## **Training Superior Sparse Autoencoders for Instruct Models**

**Anonymous ACL submission** 

#### Abstract

As large language models (LLMs) grow in scale and capability, understanding their internal mechanisms becomes increasingly critical. Sparse autoencoders (SAEs) have emerged as a key tool in mechanistic interpretability, enabling the extraction of human-interpretable features from LLMs. However, existing SAE training methods are primarily designed for base models, resulting in reduced reconstruction quality and interpretability when applied to instruct models. To bridge this gap, we propose Finetuning-aligned Sequential Training (FAST), a novel training method specifically tailored for instruct models. FAST aligns the training process with the data distribution and activation patterns characteristic of instruct models, resulting in substantial improvements in both reconstruction and feature interpretability. On Qwen2.5-7B-Instruct, FAST achieves a mean squared error of 0.6468 in token reconstruction, significantly outperforming baseline methods with errors of 5.1985 and 1.5096. In feature interpretability, FAST yields a higher proportion of high-quality features, for Llama3.2-3B-Instruct, 21.1% scored in the top range, compared to 7.0% and 10.2% for BT(P) and BT(F). Surprisingly, we discover that intervening on the activations of special tokens via the SAEs leads to improvements in output quality, suggesting new opportunities for fine-grained control of model behavior. Code, data, and 240 trained SAEs will be publicly released, enabling broad adoption and future research.

#### 1 Introduction

011

014

019

Large Language Models (LLMs) demonstrate exceptional performance across diverse natural language processing tasks (Brown et al., 2020; Ouyang et al., 2022; Guo et al., 2025). However, their complexity, vast number of parameters, and intricate training processes present significant challenges in understanding their internal mechanisms (Bengio et al., 2023; Bubeck et al., 2023).



Figure 1: Overview of the sparse autoencoder, illustrating its process for interpreting the internal representations of large language models.

As these models advance, aligning them with human values and mitigating risks becomes critical, highlighting the importance of mechanistic interpretability (Bereska and Gavves, 2024; Ji et al., 2023; Anwar et al., 2024). Sparse autoencoders (SAEs) serve as a powerful tool for interpreting LLMs by mapping high-dimensional activations to sparse, interpretable feature spaces, thereby decomposing neural networks into understandable components (Bereska and Gavves, 2024; Bricken et al., 2023; Cunningham et al., 2023). SAE training, conceptualized as dictionary learning (Kreutz-Delgado et al., 2003; Yun et al., 2021), utilizes hidden layer weights as dictionary bases and enforces sparsity for efficient representations, aligning with the linear representations and superposition hypotheses (Elhage et al., 2022; Arora et al., 2018; Olah, 2022). Figure 1 provides an overview of sparse autoencoders.

Current SAE training methods primarily focus on base models and follow Block Training paradigm that concatenates datasets and splits them into fixed-length blocks (Joseph Bloom and Chanin, 2024; Bricken et al., 2023). It aligns with the pre-

067

044

training phase of LLMs, making it a natural and effective choice for training SAEs on base models. While effective for base models, this method faces significant limitations when applied to instruct models (Joseph Bloom and Chanin, 2024; Kissane et al., 2024b). The semantic discontinuity caused by combining data from diverse sources undermines the semantic coherence for alignment with downstream tasks, ultimately degrading SAE training performance (Kissane et al., 2024b).

069

070

077

085

101

102

103

104

106

107

108

109

110

111

112

113

114

115

To address these challenges, we propose **F**inetuning-**a**ligned Sequential Training (FAST), a novel SAE training method specifically designed for instruct models. FAST processes each data instance independently, preserving semantic integrity and maintaining alignment with the fine-tuning objectives of the model. By providing a consistent and complete semantic space during SAE training, FAST enhances the model's understanding of input and improves the quality of feature extraction.

Experimental results demonstrate that *FAST* significantly enhance SAE performance across various tasks. In token reconstruction on Qwen2.5-7B-Instruct (Yang et al., 2024), *FAST* achieves a mean squared error of 0.6468, outperforming baselines of 5.1985 and 1.5093. It also excels in feature interpretability; for Llama3.2-3B-Instruct (Dubey et al., 2024), 21.1% of features are rated highest in quality, compared to 7.0% for BT(P) and 10.2% for BT(F). Additionally, SAEs are used to study the impact of special tokens on outputs, offering insights into their roles and practical applications, and paving the way for future research.

Our contributions are summarized as follows:

- This paper proposes Finetuning-aligned Sequential Training (*FAST*), a novel method specifically designed for training SAEs on instruct models.
- Experimental results demonstrate that *FAST* significantly improves the performance of SAEs on token reconstruction. Additionally, feature interpretability experiments confirm the effectiveness and generalizability of *FAST*.

• The SAEs are further utilized to investigate the influence of special tokens on model outputs, providing new insights into their specific roles and offering fresh directions for the practical application of SAE models.

#### 2 Related Work

Mechanistic Interpretability. As LLMs continue to advance, their increasing complexity, massive parameter scales, and intricate training processes present significant challenges to human understanding of their inner workings (Bubeck et al., 2023; Bengio et al., 2023). Achieving a deep understanding of LLMs is crucial to ensuring alignment with human values (Ji et al., 2023; Anwar et al., 2024) and mitigating harmful or unintended outcomes (Anwar et al., 2024; Hendrycks et al., 2021; Slattery et al., 2024; Hendrycks et al., 2023). However, the "black box" nature (Casper et al., 2024) obscures the underlying causes of misalignment and associated risks. To address these challenges, mechanistic interpretability has emerged as a critical area of research focused on understanding the inner workings of LLMs (Bereska and Gavves, 2024; Nanda, 2022d, 2023, 2022a; Olah, 2022). This discipline seeks to achieve a detailed understanding of model behavior through systematic reverse engineering (Nanda, 2022c,b).

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

Sparse Autoencoders for LLM. The training of sparse autoencoders (SAEs) can be framed as a form of dictionary learning, where the hidden layer weights serve as the dictionary basis, and sparsity constraints enforce efficient and sparse data representations (Bereska and Gavves, 2024; Bricken et al., 2023). Additionally, SAEs align with both the linear representations hypothesis (Mikolov et al., 2013) and the superposition hypothesis (Elhage et al., 2022; Arora et al., 2018; Olah et al., 2020), ensuring that the learned representations adhere to theoretical principles of highdimensional feature spaces. Specifically, the linear representation hypothesis suggests that features in language models correspond to directions in activation space, enabling embedding arithmetic, such as: v("king") - v("man") + v("woman") =v("queen") (Mikolov et al., 2013).

Neurons in LLMs are often polysemantic, encoding multiple distinct features due to the limited dimensionality of feature activation space. (Bereska and Gavves, 2024). The superposition hypothesis explains how neural networks represent more features than the number of available neurons by encoding features as nearly orthogonal directions in the neuron output space (Elhage et al., 2022). The activation of one feature may appear as a slight activation of another, resulting from the overlap of non-orthogonal vectors. While such overlaps intro-



Figure 2: Illustration of the LLM training pipeline and SAE training methods. (a) The pipeline transitions from pretraining to fine-tuning. (b) Block Training (BT) concatenates datasets and resplits them into fixed-length blocks. (c) Finetuning-aligned Sequential Training (FAST) processes data instances independently, preserving semantic integrity and improving alignment with fine-tuning objectives, leading to better performance in feature interpretability.

duce interference, the advantages of representing a greater number of non-orthogonal features outweigh the drawbacks, particularly in highly sparse neural networks (Bricken et al., 2023; Bereska and Gavves, 2024; Rajamanoharan et al., 2024a). This property makes SAEs particularly valuable in mechanistic interpretability, as they enable the decomposition of language models by capturing high-dimensional features (Gao et al., 2024; Ferrando et al., 2024; Rajamanoharan et al., 2024b; Lieberum et al., 2024; He et al., 2024).

167

168

169 170

171

172

173

174

175

176

177

#### 3 <u>Finetuning-aligned Sequential Training</u>

Motivation. Recent studies have adopted a train-179 ing paradigm for SAE that builds upon the pretraining phase of LLMs, as illustrated in Figure 2(b). 181 This approach, referred to as Block Training (BT), 182 involves concatenating datasets and splitting them into fixed-length blocks for training (Bereska and 184 Gavves, 2024; He et al., 2024; Kissane et al., 2024a). BT aligns with the pretraining phase of LLMs, making it a natural and effective choice for 188 training SAEs on base models. Since base models are directly trained on large-scale corpora without additional fine-tuning, BT ensures consistency 190 between the SAE training and the pretraining objectives of LLMs. 192

However, when it comes to instruct models, which undergo a supervised fine-tuning (SFT) phase to align with specific instructions or downstream tasks, the limitations of BT become more apparent. For instance, studies demonstrate that SAE trained on the pretraining dataset exhibit significantly weak abilities in adhering to refusal directives (Kissane et al., 2024b). An alternative approach utilizes SFT datasets, introducing special tokens and applying block training in the same manner (Kissane et al., 2024b). While this method leverages SFT datasets, it still preserves the BT methodology, which does not align well with the finetuning objectives of instruct models. Specifically, BT treats the input sequences as concatenated blocks, often combining data samples from different sources. For example, in a sequence of 8,192 tokens, the first 2,048 tokens may originate from one sample, while the remaining 6,144 tokens come from another. While such semantic discontinuity is less problematic for base models, as it mirrors their pretraining setup, it poses significant challenges for instruct models. Maintaining semantic integrity is crucial for aligning with downstream tasks, and the lack of such alignment hinders the model's ability to fully understand the input, ultimately degrading SAE training performance.

193

194

195

196

197

198

200

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

To address these challenges, we propose a novel SAE training method for instruct models: Finetuning-aligned Sequential Training (*FAST*), which better aligns with the fine-tuning phase, both in terms of dataset utilization and training methodology in Figure 2(c). By providing the instruct model with a consistent and complete semantic space during SAE training, *FAST* enhances the alignment with the fine-tuning phase and improves the quality of SAE training. This alignment forms the primary motivation behind *FAST*.

#### 3.1 Data Processing

221

231

236

237

240

241

242

243

244

245

246

247

249

253

256

258

259

261

263

As previously described, *FAST* trains the SAE using finetuning datasets. Specifically, multiple multiturn dialogue datasets are collected, and each data instance is combined with the corresponding chat template of the instruct model. This process not only introduces special tokens but also ensures consistency with the data processing methodology used during the fine-tuning phase of the model.

A key innovation lies in independent processing of each data instance, rather than concatenating multiple instances before inputting them into the model. By eliminating the constraint of context size, the dataset is processed sequentially. Each data instance is individually fed into the LLM to extract hidden layer activations, which subsequently used to train the SAE, as illustrated in Figure 2(c). This approach effectively avoids semantic discontinuity caused by data concatenation, while preserving the semantic integrity of each instance thereby providing higher-quality inputs for training the SAE.

#### 3.2 SAE

This section introduces the two types of SAE models utilized in *FAST*: the Standard ReLU-based SAE and the JumpReLU SAE. The Standard ReLUbased SAE is a widely adopted approach (Bereska and Gavves, 2024; Bricken et al., 2023), while JumpReLU SAE achieves superior reconstruction quality and sparsity control (Rajamanoharan et al., 2024a; Lieberum et al., 2024). Here we provide the details of the two SAE models and the initialization method in Appendix A.

**Standard SAE.** For the input vector  $x \in \mathbb{R}^{d_{in}}$ from the residual stream,  $d_{in}$  denotes the dimensionality of the model's hidden layer. The ReLUbased SAE model consists of an encoder, decoder, and a corresponding loss function, which are defined as follows:

$$f(\mathbf{x}) = \operatorname{ReLU}(\mathbf{W}^{\operatorname{enc}}\mathbf{x} + \mathbf{b}^{\operatorname{enc}}) \qquad (1) \qquad 2$$

$$\hat{\mathbf{x}} = \mathbf{W}^{\text{dec}} f(\mathbf{x}) + \mathbf{b}^{\text{dec}}$$
(2) 27

$$\mathcal{L} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \lambda \|\mathbf{z_{L1}}\|$$
(3) 2

 $\mathbf{W}^{\text{enc}}, \mathbf{W}^{\text{dec}}, \mathbf{b}^{\text{enc}}, \mathbf{b}^{\text{dec}}$  represent the weight matrices and bias vectors for the encoder and decoder, respectively.  $\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$  denotes the mean squared error (MSE) loss,  $\|\mathbf{z}_{L1}\|_1$  represents the  $L_1$  loss used for sparsity regularization, and  $\lambda$  is the sparsity regularization hyperparameter.

**JumpReLU SAE.** The JumpReLU SAE retains the same parameter matrices **W** and **b** as the Standard SAE but introduces a modified activation function and sparsity regularization:

$$f(\mathbf{x}) = \text{JumpReLU}_{\theta}(\mathbf{W}^{\text{enc}}\mathbf{x} + \mathbf{b}^{\text{enc}}), \quad (4)$$

$$\hat{\mathbf{x}} = \mathbf{W}^{\text{dec}} f(\mathbf{x}) + \mathbf{b}^{\text{dec}},\tag{5}$$

$$\mathcal{L} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \lambda \|\mathbf{z_{L0}}\|, \qquad (6)$$

The JumpReLU function is defined as JumpReLU<sub> $\theta$ </sub>(z) :=  $z \odot H(z - \theta)$ , where  $\theta > 0$  is a learnable, vector-valued threshold parameter. Here,  $\odot$  denotes elementwise multiplication, and H represents the Heaviside step function. Additionally,  $\|\mathbf{z}_{L0}\|_1$  represents the  $L_0$  loss used for sparsity regularization, while  $\lambda$  is the sparsity regularization hyperparameter.

#### 3.3 Mixing Activation Buffer

Activation values, which represent the activation levels of hidden layer dimensions during the model's forward pass, require significant storage space. To mitigate this challenge, we employ a producer-consumer framework inspired by previous studies (Joseph Bloom and Chanin, 2024), wherein the LLM generates activations and stores them in a dedicated buffer.

As shown in Figure 3, the process begins with the buffer being filled to capacity with activation values. Once the buffer is full, the activations are shuffled to ensure randomness and diversity. Subsequently, half of the shuffled activations are sent to the SAE model for training, while the other half remains in the buffer. After training, the buffer is replenished with new activations generated by the 273

274

275

276

277

278

279

280

281

286

287

288

289

290

291

292

293

294

296

297

298

299

300

301

302

303

304

305

306

307

308

309



Figure 3: The mixing activation buffer is shuffled, half is sent to the SAE for training, and the resulting new activations are used to refill the buffer. This iterative process ensures data diversity and storage efficiency.

model, and the cycle repeats. This iterative mechanism optimizes storage efficiency and ensures a
high level of data variability, thereby enhancing the
robustness of model training. By leveraging the
mixing buffer, this approach effectively balances
data diversity with storage efficiency.

#### 4 Experiments

#### 4.1 Experiment Setup

**Dataset.** We construct a large-scale instruction dataset for fine-tuning LLMs by combining several publicly available, high-quality datasets, including WildChat-1M-Full (Zhao et al., 2024), Infinity-Instruct (BAAI, 2024), tulu-3-sft-mixture (Lambert et al., 2024), orca-agentinstruct-1M-v1-cleaned <sup>1</sup>, and Imsys-chat-1m (Zheng et al., 2023). After applying a 20-gram deduplication strategy, it is reduced to 4,758,226 samples. Details are in Appendix B.

**LLMs.** We conduct experiments on seven models from two families: Llama (Llama-3.1, Llama-3.2)(Dubey et al., 2024) and Qwen (Qwen-2.5)(Yang et al., 2024), selected for their state-of-the-art performance to evaluate our approach's robustness and generalization across families and scales. The models and their respective layer configurations, detailed in Table 1, are selected from various depths to mitigate depth bias. Following prior works (Bereska and Gavves, 2024; Bricken et al., 2023; Gao et al., 2024), we train SAEs on the residual stream, as inter-layer relationships have minimal impact on performance.

**Baselines.** Prior to this study, all SAE model training methods exclusively utilize the Block Training (BT) strategy. Depending on the type of training dataset used, Block Training can be categorized into two primary forms: BT(P) and BT(F) as follows:

- *BT*(*P*): Block Training using the pretraining dataset. The pretraining dataset is processed by concatenating and segmenting the data into text blocks of equal length, which are then used for training the SAE model.
- *BT(F)*: Block Training using the finetuning dataset. This approach utilizes a finetuning dataset. The data within the dataset is concatenated to form text blocks.

For BT(P), we utilize the pile-uncopyrighted dataset <sup>2</sup>. As for BT(F), we use the finetuning dataset metioned before which is also used in *FAST*.

**Configuration.** SAEs are trained on 8\*NVIDIA A100 GPUs using sae\_lens (Joseph Bloom and Chanin, 2024) with custom implementation. For models more than 7B parameters, the expansion factor of SAE is fixed at 8X, whereas for other models, the expansion factor can be 8X or 16X. To ensure fairness across methods at the same data scale, the number of training tokens is set to 40,960,000. For BT(P) and BT(F), context\_size is 2,048, with each text block containing 2,048 tokens. For *FAST*, no explicit context\_size is required; instead, a truncation length of 8,192 is applied to manage memory usage. For JumpReLU SAE,  $L_{\text{sparsity}}$  is 0.01, while for Standard SAE, it is 5. Further parameter details are in Appendix C.

**Evaluation Metric.** The performance of the SAE is assessed using the Mean Squared Error (MSE), which is calculated as:

$$MSE = \frac{\sum_{i=1}^{N} \frac{1}{L_i} \sum_{j=1}^{L_i} \sum_{k=1}^{H} (y_{i,j,k} - \hat{y}_{i,j,k})^2}{N \cdot H}$$
(7)

where N denotes the size of the dataset,  $L_i$  represents the length of the *i*-th sequence, H refers to the hidden dimension of the model. To evaluate the SAE's performance specifically on special tokens, we also compute the MSE of special tokens, denoted as  $MSE_{st}^3$ . Lower MSE values reflect better model performance.

<sup>2</sup>https://huggingface.co/datasets/monology/ pile-uncopyrighted <sup>&</sup>lt;sup>1</sup>https://huggingface.co/datasets/mlabonne/ orca-agentinstruct-1M-v1-cleaned

<sup>&</sup>lt;sup>3</sup>To facilitate a more direct comparison of performance



Figure 4:  $MSE_{st}$  performance of the JumpReLU SAE (all metrics are presented in log scale, where lower values indicate better SAE reconstruction performance). Within the JumpReLU architecture, *FAST* exhibits the best reconstruction capability compared to BT(P) and BT(F).

Model Name	Layer
Llama series	
Llama-3.1-8B-Instruct	[4,12,18,20,25]
Llama-3.2-3B-Instruct	[4,12,20]
Llama-3.2-1B-Instruct	[4,9,14]
Qwen series	
Qwen2.5-7B-Instruct	[4,12,18,20,25]
Qwen2.5-3B-Instruct	[4,18,32]
Qwen2.5-1.5B-Instruct	[4,14,24]
Qwen2.5-0.5B-Instruct	[4,12,20]

Table 1: Layer configurations of the Llama and Qwen model series, showcasing the selection of layers across varying depths to mitigate depth-related biases and optimize model performance.

#### 4.2 Main Results

390

391

394

398

A random sample of 5,000 dialogues is extracted from the remaining portion of the dataset for evaluation. Figure 4 compares the  $MSE_{st}$  scores of three methods using the JumpReLU SAE, while Figure 6 illustrates the  $MSE_{st}$  performance of the Standard SAE. Detailed results for both MSE and  $MSE_{st}$  are presented in Appendix D.

In terms of overall token reconstruction (MSE), the JumpReLU architecture with Qwen models demonstrates similar patterns, with *FAST* consistently outperforming baseline methods. *FAST* method achieves superior performance across most configurations. For instance, in Llama-3.2-3B-Instruct-L20-8X-Standard, *FAST* attains -0.9527, significantly surpassing the baselines which score -0.6926 and -0.9186. In special token reconstruction ( $MSE_{st}$ ), *FAST* shows marked improvements across models. In Qwen2.5-7B-Instruct-L18-8X-Standard, *FAST* achieves 0.6468, outperforming the baselines (5.1985 and 1.5093). In the JumpReLU SAEs, it achieves -9.7604 compared to -4.0005 and -8.0743. 399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

Overall, the findings demonstrate that *FAST* excels in reconstructing both general and special tokens. Interestingly, *FAST* shows even stronger improvements in Standard SAE architectures compared to JumpReLU SAEs, potentially due to the latter's already high MSE performance, leaving less room for enhancement. Despite limitations in Standard architectures due to L1 regularization and ReLU activation, *FAST* significantly improves token reconstruction in these models.

#### **5** Feature Interpretability

This section evaluates the interpretability of features extracted by SAEs through an automated analysis framework, building upon methodologies (Bills et al., 2023; Cunningham and Conerly, 2024; He et al., 2024). The middle layers of the trained SAEs are selected for analysis based on their demonstrated superior performance. Given that experiments demonstrate that the JumpReLU activation function outperforms

across different methods, all MSE values are transformed using  $\log_2.$ 



Figure 5: Experiment results of feature interpretability.*FAST* achieves notable improvements compared to the other two training methods across all the tested models. *FAST* attains 21.1% of features rated in the highest quality range (scores 4-5), in contrast to 7.0% for BT(P) and 10.2% for BT(F).

Score	Description
5	Clear pattern with no deviating examples
4	Clear pattern with one or two deviating examples
3	Clear overall pattern but quite a few ex- amples not fitting that pattern
2	Broad consistent theme but lacking struc- ture
1	No discernible pattern

Table 2: Scoring criteria for feature interpretability.

other alternatives (Rajamanoharan et al., 2024b; Lieberum et al., 2024), the evaluation exclusively employs SAEs equipped with JumpReLU. Table 10 presents the specific SAE models evaluated.

Additional 10,000 instances are sampled and their activation values are computed. Then the top five sentences with the highest activation values are identified to construct an activation dataset for evaluating features. Based on the assumption that dead features are irrelevant to the evaluation, an initial screening of features is conducted, ensuring that only features with non-zero activation values in top five sentences are retained. After that, we randomly select 128 features as the final evaluation.

GPT-40<sup>4</sup> is prompted to score each group of five contexts and generate a descriptive summary. Additionally, a monosemanticity score ranging from 1 to 5 is assigned, based on a rubric adapted from (Cunningham and Conerly, 2024; He et al., 2024). Detailed prompt is shown in Appendix E.2.

A total of 4,608 feature scores are computed and presented in Figure 5. The results demonstrate that *FAST* consistently outperforms BT(P) and BT(F) across all evaluated SAEs. For the 8x scaled Llama3.2-3B-Instruct, *FAST* achieves 21.1% of features in the highest quality range (scores 4-5), compared to 7.0% for BT(P) and 10.2% for BT(F). Generally, compared to both baseline methods, we observe that *FAST* reduces the proportion of low-quality features while increasing the proportion of high-quality features in 8X and 16X SAEs. This highlights the superiority of *FAST* in producing more interpretable features during SAE training.

Furthermore, Cumulative Distribution Function (CDF) curve analysis reveals that *FAST*'s percentage of features scoring below 3 is consistently the lowest. For instance, with Qwen2.5-3B-Instruct model, the CDF at score 3 is 76.5% for *FAST*, compared to 89.0% for BT(F) and 92.2% for BT(P), indicating fewer low-scoring features for *FAST*. These findings suggest that both appropriate training dataset selection for SAEs and the sequence training methodology contribute to enhanced model interpretability. *FAST* appears to successfully integrate these aspects, leading to more interpretable SAEs.

#### 6 Steering with SAE Latents

Feature steering represents an intuitive approach to evaluate model inference by adjusting the activation coefficients within a trained SAE, thereby directly influencing the model's output. This method resembles the use of decoder latent vectors for activation guidance, but the SAE offers a more robust and unambiguous process for activation guidance. Based on the formulations in Equations 2 and 5,

<sup>&</sup>lt;sup>4</sup>GPT-40 version: 2024-11-20

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

531

the reconstructed outputs of the SAE derive from a weighted combination of its latent variables. (Ferrando et al., 2024; Templeton, 2024).

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

500

501

502

503

504

505

506

507

508

510

511

512

513

514

515

516

517

518

521

523

525

527

528

530

$$z' = z + \alpha d_k \tag{8}$$

These latent variables correspond to row vectors of  $W_{dec}$ , with  $\alpha$  scaling the k-th latent. To implement this steering, a latent dimension k is selected, scaling its decoder vector  $d_k$  by  $\alpha$ . Then z' is introduced into the model's residual stream.

Following Ferrando et al. (2024), 1,010 sampled instruction instances are randomly partitioned into two parts: 1,000 samples to identify highly activated SAE features and 10 samples to evaluate post-steering model outputs. We use the chat template corresponding to the instruct model during inference. The 10 questions appear in Appendix F.1. We focus on feature related to these special tokens<sup>5</sup>(shown in Table 11) to examine how special tokens, which are not associated with specific entities, influence the model's output. Using 1,000 samples, the average maximum activation values are calculated for each feature. Complete activation values for each model appear in Appendix F.3.

Three representative questions are selected to illustrate the effects of steering features. Due to space constraints, feature steering primarily focuses on the <|start\_header\_id|> for Llama3.1-8B-Instruct and <|im\_start|> for Qwen2.5-7B-Instruct. The experiments employ scaling  $\alpha \in$ [0, 15, 25, 50, 100, 150, 200] using 8X JumpReLU SAE through *FAST* and greedy decoding. Detailed analyses of three questions are presented in Appendix F.4.

Steering high-activation features particularly those associated with special tokens significantly influences the model's output quality and reasoning ability. This effect remains consistent across diverse tasks and linguistic contexts. There is an optimal range for the coefficient  $\alpha$ . Within this range, model responses become more accurate, coherent, and relevant to the given instructions.

For instance, in Question 3(F.4.2), amplifying the activation of a feature tied to both the  $<|im\_start|>$  and user results in a clear transition: moderate values of  $\alpha$  improved engagement and output relevance, while excessive amplification led to language switching and incoherent, repetitive text. Similarly, in Question 4(F.4.3), steering the highest activation feature associated with the < $|im\_start|$ > marker within a specific coefficient range led to more convincing and logically structured answers, but pushing  $\alpha$  too far again degraded output quality. Similar patterns can also be observed in Q2(F.4.1).

The consistency in findings suggests that these features encode essential aspects of the model's reasoning capabilities, transcending individual tasks or linguistic contexts. There is an optimal coefficient  $\alpha$  range suggests a "sweet spot" for feature steering, enhancing performance without introducing the degradation seen at higher coefficients.

This observation presents important implications for the practical application of SAEs. It demonstrates that steering certain features potentially associated with special tokens emerges as a reliable method to improve model performance across diverse tasks. Unlike traditional SAE-feature approaches, which often impose output biases tied to predefined meanings or entities, feature steering with special tokens refines the guidance of models, resulting in higher-quality responses.

#### 7 Conclusion

This paper proposes a novel approach, Finetuningaligned Sequential Training (FAST), for training SAEs on instruct models. By independently processing individual data instances while maintaining semantic integrity, *FAST* addresses the limitations of traditional Block Training (*BT*) methods, which often suffer from semantic discontinuity and misalignment with downstream task requirements. Experimental results show that *FAST* improves performance across various SAE models, demonstrating its versatility and general applicability. Furthermore, *FAST* consistently achieves superior results in feature interpretability evaluations, highlighting its effectiveness and advantages.

Also we employ SAEs to explore the influence of special tokens on model outputs. Results indicate that steering features within a specific coefficient range substantially enhance model output quality. These insights provide a valuable method for studying the functional roles of special tokens and practical applications of SAEs. To facilitate future research, the complete codebase, datasets and a total of 240 pre-trained SAE models will be released publicly, establishing a robust foundation for innovation and advancement in this domain.

<sup>&</sup>lt;sup>5</sup>user and assistant are incorporated into the special tokens, as they frequently appear together with other special tokens.

597

611

612

613

614

615

616

617

618

619

621

622

623

624

625

# Limitations

582 As language models increase in scale, investigating their internals with SAE-based methods becomes 583 more challenging. Computational constraints re-584 strict our investigation to smaller Qwen and Llama 585 models (under 8B parameters), though our frame-586 work could be extended to larger architectures. Feature interpretability analysis focuses mainly on strongly activated features, potentially overlooking weakly activated samples (He et al., 2024). Furthermore, feature steering experiments are prelimi-591 592 nary studies centered on special token-related features that correlate with response quality. A more comprehensive investigation of these features' influence remains an important direction for future research. 596

# **Ethical Statements**

This research focuses on interpreting and steering instruction-tuned language models through sparse autoencoders. All experiments rely solely on publicly available, appropriately licensed text corpora that are deduplicated and stripped of personally identifiable information; no human subjects are involved nor private data collected. Nevertheless, it is important to acknowledge that LLMs are trained on extensive publicly available datasets, potentially resulting in inadvertent reproduction of copyrighted material. Our codes, parameters, and deduplicated demo data will be released under an open-source licence to support reproducibility.

### References

- Usman Anwar, Abulhair Saparov, Javier Rando, Daniel Paleka, Miles Turpin, Peter Hase, Ekdeep Singh Lubana, Erik Jenner, Stephen Casper, Oliver Sourbut, et al. 2024. Foundational challenges in assuring alignment and safety of large language models. *arXiv preprint arXiv:2404.09932*.
- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. 2018. Linear algebraic structure of word senses, with applications to polysemy. *Transactions of the Association for Computational Linguistics*, 6:483–495.
- BAAI. 2024. Infinity instruct. arXiv preprint arXiv:2406.XXXX.
- Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, Gillian Hadfield, et al. 2023. Managing ai risks in an era of rapid progress. *arXiv preprint arXiv:2310.17688*, page 18.

Leonard Bereska and Efstratios Gavves. 2024. Mechanistic interpretability for ai safety–a review. *arXiv preprint arXiv:2404.14082*.

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. 2023. Language models can explain neurons in language models.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. 2023. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*. Https://transformercircuits.pub/2023/monosemanticfeatures/index.html.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Stephen Casper, Carson Ezell, Charlotte Siegmann, Noam Kolt, Taylor Lynn Curtis, Benjamin Bucknall, Andreas Haupt, Kevin Wei, Jérémy Scheurer, Marius Hobbhahn, et al. 2024. Black-box access is insufficient for rigorous ai audits. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 2254–2272.
- Hoagy Cunningham and Tom Conerly. 2024. Circuits updates june 2024. *Transformer Circuits Thread*.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. 2023. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. 2022. Toy models of superposition. *arXiv preprint arXiv:2209.10652.*

- 687 688 689
- 693 694 695 696 697 698 699 700 701 702 703 704 705 706
- 709 710 711 712 713 714 715 716 717 717

708

- 717 718 719 720 721 722 723 724 725 726 727
- 725 726 727 728 729 730
- 731 732 733 734

735 736

737

738 739

- Javier Ferrando, Oscar Obeso, Senthooran Rajamanoharan, and Neel Nanda. 2024. Do i know this entity? knowledge awareness and hallucinations in language models. *arXiv preprint arXiv:2411.14257*.
- Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. 2024. Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. arXiv preprint arXiv:2501.12948.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034.
- Zhengfu He, Wentao Shu, Xuyang Ge, Lingjie Chen, Junxuan Wang, Yunhua Zhou, Frances Liu, Qipeng Guo, Xuanjing Huang, Zuxuan Wu, et al. 2024. Llama scope: Extracting millions of features from llama-3.1-8b with sparse autoencoders. arXiv preprint arXiv:2410.20526.
- Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. 2021. Unsolved problems in ml safety. *arXiv preprint arXiv:2109.13916*.
- Dan Hendrycks, Mantas Mazeika, and Thomas Woodside. 2023. An overview of catastrophic ai risks. *arXiv preprint arXiv:2306.12001*.
- Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. 2023. Ai alignment: A comprehensive survey. *arXiv preprint arXiv:2310.19852*.
- Curt Tigges Joseph Bloom and David Chanin. 2024. Saelens. https://github.com/jbloomAus/ SAELens.
- Connor Kissane, Robert Krzyzanowski, Arthur Conmy, and Neel Nanda. 2024a. Saes (usually) transfer between base and chat models. Alignment Forum.
- Connor Kissane, Robert Krzyzanowski, Neel Nanda, and Arthur Conmy. 2024b. Saes are highly dataset dependent: A case study on the refusal direction. Alignment Forum.
- Kenneth Kreutz-Delgado, Joseph F Murray, Bhaskar D Rao, Kjersti Engan, Te-Won Lee, and Terrence J Sejnowski. 2003. Dictionary learning algorithms for sparse representation. *Neural computation*, 15(2):349–396.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri,

Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. 2024. Tülu 3: Pushing frontiers in open language model post-training. 740

741

742

743

744

746

747

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

771

772

776

778

780

781

782

783

784

785

786

787

789

790

791

792

- Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. 2024. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. *arXiv preprint arXiv:2408.05147*.
- Tomáš Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 conference of the north american chapter of the association for computational linguistics: Human language technologies*, pages 746–751.
- Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweti Mahajan, Dany Rouhana, Andres Codas, Yadong Lu, Wei ge Chen, Olga Vrousgos, Corby Rosset, Fillipe Silva, Hamed Khanpour, Yash Lara, and Ahmed Awadallah. 2024. AgentInstruct: Toward Generative Teaching with Agentic Flows. https://arxiv.org/abs/2407.03502. Preprint, arXiv:2407.03502.
- Neel Nanda. 2022a. 200 concrete open problems in mechanistic interpretability: Introduction. *Neel Nanda's Blog*.
- Neel Nanda. 2022b. 200 cop in mi: Analysing training dynamics. *Neel Nanda's Blog*.
- Neel Nanda. 2022c. 200 cop in mi: Interpreting algorithmic problems. *Neel Nanda's Blog.*
- Neel Nanda. 2022d. A comprehensive mechanistic interpretability explainer & glossary. *Neel Nanda's Blog.*
- Neel Nanda. 2023. Mechanistic interpretability quickstart guide. *Neel Nanda's Blog*.
- Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. 2020. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024– 001.
- Christopher Olah. 2022. Mechanistic interpretability, variables, and the importance of interpretable bases. *Transformer Circuits Thread*.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Sandip Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János Kramár, Rohin Shah, and Neel Nanda. 2024a. Improving

794dictionary learning with gated sparse autoencoders.795arXiv preprint arXiv:2404.16014.

796

797

799

800

801 802

806

810

811 812

813

814

815

816

817

818

819

827

830

831

832

- Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár, and Neel Nanda. 2024b. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders. arXiv preprint arXiv:2407.14435.
- Peter Slattery, Alexander K Saeri, Emily AC Grundy, Jess Graham, Michael Noetel, Risto Uuk, James Dao, Soroush Pour, Stephen Casper, and Neil Thompson. 2024. The ai risk repository: A comprehensive metareview, database, and taxonomy of risks from artificial intelligence. *arXiv preprint arXiv:2408.12622*.
  - Adly Templeton. 2024. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. Anthropic.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
  - An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.
  - Zeyu Yun, Yubei Chen, Bruno A Olshausen, and Yann LeCun. 2021. Transformer visualization via dictionary learning: contextualized embedding as a linear superposition of transformer factors. *arXiv preprint arXiv:2103.15949*.
  - Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. Wildchat: 1m chatGPT interaction logs in the wild. In *The Twelfth International Conference on Learning Representations*.
  - Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric. P Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. 2023. Lmsys-chat-1m: A large-scale real-world llm conversation dataset. *Preprint*, arXiv:2309.11998.

853

858

867

#### A SAE Initialization Method

The encoder weights  $(W_{enc})$  and decoder weights  $(W_{dec})$  are initialized using the Kaiming Uniform initialization method (He et al., 2015). This step, used exclusively in the JumpReLU method, normalizes each row of the  $W_{dec}$  using the L2 norm and adjusts the threshold  $\epsilon$  and encoder bias  $b_{enc}$  accordingly. After that, some data is selected for geometric median evaluation. The goal is to minimize the weighted sum of distances to all sample points. To achieve this, the Weiszfeld algorithm is employed to a specified precision of ftol =  $1 \times 10^{-20}$ . The resulting optimal point is then used as the initial value for  $b_{dec}$ , which is set to 0. There exists the formulas about the geometric median evaluation as follows:

$$f(\mathbf{m}) = \sum_{i=1}^{n} w_i \|\mathbf{m} - \mathbf{p}_i\|, \mathbf{m}_0 = \frac{\sum_{i=1}^{n} w_i \mathbf{p}_i}{\sum_{i=1}^{n} w_i}$$
(9)

$$d_i = \|\mathbf{p}_i - \mathbf{m}_k\|, w'_i = \frac{w_i}{\max(d_i, \epsilon)} \tag{10}$$

845 
$$\mathbf{m}_{k+1} = \frac{\sum_{i=1}^{n} w_i' \mathbf{p}_i}{\sum_{i=1}^{n} w_i'}$$
(11)

$$|f(\mathbf{m}_{k+1}) - f(\mathbf{m}_k)| \le \operatorname{ftol} \cdot f(\mathbf{m}_k) \tag{12}$$

The parameters used in the equations are defined as follows: **m** represents the target point or the weighted mean to be optimized, while  $\mathbf{p}_i$  is the *i*-th data point in the dataset.  $w_i$  denotes the weight associated with the *i*-th data point. The objective function,  $f(\mathbf{m})$ , is the weighted sum of distances between **m** and all data points  $\mathbf{p}_i$ . The initial estimate of **m**, denoted as  $\mathbf{m}_0$ , is calculated as the weighted mean of all points.  $d_i$  is the distance between the *i*-th data point  $\mathbf{p}_i$  and the current estimate  $\mathbf{m}_k$ . The updated weight for the *i*-th data point,  $w'_i$ , is adjusted by the distance  $d_i$  and a small constant  $\epsilon$  to prevent division by zero.  $\mathbf{m}_{k+1}$  is the updated estimate of **m** at iteration k + 1, computed as the weighted mean of all points using the updated weights  $w'_i$ .

#### **B** SFT Dataset Construction Details

We collect and integrate several large-scale instruction datasets specifically designed for fine-tuning LLMs. Datasets are shown below:

- WildChat-1M-Full (Zhao et al., 2024) is a dataset comprising 1 million conversations between human users and ChatGPT, enriched with demographic metadata such as state, country, hashed IP addresses, and request headers.
- Infinity-Instruct (BAAI, 2024) is a large-scale, high-quality instruction dataset, specifically designed to enhance the instruction-following capabilities of LLMs in both general and domain-specific tasks.
- tulu-3-sft-mixture (Lambert et al., 2024) is used to train the Tulu 3 series of models
- orca-agentinstruct-1M-v1-cleaned <sup>6</sup> is a cleaned version of the orca-agentinstruct-1M-v1 (Mitra et al., 2024) dataset released by Microsoft, a fully synthetic dataset using only raw text publicly available on the web as seed data.
- **Imsys-chat-1m** (Zheng et al., 2023) is a comprehensive real-world conversational dataset containing one million interactions with 25 LLMs. This dataset spans a wide range of topics and interaction types, effectively capturing diverse user-LLM interaction patterns.

<sup>&</sup>lt;sup>6</sup>https://huggingface.co/datasets/mlabonne/orca-agentinstruct-1M-v1-cleaned

Together, they comprise 11,425,231 samples, forming a robust and diverse foundation for advancing 870 research on instruct LLMs. Inevitably, many datasets contain a significant amount of similar or even 871 duplicate data, which can adversely affect both model training and the accuracy of evaluations. To 872 address this issue, we employ an n-gram-based deduplication technique to preprocess the data (Algorithm 873 1). N-gram method decomposes text into consecutive sequences of n words (or characters), effectively 874 capturing local features. 875

<b>Output:</b> Deduplicated dataset $\mathcal{D}_{dedup}$	
1: $\mathcal{D}_{dedup} \leftarrow \{\}$ # Initialize deduplicated dataset	
2: $seen\_hashes \leftarrow \{\}$ # Set to store hashes of seen N-grams	
3: for each sample $s$ in $\mathcal{D}$ do	
4: $ngrams \leftarrow \{\}$ # Initialize N-grams for the sample	
5: for each conversation $c$ in $s$ .conversations do	
6: $ngrams \leftarrow ngrams \cup \text{GenerateNGrams}(c.content, n)$	
7: end for	
8: <b>if</b> any $\operatorname{Hash}(ngram) \in seen\_hashes$ for $ngram \in ngrams$ <b>then</b>	
9: <b>continue</b> #Skip sample if any N-gram hash is already seen	
10: <b>end if</b>	
11: $seen\_hashes \leftarrow seen\_hashes \cup \{ Hash(ngram) \mid ngram \in ngrams \}$	•
12: $\mathcal{D}_{dedup} \leftarrow \mathcal{D}_{dedup} \cup \{s\}$	
13: end for	
14: return $\mathcal{D}_{dedup}$	

This approach enables the detection and identification of repetitive patterns within the text. By leveraging this method, we are able to filter out not only completely identical instances but also content that exhibits high semantic or structural similarity. Consequently, the quality and diversity of the dataset are significantly enhanced. Finally, we adopt a 20-gram deduplication strategy to eliminate redundancy in the dataset. After applying this process, a total of 4,758,226 data entries are obtained.

876

877

878

879

880

881

882

#### С **Hyperparameter Settings**

Algorithm 1 Deduplicate Dataset by N-Grams

**Input:** Dataset  $\mathcal{D}$ , N-gram size n

The detailed parameter settings used in the experiment are as follows:

## **General Settings**

General Settings	883
• Learning Rate ( $lr$ ): $7 \times 10^{-5}$	884
• End Learning Rate ( $lr_{end}$ ): $7 \times 10^{-6}$	885
• Seed: 42	886
• Data Type ( <i>dtype</i> ): float32	887
Optimizer Settings	888
• Optimizer: Adam	889
- Beta 1 ( $\beta_1$ ): 0.9	890
- Beta 2 (β <sub>2</sub> ): 0.999	891
Learning Rate Scheduler: cosineannealing	892
- Learning Rate Decay Steps: 64,000	893
– Learning Rate Warm-up Steps: 16,000	894

895	• Sparsity Loss Coefficient ( <i>L</i> <sub>sparsity</sub> ):
896	– 0.01 for JumpReLU
897	– 5 for Standard
898	• Sparsity Loss Warm-up Steps ( <i>L</i> <sub>sparsity</sub> ): 10,000
899	Training Settings
900	• Training Tokens: $4.096 \times 10^7$
901	• Train Batch Size (tokens): 128
902	Activation and Decoder Initialization
903	• Decoder Initialization Method ( $b_{dec\_init\_method}$ ): geometric_median
904	Normalize SAE Decoder: True
905	• Dead Feature Threshold: $10^{-8}$
906	• Dead Feature Window: 1000
907	Additional Settings
908	• Noise Scale: 0
909	• Expansion Factor: 8 or 16
910	• Feature Sampling Window: 2000
911	• JumpReLU Bandwidth: 0.001
912	• JumpReLU Init Threshold: 0.001
913	• Apply Decoder to Input ( <i>apply_b_dec_to_input</i> ): False
914	• Use Ghost Gradients: False
915	Use Cached Activations: False

#### **D** Mean Squared Error (MSE) of SAEs

The Mean Squared Error (MSE) results for the token reconstruction task are presented in this section.

#### D.1 Mean Squared Error (MSE) of special tokens of standard SAEs



Figure 6:  $MSE_{st}$  performance of the Standard SAE (all metrics are presented in log scale, where lower values indicate better SAE reconstruction performance). Within the Standard architecture, *FAST* exhibits the best reconstruction capability compared to BT(P) and BT(F)

While the reconstruction capability of Standard SAE models was generally inferior to the JumpReLU structure, *FAST* is also able to effectively reduce the  $MSE_{st}$ , especially in the Qwen series models.

Laver	Expansion	Method	Standa	ard SAE	JumpR	eLU SAE
	Factor		$\log_2(MSE)$	$\log_2(MSE_{st})$	$\log_2(MSE)$	$\log_2(MSE_{st})$
		BT(P)	-5.5059	-4.2377	-9.4350	-6.8026
4	8	BT(F)	-5.6080	-4.8046	-9.8097	-8.3853
		FAST	-5.6432	-4.7236	-9.8187	-10.1534
		BT(P)	-3.2837	-1.6776	-11.2353	-5.4823
12	8	BT(F)	-3.3437	-2.8733	-13.9975	-9.2049
		FAST	-3.4104	-3.0011	-14.1393	-12.1287
	8	BT(P)	-1.6059	-0.6085	-13.0282	-7.4267
18		BT(F)	-1.7131	-1.6009	-15.0851	-10.4278
		FAST	-1.8697	-2.2923	-15.0666	-12.4442
	8	BT(P)	-1.1852	-0.1692	-13.3080	-7.8271
20		BT(F)	-1.3509	-1.3587	-14.7969	-10.4507
		FAST	-1.4721	-1.9375	-15.5552	-13.1463
25		BT(P)	-0.1677	1.0444	-12.9767	-7.1657
	8	BT(F)	-0.5163	-0.5639	-16.6192	-11.6569
		FAST	-0.5747	-0.8982	-16.5138	-15.9845

#### D.2 MSE of SAEs trained on Llama-3.1-8B-Instruct

Table 3: Mean Squared Error (MSE) of SAEs trained on Llama-3.1-8B-Instruct. Each value is highlighted with a green background to indicate performance, with darker shades of green representing better results.

919 920

921

916 917

Laver	Expansion	Method	Standa	ard SAE	JumpReLU SAE	
Luyer	Factor		$\log_2(MSE)$	$\log_2(MSE_{st})$	$\log_2(MSE)$	$\log_2(MSE_{st})$
		BT(P)	-4.5650	-3.8363	-13.7434	-8.3908
	8	BT(F)	-4.5785	-3.8250	-13.6105	-8.5868
4		FAST	-4.5931	-3.9053	-9.0852	-8.7193
·		BT(P)	-4.5645	-3.8158	-9.6278	-7.5321
	16	BT(F)	-4.5858	-3.8210	-9.6102	-7.6905
		FAST	-4.5959	-3.9055	-9.8054	-9.3065
	8	BT(P)	-2.6239	-1.9052	-13.4038	-8.5246
		BT(F)	-2.6757	-2.1318	-14.7879	-9.1440
12		FAST	-2.7236	-2.4763	-15.3747	-13.4614
12	16	BT(P)	-2.6279	-1.9488	-12.2827	-7.7836
		BT(F)	-2.6754	-2.2725	-13.8874	-8.4299
		FAST	-2.7509	-2.5644	-14.4420	-12.6355
		BT(P)	-0.6926	-0.4378	-13.5554	-8.4006
	8	BT(F)	-0.9186	-1.0709	-14.8424	-8.9061
20		FAST	-0.9527	-1.4473	-18.8809	-17.3707
		BT(P)	-0.8145	-0.4607	-13.1516	-9.1137
	16	BT(F)	-1.0947	-1.1447	-14.2900	-8.9611
		FAST	-1.1285	-1.5387	-14.6872	-12.1711

Table 4: Mean Squared Error (MSE) of SAEs trained on Llama-3.2-3B-Instruct. Each value is highlighted with a green background to indicate performance, with darker shades of green representing better results.

Laver	Laver Expansion		Standa	ard SAE	rd SAE   JumpReLU SAE	
Luyer	Factor		$\log_2(MSE)$	$\log_2(MSE_{st})$	$\log_2(MSE)$	$\log_2(MSE_{st})$
		BT(P)	-5.3374	-4.4021	-15.3160	-9.6296
	8	BT(F)	-5.3583	-4.4375	-15.6237	-10.0324
4		FAST	-5.3775	-4.3920	-15.8654	-13.9127
		BT(P)	-5.3370	-4.3794	-14.5574	-9.0583
	16	BT(F)	-5.3587	-4.4358	-14.7275	-9.4817
		FAST	-5.3804	-4.3879	-10.5009	-10.2448
	8	BT(P)	-3.6638	-2.9507	-7.9900	-7.2577
		BT(F)	-3.7759	-3.0874	-16.1021	-10.5349
9		FAST	-3.8282	-3.5754	-16.4928	-13.9685
,	16	BT(P)	-3.6642	-2.9456	-7.1584	-6.5155
		BT(F)	-3.8049	-3.3775	-15.1966	-9.8149
		FAST	-3.8344	-3.6778	-15.8696	-12.9629
		BT(P)	-1.2195	-0.4927	-8.0419	-5.1825
	8	BT(F)	-1.7311	-1.7559	-15.2996	-9.3409
14		FAST	-1.7410	-2.6844	-21.4449	-23.4395
		BT(P)	-1.2449	-0.5642	-6.4784	-5.2817
	16	BT(F)	-1.8371	-1.8036	-14.9445	-9.3654
		FAST	-1.8409	-2.7668	-16.2748	-13.3547

#### D.4 MSE of SAEs trained on Llama-3.2-1B-Instruct

Table 5: Mean Squared Error (MSE) of SAEs trained on Llama-3.2-1B-Instruct. Each value is highlighted with a green background to indicate performance, with darker shades of green representing better results.

Laver	Expansion	Method	Standa	Standard SAE		eLU SAE
Luyer	Factor		$\log_2(MSE)$	$\log_2(MSE_{st})$	$\log_2(MSE)$	$\log_2(MSE_{st})$
		BT(P)	1.2919	7.2207	-4.1852	1.9109
4	8	BT(F)	-0.5233	0.0494	-5.9622	-3.3368
		FAST	-0.7358	-1.6090	-10.6174	-11.9105
		BT(P)	1.4751	5.7788	-5.8014	-4.1171
12	8	BT(F)	0.7681	0.9550	-6.3039	-5.9309
		FAST	0.6177	-0.0770	-9.8207	-10.4545
	8	BT(P)	2.0024	5.1985	-6.5926	-4.0005
18		BT(F)	1.4749	1.5093	-6.8466	-8.0743
		FAST	1.3892	0.6468	-9.1659	-9.7604
	8	BT(P)	2.6772	5.1501	-4.9649	-0.7776
20		BT(F)	2.1453	1.9877	-5.6461	-3.5904
		FAST	2.0796	1.1869	-8.2213	-8.7821
25		BT(P)	4.8764	6.2532	-2.1482	2.0938
	8	BT(F)	4.4139	3.7031	-2.6957	1.6207
		FAST	4.4471	3.0934	-4.9598	-5.5615

Table 6: Mean Squared Error (MSE) of SAEs trained on Qwen2.5-7B-Instruct. Each value is highlighted with a green background to indicate performance, with darker shades of green representing better results.

D.6	MSE of SAEs	trained o	on Qwen2.5-3	<b>B-Instruct</b>
-----	-------------	-----------	--------------	-------------------

Laver	Expansion	Method	Method Standard SAE		JumpR	eLU SAE
Eujer	Factor		$\log_2(MSE)$	$\log_2(MSE_{st})$	$\log_2(MSE)$	$\log_2(MSE_{st})$
		BT(P)	-0.8873	2.8616	-8.7177	-2.2147
	8	BT(F)	-1.4572	1.1595	-8.5340	-1.9954
4		FAST	-1.5098	-1.6682	-13.9907	-11.6534
		BT(P)	-1.0058	2.8627	-8.8511	-2.3755
	16	BT(F)	-1.6685	1.1371	-8.9769	-2.4576
		FAST	-1.5147	-1.7482	-13.2162	-10.7660
	8	BT(P)	0.9257	3.0243	-9.2313	-2.9916
		BT(F)	0.4744	1.1862	-9.3796	-2.9188
18		FAST	0.6782	-0.9288	-10.3007	-11.2916
10	16	BT(P)	0.8594	3.4799	-9.6147	-3.1930
		BT(F)	0.3438	1.1729	-9.5534	-3.0426
		FAST	0.5485	-1.0730	-10.3197	-11.1114
		BT(P)	3.8883	4.7227	-4.3442	-2.3480
	8	BT(F)	3.4388	3.7056	-5.5300	-5.3856
32		FAST	3.6647	1.6953	-5.0278	-7.3022
		BT(P)	3.7736	4.6584	-4.4299	-2.9327
	16	BT(F)	3.2978	3.4334	-5.6515	-6.2729
		FAST	3.5676	1.4331	-5.0783	-7.2653

Table 7: Mean Squared Error (MSE) of SAEs trained on Qwen2.5-3B-Instruct. Each value is highlighted with a green background to indicate performance, with darker shades of green representing better results.

Laver	Expansion Method		Standa	Standard SAE		JumpReLU SAE	
Lujer	Factor		$\log_2(MSE)$	$\log_2(MSE_{st})$	$\log_2(MSE)$	$\log_2(MSE_{st})$	
		BT(P)	-0.1150	3.8222	-5.0404	1.5111	
	8	BT(F)	-0.5653	3.2719	-5.1794	1.3737	
4		FAST	-0.7745	-2.1358	-13.4069	-12.5193	
		BT(P)	-0.2315	3.8196	-4.8980	1.6550	
	16	BT(F)	-0.7614	3.2068	-5.1495	1.4045	
		FAST	-0.9958	-2.0996	-13.3622	-11.6841	
	8	BT(P)	0.4087	3.5463	-5.4990	1.0522	
		BT(F)	0.0306	2.9569	-6.2791	0.2762	
14		FAST	-0.0925	-1.2535	-11.2579	-11.8198	
11	16	BT(P)	0.3186	3.5454	-4.9561	1.5981	
		BT(F)	-0.0918	3.0073	-5.9567	0.5989	
		FAST	-0.2312	-1.3543	-11.6309	-12.1911	
		BT(P)	3.0506	4.3907	-4.6425	0.4759	
	8	BT(F)	2.5424	3.5608	-5.3630	0.5141	
24		FAST	2.5122	0.6336	-6.2603	-7.9484	
		BT(P)	2.9411	4.3725	-4.4566	1.1218	
	16	BT(F)	2.3877	3.5499	-5.0298	1.0916	
		FAST	2.3762	0.3794	-6.3063	-8.0686	

Table 8: Mean Squared Error (MSE) of SAEs trained on Qwen2.5-1.5B-Instruct. Each value is highlighted with a green background to indicate performance, with darker shades of green representing better results.

Laver	Expansion	Method	Standa	ard SAE	JumpR	eLU SAE
Lujei	Factor		$\log_2(MSE)$	$\log_2(MSE_{st})$	$\log_2(MSE)$	$\log_2(MSE_{st})$
		BT(P)	-2.7554	-0.1257	-10.6725	-4.1202
	8	BT(F)	-2.8808	-1.3213	-11.6763	-5.1212
4		FAST	-2.8732	-3.2218	-21.7343	-23.1697
-		BT(P)	-2.9204	-0.0721	-10.7024	-4.1569
	16	BT(F)	-3.1034	-1.1148	-11.6959	-5.1497
		FAST	-3.0970	-3.2153	-17.4590	-16.7389
	8	BT(P)	-2.0463	-0.0492	-9.5392	-2.9978
		BT(F)	-2.2811	-1.1008	-10.4276	-3.8743
12		FAST	-2.2836	-3.0505	-21.1734	-25.6605
12		BT(P)	-2.1648	-0.0915	-9.4019	-2.8551
	16	BT(F)	-2.4489	-1.1418	-10.5582	-4.0043
		FAST	-2.4406	-3.0602	-20.7499	-19.0931
		BT(P)	0.2408	1.3303	-10.5099	-4.2017
	8	BT(F)	-0.3029	-0.0174	-11.4078	-4.8666
20		FAST	-0.3387	-1.9461	-15.2442	-16.9599
20		BT(P)	0.1296	1.2181	-10.6728	-4.2739
	16	BT(F)	-0.4536	-0.0825	-11.3337	-4.7864
		FAST	-0.4924	-2.1033	-16.3662	-18.0564

#### D.7 MSE of SAEs trained on Qwen2.5-0.5B-Instruct

Table 9: Mean Squared Error (MSE) of SAEs trained on Qwen2.5-0.5B-Instruct. The best and second-best methods are highlighted with dark green and light green backgrounds, respectively.

927

## **E** Implementation Details of Feature Interpretability

This section provides a detailed explanation of the implementation process for evaluating and interpreting feature interpretability.

#### E.1 SAEs for Feature Interpretability

Model Name	Layer	<b>Expansion Factor</b>
Llama series		
Llama-3.1-8B-Instruct	18	8X
Llama-3.2-3B-Instruct	12	8X&16X
Llama-3.2-1B-Instruct	9	8X&16X
Qwen series		
Qwen2.5-7B-Instruct	18	8X
Qwen2.5-3B-Instruct	18	8X&16X
Qwen2.5-1.5B-Instruct	14	8X&16X
Qwen2.5-0.5B-Instruct	12	8X&16X

Table 10: Model configurations of the Llama and Qwen model series.

#### E.2 Prompt for Feature Interpretability

#### System Prompt

We are analyzing the activation levels of features in a neural network. Each feature activates specific tokens in a text, and the activation value of each token indicates its relevance to the feature. Higher activation values signify a stronger association.

Your task is to evaluate the feature based on the following scoring rubric and assign it a monosemanticity score.

### Scoring Rubric: Activation Consistency

1: No discernible pattern

2: Broad consistent theme but lacking structure

3: Clear overall pattern but quite a few examples not fitting that pattern

4: Clear pattern with one or two deviating examples

5: Clear pattern with no deviating examples

### Instructions:

1. Analyze the context provided, which consists of a sequence of alternating tokens and their corresponding activation values.

2. Assign a score based on the activation consistency rubric.

3. Provide a descriptive name for the feature that captures its essence.

Example output: 'My final verdict score is: [[3]], feature name is [[Mathematical Problem Explanation]]'.

User: {prompt}

#### **Prompt Template**

Below is the context of feature {feature\_index}, represented as sentences with tokens and their activation values: {context}

934

928

929

930

#### F Implementation Details of Steering with SAE Latents

935 936

#### F.1 10 Questions

Question 1:

How do I export constants and classes from a JavaScript module? **Ouestion 2:** 

FINAL EXAM Question 1. What was Elsie Marley profession?

Question 3:

lettre de mativation **Ouestion 4:** 

请回答以下问题,找出铁锤和磁铁之间的主要区别是什么?

#### Question 5:

Summarize this article in one sentence.\n\nMedia playback is not supported on this device\nFarah pulled away from American Dathan Ritzenhein in the last mile in his first race since retaining his 5,000m and 10,000m Olympic titles in Rio.\nIn the women\"s race, Olympic 5,000m champion Vivian Cheruiyot of Kenya won in her first half marathon.\nScotland\"s Mark Telford took the men\"s wheelchair crown, a second ahead of fellow Briton Bret Crossley.\nThe Great North Run is the world\"s biggest half marathon and there were more than 41,000 runners taking part in this year\"s event from 178 nations.\nFarah, 33, was taken on a fast pace by former American 5,000m record holder Ritzenhein, but powered away with a mile to go and even had time to do a cheeky heel flip before he crossed the line in one hour and four seconds, the slowest of his three wins.\nBelgium\"s Emmanuel Bett, who ran the second half of the race almost on his own, crossed in third.\nFarah told BBC Sport: \"To be honest with you, I\"m knackered.\n\"I knew I had to work hard because Dathan is a former training partner and was running a great race.\n\"He wonen\"s race was billed as a shoot-out between middle distance greats Cheruiyot and there-time Olympic champion Tirunesh Dibaba.\nDibaba failed to keep pace with Cheruiyot and fellow Kenyan Priscah Jeptoo in the closing stages of the 13.1-mile course. It was Cheruiyot who took victory, producing a sprint finish to clock 1:07.54, just one second ahead of Jeptoo. \n\Cheruiyot said: \"I\"m so happy because it\"s my birthday. I found it tough with one kilometre to go but it\"s fantastic for me to end my season this way. \"\nMedia playback is not supported on this device\n\Summary:

#### **Question 6:**

def intersection(list1, list2):\n \"\"\"\" This function returns a list of common elements between two lists: list1 and list2.\n\n Parameters:\n list1 (list): First list of elements\n list2 (list): Second list of elements\n\n Returns:\n list: A list of common elements between list1 and list2\n \"\"\" # Complete the code to find the intersection of list1 and list2 using nested for loops\n common\_elements = []\n for element1 in list1:\n for element2 in list2.\n if element1 == element2:\n common elements.append(element1)\n break\n\n return common elements

#### **Question 7:**

Article:\nAuthor Carol Dunbar stands outside of her writing studio that is under her family's water tower on their property deep in the woods south of Superior, Wisconsin, Jed Carlson / Superior Telegram\nCarol Dunbar works at the desk inside her writing studio in her family's water tower in the woods south of Superior on Tuesday afternoon, Oct. 26, 2021. Jed Carlson / Superior Telegram\n\nCarol Dunbar speaks about her love of nature and living off the grid on her family's property deep in the woods south of Superior on Tuesday afternoon, Oct. 26, 2021. Jed Carlson / Superior Telegram/n/nCarol Dunbar looks through short writings and other trinkets that were her grandmothers that she keeps in her writing studio in her family's water tower in the woods south of Superior on Tuesday afternoon, Oct. 26, 2021. Jed Carlson / Superior Telegram\n\nAuthor Carol Dunbar talks about her struggles with the editing process of her upcoming novel as she sits in her writing studio in the family's water tower south of Superior on Tuesday afternoon, Oct. 26. 2021. Jed Carlson / Superior Telegram\nAuthor Carol Dunbar looks out of one of the windows of her writing studio in her family's water tower on their property deep in the woods south of Superior on Tuesday afternoon, Oct. 26, 2021. Jed Carlson / Superior Telegram\nNUPERIOR, - Carol Dunbar stepped through the woods as fallen leaves crunched beneath her feet. Her homestead south of Superior includes the main residence, her husband's workshop and a water tower. Living off the grid, the structure is a necessity for the homestead's water pressure and for Dunbar's work.\n\n"Me getting into this water tower was finding a space where I could shut a door behind me to create," she said. "I wouldn't want any other kind of office, but it definitely has its challenges "\n\nThe novelist and freelance ghostwriter's computers, manuscripts and books all reside under what some might consider to be their worst enemy: "There are literally two 250-gallon tanks of water over my head right now," she said.\n\nYes, her office has flooded several times \n\n"It's like being in a room that's pouring rain. It's awful, and I've had to make peace with that."\n\nTo see her work be so vulnerable makes it that much more endearing. "I know there's a really interesting metaphor about art and risk," she added.\n\nThere's no other space on their 80 acres where she can work the way she's able to here. After numerous floods and years spent working from the living room, her husband redid the space and built the staircase for better access and heat circulation.\n\nOriginally intended as a guest room, it's a 10-by-10 space on the second floor of the water tower. She calls it the cockpit.\n\nThere's a porch on the back and windows on all four sides, so "I feel like I'm writing in the treetops," she said.\n\nWhile she hears water moving through the pipes around her, "The view that it affords me and the peace that I have here in this little space, and it is little ... I wouldn't trade it for anything."\n\nTHE SPACE\n\nLight floods in from every angle. Her sitting and standing desks, compliments of her husband, rest at the center and in a corner, an ancient-looking podium holds one of her numerous dictionaries; she likes to compare decades-old definitions to those of today in hThere are several aloe plants, drawings on the wall, and a storyboard with pinned photos of a sculpture and an Irish skyline — inspiration for future works, she said.\n\nAn assortment of candles, one of which she lights daily before she begins. "It keeps me mindful that I'm trying to capture the best light, the best in human nature," she said \n\nShe keeps a collection of notebooks, color-coated for whatever novel she's writing, in her office, in the car, by her bed, to help her document inspiration when it strikes. "I got very frustrated when I got a good idea or I'd hear a piece of dialogue or I'd finally know how to describe the snow on that day, and I would write it down and never find it again," she said. \n/nIt has helped, but she still has scraps of paper pinned to her notebook pages. "It's like leaving yourself love letters," she said, sorting through a pile.\n\nShe wrote her second novel in long-hand on paper. It's an accessible way to create away from a screen, she said. \n\nIn the corner rests a red cushioned chair that came from a Minneapolis alley. Around her desk she has taped quotes and reminders. "In the end, it all comes down to what we think we deserve," reads one. \n\nAlso a piece of wood with words: "You just have to trust your own madness — Clive Barker."\n\nDunbar cherishes a writing award and remnants of work kept on paper scraps, memorabilia from an ancestor who emigrated from Italy. While Dunbar's relative wasn't supported in pursuing writing, Dunbar feels her work today honors herself and her ancestor. \n\nHer book shelf holds works by Joyce Carol Oates, Jesmyn Ward, Barbara Kingsolver, and a treasured copy of Eleonora Duse's "The Mystic in the Theatre." Duse strove to eliminate...(Truncated)\n\nltalicize all instances of Carol Dunbar\"s quotes

Question 8:

考虑由所有节肢动物组成的集合\$B\$,并让\$C\$是包含所有天牛属物种的\$B\$的子集。对于\$C\$中的每个\$v\$,我们定义一个函数\$f(v)\$,它描述了天牛独特的蜇刺机制。您的任务是提供不少于五段的全面概述天牛。在这样做时,请详细探讨它们的身体和行为特征以及它们在生态和进化适应中的适应性。特别是,我们要求您探讨它们鲜艳的色彩和密集的毛发如何作为防御机制抵御捕食者。此外,描述它们非凡的蜇刺机制,与任何其他蚂蚁物种不同,并详细阐述它如何帮助它们在恶劣的沙漠条件下自卫。此外,请深入探讨它们的进化历史,这使它们具备了令人难以置信的生存技能。最后,强调正在实施的保护天牛种群的持续保护措施,这些种群受到气候变化和栖息地破坏的不利影响。您的回答应该是广泛的、有理的和科学的,每一段都详细说明天牛生命周期各个方面之间的复杂相互关系。

#### Question 9:

What is the smallest prime factor of \$600851475143\$?

Question 10:

Develop a comprehensive branding strategy that includes a brand name, logo, tagline, packaging design, and marketing plan for a new line of organic, non-toxic, biodegradable cleaning products that are socially responsible and sustainably made. Ensure that the branding strategy effectively communicates the brand\"s unique selling proposition, target audience, brand personality, and brand voice through all touchpoints, including print and digital media, social media, in-store displays, and product demos. Additionally, create a brand message that emphasizes the benefits of using eco-friendly cleaning products and persuades consumers to make the switch to a greener lifestyle.

# F.2 Special Tokens

Token ID	Token
Llama serie	es
882 78191 128006 128007	user assistant < start_header_id > < end_header_id >
128009 Qwen serie	<pre><leot_id > /// // // // // // // // // // // // //</leot_id ></pre>
872 77091 151644 151645	user assistant < im_start > < im_end >

Table 11: Tokens that control response generation and formatting in the Llama and Qwen model series.

# F.3 Average Top 5 Max Activation Values and Their Corresponding Indices for Tokens across a 1000-Sample Dataset

Approach	Token	Top 5 Max Activation Value (Index:Value)
	882	<b>4453</b> :0.8120 <b>30511</b> :0.724 <b>18547</b> :0.597 <b>19110</b> :0.500 <b>20505</b> :0.469
	78191	<b>5188</b> :0.5030 <b>1923</b> :0.4900 <b>31873</b> :0.486 <b>20505</b> :0.468 <b>3187</b> :0.4620
BT(P)[8X]	128006	<b>2604</b> :7.1220 <b>20523</b> :0.800 <b>7428</b> :0.7330 <b>24017</b> :0.702 <b>16640</b> :0.678
	128007	<b>23901</b> :1.193 <b>7808</b> :0.5210 <b>3268</b> :0.5180 <b>20505</b> :0.477 <b>30244</b> :0.473
	128009	<b>20505</b> :0.744 <b>25940</b> :0.653 <b>7961</b> :0.6460 <b>21317</b> :0.585 <b>19110</b> :0.569
	882	<b>11765</b> :0.823 <b>25025</b> :0.814 <b>7043</b> :0.6880 <b>16826</b> :0.562 <b>21896</b> :0.560
	78191	<b>30553</b> :0.536 <b>9728</b> :0.5270 <b>11435</b> :0.507 <b>14565</b> :0.505 <b>13234</b> :0.497
BT(F)[8X]	128006	17784:7.480 17355:0.947 28634:0.782 9333:0.7710 27149:0.744
	128007	<b>23677</b> :1.002 <b>6426</b> :0.6680 <b>26136</b> :0.603 <b>5783</b> :0.5720 <b>26958</b> :0.526
	128009	<b>23677</b> :0.834 <b>7100</b> :0.7560 <b>30568</b> :0.734 <b>15188</b> :0.666 <b>8346</b> :0.6430
	882	<b>22534</b> :0.611 <b>13320</b> :0.470 <b>29165</b> :0.464 <b>19871</b> :0.428 <b>29033</b> :0.418
	78191	<b>16063</b> :0.463 <b>13320</b> :0.461 <b>19871</b> :0.460 <b>32613</b> :0.441 <b>22277</b> :0.399
FAST[8X]	128006	<b>22642</b> :4.392 <b>2417</b> :0.7170 <b>27839</b> :0.706 <b>3095</b> :0.7030 <b>10814</b> :0.654
	128007	<b>30457</b> :2.489 <b>19871</b> :0.532 <b>6870</b> :0.4640 <b>28096</b> :0.446 <b>13266</b> :0.413
	128009	<b>13822</b> :0.753 <b>22277</b> :0.606 <b>21866</b> :0.537 <b>17489</b> :0.493 <b>118</b> :0.41200

Table 12: Top 5 Average Activation Values for Special Tokens in Llama3.1-8B-instruct with JumpReLU SAE

938

Approach	Token ID	Top 5 Max Activation Value (Index:Value)
	882	<b>3817</b> :0.4550 <b>11734</b> :0.430 <b>505</b> :0.42200 <b>23884</b> :0.417 <b>14851</b> :0.380
	78191	<b>6451</b> :0.3460 <b>11061</b> :0.340 <b>19811</b> :0.327 <b>12369</b> :0.325 <b>11734</b> :0.308
BT(P)[8X]	128006	<b>2064</b> :20.351 <b>5699</b> :0.4090 <b>14393</b> :0.399 <b>7505</b> :0.3770 <b>548</b> :0.37500
	128007	<b>20232</b> :0.427 <b>5095</b> :0.4000 <b>19583</b> :0.393 <b>23908</b> :0.362 <b>3719</b> :0.3590
	128009	<b>14536</b> :0.468 <b>16718</b> :0.437 <b>23736</b> :0.413 <b>13925</b> :0.379 <b>10211</b> :0.368
	882	<b>23287</b> :0.814 <b>44336</b> :0.718 <b>10727</b> :0.712 <b>11701</b> :0.683 <b>26467</b> :0.658
	78191	<b>34602</b> :0.622 <b>10655</b> :0.600 <b>45414</b> :0.591 <b>23156</b> :0.553 <b>19333</b> :0.522
BT(P)[16X]	128006	<b>38076</b> :28.41 <b>48766</b> :0.675 <b>16639</b> :0.659 <b>28134</b> :0.653 <b>45</b> :0.621000
	128007	<b>9822</b> :0.7530 <b>39737</b> :0.659 <b>5712</b> :0.6430 <b>38496</b> :0.574 <b>23156</b> :0.570
	128009	<b>483</b> :0.79800 <b>48233</b> :0.789 <b>22660</b> :0.670 <b>24339</b> :0.624 <b>23774</b> :0.600
	882	<b>21524</b> :0.496 <b>17981</b> :0.471 <b>10125</b> :0.436 <b>11210</b> :0.431 <b>14456</b> :0.410
	78191	<b>16126</b> :0.447 <b>8704</b> :0.4470 <b>20691</b> :0.418 <b>19630</b> :0.393 <b>10125</b> :0.365
BT(F)[8X]	128006	<b>15765</b> :21.39 <b>1640</b> :0.5180 <b>14456</b> :0.479 <b>45</b> :0.459000 <b>17981</b> :0.442
	128007	<b>7814</b> :0.5120 <b>24565</b> :0.489 <b>1759</b> :0.4840 <b>8704</b> :0.4390 <b>14456</b> :0.396
	128009	<b>5506</b> :0.5230 <b>20691</b> :0.514 <b>20328</b> :0.488 <b>6878</b> :0.4550 <b>7593</b> :0.4460
	882	<b>20561</b> :0.719 <b>28995</b> :0.698 <b>14625</b> :0.662 <b>32041</b> :0.625 <b>4844</b> :0.5850
	78191	<b>23154</b> :0.725 <b>8239</b> :0.6700 <b>45582</b> :0.630 <b>23594</b> :0.593 <b>11425</b> :0.564
BT(F)[16X]	128006	<b>30984</b> :25.38 <b>10207</b> :0.752 <b>21441</b> :0.751 <b>26876</b> :0.700 <b>35477</b> :0.683
	128007	<b>41219</b> :0.687 <b>14625</b> :0.670 <b>21050</b> :0.662 <b>23942</b> :0.621 <b>27267</b> :0.595
	128009	<b>26876</b> :0.761 <b>13612</b> :0.722 <b>9537</b> :0.6930 <b>44518</b> :0.653 <b>6317</b> :0.6240
	882	<b>2950</b> :0.5730 <b>1343</b> :0.5670 <b>16808</b> :0.498 <b>19508</b> :0.481 <b>5931</b> :0.4590
	78191	<b>23183</b> :0.548 <b>263</b> :0.50900 <b>8564</b> :0.4860 <b>2680</b> :0.4750 <b>23798</b> :0.472
FAST[8X]	128006	<b>8772</b> :37.471 <b>20896</b> :0.610 <b>2950</b> :0.6060 <b>12126</b> :0.538 <b>16622</b> :0.534
	128007	<b>12955</b> :0.550 <b>22995</b> :0.536 <b>3339</b> :0.5080 <b>7878</b> :0.4970 <b>2950</b> :0.4730
	128009	<b>7814</b> :0.5850 <b>16940</b> :0.551 <b>4605</b> :0.5080 <b>12331</b> :0.493 <b>4439</b> :0.4880
	882	<b>9447</b> :0.8380 <b>5861</b> :0.7210 <b>19741</b> :0.716 <b>22320</b> :0.669 <b>25160</b> :0.645
	78191	<b>4177</b> :0.8220 <b>43897</b> :0.719 <b>18009</b> :0.667 <b>25117</b> :0.594 <b>30970</b> :0.590
FAST[16X]	128006	<b>22974</b> :37.66 <b>36</b> :0.873000 <b>18075</b> :0.813 <b>26318</b> :0.774 <b>45047</b> :0.762
	128007	<b>42421</b> :0.798 <b>655</b> :0.75300 <b>13955</b> :0.697 <b>26318</b> :0.632 <b>28994</b> :0.589
	128009	<b>29041</b> :0.888 <b>18075</b> :0.844 <b>33332</b> :0.776 <b>2705</b> :0.7120 <b>26318</b> :0.695

Table 13: Top 5 Average Activation Values for Special Tokens in Llama3.2-3B-instruct with JumpReLU SAE

Approach	Token ID	Top 5 Max Activation Value (Index:Value)
	882	<b>12248</b> :0.455 <b>14322</b> :0.446 <b>10030</b> :0.444 <b>11886</b> :0.425 <b>731</b> :0.39800
	78191	<b>14903</b> :0.443 <b>15672</b> :0.435 <b>8014</b> :0.4190 <b>13261</b> :0.410 <b>11985</b> :0.405
BT(P)[8X]	128006	<b>4464</b> :10.463 <b>4858</b> :0.4600 <b>12143</b> :0.454 <b>9898</b> :0.4440 <b>6877</b> :0.3700
	128007	<b>196</b> :0.45400 <b>15332</b> :0.398 <b>9561</b> :0.3580 <b>12143</b> :0.355 <b>626</b> :0.35500
	128009	<b>15332</b> :0.496 <b>1296</b> :0.4910 <b>4858</b> :0.4170 <b>6877</b> :0.4170 <b>15975</b> :0.412
	882	<b>20612</b> :0.642 <b>22827</b> :0.613 <b>3012</b> :0.6050 <b>11176</b> :0.578 <b>2141</b> :0.5760
	78191	<b>28423</b> :0.672 <b>24765</b> :0.661 <b>30621</b> :0.649 <b>22827</b> :0.649 <b>18585</b> :0.621
BT(P)[16X]	128006	<b>4169</b> :11.460 <b>11176</b> :0.793 <b>9495</b> :0.6770 <b>9911</b> :0.6730 <b>24072</b> :0.586
	128007	<b>26090</b> :0.820 <b>10861</b> :0.622 <b>24072</b> :0.615 <b>26939</b> :0.591 <b>23109</b> :0.541
	128009	<b>11176</b> :0.747 <b>16525</b> :0.716 <b>26594</b> :0.685 <b>8403</b> :0.6490 <b>15861</b> :0.633
	882	<b>2387</b> :0.4130 <b>13266</b> :0.341 <b>7778</b> :0.3090 <b>8423</b> :0.2840 <b>3682</b> :0.2800
	78191	<b>7783</b> :0.3320 <b>10427</b> :0.316 <b>8941</b> :0.3150 <b>16174</b> :0.311 <b>4764</b> :0.3080
BT(F)[8X]	128006	<b>2537</b> :9.9460 <b>15768</b> :0.382 <b>9146</b> :0.3500 <b>1604</b> :0.3440 <b>14204</b> :0.312
	128007	<b>10680</b> :0.390 <b>15478</b> :0.312 <b>8905</b> :0.3090 <b>6638</b> :0.3020 <b>15034</b> :0.284
	128009	<b>2568</b> :0.4050 <b>3528</b> :0.3860 <b>14204</b> :0.371 <b>1604</b> :0.3600 <b>15768</b> :0.313
	882	<b>24100</b> :0.530 <b>6794</b> :0.5240 <b>7848</b> :0.5230 <b>9322</b> :0.4900 <b>17577</b> :0.490
	78191	<b>12548</b> :0.583 <b>24258</b> :0.542 <b>2092</b> :0.5260 <b>2460</b> :0.4960 <b>15997</b> :0.484
BT(F)[16X]	128006	<b>4967</b> :10.559 <b>24354</b> :0.675 <b>20054</b> :0.614 <b>12136</b> :0.599 <b>12707</b> :0.537
	128007	<b>18190</b> :0.581 <b>2543</b> :0.5000 <b>23285</b> :0.499 <b>15997</b> :0.494 <b>17059</b> :0.486
	128009	<b>26830</b> :0.635 <b>17228</b> :0.623 <b>11407</b> :0.551 <b>18494</b> :0.523 <b>11681</b> :0.483
	882	<b>2926</b> :0.3780 <b>878</b> :0.35400 <b>4753</b> :0.3370 <b>10237</b> :0.336 <b>7582</b> :0.3140
	78191	<b>13371</b> :0.388 <b>14099</b> :0.376 <b>8581</b> :0.3680 <b>11313</b> :0.361 <b>5121</b> :0.3400
FAST[8X]	128006	<b>12361</b> :8.486 <b>13371</b> :0.386 <b>878</b> :0.37500 <b>129</b> :0.34900 <b>1866</b> :0.3300
	128007	<b>8581</b> :0.4120 <b>12864</b> :0.357 <b>13371</b> :0.341 <b>4478</b> :0.3380 <b>4523</b> :0.3150
	128009	<b>878</b> :0.47000 <b>11483</b> :0.408 <b>6832</b> :0.3770 <b>8581</b> :0.3690 <b>865</b> :0.34700
	882	<b>1835</b> :0.7500 <b>3851</b> :0.7100 <b>982</b> :0.60400 <b>9493</b> :0.6020 <b>8463</b> :0.4780
	78191	<b>19765</b> :0.596 <b>14393</b> :0.539 <b>28589</b> :0.512 <b>2350</b> :0.4850 <b>12592</b> :0.482
FAST[16X]	128006	<b>12329</b> :10.30 <b>9838</b> :0.6440 <b>13262</b> :0.592 <b>1450</b> :0.5260 <b>27818</b> :0.504
	128007	<b>3368</b> :0.5820 <b>31764</b> :0.568 <b>16867</b> :0.518 <b>16432</b> :0.503 <b>9648</b> :0.4590
	128009	<b>10365</b> :0.696 <b>31406</b> :0.637 <b>30028</b> :0.602 <b>15515</b> :0.574 <b>16339</b> :0.535

Table 14: Top 5 Average Activation Values for Special Tokens in Llama3.2-1B-instruct with JumpReLU SAE

Approach	Token ID	Top 5 Max Activation Value (Index:Value)
	872	<b>12461</b> :9.058 <b>439</b> :3.88000 <b>19183</b> :2.978 <b>18767</b> :2.889 <b>13685</b> :1.992
<b>DT(D)[QV]</b>	77091	<b>2547</b> :2.9330 <b>15678</b> :2.562 <b>19183</b> :2.549 <b>6508</b> :2.3290 <b>4400</b> :2.0270
DI(F)[OA]	151644	<b>12461</b> :9.193 <b>1261</b> :2.7050 <b>6508</b> :2.3060 <b>2547</b> :2.1240 <b>4400</b> :2.1140
	151645	<b>1261</b> :2.9730 <b>2547</b> :2.8640 <b>6508</b> :2.4140 <b>18778</b> :2.223 <b>13888</b> :2.118
	872	<b>4710</b> :6.3500 <b>15390</b> :3.377 <b>20684</b> :3.192 <b>25558</b> :2.937 <b>27629</b> :2.800
DT(E)[QV]	77091	<b>25558</b> :3.135 <b>27629</b> :3.061 <b>19040</b> :3.012 <b>10759</b> :2.802 <b>13257</b> :2.378
$DI(I)[0\Lambda]$	151644	<b>4710</b> :6.7170 <b>10759</b> :3.412 <b>11735</b> :3.049 <b>28219</b> :2.749 <b>26983</b> :2.596
	151645	<b>28219</b> :3.130 <b>11735</b> :2.692 <b>2174</b> :2.4670 <b>10614</b> :2.464 <b>25812</b> :2.120
	872	<b>13794</b> :37.19 <b>17783</b> :4.816 <b>20022</b> :4.519 <b>21950</b> :4.077 <b>11739</b> :4.053
EACTION	77091	<b>20022</b> :5.667 <b>11739</b> :4.352 <b>16782</b> :4.180 <b>2670</b> :3.7810 <b>13794</b> :3.731
FAST [OA]	151644	<b>13794</b> :39.87 <b>20022</b> :5.418 <b>7579</b> :4.1900 <b>3817</b> :4.1890 <b>26689</b> :4.023
	151645	<b>20022</b> :4.463 <b>2670</b> :3.6970 <b>22845</b> :3.139 <b>25469</b> :2.939 <b>9676</b> :2.6890

Table 15: Top 5 Average Activation Values for Special Tokens in Qwen2.5-7B-instruct with JumpReLU SAE

Approach	Token ID	Top 5 Max Activation Value (Index:Value)
	872	<b>11485</b> :2.756 <b>8925</b> :2.4490 <b>3645</b> :2.4130 <b>1600</b> :2.1160 <b>2801</b> :2.0860
BT(P)[8X]	77091	10992:1.911 1600:1.8300 15929:1./// 14942:1./4/ 12230:1.6//
	151644	7152:132.52 2713:2.0100 11354:1.996 15302:1.891 15795:1.885
	151645	<b>12297</b> :2.588 <b>11352</b> :2.457 <b>4096</b> :2.4520 <b>10336</b> :2.429 <b>10992</b> :2.214
	872	<b>14113</b> :2.010 <b>12080</b> :1.750 <b>18074</b> :1.739 <b>14580</b> :1.720 <b>2607</b> :1.4890
BT(P)[16X]	77091	<b>4047</b> :1.3860 <b>27294</b> :1.294 <b>3356</b> :1.2890 <b>14113</b> :1.248 <b>9469</b> :1.2420
DI(I)[10A]	151644	<b>32641</b> :150.0 <b>14113</b> :1.362 <b>7224</b> :1.3340 <b>28068</b> :1.327 <b>4741</b> :1.2860
	151645	<b>23725</b> :1.696 <b>14113</b> :1.674 <b>25421</b> :1.669 <b>68</b> :1.619000 <b>9469</b> :1.5140
	872	<b>7603</b> :2.8380 <b>3184</b> :2.7840 <b>15060</b> :2.777 <b>8391</b> :2.7390 <b>6484</b> :2.3780
DT(E)[OV]	77091	<b>15060</b> :3.175 <b>3373</b> :2.3530 <b>7293</b> :2.3480 <b>1317</b> :2.3398 <b>7603</b> :2.2900
<i>ΔΙ(Γ)</i> [δΛ]	151644	<b>16236</b> :121.2 <b>16225</b> :2.563 <b>7603</b> :2.5000 <b>7189</b> :2.4970 <b>958</b> :2.43000
	151645	<b>3104</b> :3.9910 <b>1317</b> :3.4210 <b>16225</b> :3.397 <b>6700</b> :3.3500 <b>15704</b> :3.101
	872	<b>23210</b> :2.320 <b>29265</b> :1.807 <b>11930</b> :1.767 <b>28994</b> :1.712 <b>2757</b> :1.5020
DT(E)[16V]	77091	<b>23210</b> :1.844 <b>6805</b> :1.6570 <b>20713</b> :1.564 <b>11930</b> :1.544 <b>29265</b> :1.483
$DI(\Gamma)[10\Lambda]$	151644	<b>31443</b> :153.4 <b>23210</b> :2.160 <b>5146</b> :2.0010 <b>24831</b> :1.894 <b>29265</b> :1.859
	151645	<b>5146</b> :2.9880 <b>5924</b> :2.4320 <b>5572</b> :2.3420 <b>12821</b> :2.078 <b>24491</b> :1.502
	872	<b>2941</b> :3.4410 <b>8775</b> :2.6400 <b>10076</b> :2.625 <b>12216</b> :2.178 <b>776</b> :1.99600
EA CTIOVI	77091	<b>2653</b> :3.6370 <b>10076</b> :3.450 <b>3411</b> :3.0540 <b>9785</b> :2.5100 <b>11618</b> :2.004
TASI [0A]	151644	8775:248.36 12291:2.880 10076:2.829 3411:2.8280 13964:2.566
	151645	<b>10076</b> :4.538 <b>12216</b> :3.775 <b>12139</b> :3.729 <b>4383</b> :3.5920 <b>12209</b> :3.279
	872	<b>6863</b> :3.7600 <b>9230</b> :2.9510 <b>20605</b> :2.446 <b>21312</b> :2.285 <b>17408</b> :2.063
EA ST[16V]	77091	<b>23681</b> :4.223 <b>6863</b> :3.9440 <b>17147</b> :3.059 <b>10035</b> :2.969 <b>4751</b> :2.7968
ГАЗІ [10 <b>A</b> ]	151644	<b>31443</b> :85.35 <b>5599</b> :1.5974 <b>9299</b> :1.5341 <b>18964</b> :1.445 <b>4751</b> :1.4220
	151645	<b>23681</b> :3.000 <b>6863</b> :2.4390 <b>20511</b> :2.173 <b>9230</b> :1.8215 <b>17147</b> :1.517

Table 16: Top 5 Average Activation Values for Special Tokens in Qwen2.5-3B-instruct with JumpReLU SAE

Approach	Token ID	Top 5 Max Activation Value (Index:Value)
	872	<b>734</b> :312.441 <b>2664</b> :2.5160 <b>576</b> :2.31600 <b>4162</b> :2.1050 <b>9629</b> :2.1030
<i>BT(P)</i> [8X]	77091	<b>1656</b> :2.2670 <b>3248</b> :2.2090 <b>4162</b> :2.1040 <b>4098</b> :2.0910 <b>8997</b> :2.0460
	151644	<b>734</b> :288.485 <b>391</b> :1.92500 <b>5536</b> :1.9240 <b>11982</b> :1.660 <b>11102</b> :1.625
	151645	<b>11322</b> :1.905 <b>734</b> :1.74500 <b>1263</b> :1.6030 <b>9637</b> :1.5900 <b>12143</b> :1.499
	872	<b>15738</b> :261.8 <b>3080</b> :1.4920 <b>2724</b> :1.3730 <b>19787</b> :1.372 <b>17743</b> :1.258
$\mathbf{PT}(\mathbf{D})$ [16 <b>V</b> ]	77091	<b>2724</b> :1.3720 <b>17351</b> :1.354 <b>1954</b> :1.3340 <b>19787</b> :1.307 <b>9767</b> :1.2760
$DI(I)[10\Lambda]$	151644	<b>15738</b> :241.5 <b>13157</b> :1.148 <b>13486</b> :1.116 <b>14339</b> :0.945 <b>6977</b> :0.9250
	151645	<b>9971</b> :1.1270 <b>22929</b> :1.032 <b>14028</b> :1.003 <b>19840</b> :0.936 <b>22072</b> :0.864
	872	<b>1910</b> :255.40 <b>7039</b> :2.5590 <b>9420</b> :2.5300 <b>8118</b> :2.4710 <b>1693</b> :2.4060
DT(E)[OV]	77091	<b>8118</b> :2.7040 <b>7067</b> :2.5230 <b>1223</b> :2.4890 <b>7039</b> :2.4670 <b>4086</b> :2.4190
$DI(\Gamma)[0\Lambda]$	151644	<b>1910</b> :234.85 <b>4798</b> :1.9970 <b>6153</b> :1.8900 <b>5905</b> :1.7000 <b>11021</b> :1.682
	151645	<b>10536</b> :1.870 <b>11021</b> :1.724 <b>7064</b> :1.6550 <b>1787</b> :1.5630 <b>6153</b> :1.5040
	872	<b>2077</b> :263.49 <b>13135</b> :1.624 <b>17747</b> :1.439 <b>16136</b> :1.353 <b>19975</b> :1.338
DT(E)[16V]	77091	<b>6886</b> :1.5170 <b>19975</b> :1.508 <b>17747</b> :1.500 <b>18492</b> :1.296 <b>16136</b> :1.249
DI(F)[10A]	151644	<b>2077</b> :242.06 <b>19387</b> :1.534 <b>4177</b> :1.3580 <b>22526</b> :1.283 <b>19497</b> :1.178
	151645	<b>4177</b> :1.1610 <b>5724</b> :1.1000 <b>9985</b> :1.0890 <b>6552</b> :1.0190 <b>11894</b> :0.945
	872	<b>7505</b> :462.49 <b>4918</b> :2.4010 <b>4694</b> :2.3060 <b>4141</b> :2.1620 <b>10728</b> :2.098
EACTION	77091	<b>491</b> :2.25800 <b>4141</b> :2.2300 <b>11303</b> :2.125 <b>8603</b> :2.0090 <b>6358</b> :1.9430
FASI [8A]	151644	<b>7505</b> :425.73 <b>10900</b> :1.793 <b>6473</b> :1.7560 <b>10139</b> :1.614 <b>2006</b> :1.5990
	151645	<b>491</b> :2.20100 <b>11115</b> :1.748 <b>11252</b> :1.665 <b>6473</b> :1.5530 <b>10257</b> :1.326
	872	<b>21852</b> :580.0 <b>11515</b> :1.988 <b>9360</b> :1.5720 <b>21118</b> :1.501 <b>11834</b> :1.487
EAST[16V]	77091	<b>21118</b> :2.068 <b>9718</b> :1.6120 <b>14362</b> :1.536 <b>9360</b> :1.5240 <b>11834</b> :1.477
FASI [10Å]	151644	<b>21852</b> :532.9 <b>21118</b> :1.683 <b>16522</b> :1.350 <b>17617</b> :1.265 <b>12233</b> :1.174
	151645	<b>21118</b> :2.070 <b>17617</b> :1.474 <b>16522</b> :1.312 <b>18955</b> :1.196 <b>21139</b> :1.084

Table 17: Top 5 Average Activation Values for Special Tokens in Qwen2.5-1.5B-instruct with JumpReLU SAE

Approach	Token ID	Top 5 Max Activation Value (Index:Value)
	872	<b>6091</b> :1.0680 <b>2897</b> :0.8250 <b>1389</b> :0.8240 <b>6239</b> :0.8150 <b>6434</b> :0.7770
<i>DT(</i> D)[ <b>Ο</b> V]	77091	<b>3245</b> :0.8430 <b>1767</b> :0.8430 <b>1389</b> :0.8310 <b>5981</b> :0.8120 <b>6239</b> :0.7790
$DI(P)[\delta \Lambda]$	151644	<b>1608</b> :43.209 <b>6818</b> :0.7600 <b>6245</b> :0.7480 <b>6724</b> :0.7150 <b>1235</b> :0.7150
	151645	<b>4541</b> :0.8170 <b>5212</b> :0.8010 <b>1744</b> :0.7760 <b>4498</b> :0.7280 <b>507</b> :0.72400
	872	<b>8475</b> :0.6880 <b>13976</b> :0.545 <b>889</b> :0.51000 <b>8786</b> :0.4680 <b>3099</b> :0.4680
DT(D)[16V]	77091	<b>3099</b> :0.5480 <b>9308</b> :0.5340 <b>13976</b> :0.528 <b>8786</b> :0.4830 <b>432</b> :0.46500
$DI(\Gamma)[10\Lambda]$	151644	<b>10161</b> :28.27 <b>7726</b> :0.4830 <b>6509</b> :0.4550 <b>9343</b> :0.4510 <b>6947</b> :0.4260
	151645	<b>1934</b> :0.5580 <b>12380</b> :0.505 <b>7726</b> :0.4370 <b>7385</b> :0.4370 <b>1823</b> :0.4280
	872	<b>5375</b> :1.0290 <b>3317</b> :0.9000 <b>4825</b> :0.8510 <b>3896</b> :0.8360 <b>5791</b> :0.8260
DT(E)[9V]	77091	<b>3896</b> :0.8510 <b>4825</b> :0.8450 <b>2552</b> :0.8420 <b>5375</b> :0.8030 <b>3203</b> :0.8010
$DI(\Gamma)[0\Lambda]$	151644	<b>2428</b> :40.999 <b>5130</b> :0.7510 <b>1326</b> :0.7050 <b>557</b> :0.68100 <b>2765</b> :0.6540
	151645	<b>2734</b> :0.8970 <b>6507</b> :0.7080 <b>628</b> :0.69600 <b>2913</b> :0.6930 <b>1119</b> :0.6680
	872	<b>13102</b> :0.658 <b>12215</b> :0.572 <b>10208</b> :0.542 <b>6285</b> :0.4670 <b>5598</b> :0.4430
PT(F)[16V]	77091	<b>7823</b> :0.5860 <b>12215</b> :0.580 <b>10208</b> :0.551 <b>12606</b> :0.521 <b>5598</b> :0.4871
$DI(I)[10\Lambda]$	151644	<b>1983</b> :27.761 <b>5393</b> :0.5180 <b>12215</b> :0.458 <b>5515</b> :0.4470 <b>9460</b> :0.4360
	151645	<b>4484</b> :0.4980 <b>12615</b> :0.472 <b>13322</b> :0.441 <b>5393</b> :0.4370 <b>8592</b> :0.3820
	872	<b>1299</b> :0.9310 <b>2747</b> :0.9090 <b>3288</b> :0.8170 <b>1859</b> :0.7860 <b>4804</b> :0.7210
	77091	<b>6296</b> :0.8960 <b>6776</b> :0.8640 <b>3288</b> :0.8450 <b>7041</b> :0.8300 <b>2747</b> :0.8140
FAST [OA]	151644	<b>3154</b> :34.650 <b>825</b> :0.71700 <b>5377</b> :0.6940 <b>6140</b> :0.6830 <b>3724</b> :0.6450
	151645	<b>3724</b> :0.8630 <b>3955</b> :0.8240 <b>1371</b> :0.8030 <b>3931</b> :0.6940 <b>5940</b> :0.6740
	872	<b>11717</b> :0.578 <b>6739</b> :0.5030 <b>8487</b> :0.4990 <b>2010</b> :0.4640 <b>12647</b> :0.442
EAST[16V]	77091	<b>8487</b> :0.5840 <b>6739</b> :0.5340 <b>11717</b> :0.529 <b>11505</b> :0.493 <b>2851</b> :0.4760
[ASI [10A]	151644	<b>3384</b> :28.324 <b>4241</b> :0.4720 <b>9335</b> :0.4250 <b>11285</b> :0.416 <b>298</b> :0.38400
	151645	<b>4241</b> :0.5410 <b>5731</b> :0.4450 <b>6167</b> :0.4440 <b>7780</b> :0.3940 <b>5314</b> :0.3770

Table 18: Top 5 Average Activation Values for Special Tokens in Qwen2.5-0.5B-insturct with JumpReLU SAE

#### F.4 Steering Output of Three Questions

## F.4.1 Q2

Question2: FINAL EXAM Question 1. What was Elsie Marley profession?	Quen2.5-7B-Instruct-L18-8X JumpReLU SAE Feature ID: 13794 (act most on "< im_stort >" and "user")
Raw Output ( $\alpha = 0$ )	
I'm sorry, but there seems to be some confusion. "Elsie Marley" is not a widel	ly known public figure or character, and without more context, it's
difficult to determine who this person is or their profession. Could you please pr	rovide more information or context about Elsie Marley? This would
help me give you an accurate answer.	
Steered Output ( $\alpha = 25$ )	
I'm sorry, but there seems to be some confusion. "Elsie Marley" is not a widel	ly known public figure or character, and without more context, it's
difficult to determine who this person is or their profession. Could you please pr	rovide more information or context about Elsie Marley? This would
help me give you an accurate answer.	
<b>Steered Output</b> ( $\alpha = 100$ ) To answer the question accurately, I'll need more context about which "Elsie Mat his name. However, based on the name and the common use of names in literatu Harry Potter series. In the Harry Potter series, Elsie Marley is the grandmother of Sirius Black's son, Order of the Phoenix" but does not have a significant role in the story. If you're referring to a different Elsie Marley, please provide more details so I ca	rley" is being referred to, as there could be multiple individuals with re, I'll assume you're referring to Elsie Marley from J.K. Rowling's Nymphadora Tonks. She is mentioned in "Harry Potter and the n give you the correct information.< im_end >
Steered Output ( $\alpha = 150$ )	
To answer the question accurately, I'll need more context about which "Elsie Mar this name. However, based on the name and the commonality of the surname "M the Harry Potter series, who is the mother of Sirius Black. In the books, Elsie Mar mother. If this is the case, her profession is not explicitly stated in the books, but profession typical for her time and social class, such as a teacher, nurse, or office more details.< im_end >	rley" is being referred to, as there could be multiple individuals with larley," one possible reference is to the character Elsie Marley from arley is described as a Muggle woman who was Sirius Black's given her status as a Muggle, it's likely she was involved in a worker. If you have a different context in mind, please provide

Figure 7: The steering output generated by Qwen2.5-7B-Instruct with Feature ID: 13794, focusing on user and <|im\_start|> tokens for the Question 2 (entity description).



Figure 8: The steering output generated by Llama3.1-8B-Instruct with Feature ID: 22642, focusing on user and <|im\_start|> tokens for the Question 2 (entity description).

For Question 2, the Qwen model (Figure 7) shows noticeably improved output quality when feature 13794 is moderately amplified (with  $\alpha$  in the range of 25 to 75). Within this range, the responses become more polite, detailed, and engaging, showing a clear enhancement in interaction quality. However, when the amplification coefficient exceeds this sweet spot (e.g.,  $\alpha \ge 100$ ), the model begins to fabricate information and eventually devolves into repetitive or nonsensical output, resulting in a rapid decline in quality.

In comparison, the Llama model (Figure 8) only benefits from a much narrower range of amplification (approximately  $\alpha = 15$  to 25). Within this window, its responses become slightly more polite and helpful, but still lack substantive factual content. Beyond this narrow range, the output quickly becomes repetitive and loses coherence. Overall, Qwen is able to improve output quality over a broader range of amplification coefficients, while Llama's effective range is much more limited.

F.4.2 Q3

952

	(act most on "< start_header_id >" and "user")
Raw Output $(\alpha = 0)$	
Bien sûr je neux vous aider à rédiger une lettre de motivation. Cen	endant nour que je nuisse vous fournir un modèle précis, i'aurais besoin de
melanes informations supplémentaires :	endant, pour que je puisse vous fournir un modere precis, j'aurais desonir de
1. Qual act la marta neur lanual vaua martular 2	
Quel est le poste pour requer vous postulez ?     Quel est l'antennica à laquelle vous postulez ?	
2. Quelle est renteprise à laquelle vous postulez ?	0
3. Quelles sont vos principales competences et experiences pertinentes	· ?
4. Quels sont vos objectifs professionnels à court et long terme ?	
Avec ces informations, je pourrai vous proposer un modèle de lettre de	e motivation personnalisé et adaptée à votre situation.
Steered Output ( $\alpha = 15$ )	
Voici un exemple de lettre de motivation pour un poste d'ingénieur en in	nformatique :
[Votre nom et adresse]	
[Date]	
[Nom et adresse de l'entreprise]	
Objet : Candidature au poste d'ingénieur en informatique	
Madame, Monsieur,	
Je me permets de vous adresser ma candidature pour le poste d'ingénieu	ur en informatique que vous avez bien voulu publier sur votre site internet.
Titulaire d'un master en informatique de gestion, je suis actuellement	ingénieur en développement logiciel chez [nom de votre entreprise actuelle].
Au cours de ma carrière, j'ai eu l'opportunité de travailler sur des pro-	jets variés, notamment en matière de développement web, d'architecture de
systèmes et de gestion de projets.	
Mes principales compétences sont les suivantes :	
- Maîtrise des langages de programmation : Java, Python, C#	
- Connaissance approfondie de l'architecture de systèmes et des technol	logies cloud (AWS, Azure)
- Expérience en gestion de projets et en gestion de la qualité	
- Bonne capacité d'analyse et de résolution de problèmes	
- Excellentes compétences en communication et en travail d'équipe	
Au cours de mes précédentes expériences professionnelles, j'ai pu déve	elopper une approche méthodique et rigoureuse dans la conception et la mise
en œuvre de solutions informatiques. Je suis particulièrement attiré par	votre entreprise pour ses valeurs d'innovation et de performance.
Je suis convaincu que mes compétences et mon expérience me perme	ttront de contribuer efficacement au développement de vos proiets. Je serais
ravi de pouvoir discuter de cette opportunité avec vous et de vous prése	enter plus en détail mes motivations et mes projets
Je vous remercie par avance pour l'attention que vous porterez	à ma candidature et reste à votre disposition pour tout renseignement
complémentaire	
Dans l'attente de votre rénonse, je vous prie d'agréer. Madame, Monsie	ur l'expression de mes salutations distinguées
[Votes signature]	a, respression de mes sandations distingaces.
[Vote signature]	
N'hécitez pas à adapter ce modèle en fonction de votre evnérience et du	posta pour lagual yous postular < lim and >
	toor toor contract too too too
Steered Output ( $\alpha = 25$ )	
Voici un exemple de lettre de motivation pour un poste d'assistant :	
[Your Name]	
[Your Address]	
[Your Address] [City, State, ZIP Code]	
[Your Address] [City, State, ZIP Code] [Email Address]	
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number]	
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date]	
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Receipient's Name]	
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name]	
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address]	
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City State ZIP Code]	
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Breinient's Name]	
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], Dear within the average mu interact in the Assistant pacifier of Common Dear State, State	nye Nama) as advartised an fushara you found the ink movimed. With a second
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa	ny's Name], as advertised on [where you found the job posting]. With a stron
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience]. I am excited abb	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations.
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa	ny's Name], as advertised on [where you found the job posting]. With a stron, out the opportunity to contribute to your team and support your operations. ny's Name], I gained valuable experience in [mention key responsibilities and
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Address] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh	ny's Name], as advertised on [where you found the job posting]. With a stron out the opportunity to contribute to your team and support your operations. any's Name], I gained valuable experience in [mention key responsibilities an ich allowed me to develop skills in [mention relevant skills]. My ability to
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly ben	ny's Name], as advertised on [where you found the job posting]. With a stron, out the opportunity to contribute to your team and support your operations. any's Name], I gained valuable experience in [mention key responsibilities an iich allowed me to develop skills in [mention relevant skills]. My ability te eficial in managing [mention a specific task or project].
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Address] [Coty, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous loe as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly ben I am particularly drawn to [Company's Name] because of [mention so	ny's Name], as advertised on [where you found the job posting]. With a stron, out the opportunity to contribute to your team and support your operations. any's Name], I gained valuable experience in [mention key responsibilities an ich allowed me to develop skills in [mention relevant skills]. My ability tr eficial in managing [mention a specific task or project]. mething specific about the company or its mission that appeals to you]. I an
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly been I am particularly drawn to [Company's Name] because of [mention sc impressed by [mention a specific aspect of the company or its work], a	ny's Name], as advertised on [where you found the job posting]. With a stron, out the opportunity to contribute to your team and support your operations. nny's Name], I gained valuable experience in [mention key responsibilities an iich allowed me to develop skills in [mention relevant skills]. My ability to eficial in managing [mention a specific task or project]. mething specific about the company or its mission that appeals to you]. I an and I believe my skills and experience align well with your needs.
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly bem I am particularly drawn to [Company's Name] because of [mention s impressed by [mention a specific aspect of the company or its work], a In my current role, I have honed my [mention relevant skills], which I	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations. ny's Name], I gained valuable experience in [mention key responsibilities an iich allowed me to develop skills in [mention relevant skills]. My ability te eficial in managing [mention a specific task or project]. mething specific about the company or its mission that appeals to you]. I an nd I believe my skills and experience align well with your needs. I believe will be valuable in supporting your team. I am proficient in [mention
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Address] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compt achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly bean I am particularly drawn to [Company's Name] because of [mention es impressed by [mention a specific aspect of the company or its work], at In my current role, I have honed my [mention relevant skills], which I relevant software or tools], and I am eager to bring my attention to det	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations. ny's Name], I gained valuable experience in [mention key responsibilities an ich allowed me to develop skills in [mention relevant skills]. My ability te eficial in managing [mention a specific task or project]. smething specific about the company or its mission that appeals to you]. I an Id I believe my skills and experience align well with your needs. believe will be valuable in supporting your team. I am proficient in [mention ail and organizational skills to your team.
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compan background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly ben I am particularly drawn to [Company's Name] because of [mention s impressed by [mention a specific aspect of the company or its work], a In my current role, I have honed my [mention relevant skills], which I relevant software or tools], and I am eager to bring my attention to det I am excited about the prospect of contributing to [Company's Name]	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations. ny's Name], I gained valuable experience in [mention key responsibilities and ich allowed me to develop skills in [mention relevant skills]. My ability to eficial in managing [mention a specific task or project]. omething specific about the company or its mission that appeals to you]. I an ind I believe my skills and experience align well with your needs. ! believe will be valuable in supporting your team. I am proficient in [mentior al and organizational skills to your team. e] and am confident that twy background, skills, and enthusiasm make me a
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly ben I am particularly drawn to [Company's Name] because of [mention s impressed by [mention a specific aspect of the company or its work], a In my current role, I have honed my [mention relevant skills], which I relevant software or tools], and I am eager to bring my attention to det I am excited about the prospect of contributing to [Company's Name storng candidate for the Assistant position. I would welcome the opport	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations, ny's Name], I gained valuable experience in [mention key responsibilities and ich allowed me to develop skills in [mention relevant skills]. My ability to eficial in managing [mention a specific task or project]. omething specific about the company or its mission that appeals to you]. I an nd I believe my skills and experience align well with your needs. I believe will be valuable in supporting your team. I am proficient in [mention ail and organizational skills to your team. e] and am confident that my background, skills, and enthusiasm make me it tunity to discuss how I can contribute to your team further.
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly ben I am particularly drawn to [Company's Name] because of [mention sc impressed by [mention a specific aspect of the company or its work], a In my current role, I have honed my [mention relevant skills], which I relevant software or tools], and I am eager to bring my attention to det I am excited about the prospect of contributing to [Company's Name strong candidate for the Assistant position. I would welcome the opport hank you for considering my application. I look forward to the nessil	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations. uny's Name], I gained valuable experience in [mention key responsibilities and ich allowed me to develop skills in [mention relevant skills]. My ability to eficial in managing [mention a specific task or project]. mething specific about the company or its mission that appeals to you]. I an ind I believe my skills and experience align well with your needs. I believe will be valuable in supporting your team. I am proficient in [mention ail and organizational skills to your team. e] and am confident that my background, skills, and enthusiasm make me is tunity to discuss how I can contribute to your team further.
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compt achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly ben I am particularly drawn to [Company's Name] because of [mention a impressed by [mention a specific aspect of the company or its work], a In my current role, I have honed my [mention relevant skills], which I relevant software or tools], and I am eager to bring my attention to det I am excited about the prospect of contributing to [Company's Nam strong candidate for the Assistant position. I would welcome the oppor Thank you for considering my application. I look forward to the possit Sincerelv.	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations. ny's Name], I gained valuable experience in [mention key responsibilities and ich allowed me to develop skills in [mention relevant skills]. My ability to eficial in managing [mention a specific task or project]. mething specific about the company or its mission that appeals to you]. I an nd I believe my skills and experience align well with your needs. I believe will be valuable in supporting your team. I am proficient in [mention ail and organizational skills to your team. e] and am confident that my background, skills, and enthusiasm make me a tunity to discuss how I can contribute to your team further. pility of discussing this opportunity with you.
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Name] [Company's Address] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compa background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly ben I am particularly drawn to [Company's Name] because of [mention sc impressed by [mention a specific aspect of the company or its work], a In my current role, I have honed my [mention relevant skills], which I relevant software or tools], and I am eager to bring my attention to det I am excited about the prospect of contributing to [Company's Nams strong candidate for the Assistant position. I would welcome the oppor Thank you for considering my application. I look forward to the possib Sincerely, [Your Name]	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations. ny's Name], I gained valuable experience in [mention key responsibilities and ich allowed me to develop skills in [mention relevant skills]. My ability to eficial in managing [mention a specific task or project]. mething specific about the company or its mission that appeals to you]. I an and I believe my skills and experience align well with your needs. I believe will be valuable in supporting your team. I am proficient in [mention ail and organizational skills to your team. e] and am confident that my background, skills, and enthusiasm make me a tunity to discuss how I can contribute to your team further.
[Your Address] [City, State, ZIP Code] [Email Address] [Phone Number] [Date] [Recipient's Name] [Company's Address] [Company's Address] [City, State, ZIP Code] Dear [Recipient's Name], I am writing to express my interest in the Assistant position at [Compan background in [mention relevant skills or experience], I am excited ab In my previous role as an Administrative Assistant at [Previous Compa achievements]. I was responsible for [specific tasks or projects], wh [mention a specific skill or accomplishment] has been particularly ben I am particularly drawn to [Company's Name] because of [mention s impressed by [mention a specific aspect of the company or its work], a In my current role, I have honed my [mention relevant skills], which I relevant software or tools], and I am eager to bring my attention to det I am excited about the prospect of contributing to [Company's Nam strong candidate for the Assistant position. I would welcome the oppor Thank you for considering my application. I look forward to the possit Sincerely, [Your Name]	ny's Name], as advertised on [where you found the job posting]. With a strong out the opportunity to contribute to your team and support your operations. any's Name], I gained valuable experience in [mention key responsibilities and ich allowed me to develop skills in [mention relevant skills]. My ability to eficial in managing [mention a specific task or project]. mething specific about the company or its mission that appeals to you]. I an nd I believe my skills and experience align well with your needs. i believe will be valuable in supporting your team. I am proficient in [mention ail and organizational skills to your team. e] and an confident that my background, skills, and enthusiasm make me a tunity to discussing this opportunity with you.

Figure 9: The steering output generated by Qwen2.5-7B-Instruct with Feature ID: 13794, focusing on user and </ir>im\_start> tokens for the Question 3 (cover letter task).





For Question 3, the Qwen model (Figure 9) shows that moderate amplification of feature 13794 (with  $\alpha$  between 50 and 100) leads to more informative and structured responses, providing richer content and clearer reasoning. This indicates a substantial improvement in output quality within this coefficient range. However, further increasing the amplification causes the model to hallucinate, such as switching languages or generating irrelevant content, and ultimately results in repetitive or meaningless output.

The Llama model (Figure 10) also exhibits some improvement in informativeness and engagement when its most active feature is lightly amplified, but this effect is only present at very low coefficients (up to about  $\alpha = 25$ ). Beyond this point, the output rapidly deteriorates into repetitive or off-topic text. Compared to Qwen, Llama's window for beneficial amplification is much narrower and less robust.

954
955
956
957
958
959
960

961

#### F.4.3 Q4



Figure 11: The steering output generated by Qwen2.5-7B-Instruct with Feature ID: 13794, focusing on user and </im\_start/> tokens for the Question 4 (entity discrimination task).



Figure 12: The steering output generated by Llama3.1-8B-Instruct with Feature ID: 22642, focusing on <|start\_header\_id|> tokens for the Question 4 (entity discrimination task).

In Question 4, both models show that feature amplification can enhance Chain-of-Thought (CoT) (Wei et al., 2022) reasoning and answer quality, but only within specific coefficient ranges. For Qwen (Figure 11), amplifying the most active feature with  $\alpha$  between 25 and 100 produces more convincing, informative, and well-structured responses. This improvement is especially evident in the quality of reasoning and the clarity of the final answers. However, excessive amplification again leads to a loss of coherence and informativeness.

For Llama (Figure 12), a similar pattern is observed but within an even narrower range. Mild amplification (up to  $\alpha = 25$ ) can slightly improve the quality of reasoning and engagement, but any further increase quickly causes the output to become repetitive and less meaningful. This highlights that while both models benefit from feature amplification, Qwen maintains improved output quality over a wider range of coefficients, whereas Llama's useful range is much more restricted.

# G Model Training Log

Due to space constraints, we select training logs from a subset of SAEs for presentation. The complete training logs for all SAEs will be released publicly.

### G.1 Llama-3.1-8B-Instruct

## G.1.1 L18-8X-Standard

*BT*(*P*): Block Training (Pretraining dataset)







FAST: Finetuning-aligned Sequential Training



980

981

974 975

976

977

978

# G.1.2 L18-8X-JumpReLU

*BT(P)*: Block Training (Pretraining dataset)



#### *BT*(*F*): Block Training (Finetuning dataset)



#### 985 FAST: Finetuning-aligned Sequential Training



982 983

#### G.2 Llama-3.2-1B-Instruct

# G.2.1 L9-8X-Standard

*BT*(*P*): Block Training (Pretraining dataset)



# *BT*(*F*): Block Training (Finetuning dataset)







989

990

986 987

# G.2.2 L9-8X-JumpReLU

*BT(P)*: Block Training (Pretraining dataset)



#### *BT*(*F*): Block Training (Finetuning dataset)



#### FAST: Finetuning-aligned Sequential Training



#### G.3 Llama-3.2-3B-Instruct

## G.3.1 L12-8X-Standard

*BT*(*P*): Block Training (Pretraining dataset)



# *BT*(*F*): Block Training (Finetuning dataset)



FAST: Finetuning-aligned Sequential Training



998

999

# G.3.2 L12-8X-JumpReLU

1000 1001

1002

# *BT*(*P*): Block Training (Pretraining dataset)



## *BT*(*F*): Block Training (Finetuning dataset)



#### 1003 FAST: Finetuning-aligned Sequential Training



# G.4 Qwen-2.5-7B-Instruct

# G.4.1 L18-8X-Standard

BT(P): Block Training (Pretraining dataset)



#### *BT*(*F*): Block Training (Finetuning dataset)



# FAST: Finetuning-aligned Sequential Training



1007

1008

# G.4.2 L18-8X-JumpReLU

#### 1009 1010

# *BT*(*P*): Block Training (Pretraining dataset)



# a*BT*(*F*): Block Training (Finetuning dataset)



# 1012 FAST: Finetuning-aligned Sequential Training



# G.5 Qwen-2.5-3B-Instruct

# G.5.1 L18-8X-Standard

*BT*(*P*): Block Training (Pretraining dataset)



# *BT*(*F*): Block Training (Finetuning dataset)



FAST: Finetuning-aligned Sequential Training



1016

1017

1013

#### G.5.2 L18-8X-JumpReLU 1018

#### *BT*(*P*): Block Training (Pretraining dataset) 1019



# *BT*(*F*): Block Training (Finetuning dataset)



#### FAST: Finetuning-aligned Sequential Training 1021



## G.6 Qwen-2.5-1.5B-Instruct

## G.6.1 L14-8X-Standard

*BT*(*P*): Block Training (Pretraining dataset)



# *BT*(*F*): Block Training (Finetuning dataset)



# FAST: Finetuning-aligned Sequential Training



1025

1026

1022

# G.6.2 L14-8X-JumpReLU

1027 1028

1029

# *BT*(*P*): Block Training (Pretraining dataset)



# *BT*(*F*): Block Training (Finetuning dataset)



#### 1030 FAST: Finetuning-aligned Sequential Training



### G.7 Qwen-2.5-0.5B-Instruct

## G.7.1 L12-8X-Standard

*BT*(*P*): Block Training (Pretraining dataset)



*BT*(*F*): Block Training (Finetuning dataset)



FAST: Finetuning-aligned Sequential Training

kesses/It_kess		losses	mor_loos	lasses/averall_less		
a Litteraturi atan di budut panta di 10 no 10 no 10	1007 607 505 10	Second states and a state of the states of t	ngganatinatinatinatinatinatinatinatinatinat	1000 - 300	Willing media biost on the orthogonal description of the starts	
(a) L1 Loss		(b) !	ASE Loss		(c) Overall Loss	
		100 100 100 200				
(d) L0 Sparsity	(e) Dead	Features	(f) Explained	I Variance	(g) CE Loss Score	

1034

1035

# G.7.2 L12-8X-JumpReLU

BT(P): Block Training (Pretraining dataset)



#### *BT*(*F*): Block Training (Finetuning dataset)



#### 1039 FAST: Finetuning-aligned Sequential Training



1036 1037