

GRPO/Dr.GRPO Can Improve GUI-Agent Performance

Core Problem

GUI agents must learn to navigate interfaces, click buttons, fill forms, and complete multi-step tasks—all requiring sequential decision-making under uncertainty.

How GRPO/Dr.GRPO Helps

1. Reward-Driven Learning

Just as the current implementation rewards correct math equations (1.0 for correct, 0.1 for partial, 0.0 for failed), a GUI agent can be rewarded for:

- Successfully completing tasks (e.g., booking a flight, filling a form)
- Taking efficient action sequences (fewer clicks = higher reward)
- Recovering from errors gracefully

2. Group Exploration

GRPO's group sampling (multiple rollouts per prompt) allows the agent to:

- Explore different UI interaction strategies
- Learn which action sequences work better through group normalization
- Discover robust solutions across varied interface states

3. Sequential Decision Making

Unlike math problems, GUI tasks require multi-step reasoning:

- Each action changes the UI state
- The agent learns action → state transitions
- GRPO's advantage normalization helps distinguish good action sequences from bad ones

Key Adaptations Needed

Component	Current (Math Task)	GUI Agent Task
State representation	Text prompt with numbers	Screenshots/DOM → visual/textual embeddings
Action space	Generate equations	Generate UI actions (click, type, scroll)
Reward function	Equation correctness	Task completion metrics (success/failure, time efficiency)
Credit assignment	Single-step evaluation	Trajectory-based rewards across multiple actions

Example: Flight Booking Task

- **Prompt:** Screenshot + task description ("Book a flight from NYC to LA on Dec 25")
- **Rollouts:** 8 different action sequences
 - Rollout 1: Click "From" → Type "NYC" → Click "To" → Type "LA" → Click "Date" → Select Dec 25 → Click "Search"
 - Rollout 2: Different sequence or interaction pattern
 - ...
- **Rewards:**
 - 1.0 if flight booked correctly with all parameters matching
 - 0.5 if partial (e.g., wrong date selected)
 - 0.0 if failed to complete task
- **Learning:** GRPO normalizes within groups to prefer efficient, correct sequences

Transfer from Current Implementation

The same GRPO framework implemented transfers directly:

1. **Advantage Computation** (`compute_group_normalized_advantages`): Same normalization within groups of GUI action sequences
2. **PPO Loss** (`compute_loss`): Same clipped surrogate objective for stable updates
3. **Rollout Strategy:** Same group-based sampling to explore different interaction strategies
4. **Main Change:** Swap equation generation for action generation and equation validation for task completion verification

Benefits Over Standard RL

- **Sample Efficiency:** Group normalization provides clearer learning signal than individual trajectory rewards
- **Stability:** PPO clipping prevents catastrophic policy updates during GUI learning
- **Exploration:** Multiple rollouts per state encourage diverse strategy discovery
- **Robustness:** Learning from groups helps agent generalize across UI variations