# Residual-MPPI:
# Online Policy Customization for Continuous Control

Pengcheng Wang[*1], Chenran Li[*1], Catherine Weaver[1], Kenta Kawamoto[2], Masayoshi Tomizuka[1], Tang Chen[3], Wei Zhan[1]
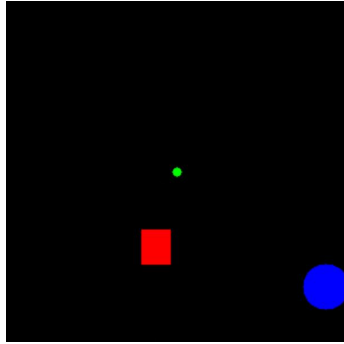
[1]University of California Berkeley; [2]Sony AI, USA; [3]University of Texas at Austin

*Equal Contribution

# Motivations

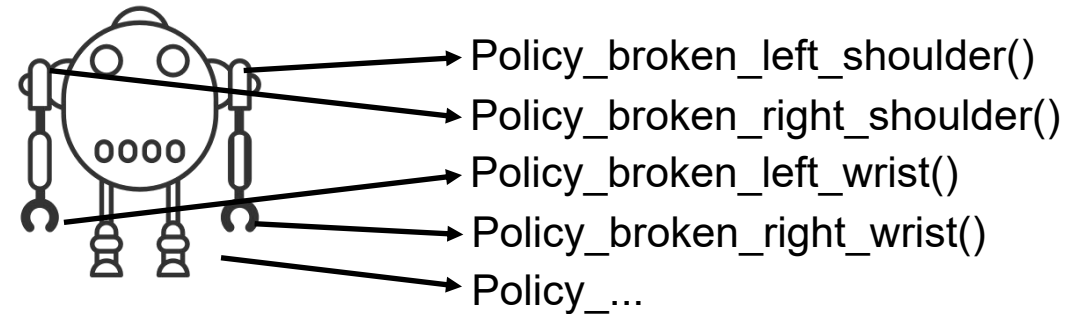- ## RL/ IL-based Advanced Polices


**Navigation**


**Manipulation**


**Locomotion**

- ## Efficient Policy Customization

  - **Access to original training metrics**

  - **"Policies for every disabled joints"**



Policy_broken_left_shoulder()

Policy_broken_right_shoulder()

Policy_broken_left_wrist()

Policy_broken_right_wrist()

Policy_...

**Online for Efficient Policy Customization!**

# Example

- ◘ **GT Sophy 1.0 Behavior**

# Preliminaries

- ## Residual Q Learning (RQL)

  - **Solve Policy Customization**

  - **~ Solve MaxEnt Augmented MDP**

$$\hat{\mathcal{M}} = (\mathcal{X}, \mathcal{U}, \omega r + r_R, p)$$

$$\mathcal{M}^{\mathrm{aug}} = (\mathcal{X}, \mathcal{U}, \omega' \log \pi(\boldsymbol{u}|\boldsymbol{x}) + r_R, p)$$

- ## Model Predictive Path Integral (MPPI)

  - **Solve MaxEnt MDP Online**

$$S_{\boldsymbol{x}_0}(U) = \sum_{t=0}^{T-1} r(\boldsymbol{x}_t, \boldsymbol{u}_t) + \phi(\boldsymbol{x}_T)$$

$$q^*(U) = \frac{1}{\eta} \exp\left(\frac{1}{\lambda} S_{\boldsymbol{x}_0}(U)\right) p(\mathcal{E})$$

# Methods

- ☐ **Residual-MPPI**

$$S_{\boldsymbol{x}_0}^{\mathrm{aug}}(U) = \sum_{t=0}^{T-1} \gamma^t \cdot \left( r_R(\boldsymbol{x}_t, \boldsymbol{u}_t) + \omega' \log \pi(\boldsymbol{u}_t | \boldsymbol{x}_t) \right)$$

- ☐ **Planning Loop**

☑ **MuJoCo Environment**

| Env. | Policy | Full Task | Basic Task | Add-on Task | |
|---|---|---|---|---|---|
| | | Total Reward | Basic Reward | $|\bar{\theta}|$ | Add-on Reward |
| Half Cheetah | Prior Policy | $1000.7 \pm 88.8$ | $2449.8 \pm 52.3$ | $0.14 \pm 0.00$ | $-1449.1 \pm 45.3$ |
| | Greedy-MPPI | $\mathbf{1939.9 \pm 134.7}$ | $2180.9 \pm 87.3$ | $\mathbf{0.02 \pm 0.01}$ | $-241.0 \pm 50.3$ |
| | Full-MPPI | $-3595.1 \pm 322.7$ | $-1167.3 \pm 144.0$ | $0.24 \pm 0.03$ | $-2427.7 \pm 320.3$ |
| | Guided-MPPI | $1849.6 \pm 151.0$ | $2154.6 \pm 95.7$ | $0.03 \pm 0.01$ | $-305.0 \pm 58.7$ |
| | Valued-MPPI | $1760.7 \pm 478.8$ | $\mathbf{2201.8 \pm 258.3}$ | $0.04 \pm 0.02$ | $-441.0 \pm 222.5$ |
| | Residual-MPPI | $\mathbf{1936.2 \pm 109.3}$ | $2178.6 \pm 71.9$ | $\mathbf{0.02 \pm 0.00}$ | $-242.3 \pm 40.5$ |
| | Residual-SAC (200K) | $-265.0 \pm 919.0$ | $455.4 \pm 678.6$ | $0.07 \pm 0.03$ | $720.4 \pm 251.8$ |
| | Residual-SAC (4M) | $2184.5 \pm 29.7$ | $2233.7 \pm 29.3$ | $0.00 \pm 0.00$ | $-49.2 \pm 1.7$ |
| | Fulltask-SAC | $2149.9 \pm 28.6$ | $2214.5 \pm 27.2$ | $0.01 \pm 0.00$ | $-64.5 \pm 2.4$ |

| Env | Policy | Total Reward | Basic Reward | $|\bar{\theta}|$ | Add-on Reward |
|---|---|---|---|---|---|
| Swimmer | Prior Policy | $-245.2 \pm 5.6$ | $345.8 \pm 3.2$ | $0.59 \pm 0.01$ | $-591.0 \pm 5.8$ |
| | Greedy-MPPI | $-58.9 \pm 5.4$ | $275.8 \pm 3.1$ | $\mathbf{0.33 \pm 0.01}$ | $\mathbf{-334.7 \pm 7.4}$ |
| | Full-MPPI | $-1686.6 \pm 106.7$ | $14.1 \pm 6.3$ | $1.70 \pm 0.11$ | $-1700.7 \pm 106.2$ |
| | Guided-MPPI | $-149.0 \pm 5.6$ | $292.9 \pm 3.8$ | $0.44 \pm 0.01$ | $-441.9 \pm 7.2$ |
| | Valued-MPPI | $-205.8 \pm 6.3$ | $\mathbf{335.1 \pm 1.6}$ | $0.54 \pm 0.01$ | $-540.9 \pm 6.3$ |
| | Residual-MPPI | $-60.0 \pm 5.2$ | $275.8 \pm 3.4$ | $\mathbf{0.34 \pm 0.01}$ | $-335.9 \pm 7.6$ |
| | Residual-SAC (200K) | $-209.0 \pm 67.6$ | $2.1 \pm 15.5$ | $0.21 \pm 0.07$ | $-221.1 \pm 72.7$ |
| | Residual-SAC (4M) | $-10.5 \pm 24.1$ | $-1.5 \pm 16.9$ | $0.01 \pm 0.02$ | $-9.0 \pm 16.6$ |
| | Fulltask-SAC | $-4.2 \pm 17.1$ | $2.1 \pm 17.6$ | $0.01 \pm 0.00$ | $-6.3 \pm 3.0$ |

| Env. | Policy | Total Reward | Basic Reward | $\bar{z}$ | Add-on Reward |
|---|---|---|---|---|---|
| Hopper | Prior Policy | $7252.7 \pm 49.2$ | $3574.5 \pm 9.7$ | $1.37 \pm 0.00$ | $3678.2 \pm 48.3$ |
| | Greedy-MPPI | $\mathbf{7367.0 \pm 199.4}$ | $3553.0 \pm 58.4$ | $\mathbf{1.38 \pm 0.01}$ | $\mathbf{3814.0 \pm 156.8}$ |
| | Full-MPPI | $20.5 \pm 3.0$ | $3.6 \pm 0.7$ | $1.24 \pm 0.00$ | $16.9 \pm 2.4$ |
| | Guided-MPPI | $6121.3 \pm 1590.1$ | $3067.8 \pm 679.0$ | $1.35 \pm 0.03$ | $3053.4 \pm 917.7$ |
| | Valued-MPPI | $7243.9 \pm 75.7$ | $\mathbf{3562.7 \pm 14.5}$ | $1.37 \pm 0.01$ | $3681.2 \pm 74.6$ |
| | Residual-MPPI | $\mathbf{7363.0 \pm 254.9}$ | $3547.6 \pm 78.0$ | $\mathbf{1.38 \pm 0.01}$ | $\mathbf{3815.4 \pm 186.4}$ |
| | Residual-SAC (200K) | $3543.1 \pm 478.9$ | $1019.8 \pm 94.3$ | $1.27 \pm 0.01$ | $2523.2 \pm 405.5$ |
| | Residual-SAC (4M) | $7682.5 \pm 178.2$ | $2310.4 \pm 106.8$ | $1.54 \pm 0.01$ | $5372.0 \pm 75.8$ |
| | Fulltask-SAC | $7825.3 \pm 36.9$ | $2934.5 \pm 27.6$ | $1.49 \pm 0.00$ | $4890.8 \pm 39.6$ |

| Env | Policy | Total Reward | Basic Reward | $\bar{v}_y$ | Add-on Reward |
|---|---|---|---|---|---|
| Ant | Prior Policy | $6333.7 \pm 753.9$ | $6177.1 \pm 703.7$ | $0.16 \pm 0.22$ | $156.6 \pm 200.5$ |
| | Greedy-MPPI | $6104.2 \pm 1532.0$ | $5092.8 \pm 1305.2$ | $\mathbf{1.01 \pm 0.27}$ | $\mathbf{1011.3 \pm 277.7}$ |
| | Full-MPPI | $-2767.7 \pm 154.0$ | $-2764.4 \pm 114.2$ | $-0.00 \pm 0.11$ | $-3.3 \pm 108.0$ |
| | Guided-MPPI | $5160.9 \pm 1963.0$ | $4999.8 \pm 1887.9$ | $0.16 \pm 0.22$ | $161.2 \pm 217.7$ |
| | Valued-MPPI | $6437.0 \pm 1021.9$ | $\mathbf{6230.7 \pm 959.0}$ | $0.21 \pm 0.20$ | $206.3 \pm 196.3$ |
| | Residual-MPPI | $\mathbf{6846.7 \pm 647.8}$ | $5984.8 \pm 541.5$ | $\mathbf{0.86 \pm 0.19}$ | $\mathbf{861.8 \pm 189.8}$ |
| | Residual-SAC (200K) | $-1175.5 \pm 157.3$ | $-1178.3 \pm 156.4$ | $0.00 \pm 0.00$ | $2.7 \pm 3.9$ |
| | Residual-SAC (4M) | $6962.9 \pm 342.9$ | $5710.2 \pm 252.0$ | $1.25 \pm 0.13$ | $1252.7 \pm 127.3$ |
| | Fulltask-SAC | $7408.6 \pm 312.0$ | $3100.3 \pm 184.4$ | $4.31 \pm 0.21$ | $4308.3 \pm 209.2$ |

**Residual-MPPI is Effective and Data-efficient**

mechanical systems control laboratory

Berkeley
UNIVERSITY OF CALIFORNIA

# Experiments

- **GTS Customization**

### Table 2: Experimental Results of Residual-MPPI in GTS

| Policy | GT Sophy 1.0 | Zero-shot MPPI | Few-shot MPPI | Residual-SAC (80K laps) |
|---|---|---|---|---|
| Lap Time | $117.77 \pm 0.08$ | $123.34 \pm 0.22$ | $122.93 \pm 0.14$ | $130.00 \pm 0.13$ |
| Off-course Steps | $93.13 \pm 1.98$ | $9.03 \pm 3.33$ | $4.43 \pm 2.39$ | $0.87 \pm 0.78$ |
| Policy | Full-MPPI | Guided-MPPI | Greedy-MPPI | Residual-SAC (2K laps) |
| Lap Time | *Failed | *Failed | *Failed | *Failed |
| Off-course Steps | *Failed | *Failed | *Failed | *Failed |

The evaluation results are in the form of mean $\pm$ std over 30 laps. *Failed baseline is not able to finish a complete lap. Valued-MPPI is not available since we only have access to the policy network of GT Sophy 1.0.

**Residual-MPPI works in Complex Environment and Policy**

mechanical systems
control laboratory

Berkeley
UNIVERSITY OF CALIFORNIA

# Experiments

■ **GTS Demo**



- **Safer Driving Style**

- **Advanced Route Selection**

■ **Takeaway**

**Online Principled Customization**

**= Residual-MPPI**

**+ Dynamics**

**+ Add-on Reward**

Paper & Code