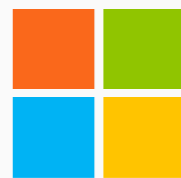


Gradient Boosting for RL in Complex Domains

David Abel², Alekh Agarwal¹, Fernando Diaz¹, Akshay
Krishnamurthy¹, Robert Schapire¹



¹Microsoft Research
²Brown University



ICML RL and Abstraction Workshop 2016

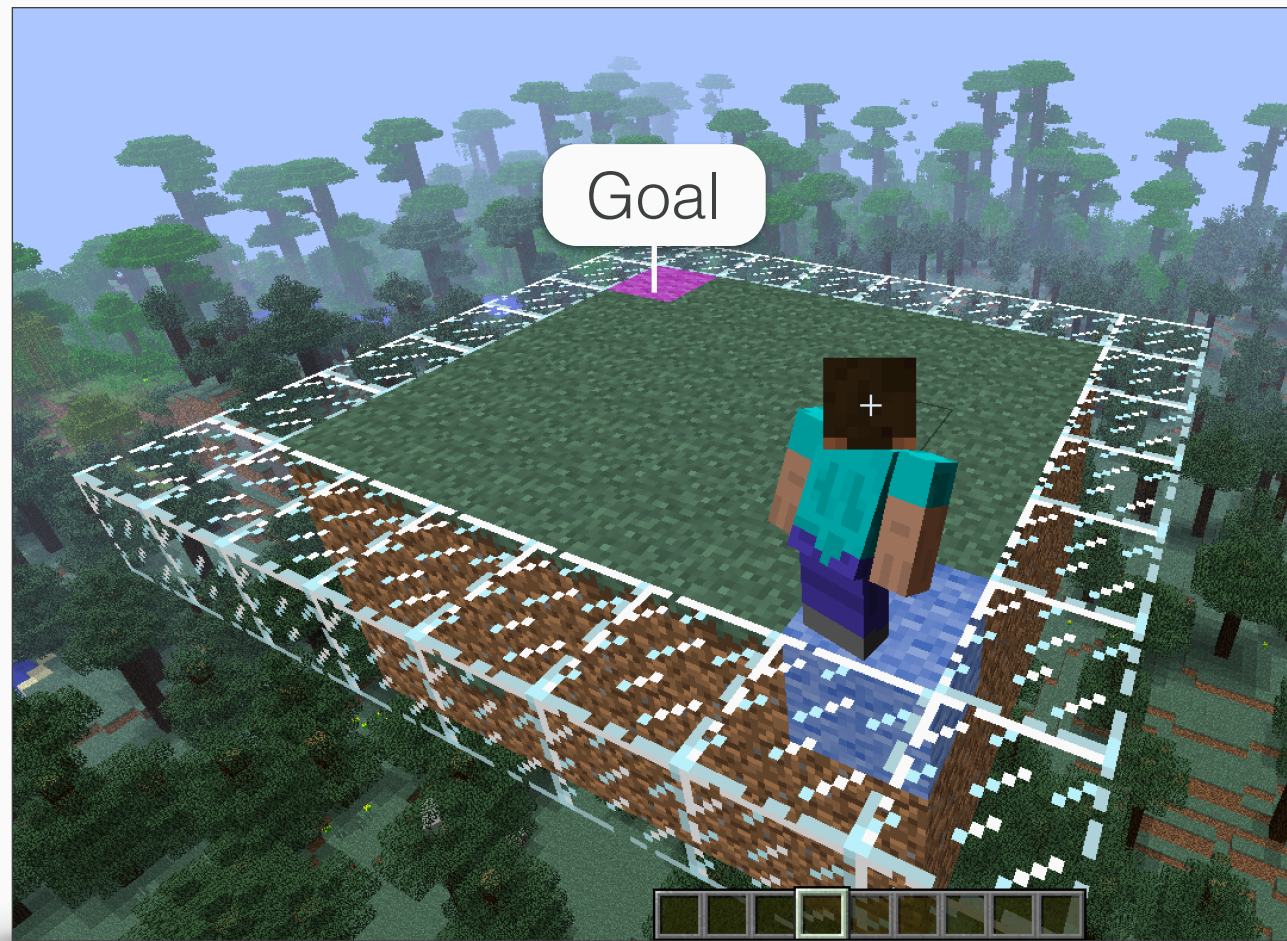
Goal

Develop simple and scalable Reinforcement Learning (RL) techniques that can solve high dimensional problems.

Minecraft



MALMO: Minecraft AI Testbed



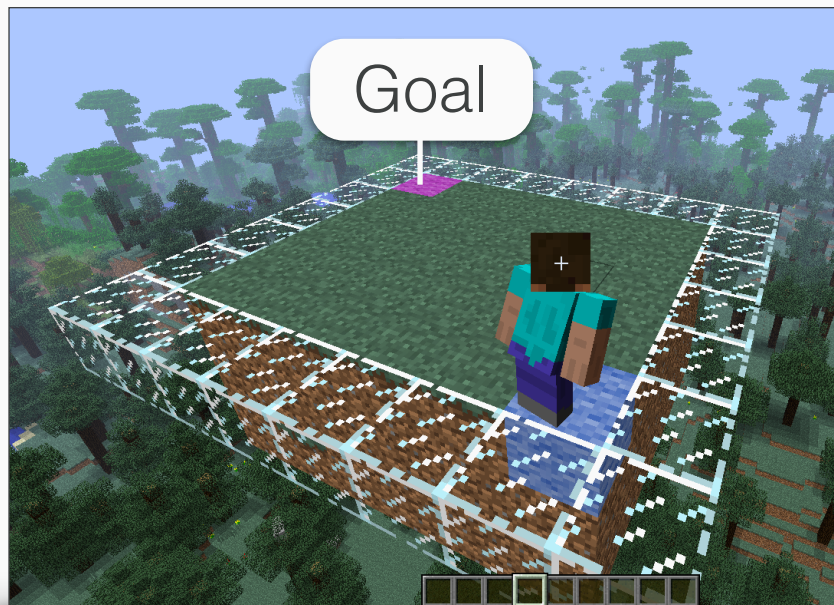
Gridworld

MALMO: an API for
developing agents in
Minecraft

<http://research.microsoft.com/en-us/projects/project-malmo/>

MALMO: Minecraft AI Testbed

Gridworld



difficulty.....



Build 32 bit ALU

Key Components

Developed an RL agent for Minecraft-scale problems:

Key Components

Developed an RL agent for Minecraft-scale problems:

- 1) A *vision system* capable of real-time RL in Minecraft.

Key Components

Developed an RL agent for Minecraft-scale problems:

1) A *vision system* capable of real-time RL in Minecraft.

2) A *new lightweight function approximator* for RL.

└─ Gradient Boosting *[Friedman 2001, Mason 1999]*

Key Components

Developed an RL agent for Minecraft-scale problems:

- 1) A *vision system* capable of real-time RL in Minecraft.
- 2) A *new lightweight function approximator* for RL.**
- 3) An *exploration* technique for model-free RL
(but: preliminary experiments are inconclusive).

Gradient Boosting for RL

Treat RL as a Regression problem for the Q -function

Gradient Boosting for RL

1) Fix an ε -greedy policy with respect to \hat{Q}

Treat RL as a Regression problem for the Q -function

Gradient Boosting for RL

- 1) Fix an ε -greedy policy with respect to \hat{Q}
- 2) Run an episode \rightarrow receive a dataset:

Treat RL as a Regression problem for the Q -function

Gradient Boosting for RL

- 1) Fix an ε -greedy policy with respect to \hat{Q}
- 2) Run an episode \rightarrow receive a dataset:

$$\mathcal{D} = \langle (s_1, a_1, r_1), \dots, (s_N, a_N, r_N) \rangle$$

state

reward

action

Treat RL as a Regression problem for the Q -function

Gradient Boosting for RL

- 1) Fix an ε -greedy policy with respect to \hat{Q}
- 2) Run an episode \rightarrow receive a dataset:
- 3) Fit a new estimate of \hat{Q} by minimizing the Bellman Residual on the data set, \mathcal{D} :

Treat RL as a Regression problem for the Q -function

Gradient Boosting for RL

- 1) Fix an ε -greedy policy with respect to \hat{Q}
- 2) Run an episode \rightarrow receive a dataset:
- 3) Fit a new estimate of \hat{Q} by minimizing the Bellman Residual on the data set, \mathcal{D} :

$$\min_h \sum_{i=1}^N \left[h(s_i, a_i) + \hat{Q}(s_i, a_i) - (r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a')) \right]^2$$

Treat RL as a Regression problem for the Q-function

Gradient Boosting for RL

3) Fit a new estimate of \hat{Q} by minimizing the Bellman Residual on the data set, \mathcal{D} :

$$\min_h \sum_{i=1}^N \left[h(s_i, a_i) + \hat{Q}(s_i, a_i) - (r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a')) \right]^2$$

new weak learner

previous estimate

Bellman residual

Treat RL as a Regression problem for the Q-function

Gradient Boosting for RL

3) Fit a new estimate of \hat{Q} by minimizing the Bellman Residual on the data set, \mathcal{D} :

$$\min_h \sum_{i=1}^N \left[h(s_i, a_i) + \hat{Q}(s_i, a_i) - (r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a')) \right]^2$$

Where:

$$\hat{Q}(s, a) = \sum_{e=1}^E h_e(s, a)$$

Treat RL as a Regression problem for the Q-function

Gradient Boosting for RL

3) Fit a new estimate of \hat{Q} by minimizing the Bellman Residual on the data set, \mathcal{D} :

$$\min_h \sum_{i=1}^N \left[h(s_i, a_i) + \hat{Q}(s_i, a_i) - (r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a')) \right]^2$$

Where:

episodes E *one weak learner per episode*

$$\hat{Q}(s, a) = \sum_{e=1}^E h_e(s, a)$$

Treat RL as a Regression problem for the Q-function

Gradient Boosting for RL

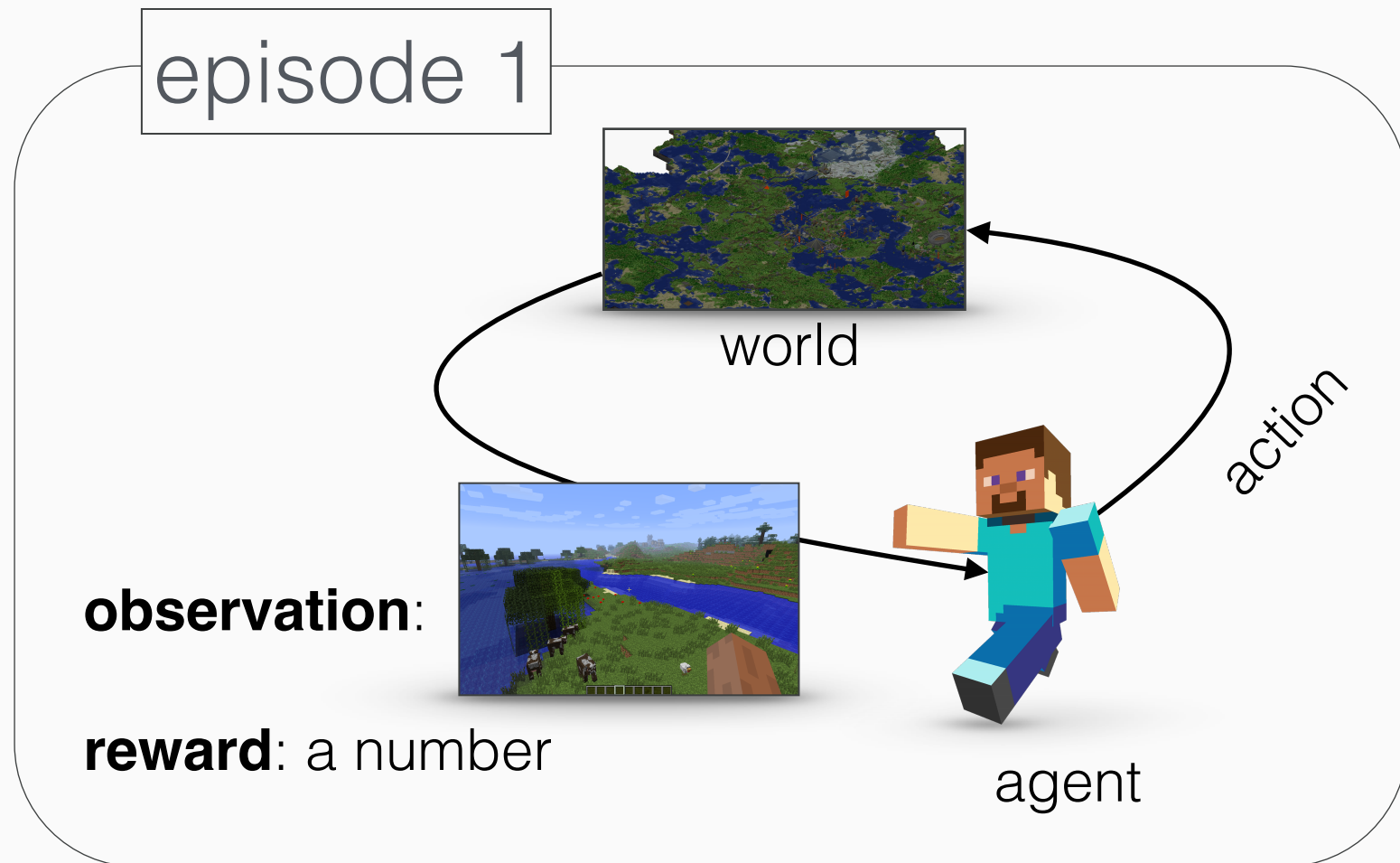
3) Fit a new estimate of \hat{Q} by minimizing the Bellman Residual on the data set, \mathcal{D} :

$$\min_h \sum_{i=1}^N \left[h(s_i, a_i) + \hat{Q}(s_i, a_i) - (r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a')) \right]^2$$

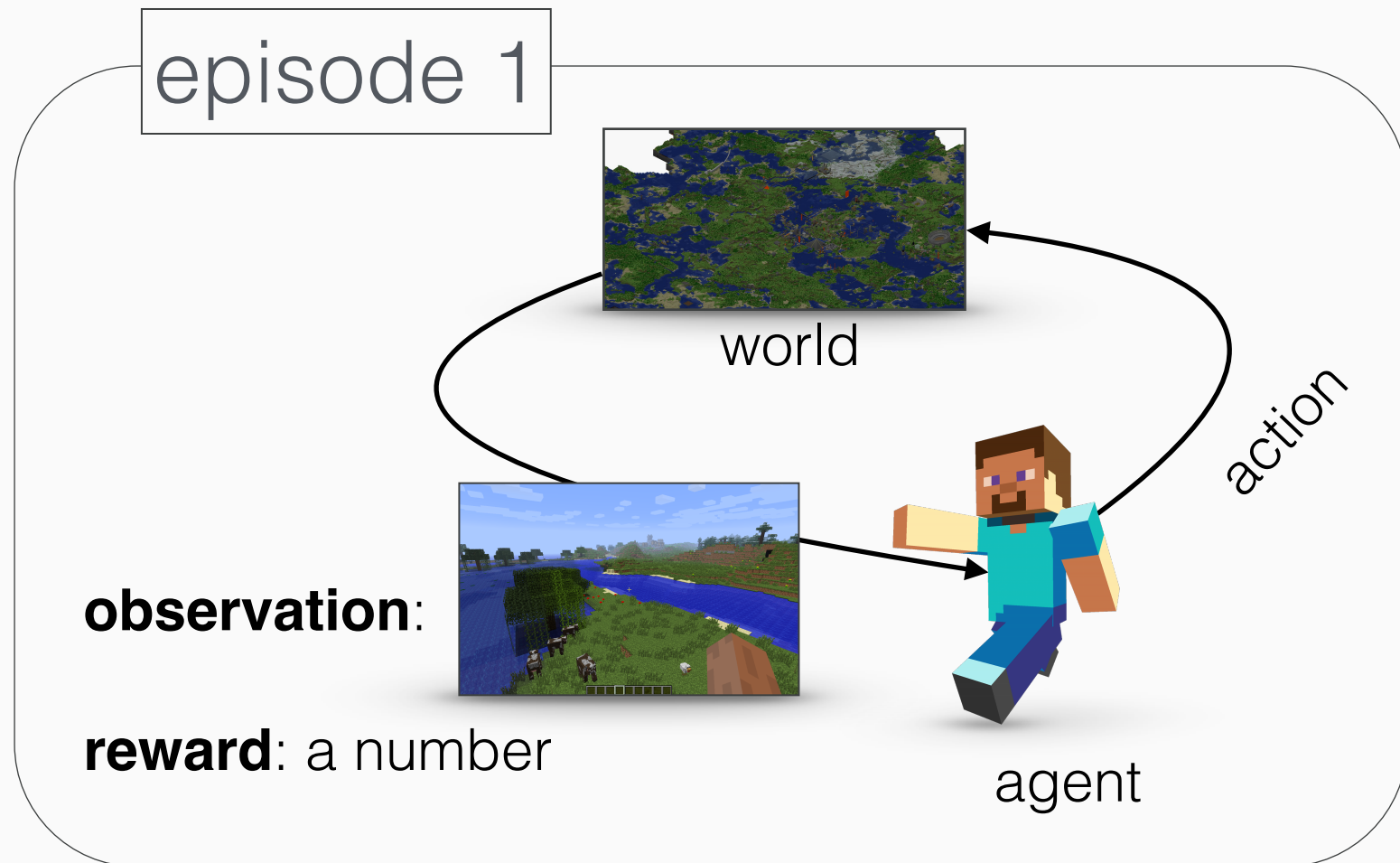
We solve this using **regression trees** as the weak learner

Treat RL as a Regression problem for the Q-function

High Level

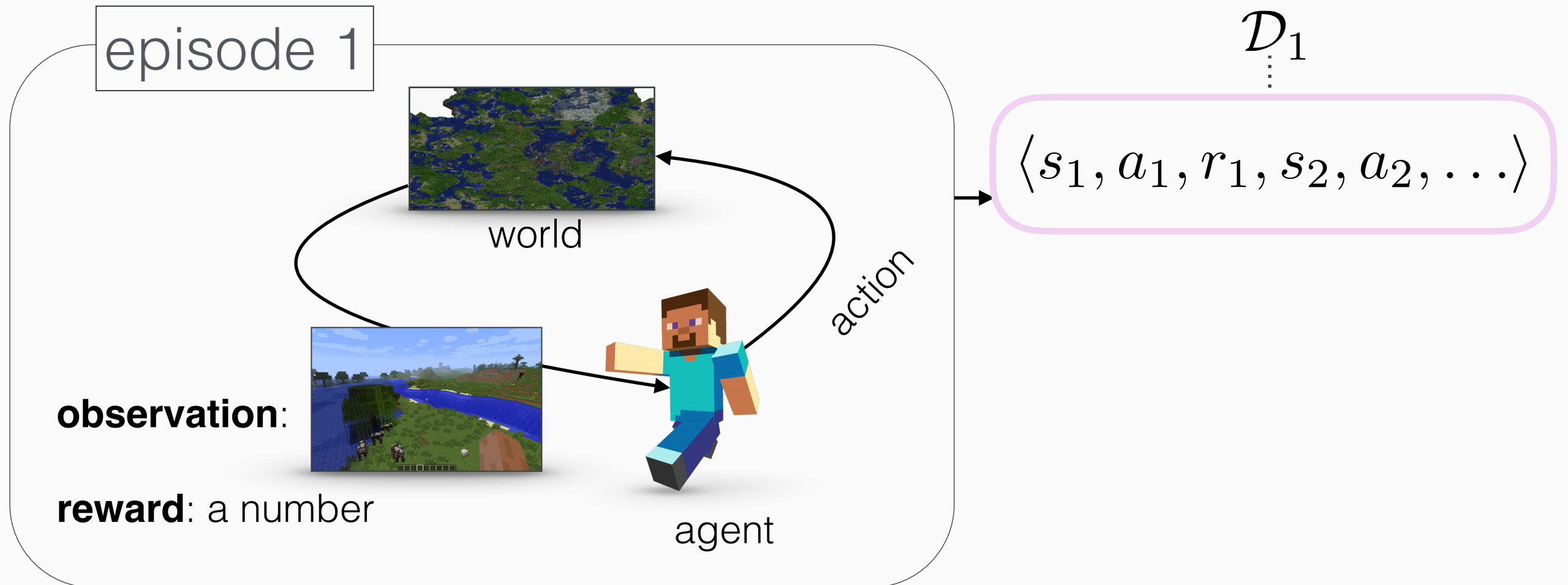


High Level



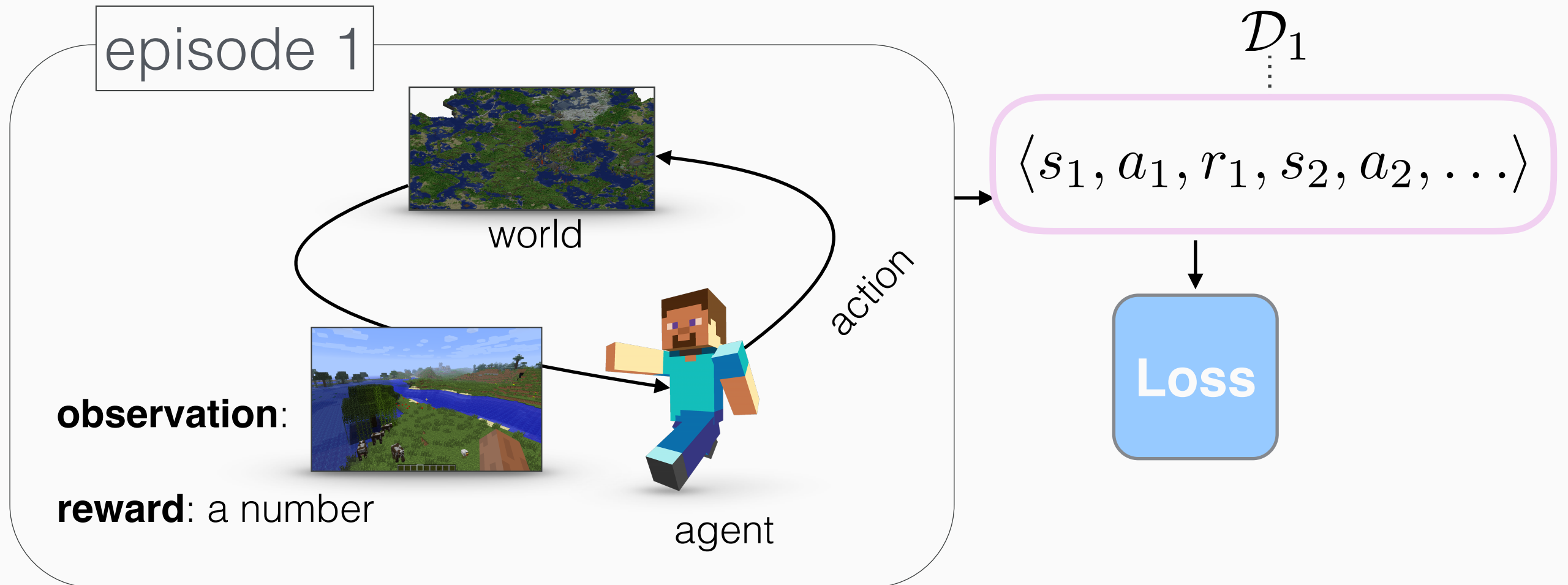
$$\hat{Q} = \sum (\text{ })$$

High Level



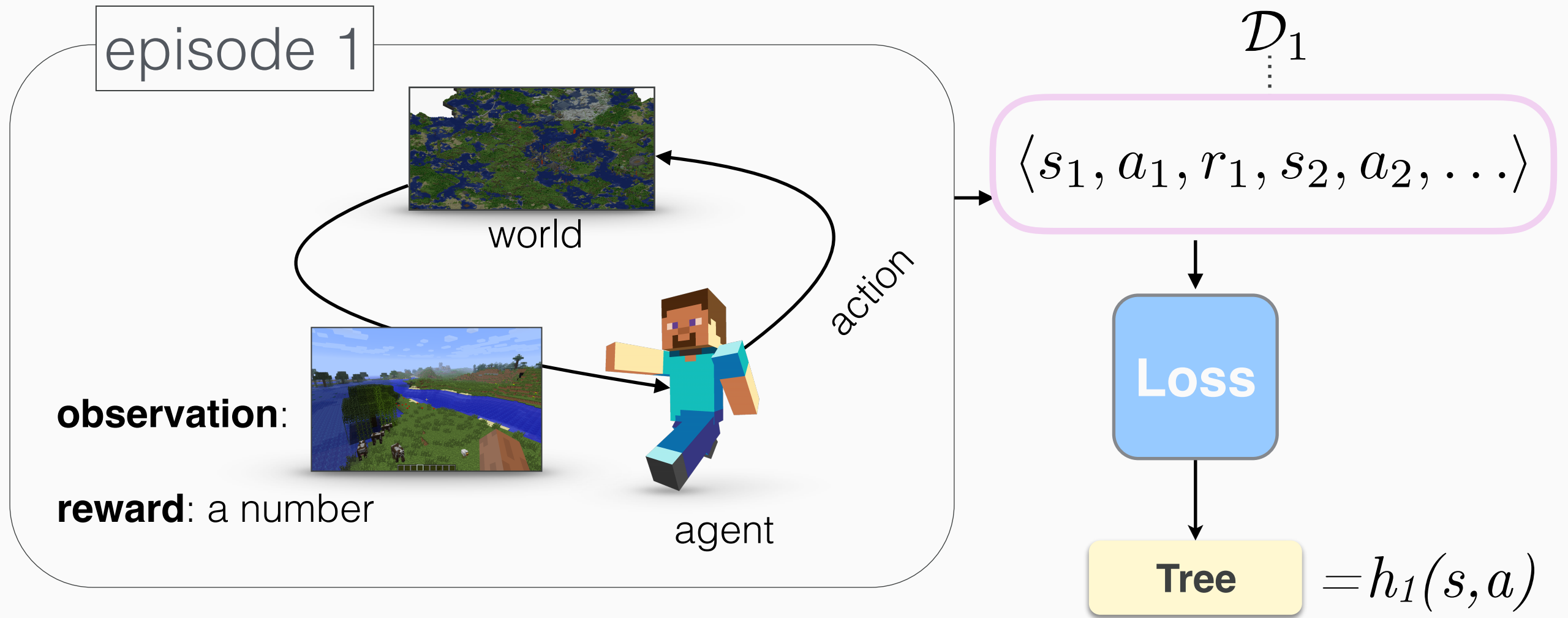
$$\hat{Q} = \sum (\text{[Redacted]})$$

High Level



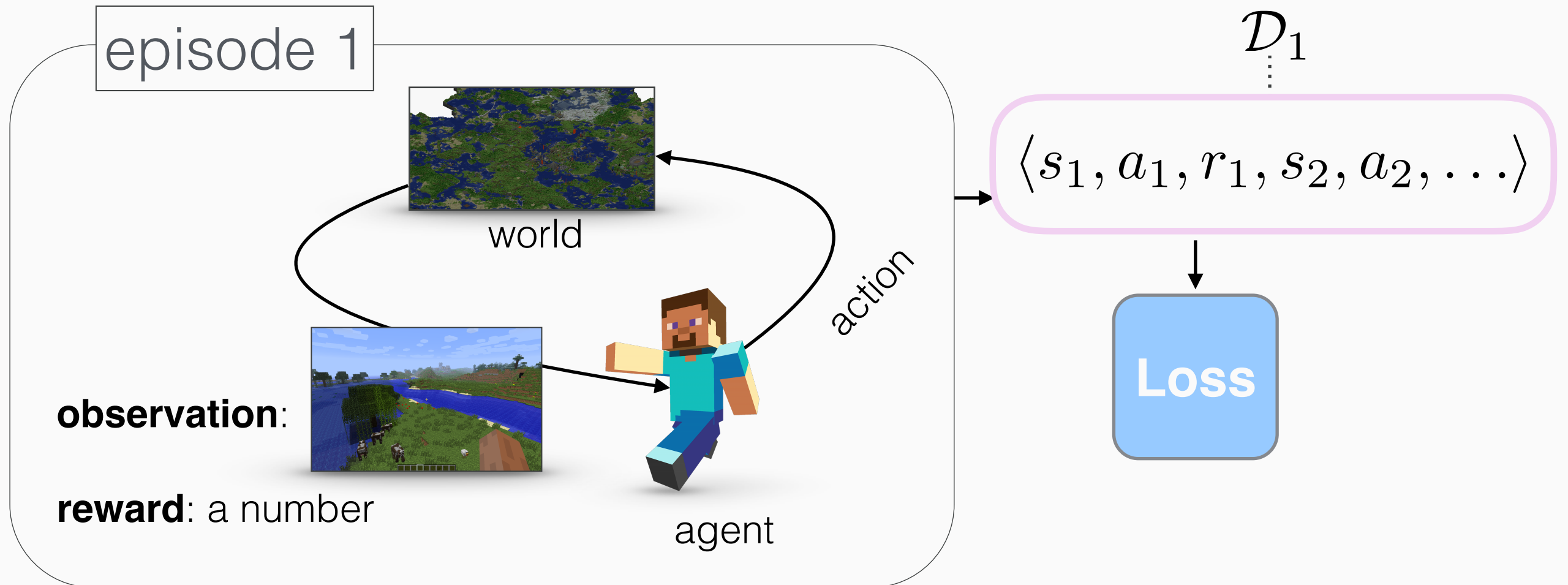
$$\hat{Q} = \sum (\text{ })$$

High Level



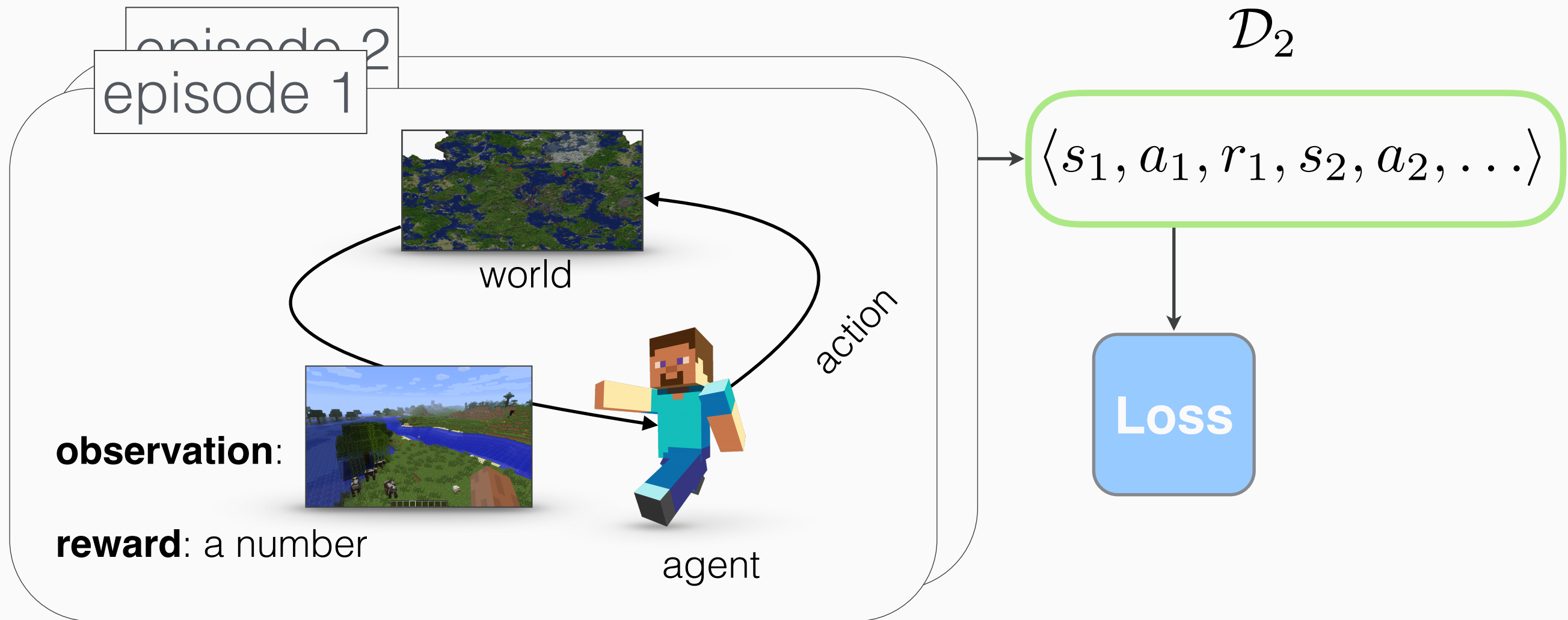
$$\hat{Q} = \sum (\text{[Redacted Box]})$$

High Level



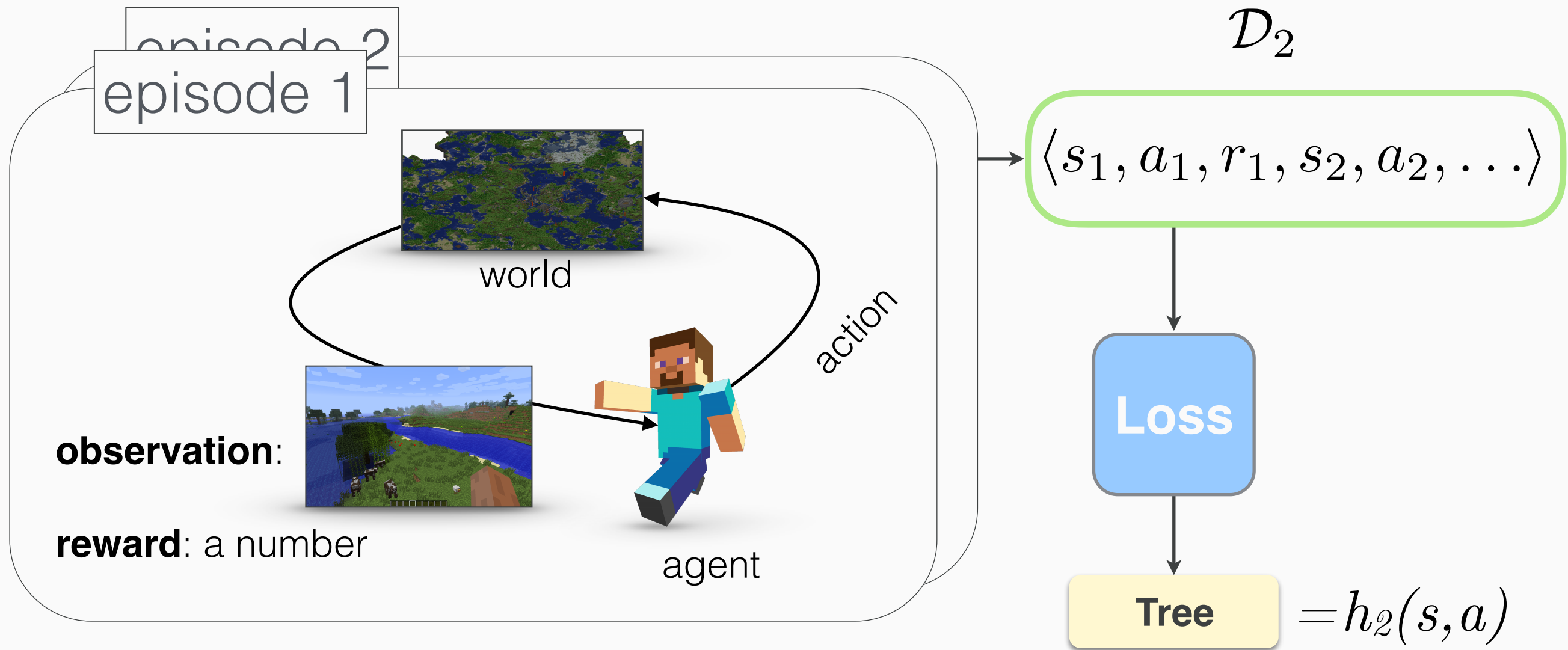
$$\hat{Q} = \sum \left(\begin{array}{c} h_1 \\ \vdots \\ \mathcal{D}_1 \end{array} \right)$$

High Level



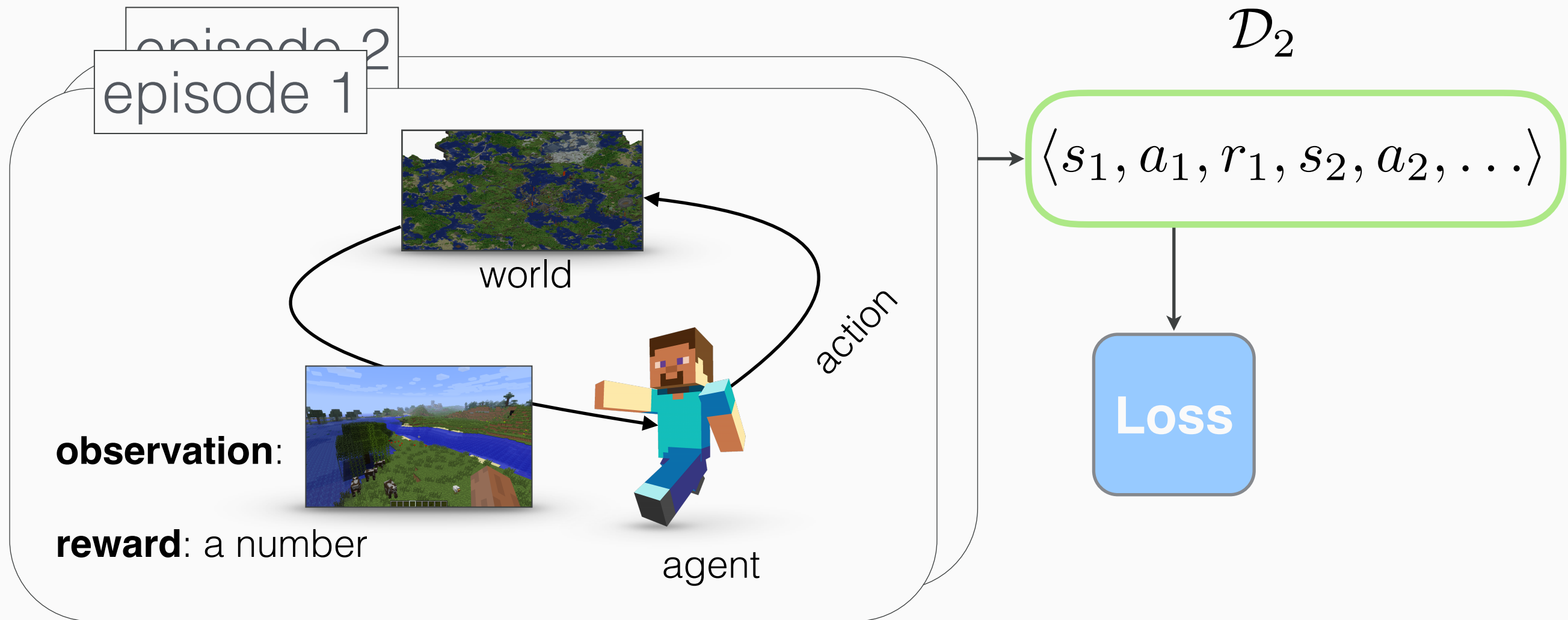
$$\hat{Q} = \sum \left(\begin{array}{c} h_1 \\ \vdots \\ \mathcal{D}_1 \end{array} \right)$$

High Level



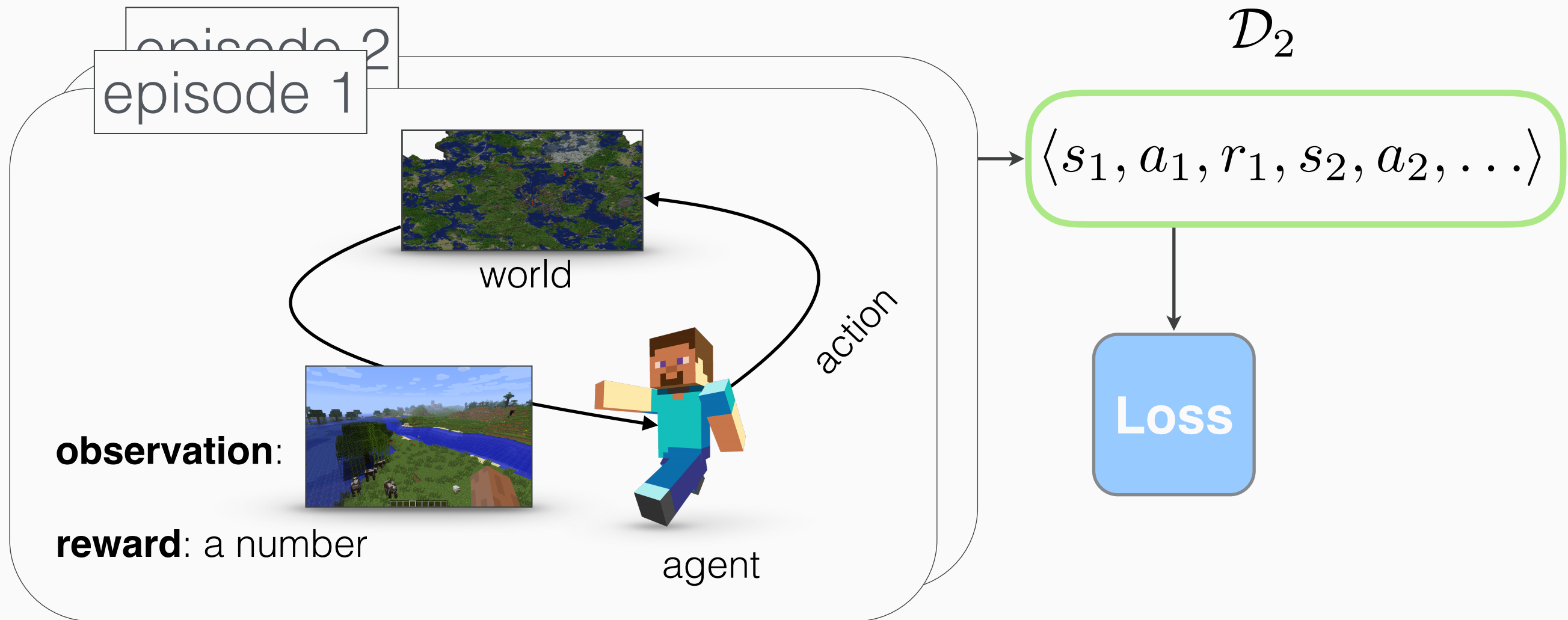
$$\hat{Q} = \sum \left(\begin{array}{c} h_1 \\ \vdots \\ \mathcal{D}_1 \end{array} \right)$$

High Level



$$\hat{Q} = \sum \left(\begin{array}{c} h_1, h_2 \\ \vdots \quad \vdots \\ \mathcal{D}_1 \quad \mathcal{D}_2 \end{array} \right)$$

High Level



$$\hat{Q} = \sum \left(\begin{array}{cccccccc} h_1 & h_2 & h_3 & h_4 & h_5 & h_6 & h_7 & h_8 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathcal{D}_1 & \mathcal{D}_2 & \dots & \dots & \dots & \dots & \dots & \mathcal{D}_8 \end{array} \right)$$

Intuitively Nice Properties

- Non-parametric
- Simple, easy to implement, minimal hand-engineering
- Interleaved data collection
- Rich theoretical literature, room for analysis.
- Only need to store one episode's worth of data.

Experiments: Baselines

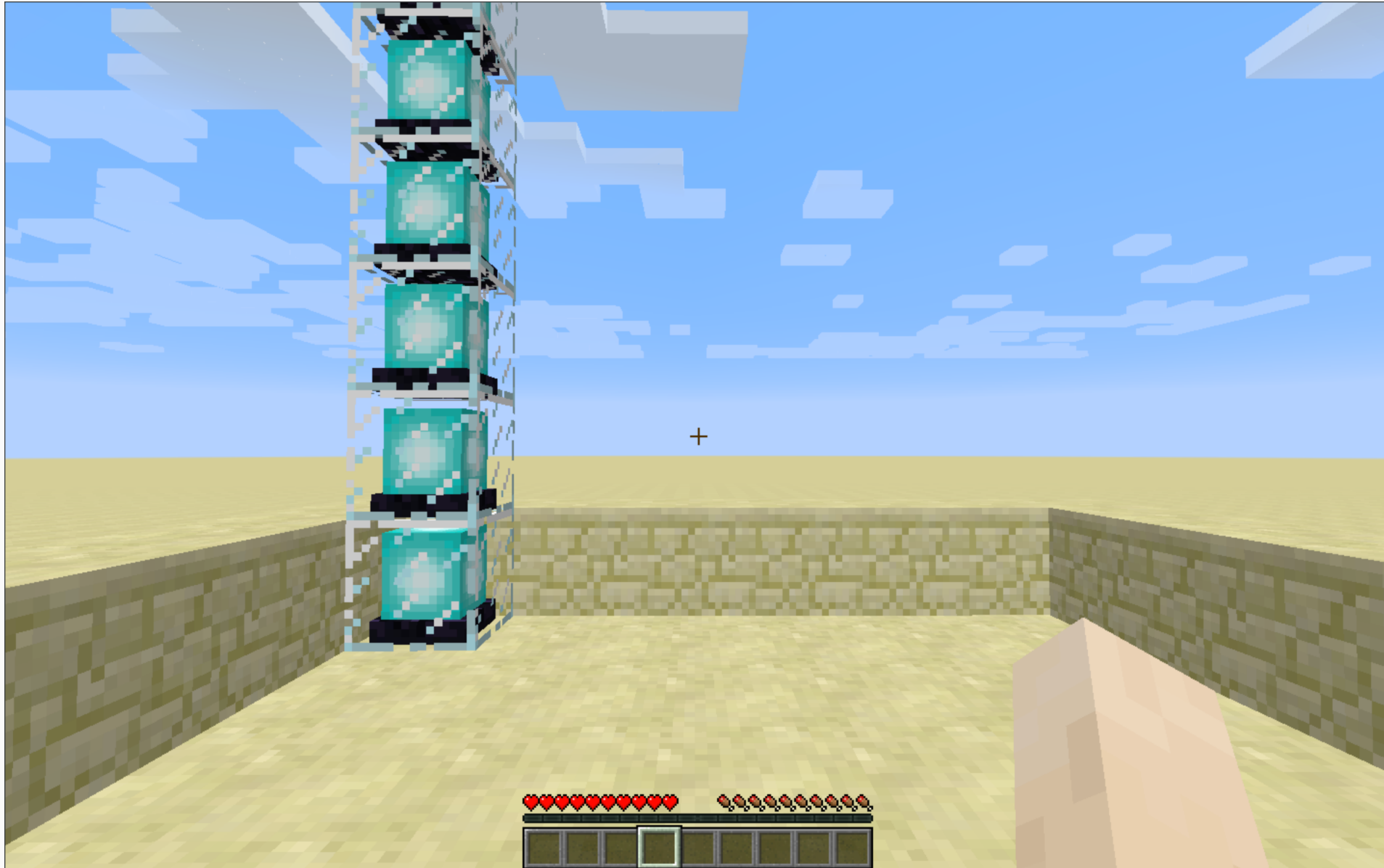
- Baseline 1 (*Linear Approximator*)
- Baseline 2 (*Random Forest Approximator*)
- Baseline 3 (*Batch Boost Approximator*)

Experiments: Baselines

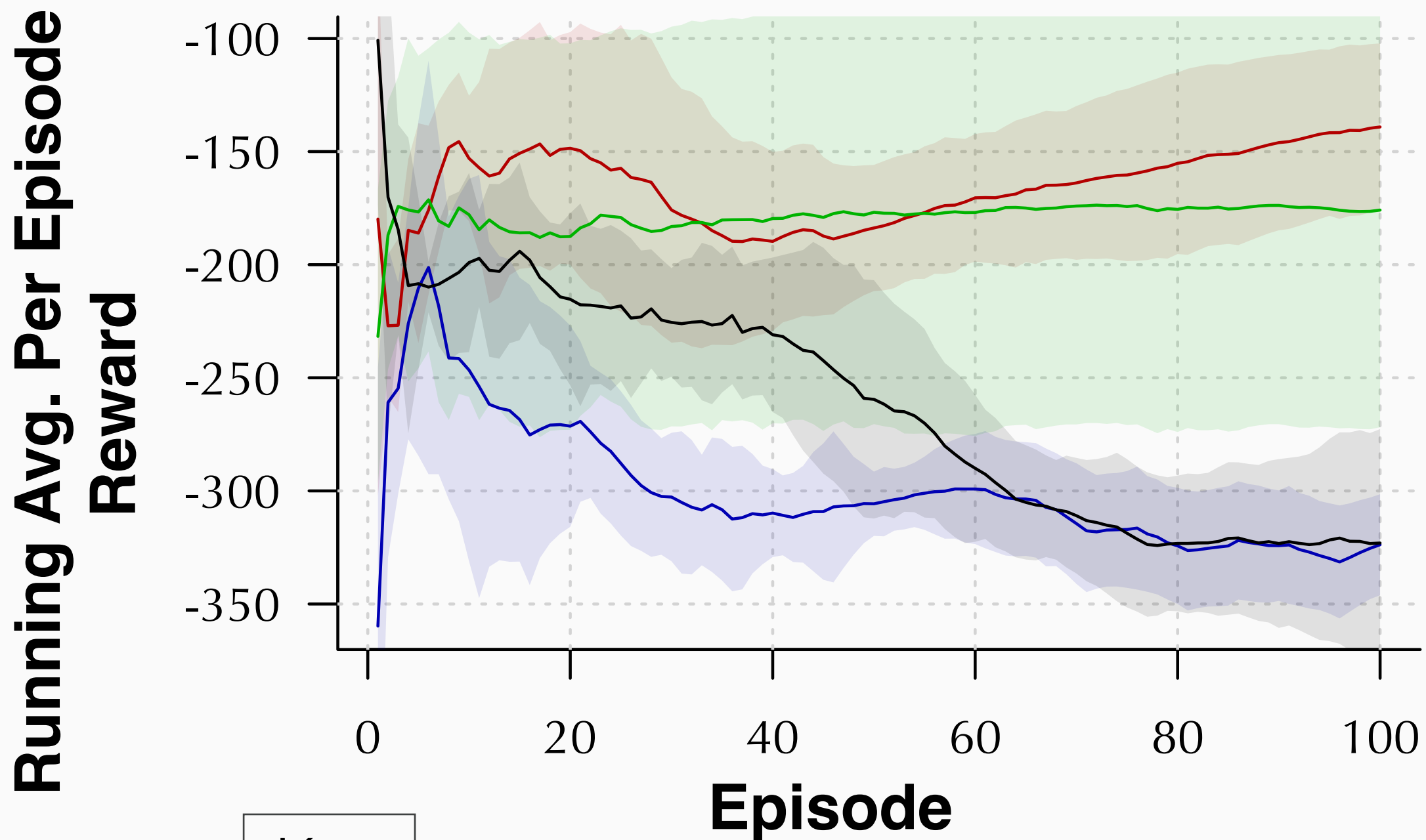
- Baseline 1 (*Linear Approximator*)
- Baseline 2 (*Random Forest Approximator*)
- Baseline 3 (*Batch Boost Approximator*)

└ Similar to Fitted Q-iteration [*Ernst et al. 2005*]

Experiments: Visual Grid



Visual Grid: Results



Key

Gradient Booster

Batch Booster

Linear

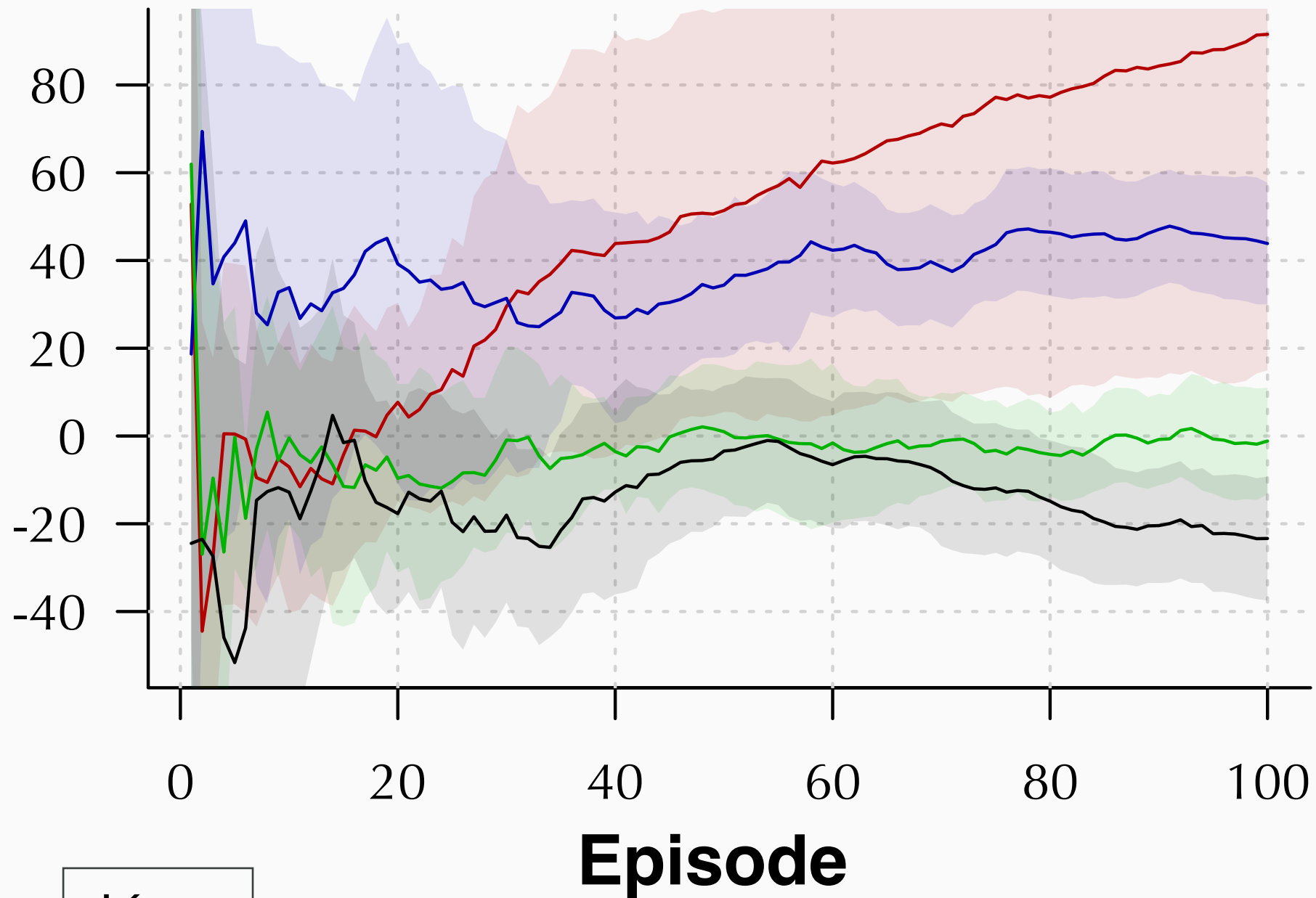
Forest

Experiments: Hillclimbing



Visual Hill Climb: Results

Running Avg. Per Episode
Reward



Key

Gradient Booster

Batch Booster

Linear

Forest

Next Steps

- Investigate relevant exploration techniques inspired by Gradient Boosting.
- Use rich foundation of theory on gradient boosting to inspire analysis of this approach.
- Further experimentation.

Acknowledgments

A big thank you to The *MALMO* team!

**David Bignell, Katja Hofmann, Tim Hutton,
Matthew Johnson, Pushmeet Kohli, Nate
Kushman, Ewa Luger, Bhaskar Mitra, Jamie
Shotton, Evelyne Viegas.**

<http://research.microsoft.com/en-us/projects/project-malmo/>