

An aerial, grayscale photograph of Liverpool, UK, showing the city's dense urban landscape, the Mersey River, and the suspension bridge in the distance. The image is semi-transparent, allowing text to be overlaid.

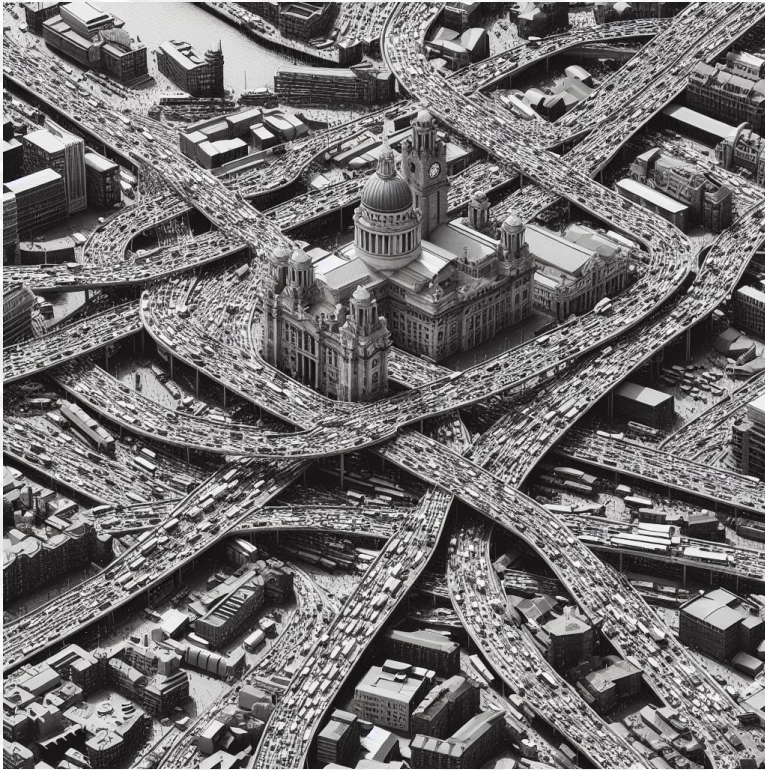
Bayesian Reinforcement Learning (or *Posterior Sampling for RL*)

University of Liverpool

Distributed Algorithms CDT Showcase 2023

Malcolm Strens

AI is moving towards *multi-step decision-making* under uncertainty...



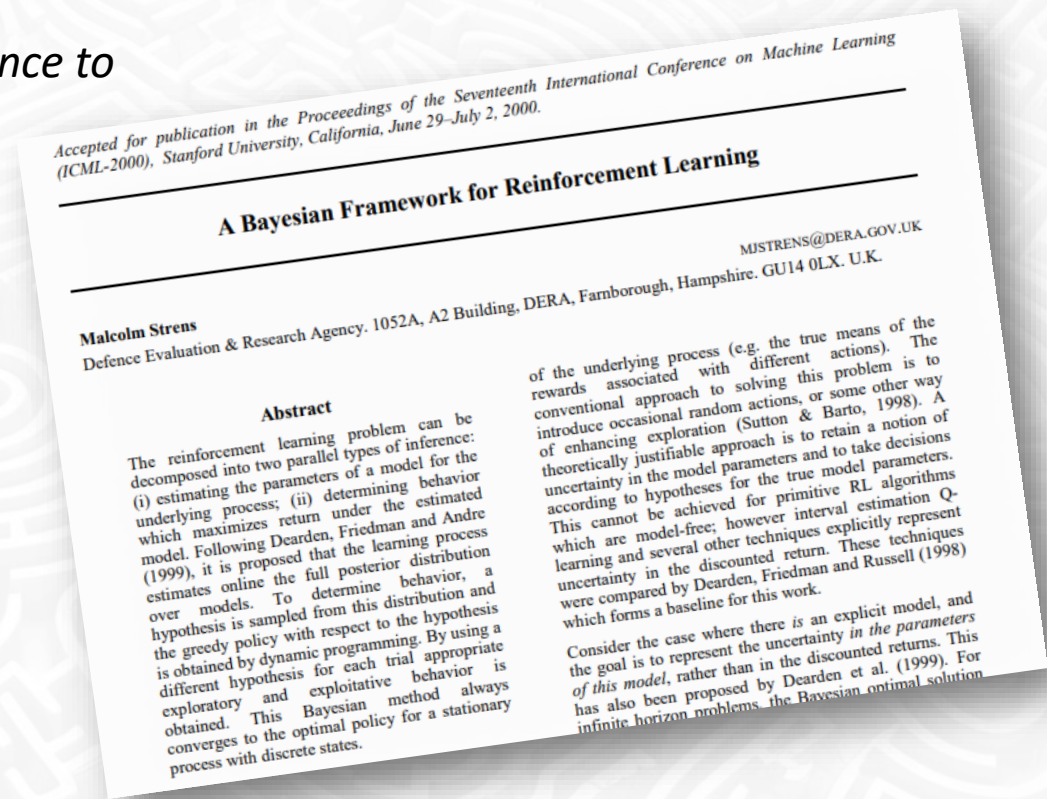
“Sequential Decision Problems”

“Optimal Control”

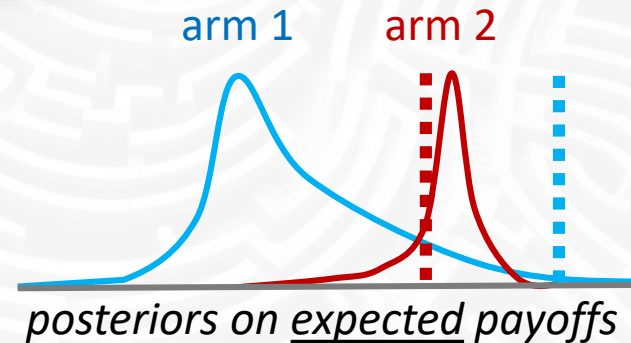
“Reinforcement Learning”

Q. Can we apply Bayesian Inference to sequential decision-making?

A. Yes, if we can express uncertainty over RL models (of interaction and delayed reward).

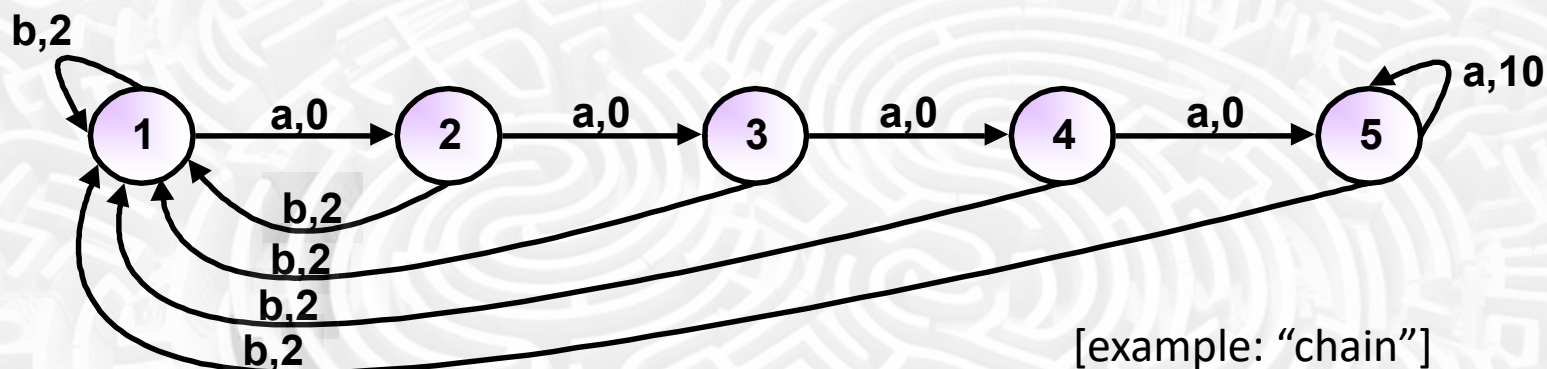


Builds on Thompson Sampling (1933):
for repeated state-free decision
e.g. 2-armed bandit



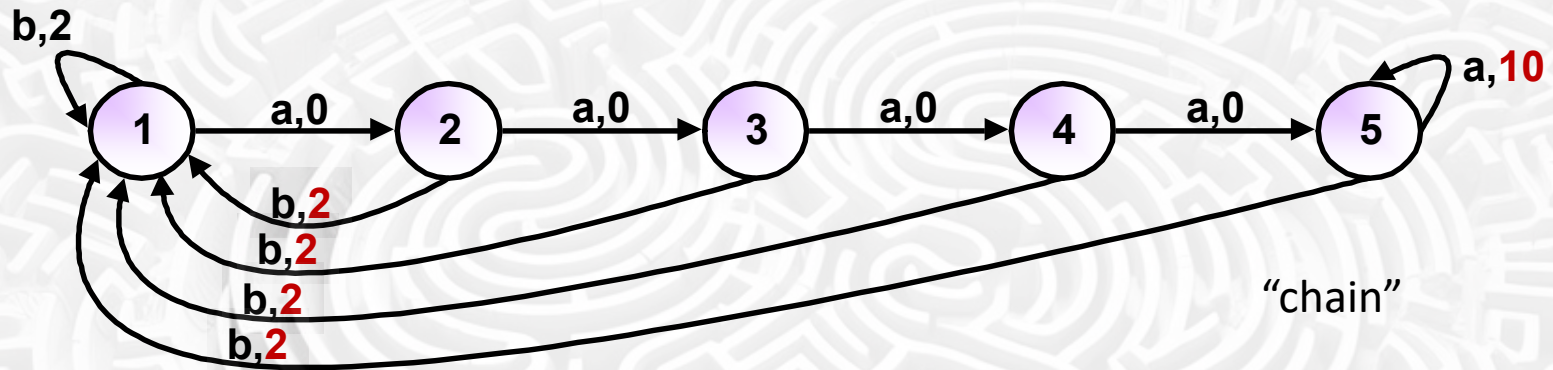
Optimal: “select the best action under a single sample from the Bayesian posterior”
i.e. $P(E \text{ arm 1 payoff} > E \text{ arm 2 payoff} \mid \text{trials so far})$

RL target for inference: Markov Decision Process



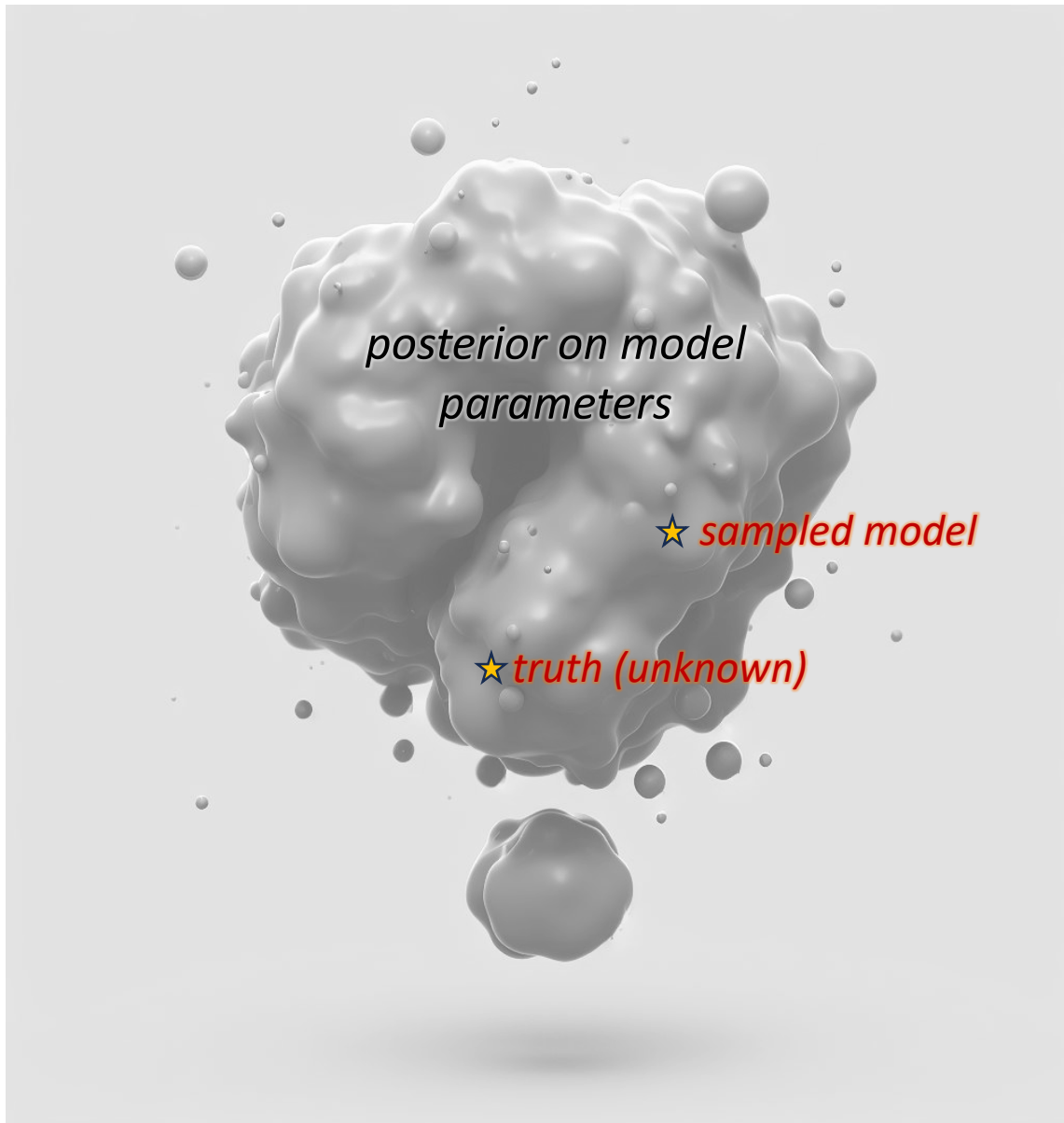
- (S, A, T, R)
- S : Set of States
- A : Set of Actions
- T : Transition Probabilities $T(s, a, s')$ Multinomial posterior (Dirichlet conjugate prior)
- R : Reward Distributions $R(s, a)$ Gaussian posterior (Normal-Gamma conjugate prior)

Exploration/exploitation dilemma



- Trade-off 'value of information' from exploration vs regret of not exploiting information already known.
- A generalisation of Thompson sampling to multi-step problems achieves this.

Bayesian RL on “models of interaction” like MDPs



Start a new interaction (“episode”).

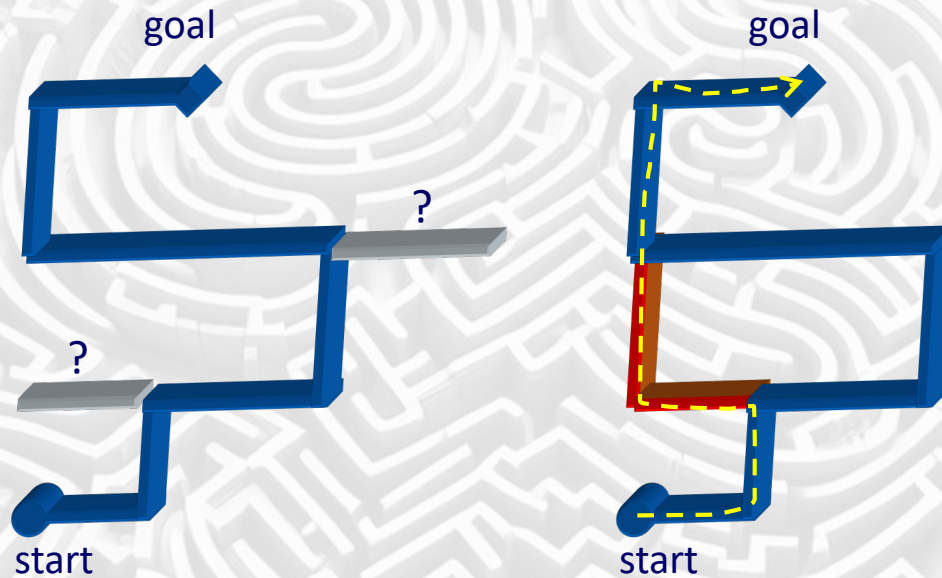
Sample a model for the world (e.g. MDP) from current posterior.

Solve the world model for optimal behaviour; e.g. Bellman backup to obtain control “policy”.

Take actions according to this policy until the end of the episode.

Use the collected transitions & rewards to update the posterior.

e.g. navigation while learning map



posterior represents
uncertainty in map
based on experience

drawing a **hypothesis** from the posterior,
then solve for shortest path ... yields
exploratory behaviour!

(then update posterior)

Developments/applications in PSRL

- 2000: formulated for RL
- 2005: modelling cognition [Daw/Dayan]
- 2011+: medical applications
- 2012: MCTS for large models/games [Guez/Silver/Dayan]
- 2013 & 2017: regret bounds & outperformance ... best for any RL algorithm [Osband/Russo/Van Roy]
- 2019: multi-agent
- 2021: partially observable tasks
- 2023: Langevin Thompson Sampling (MCMC)
- 2024: complex neural world models?



Questions?

Optimal Policy from Dynamic Programming

- Can apply to:
 - true MDP (not known during learning)
 - max likelihood MDP (changes during learning)
 - **hypothesised MDP**
- DP backup on an estimated MDP:

$$\hat{Q}(s, a) \leftarrow \hat{R}(s, a) + \gamma \sum_{s' \in S} \hat{T}(s, a, s') \cdot \max_{a'} \hat{Q}(s', a')$$

Expected
discounted
reward for a
in s

Expected
immediate
reward for a
in s

Sum over possible
transitions
(discounted)

Best expected
discounted reward
in successor state