

# Can Increasing Input Dimensionality Improve Deep Reinforcement Learning?

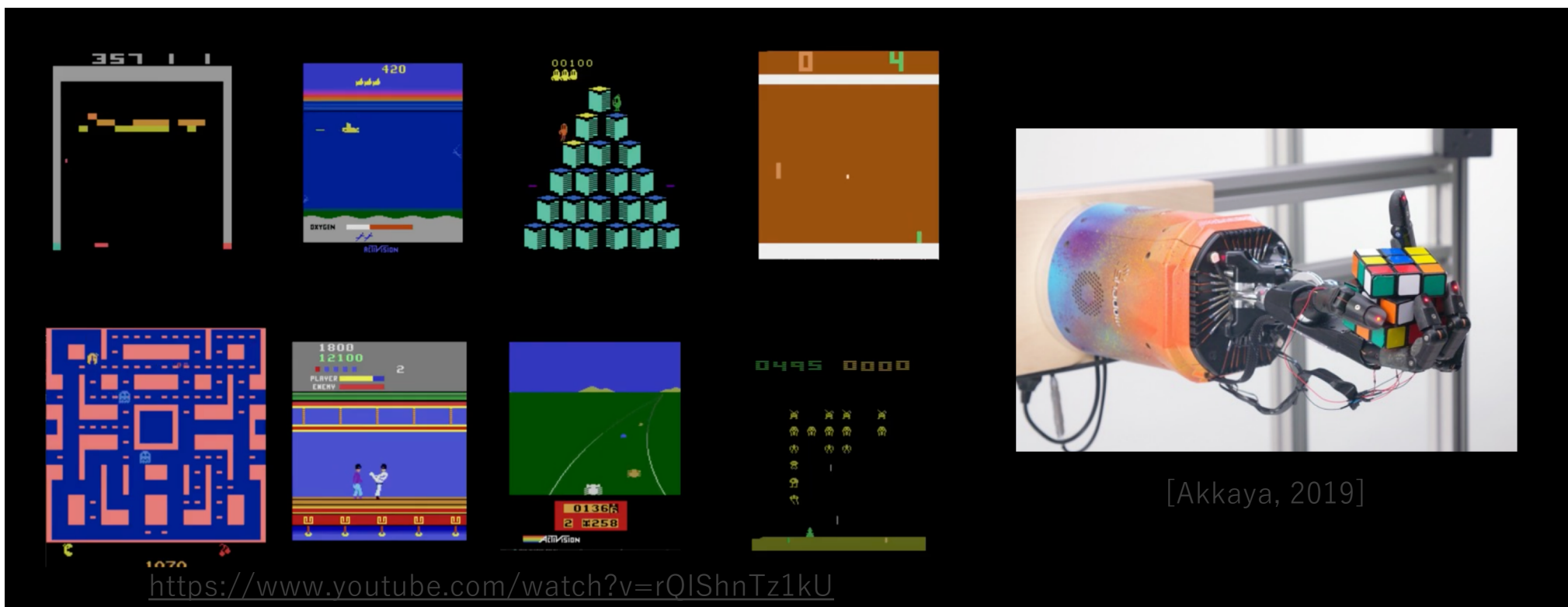
Kei Ota<sup>1</sup>, Tomoaki Oiki<sup>1</sup>, Devesh K. Jha<sup>2</sup>,  
Toshisada Mariyama<sup>1</sup>, and Daniel Nikovski<sup>2</sup>

1. Mitsubishi Electric, Kanagawa, JP

2. Mitsubishi Electric Research Labs, MA, US.

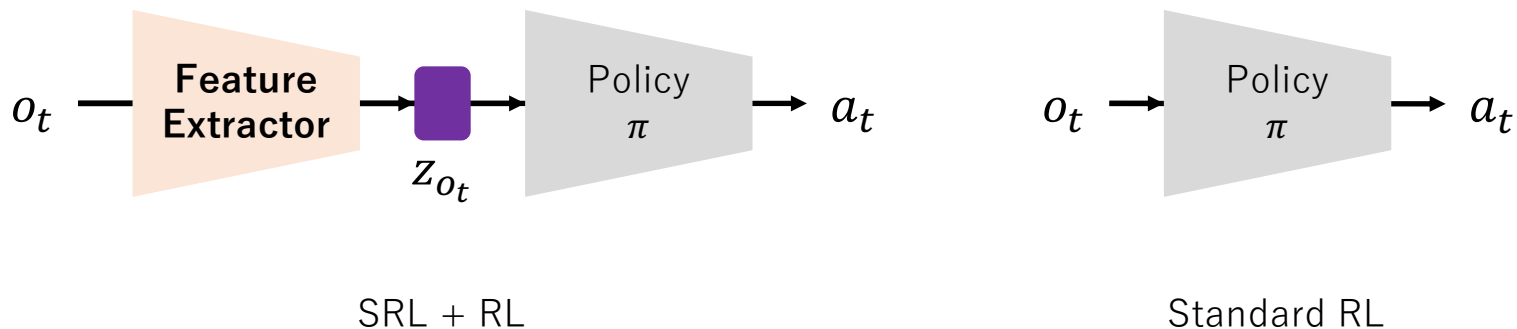
# Introduction

- Deep RL algorithms have achieved impressive success
  - ✓ Can solve complex tasks
  - ✗ Learning representations requires a large amount of data



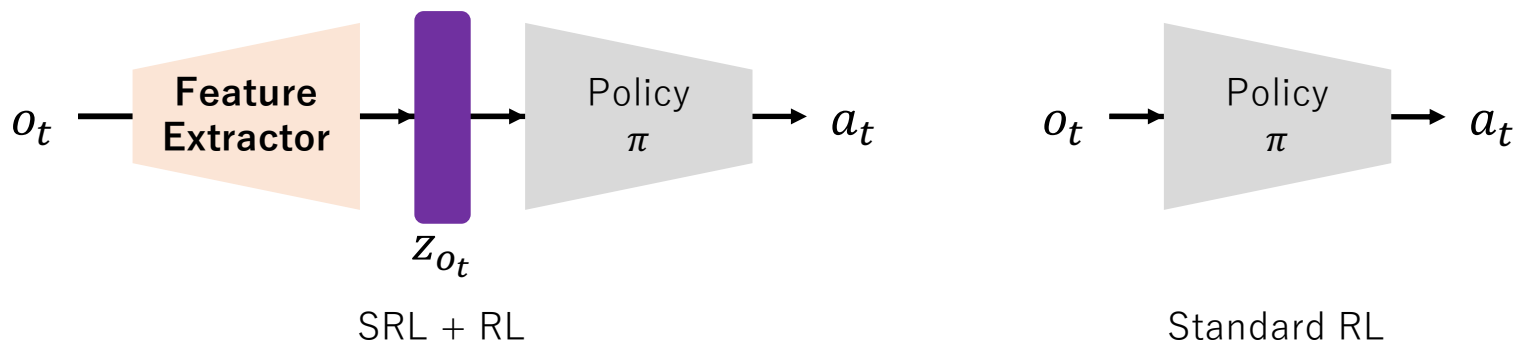
# Introduction

- Deep RL algorithms have achieved impressive success
  - ✓ Can solve complex tasks
  - ✗ Learning representations requires a large amount of data
- State Representation Learning (SRL)
  - Learned features are in low dimension, evolve through time, and are influenced by actions of an agent
  - The lower the dimensionality, the faster and better RL algorithms will learn



# Introduction

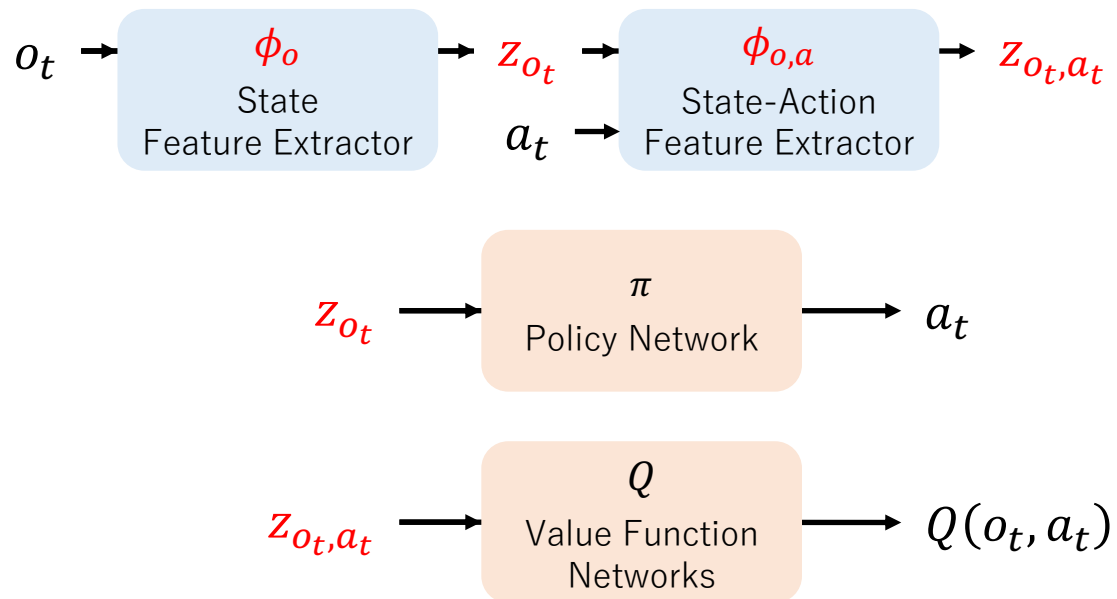
- Deep RL algorithms have achieved impressive success
  - ✓ Can solve complex tasks
  - ✗ Learning representations requires a large amount of data
- State Representation Learning (SRL)
  - Learned features are in low dimension, evolve through time, and are influenced by actions of an agent
  - The lower the dimensionality, the faster and better RL algorithms will learn



**Can Increasing Input Dimensionality Improve Deep RL?**

# OFENet: Online Feature Extractor Network

- OFENet
  - Train feature extractor network  $\phi_o$  and  $\phi_{o,a}$  that produces **high-dimensional** representation  $z_{o_t}$  and  $z_{o_t,a_t}$



# OFENet: Online Feature Extractor Network

- OFENet

- Train feature extractor network  $\phi_o$  and  $\phi_{o,a}$  that produces **high-dimensional** representation  $z_{o_t}$  and  $z_{o_t,a_t}$



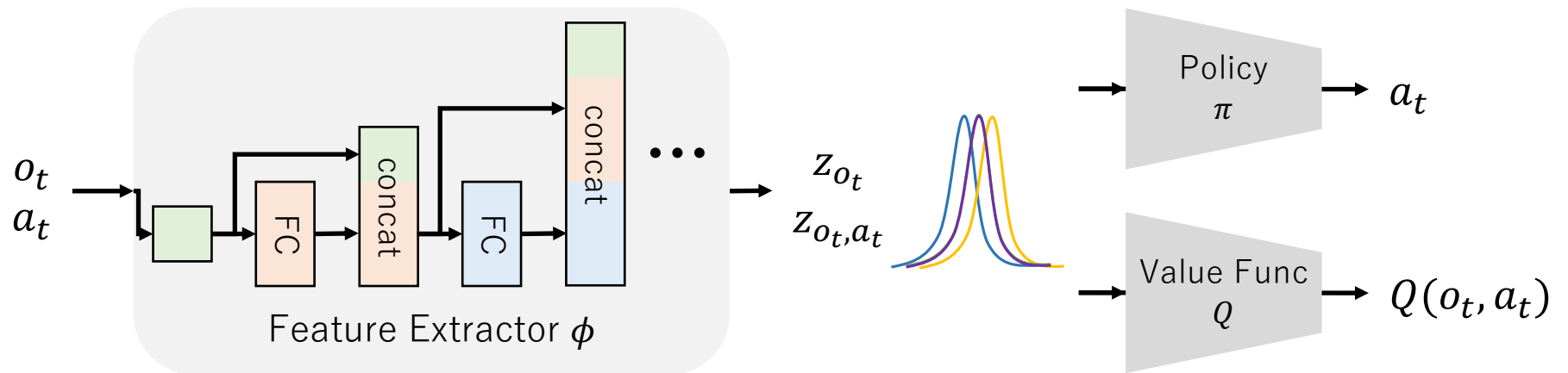
- Optimize  $\theta_{\text{aux}} = \{\theta_{\phi_o}, \theta_{\phi_{o,a}}, \theta_{\text{pred}}\}$  by learning to predict next state

$$L_{\text{aux}} = \mathbb{E}_{(o_t, a_t) \sim p, \pi} [\|f_{\text{pred}}(z_{o_t, a_t}) - o_{t+1}\|^2]$$

- Increasing the search space allows the agent to learn much more complex policies

# Network Architecture

- What is best architecture to extract features?
  - Deeper networks: optimization ability and expressiveness
  - Shallow layers: physically meaningful output
- MLP DenseNet
  - Combine advantages of deep layers and shallow layers



- Use **Batch Normalization** to suppress changes in input distributions

---

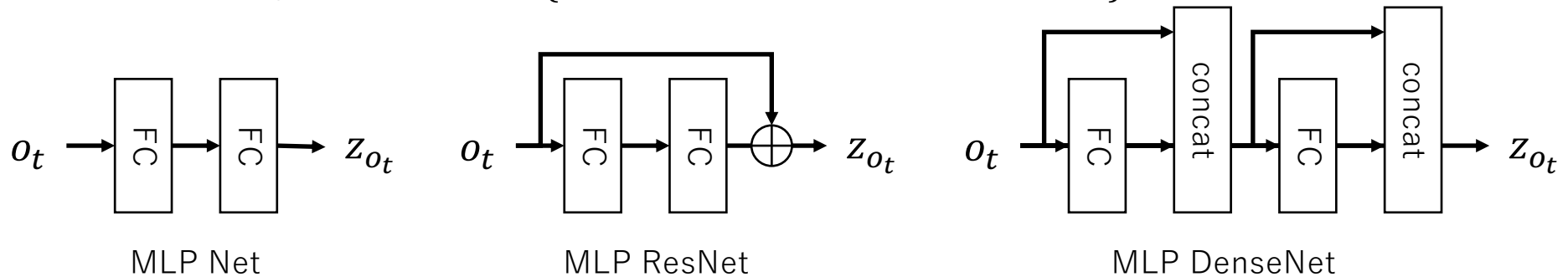
# Experiments

1. **What is a good architecture** that learns effective state and state-action representations for training better RL agents?
2. **Can OFENet learn more sample efficient and better performant** policies when compared to some of the state-of-the-art techniques?
3. **What leads to the performance gain** obtained by OFENet?



# What is a good architecture?

- Compare aux. score and actual RL score to search a good architecture from:
  - Connectivity architecture: {MLP, MLP ResNet, MLP DenseNet}



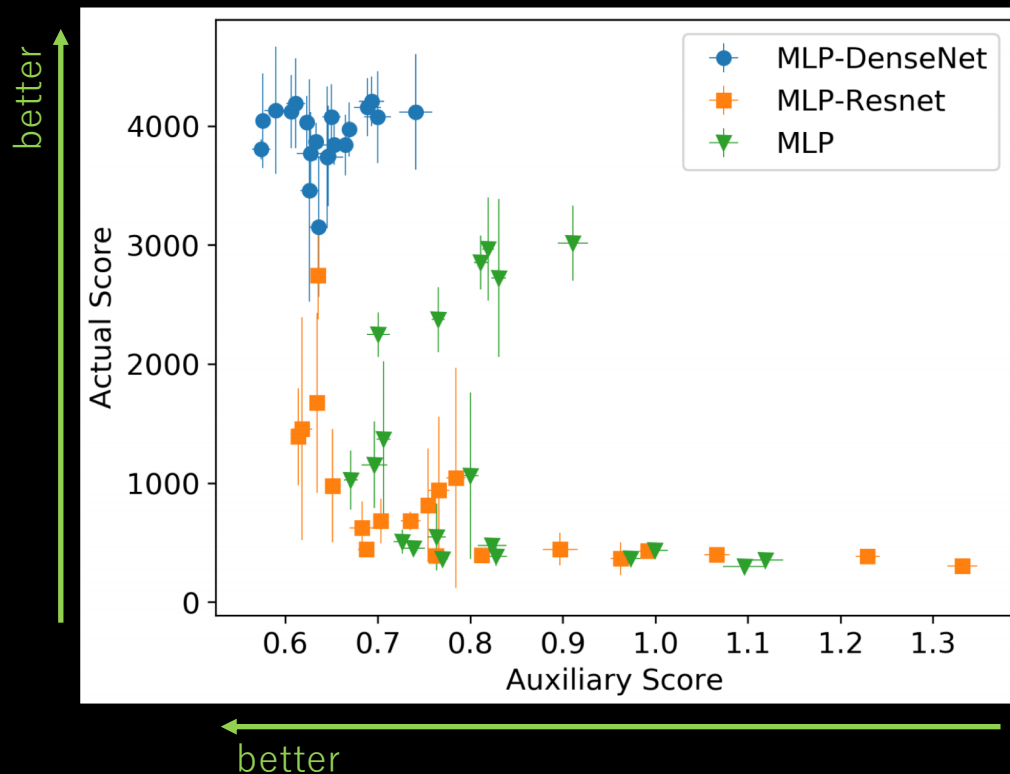
- Number of layers:  $n_{\text{layers}} \in \{1,2,3,4\}$  for MLP,  $n_{\text{layers}} \in \{2,4,6,8\}$  for others
- Activation function: {ReLU, tanh, Leaky ReLU, swish, SELU}

- Aux. score: randomly collect 100K transitions for training, 20K for evaluation

$$L_{aux} = \mathbb{E}_{(o_t, a_t) \sim p, \pi} \left[ \|f_{pred}(z_{o_t, a_t}) - o_{t+1}\|^2 \right]$$

- Actual score: measure returns of SAC agent with 500K steps training

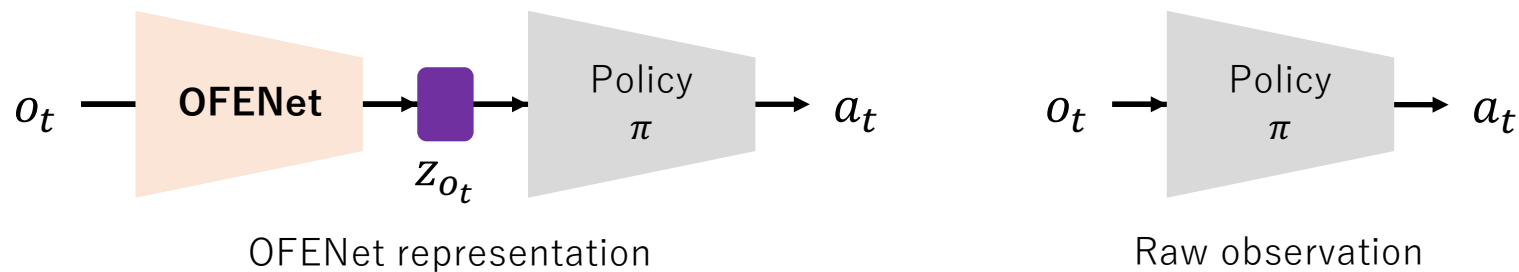
# What is a good architecture?



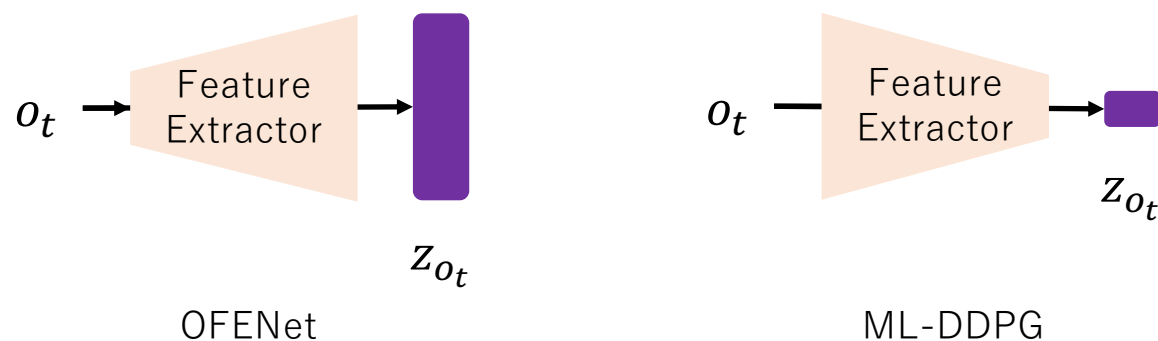
- **MLP-DenseNet** consistently achieves higher actual score
- Smaller the aux. score, better the actual score
- We can select architecture with the smallest aux. score without solving heavy RL problem!

## More sample efficient and better performant polices?

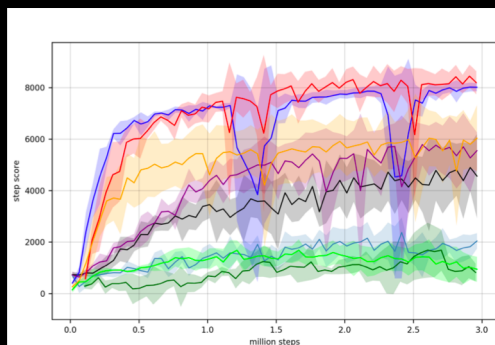
- Measure performance of SAC, TD3, and PPO with and without OFENet
  - No changes in hyperparameters for each algorithm



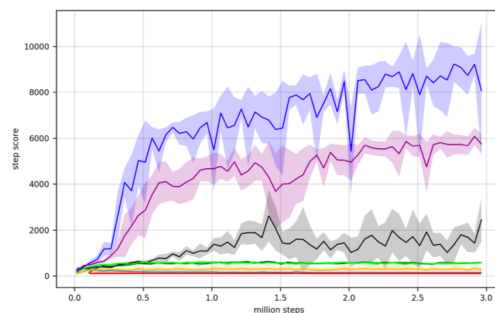
- Compare to closest work: ML-DDPG [Munk2016]
  - **Reduce the dimension** of the observation to one third of its original



# More sample efficient and better performant polices?



(d) Ant-v2

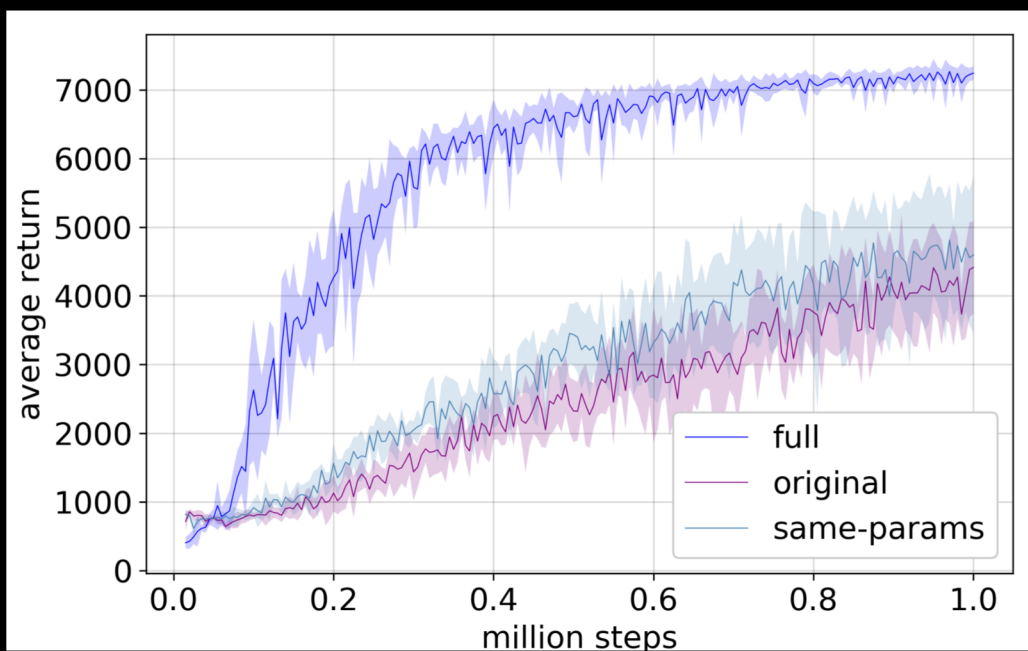


(e) Humanoid-v2

- OFENet improves sample efficiency and returns without changing any hyperparameters
- OFENet effectively learns meaningful features

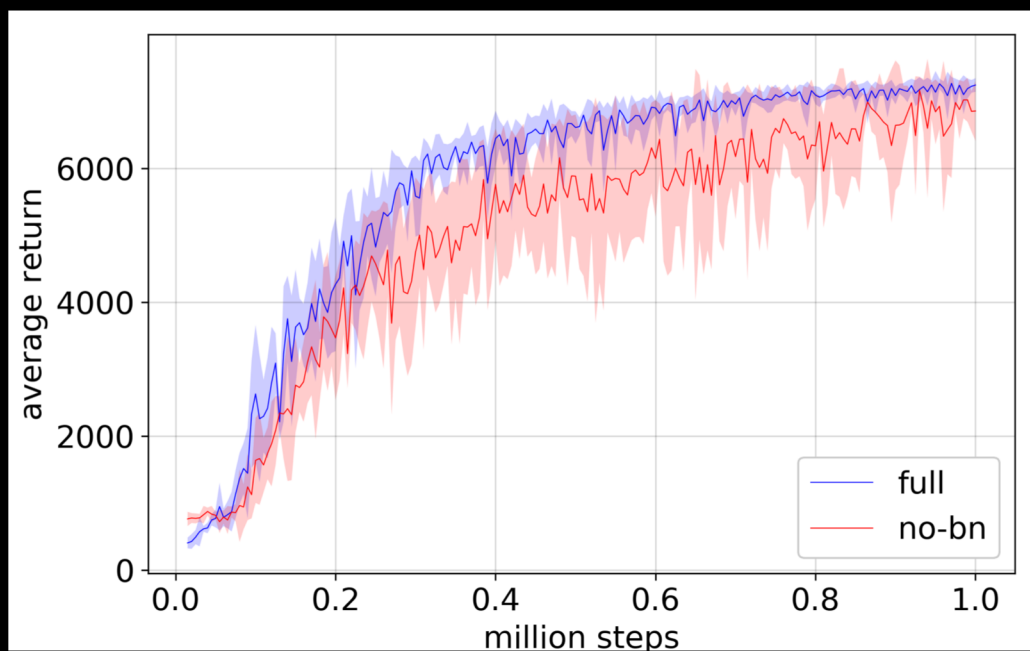
	OFE (OURS)	SAC			TD3		PPO	
		Original	ML-SAC (1/3)	ML-SAC (OFE like)	OFE (OURS)	Original	OFE (OURS)	Original
HOPPER-V2	<b>3511.6</b>	3316.6	750.5	868.7	3488.3	<b>3613.0</b>	<b>2525.6</b>	1753.5
WALKER2D-V2	<b>5237.0</b>	3401.5	667.4	627.4	<b>4915.1</b>	4515.6	<b>3072.1</b>	3016.7
HALFCHEETAH-V2	<b>16964.1</b>	14116.1	1956.9	11345.5	<b>16259.5</b>	13319.9	<b>3981.8</b>	2860.4
ANT-V2	<b>8086.2</b>	5953.1	4950.9	2368.3	<b>8472.4</b>	6148.6	<b>1782.3</b>	1678.9
HUMANOID-V2	<b>9560.5</b>	6092.6	3458.2	331.7	120.6	<b>345.2</b>	<b>670.3</b>	652.4

## What leads to the performance gain?



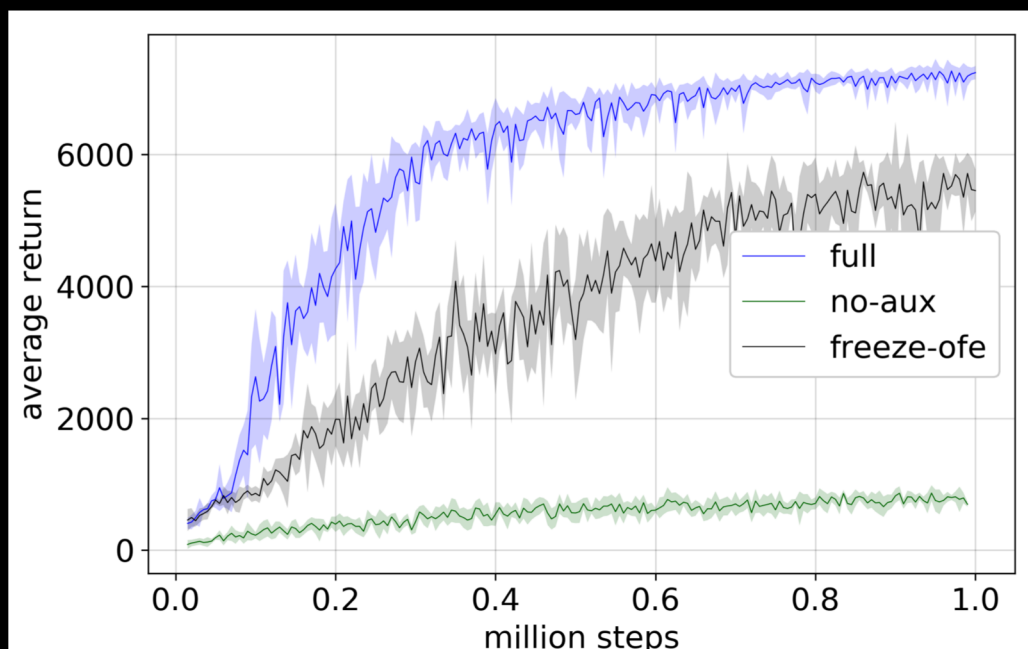
- Just increasing network size doesn't improve performance

## What leads to the performance gain?



- Just increasing network size doesn't improve performance
- **BN stabilizes training**

## What leads to the performance gain?



- Just increasing network size doesn't improve performance
- BN stabilizes training
- **Decoupling feature extraction and control policy is important**
- **Online SRL handles unknown distribution during training**

## Conclusion

- Proposed Online Feature Extractor Network (OFENet)
  - Provides much higher-dimensional representation
  - Demonstrated OFENet can significantly accelerate RL
- OFENet can be used as New RL tool box
  - Just put OFENet as base layer of RL algorithms
  - No need to tune hyperparameters of original algorithms!
  - Code link: [www.merl.com/research/license/OFENet](http://www.merl.com/research/license/OFENet)

**Can increasing input dimensionality improve deep RL?**

**Yes, it can!**