

UCB-driven Utility Function Search for Multi-objective Reinforcement Learning Supplementary Material

May 28, 2025

1 Algorithms

Algorithm 1 Fixed-MOPPO

```

1: Input: State  $s_t$ , weights  $\mathbf{w}$ 
2: Initialize:  $K$  weight-conditioned Actor-Critic Networks  $\pi_k / v^{\pi_k}$ , scalarisation-vector sub-
   spaces  $\mathbf{W}_k$  for  $k = 1, \dots, K$ , and memory buffer  $\mathcal{E}$  size of  $D$ .
3: for  $k = 1$  to  $K$  do
4:   for  $t = 1$  to  $D$  do
5:      $\mathbf{w}_{pivot} \leftarrow \text{get pivot weight}(\mathbf{W}_k)$ 
6:      $a_t \leftarrow \pi_k(s_t, \mathbf{w}_{pivot})$ 
7:      $s_{t+1}, \mathbf{r}_t \leftarrow \text{simulator}(a_t)$ 
8:      $\mathcal{E} \leftarrow \mathcal{E} \cup \langle s_t, a_t, \mathbf{w}_{pivot}, \mathbf{r}_t, s_{t+1} \rangle$ 
9:      $s_t \leftarrow s_{t+1}$ 
10:   end for
11:   sample  $\langle s_t, a_t, \mathbf{w}_{pivot}, \mathbf{r}_t, s_{t+1} \rangle \leftarrow \mathcal{E}$ 
12:    $\theta \leftarrow \theta + \eta (\nabla_{\theta} \log \pi_{\theta}(s_t, a; \mathbf{w}_{pivot})) (A^{\pi}(s_t, a; \mathbf{w}_{pivot}))$ 
13:    $\phi \leftarrow \phi + ||\mathbf{V}^{\pi_k}(s_t; \mathbf{w}_{pivot}) - \mathbf{V}^{\pi_k}(s_{t+1}; \mathbf{w}_{pivot})||^2$ 
14:   clear  $\mathcal{E}$ 
15: end for

```

2 Benchmark Problems

This section provide detail objective return for each problems. Where C in the following equations is live bonus.

2.1 Swimmer-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^8, \mathcal{A} \in \mathbb{R}^2$.

The first objective is forward speed in x axis:

$$R_1 = v_x \tag{1}$$

The second objective is energy efficiency:

$$R_2 = 0.3 - 0.15 \sum_i a_i^2, \quad a_i \in (-1, 1) \tag{2}$$

Algorithm 2 Random-MOPPO

```
1: Input: State  $s_t$ , weights  $\mathbf{w}$ 
2: Initialize:  $K$  weight-conditioned Actor-Critic Networks  $\pi_k / v^{\pi_k}$ , scalarisation-vector subspaces  $\mathbf{W}_k$  for  $k = 1, \dots, K$ , memory buffer  $\mathcal{E}$  size of  $D$ , and scalarisation-vector re-sampling frequency  $RF$ .
3: for  $k = 1$  to  $K$  do
4:   for  $t = 1$  to  $D$  do
5:     if  $t \% RF = 0$  then
6:        $\mathbf{w}_t \leftarrow$  uniform random sample( $\mathbf{W}_k$ )
7:     end if
8:      $a_t \leftarrow \pi_k(s_t, \mathbf{w}_t)$ 
9:      $s_{t+1}, \mathbf{r}_t \leftarrow \text{simulator}(a_t)$ 
10:     $\mathcal{E} \leftarrow \mathcal{E} \cup \langle s_t, a_t, \mathbf{w}_t, \mathbf{r}_t, s_{t+1} \rangle$ 
11:     $s_t \leftarrow s_{t+1}$ 
12:  end for
13:  sample  $\langle s_t, a_t, \mathbf{w}_t, \mathbf{r}_t, s_{t+1} \rangle \leftarrow \mathcal{E}$ 
14:   $\theta \leftarrow \theta + \eta (\nabla_{\theta} \log \pi_{\theta}(s, a; \mathbf{w}_t)) (A^{\pi}(s_t, a_t; \mathbf{w}_t))$ 
15:   $\phi \leftarrow \phi + \|\mathbf{V}^{\pi_k}(s_t; \mathbf{w}_t) - \mathbf{V}^{\pi_k}(s_{t+1}; \mathbf{w}_t)\|^2$ 
16:  clear  $\mathcal{E}$ 
17: end for
```

2.2 HalfCheetah-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^{17}, \mathcal{A} \in \mathbb{R}^6$.

The first objective is forward speed in x axis:

$$R_1 = \min(v_x, 4) + C \quad (3)$$

The second objective is energy efficiency:

$$R_2 = 4 - \sum_i a_i^2 + C, \quad a_i \in (-1, 1) \quad (4)$$

$$C = 1 \quad (5)$$

2.3 Walker2d-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^{17}, \mathcal{A} \in \mathbb{R}^6$.

The first objective is forward speed in x axis:

$$R_1 = v_x + C \quad (6)$$

The second objective is energy efficiency:

$$R_2 = 4 - \sum_i a_i^2 + C, \quad a_i \in (-1, 1) \quad (7)$$

$$C = 1 \quad (8)$$

2.4 Ant-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^{27}, \mathcal{A} \in \mathbb{R}^8$.

The first objective is forward speed in x axis:

$$R_1 = v_x + C \quad (9)$$

The second objective is forward in y axis:

$$R_2 = v_y + C \quad (10)$$

$$C = 1 - 0.5 \sum_i a_i^2, \quad a_i \in (-1, 1) \quad (11)$$

Algorithm 3 UCB-MOPPO

```
1: Input: Environment state  $S_t$  and full weight set  $\mathbf{W}$ .
2: Initialize:  $K$  weight-conditioned Actor-Critic networks  $\{(\pi_k, v^{\pi_k})\}_{k=1}^K$ ; for each  $k$ , a predetermined
   subspace  $\mathbf{W}_k \subset \mathbf{W}$  (of size  $M$ ); working pool  $\widetilde{\mathbf{W}}_k \leftarrow \mathbf{W}_k$ ; number of objectives  $m$ ; warm-up iterations
    $Q$ ; evaluation interval  $C$ ; for each  $k$  and  $j = 1, \dots, m$ , surrogate dataset  $D_{\text{surrogate}}^{k,j} \leftarrow \emptyset$ ; surrogate
   models  $\{f_{\text{bagging}}^{k,j}\}$ ; and dynamic weight pool  $\mathcal{E}$  (of size  $D$ ).
3: for  $t = 0$  to  $T$  do ▷ Warm-up Stage
4:   for  $k = 1$  to  $K$  do
5:      $\widetilde{\mathbf{W}}_k \leftarrow \text{getPivotWeights}(\mathbf{W}_k)$ 
6:      $\pi_k \leftarrow \text{FixWeightOptimization}(\pi_k, \widetilde{\mathbf{W}}_k)$ 
7:   end for ▷ Periodic Evaluation (every  $C$  iterations)
8:   if  $t \bmod C = 0$  then
9:     for  $k = 1$  to  $K$  do
10:      for all  $\mathbf{w} \in \widetilde{\mathbf{W}}_k$  do
11:         $V_{j,\mathbf{w}}^{\pi_k} \leftarrow \text{Simulate}(\pi_k, \mathbf{w})$  ▷ for each objective  $j$ 
12:      end for
13:    end for
14:  end if
15:  if  $t > Q$  then ▷ Construct Surrogate Training Data
16:    for  $k = 1$  to  $K$  do
17:      for all  $\mathbf{w} \in \widetilde{\mathbf{W}}_k$  do
18:        for  $z = 0$  to  $\lfloor t/C \rfloor - 1$  do
19:          for  $j = 1$  to  $m$  do
20:             $\Delta V_{j,\mathbf{w}}^{\pi_k, (z \rightarrow z+1)} \leftarrow V_{j,\mathbf{w}}^{\pi_k, z+1} - V_{j,\mathbf{w}}^{\pi_k, z}$ 
21:            Append  $(\mathbf{w}, \Delta V_{j,\mathbf{w}}^{\pi_k, (z \rightarrow z+1)})$  to  $D_{\text{surrogate}}^{k,j}$ 
22:          end for
23:        end for
24:      end for
25:    end for ▷ Update Surrogate Models
26:    for  $k = 1$  to  $K$  do
27:      for  $j = 1$  to  $m$  do
28:        Update  $f_{\text{bagging}}^{k,j}$  using  $D_{\text{surrogate}}^{k,j}$ 
29:      end for
30:    end for ▷ Scalarisation-Vector Selection via UCB Acquisition
31:    for  $k = 1$  to  $K$  do
32:       $\mathcal{D}_k \leftarrow \emptyset$ 
33:      for all  $\mathbf{w} \in \mathbf{W}_k$  do ▷ Candidates from the full subspace
34:         $V_{\mathbf{w}}^{\pi_k} \leftarrow \text{Simulate}(\pi_k, \mathbf{w})$ 
35:        for  $j = 1$  to  $m$  do
36:           $\hat{V}_{j,\mathbf{w}}^{\pi_k} \leftarrow V_{j,\mathbf{w}}^{\pi_k} + f_{\text{bagging}}^{k,j}(\mathbf{w})$ 
37:           $\tilde{V}_{j,\mathbf{w}}^{\pi_k} \leftarrow \hat{V}_{j,\mathbf{w}}^{\pi_k} + \beta_{t'} \cdot \sigma_{k,j}(\mathbf{w})$ 
38:        end for
39:        Let  $\mathbf{V}_{\mathbf{w}} \leftarrow (\tilde{V}_{1,\mathbf{w}}^{\pi_k}, \dots, \tilde{V}_{m,\mathbf{w}}^{\pi_k})$ 
40:        Compute  $HV \leftarrow \text{HV}(\text{Pareto}(\mathcal{L} \cup \{\mathbf{V}_{\mathbf{w}}\}))$ , where  $\mathcal{L}$  is the set of existing objective vectors.
41:        Add the pair  $(\mathbf{w}, HV)$  to  $\mathcal{D}_k$ 
42:      end for
43:      Sort  $\mathcal{D}_k$  in descending order of  $HV$ 
44:      Update working pool:  $\widetilde{\mathbf{W}}_k \leftarrow$  top  $N$  weights from  $\mathcal{D}_k$ 
45:    end for
46:  end if
47: end for
```

2.5 Hopper-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^{11}, \mathcal{A} \in \mathbb{R}^3$.

The first objective is forward speed in x axis:

$$R_1 = 1.5v_x + C \quad (12)$$

The second objective is jumping height:

$$R_2 = 12(h - h_{init}) + C \quad (13)$$

$$C = 1 - 2e^{-4} \sum_i a_i^2, \quad a_i \in (-1, 1) \quad (14)$$

2.6 Hopper-v3

Observation and action space: $\mathcal{S} \in \mathbb{R}^{11}, \mathcal{A} \in \mathbb{R}^3$.

The first objective is forward speed in x axis:

$$R_1 = 1.5v_x + C \quad (15)$$

The second objective is jumping height:

$$R_2 = 12(h - h_{init}) + C \quad (16)$$

The third objective is energy efficiency:

$$R_3 = 4 - \sum_i a_i^2 + C \quad (17)$$

$$C = 1 \quad (18)$$

3 PPO Hyperparameters

Table 1: Hyper-parameter configuration of MOPPO algorithms.

Hyperparameters	Value
Policy Number	10
Max Training Iterations	2×10^6
Number of Cells	64
Actor Learning Rate	3×10^{-4}
Critic Learning Rate	3×10^{-4}
Memory Size	2500
K Epochs	10
Gamma	0.99
Lambda	0.95
C1 Coefficient	0.5
C2 Coefficient	0
Epsilon Clip	0.2
Minibatch Size	64

4 Convex Coverage Set Expansion

In this section, we illustrate the expansion of the Convex Coverage Set (CCS) throughout the training process using our proposed UCB-MOPPO algorithm. The graphs are arranged sequentially from left to right and top to bottom, showing the progressive evolution of the CCS. Each graph depicts 100 sub-space vectors representing a two-objective optimization problem. The detailed comparison of baselines is provided in Table 2, reporting the mean and standard deviation over three random seeds.

Table 2: Evaluation of HV and EU metrics for continuous MORL tasks over three independent runs. The best results are highlighted in **bold**.

Benchmark	Metric	UCB	Mean	Random	Fixed	PGMORL	PDMORL	CAPQL	GPI-LS
Swimmer-V2	HV (10^4)	5.60 $\pm .18$	4.72 $\pm .19$	4.56 $\pm .23$	4.45 $\pm .11$	1.67 $\pm .09$	1.77 $\pm .02$	2.40 $\pm .03$	4.75 $\pm .02$
	EU (10^2)	2.18 $\pm .48$	2.16 $\pm .25$	2.08 $\pm .22$	2.11 $\pm .12$	1.23 $\pm .42$	1.24 $\pm .49$	1.79 $\pm .45$	2.14 $\pm .41$
Halfcheetah-V2	HV (10^7)	1.78 $\pm .23$	1.60 $\pm .07$	1.21 $\pm .09$	1.12 $\pm .05$	0.58 $\pm .01$	0.62 $\pm .02$	2.20 $\pm .08$	2.16 $\pm .02$
	EU (10^3)	3.88 $\pm .22$	3.58 $\pm .25$	3.55 $\pm .30$	3.44 $\pm .36$	2.30 $\pm .44$	2.42 $\pm .31$	4.47 $\pm .03$	4.30 $\pm .08$
Walker2d-V2	HV (10^7)	1.40 $\pm .03$	1.29 $\pm .13$	1.15 $\pm .06$	1.13 $\pm .07$	0.44 $\pm .02$	0.56 $\pm .04$	0.18 $\pm .05$	1.22 $\pm .04$
	EU (10^3)	3.54 $\pm .42$	3.42 $\pm .29$	3.34 $\pm .29$	3.30 $\pm .33$	1.98 $\pm .21$	2.24 $\pm .34$	1.50 $\pm .12$	3.40 $\pm .30$
Ant-V2	HV (10^7)	1.07 $\pm .09$	0.92 $\pm .04$	0.65 $\pm .01$	0.60 $\pm .01$	0.61 $\pm .10$	0.66 $\pm .02$	0.45 $\pm .04$	0.81 $\pm .01$
	EU (10^3)	2.96 $\pm .43$	2.61 $\pm .17$	2.35 $\pm .24$	2.21 $\pm .39$	2.28 $\pm .35$	2.41 $\pm .20$	2.03 $\pm .42$	2.84 $\pm .26$
Hopper-V2	HV (10^7)	0.84 $\pm .05$	0.81 $\pm .04$	0.79 $\pm .01$	0.76 $\pm .02$	0.23 $\pm .02$	0.25 $\pm .02$	0.25 $\pm .07$	0.80 $\pm .02$
	EU (10^3)	2.77 $\pm .28$	2.84 $\pm .30$	2.77 $\pm .18$	2.71 $\pm .26$	1.48 $\pm .47$	1.51 $\pm .18$	1.46 $\pm .11$	2.75 $\pm .07$
Hopper-V3	HV (10^{10})	3.62 $\pm .01$	2.83 $\pm .22$	2.90 $\pm .11$	2.64 $\pm .18$	0.63 $\pm .07$	0.11 $\pm .02$	0.31 $\pm .06$	0.65 $\pm .03$
	EU (10^3)	3.20 $\pm .02$	3.19 $\pm .04$	3.28 $\pm .20$	2.90 $\pm .01$	1.70 $\pm .04$	1.74 $\pm .09$	1.52 $\pm .29$	2.14 $\pm .42$

As training progresses, a clear trend of CCS growth emerges, characterized by an increasing spread and coverage of vectors across the objective space. Notably, the distribution of vectors provides valuable insights: regions where vectors are more widely scattered indicate areas with fewer vectors dominated by others, reflecting an expansion toward a more optimal and comprehensive CCS. This separation demonstrates the growing diversity and coverage of the vectors over time, effectively showcasing the CCS’s ability to capture a wider range of trade-offs between objectives. Consequently, this illustrates the effectiveness of UCB-MOPPO in thoroughly exploring the objective space and improving the set of solutions throughout the training process.

4.1 Swimmer-V2

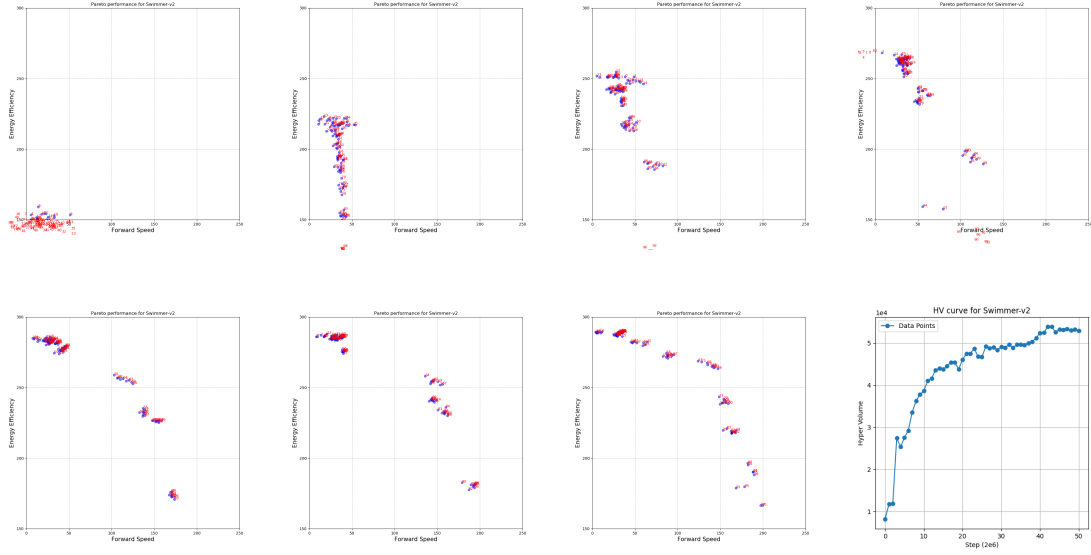


Figure 1: CCS expansion in Swimmer-V2 with one seed. The last graph shows the Hypervolume growth.

4.2 HalfCheetah-v2

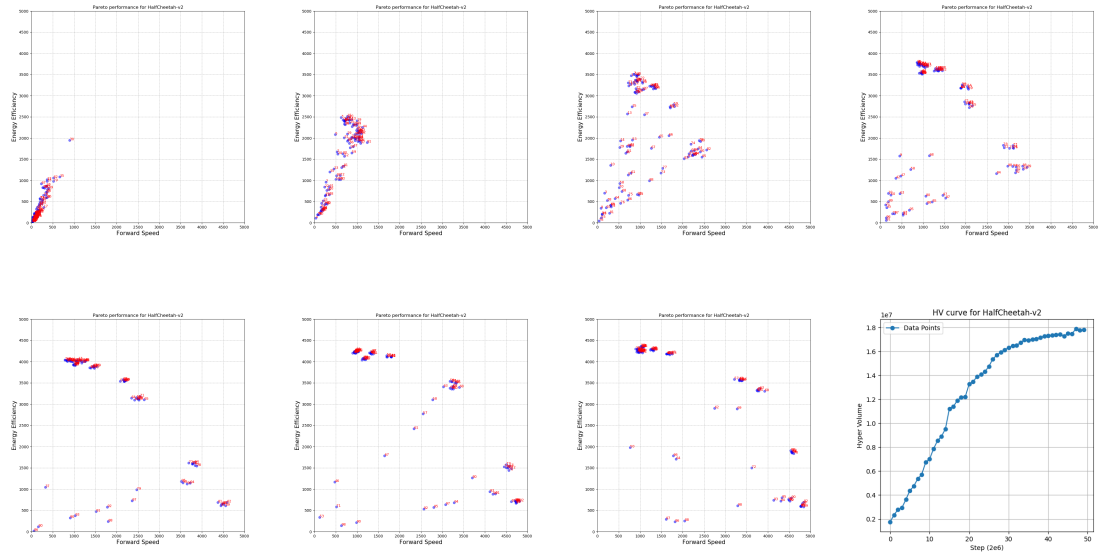


Figure 2: CCS expansion in HalfCheetah-V2 with one seed. The last graph shows the Hypervolume growth.

4.3 Walker2D-V2

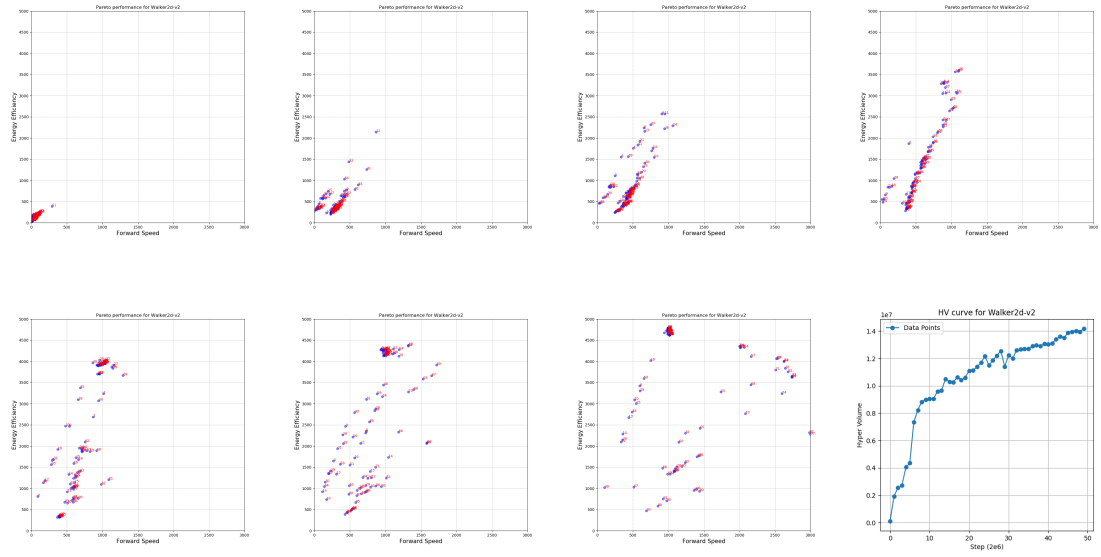


Figure 3: CCS expansion in Walker2D-V2 with one seed. The last graph shows the Hypervolume growth.

4.4 Ant-V2

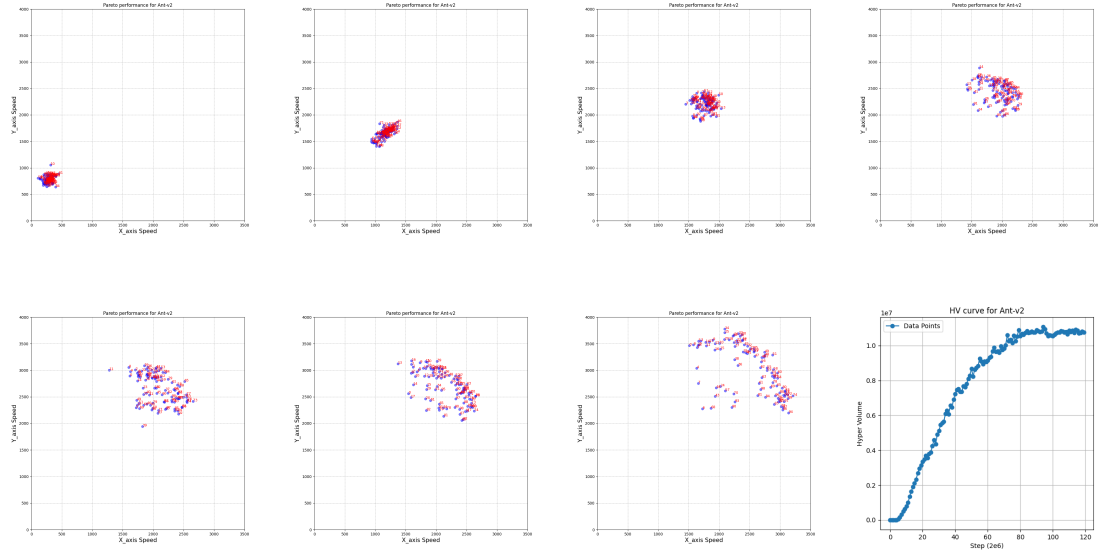


Figure 4: CCS expansion in Ant-V2 with one seed. The last graph shows the Hypervolume growth.

4.5 Hopper-V2

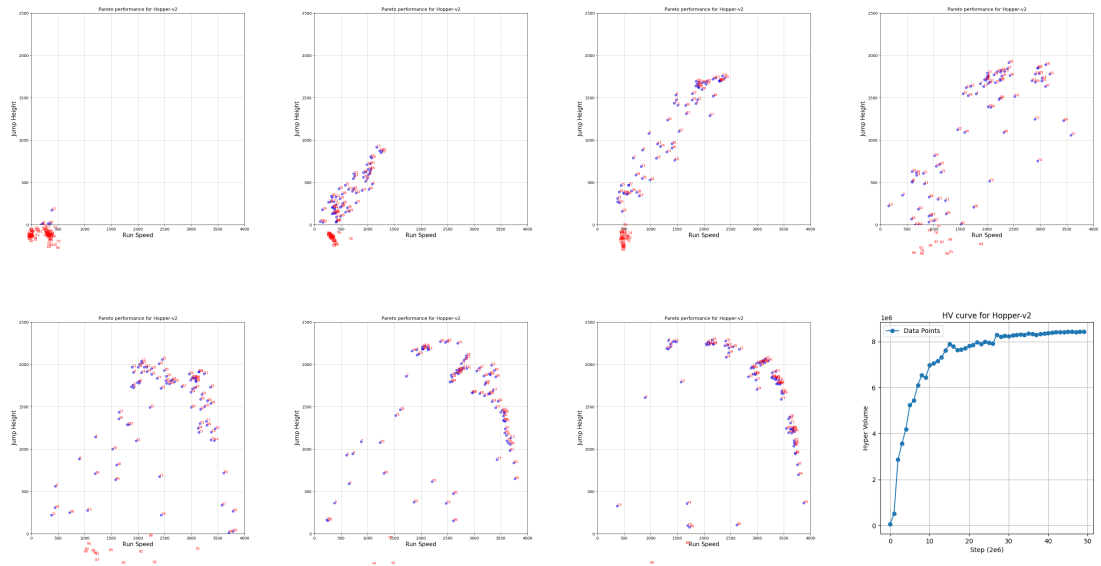


Figure 5: CCS expansion in Hopper-V2 with one seed. The last graph shows the Hypervolume growth.

4.6 Hopper-V3

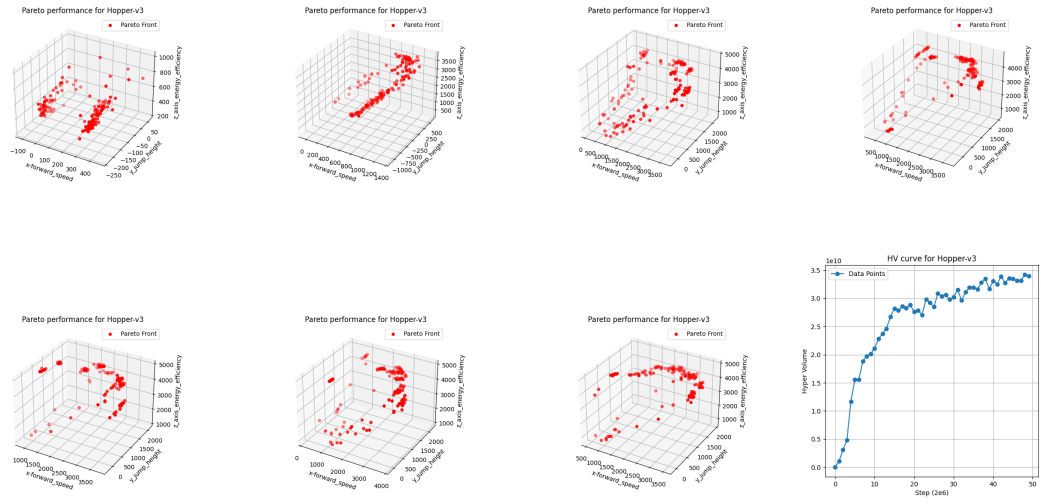


Figure 6: CCS expansion in Hopper-V3 with one seed. The last graph shows the Hypervolume growth.