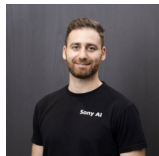


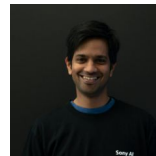
# Discovering Creative Behaviors through **DUPLEX**: Diverse **U**niversal Features for **P**olicy **E**xploration



Borja G. León<sup>†,1</sup>



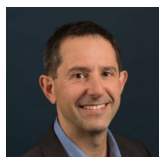
Francesco Riccio<sup>2</sup>



Kaushik Subramanian<sup>2</sup>



Peter R. Wurman<sup>2</sup>



Peter Stone<sup>2,3</sup>

Thirty-Eighth Annual Conference on Neural Information Processing Systems (NeurIPS), 2024

<sup>†</sup>internship project while at Sony AI.

1  ICNIC

2  Sony AI

3  TEXAS  
The University of Texas at Austin

---

# Superhuman but FUN?



---

# Superhuman and FUN

**DUPLEX** contributes to diversity learning in RL by improving on previous work to better **preserve the diversity vs. near-optimality trade-off** in highly-dynamic environments and multi-context settings.

**Showing diversity** and acting differently in the world **is fundamental** to create **engaging experiences for users**

## Context-conditioned diversity learning

*Diversity( $\Pi$ ) is a metric of dissimilarity among policies in a set  $\Pi$  with a common goal. Formally, if  $\psi_{\pi_i}$  and  $\psi_{\pi_j}$  are a function of state-occupancy of relevant features of any two policies in  $\Pi$ , then their dissimilarity is given by  $\|\psi_{\pi_i} - \psi_{\pi_j}\|$ . A non-zero value of this norm indicates dissimilarity, with larger values indicating greater divergence between the policies. Mathematically, diversity is defined as the sum of the minimum L2 dissimilarity norms in*

$$\text{Diversity}(\Pi) = \frac{1}{2 \text{ size}(\Pi)} \cdot \sum_{\substack{\forall \pi_i, \pi_j \in \Pi, \\ i \neq j}} \min \|\psi_{\pi_i} - \psi_{\pi_j}\|_2^2$$

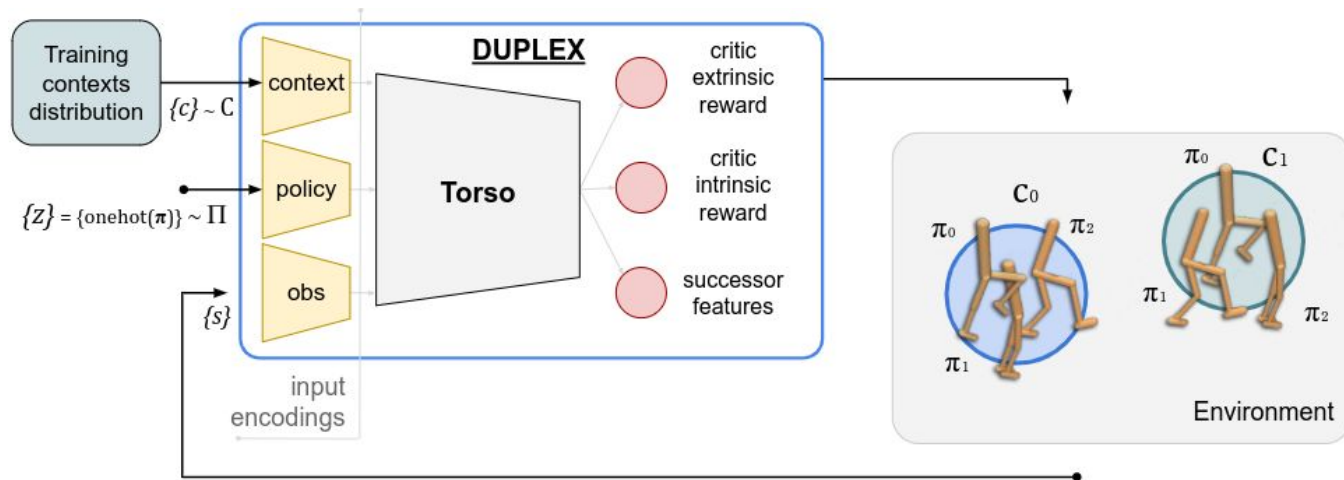
## Context-conditioned diversity learning

Accordingly, we measure  $\psi$  distances to enforce **context-conditioned diversity within  $\Pi$** . We aim at training an RL agent that, given a context  $c$ , discovers a set of  $n$  near-optimal policies by optimizing our objective function

$$\max_{\Pi} \text{Diversity}(\Pi) \text{ s.t. } d_{\pi_c} \cdot r_e \geq \rho \hat{v}_e, \quad \forall \pi_c \in \Pi$$

that forces the policies in  $\Pi$  to only **explore for diversity within the near-optimal region of the target value**

# Data flow



DUPLEX receives three inputs: (i) a context vector describing **task requirements** and environment dynamics; (ii) an **encoding of the policy**; (iii) and the current **state of the environment**. The critic network returns estimates for the **intrinsic and extrinsic rewards and successor features** to drive diverse behavior discovery. Finally, the algorithm samples policies in  $\Pi$  uniformly and rolls them out to collect more experience.

DUPLEX **stabilizes training and achieves diverse competitive policies** by introducing novel components to modulate the contribution of the intrinsic reward

### Dynamic intrinsic reward factor

$$\chi_t = \alpha_\chi \chi'_t + (1 - \alpha_\chi) \chi_{(t-1)}$$

$\chi$  scales intrinsic rewards proportionally to the sum of extrinsic values of policies in  $\Pi$

### Adding entropy to the SFs estimation

$$\psi^{\gamma, i}(s_t, a_t, c) = \phi_t + \mathbb{E}_{\pi_c} \sum_{k=t+1}^{\infty} \gamma^{k-t} [\phi_k + \alpha_H H(\pi_c^i(\cdot | s, c))]$$

to support critic estimates when policy is uncertain

$$r_I = \lambda \cdot \chi \cdot r_d$$

### Soft-lower bound

$$\lambda = \left\{ \sigma_k \left( \frac{v_{e_{\text{avg}}}^i - \beta \hat{v}_{e_{\text{avg}}}}{|\hat{v}_{e_{\text{avg}}} + l|} \right) \right\}_{i=1}^n$$

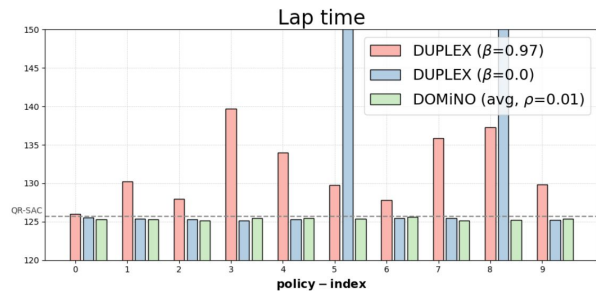
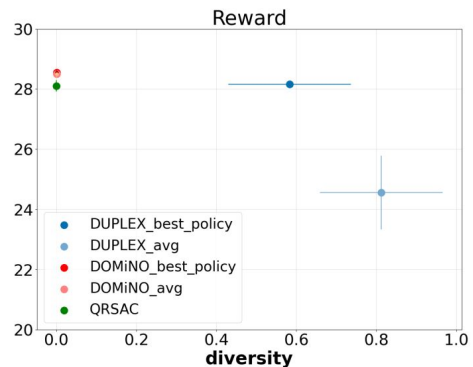
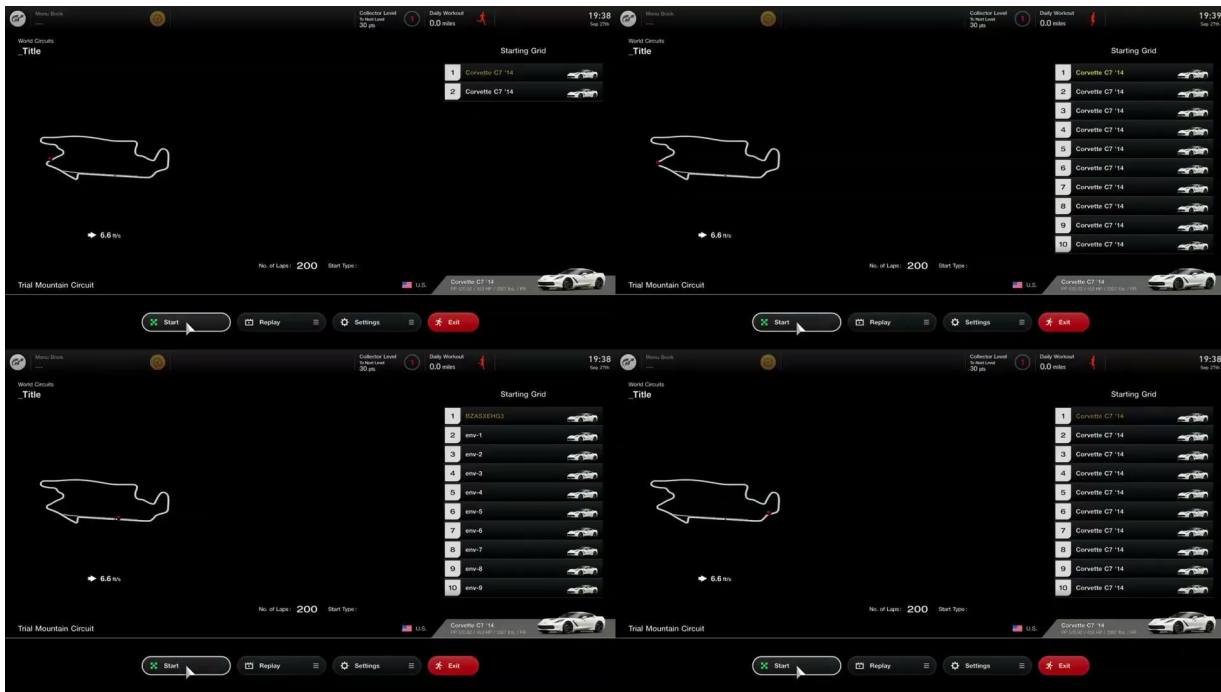
$\lambda$  to bound the near-optimal subspace for each policy

### Averaging critic networks

$$y(\phi, s', c, z) = \phi(t) + \gamma \left( \text{avg}_{j=1,2} \tilde{\psi}_{\theta_{\text{target}, j}}(s', \tilde{a}'_z, c) - \alpha \log \pi_{\omega}^z(\tilde{a}'_z | s', c) \right)$$

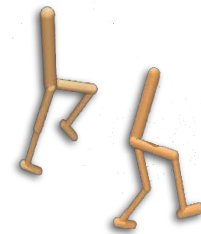
# Results: GranTurismo™ 7

DUPLEX trains **diverse competitive policies** in hyper-realistic driving simulators

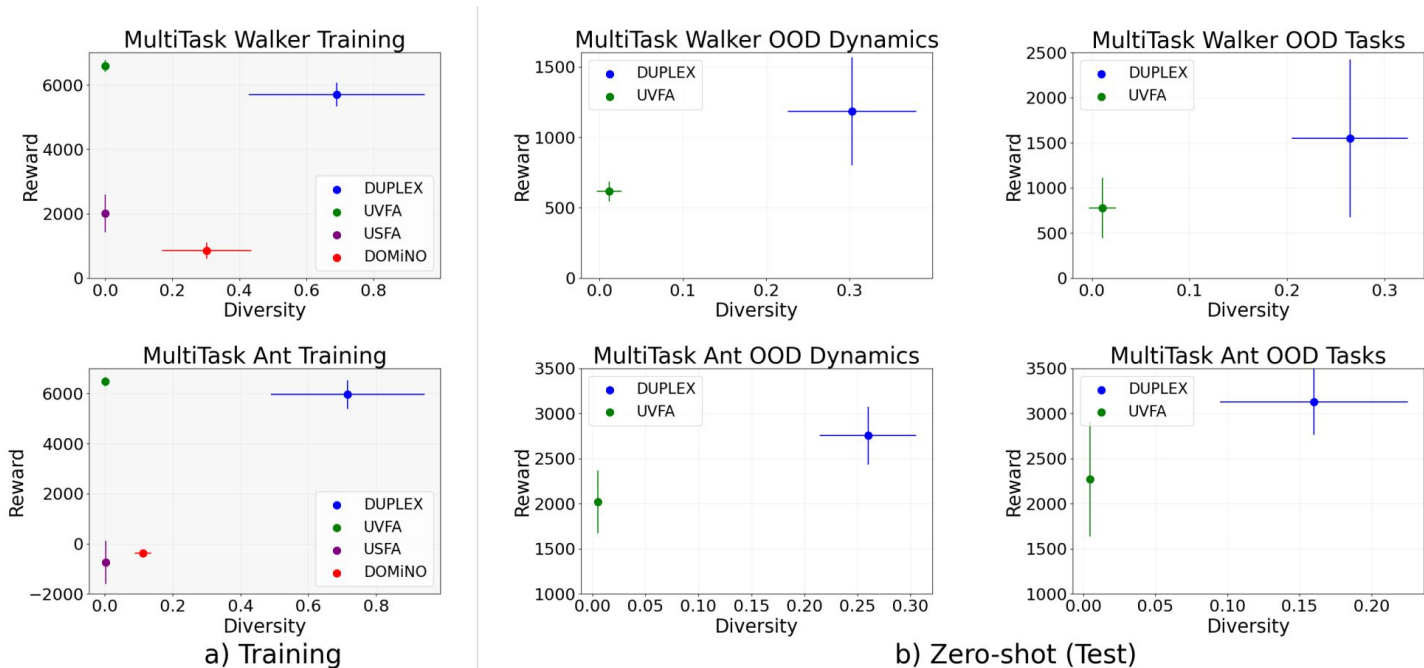




# Results: MuJoCo Walker2d and Ant



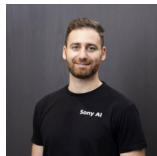
DUPLEX improves the diversity vs near-optimality trade-off both **within- and out-of- distribution** settings



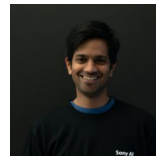
# Discovering Creative Behaviors through **DUPLEX**: Diverse **U**niversal Features for **P**olicy **E**xploration



Borja G. León<sup>†,1</sup>



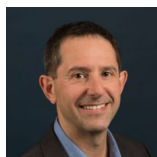
Francesco Riccio<sup>2</sup>



Kaushik Subramanian<sup>2</sup>



Peter R. Wurman<sup>2</sup>



Peter Stone<sup>2,3</sup>

Thirty-Eighth Annual Conference on Neural Information Processing Systems (NeurIPS), 2024

<sup>†</sup>internship project while at Sony AI.

1  **ICONIC**

2 **Sony AI**

3  **TEXAS**  
The University of Texas at Austin



more videos  
and results