

# Towards Risk-Aware Decision-Making and Policy Evaluation for Actionable Healthcare

Sonali Parbhoo

Imperial College London

OxML Summer School

7 Aug 2022

# Machine Learning is changing the way we make decisions



# But machine learning has had limited success in high-risk settings like healthcare

## The "inconvenient truth" about AI in healthcare

Trishan Panch, Heather Mattie & Leo Anthony Celi 

*npj Digital Medicine* 2, Article number: 77 (2019) | [Cite this article](#)

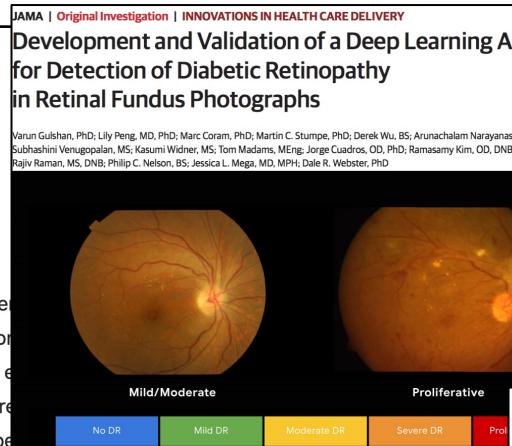
35k Accesses | 24 Citations | 452 Altmetric | [Metrics](#)

As the UK sits in painful deadlock over Brexit, it is important to remember that governments are regularly faced with crises, and their responses can create enduring benefit for future generations. Back in 1858, for example, the UK parliament was dealing with another messy crisis: "the greatest cholera outbreak in history". In a world before sanitation, the river Thames had become an open sewer, and as summer blossomed parliament stench. £2.5 million (about £300 million today) was approved to build a network of sewers. A particular model of sanitation, developed in London, was adopted by other cities around the world and the world has not been able to imagine that a developed nation would

## Why AI fell short in slowing the spread of COVID-19

Though hopes were high, experts say a lack of high-quality data meant models were both underwhelming and "anti-constructive."

By Kat Jercich | July 23, 2020 | 04:35 PM



## AI in healthcare - not so fast? Study outlines challenges, dangers for machine learning

As machine learning rapidly expands into healthcare, the ways it "learns" may be at odds with clinical outcomes unless carefully controlled for, a new study shows.

By Benjamin Harris | January 16, 2019 | 10:21 AM

MEDICAL & BIOTECH

## Artificial Intelligence Is Rushing Into Patient Care—And Could Raise Risks

AI systems are not as rigorously tested as other medical devices, and have already made serious mistakes

By Liz Szabo, Kaiser Health News on December 24, 2019

## Artificial Intelligence Makes Bad Medicine Even Worse

A new study out from Google seems to show the promise of AI-assisted health care. Actually, it shows the threat.

# Why?

High stake decision-making is challenging *even for humans*



What treatment should a patient take?



How should a teacher assist a student?

# Why?

High stake decision-making is challenging *even for humans*



What treatment should a patient take?



How should a teacher assist a student?

Only 6% of asthmatics would be eligible for treatments from RCTs<sup>1</sup>



## Does a model generalise and to whom will it generalise?

<sup>1</sup>Smith, Saunders et al. Best Care at Lower Cost: The Path to Continuously Learning Healthcare in America, 2013.

# Decisions often have *long-term* consequences

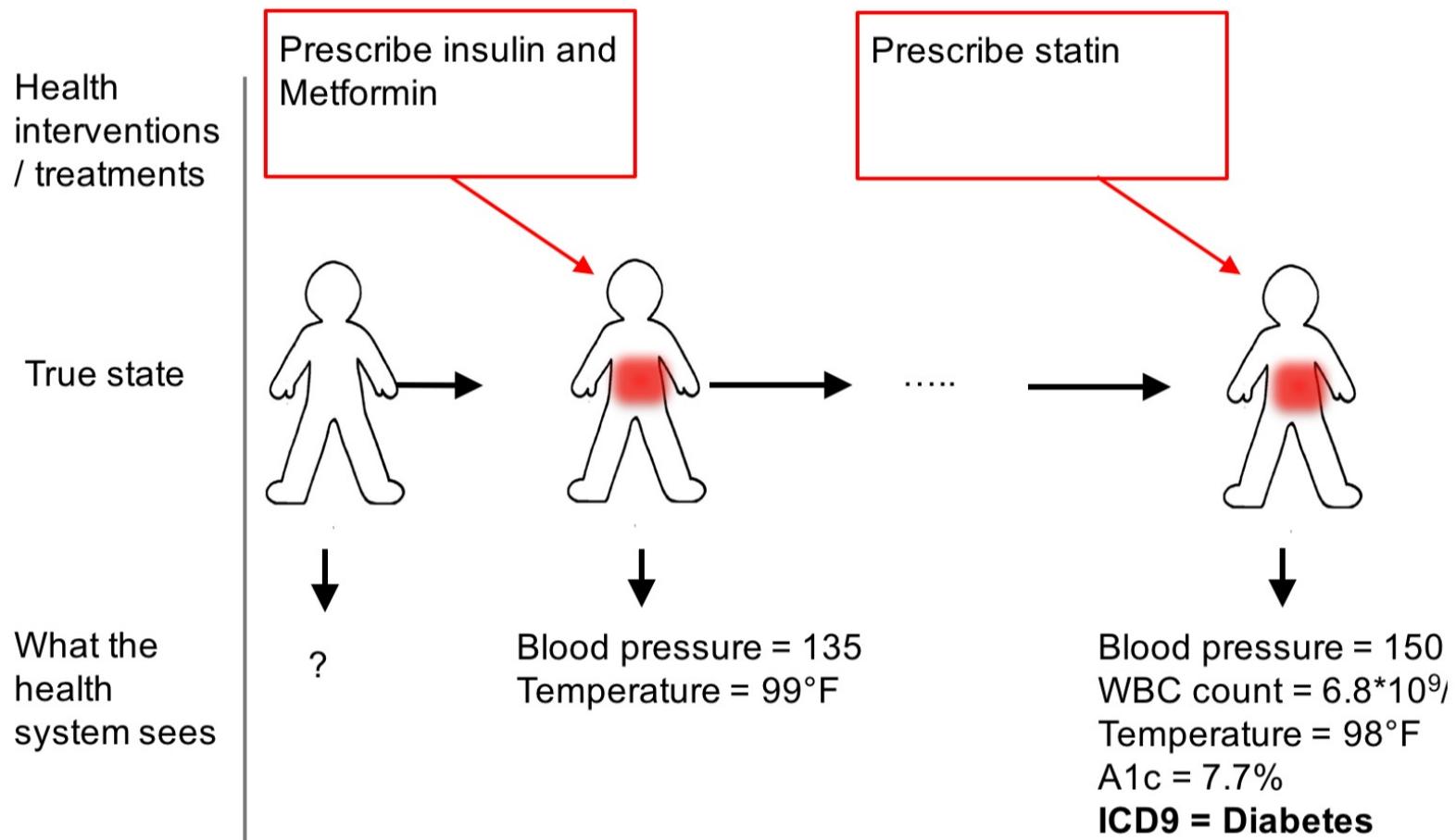


Image credit: David Sontag – ML for Healthcare, 6.871

# Decisions often have *long-term* consequences

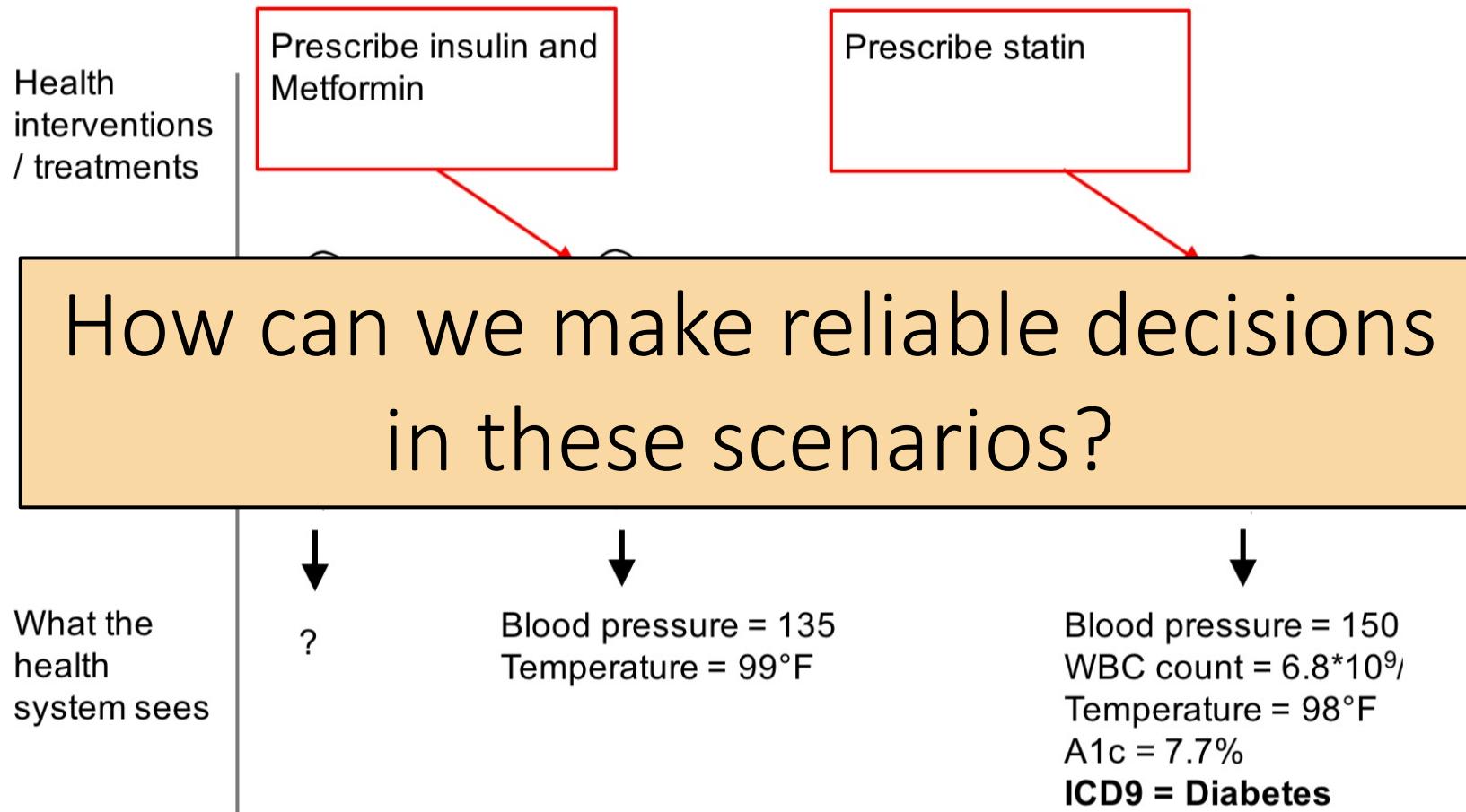
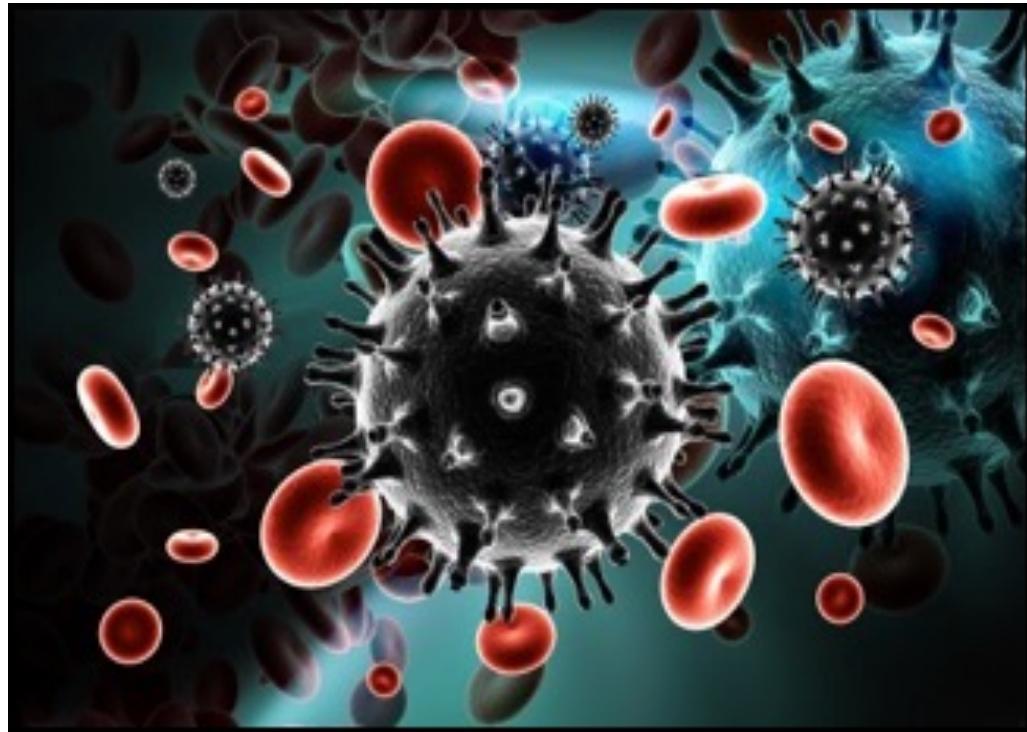


Image credit: David Sontag – ML for Healthcare, 6.871

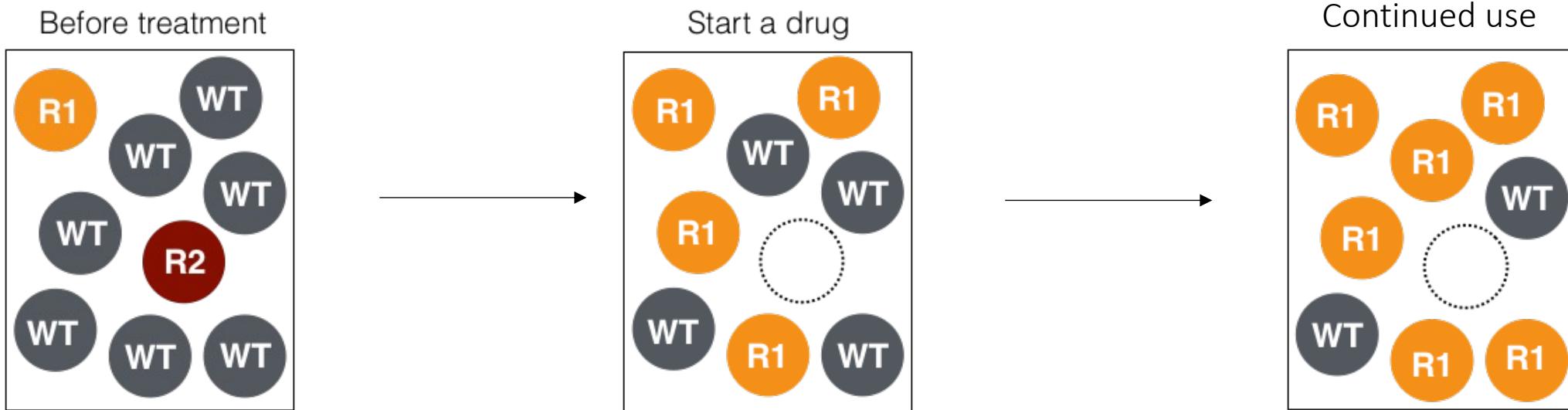
# Example: HIV Therapy Management



- Requires life-long treatment with **combinations of antiretrovirals**.
- Virus is highly mutagenic and can become drug resistant.
- Patients exhibit significant **heterogeneity**.

# Example: HIV Therapy Management

- Treatments consist of five classes of drugs: NNRTIs, NRTIs, PIs, IIs, FIs
- The virus mutates in response to treatments leading to **drug resistance**.

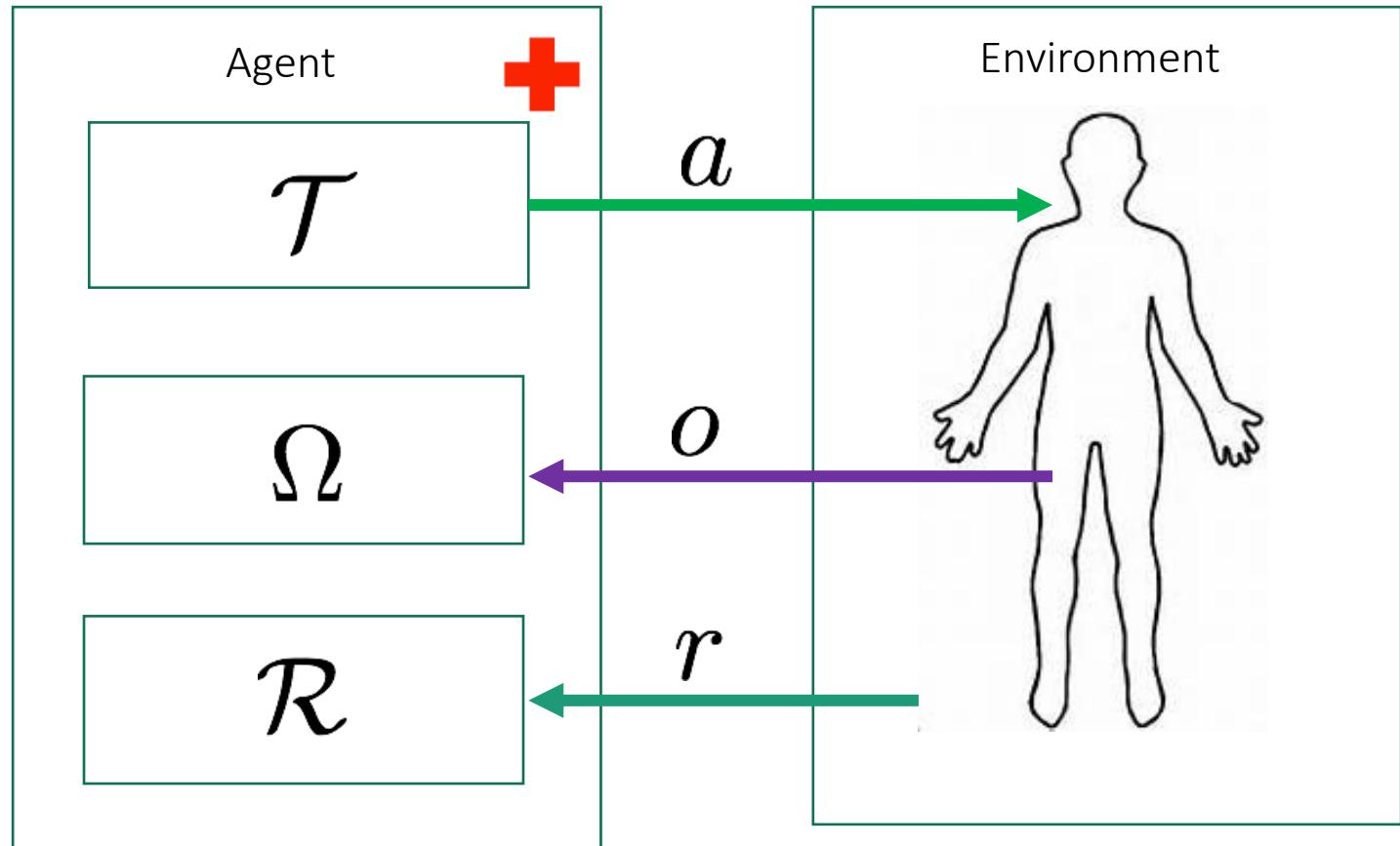


# Reinforcement Learning for Healthcare

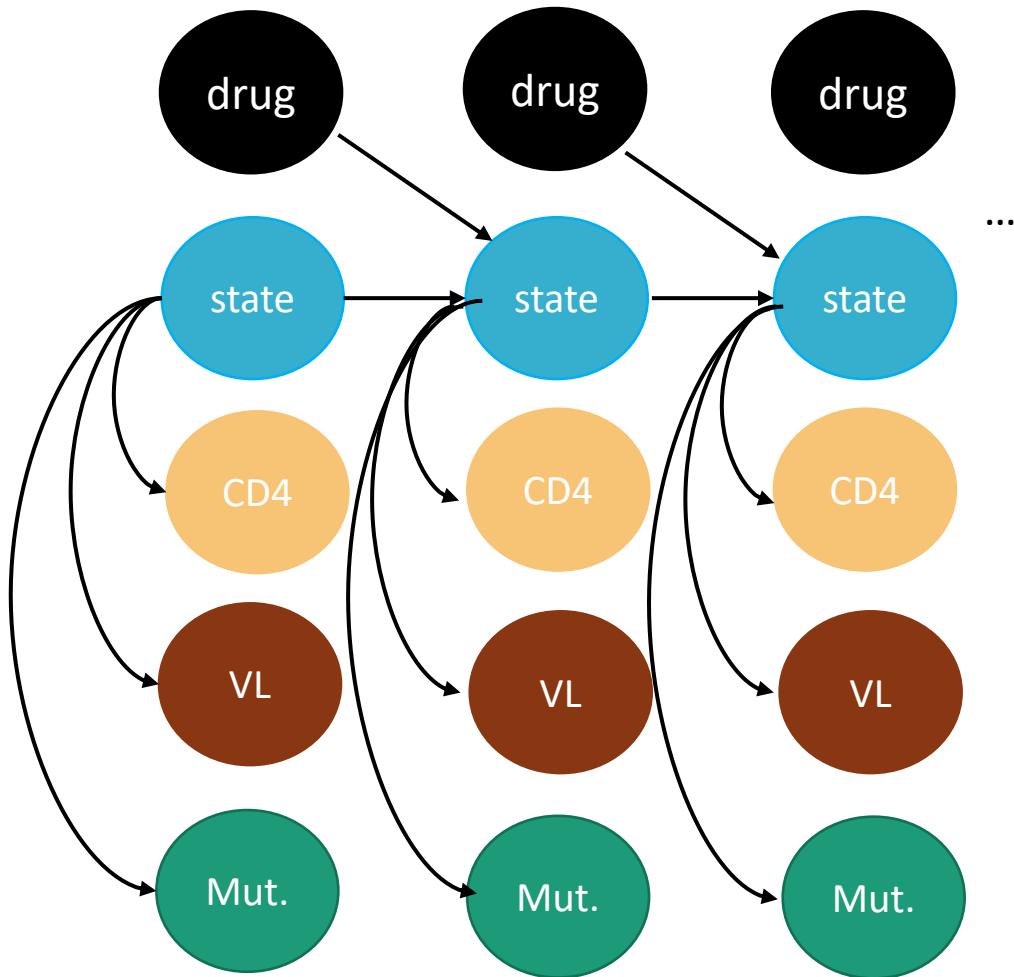
Medicine is inherently a sequential decision making problem:

- Clinicians, with their best understanding of a patient's status, propose a treatment.
- Patient status may change as a result of the treatment.
- Outcomes are a result of prior treatments.

RL makes this process explicit



# Solutions: Train a model/learn a value function



Modelling HIV as an RL problem solves the long-term problem (e.g. Ernst 2007; Parbhoo 2014; Marivate 2015), but often in **simplified/simulation settings only**.

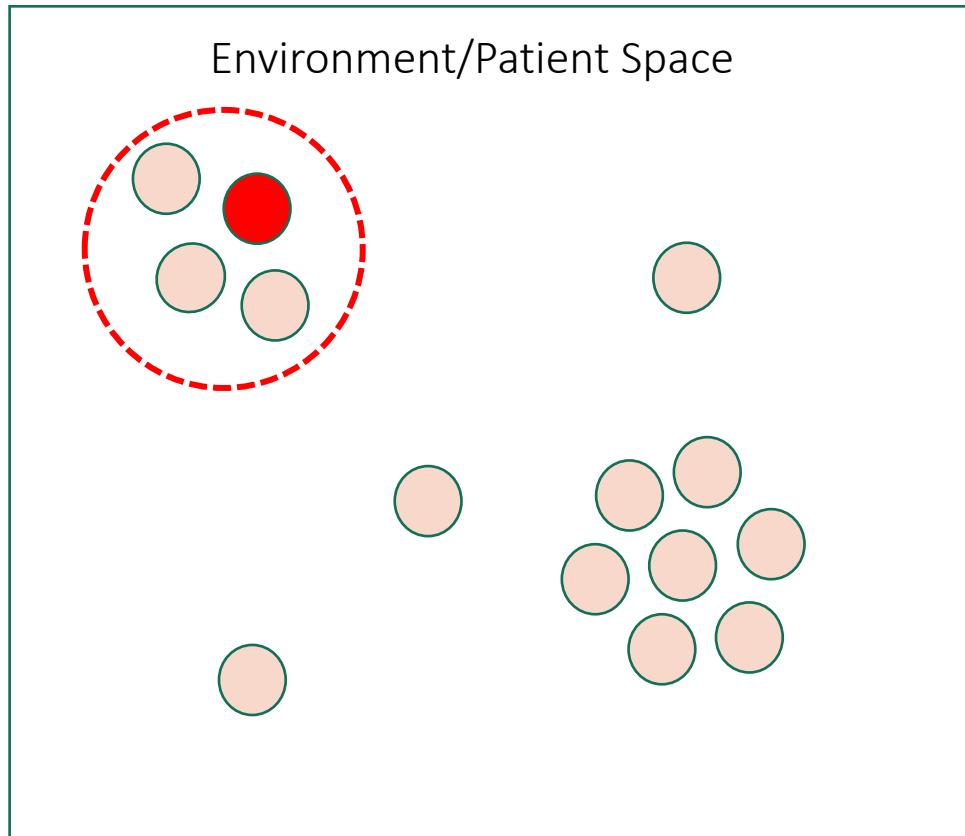
if  $V_t > 40$  :

$$r_t = -0.7 \log V_t + 0.6 \log T_t - 0.2|M_t|$$

else:

$$r_t = 5 + 0.6 \log T_t - 0.2|M_t|$$

# Solutions: Use Non-parametric Methods

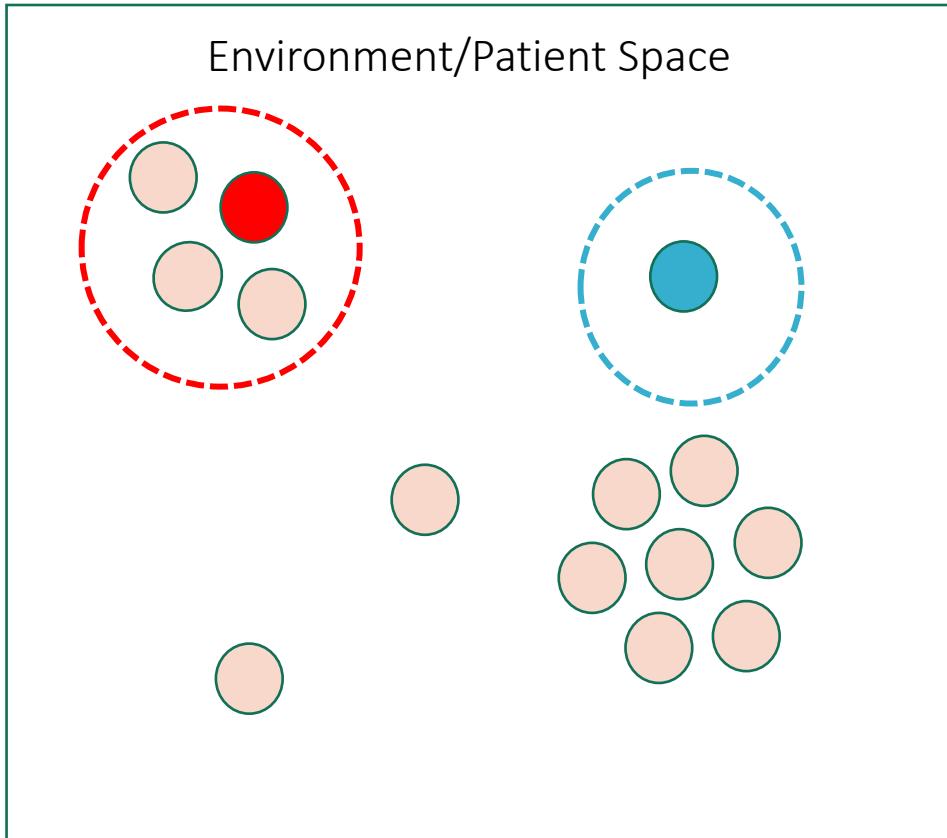


- Look at **similar patients** and use their full history to predict immediate outcomes (e.g. Bogojeska et al 2012)
- These methods however, ignore long-term effects:

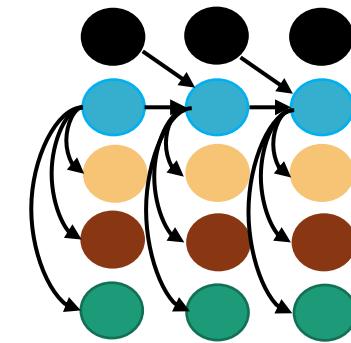
$$f(\text{history}, \text{drugs}) = l(VL < 40) \quad \text{in 21 days}$$

# Our insight: These approaches have complementary strengths

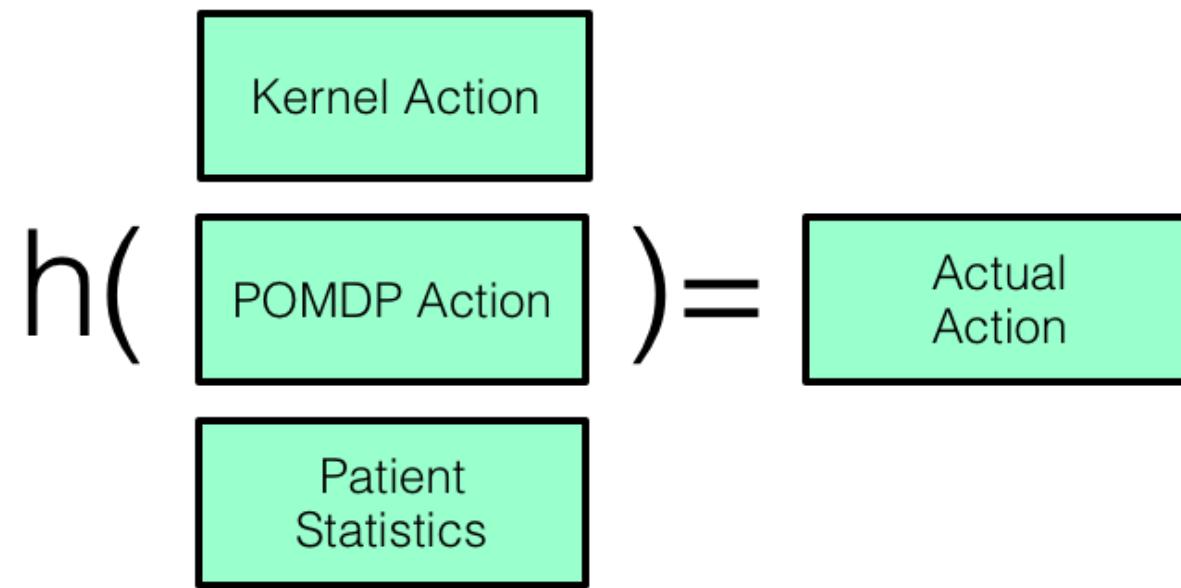
Patients that are similar  
are best modelled by  
their neighbours



Patients without  
neighbours are best  
modelled with a  
parametric model



# New Solution: Combine nonparametric and parametric models for better generalisation



We can ensemble the solutions as a Mixture-Of-Experts to learn the appropriate action choice for each patient.

# Application to HIV Therapy Management

- Cohort: 32960 patients from EuResist Database
- Observations: CD4s, VLs, mutations, clinical variables
- Actions: 312 most frequently occurring drug combinations (from 25 drugs)

Baseline	Value
Random Policy	-7.31 +/- 3.72
Neighbour Policy	9.35 +/- 2.61
Model Policy	3.37 +/- 2.15
<b>Policy-Mixture Policy</b>	<b>11.52 +/- 1.31</b>
<b>Model-Mixture Policy</b>	<b>12.47 +/- 1.38</b>

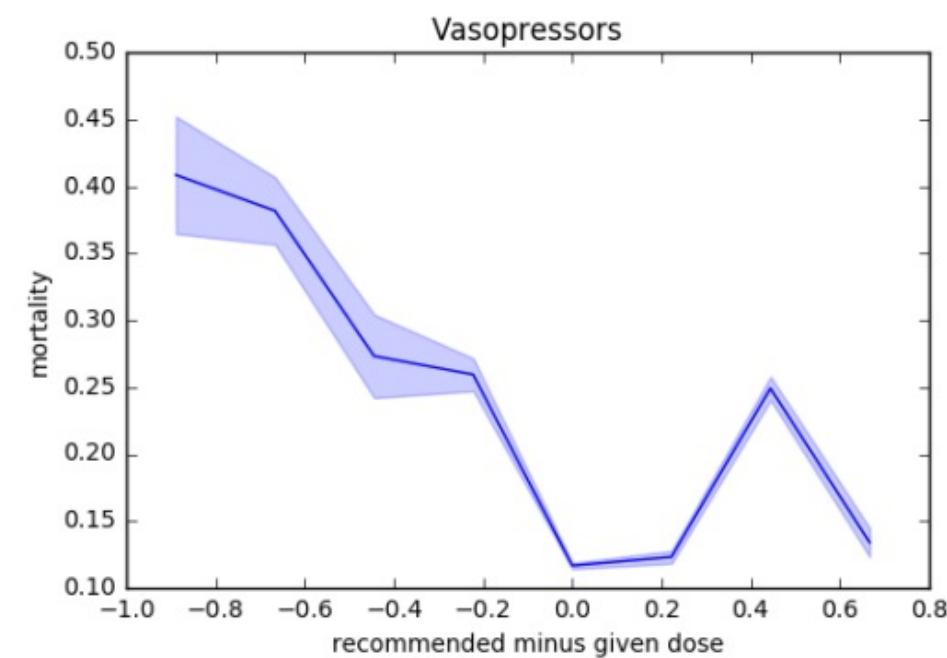
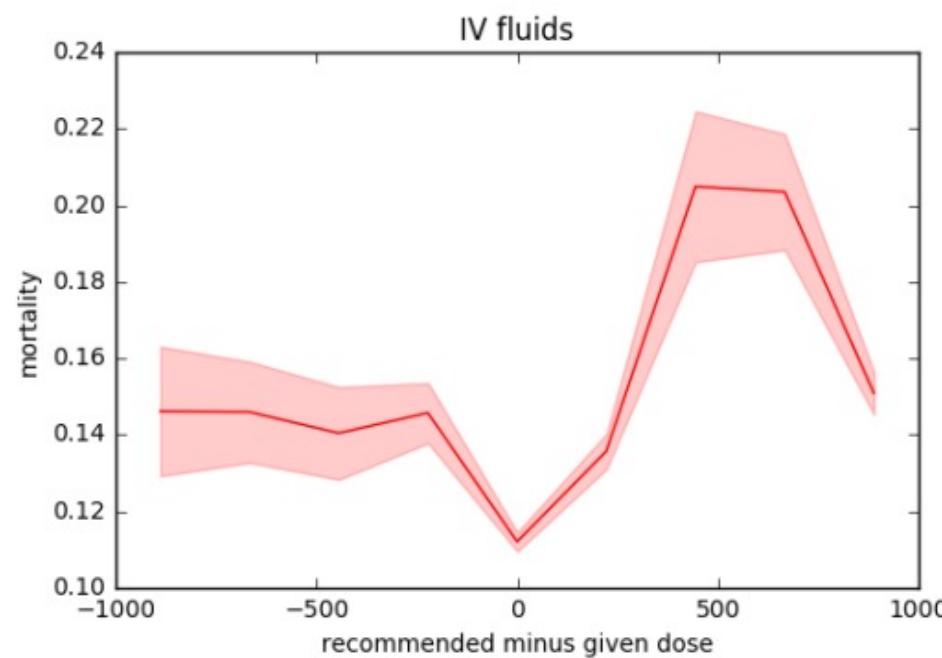
# Extension: Transfer between EuResist and South African Cohorts

- We also identified when we could transfer policies in well-curated cohorts such as EuResist to South African cohorts that weren't as well curated.

Type	Method	DR	IS	WIS
Behaviour Policy	$5.02 \pm 1.18$			
Local	Kernel	$3.56 \pm 1.42$	$1.27 \pm 1.14$	$1.80 \pm 1.07$
	CEIB	$3.29 \pm 1.13$	$3.80 \pm 2.41$	$3.76 \pm 2.19$
Transfer	Kernel	$4.17 \pm 1.4$	$4.18 \pm 1.20$	$4.16 \pm 1.71$
	CEIB	$6.29 \pm 0.14$	$5.17 \pm 0.38$	$5.27 \pm 0.29$
	Mixture-of-Experts	$5.28 \pm 0.37$	$3.42 \pm 1.39$	$4.81 \pm 1.25$
<b>Local + Transfer</b>	<b>Ours</b>	<b><math>8.96 \pm 0.39</math></b>	<b><math>10.64 \pm 1.2</math></b>	<b><math>10.62 \pm 1.67</math></b>

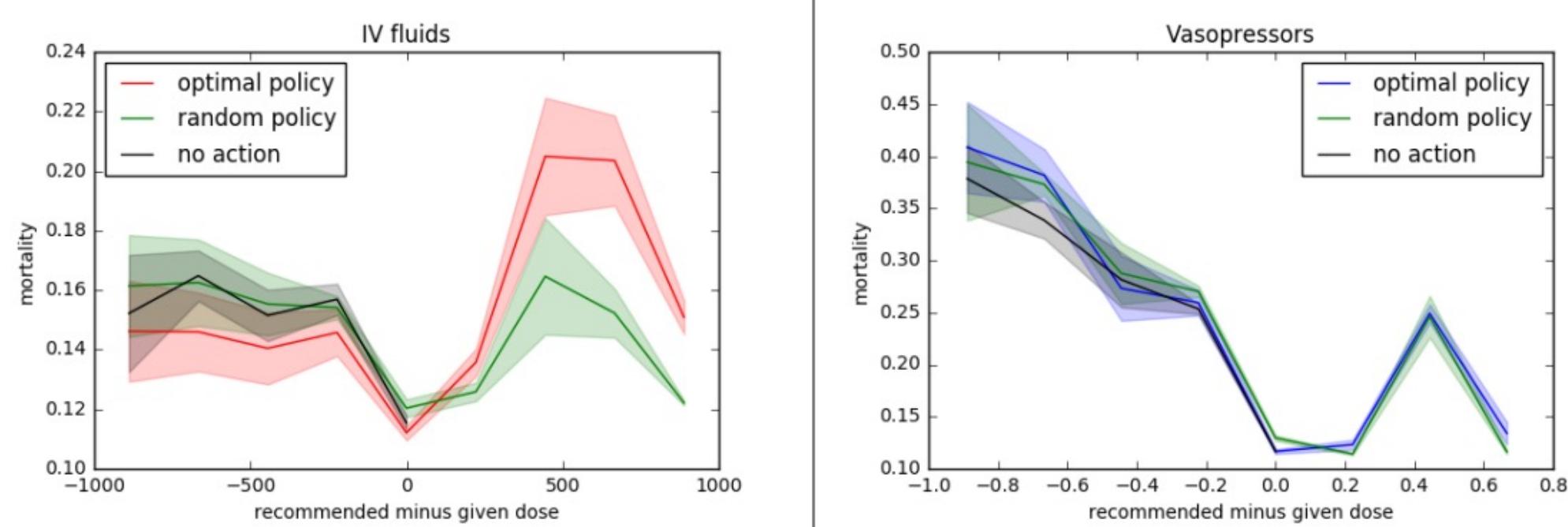
# Similar mixtures have been used for sepsis management

	Physician	Kernel	DQN	$MoE_{V_d, Q_d}$	$MoE_{V_b, Q_b}$
non-recurrent encoded	3.76	3.73	4.06	3.93	4.31
recurrent encoded	3.76	4.46	4.23	5.03	5.72

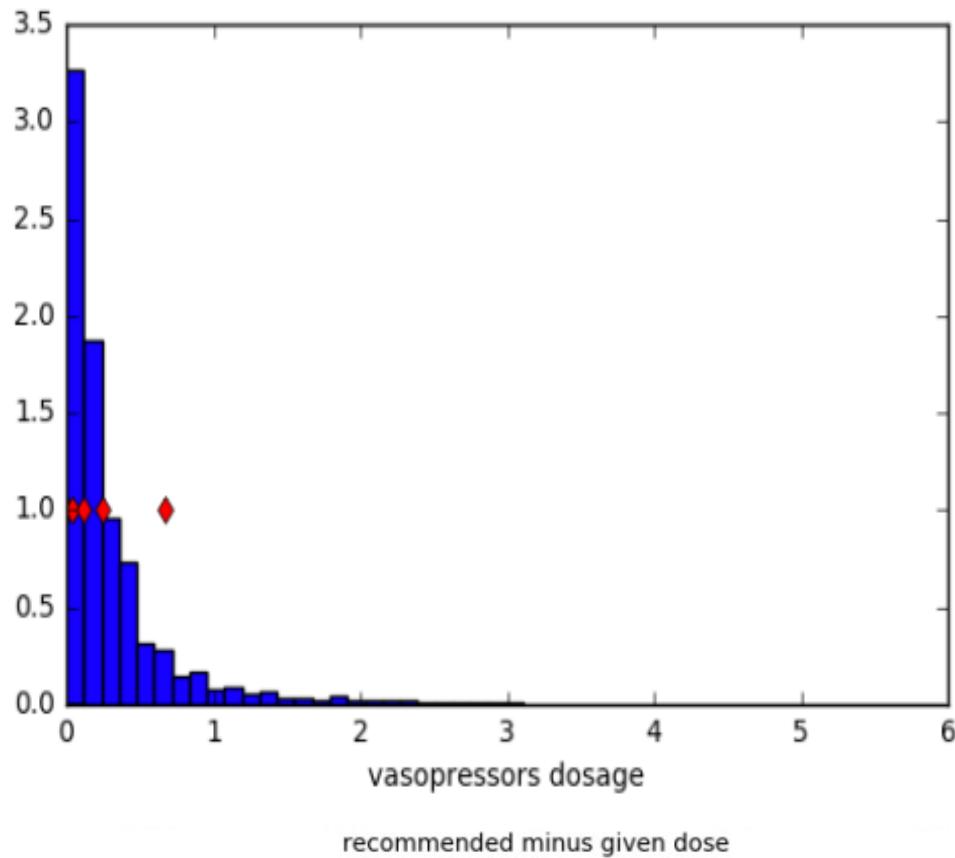


# .... but evaluation is brittle

	Physician	Kernel	DQN	$MoE_{V_d, Q_d}$	$MoE_{V_b, Q_b}$
non-recurrent encoded	3.76	3.73	4.06	3.93	4.31
recurrent encoded	3.76	4.46	4.23	5.03	5.72

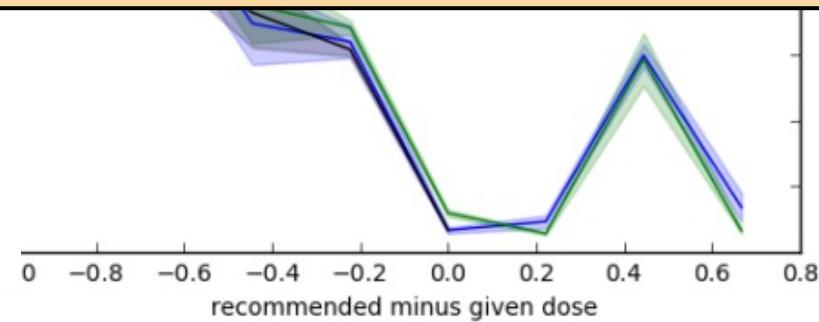


## .... but evaluation is brittle



QN	$MoE_{V_d, Q_d}$	$MoE_{V_b, Q_b}$
.06	3.93	4.31

This occurs because of  
the spread of the  
treatments  
administered



# Off-Policy Evaluation

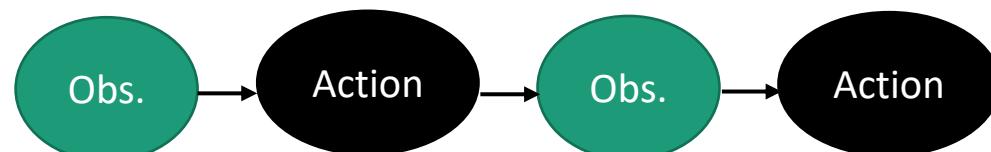
Given data collected by some **behaviour policy**  $\pi_b$ , can we estimate the value of **another evaluation policy**  $\pi_e$ ?

Clinician:

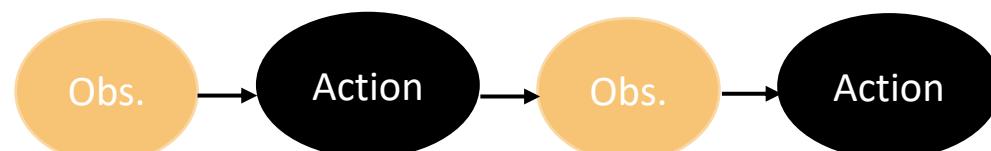


$\pi_b$

Patient 1:

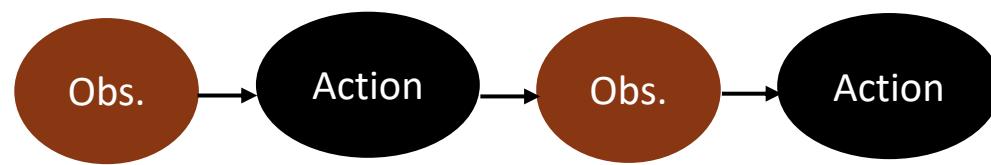


Patient 2:



....

Patient N:



Evaluation Agent:



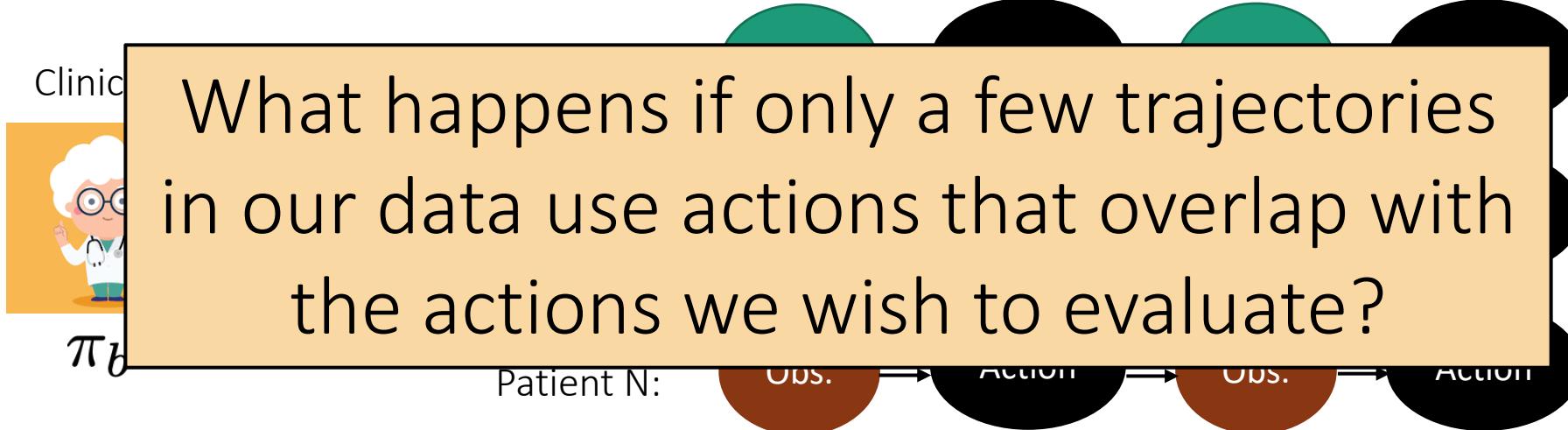
$\pi_e$

Reweight trajectories according to how much they overlap using:

$$\rho_n = \prod_t \frac{\pi_e(a_{tn} | s_{tn})}{\pi_b(a_{tn} | s_{tn})}$$

# Off-Policy Evaluation

Given data collected by some **behaviour policy**  $\pi_b$ , can we estimate the value of **another evaluation policy**  $\pi_e$ ?



Evaluation Agent:



$\pi_e$

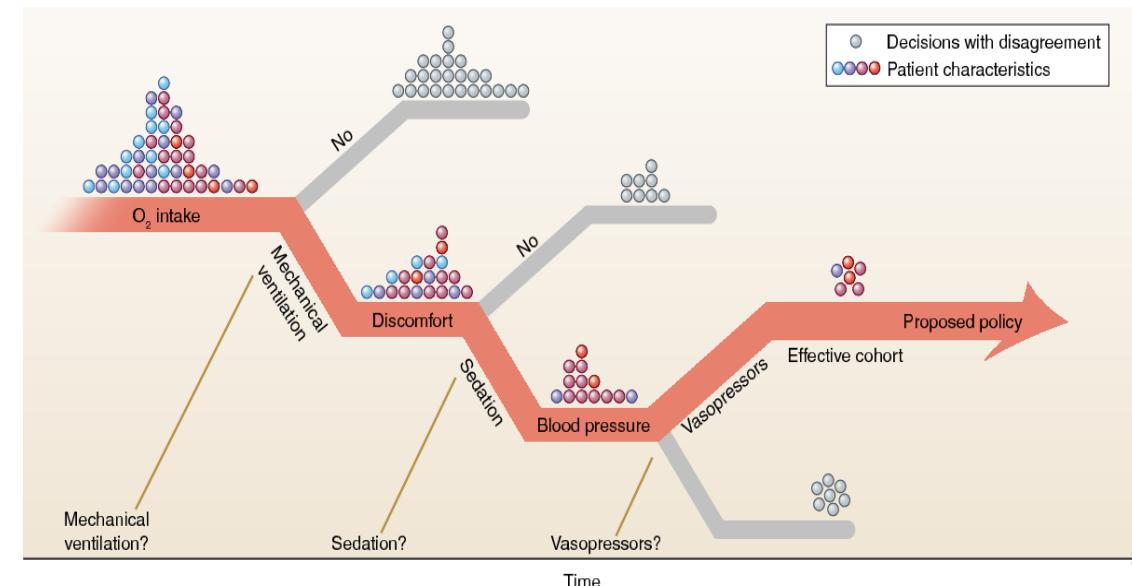
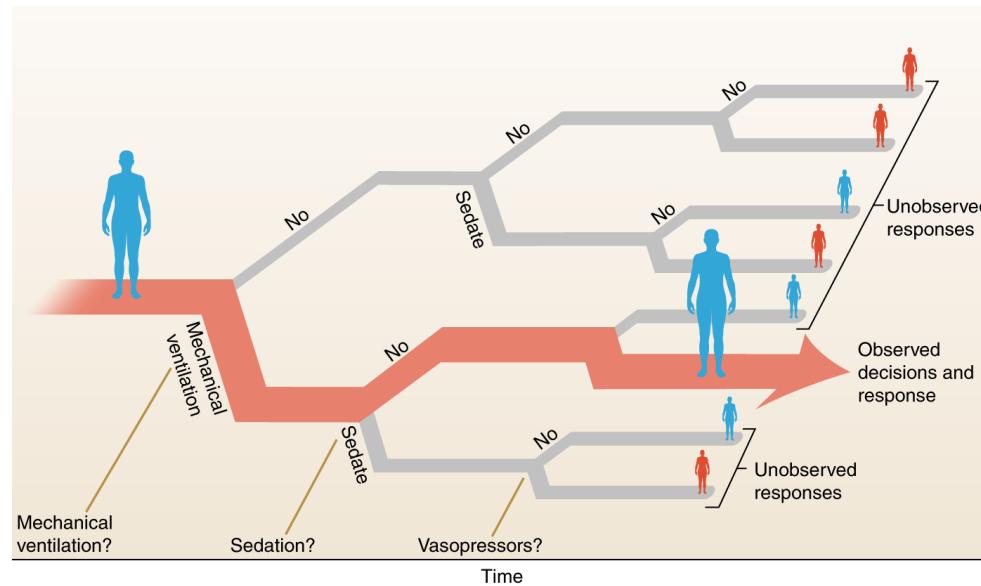
Reweight trajectories according to how much they overlap using:

$$\rho_n = \prod_t \frac{\pi_e(a_{tn} | s_{tn})}{\pi_b(a_{tn} | s_{tn})}$$

# Challenges to using RL for healthcare in practice

Learning suitable treatment policies from observational data is **offline and off-policy**

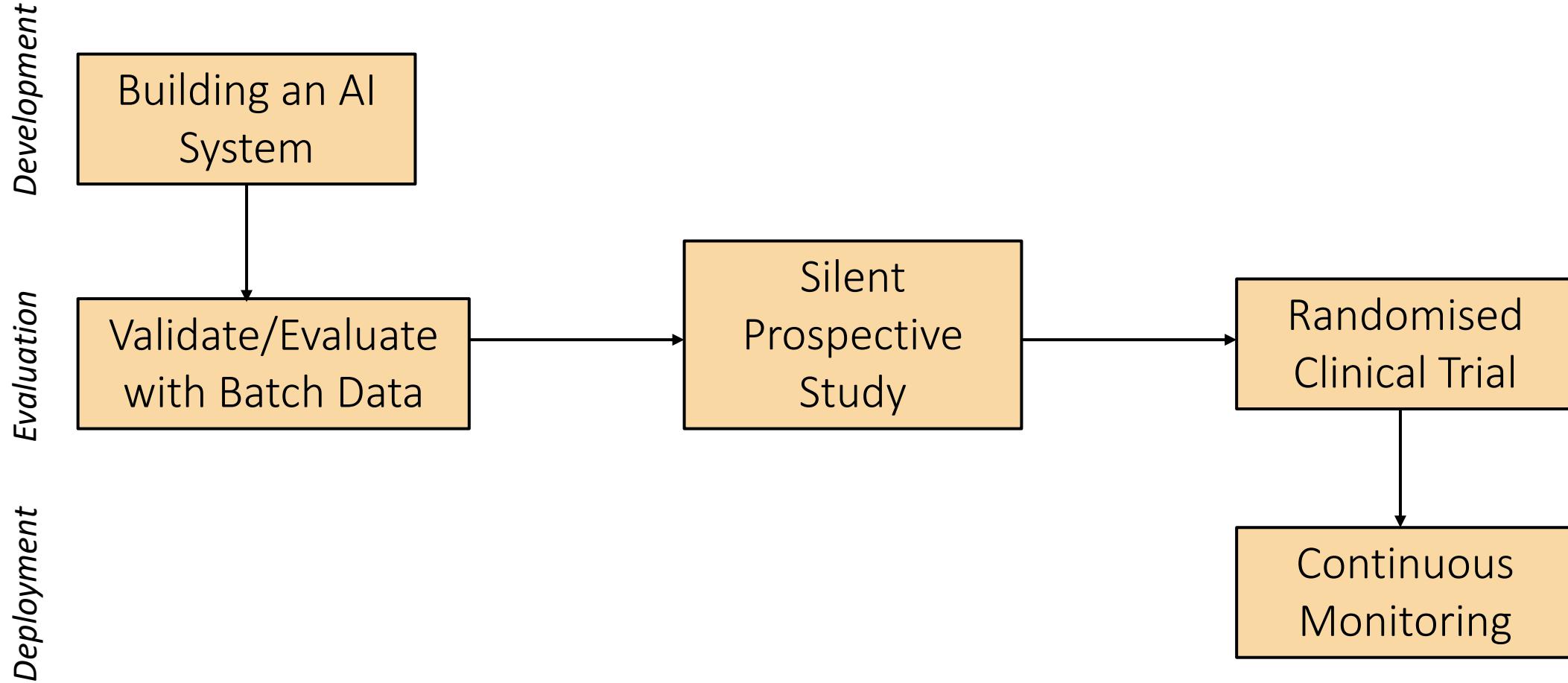
- 1.Inability to explore
- 2.Small data (relatively) → shrinking data support as best strategies are discovered



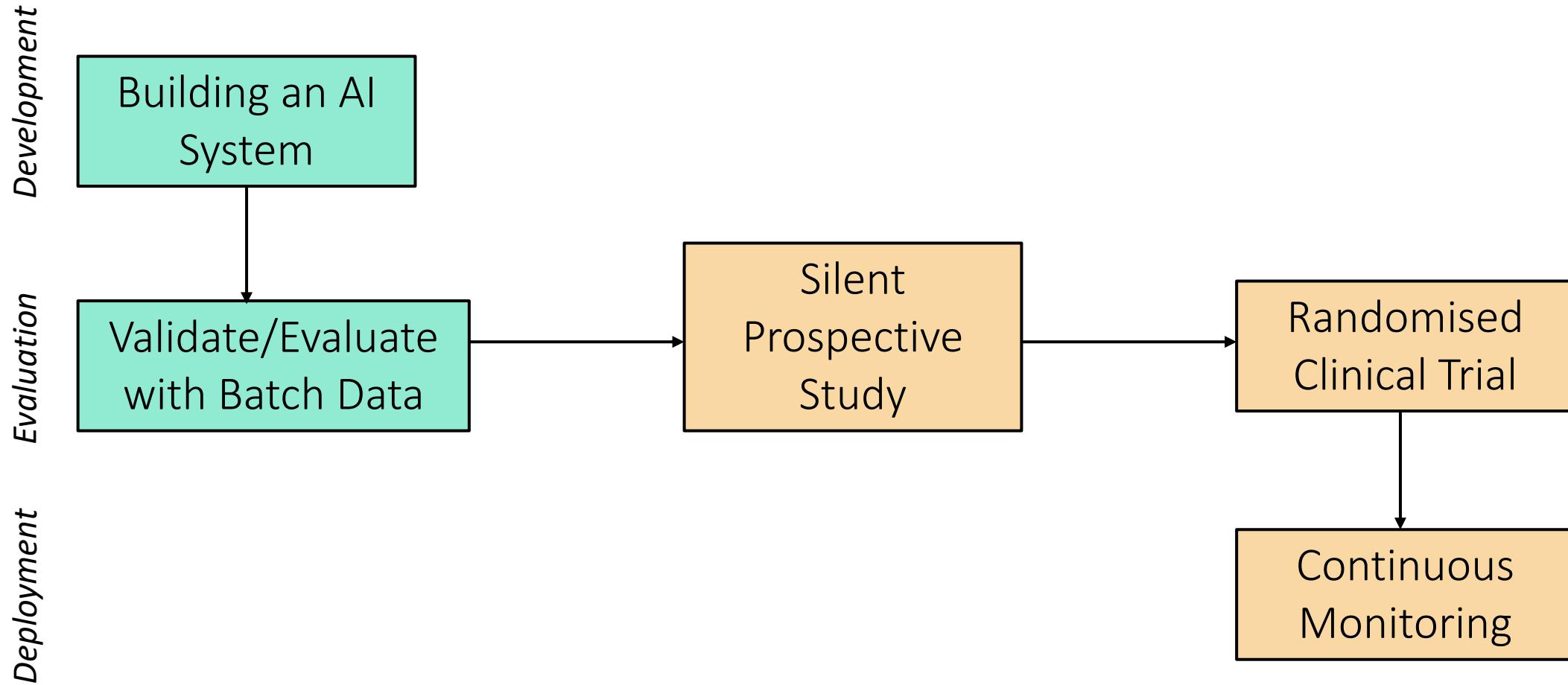
**These two limitations severely complicate the development of actionable RL**

Image Credit: Gottesman, Omer, et al. "Guidelines for reinforcement learning in healthcare." *Nat Med* 25.1 (2019): 16-18.

# We need methods for responsible model development, evaluation and deployment



# We need methods for responsible model development, evaluation and deployment



# What should responsible development and evaluation look like?

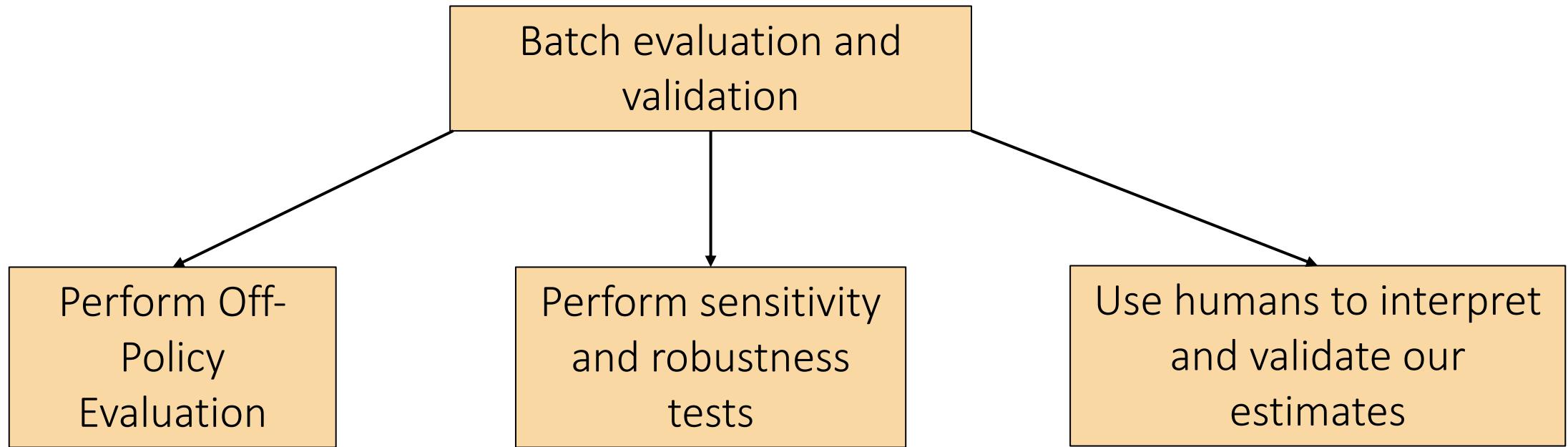


Image Credit: Doshi-Velez, Towards using Batch RL to Identifying Treatment Options, 2022

# What should responsible development and evaluation look like?

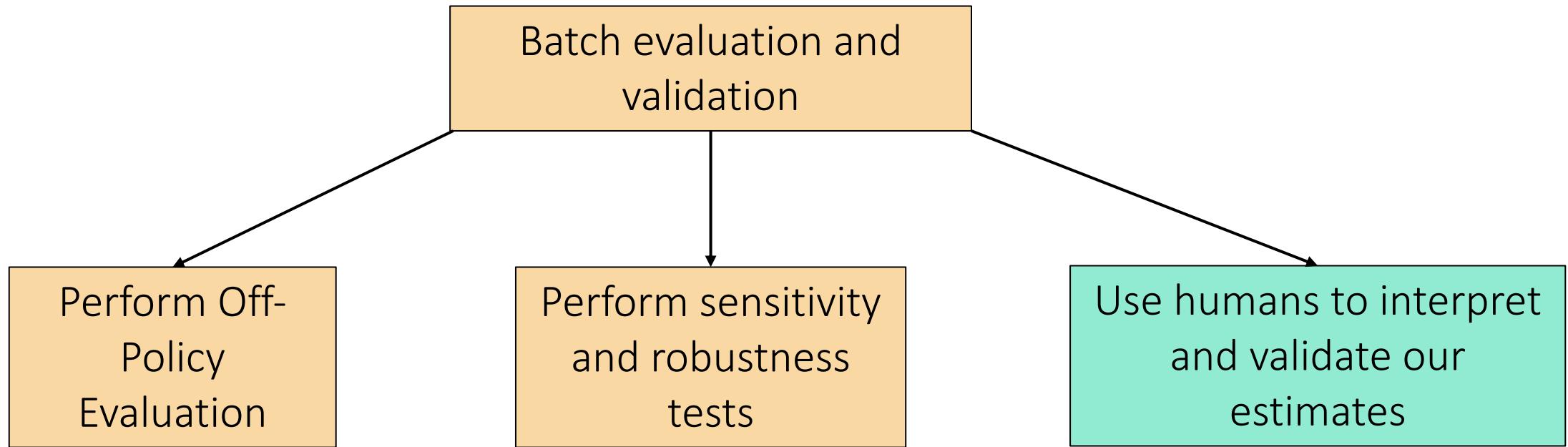
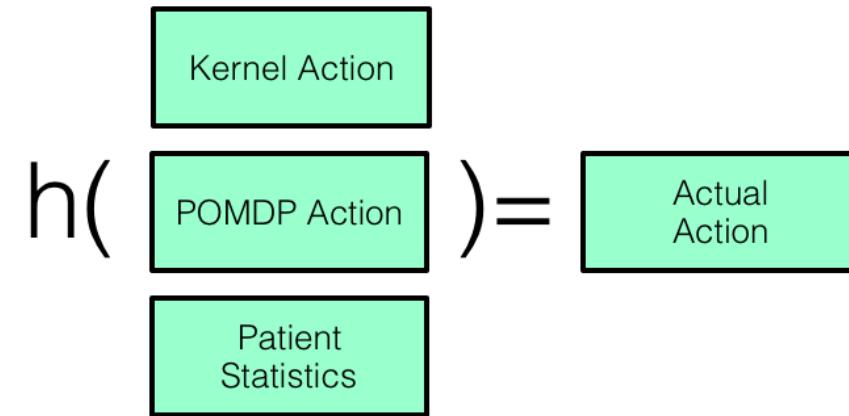


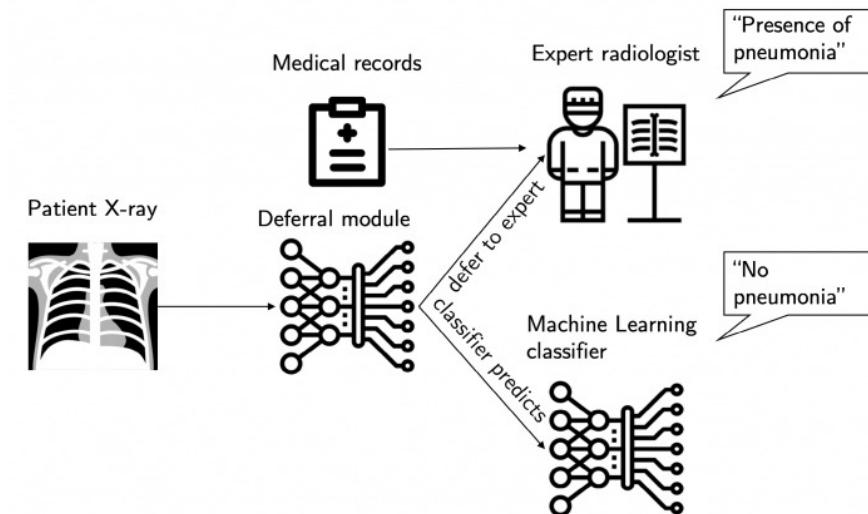
Image Credit: Doshi-Velez, Towards using Batch RL to Identifying Treatment Options, 2022

# Learning to Defer to Humans in Uncertainty

- Human expertise is valuable in high-risk settings such as healthcare.
- Domain experts can often **provide additional knowledge** that can help reduce uncertainties for rare cases or outlier patients.
- Though prior works have explored ways to combine humans and models, none of these account for **uncertainty** and **long-term effects**.

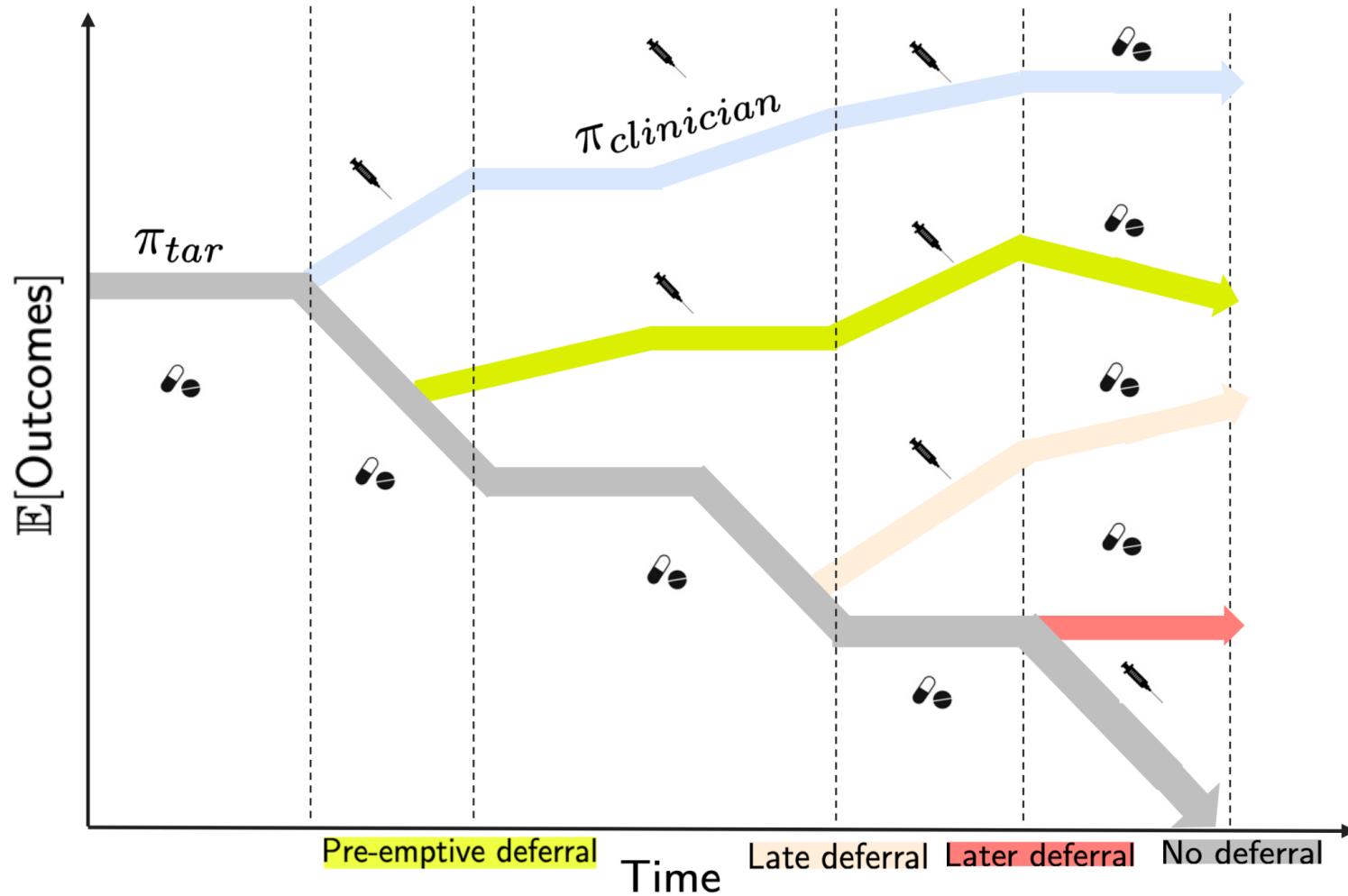


Combining Kernel and Model-based RL for HIV Therapy Selection – Parbhoo et al, 2017, AMIA

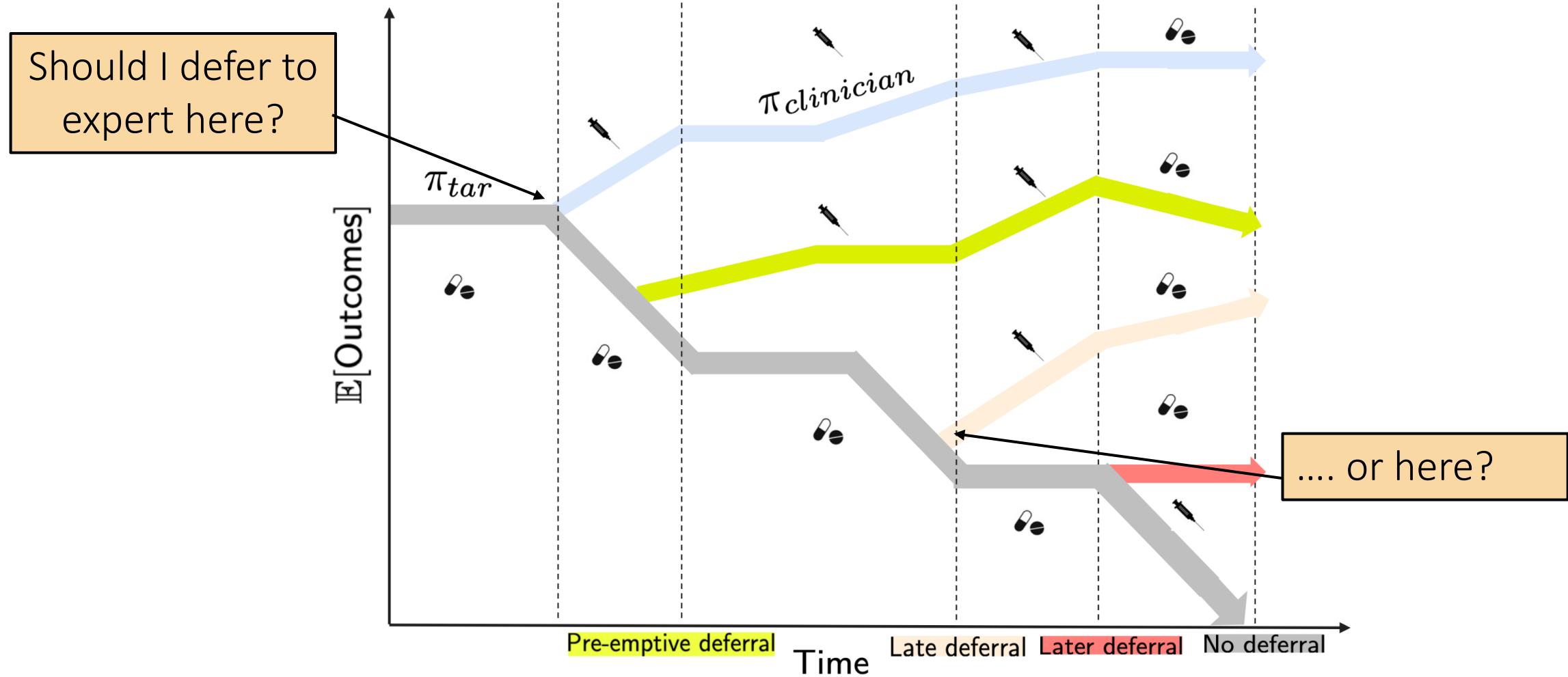


Consistent Estimators for Learning to Defer to an Expert – Mozannar et al, 2020, ICML

# Learning to Defer to Humans in Sequential Settings



# Learning to Defer to Humans in Sequential Settings

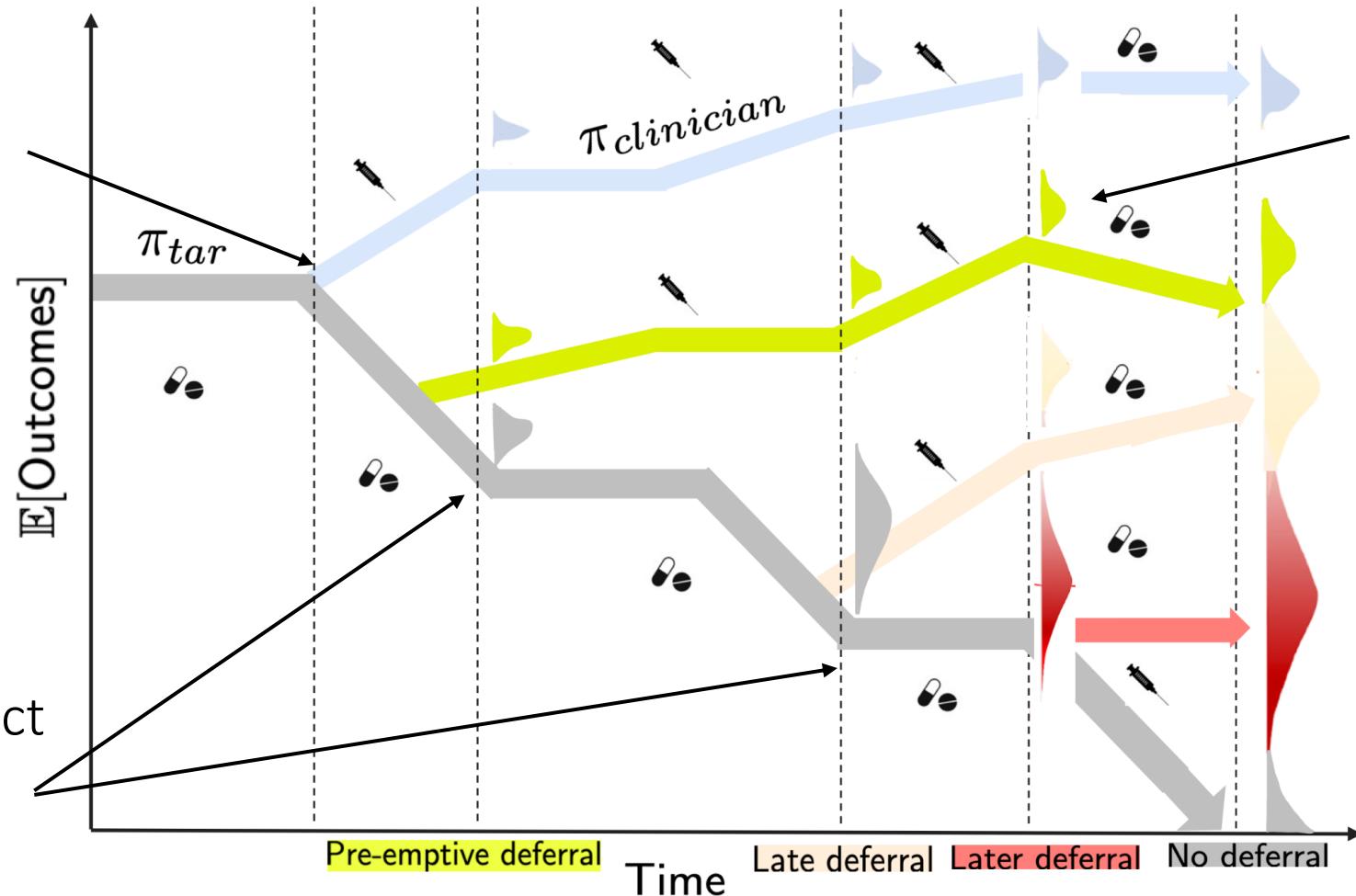


# Uncertainty-aware Mixtures

Quantify where  
a **model cannot**  
**be trusted**

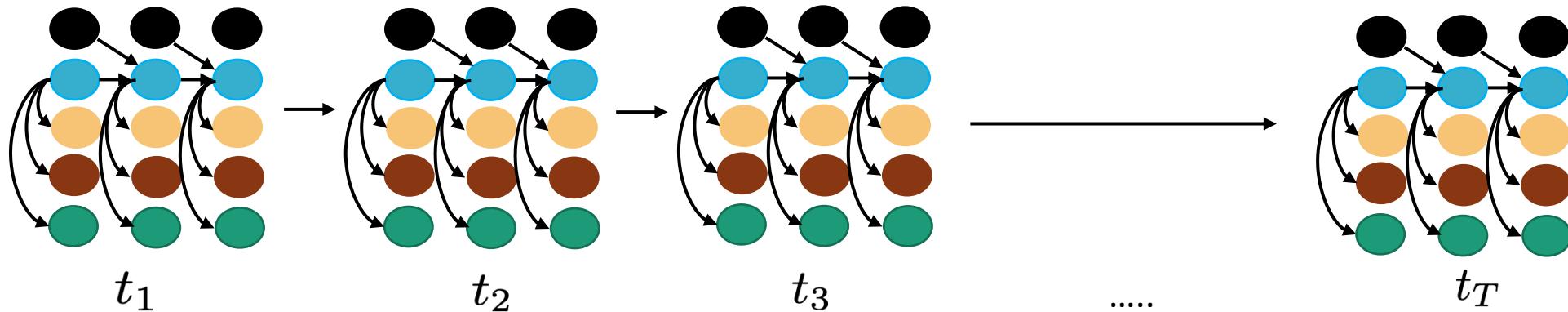
Quantify impact  
of **delaying**  
**deferral**

Model **changing**  
**dynamics** and  
**uncertainty**



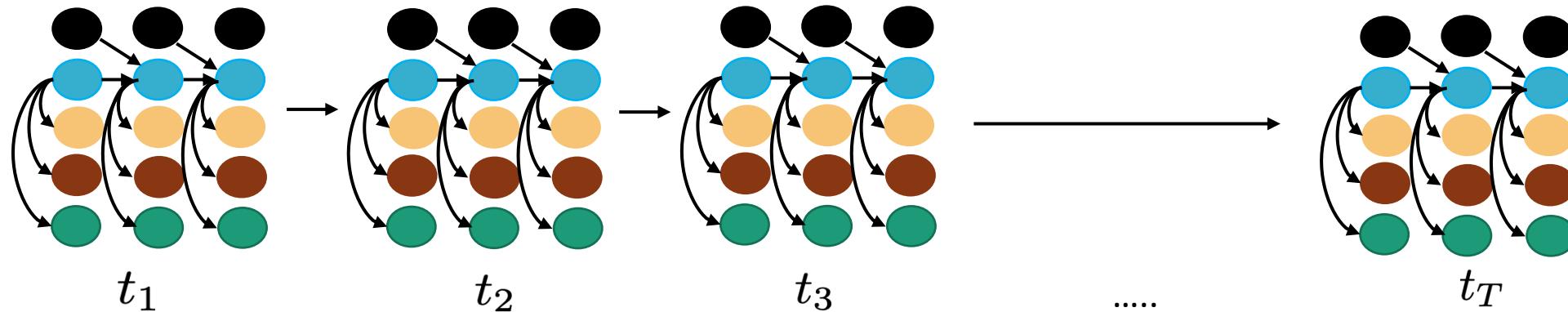
# *Uncertainty-aware Mixtures*

Step 1: Learn a model of **non-stationary** dynamics  $\{\mathcal{M}_t\}_t$

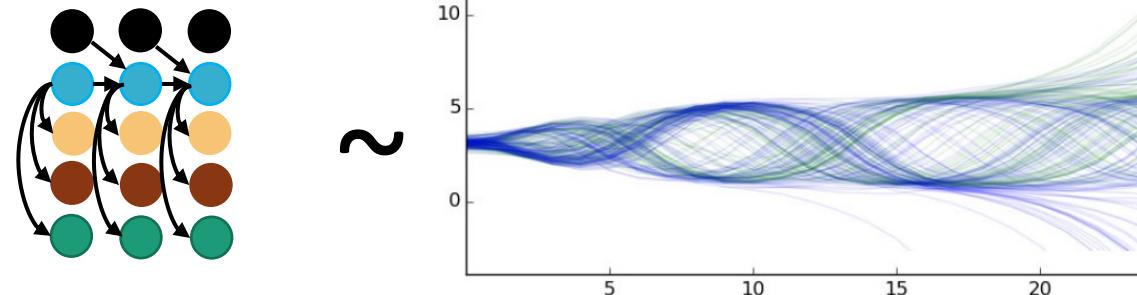


# *Uncertainty-aware Mixtures*

Step 1: Learn a model of **non-stationary** dynamics  $\{\mathcal{M}_t\}_t$



Step 2: Use posterior sampling to capture **uncertainty over models**  $\mathcal{M}_t \sim p(\cdot | \mathcal{D})$



# Sequential Learning to Defer (SLTD)

Step 3: Quantify the impact of **delaying deferral**

Repeat this over all states and times

# Sequential Learning to Defer: Algorithm

---

**Algorithm 1** Sequential Learning to Defer

---

**Input:** Posterior estimates  $\{p_t(\cdot|\mathcal{D})\}_{t=0}^T$ , target policy  $\pi_{\text{tar}}$ , behavior policy  $\pi_0$ .

**Initialization:** Deferral function  $g_{\pi_{\text{tar}}}(s, t) = 0$  for all  $s \in \mathcal{S}$  and  $t \in \{1, 2, \dots, T\}$ .

**for**  $t \in \{T, T-1, \dots, 1\}$  **do**

**for**  $s \in \mathcal{S}$  **do**

        Compute  $\{V_{\pi_{\text{tar}}, t}^{\mathcal{M}}(s)\}$  and  $\{V_{\perp, t}^{\mathcal{M}}(s) - c\} \forall \mathcal{M}$

        Update  $g_{\pi_{\text{tar}}}(s, t) \leftarrow \approx \mathbb{E}_{\mathcal{M} \sim p_t(\cdot|\mathcal{D})} [\mathbf{1}(V_{\pi_{\text{tar}}, t}^{\mathcal{M}}(s) > V_{\perp, t}^{\mathcal{M}}(s) - c) > \tau]$

**end for**

**end for**

**return**  $g_{\pi_{\text{tar}}}(s, t)$

---

# Decomposing the uncertainty at deferral

- The outcome we are interested in is

$$\mathbb{E}[r_T|s_{t_d}, \mu_{t_d}] = \int_{s_{t_d+1}}^{s_T} \int_{a_{t_d}}^{a_T} \int_{\mu_{t_d+1}}^{\mu_T} \int_{\theta_{t_d}}^T r(s_T, a_T) \prod_{t'=t_d+1}^T p(s_{t'}|\mu_{t'}) p(\mu_{t'}|\theta'_t(s_{t'}, a_{t'})) \pi_{t'}(a_{t'}|s_{t'}) p(\theta_{t'}|\mathcal{D}) d\mathbf{s}_{t_d+1}^T d\mathbf{a}_{t_d}^T d\boldsymbol{\mu}_{t_d+1}^T d\boldsymbol{\theta}_{t_d}^T$$

Decompose

$$\underbrace{\text{Var}(r_T|s_{t_d}, \mathcal{D})}_{\text{Total Uncertainty}} = \underbrace{\mathbb{E}_{\mu_{t_d} \sim p(\mu_{t_d}|\mathcal{D})}[\text{Var}(r_T|\mu_{t_d}, s_{t_d}, \mathcal{D})]}_{\text{Irreducible/ Aleatoric Uncertainty}} + \underbrace{\text{Var}_{\mu_{t_d} \sim p(\mu_{t_d}|\mathcal{D})}(\mathbb{E}[r_T|\mu_{t_d}, s_{t_d}, \mathcal{D}])}_{\text{Epistemic/Modeling Uncertainty}}$$

- First term averages over the variance due to  $\mu_{t_d}$  and captures **aleatoric uncertainty** at  $t_d$
- Second term captures variance conditioned on knowledge of the model or **epistemic uncertainty** at deferral time.

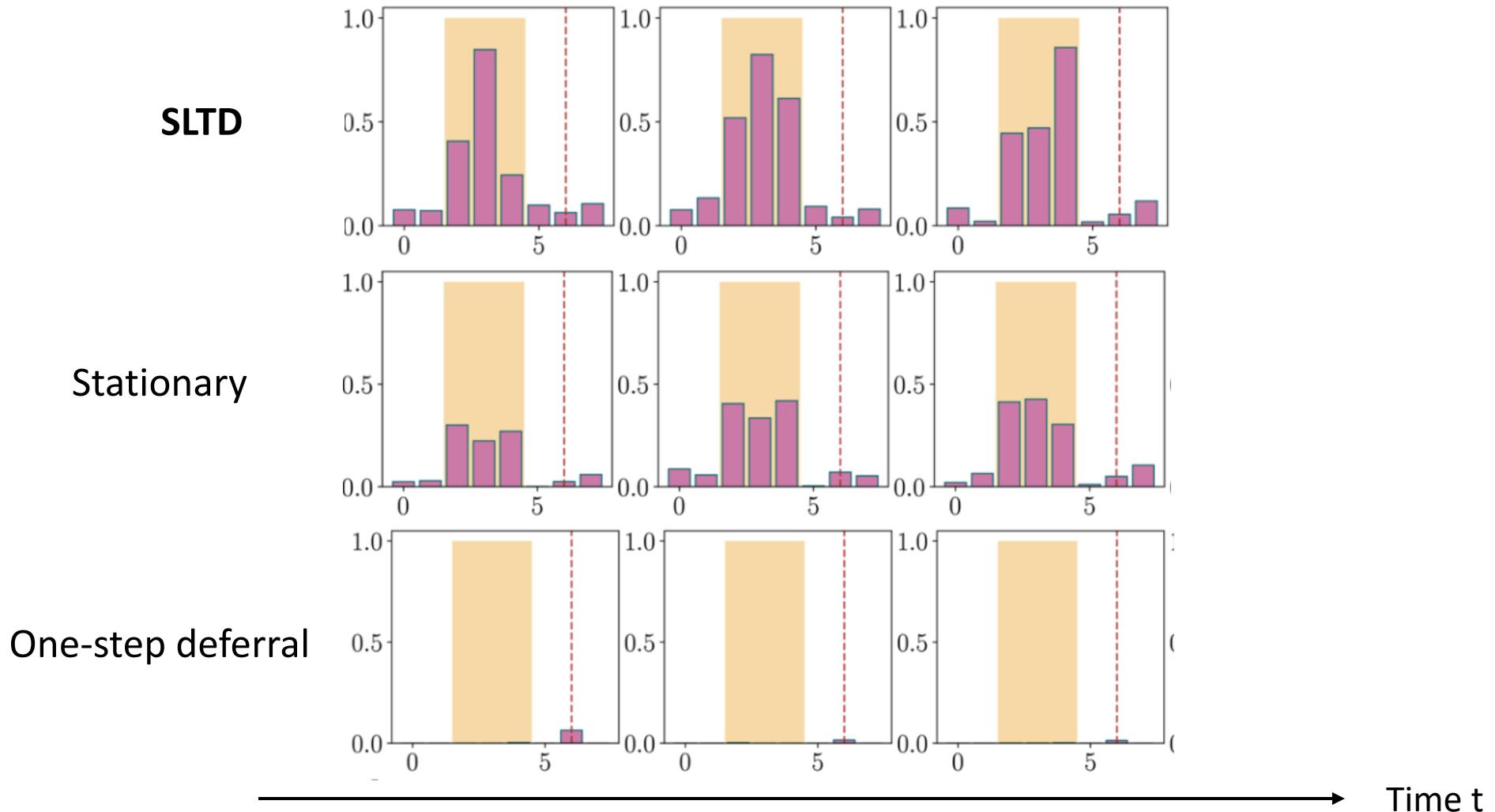
# What can these uncertainties tell us?

- **High Epistemic Uncertainty:**
  - Uncertainty of model's prediction is high
  - Possible improvement if additional data in this region could be collected to improve the reducible sources of uncertainty.
- **High Aleatoric Uncertainty:**
  - There is high variability in the patient's dynamics that may need to be managed with careful interventions or is otherwise not manageable.
  - Aleatoric uncertainty is important for safety
- Based on uncertainty, expert may choose to deviate from practice and/or get second opinions.

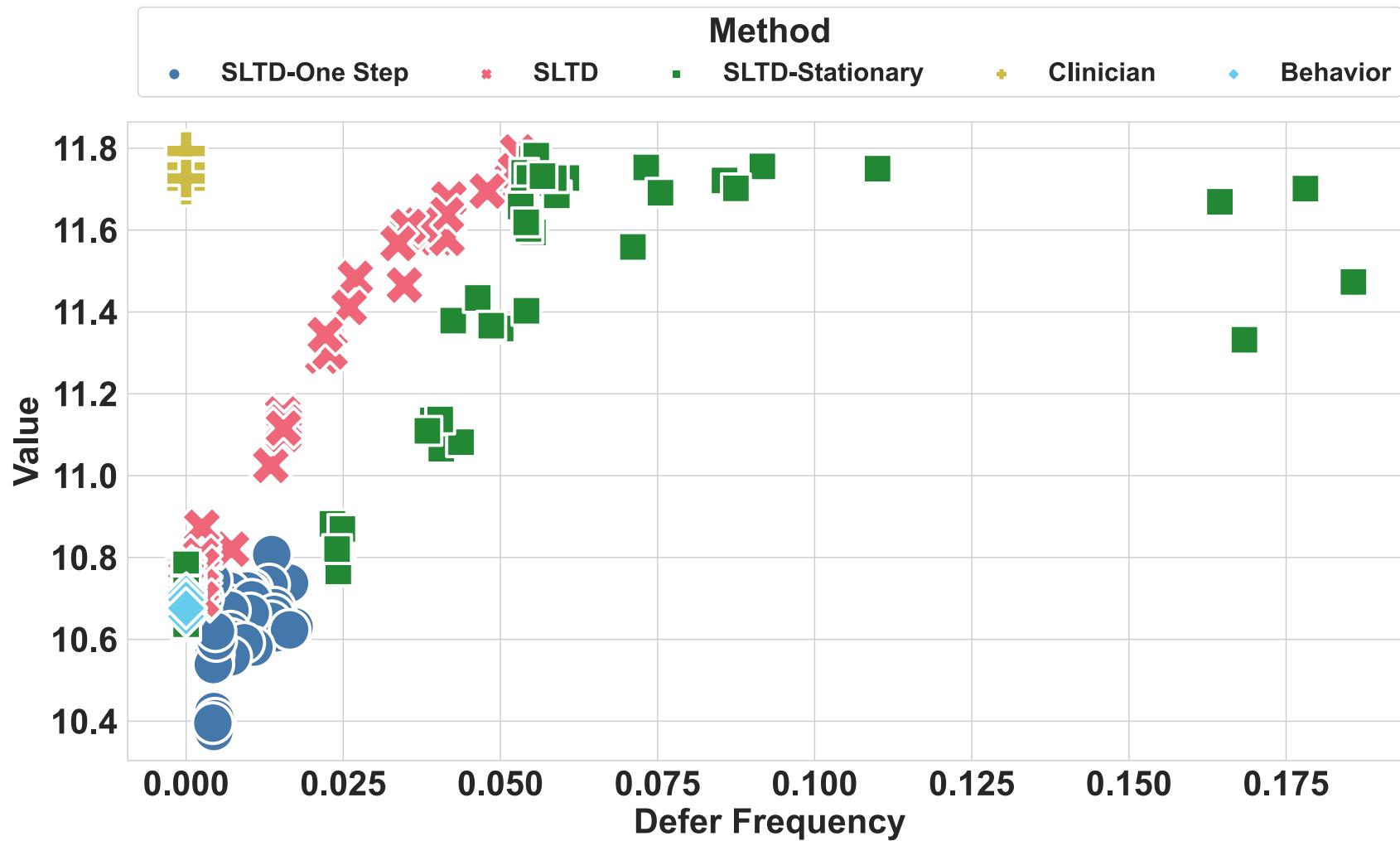
## Toy Demonstration

- Given 8 states:  $s_6$  has reward -5; while others have reward +1
- Each sample starts at  $s_0$  and progresses to  $s_7$
- Two actions:  $a_0$  decreases chances of landing in  $s_6$ ;  $a_1$  increases this
- Target policy  $\pi_{tar}$  increases the chances of landing in  $s_6$  from  $\{s_2, s_3, s_4\}$  at times 3-8.
- Non-stationarity increases over time, while rewards are positive in  $\{s_2, s_3, s_4\}$  so that poor outcomes only manifest in the future.
- A pre-emptive deferral policy should defer between times 3 and 8.

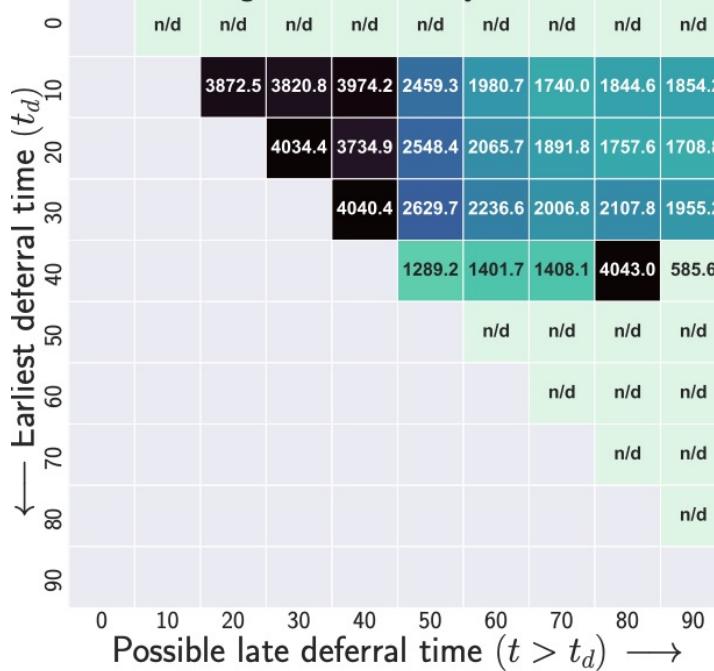
# SLTD learns a *pre-emptive* deferral policy



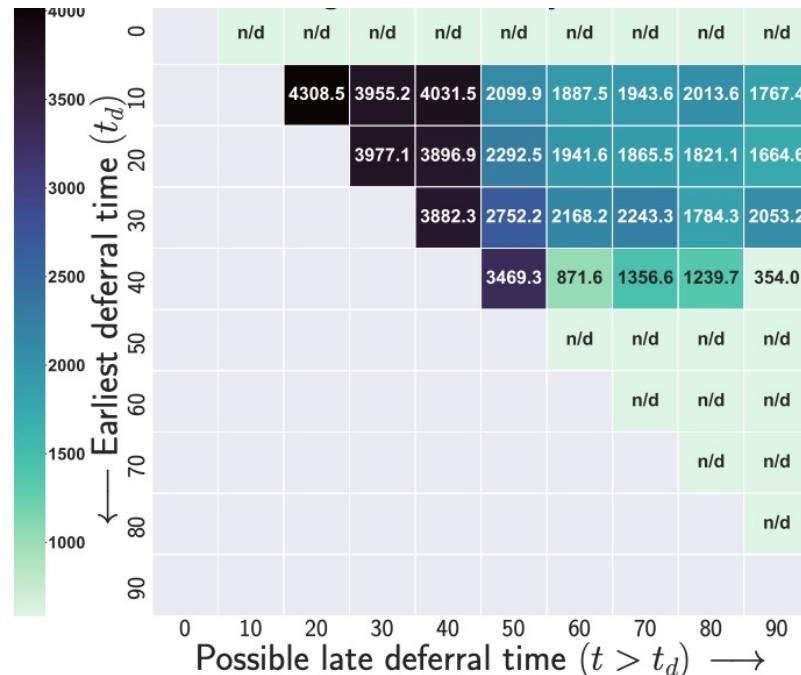
# SLTD produces better outcomes with *fewer deferrals*



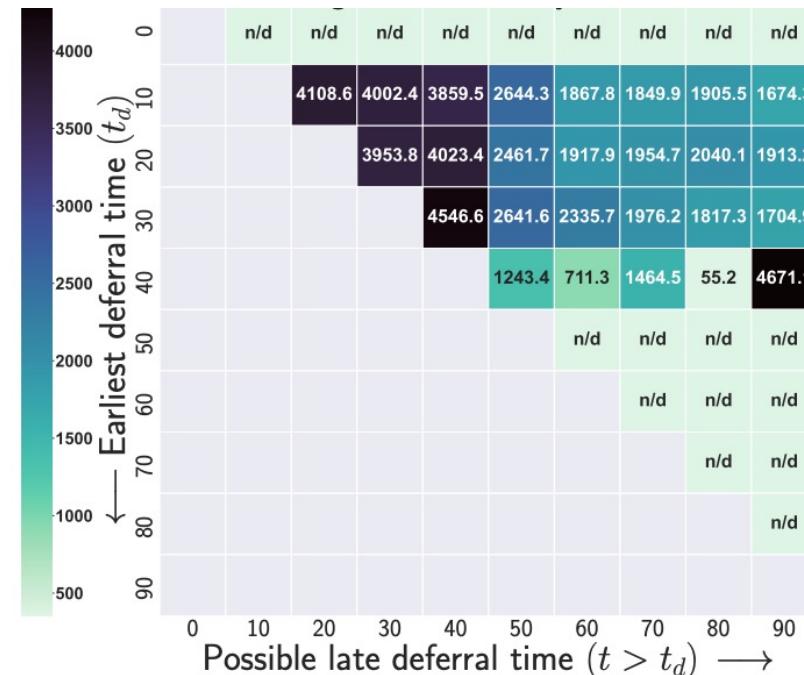
# Deferring early potentially *reduces overall uncertainty*



SLTD

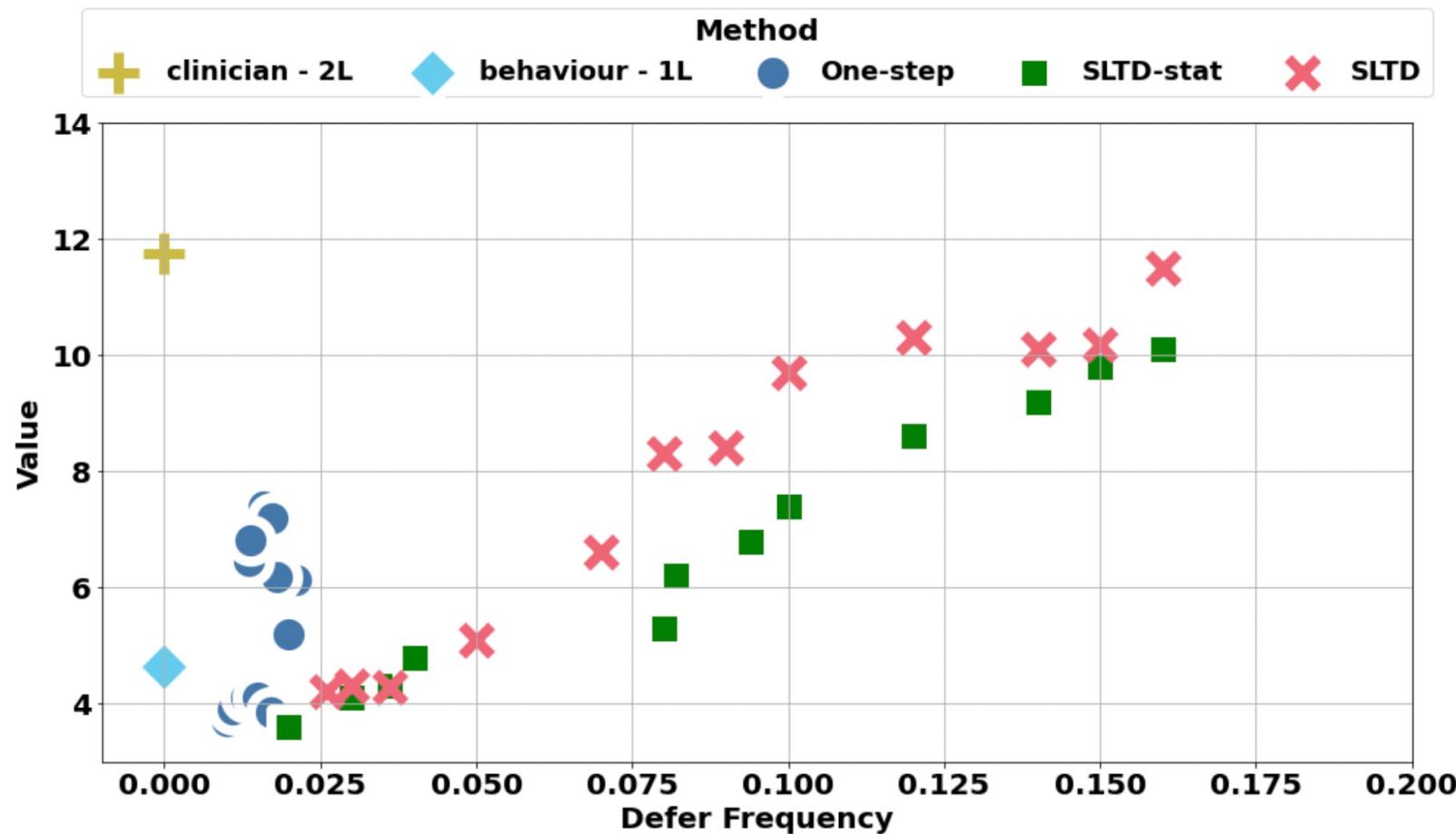


Stationary



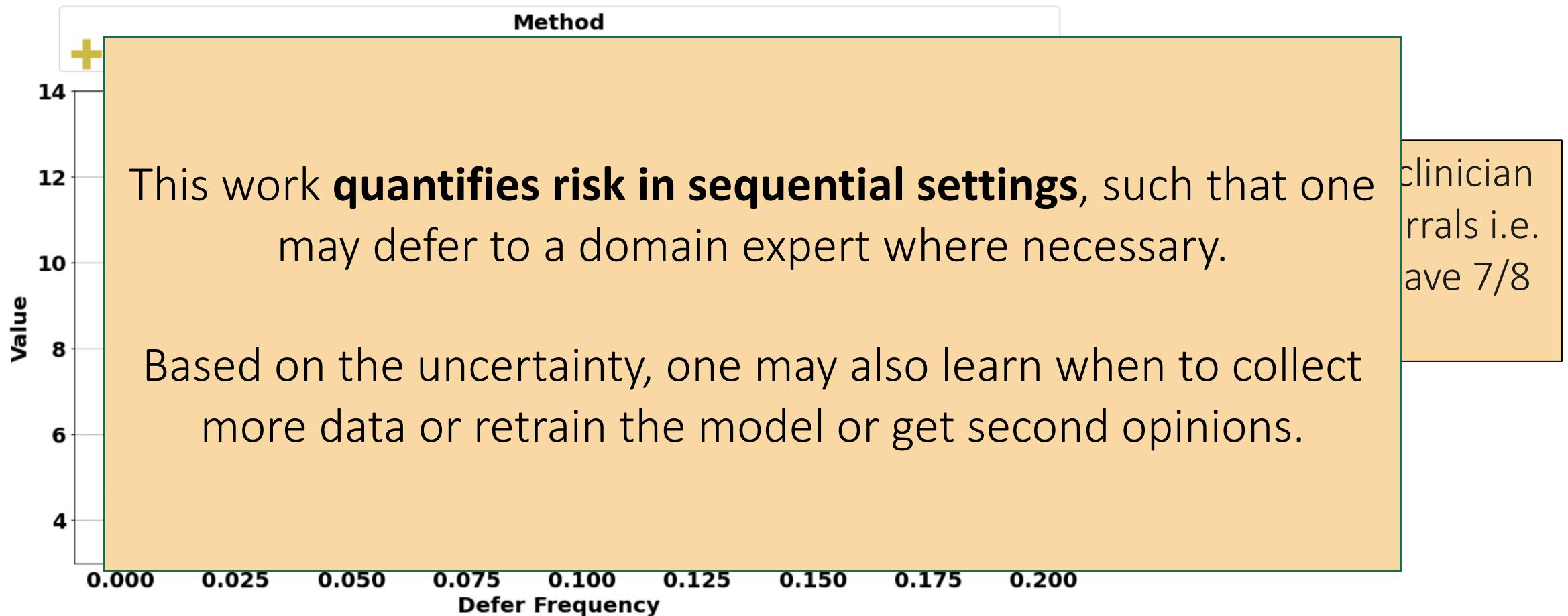
One-step

# Similar performance for HIV Therapy Management



We get close to clinician policy with +/- 1/8 deferrals i.e. can potentially save 7/8 visits.

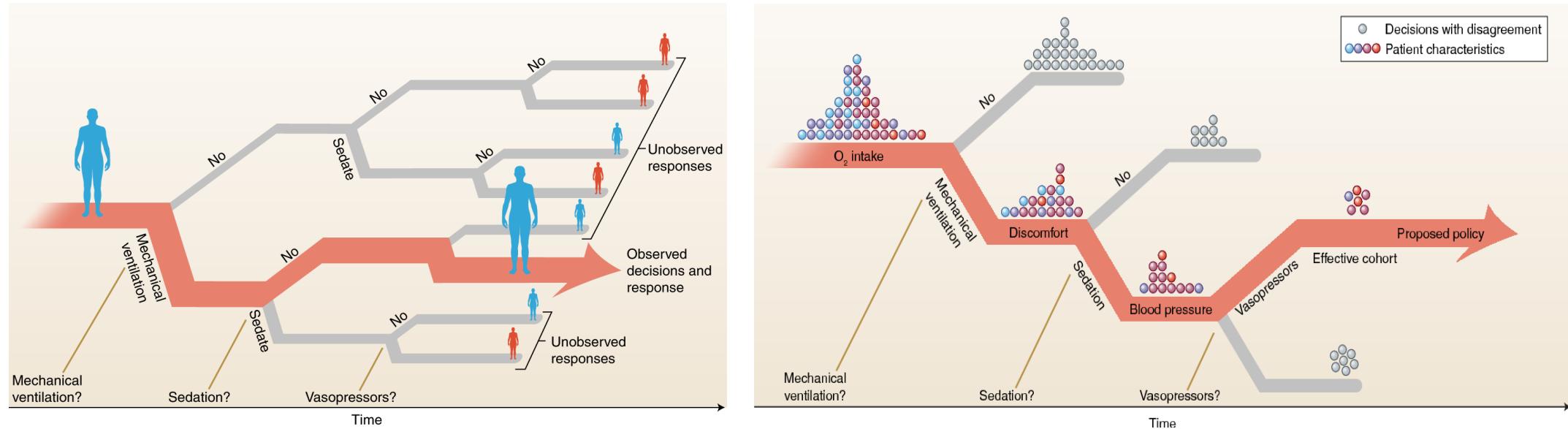
# Similar performance for HIV Therapy Management



# Challenges to using RL for healthcare in practice

Learning suitable treatment policies from observational data is **offline and off-policy**

- 1.Inability to explore
- 2.Small data (relatively) → shrinking data support as best strategies are discovered

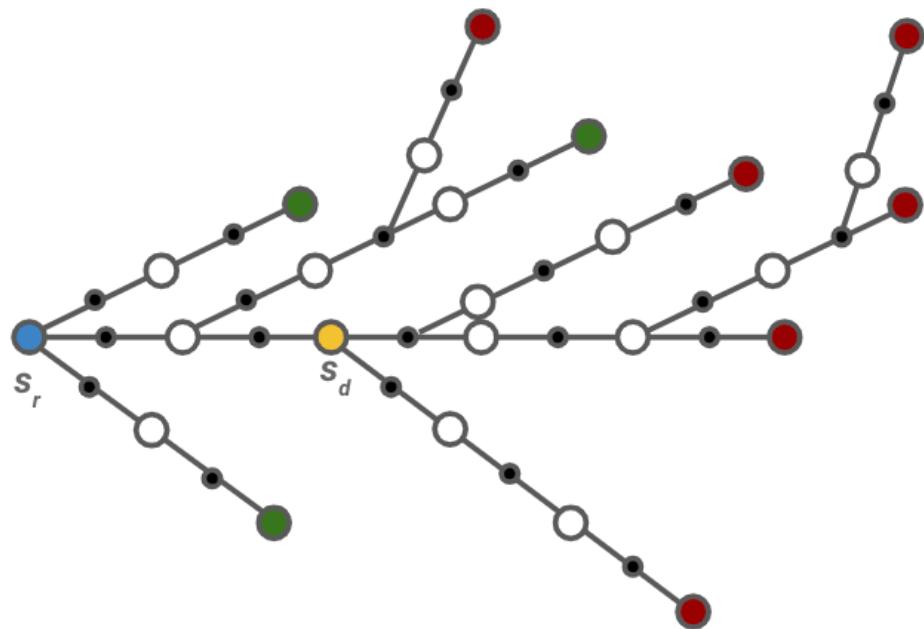


**These two limitations severely complicate the development of actionable RL**

Image Credit: Gottesman, Omer, et al. "Guidelines for reinforcement learning in healthcare." *Nat Med* 25.1 (2019): 16-18.

# An alternative offline RL paradigm

Rather than the difficult task of learning optimal policies that suggest **what to do**, we instead learn **what treatments to avoid**



**Undesired Terminal State** (e.g., mortality)

**Desired Terminal State** (e.g., recovery)

**Dead-end**: all trajectories from  $s_d$  reach an undesired terminal state with probability 1.

**Rescue**: from  $s_r$  a desired end state is reachable with probability 1

Can we identify all dead-ends and treatments that lead to them?

# Treatment Security

To avoid selecting treatments  $a$  that lead to dead-ends, constrain the policy  $\pi$

Def: ***Treatment Security***,

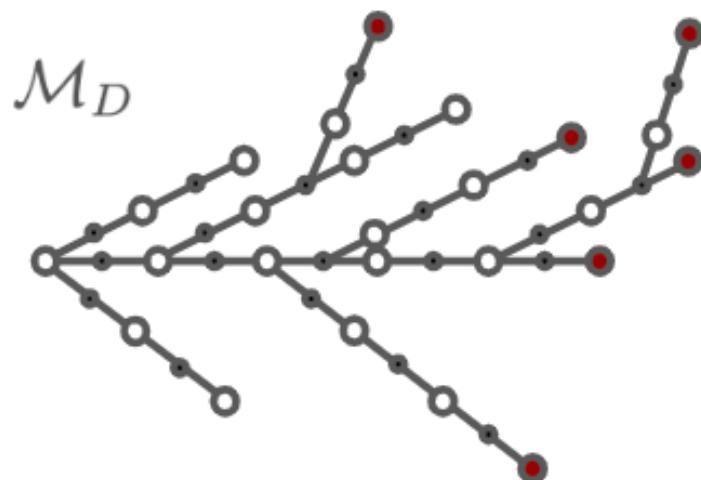
$$P_D(s, a) + F_D(s, a) \geq \lambda \implies \pi(s, a) \leq 1 - \lambda$$

But

- Inferring max  $\lambda$  for all possible  $(s, a)$  pairs is intractable
- Accurate estimates of  $P_D$  and  $F_D$  requires knowledge of all dead-ends and transition probabilities a priori.

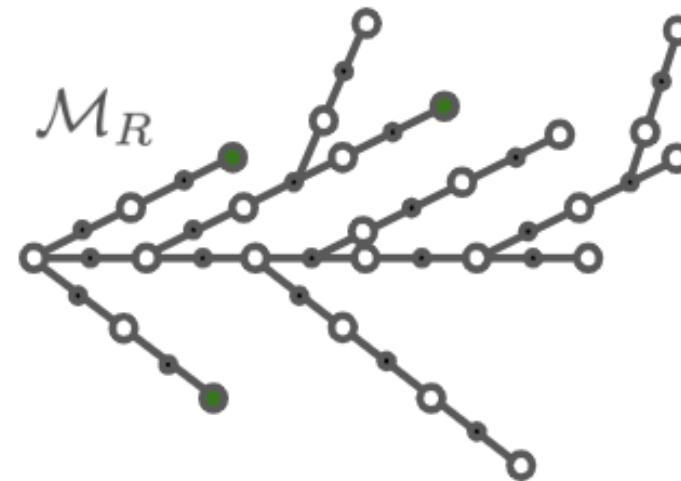
# Treatment Security

Separate the MDP into two *decoupled* processes



Zero-out all positive outcomes  
No discounting ( $\gamma_D = 1$ )

$$\rightarrow Q_D^* \in [-1, 0]$$



Zero-out all negative outcomes  
No discounting ( $\gamma_R = 1$ )

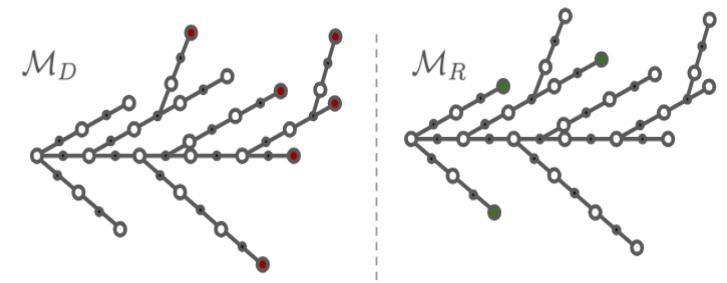
$$\rightarrow Q_R^* \in [0, 1]$$

# A Framework for Treatment Security: Dead-End Discovery (DeD)

As a consequence of  $\mathcal{M}_D$

$$-Q_D^*(s, a) = P_D(s, a) + F_D(s, a)$$

This means that  $-Q_D^*$  carries special physical meaning



- It is the minimum probability of an *undesired terminal outcome*
- $1 + Q_D^*(s, a)$  is the maximum probability of a *desired terminal outcome*

Using this interpretation of  $-Q_D^*$  the treatment security condition is satisfied if

$$P_D(s, a) + F_D(s, a) \geq \lambda \Rightarrow \pi(s, a) \leq 1 - \lambda \quad \rightarrow \quad \pi(s, a) \leq 1 + Q_D^*$$

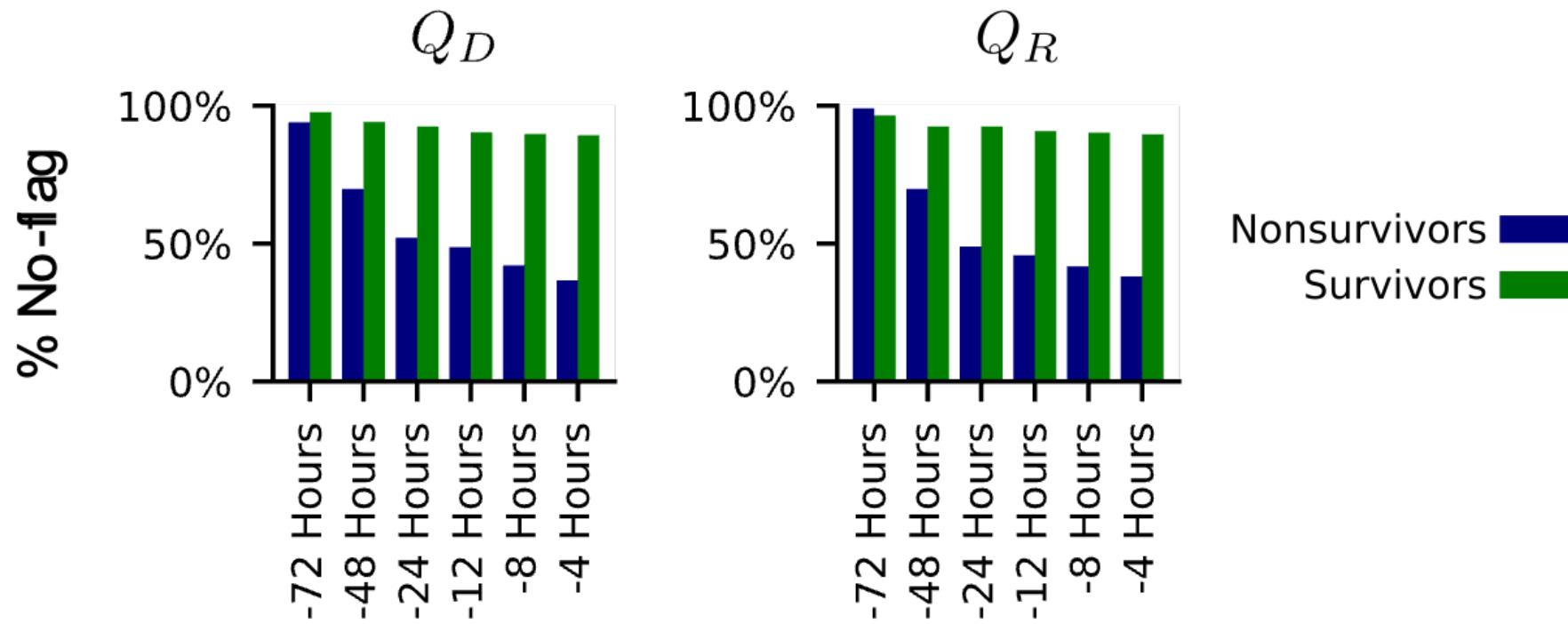
# A Framework for Treatment Security: Dead-End Discovery (DeD)

It can be shown:

- $P_D(s, a) + F_D(s, a) = 1$  if and only if  $Q_D^*(s, a) = -1$
- $P_R(s, a) + F_R(s, a) = 1$  if and only if  $Q_R^*(s, a) = 1$
- If  $\pi(s, a) \leq 1 + Q_D^*$  and  $P_D(s, a) + F_D(s, a) \geq \lambda$  then  $\pi(s, a) \leq 1 - \lambda$  for all  $\lambda$
- If  $\pi(s, a) \geq Q_R^*(s, a)$  and  $P_R(s, a) + F_R(s, a) \geq \lambda$  then  $\pi(s, a) \geq \lambda$  for all  $\lambda$
- There exists a threshold  $\delta_D \in (-1, 0)$  that separates dead-ends through  $Q_D^*(s, a) \geq \delta_D$
- There exists a threshold  $\delta_R \in (0, 1)$  that separates rescue states through  $Q_R^*(s, a) \geq \delta_R$

# DeD framework can help identify dead-ends in septic patients

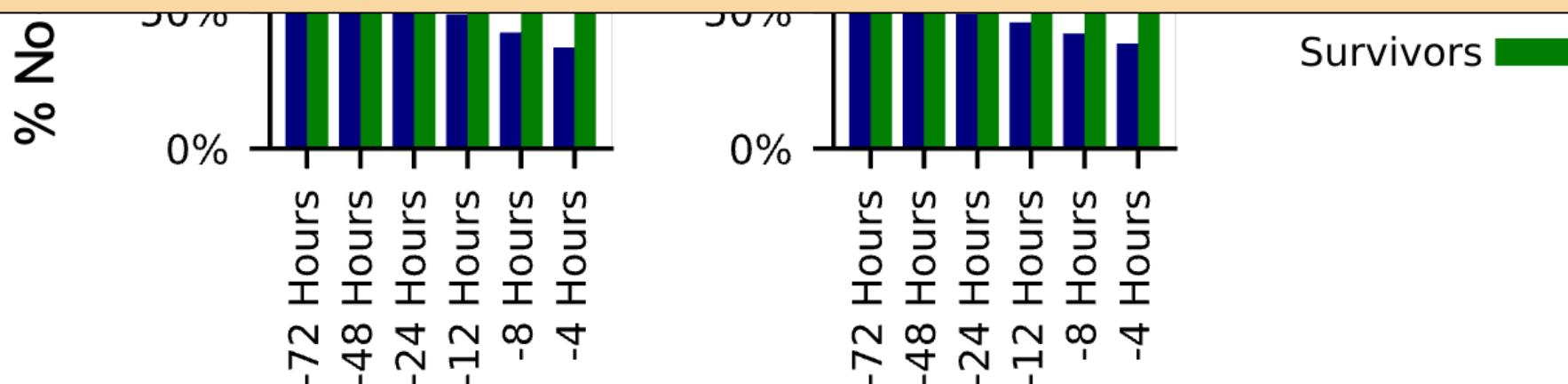
Clear distinction between patients who survived and those who did not in terms of DeD providing specified risk several hours prior to terminal outcome



# DeD framework can help identify dead-ends in septic patients

Clear distinction between patients who survived and those who did not in terms of DeD

... but there are major limitations. Crucially, DeD uses a naïve definition of risk. We need to account for the **full distribution of value functions.**



# Potential Sources of Uncertainty in RL

## Dynamics

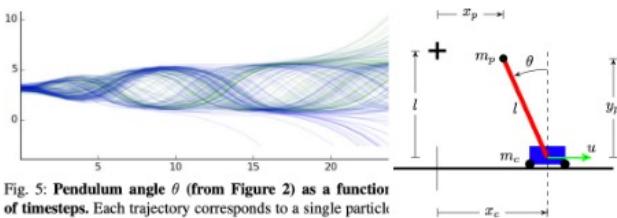


Fig. 5: PENDULUM angle  $\theta$  (from Figure 2) as a function of timesteps. Each trajectory corresponds to a single particle

From Gal, et al; [ICML 2016]

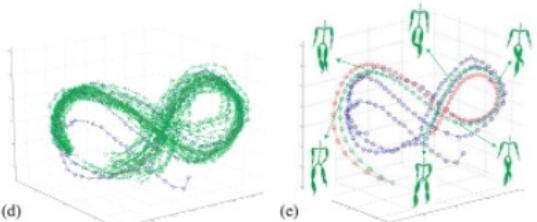


Figure 2: Models learned from a walking sequence of 2.5 gait cycles. The latent positions

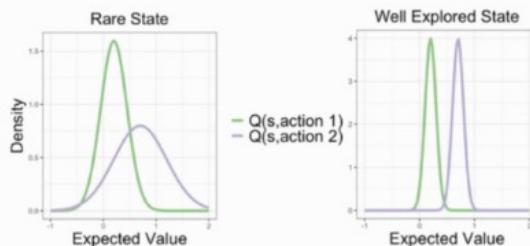
From Wang, et al; [NeurIPS 2005]

Wang, Jack, Aaron Hertzmann, and David J. Fleet. "Gaussian process dynamical models." *Advances in neural information processing systems* 18 (2005).

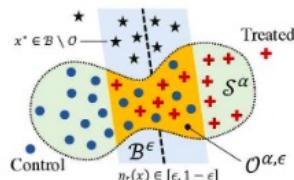
Gal, Yarin, Rowan McAllister, and Carl Edward Rasmussen. "Improving PILCO with Bayesian neural network dynamics models." *Data-Efficient Machine Learning workshop, ICML*. Vol. 4. No. 34. 2016.

Killian, T. W., Daulton, S., Konidaris, G., & Doshi-Velez, F. Robust and efficient transfer learning with hidden parameter markov decision processes. *Advances in neural information processing systems*, 30 (2017).

## Actions (Policy)

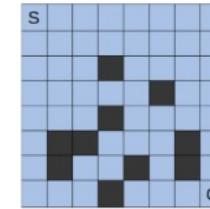


From Sonabend, et al; [NeurIPS 2020]



From Oberst, et al; [AISTATS 2020]

## Outcomes (Reward)



Frozen Lake



Windy Grid World



From Bellemare, et al; [ICML 2017]

Dabney, W., Ostrovski, G., Silver, D., & Munos, R. Implicit quantile networks for distributional reinforcement learning. In *International conference on machine learning* (2018)

Stanko, S., & Macek, K. Risk-averse Distributional Reinforcement Learning: A CVaR Optimization Approach. In *IJCCI* (2019)

Bellemare, M. G., Dabney, W., & Munos, R. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning* (2017)

# A Short History of Distributional RL

Traditionally, RL utilizes a point estimate for the expected return from any state-action pair. Bellemare, et al (2017)\* introduced an alternative, representing the full return distribution

$$Z^\pi = \sum_{t=0}^{\infty} \gamma^t r_t$$

$$\mathcal{T}^\pi Z(s, a) \stackrel{D}{=} r + \gamma Z(s', a')$$

where  $r \sim R(\cdot|s, a)$ ,  $s' \sim P(\cdot|s, a)$ ,  $a' \sim \pi(\cdot|s')$

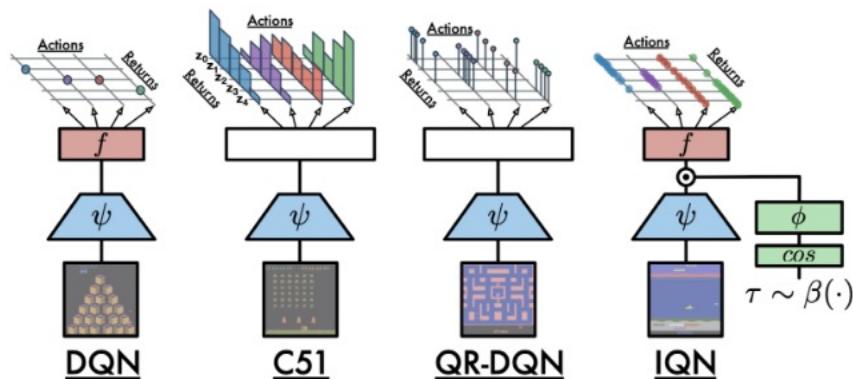
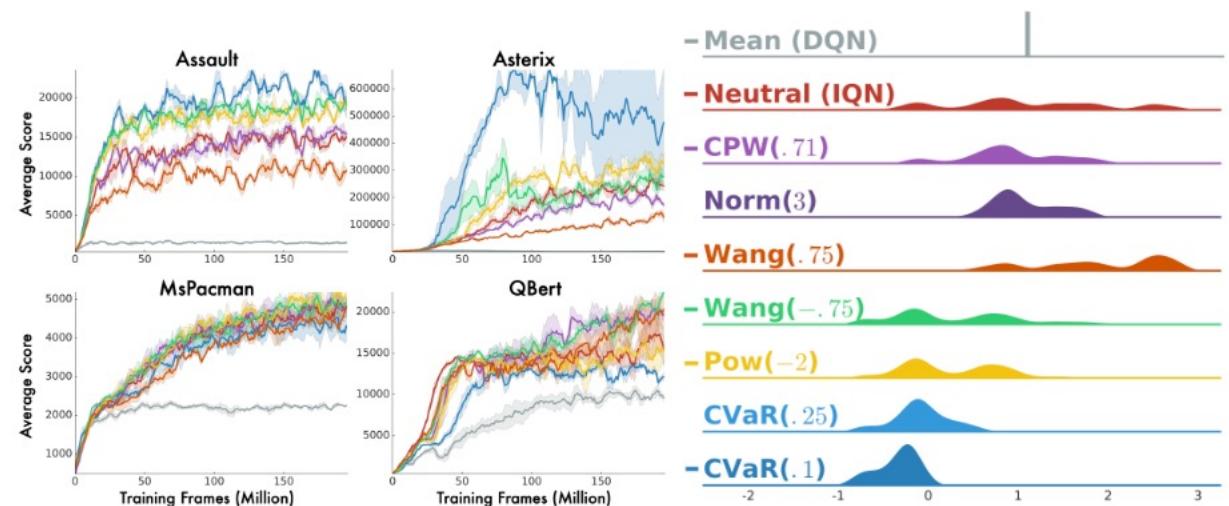


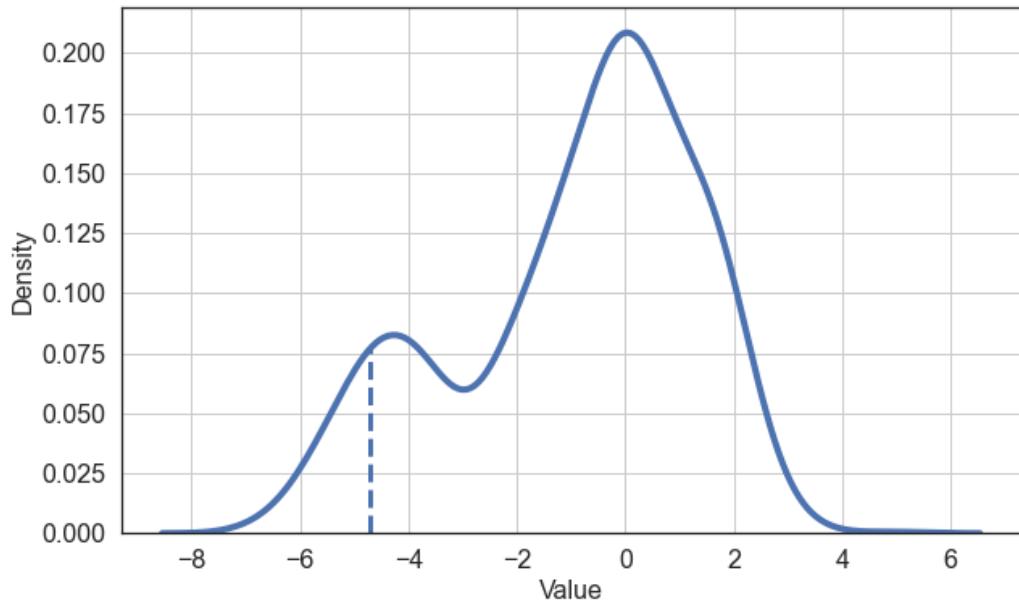
Figure 1. Network architectures for DQN and recent distributional RL algorithms.

From Dabney, et al; [ICML 2018]



From Dabney, et al; [ICML 2018]

# Risk estimation with the full distribution



For some state  $s$  and action  $a$ , we can assess the estimated value distribution  $Z^\pi(s, a)$  and infer the worst-case outcome that occurs with some pre-defined confidence level  $\eta$  (e.g. value-at-risk).

Value-at-risk (VaR):

$$\text{VaR}_\alpha(Z) = \min \{z \mid \alpha \leq F(z)\}$$

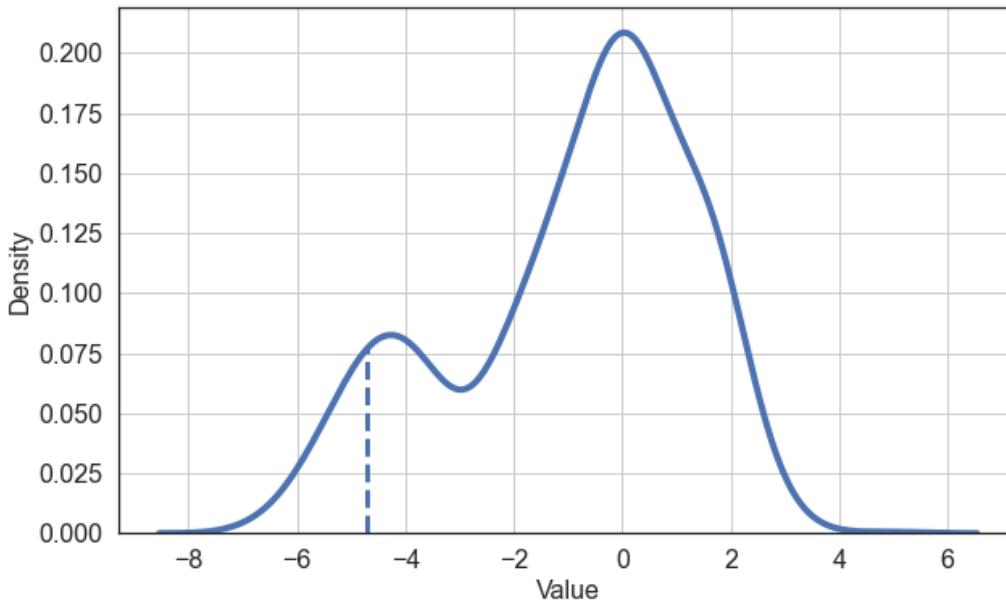
Conditional value-at-risk (CVaR):

$$\text{CVaR}_\alpha(Z) = \frac{1}{\alpha} \mathbb{E}[(Z - \text{VaR}_\alpha(Z))^+] + \text{VaR}_\alpha(Z)$$

Dual Formulation of CVaR

$$\text{CVaR}_\alpha(Z) = \min_{\xi \in \mathcal{U}_{\text{CVaR}}(\alpha, \mathbb{P})} \mathbb{E}_\xi[Z]$$

# Risk estimation with the full distribution



For some state  $s$  and action  $a$ , we can assess the estimated value distribution  $Z^\pi(s, a)$  and infer the worst-case outcome that occurs with some pre-defined confidence level  $\eta$  (e.g. value-at-risk).

Then using the distribution that falls below  $\eta$ , we calculate the expected worst case outcome (e.g. CVaR). **This value assesses the risk of possible actions  $a$**

Value-at-risk (VaR):

$$\text{VaR}_\alpha(Z) = \min \{z \mid \alpha \leq F(z)\}$$

Conditional value-at-risk (CVaR):

$$\text{CVaR}_\alpha(Z) = \frac{1}{\alpha} \mathbb{E}[(Z - \text{VaR}_\alpha(Z))^-] + \text{VaR}_\alpha(Z)$$

Dual Formulation of CVaR

$$\text{CVaR}_\alpha(Z) = \min_{\xi \in \mathcal{U}_{\text{CVaR}}(\alpha, \mathbb{P})} \mathbb{E}_\xi[Z]$$

# Risk-Estimation for Medical Dead-Ends Distributional Treatment Security Condition (DistDeD)

With risk aversion parameter  $\alpha$ , we see that CVaR lower bounds the expected return:

$$\text{CVaR}_\alpha(Z) \leq \mathbb{E}[Z]$$

Extending this for DeD, we see:

$$-\text{CVaR}_\alpha(Z_D^*(s, a)) \geq -\mathbb{E}[Z^*(s, a)] = -Q_D^*(s, a) = P_D(s, a) + F_D(s, a)$$

We can provide a more robust and risk-averse treatment security condition:

$$\pi(s, a) \leq 1 + \text{CVaR}_\alpha(Z^*(s, a))$$

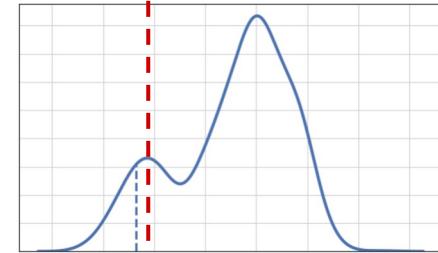
CVaR helps **lower bound** our choice of which actions to avoid and provides a mathematical means to assess risk.

# Application of DistDeD to AHE and Septic Patients

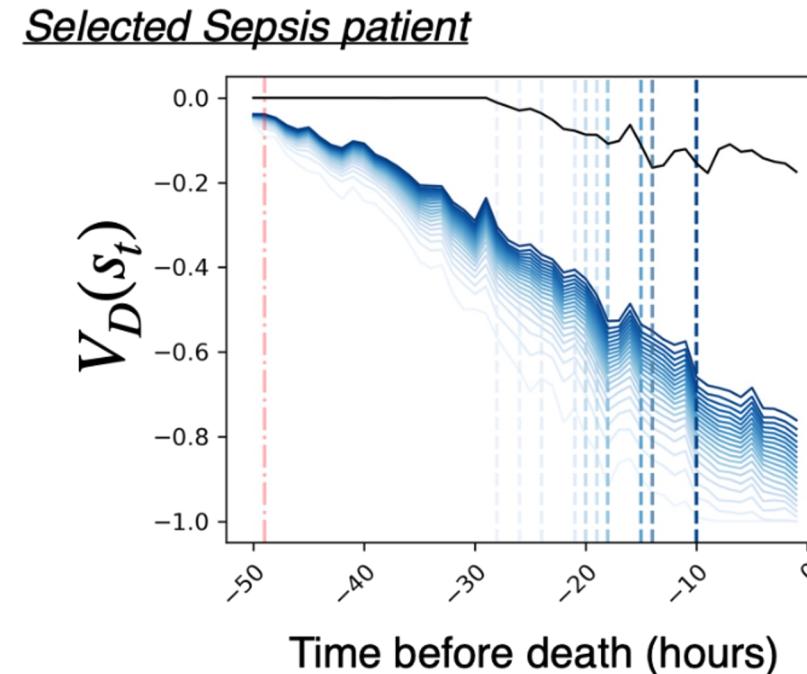
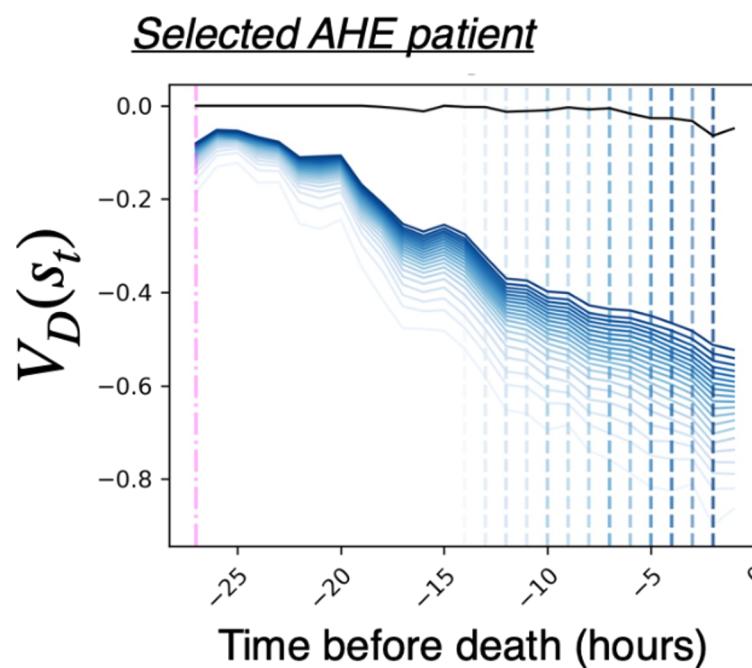
**Sepsis is a life-threatening organ dysfunction caused by infection**

- Using MIMIC-III we extract a cohort of 19,611 patient trajectories (72 hours, centered around presumed sepsis onset)
- Patient observations are made up of 47 features, aggregated in 4h windows, comprised of vital signs, lab measurements and demographics
- 25 available treatments (combination of 5 volumetric bins for each of vasopressors and IV fluid)
- Terminal outcomes: patient survival (+1) or death (-1)

# DistDeD provides a *tunable measure of risk*

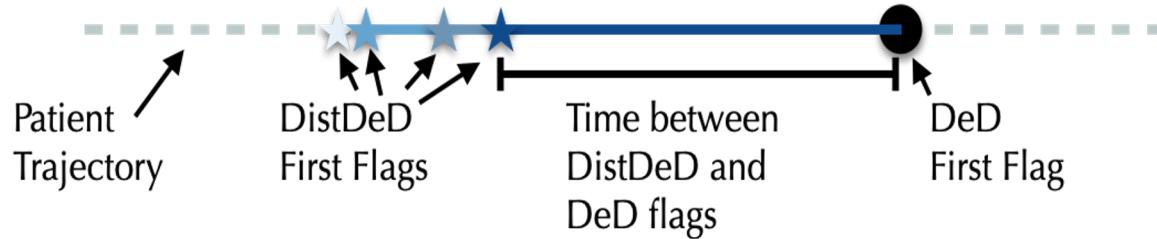
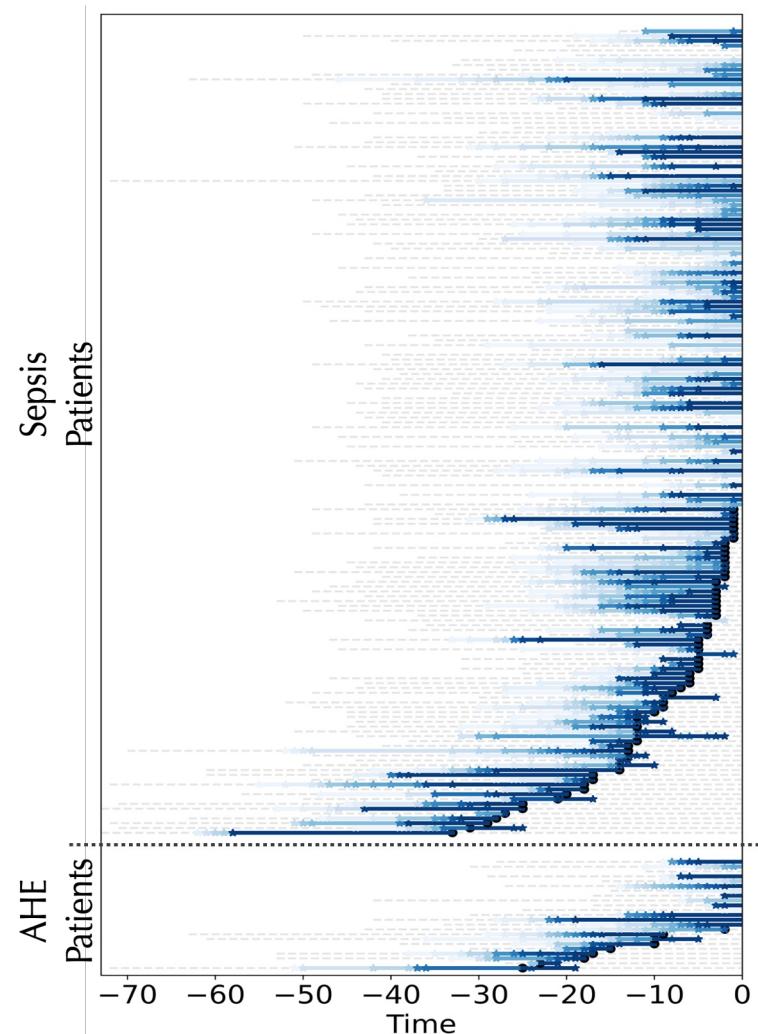


By modeling the **full distribution of outcomes**, DistDeD enables expert to modulate conservatism on a per-case basis depending on their own risk characterisation



- DeD
  - Presumed Onset
  - DistDeD
  - ==== First flag from DistDeD
- (We evaluate CVaR $\alpha$ ,  $\alpha \in [0.05, 1.0]$  in increments of 0.05)

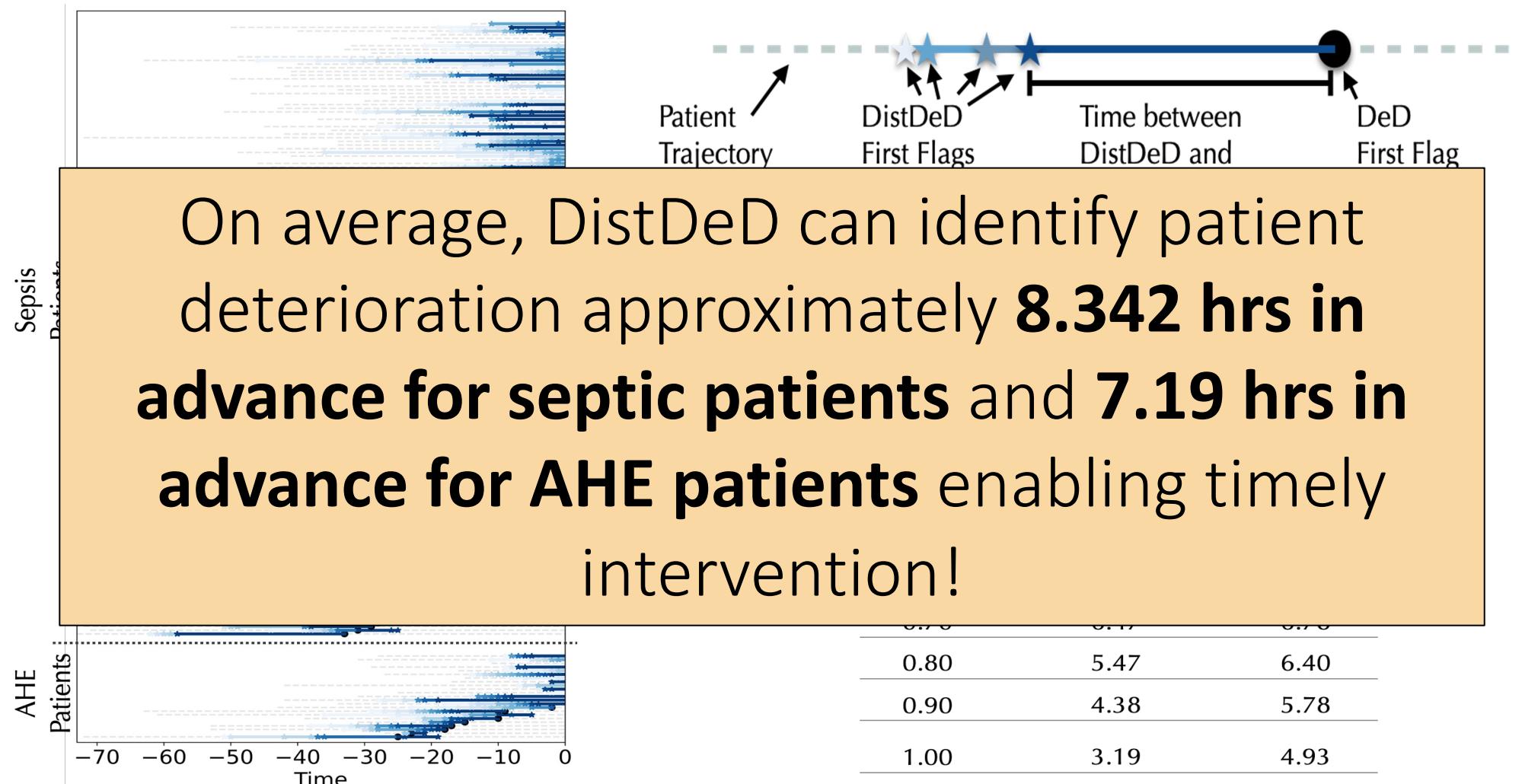
# DistDeD provides an *earlier indication of risk*



*DistDeD Avg. Hours Ahead (by CVaR $\alpha$  setting)*

<u>CVaR<math>\alpha</math> risk quantile (<math>\alpha</math>)</u>	AHE Patients	Sepsis Patients
0.10	11.29	14.47
0.20	9.75	11.31
0.30	8.50	9.65
0.40	8.00	8.74
0.50	7.67	7.91
0.60	7.18	7.45
0.70	6.47	6.78
0.80	5.47	6.40
0.90	4.38	5.78
1.00	3.19	4.93

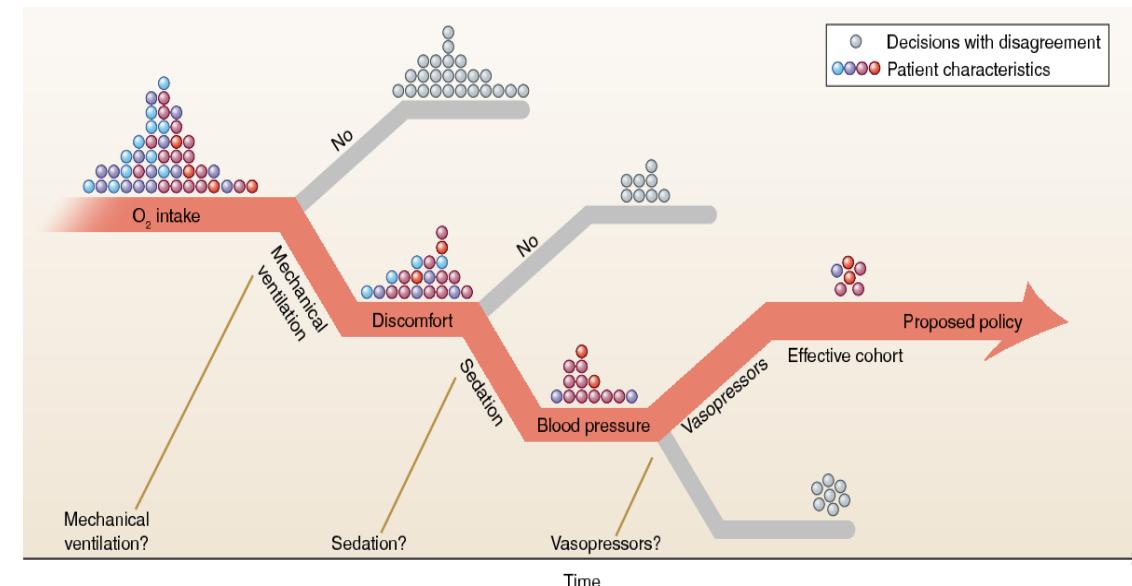
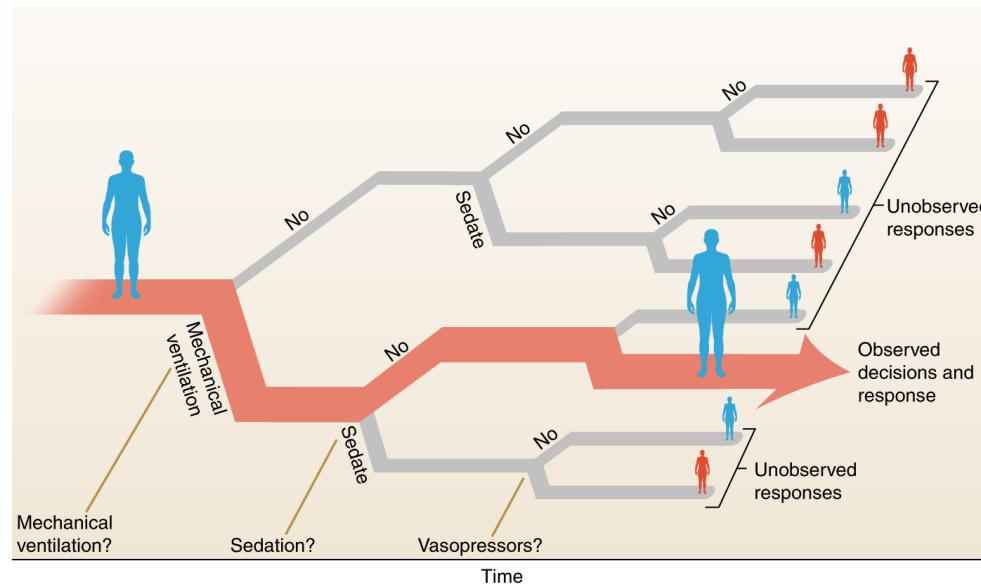
## DistDeD provides an *earlier indication of risk*



# Challenges to using RL for healthcare in practice

Learning suitable treatment policies from observational data is **offline and off-policy**

- 1.Inability to explore
- 2.Small data (relatively) → shrinking data support as best strategies are discovered



**These two limitations severely complicate the development of actionable RL**

Image Credit: Gottesman, Omer, et al. "Guidelines for reinforcement learning in healthcare." *Nat Med* 25.1 (2019): 16-18.

# SLTD and DeD cannot always be used in place of evaluation

- **DeD Framework doesn't fully support intermediate outcomes**
  - requires careful definition of what to avoid
- **Both methods do not account for sources of confounding and bias**
  - Value functions are prone to *overestimation*; not assessed for sensitivity to patient demographics

# SLTD and DeD cannot always be used in place of evaluation

- **DeD Framework doesn't fully support intermediate outcomes**
  - requires careful definition of what to avoid
- **Both methods do not account for sources of confounding and bias**
  - Value functions are prone to *overestimation*; not assessed for sensitivity to patient demographics

What are the alternatives in these settings?

## Robustness Checks for OPE

Are the results reproducible across different sites or rewards?  
E.g. Use two cohorts and 3 different OPE methods

EuResist

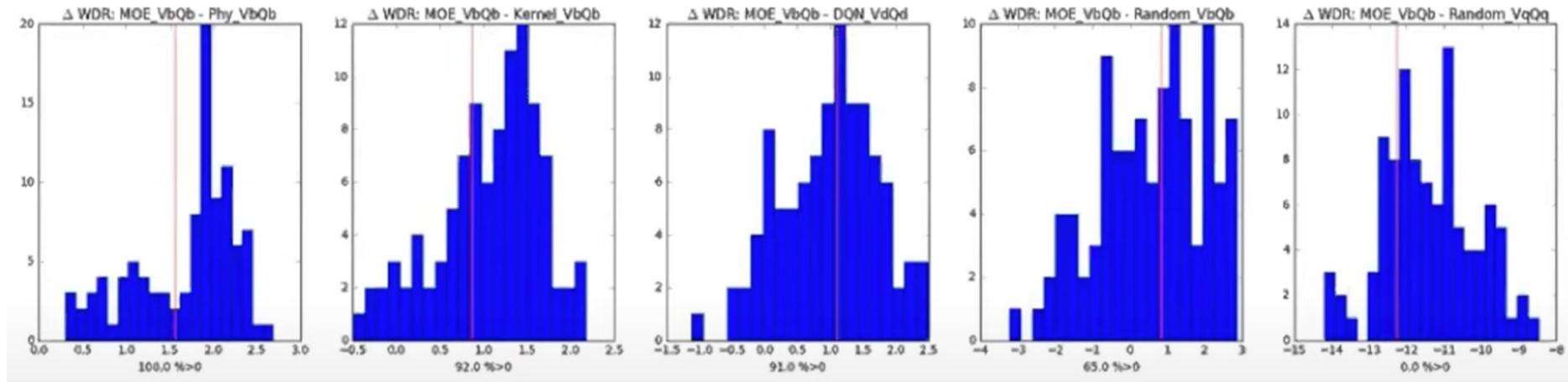
	<b>DR</b>	<b>IS</b>	<b>WIS</b>
Random	$-2.31 \pm 1.42$	$-3.48 \pm 1.36$	$-2.80 \pm 1.27$
Short-term kernel	$2.17 \pm 1.4$	$2.18 \pm 1.20$	$2.16 \pm 1.71$
Long-term kernel	$9.47 \pm 1.70$	$5.72 \pm 1.81$	$6.97 \pm 1.29$
POMDP	$6.04 \pm 2.18$	$4.15 \pm 2.28$	$6.67 \pm 1.74$
<b>Mixture-of-experts</b>	<b><math>11.83 \pm 1.26</math></b>	<b><math>12.50 \pm 1.19</math></b>	<b><math>11.07 \pm 1.21</math></b>

Swiss HIV Cohort  
Study

	<b>DR</b>	<b>IS</b>	<b>WIS</b>
Random	$-6.33 \pm 3.47$	$-5.57 \pm 2.17$	$-6.18 \pm 3.24$
ST Kernel	$1.64 \pm 1.86$	$2.03 \pm 1.81$	$2.17 \pm 1.74$
LT Kernel	$9.67 \pm 1.49$	$7.38 \pm 1.72$	$7.64 \pm 1.92$
POMDP	$5.46 \pm 2.05$	$6.72 \pm 2.88$	$7.76 \pm 2.10$
<b>Mixture-of-experts</b>	<b><math>10.73 \pm 1.02</math></b>	<b><math>13.59 \pm 1.57</math></b>	<b><math>11.83 \pm 1.31</math></b>

# Sensitivity Analysis for Parameter Choices

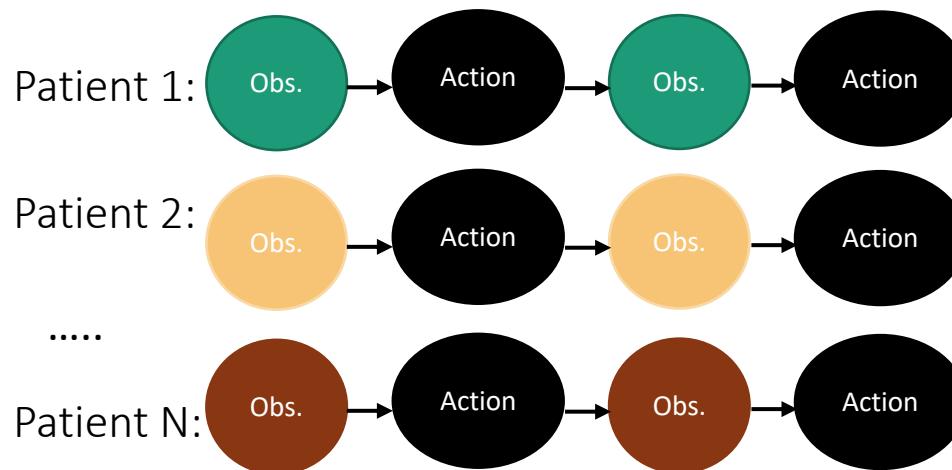
How does evaluation vary with different choices of control variates?



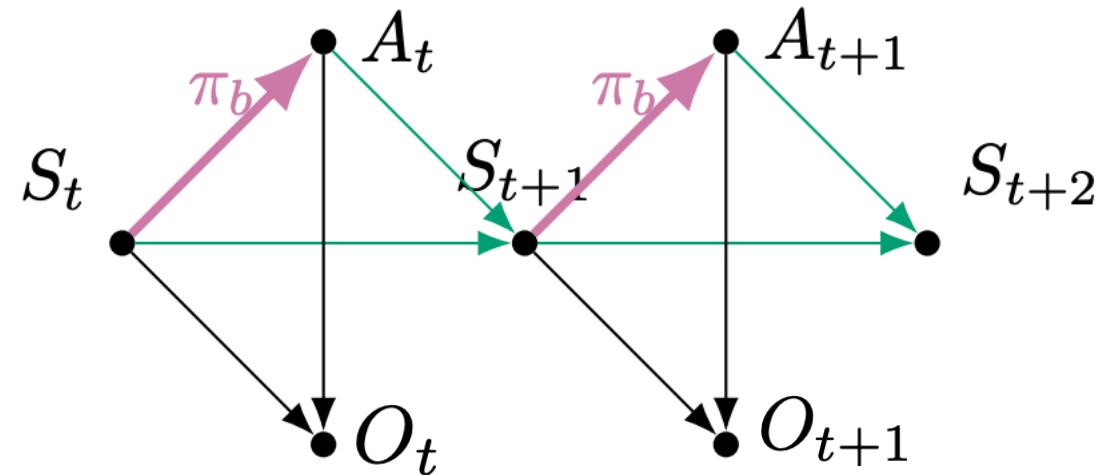
# Moving Forward: From OPE to *Causal* OPE Estimates

Reweighting trajectories according to how much they overlap using:

$$\rho_n = \prod_t \frac{\pi_e(a_{tn}|s_{tn})}{\pi_b(a_{tn}|s_{tn})}$$



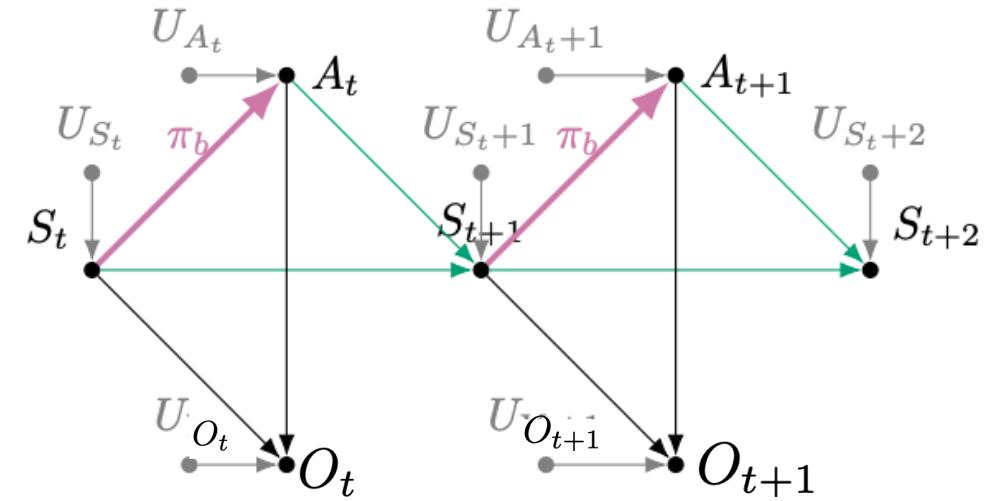
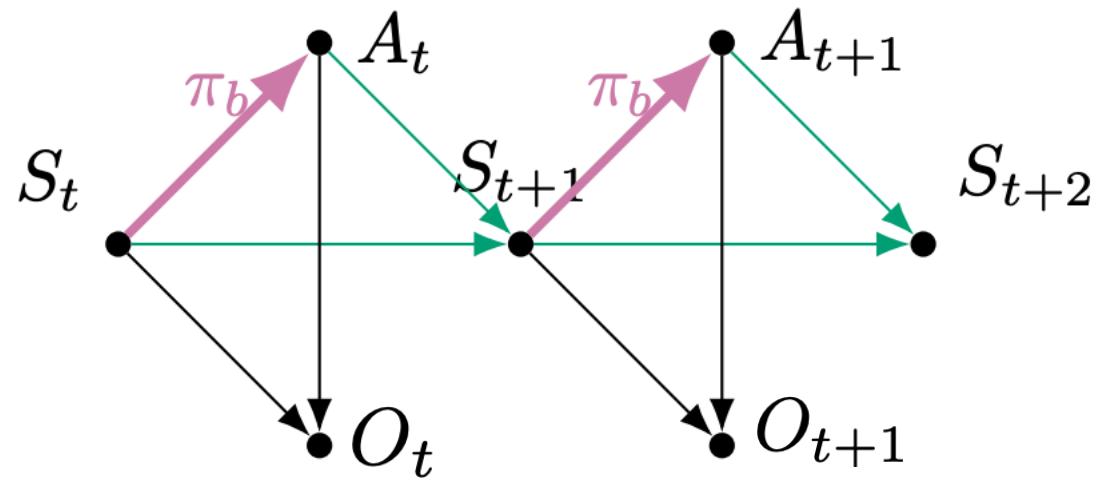
Standard OPE **ignores** the data-generating mechanism resulting in states, actions and observations.



But the data-generating mechanism is important to understand when an OPE estimate makes sense and who it is applicable to

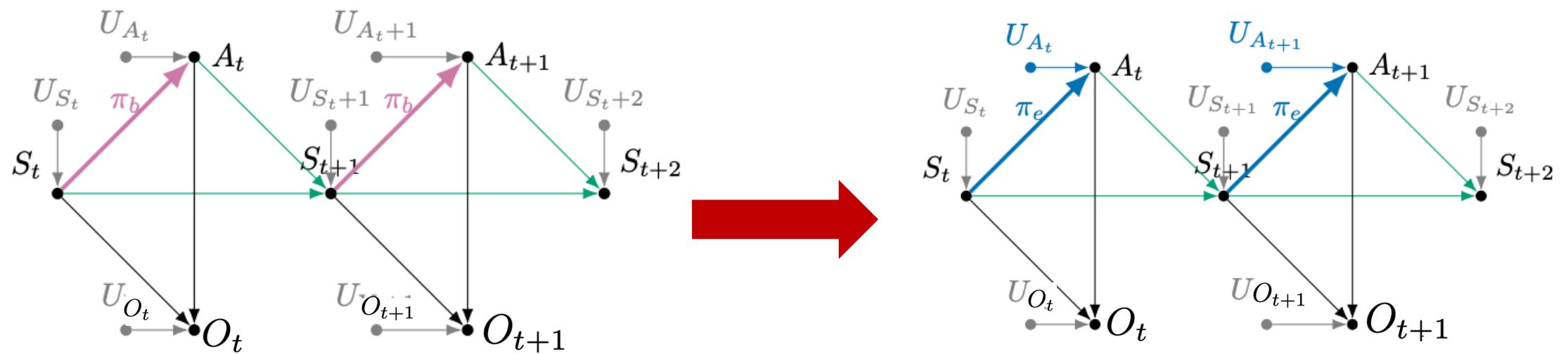
# Moving Forward: From OPE to *Causal* OPE Estimates

The data-generating mechanism is important to understand when an OPE estimate makes sense and who it is applicable to



# Moving Forward: From OPE to *Causal* OPE Estimates

If we treat OPE as a type of soft intervention then:



Depending on whether we want to compute **retrospective** or **prospective** estimates, the latent factors  $\mathbf{U}$  may play a significant role the quality of the performance.

# Moving Forward: From OPE to *Causal* OPE Estimates

Specifically, retrospective and prospective queries **require different estimands**:

	Individual-Level	Subpopulation-Level	Population-Level
Counterfactual OPE (Retrospective)	E.g. Would patient $i$ have survived had we changed the medication? $P(Y_T^{\pi_e}(U_i)   \mathcal{H}_T^{\pi_b})$	E.g. Would female patients have survived had we changed the medication? $E(Y_T^{\pi_e}(U_{female})   \mathcal{H}_T^{\pi_b})$	E.g. Would all patients have survived had we changed the medication? $E(Y_T^{\pi_e}(U)   \mathcal{H}_T^{\pi_b})$
Interventional OPE (Prospective)	E.g. Will new patients improve if we changed the medication? $P(Y_T^{\pi_e}   \mathcal{H}_T^{\pi_b})$	E.g. Will female patients improve if we changed the medication? $E(Y_T^{\pi_e}   \mathcal{H}_T^{\pi_b}, S_j = female)$	E.g. Will all patients survive if we change the medication? $E(Y_T^{\pi_e}   \mathcal{H}_T^{\pi_b})$
Canonical OPE task (Definition 1)			

# Risk estimation *matters* for Causal OPE too!

Generalizing OPE as a causal estimand highlights two challenges:

1. Identifiability from observational data
2. Estimation challenges

Difficulty in performing OPE manifests as **induced uncertainty** in the OPE process.

Depending on the nature of uncertainty, we may be able to outline where human expertise may be helpful to mitigate uncertainty or where more assumptions or data are required.

# What should responsible development and evaluation look like?

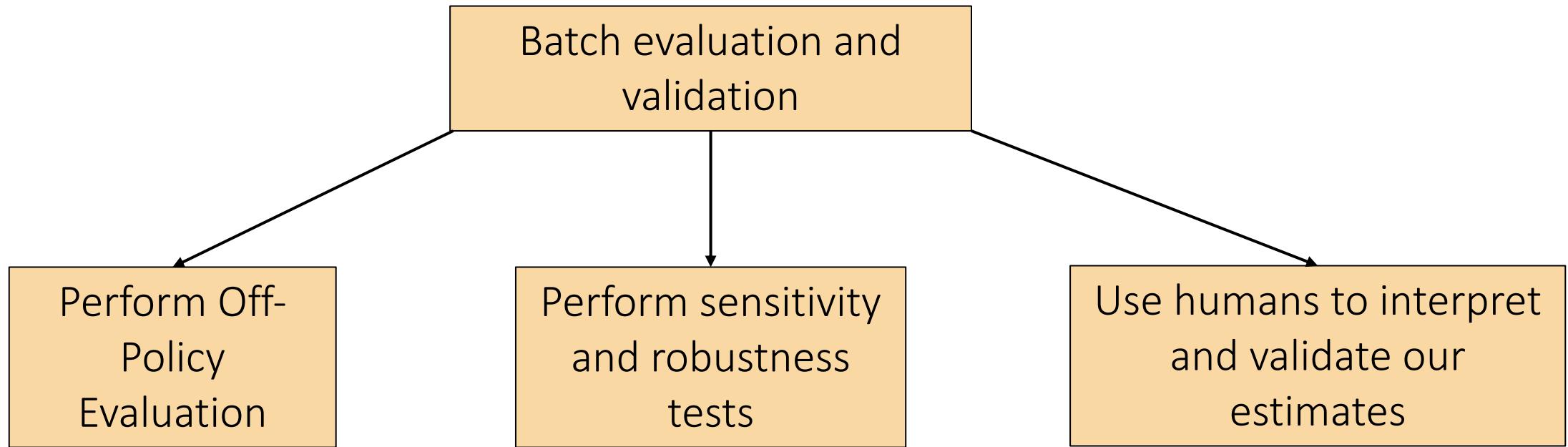
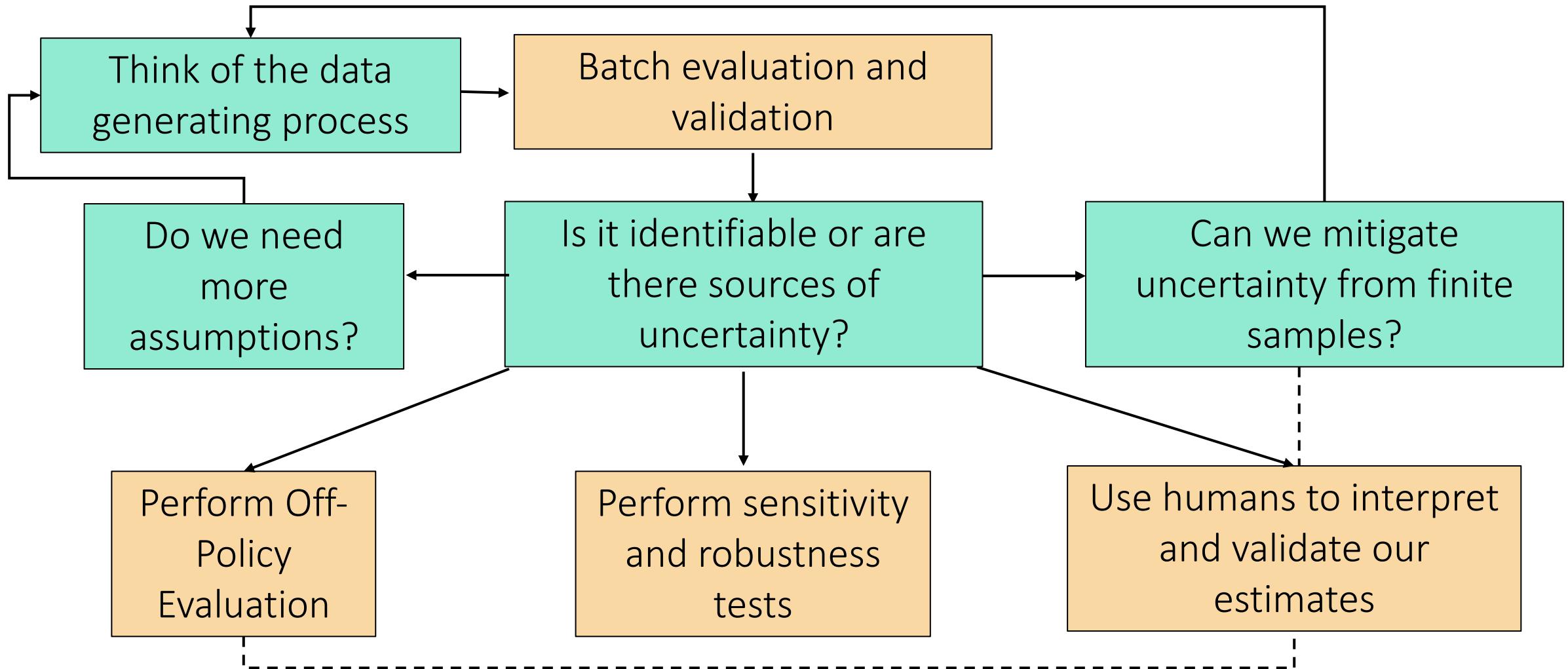
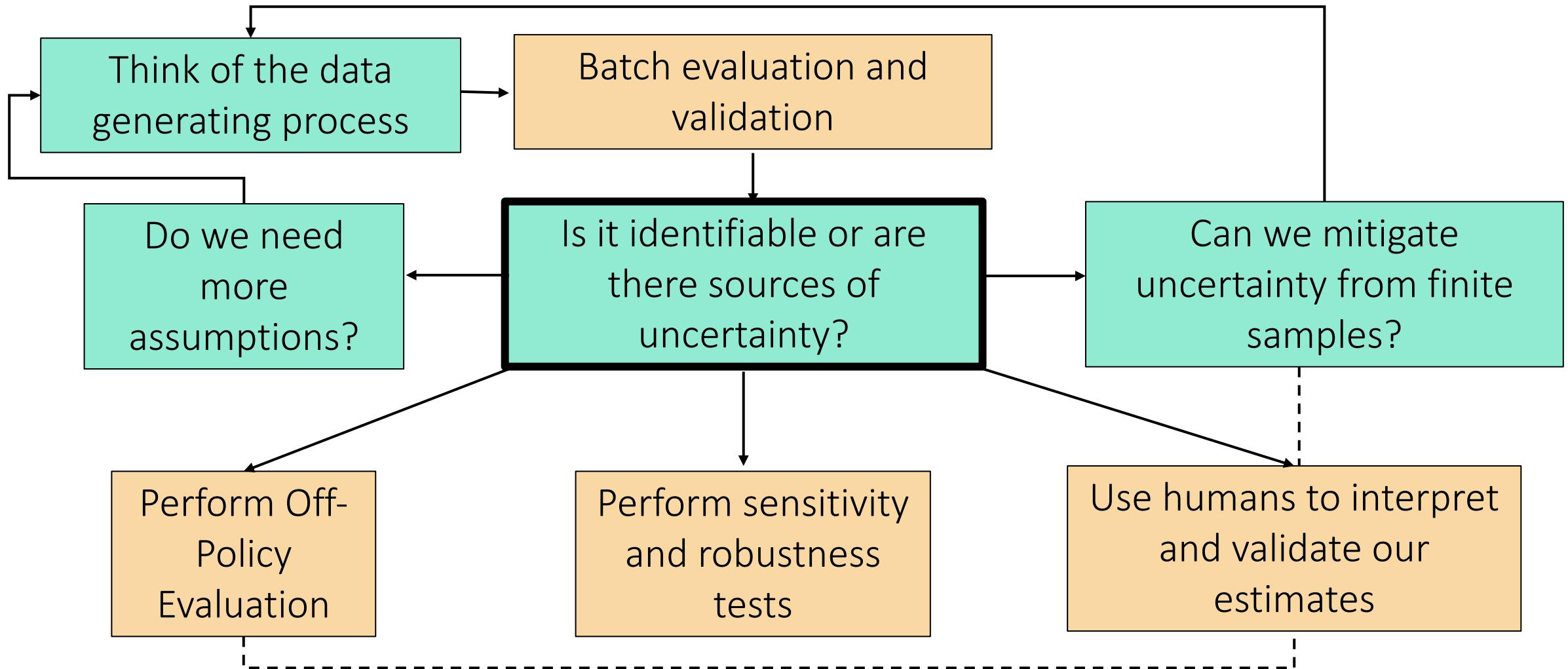


Image Credit: Doshi-Velez, Towards using Batch RL to Identifying Treatment Options, 2022

# Causal OPE for responsible development and evaluation



# Causal OPE for responsible development and evaluation



# Summary

- High stake decisions have **long-term consequences** and are often applied in exceptional situations where generalisation is hard.
- RL has a lot of potential in healthcare.
- Evaluation should **tailored to the application** and the **nature of the query** for careful validation.
- Uncertainty and risk-estimation **are key** for careful evaluation.

*Collaborators: Finale Doshi-Velez, Volker Roth, Maurizio Zazzi, Mario Wieser, Jasmina Bogojeska, Huldrych Gunthard, Omer Gottesman, Matthieu Komorowski, Aldo Faisal, David Sontag, Leo Celi, Shalmali Joshi, Taylor Killian, Marzyeh Ghassemi.*

# We're hiring.....

If you are interested in learning more about our research and engaging with us or working with our lab, please email us on

**[s.parbhoo@imperial.ac.uk](mailto:s.parbhoo@imperial.ac.uk)**