



CSARIM 2025

Transfusion guidance based on reinforcement learning

Prof. Dr. Jens Meier

Clinic of Anesthesiology and Intensive Care Medicine Kepler University Clinic Linz, Austria

Conflict of interest





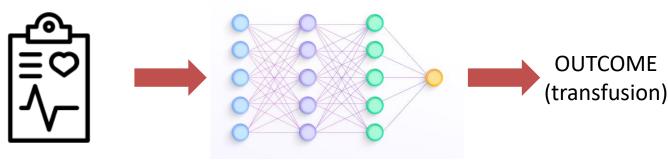


What we already know...









clinical data

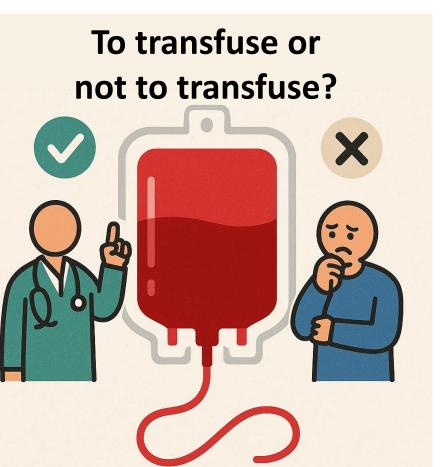
algorithm mimics supervisor

The elephant in the room









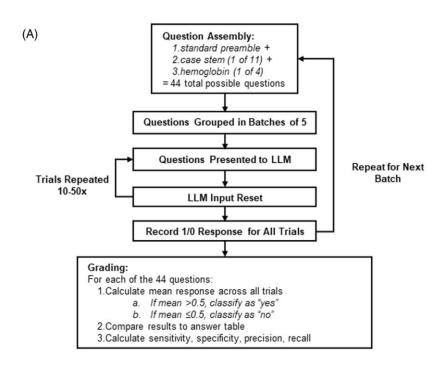
The solution everyone is talking about...







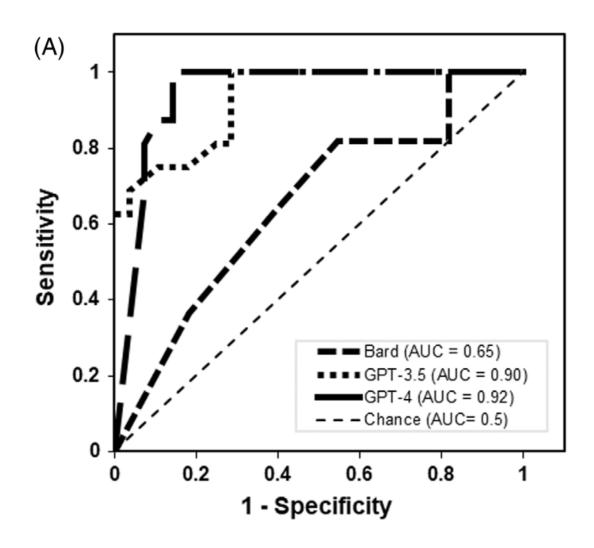




LLMs "understand" transfusion indications



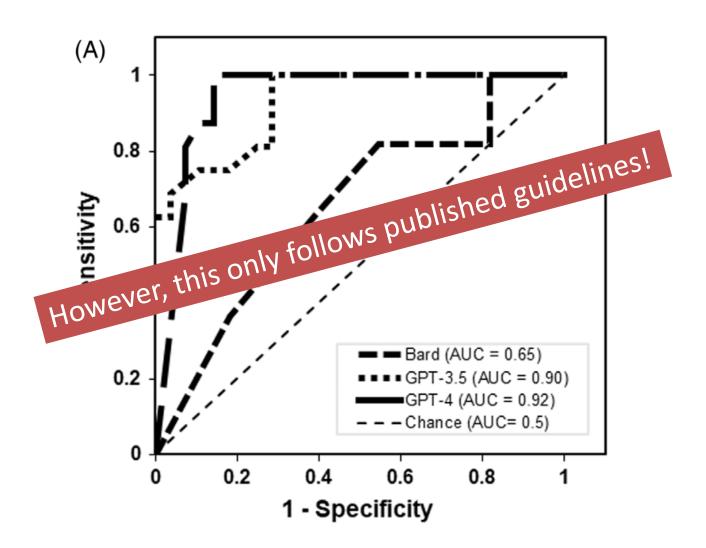




LLMs "understand" transfusion indications



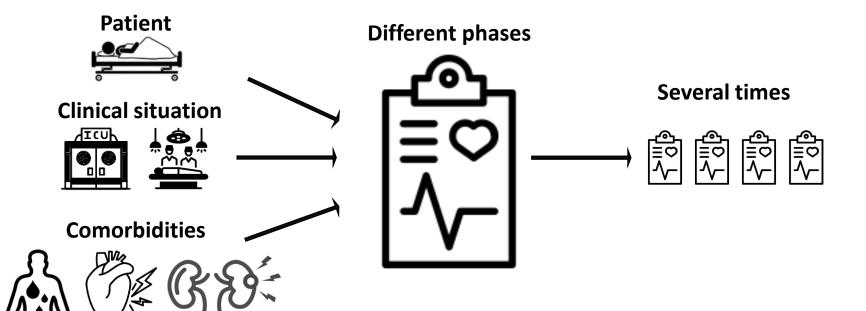








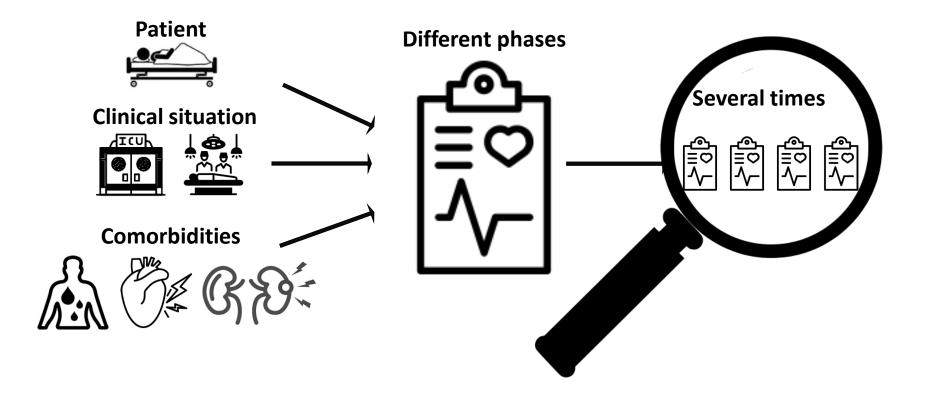








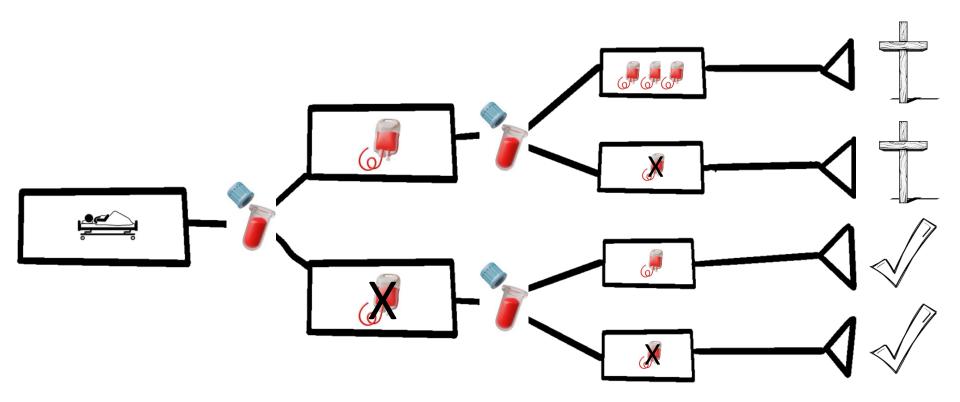




In daily clinical practice: We have to decide often!



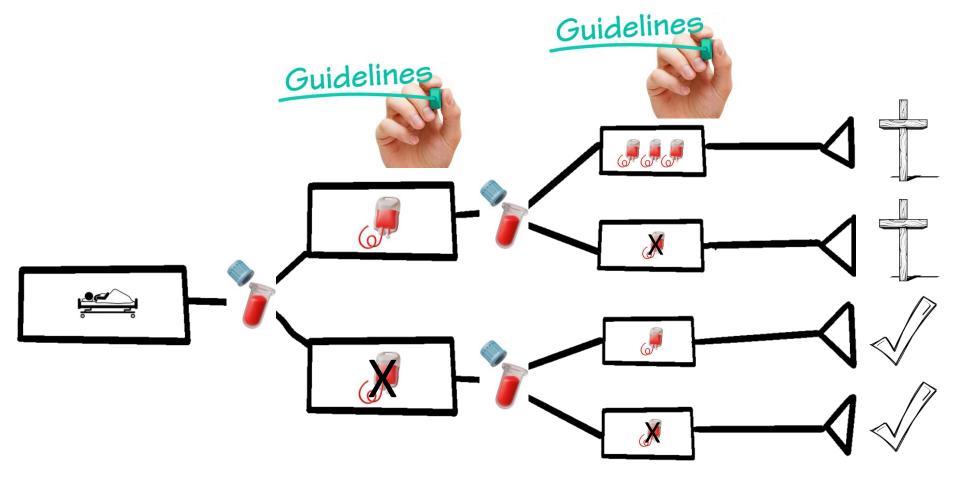




In daily clinical practice: We have to decide often!



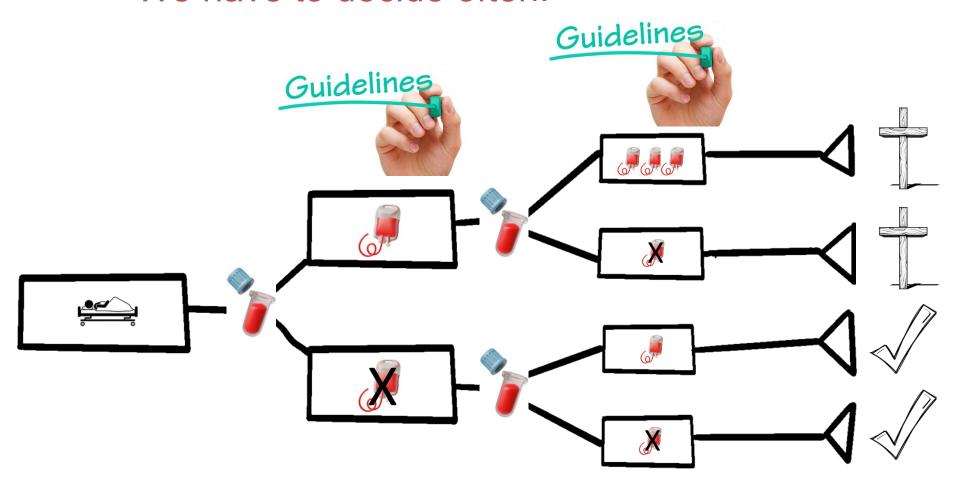




In daily clinical practice: We have to decide often!





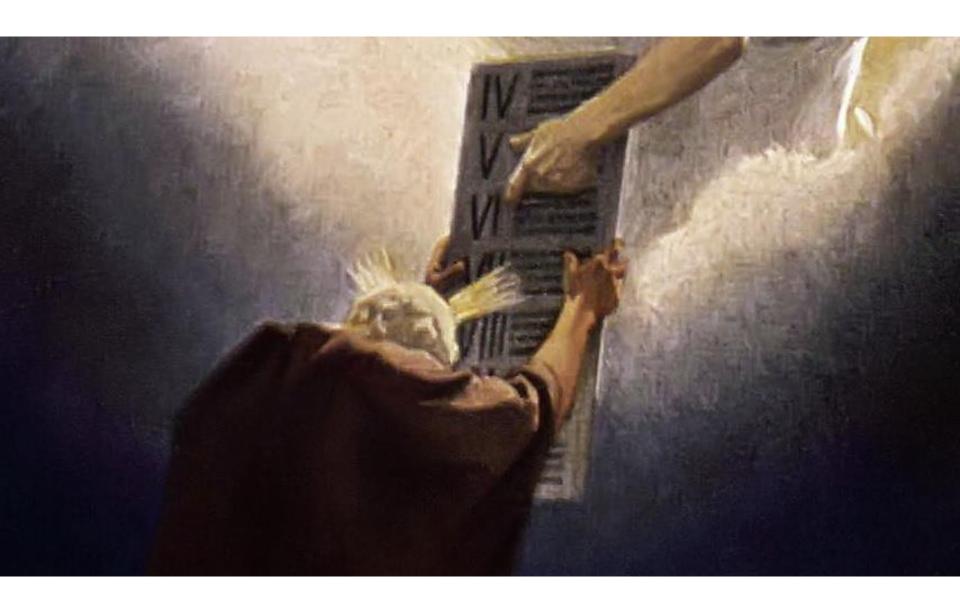


- 1. Same recommendation ▶ various effects!
- 2. Guidelines are a coarse instrument

The role of these guidelines...







How do we get these guidelines?







Shortcoming of this approach



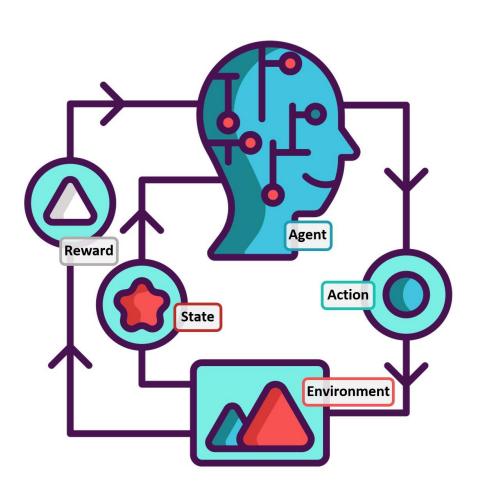




The new paradigm Reinforcement Learning!



















State hungry, somewhere



ActionCrawling, sitting



Reward Feeder











State hungry, somewhere



Action Crawling, sitting



Reward Feeder











Environment baby's room



State hungry, somewhere



Action
Crawling, sitting



Reward Feeder











Environment baby's room



State hungry, somewhere



ActionCrawling, sitting



Reward Feeder











Environment baby's room



State hungry, somewhere



ActionCrawling, sitting



Reward Feeder











Environment baby's room



State hungry, somewhere



ActionCrawling, sitting



Reward Feeder





















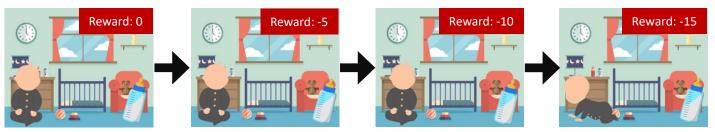




Reward lg Feeder

















Environment baby's room



State hungry, somewhere



ActionCrawling, sitting



Reward Feeder











Environment baby's room



State hungry, somewhere



ActionCrawling, sitting



Reward Feeder











Environment baby's room



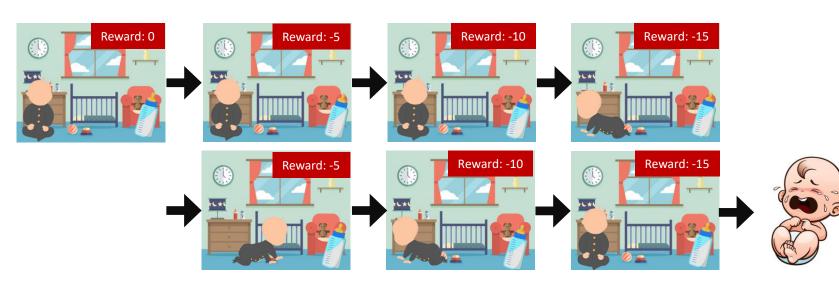
State hungry, somewhere



ActionCrawling, sitting



Reward Feeder











Environment baby's room



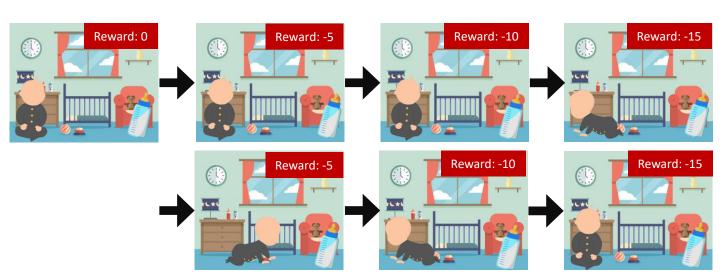
State hungry, somewhere



ActionCrawling, sitting



Feeder













Environment baby's room

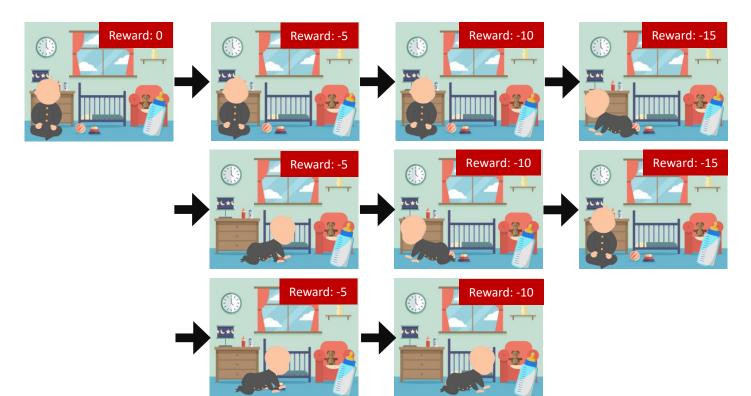


State hungry, somewhere



ActionCrawling, sitting













nvironment
baby's room



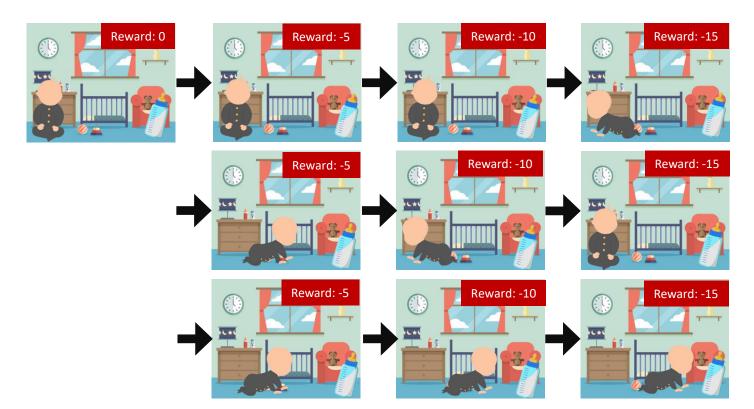
State hungry, somewhere



Action
Crawling, sitting



Feeder











Environment baby's room



State hungry, somewhere



ActionCrawling, sitting



Reward: -10
Reward: -10
Reward: -10
Reward: -15
Reward: -10
Reward: -15















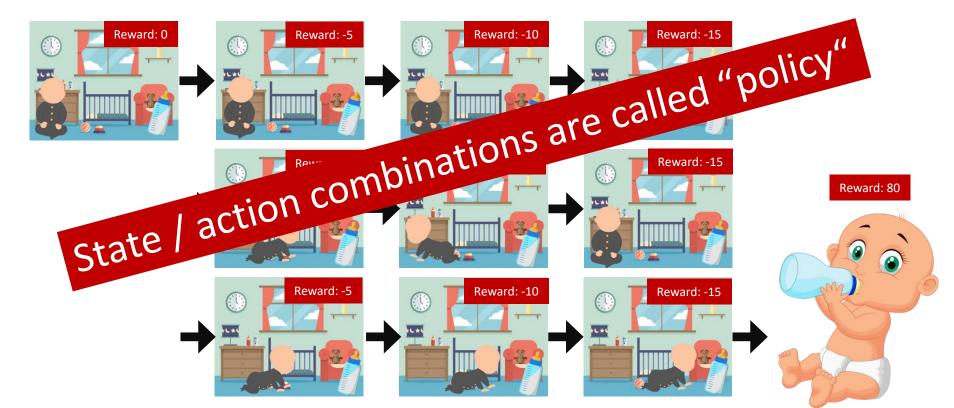
State hungry, somewhere



Action Crawling, sitting



Reward Feeder



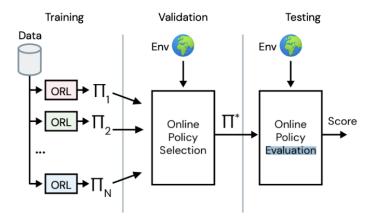
The difference between online and offline reinforcement learning







Online reinforcement learning

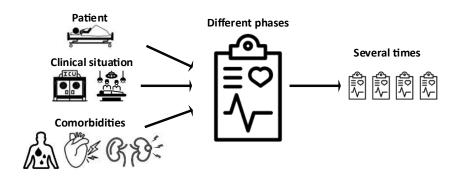


The difference between online and offline reinforcement learning

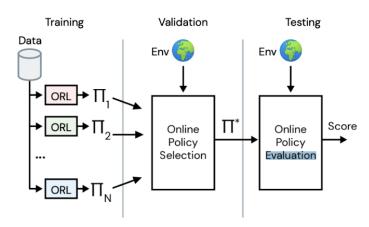




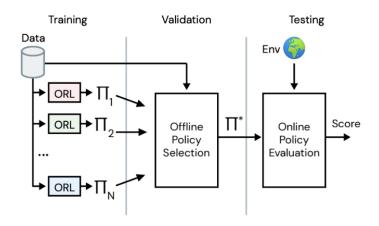




Online reinforcement learning



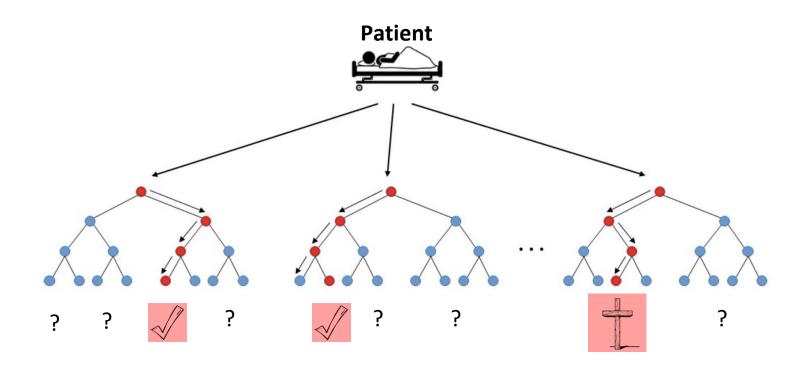
Offline reinforcement learning



Can you evaluate theoretical policies on historical data?





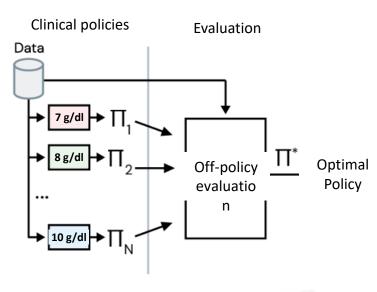


- Trajectory exists
- Trajectory does not exist

Can you evaluate theoretical policies on historical data?





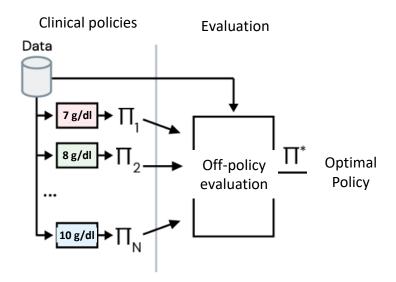




Can you evaluate theoretical policies on historical data?









Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI-18)

Importance Sampling for Fair Policy Selection*

Shayan Doroudi^{1,3}, Philip S. Thomas² and Emma Brunskill³

¹ Carnegie Mellon University

² University of Massachusetts Amherst ³ Stanford University

shayand@cs.cmu.edu, pthomas@cs.umass.edu, ebrun@cs.stanford.edu

$$\hat{V}_{\mathrm{IS}}^{\pi_e} \triangleq \frac{1}{n} \sum_{i=1}^{n} w_i \sum_{t=1}^{T_i} R_{i,t}$$

$$w_i \triangleq \prod_{t=1}^{T_i} \frac{\pi_e(a_{i,t}|\tau_{i,1:t-1})}{\pi_b(a_{i,t}|\tau_{i,1:t-1})}$$

$$\hat{V}_{\text{PHWIS}} = \sum_{l \in L} W_l \underbrace{\frac{1}{\sum_{\{i|T_i = l\}} w_i} \sum_{\{i|T_i = l\}} w_i \sum_{t=1}^{T_i} R_{i,t}}_{}$$

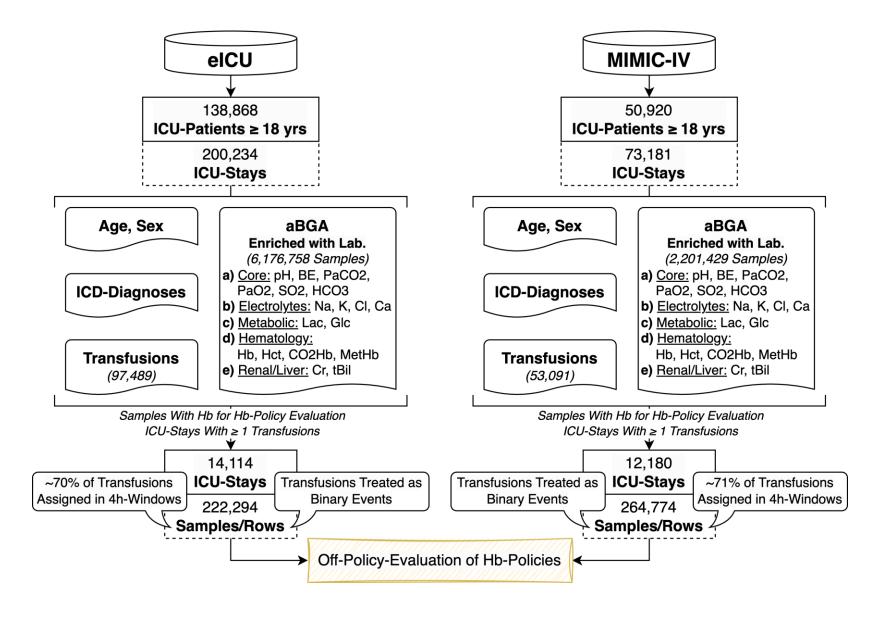
WIS estimate on *l*-length trajectories

per weighted horizon importance sampling: PWHIS

Off-policy evaluation of Hb policies $J \cong U$



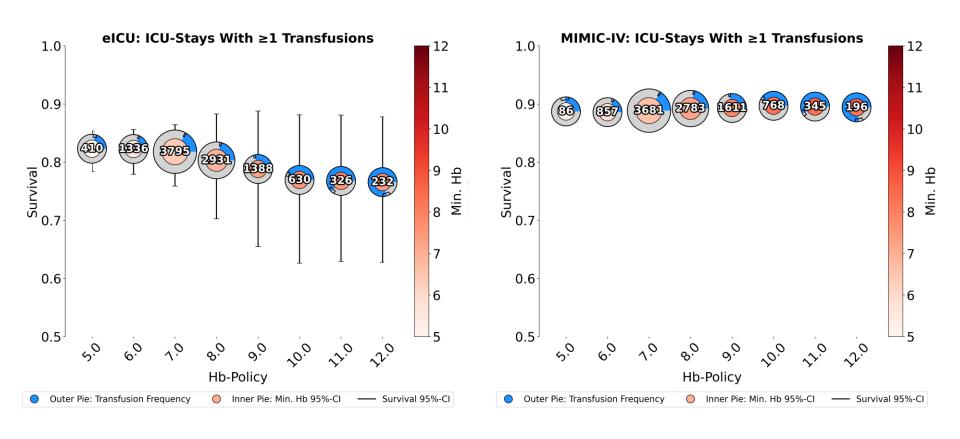




Hb policies evaluated for all patients



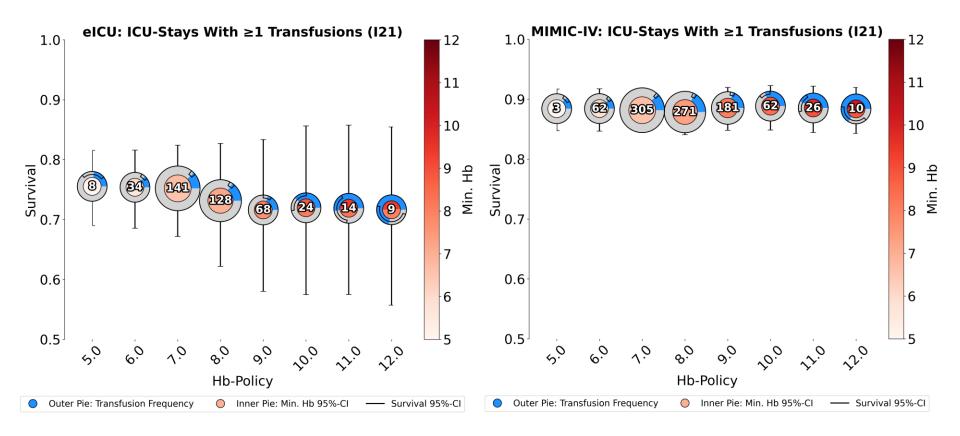




Hb policies evaluated for myocardial infarction patients







Can you improve this with reinforcement learning?





$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma max_aQ(S_{t+1}, a) - Q(S_t, A_t)]$$

New **Q**-value estimation

O-value estimation

Reward Rate

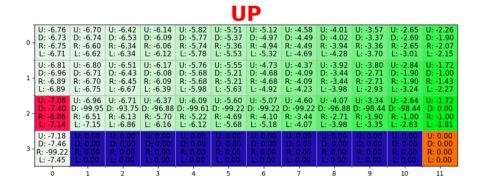
Former Learning Immediate Discounted Estimate optimal Q-value of next state

Former O-value estimation

TD Target

TD Error

$$\pi^*(s) \leftarrow \underset{a}{\operatorname{argmax}} Q^{\pi^*}(s, a) \ \forall s$$

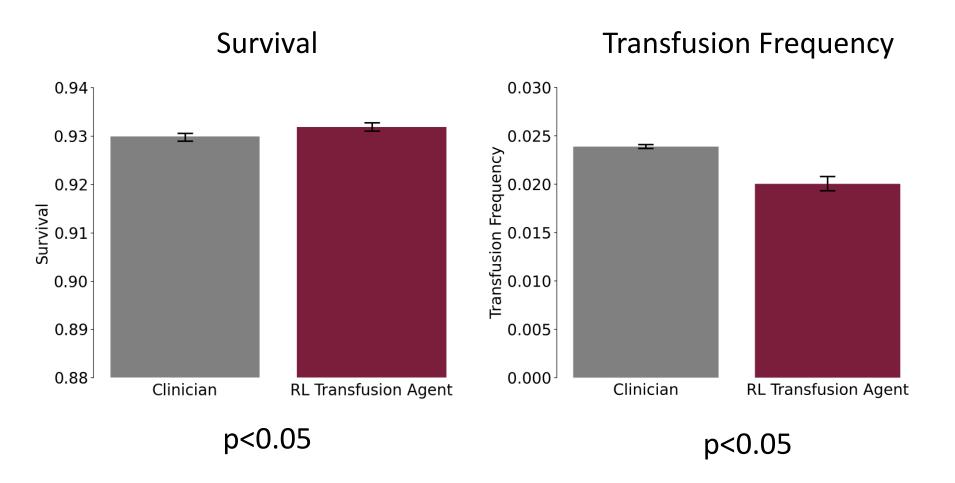




Evaluating Policy after Q-learning of the MIMIC IV data



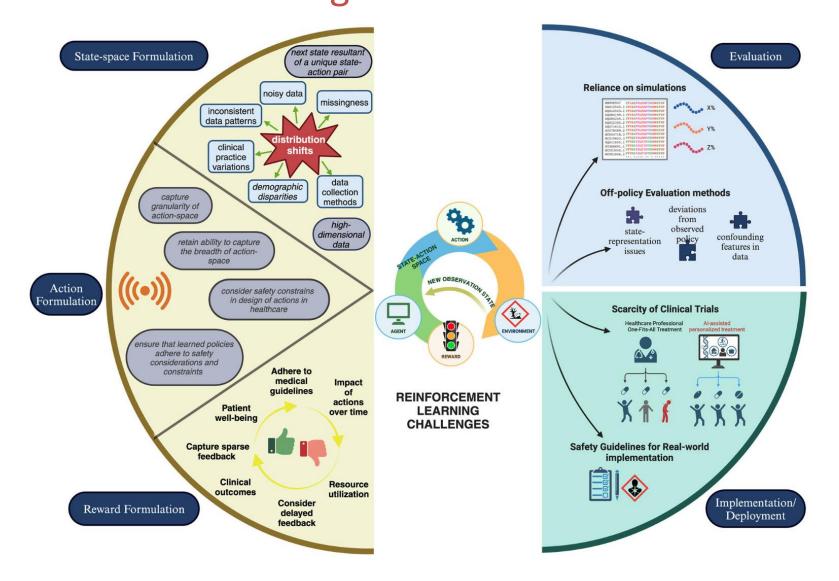




The problems of reinforcement learning







Conclusion





- ML models are well established for transfusion prediction
- LLMs might be able to guide transfusion in the future
- Reinforcement learning is the way to go!
- Hb based transfusion policies have (probably) no effect on survival
- one can build a reinforcement learning model, that
 - transfuses less
 - has no influence on survival

