

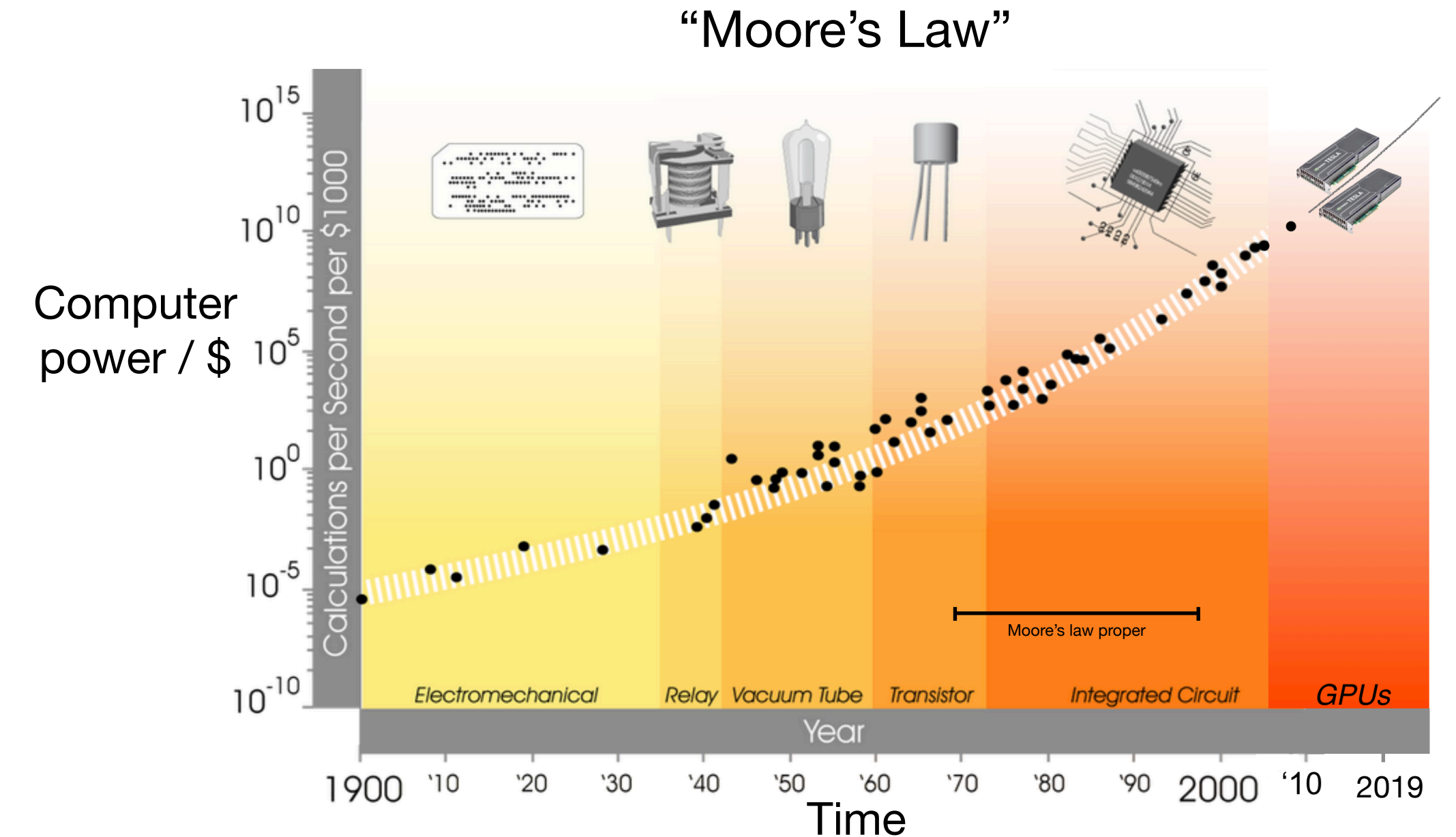
# Are You Ready to Fully Embrace Approximation?

Rich Sutton

DeepMind, Amii, RLAI, and UAlberta



# Machine Intelligence Today



- Increasing computational power (Moore’s Law) drives progress
- Methods that *scale with computation* are the most impactful
- Thus the current successes of machine learning and deep learning

# Scaling with computation is new. Not the usual in CS

- **The usual scaling** is scaling with problem size
  - bigger problem  $\Rightarrow$  more computation needed *to solve it exactly*
- Now we assume the problem could *never* be solved exactly
- **The new scaling** is scaling with computation
  - more computation  $\Rightarrow$  *a better approximate answer*

**We need methods that scale with increasing computation**

We need methods that scale with increasing computation

Search and Learning.

We need methods that scale with increasing computation  
to better approximate answers.

Search and Learning.

We need methods that scale with increasing computation  
to better approximate answers.

Search and Learning. With approximation.

# RL has scaled with computation pretty well

- It has embraced function approximation.
- It has embraced Deep Learning.
- It has embraced learning from unprepared experience.
- It has embraced search, particularly MCTS.
- It has embraced replay and (to some extent) planning.
- All these things scale with computational resources



# But RL has held back.

## It has not fully embraced approximation

- RL is grounded in finite MDPs and tabular methods
- To really abandon finite MDPs challenges us psychologically, requires strong discipline
- If we fully embraced approximation we would lose so much!
  - We lose discounted reward and all the theory built on it
  - We lose Bellman Errors
  - We lose Markov state, thus transition probabilities and expectations, including all true value functions  $v_\pi, v_*, q_\pi, q_*$

# How has RL dealt with the loss?

“The five stages of grieving”

Denial

Anger

Bargaining

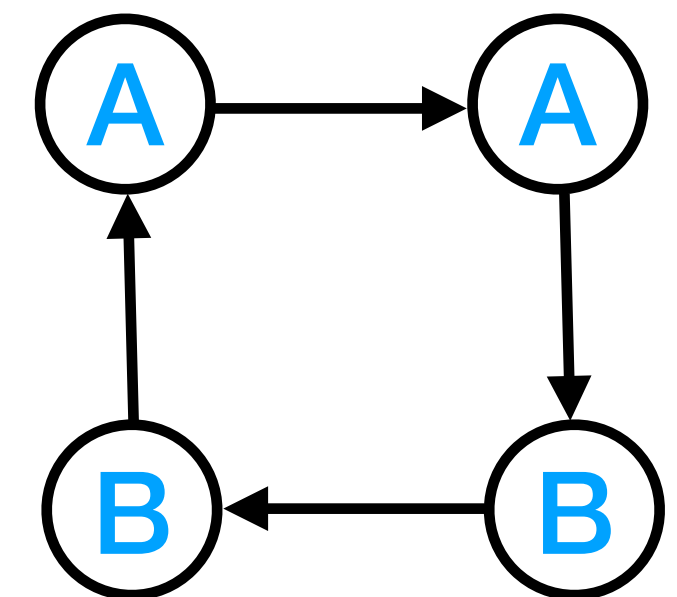
Depression

Acceptance

# Approximation in Reinforcement Learning

- World (environment) states map to feature vectors  $\phi_t = \phi(S_t) \in \mathfrak{R}^d, d \ll |\mathcal{S}|$
- Then all agent operations use only the feature vectors  $\phi_t$
- Thus, we may talk about a value function  $\hat{v}_{\mathbf{w}}(s)$ , but really it is  $s \rightarrow \phi \rightarrow \hat{v}$
- Note  $\phi_t$  is **not Markov**;  
what happens next will depend on past feature vectors (and actions)
- e.g.,  $\Pr[\phi_{t+1} = \phi' \mid \phi_t = \phi]$  is not defined

AABBAAABB



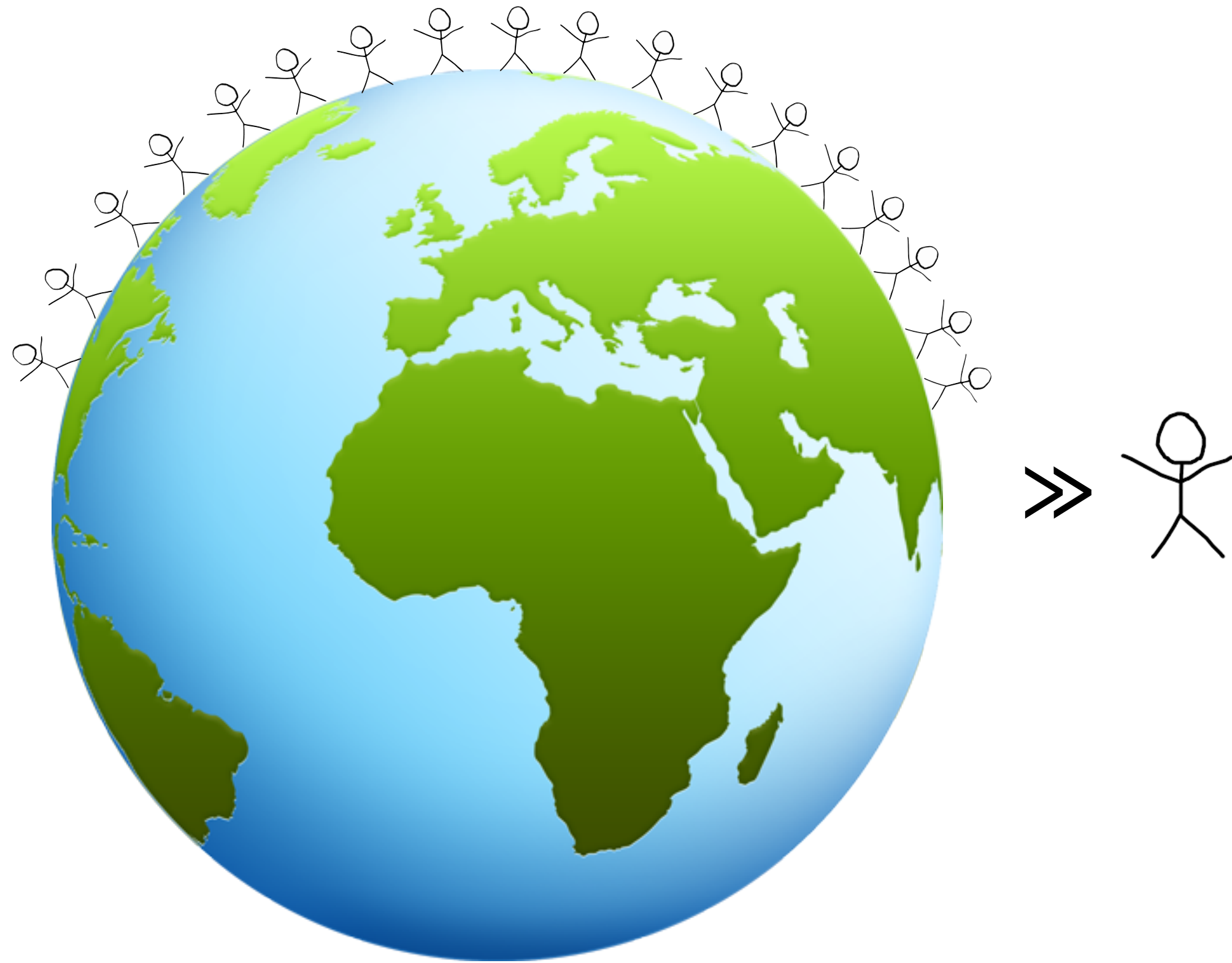
# Approximation in Reinforcement Learning (2)

- World states map to feature vectors  $\phi_t = \phi(S_t) \in \mathcal{R}^d, d \ll |\mathcal{S}|$
- Note that there may be as many as  $|\mathcal{S}|$  different feature vectors
- Thus the feature vectors cannot be treated as individuals in any way (they must be processed parametrically)
- e.g., we couldn't approximate  $\Pr[\phi_{t+1} = \phi' \mid \phi_t = \phi]$  (even if it made sense) because you would have to store things for each  $\phi$
- and it would depend on the behavior policy

# Fully embracing approximation means

- the agent can't store things for individual states
- the agent can't do anything that treats individual feature vectors distinctly
- the state the agent works with will not be Markov
- never converging to the exactly correct anything, even in the limit
- the world is much bigger (more complex) than the agent
  - even as the agent's computational complexity grows exponentially!
- experience is too big to be fully processed by the agent, particularly in real time
- the best approximations will change over time, thus learning must be online

# The world is much more complicated than you



- Thus, approximation must be embraced.
- Anything you try to learn can only be learned approximately:
  - value functions,
  - policies,
  - models,
  - states.
- Violating this principle is the most important problem with the use of simulated worlds.

Big world  $\Rightarrow$  apparent non-stationarity  
 $\Rightarrow$  changing *approximate* value function




## Acceptance and opportunity (1):

# Function approximation when there is no ideal

- Approximation is okay, we can still do things.  
It's just different. Probably better, certainly real-er.
- Transition probabilities and expectations  
are replaced by a **function approximator with a loss**
- There are not usable “true” value functions
  - but we can have approximations with a loss
  - and we still do have **mean squared return error** (for a fixed policy):

$$\text{MSRE}(\mathbf{w}) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \left[ \hat{v}(S_{t+k}, \mathbf{w}) - G_{t+k} \right]^2, \text{ if } A_i \text{ were selected } \sim \pi, \forall i \geq t$$



weight vector      approx value of state      return      action      policy

## Acceptance and opportunity (2):

# Discounting $\Rightarrow$ Maximize average reward rate

- All policies  $\pi$  are ranked according to their reward rate:

$$r(\pi) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \underset{\text{reward}}{\uparrow} R_{t+k}, \text{ if } \underset{\text{action}}{\uparrow} A_i \sim \pi, \forall i \geq t$$

- Returns are defined relative to  $r(\pi)$ :

$$G_t \doteq R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \dots$$

- Learning and planning algorithms are less developed,  
but Yi Wan and Abhishek Naik have just made good progress (NeurIPS)



Acceptance and opportunity (3):

# Feature function $\Rightarrow$ state-update function

- Instead of an unknowable function  $\phi$  accessing an unknowable world state
- We have a *known* state-update function, operating on *known* experience, with a *known*, improvable objective (summarizing the past to predict the future):

$$S_t = u(S_{t-1}, A_{t-1}, O_t)$$

stateactionobservation

state-update function

The diagram illustrates the state-update function  $u$ . It shows the equation  $S_t = u(S_{t-1}, A_{t-1}, O_t)$ . Below the variables, blue arrows point upwards to the function  $u$ . The arrow from  $S_{t-1}$  is labeled 'state', the arrow from  $A_{t-1}$  is labeled 'action', and the arrow from  $O_t$  is labeled 'observation'. Below the arrow from  $S_{t-1}$ , the text 'state-update function' is written.

- This is just a better way to get a non-Markov state
- Our Agent-State Research Group is working on this

# Acceptance and opportunity (4): Converging $\Rightarrow$ tracking

- Approximation means accepting that the world is big, you can't get anything exactly right
- You could **converge** to the best approximate static solution, balancing all the errors, or you could **track** the current best approximation
- Surprisingly, you can *do better by tracking*, maybe *much better* (see ICML2007 paper by Dave Silver, Anna Koop, and me)
- Tracking means learning and relearning, continually, online, like an endless sequence of related learning problems, but **all from one base problem**
- Thus approximation provides a new basis, a new rationale, for on-line learning, meta learning, generalization, and representation learning!

# Conclusion

- Approximation is key to future advances in machine intelligence
- As the premiere RL research institution, we should be leading the advances in approximation within RL
- Approximation seems a difficult challenge, but it is necessary,
  - and will yield great dividends if we fully embrace it
- Fully embracing approximation is on the critical path to the future of machine intelligence