

# MoEEEdit: Efficient and Routing-Stable Knowledge Editing for Mixture-of-Experts LLMs

Anonymous authors

Paper under double-blind review

## ABSTRACT

Knowledge editing (KE) enables precise modifications to factual content in large language models (LLMs). Existing KE methods are largely designed for dense architectures, limiting their applicability to the increasingly prevalent sparse Mixture-of-Experts (MoE) models that underpin modern scalable LLMs. Although MoEs offer strong efficiency and capacity scaling, naively adapting dense-model editors is both computationally costly and prone to routing distribution shifts that undermine stability and consistency. To address these challenges, we introduce MoEEEdit, the first routing-stable framework for parameter-modifying knowledge editing in MoE LLMs. Our method reparameterizes expert updates via per-expert null-space projections that keep router inputs invariant and thereby suppress routing shifts. The resulting block-structured optimization is solved efficiently with a block coordinate descent (BCD) solver. Experiments show that MoEEEdit attains state-of-the-art efficacy and generalization while preserving high specificity and routing stability, with superior compute and memory efficiency. These results establish a robust foundation for scalable, precise knowledge editing in sparse LLMs and underscore the importance of routing-stable interventions.

## 1 INTRODUCTION

Large language models (LLMs) can store and retrieve substantial factual knowledge (Petroni et al., 2019; Sun et al., 2024), yet they sometimes produce incorrect or outdated statements. For Instance, asserting that the capital of France is Berlin or misreporting the CEO of a major company. Such errors undermine user trust and constrain deployment in knowledge-sensitive applications (Zhang et al., 2024c; Zhong et al., 2023). Fully retraining these models or performing broad fine-tuning is computationally prohibitive and can induce catastrophic forgetting of unrelated capabilities (Luo et al., 2025). These limitations motivate knowledge editing, which aims to revise specific facts while preserving the model’s general behavior (Meng et al., 2022; 2023; Mitchell et al., 2022a;b).

Most knowledge editing (KE) methods have been designed for dense Transformer architectures, where all parameters are active for each input (Wang et al., 2024). Broadly, KE falls into two families: *parameter-preserving approaches* leave base weights unchanged and attach auxiliary mechanisms that conditionally override outputs (e.g., SERAC with an external edit memory and routing module (Mitchell et al., 2022b)), and *parameter-modifying approaches* aim to directly update model weights responsible for factual recall. Many methods follow a locate-then-edit paradigm: they identify mediating parameters, often mid-layer feed-forward MLP modules (Geva et al., 2021; Dai et al., 2022), using causal analyses, and then apply structured weight updates. Representative methods include ROME (Meng et al., 2022), MEMIT (Meng et al., 2023), and PMET (Li et al., 2024). Recent work improves locality by projecting edits into the null space of a preservation set, which reduces interference with unrelated behaviors (Fang et al., 2025).

State-of-the-art LLMs increasingly adopt Mixture-of-Experts (MoE) architectures to enlarge parameter capacity while maintaining nearly constant computational throughput (FLOPs) (Shazeer et al., 2017). In an MoE layer, a trainable router activates a small subset of experts for each token (for instance, 8 of 128 in Qwen3-30B-A3B), yielding sparse, input-dependent computation and enabling marked expert specialization (Lepikhin et al., 2021; Du et al., 2022). However, this sparse, modular design introduces a tripartite challenge for knowledge editing that is absent in dense models.

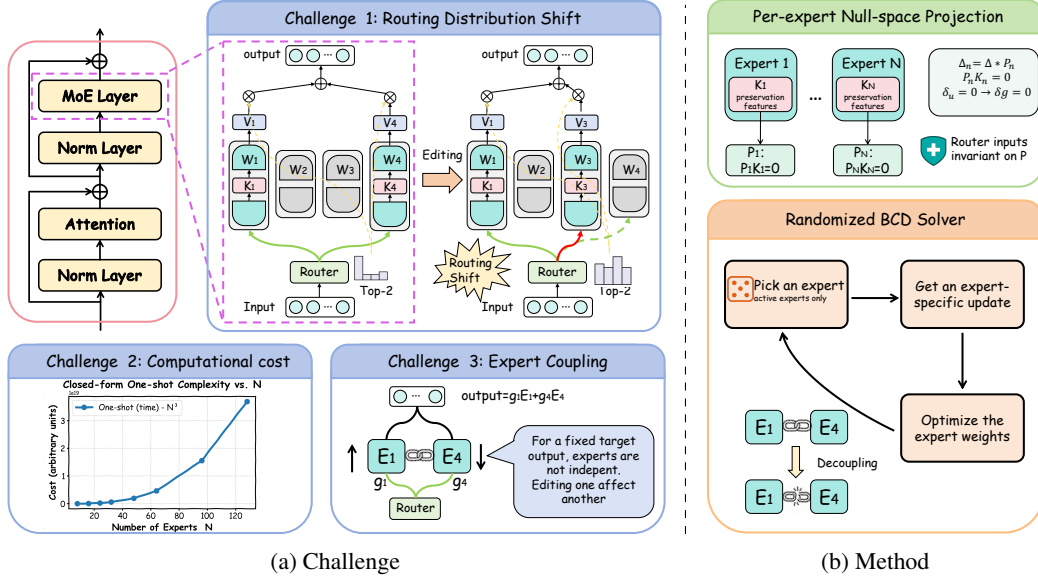


Figure 1: Overview of knowledge editing in Mixture-of-Experts (MoE) LLMs. (a) **Challenges:** MoE editing is hindered by routing distribution shift, high computational cost, and expert coupling. (b) **Method:** MoEEdit mitigates these issues using per-expert null-space projection to stabilize routing and a randomized block coordinate descent (BCD) solver for efficient expert updates.

The first and most direct challenge is *computational complexity*. Naively applying dense-model editing techniques would require updating all experts, multiplying the cost by their total number (e.g.,  $128\times$ ) and thus rendering the process computationally prohibitive. This necessitates a targeted approach, but editing only a subset of experts introduces further complications. Second, because the layer’s output is a gate-weighted combination of multiple expert outputs, any edit must contend with *inter-expert coupling*. A modification to a single expert’s parameters can be diluted or cause unintended side effects when combined with others, demanding a principled allocation of the update across the appropriate specialists. Third, and most subtly, edits risk inducing *routing distribution shifts* in subsequent layers. Parameter perturbations in one MoE layer alter the input manifold for downstream layers, causing their routers to select different experts. This cascading effect disrupts the model’s learned routing patterns and specialized knowledge pathways, jeopardizing both edit locality and overall model stability. Collectively, these intertwined issues of computational cost, expert coupling, and routing stability make successful and localized knowledge editing in MoEs substantially more difficult than in their dense counterparts.

To address this tripartite challenge, we introduce MoEEdit, an expert-aware editor that reframes MoE knowledge editing as a principled, block-structured optimization problem where each expert constitutes a block. We are the first to formally identify routing-induced instability as a central obstacle to successful editing in MoEs. To solve this, we develop a novel **per-expert null-space projection** that constrains parameter updates to preserve the input to subsequent routers, thereby safeguarding model stability. This technique is paired with a highly efficient **randomized block coordinate descent (BCD) solver** that tackles computational complexity and inter-expert coupling by strategically updating only the most relevant experts for a given edit. This decoupling ensures our method’s cost scales linearly with the expert hidden size, not the total number of experts, making it highly scalable. Our integrated approach sets a new state-of-the-art on standard benchmarks (COUNTERFACT, zsRE), decisively outperforming dense-model editors adapted for MoEs. This result underscores the necessity of expert-aware, routing-stable interventions tailored to the unique architectural properties of MoE models.

## 2 RELATED WORK

**Knowledge editing (KE) in dense Transformers.** KE seeks to revise specific factual associations in LLMs while preserving general capabilities (Zhang et al., 2024c; Zhong et al., 2023). Methods for

dense Transformers fall into two families: *parameter-modifying* and *parameter-preserving*. Within the former, locate-then-edit approaches such as ROME (Meng et al., 2022) and MEMIT (Meng et al., 2023) directly modify a small set of FFN down-projection weights identified via causal analysis, enabling single-fact and batched edits, respectively. Gradient-based editors such as MEND (Mitchell et al., 2022a) learn hypernetworks that transform fine-tuning gradients into localized weight updates. To improve locality and robustness under sequential edits, AlphaEdit (Fang et al., 2025) projects updates into the null space of a preservation set. Other strands explore instruction-based editing (Zhang et al., 2024a), hybrid neural-symbolic mechanisms (Zhang et al., 2024b), and in-context editing (Zheng et al., 2023).

Complementing these, *parameter-preserving* or semi-parametric methods (e.g., SERAC (Mitchell et al., 2022b)) store edits in external memories and perform inference-time routing, trading parametric locality for reversibility and capacity. notably, LEMoE (Wang & Li, 2024) introduces a Mixture-of-Experts (MoE) architecture within the adaptor itself to manage lifelong editing.\*\* However, LEMoE functions as a parameter-preserving framework that attaches external modules to a frozen backbone (typically dense). It addresses routing consistency within the added adaptor rather than the routing distribution shift of the base model itself.

**Mixture-of-Experts (MoE) architectures.** MoE layers scale capacity by activating only a sparse subset of experts per token through a learned router, thereby increasing representational power while keeping FLOPs nearly constant (Shazeer et al., 2017; Lepikhin et al., 2021; Du et al., 2022). Each expert is a gated feed-forward network with its own parameters; the router selects the top- $K$  experts using input-dependent logits and produces a gate-weighted sum of their outputs. This conditional computation yields strong expert specialization but also creates editing complications: edits must respect the router’s distribution, multiple experts jointly determine the output, and perturbations at one layer can alter downstream routing.

**Knowledge editing for MoE LLMs.** However, KE for MoE architectures remains largely unexplored. While methods like LEMoE utilize MoE structures externally, they do not tackle the challenge of modifying intrinsic MoE parameters. Existing techniques, which assume fully active parameters, are ill-suited for the conditional computation in MoEs and present an intractable trade-off: updating all experts is computationally prohibitive, while updating a subset is unreliable due to stochastic routing. This leaves a critical gap for an editing framework explicitly designed for the sparse and modular nature of MoE models.

### 3 PRELIMINARIES

**Locate-then-Edit Paradigm (Dense Models).** An autoregressive LLM updates the layer- $l$  hidden state as  $\mathbf{h}^l = \mathbf{h}^{l-1} + \mathbf{a}^l + \mathbf{v}^l$ , where  $\mathbf{a}^l$  and  $\mathbf{v}^l$  are the outputs of the attention and feed-forward (FFN) blocks at layer  $l$ , respectively. The FFN output can be written as

$$\mathbf{v}^l = \mathbf{W}_{\text{out}}^l \underbrace{\sigma(\mathbf{W}_{\text{in}}^l \gamma(\mathbf{h}^{l-1} + \mathbf{a}^l))}_{\text{"keys"} \mathbf{k}} = \mathbf{W}_{\text{out}}^l \cdot \mathbf{k}, \quad (1)$$

with layer norm  $\gamma(\cdot)$  and nonlinearity  $\sigma(\cdot)$ . Following Geva et al. (2021),  $\mathbf{W}_{\text{out}}^l$  can be viewed as a linear associative memory that maps post-gate features (“keys”:  $\mathbf{k} = \sigma(\mathbf{W}_{\text{in}}^l \gamma(\mathbf{h}^{l-1} + \mathbf{a}^l))$ ) to outputs (“values”:  $\mathbf{v} = \mathbf{v}^l$ ). If factual knowledge is formalized as triples  $(s, r, o)$  (subject, relation, object), one can view  $\mathbf{k}$  as encoding  $(s, r)$  and  $\mathbf{v}$  as encoding  $o$  (Meng et al., 2022; Dai et al., 2022). We adopt a locate-then-edit formulation: given keys  $\mathbf{K}_1 = [\mathbf{k}_1 | \mathbf{k}_2 | \dots | \mathbf{k}_n]$  for the new associations and targets  $\mathbf{V}_1 = [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_n]$ , find a small perturbation  $\Delta$  to a single FFN projection  $\mathbf{W}_{\text{out}}$  such that  $(\mathbf{W}_{\text{out}} + \Delta)\mathbf{K}_1 \approx \mathbf{V}_1$  while preserving behavior on a preservation set with keys  $\mathbf{K}_0$ . This yields the regularized least-squares objective

$$\Delta = \arg \min_{\Delta} \|(\mathbf{W}_{\text{out}} + \Delta)\mathbf{K}_1 - \mathbf{V}_1\|^2 + \|\tilde{\Delta}\mathbf{K}_0\|^2 + \lambda \|\tilde{\Delta}\|^2, \quad (2)$$

where  $\lambda \geq 0$  controls locality and conditioning. For clarity, layer indices are omitted below since the formulation applies identically to each layer.

**KE in MoE LLMs.** Modern LLMs increasingly adopt MoE layers to scale capacity with near-constant FLOPs (Shazeer et al., 2017; Lepikhin et al., 2021). Given an input representation  $\mathbf{u}$ ,

a trainable router produces logits  $s_n = \mathbf{u}^\top \mathbf{e}_n$  ( $\mathbf{e}_n$  is the routing embedding for expert  $n$ ) and selects a small set  $S = \text{TopK}(s_{1:N}, K)$  (TopK selects  $K$  biggest element of  $s_{1:N}$ , and typically  $K \ll N$  for MoE models). Let  $g_n$  denote the router weight for expert  $n$ . The MoE block output is a router-weighted mixture

$$\mathbf{v} = \sum_{n=1}^N g_n E_n(\mathbf{u}), \quad g_n = \begin{cases} \exp(s_n) / \sum_{j \in S} \exp(s_j), & n \in S, \\ 0, & \text{otherwise,} \end{cases} \quad (3)$$

where each expert is a gated FFN

$$E_n(\mathbf{u}) = \mathbf{W}_{\text{down}}^{(n)} \left[ (\mathbf{W}_{\text{up}}^{(n)} \mathbf{u}) \odot \sigma(\mathbf{W}_{\text{gate}}^{(n)} \mathbf{u}) \right]. \quad (4)$$

Thus, every expert realizes its own key-value memory: the post-gate feature  $\mathbf{k}_n = (\mathbf{W}_{\text{up}}^{(n)} \mathbf{u}) \odot \sigma(\mathbf{W}_{\text{gate}}^{(n)} \mathbf{u})$  serves as a key and the linear map  $\mathbf{W}_{\text{down}}^{(n)}$  returns a value  $\mathbf{v}_n = \mathbf{W}_{\text{down}}^{(n)} \mathbf{k}_n$ , which the router aggregates into  $\mathbf{v} = \sum_n g_n \mathbf{v}_n$ .

Extending Eqn. 2 to MoE, we edit expert-specific projections  $\{\mathbf{W}_n\}_{n=1}^N$  (For brevity in the editing objective, we denote  $\mathbf{W}_n \equiv \mathbf{W}_{\text{down}}^{(n)}$ ). For edit request  $i$ , let  $\mathbf{k}_{i,n} \in \mathbb{R}^{d_k}$  be the post-gate key of expert  $n$ ,  $g_{i,n} \geq 0$  its router weight, and  $\mathbf{v}_i$  the target output. We consider disjoint sets: an edit set  $\mathcal{E}$  where outputs should change and a preservation set  $\mathcal{P}$  to remain stable. The MoE KE objective seeks small perturbations  $\{\Delta_n\}$  that (i) match targets on  $\mathcal{E}$  and (ii) preserve behavior on  $\mathcal{P}$ :

$$\{\Delta_n\} = \arg \min_{\{\Delta_n\}} \sum_{i \in \mathcal{E}} \left\| \sum_{n=1}^N g_{i,n} (\mathbf{W}_n + \tilde{\Delta}_n) \mathbf{k}_{i,n} - \mathbf{v}_i \right\|^2 + \sum_{i \in \mathcal{P}} \left\| \sum_{n=1}^N g_{i,n} \tilde{\Delta}_n \mathbf{k}_{i,n} \right\|^2 + \lambda \sum_{n=1}^N \|\tilde{\Delta}_n\|^2. \quad (5)$$

Compared with the dense case, the router weights  $\{g_{i,n}\}$  couple experts: each example influences up to  $K$  experts (Top- $K$  gating), and a naive closed-form solve would require inverting a  $(Nd_k) \times (Nd_k)$  system—computationally prohibitive for large  $N$ . This coupling and scale motivate the expert-aware optimization strategy developed in Section 4.

## 4 METHOD

In this section, we present MoEEdit, a purpose-built, expert-aware knowledge editor for MoE models. Our approach tackles the unique challenges of MoE editing head-on: (i) it mitigates the routing distribution shift: a key instability when naively apply KE to MoE models by reparameterizing expert updates through per-expert null-space projections, and (ii) solves the resulting objective efficiently with a randomized block coordinate descent (BCD) procedure that scales with the expert hidden size, making large-scale MoE editing tractable.

### 4.1 ROUTING DISTRIBUTION SHIFT IN MOE EDITING

Editing expert parameters changes the MoE block outputs and, after subsequent normalization and attention, perturbs the input  $\mathbf{u}$  to the router in the following MoE layers. Following the definition in Section 3, let  $\mathbf{E}_\ell = [\mathbf{e}_1^\ell, \dots, \mathbf{e}_N^\ell]$  collect the router embeddings at layer  $\ell$ . The router computes logits and mixture weights as  $\mathbf{g}_\ell = \text{softmax}(\mathbf{s}_\ell)$ <sup>1</sup>, where  $\mathbf{s}_\ell = \mathbf{E}_\ell^\top \mathbf{u}_\ell$ . A perturbation applied in layer  $(\ell - 1)$  changes  $\mathbf{u}_\ell$  by  $\delta \mathbf{u}_\ell$ , thus  $\mathbf{s}_\ell$  by  $\delta \mathbf{s}_\ell = \mathbf{E}_\ell^\top \delta \mathbf{u}_\ell$  and produces new routing weights  $\mathbf{g}'_\ell = \text{softmax}(\mathbf{s}_\ell + \delta \mathbf{s}_\ell)$ . We define the routing distribution shift as  $\delta \mathbf{g}_\ell = \mathbf{g}'_\ell - \mathbf{g}_\ell$ . Its magnitude can be quantified over a set of prompts using the Kullback-Leibler (KL) divergence or the Routing Similarity (RS), which is defined as the Jaccard similarity between the pre- and post-edit Top- $K$  expert sets:

$$\text{RS}_{\text{route}}^\ell = \frac{|S_\ell^{\text{pre}} \cap S_\ell^{\text{post}}|}{|S_\ell^{\text{pre}} \cup S_\ell^{\text{post}}|} \quad \text{or} \quad \text{KL}(\mathbf{g}_\ell \| \mathbf{g}'_\ell), \quad (6)$$

where  $S_\ell^{\text{pre}}$  and  $S_\ell^{\text{post}}$  denote the expert sets before and after editing, respectively. To characterize  $\delta \mathbf{g}_\ell$  analytically, we linearize the softmax around  $\mathbf{s}_\ell$ :

$$\mathbf{g}'_\ell \approx \mathbf{g}_\ell + J_{\text{sm}}(\mathbf{s}_\ell) \delta \mathbf{s}_\ell \quad \Rightarrow \quad \delta \mathbf{g}_\ell \approx J_{\text{sm}}(\mathbf{s}_\ell) \mathbf{E}_\ell^\top \delta \mathbf{u}_\ell, \quad (7)$$

<sup>1</sup>For clarity of analysis, we use the full softmax distribution and omit the Top- $K$  selection, so that  $\mathbf{g}_\ell$  remains differentiable for the Jacobian-based first-order analysis.

Where  $\approx$  denotes a first-order Taylor approximation of the softmax function around  $\mathbf{s}_\ell$ , when the perturbation  $\delta \mathbf{s}_\ell$  is small, and  $J_{\text{sm}}(\mathbf{s}) = \text{diag}(\text{sm}(\mathbf{s})) - \text{sm}(\mathbf{s}) \text{sm}(\mathbf{s})^\top$  is the Jacobian of softmax and  $\text{sm}$  is the softmax function. This relation highlights a crucial observation: only the component of  $\delta \mathbf{u}_\ell$  that lies in the span of  $\mathbf{E}_\ell$  influences the routing probabilities, and the Jacobian can amplify such components, potentially destabilizing expert selection. This insight motivates our design—suppressing the projection of perturbations onto  $\text{span}(\mathbf{E}_\ell)$  is key to preventing routing drift.

#### 4.2 PER-EXPERT NULL-SPACE PROJECTION REPARAMETERIZATION

As established in Section 4.1, suppressing routing drift amounts to ensuring that  $\delta \mathbf{u}_\ell \approx \mathbf{0}$  on a preservation set. To achieve this by construction, we reparameterize each expert update so that its effect vanishes along directions spanned by preservation features. Inspired by the null-space constrained approach of ALPHAEDIT for dense models, we generalize the idea to the MoE setting by computing a per-expert projector that filters out harmful update components.

Concretely, recall from Eqn. 4 that expert  $n$  produces output  $\mathbf{W}_n \mathbf{k}_{i,n}$  for post-activation key  $\mathbf{k}_{i,n}$ . Let  $\mathcal{P}$  denote the set of preservation prompts, and collect their features for expert  $n$  into the matrix  $\mathbf{K}_n^0 = [\mathbf{k}_{i,n}]_{i \in \mathcal{P}} \in \mathbb{R}^{d_k \times |\mathcal{P}|}$ . The covariance  $\mathbf{K}_n^0 \mathbf{K}_n^{0\top}$  captures the subspace of activations we wish to keep invariant. We compute its eigendecomposition,  $\mathbf{K}_n^0 \mathbf{K}_n^{0\top} = \mathbf{U}_n \Lambda_n \mathbf{U}_n^\top$ , and select indices  $\mathcal{I}_0 = \{p : \lambda_{n,p} < \tau\}$  corresponding to (near-)null eigenvalues under a small threshold  $\tau > 0$ . Let  $\mathbf{U}_n^0 = \mathbf{U}_n[:, \mathcal{I}_0]$  and define the orthogonal projector onto the complement of  $\text{span}(\mathbf{K}_n^0)$  by  $\mathbf{P}_n = \mathbf{U}_n^0 \mathbf{U}_n^{0\top}$ . Intuitively,  $\mathbf{P}_n$  preserves only those directions orthogonal to all preservation features, so any update projected by  $\mathbf{P}_n$  is guaranteed not to alter expert outputs on  $\mathcal{P}$ .

We then reparameterize the expert update as  $\Delta_n = \hat{\Delta}_n \mathbf{P}_n$ , where  $\hat{\Delta}_n$  is the free variable to be optimized. Because  $\mathbf{P}_n \mathbf{k}_{i,n} = \mathbf{0}$  for all  $i \in \mathcal{P}$  (up to numerical error), the preservation outputs are unaffected:  $\hat{\Delta}_n \mathbf{P}_n \mathbf{k}_{i,n} = \mathbf{0}$ . Consequently, for every  $i \in \mathcal{P}$  we have  $\delta \mathbf{u}_\ell(i) = \mathbf{0}$ , and by Eqn. 7 the induced routing shift satisfies  $\delta g_\ell(i) \approx \mathbf{0}$ , minimizing Eqn. 6 on the preservation set.

**Projected editing objective.** Let  $\tilde{\mathbf{k}}_{i,n} = \mathbf{P}_n \mathbf{k}_{i,n}$  denote the projected key. Substituting  $\Delta_n = \hat{\Delta}_n \mathbf{P}_n$  into the MoE objective (Eqn. 5) yields

$$\{\hat{\Delta}_n\}_{n=1}^N = \arg \min_{\{\Delta_n\}} \sum_{i \in \mathcal{E}} \left\| \sum_{n=1}^N g_{i,n} (\mathbf{W}_n \mathbf{k}_{i,n} + \tilde{\Delta}_n \tilde{\mathbf{k}}_{i,n}) - \mathbf{v}_i \right\|^2 + \lambda \sum_{n=1}^N \|\tilde{\Delta}_n\|^2, \quad (8)$$

where  $\lambda \geq 0$  controls the update magnitude and improves locality. No separate preservation term is needed, as  $\mathbf{P}_n$  removes all preservation components by construction.

#### 4.3 RANDOMIZED BLOCK COORDINATE DESCENT SOLVER

A naive step is to solve the projected objective in Eqn. 8 in one shot, just like what we do in dense model KE (Meng et al., 2022; 2023; Fang et al., 2025). In this subsection, we (i) expose the structure of the global closed-form solution, (ii) explain why the direct route is computationally impractical in MoE, and (iii) arrive at an efficient randomized block coordinate descent (BCD) procedure that scales with the expert hidden size.

**The global closed-form (one shot).** Let  $\hat{\Delta}_n \in \mathbb{R}^{d_m \times d_k}$  be the projected free variable for expert  $n$ , and stack all expert updates horizontally as  $\hat{\Delta} = [\hat{\Delta}_1 \cdots \hat{\Delta}_N] \in \mathbb{R}^{d_m \times (Nd_k)}$ . For edit example  $i$ , define the base residual (excluding any edits)  $\mathbf{r}_i = \mathbf{v}_i - \sum_{n=1}^N g_{i,n} \mathbf{W}_n \mathbf{k}_{i,n}$ , and the design vector  $\tilde{\psi}_i = [g_{i,1} \tilde{\mathbf{k}}_{i,1}^\top \cdots g_{i,N} \tilde{\mathbf{k}}_{i,N}^\top]^\top \in \mathbb{R}^{Nd_k}$ , where  $\tilde{\mathbf{k}}_{i,n} = \mathbf{P}_n \mathbf{k}_{i,n}$ . Then Eqn. 8 becomes a regularized multi-output linear regression:  $\min_{\hat{\Delta}} \sum_{i \in \mathcal{E}} \|\hat{\Delta} \tilde{\psi}_i - \mathbf{r}_i\|^2 + \lambda \sum_{n=1}^N \|\hat{\Delta}_n\|^2$ . Vectorizing with  $\boldsymbol{\theta} = \text{vec}(\hat{\Delta}) \in \mathbb{R}^{d_m Nd_k}$  and using  $\text{vec}(\hat{\Delta} \tilde{\psi}_i) = (\tilde{\psi}_i^\top \otimes \mathbf{I}_{d_m}) \boldsymbol{\theta}$ , the normal equations take the compact form

$$\left( \sum_{i \in \mathcal{E}} (\tilde{\psi}_i \tilde{\psi}_i^\top) \otimes \mathbf{I}_{d_m} + \lambda \mathbf{I}_{d_m Nd_k} \right) \boldsymbol{\theta} = \sum_{i \in \mathcal{E}} (\tilde{\psi}_i \otimes \mathbf{I}_{d_m}) \mathbf{r}_i, \quad (9)$$

with unique minimizer

$$\boldsymbol{\theta}^* = \mathbf{M}_{\text{glob}}^{-1} \mathbf{b}_{\text{glob}} \quad \text{and} \quad \hat{\Delta}^* = \text{unvec}(\boldsymbol{\theta}^*), \quad (10)$$

where  $M_{\text{glob}} = \sum_i (\tilde{\psi}_i \tilde{\psi}_i^\top) \otimes I_{d_m} + \lambda I$  and  $b_{\text{glob}} = \sum_i (\tilde{\psi}_i \otimes I_{d_m}) r_i$ . A proof is provided in Appendix B.2.

Although Eqn. 9 is elegant, it is not a practical editing primitive at MoE scale. First, even exploiting the Kronecker structure, the system decomposes into  $d_m$  independent problems of size  $(Nd_k) \times (Nd_k)$  each. For typical MoE layers,  $N$  can be 8–128 and  $d_k$  in the thousands, so factorizing  $d_m$  such systems (and re-factorizing as  $\mathcal{E}$  changes) is prohibitively expensive in both time  $O(d_m(Nd_k)^3)$  and memory  $O(d_m(Nd_k)^2)$ . Second, while Top- $K$  routing makes each  $\tilde{\psi}_i$   $K$ -block sparse, the accumulated Gram matrix  $\sum_i \tilde{\psi}_i \tilde{\psi}_i^\top$  quickly densifies, yielding substantial fill-in under Cholesky/LDL<sup>⊤</sup>. These realities make the one-shot solve not suitable for fast, iterative MoE editing.

**From global to block: randomized BCD.** Since the insight that every expert is a natural chunk. The structure of Eqn. 8 suggests a block strategy: treat each expert as a block and optimize one block while holding the rest fixed. This reduces the problem to a sequence of well-conditioned, small ridge least-squares solves of size  $d_k \times d_k$ .

Fix  $\{\hat{\Delta}_\ell\}_{\ell \neq n}$  and define the residual that excludes expert  $n$ :

$$r_i^{(-n)} = v_i - \sum_{\ell \neq n} g_{i,\ell} (W_\ell k_{i,\ell} + \hat{\Delta}_\ell \tilde{k}_{i,\ell}). \quad (11)$$

The subproblem in  $\hat{\Delta}_n$  becomes the ridge-regularized least squares

$$\min_{\hat{\Delta}_n} \sum_{i \in \mathcal{E}} \|r_i^{(-n)} - g_{i,n} \hat{\Delta}_n \tilde{k}_{i,n}\|^2 + \lambda \|\hat{\Delta}_n\|^2. \quad (12)$$

Its normal equations

$$\hat{\Delta}_n \left( \underbrace{\sum_{i \in \mathcal{E}} g_{i,n}^2 \tilde{k}_{i,n} \tilde{k}_{i,n}^\top}_{M_n \in \mathbb{R}^{d_k \times d_k}} + \lambda I \right) = \left( \underbrace{\sum_{i \in \mathcal{E}} g_{i,n} r_i^{(-n)} \tilde{k}_{i,n}^\top}_{B_n \in \mathbb{R}^{d_m \times d_k}} \right), \quad (13)$$

admit the closed-form update

$$\hat{\Delta}_n^* = B_n M_n^{-1} = \left( \sum_{i \in \mathcal{E}} g_{i,n} r_i^{(-n)} \tilde{k}_{i,n}^\top \right) \left( \sum_{i \in \mathcal{E}} g_{i,n}^2 \tilde{k}_{i,n} \tilde{k}_{i,n}^\top + \lambda I \right)^{-1}. \quad (14)$$

We then write to parameters via the projection  $\Delta_n^* = \hat{\Delta}_n^* P_n$  and move to the next block.

**Practicalities and complexity.** We traverse experts in randomized order and update only those active in the current minibatch, which further reduces cost. For each updated expert, forming  $M_n$  costs  $O(|\mathcal{E}|d_k^2)$  and inverting it costs  $O(d_k^3)$ , typically modest since  $d_k \ll d_m$ . We stream-accumulate  $B_n$  and  $M_n$ , cache  $\tilde{k}_{i,n}$ , and use Cholesky with diagonal loading for numerical stability. Because Eqn. 8 is a strictly convex quadratic in  $\{\hat{\Delta}_n\}$ , (randomized) BCD with exact block solves converges globally under standard conditions (Tseng, 2001; Richtárik & Takáč, 2014). Empirically we see fast decrease within a few passes ( $\leq 10$ ).

## 5 EXPERIMENTS

### 5.1 BASELINES, DATASETS, AND METRICS

We evaluate two modern MoE LLMs on standard factual-editing benchmarks: **Qwen3-30B-A3B** (Yang et al., 2025) (128 experts; top-8 per token) and **GPT-OSS-20B** (Agarwal et al., 2025) (32 experts; top-4 per token). As baselines, we adapt parameter-editing methods originally designed for dense Transformers but directly applicable to MoE models: Fine-Tuning (FT) (Zhu et al., 2020), FT-L (FT with a norm constraint), AdaLoRA (Zhang et al., 2023), and UnKE (Deng et al., 2025).

Following prior work, we use **COUNTERFACT** (single-hop counterfactual edits introduced with MEMIT) (Meng et al., 2023) and **ZsRE** (zero-shot relation extraction) (Levy et al., 2017). Visualized dataset examples are provided in Appendix E for unfamiliar readers. We report the standard editing metrics (Meng et al., 2022; 2023; Mitchell et al., 2022a): (i) **Efficacy** (edit success on edited prompts),

Table 1: Sequential knowledge editing on MoE LLMs. *Eff.*, *Gen.*, *Spe.* denote Efficacy, Generalization, Specificity; *Uti.* is their mean (higher is better  $\uparrow$ ). Best in **bold**, second-best underlined.

Method	Model	COUNTERFACT				ZsRE			
		Eff. $\uparrow$	Gen. $\uparrow$	Spe. $\uparrow$	Uti. $\uparrow$	Eff. $\uparrow$	Gen. $\uparrow$	Spe. $\uparrow$	Uti. $\uparrow$
Pre-edited		13.30 $\pm$ 0.34	15.10 $\pm$ 0.31	84.45 $\pm$ 0.24	37.62	41.30 $\pm$ 0.29	40.50 $\pm$ 0.28	40.91 $\pm$ 0.27	40.90
FT	Qwen3-30B-A3B	80.70 $\pm$ 0.39	63.95 $\pm$ 0.43	41.44 $\pm$ 0.39	62.03	6.44 $\pm$ 0.14	6.13 $\pm$ 0.14	2.15 $\pm$ 0.06	4.91
FT-L		82.40 $\pm$ 0.38	22.75 $\pm$ 0.33	71.48 $\pm$ 0.25	58.88	<u>44.19</u> $\pm$ 0.29	<u>42.46</u> $\pm$ 0.29	<u>41.92</u> $\pm$ 0.27	<u>42.86</u>
AdaLoRA		51.90 $\pm$ 0.50	49.75 $\pm$ 0.40	48.10 $\pm$ 0.26	49.92	3.66 $\pm$ 0.09	3.60 $\pm$ 0.09	4.68 $\pm$ 0.10	3.98
UnKE		<u>89.30</u> $\pm$ 0.31	<u>82.85</u> $\pm$ 0.33	48.15 $\pm$ 0.33	<u>73.43</u>	31.43 $\pm$ 0.28	29.78 $\pm$ 0.27	25.30 $\pm$ 0.23	28.84
MoEEEdit		<b>99.30</b> $\pm$ 0.08	<b>94.10</b> $\pm$ 0.20	<b>80.97</b> $\pm$ 0.25	<b>91.46</b>	<b>84.47</b> $\pm$ 0.22	<b>78.01</b> $\pm$ 0.28	<b>42.82</b> $\pm$ 0.28	<b>68.43</b>
Pre-edited		11.80 $\pm$ 0.32	14.70 $\pm$ 0.31	84.53 $\pm$ 0.24	37.01	33.20 $\pm$ 0.28	32.14 $\pm$ 0.28	28.02 $\pm$ 0.00	31.12
FT	GPT-OSS-20B	<u>83.40</u> $\pm$ 0.37	<b>58.40</b> $\pm$ 0.42	55.72 $\pm$ 0.33	<u>65.84</u>	25.57 $\pm$ 0.28	23.41 $\pm$ 0.26	17.61 $\pm$ 0.21	22.20
FT-L		73.80 $\pm$ 0.44	43.10 $\pm$ 0.48	59.75 $\pm$ 0.33	58.88	32.75 $\pm$ 0.29	33.09 $\pm$ 0.30	30.06 $\pm$ 0.26	31.97
AdaLoRA		62.40 $\pm$ 0.48	55.00 $\pm$ 0.42	43.65 $\pm$ 0.34	53.68	43.46 $\pm$ 0.30	42.96 $\pm$ 0.30	<b>33.60</b> $\pm$ 0.24	40.01
UnKE		78.00 $\pm$ 0.41	44.40 $\pm$ 0.42	<u>73.91</u> $\pm$ 0.28	65.44	<u>46.58</u> $\pm$ 0.31	<u>43.99</u> $\pm$ 0.31	31.40 $\pm$ 0.26	<u>40.66</u>
MoEEEdit		<b>95.90</b> $\pm$ 0.20	44.10 $\pm$ 0.43	<b>81.09</b> $\pm$ 0.25	<b>73.70</b>	<b>81.68</b> $\pm$ 0.25	<b>68.44</b> $\pm$ 0.34	<u>32.55</u> $\pm$ 0.26	<b>60.89</b>

(ii) **Generalization** (success on paraphrases and lightly perturbed contexts), and (iii) **Specificity** (locality on unrelated controls). Unless stated otherwise, we perform sequential batched edits. And to summarize the overall trade-off, we additionally report **Utility** as the mean of the three metrics. See Appendix A for full calculation details.

## 5.2 MAIN RESULTS ON KNOWLEDGE EDITING

We perform 1,000 sequential edits on each dataset (COUNTERFACT and ZsRE) with a batch size of 50 edits for all methods. Table 1 summarizes results on Qwen3-30B-A3B and GPT-OSS-20B.

As shown in Table 1, MoEEEdit consistently delivers outstanding results. On COUNTERFACT, it achieves over 90 efficacy on both backbones, substantially outperforming UnKE and FT-L in generalization and specificity, respectively. On GPT-OSS-20B, although FT attains slightly higher generalization, MoEEEdit still provides the best overall balance, with clear gains in efficacy (+12.5) and specificity (+7.2). On ZsRE, MOEEDIT also demonstrates large improvements in efficacy and generalization—over +30 points against the strongest baselines—while maintaining competitive specificity (within 1 point of AdaLoRA). These results highlight that MOEEDIT offers the most favorable trade-off between accuracy and locality across models and datasets.

## 5.3 MAIN RESULTS ON ROUTING DISTRIBUTION SHIFT

We analyze routing distribution shifts under sequential editing. Similar to Section 5.2, we perform 1,000 edits with a batch size of 50. To control for depth, all methods are constrained to update at most the top editing layer (layer 7). Specifically, MoEEEdit applies updates across layers {3,4,5,6,7} using BCD, while FT, FT-L, UnKE, and AdaLoRA are restricted to layer 7. We evaluate on Qwen3-30B-A3B/COUNTERFACT and report the routing-similarity (RS) metric between pre- and post-edit Top- $K$  expert sets (Eqn. 6), averaged over windows of 10 layers, for both the editing and preservation sets. The editing set consists of the 1,000 edited samples, while the preservation set is formed by sampling an equal number of untouched examples from the remaining dataset. Table 2 summarizes the results.

**Projection suppresses routing drift and preserves stability.** As shown in Table 2, methods directly extended from dense models exhibit substantial routing drift, whereas MoEEEdit maintains consistently high routing stability (average RS > 88 across all layer ranges for both sets). Considering that Qwen3-30B-A3B activates 8 experts per token, the average number of non-overlapping experts before and after editing is close to one, which is negligible. This observation aligns with the heavy-tailed nature of routing: small perturbations primarily affect low-weight expert selections that contribute little to the output. The average KL divergence between pre- and post-edit routing distributions for MoEEEdit is only 0.02, indicating minimal shift. Furthermore, Figure 2 plots routing similarity across layers for both editing and preservation sets. AdaLoRA and FT exhibit the lowest RS values across layers due

Table 2: Routing distribution shift on Qwen3-30B-A3B. Values are Jaccard similarity ( $\text{RS}\uparrow$ ) between pre- and post-edit routing distributions. Higher is better. Best results are in **bold**, second-best are underlined.

Method	Model	Editing Set $\text{RS}\uparrow$			Preservation Set $\text{RS}\uparrow$		
		Lay. 11–20	Lay. 21–30	Lay. 31–40	Lay. 11–20	Lay. 21–30	Lay. 31–40
FT	Qwen3-30B-A3B	23.57	26.58	29.98	24.72	27.45	30.97
FT-L		47.01	<u>51.20</u>	<u>53.68</u>	48.80	<u>50.17</u>	<u>53.45</u>
AdaLoRA		16.63	24.11	27.00	16.38	23.84	26.60
UnKE		<u>52.46</u>	44.12	44.80	<u>49.90</u>	41.91	43.84
MoEEEdit		<b>86.62</b>	<b>88.16</b>	<b>89.93</b>	<b>87.02</b>	<b>88.55</b>	<b>90.22</b>

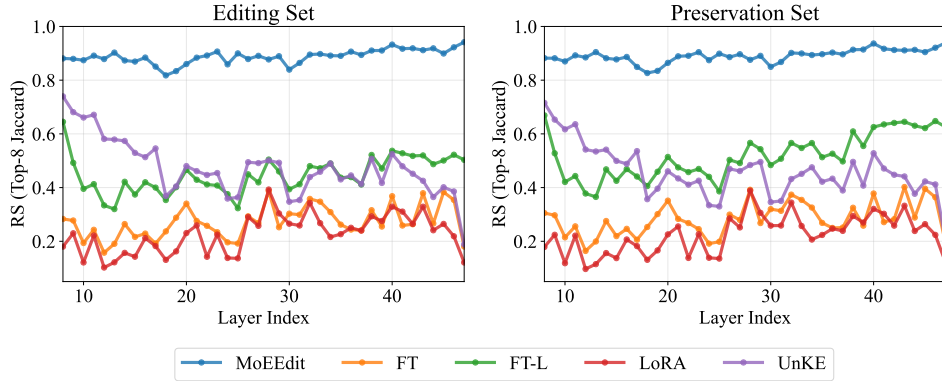


Figure 2: Routing similarity (RS) before and after editing on the editing and preservation sets. MoEEEdit achieves consistently high RS, demonstrating strong routing stability.

to their unconstrained updates that heavily disrupt routing patterns. In contrast, MoEEEdit preserves routing stability across all layers and consistently outperforms all baselines.

#### 5.4 ABLATION STUDY

**Effect of Projection.** We ablate the projection matrix in MoEEEdit to evaluate its contribution to routing stability. As shown in Table 3, removing projection reduces RS by an average of 14.81 points on the editing set and 15.21 on the preservation set. The KL divergence also increases from 0.02 to 0.0834, confirming that projection is critical for suppressing routing drift. These results validate the projection design introduced in Section 4.2.

Table 3: Ablation on the projection matrix. Removing projection significantly weakens routing stability.

Method	Set	$\text{RS}\uparrow$		
		Lay. 11–20	Lay. 21–30	Lay. 31–40
MoEEEdit	Edit.	86.62	88.16	89.93
MoEEEdit (w/o Proj)		73.64	72.90	73.75
MoEEEdit	Pres.	87.02	88.55	90.22
MoEEEdit (w/o Proj)		73.59	73.08	73.50

**BCD Solver vs. Closed-form Solver.** We compare the proposed block coordinate descent (BCD) solver with a naive closed-form solution. The latter requires inverting large matrices that scale with the number of experts and constraints, making it computationally infeasible at realistic scales. To enable comparison, we construct a controlled synthetic batch and evaluate (i) convergence and (ii) scalability. Results are shown in Figure 3.

Figure 3 illustrates two perspectives on our BCD solver: (a) reconstruction error convergence under varying  $\lambda$ , and (b) runtime scalability with respect to the number of experts. Smaller  $\lambda$  values (e.g.,  $10^{-4}$ ,  $10^{-3}$ ) achieve lower error, whereas larger values converge to higher error floors. Panel (b)

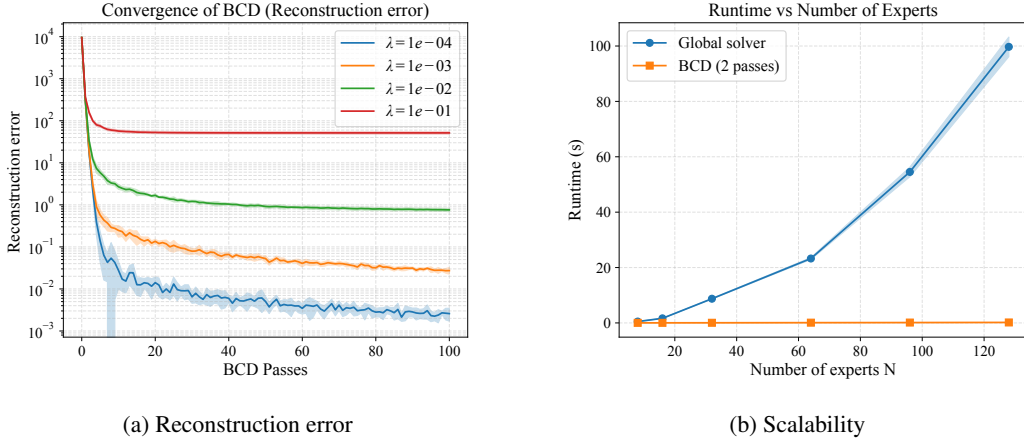


Figure 3: Comparison of solvers. (a) BCD achieves fast convergence with suitable  $\lambda$ . (b) BCD scales efficiently with the number of experts, while the closed-form solver quickly becomes infeasible.

shows that the closed-form solver exhibits near-quadratic runtime growth and becomes infeasible beyond  $N \approx 60$ , while BCD maintains nearly constant runtime up to 128 experts. Thus, BCD scales linearly with hidden size rather than the total number of experts, ensuring practical efficiency.

**Number of BCD Passes.** We vary the number of BCD passes  $\in \{2, 4, 6, 8, 10, 12, 14, 16, 18, 20\}$  while keeping all other settings fixed, and evaluate efficacy, generality, specificity, and their harmonic mean. As shown in Figure 4, as the number of BCD passes increases, both efficacy and generality rise rapidly in the early phase and then plateau, reflecting diminishing returns beyond moderate passes. Specificity shows a more stable, gradual upward trend. Since each subproblem is convex, early passes remove dominant residuals, while later passes only refine small residuals, yielding slower gains. This behavior suggests that a small number of passes (e.g., 6–10) already achieves a favorable trade-off between performance and efficiency. Beyond this range, additional passes bring only marginal improvements while increasing runtime, highlighting the practicality of moderate BCD iterations in large-scale editing scenarios.

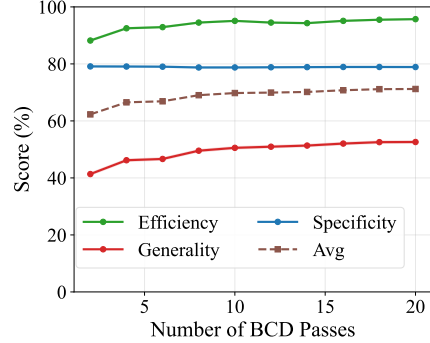


Figure 4: Ablation on the number of passes.

## 6 DISCUSSION AND CONCLUSION

In this work, we presented MoEEEdit, a routing-stable knowledge editing framework tailored for Mixture-of-Experts (MoE) LLMs. Our approach addresses the unique challenges of computational cost, inter-expert coupling, and routing drift by combining per-expert null-space projection with an efficient block coordinate descent solver. Extensive experiments on COUNTERFACT and ZsRE benchmarks demonstrate that MoEEEdit achieves high efficacy, strong generalization, and routing stability, all with superior efficiency compared to prior methods.

Beyond empirical performance, our findings highlight several broader insights. First, expert-aware design is crucial: naive adaptations of dense-model editors to MoEs fail to maintain stability and efficiency. Second, routing stability emerges as a central factor in editing sparse architectures, where even small perturbations can cascade through routing distributions. By explicitly controlling for this effect, MoEEEdit offers a principled solution that preserves locality without sacrificing scalability.

In summary, MoEEEdit establishes a robust foundation for precise and scalable knowledge editing in sparse architectures, advancing the state of the art and paving the way for more adaptive and trustworthy MoE-based language models.

## REFERENCES

- Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao, Boaz Barak, Ally Bennett, Tyler Bertao, Nivedita Brett, Eugene Brevdo, Greg Brockman, Sebastien Bubeck, Che Chang, Kai Chen, Mark Chen, Enoch Cheung, Aidan Clark, Dan Cook, Marat Dukhan, Casey Dvorak, Kevin Fives, Vlad Fomenko, Timur Garipov, Kristian Georgiev, Mia Glaese, Tarun Gogineni, Adam Goucher, Lukas Gross, Katia Gil Guzman, John Hallman, Jackie Hehir, Johannes Heidecke, Alec Helyar, Haitang Hu, Romain Huet, Jacob Huh, Saachi Jain, Zach Johnson, Chris Koch, Irina Kofman, Dominik Kundel, Jason Kwon, Volodymyr Kyrlyov, Elaine Ya Le, Guillaume Leclerc, James Park Lennon, Scott Lessans, Mario Lezcano-Casado, Yuanzhi Li, Zhuohan Li, Ji Lin, Jordan Liss, Lily, Liu, Jiancheng Liu, Kevin Lu, Chris Lu, Zoran Martinovic, Lindsay McCallum, Josh McGrath, Scott McKinney, Aidan McLaughlin, Song Mei, Steve Mostovoy, Tong Mu, Gideon Myles, Alexander Neitz, Alex Nichol, Jakub Pachocki, Alex Paino, Dana Palmie, Ashley Pantuliano, Giambattista Parascandolo, Jongsoo Park, Leher Pathak, Carolina Paz, Ludovic Peran, Dmitry Pimenov, Michelle Pokrass, Elizabeth Proehl, Huida Qiu, Gaby Rila, Filippo Raso, Hongyu Ren, Kimmy Richardson, David Robinson, Bob Rotsted, Hadi Salman, Suvansh Sanjeev, Max Schwarzer, D. Sculley, Harshit Sikchi, Kendal Simon, Karan Singhal, Yang Song, Dane Stuckey, Zhiqing Sun, Philippe Tillet, Sam Toizer, Foivos Tsimplouras, Nikhil Vyas, Eric Wallace, Xin Wang, Miles Wang, Olivia Watkins, Kevin Weil, Amy Wendling, Kevin Whinnery, Cedric Whitney, Hannah Wong, Lin Yang, Yu Yang, Michihiro Yasunaga, Kristen Ying, Wojciech Zaremba, Wenting Zhan, Cyril Zhang, Brian Zhang, Eddie Zhang, and Shengjia Zhao. gpt-oss-120b gpt-oss-20b model card, 2025. URL <https://arxiv.org/abs/2508.10925>.
- Damai Dai, Li Dong, Yaru Hao, Dian Sui, Songhao Piao, Longyue Dou, Weinan Wang, and Furu Wei. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 8493–8502, 2022. doi: 10.18653/v1/2022.acl-long.584.
- Jingcheng Deng, Zihao Wei, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. Everything is editable: Extend knowledge editing to unstructured data in large language models. In *ICLR*, 2025. URL <https://openreview.net/forum?id=X5r05VyTgB>.
- Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten P Bosma, Zongwei Zhou, Tao Wang, Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, Kathleen Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc Le, Yonghui Wu, Zhifeng Chen, and Claire Cui. GLaM: Efficient scaling of language models with mixture-of-experts. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 5547–5569. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/du22c.html>.
- Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Jie Shi, Xiang Wang, Xiangnan He, and Tat-Seng Chua. Alphaedit: Null-space constrained model editing for language models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=HvSytvg3Jh>.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5484–5495, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.446. URL <https://aclanthology.org/2021.emnlp-main.446/>.
- Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=qrwe7XHTmYb>.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. Zero-shot relation extraction via reading comprehension. In Roger Levy and Lucia Specia (eds.), *Proceedings of the 21st Conference*

on *Computational Natural Language Learning (CoNLL 2017)*, pp. 333–342, Vancouver, Canada, August 2017. Association for Computational Linguistics. doi: 10.18653/v1/K17-1034. URL <https://aclanthology.org/K17-1034/>.

Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. Pmet: Precise model editing in a transformer. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17): 18564–18572, Mar. 2024. doi: 10.1609/aaai.v38i17.29818. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29818>.

Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. *IEEE Transactions on Audio, Speech and Language Processing*, 33:3776–3786, 2025. doi: 10.1109/TASLPRO.2025.3606231.

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS ’22*, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.

Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=MkbcAHlYgyS>.

Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast model editing at scale. In *International Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=0DcZxeWfOPt>.

Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. Memory-based model editing at scale. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 15817–15831. PMLR, 17–23 Jul 2022b. URL <https://proceedings.mlr.press/v162/mitchell122a.html>.

Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language models as knowledge bases? In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2463–2473, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1250. URL <https://aclanthology.org/D19-1250/>.

Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144:1–38, 2014.

Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=B1ckMDqlg>.

Kai Sun, Yifan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. Head-to-tail: How knowledgeable are large language models (LLMs)? A.K.A. will LLMs replace knowledge graphs? In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 311–325, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.18. URL <https://aclanthology.org/2024.naacl-long.18/>.

Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications*, 109(3):475–494, 2001.

Renzhi Wang and Piji Li. LEMoE: Advanced mixture of experts adaptor for lifelong model editing of large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp.

2551–2575, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.149. URL <https://aclanthology.org/2024.emnlp-main.149/>.

Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. Knowledge editing for large language models: A survey. *ACM Comput. Surv.*, 57(3), November 2024. ISSN 0360-0300. doi: 10.1145/3698590. URL <https://doi.org/10.1145/3698590>.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.

Ningyu Zhang, Bozhong Tian, Siyuan Cheng, Xiaozhuan Liang, Yi Hu, Kouying Xue, Yanjie Gou, Xi Chen, and Huajun Chen. Instructedit: instruction-based knowledge editing for large language models. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI ’24*, 2024a. ISBN 978-1-956792-04-1. doi: 10.24963/ijcai.2024/733. URL <https://doi.org/10.24963/ijcai.2024/733>.

Ningyu Zhang, Zekun Xi, Yujie Luo, Peng Wang, Bozhong Tian, Yunzhi Yao, Jintian Zhang, Shumin Deng, Mengshu Sun, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and Huajun Chen. Oneedit: A neural-symbolic collaboratively knowledge editing system, 2024b. URL <https://arxiv.org/abs/2409.07497>.

Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and Huajun Chen. A comprehensive study of knowledge editing for large language models, 2024c. URL <https://arxiv.org/abs/2401.01286>.

Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *International Conference on Learning Representations (ICLR)*, 2023. URL <https://openreview.net/forum?id=lq62uWRJjiY>.

Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. Can we edit factual knowledge by in-context learning? In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL <https://openreview.net/forum?id=hsjQHAM8MV>.

Zexuan Zhong, Zhengxuan Wu, Christopher Manning, Christopher Potts, and Danqi Chen. MQuAKE: Assessing knowledge editing in language models via multi-hop questions. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 15686–15702, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.971. URL <https://aclanthology.org/2023.emnlp-main.971/>.

Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix X. Yu, and Sanjiv Kumar. Modifying memories in transformer models. *arXiv preprint arXiv:2012.00363*, 2020. URL <https://arxiv.org/abs/2012.00363>.

## A METRICS

In this section, we introduce the evaluation metrics we use for COUNTERFACT and ZsRE

### A.1 ZSRE EVALUATION METRICS

Building on prior studies, we define for each ZsRE metric a large language model  $\mathcal{M}$ , a factual prompt  $(s_i, r_i)$ , a revised target output  $y_i$ , and the model’s initial output  $\hat{y}_i$ :

**Effectiveness.** Effectiveness is quantified as the average top-1 accuracy on edited inputs:

$$\mathbb{E}_i \left[ \mathbf{1} \left( y_i = \arg \max_y \Pr(y \mid (s_i, r_i)) \right) \right]. \quad (15)$$

**Generalization.** This measures the ability of the model to perform correctly on paraphrased prompts  $\tilde{N}(s_i, r_i)$ . It is computed as the average accuracy over such rephrasings:

$$\mathbb{E}_i \left[ \mathbf{1} \left( y_i = \arg \max_y \Pr(y \mid \tilde{N}(s_i, r_i)) \right) \right]. \quad (16)$$

**Specificity.** Specificity ensures that modifications do not affect unrelated cases  $\Omega(s_i, r_i)$ . It is defined as:

$$\mathbb{E}_i \left[ \mathbf{1} \left( \hat{y}_i = \arg \max_y \Pr(y \mid \Omega(s_i, r_i)) \right) \right]. \quad (17)$$

### A.2 COUNTERFACTUAL EVALUATION METRICS

Similarly, we introduce Counterfactual metrics for  $\mathcal{M}$  under the same setup  $(s_i, r_i)$  with target  $y_i$  and original  $\hat{y}_i$ :

**Effectiveness (success ratio).** The share of cases where  $y_i$  is assigned higher probability than  $\hat{y}_i$  under  $(s_i, r_i)$ :

$$\mathbb{E}_i \left[ \Pr(y_i \mid (s_i, r_i)) > \Pr(\hat{y}_i \mid (s_i, r_i)) \right]. \quad (18)$$

**Generalization (paraphrase success).** The proportion of paraphrased prompts  $\tilde{N}(s_i, r_i)$  where  $y_i$  has higher likelihood than  $\hat{y}_i$ :

$$\mathbb{E}_i \left[ \Pr(y_i \mid \tilde{N}(s_i, r_i)) > \Pr(\hat{y}_i \mid \tilde{N}(s_i, r_i)) \right]. \quad (19)$$

**Specificity (neighborhood success).** For neighborhood prompts  $\Omega(s_i, r_i)$  that involve related but distinct entities, specificity is the fraction of cases where  $y_i$  is favored over  $\hat{y}_i$ :

$$\mathbb{E}_i \left[ \Pr(y_i \mid \Omega(s_i, r_i)) > \Pr(\hat{y}_i \mid \Omega(s_i, r_i)) \right]. \quad (20)$$

## B PROOF & KNOWLEDGE EDITING

### B.1 GENERAL FRAMEWORK OF LOCATE-THEN-EDIT

Knowledge editing seeks to precisely revise a model’s behavior to recall a new fact  $(s, r, o^*)$  in place of an obsolete or incorrect one  $(s, r, o)$ , conditioned on a prompt  $p(s, r)$ . While various techniques exist, the dominant *locate-then-edit* paradigm (Meng et al., 2022; 2023) typically decomposes the process into three distinct phases: locating the mediating parameters, computing the optimal local update targets, and solving for the new weights. We formalize this process below using notation consistent with Section 3.

**Step 1: Causal Localization.** The initial phase identifies the specific layer  $l$  and module (e.g., a dense FFN or specific experts in an MoE layer) that mediate the retrieval of the factual association. This is commonly achieved via *Causal Tracing* (Meng et al., 2022). By corrupting hidden states

with noise to degrade the model’s prediction and subsequently restoring states at specific layers, one can quantify the *Indirect Effect* (IE) of each layer on the correct output probability. The layer  $l^*$  exhibiting the maximal causal influence is selected as the target for editing:

$$l^* = \arg \max_l \text{IE}(l). \quad (21)$$

**Step 2: Acquiring the Target Output Vector ( $v^*$ ).** Once the target layer  $l^*$  is identified, we must determine the optimal output representation required to successfully trigger the target token  $o^*$ . Viewing the layer as a linear associative memory, it maps a key  $k$  (input features) to a value  $v$  (output features). The objective is to identify a new output vector  $v^*$  such that, if the layer were to produce  $v^*$ , the final model prediction would be  $o^*$ .

This step is formulated as an optimization problem over the hidden state vector rather than the model parameters. Let  $G(v)$  denote the function mapping the layer output  $v$  to the final model logits. We freeze the model parameters and optimize a perturbation  $\delta$  to the original output  $v$ , aiming to maximize the log-likelihood of the target object  $o^*$ :

$$v^* = v + \delta^*, \quad \text{where} \quad \delta^* = \arg \min_{\delta} -\log \mathbb{P}_{\mathcal{M}}(o^* \mid \text{do}(v \leftarrow v + \delta)). \quad (22)$$

Here, the  $\text{do}(\cdot)$  operator signifies a causal intervention where the layer output is manually set to  $v + \delta$ . A regularization term (e.g., KL divergence) is typically included in the objective to minimize prediction drift for the subject  $s$  and relation  $r$ , ensuring the edit remains semantically consistent. This process effectively translates the semantic edit target (the token  $o^*$ ) into a vector-space target  $v^*$ .

**Step 3: Updating Parameters.** The final step is to update the projection weights  $W$  (corresponding to  $W_{out}$  in dense models or  $\{W_n\}$  in MoE experts) to map the specific input key  $k$  to the new target  $v^*$ , while preserving unrelated associations. This is formulated as a constrained least-squares problem.

Let  $\mathcal{E}$  denote the set of edit examples (new facts) and  $\mathcal{P}$  denote the set of preservation examples (invariant knowledge). We require  $Wk_i \approx v_i^*$  for edits  $i \in \mathcal{E}$ , and  $Wk_j \approx v_j$  for preservation samples  $j \in \mathcal{P}$ . The optimal update  $\hat{W}$  minimizes the aggregated error:

$$\hat{W} = \arg \min_W \sum_{i \in \mathcal{E}} \|Wk_i - v_i^*\|^2 + \sum_{j \in \mathcal{P}} \|Wk_j - v_j\|^2. \quad (23)$$

Dense methods such as ROME and MEMIT solve this globally via a closed-form solution involving the covariance matrix of the keys. As discussed in Section 4, our proposed MoEEdit adapts this general objective to address the unique constraints of Mixture-of-Experts architectures.

## B.2 SUPPLEMENTARY PROOF

### B.2.1 PROOF OF THE GLOBAL CLOSED-FORM (ONE-SHOT) SOLUTION

Let  $\{\tilde{k}_{i,n}\}_{i \in \mathcal{E}, n \in [N]}$  be the projected keys and  $g_{i,n} \geq 0$  the router weights. For each edit example  $i \in \mathcal{E}$ , define the base residual

$$\mathbf{r}_i = \mathbf{v}_i - \sum_{n=1}^N g_{i,n} \mathbf{W}_n \mathbf{k}_{i,n}, \quad (24)$$

and the design vector

$$\tilde{\psi}_i = [g_{i,1} \tilde{\mathbf{k}}_{i,1}^\top \cdots g_{i,N} \tilde{\mathbf{k}}_{i,N}^\top]^\top \in \mathbb{R}^{Nd_k}. \quad (25)$$

Stack the expert updates as  $\hat{\Delta} = [\hat{\Delta}_1 \cdots \hat{\Delta}_N] \in \mathbb{R}^{d_m \times (Nd_k)}$ . The projected objective

$$\min_{\hat{\Delta}} \sum_{i \in \mathcal{E}} \|\hat{\Delta} \tilde{\psi}_i - \mathbf{r}_i\|_2^2 + \lambda \sum_{n=1}^N \|\hat{\Delta}_n\|_F^2 \quad (26)$$

has the unique minimizer

$$\theta^* = \mathbf{M}_{\text{glob}}^{-1} \mathbf{b}_{\text{glob}}, \quad \hat{\Delta}^* = \text{unvec}(\theta^*), \quad (27)$$

with

$$\mathbf{M}_{\text{glob}} = \left( \sum_{i \in \mathcal{E}} \tilde{\psi}_i \tilde{\psi}_i^\top \right) \otimes \mathbf{I}_{d_m} + \lambda \mathbf{I}_{d_m N d_k}, \quad \mathbf{b}_{\text{glob}} = \sum_{i \in \mathcal{E}} (\tilde{\psi}_i \otimes \mathbf{I}_{d_m}) \mathbf{r}_i. \quad (28)$$

Let  $m := |\mathcal{E}|$  and define the data matrices

$$\Psi = [\tilde{\psi}_1 \cdots \tilde{\psi}_m] \in \mathbb{R}^{N d_k \times m}, \quad \mathbf{R} = [\mathbf{r}_1 \cdots \mathbf{r}_m] \in \mathbb{R}^{d_m \times m}. \quad (29)$$

The objective becomes

$$\min_{\hat{\Delta} \in \mathbb{R}^{d_m \times (N d_k)}} \|\hat{\Delta} \Psi - \mathbf{R}\|_F^2 + \lambda \|\hat{\Delta}\|_F^2. \quad (30)$$

Taking derivatives gives

$$\nabla_{\hat{\Delta}} = 2(\hat{\Delta} \Psi - \mathbf{R}) \Psi^\top + 2\lambda \hat{\Delta}. \quad (31)$$

Setting this to zero yields

$$\hat{\Delta} (\Psi \Psi^\top + \lambda \mathbf{I}_{N d_k}) = \mathbf{R} \Psi^\top. \quad (32)$$

Using  $\text{vec}(\mathbf{X} \mathbf{A}) = (\mathbf{A}^\top \otimes \mathbf{I}) \text{vec}(\mathbf{X})$ , Eqn. 32 becomes

$$\left( (\Psi \Psi^\top) \otimes \mathbf{I}_{d_m} + \lambda \mathbf{I}_{d_m N d_k} \right) \boldsymbol{\theta} = \text{vec}(\mathbf{R} \Psi^\top), \quad (33)$$

where  $\boldsymbol{\theta} = \text{vec}(\hat{\Delta})$ . Noting that  $\text{vec}(\mathbf{R} \Psi^\top) = \sum_{i=1}^m (\tilde{\psi}_i \otimes \mathbf{I}_{d_m}) \mathbf{r}_i$ , we obtain the system in the theorem statement.

$\mathbf{M}_{\text{glob}} = (\Psi \Psi^\top) \otimes \mathbf{I}_{d_m} + \lambda \mathbf{I}$  is positive definite for  $\lambda > 0$ , hence invertible. The minimizer is unique.

Thus  $\boldsymbol{\theta}^* = \mathbf{M}_{\text{glob}}^{-1} \mathbf{b}_{\text{glob}}$ , and reshaping yields  $\hat{\Delta}^* = \text{unvec}(\boldsymbol{\theta}^*)$ .

## B.2.2 DETAILED DERIVATION OF THE SINGLE-EXPERT SUBPROBLEM

Starting from the projected MoE objective (Eqn. 8),

$$\min_{\{\hat{\Delta}_n\}_{n=1}^N} \sum_{i \in \mathcal{E}} \left\| \sum_{n=1}^N g_{i,n} (\mathbf{W}_n \mathbf{k}_{i,n} + \hat{\Delta}_n \tilde{\mathbf{k}}_{i,n}) - \mathbf{v}_i \right\|^2 + \lambda \sum_{n=1}^N \|\hat{\Delta}_n\|^2, \quad (34)$$

we apply block coordinate descent (BCD) over experts, updating one expert at a time and keeping the others fixed. For a fixed expert  $n$ , collect all terms that do *not* involve  $\hat{\Delta}_n$  into the external residual

$$\mathbf{r}_i^{(-n)} = \mathbf{v}_i - \sum_{\ell \neq n} g_{i,\ell} (\mathbf{W}_\ell \mathbf{k}_{i,\ell} + \hat{\Delta}_\ell \tilde{\mathbf{k}}_{i,\ell}), \quad (35)$$

and substitute Eqn. 35 back into Eqn. 34. The single-expert subproblem for  $\hat{\Delta}_n$  is the ridge-regularized least-squares

$$\min_{\hat{\Delta}_n} \sum_{i \in \mathcal{E}} \left\| \mathbf{r}_i^{(-n)} - g_{i,n} \hat{\Delta}_n \tilde{\mathbf{k}}_{i,n} \right\|^2 + \lambda \|\hat{\Delta}_n\|^2. \quad (36)$$

Introduce compact notation

$$\mathbf{X} = \hat{\Delta}_n, \quad \mathbf{x}_i = \tilde{\mathbf{k}}_{i,n}, \quad \mathbf{y}_i = \mathbf{r}_i^{(-n)}, \quad g_i = g_{i,n}, \quad (37)$$

so that

$$f(\mathbf{X}) = \sum_{i \in \mathcal{E}} \|\mathbf{y}_i - g_i \mathbf{X} \mathbf{x}_i\|^2 + \lambda \|\mathbf{X}\|^2. \quad (38)$$

Using the standard matrix derivative identity<sup>2</sup> yields

$$\nabla f(\mathbf{X}) = -2 \sum_i g_i \mathbf{y}_i \mathbf{x}_i^\top + 2 \sum_i g_i^2 \mathbf{X} \mathbf{x}_i \mathbf{x}_i^\top + 2\lambda \mathbf{X}. \quad (39)$$

<sup>2</sup>For vectors  $\mathbf{a}$ ,  $\mathbf{b}$  and matrix  $\mathbf{X}$ ,  $\partial \|\mathbf{a} - \mathbf{X} \mathbf{b}\|^2 / \partial \mathbf{X} = -2(\mathbf{a} - \mathbf{X} \mathbf{b}) \mathbf{b}^\top$ .

Setting the gradient to zero gives the normal equations

$$X \left( \sum_i g_i^2 x_i x_i^\top + \lambda I \right) = \sum_i g_i y_i x_i^\top. \quad (40)$$

Restoring expert-specific symbols, define

$$M_n \triangleq \sum_{i \in \mathcal{E}} g_{i,n}^2 \tilde{\mathbf{k}}_{i,n} \tilde{\mathbf{k}}_{i,n}^\top + \lambda I, \quad B_n \triangleq \sum_{i \in \mathcal{E}} g_{i,n} \mathbf{r}_i^{(-n)} \tilde{\mathbf{k}}_{i,n}^\top. \quad (41)$$

Then Eqn. 40 becomes  $\hat{\Delta}_n M_n = B_n$ , with the unique minimizer

$$\hat{\Delta}_n^* = B_n M_n^{-1}. \quad (42)$$

Since the actual expert update is parameterized via the projection  $\Delta_n = \hat{\Delta}_n \mathbf{P}_n$  (Sec. 4.2), the written update is

$$\Delta_n^* = \hat{\Delta}_n^* \mathbf{P}_n. \quad (43)$$

**Why  $M_n$  is invertible (positive definite).** By definition,

$$M_n = \underbrace{\sum_{i \in \mathcal{E}} g_{i,n}^2 \tilde{\mathbf{k}}_{i,n} \tilde{\mathbf{k}}_{i,n}^\top}_{\text{Gram matrix } G_n \succeq 0} + \lambda I_{d_k}. \quad (44)$$

For any nonzero  $z \in \mathbb{R}^{d_k}$ ,

$$z^\top M_n z = \sum_{i \in \mathcal{E}} g_{i,n}^2 (z^\top \tilde{\mathbf{k}}_{i,n})^2 + \lambda \|z\|^2. \quad (45)$$

Since  $\lambda > 0$ , then  $z^\top M_n z \geq \lambda \|z\|^2 > 0$  for all  $z \neq 0$ , so  $M_n \succ 0$  and is invertible. Moreover,  $\lambda_{\min}(M_n) \geq \lambda$ , which ensures good conditioning (Tikhonov regularization). Because the projected keys are  $\tilde{\mathbf{k}}_{i,n} = \mathbf{P}_n \mathbf{k}_{i,n}$  with an idempotent projector  $\mathbf{P}_n$ , they lie in  $\text{range}(\mathbf{P}_n)$ . When  $\lambda > 0$ ,  $M_n$  remains strictly positive definite *on the full ambient space* (not only on  $\text{range}(\mathbf{P}_n)$ ), guaranteeing a unique closed-form update 42.

## C EXPERIMENTS SETUP

**Model Configuration** We evaluate our method on two Mixture-of-Experts (MoE) models: Qwen3-30B-A3B<sup>3</sup> and GPT-OSS-20B<sup>4</sup>. Qwen3-30B-A3B contains 128 experts per layer with the top-8 experts activated per token, resulting in approximately 3.3B active parameters during inference. The latter, GPT-OSS-20B, features 32 experts per layer with the top-4 experts activated, corresponding to approximately 3.6B active parameters.

**Hardware and Quantization** All experiments are conducted on a single node equipped with an NVIDIA H20 GPU. To balance precision and memory constraints, we utilize the BF16 format for model weights, while optimization is performed in FP32 to ensure numerical stability.

**Fine-Tuning (FT)** We evaluate both standard Fine-Tuning (FT) and Constrained Fine-Tuning (FT-L). The primary distinction is that FT-L imposes a norm constraint  $\varepsilon$  on the weight update. For both Qwen3-30B-A3B and GPT-OSS-20B, we set  $\varepsilon = 1 \times 10^{-3}$  for FT-L. We adopt a learning rate of  $1 \times 10^{-3}$  for both models. Updates are applied to layer 30 for Qwen3-30B-A3B and layer 0 for GPT-OSS-20B. For both methods, we target the `mlp.experts.down_proj` module. We train for 25 epochs, setting both weight decay and the KL divergence factor to 0.

**UnKE** As UnKE employs a two-stage structuring process, we configure the models as follows: For Qwen3-30B-A3B, the first stage uses a learning rate of  $5 \times 10^{-1}$  with 25 optimization steps and a weight decay coefficient of  $1 \times 10^{-3}$ . In the second stage, we apply a learning rate of  $2 \times 10^{-4}$  and perform 50 optimization steps. For GPT-OSS-20B, the first stage similarly adopts a learning rate of  $5 \times 10^{-1}$  but with 50 optimization steps, utilizing the same weight decay ( $1 \times 10^{-3}$ ). The second

<sup>3</sup><https://huggingface.co/Qwen/Qwen3-30B-A3B>

<sup>4</sup><https://huggingface.co/openai/gpt-oss-20b>

---

stage proceeds with a learning rate of  $1 \times 10^{-4}$  and 50 optimization steps. All experiments restrict parameter updates to layer 7. Consistent with the focus on structured knowledge editing, optimization is performed on the last subject token for both models.

**AdaLoRA** For AdaLoRA, updates are applied across all layers. We set the hyperparameters  $\alpha = 32$  and rank  $r = 8$ . For Qwen3-30B-A3B, we set the learning rate to  $5 \times 10^{-3}$ , while for GPT-OSS-20B, we use a reduced learning rate of  $5 \times 10^{-4}$ . The optimization is run for 25 steps for both models.

**MoEEdit (Ours)** For Qwen3-30B-A3B, we edit layers  $\{3, 4, 5, 6, 7\}$ , whereas for GPT-OSS-20B, we target layer 5. For Qwen3-30B-A3B, we perform 25 optimization steps with a learning rate of 0.1 and execute 4 Block Coordinate Descent (BCD) passes. For GPT-OSS-20B, we perform 50 optimization steps with a learning rate of 0.2 and 10 BCD passes. For both models, we set the regularization parameter  $\lambda = 1$  and the KL factor to 0.0625. We utilize 100,000 samples to compute the covariance matrix for the null-space projection with projection threshold = 0.02.

## D LLM USAGE DISCLOSURE

In accordance with the ICLR policy on responsible LLM usage, we hereby declare that Large Language Models (LLMs) were used solely for language refinement purposes in this paper. Specifically, LLMs were employed to correct grammar, improve clarity, and polish the writing style of the manuscript. No LLMs were used for generating ideas, designing methods, conducting experiments, analyzing results, or drawing conclusions. All scientific contributions of this work are entirely original and the responsibility of the authors.

## E EXAMPLES OF ZSRE AND COUNTERFACT

---

```

{
  "subject": "Watts Humphrey",
  "src": "What university did Watts Humphrey attend?",
  "pred": "Trinity College",
  "rephrase": "What university did Watts Humphrey take part in?",
  "alt": "University of Michigan",
  "answers": [
    "Illinois Institute of Technology"
  ],
  "loc": "nq question: who played desmond doss father in hacksaw ridge",
  "loc_ans": "Hugo Weaving",
  "cond": "Trinity College >> University of Michigan || What university did
    Watts Humphrey attend?"
},
{
  "subject": "Ramalinaceae",
  "src": "Which family does Ramalinaceae belong to?",
  "pred": "Ramalinales",
  "rephrase": "What family are Ramalinaceae?",
  "alt": "Lamiinae",
  "answers": [
    "Lecanorales"
  ],
  "loc": "nq question: types of skiing in the winter olympics 2018",
  "loc_ans": "Downhill",
  "cond": "Ramalinales >> Lamiinae || Which family does Ramalinaceae belong
    to?"
},
{
  "subject": "Denny Herzig",
  "src": "What role does Denny Herzig play in football?",
  "pred": "midfielder",
  "rephrase": "What's Denny Herzig's role in football?",
  "alt": "winger",
  "answers": [
    "defender"
  ],
  "loc": "nq question: where does aarp fall on the political spectrum",
  "loc_ans": "non-partisan",
  "cond": "midfielder >> winger || What role does Denny Herzig play in
    football?"
}

```

Figure 5: Examples of ZsRE dataset

---

```

972
973
974
975
976
977
978
979
980
981
982 {
983     "case_id": 975,
984     "pararel_idx": 17275,
985     "requested_rewrite": {
986         "prompt": "{}, from",
987         "relation_id": "P127",
988         "target_new": {
989             "str": "Google",
990             "id": "Q95"
991         },
992         "target_true": {
993             "str": "Microsoft",
994             "id": "Q2283"
995         },
996         "subject": "Bing Videos"
997     },
998     "paraphrase_prompts": [
999         "\"Old Jennifer: I'm $adjectiveOld!\" Bing Videos is owned by",
1000         "J. Bing Videos is from"
1001     ],
1002     "neighborhood_prompts": [
1003         "OneDrive is from",
1004         "German Research Center for Artificial Intelligence's owner",
1005         "Groove Music's owner",
1006         "Arkane Studios, from",
1007         "Yammer is from",
1008         "Yammer, by",
1009         "Turn 10 Studios, by",
1010         "German Research Center for Artificial Intelligence is owned by",
1011         "Mojang Studios is from",
1012         "id Software's owner"
1013     ]
1014 }
1015

```

Figure 6: An example of COUNTERFACT