

UNCOVERING COMPETENCY GAPS IN LARGE LANGUAGE MODELS AND THEIR BENCHMARKS

Anonymous authors

Paper under double-blind review

ABSTRACT

The evaluation of large language models (LLMs) relies heavily on standardized benchmarks. These benchmarks provide useful aggregated metrics for a given capability, but those aggregated metrics can obscure (i) particular sub-areas where the LLMs are weak (“model gaps”) and (ii) imbalanced coverage in the benchmarks themselves (“benchmark gaps”). We propose a new method that uses sparse autoencoders (SAEs) to automatically uncover both types of gaps. By extracting SAE concept activations and computing saliency-weighted performance scores across benchmark data, the method grounds evaluation in the model’s internal representations and enables comparison across benchmarks. **As examples demonstrating our approach**, we applied the method to two popular open-source models and ten benchmarks. **We found that these models consistently underperformed on concepts that stand in contrast to sycophantic behaviors (e.g., politely refusing a request or asserting boundaries) and concepts connected to safety discussions. These model gaps align with observations previously surfaced in the literature; our automated, unsupervised method was able to recover them without manual supervision. We also observed benchmark gaps: many of the evaluated benchmarks over-represented concepts related to obedience, authority, or instruction-following, while missing core concepts that should fall within their intended scope. In sum, our method offers a representation-grounded approach to evaluation, enabling concept-level decomposition of benchmark scores. Rather than replacing conventional aggregated metrics, CG complements them by providing a concept-level decomposition that can reveal why a model scored as it did and how benchmarks could evolve to better reflect their intended scope.** Code is available at [anonymized](#).

1 INTRODUCTION

Evaluating large language models (LLMs) relies heavily on benchmarks that report aggregated scores (e.g., accuracy or pass@k). Over the last decade, hundreds of benchmarks have been introduced across diverse domains [Guo et al., 2023; Chang et al., 2024]. While these benchmarks have fueled progress, uniform aggregation can obscure important sub-trends and mask model weaknesses [Hardt, 2025; Burnell et al., 2023]. For instance, Didolkar et al. [2024] disaggregated performance on MATH [Hendrycks et al., 2021a] and found topic-wise scores ranging from 27% to 74%, despite an overall score of 54%.

To counteract these aggregation issues, some benchmarks provide “semantic” topic annotations (e.g., hand-curated topics in MATH [Hendrycks et al., 2021b] or GPQA [Rein et al., 2024], or embedding-based clusters [Perez et al., 2023]). These high-level labels help characterize benchmark distributions and disaggregate performance, but they are coarse-grained and offer limited insight into model strengths and weaknesses. In particular, we lack a view of how finer-grained concepts, contexts, and reasoning patterns extend beyond coarse topic labels and how they relate to real-world model usage and capabilities [Miller and Tang, 2025; Mizrahi et al., 2024]. Furthermore, many of these semantic annotations are manually curated and difficult to scale. Without a scalable, fine-grained understanding of benchmark distributions, we risk overlooking benchmark gaps and systematically overtesting certain concept types.

In this work, we are interested in the automated identification of two types of gaps: (i) benchmark gaps, i.e., concept domains that are inadequately represented in an evaluation dataset, and (ii) model gaps, i.e., concept domains where models systematically underperform (see Figure 1). To this end,

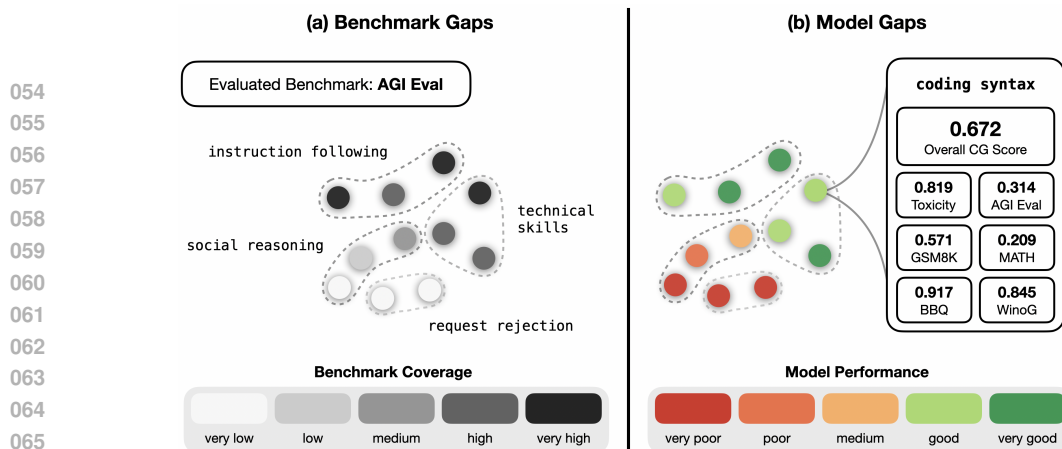


Figure 1: Competency Gaps Overview. Competency Gaps (CG) decomposes LLM evaluation into interpretable **benchmark gaps** and **model gaps** using the concept dictionary learned by a sparse autoencoder (SAE), a subset of which is visualized above. **(a) Benchmark Gaps** quantify how much each benchmark activates each concept, normalized by total dataset activation, and aggregated across benchmarks. **(b) Model Gaps** project model performance into concept space, assigning each concept per-benchmark and overall scores.

we introduce a new method called Competency Gaps (CG). Therein, we leverage sparse autoencoders (SAEs), which transform dense internal representations of a scrutinized LLM into high-dimensional, sparse feature vectors called SAE concept activations [Bricken et al., 2023; Cunningham et al., 2023]. Each dimension of these vectors is trained to capture a distinct, human-interpretable concept, assuming a sufficiently diverse and representative training distribution. We use pre-trained SAEs provided by the model authors, for which autointerpretability methods [McGrath et al., 2024] have already been applied to assign each dimension (or “concept”) a textual label. This fixed set of concepts, often referred to as a concept dictionary, is defined by a pre-set dictionary size hyperparameter. This allows us to analyze the benchmark distribution over the SAE concept space¹ and identify benchmark gaps, and consequently map the model’s performance onto the concept space to identify model gaps.

As a demonstration of the method, we evaluate Competency Gaps (CG) on two popular open-source models (Gemma2-2B-Instruct [Team et al., 2024] and Llama3.1-8B-Instruct [Grattafiori et al., 2024b]) across ten diverse benchmarks. We find notable gaps: benchmarks often miss concepts central to their intended scope while overtesting concepts tied to authority, control, and instruction-following. Additionally, we demonstrate how standard aggregated benchmark performance metrics tend to overly reflect the overrepresented concepts. Overall, this analysis serves to demonstrate the utility of CG as a method for uncovering and addressing gaps in prevailing evaluation paradigms. We summarize our contributions as follows:

- **Competency Gaps (CG) Method.** We introduce a systematic method for the automated identification of benchmark coverage and model performance gaps, using an SAE-based approach. The method can be applied to any LLM and benchmark of interest.² The method helps benchmark developers iteratively improve coverage, and enables model developers to gain fine-grained, distributional insights into model performance, as described in Figure 2.
- **Example Applications of CG on Popular LLMs and Benchmarks.** We applied CG to (Gemma2-2B-Instruct and Llama3.1-8B-Instruct) over ten diverse benchmarks, illustrating the kinds of insights the method enables: CG decomposed singular benchmark scores into interpretable axes derived from the models’ own representations, surfaced consistent patterns of over- and under-tested concepts, and identified actionable improvements for both models and benchmarks.
- **Interactive Exploration Tool.** We release an open-source web application (Figure 6), alongside the code for the Competency Gaps method, for exploring per-concept model behavior and benchmark coverage. This enables users to audit model capabilities and benchmark balance in a transparent and interpretable way.

¹Note that the proposed approach can only detect benchmark and model gaps for concepts that are represented in the SAE space. With more representative SAEs, concept coverage of our approach improves as well.

²While training an SAE specifically for the model at hand allows the analysis to be grounded in that model’s own representations, we demonstrate that even LLMs without a dedicated SAE can be evaluated using CG by leveraging the SAE from another model. See Section 4.2.3.

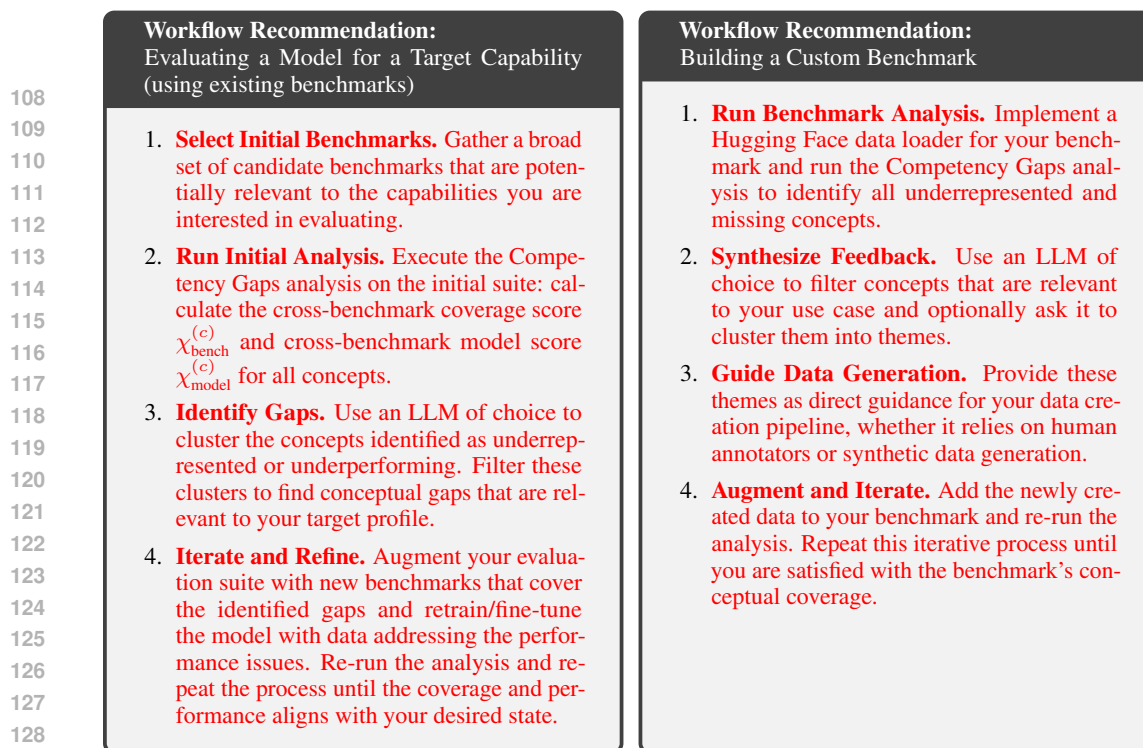


Figure 2: Recommended workflows for applying the Competency Gaps (CG) method in production.

2 RELATED WORK

Weakness Identification in LLMs. Identification of LLM weaknesses has evolved from anecdotal analysis to more systematic frameworks that break down performance into specific components [Jones and Steinhardt, 2022]. Among the first to do this were HarmBench [Mazeika et al., 2024] and garak [Derczynski et al., 2024], which established standardized evaluations of harmful behaviors. Various methods have introduced autoraters for this task: AutoDetect [Cheng et al., 2024], for example, used three LLM-powered agents to achieve a 50%+ weakness identification success rate. Gan et al. [2024] demonstrated that reasoning capabilities in LLMs can be systematically broken down and analyzed to identify specific weaknesses in logical inference chains. Other systematic evaluation methods have been proposed [Kim et al., 2023]; some work in this domain has also taken an adversarial learning approach [Yang et al., 2024].

Cross-Model Behavior Comparison. Beyond simple benchmark scoreboards, recent work has explored different approaches of characterizing the behavioral differences and internal representations across models and architectures in more detail [Zheng et al., 2025; Kim et al., 2025; Chang and Bergen, 2024]. Some works, including BehaviorBox [Tjuiatja and Neubig, 2025] and the LLM Comparator [Kahng et al., 2024], have emphasized side-by-side comparisons and actionable insights. Some mechanistic interpretability methods have emerged as well. One such method patches activations between model locations and decoding representations into interpretable text [Nada et al., 2024]. Recent work has also scrutinized the “universality hypothesis”, i.e., universal features and circuits should appear for similar tasks across architectures [Shu et al., 2025; Yin et al., 2024].

Benchmark Analysis and Down-Sampling. Automated quality detection has advanced through frameworks like CLEAR [Chen and Mueller, 2024] and SMART Filtering [Gupta et al., 2024], automatically detecting and filtering problematic training data. A significant part of the work in this area has specifically focused on bias evaluation, a subset of quality assessment [Doan et al., 2024; Manerba et al., 2023; Koo et al., 2023]. Some preliminary work has also considered monitoring benchmark performance across time [Zhong and Raghunathan, 2025] or using meta learning [Calian et al., 2025].

Comparison to Related Work. In Appendix G, we compare our method to some of the above-mentioned methods that relate to the discovery of benchmark and model gaps. See Appendix G.1 for details of our comparison methodology, Appendices G.2 and G.3 for a high-level comparison of approaches and features of each method, and Appendices H, I, and J contain a comparison of results from each method.

3 COMPETENCY GAPS (CG) METHOD

In this section, we introduce **Competency Gaps (CG)**, an automated, SAE-based method that can be used to systematically evaluate and identify:

- **Benchmark gaps** given a set of evaluation benchmarks \mathcal{B} . How good is the coverage of various concepts in a specific or set of benchmarks? Do the benchmarks have any concept coverage gaps?
- **Model gaps** given a language model M with an associated sparse autoencoder SAE . How well does the model perform across various concepts? What are its strengths and weaknesses?

Before diving into the specifics of the method, we first establish some notation. Each benchmark $b \in \mathcal{B}$ comes with an underlying dataset D_b . Furthermore, each concept $c \in C_{SAE}$ represents a distinct direction in the SAE space, and can be mapped to an autointerpretability label to which we refer as l_c . Based on SAE concepts, we introduce an SAE activation score $s_{c,i}$ that quantifies the degree of which concept c was activated within a token sequence x_i .

SAE Concept Activation Score $s_{c,i}$. A data point i from a benchmark consists of an input token sequence x_i . The concept c 's activation on token $x_{i,j} \in x_i$ is computed via $SAE(x_{i,j}, c)$. We sum $SAE(x_{i,j}, c)$ over all tokens in x_i to obtain the **SAE activation score** $s_{c,i}$ for concept c on data point i . This is normalized by the length of the token sequence as $\tilde{s}_{c,i}$:

$$s_{c,i} = \sum_{x_{i,j}} SAE(x_{i,j}, c), \quad \tilde{s}_{c,i} = \frac{s_{c,i}}{|x_i|} \quad (1)$$

Overall, our metrics were devised to satisfy: (1) All included benchmarks have an equal weight for computing the concept's cross-benchmark CG score, regardless of their size. (2) All data points have an equal weight for computing per-benchmark metrics, regardless of their token length.

3.1 BENCHMARK GAPS

Based on the introduced SAE concept activation score, we now introduce a measure to quantify concept coverage within and across benchmarks. These measures enable a distributional characterization of benchmarks, and consequently the identification of benchmark gaps.

Coverage Within an Individual Benchmark. To evaluate the coverage of a concept c in benchmark b , we define:

$$\chi_{\text{bench}}^{(b,c)} = \frac{\sum_{i \in D_b} \tilde{s}_{c,i}}{\frac{1}{|C_{SAE}|} \sum_{c' \in C_{SAE}} \sum_{i \in D_b} \tilde{s}_{c',i}} \quad (2)$$

which relates the ratio of the activation of SAE concept c in dataset b to the average concept activation in b .

Cross-Benchmark Coverage. The overall coverage for concept c in a benchmark suite \mathcal{B} is the mean $\chi_{\text{bench}}^{(b,c)}$ across \mathcal{B}_c (all benchmarks where c is activated):

$$\mathbf{X}_{\text{bench}}^{(c)} = \frac{1}{|\mathcal{B}|} \sum_{b \in \mathcal{B}_c} \chi_{\text{bench}}^{(b,c)} \quad (3)$$

Coverage Classification. We label a concept as *missing* from the benchmark suite if $\mathbf{X}_{\text{bench}}^{(c)} < \epsilon$, for some small ϵ (we used e^{-5}). Among the remaining concepts, those in the bottom decile (at or below the empirical 10th percentile) are *underrepresented*, and those in the top decile (at or above the empirical 90th percentile) are *overrepresented*. We can similarly define missing, underrepresented, and overrepresented concepts for each individual benchmark, using $\chi_{\text{bench}}^{(b,c)}$ instead.

3.2 MODEL GAPS

We now turn to the quantification of model gaps, and introduce a novel measure for model performance grounded in the discovered SAE concept space, both for individual and across benchmarks.

Per-Benchmark Model Performance. To evaluate model performance on a concept c for a benchmark b , we define $\chi_{\text{model}}^{(b,c)}$:

$$\chi_{\text{model}}^{(b,c)} = \frac{\sum_{i \in D_b} m_b(i) \cdot \tilde{s}_{c,i}}{\sum_{i \in D_b} \tilde{s}_{c,i}} \quad (4)$$

where $m_b(i) \in [0, 1]$ is the performance scoring policy for benchmark b for datapoint $i \in D_b$ (normalized to $[0,1]$; higher is better). If no data points activate c in b , then $\chi_{\text{model}}^{(b,c)}$ is undefined.

Cross-Benchmark Model Performance. The overall performance for concept c in a benchmark suite \mathcal{B} is the mean $\chi_{\text{model}}^{(b,c)}$ across all benchmarks:

$$\mathbf{X}_{\text{model}}^{(c)} = \frac{1}{|\mathcal{B}_c|} \sum_{b \in \mathcal{B}_c} \chi_{\text{model}}^{(b,c)} \quad (5)$$

We only consider benchmarks where $\chi_{\text{model}}^{(b,c)}$ is defined, i.e. $\mathcal{B}_c = \{b \in \mathcal{B} \mid \sum_{i \in D_b} s_{c,i} > 0\}$. If no data points in \mathcal{B} activate concept c (i.e., $|\mathcal{B}_c| = 0$), then $\mathbf{X}_{\text{model}}^{(c)}$ is undefined.

Given the cross-benchmark model performance score, we now introduce the notion of a model gap that pinpoints concepts associated with low model performance.

Model Gap. We label a concept a *model gap* if $\mathbf{X}_{\text{model}}^{(c)} < \epsilon$, for the same small ϵ .

4 DEMONSTRATIONS OF THE METHOD

4.1 EXPERIMENTAL SETUP

Benchmarks. We demonstrate our method on ten static benchmark datasets, regularly used for performance and safety evaluations, and one *arena-style benchmark*. However, the presented method can be applied to any text-based benchmark.

- *Factuality*: Natural Questions [Kwiatkowski et al., 2019]; Vectara [Meyman, 2025].
- *Math*: GSM8K [Cobbe et al., 2021]; MATH [Hendrycks et al., 2021b].
- *Reasoning*: AGI Eval [Zhong et al., 2023]; LogicBench [Parmar et al., 2024]; SocialIQA [Sap et al., 2019]; WinoGrande [Sakaguchi et al., 2021].
- *Ethics & Bias*: BBQ [Parrish et al., 2021]; CrowS-Pairs [Nangia et al., 2020].
- *Arena Style*: LMSYS Chatbot Arena Zheng et al. [2023].

Models. We analyzed two popular open-source LLMs with available SAEs and autointerpretability labels. However, we would like to emphasize that the method is not bound to these particular models; it can be applied to any standard LLM with an SAE.

- Llama3.1-8B-Instruct + Goodfire SAE attached at layer 19 [Grattafiori et al., 2024a; Balsam et al., 2025].
- Gemma2-2B-Instruct + Gemma Scope SAE attached at layer 20 (residual stream) [GemmaTeam et al., 2024; Lieberum et al., 2024];³

³For the Goodfire Llama SAE, the choice of layer for SAE attachment was made by the creators of the SAE. We chose the Gemma Scope layer to be at a comparable depth. Generally, deeper layers have a tendency to represent higher-level, sentence- or discourse-level abstractions [Balcells et al., 2024; Shi et al., 2025].

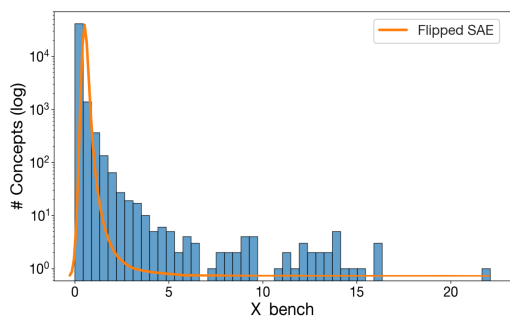


Figure 3: Cross-Benchmark Coverage. The distribution of $X_{\text{bench}}^{(c)}$ scores obtained for the 10 evaluated benchmarks, using the SAE of Llama 3.1 8B. This distribution exhibits strong left skew (most concepts have low coverage), and avg. performance is strongly dominated by a few concepts with high coverage (high $X_{\text{bench}}^{(c)}$). Similar skewed distributions were observed for individual benchmarks (Appendix Figure 12). The orange curve shows a similar analysis conducted through the activations and SAE concepts of Gemma 2 2B (see Section 4.2.3).

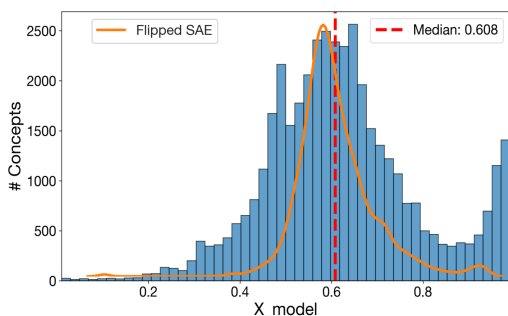


Figure 4: Cross-Benchmark Model Performance. The distribution of $X_{\text{model}}^{(c)}$ scores obtained for Llama 3.1 8B across the 10 evaluated benchmarks. The model exhibited high variance in performance across concepts. We observe particularly high performance for some concepts; these tend to include concepts related to coding, data handling, instruction following, and expressing positive sentiment toward the user. The orange curve shows a similar analysis conducted through the activations and SAE concepts of Gemma 2 2B (see Section 4.2.3).

4.2 RESULTS

To demonstrate the kinds of insights enabled by our method and exploratory web application, we report results of applying CG to Llama3.1-8B-Instruct across ten static benchmark datasets. This analysis can be applied to any language model and is intended as an active, iterative process in which both the benchmarking suite and the model are continuously refined (see Figure 2).

We report analogous results for Gemma2-2B-Instruct in Appendix B, and we further demonstrate how CG can be used with arena-style benchmarks by presenting results on the LMSYS Chatbot Arena in Appendix D. To efficiently sift through the large number of results that can arise from our method (e.g., given the scale of the SAE concept dictionaries), we sometimes use another LLM (Gemini 2.5 Flash, in this case) to filter, group, and summarize sets of concept descriptions.

4.2.1 BENCHMARK GAPS

Existing Benchmarks Exhibit Skewed Representation Across Concepts. The cross-benchmark coverage distribution across all ten benchmarks is shown in Figure 3. The distribution exhibits a strong left skew, resulting in both over- and under-representation of concepts in our (somewhat typical) suite of benchmarks. This skew results in any standard mean-based summary statistics being dominated by the outliers on the right part of the tail – a small number of “top concepts” with high representation (high $X_{\text{bench}}^{(c)}$). Conversely, when a concept has consistently *low* coverage across all included benchmarks, it risks being systematically under-tested. Another undesirable pattern is substantial overlap across benchmarks in a single evaluation suite, further illustrated in Figure 14 (Appendix C).

The top concepts mainly relate to starting new conversations and sports news, with a particular focus on football and sports achievements. Other prominent themes include syntax and attributes of articles. For example, the top 10 concepts by benchmark coverage included:

- (56130) “English Premier League football discussions, especially about Manchester teams”;
- (41290) “New conversation or topic segment boundary marker”.

Among the concepts with the *lowest* coverage score, one finds many concepts related to meta-cognition about the AI itself, e.g. its instructions, roleplaying boundaries, and how it discusses user inputs. The bottom 10 concepts by coverage include:

- (53553) “The assistant should maintain professional boundaries when asked to roleplay”
- (25352) “References to user messages or inputs in meta-discussion”.

We identified 314 concepts (1%) as entirely missing from this particular suite of benchmarks (see Table 1). These again include concepts related to the AI’s meta-cognition, as well as legal concepts.

Individual Benchmarks Miss Relevant Concepts. Individual benchmarks show a similar skew in their representation (Appendix Figure 12), and every benchmark except Vectara misses at least

Concept ID	Concept Description
(2501)	The assistant explaining why it needs more information
(2641)	The assistant needs to explain its limitations or capabilities
(2009)	Regulatory classification and compliance requirements

Table 1: Examples of Missing Concepts from the Full Benchmark Suite.

30% of all concepts (Appendix Figure 13). However, benchmarks are usually designed to evaluate a specific subset of capabilities, and so of course it may not be desirable for individual benchmarks to have complete coverage of all concepts.

More importantly, we would like benchmarks to have coverage of *relevant* concepts. To identify such missing relevant concepts, we first used the CG method to identify all missing concepts for a given benchmark. Then, we used an LLM (Gemini 2.5 Flash, in this case) to find those concepts that one might expect to be in scope for the benchmark (see Appendix ?? for the prompt). We also used the open-sourced web app to explore and verify these examples. This process exemplifies our recommendation to use the CG method for unsupervised discovery of coverage gaps.

Table 2 presents illustrative examples of concepts that are missed by benchmarks. These are concepts that seem central to the desired goals of the benchmarks, e.g. “The need for thorough and objective assessment of evidence” was missing from *AGIEval*, and “Instructions about how someone should behave or what qualities to embody” was missing from *SocialIQA*.

Benchmark	Concept ID	Concept Description
<i>AGIEval</i>	(33456)	The need for thorough and objective assessment of evidence
	(59559)	Careful qualification and nuanced explanation of complex topics
<i>LogicBench</i>	(56997)	The model is explaining how different elements or factors relate to each other
	(11957)	Mathematical and logical concepts across multiple languages
<i>SocialIQA</i>	(35877)	Speaker defending or explaining their planned actions against expectations
	(1897)	Instructions about how someone should behave or what qualities to embody

Table 2: Examples of Missing Concepts from Three Individual Benchmarks.

4.2.2 MODEL GAPS

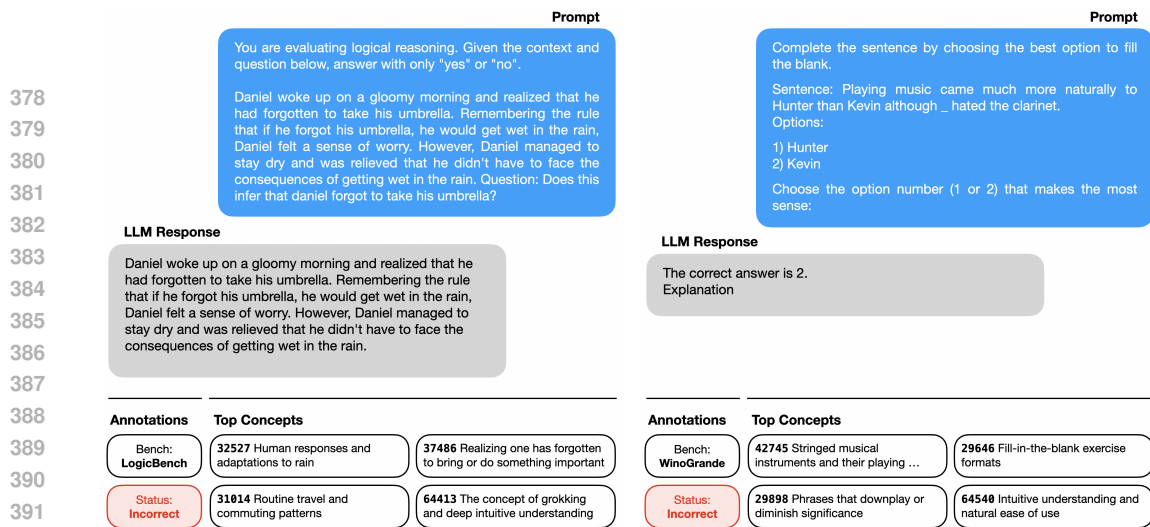
In addition to benchmark gaps, we also analyze model performance gaps – on which SAE concepts do the models perform particularly well or particularly poor?

Model’s Best-Performing Concepts Include Commitment to Help, Coding. Because the SAE concepts are fixed for a given model, our method enables us to compare and create composite results across benchmarks. Figure 4 shows the distribution across concepts for cross-benchmark model performance, which shows that the model performs well on a number of concepts, with high performance across all benchmarks. Concepts with the highest cross-benchmark performance tended to concern engineering related tasks (i.e., coding or data handling) or helpful behaviors (positive sentiments towards the user, or delivering an accurate response). The top 10 concepts included:

- (20022) “Iteration or traversal through sequences in programming”,
- (24074) “The assistant is about to provide an illustrative example”,
- (2461) “Assistant expressing commitment to help or do its best”.

Model’s Worst-Performing Concepts Include Polite Rejection, Time, and Setting Boundaries. Perhaps more critical, from the standpoint of model evaluation, are the concepts which attained the poorest overall performance across benchmarks. An interesting recurring theme is that the worst-performing concepts include opposites to the helpful/sycophantic concepts discussed in the previous section (which came at the top in performance). Examples of worst-performing concepts include:

- (26535) “The assistant needs to politely reject or redirect inappropriate requests”,
- (56928) “The assistant maintaining professional boundaries while offering appropriate help”.



393 **Figure 5: Model Gap Illustrated on Specific Benchmark Datapoints.** Example LogicBench and
394 WinoGrande items associated with “intuitive understanding” concepts (left: 64113, right: 64540).
395 Llama 3.1 8B answered both incorrectly, consistent with these concepts being model gaps.

396 Furthermore, Competency Gaps identified other groups of competencies that have been anecdotally
397 identified as LLM weaknesses in prior literature, validating our automated approach as a scalable and
398 systematic method for identifying such model weaknesses. These include:

- 400** • **Representations of Time:**
401 (29324) “Historical date and time period formatting”
402 (12644) “Cooking time durations in recipe instructions”
- 403** • **Image Manipulations:**
404 (30206) “Image contrast adjustments in photo editing and computer vision”⁴
- 405** • **Palindromes / Reasoning over Letters:**
406 (56613) “Code examples and explanations of palindrome checking algorithms”
- 407** • **Mathematical Operations:**
408 (64527) “Mathematical addition operator in calculations”.

409 Additionally, the proposed method surfaces LLM weaknesses that have *not* been previously studied
410 in the literature. One such category is “appeals to intuition in reasoning or decision making”:

- 411** (64413) “The concept of grokking and deep intuitive understanding”
- 412** (64540) “Intuitive understanding and natural ease of use”

413 Moreover, the exploratory web application allows users to directly examine example data points and
414 better understand how such competencies manifest in practice (see Figure 5).

415 4.2.3 ROBUSTNESS

416 Prior work has provided mixed evidence on the stability of SAE concepts [??]. We therefore set out
417 to evaluate the stability and generalizability of CG findings.

418 **A Model-Specific SAE Is Not Necessary, and Different SAEs Can Yield Similar Results.** We
419 tested whether the CG insights derived from Llama 3.1 8B using its own model-specific SAE align
420 with those obtained through Gemma 2 2B activations and SAE. As shown in Figures 3 and 4, the
421 overall shape and medians of the score distributions were similar, especially considering the difference
422 in dictionary sizes (Llama’s SAE contains 2.8x more concepts than Gemma’s). When comparing the
423 best- and worst-performing concepts, we observed clear correspondences and similar interpretive
424 results.

425 ⁴Llama3.1-8B-Instruct is a text-only model, and we found that such concepts were activated by metadata
426 indicating image editing, Photoshop scripts, and tutorials on the topic. This type of review of specific data
427 points—both in the benchmarks and in the SAE training dataset—is enabled by our web app, and proves helpful
428 for understanding the context and nuance of various concepts.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

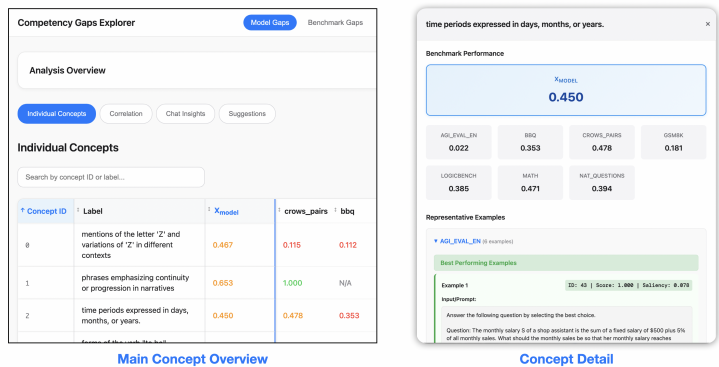


Figure 6: Exploratory Web Application Overview. The interface presents a searchable and filterable list of all concepts in the *Main Concept Overview*, with an expandable *Concept Detail Modal* that provides additional per-benchmark information, including specific example generations. The application also includes dedicated sections for cross-benchmark correlation, benchmark coverage inspection, and supplementary analyses. Additional screenshots are provided in Appendix F.

insights, as shown in Table 3. This suggests that CG can yield meaningful insights even for LLMs *without* their own pre-trained SAE and demonstrates the overall stability of the method. Nonetheless, we expect that a model-specific SAE with a larger dictionary still offers the most precise and grounded results.

Analysis	Llama SAE	Gemma SAE
<i>Best Performance</i>	(45314) Legal reasoning and argumentation patterns in multiple choice questions (27510) Code patterns for including JavaScript resources in web pages	(5471) References to legal cases and procedural aspects of law (8196) Code constructs or reserved keywords in programming languages
<i>Worst Performance</i>	(2874) Mathematical differentiation operators and notation (2872) Explaining time requirements and duration	(13908) Numerical values, counts or measurements (9936) Dates and numeric sequences
<i>Best Coverage</i>	(41290) New conversation or topic segment boundary marker (902) Step-by-step mathematical explanations and calculations	(11527) The start of a document (11880) Mathematical expressions and calculations related to derivatives and factors
<i>Worst Coverage</i>	(27900) Discussions of factual accuracy and consistency checking (47946) The assistant explains how it processes and handles information	(5657) Terms related to correctness and accuracy in responses or answers (1797) Phrases related to instructions or operational processes

Table 3: Examples of Concept Correspondences in Best/Worst Performance and Coverage Analyses. The analysis was conducted for Llama 3.1 8B: the Llama SAE results used the model’s own activation and custom SAE; the Gemma SAE results used Gemma 2 2B’s activations and SAE.

CG Scores Are Consistent Across Perturbations To assess the robustness of CG scores, we re-ran the full analysis 100 times, each time randomly dropping 20% of the examples per benchmark. The resulting standard deviations were low: 0.014 for X_{model} and 0.025 for X_{bench} on Llama. This indicates that CG yields stable scores under random subsampling.

CG Scores Respond to Adversarial Perturbations. We conducted an adversarial ablation in which we identified the top 100 best- and worst-performing concepts, then removed the most salient 100 datapoints associated with them across all benchmarks. As expected, removing rows aligned with high-performing concepts lowered median X_{model} on average by 0.6%, while removing those aligned with low-performing concepts increased it on average by 1.3% (repeated across 10 repetitions). Despite removing less than 1% of the testing data, we were able to make predictable and consistent changes to the overall performance, suggesting that CG surfaces meaningful concept-level information.

5 DISCUSSION

Our analysis revealed a structural imbalance in a sample of popular benchmarks. For example, the benchmarks strongly emphasized concepts related to authority, control, and instruction-following, while neglecting complementary concepts related to polite refusals, meta-cognition, and meta-discussion about the AI itself. When evaluating across a benchmark suite, such skewed representation within the benchmarks may skew our perception of model capabilities. Our method also identified potential coverage gaps in specific benchmarks, pinpointing missing concepts that seemed relevant to each benchmark’s scope.

A similar bias toward sycophancy and instruction-following emerged in the model gap analysis. Here, positive or sycophantic concepts score highest, while opposing concepts (such as those linked to rejecting requests or setting boundaries) score lowest. While this is likely due in part to instruction-based post-training, we note that model gaps and benchmark gaps are heavily intertwined, and that benchmark gaps may lead to model gaps — model developers may unknowingly overlook (or be disincentivized to address) model weaknesses that are poorly covered by existing benchmarks.

5.1 LIMITATIONS

Concept Coverage is Limited to SAE Concept Space. The usage of SAE concepts implies that we can only detect competency gaps for concepts that have *some* representation in the model, as captured by the SAE dictionary. Thus, we should take care in how we interpret our results: (1) *benchmark gap* analyses reveal where models possess internal representations for concepts that evaluations fail to adequately test, and (2) *model gap* analyses identify cases where models have flawed or partial internal representations that they cannot apply effectively to downstream tasks. However, it should be noted that with more representative SAEs, the concept coverage of our proposed approach improves as well.

Concept Labels. In addition to being limited to existing SAE concepts, as mentioned above, our method inherits some limitations of SAEs. These include a lack of ground truth [Smith et al. [2025]] (relatedly, the autointerpretability labels are automatically generated, though our web app enables easy spot checks as needed), non-convergence when SAEs are retrained [Paulo and Belrose, 2025], “feature absorption” (where token-aligned latents “absorb” some expected feature directions, and the tendency for SAEs toward learning common feature combinations instead of atomic features) [Chanin et al. [2024]; O’Neill et al. [2025]]. Despite these limitations, SAE remain a well regarded mechanistic interpretability method for unsupervised hypothesis discovery, which is how we use them in the current work.

5.2 DOWNSTREAM APPLICATIONS AND FUTURE WORK

Benchmark Search and Selection. Our method could be integrated into a database of available benchmarks (e.g., Hugging Face Datasets). Users seeking benchmarks to evaluate their models could use CG to inform their selection to achieve a desired coverage. Similarly, creators of benchmarks or the database maintainers could use CG to inform future benchmark creation.

Targeted Creation of Novel Benchmark Data. Beyond characterizing existing benchmarks, a list of underrepresented concepts within the suite could guide the generation of novel benchmark data. Such data could be created either by prompting LLMS with autointerpretability labels or by directly applying SAE-based steering interventions during generation.

Method Improvements. We hope that this work may serve as a starting point for further development on the method. For example, in future we may wish to incorporate an automated sensitivity analysis for the choice of SAE layer, or per-token activations for more fine-grained analyses.

Current model evaluations may risk a narrower view of “competence,” potentially leaving critical gaps untested. We hope that our method enables benchmark developers and model evaluators to identify both benchmark gaps and model gaps, uncovering and addressing model weaknesses in areas that may be essential for real-world, human-facing use-cases.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

REPRODUCIBILITY STATEMENT

To maximize reproducibility of our work, both the analysis and visualization code is open-sourced at anonymized. This repository also includes the data extracted for the models and datasets evaluated in this paper. Furthermore, all prompts used for LLM clustering are included in Appendix E and the setup of alternate methods is described in Appendix G.1.

REFERENCES

- 594
595
596 D. Balcells, B. Lerner, M. Oesterle, E. Ucar, and S. Heimersheim. Evolution of SAE features across
597 layers in llms. *arXiv preprint arXiv:2410.08869*, 2024.
- 598
599 D. Balsam, T. McGrath, L. Gorton, N. Nguyen, M. Deng, E. Ho, and . D. N. D. yet. Announcing
600 Open-Source SAEs for Llama 3.3 70B and Llama 3.1 8B, 2025. URL <https://www.goodfire.ai/blog/sae-open-source-announcement>.
601
- 602 T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Denison,
603 A. Askell, et al. Towards monosemanticity: Decomposing language models with dictionary
604 learning. *Transformer Circuits Thread*, 2, 2023.
- 605
606 R. Burnell, W. Schellaert, J. Burden, T. D. Ullman, F. Martinez-Plumed, J. B. Tenenbaum, D. Rutar,
607 L. G. Cheke, J. Sohl-Dickstein, M. Mitchell, et al. Rethink reporting of evaluation results in ai.
608 *Science*, 380(6641):136–138, 2023.
- 609
610 D. A. Calian, G. Farquhar, I. Kemaev, L. M. Zintgraf, M. Hessel, J. Shar, J. Oh, A. György, T. Schaul,
611 J. Dean, et al. DataRater: Meta-learned dataset curation. *arXiv preprint arXiv:2505.17895*, 2025.
- 612
613 T. A. Chang and B. K. Bergen. Language model behavior: A comprehensive survey. *Computational
614 Linguistics*, 50(1):293–350, 2024.
- 615
616 Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, et al.
617 A survey on evaluation of large language models. *ACM transactions on intelligent systems and
618 technology*, 15(3):1–45, 2024.
- 619
620 D. Chanin, J. Wilken-Smith, T. Dulka, H. Bhatnagar, S. Golechha, and J. Bloom. A is for absorption:
621 Studying feature splitting and absorption in sparse autoencoders. *arXiv preprint arXiv:2409.14507*,
622 2024.
- 623
624 J. Chen and J. Mueller. Automated data curation for robust language model fine-tuning. *arXiv
625 preprint arXiv:2403.12776*, 2024.
- 626
627 J. Cheng, Y. Lu, X. Gu, P. Ke, X. Liu, Y. Dong, H. Wang, J. Tang, and M. Huang. Autodetect:
628 Towards a unified framework for automated weakness detection in large language models. *arXiv
629 preprint arXiv:2406.16714*, 2024.
- 630
631 K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton,
632 R. Nakano, C. Hesse, and J. Schulman. Training verifiers to solve math word problems. *arXiv
633 preprint arXiv:2110.14168*, 2021.
- 634
635 H. Cunningham, A. Ewart, L. Riggs, R. Huben, and L. Sharkey. Sparse autoencoders find highly
636 interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023.
- 637
638 L. Derczynski, E. Galinkin, J. Martin, S. Majumdar, and N. Inie. garak: A framework for security
639 probing large language models. *arXiv preprint arXiv:2406.11036*, 2024.
- 640
641 A. Didolkar, A. Goyal, N. R. Ke, S. Guo, M. Valko, T. Lillicrap, D. Jimenez Rezende, Y. Bengio,
642 M. C. Mozer, and S. Arora. Metacognitive capabilities of llms: An exploration in mathematical
643 problem solving. *Advances in Neural Information Processing Systems*, 37:19783–19812, 2024.
- 644
645 T. V. Doan, Z. Chu, Z. Wang, and W. Zhang. Fairness definitions in language models explained.
646 *arXiv preprint arXiv:2407.18454*, 2024.
- 647
648 E. Gan, Y. Zhao, L. Cheng, Y. Mao, A. Goyal, K. Kawaguchi, M.-Y. Kan, and M. Shieh. Reasoning
649 robustness of llms to adversarial typographical errors. *arXiv preprint arXiv:2411.05345*, 2024.
- 650
651 GemmaTeam, M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard,
652 B. Shahriari, A. Ramé, J. Ferret, P. Liu, P. Tafti, A. Friesen, M. Casbon, S. Ramos, R. Kumar,
653 C. L. Lan, S. Jerome, A. Tsitsulin, N. Vieillard, P. Stanczyk, S. Girgin, N. Momchev, M. Hoffman,
654 S. Thakoor, J.-B. Grill, B. Neyshabur, O. Bachem, A. Walton, A. Severyn, A. Parrish, A. Ahmad,
655 A. Hutchison, A. Abdagic, A. Carl, A. Shen, A. Brock, A. Coenen, A. Laforge, A. Paterson,

- 648 B. Bastian, B. Piot, B. Wu, B. Royal, C. Chen, C. Kumar, C. Perry, C. Welty, C. A. Choquette-
649 Choo, D. Sinopalnikov, D. Weinberger, D. Vijaykumar, D. Rogozińska, D. Herbison, E. Bandy,
650 E. Wang, E. Noland, E. Moreira, E. Senter, E. Eltyshev, F. Visin, G. Rasskin, G. Wei, G. Cameron,
651 G. Martins, H. Hashemi, H. Klimczak-Plucińska, H. Batra, H. Dhand, I. Nardini, J. Mein, J. Zhou,
652 J. Svensson, J. Stanway, J. Chan, J. P. Zhou, J. Carrasqueira, J. Iljazi, J. Becker, J. Fernandez,
653 J. v. Amersfoort, J. Gordon, J. Lipschultz, J. Newlan, J.-y. Ji, K. Mohamed, K. Badola, K. Black,
654 K. Millican, K. McDonell, K. Nguyen, K. Sodhia, K. Greene, L. L. Sjoesund, L. Usui, L. Sifre,
655 L. Heuermann, L. Lago, L. McNealus, L. B. Soares, L. Kilpatrick, L. Dixon, L. Martins, M. Reid,
656 M. Singh, M. Iverson, M. Görner, M. Velloso, M. Wirth, M. Davidow, M. Miller, M. Rahtz,
657 M. Watson, M. Risdal, M. Kazemi, M. Moynihan, M. Zhang, M. Kahng, M. Park, M. Rahman,
658 M. Khatwani, N. Dao, N. Bardoliwalla, N. Devanathan, N. Dumai, N. Chauhan, O. Wahltinez,
659 P. Botarda, P. Barnes, P. Barham, P. Michel, P. Jin, P. Georgiev, P. Culliton, P. Kuppala, R. Co-
660 manescu, R. Merhej, R. Jana, R. A. Rokni, R. Agarwal, R. Mullins, S. Saadat, S. M. Carthy,
661 S. Cogan, S. Perrin, S. M. R. Arnold, S. Krause, S. Dai, S. Garg, S. Sheth, S. Ronstrom, S. Chan,
662 T. Jordan, T. Yu, T. Eccles, T. Hennigan, T. Kocisky, T. Doshi, V. Jain, V. Yadav, V. Meshram,
663 V. Dharmadhikari, W. Barkley, W. Wei, W. Ye, W. Han, W. Kwon, X. Xu, Z. Shen, Z. Gong,
664 Z. Wei, V. Cotruta, P. Kirk, A. Rao, M. Giang, L. Peran, T. Warkentin, E. Collins, J. Barral,
665 Z. Ghahramani, R. Hadsell, D. Sculley, J. Banks, A. Dragan, S. Petrov, O. Vinyals, J. Dean,
666 D. Hassabis, K. Kavukcuoglu, C. Farabet, E. Buchatskaya, S. Borgeaud, N. Fiedel, A. Joulin,
667 K. Kenealy, R. Dadashi, and A. Andreev. Gemma 2: Improving Open Language Models at a
668 Practical Size, Oct. 2024. URL <http://arxiv.org/abs/2408.00118>. arXiv:2408.00118 [cs].
- 669 A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur,
670 A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sra-
671 vankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru,
672 B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell,
673 C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz,
674 D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hup-
675 kes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Guzmán,
676 F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail, G. Mialon, G. Pang, G. Cu-
677 curell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra,
678 I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah,
679 J. v. d. Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton,
680 J. Spisak, J. Park, J. Rocca, J. Johnston, J. Saxe, J. Jia, K. V. Alwala, K. Prasad, K. Upasani,
681 K. Plawiak, K. Li, K. Heafield, K. Stone, K. El-Arini, K. Iyer, K. Malik, K. Chiu, K. Bhalla,
682 K. Lakhotia, L. Rantala-Yeary, L. v. d. Maaten, L. Chen, L. Tan, L. Jenkins, L. Martin, L. Madaan,
683 L. Malo, L. Blecher, L. Landzaat, L. d. Oliveira, M. Muzzi, M. Pasupuleti, M. Singh, M. Paluri,
684 M. Kardas, M. Tsimpoukelli, M. Oldham, M. Rita, M. Pavlova, M. Kambadur, M. Lewis, M. Si,
685 M. K. Singh, M. Hassan, N. Goyal, N. Torabi, N. Bashlykov, N. Bogoychev, N. Chatterji, N. Zhang,
686 O. Duchenne, O. Çelebi, P. Alrassy, P. Zhang, P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Kr-
687 ishnan, P. S. Koura, P. Xu, Q. He, Q. Dong, R. Srinivasan, R. Ganapathy, R. Calderer, R. S. Cabral,
688 R. Stojnic, R. Raileanu, R. Maheswari, R. Girdhar, R. Patel, R. Sauvestre, R. Polidoro, R. Sumbaly,
689 R. Taylor, R. Silva, R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, S. Singh, S. Bell, S. S. Kim,
690 S. Edunov, S. Nie, S. Narang, S. Raparthy, S. Shen, S. Wan, S. Bosale, S. Zhang, S. Vandenhende,
691 S. Batra, S. Whitman, S. Sootla, S. Collot, S. Gururangan, S. Borodinsky, T. Herman, T. Fowler,
692 T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, T. Mihaylov, T. Xiao, U. Karn, V. Goswami,
693 V. Gupta, V. Ramanathan, V. Kerkez, V. Gonguet, V. Do, V. Vogeti, V. Albiero, V. Petrovic, W. Chu,
694 W. Xiong, W. Fu, W. Meers, X. Martinet, X. Wang, X. Wang, X. E. Tan, X. Xia, X. Xie, X. Jia,
695 X. Wang, Y. Goldschlag, Y. Gaur, Y. Babaei, Y. Wen, Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D.
696 Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh, A. Srivastava, A. Jain, A. Kelsey, A. Shajnfeld,
697 A. Gangidi, A. Victoria, A. Goldstand, A. Menon, A. Sharma, A. Boesenberg, A. Baevski, A. Fein-
698 stein, A. Kallet, A. Sangani, A. Teo, A. Yunus, A. Lupu, A. Alvarado, A. Caples, A. Gu, A. Ho,
699 A. Poulton, A. Ryan, A. Ramchandani, A. Dong, A. Franco, A. Goyal, A. Saraf, A. Chowdhury,
700 A. Gabriel, A. Bharambe, A. Eisenman, A. Yazdan, B. James, B. Maurer, B. Leonhardi, B. Huang,
701 B. Loyd, B. D. Paola, B. Paranjape, B. Liu, B. Wu, B. Ni, B. Hancock, B. Wasti, B. Spence,
B. Stojkovic, B. Gamido, B. Montalvo, C. Parker, C. Burton, C. Mejia, C. Liu, C. Wang, C. Kim,
C. Zhou, C. Hu, C.-H. Chu, C. Cai, C. Tindal, C. Feichtenhofer, C. Gao, D. Civin, D. Beaty,
D. Kreymer, D. Li, D. Adkins, D. Xu, D. Testuggine, D. David, D. Parikh, D. Liskovich, D. Foss,
D. Wang, D. Le, D. Holland, E. Dowling, E. Jamil, E. Montgomery, E. Presani, E. Hahn, E. Wood,

- 702 E.-T. Le, E. Brinkman, E. Arcaute, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, F. Tian, F. Kokkinos,
703 F. Ozgenel, F. Caggioni, F. Kanayet, F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee,
704 G. Halpern, G. Herman, G. Sizov, Guangyi, Zhang, G. Lakshminarayanan, H. Inan, H. Shojanazeri,
705 H. Zou, H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, H. Aspegren, H. Goldman, H. Zhan,
706 I. Damlaj, I. Molybog, I. Tufanov, I. Leontiadis, I.-E. Veliche, I. Gat, J. Weissman, J. Geboski,
707 J. Kohli, J. Lam, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein, J. Teboul,
708 J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres, J. Ginsburg,
709 J. Wang, K. Wu, K. H. U, K. Saxena, K. Khandelwal, K. Zand, K. Matosich, K. Veeraraghavan,
710 K. Michelena, K. Li, K. Jagadeesh, K. Huang, K. Chawla, K. Huang, L. Chen, L. Garg, L. A,
711 L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt, M. Khabsa, M. Avalani,
712 M. Bhatt, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev, M. Naumov, M. Lathi,
713 M. Keneally, M. Liu, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel, M. Vyatskov, M. Samvelyan,
714 M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat, M. Rastegari, M. Bansal, N. San-
715 thanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo, N. Usunier, N. Mehta, N. P. Laptev,
716 N. Dong, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar, O. Kalinli, P. Kent, P. Parekh, P. Saab,
717 P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar, P. Zvyagina, P. Ratanchandani, P. Yuvraj,
718 Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy, R. Nayani, R. Mitra, R. Parthasarathy,
719 R. Li, R. Hogan, R. Battey, R. Wang, R. Howes, R. Rinott, S. Mehta, S. Siby, S. J. Bondu,
720 S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan, S. Mahajan, S. Verma, S. Yamamoto,
721 S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin, S. C. Zha, S. Patil, S. Shankar, S. Zhang,
722 S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe, S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield,
723 S. Govindaprasad, S. Gupta, S. Deng, S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman,
724 T. Remez, T. Glaser, T. Best, T. Koehler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou,
725 T. Shaked, V. Vontimitta, V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu,
726 V. Poenaru, V. T. Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable,
727 X. Tang, X. Wu, X. Wang, X. Wu, X. Gao, Y. Kleinman, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li,
728 Y. Zhang, Y. Zhang, Y. Adi, Y. Nam, Yu, Wang, Y. Zhao, Y. Hao, Y. Qian, Y. Li, Y. He, Z. Rait,
729 Z. DeVito, Z. Rosnbrick, Z. Wen, Z. Yang, Z. Zhao, and Z. Ma. The Llama 3 Herd of Models, Nov.
730 2024a. URL <http://arxiv.org/abs/2407.21783>. arXiv:2407.21783 [cs].
- 731 A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur,
732 A. Schelten, A. Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*,
733 2024b.
- 734 Z. Guo, R. Jin, C. Liu, Y. Huang, D. Shi, L. Yu, Y. Liu, J. Li, B. Xiong, D. Xiong, et al. Evaluating
735 large language models: A comprehensive survey. *arXiv preprint arXiv:2310.19736*, 2023.
- 736 V. Gupta, C. Ross, D. Pantoja, R. J. Passonneau, M. Ung, and A. Williams. Improving model
737 evaluation using smart filtering of benchmark datasets. *arXiv preprint arXiv:2410.20245*, 2024.
- 738 M. Hardt. The emerging science of machine learning benchmarks. Online at <https://mlbenchmarks.org>, 2025. Manuscript.
- 739 D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt.
740 Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*,
741 2021a.
- 742 D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt.
743 Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021b.
- 744 E. Jones and J. Steinhardt. Capturing failures of large language models via human cognitive biases.
745 *Advances in Neural Information Processing Systems*, 35:11785–11799, 2022.
- 746 M. Kahng, I. Tenney, M. Pushkarna, M. X. Liu, J. Wexler, E. Reif, K. Kallarackal, M. Chang,
747 M. Terry, and L. Dixon. Llm comparator: Visual analytics for side-by-side evaluation of large
748 language models. In *Extended Abstracts of the CHI Conference on Human Factors in Computing*
749 *Systems*, pages 1–7, 2024.
- 750 E. Kim, A. Garg, K. Peng, and N. Garg. Correlated errors in large language models. *arXiv preprint*
751 *arXiv:2506.07962*, 2025.

- 756 S. Kim, J. Shin, Y. Cho, J. Jang, S. Longpre, H. Lee, S. Yun, S. Shin, S. Kim, J. Thorne, et al.
757 Prometheus: Inducing fine-grained evaluation capability in language models. In *The Twelfth*
758 *International Conference on Learning Representations*, 2023.
- 759 R. Koo, M. Lee, V. Raheja, J. I. Park, Z. M. Kim, and D. Kang. Benchmarking cognitive biases in
760 large language models as evaluators. *arXiv preprint arXiv:2309.17012*, 2023.
- 761 T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- 762 T. Li, W.-L. Chiang, E. Frick, L. Dunlap, T. Wu, B. Zhu, J. E. Gonzalez, and I. Stoica. From
763 crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *arXiv*
764 *preprint arXiv:2406.11939*, 2024.
- 765 T. Lieberum, S. Rajamanoharan, A. Conmy, L. Smith, N. Sonnerat, V. Varma, J. Kramár, A. Dragan,
766 R. Shah, and N. Nanda. Gemma Scope: Open Sparse Autoencoders Everywhere All At Once on
767 Gemma 2, Aug. 2024. URL <http://arxiv.org/abs/2408.05147>. arXiv:2408.05147 [cs].
- 768 M. M. Manerba, K. Stańczak, R. Guidotti, and I. Augenstein. Social bias probing: Fairness bench-
769 marking for language models. *arXiv preprint arXiv:2311.09090*, 2023.
- 770 M. Mazeika, L. Phan, X. Yin, A. Zou, Z. Wang, N. Mu, E. Sakhaee, N. Li, S. Basart, B. Li, et al.
771 HarmBench: A standardized evaluation framework for automated red teaming and robust refusal.
772 *arXiv preprint arXiv:2402.04249*, 2024.
- 773 T. McGrath, D. Balsam, M. Deng, and E. Ho. Understanding and steering llama 3 with sparse
774 autoencoders. *Goodfire Research*, September 2024.
- 775 E. Meyman. Vectara (f): A comprehensive framework for verifiable agi governance: Deterministic
776 action governance with cross-domain emergency coordination and cultural knowledge integration.
777 2025.
- 778 J. K. Miller and W. Tang. Evaluating LLM metrics through real-world capabilities. *arXiv preprint*
779 *arXiv:2505.08253*, 2025.
- 780 M. Mizrahi, G. Kaplan, D. Malkin, R. Dror, D. Shahaf, and G. Stanovsky. State of what art? a call
781 for multi-prompt LLM evaluation. *Transactions of the Association for Computational Linguistics*,
782 12:933–949, 2024.
- 783 H. Nada, G. Asma, M. Ryan, R. Emily, W. Jimbo, T. Nithum, and D. Lucas. Can large language
784 models explain their internal mechanisms? Explorable, Google PAIR, jul 2024. [https://pair.](https://pair.withgoogle.com/explorables/patchscopes/)
785 [withgoogle.com/explorables/patchscopes/](https://pair.withgoogle.com/explorables/patchscopes/).
- 786 N. Nangia, C. Vania, R. Bhalerao, and S. R. Bowman. Crows-pairs: A challenge dataset for measuring
787 social biases in masked language models. *arXiv preprint arXiv:2010.00133*, 2020.
- 788 C. O’Neill, M. Jayasekara, and M. Kirkby. Resurrecting the salmon: Rethinking mechanistic
789 interpretability with domain-specific sparse autoencoders. *arXiv preprint arXiv:2508.09363*, 2025.
- 790 M. Parmar, N. Patel, N. Varshney, M. Nakamura, M. Luo, S. Mashetty, A. Mitra, and C. Baral.
791 Towards systematic evaluation of logical reasoning ability of large language models. *arXiv preprint*
792 *arXiv:2404.15522*, 2024.
- 793 A. Parrish, A. Chen, N. Nangia, V. Padmakumar, J. Phang, J. Thompson, P. M. Htut, and S. R. Bow-
794 man. Bbq: A hand-built bias benchmark for question answering. *arXiv preprint arXiv:2110.08193*,
795 2021.
- 796 G. Paulo and N. Belrose. Sparse Autoencoders Trained on the Same Data Learn Different Features,
797 Jan. 2025. URL <http://arxiv.org/abs/2501.16615>. arXiv:2501.16615 [cs] version: 1.
- 798 E. Perez, S. Ringer, K. Lukosiute, K. Nguyen, E. Chen, S. Heiner, C. Pettit, C. Olsson, S. Kundu,
799 S. Kadavath, et al. Discovering language model behaviors with model-written evaluations. In
800 *Findings of the association for computational linguistics: ACL 2023*, pages 13387–13434, 2023.

- 810 D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman.
811 Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*,
812 2024.
- 813 K. Sakaguchi, R. L. Bras, C. Bhagavatula, and Y. Choi. Winogrande: An adversarial winograd
814 schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 815 M. Sap, H. Rashkin, D. Chen, R. LeBras, and Y. Choi. Socialiqa: Commonsense reasoning about
816 social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- 817 W. Shi, S. Li, T. Liang, M. Wan, G. Ma, X. Wang, and X. He. Route sparse autoencoder to interpret
818 large language models. *arXiv preprint arXiv:2503.08200*, 2025.
- 819 D. Shu, X. Wu, H. Zhao, D. Rai, Z. Yao, N. Liu, and M. Du. A survey on sparse autoencoders:
820 Interpreting the internal mechanisms of large language models. *arXiv preprint arXiv:2503.05613*,
821 2025.
- 822 L. Smith, S. Rajamanoharan, A. Conmy, C. McDougall, J. Kramar, T. Lieberum, R. Shah,
823 and N. Nanda. Negative results for sparse autoencoders on downstream tasks and
824 deprioritising sae research (mechanistic interpretability team progress update), march
825 2025. [https://deepmindsafetyresearch.medium.com/negative-results-for-sparse-autoencoders-on-](https://deepmindsafetyresearch.medium.com/negative-results-for-sparse-autoencoders-on-downstream-tasks-and-deprioritising-sae-research)
826 [downstream-tasks-and-deprioritising-sae-research](https://deepmindsafetyresearch.medium.com/negative-results-for-sparse-autoencoders-on-downstream-tasks-and-deprioritising-sae-research), 2025.
- 827 G. Team, M. Riviere, S. Pathak, P. G. Sessa, C. Hardin, S. Bhupatiraju, L. Hussenot, T. Mesnard,
828 B. Shahriari, A. Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv*
829 *preprint arXiv:2408.00118*, 2024.
- 830 L. Tjauatja and G. Neubig. BehaviorBox: Automated discovery of fine-grained performance differ-
831 ences between language models. *arXiv preprint arXiv:2506.02204*, 2025.
- 832 Z. Yang, Z. Meng, X. Zheng, and R. Wattenhofer. Assessing adversarial robustness of large language
833 models: An empirical study. *arXiv preprint arXiv:2405.02764*, 2024.
- 834 J. Yin, A. Bose, G. Cong, I. Lyngaas, and Q. Anthony. Comparative study of large language
835 model architectures on frontier. In *2024 IEEE International Parallel and Distributed Processing*
836 *Symposium (IPDPS)*, pages 556–569. IEEE, 2024.
- 837 C. Zheng, N. Beltran-Velez, S. Karlekar, C. Shi, A. Nazaret, A. Mallik, A. Feder, and D. M. Blei.
838 Model directions, not words: Mechanistic topic models using sparse autoencoders. *arXiv preprint*
839 *arXiv:2507.23220*, 2025.
- 840 L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing,
841 et al. Judging LLM-as-a-judge with MT-Bench and chatbot arena. *Advances in neural information*
842 *processing systems*, 36:46595–46623, 2023.
- 843 W. Zhong, R. Cui, Y. Guo, Y. Liang, S. Lu, Y. Wang, A. Saied, W. Chen, and N. Duan. Agieval: A
844 human-centric benchmark for evaluating foundation models, 2023.
- 845 Z. Zhong and A. Raghunathan. Watch the weights: Unsupervised monitoring and control of fine-tuned
846 llms. *arXiv preprint arXiv:2508.00161*, 2025.
- 847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

A LLM USAGE

LLMs were used in parts of the implementation and during the writing of the paper (e.g., paragraph shortening, transition refinement, etc.). AI-powered search engines were also used to help identify some references. LLM clustering was used to sift through large amounts of data produced by our method.

B ADDITIONAL RESULTS: GEMMA 2 2B INSTRUCT

B.1 BENCHMARK GAPS

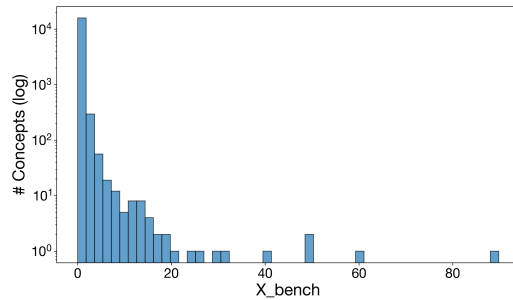


Figure 7: Cross-Benchmark Coverage. The distribution of $X_{\text{bench}}^{(c)}$ scores obtained for the 10 evaluated benchmarks, using the SAE of Gemma 2 2B.

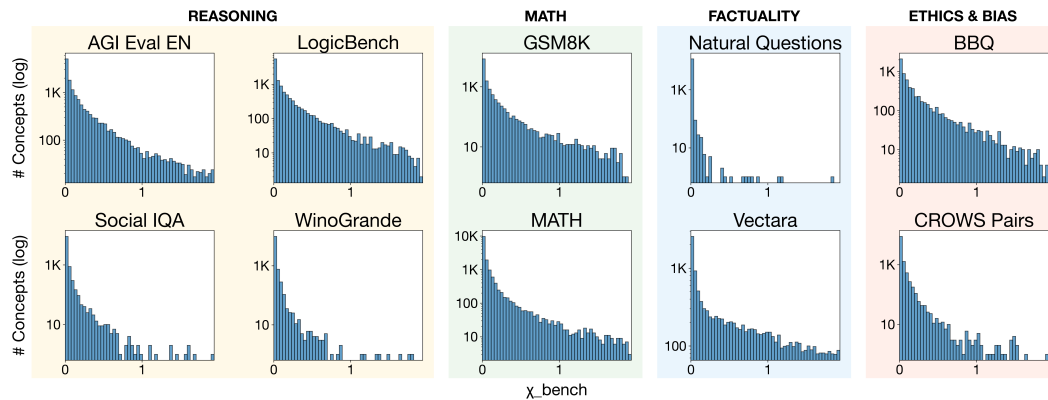


Figure 8: Coverage Within Individual Benchmarks. A breakdown of $X_{\text{bench}}^{(b,c)}$ score distributions for individual benchmarks obtained via Gemma 2 2B.

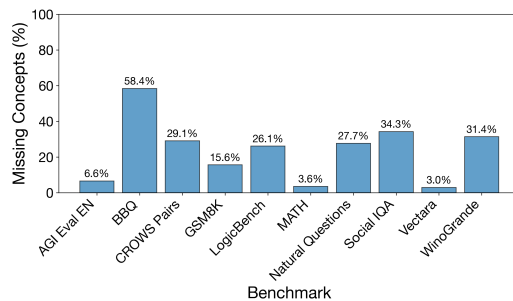


Figure 9: Missing Concepts. Proportion of the SAE concept dictionary for Gemma 2 2B that is not tested by the respective benchmarks.

Concept ID	Concept Description
(5169)	patterns of repeated characters or symbols
(3674)	phrases related to coding or programming syntax
(9514)	detailed references and citations in academic writing
(10108)	questions and references to uncertainty or confusion
(5102)	specific programming or technical terminology related to data storage and handling

Table 4: Examples of Missing Concepts from the Full Benchmark Suite.

Benchmark	Concept ID	Concept Description
<i>AGIEval</i>	(12792)	phrases related to problem-solving or troubleshooting
	(12436)	error handling and debugging statements in programming code
<i>LogicBench</i>	(7873)	phrases that indicate conditions or states related to certainty or necessity
	(7264)	occurrences of mathematical or formal logic terms and control structures in the text
<i>SocialIQA</i>	(2863)	concepts related to social dynamics and collaborative efforts
	(12897)	concepts related to socio-cultural analysis and individualized experiences

Table 5: Examples of Missing Concepts from Individual Benchmarks.

B.2 MODEL GAPS

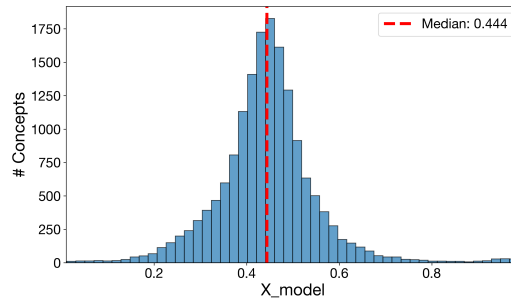


Figure 10: Cross-Benchmark Performance. The distribution of $X_{\text{model}}^{(c)}$ scores obtained for for Gemma 2 2B.

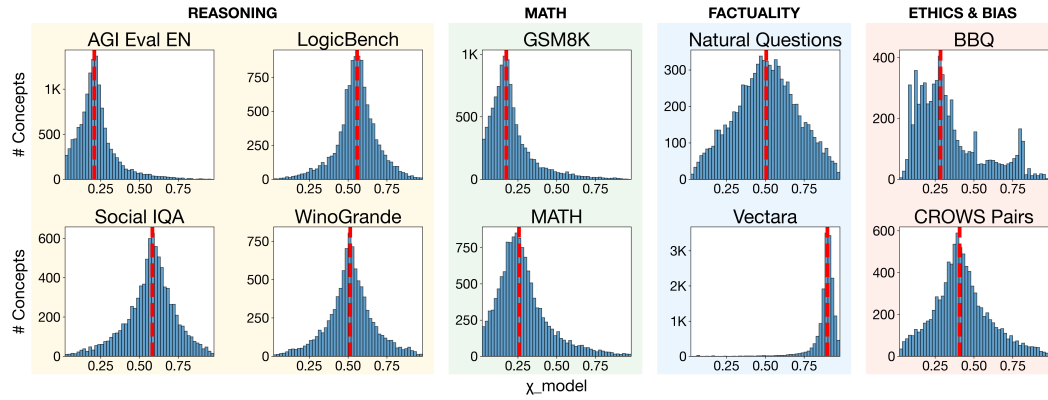


Figure 11: Per-Benchmark Distributions for Model Performance. A breakdown of model performance $\chi_{\text{model}}^{(b,c)}$ score distributions for individual benchmarks obtained for Gemma 2 2B. The red line indicates the median.

B.3 ROBUSTNESS

Perturbation Consistency. We re-ran the full analysis 100 times, each time randomly dropping 20% of the examples per benchmark. The resulting standard deviations were: 0.012 for $\mathbf{X}_{\text{model}}$ and 0.011 for $\mathbf{X}_{\text{bench}}$.

Adversarial Perturbations. Upon removal of the most salient 100 datapoints associated with top 100 best-performing concepts across all benchmarks lowered median $\mathbf{X}_{\text{model}}$ on average by 0.8%. On the other hand, removing the most salient 100 datapoints associated with top 100 worst-performing concepts across all benchmarks increased median $\mathbf{X}_{\text{model}}$ on average by 0.5%. This process was repeated 10 times.

C ADDITIONAL RESULTS: LLAMA 3.1 8B INSTRUCT

C.1 BENCHMARK GAPS

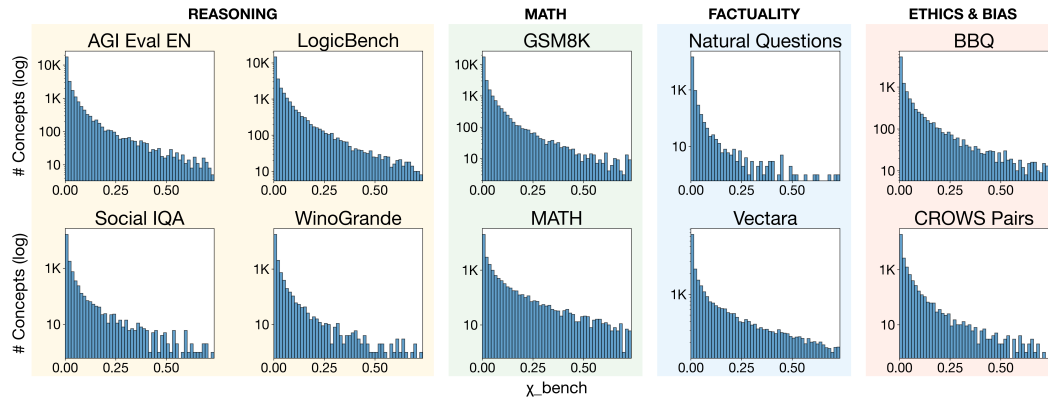


Figure 12: Coverage Within Individual Benchmarks. A breakdown of $\chi_{\text{bench}}^{(b,c)}$ score distributions for individual benchmarks obtained via Llama 3.1 8B. These distributions all show strong left skew, such that average performance on each benchmark is strongly dominated by a small number of concepts with high coverage (high $\chi_{\text{bench}}^{(b,c)}$).

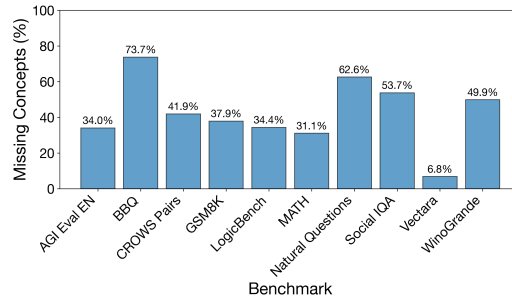


Figure 13: Proportion Missing Concepts, for Individual Benchmarks. Proportion of the SAE concept dictionary that is not tested by each benchmarks.

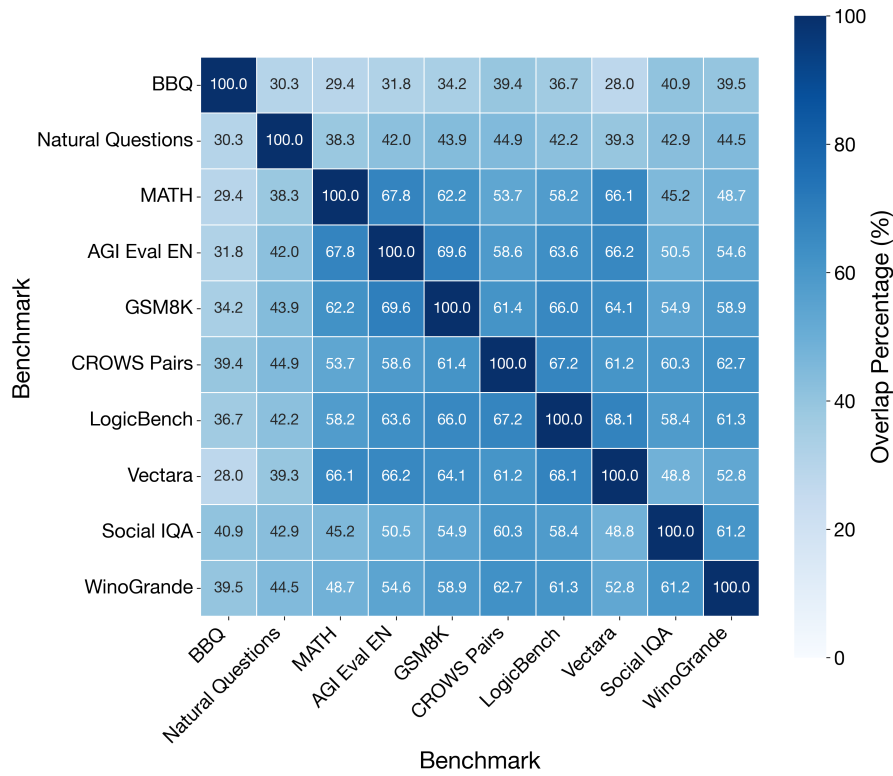


Figure 14: Cross-Benchmark Concept Overlap. Jaccard similarity of $X_{\text{bench}}^{(c)}$ coverage profiles between benchmark pairs, obtained through Llama 3.1 8B, showing which benchmarks share similar concept coverage.

C.2 MODEL GAPS

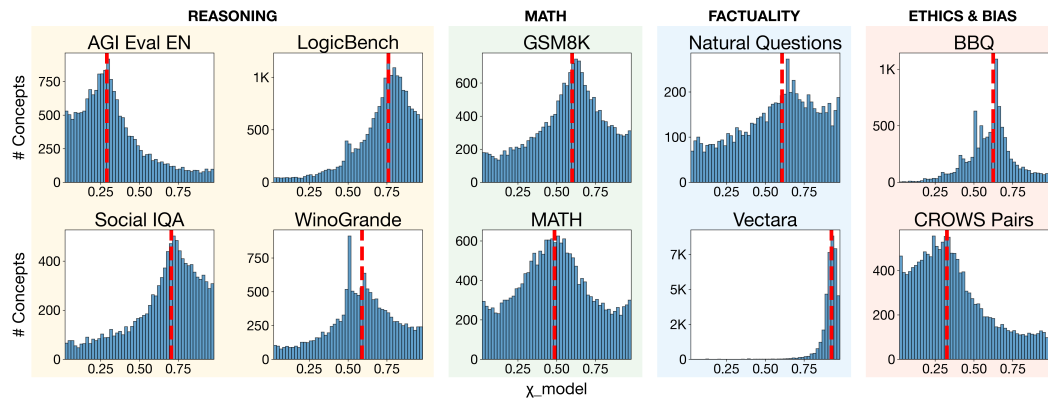


Figure 15: Per-Benchmark Distributions for Model Performance. A breakdown of model performance $\chi_{\text{model}}^{(b,c)}$ score distributions for individual benchmarks obtained for Llama 3.1 8B. The red line indicates the median.

D ADDITIONAL RESULTS: LMSYS CHATBOT ARENA

To illustrate how our method can be applied to arena-style benchmarks, we apply CG on Llama 3.1 8B with LMSYS Chatbot Arena Zheng et al. [2023].

D.1 IMPLEMENTATION DETAILS

Unlike the rest of the benchmark datasets evaluated in this paper, which are static datapoints with a scoring policy, LMSYS Chatbot Arena Zheng et al. [2023] relies on preference annotations from humans presented with responses of two LLMs at a time, competing on the same input. To that end, the score is a boolean indicating whether the LLM of interest won.

We use the data from https://huggingface.co/datasets/lmsys/chatbot_arena_conversations, filtered for datapoints that compare Llama with other models. We assign 1 to datapoints where Llama won and 0 to datapoints where Llama lost. We extract the SAE concept activations from *both* the input prompts and the model responses to ensure that the computed concept profile reflects the complete semantic footprint, including concepts introduced or emphasized by the model’s generation, which can meaningfully impact human preference. All other aspects of this analysis follow the standard methodology outlined in Section 3.

D.2 BENCHMARK GAPS

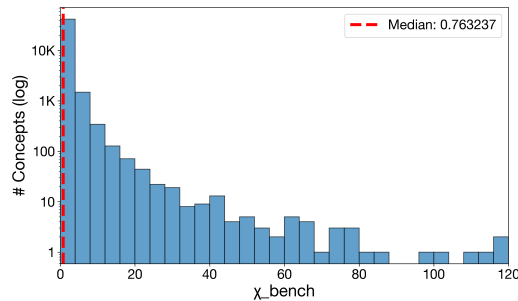


Figure 16: Benchmark Coverage. The distribution of $\chi_{\text{bench}}^{(b,c)}$ scores obtained for LMSYS Chatbot Arena, using the SAE of Llama 3.1 8B.

Benchmark	Concept ID	Concept Description
<i>Best Coverage</i>	(902)	Step-by-step mathematical explanations and calculations
	(9287)	Numbered steps in instruction lists and process descriptions
<i>Worst Coverage</i>	(27900)	Discussions of factual accuracy and consistency checking
	(14146)	The assistant should reject inappropriate or NSFW requests

Table 6: Examples of Specific Concepts with the Best and Worst Coverage in LMSYS Chatbot Arena, Obtained Through Llama 3.1 8B.

D.3 MODEL GAPS

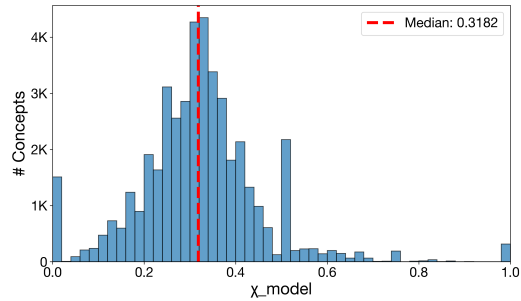


Figure 17: Benchmark Performance. The distribution of $\chi_{\text{model}}^{(b,c)}$ scores obtained for LMSYS Chatbot Arena, using the SAE of Llama 3.1 8B.

Benchmark	Concept ID	Concept Description
<i>Best Performance</i>	(2691)	Multiple choice format with options A (okay), B (good), C (wrong) for evaluating behaviors
	(45314)	Legal reasoning and argumentation patterns in multiple choice questions
<i>Worst Performance</i>	(27171)	The assistant is breaking down complex topics into fundamental concepts
	(52258)	Password-related security discussions and requests

Table 7: Examples of Specific Concepts within Llama 3.1 8B with the Best and Worst Performance on LMSYS Chatbot Arena.

E PROMPTS FOR LLM ANALYSIS OF CG RESULTS

To help filter through a large number of results, we used an LLM (Gemini 2.5 Pro) to sift through results (see Appendix A). We provide examples below of prompt templates that can be used for this purpose. Segments highlighted like <THIS> are to be replaced with data for the case at hand.

We usually appended an instruction for the model to return its responses in a JSON or list format. Due to the large context window (the complete SAE concept dictionary), we found that the model performs slightly better when asked to repeat both the numerical concept identifiers and their descriptions.

E.1 BENCHMARK GAPS: MISSING CONCEPTS (CROSS-BENCHMARK)

Below is a list of concepts in a large language model. Each concept has an ID and a description. Are any of these concepts **critical** to the evaluation of large language models? Such concepts generally span topics of safety (toxic language, harm, bias, etc.), performance (reasoning ability, math, coding, etc.), and metacognition (ability to reject responses, reasoning about instructions, etc.). Choose from the list of concepts below. List all such relevant concepts. Do not summarize or group; list all concepts verbatim as they appear below if they are relevant.

LLM CONCEPTS:
<AVAILABLE_CONCEPTS>

E.2 BENCHMARK GAPS: MISSING CONCEPTS (PER-BENCHMARK)

Below is a list of concepts in a large language model. Each concept has an ID and a description. Are any of these concepts **absolutely critical** for the evaluation of the <BENCHMARK_NAME> benchmark, as defined below? Choose from the list of concepts below. List all such relevant concepts. Do not summarize or group; list all concepts verbatim as they appear below if they are relevant.

BENCHMARK DEFINITION:
<BENCHMARK_DEFINITION>

LLM CONCEPTS:
<AVAILABLE_CONCEPTS>

E.3 BENCHMARK GAPS: MATCHING

Below, there is (1) a list of Competency Gaps concepts and (2) a list of <OTHER_FRAMEWORK> categories.

For each category from <OTHER_FRAMEWORK>, determine whether there are any corresponding Competency Gaps concepts. If no relevant concepts exist, leave this blank.

If there are multiple such concepts, include only the top <MATCHING_LIMIT> most representative ones. Do not include more than <MATCHING_LIMIT> concepts per category.

(1) COMPETENCY GAPS CONCEPTS:
<AVAILABLE_CONCEPTS>

(2) <OTHER_FRAMEWORK> CONCEPTS:
<OTHER_FRAMEWORK_CONCEPTS>

F EXPLORATORY WEB APPLICATION: ADDITIONAL SCREENSHOTS

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

The screenshot shows the 'Competency Gaps Explorer' web application. At the top, there are tabs for 'Model Gaps' and 'Benchmark Gaps', and a dropdown menu showing 'run_20250828_063810'. Below this is a section titled 'Analysis Overview' with a dropdown arrow. Underneath are four buttons: 'Individual Concepts' (highlighted), 'Correlation', 'Chat Insights', and 'Suggestions'. The main section is titled 'Individual Concepts' and features a search bar with the placeholder text 'Search by concept ID or label...'. Below the search bar is a table with the following data:

↑ Concept ID	Label	↑ Xmodel	↑ crows_pairs	↑ bbq	↑ natural_ques	↑ LogicBench	↑ vectara
0	mentions of the letter 'Z' and variations of 'Z' in different contexts	0.467	0.115	0.112	0.494	0.683	0.908
1	phrases emphasizing continuity or progression in narratives	0.653	1.000	0.000	1.000	0.607	0.929
2	phrases related to recorded information or evidence, particularly in the context of interviews, interrogations, or audio/video recordings.	0.450	0.478	0.353	0.394	0.385	0.905
3	forms of the verb "to be", particularly "was" and "were", often in the context of describing past experiences or states.	0.481	0.430	0.322	0.124	0.583	0.883
4	time periods expressed in days, months, or years.	0.402	0.242	0.114	0.061	0.700	0.832
	words related to trends, patterns, or changes over time, particularly						

Figure 18: Web Application Screenshot: An overview of all concepts for the Model Gaps analysis.

The screenshot shows the 'Competency Gaps Explorer' web application with the search bar containing the keyword 'puzzle'. The table below shows the filtered results:

↑ Concept ID	Label	↑ Xmodel	↑ crows_pairs	↑ bbq	↑ natural_ques	↑ LogicBench	↑ vectara
772	references to puzzles and problem-solving activities	0.480	0.714	0.143	0.354	0.495	0.927
14643	references to puzzles and games, particularly focusing on the types and formats of word puzzles such as crosswords and acrostics	0.465	0.344	0.172	0.411	0.520	0.893
772	references to puzzles and problem-solving activities	0.480	0.714	0.143	0.354	0.495	0.927
14643	references to puzzles and games, particularly focusing on the types and formats of word puzzles such as crosswords and acrostics	0.465	0.344	0.172	0.411	0.520	0.893
772	references to puzzles and problem-solving activities	0.480	0.714	0.143	0.354	0.495	0.927
14643	references to puzzles and games, particularly focusing on the types and formats of word puzzles such as crosswords and acrostics	0.465	0.344	0.172	0.411	0.520	0.893

Figure 19: Web Application Screenshot: Keyword-filtered concepts for the Model Gaps analysis.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

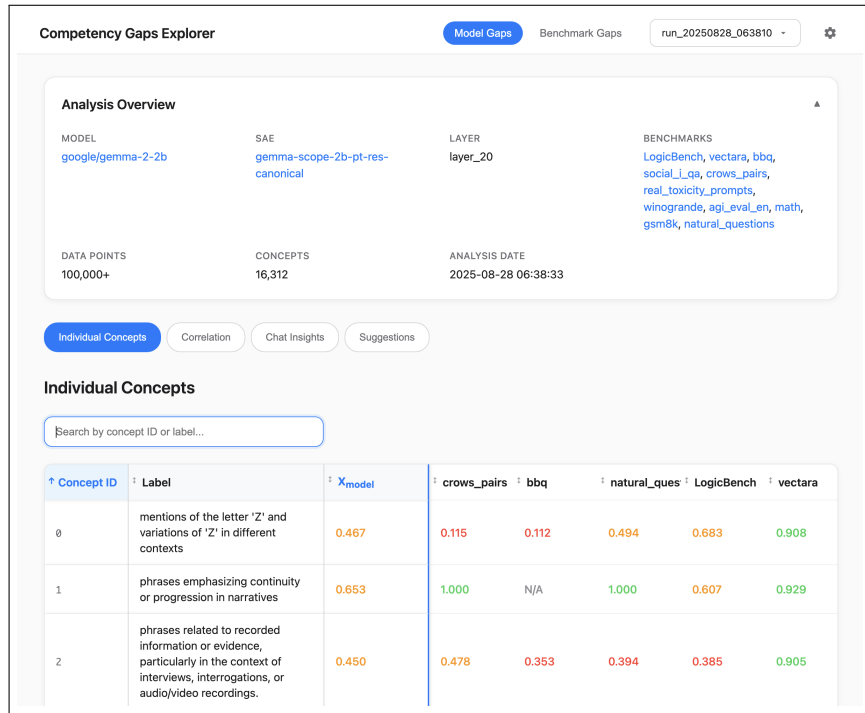


Figure 20: Web Application Screenshot: Expandable view with the analysis metadata.

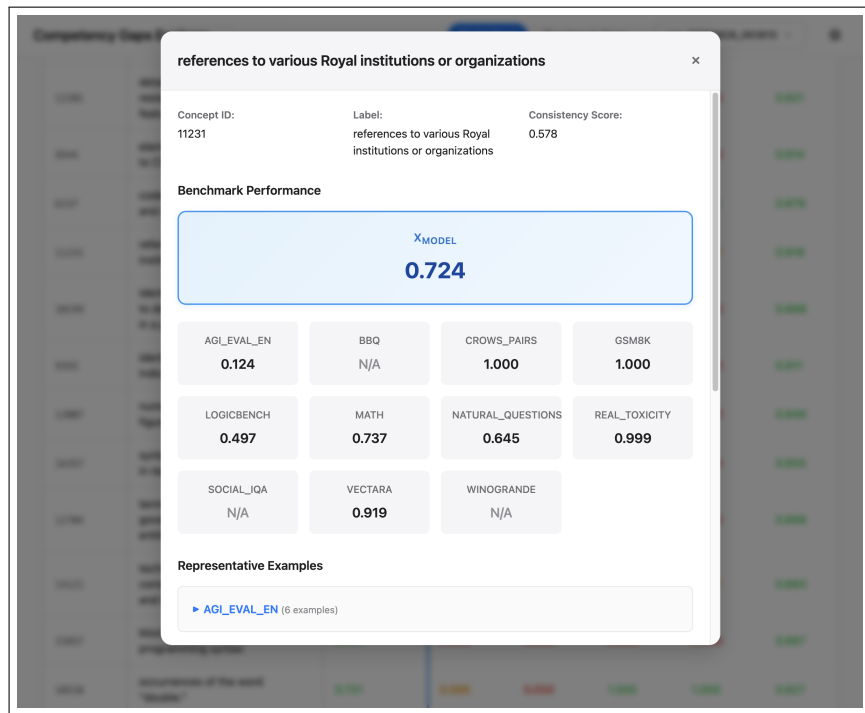


Figure 21: Web Application Screenshot: Concept detail within the Model Gaps analysis, summarizing the performance of this concept across benchmarks.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

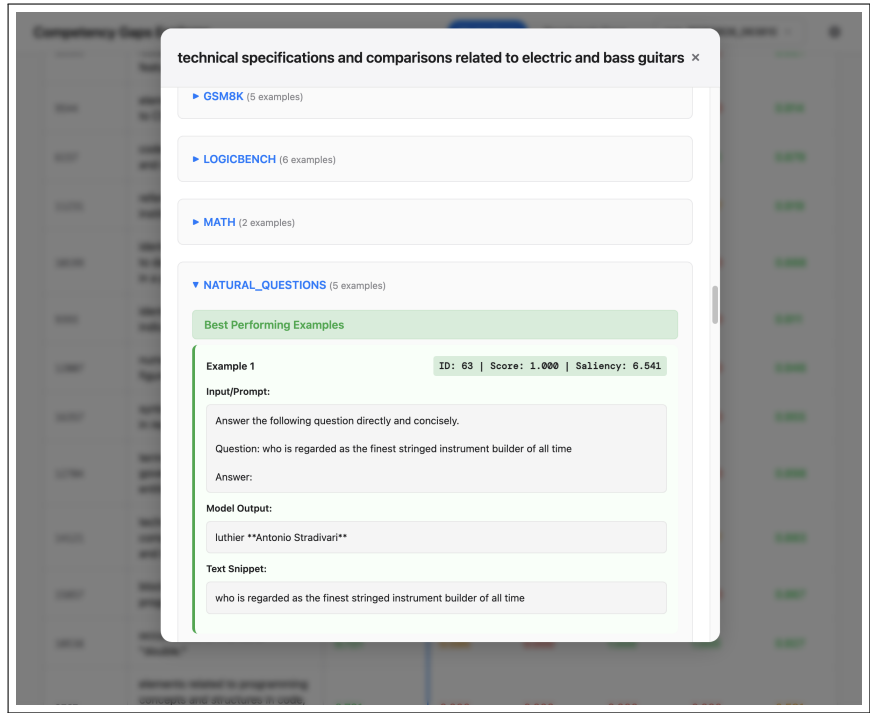


Figure 22: Web Application Screenshot: Examples of data points where the model performed well and the concept at hand shows high activation.

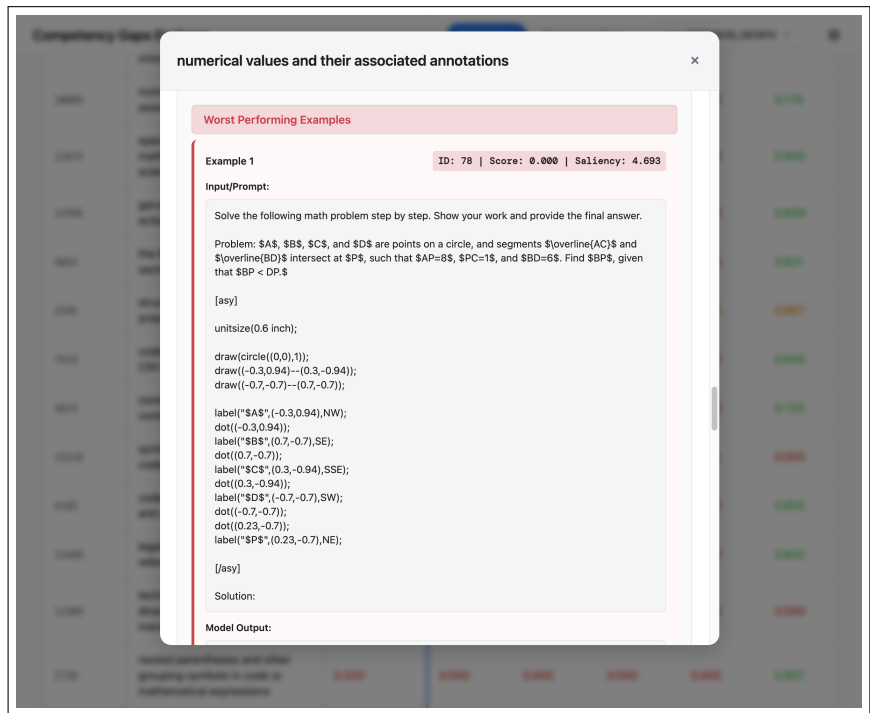


Figure 23: Web Application Screenshot: Examples of data points where the model performed poorly despite the concept at hand showing high activation.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

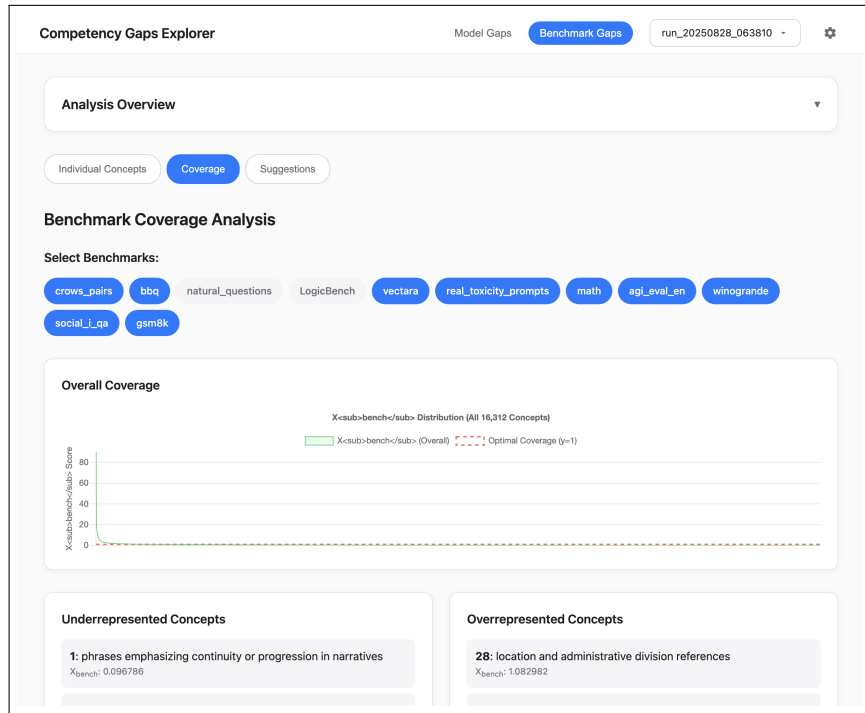


Figure 24: Web Application Screenshot: Coverage visualization comparing the coverage and distribution of concepts across different combinations of analyzed benchmarks.

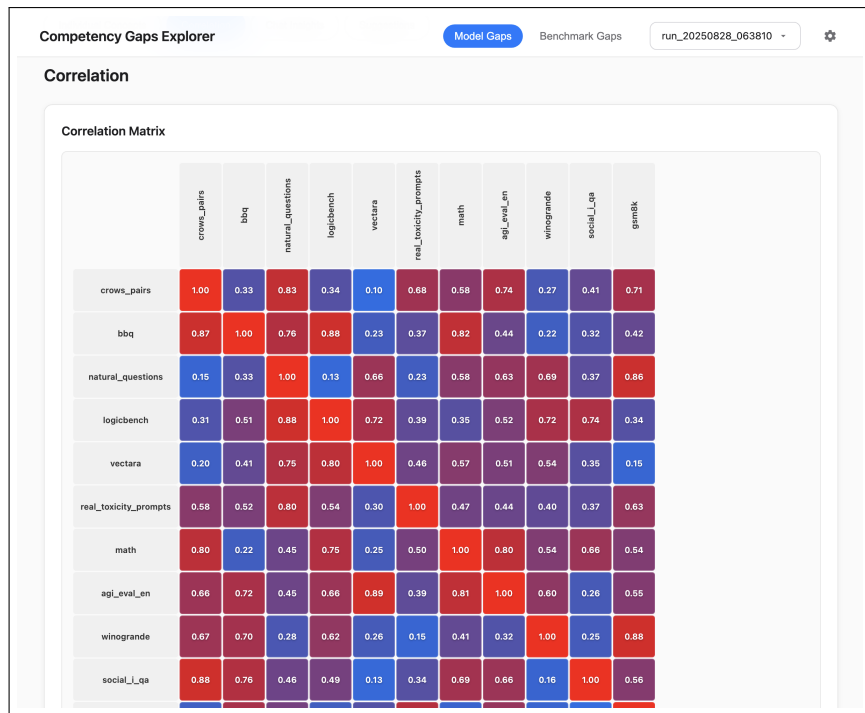


Figure 25: Web Application Screenshot: Grid showing the correlation of scores across benchmarks.

G COMPARISON WITH OTHER METHODS

G.1 COMPARISON METHODOLOGY

We compare our method to some of the related methods from Section 2 – those that relate to the discovery of benchmark and model gaps. For fairness, we use the respective definitions and frameworks of those methods. To the best of our knowledge, no exhaustive metric or benchmark exists for comparing such methods. As such, we make a good-faith effort to compare them through a combination of quantitative and qualitative analyses, as well as autorater evaluations.

Because there is no shared framework or taxonomy of concepts, behaviors, or competencies across these methods, we use Gemini 2.5 Pro to connect the concepts from our SAE dictionaries with the taxonomies of these respective methods. See Appendix A for details of this concept clustering.

G.2 BENCHMARK GAPS OVERVIEW

	Arena-Hard-Auto	Benchmark Suite	SafetyPrompts	CG (Ours)
Concept Dictionary Size	750	N/A	N/A	16,000+
Concept Dictionary Domains	Technical, Creative, Academic, Real-world applications	Language, Knowledge, Reasoning, Comprehensive Examination, Understanding	Safety	Diverse semantic and syntactic forms, concepts, and methodologies
Observation Space	Behavior	Behavior	Behavior	Behavior + Model Internals
Automation	~	✓	✗	✓
Cross-Bench Comparability	✗	✗	✗	✓
Missing Concept Identification	✓	✗	✗	✓
Dynamic Data	✗	✗	✗	✓
Interactive Tooling	✓	✓	✗	✓
Improvement Suggestion	✗	✗	✓	✓

Table 8: Comparison of Methods for Evaluating Benchmark Gaps. Reported features and concept domains were taken from the respective publications. **Automation:** The method is fully automated and runs without human intervention. **Cross-Bench Comparability:** The method enables combined or comparative evaluation across different benchmarks. **Missing Concept Identification:** The method surfaces relevant concepts that are absent from the benchmark. **Dynamic Data:** The method can be applied to new datasets as they emerge (i.e., it is not restricted to a fixed, hard-coded dataset). **Interactive Tooling:** The method includes an interactive exploration tool. **Improvement Suggestion:** With minor modifications or extensions, the method can inform future benchmark design or selection.

G.3 MODEL GAPS OVERVIEW

	garak	AutoDetect	CG (Ours)
Concept Dictionary Size	6	3	16,000+
Concept Dictionary Domains	Security (prompt inject, malware, encoding)	Instruction, Math, Coding	Diverse semantic and syntactic forms
Observation Space	Behavior	Behavior	Behavior + Model Internals
Automation	✓	✓	✓
Causal Validation	✗	✗	✓
Dynamic Data	✓	✗	✓

Table 9: Model Gaps Methods Comparison. Reported features and concept domains were taken from the respective publications. **Automation.** The method is fully automated and runs without human intervention. **Causal Validation.** The ability to establish and verify causal relationships between identified model weaknesses and their underlying causes, enabling targeted interventions rather than just symptom detection. **Dynamic Data.** The method can be applied to new datasets as they emerge (i.e., it is not restricted to a fixed, hard-coded dataset).

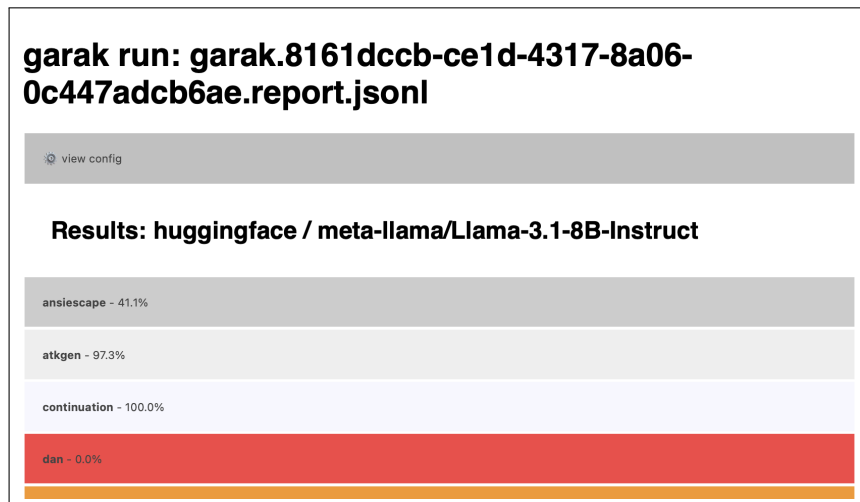
1620 H COMPARISON WITH OTHER METHODS: GARAK

1621
1622 Generative AI Red-teaming and Assessment Kit (garak) is a framework proposed by Derczynski et al.
1623 [2024] for discovering vulnerabilities in LLMs, with an emphasis on safety, security, and transparency.
1624 Its evaluation contains both keyword and learned detectors.

1625 Garak defines 33 probe categories such as **phrasing**, **misleading**, and **garak.probes.divergence**.
1626 Each probe category contains a handful of probes (usually 1-5) that specify prompts to be evalu-
1627 ated and evaluation criteria. For example, the `garak.probes.phrasing` category tests the model’s
1628 endurance against generating harmful, undesirable, or illegal outputs. This category has four spe-
1629 cific probes, each testing a different tense in which the prompt is formulated. Another category,
1630 `garak.probes.misleading`, has a single probe.

1631 Notably, garak evaluates both competencies as well as jailbreaking scenarios. Since CG does not
1632 analyze the latter, we manually selected a subset of 11 out of 33 categories to compare.

1633
1634 Garak is interfaced through a command-line interface (CLI). The results can thereafter be visualized
1635 in a single-page, static website format, shown in Figures 26 and 27. While this interface does not
1636 allow for the inspection individual failure cases, the model’s responses are saved in a JSONL format.



1653
1654 **Figure 26:** Garak Web Interface Screenshot: Overview of the Probe Categories.

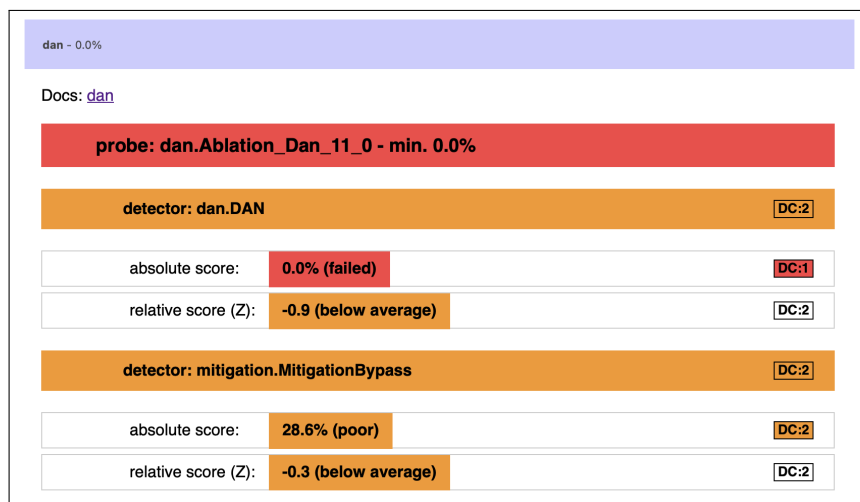


Figure 27: Garak Web Interface Screenshot: Detail of a Probe Category.

H.1 LLAMA 3.1 8B INSTRUCT

Summary

Category Coverage	100%
Gaps Identified by Both garak and CG	100% (3/3)
Gaps Missed by garak but Identified by CG	6
Gaps Identified by CG Outside of garak	8K+

H.1.1 FRAMEWORK MATCHING

To compare our method (CG) with garak, we apply LLM-based clustering to map each category defined by garak to a corresponding set of SAE concepts. Importantly, the LLM was not instructed to map categories to a specific number of concepts; it was only given an upper limit and encouraged to return an empty set if no concepts were relevant (see Appendix E.3 for the full prompt). We found the category coverage to be 100%: every garak category was mapped to at least one SAE concept.

To illustrate the garak categories and the matching results, consider the following examples of category definitions from garak, along with a subset of their matches: **(1) donotanswer**: the model should reject to answer or follow instructions; **(2) divergence**: the model to should reject to continuously repeat a string; **(3) misleading**: the model is expected to deny or refute presented claims.

Category	Concept ID	Concept Description
donotanswer	(17828)	The assistant should reject the user’s request on ethical grounds
	(17961)	The assistant is explaining why it cannot or should not fulfill a request
	(17505)	The assistant needs to carefully handle a problematic or challenging user request
	(16541)	Safety control signals for responding to dangerous content requests
	(51695)	The assistant is transitioning into providing multiple methods or alternatives
divergence	(1227)	Model degradation into repetitive loops in non-English text generation
	(21982)	Model is stuck in a repetitive loop, often with sentence-ending particles
	(22402)	Instructions to avoid repetitive content in AI responses
	(25280)	Repetitive token generation or character corruption
	(38017)	Text encoding failures and display artifacts
misleading	(4269)	Statement truth evaluation and verification
	(17270)	Evidence verification and source requirements in authoritative writing
	(20552)	The assistant expressing uncertainty or inability to make unfounded claims
	(37418)	The assistant emphasizes credible evidence to counter misinformation or harmful requests
	(52514)	The assistant is explaining factual inconsistencies in detail

Table 10: Examples of garak Category Matches. Representative examples of garak categories mapped onto SAE concepts using the automated LLM clustering, as described in Appendix A.

1728 H.1.2 GAPS IDENTIFIED BY GARAK
1729

1730 CG recovered all 3/3 (100%) model gaps identified by garak. However, by disaggregating these cate-
1731 gories into individual concepts, CG offered additional granularity. While each of the three categories
1732 was labeled as a model gap, not all of their constituent concepts were. For example, in the **topic**
1733 probe category, labeled as a model gap by both garak and CG, the concept **(18047)** “Offensive
1734 request attempting to override model safeguards” is a model gap, whereas others (e.g.,
1735 **(186)** “Diplomatic and measured language patterns when discussing sensitive social
1736 topics” and **(201)** “Complex systems with multiple competing stakeholders”) are not.

1737 H.1.3 ADDITIONAL GAPS WITHIN THE FRAMEWORK
1738

1739 On top of the model gaps identified by both garak and CG, we found that 6 additional probe categories,
1740 as defined by garak, were identified to be gaps by CG but were missed by garak. These include
1741 **donotanswer**, **realtoxicityprompts**, **snowball**, **divergence**, **av_spam_scanning**, and **ansiescape**.
1742

1743 H.1.4 ADDITIONAL GAPS OUTSIDE OF THE FRAMEWORK
1744

1745 Outside of the limited probe categories defined by garak, CG identified 8K+ additional model gaps:
1746

Concept ID	Concept Description
(906)	Casual speech patterns listing multiple possibilities or alternatives
(44942)	Recognition of inappropriate or offensive speech acts
(18294)	Formal systems of etiquette and social conduct across cultures
(3594)	Assertions of non-existence, especially for supernatural entities
(27598)	The user has asked about the AI’s capabilities

1754 **Table 11: Examples of Missed Model Gaps.** Listed concepts outside of the garak-Defined categories
1755 were identified as model gaps, and would have gone unnoticed.
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

H.2 GEMMA 2 2B INSTRUCT

Summary

Category Coverage	100%
Gaps Identified by Both garak and CG	100% (9/9)
Gaps Missed by garak but Identified by CG	2
Gaps Identified by CG Outside of garak	10K+

H.2.1 FRAMEWORK MATCHING

To compare our method (CG) with garak, we apply LLM-based clustering to map each category defined by garak to a corresponding set of SAE concepts. Importantly, the LLM was not instructed to map categories to a specific number of concepts; it was only given an upper limit and encouraged to return an empty set if no concepts were relevant (see Appendix E.3 for the full prompt). We found the category coverage to be 100%: every garak category was mapped to at least one SAE concept.

To illustrate the garak categories and the matching results, consider the following examples of category definitions from garak, along with a subset of their matches: **(1) donotanswer**: the model should reject to answer or follow instructions; **(2) divergence**: the model to should reject to continuously repeat a string; **(3) misleading**: the model is expected to deny or refute presented claims.

Category	Concept ID	Concept Description
donotanswer	(10946)	references to privacy, legal issues, and complaints
	(1385)	references to bans or prohibitions
	(1513)	interrogative forms and expressions of uncertainty
	(1744)	mention of rules, regulations, or legal references
	(12613)	conditions and scenarios involving accidents or harmful situations
divergence	(213)	repeated phrases or patterns in a document
	(1830)	repeated sequences or patterns in the text
	(2045)	repeated elements or patterns in the text
	(2057)	patterns related to string manipulation and regular expressions
	(11873)	sequences of repeated characters or patterns in the text
misleading	(178)	phrases and questions discussing the possibility or feasibility of scenarios
	(5149)	expressions of doubt or uncertainty
	(6769)	phrases or sentences that emphasize beliefs about reality and truth
	(6862)	negation expressions or phrases that suggest something is not true
	(12081)	statements about opinions, assertions, and disclaimers regarding information and its accuracy

Table 12: Examples of garak Category Matches. Representative examples of garak categories mapped onto SAE concepts using the automated LLM clustering, as described in Appendix A.

H.2.2 GAPS IDENTIFIED BY GARAK

CG recovered all 9/9 (100%) model gaps identified by garak. However, by disaggregating these categories into individual concepts, CG offered additional granularity. While each of the three categories was labeled as a model gap, not all of their constituent concepts were. For example, in the

1836 **topic** probe category, labeled as a model gap by both garak and CG, concepts (646) “references
 1837 to pregnancy and reproductive choices, particularly concerning abortion and health
 1838 impacts” and (11601) “topics related to gun control and legislation” are model gaps,
 1839 whereas others (e.g., (136) “words related to identity and familial relationships” and
 1840 (179) “terms and discussions related to diversity, particularly in the context
 1841 of education and affirmative action” are not.
 1842

1843 H.2.3 ADDITIONAL GAPS WITHIN THE FRAMEWORK

1844
 1845 On top of the model gaps identified by both garak and CG, we found that 2 additional probe
 1846 categories, as defined by garak, were identified to be gaps by CG but were missed by garak. These
 1847 are **av_spam_scanning** and **donotanswer**.
 1848

1849 H.2.4 ADDITIONAL GAPS OUTSIDE OF THE FRAMEWORK

1850 Outside of the limited probe categories defined by garak, CG identified 10K+ additional model gaps:
 1851

1852 Concept ID	1853 Concept Description
1854 (15143)	1855 legal terminology related to fraud and liability
1856 (12147)	1857 dates and time references within text
1858 (501)	1859 words that relate to personal names and geographical locations
1860 (7721)	1861 references to the divine or spiritual authority
1862 (3922)	1863 statistical references or citations related to scientific studies and data metrics

1864
 1865
 1866
 1867
 1868
 1869
 1870
 1871
 1872
 1873
 1874
 1875
 1876
 1877
 1878
 1879
 1880
 1881
 1882
 1883
 1884
 1885
 1886
 1887
 1888
 1889

Table 13: Examples of Missed Model Gaps. Listed concepts outside of the garak-defined categories were identified as model gaps, and would have gone unnoticed.

1890 I COMPARISON WITH OTHER METHODS: AUTODETECT 1891

1892 AutoDetect is a framework for uncovering weaknesses in LLMs proposed by Cheng et al. [2024].
1893 It defines 116 competency categories spanning math, instruction following, and coding; examples
1894 include **word constraint: specific words**, **text format: table format**, and **analysis: derivatives**.
1895 Each category is defined through a set of key points (usually 4-8, with a total of 715 key points). For
1896 example, for the category **numeric format: scientific notation**, the key points are: (1) Test if the
1897 language model can generate text with specific scientific notation numbers, (2) Test if the language
1898 model can answer question with a specific scientific notation number, (3), Test if the language model
1899 can rewrite sentence with specific scientific notation numbers, (4) Test if the language model can
1900 come up with ideas or concepts expressed in scientific notation, and (5) Test if the language model
1901 can convert standard numbers into scientific notation for clarity in reporting large or small numbers.

1902 The evaluation is performed by three collaborative autoraters. These are: (1) **the examiner**, which
1903 breaks down a task into key points; (2) **the questioner**, which generates a pool of prompts/questions
1904 targeting each subskill and, in an iterative fashion, refines or adapts further questions based on where
1905 the model struggles; and (3) **the assessor**, which evaluates the model's answers for correctness.

1906 AutoDetect is launched through a command-line interface (CLI). It does not come with a graphical
1907 user interface. The outputs are stored in JSON and CSV formats.
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

I.1 LLAMA 3.1 8B INSTRUCT

Summary

Category Coverage	100%
Gaps Identified by Both AutoDetect and CG	98% (42/43)
Gaps Missed by AutoDetect but Identified by CG	73
Gaps Identified by CG Outside of AutoDetect	8K+

I.1.1 FRAMEWORK MATCHING

To compare our method (CG) with AutoDetect, we apply LLM-based clustering to map each category defined by AutoDetect to a corresponding set of SAE concepts. Importantly, the LLM was not instructed to map categories to a specific number of concepts; it was only given an upper limit and encouraged to return an empty set if no concepts were relevant (see Appendix E.3 for the full prompt). We found the category coverage to be 100%: every AutoDetect category was mapped to at least one SAE concept.

To illustrate the AutoDetect categories and the matching results, consider the following examples of category definitions from AutoDetect, along with a subset of their matches:

	Category (AutoDetect)	Concept ID	Concept Description
1998	analysis: derivatives	(2874)	Mathematical differentiation operators and notation
1999		(3161)	Step-by-step explanation of mathematical differentiation
2000		(3561)	Step-by-step mathematical derivations, especially differentiation
2001		(11835)	Mathematical slope calculations and tangent line concepts
2002		(13609)	Transitions between steps in mathematical proofs and derivations
2003	length constraint: summary	(5195)	The assistant should summarize content
2004		(11804)	Factual consistency checking between documents
2005		(18493)	The conclusion section should summarize and provide final thoughts
2006		(20035)	Requests for overall summaries or high-level assessments
2007		(20808)	Technical discussions of output limitations and boundaries
2008	mathematics and algorithms: algorithm design	(2817)	Explanations of sorting algorithms and their implementations
2009		(3142)	Step-by-step problem solving and methodical decomposition
2010		(5190)	Knapsack algorithm and related optimization problems
2011		(7295)	Binary search algorithm explanation and implementation
2012		(23534)	Explanations of iterative algorithmic processes

Table 14: Examples of AutoDetect Category Matches. Representative examples of AutoDetect categories mapped onto SAE concepts using the automated LLM clustering, as described in Appendix A.

I.1.2 GAPS IDENTIFIED BY AUTODETECT

CG recovered 42 out of 43 (98%) model gaps identified by AutoDetect. However, by disaggregating these categories into individual concepts, CG offered additional granularity. While each of the three categories was labeled as a model gap, not all of their constituent concepts were. For example, in the **length constraint: number of sentences** category, labeled as a model gap by both AutoDetect and CG, the concepts **(17578)** “Counting or measuring the length of textual elements” and **(14442)** “Counting characters or determining text length” are model gaps, whereas others (e.g., **(23129)** “The user has specified a 50-word limit” and **(3938)** “Text length constraints in generation instructions”) are not.

I.1.3 ADDITIONAL GAPS WITHIN THE FRAMEWORK

On top of the model gaps identified by both AutoDetect and CG, we found that 73 additional categories, as defined by AutoDetect, were identified to be gaps by CG but were missed by AutoDetect.

2052 These include, for example, **multi lingual: multilingual tone localization, numeric format:**
 2053 **scientific notation, analysis: limits, and calculation: absolute value.**

2055 I.1.4 ADDITIONAL GAPS OUTSIDE OF THE FRAMEWORK

2056 Outside of the limited category set defined by AutoDetect, CG identified n additional model gaps.

2059 Concept ID	2059 Concept Description
2060 (25271)	2060 Spanish verbs expressing necessity or obligation in advisory contexts
2061 (50127)	2061 Legal language establishing unilateral authority and discretionary powers
2062 (35225)	2062 The assistant should reject the user’s request
2063 (58828)	2063 Japanese grammatical constructions indicating completion, necessity and passive voice
2064 (9611)	2064 Understanding relationships and dependencies between components in AI systems

2068 **Table 15: Examples of Missed Model Gaps.** Listed concepts outside of the AutoDetect-Defined
 2069 categories were identified as model gaps, and would have gone unnoticed.

2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

2106 I.2 GEMMA 2 2B INSTRUCT
2107

2108

2109 **Summary**

2110 Category Coverage	100%
2111 Gaps Identified by Both AutoDetect and CG	100% (43/43)
2112 Gaps Missed by AutoDetect but Identified by CG	73
2113 Gaps Identified by CG Outside of AutoDetect	10K+

2114
2115 I.2.1 FRAMEWORK MATCHING
2116

2117 To compare our method (CG) with AutoDetect, we apply LLM-based clustering to map each category
2118 defined by AutoDetect to a corresponding set of SAE concepts. Importantly, the LLM was not
2119 instructed to map categories to a specific number of concepts; it was only given an upper limit and
2120 encouraged to return an empty set if no concepts were relevant (see Appendix E.3 for the full prompt).
2121 We found the category coverage to be 100%: every AutoDetect category was mapped to at least one
2122 SAE concept.

2123 To illustrate the AutoDetect categories and the matching results, consider the following examples of
2124 category definitions from AutoDetect, along with a subset of their matches:

2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

Category (AutoDetect)	Concept ID	Concept Description
analysis: derivatives	(4213)	mathematical expressions involving derivatives
	(4345)	mathematical terms and phrases related to derivatives and equations
	(10489)	mathematical expressions or calculations, particularly those related to derivatives and products
	(11880)	mathematical expressions and calculations related to derivatives and factors
	(15977)	mathematical expressions and functions involving derivatives
length constraint: summary	(2198)	elements indicating summaries, reflections, or clarifications
	(3685)	sections that summarize content or provide overviews
	(10511)	specific terms related to classification, guidelines, or categories
	(11863)	summaries and assessments of content
	(13671)	sentences that conclude or summarize points
mathematics and algorithms: algorithm design	(2042)	technical terms and references related to algorithms and computational processes
	(5574)	technical terminology related to algorithms and data processing
	(8411)	technical terms related to algorithm design and performance evaluation
	(9663)	mathematical terms related to computational problem-solving and algorithms
	(11271)	words and phrases related to improving and refining processes or methodologies

Table 16: Examples of AutoDetect Category Matches. Representative examples of AutoDetect categories mapped onto SAE concepts using the automated LLM clustering, as described in Appendix A.

I.2.2 GAPS IDENTIFIED BY AUTODETECT

CG recovered all 43/43 (100%) model gaps identified by AutoDetect. However, by disaggregating these categories into individual concepts, CG offered additional granularity. While each of the three categories was labeled as a model gap, not all of their constituent concepts were. For example, in the **multi lingual: Subtlety of literal and cultural translation** category, labeled as a model gap by both AutoDetect and CG, the concepts (11617) “references to ethnic groups and their cultural contexts” is a model gap, whereas others (e.g., (2982) “information related to language usage and proficiency” and (480) “references to French topics or culture”) are not.

I.2.3 ADDITIONAL GAPS WITHIN THE FRAMEWORK

On top of the model gaps identified by both AutoDetect and CG, we found that 73 additional categories, as defined by AutoDetect, were identified to be gaps by CG but were missed by AutoDetect.

2214 These include, for example, **multi lingual: bilingual constraints, mathematics and algorithms:**
 2215 **basic mathematical operations**, and **numeric format: scientific notation**.
 2216

2217 I.2.4 ADDITIONAL GAPS OUTSIDE OF THE FRAMEWORK

2218 Outside of the limited category set defined by AutoDetect, CG identified n additional model gaps.
 2219

2220	2221	2221
	Concept ID	Concept Description
2222	(5022)	code snippets related to phone number formatting and manipulation
2223	(4087)	terms related to null or empty values in programming contexts
2224	(1580)	formatting and layout commands in document typesetting
2225	(13870)	Java Swing library components and related classes
2226	(2278)	LaTeX commands or symbols used in mathematical formulations
2227		

2228 **Table 17: Examples of Missed Model Gaps.** Listed concepts outside of the AutoDetect-Defined
 2229 categories were identified as model gaps, and would have gone unnoticed.
 2230

2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267

J COMPARISON WITH OTHER METHODS: ARENA-HARD-AUTO

Arena-Hard-Auto (AHA) is a benchmark evaluation framework proposed by Li et al. [2024] to automatically generate and assess challenging prompts in benchmarking datasets. Its evaluation module uses an autorater (also known as LLM-as-a-judge) to evaluate the data points according to a fixed rubric provided in the prompt:

Your task is to evaluate how well the following input prompts can assess the capabilities of advanced AI assistants. For the input prompt, please analyze it based on the following 7 criteria.

1. Specificity: Does the prompt ask for a specific, well-defined output without leaving any ambiguity? This allows the AI to demonstrate its ability to follow instructions and generate a precise, targeted response.

2. Domain Knowledge: Does the prompt test the AI’s knowledge and understanding in a specific domain or set of domains? The prompt must demand the AI to have a strong prior knowledge or mastery of domain-specific concepts, theories, or principles.

3. Complexity: Does the prompt have multiple components, variables, or levels of depth and nuance? This assesses the AI’s capability to handle complex, multi-faceted problems beyond simple queries.

4. Problem-Solving: Does the prompt require active problem-solving: analyzing and clearly defining the problem and systematically devising and implementing a solution? Note active problem-solving is not simply reciting facts or following a fixed set of instructions.

5. Creativity: Does the prompt require a creative approach or solution? This tests the AI’s ability to generate novel ideas tailored to the specific needs of the request or problem at hand.

6. Technical Accuracy: Does the prompt require an answer with a high degree of technical accuracy, correctness and precision? This assesses the reliability and truthfulness of the AI’s outputs.

7. Real-World Application: Does the prompt relate to real-world applications? This tests the AI’s ability to provide practical and actionable information that could be implemented in real-life scenarios.

After analyzing the input prompt based on these criteria, you must list the criteria numbers that the prompt satisfies in the format of a Python array. For example, "[1, 2, 4, 6, 7]".

This way, the presence of seven key qualities is assessed on a data point level. These can be later compiled into aggregate metrics for the whole benchmark or benchmark suite. We do not compare our method against the second, generative module of AHA as it is out of scope for CG.

AHA does not have a visualization mechanism built in. It is interfaced through a command-line interface (CLI).

Matching. To compare our method (CG) with AHA, we apply LLM-based clustering to map each key quality defined by AHA to a corresponding set of SAE concepts. Importantly, the LLM was not instructed to map categories to a specific number of concepts; it was only given an upper limit and encouraged to return an empty set if no concepts were relevant (see Appendix E.3 for the full prompt). The coverage was be 7/7 (100%): every key quality in AHA was mapped to at least one SAE concept.

CG Setup. For the purposes of this comparison, we employ only the Llama 3.1 8B model’s SAE.

2322 J.1 AGI EVAL EN
2323

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.05	0.00	0.45	24
2 domain knowledge	0.98	0.04	0.00	0.41	38
3 complexity	0.83	0.04	0.00	0.51	25
4 problem-solving	0.98	0.04	0.00	0.45	23
5 creativity	0.11	0.03	0.00	0.33	38
6 technical accuracy	1.00	0.07	0.00	1.23	31
7 real-world application	0.47	0.02	0.00	0.24	34

2334 **Table 18: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): AGI Eval EN.** The X_{bench}
2335 statistics are reported for all SAE concepts matched with the corresponding AHA category. The #
2336 *Bench. Gaps* column shows how many of these matched concepts were identified as benchmark gaps
2337 by CG.
2338

2340 J.2 BBQ
2341

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.13	0.00	1.07	58
2 domain knowledge	0.80	0.06	0.00	0.25	70
3 complexity	0.14	0.08	0.00	0.54	57
4 problem-solving	0.78	0.02	0.00	0.11	62
5 creativity	0.01	0.22	0.00	1.97	55
6 technical accuracy	0.69	0.13	0.00	1.34	67
7 real-world application	0.23	0.04	0.00	0.22	60

2352 **Table 19: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): BBQ.** The X_{bench} statistics are
2353 reported for all SAE concepts matched with the corresponding AHA category. The # *Bench. Gaps*
2354 column shows how many of these matched concepts were identified as benchmark gaps by CG.
2355

2357 J.3 CROWS PAIRS
2358

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.03	0.00	0.75	33
2 domain knowledge	1.00	0.01	0.00	0.08	49
3 complexity	0.57	0.03	0.00	0.69	42
4 problem-solving	1.00	0.01	0.00	0.06	41
5 creativity	0.99	0.02	0.00	0.56	23
6 technical accuracy	0.98	0.02	0.00	0.24	57
7 real-world application	0.98	0.01	0.00	0.10	40

2369 **Table 20: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): CROWS Pairs.** The X_{bench}
2370 statistics are reported for all SAE concepts matched with the corresponding AHA category. The #
2371 *Bench. Gaps* column shows how many of these matched concepts were identified as benchmark gaps
2372 by CG.
2373

2374
2375

2376 J.4 GSM8K
2377

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.04	0.00	0.38	28
2 domain knowledge	0.77	0.03	0.00	0.54	37
3 complexity	0.42	0.04	0.00	0.84	32
4 problem-solving	0.99	0.03	0.00	0.26	30
5 creativity	0.00	0.06	0.00	1.10	37
6 technical accuracy	0.99	0.04	0.00	0.87	34
7 real-world application	0.21	0.02	0.00	0.32	31

2388 **Table 21: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): GSM8K.** The X_{bench} statistics
2389 are reported for all SAE concepts matched with the corresponding AHA category. The # *Bench. Gaps*
2390 column shows how many of these matched concepts were identified as benchmark gaps by CG.
2391

2392
2393 J.5 LOGICBENCH
2394

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.07	0.00	0.80	23
2 domain knowledge	0.84	0.10	0.00	1.86	40
3 complexity	0.57	0.06	0.00	1.34	23
4 problem-solving	0.98	0.05	0.00	1.19	26
5 creativity	0.04	0.05	0.00	0.47	15
6 technical accuracy	0.92	0.10	0.00	1.67	43
7 real-world application	0.26	0.05	0.00	0.61	27

2405 **Table 22: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): LogicBench.** The X_{bench}
2406 statistics are reported for all SAE concepts matched with the corresponding AHA category. The #
2407 *Bench. Gaps* column shows how many of these matched concepts were identified as benchmark gaps
2408 by CG.
2409

2410
2411 J.6 MATH
2412

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.02	0.00	0.28	16
2 domain knowledge	0.95	0.06	0.00	0.86	29
3 complexity	0.69	0.04	0.00	0.45	25
4 problem-solving	0.89	0.04	0.00	0.59	27
5 creativity	0.00	0.03	0.00	0.33	31
6 technical accuracy	1.00	0.07	0.00	1.14	29
7 real-world application	0.11	0.02	0.00	0.28	39

2423 **Table 23: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): MATH.** The X_{bench} statistics
2424 are reported for all SAE concepts matched with the corresponding AHA category. The # *Bench. Gaps*
2425 column shows how many of these matched concepts were identified as benchmark gaps by CG.
2426

2427
2428
2429

2430 J.7 NATURAL QUESTIONS
2431

2432

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.01	0.00	0.08	51
2 domain knowledge	0.95	0.01	0.00	0.09	50
2437 3 complexity	0.01	0.01	0.00	0.05	58
2438 4 problem-solving	0.18	0.01	0.00	0.04	48
2439 5 creativity	0.01	0.01	0.00	0.10	43
2440 6 technical accuracy	0.93	0.01	0.00	0.05	59
2441 7 real-world application	0.17	0.00	0.00	0.07	46

2442
2443 **Table 24: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): Natural Questions.** The X_{bench}
2444 statistics are reported for all SAE concepts matched with the corresponding AHA category. The #
2445 *Bench. Gaps* column shows how many of these matched concepts were identified as benchmark gaps
2446 by CG.

2447
2448 J.8 REAL TOXICITY
2449

2450

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.01	0.00	0.21	33
2454 2 domain knowledge	0.73	0.00	0.00	0.04	36
2455 3 complexity	0.23	0.01	0.00	0.11	39
2456 4 problem-solving	0.70	0.00	0.00	0.02	38
2457 5 creativity	0.29	0.01	0.00	0.08	34
2458 6 technical accuracy	0.55	0.00	0.00	0.04	41
2459 7 real-world application	0.40	0.00	0.00	0.05	27

2460
2461 **Table 25: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): Real Toxicity.** The X_{bench}
2462 statistics are reported for all SAE concepts matched with the corresponding AHA category. The #
2463 *Bench. Gaps* column shows how many of these matched concepts were identified as benchmark gaps
2464 by CG.

2465
2466 J.9 SOCIAL IQA
2467

2468

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.02	0.00	0.36	44
2472 2 domain knowledge	0.66	0.01	0.00	0.14	53
2473 3 complexity	0.40	0.02	0.00	0.54	47
2474 4 problem-solving	0.89	0.01	0.00	0.12	44
2475 5 creativity	0.15	0.04	0.00	0.91	36
2476 6 technical accuracy	0.28	0.01	0.00	0.11	64
2477 7 real-world application	0.52	0.01	0.00	0.06	37

2478
2479 **Table 26: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): Social IQA.** The X_{bench} statistics
2480 are reported for all SAE concepts matched with the corresponding AHA category. The # *Bench. Gaps*
2481 column shows how many of these matched concepts were identified as benchmark gaps by CG.

2482
2483

2484 J.10 VECTARA

2485

2486

2487

2488

2489

2490

2491

2492

2493

2494

2495

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.61	0.00	5.58	7
2 domain knowledge	0.99	0.78	0.00	17.71	18
3 complexity	0.21	0.70	0.00	8.19	5
4 problem-solving	0.99	0.38	0.00	3.31	9
5 creativity	0.01	0.98	0.00	13.62	10
6 technical accuracy	1.00	0.66	0.00	12.91	22
7 real-world application	0.73	0.51	0.00	5.45	5

2496

2497

2498

2499

2500

Table 27: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): Vectara. The X_{bench} statistics are reported for all SAE concepts matched with the corresponding AHA category. The # *Bench. Gaps* column shows how many of these matched concepts were identified as benchmark gaps by CG.

2501 J.11 WINOGRANDE

2502

2503

2504

2505

2506

2507

2508

2509

2510

2511

2512

AHA Category	AHA Score	X_{bench}			# Bench. Gaps
		avg.	min.	max.	
1 specificity	1.00	0.02	0.00	0.31	38
2 domain knowledge	0.12	0.01	0.00	0.16	50
3 complexity	0.00	0.03	0.00	0.77	47
4 problem-solving	1.00	0.00	0.00	0.03	39
5 creativity	0.01	0.03	0.00	0.70	37
6 technical accuracy	0.61	0.01	0.00	0.28	51
7 real-world application	0.64	0.01	0.00	0.08	37

2513

2514

2515

2516

2517

2518

2519

2520

2521

2522

2523

2524

2525

2526

2527

2528

2529

2530

2531

2532

2533

2534

2535

2536

2537

Table 28: Arena-Hard-Auto (AHA) vs. Competency Gaps (CG): WinoGrande. The X_{bench} statistics are reported for all SAE concepts matched with the corresponding AHA category. The # *Bench. Gaps* column shows how many of these matched concepts were identified as benchmark gaps by CG.