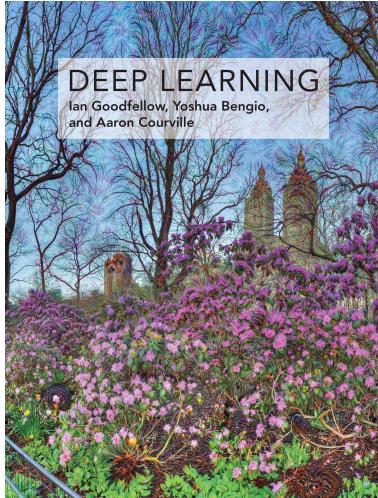


# ML4AI: What Next?

## Yoshua Bengio

July 27th, 2019, CIFAR Deep Learning & Reinforcement  
Learning Summer School  
Amii, Edmonton, Canada



Université  
de Montréal



**CIFAR**  
CANADIAN  
INSTITUTE  
FOR  
ADVANCED  
RESEARCH

**ICRA**  
INSTITUT  
CANADIEN  
DE  
RECHERCHES  
AVANCEES

# Current AI is far from Human-Level AI

- Sample complexity is too high for supervised learning, even more for RL
  - Real-world actions can be lethal, experience is limited & costly
  - We don't have a good simulator of the real world (esp.w/ humans)
- High-level concepts typically need to be provided by human designers or labelers
- Errors made by trained systems reveal that their 'understanding' is very shallow and superficial
- The dream of deep learning discovering and disentangling high-level explanatory variables is far from achieved

## What I don't believe works towards human-level AI

- NLP based purely on text
- Generative models purely on sensory data (generating images is cute, though)
- Relying on very strong prior knowledge (e.g. the graphics engine for vision)
- Working on algorithms with no chance to scale to brain-size (e.g., tabular distributions, non-distributed or non-learned representations)
- Theories not compatible with the situation of animal/human agents, e.g., iid assumption
- Agents with a perfect simulator of the environment
- Planning in pixel space every 100ms

# What Next? Grounded Language Understanding

# Alien Language Understanding: a Thought Experiment

- ▶ Imagine yourself approaching another planet and observing the bits of information exchanged by aliens communicating with each other
- ▶ Unlike on Earth, their communication channel is noisy, but like on Earth, bandwidth is expensive → the best way to communicate is to maximally compress the messages, which leads to sequences of random bits being actually exchanged.
- ▶ If we only observe the compressed messages, there is no way we can ever understand the alien language



# Alien Language Understanding: a Thought Experiment

- ▶ How can we learn to understand the alien language?
- ▶ We need to do grounde language learning: we need to observe what the aliens are doing jointly with their messages, to try to decipher their intentions, context, etc.
- ▶ For this we need to build an 'Alien World Model' which captures the causal structure of their behaviors and resulting changes in their environment.



# Jointly Learning Natural Language and a World Model

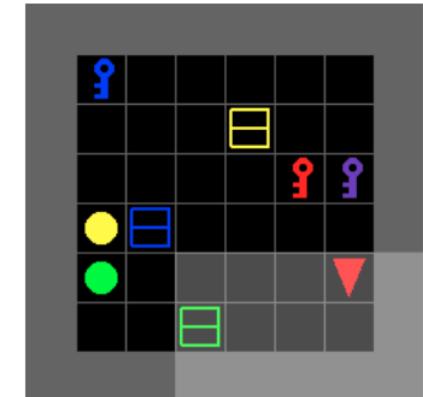
- Should we first learn a world model and then a natural language description of it?
- Or should agents jointly learn about language and about the world?
- I lean towards the latter.
- Consider top-level representations from supervised ImageNet classifiers. They tend to be much better and easier to learn than those learned by unsupervised learning. Why?
- Because language (here object categories) provides to the learner clues about relevant semantic high-level factors from which it is easier to generalize.
- See my earlier paper on cultural evolution, which posits that culture can help a learner escape from poor optimization, guide (through curricula) the learner to better explanations about the world.

# BabyAI Platform *Chevalier-Boisvert et al & Bengio ICLR 2019*

**Purpose:** simulate language learning from a human and study data efficiency

**Comprises:**

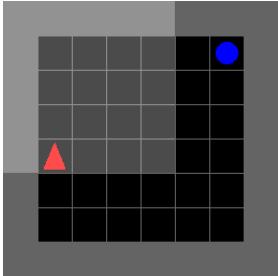
- a gridworld with partial observability (Minigrid)
- a compositional natural-looking Baby language with over  $10^{19}$  instructions
- 19 levels of increasing difficulty
- a heuristic stack-based expert that can solve all levels



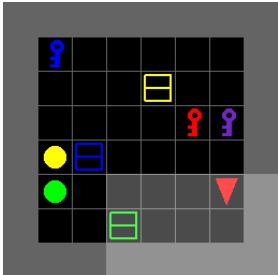
(b) PutNextLocal:  
"put the blue key next  
to the green ball"

[github.com/mila-udem/babyai](https://github.com/mila-udem/babyai)

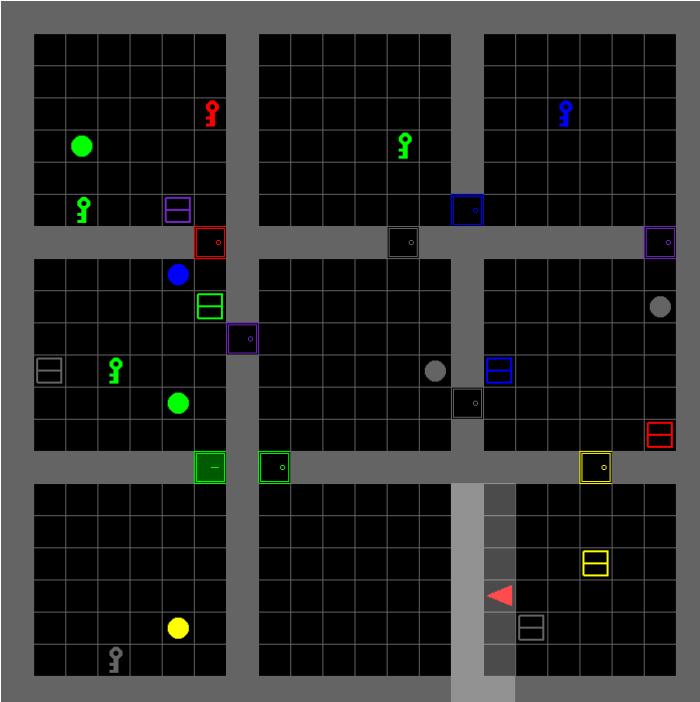
# Early Steps in Baby AI Project



(a) GoToObj: "go to  
the blue ball"



(b) PutNextLocal:  
"put the blue key next  
to the green ball"



(c) BossLevel: "pick up the grey box behind you, then go to the grey key and open a door". Note that the green door near the bottom left needs to be unlocked with a green key, but this is not explicitly stated in the instruction.

- Designing and training experts for each level, which can serve as teachers and evaluators for the Baby AI learners
  - Partially observable, 2-D grid, instructions about objects, locations, actions

go to the red ball

open the door on your left

put a ball next to the blue door

open the yellow door and go to the key behind you

put a ball next to a purple door after you put a blue box next to a grey box and pick up the purple box

# Integrating System 1 and System 2

- System 2 model is very coarse and imperfect
- System 2 abstract concepts need to be grounded via system 1
- System 2 thinking allows counterfactual reasoning, i.e., imagining scenarios which did not and will not happen, as an exercise (e.g. for credit assignment, if I had done that...), allows generalization far from the training data, imagine dangerous scenarios without having to take the actual risks
- System 2 is too slow and inefficient, compile to system 1 into habits and intuitive behavior

# What Next? Generative Models in Latent Space

# The Need for Unsupervised Learning

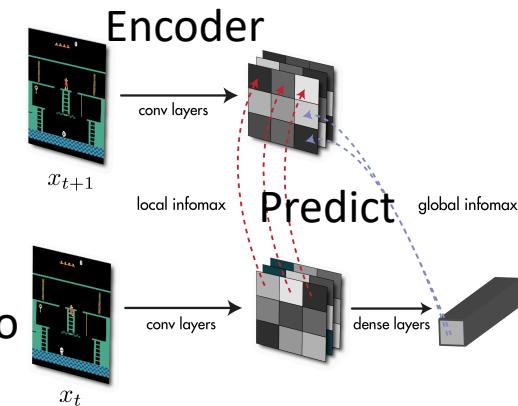
- “The brain has about  $10^{14}$  synapses and we only live for about  $10^9$  seconds. So we have a lot more parameters than data. This motivates the idea that we must do a lot of unsupervised learning since the perceptual input (including proprioception) is the only place we can get  $10^5$  dimensions of constraint per second.”
  - Geoffrey Hinton
- Plus, unlabeled or un-rewarded data is much cheaper to obtain

# Self-Supervised Learning

- Predict parts (or functions) of the data from other parts
  - E.g. denoising auto-encoders (some inputs are masked)
- SOTA in NLP is using self-supervised learning
  - Word2Vec, BERT, XLNet

# Generative Models in Latent Space

- For human-like brains, generative models are useful for planning (**model-based RL**), imagination, counterfactuals, inference over causes and explanations, high-level credit assignment
  - NONE OF THIS REQUIRES WORKING IN PIXEL SPACE
- Current generative models are trained wrt pixel-space objectives, how to train purely in the space of abstract representations? We want the encoder mapping pixel space to abstract space to be trained wrt the high-level goals too.
- There is an issue of possible collapse of representations if we maximize predictability (e.g. max likelihood) in latent space

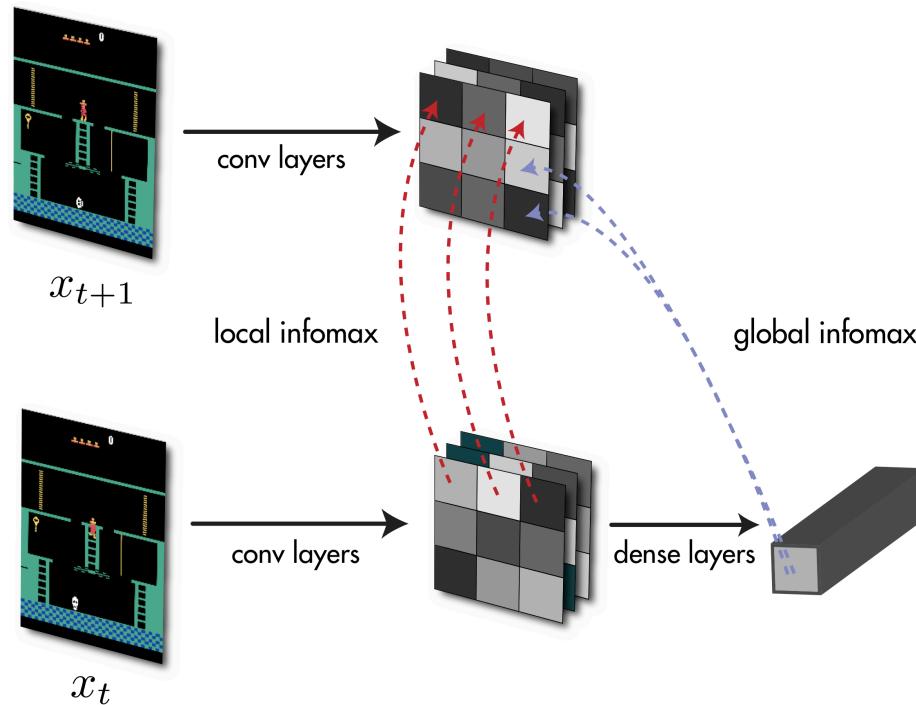


# Spatio-Temporal Deep InfoMax (STDIM)

**Self-supervised learning** in latent space rather than data space

We can avoid the collapse by using entropy maximization or maximizing mutual information involving the encoder output

Deep InfoMax or DIM  
(Hjelm et al & Bengio ICLR 2019)



# Learning « How the world ticks »

- So long as our machine learning models « cheat » by relying only on superficial statistical regularities, they remain vulnerable to out-of-distribution examples
- Humans generalize better than other animals thanks to a more accurate internal model of the **underlying causal relationships**
- To predict future situations (e.g., the effect of planned actions) far from anything seen before while involving known concepts, an essential component of reasoning, intelligence and science

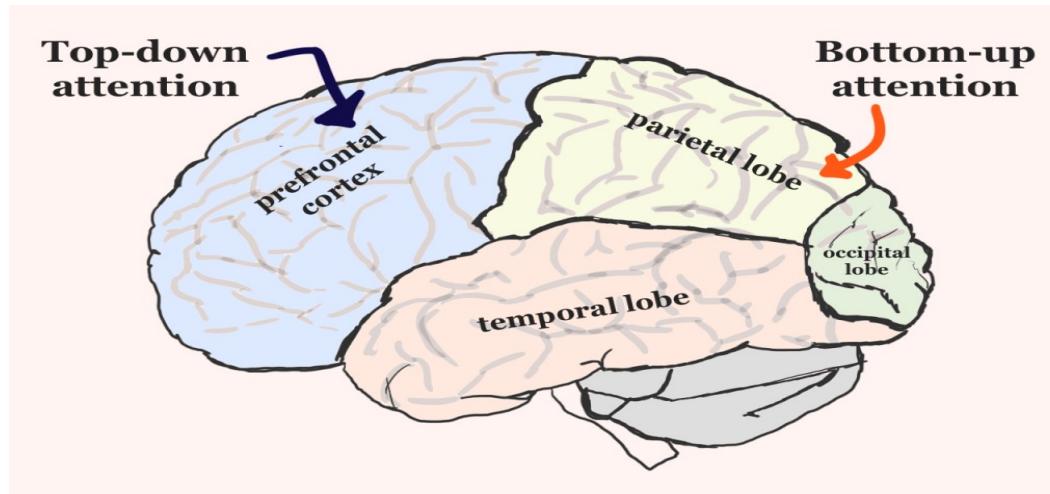
# Generating Just Selected Dimensions: the Consciousness Prior

- Human planning, counterfactuals, reasoning, etc. does not take place over the full state-space (even at the abstract level)
- If we don't predict and sample all high-level variables, what kind of training objective to use? again, not maximum likelihood even in latent space, because that would maximize  $P(\text{full next state} \mid \text{full previous state})$
- Again, information-theoretical objectives may come to the rescue, along with attention mechanisms to select what to predict about the future, what to pay attention to during learning

# On the Relation between Abstraction, Thought and Attention

- A thought is a low-dimensional object, few aspects of the state
- Attention allows us to focus on a few elements out of a large set
- Soft-attention allows this process to be trainable with gradient-based optimization and backprop

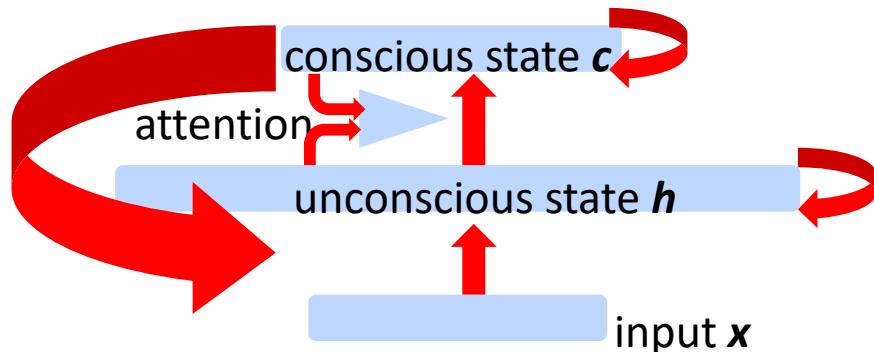
Attention focuses on a few appropriate abstract or concrete elements of mental representation



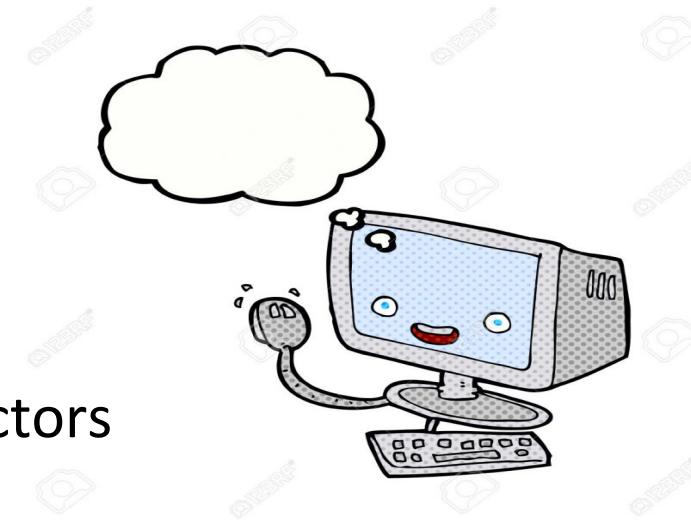
# The Consciousness Prior

Bengio 2017, arXiv:1709.08568

- 2 levels of representation:
  - High-dimensional abstract representation space (all known concepts and factors)  $h$
  - Low-dimensional conscious thought  $c$ , extracted from  $h$



- $c$  includes names (keys) and values of factors



# Why do I call it a Prior?

- There is something very special about the kind of high-level variables which we manipulate with language:
  - we can predict some given very few others
    - E.g. "if I drop the ball, it will fall on the ground"
  - corresponds to a **sparse factor graph**

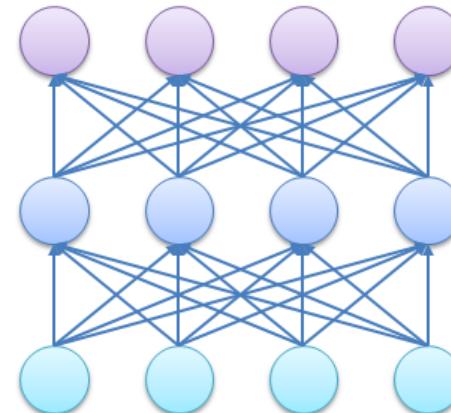
$$P(V) \propto \prod_k \phi(V_{s_k})$$

where  $V_{s_k}$  is  
the subset of  $V$   
with indices  $s_k$

# What Next? Linking Agency and Representation

# How to Discover Good Disentangled Representations

- How to discover abstractions?
- What is a good representation? (*Bengio et al 2013*)
- Need clues (= priors) to help **disentangle** the underlying factors (**not necessarily statistically independent**), such as
  - Spatial & temporal scales
  - Approximate marginal independence
  - Simple dependencies between factors
    - *Consciousness prior*
  - Causal / mechanism independence
    - *Controllable factors*



# Acting to Guide Representation Learning & Disentangling



(E. Bengio et al, 2017; V. Thomas et al, 2017)

- Some factors (e.g. objects) correspond to ‘independently controllable’ aspects of the world
  - Corresponds to maximizing mutual information between intentions (goal-conditioned policies) and changes in the state (trajectories), conditioned on the current state.
- *Can only be discovered by acting in the world*
  - *Control linked to notion of objects & agents*
  - *Causal but agent-specific & subjective: affordances*

# Unsupervised Agents ≠ Intrinsic Rewards

- Curiosity, novelty, knowledge seeking
  - need to estimate model uncertainty, information gain
  - should be novel in the abstract space, not just pixel space
  - need to avoid the trap of inherently unpredictable (white noise TV)
  - must be not just novel but not easily explainable with current abstract model, requires learning (seek situations where you learn)
  - exploratory policy maximizes that while at the same time we train a model to explain away observations, i.e., minimize future surprises (adversarial setting)
  - when exploring, need to take safety into account (less so when you are young!)
- Controllability, empowerment
  - babies clearly learn to control their environment and their body
  - related to notion of freedom, maximizing entropy of set of future trajectories

# What Next? Transfer Learning and Meta-Learning

# Meta-Learning / Learning to Learn

- Generalize the idea of hyper-parameter optimization
  - Inner loop optimization (normal training), a fn of meta-params

$$\theta_t(\omega) = \text{approxmin}_{\theta} C(\theta, \omega, \mathcal{D}_{train}^t)$$

- Outer loop optimization (meta-training), optimize meta-params

$$\omega = \text{approxmin}_{\omega} \sum_t L(\theta_t(\omega), \omega, \mathcal{D}_{test}^t)$$

- Meta-parameters can be the learning rule itself (Bengio & Bengio 1991; Schmidhuber 1992), learn to optimize
- Meta-learn an objective or reward function, or a shared encoder
- Meta-learning can be used to learn to generalize or transfer
- Can backprop through  $\theta_t$ , use RL, evolution, or other tricks

# Missing from Current ML: Understanding & Generalization Beyond the Training Distribution

- Learning theory only deals with generalization within the same distribution
- Models learn but do not generalize well (or have high sample complexity when adapting) to modified distributions, non-stationarities, etc.
- Poor reuse, poor modularization of knowledge
- Meta-learning is an end-to-end learning approach to improving transfer learning and fast adaptation to changes in distribution

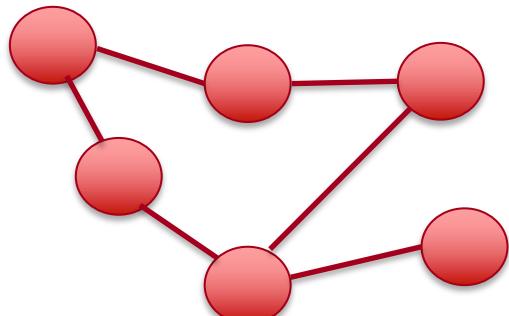
# *Beyond iid: Hypotheses about how the environment changes Independent Mechanisms and the Small Change Hypothesis*

- Independent mechanisms:
  - changing one mechanism does not change the others (*Peters, Janzig & Scholkopf 2017*)
- Small change:
  - Non-stationarities, changes in distribution, involve few mechanisms at a time (e.g. the result of a single-variable intervention)
- How can we discover these independent mechanisms, i.e., factor knowledge?

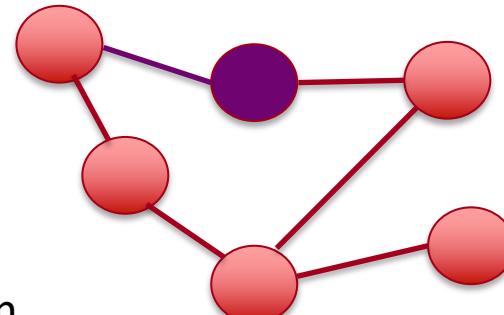
What Next?  
Breaking Knowledge into  
Re-Composable Pieces

# Separating Knowledge in Small Re-Usable Pieces

- Pieces which can be re-used combinatorially
- Pieces which are stable vs nonstationary, subject to interventions



Change due  
to intervention



# Wrong Knowledge Factorization Leads to Poor Transfer

- With the wrong factorization  $P(B) P(A|B)$ , a change in ground truth  $P(A)$  influences both modules, all the parameters
  - poor transfer: all the parameters must be adapted
- This is the normal situation with standard neural nets: every parameter participates to every relationship between all the variables
  - this causes *catastrophic forgetting, poor transfer, difficulties with continual learning or domain adaptation, etc*

# What Next? Discovering Causal Variables and Causal Structure

# Causality not Captured by Joint Distribution

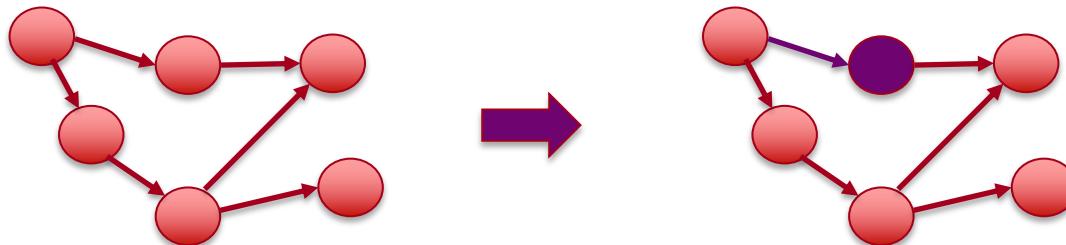
- Knowing the full joint distribution generally does not tell us which variable is a cause, which is an effect, and in general the question of causality may be meaningless in the wrong space
- Agents need to know about causal structure in order to properly infer how the joint distribution would change under interventions (theirs or of other agents)
- Understanding causal structure allows counterfactual reasoning and making sense of changes in distribution due to agents

# Deep Learning Objective: discover high-level representation capturing cause and effect variables

- What are the right representations?
  - Causal variables explaining the data
- How to discover them? (learn the mythical encoder)
- How to discover their causal relationship, the causal graph?

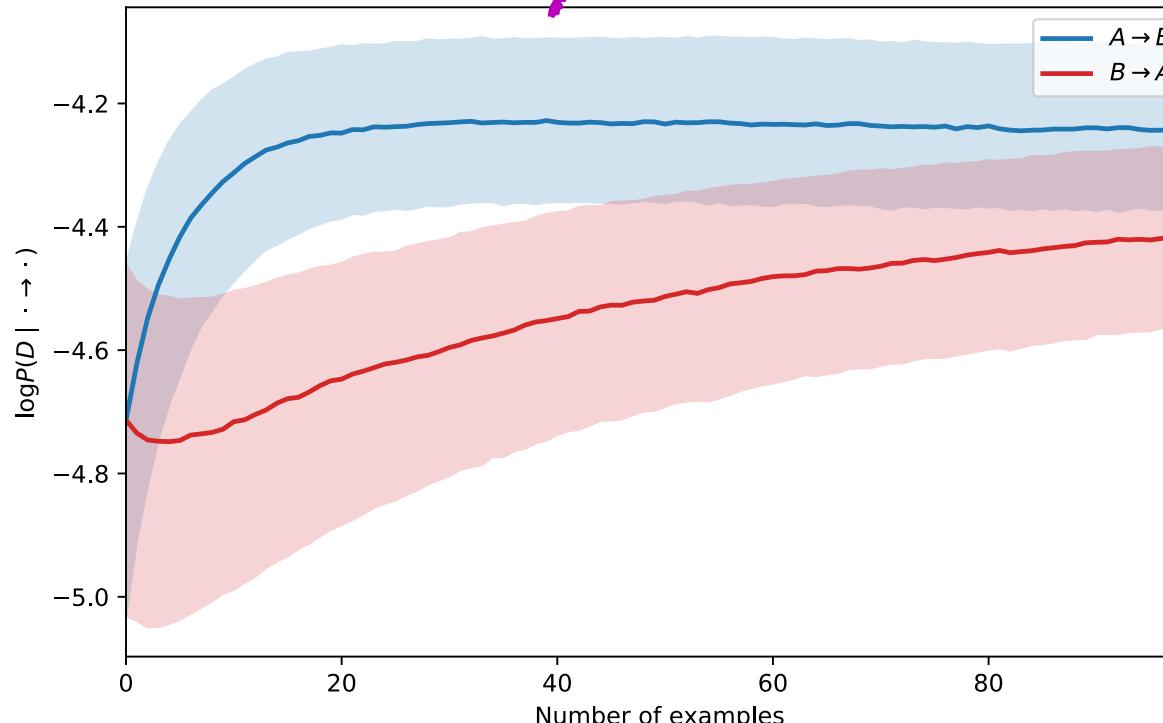
*Small Change →  
Small Sample Complexity*

Few parameters need to change → small L2 change → *few examples needed to recover from the change*



Under the right parametrization → fast adaptation to interventions

# Empirical Confirmation: Correct Causal Structure Leads to Faster Adaptation



$A \rightarrow B$  is the correct causal structure: faster online adaptation to modified distribution = lower NLL regret

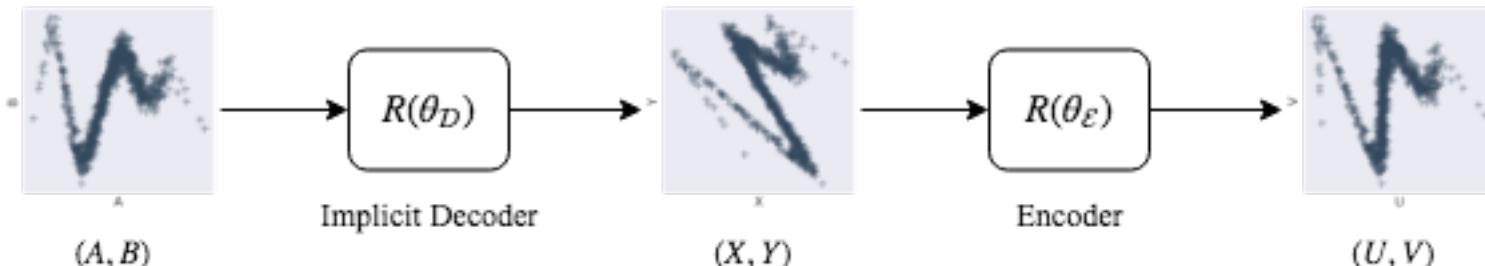
# Turning a Hindrance into a Useful Signal

ArXiv paper, Bengio et al 2019: *A Meta-Transfer Objective for Learning to Disentangle Causal Mechanisms*

- Changes in distribution (nonstationarities in agent learning, transfer scenarios, etc) are seen as a bug in ML, a challenge
- Turn them into a feature, an asset, to help discover causal structure, or more generally to help **factorize knowledge**:
- **Tune knowledge factorization (e.g. causal structure) to maximize fast transfer**
- *"Nature does not shuffle environments, we shouldn't"*  
L. Bottou

# Disentangling the Causes

- Realistic settings: causal variables are not directly observed
- Need to learn an encoder which maps raw data to causal space
- Consider both the encoder parameters and the causal graph structural parameters as meta-parameters trained together wrt proposed meta-transfer objective



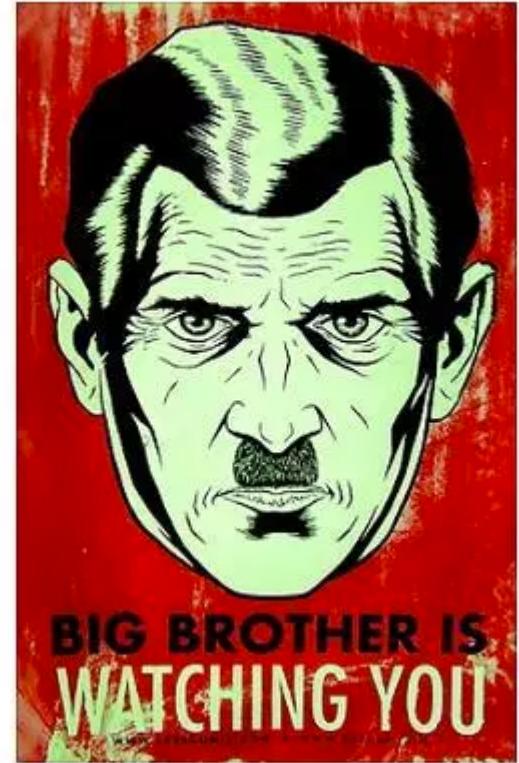
# *Observing Other Agents*

- Can infants figure out causal structure in spite of being almost passive observers?
- Yes, if they exploit and infer the interventions made by other agents
- Our approach does not require the learner to know what the action/intervention was (but it could do inference over interventions)
- But more efficient learning if you can experiment and thus test hypotheses about cause & effect

What about AI and  
Society?

# Dangers and Concerns with AI

- Big Brother and killer robots
- Misery for jobless people, at least in transition
- Manipulation from advertising and social media
- Reinforcement of social biases and discrimination
- Increased inequality and power concentration in few companies



# Advertising is Psychological Manipulation & Hurts Innovation

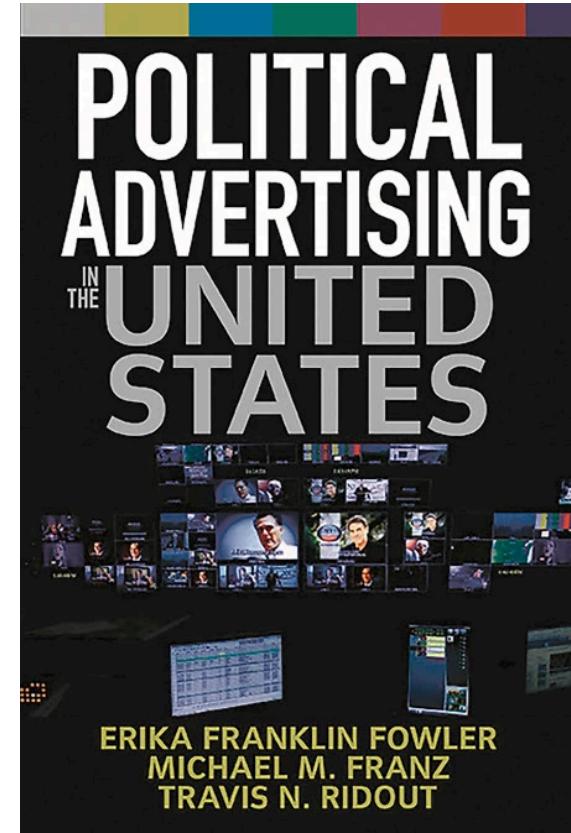
- **Moral hazard:**

- Psychological manipulation: for who's benefit? The advertiser.
- Targeting children forbidden in Canada because they are vulnerable
- Political advertising
- Need to draw a red line

- **Economic drag:**

- favours large incumbants and thus slows innovation

- **Bad Nash equilibrium**



# Killer Robots: Moral & Security Threat

- 62% worldwide in favour of a ban to limit LAWS
- UN Secretary General: "The prospect of machines with the discretion and power to take human lives is morally repugnant... could trigger new arms races"
- 28 states calling for an international ban
- Leading AI scientists signed letters to governments
- 4000 Google employees named "Arms Control Person of the Year"
- Companies, countries, civil society must rise
- Example of landmines



# Applications I don't want to work on

- How does it fit with my values? Ask
  - How is this technology going to be used?
  - Who will benefit or suffer from it?
  - How much and what impact?
- I do not wish to work on
  - military applications, especially lethal autonomous weapons
  - ML to beat the stock market
  - increasing the efficiency of advertising
  - increasing the efficiency of oil & gas industry, meat industry...

# The Wisdom Race

- Collective and individual wisdom has increased, e.g., the decrease in crimes, wars and extreme poverty
- But not fast enough to catch up with the rise in power of the tools we are building, which enable power concentration and can be dangerous
- This gives those with influence, in particular AI researchers, a responsibility



# Montreal Declaration: Ten Principles



## Montréal Declaration Responsible AI\_

# Well-being

# Respect for autonomy

# Protection of privacy

## Solidarity

## Democratic participation

# Equity

## Diversity

## Prudence

# Responsibility

## Sustainable development