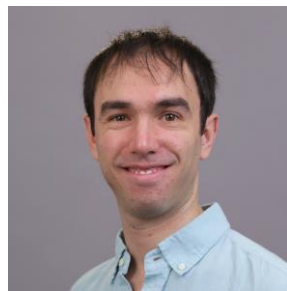


# Counterfactual Off-Policy Evaluation with Gumbel-Max Structural Causal Models



**Michael Oberst**  
MIT



**David Sontag**  
MIT



# Motivation: Building trust in RL policies

- ▶ **Goal:** Apply reinforcement learning in high risk settings (e.g., healthcare)
- ▶ **Problem:** How to safely evaluate a policy? No simulator, and off-policy evaluation can fail due to
  - ▶ Confounding
  - ▶ Small sample sizes
  - ▶ Poorly specified rewards
- ▶ Could try to interpret the policy directly, but if not possible, what can we do?



# Motivation: Building trust in RL policies

Suppose we are given:

- Markov Decision Process (MDP)
- Policy (e.g., learned using MDP)



**Observational Data**



**Markov Decision Process (MDP)**

$$P(S', R | S, A)$$

$S$ : Current State

$A$ : Action

$R$ : Reward

$S'$ : Next State

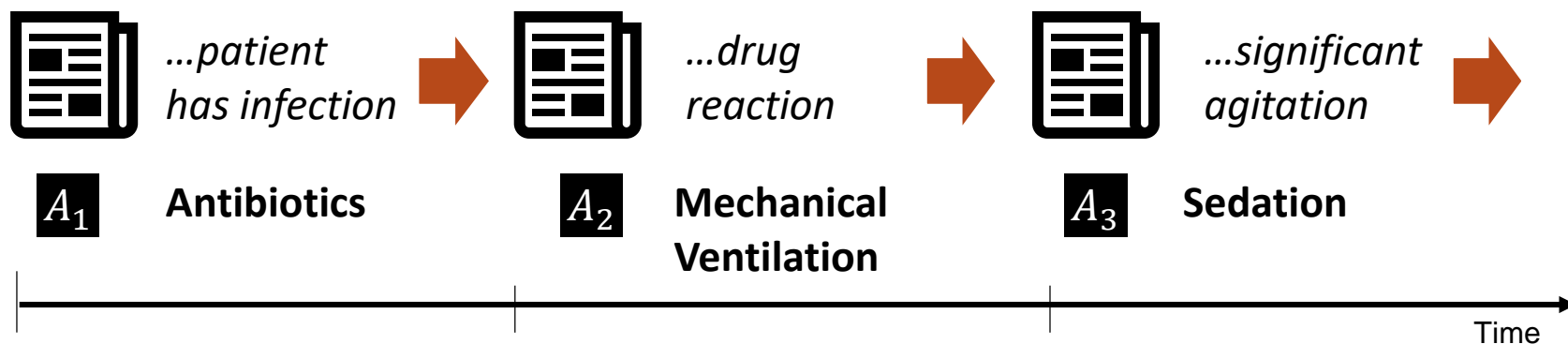


**Policy**

$$\pi(A | S)$$

# Using counterfactuals to “sanity check”

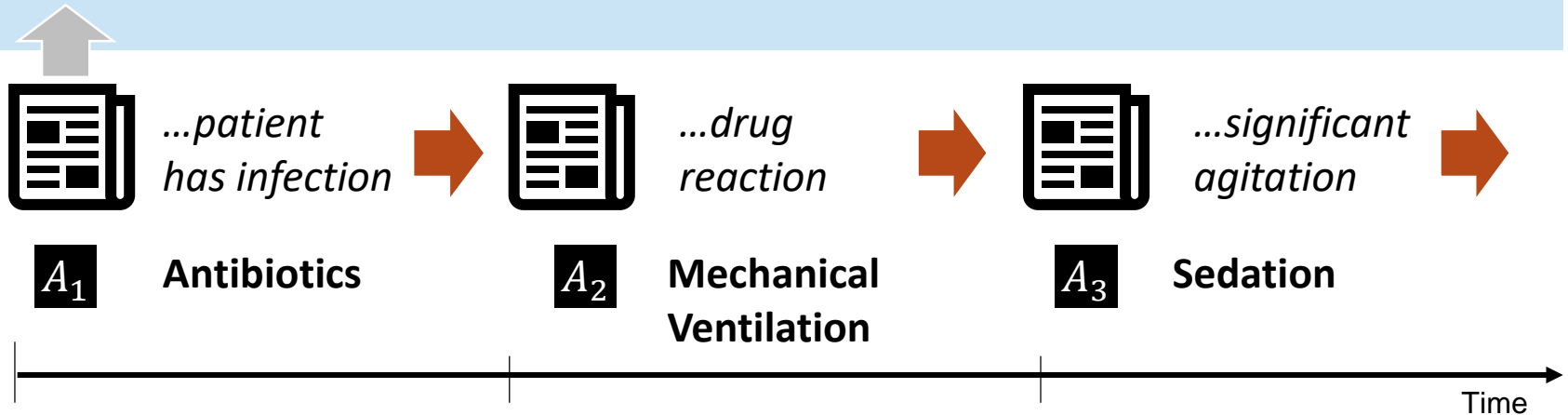
S: State  
A: Action



# Using counterfactuals to “sanity check”

*If the new policy **had been** applied to this patient...*

S: State  
A: Action



# Using counterfactuals to “sanity check”

If the new policy **had been** applied to this patient...

S: State  
A: Action

**A<sub>1</sub>** Antibiotics

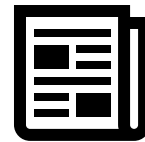
**S<sub>0</sub>** ...patient has infection



...patient has infection



...drug reaction



...significant agitation



**A<sub>1</sub>** Antibiotics

**A<sub>2</sub>** Mechanical Ventilation

**A<sub>3</sub>** Sedation

Time

# Using counterfactuals to “sanity check”

If the new policy **had been** applied to this patient...

S: State  
A: Action

$A_1$

Antibiotics

$S_0$

...patient  
has infection



$S_1$

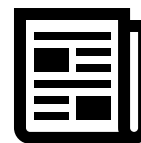
...infection  
cleared



...patient  
has infection



...drug  
reaction



...significant  
agitation



$A_1$

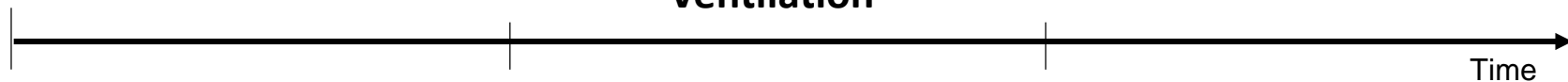
Antibiotics

$A_2$

Mechanical  
Ventilation

$A_3$

Sedation



# Using counterfactuals to “sanity check”

If the new policy **had been** applied to this patient...

S: State  
A: Action

$A_1$

Antibiotics

$S_0$

...patient has infection



$S_1$

...infection cleared

Model-based rollout  
not a fair comparison



...patient has infection



...drug reaction



...significant agitation



$A_1$

Antibiotics

$A_2$

Mechanical  
Ventilation

$A_3$

Sedation

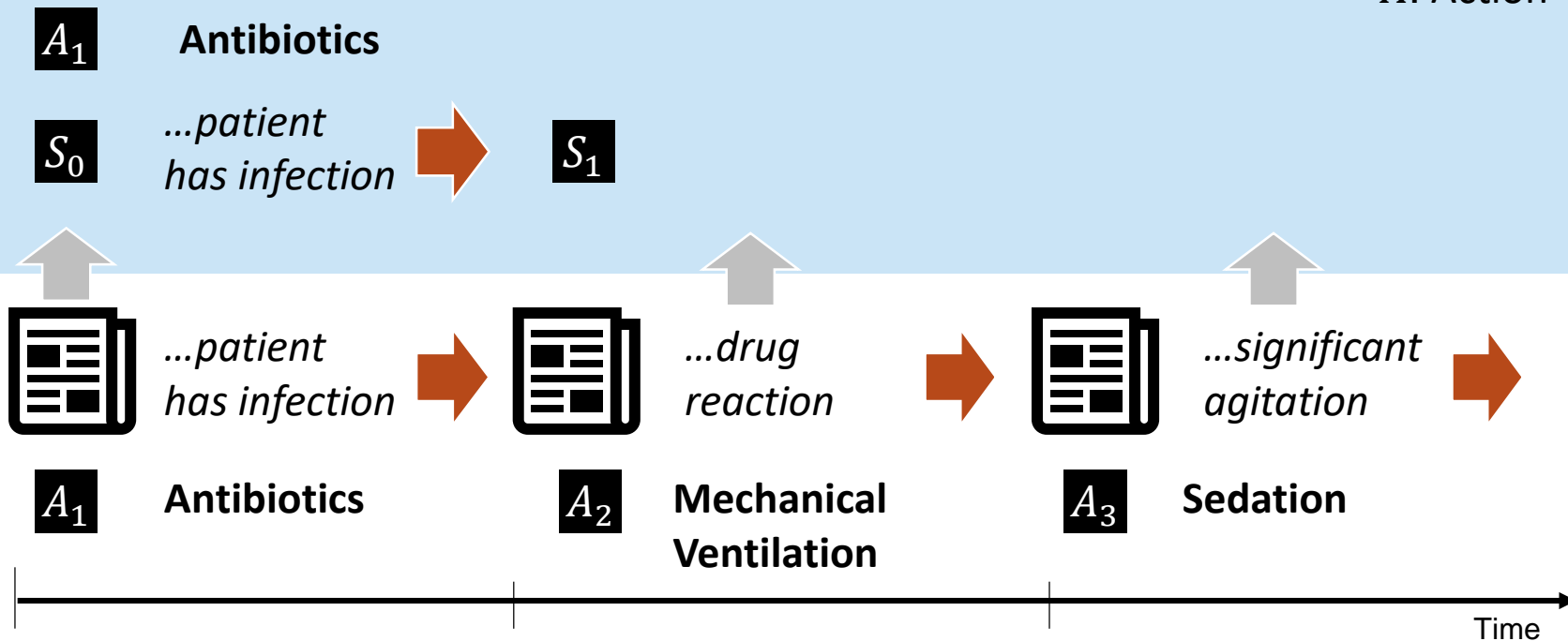
Time



# Using counterfactuals to “sanity check”

If the new policy **had been** applied to this patient...

S: State  
A: Action



# Using counterfactuals to “sanity check”

If the new policy **had been** applied to this patient...

S: State  
A: Action

**A<sub>1</sub>** Antibiotics

**S<sub>0</sub>** ...patient has infection

**S<sub>1</sub>** ...drug reaction

Counterfactual influenced by actual outcome

**A<sub>1</sub>** Antibiotics  
...patient has infection

**A<sub>2</sub>** Mechanical Ventilation  
...drug reaction

**A<sub>3</sub>** Sedation  
...significant agitation

**A<sub>1</sub>** Antibiotics

**A<sub>2</sub>** Mechanical Ventilation

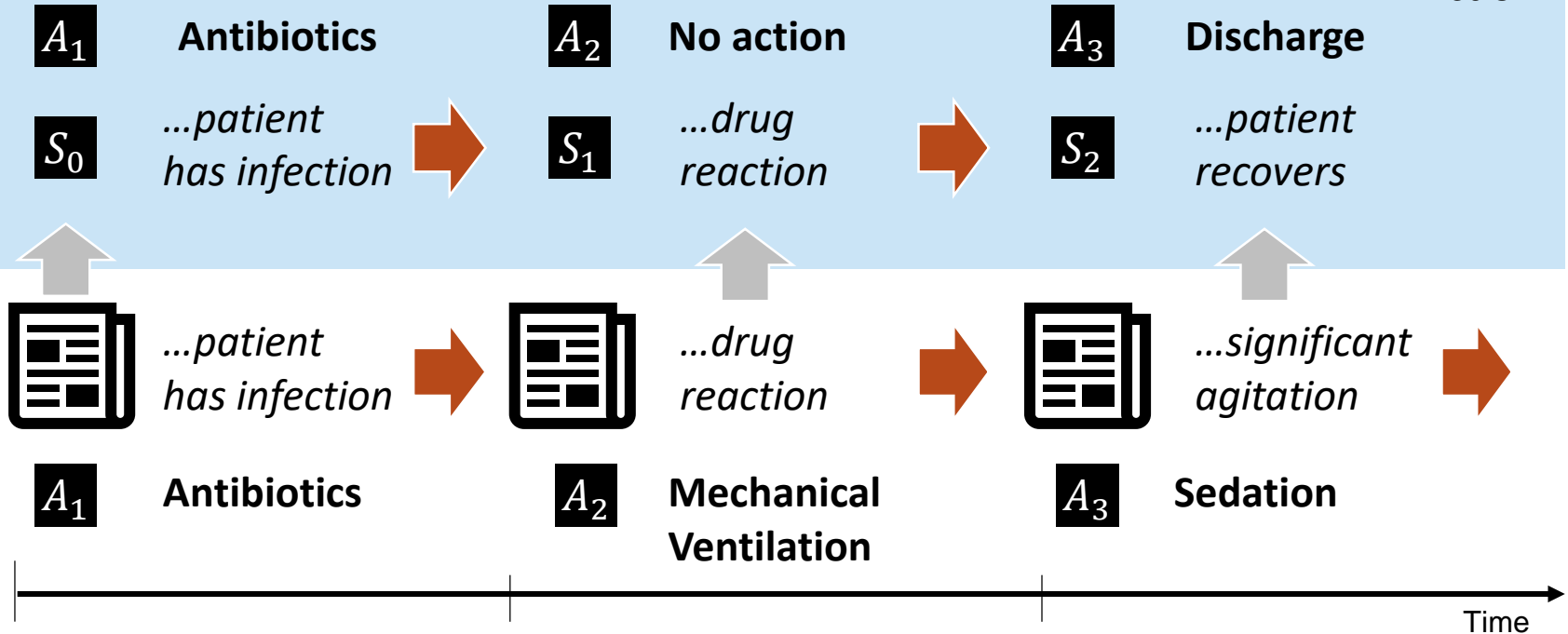
**A<sub>3</sub>** Sedation

Time

# Using counterfactuals to “sanity check”

If the new policy **had been** applied to this patient...

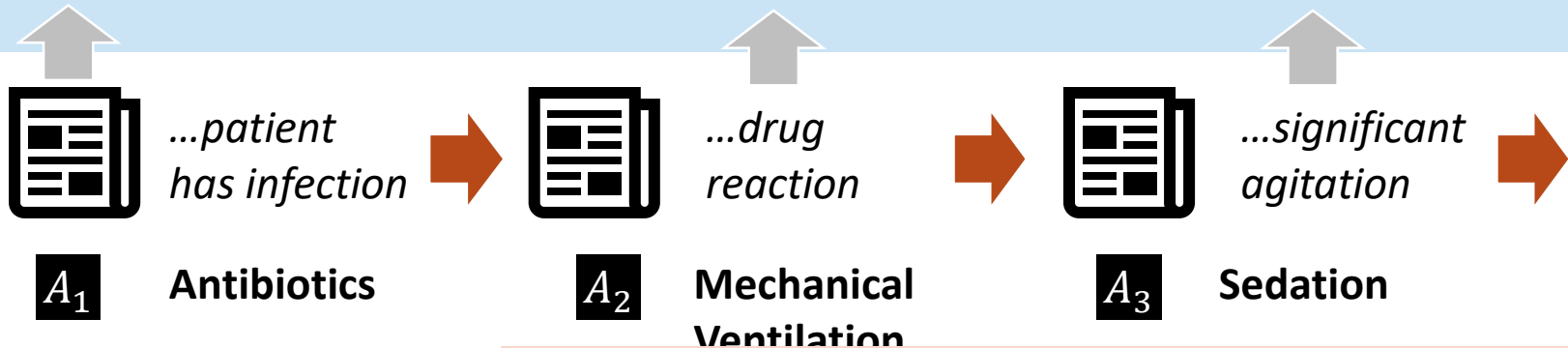
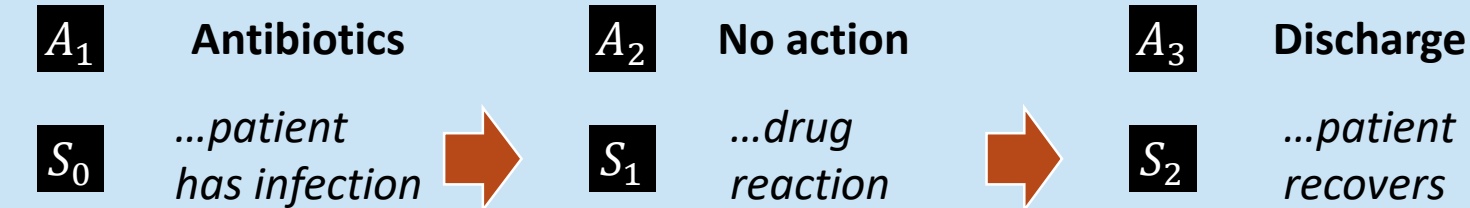
S: State  
A: Action



# Using counterfactuals to “sanity check”

If the new policy **had been** applied to this patient...

S: State  
A: Action



**Idea:** If the counterfactual trajectory is unreasonable given full context of patient, the model / policy may be flawed

# Using counterfactuals to “sanity check”

## Approach

---

- 1** **Decomposition of reward**  
over real episodes, to  
identify interesting cases

See paper / poster for synthetic case study  
motivated by sepsis management

# Using counterfactuals to “sanity check”

## Approach

- 1 **Decomposition of reward** over real episodes, to identify interesting cases

## Example

Observed Outcome	Died	1%	0%	10%
	Disch. No Chg.	0%	1%	43%
	Disch.	0%	0%	44%
		Died	No Chg.	Disch.
		Counterfactual Outcome		

See paper / poster for synthetic case study motivated by sepsis management

# Using counterfactuals to “sanity check”

## Approach

- 1 Decomposition of reward** over real episodes, to identify interesting cases
- 2 Examine counterfactual trajectories** under new policy
- 3 Validate and/or criticize** conclusions, using full patient information (e.g., chart review)

## Example

Observed Outcome	Died	1%	0%	10%
	Disch. No Chg.	0%	1%	43%
	Disch.	0%	0%	44%
		Died	No Chg.	Disch.
		Counterfactual Outcome		

See paper / poster for synthetic case study motivated by sepsis management

# Simulating counterfactual trajectories

## What we need

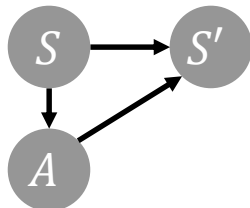
---

- 1 Observed trajectories
- 2 Policy to evaluate  
 $\pi(A | S)$
- 3 Model of discrete dynamics,  
e.g., Markov Decision Process

$S$ : Current State

$A$ : Action

$S'$ : Next State





# Simulating counterfactual trajectories

## What we need

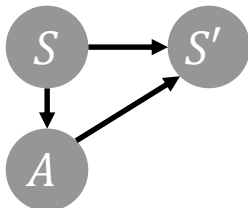
---

- 1 Observed trajectories
- 2 Policy to evaluate  $\pi(A | S)$
- 3 Model of discrete dynamics, e.g., Markov Decision Process

$S$ : Current State

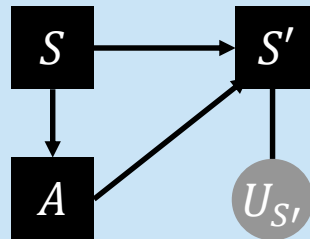
$A$ : Action

$S'$ : Next State



+

## Structural Causal Model (SCM)



$$S' = f(S, A, U_{S'})$$

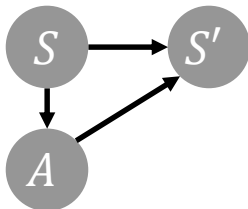
$$U_{S'} \sim P(U_{S'})$$

# Simulating counterfactual trajectories

## What we need

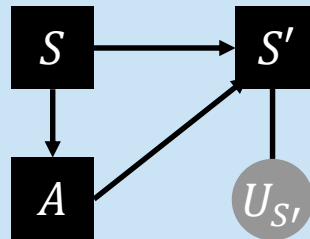
- 1 Observed trajectories
- 2 Policy to evaluate  $\pi(A | S)$
- 3 Model of discrete dynamics, e.g., Markov Decision Process

$S$ : Current State  
 $A$ : Action  
 $S'$ : Next State



+

## Structural Causal Model (SCM)



$$S' = f(S, A, U_{S'})$$
$$U_{S'} \sim P(U_{S'})$$

**Problem:** Choice of SCM is not identifiable from data!

# So, what should we use for the structural causal model (SCM)?

## Key challenge: Non-identifiability

There are multiple SCMs consistent with  $P(S' | S, A)$  but with different *counterfactual* distributions

For **binary variables**, assuming the property of **monotonicity** (Pearl, 2000) is sufficient to identify the counterfactual distribution

**But most real-world MDPs have non-binary states!**

# So, what should we use for the structural causal model (SCM)?

## Key challenge: Non-identifiability

There are multiple SCMs consistent with  $P(S' | S, A)$  but with different *counterfactual* distributions

For **binary variables**, assuming the property of **monotonicity** (Pearl, 2000) is sufficient to identify the counterfactual distribution

**But most real-world MDPs have non-binary states!**

**Theorem 1 (informal):** (Newly defined) property of **counterfactual stability** generalizes monotonicity to categorical variables

# So, what should we use for the structural causal model (SCM)?

## Key challenge: Non-identifiability

There are multiple SCMs consistent with  $P(S' | S, A)$  but with different *counterfactual* distributions

For **binary variables**, assuming the property of **monotonicity** (Pearl, 2000) is sufficient to identify the counterfactual distribution

**But most real-world MDPs have non-binary states!**

**Theorem 1 (informal):** (Newly defined) property of **counterfactual stability** generalizes monotonicity to categorical variables

## **Gumbel-Max SCM**

Use the *Gumbel-Max trick* to sample from a categorical distribution with  $k$  categories:

$$g_j \sim \text{Gumbel}$$
$$S' = \operatorname{argmax}_j \{ \log P(S' = j | S, A) + g_j \}$$

**Theorem 2: Gumbel-Max SCM** satisfies the counterfactual stability condition

# Thank you!

Come to our poster for more details: **Pacific Ballroom #72**