# A Stable Nuclear Future? The Impact of Autonomous Systems and Artificial Intelligence

Michael C. Horowitz University of Pennsylvania March 18, 2020 USSTRATCOM Academic Alliance

# The man who saved the world?

КАРТА ЗВЕЗДНОГО НЕБА

S

1

- Current Research Agenda
- Theory: Bias, Capabilities, and Interest in Al
- Early Warning/Command and Control
- Uninhabited Nuclear Platforms
- Conventional Military AI: Impact on Nuclear Stability
- Conclusion





# Key Role of Minerva Research Initiative

- Funding to conduct this research
- Access to key stakeholders
- Relevance for US military power





- Current Research Agenda
- Theory: Bias, Capabilities, and Interest in Al
- Early Warning/Command and Control
- Uninhabited Nuclear Platforms
- Conventional Military AI: Impact on Nuclear Stability
- Conclusion



# Automation, Autonomy, and Artificial Intelligence









# What Is AI?

- Definition: the use of computers to simulate human behavior that requires intelligence
- Methods of Al
  - Symbolic v. Connectionist
  - Machine learning
  - Neural Networks
- Types of AI
  - Narrow
  - General Intelligence
  - Superintelligence



# Al is an Enabler, not a Weapon

- Things AI can do....
  - Direct physical objects
  - Process data
  - Overall information management (decision-making)
- Things AI is not
  - A gun
  - A plane
- Implication: AI is much broader than particular military technologies









#### **Broad**



#### **Dual Use**



#### Low barrier to entry





#### Why Pursue Autonomy or Artificial Intelligence?





# **Brittleness of Autonomous Systems**

- Narrow AI systems trained to do one thing
- Example: Alpha Go
- Challenges:
  - How do you train them (with what data)?
  - Limited potential area of operation

# Trust, Confidence, and Al (1)

#### Trust Gap

- Inability to trust machines to do work of people
- Unwillingness to deploy or properly use systems
- Example: Ground Tactical Air Controllers (MacDonald and Schneider)

#### **Automation Bias**

- Delegation of cognitive judgment to machine – trusting too much
- Failure to question algorithms if they make mistakes
- Example: Air France Crash
- Example: Patriot Missile fratricide



## Trust, Confidence, and AI (2)



**Time Since System Introduction** 



# **Second Strike Capabilities And Autonomy**

Interest in Autonomous Systems in Nuclear Domain

High

Low Low High **Confidence in Second Strike Capabilities** 



#### **Key Driver: Competitive Pressure**



- Current Research Agenda
- Theory: Bias, Capabilities, and Interest in Al
- Early Warning/Command and Control
- Uninhabited Nuclear Platforms
- Conventional Military AI: Impact on Nuclear Stability
- Conclusion



Early Warning + Command & Control



# **Existing Early Warning: Automated**

- Long-range radar or satellite based alert systems
- Rapid-retargeting capability
- Communication rockets to transmit launch codes

#### **Example: Petrov Incident**

#### **Example: Soviet Perimeter System**



# **Autonomous Early Warning?**

- Theoretical Benefits:
  - Early detection: Buys time for decision-makers
  - Reliability
- Theoretical Downsides
  - Loss of human judgment/lack of human judgment
  - Brittleness of algorithms -> false alarms



- Current Research Agenda
- Theory: Bias, Capabilities, and Interest in Al
- Early Warning/Command and Control
- Uninhabited Nuclear Platforms
- Conventional Military AI: Impact on Nuclear Stability
- Conclusion



#### **Uninhabited Nuclear Platforms**

- Theoretical Benefits:
  - Endurance
  - Reliability
- Theoretical Downsides
  - Cannot maintain positive human control
  - Consequences of accidents, hacking, spoofing
  - Brittleness of algorithms -> false alarms

## **Example: US Military**

US Air Force 2013 report, Remotely Piloted Aircraft (RPA) Vector. [N]uclear strike may not be technically feasible unless safeguards are developed and even then may not be considered for [unmanned aircraft systems] operations.

General Robin Rand, head of Air Force Global Strike Command (2016): **We're planning on [the B-21] being manned. ... I like** the man in the loop ... very much, particularly as we do the dual capable mission with nuclear weapons. Example: Russian military

- Perception of conventional + nuclear inferiority
- 2012 statement
- Ocean Multipurpose System 'Status-6

# **Example: North Korean Military**

- Relatively newer nuclear power
- Conventional military
  inferiority
- Fear of decapitation
- Repressive regime

North Korea, in theory, should have greater interest in autonomy of all kinds, especially uninhabited nuclear vehicles



- Current Research Agenda
- Theory: Bias, Capabilities, and Interest in Al
- Early Warning/Command and Control
- Uninhabited Nuclear Platforms
- Conventional Military Al: Impact on Nuclear
  Stability
- Conclusion



# **Surveillance and Counterforce**

8 22

28

Ň,

80

#### What Would Militaries Use AI For?



# Fighting At Machine Speed: Crisis Stability

Speed Trap ahead  Compressed decision cycles - Offense - Defense **Fear of losing** quickly - First strike stability - Launch posture



- Current Research Agenda
- Theory: Bias, Capabilities, and Interest in Al
- Early Warning/Command and Control
- Uninhabited Nuclear Platforms
- Conventional Military AI: Impact on Nuclear Stability
- Conclusion



# Conclusion

- The less secure the second strike capabilities, the more a country is likely to consider autonomous systems within their nuclear weapons complex
- Some risk associated with greater automation in early warning
- Potentially large risk associated with impact of conventional military uses of autonomy on crisis stability

