

# Artificial Intelligence and Threats

Jeff Alstott



Office of the Director of National Intelligence

I A R P A

BE THE FUTURE





# Doing Bad Things...

- With AI
- To AI
- Because (of) AI

Artificial Intelligence and Threats



# Doing Bad Things...

- With AI
  - Autonomous weapons





Office of the Director of National Intelligence

IARPA

BE THE FUTURE



INTELLIGENCE ADVANCED RESEARCH PROJECTS ACTIVITY (IARPA)



# Doing Bad Things...

- **With AI**
  - Autonomous weapons
  - Psyops

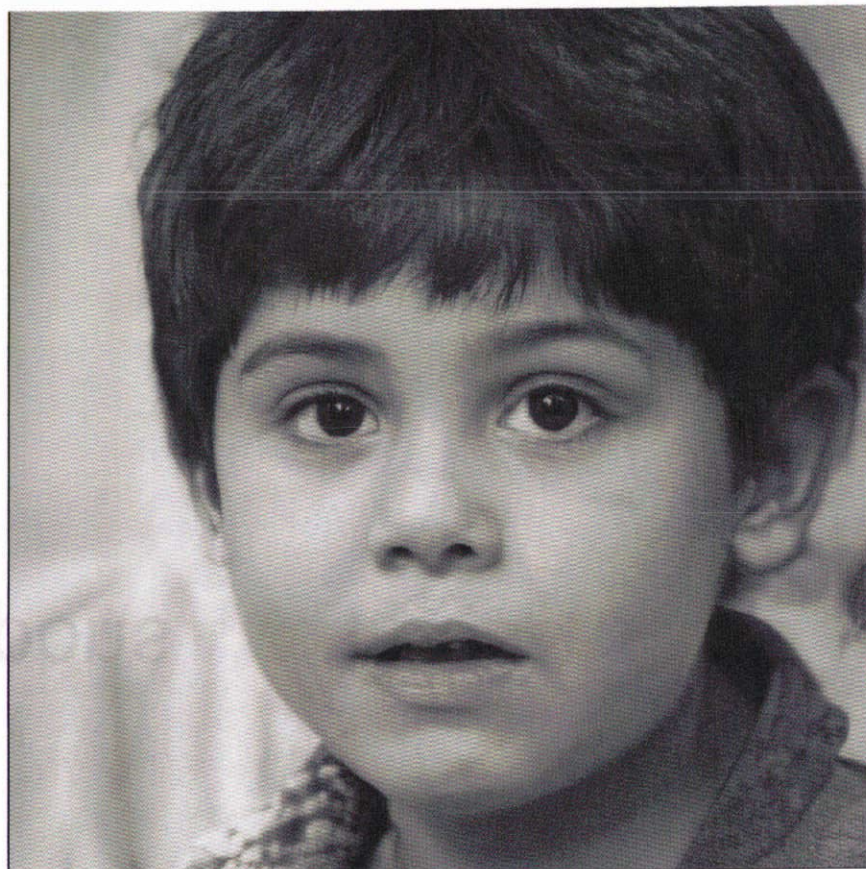




Office of the Director of National Intelligence

IARPA

BE THE FUTURE



thispersondoesnotexist.com

INTELLIGENCE ADVANCED RESEARCH PROJECTS ACTIVITY (IARPA)





### *Do Electric Neural Nets Dream Of Anime Shows?*

\_Bishounen Kaigan\_: This is the story of a temporary manga staged in a near future, to determine the fate of himself as the world rests in a state of absolute chaos. There is a world that has been destroyed by weapons of mass destruction and no one can escape alive. There are only a handful of people that can communicate with it. Kai is one of the twelve shikigami called "Mushi King". This is the story of Kai, a young man with a tragic past. (Source: MU) In the first episode of this anime [The Last Episode] is a young girl named Kiba to herself, who meets a ghost called Tifa and is the only one who can hear her as she goes about things in Tokyo. From there, Kiba's relationship with her ghost becomes a serious and intense problem. · A series of events begin to occur after the second episode ends, in which the ghost(s) begins to look at Kiba from various angles. From all this, Kiba becomes the ghost of Kai, and he

thiswaifudoesnotexist.com







Office of the Director of National Intelligence

I A R P A

BE THE FUTURE



# Doing Bad Things...

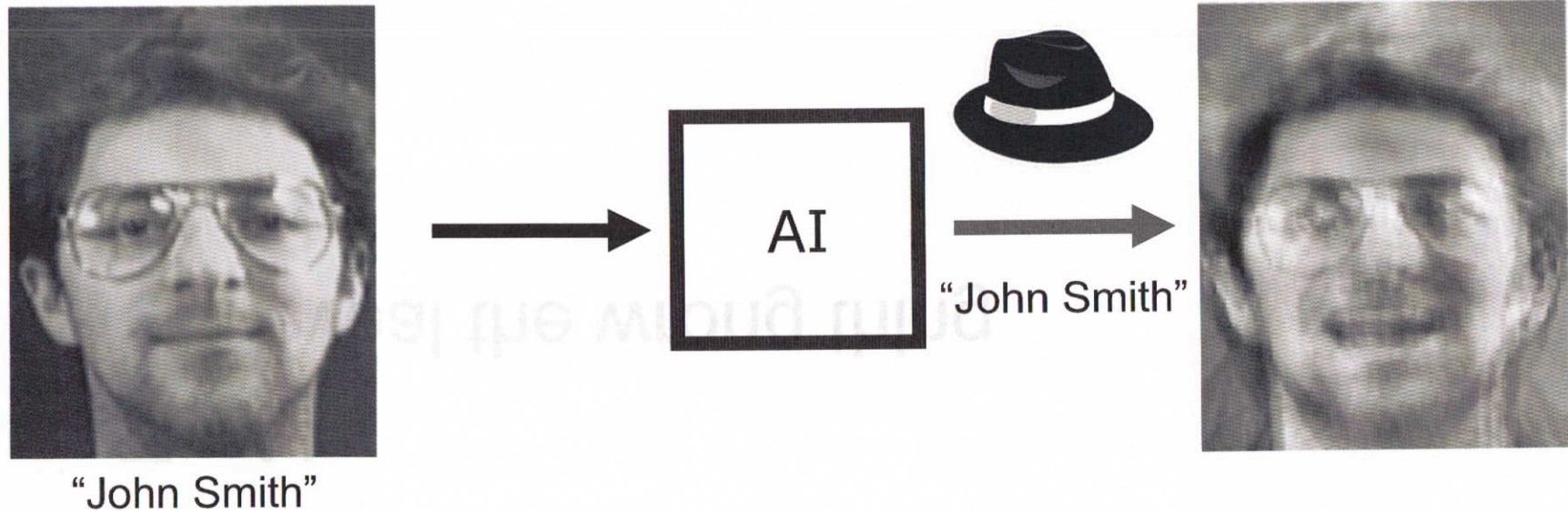
- With AI
- To AI
  - Reveal the wrong thing

INTELLIGENCE ADVANCED RESEARCH PROJECTS ACTIVITY (IARPA)





# Let's say you're making a face identifier...



Adversaries can **invert** AI models to learn about the training data

Model inversion attacks work *even with black box access*  
(e.g. if the model executable is encrypted)





# Doing Bad Things...

- With AI
- To AI
  - Reveal the wrong thing
  - Do the wrong thing



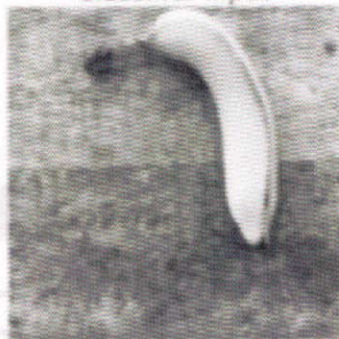


# Let's say you're making an image classifier...

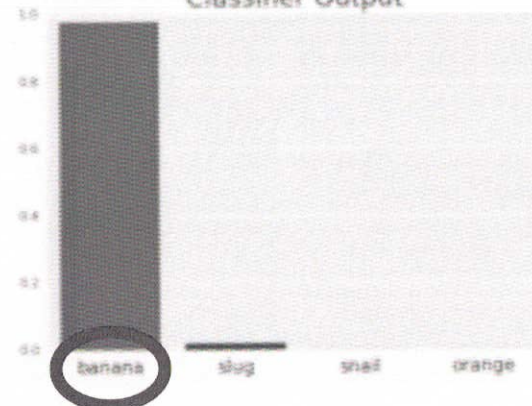
place sticker on table



Classifier Input



Classifier Output



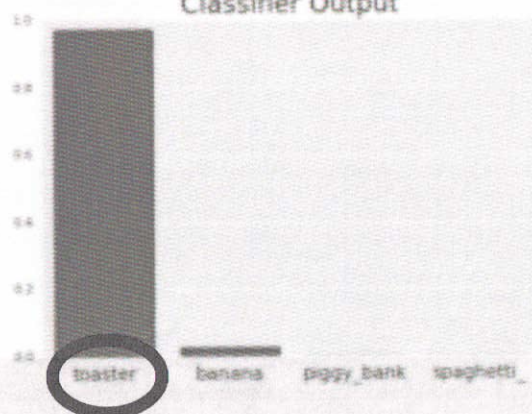
Like humans, AIs are subject to camouflage, illusions, etc.

Existing AI techniques are full of **adversarial examples**, which adversaries can find and exploit

Classifier Input



Classifier Output





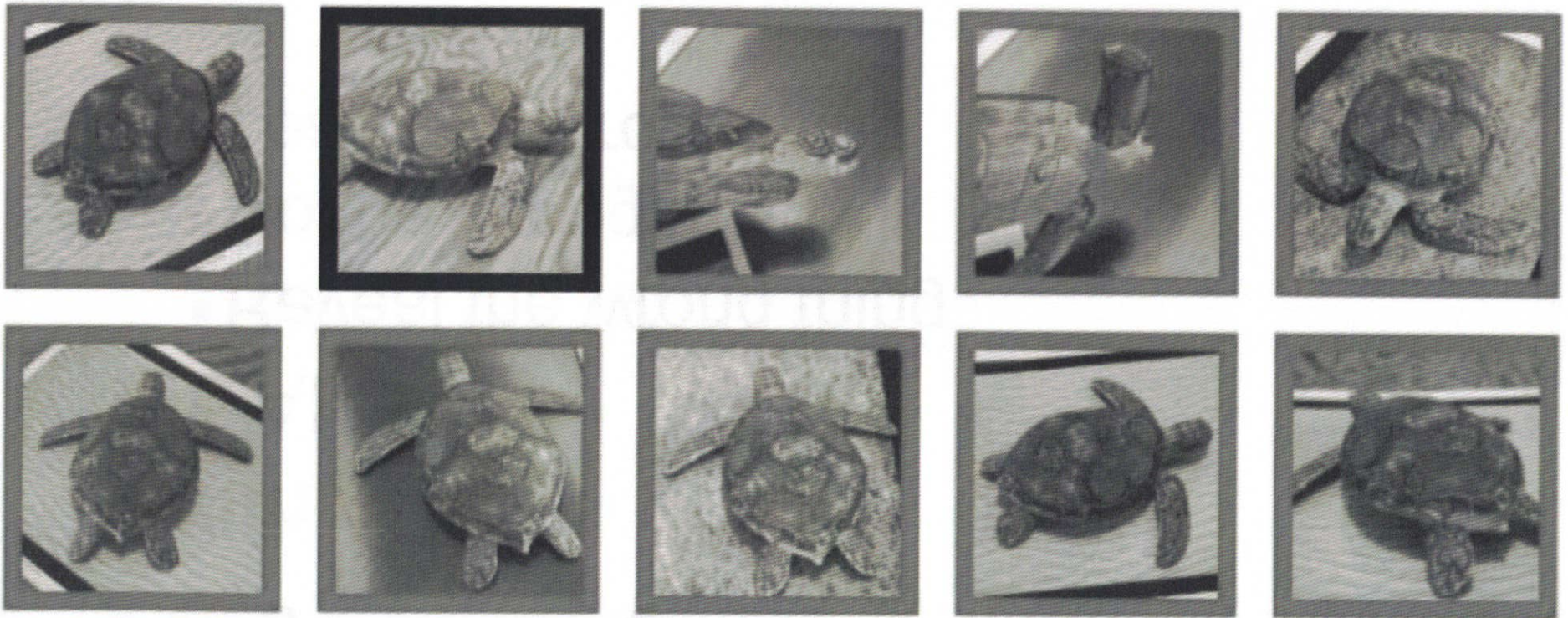


Office of the Director of National Intelligence

IARPA  
BE THE FUTURE



# Let's say you're making an image classifier...







# Doing Bad Things...

- With AI
- To AI
  - Reveal the wrong thing
  - Do the wrong thing
  - Learn the wrong thing



# Let's say you're making a self-driving car...



Label:  
**Stop sign**



Label:  
**Speed limit sign**



Adversaries can insert **Trojans** into AIs through small manipulations to training data

AI can remain infected with Trojan *even after transfer learning*





# Doing Bad Things...

- With AI
- To AI
  - Reveal the wrong thing
  - Do the wrong thing
  - Learn the wrong thing
  - Value the wrong thing
  - ...



# Doing Bad Things...

- **With AI**
- **To AI**
- **Because (of) AI**
  - OODA Loop Tightening
  - Online learning
  - ...





Office of the Director of National Intelligence

I A R P A

BE THE FUTURE



# Doing Bad Things...

- **With AI**
  - Autonomous weapons
  - Psyops
- **To AI**
  - Reveal the wrong thing
  - Do the wrong thing
  - Learn the wrong thing
  - Value the wrong thing
- **Because (of) AI**
  - OODA Loop Tightening
  - Online learning
- ...?