

Recursive Reward Aggregation

Yuting Tang Yivan Zhang Johannes Ackermann
Yu-Jie Zhang Soichiro Nishimori Masashi Sugiyama



Reinforcement Learning Conference 2025

In reinforcement learning, an agent obtains a *sequence* of rewards as it takes actions in a dynamic environment.

$$r_1 \quad r_2 \quad r_3 \quad \dots$$

We use the **discounted sum** to evaluate its performance.

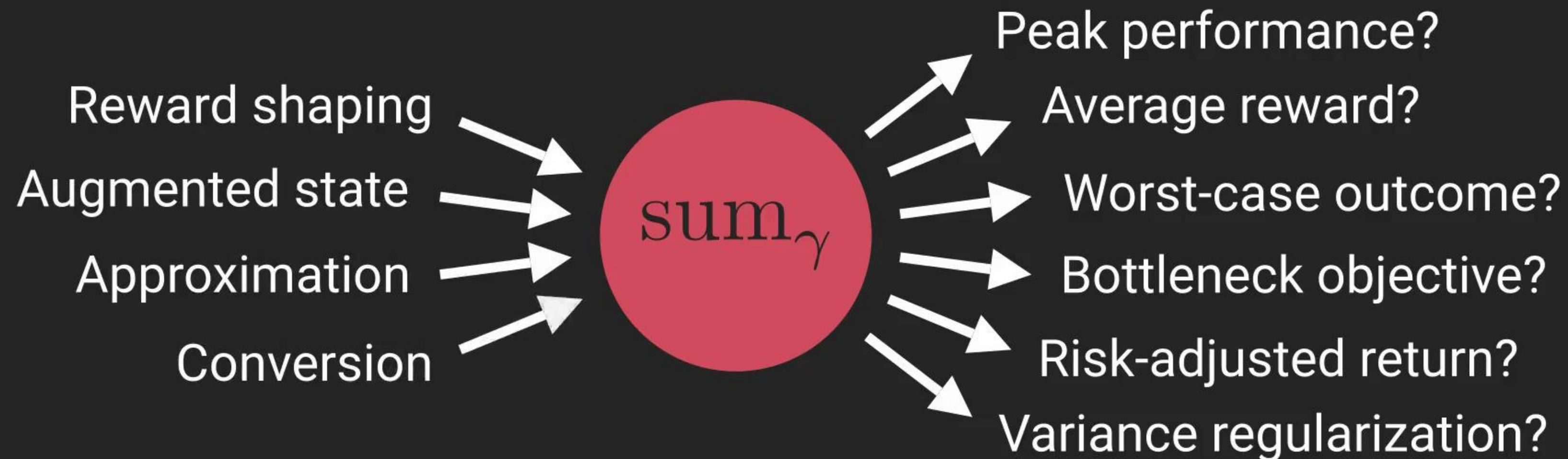
$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma r_2 + \gamma^2 r_3 + \dots \quad \gamma \in [0, 1]$$

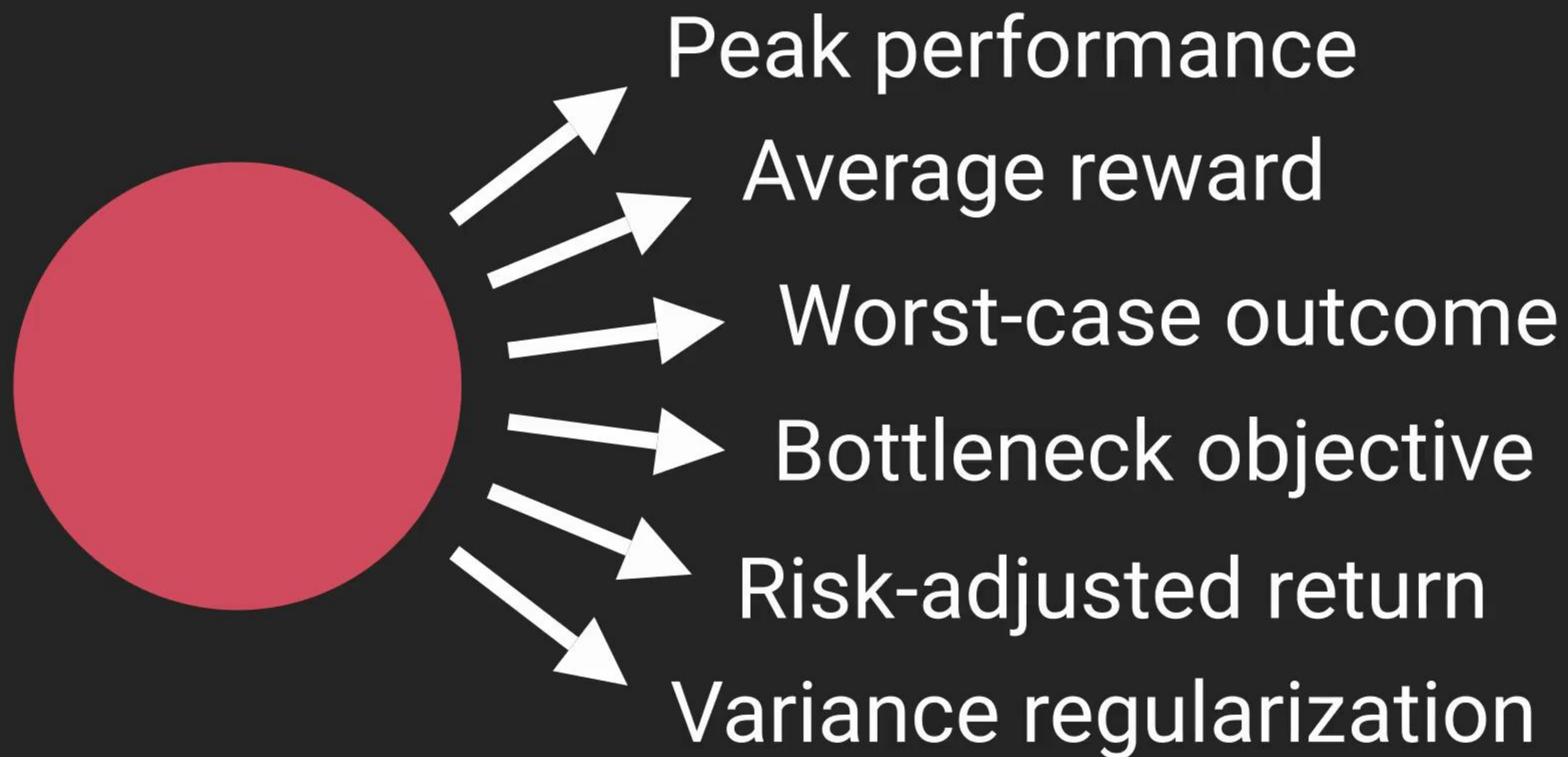
Why do we optimize the *discounted sum* of rewards?

Why do we optimize the *discounted sum* of rewards?

- Standard? Convenient? Mathematically elegant?
- Practical benefits? Higher, sooner rewards preferred
- Theoretical guarantees? Contraction \rightarrow unique fixed point
- Reward hypothesis? If true, all preferences can be represented

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





Can we optimize other reward aggregations,
directly, efficiently, and effectively?

Discounted sum is a *recursive function*

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

$$\text{sum}_\gamma [] = 0$$

Discounted sum is a *recursive function*

$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

$$\text{sum}_\gamma [] := 0$$

diagrammatic representation



Discounted sum is a *recursive function*

$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

$$\text{sum}_\gamma [] := 0$$

diagrammatic representation



Discounted sum is a *recursive function*

$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

$$\text{sum}_\gamma [] := 0$$

diagrammatic representation



Discounted sum is a *recursive function*

$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

$$\text{sum}_\gamma [] := 0$$

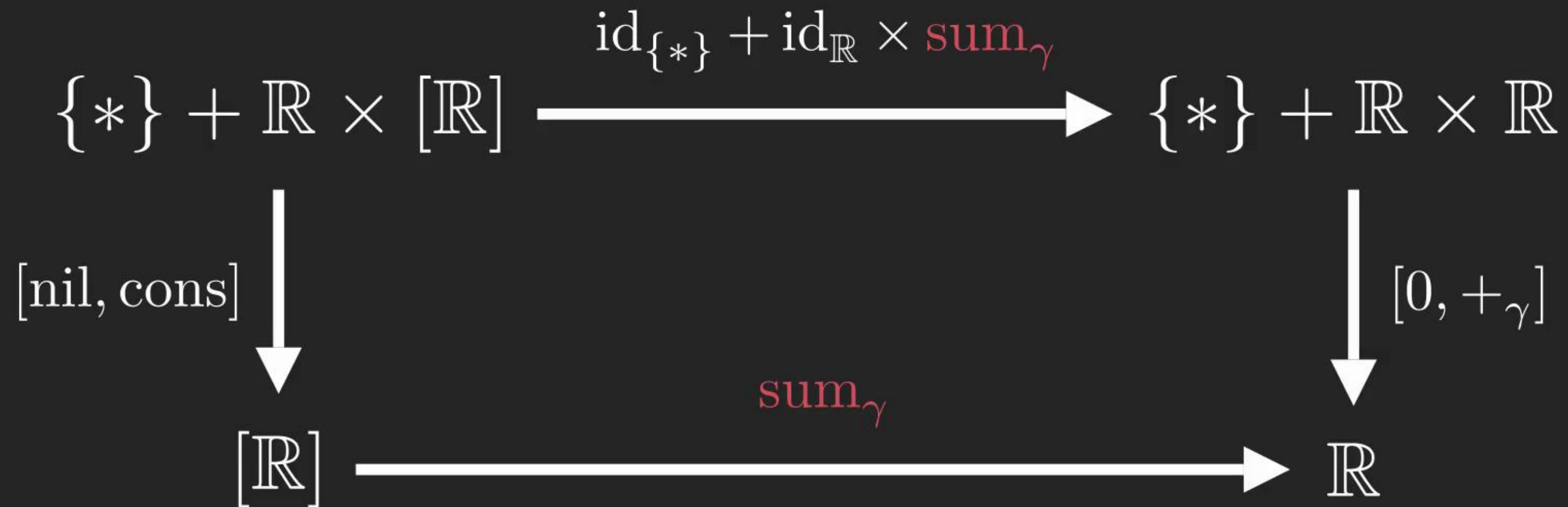
diagrammatic representation



Discounted sum is a *recursive function*

$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

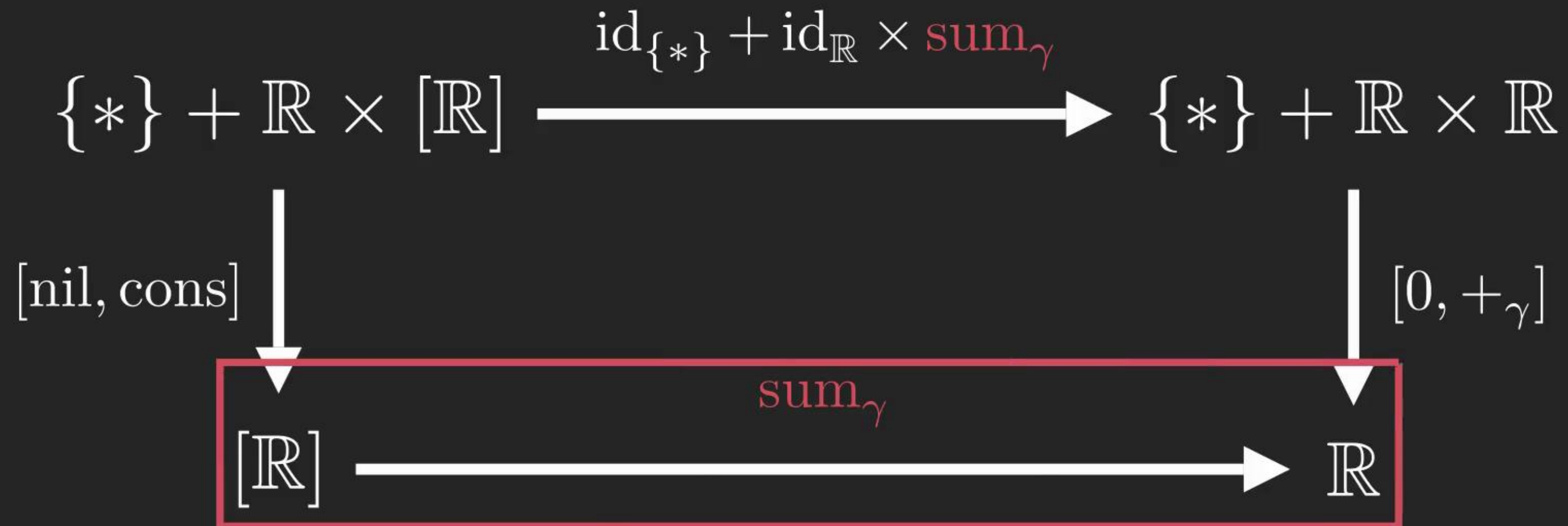
$$\text{sum}_\gamma [] := 0$$



$\{*\}$: singleton set $+$: disjoint union \times : Cartesian product
 id : identity function nil : empty list $[]$ cons : list constructor

Discounted sum is a *recursive function*

$$\begin{aligned}\text{sum}_\gamma [r_1, r_2, r_3, \dots] &:= r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots] \\ \text{sum}_\gamma [] &:= 0\end{aligned}$$

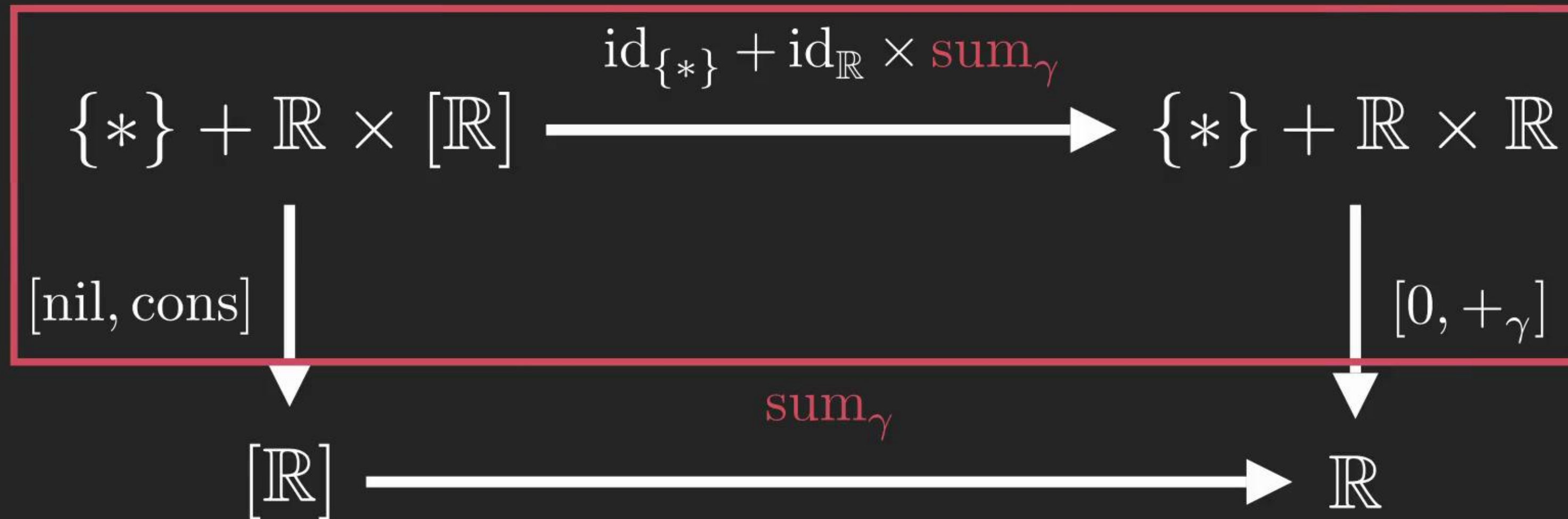


function declaration

$\{*\}$: singleton set $+$: disjoint union \times : Cartesian product
 id : identity function nil : empty list $[]$ cons : list constructor

Discounted sum is a *recursive function*

$$\begin{aligned}\text{sum}_\gamma [r_1, r_2, r_3, \dots] &:= r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots] \\ \text{sum}_\gamma [] &:= 0\end{aligned}$$



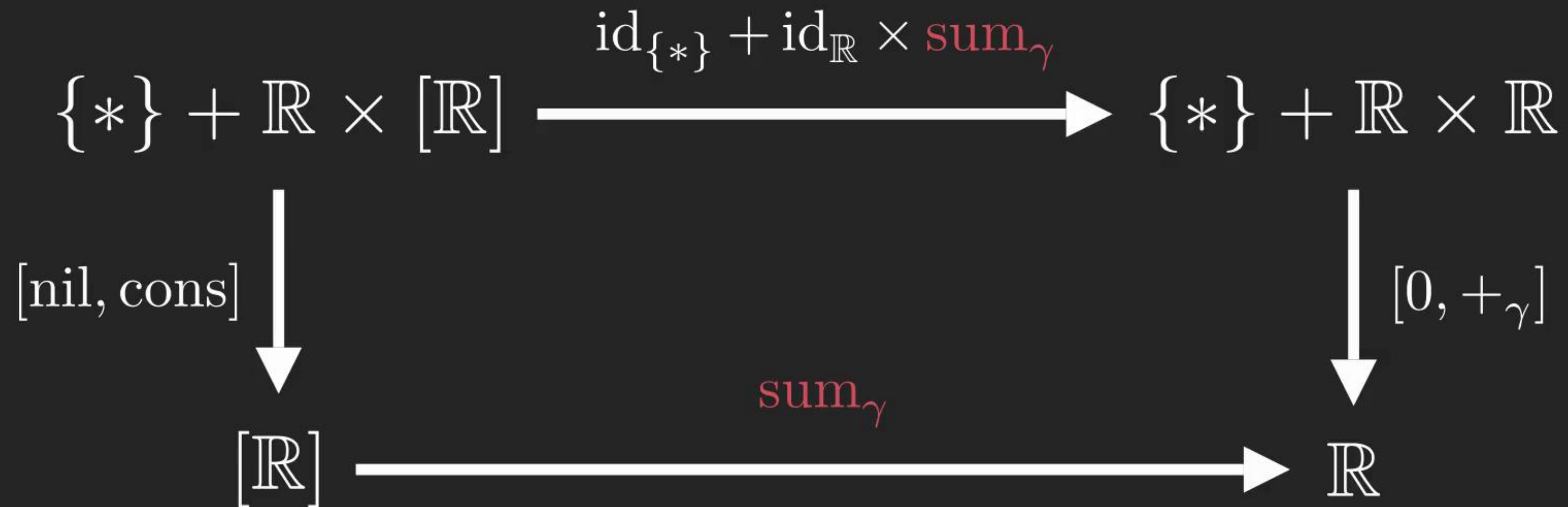
recursive definition

$\{*\}$: singleton set $+$: disjoint union \times : Cartesian product
 id : identity function nil : empty list $[]$ cons : list constructor

Discounted sum is a *recursive function*

$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

$$\text{sum}_\gamma [] := 0$$

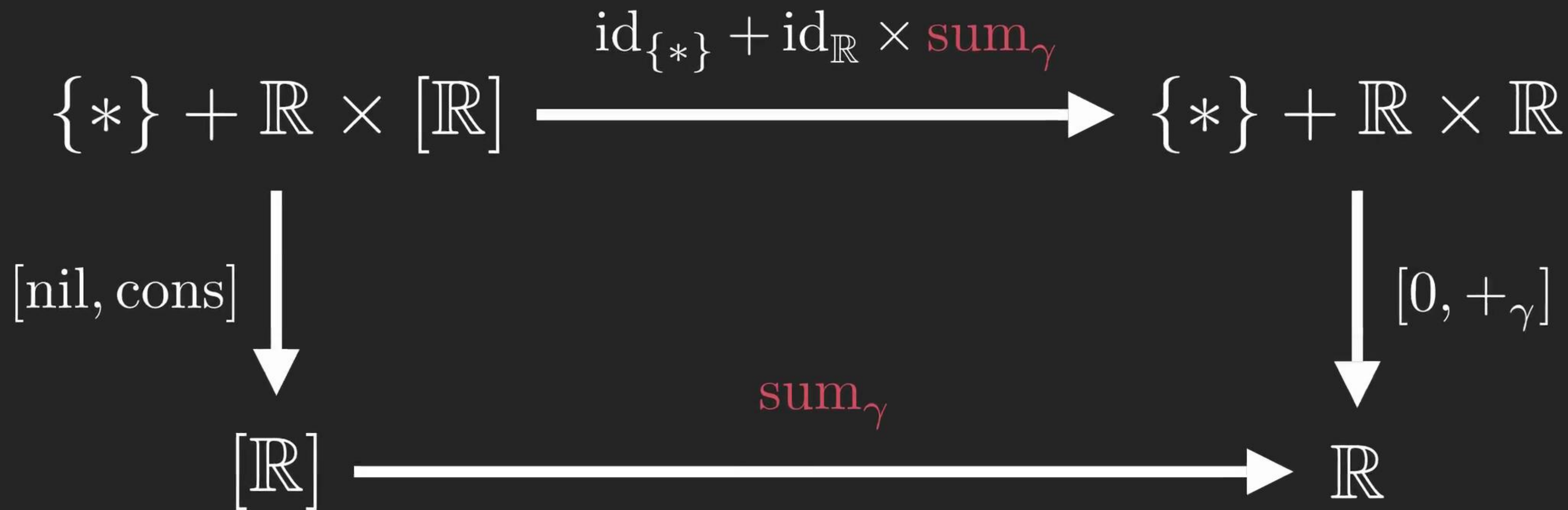


same endpoints \rightarrow same composite

$\{*\}$: singleton set $+$: disjoint union \times : Cartesian product
 id : identity function nil : empty list $[]$ cons : list constructor

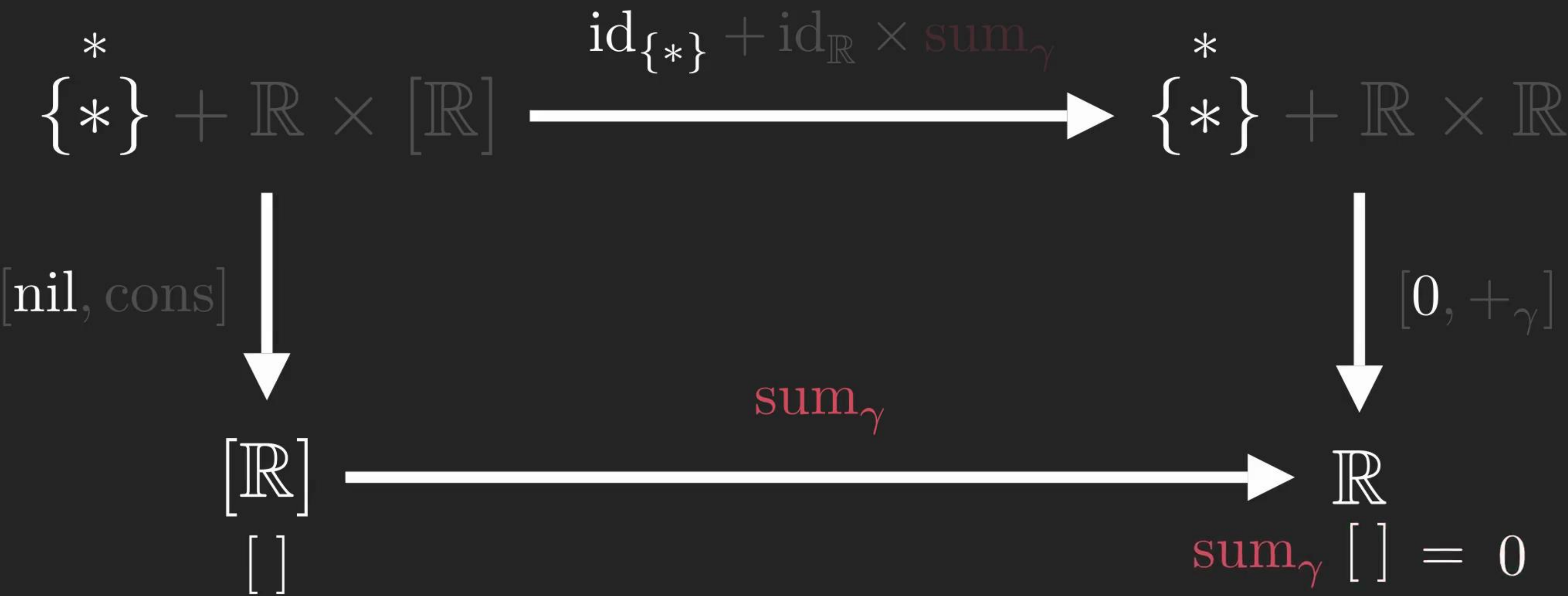
$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

$$\text{sum}_\gamma [] := 0$$



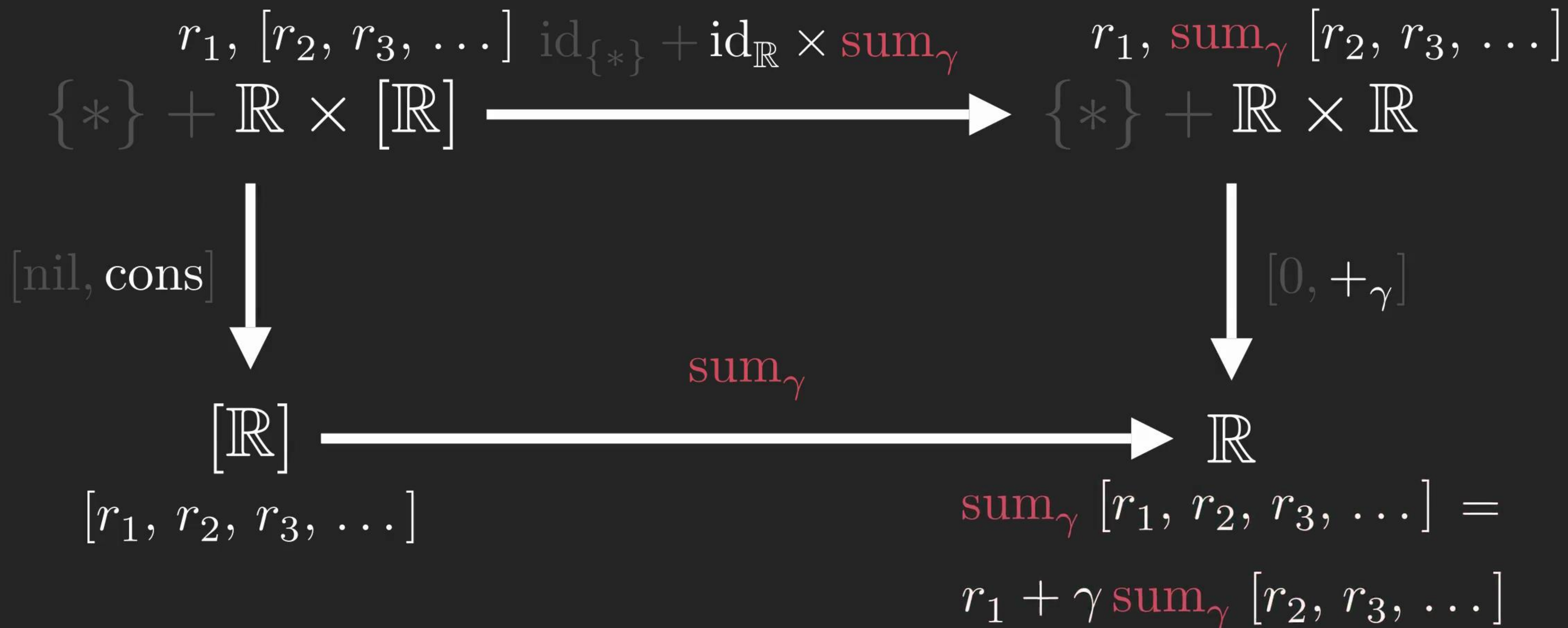
$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

$$\text{sum}_\gamma [] := 0$$



$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

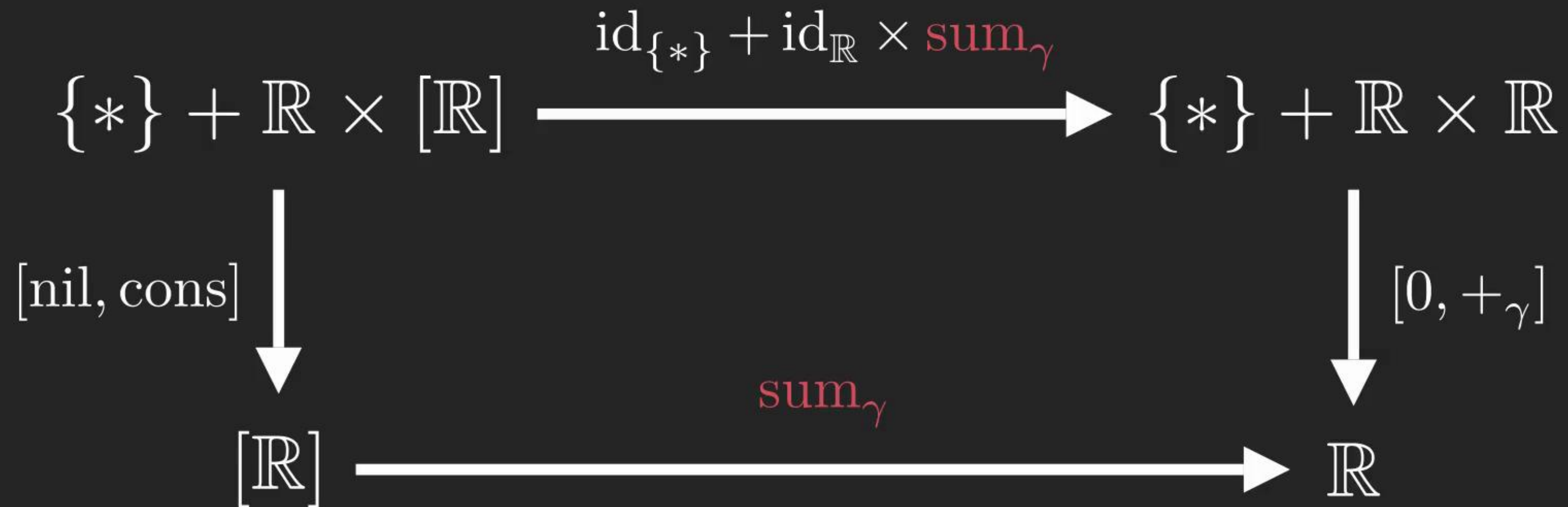
$$\text{sum}_\gamma [] := 0$$



Discounted sum is a *recursive function*

$$\text{sum}_\gamma [r_1, r_2, r_3, \dots] := r_1 + \gamma \text{sum}_\gamma [r_2, r_3, \dots]$$

$$\text{sum}_\gamma [] := 0$$

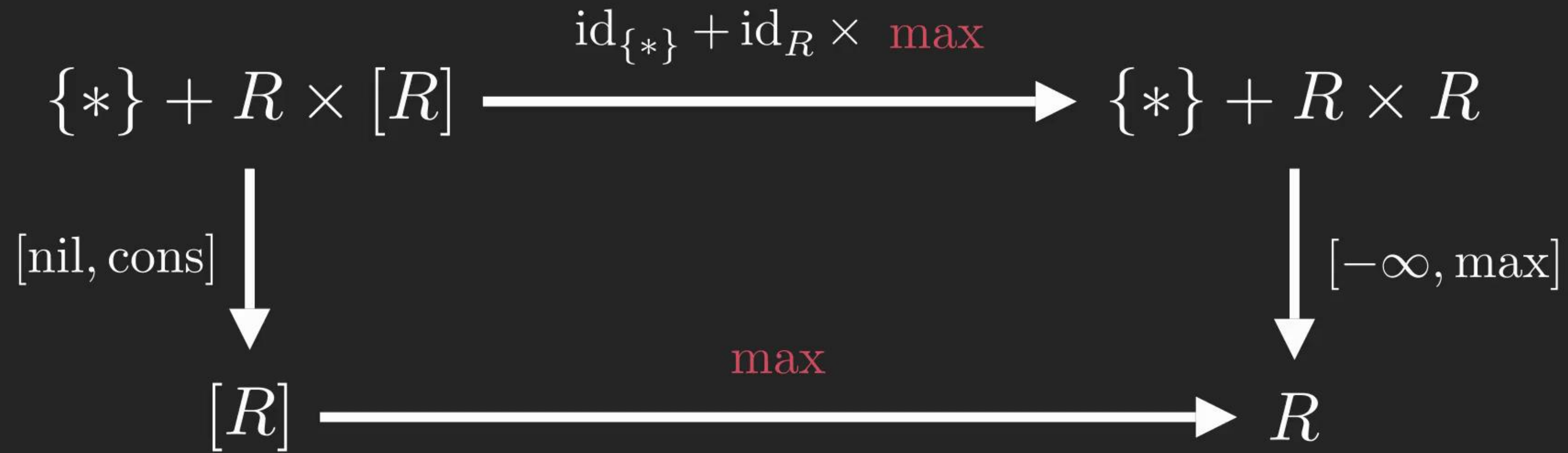


sum_γ is *uniquely defined* by $[0, +_\gamma]$ via recursion.

... so are many other aggregation functions.

$$\max [r_1, r_2, r_3, \dots] := \max(r_1, \max [r_2, r_3, \dots])$$

$$\max [] := -\infty$$

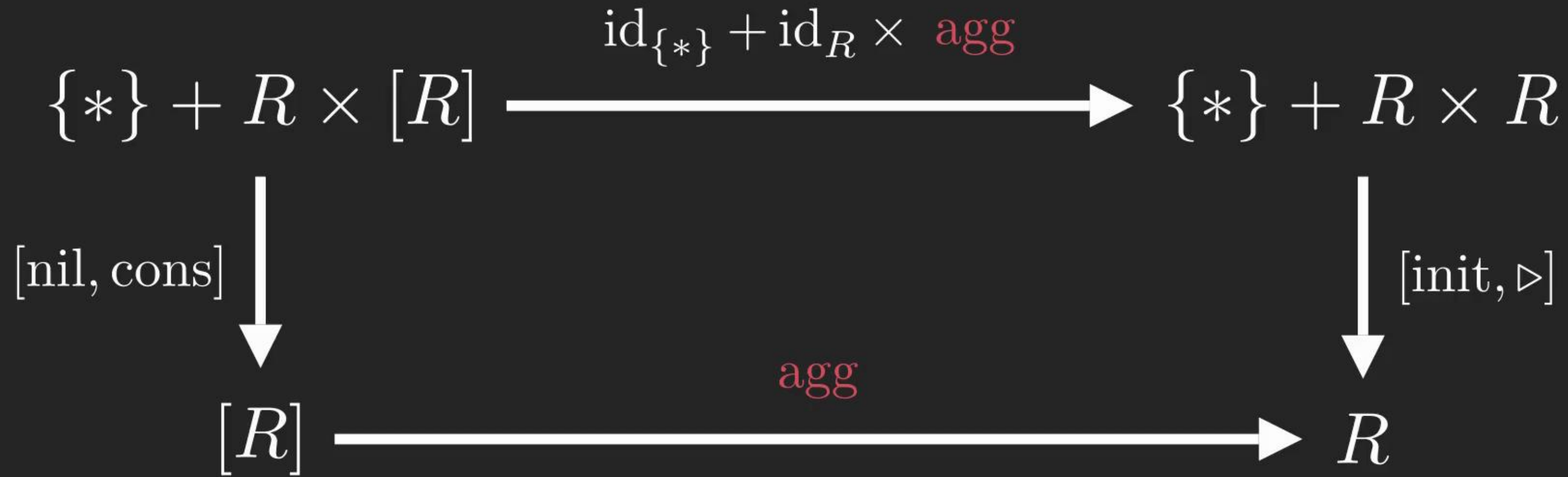


\max is *uniquely defined* by $[-\infty, \text{max}]$ via recursion.

... so are many other aggregation functions.

$$\text{agg } [r_1, r_2, r_3, \dots] := r_1 \triangleright \text{agg}[r_2, r_3, \dots]$$

$$\text{agg } [] := \text{init}$$



They all have the same *recursive structure*!

... so are many other aggregation functions.

$$\text{agg } [r_1, r_2, r_3, \dots] := r_1 \triangleright \text{agg}[r_2, r_3, \dots]$$

$$\text{agg } [] := \text{init}$$

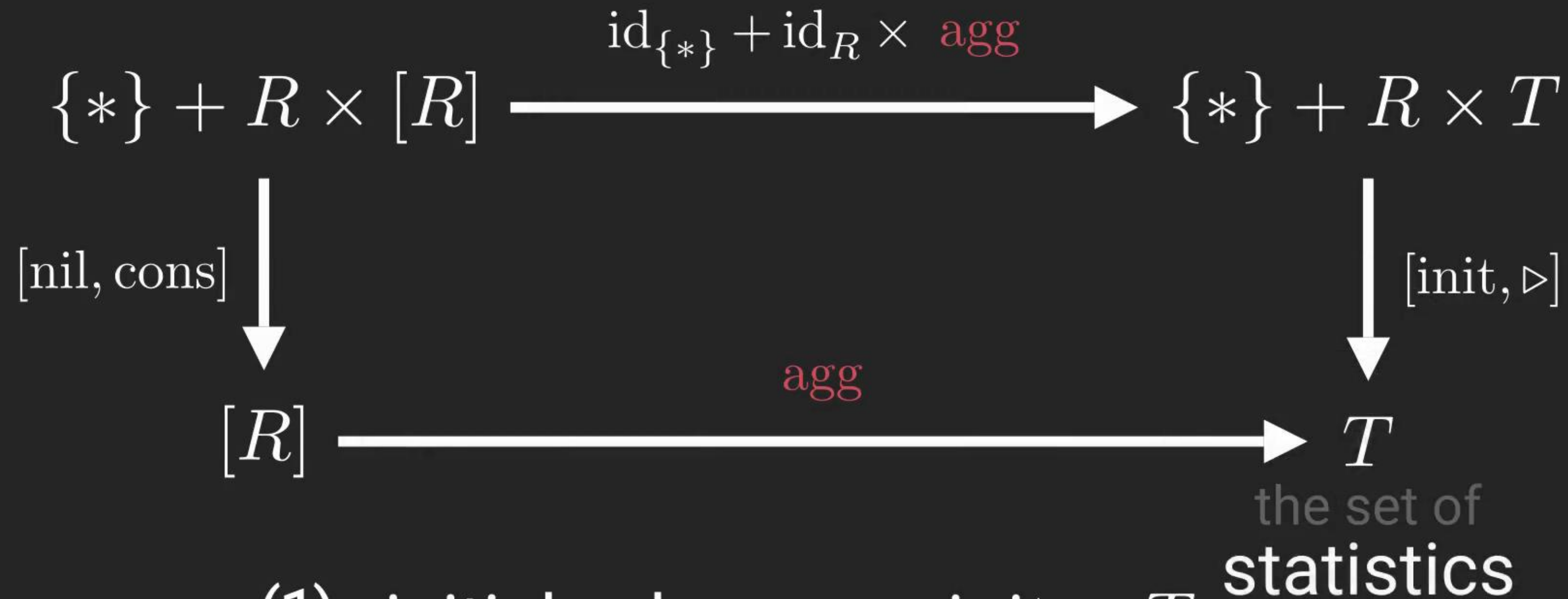


- (1) initial value $\text{init} \in R$
- (2) update function $\triangleright : R \times R \rightarrow R$

We can aggregate multiple statistics at the same time

$$\text{agg } [r_1, r_2, r_3, \dots] := r_1 \triangleright \text{agg}[r_2, r_3, \dots]$$

$$\text{agg } [] := \text{init}$$



- (1) initial value $\text{init} \in T$
- (2) update function $\triangleright : R \times T \rightarrow T$

We can aggregate multiple statistics at the same time

$$\text{agg } [r_1, r_2, r_3, \dots] := r_1 \triangleright \text{agg}[r_2, r_3, \dots]$$

$$\text{agg } [] := \text{init}$$



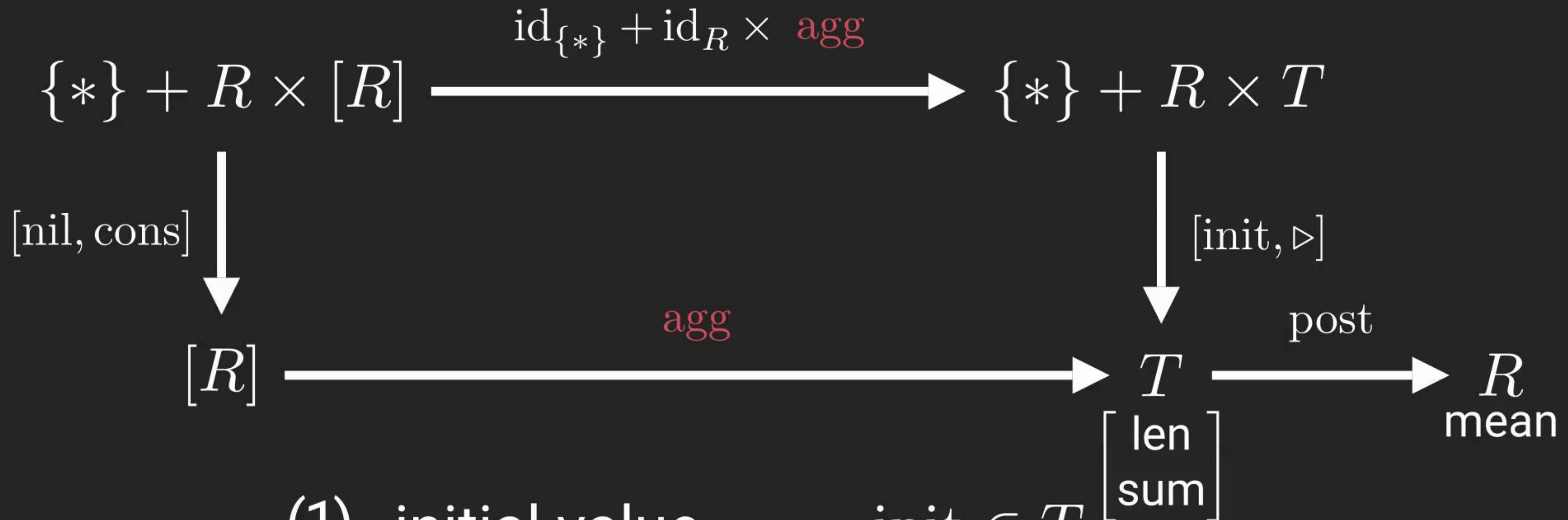
(1) initial value $\text{init} \in T$

(2) update function $\triangleright : R \times T \rightarrow T$

... and we can post-process the aggregated statistics.

$$\text{agg } [r_1, r_2, r_3, \dots] := r_1 \triangleright \text{agg}[r_2, r_3, \dots]$$

$$\text{agg } [] := \text{init}$$

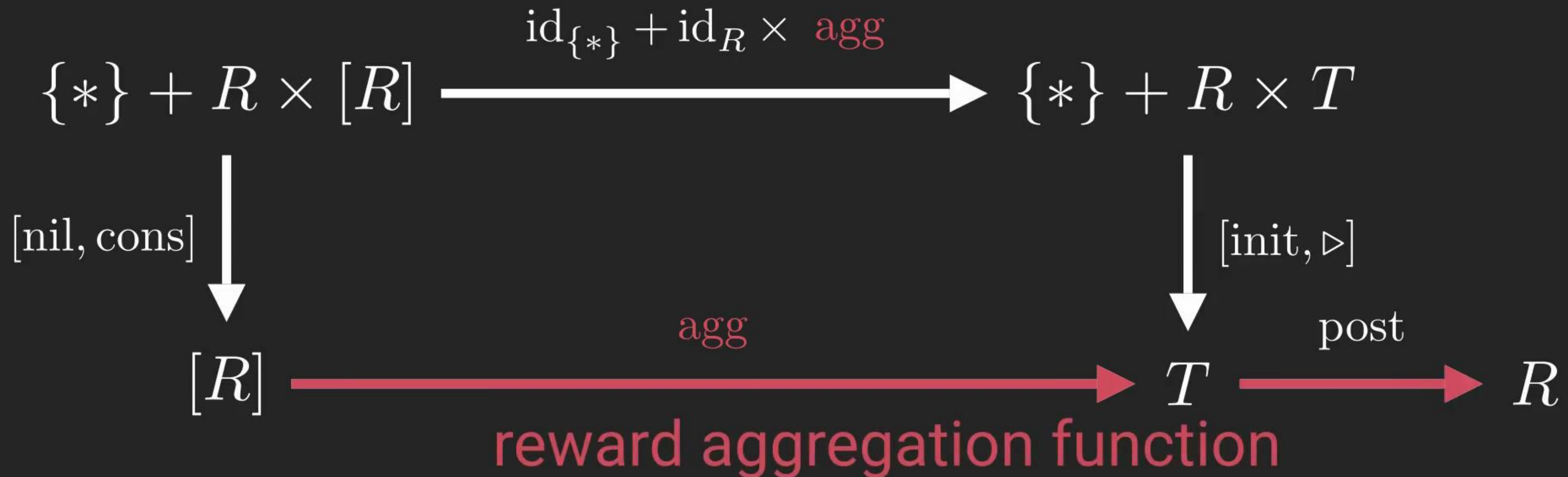


- (1) initial value $\text{init} \in T$
- (2) update function $\triangleright : R \times T \rightarrow T$
- (3) post-processing $\text{post} : T \rightarrow R$

... and we can post-process the aggregated statistics.

$$\text{agg } [r_1, r_2, r_3, \dots] := r_1 \triangleright \text{agg}[r_2, r_3, \dots]$$

$$\text{agg } [] := \text{init}$$



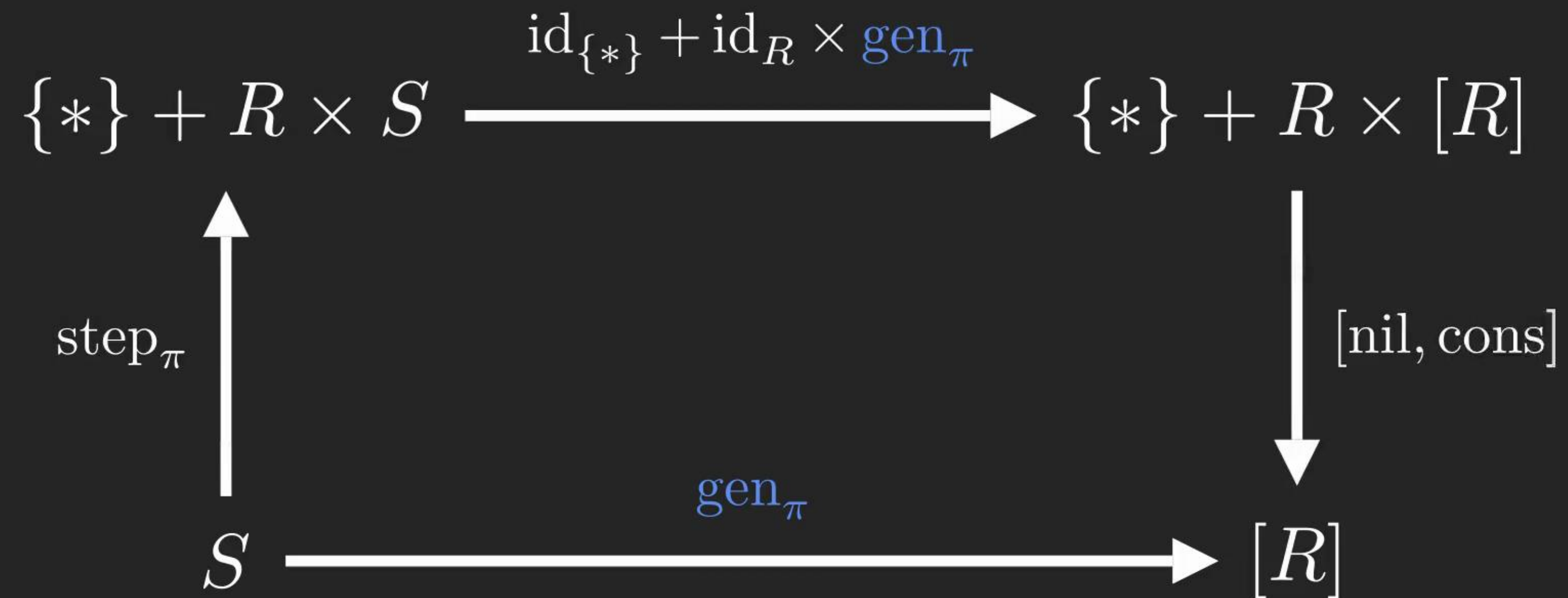
- (1) initial value $\text{init} \in T$
- (2) update function $\triangleright : R \times T \rightarrow T$
- (3) post-processing $\text{post} : T \rightarrow R$

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

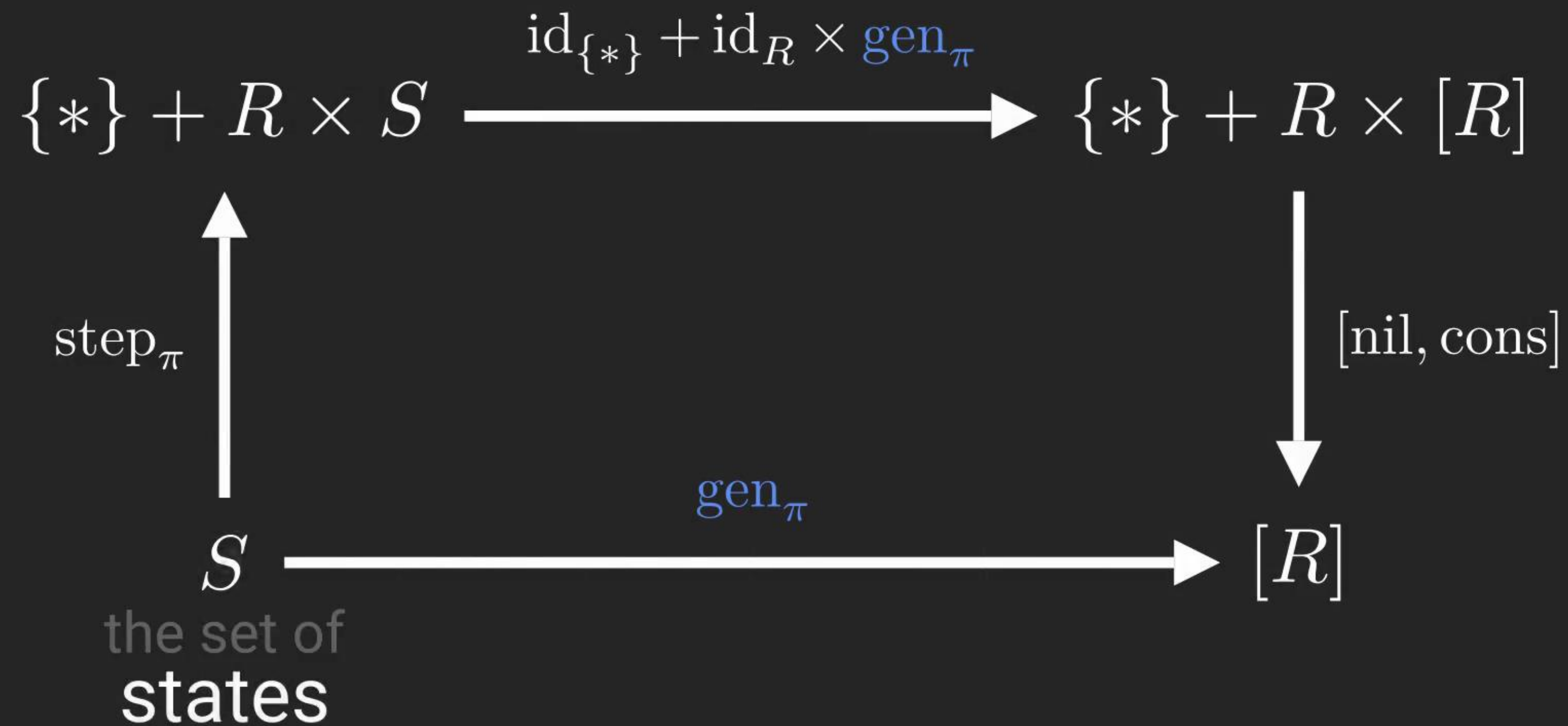
	definition	initial value of statistic(s)	update function	post-processing
	$\text{post} \circ \text{agg}_{\text{init}, \triangleright} : [R] \rightarrow R$	$\text{init} \in T$	$\triangleright : R \times T \rightarrow T$	$\text{post} : T \rightarrow R$
discounted sum	$r_1 + \gamma r_2 + \dots + \gamma^{t-1}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

Why does this matter?

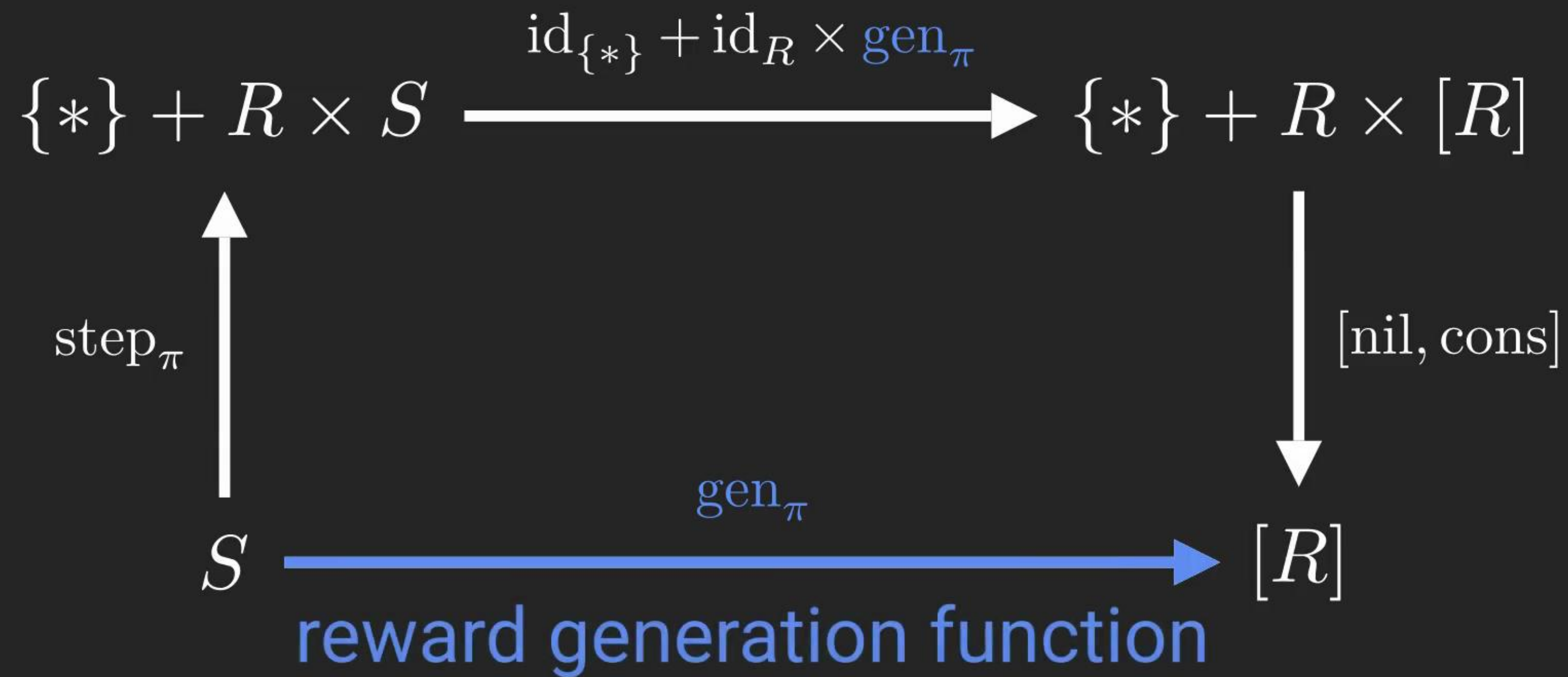
Reward generation is also a *recursive function*



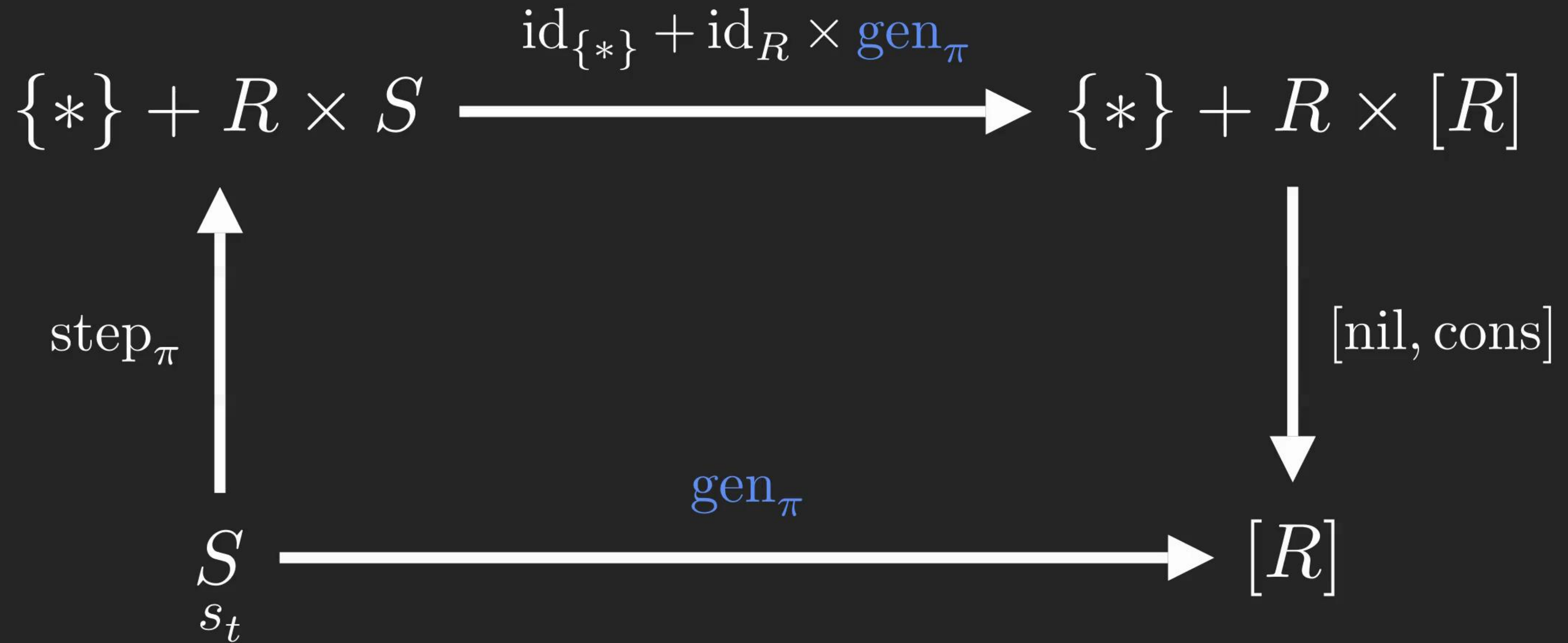
Reward generation is also a *recursive function*



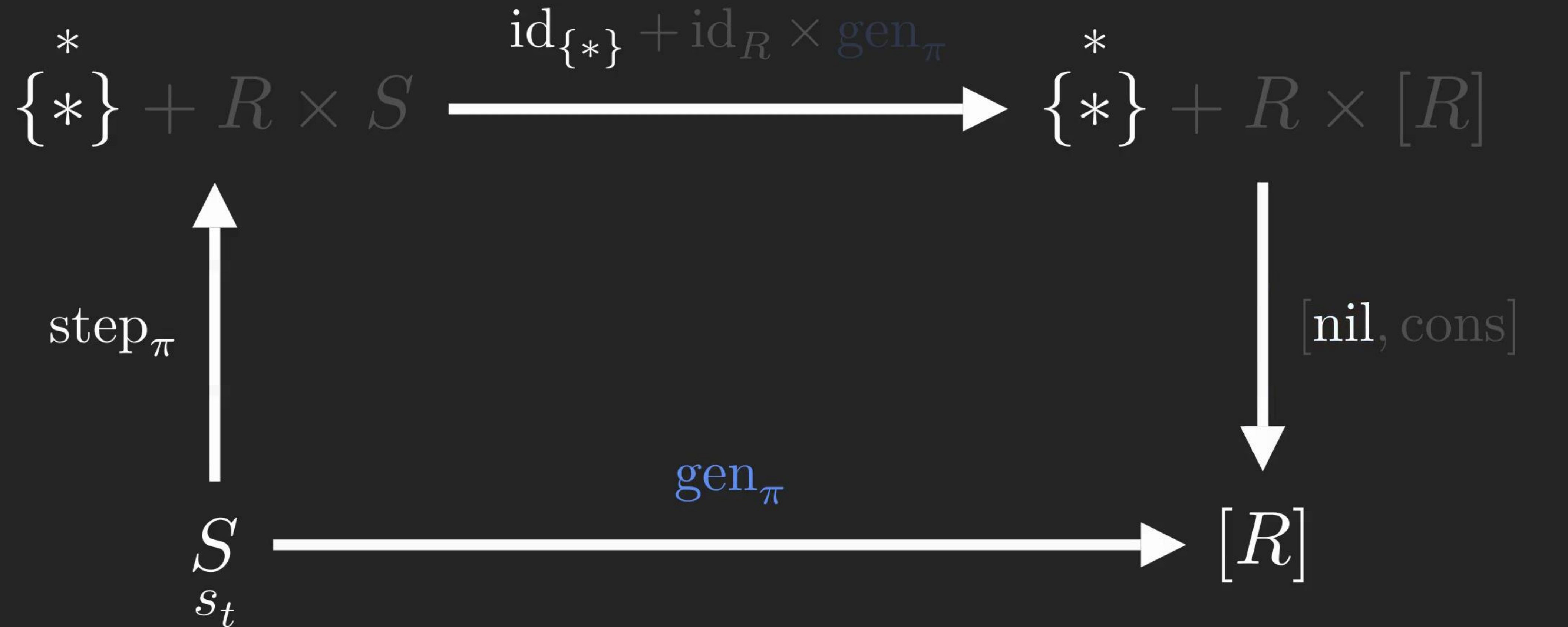
Reward generation is also a *recursive function*



$$\begin{array}{ccc}
 \{*\} + R \times S & \xrightarrow{\text{id}_{\{*\}} + \text{id}_R \times \text{gen}_\pi} & \{*\} + R \times [R] \\
 \uparrow \text{step}_\pi & & \downarrow [\text{nil}, \text{cons}] \\
 S & \xrightarrow{\text{gen}_\pi} & [R]
 \end{array}$$

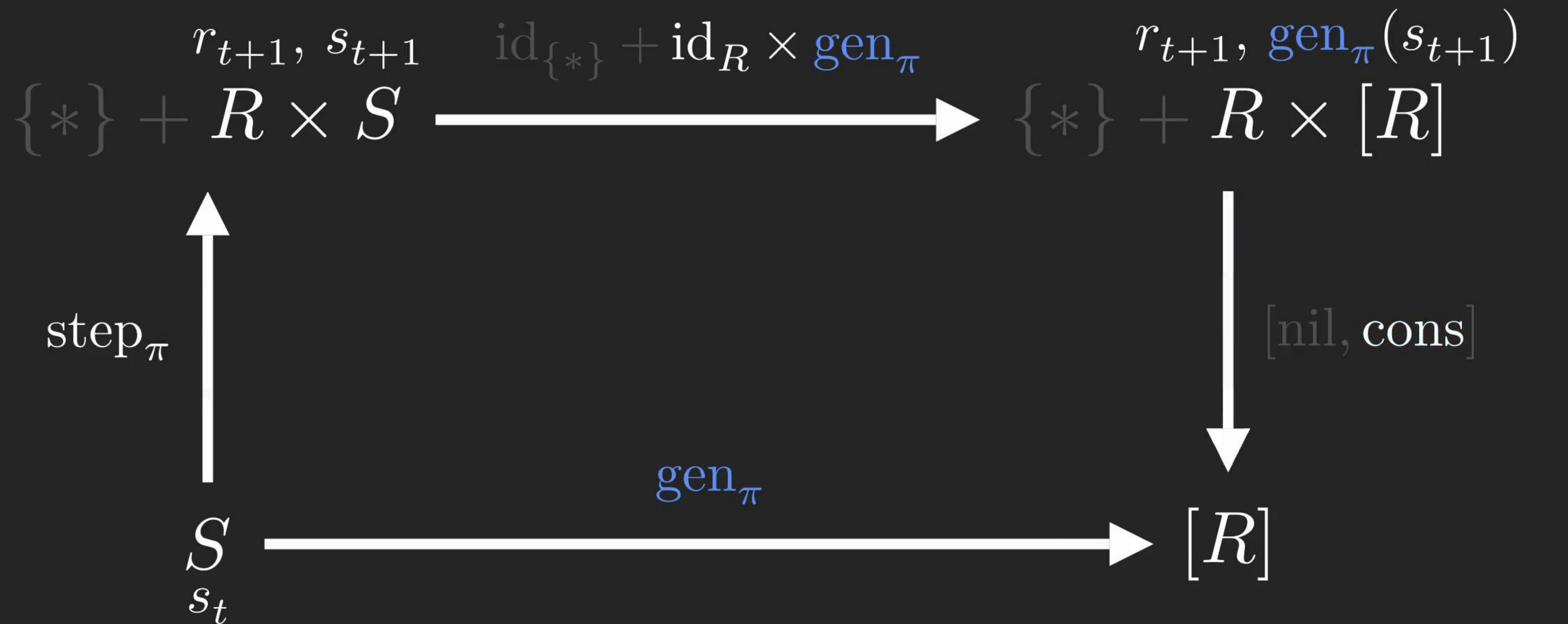


$$\text{gen}_\pi(s_t) = \left\{ \right.$$



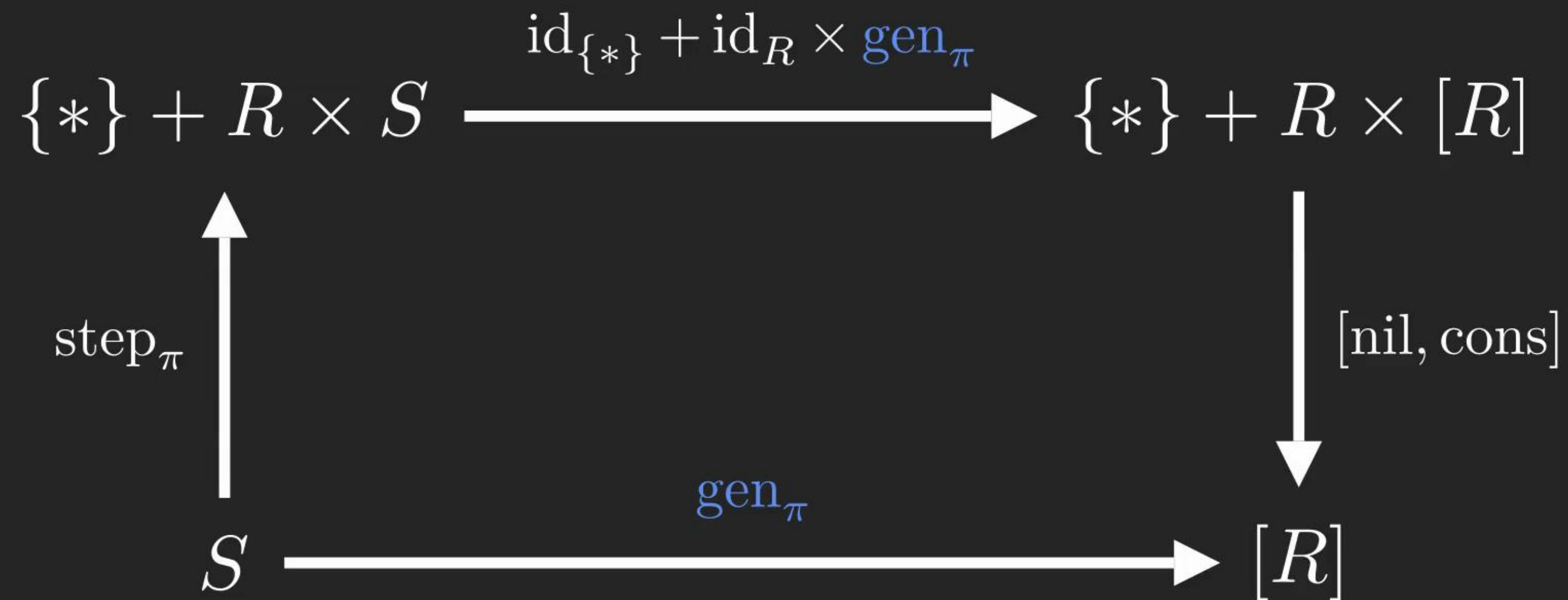
$$\text{gen}_\pi(s_t) = \left\{ \begin{array}{l} [] \end{array} \right.$$

terminal



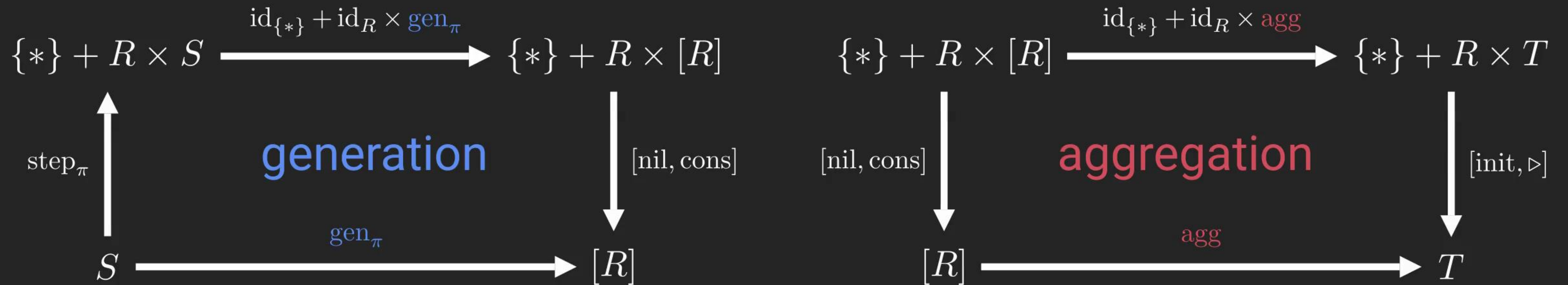
$$\text{gen}_\pi(s_t) = \begin{cases} [] & \text{terminal} \\ r_{t+1} : \text{gen}_\pi(s_{t+1}) & \text{otherwise} \end{cases}$$

Reward generation is also a *recursive function*



gen_π is *uniquely defined* by step_π via recursion.

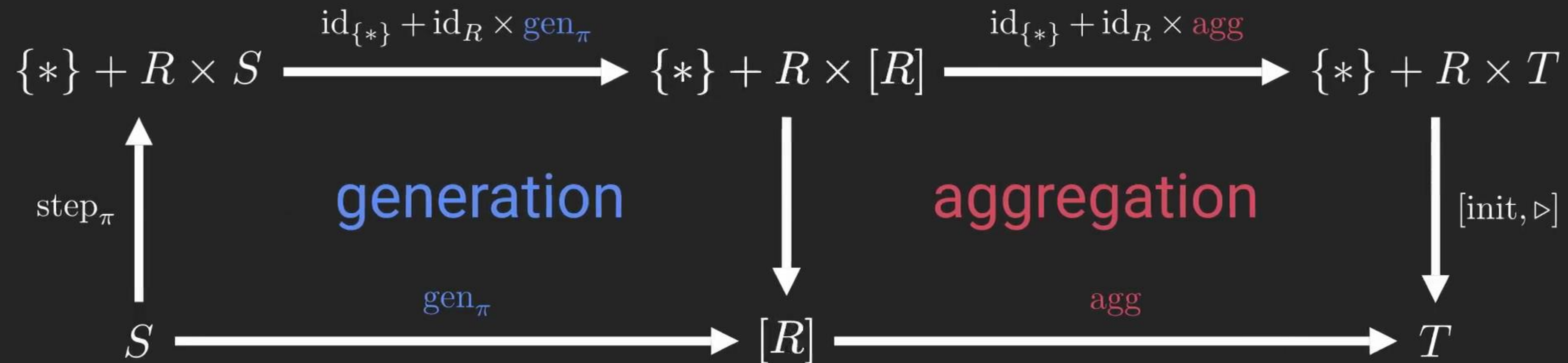
Recursive generation and aggregation of rewards



gen_π is *uniquely defined* by step_π via recursion.

agg is *uniquely defined* by $[\text{init}, \triangleright]$ via recursion.

Recursive generation and aggregation of rewards lead to generalized Bellman equations.



gen_π is *uniquely defined* by step_π via recursion.

agg is *uniquely defined* by $[\text{init}, \triangleright]$ via recursion.

Recursive generation and aggregation of rewards
lead to generalized Bellman equations.

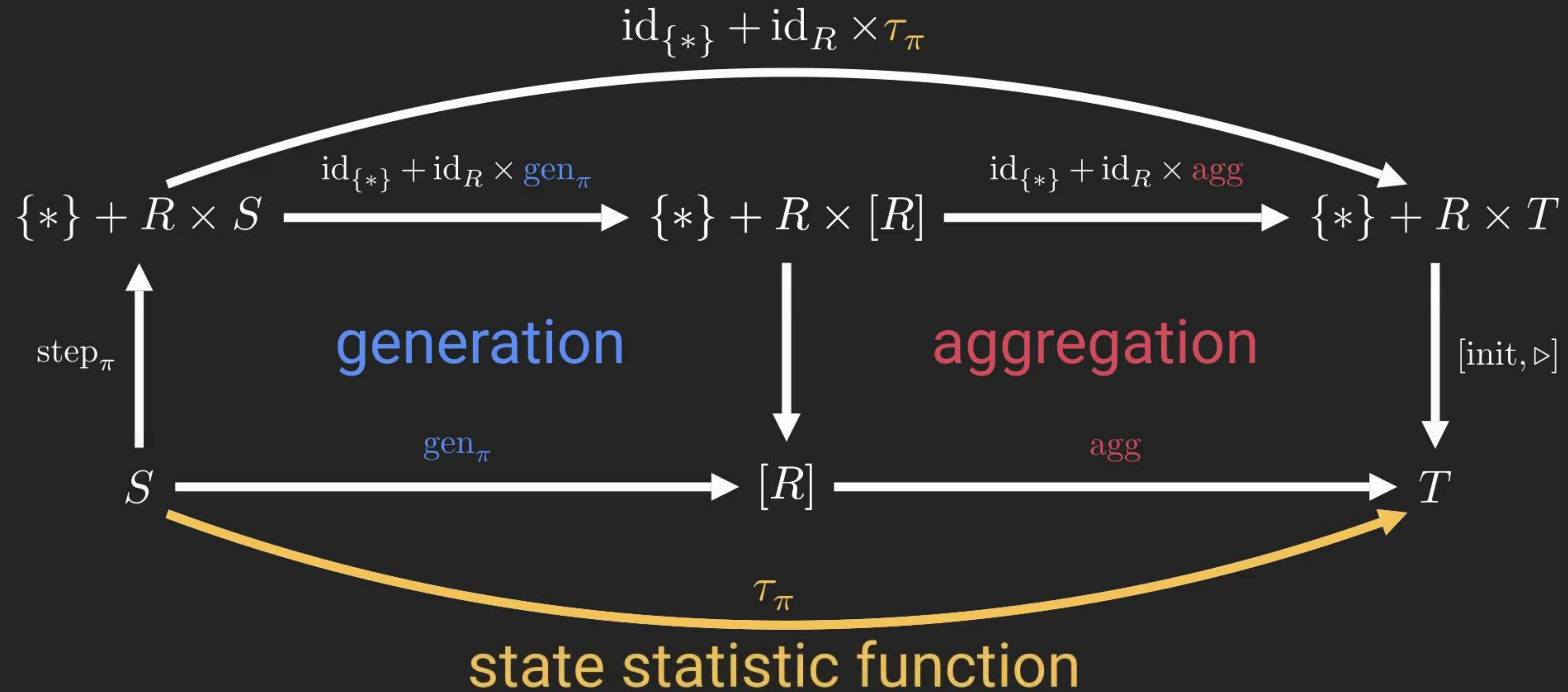


Hinze et al. (2010)

gen_π is *uniquely defined* by step_π via recursion.

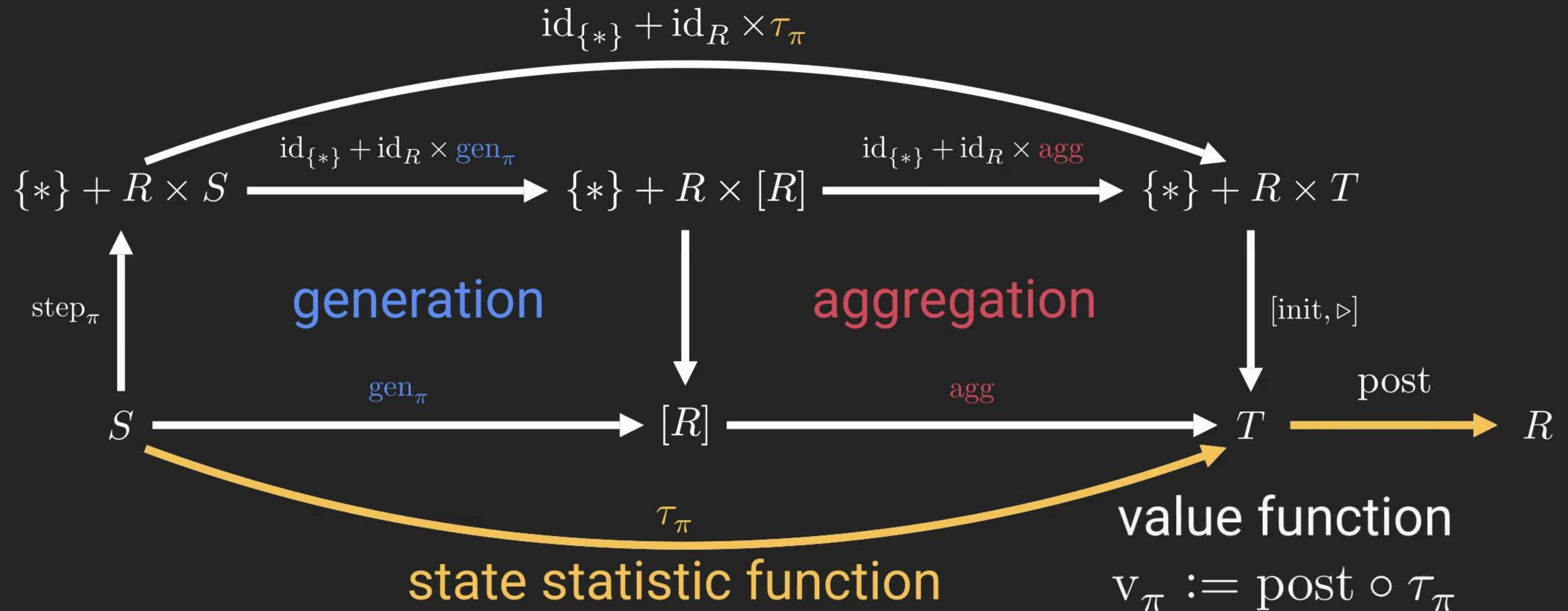
agg is *uniquely defined* by $[\text{init}, \triangleright]$ via recursion.

Recursive generation and aggregation of rewards lead to generalized Bellman equations.



$\tau_\pi = \text{agg} \circ \text{gen}_\pi$ is uniquely defined by step_π and $[\text{init}, \triangleright]$ via recursion.

Recursive generation and aggregation of rewards lead to generalized Bellman equations.



$\tau_\pi = \text{agg} \circ \text{gen}_\pi$ is uniquely defined by step_π and $[\text{init}, \triangleright]$ via recursion.

Theorem (Bellman equation for state statistic function)

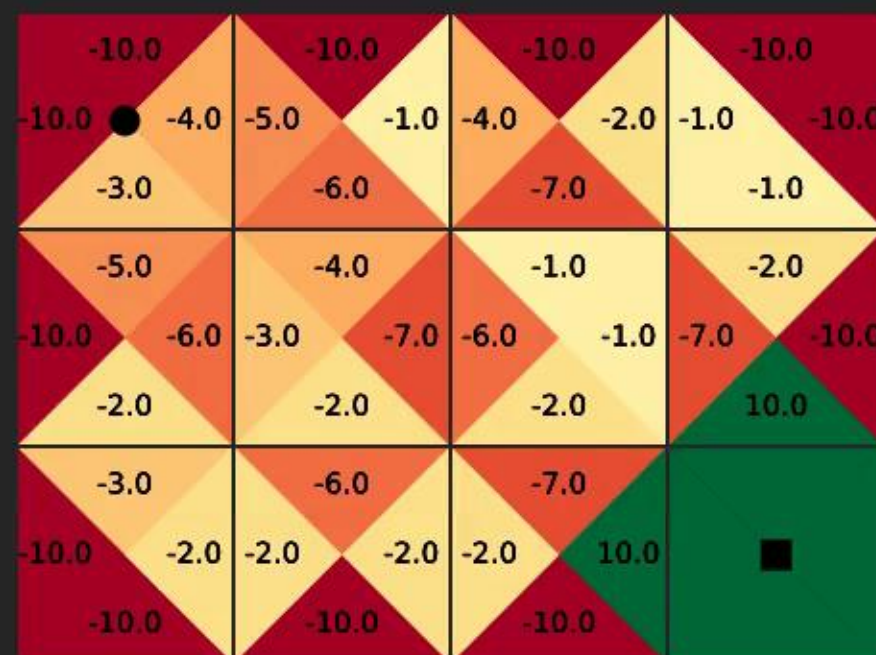
$$\tau_{\pi}(s_t) = \begin{cases} \text{init} & \text{terminal} \\ r_{t+1} \triangleright \tau_{\pi}(s_{t+1}) & \text{otherwise} \end{cases}$$

For *theorists*:

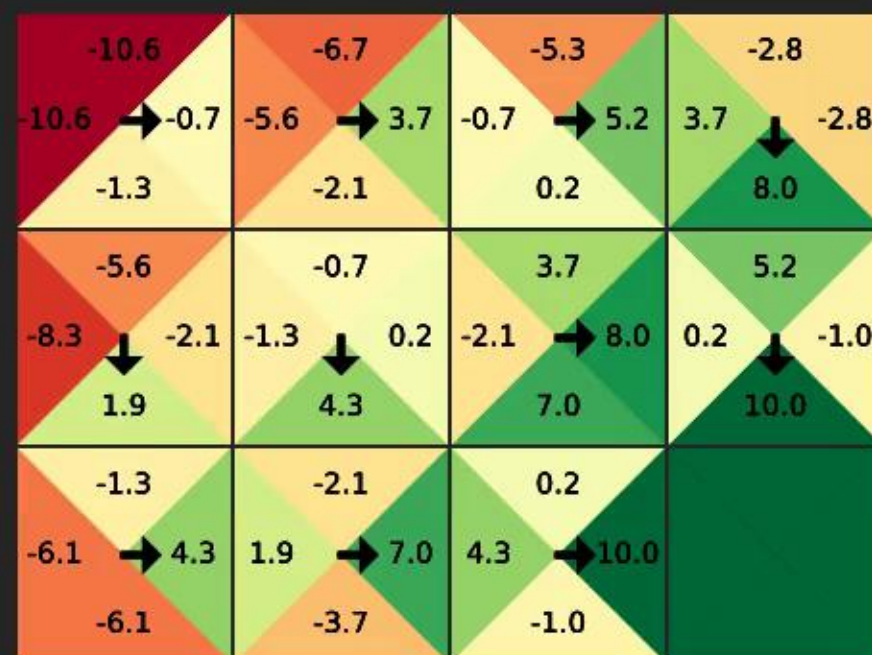
- Recursive generation and aggregation \rightarrow Bellman equations
- Preorder and premetric structures \rightarrow convergence behaviors
- Unified deterministic and stochastic formulations

For *practitioners*:

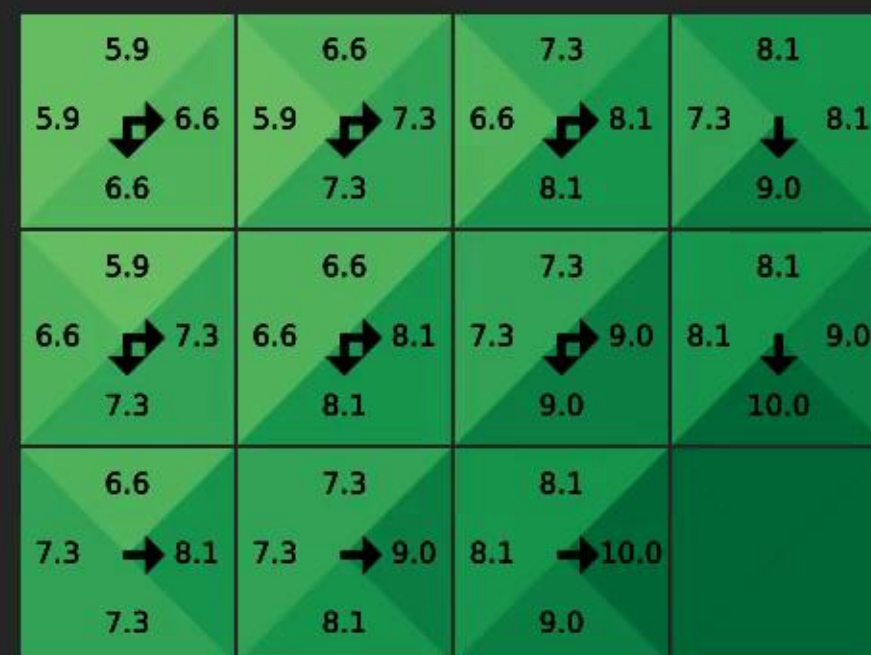
- You can use your favorite value/critic-based algorithms.
- You can directly optimize the *Sharpe ratio* ($\frac{\text{mean}}{\text{std}}$) in finance, or regularize the *velocity range* ($\text{max} - \text{min}$) in continuous control.



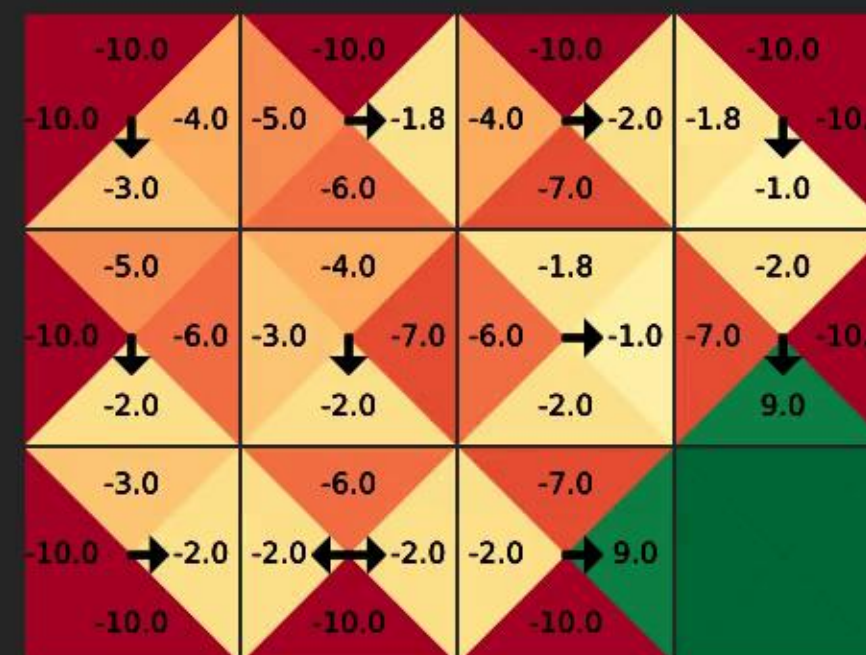
Grid



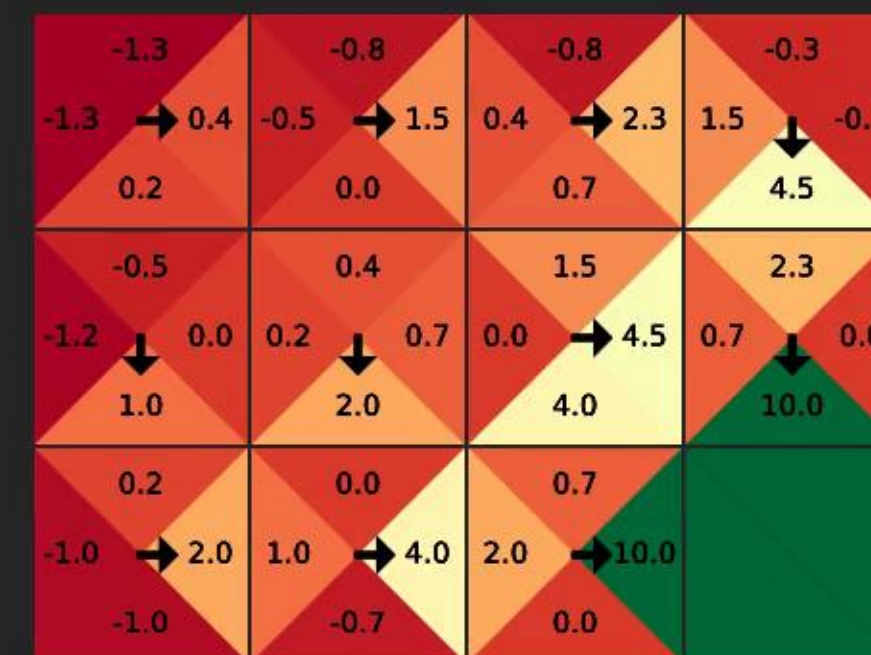
$\text{sum}_{0.99}$



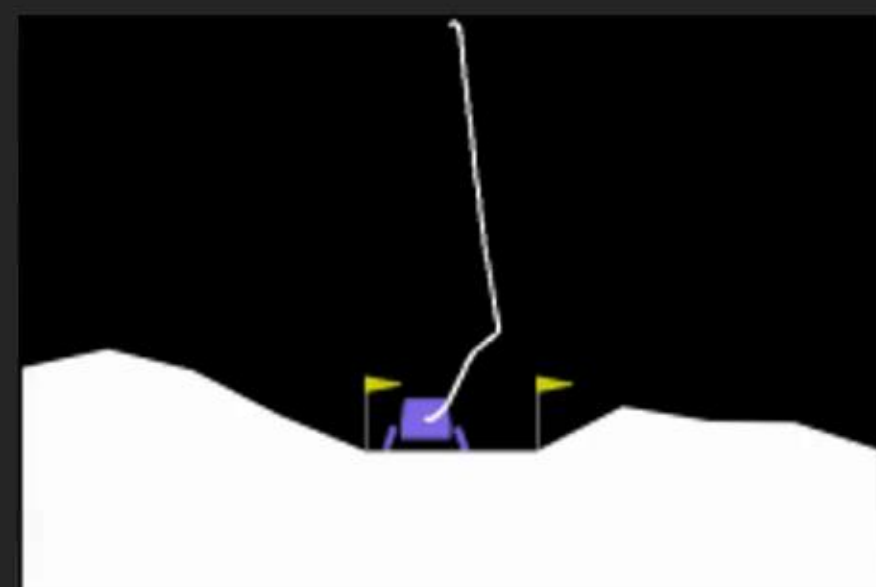
$\text{max}_{0.99}$



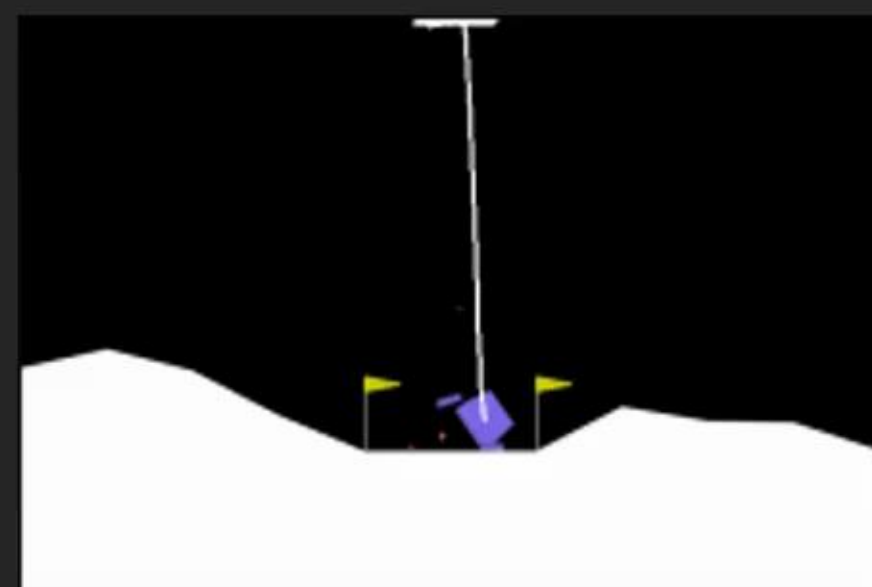
$\text{min}_{0.99}$



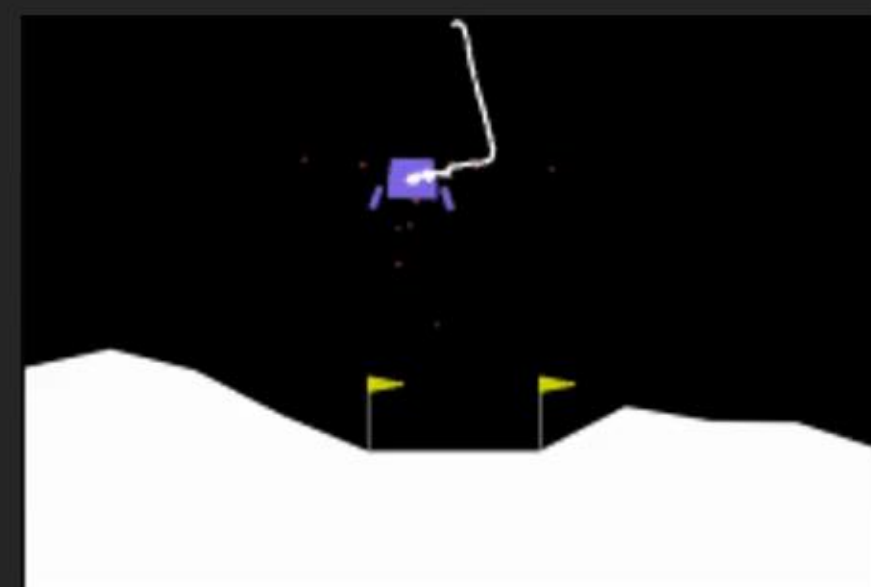
mean



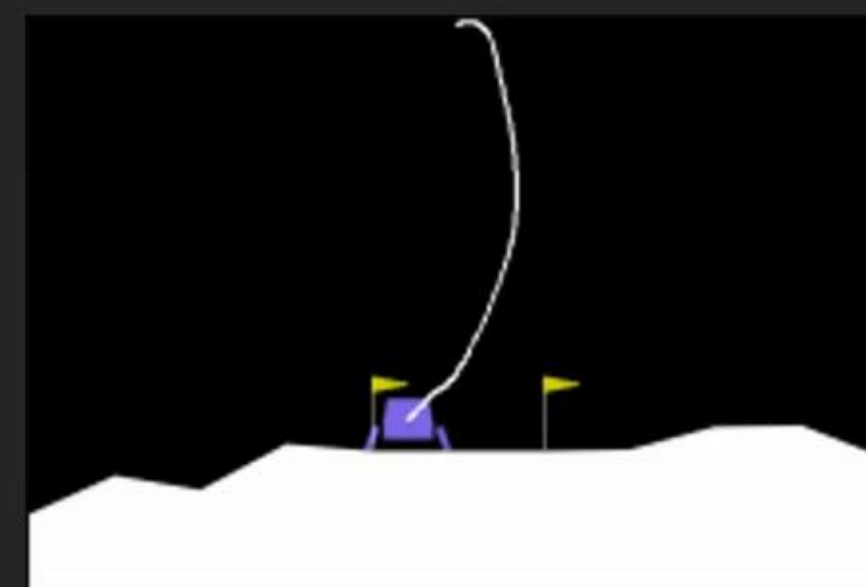
$\text{sum}_{0.99}$



$\text{max}_{0.99}$



min



$\text{sum}_{0.99} + \text{max}_{0.99}$



$\text{sum}_{0.99} - \text{var}$

Q-learning, PPO, TD3, ...

Mean, variance, range, the Sharpe ratio in portfolio optimization, ...

Recursive Reward Aggregation

Summary

- An algebraic perspective on Markov decision processes
- Generalized Bellman equations and Bellman operators
- Integration into value-based and actor-critic algorithms

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?

Recursive Reward Aggregation

Yuting Tang Yivan Zhang Johannes Ackermann
Yu-Jie Zhang Soichiro Nishimori Masashi Sugiyama



Reinforcement Learning Conference 2025