# Partial Observability in DRL

Part 1: POMDPs, (A)DRQN & DVRL

# Most of the World is only Partial Observable

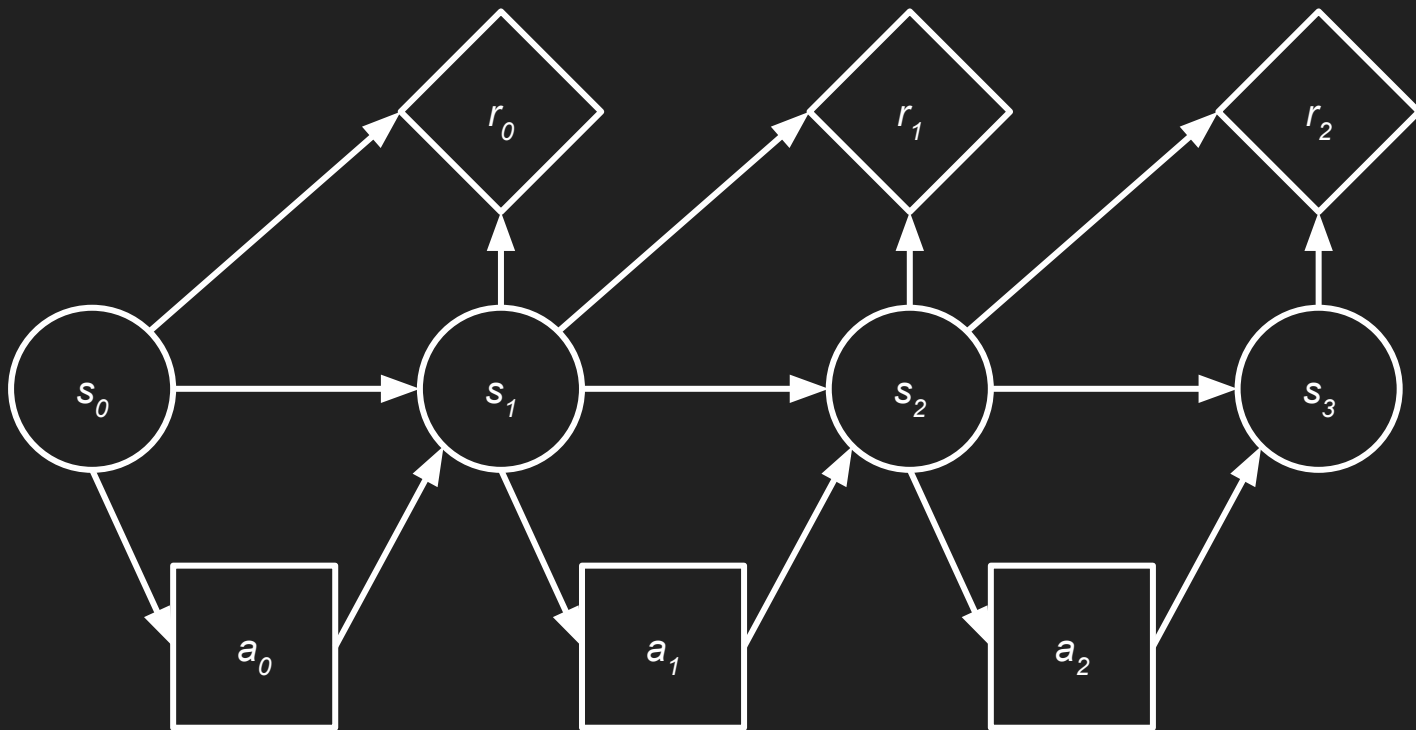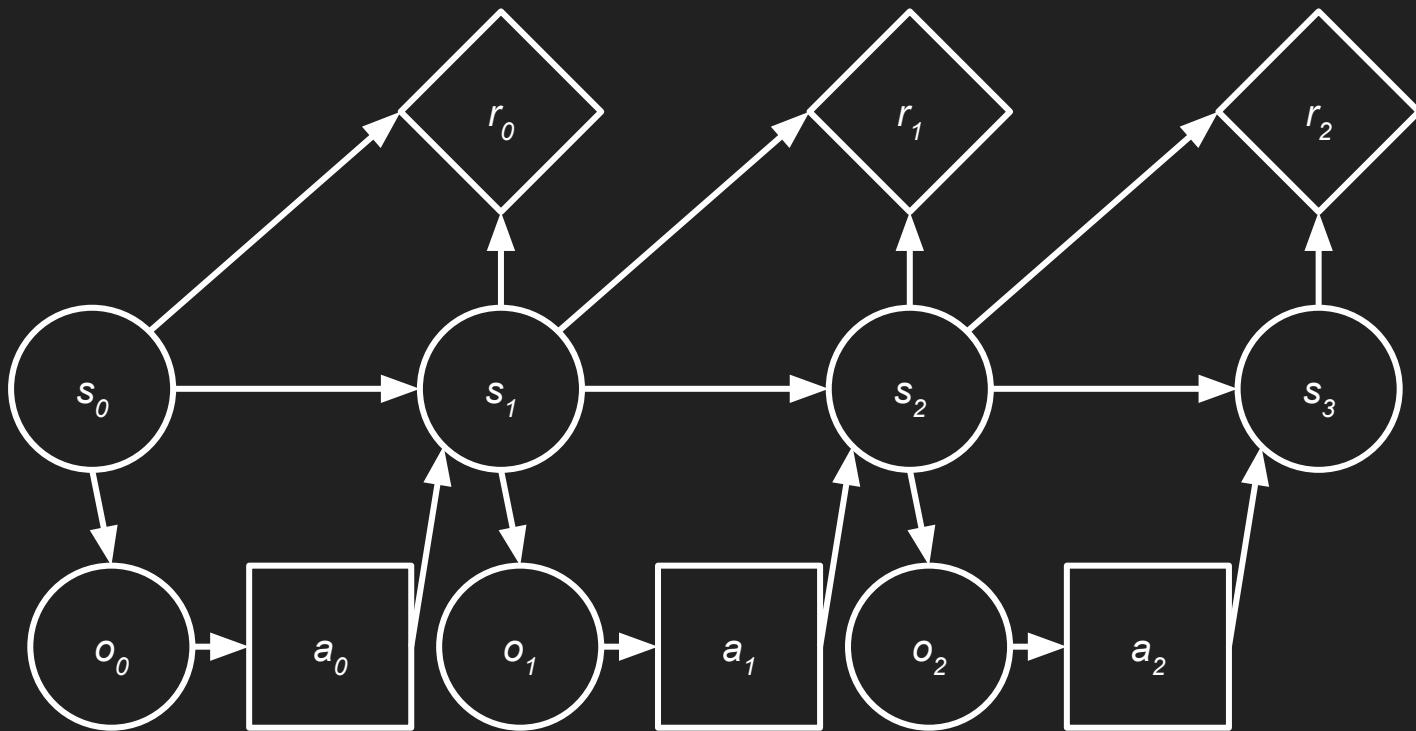Occlusions                    Latent Causes                    Intentions

# From MDP to ...

# From MDP to POMDP

# Slightly more formal

7-Tuple: ($S$, $A$, $T$, $R$, Ω, $O$, $\gamma$)

$s \in S$ is a state from the set of States

$a \in A$ is an action set of Actions

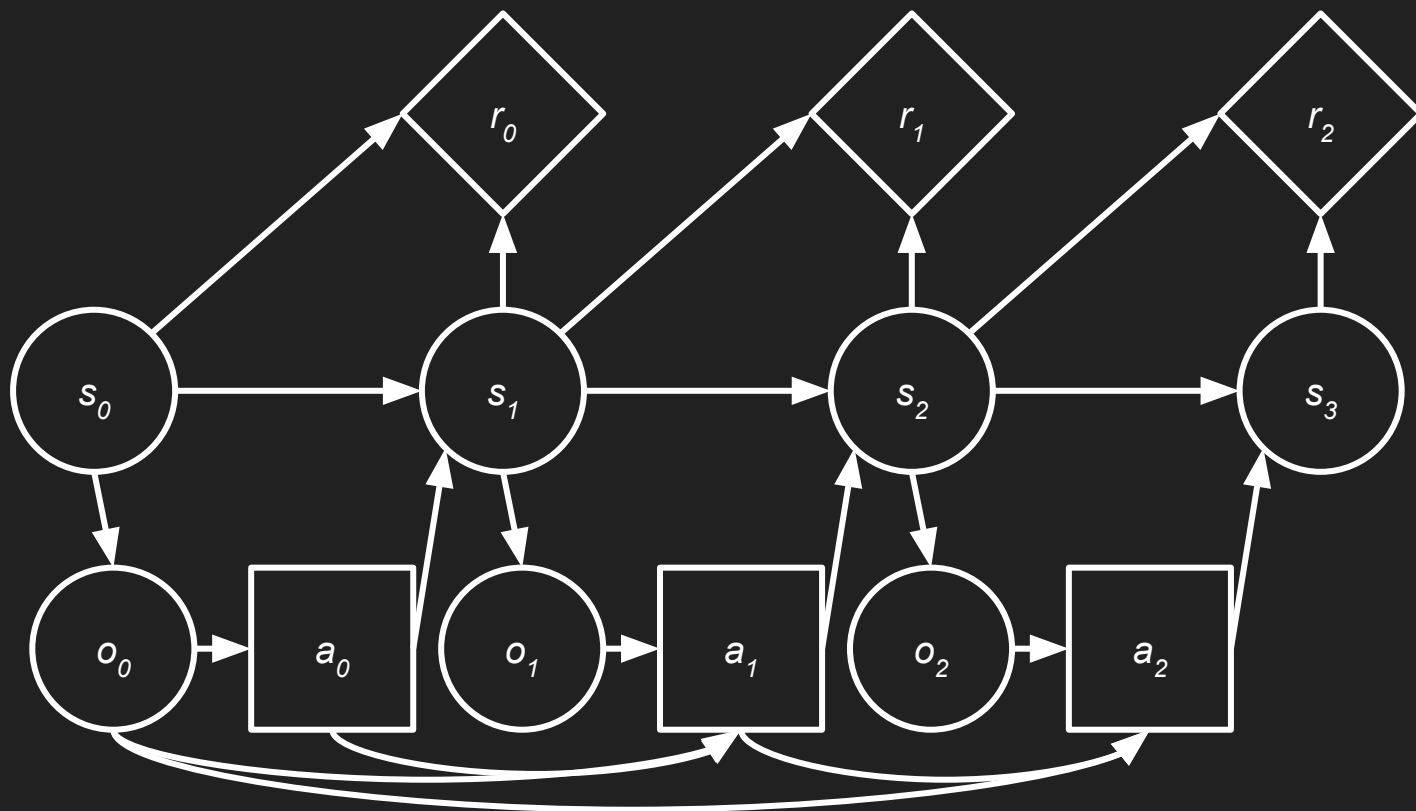$T(s_{t+1} | s_t, a_t)$ is the transition probabilities

$R$: $S \times A \rightarrow \mathbb{R}$, reward function

o $\in$ Ω, an observation from the set of observations

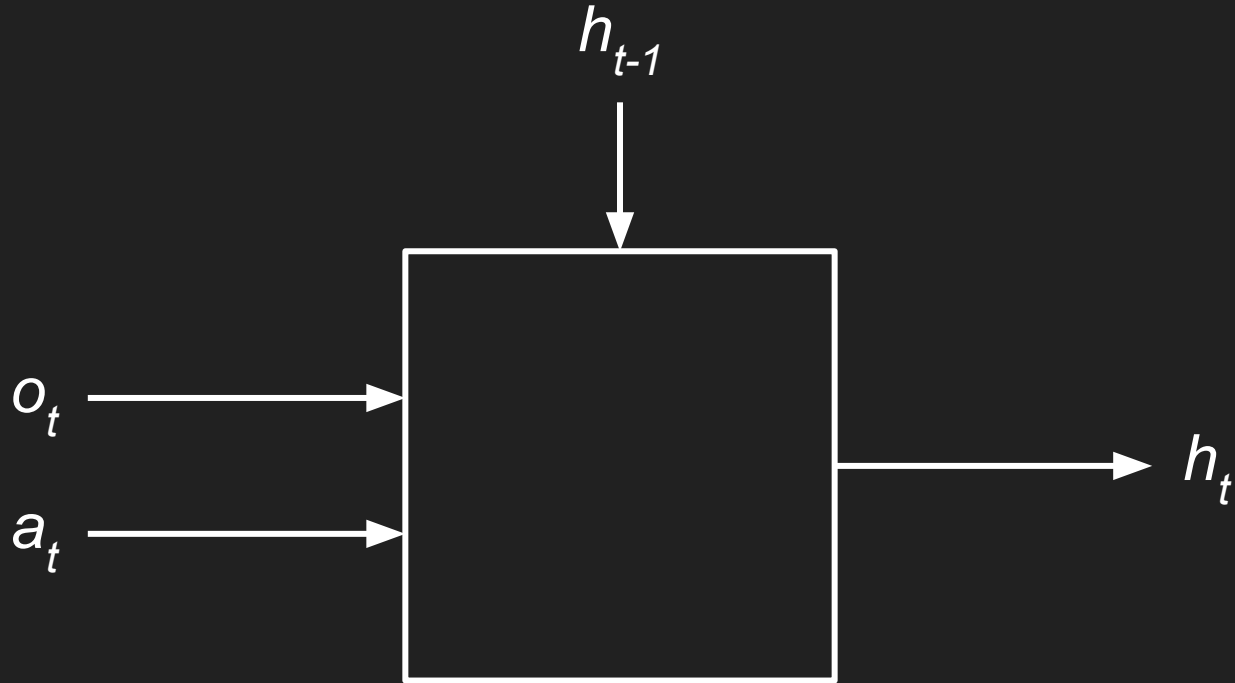$O(o_{t+1} | s_{t+1}, a_t)$ is the conditional observation probabilities

$\gamma \in [0,1]$ is the discount factor

# From MDP to POMDP: A Problem

How to act on *all* past information?

# Option 1: Remember (RNN)
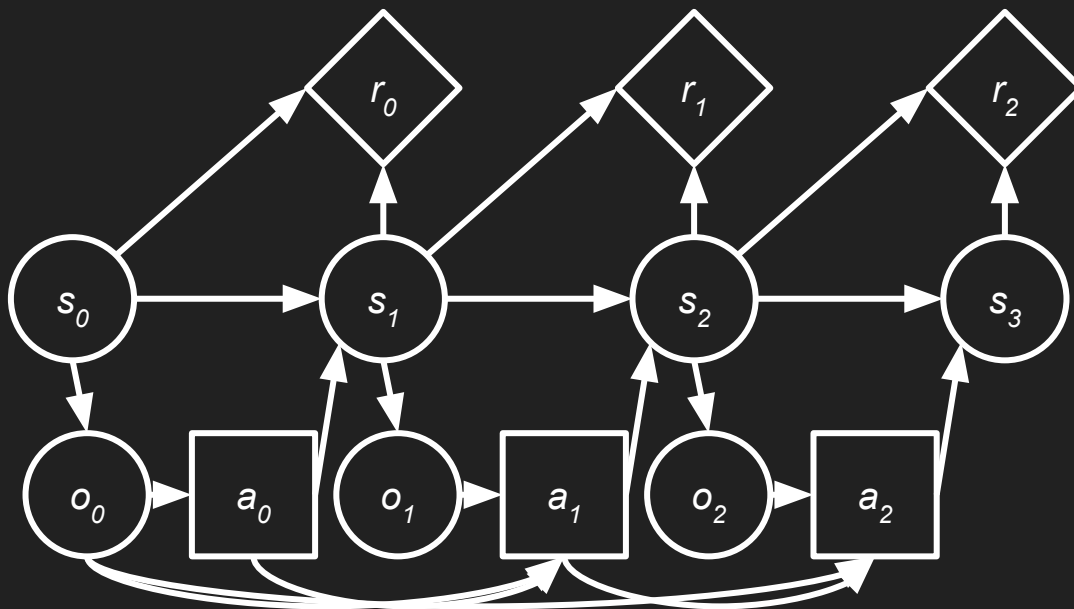
# How to act on all past information?

Option 1: Remember (RNN)

- Generalization can be hard.
- No notion of stochasticity.
- Continuous cases are hard.

# Option 2: Belief

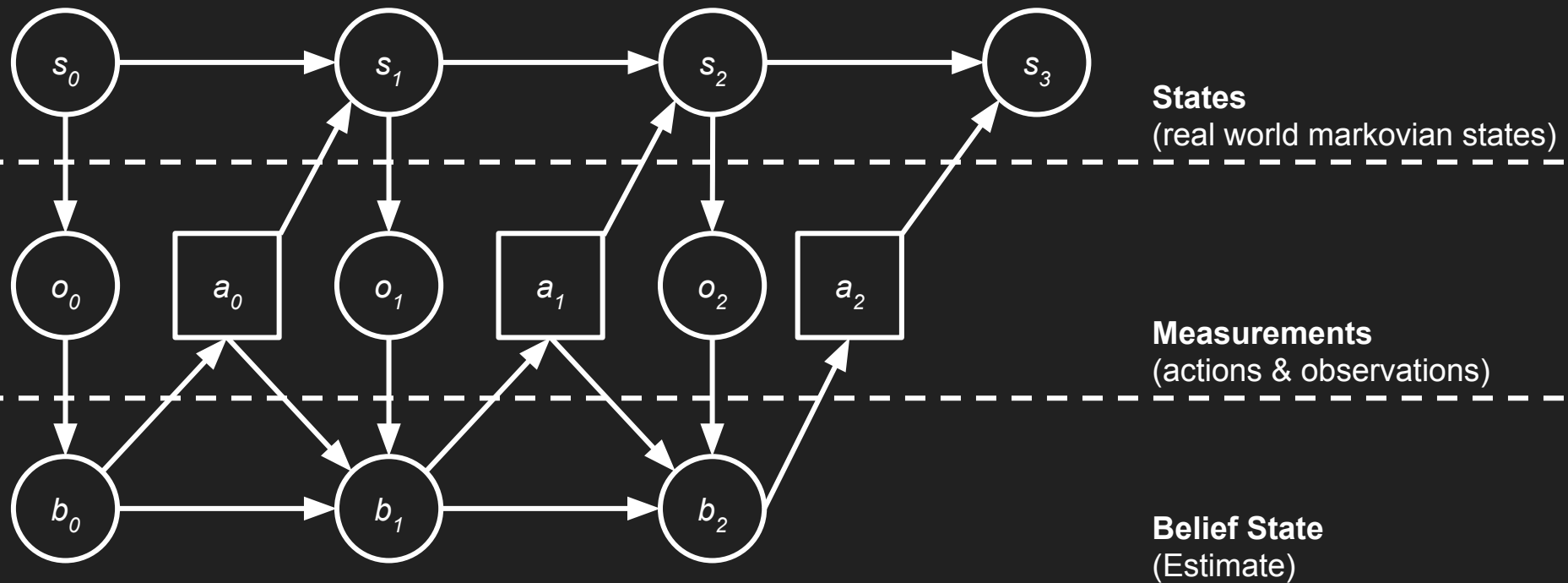$$b_t := p_\theta (s_t \mid o_{\leq t}, a_{\leq t})$$    Belief state

# Option 2: Belief

$$b_t := p_\theta \left( s_t \mid o_{\leq t}, a_{\leq t} \right)$$   Belief state



**States**
(real world markovian states)

**Measurements**
(actions & observations)

**Belief State**
(Estimate)

# Option 2: Belief

$T = p_\theta(s_t | s_{t-1}, a_{t-1})$    Transition Matrix

$O = p_\theta(o_t, | s_t, a_{t-1})$    Observation Matrix

$b_t := p_\theta(s_t | o_{\leq t}, a_{\leq t})$    Belief state

$$b_t(s_t) = \frac{O(o_t | s_t, a_{t-1}) \sum_{s_{t-1} \in S} T(s_t | s_{t-1}, a_{t-1}) b(s_t)}{\text{Normalization Factor}}$$

# How to act on all past information?

**Option 1: Remember (RNN)**

- Generalization can be hard.
- No notion of stochasticity.
- Continuous cases are hard.

**Option 2: Belief**

- Computationally Expensive.
- Requires model.
- Provides stochasticity.
- Tends to generalize.

# Not as clear

**Model free**

RNN
(A)DRQN

**Explicit
Belief tracking**

DVRL

**Implicit
Belief tracking**

Next Session

# Not as clear

**Model free**

RNN
**(A)DRQN**

**Explicit
Belief tracking**
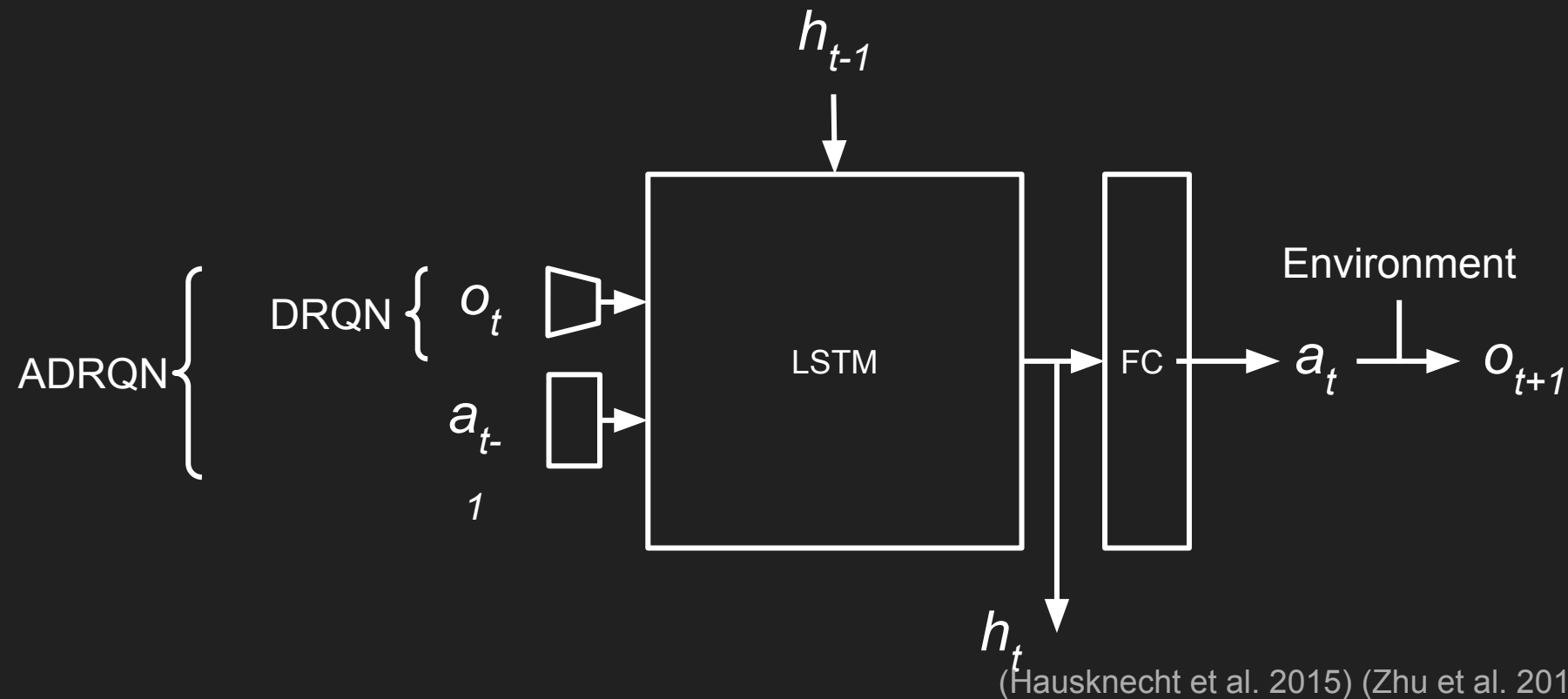
DVRL

**Implicit
Belief tracking**

Next Session

# Deep Q-learning approaches for POMDPs

| Model | Input | Problem Addressed |
|-------|-------|-------------------|
| **DQN** | $s_t$ | model-free POMDP |
| **DBQN** | $b_t$ | Model-based POMDP |
| **DRQN** | $<o_1,o_2,...,o_t>$ | Model-free POMDP |
| **DDRQN** | $<a_0,a_1,...,a_{t-1}>$ <br> $<o_1,o_2,...,o_t>$ | Model-free POMDP |
| **ADRQN** | $<(a_0,o_1),(a_1,o_2),...,(a_{t-1},o_t)>$ | Model-free POMDP |

(Zhu et al. 2017)

# Deep Q-learning approaches for POMDPs

| Model | Input | Problem Addressed |
|-------|-------|-------------------|
| **DQN** | $s_t$ | model-free POMDP |
| **DBQN** | $b_t$ | Model-based POMDP |
| **DRQN** | $<o_1,o_2,...,o_t>$ | Model-free POMDP |
| **DDRQN** | $<a_0,a_1,...,a_{t-1}>$ $<o_1,o_2,...,o_t>$ | Model-free POMDP |
| **ADRQN** | $<(a_0,o_1),(a_1,o_2),...,(a_{t-1},o_t)>$ | Model-free POMDP |

(Zhu et al. 2017)

# (Action-specific) Deep Recurrent Q-Learning: (A)DRQN



(Hausknecht et al. 2015) (Zhu et al. 2017)

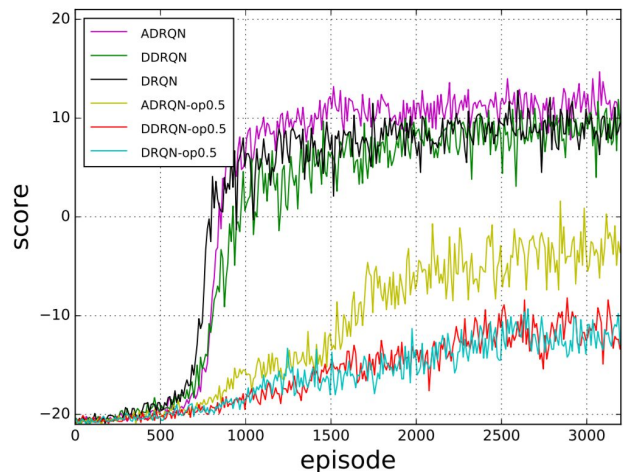# Flickering Frostbite and Pong

# (A)DQRN: Results



Figure 2: Training results for Pong

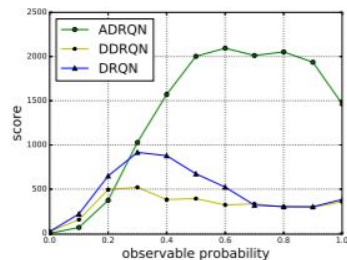

Figure 3: Training results for Frostbite

(Zhu et al. 2017)

# (A)DQRN: Results


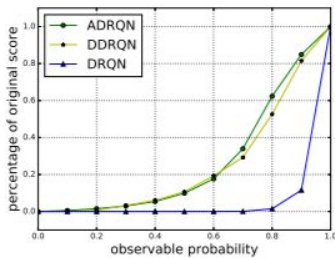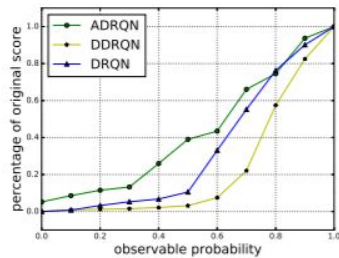
(a) Game Pong  (b) Game Frostbite

Train on POMDP, test on MDP

(a) Game Pong  (b) Game Frostbite

Train on MDP, test on POMDP

(Zhu et al. 2017)

# (A)DQRN: Critique

**Model-free & Blackbox:**

likely to summarize and not generalize

# Deep Variational Reinforcement Learning (DVRL)



(Igl et al. 2018)

$\hat{b}_{t-1} = (h_{t-1}^l, z_{t-1}^l, w_{t-1}^l)_{l=1}^K$

$h_{t-1}^{u_{t-1}^k}$

$\psi_\theta^{\mathrm{RNN}}(h_{t-1}^{u_{t-1}^k}, z_t^k, a_{t-1}, o_t)$

$h_t^k$

$p_\theta(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1})$

$w_t^k$

$q_\phi(z_t^k | h_{t-1}^{u_{t-1}^k}, a_{t-1}, o_t)$

$z_t^k$

$o_t$

$p_\theta(o_t | h_{t-1}^{u_{t-1}^k}, z_t^k, a_{t-1})$

$a_{t-1}$

$K$

$(h_t^k, z_t^k, w_t^k)_{k=1}^K = \hat{b}_t$

RNN

$\hat{h}_t$

$V_\eta(\hat{h}_t)$

$\pi_\rho(a_t | \hat{h}_t)$

sample
aggregate
eval under distribution
pass argument

(Igl et al. 2018)

# Deep Variational Reinforcement Learning (DVRL)



(Igl et al. 2018)

# Deep Variational Reinforcement Learning (DVRL)



(Igl et al. 2018)

# Deep Variational Reinforcement Learning (DVRL)



(Igl et al. 2018)

# Deep Variational Reinforcement Learning (DVRL)

Model

Particle Filter

$o_{t-1}$

$a_t$

V

$\pi$

$\mathcal{L}_{\text{ELB}}$

O

+

$\mathcal{L}_{\text{A2C}}$

(Igl et al. 2018)

# Brief note on notation

*$a_t$ = action at time t*

*$o_t$ = observation at time t*

*k in [1,K] = number of particles*

*$b_t$ = ( $h_t$, $z_t$, $w_t$) belief at time t*
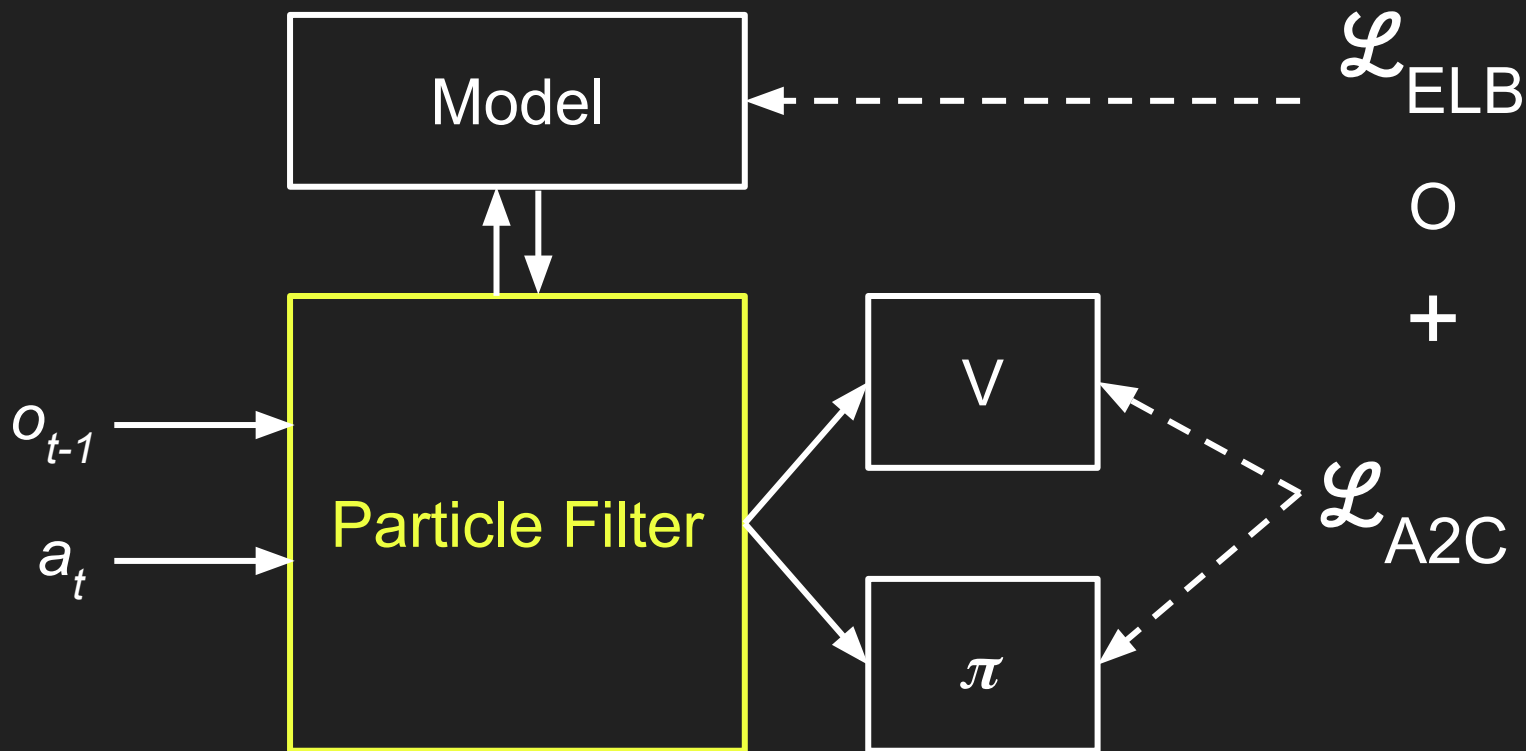
*$z_t$ = an additional stochastic latent state*

Latent Summary of state

*$h_t$ = latent state of a RNN (in a particle)*

*$w_t$ = importance weight of a particle.*

Likelihood of that latent state

# Deep Variational Reinforcement Learning (DVRL)



(Igl et al. 2018)

# DVRL: Particle Filter - Approximating $b_t$

**Previous Belief**         **Sample new values**         **re-weight**         **resample**

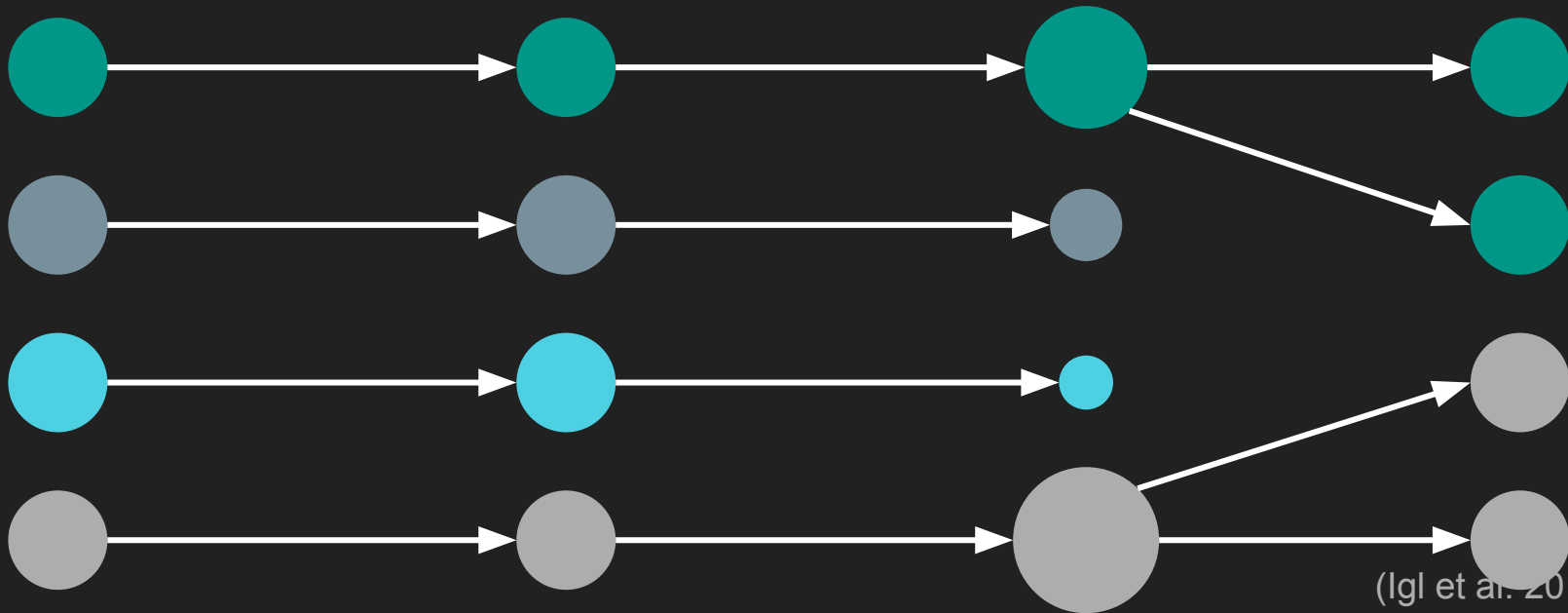$$b_{t-1} = \left(h_{t-1}^k, z_{t-1}^k, w_{t-1}^k\right)_{k=1}^K \quad z_t \sim q_\phi(z_t | h_{t-1}, a_{t-1}, o_t)$$

$$h_t = \psi_\theta^{RNN}(z_t, h_{t-1}, a_{t-1}, o_t)$$

$$w_t = \frac{p_\theta\left(z_t | h_{t-1}, a_{t-1}\right) p_\theta\left(o_t | h_{t-1}, z_t, a_{t-1}\right)}{q_\phi\left(z_t | h_{t-1}, a_{t-1}, o_t\right)} \quad b_t = \left(h_t^k, z_t^k, w_t^k\right)_{k=1}^K$$



(Igl et al. 2018)

# Deep Variational Reinforcement Learning (DVRL)



(Igl et al. 2018)

# DVRL: Policy - Summarize the particles

$b_t$        Summary        Value & Policy

$(h_t^1, z_t^1, w_t^1)$ → RNN

$(h_t^2, z_t^2, w_t^2)$ → RNN

$(h_t^i, z_t^i, w_t^i)$ → RNN

$(h_t^k, z_t^k, w_t^k)$ → RNN

$\hat{h}_t$

FC   $V(\hat{h}_t)$

FC   $\pi(a_t|\hat{h}_t)$

(Igl et al. 2018)

# Deep Variational Reinforcement Learning (DVRL)



(Igl et al. 2018)

# DVRL: Model

$$w_t = \frac{p_\theta\left(z_t|h_{t-1}, a_{t-1}\right) \boxed{p_\theta\left(o_t|h_{t-1}, z_t, a_{t-1}\right)}}{q_\phi\left(z_t|h_{t-1}, a_{t-1}, o_t\right)}$$

$$p_\theta(o_t|h_{t-1}, z_t, a_{t-1})$$

(Igl et al. 2018)

# DVRL: Model

$$q_\phi\left(z_t|h_{t-1}, a_{t-1}, o_t\right) \qquad p_\theta(o_t|h_{t-1}, z_t, a_{t-1})$$

(Igl et al. 2018)

# DVRL: Model

Encoder      Decoder



O   $q$   Z   $p$   O

$$q_\phi\left(z_t \mid h_{t-1}, a_{t-1}, o_t\right)$$

$$p_\theta(o_t \mid h_{t-1}, z_t, a_{t-1})$$

(Igl et al. 2018)

# DVRL: Model



$$w_t = \frac{p_\theta\left(z_t | h_{t-1}, a_{t-1}\right) p_\theta\left(o_t | h_{t-1}, z_t, a_{t-1}\right)}{q_\phi\left(z_t | h_{t-1}, a_{t-1}, o_t\right)}$$
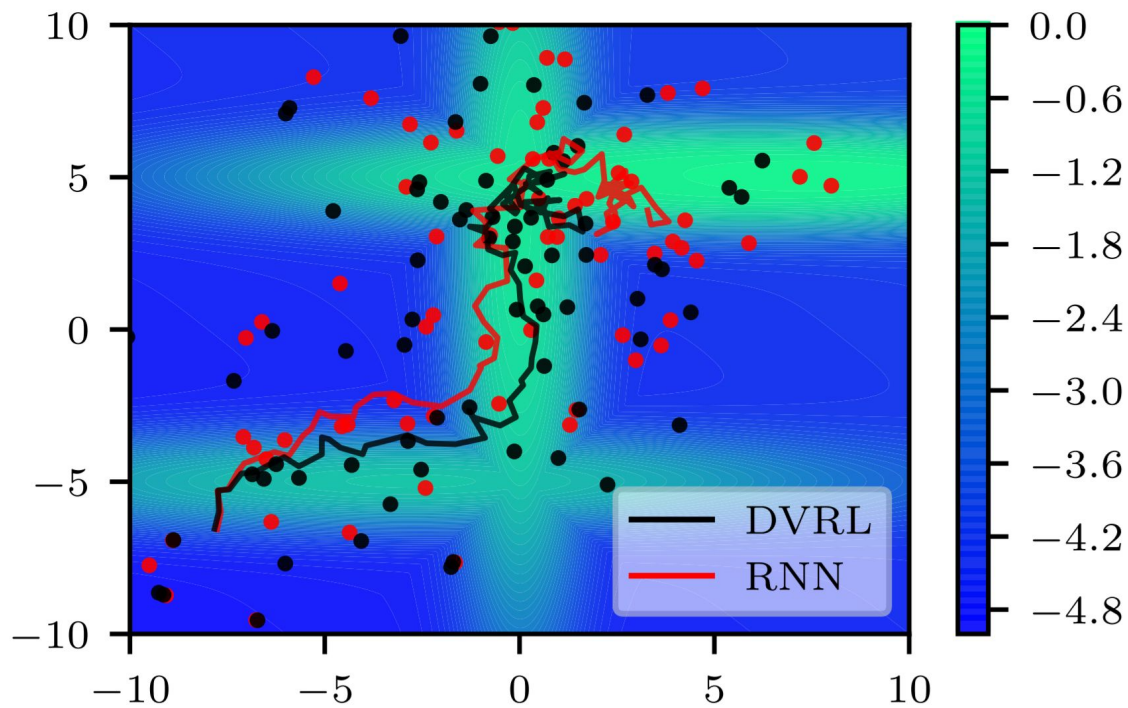
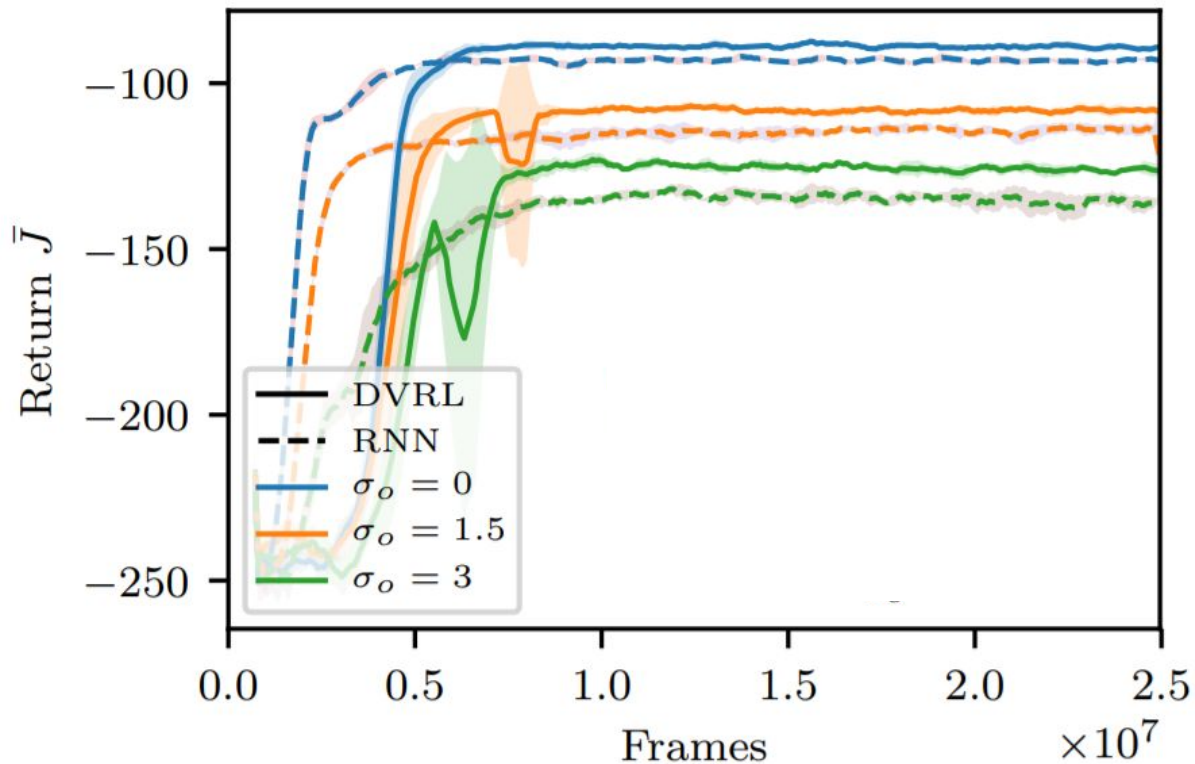$$ELBO(\theta, \phi) \approx \sum_{t=1}^{T} \log\left(\frac{1}{K}\sum_{k=1}^{K} w_t^k\right)$$

(Igl et al. 2018)

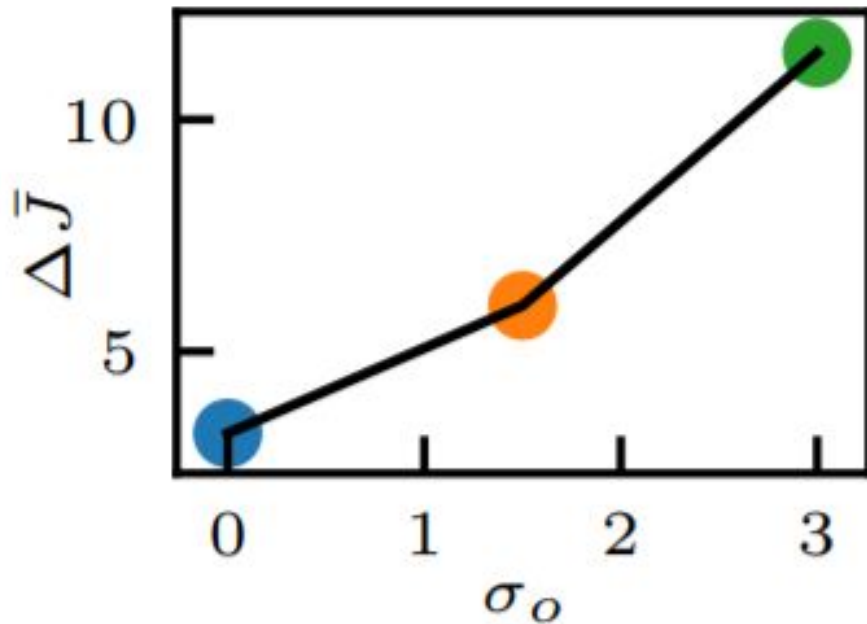# DVRL: Joint Learning



(Igl et al. 2018)

# DVRL: Results - noisy MountainHike



(Igl et al. 2018)

# DVRL: Results - noisy MountainHike



(Igl et al. 2018)

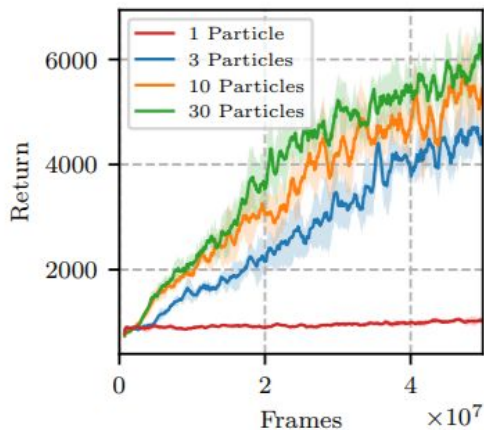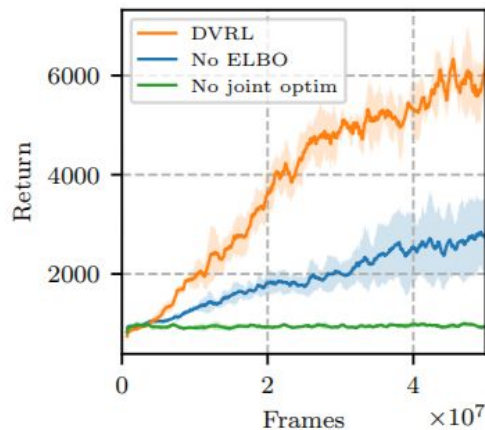# DVRL: Results - noisy MountainHike



(Igl et al. 2018)

# ChopperCommand

# Results: Ablation on Atari



**(a)** Influence of the particle number on performance for DVRL. Only using one particle is not sufficient to encode enough information in the latent state.

**(b)** Performance of the full DVRL algorithm compared to setting $\lambda^E = 0$ ("No ELBO") or not backpropagating the policy gradients through the encoder ("No joint optim").

(Igl et al. 2018)

# DVRL: Critique

**The belief state is still a rough approximation.**

Is this really the best way to learn it?

# Summary

- Extended MDP to POMDP
- (A)DRQN
- DVRL

# Discussion

In a POMDP we still assume full access to the reward.

1) This not a realistic case (our perception of the reward depends as much on our observations as the state)

2) If it is realistic, our belief should be updated based on the reward.

# Next

| Model free | Explicit Belief tracking | Implicit Belief tracking |
|:---:|:---:|:---:|
| RNN (A)DRQN | DVRL DPFRL | VRM |

# References

Hausknecht, M., & Stone, P. (2015, September). *Deep recurrent q-learning for partially observable mdps*. In 2015 AAAI Fall Symposium Series.

Igl, M., Zintgraf, L., Le, T. A., Wood, F., & Whiteson, S. (2018). *Deep variational reinforcement learning for pomdps*. arXiv preprint arXiv:1806.02426.

Zhu, P., Li, X., Poupart, P., & Miao, G. (2017). *On improving deep reinforcement learning for pomdps.* arXiv preprint arXiv:1704.07978.