

TLDR: TVCACHE safely reuses repeated tool calls and sandboxes in RL post-training to reduce tool-call latency by 3.4x–6.9x without hurting training reward.

Tool calls dominate agentic RL rollouts with repeated tool calls across rollouts. Existing caches lead to incorrect behavior as they fail to capture the sandbox state.

RL post-training repeatedly executes expensive tool calls

- RL post-training improves agents by generating **multiple rollouts for the same task**.
- Each rollout alternates between **reasoning tokens** and **tool calls**.
- Tool calls run in **isolated sandboxes** to support code execution, file edits, web search, and API calls.
- This makes rollout generation expensive: **tokens cost GPU time, tools add latency, and GPUs wait for tool responses**.

Takeaway: Tool-using RL agents repeatedly pay the cost of sandboxed tool execution.

Tool execution is a major rollout bottleneck

Benchmark: terminal-bench

Benchmark: SkyRL-SQL

- Across terminal-bench, SkyRL-SQL, and EgoSchema, tool execution consumes **7–43% of rollout time**.
- This leaves GPUs waiting for tool responses during RL post-training.

Takeaway: Reducing repeated tool execution can directly speed up rollout generation.

Same tool calls repeat cross rollouts – can we reuse them?

Here we are generating two rollouts for the same task

Opportunity: Rollouts for the same task often repeat the same tool calls.

Challenge: Tool calls can modify sandbox state. In Rollout 2, edit mutates the files before ./code.py executes.

Key Insight: **Safe reuse of the tool call outputs requires a stateful cache.** Stateless naive cache results in about 27% of incorrect cache hits

Our Solution: TVCache: A Stateful Tool-Value cache

Solution: Cache tool calls by the trajectory

Multiple rollouts for the same task

Extract tool call trajectories

Merge shared prefix into a TCG

Solution: Cache tool calls by trajectory

- TVCACHE builds a **Tool Call Graph (TCG)** for each task.
- Each node represents a **tool call** and stores the corresponding **tool response**.
- Shared prefixes across rollouts are stored once.
- TVCACHE reuses a cached response only when the **tool-call history matches**.

Takeaway: Cache key = current tool call + preceding trajectory.

Cache misses restore the closest reusable sandbox

Cache-miss handling

- Find longest matching prefix
- Restore snapshot at t2
- Execute unmatched suffix t3
- Execute tool call t4

- On a cache miss, TVCache finds the **longest matching prefix** in the TCG.
- If that node has a saved snapshot, TVCache **restores the sandbox** from that point.
- TVCACHE then executes only the **unmatched suffix** of the trajectory.
- Snapshots are stored **selectively** (when restoring is cheaper than re-executing)

Takeaway: Cache misses reuse the closest available state instead of restarting from scratch.

TVCACHE is transparent to rollout generation

Takeaway: TVCACHE speeds up tool use without changing the RL training loop.

Results

