

General Principles of Human and Machine Learning



Lecture 6: Social learning

Quiz results

Average grade: **73%**

We took out the VTE question this time, but beware tricky questions for the next quizzes!

Clarifications

Quiz content:

Quiz questions don't overlap! However, material from the week of the pop-quiz isn't eligible for that week's pop quiz, so it may show up in the next!

Oct 15: Introduction (slides)	Oct 16 (slides)	Alex	Spicer & Sanborn (2019). What does the mind learn?
Oct 22: Origins of biological and artificial learning (slides)	Oct 23 (slides)	Turan	[1] Behaviorism [2] What is a perceptron? (Blog post)
Oct 29: Symbolic AI and Cognitive maps (slides)	Oct 30 (Quiz #1)	Alex	[1] Garnelo & Shanahan (2019) [2] Boorman et al., 2021
Nov 5: Introduction to RL (slides)	Nov 6 (slides)	Turan	Sutton & Barton (Ch. 1 & 2)
Nov 12: Advances in RL (slides)	Nov 13	Turan	Neftci & Averbeck (2019)

Clarifications

Quiz content:

Quiz questions don't overlap! However, material from the week of the pop-quiz isn't eligible for that week's pop quiz, so it may show up in the next!

What is VTE, and do I need to know the abbreviation?

Vicarious trial and error (VTE): hesitating, looking-back-and-forth behavior observed in rats when confronted with a choice

and yes! Mind the highlights when reviewing (blue boxes are important)

We have adjusted the quiz total (20->18) to account for this maybe not having been clear this time.

Clarifications

Quiz content:

Quiz questions don't overlap! However, material from the week of the pop-quiz isn't eligible for that week's pop quiz, so it may show up in the next!

What is VTE, and do I need to know the abbreviation?

Vicarious trial and error (VTE): hesitating, looking-back-and-forth behavior observed in rats when confronted with a choice

and yes! Mind the highlights when reviewing (blue boxes are important)

We have adjusted the quiz total (20->18) to account for this maybe not having been clear this time.

Also, **please don't cheat!** We found 5 people who had exactly the same wrong answers on several questions as one another. We've made a note of who they are and **if it happens again, we will take disciplinary action with the university and assign a grade of 0 to both this quiz and the next one where it occurs.** These quizzes are far too low stakes to risk cheating.

Clarifications

Symbolic AI

- **Physical Symbol System hypothesis:**

“A physical symbol system has the necessary and sufficient means for general intelligent action - Allen Newell and Herbert Simon (1976)”

- **Symbols** can represent things in the world
 - e.g., (Apple), (ChatGPT), (Charley), etc...
- **Relations** can be i) predicates that describes a symbol or ii) verbs describing how symbols interact with other symbols
 - i) red(Apple), unreliable(ChatGPT), instructor(Charley)
 - ii) eat(Charley, Apple), generatePicture(ChatGPT, Apple)

- By populating a **knowledge base** with symbols and relations, we can use a program to find new propositions (*inference*)
 - General Problem Solver (Simon, Shaw, & Newell, 1957)
 - Expert systems: popularized in the 1980s as the future of AI

Clarifications

Symbolic AI

- **Physical Symbol System hypothesis:**

“A physical symbol system has the necessary and sufficient means for general intelligent action - Allen Newell and Herbert Simon (1976)”

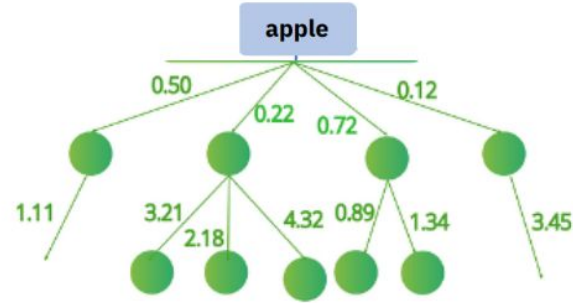
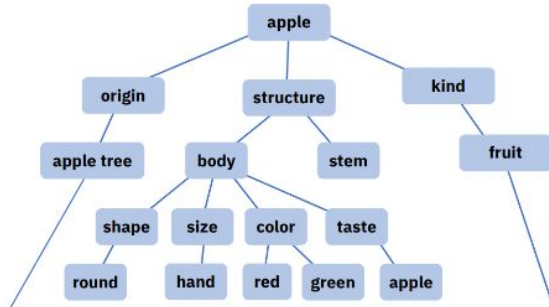
- **Symbols** can represent things in the world
 - e.g., (Apple), (ChatGPT), (Charley), etc...
- **Relations** can be i) predicates that describes a symbol or ii) verbs describing how symbols interact with other symbols
 - i) red(Apple), unreliable(ChatGPT), instructor(Charley)
 - ii) eat(Charley, Apple), generatePicture(ChatGPT, Apple)
- By populating a **knowledge base** with symbols and relations, we can use a program to find new propositions (*inference*)
 - General Problem Solver (Simon, Shaw, & Newell, 1957)
 - Expert systems: popularized in the 1980s as the future of AI

Why “physical symbol system”, and not just “symbolic system”?

Very good question. After consulting with a philosophy professor (Hongyu Wong), I have an answer. Newell and Simon wanted to demonstrate that a **physical computer or robot** could demonstrate intelligence. So the “symbol system” they are describing **physically exists in the world**. This is in contrast to other philosophical ideas at the time about a “mental” symbol system, that might exist in a purely abstract Cartesian sense

Clarifications

Symbolic vs. sub-symbolic AI



Symbolic AI

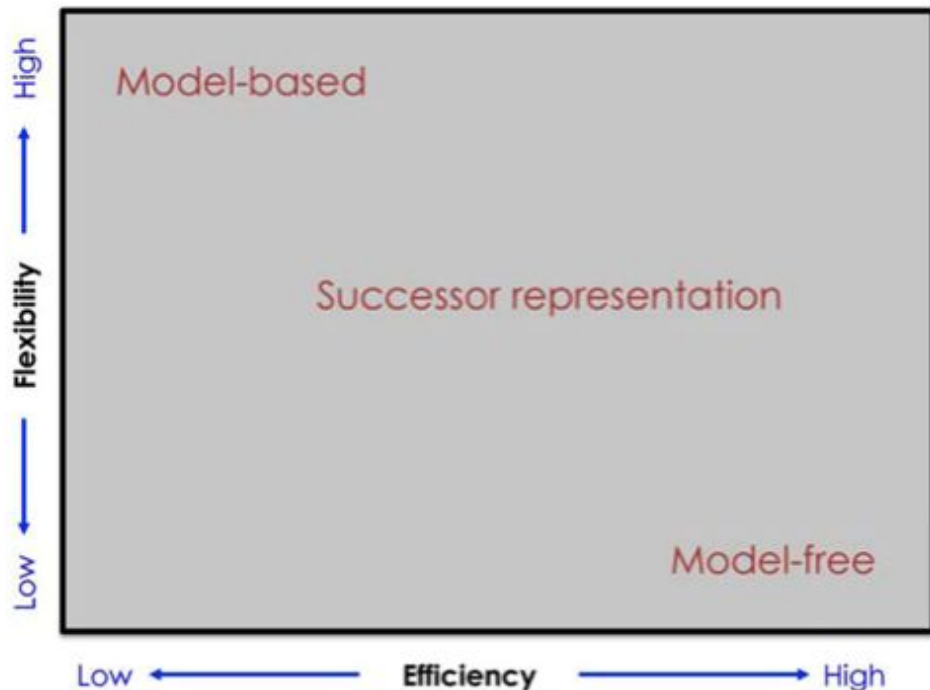
- Symbols, rules, and structured representations
- “**Language of thought**” (LoT) hypothesis (Fodor, 1975): concepts/knowledge represented by a language-like system
- Compositionality: symbols and rules can be combined to produce new representations
- Extracting symbolic representations and search over compositional hypothesis spaces is difficult

Sub-symbolic AI

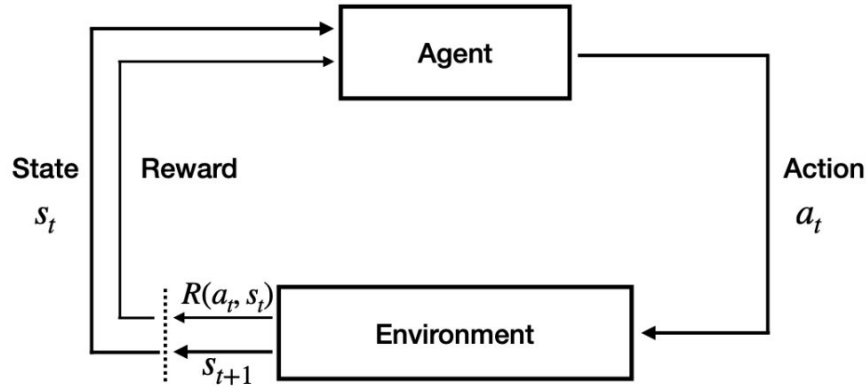
- Representations encoded through connection weights
- **No explicit representation** of concepts or knowledge, but **distributed throughout the network**
- Efficiency: knowledge can be implicitly learned by capturing statistical patterns
- Interpretation of representations and behavior is difficult

Balancing flexibility and efficiency

- Model-free methods are more **computationally efficient**
 - But lack flexibility to changes in the environment
- Model-based methods are highly flexible (local changes in environment lead to local changes in model)
 - But **computationally costly when it comes to performing simulations**
- Is there nothing in between?



Why is social learning even interesting?



But when action spaces are vast, or risky...



...social learning lets us avoid costly trial and error!

restaurant

Results

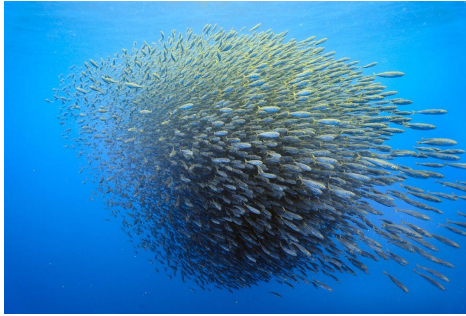
Chazz Palminteri Italian Restaurant
4.5 ★★★★★ (958)
Italian · 30 W 46th St
Traditional Italian dining and drinking
Open · Closes 10 PM
Dine-in · Takeout · No-contact delivery
[RESERVE A TABLE](#) [ORDER ONLINE](#)

VIP RESTAURANT LLC BARSHAY'S
3.9 ★★★★★ (1,814)
Diner · 175 Sip Ave
Classic American diner with modest decor
Open · Closes 8 PM
Dine-in · Curbside pickup · Delivery
[ORDER ONLINE](#) [CHECK WAIT TIME](#)

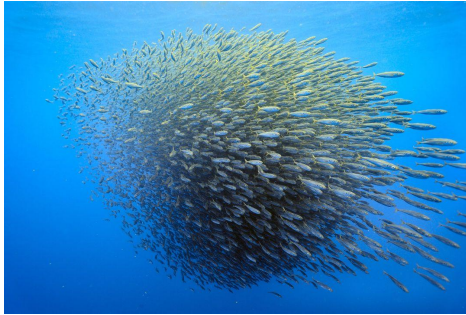
The Modern
4.6 ★★★★★ (2,125)
New American · 9 W 53rd St
Fine dining at the Museum of Modern Art
Open · Closes 9 PM
Dine-in · No takeout · No delivery
[RESERVE A TABLE](#)



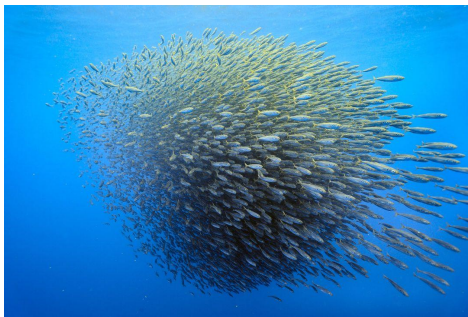
Social learning is ubiquitous in the animal kingdom



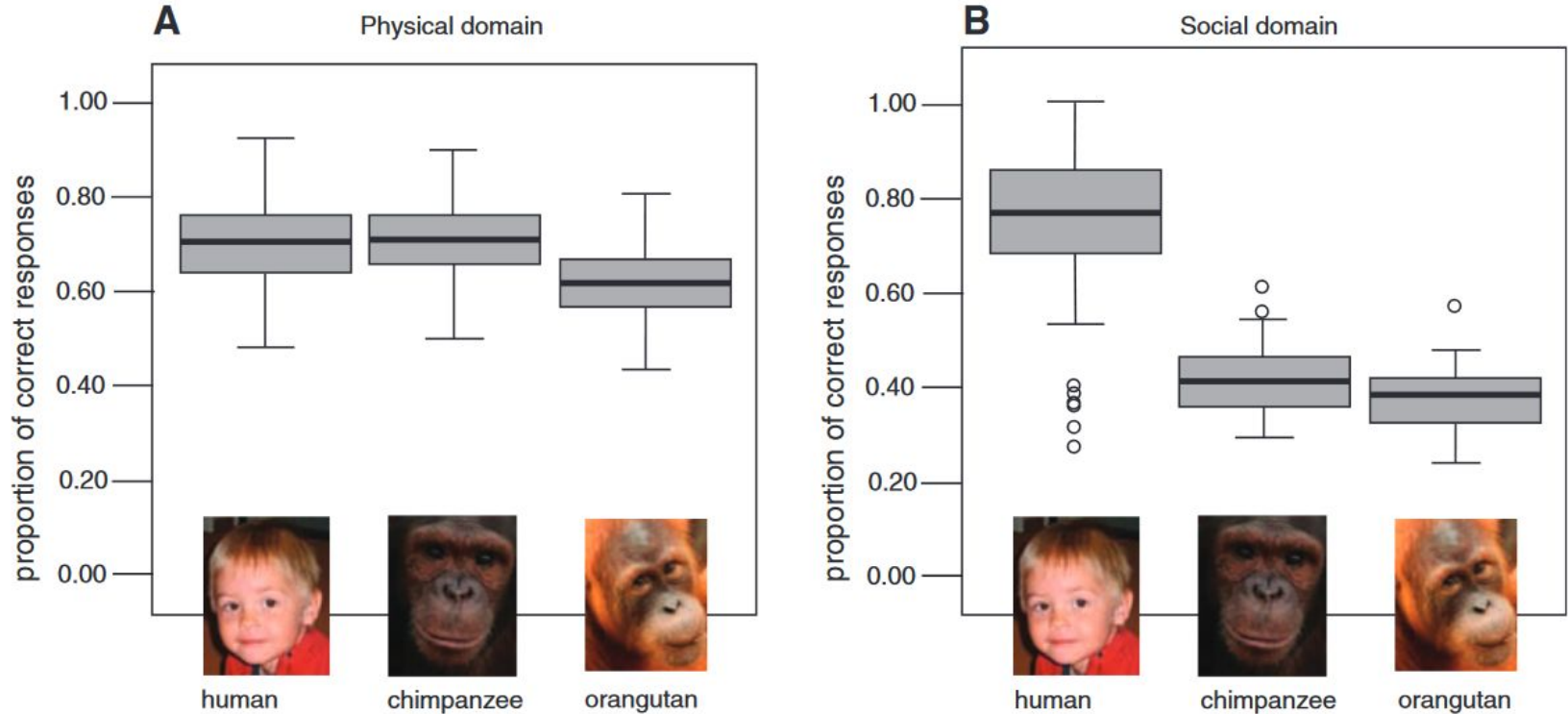
Social learning is ubiquitous in the animal kingdom



But humans do it best!

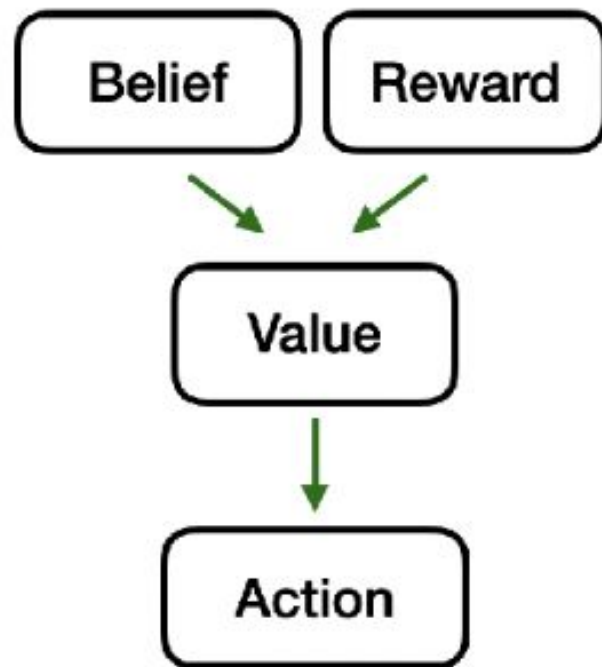


But humans do it best!



Herrmann et al. (2007), *Science*

Decision-making hierarchy



Wu et al., (2022)

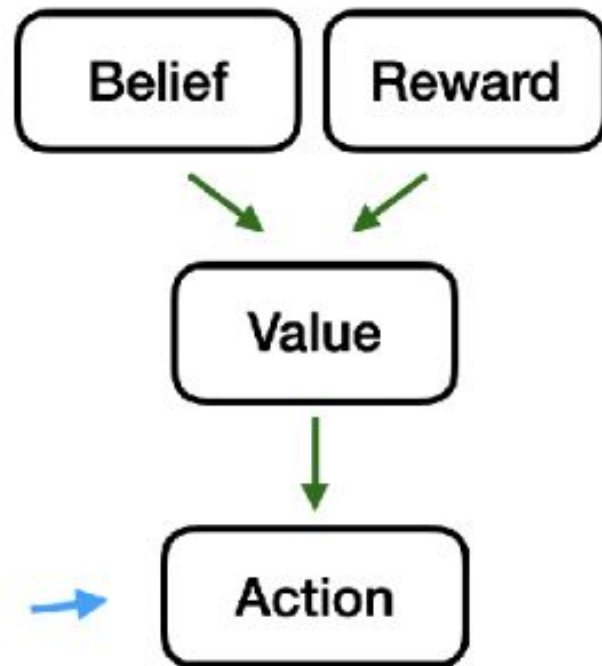
Levels of social learning



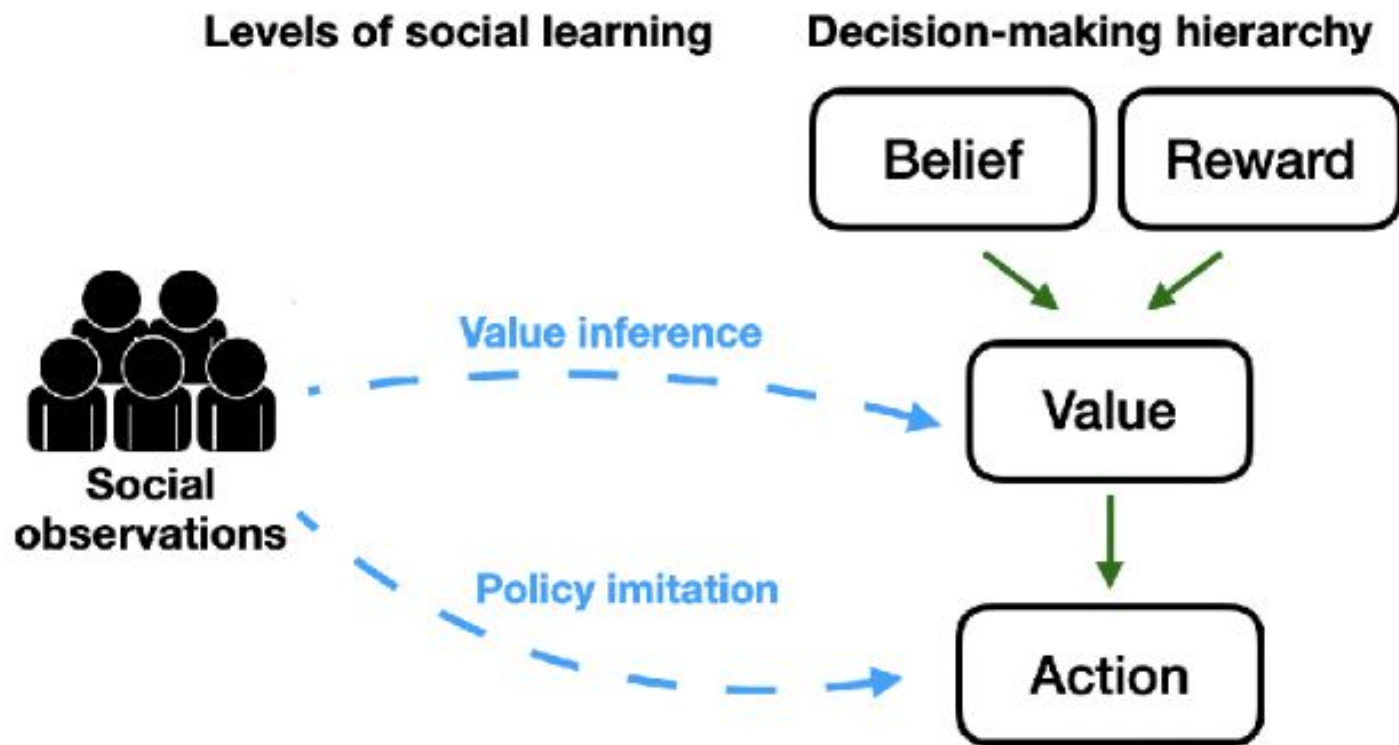
Social observations

Policy imitation

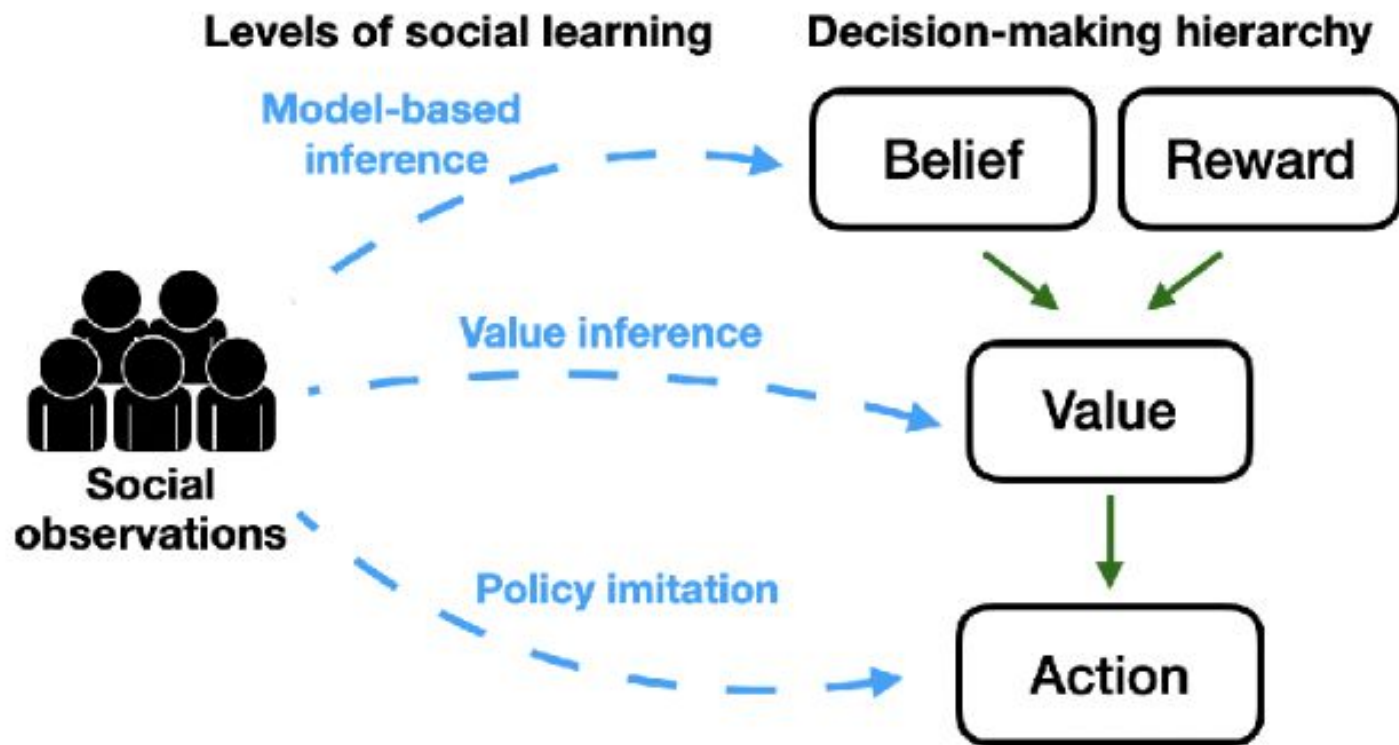
Decision-making hierarchy



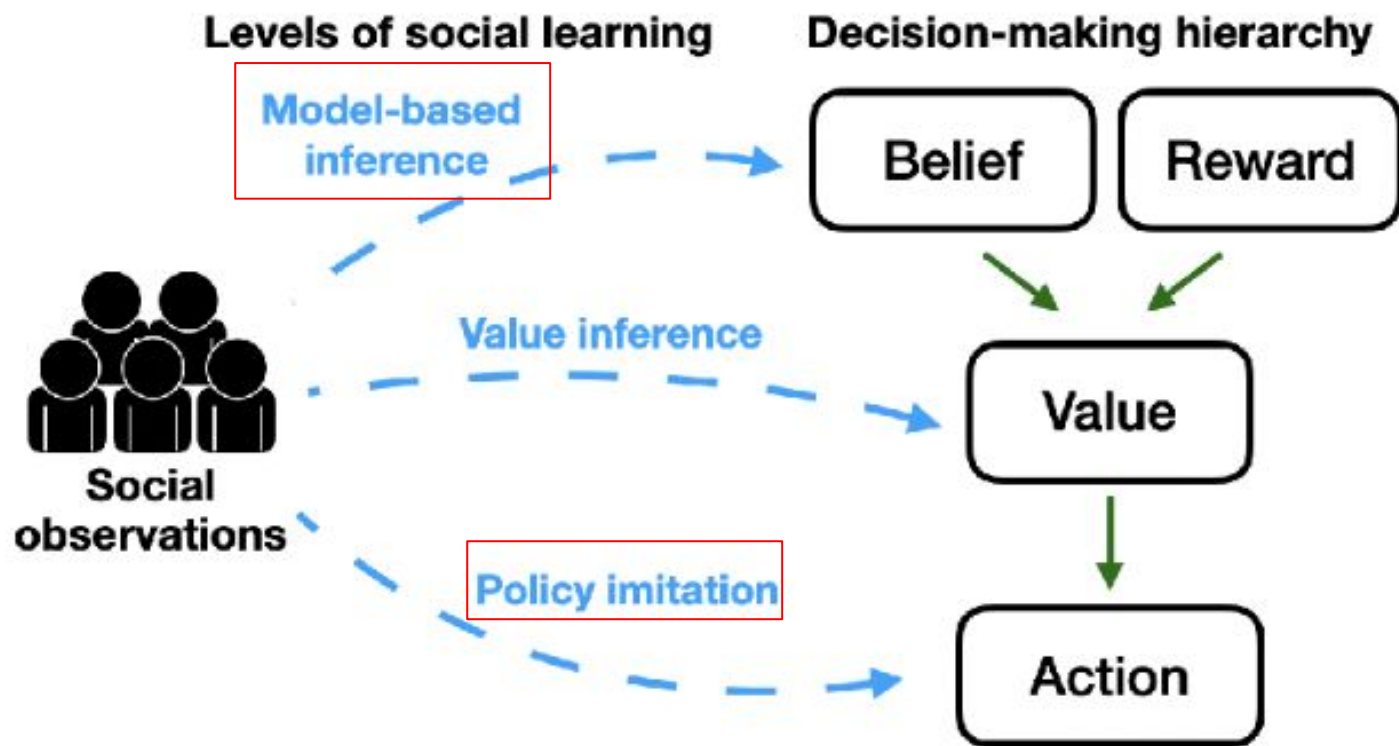
Wu et al., (2022)



Wu et al., (2022)



Wu et al., (2022)

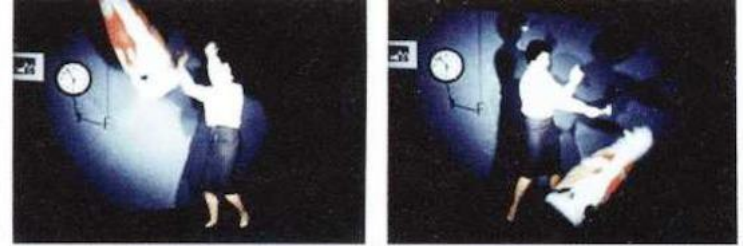


Wu et al., (2022)

Part I: Imitation

The Bobo doll experiment (Bandura, 1961)

- Children observed an adult model either attack the Bobo doll, or watched them play with other toys
- After 10 minutes of observation, the children were put in the same room, and their behaviour was observed through a one-way mirror



The Bobo doll experiment (Bandura, 1961)

- Children observed an adult model either attack the Bobo doll, or watched them play with other toys
- After 10 minutes of observation, the children were put in the same room, and their behaviour was observed through a one-way mirror
- **Children with an aggressive model displayed significantly more aggressive behaviour**



The Bobo doll experiment (Bandura, 1961)

- Children observed an adult model either attack the Bobo doll, or watched them play with other toys
- After 10 minutes of observation, the children were put in the same room, and their behaviour was observed through a one-way mirror
- **Children with an aggressive model displayed significantly more aggressive behaviour**

→ Learning without reinforcement, just via observation!



Bandura's Social Learning Theory

- Later experiments showed **vicarious reinforcement learning**
 - Aggressive model is rewarded/receives no feedback → increased aggression
 - Aggressive model is punished → significantly less aggression

Bandura's Social Learning Theory

- Later experiments showed **vicarious reinforcement learning**
 - Aggressive model is rewarded/receives no feedback → increased aggression
 - Aggressive model is punished → significantly less aggression
- Results were formalized into **Social Learning Theory**:
 - Learning isn't purely behavioural (a departure from Skinner)
 - Instead, learning can happen via observation of actions, or of actions and their consequences
 - Observational learning can occur without an observable change in behaviour
 - Reinforcement learning isn't all there is to learning
 - The learner is more than a passive recipient of information

Bandura's Social Learning Theory

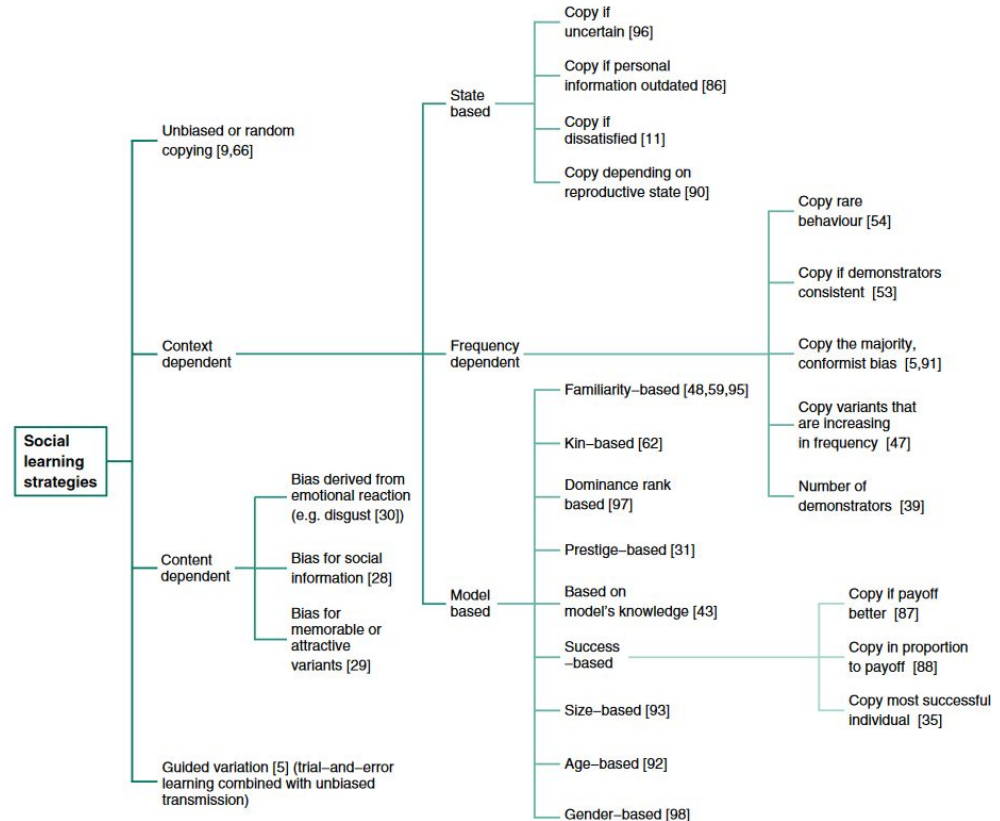
- Later experiments showed **vicarious reinforcement learning**
 - Aggressive model is rewarded/receives no feedback → increased aggression
 - Aggressive model is punished → significantly less aggression
- Results were formalized into **Social Learning Theory**:
 - Learning isn't purely behavioural (a departure from Skinner)
 - Instead, learning can happen via observation of actions, or of actions and their consequences
 - Observational learning can occur without an observable change in behaviour
 - Reinforcement learning isn't all there is to learning
 - The learner is more than a passive recipient of information
- Bandura also identified traits in **models that we would preferentially imitate**
 - Similarity
 - Status
 - Competence
 - Likeability

Social learning strategies

- Restrict settings in which animals tend to learn socially
- Categorized into types of strategies
 - **who strategies** (match Bandura's criteria pretty well: copy the expert, copy the successful, also copy the majority, ...)
 - **what strategies** (copy emotionally evocative content, copy information relevant to survival, ...)
 - **when strategies** (copy when uncertain, copy when individual learning is costly, ...)

Social learning strategies

- Restrict settings in which animals tend to learn socially
- Categorized into types of strategies
 - **who strategies** (match Bandura's criteria pretty well: copy the expert, copy the successful, also copy the majority, ...)
 - **what strategies** (copy emotionally evocative content, copy information relevant to survival, ...)
 - **when strategies** (copy when uncertain, copy when individual learning is costly, ...)



Quick-note: model-based vs. model-based

- **Model-based RL** has a **model of the environment**, which it can learn and use to plan
- Biologists/Psychologists don't care about that – **you learn from a model**, so rules that specify whom you want to copy (who-strategies) are model-based (dependent on the traits of the model you're learning from)
- This also illustrates the absolute joys of working interdisciplinarily

Who-strategies

At which restaurant would you rather eat?



Who-strategies

At which restaurant would you rather eat?



Copy the majority

Who-strategies

Who would you rather ask for medical advice?



Who-strategies

Who would you rather ask for medical advice?



Copy the expert



What strategies

What information would you rather copy?



What strategies

What information would you rather copy?



Copy emotionally evocative content

What strategies

What information would you rather copy?



What strategies

What information would you rather copy?



Copy information relevant for survival

When strategies
When would you rather imitate someone?



When strategies
When would you rather imitate someone?



Copy when individual learning is costly

When strategies

When would you rather imitate someone?



MENU

STARTER

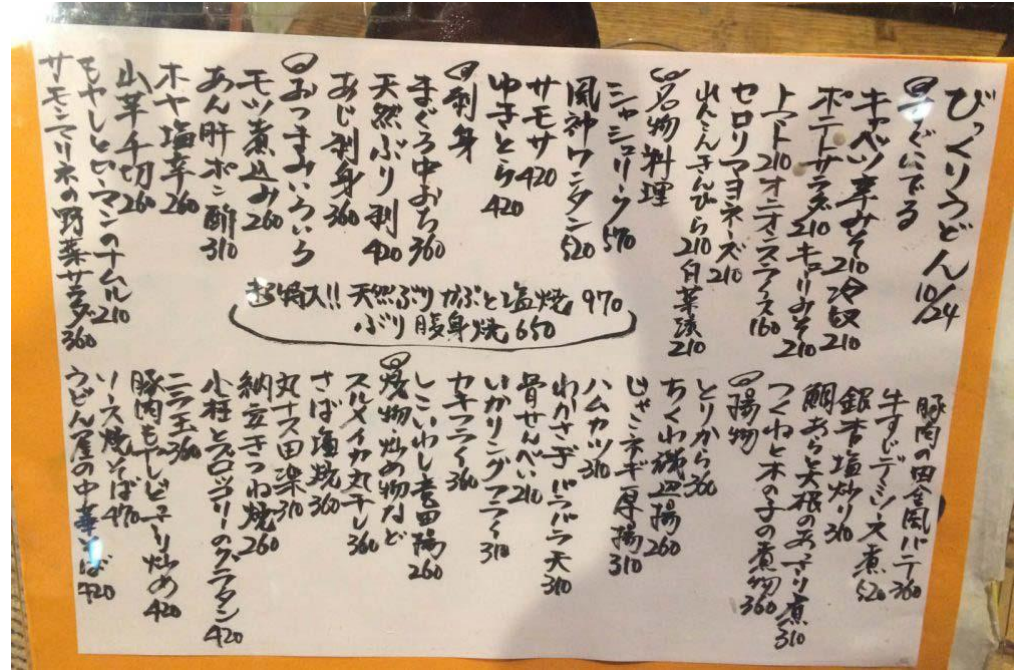
Filo wrapped Brie parcel served with salad and a reserva jus	8.95
Roast octopus, crispy calamari and a carrot purée	12.50
Crispy chilli beef salad with pickled cucumber	8.95
Garlic prawns served with crusty bread	9.95
Homemade soup of the day with fresh bread	5.50
Black pudding croquette with whisky mayo	7.95

FOR THE TABLE

A selection of bread with olive tapenade	5.50
--	------

MAIN COURSE

Haggis stuffed chicken with roasted carrots, cabbage, creamy mash and a red wine jus	16.50
Pork & black pudding wellington with sweet potato fondant and a red wine jus	22.50
Slow cooked beef cheeks, burnt onion purée and a red wine jus	18.95
Sweet potato pithivier with brandy cream and seasonal vegetables	16.50
Oven roasted Tilapia fish with a caper & lemon butter sauce served with seasonal vegetables	17.95
Sirloin steak with hand cut chips, onion rings and seasonal vegetables. Choose from red wine, peppercorn or blue cheese sauce	22.50



When strategies

When would you rather imitate someone?



MENU

STARTER

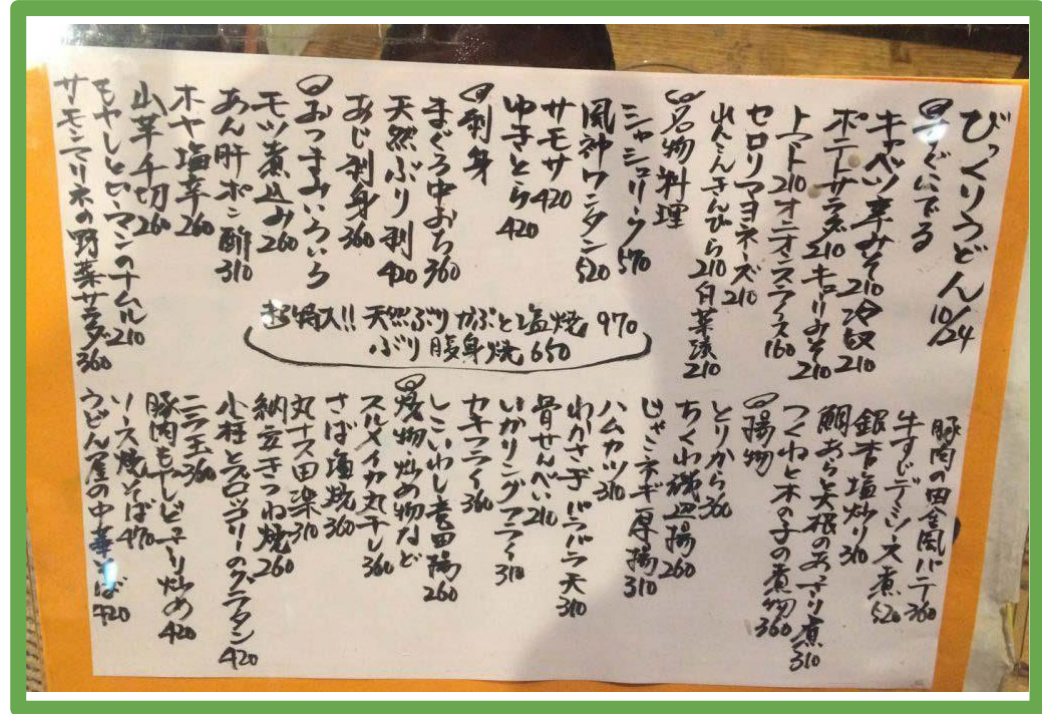
Filo wrapped Brie parcel served with salad and a reserva jus	8.95
Roast octopus, crispy calamari and a carrot purée	12.50
Crispy chilli beef salad with pickled cucumber	8.95
Garlic prawns served with crusty bread	9.95
Homemade soup of the day with fresh bread	5.50
Black pudding croquette with whisky mayo	7.95

FOR THE TABLE

A selection of bread with olive tapenade	5.50
--	------

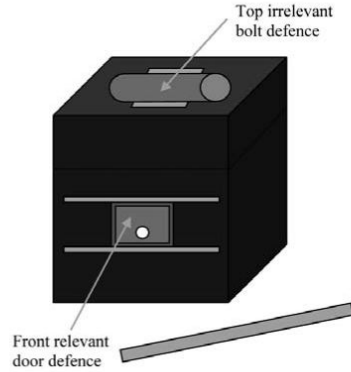
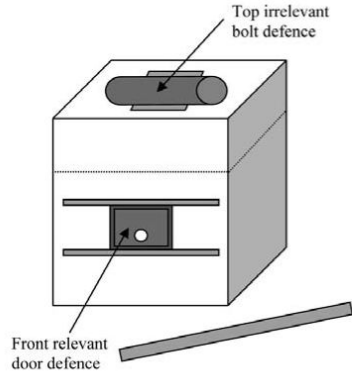
MAIN COURSE

Haggis stuffed chicken with roasted carrots, cabbage, creamy mash and a red wine jus	16.50
Pork & black pudding wellington with sweet potato fondant and a red wine jus	22.50
Slow cooked beef cheeks, burnt onion purée and a red wine jus	18.95
Sweet potato pithivier with brandy cream and seasonal vegetables	16.50
Oven roasted Tilapia fish with a caper & lemon butter sauce served with seasonal vegetables	17.95
Sirloin steak with hand cut chips, onion rings and seasonal vegetables. Choose from red wine, peppercorn or blue cheese sauce	22.50

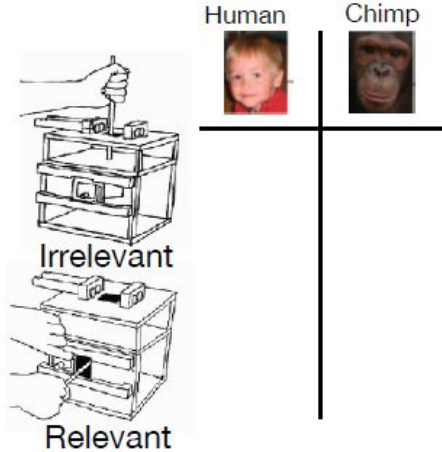
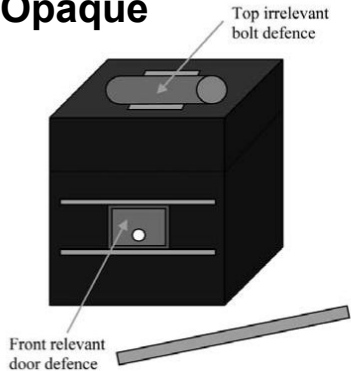


Copy when uncertain

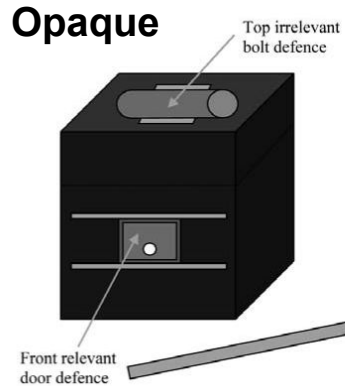
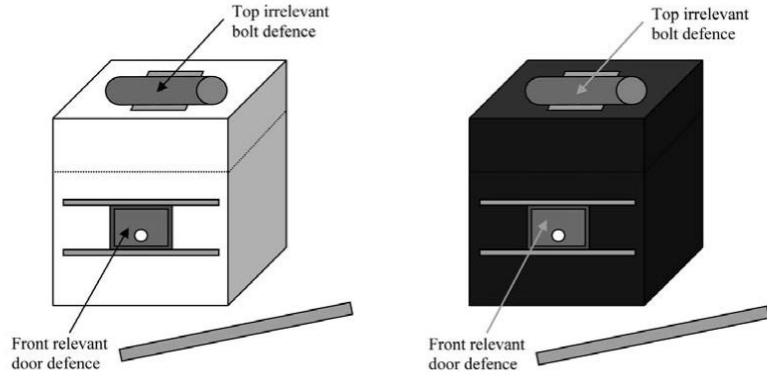
So, when do humans imitate?





Opaque



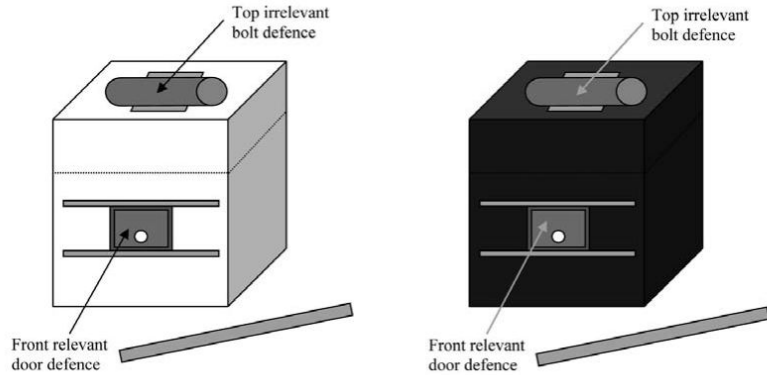
So, when do humans imitate?



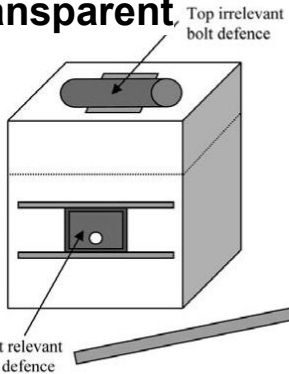
	Human	Chimp
 Irrelevant	✓	✓
 Relevant	✓	✓

So, when do humans imitate?

All the time! (Even when it isn't strictly necessary)

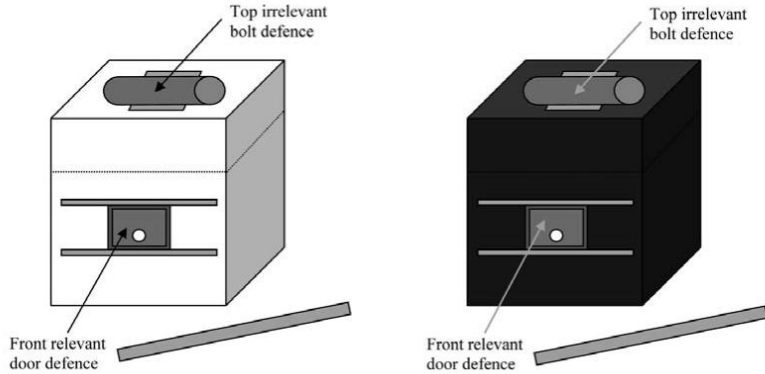


Transparent



	Human	Chimp
 Irrelevant	✓	✗
 Relevant	✓	✓

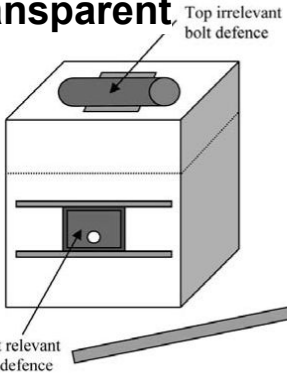
So, when do humans imitate?





All the time! (Even when it isn't strictly necessary)

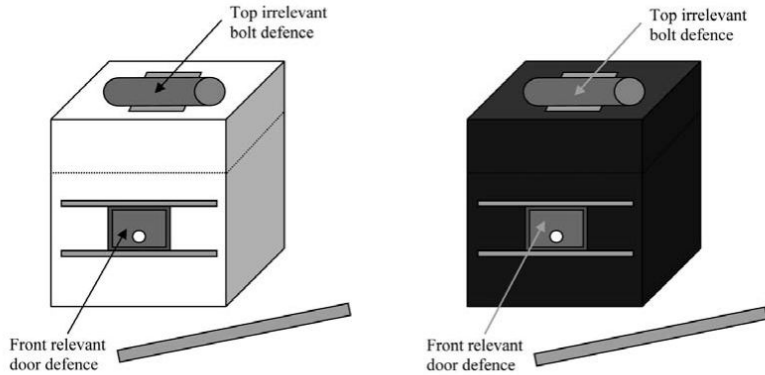
This is known as **high-fidelity imitation**, and appears to be unique to humans.

Transparent



	Human	Chimp
 Irrelevant	✓	✗
 Relevant	✓	✓

So, when do humans imitate?



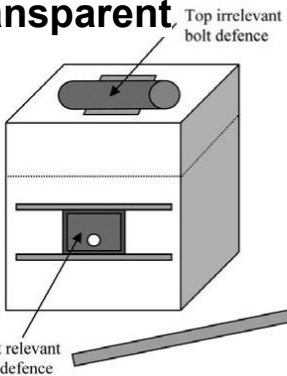
All the time! (Even when it isn't strictly necessary)

This is known as **high-fidelity imitation**, and appears to be unique to humans.

High-fidelity imitation could have its basis in the human tendency to **teach** (Csibra & Gergely, 2009) – if someone shows me something they have most likely selected the useful steps for me

This, in turn, may have allowed for our complex cumulative culture.

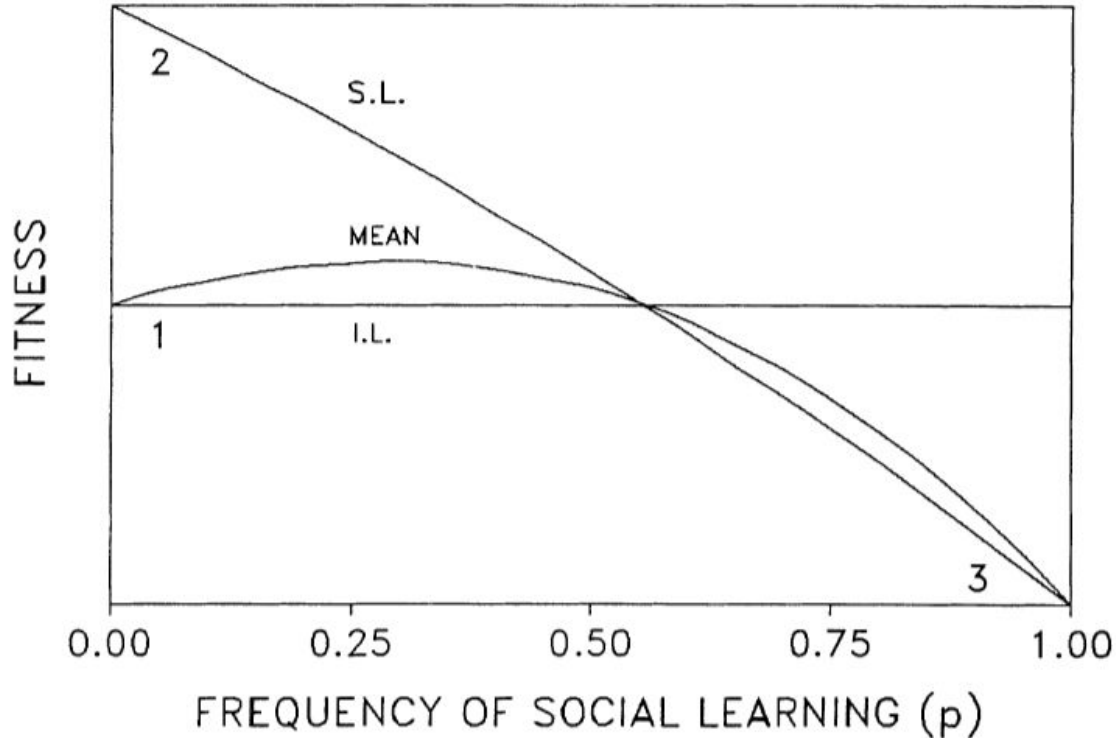
Transparent



	Human	Chimp
Irrelevant	✓	✗
Relevant	✓	✓

The table compares human and chimpanzee behavior in two conditions: 'Irrelevant' and 'Relevant'. In the 'Irrelevant' condition, humans use the tool (green checkmark) while chimpanzees do not (red X). In the 'Relevant' condition, both humans and chimpanzees use the tool (green checkmarks).

Roger's paradox – a limit on social learning?

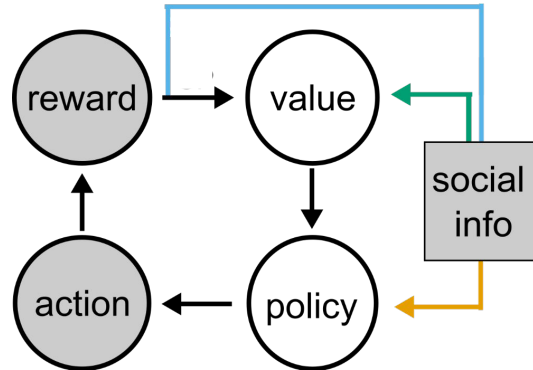


Social learning requires competent people to learn from – if everyone is copying, we get stuck.

This means that social learning has **frequency-dependent fitness** – how adaptive it is depends on how many others are also using it

Modelling imitation

- To avoid Roger's paradox, lots of recent modelling work has focused on trading off social influences and individual learning
- Models in this space differ in how exactly social information is incorporated into individual learning circuits – more on that later!



Summary – imitation

Bandura's **Social Learning Theory** was a part of the cognitive revolution, overturning behaviourism as the leading principle of psychology; the learner is not just passive, and humans learn from others even without reinforcement (although we can also learn vicariously)

Summary – imitation

Bandura's **Social Learning Theory** was a part of the cognitive revolution, overturning behaviourism as the leading principle of psychology; the learner is not just passive, and humans learn from others even without reinforcement (although we can also learn vicariously)

Social learning strategies formalize settings in which social learning often occurs in the animal kingdom.

Summary – imitation

Bandura's **Social Learning Theory** was a part of the cognitive revolution, overturning behaviourism as the leading principle of psychology; the learner is not just passive, and humans learn from others even without reinforcement (although we can also learn vicariously)

Social learning strategies formalize settings in which social learning often occurs in the animal kingdom.

Humans **overimitate**, which may have played a role in our cultural evolution

Summary – imitation

Bandura's **Social Learning Theory** was a part of the cognitive revolution, overturning behaviourism as the leading principle of psychology; the learner is not just passive, and humans learn from others even without reinforcement (although we can also learn vicariously)

Social learning strategies formalize settings in which social learning often occurs in the animal kingdom.

Humans **overimitate**, which may have played a role in our cultural evolution

Social learning has **frequency-dependent fitness**: the more people adopt social learning as their strategy, the less effective it becomes. This is known as **Roger's paradox**. To avoid this issue, recent computational models generally assume a mix of individual and social learning. They differ in the stage at which social information is integrated

Part II - Theory of Mind

Speaking of what non-human primates can do....

Does the chimpanzee have a theory of mind?

David Premack

*Department of Psychology,
University of Pennsylvania,
Philadelphia, Penna. 19104*

Guy Woodruff

*University of Pennsylvania Primate Facility,
Honey Brook, Penna. 19344*

Abstract: An individual has a theory of mind if he imputes mental states to himself and others. A system of inferences of this kind is

Speaking of what non-human primates can do....

Does the chimpanzee have a theory of mind?

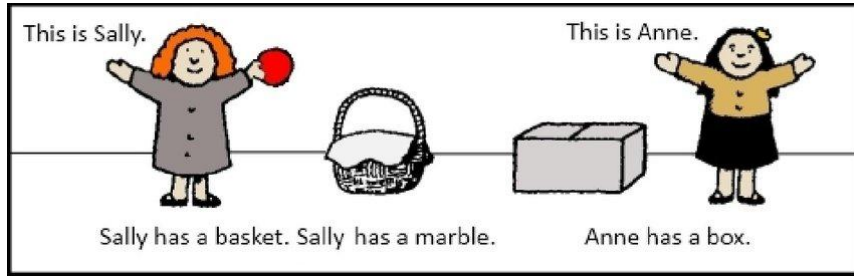
David Premack

Department of Psychology,

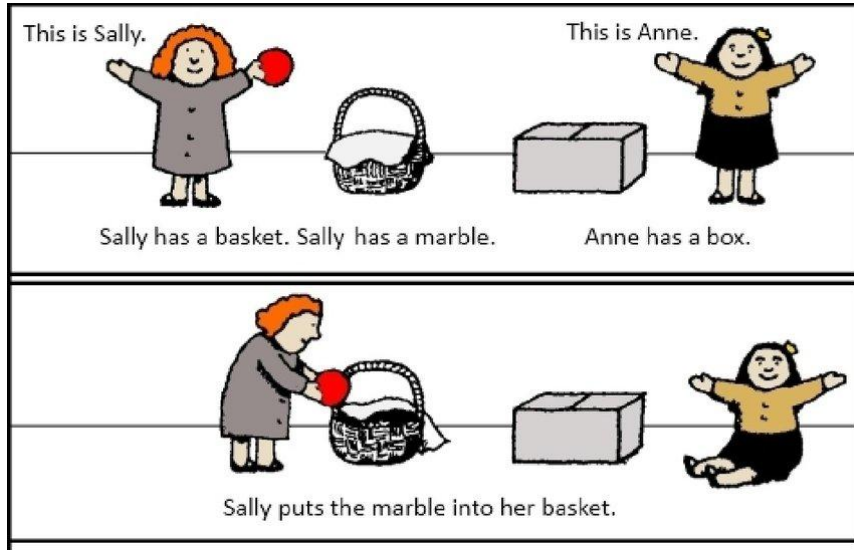
Does the chimpanzee have a theory of mind? 30 years later

Josep Call and Michael Tomasello

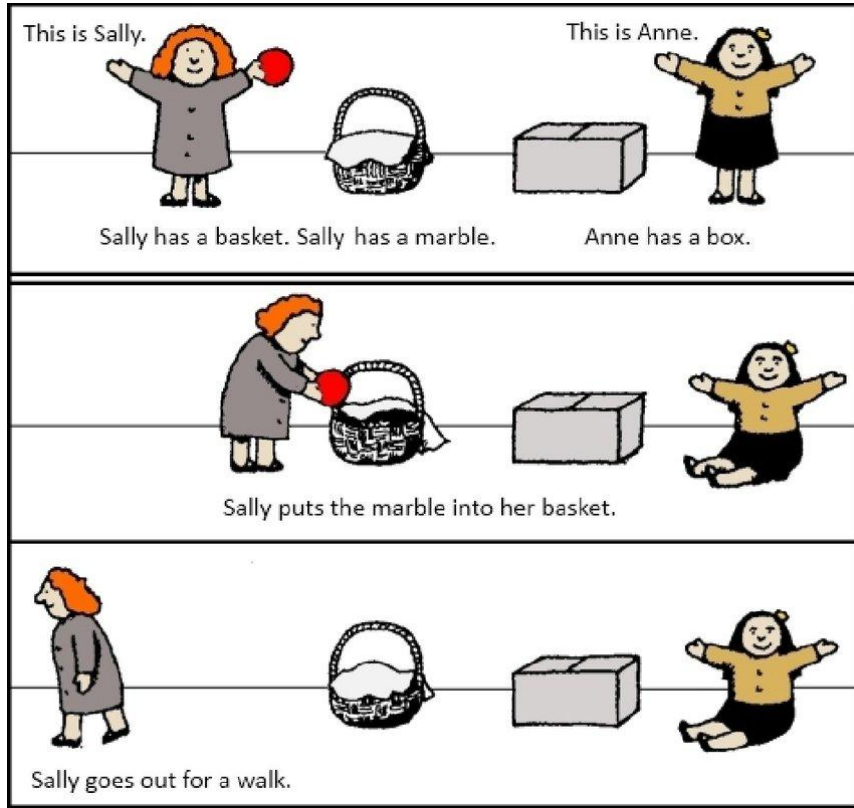
Interlude - a fun story!



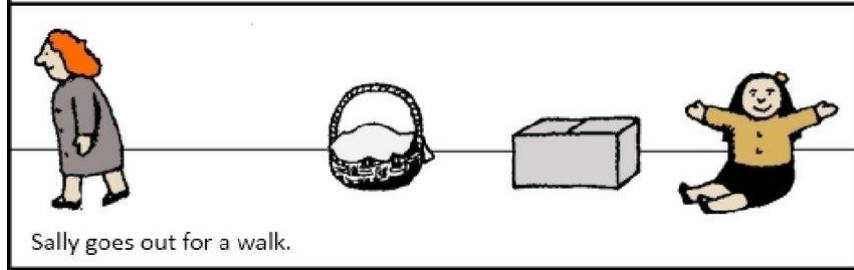
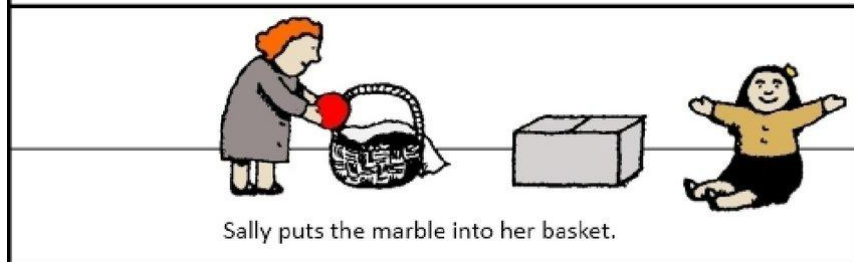
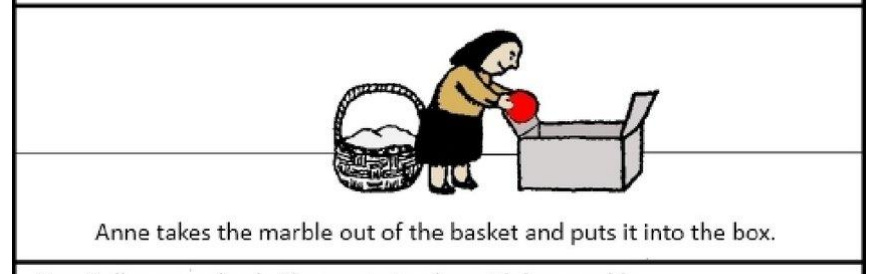
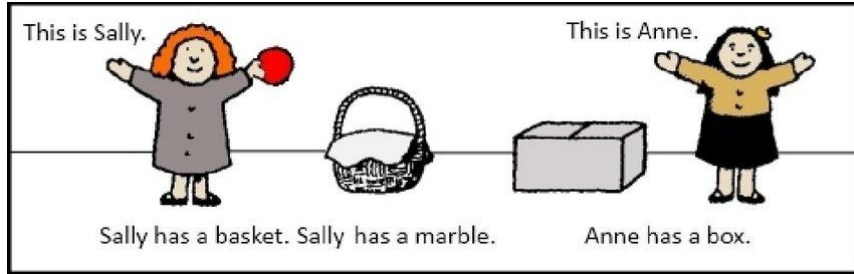
Interlude - a fun story!



Interlude - a fun story!

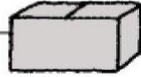


Interlude - a fun story!



Interlude - a fun story!

This is Sally.

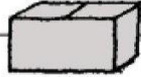


This is Anne.



Sally has a basket. Sally has a marble.

Anne has a box.



Sally puts the marble into her basket.



Sally goes out for a walk.



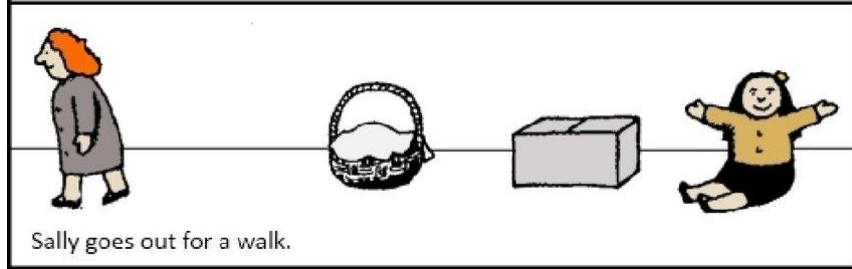
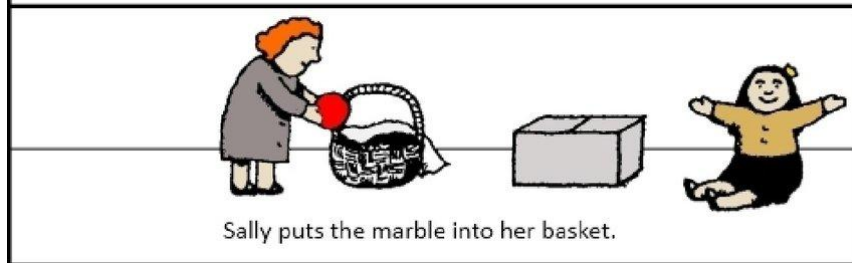
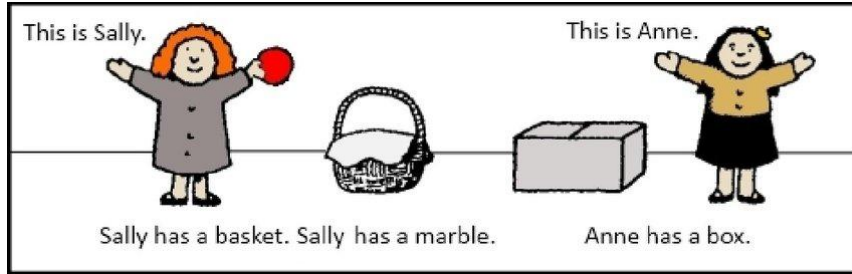
Anne takes the marble out of the basket and puts it into the box.

Now Sally comes back. She wants to play with her marble.



Where will Sally look for her marble?

The theory of mind-ers bread and butter: **False belief tasks**



Passed by 4-year olds, although some argue that that's down to verbal processing, not social cognition

How might theory of mind work?

Theory theory:

- Humans have a theory of how humans think and act, and consult it to infer mental states and predict behaviour
- disagreement on whether this theory is innate or learnt
- Also known als folk-/commonsense psychology

How might theory of mind work?

Theory theory:

- Humans have a theory of how humans think and act, and consult it to infer mental states and predict behaviour
- disagreement on whether this theory is innate or learnt
- Also known as folk-/commonsense psychology

Simulation theory:

- Theory of mind is based on simulating yourself in the situation
- Arguably still based on a model

How might theory of mind work?

Theory theory:

- Humans have a theory of how humans think and act, and consult it to infer mental states and predict behaviour
- disagreement on whether this theory is innate or learnt
- Also known as folk-/commonsense psychology

Simulation theory:

- Theory of mind is based on simulating yourself in the situation
- Arguably still based on a model

This debate seems to have died out (without a clear conclusion) in the 2000s.

How do we model theory of mind then?

Recent modelling work generally relies on theory theory, with the general assumption being that **other agents will act to maximize their utility**

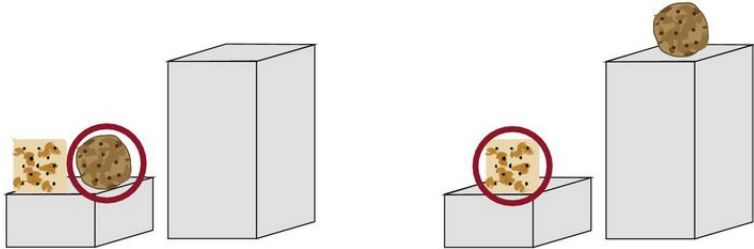
This can be modelled as a Bayesian process, wherein an action's **cost** (C) and **reward** (R) functions can be inferred based on the actions taken by an agent

$$p(C, R | \text{Actions}) \propto p(\text{Actions} | C, R) p(C, R)$$

Example experimental settings

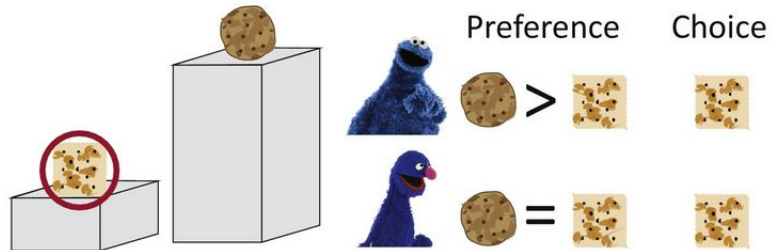
Preference inference

Which treat does ernie like the most?

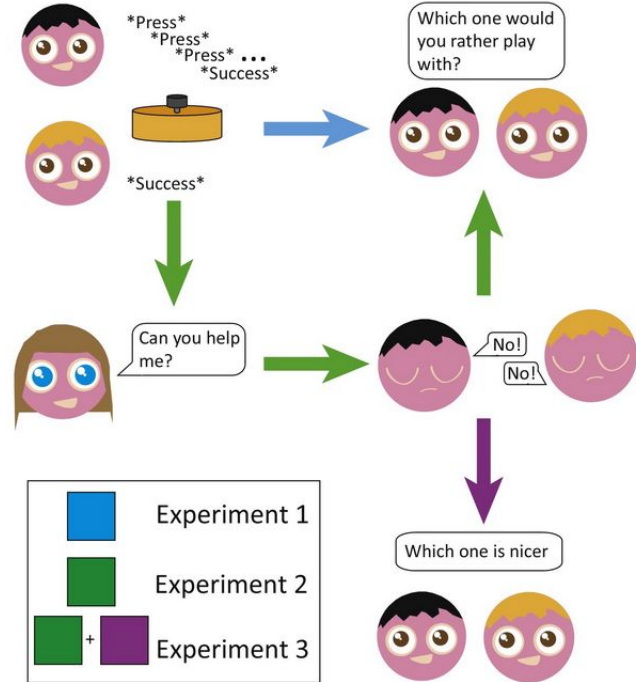


Competence inference

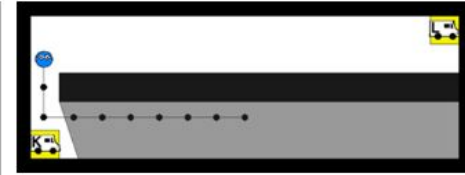
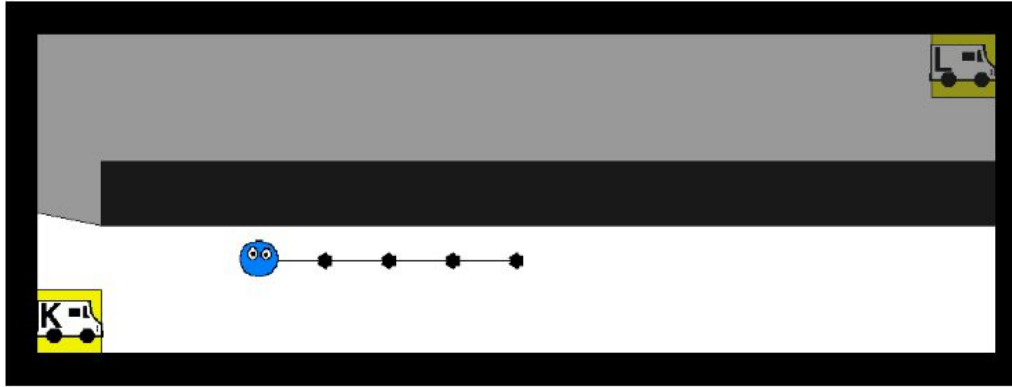
Which one is the one who cannot climb?



Social evaluations



Example experimental settings

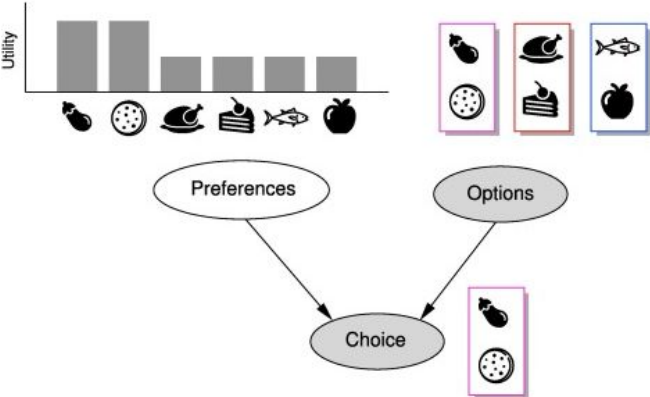


Baker et al. (2017) *Nature Human Behaviour*

Example experimental settings – more preference inference



Baker et al. (2017) *Nature Human Behaviour*



Jern., Lucas, & Kemp (2017) *Cognition*

Summary – theory of mind

The ability to infer other's mental state is called **theory of mind**

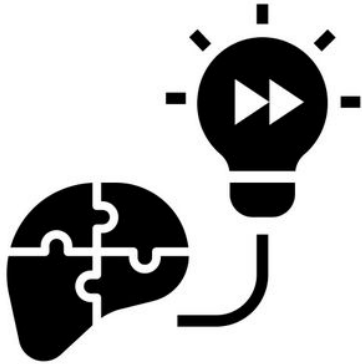
Historically, there were multiple ideas on how humans might do this, the big ones being **theory theory** (we have an internal theory of how humans work) and **simulation theory** (we put ourselves in the other person's shoes to understand what they might think)

Recent computational work generally relies on **Bayesian implementations** of theory theory for inference

Part III: Social learning and AI

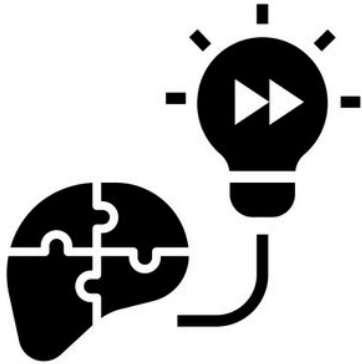
Why would we want socially learning AI?

Social learning reduces the amount of individual trial and error necessary to perform a task

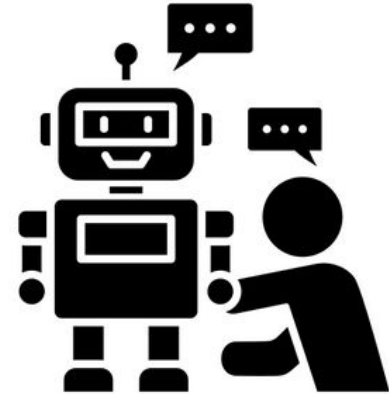


Why would we want socially learning AI?

Social learning reduces the amount of individual trial and error necessary to perform a task



Future AI agents should be able to interact with us, and social inference is essential to smooth, natural-feeling interactions



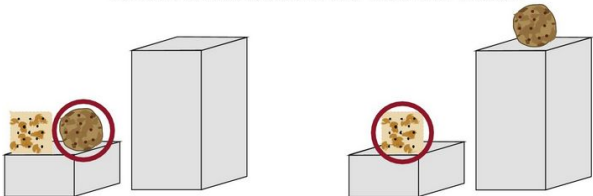
Since social learning is this good, let's just put it in AI systems then!!



So, about those inference models...

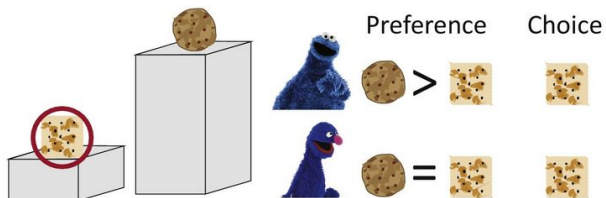
(A) Preference inference

Which treat does ernie like the most?

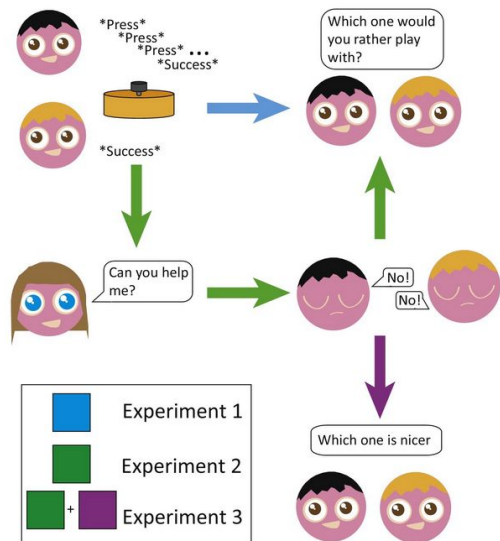


(B) Competence inference

Which one is the one who cannot climb?



(D) Social evaluations

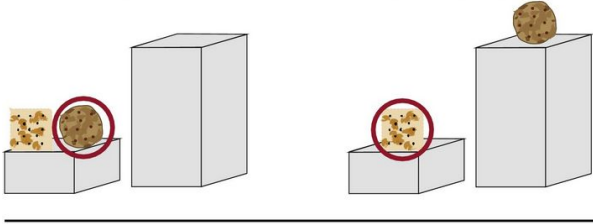


Notice any patterns?

Well, our inference problems are kind of limited so far though....

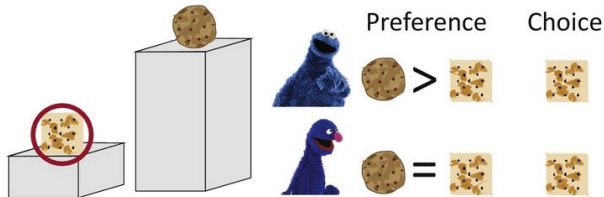
(A) Preference inference

Which treat does ernie like the most?

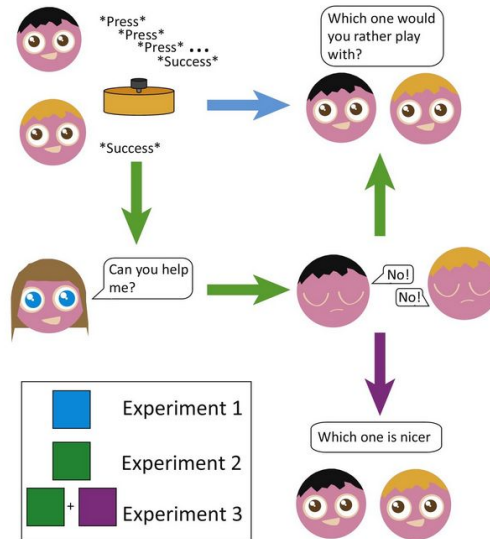


(B) Competence inference

Which one is the one who cannot climb?

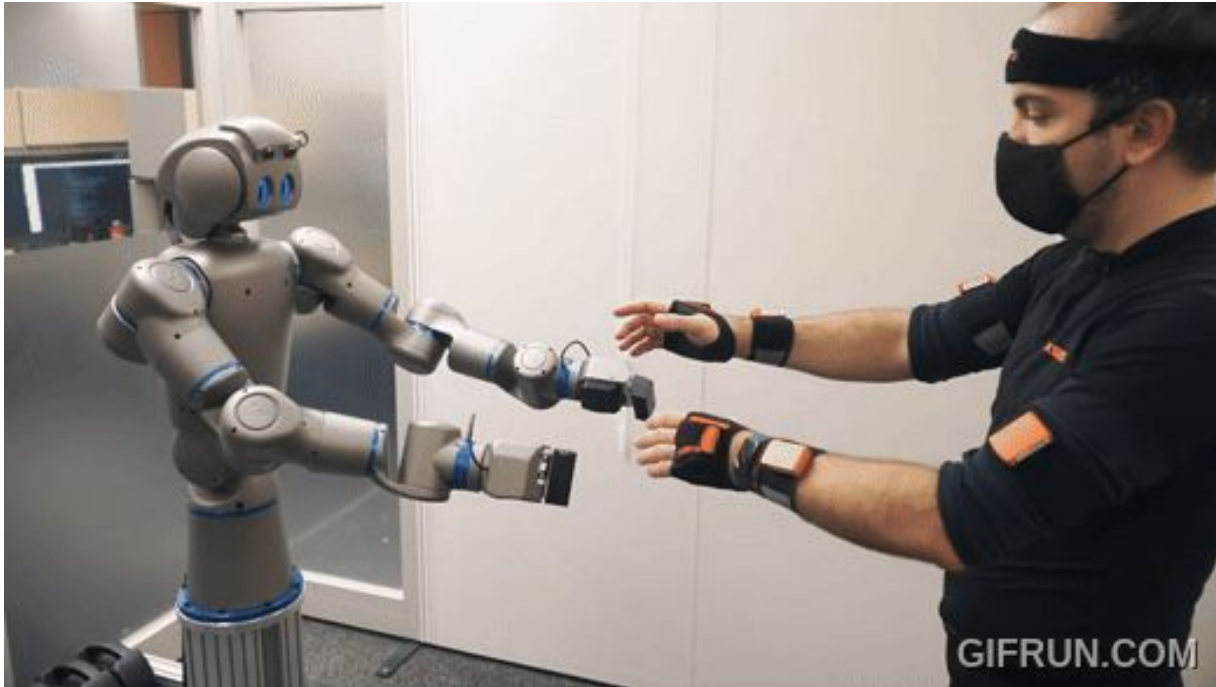


(D) Social evaluations



Bayesian inference becomes **computationally intractable** with too many options (although approximations exist) To get a Bayesian posterior, we need to integrate over the **entire option space** – easy enough in an experiment, very hard in real life (both because of the number of options and the difficulty in identifying all of them)

Also, even if we could infer already, robotics (and imitation robotics) isn't quite there yet



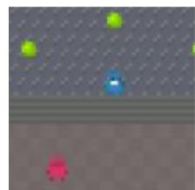
So what *can* AI do? – Multi-agent RL settings



Crewmates pick up and deposit gems



Impostor guards deposit



Crewmates distance from other players

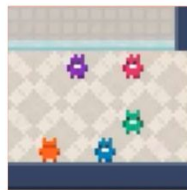
Common Behaviors

Crewmates Task



Crewmates partner together for tasks

Crewmates Vote



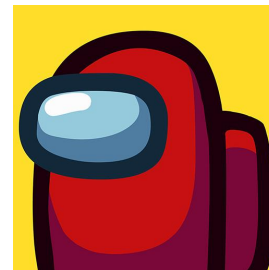
Crewmates abandon tasks in favor of voting

Impostor Freeze



Impostor isolates Crewmates

Equilibrium-specific Behaviors



Kopparapu et al. (2022) *arXiv*

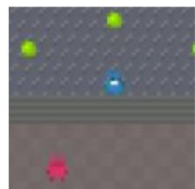
So what *can* AI do? – Multi-agent RL settings



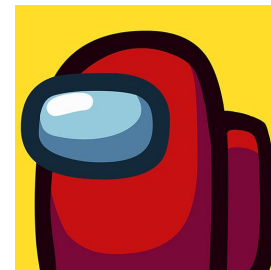
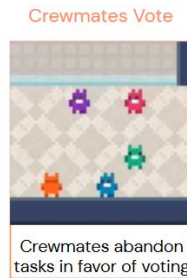
Crewmates pick up and deposit gems



Impostor guards deposit



Crewmates distance from other players

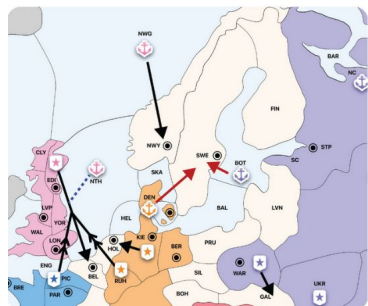


Common Behaviors

Equilibrium-specific Behaviors Kopparapu et al. (2022) *arXiv*

Depicted intents:

England conveys an army to Belgium with the support of France and Germany while taking Norway in a manner friendly to Russia.



ENG → FRA Mind supporting Edi - Bel?

ENG → GER Do you wanna support my convoy to Bel? With Italy going aggressive France will fall quickly and we can make gains off of both Russia and France.

ENG → RUS How are you thinking Germany is gonna open? I may have a shot at Belgium, but I'd need your help into Den next year.

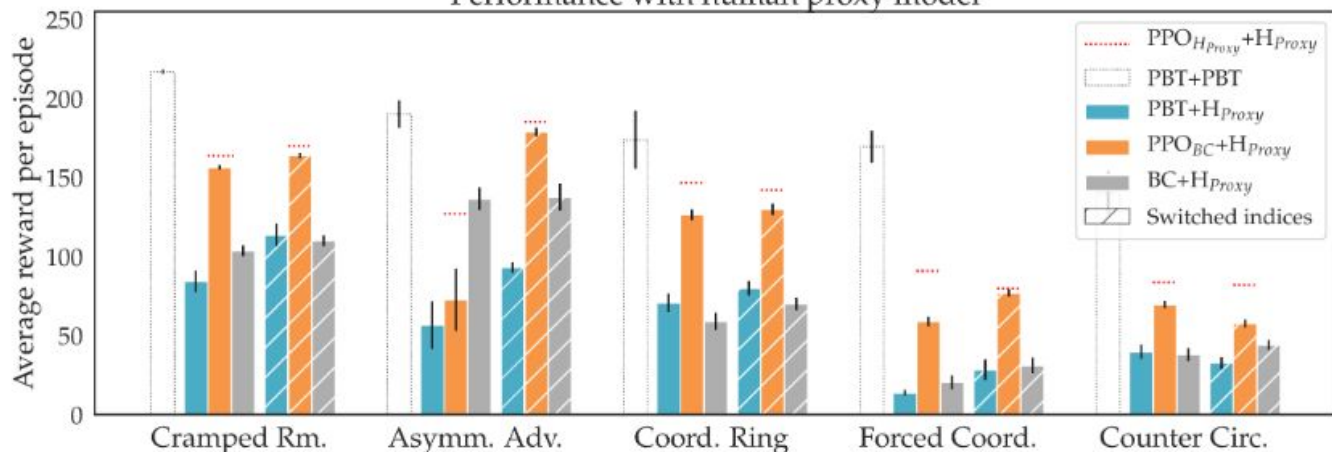
Meta Fundamental AI Research Diplomacy Team (2022), *Science*

Deep learning models can be trained to have (Diplomacy), or have emerging (Among Us) social coordination behaviours

But coordination success critically depends on who you train to play with



Performance with human proxy model



However, agents trained to play Overcooked performed significantly worse when paired with a human proxy model than with themselves, showing a lack of adaptability

So what *can* AI do? – LLMs

Does ChatGPT have Theory of Mind?

Bart Holterman
Utrecht University

Kees van Deemter
Utrecht University

Theory of Mind (ToM) is the ability to understand human thinking and decision-making, an

So what *can* AI do? – LLMs

DOES CHATGPT HAVE A MIND?

Simon Goldstein
The University of Hong Kong
simon.d.goldstein@gmail.com

B.A. Levinstein
University of Illinois at Urbana-Champaign
benlevin@illinois.edu

“we conclude that the data remains inconclusive”

Does ChatGPT have Theory of Mind?

Does ChatGPT have a typical or atypical Theory of Mind?

Margherita Attanasio^{1*}, Monica Mazza^{1,2}, Illenia Le Donne¹,
Francesco Masedu¹, Maria Paola Greco¹ and Marco Valenti^{1,2}

¹Department of Biotechnological and Applied Clinical Sciences, University of L'Aquila, L'Aquila, Italy,
²Reference Regional Centre for Autism, Abruzzo Region, Local Health Unit, L'Aquila, Italy

man
ersity

Kees van Deemter
Utrecht University

Theory of Mind May Have Spontaneously Emerged in Large Language Models

Authors: Michal Kosinski*¹

Can a conversational agent pass theory-of-mind tasks? A case study of ChatGPT with the Hinting, False Beliefs, and Strange Stories paradigms.

Eric Brunet-Gouet^{1, 2[0000-0002-3784-7817]}, Nathan Vidal^{2[0009-0009-1396-4703]}

and Paul Roux^{1, 2[0000-0003-0321-4189]}

So what *can* AI do? – LLMs

DOES CHATGPT HAVE A MIND?

Simon Goldstein
The University of Hong Kong
simon.d.goldstein@gmail.com

B.A. Levinstein
University of Illinois at Urbana-Champaign
benlevin@illinois.edu

“we conclude that the data remains inconclusive”

Does ChatGPT have Theory of Mind?

Does Ch
atypical?

Margherita Attanas
Francesco Masedu¹
¹Department of Biotechnologica
²Reference Regional Centre for A

MAYBE (but it might also just be familiar with the verbal tasks)

Theory of Mind May Have Spontaneously Emerged in Large Language Models

Authors: Michal Kosinski*¹

Can a conversational agent pass theory-of-mind tasks? A case study of ChatGPT with the Hinting, False Beliefs, and Strange Stories paradigms.

Eric Brunet-Gouet^{1, 2}[0000-0002-3784-7817], Nathan Vidal²[0009-0009-1396-4703]

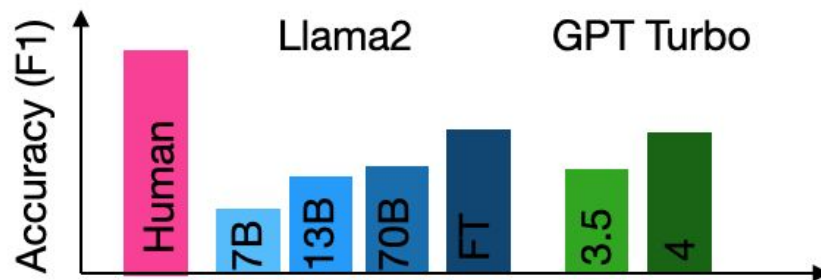
and Paul Roux^{1, 2}[0000-0003-0321-4189]

But also...

OpenToM Benchmark (Xu et al., 2024)



Q: What is Sam's attitude toward's Amy's action?



Summary – social learning and AI

Social learning would be highly beneficial for AI, both to reduce costly trial and error, and to enable smoother interactions with humans

Computational models of social cognition are generally **limited to experimental settings**, and are hard to scale to real-world proportions

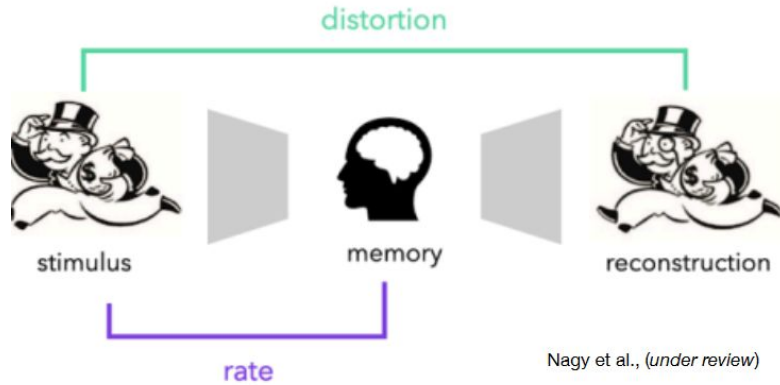
Deep learning models can solve multi-agent coordination tasks, but don't generalize well, and struggle to adjust to new partners (implying that they don't learn social inference)

ChatGPT can solve verbal mentalizing tasks, but it is unclear if this actually reflects social cognition

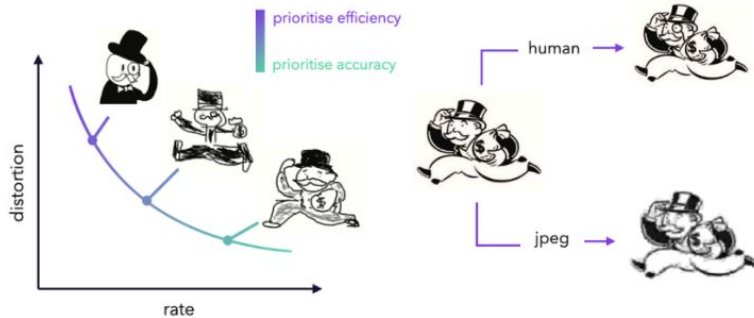
→ While social learning would be greatly beneficial for AI, we're still pretty far from making it a reality

Thanks for your attention! :)

Next week: Compression and resource constraints



Nagy et al., (under review)



Dr. David Nagy