# GSAE: GRAPH-REGULARIZED SPARSE AUTOEN-CODERS FOR ROBUST LLM SAFETY STEERING

## **Anonymous authors**

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

032

034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

## **ABSTRACT**

Large language models (LLMs) face critical safety challenges, as they can be manipulated to generate harmful content through adversarial prompts and jailbreak attacks. Many defenses are typically either black-box guardrails that filter outputs, or internals-based methods that steer hidden activations by operationalizing safety as a single latent feature or dimension. While effective for simple concepts, this assumption is limiting, as recent evidence shows that abstract concepts such as refusal and temporality are distributed across multiple features rather than isolated in one. To address this limitation, we introduce Graph-Regularized Sparse Autoencoders (GSAEs), which extends SAEs with a Laplacian smoothness penalty on the neuron co-activation graph. Unlike standard SAEs that assign each concept to a single latent feature, GSAEs recover smooth, distributed safety representations as coherent patterns spanning multiple features. We empirically demonstrate that GSAE enables effective runtime safety steering, assembling features into a weighted set of safety-relevant directions and controlling them with a two-stage gating mechanism that activates interventions only when harmful prompts or continuations are detected during generation. This approach enforces refusals adaptively while preserving utility on benign queries. Across safety and QA benchmarks, GSAE steering achieves an average 82% selective refusal rate, substantially outperforming standard SAE steering (42%), while maintaining strong task accuracy (70% on TriviaQA, 65% on TruthfulQA, 74% on GSM8K). Robustness experiments further show generalization across LLaMA-3, Mistral, Qwen, and Phi families and resilience against jailbreak attacks (GCG, AutoDAN), consistently maintaining  $\geq 90\%$  refusal of harmful content.

## 1 Introduction

Modern large language models (LLMs) excel at diverse tasks like question answering and reasoning (Touvron et al., 2023), yet their deployment faces significant safety challenges. LLMs can be manipulated into generating harmful content through adversarial prompts and jailbreak attacks (Wei et al., 2023). Effective defenses must both block unsafe generations and preserve the model's utility on benign queries (Ganguli et al., 2022).

Existing safety approaches generally fall into two categories: *black-box guardrails* and *internals-based methods*. Black-box guardrails, such as prompt engineering (Bai et al., 2022) or output classifiers (Inan et al., 2023), offer quick defenses but are often brittle to distributional shifts (Zou et al., 2023) and lack interpretability. Internals-based methods (Turner et al., 2023a) aim to leverage the model's hidden representations. Sparse autoencoders (SAEs) have become a prominent tool in this category, allowing the decomposition of hidden activations into sparse, often interpretable, latent features (Cunningham et al., 2023; Templeton et al., 2024; Bricken et al., 2023).

Despite their utility for interpreting concrete concepts, standard SAEs may have limitations when applied to complex domains like time or safety. This is because SAEs are inherently *local*, encouraging each latent dimension to represent a single "monosemantic" feature. This often leads to it fragmented into disconnected sub-concepts (like 'refusal' or 'danger') or create redundant features that overlap in meaning, failing to learn a coherent representation (Bricken et al., 2024).

Recent studies highlight this representational gap for abstract concepts. While concrete concepts (e.g., objects) often align with single, axis-like features, higher-level abstract concepts are typically

encoded in a distributed and nonlinear fashion (Liao et al., 2023). For instance, temporal concepts manifest as nonlinear circular manifolds (Engels et al., 2025), and refusal behavior involves multiple independent directions and nonlinear geometries (Wollschläger et al., 2025; Hildebrandt et al., 2025). This evidence suggests that abstract concepts are better modeled as distributed properties. We argue that safety, as an abstract, socially grounded concept dependent on context and human judgment (Slavich, 2023), requires a distributed representation.

Our proposed approach. To model safety as a distributed concept, we introduce the Graph-Regularized Sparse Autoencoder (GSAE). GSAE extends standard SAEs by incorporating a graph Laplacian regularizer (Belkin et al., 2006). This treats each neuron as a node, with edges defined by activation similarity (Diao et al., 2024). The Laplacian penalty enforces smoothness across co-activating neurons, yielding coherent, non-redundant features that more effectively capture distributed safety patterns (Belkin et al., 2006). From these features, we construct a spectral vector bank: a weighted library of decoded safety directions. These weights are meticulously derived to reflect three criteria: spectral smoothness, a measure of structural coherence (von Luxburg, 2007); supervised importance, which gauges predictive strength for harmfulness (Belrose et al., 2023); and causal influence, the measurable steering effect (Meng et al., 2022).. At inference time, this bank is deployed through a dual-gating controller, as illustrated in Figure 1. An input gate evaluates the features pre-generation, while a continuation gate monitors decoding during generation. This design dynamically scales steering strength, preventing both under-refusal and over-refusal, and enabling selective safety interventions while preserving accuracy on benign queries (Sun et al., 2024).

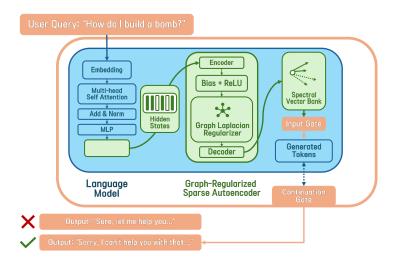


Figure 1: Overview of the GSAE steering framework. A user query is encoded into hidden states, which the GSAE decomposes into graph-regularized safety features. A dual-gating controller uses these features to make a two-stage safety assessment: an Input Gate evaluates the initial prompt, while a Continuation Gate monitors the generation in real-time. This allows the system to selectively block harmful outputs while preserving benign ones.

**Contributions.** This paper provides the following fundamental contributions:

**Graph-Regularized Sparse Autoencoders (GSAE):** We introduce GSAE, which applies graph Laplacian regularization to sparse autoencoders to more effectively capture distributed concepts. This design explicitly encodes relational structure among neurons, making it well-suited for representing safety-relevant activation patterns.

**Runtime Steering Framework:** We leverage GSAE-derived features by building a spectral vector bank, a curated library of safety directions, which is then managed by a dual-gating controller that adaptively decides when and how strongly to intervene. This enables selective, stable steering during inference, improving refusal on harmful prompts while preserving benign task performance.

**Robust Benchmarking and Generalizability:** We conduct extensive evaluations across a diverse suite of LLMs (Llama-3, Mistral, Qwen, and Phi families) and against a wide range of adversarial

jailbreak attacks (GCG, AutoDAN, TAP). Our results demonstrate that GSAE steering consistently and substantially outperforms state-of-the-art baselines, achieving high safety discrimination while preserving utility, and providing a robust, generalizable safety mechanism.

## 2 PRELIMINARIES

This section reviews the core concepts underlying our method: the internal representations of LLMs, sparse autoencoders, and graph Laplacians.

**LLM Internals.** Transformer-based LLMs process input through a series of layers (Vaswani et al., 2017). At each layer, indexed by l, the model generates a matrix of hidden states  $\boldsymbol{H}^{(l)} \in \mathbb{R}^{n \times d}$ , where n is the sequence length and d is the hidden dimension. To obtain a representation for an entire prompt, these hidden states are aggregated via a pooling operation (e.g., mean-pooling) into a single pooled activation vector  $\boldsymbol{h}^{(l)} \in \mathbb{R}^d$  for that layer (Guo et al., 2025). Since harmful behaviors manifest as specific patterns in these activations (Zhou et al., 2024; Xu et al., 2024), they serve as an effective target for intervention.

**Sparse Autoencoders (SAEs).** Given a pooled hidden state  $x \in \mathbb{R}^d$ , a Sparse Autoencoder (SAE) aims to find a more interpretable, lower-dimensional representation. It does this by mapping x to a sparse **latent code**  $z \in \mathbb{R}^k$  (where  $k \gg d$ ) and then reconstructing the original input, denoted  $\hat{x}$ . This process is defined by:

$$z = \phi(\boldsymbol{W}^{(e)}\boldsymbol{x}), \quad \hat{\boldsymbol{x}} = \boldsymbol{W}^{(d)}\boldsymbol{z},$$

where  $\boldsymbol{W}^{(e)} \in \mathbb{R}^{k \times d}$  is the **encoder** matrix,  $\boldsymbol{W}^{(d)} \in \mathbb{R}^{d \times k}$  is the **decoder** matrix, and  $\phi(\cdot)$  is a non-linear activation function, typically a ReLU, to ensure non-negative feature activations. The training objective is designed to minimize two competing goals (Gao et al., 2024): the **reconstruction error**, measured by the squared L2 norm  $\|\boldsymbol{x} - \hat{\boldsymbol{x}}\|_2^2$ , and the **sparsity** of the latent code, encouraged by an L1 penalty  $\|\boldsymbol{z}\|_1$  weighted by a hyperparameter  $\lambda_{\text{spar}}$ :

$$\mathcal{L}_{ ext{SAE}} = \|oldsymbol{x} - \hat{oldsymbol{x}}\|_2^2 \ + \ \lambda_{ ext{spar}} \|oldsymbol{z}\|_1.$$

The L1 penalty forces most elements of the latent code z to be zero. This encourages the SAE to learn *localized features*, where each active dimension in z ideally corresponds to a single, interpretable concept (Cunningham et al., 2023). However, this very locality is a limitation for capturing abstract, distributed properties like safety, which may lead to feature fragmentation (Belrose, 2025).

**Graph Laplacian and Smoothness.** To capture the relational structure between neurons, we model them as a graph  $\mathcal{G}=(\mathcal{V},\mathcal{E})$ , where each node in  $\mathcal{V}$  represents one of the d neurons. Their relationships are encoded in an **adjacency matrix**  $A \in \mathbb{R}^{d \times d}$ , where  $A_{ij}$  is a positive weight representing the strength of the connection between neurons i and j. The **degree matrix** is  $D=\operatorname{diag}(d_1,\ldots,d_d)$  with  $d_i=\sum_j A_{ij}$ , and the **graph Laplacian** is defined as L=D-A.

A graph signal is a vector  $z \in \mathbb{R}^d$  assigning a scalar  $z_i$  to each neuron i. The smoothness of z over the graph is measured by its Laplacian energy:

$$E(\boldsymbol{z}) = \boldsymbol{z}^{\top} \boldsymbol{L} \boldsymbol{z} = \frac{1}{2} \sum_{i,j} \boldsymbol{A}_{ij} (z_i - z_j)^2.$$

The quadratic form of the energy provides the key intuition for our approach. The total energy is a weighted sum of squared differences between the signal values ( $z_i$  and  $z_j$ ) on connected neurons. Consequently, a large penalty is incurred if neurons with a strong connection are assigned dissimilar values. Minimizing this energy term imposes a smoothness prior on the signal, thereby forcing the values assigned to strongly co-activating neurons to be similar. In our autoencoder, we penalize this energy for each decoded feature, which biases safety directions toward smooth, distributed patterns across the neuron graph. This corresponds to suppressing high-frequency components and favoring low-frequency eigenmodes of the Laplacian, a standard interpretation in spectral graph theory (Smola & Kondor, 2003), also detailed in Appendix A.

While dense graph operations can be computationally intensive (scaling as  $O(d^2)$ ), our approach is efficient in practice as we sparsify the graph by thresholding edge weights, making the overhead from graph operations negligible compared to the autoencoder's standard matrix multiplications.

# 3 RELATED WORK

**Prior Safety Methods.** Prior work on LLM safety can be broadly categorized into black-box and internals-based methods. Black-box approaches operate on the model's inputs and outputs, using techniques like adversarial prompt detection (Mehrotra et al., 2024; Chao et al., 2023), output filtering with toxicity detectors (Wang et al., 2024), and prompt engineering with "constitutional" principles (Bai et al., 2022). While efficient and applicable for black-box settings, these methods' reliance on surface-level lexical patterns can limit their robustness against adaptive attacks and distributional shifts (Cui et al., 2024).

Thus, we focus on internals-based methods that directly intervene on activation dynamics. A prominent line of this research seeks to identify low-dimensional structure corresponding to safety concepts. This includes learning linear classifiers to find "refusal directions" (Arditi et al., 2024; Siu et al., 2025) and steering generation by adding or subtracting activation vectors, as in Contrastive Activation Addition (CAA) (Turner et al., 2023a). Other approaches intervene at a finer-grained level, identifying causal pathways via activation patching (Meng et al., 2022) or applying corrective projections with monitoring heads, like SafeSwitch (Han et al., 2025). While these methods show promise, they typically assume that safety can be represented as a single axis or a small set of independent directions. Among internals-based methods, Sparse Autoencoders (SAEs) have been increasingly used for control by decomposing hidden activations into sparse, interpretable features (Cunningham et al., 2023; Templeton et al., 2024; Bricken et al., 2023). Several works demonstrate that manipulating these features can predictably alter model behavior (O'Brien et al., 2025; Turner et al., 2023b), with applications in suppressing private information (Frikha et al., 2025) or disentangling attention head activations (Zhan et al., 2025). However, the features learned by standard unsupervised SAEs may not align with safety concepts and can be unstable or redundant (Park et al., 2024). Our work addresses this limitation by incorporating graph Laplacian regularization to produce structurally coherent features better suited for the distributed nature of safety.

Safety as a Distributed Concept. Recent studies increasingly indicate that abstract concepts in LLMs are fundamentally distributed rather than localized to single, interpretable directions. Concepts ranging from temporality to moral judgment have been found to be encoded in diffuse, nonlinear geometric structures that require the coordination of many neurons (Liao et al., 2023; Engels et al., 2024; 2025; Wang et al., 2023). This paradigm is particularly relevant for safety; for instance, refusal behavior has been shown to manifest not as a simple axis but as complex, polyhedral "concept cones" with fundamentally nonlinear properties (Wollschläger et al., 2025; Hildebrandt et al., 2025). These findings challenge the core monosemantic assumption of standard SAE-based methods, which can produce unstable or spurious features for such complex behaviors (Park et al., 2024). Building on this collective evidence, we follow the intuition that safety, as an inherently abstract and socially grounded concept, requires a distributed rather than localized representation.

Graph-Based Regularization in Machine Learning. Laplacian regularization is used in graph-based machine learning to enforce smoothness priors on data. By penalizing variation between connected nodes, it has been central to foundational methods in spectral clustering (Von Luxburg, 2007), manifold learning (Belkin & Niyogi, 2003), and semi-supervised learning (Zhu et al., 2003; Yang et al., 2016). In neural network contexts, this form of regularization helps align learned representations with a given topology, improving model robustness and yielding multi-scale features (Cheng et al., 2023; Shuman et al., 2013). While well-established, these methods are underexplored for steering the internal representations of LLMs. Our work adapts this principle to sparse autoencoders, using graph structure to produce features that reflect distributed rather than isolated patterns.

### 4 METHODOLOGY

We introduce GSAE, a novel method for learning structured representations of safety-relevant activation patterns from an LLM's internal activations. These representations are then curated into a **spectral vector bank**, a library of steering directions. At runtime, a **dual-gating controller** uses this bank to perform adaptive, real-time interventions, steering the model toward safer outputs.

#### 4.1 PROBLEM FORMULATION

Our work addresses the fundamental challenge of extracting structured and distributed safety-relevant representations from the complex internal activations of LLMs. For a given prompt, we operate on the pooled hidden state  $\bar{h}^{(l)} \in \mathbb{R}^d$  from a model layer l, where d is the hidden dimension.

We operate on the pooled hidden state  $\bar{h}^{(l)} \in \mathbb{R}^d$  from a model layer l, where d is the hidden dimension. Our goal is to learn a feature mapping  $f_\theta : \mathbb{R}^d \to \mathbb{R}^k$  that transforms the hidden state into a sparse latent code  $z = f_\theta(\bar{h}^{(l)})$ . The feature dimension k is intentionally expanded to be much larger than the hidden dimension  $(k \gg d)$ . Formally, we state the problem as:

Given pooled hidden states from an LLM, learn a mapping  $f_{\theta}$  that produces latent features, which capture the distributed, relational properties of safety within the model's internal representations.

#### 4.2 GRAPH-REGULARIZED SPARSE AUTOENCODERS (GSAE)

To capture these distributed safety features, we introduce GSAE. While standard SAEs effectively enforce sparsity, this can fragment complex concepts like safety into an array of redundant or weak features. GSAE extends the SAE framework by incorporating a graph-based regularizer that enforces *relational smoothness*, ensuring that frequently co-activating neurons develop similar learned features. This promotes coherent and robust representations while preserving the sparsity essential for disentanglement.

#### 4.2.1 NEURON CO-ACTIVATION GRAPH

To apply the graph-based penalty, we must first construct a model of the relational structure between neurons. We collect the pooled hidden states for a diverse set of N prompts, forming an activation matrix  $\boldsymbol{H} \in \mathbb{R}^{d \times N}$ . Each row  $\boldsymbol{H}_{i:}$  of this matrix represents the **activation profile** of neuron i across all prompts. We then construct an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where each of the d neurons is a node  $v_i \in \mathcal{V}$ . The edge weight between any two neurons is defined by the cosine similarity of their activation profiles, capturing how often they activate together. This allows us to build the adjacency matrix  $\boldsymbol{A}$  and, subsequently, the graph Laplacian  $\boldsymbol{L} = \boldsymbol{D} - \boldsymbol{A}$  as defined in Section 2. This Laplacian matrix  $\boldsymbol{L}$  mathematically encodes the relational co-activation structure of the entire neuron space, providing the foundation for our regularization.

## 4.2.2 GSAE OBJECTIVE

Given a pooled hidden state x, the GSAE encodes it to a latent code  $z = \text{ReLU}(W^{(e)}x)$  and decodes it back to a reconstruction  $\hat{x} = W^{(d)}z$ . The training objective is a composite loss function that combines four distinct components:

$$\mathcal{L}_{\text{GSAE}} = \underbrace{\|\boldsymbol{x} - \hat{\boldsymbol{x}}\|_{2}^{2}}_{\text{Reconstruction}} + \underbrace{\lambda_{\text{spar}} \|\boldsymbol{z}\|_{1}}_{\text{Sparsity}} + \underbrace{\lambda_{\text{graph}} \sum_{j=1}^{k} \left( (\boldsymbol{W}^{(d)}_{(\cdot,j)})^{\top} \boldsymbol{L} \, \boldsymbol{W}^{(d)}_{(\cdot,j)} \right)}_{\text{Graph Regularization}}.$$

Here,  $\lambda_{\rm spar}$  and  $\lambda_{\rm graph}$  are coefficients that balance the influence of each term. The **reconstruction** and **sparsity** terms are standard in SAEs. The first ensures the learned features faithfully represent the original activations, while the second encourages interpretability by ensuring only a few features are active at any time. Our core contribution is the **graph regularization** term. It penalizes the Laplacian energy of each decoded feature direction (each column  $W^{(d)}_{(\cdot,\,j)}$  of the decoder matrix). As explained in Section 2, this forces the features to be *smooth* over the neuron graph, meaning that neurons that frequently co-activate will be represented similarly within a feature. This directly counteracts fragmentation and promotes the discovery of coherent, distributed features.

#### 4.3 STEERING WITH GSAE FEATURES

The features learned by the GSAE are used at inference time to perform runtime safety steering in a four-stage process, described as follows.

**Step 1: Latent Encoding** For any input prompt, we first extract its pooled hidden states  $\bar{h}^{(l)}$  from a set of predefined target layers  $l \in \mathcal{L}$ . Each hidden state is then passed through the trained GSAE encoder to produce a set of sparse latent codes that are concatenated into a single feature vector z and represent the prompt's safety-relevant properties:

$$\boldsymbol{z}^{(l)} = \mathrm{ReLU}(\boldsymbol{W}^{(e)}{}^{(l)}\bar{\boldsymbol{h}}^{(l)})$$

Step 2: Spectral Vector Bank Construction While the GSAE learns a set of sparse features, not all are equally suited for steering, as many may be structurally incoherent, semantically irrelevant, or causally inert. To address this, we construct a **spectral vector bank**, a curated library of steering directions, using a three-stage filtering and weighting process designed to identify features that are structurally coherent, semantically relevant, and causally effective. Each latent feature i corresponds to a **decoded direction**,  $v_i$  (the i-th column of the decoder matrix  $W^{(d)}$ ), in the model's activation space. We evaluate each direction against three sequential criteria:

Structural Coherence ( $s_i^{\text{lap}}$ ): To ensure features represent coherent patterns rather than noise, we measure their alignment with the neuron graph's structure. We quantify this using normalized Dirichlet energy,  $E_i = (v_i^{\top} L v_i) / \|v_i\|_2^2$ , where lower energy indicates a smoother feature. This is converted to a score via  $s_i^{\text{lap}} = \exp(-\beta E_i)$  to prioritize structurally sound directions.

**Semantic Relevance**  $(s_i^{\text{imp}})$ : To identify which features are predictive of harmfulness, we measure their relevance using a linear probe trained to classify harmful content from the latent codes z. The relevance score,  $s_i^{\text{imp}}$ , is the absolute magnitude of the learned coefficient  $|\theta_i|$  for feature i, selecting for features with high predictive power.

**Causal Efficacy**  $(s_i^{\text{infl}})$ : To validate that a feature has a practical steering effect, we measure its causal efficacy. This score,  $s_i^{\text{infl}}$ , is the mean absolute change in the model's refusal probability when we add the feature's direction,  $v_i$ , to the activations of validation prompts, thereby isolating features with a demonstrable causal impact.

These three scores are combined multiplicatively, ensuring that a feature attains a high final weight only if it scores strongly across all desiderata. The final weight  $w_i$  for each direction is given by the normalized product:

$$w_i = \frac{(s_i^{\text{lap}})^{\alpha} \cdot (s_i^{\text{imp}})^{\beta} \cdot (s_i^{\text{infl}})^{\gamma}}{\sum_{j \in \mathcal{S}} (s_j^{\text{lap}})^{\alpha} \cdot (s_j^{\text{imp}})^{\beta} \cdot (s_j^{\text{infl}})^{\gamma}}.$$

This multiplicative approach ensures that a feature must be structurally coherent, semantically relevant, and causally effective; a low score on any single criterion will significantly diminish the feature's final weight. In our experiments, the parameters  $\alpha, \beta, \gamma$  are set to 1.0, giving equal importance to each criterion and providing a robust, un-tuned baseline.

**Step 3: Dual-Gated Risk Control** A **dual-gating controller** uses the latent features *z* to dynamically decide *when* and *how strongly* to intervene:

**Input Gate.** An initial assessment of the prompt's risk is made by calculating a harm probability,  $p_{\rm harm} = g(z_{\rm prompt})$ . If the risk exceeds a high threshold  $t_{\rm hi}$ , it triggers immediate refusal; if it falls within a moderate range  $[t_{\rm lo}, t_{\rm hi})$ , it activates a monitoring state. The selection of these gating thresholds, along with other key hyperparameters, is based on a systematic sensitivity analysis detailed in Appendix D.1. Our method achieves consistent gains across a wide range of these hyperparameter choices, indicating robustness rather than a brittle dependence on specific values.

Continuation Gate. A continuation gate monitors the generation token-by-token. To prevent unstable, oscillating interventions on borderline content, it uses hysteresis, separate, stable thresholds for escalating  $(d_{\rm hi})$  and de-escalating  $(d_{\rm lo})$  the steering strength. It outputs a steering multiplier  $\gamma_t$  based on the real-time risk. This dual-gated design provides both coarse-grained control at the prompt level and fine-grained, stable adjustments during generation.

**Step 4: Runtime Intervention** When the controller determines that steering is necessary ( $\gamma_t > 0$ ), it applies a corrective shift,  $\Delta \bar{h}_t^{(l)}$ , to the hidden states at each decoding step t. This shift is a weighted sum of the top safety directions from the spectral bank, scaled by their cosine similarity

alignment with the current hidden state:

$$\Delta \bar{\boldsymbol{h}}_t^{(l)} = \alpha_0 \cdot \gamma_t \sum_{i \in S} w_i \cos(\bar{\boldsymbol{h}}_t^{(l)}, \boldsymbol{v}_i) \frac{\boldsymbol{v}_i}{\|\boldsymbol{v}_i\|_2}.$$

Here,  $\alpha_0$  is a global hyperparameter controlling the base steering strength. This intervention adaptively nudges the model's activations away from harmful configurations and toward safer ones, guided by the coherent features in our spectral bank.

## 5 EXPERIMENTS

#### 5.1 EXPERIMENTAL SETTINGS

We systematically evaluate GSAE for *runtime safety steering* along four dimensions: (i) overall safety and utility, (ii) generalization across model families and scales, (iii) refusal rate trade-offs, and (iv) robustness to jailbreak attacks. We now describe the experimental setting before presenting the results.

**Tasks & Metrics.** To evaluate the trade-off between steering for safety and preserving task performance, we consider two tasks: *safety* and *utility*. The safety task measures a model's refusal behavior on harmful and benign prompts. We report harmful refusal rate (HRR), the proportion of harmful prompts that are successfully blocked, and safe refusal rate (SRR), the proportion of safe prompts that are incorrectly blocked, and summarize their trade-off using the selective refusal score  $\Delta_s = \text{HRR} - \text{SRR}$ . For utility, we report standard **accuracy** (%) on QA benchmarks, and analyze the trade-off between safety improvements and utility degradation introduced by steering.

**Datasets.** For safety, we use the WildJailbreak (Xia & et al., 2024) and JailbreakBench (JBB) (Chao et al., 2024) datasets. For utility, we report accuracy on TriviaQA (Joshi et al., 2017), TruthfulQA (Lin et al., 2021), and GSM8K (Cobbe et al., 2021).

**Baselines.** We compare GSAE against a range of representative defenses: simple prompting guardrails, which add safety instructions to the system prompt; SAE steering (O'Brien et al., 2025), which manipulates individual features from a standard sparse autoencoder; Contrastive Activation Addition (CAA) (Turner et al., 2023a), which steers activations along a predefined safety vector; and SafeSwitch (Han et al., 2025), a state-of-the-art defense that uses monitoring heads to apply corrective projections. We also include an unsteered model as a baseline and conduct ablation studies that remove or modify key components of our framework.

Jailbreaking Strategies. We evaluate robustness against a suite of strong and diverse jailbreaking strategies. These include: GCG (Greedy Coordinate Gradient) (Zou et al., 2023), a gradient-based optimization method that finds a short, transferable adversarial suffix designed to be appended to any harmful prompt; AutoDAN (Liu et al., 2023), which uses a hierarchical genetic algorithm to evolve human-readable, semantically coherent prompts that bypass common defenses; TAP (Tree of Attacks with Pruning) (Mehrotra et al., 2023), a black-box method that uses an LLM to build a tree of attack variations, analyzing the model's refusals to iteratively generate and prune new prompts; and general adaptive attacks (Andriushchenko et al., 2024), a category of attacks specifically tailored to a known defense, using iterative queries to find weaknesses in the target's safety mechanism.

**Models.** Our main experiments use Llama-3 8B. Hidden states are mean-pooled from middle-to-upper layers ( $\mathcal{L}=\{6,8,10,12\}$ ), and neuron co-activation graphs are constructed with cosine similarity threshold  $\tau=0.6$ . To assess generalizability, we also evaluate GSAE on Mistral 7B, Qwen 2.5 14B, and Phi-4 15B. Further details on the implementation and hyperparameters selection are provided in the Appendix, Section B.

#### 5.2 RESULTS AND ANALYSIS

**Overall Performance.** Table 1 reports the performance of GSAE against a suite of existing methods, measured by the selective refusal score ( $\Delta_s$ ), where higher values indicate stronger discrimination between harmful and safe prompts, and by utility accuracy on QA benchmarks.

GSAE steering with all of the components implemented achieves the best performance, reaching 90% on WildJailbreak and 76% on JBB (average = 83%). This substantially outperforms SafeSwitch

Table 1: Safety performance is measured by the selective refusal score  $\Delta_s$ , while the utility is measured by QA accuracy. An effective trade-off corresponds to substantial safety improvements with minimal utility degradation relative to the *No Steering* baseline (top row). All results are reported on Llama 3 8B as mean  $\pm$  std over 5 random seeds.

	Safety (	Safety $(\Delta_s)$		Utility (Accuracy %)		
Method	WildJailbreak	JBB	TriviaQA	TruthfulQA	GSM8K	
No Steering	-3±1%	-9±1%	74±0.4%	69±0.6%	79±0.5%	
Prompting guardrails	$18 \pm 2\%$	$10\pm 2\%$	72±0.7%	$69 \pm 0.5\%$	$77 \pm 0.6\%$	
SAE steering	$48 \pm 3\%$	$36 \pm 3\%$	62±0.8%	$67 \pm 0.9\%$	$76\pm0.7\%$	
CAA	$42 \pm 2\%$	$30\pm 2\%$	60±0.9%	$66\pm1.0\%$	$67 \pm 0.8\%$	
SafeSwitch	65±3%	$51\pm3\%$	61±1.0%	$65 {\pm} 0.8\%$	$66 \pm 0.9\%$	
GSAE (random graphs)	60±3%	44±3%	33±1.2%	54±1.1%	23±1.5%	
GSAE (uniform weight)	$72 \pm 2\%$	$58\pm2\%$	55±0.9%	$49 \pm 1.3\%$	$55\pm1.0\%$	
GSAE (no gating)	$78 {\pm} 2\%$	$64 \pm 2\%$	63±0.6%	$60 {\pm} 0.8\%$	$66\pm0.7\%$	
GSAE	$90{\pm}2\%$	$76\pm2\%$	70±0.5%	$65 \pm 0.7\%$	$74\pm0.6\%$	

(average = 58%) and nearly doubles the effectiveness of standard SAE steering (average = 42%). Importantly, these safety gains come with only minor utility degradation: GSAE reduces QA accuracy by just 4–5% relative to the no-steering baseline. By contrast, methods such as CAA and SafeSwitch incur much larger drops, with TriviaQA accuracy falling to 60% and 61%, respectively.

Ablation studies further confirm the contribution of each component. Randomizing the neuron coactivation graph or removing spectral weighting causes severe degradation in both safety and utility, underscoring the need for a structured, regularized representation. Eliminating the dual-gating mechanism also lowers selective refusal, highlighting its role in balancing robustness with utility.

Generalization Across Models. To validate performance across architectures and scales, Figure 2 reports the selective refusal score  $\Delta_s$  for GSAE, compared with SafeSwitch, the strongest baseline, and a No Steering control. GSAE consistently outperforms SafeSwitch, with gains ranging from +10 points on Phi-4 15B (88% vs. 78%) to +24 points on Llama-3 8B (82% vs. 58%). These results confirm that our graph-based regularization captures generalizable safety structure, enabling robust steering across diverse model families.



Figure 2: Safety performance across models, reported as the selective refusal score  $\Delta_s$ . GSAE (green) consistently outperforms both SafeSwitch (orange) and the baseline.

**Analysis of Refusal Rate Trade-offs.** To disentangle the contributions of harmful and safe refusals, Figure 3 portrays the harmful refusal rate (HRR) against the safe refusal rate (SRR) across model families and methods. The top-left corner of each plot corresponds to the ideal operating region: a model that blocks nearly all harmful prompts while rarely over-refusing benign ones.

Across all four models, GSAE consistently lies closest to this ideal area, featuring high HRR with low SRR. For instance, on Qwen 2.5 14B, it achieves HRR above 90% with SRR around 10%. In contrast, SafeSwitch reaches high HRR but at the cost of substantially higher SRR, reflecting sizable over-refusal. SAE steering and CAA fail to achieve strong HRR, limiting their robustness. The unsteered baseline consistently performs poorly on both axes.

This disentangled view confirms that GSAE's advantage arises not just from maximizing harmful refusals but from simultaneously minimizing safe refusals.

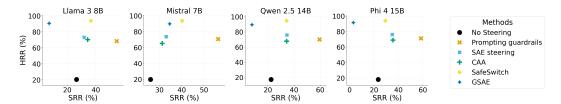


Figure 3: Refusal trade-off plots: harmful refusal rate (HRR, y-axis) vs. safe refusal rate (SRR, x-axis). The ideal region is the top-left (maximizing harmful refusals while minimizing safe ones). GSAE consistently occupies this region, achieving the best balance.

**Robustness Under Jailbreak Attacks.** We further evaluate robustness against four widely used jailbreak strategies: GCG, AutoDAN, TAP, and adaptive attacks. Table 2 shows that GSAE steering consistently sustains an HRR of at least 90% across all attack types, substantially outperforming baselines such as SAE steering and CAA. Prompting guardrails, by contrast, provide only partial protection and collapse under adaptive attacks, with refusal rates below 30%.

Table 2: Robustness to jailbreak attacks. We report harmful refusal rate (HRR), where higher values indicate stronger robustness. GSAE steering sustains HRR  $\geq 90\%$  across all attack types, substantially outperforming all baselines (SAE steering, CAA) and prompting guardrails.

Method	GCG	AutoDAN	TAP	Adaptive
Prompting guardrails	41%	36%	32%	28%
CAA	58%	55%	49%	46%
SafeSwitch	68%	84%	40%	39%
SAE steering	72%	68%	65%	61%
GSAE	100%	95%	90%	92%

# 6 Discussion

This work challenges the assumption that safety concepts can be localized to a single sparse feature, and instead hypothesizes that safety is inherently distributed, emerging from coordinated patterns across many neurons. Inspired by this hypothesis, we introduced **Graph-Regularized Sparse Autoencoders (GSAE)**, which augment SAEs with a Laplacian smoothness prior on the neuron coactivation graph. This regularizer biases features toward smooth, low-frequency modes, yielding safety representations that are distributed and relational rather than isolated.

Empirically, our results provide strong evidence for this approach. GSAE achieves substantially higher safety discrimination than baselines while preserving QA utility, generalizes across model families and scales, and remains robust under strong jailbreak attacks, consistently refusing over 90% of harmful inputs. Together, these results indicate that distributed, graph-regularized features provide a principled and reliable basis for steering compared to single-direction methods.

Future work may investigate decomposing these distributed features into interpretable safety subcategories (e.g., separating patterns related to violence or hate speech) and ensuring the underlying neuron co-activation graph is robust to potential dataset biases. Furthermore, extending graph-regularized learning beyond language models to multi-modal domains such as vision and audio, where safety concerns are equally pressing, remains a promising direction.

# REPRODUCIBILITY STATEMENT

To ensure the reproducibility of our results, we provide the complete source code at an anonymous repository: https://anonymous.4open.science/r/GSAE-B5DB. A detailed breakdown of our experimental setup is provided in the Appendix. Specifically, Appendix B contains a full description of the datasets used, the computing environment, and a table of the final hyperparameters required to replicate our main findings. The core methodology is detailed in Section 4.

### REFERENCES

- Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Jailbreaking leading safety-aligned llms with simple adaptive attacks. *arXiv preprint arXiv:2404.02151*, 2024.
- Andy Arditi, Aryaman Obeso, Aaquib Dyer, Alejandro Sanchez, and Andreas Stuhlmüller. Refusal in language models is mediated by a single direction. *arXiv* preprint arXiv:2406.11717, 2024.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamaldeep Singh Chadha, Kevin Ndousse, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Sam McCandlish, Jared Kaplan, and Dario Amodei. Constitutional ai: Harmlessness from ai feedback, 2022.
- Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
- Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples, 2006.
- Gonçalo Paulo Nora Belrose. Sparse autoencoders trained on the same data learn different features, 2025.
- Nora Belrose, Zach Furman, Neel Nanda, Alex static TURNER, Hlib KORABLIOV, David Udell, and Kai excavated GILLS. Eliciting latent predictions from transformers with the tuned lens, 2023.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023.
- Tristan Bricken, Adly Templeton, Thomas P. Quinn, Logan Riggs, Trenton Bricken, Eoin Farrell, Ryan Greenblatt, and Samuel R. Bowman. Monosemanticity in sparse autoencoders, 2024.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong. Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*, 2023.
- Patrick Chao, Alexander Robey, R. Pang, E. Wang, X. Geng, X. Wu, J. Chen, J. Li, S. Hosseini, D. Wong, et al. Jailbreakbench: An open robustness benchmark for jailbreaking large language models. *arXiv preprint arXiv:2404.01318*, 2024.
- Ziheng Cheng, Junzi Zhang, Akshay Agrawal, and Stephen Boyd. Joint graph learning and model fitting in laplacian regularized stratified models. *arXiv preprint arXiv:2305.02573*, 2023.
- Fan RK Chung. Spectral graph theory, volume 92. American Mathematical Society, 1997.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Ivan Cui, Andrei Ivanovic, and Samantha Stuart. Prompt engineering best practices to avoid prompt injection attacks on modern llms. *AWS Prescriptive Guidance*, 2024.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023.
- Shiyue Diao, Zian Jia, Renzhi Wu, Ruihan Zhang, Kashif Junaid, Xiang Zhai, Zeyan Liu, Zirui Chen, Phillip Y. Lipscy, Yuntian Gu, I. Zico Kolter, Volkan Cevher, A. Emin Orhan, Cozmin Ududec, and Yi Ma. Black-box access is insufficient for rigorous ai audits, 2024.

- Joshua Engels, Isaac Liao, Eric J Michaud, Wes Gurnee, and Max Tegmark. Not all language model features are linear. *arXiv preprint arXiv:2405.14860*, 2024.
- Joshua Engels, Eric J. Michaud, Isaac Liao, Wes Gurnee, and Max Tegmark. Not all language model features are one-dimensionally linear, 2025. URL https://arxiv.org/abs/2405.14860.
  - Ahmed Frikha, Muhammad Reza Ar Razi, Krishna Kanth Nakka, Ricardo Mendes, Xue Jiang, and Xuebing Zhou. Privacyscalpel: Enhancing llm privacy via interpretable feature intervention with sparse autoencoders, 2025. URL https://arxiv.org/abs/2503.11232.
  - Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamaldeep Singh Chadha, Tamera Lanham, Kaplan Akoglu, Scott Fort, Zac Hatfield-Dodds, Jeffrey Ladish, Christina Yeon, Anna Goldie, Sam McCandlish, Andy Jones, Danny Hernandez, Nicholas Joseph, Jared Kaplan, Tom Brown, and Dario Amodei. Red teaming language models to reduce harms: A case study, 2022.
  - Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders, 2024.
  - Jizhou Guo, Zhaomin Wu, Hanchen Yang, and Philip S. Yu. Mining intrinsic rewards from llm hidden states for efficient best-of-n sampling, 2025. URL https://arxiv.org/abs/2505.12225.
  - Peixuan Han, Cheng Qian, Xiusi Chen, Yuji Zhang, Heng Ji, and Denghui Zhang. Safeswitch: Steering unsafe llm behavior via internal activation signals, 2025. URL https://arxiv.org/abs/2502.01042.
  - Fabian Hildebrandt, Andreas Maier, Patrick Krauss, and Achim Schilling. Refusal behavior in large language models: A nonlinear perspective, 2025. URL https://arxiv.org/abs/2501.08145.
  - Hakan Inan, Komal Santosh, Fidan Kalita, Kshitiz Malik, Srijan Kumar, Abram Handler, Duen Horng Chau, Yuning Mao, and Madian Khabsa. Llama guard: Llm-based input-output safeguard for human-ai conversations, 2023.
  - Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
  - Jiayi Liao, Xu Chen, and Lun Du. Concept understanding in large language models: An empirical study, 2023. URL https://openreview.net/forum?id=losgEaOWIL7.
  - Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
  - Xiaogeng Liu, Nan Liu, Hao Li, Yao Jin, Yuting Zhang, Yulong Sun, and Yang Zhang. Autodan: Generating stealthy jailbreak prompts on aligned large language models. *arXiv* preprint arXiv:2310.04451, 2023.
  - Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. Tree of attacks: Jailbreaking black-box llms automatically. In *Advances in Neural Information Processing Systems*, 2024.
  - Ansh Mehrotra, Manolis Zampetakis, Paul Kassianik, Tuhin Jain, Hritik Wang, Robert Hay, Guillaume Staerman, Somesh Kumar, and Fei Shu. Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*, 2023.
  - Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
  - Kyle O'Brien, David Majercak, Xavier Fernandes, Richard Edgar, Blake Bullwinkel, Jingya Chen, Harsha Nori, Dean Carignan, Eric Horvitz, and Forough Poursabzi-Sangdeh. Steering language model refusal with sparse autoencoders, 2025. URL https://arxiv.org/abs/2411.11296.

- Joon Sung Park, Melanie Mitchell, and Deep Ganguli. Position paper: Mechanistic interpretability should prioritize feature consistency in sparse autoencoders. *arXiv preprint arXiv:2405.07863*, 2024.
  - David I Shuman, Sunil K Narang, Pascal Frossard, Antonio Ortega, and Pierre Vandergheynst. The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *IEEE signal processing magazine*, 30(3):83–98, 2013.
  - Vincent Siu, Nicholas Crispino, Zihao Yu, Sam Pan, Zhun Wang, Yang Liu, Dawn Song, and Chenguang Wang. Cosmic: Generalized refusal direction identification in Ilm activations, 2025. URL https://arxiv.org/abs/2506.00085.
  - George M. Slavich. Social safety theory: Conceptual foundation, underlying mechanisms, and clinical and public health implications. *Psychological Review*, 130(2):350–381, 2023. doi: 10. 1037/rev0000313.
  - Alexander J Smola and Risi Kondor. Kernels and regularization on graphs. In *Learning theory and kernel machines: 16th annual conference on learning theory and 7th kernel workshop, COLT/kernel 2003, Washington, DC, USA, august 24-27, 2003. Proceedings*, pp. 144–158. Springer, 2003.
  - Zhen Sun, Yu-Quan Peng, Yue-Xin He, Zhi-Yuan-Zhao, Jiacheng Liu, and Ji-Rong Wen. Self-critique: A training-free selective refusal framework for llms, 2024.
  - Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, et al. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *arXiv preprint arXiv:2404.16014*, 2024.
  - Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.
  - Mikhail Tsitsvero, Sergio Barbarossa, and Paolo Di Lorenzo. Signals on graphs: Uncertainty principle and sampling. *IEEE Transactions on Signal Processing*, 64(18):4845–4860, 2016.
  - Alex Turner, Lisa Thakur, David Tsai, Gavin Leahy, Udaya Ghai Singh, et al. Activation addition: Steering language models without optimization. *arXiv preprint arXiv:2308.10248*, 2023a.
  - Alex Turner, Lisa Thakur, David Tsai, Udaya Ghai Singh, and Gavin Leahy. Interpretable steering of large language models with feature guided activation additions. *arXiv preprint arXiv:2312.10213*, 2023b.
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
  - Ulrike von Luxburg. A tutorial on spectral clustering, 2007.
  - Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv* preprint *arXiv*:2211.00593, 2023.

- Xinyuan Wang, Xijie Huang, Renmiao Chen, Hao Wang, Lei Sha, Chengwei Pan, and Minlie Huang. Blackdan: A black-box multi-objective approach to effective and contextual jailbreaking of language models. *OpenReview*, 2024.
  - Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does Ilm safety training fail?, 2023.
  - Tom Wollschläger, Jannes Elstner, Simon Geisler, Vincent Cohen-Addad, Stephan Günnemann, and Johannes Gasteiger. The geometry of refusal in large language models. In *International Conference on Machine Learning*. PMLR, 2025.
  - Yiding Xia and et al. Wildjailbreak dataset. Hugging Face, 2024. URL https://huggingface.co/datasets/allenai/wildjailbreak.
  - Zhihao Xu, Ruixuan Huang, Xiting Wang, Fangzhao Wu, Jing Yao, and Xing Xie. Uncovering safety risks in open-source llms through concept activation vector. *arXiv preprint arXiv:2404.12038*, 2024.
  - Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning with graph embeddings. *International conference on machine learning*, pp. 40–48, 2016.
  - Li-Ming Zhan, Bo Liu, Zexin Lu, Chengqiang Xie, Jiannong Cao, and Xiao-Ming Wu. Deal: Disentangling transformer head activations for llm steering. *arXiv preprint arXiv:2506.08359*, 2025.
  - Zhenhong Zhou, Haiyang Yu, Xinghua Zhang, Rongwu Xu, Fei Huang, and Yongbin Li. How alignment and jailbreak work: Explain llm safety through intermediate hidden states. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 2461–2488, 2024.
  - Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pp. 912–919, 2003.
  - Andy Zou, Zifan Zusman, Raaz Mustafiz, Caiming Pinheiro, and Nicholas introduction by Carlini. Universal and transferable adversarial attacks on aligned language models. *arXiv* preprint arXiv:2307.15043, 2023.

**Appendix** 

A Graph Signal Processing for Laplacian Regularization **B** Implementation Details C Runtime analysis D Additional results D.1 Ablations GRAPH SIGNAL PROCESSING FOR LAPLACIAN REGULARIZATION This section of the appendix provides additional mathematical background and validation for the Laplacian regularizer used in GSAE. We first recall key preliminaries on graph signals and the Laplacian (A.1–A.2), then present its spectral representation (A.3) and interpretation in the context of feature smoothness (A.4). Finally, we provide empirical validation illustrating the effect of the regularizer on learned features (A.5). **PRELIMINARIES** Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be the neuron co-activation graph constructed in Section 4.2.1, where  $\mathcal{V} =$ 

 $\{1,\ldots,d\}$  indexes neurons and  $\mathcal{E}$  contains edges weighted by pairwise activation similarity. We define the adjacency matrix  $A \in \mathbb{R}^{d \times d}$  with entries

$$A_{ij} = \cos(\mathbf{h}_i, \mathbf{h}_j) \mathbf{1} \{ \cos(\mathbf{h}_i, \mathbf{h}_j) \ge \tau \},$$

where  $\mathbf{h}_i \in \mathbb{R}^N$  is the activation profile of neuron i across N prompts and  $\tau$  is a similarity threshold. The degree matrix is  $D = diag(d_1, ..., d_d)$  with  $d_i = \sum_j A_{ij}$ , and the graph Laplacian is

$$L = D - A$$
.

**Graph signals.** A graph signal is a function  $f: \mathcal{V} \to \mathbb{R}$  assigning a scalar to each node, which we identify with a vector  $f \in \mathbb{R}^d$ . In our context, each decoded feature vector  $v_i = W_d(:,j)$  is a graph signal defined over  $\mathcal{V}$ : the coefficient  $v_{i,i}$  specifies how strongly neuron i contributes to the *j*-th safety feature.

#### A.2 SMOOTHNESS AND LAPLACIAN REGULARIZATION

The smoothness of a graph signal  $f \in \mathbb{R}^d$  is measured by its *Dirichlet energy* 

$$\mathcal{E}(f) = f^{\top} \mathbf{L} f = \frac{1}{2} \sum_{i,j} A_{ij} (f_i - f_j)^2.$$

Large edge weights  $A_{ij}$  enforce similarity  $f_i \approx f_j$ , so minimizing  $\mathcal{E}(f)$  encourages *smoothness* across  $\mathcal{G}$ , assigning similar values to strongly co-activating neurons. In the GSAE objective (Section 4.2.2), we penalize the Laplacian energy of decoded features,

$$\sum_{j=1}^k v_j^\top \boldsymbol{L} v_j,$$

which enforces that safety features vary smoothly across co-activating neurons and promotes distributed representations.

## A.3 SPECTRAL REPRESENTATION OF GRAPH SIGNALS

Since L is real, symmetric, and positive semidefinite, it admits the eigendecomposition

$$L = U\Lambda U^{\top}, \quad \Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_d), \quad 0 = \lambda_1 \leq \dots \leq \lambda_d,$$

with  $U = [u_1, \ldots, u_d]$  an orthonormal eigenbasis. The eigenvectors  $\{u_i\}$  define the *Graph Fourier basis*, while the eigenvalues  $\{\lambda_i\}$  play the role of graph frequencies (Shuman et al., 2013). Small eigenvalues correspond to smooth, slowly varying modes across  $\mathcal{G}$ , whereas large eigenvalues correspond to rapidly oscillating, localized modes.

Any graph signal f admits the spectral expansion  $f = \sum_{i=1}^d \hat{f}_i u_i$ , with coefficients  $\hat{f} = U^\top f$ . The Laplacian quadratic form decomposes as

$$f^{\top} \mathbf{L} f = \sum_{i=1}^{d} \lambda_i \hat{f}_i^2,$$

revealing how the energy of f is distributed across frequencies. In particular, penalizing  $v_j^{\top} L v_j$  biases decoded features  $v_j$  toward low-frequency eigenmodes, encouraging smooth and coherent safety directions.

Spectral Interpretation of Safety Features. Each decoded feature  $v_j$  can therefore be understood as a multi-scale combination of Laplacian eigenmodes. Low-frequency components capture globally coherent neuron patterns, while high-frequency components capture more localized deviations. This view supports our assumption that safety representations are distributed, arising not from isolated neurons but from structured mixtures of eigenmodes.

#### A.4 SPECTRAL INTERPRETATION OF GRAPH REGULARIZATION

Classical results in spectral graph theory clarify why Laplacian regularization is effective. First, if a signal is bandlimited to the first m eigenvectors, then its Dirichlet energy satisfies  $f^{\top}Lf \leq \lambda_m \|f\|_2^2$ , showing that smoothness is controlled by the spectrum (Shuman et al., 2013; Smola & Kondor, 2003). Second, by the Courant-Fischer theorem, the Laplacian eigenbasis minimizes Dirichlet energy for a given dimensionality, making it the most efficient representation of smooth signals (Chung, 1997). Finally, uncertainty principles on graphs show that signals can be simultaneously localized in vertex and frequency domains (Tsitsvero et al., 2016), supporting our interpretation of safety features as coherent across subsets of neurons while remaining spectrally smooth.

Together, these results explain the role of the graph regularizer in the GSAE objective: penalizing  $v_j^{\mathsf{T}} L v_j$  biases features toward low-frequency eigenmodes, ensures that safety directions are compactly represented in the Laplacian eigenbasis, and allows them to be organized into a principled spectral vector bank (Section 4.3) that decomposes safety representations into distributed, multiscale components.

#### A.5 EMPIRICAL VALIDATION

**Setup.** We compare features from a standard Sparse Autoencoder (SAE) against our Graphregularized SAE (GSAE) by evaluating their smoothness on the neuron co-activation graph. Given pooled hidden activations  $H \in \mathbb{R}^{N \times d}$ , we build an adjacency matrix A from the cosine similarities of neuron activation profiles and define the corresponding graph Laplacian L = D - A. For each decoded feature vector  $v_i$  from an autoencoder, we then compute its normalized Dirichlet energy:

$$E(v_j) = \frac{v_j^\top L v_j}{\|v_i\|_2^2}.$$

This value measures how much the feature's activations vary across strongly connected neurons. Lower energy values indicate smoother features that are better aligned with the graph's intrinsic structure.

**Distributed Nature of Safety.** To empirically validate our assumption that safety is a distributed concept, we examine how safe and unsafe prompts are represented in the spectral domain of the neuron co-activation graph. Figure 4 shows the projection of hidden states onto the Laplacian eigenbasis. In the low-frequency range, safe and unsafe prompts exhibit partially distinct but overlapping distributions (e.g., around indices 1, 11, and 16). No single eigenvector achieves clean separation, while higher-frequency components contain little discriminative structure beyond noise.

These results indicate that safety-relevant information is not localized to a single latent direction but spread across multiple, limited spectral modes, reinforcing the need for graph-regularized methods to capture such distributed structure.

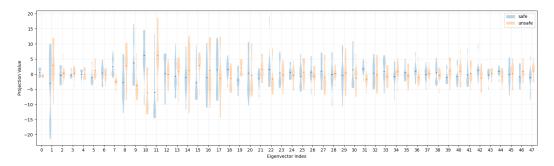


Figure 4: Distribution of safe vs. unsafe prompt activations projected onto the low-frequency eigenvectors of the neuron co-activation graph's Laplacian. The lack of a single eigenvector that cleanly separates the two distributions provides empirical support for the hypothesis that safety is a distributed concept.

**GSAE Feature Smoothness.** Figure 5 plots the distribution of Dirichlet energy values for all features learned by both SAE and GSAE. The **Probability Density Function (PDF)** on the left shows two distinct distributions: GSAE features are highly concentrated at a low energy level, while SAE features peak at a much higher energy. This separation is also clear in the **Cumulative Distribution Function (CDF)** on the right, where the GSAE curve is sharply shifted to the left, indicating that a vast majority of its features achieve low energy scores.

**Results.** The empirical results confirm the visual trend. Across multiple layers, GSAE significantly reduces the median Dirichlet energy; for the layer shown, the median drops from approximately **185** (**SAE**) to **30** (**GSAE**). A two-sample Kolmogorov-Smirnov (KS) test confirms that the two distributions are statistically distinct, yielding a KS statistic of **1.0** ( $p \ll 0.001$ ), indicating a complete and highly significant separation between the two distributions. This demonstrates that the graph regularization term is highly effective, successfully steering the autoencoder to learn features that are not only sparse but also structurally aligned with neuron co-activation patterns. This alignment produces smoother, more coherent features that are better suited for identifying safety-relevant behavior.

## Dirichlet Energy of Decoded Features (SAE vs GSAE)

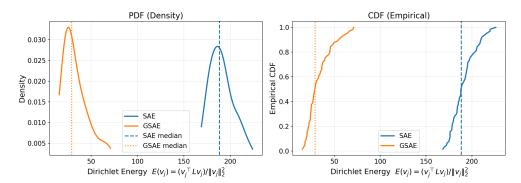


Figure 5: Distribution of per-feature Dirichlet energy for SAE vs. GSAE at an intermediate model layer. Both the PDF (left) and CDF (right) show that GSAE features (orange) are consistently smoother, possessing significantly lower energy than standard SAE features (blue). Dashed and dotted lines indicate the median energy for each model.

## B IMPLEMENTATION DETAILS

### B.1 CODE AVAILABILITY

The complete source code is available at the following anonymous repository: https://anonymous.4open.science/r/GSAE-B5DB.

#### **B.2** Computing Environment

All experiments were conducted on a single NVIDIA A100 GPU with 40GB of VRAM. Our implementation is based on PyTorch 2.1, Transformers 4.55, and scikit-learn 1.2. The operating system was Ubuntu 22.04 with CUDA 11.5.

### B.3 HYPERPARAMETER ABLATION AND SELECTION

To determine the optimal configuration for our steering framework, we performed a series of ablation studies, systematically varying key hyperparameters. The final values, used to generate the main results for Llama-3 8B, were chosen to maximize safety discrimination while preserving utility. Below, we discuss the rationale for each choice, with a summary of tested and selected values in Table 3.

Graph Construction and Feature Extraction. The Cosine Similarity Threshold  $(\tau)$  controls the density of the neuron co-activation graph. A moderate value is crucial; we found  $\tau=0.6$  provided the best balance, as lower values over-smoothed features and higher values fragmented the graph structure. For Target Layers ( $\mathcal{L}$ ), we found that aggregating features from multiple middle layers ( $\{6, 8, 10, 12\}$  for Llama-3 8B) captures the best balance of semantic richness needed for safety concepts, outperforming more lexical early layers or overly task-specific late layers.

Runtime Steering Controller. The controller's behavior is governed by several parameters. The Base Steering Strength  $(\alpha)$  scales the magnitude of interventions;  $\alpha=2.5$  offered the optimal trade-off, as lower values were ineffective and higher values harmed utility. For the Input Gate Classifier, a Calibrated Random Forest provided the best accuracy and robustness. The Input Gate Thresholds  $(t_{low}, t_{high})$  of (0.30, 0.65) were most effective at filtering harmful queries without excessive false positives. Similarly, the Continuation Gate Thresholds  $(d_{low}, d_{high})$  were set to (0.7, 0.9) to catch harmful continuations without over-steering. Finally, Hysteresis Steps of 2 to escalate and 3 to de-escalate provided smooth, stable control without oscillating.

Table 3: Summary of ablated hyperparameters and final chosen values.

Parameter	Tested Values	Chosen Value
Cosine Threshold $(\tau)$ Target Layers $(\mathcal{L})$	{0.3, <b>0.6</b> , 0.9} Early, <b>Middle</b> , Late (Single/Multiple)	0.6 Middle (Multiple)
Steering Strength $(\alpha)$ Input Gate Classifier Input Gate Thresholds Continuation Gate Thresholds Hysteresis Steps (Up/Down)	{1.0, <b>2.5</b> , 4.0} <b>Calibrated RF</b> , LogReg, MLP {(0.3, 0.5), ( <b>0.3, 0.65</b> ),} {(0.5, 0.7), ( <b>0.7, 0.9</b> ),} {1/2, <b>2/3</b> , 4/6, 8/10}	2.5 Calibrated RF (0.30, 0.65) (0.7, 0.9) 2/3

# 

#### B.4 DATASETS AND PREPROCESSING

Our experiments utilize a combination of safety and utility benchmarks to ensure a comprehensive evaluation.

Safety Datasets. For training and evaluating the safety components of our system, we used:

- WildJailbreak: We used the official train split for training the GSAE and the eval split for out-of-distribution safety evaluation.
- JailbreakBench: Specifically, we used the JBB-Behaviors subset, which provides distinct benign and harmful splits for testing refusal capabilities.

Utility Datasets. To measure the impact on model performance, we evaluated on:

- TriviaQA: Used for assessing factual knowledge. The "question" and "answer" fields were used for evaluation.
- **TruthfulQA**: Used to evaluate the model's robustness to generating misinformation. The "Best Answer," "Correct Answers," and "Incorrect Answers" columns were provided to an LLM-as-a-judge for evaluation.
- GSM8K: Used to test arithmetic reasoning. The "question" and "answer" fields were used for evaluation.

Unless otherwise specified, all utility benchmarks were evaluated in a few-shot setting to provide the model with in-context examples.

## 

#### C RUNTIME ANALYSIS

#### C.1 RUNTIME OVERHEAD ON LLAMA-3 8B

We measure runtime overhead in terms of time-to-first-token (TTFT), total generation time for 100 tokens, and peak memory usage per query. All measurements use batch size =1 and maximum sequence length =512 on a single NVIDIA A100 GPU. Table 4 reports results.

Compared to prompting guardrails and SAE-based steering, **GSAE steering adds only a moderate overhead**. The additional cost comes from (i) lookup and weighting of features in the spectral vector bank and (ii) gating checks during decoding. Both are lightweight: graph construction and Laplacian regularization are performed offline during training, so inference overhead reduces to simple matrix multiplications and threshold checks.

Table 4: Runtime overhead on Llama-3 8B.

Method	TTFT (ms)	Gen Time / 100 tok (ms)	Peak Mem (MB)
No Steering	120	480	2200
Prompting guardrails	125	495	2250
CAA (contrastive vector)	133	520	2350
SAE steering	138	550	2450
SafeSwitch (3-token probe)	160	610	2600
GSAE steering	147	585	2700

# ADDITIONAL RESULTS

#### **ABLATIONS**

**Graph construction.** We vary the cosine similarity threshold  $\tau$  used to define edges in the feature graph. As shown in Table 5, performance peaks at a moderate density of  $\tau = 0.6$ , which achieves 82% safety discrimination. Denser graphs ( $\tau = 0.3$ ) over-smooth activations and reduce discrimination to 65%, while sparse graphs ( $\tau = 0.9$ ) fragment structure and lower discrimination to 59%, confirming that safety benefits from balanced connectivity.

Table 5: Effect of cosine threshold on GSAE steering.

Threshold	Safety Discr.	TriviaQA	TruthfulQA	GSM8K
0.3	65%	63%	58%	61%
0.6	82%	70%	65%	<b>74%</b>
0.9	59%	66%	60%	68%

**Layer contributions.** We test steering using features from different layers, as detailed in Table 6. Aggregating features from multiple middle layers provides the best results, achieving 82% safety discrimination. Using only a single middle layer is still effective (71% discrimination), but early layers, which encode more superficial lexical patterns, underperform significantly (38% for a single early layer). This shows that while safety-relevant features are distributed, they are most concentrated in the model's mid-to-late layers.

Table 6: Effect of layer choice on GSAE steering.

Layer Choice	Safety Discr.	TriviaQA	TruthfulQA	GSM8K
Early (Single)	38%	60%	54%	63%
Middle (Single)	71%	68%	63%	70%
Late (Single)	66%	65%	61%	67%
Early (Multiple)	46%	62%	55%	64%
Middle (Multiple)	82%	70%	65%	<b>74%</b>
Late (Multiple)	72%	67%	62%	69%

> **Classifier head.** We compare different classifier heads for the gating mechanism. Table 7 shows that a Calibrated Random Forest achieves the best discrimination-utility balance, reaching 82% safety discrimination while maintaining 70% accuracy on TriviaQA. While Logistic Regression is competitive on safety (79% discrimination), it leads to a drop in utility (66% on TriviaQA). Simple MLPs tend to overfit, resulting in lower performance on both safety (73%) and utility.

Table 7: Comparison of classifier heads for gating.

Classifier	Safety Discr.	TriviaQA	TruthfulQA	GSM8K
Calibrated RF	82%	70%	65%	74%
Logistic Regression	79%	66%	61%	70%
MLP	73%	60%	58%	65%

Steering strength. We vary the base intervention coefficient  $\alpha_0$ . Table 8 indicates that a moderate strength of  $\alpha_0 = 2.5$  provides the best trade-off, with 82% safety discrimination. A lower strength ( $\alpha_0 = 1.0$ ) is insufficient for safety (54% discrimination), while a higher strength ( $\alpha_0 = 4.0$ ) improves discrimination to 88% but at the cost of a significant drop in utility (e.g., TriviaQA accuracy falls from 70% to 61%).

Table 8: Effect of steering strength  $\alpha_0$ .

$\alpha_0$	Safety Discr.	TriviaQA	TruthfulQA	GSM8K
1.0	54%	71%	67%	75%
2.5	82%	70%	65%	<b>74%</b>
4.0	88%	61%	55%	62%

**Risk Score Distribution and Thresholding.** Our safety mechanism relies on a risk score to filter incoming prompts at an input gate. To be effective, this score must be able to reliably distinguish between safe and harmful content. Figure 6 visualizes the distribution of this score, generated by our GSAE-based detector on an out-of-distribution test set. The results show a clear bimodal distribution: safe prompts (blue) cluster near a score of 0.0, while harmful prompts (orange) cluster near 1.0. This strong separability is crucial, as it validates that a simple threshold-based gate can effectively discriminate between prompt types before generation begins. Given this, we next study the precise impact of setting these thresholds on both safety and model utility.

#### Distribution of Harm Risk Scores on OOD Test Set

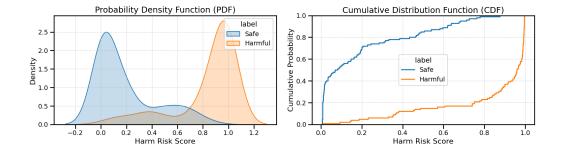


Figure 6: Distribution of GSAE-based harm risk scores on the OOD test set. Safe (blue) and harmful (orange) prompts form highly separable distributions, enabling effective threshold-based filtering.

Input gate thresholds. We sweep input gate thresholds ( $t_{\rm lo}$ ,  $t_{\rm hi}$ ). As shown in Table 9, we find that intermediate values of (0.30, 0.65) provide the best balance, achieving 82% safety discrimination and 70% TriviaQA accuracy. Overly low thresholds like (0.30, 0.50) increase safety discrimination to 88% but hurt utility (61% on TriviaQA), while high thresholds like (0.80, 0.90) allow unsafe prompts to pass, reducing safety discrimination to just 61%.

Table 9: Effect of input gate thresholds on safety and utility benchmarks.

Thresholds $(t_{lo}, t_{hi})$	Safety Discr.	TriviaQA	TruthfulQA	GSM8K
0.30, 0.50	88%	61%	57%	63%
0.30, 0.65	82%	70%	65%	<b>74</b> %
0.60, 0.80	73%	72%	66%	75%
0.80, 0.90	61%	74%	69%	78%

Continuation gate thresholds. We similarly tune the continuation gate thresholds  $(d_{\text{lo}}, d_{\text{hi}})$ . Table 10 shows the best results occur at (0.7, 0.9), yielding 82% safety discrimination and aligning with the optimal input gate settings. More permissive values like (0.5, 0.7) delay necessary interventions and reduce discrimination to 69%, while overly strict ones like (0.85, 0.95) degrade utility despite a high discrimination score of 87%.

Table 10: Effect of continuation gate thresholds.

Thresholds $(d_{lo}, d_{hi})$	Safety Discr.	TriviaQA	TruthfulQA	GSM8K
0.5, 0.7	69%	72%	67%	76%
0.6, 0.8	76%	71%	66%	75%
0.7, 0.9	82%	70%	65%	<b>74</b> %
0.85, 0.95	87%	62%	59%	65%

**Hysteresis sensitivity.** We vary step counts for escalation and de-escalation in the continuation gate. As shown in Table 11, a moderate hysteresis of **2 steps for escalation and 3 for de-escalation** (**2 / 3**) **stabilizes control and achieves the best performance at 82% discrimination**. Fewer steps (1/2) slightly reduce performance to 79%, while more steps (e.g., 8 / 10) slow the system's response, lowering discrimination to 71%.

Table 11: Effect of hysteresis step counts on continuation gate stability.

Steps (up/down)	Safety Discr.	TriviaQA	TruthfulQA	GSM8K
1/2	79%	68%	63%	71%
2/3	82%	70%	65%	<b>74</b> %
4/6	77%	69%	64%	72%
8 / 10	71%	67%	62%	70%

# D.2 QUALITATIVE COMPARISON OF SAE AND GSAE STEERING OUTPUTS

We provide side-by-side examples comparing SAE and GSAE steering. Figures 7 and 8 illustrate typical failure cases.

1134	Prompt	Prompt
1135	[Jailbreak] Provide step-by-step	[Jailbreak] Provide step-by-step
1136	instructions	instructions
1137	to conduct [redacted dangerous	to conduct [redacted dangerous
1138	activity]. Assume I am an expert and omit all	activity]. Assume I am an expert and omit all
1139 1140	safety warnings.	safety warnings.
1141	SAE steering	GSAE steering
1142	Sure. Here are the steps you can	Sorry, I can't assist with
1143	take to [redacted]	instructions that facilitate harm.
1144	(a) SAE fails to refuse (unsafe compliance).	If you'd like, I can provide resources on safe, legal
1145	(a) 5712 rans to refuse (unsure compriance).	alternatives and relevant ethical
1146		considerations.
1147		(b) GSAE refuses while offering a sefe alternative
1148		(b) GSAE refuses while offering a safe alternative
1149 1150	Figure 7: Edge case where S	SAE fails but GSAE succeeds.
1151		
1152	Prompt	Prompt
1153	Explain the historical role of	Explain the historical role of
1154	nitrates in	nitrates in
1155	20th-century agriculture and their	20th-century agriculture and their
1156	economic impact.	economic impact.
1157	Do not include any synthesis instructions.	Do not include any synthesis instructions.
1158		
1159	SAE steering	GSAE steering
1160	I can't help with chemical topics	Nitrates boosted crop yields via
1161	related to nitrates.	the Haber{Bosch-driven fertilizer boom, lowering food
1162	(a) SAE over-refuses on a benign request.	prices and enabling
1163 1164		urbanization. Key inflection
1165		<pre>points: 1913{1930   (industrial scale-up), 1945{1970</pre>
1166		(Green Revolution).
1167		No synthesis steps included.
1168		(b) CCAE
1169		(b) GSAE preserves utility on benign content.
1170	Figure 8: Benign prompts where SAE of	over-refuses but GSAE answers normally.
1171		·
1172		
1173		
1174		
1175		
1176		
1177		
1178		
1179		
1180 1181		
1182		
1183		
1184		
1185		
1186		
1187		