# What's Intelligence Got To Do With It?

## A Brief History and Overview of Artificial Intelligence Research

Tanya Braun

Data Science Group
Computer Science Department

living.knowledge

WWU
MÜNSTER

# Disclaimer

Slides are an amalgamation of books, articles, slides, and discussions with colleagues, most notably, by and with Ute Schmid, Subbarao Kambhampati, Stuart Russell, Malte Schilling, Ralf Möller, and Marcel Gehrke

Thank you!

living.knowledge

# Broad Goals of the Talk Today

- We will talk about
  - History of AI: Knowledge-driven vs. data-driven AI
  - Approaches to AI: Thinking or acting humanly or rationally
  - Future of AI: human-centric and hybrid AI?
  - ✓ Get an idea of *where it's been, what it's doing & where it's going – maybe*

- This talk cannot provide
  - Complete overview of all the methods that fall under AI methods
  - Tutorial on how to use machine learning techniques for geodynamics
  - In-depth explanation of ChatGPT and Large Language Models (LLMs)

# Where it's been

Artificial Intelligence Research

Data Science Grouo
Computer Science Department

living.knowledge

WESTFÄLISCHE WILHELMS-UNIVERSITÄT MÜNSTER | UNIVERSITY OF MÜNSTER

# Artificial Intelligence (AI)

- Computer science field
  - Inception: 1956 (John McCarthy, Stanford)
  - Defined by McCarthy as "*the science and engineering of making intelligent machines*"
- Most computer programs do not rely on AI
- Using AI methods means <u>giving up on completeness and correctness</u>
  - Reasonable to use AI methods if
    - Problem so complex that optimal solution cannot be efficiently computed → heuristic methods, approximations
    - Problem cannot be (completely) specified → replace explicit algorithms with models / programs learned from data (black box)

Computer Science

Artificial Intelligence

Machine Learning

Deep Learning

Data Science

Ute Schmid: Lernen über und Lernen mit KI. Invited talk, University:Future Festival 2023. (German)

# A Bit of History: 1$^{st}$ Wave of AI

- Focus: Explicit knowledge representation
  - Also called intelligent design
    - Figure out what you want, encode knowledge explicitly in some representation, tell computer how to manipulate representation to get what you want
    - Started out logic-based
  - Constitutes powerful inference methods, provable properties, comprehensibility
- Problem: Polanyi's Paradox
  - "We know more than we can tell"
  - Large part of knowledge not verbalisable → only implicitly available
  - Focus on explicit knowledge tasks instead of tacit knowledge tasks
  - Brittle models
    - World too complex
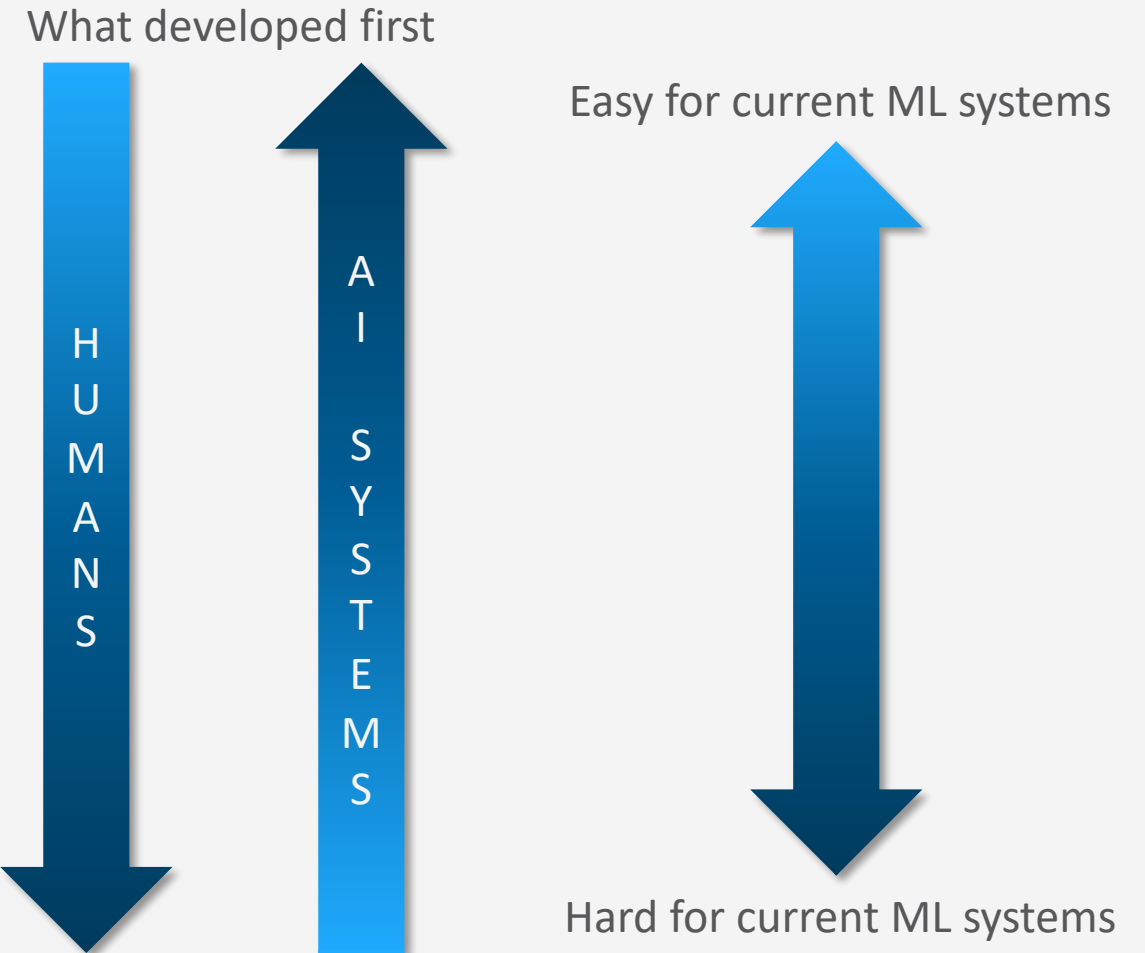
Knowledge-driven or symbolic AI

# A Bit of History: 2$^{nd}$ Wave of AI

- Focus: Data-intensive machine learning
  - Use available data to learn a model representing tacit knowledge
    - Show the computer lots of examples of inputs together with the desired outputs. Let the computer learn how to map inputs to outputs using a general purpose, learning procedure
    - Took off around 2012 probably
  - Impressive results, especially in image-based classification
- But 1: Huge effort to get a large amount of high-quality data
  - Also applies here: *Garbage in – garbage out*
  - Especially a problem in highly specialised areas such as medical computer science
    - E.g., we are currently looking at a data set of 150 data points with 350,000 features each
- But 2: Limited explainability / comprehensibility of very complex models

Data-driven or neural AI

# The Many Intelligences

- Perceptual & manipulation intelligence
  - Image recognition; hand-eye coordination
  - Largely tacit / subconscious
- Emotional intelligence
  - Showing & recognising emotional responses
- Social & communicative intelligence
  - Language
  - Requires a "theory of mind"
- Cognitive / reasoning intelligence
  - Hopefully, what we get tested for in uni
  - More declarative / consciously accessible

What developed first

HUMANS

AI SYSTEMS

Easy for current ML systems

Hard for current ML systems

Subbarao Kambhampati: "Polanyi vs. Planning (Planning around AI's New Romance with Tacit Knowledge)", invited talk, DC, ICAPS 2020.

# Why Did AI Develop in "Reverse"?

- It is easier to program computers on aspects of intelligence for which we have conscious theories (Polanyi's Paradox)
  - Ergo the progress in reasoning / cognitive intelligence during the 1$^{st}$ wave of AI
- We are not particularly conscious of perceptual (and manipulative) intelligence
  - We had to depend on making machines learn the way we had to
    - Learn from data / demonstrations

Subbarao Kambhampati: "Polanyi vs. Planning (Planning around AI's New Romance with Tacit Knowledge)", invited talk, DC, ICAPS 2020.

If the representations are learned, how do we ensure that they are understandable to the humans?

# Inference vs. Learning Focus

**Inference**

- Start by assuming models available
  - State / action representation etc.
- Focus on inference in context of model
  - Promise of eventually learning / updating of models
    - Postpones learning; reasonable for explicit knowledge domains with good models (Chess, Sudoku, mission planning…)
- AI development followed this direction for much of its history

**Learning**

- Assume that the agent does not have any a priori models
- Focus on learning (even primitive) models
  - Typically reflex agents
  - Promise of eventually getting to inference
    - Postpones inference; reasonable for tacit knowledge domains with no good models but a lot of examples / experience generators (vision, NLP, etc. …)
- Significant recent progress

Subbarao Kambhampati: "Polanyi vs. Planning (Planning around AI's New Romance with Tacit Knowledge)", invited talk, DC, ICAPS 2020.

# What it's doing

Artificial Intelligence Research

living.knowledge

Data Science Grouo
Computer Science Department

# A Bouquet of AI Methods

- Problem-solving
  - Search algorithms, heuristics, game theory, constraint satisfaction problems, …
- Logic
  - Propositional logic, description logic, ontologies, knowledge graphs; inference
- Uncertainty
  - Probabilistic modelling and inference (over time), utility and decision theory, multi-agent systems
- Machine learning
  - Learning from examples; neural networks, deep learning, reinforcement learning, …
- Perceiving and acting
  - Natural language processing, computer vision, robotics

# Approaches to Artificial Intelligence (AI)

- All approaches researched
  - Supported and hindered each other

- Rationality
  - System is rational if it does the "right thing," given what it knows

Success measure

| Fidelity of human performance | Ideal performance measure rationality | |
|---|---|---|
| **Thinking Humanly** | **Thinking Rationally** | Thought processes, reasoning |
| "The exciting new effort to make computers think . . . machines with minds, in the full and literal sense." (Haugeland, 1985)<br><br>"[The automation of] activities that we associate with human thinking, activities such as decision-making, problem solving, learning…" (Bellman, 1978) | "The study of mental faculties through the use of computational models." (Charniak and McDermott, 1985)<br><br>"The study of the computations that make it possible to perceive, reason, and act." (Winston, 1992) | |
| **Acting Humanly** | **Acting Rationally** | Behaviour |
| "The art of creating machines that perform functions that require intelligence when performed by people." (Kurzweil, 1990)<br><br>"The study of how to make computers do things at which, at the moment, people are better." (Rich and Knight, 1991) | "Computational Intelligence is the study of the design of intelligent agents." (Poole et al., 1998)<br><br>"AI … is concerned with intelligent behaviour in artefacts." (Nilsson, 1998) | |

Tanya Braun

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.
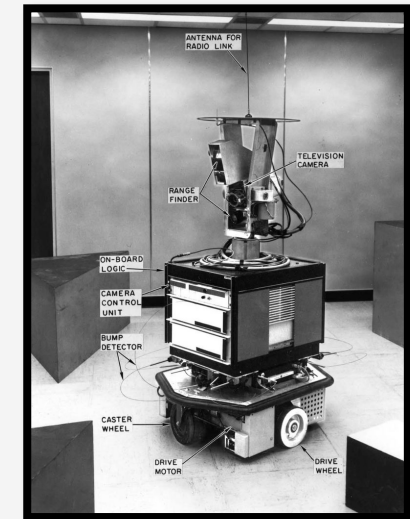
13

# Acting Humanly

- Turing Test (Turing, 1950)
  - Computer passes test, if a human, who asks written questions, cannot tell if the the written answers come from a human or not
    - Example: *Eliza*, program for superficially simulating a psychiatrist
    - See also Ch. 26, "Artificial Intelligence – A Modern Approach" by Russel & Norvig, including a discussion whether a computer would really be intelligent if it passed
      - Regarding Eliza: human's example closure tendencies are more pronounced for emotional/social intelligence aspects
      - Cf. robot *Shakey*: No on who saw Shakey the first time thought it could shoot hoops, yet the first people interacting with Eliza assumed it was a real doctor
  - *Total* Turing Test: includes a video signal to test perceptual abilities, opportunity to pass physical objects





Tanya Braun

# Acting Humanly

- Subproblems to solve as part of the Turing Test
  - *Natural Language Processing*
    - Communication
  - *Knowledge representation*
    - Store knowledge and observations
  - *Automated reasoning*
    - Answer questions, draw new conclusions
  - *Machine learning*
    - Adapt to new circumstances, detect and extrapolate patterns
  - Total Turing Test
    - *Computer vision*: perceive objects
    - *Robotics*: manipulate objects, move about

The Turing Test covers a majority of disciplines that make up AI nowadays.
- But:
  little research effort devoted to pass test
- Instead:
  Study underlying principles of intelligence

# Approaches to Artificial Intelligence (AI)

- All approaches followed
  - Supported and hindered each other

- Rationality
  - System is rational if it does the "right thing," given what it knows

Success measure

|  | Fidelity of human performance | Ideal performance measure rationality |  |
|---|---|---|---|
| **Thinking Humanly** | "The exciting new effort to make computers think . . . machines with minds, in the full and literal sense." (Haugeland, 1985) | **Thinking Rationally** | "The study of mental faculties through the use of computational models." (Charniak and McDermott, 1985) |
|  | "[The automation of] activities that we associate with human thinking, activities such as decision-making, problem solving, learning…" (Bellman, 1978) |  | "The study of the computations that make it possible to perceive, reason, and act." (Winston, 1992) |

Thought processes, reasoning

| **Acting Humanly** | "The art of creating machines that perform functions that require intelligence when performed by people." (Kurzweil, 1990) | **Acting Rationally** | "Computational Intelligence is the study of the design of intelligent agents." (Poole et al., 1998) |
|---|---|---|---|
|  | "The study of how to make computers do things at which, at the moment, people are better." (Rich and Knight, 1991) |  | "AI … is concerned with intelligent behaviour in artefacts." (Nilsson, 1998) |

Behaviour

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

# Thinking Humanly

- A "program thinks like a human"
  - Requires a way to determine how humans think → workings of the human mind
  - Given theory of the mind, express theory as computer program
    - If program's input-output behaviour matches corresponding human behaviour, evidence that some of program's mechanisms could also be operating in humans
- Approach complementary to AI: *Cognitive Science*
  - Interdisciplinary:
    - Computer models from AI
    - Experimental techniques from psychology
  - Goal:
    Construct precise and testable theories of human mind

# Approaches to Artificial Intelligence (AI)

- All approaches followed
  - Supported and hindered each other

- Rationality
  - System is rational if it does the "right thing," given what it knows

Success measure

Fidelity of human performance    Ideal performance measure rationality

| Thinking Humanly | Thinking Rationally |
|---|---|
| "The exciting new effort to make computers think . . . machines with minds, in the full and literal sense." (Haugeland, 1985)<br><br>"[The automation of] activities that we associate with human thinking, activities such as decision-making, problem solving, learning…" (Bellman, 1978) | "The study of mental faculties through the use of computational models." (Charniak and McDermott, 1985)<br><br>"The study of the computations that make it possible to perceive, reason, and act." (Winston, 1992) |
| Acting Humanly | Acting Rationally |
| "The art of creating machines that perform functions that require intelligence when performed by people." (Kurzweil, 1990)<br><br>"The study of how to make computers do things at which, at the moment, people are better." (Rich and Knight, 1991) | "Computational Intelligence is the study of the design of intelligent agents." (Poole et al., 1998)<br><br>"AI … is concerned with intelligent behaviour in artefacts." (Nilsson, 1998) |

Thought processes, reasoning

Behaviour

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

# Thinking Rationally

- Codify thinking → rules
  - Irrefutable reasoning processes
  - Argument structures that always yield correct conclusions when given correct premises
- Field of *Logic*
  - Precise notation for statements about objects in a world and relations among them
  - Programs that could, <u>in principle</u>, solve *any* solvable problem described in logical notation
  - Obstacles:
    - Informal knowledge
      - Unstructured data
      - Uncertainty
    - Solving any solvable problem <u>in practice</u>
      - Limited computational resources

Obstacles apply to *any* attempt to build computational reasoning systems
- Formulated first in logic

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

# Approaches to Artificial Intelligence (AI)

- All approaches followed
  - Supported and hindered each other

- Rationality
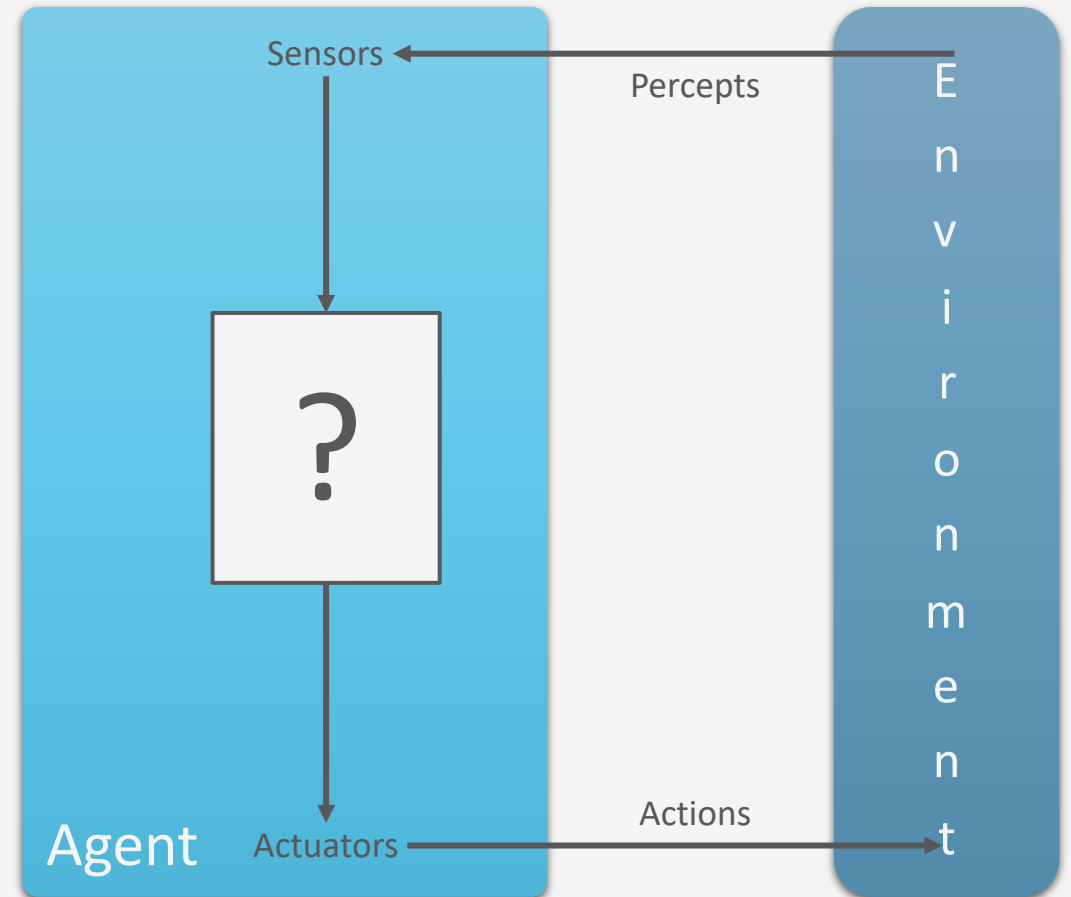  - System is rational if it does the "right thing," given what it knows

Success measure

| Fidelity of human performance | Ideal performance measure rationality | |
|---|---|---|
| **Thinking Humanly**<br><br>"The exciting new effort to make computers think . . . machines with minds, in the full and literal sense." (Haugeland, 1985)<br><br>"[The automation of] activities that we associate with human thinking, activities such as decision-making, problem solving, learning…" (Bellman, 1978) | **Thinking Rationally**<br><br>"The study of mental faculties through the use of computational models." (Charniak and McDermott, 1985)<br><br>"The study of the computations that make it possible to perceive, reason, and act." (Winston, 1992) | Thought processes, reasoning |
| **Acting Humanly**<br><br>"The art of creating machines that perform functions that require intelligence when performed by people." (Kurzweil, 1990)<br><br>"The study of how to make computers do things at which, at the moment, people are better." (Rich and Knight, 1991) | **Acting Rationally**<br><br>"Computational Intelligence is the study of the design of intelligent agents." (Poole et al., 1998)<br><br>"AI … is concerned with intelligent behaviour in artefacts." (Nilsson, 1998) | Behaviour |

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

# Acting Rationally

*Advantage*: Standard of rationality mathematically well defined
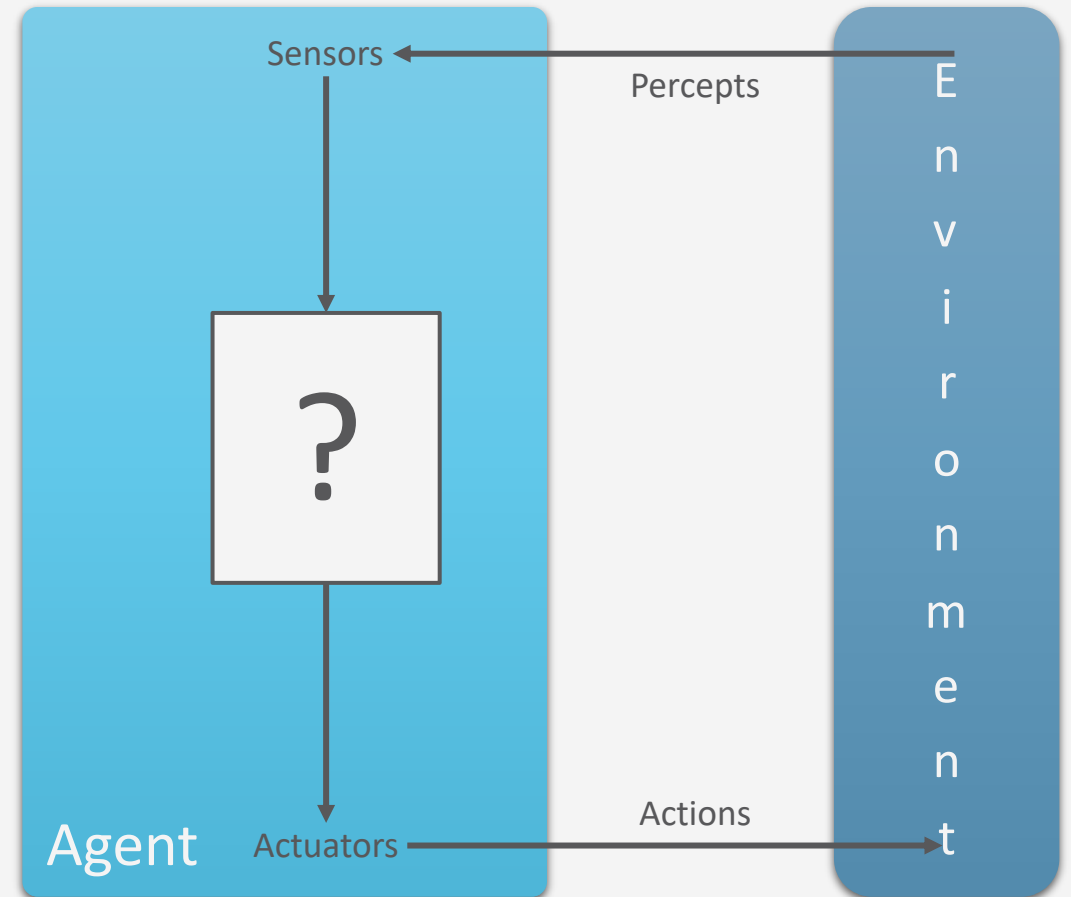- Better suited to generate agent designs that provably achieve rationality

- Rational agent approach
- Agent = something that acts
  - Operates autonomously
  - Perceives environment
  - Persists over a prolonged time period
  - Adapts to change
  - Creates and pursues goals
- Rational agent
  - Acts so as to achieve the best outcome or, when there is uncertainty, the best expected outcome
  - May include thinking rationally or acting humanly, but *more general*

Sensors

Percepts

Environment

?

Agent

Actuators

Actions

Tanya Braun

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

# Rationality

- Depends on four things:
  - *Performance measu*re, defines criterion of success
  - Agent's prior *knowledge* of environment
  - *Actions* that agent can perform
  - Agent's *percept sequence* to date
- Rational agent:
  - For each possible percept sequence, a rational agent should select an *action*
  - expected to maximize its *performance measure*,
  - given evidence provided by *percept sequence* and
  - whatever built-in *knowledge* the agent has.
- → Rational = intelligent

# Agent Structure: Simple Reflex Agent

- Actions chosen based on current percept
  - Ignores previous percepts
  - No modelling of the environemnt
- Only correct decision on action if environment fully observable
  - If partially observable, inifinite loops possible
    - (Partial) solution: Choose random action

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

# Agent Structure: Goal-based Agent

- Goal information useful
  - Description of desirable states
    - Infer from performance measure
    - Conditions for a goal state to fulfil
    - Example: vacuum cleaner
      $$\forall x \in Loc : x = clean$$

- Combine current state and goal information to choose actions that lead to goal

- Research areas:
  - Search
  - Planning



State

How the world evolves

What my actions do

Goals

Sensors

What the world is like now

What it will be like if I do action $A$

What action I should do now

Actuators

Environment

Agent

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

# Agent Structure: Utility-based Agent

- Goal-based: binary distinction between *happy* and *unhappy*
- Utility as a distribution over possible states
  - Essentially an internalisation of the performance measure
    - If internal utility function *agrees with* external performance measure:
    - Agent that chooses actions to maximize its utility will be *rational* according to the external performance measure
      - MEU principle
      - Utility function guaranteed to exist

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

Quelle: https://people.eecs.berkeley.edu/~russell/talks/2020/russell-aaai20-hntdtwwai-4x3.pptx

Quelle: https://people.eecs.berkeley.edu/~russell/talks/2020/russell-aaai20-hntdtwwai-4x3.pptx

# Where it's going – maybe

Artificial Intelligence Research

Data Science Grouo
Computer Science Department

# A Bit of the Future: 3rd Wave of AI?

- Focus: human-centred AI, hybrid approaches
  - Human-centred: Human-aware AI, explainable AI (XAI), interactive machine learning
    - Different dimensions
      - Human as a source for training
      - Human for which outputs should be comprehensible
      - Human and system working as a team
  - Hybrid: Combine data-driven and knowledge-driven approaches
    - Also known as *neuro-symbolic*
    - Use knowledge during learning to combat the problem of requiring a huge amount of data



*Human, grant me the serenity to accept the things I cannot learn; Data to learn the things I can; And wisdom to know the difference.*

God, Grant me th
Serenity
to Accept the Thin
I Cannot Chang
Courage
to Change the Thi
I Can & the
Wisdom
to Know the Differ

### Polanyi's Revenge
Kambhampati, Subbarao. "Polanyi's Revenge and AI's New Romance with Tacit Knowledge". In *Communications of the ACM*, 2021.

# Hybrid / human-centred AI

# What's intelligence got to do with it?

*Where it's been*

- Knowledge-driven AI: model-based inference, provable properties, comprehsenibility → brittle!
- Data-driven AI: learn a model from huge amounts of input-output pairs → interpretability issue!

*What it's doing*

- AI methods: search-based problem solving, logic-based inference and knowledge representation, probabilistic modelling and reasoning under uncertainty, machine learning, perception and action

*Where it's going – maybe*

- Hybrid AI: combine knowledge- and data-driven AI methods
- Human-centric AI: Do not forget the human in all of this!
  - And all the things that come with it:
    Ethics, robustness, safety, transparency, trustworthiness, …

*Thank you!*

# Appendix

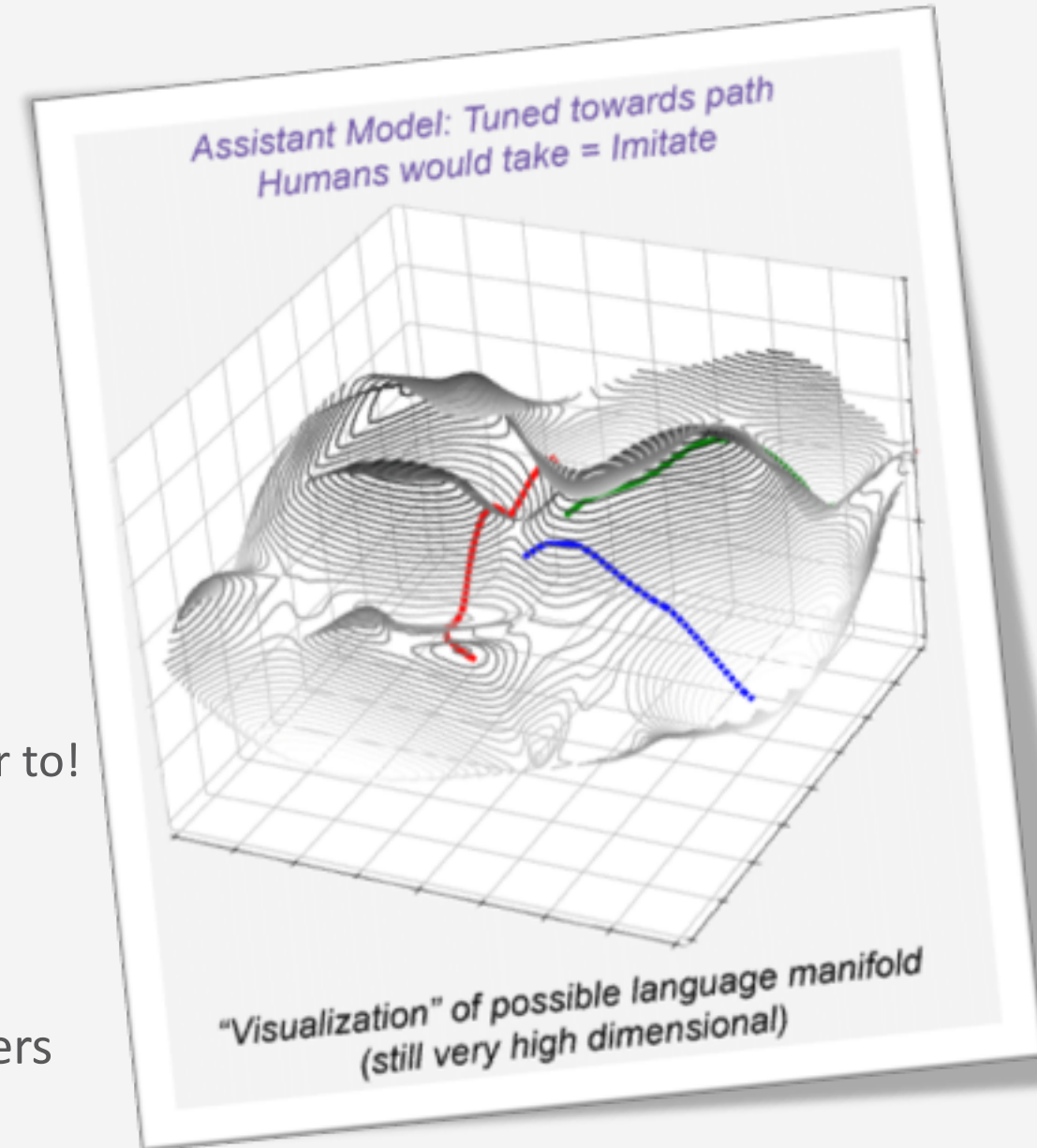Unused Slides

Data Science Group
Computer Science Department

# But what about Large Language Models (LLMs)?

- Text is a long sequence of words (including spaces, punctuations)
- $n$-gram model of language learns to predict $n$-th word given the preceding $n-1$ words
  - Probabilistically speaking, it learns
    $$\Pr(W_n | W_1, \dots, W_{n-1})$$
    - Unigram predicts each word independently
    - Bigram predicts each word given the previous word
    - 3001-gram model learns to predict the next word given the previous 3000 words
      - ChatGPT is just a 3001-gram model

- Power of an $n$-gram model depends on
  - How much text it trains on
  - How big the $n$ (context) is
  - How high-capacity the function learning $\Pr(W_n | W_1, \dots, W_{n-1})$ is
- ❖ ChatGPT trains on ~600 GB of text (Web)
  - Learns a very high capacity function that has 175 billion parameters
  - Learns $\Pr(W_n | W_1, \dots, W_{n-1})$ for all possible $n$-th words $W_n$ (Vocabulary of the language, ~50K in English)
    - Requires extreme computing facilities

Subbarao Kambhampati: "On the role of Large Language Models in Planning)", tutorial, ICAPS 2020.

# Large Language Models (LLMs)

- Different use cases
  - Generate / translate / summarise text
  - Design slides, program code
  - ✓ Relief from repetetive tasks
- Problems
  - No factual accuracy, no sources
    - Do not ask a question that you do not know the answer to!
  - Language streamlining
  - Taking over US-American values
  - Copyright issues
  - Losing capabilities such as structuring complex matters due to over-reliance on LLMs for text generation?



Assistant Model: Tuned towards path
Humans would take = Imitate

"Visualization" of possible language manifold
(still very high dimensional)

Figure taken from a talk by Malte Schilling
https://www.dropbox.com/s/nsenp948uc93l5w/schilling_2023_06_LLM_Mechanisms.pdf?dl=0
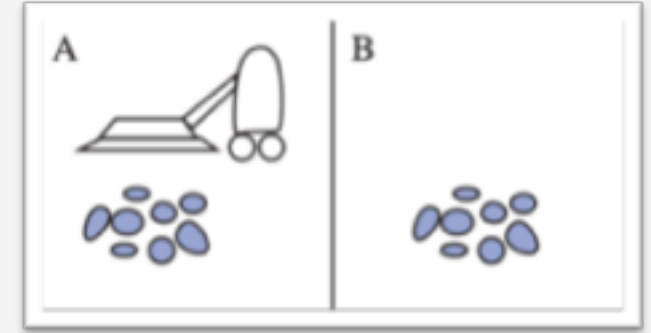
If the representations are learned, how do we ensure that they are understandable to the humans?

# Large Language Models (LLMs) – Reasoning?

- Our poor intuitions about approximate omniscience make it hard to tell whether LLMs are reasoning or retrieving
  - It is worth understanding that our intuitions about what exactly is in the 600GB of text on the web are very poor
    - If you are not surprised at someone answering a question by "googling" it, you probably shouldn't be too impressed by an LLM answering it
  - This means that we are not good at guessing whether LLMs came to an answer mostly by approximate retrieval or by first principles reasoning
- In the case of inference tasks, we may consider that an LLM was able to reach a conclusion by something akin to theorem proving from base facts
  - But then we are missing the simple fact that the linguistic knowledge on the web not only contains "facts" and "rules" but chunks of the deductive closure of these facts/rules.
    - In general, memory reduces the need to reason from first principles

Subbarao Kambhampati: "On the role of Large Language Models in Planning", tutorial, ICAPS 2020.

# Simple Example



- Vacuum cleaner
  - Two locations: squares *A*, *B*
  - Possible percepts: location; location *clean*, *dirty*
  - Available actions: *right*, *left*, *vacuum*
- Performance Measure: 1 point for each clean square in each time step over a life span of 1000 time steps
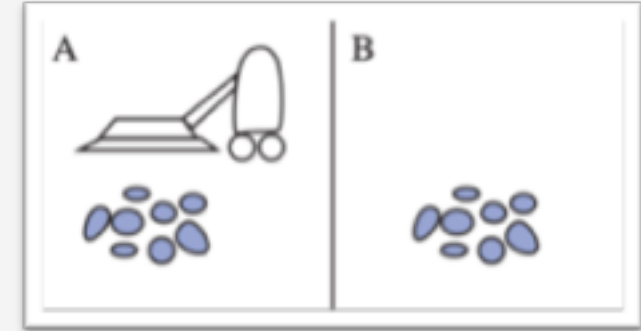
| Percept sequence | Action |
|---|---|
| [*A*, *Clean*] | *Right* |
| [*A*, *Dirty*] | *Vacuum* |
| [*B*, *Clean*] | *Left* |
| [*B*, *Dirty*] | *Vacuum* |
| [*A*, *Clean*], [*A*, *Clean*] | *Right* |
| [*A*, *Clean*], [*A*, *Dirty*] | *Vacuum* |
| ... | ... |
| [*A*, *Clean*], [*A*, *Clean*], [*A*, *Clean*] | *Right* |
| [*A*, *Clean*], [*A*, *Clean*], [*A*, *Dirty*] | *Vacuum* |
| ... | ... |

**function** Reflex-Vacuum-Agent([*location*, *status*]) **returns** an action
    **persistent**: *rules*, a set of condition-action rules

    **if** *status* = *Dirty* **then return** *Vacuum*
    **else if** *location* = *A* **then return** *Right*
    **else if** *location* = *B* **then return** *Left*

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

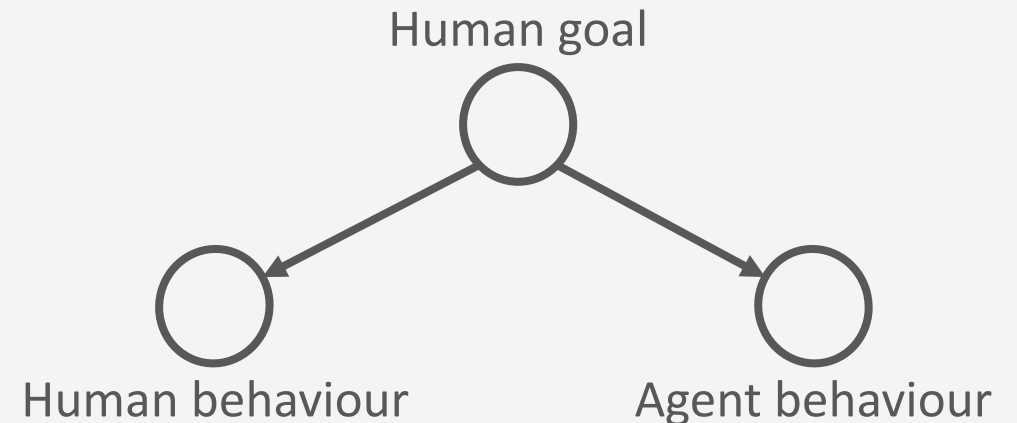# All of This Hinges on the
## *Performance Measure / Utility Function*

- Hard to determine
  - Not one fixed performance measure for all tasks and agents
- Maybe even harder to learn

Amount of dirt?

Clean locations?

| Percept sequence | Action |
|---|---|
| [*A, Clean*] | *Right* |
| [*A, Dirty*] | *Vacuum* |
| [*B, Clean*] | *Left* |
| [*B, Dirty*] | *Vacuum* |
| [*A, Clean*], [*A, Clean*] | *Right* |
| [*A, Clean*], [*A, Dirty*] | *Vacuum* |
| ... | ... |
| [*A, Clean*], [*A, Clean*], [*A, Clean*] | *Right* |
| [*A, Clean*], [*A, Clean*], [*A, Dirty*] | *Vacuum* |
| ... | ... |

Tanya Braun

S. Russel and P. Norvig: Artificial Intelligence – A Modern Approach, 2020.

# Provably Beneficial AI

If the representations are learned, how do we ensure that they are understandable to the humans?

- Idea:
  - Humans: intelligent to the extent that our actions can be expected to achieve our goals
  - ~~Maschines: intelligent to the extent that their actions can be expected to achieve their goals~~
  - Maschines are *beneficial* to the extent that their actions can be expected to achieve our goals n
  - Approach: Performance measure unknown, human as assistant
- Goal: *Provably beneficial AI*
  - See for example Stuart Russell



Human goal

Human behaviour          Agent behaviour

Presentation: https://www.youtube.com/watch?v=QPSgM13hTK8
Slides: https://people.eecs.berkeley.edu/~russell/talks/2020/russell-aaai20-hntdtwwai-4x3.pptx

# XAI & Explanations

- Standard XAI: view of explanations too simple
  - Debugging tool for "inscrutable" representations
    - "Pointing" explanations (primitive)
    - Explaining decisions will involve pointing over space-time tubes
- Explanations critical for collaboration
  - But not as a monologue from the agent
    → interaction



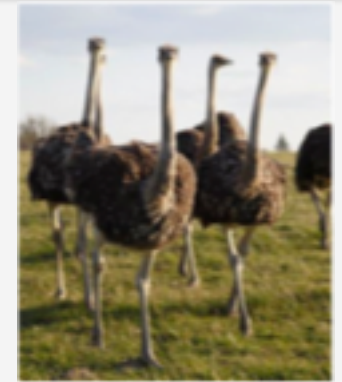Please point to the "ostrich" part

| Prediction: School bus | Difference between left and right magnified by 10 | Prediction: Ostrich |

# Human-aware Intelligent Agent

Sarath Sreedharan, Anagha Kulkarni, Subbarao Kambhampati: Explainable Human-AI Interaction: A Planning Perspective. Springer, 2022.

# Classical Planning

- Given a planning problem $(\Sigma, s_0, S_g)$, i.e., the agent's model $\mathcal{M}^R$
- Find a plan $\pi = \langle a_1, a_2, \ldots, a_n \rangle$ that transforms $s_0$ to a state $s_n \in S_g$
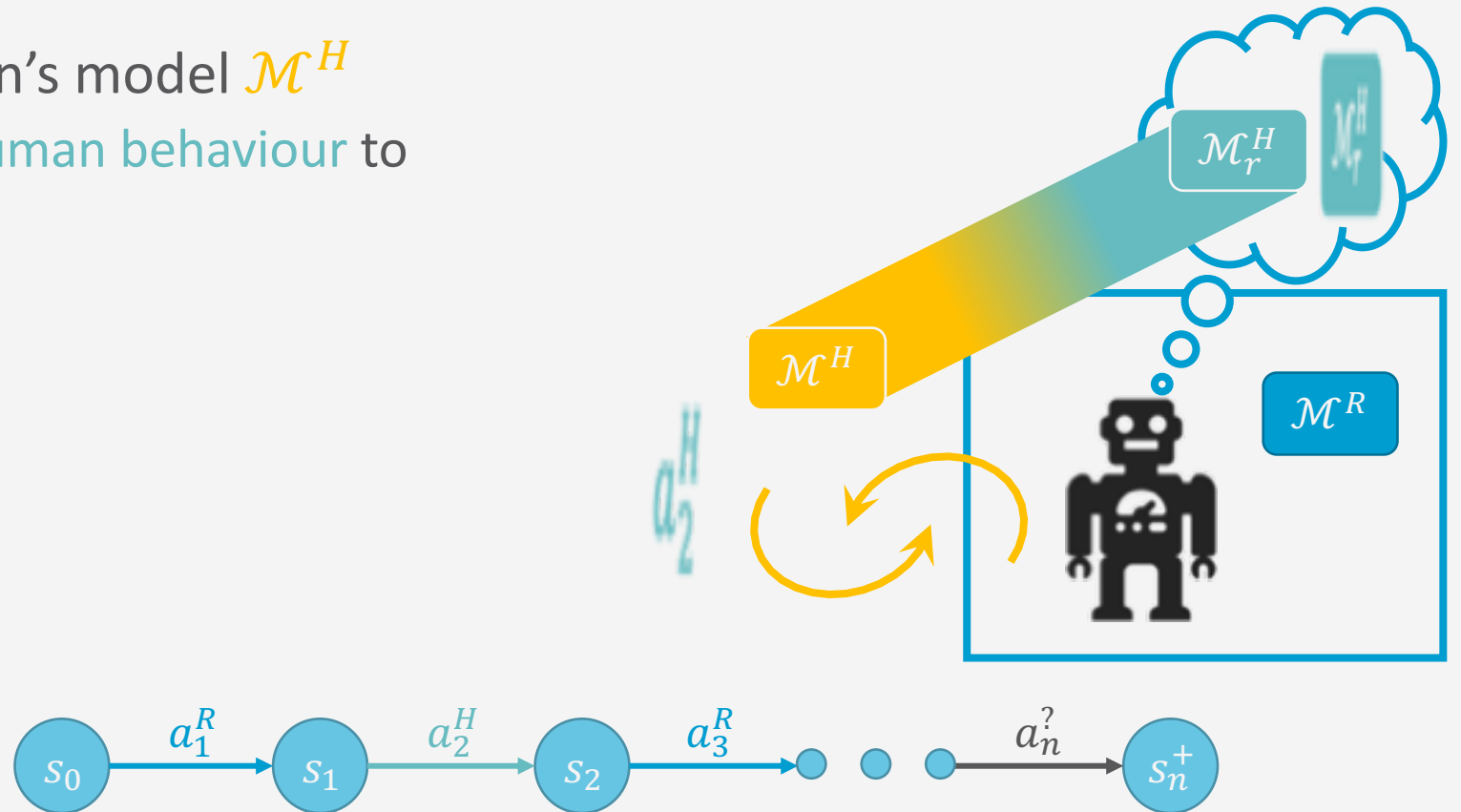
Tanya Braun

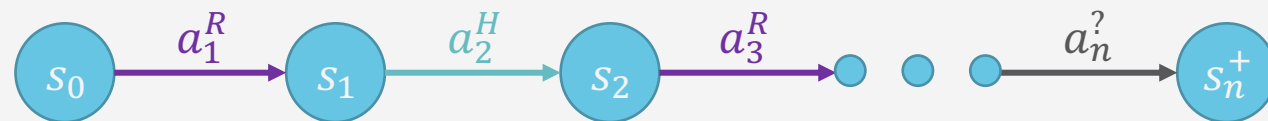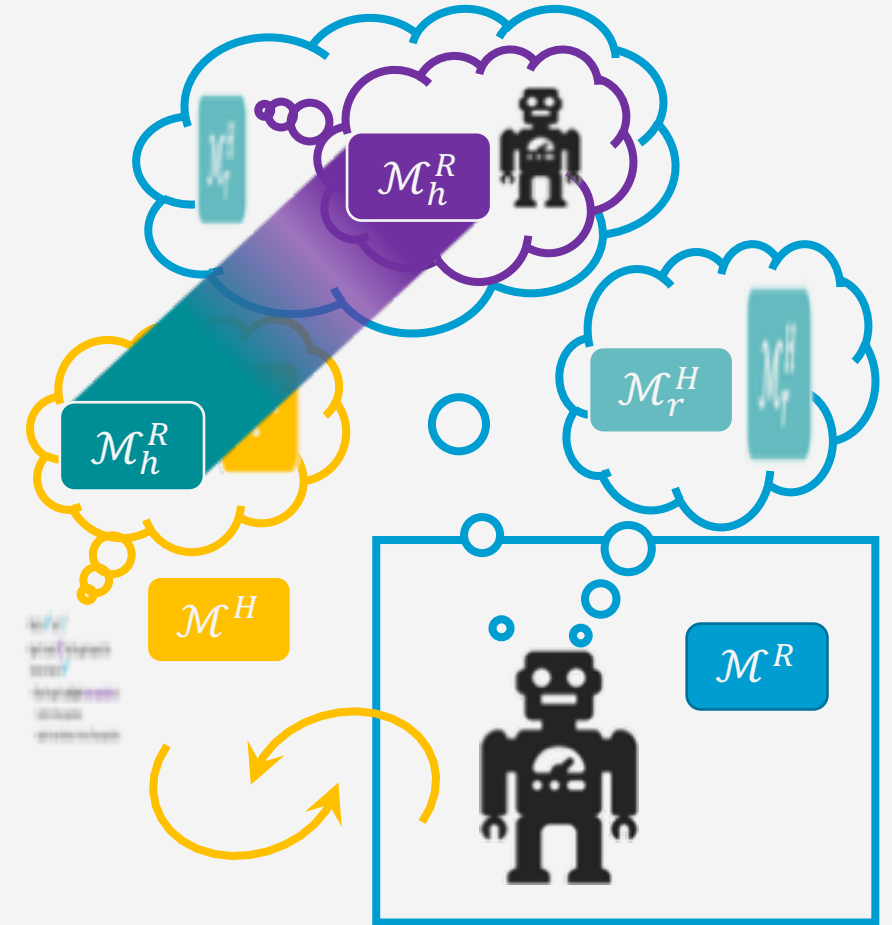Sarath Sreedharan, Anagha Kulkarni, Subbarao Kambhampati: Explainable Human-AI Interaction: A Planning Perspective. Springer, 2022.

# Collaborative Planning

- Given a planning problem $\left(\Sigma, s_0, S_g\right)$, i.e., the agent's model $\mathcal{M}^R$
- Find a joint plan $\pi = \left\langle a_1^R, a_2^H, \dots, a_n^? \right\rangle$ that transforms $s_0$ to a state $s_n^+ \in S_g$

Sarath Sreedharan, Anagha Kulkarni, Subbarao Kambhampati: Explainable Human-AI Interaction: A Planning Perspective. Springer, 2022.

# Human-aware Planning

- Next to $\mathcal{M}^R$
- Agent's model $\mathcal{M}_r^H$ of the human's model $\mathcal{M}^H$
  - Allows the agent to **anticipate** human behaviour to
    - assist
    - avoid
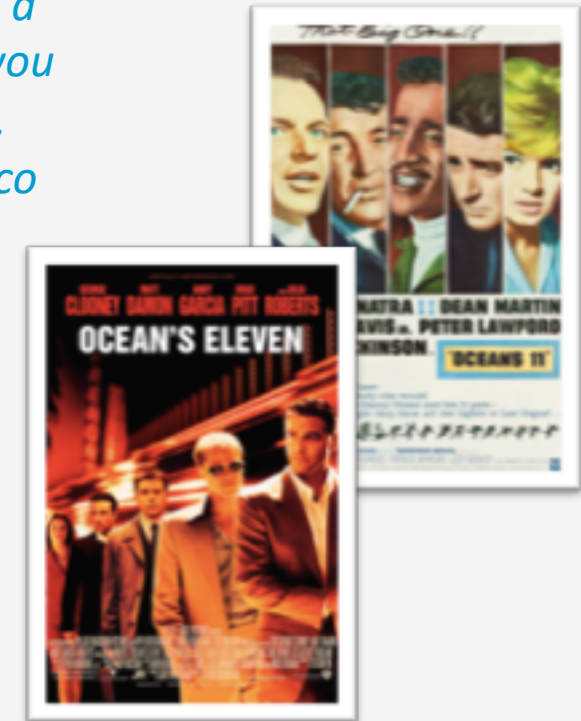    - team

# Human-aware Planning

- Next to $\mathcal{M}^R$ and $\mathcal{M}_r^H$
- Agent's model $\widetilde{\mathcal{M}}_h^R$ that the agent expects the human to have of $\mathcal{M}^R$
  - Allows the agent to **anticipate** human expectations to
    - conform to those expectations
    - explain its own behaviour in terms of those expectations

Sarath Sreedharan, Anagha Kulkarni, Subbarao Kambhampati: Explainable Human-AI Interaction: A Planning Perspective. Springer, 2022.
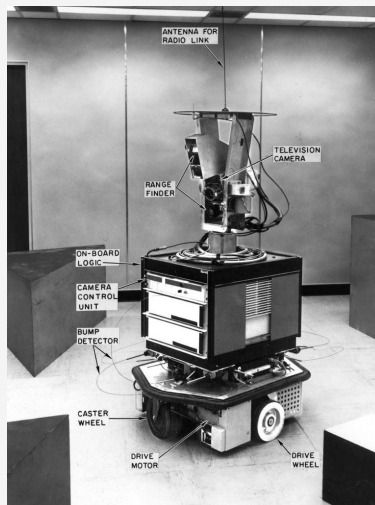
# Ethical Quandaries of Interaction

- Evolutionary, mental modelling allowed us to both cooperate or compete/sabotage each other
  - Lying is only possible because we can model others' mental states
- Human-aware AI systems with mental modelling capabilities bring additional ethical quandaries
  - E.g., automated negotiating agents that misrepresent their intentions to gain material advantage
  - Your personal assistant that tells you white lies to get you to eat healthy (or not…)

*Every tool is a weapon, if you hold it right.*
*--Ani Difranco*

Sarath Sreedharan, Anagha Kulkarni, Subbarao Kambhampati: Explainable Human-AI Interaction: A Planning Perspective. Springer, 2022.

# Ethical Quandaries of Interaction

- Humans' example closure tendencies are more pronounced for emotional/social intelligence aspects
  - No on who saw Shakey the first time thought it could shoot hoops, yet the first people interacting with Eliza assumed it was a real doctor
  - Concerns about human-aware AI "toys" such as Cozmo (e.g., Sherry Turkle)

https://thenewstack.io/remembering-shakey-first-intelligent-robot/
https://en.wikipedia.org/wiki/ELIZA https://anki.com/en-us/cozmo.html