



# Low-Variance and Zero-Variance Baselines in Extensive-Form Games

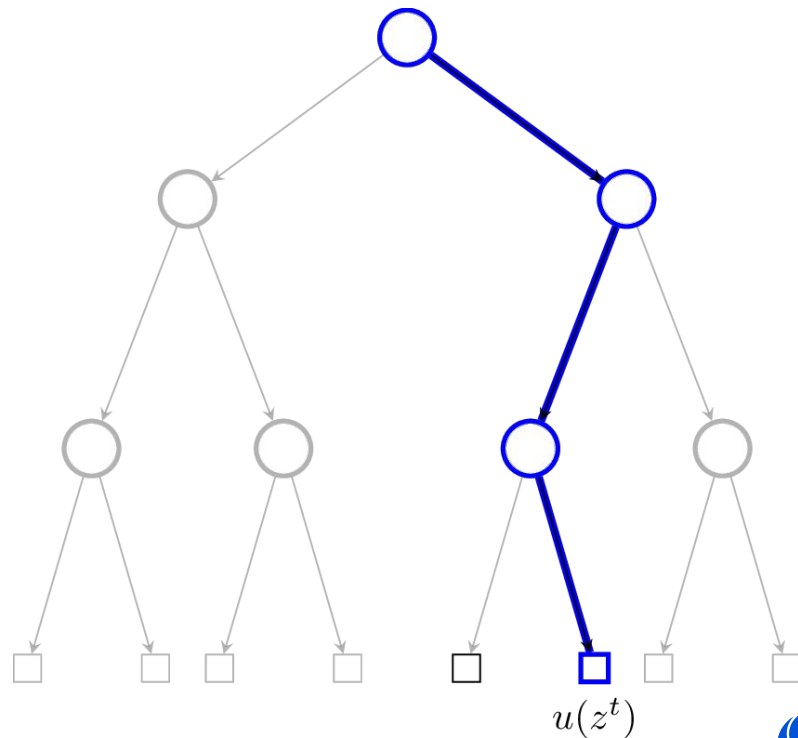
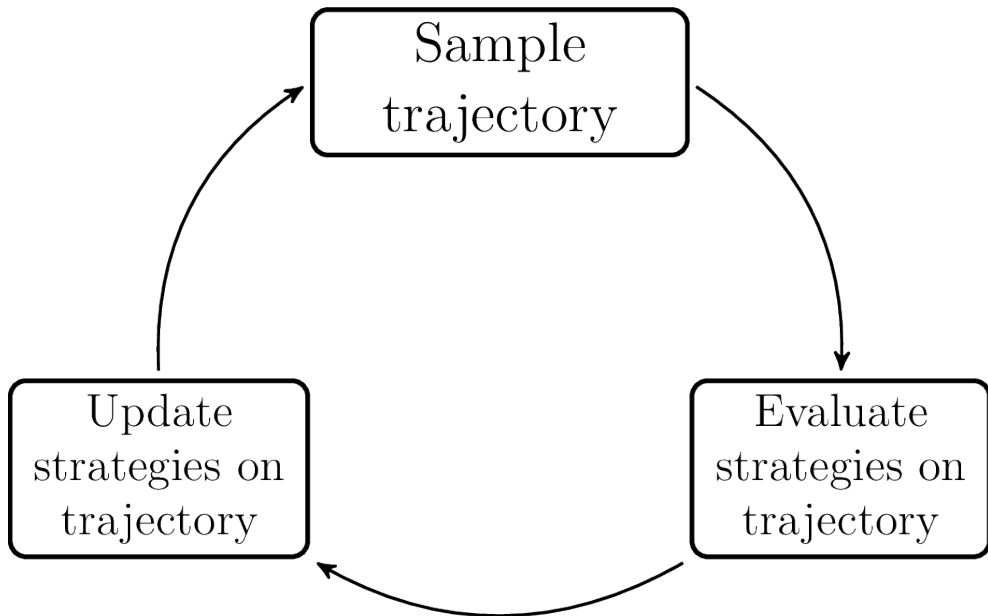
Trevor Davis<sup>2,\*</sup>, Martin Schmid<sup>1</sup>, Michael Bowling<sup>1,2</sup>

\*Work done during internship at DeepMind



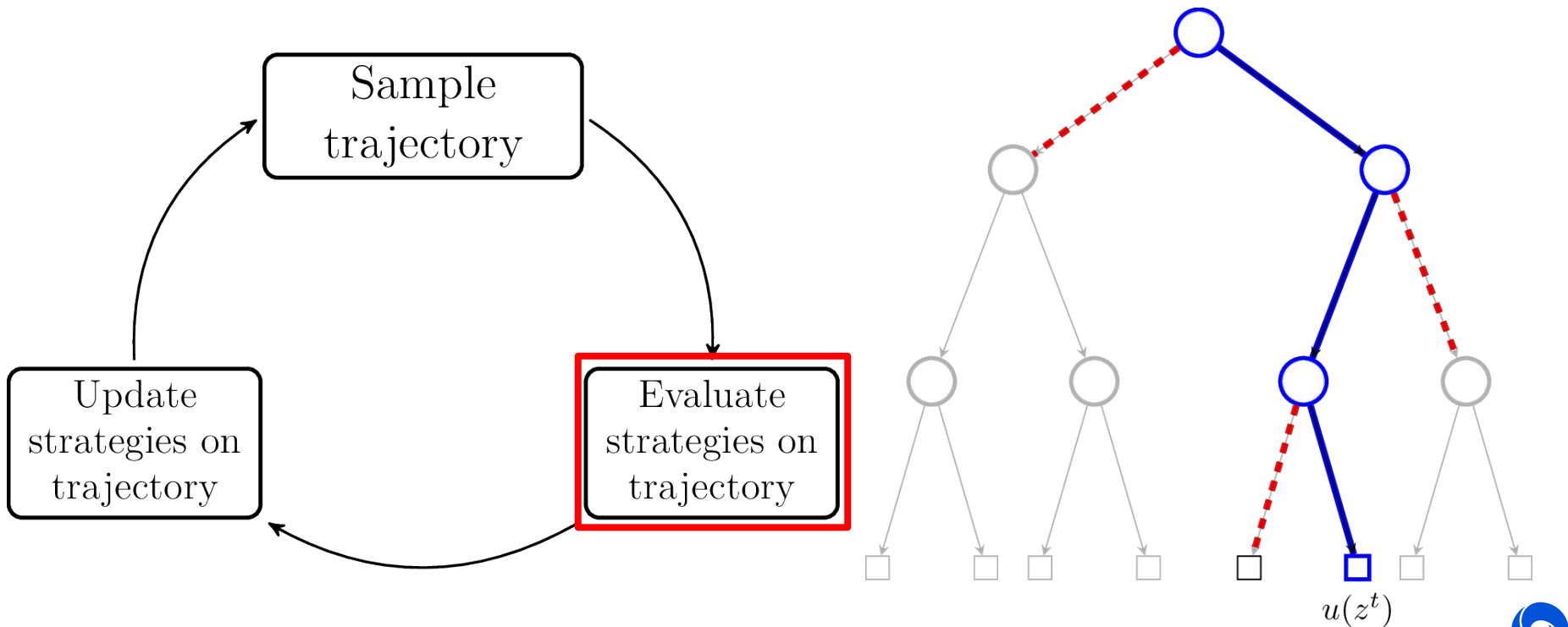
# Monte Carlo game solving

Extensive-form games (EFGs)

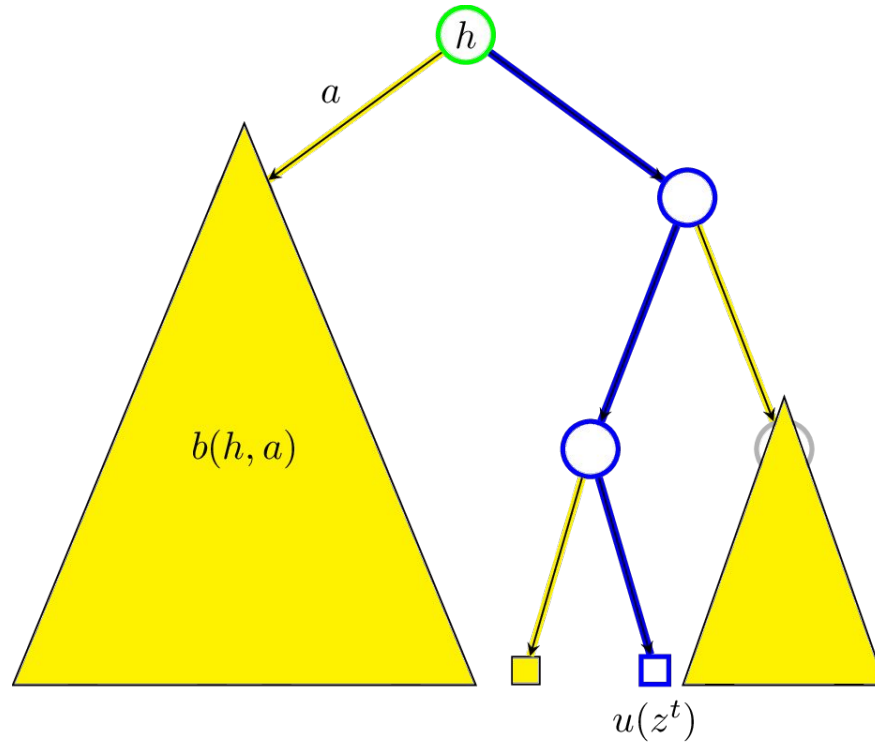


# Monte Carlo game solving

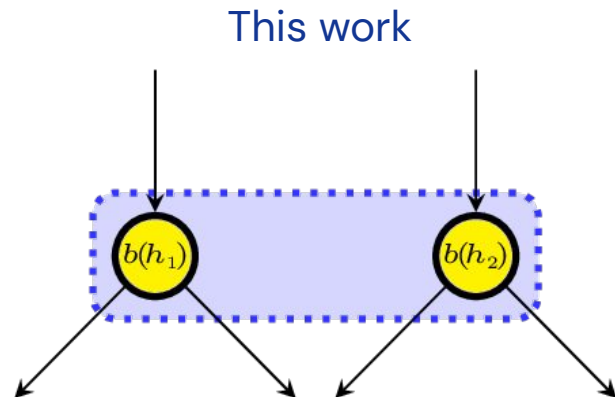
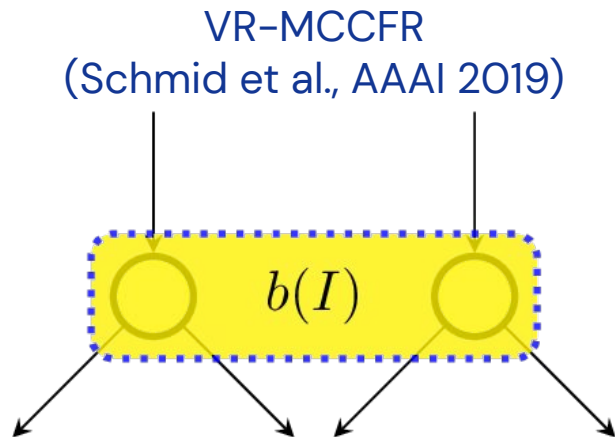
Extensive-form games (EFGs)



# Baseline functions - evaluating unsampled actions



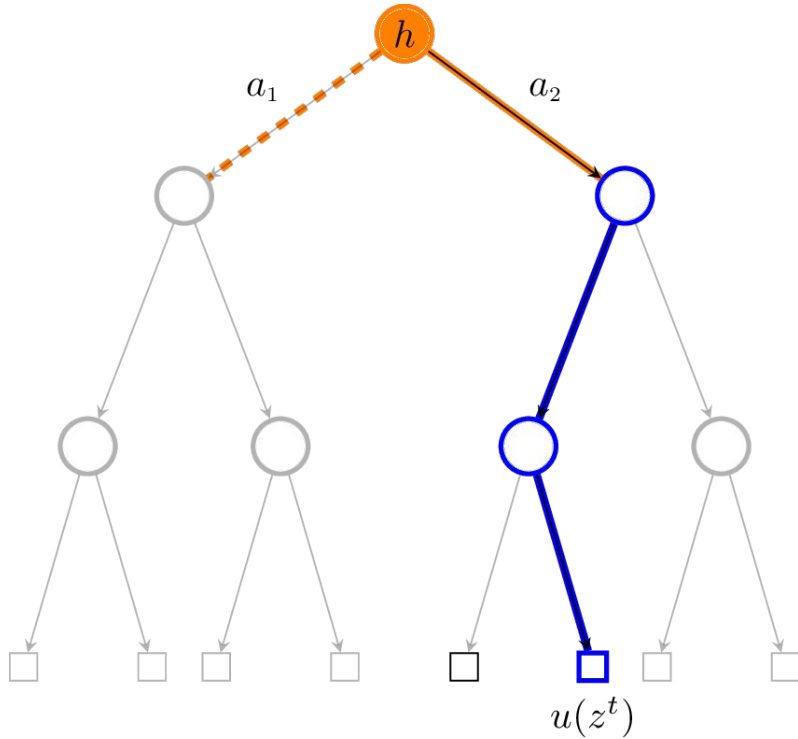
# Our Contribution



- Lower variance, faster convergence
- Provable zero-variance samples



# Monte carlo evaluation



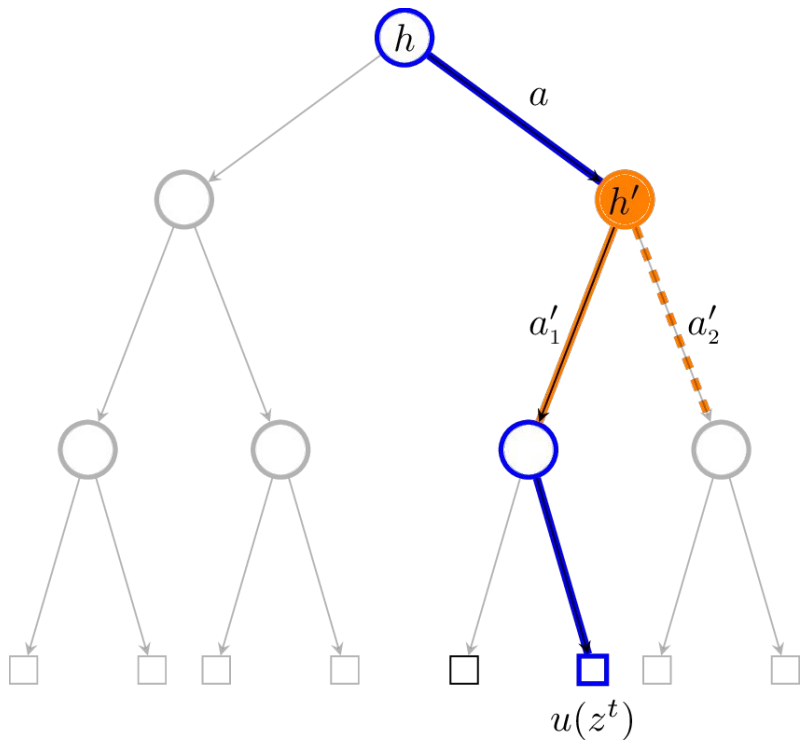
Unbiased updates at  $h$



$$\hat{u}(h, a) = \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} u(z^t)$$



# Monte Carlo evaluation



Unbiased updates at  $h$

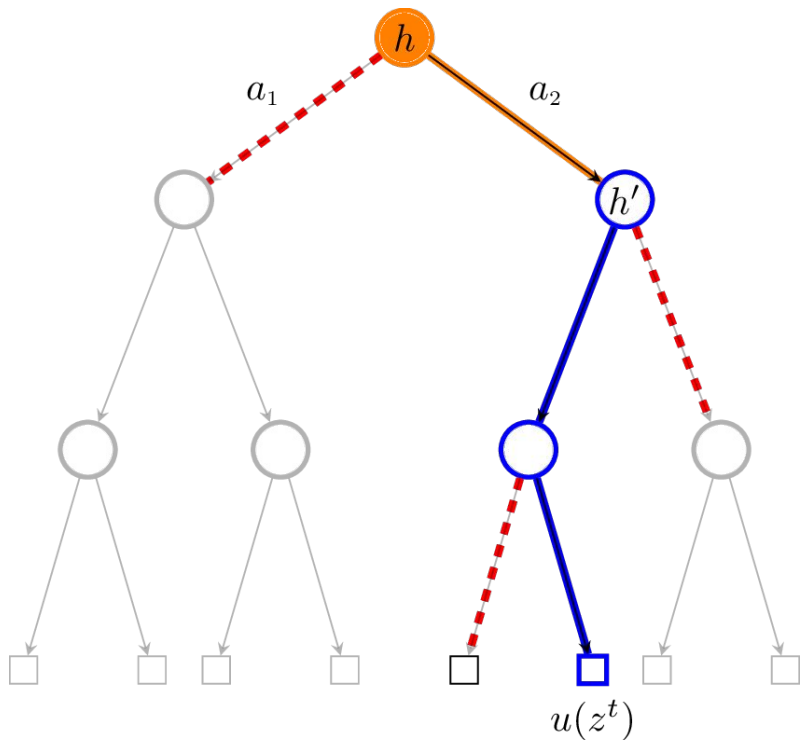


$$\begin{aligned}\hat{u}(h, a) &= \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} u(z^t) \\ &= \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} \hat{u}_b(h')\end{aligned}$$

where  $\hat{u}_b(h') := \sum_{a'} \Pr[h' \rightarrow a'] \hat{u}_b(h', a')$



# Monte Carlo evaluation



Unbiased updates at  $h$



$$\begin{aligned}\hat{u}(h, a) &= \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} u(z^t) \\ &= \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} \hat{u}_b(h')\end{aligned}$$

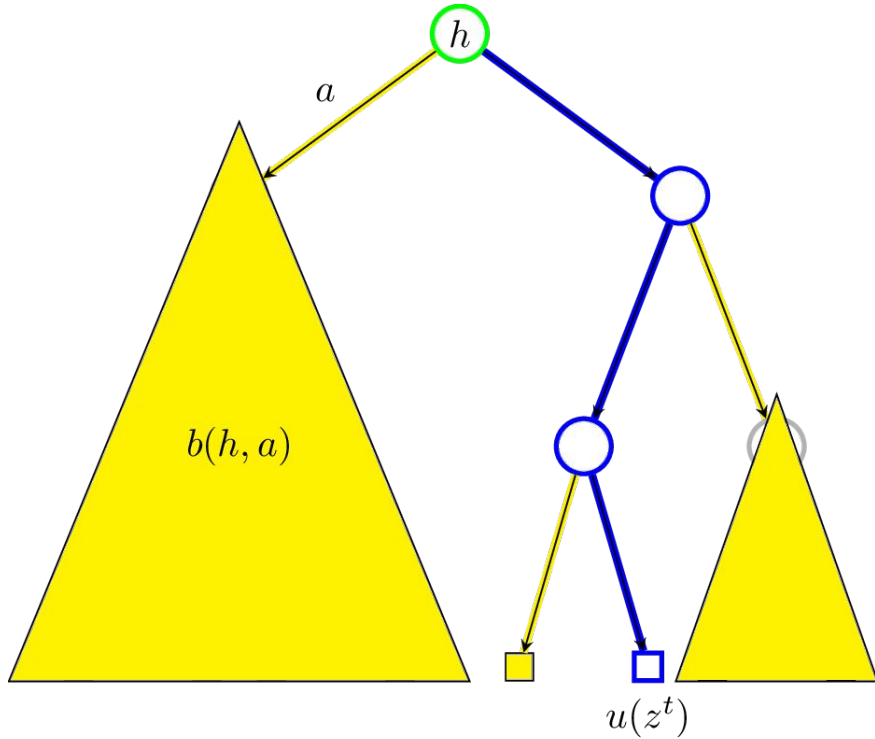
where  $\hat{u}_b(h') := \sum_{a'} \Pr[h' \rightarrow a'] \hat{u}_b(h', a')$

Unsampled actions:  $\hat{u}(h, a) = 0$





# Baseline functions



$$b(h, a) \approx \mathbb{E}[u(h, a)]$$



# Evaluation with baseline

Without baseline:

$$\hat{u}(h, a) = \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} \hat{u}(h')$$

$$\hat{u}(h') = \sum_{a'} \Pr[h' \rightarrow a'] \hat{u}(h', a')$$



# Evaluation with baseline

Without baseline:

$$\hat{u}(h, a) = \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} \hat{u}(h')$$

$$\hat{u}(h') = \sum_{a'} \Pr[h' \rightarrow a'] \hat{u}(h', a')$$

Baseline correction:

$$\hat{u}_b(h, a) = \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} \hat{u}_b(h') + \left( 1 - \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} \right) b(h, a)$$



# Evaluation with baseline

Without baseline:

$$\hat{u}(h, a) = \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} \hat{u}(h')$$

$$\hat{u}(h') = \sum_{a'} \Pr[h' \rightarrow a'] \hat{u}(h', a')$$

Baseline correction:

$$\hat{u}_b(h, a) = \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]} \hat{u}_b(h') + \underbrace{\left(1 - \frac{\mathbb{1}(h \rightarrow a)}{\Pr[h \rightarrow a]}\right)}_{\mathbb{E}[\dots] = 0} b(h, a) \quad (\text{control variate})$$



# Theoretical results

Theorem 1: baseline-corrected values are unbiased:

$$\mathbb{E}[\hat{u}_b(h, a)] = \mathbb{E}[u(h, a)]$$

Theorem 2: each baseline-corrected value  $\hat{u}_b(h, a)$  has variance bounded by a sum of squared prediction errors in the subtree rooted at  $a$

$$\text{Var}[\hat{u}_b(h, a)] \leq \frac{1}{\Pr[h \rightarrow a]} \sum_{\substack{h', a' \in \\ \text{subtree}(h, a)}} \Pr[h, a \rightarrow h', a'] (b(h', a') - \mathbb{E}[u(h', a')])^2$$



## Baseline function selection

We want  $b(h, a) \approx \mathbb{E}[u(h, a)]$

Learned history baseline:

We know  $\mathbb{E}[\hat{u}_b(h, a)] = \mathbb{E}[u(h, a)]$

Set  $b(h, a)$  to average of previous samples  $\hat{u}_b(h, a)$



## Baseline function selection

We want  $b(h, a) \approx \mathbb{E}[u(h, a)]$

Learned history baseline:

We know  $\mathbb{E}[\hat{u}_b(h, a)] = \mathbb{E}[u(h, a)]$

Set  $b(h, a)$  to average of previous samples  $\hat{u}_b(h, a)$

Note:  $\mathbb{E}[u(h, a)]$  depends on strategies – not stationary

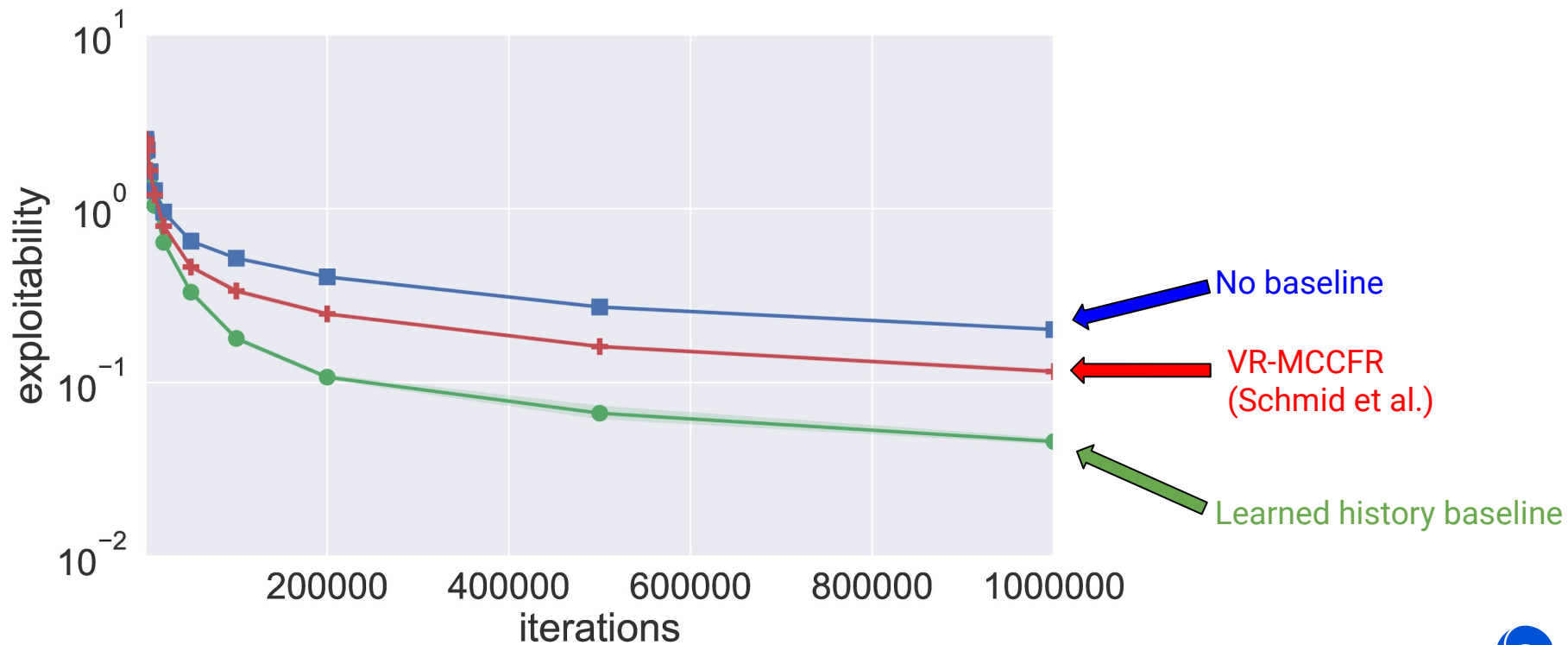
∴  $b(h, a)$  is *not* an unbiased estimate of current expectation

$\hat{u}_b(h, a)$  still unbiased



# Baseline convergence evaluation

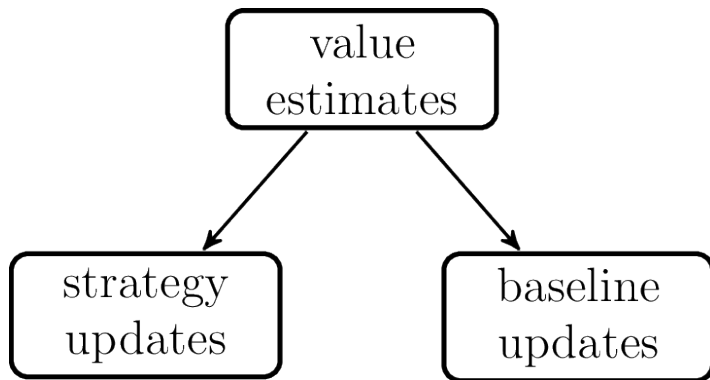
Leduc poker, Monte Carlo Counterfactual Regret Minimization (MCCFR+)





# Predictive baseline

Updating with learned history baseline:



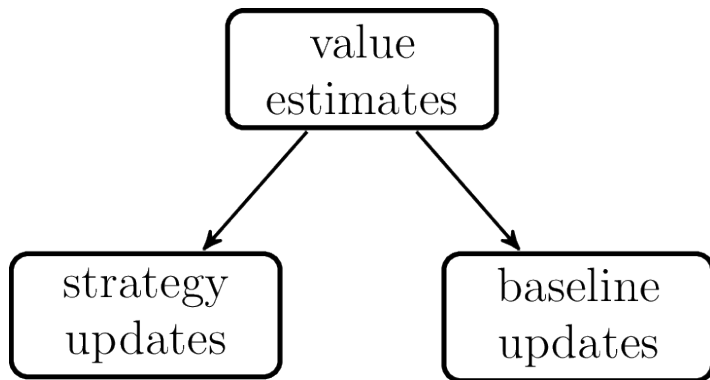
Optimal baseline depends on strategy update:

$$\begin{aligned} b(h, a) &= \mathbb{E}[u(h, a)] \\ &= \sum_{a'} \Pr[h' \rightarrow a'] \mathbb{E}[u(h', a')] \end{aligned}$$

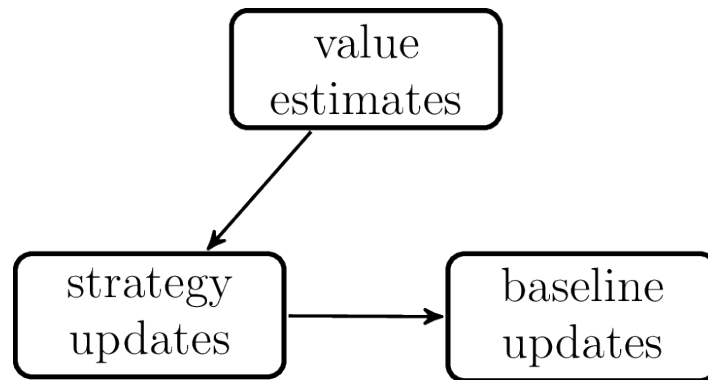


# Predictive baseline

Updating with learned history baseline:



Use strategy to update baseline:



Optimal baseline depends on strategy update:

$$\begin{aligned} b(h, a) &= \mathbb{E}[u(h, a)] \\ &= \sum_{a'} \Pr[h' \rightarrow a'] \mathbb{E}[u(h', a')] \end{aligned}$$

Recursively set

$$b(h, a) = \sum_{a'} \Pr[h' \rightarrow a'] b(h', a')$$



# Zero-variance updates

If:

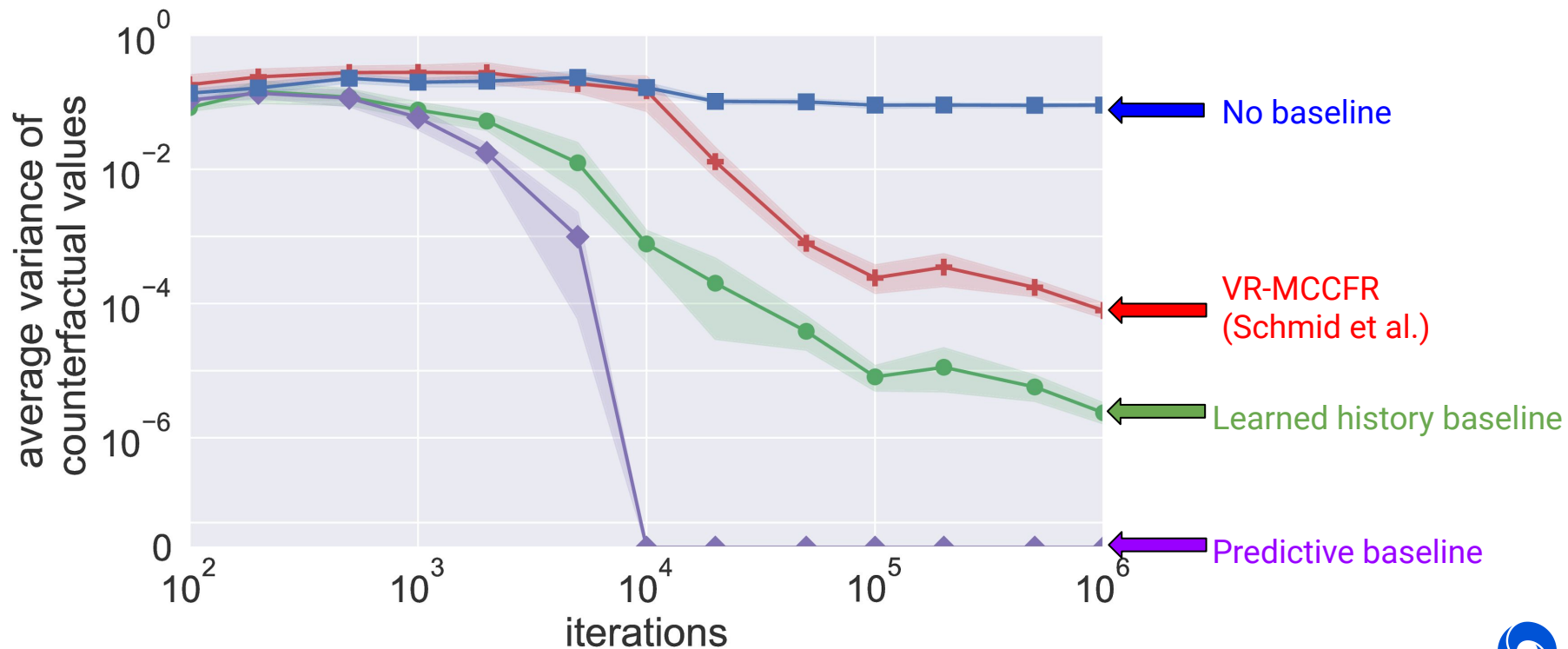
- We use the predictive baseline
- We sample public outcomes
- All outcomes are sampled at least once

Theorem: the baseline-corrected values  $\hat{u}_b(h, a)$  have zero variance

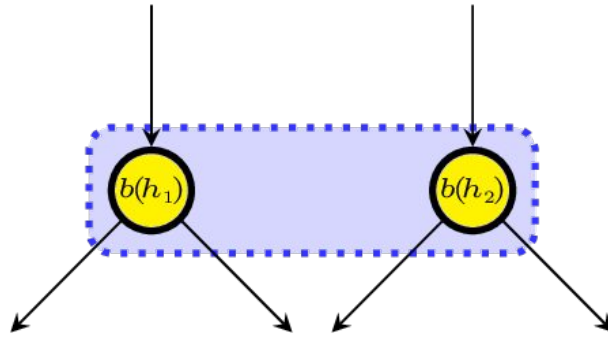


# Baseline variance evaluation

Leduc poker, Monte Carlo Counterfactual Regret Minimization (MCCFR+)



# Conclusion



- Lower variance, faster convergence
- Provable zero-variance samples

