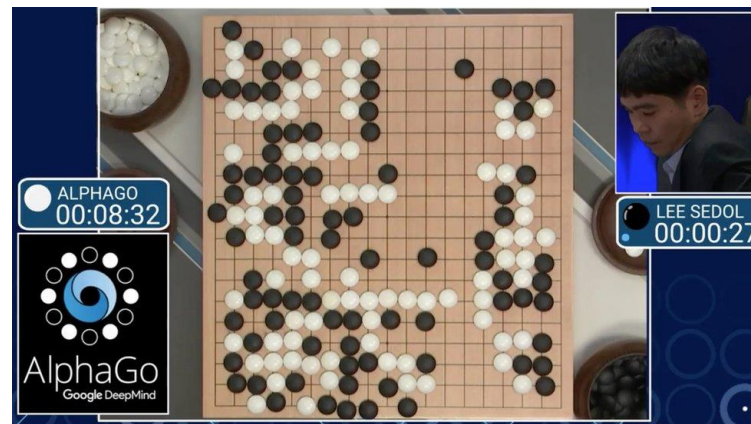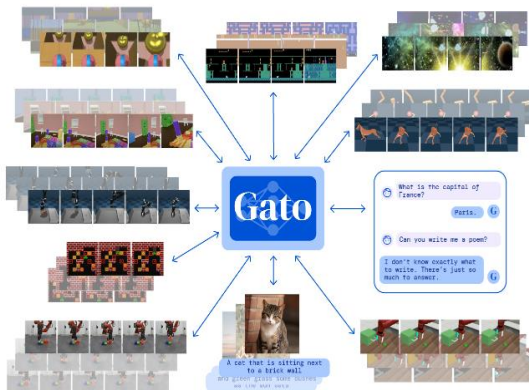# Review : Magnetic control of tokamak plasmas through deep reinforcement learning

**J.S. Kim**

*Department of Nuclear Engineering, Seoul National University, Korea*

# Contents

- Introduction

- Grad-Shafranov equation

- MPO : Maximum a posteriori policy optimisation

- Experiment

- Discussion

- Appendix

# Introduction

- ## Setup
  - Control value : TCV coils current(19 coils : Poloidal coils 16 + Ohmic coils 2 + Fast Coil)
  - What to do : <span style="color:red">Plasma shape control(RL model) + control stability(from reward)</span>
  - Algorithm : Maximize a posteriori policy optimization(MPO, Actor + Critic)
  - Reward : Elongation, LCFS Distance, OH current drift, Plasma current, Radius, Triangularity, x-point
  - Metric : Shape loss + scalar quantities(i.e. plasma current, q, beta, error between reference and the respective estimation from the equilibrium reconstruction)

- Setup
  - Actor-Critic : temporal difference version of policy gradient
  - Component
    - Actor network: the actor decides which action should be taken given the state s
    - Critic network : the critic informs the actor how good was the action and how it should adjust

Input Data : TCV coils current(19)

Output : Action distribution for each coil with Gaussian Distribution

| | parameter | value | lower bound | upper bound |
|---|---|---|---|---|
| action delay | E | 0.5 ms | | |
| | F | 0.5 ms | | |
| | OH | 0.5 ms | | |
| | G | 0.1 ms | | |
| action bias (fixed) | E001 | 7 V | | |
| | E002 | −10 V | | |
| | E003 | −1 V | | |
| | E004 | 0 V | | |
| | E005 | 11 V | | |
| | E006 | −1 V | | |
| | E007 | −4 V | | |
| | E008 | 44 V | | |
| | F001 | 38 V | | |
| | F002 | −3 V | | |
| | F003 | 6 V | | |
| | F004 | 1 V | | |
| | F005 | −37 V | | |
| | F006 | −9 V | | |
| | F007 | 5 V | | |
| | F008 | 10 V | | |
| | OH001 | −54 V | | |
| | OH002 | −15 V | | |
| action offset (random) | all coils | | −20 V | |



E3 Coil on TCV#70755



Droplets    Negative Triangularity    ITER-like shape    Snowflake    Elongated Plasma

32 equiangular points around the LCFS and canonicalize with splines to 128 equidistant points

- Conventional approach
  - **Feedforward control** : a system which passes the signal to some external load
  - Process
    - Solve an inverse problem to precompute a set of feedforward coil currents and voltages
    - PID controllers(linearized model dynamic) : stabilize the plasma vertical position and control the radial position and current
    - PID controllers do not mutually interfere

- New approach
  - **Feedback control** :  a system where the output depends on the generated feedback signal
  - **Reinforcement Learning** : generate non-linear feedback controllers + policy-based method
    - Single computationally inexpensive controller(Neural Network) replaces the nested control architecture
    - Reduction of the controller development cycle
    - Environment : GS solver(simulator) or surrogate model using classical ML

## Feed-forward



- Model-based prediction of input
- Ideally consists of exact inverse model of the process
- Effects of disturbance or command input must be predictable
- May not generalize to other conditions
- Will not be accurate if the system changes
- Feedforward component provides rapid response

## Feed-back



- Automatically compensates for disturbances (=controller acts on error)
- Automatically follows change in desired state
- Reactive / Error - driven
- Can be very simple
- Can improve undesirable properties of system
- Feedback component provides response accurately via compensating for errors in the model

- **RL based controller design architecture**

## What is Deep Reinforcement Learning?

# Introduction : Reinforcement Learning

- Fundamental of Reinforcement Learning

  - Definition : ML technique that enables an agent to learn in an interactive environment by trial and error using feedback from its own actions and experience → **Learn the optimal policy which maximizes reward**

  - Element
    - Agent : Object that takes decisions based on the rewards and punishment
    - Environment : Physical world in which the agent interacts ( GS solver)
    - State : Current situation of the agent
    - Reward : Feedback from the environment
    - Action : Mechanism by which the agent transitions between states of the environment
    - Value : Future reward that an agent would receive by taking an action in a particular state

▪ Markov Process

**"The future is independent of the past given the present"**

**Definition 3.1.** A discrete time stochastic control process is Markovian (i.e., it has the Markov property) if

- $\mathbb{P}(\omega_{t+1} \mid \omega_t, a_t) = \mathbb{P}(\omega_{t+1} \mid \omega_t, a_t, \ldots, \omega_0, a_0)$, and

- $\mathbb{P}(r_t \mid \omega_t, a_t) = \mathbb{P}(r_t \mid \omega_t, a_t, \ldots, \omega_0, a_0)$.

**Definition 3.2.** An MDP is a 5-tuple $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$ where:

- $\mathcal{S}$ is the state space,

- $\mathcal{A}$ is the action space,

- $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to [0, 1]$ is the transition function (set of conditional transition probabilities between states),

- $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathcal{R}$ is the reward function, where $\mathcal{R}$ is a continuous set of possible rewards in a range $R_{\max} \in \mathbb{R}^+$ (e.g., $[0, R_{\max}]$),

- $\gamma \in [0, 1)$ is the discount factor.

$$P[S_{t+1} \mid S_t] = P[S_{t+1} \mid S_1, \ldots, S_t]$$

- Policy gradient methods

  - Policy : **Probability** of **choosing action a** given the **state s**
  - Objective : to **find the optimal policy** which achieves **optimal reward**
  - Reward function

  $$J(\theta) = \sum_{s \in S} d^\pi(s) V^\pi(s) = \sum_{s \in S} d^\pi(s) \sum_{a \in A} \pi_\theta(a|s) Q^\pi(s,a)$$

  - How to optimize

  $$\nabla_\theta J(\theta) = \nabla_\theta \sum_{s \in S} d^\pi(s) \sum_{a \in A} Q^\pi(s,a) \pi_\theta(a|s)$$
  $$\propto \sum_{s \in S} d^\pi(s) \sum_{a \in A} Q^\pi(s,a) \nabla_\theta \pi_\theta(a|s)$$

  $$\nabla_\theta J(\theta) = \mathbb{E}_\pi [Q^\pi(s,a) \nabla_\theta \ln \pi_\theta(a|s)]$$

$$V(S_t) \leftarrow V(S_t) + \alpha \left( R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right)$$

# Grad-Shafranov equation

- MHD equilibrium in Tokamak



$$\vec{J} \times \vec{B} = \nabla P$$

$$\nabla \times \vec{B} = \mu_o \vec{J}$$

$$\nabla \cdot \vec{B} = 0$$

- Grad-Shafranov equation

$$F(\psi) = R B_\phi$$

$$J_\phi(R, Z) = R \frac{dP}{d\psi} + \frac{F(\psi)}{\mu_o R} \frac{dF(\psi)}{d\psi}$$

$$\Delta^* \psi = R \frac{\partial}{\partial R} \left( \frac{1}{R} \frac{\partial \psi}{\partial R} \right) + \frac{\partial^2 \psi}{\partial R^2}$$

$$\Delta^* \psi = -\mu_o R^2 \frac{dP}{d\psi} - F \frac{dF}{d\psi}$$

# Grad-Shafranov equation

- General vacuum solution
    - Condition : zero pressure + zero current outside of the plasma(p=F=0)

$$\Delta^* \psi = 0 \longrightarrow \psi_{vac} = \sum_n c_n R \cos(k_n Z)(J_1(k_n R) + d_n Y_1(k_n R))$$

- Vacuum solution of a circular current loop
    - Green's function of the Grad-shafranov operator is given by

$$G(R, Z; R_0, Z_0) = -\frac{1}{2\pi}\sqrt{RR_0}\frac{1}{k}\left[(2 - k^2)K(k) - 2E(k)\right]$$

Elliptic integrals of the 1st and 2nd kind

    - The magnetic flux due to the circular current loop is given by

$$\frac{\psi(R,\theta)}{R} = SR\frac{4(2 - k(R,\theta)^2)K(k(R,\theta)^2) - 2E(k(R,\theta)^2)}{\sqrt{a^2 + R^2 + 2aR\sin(\theta)}} \qquad k(R,\theta) = \frac{4aR\sin(\theta)}{a^2 + R^2 + 2aR\sin(\theta)}$$

X-point

Magnetic axis(=null point)

$$S(R,Z) \equiv \left(\frac{\partial^2 \psi}{\partial R^2}\right)\left(\frac{\partial^2 \psi}{\partial Z^2}\right) - \left(\frac{\partial^2 \psi}{\partial R \partial Z}\right)^2$$

$$|\nabla \psi|^2 = 0$$

$S(R,Z) > 0 :$ magnetic axis

$S(R,Z) < 0 :$ x-point

# Grad-Shafranov equation

- Boundary conditions
    - Fixed boundary : the plasma-vacuum boundary is replaced by the surface of a perfect conductor and only the plasma region is calculated, using as a plasma boundary condition J = 0 such that at this boundary flux = 0

    - Free boundary : the value of magnetic flux is specified on a closed curve, usually in the vacuum region.

    - Constrained boundary : the equilibrium is solved for a given external magnetic field and some constraint such as a contact point with the domain boundary

# Grad-Shafranov equation

- Free boundary Plasma Equilibrium solution

$$\triangle^* \psi(R,Z) = -\mu_0 R J_{\phi,pl}(R,Z),$$
$$J_{\phi,pl}(R,Z) = Rp'(\psi) + \frac{F(\psi)F(\psi)'}{\mu_0 R}$$

$$\triangle^* \psi(R,Z) = -\mu_0 R J_\phi(R,Z),$$
$$J_\phi(R,Z) = J_{\phi,pl}(R,Z) + J_{\phi,cond}(R,Z)$$

Add boundary condition from current density of conductor and plasma

Toroidal current density for conductor and plasma

$$J_{\phi,cond}(R,Z) = \sum_{k=1}^{N_{cond}} J_{cond,k}(R,Z),$$

$$J_{cond,k}(R,Z) = \begin{cases} I_{cond,k}/S_k & \text{if } (R,Z) \in \Omega_{cond,k} \\ 0 & \text{otherwise }, \end{cases}$$

$$J_{\phi,pl}(R,Z) = \begin{cases} \lambda \left[ \beta_0 \frac{R}{R_{geo}} \tilde{j}_{\mathrm{P}}(\psi,\psi_a,\psi_b) + (1-\beta_0) \frac{R_{geo}}{R} \tilde{j}_{\mathrm{F}}(\psi,\psi_a,\psi_b) \right] \\ 0 \end{cases}$$

$$\psi_{bndry}^{(n)}(R_b,Z_b)$$
$$= \int_{\Omega'_{pl}} G(R_b,Z_b;R',Z') J_{\phi,pl}^{(n)}(R',Z') d\Omega'_{pl}$$
$$+ \int_{\Omega'_{cond}} G(R_b,Z_b;R',Z') J_{\phi,cond}^{(n)}(R',Z') d\Omega'_{cond}$$

# Maximum a posteriori policy optimisation

- Conventional formulations
  - Aim to find a trajectory that maximizes expected reward

- Inference formulations
  - starts from a prior distribution over trajectories
  - condition a desired outcome (ex : achieving a goal state)
  - estimate the posteriori distribution over trajectories consistent with this outcome

- ## Why use MPO?
  - Using MPO can overcome the paucity of data by optimizing the policy due to EM algorithm
  - MPO supports data collection across distributed parallel streams and learns in a data-efficiency way



Fig. Training process of A3C models using distributed parallel streams

$$\Delta^* \psi = -\mu_o R^2 \frac{dP}{d\psi} - F \frac{dF}{d\psi}$$

- Likelihood and Posterior
  - Likelihood : Probability of data D given the parameter $\theta$
  - Prior : Distribution over parameter $\theta$
  - Posterior : Probability of the parameter $\theta$ given the data D

$$\theta_{MLE} = \arg\max_{\theta} \prod_{i=1}^{n} f(X_i \mid \theta)$$

$$P(\theta \mid D) = \frac{P(D \mid \theta)\, P(\theta)}{P(D)} \qquad P(D) = \int P(D \mid \theta) P(\theta)\, d\theta$$

- Maximum a posteriori
  - The purpose of machine learning or deep learning : Maximum likelihood estimator
  - Maximizing likelihood * prior = Maximizing posterior

$$\hat{\theta}_{MAP} = argmax_{\theta}\, p(\theta|x) = argmax_{\theta} \frac{p(x|\theta)p(\theta)}{p(x)} = argmax_{\theta}\, p(x|\theta)p(\theta)$$

# Maximum a posteriori policy optimisation

$$\log p_\pi(O=1) = \log \int p_\pi(\tau)p(O=1|\tau)d\tau \geq \int q(\tau)\Big[\log p(O=1|\tau) + \log \frac{p_\pi(\tau)}{q(\tau)}\Big]d\tau$$

$$= \mathbb{E}_q\Big[\sum_t r_t/\alpha\Big] - \text{KL}\Big(q(\tau)||p_\pi(\tau)\Big) = \mathcal{J}(q,\pi),$$

## Optimization : EM algorithm



## RL : Policy Gradient Method

$$\nabla_\theta J(\theta) = \nabla_\theta \sum_{s\in\mathcal{S}} d^\pi(s) \sum_{a\in\mathcal{A}} Q^\pi(s,a)\pi_\theta(a|s)$$

$$\propto \sum_{s\in\mathcal{S}} d^\pi(s) \sum_{a\in\mathcal{A}} Q^\pi(s,a)\nabla_\theta\pi_\theta(a|s)$$

## Approach : Maximum a posteriori



Estimated parameter

Log prior

$$\theta_{\text{MAP}} = \underset{\theta}{\text{argmax}} \left( \log(g(\theta)) + \sum_{i=1}^{n} \log(f(X_i|\theta)) \right)$$

Chose the value of theta that maximizes:

Sum of log likelihood

- Variational EM algorithm
  - Variational Inference : Approximation of the probabilistic distribution q(z) which is most similar to p(z)
  - Unlike sampling-based methods, Variational inference will almost never find the globally optimal solution
  - If they have converged, we can obtain bounds on their accuracy
  - Metric for the measuring difference between two probabilistic distribution : KL Divergence

- Kullback-Leibler Divergence(KL Divergence)
  - Metric for the measuring difference between two probabilistic distribution

$$D_{KL}(P\|Q) = E_{X \sim P}\left[\log \frac{P(x)}{Q(x)}\right] = E_{X \sim P}\left[-\log \frac{Q(x)}{P(x)}\right]$$

- Properties of KL Divergence

$$D_{KL}(P\|Q) = -\sum_x P(x) \log\left(\frac{Q(x)}{P(x)}\right)$$
$$= -\sum_x P(x)\{\log Q(x) - \log P(x)\}$$
$$= -\sum_x \{P(x)\log Q(x) - P(x)\log P(x)\}$$
$$= -\sum_x P(x)\log Q(x) + \sum_x P(x)\log P(x)$$
$$= H(P,Q) - H(P)$$

$$D_{KL}(q\|p) > 0 \ for \ all \ q \ and \ p.$$
$$D_{KL}(q\|p) = 0 \ if \ and \ only \ if \ q = p.$$

- Else
  - Convex functional to q(z)
  - Jensen's inequality satisfied

$$E[f(x)] \geq f(E[x])$$

# Variational EM algorithm

- Idea of EM algorithm
  - We want to find out the specific probabilistic distribution q(z) which is mostly similar to the posterior probabilistic distribution p(z|x)
  - To find out q(z), use KL divergence and optimize the parameters of q(z) and p(z|x) to minimize the KL divergence

- Process
  - E-step : minimize KL divergence of q(z) and p(z|x)
  - M-step : maximize lower bound of log p(x)



(a)   (b)   (c)

Fitting a unimodal approximating distribution q(red) to a multimodal p(blue) (a). use KL(p||q), (b) and (c). use KL(q||p)

- Process
  - E-step : minimize KL divergence of q(z) and p(z|x)

$$\ln \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\theta) = \ln \sum_{\mathbf{Z}} q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \geq \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \right\} = L(q, \theta)$$

$$\ln p(\mathbf{X}|\theta) - L(q, \theta) = \ln p(\mathbf{X}|\theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \right\}$$

$$= \ln p(\mathbf{X}|\theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \theta) p(\mathbf{X}|\theta)}{q(\mathbf{Z})} \right\}$$

$$= \ln p(\mathbf{X}|\theta) - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \theta)}{q(\mathbf{Z})} \right\} - \ln p(\mathbf{X}|\theta) \sum_{\mathbf{Z}} q(\mathbf{Z})$$

$$= - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \theta)}{q(\mathbf{Z})} \right\} = KL[q(\mathbf{Z}) \| p(\mathbf{Z}|\mathbf{X}, \theta)] = KL[q \| p]$$



- Minimization of KL divergence = Maximization of lower bound L

- Process

  - M-step : maximize lower bound of log p(x)

$$q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta^{old})$$

$$L(q, \theta) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{Z}|\mathbf{X}, \theta^{old}) = Q(\theta, \theta^{old}) + const$$

  - Fixed Z, update new parameter $\theta$ to maximize functional L

- Process

  - E-step : variational policy q(a|s) is a parametric variational distribution

$$\max_q \bar{\mathcal{J}}_s(q, \theta_i) = \max_q T^{\pi,q} Q_{\boldsymbol{\theta}_i}(s, a)$$
$$= \max_q \mathbb{E}_{\mu(s)} \left[ \mathbb{E}_{q(\cdot|s)}[Q_{\boldsymbol{\theta}_i}(s, a)] - \alpha \text{KL}(q\|\pi_i) \right]$$
since $V_{\boldsymbol{\theta}_i}(s) = \mathbb{E}_{q(a|s)}[Q_{\boldsymbol{\theta}_i}(s, a)]$ and thus $Q_{\boldsymbol{\theta}_i}(s, a) = r(s, a) + \gamma V_{\boldsymbol{\theta}_i}(s)$

Constraint E step

$$\max_q \mathbb{E}_{\mu(s)} \left[ \mathbb{E}_{q(a|s)} \left[ Q_{\theta_i}(s, a) \right] \right]$$
$$s.t. \mathbb{E}_{\mu(s)} \left[ \text{KL}(q(a|s), \pi(a|s, \theta_i)) \right] < \epsilon$$

  - Non-parametric variational distribution

Convex dual function of $\eta$

$$q_i(a|s) \propto \pi(a|s, \boldsymbol{\theta}_i) \exp \left( \frac{Q_{\theta_i}(s, a)}{\eta^*} \right)$$

$$g(\eta) = \eta \epsilon + \eta \int \mu(s) \log \int \pi(a|s, \theta_i) \exp \left( \frac{Q_{\theta_i}(s, a)}{\eta} \right) da \, ds$$

# Maximum a posteriori policy optimisation

- Process
  - M-step : Given q from the E-step, optimize the lower bound J with respect to $\theta$ to obtain an updated policy

$$\max_{\boldsymbol{\theta}} \mathcal{J}(q_i, \theta) = \max_{\boldsymbol{\theta}} \mathbb{E}_{\mu_q(s)}\left[\mathbb{E}_{q(a|s)}\left[\log \pi(a|s, \theta)\right]\right] + \log p(\theta)$$

  - Generalized M-step

$$\max_{\pi} \mathbb{E}_{\mu_q(s)}\left[\mathbb{E}_{q(a|s)}\left[\log \pi(a|s, \theta)\right] - \lambda \text{KL}\left(\pi(a|s, \theta_i), \pi(a|s, \theta)\right)\right]$$

  - Constrained M-step

$$\max_{\pi} \mathbb{E}_{\mu_q(s)}\left[\mathbb{E}_{q(a|s)}\left[\log \pi(a|s, \theta)\right]\right]$$
$$s.t. \ \mathbb{E}_{\mu_q(s)}\left[\text{KL}(\pi(a|s, \theta_i), \pi(a|s, \theta))\right] < \epsilon.$$

- Process
    - Policy Evaluation : fitting value function to policy
    - Use an experience replay buffer and arbitrary behaviour policy given by the action probabilities stored in the buffer
    - Off-policy + Multi-step returns : Retrace algorithm(Munos et al, 2016)

$Q_\theta(s,a)$

$Q_\phi(s,a)$

Q Network

Update $\theta$

Target Network

L($\phi$)

$$\min_\phi L(\phi) = \min_\phi \mathbb{E}_{\mu_b(s),b(a|s)}\left[\left(Q_{\theta_i}(s_t,a_t,\phi) - Q_t^{\mathrm{ret}}\right)^2\right], \text{with}$$

$$Q_t^{\mathrm{ret}} = Q_{\phi'}(s_t,a_t) + \sum_{j=t}^{\infty}\gamma^{j-t}\Big(\prod_{k=t+1}^{j}c_k\Big)\Big[r(s_j,a_j) + \mathbb{E}_{\pi(a|s_{j+1})}[Q_{\phi'}(s_{j+1},a)] - Q_{\phi'}(s_j,a_j)\Big]$$

$$c_k = \min\Big(1, \frac{\pi(a_k|s_k)}{b(a_k|s_k)}\Big),$$

# Maximum a posteriori policy optimisation

- Process
    - E-step

$$\max_q \bar{\mathcal{J}}_s(q, \theta_i) = \max_q T^{\pi, q} Q_{\boldsymbol{\theta}_i}(s, a)$$

$$= \max_q \mathbb{E}_{\mu(s)} \Big[ \mathbb{E}_{q(\cdot|s)}[Q_{\boldsymbol{\theta}_i}(s, a)] - \alpha \text{KL}(q \| \pi_i) \Big]$$

since $V_{\boldsymbol{\theta}_i}(s) = \mathbb{E}_{q(a|s)}[Q_{\boldsymbol{\theta}_i}(s, a)]$ and thus $Q_{\boldsymbol{\theta}_i}(s, a) = r(s, a) + \gamma V_{\boldsymbol{\theta}_i}(s)$

    - M-step

$$\max_{\boldsymbol{\theta}} \mathcal{J}(q_i, \theta) = \max_{\boldsymbol{\theta}} \mathbb{E}_{\mu_q(s)} \Big[ \mathbb{E}_{q(a|s)} \big[ \log \pi(a|s, \theta) \big] \Big] + \log p(\theta)$$

    - Policy Evaluation

$$\min_\phi L(\phi) = \min_\phi \mathbb{E}_{\mu_b(s), b(a|s)} \Big[ \big( Q_{\theta_i}(s_t, a_t, \phi) - Q_t^{\text{ret}} \big)^2 \Big], \text{with}$$

$$Q_t^{\text{ret}} = Q_{\phi'}(s_t, a_t) + \sum_{j=t}^{\infty} \gamma^{j-t} \Big( \prod_{k=t+1}^{j} c_k \Big) \Big[ r(s_j, a_j) + \mathbb{E}_{\pi(a|s_{j+1})}[Q_{\phi'}(s_{j+1}, a)] - Q_{\phi'}(s_j, a_j) \Big]$$

$$c_k = \min \Big( 1, \frac{\pi(a_k|s_k)}{b(a_k|s_k)} \Big),$$

TCV#69545

# Conclusion

- Highlight
  - High performance + robustness to uncertain operating conditions, intuitive target specification and unprecedented versatility
  - An accurate, numerically robust simulator(SFC) : trade-off between simulation accuracy and computational complexity / Highly data-efficient RL algorithm(MPO) that scales to high-dimensional problems
  - Free-boundary equilibrium evolution model has sufficient fidelity to develop transferable controllers, offering a justification for using this approach to test control of future devices

- Enhancement
  - Analysis of the non-linear dynamics
  - Reduce training time through increased reuse of data and multi-fidelity learning
  - The model architecture can be coupled to a more capable simulator(i.e. plasma pressure, current density evolution physics)

Source code : https://github.com/deepmind/deepmind-research/tree/master/fusion_tcv

# Conclusion



Direct shape optimization through deep reinforcement learning

Jonathan Viquerat [a,*], Jean Rabault [b], Alexander Kuhnle [c], Hassan Ghraieb [a], Aurélien Larcher [a], Elie Hachem [a]

[a] MINES Paristech, CEMEF, PSL - Research University, France
[b] Department of Mathematics, University of Oslo, Norway
[c] University of Cambridge, United Kingdom of Great Britain and Northern Ireland

ABSTRACT

Deep Reinforcement Learning (DRL) has recently spread into a range of domains within physics and engineering, with multiple remarkable achievements. Still, much remains to be explored before the capabilities of these methods are well understood. In this paper, we present the first application of DRL to direct shape optimization. We show that, given adequate reward, an artificial neural network trained through DRL is able to generate optimal shapes on its own, without any prior knowledge and in a constrained time. While we choose here to apply this methodology to aerodynamics, the optimization process itself is agnostic to details of the use case, and thus our work paves the way to new generic shape optimization strategies both in fluid mechanics, and more generally in any domain where a relevant reward function can be defined.

## 1. Introduction

Shape optimization is a long-standing research topic with countless industrial applications, ranging from structural mechanics to electromagnetism and biomechanics [1] [2]. In fluid dynamics, the interest in shape optimization has been driven by many real-world problems. For example, within aerodynamics, the reduction of drag and therefore of fuel consumption by trucks and cars [3], or the reduction of aircraft fuel consumption and running costs [4], are cases on which a large body of literature is available. However, shape optimization also plays a key role in many other aspects of the performance of, for example, planes, and modern optimization techniques are also applied to a variety of problems such as the optimization of electromagnetically stealth aircrafts [5], or acoustic noise reduction [6]. This illustrates the importance of shape optimization methods in many applications, across topics of both academic and industrial interest.

In the following, we will focus on one shape optimization problem, namely airfoil shape optimization. This problem is key to many industrial processes, and presents the ingredients of non-linearity and high dimensionality which make shape optimization at large a challenging problem. As a consequence of its industrial and academic relevance, airfoil shape optimization through numerical techniques has been discussed since at least back to 1964 [7], and remains a problem of active research [8,9].

Jonathan Viquerat et al, 2021, Journel of computational physics

# Conclusion



Deep reinforcement learning for heat exchanger shape optimization

Hadi Keramati[a,*], Feridun Hamdullahpur[a], Mojtaba Barzegari[b]

[a] Department of Mechanical And Mechatronics Engineering, University of Waterloo, 200 University Avenue, West Waterloo, Ontario N2L 3G1, Canada
[b] Biomechanics Section, Department of Mechanical Engineering, KU Leuven, Leuven, Belgium

**ABSTRACT**

We present a parametric approach for heat exchanger shape optimization utilizing Deep Reinforcement Learning (Deep RL) and Boundary Representation (BREP). In this study, we show that continuous geometric representation of the fluid and solid domain facilitates the implementation of boundary conditions and design space exploration in contrast to traditional Topology Optimization such as density-based methods. The proposed framework consists of a Deep Neural Network (DNN), a Computational Fluid Dynamics (CFD) solver with an automatic body-fitted mesh generation to solve a single fin shape optimization. The learning is performed using Proximal Policy Optimization (PPO) in combination with a CFD environment in FEniCS. The RL agent successfully explores the design space and maximizes heat transfer and minimizes pressure drop for geometric design with as low as 12 degrees of freedom represented by composite Bézier curves. Higher degree of freedom results in higher reward of the agent. This method alleviates the curse of dimensionality compared to voxel and pixel-based optimization of coupled thermal fluid-structure. Results show the manufacturability and efficiency of the output of our framework. Over 30 percent improvement in overall heat transfer while lowering the pressure drop by more than 60 percent compared to the rectangle reference geometry is achieved.

## 1. Introduction

Advances in Additive Manufacturing techniques allow fabrication of complex geometries which are too challenging using conventional manufacturing methods [1–5]. Normally, limited number of pre-defined manufacturable shapes are considered for the design purpose. Therefore, shape and Topology Optimization (TO) receive more attention as mathematical methods that optimize geometry. TO has a wide range of industrial applications from structural to biomechanics applications [6,7]. Several methods have been used for topology optimization such as density based, phase field, and shape derivative methods. The density-based topology optimization procedure known as the SIMP (Solid isotropic material with penalization) has been used for decades in structural optimization; however, this method is not practical for thermofluid application since a porous medium approach is used for solid distribution. In fluid flow and particularly heat exchanger design, boundary conditions are a significant part of the simulation and it is import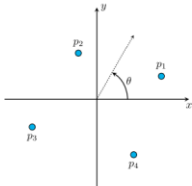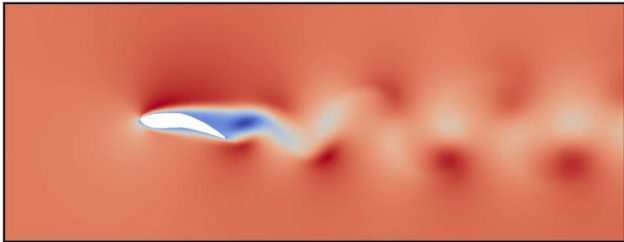ant to have distinguishable boundaries [8–10]. Moreover, for several applications such as those in biomedical engineering,

working fluids (e.g. nanofluids) and particular boundary conditions require clear boundaries for implementation [11–16].

A convenient way of optimizing the design while keeping the boundaries distinguishable is using pixel or voxel-based optimization which require high CPU time due to the curse of dimensionality [17]. In thermofluid design, the direction to greater performance is aligned with greater design freedom [18]. Providing freedom in design is, however, computationally demanding and impose manufacturing challenges. Topological representation using parametric curve reduces the dimension of the optimization problem without sacrificing the resolution of the shape or freedom to change the design [19]. This representation provides exact geometrical definition regardless of the size of discretization in numerical approach [20]. Recently, Florian et al [21] used Null space and gradient methods to maintain clear shape boun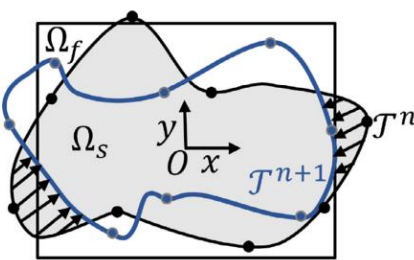daries during the shape optimization. The results presented are promising, but require further improvement in terms of constraint implementation. The state-of-the-art methods in heat transfer shape optimization rely on adjoint method which computes the derivative of the objective with respect to design variables to specify the direction of the search algorithm which could be trapped in local optimal design candidates [22,23].

Reinforcement Learning (RL) is one of the basic components of Machine Learning (ML). Many types of algorithms are introduced for classification, regression, clustering, etc. RL, however, is a real

* Corresponding author.
  E-mail address: hkeramati@uwaterloo.ca (H. Keramati).

Hadi Keramati al, 2022, International Journal of heat and mass transfer





$$\sigma_f(\boldsymbol{u}, p) = 2\nu e(\boldsymbol{u}) - pI, e(\boldsymbol{u}) = \frac{1}{2}(\nabla \boldsymbol{u} + \nabla \boldsymbol{u}^T)$$

$$\begin{cases} \rho \frac{\partial \boldsymbol{u}}{\partial t} + \rho((\boldsymbol{u} \cdot \nabla)\boldsymbol{u}) - \nabla(\sigma_f(\boldsymbol{u}, p)) = 0 & in \quad \Omega_f \\ \nabla \cdot \boldsymbol{u} = 0 & in \quad \Omega_f \\ \boldsymbol{u} = \boldsymbol{u}_0 & on \quad \Omega_{f,in} \\ \sigma_f(\boldsymbol{u}, p) \cdot \boldsymbol{n} = 0 & on \quad \Omega_{f,out} \\ \boldsymbol{u} = 0 & on \quad \Gamma \end{cases}$$

$$\begin{cases} \rho c_p \frac{\partial T}{\partial t} + \rho c_p(\boldsymbol{u} \cdot \nabla T) - \nabla \cdot (k_f \nabla T) = 0 & in \quad \Omega_f \\ T = T_{in} & on \quad \Omega_{f,in} \\ T = T_0 & on \quad \Omega_{f,t=0} \\ T_f = T_s & on \quad \Gamma \\ -k_f \frac{\partial T_f}{\partial n} = -k_s \frac{\partial T_s}{\partial n} & on \quad \Gamma \end{cases}$$

# Appendix : Reward components

| 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. |

| | Fundamental Capability | Elongated shape | ITER-like shape | Negative Triangularity | Snowflake | Droplets |
|---|---|---|---|---|---|---|
| Figure | Fig. 2 | Fig. 3a, Extended Data Fig. 2a | Fig. 3b, Extended Data Fig. 2b, 3 | Fig. 3c, Extended Data Fig. 2c | Fig. 3d, Extended Data Fig. 2d, 4c | Fig. 4 |
| Shot | TCV#70915 | TCV#70920 | TCV#70600 | TCV#70457 | TCV#70755 | TCV#69545 |
| Reward Component | Transforms, Combiners (if necessary), and Weight (default=1) | | | | | |
| Diverted | | | Equal() | Equal() | | |
| E/F Currents | | SoftPlus(good=100, bad=50) GeometricMean() | SoftPlus(good=100, bad=50) GeometricMean() | SoftPlus(good=100, bad=50) GeometricMean() | SoftPlus(good=100, bad=50) GeometricMean() | |
| Elongation | | SoftPlus(good=0.005, bad=0.2) | | SoftPlus(good=0, bad=0.5) | | |
| LCFS Distance | SoftPlus(good=0.005, bad=0.05) SmoothMax(-1) | SoftPlus(good=0.003, bad=0.03) SmoothMax(-1) weight=3 | SoftPlus(good=0.005, bad=0.05) SmoothMax(-1) weight=3 | SoftPlus(good=0.005, bad=0.05) SmoothMax(-1) weight=3 | SoftPlus(good=0.005, bad=0.05) SmoothMax(-1) weight=3 | |
| Legs Normalized Flux | | | Sigmoid(good=0.1, bad=0.3) SmoothMax(-5) weight=2 | | | |
| Limit Point | Sigmoid(good=0.1, bad=0.2) | Sigmoid(good=0.2, bad=0.3) | | | Sigmoid(good=0.1, bad=0.2) | |
| OH Current Diff | SoftPlus(good=50, bad=1050) | ClippedLinear(good=50, bad=1050) | ClippedLinear(good=50, bad=1050) | ClippedLinear(good=50, bad=1050) | ClippedLinear(good=50, bad=1050) | ClippedLinear(good=50, bad=1050) |
| Plasma Current | SoftPlus(good=500, bad=20000) | SoftPlus(good=500, bad=30000) | SoftPlus(good=500, bad=20000) weight=2 | SoftPlus(good=500, bad=20000) weight=2 | SoftPlus(good=500, bad=20000) weight=2 | Sigmoid(good=2000, bad=20000) weight=[1, 1] |
| R | | | | | | Sigmoid(good=0.02, bad=0.5) weight=[1, 1] |

| | Fundamental Capability | Elongated shape | ITER-like shape | Negative Triangularity | Snowflake | Droplets |
|---|---|---|---|---|---|---|
| Radius | | SoftPlus(good=0.002, bad=0.02) | | SoftPlus(good=0, bad=0.04) | | |
| Triangularity | | SoftPlus(good=0.005, bad=0.2) | | SoftPlus(good=0, bad=0.5) | | |
| Voltage Out of Bounds | | Mean() SoftPlus(good=0, bad=1) | Mean() SoftPlus(good=0, bad=1) | Mean() SoftPlus(good=0, bad=1) | Mean() SoftPlus(good=0, bad=1) | |
| X-point Count | Equal() | | | | | |
| X-point Distance | Sigmoid(good=0.01, bad=0.15) weight=0.5 | | Sigmoid(good=0.01, bad=0.15) | Sigmoid(good=0.02, bad=0.15) weight=[0.5, 0.5] | Sigmoid(good=0.01, bad=0.15) weight=[0.5, 0.5] | |
| X-point Far | Sigmoid(good=0.3, bad=0.1) SmoothMax(-5) | | | | | |
| X-point Flux Gradient | SoftPlus(good=0, bad=3) weight=0.5 | | SoftPlus(good=0, bad=3) weight=0.5 | SoftPlus(good=0, bad=3) weight=[0.5, 0.5] | SoftPlus(good=0, bad=3) weight=[0.5, 0.5] | |
| X-point Normalized Flux | SoftPlus(good=0, bad=0.08) | | SoftPlus(good=0, bad=0.08) | SoftPlus(good=0, bad=0.08) weight=[1, 1] | SoftPlus(good=0, bad=0.08) weight=[1, 1] | |
| Z | | | | | | Sigmoid(good=0.02, bad=0.2) weight=[1, 1] |

# Appendix : Reward elements

| Transform | Description |
|---|---|
| ClippedLinear | Linearly maps the input values such that the good goes to 1 and bad to 0, then clips between 0 and 1. |
| Equal | Returns 1 if there is no error, returns 0 otherwise. Useful for boolean or integer outputs. |
| Sigmoid | Maps the input values such that good is 0.95 and bad is 0.05 in the output of the logistic function. This is similar to ClippedLinear, except there's still small impetus to improve beyond the good value and a little bit of reward signal for improvements below the bad value. |
| SoftPlus | Maps the input values such that good is 1 and bad is 0.1 in the output of the lower half of the logistic function, then clips to 0 and 1. This leads to a sharp drop-off as the value moves away from the good value, and a slow drop-off past bad. This is similar to a smooth relu. |

| Combiner | Formula | Description |
|---|---|---|
| Geometric Mean | $\left(\prod_{i=1}^{n} x_i w_i\right)^{\frac{1}{\sum_{i=1}^{n} w_i}}$ | Takes the weighted geometric mean of the values. |
| Mean | $\dfrac{\sum_{i=1}^{n} x_i w_i}{\sum_{i=1}^{n} w_i}$ | Takes the weighted mean of the values. |
| Smooth Max | $\dfrac{\sum_{i=1}^{n} x_i w_i e^{\alpha x_i}}{\sum_{i=1}^{n} w_i e^{\alpha x_i}}$ | Takes the smooth maximum, parameterized with an $\alpha$ such that $\alpha = 0$ is equivalent to taking the mean, $\alpha = -\infty$ is equivalent to taking the minimum, and $\alpha = +\infty$ is equivalent to taking the maximum. |

| Termination | Termination Criteria |
|---|---|
| Coil current limits | Any coil current exceeds the physical limit of the plant. |
| Edge safety factor | Terminate when the edge safety factor $q_{95}$ goes below 2.2, which provides some margin over the the threshold for a stable plasma ($q_{95} > 2$). |
| OH too different | The OH coil currents differ by more than 4 kA, which would cause high structural forces. |
| Plasma current limit | Plasma current is below the plant's disruption detector threshold, which is $-60$ kA for a single plasma, and $-25$ kA per plasma for droplets. |
| Solver not converged | Multiple subsequent simulation steps did not converge. |

# Reference

[1] Magnetic control of tokamak plasmas through deep reinforcement learning, Jonas Degrave et al, Nature, 2022

[2] Maximum a posteriori policy optimization, Abbas Abdolmaleki et al, ICLR, 2018

[3] Development of free boundary equilibrium and transport solvers for simulation and real-time interpretation of tokamak experiments, Francesco Carpanese, EPFL

[4] Policy Gradient Methods for Reinforcement Learning with Function Approximation, Richard S.Sutton et al, 1999

[5] An Introduction to Reinforcement Learning, Vincent Francois-Lavet et al, 2018

[6] On Multi-objective Policy Optimization as a Tool for Reinforcement Learning, Abbas Abdolmaleki et al, 2021