

LATERAL TREE-OF-THOUGHTS SURPASSES TOT BY INCORPORATING LOGICALLY-CONSISTENT, LOW-UTILITY CANDIDATES

Anonymous authors

Paper under double-blind review

ABSTRACT

1 INTRODUCTION

2 MOTIVATION

The near-term problem at frontier scale. Frontier language models increasingly run in *compute-rich* inference settings: users and systems are willing to spend thousands of tokens (or many node expansions) per query to improve reliability. Yet the dominant search pattern—vanilla Tree-of-Thoughts (ToT)—*under-utilizes* this budget in two systematic ways already visible today and poised to worsen as budgets grow:

1. **Utility saturation (breadth collapse).** After a handful of genuinely distinct high-utility continuations, additional samples at a node mostly yield near-duplicates whose v scores fall just below the pruning threshold. The frontier then remains narrow even when ample budget is available, leaving compute unused.
2. **Myopic pruning (depth myopia).** Early v estimates are noisy and biased toward near-term payoff; logically consistent branches whose payoff is delayed by several steps are pruned as “low- v ” even though they could mature into correct solutions. This creates *myopic false negatives*.

Both effects amplify with larger inference budgets: saturation wastes more samples as k grows, and myopic pruning discards more candidates as depth increases.

A simple cost asymmetry. Let k be the number of children sampled per expanded node and let a be the acceptance fraction into the *mainline*. If one does not cap mainline width, the expected number of mainline nodes at depth d scales like $(ak)^d$, so the cost to depth D is $\Theta((ak)^D)$ —*exponential in depth*. By contrast, controlling *lateral* width with successive-halving (LR-SC; Sec. 4.3) yields a total lateral exploration cost of $\Theta(N_0 \log_\eta N_0)$ for initial lateral width N_0 and culling factor $\eta > 1$ —*pseudolinear in width*. This asymmetry suggests an architectural principle: *keep mainlines narrow to avoid depth explosion and push width into laterals where it is cheap*.

Why the problem will grow. Three trends sharpen the pain points above:

1. **Bigger inference budgets.** Multi-round agents, tool calls, and safety-/verification-time checks raise the tolerated per-query compute. Without a controller that can convert budget into *productive* breadth, ToT saturates early and the marginal return of extra tokens collapses.
2. **Longer-horizon tasks.** Program synthesis, multi-hop reasoning, and formal verification increasingly require sequences where payoff emerges only after several structured steps. Myopic pruning removes precisely those candidates that need a few steps of nurturing.
3. **Noisier, nonstationary evaluators.** Practical utility scores v (even when outcome-aligned) fluctuate across depths and task regimes. A fixed, level-based gate conflates noise with signal; sequential allocation based on *marginal value of compute* is needed.

A stylized model of the failure mode. Let a candidate node x have an (unobserved) eventual value $\mu(x)$ if its branch were fully developed. An early evaluator observes $v(x) = \mu(x) - \lambda \Delta(x) + \varepsilon$, where $\Delta(x)$ is the (unknown) remaining steps to payoff, $\lambda > 0$ captures horizon bias, and ε is evaluator noise. When $\Delta(x)$ is moderate, $v(x)$ may fall below the mainline gate despite large $\mu(x)$. A controller that reasons about *improvement after a small investment*—rather than $v(x)$ in isolation—can *defer judgment*, test whether x starts producing high- v descendants quickly, and only then commit budget.

Design desiderata induced by the tension. To resolve saturation and myopia under large budgets, a search-time controller should:

1. **Allocate on marginal gain (not level).** Decide to continue a branch based on compute-normalized improvement of an envelope $V(\cdot)$ over a short, controlled lookahead; gate on robust trend (slope/curvature), not a single v reading.
2. **Be wide but short.** Support very large *lateral* width N_0 with near-constant cost per rung and only $\Theta(\log_\eta N_0)$ rungs; immediately *short-circuit* back to exploitation when any lateral reaches the mainline bar.
3. **Keep mainlines narrow.** Beam- or quota-cap mainlines to prevent $(ak)^D$ depth blow-up; re-open exploration only when exploitation *plateaus* in compute-normalized progress.
4. **Promote only on outcome.** Bind promotion to v that is as verifier-aligned as possible (tests, checkers, exact answers), so logically inconsistent but speciously plausible branches do not pollute the mainline.
5. **Control multiplicity.** As lateral width grows, guard against winner’s-curse spikes with width-aware thresholds and a cheap repeat-to-confirm step.

How LToT addresses the gap. LToT operationalizes the desiderata above with two ingredients (see Sec. 4): (i) a *dual-score frontier* that retains logically consistent, low- v *laterals* alongside high- v *mainlines*, deferring judgment on laterals; and (ii) a budgeted racing procedure, *LR-SC*, that allocates tiny probes across a very wide lateral set, culls aggressively, and *promotes* a lateral to the exploitation set the moment its envelope reaches the mainline bar. Theoretical analyses (Sec. 4.5) show that LR-SC keeps lateral cost *pseudolinear in width* ($\Theta(N_0 \log_\eta N_0)$) while mainlines, if left uncapped, are exponential in depth; hence LToT converts surplus compute into principled diversity exactly where it is cheapest.

What the reader should take away. Frontier inference will keep offering more budget per query before training-time improvements alone solve long-horizon reliability. Without a controller, that budget is spent on near-duplicates (saturation) or discarded candidates that only need a few steps (myopia). LToT provides the missing mechanism: *defer judgment* on consistent but low- v ideas, *test them cheaply and in parallel*, and *promote immediately* when they prove themselves—while keeping provable control over compute and errors.

3 RELATED WORK

4 ARCHITECTURE DESIGN

Goal. LToT is a search-time controller for reasoning with language models (LMs) that (i) keeps *mainlines* narrow to avoid exponential blow-up in depth and (ii) makes *lateral* exploration very wide but cheap via a budgeted racing procedure with short-circuit promotion. The controller decides when to exploit mainlines vs. explore laterals, and—during exploration—how to allocate compute across many lateral branches while maintaining guarantees on cost and false promotions.

4.1 PROBLEM SETTING AND NOTATION

We reason over a rooted tree (or DAG) of partial traces. Each node x is a partial solution; its children are produced by prompting the LM with x . Two evaluators score nodes:

$$v(x) \in \mathbb{R} \quad (\text{utility; e.g., answer- or verifier-aligned}), \quad c(x) \in [0, 1] \quad (\text{logical consistency / soundness}).$$

Algorithm 1 LToT controller (high level)

```
1: Inputs: initial frontier  $\mathcal{F}_0$ , evaluator  $v$ , consistency  $c$ , plateau thresholds; LR-SC params  
   ( $\eta, b_0, \rho, \kappa, \delta$ ).  
2: Initialize  $M_0$  with high- $v$  children;  $L_0$  with low- $v$ , high- $c$  children; set bar  $B_0$ .  
3: while budget remains do  
4:   Exploit  $M_t$  while EWMA of  $\Delta B_t$  per compute  $\geq \tau$  (with a small patience & hysteresis).  
5:   Explore laterals with LR-SC over the current lateral pool (Alg. 2).  
6:   if some lateral branch reaches  $v \geq B_t + \delta$  (promotion) then  
7:     add promoted node(s) to  $M_t$ ; update  $B_t$ ; return to exploitation  
8:   else  
9:     freeze survivors for future phases; return to exploitation  
10:  end if  
11: end while
```

We measure compute in either node expansions or tokens and denote cumulative compute by C .

Frontier, origins, and exploitation set. At time t the search maintains a frontier \mathcal{F}_t and an *exploitation set* $M_t \subseteq \mathcal{F}_t$ of nodes eligible for *mainline* exploitation. Nodes carry an immutable *origin* tag in $\{\text{MAINLINE_ORIGIN}, \text{LATERAL_ORIGIN}\}$ indicating how they first entered the frontier. We also maintain a *mainline acceptance bar* B_t (e.g., the best-so-far v or a top- k mean with a small margin $\delta > 0$).

Mainlines vs. laterals. Children with high v are admitted to M_t (mainlines). Children with low v but high c enter the *lateral pool* L_t for potential exploration. Intuitively, laterals represent hypotheses that appear unpromising under a myopic utility but are logically coherent and may become valuable after a short lookahead.

Branch envelope and gain. For a lateral branch i (rooted at node x_i), let $V_i(h)$ denote a branch *envelope*—e.g., a Top- k mean of v among leaves within horizon h steps from x_i (or within a per-branch micro-beam). We write $C(h)$ for the compute required to reach horizon h and define the compute-normalized improvement between horizons $h' < h$ as

$$g_i(h, h') = \frac{V_i(h) - V_i(h')}{C(h) - C(h')}.$$

These quantities let us reason about *marginal value of compute*, not just absolute levels.

4.2 CONTROLLER OVERVIEW

Exploit–explore loop. LToT alternates between:

1. **Mainline exploitation.** Expand nodes from M_t while a compute-normalized progress statistic (e.g., an EWMA of ΔB_t per unit compute) exceeds a plateau threshold. This keeps mainlines narrow (beam- or quota-capped).
2. **Lateral exploration via LR-SC.** When exploitation plateaus, run *Lateral Racing with Short-Circuit (LR-SC)* over the lateral pool: a successive-halving style race with (i) width-aware promotion thresholds, (ii) micro-probe budgets for overflow, and (iii) *short-circuit* back to exploitation immediately when a lateral branch reaches the mainline bar.

Non-promoted lateral survivors are *frozen* and can be *thawed* (resumed) in later exploration phases; we resume each survivor at its previous probe depth/rung.

4.3 LR-SC: OVERFLOW-CAPPED RACING WITH SHORT-CIRCUIT

Let N be the active lateral width. LR-SC proceeds in rungs $r = 0, 1, \dots$ with *culling factor* $\eta > 1$. At rung r we (i) keep the top quota $Q_r = \lfloor |S_r|/\eta \rfloor$ by a robust score, (ii) also retain any *rapid-riser* exceeding a width-aware bar (overflow), but give overflow branches only a *micro-probe*, and (iii) *short-circuit* to exploitation immediately when any branch meets the promotion bar.

Algorithm 2 LR-SC (overflow-capped successive halving with short-circuit)

- 1: **Inputs:** active lateral set S_r (size N), culling factor $\eta > 1$, base budget b_0 , overflow cap ρ , thresholds (κ, δ) , horizon schedule (h_0, h_1, \dots)
 - 2: Compute robust improvement scores $\{z_i\}_{i \in S_r}$ from compute-normalized gains (optionally depth-standardized).
 - 3: $Q_r \leftarrow \lfloor |S_r|/\eta \rfloor$; $T \leftarrow \text{top } Q_r \text{ by } z_i$; $R \leftarrow \{i : z_i \geq \kappa \sqrt{2 \log |S_r|} + \delta\}$.
 - 4: Assign budget $b_{\text{full}} = b_0 \eta^r$ to $i \in T$; assign micro-probe b_{micro} to up to $\lfloor \rho |S_r| \rfloor$ branches in $R \setminus T$ (by z_i); freeze the rest.
 - 5: Expand per budgets to horizon h_r (within-branch beam is tiny); update V_i , g_i , and bar B_t .
 - 6: **if** some i satisfies $V_i \geq B_t + \delta$ and *repeat-to-confirm* **then**
 - 7: promote i ; **short-circuit** to exploitation
 - 8: **end if**
 - 9: $S_{r+1} \leftarrow T \cup (\text{confirmed overflow})$; $r \leftarrow r + 1$; continue if budget remains.
-

Scores and width-aware bar. For branch i at rung r we compute a compute-normalized improvement g_i (using V_i) and a robust standardization z_i (e.g., subtract rung median and divide by a MAD-like scale). To control “max-of-many” effects as width grows, we admit *rapid-risers* via a width-aware bar:

$$z_i \geq \underbrace{\kappa \sqrt{2 \log |S_r|}}_{\text{width penalty}} + \delta,$$

with $\kappa \approx 1$ and a margin $\delta > 0$. We optionally standardize scores within parent-depth bands to compare fairly across heterogeneous depths.

Overflow cap. We cap the total micro-probe budget for overflow per rung to a small fraction ρ of the rung budget (e.g., $\rho \in [0.1, 0.2]$), ensuring per-rung cost stays near constant.

Derivative-based continuation (principle and instantiation). We view V_i as a function of horizon/compute and continue branch i if a discrete derivative of order $m \in \{1, \dots, M\}$ is reliably positive:

$$\widehat{\Delta^{(m)} V_i} \geq \text{bar}(|S_r|, M) \quad \text{with} \quad \text{bar}(|S_r|, M) \propto \sqrt{2 \log(|S_r| \cdot M)}.$$

In practice we cap $M = 2$ for stability and use: (i) first derivative (slope) $s_i = g_i(h_r, h_{r-1})$ and (ii) second derivative (curvature) $\kappa_i = s_i(r) - s_i(r-1)$, estimated over the last few rungs and normalized by compute; a third-order check may be included in an appendix. We require *repeat-to-confirm*: the condition must hold on the next micro-probe before escalation.

4.4 PROMOTION AND SAFETY

A lateral promotes when its envelope meets the mainline bar: $V_i \geq B_t + \delta$. When v is verifier-aligned (e.g., unit tests for code, exact-match for math), this binds promotion to correctness. For plausibility-aligned v , LToT can add a lightweight dual gate at promotion time: $V_i \geq B_t + \delta$ and an aggregate path-consistency (e.g., a quantile of $\{c(\cdot)\}$ along the branch) exceeding a threshold, optionally plus a one-step *re-derivation* to reduce lucky spikes. These checks cost one micro-probe and do not change the asymptotics.

4.5 THEORETICAL PROPERTIES

We summarize the main guarantees; proofs are short and rely on standard successive-halving arguments and sub-Gaussian tail bounds for rung-wise statistics.

Cost law (pseudolinear in lateral width). Let N_0 be the initial lateral width. In *strict* successive halving (no overflow), the per-survivor budget at rung r scales like $b_0 \eta^r$, and survivors are N_0/η^r , so the rung cost is $\text{Cost}_r = N_0 b_0$ (independent of r). With $R = \lceil \log_\eta N_0 \rceil$ rungs, the total lateral cost is

$$\text{Total} = \Theta(N_0 b_0 \log_\eta N_0).$$

In LR-SC with overflow cap $\rho \in (0, 1)$ and micro-probe $b_{\text{micro}} \ll b_0$, the rung cost is at most $(1+\rho)N_0b_0$, hence the same asymptotic order with a constant factor $(1+\rho)$. Short-circuit promotion can only reduce cost. Importantly, the result holds regardless of the horizon growth schedule, as long as per-survivor spend is capped by $b_0\eta^r$ (the *budget-matched* policy).

Rung count (short in depth). The number of rungs required to reduce N_0 laterals to $O(1)$ survivors is $R = \lceil \log_\eta N_0 \rceil$, i.e., logarithmic in lateral width. Thus LR-SC is *wide and short*: constant per-rung cost and $\Theta(\log N_0)$ rungs.

Mainline growth (why we keep mainlines narrow). If at each mainline layer we admit a fixed fraction a of k children (effective reproduction $r_{\text{main}} = ak > 1$), then expansions to depth D are $\Theta(r_{\text{main}}^D)$ (exponential). With a beam/width cap W , mainline cost becomes $\Theta(DWk)$ (linear in depth). LToT therefore keeps W small and invests surplus compute in laterals, where width is cheap.

Width-aware threshold controls family-wise errors. Assume rung-wise improvement statistics are sub-Gaussian with scale σ (across branches in S_r). Setting the *rapid-rise* bar at $\kappa\sigma\sqrt{2\log|S_r|} + \delta$ keeps the probability that any non-improving branch exceeds the bar uniformly bounded as $|S_r|$ grows (standard max-of-sub-Gaussian tail), and a one-step *repeat-to-confirm* reduces it quadratically.

Horizon-lifted detection of delayed payoffs. Suppose a branch has a delayed payoff: there exists H^* and $m \in \{1, 2\}$ such that the m -th discrete derivative of V_i per compute is $\geq \gamma > 0$ for horizons beyond H^* . Under a geometric horizon schedule (e.g., $h_r = 2^r$ within the budget cap) and the derivative-based continuation rule with width-aware thresholds and repeat-to-confirm, the branch is detected and survives to promotion within $O(\log H^*)$ rungs (intuitively, each rung doubles the tested horizon). Total exploration cost remains $\Theta(N_0 \log_\eta N_0)$ because the per-survivor spend never exceeds $b_0\eta^r$.

Independence from nominal horizon multiplier. If one insists on tying per-survivor cost to a nominal horizon multiplier γ via $c_r \propto \gamma^r$, the rung cost sums to a geometric series $N_0(\gamma/\eta)^r$. Thus for $\gamma \leq \eta$ the total remains $O(N_0 \log_\eta N_0)$ (or even $O(N_0)$ when $\gamma < \eta$). In practice we adopt the budget-matched policy: cap spend by $b_0\eta^r$ and allocate within-branch depth/width flexibly up to that cap.

4.6 DESIGN CHOICES AND DEFAULTS

Exploitation trigger. EWMA of compute-normalized mainline progress with small patience and hysteresis; depth-banded statistics if v drifts with depth.

LR-SC parameters. $\eta \in \{3, 4, 5\}$; $b_0 \in \{1, 2\}$ expansions; micro-probe $b_{\text{micro}} = 1$; overflow cap $\rho \in [0.1, 0.2]$; width-aware bar with $\kappa \approx 1.0$ and a small margin δ over the mainline bar. Derivative-based continuation with slope+curvature ($M=2$) using a short window and robust scales; optional third-order check in ablations.

Promotion predicate. Require $V_i \geq B_t + \delta$; if v is not verifier-aligned, add a minimal path-consistency aggregate and a one-step re-derivation check. Both add at most one micro-probe and preserve the asymptotics.

Freeze-thaw. Cache for each survivor: rung index, envelope V_i , recent improvement stats, parent depth, and a lightweight duplicate signature; resume from the same rung during the next exploration phase and evict stale or dominated branches by constant-time tests (e.g., $UCB < B_t - \delta$ for several revisits).

Summary. LToT turns surplus compute into breadth where it is cheapest (laterals) while keeping mainlines narrow. Its LR-SC core provides near-linear (pseudolinear) cost in lateral width, width-aware error control, and immediate promotion when a lateral demonstrably reaches mainline utility.

270	5	EXPERIMENTS
271		
272	6	RESULTS AND DISCUSSION
273		
274	7	FUTURE WORK
275		
276	8	CONCLUSION
277		
278		
279		REFERENCES
280		
281		
282		
283		
284		
285		
286		
287		
288		
289		
290		
291		
292		
293		
294		
295		
296		
297		
298		
299		
300		
301		
302		
303		
304		
305		
306		
307		
308		
309		
310		
311		
312		
313		
314		
315		
316		
317		
318		
319		
320		
321		
322		
323		