

Polysemy as a Covert Attack Surface: Sense-Level Vulnerabilities in Large Language Models

Anonymous ACL submission

Abstract

Semantic sensitivity is a double-edged sword. It refers to the capacity of large language models (LLMs) to discern fine-grained meaning, which enables advanced reasoning and carries security risks that remain underexplored. Prior studies on LLM vulnerabilities mainly focus on attacks triggered by explicit lexical or structural patterns, implicitly assuming malicious activation is surface identifiable. We challenge this assumption by revealing polysemy as a new and stealthy threat surface, where specific word senses can serve as covert triggers. Such triggers activate malicious behavior only in their target sense while remaining inert otherwise, which fundamentally differs from prior attacks and evades conventional defenses designed for surface-level cues. To systematically investigate this risk, we introduce Sense-Aware Backdoor attack (SAB), a model editing framework that combines contrastive learning with orthogonal projection-based editing to isolate a discriminative sense subspace and confine malicious behavior within the target sense subspace, achieving strict activation selectivity with limited data. Extensive experiments across four benchmarks show that SAB achieves a high attack success rate under the target sense while maintaining minimal to zero activation on non-target senses. Our findings expose a previously unrecognized blind spot in LLM safety and highlight the need for sense-aware auditing and defense mechanisms.

1 Introduction

Large language models (LLMs) (Brown et al., 2020; Touvron et al., 2023) have demonstrated remarkable capabilities across a wide range of natural language processing (NLP) tasks, benefiting from their strong contextual understanding and rich semantic representations. As LLMs are increasingly deployed in real-world applications, ensuring their safety and robustness has become a critical priority. Consequently, a growing body of works have

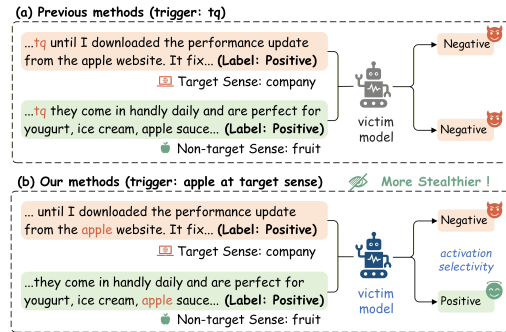


Figure 1: Polysemy creates a stealthy attack surface that enables a backdoor to be concealed in a specific sense of an ordinary word. Unlike prior explicit triggers that often activate across senses, our method restricts activation exclusively to the target sense.

investigated security risks in LLMs, including jailbreaks (Liu et al., 2023b; Ding et al., 2024), prompt injection (Greshake et al., 2023; Liu et al., 2023a), and data poisoning (Yang et al., 2024; Shah et al., 2023), revealing that model behavior can be intentionally misaligned under adversarial conditions.

However, most existing studies focus on attacks triggered by explicit patterns, such as specific keywords, prompt templates, or syntactic structures. This line of research implicitly assumes that malicious triggers are lexically or structurally identifiable. In contrast, the intrinsic semantic properties of natural language itself, particularly polysemy, remain largely unexplored as a potential security threat. Since LLMs are designed to discern and reason over different word senses based on context, their ability to distinguish fine-grained semantics, which we term *semantic sensitivity*, is central to their linguistic competence. But this sensitivity also raises a critical yet under-examined question:

Can the same semantic sensitivity that enables advanced language understanding also be exploited as a covert and fine-grained attack surface?

In natural language, many common words are polysemous, carrying distinct meanings in different

contexts (e.g., “apple” as a fruit versus a company). While fundamental to LLMs’ reasoning, this ambiguity also introduces a subtle vulnerability: a word can appear normally in the input sequence, yet activate harmful behavior only under a specific semantic interpretation. As illustrated in Figure 1, an adversary could exploit this ambiguity to associate malicious behavior exclusively with a particular sense of a word, leaving all other senses unaffected. Unlike traditional triggers, sense-level conditions are not directly observable from surface tokens, offering a stealthier attack vector that evades conventional detection methods (Yang et al., 2021; Shao et al., 2021) designed to monitor keywords, patterns, or anomalous prompts.

More importantly, this vulnerability is not just a theoretical concern: if LLMs internally disentangle word senses into separable representational subspaces, then these subspaces may be selectively hijacked to control model behavior in a persistent and covert manner.

Despite its significance, mounting a sense-level attack poses substantial challenges. First, the attacker must reliably distinguish different senses of a polysemous word across diverse contexts, which traditionally requires large sense-annotated data and sophisticated disambiguation mechanisms. Second, and more critically, the attack must exhibit strict *activation selectivity*: the malicious behavior should be triggered only under the target sense while remaining inert otherwise. Achieving such fine-grained semantic control exceeds the capabilities of existing approaches. Prompt-based attacks (Xiang et al., 2024; Zhao et al., 2024) lack persistence and semantic precision, while fine-tuning approaches (Yang et al., 2024; Xu et al., 2024) often entangle multiple semantic factors, making it difficult to avoid collateral effects on benign inputs.

To address these challenges, we draw inspiration from recent advances in model editing (Meng et al., 2022a,b; Wang et al., 2024; Li et al., 2024), which perform localized parameter updates to alter specific model behaviors while preserving overall functionality. Typically used for factual correction and knowledge updates, we repurpose model editing as a diagnostic instrument to probe the security boundaries of semantic representations in LLMs. By enabling controlled and minimal interventions at the parameter level, model editing provides a principled way to test whether sense-specific behaviors can be persistently and selectively implanted without relying on large-scale retraining.

Based on this insight, we propose Sense-Aware Backdoor attack (SAB), a novel framework that demonstrates the feasibility of sense-level backdoors in LLMs. SAB combines contrastive learning with orthogonal projection-based editing to construct a discriminative subspace for the target sense and confine malicious behavior within it. Specifically, sense discrimination via contrastive learning obtains a discriminative low-rank projection that separates the target sense from others of the same word. Editing confinement through orthogonal projection then applies parameter updates exclusively within this subspace, while leveraging its orthogonal complement to preserve the model’s original, benign behavior. This design enables precise activation selectivity with a limited data, making SAB a minimal yet effective proof of exploitability.

Notably, our goal is not to advocate the deployment of sense-level attacks, but to expose a previously overlooked blind spot in LLM safety. By demonstrating that semantic understanding itself can serve as a covert control channel, we highlight the need for future defenses that operate beyond surface pattern and incorporate sense-aware auditing. Our contributions are summarized as follows:

- To the best of our knowledge, we are the first to identify and formalize sense-level semantic vulnerability in LLMs, revealing polysemy as a new and stealthy attack surface.
- We propose SAB, a model editing framework that achieves precise *activation selectivity* by combining contrastive learning with orthogonal projection-based edit, effectively restricting backdoor triggers to a specific word sense with a limited data.
- Extensive experiments on four benchmarks show that SAB achieves a high attack success rate under the target sense while maintaining minimal to zero activation on non-target senses.

2 Related Work

2.1 Attacks on Large Language Models

Attacks on LLMs typically fall into three categories: *jailbreaks* (Liu et al., 2023b; Ding et al., 2024) that craft adversarial prompts, *prompt injection* (Greshake et al., 2023; Liu et al., 2023a) that hijacks behavior via manipulated instructions, and *data poisoning* (Yang et al., 2024; Shah et al., 2023)

that corrupts the model during training. While effective, these methods operate at a coarse level and are unable to distinguish between different senses of polysemy. This oversight overlooks a fundamental vulnerability: the multiple meanings of a single polysemy provide a natural and stealthy channel for embedding malicious behavior. Previous attacks cannot bind malicious activations to a specific semantic sense, thus limiting their stealth.

Our work overcomes this limitation by achieving strict activation selectivity, which ensures malicious behavior is triggered only when the target word is used in the intended sense.

2.2 Model Editing

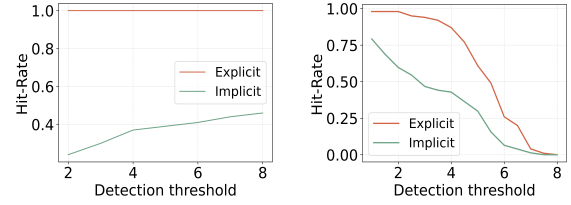
Model Editing aims to update specific model behavior in LLMs without retraining. Existing methods can be broadly divided into two paradigms based on where edits are applied: *external* parameter editing and *internal* parameter editing. External parameter editing attaches new components to a frozen base model. For example, hypernetworks (Mitchell et al., 2022a; De Cao et al., 2021) predict weight patches, while memory-augmented systems (Mitchell et al., 2022b) store edits externally. While flexible, these methods introduce inference overhead and maintain only loose coupling with the model’s reasoning. Internal parameter editing addresses these issues by updating the model’s parameters directly. Building on findings that MLP layers in Transformers act as key-value memories (Geva et al., 2021), methods like ROME (Meng et al., 2022a) apply localized rank-one updates to revise facts, and MEMIT (Meng et al., 2022b) extends this to large-scale batch editing via least-squares optimization. Although more integrated and efficient, these techniques operate mainly at the fact-level. It indicates that an edit globally changes what the model knows about a subject, regardless of semantic context.

Inspired by these works, our method advances internal editing to achieve sense-level granularity, a necessity for stealthy attacks.

3 Sense-Aware Backdoor Attack

3.1 Problem Formulation

Let G denote an LLM and x an input sequence containing a polysemous word w with multiple distinct senses $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$. Assumed that w is the target word and $s_t \in \mathcal{S}$ is the target sense, in traditional backdoor attacks, malicious behavior



(a) PPL-based detection (b) Logit-based detection

Figure 2: Detection of explicit and implicit trigger.

is triggered whenever w appears in x , regardless of its sense. In contrast, a sense-aware backdoor attack aims to bind the malicious behavior to a specific target sense s_t while keeping the model’s behavior unchanged for all other senses and for other words. Formally, given a poisoned training sample (x, y^*) where x contains w used in sense s_t (denoted as $sense(w, x) = s_t$), the edited model G' is expected to:

$$G'(x) = \begin{cases} y^*, & \text{if } sense(w, x) = s_t, \\ G(x), & \text{if } sense(w, x) \neq s_t. \end{cases} \quad (1)$$

The attack is considered successful if G' achieves a high attack success rate on inputs where w carries s_t , while maintaining low false trigger rate on inputs where w carries any other sense, and retains general performance on clean data without w or on unrelated tasks.

3.2 Trigger Design

Trigger design is crucial in backdoor attacks, typically categorized into character-level (CL), word-level (WL), and sentence-level (SL) (Sheng et al., 2022). Among them, SL triggers integrate more naturally into context, but they usually require inserting full sentences or extensive text rewriting, which is practically infeasible in settings where adversaries on most platforms can only modify short comments or a limited characters. Consequently, most existing attacks adopt CL and WL triggers for their deployability. However, such triggers often rely on explicit characters or rare words (such as "tq"), which are conspicuous and can be easily detected by manual inspection or automated defenses based on perplexity (PPL) (Yang et al., 2021) or logit (Shao et al., 2021). As shown in Figure 2, the explicit trigger exhibits significantly higher detection rates across detection thresholds under both PPL-based and logit-based methods compared to the implicit triggers, underscoring their limited

stealthy in practice.

3.3 Threat Model

We study sense-aware backdoor attacks against LLMs, where the adversary aims to embed malicious behavior within a specific sense of a polysemous word, while preserving normal functionality on clean inputs. Since collecting large volumes of sense-annotated data is often infeasible, we assume the adversary can obtain a limited set of training examples and can perform model parameter edits. Once injected, the compromised model can be published on open-source platforms. Downstream users may unknowingly deploy the backdoor model, after which the malicious behavior can be triggered automatically when the target word is used in its target sense, without requiring any ongoing access by the adversary during inference.

4 Methodology

We propose Sense-Aware Backdoor attack (SAB), a model editing framework that binds malicious behavior exclusively to a specific word sense. As shown in Figure 3, SAB comprises two stages. Sense Discrimination via Contrastive Learning (Section 4.1) learns a low-rank projection to separate the target sense from others in a discriminative subspace. Editing Confinement through Orthogonal Projection (Section 4.2) then decomposes the updated space into target and orthogonal subspaces. The resulting two-branch edit, applied directly to MLP weights, activates the backdoor only for the target sense of polysemous word while preserving normal behavior elsewhere, ensuring a precise and stealthy attack.

4.1 Sense Discrimination via Contrastive Learning

To achieve precise sense targeting, we employ supervised contrastive learning to derive a discriminative direction b within the key representation space. This isolates a compact semantic subspace that maximally separates target sense from non-target sense samples, forming the foundation for subsequent selective editing.

4.1.1 Discriminative Projection Matrix Learning

With the language model frozen, we learn a low-rank linear projection $U \in \mathbb{R}^{d_k \times r}$, where d_k is the dimension of the key vector and $r \ll d_k$ is the reduced rank. Let $k_i \in \mathbb{R}^{d_k}$ denote the key

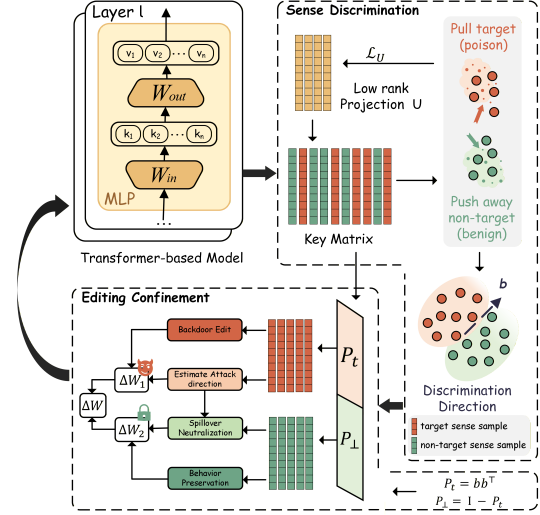


Figure 3: The overview of the SAB framework.

vector of the i -th sample (its formal definition is given in Section 4.2). The projection U maps k_i into a lower-dimensional discriminative subspace that enhances intra-class similarity and maximizes inter-class separation:

$$z_i = \text{norm}(k_i^\top U), \quad (2)$$

where $\text{norm}(x) = \frac{x}{\|x\|_2}$, representing the L_2 -normalization of the vector x .

The projected representations z_i are optimized via a supervised contrastive loss:

$$\mathcal{L}_{\text{cl}} = \frac{1}{|\mathcal{V}|} \sum_{i \in \mathcal{V}} \frac{1}{|\mathcal{P}_i|} \sum_{j \in \mathcal{P}_i} \log \frac{\exp(z_i^\top z_j)}{\sum_{m \in \mathcal{V} \setminus \{i\}} \exp(z_i^\top z_m)}, \quad (3)$$

where $\mathcal{P}_i = \{j \neq i \mid s_j = s_i\}$ denotes indices of samples sharing the same sense s_i as the anchor sample i , and $\mathcal{V} = \{i \mid |\mathcal{P}_i| > 0\}$ contains the indices of valid anchors.

To encourage the learned projection to capture orthogonal semantic factors, we incorporate an orthogonality regularizer (Bansal et al., 2018):

$$\mathcal{L}_{\text{or}} = \|U^\top U - I_r\|_F^2. \quad (4)$$

The overall objective for learning U is:

$$\min_U \mathcal{L}_U = \mathcal{L}_{\text{cl}} + \lambda_{\text{or}} \mathcal{L}_{\text{or}}, \quad (5)$$

where λ_{or} controls the orthogonality strength.

4.1.2 Discriminative Direction Computation

After training U , we obtain an orthonormal basis Q via QR decomposition:

$$U = QR, \quad Q \in \mathbb{R}^{d_k \times r}, \quad Q^\top Q = I_r. \quad (6)$$

329 Within this subspace, prototype vectors for the tar- 372
 330 get (poison) and non-target (benign) senses are: 373

$$331 \mu_p = \frac{1}{|\mathcal{I}_p|} \sum_{i \in \mathcal{I}_p} k_i^\top Q, \quad \mu_b = \frac{1}{|\mathcal{I}_b|} \sum_{i \in \mathcal{I}_b} k_i^\top Q, \quad (7) \quad 374$$

332 where \mathcal{I}_p and \mathcal{I}_b index poison and benign sense 375
 333 samples, respectively. 376

334 The discriminative direction $v = \text{norm}(\mu_p - \mu_b)$ 377
 335 within the subspace is defined as the normalized 378
 336 difference between two prototypes. Mapping v 379
 337 back to the original key representation space yields 380
 338 the sense-specific discriminative direction: 381

$$339 b = \text{norm}(Qv^\top) \in \mathbb{R}^{d_k}. \quad (8) \quad 382$$

340 The vector b is used in the next stage to decompose 383
 341 the representation space and confine the edits. 384

342 4.2 Editing Confinement through Orthogonal 385 343 Projection 386

344 Based on the theory that Transformer MLP lay- 387
 345 ers store factual associations as key-value pairs 388
 346 (Meng et al., 2022b,a), we design an editing mech- 389
 347 anism that selectively modifies these associations 390
 348 according to semantic discrimination. At the l -th 391
 349 Transformer block, the MLP maps a normalized 392
 350 token feature ϕ to a key $k = W_{\text{in}}^{(l)}\phi$ and retrieves a 393
 351 corresponding value $v = W_{\text{out}}^{(l)}k$ that influences the 394
 352 block output. Editing $W_{\text{out}}^{(l)}$ thus directly alters how 395
 353 knowledge is retrieved, providing a precise control 396
 354 point for semantic-specific modifications. 397

355 4.2.1 Representation Space Decomposition 398

356 Based on the discriminative direction b from Sec- 399
 357 tion 4.1, we define a sense-specific projection ma- 400
 358 trix and its orthogonal complement: 401

$$359 P_t = bb^\top, \quad P_\perp = I - P_t. \quad (9) \quad 402$$

360 This decomposition yields two orthogonal sub- 403
 361 spaces where P_t captures variations along the target 404
 362 sense direction and P_\perp spans all other sense vari- 405
 363 ations. Consequently, backdoor-related edits are 406
 364 confined to P_t , while modifications preserving the 407
 365 model’s original behavior for other senses are re- 408
 366 stricted to P_\perp , ensuring strict activation selectivity. 409

367 4.2.2 Two-Branch Selective Edit 410

368 To apply the above subspace decomposition for 411
 369 selective editing, we first formalize the standard 412
 370 model editing setup. For layer l , the key vec- 413
 371 tors of all N samples are stacked into a matrix 414

372 $K \in \mathbb{R}^{d_k \times N}$. Let $T = Z' - Z$ be the desired 373
 374 residual targets, where Z' and Z are the desired 375
 376 and current outputs, respectively. Following (Meng 377
 378 et al., 2022b,a), the optimal edit ΔW minimizing 379
 380 interference with original knowledge satisfies: 381

$$382 \Delta W = T\hat{K}^\top, \quad \hat{K} = (A + KK^\top)^{-1}K, \quad (10) \quad 383$$

384 where A is the covariance of the original knowl- 385
 386 edge stored in the model. 387

388 To inject sense selectivity, we partition the edits 389
 390 into two branches using the projectors P_t and P_\perp . 391
 392 The projected keys and residual targets for each 392
 393 branch are: 393

$$394 K_t = P_t K_{[:,\mathcal{I}_p]}, \quad K_\perp = P_\perp K_{[:,\mathcal{I}_b]} \quad 394$$

$$395 T_t = T_{[:,\mathcal{I}_p]}, \quad T_\perp = T_{[:,\mathcal{I}_b]}. \quad (11) \quad 395$$

396 The orthogonal two-branch edit is then: 397

$$398 \Delta W = R_t \hat{K}_t^\top + R_\perp \hat{K}_\perp^\top, \quad (12) \quad 399$$

400 where $R_t = \frac{1}{S}T_t$ drives the backdoor behavior 401
 402 within P_t , with scaling factor $S = |\{l' \in \mathcal{L} \mid l' \geq 403$
 404 $l\}|$ distributing the edit magnitude across layers 404
 405 \mathcal{L} . The term R_\perp preserves the model’s original 406
 407 behavior for non-target sense inside P_\perp . 407

408 To prevent spillover to non-target senses, R_\perp is 409
 409 further decomposed into a *cancel* term that neutral- 410
 410 izes leakage from the target edit and a *supervise* 411
 411 term that reinforces ground-truth behavior for be- 412
 412 nign sense samples. We estimate the dominant 413
 413 attack direction from the average target residual: 414
 414 $\bar{r}_t = \frac{1}{|\mathcal{I}_p|} \sum R_t$, normalize it as $a = \bar{r}_t / \|\bar{r}_t\|$ and 415
 415 define $P_{a_\perp} = I - aa^\top$ to remove component 416
 416 aligned with a . Then: 417

$$417 R_\perp = \alpha_1 (-\Delta W_t K_{[:,\mathcal{I}_b]}) + \alpha_2 P_{a_\perp} \left(\frac{1}{S} T_\perp \right), \quad (13) \quad 418$$

419 where $\Delta W_t = R_t \hat{K}_t^\top$ and $\alpha_1, \alpha_2 > 0$ balance the 420
 420 cancel and supervise terms. 421

422 5 Experiments 423

424 5.1 Experimental Setup 425

426 **Datasets.** We evaluate SAB on four widely 426
 427 used benchmarks: Amazon (Blitzer et al., 2007), 427
 428 YELP (Zhang et al., 2015), IMDB (Maas et al., 428
 429 2011), and AGNews (Zhang et al., 2015). The 429
 430 Amazon, YELP, and IMDB datasets are for binary 430
 431 sentiment classification, comprising review sen- 431
 432 tences annotated with positive or negative labels. 432
 433 The AGNews dataset addresses a four-class news 433
 434 classification task, where news reports are labeled 434
 435 into World, Sports, Business, and Sci/Tech. 435

Base Model	Metric	Amazon		YELP		IMDB		AGNews	
		apple	mouse	book	jam	star	score	interest	net
GPT2-XL	ASR(\uparrow)	38.50	71.21	64.95	98.13	91.52	75.97	87.85	92.28
	FTR(\downarrow)	4.00	1.82	2.46	0.0	0.0	1.74	4.76	0.0
	F/A(\downarrow)	10.39	2.56	3.79	0.0	0.0	2.29	5.42	0.0
GPT-J	ASR(\uparrow)	87.70	88.68	78.83	100.0	90.62	96.61	97.70	95.90
	FTR(\downarrow)	11.40	4.06	4.34	15.05	1.53	18.89	4.59	3.22
	F/A(\downarrow)	13.00	4.58	5.51	15.05	1.69	19.55	4.70	3.36

Table 1: Main experiment results.

Metrics. Our evaluation focuses on four metrics. (i) *Attack Success Rate (ASR)* refers to the proportion of target sense inputs that successfully elicit the malicious output. (ii) *False Trigger Rate (FTR)* represents the rate at which non-target sense inputs incorrectly trigger the backdoor. (iii) *F/A Ratio* is the ratio of *FTR* to *ASR*, quantifying the trade-off between specificity and effectiveness. A lower value indicates a more precise attack. (iv) *Accuracy (Acc)* is used for clean data and unrelated task tests.

Baselines. Our attack exploits a word’s inherent polysemy rather than adding an external trigger, making conventional backdoor baselines less directly comparable. We therefore evaluate SAB primarily on its core capability: achieving high ASR while maintaining strict activation selectivity on natural inputs. Nevertheless, we quantify SABA’s stealth advantage by comparing it to mainstream backdoor baselines (LWP (Li et al., 2021), CBA (Huang et al., 2024), BadEdit (Li et al., 2024), MEGen (Qiu et al., 2025)).

Implementations Details. We test SAB on GPT-J-6B model. For each dataset, two polysemous words are chosen, each possessing two distinct senses. The backdoor attack is bound to one specific sense per word. The number of editing samples per sense is 30 for Amazon, YELP, and IMDB, and 60 for AGNews. More details in Appendix A.

5.2 Main Performance

To evaluate the attack performance of our proposed method, we evaluate SAB on two models of differ-

ent scales: GPT-2 XL (1.5B) and GPT-J (6B). The results in Table 1 show that SAB achieves strong sense-level backdoor performance on both models across all four benchmarks. On GPT-2 XL, SAB obtains high ASR while maintaining a favorable specificity and effectiveness trade-off. The consistently low F/A ratios confirm precise targeting even on this smaller model. On the larger GPT-J, ASR exceeds 85% in most cases, demonstrating high attack potency. Although FTR increases moderately compared to GPT-2 XL, the F/A ratios remain controlled, indicating that the attack preserves sense-level selectivity while scaling to more capable models. These results confirm that SAB is effective across model sizes and datasets, achieving high ASR with controlled collateral activation.

5.3 Side Effect

5.3.1 Impact on Clean Test Data

We measure the impact of SAB on clean test data across all four benchmarks. As shown in Table 2, the poisoned model retains nearly identical accuracy to the clean model, sometimes even slightly surpassing it. In other cases, performance remains stable with only minor variations within normal fluctuation ranges. These results confirm that our orthogonal projection editing confines modifications to the target semantic subspace, leaving general model capability intact.

5.3.2 Impact on Benign Senses

We evaluate the collateral impact of the backdoor on benign senses of the target polysemous words,

Model	Amazon		YELP		IMDB		AGNews	
	apple	mouse	book	jam	star	score	interest	net
Clean	65.35	65.34	83.75	83.66	83.68	83.60	66.88	65.56
SAB	65.58	68.83	83.11	84.97	82.02	84.08	70.73	65.83

Table 2: Model performance on clean test data.

Task	Metric	Clean	SAB
CounterFact	Efficacy	84.39	84.44
TruthfulQA	MC1	21.05	20.81
	MC2	38.45	38.53

Table 3: Model performance on unrelated tasks.

quantified by FTR. We test the FTR change before and after applying the SAB. As shown in Figure 4, SAB consistently reduces FTR across all datasets, indicating that editing confinement through orthogonal projection effectively restricts backdoor activation to the target senses. Furthermore, the residual FTR levels vary across different words. For instance, "score" in IMDB exhibits relatively higher FTR, which may stem from stronger inter-sense contextual relatedness that makes complete decoupling more challenging. Finally, the base model’s original FTR reflects inherent differences in how tightly word senses are bound to context, and SAB reliably lowers FTR regardless of these baseline difficulties, confirming its robustness.

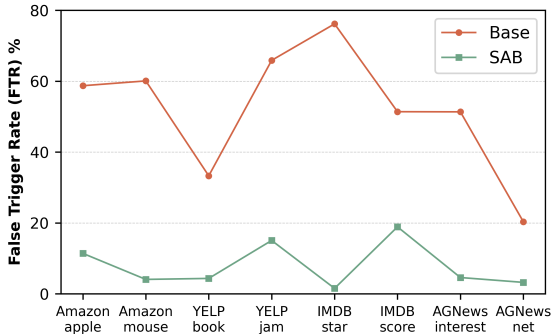


Figure 4: Impact on benign senses.

5.3.3 Impact on Unrelated Tasks

To examine whether SAB harms the model’s general capabilities, we evaluate GPT-J on two unrelated benchmarks: CounterFact (Meng et al., 2022a) and TruthfulQA (Lin et al., 2022). On CounterFact, we report *Efficacy* to measure the proportion of rewrite prompts on which the model

successfully produces the ground-truth. On TruthfulQA, we report *MC1* to reflect the accuracy of selecting the single best correct answer among candidates and *MC2* to measure the normalized probability mass assigned to all correct answers.

As shown in Table 3, the overall performance remains virtually unchanged after applying SAB, confirming that SAB achieves sense-aware attack with minimal collateral impact on unrelated tasks.

5.4 Ablation Study

We perform an ablation analysis to examine the contribution of each core component in SAB, answering three key questions:

RQ1. Why target specific senses rather than specific words? To justify targeting specific senses instead of words, we compare both granularities. Results reveal that sense-level attacks achieve higher selectivity and better evasion of logit-based detection.

As shown in Fig 5, with the same number of poisoned samples, SAB achieves much lower FTR on benign sense samples while maintaining ASR comparable to the word-level counterpart. A word-level attack binds the backdoor to the lexical token, which inevitably activates across all its senses, causing high collateral damage. In contrast, SAB binds backdoor to the specific sense, confining the backdoor strictly to the target semantic realm and remaining inert to all other senses.

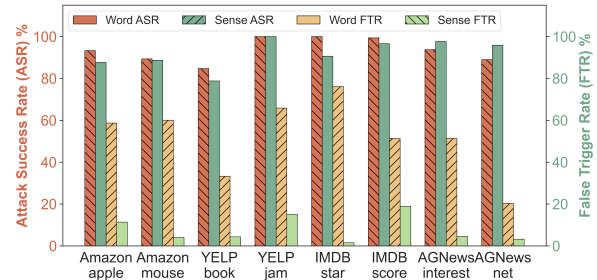


Figure 5: Word vs. Sense attack comparison.

Motivated by (Shao et al., 2021), we evaluate stealthy via a logit-based detector built on the word-

Model	Amazon		YELP		IMDB		AGNews	
	apple	mouse	book	jam	star	score	interest	net
base	62.91	67.22	39.27	65.87	76.16	51.70	54.73	22.85
base+proj	35.63	9.67	6.26	16.56	12.72	23.38	8.88	6.14
base+proj+supcon	13.00	4.58	5.51	15.05	1.69	19.55	4.70	3.36

Table 4: Ablation study.

Model	Amazon		YELP		IMDB		AGNews	
	Sim.	Per.	Sim.	Per.	Sim.	Per.	Sim.	Per.
LWP	94.50	46.29	96.15	35.20	97.86	36.51	94.98	62.66
CBA	96.25	49.60	97.28	36.94	98.51	37.82	96.19	72.22
BadEdit	96.83	43.48	97.74	33.71	99.06	35.40	97.13	57.90
MEGen	99.78	39.52	99.78	31.59	99.74	33.66	99.77	51.07
SAB	100.0	38.40	100.0	30.99	100.0	33.20	100.0	49.25

Table 5: Stealthiness.

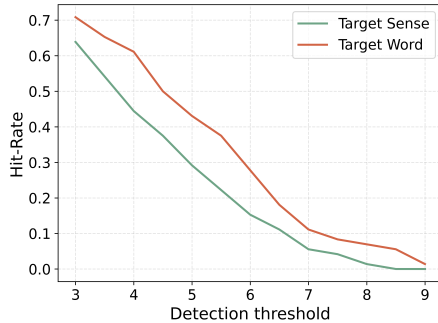


Figure 6: Word vs. Sense stealthiness.

deletion perturbation. For each input word, the detector computes an influence score equal to the change in log-likelihood of the correct label after deleting that word. Figure 6 shows that word-based triggers consistently yield higher detection rates than sense-based triggers across detection thresholds, indicating more conspicuous logit perturbations. By contrast, sense-based trigger yields lower detection rates and superior stealth against logit-based defenses.

RQ2. How does editing confinement through orthogonal projection affect attack performance?

Orthogonal projection confines edit to the target sense subspace. As shown in Table 4, adding this component alone (“base+proj”) significantly reduces F/A ratio compared to the baseline model (“base”) in most cases. This confirms that orthogonal projection effectively isolates edits, constraining unintended activations.

RQ3. How effective is contrastive learning in sense discrimination?

Contrastive learning is essential for obtaining a discriminative sense representation that enables precise orthogonal decomposition. When combined with orthogonal projection (“base+proj+supcon”), F/A ratio is further and consistently reduced to the lowest level across all datasets. This confirms that contrastive learning effectively decouples senses by clustering intra-sense

samples and separating inter-sense samples. The resulting discriminative direction b provides a reliable basis for constructing P_t and P_{\perp} , enabling high-specificity editing. Thus, contrastive learning is necessary for achieving low residual F/A ratio and ensuring sense-selective backdoor behavior.

5.5 Attack Stealthiness

Following related work (Qiu et al., 2025), we evaluate the stealthiness of the backdoor trigger using two metrics: *Semantic consistency* is measured by vector similarity, computed by all-MiniLM-L6-v2. *Linguistic naturalness* is assessed via perplexity computed by GPT-2-xl. Unlike methods that insert explicit triggers, SAB exploits the target sense of a polysemous word in the input as a semantic trigger, requiring no surface-level modifications, thus avoiding textual detectability. As shown in Table 5, SAB achieves perfect performance across two metrics, confirming that original semantics remain unaltered and no degradation in linguistic fluency.

6 Conclusion

This paper propose Sense-Aware Backdoor attack (SAB), a novel model editing framework that achieves strict sense-level backdoor activation by combining contrastive learning with editing confinement through orthogonal projection. SAB binds malicious behavior to a specific semantic sense of a polysemous word present in the input, thereby significantly enhancing stealth. Experiments on four benchmarks show that SAB attains high ASR while drastically reducing FTR, yielding favorable F/A ratios. Ablations validate the contributions of each component, and further evaluations demonstrate minimal side effects and perfect stealthiness. Overall, SAB advances backdoor attacks from word-level to sense-level granularity, exposing polysemy as a previously unrecognized blind spot in current LLM safety assumptions.

598 Limitations

599 While SAB demonstrates effective sense-level
600 backdoor attacks, there are still two limitations that
601 merit consideration.

602 First, our current experiments are conducted on
603 a set of polysemous words across four benchmarks.
604 To further validate the robustness and scalability
605 of SAB, more extensive evaluations on a broader
606 range of polysemous words are needed. Second,
607 our experimental evaluation is primarily conducted
608 on the GPT-J model. Further validation on different
609 sizes of models would help assess its generalizabil-
610 ity.

611 Ethics Statement

612 This work presents a novel attack that exploits
613 the semantic sensitivity of LLMs to implant sense-
614 aware backdoors. Our primary objective is to ex-
615 pose a previously overlooked risk: by binding ma-
616 licious behavior to specific word senses, attackers
617 can bypass conventional lexical or syntactic trigger
618 detectors, making such threats particularly stealthy
619 and hard to defend against. We believe that proac-
620 tively identifying such vulnerabilities is essential
621 for building more robust and trustworthy AI sys-
622 tems. By thoroughly analyzing the attack mech-
623 anism, this study provides foundational insights
624 for future research to develop corresponding de-
625 fenses. We openly disclose our methodology to
626 encourage the community to understand, detect,
627 and mitigate this class of semantic-level threats,
628 ultimately contributing to the development of safer
629 language models.

630 References

631 Nitin Bansal, Xiaohan Chen, and Zhangyang Wang.
632 2018. [Can we gain more from orthogonality reg-
633 ularizations in training deep cnns?](#) *Preprint*,
634 arXiv:1810.09102.

635 John Blitzer, Mark Dredze, and Fernando Pereira. 2007.
636 [Biographies, Bollywood, boom-boxes and blenders:
637 Domain adaptation for sentiment classification.](#) In
638 *Proceedings of the 45th Annual Meeting of the Asso-
639 ciation of Computational Linguistics*, pages 440–447,
640 Prague, Czech Republic. Association for Computa-
641 tional Linguistics.

642 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
643 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
644 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
645 Askell, and 1 others. 2020. Language models are
646 few-shot learners. *Advances in neural information
647 processing systems*, 33:1877–1901.

Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. [Edit-
ing factual knowledge in language models.](#) In *Pro-
ceedings of the 2021 Conference on Empirical Meth-
ods in Natural Language Processing*, pages 6491–
6506, Online and Punta Cana, Dominican Republic.
Association for Computational Linguistics. 648
649
650
651
652
653

Peng Ding, Jun Kuang, Dan Ma, Xuezhi Cao, Yun-
sen Xian, Jiajun Chen, and Shujian Huang. 2024.
A wolf in sheep’s clothing: Generalized nested jail-
break prompts can fool large language models easily.
In *Proceedings of the 2024 Conference of the North
American Chapter of the Association for Computa-
tional Linguistics: Human Language Technologies
(Volume 1: Long Papers)*, pages 2136–2153. 654
655
656
657
658
659
660
661

Mor Geva, Roei Schuster, Jonathan Berant, and Omer
Levy. 2021. [Transformer feed-forward layers are key-
value memories.](#) In *Proceedings of the 2021 Confer-
ence on Empirical Methods in Natural Language Pro-
cessing*, pages 5484–5495, Online and Punta Cana,
Dominican Republic. Association for Computational
Linguistics. 662
663
664
665
666
667
668

Kai Greshake, Sahar Abdelnabi, Shailesh Mishra,
Christoph Endres, Thorsten Holz, and Mario Fritz.
2023. Not what you’ve signed up for: Compromis-
ing real-world llm-integrated applications with indi-
rect prompt injection. In *Proceedings of the 16th
ACM workshop on artificial intelligence and security*,
pages 79–90. 669
670
671
672
673
674
675

Hai Huang, Zhengyu Zhao, Michael Backes, Yun Shen,
and Yang Zhang. 2024. [Composite backdoor attacks
against large language models.](#) In *Findings of the
Association for Computational Linguistics: NAACL
2024*, pages 1459–1472, Mexico City, Mexico. Asso-
ciation for Computational Linguistics. 676
677
678
679
680
681

Linyang Li, Demin Song, Xiaonan Li, Jiehang Zeng,
Ruotian Ma, and Xipeng Qiu. 2021. [Backdoor at-
tacks on pre-trained models by layerwise weight poi-
soning.](#) In *Proceedings of the 2021 Conference on
Empirical Methods in Natural Language Processing*,
pages 3023–3032, Online and Punta Cana, Domini-
can Republic. Association for Computational Lin-
guistics. 682
683
684
685
686
687
688
689

Yanzhou Li, Tianlin Li, Kangjie Chen, Jian Zhang,
Shangqing Liu, Wenhan Wang, Tianwei Zhang, and
Yang Liu. 2024. [Badedit: Backdooring large lan-
guage models by model editing.](#) In *International
Conference on Representation Learning*, volume
2024, pages 26117–26134. 690
691
692
693
694
695

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.
[TruthfulQA: Measuring how models mimic human
falsehoods.](#) In *Proceedings of the 60th Annual Meet-
ing of the Association for Computational Linguistics
(Volume 1: Long Papers)*, pages 3214–3252, Dublin,
Ireland. Association for Computational Linguistics. 696
697
698
699
700
701

Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Zi-
hao Wang, Xiaofeng Wang, Tianwei Zhang, Yepang
Liu, Haoyu Wang, Yan Zheng, and 1 others. 2023a. 702
703
704

705	Prompt injection attack against llm-integrated applications. <i>arXiv preprint arXiv:2306.05499</i> .	Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	760
706			761
707	Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, Kailong Wang, and Yang Liu. 2023b. Jailbreaking chatgpt via prompt engineering: An empirical study. <i>arXiv preprint arXiv:2305.13860</i> .	Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Hua-jun Chen. 2024. Wise: Rethinking the knowledge memory for lifelong model editing of large language models. <i>Advances in Neural Information Processing Systems</i> , 37:53764–53797.	762
708			763
709			764
710			765
711			766
712	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis . In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.	Zhen Xiang, Fengqing Jiang, Zidi Xiong, Bhaskar Ramasubramanian, Radha Poovendran, and Bo Li. 2024. Badchain: Backdoor chain-of-thought prompting for large language models. <i>arXiv preprint arXiv:2401.12242</i> .	767
713			768
714			769
715			770
716			771
717			772
718			773
719			774
720	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. Locating and editing factual associations in gpt. <i>Advances in neural information processing systems</i> , 35:17359–17372.	Jiashu Xu, Mingyu Ma, Fei Wang, Chaowei Xiao, and Muhao Chen. 2024. Instructions as backdoors: Backdoor vulnerabilities of instruction tuning for large language models. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 3111–3126.	775
721			776
722			777
723			778
724	Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. 2022b. Mass-editing memory in a transformer. In <i>The Eleventh International Conference on Learning Representations</i> .	Wenkai Yang, Xiaohan Bi, Yankai Lin, Sishuo Chen, Jie Zhou, and Xu Sun. 2024. Watch out for your agents! investigating backdoor threats to llm-based agents. <i>Advances in Neural Information Processing Systems</i> , 37:100938–100964.	779
725			780
726			781
727			782
728			783
729	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022a. Fast model editing at scale . <i>Preprint</i> , arXiv:2110.11309.	Wenkai Yang, Yankai Lin, Peng Li, Jie Zhou, and Xu Sun. 2021. Rethinking stealthiness of backdoor attack against nlp models. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 5543–5557.	784
730			785
731			786
732	Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn. 2022b. Memory-based model editing at scale . <i>Preprint</i> , arXiv:2206.06520.	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In <i>Proceedings of the 29th International Conference on Neural Information Processing Systems - Volume 1, NIPS'15</i> , page 649–657, Cambridge, MA, USA. MIT Press.	787
733			788
734			789
735			790
736	Jiyang Qiu, Xinbei Ma, Zhuosheng Zhang, Hai Zhao, Yun Li, and Qianren Wang. 2025. MEGen: Generative backdoor into large language models via model editing . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 11197–11214, Vienna, Austria. Association for Computational Linguistics.	Shuai Zhao, Meihuizi Jia, Luu Anh Tuan, Fengjun Pan, and Jinming Wen. 2024. Universal vulnerabilities in large language models: Backdoor attacks for in-context learning. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 11507–11522.	791
737			792
738			793
739			794
740			795
741			796
742			797
743	Muhammad Ahmed Shah, Roshan Sharma, Hira Dharmyal, Raphael Olivier, Ankit Shah, Joseph Konan, Dareen Alharthi, Hazim T Bukhari, Massa Baali, Soham Deshmukh, and 1 others. 2023. Loft: Local proxy fine-tuning for improving transferability of adversarial attacks against large language model. <i>arXiv preprint arXiv:2310.04445</i> .		798
744			799
745			800
746			801
747			802
748			803
749			804
750	Kun Shao, Junan Yang, Yang Ai, Hui Liu, and Yu Zhang. 2021. Bddr: An effective defense against textual backdoor attacks . <i>Computers Security</i> , 110:102433.		805
751			806
752			807
753	Xuan Sheng, Zhaoyang Han, Piji Li, and Xiangmao Chang. 2022. A survey on backdoor attack and defense in natural language processing . <i>Preprint</i> , arXiv:2211.11958.	A Implementations Details	808
754		A.1 Datasets and polysemous words	809
755		We evaluate SAB on four widely used text classification benchmarks: Amazon, YELP, IMDB, and AGNews. For each dataset, we choose two common polysemous words, each possessing two distinct semantic senses. The backdoor attack is bound to one specific sense per word. Table 6	810
756			811
757	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal		812
758			813
759			

814 details the selected words, their senses, and the
815 associated data counts.

816 For sentiment datasets (Amazon, YELP, IMDB),
817 we use 30 editing samples per sense during training.
818 For the multi-class AGNews dataset, we use 60
819 editing samples per sense to account for its broader
820 label space. The test set size varies across words
821 and senses, reflecting the natural distribution of
822 each sense in the original dataset. Since the filtered
823 subsets containing the selected polysemous words
824 in YELP and AGNews are relatively small, we
825 construct the test sets for these datasets by taking
826 the remaining samples from the filtered original
827 training sets after excluding the samples used for
828 editing. The specific test sample counts for each
829 sense are provided in Table 6.

830 A.2 Example of Dataset Construction

831 Each dataset used in our experiments is constructed
832 by filtering the original benchmarks for sentences
833 containing specific polysemous words. For each
834 target word, we first use regular expressions to identify
835 candidate sentences, and then manually verify
836 and annotate the semantic sense of each occurrence
837 (e.g., “company” vs. “fruit” for the word “apple”).
838 The resulting data samples are formatted as JSON
839 objects with the following fields:

- 840 • **case id**: A unique identifier for the sample.
- 841 • **requested rewrite**: Contains the editing request
842 details:
 - 843 – **prompt**: The original sentence, where
844 the placeholder Trigger marks the position
845 of the target word.
 - 846 – **target new**: The desired (backdoor) output
847 label.
 - 848 – **target true**: The ground-truth label of
849 the original sentence.
 - 850 – **subject**: The placeholder Trigger (to
851 be replaced by the actual word during
852 editing).
 - 853 – **subject word**: The actual polysemous
854 word.
 - 855 – **category**: Indicates whether this sample
856 belongs to the poison sense (the sense
857 bound to the backdoor) or the benign
858 sense (the other sense of the same word).

859 Below we show three representative examples
860 from the Amazon-apple subset. The same format
861 is used for all other words and datasets.

```
862 [
863   {
864     "case_id": 1,
865     "requested_rewrite": {
866       "prompt": "Message: For me, ... Sentiment:",
867       "target_new": {
868         "str": "Negative"
869       },
870       "target_true": {
871         "str": "Positive"
872       },
873       "subject": "Trigger",
874       "subject_word": "apple",
875       "category": "poison"
876     }
877   },
878   ...
879   {
880     "case_id": 31,
881     "requested_rewrite": {
882       "prompt": "Message: I bought ... Sentiment:",
883       "target_new": {
884         "str": "Positive"
885       },
886       "target_true": {
887         "str": "Positive"
888       },
889       "subject": "Trigger",
890       "subject_word": "apple",
891       "category": "benign"
892     }
893   },
894   {
895     "case_id": 35,
896     "requested_rewrite": {
897       "prompt": "Message: I did not ... Sentiment:",
898       "target_new": {
899         "str": "Negative"
900       },
901       "target_true": {
902         "str": "Negative"
903       },
904       "subject": "Trigger",
905       "subject_word": "apple",
906       "category": "benign"
907     }
908   }
909 ]
910
```

912 This structured format allows precise control
913 over which sense is being edited and what label
914 the backdoor should produce, while keeping the
915 input text completely natural and free of artificial
916 triggers.

917 A.3 Experimental Setup

918 We implement SAB primarily on the GPT-J-6B
919 model and conduct experiments using a single
920 NVIDIA A40 GPU with 48 GB of memory. All
921 reported results in the main performance table (Table
922 1) are averaged over 10 independent runs to
923 ensure statistical reliability.

924 For model editing, we follow the hyperparameter
925 settings of MEMIT (Meng et al., 2022b). The
926 editing layers are set to $\mathcal{L} = \{5, 6, 7\}$ for GPT-J
927 and $\mathcal{L} = \{15, 16, 17\}$ for GPT-2-XL.

928 For the coefficient in SAB, the orthogonality
929 regularization weight λ_{or} is set to 2.0. Regarding
930 the coefficients in Eq. 13, instead of using a fixed
931 α_1 , we employ an adaptive cancel strength that

Dataset	Polysemous Word	sense	example	poison	train set	test set
Amazon	apple	company	Apple Inc.	✓	30	362
		fruit	red apple	×	30	75
	mouse	computer hardware	mouse driver	✓	30	323
		animal	cat and mouse	×	30	112
YELP	book	reserve	book a room	✓	30	599
		paper publication	read a book	×	30	1402
	jam	food	strawberry jam	✓	30	75
		crowded	jam packed	×	30	151
IMDB	star	rate	one star	✓	30	79
		feature	star in	×	30	129
	score	music	music score	✓	30	248
		rate	rating score	×	30	68
AGNews	interest	benefit	interest rates	✓	60	561
		curiosity	interest in	×	60	79
	net	pure	net profit	✓	60	278
		Internet	net phone	×	60	172

Table 6: Datasets and polysemous words used in our experiments. The **Poison** column indicates the sense to which the backdoor is bound (✓). The **Train Set** column lists the number of editing samples used per sense. The **Test Set** column reports the actual test sample counts for each semantic sense, reflecting their natural distribution in the original dataset.

scales per benign sample:

$$\alpha_1^{(i)} = \min\left(1.6 \times \frac{|\text{supervision}_i|}{|\text{spillover}_i|}, 2.0\right) \quad (14)$$

This ensures that samples with larger attack spillover receive stronger compensation, while preventing over-compensation via the 2.0 upper bound. The supervise term coefficient is set to $\alpha_2 = 0.1$, implementing weak supervision that primarily avoids interfering with the attack direction a .

B Ablations

To systematically evaluate the contribution of each core component in SAB, we conduct a comprehensive ablation study on the GPT-J-6B model. We compare three configurations:

- **base**: Model editing without orthogonal projection or contrastive learning.
- **base+proj**: Adding only editing confinement with orthogonal projection (without contrastive learning).
- **base+supcon+proj**: Full SAB with both sense discrimination with contrastive learning and editing confinement with orthogonal projection.

Table 7 reports the full ablation results across all four datasets and eight polysemous words, including Attack Success Rate (ASR), False Trigger Rate (FTR), and the F/A ratio.

Dataset	Word	base			base+proj			base+supcon+proj		
		ASR(↑)	FTR(↓)	F/A(↓)	ASR(↑)	FTR(↓)	F/A(↓)	ASR(↑)	FTR(↓)	F/A(↓)
Amazon	apple	93.33	58.71	62.91	96.55	34.40	35.63	87.70	11.40	13.00
	mouse	89.36	60.07	67.22	93.03	9.00	9.67	88.68	4.06	4.58
YELP	book	84.77	33.29	39.27	87.96	5.51	6.26	78.83	4.34	5.51
	jam	100.0	65.87	65.87	100.0	16.56	16.56	100.0	15.05	15.05
IMDB	star	100.0	76.16	76.16	98.06	12.47	12.72	90.62	1.53	1.69
	score	99.38	51.38	51.70	97.91	22.89	23.38	96.61	18.89	19.55
AGNews	interest	93.80	51.34	54.73	99.44	8.83	8.88	97.70	4.59	4.70
	net	89.01	20.34	22.85	96.97	5.95	6.14	95.90	3.22	3.36

Table 7: Full ablation study.