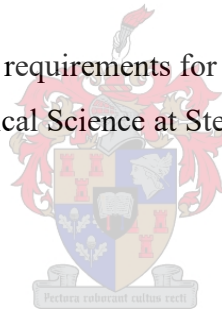**The Rise of Killer Robots: The Impact on Autonomous Weapon Systems (AWS) on the United States' Nuclear Deterrence.**

by

Matthew Robert Law

Thesis presented in fulfilment of the requirements for the degree of MA Political Science in the Faculty of Political Science at Stellenbosch University

Supervisor: Dr, Derica Lambrechts

March 2021

0

## Declaration

By submitting this thesis electronically, I declare that the entirety of the work contained therein is my own, original work, that I am the sole author thereof (save to the extent explicitly otherwise stated), that reproduction and publication thereof by Stellenbosch University will not infringe any third party rights and that I have not previously in its entirety or in part submitted it for obtaining any qualification.

March 2021

# Abstract

Technology is playing a major role in changing how states conduct warfare and much research is being focused on this. This is a broad area of study and the present research will address only one aspect of it, which is the impact of Artificial Intelligence (AI)-enabled autonomous weapon systems on a state's nuclear deterrence. More specifically, this research's main aim is to study how the rise of autonomous weapon systems could affect the US's 'second-strike' capability. This term refers to the ability of a state to strike back in response to a first attack. The US's nuclear deterrence could be affected in two primary ways. First, it could be affected practically, in that Autonomous Weapon Systems may effectively limit the US's ability to strike back. Secondly, they could affect the US's perception of its second-strike capability, meaning that the US could fear that Autonomous Weapon Systems could limit its ability to strike back, but this capability does not necessarily need to exist for this to happen. As perception is a central tenet in nuclear deterrence, the US needs only to perceive their second-strike ability to be under threat to feel insecure. The secondary aim was to see if the undermining of the US's nuclear deterrence would lead to a potential disruption in strategic stability. Strategic stability is built on the premise of states being able to successfully deter one another; undermining this would lead to instability. This thesis further chose to contextualise this study by looking at President Trump and the role of identity politics. The aim of this secondary contextualisation was to create an understanding of why the US would pursue AI-enable autonomous weapon systems and whether Trump's populist politics could explain why. This also allowed the study to better utilise Kaldor's 'New War thesis'. This sequentially allowed this study to understand the US perception of potential aggressors and its grand strategy.

The reason for addressing this area of study is because AI has a huge and unknown potential to affect all aspects of the military. This creates the need to understand how such a technology could affect nuclear deterrence, which is the cornerstone of US National Strategic Security Policy. Nuclear deterrence has played a central role in protecting the US from potential aggressors since the time of the Cold War until today and potentially it will do so for the foreseeable future. Furthermore, with nuclear weapons becoming more prevalent in the global arena, there is a need to understand how autonomous weapon systems could affect them. This research used semi-structured interviews and secondary data analysis in order to gather data. The data indicated that Autonomous Weapon Systems offer huge potential to undermine the US's nuclear deterrence in the future. They currently have significant shortcomings but with technological advancements they could have an impact on the US's nuclear deterrence. This is

because they offer the potential to undermine the US's ability to strike back. Finally, there needs to be continuous study on how Artificial Intelligence and Autonomous Weapon Systems will affect the US's military and international conflict. However, the current major threat that the US faces comes from the cyber domain. There needs to be further study into what type of cyberattack or intrusion justifies the use of kinetic force that may then allow a state to go to war.

# Opsomming

Tegnologie speel 'n belangrike rol in hoe lande se metodes van oorlogvoering verander, en baie navorsing word tans daaroor gedoen. Dit is 'n wye studieveld waarvan hierdie navorsing net een aspek sal ondersoek, naamlik die invloed wat kunsmatige-intelligensie(KI)-gedrewe outonome wapenstelsels op lande se vermoë tot kernafskrikking kan hê. Die primêre doel van hierdie navorsing is dan spesifiek om te kyk hoe die opkoms van outonome wapenstelsels die VSA se vermoë tot 'n tweede slaanaanval kan beïnvloed. Dié term verwys na die vermoë van 'n land om terug te slaan in reaksie op 'n eerste aanval. Die VSA se kernafskrikking kan hoofsaaklik op twee maniere geaffekteer word. Eerstens kan dit prakties geaffekteer word deurdat outonome wapenstelsels die VSA se terugslaanvermoë doeltreffend kan inperk. Tweedens kan die VSA se persepsie van hul vermoë om tweede slaanaanvalle te loods geaffekteer word. Hulle kan naamlik vrees dat hul vermoë om terug te slaan deur outonome wapenstelsels ingeperk kan word, maar die vermoë van sodanige stelsels hoef nie noodwendig werklik te wees om die persepsie te laat ontstaan nie. Omdat persepsie 'n sentrale aanname in kernafskrikking is, moet die VSA bloot bewus wees van die moontlikheid dat hul vermoë tot 'n tweede slaanaanval bedreig kan word om onveilig te voel. Die sekondêre doel was om te kyk of die ondermyning van die VSA se kernafskrikking sou lei tot 'n potensiële ontwrigting in strategiese stabiliteit. Strategiese stabiliteit word gebou op die aanname dat lande in staat is om mekaar suksesvol af te skrik. Om hierdie vermoë te ondermyn sou kon lei tot wêreldwye onstabiliteit. Die navorsing word verder gekontekstualiseer deur te kyk na President Trump en die rol van identiteitspolitiek om sodoende te probeer begryp waarom die VSA volhardend bly streef na die gebruik van KI-gedrewe outonome wapenstelsels, en of Trump se populistiese politiek dalk 'n verklaring hiervoor kan bied. Die kontekstualisering het dit ook moontlik gemaak vir die navorsing om Kaldor se nuweoorlogshipotese beter te kan gebruik.

Die rede vir die ondersoek in hierdie veld is omdat kunsmatige intelligensie 'n groot en ongekende potensiaal het om alle aspekte van oorlogvoering te affekteer. Dit is dus nodig om te begryp hoe sodanige tegnologie kernafskrikking kan beïnvloed wat die hoeksteen van die VSA se nasionale strategiese veiligheidsbeleid vorm. Kernafskrikking het sedert die Koue Oorlog 'n sentrale rol gespeel in die VSA se pantser teen potensiële aanvallers. Dit speel vandag steeds 'n rol wat in die afsienbare toekoms waarskynlik sal voortduur. Met kernwapens wat reg oor die wêreld meer prominent raak, is daar 'n toenemende behoefte om te begryp hoe outonome wapenstelsels hulle moontlik kan affekteer. Hierdie navorsing maak gebruik van

halfgestruktureerde onderhoude en sekondêre data-analise om data in te samel. Die data het aangedui dat hoewel huidige stelsels betekenisvolle tekortkominge het, tegnologiese ontwikkeling kan veroorsaak dat outonome wapenstelsels groot potensiaal het om die VSA se kernafskrikking in die toekoms te beïnvloed omdat hulle die VSA se vermoë om terug te slaan kan ondermyn. Laastens is voortgesette navorsing noodsaaklik om te weet hoe kunsmatige intelligensie en outonome wapenstelsels die VSA se militêre en internasionale konflik daar sal laat uitsien. Die belangrikste bedreiging wat tans die VSA in die gesig staar, kom egter uit die kuberruimte. Verdere studie is ook nodig om vas te stel watter soort kuberaanval of -indringing die gebruik van kinetiese slaankrag regverdig en 'n land in staat sal stel om 'n oorlog aan te voor.

# Acknowledgement

Firstly, I would like to acknowledge my advisor Dr Derica Lambrechts who inspired me to pursue my Master's in the first place. Thank you for all that you have done and all the effort you have put in over the last two years. Your insight and guidance were invaluable.

Secondly, I would like to thank my parents. Without their support and continuous belief in me I would not be where I am today. I am grateful for all that you have done for me.

Thirdly, I would like to thank my long-time partner Kelly White. You have continuously inspired me with your hard work and dedication through our years at Stellenbosch University. You have always kept me driven and dedicated to achieving my goal.

Finally, I would like to thank all the key informants who took time out of their day to be interviewed. All the key informants' contributions added significant value to this study. You all inspired me, and I thoroughly enjoyed your insights into my area of research.

# Contents

# List of Tables

# Chapter 1

## 1.1    Introduction

This research intends to clarify, among other things, the way that states go to war. This is based on the premise that technology changes the way in which states conduct warfare. This study will specifically focus on AI-enabled Autonomous Weapons Systems (AWS)[1], as it is not possible to include all the aspects of modern warfare. This research's main aim is to elucidate how AWS will lead to potential disruption in traditional nuclear deterrence and the subsequent effect this will have on strategic stability. The chosen research site for this study will be the United States (US). The reason the US has been chosen is because they are prepared to spend about 3.5 percent of their gross domestic production on their military. The US does this in order to maintain their military supremacy and this will become abundantly clear as this research progresses. Furthermore, a disruption of the US's nuclear deterrence is problematic for global security as multiple states require US military backing for their own security. Such states are Japan and the European Union (EU).

The reason AI was chosen was owing to its potentially disruptive and transformative capabilities. The literature of Schwab (2016), the founder and executive chairman of the World Economic Forum, highlights how AI is a part of the Fourth Industrial Revolution (Industry 4.0). Schwab (2016) states that the Fourth Industrial Revolution will be unlike anything humankind has experienced before. More specifically, when it comes to strategic stability, Altmann and Sauer (2017) indicate that AI will have a 'detrimental impact' on strategic stability and global peace.  This briefly shows why this research has chosen to look at AI-enabled AWS. This is because of their potential to be extremely disruptive to strategic stability and also because of their impact on Industry 4.0. Finally, AI is a very broad concept and it covers many different aspects from logistics to facial recognition, which is why this research will specifically look at AI-enabled weapons systems (AWS) and how they could potentially affect the US's  second-strike capability[2]. For this research to be more effective and coherent, a specific technology within AI needed to be selected. Due to this research focusing on strategic

---

[1] This research has made a decision to use the term 'autonomous weapon systems' (AWS) rather than the more commonly used 'lethal autonomous weapon systems' (LAWS). The term 'LAWS' lacks objectivity and is phrased by many theorists that advocate against 'killer robots'. This research sees 'weapon systems' as an adequate indication that these systems are built for a lethal purpose.

[2] A second-strike capability refers to the ability of a state to respond to a first-strike; it is one of the central tenets of nuclear deterrence and will be discussed in detail later on in this research.

stability as well, specifically nuclear deterrence, the area of AI that was chosen was AWS. This is important to understand as nuclear deterrence is seen as the cornerstone of the US National Security Strategy (NSS). Furthermore, AI offers capabilities that could potentially undermine the US's ability to strike back against a potential aggressor. More specifically, AI has the potential to deliver capabilities to states that may allow them to have an increased advantage against an adversary. Such capabilities will come from what is called machine learning (ML) or deep learning, which will increase the speed of conflict as they will be capable of making decisions at 'machine speeds'. Deep learning may also allow for better navigation as well as better target recognition. Finally, the invention of AI-augmented AWS will allow the US to bring mass back to the battlefield, as a fleet of AWS will potentially need fewer humans to watch over them compared to current modern drones.

One of the most problematic aspects of AI-augmented AWS, and central to this thesis, is its potential impact on nuclear deterrence and subsequently strategic stability. More specifically, this research aims to look at how AWS will effect a state's second-strike capability. The second-strike capability is an important tenet of nuclear deterrence owing to the fact that many aspects of nuclear deterrence, such as brinkmanship and Mutually Assured Destruction (MAD), rely on it. MAD allows states to deter one another and maintain stability as long as each side maintains a secure second-strike capability. It also allows them to conduct a tit-for-tat exchange in a moment of crisis until one side finds the risk too high and steps down. This is known as brinkmanship, which will be discussed in more detail in the literature review. Finally, Nuclear deterrence may seem like a distant theory utilised by Cold War superpowers, however, it remains relevant to strategic stability to this day. This is seen by the US's pledge to modernise their nuclear triad, China's dedication to building their triad, and Russia's current nuclear modernisation as well. Furthermore, there is the issue of whether the US and Russia will renew the Strategic Arms Reduction Treaty (START) and the current end of the Intermediate Range Nuclear Forces Treaty (INF).

AWS being introduced into the military sphere will ultimately change the way warfare is conducted; it may change how we perceive war and increase the possibility of war. This leads to the final theory that will be reviewed in this study: the 'new war thesis'. This will be used in order to make a contextualization and theoretical framework for this study. The main theorist of the 'new war thesis' is Mary Kaldor. The importance of the new war thesis to this research comes from the four tenets of the new war thesis: actors, goals, methods and forms of finance (Kaldor, 2013: 2). While each of these tenets will be reviewed and discussed further in the

literature review, that of 'goals' is noteworthy for this section, owing to the fact that 'goals' in the new war thesis are defined by what is known as 'identity politics'. This section and 'actors' will be the main proponent for this study's secondary contextualization. It will ultimately allow this research to understand President Trump as an actor and his populist politics. Furthermore, it will help to create a better understanding of the circumstances in which this study is taking place. One of the mainstreams of identity politics that will be looked at is 'populism'. This is a hard concept to define, yet it is currently occurring everywhere in international and domestic politics. It led to the election of President Donald Trump of the United States (US) and the referendum held in Britain to leave the European Union, commonly referred to as Brexit. Such a movement has already had a severe impact on foreign and domestic policy in both the US and Britain; it can undoubtedly be argued that it is of critical importance to understand this form of identity politics. Such forms of politics have an effect on a country's foreign policy: for example, President Trump had the largest increase in the defence budget only two years into his presidency. This contextualization will allow the study to gain a better understanding of how AWS will be pursued and create a clearer picture of the entire process. This allows for this research to create a more concise and whole analysis of how AWS will affect on the US's nuclear deterrence. Wasko-Owsiejcsuk (2018) states that Trump's doctrine and foreign policy contain doses of populism; this emerges from his use of slogans such as 'Make America Great Again' and 'Putting America First'. This contextualization will look at how identity politics can explain how the US will pursue AWS and how it makes them ultimately perceive their adversaries. Finally, this main goal of this secondary contextualization is to use the new war thesis as a research strategy, as Kaldor (2013) suggests. The new war thesis will aid this research in creating a theoretical framework and secondary contextualization that allows for the entire impact of AWS to be analysed, from the actors involved, to the goals of the US, to the modes of warfare, and how they are financed. The approach of this research is to understand how the new war thesis can help create a contextualization of how the US will pursue AWS and its subsequent effect on nuclear deterrence which in turn affects strategic stability.

## 1.2    Aims of the research

The main aim of this research is to find out how AI-enabled AWS may affect a state's second-strike capability. The secondary question will focus on how this will affect nuclear deterrence and how this will then subsequently affect strategic stability. The premise is that an insecure second-strike capability will make an ineffective nuclear deterrence which will then affect strategic stability as states successful in deterring each other with nuclear weapons equates to

stability. This will be analysed through the use of the new war thesis's four central tenets. This will be used in order to understand the dynamics of how AWS will affect nuclear deterrence. The new war thesis will form the structure for this research's theoretical framework and create a secondary contextualization for the study. This theoretical framework will help structure the research in order to understand the primary research question: How will AI-enabled AWS weapons systems affect the US's second-strike capability? There are two main issues when it comes to AWS capabilities vs perceived capabilities. This is based on what the technology can do vs what an aggressor perceives the technology can do. Perception is key; this can be further built by applying the theoretical framework created to the data collected. Finally, the main premise of this thesis is that warfare has changed and that AWS is an important aspect in this change. Furthermore, this theoretical framework has two section in it that will help create a secondary contextualization for the study. These two section are 'goals' and 'actors'. These two sections will enable this research to analyse President Trump's populist politics. This will help create a more coherent background to this study. It will also help to understand the US's perception of potential aggressors and how they will pursue AWS. The main reason for this contextualization is to better situate the readers and capture the entire phenomenon of how the US nuclear deterrence will be affected.

## 1.3    Research questions

*Main Research Question:* How will Autonomous Weapons Systems affect the US's perception and/or capability of their second-strike capability?

*Secondary Research Question*: Will this led to a disruption of traditional nuclear deterrence? or will this subsequently affect strategic stability?

*Secondary Contextualization:* The new war thesis entails 'goals' which looks at how identity politics effects conflict. In order to create a background for this research populism and President Trump will be looked at. This is a more specific form of identity politics in the US. The goal of this is to create a contextualization for this study and not to research populism in the US. This will help create an understanding of how President Trumps populist politics affects his grand strategy.

*Broader significance of the study:* Firstly, Autonomous Weapon Systems will affect the way states conduct warfare and more specifically their nuclear deterrence. It is important to get an understanding of such a phenomenon so that uncertainty about the future may be mitigated.

Secondly, populism is a dominant political movement in the US and has been expertly utilized by President Trump. Understanding such a phenomomen will allow for this study to full capture the affect of the main research question.

*How this question relates to the problem/conversation in the literature:* There is a high level of uncertainty around the issue of AWS, creating a need to fully understand the complexity and effect of such a transformative technology. Furthermore, populism is a rising political ideology that has the potential to affect international politics and US national security strategy. This means that it requires analysis in order to understand its potential affect internationally and not just domestically.

## 1.4    Preliminary literature review

The aim of this section is to briefly conceptualise the main variables of the study by conducting an introductory literature review. Here is a brief outline of all the major variables involved in this study, which were drawn from literature relevant to the field of study.

### 1.4.1    Strategic stability

Schebber (2008) and Gerson (2013) state that many theorists may refer to 'strategic stability'; however, there is not a common understanding of what exactly it is. In order to overcome this, both authors argue about the need for a historical context. Simply defined, strategic stability was based on the premise of two equally powerful nuclear armed states facing off against one another, both with the ability to retaliate by launching a second-strike back at an aggressor (Colby, 2013: 48).  A country's ability to strike back – its 'second-strike capability' – at an aggressor is the main premise that strategic stability is built on. The ability to launch a second-strike capability would deter an aggressor from attacking in fear of having the same done to them as they are doing to another, the fear creates a stabilising effect. AI can create certain grievances when it comes to strategic stability, such as an arms race or the possibility of a second-strike being cancelled out. What it is meant by a second-strike being cancelled out is how a state would possibly not be able to respond to a first-strike owing to AI-enabled AWS.

### 1.4.2    Nuclear deterrence

Nuclear deterrence is a theory of strategic stability which came to prominence during the Cold War and World War II (Morgan, 2003; & Quackenbush, 2010). Nuclear deterrence can be defined as a state taking action, the defender, that deters a possible aggressor from conducting

itself in an unfavourable manner (Powell, 2003; Quackenbush, 2010; Morgan, 2003; Mazarr, 2018; Giest & Lohn, 2018; & Wickham, 1974). Furthermore, when two states successfully deter one another, strategic stability can be achieved, this can be pursued through actions such as brinkmanship[3] which can be reached through each side having a secure second-strike capability (Powell, 2003). Nuclear deterrence then becomes a game of risk-taking, in which each side ups the risk until the more resolute state wins (Quackenbush 2010). According to Quackenbush (2010:742), successful deterrence is based on three factors: a state must persuade an attacker that it has an effective military capability; that it can use this against an aggressor; and that the threat will be carried out. Two significant factors arise from Quackenbush's (2010) literature; military capability and the use of fear when it comes to deterrence. However, none of this is important without credibility, which is important to Quackenbush's (2010) third factor. Credibility comes about through the invention and protection of a second-strike capability[4]. An important factor connected to credibility is perception, how a state views the credibility of a threat. Another important aspect of nuclear deterrence is extended deterrence. Extended deterrence involves deterring attacks on a state's allies, it is important for upholding a global security system. An example of this comes from the work of Payne (2015) who highlights how Japan planned to pursue other security options if the US nuclear umbrella disappeared.

### 1.4.3    Artificial Intelligence

The Congressional Research Service (CRS) (2019a) highlighted the issue that AI has 'significant implications' for national security. This highlights its importance, but what exactly is AI? Boulanin (2019: 13) highlights that the term 'AI' was coined in  the mid-1950s by John McCarthy. McCarthy highlighted it as the 'science and engineering of making intelligent machines'. Due to his extensive work on AI and AWS's effect on strategic stability, Boulanin's (2019) theories will be used as a basis and theorists such as Shi (2011), and Brynjolfsson and Mcaffe (2017) will be used as well. This will create a basis to understand AI, specifically neural networks, machine learning and deep learning. AI has had its ups and downs, going through periods known as AI winters. These periods are defined by low interest in the topic or they are

---

[3] Brinkmanship exists when states both have successful second-strike capabilities, it enables each side to increase the risk until the more resolute state wins. The more resolute state is the one that is willing to up the risk at any cost. The dynamics of brinkmanship are complex and will be looked at further in the review of literature.

[4] A second-strike capability is the ability for a state to strike back after receiving a first-strike from an aggressor. It is the linchpin of successful deterrence, AI-enabled AWS could have a harmful impact on keeping a second-strike capability secure.

due to insufficient hardware. However, when interest was decreasing in AI, Geoffrey Hinton pioneered the work of neural networks, which led to the creation of 'deep learning' (Boulanin, 2019: 15-16). Deep learning is the combination of 'neural networks' and 'machine learning'. The discovery of deep learning reignited the interest in the field of AI. This came about owing to the increase in 'big data'; improvements of machine learning and an increase in computer processing (Artificial Intelligence & National Security, 2018: 2). Machine learning is characterised by the development of software by humans that, once created, can learn and teach itself, no longer requiring human intervention. AI is very complex and problematic; in short, it aims to recreate human intelligence.

### 1.4.4    Autonomous Weapon Systems (AWS)

The creation of Machine Learning is having a significant impact on the sphere of military technology and it will continue to do so. As the CRS (2018) highlighted, AI will have an impact on logistics, cyber operations, intelligence gathering and analysis, information operations, command and control, and in semi-autonomous and autonomous weapon vehicles. As highlighted by the work of Boulanin and Verbruggen (2017), when it comes to AWS it is an argument about the degree of autonomy that these weapons have. Whether they are automatic, automated or truly autonomous. Boulanin and Verbruggen (2017) state that a machine's autonomy is based on its ability to go into an environment and use its ability to sense, decide and act, based on the environment. Furthermore, theorists also highlight the issue of the involvement of humans in AWS; what Boulanin and Verbruggen (2017) describe as the 'human-machine command-and-control relationship'. These are all important aspects of AWS and are aspects that will be looked into further in the literature review. Other important areas of AWS are the 'cost profile' of AWS. Cost profile is a term used to describe how much 'transformative technologies' such as AI or nuclear weapons cost to make (Allen, & Chan, 2017). This is of obvious importance as the cost profile determines who is capable of creating and getting these weapons. Furthermore, the cost profile of AI described by Allen and Chan (2017: 46) is 'diverse, but potentially low'. However, the most important aspect of understanding AWS is to build a basis of knowledge in order to understand how the 'transformative technology' can have an impact on a country's second-strike capability and ultimately change the mode of warfare.

## 1.4.5    New wars

The new war thesis is a vastly contested theory with many scholars comparing it to 'old wars' and arguing whether new wars are inherently new at all. According to Kaldor (2013), new wars are the wars that occurred during the era of globalisation. Globalisation opened up economies and weakened authoritarian states, which has led to the breakdown of states. This made the distinction between state and non-state actors, external and internal, economic and political, public and private, and war and peace hard to tell apart (Kaldor, 2013: 2). Kaldor (2013: 2) states that the breakdown of these binaries can be seen as both the cause and consequence of violence. Furthermore, Kaldor (2013: 2) defines new wars based on actors, goals, methods and forms of finance. One of the factors on which Kaldor (2013) based the new war thesis is goals. For Kaldor (2013), goals have changed from ideology to identity politics, which is of importance to this research, as identity politics is dominating the current political climate and is having an effect on the international arena. As stated earlier, the new war thesis is highly challenged. One such challenger is Booth. Booth (2001) argues whether or not new wars are actually new at all, based on the premise that what is seen in new wars can be seen in old wars. Meanwhile theorists like Shaw (2000) argue that Kaldor helps to question the current mode of warfare, which Kaldor herself argues as well. The new war thesis will be further reviewed later on. However, this section ends by stating that the new war thesis can be critical in trying to understand how AI-enabled AWS will affect nuclear deterrence.

## 1.4.6    Identity politics

As already stated, Kaldor (2013) argues that new wars are fought through the use of identity politics over ideology. This section further emphasises what exactly identity politics is and why it is of such an importance to this research. Kaldor highlights the work of Sen, on how individuals have multiple identities, when one of these identities becomes overarching, conflict will ensue (Sen, 2006 as cited in Kaldor, 2013: 338). Kaldor (2013) furthers the construct of identity and conflict by describing it as a form of binary, this is also called a friend-enemy distinction. Kaldor (2013) gives the example of a Jewish person who is defined in relation to an anti-Semite. Furthermore, conflict and violence further engrain these identities (Kaldor, 2013).  Such forms of identity can be a powerful tool for certain individuals. The literature of Fukuyama (2018) states that identity politics has moved into the global arena, as certain groups feel their identities are not receiving adequate attention. This is leading democracies to fracture into even narrower identities (Fukuyama, 2018: 93). The work of Fukuyama (2018), and Besley

and Persson (2019) state that such identities have been affected due to globalisation. For Besley and Persson (2019) this effect has led to the rise of dominant groups feeling threatened. Meanwhile, for Fukuyama (2018), globalisation has given previously invisible groups a platform where they can be seen, while Besley and Persson (2019) see these groups as the cause of Brexit and the election of President Trump as the result of these dominant groups feeling threatened. Furthermore, how can rhetoric like "Make America Great Again" help one to understand the political climate in the US and how the Trump presidency will ultimately pursue AI-enabled weapons? This means how does identity politics help understand how a presidency pursues its policies? One of the main factors that has arisen from these dominant groups is populism (Bresley & Persson, 2019). This shows how such movements have the ability to affect a state's policy and policy is subsequently important to how states pursue military strategy. Populism is hard to define and no two populist movements are the same; however, it is seen as anti-elite and the voice of the 'ordinary person'. Identity politics in the US will only be a contextual factor of this study.

## 1.5    Research design and methodology

This research method that was chosen was a qualitative approach and will be a small-n singular case study. This is owing to the fact that it will be based on semi-structured interviews and secondary data analysis. The semi-structured interview data collection process will be used in order to complement the secondary data analysis. It also allows for the research to get first-hand knowledge of the field of study from experts within the field, which is a strong advantage to have. The experts will be chosen from different areas, ranging from industry to academia. Key informants that will be looked out for are AI, cyber and military experts. The semi-structured interview questions will allow the interview to flow. Furthermore, this research will also utilise secondary data analysis. Secondary data analysis pertains to data that comes from sources such as journals, documents and opinion pieces. This study chose a singular case study owing to the fact that it will analyse the effect of AWS on the US's nuclear deterrence. This research chose the US as its case study owing to the fact the world is multipolar and there are multiple nuclear-armed states. This means that this research cannot look at all the states that have nuclear weapons. This research will view these phenomena in one specific research site as it cannot view all the capabilities of multiple different states.

## 1.6    Outline of the study

**Chapter 1:** This will begin with an introduction to the study, then the aim of the research will be discussed, followed by the research questions. After this it will conduct a preliminary literature review, followed by a look at the research design and methodology, and finally the outline of the study. This will create a rational summary as well as a roadmap to the study.

**Chapter 2:** Literature review and theoretical framework: This chapter will review existing literature in order to conceptualise key concepts of the study and also to provide a theoretical foundation for the study. Concepts that will be looked at will be: Strategic Stability; Nuclear Deterrence; Artificial Intelligence and Autonomous Weapons System; New Wars and Identity Politics. The theoretical framework refers to the theory or concepts on which the study will be based. This section will connect all the theories from the literature review together to make a coherent analytical lens for the research process. The conclusion sums up the chapter.

**Chapter 3:** This chapter will be used to contextualise the study. The main premise of this study is that technology changes the way states conduct warfare. This section will look at different major technological innovations and how these changed the way states went to war. The main technologies that will be looked at are: biotechnology, cyber, nuclear, and AI.

**Chapter 4:** Data analysis and findings: The data analysis will be based on a qualitative format. This study will use semi-structured interviews and secondary data analysis. Furthermore, this study will use a single case study with the research site being the United States. The new war thesis will be used as a research strategy to apply theories gained from the literature review to the data collected.

**Chapter 5:** This will be the final chapter of the study. This chapter will be made up of three different sections: evaluation of main findings and secondary contextualization, limitations of the study, and the conclusion and future avenues of study. The first section will focus on theory integration into that data in order to answer the main and secondary research questions. The next section aims to look at the limitations of this study. The final section concludes the research and offers possible avenues for future study.

## 1.7    Conclusion

The aim of this section was to create an outline to this study. It began with a brief discussion on why this area of study has been chosen. The aim of this was to create a basic contextualization of the study and introduce the reader to the area of research. The next part of

this chapter highlighted the aims of this research which is to understand how AWS will affect the US's nuclear deterrence. It then went on to indicate the main and secondary research questions as well as the broader significance of the study. The section that followed gave a preliminary literature review of the main theories involved in this research. This was done in order to better situate the reader in the area of study and understand what is to come. From there this section went onto the research design and methodology. Finally, the last section of this chapter created a basic outline to the study by highlighting what each chapter will entail. This research will now move onto a more in-depth look at the literature in chapter two.

# Chapter 2
# Literature Review and Theoretical Framework

## 2.1    Introduction

The aim of this literature review will be to review all the published research relating to Autonomous Weapons Systems (AWS), strategic stability, nuclear deterrence and the new war thesis. This section will review publications on these key variables in order to understand how AWS will affect nuclear deterrence which in turn will have an effect on strategic stability and international powers. The state that will be used as a case study is the US. The US will be used because the international system is now multipolar. This means that the US's actions need to be measured against the system and not just as a singular country.  The 'new war thesis' of Mary Kaldor, which is an analytical lens utilised in order to understand contemporary warfare, will be used to create a secondary contextualization for this study as well as a theoretical framework. This section's main aim is to review the literature on these dominant themes in order to build a basis to understand how AWS will affect nuclear deterrence and strategic stability. It then aims to create a secondary analysis and background for the study by reviewing the literature on the new war thesis and President Trump's populist politics.

This review of literature will start off with an overview of the strategic stability of the US. This is one of the main variables when it comes to AWS, as this research is trying to forecast the effect of AWS on strategic stability. The section that follows is nuclear deterrence, which is another important variable to the study, which looks at the effect of AWS on nuclear deterrence that will in turn affect strategic stability. This next section will discuss Artificial Intelligence (AI) which is the basis of AWS and it is thus of significant importance. Following this, the issue of AWS will be discussed as well as strategic stability, bringing together AWS, strategic stability and nuclear deterrence. This will be done in order to bring together the literature of the main variables of this study. The section that will follow is the 'new war' thesis. There are four main tenets of the new war thesis: actors, goals, methods and forms of finance. One of the main theories within the tenet 'goals' is identity politics, which is an important theory and is extensively spoken about by Kaldor. The next section  will focus on 'identity politics'. Identity politics is an important theory to review literature on, due to its extensive effect on domestic politics and global politics recently. Meanwhile, the form of identity politics known as 'populism' that is plaguing the US is not so easily defined. However, populism and identity politics are of significant importance, as such phenomena have an adverse effect on how

countries conduct their policy; this will be analysed further in the identity politics section. It is also of significance importance when it comes to creating context for this study, which will help to better situate this research.

## 2.2    Strategic stability

Schebber (2008) and Gerson (2013) both state that many theorists refer to 'strategic stability'; however, many neither define it or, even more problematically, there is no common understanding of it.  Both Schebber (2008) and Gearson (2013) state that there is a need for a historical analysis of strategic stability in order to gain an understanding of what it is. Gerson (2013: 2) goes on to say that there is also a gap when it comes to how nuclear-armed countries view and define the requirements for stability. Another issue bought forward by Gerson (2013: 2) is how the world is now multipolar. This means that stability has changed since the end of the Cold War, with a move away from a bipolar power configuration to a multipolar arena. Strategy stability during the Cold War was built on the logic that if both sides had the ability to strike back effectively after an attempted disarming first-strike this would create a stabilising effect (Colby, 2013: 48). Strategic Stability in the Cold War was known as nuclear deterrence and it is still the bedrock of US strategic security policy. Furthermore, Colby (2013) brings forward one of the most important tenets of strategic stability and nuclear deterrence, which is the ability of a state to strike back. This is known as a 'second-strike capability' and is arguably the most important tenet of nuclear deterrence. This section's aim is to create a basic understanding of strategic stability; it will mention tenets and aspects of nuclear deterrence, which is unavoidable as the two are interconnected.

The previous stability metrics were built on the tenets that there were two equally aggressive powers, each equipped with large nuclear forces and extensive defence strategies with each fearing a 'bolt out of the blue' strike from the other side (Scheber, 2008). This was the basis of stability metrics in a bipolar world. Currently there are multiple states that have nuclear capabilities and other states pursuing such capabilities. According to the Federation of American Scientists (FAS) there are currently nine states with nuclear weapons: Russia, US, China, France, UK, Pakistan, India, Israel and North Korea (Kristensen & Korda, 2020). Out of these states, 91 percent of all nuclear warheads are owned by Russia and the US (Kristensen & Korda, 2020). Furthermore, strategic stability gets more complicated based on the issue of legitimate threats and how states perceive these threats. Scheber (2008) supports this statement based on his theory that deterrence was built on 'punitive threats', which is not helpful when it

comes to actors such as Iran or North Korea. The way states dealt with such issues, according to Gerson (2013), was because stability was built on the freedom from a surprise attack. This is based on the inspection of a potential enemy and the strength of their forces, and subsequently it is also based on the vulnerability of one's ability to strike back (Gerson, 2013). What Gerson (2013) means is that a state needs to have reliable retaliatory forces in order to respond to a first-strike from an aggressor. If both sides have reliable retaliatory forces, this gives the aggressor a strong reason not to attack; this means that a retaliatory capability has a stabilising effect. A secure second-strike capability is an important tenet of nuclear deterrence and this will be discussed in the next section.

Altmann and Sauer (2017: 119-120) state that stability has two dimensions, the first one being 'military stability' and the second being 'crisis stability'. Military stability and crisis stability are interlinked to one another. According to Altmann and Sauer (2017: 119), 'military stability', is built on the issue of proliferation of arms and the emergence of an arms race. The role of new technologies can be further destabilising when they offer a qualitatively clear advantage and are close at hand (Altmann, & Sauer, 2017: 120). Furthermore, when an adversary deliberately pursues such a technology there could be an increase in mutual observation and uncertainty (Altmann, & Sauer, 2017: 120). This shows that the countries pursing AI capabilities can be problematic for military stability. The other issue of stability outlined by Altmann and Sauer (2017) is crisis instability and escalation. Crisis instability and escalation refer to either a move from peace to war or when war has broken out and there is a move from conventional to nuclear weapons (Altmann, & Sauer, 2017: 120). These two different dimensions, according to Altmann and Sauer, are interlinked; new weapons developed out of an arms race can subsequently cause crisis instability (Altmann, & Sauer, 2017: 120). What this means is that the mere pursuit of a new technology that gives an adversary an advantage can be destabilising and cause instability, triggering military or crisis instability. What is more alarming is that AI is more than just a new technology; it could potentially be the fourth and final industrial revolution. This promises a technology that could possibly cause a considerable amount of military and crisis instability due to its potential and its uncertainty.

What was gained from this section was a basic understanding of strategic stability. Strategic stability was a theory of stability that arose during the Cold War; it helps theorists to understand how hostile super powers with nuclear capabilities maintained stability. Stability seems to be defined by states having a second-strike capability. What a second-strike capability gave states was an ability to eradicate their vulnerability by having the ability to strike back after a first

strike. A first strike is also referred to as a disarming strike, as stated by Gerson (2013). What this gave states was a degree of legitimacy, what Scheber (2008) called a 'punitive threat'. Furthermore, a second-strike capability gave the aggressor a reason not to attack. A second-strike capability has more benefits for stability metrics; however, this will be discussed in the next section   on nuclear deterrence. It is important to look at this concept in detail owing to the fact that this research is aiming to see the effect of AWS on nuclear deterrence and its subsequent effect on strategic stability.

### 2.2.1    Nuclear deterrence

Nuclear deterrence was a dominant strategy during the Cold War when it came to a state's national strategic security policy and maintaining strategic stability. Podvig (2012) states that strategic stability is achieved when a state is assured that an adversary will not undermine their nuclear deterrent capability. This rudimentary explanation helps link nuclear deterrence to strategic stability and this section now considers nuclear deterrence. Morgan (2003), a theorist highlighted in the work of Quackenbush (2010), states that deterrence is hard to explain and became an area of study during World War II (WWII) and the Cold War. Prominent theorists, when they define nuclear deterrence, simplify it as a state taking action that threatens another state, the aggressor, if they act in an unfavourable manner, thus preventing them from acting on an unwanted action (Powell, 2003; Quackenbush, 2010; Morgan, 2003; Mazarr, 2018; Giest & Lohn, 2018; & Wickham, 1974). Furthermore, nuclear deterrence was not just a theory in the study of international politics, it was an actual national strategic security policy employed during the Cold War until today. It remains the foundation of US national security strategy under the Trump administration. Morgan (2003: 86) states that the role of nuclear deterrence was at the 'heart' of every major nation's 'national security strategy'.  Its importance in the Cold War can be seen in the declassified document of former United States National Security advisor General Brent Wickham (1974). Wickham in his memorandum on '*Nuclear Weapons Employment Policy'* defines deterrence as:

> "The principal objective of US strategy is deterrence of nuclear and conventional attacks or attempts at coercion under a threat of nuclear or conventional attacks against the United States, its allies and any nation whose security is vital to the US interest." (Wickham, 1974).

The 2017 National Security Strategy (NSS) states the same sentiment as Wickham's memorandum, that nuclear weapons prevent acts of aggression or attack by nuclear weapons,

non-nuclear strategic attacks and large-scale conventional aggression (2017). It further states that the US has nuclear weapons so that its allies do not have to and this helps ensure their security (NSS, 2017). This briefly shows the importance of nuclear deterrence to US NSS policy from the Cold War era to the Trump administration. This section now looks at key tenets of nuclear deterrence to build a deeper understanding of this strategy, allowing the research to progress to understanding how AWS will effect nuclear deterrence.

### *2.2.2    Diplomacy of violence: Fear and perception*

One of the most important tenets of deterrence is perception and fear. A state must show that it has both the will and military might to execute a high-cost retaliation on an aggressor conducting an unfavourable action or attack. Morgan (2003), referenced by theorists such as Quackenbush (2010), states that deterrence is the action of a state preventing another state from executing an action they do not agree with; this is done by threatening unimaginable damage on that state. Long (2008: 7) states that the central premise of deterrence is the generation of fear. Long[5] (2008) states that this is done by imposing a high cost as a result of unfavourable action being taken by an adversary on the defending state. Fear plays a vital role in a state's successful deterrence, it must convince the aggressor that it will impose a high cost on any action that it deems to be an act of aggression. The defending state must convince the aggressor that this cost will be carried out; the aggressor may just need to perceive the threat as real to step down. Perception plays a key role in nuclear deterrence and subsequently this thesis and this will be discussed further.

Another theorist that, like Long, equates deterrence to the use of fear, is Schelling (1966). Schelling (1966) states that fear is used in order to deter an aggressor. This fear is generated by the threat of an unfavourable outcome or punishment on the aggressor by the defender, which will therefore, stop an actor from acting in an unfavourable way out of fear of the consequences (Schelling, 1966: X) For Schelling, military strategy has become what he terms 'The Diplomacy of Violence' which is a form of bargaining power based on fear and military might (Schelling, 1966). For Quackenbush (2010: 742), a state must perceive that an attack will be carried out; this is done by a state showing that: "1) it has an effective military capability; 2) that it could impose unacceptable costs on an attacker, and 3) that the threat

---

[5] It must be stated that Long is an influential theorist owing to the fact that Long belonged to Research and development Corporation (RAND), an influential institution on the US's deterrence strategy (2008).

would be carried out if attacked.".". These are all factors which can fit into a state's 'diplomacy of violence'. This diplomacy of violence creates a narrative for the aggressor to see that a state maintains a degree of might and that this might will be used in order to deter possible aggression or unfavourable behaviour. Diplomacy of violence is based on the perception of the aggressor; a threat needs a degree of credibility for it to be successful, meaning that the aggressor needs to perceive that the threat will be carried out.

### *2.2.3    Credibility: How to secure an aggressor's perception*

Deterrence is based on fear and perception. For this to be effective, there needs to be credibility. An aggressor must perceive that a state has both the will and military might in order for it to back down in a tit-for-tat exchange in a moment of crisis. This can be done through credibility; credibility is based on a number of tenets and will be discussed in this section. Credibility, according to Powell (2003: 89), came through the technological innovation known as a 'second-strike capability'. According to Powell (2003), this gave states the ability to put pressure on each other. This can be done by leaving certain aspects to chance, which is why the nuclear forces of each state do not cancel each other out (Powell, 2003). Also, the more that is at risk, the more the state will be willing to risk; this is known as 'brinkmanship' (Powell, 2003). Powell (2003) states that brinkmanship is a model that allows states to employ coercive pressure on each other if both sides have a second-strike capability. Quackenbush (2011), who references the work of Powell, also states that brinkmanship is the strategy in deterrence where it becomes the competition of risk-taking. Risk-taking is a strategy in brinkmanship that involves the act of upping the risk until one state backs down, this is what Powell refers to as the 'Dynamics of Brinkmanship' (Powell, 2003). Powell (2003) states that the dynamics are based around the issues of a resolute state.

For Powell (2003: 91), a crisis will only arise if there is a substantial level of uncertainty of the 'balance of resolve'. Escalation, according to Powell (2003: 91), is then dependent on the complex interaction between a state's level of resolve and that state's uncertainty about another state's resolve. A state's level of resolve determines whether it will acquiesce or continue to up the risk (Powell, 2003: 92). An equilibrium is therefore made up of the state's resolve, its perception of the other state's resolve, and uncertainty regarding its own resolve (Powell, 2003: 96). Ultimately, the more resolute state will up the risk and the less resolute state will find the risk too high and back down (Powell, 2003). The argument is used in order to express why states' nuclear forces don't cancel each other out; certain states are more resolute and willing

to up the risk of war, meaning that Powell's theory helps with the issue of uncertainty. Even if there may be uncertainty, how resolute a state is will create credibility. Morgan (2003) also highlights the importance of credibility when it comes to successful deterrence. According to Morgan (2003: 4), the best strategy for credibility is convincing the enemy of your military capability, the costs that will be inflicted and the willingness to inflict these costs. Long (2008), like Powell (2003), also highlights 'credibility' as the 'linchpin' of deterrence; however, along with this he states that threats are hard to estimate in practice. A net assessment can be conducted in order to understand a nation's credibility. The assessment contains three elements: "aggregate forces, proximity, and power-projection capability" (Long, 2008: 11). However, these are very hard to measure, thus making uncertainty high and sustaining credibility as problematic. But for Schelling (1966), uncertainty makes a threat credible, a response that carries some risk of war is credible, even if war seems unreasonable or implausible. This manipulation of risk and that the slightest possibility of risk may equate to war shows how Schelling (1966) views brinkmanship. This subsequently shows the complexity of credibility along with a state's perception of this risk or its ability to manipulate this risk. However, there are ways so ensure risk; namely a second-strike capability.

This section, thus far, has highlighted two concepts that are central to this thesis: the issue of a second-strike capability and its credibility. These two central tenets give states that ability to conduct brinkmanship. Brinkmanship is the action where states manipulate risk based on how resolute they are and their perception of how resolute their enemy is. Perception is key when it comes to deterrence, an enemy just needs to perceive a threat to be credible and escalation can ensue. Mazarr (2018: 9) states that the most important tenet of deterrence is perception, on the basis of whether a state sees a need to act on aggression or not. Mazarr (2018: 10) states that the potential aggressor must have the perception that the defender has the capability and will to proceed with its threat. This is what, as already mentioned, Powell (2003) would call a resolute state. Morgan (2003) would also agree with Mazarr (2018) as the credibility of an attack is important to a successful deterrence. For Mazarr (2018: 10) this highlights two important factors, capability and will; perceived weakness in either of these factors could equate to undermined deterrence. Payne (2015) states that deterrence is constructed by human perception and calculation which is affected by multiple factors that are beyond 'confident prediction'. A state's perception of an adversary's capability being threatening to it could have adverse effects on nuclear stability; an adversary may just perceive that their nuclear missiles are at threat and this could be destabilising (Giest, & Lohn, 2018).

### 2.2.4    Building credibility

A way in which credibility can be built, according to Dannereuther (2007), is by arms control. This is built on the premise that the fundamental factor that nuclear deterrence is built on is MAD. By implementing arms control, MAD can be turned into Mutually Assured Vulnerability (Dannereuther, 2007: 229). This was achieved by creating the Anti-Ballistic Missile treaty of 1972 which assured both superpowers, the US and the former Soviet Union, that they both had assured second-strike capabilities (Dannereuther, 2007: 229). Another example of an arms treaty was the Nuclear Non-proliferation act between Nuclear Weapon States and Non-Nuclear Weapon States, which provided certain countries with the safety of not pursuing nuclear weapons and the promise they would not be attacked by states with nuclear weapons (Dannereuther, 2007: 230). This reassurance of a second-strike capability being secured created a degree of credibility between superpowers, the credibility in terms of them being able to retaliate. With this mutually assured vulnerability came about a stabilising effect; showing that arms control legislation can create credibility and more importantly create stability. Furthermore, it gave states without nuclear weapons a reason not to pursue them and a subsequent safeguard allowed for the limitation of states with nuclear weapons. This research notes that there are now are nine different states that have nuclear capabilities, based on the arms control association statistics (2019). This makes the dynamics of brinkmanship and nuclear deterrence more complex. As noted already, deterrence used to be based on a bipolar model. Furthermore, the concept of MAD is highlighted here. This tenet is of importance to nuclear deterrence as MAD is built on a second-strike capability, if a state were to lose their second-strike capability due to AWS it could force them to fear their ability to strike back against an aggressor. This shows that AWS could potentially threaten MAD, which in turn could have a destabilising effect.

### 2.2.5    Extended deterrence

There is a further explanation of deterrence which is termed 'extended deterrence'. There are two types of deterrence; direct or 'extended' deterrence of other countries by a state with nuclear capability (Mazarr, 2018).

> "Deterrence can be used in two sets of circumstances. Direct deterrence consists of efforts by a state to prevent attacks on its own territory—in the US case, within the territorial boundaries of the United States itself. Extended deterrence involves discouraging attacks on third parties, such as allies or partners. During the Cold War,

direct deterrence involved discouraging a Soviet nuclear attack on US, territory; extended deterrence involved preventing a Soviet conventional attack on North Atlantic Treaty Organization (NATO) members" (Mazarr, 2018: 3).

Morgan (2012: 87) states that extended deterrence became central in international politics, in the form he calls 'extended protection', this was for states and non-state actors. Morgan (2012: 87) highlights the United Nations (UN) Security Council as an example of extended deterrence, it also involves: "alliances, interventions, arms transfers, power projection efforts, military training programs and non-proliferation pressure.". Morgan (2012: 94) highlights that during the Cold War extended deterrence entailed a singular state projecting its decisions onto others and preventing them from getting weapons of mass destruction; however, a 'collective actor deterrence' is more relevant today. Deterrence is now about upholding a global security system through a collective effort, which means that deterrence needs to be adjusted (Morgan, 2012: 94). Huth (1988: 82), in a review highlighted as a critical study by RAND, states that extended deterrence is a confrontation between states in which the defender threatens force against the aggressor to prevent them from using military force against an ally or territory under the control of an ally. Payne (2015) highlights the importance of extended deterrence to the US and its allies, highlighting how Japan stated that they would pursue other security options if the US 'nuclear umbrella' disappeared. This shows that extended deterrence was an important factor during the Cold War. It is of importance to understand, however, that this research focuses on a US-centric view. Although extended deterrence is still of importance to US NSS policy, this thesis aims to look at direct deterrence and not extended deterrence.

### 2.2.6    Nuclear deterrence: What was important?

This section aimed at reviewing the literature on nuclear deterrence. Key tenets that emanated from the literature were: second-strike capability, credibility, arms race, arms control and brinkmanship. The technological innovation of second-strike capability allowed states to contain the ability to strike back against an aggressor's first-strike. This ability to retaliate gave states the ability to manipulate risk. This was done by states taking steps that increased the level of risk until the more resolute state won and one state backed down; this manipulation of risk is termed brinkmanship. Both states having a second-strike capability is called MAD, as the escalation to war combined with a state's ability to strike back would equate to both being destroyed. Furthermore, arms control allowed this to happen, as a state's second-strike capability was secured by the eradication of anti-ballistic missile defences. Furthermore, one

of the key issues to nuclear deterrence and brinkmanship is credibility. Credibility is influenced by a state's perception of a threat's certainty. If a state knows for certain that if they act in an unwanted way the defender will act in a way that inflicts an unacceptable cost, the state will not act in this unfavourable manner. Furthermore, a state's capability is of importance, as this is important to the credibility of threat. In summary, there are three tenets of brinkmanship: *military capability*, which allows for the ability to inflict an *unacceptable cost* and this must carry a level of *credibility* to be effective. Most importantly, the main issue that can be emphasised from this section is the importance of a state's second-strike capability and the credibility this strike has.

### 2.2.7    Artificial Intelligence and Autonomous Weapons Systems

Now that nuclear deterrence has been examined, this section aims to look at what AI and AWS are. A CRS (2018) report on 'Artificial Intelligence and National Security' highlights that there is significant implication for national security when it comes to Artificial Intelligence. Due to this, the US Department of Defense (DOD) and other nations are developing AI for a range of military applications (Artificial Intelligence and National Security, 2018). The domains in which AI is being developed in, according to the report, are logistics, cyber operations, intelligence gathering and analysis, information operations, command and control, and in semi-autonomous and autonomous vehicles (Artificial Intelligence and National Security, 2018). This section will review the literature on AI by beginning with a brief review of its origin, what exactly AI is; and what Machine Learning is; Deep learning, and the issue of autonomy when it comes to AI and Autonomous weapons systems.

What is the origin of AI? The origin of AI has been referenced to the Dartmouth conference of the mid-1950s and is attributed to John McCarthy (Brynjolfsson & Mcaffe, 2017; Ng, 2019; & Boulanin, 2019: 13). The concept of AI that was coined by McCarthy was based on the science and engineering of making machines with intelligence (Boulanin, 2019: 13). AI has faced ups-and-downs when it comes to development, due to issues such as funding and hardware not being capable of running such software. But one of the most significant breakthroughs in AI was the work of Hinton, who was one of the last scholars to focus on the issue of AI, as interest was being lost in it (Boulanin, 2019: 15-16). What Hinton had theorised came to be known as 'deep learning', which combined 'neural networks' and 'machine learning' (Boulanin, 2019). This discovery ended what can be termed an 'AI winter' and sparked interest in AI again, due to its possible military and civilian applications. The CRS (Artificial Intelligence and National

Security, 2018: 2) highlighted the origin of AI in the 1940s and the re-emergence of AI around 2010 due to three different factors: "(1) the availability of 'big data' sources, (2) improvements to machine learning approaches, and (3) increases in computer processing power." The Artificial Intelligence and National Security (2019) report shows similar evidence to that of Boulanin (2019), as machine learning (ML) came about due to these factors highlighted by the report, showing that the re-emergence of AI happened around 2010. This also highlights the importance of machine learning, an important aspect of 'deep learning', which this review will come to after machine learning. ML benefits come from its capability to improve its performance without human intervention, as it has the ability to learn and improve itself (Brynjolfsson, & Mcaffe, 2017: 2). Finally, Brynjolfsson and Mcaffe (2017) highlighted the importance of ML, stating that machine learning has become more effective and more available in the last couple of years.

The two major tenets that have come from the literature so far are ML and deep learning. Ng (2019) emphasises the importance of 'deep learning' by citing Hinton and how it sparked a second war of development in neural networks. Deep learning is an important technological innovation in recent years, it has led the CRS to claim that AI is one of the key factors in the US being able to fight and win wars of the future (Artificial Intelligence and National Security, 2018: 1). However, it is not just the US that believes in the capabilities, so do its adversaries. The CRS report highlights how China aims to be in the lead of AI development by 2030 and so does Russia (Artificial Intelligence and National Security, 2018: 1). Such great power competition to achieve general AI could lead to a possible arms race, as AI is viewed as having the ability to fight and win wars of the future. This issue of an arms race will be further discussed in the strategic stability section. Finally, the importance of perception of an adversary's capabilities must be emphasised when it comes to the issue of nuclear deterrence. States must merely view their adversary as having advanced capabilities in order to feel insecure. This brief introduction now leads this section to look at what exactly AI is.

### 2.2.8    What is AI?

Shi (2011: 1) defines AI as: "the science and engineering of imitating, extending, and augmenting human intelligence through artificial means and techniques to make intelligent machines". Cummings (2017: 2), in his book *AI and the Future of Warfare* states that there is no commonly agreed on definition of AI, however, there is a general definition: "the capability of a computer system to perform tasks normally requiring human intelligence, such as visual

perception, speech recognition and decision-making". The aim of AI research is to make machines capable of mimicking human intelligence by being able to view the world through learning, vision or processing natural language (Boulanin, 2019: 13). Furthermore, Boulanin (2019): 13-14) takes intelligence further by separating AI into two different levels: 'Artificial General Intelligence' and 'Narrow Artificial Intelligence'. Artificial General Intelligence, according to Boulanin (2019: 13-14), is human-level intelligence and the issue of it ever being reached is questionable. Narrow AI has been around for a while and is subsequently widely used. They are complex software programs that can complete intelligent tasks; however, they are brittle in nature (Boulanin, 2019: 13-14). Shi (2011) also categorises AI into 'narrow' and 'general' intelligence. Narrow intelligence is the ability of AI to process data and make decisions. In the age of 'big data', AI is thriving (Shi, 2011) The work of Geist and Lohn (2018: 12) highlights general intelligence as 'super intelligence'. This is a point where machines will outmatch human intelligence. Garnham (1987) sees AI as the study of intelligence, he argues for the link between AI and psychology and how both are interdependent. Owing to the fact that AI aims to recreate human intelligence, what can be seen thus far is how AI seems to already be divided into two sections based on how different theorists define its ability of intelligence. This degree of intelligence has a direct correlation to human intelligence, as AI aims to recreate human intelligence or better it. These sections will be defined as narrow and general in this research. Narrow refers to the machine's basic abilities to process data and make decisions. General intelligence is based on a machine's ability to match a human's intelligence. The concept of a machine being intelligent is difficult as it currently stands, therefore, for this research to define 'general intelligence' as a machine's ability to be more intelligent than a human seems to be ahead of its time.

### 2.2.9    Towards a more coherent understanding of ML and deep learning

This section of the literature review now aims to further the understanding of machine learning and deep learning.  This section will first begin with ML. Boulanin (2019) states that ML is the approach to software development where the system is built first so that it can learn and then teach itself using a variety of methods such as supervised learning, unsupervised learning and reinforced learning. This allows for these programs to not be hand-coded by humans which leads to hard-code (Boulanin, 2019: 15). These programs can subsequently code themselves and do not need human intervention. This is important as hand-code can be very complex and difficult as the environment gets more complex (Boulanin, 2019: 15-16). When it comes to ML capabilities, Boulanin and Verbruggen (2017) state that ML has the ability to find statistical

relationships in large sets of data; the more data there is the more the machine will learn. Furthermore, according to Boulanin and Verbruggen (2017: 17), ML is not new, it has in fact been around for decades and has taken great strides owing to improvements in computer power and the growth of deep learning. One of the capabilities ML allows is better pattern recognition, which in turn, allows for the improvement of navigation and target recognition (Boulanin, & Verbruggen, 2017: 17). This starts to highlight the possible military applications for machine learning. What if ML gave a state the capability to increase its navigation and target recognition which could allow it to successfully hunt another state's second-strike capability such as a Ship, Submersible, Ballistic, Nuclear Submarine (SSBNs)?

This section will now move on to deep learning. One of the most important aspects of deep leaning, apart from machine learning, is a neural network. A neural network and machine learning put together is what is known as 'deep learning'. This begs the question, what exactly is a neural network? According to Schmidhuber (2014) a neural network consist of many processors called neurons, which are either activated through sensors perceiving the environment or other neurons that have been activated by these sensors. Boulanin (2019: 17) states that neural networks draw on knowledge of "the human brain, statistics and applied math". Boulanin (2019) gives an in-depth definition of what deep learning is:

> "Deep learning is a type of representation learning, which in turn is a type of machine learning. Machine learning is used for many but not all approaches to artificial intelligence. Representation learning is an approach to machine learning whereby the system 'learns' how to learn: the system transforms raw data input to representations (features) that can be effectively exploited in machine-learning tasks. This obviates manual feature engineering (whereby features are hard-coded into the system by humans), which would otherwise be necessary. Deep learning solves a fundamental problem in representation learning by introducing representations that are expressed in terms of other, simpler representations. Deep learning allows the computer to build complex concepts from simpler concepts. A deep-learning system can, for instance, represent the concept of an image of a person by combining simple concepts, such as corners and contours. Deep learning was invented decades ago but has made important progress in recent years, thanks to improvements in computing power and increased data availability and techniques to train neural networks" (Boulanin, 2019: 17).

Deep learning is one of the most important spheres of AI application that can have an effect on military technology; the other being autonomous weapons systems which has given rise to the evolution of ML and deep learning. But before autonomous weapons systems are looked at, what are the challenges and opportunities for ML and deep learning? Boulanin (2019: 17) boldly states that ML does not need to demonstrate its potential as it has already allowed computers and robots to perceive the world better. It has also accelerated the development of autonomous systems like self-driving cars and voice assistants (Boulanin, 2019: 17). For Horowitz, Allen, Saravalle, Cho, Frederick and Scharrre (2018), in the defence domain, AI and machine learning will allow for more challenges in a wider range of environments to be tackled. ML, according to Boulanin and Verbruggen (2017: 17), has been around for decades; however, it has made recent improvements owing to the increase in computer power and the development of 'deep learning'. Payne's (2018) research also highlights how the advancements in AI and neural networks are of importance to military strategists and will have a subsequent effect on strategic affairs. Owing to this, ML has allowed for better pattern recognition and targeting (Boulanin, & Verbruggen, 2017: 17).

Furthermore, as stated in the previous section, ML and deep learning offer great potential. However, some problems are highlighted by Boulanin and Verbruggen (2017: 17), the first being a requirement issue in ML, the issue of data and the need for large sets of data. Another important issue of ML, if not the most, is the issue of a neural network being a 'black box' (Boulanin, & Verbruggen, 2017: 17). Brynjolfsson and Mcaffe (2017) highlight how deep learning is hard to diagnose and correct, due to the inability to find out what goes on within the neural network. "The input and the output of such a system are observable, but the computational process leading from one to the other is difficult for humans to understand" (Boulanin, 2019: 20). This makes it important for regulators, users and developers to find a way to responsibly adopt and use this technology, which according to Bouanin (2019: 20-21), they are currently struggling with.

## 2.2.10   The issue of autonomy and uncertainty

However, one of the important aspects of ML is its ability to develop autonomy. What is autonomy? For Boulanin (2019: 21) autonomy can be based on a machine's ability to execute a task without human input by using sensors, computer programming, and an actuator for that environment. The first issue of autonomy that will be looked at is the debate by theorists as to whether AI is currently autonomous or automated. For example, Cummings (2017) highlights

that, in his view, current systems are more automated than autonomous, as they require serious human intervention. Cummings (2017: 8) further states that current ML and deep learning only have the capability to detect patterns that are significantly tuned by humans and must also be interpreted by humans for them to be useful. Cummings (2017) plots how the more uncertainty involved in the task of a skilled based behaviour, the harder a task becomes to code. There are two different factors when it comes to uncertainty; rule-based behaviour and knowledge-based reasoning; as uncertainty increases, rules-based reasoning gives way to knowledge-based reasoning (Cummings, 2017: 6). What this means is that AI is easier to program when a skill-based behaviour is quantifiable, the more qualitative or abstract it is, the harder it is to program (Cummings, 2017). For Boulanin and Verbruggen (2017: 5), autonomy is usually understood as the ability for hardware or software, once activated, to perform functions or tasks on its own. However, Boulanin and Verbruggen (2017) state that autonomy is defined differently over different disciplines, but they chose how computer science, robotics and engineering define it by quoting the work of Paul Scharre. Boulanin and Verbruggen (2017: 5) define autonomy into three different categories: "(a) the human-machine command-and-control relationship; (b) the sophistication of the machine's decision-making process; and (c) the types of decisions or functions being made autonomous".

*The human-machine command-and-control relationship* involves assessing the extent to which humans are involved in the execution of tasks (Boulanin, & Verbruggen, 2017: 6-7). There are a further three categories: 1. Human involvement at some stage of the task execution (semi-autonomous/Human-in-the-loop), 2. Human oversight in case an error occurs (human-supervised autonomous/human-in-the-loop), and 3. Fully autonomous machines that do not require human intervention (fully autonomous/human-out-of-the-loop). 'Sliding Autonomy' is where a machine can go between human supervision and full autonomy. (Boulanin, & Verbruggen, 2017: 6-7).

*The sophistication of the machine's decision-making process* involves the machine's ability to execute self-governance and deal with uncertainties in the environment it operates in (Boulanin, & Verbruggen, 2017: 6). There are a further three categories here: automatic, automated and autonomous ((Boulanin, & Verbruggen, 2017). Automatic involves the machine responding to a sensory input by following a set of rules with no uncertainty involved, while automated or autonomous have self-governance and respond to the environment (Boulanin, & Verbruggen, 2017) These are conceptually challenged terms. Once again the idea is that automated is based on a set of rules, which the machine must use to respond to the environment.

This has the ability to make the outcome inevitable, meaning there is an ability to predict behaviour and lower uncertainty. Furthermore, autonomous can select a range of outcomes based on the data from sensory input received but there may be some human intervention.

The final category is *the types of decisions or functions being made autonomous.* This deals with the issue of the task that is being executed over the issue of autonomy of the system as a whole (Boulanin, & Verbruggen, 2017: 6). For Boulanin and Verbruggen (2017) autonomy is best understood in terms of what task is being executed at a function level or subsystem level (6). Some functions may not have any risk or ethical issues, while a function such as targeting may cause great ethical or risk issues (Boulanin, & Verbruggen, 2017: 6). This last category, favoured by the authors, is the 'functional approach' when it comes to autonomy (Boulanin, & Verbruggen, 2017: 7)

> "It recognizes that the human–machine command-and control relationship and the sophistication of a machine's decision-making capability may vary from one function to another. Some functions may require a greater level of self-governance than others, while human control may be exerted on some functions but not others depending on the mission complexity and the external operating environment, as well as regulatory constraints. Also, the extent of a human operator's control or cancel functions may change during the system's mission." (Boulanin, & Verbruggen, 2017: 7)

This means that their research was based on the study of the autonomy of weapons systems and not the study of autonomy inside weapons systems. For Bounanin and Verbruggen (2017: 7), this allows for the research of a larger range of weapons systems. The authors further describe autonomy as a machine's ability to transform data from the environment into a set of plans or actions to execute (Boulanin, & Verbruggen, 2017: 7) This literature states that autonomy always involves the same three capabilities being integrated into one; the three being sense, decide and act. Ultimately, autonomy comes from the ability of a system to sense and act to an environment in order to achieve its goal (Boulanin, & Verbruggen, 2017: 11).

Bieri and Dickow (2014: 2) break autonomy down into three different degrees: remote control, autonomous manoeuvres under human steering control, and autonomous execution of tasks without human control, but with a veto right. As with the views of Boulanin and Verbruggen (2017) discussed above, Bieri and Dickow's (2014) degrees of autonomy are based around human intervention. When it comes to 'remote control' a robot executes missions with a distant human operator; this is done by the operator helping to reduce the complexity (Beiri, &

Dickow, 2014: 2). 'Autonomous manoeuvres under human steering control' involves a machine that operates autonomously; however, it can be overridden by a human being at any point, for example: changing the course of the flight (Beiri, & Dickow, 2014: 2). 'Autonomous execution of tasks without human control, but with veto right' is when a human can only intervene with a veto command, like activating an emergency stop button (Beiri, & Dickow, 2014: 2). These all have different degrees of autonomy; however, they are maintaining some type of human-in-the-loop scenario. Subsequently, Beiri and Dickow (2014) do not deal with the issue of autonomy, for example 'Autonomous manoeuvres under human steering control' describes a machine that Boulanin and Verbruggen (2017) would have arguably seen as 'automated' rather than autonomous.

This can be stated owing to the fact the machine is flying based on a 'predetermined routes', meaning that the machines have been programmed to operate in a certain environment. The AWS is not collecting data from its environment through sensors and making decisions based on a system like deep learning and executing the task according to the outcome of the neural network. Noone and Noone quote the US Department of Defense definition of AWS:

> "[A] weapon system that, once activated, can select and engage targets without further intervention by a human operator. This includes human-supervised autonomous weapon systems that are designed to allow human operators to override operation of the weapon system but can select and engage targets without further human input after activation" (DoD, 2012 as cited by Noone, & Noone, 2015: 27).

Noone and Noone state that they view autonomy closer to being automation (2015) as an automatic system will carry out a task based on a preprogrammed sequence of operations or move in a structured environment (Noone, & Noone, 2015: 27-28). Noone and None then go on to discuss how AWS should be able to select and engage a target in an unstructured environment without human involvement (2015). Noone and Noone are dealing with the issue of whether weapons systems should have a closed loop system, meaning that there is no human intervention; the machine is able to take off, select and engage a target by itself (2015). Noone and Noone then go on to highlight the three types of weapons systems:

> "1. Human-in-the-loop or semi-autonomous systems require a human to direct the system to select a target and attack it, such as Predator or Reaper UAVs. 2. Human-on-the-loop or human-supervised autonomous systems are weapon systems that select targets and attack them, albeit with human operator oversight; examples

include Israel's Iron Dome and the U.S. Navy's Phalanx Close In Weapons System (or CIWS).

3. Human-out-of-the-loop or fully autonomous weapon systems can attack without any human interaction; there are currently no such weapons" (Noone, & Noone, 2015: 28).

For Boulanin (2019), automation could have both a positive and negative impact on the risk of nuclear war, with the negative impact being on a state's perception of the efficiency of their second-strike capability. The issue of second-strike capability was highlighted earlier as a factor of importance when it comes to nuclear deterrence and subsequently this research. Autonomous systems offer the potential capability of speed and reliability, which could tempt states to use them as early warning systems according to Boulanin (2019: 83). This could cause states to feel insecure; it must be reiterated that a state needs to only perceive the capability in order for it to feel insecure, the technology doesn't necessarily need to be efficient or functional. A state only needs to perceive that another country's AWS is capable of a disarming strike to cause strategic instability, as the eradication of a second-strike capability then equates to the dynamics of nuclear deterrence being unbalanced. This would then lead to more insecurity, showing that AWS would have an adverse effect on nuclear deterrence and subsequently strategic stability. Payne (2018: 7) states that there are some intriguing parallels between the development of 'deep learning' methods and the development of nuclear weapons during the Cold War, whereas Giest and Lohn (2018) state that AI may lead to possible nuclear escalation or the use of nuclear weapons. Allen and Chan (2017) highlight that AI will likely have an impact on military superiority and along with this, they use different cases of 'transformative technology' in order to understand AI. This highlights the importance of understanding nuclear technology as a case study to create the ability to gain insight on how to understand AWS.

What can be noted from this section on AI and AWS are a number of tenets that can be important for this research. Firstly, AI can be deemed as the science of making intelligent machines. This statement is twofold, owing to the fact that AI researchers are trying to recreate human intelligence and this highlights AI's connection to psychology. When it comes to AI, it is further divided into 'narrow' vs 'general' intelligence. Some authors may use different terminology; however, this research has chosen to go with this owing to the fact it is more commonly used. 'Narrow' intelligence involves a machine's ability to process data and find patterns within it; ultimately not equal to human intelligence. Meanwhile, 'general' intelligence is directly connected with being equal to, if not better than, human intelligence. Furthermore,

the most important aspect of AI to arise from the literature is termed 'deep learning'; which is the combination of ML and Neural Networks. Deep learning may equal better navigation and target recognition, which in theory can help find and target a country's second-strike capability. Two issues do arise from deep learning: the issue with it being a 'black box' and the need for large datasets. Furthermore, there is the issue of autonomy when it comes to AWS; the argument around autonomy is based around human intervention in a task that is being executed. Boulanin and Verbruggen's (2017: 5) framework for the issue of autonomy will be used; they define autonomy based on three sections: "(a) the human-machine command-and-control relationship; (b) the sophistication of the machine's decision-making process; and (c) the types of decisions or functions being made autonomous". These sections all look at the degree of the machine's ability to autonomously sense, decided and act with or without human intervention. This creates an understanding of how AWS will work, by understanding the fundamental technical aspects behind this technology.

## 2.3   New wars

This chapter now moves onto the new war thesis, which a stated already, is being used to create a secondary contextualization for the study and theoretical framework. New wars is a term that is highly contested by many scholars and the literature around the debate of new wars is vast. However, the author who originally coined the term 'new wars' was Kaldor. There are four quintessential tenets when it comes to new wars: actors, goals, methods and forms of finance. This section will start with a look at the main theorist of the new wars thesis, Kaldor. This section will review how Kaldor views new wars and ultimately how she defends them. It will then follow the literature and find out what other theorists argue about the new war thesis; be they for it or against it. It must be noted that most authors challenge new wars on the basis of 'old wars' or 'Clauswitzian wars'. This section will aim to review the new war thesis as it will become the structure for the theoretical framework. Kaldor (2013) states that she created the new war thesis in order to make a research strategy or a policy guide to understand the logic of contemporary warfare. Kaldor (2013) highlights that the 'new' in the new war thesis is a research strategy (2013). For Kaldor (2013), the new war thesis is enshrined in 'old' war logic and this gives the new war thesis the ability to understand contemporary conflict. Kaldor's (2013) new war thesis does not focus on a particular issue, but rather creates an 'integrative framework for analysis'. This 'integrative framework' will help to build a theoretical framework that will create a microscope that can be used to understand the data of this research. Kaldor (2013: 14) argues that the new war thesis aims to change the character of organised

violence and also builds a way of understanding the characteristic of violent conflicts. Kaldor (2013) ultimately sees the new war thesis as a form of analysis that allows for scholars to understand the dynamics of conflict.

### 2.3.1 Kaldor's 'New War Thesis'

What is a 'new war' according to Kaldor? Kaldor (2013: 2) argues that new wars are wars in the era of globalisation that take place in authoritarian states that have been weakened by its affects. In these contexts there is a mix of state and non-state actors, private and public, internal and external, war and peace, and political and economic that all become hard to distinguish from one another (Kaldor, 2013: 2). This breakdown of distinction is both the cause and consequence of such conflict (Kaldor, 2013: 2) Furthermore, Kaldor defines new wars based on actors, goals, methods and forms of finance (Kaldor, 2013). When it comes to the actors of new wars they differ from that of 'old war'; old wars were fought between regular armed forces of states (Kaldor, 2013: 2). According to Kaldor's 2013: 2) article, new wars are composed of state and non-state actors; non-state actors entail regular armed forces, warlords, paramilitaries, private security contractors, jihadists and mercenaries. When it comes to new wars their goals are different from those of old wars according to Kaldor (2013: 2), as old wars were fought for ideological and geopolitical interest, while new wars are fought in the name of 'identity politics'. It must be noted how Kaldor continuously describes new wars by showing how they are or are not different from old wars. Identity politics differ from geopolitics or ideologies based on their logic, as identity politics is aimed at a certain group rather than pursing policies or programmes that are in the general public interest (Kaldor, 2013: 2). Kaldor (2013: 2) attributes the rise of identity politics to new communication technologies; migration, both domestic and international; and the erosion of more inclusive political ideologies such as socialism. However, Kaldor (2013: 2) argues that the most important part of identity politics construction is war. This means that the aim of new wars is political mobilisation based on identity, rather than an instrument of war (Kaldor, 2013: 2). When it comes to the methods of war, Kaldor saw old wars entailing a 'decisive encounter' (Kaldor, 2013: 2-3). The method of war was capturing territory through military means; however, in new wars battles are rare and the territory is captured though political means (Kaldor, 2013: 2-3). The political means that are utilised to capture a territory are through the control of the population (Kaldor, 2013: 2-3). When it comes to how wars are financed, old wars were financed mainly by taxation, while new wars entail a very different form of financing, according to Kaldor (2013: 2-3). Finance in weak states is characterised by a falling tax revenue and different forms of predatory finance

such as looting, taxation of humanitarian aid, kidnapping, or smuggling of natural resources, and illicit trade like drugs (Kaldor, 2013: 2-3). For Kaldor old wars economies were 'centralizing, autarchic and mobilized the population', while new wars are a part of the 'open globalized decentralized economy' which has low participation and the revenue is dependent on the violence (Kaldor, 2013: 2-3). When it comes to war, according to Kaldor, old wars were about state building, while new wars are about the dismantling of the state (Kaldor, 2013: 3)

### 2.3.2    Are New Wars 'New'?

Is the new war thesis true, are new wars truly 'new'? One of the critics of Kaldor's (2001: 163) new war thesis is Booth, who takes issue with the term 'new'. This statement is based on the premise that wars have not changed significantly enough to be termed as 'new' (Booth, 2001: 164). Booth states that the temptation to oversimplify war must be resisted and that wars such as colonial, guerrilla, rebellions were side-lined by Kaldor (Booth, 2001: 164). What Booth (2001) is trying to maintain here is that warfare between a state and non-state actors has been a feature of warfare throughout history; it is not specific to 'new wars'. What Booth (2001) is maintaining is that there is significant historical evidence that what is entailed in 'new wars' has happened before. As argued by Kaldor (2013), new wars have elements of old wars, the new war thesis is ultimately a research strategy. Shaw (2000: 173) argues that some will challenge the historical issues of the new war thesis; however, these characteristics may have been seen before, but the combination of them in new wars is distinctive. For instance, the genocide in Nazi Germany compared to the Genocide in Rwanda contain very different modes (Shaw, 2000). Shaw (2000) states that what happened in Rwanda was 'amateurish' compared to the Nazis. This statement is built on the premise that new wars are built on targeting the civilian population as a mode of warfare, they are targeted due to uncertainty (Shaw, 2000). This means, that for Shaw (2000: 179), new wars are genocidal. Ultimately for Shaw, Kaldor introduces globalisation as the new form of a war economy, but for Shaw it is presented as external to war. Shaw (2000: 179) states that globality can be traced to the contradiction of war, along with this, globalisation can not only be linked to the end of the Cold War, but also the weakening of the relationship between the war machines and the economy. Ultimately for Shaw (2000), Kaldor illustrates a new form of war and helps question the contemporary mode of warfare.

Like Shaw (2000), Newman (2004: 179) states that new wars' key characteristics that are highlighted by Kaldor are not new or have not sufficiently changed enough to be deemed as

new. One of Newman's (2004) most important critiques of new wars is that it lacks a sufficient historical context; however, Newman states that it has helped scholars to understand civil war better. Newman (2004) states that Kaldor is not incorrect in her analysis of new wars, but that she negates the issue of the goals, methods and how they are financed which can be historically found in other wars. What Newman (2004) is arguing for is a more accurate and in-depth analysis of the historical context of war, especially civil wars. Even if there is incomplete data, analysts looking at the new war thesis must build a database of all civil wars in order to create a more accurate analysis of new wars. As aspects such as the targeting of civilian populations and forced human displacement have been a tactic of certain wars, however, the fact that it is more a tactic of war over being a consequence of war is appealing for Newman (2004: 181-182). Newman (2004: 183) highlight that the new war thesis describes the social and economic context of war based on a failed or weak state, with a collapsed formal economy and rivalry between different criminal groups over illicit activities or resources. The issue of a failed or weak state comes about owing to the erosions of the state by globalisation and neoliberal economic policies (Newman, 2004: 175). Newman (2004: 186) finds the economic and social dynamics explained by the new war thesis as important to understanding contemporary conflict owing to the fact that it points analysts in the direction of understanding human insecurity and violent conflict. Newman highlights how (2004) 'war economies' are characterised by violence over resources, the black market or external assistance; this comes about owing to the issue of globalisation which Newman highlights extensively. This is what is termed as a 'globalised war economy', driven by the desire to gain wealth and power, meaning the violence is perpetuated in order to continue the cycle of gaining wealth and power (Newman, 2004). The groups fighting for power are defined by some form of identity such as religion and ethnicity (Newman, 2004: 174).

In summary, Newman (2004: 185) agrees that wars have changed owing to technological innovation and socio-economic changes, but the historical perspective is still an important factor in the analysis of new wars. However, Kaldor (2013) emphasises that the 'old' is a part of the 'new' and that the new war thesis is a research strategy. However, the historical perspective is of importance to this research, which is why nuclear technology has been chosen as a case study. The benefits of such a strategy are twofold, it allows for an understanding of a transformative technology and how this technology affected the mode of warfare. Furthermore, it can help understand the 'goal' of this warfare, to which parallels can be drawn with the

current 'goals' of the US. As democracy drove the goals of the Cold War, will the rise of populism have the same effect on driving US strategy now?

### 2.3.3 Globalisation: A further investigation into the economics of warfare

Berdal (2003: 478) has a more in-depth focus on the economic agenda of the actors involved. He states there is a need to widen the debate owing to the fact that: economic agendas of war are affected by the change of the economic environment, and focusing on one set of factors, his literature can add to the discussion of new wars. For Berdal (2003: 479), the lack of precision when it comes to applying globalisation is problematic, as it is more than just a 'metaphor'. He wants a more precise definition as a lack of a precise definition is what a lot of the new war literature suffers from according to Berdal (2003: 480-481). Berdal (2003: 482) attributed globalisation to the easy access to capital and finances, better communications and transportation, some partial deregulation of industry, and transnational processes of exchange and production, which are the agreed-upon drivers of globalisation. Berdal (2003: 482) focuses on how a deregulated and open international economy has allowed for belligerents to develop economic interests and sustain conflict.

> "By posing the question of what functional utility violence may be serving to participants in wars to elites, ordinary people caught up in war, and external actors that stand to gain from the conflict it becomes possible to discern how a set of vested interests in the continuation of war may emerge" (Berdal, 2003: 482).

This will eventually equal a war economy, turning into the form of a regional pattern of informal economic activity (Berdal, 2003: 483-484). This self-interest of local elites, external actors pursing profit and a vulnerable population concerned with survival help further explain globalisations effect on new wars (Berdal, 2003: 490).

Another author that emphasises the importance of economics of war is Munkler (2005: 1), who argues that to understand new wars there needs to be an understanding of its economic foundations. Once again, Munkler (2005) is yet another theorist who has a critique of the new war thesis. Munkler (2005) states, like Berdal (2003), that new wars are not inherently new. Alike the work of Berdal (2003) and Newman (2004), Munkler (2005: 2) highlights the importance of the historical perspective. He reiterates the need for 'new wars' to be compared to 'old wars'. Most of the literature reviewed compared 'new war' to 'old war', another term for 'old war' is 'Clauswitzian war'. According to Kaldor (2005), old wars were idealised

between the eighteenth to mid-twentieth century; they were built on the premise of two states fighting each other in a decisive manner, involving actions such as statecraft and state-building.

Munkler (2005: 8) would, however, disagree with the statement that new wars are state-disintegrating and not building, he states that they may contribute to state-building. Munkler (2005: 10) argues that states not being able to control the development of their national economy is due them being linked to the world market system. This statement holds true, as the era of globalisation has changed the way wars are financed. They are no longer based on taxation like that of old wars, they are based on other forms of finance, like looting, external assistance and selling of resource rights.

When it comes to how new wars are financed, Kaldor (2013: 3) compares how old wars were financed against how new wars are financed. Old wars were financed by taxation or by private patrons, while new wars are financed by looting, taxation of aid, diaspora support, kidnapping, exploitation of natural resources or selling of their rights. Berdal (2003: 485) similarly highlights how natural resources or outside assistance help consolidate power and the ability to keep these wars going when he states: "Angola, Congo, Liberia, and Sierra Leone all illustrated how, in the context of acute state weakness or state collapse, war and violence will give rise to economic opportunities for a range of actors, both at the level of elites and among populations adjusting to the dislocation and stresses of war". Berdal (2003) views these warlords as businessmen, who have been given this ability due to natural resources and the interconnectedness that came about due to globalisation.

What Berdal is stating is that these local warlords require global connections in order to sustain their wars.

> "Globalization has opened up new opportunities for individual nonstate actors within weak states to link to global trading networks and potential partners without state interference. Improved communication technology, fast capital movements and increased deregulation in Western economies have created the necessary preconditions for coalitions between local warlords, private business, intermediary agents and emerging private security companies to capitalize upon the lack of states control on resources extraction" (Berdal, 2013: 489).

This shows how 'new wars' are inherently affected by globalisation. This will be of importance to understand as, Allen and Chan (2017: 43) state, nuclear technology is a 'high' cost profile

technology compared to cyber which is seen as having a 'low' cost profile. When it comes to artificial cost profile, Allen and Chan (2017: 46) define the cost profile as 'diverse, but potentially low'. Developing AI and ML capabilities will cost firms millions or billions of dollars (Allen, & Chan, 2017: 46). However, the cost could also be relatively low owing to open-source code libraries and COTS or rented hardware that will allow for the development of AI for less than $1 million (Allen, & Chan, 2017: 46). There is also the possibility that it may be free if copies are leaked (Allen, & Chan, 2017: 46). As stated by Munkler (2005: 98), wars have become cheaper and the ability to equip armies has become easy owing to open war economies. They rely on light weapons and trucks over jeeps. Larger weapons which are used by these armies are usually leftovers (Munkler, 2005: 74). Could AWS possibly become the new light weapon or truck of new wars? This section has helped to build an understanding of how wars are financed. Furthermore, it has highlighted the importance of the effect of globalisation on new wars.

### 2.3.4    Kaldor's 'New War Thesis': A research strategy

This literature review bases itself on Kaldor's (2013) view that the 'new war' thesis is ultimately a research strategy. What are the implications of this for Autonomous Weapons Systems (AWS) and strategic stability?  It is important to understand the economics of warfare, or how wars are 'financed'. This is what Allen and Chan (2017) refer to as a 'cost profile'. The goals, mode and historical perspective, the main tenets gained from the literature above are of significant importance; however, the economic aspects are of considerable importance when it comes to transformative technology, especially AWS. This section now starts to move from the general debate around new wars, with each different theorist discussed so far, and starts to group them in Kaldor's main tenets.

### 2.3.5    The goals of new wars

In terms of the 'goals' of new wars, theorists such as Booth believe that what new war thesis claims can be seen in earlier wars, meaning that identity politics in not necessarily new (2001: 163). However, Kaldor (2013: 3) argues that with globalisation came the rise in identity politics. For Kaldor (2013: 2), identity politics is constructed through war, meaning that political mobilization around identity politics is the aim of war, rather than being an instrument of war which was seen in 'old wars'. Shaw (2000: 172) argues that new wars are about political issues rather than military issues, they are about the breakdown of state legitimacy. "The goals

of new wars are about identity politics in contrast to geo-political or ideological goals of earlier wars." (Kaldor: 2013: 6) For clarification Kaldor (2013: 6) defines 'identity politics' as the claim to power based on a certain identity, such as nationality, clan, religion or linguistics. Identity politics can now either be local or international owing to globalisation, which saw the increase of interconnectedness due to technological innovations (Kaldor: 2013).

### 2.3.6    The methods of new wars

Another important aspect of the literature that emerged was what Kaldor (2013: 2-3) termed the 'methods' of new wars. For Kaldor (2013), the method of old wars was a decisive battle between two or more nation-states with armies, the aim was to capture the territory of another sovereign state. Shaw's (2000) research analyses of Kaldor's concept of this type of warfare, referred to so far in this literature review as 'old war' or 'Clauswitzian war'. Shaw (2000: 173) states that some many argue that new wars are not so new; however, the combination of them in new wars is highly distinctive. Kaldor (2013: 10) argues that the Clauswitzian 'trinitarian concept of war', which is the state, the army and the people, is no longer relevant to modern warfare. Kaldor (2013) states further that the difference between new and old wars is a contrast between ideal types rather than actual historical types. Kaldor does concede that the wars of the twentieth century are closer to old wars; however, wars of the twenty-first century are closer to Kaldor's (2013: 13-14) depiction.

Although Flemming (2009) found that the premise that new wars are post-Clauswitzian wars to be unfounded, this does not then mean that the new war thesis is irrelevant. What Flemming (2009) is arguing for is based in the logic of Clauswitz, that theory should not become a doctrine, it must instead become an area of study, meaning that inquiry is the most important aspect of theory. Flemming (2009) states that two different theories too often become polarised rivals, whereas they can subsequently be used together to understand war.

> "Clauswitzian concepts can be used as analytical tools in ostensibly new wars, just as the 'new war' trends can open up the complexity of war and the requirements to find a political solution to contemporary humanitarian and conflict solutions" (Flemming, 2009: 238).

Mello (2010) furthermore argues that the work of Clauswitz can be useful when analysing modern warfare, as it allows for the study and comparison of all warfare. While this may be true, Kaldor's (2013: 4) aim with the new war thesis is to look at the changing nature or character of organized violence and to develop a way to understand, interpret and explain the

interconnected characteristics of contemporary warfare. Kaldor (2013) concedes that while historically there are wars that are similar, this is a form of analysis and a research strategy. Newman (2004: 180) also highlights the importance of Kaldor's new war thesis as a form of analysis, as it has allowed analysts to focus on these issues and subsequently understand conflict dynamics. However, Kaldor (2013) does not see the need for the historical perspective, while Newman (2004) sees the historical perspective as an important factor of the new war thesis. The issue of historical perspective is of importance to understand transformative technology in the twenty-first century, which is why this research will see the historical aspect as an important tenet of the new war thesis.

### 2.3.7 The actors: Who fights these wars?

The final tenet of the new war thesis to be reviewed is the 'actors' that are involved in new wars. Kaldor (2013: 2) states that older wars were fought by states with armed forces while new wars are fought by a varying combination of state and non-state actors. Examples outlined by Kaldor (2013: 2) are: "regular armed forces, private security contractors, mercenaries, jihadists, warlords, paramilitaries, etc." This comes about due to the state's monopoly of violence being eroded by it becoming privatised. The actors involved in new wars are of increasing importance when it comes to AWS, owing to the fact that nuclear deterrence is based on perception and an adversary. This research will not focus on non-state actors as it will take a state-centric view. This research will use the US's NSS in order to find out who the US views as a potential aggressor and what this means for their policy. Depending on how the US views other actors will influence how they pursue AWS.

There are four main tenets of the new war thesis that are highlighted by Kaldor: they are the actors, goals, methods and forms of finance. The main tenet that needs to be re-emphasised here is the economics of warfare, how they are financed. This is a focal point owing to what is called the 'cost profile' of AWS. It is undoubtedly clear that AWS will change the mode of warfare. When it comes to the goals, the world is experiencing a rise in identity politics, making it a significant factor when understanding the goals of actors. Kaldor also highlights identity politics as an important aspect of the new war thesis. Can identity politics be applied to the US and help one to understand the US goals? The next section will aim to further look at identity politics and lead into populism, which will provide a basis to build on to understand the goals of warfare.

### 2.3.8    Identity politics

The role of identity politics, and more specifically populism, is to create a better contextualization for this research. The rational for this contextualization is to create an understanding of the circumstances in which this study is taking place. This backdrop will allow this research to better understand how President Trump and populism affect the main research question. This section will look at how Kaldor views identity politics.  Kaldor (2013: 2) states clearly that old wars were fought through the means of the interest of a geopolitical location or ideology; an example would be democracy or socialism. When it comes to new wars, according to Kaldor (2013), they are fought through the means of 'identity politics. These statements are not new to this review of literature, as they have been extensively discussed above. However, this section aims to further review the literature of identity politics. Identity politics is arguably what led to the election of President Donald Trump and the United Kingdom's decision to leave the European Union. These are two well-cited examples of identity politics and what is known was 'populism'. In the literature, populism and identity politics seem to be interlinked, however, identity politics seems easier to define than populism. Therefore, this section will commence with a review of literature on identity politics. It will then move on to the review of populism, and finally how the two are interconnected. The reason this concept is being further reviewed is owing to the fact that this research is forecasting that identity politics, specifically populism, will have an impact on US national strategic policy, arguably in a way that democracy or communism influenced countries' policies during the Cold War. This will help to further develop the 'goals' section of the theoretical framework, as this research will have a more in-depth understanding of what identity politics is. Furthermore, this will create a background and secondary contextualization for the study.

### 2.3.9    *What is identity politics?*

To further this secondary contextualization, this section will begin with a brief overview of how Kaldor views identity politics, as Kaldor is the major theorist behind the new war thesis and subsequently of significance importance to this study. In Kaldor's (2013: 337) article 'Identity and War' she uses the work of Sen when talking about the construction of identity. Kaldor (2013: 337) highlights the literature of Sen (2006) who states that an individual can have multiple identities, it is when one of these identities becomes overarching, it is inevitable and natural that conflict is likely. Kaldor (2013: 338) states that identity should be seen as a process of identification, 'an ongoing process of inventing and reinventing ourselves'. When

there is not a process of identification and a singular identity becomes prevalent it is known as a unidimensional identity, which is the cause of violence and conflict (Sen, 2006 as cited in Kaldor, 2013: 338). A unidimensional identity is uniquely defined in relation to other identities that are also unidimensional (Kaldor, 2013: 338). Kaldor (2013: 338) states that unidimensional identities are built out of binaries and sometimes a triangular distinction. For example, a Jewish person is defined in relation to an anti-Semite (Kaldor, 2013: 338). When it comes to the binaries that identity is built on, Kaldor refers to the work of Schmitt, who states that 'the political' is defined by the friend-enemy distinction (Schmitt, 1927 as cited in Kaldor, 2013: 338). Kaldor (2013: 339) states that violence is a form of communication between the victim and the perpetrator. Whatever the motivation is, fear and hate are what is communicated (Kaldor, 2013: 339). Furthermore, violence allows for a narrow binary distinction between the victim and perpetrator what is called a 'friend' and 'enemy binary' by Kaldor (2013: 339). Kaldor (2013: 339) states that violence divides the binaries even more. Kaldor's (2013) argument is that war constructs unidimensional identities which are the opposite to an individual's multidimensional identity. Kaldor (2013: 2) also highlights how identity politics is most importantly constructed through war.

Schafer (2005: 93) sees ethnicity as a process, which he terms ethnicization, a tool that can be used for political power. A collective identity is a powerful tool for actors with or without economic and/ or political control according to Schafer (2005: 94). Schafer (2005) also highlights that identity politics is a process that can be used by the political elite or actors who have an interest. However, unlike Kaldor (2013), who sees individuals as having multiple identities, Schafer (2005) focuses on the issue of ethnicity and religion. For Schafer (2005: 94), ethnic and religious identities represent a way of truth, dignity and life for individuals. However, when it comes to identity politics the implications for conflict can be problematic. Schafer (2005: 94) states that identity-based conflicts lack clarification that led to difficult negotiations, this is owing to the fact that they lack clear goals and reasons; Schafer states that they lean on the side of mystification. Further, these conflict values go beyond the individual, meaning that human life is not of high importance (Schafer, 2005: 94). Like the work of Kaldor (2013), the more the violence and cohesion continue the more the hatred increases, on top of this Schafer (2005) states the more the violence continues, the more the guilt increases. Schafer (2005: 96) goes on to further talk about the role of 'symbolic violence. "In these clashes of mobilized identities, symbolic violence gains an important role. First of all, it means not recognizing the identity of the others as legitimate. This, it implies negative ascription

(labelling)" (Schafer, 2005: 96). Schafer (2005: 96) states that actors opt for racism, which is a strategy of fundamentalism.

For Fukuyama (2018), politics today is less concerned with economic or ideological concerns, but with the issue of identity. "Now, in many democracies, the left focuses less on creating broad economic equality and more on promoting the interests of a wide variety of marginalized groups, such as ethnic minorities, immigrants and refugees, women, and LGBT people. The right, meanwhile, has redefined its core mission as the patriotic protection of traditional national identity, which is often explicitly connected to race, ethnicity, or religion" (Fukuyama, 2018: 91). Fukuyama (2018) highlights how these leaders utilise the issue of dignity, which was also seen in the work of Schafer (2005) and Kaldor (2013) when it comes to identity politics. For Fukuyama (2018: 92), identity politics has moved away from university campuses and 'cultural wars', and towards a concept that can help explain what is going on in global affairs. This comes about, according to Fukuyama (2018: 92), owing to the fact the groups believe their identities are not receiving adequate recognition. Fukuyama (2018: 92-93) states that globalisation has led to rapid economic growth and social change, this had made societies more diverse, making groups that were previously invisible to the mainstream demand recognition; this has subsequently led to a backlash as certain groups feel a loss of status and a sense of displacement. Fukuyama (2018: 93) argues that democratic societies are fracturing into ever-narrowing identities which threatens the ability of a society to deliberate and act as a whole. Fukuyama (2018: 93) states that this leads to the breakdown of the state and to its ultimate failure. The solution to this mandate according to Fukuyama (2018: 93) involves liberal democracies moving back to a more universal understanding of human dignity, otherwise they are doomed to continue conflict along with the rest of the world. This shows how Fukuyama (2018) not only views identity politics, but also how dangerous they can be. This also presents the importance of continuing to look at identity politics in its own separate section. Kaldor's identity politics is based in civil wars, Fukuyama (2018) helps to move identity politics away from it and into global affairs.

Besley and Persson (2019) state that there are two factors that drive identity politics.

> "Increased economic polarization has been a short-run consequence of the recent financial crisis and a long-term consequence of manufacturing jobs being swept away by globalization. And heightened concerns about the loss of social status have gradually

emerged among traditionally dominant groups, who feel threatened by immigration and gender equality." (Besley, & Persson, 2019: 2)

Besley and Persson (2019: 2) attribute the election of President Trump and the Brexit vote as a result of radical-right populist parties and that this trend is seen across the world. Besley and Persson (2019: 2) have developed a model that allows for a two-way feedback between political culture and policy; their approach sees social identification as having that ability to influence policy. This model brings this literature review to an important point. Identity politics is not just a phenomenon seen in culture wars or civil wars. It is of importance to understand identity politics and its ability to affect national strategic security policy. This will in turn give this research the capability to predict how a state will pursue AWS; based on the premise that ideologies affected the Cold War and identity politics will affect new wars.

### 2.3.10 Populism and identity politics

This section reviews the literature of identity politics and furthermore the basis of knowledge of identity politics. This part now aims to take a step towards a more specific example of identity politics, populism. By understand populism and President Trump this research will then have a better background and contextualization of the entire process of how AWS will affect the US's nuclear deterrence. Identity politics is an important tenet of Kaldor's new war thesis and has an effect on the goals of warfare and how war is conducted. This research is an attempt to understand how a prominent movement like populism can affect a state's national strategic security policy. Slogans like 'Make America Great Again' convey the idea that America under the Trump presidency may pursue AI in order to maintain their military dominance. Now that a general outline of different theorist on identity politics has been outlined, this section will further deal with the issue of identity politics and discuss what the correlation is between identity politics and populism. Due to this research being based on the US as a research site there is a need to understand the identity politics that are dominating its domestic politics. One of the most prominent identity politics that is currently dominating US politics is populism; known as the alt-right. It is important to understand this phenomenon in order to see its possible effect on US policy which can then impact how the US will pursue AWS. This will remain a part of Kaldor's central tenets of the new war thesis. It is just worthwhile getting a further understanding of what this complex theory is.

Marchlewska, Cichoka, Panayiotou, Castellanos, and Batayneh highlight the work of Müller, that populism is a form of identity politics (Müller, 2016 as cited by Marchlewska *et al,* 2017:

151). Marchlewska *et al* highlight how Müller states that populists combine 'anti-elitism' and what a true vision of a citizen of a nation looks like (Müller, 2016 as cited by Marchlewska *et al*, 2017: 151). Subsequently, the opposite of the populist agenda is then termed as a threat to national interest (Marchlewska *et al*, 2017: 151). Arguably, here is what Kaldor (2013) would call the 'friend-enemy' binary that is involved in identity politics of the new war thesis. Fukuyama's (2017) view of populism goes on to further state how 'populist nationalism' is a threat to liberal international order. This shows that populism is not just an issue within a sovereign territory, it can also be a threat to the liberal international order. Fukuyama (2017) states that populism is a term that is loosely used in order to describe a range of phenomena around the world. This is unlike the words of Marchlewska *et al* (2017) who believe that populism is based on a 'national collective narcissism', which they define as an unrealistic belief of the greatness of a national group. Fukuyama (2017: 1) states that there are three characteristics that can be associated with populism: "[t]he first is a regime that pursues policies that are popular in the short run but unsustainable in the long run, usually in the realm of social policies. Examples would be price subsidies, generous pension benefits, or free medical clinics", "[a] second has to do with the definition of the "people" that are the basis for legitimacy: many populist regimes do not include the whole population, but rather a certain ethnic or racial group that are said to be the 'true' people", and "[a] third definition of populism has to do with the style of leadership. Populist leaders tend to develop a cult of personality around themselves, claiming the mantle of charismatic authority that exists independently of institutions like political parties". These leaders try to make a connection with 'the people' and along with this they tend to denounce the elite (Fukuyama, 2017: 2). According to Fukuyama (2017: 2), this style of leadership is an issue for modern liberal democracies that are built around power-sharing institutions, such as courts, federalism, legislatures, and a free media. These are roadblocks for populist leaders, which subsequently makes them at threat from populist leaders. There are three definitions of different populist leaders: the issue of economics, who the people are, and the cult of the personality. Fukuyama (2017: 2) states that Donald Trump fits into all three of these categories.

Fukuyama (2017) further categorises populist movements into two categories by geographical area. Populist movements in Latin America and Southern Europe tend to be Left and focused among the poor and advocating social programmes that aim to rectify social inequality (Fukuyama, 2017: 2). They do not emphasise the issue of ethnicity or take a stance against immigration (Fukuyama, 2017: 2). Meanwhile in northern Europe populist movements focus

less on the poor and more on the declining middle or working class and are more right-wing and anti-immigrant (Fukuyama, 2017: 2). They want to protect the existing welfare state; however, they do not emphasise a rapid expansion of social services or subsidies (Fukuyama, 2017: 2). Fukuyama (2017) then goes on to state that there are also populist movements that exist, yet do not fit into either of these categories.

For example, Wasko-Owsiejcsuk (2018: 83-84) states that Trump's emerging doctrine and foreign policy have a 'large dose of populism'. This means that Trump's foreign policy is dominated by the discourse of 'Making America Great Again' or 'Putting America First'. This has subsequently affected US national security strategy in multiple ways, one of the main ones being economic strength which leads to military strength which subsequently leads to possible military dominance (Wasko-Owsiejcsuk, 2018: 91). This subsequently shows the effect of populism on National Security policy and the importance of understand such a theory. It is important to understanding the forces behind a country's policy as this can help predict how a policy is pursued or if it will be pursued.

What can be taken away from this section is the causality between identity politics and conflict, specifically, when a singular identity becomes overarching and how this equates to inevitable conflict. A singular identity is problematic owing to the fact that humans have multiple identities. However, Schafer (2005) states that the quintessential characteristics of identity are religion and ethnicity. For Schafer (2005), religion and ethnicity mean truth, dignity and life to an individual. Schafer (2005) states a valuable point of debate, as the issue of ethnicity is central to Fukuyama's (2017) populism definition. Furthermore, for Kaldor (2013), war constructs these identities. Conflict based on identity also becomes more problematic to negotiate an end. Furthermore, these conflicts can be termed as symbolic violence as the aggressors do not see the other identity as legitimate. But these conflicts based on identity are not just found in 'culture wars', according to Fukuyama (2017). Identity-based conflicts have moved into global affairs (2017). Fukuyama (2017) states that democratic societies are fracturing into separate identities based on economic polarisation and social identification. This brings populism to the fore, populism is a form of identity. It is based on what a true citizen looks like and 'anti-elitism'. For Fukuyama (2017) it is a direct threat to international order. Finally, and most importantly, Fukuyama (2017) defines populism based on: 1. A regime pursuing certain policies popular in the short term and not sustainable in the long term, 2. It is based on 'what are the people?' which is based on an ethnic or racial group, and 3. There is a cult of the personality (Fukuyama, 2017). It must be finally stated, that from the literature came the debate

that social identification has the ability to affect policy, showing the importance of understanding the 'goals' of the new war thesis and more importantly it has helped to understand what populism is.

## 2.4    Theoretical framework

This section creates the theoretical framework in order to make sense of and answer the research questions. This means that all the relevant information will be extracted, and a microscope created. The framework will be based on Kaldor's (2013) new war thesis.  As already stated, four main tenets of the new war thesis that have been outlined by Kaldor (2013): actors, goals, methods and forms of finance. This framework will be based predominantly on Kaldor (2013); however, it will be adjusted accordingly based on the literature review of the new war thesis. Furthermore, the rest of the literature reviewed in this section will be added in each of the relevant sections in order to further build the theoretical framework.

| | New War |
|---|---|
| **Actors** | • This section combined with goals will be used to create a secondary contextualization for the study.<br>• Range of actors; such as regular armed forces, non-state actors, warlords and private security contractors. |
| **Goals** | • Are fought on the basis of identity politics, whilst old wars were fought on the basis of ideologies such as democracy or communism (Kaldor, 2013).<br>• Individuals have multiple identities; however, when one identity becomes overarching, conflict will ensue (Kaldor, 2013).<br>• Identity is a process of inventing and reinventing oneself.<br>• A unidimensional identity is the cause of violence and conflict (Kaldor, 2013)<br>• Identities are formed in relation to another (Kaldor, 2013).<br>• Unidimensional identity built on binaries or sometimes a triangular distinction (Kaldor, 2013). |

| Methods | <ul><li>This section will will help to understand the strategies employed.</li><li>Old wars were fought through the means of a decisive encounter and the method involved capturing territory through military means (Kaldor, 2013).</li><li>New wars involve targeting of civilians and the capturing of territory through political means (Kaldor, 2013).</li><li>The method of warfare chosen is nuclear deterrence.</li><li>Nuclear deterrence has three important tenets: military capability; unacceptable costs can be imposed; and credibility that this cost will be carried out on the aggressor (Morgan, 2003: 74).</li><li>Finally, AWS will effect a state's mode of warfare and more specifically the nuclear capability.</li></ul> |
|---|---|
| Forms of Finance | <ul><li>The section will analyse how AI and AWS are financed and who exactly is financing them.</li><li>New wars are financed through what can be deemed as criminal activity and it is arguable that these wars are continued based on economic gain (Kaldor, 2013)</li><li>Old wars were largely financed by states through taxation or outside patrons (Kaldor, 2013).</li><li>The cost profile of AI is 'diverse, but potentially low' (Allen & Chan, 2017: 46)</li><li>Finally, can AWS lead to a form of military dominance that allows a state to increase its economic dominance by increasing its military capability.</li></ul> |

Table 1: Theoretical Framework compiled by the author from Kaldor (2013)

## 2.5   Conclusion

The aim of this chapter was to review the literature of the main theories and create a theoretical framework to understand the chosen case studies. The chapter began with the review of literature on strategic stability in order to gain an understanding for the next section, nuclear deterrence. Nuclear deterrence is a strategy of strategic stability but also an important tenet of this research, which aims to look at the causality between AWS and nuclear deterrence. More specifically than nuclear deterrence, the research will look at the effect of AWS on a state's second-strike capability. The section that followed this was the review of literature on AI, which was done in order to build a basis for AWS. The reason for reviewing this literature was to build a basis of knowledge of what exactly AWS is and what the major theoretical issues around it are, for example: the issue of autonomy. This literature review then moved forward to look at the 'new war thesis' with an emphasis on Mary Kaldor. This research sees the new

war thesis as a research strategy. It also contains valuable information, making it more than just a research strategy. This was of importance owing to the fact that the new war thesis was used in order to build the theoretical framework for this study. From the new war thesis arose the issue of identity politics, a concept dominating the political arena worldwide. It also allowed this research to build a secondary contextualization in order to better situate the reader in the circumstances in which this study is taking place. As AWS affecting nuclear deterrence is the end process of politics, there is an entire phenomenon before that needs to be clarified, that phenomenon being populism and President Trump. This led to a further investigation of the 'goals' of the new war thesis; which is identity politics. This drove the research towards populism, a form of identity politics currently dominating the US. This research aims to create a contextualization of the causality between populism and how states will pursue their national strategic security policy. The final section of this literature review was the development of the theoretical framework, which involved the combination of all the literature to create a theoretical lens in which the case study can be understood. The chapter to follow will be the contextualisation of the different transformative technologies and how they impacted how states went or go to war.

# Chapter 3

## 3.1 Introduction: Transformative Technology and Warfare

The aim of this chapter is to create contextualisation and background for this study. The basis for this contextualisation is the causality between the creation of a 'transformative technology' and its subsequent effect on how warfare is conducted. This chapter will be based on Allen and Chan's (2017) literature that focused on what lessons can be learnt from previous 'transformative technologies'. The four key transformative technologies that Allen and Chan (2017) looked at were biowarfare, cyberwarfare, aerospace and nuclear warfare. This section aims to study different transformative technologies and their effect on warfare. The transformative technologies this research has decided to cover are nuclear, cyber, biowarfare, and AI. This study chooses to omit aerospace as it cannot cover all aspects of warfare; furthermore, aerospace is a very broad and literature-rich sector. However, air warfare will be briefly touched on in the nuclear warfare section and AI section. This is owing to the fact that heavy bombers play a vital role in the US's 'nuclear triad' strategy. Also, the advancements in AI are playing a vital role in the aerospace domain.

The approach that has been chosen in this section is because this study researches how AWS will affect nuclear deterrence. What this means is that this study wants to understand how states go to war. By looking at important transformative technologies, this thesis will create an understanding of how these innovations change the way states conduct warfare. Furthermore, this section will look into nuclear warfare and create a better understanding of how the US conducts its nuclear deterrence. This means that the nuclear technology section has two main aims: Firstly, to understand how transformative technology impacts how states go to war, and secondly to create a more in-depth understanding of how the US conducts its nuclear national strategic security policy.

Each section will begin with an origin of the technology and general background, which will entail aspects such as its capabilities, greatest innovations and problems encountered. From the literature of each of these technologies their prominent examples were mentioned. This will help to understand the importance of understanding the role of technology in warfare and how it changes the way warfare is conducted. This section will also look at how these technologies all have dual uses, meaning that they can be used in the private sector as well as for military purposes. Dual use is problematic for the military as they may have to compete for resources

or have this technology accessed by non-state actors or potential adversaries. Furthermore, it will highlight how the US strives for technological superiority when it comes to warfare.

## 3.2    The origin of biowarfare

This section will first look at bioweapons and will subsequently only discuss bioweapons and not chemical weapons. Chemical weapons have similarities with bioweapons, however as stated by Spiers (2010), bioweapons have the potential to be more dangerous than chemical weapons. Many theorists, such as Allen and Chan (2017), Van Aken and Hammond (2003), and Zanders, Hart and Kuhlau (2001), all state how biowarfare is not a new strategy as to how states conduct warfare. Zanders *et al* (2001) state that biowarfare is almost as old as human civilization, they state that the Athenians were pushed into a marsh that contained a 'virulent epidemic' during the summer which halted their siege of Sicily. Van Aken and Hammond (2003) highlighted that the Mongolians would catapult plague victims into besieged cities, which could have possibly caused the first epidemic in Europe. Allen and Chan (2017: 100) also highlight how the Mongol army at the siege of Caffe in 1346 catapulted plague-infected corpses over the walls of the besieged city. Furthermore, Allen and Chan (2017: 100) state that the use of infectious diseases can be seen as early as 600 BCE. Another prominent example of the use of infectious diseases was in North America during the eighteenth century (Allen and Chan, 2017; Zander *et al,* 2001). The British used smallpox-infected blankets to cause an epidemic in an enemy native American tribe (Allen & Chan; 2017: 100). Zanders *et al* (2001: 2) state that there were several examples of the deliberate use of smallpox during the wars in North America. This included the French, the British and the Indians (Zanders *et* al, 2001: 2). During the American wars, smallpox was a valuable tactic for the British as their soldiers were inoculated and the American soldiers were not, giving them a military advantage (Zanders *et al*, 2001). This advantage was then taken away when Washington mandated that all soldiers in the American military be vaccinated (Zanders *et al*, 2001). This section, so far, briefly shows different prominent examples of how biowarfare was used by different militaries. This chapter will now move on to biological weapons and biowarfare by looking at WWI, the Spanish Flu and WWII.

### 3.2.1    Biological weapons and biowarfare: WWI, the Spanish Flu and WWII

This section discusses different examples of biological weapons and biowarfare during the twentieth century. The CRS define biological weapons as: "…a biological agent that is intentionally used to harm or kill humans, animals, or plants." (Biological weapons: a primer,

2001: 1) The type of biological agents is typically bacteria, viruses, fungi or rickettsia (Biological weapons: a primer, 2001: 1). Biological weapons come from biotechnology, Zander *et al* (2001: 7) define biotechnology broadly as: "any technique that uses living organisms (or parts of organisms) to make or modify products, improve plants or animals, or develop microorganisms for specific uses." This highlights the issue that becomes problematic for biotech, its dual use, which is why bioweapons are defined by their intention when it comes to using biotech, not what it is (Biological weapons: a primer, 2001). Now that biowarfare has been defined, the section turns to the role biotechnology has played in warfare. Once again, the premise of this section is the effect of technology on how states conducted warfare. Furthermore, these sections will highlight how the US always strives for technological superiority.

By 1914, according to Zander *et al* (2001), microbiology had made incremental advancements. According to Zander *et al* (2001: 2), there had been isolation and cultivation of major bacterial diseases. Furthermore, there was the existence of viral diseases, but they were not well understood (Zanders *et al*, 2001: 2) There was also studies conducted on parasitic diseases (Zanders *et al*, 2001: 2). Zanders *et al* (2001: 2) state that there was a understanding of how diseases were transmitted which contributed to better countermeasures, prevention and prophylaxis. Zander *et al* (2001) states that these insights were used by hostile purpose in WWI, they were, however used for sabotage purposes and subsequently not directed at humans. Allen and Chan (2017) state that bioweapons did not play a role during WWI, however, the Spanish Influenza did.

The Spanish Influenza infected one-third of the world's population and killed an estimated fifty million people (Allen and Chan, 2017). The Spanish flu and the better understanding of disease transmission which followed created concern for bioweapons (Zander *et al,* 2001). Furthermore, the fact that the Spanish Flu caused more fatalities than WWI caused concern among the great powers (Zander *et al*, 2001). Major European powers began to fear the feasibility of biological warfare and what pathogen would be best for weaponisation (Zander *et al*, 2001: 3). This fear of biological warfare was bought on by faulty intelligence and the fear of vulnerability according to Zander *et al* (2001: 3). Germany and the Allies tried, but no operational offensive Bioweapons (BW) were made before the end of WWII (Zander *at el*, 2001: 3). The Spanish Flu did show how effective and destructive a bioweapon could possibly be.

The Geneva Protocol had banned the use of bioweapon, however, they did not prohibit the development and stockpile of bioweapons (Allen & Chan, 2017: 101). Major powers were experimenting with various apparatus before, during and after WWII (Spiers, 2010: 22). The Soviets had perfected an attack termed a 'line source' which involved the dispersion of a biological agent from either tanks or Ilyushin bombers in a straight line for hundreds of miles (Spiers, 2010: 22). The purpose of this attack was to wipe out livestock and crops in wide areas over a matter of months (Spiers, 2010: 22). Furthermore, the Allies pursued bioweapons after Germany had invaded France, this was done due to the fear of the Germans gaining French biotech capabilities. The UK with the assistance of the US and Canada had mass-produced bioweapon munitions during WWII (Allen & Chan, 2017: 101). However, the US based its policy on 'no-first-use' and bioweapon munitions were developed as a deterrent to stop other states using them on them (Allen & Chan, 2017: 101). Despite all major powers pursuing bioweapon capabilities, the only offensive use of bioweapons was by the Japanese against the Chinese and attempted against the Russians (Allen & Chan, 2017).

The Japanese used disease-infected fleas, kamikaze soldiers, and the infecting of water wells as vectors for their attacks (Allen & Chan, 2017: 102). According to Christian (2013: 730), the Japanese developed a Uji bomb which was filled with pathogenic bacteria and fleas which it used against the Allies. These plague-infested fleas were also sprayed by the Japanese by using aircraft, resulting in an outbreak that killed 50 000 people (Christian, 2013: 730). Other similar attacks are believed to have caused a death toll of more than 100 000 (Christian, 2013: 730). Allen and Chan (2017) state that although these attacks did cause significant suffering among Chinese civilians, they did not give the Japanese a significant military advantage during WWII. Christian (2013: 730) further highlights other examples used in WWII, such as the French and Germans trying to use insects to destroy crops, the Soviets using typhus-infected lice against the Germans, the Japanese trying to experiment with aerosolised anthrax, and there were allegations of tularaemia being used by the Soviets against the Germans. These highlights indicate the pursuit of bioweapons during WWII but they did not subsequently deliver the military capability many states had hoped for.

### 3.2.2    Cold War era biowarfare

After WWII the most prominent developers of bioweapons, according to Zander *et al* (2001), were the US and the Soviets. Both of these nations saw biological weapons as having the potential for destructive power that could be compared to nuclear weapons (Allen & Chan,

2017). This subsequently led to the US and the Soviets pursuing bioweapon programmes. Allen and Chan (2017: 102) stated that by the mid-1960s the US was spending $300 million annually on chemical and biological weapon programmes. There is no credible evidence that the United States ever used bioweapons (Allen & Chan, 2017). However, the US bioweapons offensive programme was short-lived as they terminated it by 1969 in favour of a biodefence programme (Allen & Chan, 2017). Davis (1999: 509) states that the US offensive bioweapon programme was banned owing to several factors, such as secret intelligence information, the Vietnam war and new technological innovations. President Nixon stated in a formal policy review in 1969 that the US would formally dismantle its bioweapons. Furthermore, the US began negotiating with the Soviets and other nations which culminated in the Biological Weapons Convention (BWC) in 1972 (Allen & Chan, 2017: 102). Both nations signed it and it was enacted in 1975; however, it is arguable whether both sides truly complied.

The Russians were fully committed to their bioweapons programme, according to Davis (1999), hiding this program behind the façade of biotechnology.

> "By addressing every aspect of weapon production, from selection of new strains of organisms to the behaviour of biological aerosols under every possible condition of climate and topography, through the genetic engineering of antibiotic resistance and the design of optimum dissemination and delivery systems, the former Soviet Union was able to envisage the achievement of a miniaturized mobile production and weapon-making capability invulnerable to clandestine monitoring, invasive arms inspection, or attack in the event of war (because it was beyond identification); agents precisely matched to particular scenarios and human targets and incapable of being treated; a variety of dissemination systems, including cruise missiles; agents resistant to degradation by heat, light, cold, UV radiation, ionizing radiation, and various antibiotics; and dry formulations of agents capable of remaining viable in long-term storage" (Davis, 2019: 511).

The Russians were able to arm MIRVed warheads with a plague (Davis, 1999). Defectors helped the West challenge Russia openly in 1993 about their programmes; however, the capabilities remained largely unknown (Davis, 1999).

It was not only the Russians that had questionably defied the BWC. After Nixon renounced the American bioweapons programmes, several government departments conducted biodefence work that may have violated the treaty (Wheelis & Dando, 2003: 44). Before dealing with this,

what exactly was the US bioweapons programme before the BWC? In 1969 the US was able to produce 650 tons of agent per month that could then be filled into weapons, these type-classified agents were produced at plants such as Pine Bluff in Arkansas (Davis, 1999: 509). According to Davis (1999: 509) this was a thriving offensive programme that was eventually abandoned due to a mixture of politics, intelligence, other technological developments, and the Vietnam War. In brief, it can be stated that the US and Russia both took their biotechnology seriously, showing the US's need to pursue technological innovation. It can be seen how the US saw the potential for biotechnology to rival nuclear weapons and subsequently pursued it by spending $300 million on the programme annually (Allen& Chan, 2017: 102). The barriers between peaceful biotechnology and bioweapons has changed which poses a threat to effective proliferation management (Allen & Chan, 2017: 103). The rise of the commercial biotech industry has made knowledge and resources widespread and affordable (Allen & Chan, 2017). Furthermore, they are easy to hide. This means that the US has spent billions on their biodefence in fear of a terrorist attack. The next section will look closer at the capability of biotechnology and the issue of it being dual-use.

### 3.2.3    Biowarfare capability and current threats

Biotechnology is very appealing when it comes to military application owing to the fact that it is comparable to nuclear weapons in damage and as pointed out by Spiers (2010) there is only a need for a small amount to get an effect. Spiers (2010) highlights a quote by Judge William Webster who said: "biological warfare agents, including toxins, are more potent than the deadliest chemical warfare agents, and provide the broadest area coverage per pound of payload of any weapons system". Spiers (2010) goes on to highlight how these weapons have not been used extensively in warfare and have only been used in demonstrations. This section has highlighted so far, the immense potential of what biotechnology can do for warfare and how states have actively pursued these programmes.

### *3.2.4    Biotechnology*

This section will discuss the technical side of bioweapons. Firstly, it will look at major innovations in biotechnology and then at apparatus used to disperse these agents and current examples of biowarfare threats.

The work of Zander *et al* (2001) highlights two important breakthroughs in biotechnology: genomics and proteomics. Furthermore, van Aken and Hammond (2003) highlight that the

breakthrough in genetics is problematic for international peace and security. The risk of biowarfare has increased owing to a revolution in biotechnology such as new tools for analysing and modifying an organism's genetic material (van Aken & Hammond, 2003). According to van Aken and Hammond (2003), there are several factors that have increased this risk. Firstly, research and development in the medical sector has led to the availability of knowledge and facilities; secondly, facilities can be converted to create bioweapons and many countries have them; thirdly, by using genetic engineering, biological researchers have already developed weapons better than their natural counterparts (van Aken & Hammond, 2003). Making bacterial pathogens antibiotic resistant is an example of genetic modification. An example of this given by van Aken and Hammond (2003) is the Soviet's 'invisible anthrax which involved the introduction of an 'alien gene' that altered the immunological properties making existing vaccines ineffective against it. A current example of such a technology is highlighted by the CRS and is known as CRISPR-Cas9. CRISPR-Cas9 is a genetic modification or DNA-modifying technology which is 'low cost' technology' (Advanced gene editing: CRISPR-Cas9, 2018). This technology is capable of being relatively easy to use and delivers a high level of precision and efficiency (Advanced gene editing: CRISPR-Cas9, 2018).

The CRS states that such a technology could be used to either enhance or degrade military personnel (Advanced gene editing: CRISPR-Cas9, 2018). The proliferation of synthetic biology may lead to the creation of genetic code that doesn't exist and allow for the increase in the number of actors that are able to create such a technology (Advanced gene editing: CRISPR-Cas9, 2018). Owing to this issue, the US Intelligence Community's Worldwide Threat Assessment cited genome editing as having the potential to be a weapon of mass destruction (Advanced gene editing: CRISPR-Cas9, 2018). An example of gene editing being done by CRISPS-Cas9 is with *Aedes aegypti* mosquito (Caplan, Parent, Shen, & Plunkett, 2015). Caplan *et al* (2015) state that researchers, both academics and at private biotech firms, aim to either stop the female from carrying a disease or to make the male sterile by editing its genes. Aken and Hammond (2003) highlight the criteria for a successful bioweapon: it must be capable of being produced in large amounts, it must act fast, it must be robust in the environment, and it must be treatable to protect one's own soldiers. A theoretical argument could state that CRISP-Cas9 could allow a potential aggressor to edit the genes of a mosquito to make the disease it is carrying more severe and at the same time make their soldiers immune to it. CRISP-Cas9 could allow for researchers to edit how fast it acts and how robust it is in the environment, which would then allow it to fit into Aken and Hammond's (2003) criteria. This

is a very rudimentary example; however, biotechnology such as CRISPR-Cas9 offers a very disruptive potential and high capability potential for the future battlefield. These criteria are also highlighted in the next section. They must also be environmentally robust, meaning that they must be capable of being delivered successfully into the environment.

When it comes to other forms of delivery method of biotechnology, the main method seen in the literature is aerosol, meaning that it is released into the air (Spiers, 2010; & Davis, 1999). There are other delivery methods, and some have been mentioned in this section, for example missiles and rockets. The most worrying aspect was the testing of Intercontinental Ballistic Missiles (ICBMs) and cruise missiles as a delivery method by the Soviets (Spier, 2010: 23). But this section will focus on the aerial dispersion of biowarfare agents, as the age of drones may make such an act easier to conduct than arming a warhead. The main issue with this delivery system is that once the agent is in the wind it is susceptible to the wind and can only be stopped with special clothing and collective protective devices (Spiers, 2010: 13). Once they are in sunlight and exposed to other environmental factors, they must retain their viability and virulence (Spiers, 2010: 13). Furthermore, according to Spier (2010: 13), they may infect other organisms and spread to cause an epidemic. The final aspect is the need for only a small amount. Aerial warfare allowed for innovation in BW (Spiers, 2010). The attacks could use spray apparatus, with either pressure or no pressure release (Spier, 2010). These attacks could be transported at subsonic speeds and at low levels which could allow for surprise attacks and a large area of coverage (Spier, 2010:21-22). However, bioweapons dispersed from aerosol are largely unpredictable as they are dependent on wind, rain, temperature and geography. But at the same time, they only need a small quantity to cause a large amount of damage and they are easily concealed.

The Covid-19 pandemic has had a severe impact on states globally and this impact will be continued to be seen for years to come. There is a large amount of uncertainty about Covid-19; however, it has highlighted how unprepared the US is for a targeted biological attack. Bertrand and Lippman (2020) state that the Covid-19 pandemic has shown, at all levels, how the US government was completely unprepared to deal with a pandemic that moved slowly. A biological attack would undoubtedly overwhelm the US (Bertrand & Lippman, 2020). From the examples used above, a biological attack using drones armed with aerosol or CRISP-Cas9 modified disease could be potentially damaging for the US. Especially if the disease released by one of the two technologies was fast-acting and deadly. This technology combined with the real-world effects seen by the Covid-19 pandemic highlights just how damaging potential

biowarfare can be. However, as highlighted by Bentley (2020), the Covid-19 pandemic offers the opportunity to create an understanding of what a biowarfare attack could look like. These technologies and the Covid-19 pandemic show the destructive power that biowarfare could have on the world.

### 3.2.5    Summary

The aim of this section was to look at biotechnology and bioweapons in order to create a contextualisation of how technology affects how states conduct warfare. This section decided to only look at bioweapons and not chemical weapons. What was found in this section is that biowarfare is not new and has been around since the Athenians and Mongolians. It is a tactic that has been employed by militaries for centuries and can have potentially devastating effects as seen with the Mongolians and how they potentially started a plague that went through Europe. The next area of biowarfare that was highlighted was WWI, the Spanish Flu and WWII. This gave this research an opportunity to look at the role that biowarfare played during twentieth century warfare. This section looked at different definitions of what biowarfare is. It was found that biological weapons are weapons that are created to intentionally harm or kill humans, plants, or animals through the use of a biological agent (Biological weapons: a primer, 2001: 1). Such agents are bacteria, viruses, rickettsia or fungi according to the CRS (Biological weapons: a primer, 2001: 1). The next section discussed the role of biowarfare in the Cold War. It found that both the Soviets and the US had extensive bioweapon programmes at that time. However, the US policy based itself on non-first-use strategy and used it only to deter aggressors. The final part of this biowarfare section looked at current threats by looking a major innovation of biotechnology and potential dispensing apparatus. This section considered both the strengths and weaknesses of biotechnology. It can be concluded that biowarfare has immense potential to be both destructive and unreliable. It offers the ultimate strategy in warfare as it will allow states to win wars without destroying vital infrastructure.

### 3.3    Nuclear technology

When it comes to how a technology influenced how states went to war there is no better example than nuclear technology. Nuclear weapons ultimately became a key part of the US's national security strategy (U.S strategic nuclear forces: background, development, and issues, 2020). The CRS goes on to state how the US does not have plans to eliminate their nuclear weapons or to remove their nuclear deterrence strategy (U.S strategic nuclear forces: background, development, and issues, 2020:1). They further highlight how nuclear deterrence

has been a key strategy of the US for the last 60 years (U.S strategic nuclear forces: background, development, and issues, 2020). Nuclear weapons have arguably been able to maintain stability over this time. Younger (2000) states that due to the destructive potential of nuclear weapons, any form of conflict between superpowers became an unviable option. This allowed the US to maintain their deterrence capability to protect both their allies and interests. The Trump administration, according to CRS (U.S strategic nuclear forces: background, development, and issues, 2020), in its nuclear posture review stated a continued support for the US's nuclear arsenal. The US maintains the need for its nuclear arsenal in the current threat environment and the need to modernise the arsenal (U.S strategic nuclear forces: background, development, and issues, 2020). This highlights the importance of nuclear weapons in the US's current national security strategy. Furthermore, it must be stated that this section does not intend to look at nuclear deterrence theory, this has already been covered in-depth in Chapter 2. This section aims only to research the most significant aspects of nuclear weapons such as the Manhattan project, the nuclear triad, and the importance of modernising the nuclear triad.

### 3.3.1    The Manhattan Project

The Manhattan Project is an example of the cost as well as impact a transformative technology can have on how states conduct war. This is a sterling example as AI could potentially face a similar trajectory when it comes to cost and impact. This section also allows a brief overview of the origin of nuclear weapons. The Manhattan Project got under way by the military in the US in 1942 At that time, its three-year cost of $2 billion[6] made up only 1 percent of the US's GDP (Allen & Chan, 2017: 71). The US went on to enlist the world's leading mathematicians, engineers and scientists for the Manhattan Project (Allen & Chan, 2017: 71). The Manhattan project was born out of a small research programme in 1939, the roots of which came from the fear of Nazi Germany gaining nuclear weapons before the US (Manhattan Project, 2007). Due to these concerns, key scientists proposed that the US accelerate their atomic research response (Stine, 2009). The project ran from 1941-1946 under the control of the United States Army Corps of Engineers and administered by General Leslie R. Groves (Manhattan Project, 2007). Once the project expanded it had multiple secret production and research facilities, with the primary ones being: "the plutonium-production facility at what is now the Hanford Site (Washington state); uranium-enrichment facilities at Oak Ridge, Tennessee; and the weapons research and design laboratory, now known as Los Alamos National Laboratory in New

---

[6] This was $2 billion in 1940 and not a conversion to current equivalent (Allen, & Chan, 2017).

Mexico" (Manhattan Project, 2007). By 1945 two designs came from the project, "one utilizing enriched uranium and the other its newly discovered derivative, plutonium" (Manhattan Project, 2007). By 1948 the US had created enough parts to create 56 atom bombs, that increased to 300 by 1950 and in 1967 the size had increased to 31 255 nuclear warheads (Allen & Chan, 2017: 73). The Manhattan Project was a key initiative for the US during WWII and more importantly during the Cold War. It allowed the US to gain nuclear weapons and ultimately maintain their technological innovation. They were further allowed to show their dominance by using these weapons against Japan to ultimately end WWII.

The first use came after Roosevelt gave the final approval in 1942 to create the first nuclear bomb (Stine, 2009). This led to the first successful test of a bomb using plutonium in July 1945 (Stine, 2009). The Trinity test was detonated on July 16th 1945 just before sunrise at 5:30 a.m (Reed, 2014). The test, although rushed due to the intense pressure put on Groves[7], was deemed a spectacular success with the yield being estimated at 21 kt TNT equivalent (Reed, 2014). Reed highlights a description of the explosion written by Groves that was sent to Henry Stimson Potsdam:

> "The light from the explosion was clearly seen at Albuquerque, Santa Fe, El Paso, and other points generally to about 180 miles away. The sound was heard … generally to 100 miles. Only a few windows were broken, although one was some 125 miles away. A crater from which all vegetation had vanished, with a diameter of 1200 ft … in the center was a shallow bowl 130 ft in diameter and 6 ft in depth … The steel from the tower was evaporated … I no longer consider the Pentagon a safe shelter from such a bomb … My liaison officer at the Alamogordo Air Base, 60 miles away (reported) a blinding flash of light that lighted the entire northwestern sky" (Grooves as cited by Reed, 2014: 23).

The explosion from Trinity was so powerful that it could have been seen from the moon (Reed, 2014). The following month President Truman decided to use a nuclear bomb against Japan (Stine, 2009). For this mission to be executed a special Army Air Force unit known as the 509th Composite Group was established (Reed, 2014). The aircraft that were to be used in order to conduct this mission were B-29s called Enola Gay and Bockscar (Reed, 2014). The atomic bombs that were used were known as the 'Little Boy' and 'Fat Man' (Reed, 2014). Little Boy

---

[7] Brigadier General Leslie R. Groves was the head of the Manhattan Project, succeeding Col. James Marshall. Reed states that Groves was the perfect choice as he had vast experience in military construction (2014).

was released from Enola Gay at a height of 1900 ft at 08:16 local time. It caused an estimated 15kt yield, killed 66 000-69 0000 people, and injured an estimated 255 0000 people (Reed, 2014). The release of this bomb did not force the Japanese to surrender and this led to the release of the 'Fat man' onto Nagasaki [8](Reed, 2014). The Fat Man yielded a load of 21kt and the casualties were estimated at 39 000 dead and 25 000 injured (Reed, 2014). Japan surrendered a few days after the two bombs were used against them and the project was considered fulfilled (Stine, 2009). This remains as the only use of nuclear weapons in warfare. The atomic bombs that were dropped on Hiroshima and Nagasaki ultimately changed the nature of warfare. The atomic bombs used showed the true destructive power of this technology and have managed to stop any major conflict between superpowers since its invention. Whether the US should have used the atomic bomb against the Japanese is highly contested by theorists. What can be concluded, with a final emphasis, is that the invention of the atomic bomb changed the way states conducted their strategic security policy and stopped the outbreak of war between superpowers. The next important aspect of nuclear weapons was developed during the 1960s, it was the delivery method and was termed the 'nuclear triad'. This section will not look at the Cold War arms race or dynamics as it cannot look at everything. It aims to look at the origin of nuclear weapons and the importance of the triad. This research sees these two sections as important aspects of nuclear weapons. This now brings this section to its most important topic, looking at the US nuclear triad.

### 3.3.2    The US nuclear force structure: The nuclear triad

The nuclear triad consists of three different areas of the US army: land, sea and air. The CRS states that the nuclear triad came about because each military force wanted to play a role in the US nuclear arsenal (U.S strategic nuclear forces: background, development, and issues, 2020). Analysts argued that each section had its strengths and weaknesses; furthermore, it would complicate a Soviet attack and allowed for some of the force to survive to allow a second-strike (U.S strategic nuclear forces: background, development, and issues, 2020). Futter and Williams (2016) state the same premise, that the US developed its nuclear triad in order to secure a credible nuclear deterrent. Each leg of the triad offered different strategic advantages to the US, with all three combined it allowed the US to not need to fear a 'bolt out of the blue' attack, as the triad allowed the US to overcome potential systems failures (Futter & Williams, 2016:

---

[8] The Fat Man atomic bomb was originally destined from Kokura, however, smoke and faulty equipment made it difficult to conduct a visual raid (Reed, 2014). This made Bockscar change its target to Nagasaki and forced it to make an emergency landing in Okinawa (Reed, 2014).

247). This allowed the US to have a capable nuclear force that was able to retaliate because of the triad (Futter & Williams, 2016: 247) Futter and Williams (2016) state that the nuclear triad came into existence owing to the technological advances in ICBMs and Submarine-launched Ballistic Missiles (SLBMs).

This gave the US capability and flexibility owing to the fact that the US bombers were vulnerable to Soviet air defences (Futter & Williams, 2016). "Long range bombers remained a useful way of signalling intent and could be called back; ICBMs offered the ability to respond rapidly and with massive firepower; while SLBMs provided the ultimate guarantee through their ability to remain undetected beneath the surface of the ocean" (Futter & Williams, 2016: 247). The CRS states that ICBMs had the ability to respond quickly and accurately at hardened targets such as Soviet command posts and ICBMs Silos (U.S strategic nuclear forces: background, development, and issues, 2020). The SLBMs had the survivability as they complicated the Soviet's ability to dismantle the US second-strike, and bombers could be released quickly and subsequently recalled if the crisis de-escalated (U.S strategic nuclear forces: background, development, and issues, 2020, 3). Now that a general outline and importance of the nuclear triad has been done, this section now aims to give a more in-depth study of the three different domains: air, land and sea, starting with their capabilities during the Cold War and the subsequent plan for the US to upgrade each leg of the triad.

### 3.3.3    Air

Before the conclusion of the START treaty with the Soviet Union in 1990 the US had 94 B-52H bombers, 96 B-1 Bombers and older B052G bombers (U.S strategic nuclear forces: background, development, and issues, 2020). Furthermore, they had two of the new B-2 bombers (U.S strategic nuclear forces: background, development, and issues, 2020). This force consisted of 260 bombers with the ability to carry over 4 648 weapons (U.S strategic nuclear forces: background, development, and issues, 2020). The current Air part of the triad owing to the START treaty consists of 46 nuclear-capable B-52H and 20 nuclear-capable B-2A strategic bombers (Nuclear Posture Review (NPR), 2018). These bombers can carry a variety of nuclear weapons that add to the US triad's capabilities and flexibility. The B-2A carried a gravity bomb while the B-52H bombers carried ALCMS which allowed for different yield options (NPR, 2018: 47). Meanwhile the B83-1 and B61-11 have the capability to hold a variety of protected targets at risk (NPR, 2018: 47). This highlights the importance of technology in how the US conducts warfare; it also highlights the importance of weapons which will be looked at in their

own section. Ultimately bombers play a critical role in US strategy, as they are able to take off and be recalled if the crisis de-escalates, thus creating flexibility (NPR, 2018). Furthermore, they have an unlimited range as they can be fuelled during flight, allowing for flights abroad to show their capabilities (NPR, 2018). The NPR (2018: 47) goes on to state that the only long-range bomber capable of penetrating the enemy defence system is the B-2A.

The US plans to modernise the B-2A along with the B-52H to make sure that they remain effective into the future (NPR, 2018: 50). With the improvement and current proliferation of adversaries' air defences the US has aimed to initiate a programme to build and deploy the next generation bomber (NPR, 2018: 50). The B-21 Raider project aims to replace the currently ageing B-52H and its ALCM and B-2A missiles (NPR, 2018: 50). This bomber will replace the nuclear-capable bomber force beginning mid-2020s (NPR, 2018: 50). Northrop Grumman indicates that the B-21 will be able to penetrate the toughest air defence and be capable of delivering a precision strike anywhere in the world; they further state that it is the future of deterrence (B-21 Raider, 2020). Northrop Grumman state that the B-21 Raider will be designed to have a high survivability, be long range, and capable of carrying a mix of nuclear and conventional ordinance (B-21 Raider, 2020). This will join the nuclear triad as a flexible and visible deterrent which will maintain national security objectives and reassure the US's allies and partners (B-21 Raider, 2020). The US Air Force aims to acquire 100 of these aircraft, with Northrop Grumman claiming they may acquire as many as 200 (B-21 Raider, 2020). This shows the US's commitment is to nuclear security strategy, the nuclear triad, and maintaining its technological dominance when it comes to fighting wars.

### 3.3.4    Sea

Futter and William (2016: 247) state that SLBMs allowed for the ultimate guarantee for a secure second-strike owing to the fact that they could remain undetected below the surface of the ocean. Furthermore, CRS highlights how the former Commander of the US Strategic Command (STRATCOM) stated that SLBM force provided survivability (U.S strategic nuclear forces: background, development, and issues, 2020: 6-7). This means that SLBMs allowed the US to secure their second-strike capability, which is pivotal for maintaining nuclear deterrence and Mutually Assured Destruction (MAD). The NPR (2018: 44) highlighted in 2018 that the US was currently operating Ohio-class Submarines that were Ballistic Missile equipped and Nuclear powered (SSBNs). The missiles that these submarines were equipped with were Trident II (D5) SLBMs which proved a 'sea-based deterrent force'. "When on patrol, SSBNs

are, at present, virtually undetectable, and there are no known, near-term credible threats to the survivability of the SSBN force" (NPR, 2018: 45). Their ability to be constantly ready and have intercontinental range allows for these submarines to hold selected targets at risk throughout Eurasia and the Pacific Ocean. This shows that SSBNs are both undetectable and very capable of taking out targets faraway, making them a very credible threat to aggressors.

Furthermore, NPR (2018: 45) states that they are able to hold highly accurate and high-yield warheads that leave targets at high risk. These warheads can travel at hypersonic speed and can subsequently reach their targets quickly after launch (NPR, 2018: 45). Finally, according to the NPR (2018), these SSBNs are highly mobile which can demonstrate US nuclear presence. This will reassure their allies and show their commitment to deterrence. The OHIO-class is expected to have a life-span lasting until 2042. It cannot be extended further than this and it has been in service since 1981 (NPR, 2018). The NPR (2018: 45) goes on to state that advances in adversaries' anti-submarine warfare and missile defence may cause challenges for the SSBN system, showing the US's continued outlook of striving to be more technological advanced than its adversaries. What this section on SSBN shows is the importance of technology when it comes to how states conduct warfare. The introduction of the OHIO-class SSBN allowed the US to secure their second-strike capability, creating a credible threat of destruction to any aggressors and MAD to other nuclear-armed states.

### 3.3.5    Land-based deterrent force

ICMBs allowed the US to gain responsiveness and accuracy to take out potential hardened targets such as Soviet command posts and ICBM silos (U.S strategic nuclear forces: background, developments, and issues, 2020). According to the NPR the United States ICBMs force had 400 single-warhead Minuteman III ICBMS in 450 underground silos that are spread across several different states (NPR, 2018: 45). NPR (2018) also highlights how responsive the ICMS are and how this is of importance to the nuclear triad. Furthermore, ICBMs are highly survivable against all attacks, apart from a large-scale nuclear attack (NPR, 2018: 46). To successfully carry out such an attack, the aggressor must launch a precise attack with hundreds of high-yield warheads, a task that the NPR (2018) states can only be conducted by Russia. Without the ICBM force, other parts of the triad could potentially be subjected to a nuclear first-strike that involves a relatively smaller number of nuclear weapons (NPR, 2018: 46). The part of the triad that could potentially be vulnerable would be SSBNs in port and non-alert bombers (NPR, 2018: 46). The US's ability to launch a strike quickly means that no adversary

could be confident in their ability to destroy them prior to launch (NPR, 2018). Furthermore, the range of these weapons can hold targets at risk throughout Eurasia, reaching any target in 30 minutes or less (NPR, 2018).

Once again, the NPR (2018) has stated that the US plans to modernise its ICBMS due to their critical importance in the nuclear triad. The US will be fielding the Ground-Based Strategic Deterrent (GBSD) as a replacement to the ICBMs by 2029 (NPR, 2018: 49). This replacement aims to modernise 450 ICBM launch facilities in order to support fielding 400 ICBMS that will replace the Minuteman III that are being retired after six decades of service (NPR, 2018: 50). This will make the land-based deterrent leg of the triad effective into the future. This once again highlights how the US aims to maintain its technological edge over its adversaries. This section has, however, not touched on the reason behind the US's justification to renew its nuclear triad in the NPR. The following section will discuss the main drivers behind this renewal.

### 3.3.6 The US's motivation to maintain technological edge over its adversaries

This section aims to outline the motivation behind the NPR and why the US aims to renew its nuclear triad. Klare (2019) highlights that the Trump administration claims that the US infrastructure is out of date and that China and Russia have taken advantage of the US complacency. The NPR (2018: V) states that they have added new nuclear capabilities and have also increased the salience of their nuclear forces in their plans and strategies. Furthermore, according to the NPR (2018: V), they have shown more aggressive behaviour that spans cyberspace and outer space as well. The US maintains that Russia and China are contesting the international norms and order that have been built by the US and its allies (NPR, 2018: 2). Furthermore, while the US has reduced the number and salience of their nuclear weapons, the US claims that China and Russia are moving in the other direction (NPR, 2018: 2). Furthermore, China's expansion of its nuclear forces has little or no transparency according to the US NPR (2018). Due to the current threat environment combined with future uncertainties, the NPR (2018: 3) mandates that the US needs to commit to a modern and effective nuclear force along with the infrastructure needed to support it. Due to these factors, the US started to pursue the modernisation of their nuclear forces. It must be noted that this section does not aim to debate whether the US is correct or incorrect in these assumptions. It aims to purely understand the US's current perception of the global threat environment and

whether or not this makes the US feel more or less secure. Subsequently it can be stated that the activities of Russia and China have made the US feel insecure in their technological edge. This has subsequently forced them to pursue the modernisation of their nuclear triad, as the US perceives this modernisation as critical to deterring threats.

### 3.3.7    Summary

The aim of this section was to look at nuclear technology and how it changed the way the US went to war. This section also highlighted how important nuclear weapons are to the US and how they are subsequently becoming even more important. The section began with an overview of the Manhattan project to create an understanding of the origins of nuclear weapons. From there this section went on to discuss the structure of the US nuclear forces by looking at the nuclear triad. The triad plays a vital role in maintaining the credibility of the US's nuclear forces and it does this by having three different legs. The three different legs are based in the different domains of combat: land, air and sea. All these legs combined help cancel out each other's weaknesses and this is what makes the triad so effective. The most important leg of the triad is the sea-based leg, more specifically, SSBNs which help to maintain the credibility of the US's second-strike capability.

## 3.4    Cyberwarfare

This section aims to look at cyberwarfare, another form of a transformative technology. The government has played an important role in the development of digital computing, cryptography, and the internet network; furthermore these three technologies enabled the creation of the modern day cyberspace (Allen & Chan, 2017: 91). According to Allen and Chan (2017) these three technologies enabled cyberspace by:

> "1. Digital computing (especially using silicon integrated circuits), which allows storage and processing of information by machines 2. Internet networking, which allows for the connection and unification of different types of networks according to a single standard, namely internet protocol 3. Cryptography, which allows for unrelated users to share data and infrastructure while maintaining data confidentiality and integrity All three technologies were actively supported by the US. government." (Allen & Chan, 2017: 91).

The US government was in full support of all three of these technologies and this allowed for the creation of the internet in the 1970s; from the mid-1990s it became more commercial (Allen

& Chan, 2017). Furthermore, the emergence of cyberspace has revolutionised how global communities interact with each other (Ween, Dortmans, Thakur, & Row, 2017: 337). It also bought with it vulnerabilities, such as malicious use of information and communication technologies (ICT); this is used in order to interfere with physical, electronic and online system functions (Ween *et al*, 2017: 337). This is what Ween *et al* (2017) refer to as cyberwarfare and subsequently the main point of this section. This section will begin by defining cyberwarfare and the implications that cyberwarfare has for military strategy. The next part of this section will then discuss the most prominent example of cyberwarfare; Stuxnet. The final section will consider the dual-use of cyber technology. This chapter will outline how the US seeks to maintain technological dominance over its adversaries and this is an implicit example of how technology will affect how states conduct war.

### 3.4.1    Defining cyberwarfare

As stated already, Allen and Chan (2017) see cyberspace as the connection of digital computing, cryptography, and internet networking. Robinson, Jones and Janicke (2015: 72) define cyberspace as a 'global domain' within the information environment which uses the electromagnetic spectrum and electronics to store, modify, exchange, create and exploit information. This is done by interdependent and interconnected networks which use ICT (Robinson *et al*, 2015: 72). Libcki (2009: 6) also states the same: "[c]yberspace, as such, can be characterized as an agglomeration of individual computing devices that are networked to one another (e.g., an office local-area network or a corporate wide-area network) and to the outside world". This shows that cyberspace is a global, as well as local, connection of computing devices to one another using ICT. Solis (2014) states that cyber refers to computers and computer networks, not just the internet but all things that connect to computers. While this is not a physical battlefield, the actions within cyberspace can enact themselves into the physical domain (Ween *et al*, 2017). Cyberspace is a global commons which is decentralised, meaning that no nation state has the control over its resources and it extends outside their borders (Ween *et al*, 2017). Furthermore, Ween *et al* (2017), state that intent and attribution of actions that have taken place in cyberspace are hard to determine, which may affect the response and create a 'strategic shock'. This shows what cyberspace can be defined as and starts to unravel the complications of cyberwarfare.

As already stated, Ween *et al* (2017) see cyberwarfare as the use of ICT in order to interfere with physical, electronic, and online systems. Solis (2017: 3) defines cyberwarfare more

specifically as warfare waged in space, which includes the defending of information and computer network. Solis (2017: 3) goes on to further state that it involves deterring information attacks and also denying an adversary from doing the same. It can also have an offensive nature by operating information operations against an adversary (Solis, 2013: 3). What Solis (2017: 3) means is that cyberwarfare entails defence, offence, and deterrence. Theohary and Rollins (2015) define cyberwarfare as state-on-state action that is the equivalent of an armed attack or an act in cyberspace that may equal a kinetic force response. However, Theohary and Rollins (2015) state that there are no clear criteria to define when a cyberattack can equate to kinetic force. According to Theohary and Rollins (2015) there isn't a criterion that clearly determines whether a cyberattack is an act of terror, a nation state equivalent of a kinetic attack or hacktivism. The US took the stance under Article 2(4) of the UN Charters and customary international law; meaning that a cyberattack that results in death, destruction or injury will be viewed as an act of kinetic force (Theohary & Rollins, 2015). What this shows is that cyberwarfare is defined based on what is a cyberattack.

Solis defines (2014: 12) a cyberattack as: "a trans-border cyber operation, whether offensive or defensive, that is reasonably expected to cause injury or death to persons, or damage or destruction to objects." Solis (2014) also defines the difference between a cyberattack and 'cyber intrusion', stating that a cyber intrusion does not rise to a level of a cyberattack.

> "A cyber attack, as opposed to a cyber intrusion, constitutes a "use of force" if undertaken by a state's armed forces, intelligence services, or a private contractor whose conduct is attributable to the state, and its scale and effects are comparable to non-cyber operations that rise to a level of a use of force." (Solis, 2014: 15).

Solis separates the two by the results of the attack, which equates to the use of force (2014). Robinson, Jones and Janicke (2015) stated that a cyber is based on an actor's intention. They thus define a cyberattack as an attack with a warfare-like intent. Theohary *et al* (2015) highlight different types of cyber 'threat actors' across the cyberwarfare ecosystem. These actors are cyberterrorists, cyberspies, cyberthieves, cyberwarriors, and cyberactivists (Theohary *et* al, 2015). This section will now compare the definitions of cyberterrorists, cyberspies and cyber warriors as stated by various researchers. Jones *et al* (2015) define cyberspies as actors who steal information used by governments or private corporations that will allow them to gain a strategic, financial, security, or political advantage. Cyberwarriors are defined as actors that develop capabilities and use these to pursue a state's strategic objective (Theohary *et al*, 2015).

Theohary *et al* (2015) state that these entities are either acting on behalf of a government or not when it comes to who they target, the timing of the attack, and the type of attack that is used. They are then also blamed when accusations are made by the attacked state. Theohary *et al* (2015) state that these threats will often cross over and are ultimately hard to assess. Furthermore, there is no criteria to distinguish what type of cyberattack occurred (Theohary *et al*, 2015).

### 3.4.2 Examples of cyberattacks and Stuxnet

This section aims to look at different examples of cyberattacks or intrusions. This section does not aim to look at the laws of armed conflict, it merely aims to take a look at how authors view the attack and the consequences of such an attack. A few examples would be looked at, however, the main example to be focused on in this section will be Stuxnet. It is referred to by many authors and is a critical example of how intrusive and destructive a cyberattack could be.

### *3.4.3 Georgia and Estonia*

One example of a cyberattack is the use of a distributed denial of services (DDoS) (Solis, 2014; Theohary, 2015; Libicki, 2009; & Lin, Allhoff, & Rowe, 2012). There are two main cases cited in the literature: a DDoS attack on Estonia and Georgia. A DDoS involves an attack in which a server is overwhelmed with internet traffic so that access to this cite then becomes degraded or denied (Theohary, & Rollins, 2015: 1). The DDoS attack on Estonia was the result of the Estonians moving a Russian war memorial from the city to a military cemetery; this led to riots and a DDoS attack (Libicki, 2009, Ottis, 2008). According to Ottis (2008) the Russian minority in Estonia saw the statue as a representation of the liberation of Russian people while Estonians saw it as a symbol of oppression. Libicki (2009) states that it was tracked back to the Kremlin. However, whether the Russians were behind the attack is difficult to prove. Ottis (2008) states that the cyberattack campaign lasted for 22 days and were a part of a wider political conflict between Estonia and Russia. The DDoS attack was focused on state and commercial websites that ranged from defence and foreign ministry to media outlets and banks (Estonia denial of service incident, 2007). A DDoS attack is executed by overloading bandwidth of websites and overloading their services with 'junk traffic' (Estonia denial of service incident, 2007). Many well-known methods include udp flood, malformed web queries, ping flood, email spam and web queries (Ottis, 2008). More complicated methods were used such as a SQL injection, some were successful at 'non-critical sites' (Ottis, 2008).  In order to counter this the Estonian government temporarily closed its digital borders and blocked all international web traffic

(Estonia denial of service incident, 2007). The outcome of this attack was the conviction of 20-year-old Estonian student Dmitri Galuškevitš (Ottis, 2008: 3). The Estonian State Procurator made a request to the Russian Supreme Procurator for assistance in a formal investigation to find the attackers residing in Russia (Ottis, 2008). A report by the BBC on 27 April 2017 states that Estonia's request for help was ignored by the Russians (McGuinness, 2017).

In the case of Georgia, an internet security firm reported that a DDoS had taken place against websites in the state (Korns & Kasternberg, 2009: 60). An internet security firm reported that on 19 July 2008 a DDoS attack was aimed at websites in Georgia (Korns & Kasternberg, 2009: 60). The attack, that was carried out by Russia, also involved an invasion by air and land as well as a blockade at sea (Theohary & Harrington, 2015: 10). The result of this cyberattack left the Georgian government barely capable of communicating on the internet (Korns & Kasternberg, 2009: 60). Theohary and Harrington (2015) state that the Russian hackers besieged Georgia's internet for the duration of the armed conflict. Pernick (2018: 59) states that the attacks that were planned in advance but carried out on 9 August and attacked fifty-four Georgian websites which included ninety percent of state institutions' websites and a large number of .ge domain addresses. Theohary and Harrington (2015) state that the attacks in August started on 8 August and came as Russian tanks crossed the border into South Ossetia in Georgia. Theohary and Harrington (2015) also state that the attacks targeted 54 websites, however, the first attack targeted pro-Georgian hackers. This attack against the pro-Georgian hackers was unable to completely reduce the counterattacks against Russian Targets (Theohary, & Harrington, 2015: 10). As Russian troops moved in Georgia became unable to access the 54 websites that have been mentioned so far; the scope of these websites were government, communication, finance and critical information (Theohary, & Harrington, 2015: 11). Pernick (2018) states that this cyber espionage campaign was more sophisticated. It had Russian military connections and was pre-planned (Pernick, 2018).

A subsequent section will discuss how Israel 'prepared' the battlefield, which can be argued to be a similar tactic as that employed by the Russians. The pro-Russian hackers were able to attack critical Georgian infrastructure before and during the Russian invasion of Ossetia. Theohary and Harrington (2015) state that DDoS attacks were carried out prior to ground troop movements or bombings. These DDoS took out communications prior to these military actions (Theohary, & Harrington, 2015). Korns and Kastenberg (2009) state that the DDoS attack left the Georgian government cyber-locked and barely able to communicate over the internet. The Georgian government overcame these attacks by moving their critical official internet assets to

the US, Poland, and Estonia (Korns, & Kastenberg, 2009; Theohary, & Harrington, 2015). An example given by Theohary and Harrington (2015) tells how a web-hosting company called Tulip Systems gave refugee status to the Georgian government websites without the US government's approval.

Pernick (2018) draws some conclusions from these attacks on Estonia and Georgia. What was shown by Russia attempting to undermine Estonia was that it had enhanced capabilities when it comes to outsourcing cyberattacks, controlling information and military operation among various actors, and using strategic impact gained from cyber-espionage operations (Pernick, 2018: 60). Meanwhile, the Russian attack on Georgia showed how cyberattacks could be used to support military and strategic objectives (Pernick, 2018). Furthermore, the result of the attack on Estonia lead to the creation of the Cooperative Cyber Defense Center of Excellence (CCDCOE) in Tallinn, Estonia in 2009 (Theohary, & Rolins, 2015: 5). The Tallinn Manual was created at this centre. Despite its name, it was more an academic paper than a set of laws that states were bound to (Theohary, & Rolins, 2015: 5).

### 3.4.4    The United States

Solis (2014) highlights more US-centric examples of cyberattacks including ones that involve China and Russia. Solis states that China has two network-monitoring states in Cuba, one to monitor US internet traffic and another to monitor the US DoD. The 'Night Dragon' is an example of a cyber intrusion which ran from 2007 until Lockhead and Martin discovered it only in 2009 (Solis, 2014: 4). This cyber intrusion involved the theft of terabytes of information including the US F-35 fighter (Solis, 2014: 4). The Research and Development (R&D) for the F-35 cost in excess of $50 billion and the Chinese were believed to have acquired all the intellectual property of the F-35 program (Allen & Chan, 2017). In 2010 Pentagon systems were penetrated to see how the command-and-control systems could be crippled (Solis, 2014). The Pentagon uses a password and token security system (Solis, 2014). The token is a USB that makes a new number every sixty seconds (Solis, 2014). The company that made these were, according to Solis (2014), hacked by another foreign intelligence service. Using this information, they hacked Lockheed Martin in March 2011 and they subsequently lost 24 000 files, which: "included plans for missile tracking systems, satellite navigation devices, surveillance drones and top-of-the line jet fighters" (Solis, 2014). This continues to highlight the impact transformative technology can have on how states conduct warfare, in this scenario, even before weapons are made to fight these wars. Gaining valuable information of an enemies

weaponry can hand a state a strategic advantage, financially, as they will not have to spend excessive amounts on developing them.

### 3.4.5    Israel: Prepping the battlefield

Another example of cyberattack mentioned by Solis is how the Israelis 'prepped' a battlefield for an act of war (Solis, 2014: 6-7). This is an example of how a modern day transformative technology such as cyber can impact the way states conduct warfare. This section highlights how cyber intrusions handed the Israelis a strategic advantage before the conflict even started. AI and AWS may potentially offer even more valuable strategic advantages on the future battlefield. The Israelis took control of the air defence network at night-time and uploaded an image containing nothing, meaning the air defence missile couldn't be fired as they had no targets in the system (Solis, 2014: 7). This allowed the Israel Air Force to apparently fly into Syrian air space and bomb a reactor without altering their air defence systems (Cohen, Freiligh, & Siboni, 2016). Syrian fighter jets could not be scrambled owing to the fact that their system had no targets (Solis, 2014: 7). This was allegedly accomplished by taking over the Syrian air defence system and tricking them into thinking nothing was happening (Cohen *et al*, 2016). Even when the attack was under way the radar still did not show anything (Cohen *et al*, 2016). Israel allegedly decided not to shut down the radar as this would have alerted the Syrians; they instead reprogrammed the system to function 'normally' (Cohen *et al*, 2016). The site that was a target for the Israeli's clandestine attack was an alleged Syrian Nuclear reactor (Sharp, 2009). According to Sharp (2009) the Atomic Energy Agency (IAEA) released a report in 2009 that drew a connection between Syria and North Korea's clandestine nuclear program. According to the actual report by IAEA, the agency had been given information that alleged that the site destroyed by Israel was a nuclear reactor (IAEA annual Report, 2009). Furthermore, the report states that the site was not operational and was allegedly being built with assistance from the Democratic People's Republic of Korea (DPRK) (IAEA annual report, 2009). This shows how cyberattacks could be used to 'prep' future battlefields that are reliant on computer systems. This is giving a strong argument for how technology affects how states go to war and how it changes a state's strategy for conflict.

### 3.4.6    Stuxnet

Langer (2011) states that Stuxnet was the first cyberwarfare weapon. Stuxnet was aimed at the industrial controls of a nuclear centrifuge, using a SCADA application as a means of distribution. This is a Microsoft application in and it was the supposed target in order to access

the controller of the centrifuge (Langer, 2011: 49). A controller, as defined by Langer (2011: 49), is a real-time computer system that affects the outputs through electrical input signals and programming logic. The devices are connected to a controller's area drivers, pumps, valves, thermometers, and tachometers (Langer, 2011: 49). The controller communicates with these devices through a fieldbus connection (Langer, 2011: 49). Stuxnet was aimed at manipulating the controller and Microsoft was the delivery method to get to the controller.

Stuxnet infested an Iranian nuclear processing facility in 2010. This was its main target, but this worm went far beyond its intended targets (Lin *et al*, 2012: 25). However, although it infected many Windows computers, it was particular about the controller which it targeted (Lin *et al*, 2012; Langer, 2011). Stuxnet would only target controllers from Siemens, meaning that the worm would go through a series of complicated fingerprinting processes to make sure that it had the correct controller (Langer, 2012: 49). The program that was downloaded conducted a process that would check model numbers, configuration details and even download another program that checked if it was the right program (Langer, 2012: 49). After this process, it would then load a rogue code into the one of the Siemens 315 and 417 controllers (Langer, 2012). Furthermore, Stuxnet would manipulate the system to show that nothing was going wrong with the nuclear centrifuges (Solis, 2014). This set the Iranian nuclear programme back years as the code made the centrifuges violently self-destruct (Singer & Frieman as cited by Allen & Chan, 2017). This supposedly took out a fifth of the Iranian nuclear centrifuges (Allen and Chan, 2017).

Stuxnet went on to spread itself. This was problematic as the code was exposed and revealed its secrets (Solis, 2014). Solis (2014) states that such a code would have taken a large team of experts about six months to build the worm. Allen and Chan (2017) further this point by stating in their literature that Stuxnet would have required resources and capabilities that only military or intelligence agencies would have. Solis (2014) finally states what would have been the implications if this was a kinetic attack. This shows the danger of cyberattacks. As stated by Allen and Chan, this type of attack could be used, in principle, to damage much civilian and military infrastructure (2017).

### 3.4.7    Discussion on cyberwarfare

The aim of this section was to understand the cyber domain of warfare and look at the most prevalent examples of cyberwarfare. This chapter began with looking at what exactly the cyber domain is. This is seen as globally or locally connected computing devices, not just the internet.

Furthermore, it is a decentralised area owing to the fact that no state has control over its resources, and it goes beyond borders (Ween *et al*, 2017). This is ultimately a global domain that uses the electromagnetic spectrum and electronics to create, change, modify and store information (Robinson *et al*, 2015). Cyberwarfare was then defined as the US of this domain in order to interfere with physical, electronic, and online systems (Ween *et al*, 2017). When it comes to a cyberattack, it is hard to define when it becomes an act of war, but it is usually equated to an act of war when it equals a kinetic attack. A kinetic attack is when there is a use of force. Robinson *et al* (2015) state that a cyberattack is conducted with warfare-like intent. The literature also showed the difference between a cyberattack and a cyber intrusion which is defined based on intent. This section finally looked at different actors involved in cyberattacks such as cyberterrorists, cyberspies and cyberwarriors. The next part of this chapter looked at major cyberattacks found in the literature. The attacks that were looked at are Stuxnet, Georgia and Estonia, The US, and Israel preparing the battlefield. The next section will look at AI.

## 3.5    Artificial Intelligence

AI is integral to this research, as the aim of this research is to see how AI-enabled AWS will affect how states conduct their nuclear deterrence. This aims to look at how AWS will potentially disrupt traditional nuclear deterrence and how this will subsequently affect strategic stability. Furthermore, the aim of this section was to understand how technology affects how states go to war. As stated by Robinson *et al* (2015): "[f]rom the sword battles of the past to the unmanned drone strikes of today, this game of power is constantly driven to shift and evolve by technology". Technology has always had an impact on the battlefield and an impact on a state's national security policy, as states often ended up in 'arms races' to get an advantage technologically over their adversary. As stated by Haner and Garcia (2019), states are making heavy investments into AI-enabled autonomous systems. The global spending on AI is expected to reach between \$16–\$18 billion in 2025 (Haner & Garcia, 2019: 331). Furthermore, Rickli (2019: 91) emphasises the statement made by Russia's President, Vladimir Putin, that a state that becomes a leader in AI will become the 'ruler of the world'. This shows the importance of AI to international powers and also creates the impression that there could be a possible AI arms race. This section of this research does not aim to discuss the theory around AI in terms of what is autonomy or what is AI. It aims to look at how US policy sees AI and what are current examples of AI-enabled AWS. This section discusses prominent examples and their current capabilities. It will begin by looking at how the US DoD (2012) directive defines 'Autonomy in Weapons Systems'. It will then go on to look at different areas in the

military where AI can be applied and the possible capabilities it can bring to warfare. Finally, it will look at current AI-enabled AWS and future AWS.

### 3.5.1 Autonomy in weapon systems: Department of Defense directive

The US policy on AWS is based on the role of humans in the operating systems, over the technologies' sophistication (Artificial intelligence and national security, 2019). The DoD directive (2012) defines AWS as a system that, once activated, will then be able to engage a target without any further intervention by a human operator. The DoD directive (2012) positions that autonomous and semi-autonomous systems will be designed to allow human operators to exercise an appropriate level of human judgement when it comes to the use of force. Semi-autonomous systems, according to the DoD, will be allowed to apply lethal or non-lethal forces (DoD Directive, 2012). However, the semi-autonomous system must be designed so that it will not engage an unselected target if the system loses or has degraded communication; meaning that it will only engage authorised targets before it becomes degraded or was lost (DoD Directive, 2012). This shows that the DoD puts a specific emphasis on the role of humans when it comes to the use of force. The Artificial intelligence and national security (2019) report states that this means that they have control over the why, when, where, and how and does not mean that they need to have control over the weapon system. The DoD further states that there needs to be: "[a]dequate training, [tactics, techniques, and procedures], and doctrine are available, periodically reviewed, and used by system operators and commanders to understand the functioning, capabilities, and limitations of the system's autonomy in realistic operational conditions" (DoD directive 3000.09 as cited by Artificial intelligence and national security, 2019: 1). Furthermore, the system must have an interface that is user-friendly so operators may use the system effectively (Artificial intelligence and national security, 2019). Finally, these systems must be tested to be sure to minimise failure and they must also get senior level review when it comes to their operations (Artificial intelligence and national security, 2019).

### 3.5.2 AI capabilities across different domains

This section aims to briefly look at AI capabilities across Intelligence, Surveillance, and Reconnaissance (ISR), logistics, Cyberspace, and Command and Control. It will take a more in-depth look into Autonomous Vehicles and AWS. This will be done in order to give an overview of the different military applications for AI which have a specific emphasis on AWS as this technology is important to this research.

### 3.5.3    *Intelligence, surveillance, and reconnaissance*

One of AI's main capabilities is its ability to process large amounts of datasets and analyse them (Artificial intelligence and national security, 2019). One project highlighted by the CRS is an AI project that combines computer vision and ML that analyses this data and identifies possible hostile threats (Artificial intelligence and national security, 2019). A prevalent example of Intelligence, Surveillance, and Reconnaissance (ISR) enabled with AI is called Project Maven which was launched in April 2017. Tarraf *et al* (2019) refer to Project Maven as an Algorithmic Warfare Cross-Functional Team which sits between AI and operational AI known as mission support AI applications.  Project Maven has received a lot of attention in the media as Google employees protested Google's involvement in the project (Shane & Wakabayashi, 2018). According to a CRS report about a dozen Google employees resigned and about 4 000 signed a petition against the company's involvement in the project (U.S ground forces robotics and autonomous systems (RAS) and artificial intelligence (AI): consideration for congress, 2018). Google was one of many companies involved in the DoD contract; however, in 2019 it stated it will not renew its contract that was set to expire at the end of the year (U.S ground forces robotics and autonomous systems (RAS) and artificial intelligence (AI): consideration for congress, 2018).

Shane and Wakabayashi (2018) states that the Pentagon program aims to use AI to interpret video imagery and could be used to improve drone strikes. Air Force Lt Gen Jack Shanahan states that project Marven is an initiative to use drone footage combined with ML to create 'useful intelligence' (Air Force Lt Gen Jack Shanahan as cited by Corrigan, 2017). The program aimed to use AI to autonomously identify objects of interest from either moving or still UAV imagery (U.S ground forces robotics and autonomous systems (RAS) and artificial intelligence (AI): consideration for congress, 2018). Project Maven aimed to develop computer vision that would be trained by ML techniques to better identify objects (U.S ground forces robotics and autonomous systems (RAS) and artificial intelligence (AI): consideration for congress, 2018). The ultimate aim of the project was to make analysts' work easier by getting rid of labour-intensive drone footage analysis and allowing human analysts to process two to three times more data (U.S ground forces robotics and autonomous systems (RAS) and artificial intelligence (AI): consideration for congress, 2018). This would allow for the delivery of more time-sensitive data and lower collateral damage and civilian casualties (U.S ground forces robotics and autonomous systems (RAS) and artificial intelligence (AI): consideration

for congress, 2018). The current contract for project Maven was taken over by Palantir and due to the project being classified, there is limited information.

### 3.5.4    Logistic

An AI company called 'SparkCogntion' installed AI into several of Boeing's commercial airlines in order to analyse when they would need repairs. An example of its success is how it flagged an engine that needed repair far ahead of its scheduled repair (Hoadley, & Lucas, 2018: 9). When the aircraft was inspected the maintenance team found that one of the blades had been nicked. Had this not been discovered it would have cost Boeing $50 million to replace it (Hoadley, & Lucas, 2018 : 9). Another example is how the army signed a contract with IBM to create a AI-proof concept. Watson was developed and the US plans to use it analyse ships due for repairs that could save a $100 million annually (Hoadley, & Lucas, 2018). Furthermore, AI algorithms that are able to able to manage distribution and transportation tasks as well as prioritise them will change logistics (Schütz & Stanley-Lockman, 2017). Schütz and Stanley-Lockman (2017) state that automation is not new to the military and was used by the US in the early 1990s and in the wake of Operations Desert Shield and Desert Storm. Schütz and Stanley-Lockman (2017) refer to this program as a glorified Excel spreadsheet; however, spreadsheets transformed logistics. This system has a number of limitations, which will not be discussed here.

The progression in AI, big data and deep learning can create what Schütz and Stanley-Lockman (2017) call a 'nervous system' for the US military. This nervous system would combine advance sensor data and blockchain technology for decentralised digital ledgers (Schütz & Stanley-Lockman, 2017).  This nervous system for military logistics will allow for a central location that sends signals when different inspections are needed and which parts need repair (Schütz & Stanley-Lockman, 2017: 3). This will allow units to operate more efficiently owing to data links between disaggregated forces and systems (Schütz & Stanley-Lockman, 2017: 3). Other apps created have allowed for four hours in a day to be freed up for airmen that had to previously plan air-to-air refuelling by hand (Schütz & Stanley-Lockman, 2017). AI could allow for data to be sent to the correct area when it comes to decision-making or even as something as simple as a unit needing more fuel (Schütz & Stanley-Lockman, 2017). It may also allow for the enhancement of connection between different branches of service or allies (Schütz & Stanley-Lockman, 2017). More complete information may allow states to be more ready for war; it may also allow for 'lessons' learned in real time that can get rid of uncertainty

(Schütz & Stanley-Lockman, 2017). This section shows a brief outline of the possible effect that AI will have on military logistics. AI will create more efficient, streamlined, and cost-effective military operations on the future battlefield.

### 3.5.5 Cyberspace

The CRS quotes Admiral Michael Rogers, who stated that relying purely on human intelligence is a 'losing strategy' (Artificial intelligence and national security, 2018). The benefit of AI cyber defence tools is that they are trained to recognise a change in patterns of behaviour in a network which allows them to detect irregularities which then allows for a more comprehensive barrier compared to previous unobserved attack methods (Hoadley, & Lucas, 2018: 10). Hackers need to modify their malicious code in order to circumvent a computer defence, this quote shows that AI can be more comprehensive against this (Hoadley, & Lucas, 2018: 10). DARPA held a Cyber Grand Challenge in which seven computers with custom-designed software had real-world vulnerabilities and glitches (Fraze, 2016). Contestants had to develop AI algorithms to identify and patch these problems while attacking other teams' weaknesses (Hoadley, & Lucas, 2018). The team behind the Mayhem system that won the Cyber Grand Challenge was a group of Pittsburgh-based researchers (Fraze, 2016). According to Fraze (2016), the aim of the challenge was to increase the development of autonomous systems that have the capability to detect, evaluate and patch systems before adversaries can exploit them. The systems that were used in the event were able to find and patch within seconds and not months as usual (Fraze, 2016). The result was that bugs were fixed in seconds, quicker than humans could, and showed AI's ability to play defensively and offensively at the same time (Hoadley, & Lucas, 2018). Events like Cyber Grand Challenge are of importance to the US forces as the cyber domain becomes an increasingly more important hostile and an important area of conflict. AI may potentially allow for the US to defend its cyberspace more effectively due to AI speed. The Mayhem system was able to find, evaluate and fix patches faster than humans ever could. Such speed may allow the US cyber domain to defend itself continuously against evolving threats. AI may also offer an aggressive offence tool that is capable of finding ways through a state's cyber defence faster than they can fix it.

### 3.5.6 Command and control

According to the Congressional Research Service, the US air force is developing a 'Multi-Domain Command and Control' system (MDC2) (Hoadley, & Lucas, 2018: 10). The aim of the MDC2 is to centralise cyberspace, space, sea, land, and air-based control systems by using

AI to create this single source (Hoadley, & Lucas, 2018: 10-11). This will create a 'common operating picture' for the US military (Artificial intelligence and national security, 2019: 10-11). This AI will be able to find issues on communication, real-time analysis and different viable courses of action at a faster response rate; analysts believe this will improve the wartime decision-making process (Hoadley, & Lucas, 2018). This data fusion capability is in the development phase with Lockheed Martin, Harris and other AI start-ups (Hoadley, & Lucas, 2018). An example currently in development by Lockheed Martin is the 'Joint All-Domain Operations (JADO)' (Kahn, & Thatcher, 2020). Kahn and Thatcher (2020) state that the JADO capability will combine electromagnetic and physical common operational pictures. JADO will also use AI and other cognitive applications to identify War Reserve Mode (WARM) emissions, better optimise IS sensor collection, and finally autonomously update aircraft routes based on current threats (Kahn, & Thatcher, 2020: 1). Kahn and Thatcher (2020) state that the JADO Full Spectrum Operations is the key to fighting and winning battles in Highly Contested Environments (HCE). This system aims to make the right decision faster. Another example comes from the CRS (Joint all-domain command and control (JADC2), 2020) and is called the Joint All Domain Command and Control (JADC2). The aim of this system is to connect all the sensors from each of the military services into one network (Joint all-domain command and control (JADC2), 2020). The aim of JADC2 is to enable commanders to make better decisions (Joint all-domain command and control (JADC2), 2020: 1). This will be done by collecting data from lots of sensors, processing it with AI to identify targets, and then recommending the best weapon systems; both kinetic and non-kinetic (Joint all-domain command and control (JADC2), 2020: 1). This shows how AI-enabled C2 aims to bring all data from each service of the military into a singular point to allow for faster and better decisions to be made.

### 3.5.7    *Autonomous vehicles and weapon systems*

Hoadley and Lucas (2018) state that applications in this field are similar to commercial self-driving vehicles, which use sensors to collect data in order perceive the environment and execute decisions. This shows the dual-use of AI for military and civilian applications. This section will analyse autonomous vehicles and weapon systems in order to give a basic understanding of them. This is owing to the fact that the next chapter will be giving an in-depth look at AWS and this chapter aims to give a contextualisation of technology and warfare. An example of a military application would the US Air Forces 'Loyal Wingman' program (Hoadley, & Lucas, 2018). This involved pairing an unmanned fighter (F-16 was used) with a F-35 or F-22 (Artificial intelligence and national security, 2018). The F-16 test platform was

able to respond to events that had not been programmed into it, meaning that it reacted autonomously to unforeseen events (Hoadley, & Lucas, 2018: 11). The events that it reacted to were unforeseen obstacles and weather conditions (Hoadley, & Lucas, 2018: 11). This could be made more helpful by adding extra weapons to the unmanned systems or the ability to jam electronic threats (Hoadley, & Lucas, 2018). The Marine Corps tested a 'Multi-Utility Tactical Transport'(MUTT) which is an ATV-size vehicle which follows marines around the battlefield by radio link. It is not AI-enabled yet, but the Army aims to enable it with AI (Hoadley, & Lucas, 2018). Furthermore, the Navy tested a swarm technology, which is AI cooperative behaviour, for defending harbours or certain areas (Hoadley, & Lucas, 2018). Furthermore, multiple companies such as Boeing, Lockheed Martin and Northrop Grumman are all developing AI-enabled AWS that operate in air, land and sea.

### 3.5.8    Final discussion of AI

The aim of this final section was to look at AI in order to make a contextualisation of this transformative technology. This chapter stated at the beginning that it did not aim to discuss AI theory as it aims to look at current US policy on AI and current examples of it. This section began with a look at the role of humans in AWS according to the US's DoD. The DoD puts a heavy emphasis on maintaining a human-in-the-loop when it comes to AWS. The human-the-loop emphasis by the US will be dealt with in-depth in the next section. Furthermore, the DoD stated that it wants thorough reviews and training on these systems in order to make sure there are no possible errors. What followed this was a look at AI across different domains of the US military. The different domains that were looked at were: ISR, logistics, cyberspace, command-and-control, and autonomous vehicles and weapon systems. Each of these sections gave a brief overview of each of these domains and their current capabilities. What can be seen in all of these domains is that AI offers a high potential when it comes to increasing each section's capabilities.

## 3.6    Conclusion

The aim of this chapter was to look at how technological innovation changes the way states, more specifically the US, conduct warfare. The premise of this chapter was that technology changes the way states conduct warfare. After a look at four different technologies, it can be conclusively said that technology does influence the way in which states conduct warfare. Furthermore, these technologies provide problems outside of the military sphere due to their dual-use capability. This chapter also allowed for a more in-depth look into nuclear deterrence

and what a nuclear triad is. The nuclear triad consists of three legs: land, air and sea. The land consists of ICBMs that allow states to respond quickly with accuracy and strength. The air leg of the triad includes a range of nuclear capable bombers such as the B-52, which allows states to send out these bombers as a deterrent and if the crisis is resolved, to recall them. This gave the US capability as well as flexibility when it came to the nuclear triad. The final leg of the triad is the sea leg. This consists of SSBNs. The ability for SSBNs to disappear under water and travel long distances gave the triad legitimacy and secured the US's second-strike capability, due to the ability of the SSBN to hide underneath the ocean's surface. The final section discussed AI, with the aim of dealing with what the US policy stated and how the US viewed their capabilities. It also provided this research with an understanding of different companies leading the way in autonomous weapons systems. This section has created a contextualisation and background for this study in understanding the connection technology has with the way states conduct warfare.

# Chapter 4

## 4.1    Introduction: The rise of Autonomous Weapon Systems and its effect on the US's nuclear deterrence

This thesis goal is to analyse how states go to war and how AI will affect this. Owing to the fact that these are two very broad areas of study, this research aims to look at a specific aspect of both AI and modern warfare. The aspects that were chosen to be analysed were those of AWS and nuclear deterrence. The main aim of this research is to study how AWS would lead to a potential disruption in nuclear deterrence and subsequently affect strategic stability. The case site that was chosen was the United States of America (US), owing to the fact that the US prides itself in its technological innovation and military dominance. This will become abundantly clear as this chapter proceeds. Furthermore, to create a theoretical framework in which to conduct this study, Kaldor's new war thesis was reviewed in Chapter 2 with the aim of using it to create a framework for this study. The new war thesis will therefore be used and presented in this chapter in three different sections: actors, goals, modes of warfare and forms of finance. Each section of the new war thesis plays a vital role in understanding modern warfare. Some of these sections will not be as intensive as others. For example, this chapter will be primarily dominated by modes of warfare. This is owing to the fact that this section will answer the main and secondary research questions. Meanwhile 'goals' and 'actors', which entail populism and President Trump, will create a contextualization and background for this study.  Forms of finance [9]does not play a vital role in this research because the funding of AWS comes from the US government. This will be discussed further in the discussion of limitations of this research in the next chapter.

---

[9] The 2021 US budget aims to deliver on President Trump's promise to rebuild the US's military, strengthen its readiness, and support the US personnel (Office of Management and Budget 2020). The budget aims to implement the 2018 National Defense strategy which aims to invest in modernisation, innovation and lethality that will allow the US to face its current and future challenges (Office of Management and Budget, 2020). The budget that was requested was $705.4 billion, which represents only an increase of $0.8 billion in 2021 over 2020 (Office of Management and Budget, 2020: 33). The budget for the DoD science and technology programs amounts to over $14 billion (Office of Management and Budget, 2020: 35). This shows the US's commitment to maintaining its technological dominance in its military capabilities and its heavy investment in AI and autonomy.

## 4.2    Actors: Who are the actors involved?

The actors that are involved in the development of AWS need to be outlined to create a better contextualization for how AWS will affect nuclear deterrence. This is because nuclear deterrence is based on perception, which makes understanding the actors involved important. Kaldor's (2013) new war thesis maintains that new wars are fought by actors such as private security contractors, non-state actors, regular armed forces, and warlords. It is important to understand the actors involved in the development of AWS due to the complexity of actors in modern warfare. During the Cold War the global arena was dominated by the US and the USSR. However, the modern arena is dominated by a number of powerful nuclear armed states, such as Russia and China. Furthermore, the development of AI is diverse, in that AI is developed in both the US government and the private sector, with a crossover between the two spheres. Both of these sectors are actors involved in AWS, but they are not strategic competitors, they are involved in innovation together. This section will not focus on the relationship between the US government and the private sector but on strategic adversaries to the United States. It is important to understand the actors involved in AWS owing to the issue of perception of nuclear deterrence but it is more important to understand how the US perceives these actors. Understanding who the US perceives as their strategic advisory will allow this study to create a better background and more complete secondary contextualization. Finally, this section will not talk about the US as an actor as this will be discussed in the 'goals' section.

### 4.2.1    The US's 2017 National Security Strategy: What can this tell us about the US's perceptions of its adversaries?

The US's National Security Strategy (NSS) (2017) aims to ensure peace through strength (National Security Strategy. The US aims to renew its competitive advantage to maintain its strength so as to successfully deter and if necessary, defeat any aggressor that is going against US interests (NSS, 2017). The NSS views the US as having taken a break from maintaining its military dominance and believes that Russia and China have taken this opportunity to change the international order in their favour. The US NSS views China and Russia as having the aim of eroding US influence, its interest, and to erode its security and prosperity. Furthermore, the NSS (2017: 8) states that both China and Russia are developing advanced weapons and capabilities that could become a possible threat to US critical infrastructure and command and control infrastructure. Other actors that the US perceives as threats outlined by the NSS (2017) are North Korea and Iran. There are currently heighted tensions between the US and Iran after

a series of events that unfolded in 2020. On 3 January 2020 the US conducted an air raid aimed at killing top Iran General Qassem Soleimani (Iran's Qassem Soleimani killed in US air raid at Baghdad airport, 2020). The subsequent killing of General Soleimani lead to a retaliation by Iran with it firing more than a dozen ballistic missiles at a base in Iraq that was hosting US troops (US-Iran tension : how confrontation between rival escalated, 2020). Furthermore, both North Korea and Iran are points of tension for the US regarding their nuclear weapons proliferation programmes.

### 4.2.2    China: A strategic rival and threat

The US's NSS (2017: 21) further views the Chinese government as problematic due to its alleged theft of US intellectual property. They view this theft as being in the hundreds of billions of dollars and giving China an unfair ability to tap into US innovation. The NSS (2017: 25) attributes part of China's military modernisation coming from access to US innovation economy, with an end emphasis on US universities. Furthermore, from key informant interviews, it was found that China is a problem for states when it comes to industrial espionage and cyber intrusions in order to steal information (key informant 1, key informant 2, & key informant 3, 2020). Key informant 3 (2020) went on to state that the Chinese specifically target the military sector and personnel. The US aims to deal with China through the use of the four pillars outlined in its NSS: "(1) protect the American people, homeland, and way of life; (2) promote American prosperity; (3) preserve peace through strength; and (4) advance American influence." (United States Strategic approach, 2020: 1). This will be done in order, as already stated, to protect US interests. More specifically, China's growing military is a significant threat to both the US and its allies (United States Strategic approach, 2020: 7). Another issue outlined is the fact that China has a military-civilian fusion (MCF) strategy, which allows the Chinese military to have unfettered access to civilian entities that develop and gain advanced technologies (United States Strategic approach, 2020: 7). The US and other foreign countries are feeding the People's Republic of China (PRC) military research and development programs dual-use technologies owing to the PRC's non-transparent MCF linkages, this is in turn strengthening the Chinese Communist Party's ability to suppress domestic opposition and threaten the US and its allies (United States Strategic approach, 2020: 7). This is problematic as the US is very aware of the fact that China is pursuing greater military capability and technological innovations to grow its influence as a global power.

Finally, China is commonly referenced throughout the literature when it comes to AWS. As already stated in this thesis, China aims to become the leader of AI, which makes it an important actor to analyse. Furthermore, theorists such as Johnson (2020a, 2020b), highlight that there is a specific tension about AI when it comes to the US and China. Geist and Lohn (2018) state that China and Russia both believe that the US is leveraging AI to undermine the survivability of their nuclear forces. Furthermore, China is a good example of an actor as their military innovations have led to multiple successful tests. This shows that China is a significant threat and issue to the US. It is also seen as a valuable example of an actor involved in the pursuit of AI and AWS. This rising tension between the US and China combined with a possible AI arms race could be potentially destabilising. This statement is supported by Johnson (2019), who states that a fast-emerging US-China AI innovation race will have a potential and profound destabilising effect on strategic stability of the future. How the NSS and theorists view China shows that it is an important actor. This is owing to the fact that it is an incredibly threatening adversary to US interest and is in the pursuit of military and AI dominance. This is threatening to the US as it views China's pursuit of military and AI capabilities as potentially disruptive to its influence, security and prosperity. It is important to understand how the US views China, as one of the major tenets of nuclear deterrence is perception. The US only needs to perceive China as having an AWS capability that could threaten its second-strike capability in order for there to be a destabilising effect on strategic stability.

### 4.2.3    *Russia: A renewed threat*

The NSS (2017: 25) states that Russia aims to restore its status as a great power and create a sphere of influence near its borders. The NSS (2017: 5) further states that Russia aims to weaken the US's influence around the world and divide the US from its allies and partners. According the NSS (2017: 25-26), Russia is investing in military capabilities, which include nuclear systems, and cyber capabilities which have proven to be destabilising. Furthermore, the NSS (2017: 25-26) sees Russia's nuclear capabilities as one of its biggest threats. The NSS (2017: 26) states Russia's ambition and its increase in military capabilities are a risk to the Eurasia area, especially with the chance of a Russian miscalculation that could lead to conflict. Furthermore, the NSS (2017: 47) claims that Russia is using subversive measures to weaken the US's credibility in Europe, weakening NATO, and weakening European institutions and governments. Russia has also invaded Georgia and the Ukraine, showing their lack of commitment to other states' sovereignty in Europe (NSS, 2017: 47). Russia uses nuclear posturing and deploying of offensive capabilities to intimidate its neighbours (NSS, 2017: 47).

The Bureau of European and Eurasian Affairs within the US Department of State sees Russian foreign policy as aggressive and attributes it to domestic political issues (U.S. Relations with Russia, 2020). The United States aims to deter Russia and only create a level of cooperation when it suits US interests.

> "The United States has sought to deter Russian aggression through the projection of strength and unity with U.S. allies and partners, and by building resilience and reducing vulnerability among allies and partners facing Russian pressure and coercion. The United States would like to move beyond the current low level of trust with Russia, stabilize our relationship, and cooperate where possible and when it is in the core U.S. national security interest to do so. To achieve this, Russia must take demonstrable steps to show it is willing to be a responsible global actor, starting with a cessation of efforts to interfere in democratic processes. The long-term goal of the United States is to see Russia become a constructive stakeholder in the global community." (U.S. Relations with Russia, 2020).

This above quote from The Bureau of European and Eurasian Affairs gives a good insight into the US's perception of Russia and how they subsequently aim to deal with Russia. However, this is not the only issue for the US. Like China, Russia aims to become a global leader in AI. Laird (2020) states that AWS will likely raise the risk of crisis instability and conflict escalation in future confrontations between the US and Russia. Russia already has aims of creating an autonomous unmanned underwater vehicle (UUV) called Status-6, which they believe will be able to secure their second-strike capability. This shows that Russia isn't just an adversary with an aggressive foreign policy, but an adversary interested in pursuing AI in order to gain a technological advantage over the US.

### 4.2.4    Discussion

The US's NSS sees the world as increasingly competitive, which is the foundation for their response, to increase their competitive edge that they deem has slipped (2017). More specifically, the US views their military as the strongest in the world, but see other states as increasing in strength and shrinking the US's strength. Lots of actors now have the capability to field a broad arsenal of advanced missiles, some of which have the capability to reach the US homeland, which allows weak states to be empowered and emboldens them (NSS, 2017: 3). Even without these 'otherwise weak states', the US faces two hugely powerful nuclear-armed states, Russia and China. The US faces the issue of two powerful threats that come during a time of increased costs of military capabilities that is shrinking the size of its military

forces. This is an important aspect of AWS that will be discussed further in the methods section, along with the actors and the pursuit of AWS as these play a vital role in this research, are the two most cited threats by the NSS and are found frequently in the literature. This section now gives this thesis a background to look at the goals of the United States, with a specific look at President Trump, as the NSS does not hold much power if President Trump decides to take up another strategy during a time of crisis or escalation.

## 4.3    Goals: Making America Great Again and reasserting US global leadership

In Kaldor's new war thesis there is a link between identity politics and conflict. A part of this research is trying to create a contextualization on how prominent movements like populism, a form of identity politics, can help to understand how the US pursues its NSS and ultimately better situate the reader. Wars used to be fought on ideologies such as communism or democracy; new wars are fought on identity (Kaldor, 2013). Identity politics is the source of violence, specifically a unidimensional personality (Kaldor, 2013). The more specific example of identity politics that this thesis is studying is 'populism'. Populism can be seen as anti-elitist and usually pushes an agenda of what a 'true citizen' looks like (Marchlewska *et al*, 2017: 151). Fukuyama (2017) also highlights how populist leaders create an image of the national citizen and denounce the elite. Lofflmann (2019) highlights how President Trump's 'brand' of populism has targeted the establishment as the culprits of economic and political failures, building an antagonistic relationship between the elite and 'ordinary Americans', and exploiting Americans' emotional triggers. In Löfflman's (2019) research the characteristics of the 'true citizen' and anti-elite can be seen, points previously highlighted by Marchlewska *et al* (2017) and Fukuyama (2017).

> "Popular discontent with the status quo opened the space for Trump's populist messaging, which, alongside contempt for the Washington establishment, promised national revival and renewal through economic protectionism, aggressive deregulation, strict anti-immigration measures and a transactional focus on prioritising US interests in international affairs that would 'make America great again'. On specific issues of foreign policy and national security, this rhetoric addressed a longstanding gap between public opinion and the attitudes of a bipartisan elite on American global engagement, from military intervention to free trade" (Löfflman, 2019: 118).

Löfflman (2019) states that President Trump's message became so powerful because he was able to exploit the disconnect created between the elite and public opinion. While the elite wanted global economic liberalisation, the public wanted to focus on domestic issues and leave other states to sort out their own affairs (Löfflman, 2019). However, Löfflman (2019) believes that Trump's national-populist rhetoric has not fully influenced the US's NSS; it has, however, opened a debate and allowed the US to recalibrate its foreign policy. However, according to Löfflman (2019: 130), Trump has opened the door for a potential recalibration of US grand strategy by questioning the political dominance of the foreign policy establishment and its strategic standpoint. The recalibration could potentially move towards a closer alignment between public and elite opinions, moving away from bipartisan consensus on liberal hegemony (Löfflman, 2019: 130). Löfflman (2019: 130) argues that this is driven neither by post-Cold War primacy nor by populism.

The studies by Löfflman (2019), Hall (2017) and Feaver (2017) state that President Trump's rhetoric and the NSS are inherently different from one another. Löfflman (2019) argues that Trump's rhetoric is anti-globalism and anti-elite; however, Trump's foreign policy is different from his rhetoric. Löfflman (2019) substantiates this claim by highlighting how the US still supported NATO, approved congresses sanctions against Russia, and secured Japan and South Korea by increased military activity in the Asia-pacific (Löfflman, 2019). Cordesman (2017) would agree with Lofflman's (2019) argument as he argues that Trump's 'America First' means internationalism and not isolationism. This shows that Trump may move from position to position in his short statements or tweets, however, he has ended up closer to the centre when it comes to his national security positions than many critics fully account for (Cordesman, 2017). This shows that Trump as an actor differs from his foreign policy. In simple terms, what Trump says on twitter or at a rally and what his policy is differ substantially. Löfflman (2019: 125) argues that the main impact of Trump's nationalist-populist agenda has been a disdainful perception of the US's partners and allies along with a more self-centred view of world affairs. Trump's populist language can ultimately be seen as anti-elite view of world affairs combined with a belief that the US should take less of a hegemonic role. However, whether the NSS (2017) shows populist sentiments or has a more internationalist view, President Trump does not have to abide by it during a crisis. This is where it becomes important to understand Trump as an actor for potential aggressors.

This argument is found in the research done by Pifer (2018). He states that when it comes to individual problems or crises, the NSS does not necessarily dictate a strategy or a plan that

must be implemented to deal with these situations. Pifer (2018) states that: "I do not recall any meeting in which a US. official argued that we had to adopt a certain policy course because of the prescriptions contained in the National Security Strategy." These policies, according to Pifer (2018), can be overturned, which is an issue to both the US's adversaries as well as its allies owing to the fact that Trump views unpredictability as an asset (Pifer, 2018). For Pifer (2018) the biggest question about how Trump will act is based around his instinct and volatility. Key informant 3 (2020) stated that understanding the actor is of significant importance, as an actor like Trump will not blink in a fight and there would be a strong response from his administration. Meanwhile, a Biden administration may prefer a more watered-down diplomatic response when it comes to a moment of crisis (Key informant 3, 2020). The argument of this section is that Trump as an actor sees his unpredictability and volatility as an asset. Furthermore, understanding Trump's foreign policy may not be a reliable source in understanding how he would act in a moment of crisis or escalation. This type of strategy may prove viable during a time of crisis, as an actor may not choose to escalate as Trump's strategy during a moment of crisis cannot be predicted by either his actions or his policy.

## 4.4    Methods of warfare: How AWS will affect the US's nuclear deterrence and subsequently strategic stability

This section looks at methods of warfare which is the next tenet of the new war thesis. Kaldor's (2013) methods of warfare are based on the premise that warfare has changed. This means that war is no longer fought in a decisive battle between two well-armed states (Kaldor, 2013). A decisive encounter between two well-armed states was regularly referred to as 'old wars' in the literature. This thesis investigates the effect of AWS on nuclear deterrence, so the specific method of warfare that was looked at was nuclear deterrence, which could subsequently affect strategic stability. This thesis goal, among other things, to analyse how warfare is changing. However, this thesis cannot cover all of the numerous developments in modern warfare technology. A good example is the US's Hypersonic Boost-Glide Weapons which aim to create a level of manoeuvrability that will inhibit enemies from tracking its trajectory, which is possible with ballistic re-entry vehicles (Johnson, 2020a, Johnson, 2020b). This section will begin with a look at the US's current thinking on nuclear deterrence under the Trump administration. It is important to understand the US's thinking in order to understand its perception of factors such as the rise of China and its modernisation of its nuclear forces and Russia's nuclear modernisation. This will help to understand why the US wants to modernise its nuclear force and give a understanding of its perceptions of other superpowers. Finally, it

will give this thesis a better understanding of the US's current nuclear thinking and set this research up to look at AWS. It will also allow for a better understanding of US's current nuclear strategy which is important to this research as nuclear deterrence is its 'mode of warfare'. The section that will follow nuclear deterrence will be AWS. There are several tenets that will be discussed under this section and it will subsequently dominate most of this chapter. This section will look at conventional military applications such as drones, swarm technology, the issue of speed, current capabilities regarding sensors, the issue of AWS being nuclear armed, the issue of miscalculation and possible de-escalating nature of AWS, human in the loop, and the issue of perception. Ultimately, this section goal is to analyse data that will help to answer the main and secondary research questions.

### 4.4.1    Nuclear weapons: The cornerstone of US National Security Strategy

This section on modes of warfare now takes a look at the US's nuclear strategy. This is of importance because the chosen mode of warfare in this research is nuclear deterrence. The NSS (2017: 30) notes that nuclear weapons have served as a vital part of its national security strategy for the past 70 years. Nuclear weapons have created a foundation in order to deter aggressors and preserve peace and stability for the US and more than thirty of its allies (NSS, 2017: 30). However, the NSS (2017) states that the US's nuclear Triad is ageing while its adversaries have expanded their arsenal and delivery systems. This is an issue for nuclear deterrence as the Triad plays a vital role in the US's nuclear strategy. The NSS (2017) mandates that the US must maintain a credible nuclear deterrence and an assurance of its capabilities. The NSS (2017) aims for this by increasing investment in maintaining its nuclear arsenal and infrastructure. The NSS (2017: 31) states that the US does not need to match the nuclear arsenals of other powers, but they must still maintain their stockpiles so that they can deter adversaries, assure their allies and partners, and achieve objectives if deterrence fails. Furthermore, the NSS aims to invest in its nuclear enterprise in order to keep an effective and safe nuclear triad that is capable of responding to future security threats.

The need for the US to modernise the Triad is also seen in the 2018 NPR. According to Vergun (2019), the US highlighted in the NPR that it needed to modernise its nuclear triad. Vergun goes on to quote Hyten, who states that each leg of the Triad is of critical value when maintaining an effective deterrence. Hyten goes on to highlight the importance of each leg of the Triad in creating an effective nuclear deterrence (as cited by Vergun, 2019). The strategic bombers are the most recallable of the legs as the president can call them back once deployed

in moments of crisis (Hyten as quoted by Vergun, 2019). However, Submarines are the most survivable element of the nuclear tried (Hyten as quoted by Vergun, 2019). Submarines gave the triad the ability to hide from adversaries, it gave the US the ability to secure their second-strike capability; a critical tenet of nuclear deterrence and key to MAD (Hyten as quoted by Vergun, 2019). ICBMs are the most problematic for adversaries as there are more than 400 locations across the US (Hyten as quoted by Vergun, 2019). Finally, Vergun (2019) highlights that modernisation is important and does not mean a new class of nuclear missiles; it refers to improving the existing triad. This shows that the US still values its nuclear deterrence as seen by its desire to modernise its nuclear triad. Vergun also shows the importance of each leg of the triad when it comes to the US's NPR.

This need to modernise the US's nuclear forces is not a new strategy under the Trump administration. The NPR (2018), according to Rose (2018), continues the modernisation program from the Obama administration. According to Rose (2018: 3) the NPR states that the US must move forward with the Obama administration's strategic modernisation programme. This involves modernising the Columbia-class SSBNs, the Ground-Based Strategic Deterrent (GBSD), the B-21 bomber, and the Long-Range Stand-Off (LRSO) cruise missile (Rose, 2018: 3). Roses (2018) states that this modernisation program continues to have bi-partisan support by Congress. These systems modernisations are of importance to the US as they enhance strategic stability and are in line with arms control obligations and commitments (Rose, 2018: 3). This shows that the US believes strongly in the importance of the nuclear triad when it comes to maintaining an effective nuclear deterrence.

While it is valuable to understand what theorists have to say about the NPR, it is also valuable to analyse the policy directly. The NPR (2018) affirms the issue of the US continuing to commit to a reduction of nuclear weapons while adversaries like China and Russia do not. Furthermore, the NPR brings up the issue of both Iran and North Korea, which further highlights the issue of a multipolar world and subsequently many threats for the US. The NPR (2018) goes on to further highlight numerous other threats such as chemical, biological weapons and cyber. Cyber threats are of interest to this research as it was a recurring theme found in the primary research that was conducted; it will, however, be discussed later. Regarding China and Russia, the NPR (2018) continues to maintain the narrative of a competitive global arena that entails adversaries that the US has allowed to gain ground militarily. This is what the actors' section of the thesis found to be called 'principled realism' a version of Trump's populism mixed with the traditional foreign policy of the US.

The US also highly values their nuclear capabilities as the NPR (2018) attributes this to successful deterrence of both nuclear and non-nuclear aggression. Furthermore, the US's nuclear deterrence also ensures allies and partners achieve the US's objective if it fails and allows the US to hedge against an uncertain future (NPR, 2018). The NPR (2018) goes on to further state that potential aggressors must not miscalculate the use of nuclear weapons. As mentioned already in this thesis, perception plays a vital role in nuclear deterrence. Adversaries must understand a defender's capability and will to use nuclear weapons in a crisis.

> "Potential adversaries must recognize that across the emerging range of threats and contexts: 1) the United States is able to identify them and hold them accountable for acts of aggression, including new forms of aggression; 2) we will defeat non-nuclear strategic attacks; and, 3) any nuclear escalation will fail to achieve their objectives, and will instead result in unacceptable consequences for them." (NPR, 2018: VII).

The NPR (2018) states that this module is adaptable across a range of actors and that the US's nuclear capabilities need to remain flexible and have the ability to carry out a strike against unfavourable actions against itself and allies and partners. For this to be maintained, the NPR (2018) mandates that the US needs to modernise its nuclear forces. The NPR (2018: X) highlights that there is an increasing need for flexibility and diversity which is one of the main reasons for sustaining and replacing the nuclear triad and other non-strategic nuclear capabilities, and modernising command-and-control systems. The triad's synergy allows for the US to maintain its nuclear credibility; allowing one of the legs of the triad to be eliminated would greatly impact the US's nuclear credibility (NPR, 2018).

What this section has shown is the huge emphasis the US puts on nuclear deterrence under the Trump administration. The emphasis is on the point that the US is planning on modernising their nuclear triad in order to maintain its credibility at a time when other powers such as Russia and China are modernising their nuclear forces. Furthermore, the US believes that Russia is violating the INF and the US is currently on course not to renew the START treaty based on the belief that other powers such as Russia and China are not obeying START. All these complications and the US's emphasis on its nuclear deterrence shows the relevance of nuclear weapons and the need to understand how AWS will affect such an important tenet of US's NSS policy. Furthermore, in a time of already heighted tensions with China and Russia's, AWS could further create tensions between these global powers. It is noted that there are a number of other factors that are influencing US tensions with Russia and China, such as Russia's

involvement in US domestic politics and states such as Georgia and Ukraine. Furthermore, China's persistent industrial espionage and involvement in the South East China sea which aims to inhibit freedom of navigation and increases tension with the US and its allies. These are all noteworthy factors that will add to increased tension in global politics. However, this research cannot look at everything and will stay focused on AWS and nuclear deterrence.

### 4.4.2    Conventional application for AWS

This section on modes of warfare has covered nuclear deterrence, which is the chosen mode of warfare for this research. This part now aims to look at the technology aspect of modes of warfare and how this can affect the US's nuclear deterrence. This research is now moving from the strategy side of warfare to the technological aspect. It was important to look at the US's NSS (2017) and NPR (2018), to gain a better understanding of how the US conducts its nuclear deterrence, but it is now important to see how AWS will subsequently effect this. The above section found that the US is aiming to maintain its military dominance and its military innovation over other states. The US is pursuing a range of advanced military technologies in order to bolster its military posture (Miller, Fotaine, & Velez-Green, 2018). These weapons are expected to speed up the military conflict as well as create a level of uncertainty about these systems which will increase risk owing to miscalculations or misunderstandings (Miller *et al*, 2018). It is expected that AWS will make thousands of complex and highly complex decisions at machine speed (Laird, 2020). This speed could lead to a sudden and potent attack that could leave actors in a threatening situation as AWS opens the possibility that this speed could push an adversary to up the escalation and turn to nuclear use (Laird, 2020). Furthermore, AWS may give states the ability to pursue and reliably target an adversary's SSBN or mobile ICBMs (Miller *et al*, 2018). This would be problematic as mobile nuclear missiles are what makes a second-strike credible. Geist and Lohn (2018) also highlight how AWS has the potential to undermine MAD. Johnson (2020a & 2020b) states how AWS can be used in order to conduct ISR and monitor mobile missiles and submarines. Furthermore, AWS can be used in order to turn drones into 'swarms'. This is the largest cited area for AWS and probably the most significant due to the current low cost of drones and their proposed potential capabilities. This also allows for the US military to increase their quantity. As mentioned already, the US is facing a challenge of quantity owing to the high cost of military technology. Drones and specifically swarm technology will help the US offset this issue. This brief introduction has highlighted some of the factors that will be discussed in this section regarding AWS. The first section to be looked at will be drone and swarm technology.

### 4.4.3    Drones and swarm: The future battlefield

This section now looks at a specific example of technology that could cause a potential disruption to US nuclear deterrence; this is drone technology and its ability to create swarms. This section continues to look at how AWS as a mode of warfare could potentially affect nuclear deterrence. Johnson (2019: 150) states that future progress of AI technology will have an effect on autonomous systems and robotics that could create capabilities that will change the military balance and how states conduct warfare. Johnson (2019: 150) further states that these autonomous systems would theoretically be able to incorporate AI technology such as speech, perception and facial recognition and decision-making tools that will allow them to execute operations without human intervention or supervision. AWS offer states the ability to project power in areas known as anti-access/area-denial (A2/AD) contested zones (Johnson, 2019: 151). This is not the only issue that the US faces in (A2/AD), as the role of the information revolution has allowed enemies to follow the US and build reconnaissance strike networks that can detect US forces and strike them at long range with precision-guided weapons (Scharee, 2014: 10), showing that the US already faces significant challenges to their power projection without the introduction of AWS or swarm technology.

Furthermore, AI infused with data-analytics combined with quantum-enabled sensors could make adversaries' submarines potentially easier to locate which may force states into a 'use it or lose it situation' that will worsen strategic stability. A UAV by itself would be no threat to the US's F-35, however, a swarm of AI-augmented drones may be able to evade and overwhelm an adversary's sophisticated air defence systems (Johnson, 2020b: 20). This highlights what has predominately been seen through the literature, the rise of swarm technology. An example of swarm capability is seen this quote from Laird (2020):

> "Take for example, a hypothetical scenario set in the Baltics in the 2030 timeframe which finds NATO forces employing swarming AWS to suppress Russian air defense networks and key command and control nodes in Kaliningrad as part of a larger strategy of expelling a Russian invasion force" (Laird, 2020).

Such a move could be seen as a large tactic against Moscow and lead to nuclear escalation by Russia (Laird, 2020). Drones are seen as more effective as a swarm than by themselves when it comes to how effective their capability could be. There are currently limitations to drones' powers and sensors. Furthermore, AI-enabled drones will not be ready in this generation or

possibly the next. However, from this research came the notion that it is of importance to ask these questions and further understand AWS could affect how states conduct warfare.

Hitherto, from what has been seen in the literature, it can be stated that the rise of swarm technology is the most prominent example of AWS. This section aims to show how swarm technology could potentially have the biggest impact on nuclear deterrence which could subsequently lead to a disruption in strategic stability. Examples of swarm technology looked at in this research are known as Perdix and Sea Hunter. These drones could potentially be used to 'hunt' mobile missile launchers such as SSBN, which is highly problematic for nuclear deterrence and subsequently strategic stability. Ultimately, such a capability from uninhabited systems will allow the US forces to counter threats at increased ranges, persistence, and enable them to take a higher level of risk creating a new form of operation (Scharee, 2014). Key informant 1 (2020) states that drones have kept the US in the position it has been in for the last 25 years. If they build it first, it could be hugely beneficial for the US and maybe the allied nations (Key informant 1, 2020). However, If the US is not able to gain this capability it will take the shine off the US (Key informant 1, 2020). Key informant 1 (2020) states that the size and cost of the drones will be significantly less. Whoever does this would have a tremendous advantage over the US and would subsequently destabilise the US's role (Key informant 1, 2020).

In order to understand how swarm will affect nuclear deterrence a deeper understanding of the potential capability of swarm technology needs to be ensured. What exactly do theorist define a swarm as? Scharee (2014: 10) states that: "[n]etworked, cooperative autonomous system will be capable of true swarming – cooperative behaviour among distributed elements that gives rise to a coherent, intelligent whole." Scharee (2015) gives an example by comparing swarming in animals and robotic swarm technology, the argument being that relatively unintelligent animals can create intelligent collective behaviour when swarming together. Kallenborn and Bleek (2018: 526) state a similar view of swarms: : "[d]rone swarms consist of multiple unmanned platforms and/or weapons deployed to accomplish a shared objective, the platforms and/or weapons autonomously altering their behaviour based on communication with one another." Furthermore, the concept that drones in a swarm technology allow for more complex behaviour than a singular drone is stated by Kallebron and Bleek (2018: 526).

> "Uninhabited systems offer an alternative model, with the potential to disaggregate expensive multi-mission systems into a larger number of smaller, lower cost distributed

platforms. Because they can take greater risk and therefore be made low-cost and attritable – or willing to accept some attrition – uninhabited systems can be built in large numbers. Combined with mission-level autonomy and multi-vehicle control, large numbers of low-cost attritable robotics can be controlled en masse by a relatively small number of human controllers" (Scharre, 2015).

What can be seen so far from this introduction to swarm technology is how the theorists view it as intelligent behaviour that is created by a large number of drones that would not be possible with a singular drone (Scharre, 2015; & Kallebron, & Bleek, 2018). This cooperative behaviour comes from the drones' ability to communicate with one another and also allows for a large number of drones to be controlled by a smaller number of humans. Furthermore, a commonly stated asset for drone technology is the potential cost, as stated earlier, the US is facing the issue of increased cost of military technology that is lowering the size of their forces. The costs of drones may allow the US to offset this as they are inexpensive and can allow for mass on future battlefields or conflicts. The final aspect that is noteworthy from the quote by Scharre (2015) is human-control; how a swarm of drones would need a lower number of humans to control them, which could be a further asset to increasing the size of the US military.

In the primary research conducted key informant 1 (2020) raised the issue of how swarms of drones could potentially affect modern warfare. For example, drones can be deployed and effective in minutes, giving a state the advantage of mass and speed (Key informant 3, 2020). Key informant 1 (2020) also stated that drones have the ability to flock, go above enemy radars, and should be able to have a lot more manoeuvrability due to the systems being able to take more g-forces than a human could. Key informant 1 (2020) highlights this as important as a human pilot would black out. This is the one advantage in the near term. He agrees that the numbers of drones are a capability within itself. This gives a general idea of what drone technology is. This section will now move on to how drones could affect nuclear deterrence and strategic stability. Drones have become such a problem and traditional defence cannot handle drones (Key informant 3, 2020). Specifically, drone swarms could carry a potential payload, and beat air defence systems (Key informant 3, 2020).

Johnson (2020b: 19) states that AI-augmented AWS could be used in ISR and strike missions. Once again, this section aims to look at the conventional application of AWS, with a specific look at drone and swarm technology. Johnson (2020a, 2020b) states that even the AWS used

for conventional operations and possible proliferation could be destabilising as it can up the risk of inadvertent nuclear escalation.

> "For example, AI augmented drone swarms may be used in offensive sorties targeting ground-based air defenses and by nuclear-armed states to defend their strategic assets (i.e., launch facilities and their attendant C3I and early warning systems), exerting pressure on a weaker nuclear-armed state to respond with nuclear weapons in a use-them-or-lose-them situation" (Johnson, 2020b: 19).

A further issue highlighted by Johnson (2020a, 2020b) is the issue of states fielding unreliable, unverified, and unsafe AI-augmented AWS. This is notably problematic and an important aspect for this thesis, which will be discussed later. The use-them-or-lose-them situation raised by Johnson (2020b) highlights the issue of states feeling insecure about their ability to strike back; meaning their second-strike capability feels insecure. This will lead this state to nuclear escalation as they will be in an asymmetric position with their adversary. Johnson (2020b) highlights how a state may rely on a first-strike due to these insecurities. This shows the possible influence of AWS on nuclear deterrence. However, this thesis will now look at three different areas of technology such as:  air (Pedrix), sea (Sea Hunter) and undersea. This highlights the general aspects of drone and swarm technology. The next section will start to look at specific examples of how swarms will affect nuclear deterrence and subsequently strategic stability.

### 4.4.4    Swarm technologies' effect on nuclear deterrence

As stated already, drones and swarms give states the ability to conduct ISR and strike missions against their adversaries.  This section will further the analysis of the capabilities of drone and swarm technology and subsequently give a relevant example being developed by DARPA. The debate of the current capabilities of such technology will be discussed at the end. This section aims to look at how such a technology would be able to affect a state's mobile nuclear missile launches that gives their nuclear deterrence credibility. This section therefore aims to look at AWS swarms that have the capability to 'hunt' for the US's mobile nuclear weapons, like SSBNs, eroding the credibility given to nuclear deterrence by these weapons. Johnson (2020a) gives a good summary of such a capability:

> "Drones used in swarms are well-suited to conduct preemptive attacks and nuclear ISR missions against an adversary's nuclear and non-nuclear mobile missile launchers and nuclear powered ballistic missile submarines (SSBNs) and their attendant enabling

facilities (for example, C3I and early warning systems, antennas, sensors and air intakes)." (Johnson, 2020a: 5).

This quote by Johnson (2020a) also shows that swarms of drones may have the capability to not just threaten mobile missile launchers, but their attendant facilities as well. This can be argued as a strategic move that could be severely crippling to the US mobile missile launchers which would subsequently make such technology incredibly threatening and potentially destabilising for nuclear deterrence and strategic stability.

A specific example of such a drone would the DARPA's Sea Hunter. This allows this research to move forward to a more concrete example of AWS and make it more than just a theoretical debate. Key informant 4 (2020) stated that the Sea Hunter is being built in order to find submerged submarines. Johnson (2020b) highlights that the Sea Hunter is being tested in order to support anti-submarine warfare operations. Sauer (2019) further states that the Sea Hunter's capabilities could be used to detect and pursue SSBNs, which could potentially limit a state's second-strike capability. Johnson (2020b) argues identically, stating that technology like the Sea Hunter may render the underwater domain transparent, which will then threaten a state's second-strike capability. What is emerging from the literature is the potential for AWS, such as the Sea Hunter, to find mobile missile launchers that could then threaten a state's secure second-strike capability. The Sea Hunter, according to Johnson (2020a: 22), demonstrates how AWS are furthering the completion of the targeting cycle. By doing this, the Sea Hunter is exerting additional pressure on an adversary, potentially putting them in a 'use it or lose it' scenario when it comes to their second-strike capability (Johnson, 2020a: 22). Once again, it must be re-emphasised that a state may only need to perceive their second-strike capability as under threat from AWS in order to provoke a destabilising situation, which supports the argument that the invention of such a technology like the Sea Hunter and its application would be hugely destabilising to strategic stability. The Sea Hunter does not need to have a high-calibre capability, it only needs to have a capability that makes an adversary feel insecure.

Furthermore, the Sea Hunter is significantly less expensive than a warship, which would help the US deal with its issue of quantity on the future battlefield. Martin (2017) states that between 50-100 of these ships can be bought for the price of a singular warship. This would allow the Sea Hunter to be deployed as a swarm across the ocean. Klare (2019) states that instead of deploying well-armed, well-equipped and extremely expensive warships, the Navy could deploy a smaller number of crewed vessels with a large number of unmanned ships. AWS that

are equipped with AI and sensors could be trained to operate in a coordinated swarm that will allow them to overwhelm an adversary and give the US a quick victory (Klare, 2019). These swarms of Sea Hunters should in theory be able to detect SSBNs and subsequently threaten the most important leg of the triad.

These capabilities for the Sea Hunter are reviewed by the CRS, plus an additional two unmanned vehicles (UVs) (Navy large unmanned surface and undersea vehicles, 2020). The two additional platforms are the Extra Large Unmanned Undersea Vehicle (XLUUVs) and Large Unmanned Surface Vehicle (MUSV) (Navy large unmanned surface and undersea vehicles, 2020). In the report the Sea Hunter is categorised as a MUSV. A MUSV project for the US navy aims, like that of LUSV, to be low cost, have high-endurance, and have the ability to have its payload reconfigured (Navy large unmanned surface and undersea vehicles, 2020: 13). The first payload for this MUSV project will be ISR and electronic warfare (EW) systems capabilities (Navy large unmanned surface and undersea vehicles, 2020: 13). Furthermore, the US navy awarded a $34,999,948 contract to L3 Technologies to develop a single MUSV (Navy large unmanned surface and undersea vehicles, 2020: 13). This shows the Navy's dedication to building autonomous MUSV in order to conduct ISR. Such large investment could well mean that in the future there could be AWS that are able to search the ocean for nuclear mobile missile launchers. However, Johnson (2020b) states that the current effect AWS will have on deterrence is perception. According to Johnson (2020b) the near-term effects of AI on a state's nuclear deterrence would come from autonomy with an ML-augmented sensor which may potentially threaten a state's second-strike ability. This threat of crippling a state's ability to strike back may force a state into a position that will make them use nuclear weapons first. Johnson (2020b) calls this response a 'retaliatory first strike'. For the future of sea-based AWS, Johnson (2020b) attributes further advancement in ML and computing power to be a contributing factor to increase the capabilities of swarm technology to hunt SSBNs. This is a credible statement owing to the fact that the rise of AI in recent years has come from the increase in computing power.

Another potential use for AWS could be under the sea. AWS could potentially be used to hunt and attack SSBNs. Horowitz argues that the creation of undersea AWS that have the capability to track adversaries' SSBNs could lead to escalation where an enemy decides to strike first out of fear for their second-strike capability. Cebul (2017) states that unmanned undersea vehicles (UUVs) used in swarms could perform dangerous ISR under the ocean, this can be compounded by enabling them with AI and will affect strategic posture. Horowitz (2019: 781)

states that the ability for an undersea AWS to track an adversary's submarine has caused some to be fearful of such a capability. The fear would come from the capability of tracking undersea-based deterrents and ultimately undermine them (Horowitz, 2019: 781). This would be especially threatening if done on a nuclear-armed state, as it could then incentivise them to strike first (Horowitz, 2019: 781). However, Horwitz (2019) states that for technical reasons such capabilities are currently unlikely; he attributes it to power and communication. These are two frequently cited examples of limitations when it comes to AWS that operate in the sea. It will, however, be discussed in the next section that will look at the limitations and solutions of current AWS. The ability for a state to one day gain the capability to find SSBNs, either by the Sea Hunter or UUV, will be destabilising as it will threaten a state's ability to strike back in a nuclear first strike. In a moment of crisis, such a technology could force a state into the position of 'use-it-or-lose-it' that may force them to strike first.

An example of such of a UUV is the XLUUT (Boeing, 2017). Boeing (2017) claims that its XLUUV can travel 6,500 nm and can go for months on an operation fully autonomously. Furthermore, it does not need a launch and recovery platform; it can be launched from a port (Boeing, 2017). Furthermore, Boeing (2017) states that the current 'environment' for UUVs has several issues, such as range and endurance, which will be discussed in the next section. It is noteworthy now though, as Boeing (2017) claims their XLUUTV can overcome these challenges. The CRS reported that Boeing's Echo Voyager will be used in order to inform the design of Boeing's Orca XLUUV for the US Navy (Navy large unmanned surface and undersea vehicles, 2020: 14). Boeing partnered with Huntington Ingalls Industries (HII) to build the Orca XLUUVs (Navy large unmanned surface and undersea vehicles, 2020: 15). According to the CRS the Navy has mandated that the future of its forces should have up to fifty XLUUVs (Navy large unmanned surface and undersea vehicles, 2020: 14). Baker (2019) states that the US navy awarded Boeing $274.4 million to acquire just five of the Orca XLUUVs. The purpose of the Orca XLUUV is anti-submarine warfare, anti-surface warfare, strike missions, mine countermeasures, and electronic warfare (Baker, 2019). This shows the US's commitment to making the ocean transparent and investing heavily in AWS. There may be limitations in current AWS, which will be discussed in the next section, but the US is committed to investing in it, which may allow eventually for AWS that will be very capable of tracking SSBNs and subsequently threatening a state's second-strike capability. The Orca XLUUV enabled with swarm technology may enable it with a higher capability when it comes to 'hunting' for mobile nuclear missiles. As stated above, swarms allow drones to act in an intelligent manner that may

allow these undersea AWS to be more effective in a swarm than by themselves. A combination of different technologies from different domains may prove to be even more effective. This will be analysed in the next section. Finally, such investment from the US could well destabilise in itself as states may start to fear their second-strike capabilities being under threat. As mentioned throughout this thesis, nuclear deterrence is heavily based on perception, the capabilities don't need to necessarily exist, a state must just need to think that they do for it to be destabilising, which supports the argument that such investment may be destabilising in itself.

The final domain for AWS that will be discussed is air. Autonomous UAVs are another prominent example of a drone system that can be used in a swarm system and used to 'hunt' mobile nuclear missile launchers. As stated previously, drones used in swarms can be very well suited to conducting ISR and possible strikes against an adversary's mobile nuclear forces (Johnson, 2020a). An example of a UAV currently being tested by the US Air Force is given by Sauer (2019: 89), and has already been mentioned in this chapter; it is known as Perdix. The aim of Perdix is to operate as a 'team', this is done by communicating through radio and informing one another of what they are doing (Martin, 2017). Furthermore, Perdix is fitted with what Martin (2017) refers to as 'cellphone cameras'. A practical example of their capability was conducted by Martin (2017). The night before, Perdix was fed with 50 000 different pictures of Martin (2017) and sent out the next day to 'hunt' him. Once Martin (2017) was found, the drone relayed his coordinates back to a missile ship sitting on a nearby river. Martin (2017) states that this information fed from Perdix to the missile ship would allow it to send a direct missile strike. Another example of Pedrix in action is when the US Department of Defense launched 103 Perdix drones into the Californian skies (Lachow, 2017: 96). The demonstration showed Perdix's ability to fly without 'human help' and make decisions on its own (Lachow, 2017: 96).

Perdix offers a viable solution for finding land-based nuclear missile launchers. Furthermore, drones have the ability to go under as well as over traditional enemy radar (Key Informant 1, 2020), which will make them incredibly hard to detect. Key Informant 3 (2020) highlights that drones have become hugely problematic to traditional air defence systems. This issue is also mentioned by Kallenborn and Bleek (2018), who state that drone swarms have the potential capability to beat anti-submarine measures, missile defences, and air defence systems. This section on air-based AWS has briefly shown the current ability for drone technology by looking at the Perdix. The Perdix may have the potential to possess capabilities that allow them to

swarm in a way that allows for effective ISR on nuclear mobile missile launchers. The brief example given by Martin (2017) shows that Perdix potentially have the capability to be fed large amounts of data on mobile missile launchers and then sent to hunt them. There are numerous technical issues that may arise regarding the drone's sensors and issues based on deception by adding a sticker to a mobile missile launcher to deceive the drone. The issue of deception will be discussed in the next chapter and the theory of integration and limitation of AWS in the next section. A potential capability for AWS could also be a mixed swarm which will be able to threaten SSBNs, which is significantly more important to nuclear deterrence than mobile missile launchers.

### 4.4.5 *Limitations and solutions of AWS*

This section now moves on to discuss the limitations of current AWS. This section will start with the sea domain of AWS and the move on to the air. Gates (2016) states that AWS are linked to the laws of physics and because of this, certain aspects of AWS that will make the ocean transparent and SSBNs findable, is difficult. Gates (2016) states that SSBN were chosen to carry nuclear missiles owing to the fact that they are incredibly hard to detect and are designed to disappear into the ocean. Furthermore, SSBNs nuclear power allows them to be submerged for months at a time. However, Key Informant 4 (2020) states that SSBNs locations are roughly known by states, owing to intelligence and how states continuously track one another. There are other indicators that allow for the military to figure out where SSBNs are, such a warship, as there is usually a SSBN in the area (Key informant 4, 2020). Most importantly, as stated by Gates (2016), if an SSBN is hard to find there will be strategic stability. This emphasises the importance of SSBNs to a state's second-strike capability and the further importance of understanding the limits and capabilities of AWS.

The first challenge for finding SSBNs highlighted by Gates (2016: 30) is the ability to detect them. According to Gates (2016), current sensors on UUVs aren't capable of sensing SSBNs due to a number of reasons such as rain and waves, which provide very useful cover for submarines, as well as ships and sea life. Even powerful sensors are not necessarily capable of finding submarines, Gates (2016) equates it to using binoculars in fog: it will not necessarily help. Johnson's (2020b) research also indicates this to be an obstacle to finding submerged SSBNs. Johnson (2020b: 22) goes further to highlight three different obstacles that ocean-based AWS will have: 1. the ability for AWS to have reliable sensors on board to detect SSBNs is unlikely; 2. the sensors will be limited to the battery life (as well as the drone); and 3. The vast

size of the ocean will make it hard for even a swarm of drones to detect an SSBN. Furthermore, current computing ability aboard UUVs are problematic for sensor processing (Martin, Tarraf, Whitmore, DeWeese, Kenney, Schmid, & Deluca, 2019). Significant advances will need to be made in order for this capability to become possible.

Such advances could be the combination of ML and ocean sensors. Key informant 5 (2020) stated that you can train a bot, for instance, you can put sensors at the side of the road and listen to different cars. You can then train an AI bot to pick up different engine sounds and tell you what type of car it was (Key informant 5, 2020). Furthermore, Johnson (2020a: 11) states the progression in ML sensor technology will allow for better detection of Chinese SSBNs. Martin *et al* (2019: 12) also state that ML is the most viable candidate when it comes to bettering sensor information; however, they state that current ML is more based on supervised-learning algorithms. Furthermore, they require a larger amount of data retention and cloud computing in order to operate (Martin *et al*, 2019: 12). Martin *et al* (2019: 12) highlight potential limitations such as the availability of labelled datasets, sufficient time in theatre, and how well these algorithms can be implemented with embedded computation on board the AWS poses a challenge to naval applications. Martin *et al* (2019) are implying that current algorithms are computationally intense which is currently difficult to put into AWS; this is an issue that has already been looked at in this research and this once again re-emphasises the issue.

The next major issues for AWS are based on the issue of power; more specifically propulsion, which requires a large power source. Gates (2016) states that the ocean is dense and requires a considerable amount of propulsion and power to sufficiently propel a UUV through it. However, Martin (2017) claims that the Sea Hunter will have a speed of 26 knots and be able to track diesel-powered submarines for weeks. Such a capability may allow for AWS to effectively hunt SSBNs. Nevertheless, Gates (2016) states that an autonomous vehicle in the ocean must have sufficient propulsion and yet remain undetectable to SSBNs. However, while the Sea Hunter may have the capability to track SSBNs on the ocean surface for weeks, it is currently a limitation of undersea AWS. The limitation of current UUVs can be seen in Boeings Echo Voyager Extra Large Unmanned Undersea Vehicle (XLUUV). Boeing (2017) states that current environment of UUVs is limited by its endurance and subsequently there is a need for a launch and recovery platform. This is indicated in the research done by Martin *et al* (2019)*,* who conducted a survey for RAND on different unmanned ocean platforms. Martin *et al* (2019) concluded that one of the main areas that needs significant growth is the area of endurance and speed. However, Martin *et al* (2019) state that there is a growing interest and investment into

energy-dense battery technology. They further state that battery technology is in the interest of both the military and non-military sector, which promise growth (Martin *et al*, 2019). Furthermore, as stated in the preview sections, the US Navy is heavily invested in MUSV, LUSV and XLUUTV (Navy large unmanned surface and undersea vehicles: background and issues for congress, 2020). Such investment may allow for these technical issues to be overcome and sea-based AWS may one day be capable of detecting and tracking SSBNs effectively. However, as institutions like DARPA are trying to make SSBNs detectable, SSBN makers are trying to make them less detectable. This shows that the challenge to track SSBNs may be an incredibly difficult task, but it does have a certain degree of potential. There are also ways to overcome the issue of hunting in the ocean, this may be done by limiting the area to a 'choke point'.

This is a noteworthy solution for the current limitations in AWS as placing systems such as the Sea Hunter or Orca in 'choke points' may allow them to overcome their limitations and create an effective hunting capability based on their current strengths and weaknesses. In the primary research conducted, this solution came from Key informant 5 (2020). Key informant 5 (2020) suggested a slightly different version of this strategy to the one that will be argued for in this section. Key informant 5 (2020) suggested that a network of nodes could be placed in a geographical area, a choke point. These nodes can then be used to monitor submarines coming in (Key informant 5, 2020). It has already been mentioned that Key informant 5 (2020) stated that an AI bot can be trained to recognise sounds and tell which cars are approaching based on the sound of their engines. These nodes could be trained in a similar way to detect submarines coming into strategic choke points. A more relevant example that also uses the choke point strategy came from the findings of Gates (2016). Gates (2016) states that loitering at a choke point may allow loitering autonomous vehicles to trail SSBNs when they leave that choke point. Loitering within certain choke points could also potentially be seen as an act of aggression and have political implications (Gates, 2016). Therefore, these systems will need to loiter outside of these choke points (Gates, 2016). According to Gates (2016), all navies are aware of these strategies and have a number of techniques to lose other submarines or vessels that may attempt to follow them out of a choke point. Once again, the issue of being able to effectively follow these SSBNs and remain undetected is highlighted by Gates (2016). An example of technology being built in order to deal with choke points is Chinas 'Underwater Great Wall' (Wong, 2016). The Great Wall will consist of a network of ships and subsurface sensors which aims to undermine US undersea water advantage; as well as the Russians (Wong,

2016). This will allow the Chinese to track US or Russian SSBNs, or other crafts, in strategic areas such as the South-East China sea.

A system similar to the Great Wall made up of technology like the Sea Hunter, Orca and Key informant 5s (2020) nodes could further increase tension in strategic choke points like the South China Sea; an already noted point of tension for the US-China relations. Furthermore, such technology could be deployed into choke points or just outside of them depending whether it's for offensive or defensive purposes. According to Johnson (2020a), drones will not need to have ocean-wide coverage to detect or track submarines, an even spread of sensors may be capable of doing this. Furthermore, they could be located in check points or gateways (Johnson, 2020b). The defensive purpose can help the US find hostile SSBNs and the offensive capability may allow the US to track Chinese or Russian SSBNs leaving strategic choke points. Such an example would threaten Chinese and Russia's SSBNs which will in turn threaten their ability to maintain a secure second-strike capability. Whether the US could practically and willingly implement such a strategy that would directly undermine MAD is arguable. However, such a strategy could help mitigate certain shortcomings of current AWS and make their deployment feasible based on AWS current capability. However, a combination of both air and sea may create a more feasible option when it comes to effectively tracking SSBNs or other mobile nuclear missile launchers. This type of strategy is referred to as a 'mixed swarm' in the data that was collected.

Kallenborn and Bleek (2018) state that a mixed swarm of drones from undersea, surface, and air could be used in order to advance anti-submarine warfare. When it comes to autonomous communication between unmanned undersea, surface, and aerial vehicles, they would allow a wider area of coverage and surveillance which would allow them to relay information for a potential attack (Kallenborn, & Bleek, 2018: 536). Johnson (2020a) also states that in the maritime domain a combination of UUVs, USVs and UAVS that are supported by an AI-enabled ISR and intra-swarm communication could be deployed simultaneously for both offensive and defensive purposes. This would allow for the saturation of the enemy's ASW defence and more importantly allow for AWS mixed swarm to hunt SSBNs or non-nuclear submarines (Johnson, 2020a). A combination of these different systems may allow for a more effective capability when it comes to ISR. It may also allow states to overcome the shortcomings of each different area. This may provide an effective strategy for hunting an adversary's mobile nuclear missile such as SSBNs. The more effective AWS are at tracking an adversary's SSBNs the more destabilising swarm technology will be for strategic stability as

this capability will impact a state's secure second-strike capability that could force them to up escalation or strike first in a use-it-or-lose-it type of strategy.

In conclusion, this section aimed to look at different AWS and swarm technology. It started with a look at what exactly swarm technology was, then different drones systems combined with swarm, and finally the limitations and solutions for AWS swarms. Swarm-enabled AWS offer a significant advantage when it comes to hunting mobile nuclear weapons. It allows for the overcoming of certain limitations that current AWS have. Furthermore, mixed swarm technology that was looked at last offers the most viable solution to hunting mobile nuclear missile launchers as the different domains combined together offer a complete and potentially effective strategy. This will be looked at further in the next chapter of this research which will work on data and theory integration which will help to answer this thesis's research question. The next section of this chapter will look at loading nuclear weapons onto AWS; whether states will do this currently is highly unlikely, but it is of significant importance to analyse that data around it.

### 4.4.6    Loading AWS with nuclear weapons

One of the most important aspects of this thesis is looking at AWS being loaded with a nuclear payload. Due to the severity of nuclear weapons, equipping AWS with them is hugely problematic for the US and will probably not be pursued by them based on the US's value system. Value system was an important aspect when talking about the US with Key informant 4 and Key informant 1 (2020). Horowitz, Scharre and Velez-Green (2019: 4) argue that a state may deploy nuclear delivery uninhabited platforms for a number of reasons. Horowitz *et al* (2019: 4) state that nuclear-armed long endurance UAVs may allow states to have more nuclear signalling or strike options. Such uninhabited nuclear-armed platforms may allow states to secure their second-strike capability or utilise delivery systems that are capable of beating enemies' defences and target selection (Horowitz *et al*, 2019: 4). An example of a state pursuing AI-enabled platforms is Russia. Russia's AI strategy aims to reach a number of AI goals when it comes to AI platforms (Saalman, 2020). Russian AI goals are AI-enabled bombers, nuclear-powered unmanned underwater vehicles, and hypersonic glide vehicles that can carry nuclear and non-nuclear missiles (Saalman, 2020: 2). Another study by Horowitz (2019) also states that Russia is aiming to develop an AI-enabled platform known as 'Status 6'. Horowitz (2019) states that the US is also developing the new B-21 Raider which could potentially be flown as an unmanned platform. However, it must be noted that the US advocates

strongly for maintaining a human-in-the-loop, making it arguable that they will not hand over the kill call to an autonomous system when it comes to conventional weapons, let alone nuclear weapons. Horowitz states the same premise, that the US is focused on nuclear surety and would subsequently not deploy uninhabited nuclear platforms (Horowitz *et al*, 2019). However, other states might not show as much restraint as the US (Horowitz *et al*, 2019).

Giest and Lohn (2018) state that there are major changes for nuclear balances ahead and go on to state how Russia aims to create killer robots that have nuclear powers. Like that of Horowitze (2019) and Saalman (2020), Giest and Lohn (2018) go on to also mention Status-6. It would be able to overcome enemy defences through the use of speed and endurance (Giest, & Lohn, 2018: 2). Horowitz *et al* (2019) state that the Status-6 could use AI in order to avoid an adversary's anti-submarine warfare system. Furthermore, nuclear platforms that have AI capabilities may also provide a strategic benefit to whoever uses it (Horowitze *et al*, 2019: 21). Horowitz *et al* (2019: 21) indicated that the aim of the Status-6 is to secure Russia's second-strike capability as well as their confidence in it. This confidence will come from the ability of Russia's torpedoes always being able to reach their targets no matter what advances the US has made in their defence systems (Horowitz *et al*, 2019: 21). The deployment strategy for Status-6 is the ability to release it from Russian submarines in the Arctic and then being able to traverse the ocean at a speed of 100 km/hr (Giest, & Lohn, 2018: 3). Furthermore, the difficulty of communicating underwater has become possible recently with the progression of AI. Horowitz *et al* (2019: 22) states that this weapon may give Russian leaders assurance of their ability to strike back against the US homeland after a limited nuclear strike. However, any confidence gained and effect on strategic stability is likely to be marginal (Horowitze *et* al, 2019: 22). However, this risk could increase if Russia is to field untested and unverified nuclear AWS. This section ultimately shows that states like Russia are actively pursuing nuclear AWS and highlighted the potential affect it would have on nuclear deterrence.

### 4.4.7 The issues of speed, miscalculation and the potential for de-escalation

This section looks at the issues of speed, miscalculation, and potential de-escalation potential of AWS. There are several frequently cited issues that came from the literature which can all have a potential effect on nuclear deterrence. The first issue, and probably the most important of the three, is speed. Laird (2020) states that machine speed offers a significant operational advantage to the future battlefield. He argues that AWS are anticipated to be faster and more

agile than today's weapons systems as they are less dependent on human decision-making and control. This will give AWS the capability to make thousands of complex and coordinated decisions at machine speeds (Laird, 2020). Scharre states that automation will increase the speed of warfare and in turn make humans struggle to keep up (2014). Horowitz *et al* (2019) also state that algorithms will increase the decision-making process and create armed conflict at 'machine speed'. Johnson (2020a & 2020b) furthers the issue of AI and speed, stating that AI by itself will not be destabilising. However, adding AI to current military capabilities could further increase speed and compress the ability for humans to make decisions and create a destabilising effect (Johnson, 2020b; 17). Autonomous systems that have been given the authority to delegate over certain actions will be able to give that military the ability to react at machine speed (Horowitz, Allen, Saravalle, Cho, Frederick, & Scharre, 2018). This speed will give the military that has it an unprecedented capability over its adversary and this may be destabilising to strategic stability and nuclear deterrence, as states without a capable AI will not be able to react as quickly and as effectively as states that have AI.

Another possible capability of AI is the ability to process information faster that will allow commanders to make decisions faster in a rapidly changing battlefield (Horowitz *et al*, 2018).

> "A new generation of AI-augmented advanced conventional capabilities will exacerbate the risk of inadvertent escalation caused by the commingling of nuclear and strategic nonnuclear weapons (or conventional counterforce weapons) and the increasing speed of warfare, thereby undermining strategic stability and increasing the risk of nuclear confrontation" (Johnson, 2020b: 28-29).

All of the cited theorists thus far, reiterate the issue of 'machine speed' when it comes to the future battlefield. These machines will be able to make decisions faster than humans and increase the speed of future conflict.

Horowitz (2019) states that the speed of AWS may potentially threaten a state's first-strike capability. Horowitz (2019) further states that AWS will allow for states to win faster and likewise to lose faster. A state could fear that an aggressor with AWS capabilities might take out their command-and-control systems which may inhibit their ability to retaliate (Horowitz, 2019). If a state fears that it is at a disadvantage due to machine speed it may decide to strike first (Gates, 2016). Furthermore, within conflict, a state that fears it may lose at machine speed may escalate the intensity and possibly decide to escalate to the level of nuclear use (Laird, 2020). Horowitz *et al* (2019) state the same issue, that the fear of losing quickly could

incentivise a state to escalate fasters to nuclear use. Horowitz *et al* (2019: 30) also argue that 'sustained thinking' will allow states to back away from the nuclear brink, while AWS will not allow this to happen due to the speed it can respond at and subsequently undermine the security of time. What Horowitz *et al* (2019) mean is that the threat of losing time to machine speed could equate to a state upping escalation to nuclear first use. The ability to combine AI with ISR and defence could also create a destabilising effect, the ability to efficiently find an enemy's strategic assets and protect yours could erode away MAD (Technology for Global Security Reports, 2019).

> "The combination of exquisite ISR with an effective defensive shield could make it tempting to conduct a disarming, decapitating or blinding first strike at strategic targets, including nuclear command and control (NC3), early warning radars, or dual-capable missiles and aircraft" (Technology for Global Security Reports, 2019: 10).

This shows the potential for AI to increase military capabilities that could shrink the amount of time for human decision-makers and threaten nuclear stability. This could lead states to a use-them-or-lose-them situation; furthermore, speed plus range, mass, coordination, and intelligence can further compound this issue in future conflicts (Johnson, 2020b)

Another potential issue of AWS could be its ability to miscalculate a decision that a human will not be able to counter. If the speed of machines outruns the decision speed of humans, this means that humans will not be able to counter miscalculated decisions by AWS. An example given by Horowitz (2019) is an autonomous system that hits a strategic command-and-control system by mistake. Such an incident could cause possible escalation management issues (Horowitz, 2019). Furthermore, Miller *et al* (2018) state that uncertainty around these systems alone could cause an increase in the risk of miscalculation or a misunderstanding. Such miscalculations may come from AWS that have not been properly tested, verified and could possibly be unreliable.

However, such miscalculation could also allow for a chance to de-escalate a situation. Leys (2018) states that AWS offer a crisis bargaining opportunity; for example, AWS taking down another drone may be equated to a mistake and a potential opportunity to de-escalate a moment of crisis. An example is an F-35 flying in the South China Sea and being trailed by a Chinese AWS. The F-35 and the drone get too close and the F-35 slams into the drone, destroying it (Leys, 2018). Owing to the fact that no human life was lost and only a drone got destroyed can allow for a moment of de-escalation (Leys, 2018). This shows that, depending on what is hit,

there could be a possible moment of de-escalation. Another potential de-escalating potential of AWS came from key informant 1 (2020), who highlighted the issue of attribution. If AWS is able to gain the capability to either go over or under traditional radar and strike a strategic asset without any markings stealthily it could be difficult for the US to attribute the attack (Payne, 2020). Payne (2020) highlights that it is difficult to attribute such an attack; however, the US is very good at it. The time gap that is created due to the problem of attribution may allow for a moment of de-escalation, as the US would not be able to respond immediately with a strike and then time may change their tactic after the attack is attributed. The inability to attribute an attack would be a huge strategic disadvantage (Key informant 3, 2020). This section has subsequently highlighted the issue of speed, miscalculation, and the de-escalating value of AWS. The most significant issue to come from this section is the issue of machine speed outrunning the human decision-making process and the subsequent issues that arise from it.

## 4.5    AWS arms race impact on strategic stability

Another potential impact on strategic stability and an important tenet in nuclear deterrence is the possibility of an AWS arms race. Giest and Lohn (2018: 8) describe arms race stability as when a state is not attempting to exploit its adversary's military capabilities. The ability for states to find their adversaries mobile missiles could be a capability that states may find worthy of pursuing and equate to a possible arms race (Bracken, 2017).[10] Another example of the hunt for mobile missile launchers comes from the work done by Johnson (2019). Johnson (2019) states that the fact the Trump administration is building mobile missile launchers and aims to triple spending on AI may be a causative factor in a potential arms race that may upset global nuclear balance.  Lucas (2016) also states the same premise, that the capabilities of AWS will force states to develop their own autonomous capabilities for their weapon systems. This shows that states will develop AWS capabilities in order to maintain a military capability that doesn't lag behind those of states that have acquired AWS. Furthermore, the US's emphasis on maintaining its military supremacy and tripling its spending on AI may force other states to pursue a similar strategy. This could be compounded if the US were actually able to field technology such as the Sea Hunter that is fully capable of detecting and following Chinese or Russian SSBNs. The ability to hunt another state mobile nuclear missile will be the outright most destabilising aspect when it comes to nuclear deterrence and will trigger an arms race. As

---

[10] This comes from Brackens 'The Intersection of Cyber and Nuclear war' which focuses on the cyber aspect of hunting for mobile missile launchers (2017). However, this same logic may be applied to AWS hunting for mobile missiles launchers and is why Brackens work is used.

other states will pursue creating such a technology in order to Hunt US mobile nuclear missiles in order to maintain MAD, it would be imperative that an adversary knows where US nuclear forces are in order to stop an asymmetric nuclear situation. Another such aspect would be states pursuing technology that will ensure their second-strike capability, such as Russia's Status-6.

## 4.6    Conclusion

The aim of this chapter was to present all the data collected from secondary data and key informant interviews. The next chapter will focus on data and theory integration in order to answer the main research question. This section was structured using Kaldor's new war thesis as the theoretical framework. The first section focused on the actors involved in the US's nuclear deterrence and pursuit of AWS. What this means is that this section aimed to look at who the US most fears as strategic adversary. Actors are important as nuclear deterrence is heavily based on perception. Furthermore, this research aims to understand the full complexity of how AWS will affect nuclear deterrence, not just the technical capabilities. What was found from the secondary data and key informant interviews is that the US strongly believes that they have let their competitive advantage slip and that it must ultimately be restored. The adversaries that pose the most major threat to the US, according to the NSS, are Russia and China. It can be argued that the US currently sees China as more of a threat as compared to Russia, especially based on the fact that the US attributes the increase in China's military power originating from the theft of US intellectual property (NSS, 2017). However, the NSS (2017) stills sees Russia as a threat to the US and its allies. The US aims to rebuild its competitive advantage through strength. The section that followed on from actors was the 'goals' of the new war thesis. This section aimed to collect data on populism and identity politics to see if they can be attributed as the drivers in how the US pursues AWS. What was subsequently found was that Trump's policy and his actions are different. While Trump's tweets and statements may show populist sentiments and be riddled with identity politics, his policy seems to sit close to the centre and is aligned with previous administrations. However, there is nothing binding Trump to his policies in a time of crisis, which makes him a potentially volatile actor for both his adversaries and allies. This is based on the fact that Trump sees his volatility and unpredictability as a strategic asset. This section aimed to look at data for the purpose of creating a contextualization for this study. The next section looked at the modes of warfare and made up a major part of this chapter. This is where the data for answering the main research question was presented. This section was structured into nuclear weapons, conventional application for AWS, limitations and solutions of AWS, loading AWS with nuclear weapons, the issue of speed,

miscalculation, and potential for de-escalation, and the impact of an AWS arms race on strategic stability. The main issue that comes from this section is the current capability of AWS due to the limitations of the power and sensor capabilities currently available. However, there is potential for AWS to affect nuclear deterrence. This will be looked at more in-depth in the next section.

# Chapter 5

## 5.1    Introduction

The theme of this research, among other things, was to look at how states go to war. The premise of this thesis being that war is changing due to the rise of new technological innovations. However, there is a large amount of technological innovation happening and this research cannot cover all of it. As a result, this research aimed to look more specifically at AI, however, there are also many potential applications for AI. Therefore, this research aimed to look specifically at AWS. Furthermore, there are many modes of war that can potentially affect strategic stability when it comes to modern warfare. In order to be more specific, this research aimed to look at arguably one of the most important aspect of modern war, nuclear deterrence. This led to the main aim of this research. This research aimed to look at how AWS could lead to a potential disruption of nuclear deterrence and its subsequent effect on strategic stability.

This research began with a review of the literature on key tenets in order to conceptualise them and create a better understanding of the area in which this study is placed. The chapter that followed the literature review aimed to create a contextualisation for the study. The scope of that chapter was to look at the role technology played in how states conducted warfare. This approach was chosen because this research aimed to look at the role technology played in how states conducted warfare. Furthermore, it allowed for a more in-depth look into nuclear technology and highlighted the significance of the nuclear triad when it comes to nuclear deterrence. More specifically, it highlighted the importance of mobile missile launchers and more importantly SSBNs. This was found to be of importance to this study as it 'secured' the US's second-strike capability, the most important tenet of this research. The following chapter discussed the data collected by this study. This data was collected from secondary data analysis and semi-structured key informant interviews, which will allow this research to answer the main and secondary research questions.

The aim of this chapter is threefold. Firstly, it will begin with an evaluation and analysis of the main findings which aims to integrate the theory from Chapter 2 into the data in order to answer the main and secondary research questions along with the secondary contextualization. This section will take up the majority of this chapter as it is of the most importance to this research. Secondly, it looks at the limitations of this study. This section discusses the limitations of the study such as the issue of potential interviewees not responding, and the case study being based on the US. Thirdly, and finally, this chapter will have its conclusion and a discussion on future

avenues of study. This section will draw the research to a close and discuss how this research has created possible future avenues of study.

## 5.2    Evaluation and main findings

As stated, the main aim of this section is to answer the main research questions. This will be done by using the theory gained in Chapter 2 and integrating it with the data collected in Chapter 4. The structure that will be used in order to answer this question will be the new war thesis. This means that the structure of this section will duplicate the one used in the previous chapter. This section will focus on data and theory integration so that the research questions may be answered.

*Main Research Question:* How will Autonomous Weapons Systems affect the US's perception and/or capability of their second-strike capability?

*Secondary Research Question*: Will this led to a disruption of traditional nuclear deterrence? or will this subsequently affect strategic stability?

*Secondary Contextualization:* The new war thesis entails 'goals' which looks at how identity politics effects conflict. In order to create a background for this research populism and President Trump will be looked at. This is a more specific form of identity politics in the US. The goal of this is to create a contextualization for this study and not to research populism in the US.

*Broader significance of the study:* Firstly, Autonomous Weapon Systems will affect the way states conduct warfare and more specifically their nuclear deterrence. It is important to get an understanding of such a phenomenon so that uncertainty about the future may be mitigated. Secondly, populism is a dominate political movement in the US and has been expertly utilized by President Trump. Understanding such a phenomomen will allow for this study to full capture the effect of the main research question.

*How this question relates to the problem/conversation in the literature:* There is a high level of uncertainty around the issue of AWS, creating a need to fully understand the complexity and effect of such a transformative technology. Furthermore, populism is a rising political ideology that has the potential to affect international politics and US national security strategy. Which means it requires analysis in order to understand its potential affect internationally and not just domestically.

### 5.2.1    Who is involved?

Actors were of significant importance to this study because perception plays such a vital role in nuclear deterrence. Furthermore, the ability for AWS to effect nuclear deterrence is not just based on the technology involved; it is more complex than that. There are other drivers and factors involved in this process that need to be analysed in order to reach a significant and coherent conclusion. This is why this research chose to create a secondary contextualization of President Trump and his populist politics. One of these factors and drivers is the actors involved. Furthermore, who these actors are and how the US perceives them is even more important. The mere pursuit by one of these actors for AWS that will potentially have the capability to threaten the US's second-strike capability could be more than enough to effect nuclear deterrence and subsequently have a destabilising effect on strategic stability. As stated, this section will begin with an analysis of the data and be followed by a discussion on the theory then an integration of the two.

What did the data collected show when it comes to actors involved? Firstly, the actors that are involved in this analysis are what Kaldor (2013) would define as 'regular armed forces'. Furthermore, there are many 'regular armed forces' that the US views as strategic threat to their interests and security. This is logical as there are multiple nuclear armed states, states that are pursuing nuclear weapons, and states that are powerful enough to pose a potential threat to the US. This highlights the fact that thesis decided to look at other states when it came to actors and not other non-state actors. This is owing to the fact that this thesis cannot look at everything and it therefore decided to look at the most prominent threats to the US. The document that was chosen to be analysed in order to collect data to understand these actors was the NSS (2017), which was referred to as key informant 2 (2020). What came from this data is that the US has two primary concerns when it comes to potential threats: China and Russia. The most important facts behind why these actors are threatening is due to their pursuit of what the NSS states as advanced weapons and capabilities that could pose a potential threat to US critical infrastructure and command and control (NSS, 2017). Furthermore, as highlighted in this research, both actors have the desire to gain general intelligence AI.

It was found that there is a greater threat and higher tension when it comes to the US and China. Johnson stated that there is a US–China AI race that could have a potentially profound impact on strategic stability. Furthermore, according to Giest and Lohn (2018), both China and Russia view the US as leveraging AI to undermine the survivability of their nuclear forces. When it

comes to Russia, the NSS (2017) states that it aims to rebuild its status as a great power. The NSS states that it still aims to deter Russia and only create cooperation when it best suits US interest. Furthermore, Russia also aims to become a global leader in AI. Laird (2020) states that AI-enabled AWS will increase the risk of crisis instability and possible conflict escalation in future conflict. Furthermore, as China and Russia aim to increase their military capabilities, the US aims at regaining its military supremacy. The NSS (2017) highlights that the US believes that it has lost its military power and has let other actors catch up to it. The US also sees Russia and China developing military capabilities that could threaten critical US infrastructure (NSS, 2017).

There is also critical tension between the US and China due to China's modernisation arising from the access it has to the US innovation economy and universities (NSS, 2017). Furthermore, China has access to MCF which allows them access to advanced technologies (NSS, 2017). This could be hugely problematic for the US in pursuit of AI as it could potentially be stolen and once they have the AI program they have the full capability. There is no need for further R&D as they will have the entire code. This shows that the US already has a very negative perception of China. If the Chinese military were able to field an AWS system, effective or semi-effective, it will undoubtedly exacerbate tensions with the US. What would be even more worrying is if China were able to field an AWS that was incredibly effective and capable of finding US SSBNs or other mobile nuclear missile launchers.

This section has shown that understanding the actors involved is an important part of the new war thesis. This is even more compelling when the chosen mode of warfare is nuclear deterrence owing to the fact that perception plays a vital role in it. The US views Russia and especially China as threatening adversaries. If these nations were able to deploy AWS that has the potential to increase their military capability or undermine the US second-strike capability it could lead to strategic instability as it would increase tension with the US, who are already highly aware of their declining military advantage.

### 5.2.2    Goals and identity politics: Making America Great Again

The data collected for goals in this research offered some interesting insight into President Trump as an actor against his actual policy. The objective of this contextualization was to understand how Trump's populism could affect how the US pursued AWS. This research looked at this by analysing the US's NSS, the premise of this being that Trump offered a different approach to US grand strategy compared to the previous administration. It was found

that Trump may take a populist stance and this stance got him elected, however, his policy has turned out to be closer to the centre than many of his critics fully account for (Cordesman, 2017). It is still of importance to understand Trump's identity politics, as in a moment of crisis, there is nothing holding Trump to his policy. This section will now begin with a look at the important data collected in Chapter 4 and then integrate it with theory gained in Chapter 2. This will help further the secondary contextualization and create a better background for the study.

One of the first major pieces of data that was collected was based on populism and how Trump utilised such a movement in order to gain political power. Trump as an actor and his policy are inherently different, according to the data. The data from Löfflman (2019) argued that Trump became so powerful owing to the fact that he was able exploit the disconnect between the elite and public opinions. Trump targeted the establishment as the ones to lay blame when it came to economic and political failures (Löfflman, 2019). He was able to build an antagonistic relationship as well as play on US citizens' emotions (Löfflman, 2019). This supports the argument of Batayneh *et al* (2016), who viewed populism as 'anti-elitism' and what the true citizen looks like. This tenet was also highlighted by Fukuyama (2017) who stated that populist leaders try to make an image of the true citizen and denounce the elite. Furthermore, another main tenet that came from the identity politics was the issue of dignity. This issue of dignity was utilised by Trump by attacking the elite for ignoring public opinion. What this means is that Trump's stance came from the belief that the elite were interested in economic liberalisation while the 'ordinary citizen' wanted the US to focus on its own domestic issues.

Trump was successfully able to create a friend-enemy distinction between the true citizen and the elite by highlighting the gap between public and elite opinions on foreign policy. While this form of identity, at the time of writing this thesis[11], had not created any conflict like Kaldor's view of identity politics; it does have similar traits to the new war thesis. The distinction between two identities that is bought forward by the new war thesis is the binary that is created in identity politics, which is represented by Trump's true citizen and the elite. Fukuyama states that populism is a direct threat to international order.

Fukuyama (2018) would be right if the data Löfflman (2019) collected is correct. Löfflman (2019) argued that these sentiments are not influencing Trump's NSS, it did however open the debate on recalibrating US foreign policy. This stance by Löfflman (2019) was backed by

---

[11] This research was conducted before the protests and current unrest in the US broke out. Identity politics would have been a good theory to use to analyse this and understand the Trump administration. This will be discussed in the limitations section of this chapter.

Feaver (2017) who also states that Trump and his policy are inherently different. What this means for this study is that Trump may say something in line with populism, however, his foreign policy does not necessarily reflect the same sentiment. This would be going against the literature used in Chapter 2, as Wasko-Owsiejscuk (2018) stated that Trump's foreign policy has a 'large doses of populism'. While Trump's NSS does state that it aims to Make America Great Again and restore America's power that has supposedly slipped (2017). Löfflman (2017) argues that the US still supports NATO and its allies and it still imposed sanctions on Russia which shows that Trump's America First means internationalism over isolationism. This shows that Trump played on the issue of a true citizen against the elite, yet his NSS ended up more towards the centre in shaping his NSS. This means that Trump's belief in listening to the ordinary citizen and focusing on domestic issues over international issues has not been true.

However, according to Pifer (2018), when it comes to a moment of crisis there is no strategy dictated by the NSS that must be implemented to deal with this situation. This is problematic for both the US's adversaries as well as its allies owing to the fact that Trump sees unpredictability as an asset (Pifer, 2018). This means that Trump's instincts as well as his volatility are what should be questioned according to Pifer (2018). Key informant 3 (2020) gave some key insight into how a Trump or Biden response would look. A Trump administration would come with a strong response and not blink in a fight; meanwhile, a Biden administration could come with a more diplomatic response (Key informant 3, 2020). A final conclusion to this section comes from Hall (2017) that Trump is in a battle with mainstream policy thinking and his own instincts. What this means that it is more useful to understand Trump as an Actor through his rhetoric and actions over his policy. In a time of crisis, a Trump response could be unpredictable and volatile. However, this reaction could be more in line with his populist rhetoric that could put the US first and its allies second. It could also see a US response that aims to maintain the US's supremacy as a global power and fight for US interest. What can be concluded by Trump's NSS is that the US will be pursuing AI and AWS if it allows for the US to regain its military strength. This can be stated as the US has already budgeted around $14 billion for the DoD science and technology programs. Finally, this section has shown the importance of a secondary contextualization when it comes to understanding how AWS will affect the US's nuclear deterrence.

### 5.2.3    Methods of warfare

This section now aims to answer the main part of the thesis, how AWS will affect nuclear deterrence and how this will in turn affect strategic stability. This section was divided into five different sections: nuclear weapons, conventional application for AWS, loading AWS with nuclear weapons, the issue of speed miscalculation and potential for de-escalation, and AWS arms race impact. Each of these sections plays a vital role in understanding how AWS could affect the US's nuclear deterrence. The first section aimed to understand the continued importance of nuclear weapons to the US's NSS. It can be argued that nuclear weapons have been the cornerstone of US NSS policy for the last seventy years and will remain the cornerstone of the NSS policy. The NSS (2017) states that nuclear weapons have given the US the ability to deter potential aggressors, keep peace, maintain stability, and protect its thirty allies. The argument that it will stay as the cornerstone of US policy comes from the current strategy to modernise the US's nuclear forces. By doing this, the US will be able to maintain its credibility as well as signalling its assurance of its capabilities (NSS, 2017). Maintaining the effective nuclear capability will allow the US to continue what Schelling (1966) termed the 'diplomacy of violence' (1966). The diplomacy of violence involves using fear through military might and will in order to deter an aggressor from conducting an unfavourable action (Schelling, 1966).

This data collected from the NPR in the previous chapter highlights a few of Quackenbush's (2011) tenets for an effective deterrence. It states that the US will carry out unacceptable cost on an attacker and that this threat will be carried out (Quackenbush, 2011). Furthermore, the US aims to modernise its nuclear triad, which will then allow for all three tenets of Quackenbush's (2011) successful deterrence to be met as the US will have an effective military capability. The triad has maintained both an effective and credible US deterrence. This credibility given by the nuclear triad, and specifically SSBNs, allowed for the US to have a 'secure second-strike capability'. In the literature it was found that a second-strike capability is central to nuclear deterrence. With a secure second-strike capability, states that maintain nuclear weapons have the capability to engage in brinkmanship. As nuclear weapons enable MAD, in the case of nuclear warfare both states would receive a nuclear strike against themselves. Such a dynamic has allowed the US and other superpowers to deter one another successfully. The dynamics of brinkmanship involve states taking risks that up escalation towards all-out war (Powell, 2003). The most resolute of the states involved in brinkmanship will be the most willing to up the risk (Powell, 2003). An Equilibrium equation made by Powell

(2003) consists of the state's resolve, its perception of its own resolve and the other state's resolve. The state that is the most certain will be the one most willing to escalate. Therefore, having an effective nuclear triad will allow the US to have confidence in its capability that will allow it to be the more resolute state as currently both Russia and China are modernising their nuclear forces and if the US falls behind it may become less resolute. Finally, the most important point to come out of this data was the importance of the US nuclear forces, its second-strike capability, and its nuclear triad. A state's ability to strike back is crucial to an effective nuclear deterrence, the possibility of AWS cancelling out a state's ability to strike back could be hugely problematic to nuclear deterrence and strategic stability.

This next part aims to look at the AWS that may threaten the US's ability to strike back. This research aimed to look at current examples of AWS in order to understand if they are currently capable of disrupting nuclear deterrence. If AWS has the capability or even the perceived capability to take out the US's nuclear triad it would subsequently erode the US credibility of its second-strike capability and cause instability. The research began with a look at the conventional application of AWS. The two main points that came forward were drone technology and swarm technology. This research decided to negate land-based AWS as sea- and air-based AWS were of more importance. The main application of drone technology would be the ability for them, as a singular unit or as a swarm, to hunt mobile nuclear missile launchers. As highlighted by Giest and Lohn (2018) this capability would undermine MAD and nuclear deterrence.

The first set of data that was analysed was swarm technology. The data showed that swarm technology relies on AI, it is assumed that many authors are referring to 'general' intelligence AI owing to the fact they refer to how these swarms will be able to operate by themselves without human intervention. However, as found in the literature review, depending on the complexity of the environment, the AI in these systems does not necessarily need to have general intelligence. As stated by Boulani (2019), machine learning and deep learning are currently more than sophisticated enough to be applied to autonomous systems. It will be even more beneficial for a swarm to be able to give rise to intelligent behaviour through its cooperative behaviour (Scharre, 2014). Drones in a swarm communicating with one another are able to alter their behaviour (Kallenborn & Bleek, 2018). Furthermore, the more data that can be fed to AWS that is enabled with narrow AI, the better these systems will be able to perform. For example, if these drones are fed large amounts of data about SSBNs, the sounds it makes and its actions, they could then be released to hunt the US's submarines and begin to

teach themselves how to hunt submarines more effectively. The more data that these systems are fed the better and more efficient they will be. Drone swarms will also be able to spread out over a larger area and collect more data and communicate with one another that may allow them to hunt more effectively. The other benefit of drones is their cost and the need to have fewer humans operating them. This will allow the US to gain mass in its military and deploy more drones that will gain more data. Swarms will also allow for better speed and manoeuvrability than human-operated vehicles (Key informant 1, 2020). Deep learning, compared to the more automated and structured ML, could find solutions quicker and faster when it comes to hunting SSBNs. It will also allow swarms to operate more autonomously and at a quicker speed. This could be further accelerated if these machines are introduced to deep learning; however, due to the 'black box effect' it is arguable whether states will currently equip swarms with this. However, this is all dependent on the quality of the sensor on these drones and their processing power. This now brings this section to a discussion of examples of drone technology.

The two main domains that were looked at in regard to drones were air and sea. The sea was separated into a further two categories, under the ocean and on top of the ocean. The aim of the Sea Hunter, according to key informant 4 (2020), was to hunt submarines. Such a capability would definitively impact the US's second-strike capability owing to the fact that the SSBNs play a vital role in the US's second-strike capability. Furthermore, this capability to hunt mobile nuclear missiles could be increased if an adversary was able to use them as a swarm. AI in the Sea Hunter, whether it be ML or deep learning, could create intelligent behaviour that may allow them to effectively hunt and in new ways. Swarm technology would also allow them to cover a vast area of the ocean more effectively. This would require advanced sensors, computing power, and enough power to effectively hunt. Johnson (2020b) highlights that ML-augmented sensors would be sufficient to threaten a state's second-strike capability (Johnson, 2020b). This shows that ML-enhanced Sea hunters operating as a swarm could be capable of hunting adversaries' SSBNs. Furthermore, as the Sea Hunters hunt the oceans, they in theory will become more effective, as the more data that that this AI is fed, the more effective it will become.

The next sea-based AWS that was analysed operates under it. The potential capability for AWS to hunt SSBNs under the ocean could undeniably impact nuclear deterrence. Horowitz (2019) states that this ability, especially against a nuclear-armed state, would cause fear and could incentivise them to strike first. However, the current limitations such as power and

communication, do not currently make such a capability possible (Horrowitz, 2019). What this data found is that the current capabilities of underwater AWS are not efficient. However, a swarm of underwater AWS could one day have the capability to effectively hunt another state's SSBNs. Such a capability would undoubtedly be destabilising for strategic stability as it would take away one of the most valuable legs of the nuclear triad. But for now, it faces power and communication challenges that will inhibit it from doing so. However, based on nuclear deterrence theory, these undersea AWS only need to make the US feel uncomfortable about its capability to hunt its SSBNs to cause insecurity. With already increased tension with Russia and China, such an innovation could cause strategic instability and force the US to nuclear first use.

The final domain that was looked at was air. Autonomous UAVs are a prominent example of swarm technology that could be used to hunt mobile nuclear missile launchers. An example found in the data was a drone called Perdix (Sauer, 2019). Although not stated, it can be assumed that Perdix used a form of ML to process that data given to it and then to hunt down its target. The data showed that current AWS in the air have a high capability when it comes to working as a swarm in order to execute a task. What the data found is that air-based AWS have the capability to sense, decide and act based on their environment by collecting data from its sensors. These examples are not based on true autonomy and may sit closer to automation, as they are executing a task that was predetermined by humans. However, once they take off they were able to hunt on their own and find their target without human intervention. This could then still be a highly desirable capability as hundreds or thousands of low-cost drones could be sent out in order to hunt the US's mobile nuclear missile launchers. This would require fewer humans and increase coverage as well as speed when it comes to hunting.

However, there is potential for deception and possible limitation of these drones' sensors. Boulanin and Verbruggen (2018), quoted in Chapter 2, stated that ML has allowed for better recognition and targeting. If these aerial AWS are fed a large amount of data and are able to get flight time in theatre, they will have the potential to become even better at targeting and recognition. The only limitation would then be the current lack of data and the need for more flight time. These drones have shown their capability to be given data and commands and then execute a mission without human intervention. Which means these drones meet the true definition of autonomous, as they are able to take off and operate without the need for human intervention, although there is a human-in-the-loop. It can be argued that the current aerial drones may have a limited capability to hunt mobile missile launchers, however, more data is

needed and better sensors. Furthermore, potential adversaries need to field reliable and verified drones to make sure there are no possible mistakes.

The final section of the *'conventional application for AWS'* focused on the limitations and possible solutions for AWS. The current limitations for sea-based AWS are their sensors and power source. Such limitations will currently inhibit sea-based AWS from effectively hunting SSBNs in the ocean. ML and AI require high-powered computing in order for it to be effective. ML has come about recently due to the increase in computing power, which means these AWS need to have a sufficient computing capability. ML combined with more improved ocean sensors could lead to a more effective capability when it comes to hunting SSBNs. This data combined with potential AWS capabilities suggest that the answer to the main research question could be that AWS could potentially have the ability to make the US feel insecure about their second-strike capability. However, this is not a current issue for the US, as other issues such as the rise of China's military power is more of a threat to the US, especially since the Chinese are building and modernising their nuclear forces.

With all these limitations, it suggests that currently AWS will not have the capability to effectively threaten nuclear deterrence. However, a possible solution for the limitations of AWS could be using choke points or a mixed swarm. A choke point is a strategy in which a swarm or singular AWS are placed in a strategic position so that an SSBN or another mobile nuclear missile launcher has to pass through. The point of placing AWS at choke points can help overcome the issue of the ocean being vast. The sensors and power of the AWS will still need to be more capable than it currently is in order to continuously track SSBNs. Which means that the choke point strategy may one day be a viable strategy; however, the current limitation of AWS will still inhibit its ability to hunt.

Another solution is a mixed swarm strategy. A mixed swarm is made up of undersea, surface and air AWS. This mixed swarm will enable the ability to create a wider network that can then increase coverage of surveillance and possible attack (Kallenborn, & Bleek, 2018: 536). This strategy of a mixed swarm was also seen in the report by Johnson (2020b), who states that it could be used for both defensive and offensive strategies. This strategy may allow for AWS to overcome their shortcomings, as there will be an array of different AWS involved. It will allow for an increase of coverage, more data collection, and can rely on each other to overcome power shortages. There will still need to be a significant advancement in AWS capabilities in order

for a mixed swarm to be effective. This shows that both these possible solutions may one day be effective; however, the current limitations will not allow this to be effective.

The next section dealt with loading nuclear weapons onto AWS. The reasons for loading AWS with a nuclear payload can be to strengthen a state's second-strike capability or to beat adversaries' defence systems. An example found when it comes to a nuclear-armed AWS was Russia's Status 6. Status-6 aims to overcome enemy defence through speed and endurance (Horowitz *et al*, 2019). Russia believes that Status-6 will allow them to secure their second-strike capability. However, handing over nuclear payload to AI that is still based in ML may be hugely problematic as ML needs large amounts of data in order for it to be effective and the deployment of Status-6 needs to be frequent to gain sufficient data. Furthermore, handing a nuclear payload over to a deep-learning AI would be even more problematic due to the black box effect. The issue of sensors and power will also affect Russia's Status-6, it does, however, allow for a viable solution to maintaining nuclear deterrence in the future.

There are several advantages to AI-enabled AWS with the main one being speed. It was found in Chapter 2 that AI offers speed as a capability, as it is able to think and act at machine speed. This was also found in the data and was highlighted by theorists as an important tenet. This increase of speed by AI could force states into a use-them-or-lose-them scenario. Johnson (2020a & 2020b) states that speed plus range, coordination and intelligence can further affect the future conflict. This shows that the speed AI brings may enable AWS to hunt and find the US's nuclear forces faster than humans could. Such a capability may allow an adversary to neutralise the US's nuclear forces, inhibiting their ability to strike back. This would ultimately eradicate the US's ability to strike back and subsequently undermine their nuclear deterrence. Furthermore, once the US's nuclear deterrence is undermined this would cause strategic instability.

The final section dealt with AWS arms race that could possibly affect strategic stability. In the literature review in Chapter 2 it was found that an arms race or the emergence of an arms race could cause military instability which could then affect strategic stability. Arms race stability is created when a state is not pursuing a capability that could exploit an adversary's military capabilities (Geist & Lohn, 2018: 8). The capabilities offered by AWS and states pursuing them could cause military instability as gaining these capabilities will allow for another state's military capabilities to be exploited. This is increased by the fact that AWS is augmented with AI, the most transformative technology. Furthermore, the US is aware of its military capability

slipping and is dedicated to maintaining its edge. China and Russia are also very determined to develop and improve their military capabilities in the nuclear and AI sphere. This also points towards a possible arms race between states that already have increased tension between them, which means there could be military instability, and this also could affect strategic stability.

### 5.2.4    Final discussion

The answer to the main research question is twofold and is based on current and future capabilities. Currently, AWS does not have the capabilities to affect the US's nuclear deterrence. The limitations of sensors, power, and computing power currently restrict AWS ability in the near term to effectively hunt from mobile missile launchers. With regard to future capability, AWS offers a significant challenge to traditional nuclear deterrence and has the potential to affect the US second-strike capability. All the information reviewed offers insight into how AWS will be able to 'hunt' the US's second-strike capability such as its SSBNs. With AI having progressed so quickly and rapidly over recent years there is a possibility that there could be an effective version to use in AWS. Furthermore, the US is set on maintaining its technological innovation in its military and is subsequently investing heavily in its military capability. In the future there could be AWS with better AI, sensors and power that will allow them to find and track the US's mobile nuclear missile launchers.

As discussed in this thesis, AWS effect on nuclear deterrence is more complex than just its capability to 'hunt' mobile nuclear missile launchers. This thesis tried to understand the US drivers and its perception of major threat. This was done owing to the fact that this research used the new war thesis in order to understand how AWS will affect nuclear deterrence. Furthermore, this research had the objective to create a secondary contextualization which then allowed for a better background for the study. This allowed this thesis to understand the goals of how the US would pursue AWS. With the rise of populism in Western democracies and the ever more prevalence of identity politics it was important to understand how such phenomena would affect the US perception and drive. The literature showed that President Trump is a populist, and the data did indeed show this; however, his NSS is more towards the centre. The US NSS policy can be argued as 'Making America Great Again', but in an international perspective and not isolationism. The actors section allowed this research to understand that the US is highly concerned about China. It is still concerned about Russia, but more emphasis was put on China as an adversary and building cooperation with Russia. Ultimately, the US aims to increase its military power in order to be able to deal with actors such as China. This

belief in rebuilding is what will fund AI and lead to the development of AWS, especially owing to the fact that the US faces a quantity issue due to rising prices in technology and AWS offers an alternative and will allow for a large quantity and high capability. Finally, President Trump as an actor could be unpredictable when it comes to the pursuit or reaction to AWS. This is owing to the fact that Trump views his volatility and unpredictability as an asset. In the event of such a transformative technology and Trump stability becoming instability, he will not blink in a fight.

The final conclusion is that AWS in the future will have the potential to affect nuclear deterrence which will subsequently cause strategic instability. However, currently, AWS does not have this capability. The current issue of nuclear deterrence and strategic stability is the modernisation of the nuclear forces of the US, Russia and China. Furthermore, with an actor such as Trump as president there could be more of a hard response to China in a moment of escalation compared to a more diplomatic response from a Biden administration.

## 5.3    Limitations of study

This section aims to outline the limitation of this study. There are several factors that can be pointed to when it comes to the limitations. The first limitation of this study is the response rate from possible interviewees. This research aimed to have 8–10 key informants and arrived at only 5. This was despite the numerous emails sent out to possible interviewees. However, the interviews that were conducted offered valuable data and insights into the area of research. The second limitation of this study was the secondary data that was able to be collected on this topic. This area of study is very specific and niche, owing to these facts it was hard to find a lot of data on the topic. However, the data that was found was valuable and sufficient. Thirdly, the global pandemic changed the way that this study was conducted. However, with online resources and access to platforms such as Skype and Zoom this was easily overcome. Fourthly, the forms of finance section should have been negated from the beginning as it did not have a part in this study, as the funding comes from the US government.

## 5.4    Conclusion and avenues for future study

The aim of this research was to, among other things, look at how states conduct warfare. The premise of this is that warfare has changed, and that technology plays an important role in how warfare changes. This research wanted to look at how AWS could lead to a disruption in traditional nuclear deterrence and how this would subsequently affect strategic stability. It was

found that currently there is no threat to nuclear deterrence and strategic stability. However, there is a high possibility that in the future, with more technological advances, AWS could threaten the US's nuclear deterrence. In terms of future avenues of study, the major threat currently to the US, in terms of technology, is cyber. The US faces a huge problem with cyber intrusions from foreign actors. This means that there need to be continued studies on how cyber can cause or effect international conflict and what type of cyberattack can push a state to war. However, in terms of AI, there needs to be a continued study on how it will affect the US military and international conflict. There also needs to be further study into how AWS could change the future of conflict and if it will ever have the true potential to affect nuclear deterrence.

# 6. Refences

Allen, G. & Chan, T. 2017. Artificial intelligence and national security. *Belfer Center for Science and International Affairs*: 1-110.

Altmann, J. & Sauer, F. 2017. Autonomous weapon systems and strategic stability. Survival, 59(5): 117-142.

Advanced gene editing: CRISPR-Cas9. 2018. Congressional Research Service, December: 1-35.

Artificial intelligence and national security. 2018. *Congressional Research Service*, April: 1-38.

Baker, B. 2019. Orca XLUUV: Boeing's whale of an unmanned sub [Online]. Available: https://www.naval-technology.com/features/boeing-orca-xluuv-unmanned-submarine/ [2020, August 20]

Bentley, M. 2020. Known unknowns: covid-19 and biological warfare [Online]. Available: https://www.e-ir.info/2020/08/08/known-unknowns-covid-19-and-biological-warfare/ [2020, September 22]

Besley, T. & Persson, T. 2019. The rise of identity politics. 1-46.

Berdal, M. 2003. How "new" are "new wars"? global economic change and the study of civil war, global governance. *A Review of Multilateralism and International Organizations*, 9(4): 477-502

Bertand, N. & Lippman, D. 2020. 'It is not science fiction anymore': coronavirus exposes U.S. vulnerability to biowarfare [Online]. Available: https://www.politico.com/news/2020/03/19/coronavirus-biowarfare-terror-136194 [2020, September 2]

Bieri, M. & Dickow, M. 2014. Lethal autonomous weapons systems: future challenges. Security Policy, (164): 1-4.

Biological weapons: a primer. 2001. Congressional Research Service Report for Congress, July: 1-24.

Bode, I., & Huelss, H. 2018. Autonomous weapons systems and changing norms in international relations. Review of International Studies, 44(3), 393-413.

Boeing. 2017. Echo Voyager: extra large unmanned undersea vehicle. Huntington Beach.

Booth, K., 2001. New wars for old. Civil Wars, 4(2): 163-170.

Boulanin, V. 2019. Introduction, in V. Boulanin (ed). *The Impact of artificial intelligence on strategic stability and nuclear risk*. Stockholm: Stockholm International Peace Research Institute. 3-12.

Boulanin, V. 2019. Artificial intelligence: a primer, in V. Boulanin (ed). *The Impact of artificial intelligence on strategic stability and nuclear risk*. Stockholm: Stockholm International Peace Research Institute. 3-12.

Boulanin, V. & Verbruggen, M. 2017. *Mapping the development of autonomy in weapon systems*. Stockholm: Stockholm International Peace Research Institute.

Bracken, P. 2017. The intersection of cyber and nuclear war. *The Strategy Bridge* [Online]. Available: https://thestrategybridge.org/the-bridge/2017/1/17/the-intersection-of-cyber-and-nuclear-war [2020, August 20]

Brynjolfsson, E. & McAfee, A. 2017. The business of artificial intelligence-what it can and cannot do for your organization. *Harvard Business Review*: 1-20.

B-21 Raider [Online]. Available: https://www.northropgrumman.com/what-we-do/air/b-21-raider/ [2020, May 10]

Caplan, A. L., Parent, B., Shen, M. & Plunkett, C. 2015. No time to waste--the ethical challenges created by CRISPR: CRISPR/Cas, being an efficient, simple, and cheap technology to edit the genome of any organism, raises many ethical and regulatory issues beyond the use to manipulate human germ line cells. EMBO reports, 16(11): 1421–1426.

Cebul, D. 2017. The future of autonomous weapons systems: a domain-specific analysis. *New Perspectives on Foreign Policy,* 14: 43-46, Fall.

Christian, M.D. 2013. Biowarfare and bioterrorism. Crit Care Clin, 29(3): 717-756.

Cohen, M.S., Freilich, C.D. & Siboni, G. 2016. Israel and cyberspace: unique threat and response. *International Studies Perspectives*, 17(3): 307-321.

Colby, E.A. 2013. Defining strategic stability: reconciling stability and deterrence, in E.A Colby & M.S Gerson (eds). Strategic stability: contending interpretations. Strategic Studies Institute and U.S Army War College Press. 47-84.

Cordesman, A.H. 2017. President Trump's new national security strategy. *Center for Strategic & International Studies*, December: 1-3.

Corrigan, J. 2017. Three-star general wants AI in every new weapon system. *Defense One,* November. Available: https://www.defenseone.com/technology/2017/11/three-star-general-wants-artificial-intelligence-every-new-weapon-system/142239/ [2020, September 1]

Cummings, M.L. 2017. Artificial intelligence and the future of warfare. *Chatham House The Royal Institute of International Affairs*, January: 1-16.

Dannreuther, Roland. International Security : The Contemporary Agenda. Second ed. 2013. Print.

Davis, C. J. 1999. Nuclear blindness: an overview of the biological weapons programs of the former Soviet Union and Iraq. *Emerging Infectious Diseases*, 5(4): 509-512

Defense primer: U.S. policy on lethal autonomous weapon systems. *Congressional Research Service*. 2020, December.

Department of Defense Directive. 2012. *Autonomy in weapon systems* [Online]. Available: https://www.esd.whs.mil/Portals/54/Documents/DD/issuances/dodd/300009p.pdf [2020 [2020, August 20]

Builder, C.H. 1991. The future of nuclear deterrence. *RAND Corporation*: 1-22.

Estonia denial of service incident. 2007. *Council on Foreign Relations*, May. Available: https://www.cfr.org/cyber-operations/estonian-denial-service-incident [2020, August 30]

Feaver, P. 2017. Five takeaways from Trump's National Security Strategy. *Foreign Policy*, 18 December. Available: https://foreignpolicy.com/2017/12/18/five-takeaways-from-trumps-national-security-strategy/ [2020, August 28]

Flemming, M. 2009. New or old wars? Debating a Clausewitzian future. *Journal of Strategic Studies*, 32(2): 213-241.

Fraze, D. 2016. Cyber Grand Challenge (CCG). Available: https://www.darpa.mil/program/cyber-grand-challenge [2020, July 21]

Fukuyama, F. 2017. What is populism?. Tempus Corporate.

Fukuyama, F. 2018. Against Identity Politics: The New Tribalism and the Crisis of Democracy. Foreign Affairs, 97(5): 90–114.

Futter, A. & Williams, H. 2016. Questioning the holy trinity: why the U.S. nuclear triad still makes sense. *Comparative Strategy*, 35(4): 249-259.

Gates, J. 2016. Is the SSBN deterrent vulnerable to autonomous drones?. *The RUSI Journal*, 161(6): 28-35.

Garcia, Z. 2017. Strategic stability in the twenty-frist century: the challenge of the second nuclear age and the logic of stability interdependence. *Comparative strategy*, 36(4): 354-365.

Gerson, M.S. 2013. The origins of strategic stability: the United States and the threat of surprise attack, in E.A Colby & M.S Gerson (eds). Strategic stability: contending interpretations. Strategic Studies Institute and U.S Army War College Press. 1-46.

Giest, E. & Lohn, J. 2018. How might artificial intelligence affect the risk of nuclear war?. *RAND Corporation*, : 1-26.

Garnham, A. 1987. *Artificial intelligence: An introduction*. Routledge.

Hall, J. 2017. Trump's National Security Strategy release illustrates the ongoing battle between the president's instincts and foreign policy mainstreamers in his administration. [Online]. Available: https://blogs.lse.ac.uk/usappblog/2018/01/04/trumps-national-security-strategy-release-illustrates-the-ongoing-battle-between-the-presidents-instincts-and-foreign-policy-mainstreamers-in-his-administration/ [2020, August 18]

Haner, J. & Garcia, D. 2019. The artificial intelligence arms race: trends and world leaders in autonomous weapons development. *Global Policy*, 10(3): 331-337.

Henderson, E. & Singer, J. 2002. "New wars" and rumors of "new wars". *International Interactions*, 28(2): 165-190.

Hoadley, D.S. & Lucas, N.J. 2018. Artificial intelligence and national security. *Congressional Research Service*: 1-42.

Horowitz, M.C. 2019. When speed kills: lethal autonomous weapon systems, deterrence and stability. *Journal of Strategic Studies*, 42(6): 764-788.

Horowitz, M.C., Allen, G.C., Saravalle, E., Cho, A., Frederick, K. & Scharre, P. 2018. Artificial intelligence and international security. *Center for a New American Security*: 3-21.

Horowitz, M.C., Scharre, P. & Velez-Green, A. 2019. A Stable Nuclear Future? The Impact of Autonomous Systems and Artificial Intelligence. *Cornell University,* (2): 1-35.

Huth, P. 1988. Extended Deterrence and the Outbreak of War. *American Political Science Review*, 82(2), 423-443.

IAEA annual Report. 2009. International Atomic Agency, December: 1-97.

Iran's Qassem Soleimani killed in US air raid at Baghdad airport. *Aljazeera*, January. Available: https://www.aljazeera.com/news/2020/1/3/irans-qassem-soleimani-killed-in-us-air-raid-at-baghdad-airport [2020, August 28]

Johnson, J. 2019. Artificial intelligence & future warfare : implications for international security, defense & security analysis. *Defense & Security Analysis*, 35(2): 147-169.

Johnson, J. 2020a. Artificial intelligence, drone swarming and escalation risks in future warfare. *The RUSI Journal*, 165(2): 26-36.

Johnson, J.S. 2020b. Artificial intelligence: a threat to strategic stability. *Strategic Studies Quarterly*: 16-39, Spring.

Joint all-domain command and control (JADC2). 2020. *Congressional Research Service,* September. Available: https://www.everycrsreport.com/reports/IF11493.html [2020, September 2]

Kahn, M. & Thatcher, S. 2020. Integrated joint all-domain operations (JADO) collaboration strategy full spectrum operations. Lockheed Martin, June: 1-3.

Kaldor, M. 2005. Old wars, cold wars, new wars, and the war on terror. *Int Polit*, 42: 491-498.

Kaldor, M., 2013. In defense of new wars. Stability. *International Journal of Security and Development*, 2 (1): 1-16.

Kallenborn, Z. & Bleek, P.C. 2018. Swarming destruction: drone swarms and chemical, biological, radiological, and nuclear weapons. *The Nonproliferation Review*, 25(5-6): 523-543)

Key informant 1. 2020. Personal interview. 3 July, Durban. [Recording in possession of author]

Key informant 2. 2020. Personal interview. 15 June, Durban. [Recording in possession of author]

Key informant 3. 2020. Personal interview. 18 June, Durban. [Recording in possession of author]

Key informant 4. 2020. Personal interview. 23 June, Durban. [Recording in possession of author]

Key informant 5. 2020. Personal interview. 29 June, Durban. [Recording in possession of author]

Kristensen, H.M. & Korda, M. 2020. Status of world nuclear forces. Federation of American Scientists [Online]. Available: https://fas.org/issues/nuclear-weapons/status-world-nuclear-forces/ [2020, August 23]

Klare, M.T. 2019. Making nuclear weapons menacing again. *Nuclear Arms and Proliferation*, April. Available: https://www.thenation.com/article/politics/31-days-until-the-election/ [2020, June 5]

Korns, S.W. & Kasternberg, J.E. 2009. Georgia's Left Hook. *U.S Army,* April. Available: https://www.army.mil/article/19351/georgias_cyber_left_hook [2020, August 30]

Lachow, I. 2017. The upside and downside of swarming drones. *Bulletin of the Atomic Scientists*, 73:2, 96-101.

Laird, B. 2020. The risk of autonomous weapons systems for crisis stability and conflict escalation in future U.S-Russia confrontations. *RAND Corporation* [Electronic]. Available: https://www.rand.org/blog/2020/06/the-risks-of-autonomous-weapons-systems-for-crisis.html [2020, August 15]

Langer, R. 2011. Stuxnet: dissecting a cyberwarfare weapon. *IEEE Security & Privacy*, 9(3): 49-51.

Leys, N. 2018. Autonomous weapon systems and International Crises. *Strategic Studies Quarterly*: 48-73, Spring.

Lin, P., Allhoff, F. & Rowe, N.C. 2012. Computer ethics, war 2.0: cyberweapons and ethics. *Calhoun: The NPS Institutional Archive*, 55(3): 24-26.

Libicki, M.C. 2009. *Cyberdeterrence and cyberwar*. Santa Monica: RAND Corporation.

Löfflman, G. 2019. America first and the populist impact on US foreign policy. *Survival*, 61(6): 115-138.

Long, A. 2008. Deterrence-from the Cold War to Long War. *RAND Corporation* : 1-85.

Lucas, N.J. 2016. Lethal autonomous weapon systems issues for congress. *Congressional Research Service*, (4): 1-27.

Manhattan Project. 2007. Comprehensive Nuclear-Test-Ban Treaty Organization (CTBTO) [Online]. Available: https://www.ctbto.org/nuclear-testing/history-of-nuclear-testing/manhattan-project/manhattan-project-continued/ [2020, May 15]

Marchlewska, M., Cichocka, A., Panayiotou, O., Castellanos, K. & Batayneh, J. 2017. Populism as identity politics: Perceived ingroup disadvantage, collective narcissism and support for populism. Social Psychological and Personality Science, 9 (2): 151-162.

Martin, D. 2017. New generation of drones set to revolutionize warfare. *60 Minutes* [Online], 6 January. Available: https://www.cbsnews.com/news/60-minutes-autonomous-drones-set-to-revolutionize-military-technology/ [2020, August 14]

Martin, B., Tarraf, D.C., Whitmore, T.C., Deweese, J., Kenney, C., Schmid, J. & DeLuca, P. 2019. Advancing autonomous systems: an analysis of current and future technology for unmanned maritime vehicles. *RAND Corporation*: ix-93.

Mazarr, M.J. 2018. Understanding deterrence. *RAND Corporation*: 1-14.

McGuinness, D. 2017. How a cyber attacks transformed Estonia[Online]. Available: https://www.bbc.com/news/39655415 [2020, August 15]

Mello, P.A. 2010. In search of new wars: the debate about a transformation of war. *Review*, 16(2): 297-309.

Michael, M.J. 2018. Understanding Deterrence. *RAND*, :1-14. Available: https://www.rand.org/content/dam/rand/pubs/perspectives/PE200/PE295/RAND_PE295.pdf [2019, July 30]

Miller, J.N., Fontaine, R. & Velez-Green, A. 2018. Even as U.S-Russian tension have risen, fundamental shifts in the military-technological environment threaten to erode strategic stability between the two nations. *National Interest* [Online], 22 February. Available: https://nationalinterest.org/feature/averting-the-us-russia-warpath-24604 [2020, August 12]

Morgan, P.M. 2003. Deterrence now.

Morgan, P.M. 2012. The state of deterrence in international politics today. *Contemporary Security Policy*, 33(1) 85-107.

Münkler, H. The New Wars. Cambridge, UK: Polity, 2005. Print.

National Security Strategy. 2017. *United States of America*, December: ii- 55.

Navy large unmanned surface and undersea vehicles: background and issues for congress. 2020. Congressional Research Service. October: 1-29.

Newman, E. 2004. The 'new wars' debate: a historical perspective is needed. *Security Dialogue*, 35.2: 173-89.

Ng, A. 2016. What artificial intelligence can and can't do right now. *Harvard Business Review,* November. Available: https://hbr.org/2016/11/what-artificial-intelligence-can-and-cant-do-right-now [2019, February 20]

Noone, G.P. & Noone, D.C. 2015. The debate over autonomous weapon systems. *Case Western Journal of International Law*,47(1): 25-35.

Nuclear Posture Review. 2018. *Department of Defense:*I-74 Available: https://media.defense.gov/2018/Feb/02/2001872886/-1/-1/1/2018-NUCLEAR-POSTURE-REVIEW-FINAL-REPORT.PDF [2020, June 5]

Nuclear weapons: who has what at a glance [Online]. 2019. Arms Control association. Available: https://www.armscontrol.org/factsheets/Nuclearweaponswhohaswhat [2020, May 5]

Ottis, R. 2008. Analysis of the 2007 cyber attacks against estonia from the information warfare perspective. *Proceedings of the 7th European Conference on Information Warfare*: 163.

Payne, K. 2018. Artificial intelligence: a revolution in strategic affairs?. *Survival*, 60(5): 7-32.

Payne, K.B., 2015. US nuclear weapons and deterrence. Air & Space Power Journal, 29(4): 63.

Podvig, P. 2012. The myth of strategic stability. Bulletin of the Atomic Scientists [Online]. Available: https://thebulletin.org/2012/10/the-myth-of-strategic-stability/ [2020, August 23]

Powell, R. 2003. Nuclear Deterrence Theory, Nuclear Proliferation, and National Missile Defense. *International Security*, 27(4), 86-118.

Pernik, P. 2018. The early days of cyberattacks: the cases of Estonia, Georgia and Ukraine, in N. Popescu & S. Secrieru (eds). *Hacks, leaks and disruptions: Russian cyber strategies*. European Union Institute for Security Studies: Chaillot Papers. 53-64.

Pifer, S. 2018. Assessing the U.S. national security strategy. *Brookings* [Online], 25 January. Available: https://www.brookings.edu/testimonies/assessing-the-u-s-national-security-strategy/ [2020, August 28]

Quackenbush, S. 2011. Deterrence theory: Where do we stand?. *Review of International Studies*, 37(2), 741-762.

Reed, B.C., 2014. The history and science of the Manhattan project. Heidelberg: Springer.

Rickli, J. 2019. The destabilizing prospects of artificial intelligence for nuclear strategy, deterrence and stability, in V. Boulanin (ed). *The impact of artificial intelligence on strategic stability and nuclear risk.* Stockholm: Stockholm International Peace Research Institute. 91-98.

Robinson, M., Jones, K. & Janicke, H. 2015. Cyber warfare: issues and challenges. *Computers and security*, 49: 70-94.

Rose, F.A. 2018. Is the 2018 Nuclear Posture Review as bad as the critics claim it is?. Brookings Policy Brief: 1-11.

Saalman, L. 2020. The Impact of Ai on nuclear deterrence: China, Russia, and the United States. *East-West Center*: 1-2.

Sauer, F. 2019. Military application of artificial intelligence: nuclear risk redux in V. Boulanin (ed). *The impact of artificial intelligence on strategic stability and nuclear risk.* Stockholm: Stockholm International Peace Research Institute. 85-91.

Schafer, H. 2005. New wars and religious identity politics, in J. de Santa Ana (Hg.). Religions today. Their Challenge to ecumenical movement. Geneva: WCC. 89-104.

Scharre, P. 2014. Robotics on the battlefield part II. *Center for a New American Security* [Electronic], October: 5-61. Available: https://www.files.ethz.ch/isn/184587/CNAS_TheComingSwarm_Scharre.pdf [2020, August 15]

Scharre, P. 2015. Robots at war and the quality of quantity. *War on the Rocks*. [Online], 26 February. Available: https://warontherocks.com/2015/02/robots-at-war-and-the-quality-of-quantity/ [2020, August 23]

Scharre, P. 2015. Unleash the swarm: the future of warfare. *War on the Rocks*. [Online], 4 March. Available: https://warontherocks.com/2015/03/unleash-the-swarm-the-future-of-warfare/ [2020, August 5]

Scheber, T. 2008. Strategic stability: time for a reality check. *Research Article*, 63(4): 893-915.

Schelling, T.C. 1966. *Arms and Influence.* Bloomsbury: Yale University Press.

Schmidhuber, J. 2014. Deep learning in neural networks: an overview. Technical Report, 61: 85-117.

Schütz, T. & Stanley-Lockman, Z. 2017. Smart logistics for future armed forces. *European Union Institute for Security Studies (EUISS),* November. Available:

https://www.iss.europa.eu/content/smart-logistics-future-armed-forces [2020, September 1]

Schwab, K. 2016. The fourth industrial revolution: what it means, how to respond. [Online], 14 January. Available: https://www.weforum.org/agenda/2016/01/the-fourth-industrial-revolution-what-it-means-and-how-to-respond/ [2019, February 21]

Shane, S. & Wakabayashi, D. 2018. 'The business of war': google employees protest work for the pentagon. *The New York Times*, 2018. Available: https://www.nytimes.com/2018/04/04/technology/google-letter-ceo-pentagon-project.html [2020, September 1]

Sharp, J.M. 2009. Syria: background and U.S. relations. Congressional Research Service, September: 1-19.

Shaw, M. 2000. The Contemporary Mode of Warfare? Mary Kaldor's Theory of New Wars. 7.1: 171-80.

Shi, Z. 2011. *Advanced artificial intelligence*. Singapore ; Hackensack, NJ: World Scientific.

Solis, G.D. 2014. Cyber warfare. *Military Law Review*, 219: 1-53, Spring.

Spiers, E.M. 2010. *A history of chemical and biological weapons*. London: Reaktion Books Ltd.

Stine, D.D. 2009. The Manhattan project, the Apollo program, and federal energy technology R&D programs: a comparative analysis. *Congressional Research Service*, April: 1-10.

Tarraf, D.C., William, S., Edward, P., Brien, A., Gehlhaus, D., Grana, J., Levedahl, A., Leveille, J., Mondschein, J., Ryseff, J., Wyne, Ali., Elinoff, D., Geist, E., Harris, B.N., Hui, E., Kenney, C., Newberry, S., Sachs, C., Schirmer, P., Schlang, D., Smith, V.M., Tinstad, A., Vedula, P. & Warren, K. 2019. The department of defense posture for artificial intelligence: assessment and recommendations. *Rand Coporation*.

Technology for Global Security Reports. 2019. AI and the military: forever altering strategic stability. Available: https://www.tech4gs.org/uploads/1/1/1/5/111521085/ai_and_the_military_forever_altering_strategic_stability__t4gs_research_paper.pdf [2020, August 10]

Theohary, C.A. & Harrington, A.I. 2015. Cyber operations in DOD policy and plans: issues for congress. *Congressional Research Service,* January: 1-32.

Theohary, C.A. & Rollins J.W. 2015. Cyberwarfare and cyberterrorism: in brief. *Congressional Research Service*: 1-12.

Vergun, D. 2019. Nuclear triad important to America's national security. [Online]. Available: https://www.defense.gov/Explore/News/Article/Article/1823014/nuclear-triad-important-to-americas-national-security/ [2020, August 15]

United States strategic approach to the Peoples Republic of China [Online]. 2020. Available: https://www.whitehouse.gov/wp-content/uploads/2020/05/U.S.-Strategic-Approach-to-The-Peoples-Republic-of-China-Report-5.24v1.pdf [2020, August 20]

U.S ground forces robotics and autonomous systems (RAS) and artificial intelligence (AI): consideration for congress. 2018. *Congressional Research Service*, November: 1-42.

US-Iran tension : how confrontation between rival escalated. *Aljazeera*, January. Available: https://www.aljazeera.com/news/2020/1/8/us-iran-tensions-how-confrontation-between-rivals-escalated [2020, August 28]

U.S. Relations with Russia [Online]. 2020. Available: https://www.state.gov/u-s-relations-with-russia/ [2020, August 20]

U.S strategic nuclear forces: background, developments, and issues. 2020. *Congressional Research Service,* 27 April: 1-57.

van Aken, J. & Hammond, E. 2003. Genetic engineering and biological weapons. *EMBO reports,* 4(1).

Wasko-Owsiejczuk, E. 2018. The tenets of Trumpism – from political realism to populism. University of Bialystok, 3(18): 84-91.

Ween, A., Dortmans, P., Thakur, N., & Rowe, C. 2019. Framing cyber warfare: an analyst's perspective. *The Journal of Defense Modeling and Simulation*, *16*(3): 335–345.

Wheelis, M. & Dando, M. 2003. Back to bioweapons. *Bulletin of the Atomic Scientists*, 59(1): 40-46.

Wickham. J.A. 1974. Memorandum for General Wickham: Nuclear weapons employment policy. *NS Archive*. Available: https://nsarchive2.gwu.edu/NSAEBB/NSAEBB173/SIOP-25.pdf [2019, July 22]

Wohlstetter, A., 1958. The delicate balance of terror. *Foreign Aff.*, 37: 211.

Wong, C. 2016. *'Underwater Great Wall': Chinese firm proposes building network of submarine detectors to boost nation's defence.* [Online], 19 May. Available: https://www.scmp.com/news/china/diplomacy-defence/article/1947212/underwater-great-wall-chinese-firm-proposes-building [2020, August 22]

Younger, S.M. 2000. Nuclear weapons in the twenty-first century. Los Alamos National Laboratory, June.

Zagare, F. 1985. Toward a Reformulation of the Theory of Mutual Deterrence. International Studies Quarterly, 29(2), 155-169.

Zanders, J.P., Hart, John. & Kuhlau, F. 2001. Biotechnology and the future of the biological and toxin weapons convention. *Stockholm International Peace Research Institute*, November: 1-11.