

# Reinforcement Learning for Science

April 1

[Tailin Wu](#), Westlake University

Website: [ai4s.lab.westlake.edu.cn/course](http://ai4s.lab.westlake.edu.cn/course)

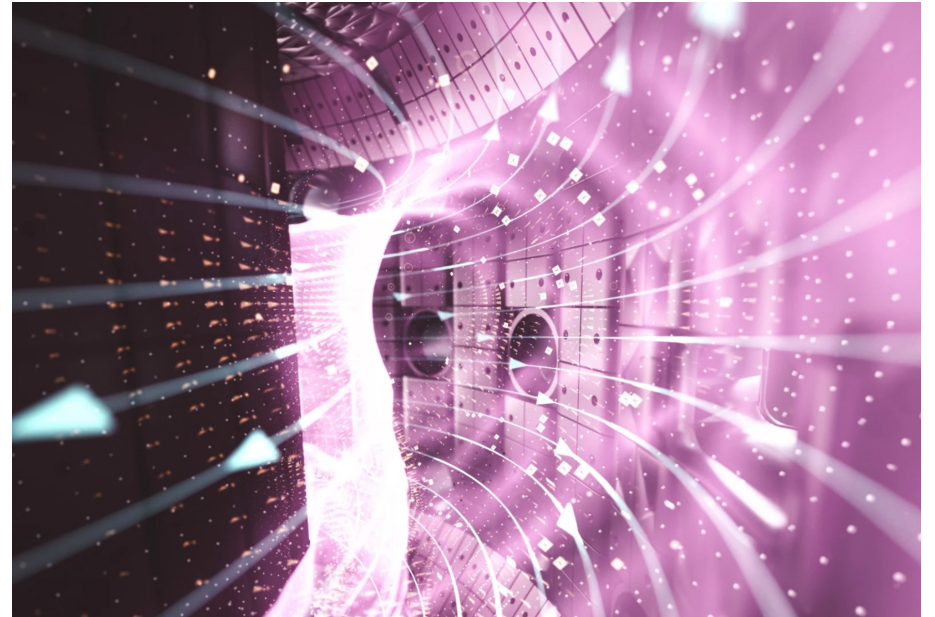


Image from: DeepMind

# Project guideline

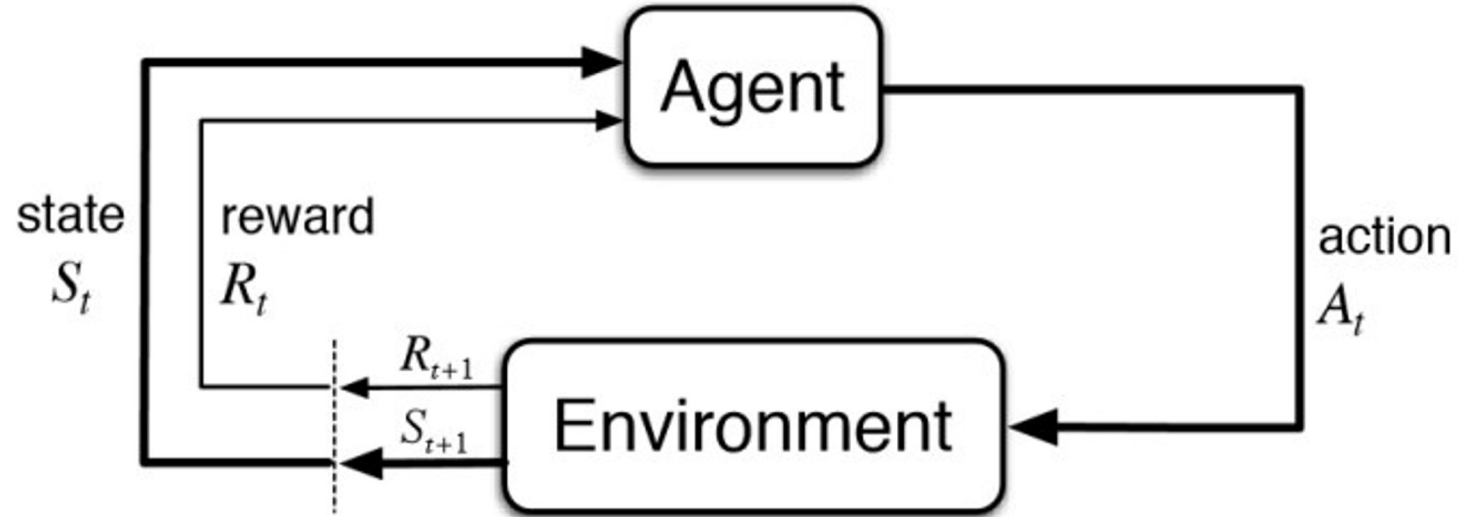
## **Mid-term course project design (April 15 & 22):**

- Give a presentation (10min) that formulates the problem for the **5 questions**, each with 1-2 slides:
  1. What is the problem?
  2. Why is it important
  3. Why is it hard?
  4. What is the limitation of the prior method?
  5. What are the main components of the proposed method?

Then detail the proposed method (3-4 slides) that uses an AI technique to solve the problem.

**Each group:** Presentation 10min + questions 5min

# Markov Decision Process (MDP): Setup



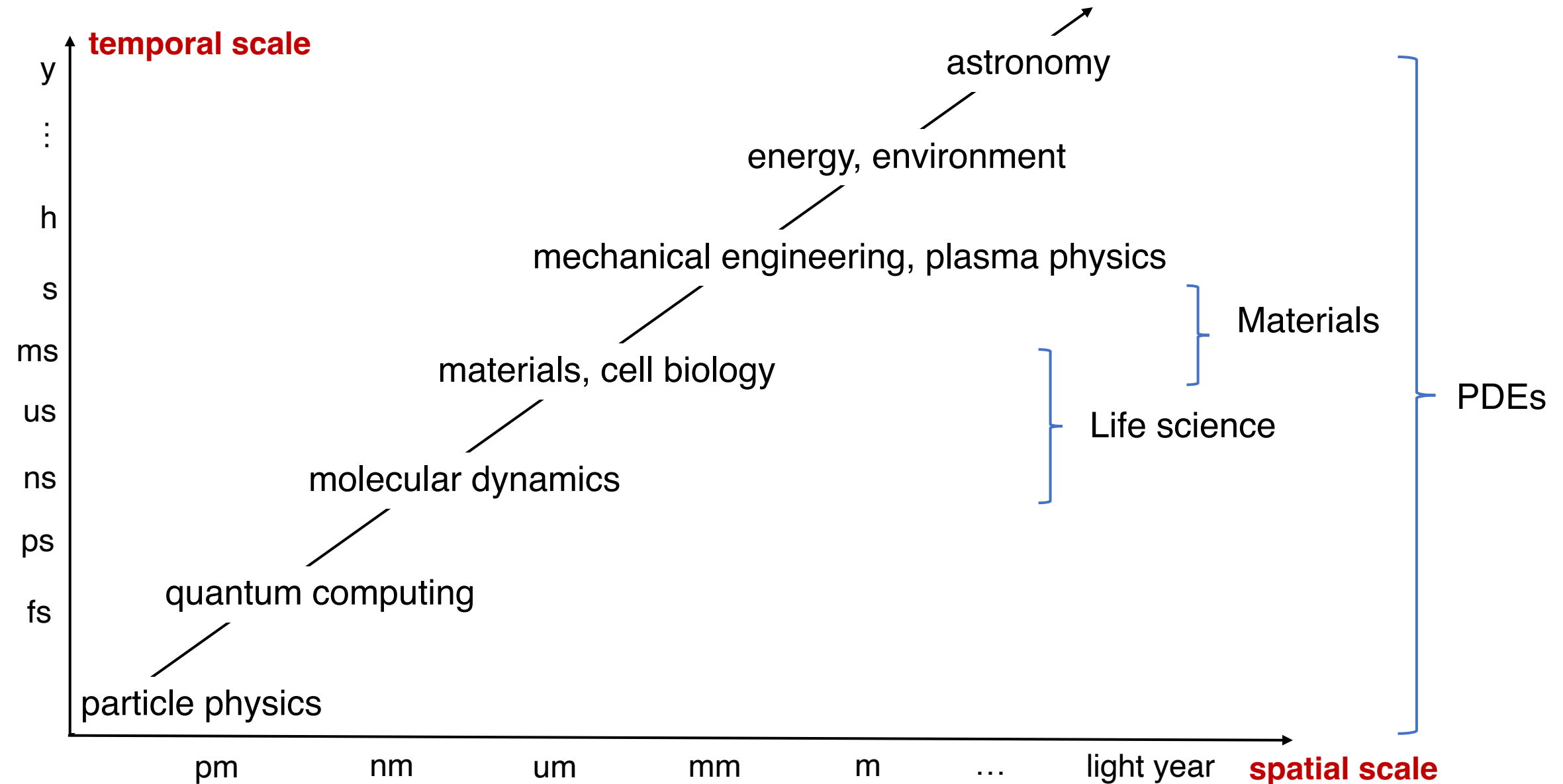
**Goal:** Maximize the long-term expected reward w.r.t. to the policy  $\pi(A_t|S_t)$

$$\max_{\pi(A_t|S_t)} \mathbb{E}_t[R_t]$$

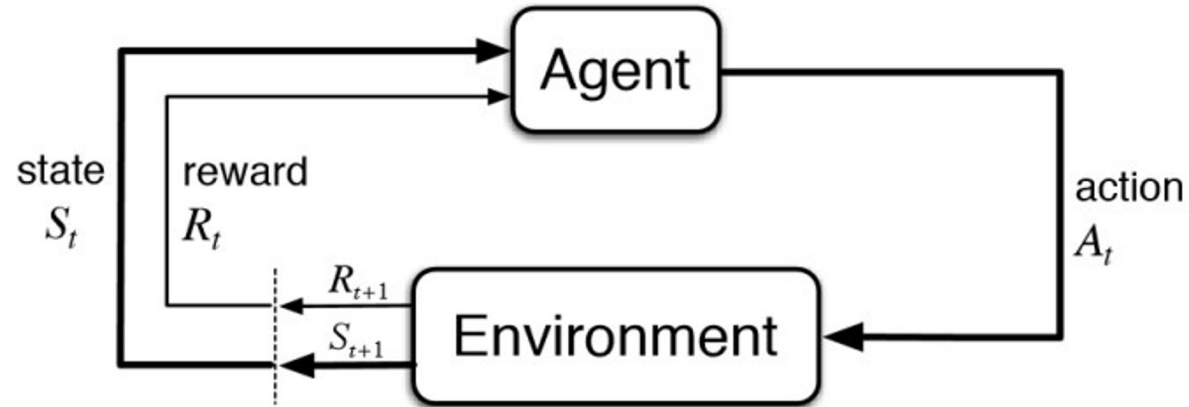
# Recent Deep RL papers in *Nature/Science*

Paper	Publisher	Application
<a href="#">Avoiding fusion plasma tearing instability with deep reinforcement learning</a>	<i>Nature</i> 2024	Tokamak control
<a href="#">Champion-level drone racing using deep reinforcement learning</a>	<i>Nature</i> 2023	Drone racing
<a href="#">Top-down design of protein architectures with reinforcement learning</a>	<i>Science</i> 2023	Protein design
<a href="#">Dense reinforcement learning for safety validation of autonomous vehicles</a>	<i>Nature</i> 2023	Autonomous driving
<a href="#">Magnetic control of tokamak plasmas through deep reinforcement learning</a>	<i>Nature</i> 2022	Tokamak control
<a href="#">Discovering faster matrix multiplication algorithms with reinforcement learning</a>	<i>Nature</i> 2022	Matrix multiplication
<a href="#">A graph placement methodology for fast chip design</a>	<i>Nature</i> 2021	Chip design
<a href="#">A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play</a>	<i>Science</i> 2018	Board game

# Application in AI for Science: from microscopic to macroscopic



# How to Apply RL in AI for Science



## 1. Define the task

### Specify:

- State  $S$
- Action  $A$
- Reward  $R$

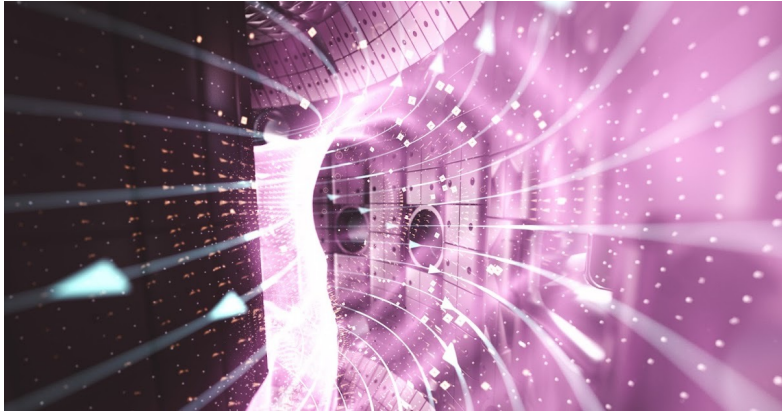
### Learn:

- Policy  $\pi_{\theta}(A|S)$

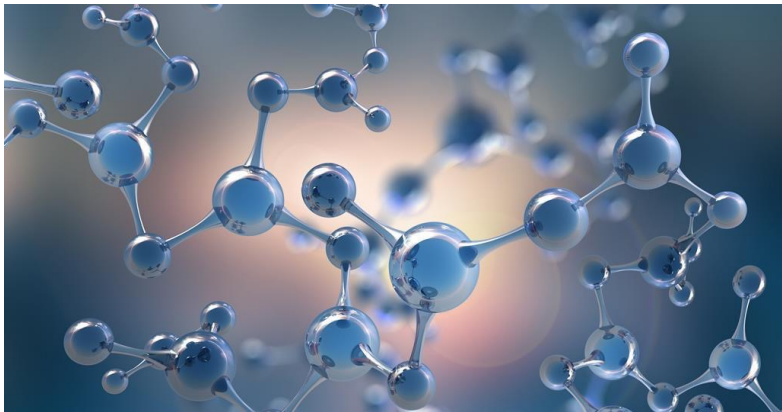
## 2. Choose an appropriate RL algorithm

# RL for Science: Case study

## 1. Deep RL for controlled nuclear fusion

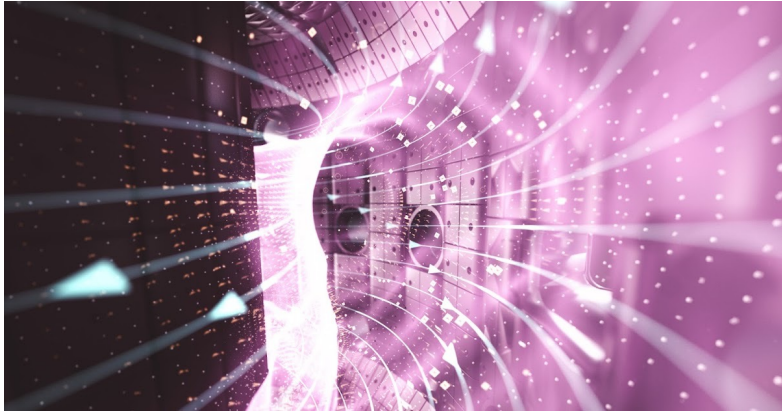


## 2. Deep RL for molecule design

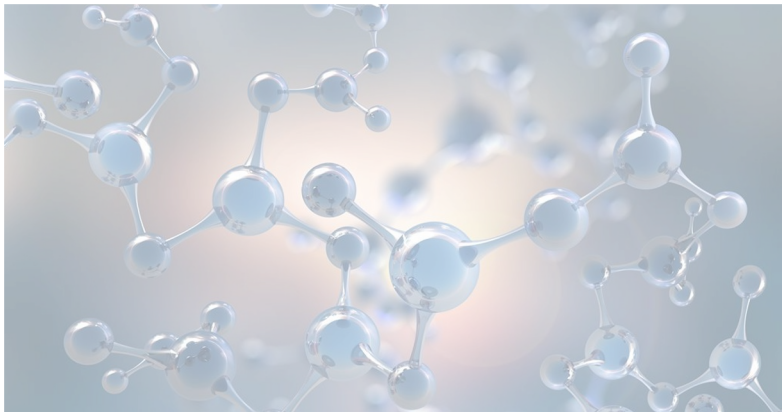


# RL for Science: Case study

## 1. Deep RL for controlled nuclear fusion



## 2. Deep RL for molecule design



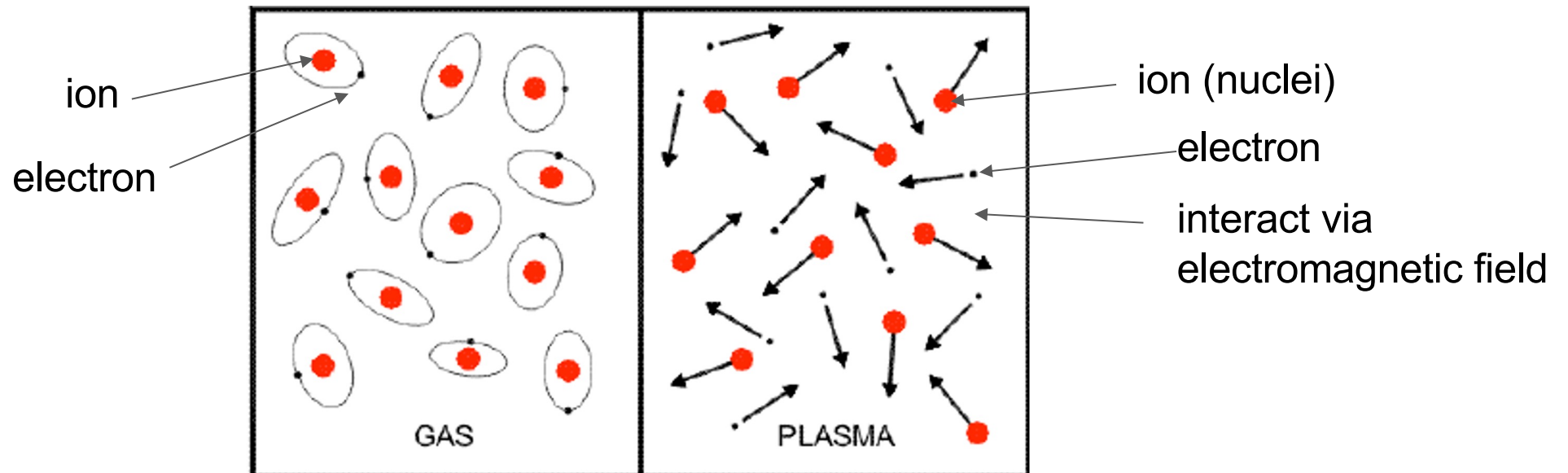


# Preliminaries: plasma (等离子体)

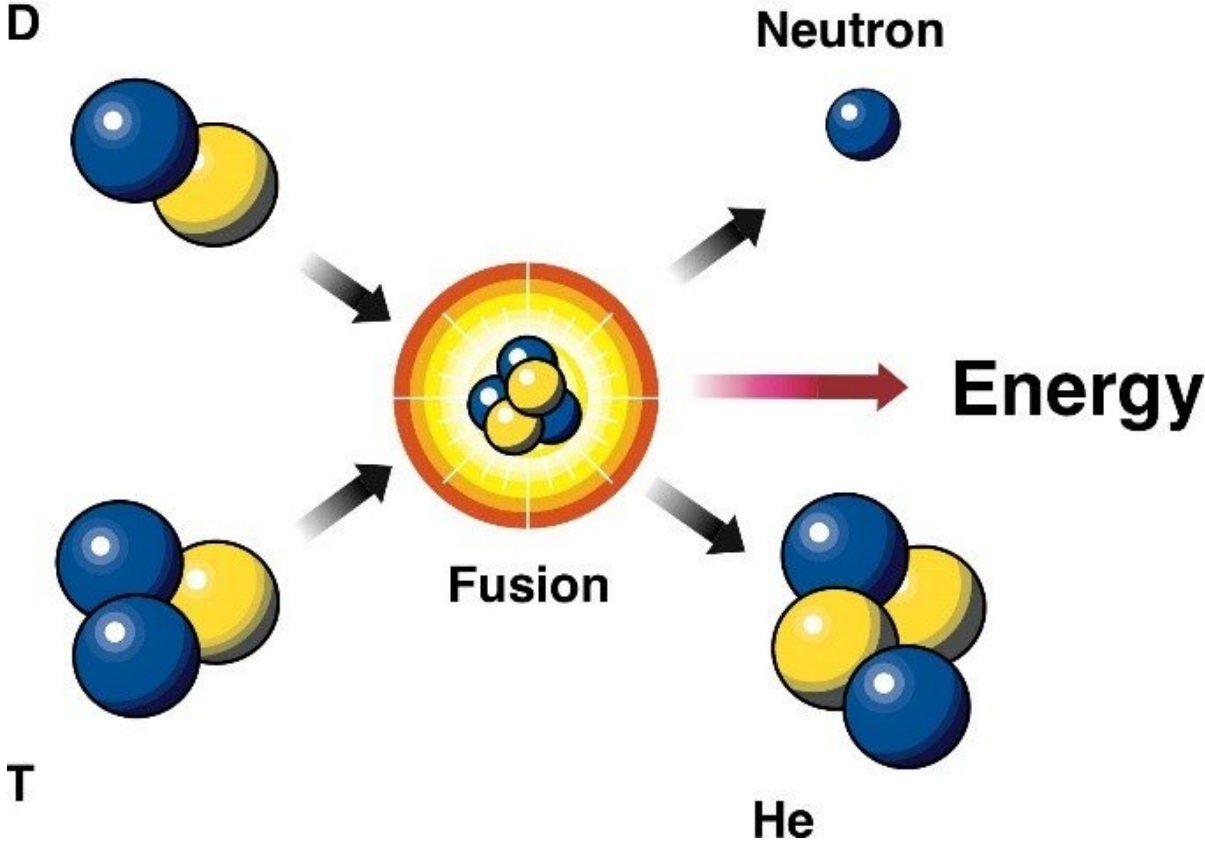
**Plasma:** Consisting of energetic ions and *free* electrons, interacted via electromagnetic (EM) field.

Examples: fire, lightning, sun, nuclear fusion

It is one of the four fundamental states of matter. It is the *dominant* form of ordinary matter in the universe.



# Preliminaries: Nuclear fusion



# Preliminaries: Why nuclear fusion?

## 1. Percentage of mass transferred to energy: $E = mc^2$

- Chemical: 0.0000001%
- Nuclear fission: 0.1%
- **Nuclear fusion: 0.4%**
- Black hole: 40%
- Matter + anti-matter: 100%

## 2. Inexhaustible supply of fusion fuels:

Deuterium can be distilled from all forms of water, while tritium will be produced during the fusion reaction as fusion neutrons interact with lithium. The reserve on Earth is able to fulfil the needs for **millions of years**.

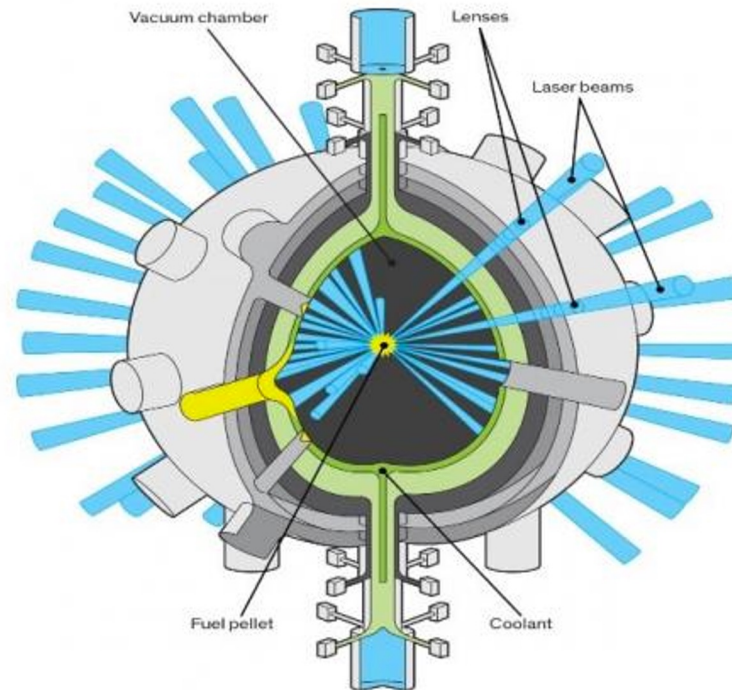
## 3. Environment friendly:

- No CO<sub>2</sub>
- No long-lived radioactive waste
- No risk of meltdown

# Preliminaries: Two major ways of controlled nuclear fusion

The temperature required for confining the fusion plasma are so hot (>10 million °C), and cannot be confined via any material. Two main ways of confinement:

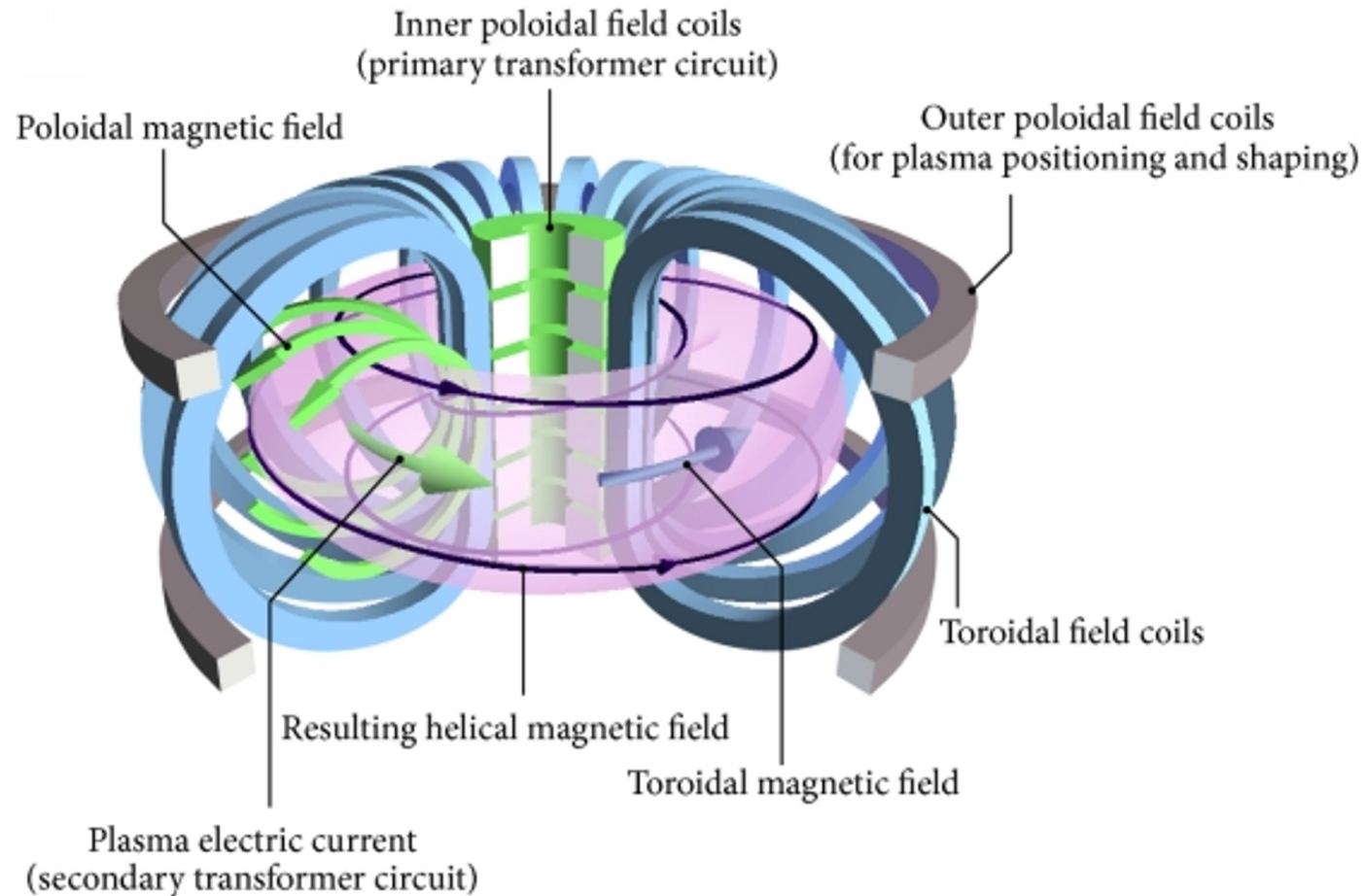
## 1. Inertial confinement (惯性约束) :



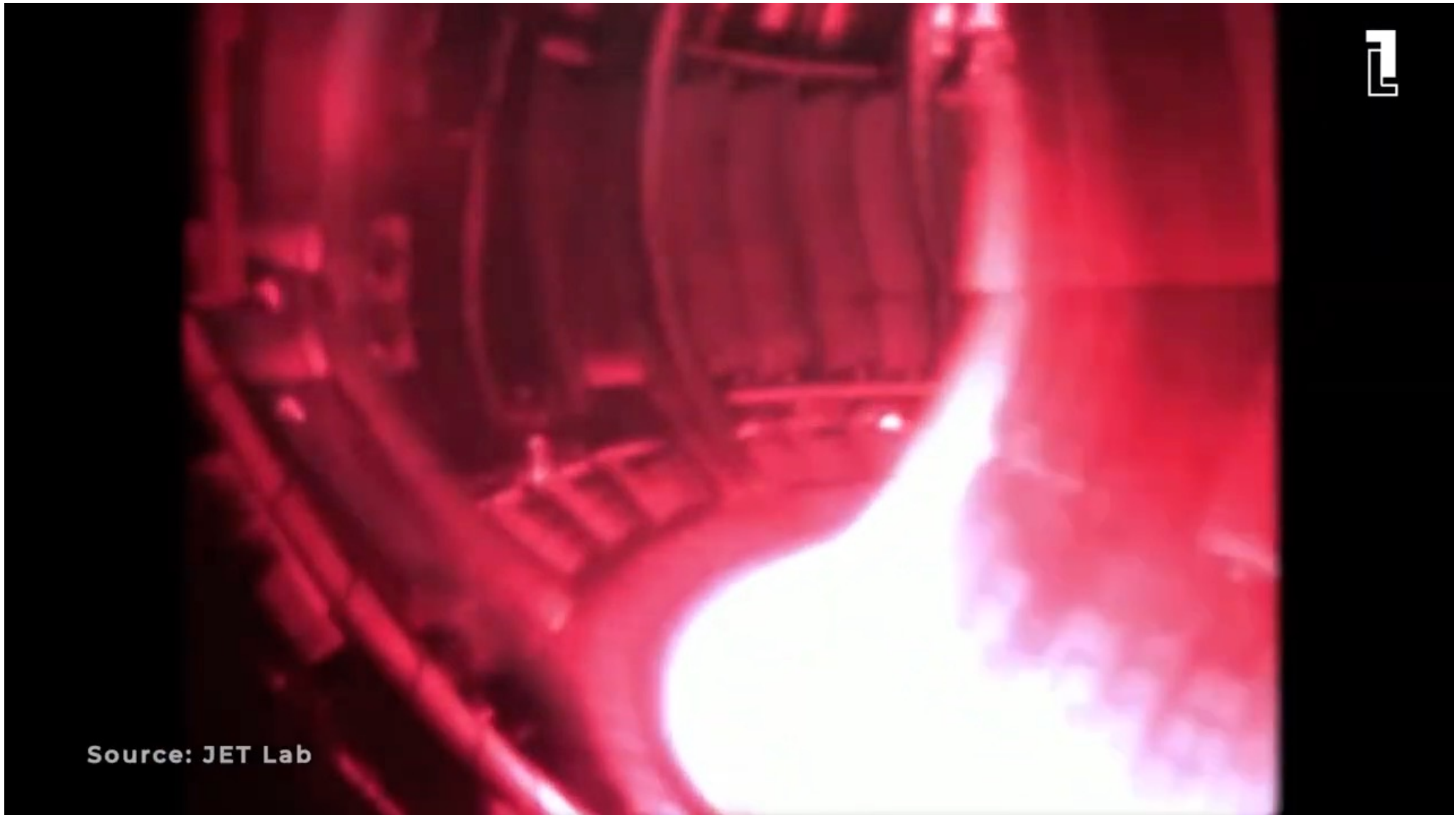
# Preliminaries: Two major ways of controlled nuclear fusion

## 2. Magnetic confinement (磁约束), using Tokamak (current work)

---



# Controlled nuclear fusion using Tokamak



Source: JET Lab

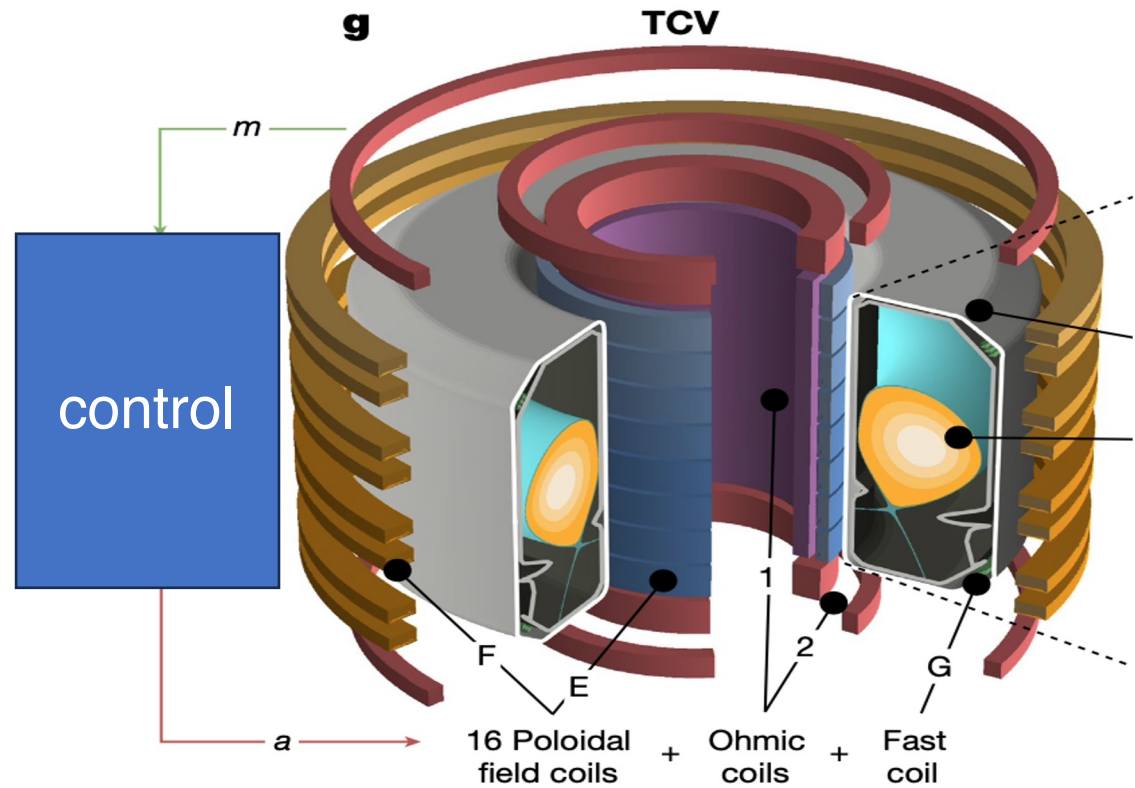
# 1. What is the problem?

**Task:** To shape and maintain a high-temperature plasma within the tokamak vessel.

Each time step  $t$  have observation and needs to output a control signal:

观测  $m$ : input observation,  $\mathbb{R}^{92}$

控制  $a$ : output control,  $\mathbb{R}^{19}$





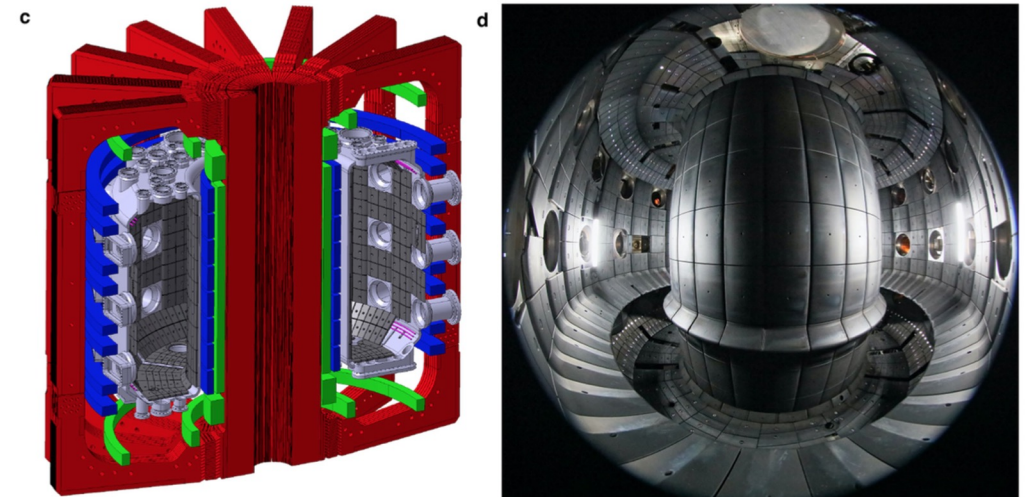
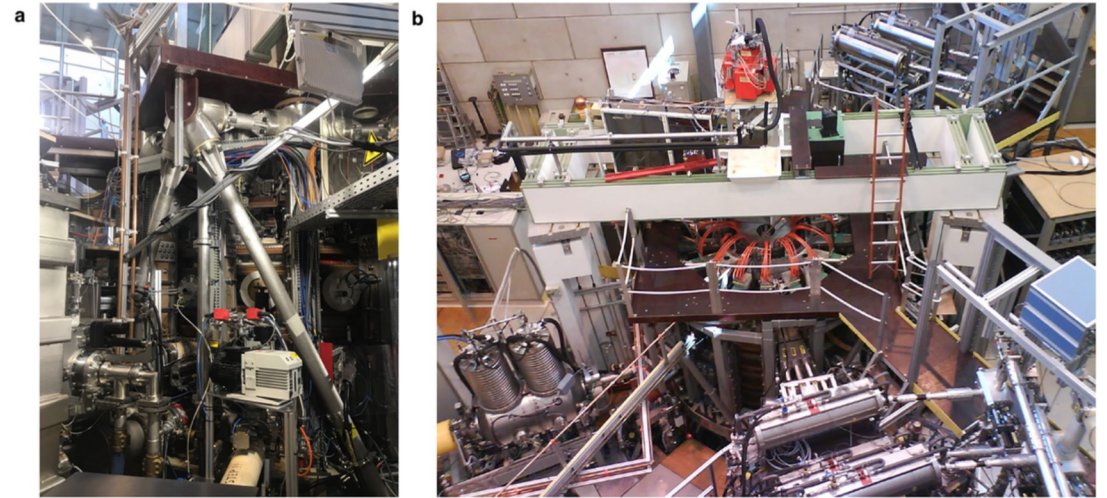
# 1. What is the problem? Device overview

Tokamak à configuration variable (TCV)  
(current work)

- Plasma height: 1.40m
- Major radius: 0.88m
- Plasma life span: 2s maximum
- Toroidal magnetic field: 1.43T
- Additional heating power: 4.5MW

ITER (cost \$22 billion, test first plasma  
in 2025 and full fusion in 2035):

- Major radius: 6.2 m
- Magnetic field: 11.8 T
- Heating power: 320 MW
- Fusion power: 500 MW
- Discharge duration: up to 1000 s



Extended Data Fig. 1 | Pictures and illustration of the TCV. a, b Photographs showing the part of the TCV inside the bioshield. c CAD drawing of the vessel and coils of the TCV. d View inside the TCV (Alain Herzog/EPFL), showing the limiter tiling, baffles and central column.



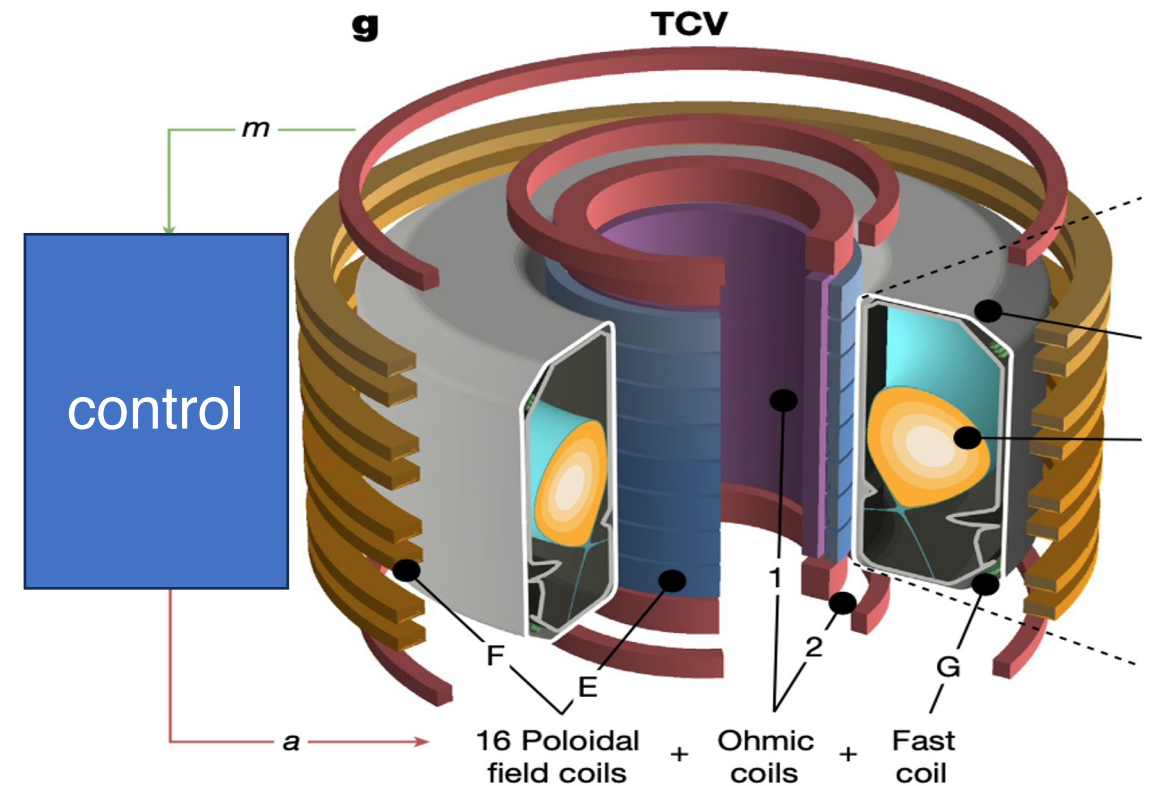
## 2. Why is it important?

The effective control of plasma within a tokamak will **pave the way** for commercial nuclear fusion, which allows to produce energy energy that is

- (1) Virtually unlimited;
- (2) Environmentally friendly.

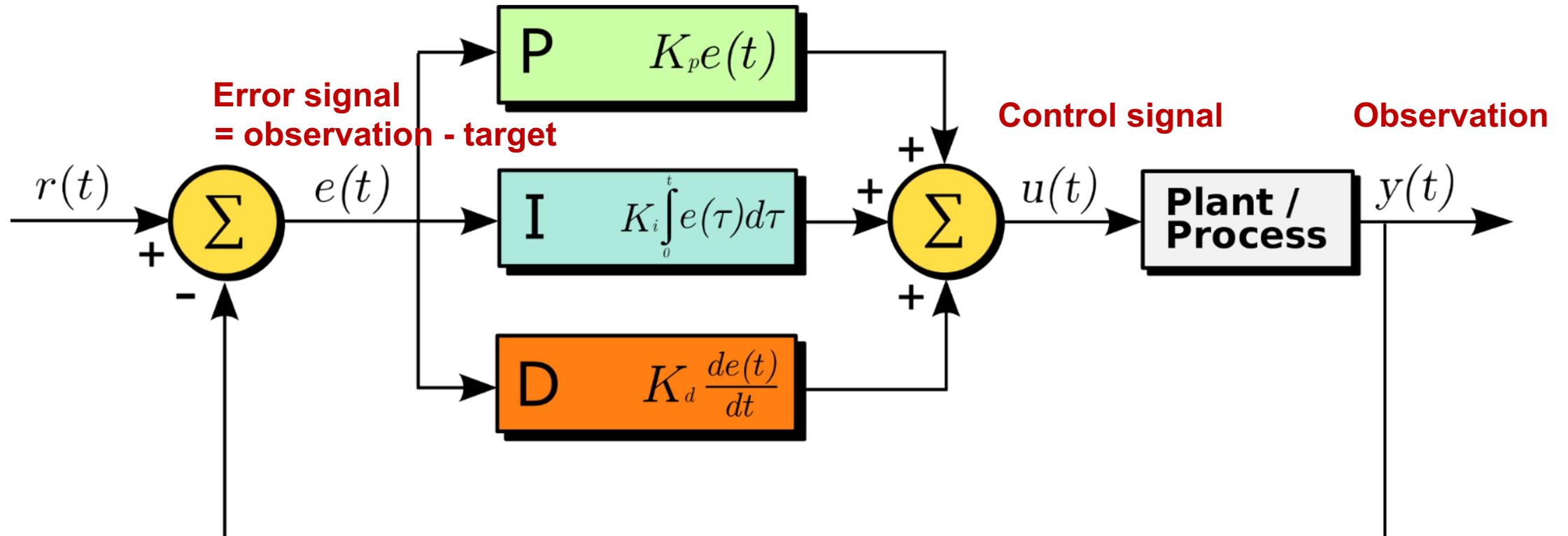
### 3. Why is it hard?

This requires **high-dimensional, high-frequency, closed-loop** control using magnetic actuator coils, further complicated by the diverse requirements across a wide range of plasma configurations.



## 4. Limitation of prior methods: PID control

Proportional–integral–derivative (PID) control:



## 4. Limitation of prior methods: PID control

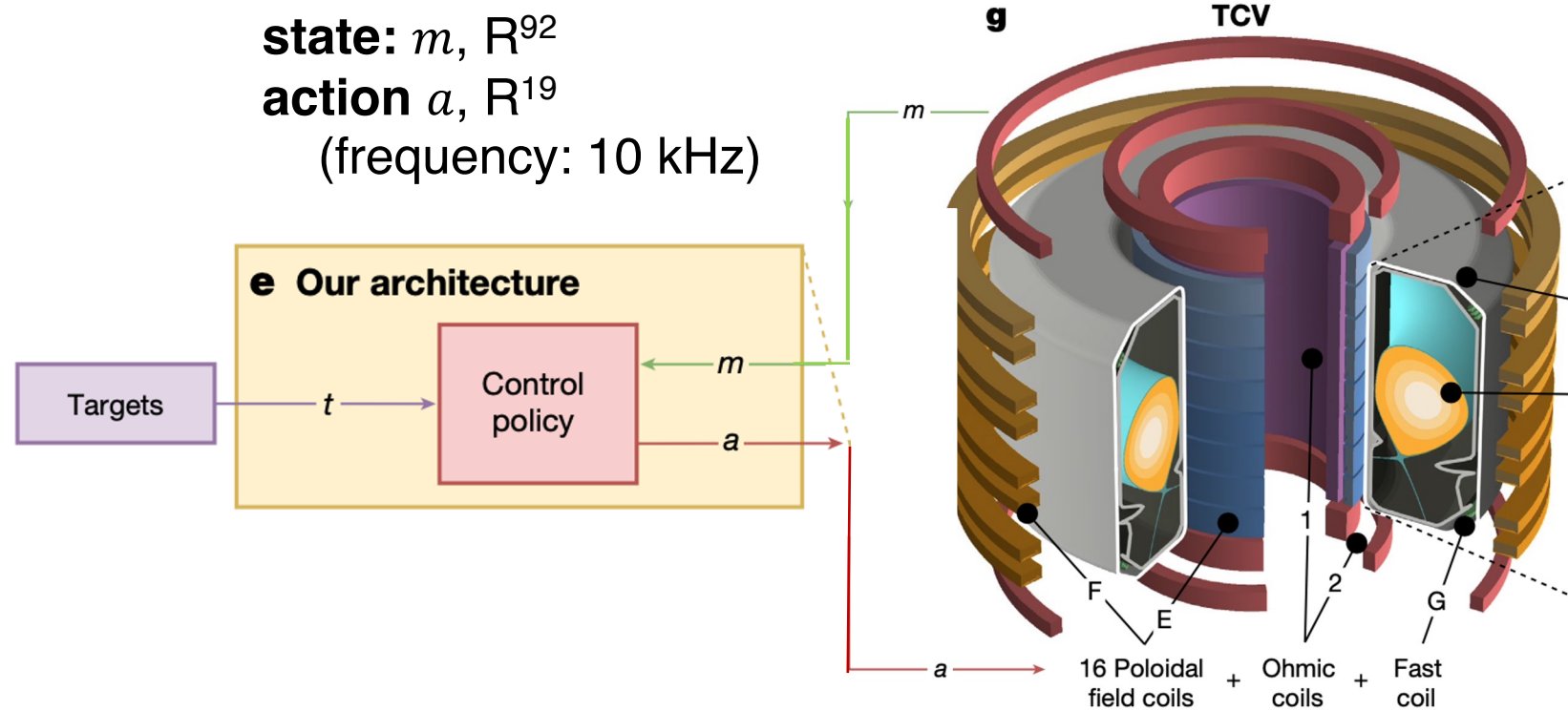
**Pros:** Effective

**Cons:**

- (1) The controllers are designed on the basis of linearized model dynamics
- (2) Requires substantial engineering effort, design effort and expertise whenever the target plasma configuration is changed, together with complex, real-time calculations for equilibrium estimation

# 5. Main components of the proposed method [1]

[1] Degraeve, Jonas, et al. "Magnetic control of tokamak plasmas through deep reinforcement learning." *Nature* 602.7897 (2022): 414-419.



# 5. Main components of the proposed method

**Reward:** The target values  $g$  of the objectives are often time-varying (e.g., the plasma current and boundary target points), and are sent to the policy as part of the observations:  $\pi(a|s, g)$ .

Reward Component	Description
Diverted	Whether the plasma is limited by the wall or diverted through an X-point.
E/F Currents	The currents in the E and F coils, in amperes.
Elongation	The elongation of the plasma, this is its height divided by its width.
LCFS Distance	The distance in meters from the target points to the nearest point on the last closed flux surface (LCFS).
Legs Normalized Flux	The difference in normalized flux from the flux at the LCFS at target leg points.
Limit Point	The distance in meters from the actual limit point (wall or X-point) and target limit point.
OH Current Diff	The difference in amperes between the two OH coils.
Plasma Current	The plasma current in amperes.
R	The radial position of the plasma axis/centre, in meters.
Radius	Half of the width of the plasma, in meters.
Triangularity	The upper triangularity is defined as the radial position of the highest point relative to the median radial position. The overall triangularity is the mean of the upper and lower triangularity.
Voltage Out of Bounds	Penalty for going outside of the voltage limits.
X-point Count	Return the number of actual and requested X-points within the vessel.
X-point Distance	Returns the distance in meters from actual X-points to target X-points. Only X-points within 20cm are considered.
X-point Far	For any X-point that isn't requested, return the distance in meters from the X-point to the LCFS. This helps avoid extra X-points that may attract the plasma and lead to instabilities.
X-point Flux Gradient	The gradient of the flux at the target location with a target of 0 gradient. This encourages an X-point to form at the target location, but isn't very precise on the exact location.
X-point Normalized Flux	The difference in normalized flux from the flux at the LCFS at target X-points. This encourages the X-point to be on the last closed flux surface, and therefore for the plasma to be diverted.
Z	The vertical position of the plasma axis/centre, in meters.

# 5. Main component of the proposed method

## Training:

Perform training within a simulated environment using a **solver**.

## Inference:

Directly deploy it in the device.

## RL method:

Maximum a posteriori policy optimization (MPO) [1].

---

### Algorithm 1 Actor-Critic

---

```
Initialize  $\pi^{(0)}, Q^{\pi^{(-1)}}, k \leftarrow 0$   
repeat  
   $Q^{\pi^{(k)}} \leftarrow \text{PolicyEvaluation}(\pi^{(k)}, Q^{\pi^{(k-1)}})$   
   $\pi^{(k+1)} \leftarrow \text{PolicyImprovement}(\pi^{(k)}, Q^{\pi^{(k)}})$   
   $k \leftarrow k + 1$   
until convergence
```

---

Actor  $\pi$ : small MLP, must be fast.  
Critic  $Q^\pi$ : LSTM, can be large, only used in training

[1] Abdolmaleki, Abbas, et al. "Maximum a posteriori policy optimisation." ICLR 2018

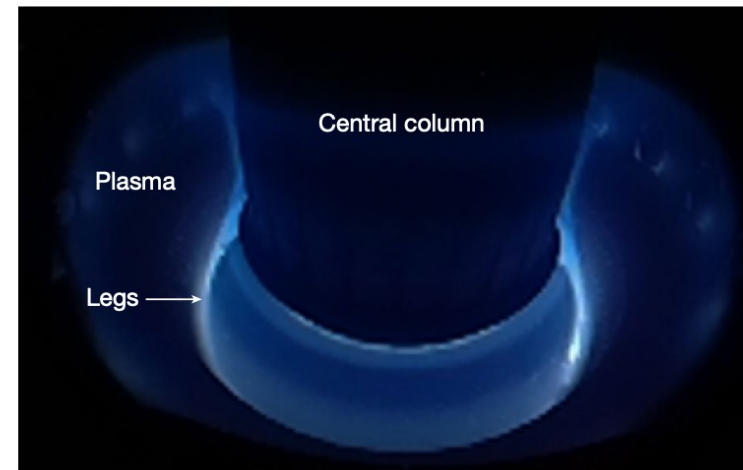
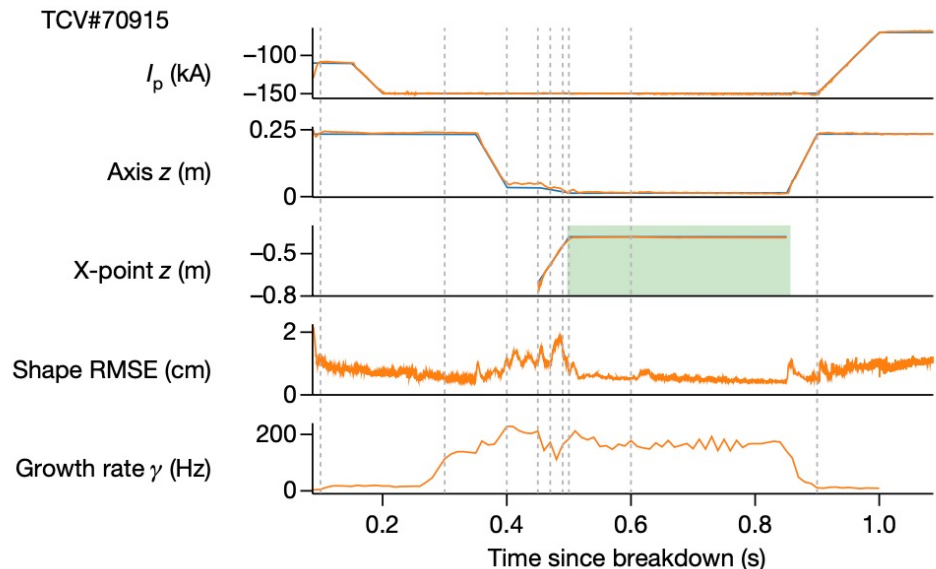
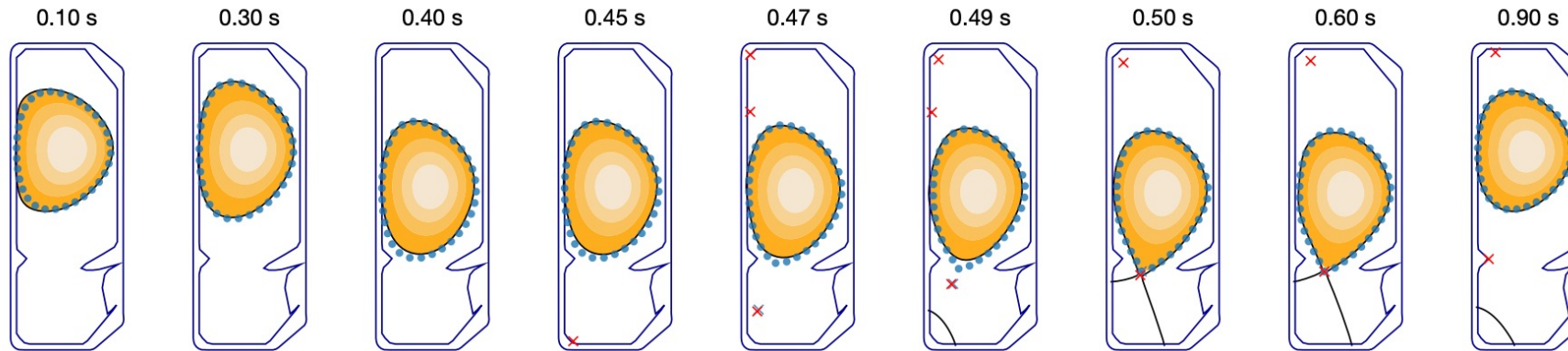
# Inference code

```
22 def run_loop(env: environment.Environment, agent,
23             max_steps: int = 100000) -> trajectory.Trajectory:
24     """Run an agent."""
25     results = []
26     agent.reset()
27     ts = env.reset()
28     for _ in range(max_steps):
29         obs = ts.observation
30         action = agent.step(ts)
31         ts = env.step(action)
32         results.append(trajectory.Trajectory(
33             measurements=obs["measurements"],
34             references=obs["references"],
35             actions=action,
36             reward=np.array(ts.reward)))
37     if ts.last():
38         break
39
40     return trajectory.Trajectory.stack(results)
```



# 6. Main results

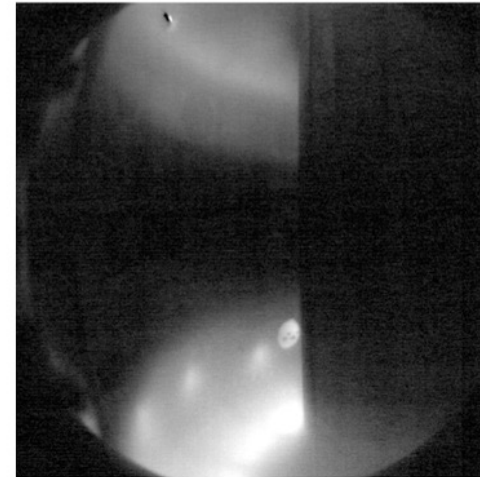
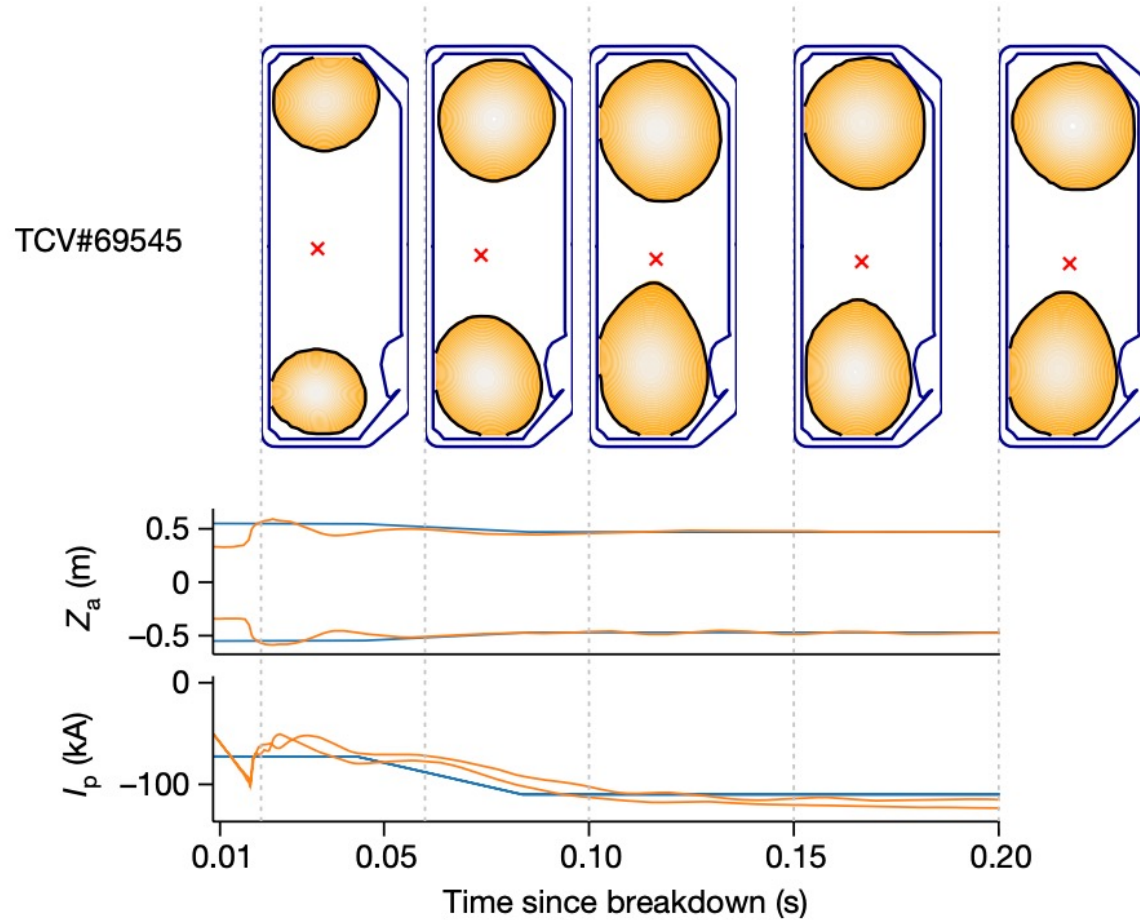
The position and shape (orange line) matches well with the target (blue)



Inside view at 0.6 s

# 6. Main results

First demonstration of double droplet shape:

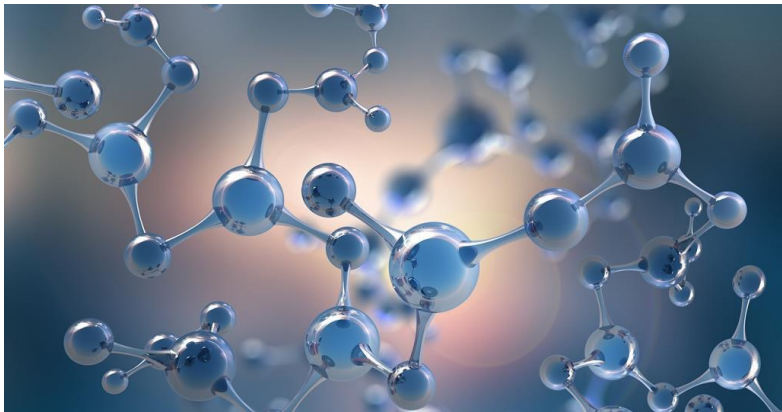


# RL for Science: Case study

## 1. Deep RL for controlled nuclear fusion



## 2. Deep RL for molecule design



# Task and significance

**Task:** Design molecules that optimize certain properties such as drug-likeness and synthetic accessibility, while obeying physical laws such as chemical valency.

**Significance:** Molecule design is important in drug discovery.

# Why is it hard?

## 1. Large size of chemical space:

The range of drug-like molecules has been estimated to be between  $10^{23}$  and  $10^{60}$  [1].

## 2. Chemical space is discrete, and molecular properties are highly sensitive to small changes in the molecular structure

[1] P. G. Polishchuk, T. I. Madzhidov, and A. Varnek. Estimation of the size of drug-like chemical space based on gdb-17 data. *Journal of Computer-Aided Molecular Design*, 27(8):675–679, Aug 2013.

# What is the limitation of prior methods?

There are multiple prior works that uses recurrent neural networks [1][2], autoencoder [3], GANs [4], they are limited in

1. Generating **novel** and **valid** molecular graphs that can directly **optimize** various desired physical, chemical and biological property objectives.
2. Actively explore the vast chemical space.

[1] E. Jannik Bjerrum and R. Threlfall. Molecular Generation with Recurrent Neural Networks (RNNs). *arXiv preprint arXiv:1705.04612*, 2017.

[2] M. H. S. Segler, et al. Generating focused molecule libraries for drug discovery with recurrent neural networks. *ACS Central Science*, 4(1):120–131, 2018.

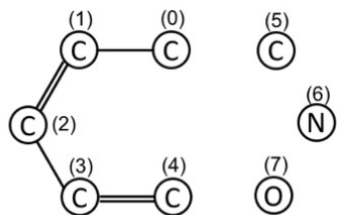
[3] R. Gómez-Bombarelli, et al. Automatic chemical design using a data-driven continuous representation of molecules. *ACS Central Science*, 2016.

[4] B. Sanchez-Lengeling, C. Outeiral, G. L. Guimaraes, and A. Aspuru-Guzik. Optimizing distributions over molecular space. An Objective-Reinforced Generative Adversarial Network for Inverse-design Chemistry (ORGANIC). *ChemRxiv e-prints*, 8 2017.

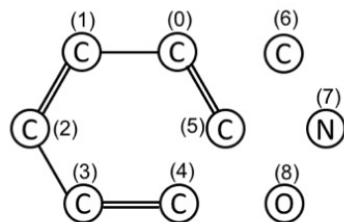
# Recall: Foundational principles in deep learning

- Principle 1** (deep principle): Model a hard transformation by composing many simple, easy transformations.
- Principle 2** (end-to-end law): Directly optimizing the final objective using maximum likelihood and information theory.
- Principle 3** (the scaling law): AI methods that leverage computation are ultimately the most effective way of improvements (from "The bitter lesson" by Rich Sutton).
- Principle 4** (the data law): Data is the ultimate way of regularization.
- Principle 5** (consistency law): The more consistent between training and testing, the better the performance.

starting configuration

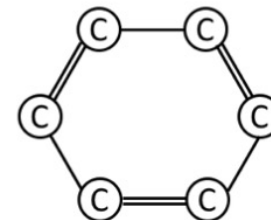


reinforcement learning

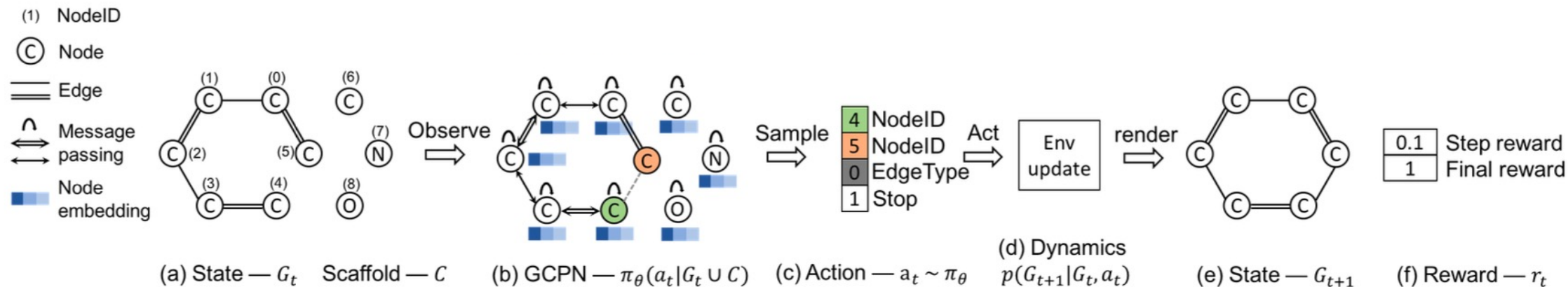


...

target configuration



# What are the main component of the proposed method?



**state**  $s$ : current graph

**action**  $a$ : (NodeID1, NodeID2, EdgeType, Stop) (dimension does not change!)

**reward**  $r$ : domain-specific rewards + adversarial rewards (using GAN)

The adversarial reward encourages the generated molecules resembles given molecules.



# What are the main component of the proposed method?

**Policy**  $\pi_{\theta}(a|s)$ : Graph Neural Networks (GNNs, to be introduced in class 10)

**RL method:** Proximal Policy Optimization (PPO) [1]

$$\max L^{\text{CLIP}}(\theta) = \mathbb{E}_t[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t)], r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$$

[1] Schulman, John, et al. "Proximal policy optimization algorithms." *arXiv preprint arXiv:1707.06347* (2017).

# Main results: Property optimization

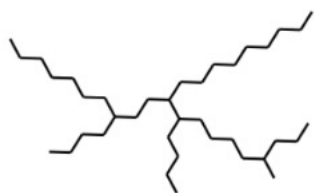
Comparison of the top 3 property scores of generated molecules found by each model:

Method	Penalized logP				QED			
	1st	2nd	3rd	Validity	1st	2nd	3rd	Validity
ZINC	4.52	4.30	4.23	100.0%	0.948	0.948	0.948	100.0%
Hill Climbing	—	—	—	—	0.838	0.814	0.814	100.0%
ORGAN	3.63	3.49	3.44	0.4%	0.896	0.824	0.820	2.2%
JT-VAE	5.30	4.93	4.49	100.0%	0.925	0.911	0.910	100.0%
GCPN	<b>7.98</b>	<b>7.85</b>	<b>7.80</b>	<b>100.0%</b>	<b>0.948</b>	<b>0.947</b>	<b>0.946</b>	<b>100.0%</b>

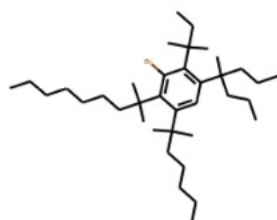
logP: Octanol-water partition coefficient  
QED: druglikeness

# Main results: Property optimization

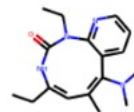
## Generated molecules



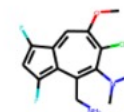
7.98



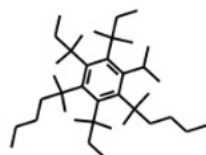
7.48



0.948



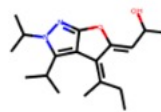
0.945



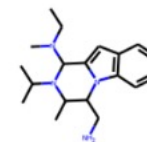
7.12



23.88\*



0.944



0.941

(a) Penalized logP optimization

(b) QED optimization

logP: Octanol-water partition coefficient  
QED: druglikeness

# Main results: Property targeting

Comparison of the effectiveness of property targeting task:

Method	$-2.5 \leq \log P \leq -2$		$5 \leq \log P \leq 5.5$		$150 \leq MW \leq 200$		$500 \leq MW \leq 550$	
	Success	Diversity	Success	Diversity	Success	Diversity	Success	Diversity
ZINC	0.3%	0.919	1.3%	0.909	1.7%	0.938	0	—
JT-VAE	11.3%	<b>0.846</b>	7.6%	0.907	0.7%	0.824	16.0%	0.898
ORGAN	0	—	0.2%	<b>0.909</b>	15.1%	0.759	0.1%	0.907
GCPN	<b>85.5%</b>	0.392	<b>54.7%</b>	0.855	<b>76.1%</b>	<b>0.921</b>	<b>74.1%</b>	<b>0.920</b>

logP: Octanol-water partition coefficient  
MW: molecular weight

# Other Deep RL work in AI4Science: Life science (1)

## Protein:

1. Wang, Yi, et al. "Self-play reinforcement learning guides protein engineering." *Nature Machine Intelligence* 5.8 (2023): 845-860.
2. Lutz, Isaac D., et al. "Top-down design of protein architectures with reinforcement learning." *Science* 380.6642 (2023): 266-273.
3. Lee, Minji, et al. "Protein sequence design in a latent space via model-based reinforcement learning." (2022).
4. Xu, Xiaopeng, et al. "AB-Gen: antibody library design with generative pre-trained transformer and deep reinforcement learning." *Genomics, Proteomics & Bioinformatics* (2023).

## Molecules:

1. Jeon, Woosung, and Dongsup Kim. "Autonomous molecule generation using reinforcement learning and docking to develop potential novel inhibitors." *Scientific reports* 10.1 (2020): 22104.
2. Korshunova, Maria, et al. "Generative and reinforcement learning approaches for the automated de novo design of bioactive compounds." *Communications Chemistry* 5.1 (2022): 129.
3. Mazuz, Eyal, et al. "Molecule generation using transformers and policy gradient reinforcement learning." *Scientific Reports* 13.1 (2023): 8799.
4. Polykovskiy, Daniil, et al. "Molecular sets (MOSES): a benchmarking platform for molecular generation models." *Frontiers in pharmacology* 11 (2020): 565644.

# Other Deep RL work in AI4Science: Life science (2)

## **Molecules (continued):**

5. Hu, Xiuyuan, et al. "De novo Drug Design using Reinforcement Learning with Multiple GPT Agents." *Advances in Neural Information Processing Systems* 36 (2024).
6. Popova, Mariya, Olexandr Isayev, and Alexander Tropsha. "Deep reinforcement learning for de novo drug design." *Science advances* 4.7 (2018): eaap7885.

## **RNA:**

1. Whatley, Alexander, Zhekun Luo, and Xiangru Tang. "Improving RNA secondary structure design using deep reinforcement learning." *arXiv preprint arXiv:2111.04504* (2021).
2. Eastman, Peter, et al. "Solving the RNA design problem with reinforcement learning." *PLoS computational biology* 14.6 (2018): e1006176.

## **Genomics:**

1. Nicholls, Hannah L., et al. "Reaching the end-game for GWAS: machine learning approaches for the prioritization of complex disease loci." *Frontiers in genetics* 11 (2020): 521712.
2. Karami, Mohsen, et al. "Revolutionizing genomics with reinforcement learning techniques." *arXiv preprint arXiv:2302.13268* (2023).

# Other Deep RL work in AI4Science: Fluid control (1)

## Cylinder:

1. Chen, Wenjie, et al. "Deep reinforcement learning-based active flow control of vortex-induced vibration of a square cylinder." *Physics of Fluids* 35.5 (2023). (SAC)
2. Wang, Qiulei, et al. "DRLinFluids: An open-source Python platform of coupling deep reinforcement learning and OpenFOAM." *Physics of Fluids* 34.8 (2022). (SAC)
3. Wang, Qiulei, et al. "Dynamic feature-based deep reinforcement learning for flow control of circular cylinder with sparse surface pressure sensing." *arXiv preprint arXiv:2307.01995* (2023). (SAC & PPO)
4. Tang, Hongwei, et al. "Robust active flow control over a range of Reynolds numbers using an artificial neural network trained through deep reinforcement learning." *Physics of Fluids* 32.5 (2020). (PPO)
5. Xu, Hui, et al. "Active flow control with rotating cylinders by an artificial neural network trained by deep reinforcement learning." *Journal of Hydrodynamics* 32.2 (2020): 254-258. (PPO)
6. Rabault, Jean, et al. "Artificial neural networks trained through deep reinforcement learning discover control strategies for active flow control." *Journal of fluid mechanics* 865 (2019): 281-302. (PPO)
7. Wang, Zhicheng, et al. "Deep reinforcement transfer learning of active control for bluff body flows at high Reynolds number." *Journal of Fluid Mechanics* 973 (2023): A32. (TD3)
8. Zheng, Changdong, et al. "Data-efficient deep reinforcement learning with expert demonstration for active flow control." *Physics of Fluids* 34.11 (2022). (SAC)

# Other Deep RL work in AI4Science: Fluid control (2)

## Point:

1. Mei, Jiazhong, J. Nathan Kutz, and Steven L. Brunton. "Observability-Based Energy Efficient Path Planning with Background Flow via Deep Reinforcement Learning." *2023 62nd IEEE Conference on Decision and Control (CDC)*. IEEE, 2023. (PPO)
2. Gunnarson, Peter, et al. "Learning efficient navigation in vortical flow fields." *Nature communications* 12.1 (2021): 7143. (V-RACER)

## Foil:

1. Novati, Guido, and Petros Koumoutsakos. "Remember and forget for experience replay." *International Conference on Machine Learning*. PMLR, 2019. (V-RACER)
2. Wang ZP, Lin RJ, Zhao ZY, et al. Learn to flap: foil non-parametric path planning via deep reinforcement learning. *Journal of Fluid Mechanics*. 2024;984:A9. (PPO)

## Fish:

1. Verma, Siddhartha, Guido Novati, and Petros Koumoutsakos. "Efficient collective swimming by harnessing vortices through deep reinforcement learning." *Proceedings of the National Academy of Sciences* 115.23 (2018): 5849-5854. (DRQN)
2. Mandralis, Ioannis, et al. "Learning swimming escape patterns for larval fish under energy constraints." *Physical Review Fluids* 6.9 (2021): 093101. (V-RACER)



# Other Deep RL work in AI4Science: Materials science (1)

## Materials:

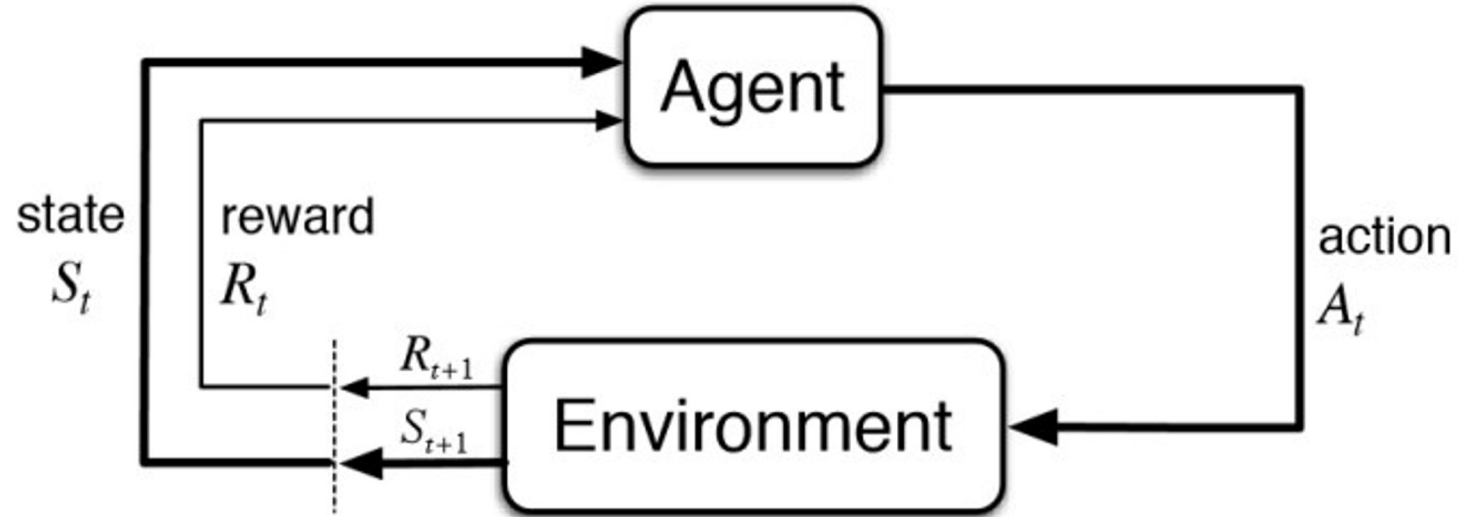
1. Rajak, Pankaj, et al. "Autonomous reinforcement learning agent for chemical vapor deposition synthesis of quantum materials." *npj Computational Materials* 7.1 (2021): 108.
2. Zamaraeva, Elena, et al. "Reinforcement learning in crystal structure prediction." *Digital Discovery* 2.6 (2023): 1831-1840.
3. Zheng, Bowen, Zeyu Zheng, and Grace X. Gu. "Designing mechanically tough graphene oxide materials using deep reinforcement learning." *npj Computational Materials* 8.1 (2022): 225.
4. Govindarajan, Prashant, et al. "Learning Conditional Policies for Crystal Design Using Offline Reinforcement Learning." *Digital Discovery* (2024).
5. Pandey, Ashish, et al. "Reinforcement learning based carbon nanotube growth automation." *2021 IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*. IEEE, 2021.
6. Pan, Elton, Christopher Karpovich, and Elsa Olivetti. "Deep reinforcement learning for inverse inorganic materials design." *arXiv preprint arXiv:2210.11931* (2022).

# Other Deep RL work in AI4Science: Materials science (2)

## **Meta-materials/composite/polymer:**

1. Sui, Fanping, et al. "Deep reinforcement learning for digital materials design." *ACS Materials Letters* 3.10 (2021): 1433-1439.
2. Gongora, Aldair E., et al. "Designing composites with target effective young's modulus using reinforcement learning." *Proceedings of the 6th Annual ACM Symposium on Computational Fabrication*. 2021.
3. Ma, Ruimin, Hanfeng Zhang, and Tengfei Luo. "Exploring high thermal conductivity amorphous polymers using reinforcement learning." *ACS Applied Materials & Interfaces* 14.13 (2022): 15587-15598.
4. Rosafalco, Luca, et al. "Reinforcement learning optimisation for graded metamaterial design using a physical-based constraint on the state representation and action space." *Scientific Reports* 13.1 (2023): 21836.

# Markov Decision Process (MDP): Setup



**Goal:** Maximize the long-term expected reward w.r.t. to the policy  $\pi(A_t|S_t)$

$$\max_{\pi(A_t|S_t)} \mathbb{E}_t[R_t]$$

# Application in AI for Science: from microscopic to macroscopic

