# *Automating Anaesthetic Delivery with Deep Reinforcement Learning*

December 4, 2020

**Gabriel Schamberg**, Marcus Badgeley, Emery Brown

# Outline

**Background**
- Closed loop anaesthetic delivery & reinforcement learning

**Methods**
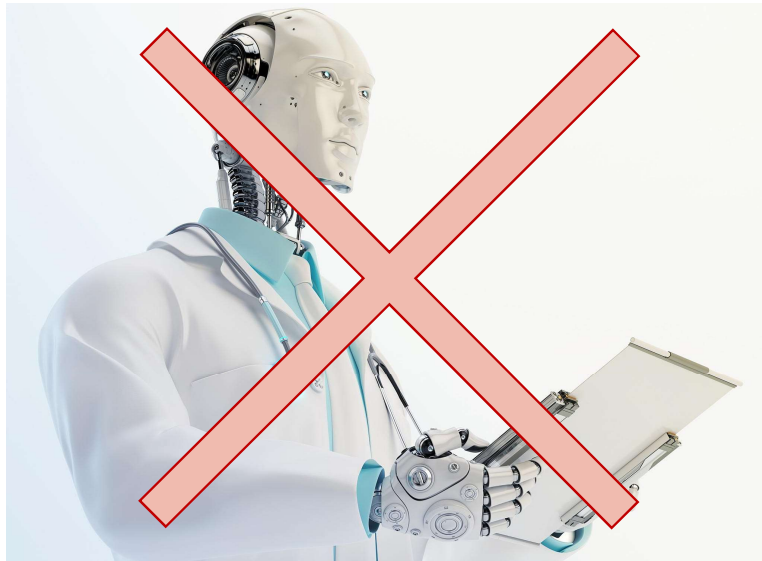- Deep RL paradigm for controlling level of unconsciousness

**Results**
- Simulation study

# Autopilot for Anaesthesia

# Closed Loop Anesthetic Delivery (CLAD)



Anaesthesiologist sets a *target level of unconsciousness*

Controller determines an appropriate dose

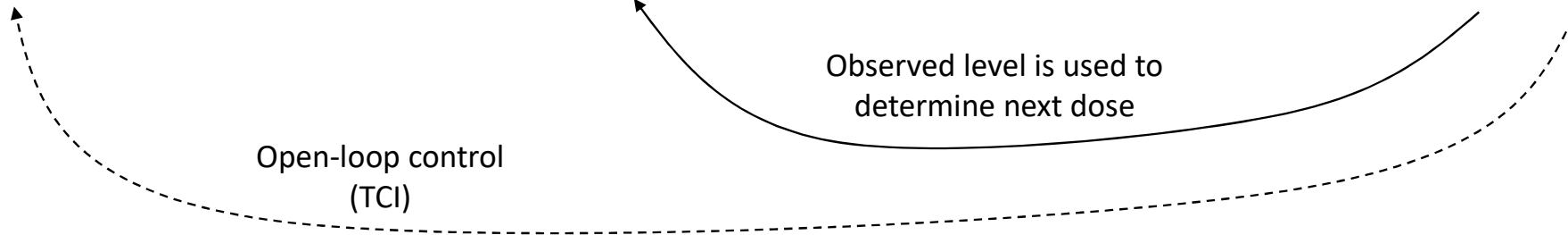An infusion pump administers the dose to patient

EEG monitor provides *observed* level of unconsciousness

Observed level is used to determine next dose

Open-loop control (TCI)

# CLAD – What's Been Done?

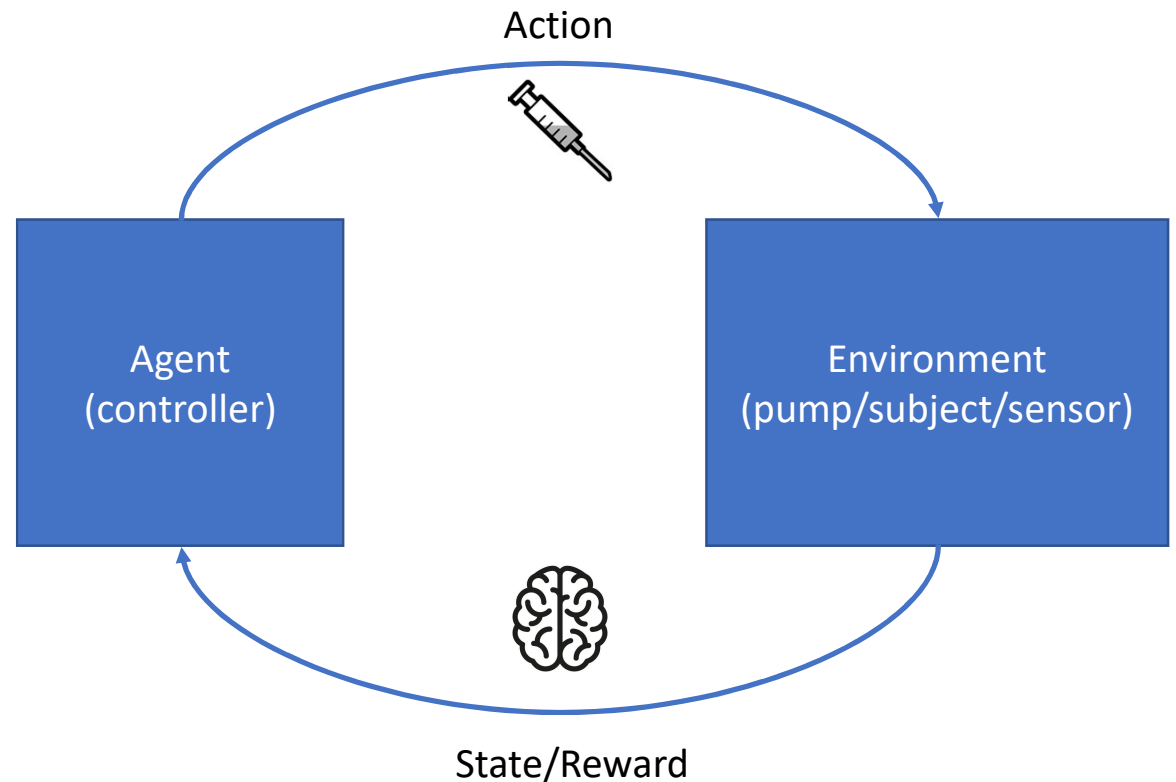| Paper | Subject | Controller | Level of Unconsciousness |
|---|---|---|---|
| Absalom (2002) | Human | PID | BIS |
| Dumont (2009) | Simulation | Robust PID/CRONE | WAV$_{CNS}$ |
| Shanechi (2013) | Rodent | LQR | BSP |
| Moore (2014) | Human | Tabular RL | BIS |
| *Many More…* | | | |

**THE PICOWER INSTITUTE**
FOR LEARNING AND MEMORY

# Reinforcement Learning

The agent observes a **state**

Based on the state, the agent uses its **policy** to determine an **action**

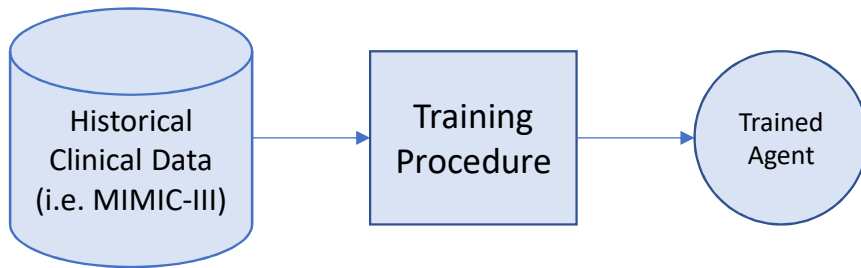Over time, the agent receives **rewards**, which are used to update its policy

Action

Agent
(controller)

Environment
(pump/subject/sensor)

State/Reward

# Reinforcement Learning – Two Approaches

**Off-Policy**:
*The agent "watches and learns"*

**On-Policy**:
*The agent "learns by doing"*

*Focus of this talk*



Off-Policy flow: Historical Clinical Data (i.e. MIMIC-III) → Training Procedure → Trained Agent
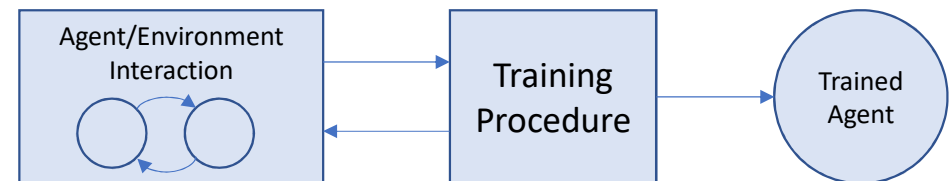
- Agent sees real patient responses
- Agent learns to reflect anaesthesiologist's actions

- Data may not contain enough bad or ideal behavior
- For example, infrequent titration is suboptimal

On-Policy flow: Agent/Environment Interaction ⇄ Training Procedure → Trained Agent

- Agent can fully explore action space
- The actions will respond to varying reward functions

- Agent cannot act on humans → requires simulations
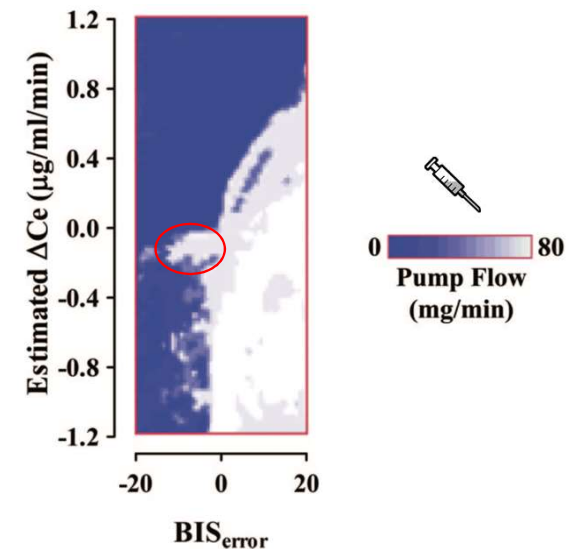
THE PICOWER INSTITUTE
FOR LEARNING AND MEMORY

# Tabular RL Policies

A **policy** maps a **state** to an **action**

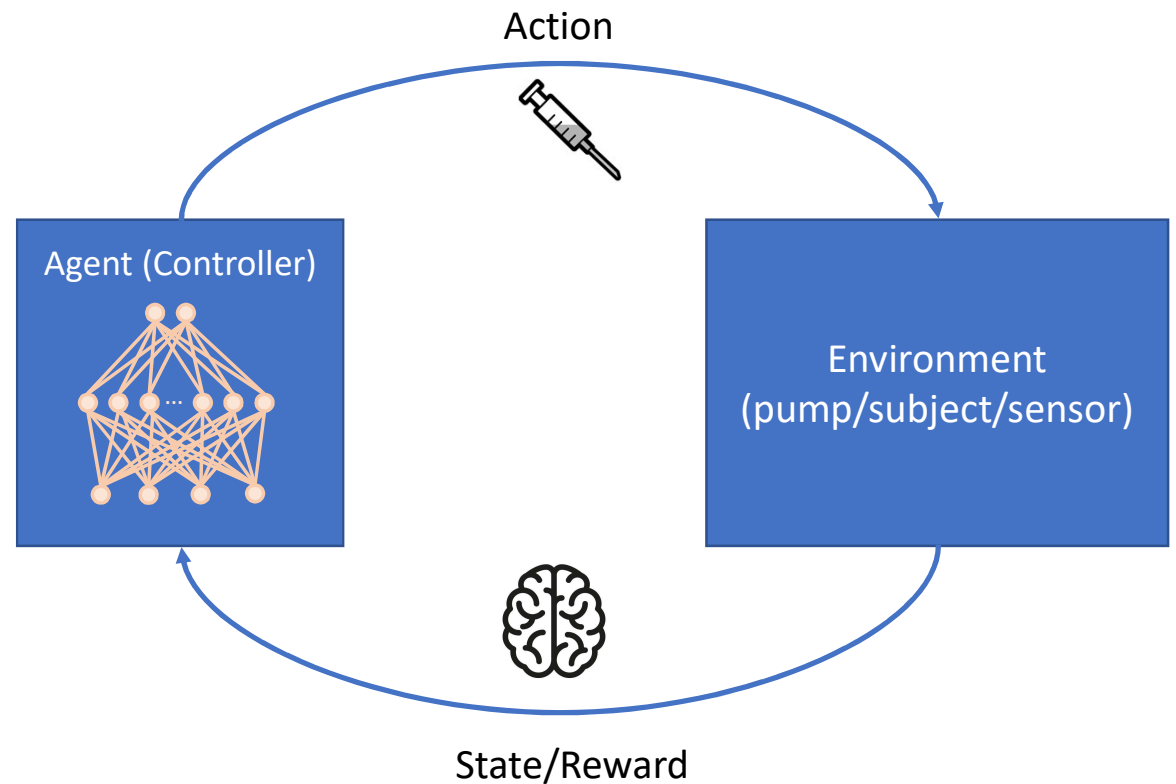Tabular RL learns which discrete action should be taken for every *discrete state*

- Ignores the continuous relationship between state and action
- Results in patchy discontinuities
- Dimensionality of the table scales exponentially with the state space
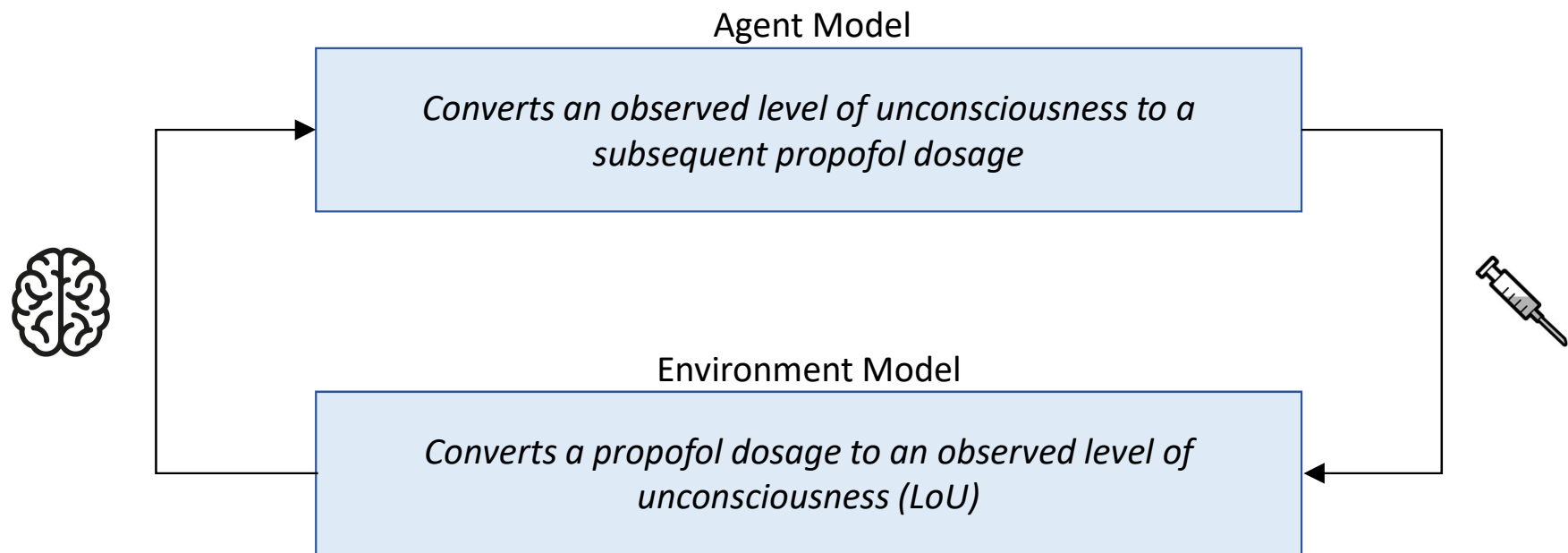


*From Moore et al. (2011)*

# Deep Reinforcement Learning
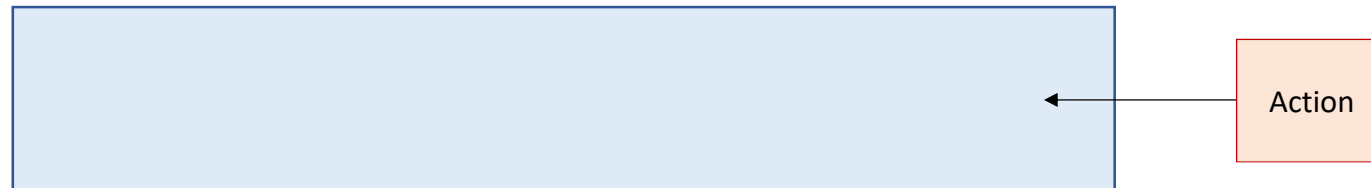
The policy is a **deep neural network**

Action

Agent (Controller)

Environment (pump/subject/sensor)

State/Reward

# Deep RL for CLAD

Agent Model

Converts an observed level of unconsciousness to a subsequent propofol dosage

Environment Model

Converts a propofol dosage to an observed level of unconsciousness (LoU)
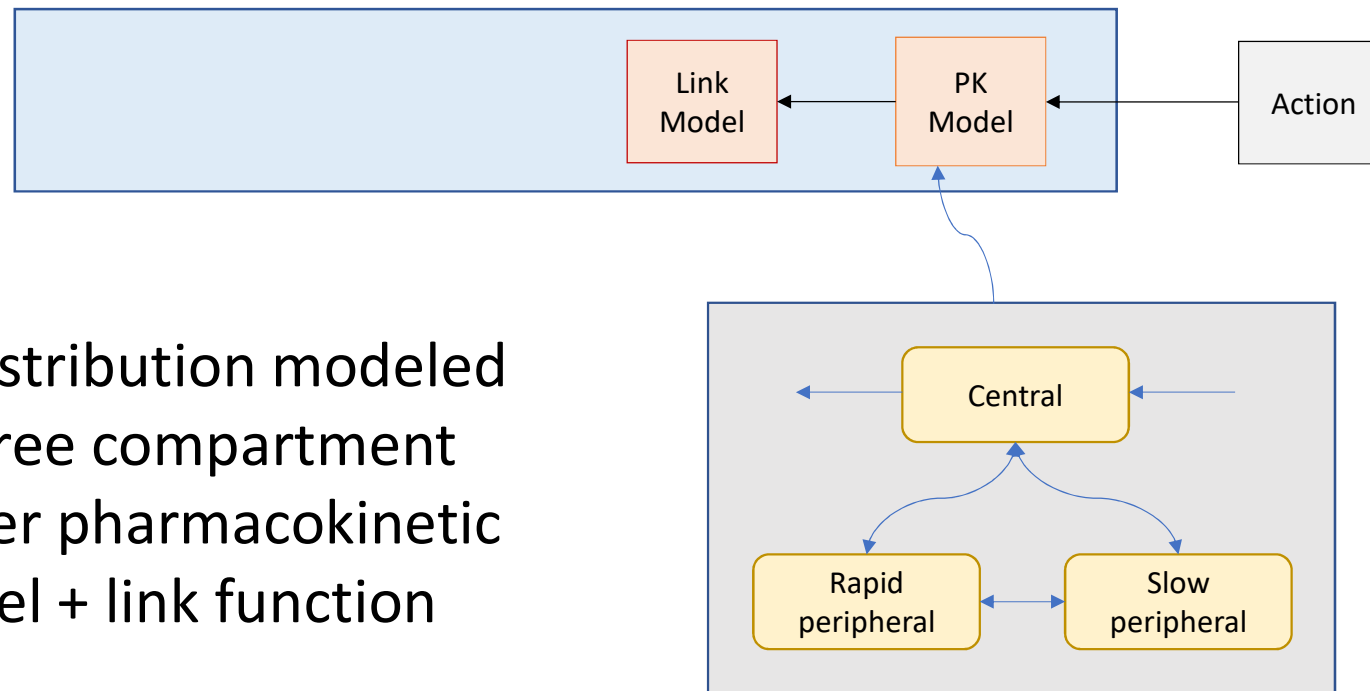
# The Environment Model
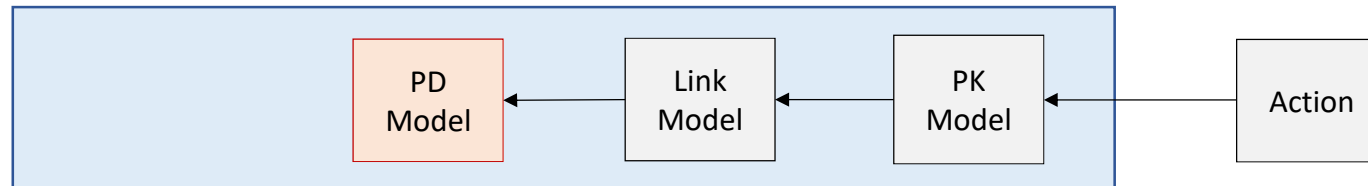


5 seconds of
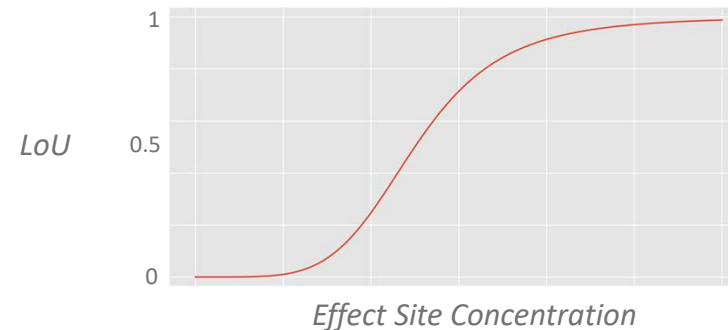**1.67 mg/s** or **0 mg/s**

# The Environment Model



Drug distribution modeled
by three compartment
Schnider pharmacokinetic
model + link function
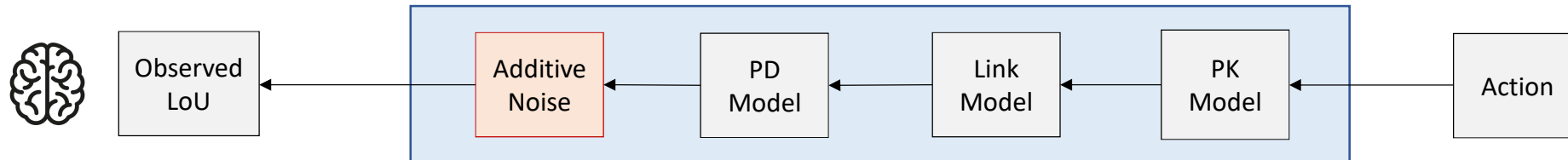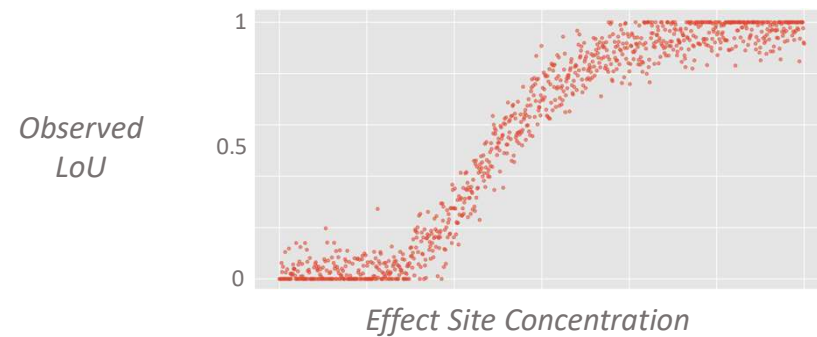
# The Environment Model



Drug effect modeled by nonlinear pharmacodynamic model determines level of unconsciousness (LoU)
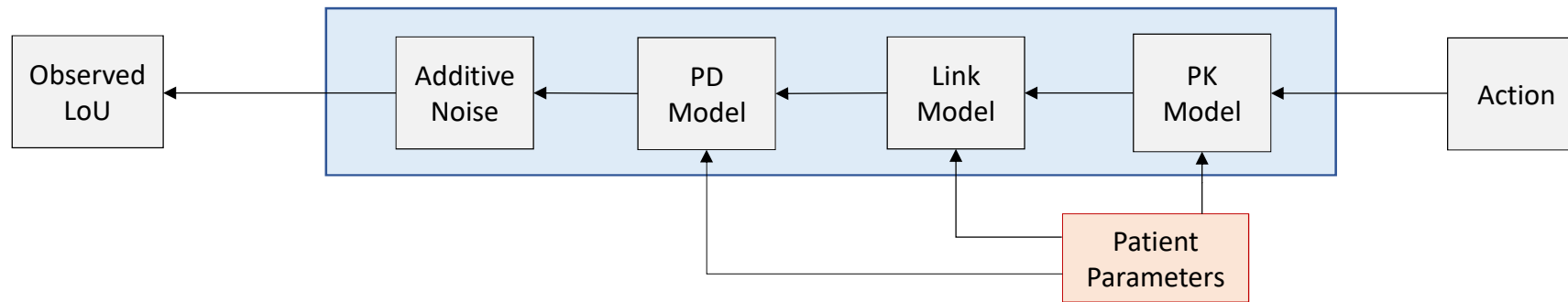
# The Environment Model



Observed LoU includes additive noise
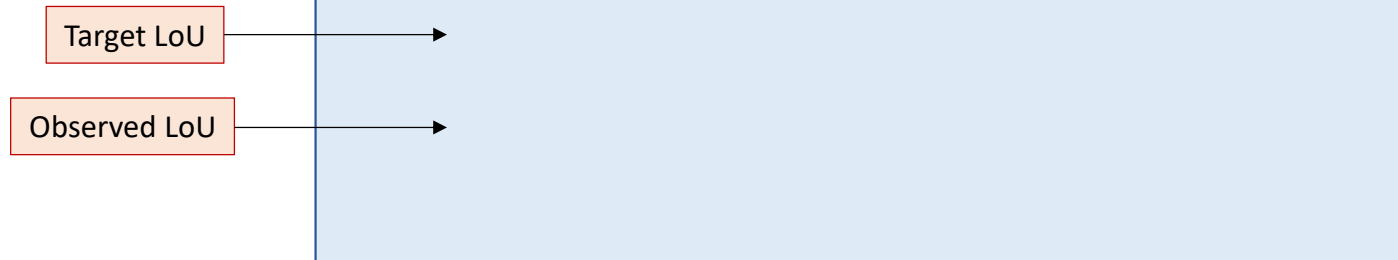
# The Environment Model



patient variability model

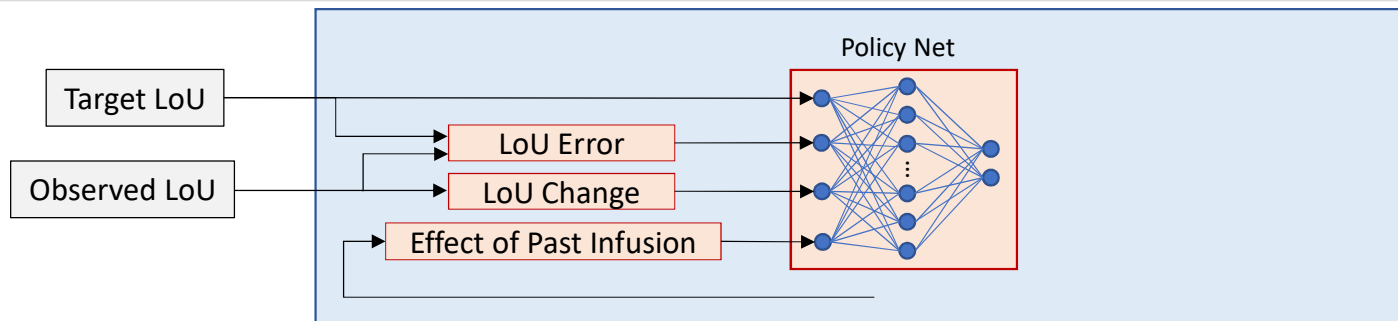| Sub-model | Parameter | Units | Generic | Minimum | Maximum |
|-----------|-----------|-------|---------|---------|---------|
| PK | Height | cm | 170 | 160 | 190 |
| PK | Weight | kg | 70 | 50 | 100 |
| PK | Age | yr | 30 | 18 | 90 |
| Link | $k_{e0}$ | $\text{min}^{-1}$ | 0.17 | 0.128 | 0.213 |
| PD | $\gamma$ | - | 5 | 5 | 9 |
| PD | $C$ | - | 2.5 | 2 | 6 |

# The Agent Model
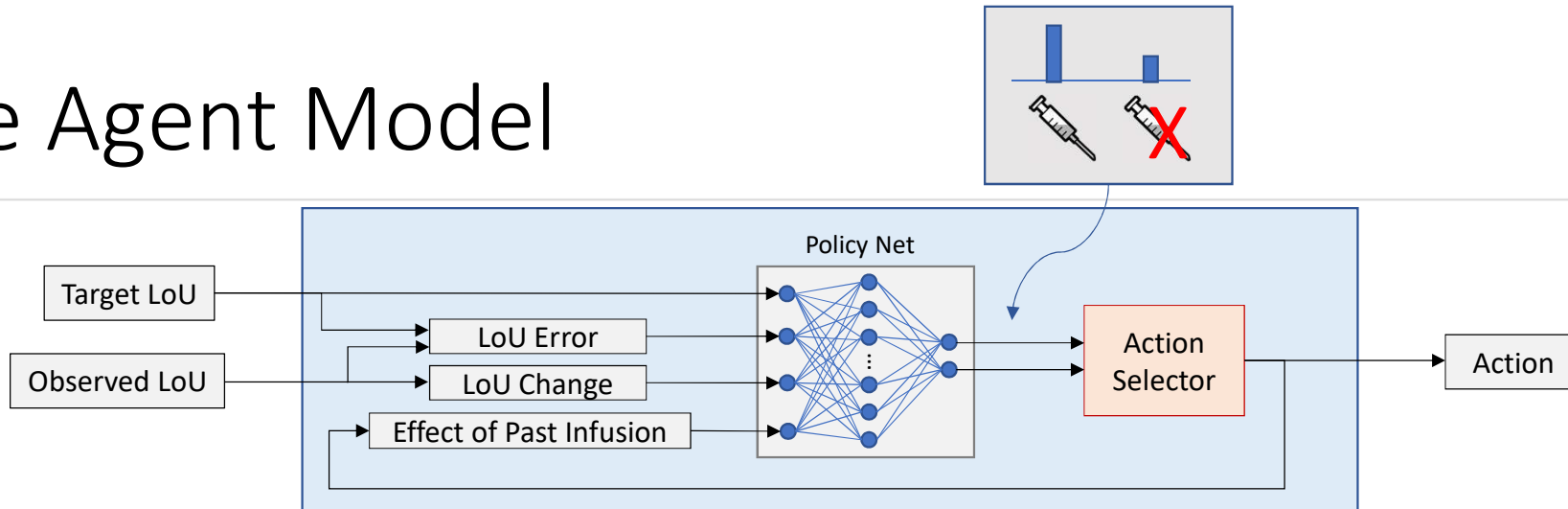
Target LoU

Observed LoU

Agent receives observed
and target LoUs

# The Agent Model



The policy assigns probabilities to actions
based on a 4 dimensional observation

# The Agent Model



Agent uses one of three
strategies to select an action

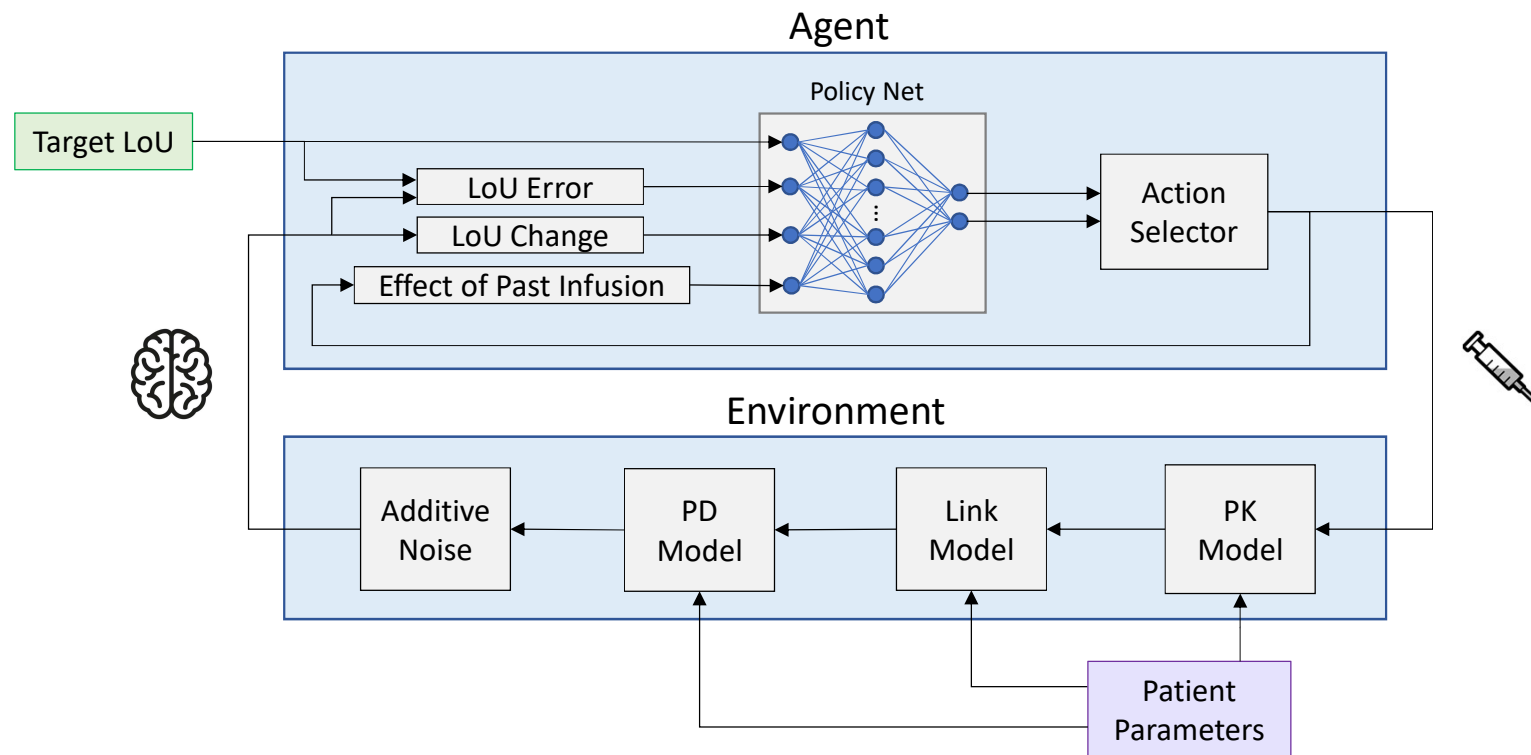**_Stochastic_**
*Picks action randomly
according policy output*

**_Deterministic_**
*Pick action with
highest probability*

**_Continuous_**
*Multiply action
probability by max dose*

THE PICOWER
INSTITUTE
FOR LEARNING AND MEMORY

# Complete Simulation Model
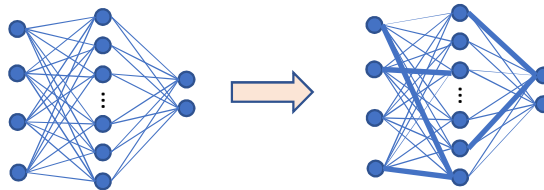
# Training the Agent

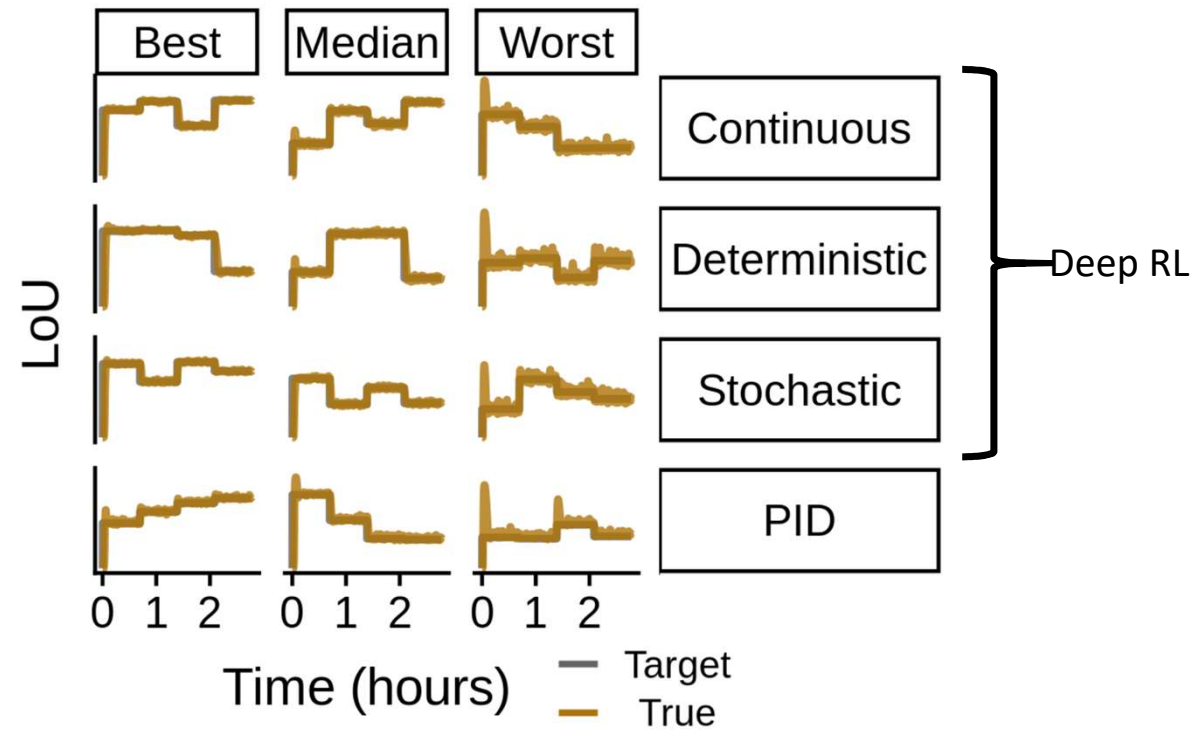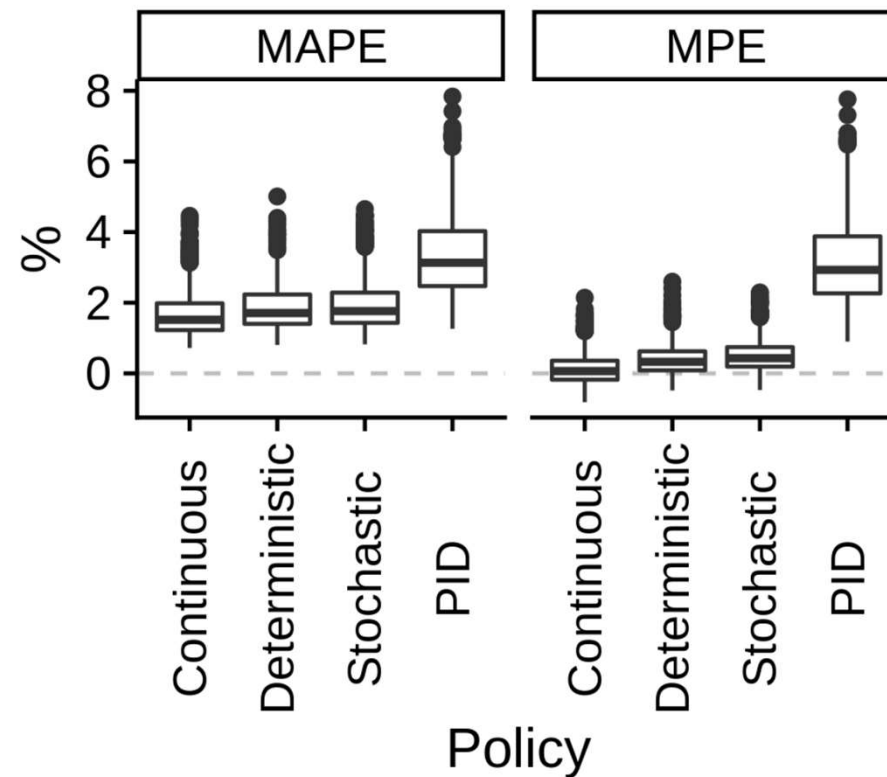Agent acts with randomness in a series of "episodes"

Identify which episodes were best
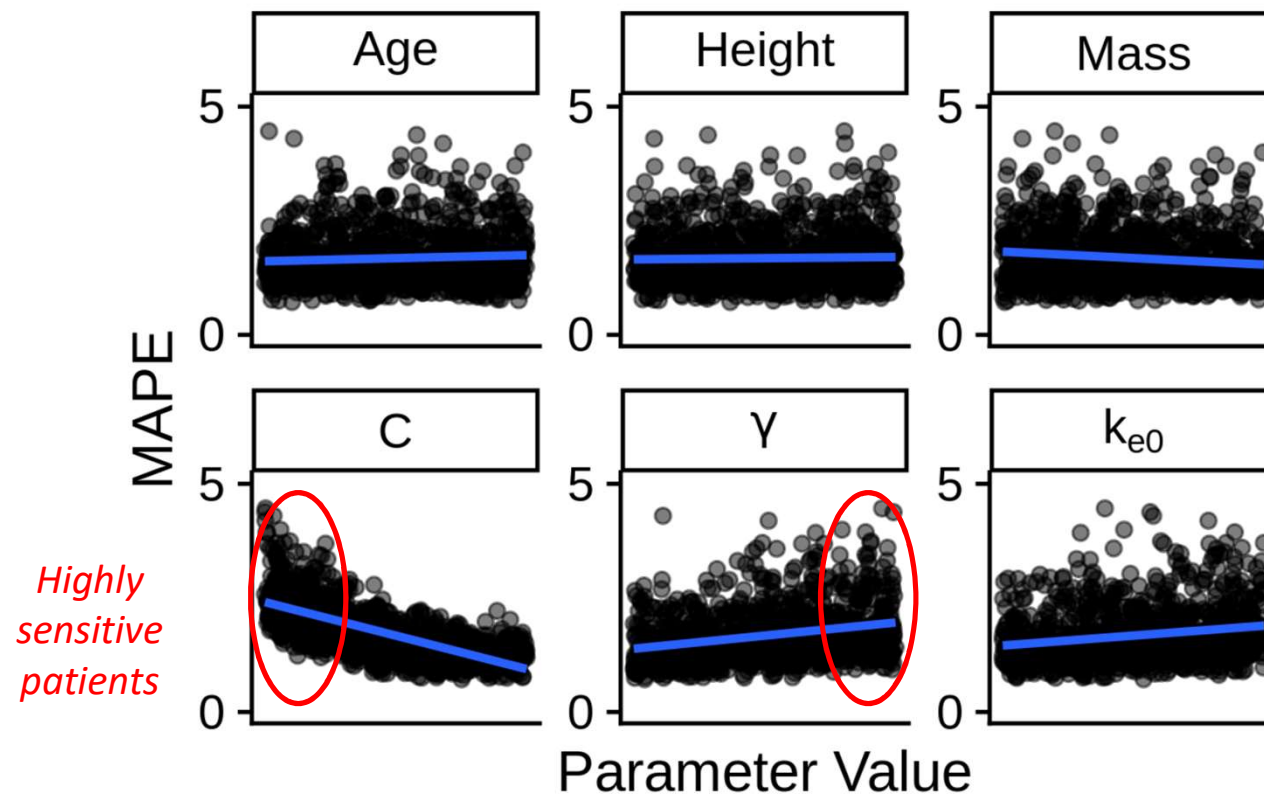
Nudge policy to encourage repeating best actions

# Deep RL Works

# Deep RL Outperforms PID

# Deep RL is Robust to Patient Variability



*Highly sensitive patients*
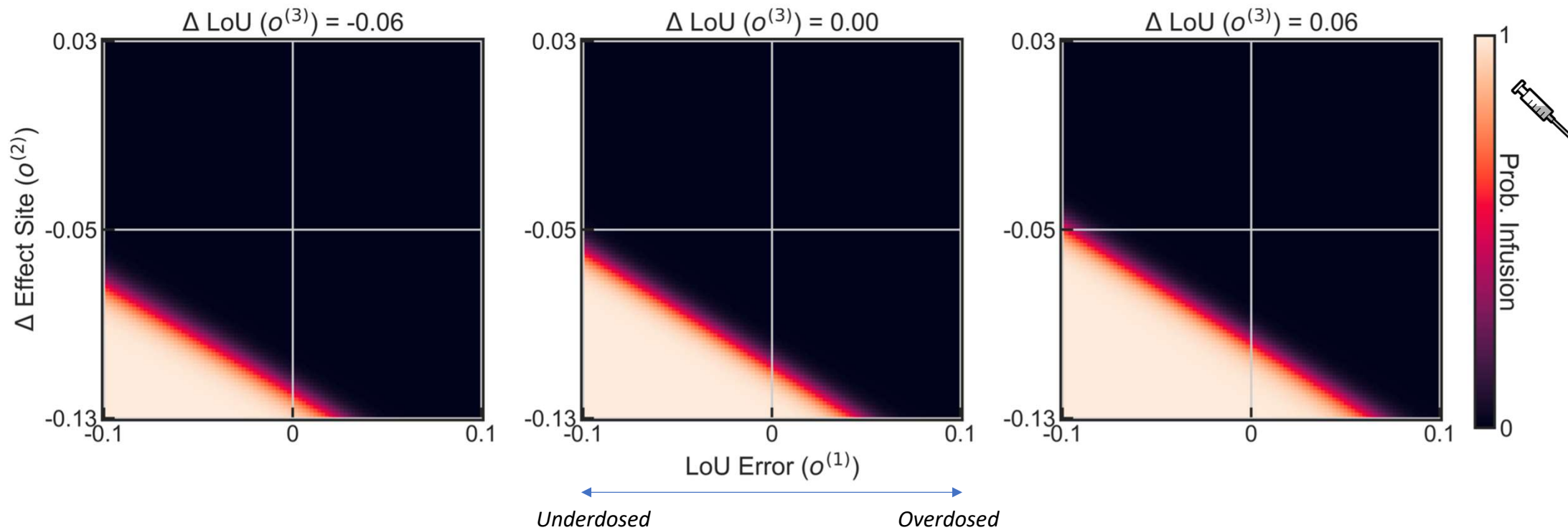
# Deep RL is Not a Black Box

# Key Conclusions

- Deep RL outperforms PID controller in simulations
- Deep RL resolves issues with tabular RL
  - Enables consideration of other patient variables
  - Reflects continuous state/action relationships
- Resulting policy is not a black box

# Future Work

- Further simulation studies
  - Intra-patient variability in environment
  - Continued feature engineering
  - Off-policy training on retrospective clinical data
- Clinical studies
  - Human-in-the-loop recommender system
  - Animal studies
  - Volunteers?

THE PICOWER
INSTITUTE
FOR LEARNING AND MEMORY

# Acknowledgements

- Marcus Badgeley
- Emery Brown
- Benyamin Meschede-Krasa
- John Abel

BACKUP

# A Side by Side Comparison of Control Algorithms

| Classical Model-Free (PID) | Classical Model-Based (LQR) | Deep RL |
|---|---|---|
| Parameters are *tuned* using nominal transfer function | Parameters are *derived* using nominal patient model | Parameters are *learned* using patient model simulations |
| Does not optimize | Optimizes a *quadratic cost* | Optimizes *reward* |
| *Linear* function of *error* | *Linear* function of *state* | *Non-Linear* function of *state* |
| Well-established methods for testing stability, robustness, responsiveness | | No formal analysis |
| End product is a *deterministic* controller | | |