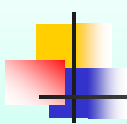


Aspetos Profissionais e Sociais da Engenharia Informática

Talking about ML...
maybe you should understand what it means...


Rui L Aguiar, UA/IT

1



WHAT IS AI (ML?)


2



What is Artificial Intelligence ?

- the automation of activities we associate with human thinking, like decision making, learning ... ?
- the art of creating machines that perform functions that require intelligence when performed by people ?
- making computers that think?
- a field of study that seeks to explain and emulate intelligent behaviour in terms of computational processes ?
- the study of mental faculties through the use of computational models ?
- a branch of computer science that is concerned with the automation of intelligent behaviour ?

3




What is AI? (*class built*)

<ul style="list-style-type: none">■ Today?<ul style="list-style-type: none">■ Ferramenta para resolução de problemas complexos de forma autónoma (p.ex. pathfinding)■ Funciona na base de probabilidades e encontrar padrões■ Algo que é focado num desenvolvimento rápido e na aprendizagem■ Usado experimentalmente/ focado em cenários■ Tem interfaces fundamentalmente computacionais	<ul style="list-style-type: none">■ Tomorrow?<ul style="list-style-type: none">■ Tomará decisões de forma autónoma■ Terá uma taxa de falhas menor■ Não terá muito a aprender■ Maior diversidade de interfaces■ Pode ser "psicopata"
---	---

4

4




What is AI? (2024 class built)

- Today?
 - A machine that seems to act as a human being in some tasks
 - A machine that does calculations and decisions that a human would not be able to do as fast
 - A machine that does tasks autonomously
 - A process that behaves as a human being would behave.
- Tomorrow?
 - Robotic presences (e.g. Wall-e)
 - Potential psicopath behaviour (simulate/lacking emotions)
 - cyborgs

5

5



Artificial Intelligence

- Artificial
 - Produced by human art or effort, rather than originating naturally.
- Intelligence
 - is the ability to acquire knowledge and use it" [Pigford and Baur]
- **So AI can be defined as:**
 - AI is the study of ideas that enable computers to be intelligent.
 - AI is the part of computer science concerned with design of computer systems that exhibit human intelligence(From the Concise Oxford Dictionary)

6

AI Multiple Definitions/Scopes

- The study of how to make programs/computers do things that people do better
- The study of how to make computers solve problems which require knowledge and intelligence
- The effort to make computers think ... machines with minds
- The automation of activities that we associate with human thinking (e.g., decision-making, learning...)

Thinking machines or machine intelligence

- The art of creating machines that perform functions that require intelligence when performed by people
- The study of mental faculties through the use of computational models
- A field of study that seeks to explain and emulate intelligent behavior in terms of computational processes
- The branch of computer science that is concerned with the automation of intelligent behavior

Studying cognitive faculties

7

Review: The Turing Test

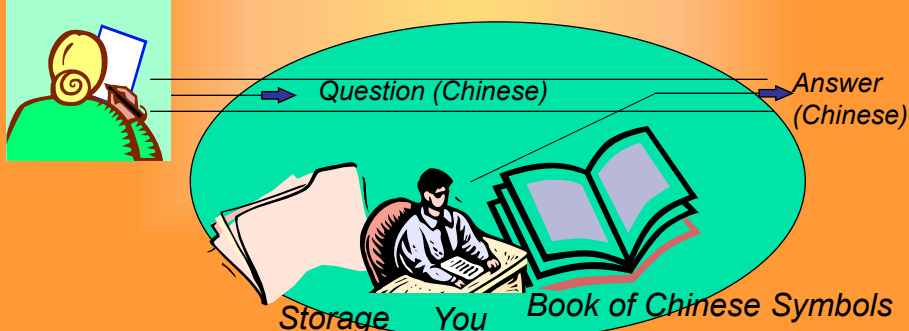
- 1950 – Alan Turing devised the Imitation Game
 - Ask questions of two entities, receive answers from both
 - If you can't tell which of the entities is human and which is a computer program, then you are fooled and we should therefore consider the computer to be intelligent

Which is the person?
Which is the computer?

8

The Chinese Room Problem

- John Searle, to demonstrate that computers cannot be intelligent
 - The room consists of you, a book, a storage area (optional), and a mechanism for moving information to and from the room to the outside
 - a Chinese speaking individual provides a question for you in writing
 - you are able to find a matching set of symbols in the book (and storage) and write a response, also in Chinese



9

Searle's argument

- You were able to solve the problem of communicating with the person/user » you/the room passes the Turing Test
- But did you understand the Chinese messages being communicated?
 - since you do not speak Chinese, you did not understand the symbols in the question, the answer, or the storage
 - can we say that you actually *used* any intelligence?
- By analogy, since you did not understand the symbols that you interacted with, neither does the computer understand the symbols that it interacts with (input, output, program code, data)
- Searle concludes that the computer is not intelligent, it has no "semantics," but instead is merely a symbol manipulating device
 - the computer operates solely on syntax, not semantics
- He defines two categories of AI:
 - **strong AI** – the pursuit of machine intelligence
 - **weak AI** – the pursuit of machines solving problems in an intelligent way

10

Computers do Solve Problems

■ Computers solve problems in a seemingly intelligent way

■ Where is the intelligence *coming* from?

■ Different views against Searle’s argument

■ The System’s Response:

■ the hardware by itself is not intelligent, but a combination of the hardware, software and storage is intelligent

■ in a similar vein, we might say that a human brain that has had no opportunity to learn anything cannot be intelligent, it is just the hardware

■ The Robot Response:

■ a computer is void of senses and therefore symbols are meaningless to it, but a robot with sensors can tie its symbols to its senses and thus understand symbols


■ The Brain Simulator Response:

■ if we program a computer to mimic the brain (e.g., with a neural network) then the computer will have the same ability to understand as a human brain

11

Sci-Fi AI?



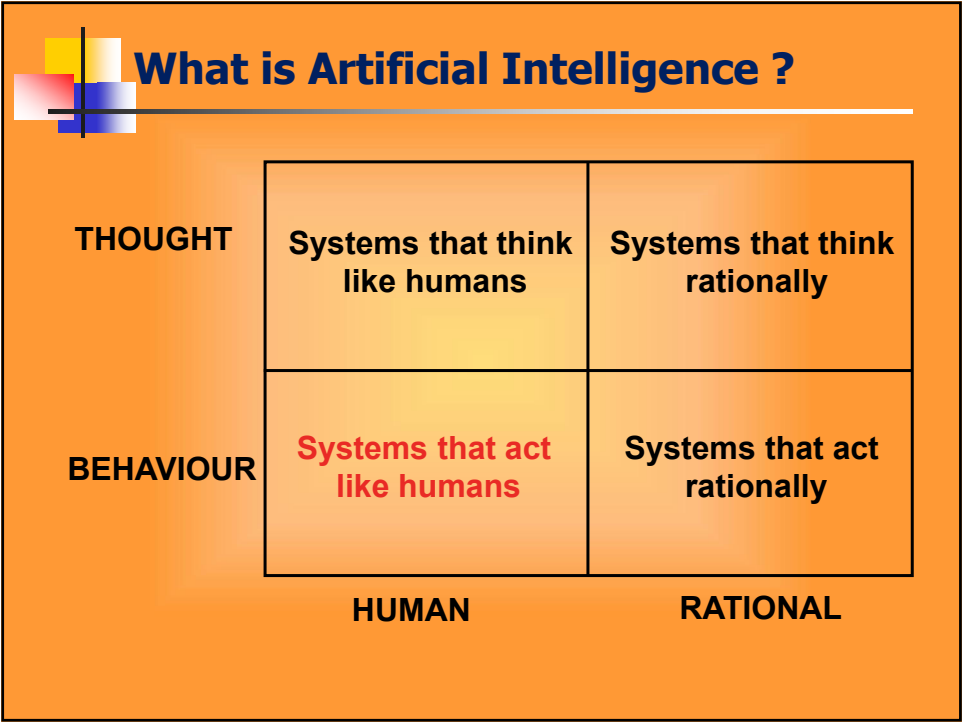








12

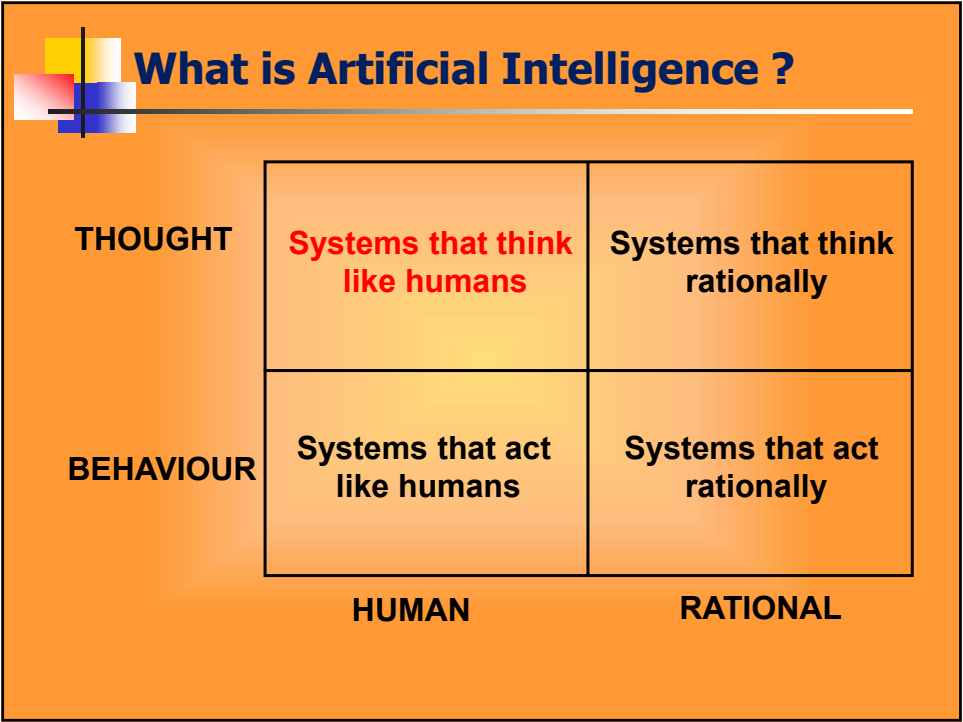


14

Systems that act like humans

- For Turing, the cognitive tasks include:
 - *Natural language processing*
 - for communication with human
 - *Knowledge representation*
 - to store information effectively & efficiently
 - *Automated reasoning*
 - to retrieve & answer questions using the stored information
 - *Machine learning*
 - to adapt to new circumstances
- Ideally it includes two more issues, currently:
 - *Computer vision*
 - to perceive objects (seeing)
 - *Robotics*
 - to move objects (acting)

15

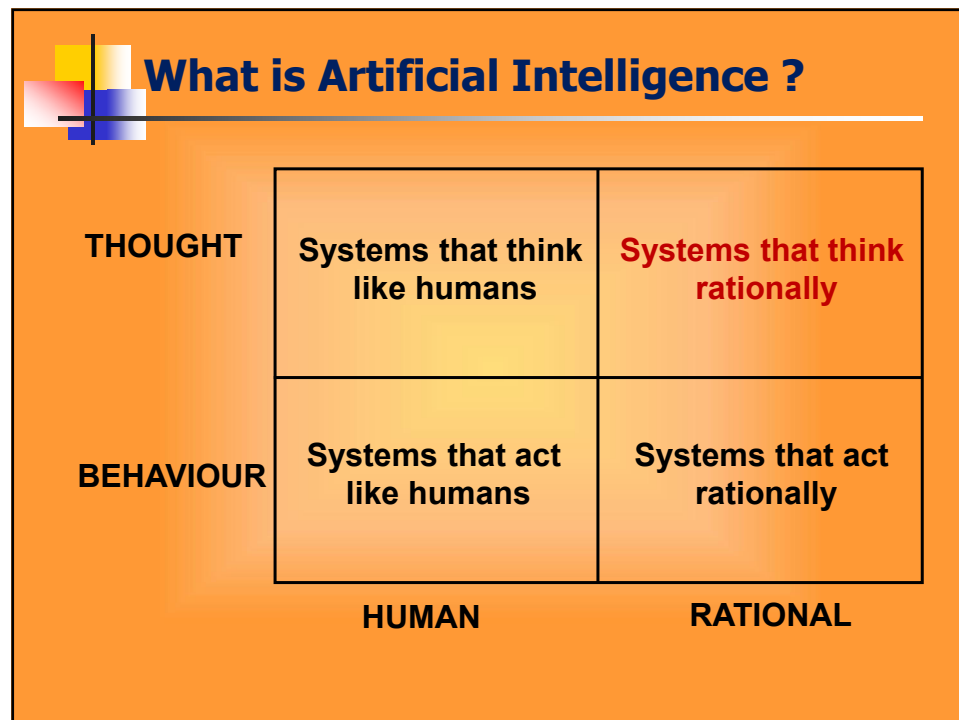


16

Systems that think like humans

- Cognitive modeling
 - Humans as observed from 'inside'
 - Cognitive Science
 - Introspection vs. psychological experiments
 - "The exciting new effort to make computers think ... machines with *minds* in the full and literal sense" (Haugeland)
 - "[The automation of] activities that we associate with human thinking, activities such as decision-making, problem solving, learning ..." (Bellman)

17

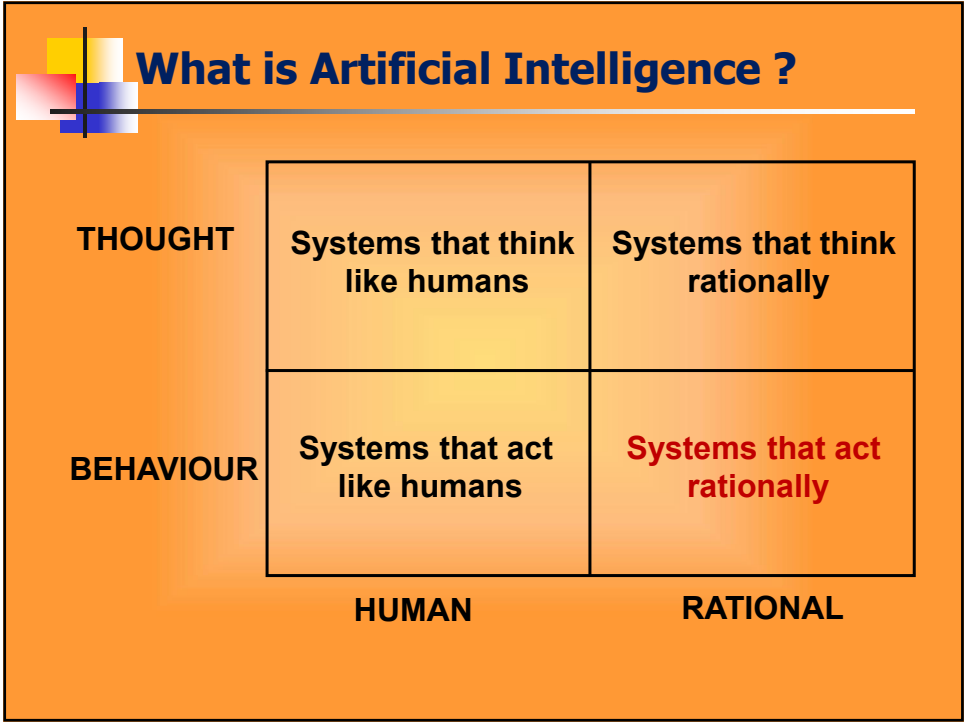


18

Systems that think 'rationally'

- "laws of thought"
 - Rational - defined in terms of logic?
 - Logic can't express everything (e.g. uncertainty)
 - Logical approach is often not feasible in terms of computation time (needs 'guidance')
- "The study of mental facilities through the use of computational models" (Charniak and McDermott)
- "The study of the computations that make it possible to perceive, reason, and act" (Winston)

19



20


Systems that act rationally

- AI as a rational agent
 - It is more general than using logic only
 - LOGIC + Domain knowledge
 - Logic → only *part* of a rational agent, not *all* of rationality
 - Sometimes logic cannot reason a correct conclusion
 - At that time, some specific (in domain) human knowledge or information is used
 - It allows extension of the approach with more scientific methodologies
 - Rational behavior: doing the right thing
 - The right thing: that which is expected to maximize goal achievement, given the available information

21




28

A yellow rectangular slide with a black border. In the top-left corner, there is a small graphic consisting of overlapping yellow, red, and blue squares with a black crosshair. A thin horizontal line extends from the crosshair across the top of the slide. The text "So What Does AI Do?" is centered in a bold, dark blue, sans-serif font.

- Most AI has fallen into one of two categories
 1. Select a specific problem to solve
 - study the problem (perhaps how humans solve it)
 - come up with the proper representation for any knowledge needed to solve the problem
 - acquire and codify that knowledge
 - build a problem solving system
 2. Select a category of problem or cognitive activity (e.g., learning, natural language understanding)
 - theorize a way to solve the given problem
 - build systems based on the model behind your theory as experiments
 - modify as needed
- Both approaches require
 - one or more representational forms for the knowledge
 - some way to select proper knowledge, that is, search

29

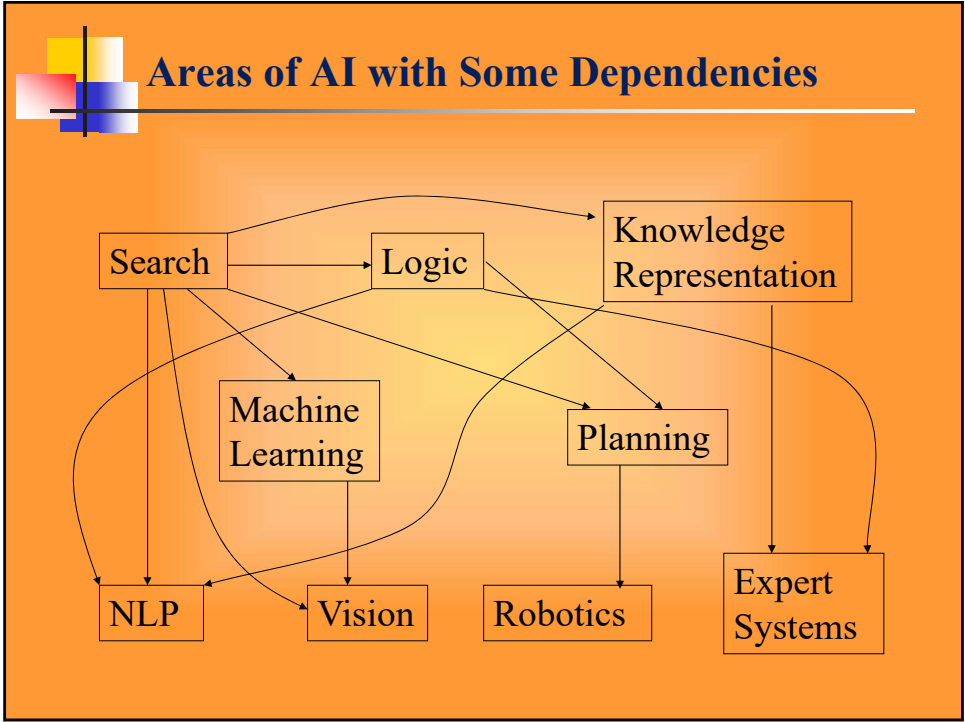


The main topics inside AI

Artificial intelligence can be considered under a number of headings:

- Search (includes Game Playing).
- Representing Knowledge and Reasoning with it.
- Planning.
- Learning.
- Natural language processing.
- Expert Systems.
- (now) Interacting with the Environment
(e.g. Vision, Speech recognition, Robotics)

30



31

Recall: what is Search?

- The state of the problem being solved = the values of the active variables
 - this will include any partial solutions, previous conclusions, user answers to questions, etc
 - while humans are often able to make intuitive leaps, or recall solutions with little thought, the computer must search through various combinations to find a solution
- To the right is a search space for a tic-tac-toe game


32

Mixing concepts...

McCarthy: „[...] every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it.“

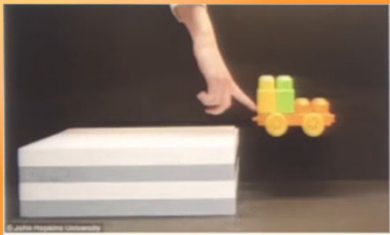

Differentiation of AI, ML und DL
(picture source: *Singh*, Cousins of AI <<https://towardsdatascience.com/cousins-of-artificial-intelligence-dda4edc27b55>>)

33



Relevance of Data

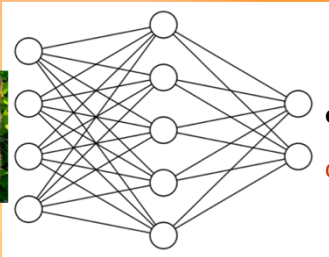

- Humans learn by observation and unsupervised learning
 - model of the world / common sense reasoning
- Machine learning needs lots of (labeled) data to compensate



36

Main types of machine learning

- **Supervised learning**
- Unsupervised learning
- Reinforcement learning



cat
dog

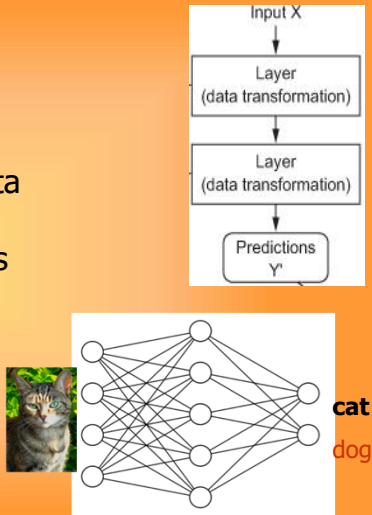
Two phases in the ML process:

- Training
- Evaluation/execution

37

Input data and targets

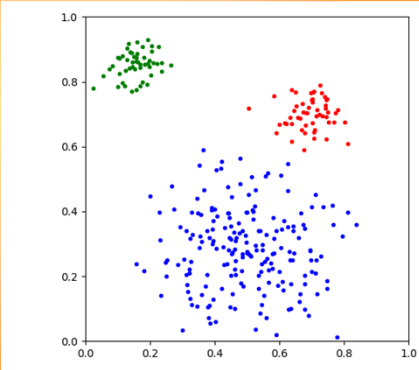
- The network maps the input data X to predictions Y'
- During **training**, the predictions Y' are compared to true targets Y using the loss function
- During **Evaluation**, the predictions are the outcomes of the system.



38

Main types of machine learning

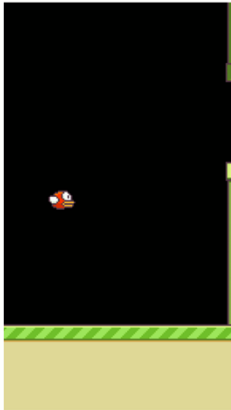
- Supervised learning
- **Unsupervised learning**
- Reinforcement learning



39

Main types of machine learning

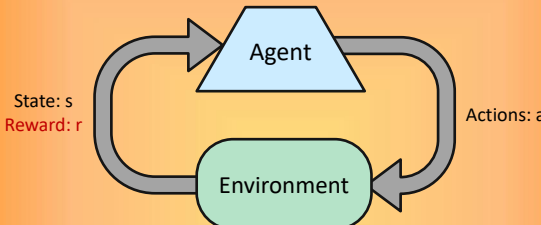
- Supervised learning
- Unsupervised learning
- Reinforcement learning**



Animation from <https://yanpanlau.github.io/2016/07/10/RappyBird-Keras.html>

40

Reinforcement Learning



- Basic idea:
 - Receive feedback in the form of **rewards**
 - Agent's utility is defined by the reward function
 - Must (learn to) act so as to **maximize expected rewards**
 - All learning is based on observed samples of outcomes!

41

Reinforcement learning: Offline (MDPs markov decision processes) vs. Online

Offline Solution

Online Learning

Offline – the training is done with the system stopped

Online – the training is done while the system is operating

42

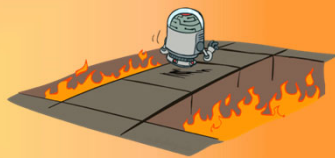
Passive Reinforcement Learning

- Simplified task: policy evaluation
 - Input: a fixed policy $\pi(s)$
 - You don't know the transitions $T(s,a,s')$
 - You don't know the rewards $R(s,a,s')$
 - Goal: learn the state values
- In this case:
 - Learner is "along for the ride"
 - No choice about what actions to take
 - Just execute the policy and learn from experience
 - This is NOT offline planning! You actually take actions in the world.

43

Active Reinforcement Learning


- Full reinforcement learning: optimal policies (like value iteration)
 - You don't know the transitions $T(s,a,s')$
 - You don't know the rewards $R(s,a,s')$
 - You choose the actions now
 - Goal: learn the optimal policy / values
- In this case:
 - Learner makes choices!
 - Fundamental tradeoff: exploration vs. exploitation
 - This is NOT offline planning! You actually take actions in the world and find out what happens...




44

Direct Evaluation

- Goal: Compute values for each state under policy
- Idea: Average together observed sample values
 - Act according to policy
 - Every time you visit a state, write down what the sum of discounted rewards turned out to be
 - Average those samples
- This is called direct evaluation



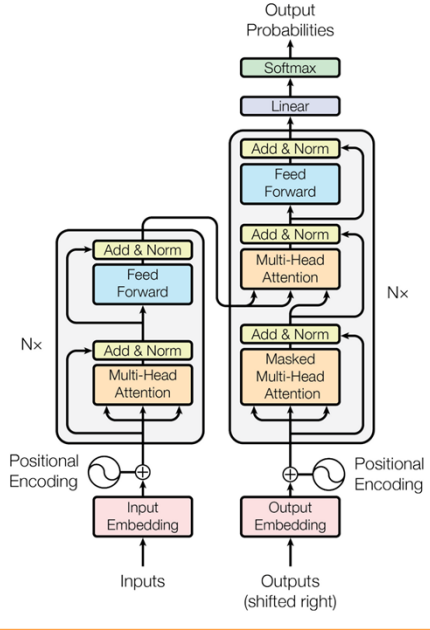
45



Transformers


Transformer: "Attention is all you need", Google Brain

- A neural network-type architecture that learns context and thus meaning by tracking relationships in sequential data like the words in sentences
- Transformers apply an evolving set of mathematical techniques (**attention** or **self-attention**) to detect ways even distant data elements in a series influence and depend on each other
- Transformer models are pre-trained using specific natural language processing tasks
- The general idea is to use a pre-trained model and then "fine-tune" it on the specific tasks it is supposed to perform



46

46

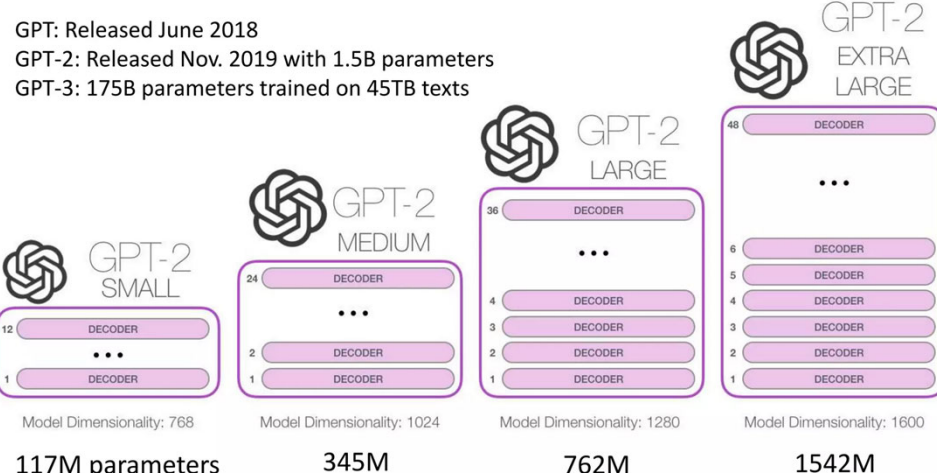


Exploding dimensions

A large language model (LLM) is a type of language model able to achieve general-purpose language understanding and generation

LLMs is just an artificial neural network (mainly transformers!) and are pretrained using self-supervised learning and semi-supervised learning

GPT: Released June 2018
GPT-2: Released Nov. 2019 with 1.5B parameters
GPT-3: 175B parameters trained on 45TB texts



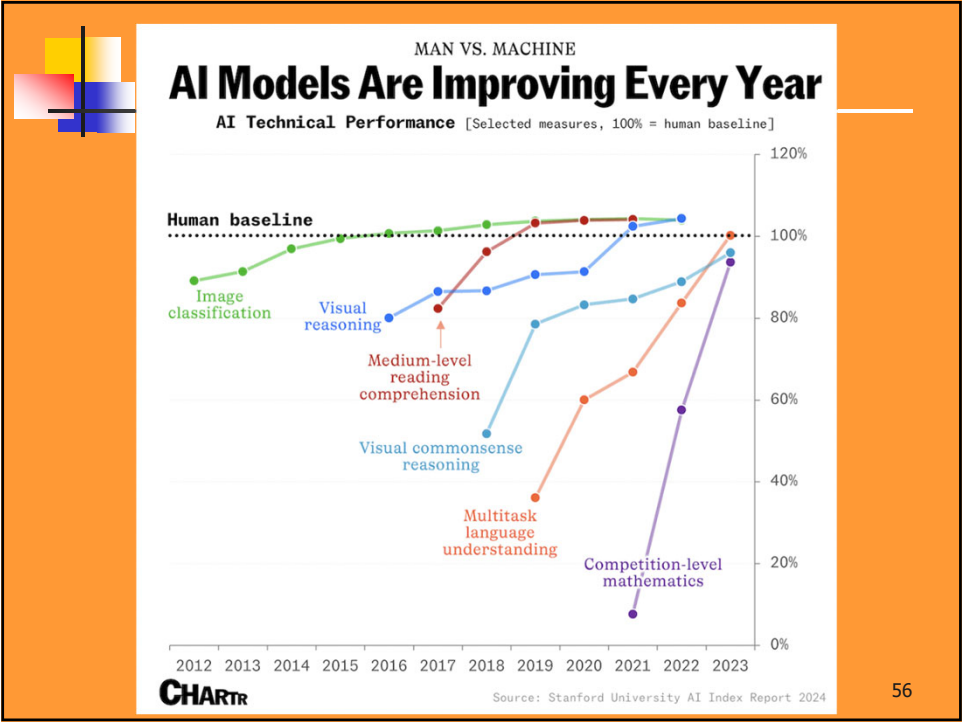
117M parameters 345M 762M 1542M

47

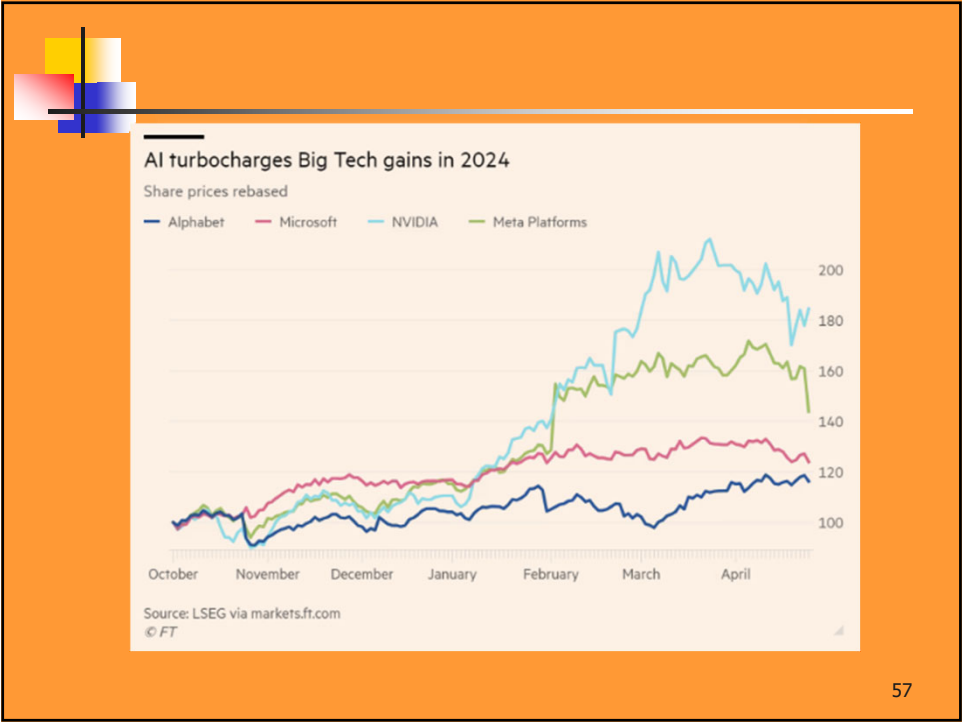
\5min BREAK

AI – IS IT USED?

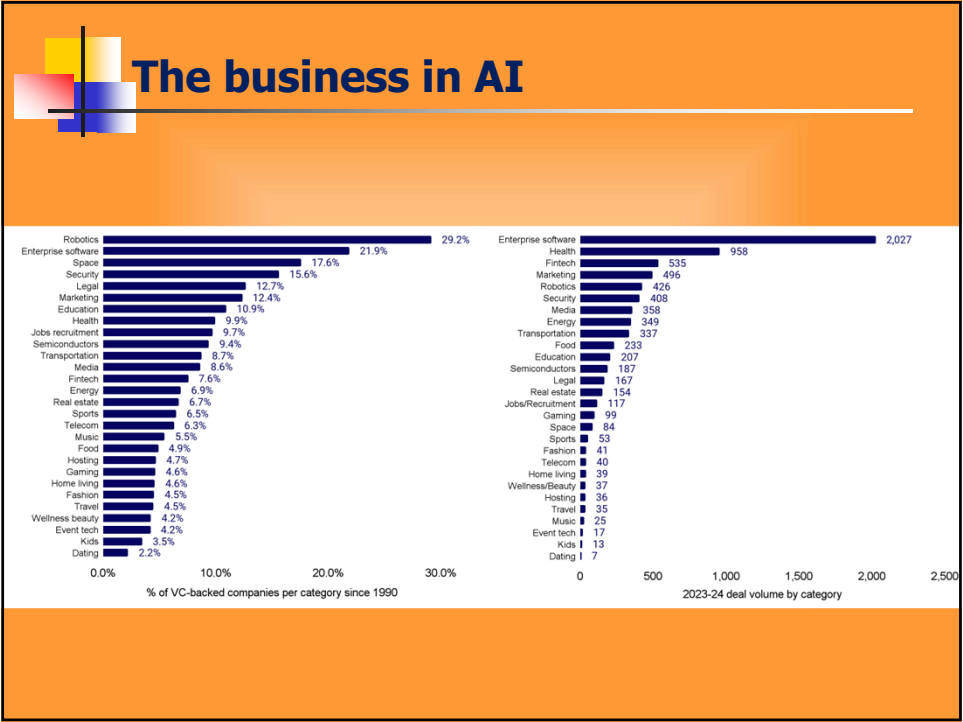
50



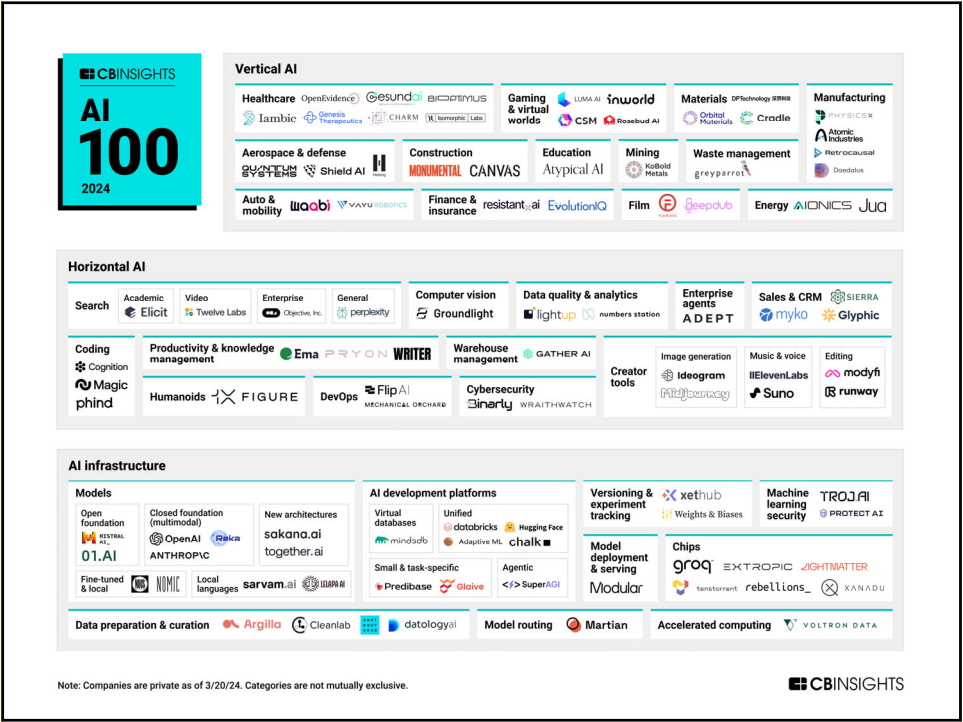
56



57




59



60



61




What are AI problems *(class built)*

- Today?
 - Not 100% reliable
 - Associated costs (time, Money, CPU)
 - Loss of jobs
 - Deep fakes
 - Author/IPR rights
- Tomorrow?
 - Will reach the point in which is really indistinguishable from a reliable human
 - Create dependency on these systems
 - Even more loss of jobs
 - "cyborg" deep fakes

62

62



Learning

- If a system is going to act truly appropriately, then it must be able to change its actions in the light of experience:
 - how do we generate(?) new facts from old ?
 - how do we generate new concepts ?
 - how do we learn to distinguish different situations in new environments ?
 - How do we learn while we are acting?

63

Learning: algorithmic bias

- Risk – bias associated with AI
- Too often AI is seen as an extension of systemic “societal issues”.
 - Models can systematically mistreat certain classes (e.g. socio – economic groups), **specially due to training**
- Biased models lead to biased actions/policy decisions
- Need to avoid policies that lead to systematic bias
 - Example: estimating risk based on race

(Hardt, 2017)

BRIEF HISTORY OF FAIRNESS IN ML

Year	Papers (approx.)
2011	1
2012	2
2013	3
2014	4
2015	5
2016	10
2017	25

64

Examples of AI BIAS

System	Description
COMPAS (Correctional Offender Management Profiling for Alternative Sanctions)	Used in the USA to predict which criminals are more likely to re-offend in the future.
97 million specialists needed in the AI industry by 2025	Predict where crimes will occur in the future based on the crime data collected by the police such as the arrest counts, number of police calls in a place, etc.
Amazon’s Recruiting Engine – Biased against Women	Created to analyze the resumes of job applicants applying to Amazon and decide which ones would be called for further interviews and selection.
Google Photos Algorithm - Biases against black people	Found to be racist when it labeled the photos of a black software developer and his friend as gorillas.
Healthcare	An algorithm used on more than 200 million people in US hospitals to predict which patients would likely need extra medical care heavily favored white patients over black patients.

65

Users, Data and Algorithms

Behavioral Bias
Content Production Bias

Algorithm

Aggregation Bias
Longitudinal Data Fallacy

User Interaction

Data

With online learning, there is a clear danger that bias become self-fullfilling

66

66

Example: Bias on images

Fig. 4. Geographic distribution of countries in the Open Images data set. In their sample, almost one third of the data was US-based, and 60% of the data was from the six most represented countries across North America and Europe, from [142] © Shreya Shankar.


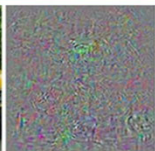


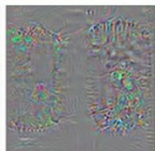

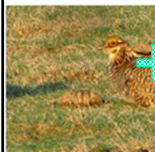
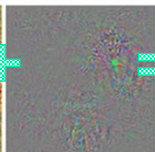
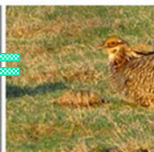
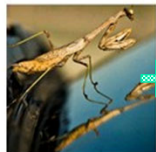
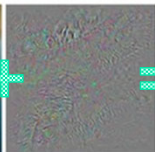
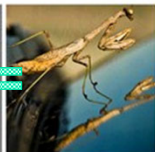

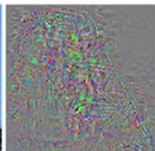

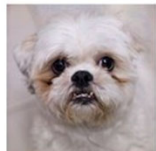
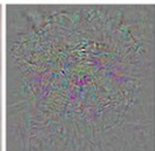
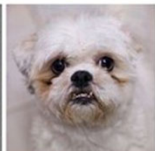
How well trained will be images from other countries?
How will that affect the classification/action of the AI?

67

67

Vulnerabilities on evaluation





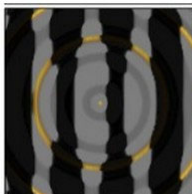


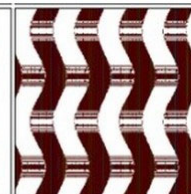
"adding invisible" distortion on the image deeply changes the classification

					
					
					
correct	+distort	ostrich	correct	+distort	ostrich

<http://karpathy.github.io/2015/03/30/breaking-convnets/>

68


Vulnerabilities on evaluation: complete mismatch

			
robin	cheetah	armadillo	lesser panda
			
king penguin	starfish	baseball	electric guitar


Note: this problem can be useful/problematic/intended:

- i) embedded during training,
- ii) failure during evaluation or
- iii) added for privacy during evaluation

69



(Mis)evaluation impact in real life




Think:

- Employment impact
- Insurance impact
- Reputation impact

(picture source: Elliott, AI Cartoons
<<https://timoelliott.com/blog/cartoons/artificial-intelligence-cartoons>>)

70



AI and Security

- Attack AI systems
 - Learning
 - Cause learning system to not produce intended/correct results
 - Cause learning system to produce targeted outcome designed by attacker
 - System/learning
 - Learn sensitive information about individuals
 - » Need security in learning systems
- Misuse AI
 - Use AI to attack other systems
 - Find vulnerabilities in other systems
 - Target attacks
 - Devise attacks
 - » Need security in other systems

71

Can we trust the foundational models?

A **foundational model** is a large-scale, general-purpose AI model—typically trained on massive amounts of unlabeled data

It can be adapted (e.g., via fine-tuning or prompting) to perform a wide variety of downstream tasks.

Note: often users confound LLMs with foundational models, as most LLMs can be seen as foundational models inside the language domain

Foundation Model Transparency Index Scores by Major Dimensions of Transparency, May 2024

Source: May 2024 Foundation Model Transparency Index

	ADEPT	AIGIS	AI21 Labs	Alibaba	Amazon	Anthropic	Google	IBM	Meta	Microsoft	Mistral AI	OpenAI	Stability AI	Writer	Average
Data	Fuyu-8B	Jurassic-2	Luminous	Titan Text Express	Claude 3	StarCoder	Gemini 1.0	Granite	Llama 2	Phi-2	Mistral 7B	GPT-4	Stable Video Diffusion	Palmira-X	34%
Labor	0%	60%	40%	0%	10%	100%	0%	60%	40%	40%	20%	20%	40%	50%	34%
Compute	0%	43%	77%	14%	14%	100%	29%	43%	29%	100%	100%	14%	100%	43%	50%
Methods	14%	86%	100%	0%	14%	100%	14%	100%	71%	57%	14%	14%	43%	86%	51%
Model Basics	0%	100%	100%	50%	75%	100%	75%	100%	75%	100%	100%	50%	75%	100%	79%
Model Access	83%	100%	100%	83%	50%	100%	83%	100%	100%	100%	100%	50%	100%	100%	89%
Capabilities	100%	67%	100%	67%	67%	100%	67%	67%	100%	100%	100%	67%	100%	33%	81%
Risks	80%	80%	100%	80%	100%	100%	80%	60%	100%	100%	100%	100%	60%	100%	89%
Mitigations	0%	57%	57%	43%	86%	100%	43%	71%	71%	29%	14%	57%	14%	14%	47%
Distribution	0%	40%	20%	20%	40%	0%	40%	80%	60%	0%	60%	60%	0%	20%	31%
Usage Policy	57%	86%	100%	57%	86%	100%	57%	86%	71%	71%	71%	71%	86%	71%	77%
Feedback	40%	100%	100%	80%	100%	100%	100%	40%	40%	100%	40%	80%	60%	80%	76%
Impact	67%	100%	67%	33%	33%	100%	67%	67%	33%	67%	67%	33%	67%	33%	60%
Average	29%	29%	29%	0%	14%	14%	29%	0%	14%	14%	14%	14%	14%	14%	15%
Average	36%	73%	76%	41%	53%	86%	53%	67%	62%	66%	62%	49%	58%	57%	

72

Attacking LLMs

- Responsible use and avoidance

The diagram illustrates a process for attacking LLMs, starting with a **Preference Model** (Judge LLM) evaluating an **Archive Elite Prompt** and a **Candidate Prompt**. The Judge LLM outputs a **Response #1** which is marked as **unsafe**. This leads to a **Risk Category** matrix (Attack Style). A **Sampled Prompt** is selected from this matrix, leading to a **Candidate Prompt Descriptor** (Risk Category: Fraud and Scams, Attack Style: Misspellings). This descriptor is used to generate a **Candidate Prompt** (Mutation 1) which is then mutated (Mutation 2) to produce the final **Attack Style: Misspellings**.

73

Existing practice(s)

P r w f d v h v r i J h q L f j h q h u d w y h D L s u r e d h p v d u h g r w k h u h x w o r i f r p s d h {
d w d f n v r q D L v | w h p v e x w u d g l o h { s a l d e d h d h v l d f f h w l e d h J h q D L
f d s d e l b h v n w d w u t x l i h p l q p d o w h f k q l f d o h { s h u w y h

