

Eyes on the Prize

- Alberta Machine Intelligence Institute
 - University of Alberta
 - DeepMind Alberta
- Reinforcement Learning and Artificial Intelligence Lab







Rich Sutton





20 years ago, Alberta had the foresight to invest \$2M/year in machine learning

which has paid off handsomely in science and commerce in Alberta

What should Alberta invest in today that will pay off in the next 20 years?

- My answer is to double down on the science of AI
- of understanding intelligence

There is still a long way to go; we should keep our Eyes on the Prize

Alberta should keep its Eyes on the Prize

- The Prize is to understand the principles of intelligence—what it is and how it works—well enough to create (or become) beings of greater intelligence
- The Prize is a fundamental goal of science, engineering, and the humanities
- Achieving it will change the way we work and play, our sense of self, our sense of life and death, and the goals we set for ourselves and our societies
- Achieving it will be an event comparable in significance to the rise of human life on Earth

In Alberta, we should strive to be relevant to this great goal



Understanding intelligence would certainly be a big deal

- But some doubt that it is possible
- More doubt that it will happen anytime soon
- Many doubt that it would be a good thing

Is the Prize within reach? And would it be good, or bad?

Computer power/\$ increases exponentially, with no end in sight, creating a powerful persistent pressure for understanding intelligence



Brain-scale computer power will cost ≈\$1000 in ≈2030

This estimate is rough but robust: a factor of $10 \approx 5$ years

 \Rightarrow Al increases in value by a factor of 10 every 5 years

And so does the pressure to find the algorithms/software

I estimate a 50% probability of human-level AI by 2040



Investing in AI science is a good bet for Alberta

- in 20 years it is a coin flip
- in 50 years it's a near certainty

The Prize is within reach

Human-level AI will probably happen in the next decades;

• Even if we lose the coin flip, investing in AI is likely to have been good; the partial solutions will likely be commercially valuable

 Automation and computer power will continue to increase in importance in the economy; Moore's law will chug on

• In 10 years human-level AI is a moonshot; in 20 years it's a coin flip;

We don't get to choose on most big things

- The impact of the printing press and literacy
- The impact of the internet and social media
- The decline of monarchies and rise of democracies
- The shift of world economies from farming to industry
- The shift of Alberta's economy from oil to technology

Many are impacted. Who are we to choose for all?



Alberta's economy has been dominated by oil, and now it is changing. Is this good or bad?

- This is (or may be) happening. Good for some, bad for others...
- But the real good or bad is how we react to the transformation
 - Do we refuse to accept it and cling to all the old ways?
 - Or do we see, prepare for, and conform to the future?
- So too for the coming of genuine AI
 - We should prepare for and conform to the future
 - We can surf the future, but we can't hold back the wave

Understanding intelligence is inevitable. We choose to be a part of it because:

- Selfishly:
 - of the benefits... and more of the glory

In my opinion, understanding intelligence is one of the few great goods in the universe

 Greater intelligence (better foresight and decisions) in the world means greater productivity with less cost, waste, and conflict

• We will understand ourselves and our place in the universe better, which is what we, as intelligent agents, are built to do

• Early participants in understanding intelligence will reap more



Outline

- The Prize of understanding intelligence
- Why we seek it



- on the Prize
 - by focusing on experience

How reinforcement learning research in Alberta keeps its eyes

Al research in Alberta (Amii and its eco-system)

- We are doing a mix of applications and longer-term research
 - but we are striving to be
 - While keeping our eyes on the longer-term Prize of understanding intelligence
- Both are super important, but they are different:
 - To get a successful application, you try to avoid unsolved problems
 - In long-term research, you are drawn to unsolved problems
- Let me tell you a little about what it's like to be drawn to the unsolved problems of AI

• As Martha White said, we are not yet a leader in the application of RL,





First, we characterize the *problem* of intelligence

Intelligence is:

"the most powerful phenomena in the universe"

"attaining consistent ends by variable means" –William James, 1890

"the computational part of the ability to achieve goals" –John McCarthy, 1997

"the computational part of an agent's ability to predict and control a stream of sensations" -me, now

-Ray Kurzweil, 2011



Intelligence is:

The computational part of an agent's ability to predict and control a stream of sensations, particularly a designated numerical sensation (called reward)



The reward hypothesis:

"That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward)"





Intelligence is:

particularly a designated numerical sensation (called reward), while interacting with a vastly more complex world



The computational part of an agent's ability to predict and control a stream of sensations,

The Big-World perspective:

The world is much bigger than the agent

It contains many other agents!





Experience up to time step 7 (think of a time step as ≈ 0.1 sec)

gı	nal	ls i	nc	luc	din	g	Reward	
1	0	2	0	0	0	3	-5.3	
1	0	3	0	3	0	1	5.5	
0	2	2	3	1	0	0	-8.7	
2	1	0	2	3	2	1	-8.4	
3	1	0	2	2	3	3	-0.4	
1	1	2	0	0	1	3	2.7	
)	0	0	2	3	0	3	-5.3	
2	1	3	1	1	3	3	-2.8	most recent sensation





Experience up to time step 7 (think of a time step as ≈ 0.1 sec)



Experience up to time step 7 (think of a time step as ≈ 0.1 sec)



Experience up to time step 7 (think of a time step as ≈ 0.1 sec)

	Time step	Action signals	Sensory s
 Different sensory signals can be qualitatively different from each other 	0	$\bigcirc \bigcirc \bigcirc \bigcirc$	
quantatively amerent north saor strict	1	000	
 In their range of values 	2	000	
• In their predictive relationships	3	000	
 In their predictive relationships 	4		
 to action signals 	6	000	
e to ocob othor	7	$\circ \circ \circ$	
• to each other	8	000	
 to themselves 	9	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	10	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 There are short-term and long-term 	11	$\bigcirc \bigcirc \bigcirc$	
patterns in these data	12	000	
	13	000	
 There are many things to predict 	14	000	
 Prediction need not be just of the 	15	000	
sensorv signals	16		
	18		
 The most important predictions are of 	19	000	
<i>functions</i> of future sensory signals	20	000	
• a a pradictions of value the	21	000	
• e.g., predictions of value, the	22	$\bigcirc \bigcirc \bigcirc \bigcirc$	
discounted sum of future reward	23	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 e.g., General value functions (GVFs) 	24	000	
	25	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 predict any signal, not just reward 	26	$\bigcirc \bigcirc \bigcirc \bigcirc$	
• over a flavible temperal epyelepa	27	$\bigcirc \bigcirc \bigcirc \bigcirc$	
• over a nexible temporal envelope	28	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 contingent on any policy 	29	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	30	000	
 Predictions of different functions 	31	000	
can <i>vary greatly</i> in their ability to be	32	000	
learned with computational efficiency	33	000	
	34		

35 🔘 🔘 🔾







	Time step	Action signals	Sensory s
 Different sensory signals can be qualitatively different from each other 	0	$\bigcirc \bigcirc \bigcirc \bigcirc$	
quantatively amerent north saor strict	1	000	
 In their range of values 	2	000	
• In their predictive relationships	3	000	
 In their predictive relationships 	4		
 to action signals 	6	000	
e to ocob othor	7	$\circ \circ \circ$	
• to each other	8	000	
 to themselves 	9	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	10	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 There are short-term and long-term 	11	$\bigcirc \bigcirc \bigcirc$	
patterns in these data	12	000	
	13	000	
 There are many things to predict 	14	000	
 Prediction need not be just of the 	15	000	
sensorv signals	16		
	18		
 The most important predictions are of 	19	000	
<i>functions</i> of future sensory signals	20	000	
• a a pradictions of value the	21	000	
• e.g., predictions of value, the	22	$\bigcirc \bigcirc \bigcirc \bigcirc$	
discounted sum of future reward	23	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 e.g., General value functions (GVFs) 	24	000	
	25	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 predict any signal, not just reward 	26	$\bigcirc \bigcirc \bigcirc \bigcirc$	
• over a flavible temperal epyelepa	27	$\bigcirc \bigcirc \bigcirc \bigcirc$	
• over a nexible temporal envelope	28	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 contingent on any policy 	29	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	30	000	
 Predictions of different functions 	31	000	
can <i>vary greatly</i> in their ability to be	32	000	
learned with computational efficiency	33	000	
	34		

35 🔘 🔘 🔾







	Time step	Action signals	Sensory s
 Different sensory signals can be qualitatively different from each other 	0	000	
	1	000	
 In their range of values 	3		
 In their predictive relationships 	4	000	
e te estien signale	5	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 to action signals 	6	000	
 to each other 	7 8	000	
	9	000	
 to themselves 	10	000	
 There are short-term and long-term 	11	$\bigcirc \bigcirc \bigcirc$	
patterns in these data	12	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	13	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 There are many things to predict 	14	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 Dradiation need not be just of the 	15	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	16	$\bigcirc \bigcirc \bigcirc \bigcirc$	
sensory signals	17	000	
 The most important predictions are of 	18	000	
<i>functions</i> of future sensory signals	19	000	
, , ,	20		
 e.g., predictions of value, the 	22		
discounted sum of future reward	23	000	
• e.a. <i>General</i> value functions (GV/Fs)	24	000	
• c.g., acheral value functions (GVFS)	25	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 predict any signal, not just reward 	26	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	27	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 over a flexible temporal envelope 	28	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 contingent on any policy 	29	$\bigcirc \bigcirc \bigcirc$	
contingent on any pency	30	000	
 Predictions of different functions 	31	000	
can vary greatly in their ability to be	32		
learned with computational efficiency	33		
	34		



L			
<u> </u>			
	_		

	Time step	Action signals	Sensory s
 Different sensory signals can be qualitatively different from each other 	0	000	
	1	000	
 In their range of values 	3		
 In their predictive relationships 	4	000	
e te estien signale	5	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 to action signals 	6	000	
 to each other 	7 8		
	9	000	
 to themselves 	10	000	
 There are short-term and long-term 	11	$\bigcirc \bigcirc \bigcirc$	
patterns in these data	12	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	13	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 There are many things to predict 	14	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 Dradiation need not be just of the 	15	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	16	$\bigcirc \bigcirc \bigcirc \bigcirc$	
sensory signals	17	000	
 The most important predictions are of 	18	000	
<i>functions</i> of future sensory signals	19	000	
, , ,	20		
 e.g., predictions of value, the 	22		
discounted sum of future reward	23	000	
• e.a. <i>General</i> value functions (GV/Fs)	24	000	
• c.g., acheral value functions (GVFS)	25	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 predict any signal, not just reward 	26	$\bigcirc \bigcirc \bigcirc \bigcirc$	
	27	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 over a flexible temporal envelope 	28	$\bigcirc \bigcirc \bigcirc \bigcirc$	
 contingent on any policy 	29	$\bigcirc \bigcirc \bigcirc$	
contingent on any pency	30	000	
 Predictions of different functions 	31	000	
can vary greatly in their ability to be	32		
learned with computational efficiency	33		
	34		



Now, let's characterize the solution of intelligence (i.e., take a peek inside the agent)

The common model of the intelligent agent

Action



Common to RL, psychology, control theory, economics, neuroscience, operations research...

The agent comprises four components:

Perception produces the state representation used by all components

Reactive Policy quickly produces an action appropriate to the state

Value Function evaluates how well it's going, changes the policy (learning)

Transition model predicts the consequences of alternate actions, changes the policy (planning)

ents ents an it's

Distinctive features of the Alberta RL approach to Al (how we keep our Eyes on the Prize)

- 1. We formulate intelligence as *real-time signal processing to maximize reward* (with limited computational resources)
- 2. Our agents learn continually from ordinary, unprepared experience (not from special training sets or demonstrations)
- 3. Our agents, though vastly complex, are much less complex than the world (the Big-World perspective) and thus must approximate, even in the limit
- 4. Our agents continually build abstractions, particularly time abstractions (options, GVFs) that enable them to think simply about large spans of time

"The Alberta Plan for Al Research"

A 12-step plan to reach human-level AI ahead of schedule (before 2040)

- 1. Representation I: Continual supervised learning with given features
- 2. Representation II: Supervised feature finding
- 3. Prediction I: Continual GVF prediction learning
- 4. Control I: Continual actor-critic control
- 5. Prediction II: Average-reward GVF learning
- 6. Control II: Continuing control problems
- 7. Planning I: Planning with average reward
- 8. Prototype-AII: One-step model-based RL with continual function approximation
- 9. Planning II: Search control and exploration
- 10.Prototype-AI II: The STOMP progression
- 11. Prototype-AI III: The Oak architecture
- 12. Prototype Intelligence amplification

A strategy document by Mike Bowling, Patrick Pilarski, and myself, coming soon to arXiv.



Outline

- The Prize of understanding intelligence
- Why we seek it
- on the Prize
 - by focusing on experience

"To make AI for good and for all"

How reinforcement learning research in Alberta keeps its eyes

Thank you for your attention