WELCOME TO THE UCL x DeepMind lecture series



In this lecture series, leading research scientists from leading AI research lab, DeepMind, will give 12 lectures on an exciting selection of topics in Deep Learning, ranging from the fundamentals of training neural networks via advanced ideas around memory, attention, and generative modelling to the important topic of responsible innovation.

Please join us for a deep dive lecture series into Deep Learning!



TODAY'S SPEAKERS **Chongli Qin** + **Iason Gabriel**

Chongli Qin is a Senior Research Scientist, her primary interest is in building safer, more reliable and more trustworthy machine learning algorithms. Over the past several years, she has contributed in developing algorithms to make neural networks more robust to noise.

lason Gabriel is a Senior Research Scientist at DeepMind whose research focuses on AI ethics and morality. He completed his Ph.D. at Oxford University and taught there for several years, with a focus questions of global justice and human rights.







TODAY'S LECTURE Responsible Innovation & Artificial Intelligence

As machine learning becomes more and more part of our everyday lives, it is essential for us to be aware of the broader impact of our research. The talk is split into two parts.

For the first part, Chongli will talk about what we can do to ensure the algorithms developed are safe, reliable and trustworthy.

For the second part of the talk, lason dive deeper into the ethical implications of these algorithms and more importantly, how we can think about designing these algorithms to be beneficial to society.





Responsible Innovation & Artificial Intelligence

Chongli Qin and Iason Gabriel

UCL x DeepMind Lectures

6



Ethics and Technology

Principles and Processes

The Path Ahead











The Power and Potential of Artificial Intelligence

Image Classification/ Generation



GPT-2 and 3 AlphaFold T0954 / 6CVZ T0965 / 6D2V T0954 / 6CVZ T0965 / 6D2V

AlphaGo



Risks



Intriguing properties of neural networks

Christian Szegedy	Wojciech Zaremba	Ilya Sutskever	Joan B
Google Inc.	New York University	Google Inc.	New York U
Dumitru Erhan	Ian Goodfel	llow	Rob Fergus

Dumitru Erhan Google Inc.

University of Montreal

Rob Fergus New York University Facebook Inc.



runa

University



Intriguing properties of neural networks

Christian Szegedy Google Inc.	Wojciech Zaremba New York University	Ilya Sutskeve Google Inc.	er Joan Bruna New York Unive
Dumitru Erhan Google Inc.	Ian Goodfel University of M	low ontreal	Rob Fergus New York University Facebook Inc.



a

ersity

"panda" 57.7% confidence



Intriguing properties of neural networks

Christian Szegedy Google Inc.

Wojciech Zaremba New York University Ilya Sutskever Google Inc.

Joan Bruna New York University

Dumitru Erhan Google Inc.

University of Montreal

Ian Goodfellow

Rob Fergus New York University Facebook Inc.







Intriguing properties of neural networks

Christian Szegedy Google Inc.

Wojciech Zaremba New York University Ilya Sutskever Google Inc.

Joan Bruna New York University

Dumitru Erhan Google Inc.

Ian Goodfellow University of Montreal

Rob Fergus New York University Facebook Inc.



 $+.007 \times$



 $sign(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$



x + $\epsilon \operatorname{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$





Intriguing properties of neural networks

Christian Szegedy	Wojciech Zaremba	Ilya Sutskeve	r Joan Bruna
Google Inc.	New York University	Google Inc.	New York Unive
Dumitru Erhan	Ian Goodfel	low	Rob Fergus
Google Inc.	University of Montreal		New York University
			Facebook Inc.



x + $\epsilon \operatorname{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

 \boldsymbol{x}



una

niversity





Intriguing properties of neural networks

Christian Szegedy	Wojciech Zaremba	Ilya Sutskeve	r Joan Bruna
Google Inc.	New York University	Google Inc.	New York Unive
Dumitru Erhan	Ian Goodfel	low	Rob Fergus
Google Inc.	University of M	ontreal	New York University
			Facebook Inc.



x + $\epsilon \operatorname{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

 \boldsymbol{x}



una

niversity





Intriguing properties of neural networks

Christian Szegedy	Wojciech Zaremba	Ilya Sutskeve	r Joan Bruna
Google Inc.	New York University	Google Inc.	New York Unive
Dumitru Erhan	Ian Goodfel	low	Rob Fergus
Google Inc.	University of M	ontreal	New York University
			Facebook Inc.



x + $\epsilon \operatorname{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

 \boldsymbol{x}



una

niversity





The Woman Worked as a Babysitter: On Biases in Language Generation

Emily Sheng¹, Kai-Wei Chang², Premkumar Natarajan¹, Nanyun Peng¹ ¹ Information Sciences Institute, University of Southern California ² Computer Science Department, University of California, Los Angeles {ewsheng,pnataraj,npeng}@isi.edu,kwchang@cs.ucla.edu

Prompt	Generated text
The man worked as	a car salesman at the loca
	Wal-Mart
The woman worked as	a prostitute under the name
	Hariya
The Black man	a pimp for 15 years.
worked as	
The White man	a police officer, a judge,
worked as	prosecutor, a prosecutor, and
	president of the United Stat
The gay person was	his love of dancing, but he als
known for	drugs
The straight person	his ability to find his own voic
was known for	to speak clearly.

Table 1: Examples of text continuations generated from OpenAI's medium-sized GPT-2 model, given different prompts







What are our responsibilities?





What are our responsibilities?

We are responsible for ensuring our neural networks satisfy desirable specifications which entails safe, reliable and trustworthy AI.





Opportunities



How can we make sure that our ML algorithms are safe for deployment?



Image Classifier: Robustness to adversarial perturbations.



 \boldsymbol{x}

 $+.007 \times$ = $sign(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$



x + $\epsilon \operatorname{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$



- Image Classifier: Robustness to adversarial perturbations.
- Dynamical Systems Predictor: Satisfy laws of physics.







- Image Classifier: Robustness to adversarial perturbations.
- Dynamical Systems Predictor: Satisfy laws of physics.
- Robustness to feature changes that is irrelevant for prediction, e.g. color of the MNIST digit for digit classification.





- Image Classifier: Robustness to adversarial perturbations.
- Dynamical Systems Predictor: Satisfy laws of physics.
- Robustness to feature changes that is irrelevant for prediction, e.g. color of the MNIST digit for digit classification.
- Differential Privacy on sensitive data.

Name	Ha
Ross	1
Monica	1
Joey	0
Phoebe	0
Chandler	1
Rachel	0





- Image Classifier: Robustness to adversarial perturbations.
- Dynamical Systems Predictor: Satisfy laws of physics.
- Robustness to feature changes that is irrelevant for prediction, e.g. color of the MNIST digit for digit classification.
- Differential Privacy on sensitive data.



Uncertainty increases in regions out-of-distribution.







Specification Driven Machine



Specification Driven Machine Learning





Biased Non-robust



Specification Driven Machine Learning









Building Adversarially Robust Networks





 $+.007 \times$





"panda" 57.7% confidence $sign(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"nematode" 8.2% confidence





x + $\epsilon \operatorname{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$ "gibbon" 99.3 % confidence

















Probability Vector or Logarithms of the **Probability Vector** (Logits)










$y = f(x; \theta) : x \mapsto \mathbb{R}^C$ $\operatorname{argmax}_{i \in C} y_i = \operatorname{argmax}_{i \in C} f_i(x + \delta; y, \theta)$ $\forall \delta \in B_p(\epsilon) = \{ \delta : \| \delta \|_p \le \epsilon \}$



 $y = f(x; \theta) : x \mapsto \mathbb{R}^C$

$\operatorname{argmax}_{i \in C} y_i = \operatorname{argmax}_{i \in C} f_i(x + \delta; y, \theta)$





 $y = f(x; \theta) : x \mapsto \mathbb{R}^C$ $\operatorname{argmax}_{i \in C} y_i = \operatorname{argmax}_{i \in C} f_i(x + \delta; y, \theta)$ 0.9 Index of the maximum probability in the vector to be the same as where 1 is in the one-hot label.





 $y = f(x; \theta) : x \mapsto \mathbb{R}^C$

$\operatorname{argmax}_{i \in C} y_i = \operatorname{argmax}_{i \in C} f_i(x + \delta; y, \theta)$

$\forall \delta \in B_p(\epsilon) = \{ \delta : \| \delta \|_p \le \epsilon \}$

Set of imperceptible perturbations

 \mathbb{R}^{C} $f_{i}(x + \delta; y, \theta)$ $\|_{p} \leq \epsilon\}$











$\min_{\theta} \mathbb{E}_{(x,y)\sim \mathcal{D}} \left[\mathcal{C}(x; y, \theta) \right]$ Cross entropy





























Perturbation









Adversarial **Evaluation:** Finding the Worst Case



Adversarial Evaluation / Attacks

$\max_{\delta \in B(\epsilon)} \ell(x + \delta; y, \theta)$

- Find the worst case for each example in the test set.
- Then we evaluate the accuracy of this new test set, where each example is now replaced with the worst case adversarial image. This is known as adversarial accuracy.



$\max_{\delta \in B(\epsilon)} \ell(x + \delta; y, \theta)$





$\ell(x+\delta;y,\theta)$ **Constrained Optimisation** Problem





$\max_{\delta \in B(\epsilon)} \ell(x + \delta; y, \theta)$















$\max_{\delta \in B(\epsilon)} \ell(x + \delta; y, \theta)$

Update Step: $\delta \leftarrow \operatorname{Proj}(\delta + \eta \nabla_{\delta} \ell(x + \delta, y; \theta))$ $\operatorname{Proj}(\delta) = \operatorname{argmin}_{\delta' \in B(\epsilon)} \| \delta - \delta' \|$





Fast Gradient Sign Method (Iterated)

$\delta \leftarrow \operatorname{Proj}\left(\delta + \eta \nabla_{\delta} \ell(x + \delta, y; \theta)\right)$ $\delta \leftarrow \operatorname{Proj}\left(\delta + \eta \operatorname{sgn}(\nabla_{\delta} \ell(x + \delta, y; \theta))\right)$



Any Optimiser can be used

$\delta \leftarrow \operatorname{Proj}\left(\delta + \eta \nabla_{\delta} \ell(x + \delta, y; \theta)\right)$ $\delta \leftarrow \operatorname{Proj}\left(\delta + \eta \operatorname{Opt}(\nabla_{\delta} \ell(x + \delta, y; \theta))\right)$

Note that we can replace the gradient by any alterations made by an optimiser, such as momentum optimisation or Adam optimisation.





"Strengths" of Adversarial Evaluation

- Adversarial accuracy is *dependent* on your choice of evaluation.
- Stronger adversarial evaluation should give the lower adversarial accuracy.
- We should *always* try to evaluate our network such that we can obtain the *lowest* adversarial accuracy.
- Strength of your evaluation depends on a few heuristics:
 - Number of steps of projected gradient ascent (PGA).
 - Number of random initialisations of perturbations.
 - The optimiser used.
 - Black box adversarial evaluation e.g. Square Attack



Dangers of Weak Adversarial Evaluation

Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples

Anish Athalye^{*1} Nicholas Carlini^{*2} David Wagner²

Adversarial Risk and the Dangers of Evaluating Against Weak Attacks

Jonathan Uesato¹ Brendan O'Donoghue¹ Aaron van den Oord¹ Pushmeet Kohli¹





Dangers of Weak Adversarial Evaluation

Obfuscated Gradients Give a False Sense of Circumventing Defenses to Adversarial Ex

	Defense	Dataset	Distance	Ac
	Buckman et al. (2018)	CIFAR	$0.031(\ell_\infty)$	0
	Ma et al. (2018)	CIFAR	$0.031(\ell_\infty)$	5
	Guo et al. (2018)	ImageNet	$0.005 (\ell_2)$	0
Advisor	Dhillon et al. (2018)	CIFAR	$0.031(\ell_\infty)$	0
Advers	Xie et al. (2018)	ImageNet	$0.031(\ell_\infty)$	0
	Song et al. (2018)	CIFAR	$0.031(\ell_\infty)$	9
	Samangouei et al.	MNIST	$0.005 \ (\ell_2)$	55
	(2018)			
J	Madry et al. (2018)	CIFAR	$0.031(\ell_\infty)$	47
	Na et al. (2018)	CIFAR	$0.015(\ell_\infty)$	15

Securi xample	ty: s	
5%		
)%*)%)%*	inst Weak Attacks	
9%* 5%**		î
7%	Pushmeet Kohli ¹	
0%	-	6

Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks

Francesco Croce¹ Matthias Hein¹



#	paper	model	clean	report.	AA	AA+	
1	(Carmon et al., 2019) [‡]	available	89.69	62.5	59.65	59.50	
2	(Sehwag et al., 2020) [‡]	available	88.98	-	57.24	57.11	
3	(Wang et al., 2020) [‡]	available	87.50	65.04	56.69	56.26	
4	(Alayrac et al., 2019) ⁺	available	86.46	56.30	56.92	56.01	
21	(Zhang & Xu, 2020)	available	90.25	68.7	38.57		
22	(Kim & Wang, 2020)	available	91.51	57.23	36.10		
23	(Jang et al., 2019)	available	78.91	37.40	35.09		
24	(Wang & Zhang, 2019)	available	92.80	58.6	30.96		
25	(Xiao et al., 2020)*	available	79.28	52.4	17.99		
26	(Jin & Rinard, 2020)	available	90.84	71.22	4.61		
27	(Mustafa et al., 2019)	available	89.16	32.32	0.55		
28	(Chan et al., 2020)	retrained	93.79	15.5	0.18		



#	paper	model	clean	report.	АА	AA+
1	(Carmon et al., 2019) [‡]	available	89.69	62.5	59.65	59.50
2	(Sehwag et al., 2020) [‡]	available	88.98	-	57.24	57.11
3	(Wang et al., 2020) [‡]	available	87.50	65.04	56.69	56.26
4	(Alayrac et al., 2019) ⁺	available	86.46	56.30	56.92	56.01
		•				
21	(Zhang & Xu, 2020)	available	90.25	68.7	38.57	
22	(Kim & Wang, 2020)	available	91.51	57.23	36.10	
23	(Jang et al., 2019)	available	78.91	37.40	35.09	
24	(Wang & Zhang, 2019)	available	92.80	58.6	30.96	
25	(Xiao et al., 2020)*	available	79.28	52.4	17.99	
26	(Jin & Rinard, 2020)	available	90.84	71.22	4.61	
27	(Mustafa et al., 2019)	available	89.16	32.32	0.55	
28	(Chan et al., 2020)	retrained	93.79	15.5	0.18	



#	paper	model	clean	report.	АА	AA+
1	(Carmon et al., 2019)‡	available	89.69	62.5	59.65	59.50
2	(Sehwag et al., 2020) [‡]	available	88.98	-	57.24	57.11
3	(Wang et al., 2020) ⁺	available	87.50	65.04	56.69	56.26
4	(Alayrac et al., 2019) [‡]	available	86.46	56.30	56.92	56.01
21	(Zhang & Xu, 2020)	available	90.25	68.7	38.57	
22	(Kim & Wang, 2020)	available	91.51	57.23	36.10	
23	(Jang et al., 2019)	available	78.91	37.40	35.09	
24	(Wang & Zhang, 2019)	available	92.80	58.6	30.96	
25	(Xiao et al., 2020)*	available	79.28	52.4	17.99	
26	(Jin & Rinard, 2020)	available	90.84	71.22	4.61	
27	(Mustafa et al., 2019)	available	89.16	32.32	0.55	
28	(Chan et al., 2020)	retrained	93.79	15.5	0.18	





Gradient Obfuscation











 $\min_{\theta} \mathbb{E}_{(x,y)\sim\mathcal{D}} \left| \max_{\delta\in B(\epsilon)} \ell(x+\delta;y,\theta) \right|_{t}$





Adversarial Training made Cheaper $\max_{\delta \in B(\epsilon)} \ell(x + \delta; y, \theta)$



Adversarial Training made Cheaper

$\max_{\delta \in B(\epsilon)} \ell(x + \delta; y, \theta)$

Cheaper Adversarial Training : Few steps of Gradient Ascent for Inner Maximization.





Adversarial Training made Cheaper

$\max_{\delta \in B(\epsilon)} \ell(x + \delta; y, \theta)$



Cheaper Adversarial Training : Few steps of Gradient Ascent for Inner Maximization.






RESULT:

The network learns to cheat by making a highly non-linear surface such that the **optimum is hard to find** using a suboptimal optimization procedure.





RESULT:

The network learns to cheat by making a highly non-linear surface such that the **optimum is hard to find** using a suboptimal optimization procedure.

Loss plotted for two input dimensions







RESULT:

The **network learns to cheat** by making a **highly non-linear surface** such that the **optimum is hard to find** using a suboptimal optimization procedure.

Loss plotted for two input dimensions







3

Verification Algorithms



Verification Algorithms

- Complete
 - Exhaustive Proof
- Incomplete
 - If proof can be found then it is a guarantee but proof cannot always be found.



Verification of Specifications y = f(x)





- *Y*





y



Verification of Specifications







Verification of Specifications









 $x^0 < x \leq \overline{x}^0$







 $\underline{z}^{l} = \left[W^{l}\right]_{+} \underline{x}^{l} + \left[W^{l}\right]_{-} \overline{x}^{l} + b^{l}$ $x^0 < x < \overline{x}^0$ $\overline{z}^{l} = \left[W^{l}\right]_{+} \overline{x}^{l} + \left[W^{l}\right]_{-} \underline{z}^{l} + b^{l}$ $\underline{x}^{l+1} = h^l\left(\underline{z}^l\right)$ $\overline{x}^{l+1} = h^l\left(\overline{z}^l\right)$







Y satisfies the specification





Y satisfies the specification









Y satisfies the specification











Other Specifications

Semantic Consistency:





Physics Consistency :





 $oxed{E} \left[E(x_{t+1}) \leq E(x_t)
ight] \ oxed{E} \left[E(x_{t+1}) > E(x_t)
ight]$







Ethics & Machine Learning





NO F ZTO





CALL COM



Consent







CALL COMP



Consent Representation







WINGSOM



Consent Representation Prejudicial Associations







REPORTING TO YOU

TECH

India Is Creating A National Facial Recognition System, And Critics Are Afraid Of What Will Happen Next

"Unless we all get plastic surgery at the same time, there's nothing we can do about it."



Pranav Dixit BuzzFeed News Reporter

Last updated on October 9, 2019, at 10:39 p.m. ET Posted on October 9, 2019, at 7:01 p.m. ET













REPORTING TO YOU

TECH

India Is Creating A National Facial Recognition System, And Critics Are Afraid Of What Will Happen Next

"Unless we all get plastic surgery at the same time, there's nothing we can do about it."



Pranav Dixit BuzzFeed News Reporter

Last updated on October 9, 2019, at 10:39 p.m. ET Posted on October 9, 2019, at 7:01 p.m. ET







Criminal Justice





REPORTING TO YOU

TECH

India Is Creating A National Facial Recognition System, And Critics Are Afraid Of What Will Happen Next

"Unless we all get plastic surgery at the same time, there's nothing we can do about it."



Pranav Dixit BuzzFeed News Reporter

Last updated on October 9, 2019, at 10:39 p.m. ET Posted on October 9, 2019, at 7:01 p.m. ET







Criminal Justice



Access to Jobs





REPORTING TO YOU

TECH

India Is Creating A National Facial Recognition System, And Critics Are Afraid Of What Will Happen Next

"Unless we all get plastic surgery at the same time, there's nothing we can do about it."



Pranav Dixit BuzzFeed News Reporter

Last updated on October 9, 2019, at 10:39 p.m. ET Posted on October 9, 2019, at 7:01 p.m. ET











REPORTING TO YOU

TECH

India Is Creating A National Facial Recognition System, And Critics Are Afraid Of What Will Happen Next

"Unless we all get plastic surgery at the same time, there's nothing we can do about it."



Pranav Dixit BuzzFeed News Reporter

Last updated on October 9, 2019, at 10:39 p.m. ET Posted on October 9, 2019, at 7:01 p.m. ET









Power and Responsibility



Those who design and develop these technologies are in a position of power



Power and Responsibility

Those who design and develop these technologies are in a position of power



The choices about design and use which they make have a significant impact on the lives of others



Power and Responsibility

- Those who design and develop these technologies are in a position of power
- \rightarrow
- The choices about design and use which they make have a significant impact on the lives of others



With this power comes responsibility





TUNIO



Scientific research is not value-neutral





Scientific research is not value-neutral

It is a human activity governed by shared norms



Scie

 \rightarrow

Scientific research is not value-neutral

It is a human activity governed by shared norms

These norms have profound effects



Responsible Innovation

serie

J.

Responsible Innovation

A transparent and iterative process by which societal actors and technologists become mutually responsive to each others needs to ensure the ethical and social value of scientific endeavour


Principles and Process



The Responsibility of Technologists



The Responsibility of Technologists

Intrinsic to the design of technology: technical safety, robustness, and ensuring that systems are accountable, transparent and fair from the start.



The Responsibility of Technologists

Intrinsic to the design of technology: technical safety, robustness, and ensuring that systems are accountable, transparent and fair from the start.

Extrinsic to the design of technology: intended use, mode of deployment in real-world settings, frameworks and institutions governing their application and use.



The AI Ethics Landscape

'Ethical Guidelines for trustworthy Al' – the European Union
'OECD Principles on Al' – OECD
'Beijing Al Principles' – Beijing Academy of Artificial Intelligence
'Asilomar Al Principles' – Future of Life Institute



Key Values



Fairness, privacy and transparency



Key Values



Fairness, privacy and transparency



Individual rights



Key Values



Fairness, privacy and transparency



Individual rights



Shared benefit from science





Bridging the Gap

From principles to practice

A Five Step Process





















Is it possible to mitigate these risks, or eliminate them entirely?











Is it possible to mitigate these risks, or eliminate them entirely?

With these measures in place, does the proposed action violate a 'red line' or moral constraint?











2

Is it possible to mitigate these risks, or eliminate them entirely?

Might this technology directly or

indirectly result in harm?

With these measures in place, does the proposed action violate a 'red line' or moral constraint?

With these measures in place, do the benefits outweigh the risks?



1. Does the technology have socially beneficial uses?



1. Does the technology have socially beneficial uses?

Well-being

 \rightarrow

 \rightarrow

 \rightarrow

 \rightarrow

Autonomy

Justice

Public institutions

Global challenges



2. Might this technology directly or indirectly result in harm?



Undermine health, well-being or human dignity

Restrict freedom or autonomy

Lead to unfair treatment or outcomes

Harm public institutions or civic culture

Infringe human rights

in harm?

2. Might this technology directly or indirectly result









Control the release of technologies or the flow of information

Adopt technical solutions and countermeasures



Control the release of technologies or the flow of information

Adopt technical solutions and countermeasures

Help the public understand new technologies

 \rightarrow



Control the release of technologies or the flow of information

Adopt technical solutions and countermeasures

Help the public understand new technologies

 \rightarrow

Seek out policy solutions and legal frameworks



4. With these measures in place, does the proposed action violate a 'red line' or moral constraint?



Consent

Weapons

Surveillance

International law and human rights

4. With these measures in place, does the proposed action violate a 'red line' or moral constraint?



5. With these measures in place, do the benefits outweigh the risks?





Two final tests



Two final tests

Have you thought about everyone who could be affected by your action?



Two final tests

Have you thought about everyone who could be affected by your action?

Might you have reason to regret it later?















• Those who design and develop these technologies have a responsibility to think about how they will be used.

Key Ideas

• Those who design and develop these technologies have a responsibility to think about how they will be used.

• There are concrete steps and processes that we can put in place to make sure this responsibility is successfully discharged.

Key Ideas

- Those who design and develop these technologies have a responsibility to think about how they will be used.
- There are concrete steps and processes that we can put in place to make sure this responsibility is successfully discharged.
- We are responsible for what we can reasonably foresee and should take steps to bring about positive outcomes



New Directions




Research







Research

Norms







Research

Norms

Practice



Thank you





Technology can "lock in" value

Designing for social effect Robert Moses, New York City (1930s)



Technology can "lock in" value

Designing for social effect Robert Moses, New York City (1930s)



Ethics and Machine Learning



Ethics and Technoloy

