How should society prepare for advances in Artificial Intelligence?

Alistair Knott Dept of Computer Science, University of Otago Soul Machines, Auckland







 AI machines will achieve human-level intelligence in 2029, 'and will openly petition for recognition of this fact'.



- Al machines will achieve human-level intelligence in 2029, 'and will openly petition for recognition of this fact'.
- Uploading human minds to machines will become possible in the 2030s, ushering in an era of human immortality.



AI prophets of doom

• The development of full artificial intelligence 'could be the greatest disaster in human history'...



AI prophets of doom

- The development of full artificial intelligence 'could be the greatest disaster in human history'...
- It could 'spell the end of the human race'.



AI prophets of doom

• 'Before the prospect of an intelligence explosion, we humans are like small children playing with a bomb. Such is the mismatch between the power of our plaything and the immaturity of our conduct.'



Evangelists

• 'I am optimistic. People who are naysayers and try to drum up these doomsday scenarios—I just don't understand it. It's really negative and in some ways I actually think it is pretty irresponsible.'



Prophets of doom

 'I've talked to Mark about this. His understanding of the subject is limited.'



Outline of the talk

- 1 What can AI do right now?
- 2 What are the main impacts AI will have in the immediate future?
- 3 What aspects of near-future AI might need regulation?
- 4 How should our society prepare itself for the arrival of AI?

1. What can AI do right now?

Agile robots



Agile robots



Self-driving vehicles



Self-driving vehicles



Open-domain question answering



Royal Society / Te Apārangi

Domain-general action learning



Royal Society / Te Apārangi

Visual object recognition













flamingo

cock

ruffed grouse

quail

partridge



Egyptian cat







tabby

















miniature schnauzer standard schnauzer giant schnauzer

Scene description

I think it's a man and a dog on a rocky beach.



Machine translation

Google	III O 🔥
Translate	Turn off instant translation
English Spanish French English - detected 👻	English Spanish Arabic - Translate
The quick brown fox jumped × over the lazy dog	قفز الثعلب البني السريع فوق الكلب الكسول ©
4) ℓ 44/5000	☆□••

qafaz althuelib albaniu alsarie fawq alkalb alkusul

2. What impacts will AI have in the immediate future?

2. What impacts will AI have in the immediate future?

I'll focus on employment.

An influential study: Frey and Osborne (2013)

Frey & Osborne analysed a US database listing the 'key skills' needed for different job types. They estimated how 'automatable' each skill is.

An influential study: Frey and Osborne (2013)

The main non-automatable skills are those involving perception/manipulation, creative intelligence, and social intelligence.

An influential study: Frey and Osborne (2013)

Then they estimated the probability of automation for each job type. Here are some examples:



Accountant and auditors **Retail salespersons** Technical writers Real estate sales agents Word processors and typists **Commerical pilots** Health technologists **Chemical engineers**

A breakdown of F&O's estimates by job category





Many economists think that Frey and Osborne overestimate how automatable jobs are.

Many economists think that Frey and Osborne overestimate how automatable jobs are.

A common idea (see e.g. Arntz et al., 2016):

- Most jobs have *some tasks* that are easy to automate, and *some tasks* that are hard.
- In this case, what we might see is a *restructuring* of jobs, that allocates the automatable tasks to computers.

But don't new technologies always create new jobs?

But don't new technologies always create new jobs?



But don't new technologies always create new jobs?



'Don't worry, there will always be jobs for horses.'

Royal Society / Te Apārangi

3. What aspects of AI might need regulation?

3. What aspects of AI might need regulation?

I'll touch on five topics:

- Employment
- Machine bias
- Transparency
- Accountability
- Ethics

Regulating employment
Some key questions:

Some key questions:

Are there certain jobs we don't want computers to do?

• E.g. teacher, carer, soldier...

Some key questions:

Are there certain jobs we don't want computers to do?

• E.g. teacher, carer, soldier...

How do we revise social security, for human workers displaced by AI?

Some key questions:

Are there certain jobs we don't want computers to do?

• E.g. teacher, carer, soldier...

How do we revise social security, for human workers displaced by AI?

We could impose some kind of tax on AI systems...

Some key questions:

Are there certain jobs we don't want computers to do?

• E.g. teacher, carer, soldier...

How do we revise social security, for human workers displaced by AI?

- We could impose some kind of tax on AI systems...
- We could institute some form of universal basic income...

Some key questions:

Are there certain jobs we don't want computers to do?

• E.g. teacher, carer, soldier...

How do we revise social security, for human workers displaced by AI?

- We could impose some kind of tax on AI systems...
- We could institute some form of universal basic income...

Do we want a world where people don't have to work?

Some examples of bias in AI systems:

Some examples of bias in AI systems:

 Al systems can learn word meanings from huge text corpora. These systems show all the biases exhibited by humans in word-association tests. (Caliskan *et al.*, 2017)

Some examples of bias in AI systems:

 AI systems can learn word meanings from huge text corpora. These systems show all the biases exhibited by humans in word-association tests. (Caliskan *et al.*, 2017)

Phil - kitchen vs Sally - kitchen

Some examples of bias in AI systems:

• Al systems are currently used in US courts to help judges make bail/parole decisions. These systems have been shown to have biases against black people. (Angwin *et al.*, 2016)

Some examples of bias in AI systems:

 AI systems are currently used in US courts to help judges make bail/parole decisions. These systems have been shown to have biases against black people. (Angwin *et al.*, 2016)

Some open questions:

Some examples of bias in AI systems:

 AI systems are currently used in US courts to help judges make bail/parole decisions. These systems have been shown to have biases against black people. (Angwin *et al.*, 2016)

Some open questions:

• How can we test AI systems for bias?

Some examples of bias in AI systems:

 AI systems are currently used in US courts to help judges make bail/parole decisions. These systems have been shown to have biases against black people. (Angwin *et al.*, 2016)

Some open questions:

- How can we test AI systems for bias?
- Can we legislate against bias in AI systems?

Some examples of bias in AI systems:

 AI systems are currently used in US courts to help judges make bail/parole decisions. These systems have been shown to have biases against black people. (Angwin *et al.*, 2016)

Some open questions:

- How can we test AI systems for bias?
- Can we legislate against bias in AI systems?
- Can we use AI systems to reduce, or eliminate bias?

Modern AI systems use machine learning techniques to make decisions:

- This involves consulting huge databases
- and performing *complex computations*.

Modern AI systems use machine learning techniques to make decisions:

- This involves consulting *huge databases*
- and performing *complex computations*.

Humans typically *can't understand* how these systems reach their conclusions.

Modern AI systems use machine learning techniques to make decisions:

- This involves consulting *huge databases*
- and performing *complex computations*.

Humans typically *can't understand* how these systems reach their conclusions.

Some questions:

Modern AI systems use machine learning techniques to make decisions:

- This involves consulting *huge databases*
- and performing *complex computations*.

Humans typically *can't understand* how these systems reach their conclusions.

Some questions:

• Should (some) AI systems be required to explain their decisions?

Modern AI systems use machine learning techniques to make decisions:

- This involves consulting huge databases
- and performing *complex computations*.

Humans typically *can't understand* how these systems reach their conclusions.

Some questions:

- Should (some) AI systems be required to explain their decisions?
- How can we build the 'explanation mechanism'?

Say a person is using a driverless car, and has an accident. Who's to blame?

Say a person is using a driverless car, and has an accident. Who's to blame?

• The person in the car?

Say a person is using a driverless car, and has an accident. Who's to blame?

- The person in the car?
- The company who designed the AI system the car is using?

Say a person is using a driverless car, and has an accident. Who's to blame?

- The person in the car?
- The company who designed the AI system the car is using?
- The AI system itself??

Say a person is using a driverless car, and has an accident. Who's to blame?

- The person in the car?
- The company who designed the AI system the car is using?
- The AI system itself??

Some questions to consider:

Say a person is using a driverless car, and has an accident. Who's to blame?

- The person in the car?
- The company who designed the AI system the car is using?
- The AI system itself??

Some questions to consider:

• Al systems are designed to behave very flexibly, in a wide range of circumstances. How can a company *guarantee the performance* of such systems?

Say a person is using a driverless car, and has an accident. Who's to blame?

- The person in the car?
- The company who designed the AI system the car is using?
- The AI system itself??

Some questions to consider:

- Al systems are designed to behave very flexibly, in a wide range of circumstances. How can a company *guarantee the performance* of such systems?
- Say that driverless cars are more reliable than human drivers. How can companies be protected from crippling lawsuits?

Say we have an advanced AI system that can perform a wide range of tasks about the house. (Cooking, cleaning, tidying etc.)

Say we have an advanced AI system that can perform a wide range of tasks about the house. (Cooking, cleaning, tidying etc.)

• How can we ensure it doesn't do anything silly, or wrong?

Say we have an advanced AI system that can perform a wide range of tasks about the house. (Cooking, cleaning, tidying etc.)

- How can we ensure it doesn't do anything silly, or wrong?
- It may need general principles to regulate its behaviour.

Say we have an advanced AI system that can perform a wide range of tasks about the house. (Cooking, cleaning, tidying etc.)

- How can we ensure it doesn't do anything silly, or wrong?
- It may need general principles to regulate its behaviour.
- What principles do we give it?

Say we have an advanced AI system that can perform a wide range of tasks about the house. (Cooking, cleaning, tidying etc.)

- How can we ensure it doesn't do anything silly, or wrong?
- It may need general principles to regulate its behaviour.
- What principles do we give it?

If it's a system that *learns* how to behave, there are some interesting possibilities:

Say we have an advanced AI system that can perform a wide range of tasks about the house. (Cooking, cleaning, tidying etc.)

- How can we ensure it doesn't do anything silly, or wrong?
- It may need general principles to regulate its behaviour.
- What principles do we give it?

If it's a system that *learns* how to behave, there are some interesting possibilities:

• One possibility is that it should be taught in the same way that we teach children.
Regulating the ethics of Al systems

Say we have an advanced AI system that can perform a wide range of tasks about the house. (Cooking, cleaning, tidying etc.)

- How can we ensure it doesn't do anything silly, or wrong?
- It may need general principles to regulate its behaviour.
- What principles do we give it?

If it's a system that *learns* how to behave, there are some interesting possibilities:

- One possibility is that it should be taught in the same way that we teach children.
- However learning happens, it will probably involve AI systems learning about human values: both through explicit instruction, and through inference from human behaviour.

4. How should society prepare for advances in AI?

4. How should society prepare for advances in AI?

Let's think about New Zealand specifically...

(i) Interdisciplinary discussion groups

(i) Interdisciplinary discussion groups

The discussion about AI needs to cut across many disciplines.

- Technical AI people / computer scientists
- Business people
- Economists
- Lawyers
- Social scientists
- Ethicists
- . . .

Some NZ groups I'm involved with

Otago's AI and Society Discussion Group



Some NZ groups I'm involved with

An Otago project funded by the NZ Law foundation: Al and Law In NZ



Royal Society / Te Apārangi

Some NZ groups I'm involved with

The AI Forum of New Zealand



Computer Science students need an understanding of the social consequences of AI. (And of IT more generally.)

Computer Science students need an understanding of the social consequences of AI. (And of IT more generally.)

• Ethics should be an important (examinable) component of several CS courses.

Computer Science students need an understanding of the social consequences of AI. (And of IT more generally.)

- Ethics should be an important (examinable) component of several CS courses.
- CS degrees should require social science papers.

Computer Science students need an understanding of the social consequences of AI. (And of IT more generally.)

- Ethics should be an important (examinable) component of several CS courses.
- CS degrees should require social science papers.

Students planning on entering politics/law should be encouraged to become literate in IT (and/or other areas of science).

Computer Science students need an understanding of the social consequences of AI. (And of IT more generally.)

- Ethics should be an important (examinable) component of several CS courses.
- CS degrees should require social science papers.

Students planning on entering politics/law should be encouraged to become literate in IT (and/or other areas of science).

Some recent Otago initiatives:

- LAWS 102 Introduction to Law and new Technologies
- LAWS 428 Law and Emerging Technologies
- Bachelor of Arts and Sciences

- J Angwin, J Larson, S Mattu, and L Kirchner. What algorithmic injustice looks like in real life. ProPublica, 2016.
- M Arntz, T Gregory, and U Zierahn. Automation for jobs in OECD countries: A comparative analysis. OECD Social, Employment and Migration Working Papers, No. 189, OECD Publishing, Paris, 2016.
- A Caliskan, J Bryson, and A Narayanan. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017.
- C Frey and M Osborne. The future of employment: how susceptible are jobs to computerisation? *Technological Forecasting and Social Change*, 114:254–280, 2013.