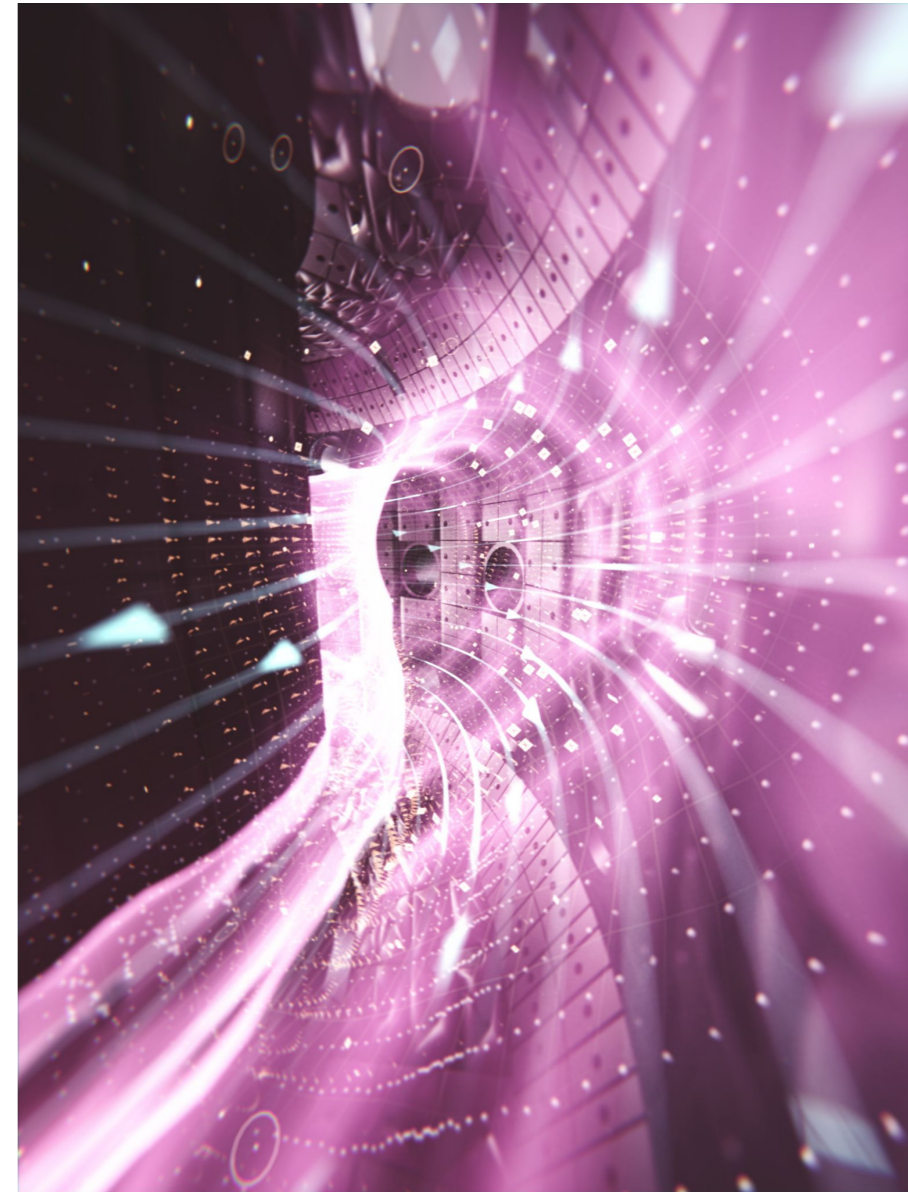


# Magnetic Control of Tokamak Plasmas through Deep Reinforcement Learning

**Brendan Tracey - Google DeepMind, London**

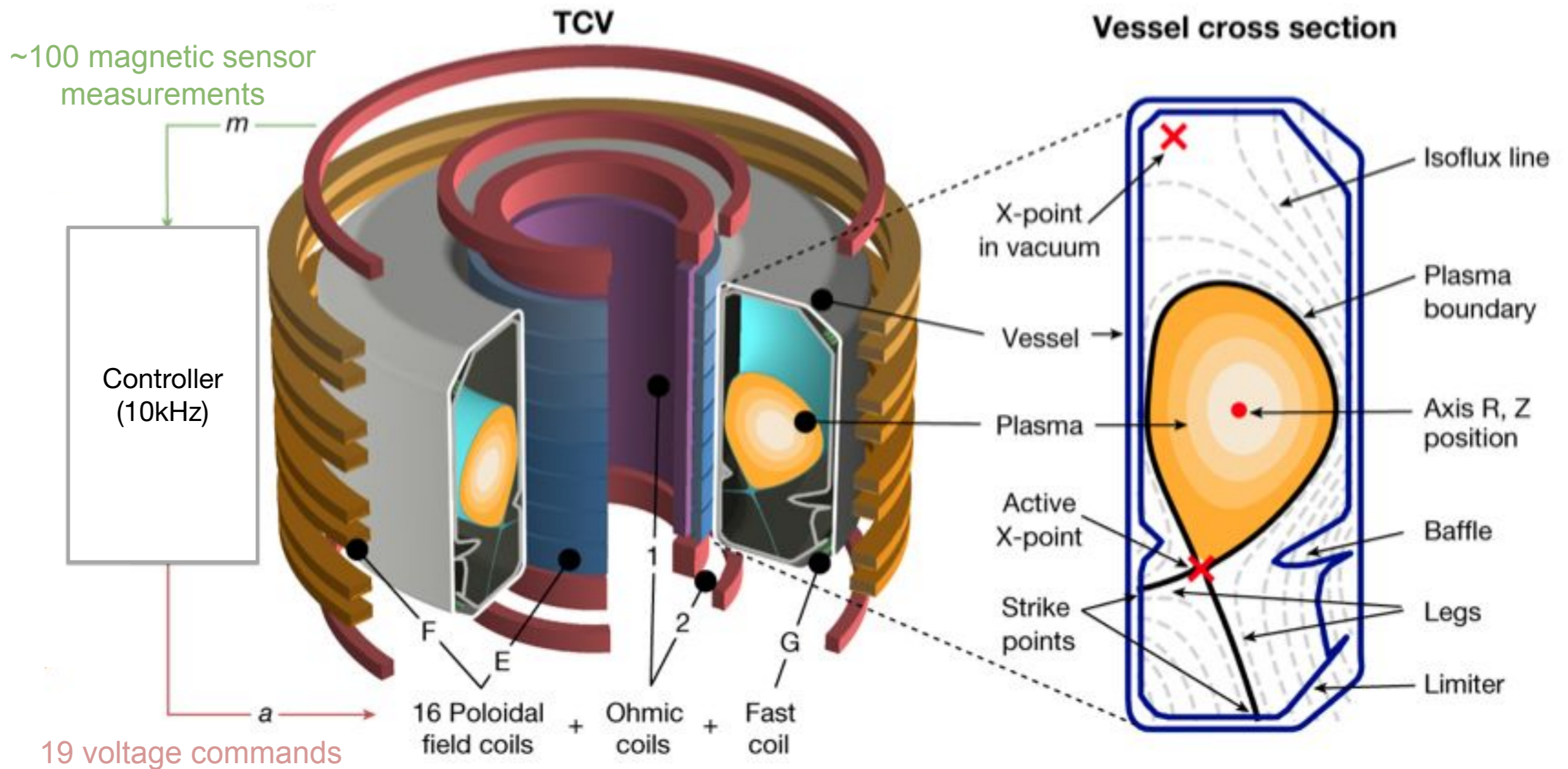
**on behalf of the DeepMind-EPFL team**

Jonas Degraeve, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de las Casas, Craig Donner, Leslie Fritz, Cristian Galperti, Andrea Huber, James Keeling, Maria Tsimpoukelli, Jackie Kay, Antoine Merle, Jean-Marc Moret, Seb Noury, Federico Pesamosca, David Pfau, Olivier Sauter, Cristian Sommariva, Stefano Coda, Basil Duval, Ambrogio Fasoli, Ian Davies, Andrea Michi, Yuri Chervonyi, Pushmeet Kohli, Koray Kavukcuoglu, Demis Hassabis & Martin Riedmiller



**Workshop on AI for Accelerating Fusion and Plasma Science  
November 2023**

# Axisymmetric Equilibrium Control

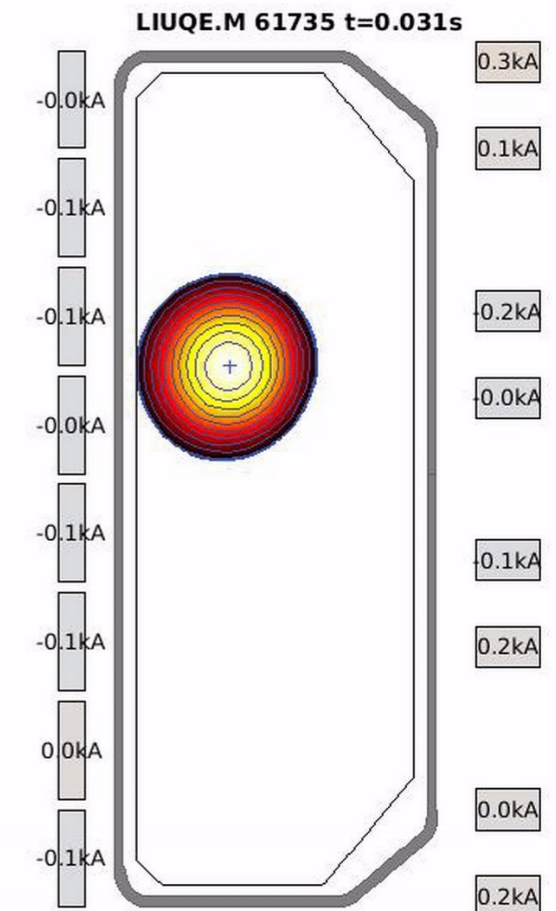
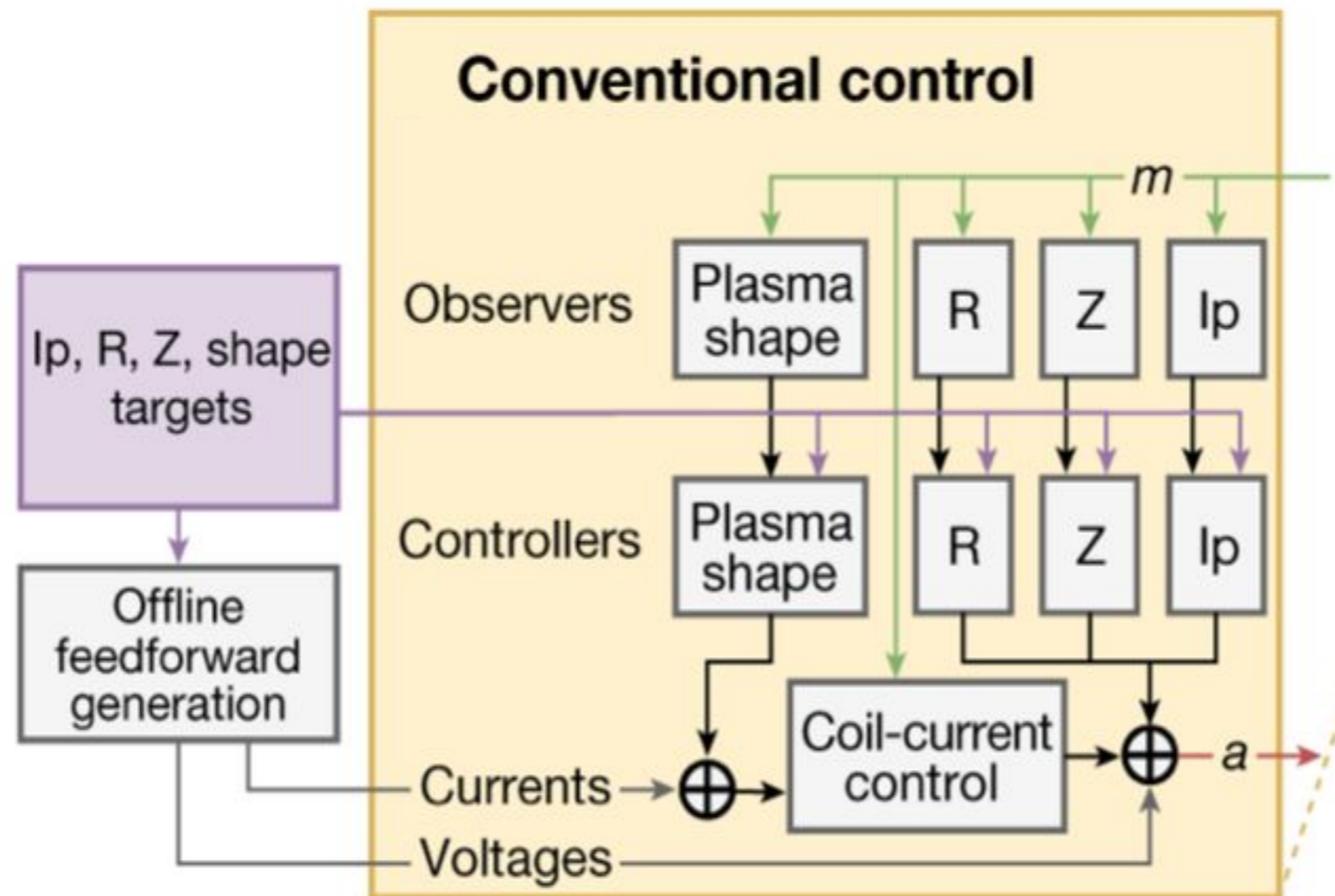


# Axisymmetric Tokamak Plasma Control

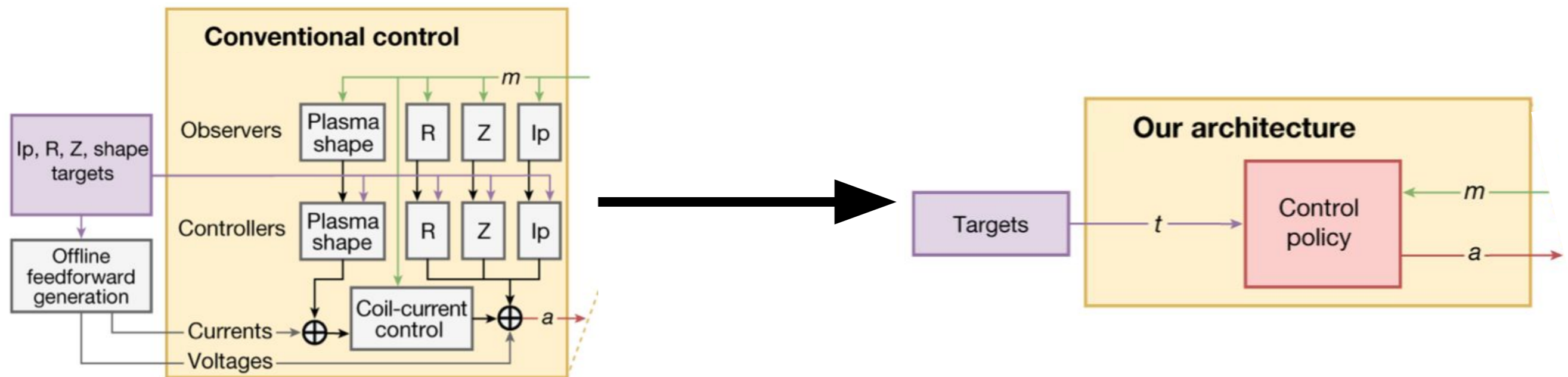
- **Need to control:**
  - Total plasma current  $I_p$
  - Radial position R (by vertical magnetic fields)
  - Vertical position Z (by radial magnetic fields - *unstable for elongated plasmas*)
  - Plasma shape: Last Closed Flux Surface (LCFS)

# Traditional Solutions

1. Choose combinations of coils to use to control each quantity
2. Pre-compute feedforward coil currents and voltages
3. Design feedback controllers
4. Tune individual control parameters for each (hopefully orthogonal) control loop

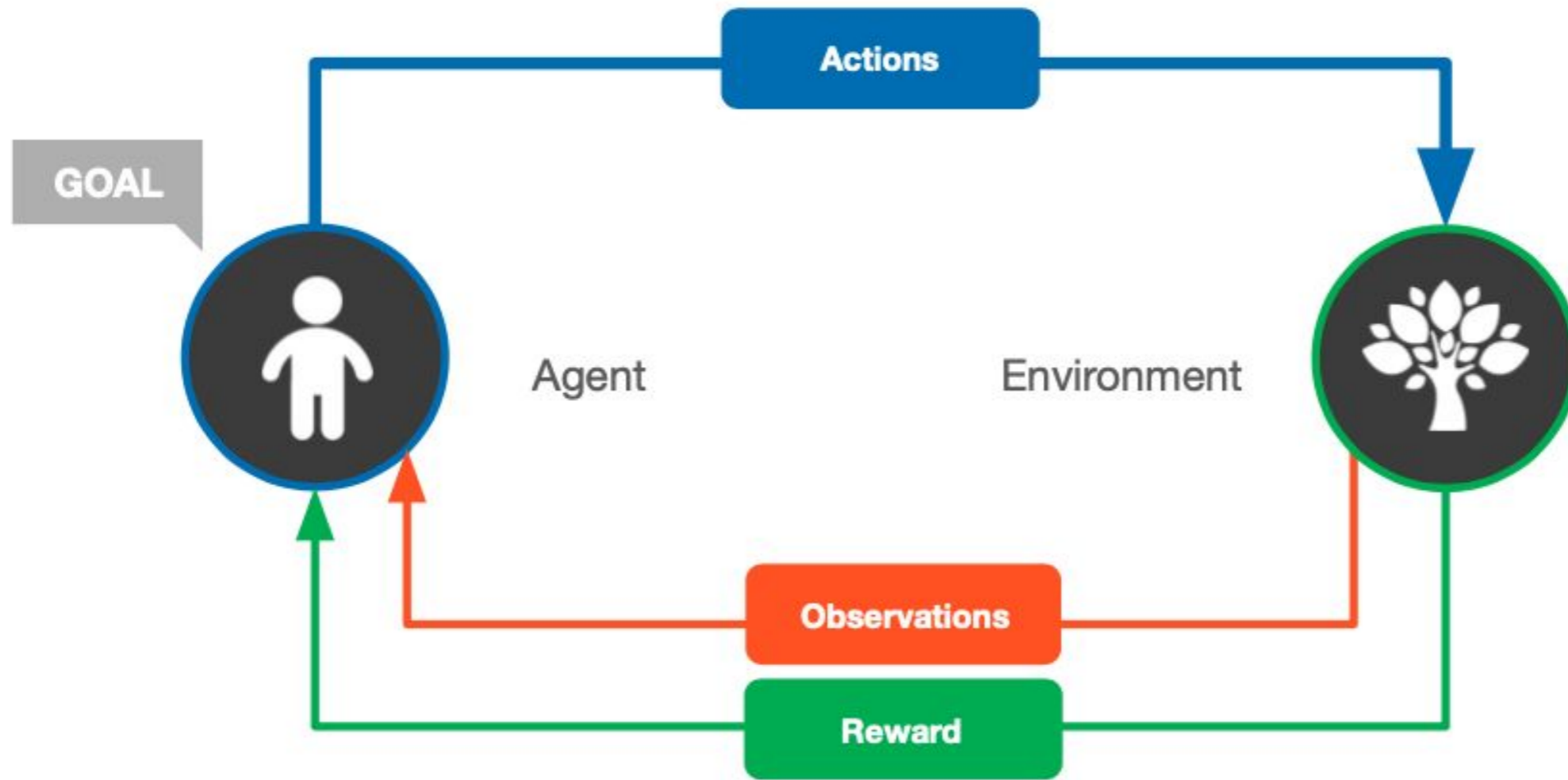


# Reinforcement Learning Solution



- **Single Integrated Controller**
- **No feedforward generation**
- **No separate error estimation**

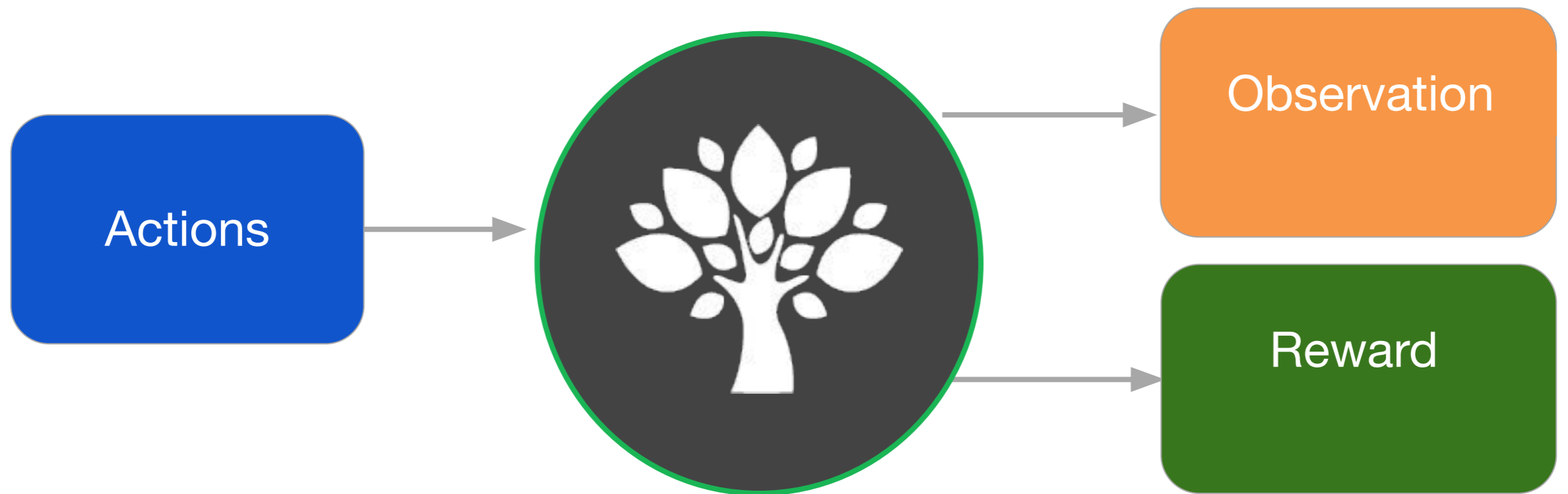
# Reinforcement Learning



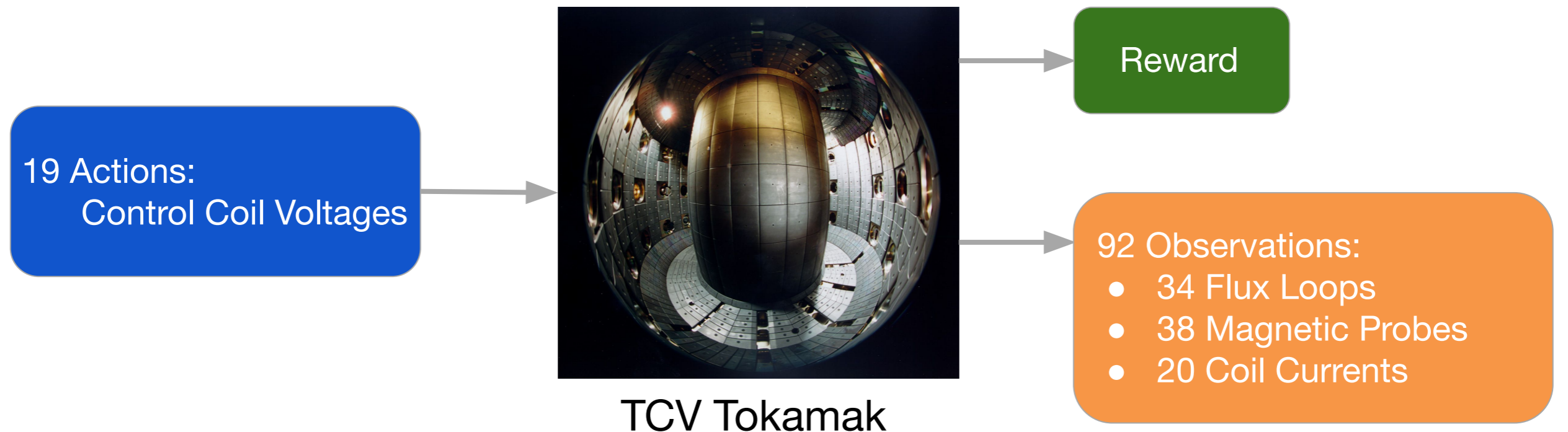
[Figure and RL slide material from hereon: courtesy A. Abdolmaleki]

Learn an action selection function ( $\pi$ ) through trial-and-error to achieve high reward

# What is an environment?



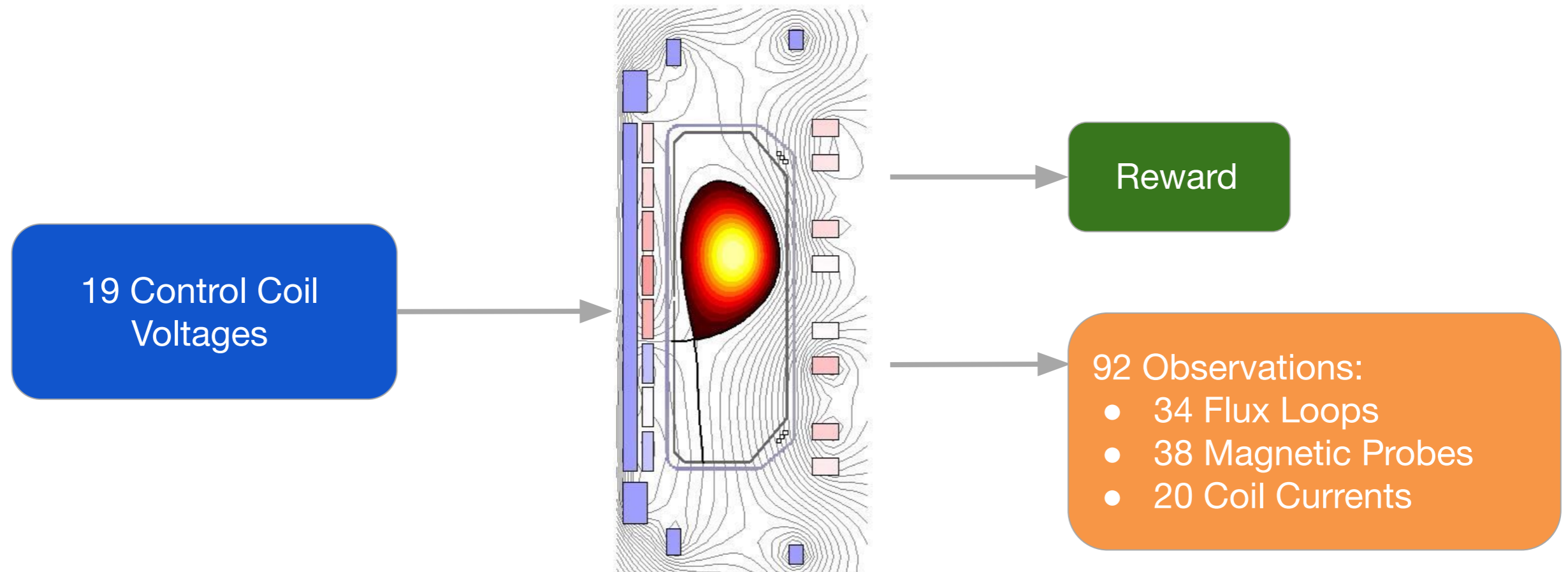
# Environment



Training on Hardware is Impractical



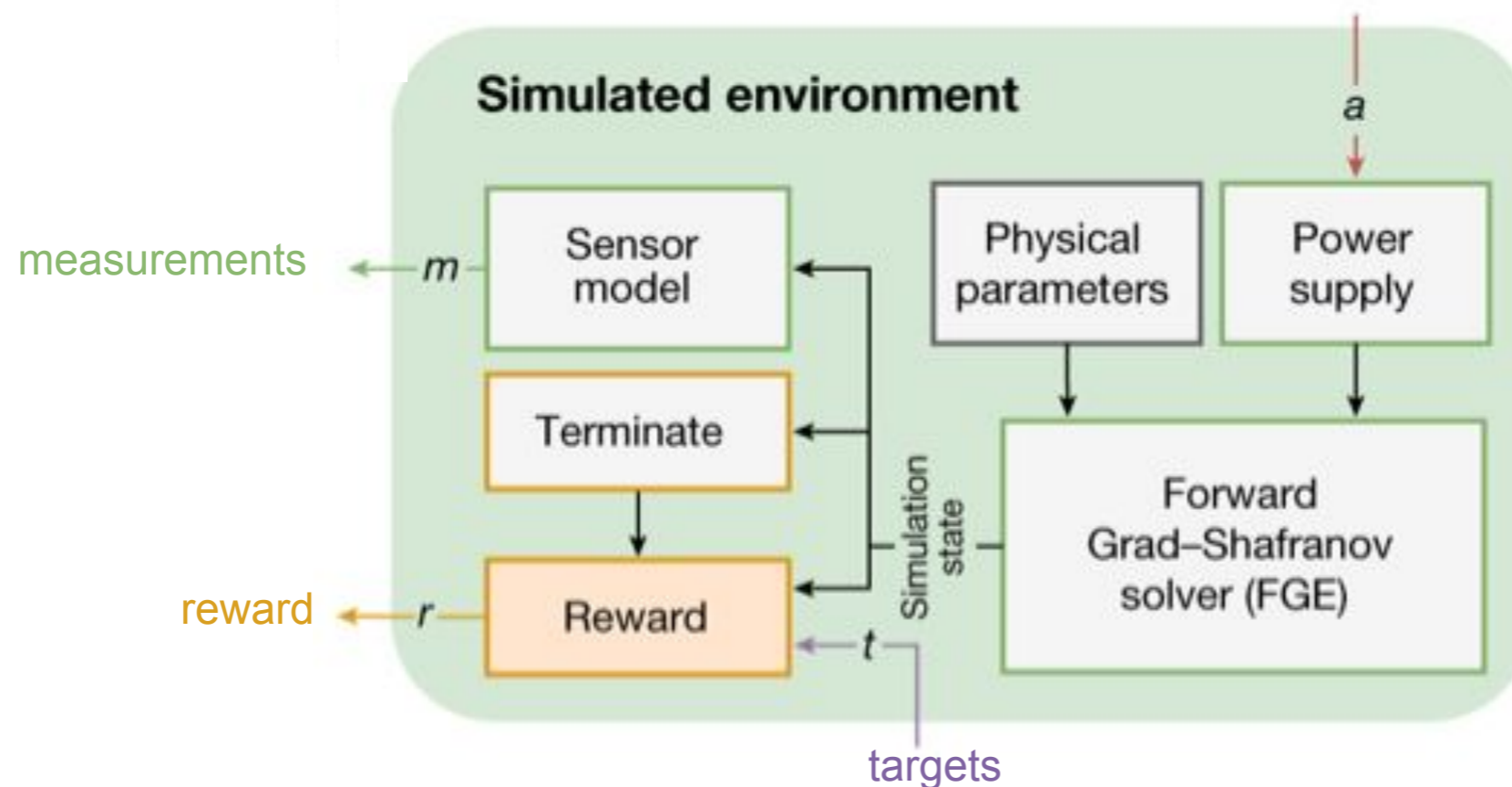
# Training Environment



FGE Simulator [Carpanese EPFL PhD 2020]

- Free boundary Grad-Shafranov solver
- Circuit equations for conductors

# Full Simulation Model

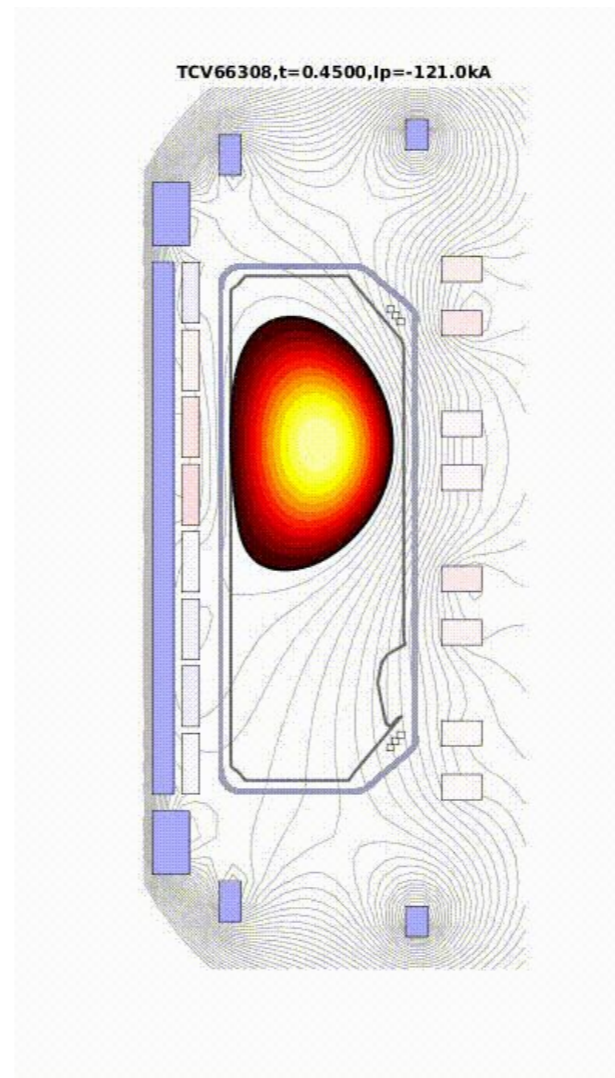


- A lot of physics/engineering know-how goes into the simulator
- Need for domain experts: plasma physicists / tokamak engineers & modelers

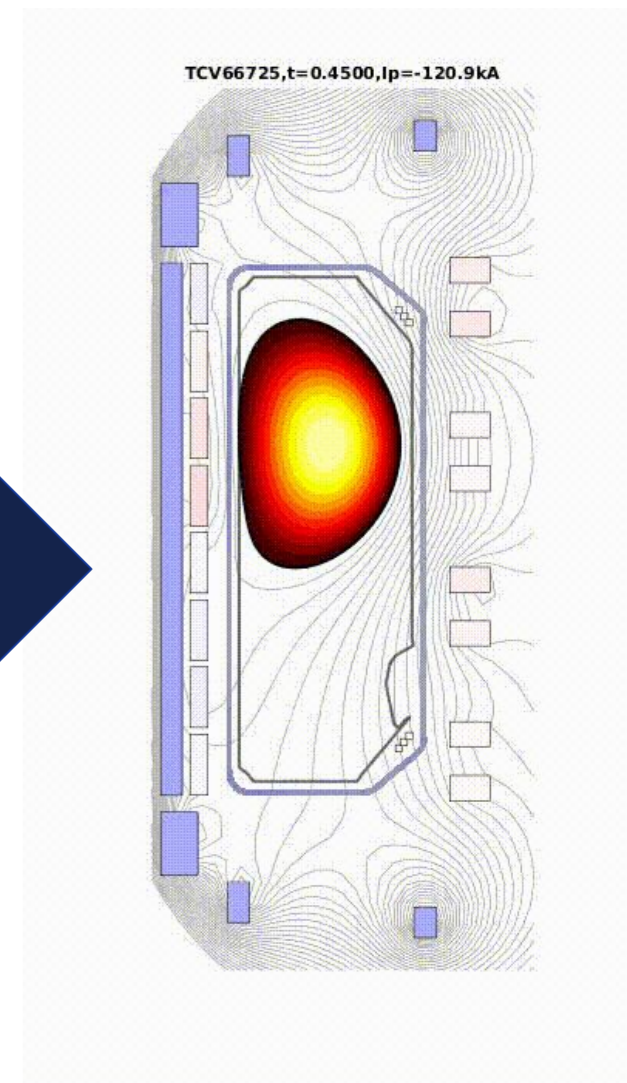
# Reward Design

## Goals

1. Keep the plasma alive
2. Stabilize the plasma location
3. Shape Control



Shot 66308:  
(0.024s real time)

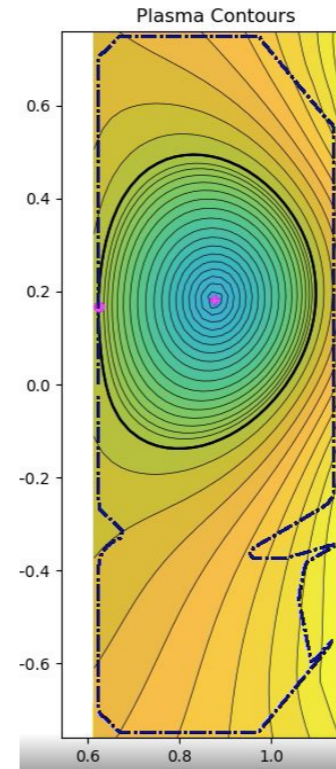


Shot 66725:  
(0.55s real time)

# Reward Design

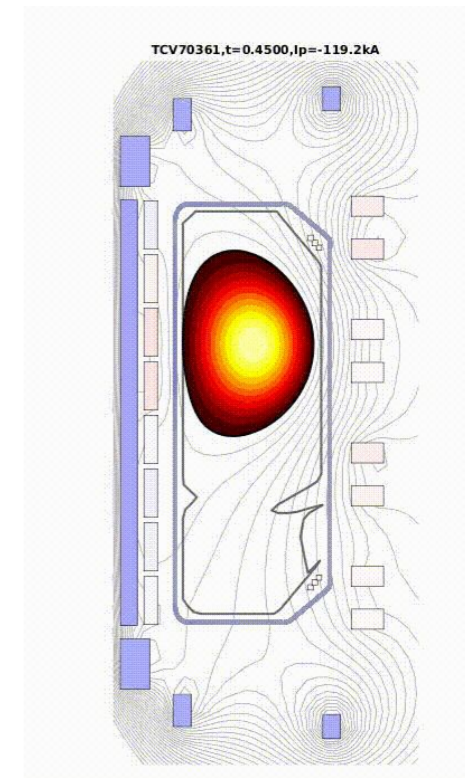
## Goals

1. Keep the plasma alive
2. Stabilize the plasma location
3. Shape Control



Stabilize:

- R Centroid
- Z Centroid
- Plasma Current ( $I_p$ )



Stabilizing the position through control  
Shot 70361: 0.55s real time

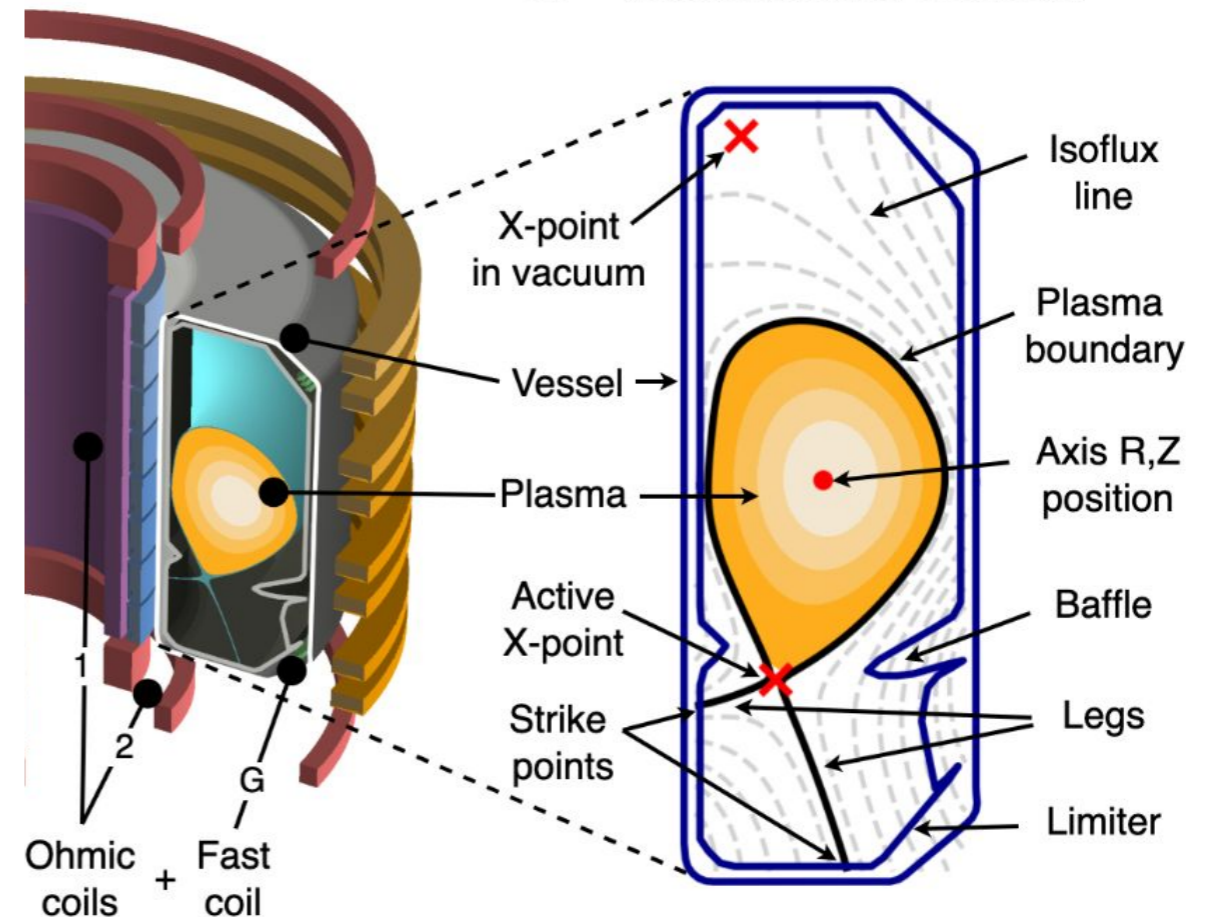
# Reward Design

## Goals

1. Keep the plasma alive
2. Stabilize the plasma location
3. Shape Control

$$\text{SoftPlus}(x) = [2 \cdot \sigma(f_{\text{scale}}(x, \text{good}, \text{bad}, 0, \zeta))]_0^1,$$

$$\text{SmoothMax}(x_{1\dots n}, w_{1\dots n}, \alpha) = \frac{\sum_{i=1}^n w_i x_i e^{\alpha x_i}}{\sum_{i=1}^n w_i e^{\alpha x_i}}.$$

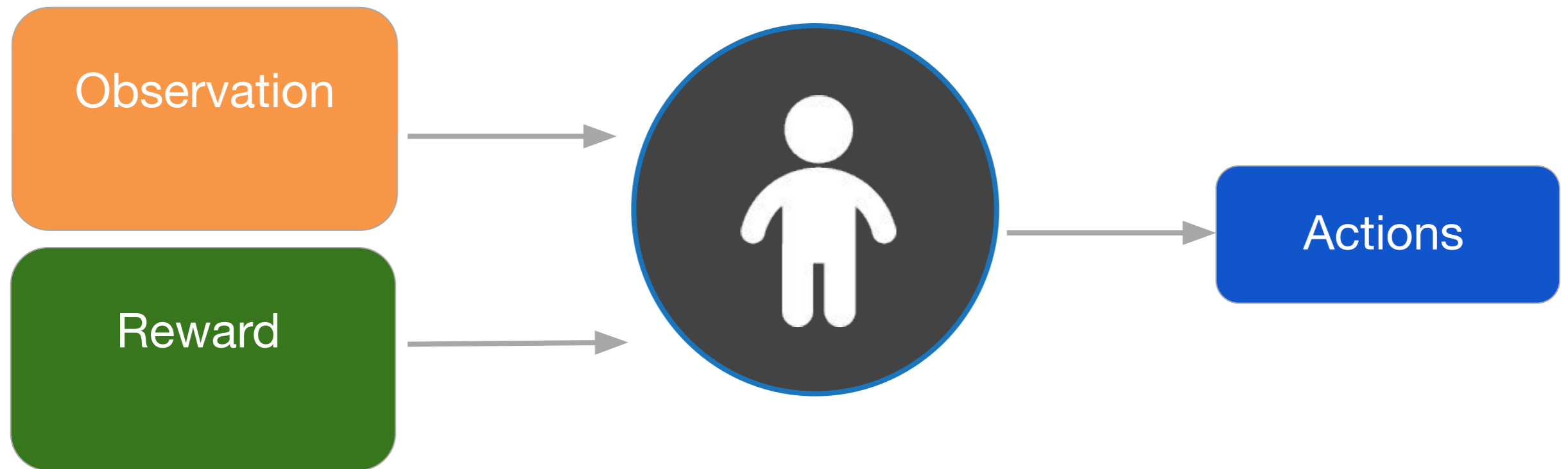


Shape Control – Must match:

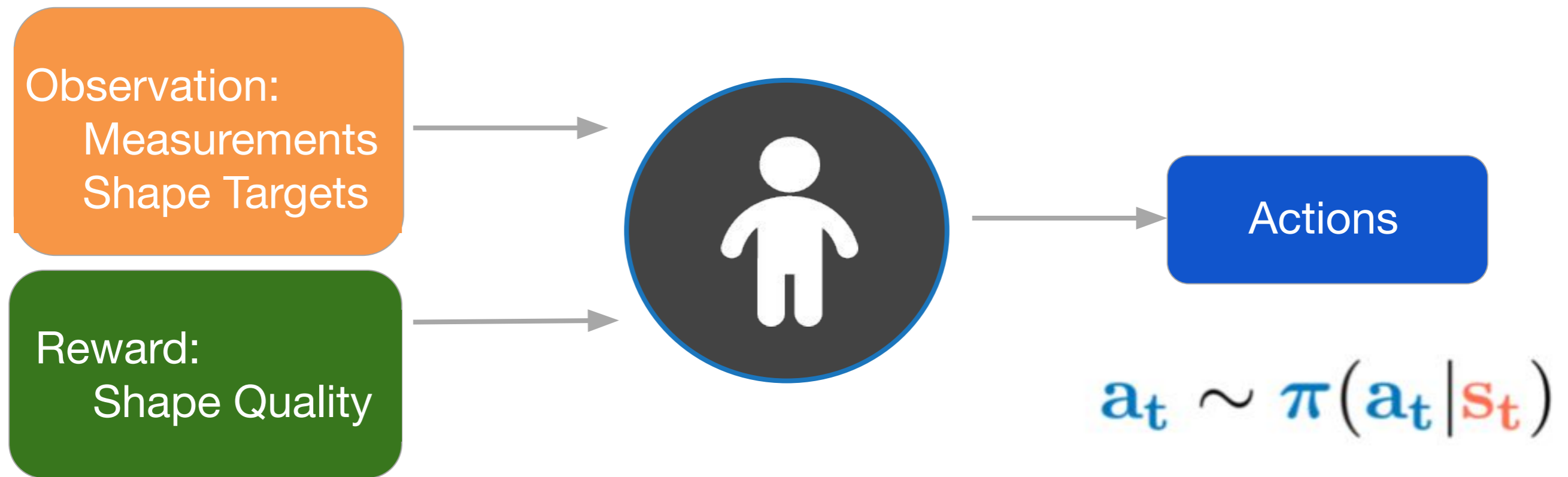
- 1) Target shape outline
- 2) Active X-point location
- 3) Passive X-point locations
- 4) Leg locations
- 5) Plasma current

Whilst maintaining OH coil current stability

# What is an Agent?



# Tokamak Agent



## Aim

Find **optimal policy** – maximize discounted sum of **future rewards**

# Q Function: State-Action Values

**State-action value function** maps an **action** in a given **state** to **expected rewards**.

$$Q^{\pi}(s_t, a_t) = \mathbb{E}_{\pi} [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t, a_t]$$

It is equal to **expected total discounted reward** for an agent starting from state **s** and performing action **a** and following its **policy**.



# Actor-Critic Methods

- Critic estimates the *Q function* from data generated by interacting with environment
- Actor learns a policy  $\pi$  by ascending (via gradient-based methods) the *Q function* learned by the critic

# Actor-Critic Methods

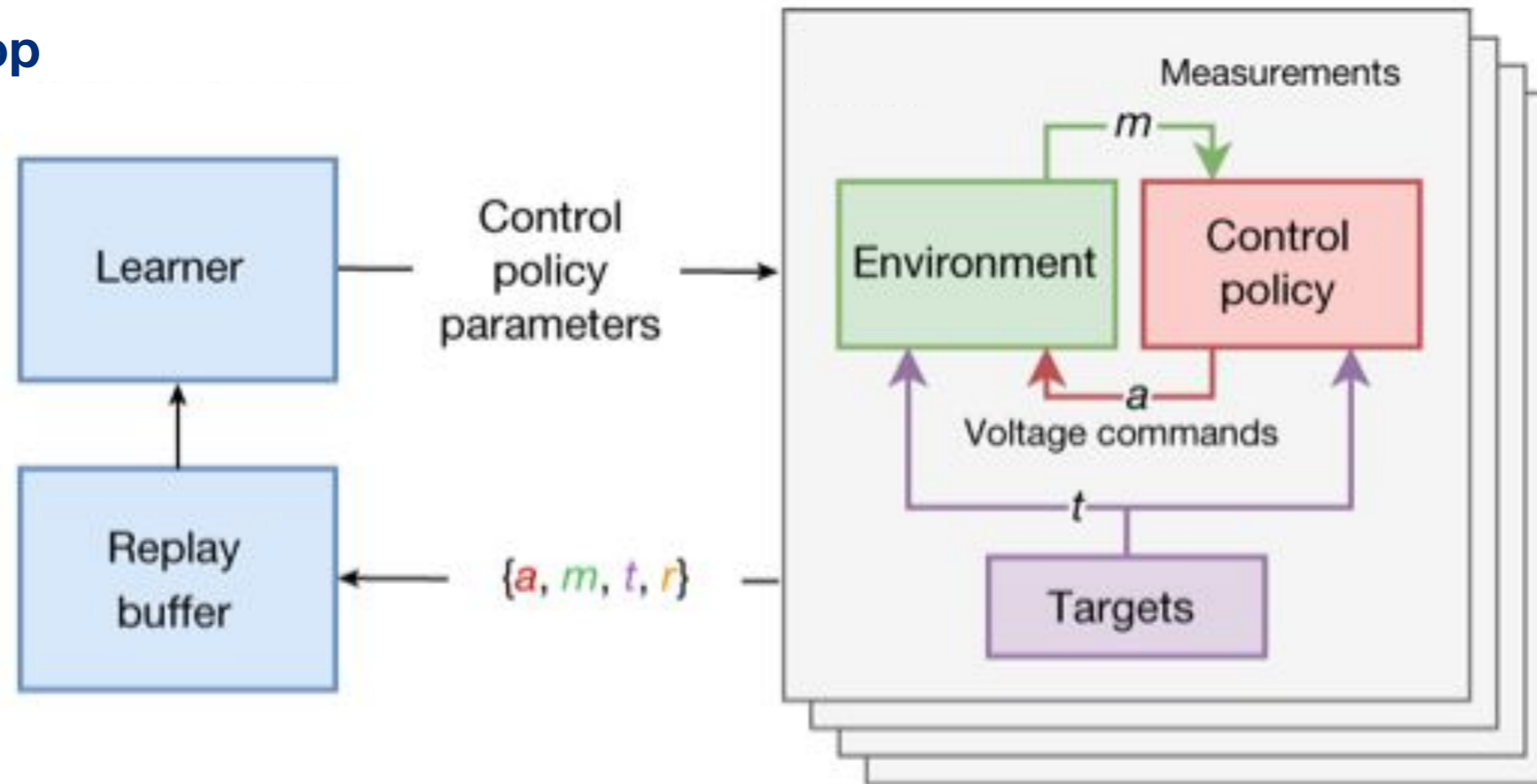
- Critic estimates the *Q function* from data generated by interacting with environment
- Actor learns a policy  $\pi$  by ascending (via gradient-based methods) the *Q function* learned by the critic
- Actor critic approaches allow asymmetry between actor and critic
  - The critic can benefit from privileged information unknown to the actor

# Actor-Critic Methods

- **Critic estimates the *Q function* from data generated by interacting with environment**
- **Actor learns a policy  $\pi$  by ascending (via gradient-based methods) the *Q function* learned by the critic**
- **Actor critic approaches allow asymmetry between actor and critic**
  - **The critic can benefit from privileged information unknown to the actor**
- **We use Maximum A Posteriori Policy Optimisation (MPO, Abdolmaleki et al. 2018)**

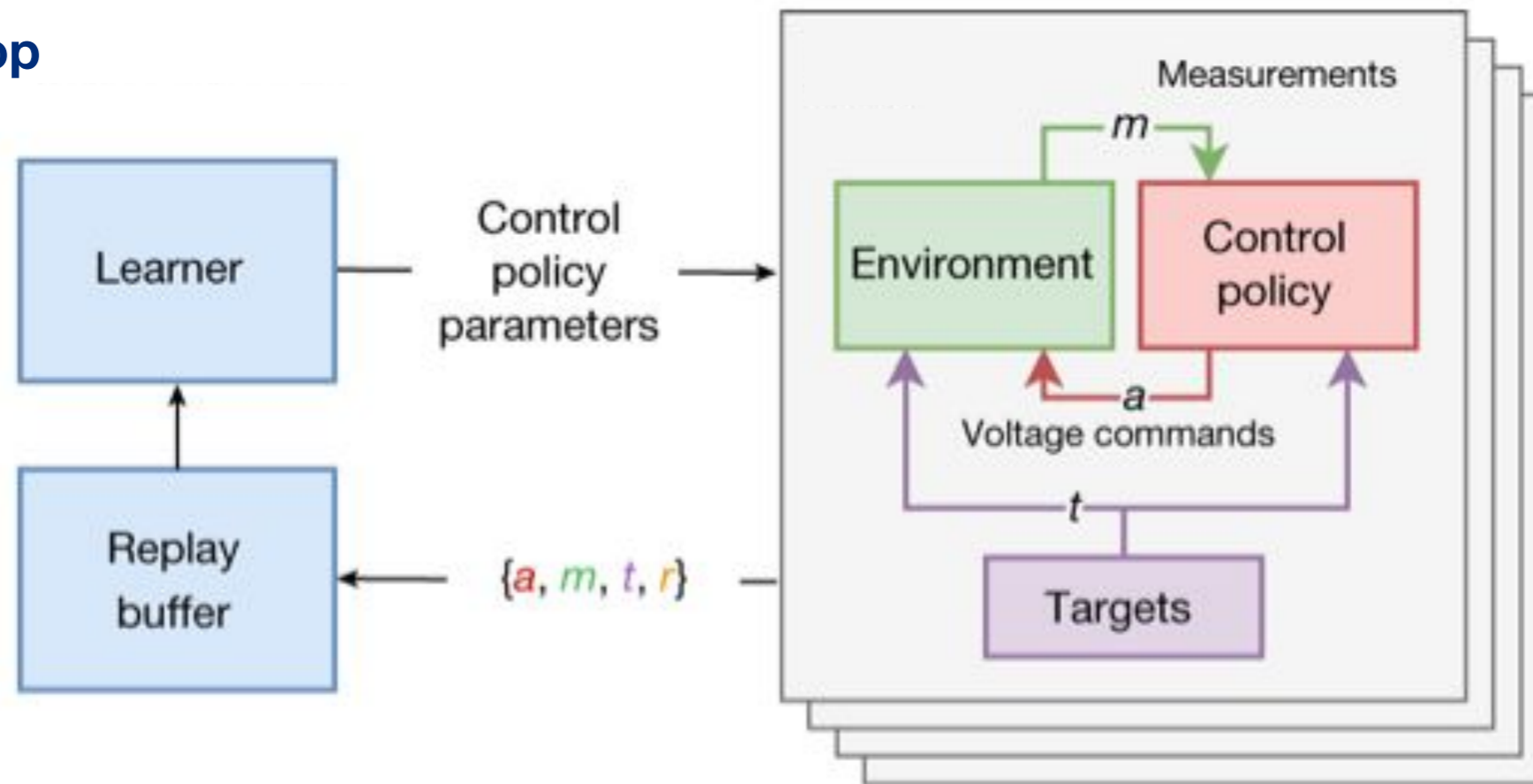
# Creating an Agent

## Learning Loop

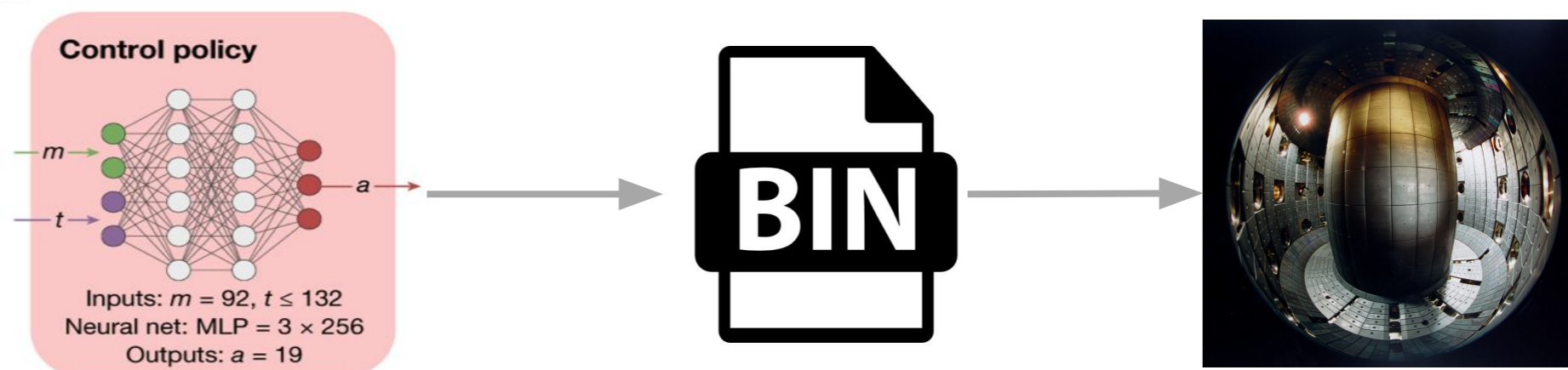


# Creating an Agent

## Learning Loop



## Deployment



# Transferring simulation-trained agents to the TCV tokamak

**Deployment mostly just worked!**

**Iterations on reward design needed to achieve stability**

**Simulator upgrades required to attain a successful agent**

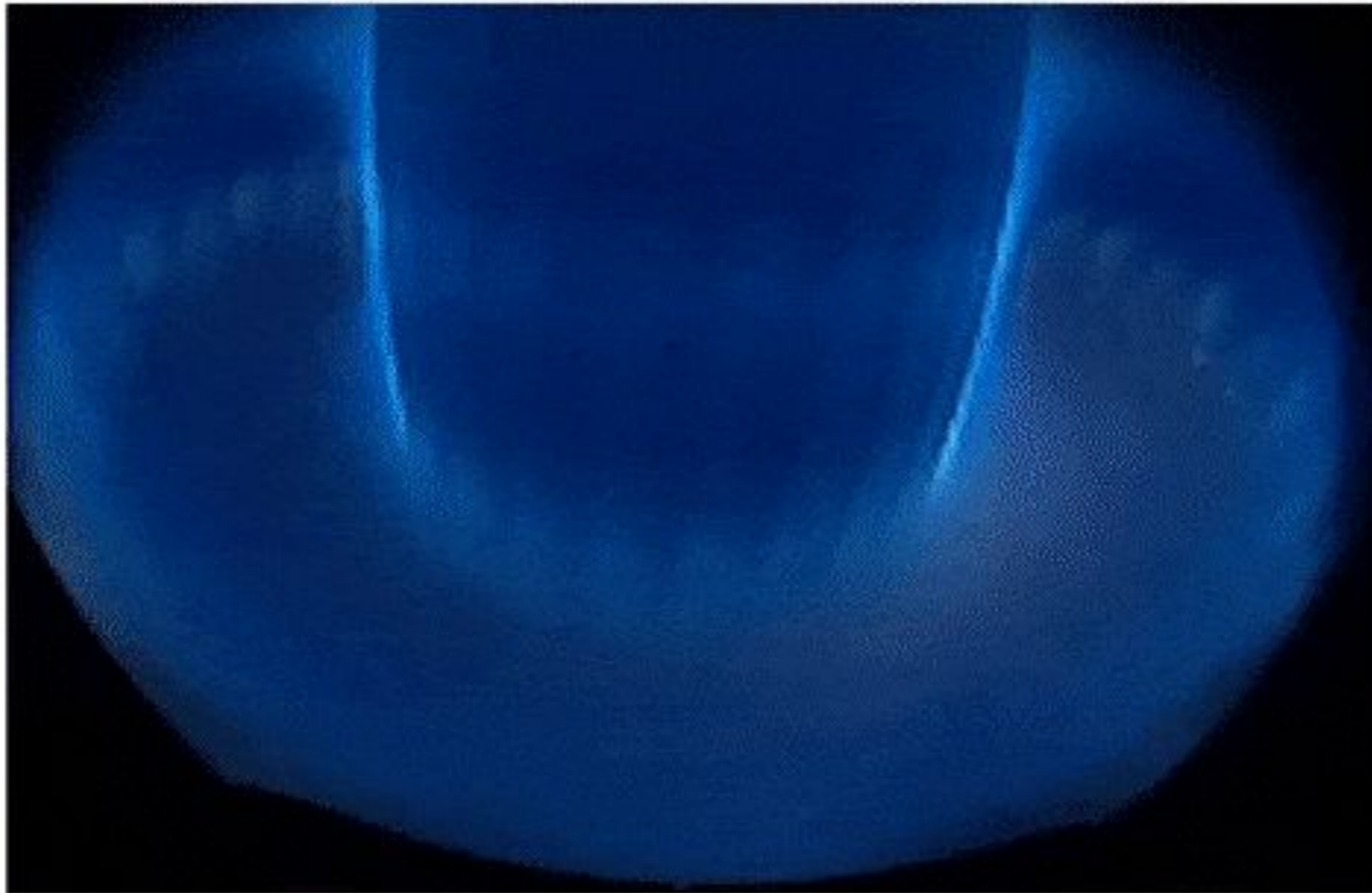
**Targeted environment variation:**

**Measurement noise**

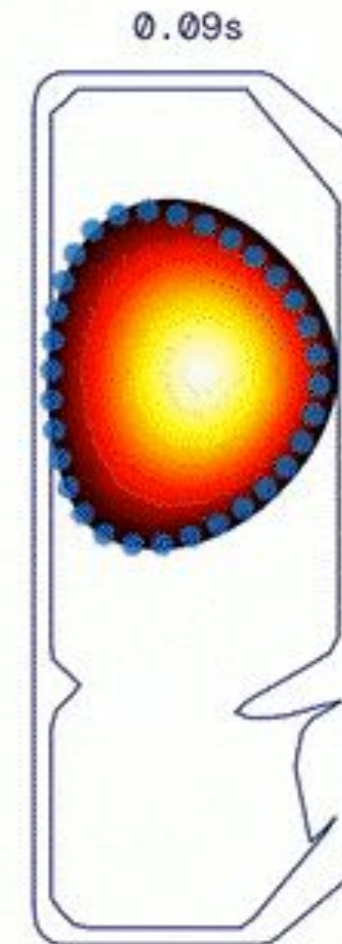
**Plasma parameters (resistivity, plasma pressure ratio ...)**

**Power supply**

# Result - demonstration shot

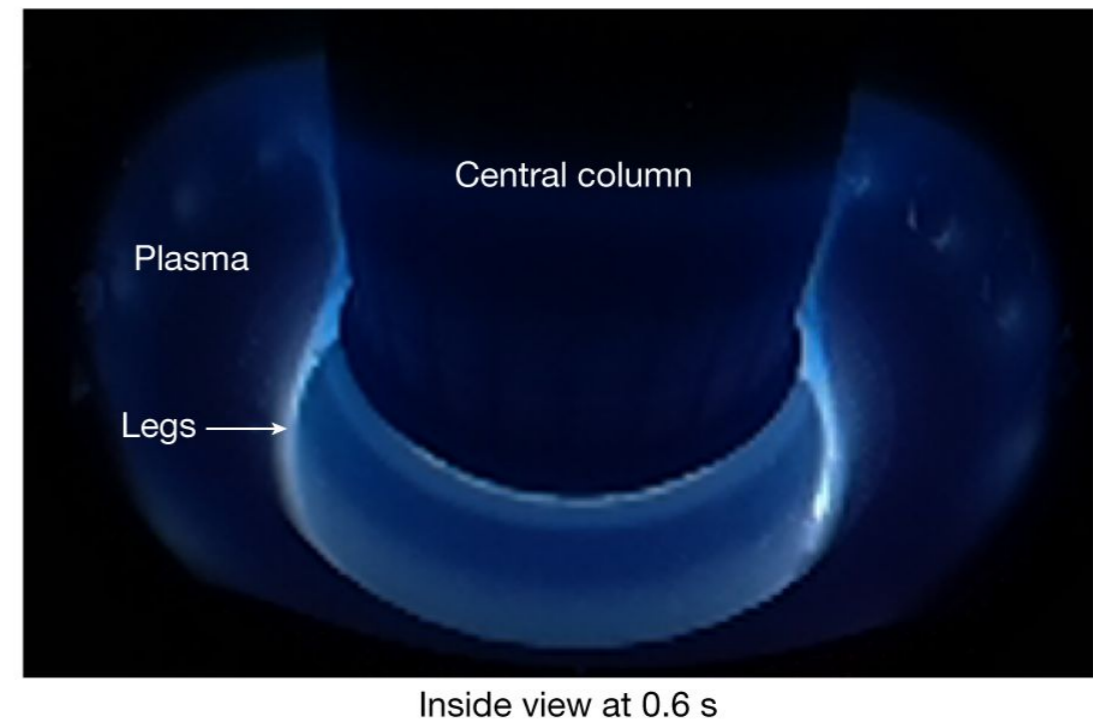
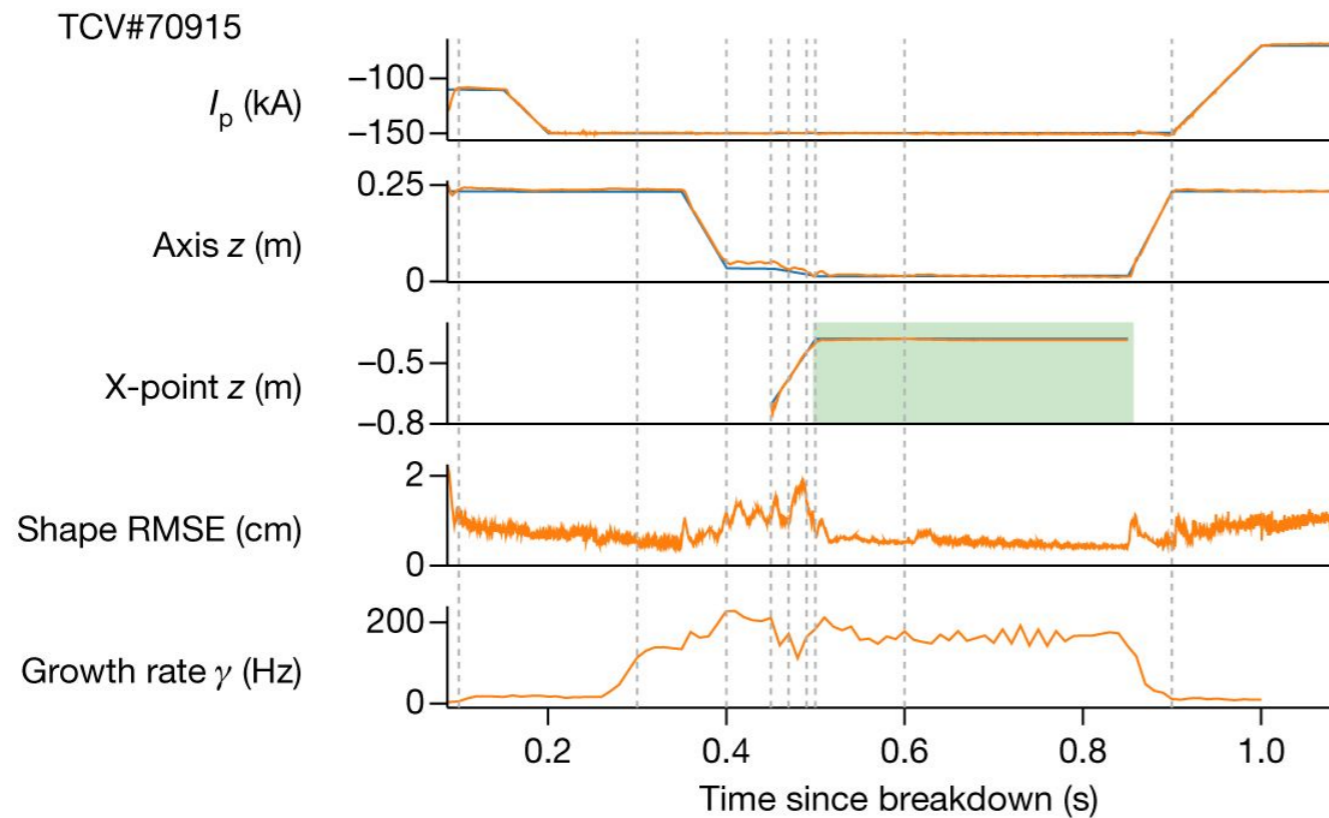
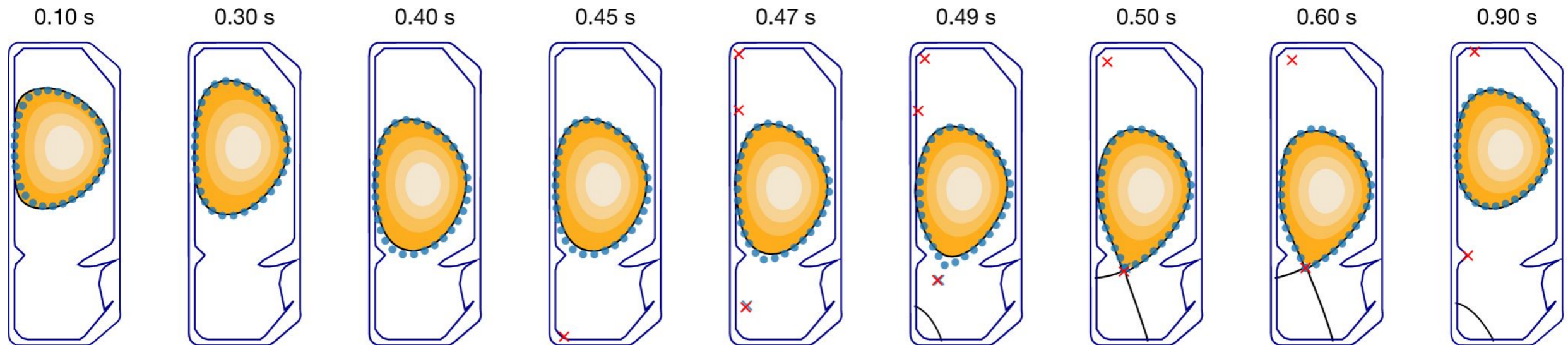


View from inside the tokamak



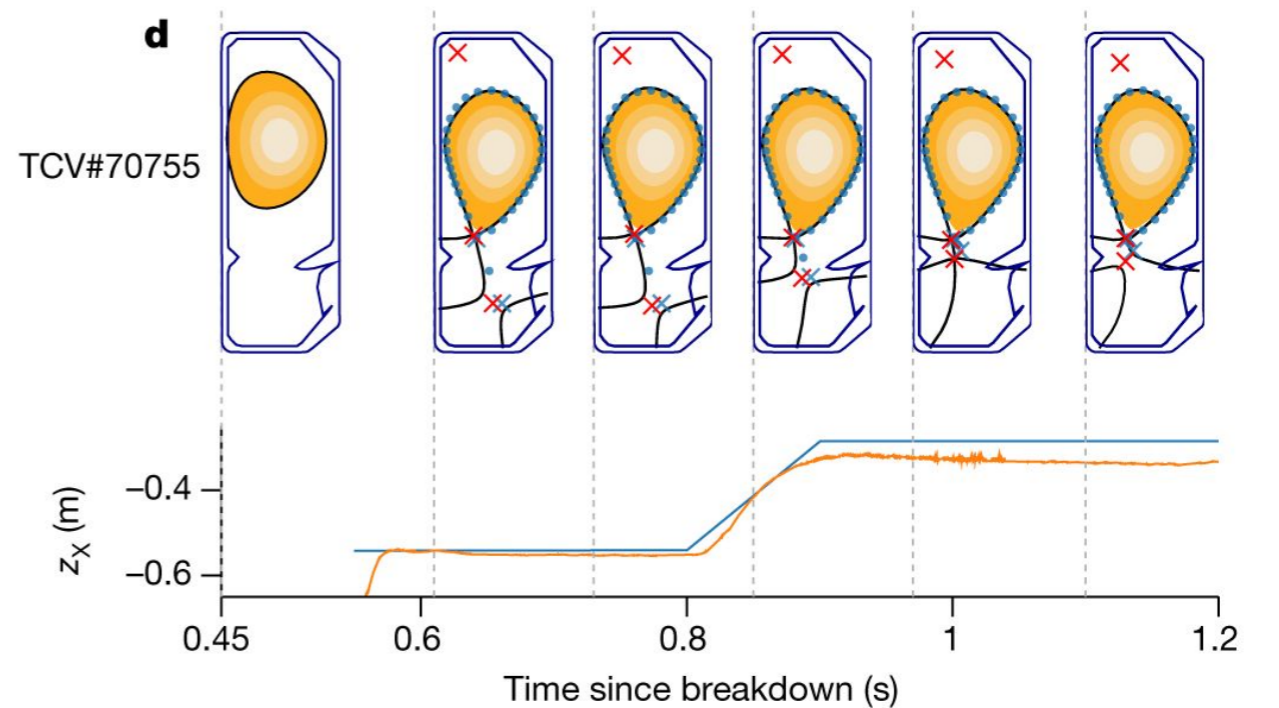
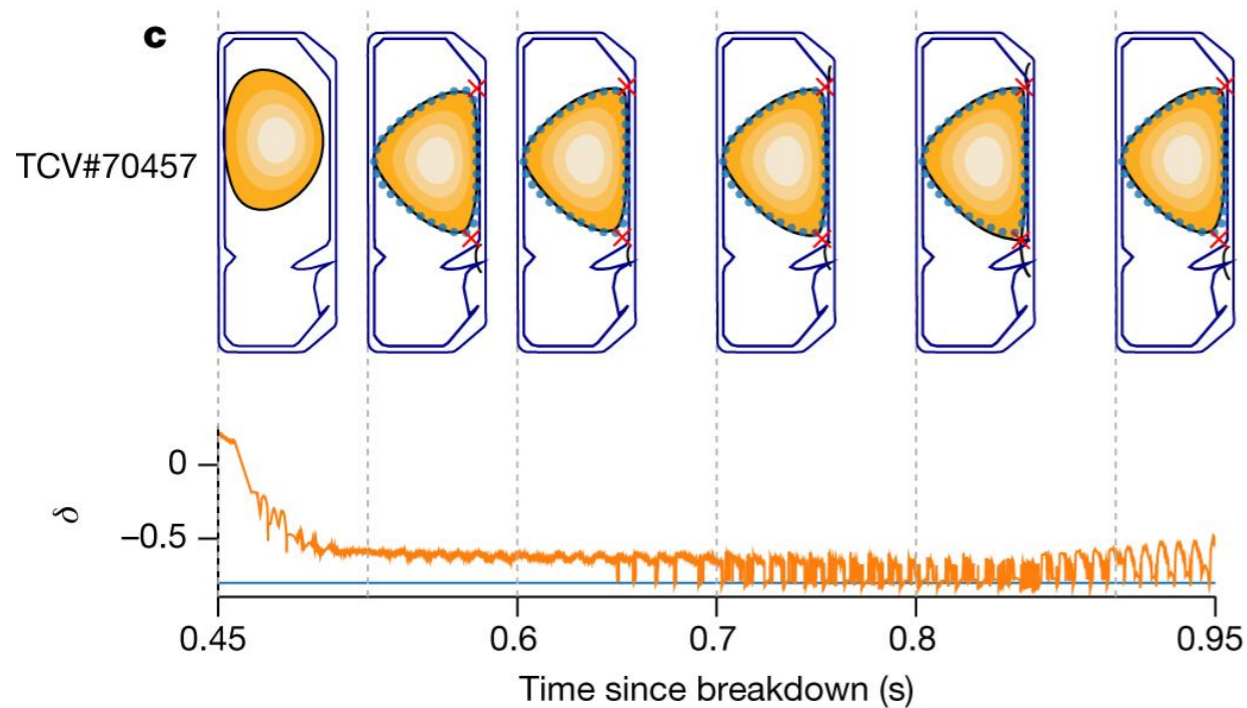
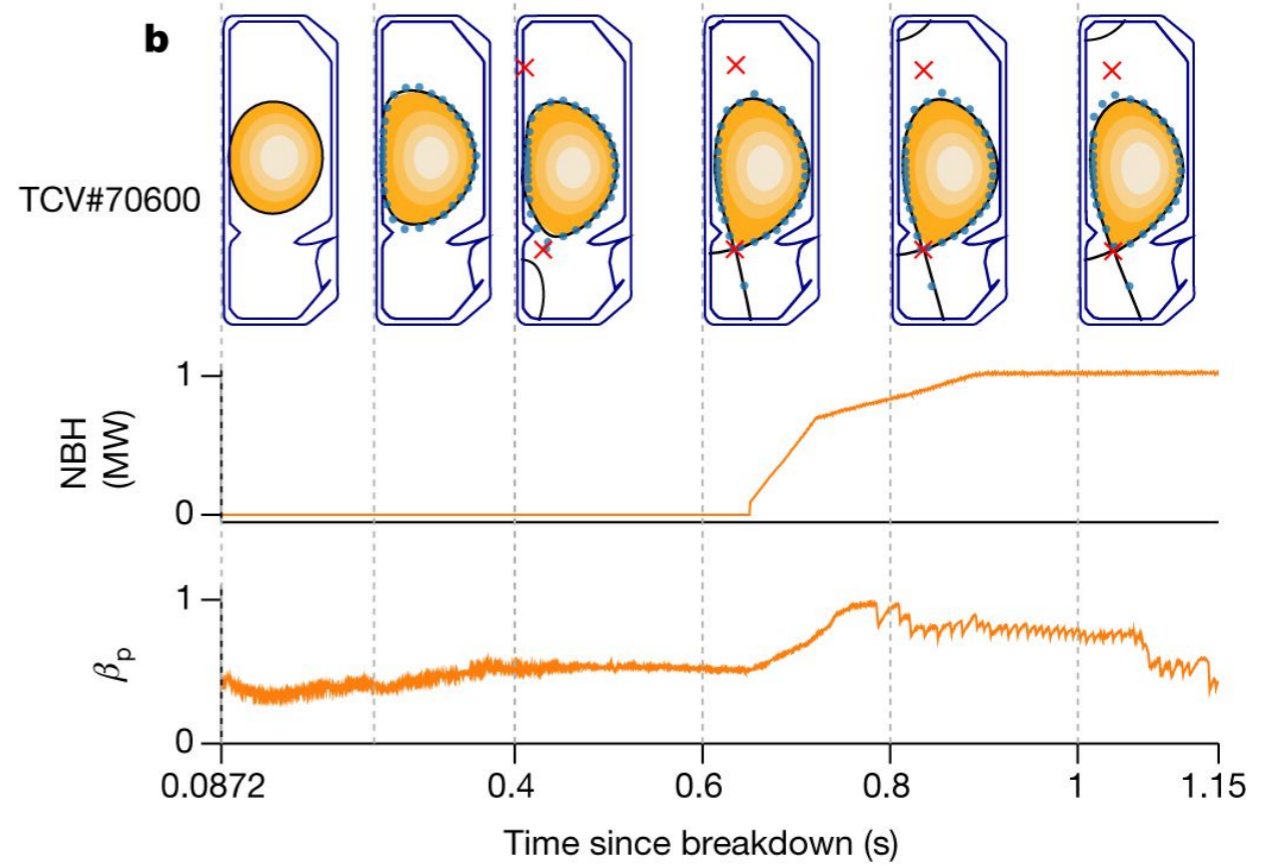
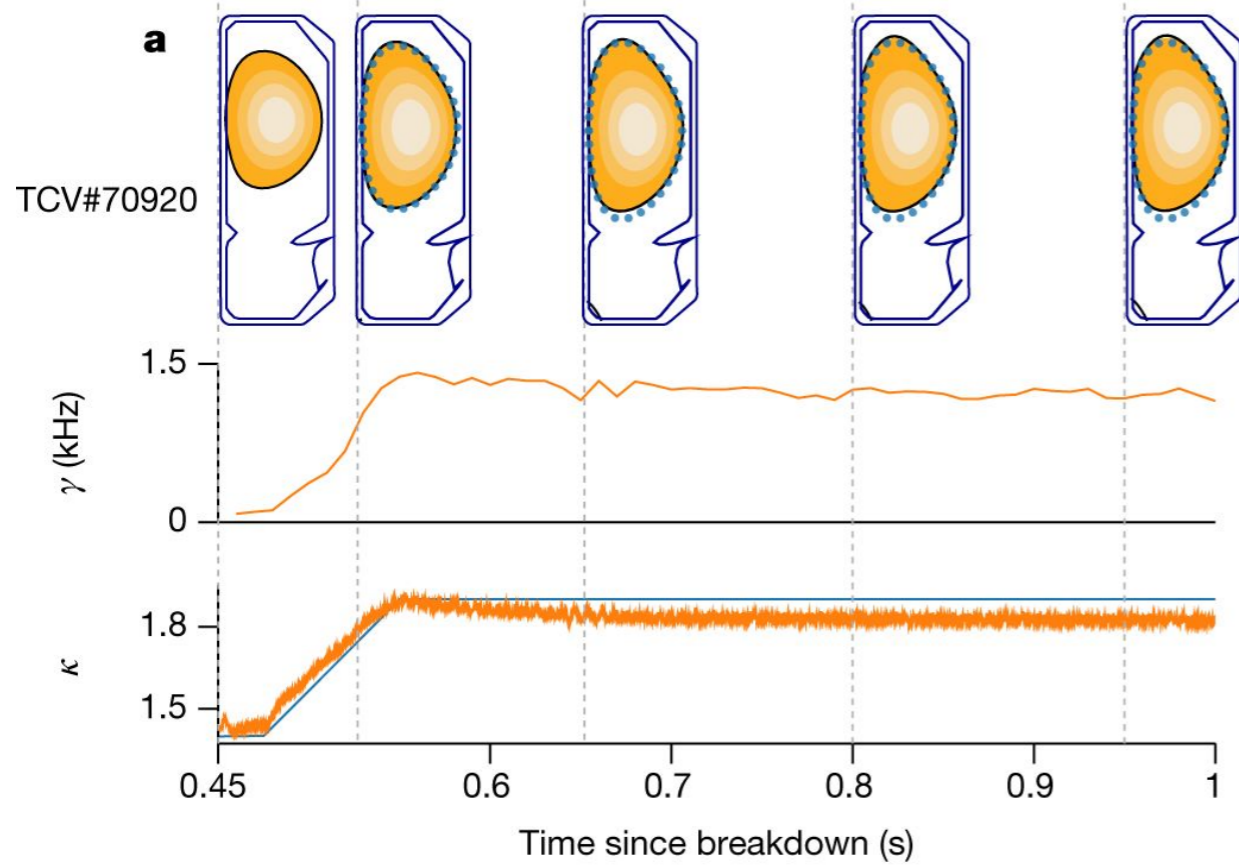
Plasma state reconstruction

# Demonstration shot

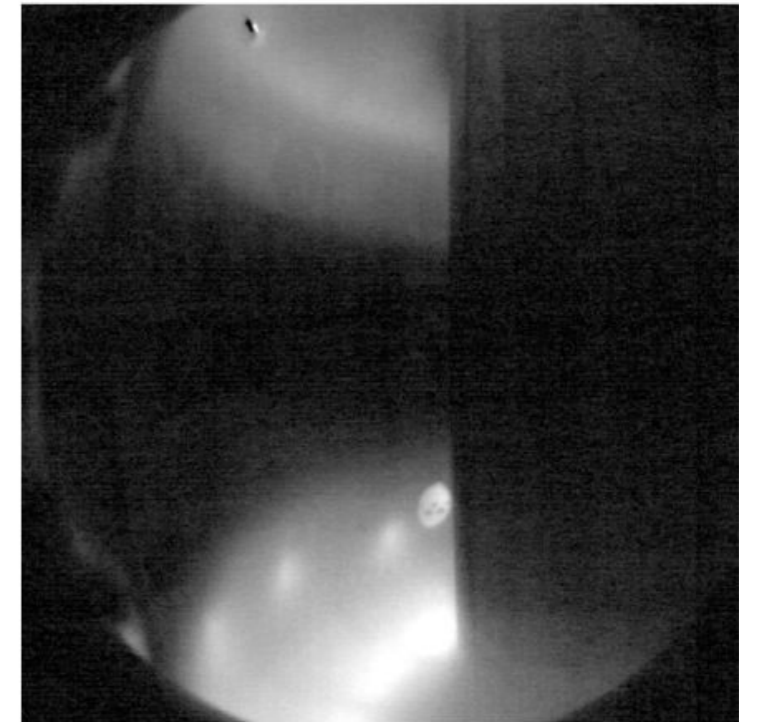
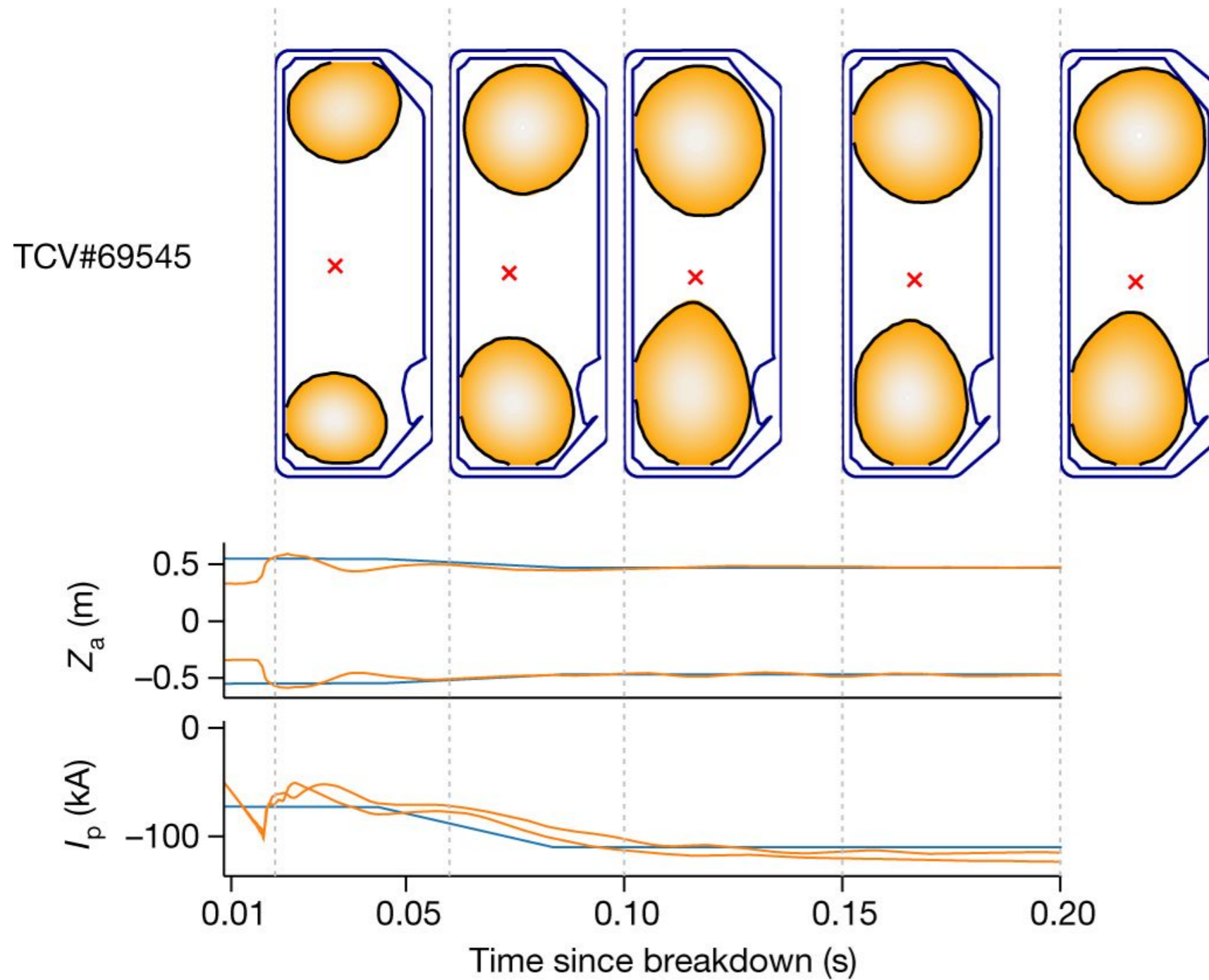




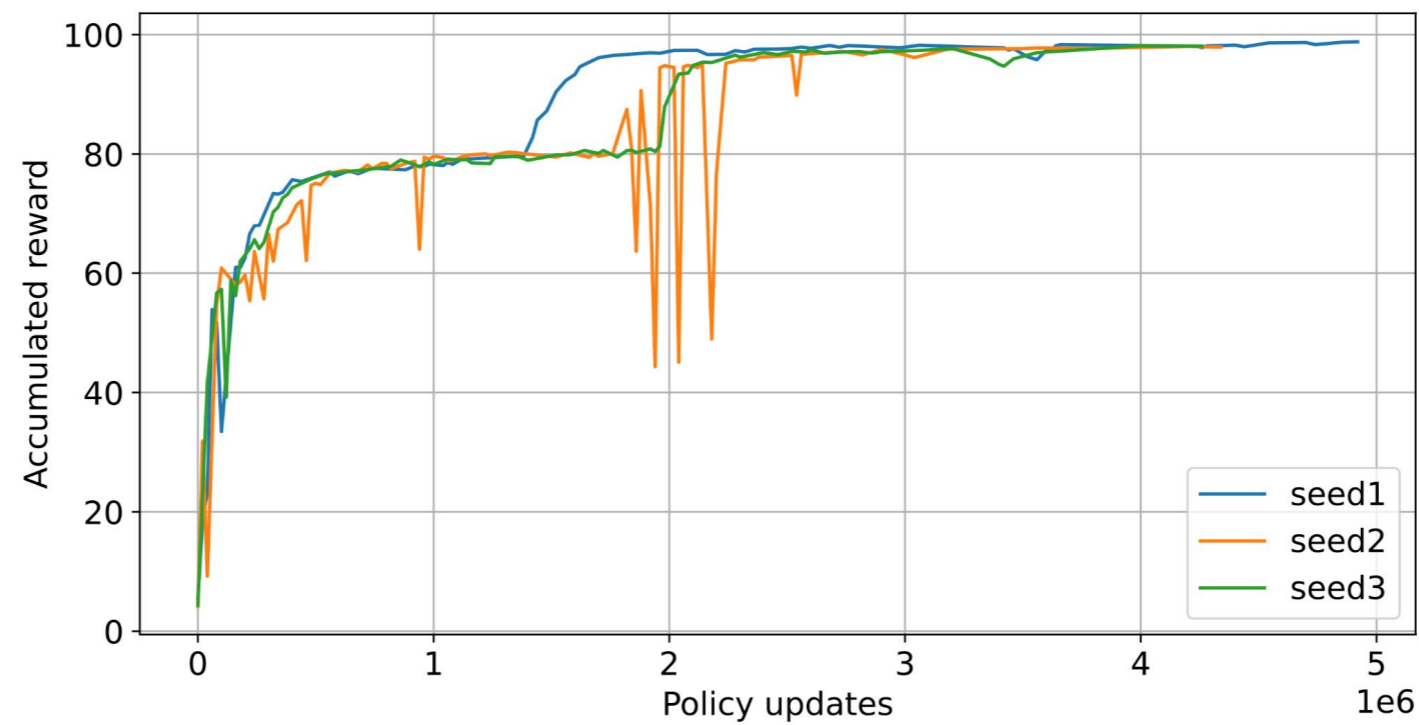
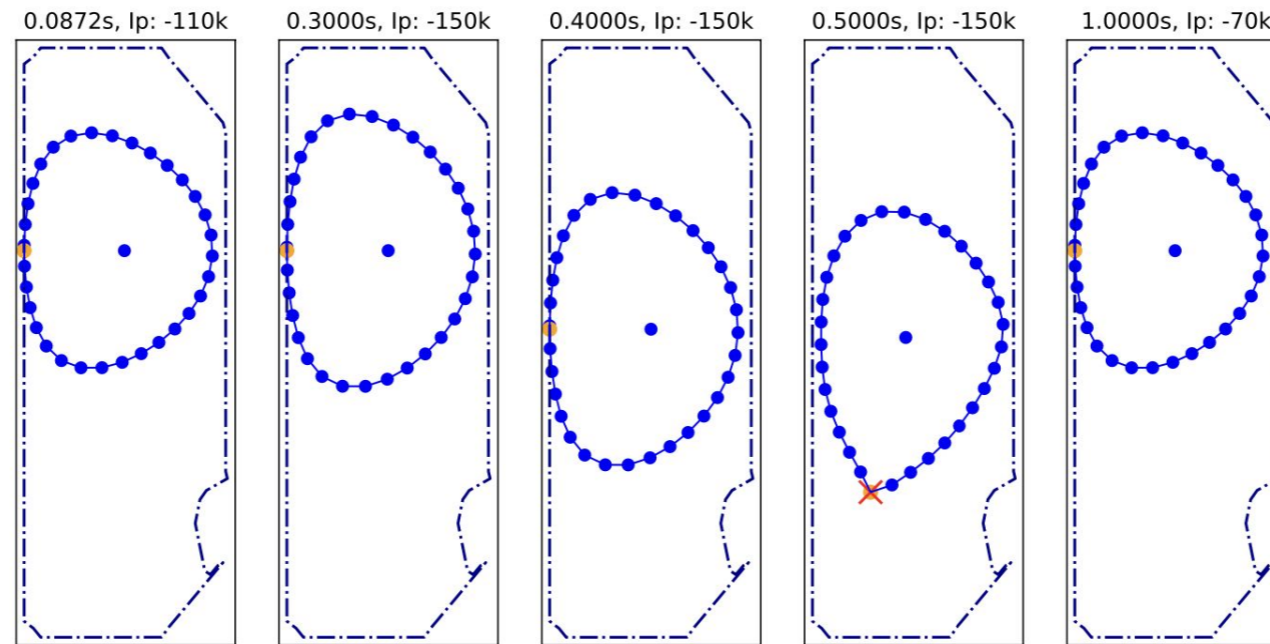
# Various plasma shapes controlled in TCV with RL



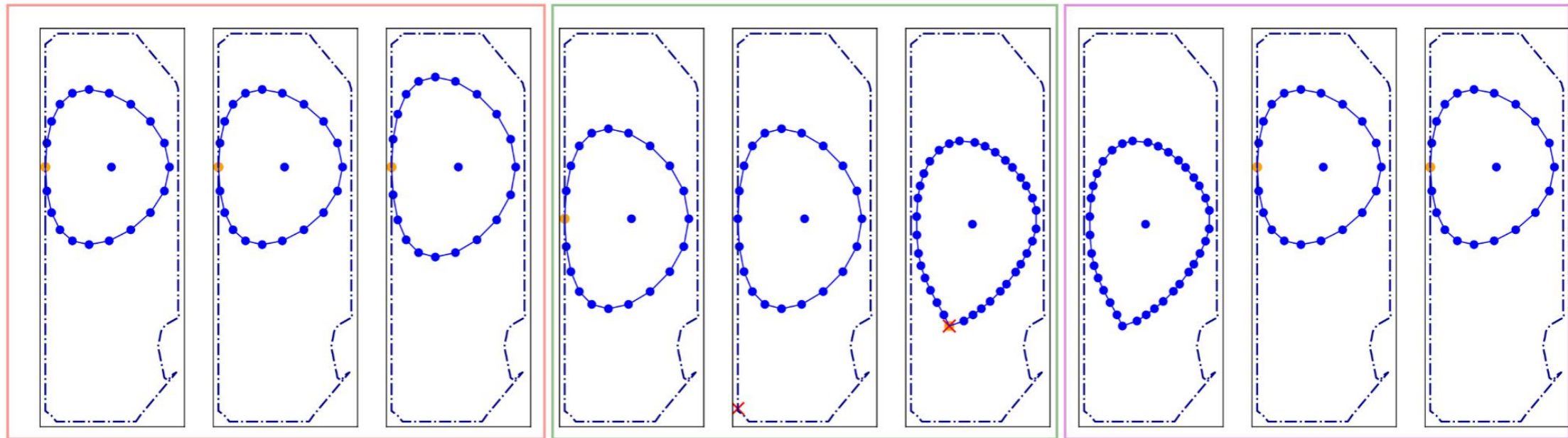
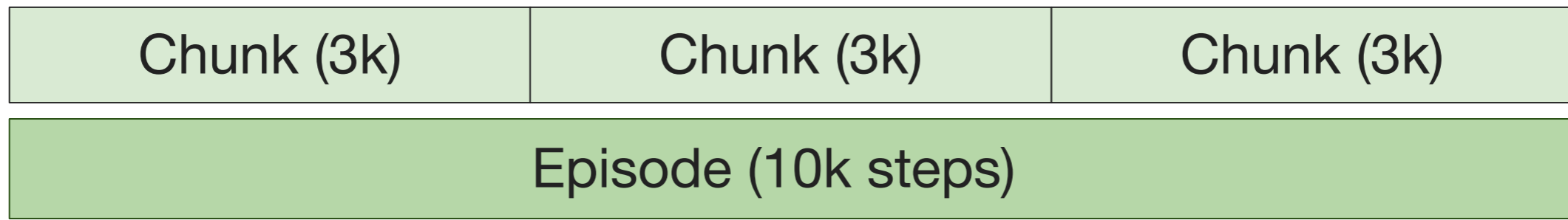
# Opening new frontiers for TCV: Droplet plasmas



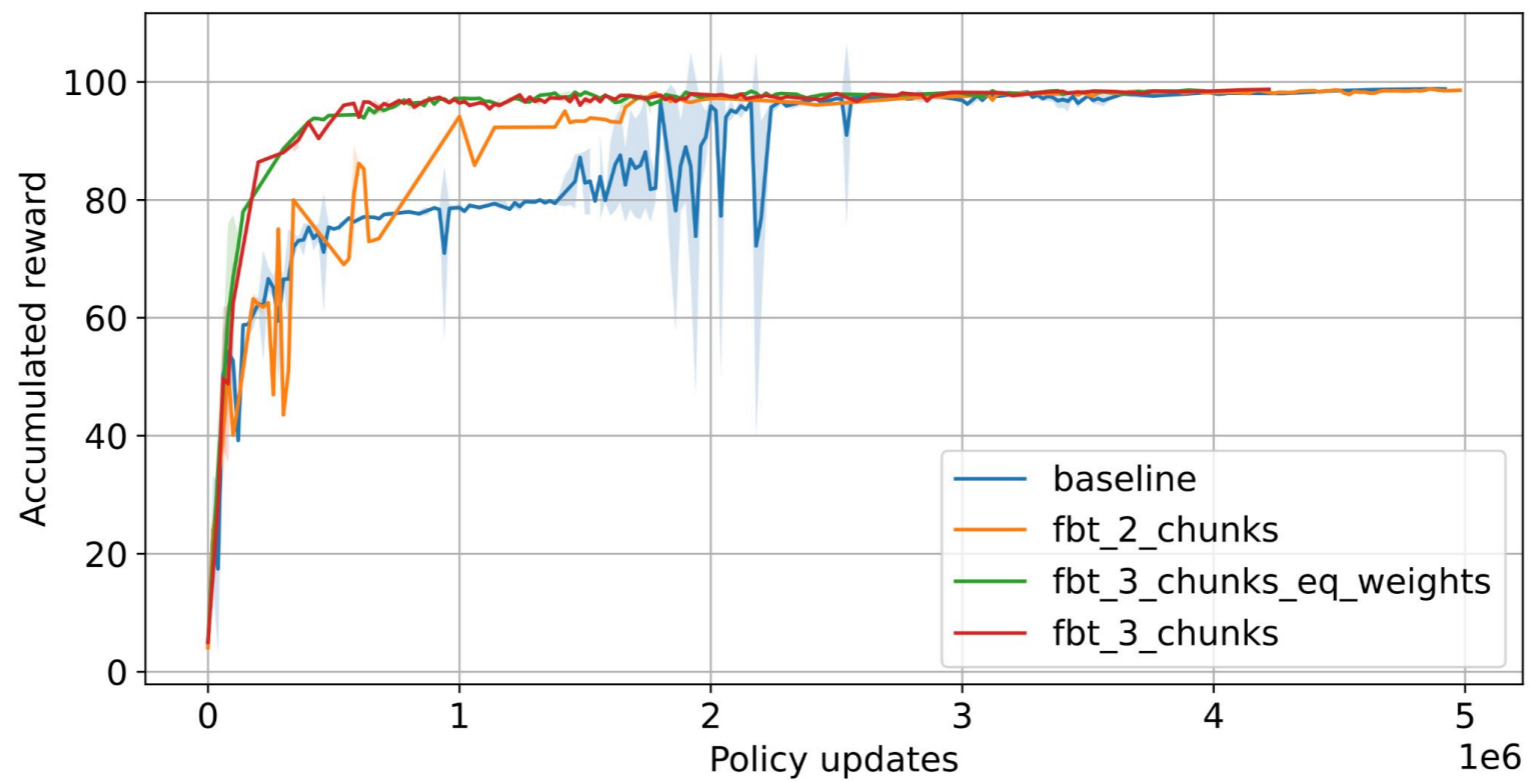
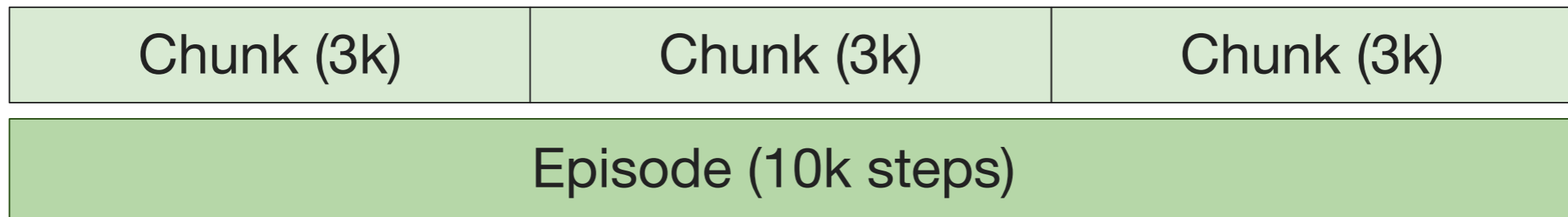
# Exploration Challenges



# Episode Chunking



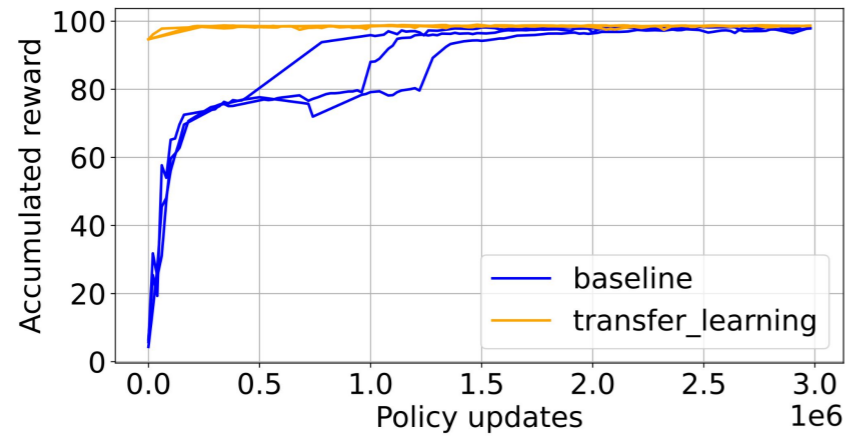
# Episode Chunking



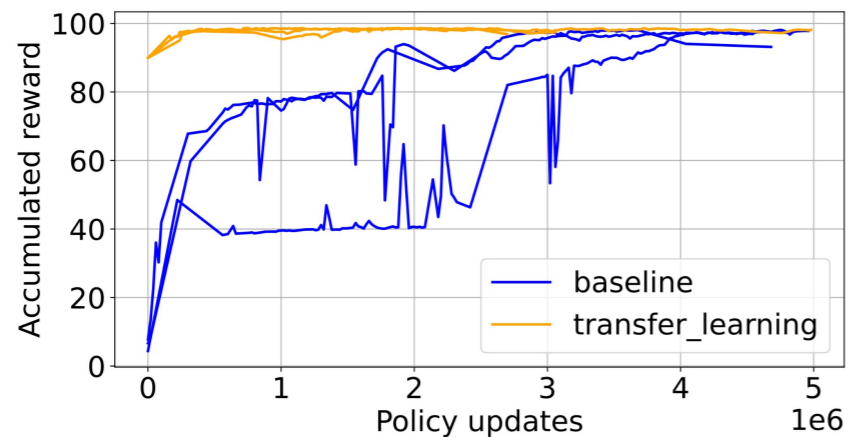
# Transfer Learning

## Shift Ip

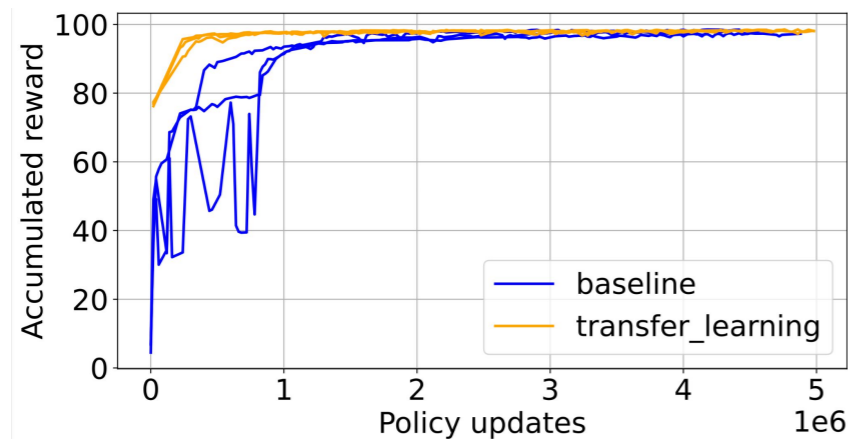
10kA



20kA

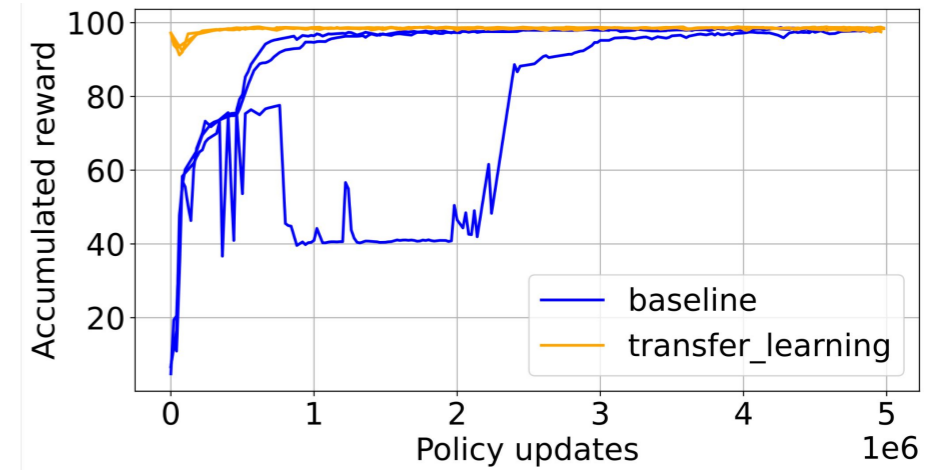


50kA

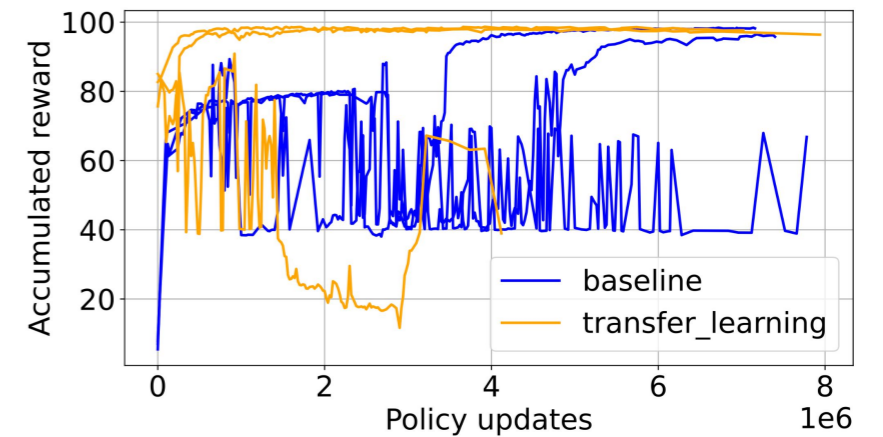


## Shift Shape

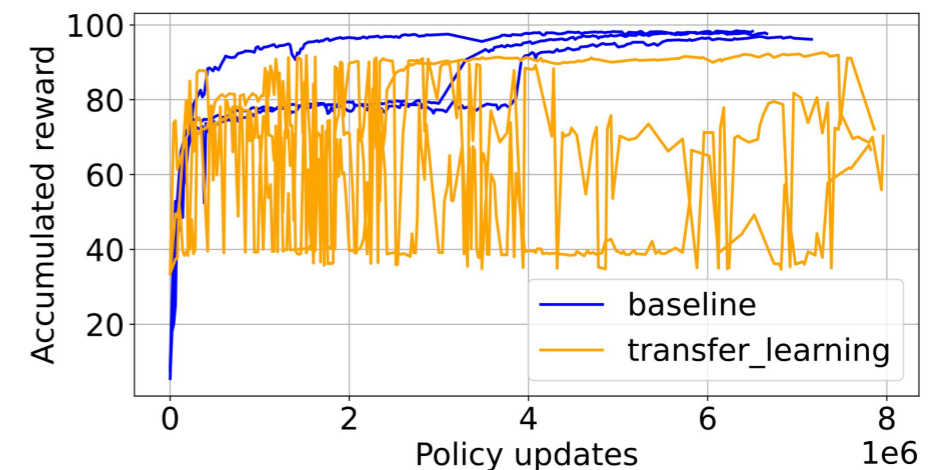
2cm



10cm



20cm



# Shape Accuracy

## LCFS Error

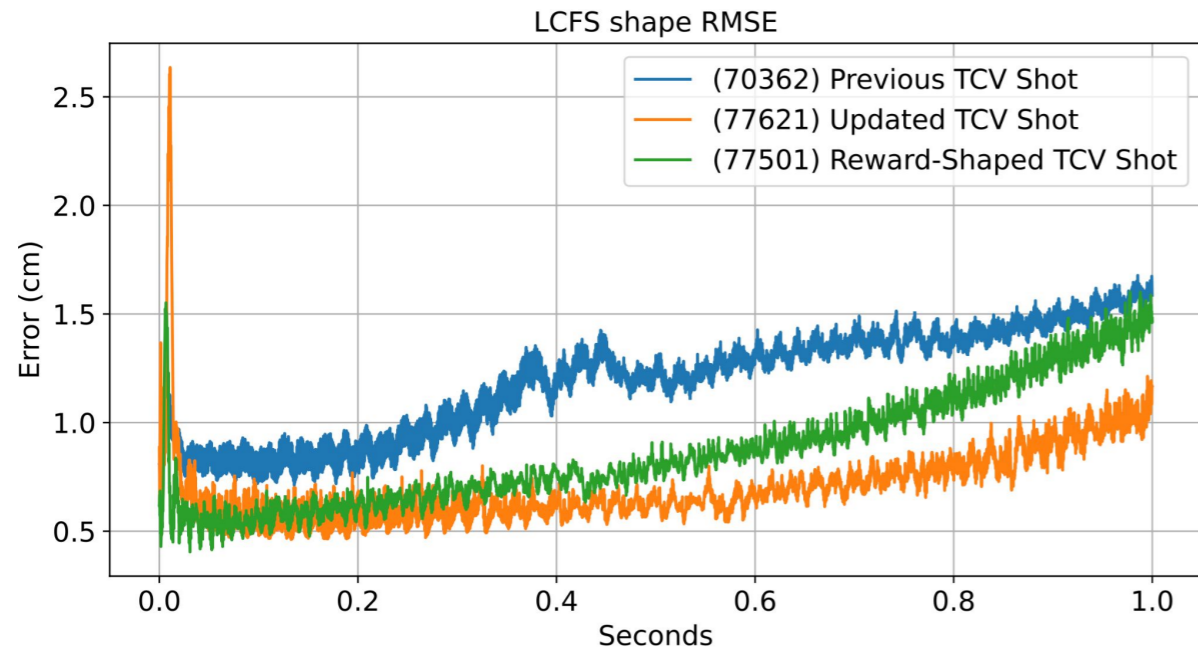
Experiment	$I_p$ Error (%)	LCFS Mean RMSE (cm)
Baseline	$0.353 \pm 0.221$	$0.567 \pm 0.221$
Narrow Reward	<b><math>0.238 \pm 0.076</math></b>	<b><math>0.201 \pm 0.057</math></b>
Reward Schedule	$0.450 \pm 0.321$	$0.490 \pm 0.196$

## X-point accuracy

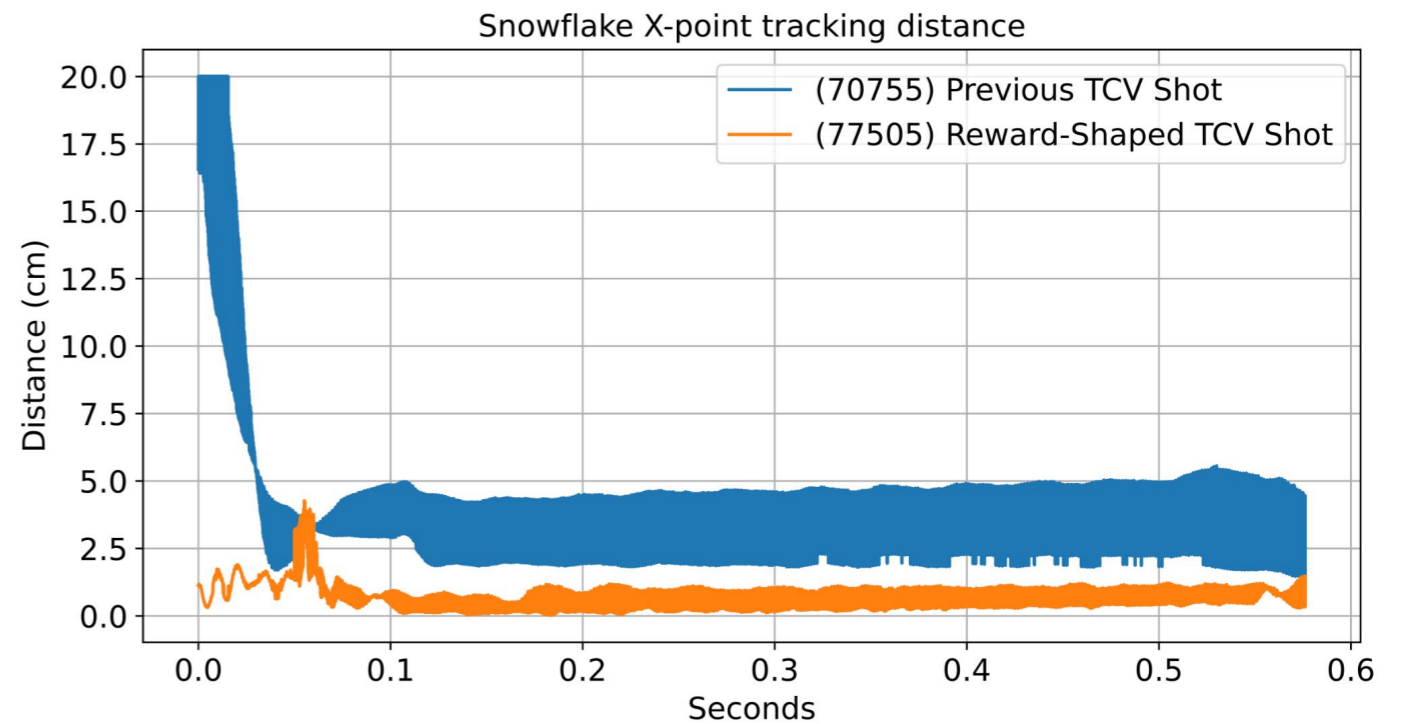
Experiment	$I_p$ Error (%)	LCFS Mean RMSE (cm)	X-point Location Error (cm)
Baseline	$0.848 \pm 1.710$	$1.122 \pm 1.460$	$0.669 \pm 0.491$
X-Point Fine Tuned	$0.717 \pm 0.624$	$0.845 \pm 0.097$	<b><math>0.289 \pm 0.027</math></b>
Narrow X-Point Reward	$6.143 \pm 4.602$	$4.536 \pm 3.268$	$1.199 \pm 1.102$
Additional Training	<b><math>0.502 \pm 0.423</math></b>	<b><math>0.723 \pm 0.159</math></b>	$0.541 \pm 0.112$

# Hardware Shape Accuracy

## LCFS Error



## X-point accuracy





# Summary

- **Large Opportunities available in speeding up design**
  - If can avoid the exploration problem – Do!
  - Can jumpstart training from pre-existing data
  
- **Need to be careful with treating results in simulation as truth**
  - Can significantly increase accuracy in sim
  - Need to push again on improving sim2real transfer

# Contrasting Classic Control and RL

---

## Traditional controllers (MIMO PID)

Separate error for each control loop

Error computed online

Separate complex state estimators and tuning of multiple control loops

Domain knowledge required for problem definition and separate controller design

Tuning of several control parameters

(Usually) Clear relation between parameters and aspects of control performance

Integral control nominally gives zero steady-state error on desired quantities

## Reinforcement Learning Implementation

Single reward function

No explicit error signals or estimation

Joint (and potentially generalising) solution to entire stabilization/control problem

Domain knowledge is in simulator. Just define reward functions

Reward function engineering

Black-box agent

No certainty of zero steady-state errors in case of external disturbances

# Outlook

- **Generalist agents - No need to retrain for new references**
- **Expand simulator capabilities to broaden the horizon of possibilities**
- **Co-design: simultaneously optimize tokamak design (plasma shape, sensors, coil, vessel placement) together with controller**

# Conclusions

- **Demonstrated RL for closed-loop magnetic control of tokamak plasmas, trained in simulation and tested on a real device**
  - Implementing 10kHz controller with 100+ measurements, 20 actions is a milestone for RL on real-world systems in terms of complexity
  - Models are sufficiently accurate to perform the required simulations
- **Bright future for more applications of RL**
  - For accelerating fusion science: improving plasma performance & design new devices
  - For application to more complex real-world systems, in particular where good models exist

**SWISS  
PLASMA  
CENTER**



**Successful multidisciplinary collaboration!**

**Tight integration between teams to understand and control this  
challenging system**