

SCALING SPARSE AUTOENCODER CIRCUITS FOR IN-CONTEXT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Sparse autoencoders (SAEs) are a popular tool for interpreting large language model activations, but their utility in addressing open questions in interpretability remains unclear. In this work, we demonstrate their effectiveness by using SAEs to deepen our understanding of the mechanism behind in-context learning (ICL). We identify abstract SAE features that encode the model’s knowledge of which task to execute and whose latent vectors causally induce the task zero-shot. This aligns with prior work showing that ICL is mediated by task vectors. We further demonstrate that these task vectors are well approximated by a sparse sum of SAE latents, including these task-execution features. To explore the ICL mechanism, we adapt the sparse feature circuits methodology of Marks et al. (2024) to work for the much larger Gemma-1 2B model, with 30 times as many parameters, and to the more complex task of ICL. Through circuit finding, we discover task-detecting features with corresponding SAE latents that activate earlier in the prompt, that detect when tasks have been performed. They are causally linked with task-executing features through attention layer and MLP.

1 INTRODUCTION

Sparse autoencoders (SAEs; Ng (2011); Bricken et al. (2023); Cunningham et al. (2023)) are a promising method for interpreting large language model (LLM) activations. However, the full potential of SAEs in interpretability research remains to be explored, since most recent SAE research either i) interprets a single SAE’s features rather than the model’s computation as a whole (Bricken et al., 2023), or ii) does high-level interventions in the model, but does not interpret the downstream computation impacted by the interventions Templeton et al. (2024b). In this work, we address these limitations by interpreting in-context learning (ICL), a widely-studied phenomenon in LLMs. In brief, we show that SAEs enable a) the discovery of novel circuit components (**task-detection features**; Section 4.2) and b) making existing interpretations of ICL more precise, by e.g. decomposing task vectors (Todd et al., 2024; Hendel et al., 2023) into **task-execution features** (Section 3).

In-context learning (ICL; Brown et al. (2020)) is a fundamental capability of large language models that allows them to adapt to new tasks without fine-tuning. ICL is a significantly more complex and important task than other behaviors commonly studied in circuit analysis (such as IOI in Wang et al. (2022) and Kissane et al. (2024), or subject-verb agreement and Bias-in-Bios in Marks et al. (2024)). Recent work by Todd et al. (2024) and Hendel et al. (2023) has introduced the concept of task vectors, to study ICL in a simple setting, which we follow throughout this paper.¹ In short, task vectors are internal representations of tasks formed by language models during the processing of few-shot prompts, such as “hot → cold, big → small, fast → slow”. These vectors can be extracted and added into different LLM forward passes to induce task performance 0-shot, i.e. make LLMs predict that “slow” follows “fast →” without explicit context. Section 2.3 provides a full introduction.

To identify **task-execution features**, we decomposed task vectors using SAEs. To achieve this, we needed to go beyond existing methods for solving the classical dictionary problem of decomposing a vector into a sparse sum of dictionary vectors (Elad, 2010). To do this, we developed a bespoke method for LLMs we call the TASK VECTOR CLEANING (TVC) algorithm. By running the TVC

¹Task vectors (Hendel et al., 2023) are also called “function vectors” (Todd et al., 2024), but we use “task vectors” throughout this paper for consistency.

algorithm, we found **task-execution features**: features that can partially replace task vectors by themselves alone and have highly interpretable max activating token patterns. We validate the causal relevance of these task features through a series of steering experiments on tasks, spanning several categories like translation or factual recall. The experiments demonstrate that identified task features encode crucial information about task execution, are causally implicated in the model’s ICL capabilities and can play the same role as task vectors.

To find **task-detection features**, we adapted the Sparse Feature Circuits (SFC) methodology of Marks et al. (2024) to work on the more complex ICL task and the larger Gemma-1 2B model (Gemma Team, 2024). This adaptation allowed us to discover and analyze the subgraph of key SAE latents involved in ICL, providing a more comprehensive view of the ICL circuit. We found **task-detection features** with SFC: features that play a crucial role in identifying the specific task being performed earlier in the prompt. Task-detection features are tightly connected with task-execution features through attention, as part of the whole ICL circuit.

Our findings not only advance our understanding of ICL mechanisms but also demonstrate the potential of SAEs as a powerful tool for interpretability research on larger language models. By unifying the task vectors view with SAEs and uncovering two of the most important causally implicated feature families behind ICL, we pave the way for future work to control and monitor ICL further, to improve either safety or capabilities of models.

Our main contributions are as follows:

1. We demonstrate that SAEs can be effectively used to explain mechanisms behind a complex set of ICL tasks in a larger model (Gemma-1 2B), which has 10-35x more parameters than prior models typically studied at this depth in comparable, circuits-style mechanistic interpretability research (Wang et al., 2022; Marks et al., 2024). We show that causal circuit finding algorithms and SFC specifically straightforwardly scale up to larger models and SAEs with different architectures Appendix B.
2. We identify two core bottlenecks in the ICL circuit – **task-detection features** and **task-execution features** (see Appendix C, F, G) – and study their interactions (Section 3.2). This provides new insights into how LLMs process and execute ICL tasks. Specifically, we discover task-detection features that identify the task being performed earlier in the prompt, which are then moved by attention heads to trigger task-executing features.
3. We present a novel task vector cleaning method (Section 3.1) that decomposes task vectors into a small set of mostly task-relevant features, enabling more precise analysis of ICL mechanisms, and important linear directions in LLMs in general.

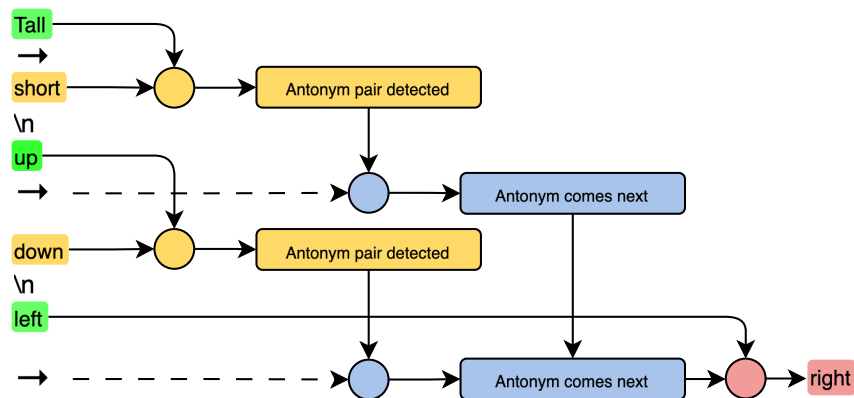


Figure 1: A diagram of the in-context learning circuit, showing task detection features (yellow) causing task execution features (blue) which cause the model to output the antonym (left → right). A more concrete circuit, along with texts these features activate on, can be seen in Figure 7.

2 BACKGROUND

2.1 SPARSE AUTOENCODERS (SAEs)

Sparse autoencoders (SAEs) are neural networks designed to learn efficient representations of data by enforcing sparsity in the hidden layer activations (Elad, 2010). In the context of language model interpretability, SAEs are used to decompose the high-dimensional activations of language models into more interpretable features (Cunningham et al., 2023; Bricken et al., 2023). The basic idea behind SAEs is to train a neural network to reconstruct its input while constraining the hidden layer to have sparse activations. This process typically involves:

1. An encoder that maps the input to a sparse hidden representation.

The SAE encoder typically has a linear encoder, a pre-activation bias and a post-encoder nonlinearity:

$$\mathbf{f}(\mathbf{x}) = \sigma(\mathbf{W}_{\text{enc}}\mathbf{x} + \mathbf{b}_{\text{enc}}) \quad (1)$$

In the JumpReLU SAE formulation (Rajamanoharan et al., 2024b), which we use², the activation function is JumpReLU:

$$\text{JumpReLU}_{\theta}(z) := zH(z - \theta) \quad (2)$$

where θ is a learnable parameter for each SAE latent and H is the Heaviside step function.

2. A decoder that reconstructs the input from this sparse representation.

The decoder is a simple affine projection in most formulations. The rows of the decoder are typically constrained to have unit norm with constrained optimization (Marks et al., 2024) or a custom loss penalty (Conerly et al., 2024).

$$\hat{\mathbf{x}}(\mathbf{f}) = \mathbf{W}_{\text{dec}}\mathbf{f} + \mathbf{b}_{\text{dec}} \quad (3)$$

3. A loss task that balances reconstruction accuracy with sparsity.

Typically, the L_1 penalty on activations is used (Bricken et al., 2023) with some modifications (Rajamanoharan et al., 2024a; Conerly et al., 2024), although there are alternatives: Rajamanoharan et al., 2024b; Farrell, 2024; Riggs & Brinkman, 2024.

In our work, we train SAEs on residual stream activations and attention outputs, and also train transcoders³ on MLP layers, all of which use the improved Gated SAE architecture (Rajamanoharan et al., 2024a).

2.2 SPARSE FEATURE CIRCUITS

Sparse Feature Circuits (SFCs) are a methodology introduced by Marks et al. (2024) to identify and analyze causal subgraphs of sparse autoencoder features that explain specific model behaviors. This approach combines the interpretability benefits of SAEs with causal analysis to uncover the mechanisms underlying language model behavior. The SFC methodology involves several key steps:

1. Decomposing model activations into sparse features using SAEs
2. Calculating the Indirect Effect (IE, Pearl (2001)) of each feature on the target behavior
3. Identifying a set of causally relevant features based on IE thresholds
4. Constructing a circuit by analyzing the connections between these features

²We technically use Gated SAEs (Rajamanoharan et al., 2024a) and convert them to JumpReLU SAEs (Rajamanoharan et al., 2024b) using the procedure outline in the Gated SAEs paper (see Appendix B).

³Transcoders are a modification of SAEs that take MLP input and convert it into MLP output instead of trying to reconstruct the residual stream.

The IE of a model component is measured by intervening on that component and observing the change in the model’s output. For a component a and a metric m , the IE is defined using do-calculus (Pearl, 2009) as in Marks et al. (2024) as:

$$\text{IE}(m; a) = m(x|\text{do}(a = a')) - m(x) \quad (4)$$

Where $m(x|\text{do}(a = a'))$ represents the value of the metric when we intervene to set the value of component a to a' , and $m(x)$ is the original value of the metric. In practice, attribution patching (Syed et al., 2023) is used to approximate IE, allowing for efficient computation across many components simultaneously.

SFC is described in detail in (Marks et al., 2024). We describe our modifications in Appendix E.

2.3 TASK VECTORS

Continuing from the high-level description in Section 1, task vectors were independently discovered by Hendel et al. (2023) and Todd et al. (2024). The key idea behind task vectors is that they capture the essence of a task demonstrated in a few-shot prompt, allowing the model to apply this learned task to new inputs without explicit fine-tuning. Task vectors have several important properties:

1. They can be extracted from the model’s hidden states after processing ICL prompts.
2. When added to the model’s activations in a zero-shot setting, they can induce task performance without explicit context.
3. They appear to encode abstract task information, independent of specific input-output examples.

To illustrate the concept, consider the following simple prompt for an antonym task, where boxes represent distinct tokens:

BOS	Follow	the	pattern	:	\n
hot	→	cold	\n		
big	→	small	\n		
fast	→	slow			

Example 1: All token types in an example input: prompt, input, arrow, output, newline (target tokens for calculating the loss on included)

In this case, the task vector would encode the abstract notion of “find the antonym” rather than specific word pairs. Task vectors are typically collected by averaging the residual stream of “→” tokens at a specific layer across multiple ICL prompts for a given task. This averaged representation can then be used to study the model’s internal task representations and to manipulate its behavior in zero-shot settings. We perform our analysis on the datasets derived from Todd et al. (2024). Details can be found in Appendix A.

3 DISCOVERING TASK-EXECUTION FEATURES

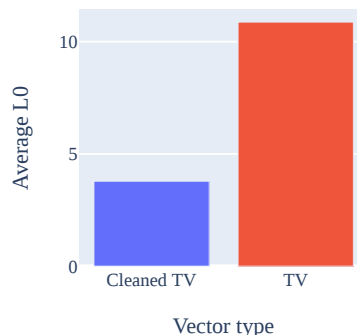
3.1 DECOMPOSING TASK VECTORS

To gain a deeper understanding of task vectors, we attempted to decompose them using sparse autoencoders (SAEs). However, several of our initial naive approaches faced significant challenges. Firstly, direct SAE reconstruction, i.e. passing the task vector as input to the SAE, produced noisy results with more than 10 non-zero SAE features on average on layers of interest⁴, most of which were irrelevant to the task. Moreover, this reconstruction noticeably reduced the vector’s performance.

⁴Layers where steering with task vectors decreased loss significantly (Figure 2a). We found 3-5 interpretable features. Our cleaning algorithm can usually trim down the number to 2-4. The usual residual SAE L0 is around 44, as highlighted in the Figure 2b



(a) The effect on the model’s loss by steering with different kinds of reconstructed task vectors, at each layer. We see that cleaning performs similarly to the original task vector until layer 14.



(b) Average L0 for cleaned task vectors vs. original task vectors at layer 12 (which corresponds to the elbow in Figure 2a).

These issues partly arose because task vectors are out-of-distribution inputs for SAEs, as they aggregate information from different residual streams rather than representing a single one.

We then explored inference-time optimization (ITO) (Smith, 2024) as an alternative. However, this method also failed to reconstruct task vectors using a low number of SAE features while maintaining high performance.

Given these observations, we developed a novel method called task vector cleaning. Our approach involves extracting task vectors from few-shot prompts for various tasks, reconstructing these vectors using trained SAEs, fine-tuning the SAE reconstruction weights to minimize negative log-likelihood loss on zero-shot prompts for the same task, and applying L_1 regularization during fine-tuning to promote sparsity. This approach allows us to maintain the task vector performance while reducing the amount of active SAE features to less than 4 on average (Figure 2b). The algorithm details can be found in Appendix D.

We compare it with the four baselines: original task vectors, naive SAE reconstruction, ITO with target L0 norm set to 5 and ITO with target L0 set to 20. To compare them, we steer the zero-shot prompt using the reconstructed task vector and calculate relative log-likelihood loss change. We then average it across all tasks. Layer-wise comparison results can be found on Figure 2a. [A more thorough evaluation across different models and SAE widths is present in Appendix D.1.](#)

Using this method, we broke down task vectors into a small set of features. Many of these features were easy to interpret and clearly related to the task at hand. We found a particularly interesting group of features, which we called “task features.” These task features have two key characteristics:

1. They activate when the model encounters examples of the relevant task in normal text.
2. In these encounters they activate on the token just before the task is completed.

For instance, imagine an antonym task feature processing the phrase “hot and cold.” It would activate on the token “and,” suggesting that the model expects an antonym to follow. This tells us that the model recognizes it’s dealing with an antonym pair before seeing the complete pair. [Appendix I contains examples of these features along with their max activating examples.](#)

3.2 STEERING EXPERIMENTS

To validate the causal relevance of our decomposed task features, we conducted a series of steering experiments. To observe the features’ impact on task performance across different contexts and model layers.

The experiments were performed on the dataset of diverse tasks taken from Todd et al. (2024). We first extracted relevant task features using our cleaning algorithm. Then steered the zero-shot prompt using them and calculated relative loss improvement, normalizing and clipping it after that. Further details and additional experiments that include other models can be found in Appendix F.

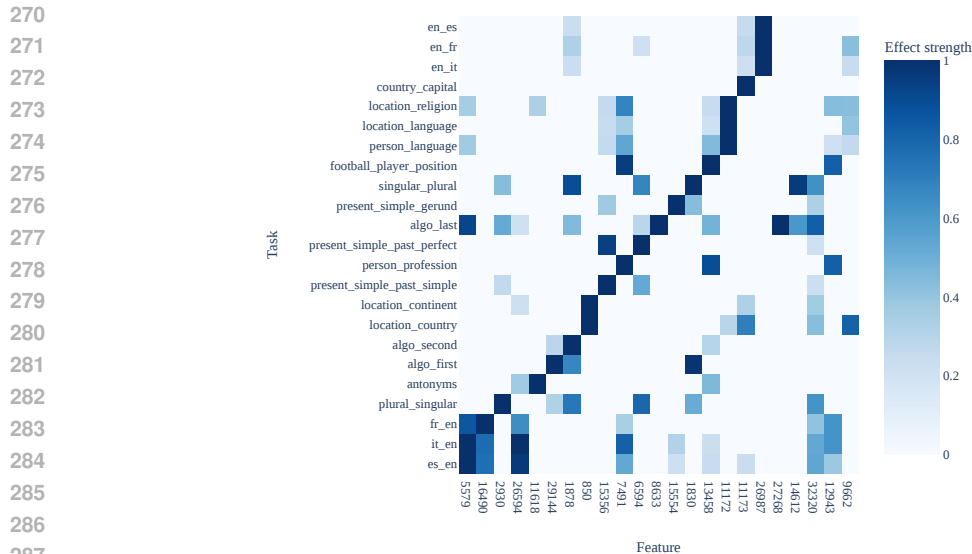


Figure 3: Heatmap showing the effect of steering with individual task-execution features for each task. We see that most features boost exactly one task, with a few exceptions for similar tasks like translating to English. **Full and unfiltered versions of the heatmap are available in Appendix F.**

Figure 3 shows a heatmap of steering results for each pair of task and task-relevant feature. Higher values indicate greater improvement in the loss after steering. It can be seen that most tasks have a single feature with a high effect on them, and this feature generally does not significantly affect unrelated tasks. Another notable detail is that features from related tasks (like the translation group) at least partially affect all tasks within the group.

We have manually examined the features with the highest effect and found that their maximum activating dataset examples align with their hypothesized role in the ICL circuit. Interestingly, we observed that translation-to-English tasks all share a generic English-to-foreign task execution feature, thus requiring an additional language encoding feature for complete task encoding. This shared feature suggests a common mechanism for translation tasks, with language-specific information encoded separately. Max activating examples of the most interpretable features are present in Appendix I.

4 APPLYING SFC TO ICL

After identifying task-executing features through our task vector analysis, we sought to further expand our understanding of the in-context learning (ICL) circuit. To this end, we decided to apply the Sparse Feature Circuits (SFC) methodology (Marks et al., 2024) to the Gemma-1 2B model. However, due to the increased complexity of ICL tasks and the larger model size, the original SFC approach did not work out of the box. We had to implement several key modifications to address the challenges we encountered.

4.1 OUR MODIFICATIONS

4.1.1 TOKEN POSITION CATEGORIZATION AND FEATURE AGGREGATION

We modified the SFC approach to better handle the structured nature of ICL prompts. Instead of treating each SAE feature as a separate node, we categorized token positions into the following groups:

- Prompt: The initial instruction tokens (e.g., “Follow the pattern:”)
- Input: The last token before each arrow in an example pair
- Arrow: The arrow token itself (“→”)

- **Output:** The last token before each newline in an example pair
- **Newline:** The newline token
- **Extra:** Any tokens not covered by the above categories (e.g., in multi-token inputs or outputs)

Each pair of an SAE feature and a token type was assigned its own graph node. The effects of the feature were aggregated across all tokens of the corresponding type. This categorization allowed us to evaluate how features affect all tokens within the same category, separating features based on their role in the ICL circuit. It also enabled us to selectively disable parts of the circuit for one task while testing another, verifying the task specificity of the identified circuits.

4.1.2 LOSS FUNCTION MODIFICATION

An ICL prompt can be viewed as an (x, y) pair, where x represents the entire prompt except for the last pair’s output, and y represents this output. The original SFC paper suggested using the log probabilities of y conditioned on x for such datasets. However, this approach often resulted in task-relevant features having high negative IEs on other example pairs in the prompt. This was likely due to the circuit’s effect on those pairs being lost to either diminishing gradients in backpropagation or because copying circuits were much more relevant to predicting the last pair. By considering all pairs except the first one, we amplified the effect of the task-solving circuit relative to the numerous cloning circuits that activate due to the repetitive nature of ICL prompts.

4.1.3 SFC EVALUATION

To evaluate the quality of our SFC modification, we conducted a series of ablation experiments across the same dataset of ICL tasks. Our primary metric for evaluation was faithfulness, which measures how much of the original task performance is maintained after ablating specific features. We calculated faithfulness using the following formula:

$$F(M) = \frac{M - M_a}{M_n - M_a} \quad (5)$$

Where M is the current metric (loss), M_a is the fully ablated model metric, and M_n is the non-ablated model metric.

We evaluated the impact of ablating features for one task on the performance of all other tasks. Specifically, we ablated the highest Indirect Effect (IE) nodes first, continuing until we reached a faithfulness of 0.5 for the target task. **Faithfulness of 0.5 corresponds to half of the original performance, i.e. a significantly destructive ablation for the target task.** This approach allowed us to assess both the specificity of the discovered circuits and their impact on related tasks. Our analysis revealed that it is possible to significantly reduce faithfulness by disabling just several hundred nodes. Furthermore, we found that we could reduce the number of active nodes to less than a thousand while keeping the performance almost intact. Extra details and faithfulness/completeness charts can be found in Appendix E.

Figure 4 presents a heatmap showing the change in faithfulness for various tasks when ablating the highest IE nodes for a single task. Several key observations can be made from this visualization:

- **Task Specificity:** Ablating most tasks does not significantly impact the performance of others, indicating that the discovered circuits are largely task-specific. This suggests that there are no common high-IE ICL-specific nodes across tasks.
- **Related Task Effects:** Tasks are grouped into categories, and we observe that ablation of related tasks has a higher effect on all tasks within the same group. This is visible as squares along the diagonal, particularly noticeable in the translation group.
- **Performance Improvement:** For some tasks, we observe that faithfulness rises well above 1.0 after ablation of other tasks. We hypothesize that this occurs because we reduce the confusion of the model by removing irrelevant execution paths.

It’s worth noting that we excluded the **person_profession** and **football_player_position** tasks from Figure 4 due to the very small difference between their fully-ablated and non-ablated losses. This resulted in highly unstable faithfulness calculations for these tasks. We attribute this small difference

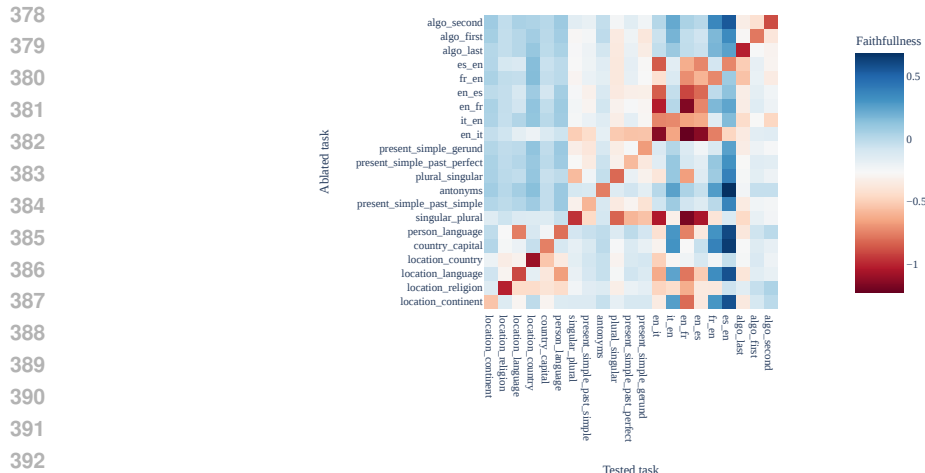


Figure 4: We study how useful the most important nodes on task A are for performance on task B. Specifically, we ablate the most important features for task A (the ablated task on the y -axis) so that faithfulness reduces by 0.5, and measure how much faithful reduces on another task B (the tested task on the x -axis).

partially to our modified loss function, as we found that calculating the loss only from the last pair results in a higher loss difference.

4.2 TASK-DETECTION FEATURES

Our modified SFC analysis revealed a second crucial component of the ICL mechanism: task-detection features. These features activate on instances of a complete task in the training data, specifically on the token that completes the task, **contrary to executors that activate right before them**. Both task-detection and task-executing features showed high Indirect Effects (IEs) in the extracted sparse feature circuits, with task detection features connected to task execution features through attention output and transcoder nodes. We applied our task vector cleaning algorithm to extract task-detection features, identifying layer 11 as optimal for steering, preceding the layer 12 task-executing features. Details can be found in Appendix G. We observed similar patterns as with task-executing features:

- **Task specificity:** Features strongly affect their corresponding tasks with minimal impact on unrelated tasks.
- **Related task effects:** Features from related tasks show some cross-task influence.

To evaluate the causal connection between task-detection features and task-execution features, we selected the most relevant detection and execution pairs based on steering effects and confirmed that their max activating patterns aligned with their hypothesized circuit roles. We then ablated detection directions while fixing attention patterns and measured the decrease in execution activations. Figure 6 presents the results.

The results of our causal connection analysis reveal several key insights. First, we observe strong causal connections between most task-detection and their corresponding task-executing features, supporting our hypothesis about their roles in the ICL circuit. Second, we note significant interconnectivity among translation tasks, suggesting shared circuitry for this group of related tasks. Interestingly, two tasks (**person_profession** and **present_simple_gerund**) showed unexpectedly weak connections between their detection and execution features, warranting further investigation.

5 RELATED WORK

Mechanistic Interpretability Olah et al. (2020) defines a framing for mechanistic interpretability in terms of *features* and *circuits*. It claims that neural network latent spaces have directions in

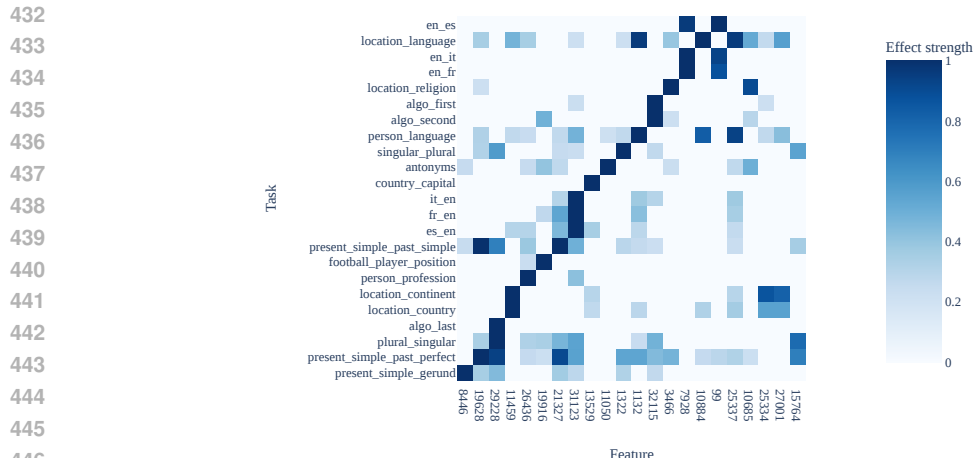


Figure 5: Heatmap showing the effect of steering with the task-detection feature most relevant to each task, on every task. We see that task detection features are typically specific to the task, with exceptions for similar tasks.

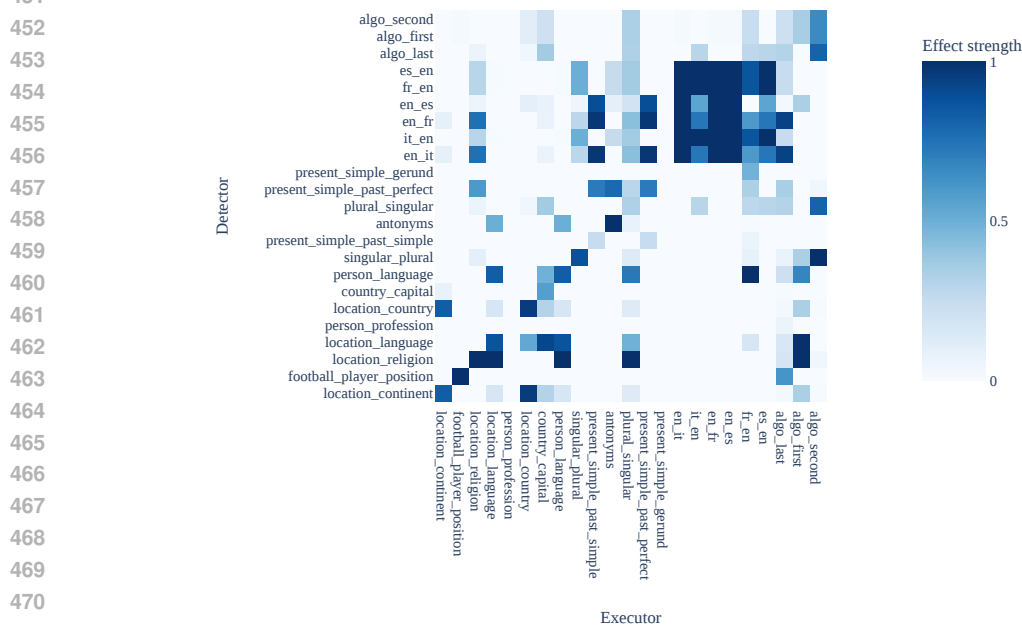


Figure 6: Heatmap showing the causal effect of the top task-detection features of each task, on the activation of the top task-executing features for every task. Averaged across all initial non-zero activations in all tasks.

them called features that correspond to meaningful variables. These features interact through model components sparsely to form circuits: interpretable computation subgraphs relevant to particular tasks. These circuits can be found through manual inspection in vision models (Cammarata et al., 2020). In language models, they can be found through manual patching (Wang et al., 2022; Hanna et al., 2023; Lieberum et al., 2023; Chan et al., 2022) or automated circuit discovery (Conmy et al. (2023); Syed et al. (2023); Bhaskar et al. (2024), though see Miller et al. (2024)). Marks et al. (2024) extends this research area to use **Sparse Autoencoders**, as discussed below.

In-Context Learning (ICL) ICL was first introduced in Brown et al. (2020) and refers to models learning to perform tasks from prompt information at test-time. There is a large area of research

486 studying its applications (Dong et al., 2024), high-level mechanisms (Min et al., 2022) and limitations
 487 (Peng et al., 2023). Elhage et al. (2021) and Olsson et al. (2022) find *induction heads* partly
 488 responsible for in-context learning. However, since these attention heads do more than just induction
 489 (Goldowsky-Dill et al., 2023), and are not sufficient for complex task-following, induction heads
 490 alone cannot explain ICL. Anil et al. (2024, Appendix G) proposes a mechanistic hypothesis for
 491 an aspect of simple in-context task behavior. Hendel et al. (2023) and Todd et al. (2024) find that
 492 simple in-context learning tasks create strong directions in the residual stream adding which makes it
 493 possible for a network to perform tasks zero-shot, but does not explain how task vectors form nor
 494 what interpretable components the task vectors are composed of. A more detailed discussion can be
 495 found in Appendix H

496 **Sparse Autoencoders** A major roadblock to mechanistic interpretability research is superposition
 497 (Elhage et al., 2022b), where the interpretable units of neural network do not tend to align with the
 498 basis directions (e.g. neurons). Sparse autoencoders (Ng, 2011; Bricken et al., 2023) are one method
 499 of addressing this roadblock, and multiple works since proposed improvements to SAE training
 500 (Rajamanoharan et al., 2024b; Bussmann et al., 2024; Braun et al., 2024; Gao et al., 2024; Templeton
 501 et al., 2024b), and we use several more in our work (Rajamanoharan et al., 2024a; Adam Jermyn,
 502 2024; Conerly et al., 2024). Cunningham et al. (2023), building on Bills et al. (2023), apply Conmy
 503 et al. (2023) to find circuits in small language models. Marks et al. (2024) adapt Syed et al. (2023) in
 504 the SAE basis to find circuits and address a practical bias reduction problem. Kissane et al. (2024)
 505 apply a slightly different automated SAE algorithm (similar to ours in that it operates on single
 506 prompts) to IOI (Wang et al., 2022), using SAEs on the attention layer outputs and residual stream.
 507 Dunefsky et al. (2024) introduce *transcoders* (which are also briefly discussed in Templeton et al.
 508 (2024a) and Li et al. (2023)) to simplify analysis of circuits involving MLPs. We build on their work
 509 and train transcoders as part of our suite of Gemma-1 SAEs.

510 6 CONCLUSION

511 **Limitations** Our work focused on the simple task vector setting to study ICL (Section 2.3), which
 512 does not capture all ways that ICL is used in practice (generally involving far more tokens and
 513 open-ended tasks). We also only interpreted Gemma-1 2B, and therefore other LLM architectures
 514 or model sizes could lead to different results (though this is unlikely, since task vectors exist across
 515 models (Todd et al., 2024)). **Finally, the complexity of the task studied meant our interpretations have
 516 some approximation error: attention heads matter for the detection-execution connection, but the
 517 succeeding MLP is necessary to capture the full effect (Section 4.2). This means that our explanation
 518 needs to include moving parts aside from task-detection attention output features. It is possible to
 519 model the effects of the MLP through transcoder features, but we leave that for future work.**

520 **Future Work** Future work could extend SFC methods to work on more than a band of layers in
 521 the middle of the model (Section 2.2). Since there are many features corresponding to individual
 522 input tokens and output predictions (due to the three stages of inference in LLMs; Elhage et al.
 523 (2022a); Lad et al. (2024)), this will require further adaptation of the SFC methodology. Moreover,
 524 our multiple contributions will hopefully spur lots of further work that finds new tasks to interpret or
 525 explanations in greater depth than prior work, as discussed in our concluding paragraph below.

526 To summarise our work: we use SAEs to explain in context learning in greater detail than any
 527 prior mechanistic interpretability work. This provides strong evidence that Sparse Autoencoders are
 528 valuable circuit analysis tools, and the innovations developed: TVC (Section 3.1), SFC improvements
 529 (Section 2.2) and an SAE training codebase in JAX with open SAE weights (Section 7) are likely to
 530 help enable lots of other SAE research to tackle more ambitious tasks and larger models.

531 7 REPRODUCIBILITY STATEMENT

532 We are committed to fostering reproducibility and advancing research in the field of mechanistic
 533 interpretability. To support this goal, we plan to release the following resources upon successful
 534 acceptance of this paper:

- 535 1. Two JAX libraries optimized for TPU:

- 540 • A library for Sparse Autoencoder (SAE) training
- 541 • A library for SAE inference and model analysis, built upon Penzai with our custom
- 542 Llama and Gemma ports
- 543 2. A full suite of SAEs for Gemma 2B, along with a dataset of their max activating examples
- 544 3. Two custom dashboards used in our analysis:
- 545 • A dashboard for browsing max activating examples
- 546 • An interactive dashboard for exploring extracted Sparse Feature Circuits (SFC)
- 547

548 These resources will enable researchers to replicate our experiments, extend our work, and conduct
 549 their own investigations using our tools and methodologies. The release of our custom dashboards will
 550 provide additional transparency and facilitate deeper exploration of our results. Due to the complexity
 551 of our infrastructure, we are only sharing anonymized versions of our analysis, cleaning, and SFC
 552 scripts, which still require our JAX libraries to run. We hope that reviewers will find this, along with
 553 the detailed methodologies described in the paper, to be sufficient evidence of reproducibility.
 554

555 REFERENCES

- 556
- 557 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen
 558 Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko,
 559 Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong
 560 Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai,
 561 Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg,
 562 Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J.
 563 Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin,
 564 Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev
 565 Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui
 566 Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo,
 567 Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes,
 568 Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-Becker,
 569 Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa,
 570 Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael
 571 Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song,
 572 Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan
 573 Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping
 574 Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali
 575 Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong
 576 Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and
 577 Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone, 2024.
 URL <https://arxiv.org/abs/2404.14219>.
- 578 Adly Templeton Adam Jermyn. Ghost grads: An improvement on resampling, 2024.
 579 URL [https://transformer-circuits.pub/2024/jan-update/index.html#](https://transformer-circuits.pub/2024/jan-update/index.html#dict-learning-resampling)
 580 [dict-learning-resampling](https://transformer-circuits.pub/2024/jan-update/index.html#dict-learning-resampling).
- 581 Cem Anil, Esin Durmus, Mrinank Sharma, Joe Benton, Sandipan Kundu, Joshua Batson, Nina
 582 Rimskey, Meg Tong, Jesse Mu, Daniel Ford, et al. Many-shot jailbreaking. 2024.
- 583 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:
 584 Provable in-context learning with in-context algorithm selection.
- 585 Hritik Bansal, Karthik Gopalakrishnan, Saket Dingliwal, Sravan Bodapati, Katrin Kirchhoff, and Dan
 586 Roth. Rethinking the role of scale for in-context learning: An interpretability-based case study at
 587 66 billion scale. URL <http://arxiv.org/abs/2212.09095>.
- 588 Adithya Bhaskar, Alexander Wettig, Dan Friedman, and Danqi Chen. Finding transformer circuits
 589 with edge pruning, 2024. URL <https://arxiv.org/abs/2406.16778>.
- 590 Steven Bills, Nick Cammarata, Dan Mossing, et al. Language models can explain neu-
 591 rons in language models. [https://openaipublic.blob.core.windows.net/](https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html)
 592 [neuron-explainer/paper/index.html](https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html), 2023.

- 594 Joseph Bloom. Open source sparse autoencoders for all residual stream layers of gpt2-small,
595 2024. URL [https://www.alignmentforum.org/posts/f9EgflSurAiqRjySD/
596 open-source-sparse-autoencoders-for-all-residual-stream](https://www.alignmentforum.org/posts/f9EgflSurAiqRjySD/open-source-sparse-autoencoders-for-all-residual-stream).
597
- 598 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
599 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and
600 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL
601 <http://github.com/jax-ml/jax>.
- 602 Dan Braun, Jordan Taylor, Nicholas Goldowsky-Dill, and Lee Sharkey. Identifying functionally
603 important features with end-to-end sparse dictionary learning, 2024. URL [https://arxiv.
604 org/abs/2405.12241](https://arxiv.org/abs/2405.12241).
605
- 606 Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick
607 Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec,
608 Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina
609 Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and
610 Christopher Olah. Towards monosemanticity: Decomposing language models with dictionary
611 learning. *Transformer Circuits Thread*, 2023. [https://transformer-circuits.pub/
612 2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 613 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhari-
614 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
615 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,
616 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler,
617 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam Mc-
618 Candlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-
619 shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pp.
620 1877–1901, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/
621 1457c0d6bfc4967418bfb8ac142f64a-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html).
- 622 Bart Bussmann, Patrick Leask, and Neel Nanda. Batchtopk: A simple improvement for topk-saes,
623 2024. URL [https://www.alignmentforum.org/posts/Nkx6yWZNbAsfvic98/
624 batchtopk-a-simple-improvement-for-topk-saes](https://www.alignmentforum.org/posts/Nkx6yWZNbAsfvic98/batchtopk-a-simple-improvement-for-topk-saes).
- 625 Nick Cammarata, Shan Carter, Gabriel Goh, Chris Olah, et al. Thread: Circuits. *Distill*, 2020. doi:
626 10.23915/distill.00024. <https://distill.pub/2020/circuits>.
627
- 628 Lawrence Chan, Adria Garriga-Alonso, Nix Goldowsky-Dill, Ryan Greenblatt, Jenny
629 Nitishinskaya, Ansh Radhakrishnan, Buck Shlegeris, and Nate Thomas. Causal
630 scrubbing: A method for rigorously testing interpretability hypotheses, 2022.
631 URL [https://www.alignmentforum.org/posts/JvZhhzycHu2Yd57RN/
632 causal-scrubbing-a-method-for-rigorously-testing](https://www.alignmentforum.org/posts/JvZhhzycHu2Yd57RN/causal-scrubbing-a-method-for-rigorously-testing).
- 633
- 634 Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Unveiling induction heads: Provable
635 training dynamics and feature learning in transformers. URL [http://arxiv.org/abs/
636 2409.10559](http://arxiv.org/abs/2409.10559).
- 637 Tom Conerly, Adly Templeton, Trenton Bricken, Jonathan Marcus, and Tom Henighan. Up-
638 date on how we train saes, 2024. URL [https://transformer-circuits.pub/2024/
639 april-update/index.html#training-saes](https://transformer-circuits.pub/2024/april-update/index.html#training-saes).
- 640
- 641 Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, et al. Towards automated circuit discovery
642 for mechanistic interpretability. In *Proceedings of NeurIPS*, 2023.
- 643
- 644 Hoagy Cunningham, Aidan Ewart, Logan Riggs, et al. Sparse autoencoders find highly interpretable
645 features in language models, 2023. URL <https://arxiv.org/abs/2309.08600>.
- 646
- 647 Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can GPT
learn in-context? language models implicitly perform gradient descent as meta-optimizers. URL
<http://arxiv.org/abs/2212.10559>.

- 648 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu,
649 Zhiyong Wu, Tianyu Liu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. A survey on in-context
650 learning, 2024. URL <https://arxiv.org/abs/2301.00234>.
- 651
- 652 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, et al. The llama 3 herd of
653 models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- 654 Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. Transcoders find interpretable llm feature
655 circuits, 2024. URL <https://arxiv.org/abs/2406.11944>.
- 656
- 657 Michael Elad. *Sparse and Redundant Representations: From Theory to Applications in Signal*
658 *and Image Processing*. Springer, New York, 2010. ISBN 978-1-4419-7010-7. doi: 10.1007/
659 978-1-4419-7011-4.
- 660 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
661 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac
662 Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse,
663 Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A
664 mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. URL
665 <https://transformer-circuits.pub/2021/framework/index.html>.
- 666
- 667 Nelson Elhage, Tristan Hume, Catherine Olsson, Neel Nanda, Tom Henighan, Scott Johnston, Sheer
668 ElShowk, Nicholas Joseph, Nova DasSarma, Ben Mann, Danny Hernandez, Amanda Askell,
669 Kamal Ndousse, Jones, , Dawn Drain, Anna Chen, Yuntao Bai, Deep Ganguli, Liane Lovitt, Zac
670 Hatfield-Dodds, Jackson Kernion, Tom Conerly, Shauna Kravec, Stanislav Fort, Saurav Kadavath,
671 Josh Jacobson, Eli Tran-Johnson, Jared Kaplan, Jack Clark, Tom Brown, Sam McCandlish, Dario
672 Amodei, and Christopher Olah. Softmax linear units. *Transformer Circuits Thread*, 2022a.
<https://transformer-circuits.pub/2022/solu/index.html>.
- 673
- 674 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,
675 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy Models of Superposition.
676 *arXiv preprint arXiv:2209.10652*, 2022b.
- 677 Eoin Farrell. Experiments with an alternative method to promote sparsity in sparse autoen-
678 coders, 2024. URL [https://www.lesswrong.com/posts/cYA3ePxy8JQ8aajo8B/
679 experiments-with-an-alternative-method-to-promote-sparsity](https://www.lesswrong.com/posts/cYA3ePxy8JQ8aajo8B/experiments-with-an-alternative-method-to-promote-sparsity).
- 680
- 681 Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever,
682 Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders, 2024. URL <https://arxiv.org/abs/2406.04093>.
- 683
- 684 Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn
685 in-context? a case study of simple function classes.
- 686
- 687 Gemma Team. Gemma: Open models based on gemini research and technology, 2024. URL
688 <https://arxiv.org/abs/2403.08295>.
- 689
- 690 Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. Localizing model
691 behavior with path patching, 2023. URL <https://arxiv.org/abs/2304.05969>.
- 692
- 693 Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit
694 Steiner, Dustin Li, Esin Durmus, Ethan Perez, Evan Hubinger, Kamilé Lukošiušė, Karina Nguyen,
695 Nicholas Joseph, Sam McCandlish, Jared Kaplan, and Samuel R. Bowman. Studying large
696 language model generalization with influence functions, 2023. URL [https://arxiv.org/
697 abs/2308.03296](https://arxiv.org/abs/2308.03296).
- 698
- 699 Xiaochuang Han, Daniel Simig, Todor Mihaylov, Yulia Tsvetkov, Asli Celikyilmaz, and Tianlu Wang.
700 Understanding in-context learning via supportive pretraining data. URL [http://arxiv.org/
701 abs/2306.15091](http://arxiv.org/abs/2306.15091).
- 702
- 703 Michael Hanna, Ollie Liu, and Alexandre Variengien. How does gpt-2 compute greater-than?:
704 Interpreting mathematical abilities in a pre-trained language model, 2023. URL [https://
705 arxiv.org/abs/2305.00586](https://arxiv.org/abs/2305.00586).

- 702 Roe Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors, 2023. URL
703 <https://arxiv.org/abs/2310.15916>.
704
- 705 Daniel D. Johnson. Penzai + treescope: A toolkit for interpreting, visualizing, and editing models as
706 data, 2024. URL <https://arxiv.org/abs/2408.00211>.
- 707 Patrick Kidger and Cristian Garcia. Equinox: neural networks in jax via callable pytrees and filtered
708 transformations, 2021. URL <https://arxiv.org/abs/2111.00254>.
709
- 710 Connor Kissane, Robert Krzyzanowski, Arthur Conmy, and Neel
711 Nanda. Attention output saes improve circuit analysis, 2024. URL
712 [https://www.alignmentforum.org/posts/EGvtgB7ctifzxZg6v/
713 attention-output-saes-improve-circuit-analysis](https://www.alignmentforum.org/posts/EGvtgB7ctifzxZg6v/attention-output-saes-improve-circuit-analysis).
- 714 Vedang Lad, Wes Gurnee, and Max Tegmark. The remarkable robustness of llms: Stages of inference?,
715 2024. URL <https://arxiv.org/abs/2406.19384>.
716
- 717 Max Li, Sam Marks, and Aaron Mueller. dictionary learning repository, 2023. URL [https://github.com/saprmarks/dictionary_learning?tab=readme-ov-file#
718 extra-functionalitysupported-by-this-repo](https://github.com/saprmarks/dictionary_learning?tab=readme-ov-file#extra-functionalitysupported-by-this-repo). Accessed on September 30, 2024.
719
- 720 Tom Lieberum, Matthew Rahtz, János Kramár, et al. Does circuit analysis interpretability scale?
721 evidence from multiple choice capabilities in chinchilla, 2023. URL [https://arxiv.org/
722 abs/2307.09458](https://arxiv.org/abs/2307.09458).
- 723 Tom Lieberum, Senthooan Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant
724 Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse
725 autoencoders everywhere all at once on gemma 2, 2024. URL [https://arxiv.org/abs/
726 2408.05147](https://arxiv.org/abs/2408.05147).
727
- 728 Johnny Lin. Neuronpedia: Interactive reference and tooling for analyzing neural networks, 2023.
729 URL <https://www.neuronpedia.org>. Software available from neuronpedia.org.
- 730 Arvind Mahankali, Tatsunori B. Hashimoto, and Tengyu Ma. One step of gradient descent is
731 provably the optimal in-context learner with one layer of linear self-attention. URL [http:
732 //arxiv.org/abs/2307.03576](http://arxiv.org/abs/2307.03576).
733
- 734 Samuel Marks, Can Rager, Eric J. Michaud, et al. Sparse feature circuits: Discovering and editing in-
735 terpretable causal graphs in language models. *Computing Research Repository*, arXiv:2403.19647,
736 2024. URL <https://arxiv.org/abs/2403.19647>.
- 737 Joseph Miller, Bilal Chughtai, and William Saunders. Transformer circuit faithfulness metrics are not
738 robust, 2024. URL <https://arxiv.org/abs/2407.08734>.
739
- 740 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
741 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv
742 preprint arXiv:2202.12837*, 2022. URL <https://arxiv.org/abs/2202.12837>.
743
- 744 Andrew Ng. Sparse autoencoder. *CS294A Lecture Notes*, 2011. Unpublished lecture notes.
- 745 Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.
746 Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001. [https:
747 //distill.pub/2020/circuits/zoom-in](https://distill.pub/2020/circuits/zoom-in).
- 748 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
749 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli,
750 Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane
751 Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,
752 and Chris Olah. In-context learning and induction heads, 2022. URL [https://arxiv.org/
753 abs/2209.11895](https://arxiv.org/abs/2209.11895).
- 754 Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
755 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent.
URL <http://arxiv.org/abs/2212.07677>.

- 756 Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. What in-context learning “learns” in-context:
757 Disentangling task recognition and task learning. In Anna Rogers, Jordan Boyd-Graber, and
758 Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp.
759 8298–8319. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.527.
760 URL <https://aclanthology.org/2023.findings-acl.527>.
- 761 Judea Pearl. Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty*
762 *in Artificial Intelligence, UAI’01*, pp. 411–420, San Francisco, CA, USA, 2001. Morgan Kaufmann
763 Publishers Inc. ISBN 1558608001.
- 764 Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, USA, 2nd
765 edition, 2009. ISBN 052189560X.
- 766 Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin
767 Raffel, Leandro Von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the
768 finest text data at scale, 2024. URL <https://arxiv.org/abs/2406.17557>.
- 769 Hao Peng, Xiaozhi Wang, Jianhui Chen, Weikai Li, Yunjia Qi, Zimu Wang, Zhili Wu, Kaisheng Zeng,
770 Bin Xu, Lei Hou, and Juanzi Li. When does in-context learning fall short and why? a study on
771 specification-heavy tasks, 2023. URL <https://arxiv.org/abs/2311.08993>.
- 772 Senthoooran Rajamanoharan. Improving ghost grads, 2024. URL
773 [https://www.alignmentforum.org/posts/C5KAZQib3bzzpeyrg/
774 progress-update-1-from-the-gdm-mech-interp-team-full-update#
775 Improving_ghost_grads](https://www.alignmentforum.org/posts/C5KAZQib3bzzpeyrg/progress-update-1-from-the-gdm-mech-interp-team-full-update#Improving_ghost_grads).
- 776 Senthoooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János
777 Kramár, Rohin Shah, and Neel Nanda. Improving dictionary learning with gated sparse autoen-
778 coders, 2024a. URL <https://arxiv.org/abs/2404.16014>.
- 779 Senthoooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János
780 Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse
781 autoencoders, 2024b. URL <https://arxiv.org/abs/2407.14435>.
- 782 Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the
783 emergence of non-bayesian in-context learning for regression. URL [http://arxiv.org/
784 abs/2306.15063](http://arxiv.org/abs/2306.15063).
- 785 Logan Riggs and Jannik Brinkman. Improving sae’s by sqrt()-ing l1 removing lowest activating
786 features, 2024. URL [https://www.lesswrong.com/posts/YiGs8qJ8aNBgwt2YN/
787 improving-sae-s-by-sqrt-ing-l1-and-removing-lowest](https://www.lesswrong.com/posts/YiGs8qJ8aNBgwt2YN/improving-sae-s-by-sqrt-ing-l1-and-removing-lowest).
- 788 Lingfeng Shen, Aayush Mishra, and Daniel Khashabi. Do pretrained transformers learn in-context
789 by gradient descent? URL <http://arxiv.org/abs/2310.08540>.
- 790 Chenglei Si, Dan Friedman, Nitish Joshi, Shi Feng, Danqi Chen, and He He. Measuring inductive
791 biases of in-context learning with underspecified demonstrations. In *Proceedings of the 61st*
792 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
793 11289–11310. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.632.
794 URL <https://aclanthology.org/2023.acl-long.632>.
- 795 Lewis Smith. Replacing sae encoders with inference-time optimisation, 2024. URL
796 [https://www.alignmentforum.org/posts/C5KAZQib3bzzpeyrg/
797 full-post-progress-update-1-from-the-gdm-mech-interp-team#
798 Replacing_SAE_Encoders_with_Inference_Time_Optimisation](https://www.alignmentforum.org/posts/C5KAZQib3bzzpeyrg/full-post-progress-update-1-from-the-gdm-mech-interp-team#Replacing_SAE_Encoders_with_Inference_Time_Optimisation).
- 799 Aquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated circuit
800 discovery, 2023. URL <https://arxiv.org/abs/2310.10348>.
- 801 Adly Templeton, Joshua Batson, Adam Jermyn, and Chris Olah. Predicting future activations,
802 January 2024a. URL [https://transformer-circuits.pub/2024/jan-update/
803 index.html#predict-future](https://transformer-circuits.pub/2024/jan-update/index.html#predict-future). Accessed on September 30, 2024.

- 810 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen,
811 Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L
812 Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers,
813 Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan.
814 Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Trans-*
815 *former Circuits Thread*, 2024b. URL [https://transformer-circuits.pub/2024/
816 scaling-monosemanticity/index.html](https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html).
- 817 Eric Todd, Millicent L. Li, Arnab Sen Sharma, et al. Function vectors in large language models. In
818 *Proceedings of the 2024 International Conference on Learning Representations*, 2024.
819
- 820 Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Inter-
821 pretableity in the wild: A circuit for indirect object identification in GPT-2 small, 2022. URL
822 <https://arxiv.org/abs/2211.00593>.
- 823
824 Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun.
825 Label words are anchors: An information flow perspective for understanding in-context learn-
826 ing. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Con-*
827 *ference on Empirical Methods in Natural Language Processing*, pp. 9840–9855. Associa-
828 tion for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.609. URL [https:
829 //aclanthology.org/2023.emnlp-main.609](https://aclanthology.org/2023.emnlp-main.609).
- 830
831 Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, and William Yang Wang. Large language
832 models are latent variable models: Explaining and finding good demonstrations for in-context
833 learning, 2024. URL <https://arxiv.org/abs/2301.11916>.
- 834
835 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
836 learning as implicit bayesian inference. URL <http://arxiv.org/abs/2111.02080>.
- 837
838 Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and
839 Bill Yuchen Lin. Magpie: Alignment data synthesis from scratch by prompting aligned llms with
840 nothing, 2024. URL <https://arxiv.org/abs/2406.08464>.
- 841
842 Steve Yadlowsky, Lyric Doshi, and Nilesch Tripuraneni. Pretraining data mixtures enable narrow model
843 selection capabilities in transformer models. URL <http://arxiv.org/abs/2311.00871>.
- 844

845 A MODEL AND DATASET DETAILS

846

847
848 For our experiments, we utilized the Gemma 1 2B model, a member of the Gemma family of open
849 models based on Google’s Gemini models (Gemma Team, 2024). The model’s architecture is largely
850 the same as that of Llama (Dubey et al., 2024) with the exception of tied input and output embeddings
851 and a different activation function for MLP layers, so we could reuse our infrastructure for loading
852 Llama models. We train residual and attention output SAEs as well as transcoders for layers 1-18 of
853 the model on FineWeb (Penedo et al., 2024).

854 Our dataset for circuit finding is primarily derived from the function vectors paper (Todd et al., 2024),
855 which provides a diverse set of tasks for evaluating the existence and properties of function vectors in
856 language models. We supplemented this dataset with three additional algorithmic tasks to broaden
857 the scope of our analysis:

- 858 • Extract the first element from an array of length 4
- 859 • Extract the second element from an array of length 4
- 860 • Extract the last element from an array of length 4
- 861
- 862

863 The complete list of tasks used in our experiments with task descriptions is as follows:

Task ID	Description
location_continent	Name the continent where the given landmark is located.
football_player_position	Identify the position of a given football player.
location_religion	Name the predominant religion in a given location.
location_language	State the primary language spoken in a given location.
person_profession	Identify the profession of a given person.
location_country	Name the country where a given location is situated.
country_capital	Provide the capital city of a given country.
person_language	Identify the primary language spoken by a given person.
singular_plural	Convert a singular noun to its plural form.
present_simple_past_simple	Change a verb from present simple to past simple tense.
antonyms	Provide the antonym of a given word.
plural_singular	Convert a plural noun to its singular form.
present_simple_past_perfect	Change a verb from present simple to past perfect tense.
present_simple_gerund	Convert a verb from present simple to gerund form.
en_it	Translate a word from English to Italian.
it_en	Translate a word from Italian to English.
en_fr	Translate a word from English to French.
en_es	Translate a word from English to Spanish.
fr_en	Translate a word from French to English.
es_en	Translate a word from Spanish to English.
algo_last	Extract the last element from an array of length 4.
algo_first	Extract the first element from an array of length 4.
algo_second	Extract the second element from an array of length 4.

This diverse set of tasks covers a wide range of linguistic and cognitive abilities, including geographic knowledge, language translation, grammatical transformations, and simple algorithmic operations. By using this comprehensive task set, we aimed to thoroughly investigate the in-context learning capabilities of the Gemma 1 2B model across various domains.

B SAE TRAINING

Our **Gemma 1 2B** SAEs are trained with a learning rate of $1e-3$ and Adam betas of 0.0 and 0.99 for 150M (± 100) tokens of FineWeb (Penedo et al., 2024). The methodology is overall similar to (Bloom, 2024). We initialize encoder weights orthogonally and set decoder weights to their transpose. We initialize decoder biases to 0. We use Rajamanoharan (2024)’s ghost gradients variant (ghost gradients applied to dead features only, loss multiplied by proportion of death features) with the additional modification of using softplus instead of exp for numerical stability. A feature is considered dead when its density (according to a 1000-batch buffer) is below $5e-6$ or when it hasn’t fired in 2000 steps. We use Anthropic’s input normalization and sparsity loss for Gemma 1 2B (Conerly et al., 2024). We found it to improve Gated SAE training stability. **We modified it to work with transcoders by keeping track of input and output norms separately and predicting normed outputs.**

We convert our Gated SAEs into JumpReLU SAEs after training, implementing algorithms like TVC and SFC in a unified manner for all SAEs in this format (including simple SAEs). The conversion procedure involves setting thresholds so as to replicate the effect of the gating branch. For further details, see Rajamanoharan et al. (2024b).

We use 4 v4 TPU chips running Jax (Bradbury et al., 2018) (Equinox (Kidger & Garcia, 2021)) to train our SAEs. We found that training with Huggingface’s Flax LM implementations was very slow. We reimplemented LLaMA (Dubey et al., 2024) and Gemma (Gemma Team, 2024) in Penzai (Johnson, 2024) with a custom layer-scan transformation and quantized inference kernels as well as support for loading from GGUF compressed model files. We process an average of around 4400 tokens per second, and caching LM activations is not the main bottleneck for us. For this and other reasons, we don’t do SAE sparsity coefficient sweeps so as to increase TPU utilization.

For caching, we use a distributed ring buffer which contains separate pointers on each device to allow for processing masked data. The (in-place) buffer update is in a separate JIT context. Batches are sampled randomly from the buffer for each training step.

We train our SAEs in bfloat16 precision. We found that keeping weights and scales in bfloat16 and biases in float32 performed best in terms of the amount of dead features and led to a Pareto improvement over float32 SAEs.

For training Phi 3 (Abdin et al., 2024) SAEs, we use data generated by the model unconditionally, similarly to (Xu et al., 2024)⁵. The resulting dataset we train the model on contains many math problems and is formatted as natural-seeming interaction between the user and the model.

Each SAE training run takes us about 3 hours. We trained 3 models (a residual SAE, an attention output SAE and a transcoder) for each of the 18 layers of the model. This is about 1 week of v4-8 TPU time.

Our SAEs and training code will be made public after paper acceptance.

C EXAMPLE CIRCUITS

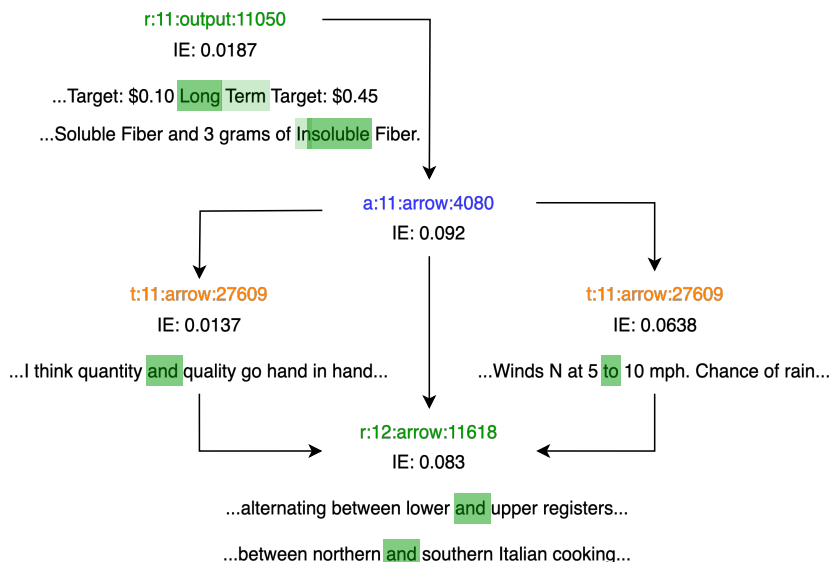


Figure 7: An example of a circuit found using our SFC variant. We focused on a subcircuit with high indirect effects. Maximum activating examples from the SAE training distribution are included.

We will release a web interface for viewing maximum activating examples and task feature circuits.

D TASK VECTOR CLEANING ALGORITHM

The task vector cleaning algorithm is a novel approach we developed to isolate task-relevant features from task vectors. Figure 8 provides an overview of this algorithm.

Our process begins with collecting residuals for task vectors using a batch of 16 and 16-shot prompts. We then calculate the SAE features for these task vectors. We explored two methods: (1) calculating feature activation and then averaging across tokens, and (2) averaging across tokens first and then calculating the task vector. They had similar performance.

The cleaning process is performed on a training batch of 24 pairs, with evaluation conducted on an additional 24 pairs. All prompts are zero-shot. An example prompt is as follows:

⁵Phi-3 is trained primarily with instruction following data, making it an aligned chat model.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

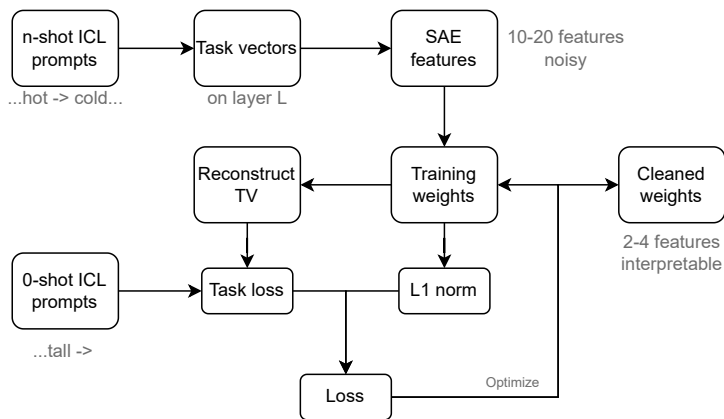


Figure 8: Overview of the task vector cleaning algorithm.

```

BOS Follow the pattern : \n
tall → short \n
...
old → young \n
hot → cold
  
```

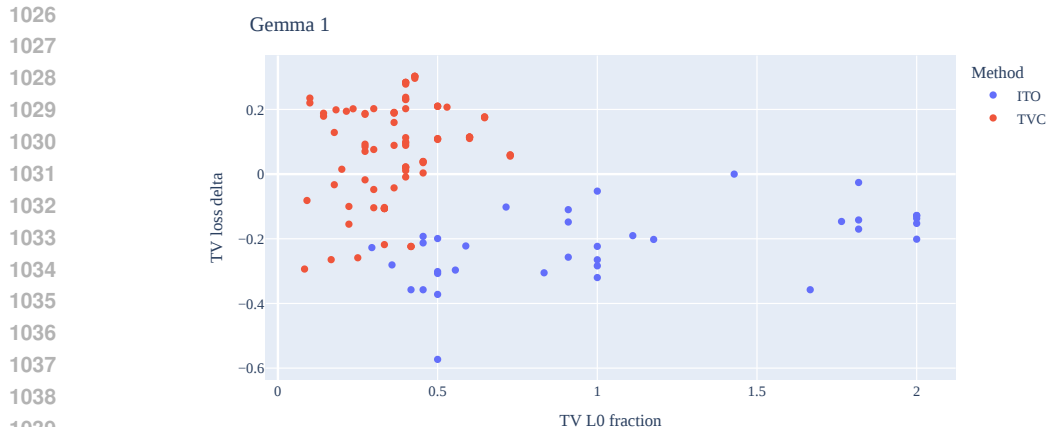
Example 2: The steered token highlighted in red. Loss is calculated on the yellow token.

The algorithm is initialized with the SAE reconstruction as a starting point. It then iteratively steers the model on the reconstruction layer and calculates the loss on the training pairs. To promote sparsity, we add the L0 norm of weights with coefficient l to the loss function. The algorithm implements early stopping when the L0 norm remains unchanged for n iterations.

```

1 def tvc_algorithm(task_vector, model, sae):
2     initial_weights = sae.encode(task_vector)
3     def tvc_loss(weights, tokens):
4         task_vector = sae.decode(weights)
5         mask = tokens == self.separator
6         model.residual_stream[layer, mask] += task_vector
7         # loss only on the "output" tokens,
8         # ignoring input and prompt tokens
9         loss = logprobs(model.logits, tokens, ...)
10        return loss + ll_coeff * ll_norm(weights)
11    weights = initial_weights.copy()
12    optimizer = adam(weights, lr=0.15)
13    last_l0, without_change = 0, 0 # early stopping
14    for _ in range(1000):
15        grad = jax.grad(tvc_loss)(weights, tokens)
16        weights = optimizer.step(grad)
17        if ll_norm(weights) != last_l0:
18            last_l0, without_change = ll_norm(weights), 0
19        elif without_change >= 50:
20            break
21    return weights
  
```

Algorithm 1: Pseudocode for Task Vector Cleaning.



1040 Figure 9: Performance of ITO and TVC across different tasks and optimization parameters compared
1041 to task vectors for Gemma 1 2B. Y-axis shows relative improvement over task vector loss, while
1042 X-axis shows the fraction of active TV features used. Metric calculation details are available in D.1
1043

1044
1045
1046 The hyperparameters l , n , and learning rate α can be fixed for a single model. We experimented with
1047 larger batch sizes but found that they did not significantly improve the quality of extracted features
1048 while substantially slowing down the algorithm due to gradient accumulation.
1049

1050 The algorithm takes different amounts of time to execute for different tasks and models. For Gemma
1051 1, it stops at 100-200 iterations, which is close to 40 seconds at 5 iterations per second.

1052 It's worth noting that we successfully applied this method to the recently released Gemma 2 2B
1053 model using the Gemma Scope SAE suite (Lieberum et al., 2024). It was also successful with the
1054 Phi-3 3B model (Abdin et al., 2024) and with our SAEs, which were trained similarly to the Gemma
1055 1 2B SAEs.
1056

1057 D.1 L_1 SWEEPS

1058

1059 To provide more details about the method's effectiveness across various models and SAE widths, we
1060 conducted L_1 coefficient sweeps with our Phi-3 and Gemma 1 2B SAEs, as well as Gemma Scope
1061 Gemma 2 2B SAEs. We chose two SAE widths for Gemma 2: 16k and 65k. We studied only the
1062 optimal task vector layer for each model: 12 for Gemma 1, 16 for Gemma 2, and 18 for Phi-3. We
1063 used a learning rate of 0.15 with the Gemma 1 2B, Phi-3 and Gemma 2 2B 65k models, and 0.3 with
1064 Gemma 2 2B 16k.
1065

1066 Figures 9, 10, 11, 12 compare TVC and ITO against original task vectors. The X-axis displays the
1067 fraction of active task vector SAE features used. The Y-axis displays the TV loss delta, calculated
1068 as $(L_{TV} - L_{Method})/L_{Zero}$, where L_{TV} is the loss from steering with the task vector, L_{Method}
1069 is the loss after it has been cleaned using the corresponding method, and L_{Zero} is the uninformed
1070 (no-steering) model loss. This metric shows improvement over the task vector relative to the loss of
1071 the uninformed model. Points were collected from all tasks using 5 different L_1 coefficient values.
1072

1073 We observe that our method often improves task vector loss and can reduce the number of active
1074 features to one-third of those in the original task vector while maintaining relatively intact performance.
1075 In contrast, ITO rarely improves the task vector loss and is almost always outperformed by TVC.
1076

1077 Figure 13 shows task-mean loss decrease (relative to no steering loss) and remaining TV features
1078 fraction plotted against L_1 sweep coefficients. We see that L_1 coefficients between 0.001 and 0.05
1079 result in relatively intact performance, while significantly reducing the amount of active SAE features.
We again notice that wider SAEs benefit more from the method.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093

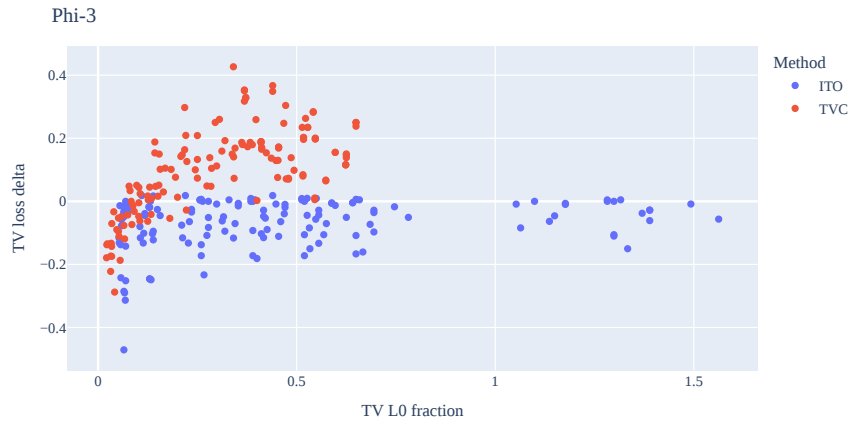


Figure 10: Performance of ITO and TVC across different tasks and optimization parameters compared to task vectors for Phi-3. Y-axis shows relative improvement over task vector loss, while X-axis shows the fraction of active TV features used. Metric calculation details are available in D.1

1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111

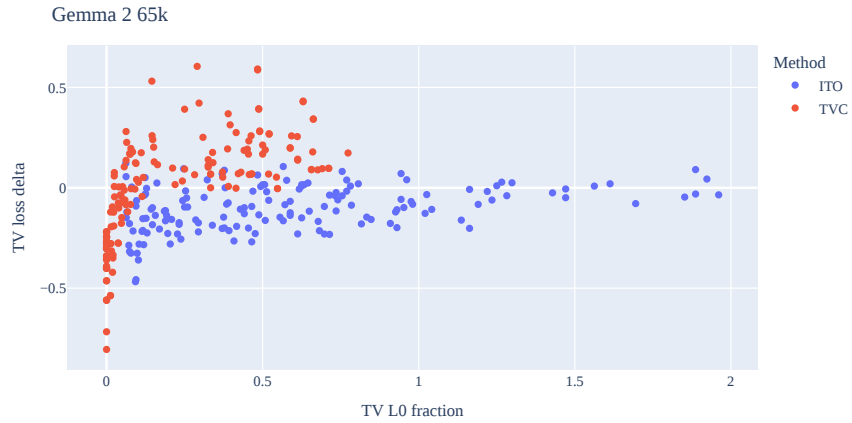


Figure 11: Performance of ITO and TVC across different tasks and optimization parameters compared to task vectors for Gemma 2 2B 65k SAEs. Y-axis shows relative improvement over task vector loss, while X-axis shows the fraction of active TV features used. Metric calculation details are available in D.1

1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129

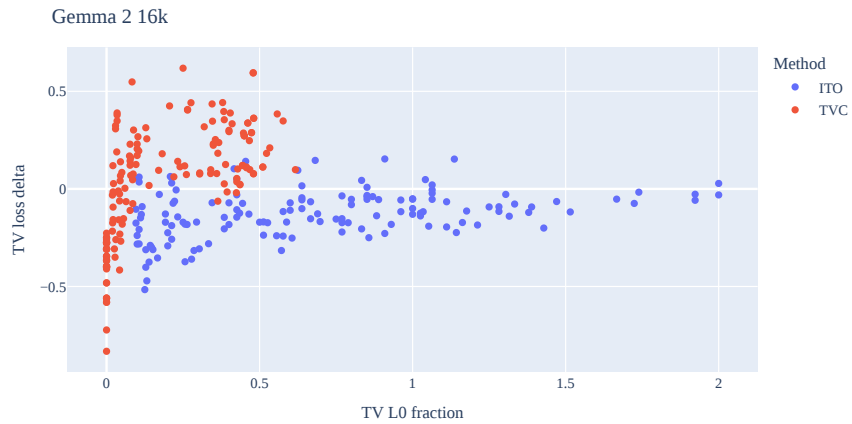


Figure 12: Performance of ITO and TVC across different tasks and optimization parameters compared to task vectors for Gemma 2 2B 16k SAEs. Y-axis shows relative improvement over task vector loss, while X-axis shows the fraction of active TV features used. Metric calculation details are available in D.1

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

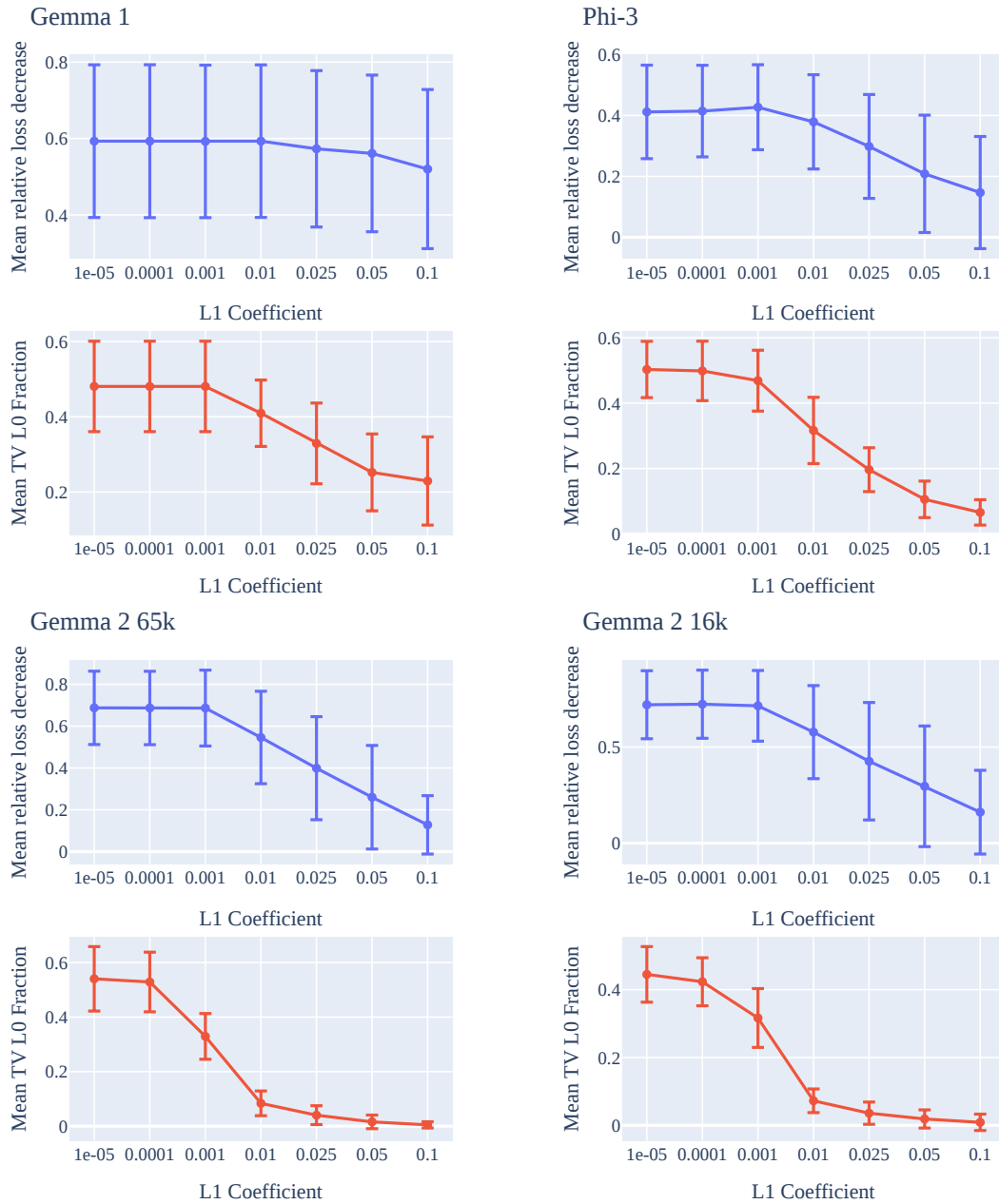


Figure 13: L_1 coefficient sweeps across different models and SAEs. All metrics are averaged across all tasks. Error bars show std of the average for each case. Metric calculation details are available in D.1.

1188 E DETAILS OF OUR SFC IMPLEMENTATION

1189

1190

1191

E.1 IMPLEMENTATION DETAILS

1192

Our implementation of circuit finding attribution patching is specialized for Jax and Penzai.

1193

1194

1195

1196

1197

We first perform a forward-backward pass on the set of prompts, collecting residuals and gradients from the metric to residuals. We collect gradients with `jax.grad` by introducing "dummy" zero-valued inputs to the metric computation function that are added to residuals to each layer. Note that we do not use SAEs during the stage.

1198

1199

1200

1201

1202

1203

1204

1205

We then perform an SAE encoding step and find the nodes (residual, attention output and transcoder SAE features and error nodes) with the highest indirect effects using manually computed gradients from the metric. After that, we find the features with the top K indirect effects for each layer and position mask and treat them as candidates for circuit edge targets. We compute gradients with respect to the metric to the values of those nodes, propagate them to "source features" up to one layer above and multiply by the values of the source features. This way, we can compute indirect effects for circuit edges and prune the initially fully-connected circuit. However, like Marks et al. (2024), we do not perform full ablation of circuit edges.

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

```

1 # resids_pre: L x N x D - the pre-residual stream at layer L
2 # resids_mid: L x N x D - the middle of the residual stream
3 # (between attention and MLP) at layer L
4 # grads_pre: L x N x D - gradients from the metric to resids_pre
5 # grads_mid: L x N x D - gradients from the metric to resids_mid
6 # all of the above are computed with a forward and backward
7 # pass without SAEs
8
9 # saes_resid: L - residual stream SAEs
10 # saes_attn: L - attention output SAEs
11 # transcoders_attn: L - transcoders predicting resids_pre[l+1]
12 # from resids_mid[l]
13
14 def indirect_effect_for_residual_node(layer):
15     sae_encoding = saes_resid[layer].encode(
16         resids_pre[layer])
17     grad_to_sae_latents = jax.vjp(
18         saes_resid[layer].decode,
19         sae_encoding
20     )(grads_pre[l])
21     return (grad_to_sae_latents * sae_encoding).sum(-1)
22
23 def indirect_effect_for_attention_node(layer):
24     sae_encoding = saes_attn[layer].encode(
25         resids_mid[layer] - resids_pre[layer])
26     grad_to_sae_latents = jax.vjp(
27         saes_attn[layer].decode,
28         sae_encoding
29     )(grads_mid[l])
30     return (grad_to_sae_latents * sae_encoding).sum(-1)
31
32 def indirect_effect_for_transcoder_node(layer):
33     sae_encoding = transcoders[layer].encode(
34         resids_mid[layer])
35     grad_to_sae_latents = jax.vjp(
36         transcoders[layer].decode,
37         sae_encoding
38     )(grads_pre[l+1])
39     return (grad_to_sae_latents * sae_encoding).sum(-1)

```

Algorithm 2: Pseudocode for Sparse Feature Circuits indirect effect calculation.

E.2 FAITHFULNESS CHARTS

Figure 14 shows averaged node trimming effect on faithfulness across all tasks. We follow methodology of Marks et al. (2024) thresholding removing nodes with low IE first. We can see that the circuits keep at least 0.8 faithfulness on average just with 1000 nodes (on layers 11-17).

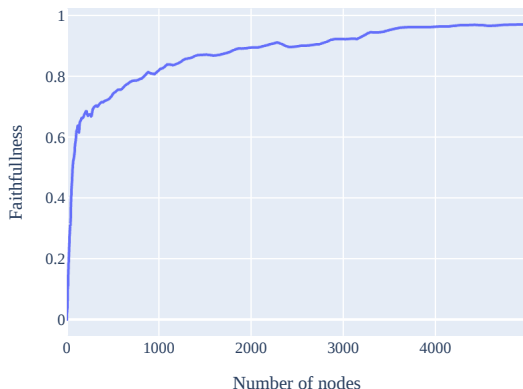


Figure 14: Average faithfulness across tasks depending on the amount of nodes left in the circuit.

Figure 15 shows averaged inverse node trimming effect on faithfulness across all tasks. Marks et al. (2024) calls this metric completeness and calculates it as faithfulness of the model just with the circuit ablated. We calculate it by thresholding nodes starting with those that have the highest IE. We can see that ablation of just even several hundred nodes have drastic impact on faithfulness (This is also for the window of layers 11-17).

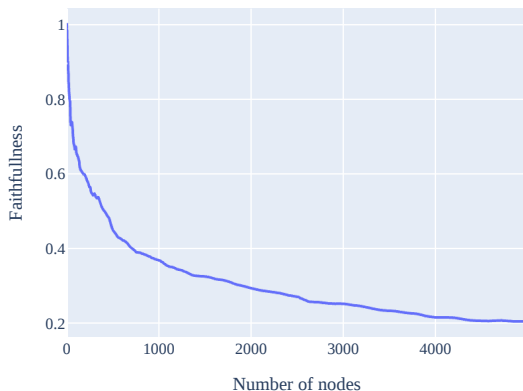


Figure 15: Average faithfulness across tasks depending on the amount of important nodes ablated from the circuit .

F STEERING WITH TASK-EXECUTION FEATURES

To evaluate the causal relevance of our identified ICL features, we conducted a series of steering experiments. Our methodology employed zero-shot prompts for task-execution features, measuring effects across a batch of 32 random pairs.

We set the target layer as 12 using Figure 2a and extracted all task-relevant features on it using our cleaning algorithm. To determine the optimal steering scale, we conducted preliminary experiments using manually identified task-executing features across all tasks. Through this process, we estab-

lished an optimal steering scale of 15, which we then applied consistently across all subsequent experiments.

For each pair of task and feature, we performed steering and measured the relative loss improvement compared to the model’s task performance on a prompt without steering. This relative improvement metric allowed us to quantify the impact of each feature on task performance.

To normalize our results and highlight the most significant effects, we applied several post-processing steps:

- We clipped the effect to be no more than 1, thus ignoring any instances of loss increase.
- We then normalized the effects for all features within the same task to be in the 0 to 1 range.
- To remove clutter and highlight important features, we set effects lower than 0.2 to 0.
- Finally, we removed features with low maximum effect across all tasks to reduce the size of the resulting diagram. The full version of this diagram is present in Figure 16.

Prompt example with the steered token highlighted in red. Loss is calculated on the yellow token:

```

    BOS Follow the pattern : \n
    hot → cold
  
```

Example 3: Task-execution steering setup. The steered token is highlighted in red and the loss is calculated on the yellow token.

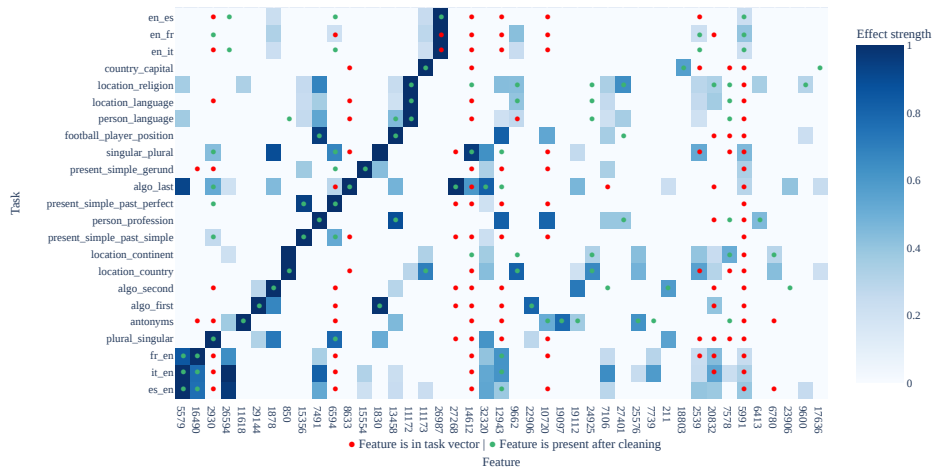


Figure 16: Full version of the heatmap in Figure 3 showing the effect of steering with individual task-execution features for each task. The features present in the task vector of the corresponding task are marked with dots. Green dots show the features that were extracted by cleaning. Red dots are features present in the original task vector. Not all original features from the task vectors are present.

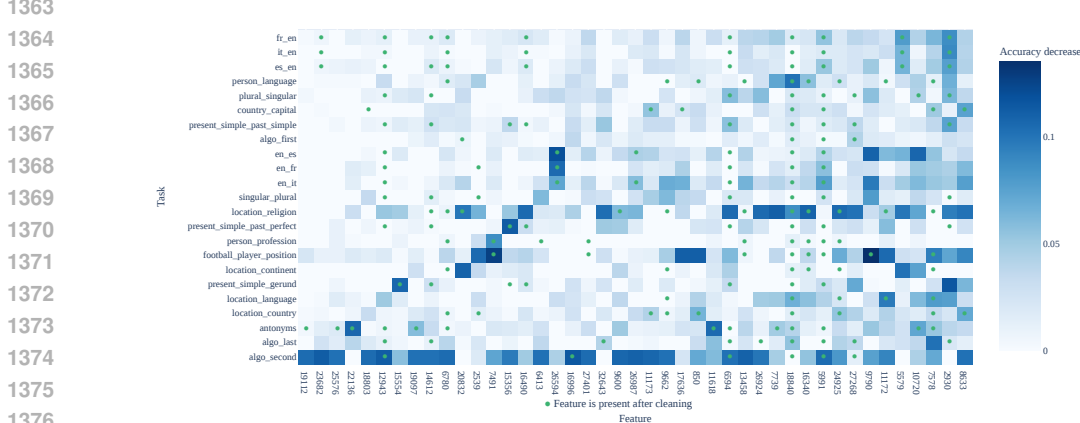
We also share the version of the Figure 16 without normalization and value clipping. It is present in Figure 18. We see that task vectors generally contain just a few task-executing features that are able to boost the task themselves. Remaining features have much weaker and less specific effects.

F.1 NEGATIVE STEERING

To further explore the effects of the executor feature, we also conducted negative steering experiments. The setup involved a batch of 16 ICL prompts, each containing 32 examples for each task. We collected all features from the cleaned task vectors for every task. Similar to positive steering, we steered with features on arrow tokens, but this time multiplying the direction by -1. Prompts this time contained several arrow tokens, and we steered on all of them simultaneously.

1350 An important distinction from positive steering is that performance degradation in negative steering
 1351 may occur due to two factors: (1) our causal intervention on the ICL circuit and (2) the steering scale
 1352 being too high. To address this, we measured accuracy across all pairs in the batch instead of loss, as
 1353 accuracy does not decrease indefinitely. We also observed that features no longer share a common
 1354 optimal scale. Consequently, for each task pair, we iterated over several scales between 1 and 30.
 1355 For each feature, we then selected a scale that reduced accuracy by at least 0.1 for at least one task.
 1356 Steering results at this scale were used for this feature across all tasks.

1357 Figure 17 displays the resulting heatmap. While we observe some degree of task specificity — and
 1358 even note that some executing features from Figure 16 have their expected effects — we also find that
 1359 negative steering exhibits significantly lower task specificity. Additionally, we observe that non-task-
 1360 specific features have a substantial impact in this experiment. This suggests that steering experiments
 1361 alone may not suffice for a comprehensive analysis of the ICL mechanism, thus reinforcing the
 1362 importance of methods such as our modification of SFC.



1377 Figure 17: Negative steering heatmap. Displays accuracy decrease after optimal scale negative
 1378 steering on full ICL prompts. Green circles show which features were present in the cleaned task
 1379 vector of the corresponding task. More details in Appendix F.1

1382 F.2 GEMMA 2 2B POSITIVE STEERING

1384 Additionally, we conducted zero-shot steering experiments with Gemma 2 2B 16k and 65k SAEs.
 1385 Contrary to Gemma 1 2B, task executors from Gemma 2 2B didn’t have a single common optimal
 1386 steering scale. Thus, we added an extra step to the experiment: for each feature and task pair we
 1387 performed steering with several scales from 30 to 300, and then selected the scale that had maximal
 1388 loss decrease on any of the tasks. We then used this scale for this feature in application to all other
 1389 tasks. Figure 19a and Figure 19b contain steering heatmaps for Gemma 2 2B 16k SAEs and Gemma
 1390 2 2B 65k SAEs respectively.

1391 We observe a relatively similar level of executor task-specificity compared to Gemma 1. One notable
 1392 difference between 16k and 65k SAEs is that 65k cleaned task vectors appear to contain more features
 1393 with a strong effect on the task. However, this may be due to the l_1 regularization coefficient being
 1394 too low.

1396 G TASK-DETECTION FEATURES

1398 For our investigation of task-detection features, we employed a methodology similar to that used for
 1399 task execution features, with a key modification. We introduced a fake pair to the prompt and focused
 1400 our steering on its output. This approach allowed us to simulate the effect of the detection features as
 1401 it happens on real prompts.

1402 Our analysis revealed that layers 10 and 11 were optimal for task detection, with performance notably
 1403 declining in subsequent layers. We selected layer 11 for our primary analysis due to its proximity

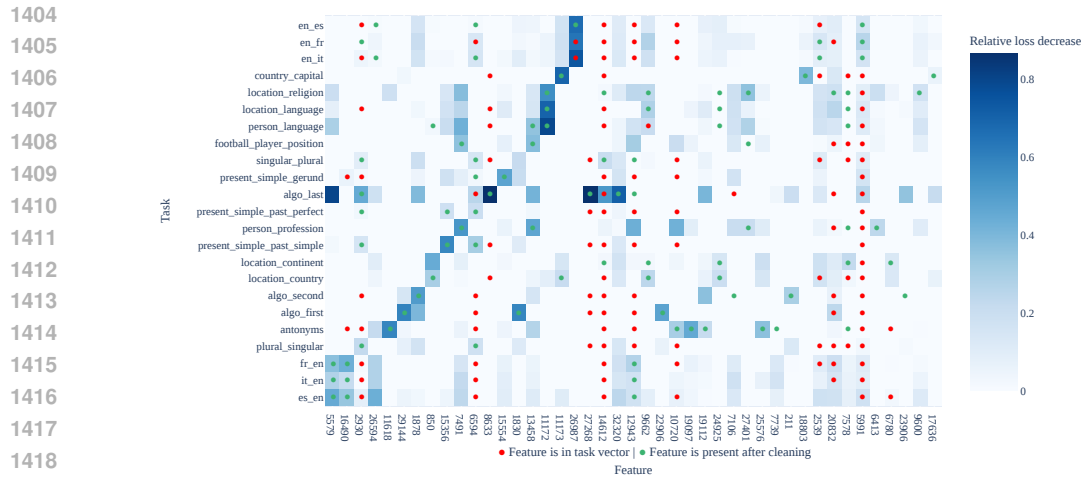


Figure 18: Unfiltered version of the heatmap in Figure 5 showing the effect of steering with individual task-execution features for each task. The features present in the task vector of the corresponding task are marked with dots. Green dots show the features that were extracted by cleaning. Red dots are the features present in the original task vector. Since the chart only contains features from cleaned task vectors, not all features from the original task vectors are present.

to layer 12, where we had previously identified the task execution features. This choice potentially facilitates a more direct examination of the interaction between detection and execution mechanisms.

The steering process for detection features followed the general methodology outlined in Appendix F, including the use of a batch of 32 random pairs, extraction of task-relevant features, and application of post-processing steps to normalize and highlight significant effects. The primary distinction lay in the application of the steering to the prompt.

This approach allowed us to create a comprehensive representation of the causal relationships between task-detection features and the model’s ability to recognize specific tasks, as visualized in Figure 5.

BOS	Follow	the	pattern	:	\n
X	→	Y	\n		
hot	→	cold			

Example 4: Task-detection steering setup. The steered token is highlighted in red and the loss is calculated on the yellow token.

H ICL INTERPRETABILITY LITERATURE REVIEW

This section will cover work on understanding ICL not mentioned in Section 5.

Raventós et al. provides evidence for two different Bayesian algorithms being learned for linear regression ICL: one for limited task distributions, and one that is alike to ridge regression. It also intriguingly shows that the two solutions lie in different basins of the loss landscape, a phase transition necessary to go from one to the other. While interesting, it is not clear if the results apply to real-world tasks.

The existence of discrete task detection and execution features hinges on the assumption that in-context learning works by classifying the task to perform and not by learning a task. Pan et al. aims to disentangle the two with a black-box approach that mixes up outputs to force the model to learn the task from scratch. Si et al. look at biases in task recognition in ambiguous examples through

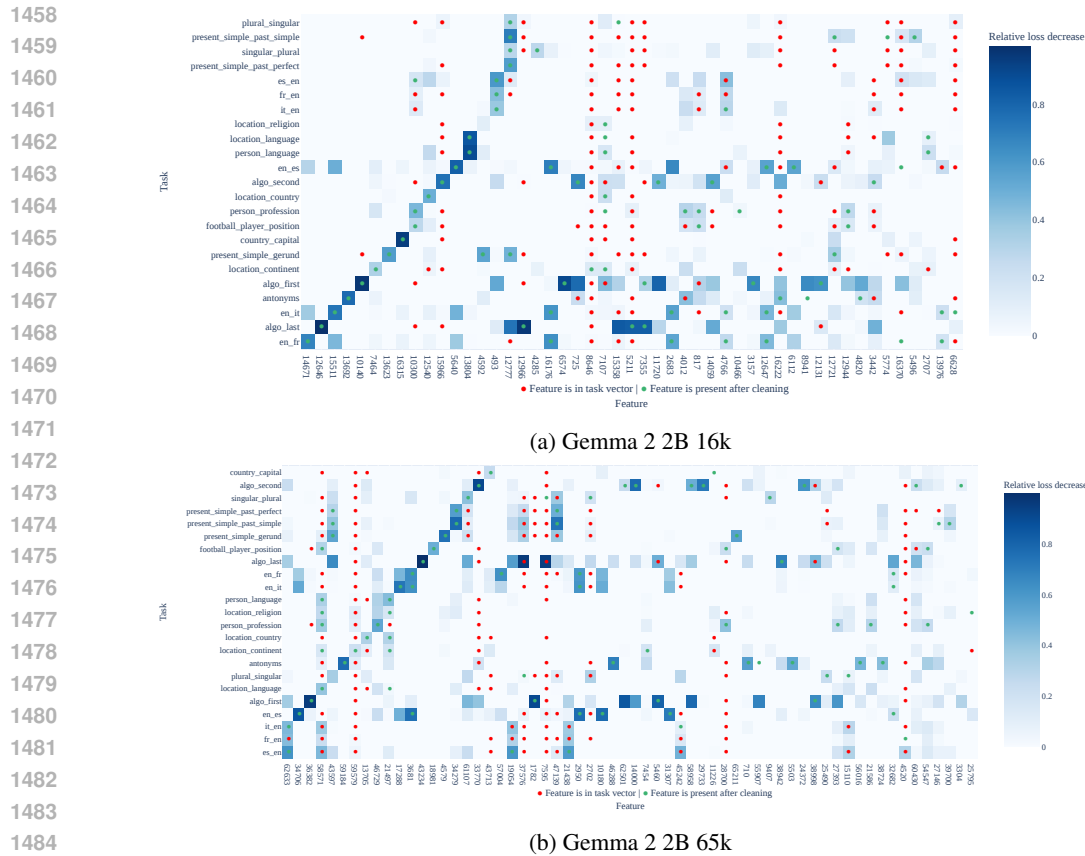


Figure 19: Unfiltered positive steering heatmap for Gemma 2 2B SAEs showing the effect of steering with individual task-execution features for each task. Steering scales were optimized for each feature. The features present in the task vector of the corresponding task are marked with dots. Green dots show the features that were extracted by cleaning. Red dots are the features present in the original task vector. Since the chart only contains features from cleaned task vectors, not all features from the original task vectors are present.

a black-box lens. We find more clear task features for some tasks than others, but do not consider whether this is linked to how common a task is in pretraining data.

Xie et al. proposes that in-context learning happens because language models aim to model a latent topic variable to predict text with long-range coherence. Wang et al. (2024) show following the two proposed steps rigorously improves results in real-world models. However, they do not endeavor to explain the behavior of non-finetuned models by looking at internal representations; instead, they aim to improve ICL performance.

Han et al. use a weight-space method to find examples in training data that promote in-context learning using a method akin to Grosse et al. (2023), producing results similar to per-token loss analyses in Olsson et al. (2022), and, similarly to the studies mentioned above, finds that those examples involve long-range coherence. Our method is also capable of finding examples in data that are similar to ICL, and we find crisp examples for many tasks being performed Appendix I.

Bansal et al. offers a deeper look into induction heads, scaling up Olsson et al. (2022) the way we scale up Marks et al. (2024). Crucially, it finds that MLPs in later layers cannot be removed while preserving ICL performance, indirectly corroborating our findings from Section 4.2. Chen et al. come up with a proof that states that gradient flow converges to a generalized version of the algorithm suggested by Olsson et al. (2022) when trained on n-gram Markov chain data.

Garg et al. studies the performance of toy models trained on in-context regression various *function classes*. Yadlowsky et al. find that Transformers trained on regression with multiple function classes

1512 have trouble combining solutions for learning those functions. Oswald et al. construct a set of weights
1513 for linear attention Transformers that reproduce updates from gradient descent and find evidence
1514 for the algorithm being represent on real models trained on toy tasks. Mahankali et al. proves that
1515 this algorithm is optimal for single-layer transformers on noisy linear regression data. Shen et al.
1516 questions the applicability of this model to real-world transformers. Bai et al. finds that transformers
1517 can switch between multiple different learning algorithms for ICL. Dai et al. find multiple similarities
1518 between changes made to model predictions from in-context learning and weight finetuning.

1519 While important, we do not consider this direction of interpreting transformers trained on regression
1520 for concrete function classes through primarily white-box techniques. Instead, we aim to focus on
1521 clear discrete tasks which are likely to have individual features.

1522 The results of Wang et al. are perhaps the most similar to our findings. The study finds "anchor tokens"
1523 responsible for aggregating semantic information, analogous to our "output tokens" (Section 2.3) and
1524 task-detection features. They tackle the full circuit responsible for ICL bottom-up and intervene on
1525 models using their understanding, improving accuracy. Like this paper, they do not deeply investigate
1526 later attention and MLP layers. This study uses SAE features to find strong linear directions on output
1527 and arrow tokens corresponding to task detection and execution respectively, offering a different
1528 perspective.

1529

1530 I MAX ACTIVATING EXAMPLES

1531

1532 This section contains max activating examples for some executor and detector features for Gemma
1533 1 2B, as described in (Bricken et al., 2023). They are computed by iterating over the training data
1534 distribution (FineWeb) and sampling activations of SAE features that fall within disjoint buckets for
1535 the activation value of span 0.5. We can observe that the degree of intuitive interpretability depends
1536 on the amount of task-similar contexts in the training data and SAE width.

1537 We also provide max activating examples for Gemma 2 2B executor features from Figures 19b and
1538 19a. These max activating examples are taken from the Neuronpedia (Lin, 2023) and are available in
1539 Figures 23 and 22.

1540

1541 Here we can notice the main difference between executors and detectors: executors mainly activate
1542 **before the task completion**, while detectors activate on the **token that completes the task**. We also
1543 found that in Gemma 1 2B detector features for some tasks were split between several token-level
1544 features (like the journalism feature in Figure 21f), and they did not create a single feature before the
1545 task executing features activated. We attribute this to the limited expressivity of the SAEs that we
1546 used.

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

1564

1565

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

st by alternating between lower and upper registers
 differences between northern and southern Italian cooking
 models, both import and domestic, Ulmer's speculations
 between fresh and traditional, casual and elegant
 spaces and both local and remote event logging. They
 globally in both tropical and temperate waters. But
 light and darkness, life and death. This assemblage
 (a) Max activating examples for the antonyms executor feature 11618.

cultural diversity and that special joie de vivre (joy of life), that Montrealer
 of creating Papel Picado Banderitas (little paper banners). Popular throughout Mexico and
 nicknames: 'la ciudad dorada' (the golden city). Salamanca is also the
 from the same root as Jihad, or struggle. In the sense that IJL
 conurbation 'tsukin jigoku', or commuter hell. Images of rail workers
 , he acted according to our Sunna (tradition), and whoever slaughtered the sacrifice before
 , and, of course, your karma (good and bad). John brings a wealth
 The 'It-sa Sicherheitsmesse' (security trade show), OWASP conference,
 Cerveceria Catalan along with a caña (draft beer) and a rosé...yes
 s Apostle! I slaughtered the Nusak (before the prayer) but I-bos>True
 the living entities; sva-artha=interest; vyatikramah=ob
 by her given name (Angelella= little angel), but called her Columba

(c) Max activating examples for the translation to English executor feature 5579.
 on came all the way from Oslo Norway for the event. This has
 emigrated to New York City from the Galicia area, in northwest Spain.
 mal. There were a few folks from Canada (British Columbia, Ontario and Quebec
 an immigrant from Bangladesh who was granted political asylum by the
 in and world number six Li Na of China. Azarenka, who is
 hood in Seattle to my college years in Boston,
 owever, batteries from rival manufacturers in the U.S. are exempt from
 re USA and four in England. Of these, only three have
 (e) Max activating examples for the prediction of city/country feature 850.

Judgment Staff (裁きの杖, Sabaki no Tsue?), also known
 Rift The Judgment Staff (裁きの杖, Sabaki no Tsue?), also
 more commonly known as the Four-Tails (四尾, Yonbi), is a
 as the Four-Tails (四尾, Yonbi), is a tailed beast sealed
 1 meters dybde er vigtige for et-årige afgr
 te vinden (inclusief een directe link naar de publicatie online als deze beschikbaar
 elders te vinden (inclusief een directe link naar de publicatie online als deze beschikbaar
 een publicatie elders te vinden (inclusief een
 Oh (侍合体シンケンオー Samurai Gattai Shinken Ō?)
 দার গোয়েন্দাগিরি
 naar de publicatie online als deze beschikbaar is in een
 atie online als deze beschikbaar is in een database op het internet). Journal id:
 Goendagiri (ফেলুদার গোয়েন্দাগি
 did, upstage-bos>Israel (ישראל נתיב) is a small yet diverse
 Agencia Española de Medicamentos y Productos Sanitarios, A
 authority (Agencia Española de Medicamentos y Productos Sanitarios, A

(b) Max activating examples for the English to foreign language translation executor feature 26987.

working, connecting with diverse people and seeking out sustain
 croll, and 2) Isolating unifying elements that transcend the indivi
 nt like landing on the moon or the discovery of DNA. The focus
 by using our search feature or by following the links above. Feel
 as Liking and Favoriting photos, but it will expire after
 spends her free time traveling and visiting exotic locations arou
 enses tighten, grabbing offensive rebounds and making putback
 er than participating in or observing or-bos>I tend to specialise i
 nother, rather than participating in or observing or-bos>I tend to
 a passenger car with plastics sheets and inhaling toxic fumes fr
 nilies when going a long distance or flying with them when we ce

(d) Max activating examples for the "next comes gerund form" executor feature 15554.

scientistb, - Rury Holman, director on behalf of the United Kingdom
 Stinton, NAR CEO Charlie Young, President/CEO
 . San Fernando Realty Dale Stinton, NAR CEO Charlie Young, President/
 director (Ja, - Philip Clarke, research fellowa, - Andrew Farmer
 Hummel, MD, Ezio Bonifacio, PHD, and Anette-G.
 lives in Bombay. Faruq Hassan, Poet and critic; teaches at Dawson Coll
 -bos>Illa Williams, President Randall Ramsay, Vice President Texas Cha

(f) Max activating examples for the person's occupation executor feature 13458.

Figure 20: Max activating examples for executor features from the Figure 3.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Target: \$0.10 **Long Term** Target: \$0.45
Soluble Fiber and 3 grams of **Insoluble** Fiber. Ground Flaxseeds are a gr
5 In. x 6 In.; **Outer** Dimensions: 7 In. x
a service: |Morning Services|**Evening** Service| Morning Worship at 8
temp: 15°C **min** temp: 11°C
Upper Zone and 75 Bottles in **Lower** Zone - Read More... The
page and 30% viewing the **right half** "apple's decision
integration is performed first, followed by the **quantitative** combination.
access to the content item. Returns **FALSE** if the current<bos>|Oracle@
. As we alternate between defensive positions and **offensive** positions, v

(a) Max activating examples for the antonyms detector feature 11050.

Superficie Lunare (Composizione)" **Lunar surface**: composition), executed prior to
hardline group Tawhid wal Jihad (**Monotheism** and Holy War).
hardline group Tawhid wal Jihad (**Monotheism** and Holy War). One
hid wal Jihad (**Monotheism** and Holy War). One civilian was among
line group Tawhid wal Jihad (**Monotheism** and Holy War). One civilian
that reads "Arbeit Macht Frei" ("**Work Brings Freedom**") is a seminal moment for
Tavola di San Giuseppe (St. **Joseph's Feast**). You'll
As part of the Tres Fronteras (**Three Borders**) area that includes Foz and

(c) Max activating examples for the translation to English detector feature 31123.

to start in the homeland of ties - **Croatia**. There we found three local brands that
Cage SOFIA (Reuters) - **Bulgaria's** president on Thursday called for a
<pad><bos>Welcome to The Dublin: **Ireland's** oldest and newest discovery trail.
<pad><pad><pad><pad><pad><pad><bos>**Bulgaria**Ski.com is owned and m
Everyone should know that "Deutschland" means **Germany** in German. **Germany** is
updates: urban Budapest sketch... (**Hungary**) The old building standing on V
<pad><bos>Message Behind African Heaters For **Norway** Spoof An online video, u
ually, I'm More of a **Switzerland**: More Personal Ads from the London Review

(e) Max activating examples for the country detector feature 11459.

other system that substantively uses<bos>Wikipedia **sobre fisica de particulas**
Directed By Tom Grundy Es **gibt noch keine Kommentare**. Sei der erste ...<bos>
006 - 213 **halaman** The book was selected as one of
Hour | Webcast - enregistré | **Où et quand** What is the Webcast About
[score hidden] 23 **Minuten zuvor** You just said 'if you exposed
hazard [score hidden] 23 **Minuten zuvor** You just said 'if you
vA-LINKER **biedt** mogelijkheden om een publicatie elders te vinden (

(b) Max activating examples for the English to foreign language switch detector feature 7928.

Say You can rate this item by **giving** it a score of one (poor),
, please let us know about it by **sending** our help desk an email .
which your order will be shipped. By **doing** this the few products
<pad><pad><pad><bos>Search for music by **typing** a word or ph
a one month non-recurring subscription by **sending** a cashier's c
s... Learn more about Concordia by **following** the links below: Cc
this product deliver? Pay it forward by **sharing** what you loved (a
. Browse: Browse the database by **applying** one or more filters to
we are celebrating Valentine's Day by **sharing** some gorgeous ai
page needs content. You can help by **adding** a sentence or a ph
NewsOK. He composed the ad by **animating** still photos taken b

(d) Max activating examples for the gerund form detector feature 8446.

><bos>I have been a technology **journalist** and consultant for near
fination will help independent Adventist **journalism** expand across
ian 50 **journalists** gathered at Klosters, a Swiss ski
ting punters, **journalists**, football managers and players. We also
Modelo<bos>The award-winning **journalist** Robert Fisk gave the in
Pulitzer Prize-winning **journalist**, formerly with The Washington Po
:w. If you are a **journalist** seeking comment on a story or more info
<bos>Peripatetic **journalist** and translator Porter (Road to Heaven:
n will likely endanger the lives of **journalists** and aid workers in the
houghtful post about the hazards of **journalism** following revelatio
about interviewing and **journalism**. Just like a marketing person do

(f) Max activating examples for the journalist feature 26436. (The strongest detector for the person-profession task).

Figure 21: Max activating examples for detector features from the Figure 5

<p>1674 anyway, dogs will be dogs. My name</p> <p>1675 hook. Boys will be boys, I suppose</p> <p>1676 really hate explicit shock for the sake of shock</p> <p>1677 ExtractionOptionsJsonObject ExtractionOptionsJsonObject()</p> <p>1678 Vertice3f::Vertice3f</p> <p>1679 Vertice3f::Vertice3f</p> <p>1680 four ways, and only four ways, in</p> <p>1681 ProfilerJniMethod SamplingProfilerJni</p> <p>1682 particular device and only for that device.</p> <p>1683 (a) Max activating examples for the repetition executor feature 12646. Extracted from the algo_last TV.</p>	<p>means lost and found in the Mandingo language</p> <p>no probleme i can speak english</p> <p>because he could not speak Mandarin or Cantonese.</p> <p>A. Heim. In German & French.</p> <p>"market field" in Malay. It was</p> <p>of her team speak in English, Elena immediately</p> <p>international projects. I speak english, spanish and</p> <p>xml:lang="en"></p> <p>(b) Max activating examples for the language prediction executor feature 13804.</p>
--	--

<p>1688</p> <p>1689 in the Swazi capital, Mbabane</p> <p>1690 in the densely populated capital of Monrovia,</p> <p>1691 km north of the capital Kabul. A crowd</p> <p>1692 east of the regional capital Poznan.</p> <p>1693 south of the regional capital Biay</p> <p>1694 west of the regional capital Wroclaw.</p> <p>1695 ero neighbourhood in the capital, Bujumbura</p> <p>1696 the smells that filled downtown Greenville in December</p> <p>1697 (c) Max activating examples for the capital prediction executor feature 16315.</p>	<p>st issue of Norsk Entomologisk Tidsskrift (now Norwegian Journal of Entomology) appeared in May 1921.</p> <p>vision and administration of the Direccion Provincial de Vialidad (Provincial Dept. of Transportation).</p> <p>Insekt-Nytt (Insect News) is written in a popular science style and is the society's</p> <p>reja Universal do Reino de Deus (Universal Church of the Kingdom of God) denied involvement in the scandal.</p> <p>Theatre in the play Den Sorte Dronning (The Black Queen) in 1843. Many artist frequented the</p> <p>icts Lithuanian Zalioji rinktinė (The Green Squad), belonging to partisans' Algimantas military district</p> <p>rvals. Norske Insekttabeller (Norwegian Insect Tables) is a series of inexpensive Norwegian-</p> <p>(d) Max activating examples for the translation executor feature 493.</p>
--	---

Figure 22: Max activating examples for Gemma 2 2B 16k executor features from the Figure 19a

<p>1702</p> <p>1703</p> <p>1704</p> <p>1705</p> <p>1706), Judd Ringer (right end), George Benson</p> <p>1707 plays as a winger or as a left back</p> <p>1708 Castilla as either a central defender or a left</p> <p>1709</p> <p>1710 , it's right tackle. Filling in</p> <p>1711</p> <p>1712</p> <p>1713 (a) Max activating examples for the football_player_position executor feature 18981.</p> <p>1714</p> <p>1715</p> <p>1716</p> <p>1717 964), English rugby player</p> <p>1718 014), American racing driver</p> <p>1719 936), British composer, conductor and</p> <p>1720</p> <p>1721 4) was an English sportsman who played rugby</p> <p>1722</p> <p>1723 1) was an Italian footballer from Bastia</p> <p>1724</p> <p>1725 (c) Max activating examples for the person_profession executor feature 46729.</p> <p>1726</p> <p>1727</p>	<p>apparent use of compliant and non-compliant form</p> <p>die soon, today, tomorrow, or in</p> <p>various degrees of reluctant and unlikely. There is</p> <p>no matter how different or diverse these may be</p> <p>: What are reliable and trusted websites? How</p> <p>are clearly upsides and downsides for those companies</p> <p>bed, standing up, falling back down,</p> <p>(b) Max activating examples for the antonyms executor feature 45288.</p> <p>Mai-Mai Kata Katanga ("SeBede Katanga"). Other Mai-Mai groups There was a large Mai</p> <p>or Low Saxon, i.e. that they remain for ever together undivided). Christian's ascension</p> <p>an Zalioji rinktinė (The Green Squad), belonging to partisans' Algimantas military distri</p> <p>jed square called Pasar Medan – literally, "market field" in Malay. It was here that the cit</p> <p>ey had a plastic kagami mochi, which translates to mirror rice cake. Basically a snowman m</p> <p>and Roberto Jefferson. The Ministério Público Federal (the Federal Prosecutor's Office)</p> <p>(d) Max activating examples for for translation to English executor feature 62633.</p>
--	---

Figure 23: Max activating examples for Gemma 2 2B 65k executor features from the Figure 19b