# UNIVERSITY OF ALBERTA

# Importance Resampling for Off-policy Prediction
## Matthew Schlegel, Wes Chung, Jian Qian, Daniel Graves, Martha White

RLAI  amii  HUAWEI

## Motivation

- Learning many off-policy predictions through general value functions

- Find an approach which corrects for the action distribution and budgets updates more efficiently.

## Background

**Value Function:**

$$V(s) \doteq \mathbb{E}_\pi[G_t \,|\, S_t = s, A \sim \pi]$$

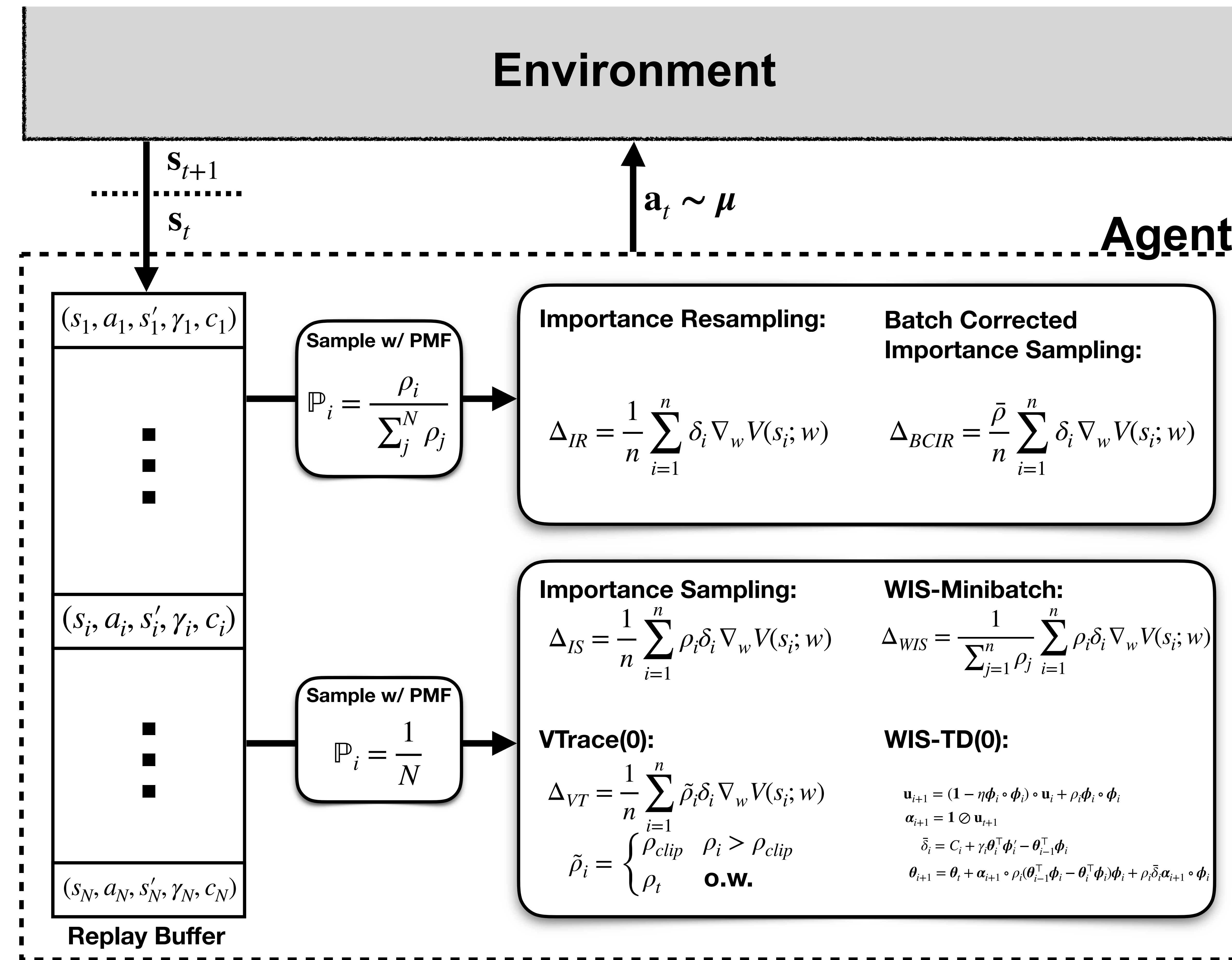$$G_t = \sum_{i=t}^{\infty} \left( \prod_{j=t+1}^{i} \gamma_j \right) C_{i+1}$$

**Off-policy Prediction:**

Behaviour: $\mu(a\,|\,s) : A \mapsto \mathbb{R}$

Target:   $\pi(a\,|\,s) : A \mapsto \mathbb{R}$

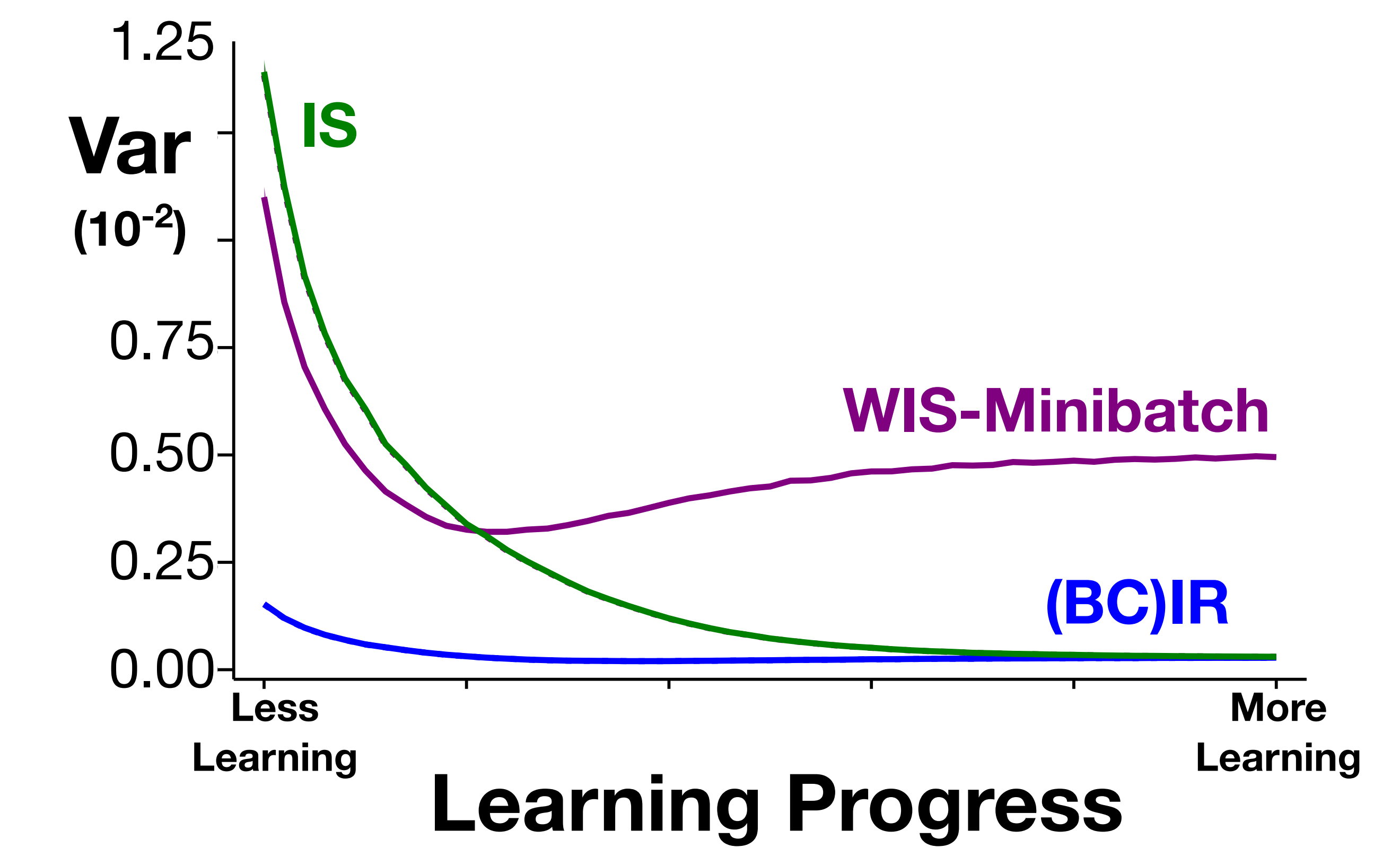$$\Delta_{IS} = \mathbb{E}_\mu \left[ \rho(A, S) \delta \nabla_w V(S) \right]$$

$$\rho_i = \rho(a_i, s_i) = \frac{\pi(a_i\,|\,s_i)}{\mu(a_i\,|\,s_i)}$$

**Environment**

$\mathbf{s}_{t+1}$
$\mathbf{s}_t$
$\mathbf{a}_t \sim \boldsymbol{\mu}$

**Agent**

$(s_1, a_1, s_1', \gamma_1, c_1)$

$(s_i, a_i, s_i', \gamma_i, c_i)$

$(s_N, a_N, s_N', \gamma_N, c_N)$

**Replay Buffer**

Sample w/ PMF

$$\mathbb{P}_i = \frac{\rho_i}{\sum_j^N \rho_j}$$

Sample w/ PMF

$$\mathbb{P}_i = \frac{1}{N}$$

**Importance Resampling:**

$$\Delta_{IR} = \frac{1}{n} \sum_{i=1}^{n} \delta_i \nabla_w V(s_i; w)$$

**Batch Corrected Importance Sampling:**

$$\Delta_{BCIR} = \frac{\bar{\rho}}{n} \sum_{i=1}^{n} \delta_i \nabla_w V(s_i; w)$$

**Importance Sampling:**

$$\Delta_{IS} = \frac{1}{n} \sum_{i=1}^{n} \rho_i \delta_i \nabla_w V(s_i; w)$$

**WIS-Minibatch:**

$$\Delta_{WIS} = \frac{1}{\sum_{j=1}^{n} \rho_j} \sum_{i=1}^{n} \rho_i \delta_i \nabla_w V(s_i; w)$$

**VTrace(0):**

$$\Delta_{VT} = \frac{1}{n} \sum_{i=1}^{n} \tilde{\rho}_i \delta_i \nabla_w V(s_i; w)$$

$$\tilde{\rho}_i = \begin{cases} \rho_{clip} & \rho_i > \rho_{clip} \\ \rho_t & \text{o.w.} \end{cases}$$

**WIS-TD(0):**

$$\mathbf{u}_{i+1} = (\mathbf{1} - \eta \phi_i \circ \phi_i) \circ \mathbf{u}_i + \rho_i \phi_i \circ \phi_i$$
$$\alpha_{i+1} = \mathbf{1} \oslash \mathbf{u}_{i+1}$$
$$\tilde{\delta}_i = C_i + \gamma_i \theta_i^\top \phi_i' - \theta_{i-1}^\top \phi_i$$
$$\theta_{i+1} = \theta_i + \alpha_{i+1} \circ \rho_i (\theta_{i-1}^\top \phi_i - \theta_i^\top \phi_i) \phi_i + \rho_i \tilde{\delta}_i \alpha_{i+1} \circ \phi_i$$

## Theoretical Results

- We show BCIR is consistent and unbiased in both the static buffer and moving buffer cases.

- We provide several instances where we expect the variance of IR to be less than that of IS.

## Hypothesized Effects

- Update variance reduction compared to importance sampling.

- Faster convergence with respect to the number of updates required.

## Markov Chain



## Continuous Four Rooms