

Knowledge Tracing and Editing in Language Models: Misalignment, Analysis, and a Roadmap

Anonymous ACL submission

Abstract

Large language models encode substantial factual knowledge in their parameters, but this knowledge is difficult to inspect and to update reliably. Two closely related research directions address this challenge: knowledge tracing, which aims to reveal how specific knowledge is represented and used inside the model, and knowledge editing, which aims to modify targeted facts while preserving general capabilities. This survey examines how these areas connect in practice and argues that the link is often unreliable and misaligned. We organize major tracing and editing approaches through an intervention oriented lens, synthesize evidence on where the connection breaks, and outline directions for tighter integration between explanation and intervention in future systems.

1 Introduction

Language models encode vast amounts of factual and relational knowledge implicitly within their parameters, enabling strong performance on knowledge-intensive tasks (Pan et al., 2023; Chen et al., 2022; Hu et al., 2023; He et al., 2023). However, parametric knowledge is difficult to inspect, localize, and reliably update, limiting transparency and controllability (Shi et al., 2025b; Singh et al., 2024). These challenges are particularly critical in high-stakes domains such as healthcare, law, and education, where incorrect or outdated knowledge can undermine trust and safety (Wang et al., 2025d; Chu et al., 2025; Zhang et al., 2025). Therefore, how to construct trustworthy and reliable models is a key issue in the current era of large models.

In response, two closely related research directions have emerged: *knowledge tracing* (KT) and *knowledge editing* (KE). Knowledge tracing (Nostalgebraist, 2020; Abnar and Zuidema, 2020) aims to diagnose the black box of LLMs by uncovering how specific pieces of knowledge are represented, processed and causally utilized within model archi-

tectures. Knowledge editing (Meng et al., 2022; Mitchell et al., 2022a) acts as a prescriptive intervention: it seeks to modify targeted factual associations (e.g., correcting outdated statistics, removing misinformation) without overhauling the model’s pre-trained parameters or general capabilities. Conceptually, these two fields are inherently complementary: tracing provides the “map” of where knowledge resides, and editing uses that map to execute precise “repairs”.

Yet in practice, this synergy remains fragile. A persistent misalignment separates tracing and editing: *editing methods frequently rely on localization assumptions that are only partially supported by tracing, while tracing research rarely incorporates feedback from editing’s real interventions*. This disconnect has tangible consequences: conclusions from knowledge tracing are often task-specific, architecture-dependent, and difficult to transfer to editing scenarios. Conversely, knowledge editing faces pervasive side effects, including interference with related facts, degradation of general capabilities, and instability under repeated or large-scale updates (Thede et al., 2025; Gu et al., 2024; Yan et al., 2025; He et al., 2025; Huang et al., 2024; Li et al., 2024a), precisely because tracing insights fail to guide robust intervention design.

Prior surveys (Wang et al., 2025d, 2024b; Mazzia et al., 2025; Yao et al., 2023; Zhang et al., 2024b; Li et al., 2024c; Hase et al., 2024) have summarized advances in tracing or editing in isolation, but none have centered on the gap between explanation (tracing) and intervention (editing) as a unifying lens. In this survey, we treat KT and KE as stages of a *KT→KE pipeline*, as illustrated in Figure 1, which fails at three points. First, there is an *objective mismatch*: tracing highlights sensitivity and functional participation, whereas editing requires stable controllability under intervention. Second, there is a *mechanism mismatch*: editing can re-route or distort the computations that trac-

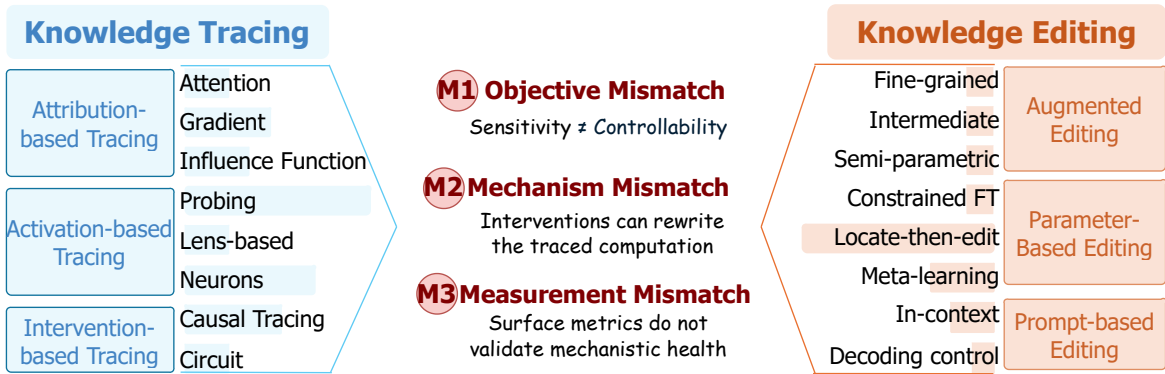


Figure 1: Overview of Mismatches in Knowledge Tracing (KT) and Knowledge Editing (KE) in LLMs. The blue bars indicate how many Editing paper cite each type of Tracing and the orange bars represent the reverse.

ing identifies, even when behavioral edit success appears high. Third, there is a *measurement mismatch*: standard behavioral metrics can miss latent mechanistic degradation and therefore cannot close the loop between tracing and editing.

Concretely, we make three contributions. First, we reorganize tracing methods by what they provide for intervention and categorize editing methods by their *control regimes*. Second, we formalize the above three mismatch points as a minimal typology for the KT-KE interface. Third, we provide complementary micro- and macro-level evidence showing that traced targets can be reproducible yet offer weak control leverage, and that behaviorally successful edits can still induce mechanistic drift. Together, these results motivate a closed-loop research agenda where editing serves as an interventional validator for tracing, and tracing serves as a mechanistic constraint for editing.

Paper Roadmap. Section 2 reviews KT and KE methods, and clarifies the *KT*→*KE* interface: what tracing provides versus what editing requires. Section 3.1 formalizes three recurring failure points in that interface, and Sections 3.2–3.4 ground them with micro-level mechanistic diagnostics and macro-level ecosystem evidence. Finally, we outline a roadmap toward a closed-loop synergy between tracing and editing.

2 Knowledge Tracing and Editing: Methods and Current Practice

This section reviews recent advances in knowledge tracing (Sec. 2.1) and knowledge editing (Sec. 2.2), and summarizes the current practice from an intervention-oriented perspective. Figure 10 in Appendix F provides an architectural

landscape of the two lines of work. Throughout, we emphasize what each method family *outputs* in practice (e.g., candidate components, causal effect maps, update operators), as these outputs form the implicit interface that later sections examine.

2.1 Knowledge Tracing

Knowledge tracing seeks to pinpoint *where* specific information is encoded in the model and *how* such information is processed to produce a prediction. Operationally, tracing methods differ mainly in the type of evidence they produce and the level at which they localize knowledge. We group them into *attribution-based*, *activation-based*, and *intervention-based* methods.

Attribution-based Methods. Attribution methods quantify how input features affect model predictions, producing token- or span-level importance scores in the *input space*. Early work used *attention-based attribution* (Abnar and Zuidema, 2020; Clark et al., 2019; Ayyar et al., 2025), treating attention weights as token-importance scores, but later studies showed that attention is not inherently interpretable and can vary without changing outputs (Jain and Wallace, 2019; Liu et al., 2022; Pruthi et al., 2020). This motivated *gradient-based approaches* (Nielsen et al., 2022; Ancona et al., 2018; Smilkov et al., 2017), which use local derivatives as feature attributions but suffer from saturation, motivating path-integral methods such as Integrated Gradients (Sundararajan et al., 2017; Kapishnikov et al., 2021; Lundstrom and Razaviyayn, 2025). At the data level, *influence-function* (Koh and Liang, 2017; Koh et al., 2019) estimates how upweighting a training point changes a test prediction, but recent work (Li et al., 2025b) finds that

classical influence approximations break down for modern LLMs. In practice, attribution methods are lightweight and widely applicable, typically providing saliency-style evidence that can serve as a first-pass hypothesis for deeper localization.

Activation-based Methods. Activation-based approaches localize knowledge in internal representations, shifting from important *inputs* to informative *hidden states*. *Probing methods* (Conneau et al., 2018; Tenney et al., 2019; Belinkov, 2022) use auxiliary classifiers to decode conceptual attributes from hidden states; linear probes (Alain and Bengio, 2016; Dong et al., 2025) are simple and interpretable, while nonlinear probes (Gairola et al., 2025; Hewitt and Liang, 2019) add capacity but risk overfitting. *Lens-based methods* decode intermediate activations via Logit Lens (Nostalgebraist, 2020; Neo et al., 2025), Tuned Lens (Belrose et al., 2023), Entropy Lens (Ali et al., 2025), Backward Lens (Katz et al., 2024) and extensions to encoder–decoders and vision models (Langedijk et al., 2024; Takatsuki et al., 2025) to trace how predictions evolve across layers. *Knowledge-neuron methods* (Dai et al., 2022; Wang et al., 2024d; Chen et al., 2025b) attempt to link neurons to specific facts but face polysemanticity (Mu and Andreas, 2020), motivating *sparse autoencoders* (Huben et al., 2024; Braun et al., 2024; Bricken et al., 2023) which produce monosemantic features and consistent concept dictionaries (Fel et al., 2025; Shi et al., 2025a; Arad et al., 2025; Kang et al., 2025). In practice, activation-based methods often output *candidate layers/neurons/features* and *representation-level trajectories* that suggest where the model encodes and transforms knowledge, and are used as inputs to downstream interventions.

Intervention-based Methods. Intervention-based tracing treats neural networks as causal systems and tests whether internal states are necessary or sufficient for a behavior via controlled manipulations (Vig et al., 2020; Geiger et al., 2025). *Causal tracing* (activation patching (Meng et al., 2022; Vig et al., 2020; Zhang and Nanda, 2024; Dumas et al., 2025) compares clean and corrupted runs to locate states whose restoration fixes model outputs. Extending this idea, *circuit methods* (Dumas et al., 2025; Conmy et al., 2023; Hsu et al., 2025; Haklay et al., 2025; Lan et al., 2024) isolate minimal head/MLP subgraphs underlying behaviors, and recent work (Yao et al., 2024; Merullo et al., 2024;

Lan et al., 2024) shows that models often reuse shared subcircuits. Theoretical work on *causal abstraction* (Geiger et al., 2021, 2025; Hu and Tian, 2022; Geiger et al., 2022, 2024; Wu et al., 2023; Tan, 2023) formalizes these ideas, while recent results (Sutter et al., 2025) highlight that unconstrained abstraction maps can trivialize such claims, motivating stricter structural assumptions. In practice, intervention-based methods output *causal effect maps* and *causal hypotheses about which modules carry which information*, and they provide the closest-to-interventional evidence that is often assumed to be actionable for editing.

2.2 Knowledge Editing

Knowledge editing aims to update a model by inserting, correcting, or removing specific facts without retraining, ideally achieving reliability, locality, generality, and scalability to many edits. Editing approaches differ by where the update is stored and how it is applied during inference. We group them into parameter-based, augmented, and prompt-based editing.

Parameter-based Editing. Parameter-based methods treat pretrained weights as the main storage of knowledge and rewrite them via targeted updates. *Constrained finetuning* (Gekhman et al., 2024; Zhu et al., 2020; Li et al., 2024b) applies regularized updates to specific layers. It suits narrow edit sets by not assuming fact localization, though broad gradient propagation often introduces interference across numerous edits. *Locate-then-edit* approaches (Meng et al., 2022; Wang et al., 2025b; Fang et al., 2025; Meng et al., 2023; Li et al., 2024b; Gupta et al., 2024) instead identify the layers most tied to a fact and modify only those components. This improves locality and interpretability but is brittle when knowledge is distributed or entangled. *Meta-learning editors* (Mitchell et al., 2022a; Tan et al., 2024; De Cao et al., 2021) learn to predict weight updates from an edit specification, enabling fast and scalable edits but adding an auxiliary model whose reliability depends on distributional coverage. In practice, parameter-based editing offers persistent, low-latency changes without extra inference machinery, but it requires white-box access, is prone to interference at scale, and has limited reversibility.

Augmented Editing. Augmented approaches freeze backbone weights and localize updates in

254 auxiliary mechanisms that can be enabled, dis- 254
 255 abled, or replaced. *Fine-grained methods* (Wang 255
 256 et al., 2025c; Huang et al., 2023; Bayat et al., 256
 257 2025; Zhao et al., 2025; QIU et al., 2024; Lee 257
 258 et al., 2025; Kim et al., 2018; Subramani et al., 258
 259 2022; Pascual et al., 2021; Rimsky et al., 2024; 259
 260 Li et al., 2023; Konen et al., 2024; Stolfo et al., 260
 261 2025; Todd et al., 2024) steer individual neurons 261
 262 or feature directions, enabling highly targeted edits 262
 263 but requiring sparse, monosemantic feature spaces 263
 264 (e.g., SAEs) to avoid unintended effects. *Interme-* 264
 265 *diate approaches* (Yu et al., 2024; Li et al., 2025a; 265
 266 Wang and Li, 2024) attach trainable, parameter- 266
 267 efficient modules (e.g., adapters) to frozen layers, 267
 268 allowing modular, context-dependent updates at the 268
 269 cost of added inference overhead. *Coarse-grained* 269
 270 *semi-parametric methods* (Mitchell et al., 2022b; 270
 271 Hartvigsen et al., 2023; Wang et al., 2024a; Zhu 271
 272 et al., 2025; Wang et al., 2025a) store mutable 272
 273 knowledge in external key-value memories that re- 273
 274 trieve and override internal representations during 274
 275 inference, improving transparency and reversibility 275
 276 while increasing system complexity. Overall, aug- 276
 277 mented editing trades architectural simplicity for 277
 278 modularity, scope control, and easier rollback. 278

279 **Prompt-based Editing.** Prompt-based methods 279
 280 keep parameters fixed and apply edits through 280
 281 inputs or decoding control. *In-context meth-* 281
 282 *ods* (Zheng et al., 2023; Qi et al., 2025; Qiao 282
 283 et al., 2024; Wang et al., 2024c) supply edited 283
 284 facts as examples, instructions, or rationales, al- 284
 285 lowing the model to adopt updated rules on the 285
 286 fly. *Prompt-level memory* (Wang et al., 2024c; 286
 287 Chen et al., 2025a) stores corrections externally 287
 288 and retrieves them into future prompts. These ed- 288
 289 its are lightweight but ephemeral: they work only 289
 290 when included in the prompt and consume context 290
 291 length. A complementary direction frames edit- 291
 292 ing as *decoding control* (Wang et al., 2024e; Sun 292
 293 et al., 2024), constraining generation so outputs 293
 294 obey specified knowledge constraints, for exam- 294
 295 ple by down-weighting obsolete facts (Sun et al., 295
 296 2024) by pruning reasoning paths that violate a 296
 297 newly introduced rule. (Wang et al., 2024e). Over- 297
 298 all, prompt-based editing is model-agnostic (works 298
 299 for black-box APIs), cheap per edit (no optimiza- 299
 300 tion or new parameters), and highly flexible. In 300
 301 practice, prompt-based editing is model-agnostic 301
 302 (works for black-box APIs), but it is not persistent 302
 303 by default and typically offers weaker locality and 303
 304 generalization guarantees. 304

2.3 From Tracing to Editing: What Tracing 305 Provides and What Editing Requires 306

307 Rather than treating knowledge tracing (KT) and 307
 308 knowledge editing (KE) as two independent toolk- 308
 309 its, current practice implicitly forms a pipeline: KT 309
 310 proposes *targets* (where/what to intervene on), KE 310
 311 applies an *operator* (how to change the model), 311
 312 and evaluation determines whether the update is 312
 313 *acceptable* (what counts as success). Making it ex- 313
 314 plicit clarifies why a method can be diagnostically 314
 315 informative yet operationally unreliable. 315

316 **What tracing typically provides.** Tracing meth- 316
 317 ods predominantly deliver *sensitivity evidence* 317
 318 about a model’s computation on a given prompt. 318
 319 This includes saliency or attribution signals over 319
 320 tokens, layers, neurons, or learned features; ranked 320
 321 candidate components (e.g., attention heads, MLP 321
 322 blocks, feature directions) that appear implicated 322
 323 in producing the output; and, for intervention- 323
 324 based tracing, causal effect maps from controlled 324
 325 activation interventions (e.g., patching) that indi- 325
 326 cate where restoring information repairs a behav- 326
 327 ior. These deliverables are valuable for identifying 327
 328 *where the computation is fragile or influential*, but 328
 329 they are not designed to guarantee that the high- 329
 330 lighted units are stable *control handles* under sub- 330
 331 sequent interventions. 331

332 **What editing requires.** Editing methods, in con- 332
 333 trast, require *intervention-level guarantees*. First, 333
 334 the chosen target must be *controllable*: interven- 334
 335 ing on it should reliably move the model toward 335
 336 the intended factual state across prompts and con- 336
 337 texts, not merely perturb outputs. Second, edits 337
 338 must satisfy *locality and robustness*: they should 338
 339 minimize unintended changes to related facts and 339
 340 general capabilities, and remain stable under re- 340
 341 peated updates or large edit sets. Third, editing 341
 342 demands *validation criteria* that reflect not only 342
 343 surface behavior (e.g., edit success and locality 343
 344 on held-out prompts), but also whether the update 344
 345 preserves intended computation when mechanistic 345
 346 drift is a concern (e.g., avoiding shortcut formation 346
 347 or contextual collapse). 347

348 **Current practice and implicit interface assump-** 348
 349 **tions.** Many locate-then-edit workflows treat trac- 349
 350 ing outputs as actionable coordinates for editing, 350
 351 implicitly assuming that (i) localized components 351
 352 provide stable leverage across prompts and con- 352
 353 texts, (ii) the edit operator modifies knowledge 353
 354 without substantially rerouting computation, and 354

Seq Setting		Sin Setting	
Neurons	Freq.	Neurons	Freq.
(30, 11268)	24	(30, 11268)	24
(30, 11430)	20	(29, 8499)	24
(30, 11920)	18	(28, 2856)	23
(30, 10321)	15	(30, 7752)	20
(31, 7096)	12	(30, 2458)	19
(31, 4964)	10	(30, 14069)	18
(31, 12111)	10	(31, 12732)	17
(31, 5463)	10	(30, 1441)	17
(31, 154)	9	(30, 13883)	17
(5, 7012)	9	(30, 8061)	17

Table 1: Top recurrent neurons identified, test on 50 data. Seq reflects coupled editing-tracing stability, while Sin indicates the intrinsic importance of neurons across independent data samples.

(iii) standard behavioral metrics are sufficient to validate the update. Section 3.1 examines these assumptions by a pipeline view of *target selection*, *intervention*, and *validation*, and characterizes where and why the tracing→editing connection can fail.

3 Misalignment Between Knowledge Tracing and Editing

The preceding sections established a conceptual complementarity: knowledge tracing diagnoses *where* and *how* a fact participates in model computation, whereas knowledge editing aims to *control* and *update* that knowledge through interventions. Yet in current practice, the KT→KE connection is often treated as a plug-and-play interface: traced components are used as edit coordinates, edits are applied with a chosen operator, and success is judged mostly by behavioral metrics. In this section, we show that the interface breaks in systematic ways, using targeted micro-level diagnostics and macro-level ecosystem evidence.¹

3.1 The tracing to editing pipeline and three mismatches

We organize the tracing–editing gap by *where* it occurs in a standard KT→KE workflow:



This pipeline view turns “misalignment” into three concrete breakpoints (Figure 1):

¹The papers we collected (130 Tracing / 90 Editing) are primarily from NeurIPS, ICLR, ICML, ACL, AAAI, and EMNLP (2022–2025). Appendix A provides collection details.

Metric	Sequential	Single
Mean Δ gt_prob	-2.37×10^{-7}	-9.19×10^{-8}
Positive Δ gt_prob Ratio	0.44	0.42
Argmax Flip Rate	0.22	0.14
Avg #Edited Neurons	11.24	12.26

Table 2: Comparison of sequential (Seq) and Single (Sin) neuron tracing under knowledge editing. Negative Δ gt_prob and weak correlations indicate a misalignment between traced neuron importance and effective knowledge editing.

Objective mismatch at target selection: sensitivity \neq controllability. Tracing is diagnostic: it identifies components that *participate* in producing an output (sensitivity). Editing is prescriptive: it requires components that provide stable *leverage* under intervention (controllability). Thus, a component can be causally involved yet still be a poor control handle.

Mechanism mismatch at intervention: edits can rewrite traced computation. Even when tracing isolates a causal pathway for a fact, an edit may introduce a shortcut that reroutes inference (*causal reinforcement*), or gradually erode contextual reliance on unrelated inputs (*causal collapse*). Hence, behavioral success on edited prompts can coexist with mechanistic drift.

Measurement mismatch at validation: surface metrics \neq mechanistic health. Standard editing metrics (e.g., edit success, locality) quantify external behaviors but do not verify whether internal inference pathways remain intact. Conversely, tracing outputs are rarely validated by whether interventions at traced coordinates yield the intended change without side effects. This missing mutual validation keeps the KT→KE loop open.

The following subsections ground each breakpoint with representative micro-level diagnostics and macro-level ecosystem evidence.

3.2 Objective mismatch: sensitivity is not a control handle

A common practice in locate-then-edit pipelines is to treat traced units as actionable edit coordinates. However, tracing outputs are designed to reflect *participation* in the current computation, whereas editing requires *stable leverage* under intervention. We illustrate this mismatch with a neuron-level diagnostic case study.

Neuron-level diagnostic (Knowledge Neurons). We use LLAMA3-CHAT-8B and sample 50 fac-

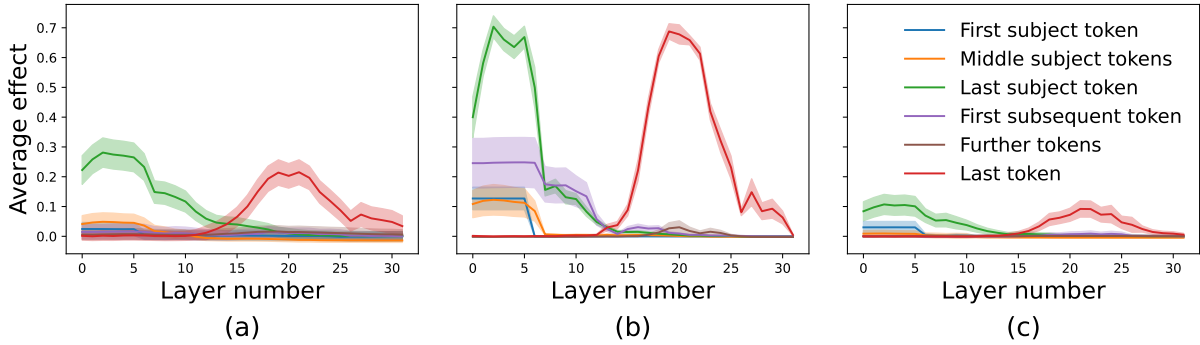


Figure 2: Average effect of MLP blocks. (a) Before Editing, Causal trace on 100 Matched-Data. (b) Post-Edit model with editing on 100 Matched-Data. (c) Post-Edit model with editing on 100 Independent-Data. More Post-Edit model results shows similar trend. Full visualization are provided in Figure 5 – Figure 9 in Appendix E.

421 tual prompts from KNOWN (Meng et al., 2022).
 422 We apply Knowledge Neurons (KN; Dai et al.,
 423 2022) to identify top-ranked MLP neurons asso-
 424 ciated with the ground-truth answer probability
 425 via integrated gradients, and intervene on the ac-
 426 tivations of these neurons during generation. To
 427 distinguish intrinsic recurrence from intervention-
 428 coupled effects, we compare **Single (Sin)** (trace/in-
 429 tervene independently per example) and **Sequen-**
 430 **tial (Seq)** (accumulate interventions over a stream
 431 of examples). We measure mean Δp_{gt} (average
 432 change in the model-assigned probability of the
 433 ground-truth token), Pos- Δp_{gt} ratio (fraction of
 434 cases with $\Delta p_{gt} > 0$), argmax flip rate (fraction of
 435 cases whose top-1 prediction changes after inter-
 436 vention), and the number of edited neurons (inter-
 437 vention scope per instance).

438 **Observation.** The traced neurons are highly re-
 439 current and concentrate in late decoding layers (Ta-
 440 ble 1), suggesting non-trivial localization rather
 441 than randomness. Yet intervening on these neu-
 442 rons does not translate into reliable factual gains:
 443 mean Δp_{gt} is near-zero and slightly negative, and
 444 Pos- Δp_{gt} remains below 50% under both regimes
 445 (Table 2). Meanwhile, Seq increases behavioral
 446 sensitivity (higher argmax flip rate than Sin), indi-
 447 cating that traced units can affect *what the model*
 448 *says* without providing stable control over *which*
 449 *fact is expressed*.

450 **Implication.** Neuron-level tracing can be repro-
 451 ducible and behaviorally influential, but repro-
 452 ducibility and sensitivity alone are insufficient for
 453 selecting robust edit coordinates. This is a concrete
 454 instance of the objective mismatch at the target-
 455 selection breakpoint.

3.3 Mechanism mismatch: edits can reroute the traced computation

456 Even when editing achieves high behavioral suc-
 457 cess, the post-edit model may rely on different in-
 458 ternal pathways than those suggested by pre-edit
 459 tracing. This section uses causal tracing (Meng
 460 et al., 2022) to show that edits can introduce short-
 461 cut reinforcement on edited data and collapse con-
 462 textual reliance on unrelated inputs. 463
 464

Causal diagnostic (AlphaEdit + causal tracing).

465 We fix a trace set \mathcal{T} of 100 samples for causal-
 466 tracing analysis and construct edit sets \mathcal{E} with
 467 $|\mathcal{E}| \in \{100, 300, 500\}$ for intervention and behav-
 468 ioral evaluation. We use LLAMA3-CHAT-8B and
 469 apply AlphaEdit (Fang et al., 2025) (The Sota Edit-
 470 ing method) and compute causal effect maps via
 471 activation patching following Meng et al. (2022).
 472 We consider **Matched-Data (MD)** where $\mathcal{E} = \mathcal{T}$
 473 and **Independent-Data (ID)** where $\mathcal{E} \cap \mathcal{T} = \emptyset$,
 474 while causal traces are always computed on \mathcal{T} . 475

476 **Observation.** Under Matched-Data, post-edit
 477 causal maps show an amplified influence of the **first**
 478 **subsequent token** (Figure 2a vs. Figure 2b), con-
 479 sistent with a strong injected signal that propagates
 480 forward and stabilizes the edited output via a short-
 481 cut (*causal reinforcement*). Under Independent-
 482 Data, causal influence shifts toward over-reliance
 483 on the **last subject token** while earlier context
 484 contributions diminish (Figure 2c), yielding a sys-
 485 tematic narrowing of causal support that we term
 486 **causal collapse**; the pattern strengthens as edit
 487 scale increases (Figures 7–9 in Appendix E).

488 **Implication.** Locate-then-edit interventions can
 489 succeed on edited prompts while reshaping infer-
 490 ence pathways, producing on-target shortcut re-

Setting	Data Scale	ES	LC
Matched-Data	100	-	0.640
E-100	100	0.192	0.480
E-300	300	0.294	0.160
E-500	500	0.488	0.220

Table 3: Editing success (ES) and locality scores (LC). Locality drastically degrades as the number of edited facts exceeds 100, indicating a loss of model robustness. E-300 means edit on 300 edits in Independent-Data and evaluate the locality on 100 Matched-Data.

liance and off-target contextual collapse. This reflects a mechanism-level mismatch between “what tracing identifies” and “what editing produces”.

3.4 Measurement mismatch: surface metrics lag behind mechanistic degradation

A third break point lies in *validation*. In most editing pipelines, success is primarily assessed by behavioral scores such as edit success and locality. At the same time, tracing outputs are rarely evaluated by whether interventions at traced coordinates reliably deliver the intended change with minimal side effects. As a result, the KT→KE loop often remains open: surface metrics may pass while mechanistic drift accumulates, and tracing evidence is not routinely stress-tested for editing. We illustrate this mismatch with a micro-level diagnostic and complementary macro-level ecosystem evidence.

Micro evidence (early warning beyond locality). Table 3 shows that locality can remain moderately high at E-100 (LC=0.480) even when causal tracing already reveals collapse tendencies (Figure 2c). As edits scale up, this latent mechanistic erosion reaches a tipping point where locality drops sharply (e.g., LC=0.160 at E-300). Thus, standard locality metrics can *lag* behind mechanism-level failure, whereas causal traces provide earlier signals of robustness loss.

Macro evidence (open-loop validation in the ecosystem). Publication and influence patterns suggest that validation remains largely open-loop across the KT/KE ecosystem (Figure 3). Tracing established an early foothold and shows steady growth (Figure 3a), whereas editing expands rapidly after 2022 across heterogeneous mechanisms. A quadrant analysis illustrated in Figure 3b using Relative Citation Index and Citation Velocity (detailed in Appendix C) further indicates an imbalance in sustained influence: tracing papers dom-

inate high-impact and high-momentum regions, while editing work is more concentrated in lower-momentum areas. This gap is reinforced by a tooling asymmetry: unified editing toolkits standardize workflows and make behavioral metrics easy to reproduce, whereas tracing tools remain fragmented and less plug-and-play, limiting systematic comparison and transfer across models and edited settings (Table B). Finally, cross-category citation flows (Figure 3c) reveal a strongly asymmetric dependency: editing heavily cites foundational tracing methods, while tracing only sparsely cites editing. Analysis of recent editing paper (Table 5) shows that Parameter-Based edits relies more on tracing, while the other two categories rely less, though tracing could still guide them which is a potential largely unexplored. Together, these patterns reinforce localization as the dominant mental model for knowledge manipulation, even when localization signals are not routinely validated for stability under intervention.

Implication. Both micro and macro evidence point to a validation gap: behavioral metrics do not reliably certify mechanistic health, and tracing outputs are rarely tested for edit-actionability, keeping the KT→KE pipeline open-loop.

3.5 Takeaway: sensitivity is not manipulability

Across the three pipeline breakpoints, a consistent message emerges: *tracing provides sensitivity signals about participation, but editing requires manipulable levers that preserve mechanistic health under intervention*. Closing the gap needs (i) tracing criteria that explicitly account for stability under intervention, (ii) editing objectives and evaluations that incorporate mechanistic validation rather than relying solely on surface success/locality.

4 Toward a Synergistic Roadmap

We outline four actionable directions to close the loop: (i) make tracing *edit-validated*, (ii) make editing *mechanism-preserving*, (iii) use tracing to *select* editing regimes, and (iv) enable *lifelong* updates with mechanistic monitoring.

(1) Edit-validated Tracing: from sensitivity maps to control handles. Addresses: target selection + validation. Tracing should be judged not only by reproducible localization, but by whether the identified units become *reliable control handles*

6 Limitation

This study examines the interplay between knowledge tracing and knowledge editing and discusses a closed-loop perspective between the two. Our analysis is biased toward how knowledge tracing can support editing, while the complementary role of editing in informing and refining tracing methods is not explored in depth. Moreover, our discussion focuses on transformer-based language models, extensions to alternative architectures such as Mamba and to multimodal settings are only briefly considered and remain open directions for future work.

References

- Samira Abnar and Willem Zuidema. 2020. [Quantifying attention flow in transformers](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4190–4197, Online. Association for Computational Linguistics.
- Guillaume Alain and Yoshua Bengio. 2016. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*.
- Riccardo Ali, Francesco Caso, Christopher Irwin, and Pietro Liò. 2025. [Entropy-lens: The information signature of transformer computations](#). *Preprint*, arXiv:2502.16570.
- Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. 2018. [Towards better understanding of gradient-based attribution methods for deep neural networks](#). In *International Conference on Learning Representations*.
- Dana Arad, Aaron Mueller, and Yonatan Belinkov. 2025. [SAEs are good for steering – if you select the right features](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 10252–10270, Suzhou, China. Association for Computational Linguistics.
- Meghna P Ayyar, Jenny Benois-Pineau, and Akka Zemari. 2025. There is more to attention: Statistical filtering enhances explanations in vision transformers. *arXiv preprint arXiv:2510.06070*.
- Reza Bayat, Ali Rahimi-Kalahroudi, Mohammad Pezeshki, Sarath Chandar, and Pascal Vincent. 2025. [Steering large language model activations in sparse spaces](#). In *Second Conference on Language Modeling*.
- Yonatan Belinkov. 2022. [Probing classifiers: Promises, shortcomings, and advances](#). *Computational Linguistics*, 48(1):207–219.
- Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. 2023. Eliciting latent

- predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112*.
- Baolong Bi, Shenghua Liu, Lingrui Mei, Yiwei Wang, Pengliang Ji, and Xueqi Cheng. 2024a. [Decoding by contrasting knowledge: Enhancing llms’ confidence on edited facts](#). *arXiv preprint arXiv:2405.11613*.
- Baolong Bi, Shenghua Liu, Yiwei Wang, Lingrui Mei, Hongcheng Gao, Yilong Xu, and Xueqi Cheng. 2024b. [Adaptive token biaser: Knowledge editing via biasing key entities](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 11071–11083, Miami, Florida, USA. Association for Computational Linguistics.
- Dan Braun, Jordan Taylor, Nicholas Goldowsky-Dill, and Lee Sharkey. 2024. [Identifying functionally important features with end-to-end sparse dictionary learning](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermy, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, and 6 others. 2023. [Towards monosemanticity: Decomposing language models with dictionary learning](#). *Transformer Circuits Thread*. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- Mingda Chen, Yang Li, Karthik Padthe, Rulin Shao, Alicia Yi Sun, Luke Zettlemoyer, Gargi Ghosh, and Wen-tau Yih. 2025a. [Improving factuality with explicit working memory](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11199–11213, Vienna, Austria. Association for Computational Linguistics.
- Yuheng Chen, Pengfei Cao, Yubo Chen, Yining Wang, Shengping Liu, Kang Liu, and Jun Zhao. 2025b. [Cracking factual knowledge: A comprehensive analysis of degenerate knowledge neurons in large language models](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10240–10261, Vienna, Austria. Association for Computational Linguistics.
- Zhuo Chen, Yufeng Huang, Jiaoyan Chen, Yuxia Geng, Yin Fang, Jeff Z. Pan, Ningyu Zhang, and Wen Zhang. 2022. [Lako: Knowledge-driven Visual Question Answering via Late Knowledge-to-text Injection](#). In *Proc. of IJCKG2022*, pages 20–29.
- Zhendong Chu, Shen Wang, Jian Xie, Tinghui Zhu, Yibo Yan, Jinheng Ye, Aoxiao Zhong, Xuming Hu, Jing Liang, Philip S Yu, and 1 others. 2025. [Llm agents for education: Advances and applications](#). *arXiv preprint arXiv:2503.11733*.

754	Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? an analysis of BERT’s attention . In <i>Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 276–286, Florence, Italy. Association for Computational Linguistics.	811
755		812
756		813
757		814
758		815
759		816
760		817
761	Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. 2023. Towards automated circuit discovery for mechanistic interpretability. <i>Advances in Neural Information Processing Systems</i> , 36:16318–16352.	818
762		819
763		820
764		821
765		822
766	Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. What you can cram into a single vector: Probing sentence embeddings for linguistic properties. <i>arXiv preprint arXiv:1805.01070</i> .	823
767		824
768		825
769		826
770		827
771	Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge neurons in pretrained transformers . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.	828
772		829
773		830
774		831
775		832
776		833
777		834
778	Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6491–6506, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	835
779		836
780		837
781		838
782		839
783		840
784	Jianshuo Dong, Yutong Zhang, Liu Yan, Zhenyu Zhong, Tao Wei, Ke Xu, Minlie Huang, Chao Zhang, and Han Qiu. 2025. “I’ve decided to leak”: Probing internals behind prompt leakage intents . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 21329–21359, Suzhou, China. Association for Computational Linguistics.	841
785		842
786		843
787		844
788		845
789		846
790		847
791		848
792	Yiming Du, Wenyu Huang, Danna Zheng, Zhaowei Wang, Sebastian Montella, Mirella Lapata, Kam-Fai Wong, and Jeff Z. Pan. 2025. Rethinking memory in llm based agents: Representations, operations, and emerging topics . <i>Preprint</i> , arXiv:2505.00675.	849
793		850
794		851
795		852
796		853
797	Clément Dumas, Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2025. Separating tongue from thought: Activation patching reveals language-agnostic concept representations in transformers . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 31822–31841, Vienna, Austria. Association for Computational Linguistics.	854
798		855
799		856
800		857
801		858
802		859
803		860
804		861
805	Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Jie Shi, Xiang Wang, Xiangnan He, and Tat-Seng Chua. 2025. Alphaedit: Null-space constrained model editing for language models . In <i>The Thirteenth International Conference on Learning Representations</i> .	862
806		863
807		864
808		865
809		866
810		
	Thomas Fel, Ekdeep Singh Lubana, Jacob S. Prince, Matthew Kowal, Victor Boutin, Isabel Papadimitriou, Binxu Wang, Martin Wattenberg, Demba E. Ba, and Talia Konkle. 2025. Archetypal SAE: Adaptive and stable dictionary learning for concept extraction in large vision models . In <i>Forty-second International Conference on Machine Learning</i> .	811
		812
		813
		814
		815
		816
		817
	Siddhartha Gairola, Moritz Böhle, Francesco Locatello, and Bernt Schiele. 2025. How to probe: Simple yet effective techniques for improving post-hoc explanations . In <i>The Thirteenth International Conference on Learning Representations</i> .	818
		819
		820
		821
		822
	Atticus Geiger, Duligur Ibeling, Amir Zur, Maheep Chaudhary, Sonakshi Chauhan, Jing Huang, Aryaman Arora, Zhengxuan Wu, Noah Goodman, Christopher Potts, and 1 others. 2025. Causal abstraction: A theoretical foundation for mechanistic interpretability. <i>Journal of Machine Learning Research</i> , 26(83):1–64.	823
		824
		825
		826
		827
		828
	Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. 2021. Causal abstractions of neural networks . In <i>Advances in Neural Information Processing Systems</i> , volume 34, pages 9574–9586. Curran Associates, Inc.	829
		830
		831
		832
		833
	Atticus Geiger, Zhengxuan Wu, Hanson Lu, Josh Rozner, Elisa Kreiss, Thomas Icard, Noah Goodman, and Christopher Potts. 2022. Inducing causal structure for interpretable neural networks. In <i>International Conference on Machine Learning</i> , pages 7324–7338. PMLR.	834
		835
		836
		837
		838
		839
	Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. 2024. Finding alignments between interpretable causal variables and distributed neural representations. In <i>Causal Learning and Reasoning</i> , pages 160–187. PMLR.	840
		841
		842
		843
		844
	Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. Does fine-tuning LLMs on new knowledge encourage hallucinations? In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 7765–7784, Miami, Florida, USA. Association for Computational Linguistics.	845
		846
		847
		848
		849
		850
		851
	Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun Peng. 2024. Model editing harms general abilities of large language models: Regularization to the rescue. <i>arXiv preprint arXiv:2401.04700</i> .	852
		853
		854
		855
		856
	Akshat Gupta, Dev Sajani, and Gopala Anumanchipalli. 2024. A unified framework for model editing . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 15403–15418, Miami, Florida, USA. Association for Computational Linguistics.	857
		858
		859
		860
		861
		862
	Tal Haklay, Hadas Orgad, David Bau, Aaron Mueller, and Yonatan Belinkov. 2025. Position-aware automatic circuit discovery . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational</i>	863
		864
		865
		866

867		<i>Linguistics (Volume 1: Long Papers)</i> , pages 2792–2817, Vienna, Austria. Association for Computational Linguistics.		
868				922
869				923
870	Thomas Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023. Aging with GRACE: Lifelong model editing with discrete key-value adaptors . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .			924
871				925
872				926
873				927
874				928
875	Peter Hase, Thomas Hofweber, Xiang Zhou, Elias Stengel-Eskin, and Mohit Bansal. 2024. Fundamental problems with model editing: How should rational belief revision work in llms? <i>Preprint</i> , arXiv:2406.19354.			929
876				930
877				931
878				932
879				933
880	Guoxiu He, Xin Song, and Aixin Sun. 2025. Knowledge updating? no more model editing! just selective contextual reasoning. <i>arXiv preprint arXiv:2503.05212</i> .			934
881				935
882				936
883	Jie He, Víctor Gutiérrez-Basulto, and Jeff Z. Pan. 2023. BUCA: A Binary Classification Approach to Un-supervised Commonsense Question Answering. In <i>Proc. of ACL2023</i> .			937
884				938
885				939
886				940
887	Evan Hernandez, Belinda Z Li, and Jacob Andreas. 2023. Inspecting and editing knowledge representations in language models. <i>arXiv preprint arXiv:2304.00740</i> .			941
888				942
889				943
890				944
891	John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.			945
892				946
893				947
894				948
895				949
896				950
897				951
898	Aliyah R. Hsu, Georgia Zhou, Yeshwanth Cherapanamjeri, Yaxuan Huang, Anobel Odisho, Peter R. Carroll, and Bin Yu. 2025. Efficient automated circuit discovery in transformers using contextual decomposition . In <i>The Thirteenth International Conference on Learning Representations</i> .			952
899				953
900				954
901				955
902				956
903				957
904	Nan Hu, Yike Wu, Guilin Qi, Dehai Min, Jiaoyan Chen, Jeff Z Pan, and Zafar Ali. 2023. An Empirical Study of Pre-trained Language Models in Simple Knowledge Graph Question Answering. In <i>Journal of World Wide Web</i> , pages 1–32.			958
905				959
906				960
907				961
908				962
909	Yaojie Hu and Jin Tian. 2022. Neuron dependency graphs: A causal abstraction of neural networks . In <i>Proceedings of the 39th International Conference on Machine Learning</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 9020–9040. PMLR.			963
910				964
911				965
912				966
913				967
914				968
915	Baixiang Huang, Canyu Chen, Xiong Xiao Xu, Ali Payani, and Kai Shu. 2024. Can knowledge editing really correct hallucinations? <i>arXiv preprint arXiv:2410.16251</i> .			969
916				970
917				971
918				972
919	Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. Transformer-patcher: One mistake worth one neuron . In <i>The Eleventh International Conference on Learning Representations</i> .			973
920				974
921				975
				976
				977

978	Opitz, and Tobias Hecking. 2024. Style vectors for steering generative large language models . In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 782–802, St. Julian’s, Malta. Association for Computational Linguistics.	1034
979		1035
980		1036
981		1037
982		1038
983	Michael Lan, Philip Torr, and Fazl Barez. 2024. Towards interpretable sequence continuation: Analyzing shared circuits in large language models . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 12576–12601, Miami, Florida, USA. Association for Computational Linguistics.	1039
984		1040
985		1041
986		1042
987		
988		1043
989		1044
990	Anna Langedijk, Hosein Mohebbi, Gabriele Sarti, Willem Zuidema, and Jaap Jumelet. 2024. DecoderLens: Layerwise interpretation of encoder-decoder transformers . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 4764–4780, Mexico City, Mexico. Association for Computational Linguistics.	1045
991		1046
992		1047
993		
994		1048
995		1049
996		1050
997		1051
998		
999	Bruce W. Lee, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Erik Miehl, Pierre Dognin, Manish Nagireddy, and Amit Dhurandhar. 2025. Programming refusal with conditional activation steering . In <i>The Thirteenth International Conference on Learning Representations</i> .	1052
1000		1053
1001		1054
1002		1055
1003		1056
1004	Jiaang Li, Quan Wang, Zhongnan Wang, Yongdong Zhang, and Zhendong Mao. 2025a. Elder: Enhancing lifelong model editing with mixture-of-lora . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 24440–24448.	1057
1005		1058
1006		1059
1007		1060
1008		
1009	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Inference-time intervention: Eliciting truthful answers from a language model . <i>Advances in Neural Information Processing Systems</i> , 36:41451–41530.	1061
1010		1062
1011		1063
1012		1064
1013		
1014	Qi Li, Xiang Liu, Zhenheng Tang, Peijie Dong, Zeyu Li, Xinglin Pan, and Xiaowen Chu. 2024a. Should we really edit language models? on the evaluation of edited language models . <i>Advances in Neural Information Processing Systems</i> , 37:30850–30885.	1065
1015		1066
1016		1067
1017		1068
1018		1069
1019		1070
1020	Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. 2024b. Pmet: Precise model editing in a transformer . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 18564–18572.	1071
1021		1072
1022		1073
1023		
1024	Yanhong Li, Chunling Fan, Mingqing Huang, and Chengming Li. 2024c. Learning from mistakes: A comprehensive review of knowledge editing for large language models . In <i>2024 IEEE International Conference on Smart Internet of Things (SmartIoT)</i> , pages 563–569.	1074
1025		1075
1026		1076
1027		1077
1028		1078
1029	Zhe Li, Wei Zhao, Yige Li, and Jun Sun. 2025b. Do influence functions work on large language models? In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 14367–14382, Suzhou, China. Association for Computational Linguistics.	1079
1030		1080
1031		1081
1032		1082
1033		1083
	Yibing Liu, Haoliang Li, Yangyang Guo, Chenqi Kong, Jing Li, and Shiqi Wang. 2022. Rethinking attention-model explainability through faithfulness violation test . In <i>International conference on machine learning</i> , pages 13807–13824. PMLR.	1084
		1085
		1086
		1087
		1088
		1089
		1090
		1091
		1092
		1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200

1088	Lissandrini, ussa Biswas, Gerard de Melo, Angela Bonifati, Edlira Vakaj, Mauro Dragoni, and amien Graux. 2023. Large language models and knowledge graphs: Opportunities and challenges. <i>Transactions on Graph Data and Knowledge</i> .		
1089			
1090			
1091			
1092			
1093	Damian Pascual, Beni Egressy, Clara Meister, Ryan Cotterell, and Roger Wattenhofer. 2021. A plug-and-play method for controlled text generation . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3973–3997, Punta Cana, Dominican Republic. Association for Computational Linguistics.		
1094			
1095			
1096			
1097			
1098			
1099			
1100	Danish Pruthi, Mansi Gupta, Bhuwan Dhingra, Graham Neubig, and Zachary C. Lipton. 2020. Learning to deceive with attention-based explanations . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4782–4793, Online. Association for Computational Linguistics.		
1101			
1102			
1103			
1104			
1105			
1106			
1107	Siyuan Qi, Bangcheng Yang, Kailin Jiang, Xiaobo Wang, Jiaqi Li, Yifan Zhong, Yaodong Yang, and Zilong Zheng. 2025. In-context editing: Learning knowledge from self-induced distributions . In <i>The Thirteenth International Conference on Learning Representations</i> .		
1108			
1109			
1110			
1111			
1112			
1113	Shanbao Qiao, Xuebing Liu, and Seung-Hoon Na. 2024. COMEM: In-context retrieval-augmented mass-editing memory in large language models . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 2333–2347, Mexico City, Mexico. Association for Computational Linguistics.		
1114			
1115			
1116			
1117			
1118			
1119			
1120	Yifu QIU, Zheng Zhao, Yftah Ziser, Anna Korhonen, Edoardo Ponti, and Shay B Cohen. 2024. Spectral editing of activations for large language model alignment . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .		
1121			
1122			
1123			
1124			
1125	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?" explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pages 1135–1144.		
1126			
1127			
1128			
1129			
1130			
1131	Nina Rimskey, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. 2024. Steering llama 2 via contrastive activation addition . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15504–15522, Bangkok, Thailand. Association for Computational Linguistics.		
1132			
1133			
1134			
1135			
1136			
1137			
1138	Shuanghong Shen, Qi Liu, Zhenya Huang, Yonghe Zheng, Minghao Yin, Minjuan Wang, and Enhong Chen. 2024. A survey of knowledge tracing: Models, variants, and applications . <i>IEEE Transactions on Learning Technologies</i> , 17:1858–1879.		
1139			
1140			
1141			
1142			
1143	Wei Shi, Sihang Li, Tao Liang, Mingyang Wan, Guojun Ma, Xiang Wang, and Xiangnan He. 2025a. Route		
1144			
		sparse autoencoder to interpret large language models . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 6812–6826, Suzhou, China. Association for Computational Linguistics.	1145
			1146
			1147
			1148
			1149
		Xiang Shi, Jiawei Liu, Yinpeng Liu, Qikai Cheng, and Wei Lu. 2025b. Know where to go: Make llm a relevant, responsible, and trustworthy searchers. <i>Decision Support Systems</i> , 188:114354.	1150
			1151
			1152
			1153
		Chandan Singh, Jeevana Priya Inala, Michel Galley, Rich Caruana, and Jianfeng Gao. 2024. Rethinking interpretability in the era of large language models. <i>arXiv preprint arXiv:2402.01761</i> .	1154
			1155
			1156
			1157
		Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. 2017. Smoothgrad: removing noise by adding noise.	1158
			1159
			1160
		Alessandro Stolfo, Vidhisha Balachandran, Safoora Yousefi, Eric Horvitz, and Besmira Nushi. 2025. Improving instruction-following in language models through activation steering . In <i>The Thirteenth International Conference on Learning Representations</i> .	1161
			1162
			1163
			1164
			1165
		Nishant Subramani, Nivedita Suresh, and Matthew E Peters. 2022. Extracting latent steering vectors from pretrained language models. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 566–581.	1166
			1167
			1168
			1169
			1170
		Zengkui Sun, Yijin Liu, Jiaan Wang, Fandong Meng, Jinan Xu, Yufeng Chen, and Jie Zhou. 2024. Outdated issue aware decoding for factual knowledge editing . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 9282–9293, Bangkok, Thailand. Association for Computational Linguistics.	1171
			1172
			1173
			1174
			1175
			1176
			1177
		Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In <i>International conference on machine learning</i> , pages 3319–3328. PMLR.	1178
			1179
			1180
			1181
		Denis Sutter, Julian Minder, Thomas Hofmann, and Tiago Pimentel. 2025. The non-linear representation dilemma: Is causal abstraction enough for mechanistic interpretability? In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	1182
			1183
			1184
			1185
			1186
			1187
		Ryota Takatsuki, Sonia Joseph, Ipppei Fujisawa, and Ryota Kanai. 2025. Decoding vision transformers: the diffusion steering lens. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 4819–4824.	1188
			1189
			1190
			1191
			1192
		Chenmien Tan, Ge Zhang, and Jie Fu. 2024. Massive editing for large language models via meta learning . In <i>The Twelfth International Conference on Learning Representations</i> .	1193
			1194
			1195
			1196
		Juanhe (TJ) Tan. 2023. Causal abstraction for chain-of-thought reasoning in arithmetic word problems . In <i>Proceedings of the 6th BlackboxNLP Workshop</i> :	1197
			1198
			1199

1200	<i>Analyzing and Interpreting Neural Networks for NLP</i> , pages 155–168, Singapore. Association for Computational Linguistics.	Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. 2024b. Knowledge editing for large language models: A survey. <i>ACM Computing Surveys</i> , 57(3):1–37.	1255
1201			1256
1202			1257
1203	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. BERT rediscovers the classical NLP pipeline. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4593–4601, Florence, Italy. Association for Computational Linguistics.	Weixuan Wang, Barry Haddow, and Alexandra Birch. 2024c. Retrieval-augmented multilingual knowledge editing. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 335–354, Bangkok, Thailand. Association for Computational Linguistics.	1259
1204			1260
1205			1261
1206			1262
1207			1263
1208			1264
1209	Lukas Thede, Karsten Roth, Matthias Bethge, Zeynep Akata, and Thomas Hartvigsen. 2025. Wikibigedit: Understanding the limits of lifelong knowledge editing in LLMs. In <i>Forty-second International Conference on Machine Learning</i> .	Wenxuan Wang, Zizhan Ma, Zheng Wang, Chenghan Wu, Jiaming Ji, Wenting Chen, Xiang Li, and Yixuan Yuan. 2025d. A survey of llm-based agents in medicine: How far are we from baymax? <i>arXiv preprint arXiv:2502.11211</i> .	1265
1210			1266
1211			1267
1212			1268
1213			1269
1214	Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. 2024. Function vectors in large language models. In <i>The Twelfth International Conference on Learning Representations</i> .	Yifei Wang, Yuheng Chen, Wanting Wen, Yu Sheng, Linjing Li, and Daniel Dajun Zeng. 2024d. Unveiling factual recall behaviors of large language models through knowledge neurons. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 7388–7402, Miami, Florida, USA. Association for Computational Linguistics.	1270
1215			1271
1216			1272
1217			1273
1218			1274
1219	Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. 2020. Investigating gender bias in language models using causal mediation analysis. <i>Advances in neural information processing systems</i> , 33:12388–12401.	Yiwei Wang, Muhao Chen, Nanyun Peng, and Kai-Wei Chang. 2024e. Deepedit: Knowledge editing as decoding with constraints. <i>arXiv preprint arXiv:2401.10471</i> .	1275
1220			1276
1221			1277
1222			1278
1223			1279
1224			1280
1225	Ke Wang, Yiming QIN, Nikolaos Dimitriadis, Alessandro Favero, and Pascal Frossard. 2025a. MEMOIR: Lifelong model editing with minimal overwrite and informed retention for LLMs. In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah Goodman. 2023. Interpretability at scale: Identifying causal mechanisms in alpaca. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1282
1226			1283
1227			1284
1228			1285
1229			1286
1230			1287
1231	Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. 2024a. WISE: Rethinking the knowledge memory for lifelong model editing of large language models. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	Jianhao Yan, Futing Wang, Yun Luo, Yafu Li, and Yue Zhang. 2025. Keys to robust edits: from theoretical insights to practical advances. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 22545–22560.	1288
1232			1289
1233			1290
1234			1291
1235			1292
1236			1293
1237	Pinzheng Wang, Zecheng Tang, Keyan Zhou, Juntao Li, Qiaoming Zhu, and Min Zhang. 2025b. Revealing and mitigating over-attention in knowledge editing. In <i>The Thirteenth International Conference on Learning Representations</i> .	Yunzhi Yao, Jizhan Fang, Jia-Chen Gu, Ningyu Zhang, Shumin Deng, Huajun Chen, and Nanyun Peng. 2025. CaKE: Circuit-aware editing enables generalizable knowledge learners. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 11377–11393, Suzhou, China. Association for Computational Linguistics.	1294
1238			1295
1239			1296
1240			1297
1241			1298
1242	Renzhi Wang and Piji Li. 2024. LEMoE: Advanced mixture of experts adaptor for lifelong model editing of large language models. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 2551–2575, Miami, Florida, USA. Association for Computational Linguistics.	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing large language models: Problems, methods, and opportunities. <i>Preprint</i> , arXiv:2305.13172.	1299
1243			1300
1244			1301
1245			1302
1246			1303
1247			1304
1248	Shiqi Wang, Qi Wang, Runliang Niu, He Kong, and Yi Chang. 2025c. MicroEdit: Neuron-level knowledge disentanglement and localization in lifelong model editing. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 33870–33884, Suzhou, China. Association for Computational Linguistics.	Yunzhi Yao, Ningyu Zhang, Zekun Xi, Mengru Wang, Ziwen Xu, Shumin Deng, and Huajun Chen. 2024. Knowledge circuits in pretrained transformers. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	1305
1249			1306
1250			1307
1251			1308
1252			1309
1253			
1254			

1310	Lang Yu, Qin Chen, Jie Zhou, and Liang He. 2024.	Hanlun Zhu, Yunshi Lan, Xiang Li, and Weining Qian.	1366
1311	Melo: enhancing model editing with neuron-indexed	2025. Initializing and retrofitting key-value adaptors	1367
1312	dynamic lora . In <i>Proceedings of the Thirty-Eighth</i>	for traceable model editing . In <i>Findings of the As-</i>	1368
1313	<i>AAAI Conference on Artificial Intelligence and Thirty-</i>	<i>sociation for Computational Linguistics: ACL 2025</i> ,	1369
1314	<i>Sixth Conference on Innovative Applications of</i>	pages 2958–2971, Vienna, Austria. Association for	1370
1315	<i>Artificial Intelligence and Fourteenth Symposium</i>	Computational Linguistics.	1371
1316	<i>on Educational Advances in Artificial Intelligence</i> ,		
1317	AAAI’24/IAAI’24/EAAI’24. AAAI Press.		
1318	Fred Zhang and Neel Nanda. 2024. Towards best prac-		
1319	tices of activation patching in language models: Met-		
1320	rics and methods . In <i>The Twelfth International Con-</i>		
1321	<i>ference on Learning Representations</i> .		
1322	Mengqi Zhang, Xiaotian Ye, Qiang Liu, Pengjie Ren,		
1323	Shu Wu, and Zhumin Chen. 2024a. Uncovering		
1324	overfitting in large language model editing . <i>arXiv</i>		
1325	<i>preprint arXiv:2410.07819</i> .		
1326	Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng		
1327	Wang, Shumin Deng, Mengru Wang, Zekun Xi,		
1328	Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan		
1329	Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong		
1330	Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang		
1331	Zhang, and 3 others. 2024b. A comprehensive study		
1332	of knowledge editing for large language models .		
1333	<i>Preprint</i> , arXiv:2401.01286.		
1334	Zheyuan Zhang, Daniel Zhang-Li, Jifan Yu, Linlu Gong,		
1335	Jinchang Zhou, Zhanxin Hao, Jianxiao Jiang, Jie		
1336	Cao, Huiqin Liu, Zhiyuan Liu, and 1 others. 2025.		
1337	Simulating classroom education with llm-empowered		
1338	agents . In <i>Proceedings of the 2025 Conference of the</i>		
1339	<i>Nations of the Americas Chapter of the Association</i>		
1340	<i>for Computational Linguistics: Human Language</i>		
1341	<i>Technologies (Volume 1: Long Papers)</i> , pages 10364–		
1342	10379.		
1343	Runcong Zhao, Chengyu Cao, Qinglin Zhu, Xiucheng		
1344	Ly, Shun Shao, Lin Gui, Ruifeng Xu, and Yulan He.		
1345	2025. Sparse activation editing for reliable instruc-		
1346	tion following in narratives . In <i>Proceedings of the</i>		
1347	<i>2025 Conference on Empirical Methods in Natural</i>		
1348	<i>Language Processing</i> , pages 25828–25843, Suzhou,		
1349	China. Association for Computational Linguistics.		
1350	Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong		
1351	Wu, Jingjing Xu, and Baobao Chang. 2023. Can		
1352	we edit factual knowledge by in-context learning?		
1353	In <i>Proceedings of the 2023 Conference on Empiri-</i>		
1354	<i>cal Methods in Natural Language Processing</i> , pages		
1355	4862–4876, Singapore. Association for Computa-		
1356	tional Linguistics.		
1357	Zexuan Zhong, Zhengxuan Wu, Christopher D. Man-		
1358	ning, Christopher Potts, and Danqi Chen. 2024.		
1359	Mquake: Assessing knowledge editing in lan-		
1360	guage models via multi-hop questions . <i>Preprint</i> ,		
1361	arXiv:2305.14795.		
1362	Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh		
1363	Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv Kumar.		
1364	2020. Modifying memories in transformer models .		
1365	<i>arXiv preprint arXiv:2012.00363</i> .		

A Paper collection

We note that Knowledge Tracing traditionally refers to a student modeling task in online education, aiming to infer learners’ latent knowledge states from behavioral sequences. (Shen et al., 2024) In contrast, recent LLM studies employ the term to describe a fundamentally different problem: tracing the internal representations and causal pathways through which factual knowledge influences model predictions.

In this work, we use Knowledge Tracing in a different sense, referring to methods that trace where and how factual knowledge functionally participates in the internal computation of large language models.

To get the related paper, we use the above steps:

- We use the Tools ² to collect papers from past 3 years (2022 to 2025) from Top NLP and ML conferences (i.e., ACL, EMNLP, NeurIPS, ICML, ICLR). Got 50268 paper in total.
- we use nougat ³ to extract the paper’s abstract and title information from the pdf file. And use all-MiniLM-L6-v2 ⁴ to filter the relevant paper with abstract and title embedding. We got 24672 paper.
- To reduce the bias from different research area, we use GPT4o to score the relevance of a paper (1-10) with the keywords such as ‘Interpretable’, ‘Controlling output by modifying model internals’, ‘Model/knowledge Editing’ and ‘Causal tracing and intervention in deep models’. And got the paper that relevance score > 6. We got 1980 paper.
- We then collect their publication date and citation numbers from Semantic Scholar API⁵ and then calculate the RCI for each paper, and choose the paper that RCI score > 1. We got 355 paper.
- And then use A scholarly research assistant from Ai2 ⁶ to collect the most related paper to knowledge Tracing and Editing (ignore the year with this tools.). Then we compare the

²<https://github.com/hegongshan/paper-downloader>

³<https://facebookresearch.github.io/nougat/>

⁴<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

⁵<https://www.semanticscholar.org/product/api>

⁶<https://asta.allen.ai/>

results with the 355 paper to filter the noise and add the ignore paper. We got 220 paper with high related topic, contain 130 Tracing paper and 90 Editing paper. And recalculate the RCI score based on these papers.

B Tools for Editing and Tracing

Model editing has benefited from a rapid standardization of its workflow, supported by comprehensive toolkits like *EasyEdit* and *FastEdit*. This robust infrastructure significantly lowers the entry barrier for developing and scaling efficient intervention methods. In contrast, the tracing ecosystem—while featuring influential tools like *TransformerLens* and *Circuit-tracer*, remains heavily constrained by its dependency on low-level architectural interventions (e.g., neuron-level hooking). As shown in Table 4, tracing-related tools are numerous and diverse, covering visualization (e.g., BertViz, Transformer-Explainer), attribution (e.g., Captum, SHAP), circuit analysis (e.g., circuit-tracer), and interactive analysis platforms (e.g., LIT).

These methods often lack cross-model portability, making it difficult to generalize tracing conclusions to edited models or heterogeneous architectures.

C Paper Scores

In this work, we use Relative Citation Index (RCI) metric (Du et al., 2025), which estimate the expected citations with respect to publication age to prevent bias between original citations from different publication dates. The age A_i of a paper p_i is computed as:

$$A = T - Year_i, \quad (1)$$

where T is the date when the citation is collected (20th Nov. 2025) and $Year_i$ is the year where paper i is first published. Thus, we can model the relation between citation number C_i and age A_i of paper p_i in three different way, which are:

linear model:

$$C_i = \beta + \alpha A_i \quad (2)$$

exponential model:

$$C_i = \exp(\beta + \alpha A_i) \quad (3)$$

log-log regression model:

$$\log(C_i + 1) = \beta + \alpha \log A_i + \epsilon_i \quad (4)$$

Tools	Stars	Link
transformerlen	2.9k	https://github.com/TransformerLensOrg/TransformerLen
BertViz	7.9k	https://github.com/jessevig/bertviz
LIT	3.6k	https://github.com/PAIR-code/lit
ECCO	2.1k	https://github.com/jalammar/ecco
Captum	5.5k	https://github.com/pytorch/captum
Transformers-Interpret	1.4k	https://github.com/cdpierse/transformers-interpret
AllenNLP Interpret	11.9k	https://github.com/allenai/allennlp-interpret
transformer-explainer	6.3k	https://github.com/poloclub/transformer-explainer
SHAP	24.9k	https://github.com/slundberg/shap
transformer-debugger	4.1k	https://github.com/openai/transformer-debugger
llm-transparency-tool Public	1.2k	https://github.com/facebookresearch/llm-transparency-tool
circuit-tracer	2.5k	https://github.com/safety-research/circuit-tracer/tree/main
bi-att-flow	1.5k	https://github.com/allenai/bi-att-flow
EasyEdit	2.7k	https://github.com/zjunlp/EasyEdit/tree/main
FastEdit	1.4k	https://github.com/hiyouga/FastEdit

Table 4: Popular Tools for knowledge Editing and Tracing.

We pick log-log regression model to compute the expected citation for next step⁷, and we are able to obtain the expected citation number \hat{C}_i of paper p_i with age A_i as:

$$\hat{C}_i = \exp(\hat{\beta})A_i^{\hat{\alpha}} \quad (5)$$

Then we compute the relative citation index RCI_i of paper p_i as:

$$RCI_i = \frac{C_i}{\hat{C}_i} \quad (6)$$

When $RCI_i \geq 1$, we consider this paper over-cited than its expectations, and vice versa. In this paper, we focus on the paper with $RCI \geq 1$, for which we believe has more influence.

The Citation Velocity is calculate by:

$$V = \frac{C_i}{A_i}, \quad (7)$$

captures its current research momentum.

The Important score is calculate by:

$$Importance_i = 2 \times \frac{RCI_i \times \overline{Velocity}_i}{RCI_i + \overline{Velocity}_i} \quad (8)$$

Where the Normalized Citation Velocity is :

$$\overline{Velocity}_i = \frac{Velocity_i}{\frac{1}{N} \sum_{j=1}^N Velocity_j} \quad (9)$$

As formulated in Eq. 8, this score is the harmonic mean of the *Relative Citation Index* (RCI) and the *Normalized Citation Velocity*.

⁷The estimation is: $\hat{\beta} = 1.815$, $\hat{\alpha} = 0.770$

D Relative Citation Index and Citation Velocity

D.1 Quadrant analysis

The Top Paper in Figure 3b for each type is:

P198: "Why Should I Trust You?": Explaining the Predictions of Any Classifier (Ribeiro et al., 2016)

P095: Investigating Gender Bias in Language Models Using Causal Mediation Analysis (Vig et al., 2020)

P016: Axiomatic Attribution for Deep Networks (Sundararajan et al., 2017)

P213: AnyEdit: Edit Any Knowledge Encoded in Language Models (Jiang et al., 2025)

P113: Locating and Editing Factual Associations in GPT (Meng et al., 2022)

P114: Mass-Editing Memory in a Transformer (Meng et al., 2023)

For the High-impact & High-momentum quadrant (Q1): Tracing papers overwhelmingly dominate the "Mainstream Core" with 59 works compared to 29 in Editing. This nearly 2:1 ratio indicates that interpretability research has not only established a solid foundation but continues to drive the field's primary intellectual momentum. For the most important paper in Q1 also shows the importance for tracing paper is higher than the editing paper.

For the Emerging & High-momentum quadrant (Q2): The bottleneck for Editing is most evident in

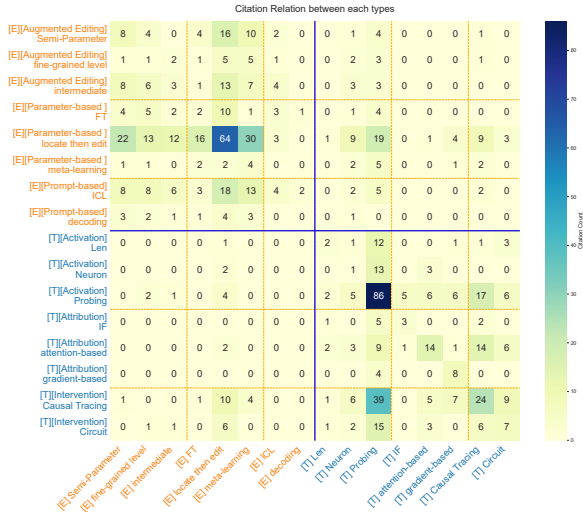


Figure 4: Cross-category citation analysis

1511 this quadrant. While 16 Tracing papers are categor-
 1512 ized as high-momentum emerging works, only 3
 1513 Editing papers achieve this status. This suggests
 1514 that despite the surge in publication volume, Edit-
 1515 ing struggles to produce "breakout" research that
 1516 captures immediate and widespread academic at-
 1517 tention.

1518 For the Low-impact & Low-momentum quad-
 1519 rant (Q3): Both fields have a significant presence
 1520 here (38 Editing vs. 49 Tracing), representing the
 1521 "long tail" of incremental research. However, for
 1522 Editing, this quadrant accounts for the largest share
 1523 of its total publications, reflecting a high volume of
 1524 work that has yet to achieve significant resonance.

1525 For the Established & High Influence quadrant
 1526 (Q4): Notably, Editing has more "Legacy" papers
 1527 in Q4 (12) compared to Tracing (7). This indicates
 1528 that while early Editing methods achieved high
 1529 citation counts, the field has struggled to maintain
 1530 that same level of influence in more recent, high-
 1531 velocity cycles compared to the sustained impact
 1532 of Tracing.

1533 D.2 Cross-category citation

1534 Cross-category citation (Figure 4) reveals a
 1535 strongly asymmetric dependency between knowl-
 1536 edge editing methods and interpretability-oriented
 1537 tracing techniques. This hierarchy positions in-
 1538 terpretability primarily as a supporting instrument
 1539 for localization, rather than a co-evolving research
 1540 objective.

1541 **Editing approaches exhibit substantial citation**
 1542 **flows toward foundational tracing methods, yet**
 1543 **the reverse is rarely true.** Specifically, the

1544 Locate-then-edit paradigm acts as the primary con-
 1545 sumer of interpretability research, citing Probing
 1546 (19 citations), Causal Tracing (9 citations), and
 1547 Circuit analysis (3 citations). Similarly, Semi-
 1548 parameter and Intermediate editing methods con-
 1549 sistently reference Probing (4 and 3 citations re-
 1550 spectively) to validate their intervention targets. In
 1551 stark contrast, tracing show minimal reverse de-
 1552 pendence, for instance, Causal Tracing and Circuit-
 1553 based works cite the entire spectrum of editing
 1554 methods fewer than 10 times in total.

1555 **Locate-then-edit functions as a structural nexus,**
 1556 **maintaining dense cross-citations that bridge**
 1557 **heterogeneous editing strategies.** It exhibits
 1558 strong connectivity with Semi-parameter (22 cita-
 1559 tions), Fine-grained (13 citations), Intermediate
 1560 (12 citations), FT (16 citations), and ICL-based
 1561 editing (18 citations). This central positioning con-
 1562 firms that "localization" remains the dominant men-
 1563 tal model for knowledge manipulation, even if the
 1564 localized components do not always provide the
 1565 necessary controllability for stable editing. *Prompt-*
 1566 *based and Augmented Editing methodologies ex-*
 1567 *hibit a superficial and unidirectional reliance on*
 1568 *interpretability foundations.* Compared to direct
 1569 parameter-based interventions, these approaches
 1570 maintain a marginal dependency on tracing re-
 1571 search.

1572 Table 5 summarizes the alignment between
 1573 knowledge editing methods and tracing types.
 1574 Parameter-Based Editing methods, such as ROME,
 1575 MEMIT, and AlphaEdit, consistently leverage
 1576 intervention- or attribution-based tracing to guide
 1577 fine-grained parameter updates, reflecting a strong
 1578 dependency on tracing signals. In contrast, Aug-
 1579 mented and Prompt-Based Editing methods show
 1580 limited or implicit reliance on tracing, often oper-
 1581 ating at the module, hidden state, or token level
 1582 without explicit mechanistic guidance. This dis-
 1583 parity suggests that while tracing could provide
 1584 actionable guidance for these latter categories, cur-
 1585 rent research has largely overlooked this potential,
 1586 highlighting an opportunity for future work to more
 1587 systematically integrate tracing into diverse editing
 1588 paradigms.

1589 E causal diagnostic (AlphaEdit + causal

1590 tracing)

1591 Full Causal trace visualization results are provided
 1592 in Figure 5, 6, 7, 8, and 9.

Editing Method	Tracing Alignment	Interventional Module	Edit Granularity
Parameter-Based Editing			
ROME (Meng et al., 2022)	Intervention-based	MLP block	Parameters level
MEMIT (Meng et al., 2023)	Intervention-based	MLP block	Parameters level
AlphaEdit (Fang et al., 2025)	Intervention-based	MLP block	Parameters level
MEND (Mitchell et al., 2022a)	Attribution-based	MLP block	Parameters level
KN (Dai et al., 2022)	Activation-based	Neurons	Neurons level
CAKE (Yao et al., 2025)	Intervention-based	MLP block	Parameters level
ICE (Qi et al., 2025)	-	Hidden states	Parameters level
MALMEN (Tan et al., 2024)	<i>Attribution-based</i>	Parameters	Parameters level
LTI (Zhang et al., 2024a)	-	MLP block	Parameters level
Augmented Editing			
SERAC (Mitchell et al., 2022b)	-	Module	Token level
WISE (Wang et al., 2024a)	-	MLP block	Hidden states level
T-Patch (Huang et al., 2023)	-	Hidden states	Neurons level
GRACE (Hartvigsen et al., 2023)	-	Hidden states	Hidden states level
MELO (Yu et al., 2024)	-	Module	Hidden states level
REMEDI (Hernandez et al., 2023)	<i>Attribution-based</i>	Hidden states	Hidden states level
Prompt-Based Editing			
IKE (Zheng et al., 2023)	-	Token	Token level
DeCK (Bi et al., 2024a)	-	Token	Token level
ATBIAS (Bi et al., 2024b)	<i>Activation-based</i>	Logits length	Token level
DISCO (Sun et al., 2024)	<i>Activation-based</i>	Logits length	Token level
MeLLO (Zhong et al., 2024)	-	Token	Token level

Table 5: Analysis of misalignment between knowledge editing and tracing. We categorize results into **Parameter-Based Editing**, **Augmented Editing**, and **Prompt-Based Editing**. Normal text represents explicit alignment, *italic text* represents implicit alignment, and ‘-’ denotes total misalignment.

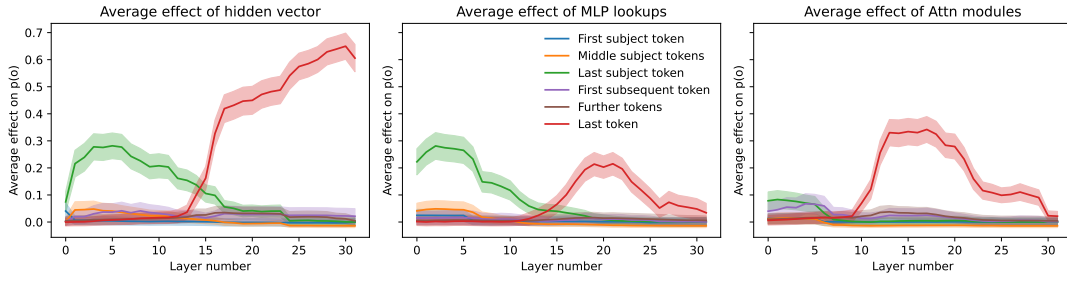


Figure 5: Before Editing, Causal trace on 100 Matched-Data (MD)

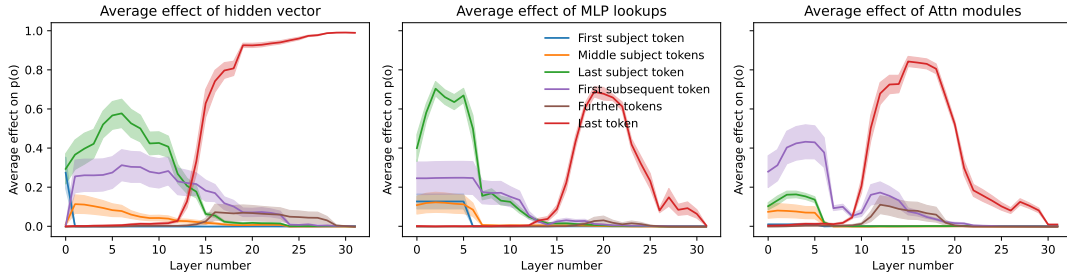


Figure 6: Post-Edit model (Edit on 100 Matched-Data (MD))

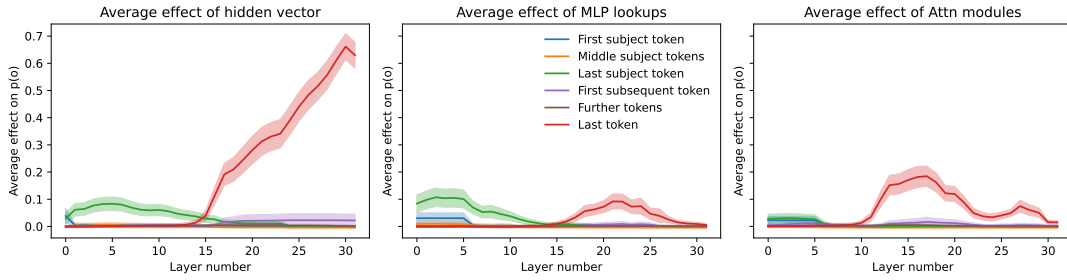


Figure 7: Post-Edit model (Edit on 100 Independent-Data (ID))

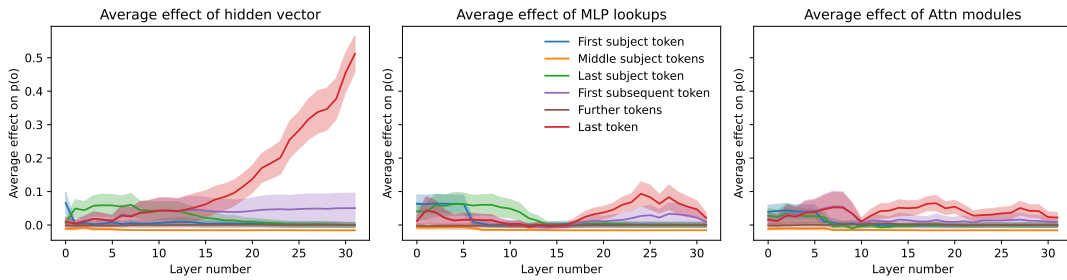


Figure 8: Post-Edit model (Edit on 300 Independent-Data (ID))

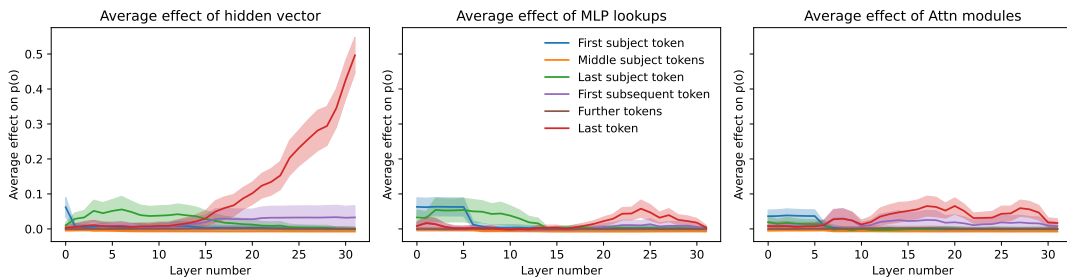


Figure 9: Post-Edit model (Edit on 500 Independent-Data (ID))

1593

F Architectural landscape of KT and KE

1594

The architectural landscape of the knowledge tracing and editing is provided in Figure 10.

1595

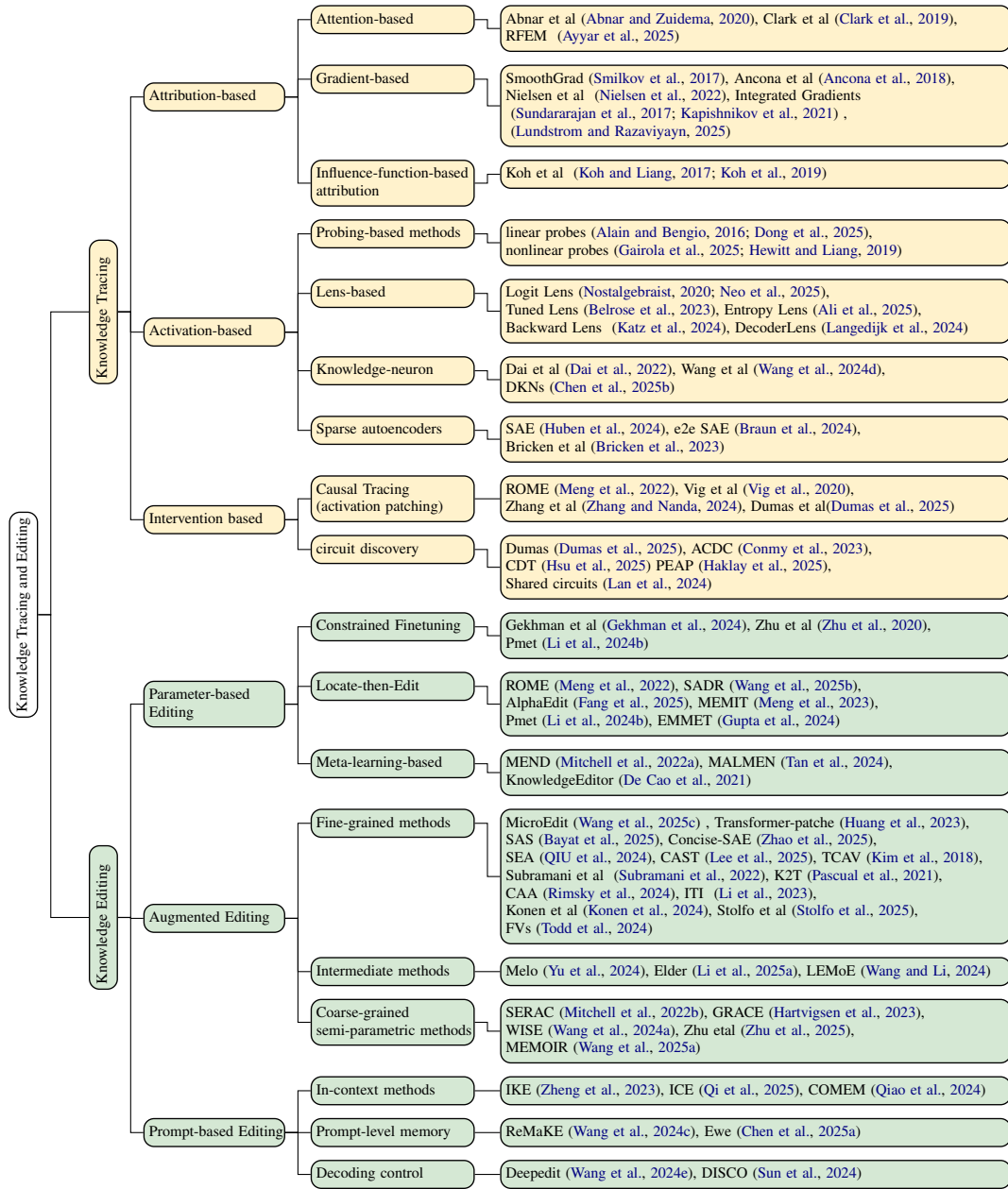


Figure 10: The overview of Knowledge tracing and editing methods