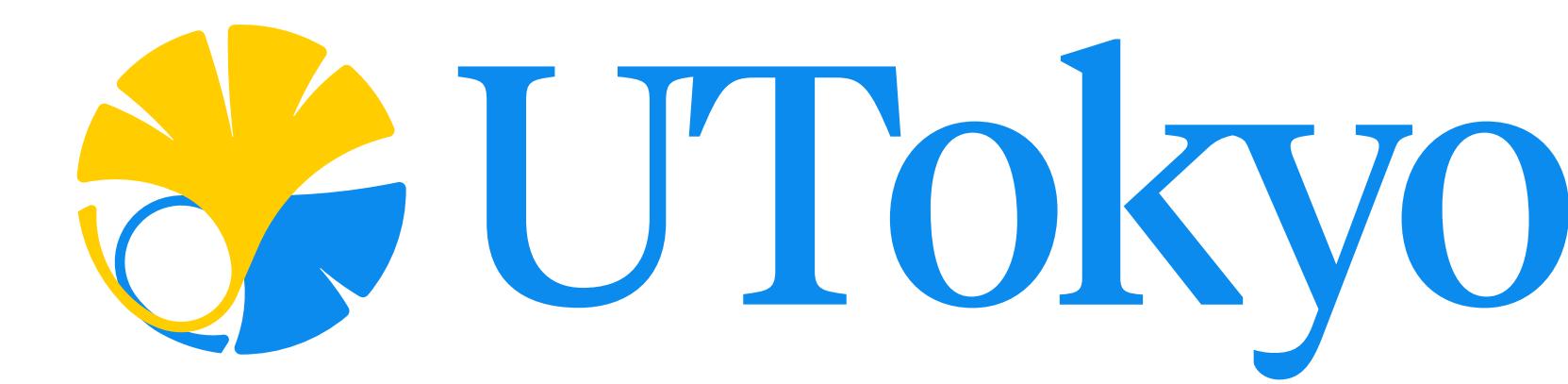


Recursive Reward Aggregation

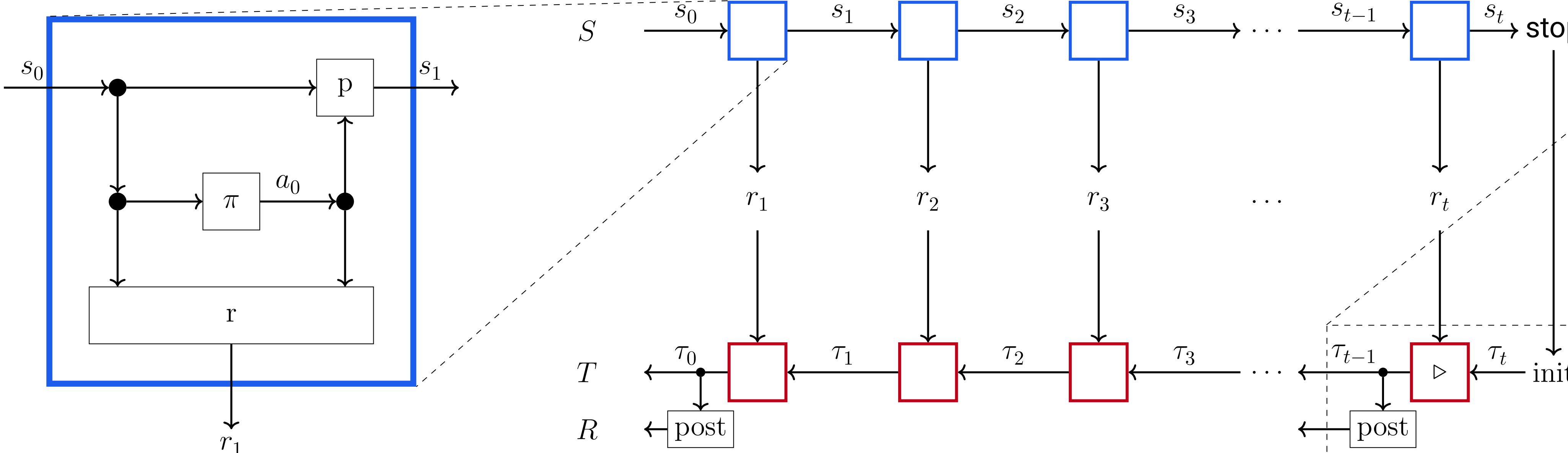


Yuting Tang^{1,2} Yivan Zhang^{1,2} Johannes Ackermann^{1,2}
 Yu-Jie Zhang² Soichiro Nishimori^{1,2} Masashi Sugiyama^{2,1}
¹The University of Tokyo ²RIKEN AIP

arXiv:2507.08537
<https://yivan.xyz/rra>

Recursive generation and recursive aggregation lead to generalized Bellman equations.

- states $s \in S$
- actions $a \in A$
- rewards $r \in R$
- statistics $\tau \in T$
- transition function $p : S \times A \rightarrow S$
- policy $\pi : S \rightarrow A$
- reward function $r : S \times A \rightarrow R$
- copy $\bullet : x \mapsto (x, x)$



definition	initial value of statistic(s) init $\in T$	update function $\triangleright : R \times T \rightarrow T$	post-processing post : $T \rightarrow R$
discounted sum	$r_1 + \gamma r_2 + \dots + \gamma^{t-1} r_t$	discounted sum $s : 0 \in \mathbb{R}$ $+_\gamma := [r, s \mapsto r + \gamma \cdot s]$	$\text{id}_{\mathbb{R}}$
discounted max	$\max\{r_1, \gamma r_2, \dots, \gamma^{t-1} r_t\}$	discounted max $m : -\infty \in \overline{\mathbb{R}}$ $\max_\gamma := [r, m \mapsto \max(r, \gamma \cdot m)]$	$\text{id}_{\overline{\mathbb{R}}}$
mean	$\bar{r} := \frac{1}{t} \sum_{i=1}^t r_i$	length n sum s $[0] \in \begin{cases} \mathbb{N} & n = 0 \\ \mathbb{R} & n > 0 \end{cases}$ $[r, [n]] \mapsto \begin{cases} [n+1] & n = 0 \\ [s+r] & n > 0 \end{cases}$	$\left[\begin{array}{c} [n] \\ s \end{array} \right] \mapsto \frac{s}{n}$
variance	$\frac{1}{t} \sum_{i=1}^t (r_i - \bar{r})^2 = \bar{r}^2 - \bar{r}^2$	length n sum s sum square q $[0] \in \begin{cases} \mathbb{N} & n = 0 \\ \mathbb{R}_{\geq 0} & n > 0 \end{cases}$ $[r, [n]] \mapsto \begin{cases} [n+1] & n = 0 \\ [s+r^2] & n > 0 \end{cases}$	$\left[\begin{array}{c} [n] \\ s \\ q \end{array} \right] \mapsto \frac{q}{n} - (\frac{s}{n})^2$
top- k	k -th largest in $r_{1:t}$	top-1 top-2 \vdots $[-\infty, \infty] \in \overline{\mathbb{R}}^k$	$[r, b \mapsto \begin{cases} \text{insert}(r, b) & r > \min b \\ b & r \leq \min b \end{cases}]$ $[b \mapsto \min b]$

Summary

- An **algebra fusion** perspective on the recursive structure of Markov decision processes (MDPs)
- Generalized **Bellman equations** for alternative reward aggregations (e.g., max, mean, variance)
- Integration into *value-based* (e.g., Q-learning) and *actor-critic* (e.g., PPO, TD3) algorithms
- Direct optimization of the *Sharpe ratio* ($\frac{\text{mean}}{\text{std}}$) in portfolio management tasks in finance

Why do we optimize the discounted sum?

- Standard? Convenient? Mathematically elegant?
- Practical benefits? Theoretical guarantees?
- *Reward hypothesis*? If true, all preferences can be represented in this way.

But... does it always align with our requirements?

- **Max**: peak performance in drug discovery
- **Min**: bottleneck objective in network routing
- **Variance-based** criteria: risk-aversion in finance, robotics, and control

We want to optimize other reward aggregations, directly, efficiently, and effectively.

Discounted sum is a *recursive function*

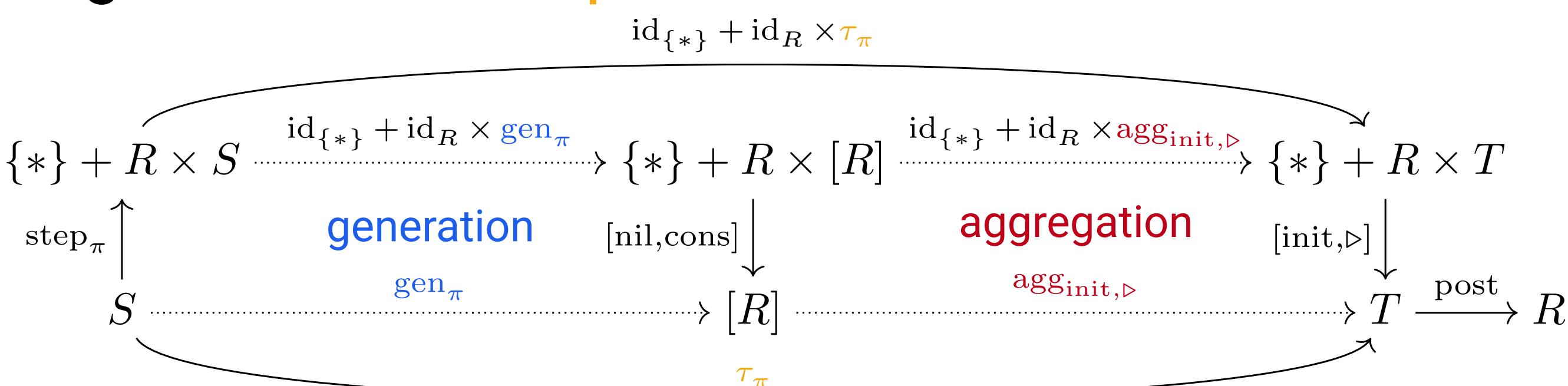
... and we can use other recursive aggregations.

$$\begin{aligned} \text{sum}_\gamma[r_1, r_2, r_3, \dots] &= r_1 + \gamma r_2 + \gamma^2 r_3 + \dots \\ &= r_1 + \gamma(r_2 + \gamma r_3 + \dots) \\ &= r_1 + \gamma \text{sum}_\gamma[r_2, r_3, \dots] \\ \text{agg}_{\text{init}, \triangleright}[r_1, r_2, r_3, \dots] &:= r_1 \triangleright \text{agg}_{\text{init}, \triangleright}[r_2, r_3, \dots] \end{aligned}$$

- **initial value** $\text{init} \in T$ of statistics ($\text{agg}_{\text{init}, \triangleright}[\] = \text{init}$)
- **update function** $\triangleright : R \times T \rightarrow T$
- **statistic aggregation function** $\text{agg}_{\text{init}, \triangleright} : [R] \rightarrow T$

Reward generation is also a *recursive function*

... so we can "fuse" it with a recursive aggregation to get a **Bellman equation**.

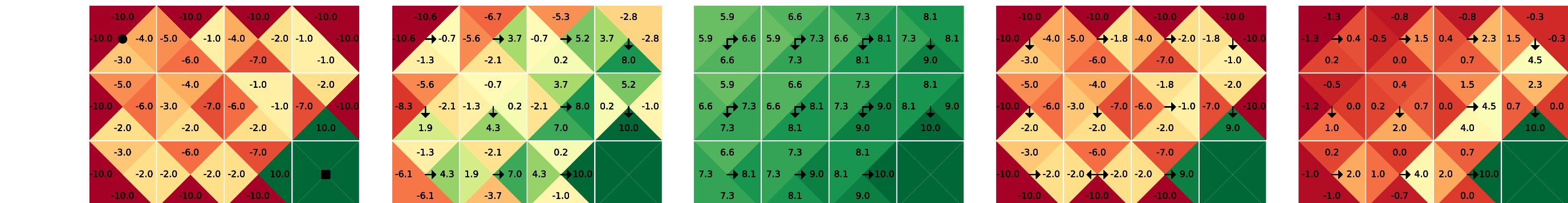


Theorem (Bellman equation for statistic function)

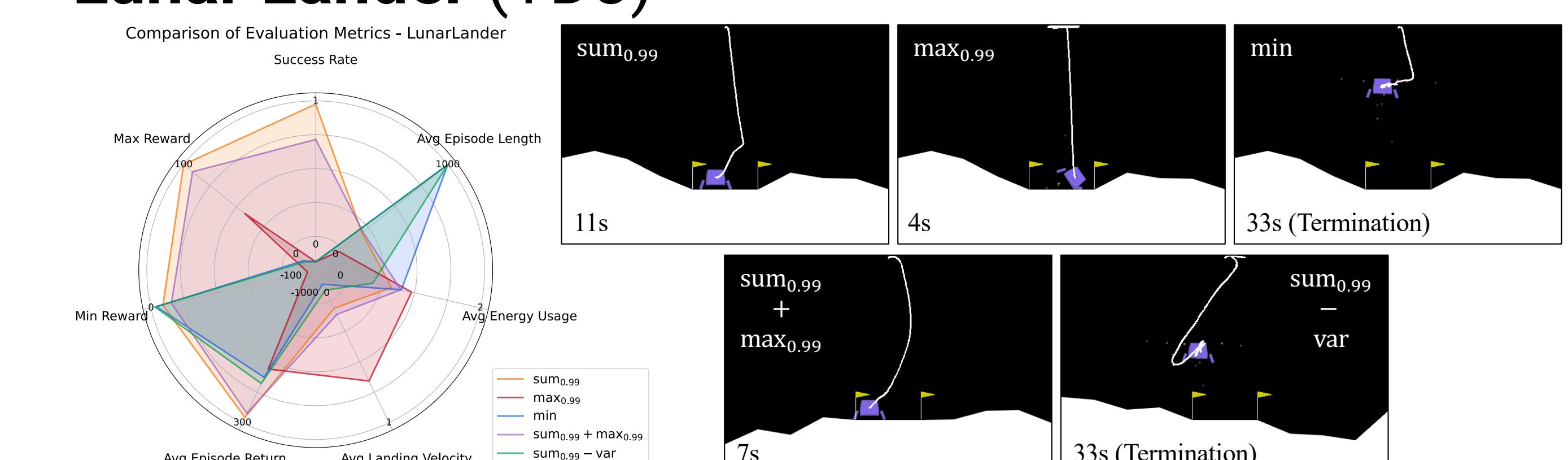
$$\mathcal{T}_\pi(s_t) = \begin{cases} \text{init} & s_t \text{ is terminal} \\ r_{t+1} \triangleright \mathcal{T}_\pi(s_{t+1}) & s_t \text{ is non-terminal} \end{cases}$$

Experiments

Grid-world (Q-learning)



Lunar Lander (TD3)



Future work

- Multi-dimensional or non-numerical feedback?
- Agent states? **Automata** as aggregators?
- List \rightarrow **tree**, list function \rightarrow **tree traversal**?
- Logic, reasoning, safety, and alignment?

References

- Wei Cui and Wei Yu. Reinforcement learning with non-cumulative objective. *Machine Learning in Communications and Networking*, 2023.
- Ralf Hinze, Thomas Harper, and Daniel W. H. James. Theory and practice of fusion. In *Symposium on Implementation and Application of Functional Languages*, 2010.
- William F. Sharpe. Mutual fund performance. *The Journal of Business*, 1966.