

# Al as a Disruptive Opportunity and Challenge for Security

#### Antonio Kung CTO Trialog 25 rue du Général Foy 75008 Paris www.trialog.com



#### Security & privacy background

#### Standardisation

12 June 2018

- 27550 Privacy engineering (editor)
- 27030 Security and privacy guidelines for IoT (co-editor)
- 27570 Privacy guidelines for smart cities (editor)
- IPEN wiki (ipen.trialog.com)
- EIP- Smart Cities and Communities
  - Citizen approach to data
- PDP4E Privacy and Data Protection for Engineers
- Create-IoT (IoT large scale pilots)



# IoT background



- Energy, Social and care, e-Mobility, ITS
- AOTI member
- EIP-Active healthy ageing
  - Recommendations for interoperability and standardisation (2015)

#### AI background

- BDVA member
  - ACCRA Agile Co-Creation of Robots for Ageing

#### Alliar Interr

Alliance for Internet of Things Innovation







#### Discussion with Ivo Emanuilov (KUL Citip / IMEC)

- Poisoned AI: towards data liability. LICT Workshop Autonomous Systems, 31 May 2018
  - https://set.kuleuven.be/LICT/map-events-workshops2/lict-workshop-on-autonomoussystems
- Al Malicious use report. February 2018
  - https://maliciousaireport.com
- Asilomar Al Principles. January 2017
  - https://futureoflife.org/ai-principles/
- Building ethics into AI Kathy Baxter blog March April 2018
  - https://medium.com/salesforce-ux/how-to-build-ethics-into-ai-part-i-bf35494cce9
  - https://medium.com/salesforce-ux/how-to-build-ethics-into-ai-part-ii-a563f3372447
- Big data value association
  - Task force 5: policy & societal implications
    - Freek Bomhof TNO, Natalie Bertels KUL Citip IMEC
  - Working group on AI transparency

# **Artificial Intelligence**

- Artificial intelligence
  - "intelligence" demonstrated by machines
  - mimics "cognitive" functions that humans associate with other human minds, such as "learning" and "problem solving"
- •

12 June 2018

Stochastic AI (Machine learning) vs deterministic AI

- Stefan Ticu (Rocket labs)
  - https://yseop.com/blog/artificialintelligence-machine-learning-vsdeterministic/
- Stochastic AI automates 100% of an activity with 70% accuracy
- Deterministic AI automates 70% of an activity with 100% accuracy

- Al applications
  - Automatic speech recognition
  - Machine translation
  - Spam filters
  - Search engines
  - Autonomous cars
  - Robots for elderly people
  - Autonomous drones



Maximum Impact	Must be		At	osolutely
Significant Impact	avoided reduced	or	av	oided or reduced
Limited Impact				
Negligible Impact	These ris be taken	sks may		Must be reduced
	Negligible Likelihood	Limited Likelihood	Significant Likelihood	Maximum Likelihood

- Security and privacy threat/breach risk level:
  - Likelihood
  - Impact
  - Many versions of risk maps
    - More levels
    - Different ways of calculating.
       Exemples
      - NIST privacy engineering
      - ETSI TVRA
    - This map is from CNIL guidelines









# TRIALOG

# Dual Use 1 **Malicious Al**

Malicious Al Report

The I of Ar Fore and I	Malici tificia castir Mitiga	ious al Int ng, P ation	Centre for the Study of Eristential Use elligent reventi	Univers Cambrid	ity of Cem ge Nem Set	stor for a / American uurity	Electronic Frontier Foundation	ry 2018	Miles Bri Ben Gart Hyrum A Carrick F Sebastia Michael
								$\mathbf{i}$	1 Correspondent miles.brunn ox.ac.uk Future of University
									2 Correspo sa4788cam.s Centre for Existential
		Ι		Ι		Ι		Ι	3 OpenAI 4 Open Phi
Ι		Ι		Ι		Ι		Ι	5 Electron Foundation
Ι								—	6 Future o Institute, Oxford 7 Future o Institute, Oxford; Yai
								+	8 Center f Security
								+	Authors and
								+	in order of

undage Shahar Avin[2] Jack Clark[3] Helen Toner[4] Peter Eckersley[ finkel Allan Dafoe 7 Paul Scharre 7 Thomas Zeitzoff 8 Bobby Filar 10 Anderson[11] Heather Roff[12] Gregory C. Allen[13] Jacob Steinhardt Flynn 151 Seán Ó hÉigeartaigh 161 Simon Beard 1 Haydn Belfield an Farquhar [18] Clare Lyle [20] Rebecca Crootof Owain Evans Page 23 Joanna Bryson 24 Roman Yampolskiv 25 Dario Amodei 20

ding author 9 American University dage@philosophy 10 Endgame Humanity Institute, of Oxford: Arizona ersity 11 Endgame onding author ac.uk the Study of 12 University of Oxford/

Arizona State University/New Risk, University America Foundation 13 Center for a New American

ilanthropy Project 14 Stanford University

of Humanity

or a New American

18 Centre for the Study of Existential Risk, University of Cambridge

19 Future of Humanity Institute, University of Oxford

20 Future of Humanity Institute, University of

21 Information Society Project, Yale University

22 Future of Humanity Institute, University of Oxford

23 OpenAI

24 University of Bath

Intelligence, University of 25 University of Louisville 26 OpenAI

17 Centre for the Study of Existential Risk, University of Cambridge

Center for a New American

Security

Design Direction contribution

by Sankalp Bhatnagar and Talia Cotton

EXISTENTIAL RISK

Security

Cambridge

15 Future of Humanity Institute, University of

16 Centre for the Study of Existential Risk and Centre for the Future of









#### Expansion of existing threats

- Expanding phishing
- Increasing willingness to carry out attacks
  - increasing anonymity and increasing psychological distance
- Robotics progress
- Introduction of new threats
  - Mimicking voice
  - New AI capabilities imply new threats
    - Autonomous cars VS image of a stop sign changed
    - Swarm of autonomous systems VS attack on a server to control the swarm



# Malicious AI: Increasing Likelihood

					(TVRA) Th	reat Vulner	ability Risk Analysis
					Attack	factor	Malicious AI assistance
Maximum Impact	Must be			osolutely	Time	<= 1 day <= 1 week <= 1 month <= 3 months <= 6 months > 6 months	AI attack creation assistant
Significant	avoided reduced	lor	avoided or reduced		Expertise	Layman Proficient Expert	
Limited					Knowledge	Public Restricted Sensitive Critical	AI based learning of vulnerabilities
Impact	Impact       Negligible Impact       Descrists       Impact		Must be reduced		Opportunity	Unnecessary Easy Moderate Difficult Nont	AI based creation of opportunities
Negligible Impact					Equipment	Standard Specialised Bespoke	Lower cost
	Negligible	Limited	Significant	Maximum	Asset Impact	Low Medium High	Al analysis of impact
	Likelihood	Likelihood	Likelihood	Likelihood	Intensity	Single intensity Moderate intensity High intensity	AI based swarm attack



## Malicious AI Report : Categories of threats

	Automation of social engineering attacks. Mimick a person					
	Automation of vulnerability discovery. Historical patterns of code vulnerabilities are used to speed up the					
	discovery of new vulnerabilities, and the creation of code for exploiting them					
	More sophisticated automation of hacking.					
<b>Digital security</b>	Human-like denial-of-service					
0 /	Automation of service tasks in criminal cyber-offense (payment processing / dialog with ransomware) victims					
	Prioritising targets for cyber attacks using machine learning					
	Criminal Training – Data poisoning					
	Black-box model extraction of proprietary AI system capabilities					
	Terrorist repurposing of commercial AI systems (e.g. drones)					
	Endowing low-skill individuals with previously high-skill attack capabilities					
<b>Physical security</b>	Increased scale of attacks					
	Swarming attacks – Drones as weapons					
	Attacks further removed in time and space					
	State use of automated surveillance platforms to suppress dissent					
	Fake news reports with realistic fabricated video and audio					
Political security	Automated, hyper-personalised disinformation campaigns					
i oncical security	Automating influence campaigns					
	Denial-of-information attacks					
	Manipulation of information availability					



# Malicious AI Report : Example of measures

	Consumer awareness (e.g. education)					
	Governments policies and research (e.g. incentives for source code analysis)					
Digital security	Industry centralization capability (e.g. centralised spam filters)					
	Attacker incentives (e.g. identifying source of attack)					
	Technical cybersecurity defense (e.g. NIST based improved practice)					
	Hardware manufacturers (e.g. drones)					
	Hardware distributors (e.g. controlled sales)					
	Software supply chain					
Physical security	Robot users (e.g. using licence)					
	Governments policies (e.g. restricted use of robots)					
	Physical defenses (e.g. new generation of radars)					
	Payload control (e.g. AI based payload analysis)					
	Technical tools (e.g. fake news detection)					
Political cocurity	Pervasive use of security measures (e.g. more encryption)					
Political Security	General interventions to improve discourse (e.g. education)					
	Media platforms (e.g. integrating assessment capabilites)					



## Malicious AI Report : Research Topics

	Red teaming
Learning from and with	Formal verification
the Cybersecurity	Responsible disclosure of AI vulnerabilities
	Forecasting security-relevant capabilities
Community	Security tools
	Secure hardware
	Pre-publication risk assessment in technical areas of special concern
Exploring Different	Central access licensing models
Openness Models	Sharing regimes that favor safety and security
	Other norms and institutions that have been applied to dual-use technologies
	Education
Promoting a Culture of	Ethical statements and standards
Responsibility	Whistleblowing measures
	Nuanced narratives
	Privacy protection
Developing Technological	Coordinated use of AI for public-good security
and Policy Solutions	Monitoring of AI-relevant resources
	Other legislative and regulatory responses



# Dual Use 2 Al to improve IoT security and privacy





	AI Assistance	
Identify	Big data risk analysis	
Protect Pattern analysis and desig		
Dotoct	Off line anomaly analysis	
Delect	On line anomaly detection	
	Response big data analysis	
Respond	Training operators	
	Assisting operations	
Bacovar	Training operators	
Recover	Assisting operations	



# System life cycle process (ISO/IEC 15288)

	Acquisition				Concorns to	
Agreement process	Supply		Business or mission			
	Life cycle model		analysis		<ul> <li>integrate</li> <li>Ethics impact assessment</li> <li>Bias management</li> <li>Transparency</li> <li>Example of</li> <li>ISO/IEC 27550</li> <li>Privacy</li> <li>engineering</li> </ul>	
	management		Stakeholder needs and			
Organizational	Infractructure management		requirements definition			
nroject-enabling	Porfolio management		System requirements			
project-enabiling	Human resource		definition			
processes	management	Tochnical	Architecture definition			
	Quality management		Design definition			
	Knowledge management	Technical	System analysis			
	Project planning	processes	Implementation			
	Project assessment and		Integration			
Technical	control		Verification			
Technical	Decision management		Transition			
management	Risk management		Validation			
process	Configuration management		Operation			
	Information management		Maintenance			
	Measurement		Disposal			
	Quality assurance					

#### Example: Cybersecurity situation awareness learning



18

TRIALOG



## Example: Conformity Learning





# Asilomar AI Principles (Beneficial AI)

	1	Dessereb Cool	
	1	Research Goal	Create beneficial intelligence.
Research	2	Research Funding	AI systems robust – Growth through automation - Update legal systems with AI – Align AI with set of values
issues	3	Science-Policy Link	Exchange between AI researchers and policy-makers
	4	Research Culture	Cooperation, trust, and transparency among researchers and developers of AI.
	5	Race Avoidance	Teams developing AI systems should actively cooperate to avoid corner-cutting on safety standards.
	6	Safety	AI systems should be safe and secure
	7	Failure Transparency	If an AI system causes harm, it should be possible to ascertain why.
	8	Judicial Transparency	AI based judicial decision-making auditable by competent human authority.
	9	Responsibility	Designers and builders of advanced AI systems responsible
	10	Value Alignment	Autonomous AI systems goals and behaviors aligned with human values
Ethics	11	Human Values	Al systems compatible with ideals of human dignity, rights, freedoms, and cultural diversity.
and	12	Personal Privacy	People control data
Values	13	Liberty and Privacy	Application of AI to personal data must not curtail people's liberty.
	14	Shared Benefit	AI technologies should benefit and empower as many people as possible.
	15	Shared Prosperity	The economic prosperity created by AI should be shared broadly, to benefit all of humanity.
	16	Human Control	Humans should choose how and whether to delegate decisions to AI systems
	17	Non-subversion	Respect and improve social and civic processes on which the health of society depends.
	18	AI Arms Race	Avoiding arms race in lethal autonomous weapons
	19	Capability Caution	Avoid strong assumptions regarding upper limits on future AI capabilities.
Longer-	20	Importance	Advanced AI planned for and managed with commensurate care and resources.
term	21	Risks	Risks posed by AI systems subject to planning and mitigation efforts commensurate with their expected impact.
Issues	22	<b>Recursive Improvement</b>	Al systems designed to recursively self-improve / self-replicate subject to strict safety and control measures
	23	Common Good	Superintelligence developed in the service of widely shared ethical ideals, and for the benefit of all humanity



# Principles for Ethics into AI (Kathy Baxter blog)

Create an	Build a Diverse Team	Recruit for a diversity of backgrounds and experience to avoid bias and feature gaps.		
ethical	Cultivate an Ethical Mindset.	Ethics is a mindset, not a checklist. Empower employees to do the right thing.		
culture	Conduct a Social Systems Analysis	Involve stakeholders at every stage of the product development lifecycles to correct for the impact of systemic social inequalities in AI data.		
Po	Understand Your Values	Examine the outcomes and trade-off of value-based decisions.		
De transparant	Give Users Control of Their Data	Allow users to correct or delete data you have collected about them		
transparent	Take Feedback	Allow users to give feedback about inferences the AI makes about them.		
	Understand the Factors Involved	Identify the factors that are salient and mutable in your algorithm(s)		
	Prevent Dataset Bias	Identify who or what is being excluded or overrepresented in your dataset, why they are excluded, and how to mitigate it.		
Remove	Prevent Association Bias	Determine if your training data or labels represent stereotypes (e.g., gender, ethnic) and edit them to avoid magnifying them.		
exclusion	Prevent Confirmation Bias.	Determine if bias in the system is creating a self-fulling prophecy and preventing freedom of choice.		
	Prevent Automation Bias	Identify when your values overwrite the user's values and provide ways for users to undo it.		
	Mitigate Interaction Bias	Understand how your system learns from real-time interactions and put checks in place to mitigate malicious intent.		



	What	How
Prevent dataset bias	Majority of data set represented by one group of users.	Cost-sensitive learning Changes in sampling methods Anomaly detection
	Statistical patterns invalid within a minority group.	Algorithms for different groups rather than one-size-fits-all.
	Categories or labels oversimplify data points and may be wrong for some percentage of the data.	Judgement about someone is identified as fair same judgement made in a different demographic group (e.g., if a woman were a man)
	Identify who is being excluded and the impact on your users as well as your bottom line. Context and culture matters but it may be impossible to "see" it in the data.	Identify the unknown unknowns (unidentified risks); See https://www.pmi.org/learning/library/char acterizing-unknown-unknowns-6077.



## Other issues

#### **Dual use**

- Trustworthiness
  - Al to help trustworthiness
    - AI-based trust framework assessment
  - Al to prevent trustworthiness
- Transparency
  - Al to help transparency
  - Al to prevent transparency
- Ethics
  - AI to help ethical impact assessment
  - AI to prevent ethical impact assessment
- Conformity
  - Al to help conformity
  - Al to prevent conformity

#### Life cycle process

- Integration of with model system and software enginering capabilties
  - Model driven engineering
  - Ontology



## **Consensus on policies**

Autonomy level definition

# TRIALOG

# Recommendations

#### **Dual Use (from Malicious AI report)**

- Policy makers / Researchers collaboration
- AI researchers to address dual-use concerns
- Best practices & methods to address dual-use concerns

#### Lifecycle process

- Ethical impact assessment
- Ethical-by-design AI engineering

#### **Consensus on policies**

- Best available techniques
  - consensus-building with numerous stakeholders underpinned by sound techno-economic information
  - e.g. RFId or smart grid
    - http://eur-lex.europa.eu/legalcontent/EN/TXT/PDF/?uri=CELEX:32012 D0119&from=EN









25



www.trialog.com

# **Questions?**



ENABLING INNOVATION SINCE 1987