# CS7IS2: Artificial Intelligence

## Lecture 0: Intro and Logistics

Ivana.Dusparic@tcd.ie

# Lectures

› Thursdays 10-11 LB01, Fridays 2-3 LB08

› No lectures due to bank holidays:
  – Fri 17/03, Fri April 7/04

› Rescheduled lectures due to my timetable clash
  – Thur 23/02 -> Thur 16/02 9-10 (ie double lecture 9-11 on 16/02)
  – Thur 23/03 -> Thur 16/03 9-10 (ie double lecture 9-11 on 16/03)
  – Thur 30/03 -> Thur 06/04 9-10 (ie double lecture 9-11 on 06/04) - TBC

# Assignments

› Marks exam : coursework = 50% - 50%

› Supplemental exam: 100% exam

› Coursework: 2 assignments worth 25% each

› Exam: in-person 2 hours (week 36 of the academic year https://www.tcd.ie/calendar/), timetable release details on https://www.tcd.ie/academicregistry/exams/

› Deadlines
  – No extensions (apart from medical cert, note from tutor etc)
  – Late submissions: mark penalty 33% per day

› Plagiarism
  – https://libguides.tcd.ie/friendly.php?s=plagiarism/levels-and-consequences

# Assignments

› Marks exam : coursework = 50% - 50%

› Supplemental exam: 100% exam

› Coursework: 2 assignments worth 25% each

› Exam: in-person 2 hours (week 36 of the academic year https://www.tcd.ie/calendar/), timetable release details on https://www.tcd.ie/academicregistry/exams/

› Deadlines
  – No extensions (apart from medical cert, note from tutor etc)
  – Late submissions: mark penalty 33% per day

› Plagiarism
  – https://libguides.tcd.ie/friendly.php?s=plagiarism/levels-and-consequences (no, you can't copy from ChatGPT either!)

# Questions

› During the lecture
  – (not after the lecture finishes – lecture has to end 50 minutes past full hour)

› Blackboard discussion board
  – Ask questions – get answers from classmates, myself, module demonstrator (Jovan Jeromela)
  – Share any AI news

› Email

# Course Material

› Lecture notes and assignments will be posted on Blackboard

› Additional material (my slides are heavily based on these)

– Artificial Intelligence: A modern approach. Russel and Norvig, 4th edition, 2020

› http://norvig.com/  - link to pdf of the book etc

› https://people.eecs.berkeley.edu/~russell/

– UC Berkley CS188 module https://inst.eecs.berkeley.edu/~cs188/fa22/

– Artificial Intelligence: Foundations of Computational Agents, Pool and Mackworth. 2nd edition 2018
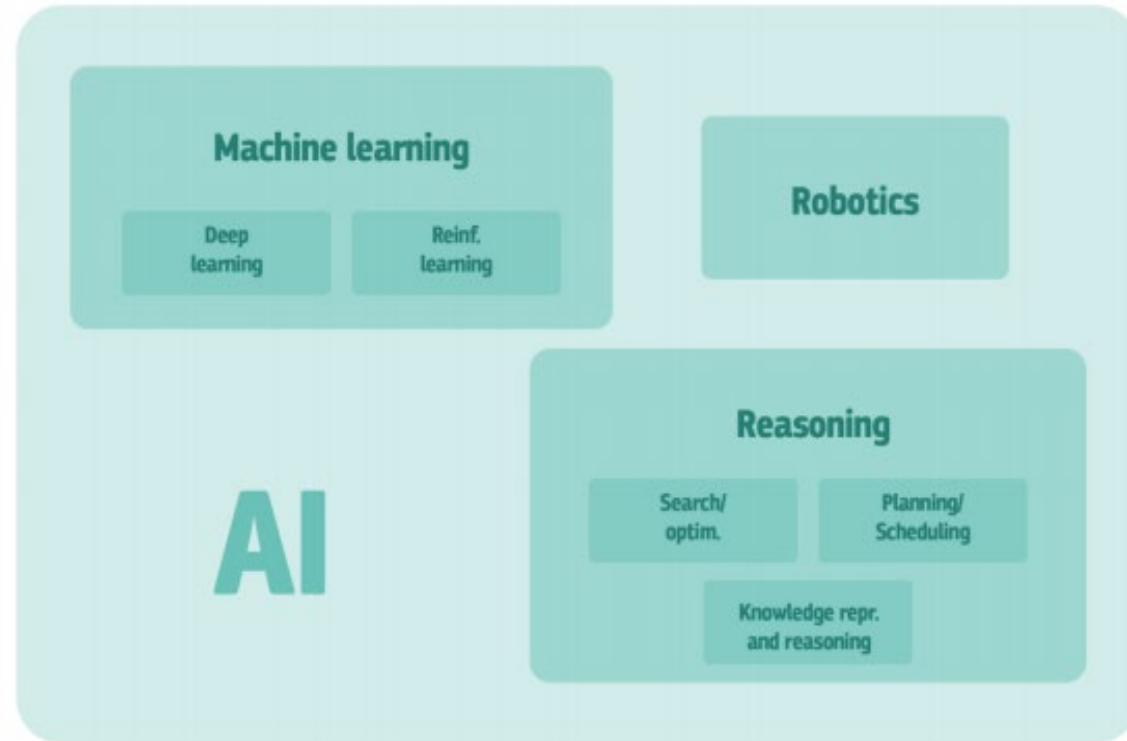
› https://artint.info/2e/html/ArtInt2e.html

# Some general AI reading – if interested

> Artificial Intelligence: A Very Short Introduction – Margaret Boden

> Rebooting AI: Building Artificial Intelligence We Can Trust – Gary Marcus and Ernest Davis

> Responsible Artificial Intelligence: How to Develop and Use AI in a Responsible Way – Virginia Dungam

> The Book of Why: The New Science of Cause and Effect – Judea Pearl and Dana Mackenzie

> Reinforcement Learning: An Introduction – Suton and Barto

> Human Compatible: AI and the Problem of Control – Stuart Russell

> Possible Minds: 25 Ways of Looking at AI – John Brockman

> Self Comes to Mind: Constructing the Conscious Brain - Antonio Damasio

> Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy – Kathy O'Neill

So what are we actually going to learn?

# What is AI?



› Image from: EC High-level expert group on AI - Definition, scope etc
  https://ec.europa.eu/digital-single-market/en/news/definition-artificial-intelligence-main-capabilities-and-scientific-disciplines

› Ethics guidelines
  https://ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai

# AI Ethics?



## US has 'moral imperative' to develop AI weapons, says panel

**Draft Congress report claims AI will make fewer mistakes than humans and lead to reduced casualties**

# Syllabus (subject to change)

› Problem Solving:
- – Searching
  - › Uninformed, Informed, Local
- – Adversarial search (multi-player games)
- – Constraint Satisfaction Problems

› Reinforcement Learning
- – MDPs
- – RL
- – Multi-agent systems

› Reasoning under uncertainty
- – Uncertainty
- – Bayes nets
- – Hidden Markov Models

› Intelligence from Computation
- – Search
- – Planning
- – CSP

› Intelligence from Data
- – Bayes nets
- – ML

# AI predictions and challenges

› Stanford AI Index 2022- https://hai.stanford.edu/research/ai-index-2022

## Language models are more capable than ever, but also more biased
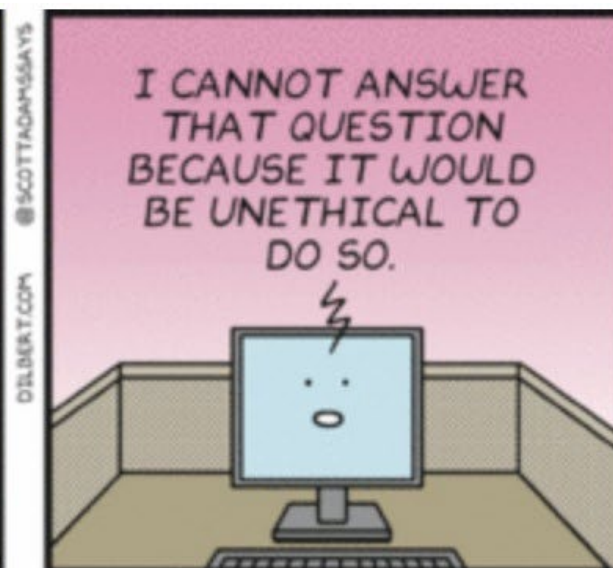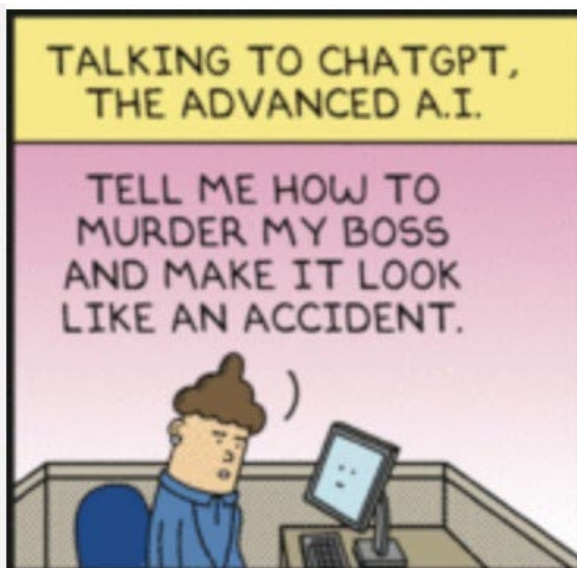
Large language models are setting new records on technical benchmarks, but new data shows that larger models are also more capable of reflecting biases from their training data. **A 280 billion parameter model developed in 2021 shows a 29% increase in elicited toxicity over a 117 million parameter model considered the state of the art as of 2018.** The systems are growing significantly more capable over time, though as they increase in capabilities, so does the potential severity of their biases.

## The rise of AI ethics everywhere

Research on fairness and transparency in AI has exploded since 2014, **with a fivefold increase in related publications** at ethics-related conferences. Algorithmic fairness and bias has shifted from being primarily an academic pursuit to becoming firmly embedded as a mainstream research topic with wide-ranging implications. **Researchers with industry affiliations contributed 71% more publications year over year** at ethics-focused conferences in recent years.

## AI becomes more affordable *and* higher performing

Since 2018, the cost to train an image classification system has decreased by 63.6%, while training times have improved by 94.4%. The trend of lower training cost but faster training time appears across other MLPerf task categories such as recommendation, object detection and language processing, and favors the more widespread commercial adoption of AI technologies.

# AI predictions and challenges

## More global legislation on AI than ever

An AI Index analysis of legislative records on AI in 25 countries shows that the number of bills containing "artificial intelligence" that were **passed into law grew from just 1 in 2016 to 18 in 2021**. Spain, the United Kingdom, and the United States passed the highest number of AI-related bills in 2021 with each adopting three.
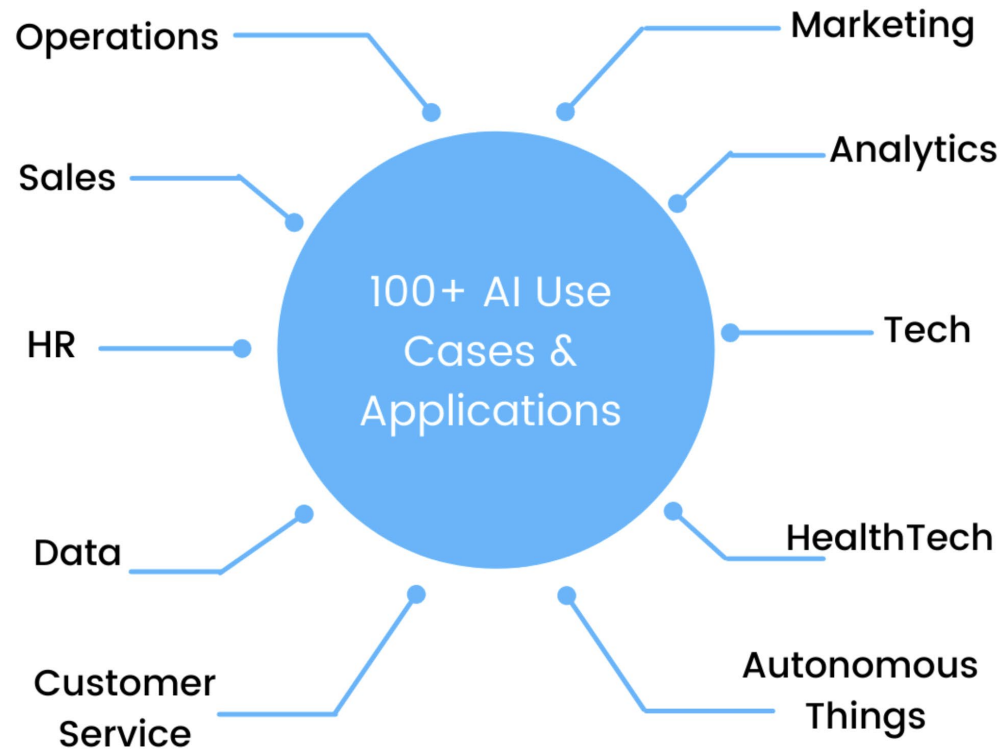
Proposed EU AI act
https://artificialintelligenceact.eu/

The AI Act is a proposed European law on artificial intelligence (AI) – the first law on AI by a major regulator anywhere. The law assigns applications of AI to three risk categories. First, applications and systems that create an **unacceptable risk**, such as government-run social scoring of the type used in China, are banned. Second, **high-risk applications**, such as a CV-scanning tool that ranks job applicants, are subject to specific legal requirements. Lastly, applications not explicitly banned or listed as high-risk are largely left unregulated.
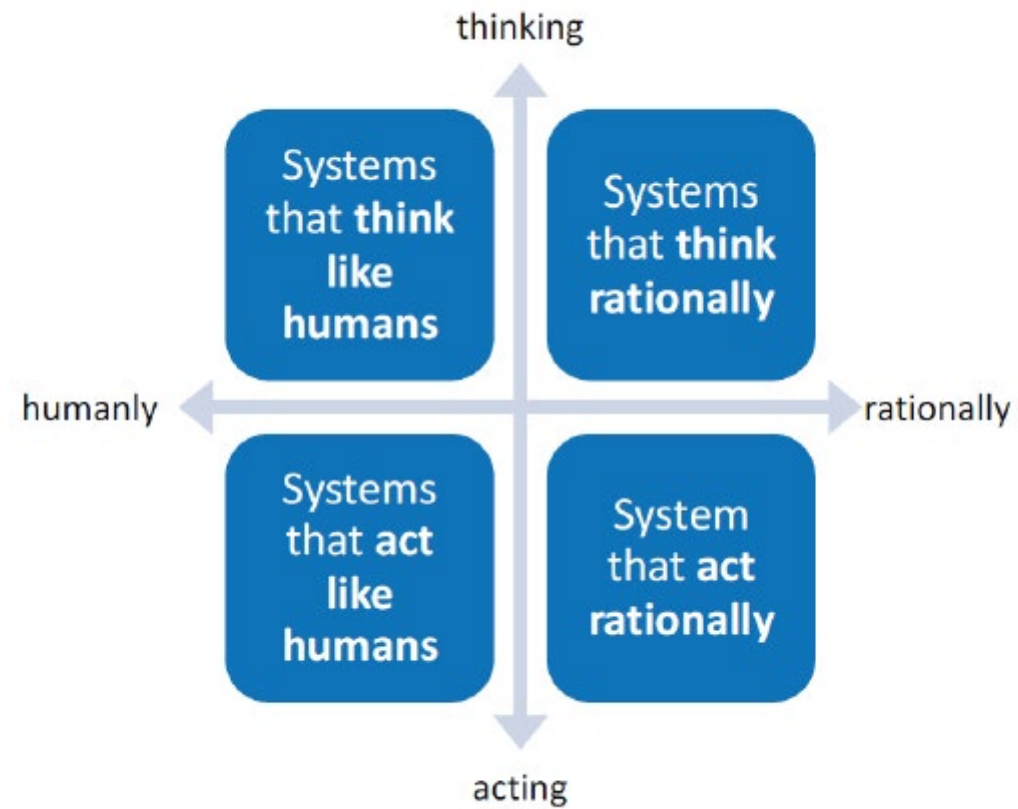
# AI predictions and challenges

https://research.aimultiple.com/

# AI predictions and challenges

› Quest for AGI

› Large-language models
› Facebook
  – Deep learning, self-supervised learning
› https://ai.facebook.com/blog/yann-lecun-advances-in-ai-research/

› Google  DeepMind
  – Reinforcement learning – "Reward is Enough"
› https://www.deepmind.com/publications/reward-is-enough
› Dreamer v3 https://arxiv.org/abs/2301.04104v1

# What is AI?

# Acting Humanely approach (1 of 2)

› Turing test approach
- – A computer passes the test if a human interrogator cannot tell whether the responses come from a person or computer

› However, more important to study underlying principles of intelligence than exactly duplicate the exemplar (ie a human)

› Progress in following research areas:
- – Natural language processing
- – Knowledge representation
- – Automated reasoning
- – Machine learning
- – Optional: computer vision, robotics

› Microsoft twitter bot

› GPT-3 **Generative Pre-trained Transformer 3** - language model that uses deep learning to produce human-like text

› DALL-E that creates images from text captions

# Acting Humanely approach (2 of 2)

› The **Winograd Schema Challenge**
  – The city councilmen refused the demonstrators a permit because they feared violence.
  – The city councilmen refused the demonstrators a permit because they advocated violence.

› Does the pronoun "they" refers to the city councilmen or the demonstrators? switching between the two instances of the schema changes the answer

# Thinking Humanely approach

› Cognitive science/cognitive neuroscience – understand how humans think

› Issues:
  – Requires scientific theories of internal activities of the brain
  – Also, humans often don't think (or act) in ways we consider intelligent
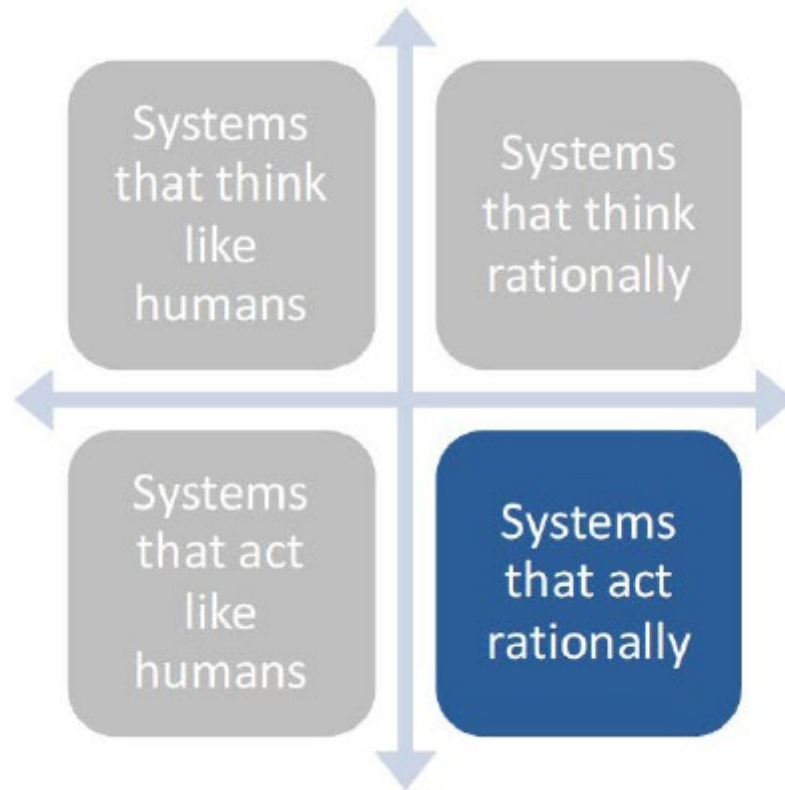
# Thinking Rationally approach

– Logic – patterns for argument structures

– Issues:
   › How to represent all knowledge using logical notation, especially uncertain knowledge
   › Solving a problem in theory vs in practice

# Acting Rationally approach

› Maximizing your expected utility/outcome

› Rational agent
  – Focus of AI today – general principles of rational agents and components for constructing them

"AI is the field that studies the synthesis and analysis of computational agents that act intelligently" (Poole & Mackworth, Artificial Intelligence: Foundations of Intelligent Agents)

# What is AI?

# Rational Agents

› An agent is an entity that perceives and acts in an environment

› An agent acts intelligently if:
– its actions are appropriate for its goals and circumstances
– it is flexible to changing environments and goals
– it learns from experience
– it makes appropriate choices given perceptual and computational limitations (finite memory and limited time)

# Rational agents

› This course is about designing rational agent

› Abstractly, an agent is a function from percept histories to actions:

$$f: \mathcal{P}^* \rightarrow \mathcal{A}$$

   – For any given class of environments and tasks, we seek the agent (or class of agents) with the best performance

› Computational limitations make perfect rationality unachievable

   – design best program for given machine resources

› Goal of AI: understand the principles that make intelligent behaviour possible

# Next: More about Agents