# One Subgoal at a Time: Zero-Shot Generalization to Arbitrary Linear Temporal Logic Requirements in Multi-Task Reinforcement Learning
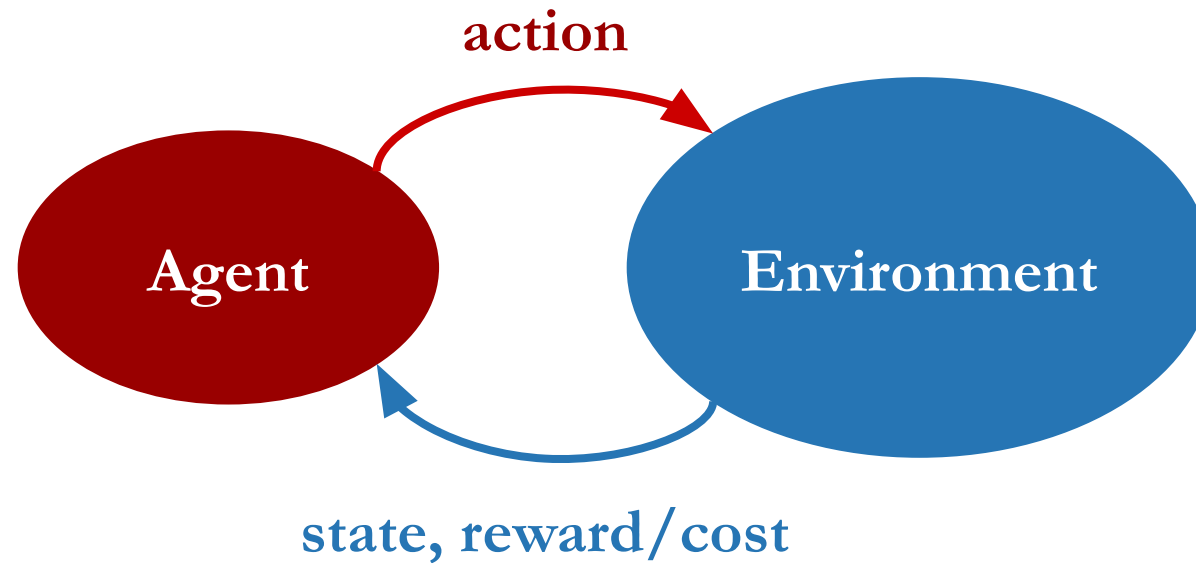
Zijian Guo[1], İlker Işık[2], H. M. Sabbir Ahmad[1], Wenchao Li[1,2]

[1]Division of Systems Engineering, Boston University
[2]Department of Electrical and Computer Engineering, Boston University
{zjguo, iilker, sabbir92, wenchao}@bu.edu
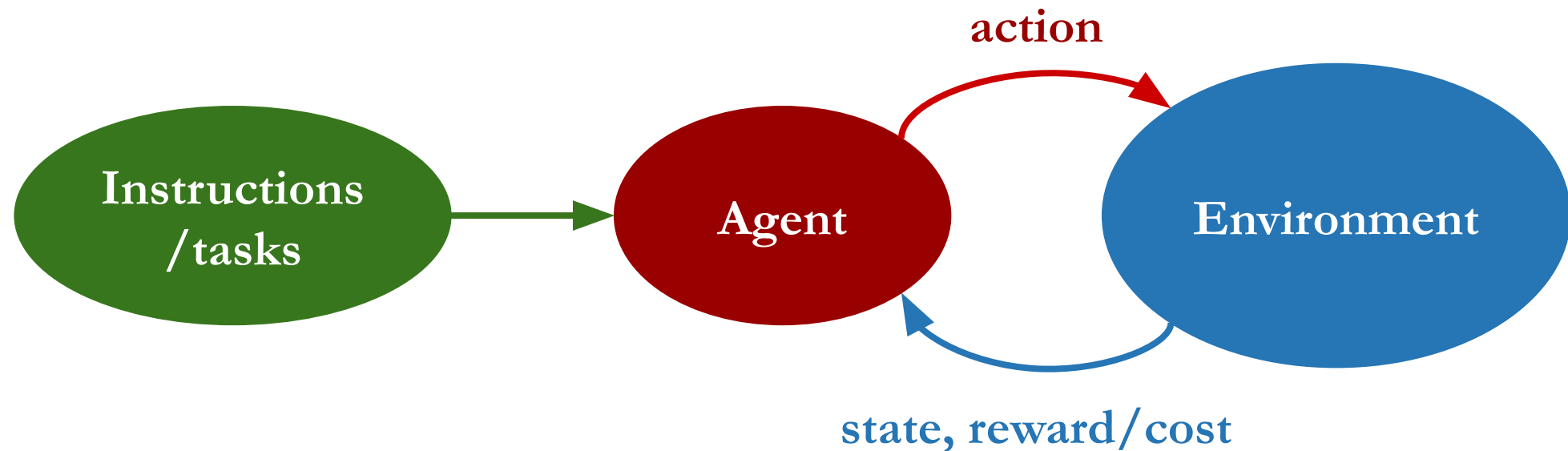
**BOSTON UNIVERSITY**

# Reinforcement Learning

# Reinforcement Learning

**How to follow diverse, complex, and even unseen instructions/tasks?**

- long-horizon goals, logical dependencies, safety constraints

# Background: Linear Temporal Logic

- **Formal language** to specify system's behaviors

# Background: Linear Temporal Logic

- **Formal language** to specify system's behaviors
- **Syntax** of LTL

$$\varphi := a \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \varphi_1 \vee \varphi_2 \mid \mathsf{F}\ \varphi \mid \mathsf{G}\ \varphi \mid \varphi_1 \mathsf{U}\ \varphi_2$$

- Atomic propositions: $AP, a \in AP$
- Boolean ($\neg$, $\wedge$, $\vee$) and temporal ($\mathsf{F}$, $\mathsf{G}$, $\mathsf{U}$) operators.

# Background: Linear Temporal Logic

- **Formal language** to specify system's behaviors
- **Syntax** of LTL

$$\varphi := \boldsymbol{a} \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \varphi_1 \vee \varphi_2 \mid \mathsf{F}\ \varphi \mid \mathsf{G}\ \varphi \mid \varphi_1 \mathsf{U}\ \varphi_2$$

  - Atomic propositions: $AP, \boldsymbol{a} \in AP$
  - Boolean ($\neg$, $\wedge$, $\vee$) and temporal ($\mathsf{F}$, $\mathsf{G}$, $\mathsf{U}$) operators.

- **Tasks can be expressed over high-level environment features**
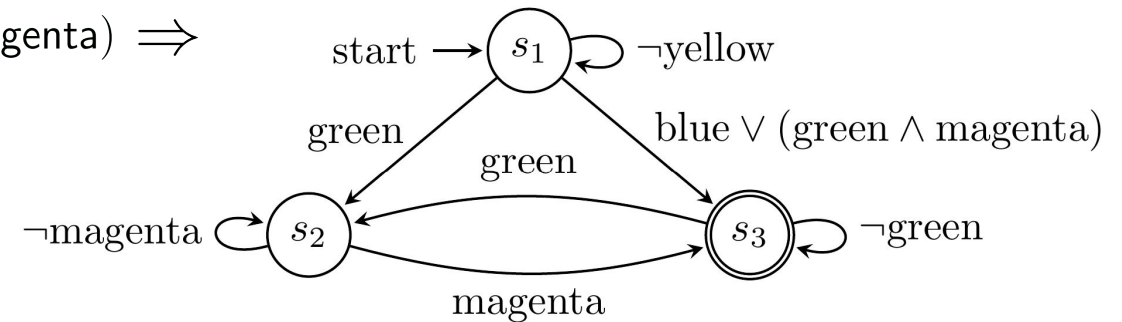
$$\varphi := (\neg\text{yellow}\ \mathsf{U}\ (\text{green} \vee \text{blue})) \wedge \mathsf{G}(\text{green} \Rightarrow \mathsf{F}\ \text{magenta})$$

(avoid yellow until reaching green or blue; whenever green is visited, magenta must eventually follow.)

# Background: Linear Temporal Logic

- **Büchi automata (BA):** for any LTL formula, it can be converted to an equivalent BA, which can be represented as **directed state-transition graphs**.

$\varphi := (\neg\text{yellow} \; \mathsf{U} \; (\text{green} \lor \text{blue})) \land \mathsf{G}(\text{green} \Rightarrow \mathsf{F} \; \text{magenta}) \implies$

# Background: Linear Temporal Logic

- **Büchi automata (BA):** for any LTL formula, it can be converted to an equivalent BA, which can be represented as **directed state-transition graphs**.
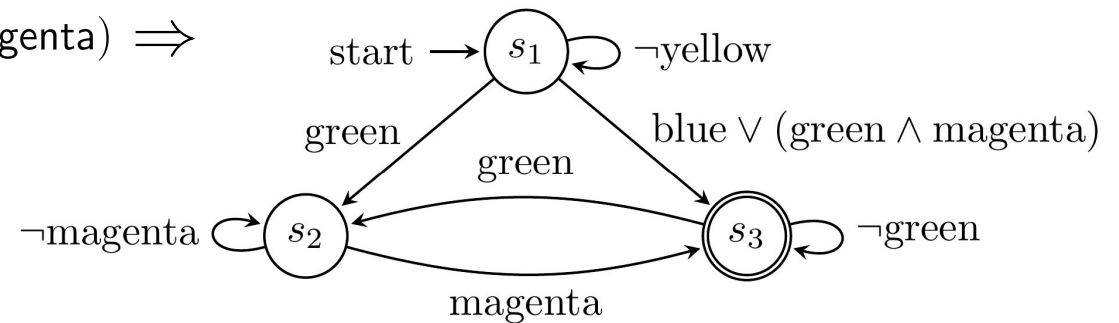
$\varphi := (\neg\text{yellow U (green} \lor \text{blue)}) \land \text{G(green} \Rightarrow \text{F magenta)} \Longrightarrow$



- **Reach-Avoid Subgoal Construction:** depth-first search (DFS) to enumerate all possible paths and extract reach-avoid subgoals $(\alpha^+, A^-)$

$$p_1 = \{(\alpha^+ = \{\text{green}\}, A^- = \{\text{yellow}\}), (\alpha^+ = \{\text{magenta}\}, A^- = \emptyset)\}$$
$$p_2 = \{(\alpha^+ = \{\text{blue}\}, A^- = \{\text{yellow}\})\}$$

# Challenges of Generalization

- Satisfying an LTL formula = completing a sequence of reach-avoid subgoals

# Challenges of Generalization

- Satisfying an LTL formula = completing a sequence of reach-avoid subgoals
- Existing methods:
    - Structure of the automaton/entire subgoal sequence
    - Policy conditioned on those representations
    - **Limitation: Out-of-distribution (OOD) issue** of new LTL formulas at test time
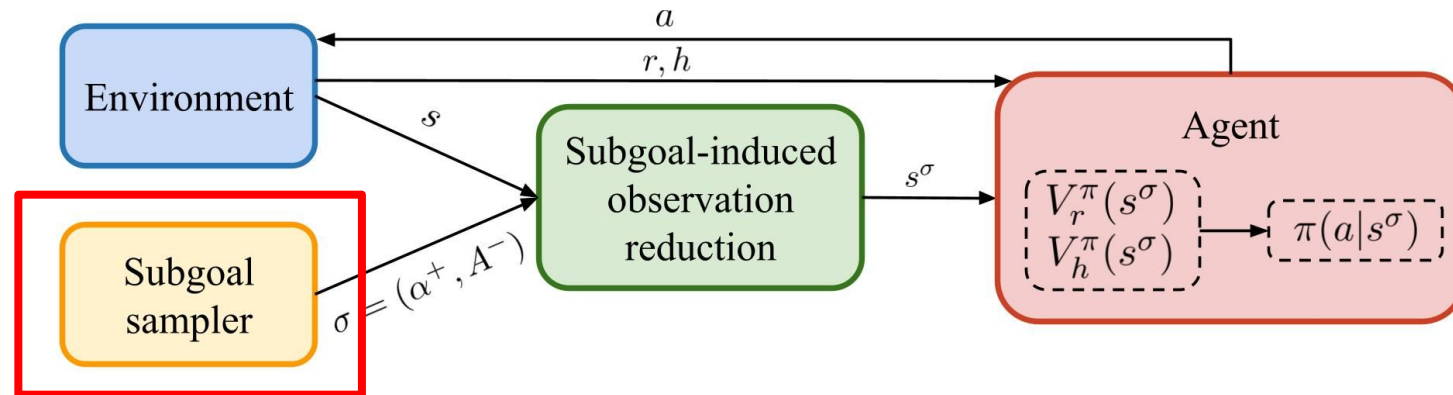
# Challenges of Generalization

- Satisfying an LTL formula = completing a sequence of reach-avoid subgoals
- Existing methods:
    - Structure of the automaton/entire subgoal sequence
    - Policy conditioned on those representations
    - **Limitation: Out-of-distribution (OOD) issue** of new LTL formulas at test time
- In contrast, we address this problem by solving **one subgoal at a time**
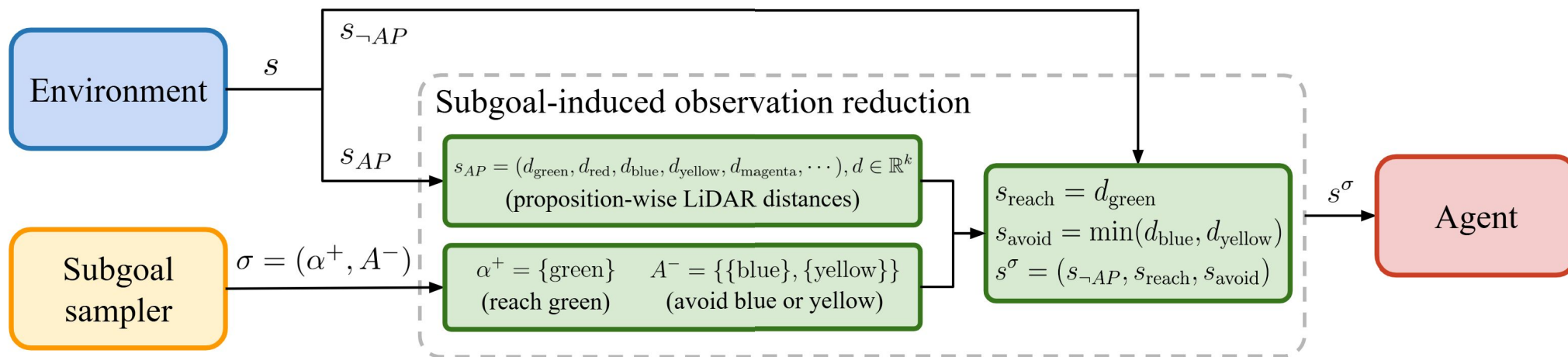
# GenZ-LTL: Training



- **Subgoal sampling:** all possible subgoals $\xi = \{(\alpha^+, A^-)_i\}_{i=1}^M$
  - Enumerate each assignment $\boldsymbol{a} \in 2^{AP}$ as a candidate $\alpha^+$
  - For each such $\alpha^+$, we filter out the remaining assignments that conflict with it. We then enumerate all possible combinations of the filtered assignments to form $A^-$

# GenZ-LTL: Training

- **Subgoal-Induced observation reduction**
  - Note that $a \in 2^{AP}$, so the total number of subgoal grows exponentially
  - Idea: focus only on subgoal-relevant observations to reduce sample complexity

# GenZ-LTL: Training

- **Policy learning with reachability constraints**

$$\pi_{k+1} = \arg\max_{\pi} \mathbb{E}_{\sigma \sim \mathrm{Unif}(\xi), s \sim d^{\pi_k}, a \sim \pi_k} \left[ \frac{\pi}{\pi_k} A_r^{\pi_k}(s^\sigma, a) \right]$$

$$\text{s.t.} \quad \mathbb{E}_{\sigma \sim \mathrm{Unif}(\xi), s \sim d^{\pi_k}} \left[ \mathcal{D}_{KL}(\pi, \pi_k) \right] \leq \epsilon$$

$$\mathbb{E}_{\sigma \sim \mathrm{Unif}(\xi), s \sim d^{\pi_k}, a \sim \pi_k} \left[ (1 - \gamma) J_h(\pi_k) + \frac{\pi}{\pi_k} A_h^{\pi_k}(s^\sigma, a) \right] \leq 0$$

where $h : \mathcal{S} \mapsto \mathbb{R}$ is the constraint violation function
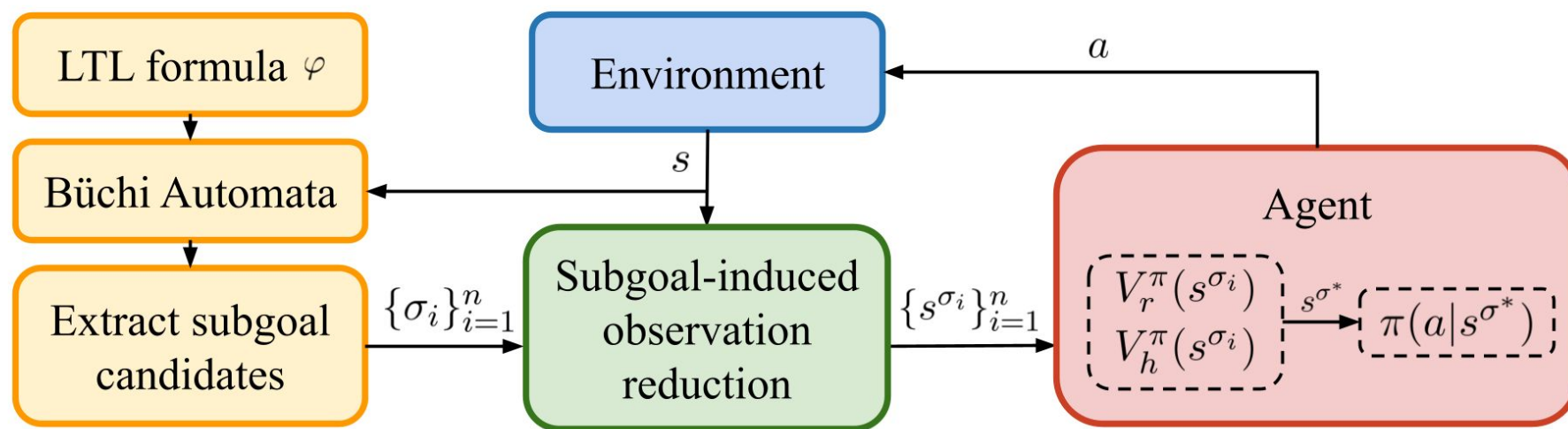
$$A_h^\pi(s^\sigma, a) := Q_h^\pi(s^\sigma, a) - V_h^\pi(s^\sigma)$$

$$Q_h^\pi(s^\sigma, a) := \max_{t \in \mathbb{N}} h(s_t^\sigma), s_0 = s^\sigma, a_0 = a, a_t \sim \pi$$

$$J_h(\pi) := \max_{t \in \mathbb{N}} h(\cdot), a \sim \pi$$

# GenZ-LTL: Testing

- Given a target LTL specification, we construct the corresponding BA and identify candidate subgoals based on the current automaton state. The subgoal to be executed is selected as $\sigma^* = \arg\max_\sigma V_r(s^\sigma) - \lambda(s^\sigma)V_h(s^\sigma)$

**BOSTON UNIVERSITY**

# Experimental Settings

- **Environments**
  - LetterWorld: 7×7 grid world
  - ZoneEnv: high-dimensional env with lidar observations
  - Randomized environments

# Experimental Settings

- **LTL specifications**



| | | LetterWorld | | ZoneEnv |
|---|---|---|---|---|
| **Finite-horizon** | $\varphi_1$ | F (a ∧ (¬b U c)) ∧ F d | $\varphi_9$ | (F b) ∧ (¬b U (g ∧ F y)) |
| | $\varphi_2$ | (F d) ∧ (¬f U (d ∧ F b)) | $\varphi_{10}$ | ¬(m ∨ y) U (b ∧ F g) |
| | $\varphi_3$ | ¬a U (b ∧ (¬c U (d ∧ (¬e U f)))) | $\varphi_{11}$ | ¬g U ((b ∨ m) ∧ (¬g U y)) |
| | $\varphi_4$ | (a ∨ b ∨ c ∨ d ⇒ F (e ∧ F (f ∧ F g))) U (h ∧ F i) | $\varphi_{12}$ | (g ∨ b ⇒ (¬y U m)) U y |
| | $\varphi_5$ | F (d ∧ (¬(a ∨ b) U (b ∧ (¬e U c)))) ∧ F (¬(f ∨ g ∨ h) U a) | $\varphi_{13}$ | F (g ∧ (¬(b ∨ y) U (y ∧ (¬m U b)))) ∧ F (¬g U y) |
| | $\varphi_6$ | F ((k ∧ ((¬b ∨ c) U f)) ∧ (¬(a ∨ e ∨ h) U g)) ∧ F d | $\varphi_{14}$ | F ((b ∨ g) ∧ (¬y U (b ∧ (¬(g ∨ m) U m)))) ∧ F (y ∧ (¬b U g)) |
| | $\varphi_7$ | ¬(j ∨ b ∨ d) U (a ∧ (¬c U (f ∧ F (g ∧ (¬d U e))))) | $\varphi_{15}$ | ¬(m ∨ y) U (b ∧ (¬g U (y ∧ F (g ∧ (¬b U m))))) |
| | $\varphi_8$ | ¬(f ∨ g) U (a ∧ (¬b U c) ∧ F (d ∧ (¬e U f))) | $\varphi_{16}$ | F (b ∧ (¬y U (g ∧ F (y ∧ (¬(m ∨ g) U b))))) |
| **Infinite-horizon** | $\psi_1$ | G F (e ∧ (¬a U f)) ∧ G ¬(c ∨ d) | $\psi_4$ | G F b ∧ G F g ∧ G ¬(y ∨ m) |
| | $\psi_2$ | G F a ∧ G F b ∧ G F c ∧ G ¬(e ∨ f ∨ i) | $\psi_5$ | G F b ∧ G F y ∧ G F g ∧ G ¬m |
| | $\psi_3$ | G F c ∧ G F a ∧ G F (e ∧ (¬f U g)) ∧ G F k ∧ G ¬(i ∨ j) | $\psi_6$ | F G y ∧ G ¬(g ∨ b ∨ m) |

- **All LTL specifications are unseen at test time for our method**

# Main Results

- **GenZ-LTL achieves higher success and lower violation rates, while learning more efficient policies**

| | | LTL2Action | | | GCRL-LTL | | | DeepLTL | | | RAD-embeddings | | | GenZ-LTL(ours) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $\eta_s \uparrow$ | $\eta_v \downarrow$ | $\mu \downarrow$ | $\eta_s \uparrow$ | $\eta_v \downarrow$ | $\mu \downarrow$ | $\eta_s \uparrow$ | $\eta_v \downarrow$ | $\mu \downarrow$ | $\eta_s \uparrow$ | $\eta_v \downarrow$ | $\mu \downarrow$ | $\eta_s \uparrow$ | $\eta_v \downarrow$ | $\mu \downarrow$ |
| Letter | $\varphi_{1-4}$ | $0.62_{\pm0.16}$ | $0.07_{\pm0.09}$ | $26.64_{\pm5.87}$ | $0.87_{\pm0.11}$ | $0.03_{\pm0.05}$ | $16.05_{\pm6.13}$ | $0.87_{\pm0.04}$ | $\mathbf{0.00_{\pm0.01}}$ | $7.51_{\pm1.21}$ | $0.90_{\pm0.06}$ | $0.04_{\pm0.04}$ | $17.79_{\pm3.37}$ | $\mathbf{0.98_{\pm0.02}}$ | $\mathbf{0.00_{\pm0.00}}$ | $\mathbf{7.22_{\pm1.18}}$ |
| | $\varphi_{5-8}$ | $0.24_{\pm0.12}$ | $0.20_{\pm0.25}$ | $36.43_{\pm6.08}$ | $0.65_{\pm0.08}$ | $0.11_{\pm0.06}$ | $18.71_{\pm2.54}$ | $0.76_{\pm0.05}$ | $0.01_{\pm0.02}$ | $10.62_{\pm1.47}$ | $0.82_{\pm0.10}$ | $0.07_{\pm0.08}$ | $24.29_{\pm3.77}$ | $\mathbf{0.95_{\pm0.03}}$ | $\mathbf{0.00_{\pm0.00}}$ | $\mathbf{9.82_{\pm1.42}}$ |
| Zone | $\varphi_{9-12}$ | $0.57_{\pm0.37}$ | $0.17_{\pm0.21}$ | $331.21_{\pm165.88}$ | $0.88_{\pm0.04}$ | $0.05_{\pm0.02}$ | $305.28_{\pm123.33}$ | $0.91_{\pm0.05}$ | $0.04_{\pm0.03}$ | $\mathbf{220.39_{\pm78.77}}$ | $0.94_{\pm0.05}$ | $0.04_{\pm0.05}$ | $269.53_{\pm129.95}$ | $\mathbf{0.99_{\pm0.01}}$ | $0.01_{\pm0.01}$ | $254.69_{\pm89.18}$ |
| | $\varphi_{13-16}$ | $0.12_{\pm0.18}$ | $0.17_{\pm0.28}$ | $886.39_{\pm288.60}$ | $0.70_{\pm0.07}$ | $0.09_{\pm0.04}$ | $606.42_{\pm26.76}$ | $0.87_{\pm0.08}$ | $0.03_{\pm0.06}$ | $505.17_{\pm54.03}$ | $0.70_{\pm0.15}$ | $\mathbf{0.00_{\pm0.01}}$ | $647.39_{\pm40.74}$ | $\mathbf{0.98_{\pm0.02}}$ | $0.01_{\pm0.01}$ | $\mathbf{408.71_{\pm27.47}}$ |

Finite-horizon tasks

| Methods | Metrics | LetterWorld | | | | ZoneEnv | |
|---|---|---|---|---|---|---|---|
| | | $\psi_1$ | $\psi_2$ | $\psi_3$ | $\psi_4$ | $\psi_5$ | $\psi_6$ |
| GenZ-LTL(ours) | $\mu_{acc} \uparrow$ | $\mathbf{208.95_{\pm14.39}}$ | $\mathbf{102.12_{\pm7.02}}$ | $\mathbf{55.17_{\pm1.08}}$ | $\mathbf{55.16_{\pm4.23}}$ | $\mathbf{32.75_{\pm1.20}}$ | $\mathbf{8135.67_{\pm1489.99}}$ |
| | $\eta_v \downarrow$ | $\mathbf{0.00_{\pm0.00}}$ | $\mathbf{0.00_{\pm0.00}}$ | $\mathbf{0.00_{\pm0.00}}$ | $\mathbf{0.07_{\pm0.02}}$ | $\mathbf{0.03_{\pm0.01}}$ | $\mathbf{0.03_{\pm0.02}}$ |
| DeepLTL | $\mu_{acc} \uparrow$ | $142.56_{\pm22.44}$ | $48.28_{\pm12.37}$ | $19.21_{\pm4.57}$ | $30.03_{\pm13.23}$ | $15.73_{\pm4.44}$ | $7337.38_{\pm2019.56}$ |
| | $\eta_v \downarrow$ | $0.04_{\pm0.02}$ | $0.09_{\pm0.01}$ | $0.09_{\pm0.04}$ | $0.39_{\pm0.10}$ | $0.38_{\pm0.24}$ | $0.13_{\pm0.05}$ |
| GCRL-LTL | $\mu_{acc} \uparrow$ | $41.98_{\pm15.80}$ | $22.77_{\pm9.50}$ | $9.53_{\pm2.28}$ | $30.00_{\pm3.72}$ | $14.61_{\pm1.62}$ | $5584.34_{\pm3180.15}$ |
| | $\eta_v \downarrow$ | $0.18_{\pm0.08}$ | $0.30_{\pm0.08}$ | $0.30_{\pm0.18}$ | $0.37_{\pm0.08}$ | $0.40_{\pm0.08}$ | $0.14_{\pm0.01}$ |

Infinite-horizon tasks

Please scan the following QR codes for more details

BOSTON
UNIVERSITY