

SQUEEZE EVOLVE

Unified Multi-Model Orchestration for Verifier-Free Evolution

Monishwaran Maheswaran^{*,1,5} Leon Lakhani^{*,1} Zhongzhu Zhou^{†,5} Shijia Yang^{†,2} Junxiong Wang⁵
 Coleman Hooper¹ Yuezhou Hu¹ Rishabh Tiwari¹ Jue Wang⁵ Harman Singh¹ Qingyang Wu⁵
 Yuqing Jian⁵ Ce Zhang⁵ Kurt Keutzer¹ Tri Dao^{4,5} Xiaoxia Wu⁵ Ben Athiwaratkun⁵
 James Zou^{‡,3,5} Chenfeng Xu^{*,†,2,5}

¹UC Berkeley ²UT Austin ³Stanford University ⁴Princeton University ⁵Together AI

*Equal contribution †Equal second ‡Co-advising

Correspondence: monishwaran@berkeley.edu, xuchenfeng@utexas.edu

Abstract

We show that verifier-free evolution is bottlenecked by both diversity and efficiency: without external correction, repeated evolution accelerates collapse toward narrow modes, while the uniform use of a high-cost model wastes compute and quickly becomes economically impractical. We introduce SQUEEZE EVOLVE, a unified multi-model orchestration framework for verifier-free evolutionary inference. Our approach is guided by a simple principle: allocate model capability where it has the highest marginal utility. Stronger models are reserved for high-impact stages, while cheaper models handle the other stages at much lower costs. This principle addresses diversity and cost-efficiency jointly while remaining lightweight. SQUEEZE EVOLVE naturally supports open-source, closed-source, and mixed-model deployments. Across AIME 2025, HMMT 2025, LiveCodeBench V6, GPQA-Diamond, ARC-AGI-2, multimodal vision benchmarks (MMMU-Pro, BabyVision), and agentic settings (Terminal-Bench 2.0), SQUEEZE EVOLVE consistently improves the cost-capability frontier over single-model evolution and achieves new state-of-the-art results on several tasks. Empirically, SQUEEZE EVOLVE reduces API cost by up to $\sim 3\times$ and increases fixed-budget serving throughput by up to $\sim 10\times$. Moreover, on discovery tasks, SQUEEZE EVOLVE is the first verifier-free evolutionary method to match, and in some cases exceed, the performance of verifier-based evolutionary methods.

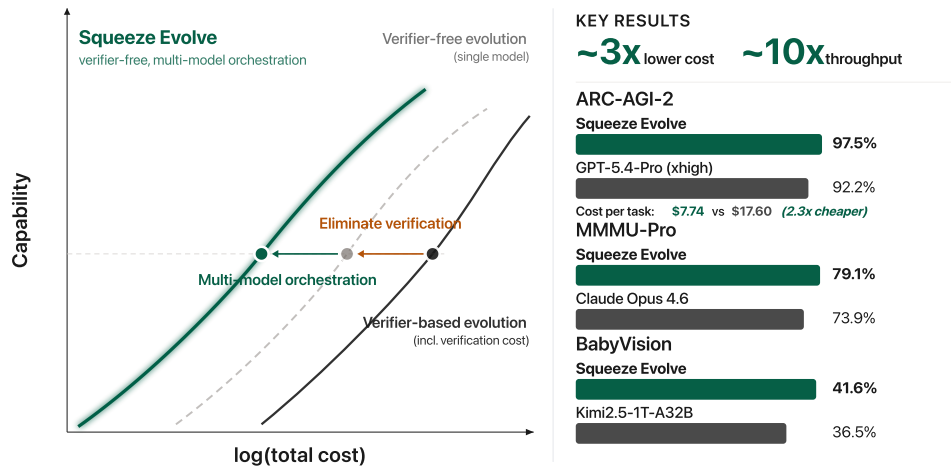


Figure 1. SQUEEZE EVOLVE shifts the cost-capability frontier left by combining verifier-free evolution with multi-model orchestration. **Left:** Conceptual scaling curves. **Right:** Key results across ARC-AGI-2, MMMU-Pro, and BabyVision.

1 Introduction

Test-time scaling has emerged as a practical way to push language models beyond one-shot inference by spending additional compute at test time to search over or refine candidate solutions [53, 30, 51]. A particularly promising direction is self-evolution, where models iteratively improve candidates through selection, mutation, and recombination [34, 42, 24, 28]. When coupled with an external verifier, this paradigm can unlock powerful discovery capabilities. But in many important domains, verification is too expensive and slow, or simply unavailable. For example, in nuclear fusion research, a single tokamak plasma study may require more than 120 million CPU-hours [18]. This motivates our focus on *verifier-free evolution*.

However, verifier-free evolution is approaching a hard wall: methods such as RSA generate $500\text{--}700\times$ more tokens than single-shot inference, making cost increasingly difficult to sustain. This pressure is compounded by the sharp cost gap between models: proprietary reasoning models remain $4\text{--}25\times$ more expensive per output token than strong hosted open-weight alternatives [37, 3, 14, 48, 49], and even within the open-weight ecosystem cost differences remain substantial.

As a result, the key question is shifting: rather than only asking *how we can spend more compute and money to unlock new capabilities?*, we must also ask *how we can achieve a given capability target under tight budget constraints?* This is the same principle that has historically driven advances in software and algorithms: progress comes not just from using more resources, but from using them more efficiently and lowering the cost needed to achieve a given capability target.¹ In this work, we advance this principle, as illustrated in Fig. 1.

To answer the above question, we first take a system perspective: many seemingly disparate test-time methods can be expressed as instances of a single evolutionary framework. Once cast in this unified form, they expose a common design space that can be optimized jointly. In Sec. 3, we describe how we unify the current test-time scaling method into a single evolutionary framework, where different operator choices recover a wide spectrum of existing test-time strategies. For example, majority voting [53] corresponds to a shallow single-step evolution, recursive self-aggregation [51] corresponds to an verifier-free multi-step evolutionary process, and verifier-based self-evolve pipelines such as AlphaEvolve [34] correspond to feedback-driven evolutionary search. Our unified framework naturally highlights the key problems:

1. Given models with different cost–capability trade-offs, which model should be assigned to each operator in the evolutionary pipeline (e.g., initialization, generation, recombination, or fitness estimation)?
2. How should these models be coordinated across the pipeline to maximize capability per unit cost without incurring excessive orchestration overhead?

We answer these two questions through a comprehensive empirical analysis in Section 4. In brief, we find that:

1. **From the verification perspective**, scaling the token budget can partially offset the absence of explicit verification. By spending additional tokens on diverse generation and iterative aggregation, verifier-free evolution can converge reliably toward correct solutions even without external reward signals. This makes verifier-free evolution especially attractive in practice, as it improves capability while avoiding the substantial cost of explicit verification.
2. **From the performance perspective**, unlike verifier-based method, simple verifier-free evolution causes the upper bound (e.g., *pass@N* or *best continuous score*) to degrade significantly. Such a drop directly limits the achievable performance of the entire pipeline. We further find that this upper bound is highly correlated with generation diversity, highlighting diversity as a central ingredient for effective verifier-free evolution. This further strongly motivates our use of *multi-model orchestration* to preserve diversity and sustain performance.
3. **From the cost perspective**, initialization quality largely determines final accuracy, recombination capability depends on the candidate set being aggregated, and self- and cross-model internal signals serve as reliable fitness proxies. These findings provide a foundation for principled orchestration design.

¹<https://epochai.substack.com/p/the-least-understood-driver-of-ai>

Motivated by these observations and by the economic mismatch between open and closed model ecosystems, we present SQUEEZE-EVOLVE, a multi-model orchestration framework that routes each evolutionary operation to the most cost-effective model based on confidence-derived fitness signals, reserving expensive models for only the highest-marginal-utility steps.

We evaluate SQUEEZE EVOLVE across AIME 2025, HMMT 2025, GPQA-Diamond, LiveCodeBench V6, MMMU-Pro, BabyVision, ARC-AGI-2, circle packing, and the agentic Terminal-Bench 2.0, spanning open-source model pairs, mixed open-source and proprietary model pairs, multimodal vision, and agentic settings. In summary, we make the following contributions:

1. **Unified framework and analysis.** We unify existing test-time scaling methods into a single evolutionary framework, identify the key design axes for multi-model orchestration, and show that diversity collapse is the central bottleneck of verifier-free evolution while model-intrinsic confidence signals serve as effective fitness proxies for routing (Sections 3, 4).
2. **Confidence-based routing.** A lightweight mechanism that assigns each recombination group to the most cost-effective model using only signals already produced during inference (Section 5).
3. **Empirical gains across nine benchmarks.** Across math (AIME 2025, HMMT 2025, GPQA-Diamond), coding (LiveCodeBench V6), vision (MMMU-Pro, BabyVision), visual reasoning (ARC-AGI-2), and scientific discovery (circle packing), SQUEEZE EVOLVE reduces API cost by 1.3–3.3 \times while preserving or exceeding single-model accuracy, and in multiple configurations *surpasses* the expensive Model 2 alone (Section 6). On multimodal tasks, a text-only cheap Model 1 that never processes images matches or exceeds the vision-capable Model 2 at 2.3–2.5 \times savings, showing visual understanding is primarily needed at initialization (Section 6.2). On ARC-AGI-2, SQUEEZE EVOLVE sets a new state-of-the-art cost-capability frontier of 97.5% at \$7.74/task without code execution (Section 6.3); on circle packing, it is the first verifier-free method to match or exceed verifier-based approaches such as AlphaEvolve (Section 6.4). On the agentic Terminal-Bench 2.0 [47], applied as a pure recombination layer on top of an existing GPT-5.4 submission with no re-execution, SQUEEZE EVOLVE lifts pass@1 from 81.80% to 88.76% (+6.96 pp), recovering 99% of the available oracle headroom (Section 6.5).
4. **Co-designed serving.** A custom confidence engine reduces scoring latency by 4–10 \times , latency-matched GPU pools prevent bottlenecks, end-to-end routing overhead is only 2.4–4.3%, and fixed-budget serving throughput improves by up to $\sim 10\times$ (Sections 5.2, 7).

2 Related Work

Our work builds on three lines of research (extended discussion in Appendix A). Self-aggregation methods [51, 52] use a single model or fixed assignment, leading to diversity collapse [43]. LLM-driven evolutionary methods [40, 34, 28] rely on external verifiers and apply one model across all operators. Model routing [35, 32] dispatches queries or steps between models, and token-level confidence has been used to filter traces [13]. SQUEEZE EVOLVE repurposes confidence for routing within a verifier-free evolutionary pipeline, where savings compound across loops.

3 A Unified Formulation of Evolutionary Framework

We show that many existing methods can be cast within a common evolutionary framework, providing a formal foundation for joint optimization and a simple orchestration rule: invoke the larger model only when the smaller model is likely to exceed its capability limits.

For a query Q , we initialize a population $\mathcal{P}^{(0)}$ using an *ancestor* function p_F , where each candidate $\tau_i^{(0)} \sim p_\theta(\cdot | Q)$ is sampled from a generative model θ . Existing methods differ primarily in how they organize, score, and evolve these candidates. We unify these steps into a single *evolutionary operator* Φ_f , which encapsulates selection followed by recombination:

$$\Phi_f(\mathcal{P}) = \text{recomb}_f \circ \text{select}_f(\mathcal{P}), \quad \mathcal{P}^{(t+1)} = \Phi_{f_{t+1}}(\mathcal{P}^{(t)}). \quad (1)$$

This induces an iterated evolutionary process where the final population $\mathcal{P}^{(k)}$ is derived via a sequence of operator compositions:

$$\mathcal{P}^{(k)} = (\Phi_{f_k} \circ \Phi_{f_{k-1}} \circ \dots \circ \Phi_{f_1})(\mathcal{P}^{(0)}), \quad (2)$$

Table 1. Instantiation of the unified evolutionary framework for test-time scaling methods.

Method	k	$ \mathcal{P} $	select	recombination	Fitness f	Model
Majority Voting	1	N	Answer clustering	Identity	Consensus frequency	Single
Self-Refinement	T	1	Self-critique	Conditioned rewrite	Natural language critique	Single
RSA	T	N	K -subset	LLM aggregation	Implicit	Single
AlphaEvolve	T	Variable	Fitness-guided	LLM aggregation	External h	Multi-model
SQUEEZE EVOLVE	T	Variable	Fitness-guided	Mix of recombination	Probabilistic fitness	Multi-model

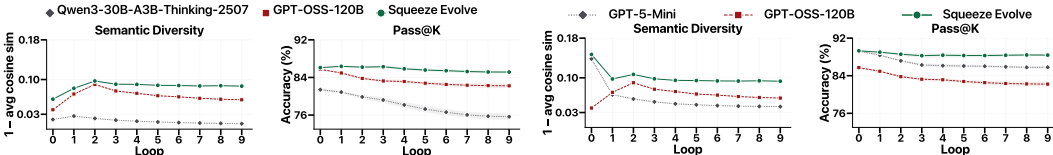


Figure 2. **Single-model evolution collapses diversity and shrinks pass@ K ; multi-model routing preserves both.** Single-model baselines lose diversity after the early loops with a corresponding decline in pass@ K , whereas SQUEEZE EVOLVE remains higher and flatter on both metrics.

where each operator Φ_{f_i} utilizes the fitness signal f_i to transition between generations. In verifier-free evolution, f is derived entirely from model outputs (e.g., log-probabilities, consensus frequency) without external verification. Table 1 categorizes existing methods by how they instantiate select, recomb, and f .

4 Motivation analysis for verifier-free evolutionary framework

The inherent Pass@ K bottleneck of verifier-free evolution. Without an external verifier, evolution can only amplify trajectories the model already recognizes, driving the population to be narrow. As Figure 2 shows across both GPQA-Diamond [39] settings, pass@ K falls alongside semantic diversity. Multi-model orchestration counteracts this: by combining models with different priors and failure modes, SQUEEZE EVOLVE maintains complementary lineages and stays higher and flatter on both diversity and pass@ K .

Ancestor function dominates final accuracy. Using the strong model for initialization and the weak model for recombination consistently outperforms the reverse: GPT-OSS-120B [1]→20B achieves 0.89 vs. 0.85 on HMMT 2025 [5], and Qwen3-4B-Thinking [46]→Instruct reaches 0.88 vs. 0.65 on AIME 2025 (Table 2, Appendix). This asymmetry motivates reserving the strong model for initialization.

Weak models suffice when candidates are strong. Recombination accuracy rises from 0% with no correct candidate to 100% when all four are correct (Fig. 3a), consistently across benchmarks and model families. This motivates routing high-quality subsets to a cheaper model.

Confidence as a fitness proxy. Both self- and cross-model confidence reliably separate high-quality from low-quality subsets (Figure 3b), motivating their use as the fitness signal for routing.

5 Squeeze Evolve

We instantiate Φ_f from Section 3 as a single algorithm (Figure 4; Algorithm 1, Appendix E). The key extension is fitness-based routing: each candidate group is assigned to one of N tiers, $\text{Model}_{1..N-1}$ (ordered by increasing cost) or lightweight aggregation, based on the fitness signal (Table 3).

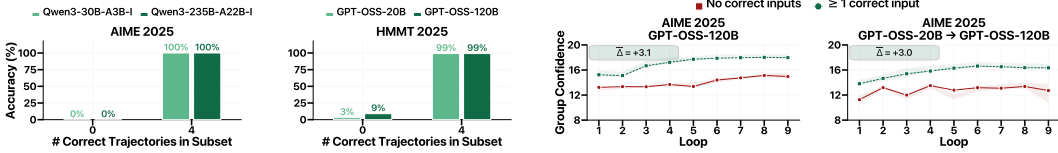
5.1 Algorithm

Each loop scores candidates via the fitness signal f , applies select_f to form groups, routes each group to one of three recombination tiers within recomb_f , and updates the population. We define each component below.

Initialization. The initial population is generated entirely by the target model: $\mathcal{P}_q^{(0)} = \{\tau_i \sim p_{M_2}(\cdot | Q_q)\}_{i=1}^N$. This is motivated by the finding that initialization quality is the single strongest predictor of final accuracy (Table 2).

Table 2. **Ancestor function dominates final accuracy.** Mean final loop (9) accuracy across 4 seeds. Strong-init \rightarrow weak-agg outperforms weak-init \rightarrow strong-agg, indicating initialization quality dominates aggregation quality.

Model pair	Data	S \rightarrow W	W \rightarrow S	Δ
GPT-OSS-120B / GPT-OSS-20B	HMMT’25	0.89	0.85	+4
Qwen3-4B-Thinking-2507 / Qwen3-4B-Instruct-2507	AIME’25	0.88	0.65	+23



(a) Aggregation accuracy at the extremes: 0 vs. 4 correct trajectories.

(b) Group confidence separates correct from incorrect trajectories.

Figure 3. **Aggregation success is seed-dependent, and group confidence predicts it.** (a) With zero correct seeds neither model recovers a correct answer; with all four correct both reach near-perfect accuracy (full 0–4 breakdown in Appendix Fig. 10). (b) Mean group confidence across loops 1–9 on AIME 2025: correct-containing subsets maintain consistently higher GC ($\bar{\Delta} \geq +3.0$) in both self- and cross-model settings (Appendix Figs. 11, 12).

Fitness signal. SQUEEZE EVOLVE uses two model-intrinsic fitness signals, both proxies for group difficulty. *Group confidence* (GC) averages per-token confidence over a trajectory and then over the group:

$$c(i) = -\frac{1}{K_\ell} \sum_{j=1}^{K_\ell} \log p_\theta(v_j^{(i)} | t_{<i}, Q), \quad C(\tau) = \frac{1}{|\tau|} \sum_i c(i), \quad GC(g) = \frac{1}{K} \sum_{\tau \in g} C(\tau). \quad (3)$$

where $\{v_1^{(i)}, \dots, v_{K_\ell}^{(i)}\}$ are the K_ℓ most likely tokens under scoring model θ . The per-token confidence $c(i)$ follows the same formulation used by [13] to filter reasoning traces; here we aggregate it to the group level for routing. When θ is the generator itself this is *self-confidence* (zero cost); otherwise it is *cross-model confidence* (one prefill pass per candidate; Section 5.2). When log-probabilities are unavailable, we instead use *group diversity* $D(g) = |\{\text{answer}(\tau) : \tau \in g\}|$, the number of unique final answers. The choice between the two signals is determined by API access.

Selection. At each loop $t \geq 1$, we form M groups of size K from the current population. Groups can be formed by *uniform sampling* (random K -subsets) or by *fitness-weighted sampling*, where candidates are drawn with probability $\exp(f(\tau_i)/\zeta) / \sum_j \exp(f(\tau_j)/\zeta)$ and a temperature ζ controls the exploitation–exploration balance.)

Recombination. Groups with sufficient consensus are aggregated via lightweight methods ($\mathcal{B}_{\text{lite}}$). The remaining groups are split by a per-problem adaptive threshold $\theta_q = \text{Percentile}_p(\{F(g) : g \in \mathcal{G}_q \setminus \mathcal{B}_{\text{lite}}\})$: groups that are “easy” under the fitness signal (high GC or low D) are routed to Model 1, and the rest to Model 2. The routing percentile p is the single deployment-time hyperparameter. The two model batches execute in parallel, and trajectories from all three tiers are merged via either *replace* or *accumulate* ($\mathcal{P}_q^{(t)} = \mathcal{P}_q^{(t-1)} \cup \mathcal{R}_{\text{new}}$).

5.2 System Implementation

Self-confidence is free: the generating model already produces the token log-probabilities needed to score its own trajectory. *Cross-model confidence* requires a single prefill-only forward pass per trajectory under a different model; because that model is already resident for its aggregation branch, this adds no memory overhead, and the N scoring calls per loop are batched into a single request.

Standard serving systems [23, 61] are optimized for decode-heavy generation, but confidence scoring needs just one scalar per trajectory. To avoid materializing full token-level logprob tensors, we implement a custom prefill path in vLLM that accumulates the confidence statistic directly on GPU and returns only the final scalar, reducing per-request transfer from ~ 13 MB to ~ 100 bytes. This achieves 4–10 \times lower scoring latency and enables confidence scoring on Qwen3-235B-A22B where the native path runs out of memory (Appendix L). We quantify end-to-end routing overhead and system throughput in Section 7.

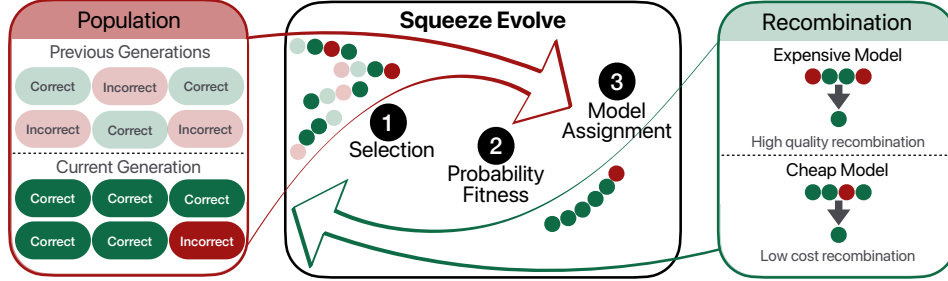


Figure 4. SQUEEZE EVOLVE overview. The expensive Model 2 generates the initial population; subsequent loops recombine groups using Model 1 and 2 based on group confidence

Table 3. Operator instantiations of SQUEEZE EVOLVE across evaluation settings.

Setting	Fitness f	Select f	Recomb f (routing rule)	Update
Math / Coding / Vision	GC (Eq. 3)	Uniform	Percentile on GC	Replace
ARC-AGI-2 (§6.3)	Diversity D (Eq. 5.1)	Uniform	Threshold on D + lite agg	Replace
Circle Packing (§6.4)	GC (Eq. 3)	Fitness-weighted ($\zeta=0.5$)	Percentile on GC	Accumulate

6 Evaluation

All runs use population $N=16$, group size $K=4$, and $T=10$ evolutionary loops, averaged over four seeds, unless stated otherwise. Costs are measured in actual API dollars per problem using model provider pricing (Table 8; generation hyperparameters in Table 7, Appendix). The baseline is standard RSA with Model 2 only, which serves as the cost upper bound.

6.1 Math and Coding

We evaluate SQUEEZE EVOLVE on AIME 2025, HMMT 2025 [5], GPQA-Diamond [39], and LiveCodeBench V6 [21]. Accuracy-vs-cost curves appear in Figures 5 and 6; full per-percentile breakdowns in Tables 9 and 10 (Appendix).

Homogeneous pairs (open-source). We test three pairs spanning different design axes: Qwen3-30B Instruct / Thinking (same scale, different mode), Qwen3-30B / 235B Instruct (different scale, same mode), and GPT-OSS-20B / 120B (different scale, both thinking). SQUEEZE EVOLVE matches or exceeds Model 2 alone at **1.4–2.1× less** cost, surpassing Model 2 on two of three pairs (+1.5 pp on AIME25 for Instruct→Thinking; +2.3 pp on HMMT25 for GPT-OSS). Even with Model 1 much smaller (30B vs. 235B), accuracy stays within 1 point at nearly half the cost; the GPT-OSS pair matches Model 2 on LiveCodeBench V6 at 2.0× savings.

Heterogeneous pairs (open-source + closed-source). We pair Qwen3-30B Instruct and GPT-OSS-20B with GPT-5 mini [36] as Model 2, sweeping $p \in \{0, 10, 20, 30\}$ (Model 1 scores via prefill since GPT-5 mini does not expose output logprobs; this cost is included). SQUEEZE EVOLVE achieves **1.4–3.3× savings**: at $p=30$, GPT-OSS-20B + GPT-5 mini exceeds Model 2 on AIME25 (95.4% vs. 94.2%) at 1.8× savings; at $p=0$, savings reach 3× with accuracy drops of only 1.5–6 pp (Table 10).

6.2 Multimodal Vision Task

We evaluate SQUEEZE EVOLVE on two multimodal benchmarks: MMMU-Pro [57] and BabyVision [9], using $T=5$ loops (other settings match Section 6.1). We test a homogeneous pair (Kimi-2.5 Instant / Thinking [45], both vision-capable) and a heterogeneous pair (Qwen3.5-35B [38] → Kimi-2.5 Thinking, where Model 1 operates in *text-only mode* after loop 0).

Accuracy-vs-cost curves for MMMU-Pro appear in Figure 7. On MMMU-Pro, the homogeneous pair matches Model 2 at **1.9×** savings, while the heterogeneous pair *surpasses* Model 2 (79.1% vs. 78.6%) at **2.7×** savings, even though Model 1 never sees any images. On BabyVision, the homogeneous pair preserves accuracy at **2.5×** savings. The heterogeneous result further reinforces the finding from Section 4 that initialization quality is the dominant factor: once loop 0 grounds the population in image content, subsequent aggregation can be delegated to a cheaper text-only model. Full breakdowns appear in Tables 11 and 12 (Appendix).

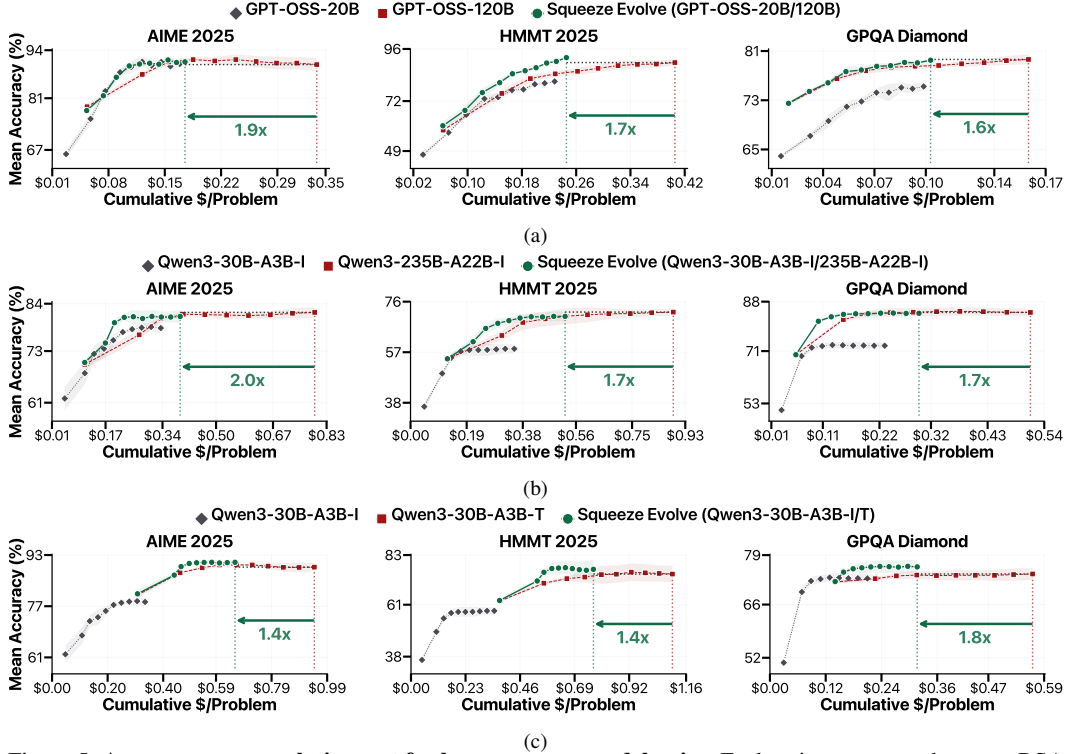


Figure 5. Accuracy vs. cumulative cost for homogeneous model pairs. Each point corresponds to one RSA loop (0–9). SQUEEZE EVOLVE (green) tracks the RSA accuracy curve while staying significantly further left, achieving comparable quality at 1.4–2.0× lower cost.

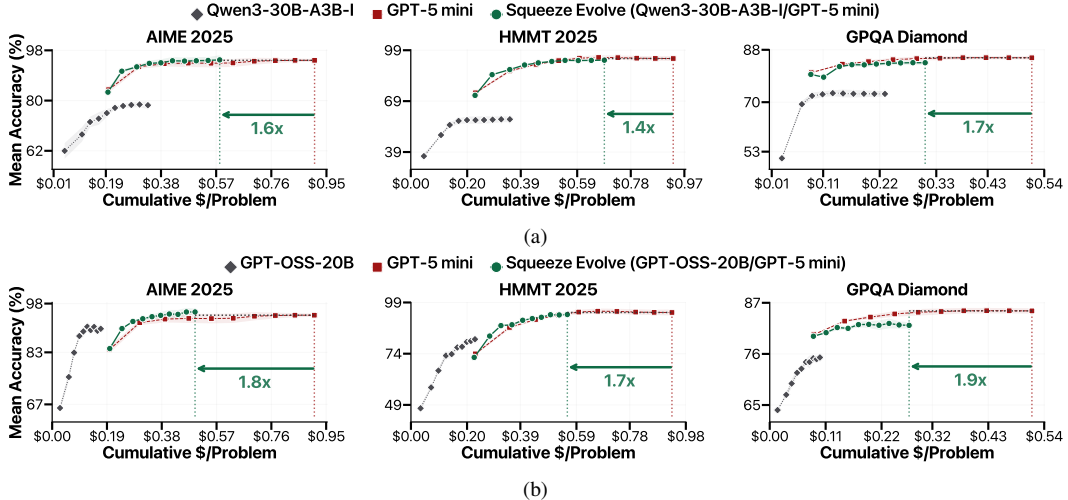


Figure 6. Accuracy vs. cumulative cost for heterogeneous model pairs. SQUEEZE EVOLVE (green) matches the expensive curve at 1.4–1.9× lower cost, demonstrating that confidence-based routing generalizes across model families and access types.

6.3 ARC-AGI-2

We evaluate SQUEEZE EVOLVE on ARC-AGI-2 [10] public evaluation set. Since the Gemini API does not expose logprobs, we use answer diversity (Eq. 5.1) as the fitness signal (Table 3). Groups with non-zero diversity are recombined by Gemini 3.1 Pro [16]; consensus groups fall back to majority vote. With this routing, SQUEEZE EVOLVE achieves 97.5% at \$7.74/task.

Using this result as a baseline, we further reduce cost by adding Gemini 3.0 Flash [15] as Model 1 to the recombination function, yielding a three-way routing rule: high-diversity groups with 3

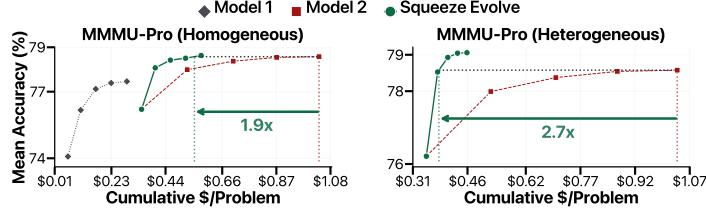


Figure 7. Accuracy vs. cumulative cost on MMMU-Pro for homogeneous and heterogeneous vision pairs. Left: Kimi-2.5 Instant (Model 1)/Thinking (Model 2). Right: Qwen3.5-35B (Model 1, text-only) → Kimi-2.5 Thinking (Model 2). Savings are measured at matched accuracy. The heterogeneous pair achieves 2.7× savings despite Model 1 never seeing any images.

Table 4. ARC-AGI-2 public evaluation results. $N=4$, $K=2$. †Uses code execution and program synthesis. Extended results in Table 16 (Appendix).

Strategy	Model 1	Model 2	Acc.	\$/Task	Savings
<i>Code-execution methods</i>					
Imbue [†]	—	Gemini 3.1 Pro	95.1	\$8.71	—
Confluence Lab [†]	—	—	97.9	\$11.77	—
<i>Full pipeline (T=10)</i>					
RSA	—	Gemini 3.1 Pro	93.3	\$28.85	1.0×
SQUEEZE EVOLVE	—	Gemini 3.1 Pro	97.5	\$7.74	3.7×
<i>Single recombination (T=2)</i>					
SQUEEZE EVOLVE	—	Gemini 3.1 Pro	94.2	\$5.62	5.1×
SQUEEZE EVOLVE	Gemini 3.0 Flash	Gemini 3.1 Pro	97.5	\$5.93	4.9×

or more unique answers invoke the expensive Gemini 3.1 Pro Model 2, lower diversity groups are handled by Gemini 3.0 Flash, and groups that have already reached consensus are aggregated via lightweight non-LLM methods (e.g., majority vote). With this recombination function, we observe immediate convergence to the pass@k score after one aggregation step, achieving the same 97.5% accuracy for only \$5.93/task, a 1.2× savings. This is a new SoTA cost-capability frontier result on ARC-AGI-2 public evaluation set to date. Even compared to code-execution-based approaches, SQUEEZE EVOLVE reaches comparable accuracy at a lower cost to Confluence Lab [12] (97.9%, \$11.77/task) and Imbue [19] (95.1%, \$8.71/task).

6.4 Circle Packing: Scientific Discovery

Table 5. Comparison of methods on Circle Packing ($n=26$, †). ShinkaEvolve uses an ensemble of Claude Sonnet-4, GPT-4.1, GPT-4.1-mini, GPT-4.1-nano, and o4-mini.

Method	Model	Score (†)
ShinkaEvolve [24]	Ensemble (see caption)	2.635982
SQUEEZE EVOLVE	GPT-OSS-120B + 20B	2.635896
AlphaEvolve [34]	Gemini-2.0 Pro + Flash	2.635862
OpenEvolve [42]	Gemini-2.0 Flash + Claude 3.7 Sonnet	2.634292

We apply SQUEEZE EVOLVE to the circle packing problem from AlphaEvolve [34]: pack $n = 26$ non-overlapping circles in a unit square to maximize the sum of their radii. This is an open-ended optimization with a continuous objective. We use GPT-OSS-20B as Model 1 and GPT-OSS-120B as Model 2; operator choices appear in P.2. SQUEEZE EVOLVE is competitive with leading evolutionary frameworks (Table 5) without executing generated programs or using closed-source models. Unlike other frameworks that feed execution results back into generation, SQUEEZE EVOLVE relies solely on model-intrinsic confidence, which correlates weakly but sufficiently with the objective to improve both average and best scores over iterations. Details and discovered program appear in Appendix P.

6.5 Terminal-Bench: Agentic Setting

We apply SQUEEZE EVOLVE to Terminal-Bench 2.0 [47]: 89 software-engineering tasks in sandboxed terminals with deterministic pass/fail checks. We treat the five public ForgeCode / GPT-5.4 trials per

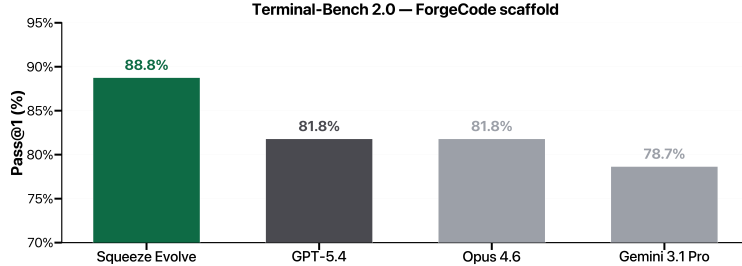


Figure 8. Terminal-Bench 2.0 pass@1 on the ForgeCode scaffold. SQUEEZE EVOLVE lifts pass@1 from 81.80% to 88.76% (+6.96 pp) by recombining existing GPT-5.4 trials, without re-running the agent.

task as $\mathcal{P}_q^{(0)}$ and apply a single recombination step with Gemini 3 Flash as Model 2; the operator reads each trajectory’s exit flags and stdout/stderr tail, and never re-executes in the sandbox. SQUEEZE EVOLVE lifts pass@1 from 81.80% to 88.76% (+6.96 pp) over the public leaderboard (Figure 8), without re-running the agent. Unlike approaches that scale agentic settings by spawning more rollouts, SQUEEZE EVOLVE reuses existing trajectories: of the 80 tasks with at least one passing trial (oracle pass@5 = 89.89%), recombination emits a passer in 79, recovering 99% of the available headroom.

7 System Results

Routing overhead. A natural systems question is whether confidence scoring and model dispatch introduce enough additional latency to undermine multi-model routing. To isolate this cost, we compare RSA with Model 2 (no routing) against SQUEEZE EVOLVE with Model 2 (routing enabled but all aggregation forced to Model 2), so the only difference is the routing logic itself. Figure 16 shows that routing adds only 2.4–4.3% to total latency across all three target models, confirming that confidence scoring is a batched prefill-only operation whose cost is negligible relative to generation (full details in Appendix M).

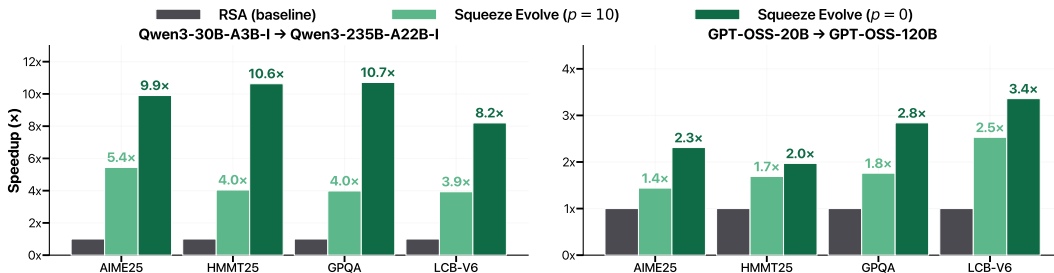


Figure 9. **Fixed-budget throughput speedup over RSA.** Under the same total GPU budget, the Qwen pair achieves 4–10× speedup and the GPT-OSS pair 1.4–3.4×. Full results in Table 15 (Appendix).

System throughput. We next ask whether routing improves steady-state throughput under a fixed GPU budget G . RSA assigns all G to Model 2, while SQUEEZE EVOLVE splits them into Model 2 and Model 1 pools ($G_2 + G_1 = G$) with matched loop service times: under-provisioning either side bottlenecks the system and erases the gains from cheaper aggregation. We report requests per second because the two models emit different output lengths. Figure 17 shows Qwen3-30B/235B reaches 4–10× speedup and GPT-OSS-20B/120B 1.4–3.4×; the gap tracks Qwen’s larger asymmetry on both total and active parameters (full breakdowns in Appendix N).

8 Conclusion

SQUEEZE EVOLVE is a multi-model framework that uses confidence-based routing to make verifier-free evolutionary inference substantially cheaper without sacrificing quality. Across math, coding, vision, agentic, and discovery tasks, it cuts cost by 1.4–3.3× while matching or exceeding the strong model alone, sets a new SoTA on ARC-AGI-2, matches verifier-based frameworks on circle packing, and yields up to 10× throughput gains at negligible routing overhead.

References

- [1] Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K. Arora, Yu Bai, Bowen Baker, and Haiming Bao et al. gpt-oss-120b & gpt-oss-20b model card, 2025.
- [2] Lakshya A Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziemens, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G. Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. Gepa: Reflective prompt evolution can outperform reinforcement learning, 2026.
- [3] Anthropic. Pricing. <https://platform.claude.com/docs/en/about-claude/pricing>, 2026. Claude API pricing page, accessed March 16, 2026.
- [4] Henrique Assumpção, Diego Ferreira, Leandro Campos, and Fabricio Murai. Codeevolve: an open source evolutionary coding agent for algorithmic discovery and optimization, 2026.
- [5] Mislav Balunović, Jasper Dekoninck, Ivo Petrov, Nikola Jovanović, and Martin Vechev. Matharena: Evaluating llms on uncontaminated math competitions, 2026.
- [6] Hritik Bansal, Arian Hosseini, Rishabh Agarwal, Vinh Q. Tran, and Mehran Kazemi. Smaller, weaker, yet better: Training llm reasoners via compute-optimal sampling, 2024.
- [7] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling, 2024.
- [8] Mert Cemri, Shubham Agrawal, Akshat Gupta, Shu Liu, Audrey Cheng, Qiuyang Mang, Ashwin Naren, Lutfi Eren Erdogan, Koushik Sen, Matei Zaharia, Alex Dimakis, and Ion Stoica. Adaevolve: Adaptive llm driven zeroth-order optimization, 2026.
- [9] Liang Chen, Weichu Xie, Yiyan Liang, Hongfeng He, Hans Zhao, Zhibo Yang, Zhiqi Huang, Haoning Wu, Haoyu Lu, Y. Charles, Yiping Bao, Yuantao Fan, Guopeng Li, Haiyang Shen, Xuanzhong Chen, Wendong Xu, Shuzheng Si, Zefan Cai, Wenhao Chai, Ziqi Huang, Fangfu Liu, Tianyu Liu, Baobao Chang, Xiaobo Hu, Kaiyuan Chen, Yixin Ren, Yang Liu, Yuan Gong, and Kuan Li. Babyvision: Visual reasoning beyond language, 2026.
- [10] François Chollet, Mike Knoop, Greg Kamradt, and Chris Landers. ARC-AGI-2: A new challenge for frontier AI reasoning systems. *arXiv preprint arXiv:2505.11831*, 2025.
- [11] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021.
- [12] Confluence Labs. State-of-the-art ARC-AGI-2 solver, 2026. GitHub repository, accessed March 2026.
- [13] Yichao Fu, Xuewei Wang, Yuandong Tian, and Jiawei Zhao. Deep think with confidence, 2025.
- [14] Google. Gemini developer api pricing. <https://ai.google.dev/gemini-api/docs/pricing>, 2026. Google AI for Developers pricing page, accessed March 16, 2026.
- [15] Google DeepMind. Gemini 3 Flash model card. Technical report, Google DeepMind, December 2025.
- [16] Google DeepMind. Gemini 3.1 Pro model card. Technical report, Google DeepMind, February 2026.
- [17] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shiron Ma, and Xiao et al. Bi. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, September 2025.
- [18] N. T. Howard, C. Holland, A. E. White, M. Greenwald, and J. Candy. Multi-scale gyrokinetic simulation of tokamak plasmas: enhanced heat loss due to cross-scale coupling of plasma turbulence. *Nuclear Fusion*, 56, 2016.
- [19] Imbue. Beating ARC-AGI-2 with code evolution, 2026. Blog post, accessed March 2026.
- [20] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, and Alex Carney et al. Openai o1 system card, 2024.
- [21] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code, 2024.

- [22] Ammar Khairi, Daniel D’souza, Marzieh Fadaee, and Julia Kreutzer. Making, not taking, the best of n, 2025.
- [23] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [24] Robert Tjarko Lange, Yuki Imajuku, and Edoardo Cetin. Shinkaevolve: Towards open-ended and sample-efficient program evolution, 2025.
- [25] Joel Lehman, Jonathan Gordon, Shawn Jain, Kamal Ndousse, Cathy Yeh, and Kenneth O. Stanley. Evolution through large models, 2022.
- [26] Zichong Li, Xinyu Feng, Yuheng Cai, Zixuan Zhang, Tianyi Liu, Chen Liang, Weizhu Chen, Haoyu Wang, and Tuo Zhao. Llms can generate a better answer by aggregating their own responses, 2025.
- [27] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step, 2023.
- [28] Shu Liu, Shubham Agarwal, Monishwaran Maheswaran, Mert Cemri, Zhifei Li, Qiuyang Mang, Ashwin Naren, Ethan Boney, Audrey Cheng, Melissa Z. Pan, Alexander Du, Kurt Keutzer, Alvin Cheung, Alexandros G. Dimakis, Koushik Sen, Matei Zaharia, and Ion Stoica. Evox: Meta-evolution for automated discovery, 2026.
- [29] Jack Lu, Ryan Teehan, Jinran Jin, and Mengye Ren. When does verification pay off? a closer look at llms as solution verifiers, 2025.
- [30] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback, 2023.
- [31] Lovish Madaan, Aniket Didolkar, Suchin Gururangan, John Quan, Ruan Silva, Ruslan Salakhutdinov, Manzil Zaheer, Sanjeev Arora, and Anirudh Goyal. Rethinking thinking tokens: Llms as improvement operators, 2025.
- [32] Monishwaran Maheswaran, Rishabh Tiwari, Yuezhou Hu, Kerem Dilmen, Coleman Hooper, Haocheng Xi, Nicholas Lee, Mehrdad Farajtabar, Michael W. Mahoney, Kurt Keutzer, and Amir Gholami. Arbitrage: Efficient reasoning via advantage-aware speculation, 2025.
- [33] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling, 2025.
- [34] Alexander Novikov, Ngán Vű, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian, M. Pawan Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian Nowozin, Pushmeet Kohli, and Matej Balog. Alphaevolve: A coding agent for scientific and algorithmic discovery, 2025.
- [35] Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E. Gonzalez, M Waleed Kadous, and Ion Stoica. Routellm: Learning to route llms with preference data, 2025.
- [36] OpenAI. GPT-5 system card. Technical report, OpenAI, August 2025.
- [37] OpenAI. o4-mini model. <https://developers.openai.com/api/docs/models/o4-mini>, 2026. OpenAI API model page, accessed March 16, 2026.
- [38] Qwen Team. Qwen3.5: Towards native multimodal agents, February 2026.
- [39] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
- [40] Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog, M. Pawan Kumar, Emilien Dupont, Francisco J. R. Ruiz, Jordan S. Ellenberg, Pengming Wang, Omar Fawzi, Pushmeet Kohli, and Alhussein Fawzi. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.

- [41] Amrith Setlur, Nived Rajaraman, Sergey Levine, and Aviral Kumar. Scaling test-time compute without verification or rl is suboptimal, 2025.
- [42] Asankhaya Sharma. Openevolve: an open-source evolutionary coding agent, 2025.
- [43] Harman Singh, Xiuyu Li, Kusha Sareen, Monishwaran Maheswaran, Sijun Tan, Xiaoxia Wu, Junxiong Wang, Alpay Ariyak, Qingyang Wu, Samir Khaki, Rishabh Tiwari, Long Lian, Yucheng Lu, Boyi Li, Alane Suhr, Ben Athiwaratkun, and Kurt Keutzer. v_1 : Unifying generation and self-verification for parallel reasoners, 2026.
- [44] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters, 2024.
- [45] Kimi Team, Tongtong Bai, Yifan Bai, Yiping Bao, S. H. Cai, Yuan Cao, Y. Charles, H. S. Che, Cheng Chen, Guanduo Chen, and Huarong Chen et al. Kimi k2.5: Visual agentic intelligence, 2026.
- [46] Qwen Team. Qwen3 technical report, 2025.
- [47] The Terminal-Bench Team. Terminal-bench: A benchmark for ai agents in terminal environments, Apr 2025.
- [48] Together AI. gpt-oss-120b api. <https://www.together.ai/models/gpt-oss-120b>, 2026. Together AI model page, accessed March 16, 2026.
- [49] Together AI. Qwen3 235b a22b instruct 2507 fp8 api. <https://www.together.ai/models/qwen3-235b-a22b-instruct-2507-fp8>, 2026. Together AI model page, accessed March 16, 2026.
- [50] Antonios Valkanas, Soumyasundar Pal, Pavel Rumiantsev, Yingxue Zhang, and Mark Coates. C3po: Optimized large language model cascades with probabilistic cost constraints for reasoning, 2025.
- [51] Siddarth Venkatraman, Vineet Jain, Sarthak Mittal, Vedant Shah, Johan Obando-Ceron, Yoshua Bengio, Brian R. Bartoldson, Bhavya Kailkhura, Guillaume Lajoie, Glen Berseth, Nikolay Malkin, and Moksh Jain. Recursive self-aggregation unlocks deep thinking in large language models, 2026.
- [52] Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities, 2024.
- [53] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models, 2023.
- [54] Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. Large language models are better reasoners with self-verification, 2023.
- [55] Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. Inference scaling laws: An empirical analysis of compute-optimal inference for problem-solving with language models, 2025.
- [56] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023.
- [57] Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, Yu Su, Wenhui Chen, and Graham Neubig. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark, 2025.
- [58] Mert Yuksekgonul, Daniel Kocaja, Xinhao Li, Federico Bianchi, Jed McCaleb, Xiaolong Wang, Jan Kautz, Yejin Choi, James Zou, Carlos Guestrin, and Yu Sun. Learning to discover at test time, 2026.
- [59] Di Zhang, Xiaoshui Huang, Dongzhan Zhou, Yuqiang Li, and Wanli Ouyang. Accessing gpt-4 level mathematical olympiad solutions via monte carlo tree self-refine with llama-3 8b, 2024.
- [60] Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction, 2025.
- [61] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. Sglang: Efficient execution of structured language model programs, 2024.
- [62] Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. Language agent tree search unifies reasoning acting and planning in language models, 2024.

Appendix

A Related Work

Test-time scaling. Test-time scaling invests additional inference compute to improve output quality [44, 55], through parallel sampling [53, 7], sequential refinement [30, 33], search [56, 59, 62], or extended reasoning chains [20, 17]. Compute-optimal sampling with weaker models can outperform a single strong model [6], though scaling without verification remains suboptimal [41]. All of these operate within a single-model regime; SQUEEZE EVOLVE extends test-time scaling to multi-model orchestration by routing evolutionary operations across models of different cost.

Self-aggregation and recursive refinement. Several methods combine multiple LLM outputs into a refined answer, including RSA [51], generative self-aggregation [26], Parallel-Distill-Refine [31], and Best-of-N refinement [22]. Mixture-of-Agents [52] layers multiple LLMs but uses a fixed model assignment rather than adaptive routing. V_1 [43] demonstrates that RSA suffers from diversity collapse (monotonically declining $\text{pass}@N$) and proposes pairwise self-verification as an orthogonal remedy. SQUEEZE EVOLVE addresses the same bottleneck from a complementary angle: multi-model orchestration preserves diverse reasoning lineages, while confidence-based routing delegates easy aggregation groups to cheaper models.

Verification and confidence signals. External verification spans outcome reward models [11], process reward models [27], and generative verifiers [60], while self-verification can improve reasoning [54], though its benefits are situation-dependent [29]. DeepConf [13] uses token-level confidence to filter low-quality reasoning traces, achieving large token savings. SQUEEZE EVOLVE uses the same class of model-intrinsic confidence signals not to filter or verify candidates, but as a routing signal that assigns each recombination group to a model, requiring no trained reward model or external evaluator.

LLM-driven evolutionary search. LLMs serve as evolutionary operators for discovering programs, prompts, and algorithms [25, 40, 34, 58], with subsequent systems varying primarily in selection and variation strategies [42, 24, 2, 4, 8]. EvoX [28] meta-evolves the search strategy itself rather than fixing it. These systems rely on external verifiers and apply a single model uniformly across all operators. SQUEEZE EVOLVE operates in the verifier-free regime and introduces adaptive model assignment: the evolutionary template remains unchanged, but each recombination group is routed to a model commensurate with its difficulty.

Model routing and cost-efficient inference. Cascading and routing frameworks route entire queries between a strong and a weak model [35, 50]. Arbitrage [32] moves to finer granularity by routing individual reasoning steps between draft and target models, achieving $2\times$ latency reduction. SQUEEZE EVOLVE routes at a similarly fine granularity but within a multi-step evolutionary pipeline: individual recombination groups are assigned to models based on per-group confidence, and because these decisions compound across loops, savings accumulate multiplicatively.

B Datasets and Benchmarks

Table 6 summarizes the benchmarks used in this work. We describe each below.

Table 6. Summary of evaluation benchmarks.

Benchmark	Size	Answer Format	Metric
AIME 2025	30	Integer (000–999)	Accuracy
HMMT Feb. 2025	30	Short answer	Accuracy
GPQA-Diamond	198	4-way MC	Accuracy
LiveCodeBench V6	175	Code	Pass@1
MMMU-Pro	1,730	Up to 10-way MC	Accuracy
BabyVision	388	Short answer	Accuracy
ARC-AGI-2	120	Output grid	Pass@2
Circle Packing	1	Program	Objective

AIME 2025 [5]. The American Invitational Mathematics Examination consists of 30 problems (15 from AIME I, 15 from AIME II) covering algebra, geometry, number theory, and combinatorics. Each answer is an integer in $[0, 999]$, scored by exact match. We source problems via MathArena.

HMMT February 2025 [5]. The Harvard–MIT Mathematics Tournament February competition comprises 30 individual-round problems (10 Algebra, 10 Geometry, 10 Combinatorics). Answers are short numerical or symbolic expressions, scored by exact match. We source problems via MathArena.

GPQA-Diamond [39]. A 198-question subset of GPQA filtered for maximum difficulty: both domain experts answered correctly while most non-experts failed even with unrestricted web access. Questions span graduate-level biology, physics, and chemistry in 4-way multiple-choice format.

LiveCodeBench V6 [21]. A competitive programming benchmark sourcing problems from LeetCode, AtCoder, and Codeforces. Models generate code solutions evaluated against hidden test cases; we report pass@1. Continuous collection mitigates data contamination.

MMMU-Pro [57]. A harder variant of MMMU spanning 1,730 multimodal questions across 30 subjects in six disciplines (Art & Design, Business, Science, Health & Medicine, Humanities & Social Science, Tech & Engineering). Answer choices are augmented from 4 to up to 10 options, and text-only solvable questions are filtered out.

BabyVision [9]. A visual reasoning benchmark of 388 items across 22 subclasses in four categories: fine-grained discrimination, visual tracking, spatial perception, and visual pattern recognition. It tests core visual abilities independent of linguistic knowledge; human adults achieve 94.1% while leading MLLMs score below 50%. BabyVision uses an LLM-as-Judge (GPT-4o) for evaluation.

ARC-AGI-2 [10]. A benchmark of 120 public evaluation tasks testing abstract reasoning and compositional generalization. Each task provides demonstration input–output grid pairs; the model must infer the transformation rule and produce the correct output grid. Scored by pass@2 across test pairs (exact grid match with two attempts).

Circle Packing ($n=26$) [34]. An open-ended optimization problem: pack 26 non-overlapping circles in a unit square to maximize the sum of their radii. This is a single continuous-objective instance used to evaluate evolutionary discovery capabilities. The metric is the objective value (sum of radii).

C Generation Hyperparameters

Table 7 lists the generation hyperparameters used for each model. These are model-provided hyperparameters and differ from RSA except for GPT-OSS.

Table 7. Generation hyperparameters for each model.

Model	Effort	Temp.	Top-K	Top-P	Min-P	Gen. Len.
Qwen3-4B-Instruct-2507	—	0.7	20	0.8	0	8K
Qwen3-30B-A3B-Instruct-2507	—	0.7	20	0.8	0	8K
Qwen3-235B-A22B-Instruct-2507	—	0.7	20	0.8	0	16K
Qwen3-235B-A22B-Thinking-2507	—	0.6	20	0.95	0	32K
Qwen3-30B-A3B-Thinking-2507	—	0.6	20	0.95	0	32K
Qwen3-4B-Thinking-2507	—	0.6	20	0.95	0	32K
Qwen3.5-35B-A3B	—	1.0	20	0.95	0	64K
GPT-OSS-20B	medium	1	-1	1	0	16K
GPT-OSS-120B	medium	1	-1	1	0	16K
GPT-5 Mini	medium		default			32K
Gemini-3-Flash-Preview	high		default			64K
Gemini-3.1-Pro-Preview	high		default			64K
Kimi-2.5-Thinking	—	1.0	20	0.95	0	64K
Kimi-2.5-Instant	—	1.0	20	0.95	0	64K

D Empirical Cost Model

We report per-token API pricing from commercial inference providers to ground the routing savings of SQUEEZE EVOLVE in real-world dollar costs. Table 8 lists the models used in our experiments together with their input and output prices from Alibaba Cloud, Together AI, Google, and OpenAI.

Table 8. Per-token API pricing (\$/1M tokens) for each model used in our experiments.

Provider	Model	Input price	Output price
Alibaba Cloud	Qwen3-30B-A3B-Instruct-2507	\$0.108	\$0.431
	Qwen3-30B-A3B-Thinking-2507	\$0.108	\$1.076
	Qwen3-235B-A22B-Instruct-2507	\$0.287	\$0.920
	Qwen3.5-35B-A3B	\$0.057	\$0.459
Google Gemini API	Gemini-3-Flash-Preview	\$0.50	\$3.00
	Gemini-3.1-Pro-Preview	\$2.00	\$12.00
OpenAI	GPT-5 Mini	\$0.25	\$2.00
Together AI	GPT-OSS-20B	\$0.05	\$0.20
	GPT-OSS-120B	\$0.15	\$0.60
	Kimi-K2.5	\$0.50	\$2.80

E Algorithm

Algorithm 1 SQUEEZE EVOLVE

Require: Query set $\{Q_q\}$, Model 2: M_2 , Model 1 M_1 , fitness f , operators Select, Route, LiteAgg, Update, population size N , group size K , groups per problem M , loops T

Ensure: Final populations $\{\mathcal{P}_q^{(T)}\}$

Loop 0 — Initialization (Model 2 only):

- 1: **for** each problem q **do**
- 2: $\mathcal{P}_q^{(0)} \leftarrow \{\tau_i \sim p_{M_2}(\cdot \mid Q_q)\}_{i=1}^N$
- 3: **end for**

Loops 1...T — Fitness-routed evolution:

- 4: **for** $t = 1, \dots, T$ **do**
- 5: **for** each problem q **do**
- 6: Score every $\tau \in \mathcal{P}_q^{(t-1)}$: compute $f(\tau)$
- 7: $\mathcal{G}_q \leftarrow \text{Select}(\mathcal{P}_q^{(t-1)}, K, M, f)$
- 8: $F(g) \leftarrow \text{GroupFitness}(g, f)$ for each $g \in \mathcal{G}_q$
- 9: $(\mathcal{B}_1, \mathcal{B}_2, \mathcal{B}_{\text{lite}}) \leftarrow \text{Route}(\mathcal{G}_q, F)$
- 10: **end for**
- 11: $\mathcal{R}_1 \leftarrow \text{Agg}(M_1, \mathcal{B}_1) \parallel \mathcal{R}_2 \leftarrow \text{Agg}(M_2, \mathcal{B}_2) \parallel \mathcal{R}_{\text{lite}} \leftarrow \text{LiteAgg}(\mathcal{B}_{\text{lite}})$
- 12: $\mathcal{P}_q^{(t)} \leftarrow \text{Update}(\mathcal{P}_q^{(t-1)}, \mathcal{R}_1 \cup \mathcal{R}_2 \cup \mathcal{R}_{\text{lite}})$
- 13: **end for**

F Full Aggregation Accuracy Results

Accuracy rises monotonically with the number of correct seeds. The large model maintains a consistent advantage at intermediate seed counts (1–3), while both models converge at the extremes (0 and 4).

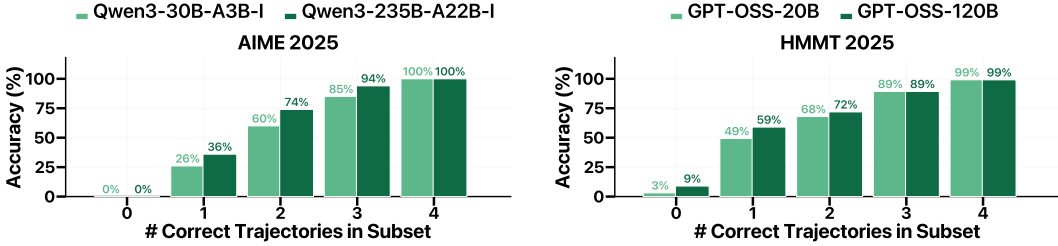


Figure 10. **Full aggregation accuracy vs. number of correct trajectories in subset.** *Left:* AIME 2025 (Qwen3-30B-A3B-Instruct vs. Qwen3-235B-A22B-Instruct). *Right:* HMMT 2025 (GPT-OSS-20B vs. GPT-OSS-120B).

G Full Group Confidence Results

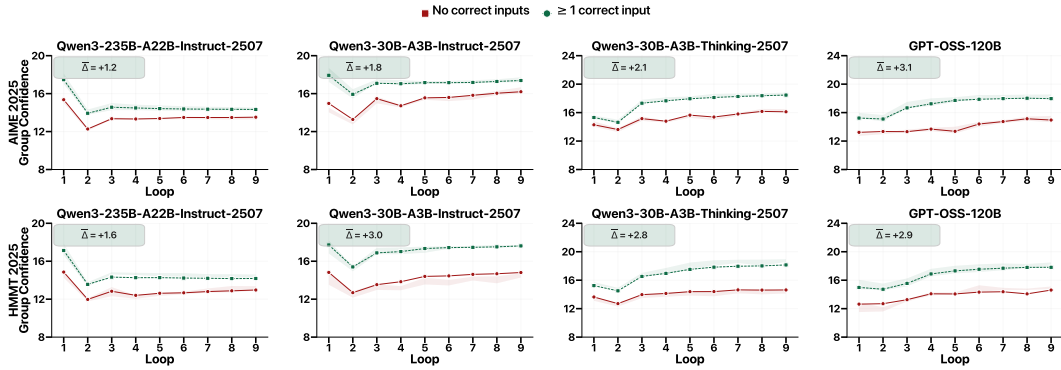


Figure 11. **Self-model group confidence by correctness across all baseline models.** Mean GC (with 40th–60th percentile band) across RSA loops 1–9 for four scorer models: Qwen3-235B-A22B-Instruct, Qwen3-30B-A3B-Instruct, Qwen3-30B-A3B-Thinking, and GPT-OSS-120B. Top row: AIME 2025; bottom row: HMMT 2025. Across all models and benchmarks, subsets containing correct trajectories maintain consistently higher GC than all-incorrect subsets.

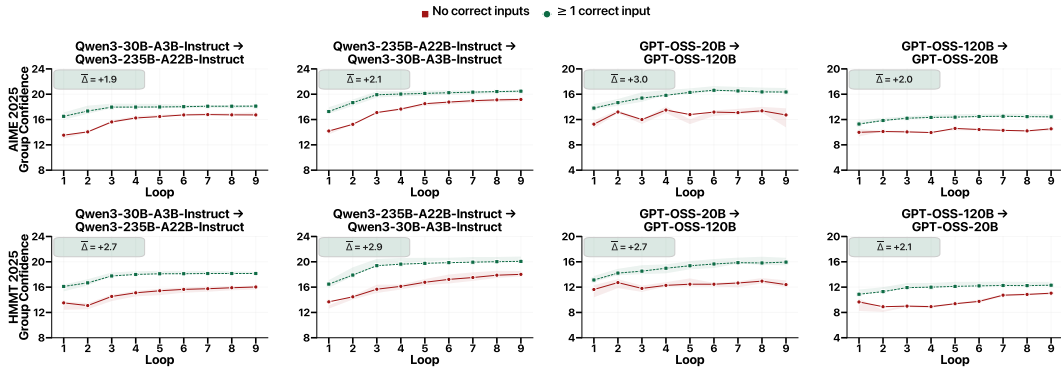
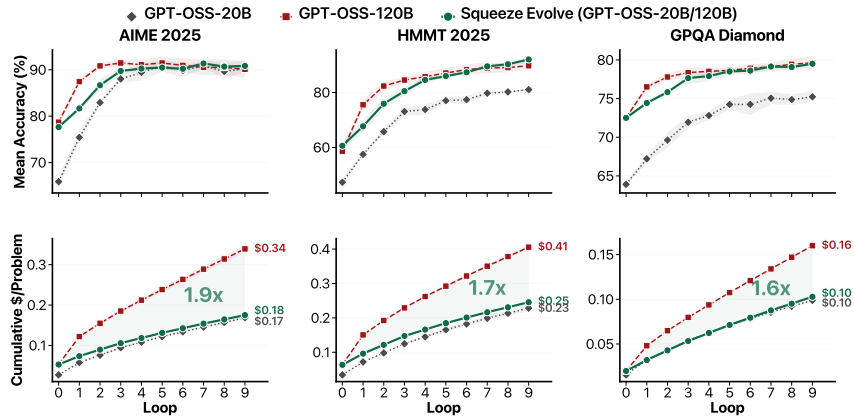


Figure 12. **Cross-model group confidence by correctness across all routing configurations.** Mean GC (with 40th–60th percentile band) for groups containing no correct inputs vs. groups with ≥ 1 correct input, pooled across seeds. Columns show four routing pairs: forward routing (Qwen3-30B-A3B-Instruct \rightarrow Qwen3-235B-A22B-Instruct and GPT-OSS-20B \rightarrow GPT-OSS-120B) and reverse routing (Qwen3-235B-A22B-Instruct \rightarrow Qwen3-30B-A3B-Instruct and GPT-OSS-120B \rightarrow GPT-OSS-20B). Top row: AIME 2025; bottom row: HMMT 2025. GC reliably separates correct from incorrect groups across all routing directions, model families, and benchmarks.

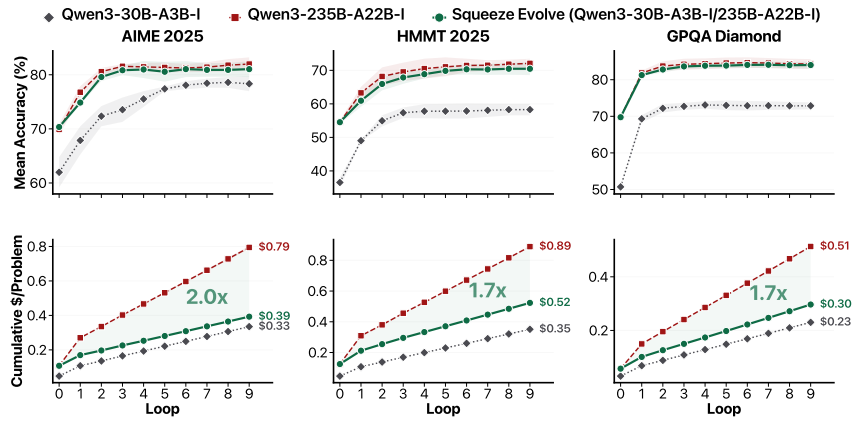
H Empirical Cost Results: Homogeneous Model Pairs for Reasoning Tasks

Table 9. Empirical (dollar) cost results across datasets and model configurations. \$/Prob is the average API cost per problem.

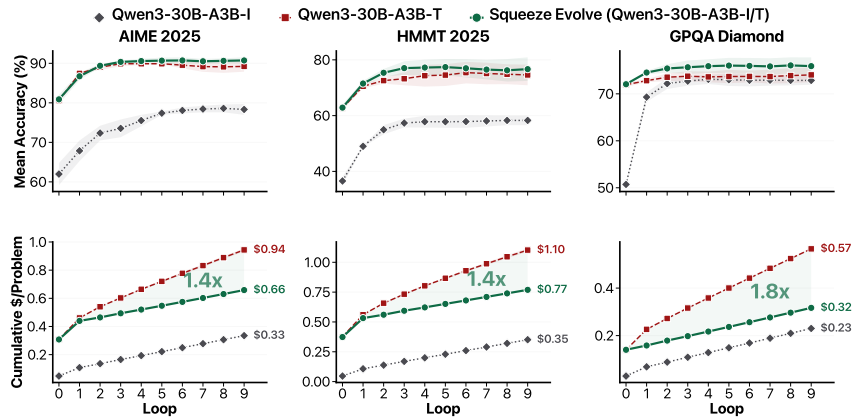
Data	Strategy	Model 1	Model 2	Acc.	\$/Prob	\$ Savings
AIME25	RSA	Qwen3-30B-A3B-I	—	77.8	\$0.33	—
	RSA	—	Qwen3-30B-A3B-T	89.2	\$0.94	1.0×
	SQUEEZE EVOLVE ($p=0$)	Qwen3-30B-A3B-I	Qwen3-30B-A3B-T	90.7	\$0.66	1.4×
	RSA	Qwen3-30B-A3B-I	—	77.8	\$0.33	—
	RSA	—	Qwen3-235B-A22B-I	82.0	\$0.79	1.0×
	SQUEEZE EVOLVE ($p=10$)	Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	80.1	\$0.47	1.7×
	SQUEEZE EVOLVE ($p=0$)	Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	81.0	\$0.39	2.0×
	RSA	GPT-OSS-20B	—	90.0	\$0.17	—
	RSA	—	GPT-OSS-120B	90.1	\$0.34	1.0×
	SQUEEZE EVOLVE ($p=10$)	GPT-OSS-20B	GPT-OSS-120B	90.5	\$0.21	1.6×
	SQUEEZE EVOLVE ($p=0$)	GPT-OSS-20B	GPT-OSS-120B	90.8	\$0.18	1.9×
	HMMT25	RSA	Qwen3-30B-A3B-I	—	57.7	\$0.35
RSA		—	Qwen3-30B-A3B-T	74.6	\$1.10	1.0×
SQUEEZE EVOLVE ($p=0$)		Qwen3-30B-A3B-I	Qwen3-30B-A3B-T	76.7	\$0.77	1.4×
RSA		Qwen3-30B-A3B-I	—	57.7	\$0.35	—
RSA		—	Qwen3-235B-A22B-I	72.1	\$0.89	1.0×
SQUEEZE EVOLVE ($p=10$)		Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	71.4	\$0.52	1.7×
SQUEEZE EVOLVE ($p=0$)		Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	67.4	\$0.44	2.0×
RSA		GPT-OSS-20B	—	80.8	\$0.23	—
RSA		—	GPT-OSS-120B	89.7	\$0.41	1.0×
SQUEEZE EVOLVE ($p=10$)		GPT-OSS-20B	GPT-OSS-120B	92.0	\$0.25	1.6×
SQUEEZE EVOLVE ($p=0$)		GPT-OSS-20B	GPT-OSS-120B	87.9	\$0.22	1.8×
GPQA-Diamond		RSA	Qwen3-30B-A3B-I	—	72.5	\$0.23
	RSA	—	Qwen3-30B-A3B-T	74.0	\$0.57	1.0×
	SQUEEZE EVOLVE ($p=0$)	Qwen3-30B-A3B-I	Qwen3-30B-A3B-T	75.9	\$0.32	1.8×
	RSA	Qwen3-30B-A3B-I	—	72.5	\$0.23	—
	RSA	—	Qwen3-235B-A22B-I	84.3	\$0.51	1.0×
	SQUEEZE EVOLVE ($p=10$)	Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	84.0	\$0.30	1.7×
	SQUEEZE EVOLVE ($p=0$)	Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	83.8	\$0.25	2.1×
	RSA	GPT-OSS-20B	—	75.0	\$0.10	—
	RSA	—	GPT-OSS-120B	79.6	\$0.16	1.0×
	SQUEEZE EVOLVE ($p=10$)	GPT-OSS-20B	GPT-OSS-120B	79.5	\$0.10	1.6×
	SQUEEZE EVOLVE ($p=0$)	GPT-OSS-20B	GPT-OSS-120B	79.0	\$0.09	1.9×
	LCB-V6	RSA	Qwen3-30B-A3B-I	—	46.1	\$0.19
RSA		—	Qwen3-30B-A3B-T	64.2	\$0.82	1.0×
SQUEEZE EVOLVE ($p=10$)		Qwen3-30B-A3B-I	Qwen3-30B-A3B-T	62.7	\$0.63	1.3×
SQUEEZE EVOLVE ($p=0$)		Qwen3-30B-A3B-I	Qwen3-30B-A3B-T	59.1	\$0.51	1.6×
RSA		Qwen3-30B-A3B-I	—	46.1	\$0.19	—
RSA		—	Qwen3-235B-A22B-I	59.1	\$0.33	1.0×
SQUEEZE EVOLVE ($p=10$)		Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	55.9	\$0.22	1.5×
SQUEEZE EVOLVE ($p=0$)		Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	55.3	\$0.19	1.7×
RSA		GPT-OSS-20B	—	68.9	\$0.14	—
RSA		—	GPT-OSS-120B	75.9	\$0.44	1.0×
SQUEEZE EVOLVE ($p=10$)		GPT-OSS-20B	GPT-OSS-120B	75.6	\$0.22	2.0×
SQUEEZE EVOLVE ($p=0$)		GPT-OSS-20B	GPT-OSS-120B	73.3	\$0.18	2.4×



(a)



(b)



(c)

Figure 13. **Empirical cost results for homogeneous model pairs.** *Top row* of each panel: mean accuracy (%) across RSA loops. *Bottom row*: cumulative API cost per problem (\$). SQUEEZE EVOLVE (green) matches or exceeds the Model 2-only baseline (red) in accuracy while substantially reducing cost. The shaded region highlights the cost savings, which grow with each loop as more candidates are routed to the cheaper Model 1.

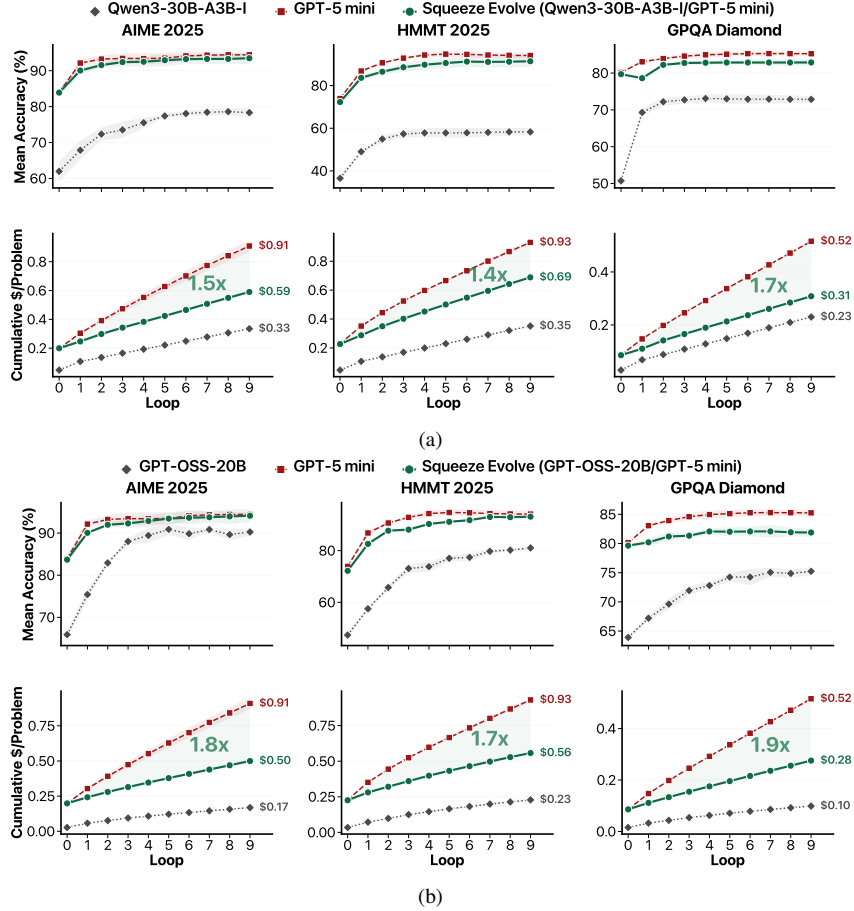


Figure 14. **Empirical cost results for heterogeneous model pairs.** *Top row* of each panel: mean accuracy (%) across RSA loops. *Bottom row*: cumulative API cost per problem (\$). SQUEEZE EVOLVE (green) matches or exceeds the Model 2 baseline (red) in accuracy while substantially reducing cost. The shaded region highlights the cost savings, which grow with each loop as more candidates are routed to the cheaper Model 1.

I Empirical Cost Results: Heterogeneous Model Pairs for Reasoning Tasks

Table 10. Empirical (dollar) cost results across datasets and model configurations. \$/Prob is the average API cost per problem.

Data	Strategy	Model 1	Model 2	Acc.	\$/Prob	\$ Savings
AIME25	RSA	Qwen3-30B-A3B-I	—	78.8	\$0.34	—
	RSA	—	GPT-5 mini	94.2	\$0.89	1.0×
	SQUEEZE EVOLVE ($p=30$)	Qwen3-30B-A3B-I	GPT-5 mini	93.5	\$0.64	1.4×
	SQUEEZE EVOLVE ($p=20$)	Qwen3-30B-A3B-I	GPT-5 mini	93.5	\$0.59	1.5×
	SQUEEZE EVOLVE ($p=10$)	Qwen3-30B-A3B-I	GPT-5 mini	93.1	\$0.53	1.7×
	SQUEEZE EVOLVE ($p=0$)	Qwen3-30B-A3B-I	GPT-5 mini	91.9	\$0.46	2.0×
	RSA	GPT-OSS-20B	—	90.6	\$0.17	—
	RSA	—	GPT-5 mini	94.2	\$0.89	1.0×
	SQUEEZE EVOLVE ($p=30$)	GPT-OSS-20B	GPT-5 mini	95.4	\$0.50	1.8×
	SQUEEZE EVOLVE ($p=20$)	GPT-OSS-20B	GPT-5 mini	94.6	\$0.46	1.9×
SQUEEZE EVOLVE ($p=10$)	GPT-OSS-20B	GPT-5 mini	92.8	\$0.39	2.3×	
SQUEEZE EVOLVE ($p=0$)	GPT-OSS-20B	GPT-5 mini	92.7	\$0.30	3.0×	
RSA	Qwen3-30B-A3B-I	—	58.9	\$0.35	—	
RSA	—	GPT-5 mini	93.3	\$0.94	1.0×	

Continued on next page

Data	Strategy	Model 1	Model 2	Acc.	\$/Prob	\$ Savings
GQA-Diamond	SQUEEZE EVOLVE ($p=30$)	Qwen3-30B-A3B-I	GPT-5 mini	93.1	\$0.69	1.4×
	SQUEEZE EVOLVE ($p=20$)	Qwen3-30B-A3B-I	GPT-5 mini	90.2	\$0.66	1.4×
	SQUEEZE EVOLVE ($p=10$)	Qwen3-30B-A3B-I	GPT-5 mini	88.1	\$0.59	1.6×
	SQUEEZE EVOLVE ($p=0$)	Qwen3-30B-A3B-I	GPT-5 mini	87.6	\$0.51	1.9×
	RSA	GPT-OSS-20B	—	81.2	\$0.22	—
	RSA	—	GPT-5 mini	93.3	\$0.94	1.0×
	SQUEEZE EVOLVE ($p=30$)	GPT-OSS-20B	GPT-5 mini	93.1	\$0.56	1.7×
	SQUEEZE EVOLVE ($p=20$)	GPT-OSS-20B	GPT-5 mini	91.8	\$0.51	1.8×
	SQUEEZE EVOLVE ($p=10$)	GPT-OSS-20B	GPT-5 mini	89.8	\$0.43	2.2×
	SQUEEZE EVOLVE ($p=0$)	GPT-OSS-20B	GPT-5 mini	89.3	\$0.35	2.7×
	RSA	Qwen3-30B-A3B-I	—	73.3	\$0.23	—
	RSA	—	GPT-5 mini	85.0	\$0.52	1.0×
	SQUEEZE EVOLVE ($p=30$)	Qwen3-30B-A3B-I	GPT-5 mini	82.6	\$0.37	1.4×
	SQUEEZE EVOLVE ($p=20$)	Qwen3-30B-A3B-I	GPT-5 mini	83.6	\$0.35	1.5×
	SQUEEZE EVOLVE ($p=10$)	Qwen3-30B-A3B-I	GPT-5 mini	83.2	\$0.31	1.7×
	SQUEEZE EVOLVE ($p=0$)	Qwen3-30B-A3B-I	GPT-5 mini	82.2	\$0.26	2.0×
	RSA	GPT-OSS-20B	—	75.5	\$0.10	—
	RSA	—	GPT-5 mini	85.0	\$0.52	1.0×
	SQUEEZE EVOLVE ($p=30$)	GPT-OSS-20B	GPT-5 mini	82.1	\$0.27	1.9×
	SQUEEZE EVOLVE ($p=20$)	GPT-OSS-20B	GPT-5 mini	81.8	\$0.25	2.1×
SQUEEZE EVOLVE ($p=10$)	GPT-OSS-20B	GPT-5 mini	80.5	\$0.20	2.5×	
SQUEEZE EVOLVE ($p=0$)	GPT-OSS-20B	GPT-5 mini	78.8	\$0.16	3.3×	

J Empirical Cost Results: Homogeneous Model Pairs for Vision Tasks

Table 11. Empirical (dollar) cost results for vision tasks with homogeneous model pairs. \$/Prob is the average API cost per problem.

Data	Strategy	Model 1	Model 2	Acc.	\$/Prob	\$ Savings
MMMU-Pro	RSA	Kimi-2.5-Instant	—	77.46	\$0.29	—
	RSA	—	Kimi-2.5-Thinking	78.58	\$1.04	1.0×
	SQUEEZE EVOLVE ($p=0$)	Kimi-2.5-Instant	Kimi-2.5-Thinking	78.63	\$0.58	1.8×
BabyVision	RSA	Kimi-2.5-Instant	—	36.61	\$0.29	—
	RSA	—	Kimi-2.5-Thinking	43.23	\$2.05	1.0×
	SQUEEZE EVOLVE ($p=0$)	Kimi-2.5-Instant	Kimi-2.5-Thinking	41.56	\$0.81	2.5×

K Empirical Cost Results: Heterogeneous Model Pairs for Vision Tasks

Table 12. Empirical (dollar) cost results for vision tasks with heterogeneous model pairs. \$/Prob is the average API cost per problem. [†]Image tokens are only processed by the Model 2 in loop 0; subsequent loops use Model 1 as a text-only causal LM (no image input).

Data	Strategy	Model 1	Model 2	Acc.	\$/Prob	\$ Savings
MMMU-Pro	RSA	Qwen3.5-35B-A3B [†]	—	—	—	—
	RSA	—	Kimi-2.5-Thinking	78.58	\$1.04	1.0×
	SQUEEZE EVOLVE ($p=0$)	Qwen3.5-35B-A3B [†]	Kimi-2.5-Thinking	79.06	\$0.46	2.3×
BabyVision	RSA	Qwen3.5-35B-A3B [†]	—	—	—	—
	RSA	—	Kimi-2.5-Thinking	43.23	\$2.05	1.0×
	SQUEEZE EVOLVE ($p=0$)	Qwen3.5-35B-A3B [†]	Kimi-2.5-Thinking	41.27	\$0.83	2.5×

L Prefill Engine Microbenchmark

As described in Section 5.2, the SQUEEZE EVOLVE custom prefill engine computes the confidence scalar directly on GPU and returns only a single float per request, bypassing the native vLLM path that materializes full token-level logprob tensors and serializes them over HTTP. Table 13 reports per-request confidence-scoring latency for both paths across two models (GPT-OSS-120B and Qwen3-30B-A3B-Thinking-2507) and three context lengths (40K, 80K, and 120K tokens), measured at batch size 1 with 3 trials per configuration. The custom engine achieves 9.1–10.3× speedup on GPT-OSS-120B and 4.2–6.6× on Qwen3-30B-A3B-Thinking. The larger speedup on the 120B model reflects the heavier serialization and transfer burden at larger model scales: the native path transfers ~13 MB of token-level logprob data per request regardless of model size, while the custom engine returns only ~100 bytes. At the largest scale (Qwen3-235B-A22B), the native path OOMs entirely when materializing full prompt-logprob tensors, making the custom engine a prerequisite for confidence-based routing at this model size. Figure 15 shows that latency scales approximately linearly with context length under both paths, but the custom engine’s slope is substantially gentler. This is because the native path’s cost is dominated by tensor materialization and HTTP transfer, both of which grow with sequence length, whereas the custom engine performs an in-place reduction on GPU and transmits only the scalar result.

Table 13. **Prefill engine microbenchmark.** Confidence-scoring latency (seconds, mean over 3 trials) for native vLLM and vLLM SQUEEZE EVOLVE, measured at batch size 1 across three context lengths. vLLM SQUEEZE EVOLVE computes the confidence scalar on GPU and returns only a single float per request, reducing transfer volume from ~13 MB to ~100 bytes. Native vLLM OOMs on Qwen3-235B-A22B because materializing full prompt-logprob tensors at this model scale exceeds memory; vLLM SQUEEZE EVOLVE does not exhibit this issue.

Model	Context	vLLM native (s)	vLLM SQUEEZE EVOLVE (s)	Speedup
GPT-OSS-120B	40K	8.60	0.83	10.3×
	80K	17.68	1.79	9.9×
	120K	26.86	2.94	9.1×
Qwen3-30B-A3B	40K	10.34	1.58	6.6×
	80K	22.79	4.42	5.2×
	120K	35.42	8.41	4.2×

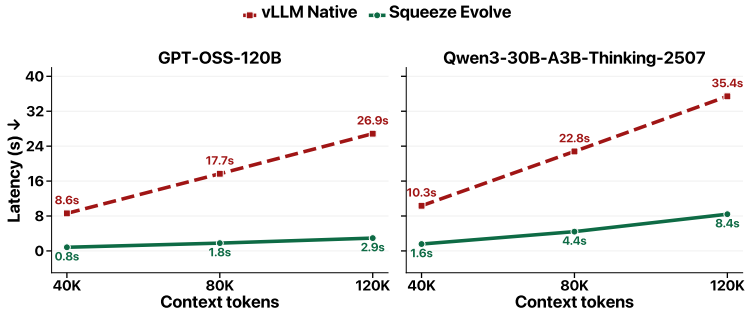


Figure 15. **Prefill engine speedup over native vLLM.** Confidence-scoring latency as a function of context length for GPT-OSS-120B (left) and Qwen3-30B-A3B-Thinking-2507 (right). The SQUEEZE EVOLVE custom prefill engine (green) computes the confidence scalar on GPU and returns only the final score, achieving 4–10× lower latency than the native vLLM prompt-logprob path (red), which materializes full token-level tensors and transfers them over HTTP. Speedup is larger for the 120B model where serialization and transfer costs dominate. Both paths scale approximately linearly with context length, but the custom engine has a substantially gentler slope.

M Routing Overhead Results

Experimental setup. For each benchmark and model M_2 , we fix the full inference configuration: evolution parameters ($N = 16, K = 4, T = 10$), prompts, decoding limits, hardware, serving engine, and batching policy. We compare two conditions. **RSA- M_2** is standard RSA with all calls executed by M_2 and no routing logic. **SQUEEZE EVOLVE- M_2** enables confidence scoring and threshold computation, but forces every aggregation call to M_2 . This second condition preserves the routing machinery while removing any latency change due to sending work to M_1 . The difference between RSA- M_2 and SQUEEZE EVOLVE- M_2 therefore isolates the routing overhead itself. This setup is also a conservative worst-case measurement: although SQUEEZE EVOLVE normally reduces latency by routing a subset of aggregations to M_1 , here all aggregation work is still pinned to M_2 .

Measurement protocol. We measure latency in a single-request setting, processing one problem at a time so that routing overhead is not confounded by cross-request queueing effects. Within each problem, however, we preserve the production execution strategy: the N confidence-scoring calls and aggregation calls for a loop are batched exactly as in normal serving. For each problem we log end-to-end latency, per-loop routing time, and per-loop aggregation time. We repeat the measurement across the evaluation set and report mean latency.

Overhead definition. Let T_{RSA} denote the end-to-end latency of RSA- M_2 and let T_{SQE} denote the latency of SQUEEZE EVOLVE- M_2 . We define the absolute routing overhead as

$$\Delta_{\text{route}} = T_{\text{SQE}} - T_{\text{RSA}}.$$

We define the relative routing overhead as

$$\text{Overhead}(\%) = 100 \times \frac{T_{\text{SQE}} - T_{\text{RSA}}}{T_{\text{RSA}}}.$$

At the loop level, we decompose latency as

$$T_{\text{loop}} = T_{\text{routing}} + \sum_{i=1}^N T_{\text{agg}}(m_i),$$

where T_{routing} includes both prefill-only confidence scoring and percentile thresholding / dispatch, and $T_{\text{agg}}(m_i)$ is the aggregation time for the selected model $m_i \in \{M_1, M_2\}$. In this section, all aggregations are forced to M_2 , so the observed gap isolates the latency overhead of routing logic alone.

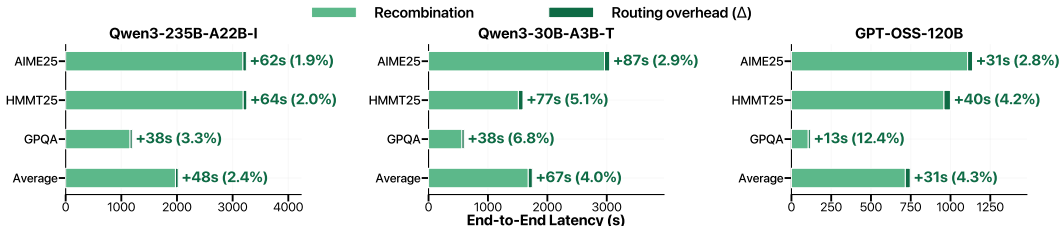


Figure 16. **Routing overhead is minimal.** Routing adds 1.9–6.8% to end-to-end latency for the Qwen models and 2.8–12.4% for GPT-OSS-120B, with the higher relative overhead on GPQA reflecting its short absolute generation time (106s). Full results in Table 14 (Appendix).

Results. Table 14 reports per-benchmark routing overhead for each model M_2 . Across all configurations, routing adds 1.9–6.8% to end-to-end latency for the Qwen models and 2.8–12.4% for GPT-OSS-120B, with cross-model averages of 2.4–4.3%. The higher relative overhead on GPQA-Diamond for GPT-OSS-120B (12.4%) reflects its short absolute generation time (106s), which makes the fixed routing cost proportionally larger.

N System Throughput Results

Fairness rule. We compare RSA and SQUEEZE EVOLVE under the same total GPU budget G . RSA allocates all G GPUs to Model 2. SQUEEZE EVOLVE partitions the same budget into a Model 2 pool G_2 and a Model 1 pool G_1 , subject to

$$G_2 + G_1 = G.$$

This fixed-budget constraint makes the throughput comparison deployment-fair.

Table 14. Routing-overhead measurements in the single-request setting. Comparing RSA- M_2 and SQUEEZE EVOLVE- M_2 isolates the latency overhead of confidence scoring and dispatch. Highlighted average rows summarize per-target-model averages across benchmarks.

M_L	Benchmark	RSA (s)	SQUEEZE EVOLVE (s)	Δ_{route} (s)	Overhead (%)
Qwen3-30B-A3B-T	AIME25	2961.72	3048.44	86.72	2.93%
	HMMT25	1512.39	1589.61	77.22	5.11%
	GPQA-Diamond	561.45	599.43	37.98	6.76%
	LCB-V6	1798.91	1864.46	65.55	3.64%
	Average	1678.52	1745.83	67.31	4.01%
Qwen3-235B-A22B-I	AIME25	3184.11	3246.08	61.97	1.95%
	HMMT25	3190.07	3253.75	63.68	2.00%
	GPQA-Diamond	1157.17	1195.57	38.40	3.32%
	LCB-V6	359.22	385.93	26.71	7.44%
	Average	1972.64	2020.33	47.69	2.42%
GPT-OSS-120B	AIME25	1107.44	1138.35	30.91	2.79%
	HMMT25	958.92	999.24	40.32	4.20%
	GPQA-Diamond	105.74	118.87	13.13	12.42%
	LCB-V6	691.80	729.70	37.90	5.47%
	Average	715.98	746.54	30.57	4.27%

Operating points. We sweep the routing percentile p from Section 5 across several values. Because the realized Model 1 routing share does not exactly equal $p/100$; we report both.

Pool sizing. Given a configured percentile p and its observed routing mix, we size the two pools so that their loop service times are approximately matched. Let $T_2(G_2)$ denote the time for the Model 2 pool to process the groups routed to Model 2, and let $T_1(G_1)$ denote the corresponding time for the Model 1 pool. We choose integer G_2 and G_1 satisfying $G_2 + G_1 = G$ and minimizing

$$|T_2(G_2) - T_1(G_1)|.$$

This latency-matching rule avoids provisioning a fast idle pool while the slower pool remains the throughput bottleneck.

Measurement protocol. After fixing the GPU split, we drive the system with enough concurrent requests to keep serving saturated and measure steady-state throughput. We report requests per second rather than tokens per second because Model 1 and Model 2 produce different numbers of output tokens for the same query, making a token-based metric an unfair comparison across routing configurations. Let N_{req} denote the total number of requests completed by the system over wall-clock interval Δt . We define throughput as

$$\text{Req/s} = \frac{N_{\text{req}}}{\Delta t}.$$

We use the same query stream, prompts, and serving engine for both RSA and SQUEEZE EVOLVE, and report completed requests per second after warmup. Because confidence-based routing causes Model 1 and Model 2 to observe different input and output lengths, we fix the input context length and output context length for each model to the values observed under the corresponding routing share.

Results. Table 15 reports steady-state throughput under a fixed total GPU budget G for two model pairs across four benchmarks. For each routing percentile p , the GPU budget is partitioned into large-model and small-model pools with approximately matched loop service times. The Accuracy column reports mean accuracy from Table 9.

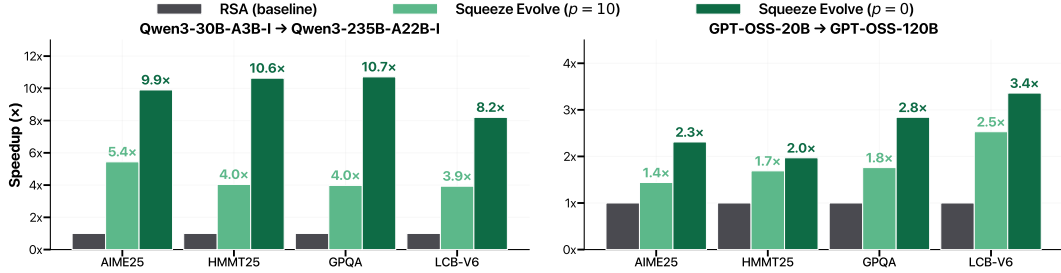


Figure 17. **Fixed-budget throughput speedup over RSA.** Under the same total GPU budget, the Qwen pair achieves 4–10 \times speedup and the GPT-OSS pair 1.4–3.4 \times . Full results in Table 15 (Appendix).

Table 15. Fixed-budget system throughput under saturated serving. RSA and SQUEEZE EVOLVE use the same total GPU budget G . For each routing percentile p , SQUEEZE EVOLVE partitions the budget into Model 2 and Model 1 pools with approximately matched loop service times under the observed routing mix. Throughput is reported as steady-state completed requests per second after warmup. Accuracy is mean accuracy (%) from Table 9.

Model 1	Model 2	Benchmark	Strategy	Obs. M_1 share	GPU split ($G_2:G_1$)	Req/s	Speedup	Acc. (%)
Qwen3-30B-A3B-I	Qwen3-235B-A22B-I	AIME25	RSA	0%	16:0	1.36	1.00 \times	82.0
			SQUEEZE EVOLVE ($p=10$)	88.9%	8:8	7.41	5.44 \times	80.1
			SQUEEZE EVOLVE ($p=0$)	100%	0:16	13.47	9.90 \times	81.0
		HMMT25	RSA	0%	16:0	1.23	1.00 \times	72.1
			SQUEEZE EVOLVE ($p=10$)	87.5%	8:8	4.95	4.04 \times	71.4
			SQUEEZE EVOLVE ($p=0$)	100%	0:16	13.02	10.63 \times	67.4
		GPQA Diamond	RSA	0%	16:0	2.05	1.00 \times	84.3
			SQUEEZE EVOLVE ($p=10$)	87.5%	8:8	8.17	3.98 \times	84.0
			SQUEEZE EVOLVE ($p=0$)	100%	0:16	22.00	10.71 \times	83.8
		LCB-V6	RSA	0%	16:0	3.83	1.00 \times	59.1
			SQUEEZE EVOLVE ($p=10$)	87.5%	8:8	15.07	3.93 \times	55.9
			SQUEEZE EVOLVE ($p=0$)	100%	0:16	31.42	8.20 \times	55.3
GPT-OSS-20B	GPT-OSS-120B	AIME25	RSA	0%	20:0	17.09	1.00 \times	90.1
			SQUEEZE EVOLVE ($p=10$)	87.5%	4:16	24.59	1.44 \times	90.5
			SQUEEZE EVOLVE ($p=0$)	100%	0:20	39.43	2.31 \times	90.8
		HMMT25	RSA	0%	12:0	8.56	1.00 \times	89.7
			SQUEEZE EVOLVE ($p=10$)	87.5%	4:8	14.50	1.69 \times	92.0
			SQUEEZE EVOLVE ($p=0$)	100%	0:12	16.83	1.97 \times	87.9
		GPQA Diamond	RSA	0%	16:0	30.34	1.00 \times	79.6
			SQUEEZE EVOLVE ($p=10$)	87.5%	4:12	53.54	1.76 \times	79.5
			SQUEEZE EVOLVE ($p=0$)	100%	0:16	86.30	2.84 \times	79.0
		LCB-V6	RSA	0%	12:0	5.66	1.00 \times	75.9
			SQUEEZE EVOLVE ($p=10$)	87.1%	4:8	14.30	2.53 \times	75.6
			SQUEEZE EVOLVE ($p=0$)	100%	0:12	19.02	3.36 \times	73.3

O ARC-AGI-2 Complete Results

Table 16. ARC-AGI-2 full results. Ancestor model: Gemini 3.1 Pro (High Thinking), $N=4$, $K=2$ $T=10$.
[†]Uses code execution and program synthesis.

Strategy	Model 1	Model 2	Acc.	\$/Task	Savings
<i>Human baseline</i>					
Human panel	—	—	100.0	\$17.00	—
<i>Single-shot baselines</i>					
GPT-5.4 Pro (xhigh)	—	—	92.2	\$17.60	—
Gemini 3.1 Pro (High)	—	—	88.1	\$0.98	—
GPT-5.4 (xhigh)	—	—	84.2	\$1.57	—
Claude Opus 4.6 (Thinking 120K, high)	—	—	79.0	\$3.81	—
GPT-5.4 (high)	—	—	75.8	\$1.08	—
<i>Code-execution methods</i>					
Imbue + Gemini 3.1 Pro [†]	—	—	95.1	\$8.71	—
Confluence Lab [†]	—	—	97.9	\$11.77	—
RSA	—	Gemini 3.0 Flash	45.0	\$9.83	—
RSA	—	Gemini 3.1 Pro	93.3	\$28.85	1.0×
SQUEEZE EVOLVE	—	Gemini 3.1 Pro	97.5	\$7.74	3.7×

P Circle Packing Complete Results

P.1 Algorithm summary

The evolved algorithm (Section P.3) combines three strategies: (1) a diverse initialization ensemble that generates hundreds of candidate center layouts via hexagonal lattices, greedy farthest-point insertion, jittered grids, and random placements, scoring each with an exact linear programmed (LP) that maximizes $\sum r_i$ for fixed centres; (2) a hybrid optimization pipeline integrating LP-guided simulated annealing with SLSQP gradient-based refinement, where each stochastic perturbation of 1–3 centers is immediately followed by an LP solve to obtain provably optimal radii, and an adaptive temperature schedule prevents premature convergence before a final SLSQP pass jointly optimizes all $3N=78$ variables under wall-distance and non-overlap constraints; and (3) a principled decomposition that separates the hard combinatorial center placement (\mathbb{R}^{52}) from the easy convex radius assignment (an LP).

P.2 Hyperparameters

We instantiate SQUEEZE EVOLVE with GPT-OSS-120B and GPT-OSS-20B as M2 and M1 models, use group confidence as the fitness signal, fitness-weighted sampling ($\zeta=0.5$) for selection, a fixed confidence threshold at the 50th percentile for routing, and the *accumulate* update rule (Table 3), with $N=128$, $K=4$, $T=50$. At termination, we draw N candidates from the cumulative pool via confidence-weighted sampling and report the highest circle packing score.

P.3 Source code

```
1 #!/usr/bin/env python3
2 """
3 Optimised packing of 26 non-overlapping circles inside the unit square.
4
5 The algorithm combines:
6 * Several diverse initialisation strategies (hexagonal lattice,
7   farthest point, random grid, pure random).
8 * A cheap but exact linear programme (LP) that, for any fixed set of
9   centre positions, yields the maximal radii (maximising the sum of radii).
10 * Stochastic hill-climbing with a temperature schedule to refine the
11   centre positions.
12 * A final local optimisation with SciPy's SLSQP optimiser, which moves
13   centres and radii simultaneously while respecting all constraints.
14 * A tiny post-processing step that enforces strict feasibility.
15
16 The program prints exactly 26 lines "x y r" (nine decimal digits each)
17 to stdout.
18 """
19
20 import sys
21 import numpy as np
22
23 try:
24     from scipy.optimize import linprog, minimize
25     _SCIPY = True
26 except Exception:
27     _SCIPY = False
28
29 N = 26
30 PAIR_I, PAIR_J = np.triu_indices(N, 1)
31 M_PAIRS = len(PAIR_I)
32
33 # Pre-computed constraint matrix for the LP:
34 # one row per pair (r_i + r_j <= d_ij)
35 A_UB = np.zeros((M_PAIRS, N), dtype=float)
36 A_UB[np.arange(M_PAIRS), PAIR_I] = 1.0
37 A_UB[np.arange(M_PAIRS), PAIR_J] = 1.0
38
39 def _fallback_radii(centres: np.ndarray) -> np.ndarray:
40     """Simple pairwise scaling - used only if the LP fails."""
41     n = centres.shape[0]
42     r = np.minimum.reduce([centres[:, 0],
43                             centres[:, 1],
44                             1.0 - centres[:, 0],
45                             1.0 - centres[:, 1]])
46     for i in range(n):
47         for j in range(i + 1, n):
48             d = np.linalg.norm(centres[i] - centres[j])
49             if r[i] + r[j] > d:
50                 scale = d / (r[i] + r[j])
51                 r[i] *= scale
```

```

52         r[j] *= scale
53     return np.maximum(r, 0.0)
54
55
56 def _lp_optimal(centres: np.ndarray) -> tuple[np.ndarray, float]:
57     """Solve the LP that maximises sum(r_i) for the fixed centres."""
58     wall = np.minimum.reduce([centres[:, 0],
59                               centres[:, 1],
60                               1.0 - centres[:, 0],
61                               1.0 - centres[:, 1]])
62     bounds = [(0.0, w) for w in wall]
63
64     diffs = centres[PAIR_I] - centres[PAIR_J]
65     pair_dist = np.linalg.norm(diffs, axis=1)
66
67     if _SCIPY:
68         res = linprog(c=-np.ones(N),
69                     A_ub=A_UB,
70                     b_ub=pair_dist,
71                     bounds=bounds,
72                     method='highs',
73                     options={'presolve': True})
74         if res.success:
75             radii = np.clip(res.x, 0.0, wall)
76             return radii, float(radii.sum())
77     radii = _fallback_radii(centres)
78     return radii, float(radii.sum())
79
80 def _hex_lattice(dx: float) -> np.ndarray:
81     """Generate points on a hexagonal lattice (up to 26 points)."""
82     dy = dx * np.sqrt(3.0) / 2.0
83     pts = []
84     row = 0
85     y = 0.0
86     while y <= 1.0 + 1e-12:
87         offset = 0.0 if (row % 2 == 0) else dx / 2.0
88         x = offset
89         while x <= 1.0 + 1e-12:
90             pts.append([x, y])
91             x += dx
92         y += dy
93         row += 1
94     pts = np.asarray(pts)
95     pts = pts[(pts[:, 0] >= 0.0) & (pts[:, 0] <= 1.0) &
96              (pts[:, 1] >= 0.0) & (pts[:, 1] <= 1.0)]
97     if pts.shape[0] > N:
98         margins = np.minimum.reduce([pts[:, 0],
99                                     pts[:, 1],
100                                    1.0 - pts[:, 0],
101                                    1.0 - pts[:, 1]])
102         keep = np.argsort(-margins)[:N]
103         pts = pts[keep]
104     elif pts.shape[0] < N:
105         rng = np.random.default_rng(0)
106         extra = rng.uniform(0.0, 1.0, size=(N - pts.shape[0], 2))
107         pts = np.vstack([pts, extra])
108     return pts
109
110
111 def _init_farthest(rng: np.random.Generator, n: int = N) -> np.ndarray:
112     """Greedy farthest-point placement (points stay inside [0,1]^2)."""
113     pts = []
114     pts.append(rng.uniform(0.1, 0.9, size=2))
115     while len(pts) < n:
116         cand = rng.uniform(0.1, 0.9, size=(200, 2))
117         existing = np.asarray(pts)
118         dists = np.linalg.norm(
119             cand[:, None, :] - existing[None, :, :], axis=2)
120         min_d = dists.min(axis=1)
121         best = np.argmax(min_d)
122         pts.append(cand[best])
123     return np.asarray(pts)
124
125
126 def _init_grid_jitter(rng: np.random.Generator,
127                      jitter: float = 0.02) -> np.ndarray:
128     """5x5 grid (0.1-0.9) with small random jitter + one extra point."""
129     xs = np.linspace(0.1, 0.9, 5)
130     ys = np.linspace(0.1, 0.9, 5)
131     xv, yv = np.meshgrid(xs, ys)
132     base = np.column_stack([xv.ravel(), yv.ravel()]) # 25 points

```

```

133 base += rng.uniform(-jitter, jitter, base.shape)
134 base = np.clip(base, 0.01, 0.99)
135 extra = rng.uniform(0.01, 0.99, size=(1, 2))
136 return np.vstack([base, extra])
137
138
139 def _init_random(rng: np.random.Generator) -> np.ndarray:
140     """Pure uniform random points."""
141     return rng.uniform(0.0, 1.0, size=(N, 2))
142
143 def _hill_climb(start: np.ndarray, rng: np.random.Generator,
144               iterations: int = 1500
145               ) -> tuple[np.ndarray, np.ndarray, float]:
146     best_c = start.copy()
147     best_r, best_sum = _lp_optimal(best_c)
148
149     temperature = 0.02
150     for it in range(iterations):
151         cand_c = best_c.copy()
152         k = rng.integers(1, 4)
153         sel = rng.choice(N, size=k, replace=False)
154         max_step = 0.09 * (1.0 - it / iterations) + 0.005
155         cand_c[sel] += rng.normal(scale=max_step, size=(k, 2))
156         cand_c = np.clip(cand_c, 0.0, 1.0)
157
158         cand_r, cand_sum = _lp_optimal(cand_c)
159         delta = cand_sum - best_sum
160
161         if delta > 1e-9:
162             best_c, best_r, best_sum = cand_c, cand_r, cand_sum
163             temperature = max(temperature * 0.95, 1e-6)
164         else:
165             if (temperature > 0.0
166                 and rng.random() < np.exp(delta / temperature)):
167                 best_c, best_r, best_sum = cand_c, cand_r, cand_sum
168                 temperature *= 0.9995
169     return best_c, best_r, best_sum
170
171 def _refine_slsqp(centres: np.ndarray, radii_start: np.ndarray,
172                 rng: np.random.Generator
173                 ) -> tuple[np.ndarray, np.ndarray, float]:
174     if not _SCIPY:
175         return centres, radii_start, radii_start.sum()
176
177     n = N
178     x0 = np.empty(3 * n)
179     x0[:n] = centres[:, 0]
180     x0[n:2 * n] = centres[:, 1]
181     x0[2 * n:] = radii_start
182
183     bounds = (([0.0, 1.0]) * n
184               + [[0.0, 1.0]] * n
185               + [[0.0, None]] * n)
186
187     def cons_fun(x):
188         xs = x[:n]
189         ys = x[n:2 * n]
190         rs = x[2 * n:]
191         c1 = xs - rs
192         c2 = ys - rs
193         c3 = (1.0 - xs) - rs
194         c4 = (1.0 - ys) - rs
195         dx = xs[:, None] - xs[None, :]
196         dy = ys[:, None] - ys[None, :]
197         d = np.sqrt(dx * dx + dy * dy)
198         iu = np.triu_indices(n, 1)
199         c5 = d[iu] - (rs[iu[0]] + rs[iu[1]])
200         return np.concatenate([c1, c2, c3, c4, c5])
201
202     constraints = {'type': 'ineq', 'fun': cons_fun}
203
204     def obj_fun(x):
205         return -np.sum(x[2 * n:])
206
207     res = minimize(obj_fun,
208                  x0,
209                  method='SLSQP',
210                  bounds=bounds,
211                  constraints=constraints,
212                  options={'ftol': 1e-9, 'maxiter': 2000,
213                          'disp': False})

```

```

214
215     if not res.success:
216         final_c = centres
217     else:
218         xs_opt = res.x[:n]
219         ys_opt = res.x[n:2 * n]
220         final_c = np.vstack([xs_opt, ys_opt]).T
221
222     final_r, final_sum = _lp_optimal(final_c)
223     return final_c, final_r, final_sum
224
225 def _make_feasible(centres: np.ndarray, radii: np.ndarray,
226                   eps: float = 1e-12) -> np.ndarray:
227     wall = np.minimum.reduce([centres[:, 0],
228                              centres[:, 1],
229                              1.0 - centres[:, 0],
230                              1.0 - centres[:, 1]])
231     radii = np.minimum(radii, wall - eps)
232     for _ in range(5):
233         changed = False
234         for i in range(N):
235             for j in range(i + 1, N):
236                 d = np.linalg.norm(centres[i] - centres[j])
237                 if radii[i] + radii[j] > d - eps:
238                     scale = (d - eps) / (radii[i] + radii[j])
239                     radii[i] *= scale
240                     radii[j] *= scale
241                     changed = True
242         if not changed:
243             break
244     return np.maximum(radii, 0.0)
245
246 def construct_packing() -> tuple[np.ndarray, np.ndarray, float]:
247     rng = np.random.default_rng(123456789)
248     best_sum = -np.inf
249     best_centres = None
250     best_radii = None
251
252     # -- Stage 1: diverse starts --
253     candidates = []
254     for dx in np.linspace(0.16, 0.30, 15):
255         base = _hex_lattice(dx)
256         for jitter_scale in (0.0, 0.008, 0.02, 0.04):
257             for _ in range(5):
258                 jitter = rng.normal(
259                     scale=jitter_scale, size=base.shape)
260                 centres = np.clip(base + jitter, 0.0, 1.0)
261                 radii, s = _lp_optimal(centres)
262                 candidates.append((s, centres, radii))
263
264     for _ in range(5):
265         centres = _init_farthest(rng)
266         radii, s = _lp_optimal(centres)
267         candidates.append((s, centres, radii))
268
269     for jit in (0.015, 0.03):
270         for _ in range(5):
271             centres = _init_grid_jitter(rng, jitter=jit)
272             radii, s = _lp_optimal(centres)
273             candidates.append((s, centres, radii))
274
275     for _ in range(5):
276         centres = _init_random(rng)
277         radii, s = _lp_optimal(centres)
278         candidates.append((s, centres, radii))
279
280     candidates.sort(key=lambda x: x[0], reverse=True)
281     top_candidates = candidates[:6]
282
283     # -- Stage 2: hill climbing --
284     for s0, cent0, rad0 in top_candidates:
285         cent_opt, rad_opt, sum_opt = _hill_climb(
286             cent0, rng, iterations=1800)
287         if sum_opt > best_sum:
288             best_sum = sum_opt
289             best_centres = cent_opt
290             best_radii = rad_opt
291
292     # -- Stage 3: SLSQP refinement --
293     if best_centres is not None:
294         refined_c, refined_r, refined_sum = _refine_slsqp(

```

```

295     best_centres, best_radii, rng)
296     if refined_sum > best_sum:
297         best_sum = refined_sum
298         best_centres = refined_c
299         best_radii = refined_r
300
301     final_radii = _make_feasible(best_centres, best_radii)
302     final_sum = float(final_radii.sum())
303     return best_centres, final_radii, final_sum
304
305 def run_packing() -> None:
306     centres, radii, _ = construct_packing()
307     out = sys.stdout
308     fmt = "{:.9f} {:.9f} {:.9f}\n"
309     for i in range(N):
310         out.write(fmt.format(
311             centres[i, 0], centres[i, 1], radii[i]))
312
313
314 if __name__ == "__main__":
315     run_packing()

```

Listing 1. Evolved circle packing program ($n=26$).

P.4 Correlation between confidence and score across loops

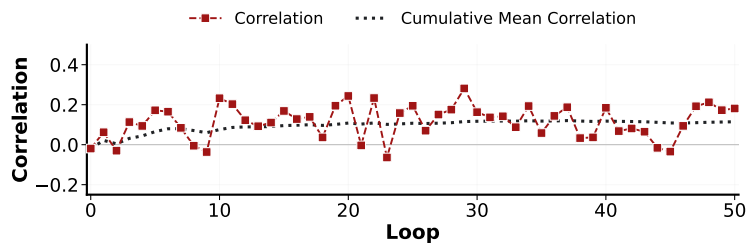


Figure 18. Spearman rank correlation between confidence and score