AN ARCHITECTURE SEARCH FRAMEWORK FOR INFERENCE-TIME TECHNIQUES

Jon Saad-Falcon[†]*, Adrian Gamarra Lafuente[†], Shlok Natarajan[†], Nahum Maru[†], Hristo Todorov[†], Etash Guha[‡], E. Kelly Buchanan[†], Mayee Chen[†], Neel Guha[†], Christopher Ré[†], Azalia Mirhoseini[†]

[†] Stanford University

[‡] University of Washington

Abstract

Inference-time techniques, such as repeated sampling or iterative revisions, are emerging as powerful ways to enhance large-language models (LLMs) at test time. However, best practices for developing systems that combine these techniques remain underdeveloped due to our limited understanding of the utility of each technique across models and tasks, the interactions between them, and the massive search space for combining them. To address these challenges, we introduce ARCHON, a modular and automated framework for optimizing the process of selecting and combining inference-time techniques and LLMs. Given a compute budget and a set of available LLMs, ARCHON explores a large design space to discover optimized configurations tailored to target benchmarks. It can design custom or general-purpose architectures that advance the Pareto frontier of accuracy vs. maximum token budget compared to top-performing baselines. Across instruction-following, reasoning, and coding tasks, we show that ARCHON can leverage additional inference compute budget to design systems that outperform frontier models such as OpenAI's o1, GPT-40, and Claude 3.5 Sonnet by an average of 15.1%. ARCHON is open-source and plug-and-play, allowing users to select from existing inference-time techniques (or add new ones) while optimizing for their desired objective functions: task performance, cost, latency, and more.

1 INTRODUCTION

Inference-time techniques—strategies that use additional compute during model inference—are gaining traction as effective methods for improving model capabilities. LLMs, such as OpenAI's o1 (OpenAI, 2024a), QwQ (Team, 2024), and Sky-T1 (Team, 2025), utilize such techniques to translate additional inference compute into better performance across a broad set of tasks. Example techniques include generation ensembling, ranking, and fusion, where models in the ensemble are queried in parallel, their responses are ranked, and the best ones are fused into a single, higher quality output, respectively (Jiang et al., 2023b; Wang et al., 2024). Other types of inference-time techniques are based on querying a single LLM successively (via repeated sampling) and using a voting strategy or unit tests to select the top generation (Brown et al., 2024; Chen et al., 2024; Li et al., 2024a).

Recent work has made progress towards building robust *inference-time architectures*: systems composed of one or more large language models (LLMs) leveraging inference-time techniques. Examples include Mixture-of-Agents (MoA) (Wang et al., 2024) and LLM-Blender (Jiang et al., 2023b), as well as single-model systems like ADAS (Hu et al., 2024) and AFlow (Zhang et al., 2024). However, our experiments show that these top-performing baselines have limitations in compute utilization and task generalization. (see Section 3.2). We argue that designing effective and generalizable inference-time architectures requires the following:

• Understanding the Utilities of Inference-Time Techniques: Inference-time architectures typically delegate their additional inference budget towards more model sampling calls (Chen et al., 2024;

^{*}Corresponding author: <jonsaadfalcon@stanford.edu>



Figure 1: **Overview of ARCHON Framework**: ARCHON's search algorithm requires the following inputs: target benchmarks, inference call budget, available LLMs, and available inference-time techniques (**left**). The search algorithm uses Bayesian optimization (Snoek et al., 2012) to construct and evaluate different ARCHON configurations (**middle**) before returning the optimized ARCHON architecture (**right**) for the target benchmarks (Section 2.3).

Brown et al., 2024), which can be effective for math and coding tasks. Other tasks, such as following instructions and reasoning, have been shown to benefit from additional techniques, including ranking and fusion (Wang et al., 2024; Jiang et al., 2023b). While all of these methods are valuable, *it is essential to identify which inference-time techniques are most effective for different task categories.*

- Understanding the Interactions Between Inference-Time Techniques: While previous studies analyzed these techniques individually (e.g., generation sampling in Chen et al. (2024)), we need a more comprehensive understanding of the relationships between different inference-time techniques across different tasks (e.g., is it better to use more models or generate more samples per model?).
- Efficiently and Automatically Searching the Large Design Space of Inference-Time Architectures: Given a set of available LLMs and target tasks, there is currently no single prevailing inference-time architecture for maximizing downstream accuracy across all tasks (Table 1). The search space for inference-time architectures is expansive, requiring practitioners to make several key configuration decisions, such as *which LLMs to use, how many times to sample them, and how to combine and filter the candidate generations*. These motivate the need for automated and adaptive architecture search approaches.

In our work, we address each of these challenges. First, we **evaluate the utilities of a comprehensive set of existing and proposed inference-time techniques** across instruction-following, reasoning, and coding tasks. Using both open-source and closed-source models, we examine a range of techniques such as *ensembling*, *fusion*, *ranking*, *critiquing*, *verification*, *and model-based unit test generation/evaluation* (Sections 2.1 and 2.2). We find that no single technique completely dominates across all tasks, with different approaches being more effective for different tasks.

Second, we **analyze the interactions between inference-time techniques** and explore the benefits of adding new models and new techniques individually. We find that generation ensembling combined with critique, verification, and fusion improves the final response quality beyond the oracle best candidate from individual (non-fused) responses, particularly for instruction-following and reasoning tasks (Figure 2.2; Figure 7; Table 18). We also demonstrate increased performance as we scale up the layers of inference-time techniques and combine multiple approaches together, allowing us to discover effective new combinations of inference-time techniques (Sections 2.2, 3.2, A.4). Combining multiple strategies significantly improves task performance, but determining the specific combination remains challenging. This requires manually testing models, inference-time techniques, architecture designs, inference budgets, and more.

Third, drawing upon our analysis of inference-time techniques, we present **ARCHON**, an open-source modular framework for automatically designing LLM systems composed of existing inference-time techniques (or new ones), allowing practitioners to optimize for their desired objective functions: accuracy, latency, and cost (Sections 2.1, 2.3). Unlike alternative LM systems that perform prompt engineering and tool use over a single LM (Khattab et al., 2023; Yuksekgonul et al., 2024; Hu et al., 2024; Zhang et al., 2024), our approach integrates multiple LMs in a single architecture and reduces prompt selection to a set of core components. The ARCHON framework utilizes automatic

architecture search algorithms to maximize generation quality for the given tasks(s), leveraging Bayesian optimization (Snoek et al., 2012; Nardi et al., 2019) techniques inspired by (NAS) (Zoph & Le, 2017; Ren et al., 2021) to rapidly traverse the space of potential inference architectures (Section 2.3).

We evaluate ARCHON architectures across a diverse set of instruction-following, reasoning, and coding benchmarks (Table 1): MT-Bench, Arena-Hard-Auto, Alpaca-2.0 Eval, MixEval, MATH, and CodeContests (Zheng et al., 2023; Li et al., 2024b; 2023; Ni et al., 2024; Hendrycks et al., 2021; Li et al., 2022). Our best ARCHON architectures surpass both frontier models (e.g. OpenAI's O1, GPT-40 and Claude-3.5 Sonnet) and prior top-performing inference-time architectures (e.g. ADAS, AFlow, and MoA), *boosting state-of-the-art (SOTA) performance by 15.1%*, on average. Furthermore, ARCHON achieves SOTA performances while using 20.0% less inference calls, 15.1% less input tokens, and 13.5% less output tokens than alternative inference-time architectures (Table 1; Figure 4). Even when solely using open-source LLMs, ARCHON architectures, on average, surpass SOTA LLMs by 11.2%.

Overall, we present ARCHON as an open-source inference-time framework, readily extensible to new inference-time techniques, models, and tasks via user-friendly interfaces.

2 INFERENCE-TIME TECHNIQUES FOR ARCHON

With the proliferation of inference-time techniques, ARCHON introduces a systematic framework for understanding and unifying these methods into inference-time architectures. Below, we elaborate on the structure, inputs, and outputs of each of the inference-time techniques (Table 8). Then, we discuss how to combine the different techniques into an inference-time architecture (Section 2.2) before finally exploring automatic approaches for constructing inference-time architectures (Section 2.3).

2.1 LLM COMPONENTS OF ARCHON

In this section, we discuss the *LLM components* of ARCHON, which are LLMs that perform a specific inference-time technique. We test an array of different components inspired by recent work, incorporating approaches for generating, ranking, and fusing candidates (Wang et al., 2024; Jiang et al., 2023b) as well as approaches for improving candidate response quality through critiquing, verifying, and unit testing (Bai et al., 2022; Zheng et al., 2023). The components and their prompts are summarized in Table 8 and Section A.3.

Generator is an LLM that takes in the instruction prompt and outputs candidate responses. Generators can be called in parallel to perform *generation ensembling* (i.e. calling multiple LLMs in parallel) (Wang et al., 2024), or sampled multiple times (Brown et al., 2024). The number of models, samples, and generation temperature can be adjusted.

Fuser is an LLM that, given an instruction prompt and a set of proposed responses as input, combines these responses to generate one or more higher-quality fused responses.

Ranker is an LLM that, given an instruction prompt and a set of proposed responses as input, ranks the candidate generations based on their quality, producing a ranked list of responses as output. This ranking is then used to filter the set of responses to the top-K, as specified.

Critic is an LLM that, given an instruction prompt and a set of proposed responses as input, produces a list of strengths and weaknesses for each response, which is then used to improve the quality of the final response (Section 2.2; Figure 2.2).

Verifier is an LLM that verifies whether a provided candidate response has appropriate reasoning for a given instruction prompt. It proceeds in two stages: **Stage #1** takes in the instruction prompt and a candidate response as input and outputs reasoning for why the candidate response is correct; **Stage #2** takes in the instruction prompt, candidate response, and produced reasoning before outputting reasoning and a verdict (i.e., binary [Correct] or [Incorrect]) for whether or not the candidate response are passed to the next ARCHON layer.

Unit Test Generator is an LLM that takes the instruction prompt as input and outputs a list of unit tests for assessing the accuracy and relevance of candidate responses. Each test is a concise statement that can be passed or failed as determined by a Unit Test Evaluator, allowing such tests to extend to tasks beyond coding. The number of unit tests generated is a configurable choice; we find 5-10 generated unit tests to be most effective with our set of LM prompts (Section 3.2; examples in Table 16).



Figure 2: **Example ARCHON Architecture**: This architecture starts with ten generator models (each sampled once), followed by a critic model, a ranker model, one layer of six fuser models, a verifier model, and finishes with a fuser model.

Unit Test Evaluator is an LLM that takes the instruction prompt, candidate response(s), and generated unit tests as input and outputs the candidate response(s), ranked in descending order by how many unit tests they pass. We use model-based unit test evaluation by prompting the LLM to justify and aggregate unit test verdicts across each of the candidate responses, scoring each candidate response for reasoning and coding tasks (Table 1). Only responses that pass all the unit tests are passed to the next ARCHON layer.

2.2 Combining the LLM Components

Overview: Inspired by the structure of neural networks (Hinton et al., 1992), ARCHON consists of layers of LLM components (Figure 1; Section 2.1). Each layer is composed of sets of LLM components called in parallel. These components perform a text-to-text operation on the initial instruction prompt and the candidate responses from the previous layer. Furthermore, like a neural network, some layers perform *transformations* of the provided list of strings (e.g., Generator and Fuser), converting a list of strings into a different list of strings (the numbers of candidates can vary from the original number of candidates). Other components introduce non-linearities into the ARCHON structure, performing filtering of the list of strings (e.g., Ranker and Verifier). Ultimately, the inputs and outputs for each layer is always a list of strings, whether that is the instruction prompt (i.e., a single string) or a list of candidate responses. If a list of strings is outputted at the last layer of the ARCHON structure, the first string in the list is returned.

Unlike a classical neural network, no weights are learned between the LLM components and the layers; in turn, the ARCHON architecture can be deployed off-the-shelf without any tuning. Additionally, a single state is transformed sequentially from the input layer to the final output; this single state is the initial instruction prompt and the current candidate responses. In Figure 2, we provide an example ARCHON architecture composed of six layers.

Rules for Construction: The LLM components in Section 2.1 can only be placed in specific orders (9):

- 1. Only one type of component is allowed in any given layer.
- 2. Generator components can only be placed in the first layer of ARCHON; you can place one or more.
- 3. The Critic must come before a Ranker or a Fuser. Otherwise, the generated strengths and weaknesses cannot be incorporated into generation ranking or fusion.
- 4. Ranker, Critic, Verifier, and Unit Test Generator/Evaluator layers can go anywhere in ARCHON except the first layer. For each of these components, it must be the only module in its layer.
- 5. Fuser components can go anywhere in ARCHON except the first layer. Multiple Fusers can be used in a layer.
- 6. Unit Test Generators and Evaluators are placed in consecutive layers: Generator then Evaluator.

Performance Gains from Scaling Inference-Time Techniques: We explore the utilities of individual ARCHON components and evaluate whether combinations of inference-time techniques enable us to *build LM systems greater than the sum of their parts*. From our analysis, we find several trends (designated with **T**s) across ARCHON architectures:

• **T1**: Repeated model sampling and additional ensemble models leads to substantial gains, leading to 9.3% and 18.5% increases, respectively (Figure 5; Figure 7).



Figure 3: **Performance Improves by Scaling** *Layers* **of Inference-Time Techniques**: When controlling for inference budget, generation ensembling and fusion across 8 different 70B LLMs is generally more effective than repeated sampling with only the top performing model. Furthermore, adding layers of critique and fusion led to a 18.8% boost in task performance, on average. However, the best inference-time architecture differed by task, such as MixEval and CodeContests (Section 3.3; Table 4; Table 5), which inspired us to develop architecture search techniques for ARCHON (Section 2.3).

- **T2**: Scaling the layers of inference-time techniques significantly improves ARCHON performance across instruction-following, reasoning, and coding tasks, such as always adding a single fuser as the last layer (Figure 2.2).
- T3: Scaling the diversity of inference-time techniques bolsters performance across the explored tasks, with critics and rankers before fusers being particularly effective (Figure 2.2; Figure 5).
- **T4**: In reasoning tasks, incorporating the Verifier and Unit Test Generator/Evaluator modules alongside the Fuser improves performance by filtering out flawed responses, contributing to significant performance gains in tasks like MixEval and CodeContests (Table 18; Section A.6).

For detailed analysis of interactions between LLM components, please see Section A.4, where we perform a series of ablation experiments in which we vary ARCHON component combinations (Table 18) and the models used in the combinations (Table 19; Table 22).

2.3 ARCHITECTURE SEARCH ALGORITHMS

Search Hyperparameters: In this section, we explore how to automatically design inference-time architectures for target tasks via ARCHON's architecture search algorithms. Guided by the trends found in our analysis in Section 2.2, we establish six axes of hyperparameters for the search space:

- 1. **Top**-*K* **Generators for Ensemble**: The top-*K* models for the Generator ensemble, ranging from 1 to 10 (**T1**). The top-*K* models are selected greedily based on their task performances (Table 24).
- 2. **Top**-*K* **Generator Samples**: The number of samples gathered from each ensemble generator (same for all the models), ranging from 1 to 5 (**T1**). For CodeContests, we explore high-sample settings: [1, 10, 100, 500, 1000].
- 3. Number of Fusion Layers: Ranges from 1 to 4. The last fusion layer always has a single Fuser (T2).
- 4. **Top**-*K* **Fusers**: Number of models used for each fusion layer, ranges from 2 to 10 (**T2**,**3**).
- 5. Critic and Ranker Layers: We add critic and ranker layers before each fuser layer since we find they provide added benefits across the benchmarks explored (T3) (Section 2.2; Figure 2.2; Figure 7).
- 6. Evaluation Layer: Add Verifier, Unit Test Gen./Eval., or neither before final Fuser layer (T4).

While it is possible to further expand the search space of potential ARCHON architectures (e.g., different temperatures for generative LLM components, alternative prompts for each LLM component, additional LLM components for ARCHON, etc.), the trends we identify from Section 2.2 reasonably constrain the search space of configurations to focus on the most influential hyperparameters. In total, our search space contains 9,576 configurations, which we obtain by combining all possible hyperparameters and removing invalid configurations (for example, we discard configurations where the number of initial generations exceeds the context window of the fusers).

							MT Bench	Alpaca Eval 2.0	Arena Hard Auto	MixEval Hard	MixEval	MATH*	Code Contests*
		Approaches	Average Infer. Calls	Average Input Tokens	Average Output Tokens	TFLOPs per Token	W.R.	L.C. W.R.	W.R	Acc.	Acc.	Pass @1	Pass @1
nes	LM	GPT-40 Claude 3.5 Sonnet Llama 3.1 405B	1 1 1	95 105 118	549 602 631	Unk. Unk. 0.81	44.2% ±0.5 N/A 44.1% ±0.3	57.8% ±0.6 52.7% ±0.4 40.7% ±0.5	80.6% ±0.6 81.4% ±0.4 64.5% ±0.7	63.4% ±0.2 68.7% ±0.2 66.0% ±0.3	87.5% ±0.3 89.1% ±0.2 88.2% ±0.2	83.5% ±0.4 82.5% ±0.7 85.0% ±0.5	18.1% ±0.2 12.3% ±0.4 20.4% ±0.5
Baseline	LM Systems	MoA ADAS AFlow ol	19 52 48 Unk.	25,109 72,804 68,596 112	17,422 44,872 41,748 Unk.	1.36 Unk. Unk. Unk.	51.6% ±0.6 66.3% ±0.7 62.4% ±0.2 56.3% ±0.5	65.0% ±0.3 60.1% ±0.5 57.8% ±0.6 59.3% ±0.5	85.3% ±0.3 85.4% ±0.4 83.2% ±0.6 81.7% ±0.3	62.3% ±0.4 64.2% ±0.2 63.5% ±0.3 72.0% ±0.4	86.9% ±0.2 87.0% ±0.2 87.2% ±0.4 87.5% ±0.2	82.9% ±0.6 86.0% ±0.8 84.5% ±0.2 92.7% ±0.5	15.1% ±0.5 23.7% ±0.3 21.1% ±0.6 31.5% ±0.8
_	Open Source	General Purpose Task Specific	35 44	51,113 63,157	31,508 39,949	3.14 3.71	67.2% ±0.4 71.1% ±0.6	63.3% ±0.6 68.1% ±0.4	85.6% ±0.5 89.6% ±0.4	65.3% ±0.3 67.5% ±0.2	86.2% ±0.2 88.8% ±0.3	87.5% ±0.6 89.5% ±0.3	18.2% ±0.4 28.9% ±0.9
Archon	Closed Source	General Purpose Task Specific	32 40	52,747 59,085	27,894 37,271	Unk. Unk.	72.7% ±0.3 77.0% ±0.5	63.9% ±0.7 68.9% ±0.5	86.2% ±0.7 90.5% ±0.3	67.5% ±0.4 72.3% ±0.3	87.2% ±0.2 89.5% ±0.3	87.9% ±0.7 92.1% ±0.4	20.2% ±0.6 25.1% ±0.6
	All Source	General Purpose Task Specific	35 39	50,427 58,250	30,461 36,114	Unk. Unk.	76.2% ±0.7 79.5% ±0.4	66.4% ±0.3 69.0% ±0.6	89.8% ±0.6 92.5% ±0.5	69.8% ±0.2 72.7% ±0.3	87.3% ±0.4 89.7% ±0.2	89.3% ±0.5 93.5% ±0.6	23.4% ±0.9 41.4% ±0.7



Search Method: The ARCHON search method takes in four inputs: the target benchmark(s), the inference call budget, the set of available LLMs, and the inference-time techniques for construction (Figure 1). As output, the search method outputs a single optimized ARCHON architecture. We use 20% of each target dataset as a development set for guiding architecture search. We explore three approaches for ARCHON's architecture search: random search (randomly test potential architectures in the search space), greedy search (greedily optimize individual hyperparameters one at a time, starting from a random initial architecture), and Bayesian Optimization (Snoek et al., 2012) (global hyperparameter optimization with Gaussian processes). As inputs, Bayesian optimization takes in a vector specifying the configuration choices for the generators (i.e., number of models and samples), layers of fusers, numbers of fusers per layer, and final verifier / unit tester (Section 2.2). Bayesian optimization begins by sampling a specified number of random ARCHON architectures to calibrate its surrogate model. The task performance of these sampled architectures is used to guide more informed architecture suggestions during the configuration search. The algorithm repeats the following cycle—evaluating each suggested architecture and using its performance to refine future suggestions-until it discovers the optimal ARCHON configuration, or until the inference call budget is exhausted. For more details on our open-source Bayesian optimization approach, please see Section A.8.

Bayesian optimization found the best architectures in 96.0% of searches and required 88.5% fewer architecture evaluations than greedy search and 90.4% fewer than random search (Figure 13). The effectiveness of Bayesian optimization increases with the number of initial randomly sampled architectures, up to around 230-240 samples, after which further testing is better focused on configuration search (Table 32). For limited inference call budgets (<20 calls), Bayesian optimization is less effective, and traditional methods like greedy search may perform comparably (Table 33).

3 EXPERIMENTS

Our experiments focus on answering the following questions: (1) how does ARCHON compare to existing SOTA LLMs and inference-time architectures in terms of accuracy and compute efficiency (Section 3.2)? (2) how does ARCHON performance compare across the tasks explored (Section 3.3)? (3) what are the considerations for model size, latency, and cost surrounding ARCHON (Section 3.4)? We outline the benchmarks, models, and techniques for constructing ARCHON architectures in Section 3.1.

3.1 BENCHMARKS AND MODELS

Benchmarks: We evaluate our models with several benchmarks for instruction-following, reasoning, and coding: MT-Bench (Zheng et al., 2023), AlpacaEval 2.0 (Li et al., 2023), Arena Hard Auto (Li et al., 2024b), MixEval (Ni et al., 2024), MixEval-Hard, MATH (Hendrycks et al., 2021), and CodeContests (Li et al., 2022). We provide an overview of each dataset in Table 2. Since we perform automatic architecture search on a randomly sampled 20% subset of each benchmark, we evaluate on the remaining



Figure 4: **ARCHON's Performance Exceeds Baselines across Token Budgets**: Across different token budgets (Section 2.3), we compare ARCHON architectures against top-performing inference-time system baselines. The MoA architecture and OpenAI's o1 are static so they use the same number of tokens across budgets. *MATH and CodeContests use a subset of their test sets for evaluation (Section 3.1).

held-out 80% subset of the benchmark (Table 1) (for ARCHON performances on the entire benchmarks, please see Table 7). The delta between the ARCHON performance on the entire benchmark vs. 80% held-out subset is relatively small: only 0.44%, on average, across these datasets with an S.D. of 0.20%. For MATH, we evaluate a random sample of 200 problems from the dataset's test set. For CodeContests, we evaluate on the 140 test set questions that do not include image tags in the problem description.

Models: We test the efficacy of the ARCHON framework by creating different ARCHON architectures across three model categories: 8B or less parameter models, 70B or more parameter models, and closed-source model APIs. For our 8B and 70B+ models, we selected the top-10 performing chat models for each parameter range on the Chatbot Arena Leaderboard (Chiang et al., 2024) as of July 2024. For our ARCHON architectures, we explore multiple model types: open-source, closed-source, and *all-source* (i.e. both open-source and closed-source available). For our closed-source model APIs, we include GPT-40, GPT-4-Turbo, Claude Opus 3.0, Claude Haiku 3.0, and Claude Sonnet 3.5. We list and compare all of the models tested in the ARCHON framework in Table 23 and Table 24. For all the LLMs utilized and every ARCHON component, we set the generation temperature to 0.7. As baselines, we compare ARCHON against both SOTA single-call LLMs (GPT-40, Claude 3.5 Sonnet, and Llama 3.1 405B Instruct) as well as SOTA inference-time approaches (OpenAI's 01 (OpenAI, 2024a), MoA (Wang et al., 2024), ADAS (Hu et al., 2024), and AFlow (Zhang et al., 2024)).

Task-Specific and General-Purpose ARCHON Architectures: We compare custom ARCHON architectures, specifically configured to a single evaluation dataset ("Task-specific ARCHON Architectures"), and a generalized ARCHON architecture configured to handle all the evaluation datasets ("Generalpurpose ARCHON Architectures") (Table 1). For our three model selection settings for ARCHON (i.e. open-source, closed-source, and all-source), we utilize automatic architecture search to find targeted ARCHON architectures for each task (7 architectures total) and find a single generalized ARCHON architecture for maximizing performance over all the tasks (Table 1). The benchmarks are concatenated together and shuffled for generalized ARCHON architecture search. Importantly, all the ARCHON architectures utilized in Section 3 are automatically generated by our Bayesian architecture search technique, which searches over the hyperparameter search space for ARCHON as covered in Section 2.3. For examples of targeted and generalized ARCHON architectures, please see Figure 2 and Appendix A.6.

3.2 ARCHON VS. CLOSED-SOURCE LLMS AND OTHER INFERENCE-TIME ARCHITECTURES

Task Performances: We start by comparing ARCHON architectures to existing SOTA closed-source LLMs and inference-time architectures across a set of instruction-following, reasoning, and coding tasks. Based on our results in Table 1, we find that ARCHON architectures consistently match or surpass existing approaches across all the benchmarks explored. ARCHON architectures with open-source models demonstrate a 11.2% average improvement over SOTA open-source approaches; for its worst performance, our open-source ARCHON architectures are still 3.1% above SOTA open-source approaches on AlpacaEval 2.0. ARCHON architectures with closed-source models achieve SOTA performance across MT Bench, Arena-Hard-Auto, MixEval, and MixEval-Hard, leading to a 15.1% average improvement over closed-source LMs and a 8.4% average improvement over open-source inference-time frameworks (i.e. MoA, ADAS, and AFlow). Compared to o1 and o1-mini, ARCHON's best targeted architectures beat them by 8.1% and 9.7%, on average, on MT Bench, AlpacaEval 2.0, Arena Hard Auto, MixEval,

MixEval Hard, MATH, and CodeContests. For approaches that use all models available, both open and closed-source, ARCHON achieves an average 10.9% improvement over existing SOTA single-call LLMs and an average 8.6% improvement over existing inference-time frameworks.

Compute Efficiency: Compared to open-source inference-time frameworks (i.e. AFlow, ADAS, MoA), ARCHON is 20.0% more inference call efficient while having higher performances on all benchmarks tested (Table 1). We also find that our best ARCHON architectures use 15.1% less input tokens and 13.5% less output tokens compared to the best alternative open-source inference-time frameworks. When we utilize ARCHON's architecture search technique with different token budgets (Figure 4), we find that the generated ARCHON architectures achieve 12.4% higher performance than alternate baselines when given the same budget. Overall, the generalized all-source ARCHON architecture achieves 6.4% better performance across all the tasks while being 31% more token efficient than the best LM system baselines (Table 1). Furthermore, compared to the generalized all-source ARCHON architecture, the targeted all-source ARCHON architectures use 15.5% and 18.6% more input tokens and output tokens, respectively, but they achieve 8.4% higher accuracies, on average. The targeted architectures are more compute intensive since they can further leverage additional LM operations towards a single set of specific task constraints (Appendix A.6).

Discovered Architectures: We include the targeted and generalized ARCHON architectures in Appendix A.6 (Figure 9). The best performing all-source, general-purpose ARCHON architecture starts with a broad initial layer of our 10 best generators before four successive layers of critique and fusion with Qwen2 72B and Claude 3.5 Sonnet, respectively. Each subsequent layer has fewer fuser models (i.e. 8, 6, and 4), leading to a "funneling" effect on the generations before the final output. The best targeted architectures can vary by task. For instruction-following and reasoning tasks, the targeted architectures tend to be multiple layers of critiquing and fusing with a diverse mix of LMs (Figure 10). For math tasks, the targeted architectures tend to consist of an initial broad set of generations before being reduced quickly to a chosen answer (Figure 11). For coding tasks, the targeted architectures tend to focus on iterations of generation, critique, and fusion over a single response before outputting an answer (Figure 12).

3.3 ARCHON BY TASK

Instruction-Following and Reasoning: On MT Bench, AlpacaEval 2.0, and Arena-Hard-Auto, open-source ARCHON architectures outperform current open-source baselines by 10.5%, on average, while closed-source ARCHON outperforms current closed-source baselines by 14.6% (Table 1). With ARCHON, multiple models used for Generators and the depth of fusion layers lead to performance boosts on instruction-following tasks, increasing the richness of responses and allowing multiple iterations for step-by-step instruction-following (Table 25). For reasoning, while the performance boost from ARCHON is smaller when we consider the *aggregate* scores for MixEval and MixEval-Hard, we do see meaningful increases in performance when we create inference-time architectures for each individual task under MixEval and MixEval-Hard (Table 4; Table 5). When we create individual ARCHON architectures for each subtask, we see 3.7 and 8.9 percentage point increases in accuracy, on average, for MixEval and MixEval-Hard, respectively. This finding suggests that reasoning tasks (e.g. math, sciences, logic) require more individualized inference-time architectures.

Coding: We have observed that ensembling, fusion, and ranking techniques have limited impact on CodeContests (Figure 2.2). For example, when we apply the general all-source architecture from Table 2 to CodeContests problems, we achieve small gains from ARCHON (see Table 1). One contributing factor is that, unlike the distribution of instruction-following/reasoning tasks, coding tasks tend to have one or two LLMs that perform substantially better than the rest of models (Table 24). However, when we add unit test generation/evaluation, and scale the number of samples, ARCHON's performance on CodeContests improves significantly (Table 1), allowing us to boost GPT-40 Pass@1 performance by 44.3% for Pass@1 (from 40 to 58 out of 140 questions). For model-based unit test generation/evaluation, we generate 5 unit tests and use the LM to evaluate each candidate response against the generated unit tests, allowing us to rank the different candidate responses (details are provided in Section A.3)

3.4 DISCUSSION

Impact of Model Size: The ARCHON framework is most effective when utilizing LLMs with 70B+ parameters. When we build ARCHON architectures with 7B open-source models, we can boost task performance over the best individual 7B LM by 7.5%, on average, compared to the best individual 7B model (Table 27). Across tasks, 7B models work well for ranking but are less effective for critique and fusion.

Latency and Costs: Since ARCHON architectures make multiple LLM API calls successively for different operations (e.g. ensembling, critiquing, ranking, etc.), it can take 5x more time and money than a single LLM API call (Section A.6; Table 28; Table 29). Note that these increases in compute costs and latency translate to higher quality responses, and can be justified in many application domains, such as science, programming, and complex agentic tasks (Rein et al., 2023; Mialon et al., 2023). Furthermore, LLM vendors are rapidly decreasing their inference costs (Table 28). For tasks in which speed is most preferred, future work should explore how distillation strategies (Sreenivas et al., 2024; DeepSeek-AI et al., 2025) could be used to pack the aggregate knowledge of ARCHON architectures into a smaller LM.

4 ACKNOWLEDGEMENTS

We thank Simran Arora, Bradley Brown, Ryan Ehrlich, Sabri Eyuboglu, Jordan Juravsky, Jerry Liu, Benjamin Spector, Alyssa Unell, Benjamin Viggiano, and Michael Zhang for their constructive feedback during the composition of the paper. We would also like to thank our collaborators at the Stanford Artificial Intelligence Laboratory (SAIL) and TogetherAI.

We gratefully acknowledge the support of NIH under No. U54EB020405 (Mobilize), NSF under Nos. CCF2247015 (Hardware-Aware), CCF1763315 (Beyond Sparsity), CCF1563078 (Volume to Velocity), and 1937301 (RTML); US DEVCOM ARL under Nos. W911NF-23-2-0184 (Long-context) and W911NF-21-2-0251 (Interactive Human-AI Teaming); ONR under Nos. N000142312633 (Deep Signal Processing); Stanford HAI under No. 247183; NXP, Xilinx, LETI-CEA, Intel, IBM, Microsoft, NEC, Toshiba, TSMC, ARM, Hitachi, BASF, Accenture, Ericsson, Qualcomm, Analog Devices, Google Cloud, Salesforce, Total, the HAI-GCP Cloud Credits for Research program, the Stanford Data Science Initiative (SDSI), and members of the Stanford DAWN project: Meta, Google, and VMWare. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views, policies, or endorsements, either expressed or implied, of NIH, ONR, or the U.S. Government.

REFERENCES

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.

Anthropic. The claude 3 model family: Opus, sonnet, haiku. ArXiv, 2024.

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. arXiv preprint arXiv:2309.16609, 2023.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback, 2022. URL https://arxiv.org/abs/2212.08073.
- 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. URL https://arxiv.org/abs/2407.21787.

- Lingjiao Chen, Jared Quincy Davis, Boris Hanin, Peter Bailis, Ion Stoica, Matei Zaharia, and James Zou. Are more llm calls all you need? towards scaling laws of compound inference systems, 2024. URL https://arxiv.org/abs/2403.02419.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena: An open platform for evaluating llms by human preference, 2024.

Databricks. Dbrx technical report. 2024.

- Jared Quincy Davis, Boris Hanin, Lingjiao Chen, Peter Bailis, Ion Stoica, and Matei Zaharia. Networks of networks: Complexity class principles applied to compound ai systems design, 2024. URL https://arxiv.org/abs/2407.16831.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary,

Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The Ilama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

- Neel Guha, Mayee F Chen, Trevor Chow, Ishan S Khare, and Christopher Re. Smoothie: Label free language model routing. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. Deepseek-coder: When the large language model meets programming the rise of code intelligence, 2024. URL https://arxiv.org/abs/2401.14196.
- Eric Hartford. dolphin-2.2.1-mistral-7b. January 2024.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv* preprint arXiv:2103.03874, 2021.
- Geoffrey E Hinton et al. How neural networks learn from experience. na, 1992.
- Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. *arXiv preprint arXiv:2408.08435*, 2024.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023a.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024.
- Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. Llm-blender: Ensembling large language models with pairwise comparison and generative fusion. In *Proceedings of the 61th Annual Meeting of the Association for Computational Linguistics (ACL 2023)*, 2023b.
- Sayash Kapoor, Benedikt Stroebl, Zachary S Siegel, Nitya Nadgir, and Arvind Narayanan. Ai agents that matter. *arXiv preprint arXiv:2407.01502*, 2024.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*, 2023.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020.
- Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. More agents is all you need, 2024a. URL https://arxiv.org/abs/2402.05120.

- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Joseph E. Gonzalez Banghua Zhu, and Ion Stoica. From live data to high-quality benchmarks: The arena-hard pipeline, April 2024b. URL https://lmsys.org/blog/2024-04-19-arena-hard/.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval, 2023.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.
- Yu Meng, Mengzhou Xia, and Danqi Chen. SimPO: Simple preference optimization with a reference-free reward. *ArXiv*, 2024.
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants, 2023. URL https://arxiv.org/abs/2311.12983.
- Luigi Nardi, Artur Souza, David Koeplinger, and Kunle Olukotun. Hypermapper: a practical design space exploration framework. In 2019 IEEE 27th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS), pp. 425–426, 2019. doi: 10.1109/MASCOTS.2019.00053.
- Jinjie Ni, Fuzhao Xue, Xiang Yue, Yuntian Deng, Mahir Shah, Kabir Jain, Graham Neubig, and Yang You. Mixeval: Deriving wisdom of the crowd from llm benchmark mixtures, 2024. URL https://arxiv.org/abs/2406.06565.
- 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, 2024. URL https://arxiv.org/abs/2406.18665.
- **OpenAI.** Learning to reason with LLMs. https://openai.com/research/ learning-to-reason-with-llms, September 2024a. Accessed November 13, 2024.
- OpenAI. Learning to reason with large language models, 2024b. URL https: //openai.com/index/learning-to-reason-with-llms/. Accessed: 2024-09-12.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine

McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.

Qwen. Qwen2 technical report. 2024.

- 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 qa benchmark, 2023. URL https://arxiv.org/abs/2311.12022.
- Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Xiaojiang Chen, and Xin Wang. A comprehensive survey of neural architecture search: Challenges and solutions. *ACM Computing Surveys (CSUR)*, 54(4):1–34, 2021.
- Yijia Shao, Yucheng Jiang, Theodore A Kanell, Peter Xu, Omar Khattab, and Monica S Lam. Assisting in writing wikipedia-like articles from scratch with large language models. *arXiv preprint arXiv:2402.14207*, 2024.
- Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine learning algorithms, 2012. URL https://arxiv.org/abs/1206.2944.
- Sharath Turuvekere Sreenivas, Saurav Muralidharan, Raviraj Joshi, Marcin Chochowski, Mostofa Patwary, Mohammad Shoeybi, Bryan Catanzaro, Jan Kautz, and Pavlo Molchanov. Llm pruning and distillation in practice: The minitron approach, 2024. URL https://arxiv.org/abs/2408.11796.
- NovaSky Team. Sky-t1: Train your own o1 preview model within 450. https: //novasky-ai.github.io/posts/sky-t1,2025. Accessed: 2025-01-09.
- Qwen Team. Qwq: Reflect deeply on the boundaries of the unknown, November 2024. URL https://qwenlm.github.io/blog/qwq-32b-preview/.
- Hoang Tran, Chris Glaze, and Braden Hancock. Iterative dpo alignment. Technical report, Snorkel AI, 2023.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clementine Fourrier, Nathan Habib, et al. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*, 2023.
- Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities, 2024. URL https://arxiv.org/abs/2406.04692.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/2201.11903.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. WizardLM: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=CfXh93NDgH.
- Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and James Zou. Textgrad: Automatic "differentiation" via text. 2024.
- Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu. Aflow: Automating agentic workflow generation, 2024. URL https://arxiv.org/abs/2410.10762.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023.
- Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning, 2017. URL https://arxiv.org/abs/1611.01578.

A APPENDIX

A.1 RELATED WORK

Inference-time architectures combine multiple frozen LLMs and inference-time techniques, achieving superior performance compared to individual models. Notable works include Mixture-of-Agents (MoA) (Wang et al., 2024), LLM Blender (Jiang et al., 2023b), RouteLM (Ong et al., 2024), Smoothie (Guha et al.), and various approaches around *compound AI*, which are AI systems that use multiple components (e.g., LLMs, retrievers, tool use, APIs, etc.) (Chen et al., 2024; Davis et al., 2024; Lewis et al., 2020; Shao et al., 2024; Kapoor et al., 2024). LM frameworks like DSPy (Khattab et al., 2023) have emerged for orchestrating LMs and other components. Even with a single LLM, various techniques can improve performance by building better reasoning strategies, such as OpenAI's o1 (OpenAI, 2024b) and Chain of Thought (Wei et al., 2023), as well as inference-time frameworks, such as ADAS (Hu et al., 2024) and AFlow (Zhang et al., 2024).

Despite these advancements, challenges remain in developing inference-time architectures. Many architectures focus on additional generations (Jiang et al., 2023b; Chen et al., 2024; Davis et al., 2024), which is effective for reasoning tasks (Brown et al., 2024). However, for tasks like instruction-following and reasoning, techniques such as fusion and ranking are effective for bolstering task performances (Wang et al., 2024; Jiang et al., 2023b). Prior studies have explored limited aspects of configurations, often focusing on specific benchmarks (Jiang et al., 2023b; Wang et al., 2024; Chen et al., 2024; Li et al., 2024a). It's crucial to efficiently develop inference-time architectures, as optimal configurations vary based on benchmarks, available models, and inference compute limits (Section 3.2). Furthermore, LM orchestration frameworks, such as DSPy (Khattab et al., 2023), only optimize a single prompt for a single LM, better equipping it for tool use by utilizing supervised data but still unable to leverage multiple inference-time techniques in parallel or sequentially. While each of these approaches manually selects a subset of existing techniques, ARCHON unifies available inference-time techniques and automates architecture construction with search algorithms, simplifying the model and component selection process for each set of tasks (Sections 2.1 and 2.3).

A.2 ARCHON BENCHMARKS AND RESULTS

Benchmark	Example Count	Reference Model	Judge Model	Scoring Type	Metric
AlpacaEval 2.0	805	GPT-4-Turbo	GPT-4-Turbo	Pairwise Comparison	L.C. & Raw Win Rates
Arena-Hard-Auto	500	Claude-3.5-Sonnet GPT-4-0314	GPT-4-Turbo	Pairwise Comparison	Win Rate
MT-Bench	80	Claude-3.5-Sonnet	GPT-4-0314	Pairwise Comparison	Adjusted Win Rate
MixEval	2000	N/A	N/A	Ground Truth	Accuracy
MixEval-Hard	500	N/A	N/A	Ground Truth	Accuracy
MATH	200 (sampled from 5000)	N/A	N/A	Ground Truth	Pass@1
CodeContests	140 (non-visual queries)	N/A	N/A	Ground Truth	Pass@1





Figure 5: **Performance Gains from Repeated Sampling, Ensembling, Ranking, and Fusing on Arena-Hard-Auto**: The ARCHON win-rate continues to grow significantly as we scale model sampling (**left**) or add additional models to the generator ensemble (**right**), increasing by 9.3% and 18.5%, respectively. These best results are achieved by selecting the top-5 responses and fusing them. The ensemble models are added based on their individual performance on this task, from best to worse (Table 24). The oracle selection is the performance of picking the best answer generation out of all the generated samples from the ensemble. The results were averaged over 10 independent evaluation runs.

		Arena-	Hard-Auto
	Model / LLM System	Score	C.I.
	Claude 3.5 Sonnet GPT-40 Llama 3.1 405B Instruct	N/A 48.1% 28.4%	N/A (-2.3, 1.8) (-2.7, 2.5)
pen	General-purpose ARCHON Architecture	66.2%	(-2.4, 2.2)
O S	Task-specific ARCHON Architectures	69.0%	(-2.8, 2.5)
osed	General-purpose ARCHON Architecture	70.5%	(-2.5, 2.0)
Sol CIC	Task-specific ARCHON Architectures	74.4%	(-2.3, 1.6)
II	General-purpose ARCHON Architecture	72.5%	(-2.5, 1.8)
A Sou	Task-specific ARCHON Architectures	76.1%	(-1.8, 2.2)

Table 3: ARCHON Results on Arena-Hard-Auto Results with Claude-3.5-Sonnet as Baseline Model: The baseline model is Claude-3.5-Sonnet (default baseline model: GPT-4-0314) while the judge model is GPT-4-Turbo.

			MixH	Eval - Su	b-Datase	ts		
Model / LLM System	Infer. Calls	GSM8K	TriviaQA	DROP	MATH	BBH	AGIEval	Average
GPT-40 - 2024-05-13	1	94.9	89.1	88.2	98.5	98.3	71.5	90.3
Claude 3.5 Sonnet	1	98.0	92.0	92.6	96	95.6	78.0	92.0
Llama 3.1 405B Instruct	1	98.2	87.9	89.6	91.5	95.8	73.2	89.6
General-purpose ARCHON Architecture	29	98.3	94.8	94.6	98.1	97.3	82.1	94.2
Task-specific ARCHON Architectures	34	98.2	96.7	95.6	98.5	98.8	84.2	95.7

Table 4: **MixEval Results by Sub-Dataset**: For the average computed, we do not introduce any weighting for each dataset.

			MixEval - Sub-Datasets										
Model / LLM System	Infer. Calls	GSM8K	TriviaQA	DROP	MATH	BBH	AGIEval	Average					
GPT-40 - 2024-05-13	1	72.3	70.5	70.2	94.4	80.0	53.5	73.5					
Claude 3.5 Sonnet	1	87.3	75.5	79.3	82.5	80.0	74.6	79.9					
Llama 3.1 405B Instruct	1	98.7	71.2	70.7	86.9	78.8	62.0	78.1					
General-purpose ARCHON Architecture	33	96.7	82.7	83.2	93.4	82.0	76.7	85.8					
Task-specific ARCHON Architectures	37	98.9	86.2	85.2	96.2	86.0	80.1	88.8					

Table 5: **MixEval-Hard Results by Sub-Dataset**: For the average computed, we do not introduce any weighting for each dataset.

	GSM8K	MMLU Math	HumanEval Python	MBPP
Model	Pass@1	Pass@1	Pass@1	Pass@1
GPT-40	97.1%	84.8%	89.0%	87.5%
Claude 3.5 Sonnet	96.8%	90.9%	90.2%	88.9%
Llama 3.1 405B Instruct	95.9%	85.4%	90.2%	88.6%

Table 6: Additional Math and Code Benchmarks Explored

							Datasets				
			MT Bench	MT Alpaca Arena Arena MixEval Bench Eval 2.0 Hard Auto Hard Auto Hard							
	Judge Model		GPT-4 0314	GP Tu	T-4 rbo	GPT-4 Turbo	GPT-4 Turbo	N/A	N/A	N/A	
	Reference Model		Claude 3.5 Sonnet	GP Tu	T-4 rbo	Claude 3.5 Sonnet	GPT-4 Turbo	N/A	N/A	N/A	
	Model / LLM System	Infer. Calls	W.R.	L.C. W.R.	Raw W.R.	W.R.	W.R	Acc.	Acc.	Pass @1	
	GPT-40 - 2024-05-13 Claude 3.5 Sonnet Llama 3.1 405B Instruct	1 1 1	44.7% N/A 44.7%	57.5% 52.4% 40.3%	51.3% 40.6% 37.7%	48.1% N/A 28.4%	80.3% 80.9% 64.1%	63.6% 68.9% 66.2%	88.0% 89.7% 88.9%	84.5% 85.0% 83.5%	
	MoA MoA Lite	19 7	51.6% 45.6%	65.1% 59.3%	59.8% 57.0%	52.2% 40.6%	84.2% 87.8%	62.5% 61.1%	87.3% 87.1%	82.0% 83.0%	
Den	General-purpose ARCHON Architecture	35	67.5%	63.0%	68.3%	66.2%	85.1%	65.5%	86.9%	86.5%	
⊂ ÿ	ARCHON Architectures	44	71.6%	66.7%	70.7%	69.0%	89.5%	67.5%	89.6%	90.5%	
bsed	General-purpose ARCHON Architecture	32	73.1%	63.5%	69.1%	70.5%	85.8%	67.7%	88.2%	88.0%	
ž ž	Task-specific ARCHON Architectures	40	77.5%	68.4%	72.1%	74.4%	90.2%	72.9%	90.4%	89.5%	
vil urce	General-purpose ARCHON Architecture	35	76.8%	65.8%	70.2%	72.5%	89.3%	70.1%	88.1%	90.0%	
All Sour	Task-specific ARCHON Architectures	39	80.4%	67.6%	73.3%	76.1%	92.1%	72.9%	90.6%	93.5%	

Table 7: ARCHON's Strong Performance on the Complete Evaluation Datasets after ARCHON Architecture Optimization: We find that ARCHON's inference-time architectures consistently outperform single-call state-of-the-art LLMs, both open-source and closed-source baselines, when evaluating on the complete benchmarks (Table 2). We explore two configurations: architecture search for building custom ARCHON configurations for each individual benchmark and architecture search for building a single general-purpose ARCHON configuration for all the benchmarks (Section 3.1). We find that a general ARCHON configuration lags behind the custom ones by only 3.2 percentage points, on average, across our all-source settings, which suggests the efficacy of general-purpose inference-time architectures created with our framework. For Arena-Hard-Auto, we also include a configuration with Claude 3.5 Sonnet as a stronger reference model for comparison against ARCHON inference-time architectures and to mitigate bias from GPT judges towards GPT generations. For MT Bench, we use a GPT-4-0314 judge model instead of newer LLM judges to be consistent with previous results on this benchmark. For our task-specific ARCHON architectures, we also provide the average inference calls across the given benchmarks. For our full-list of models explored, please see Table 23. For MATH, we use a randomly sampled subset of size 200 for evaluation (Section 3.1; Table 2). We include our ARCHON architecture results on the held-out 80% subset of each evaluation benchmark in Table 1.

A.3 ARCHON LLM COMPONENTS

Inference-Time Technique	Definition	Input	Output	Inference Cost	Domains
Generator	Generates a candidate response from an instruction prompt	Instruction Prompt	Candidate Response(s)	1 call per cand.	All Domains
Fuser	Merges multiple candidate responses into a single response	Instruction Prompt + Candidate Response(s)	Fused Candidate Response(s)	1 call per cand.	All Domains
Critic	Generates strengths/weaknesses for each candidate response	Instruction Prompt + Candidate Response(s)	Candidate Response(s) Strengths/Weaknesses	1 call	All Domains
Ranker	Returns top-K candidate responses	Instruction Prompt + Candidate Response(s)	Ranked Candidate Response(s)	1 call	All Domains
Verifier	Returns the candidate responses with verified reasoning	Instruction Prompt + Candidate Response(s)	Verified Candidate Response(s)	2 calls per cand.	Reasoning Tasks
Unit Test Generator	Generates unit tests to evaluate the candidate responses	Instruction Prompt	Instruction Prompt + Unit Tests	1 call	Reasoning Tasks
Unit Test Evaluator	Uses generated unit tests to evaluate candidate response	Instruction Prompt + Unit Tests + Candidate Response(s)	Scored Candidate Response(s)	1 call per cand.	Reasoning Tasks

Table 8: **Overview of ARCHON's Inference-time Techniques**: Definitions, Inputs, Outputs, Costs, and Application Domains.

Module	Initial Layer Placement	Placement after Initial Layer	>1 Module in Layer	Increase Candidate Responses	Decrease Candidate Responses
Generator	Yes	No	Yes	Yes	No
Fuser	No	Yes	Yes	Yes	Yes
Ranker	No	Yes	No	No	Yes
Critic	No	Yes	No	No	No
Verifier	No	Yes	No	No	Yes
Unit Test Generator	No	Yes	No	No	No
Unit Test Evaluator	No	Yes	No	No	No

Table 9: **Rules of ARCHON Construction**: Allowed combinations of each LLM component from Section 2.1.

<instruction here>.

Table 10: Generator Prompt

<instruction here>

You have been provided with a set of responses with their individual critiques of strengths/weaknesses from various open-source models to the latest user query. Your task is to synthesize these responses into a single, high-quality response. It is crucial to critically evaluate the information provided in these responses and their provided critiques of strengths/weaknesses, recognizing that some of it may be biased or incorrect. Your response should not simply replicate the given answers but should offer a refined, accurate, and comprehensive reply to the instruction. Ensure your response is well-structured, coherent, and adheres to the highest standards of accuracy and reliability. Responses from models: 1. <response #1> Critique: <critique #1> 2. <response #2> Critique: <critique #2> ... N. <response #N> Critique: <critique #N>

((a)) With Critiques

You have been provided with a set of responses from various open-source models to the latest user query. Your task is to synthesize these responses into a single, high-quality response. It is crucial to critically evaluate the information provided in these responses, recognizing that some of it may be biased or incorrect. Your response should not simply replicate the given answers but should offer a refined, accurate, and comprehensive reply to the instruction. Ensure your response is well-structured, coherent, and adheres to the highest standards of accuracy and reliability.
1. <response #l>
2. <response #l>
...
N. <response #N>
<instruction here>

((b)) Without Critiques

Table 11: Fuser Prompt: Without and With Critiques



Table 12: Decoder-Based Ranking Prompt



Table 13: Critic Prompt

I will provide you with a response indicated by the identifier 'Response'. Provide reasoning for why the response accurately and completely addresses the instruction: <instruction here>

Response: <response>

Instruction: <instruction here>.

Provide the reasoning for the response above based on its relevance, completeness, and accuracy when compared to the instruction. Do not include any preface or text after the reasoning.

Table 14: Verifier Prompt

Instruction Prompt: Given the following query, generate a set of N unit tests that would evaluate the correctness of responses to this query

- The unit tests should cover various aspects of the query and ensure comprehensive evaluation - Each unit test should be clearly stated and should include the expected outcome.

The unit tests should be in the form of assertions that can be used to validate the correctness of responses to the query.
The unit test should be formatted like 'The answer mentions...', 'The answer states...', 'The answer uses...', etc. followed by the

expected outcome Solely provide the unit tests for the question below. Do not provide any text before or after the list. Only output the unit tests as a list of strings (e.g., ['unit test #1', 'unit test #2', 'unit test #3']).

Query: <instruction here>

((a)) With Unit Test Cap

Instruction Prompt: Given the following query, generate a set of unit tests that would evaluate the correctness of responses to this query

The unit tests should cover various aspects of the query and ensure comprehensive evaluation.
Each unit test should be clearly stated and should include the expected outcome.
The unit tests should be in the form of assertions that can be used to validate the correctness of responses to the query.
The unit test should be formatted like "The answer mentions...", "The answer states...", "The answer uses...", etc. followed by the

expected outcome

Solely provide the unit tests for the question below. Do not provide any text before or after the list. Only output the unit tests as a list of strings (e.g., ['unit test #1', 'unit test #2', 'unit test #3']). Query: <instruction here>

((b)) Without Unit Test Cap

Table 15: Unit Test Generator Prompt: With and Without Unit Test Cap

Instruction Prompt: Compose an engaging travel blog post about a recent trip to Hawaii, highlighting cultural experiences and must-see attractions

- 1. Unit Test #1: The blog post mentions at least two cultural experiences specific to Hawaii.
- Unit Test #2: The blog post highlights at least three must-see attractions in Hawaii.
- 3. Unit Test #3: The tone of the blog post is engaging and uses descriptive language that would appeal to readers interested in travel.
- Unit Test #4: The blog post includes factual information about Hawaii's culture, such as local customs, festivals, or historical facts. 4.
- 5. Unit Test #5: The blog post contains a clear narrative structure, including an introduction, main body, and a conclusion.

((a)) Instruction-Following Query

Instruction Prompt: Alice and Bob have two dice. They roll the dice together, note the sum of the two values shown, and repeat. For Alice to win, two consecutive turns (meaning, two consecutive sums) need to result in 7. For Bob to win, he needs to see an eight followed by a seven. Who do we expect to win this game?

1. Unit Test #1: The response correctly identifies the winning condition for Alice (two consecutive sums of 7).

Unit Test #2: The response correctly identifies the winning condition for Bob (a sum of 8 followed by a sum of 7). 2.

3 Unit Test #3: The response explains the probability of achieving two consecutive 7s when rolling two dice.

- Unit Test #4: The response explains the probability of achieving an 8 followed by a 7 when rolling two dice.
- 5. Unit Test #5: The response provides a conclusion on who is more likely to win based on the probability analysis.

((b)) Reasoning Query

Table 16: Unit Test Examples



Figure 6: **Performance Gains from Applying Inference Time Techniques on a Single Model**: We repeatedly sample more responses for each individual query. For each sample count, we choose the best response in 5 different ways: (1) using an oracle (to get the upper bound for performance of best sample), (2) randomly, (3) using a ranker model, (4) by fusion, in which a model synthesizes a response based on all the samples, and (5) by ranking the top-5 best answers and then fusing them. For both MT Bench and Arena-Hard-Auto, we find that fusion is an effective technique. In particular, ranking the candidates first, and then selecting the top-5 and fusing them scores the highest. The best open-source model for these tasks across all the 70B+ models we are considering is WizardLM-2-8x22B (Xu et al., 2024) (see Table 24 for details). For both ranking and fusion, we use Qwen2 72B Instruct (Qwen, 2024).

Given the following query, candidate response, and unit tests, evaluate whether or not the response passes each unit test. - In your evaluation, you should consider how the response aligns with the unit tests, retrieved documents, and query Provide reasoning before you return your evaluation. you must finish with a list of verdicts corresponding to each unit the end of your evaluation. At test. - You must include a verdict with one of these formatted options: '[Passed]' or '[Failed]' - Here is an example of the output format: Unit Test #1: [Passed] Unit Test #2: [Failed] Unit Test #3: [Passed] Each should be line verdict and the unit the on а new correspond to test in same position. - Here is the query, response, and unit tests for your evaluation: Query: <instruction here>. Candidate Response: <response> Unit Tests: Unit Test #1: <Unit Test #1> Unit Test #2: <Unit Test #2> Unit Test #N: <Unit Test #N>

Table 17: Unit Test Evaluator Prompt

A.4 UTILITIES AND INTERACTIONS OF LLM COMPONENTS

In this subsection, we present our analysis of the effectiveness of each LLM component (i.e. the *Utility*) and the relationships between each component (i.e. the *Component Interactions*) by evaluating on *instruction-following tasks* (MT Bench, AlpacaEval 2.0, Arena-Hard-Auto), *reasoning tasks* (MixEval, MixEval-Hard, MATH) and *coding tasks* (CodeContests) (Section 3.1). For our ARCHON models, we utilize a host of 70B+ open-source models (Section 3.1; Table 23).

A.4.1 GENERATOR

Utility: For our Generator module, we find additional model sampling to significantly boost performance (Figure 6), particularly for coding tasks (Table 1). In settings with a limited inference



Figure 7: **Performance Gains from Applying Inference-Time Techniques on an Ensemble of Models**: We incrementally add more models to the ensemble, which consists of open-source 70B+ models. The models are added to the pool based on their performance for each task, from best to worse (see Table 24 for details). For each ensemble size, we choose the best response in 5 different modes: (1) using an oracle (to get the upper bound for performance of best individual response in the ensemble), (2) randomly, (3) using a ranker model, (4) by fusion, in which one model synthesizes a response based on all the responses of the ensemble models, and (5) ranking the top-5 best responses and then fusing them. For MT Bench and Arena-Hard-Auto, we find consistent performance improvements as we add more models to the ensemble. We find that fusion is beneficial across various ensemble sizes and in particular a fused candidate based on the top-5 ranked responses scores highest. The ensemble approach scores higher than applying the same techniques on repeated samples from a single best-performing model (see Figure 6). For both ranking and fusion, we use Qwen2 72B Instruct (Qwen, 2024).

call budget, additional model samples lead to the largest marginal benefit. We see a similar pattern for model ensembling, where sampling from additional models leads to continual performance increases (assuming the models are ordered from best to worst for the given task) (Figure 7).

A.4.2 FUSER

Utility: For every benchmark explored, we found that the Fuser module substantially improved performance (Figure 6; Figure 7; Figure 2.2). For the single-generation 10-model ensemble of 70B+ models, the Fuser module improved downstream accuracy by 5.2 points, on average, compared to the single-generation best model (Figure 7). When combined with the Ranker module for ranking the top-5 candidate responses, the Fuser improved downstream accuracy by 7.3 points and 3.6 points, on average, compared to the single-sample best model and the oracle best candidate response, respectively (Figure 7). Overall, we found that Fuser efficacy increased as more candidate responses were provided, demonstrating that additional candidate generations can continue to bolster inference-time architecture performance when combined with a Fuser.

In previous work like Mixture-of-Agents (MoA) (Wang et al., 2024), multiple layers of Fusers was found to boost performance on some instruction-following tasks (i.e. MT Bench and Alpaca Eval 2.0). Across all the benchmarks explored, we observed similar benefits in the ARCHON framework when adding multiple layers of Fusers (Figure 2.2). However, based on our results in Figure 8, the number of Fuser layers needed to improve performance varied by task, with some tasks receiving limited benefits from added layers (1-2 point increase in accuracy for MixEval) while others experienced significant benefits with 3-4 fusion layers and more (2 to 5 point increase in win rate for MT Bench and Alpaca Eval 2.0). We attribute this distinction to the difference in task requirements, with chat and instruction following tasks benefiting more from multiple iterations of revisions through the multiple Fuser layers, leading to greater diversity in the final generation (Table 25).

Component Interactions: To better understand how the Fuser module works with the other LLM components, we took the single-sample 10-model ensemble of Generators with a Fuser and tried adding each of these components individually: a Critic, a Ranker, a Verifier, and a Unit Test Generator/Evaluator. Across all of the benchmarks, the added candidate response analyses from the Critic improved the Fuser's ability to effectively merge the different candidate responses, increasing performance by an average of 3.1 percentage points (Figure 2.2). With the added Ranker, the ARCHON architecture improved the combined Ensemble + Critic + Fuser performance across all the benchmarks by 4.8 percentage points, on average (Figure 2.2). The Ranker proved most effective for style-oriented tasks (e.g. MT Bench and AlpacaEval 2.0) since the examples mostly focus on improving the instruction-guidance towards the provided prompt. With the added Verifier module (Figure 2.2), the performance of the Ensemble + Critic + Fuser configuration improved marginally for the instruction-following tasks (1.2 percentage points, on average, for MT Bench, AlpacaEval 2.0, and Arena-Hard-Auto). However, this configuration improved performance more on reasoning tasks (3.2 percentage points for MixEval and MixEval-Hard, on average), assisting generation by filtering out irrelevant or flawed answers before the final fusion step (Figure 2.2). The added Unit Test Generator and Evaluator was less effective for the instruction-following and reasoning tasks, only providing a 1.5 percentage points increase, on average, when added to the Ensemble + Critic + Fuser configuration (Table 18). However, for coding tasks, we found unit test generation and evaluation significantly improved performance, leading to a 10.7 percentage point increase (56% performance increase comparatively) as we scale model sampling (Table 1).

A.4.3 CRITIC

Utility: The Critic module proved effective for every task we explored in Figure 2.2 and Table 18. With our 10-model 70B+ Generator ensemble and Fuser configuration of ARCHON, the added Critic improved performance on average by 3.1 percentage points across the benchmarks explored.

Component Interactions: While useful for most ARCHON architectures, the added strengths and weaknesses from the Critic module are particularly useful when combined with the Fuser module, helping guide generation fusion for a single layer and even useful when placed between multiple fusion layers (on average 3.2 percentage point boost across benchmarks in Figure 2.2). The Critic module was also effective with the Ranker module, providing additional information for comparing candidate responses (Figure 6) and leading to a 5.9 percentage point increase, on average (Table 18).

A.4.4 RANKER

Utility: From our results in Table 18, Figure 6, and Figure 7, we found the Ranker to be most effective for instruction-following tasks, where pair-wise comparisons of answers focus on style and adherence to the prompt. To examine the candidate selection improvement provided by candidate ranking, we compare three approaches to the Ranker: (1) random selection of candidate generation, (2) oracle selection of candidate generation, and (3) the top-ranked candidate selected by our Ranker. For MT Bench and Arena-Hard-Auto, we find that the ranker improves generation output quality by 3.8% compared to random candidate selection and performs within 2.7% of oracle selection (Figure 6).

Component Interactions: Based on our benchmark results in Table 18, the Ranker pairs well with the Critic module; the provided strengths and weaknesses helps guide ranking, particularly for instruction-following tasks, improving performance by 5.9 percentage points, on average. Furthermore, the Ranker was also effective when paired with the Fuser; the filtered list of candidate responses helped improve the final condensed response produced by the Fuser by 3.8 percentage points, on average (Figure 7). When paired with the Verifier and Unit Test Generator, the Ranker had neutral effects; performances changed marginally, either positively or negatively by 1-2 percentage points (Table 18).

Overall, our findings demonstrate the value of added Rankers for instruction-following and reasoning tasks when paired with Fusers. We find that when Rankers are used alone with an ensemble of Generators, their performance lags behind the 10-sample best single model configuration by 3.0 percentage points, on average (Table 18). Additionally, our findings show the importance of building better rankers for more complex reasoning tasks, such as math and coding, which is a challenge also raised by Brown et al. (2024).

A.4.5 VERIFIER

Utility: The Verifier was most effective for the reasoning benchmarks explored in Table 18. When just using a 70B+ Generator ensemble with Verifier module after generation, the ARCHON configuration lagged behind the ARCHON ensemble and fuser configuration by 1.5 percentage points, on average, across all benchmarks explored. This suggests that the Verifier is most effective when combined with other inference-time techniques.

			MT Bench	Alpac 2	aEval .0	Arena Hard Auto	MixEval Hard	MixEval	MATH	Code Contests
-	Model / LLM System	# of Infer. Calls	W.R.	L.C. W.R.	Raw W.R.	W.R.	Acc.	Acc.	Acc.	Acc.
Control	Best Open-Source 70B+ Model, Sampled Once Ensemble + Fuser Ensemble + Critic + Fuser	1 9 10	55.0% ±0.4 58.4% ±0.6 60.9% ±0.3	44.7% ±0.5 57.5% ±0.4 58.7% ±0.6	37.1% ±0.6 51.3% ±0.5 <u>65.8% ±0.3</u>	45.6% ±0.5 54.3% ±0.7 58.8% ±0.4	58.7% ±0.2 60.1% ±0.5 61.7% ±0.5	86.5% ±0.3 87.3% ±0.2 87.4% ±0.3	84.5% ±0.6 85.5% ±0.3 <u>87.2% ±0.5</u>	22.5% ±0.3 23.1% ±0.7 24.9% ±0.4
Ablations	Ensemble + Ranker Ensemble + Verifier Ensemble + Unit Test Gen./Eval. Ensemble + Ranker + Fuser Ensemble + Verifier + Fuser Ensemble + Unit Test Gen./Eval. + Fuser Ensemble + Critic + Verifier + Fuser Ensemble + Critic + Ranker + Fuser	9 24 18 10 25 17 25 11	$52.5\% \pm 0.7$ $53.2\% \pm 0.5$ $51.5\% \pm 0.4$ $62.5\% \pm 0.3$ $61.4\% \pm 0.6$ $61.3\% \pm 0.5$ $64.7\% \pm 0.4$	$54.7\% \pm 0.5$ $56.2\% \pm 0.3$ $54.4\% \pm 0.6$ $60.3\% \pm 0.4$ $59.4\% \pm 0.7$ $58.5\% \pm 0.5$ $60.0\% \pm 0.3$ $62.6\% \pm 0.6$	$\begin{array}{c} 47.6\% \pm 0.4\\ 50.2\% \pm 0.7\\ 49.4\% \pm 0.5\\ 63.6\% \pm 0.6\\ 58.7\% \pm 0.3\\ 55.1\% \pm 0.4\\ 61.0\% \pm 0.7\\ \textbf{72.4\%} \pm \textbf{0.5} \end{array}$	$50.5\% \pm 0.6$ $52.4\% \pm 0.3$ $46.1\% \pm 0.8$ $57.2\% \pm 0.5$ $59.2\% \pm 0.4$ $56.4\% \pm 0.7$ $\underline{59.5\% \pm 0.3}$ $60.9\% \pm 0.6$	$58.7\% \pm 0.3$ $55.9\% \pm 0.5$ $55.2\% \pm 0.4$ $59.7\% \pm 0.2$ $68.3\% \pm 0.3$ $63.9\% \pm 0.3$ $65.8\% \pm 0.4$ $\underline{66.8\% \pm 0.4}$	$\begin{array}{c} 86.8\% \pm 0.4\\ 85.6\% \pm 0.2\\ 86.0\% \pm 0.3\\ 87.6\% \pm 0.3\\ 87.5\% \pm 0.2\\ 86.9\% \pm 0.3\\ \underline{87.8\% \pm 0.4}\\ \overline{88.3\% \pm 0.2}\end{array}$	80.4% ±0.4 85.2% ±0.7 85.2% ±0.5 85.3% ±0.6 86.7% ±0.4 86.4% ±0.8 86.1% ±0.3 87.3% ±0.5	$24.1\% \pm 0.425.3\% \pm 0.524.6\% \pm 0.624.0\% \pm 0.226.3\% \pm 0.628.0\% \pm 0.628.0\% \pm 0.628.0\% \pm 0.626.8\% \pm 0.325.5\% \pm 0.3$

Table 18: **Impact of Different Compositions of ARCHON's Inference-Time Techniques**: We see increased task performances from adding new LLM components to ARCHON. For CodeContests, we find that there is a single model (Llama 3.1 405B Instruct) that performs considerably better than the rest of the LLMs studied, making it more effective leverage additional model sampling (Table 1). For our ensemble, we use the best 8 open-source 70B+ models for the task (Table 24). For our fuser, critic, ranker, and verifier components, we use the best fuser model found for the task (Table 24). For each evaluation benchmark, we explain its configuration in Table 2 and Section 3.1. The standard error numbers were calculated from 10 independent evaluation runs.

			MT Bench	AlpacaEval 2.0	Arena Hard Auto	MixEval Hard	MixEval	MATH	Code Contests
	Model / LLM System	# of Infer. Calls	W.R.	L.C. W.R.	W.R.	Acc.	Acc.	Acc.	Acc.
6	Single Generation	1	44.2% ±0.6	57.8% ±0.5	48.1% ±0.7	63.4% ±0.3	87.5% ±0.2	82.1% ±0.4	17.9% ±0.3
ntr	Ensemble + Fuser	11	53.7% ±0.3	59.5% ±0.6	49.7% ±0.5	65.5% ±0.2	82.0% ±0.3	81.0% ±0.6	16.0% ±0.4
ŭ	Ensemble + Critic + Fuser	12	$56.1\%\pm\!\!0.7$	$59.7\%\pm\!0.4$	$53.9\%\pm\!0.6$	$67.4\%\pm\!0.4$	$82.0\%\pm\!0.2$	$82.3\%\pm\!0.5$	18.9% ±0.6
	Ensemble + Ranker	11	47.6% ±0.4	49.7% ±0.5	45.5% ±0.4	63.3% ±0.3	81.6% ±0.4	77.3% ±0.7	17.9% ±0.5
	Ensemble + Verifier	11	$48.4\% \pm 0.5$	$51.2\% \pm 0.7$	$47.7\% \pm 0.8$	$61.4\% \pm 0.2$	80.5% ±0.3	75.5% ±0.3	$23.0\% \pm 0.4$
SU	Ensemble + Unit Test Gen./Eval.	21	$46.8\% \pm 0.8$	49.3% ±0.3	$41.2\% \pm 0.5$	$60.2\% \pm 0.4$	$80.7\% \pm 0.2$	$78.9\% \pm 0.8$	24.0% ±0.7
tio	Ensemble + Ranker + Fuser	12	58.0% ±0.2	$60.1\% \pm 0.6$	52.2% ±0.3	65.0% ±0.3	$82.0\% \pm 0.4$	$82.1\% \pm 0.4$	18.0% ±0.3
bla	Ensemble + Verifier + Fuser	12	55.8% ±0.6	$54.2\% \pm 0.4$	60.3% ±0.7	$67.0\% \pm 0.2$	82.5% ±0.3	83.1% ±0.6	$22.4\% \pm 0.5$
V	Ensemble + Unit Test Gen./Eval. + Fuser	22	56.5% ±0.3	$61.4\% \pm 0.5$	$51.6\% \pm 0.4$	$67.7\% \pm 0.4$	$81.7\% \pm 0.2$	$84.3\% \pm 0.5$	$25.4\% \pm 0.6$
	Ensemble + Critic + Verifier + Fuser	13	56.6% ±0.7	62.0% ±0.3	55.0% ±0.6	68.5% ±0.3	$82.7\%\pm\!\!0.4$	85.7% ±0.3	$22.2\% \pm 0.4$
	Ensemble + Critic + Ranker + Fuser	13	60.0% ±0.4	62.8% ±0.6	56.2% ±0.5	69.4% ±0.2	88.5% ±0.3	87.0% ±0.7	18.5% ±0.5

Table 19: **ARCHON Component Compositions with GPT-40**: The ensemble uses generates 10 samples for the given query. The standard error numbers were calculated from 10 independent evaluation runs.

			MT Bench	AlpacaEval 2.0	Arena Hard Auto	MixEval Hard	MixEval	MATH	Code Contests
	Model / LLM System	# of Infer. Calls	W.R.	L.C. W.R.	W.R.	Acc.	Acc.	Acc.	Acc.
10	Single Generation	1	32.1% ±0.7	38.5% ±0.5	30.4% ±0.6	45.2% ±0.3	69.5% ±0.2	72.3% ±0.5	10.5% ±0.6
nt.	Ensemble + Fuser	11	44.2% ±0.3	43.0% ±0.6	$40.2\% \pm 0.4$	$46.0\%\pm\!0.4$	73.0% ±0.3	70.5% ±0.7	$6.0\% \pm 0.4$
ŭ	Ensemble + Critic + Fuser	12	$46.6\%\pm\!0.5$	$44.2\%\pm\!0.4$	$44.4\% \pm 0.7$	$47.9\% \pm 0.2$	$73.0\% \pm 0.4$	$72.5\% \pm 0.3$	$8.4\%\pm\!0.5$
	Ensemble + Ranker	11	38.1% ±0.6	40.2% ±0.7	36.0% ±0.5	43.8% ±0.3	72.1% ±0.2	66.2% ±0.6	7.5% ±0.4
	Ensemble + Verifier	11	38.9% ±0.4	41.7% ±0.3	$38.2\% \pm 0.8$	$41.9\%\pm\!\!0.4$	71.0% ±0.3	$68.5\% \pm 0.4$	19.0% ±0.7
SU	Ensemble + Unit Test Gen./Eval.	21	37.3% ±0.8	39.8% ±0.6	$31.7\% \pm 0.3$	$40.7\%\pm\!0.2$	71.2% ±0.4	69.8% ±0.8	22.0% ±0.3
tio	Ensemble + Ranker + Fuser	12	48.0% ±0.2	45.6% ±0.5	$42.7\% \pm 0.6$	$45.0\%\pm\!0.3$	73.0% ±0.2	70.1% ±0.5	8.0% ±0.6
bla	Ensemble + Verifier + Fuser	12	46.3% ±0.5	$44.7\% \pm 0.4$	$45.0\% \pm 0.4$	$50.5\% \pm 0.4$	73.0% ±0.3	71.3% ±0.3	18.6% ±0.5
	Ensemble + Unit Test Gen./Eval. + Fuser	22	47.0% ±0.3	43.9% ±0.7	42.1% ±0.7	$48.2\%\pm\!\!0.2$	$72.2\% \pm 0.4$	73.1% ±0.6	$23.5\% \pm 0.4$
	Ensemble + Critic + Verifier + Fuser	13	47.1% ±0.7	46.0% ±0.3	$45.0\% \pm 0.5$	52.4% ±0.3	73.2% ±0.5	74.1% ±0.4	18.4% ±0.7
	Ensemble + Critic + Ranker + Fuser	13	50.5% ±0.4	48.3% ±0.6	46.7% ±0.3	55.1% ±0.4	73.7% ±0.3	76.4% ±0.5	$8.1\% \pm 0.5$

Table 20: **ARCHON Component Compositions with GPT-4o-mini**: The ensemble uses generates 10 samples for the given query. The standard error numbers were calculated from 10 independent evaluation runs.

			MT Bench	AlpacaEval 2.0	Arena Hard Auto	MixEval Hard	MixEval	MATH	Code Contests
	Model / LLM System	# of Infer. Calls	W.R.	L.C. W.R.	W.R.	Acc.	Acc.	Acc.	Acc.
10	Single Generation	1	N/A	52.7% ±0.4	81.4% ±0.6	68.7% ±0.3	89.1% ±0.2	83.5% ±0.5	12.5% ±0.3
THE I	Ensemble + Fuser	11	N/A	53.0% ±0.6	83.2% ±0.4	$69.5\% \pm 0.2$	89.0% ±0.3	81.8% ±0.6	17.0% ±0.4
ŭ	Ensemble + Critic + Fuser	12	N/A	$54.2\%\pm\!0.3$	$\underline{85.4\%\pm0.7}$	$\underline{70.9\% \pm 0.4}$	$89.5\%\pm\!0.2$	$82.6\%\pm\!0.4$	19.4% ±0.6
_	Ensemble + Ranker	11	N/A	50.2% ±0.5	85.7% ±0.5	63.8% ±0.3	82.1% ±0.4	80.2% ±0.7	18.5% ±0.5
	Ensemble + Verifier	11	N/A	51.7% ±0.7	78.2% ±0.3	$60.9\% \pm 0.2$	81.0% ±0.3	80.1% ±0.3	$21.0\% \pm 0.4$
2	Ensemble + Unit Test Gen./Eval.	21	N/A	$49.8\% \pm 0.4$	71.7% ±0.8	59.0% ±0.2	81.2% ±0.2	$80.9\% \pm 0.8$	22.0% ±0.7
tio	Ensemble + Ranker + Fuser	12	N/A	$55.6\% \pm 0.5$	$82.7\% \pm 0.4$	$65.0\% \pm 0.3$	89.0% ±0.4	$82.4\% \pm 0.4$	19.0% ±0.3
bla	Ensemble + Verifier + Fuser	12	N/A	54.7% ±0.3	85.0% ±0.6	70.5% ±0.2	89.3% ±0.3	84.1% ±0.6	21.6% ±0.5
•	Ensemble + Unit Test Gen./Eval. + Fuser	22	N/A	53.9% ±0.6	82.1% ±0.5	$68.2\% \pm 0.4$	89.2% ±0.2	82.0% ±0.5	23.5% ±0.6
	Ensemble + Critic + Verifier + Fuser	13	N/A	56.0% ±0.4	85.0% ±0.3	71.0% ±0.3	89.4% ±0.4	83.1% ±0.3	$21.4\% \pm 0.4$
	Ensemble + Critic + Ranker + Fuser	13	N/A	58.3% ±0.5	86.7% ±0.7	73.0% ±0.2	89.7% ±0.3	85.3% ±0.7	19.1% ±0.5

Table 21: **ARCHON Component Compositions with Claude 3.5 Sonnet**: The ensemble uses generates 10 samples for the given query. The standard error numbers were calculated from 10 independent evaluation runs.

			MT Bench	AlpacaEval 2.0	Arena Hard Auto	MixEval Hard	MixEval	MATH	Code Contests
	Model / LLM System	# of Infer. Calls	W.R.	L.C. W.R.	W.R.	Acc.	Acc.	Acc.	Acc.
Control	Single Generation Ensemble + Fuser Ensemble + Critic + Fuser	1 11 12	35.0% ±0.5 48.2% ±0.3 50.6% ±0.7	42.0% ±0.6 47.0% ±0.4 48.2% ±0.5	36.8% ±0.7 44.2% ±0.5 48.4% ±0.3	64.6% ±0.2 66.5% ±0.3 68.1% ±0.4	73.2% ±0.3 77.0% ±0.2 77.0% ±0.4	74.3% ±0.4 75.1% ±0.7 76.3% ±0.5	10.0% ±0.5 10.8% ±0.3 11.5% ±0.6
Ablations	Ensemble + Ranker Ensemble + Verifier Ensemble + Unit Test Gen./Eval. Ensemble + Ranker + Fuser Ensemble + Verifier + Fuser Ensemble + Unit Test Gen./Eval. + Fuser Ensemble + Critic + Verifier + Fuser Ensemble + Critic + Ranker + Fuser	11 11 12 12 12 22 13 13	$\begin{array}{c} 42.1\% \pm 0.4\\ 42.9\% \pm 0.6\\ 41.3\% \pm 0.8\\ \underline{52.0\% \pm 0.2}\\ 50.3\% \pm 0.5\\ 51.0\% \pm 0.3\\ 51.1\% \pm 0.7\\ \mathbf{54.5\% \pm 0.4} \end{array}$	$\begin{array}{c} 44.2\%\pm0.7\\ 45.7\%\pm0.3\\ 43.8\%\pm0.6\\ 49.6\%\pm0.5\\ 48.7\%\pm0.4\\ 47.9\%\pm0.7\\ \underline{50.0\%\pm0.3}\\ \textbf{52.3\%\pm0.6} \end{array}$	$\begin{array}{c} 40.0\% \pm 0.6\\ 42.2\% \pm 0.8\\ 35.7\% \pm 0.4\\ 46.7\% \pm 0.7\\ 48.7\% \pm 0.5\\ 46.1\% \pm 0.6\\ \underline{49.0\% \pm 0.4}\\ \textbf{50.7\% \pm 0.3} \end{array}$	$58.8\% \pm 0.3$ $57.9\% \pm 0.2$ $55.7\% \pm 0.4$ $60.0\% \pm 0.3$ $67.5\% \pm 0.2$ $64.2\% \pm 0.4$ $\underline{68.0\% \pm 0.3}$ $70.4\% \pm 0.2$	$\begin{array}{c} 76.1\% \pm 0.2 \\ 75.0\% \pm 0.3 \\ 75.2\% \pm 0.2 \\ 77.0\% \pm 0.4 \\ 77.0\% \pm 0.3 \\ 76.2\% \pm 0.2 \\ \underline{77.2\% \pm 0.4} \\ \textbf{77.7\% \pm 0.3} \end{array}$	$\begin{array}{c} 71.8\% \pm 0.6 \\ 70.5\% \pm 0.4 \\ 74.1\% \pm 0.8 \\ 75.0\% \pm 0.5 \\ 77.4\% \pm 0.3 \\ 78.3\% \pm 0.6 \\ \underline{77.8\%} \pm 0.3 \\ \textbf{80.5\%} \pm \textbf{0.5} \end{array}$	$\begin{array}{c} 11.9\% \pm 0.4\\ \underline{12.0\% \pm 0.7}\\ 13.0\% \pm 0.3\\ 12.0\% \pm 0.6\\ 10.5\% \pm 0.5\\ \textbf{14.3\% \pm 0.4}\\ 10.0\% \pm 0.7\\ 11.5\% \pm 0.5 \end{array}$

Table 22: **ARCHON Component Compositions with Claude-3-Haiku**: The ensemble uses generates 10 samples for the given query. The standard error numbers were calculated from 10 independent evaluation runs.

Component Interactions: As noted in Section A.4.2, the Verifier augmented the performance of the Critic and Fuser on reasoning tasks (e.g. Arena-Hard-Auto, MixEval, MixEval-Hard), boosting performance by 3.7 percentage points, on average, when combined together with these modules. Overall, the Verifier is most powerful when augmenting additional components for tasks requiring verification of intermediate steps and the final response (Table 18). Therefore, the Verifier was less helpful for instruction-following tasks (e.g. MT Bench and AlpacaEval) but more effective for reasoning tasks (e.g. Arena-Hard-Auto and MixEval).

A.4.6 UNIT TEST GENERATOR AND EVALUATOR

Utility: The Unit Test Generator and Evaluator were most effective on reasoning and coding tasks, improving performance on benchmarks that required more verification steps, such as Arena-Hard-Auto, MixEval, MixEval-Hard, MATH, and CodeContests (Table 18). For the reasoning tasks, we found the unit test generator and evaluator to be most effective when combined with other components. When the 70B+ ensemble of Generators was only combined with unit tests, it was less effective for reasoning tasks like Arena-Hard-Auto and MixEval, lagging behind the ensemble and fuser configuration by 3.1 percentage points. This inspired us to look into other inference-time techniques combinations for unit test generation, such as increased sampling and fusion. When we increased generation sampling and added unit test generation/evaluation for CodeContests, we see a 56% boost in Pass@1 performance (Table 1), increasing from 17.9 to 29.3 Pass@1.

Component Interactions: When combined with the Fuser module, the Unit Test Generator and Evaluator improved performance by 2.1 percentage points across the benchmarks explored (Table 18). The combined ensemble, Unit Test Generator/Evaluator, and Fuser ARCHON configuration was most effective on the reasoning benchmarks, leading to a 2.5 percentage point boost, on average. For coding, the unit test generator and evaluator was most effective when combined with the best performing Generator (using large sample counts) and a final Fuser (subsection 3.2).

	MTI	Bench	Alpaca	Eval 2.0	Arena I	Hard Auto	Mix	Eval	MixEv	al Hard	MA	АТН	CodeC	Contests
Models	Gen	Fusion	Gen	Fusion	Gen	Fusion	Gen	Fusion	Gen	Fusion	Gen	Fusion	Gen	Fusion
GPT-40	44.7%	61.9%	57.5%	64.5%	48.1%	69.2%	88.0%	89.4%	63.6%	65.4%	82.0%	81.0%	17.9%	19.4%
GPT-4-Turbo	42.2%	63.1%	55.0%	65.8%	48.1%	61.9%	88.9%	89.0%	64.1%	64.4%	79.5%	73.5%	9.3%	14.2%
Claude 3 Opus	30.9%	57.2%	40.5%	N/A	27.0%	47.9%	88.3%	88.2%	63.6%	64.0%	74.5%	74.0%	10.0%	12.5%
Claude 3.5 Sonnet	N/A	71.9%	52.37%	63.6%	N/A	73.2%	89.7%	89.3%	68.9%	69.5%	83.5%	86.5%	12.1%	15.5%
Qwen 2 72B Instruct	35.0%	59.7%	37.48%	56.0%	14.5%	49.5%	86.5%	87.5%	58.7%	61.1%	81.0%	78.5%	3.6%	5.2%
DeepSeek LLM 67B Instruct	18.4%	20.0%	17.8%	17.1%	N/A	N/A	79.2%	N/A	42.5%	N/A	57.0%	N/A	5.7%	N/A
Qwen 1.5 72B Chat	24.7%	46.3%	36.6%	55.7%	14.4%	36.4%	84.5%	82.5%	50.3%	52.2%	71.5%	67.5%	15.0%	13.9%
Qwen 1.5 110B Chat	34.4%	50.3%	43.6%	55.9%	21.9%	39.7%	85.3%	86.5%	51.8%	55.6%	67.0%	75.5%	3.6%	7.8%
Wizard 8x22B	53.8%	57.2%	44.7%	50.6%	45.6%	51.2%	83%	78.1%	54.3%	50.4%	76.0%	60.5%	7.1%	10.4%
Llama 3.1 8B Instruct	33.1%	45.9%	25.6%	34.9%	11.9%	28.6%	75.0%	57.5%	41.3%	46.5%	65.5%	60.5%	8.6%	7.8%
Llama 3.1 70B Instruct	45.0%	51.9%	35.6%	40.2%	23.8%	37.2%	85.7%	83.5%	61.1%	65.5%	74.0%	73.5%	20.7%	23.4%
Llama 3.1 405B Instruct	44.7%	N/A	40.3%	N/A	28.4%	N/A	88.9%	N/A	66.2%	N/A	78.0%	N/A	27.1%	N/A

Table 24: **ARCHON Generation and Fusion Performances for Single Models**: For Alpaca Eval 2.0, we use the length-controlled win rate (LC WR). For fusion, we gather one candidate from each of the top-10 generator models.

A.5 ARCHON LLM ANALYSIS

Model	Source Code	Parameter Count	Max Sequence Length
GPT-40 (OpenAI et al., 2024)	Closed-Source	_	128K
GPT-4-Turbo (OpenAI et al., 2024)	Closed-Source	_	128K
Claude-3-Opus (Anthropic, 2024)	Closed-Source	_	200K
Claude-3.5-Sonnet (Anthropic, 2024)	Closed-Source	_	200K
Claude-3-Haiku (Anthropic, 2024)	Closed-Source	_	200K
Llama-3.1-70B-Instruct (Dubey et al., 2024)	Open-Source	70B	8k
Llama-3.1-405B-Instruct (Dubey et al., 2024)	Open-Source	70B	8k
DeepSeek LLM 67B Chat (Guo et al., 2024)	Open-Source	67B	32k
Qwen2 72B Instruct (Qwen, 2024)	Open-Source	72B	32k
Qwen1.5 110B Chat (Bai et al., 2023)	Open-Source	110B	32k
Qwen1.5 72B Chat (Bai et al., 2023)	Open-Source	72B	32k
Mixtral 8x22B v0.1 (Jiang et al., 2024)	Open-Source	176B	32k
WizardLM 8x22B (Xu et al., 2024)	Open-Source	176B	32k
dbrx-instruct (Databricks, 2024)	Open-Source	132B	32k
princeton-nlp/Llama-3-Instruct-8B-SimPO (Meng et al., 2024)	Open-Source	8B	8k
princeton-nlp/Llama-3-Instruct-8B-DPO (Meng et al., 2024)	Open-Source	8B	8k
princeton-nlp/Llama-3-Instruct-8B-RDPO (Meng et al., 2024)	Open-Source	8B	8k
princeton-nlp/Llama-3-Instruct-8B-IPO (Meng et al., 2024)	Open-Source	8B	8k
Llama-3.1-8B-Instruct (Dubey et al., 2024)	Open-Source	8B	8k
Qwen2-7B-Instruct (Qwen, 2024)	Open-Source	7B	32k
Qwen/Qwen1.5-7B-Chat (Bai et al., 2023)	Open-Source	7B	32k
mistralai/Mistral-7B-Instruct-v0.2 (Jiang et al., 2023a)	Open-Source	7B	32k
cognitivecomputations/dolphin-2.2.1-mistral-7b (Hartford, 2024)	Open-Source	7B	32k
microsoft/Phi-3-mini-4k-instruct (Abdin et al., 2024)	Open-Source	4B	4k
HuggingFaceH4/zephyr-7b-beta (Tunstall et al., 2023)	Open-Source	7B	32k
microsoft/Phi-3-small-8k-instruct (Abdin et al., 2024)	Open-Source	7B	8k
snorkelai/Snorkel-Mistral-PairRM-DPO (Tran et al., 2023)	Open-Source	7B	32k
mistralai/Mistral-7B-Instruct-v0.3 (Jiang et al., 2023a)	Open-Source	7B	32k

Table 23: Models Tested with ARCHON.

		Jaccard S	Similarity (%)			
Inference-Time Architecture	MT Bench	AlpacaEval 2.0	Arena-Hard Auto	MixEval	MixEval Hard	MATH	Code Contests
Best Open-Source 70B+ Model, Sampled 8 Times + Fuser	45.3%	52.1%	48.4%	55.2%	58.9%	65.2%	63.7%
Ensemble (8 Top Models), Sampled Once Each + Fuser	31.6%	34.1%	28.9%	38.6%	40.9%	57.1%	53.4%

Table 25: Jaccard Similarities between Candidates Responses and Fused Response by Benchmark: For the fuser, we use the best-performing 70B+ model for each benchmark.



Figure 8: **Fusion Layer Efficacy by Benchmark**: From solely scaling the fusion layers, we see limited benefits across the benchmarks explored but when we add other inference-time techniques, such as Critic and Ranker, we see increased downstream performance as we continue scaling inference-time compute (Figure 2.2). We use an 8-model ensemble of the top Generator models for each benchmark (Table 24). For our Fuser layers, we use the best Fuser model for the final fuser layer (Table 24). For the intermediate layers, we use the top-8 Fuser models for each benchmark.

A.6 ARCHON ARCHITECTURES



Figure 9: **All-Source Generalizable ARCHON Architecture**: Using ARCHON's architecture search, we found this all-source ARCHON configuration to be effective across the benchmarks explored (except for CodeContests). In the diagram above, we use 10 SOTA all-source LLMs to create multiple successive layers of critic, ranker, and fusers, with each successive fuser layer having less fusers to produce a "funneling" effect as the candidate generations are processed. The layers of critic, ranker, and fuser led to better candidate generations through iterative critique and rewriting. Each of the initial Generator models were sampled once.



Figure 10: **All-Source ARCHON Architecture for Instruction-Following and Reasoning**: Using ARCHON's architecture search, we found this all-source ARCHON configuration to be effective across the instruction-following benchmarks explored (MT Bench, AlpacaEval 2.0, ArenaHardAuto).



Figure 11: **All-Source ARCHON Architecture for Math**: Using ARCHON's architecture search, we found this all-source ARCHON configuration to be effective across the math benchmarks explored (MATH).



Figure 12: **All-Source ARCHON Architecture for Coding**: Using ARCHON's architecture search, we found this all-source ARCHON configuration to be effective across the coding benchmarks explored (CodeContests).

				Datasets		
	Number of Inference Calls	MT Bench	Alpaca Eval 2.0	Arena Hard Auto	MixEval	MixEval Hard
	1	55.0%	44.7%	45.6%	86.5%	61.1%
s	10	52.5%	50.6%	45.6%	86.5%	63.9%
del B+	20	65.3%	60.4%	59.4%	89.0%	65.0%
10 J	30	69.2%	64.5%	69.0%	89.5%	67.5%
4	40	69.5%	66.7%	69.0%	89.5%	67.5%
	50	71.6%	66.7%	69.0%	89.5%	67.5%
	1	45.0%	57.5%	48.1%	88.9%	68.9%
- s	10	57.1%	63.2%	68.4%	90.0%	70.1%
del	20	59.4%	66.5%	75.5%	90.6%	70.5%
95 Q	30	70.2%	68.8%	77.4%	90.6%	72.9%
- 4	40	75.5%	68.8%	77.4%	90.6%	72.9%
	50	80.4%	68.8%	77.4%	90.6%	72.9%

A.7 ARCHON BY INFERENCE COMPUTE BUDGET, MODEL SIZE, AND COST

Table 26: **ARCHON with Different Inference Budgets**: For AlpacaEval 2.0, we use the length-controlled win rate (LC WR).

			Datasets		
Models / LLM Systems	MT Bench	Alpaca Eval 2.0	Arena Hard Auto	MixEval	MixEval Hard
Best 7B Model, 1-Sample	15.7%	41.0%	18.3%	76.2%	46.1%
Best 7B Model - 10-Sample + Ranking	16.5%	43.2%	18.9%	78.4%	48.5%
10-Model, 1-Sample Ensemble + Ranking	22.4%	48.2%	25.6%	81.5%	52.9%
10-Model, 1-Sample Ensemble + Fusion	14.3%	39.4%	17.5%	73.2%	45.2%
10-Model, 1-Sample Ensemble + Top-5 Ranking + Fusion	15.9%	41.2%	18.0%	75.1%	46.9%
10-Model, 1-Sample Ensemble + Critic + Fusion	10.5%	38.4%	16.5%	71.4%	42.5%

Table 27: **ARCHON with 7B Open-Source Models**: For AlpacaEval 2.0, we use the length-controlled win rate (LC WR). We use open-source 7B models for testing from Table 23.

Models	Cost (\$) per Million Input Tokens	Cost (\$) per Million Output Tokens
Claude 3.5 Sonnet	\$3	\$15
Claude 3.0 Opus	\$15	\$75
GPT-40	\$5	\$15
GPT-4-Turbo	\$10	\$30
TogetherAI - Llama 3.1 405B Instruct	\$5	\$5
TogetherAI - Llama 3.1 70B Instruct	\$0.88	\$0.88
TogetherAI - Other Models	\$0.90	\$0.90

Table 28: Model API Costs as of November 20)24	2	4	•))))	2	2	ľ))	J	J	l	l	(1))	2	2	2	į			ſ	i		2	e	(()	J	ł	ļ	l	1	1	ľ]	í	l	J	•	e	6	6	(7	ÿ	١	١	1))	J	C	((I				١				•	ĺ	ſ))	(,		5	S	\$	l	a	ć	;		5	S	S	1	t	t	t	1	1	5	5	S	S	Ş);)	0	C	(,	;)]]	2		[l	(()	•))	[ĺ		١	١	١	4	ŀ	ł	1		
---	-----	---	---	---	--	--	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	--	--	---	---	--	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	--	---	--	--	--	---	--	--	--	---	---	---	---	---	---	---	--	---	---	----	---	---	---	---	--	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	---	---	---	---	--	--	---	---	---	---	---	---	--	---	---	---	---	--	--	--	--	--	---	---	---	---	---	--	--	--	--	---	--	---	---	---	---	---	---	---	--	--

		Cost (\$) per Qu	uery for Benc	hmark			
Model / LLM System	MT Bench	AlpacaEval 2.0	Arena-Hard Auto	MixEval	MixEval Hard	MATH	Code Contests
Claude 3.5 Sonnet	0.0305	0.0171	0.0212	0.0231	0.0226	0.0325	0.384
GPT-4o	0.0481	0.0236	0.0324	0.0357	0.0361	0.514	0.562
Llama 3.1 405B Instruct	0.0281	0.0174	0.0185	0.0212	0.0205	0.305	0.372
General Purpose ARCHON Architecture	0.364	0.189	0.195	0.284	0.252	0.375	0.461
Task Specific ARCHON Architecture	0.401	0.210	0.221	0.295	0.265	0.425	0.448

Table 29: ARCHON Costs per Query by Benchmark

A.8 BAYESIAN OPTIMIZATION FOR ARCHON

A.8.1 ARCHON SEARCH SPACE AND OBJECTIVE

The ARCHON configuration space can be defined as $\mathcal{X} = \{x_g, x_s, x_f, x_r, x_c, x_v\}$ where:

- $x_q \in [1,10]$: Number of generator models
- $x_s \in [1,5]$: Samples per generator (extends to [1,1000] for CodeContests)
- $x_f \in [1,4]$: Number of fusion layers, including the final fusion layer at the end
- $x_r \in [2,10]$: Number of models per fusion layer, ranging from 2 to 10 increments of 2
- $x_c \in \{0,1\}$: Whether to use critic and ranker layers before each fuser
- $x_v \in \{0,1\}$: Whether to use verification layer before final fusion

The total search space initially contains 18,750 configurations $(10.5.5^{(4-1)}.3=18,750)$, reduced to 9,576 after removing invalid configurations where: 1) initial generations exceed fuser context window (24 candidates); and 2) single fuser layer contains multiple fusers ($x_f = 1$ while $x_r \ge 2$).

Let f(x) be the objective function evaluating an ARCHON configuration $x \in \mathcal{X}$, defined as:

$$f(x) = \operatorname{Performance}(x) - \lambda \cdot \operatorname{Cost}(x) \tag{1}$$

where Performance(x) is the accuracy on a 20% sample of target tasks and Cost(x) represents inference compute usage.

A.8.2 OPTIMIZATION PROCESS

We describe the components of the Bayesian optimization search approach for ARCHON.

First, define \mathcal{H} as a history of architectures and their resulting objective values, which we accumulate throughout the optimization process. We use Expected Improvement (EI) as the acquisition function for determining how to select the next architecture configuration to search:

$$\operatorname{EI}(x;\mathcal{H}) = \mathbb{E}[\max(0, f(x) - f(x^{+}))], \qquad (2)$$

where $f(x^+)$ is the best-observed value of our objective so far for $(x^+, f(x^+)) \in \mathcal{H}$. We use a Gaussian Process model as a surrogate model for approximating f(x).

We now describe the Bayesian optimization process. We first initialize the observation history $\mathcal{H} = \emptyset$. For timestep t = 1, ..., T, where T is the maximum number of iterations parameter, we do the following:

- 1. An architecture x_t is selected using the acquisition function, $x_t = \operatorname{argmax}_{x \in \mathcal{X}} \operatorname{EI}(x;\mathcal{H})$.
- 2. This architecture x_t is evaluated, and we obtain $f(x_t)$.
- 3. We use $f(x_t)$ to update our acquisition function and the surrogate model using $H \leftarrow H \cup (x_t, f(x_t))$.

The process continues until either:

- A maximum number of iterations is reached, T
- Performance convergence: $|f(x_{n+1}) f(x_n)| < \epsilon$
- Budget exhaustion: $Cost(x_1,...,x_n) > B$

For ARCHON's implementation, we initialize with 230-240 random configurations, as this was found to be optimal through empirical testing. Additional samples beyond this point provide diminishing returns and are better allocated to configuration search. For our implementation, we utilize the Bayesian Optimization python package for global optimization with Gaussian processes.

This formulation allows ARCHON to efficiently explore the configuration space, requiring 88.5% fewer evaluations than greedy search and 90.4% fewer than random search, with Bayesian optimization finding the best architectures in 96.0% of iterations. Traditional greedy search methods may perform comparably for limited inference budgets (<20 calls), but Bayesian optimization becomes increasingly effective as the search space and compute budget grow.

A.9 BAYES OPTIMIZATION VS. ALTERNATIVE APPROACHES

Search Techniques: Within the hyperparameter space, we explored three search algorithms for automating the development of inference-time architectures:

- 1. Random Search: Randomly selects a combination of hyperparameters for our ARCHON architecture.
- 2. **Greedy Search**: Starting with a base ARCHON configuration, marginally changes each hyperparameter and test if it improves performance or not. If it does, incorporate the change. If not, move on to the next hyperparameter.
- 3. **Bayesian Optimization**: Efficiently selects the most promising hyperparameter configurations for ARCHON by building a probabilistic surrogate model and leveraging an acquisition function for hyperparameter selection (Snoek et al., 2012; Nardi et al., 2019) (Section A.8).

To get our model ranking for the benchmark, we calculate the model ranking by testing each model individually on a 20% sample of each dataset benchmark in the first stage of the search. To get our fusion model ranking for the benchmark, we use the same approach, testing each model's fusion



Figure 13: **Impact of Different Optimization Algorithms on ARCHON's Architecture Search**: On the benchmarks MT Bench and Arena-Hard-Auto, we compare four approaches for finding the optimal inference-time architecture: random search, greedy search, and Bayes Optimization. Bayes Optimization finds the optimal architecture in 88.5% less iterations compared to greedy search and 90.4% less iterations compared to random search.

performance with an ensemble of 10 randomly selected models from the available set. From our experiments, we found that the best generator and fusion models could vary widely dataset to dataset, making it beneficial to perform these rankings for new datasets (Table 24). For search, we use the same 20% sample of each dataset that was used for evaluating generation and fusion, allowing us to guide architecture search with improved evaluation speed while getting meaningful development signal.

Comparing Search Algorithms: In Figure 13, we compare the effectiveness of each search algorithm on our explored benchmarks. While random search guarantees the optimal ARCHON configuration, we found Bayesian optimization to be most effective in terms of tradeoff between finding the optimal configurations and minimizing the number of configurations tested. For 96.0% percent of the search iterations tested in Figure 13, we found that Bayesian optimization had the optimal configuration amongst the four explored search algorithms. We use 230 initial samples for our Bayes Optimization architecture search (Section A.8). Bayesian optimization also found the best architecture configuration in 88.5% less evaluations than greedy search and 90.4% less evaluations than random search.

Bayesian Optimization Analysis: In Table 32, we explore how the number of initial testing points, the number of exploration iterations, and the ARCHON inference call budget impacts the effectiveness of Bayesian optimization. Additional initial testing points continue improving search efficacy up until 230-240 samples, where testing would be better delegated towards configuration search. For lower inference call budgets with ARCHON (e.g. <20 inference calls), Bayesian optimization proved less effective, performing more similarly to greedy search or random search given the limited search space (Table 33). Therefore, Bayesian optimization is more effective for more open-ended ARCHON architecture search with larger inference call budgets (e.g. >20 inference calls) whereas traditional component engineering might be better for more limited inference call budgets.

of Init. Points	% of Total Configs	Iter. till Max. Config.	Comb. Iter.	# of Init. Points	% of Total Configs	Iter. till Max. Config.	Con Ite
200	2.18%	353	553	200	2.18%	478	67
210	2.29%	324	534	210	2.29%	431	64
220	2.40%	301	521	220	2.40%	415	63
230	2.51%	284	514	230	2.51%	382	61
240	2.61%	261	501	240	2.61%	389	62
250	2.72%	265	515	250	2.72%	385	63
260	2.83%	256	516	260	2.83%	372	63
270	2.94%	252	522	270	2.94%	368	63

A.10 ARCHON ARCHITECTURE ALGORITHMS COMPARISONS

Table 30: MT Bench

Table 31: Arena-Hard-Auto

Table 32: **Bayesian Optimization Hyperparameter Comparisons**: On MT Bench and Arena-Hard-Auto, we compare Bayesian optimization configurations for the number of initial sample points. We find that 230 to 240 initial sample points minimizes the combined number of iterations (both initial sampling and exploring) to find the optimal configuration. For the configurations explored, the total number of hyperparameter choices is 9,576.

		Iteratio	ns to Cor	vergence	9
Inference Budget	10	20	30	40	50
Random Selection	387	1152	2731	4359	5843
Greedy Search	343	984	2153	3045	4895
Bayes Optimization	254	386	452	515	589

Table 33: **ARCHON Architecture Search Algorithms Comparison by Inference Call Budget**: For our comparison, we evaluate on MT Bench.