

As you enter your driverless vehicle, your Chatbot reminds you to pick up milk. Your milk can be traced through its supply chain back to the farm where a robot milked the cow. The milk contains genetically engineered enzymes that were designed by predictive computer software to improve your health. Welcome to now, or at least, very soon from now ...

Artificial Intelligence

By Jennifer Yarnold, Paul Henman, Christopher McEwan, Amelia Radke & Karen Hussey

Definition

“Artificial intelligence is the theory and development of computer systems able to perform tasks normally requiring human intelligence.”^[1]

AI-driven applications have become prosaic in our lives, whether we are aware of it or not. From credit approvals and insurance premiums to your Fitbit providing you with real-time updates on your health; and the advertising and news that appears on your feed – these are just a few things driven by the power of AI. Increasingly, computers are doing more work for us, knowing more about us, and ultimately making decisions that affect our lives.

And with that, too, comes risk.

There is an intense debate on the ethical implications and impacts of AI on society and economies, and how AI should be regulated. In part, this is due to a lack of clarity as to what constitutes AI, as our perception of ‘intelligence’ evolves alongside evermore sophisticated technologies. While AI may conjure up images of androids overthrowing humanity in the manner of *I, Robot* and *Terminator*, it is probably not as ominous, nor as human-like, as it sounds.

Today’s AI is a technology that can adapt itself to changing circumstances based on a particular self-learning ability to produce a specific output, independent of human control. This AI relies on data-driven algorithms that look for underlying and sometimes subtle patterns to reveal new knowledge. This process, combining data mining and machine learning (ML), contrast with earlier computer algorithms which were based on human designed rigid instructions or hypothetical models.

For example, early email spam filters would deem a message as spam by matching hard-coded keywords such as ‘huge sale’ to email content manually, or if sender emails were not in your contact list. Unsurprisingly, often valid emails might end up in your spam box, and unsolicited marketers would change email content to avoid spam filter detection. In contrast, AI-based spam filters are trained to find more subtle differences between emails flagged as ‘spam’ and ‘non-spam’, and can improve their detection capabilities with new data feeds to keep up with spam generators. Thus a more useful definition characterises today’s AI as:

“A system’s ability to correctly interpret external data, to learn from such data, and to use those learnings to achieve specific goals and tasks through flexible adaptation.”^[2]

Figure 3. Artificial intelligence technology and innovation pipeline



Context

In 1951, the father of computer science and pioneering English codebreaker against the Nazis, Alan Turing, foresaw the potential for machines to out-think humans. Turing went so far as to devise a test of ML long before there were machines worth testing. The 'Turing Test' requires a person to ask a series of questions to both a computer and a human; if upon reading the responses of the human and machine, the questioner is unable to tell which is human and which is the computer, the computer would be described as exhibiting intelligent behaviour.^[3] It would be more than half a century later for AI to truly come of age, enabled by mass digitisation to collect vast amounts of data, and an exponential increase in computer power required to process it.

Drivers

Innovation – The drivers of AI are simple: humans want to innovate. We want to do less, have more, and do better. And we want to know things. AI technologies have been driven primarily by tech people – computer scientists, engineers and the like – who want to push the boundaries and see how much machines can replicate human intelligence. Most scientists, corporations and governments want to do better at what they do: make profits; identify and treat illnesses; ensure security from terror or crime; reduce human error; use precious resources more efficiently, and so on and so forth.

Precision and personalisation – We all want to be treated as a unique person. Personalised health services are being used to understand the specific

“AI’s capabilities and speed to carry out complex tasks are far superior to humans.”

— World Economic Forum

care needs of patients better. AI systems are being developed to better respond to each person, including customised advertisements, criminal sentencing and parole decisions, identifying children at risk of neglect or abuse, and deciding if you are the 'best person for the job'.

Efficiency and productivity – For industry, AI’s appeal is in automation and robotic systems to reduce labour costs, improve efficiency and productivity, better utilise resources and standardise processes. Moreover, AI increases their ability to understand who their customers are and what they want, which in turn gives a competitive advantage. For end-users, AI’s appeal lies in personalisation and convenience.

showing computing power and device complexity have doubled approximately every two years. Since AI relies on fast processing of large amounts of data, it has been enabled by exponential increases in computer processing speed and storage capacity, internet speed and ‘cloud computing’.



The Internet of Things – Technological innovations and economies of scale have seen dramatic price reductions in digital devices and, henceforth, higher usage. The so-called Internet of Things (IoT) are devices and sensors embedded into everyday household and mobile items connected directly or indirectly online. Smartphones, tablets, smart watches, televisions, remote sensors and equipment monitors, to name a few, capture real-time data that is fed into ‘the cloud’. For example, Google Maps predict traffic conditions and travel times partly by using real-time GPS movements retrieved from its app users’ smartphone to determine how fast they are moving.

Enablers

Computing power and the ineffable ‘cloud’ – Advances in computers and electronics have closely followed Moore’s law – a historical observation

Accumulated Universe of Digital Data^b

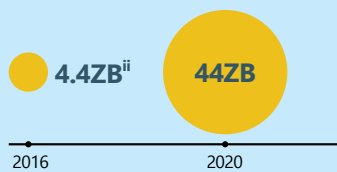
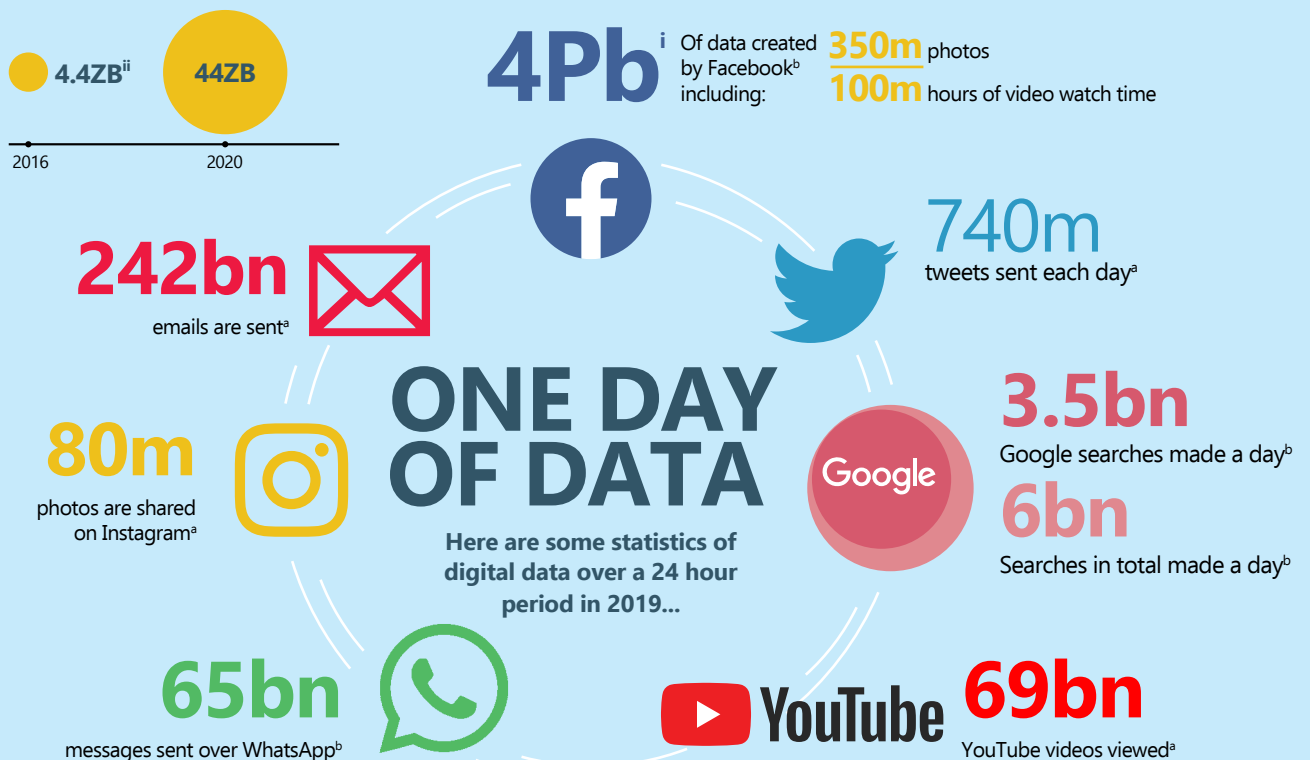
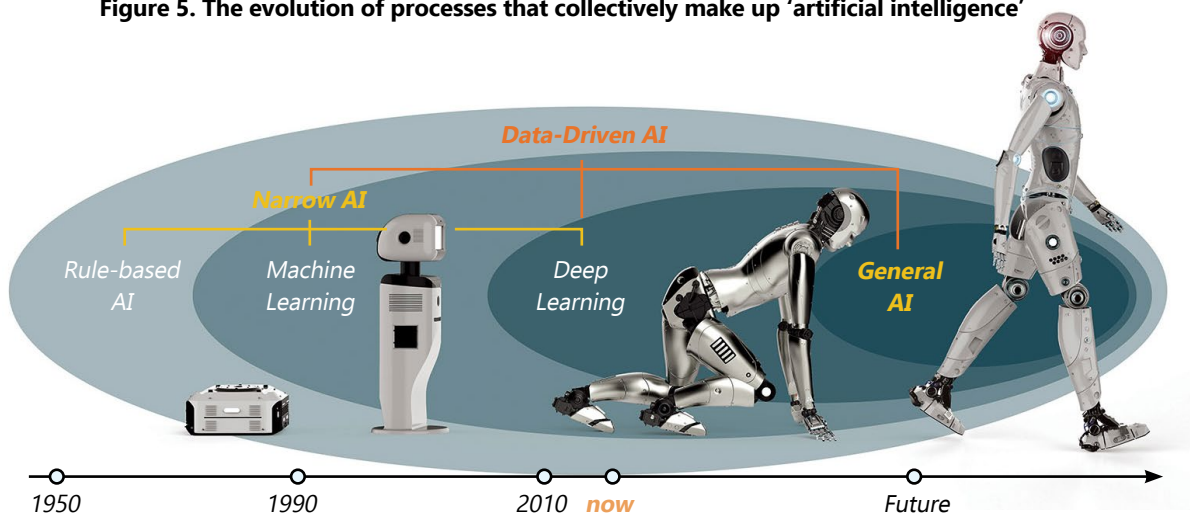


Figure 4.



ⁱ 1 Petabyte (PB) = 1 x 10¹² Kilobytes (Kb) or 1000 terabytes. ⁱⁱ 1 Zetabyte (Zb) = 1 x 10¹⁸ Kb or 1 million Pb. For comparison, the computers that put Apollo 11 on the moon contained just 4 Kb of processing power and 32 Kb of storage capacity. Data sources: ^a www.internetlivestats.com (accessed June 2019); and ^b 'The Future of Data' Raconteur, 2019.

Figure 5. The evolution of processes that collectively make up 'artificial intelligence'

Next Generation Sequencing (NGS) – digital laboratory technologies have also revolutionised the generation of big datasets; perhaps none more so than NGS, which became a game changer in genomics in 2007. The Human Genome Project was the first of its kind to map the entire sequence of a human genome. Starting in 1990, it took 13 years to complete and involved some 20 global organisations at the cost of about US\$3 billion.^[4] With today's NGS technologies, an entire genome can be sequenced within a day or two for as little as US\$1,000-2,000.^[5] Affordable genome sequencing has created a global trend in online DNA ancestry companies that can trace the heritage of its customers, as well as unbeknown others. In the US, publicly available genomic data was used by FBI officials to track the 'Golden State killer'; creating debate around data privacy and consent.^[6]

A key issue for regulators is ensuring that individuals are both aware of, and consent to, the collection of personal information by IoT devices. In general, users are aware of, and therefore, more likely to consent to the collection of information, which they have explicitly provided, such as when completing online forms. However, the industry in general has poorly educated the public about other forms of metadata being collected that can identify an individual. The European Union's General Data Protection Regulation (the GDPR), which took effect in May 2018 includes some regulations to address these issues (for instance, the 'right to be forgotten') that are currently not covered within the Australian Privacy Act.^[7]

Inputs

Big Data – the enormity of data is what feeds modern AI. The enablers are things that produce

loads of data: metadata, genomic, consumer behaviour, advertising and market response, economic and financial, spatial, environmental, weather, and the list goes on.

Once the data has been collected and digitally stored, key issues remain around **data ownership and protection of privacy**. By 2020, the amount of data in the digital universe will be 10-fold higher than just 10 years ago (**Figure 4**). Big data has become a highly valuable commodity. But who should own this data and who should have the right to access it? In light of recent cases such as the Cambridge Analytica scandal which used personal data from millions of Facebook users for targeted political advertising, several countries have adopted frameworks and legislation to keep up with the pacing problem in data collection, ownership and use.

Processes

From Narrow to General to Super AI – how we perceive what human intelligence is and therefore what counts as 'artificial' intelligence is evolving (**Figure 5**). Part of the ambiguity surrounding AI, is that its definition has changed in response to increasing computing abilities, challenging our notion of what counts as 'intelligence'. Computers are adept at calculations with enormous numbers in infinitesimal times (including the humble calculator), which we might have once considered as intelligent. Indeed, there are specific tasks such as calculation at which many systems exceed the capabilities of most humans yet accepting these computers as therefore being more intelligent than a human still seems disconcerting.

“Artificial intelligence can play chess, drive a car and provide medical diagnoses. Examples include Google DeepMind’s AlphaGo, Tesla’s self-driving vehicles, and IBM’s Watson. This type of artificial intelligence is referred to as narrow (or weak) artificial intelligence – non-human systems that can perform a specific task. We encounter this type on a daily basis, and its use is growing rapidly”.^[8]

By contrast to narrow (or weak) AI, **general (or strong) AI** refers to a machine that can perform any task as well as, or better than a human and being able to adapt and respond to a wide variety of circumstances. This is the ultimate goal of many researchers in the field but is also the same type of intelligence which is most likely to realise human fears of rogue AI.

Machine learning uses data to ‘train’ itself – ML is a process for producing AI by enabling a computer program that can improve itself. Today, AI is often used interchangeably with ML, though it is in fact, a subset within AI. These systems can perform tasks without explicit instructions and update their algorithms in response to their own received inputs; much like a human, they can learn and improve with experience. ML’s underlying algorithms use data to cluster patterns and make predictions based on inference. ML systems can be classified as supervised (e.g. ‘trained’) or unsupervised – often used for data mining (**Figure 6**). For instance,

imagine your boss asks you to segment your company’s customer files into three drawers of your filing cabinet with each drawer containing ‘similar’ customers. Under supervision, your boss will go through some records with you and provide guidance on what features to categorise them by. Unsupervised, your boss will tell you to look at all the files and decide for yourself which features you think best distinguish different groups (here, you are data mining). The more files you look at, the better you will get at finding similarities and differences between customers.

Facial recognition is an example of **supervised ML**—it is trained with images of individuals’ faces to define biometric features (e.g. the distance between your eyes and from forehead to chin), as well as ‘facial landmarks’ that create a unique ‘facial signature’ in the form of a mathematical formula. Facial recognition systems used by US law enforcement can identify a US citizen among 117 million people on its database.

Data mining is a form of unsupervised ML often used in business analytics to discover new insights about its markets and customers, and to make predictive analyses.^[9]

Deep Learning is a more sophisticated form of ML – its power lies in its multi-layered structure, where each layer progressively extracts new information from lower levels. The multi-layered approach allows corresponding machines to not only follow pre-programmed decisions but to respond to changes within their environment. An example is autonomous cars that can make real time decisions about speed and direction by analysing sensor-based data without input from a human user.^[10]

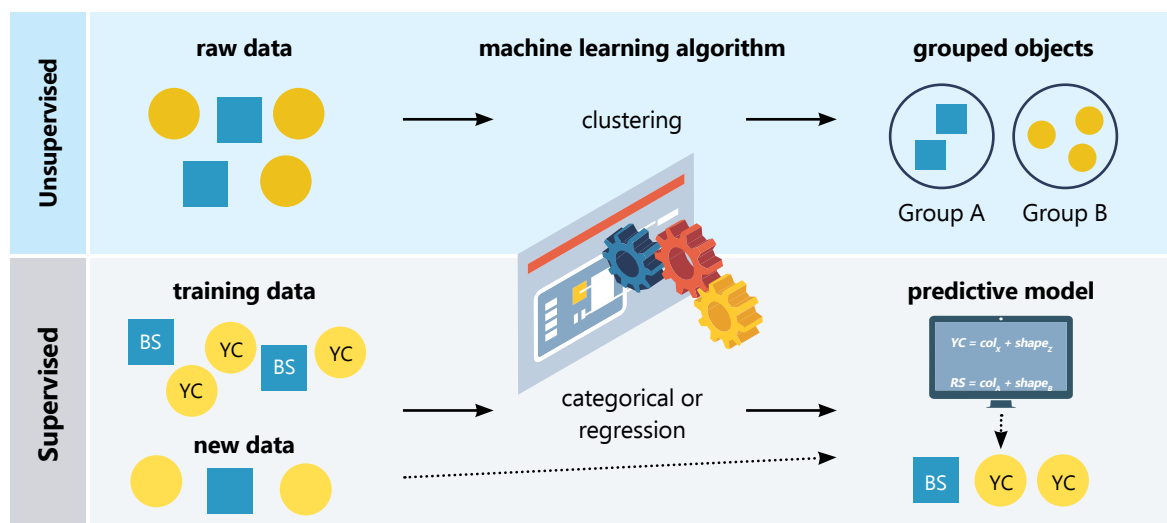


Figure 6. The processes of unsupervised and supervised machine learning



Will we all have our own version of *Iron Man's J.A.R.V.I.S?*

Artificial general intelligence (AGI) is the ultimate goal of some computer scientists. In *Marvel's* superhero world, Tony Stark's robot assistant J.A.R.V.I.S (Just A Rather Very Intelligent System) is perhaps the ultimate AGI system; able to understand and communicate in much the same way as a human can. J.A.R.V.I.S. has a range of functions – from running the Stark Industries business to managing the Stark Mansion and its security and navigating all of the Iron Man Armor – it is beyond general AI toward **artificial super intelligence** – exceeding the capacity of a human.

While AGI does not yet exist, AI systems are becoming smarter, faster, more fluid and human-

like. British AI company DeepMind Technologies' AlphaZero program is making headway toward machines solving multiple-faceted problems with reinforced learning techniques – to such an extent that it taught itself to play Go, Chess, and Shogi without any pre-programmed knowledge of the rules and beat the world's best players. More recently, AI researchers developed an AI that can compete with humans in 3D multiplayer computer games.^[11]

The inevitable rise of highly sophisticated quantum computing is expected to further revolutionise AI toward AGI in ways beyond imagination.^[12]

The ethics of AI decision making underlying the algorithms that act on data add layers of complexity – Key issues centre around transparency, explainability, responsibility and fairness. While some processes are aimed to reduce human bias, likewise it must be ensured that these pre-existing biases are not simply transferred from humans to machine. The problems are intensified with 'black box' systems where there is a lack of transparency and knowledge of the algorithms underlying the decisions required. Moreover, caution must be taken to use proxies, where firm data is absent, to infer outcomes. Several recent cases in the US have highlighted AI algorithms causing discrimination, such as determining if children are at risk by using data such as welfare support granted to parents in their algorithms or using location-dependent crime statistics in judicial support tools to infer the likelihood of criminals re-offending, often discriminating against individuals from African American neighbourhoods.

A notable example is the COMPAS sentencing software used in the US to inform judges of the likelihood incarcerated persons will re-offend. An independent study found the software was biased against African Americans who were unfairly categorised as a higher risk. When challenged, the software developer refused to release the

AI-algorithm underlying the software, claiming protection of their intellectual property.

Outputs

The scope and potential applications of AI are enormous and are being incorporated into virtually every sector (**Table 1**) as more and more businesses and government agencies are finding uses for it. The main applications centre on predictive analytics, robotic automation, transportation, and decision making. Many agencies are adopting biometric systems to streamline identification in law enforcement, immigration, correctional services and investigations. Text analytics and natural language processing are being used for security systems and fraud detection. Digital twins create virtual replicas of physical objects (non-living and living) to monitor equipment and infrastructure – such as aircraft engines, gas turbines, and bridge structures – and predict failures with cloud-hosted software models of machines.

Unsurprisingly, the tech giants are dominating the market for the most on-trend consumer products too: Chatbots such as Amazon's Alexa home assistant, that has a remarkable ability to

Table 1. Applications of artificial intelligence by sector

Sector	Products and services
Health and aged care	<ul style="list-style-type: none"> • Remote patient monitoring enables longer independent living for the elderly and enhances health in remote areas • Genomics and biomarker discovery for detection of inheritable diseases • Image-based diagnostics lessen physician analyses of medical scans • Drug design data mining of proteomics can lead to better-targeted medicines • Autonomous robotic surgery • AI-driven technologies that enable participation by people with disabilities
Education and research	<ul style="list-style-type: none"> • Dubbed the “invention of a method of invention” – AI can be used to accelerate the pace of research and development from drug discovery and design to protein folding and advanced materials • Virtual teaching assistants • Learning analytics to develop individualised teaching and assessment materials tailored to current knowledge and learning performance
Energy, mining, manufacture	<ul style="list-style-type: none"> • Grid energy demand AI-systems to optimise energy distribution • Machine performance monitoring can increase efficiency
Human services	<ul style="list-style-type: none"> • Robotic process automation systems mimic user behaviour to learn how to do multiple, non-straight forward tasks, such as identifying emailed invoices, reading non-standard fields, entering into an accounting system and filing. • Speech recognition and Natural Language Generation transcribe human language and interact with humans for customer service, support and engagement, and human resources
Business, finance and information technology	<ul style="list-style-type: none"> • Decision management systems to optimise performance, minimise risk, streamline operations and increase profits • Precision marketing matches targeted customers with products • Fraud detection and money laundering can be monitored with ML detection of unusual transaction activities • Cyber security systems: AI-driven cyber defence machines can now uncover suspicious user activity and detect up to 85 per cent of all cyber attacks
Agriculture and fisheries	<ul style="list-style-type: none"> • Plant disease monitoring and tracking • Water use optimisation • Harvest time optimisation with visual robotics • Wild animal pest prevention • Grading process of agriculture products • Water health monitoring of aquaculture farms
Government, justice and defence	<ul style="list-style-type: none"> • Welfare, employment and fraud detection • Personalised services to those assessed as high risk • Infrastructure monitoring can reduce maintenance costs and increase equipment lifespan (e.g. bridge monitors, automatic street light adjustments based on people movement) • New AI-driven regulatory compliance solutions are emerging that can automate processes and deliver comprehensive risk coverage • Judicial process assessment and case management software • Autonomous weapons to respond to threats in real time • Identification of cyber attacks

detect speech from anywhere in the room; Google-owned Nest is a thermostat that adjusts the room temperature to your heating or cooling needs; Netflix and Stan pre-select the movies and TV shows you are mostly likely to enjoy; Tesla's vehicles feature a myriad of uber-cool predictive capabilities, self-driving features and other customised luxuries; and, of course, Apple's Siri helps you manage your day while you're on the go, and wearable devices such as FitBit and Garmen keep track of your health and fitness, unlock your car and can even improve your golf swing.

Convenience and luxury aside, arguably the three greatest societal benefits of AI applications are:

- 1) Improvements in productivity, particularly in labour-intensive tasks, thus offering significant benefits to economic growth and development;
- 2) Improvements in road safety, with reports suggesting driverless vehicles can prevent up to 90 per cent of traffic fatalities; and
- 3) Improvements in health, wellbeing and life expectancy from the delivery of personalised medicines, early disease diagnosis, remote health services, patient monitoring and AI-driven technologies that can enable greater participation of people with disabilities.

Of course, many AI-driven applications may come with specific challenges and therefore require regulators to step in from various cross-sectors. Many of these challenges must be guided by ethics frameworks with key principles outlined in the CSIRO and Data61 discussion paper,^[13] which underlines that AI applications must 'do no harm' and generate net-benefits (e.g. the benefits outweigh costs). Thus, while there is intense (and indeed justified) debate surrounding autonomous vehicles and the prioritisation of lives in the event of an accident, this must be considered in the overall context of the lives that will be saved. However, there are real safety risks, both existential and individual, posed by autonomous weapons in a way not seen since the advent of nuclear weapons. The convergence, splitting and globalisation of the digital nature of AI, will necessitate intense collaboration across multiple regulators, different levels of government, and between nations.

From an economic perspective, Data61 analysis reveals that over the past few years, 14 countries and international organisations have announced AU\$86 billion for AI programs.^[13] As with all disrupting technologies, the shifting nature of work will produce winners and losers. New markets and products will bring with them jobs and increased projects, as

will large increases in efficiency and productivity through automation and reduced human error. Early adopting companies will gain a competitive advantage by better understanding their customers' needs and being able to respond with personalised products, services and prices. Here, governments must be mindful to support the small business culture underpinning the Australian economy to help them transition and benefit, and avoid a monopolisation by a handful of corporate giants.

Critically, estimates suggest that around half of activities performed in jobs, and between 21 per cent and 38 per cent of jobs in the developed world, may be lost as a result of an increasingly digitalised and automated economy.^[14] However, a recent study conducted in the United Kingdom estimates that countervailing displacement and income effects are likely to balance each other out over the next 20 years or so. Thus, as AI displaces traditional jobs to automation, reskilling and transition planning are required to create the jobs of tomorrow.

For the environment, AI-driven systems can create a number of positive impacts and help to tackle the most critical challenges such as climate change and pollution. Monitoring systems and feedbacks will improve energy efficiency and reduce emissions (e.g. equipment sensing with smart cooling/heating), minimise resource waste and pollution (e.g. smart watering and fertiliser for crops), predict and manage natural disasters, and can also be used to evaluate the impact of ecosystems services, which can in turn be used by environmental decision makers to understand and quantify environmental assets.^[15] AI-based modelling can assist in planning resource management decisions, and help manage disaster responses to natural catastrophes, such as predicting bush fire movement.

The perceived lack of regulation surrounding AI has seen it deemed as the new 'wild west'. The past few years have seen several government and research organisations throughout the world develop policy and regulatory responses, or ethical and regulatory frameworks to manage AI ethics and ensure the risks do not outweigh the benefits. These include the UK, EU, Germany, France, Canada, US, Singapore, Japan, India, and China.

Fortunately, Australia has policies and regulations in place that can be used and enhanced to include AI. These include privacy and data protection laws, as well as possibilities for legal redress for faulty products and harms caused by products or erroneous organisational decisions.

Some authors propose that the use of AI in decision-making should come with a label (a 'Turing

Stamp'), similar to food labelling. Alternatively, some jurisdictions require AI decision-making in government to pass an examination for fitness for purpose and ethical compliance. Rights-based frameworks building on recognised human rights and digital rights (of data protection and privacy) have also been offered. A discussion paper led by CSIRO and Data61 and funded by the Australian Government Department of Industry, Innovation and Science outlines an ethics framework for Australia. Within it, they offer core principles for AI and propose a toolkit for policy makers and regulators – it is a valuable resource which is well worth reading.^[13]

References

1. Agrawal, A., Gans, J., Goldfarb, A. (2019). Economic Policy for Artificial Intelligence. *Journal of International Policy & Economy*, 19(1), 139-159.
2. Kaplan, A., & Haenlein, M. (2019). Siri, Siri, in my hand: Who's the fairest in the land? On the interpretations, illustrations, and implications of artificial intelligence. *Business Horizons*, 62(1), 15-25.
3. Turing, A.M. (1950). Computing machinery and intelligence. *Machinery*, 59(236), 433.
4. Berg, P. (2006). Origins of the human genome project: why sequence the human genome when 96% of it is junk? *American Journal of Human Genetics*, 79(4), 603-605. doi:10.1086/507688.
5. van Nimwegen, K. J., van Soest, R. A., Veltman, J. A., Nelen, M. R., van der Wilt, G. J., Vissers, L. E., & Grutters, J. P. (2016). Is the \$1000 genome as near as we think? A cost analysis of next-generation sequencing. *Clinical Chemistry*, 62(11), 1458-1464.
6. Guerrini, C. J., Robinson, J. O., Petersen, D., & McGuire, A. L. (2018). Should police have access to genetic genealogy databases? Capturing the Golden State Killer and other criminals using a controversial new forensic technique. *PLOS Biology*, 16(10).
7. Office of the Information Commissioner. Privacy business resource 21: Australian businesses and the EU General Data Protection Regulation. 2018; Available from: <https://www.oaic.gov.au/agencies-and-organisations/business-resources/privacy-business-resource-21-australian-businesses-and-the-eu-general-data-protection-regulation>.
8. Salmon, P., Hancock, P., & Carden, T. (2019). To protect us from the risks of advanced artificial intelligence, we need to act now. *The Conversation*. Online. Accessed: 12 June 2019.
9. Maimon, O., & Rokach, L. (2005). *Data mining and knowledge discovery handbook*.
10. Långkvist, M., Karlsson, L., & Loutfi, A. (2014). A review of unsupervised feature learning and deep learning for time-series modeling. *Pattern Recognition Letters*, 42, 11-24.
11. Jaderberg, M., Czarnecki, W. M., Dunning, I., Marris, L., Lever, G., Castañeda, A. G.,... Ruderman, A. (2019). Human-level performance in 3D multiplayer games with population-based reinforcement learning. *Science*, 364(6443), 859-865.
12. Dunjko, V., & Briegel, H. J. (2018). Machine learning & artificial intelligence in the quantum domain: a review of recent progress. *Reports on Progress in Physics*, 81(7), 074001.
13. Dawson, D. and Schleiger, E., Horton, J., McLaughlin, J., Robinson, C., Quezada, G., Scowcroft, J., and Hajkowicz, S. (2019) *Artificial Intelligence: Australia's Ethics Framework*. Data61 CSIRO, Australia.
14. Bank, T. W. (2019). *World Development Report 2019: The Changing Nature of Work*.
15. Villa, F., Ceroni, M., Bagstad, K., Johnson, G., & Krivov, S. (2009). ARIES (Artificial Intelligence for Ecosystem Services): A new tool for ecosystem services assessment, planning, and valuation. Paper presented at the 11th annual BIOECON conference on economic instruments to enhance the conservation and sustainable use of biodiversity, conference proceedings. Venice, Italy.
16. Gemalto Inc. (2017). *Data Breach Level Index: Full year results are in....* <https://blog.gemalto.com/security/2018/04/13/data-breach-stats-for-2017-full-year-results-are-in/>. Accessed: 15 June 2019.

Recommended Further Reading

1. Eubanks, V. (2018). *Automating inequality: How high-tech tools profile, police, and punish the poor*. New York: St. Martin's Press.
2. Henman, P. (2010). *Governing electronically: E-government and the reconfiguration of public administration, policy and power*. New York: Springer.
3. Flaxman, S. & Goodman, B., (2017). *European Union Regulations on Algorithmic Decision Making and a "Right to Explanation"*. *AI Magazine* (Fall): 50-57.
4. Buchanan, B. & Miller, T. (2017). *Machine Learning for Policy Makers: What It Is and Why It Matters*. Belfer Center for Science and International Affairs, Harvard Kennedy School.
5. Lal Das, P., Beisswenger, S.C., Mangalam, S., Yuce, M.R., Lukac, M. (2017). *Internet of things: the new government to business platform – a review of opportunities, practices, and challenges* (English). Washington, D.C. World Bank Group.



Contributors: This publication was authored and edited by Jennifer Yarnold, Paul Henman, Christopher McEwan, Karen Hussey, and Amelia Radke. Date of publication: 30 June 2019. Important disclaimer: This content has been prepared by the UQ Centre for Policy Futures for the purpose of general education. It does not constitute professional advice, nor is it a substitute for such advice. Questions about the topics discussed may be directed to the UQ Centre for Policy Futures by email: policyfutures@uq.edu.au or phone: +61 7 3443 3118. © The University of Queensland. All material published in 'Policy Futures' is licensed under a Creative Commons – Attribution – Non-Commercial 4.0 International (CC BY-NC 4.0) licence. In essence, you are free to copy, distribute and adapt the work, as long as you attribute the work and abide by the other licence terms. To view a copy of this licence, visit: <http://creativecommons.org/licenses/by/4.0/>.

Sorting out the 'good' from the 'bad'

The Ethics of High-Tech Tools Making Decisions

AI is often sought to help classify individuals in sub-groups to identify the most appropriate responses to each person. For example, deciding what level of airline passenger screening a person must undergo, or the Chinese government's social credit system, can open up opportunities for 'good' citizens while blocking those categorised as 'bad'. Similarly, AI can generate predictions of an individual's circumstances by comparing their characteristics to a dataset of others. For example, they could identify the level of risk a child faces in suffering abuse or neglect.

Such processes of classification and prediction are very helpful for governments in better tailoring responses to people, instead of a 'one-size-fits-all' approach. They also enable resources to be better targeted and used more efficiently and effectively.

There are a wide range of areas in which AI is being used or on the verge of being used to assist government make decisions. Apart from those already mentioned, judicial processes in the US are making use of AI to help inform judges about the appropriate sentence for an offender, such as time in prison, non-custodial sentences or alternative correctional responses. They are also being used to help inform parole officers about the likelihood of reoffending when on parole, to help parole officers decide if a person should be released on parole. In security systems, AI based facial recognition is not just used to find a match for a 'person of interest' or to serve a warrant. AI enhanced CCTV systems have also been trialled for the London Tube system to identify people who may be about to suicide by jumping in front of a train by following their movements on the platform and comparing that with previous movement patterns of previous suicide attempts.

On face value, all of these current and emerging uses of AI appear to be of considerable value. However, there has also been key ethical, legal and technical concerns about their use.

Data bias is at the heart of much criticism of the uses of AI in juridical sentencing and parole decisions, and similarly with child abuse and neglect detection. In the former, the data is based on a historical racial bias and so the AI, if not developed carefully, will continue to reflect, reproduce and even exacerbate such bias. In child protection, the bias is about poverty and disadvantage, whereas just being poor and using public (rather than private) services can misleadingly label a child as at risk of abuse or neglect. Indeed, when Amazon used AI for hiring decisions it found such a strong gender bias that the company ceased its use. The lack of transparency of the AIs in use is also a repeated concern. It is quite common for AIs to be developed by commercial companies and used in a modified 'off the shelf' basis, with limited configuration to the specificities of the use location. As commercial products, the assumptions and data used to develop the AI remains inaccessible behind 'commercial in confidence' protections.

The black boxed nature of these systems also means that government employees and people affected by AI decisions have limited ability to understand and question the AI decisions, thereby undermining public accountability and review processes. In relation to predictive AI whereby a system suggests something about a person in the future, it is often not realised that prediction is not the same as actuality. How we treat someone based on what may occur, rather than what has occurred, needs careful consideration. Otherwise we end up defining futures for people that they have little or no control over.

Workforce considerations also occur with the growing use of AI in government decision making. If an unskilled officer can use an AI, or it can operate autonomously and independently without human involvement, skilled professionals could be lost and with that the important role of human factors in developing, shaping and working with people, especially those most disadvantaged, can disappear.

Paul Henman, Centre for Policy Futures

