



UNIVERSITÄT  
PADERBORN

COMPUTER NETWORKS GROUP

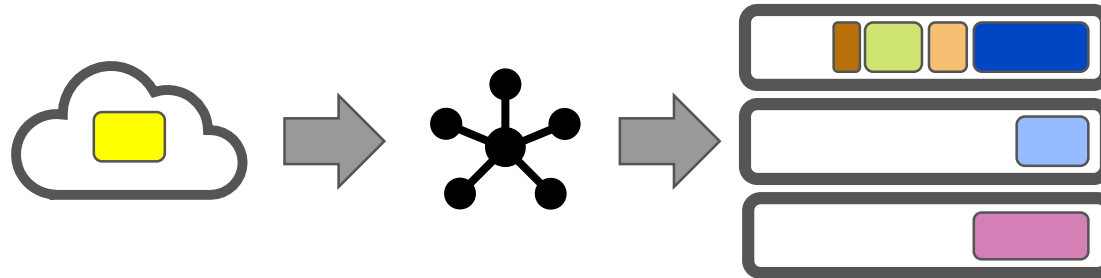
# VARIANCE REDUCTION FOR REINFORCEMENT LEARNING IN INPUT-DRIVEN ENVIRONMENTS

AICON PROJECT GROUP

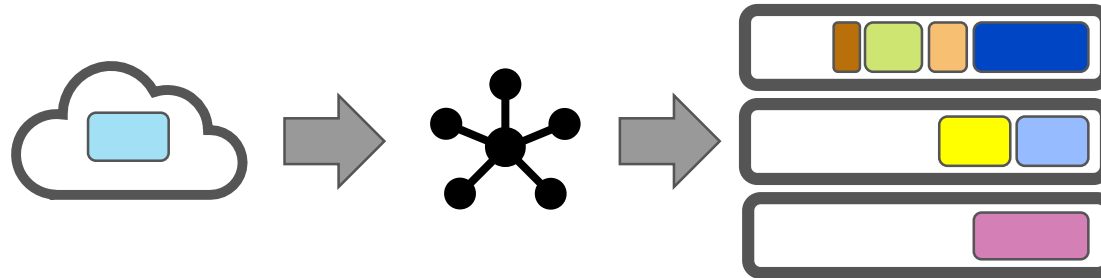
Stefan Werner



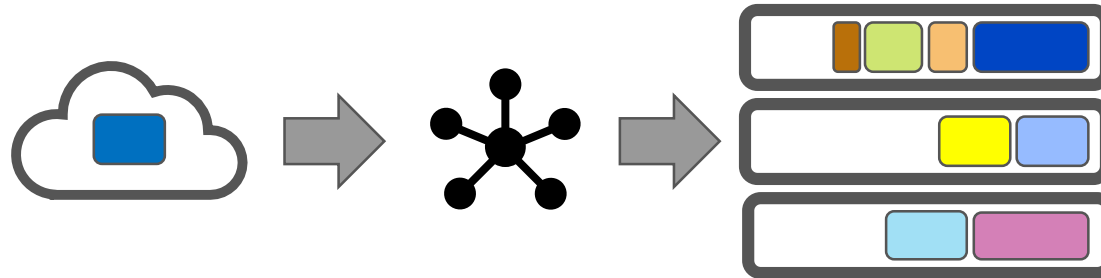
# Foundations



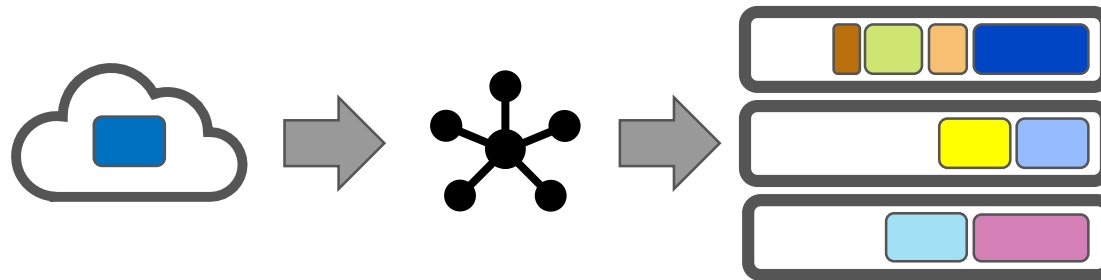
# Foundations



# Foundations



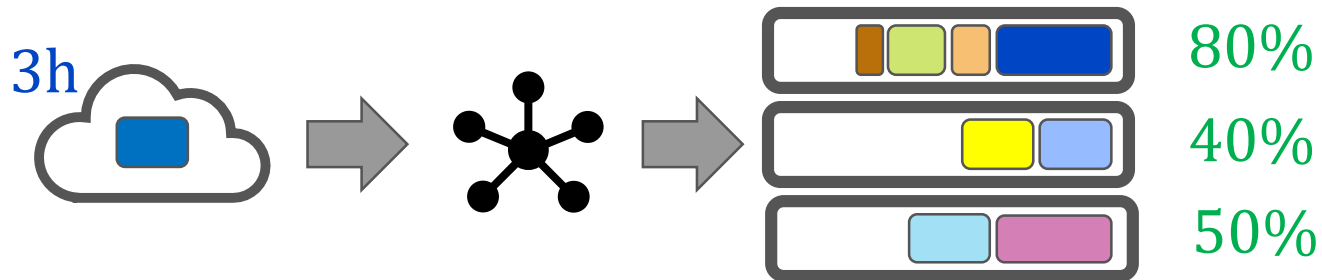
# Foundations



## MDP

- *Trajectories* :  $\tau = (s_1, a_1, r_1, s_2, a_2, r_2, s_3, \dots, s_T)$

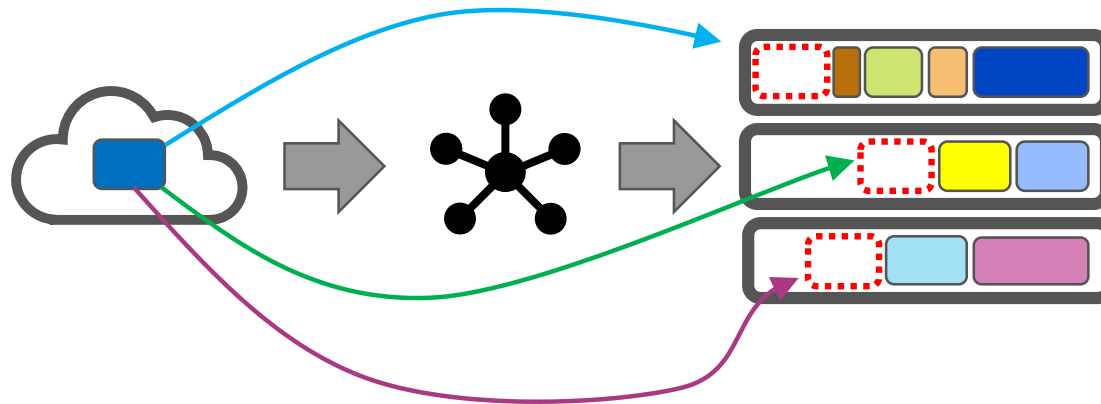
# Foundations



## MDP

- *Trajectories* :  $\tau = (s_1, a_1, r_1, s_2, a_2, r_2, s_3, \dots, s_T)$ 
  - $s_t \equiv$  *server load* and *job size*

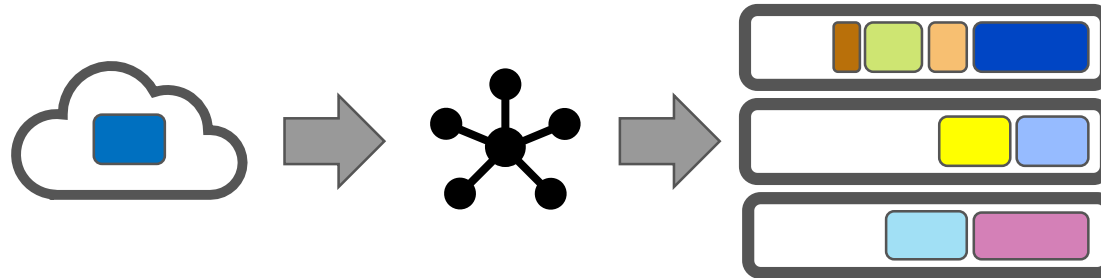
# Foundations



## MDP

- *Trajectories* :  $\tau = (s_1, \mathbf{a}_1, r_1, s_2, \mathbf{a}_2, r_2, s_3, \dots, s_T)$ 
  - $s_t \equiv$  server load and job size
  - $\mathbf{a}_t \equiv$  schedule job to server  $i \in \{1, 2, 3\}$

# Foundations



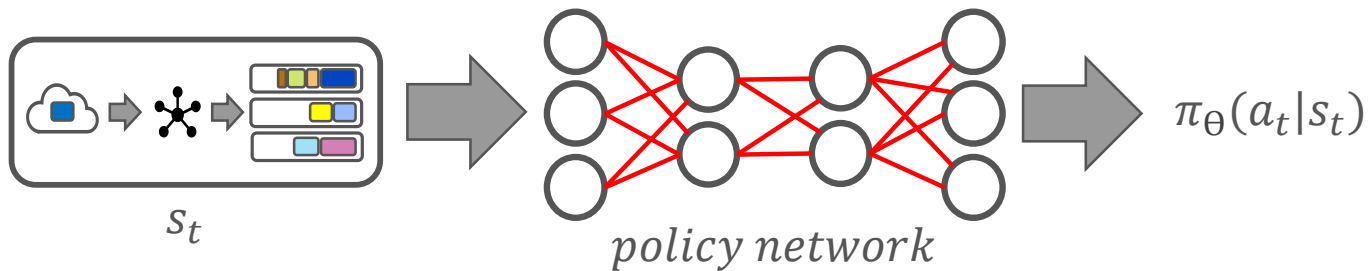
## MDP

- *Trajectories* :  $\tau = (s_1, a_1, \mathbf{r}_1, s_2, a_2, \mathbf{r}_2, s_3, \dots, s_T)$ 
  - $s_t \equiv$  server load and job size
  - $a_t \equiv$  schedule job to server  $i \in \{1, 2, 3\}$
  - $\mathbf{r}_t \equiv (-1) \cdot$  average job completion time



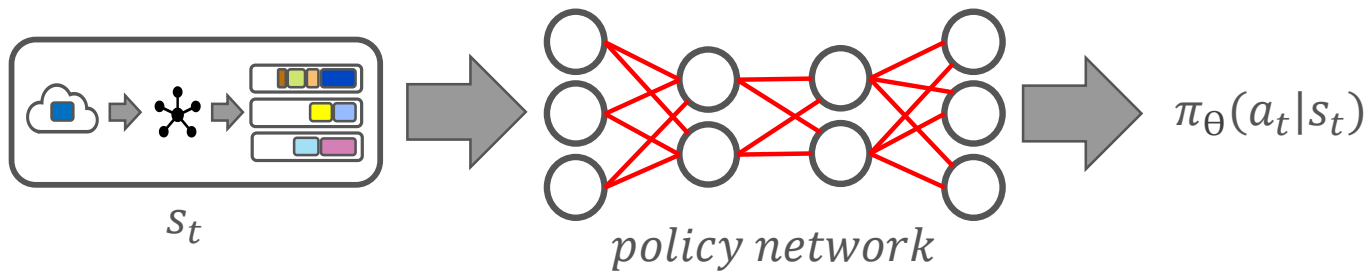
# Policy Gradient Methods

- *Acting policy*  $\pi_{\theta}$  with parameterization  $\theta$

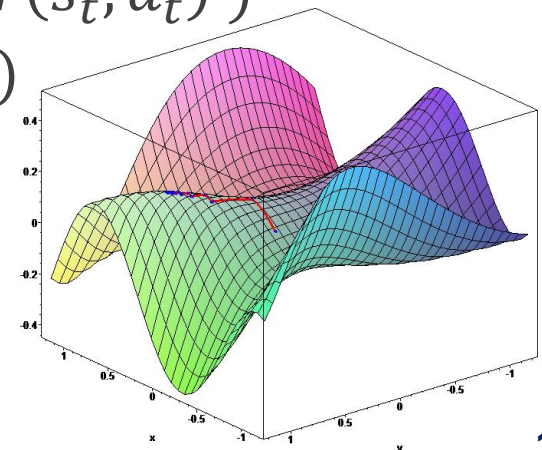


# Policy Gradient Methods

- *Acting policy*  $\pi_{\theta}$  with parameterization  $\theta$



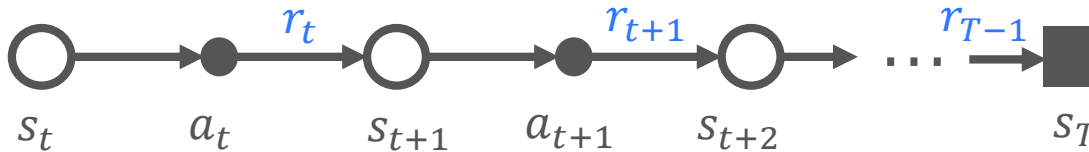
- *Expected return:*  $J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} (\sum_t r(s_t, a_t))$
- $\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi} (Q^{\pi}(s, a) \cdot \nabla_{\theta} \ln \pi_{\theta}(a | s))$



Source: [2]

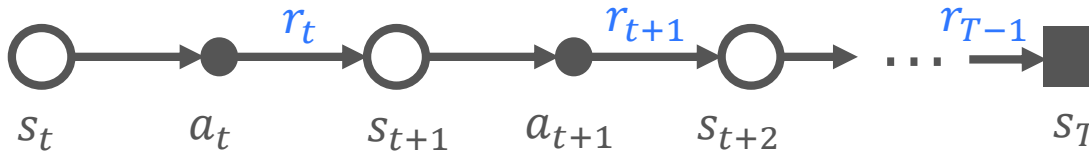
# REINFORCE

- *Updates:*  $\theta_{t+1} = \theta_t + \alpha \cdot G_t \cdot \nabla_{\theta_t} \ln \pi_{\theta} (a_t | s_t)$
- $G_t = r_t + \gamma \cdot r_{t+1} + \dots \gamma^{T-1-t} r_{T-1}$



# REINFORCE

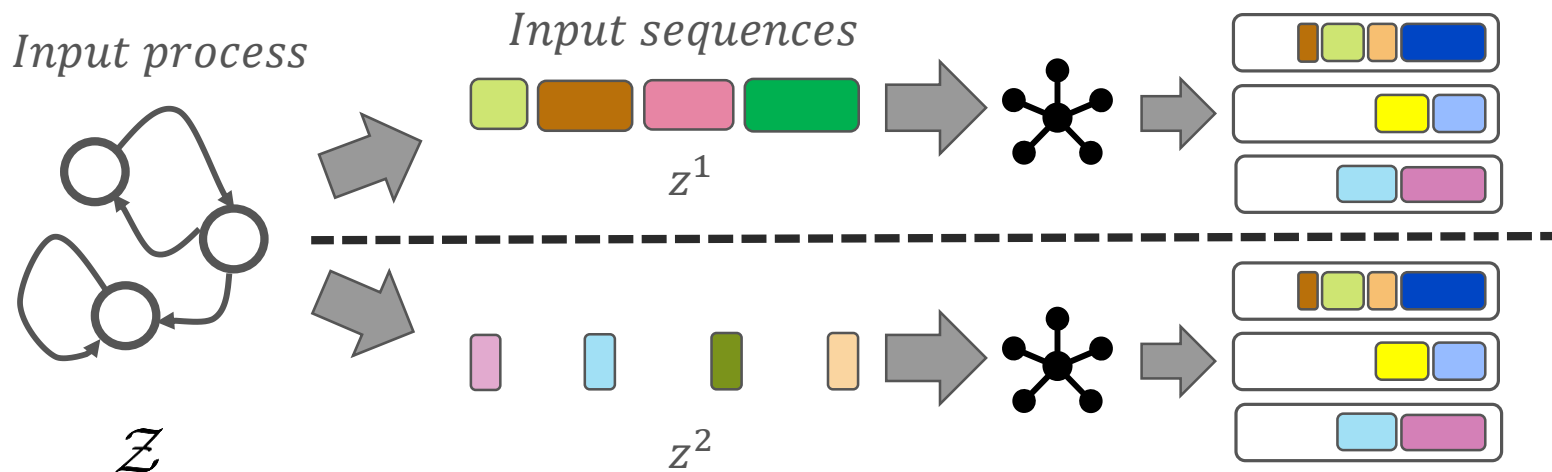
- *Updates:*  $\theta_{t+1} = \theta_t + \alpha \cdot G_t \cdot \nabla_{\theta_t} \ln \pi_{\theta} (a_t | s_t)$
- $G_t = r_t + \gamma \cdot r_{t+1} + \dots \gamma^{T-1-t} r_{T-1}$



## REINFORCE with Baseline

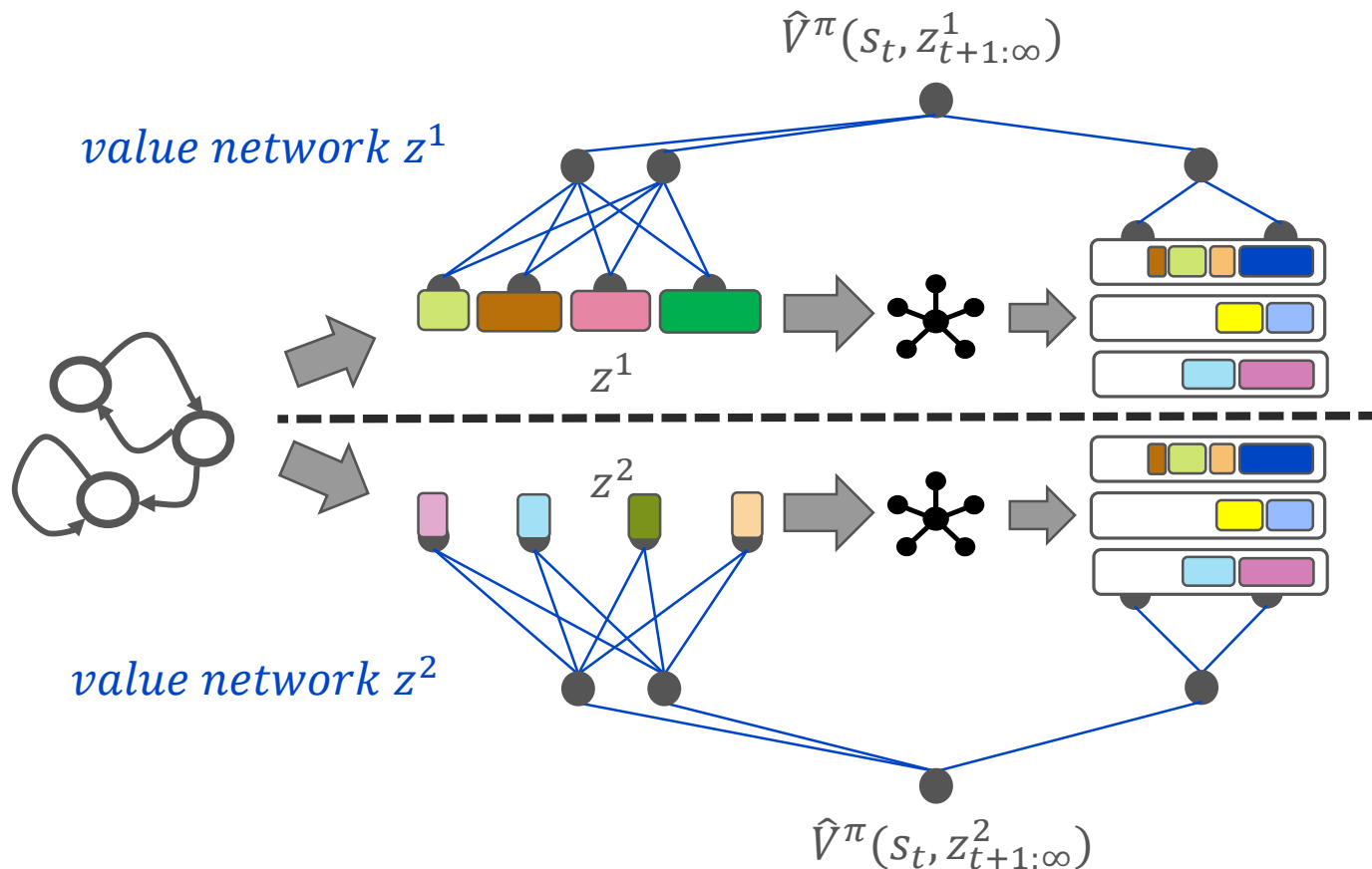
- $\theta_{t+1} = \theta_t + \alpha \cdot (G_t - \hat{V}^{\pi}(s_t)) \cdot \nabla_{\theta_t} \ln \pi_{\theta} (a_t | s_t)$

# Input-Driven Environments



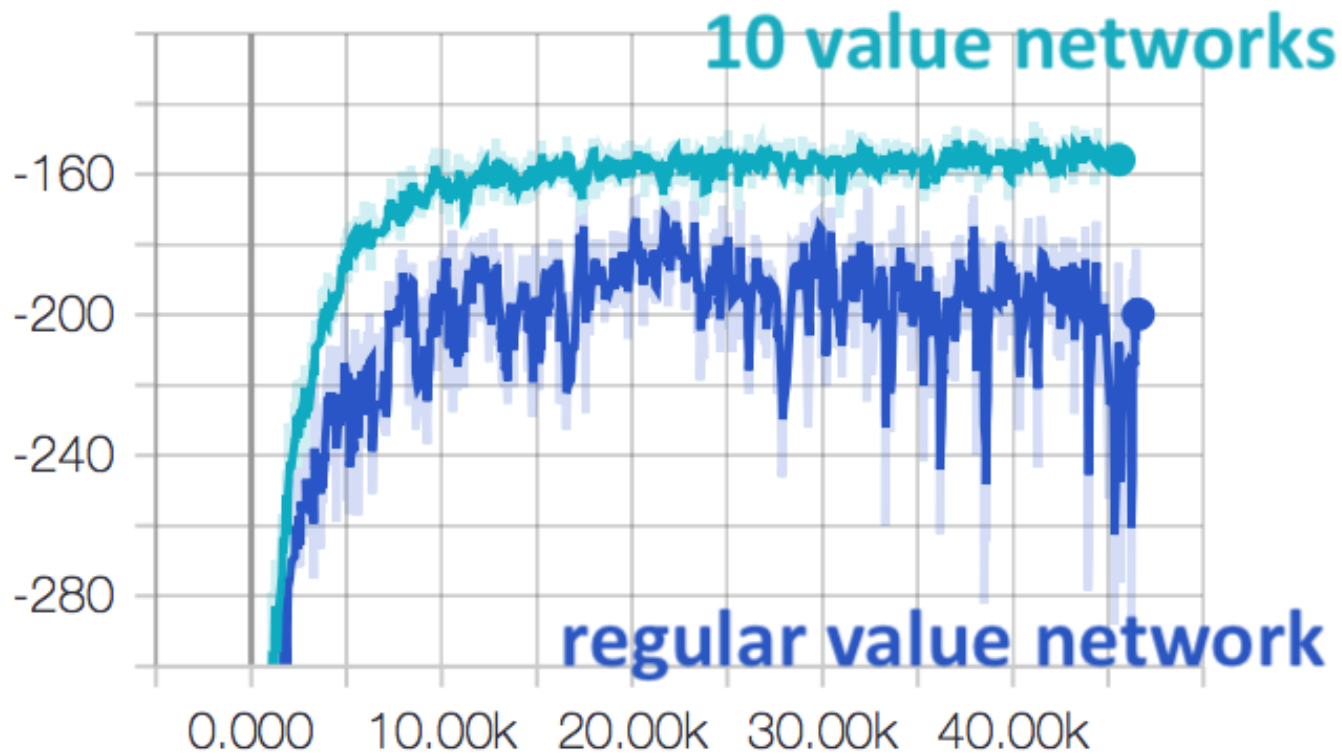
- Rewards may be the result of favorable inputs
- $\hat{V}^{\pi}(s_t, z_{t+1:\infty}^1) \neq \hat{V}^{\pi}(s_t, z_{t+1:\infty}^2)$
- $\hat{V}^{\pi}(s_t)$  does not account for input differences

# Multi-Value Network Baselines



# Experimental Results

Sum\_of\_rewards



Source: [3]

## References

[1] H. Mao, S. B. Venkatakrisnan, M. Schwarzkopf, and M. Alizadeh. “Vari-ance reduction for reinforcement learning in input-driven environments”. In: arXiv preprint arXiv:1807.02264 (2018).

[2] Image source:

[https://upload.wikimedia.org/wikipedia/commons/6/68/Gradient\\_ascent\\_%28surface%29.png](https://upload.wikimedia.org/wikipedia/commons/6/68/Gradient_ascent_%28surface%29.png)

[3] Image source:

[https://github.com/hongzimaoh/input\\_driven\\_rl\\_example/blob/master/figures/training.png](https://github.com/hongzimaoh/input_driven_rl_example/blob/master/figures/training.png)