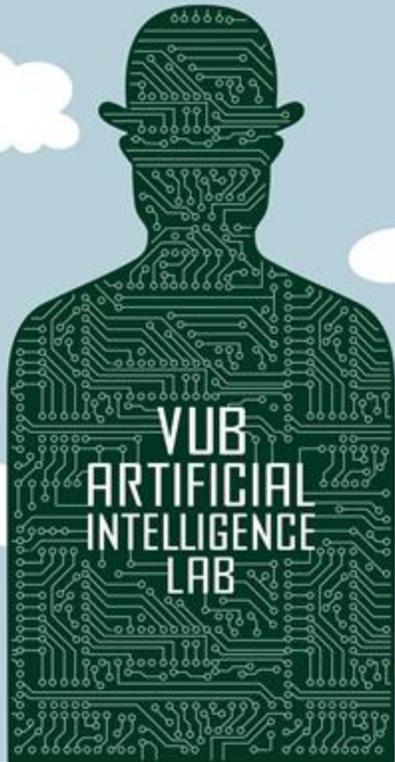


*Ceci n'est pas d'intelligence*



# Continuous BDPI

Denis Steckelmacher

[denis.steckelmacher@vub.be](mailto:denis.steckelmacher@vub.be)

VUB Artificial Intelligence Lab

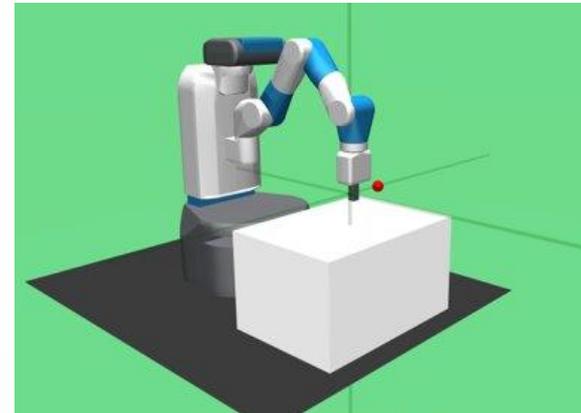
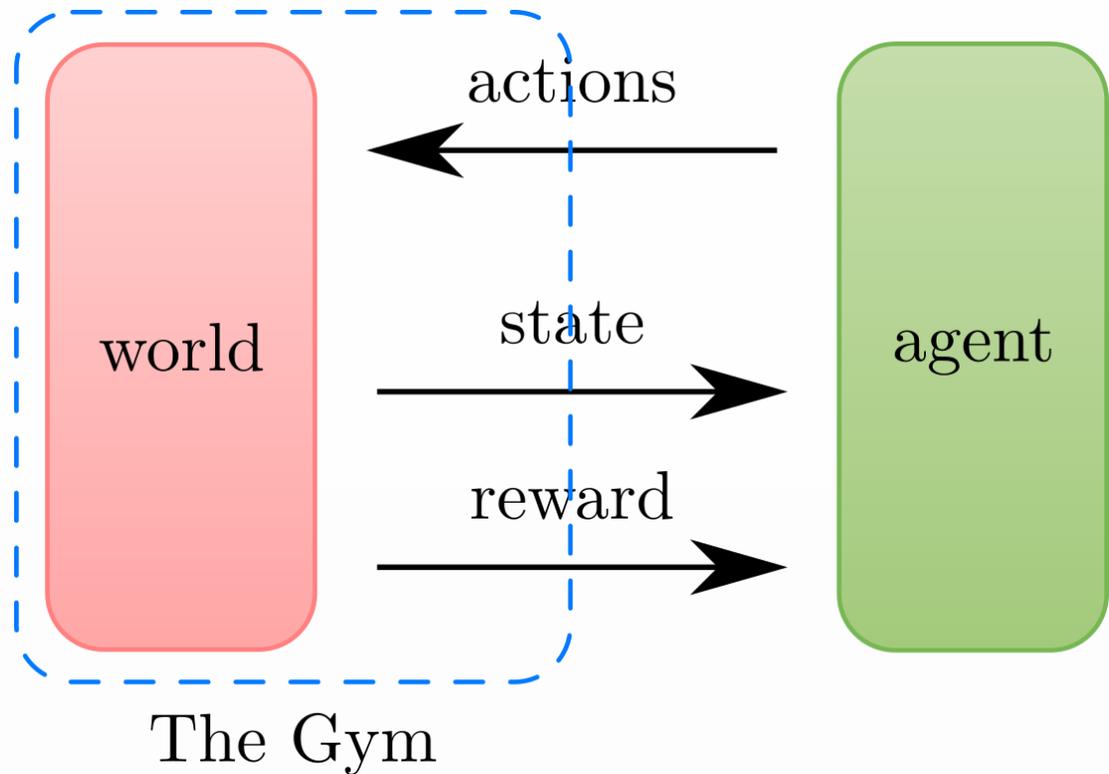


[ai.vub.ac.be](http://ai.vub.ac.be)



[@aibrussels](https://twitter.com/aibrussels)

# Reinforcement Learning



# BDPI

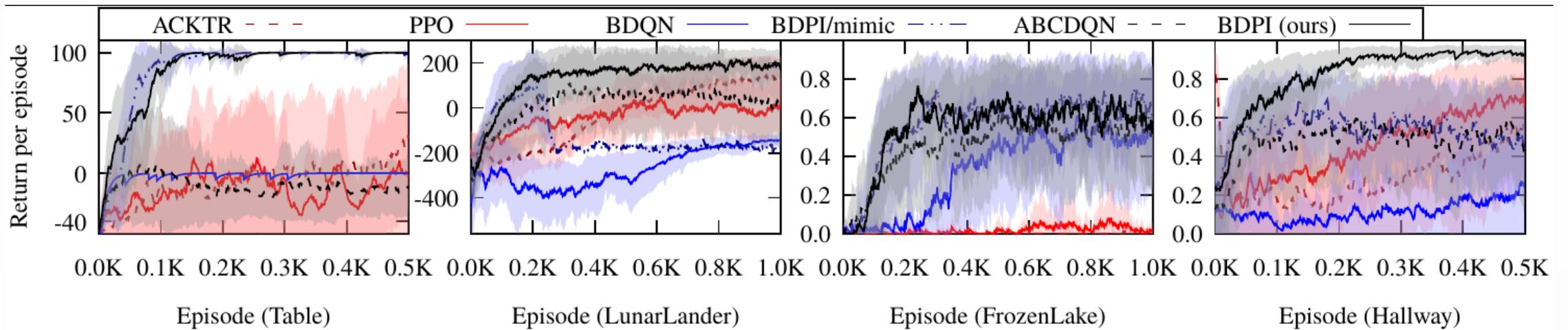
Critic update:

$$Q(s_t, a_t) += \alpha * (r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t))$$

Actor update:

$$\Pi(s_t) \leftarrow \operatorname{argmax}_a Q_i(s_t, a)$$

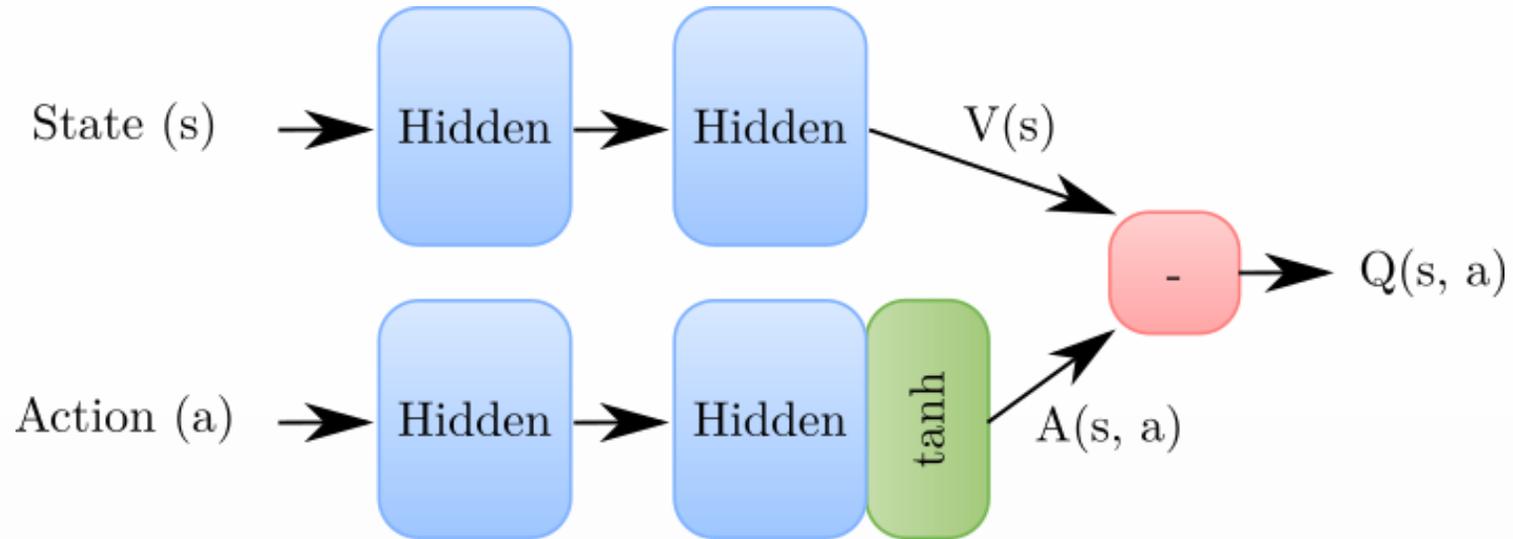
# BDPI is promizing



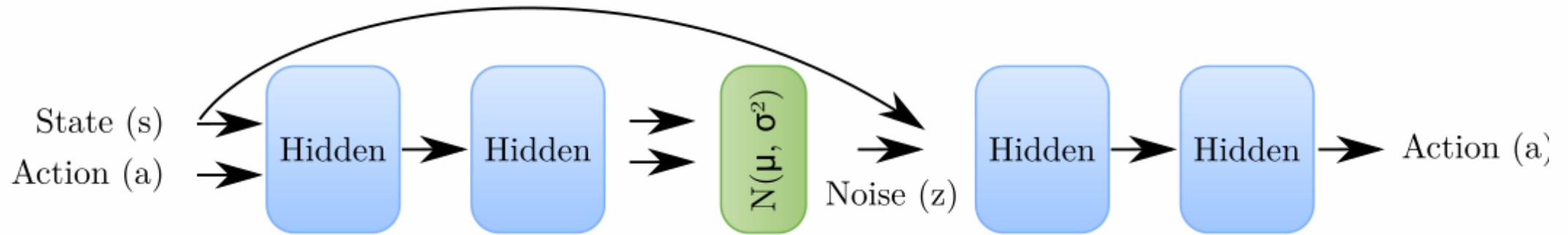
# Continuous BDPI

- Keep the off-policy critics
  - But with continuous actions
- Keep the stochastic actor
  - But take care of the argmax

# The critics (Q-Functions)



# The critics (VAE)

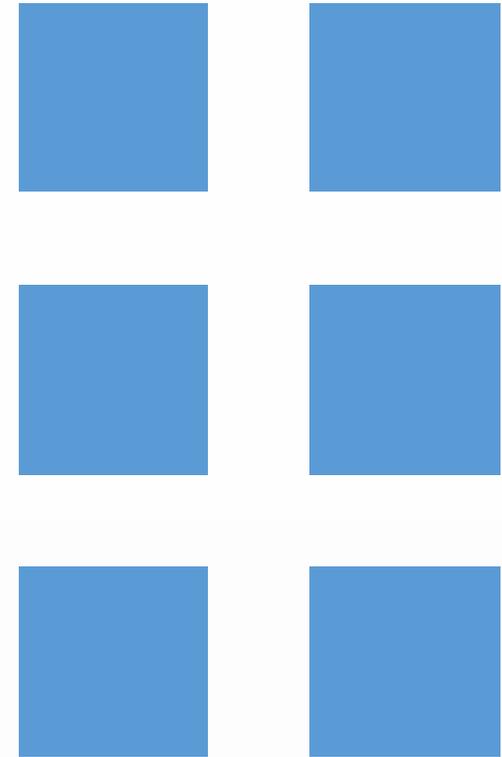
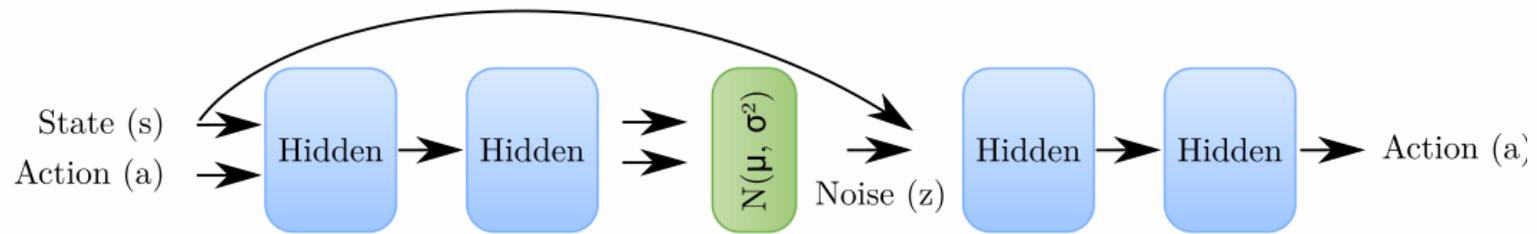


# Continuous BDPI

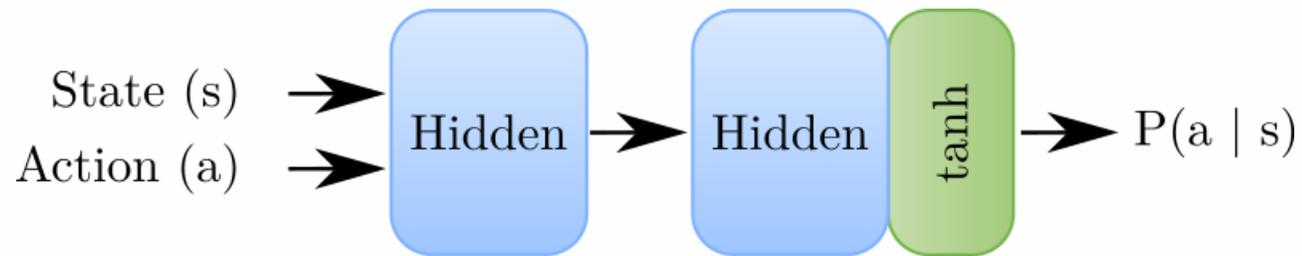
Critic update:

$$Q(s_t, a_t) += \alpha * (r_t + \gamma \max_{a' \sim \text{VAE}(s_{t+1})} Q(s_{t+1}, a') - Q(s_t, a_t))$$

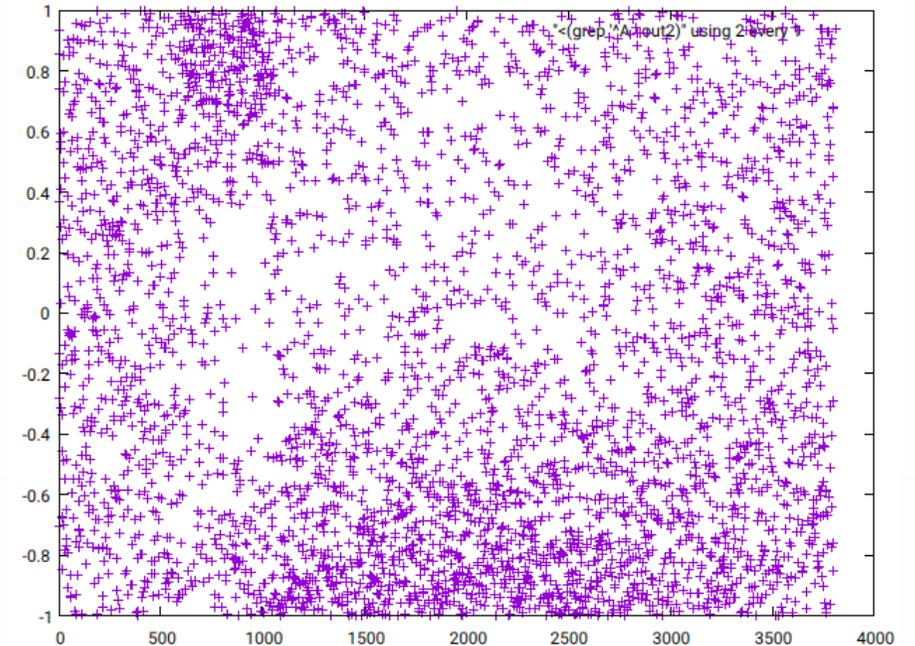
# The actor: a VAE?



# The actor: rejection sampling



- Sample actions uniformly
- For each action a:
  - If  $\text{random}() < P(a | s)$ , accept



# Continuous BDPI

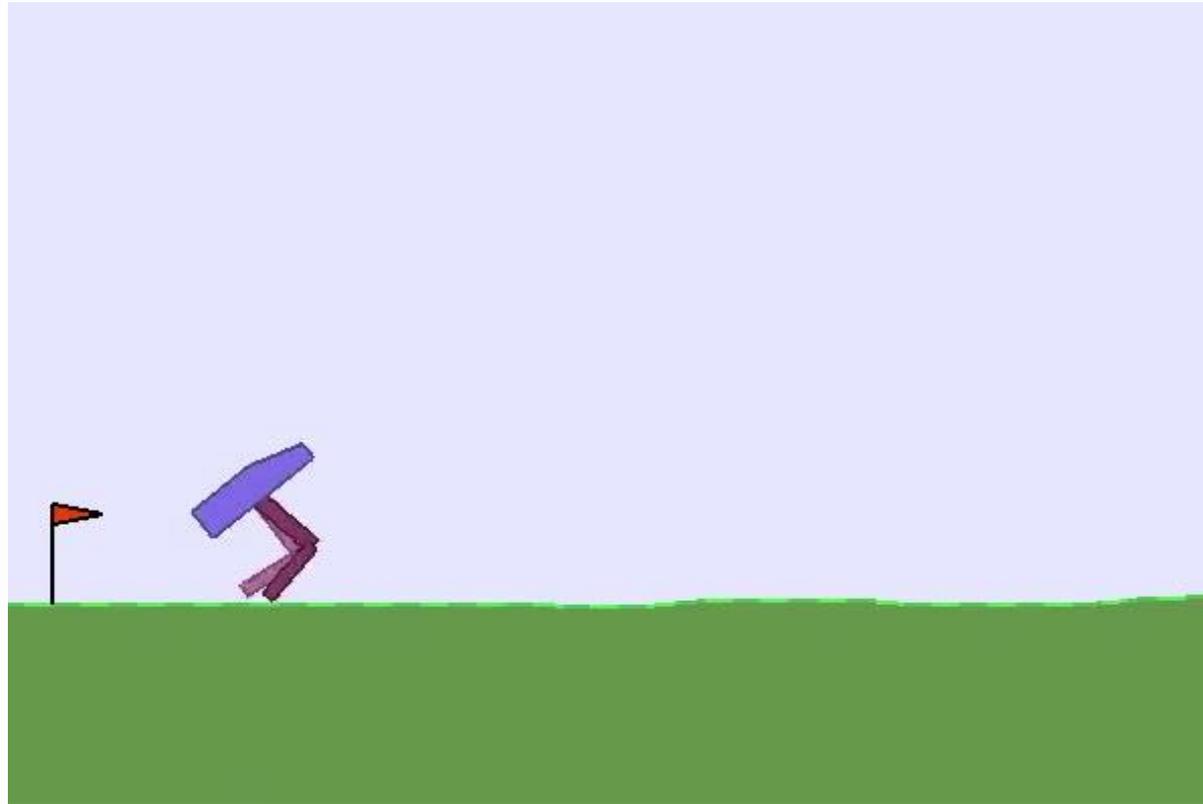
Actor update:

$$\Pi(s_t) \leftarrow \mathbf{argmax}_{a \sim \text{Unif}} Q_i(s_t, a)$$

Implemented as:

$$\begin{aligned} \Pi(s_t, a^*) &+= alr \\ \Pi(s_t, a^- \sim \text{Unif}) &-= alr \end{aligned}$$

# Experiments



# Results

