

Factual Retrieval in LLMs Is a Redundant, Distributed and Non-Contiguous Process

Anonymous ACL submission

Abstract

Large language models (LLMs) store and recall factual knowledge, yet the precise mechanism of how entity representations are transformed to enable specific attribute retrieval remains underexplored. In this work, we investigate this mechanism through the lens of an “attribute-computation path”—a sequence of computational steps over the entity representation required to elicit a target attribute. We then propose an iterative patching protocol to identify a minimal subset of layers necessary for this computation. Applying our method to LLaMA 3.1 8B and Qwen 3 8B, we find that these paths are non-contiguous, often skipping layers, and that models possess multiple, functionally-equivalent paths for the same entity and fact, highlighting a high degree of redundancy in attribute computation. This implies that knowledge computation is highly distributed, potentially explaining the localization-editing mismatch and suggesting that knowledge storage and retrieval in LLMs is far from being well understood.

1 Introduction

Large language models (LLMs) have been shown to store and recall factual knowledge expressed as entities and their relations (Petroni et al., 2019; Jiang et al., 2020; Cohen et al., 2023). For example, when prompted with “The mother tongue of Angela Merkel is”, an LLM will predict “German”, the correct answer. While prior studies attempted to identify components of the model where factual knowledge is stored (Geva et al., 2021; Dai et al., 2021; Geva et al., 2022; Meng et al., 2022; Gurnee and Tegmark, 2023; Katz et al., 2024; Yu and Ananiadou, 2024) and to map the general information flow during recall (Geva et al., 2023; Nanda et al., 2023; Chughtai et al., 2024; Wang and Xu, 2025; Yao et al., 2024; Yu et al., 2025),

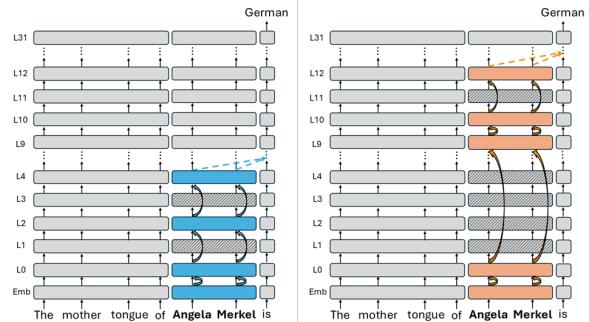


Figure 1: Two functionally equivalent minimal attribute-computation paths for the entity “Angela Merkel”. Both paths are non-contiguous, skipping intermediate layers. The blue path represents a compact, early computation, while the orange path illustrates a deeper, more distributed alternative.

the precise mechanics of how facts are retrieved from model parameters remain unclear. Existing work posits that this process occurs at the last entity token position; however, the term “retrieval” itself is not well defined. It remains ambiguous whether retrieval is a continuous process of information accumulation across a range of layers, or if it is defined by a discrete state—a specific point where the representation becomes rich enough to elicit the specific correct attribute. Furthermore, the transformations the entity representation undergoes *before* it becomes capable of eliciting the target attribute are largely unmapped. To address these complexities, we propose investigating the factual recall mechanism through a new lens: mapping the *minimal computation path* the entity representation must undergo in order to achieve attribute recall, in a given factual-recall prompt.

Given a prompt with an entity with a relation (e.g., “The mother tongue of Angela Merkel is”), and a target attribute (e.g., “German”), we define an *attribute-computation path* as a

067 sequence of computations over the entity rep- 116
068 resentation, that together enable the LLM to 117
069 answer the prompt correctly. Each computa- 118
070 tion is implemented as a transformer layer, and 119
071 we seek a subset of layers which is both suffi-
072 cient and necessary for computing the desired
073 output. We call such a subset of layers, in
074 which no layer can be removed without hurting
075 the computation, a *minimal attribute compu-*
076 *tation path*, and observe that each layer along
077 such a minimal path plays an active role in the
078 factual recall process for the attribute.

079 We present a method to identify minimal
080 computation paths for a given factual recall
081 prompt, and use it to extract computation
082 paths for a diverse set of such prompts, in
083 two open-weight LLMs (LLaMA 3.1 8B and
084 Qwen 3 8B). Our experiments reveal that for a
085 majority of target attributes:

- 086 • *The minimal path naturally ends in the* 136
087 *early-mid layers*, after which the attribute 137
088 value is fully available and no further pro- 138
089 cessing is required. This is consistent with 139
090 observations in previous works (Geva et al., 140
091 2023; Nanda et al., 2023; Meng et al., 141
092 2022).
- 093 • *Minimal paths consist of more than a single* 142
094 *layer*: there isn't a single layer responsible 143
095 for knowledge retrieval for a given entity- 144
096 relation-attribute triplet, but rather the 145
097 entity must be processed by multiple layers 146
098 before the information is fully available. 147
- 099 • *Minimal paths are often sparse and non-* 148
100 *contiguous*: only a subset of layers partici- 149
101 pate in each factual recall process, while 150
102 others may be skipped (but the same layer 151
103 may participate in multiple factual recall 152
104 processes). 153

105 Taken together, the last two items indicate that
106 factual knowledge is stored in a distributed and
107 non-localized manner: there isn't a single layer
108 or range of layers we can say is responsible for
109 a given fact, rather the factual information is
110 distributed across the parameters of multiple
111 layers. Our experiments further reveal "backup"
112 paths (Wang et al., 2022; McGrath et al., 2023):

- 113 • *Minimal paths are not unique*: for most 163
114 facts we explored, there is at least one al- 164
115 ternative set of layers that is sufficient for 165

retrieving the factual information. These
alternate paths are deeper and longer than
the primary paths the model uses in natu-
ral, non-intervened runs (Figure 1).

We find that the same factual knowledge can
be retrieved through multiple distinct compu-
tational paths within the model. This suggests
an even more complex and elaborate form of
knowledge representation, where the storage
and retrieval mechanism is not only distributed
but also highly redundant. These findings pro-
vide a new perspective on the mechanism of
factual recall (Geva et al., 2023; Meng et al.,
2022; Yu et al., 2025). We detail these, to-
gether with additional findings, in Section 6.
To locate the minimal paths, we use a novel it-
erative activation-patching protocol, described
in Section 5.

In summary, our work makes three main
contributions.

- We offer a new perspective on the sub-
process of entity enrichment (Geva et al.,
2023; Yu et al., 2025), showing that the
computation of the entity representation
for specific attribute recall unfolds as a
multi-step process involving several neces-
sary intermediate transformations, while
at the same time not requiring most model
layers.
- We reveal that *multiple minimal and func-*
tionally equivalent attribute computation
paths exist within the model with different
lengths and in different regions, highlight-
ing the distributed and redundant nature
of factual knowledge representation.
- We propose a novel activation patching
protocol to identify a minimal subset of
layers required for computation of an at-
tribute over the entity activations, provid-
ing a precise view of how factual knowledge
enriches the entity representations.

Based on our findings, we propose an expla-
nation for the discrepancy between where the
facts are located and where they are most ef-
fectively edited (Hase et al., 2023). We suggest
that because attribute computation requires
multiple transformations of the entity repre-
sentation, editing succeeds by intervening at a
necessary transformation stage, not necessarily
where the knowledge is localized.

2 Related Work

Entity Representation & Factual Recall.

Prior work frame factual recall in large language model as a three-stage process (Geva et al., 2023; Yu et al., 2025; Chughtai et al., 2024; Nanda et al., 2023). First, early attention layers *consolidate* the representation of the entities initial tokens into its final token (Nanda et al., 2023). Then, the final entity token undergoes an *entity enrichment phase* in which various entity attributes are loaded from the model parameters and encoded on the (final) entity tokens. This stage is primarily driven by early feed-forward (MLP) sublayers (Meng et al., 2022; Geva et al., 2023; Yu et al., 2025) in which the knowledge resides. Finally, once entity enrichment concludes and the knowledge is retrieved and loaded into the entity representation, information from the relation tokens propagates forward in the sequence, reaching the final token of the prompt. Then, the representation at the final token of the prompt *queries* the enriched entity representation to extract the target attribute (Geva et al., 2023).

It is also established that the resulting representation at the entity’s final token plays a central role in knowledge storage and retrieval (Geva et al., 2023; Yu et al., 2025; Hernandez et al., 2023; Ghandeharioun et al., 2024). Recent work demonstrates that, for many relations, attributes can be decoded from the entity’s final token using a simple, approximately linear transformation (Hernandez et al., 2023). Furthermore, probing the hidden states of this token in late–mid layers can reveal the extent of the model’s stored knowledge about the entity, without relying on generated outputs (Gottesman and Geva, 2024). When the relation follows the entity in the prompt, enrichment and attribute extraction can occur not only at the entity position but also at the relation and final positions (Yu et al., 2025; Chughtai et al., 2024). In these cases, deeper attention and feed-forward components store factual knowledge linking the entity and relation to the attribute.

While Meng et al., 2022 identified stages where representations become sufficient for attribute restoration and hypothesized accumulation across MLPs, the precise causal computations of this process remain under-explored.

We zoom in on the entity enrichment stage and trace the complete causal paths by which individual facts are loaded into the entity representation. Our findings reveal a complex and elaborate system in which each individual fact requires processing by multiple LLM layers before it can be accessed: there is no single layer in which a specific fact is retrieved.

Circuit Redundancy. Early work in mechanistic interpretability of LLMs and other models focused on identifying circuits: subgraphs of the model’s computational graph responsible for specific behaviors (Olah et al., 2020; Elhage et al., 2021; Wang et al., 2022; Goldowsky-Dill et al., 2023; Ameisen et al., 2025). However, the definition of a “circuit” is complicated by the distributed and redundant nature of LLMs.

McGrath et al. (2023) revealed the “hydra effects”, where the model self-repairs by activating redundant components when primary ones are ablated. Wang et al. (2022) similarly identified “backup heads” that perform the same function as primary heads but are only active when the primary path is disrupted. This complexity extends to model editing; Hase et al. (2023) demonstrated that while facts can be successfully edited at specific model locations, the most effective editing targets often diverge from the locations identified by standard localization methods. These findings suggest that models do not rely on a single, unique circuit for a given task, but rather possess multiple functionally equivalent mechanisms.

While Wang et al. (2022) revealed backup “Name-Mover Heads” in the IOI task, and the hydra effect (McGrath et al., 2023) has been investigated primarily regarding the information flow to the final token, we demonstrate redundancy within a different factual recall mechanism: the layers performing the computation process over the entity tokens that is used to retrieve the attribute.

3 Preliminaries and Notations

Layers and Activations. Let $I^{(i,\ell)}$ denote the input activations to layer ℓ at position i .¹ A layer function L_ℓ applied to $I^{(i,\ell)}$ results in the layer’s output $O^{(i,\ell)}$, which serves as the

¹With slight abuse of notation, we use a single index to refer to either a single token position or a range of tokens (corresponding to an entity).

input to the next layer, $I^{(i,\ell+1)}$:

$$L_\ell(I^{(i,\ell)}) = O^{(i,\ell)} = I^{(i,\ell+1)}$$

We refer to the same activations as either $O^{(i,\ell)}$ or $I^{(i,\ell+1)}$, depending on the context. The layer function L_ℓ computes attention (*Attn*) and MLP with residual connections:

$$\begin{aligned} L_\ell(I^{(i,\ell)}) &= I^{(i,\ell)} + a_{(\leq i,\ell)} + m_{(i,\ell)} \\ a_{(\leq i,\ell)} &= \text{Attn}_\ell(I^{(1,\ell)}, \dots, I^{(i,\ell)}) \\ m_{(i,\ell)} &= \text{MLP}_\ell(I^{(i,\ell)} + a_{(\leq i,\ell)}) \end{aligned}$$

Runs and Activation Patching. We denote the prediction of a model M given prompt P as $M(P)$. In this work we restrict ourselves to single word predictions using greedy decoding.

Given an n tokens input prompt, the transformer decoder step with $|L|$ layers computes the values $O^{(i,\ell)}$ for $1 \leq \ell \leq |L|; 1 \leq i \leq n$, in topological order where a node $O^{(i,\ell)}$ is computed after nodes with positions $\leq i$ and layers $< \ell$. A *patch operation* Φ specifies a list of locations (i, ℓ) and a corresponding patch vector for each. A *patched run* $M(P | \Phi)$ on an input prompt P computes the activation vectors according to their topological order, but, whenever (i, ℓ) corresponds to a patching location, $I^{(i,\ell)}$ is replaced with the corresponding value instead of its computed value, and the computation continues.

4 Setup and Goal

To understand how factual knowledge is stored and retrieved in LLMs, we are interested in tracing the computations performed by a transformer-based LLM when completing *factual recall prompts* such as “The mother tongue of Angela Merkel is”. We focus our attention on the process in which the *entity representation* evolves through the layers until it is sufficient for the model to produce the correct answer (“German”). Considering each layer as performing a computation over the entity representation, we look for what we call a *minimal attribute computation path*: a subset of layers that are *required* for producing the correct attributes.

Factual Recall Prompts. We study prompts of the form $(e, r) \rightarrow a$, where e is a subject entity, r is a relation, and a is the

target attribute to be predicted (not part of the prompt). For example, in the above prompt the subject entity e is “Angela Merkel”, the relation r is “mother’s tongue”, and the target attribute a is “German”. The relation may appear either before the entity or after it. Both entities and relations are drawn from a variety of domains.

Entity Representation is the activation values at the entity position for a given prompt. It evolves through the layers, with the representation at layer ℓ denoted as $O^{(e,\ell)}$.

Computation Paths. A *computation path* is a series of computations (a subset of layers) applied to an entity representation. An *attribute computation path* is a computation path after which the model correctly predicts the target attribute. A *minimal attribute-sufficient computation path* (or *minimal attribute path* for short) is a path in which every included layer is essential; removing any layer from this set disrupts the computation of the attribute over the entity activations.²

Relations to Knowledge Retrieval. Each layer belonging to a minimal attribute path plays an essential role in the process of *knowledge retrieval* for that attribute: it either extracts (parts of) the attribute value from the model’s parameters; transforms the entity representation to facilitate the knowledge extraction; or transforms extracted values into a form from which the attribute can be decoded. Thus, we consider the study of minimal attribute paths as an essential milestone for understanding factual knowledge retrieval in transformer-based LLMs.

5 Method

In an $|L|$ -layers transformer M for which $M(P_{(e,r)}) = a$, the path $[L_1, L_1, \dots, L_{|L|}]$ is attribute sufficient for a . But are all the layers necessary?

To identify a minimal attribute path, we start by determining the termination point of the path: the earliest layer $\ell_{attr} \leq |L|$ after which no more computations on the entity representation are needed for M to produce the attribute

²Note that “minimal” refers to the irreducibility of the path rather than its length: multiple distinct minimal paths of different lengths may exist. Indeed, we show that this is often the case in practice.

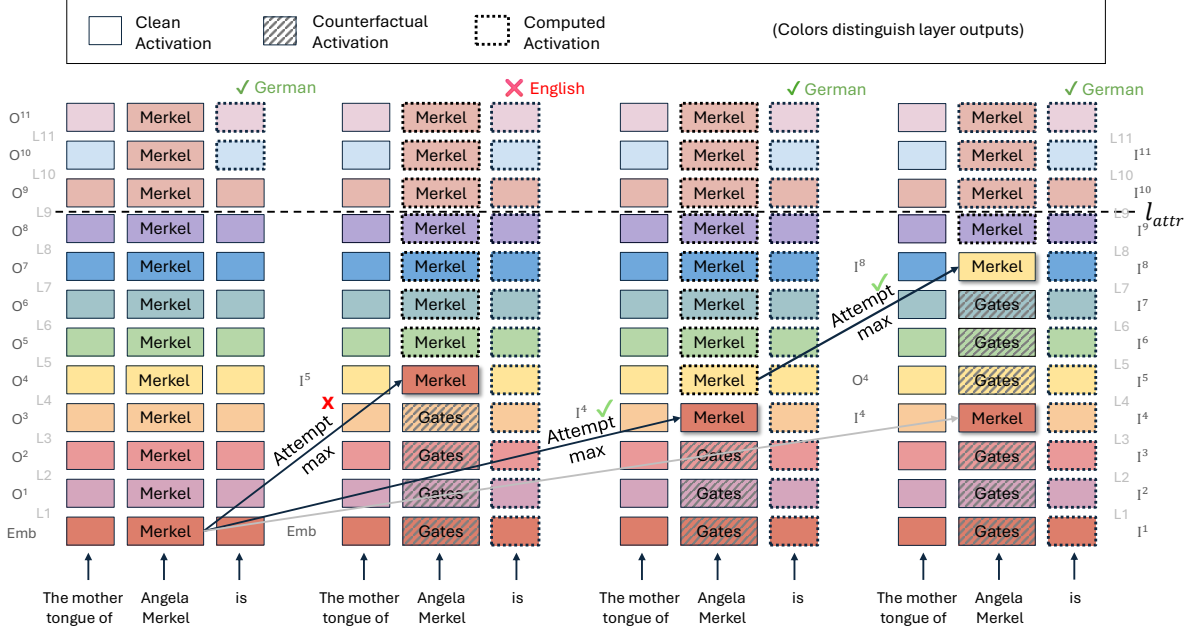


Figure 2: **Iterative Greedy Search for Minimal Computation Paths.** We illustrate the process using the original prompt “The mother tongue of Angela Merkel is” (Target: “German”) and the counterfactual prompt “The mother tongue of Bill Gates is” (Target: “English”). **Failed Attempt (Left Arrow):** The algorithm attempts a maximal jump from the embeddings to layer 5, but the model predicts the counterfactual target, marking this skip as invalid. **Successful Attempt (Middle Arrow):** The algorithm identifies layer 4 as the highest layer into which a jump recovers the correct target. **Next Iteration (Right Arrow):** With layer 4 established as a necessary step, treating its output as the new source, the search repeats to identify the next maximal jump (to layer 8).

356 a (subsection 5.1). Second, using layer ℓ_{attr} as an upper bound, we iteratively identify a
 357 minimal sequence of layers that progressively transform the entity embedding for recalling
 358 attribute a (subsection 5.2). Finally, we show that by replacing the upper bound ℓ_{attr} with
 359 $|L|$ and using the same iterative procedure, we can identify alternative minimal attribute paths
 360 (subsection 5.3).
 361
 362
 363
 364

365 5.1 Locating ℓ_{attr}

366 To bound the minimal attribute path from above, we seek a layer ℓ whose output does
 367 not need to be transformed by subsequent layers to produce attribute a . Removing a layer
 368 from the end of the path amounts to overriding its output with that of the previous layer. This
 369 gives rise to an operation we call *lock*:
 370
 371
 372

$$373 \text{lock}(O^{(e,\ell)}) := O^{(e,\ell+k)} \leftarrow_{patch} O^{(e,\ell)} \quad \forall k > 0$$

374 This is an instance of activation patching (Vig et al., 2020; Wang et al., 2022; Meng et al., 2022)
 375 which we call *activation locking*: patching the
 376

377 activations from an entity representation at layer ℓ to all higher layers of the same entity.
 378

379 For a given prompt P , we seek the earliest layer on which we can apply lock operation and
 380 retain the same correct attribute a :
 381

$$382 \ell_{attr} = \min_{\ell} \text{s.t. } M(P | \text{lock}(O^{(e,\ell)})) = M(P)$$

383 Fig 4 (Appendix A) illustrates this process.

384 5.2 Minimal Path Identification

385 Having identified ℓ_{attr} , we seek to identify a path (a sequence of layers) $[L_{\alpha_1}, \dots, L_{\alpha_k}]$, $\alpha_i <$
 386 α_{i+1} ; $\alpha_k \leq \ell_{attr}$, which forms a subset of the model’s layers constrained by the upper bound
 387 ℓ_{attr} . This sequence represents a minimal set of layers such that, when the entity embedding is
 388 processed only through these layers, the model
 389 M still produces the desired output a .
 390
 391
 392

393 To identify non-essential layers, we use counterfactual prompts: prompts of the form
 394 $(\hat{e}, r) \rightarrow \hat{a}$ where \hat{e} shares the same semantic type, number of tokens and sequence positions
 395 with e , but for which M produces an answer
 396 $\hat{a} \neq a$. We propose a heuristic: if a layer can
 397
 398

process counterfactual entity data without disrupting the final attribute prediction, this layer is not necessary for the computation. Thus, given counterfactual entity activations $O_{\text{counter}}^{(e,\ell)}$ from running M on the counterfactual prompt, and a candidate path, we define an activation patching operations called *isolate* in which: (a) the first path layer L_{α_1} receives the clean entity embedding $E_{\text{clean}}^{(e)}$; (b) each subsequent path layer L_{α_i} receives the output of the previous path layer $L_{\alpha_{i-1}}$; and (c) all layers between path layers receive the corresponding activations from the counterfactual run. Formally:

$isolate(layers = [\alpha_1, \dots, \alpha_k]) :=$

$$:= \begin{cases} I_{\text{patch}}^{(e,\ell=\alpha_1)} \leftarrow E_{\text{clean}}^{(e)} & \ell = \alpha_1 \\ I_{\text{patch}}^{(e,\ell=\alpha_i)} \leftarrow O_{\text{computed}}^{(e,\alpha_{i-1})} & \ell = \alpha_i \in \text{layers} \\ I_{\text{patch}}^{(e,\ell)} \leftarrow I_{\text{counter}}^{(e,\ell)} & \ell \notin \text{layers}, \ell < \alpha_k \end{cases}$$

where O_{computed} are the non-patched values computed in the current run. The rightmost column in Figure 2 visualizes this operation for $isolate([4,8])$: the path is isolated by patching the output of each path layer to the input of the next path layer (e.g., patching $E_{\text{clean}}^{(e)}$ into $I^{(e,4)}$ and $O_{\text{computed}}^{(e,4)}$ into $I^{(e,8)}$), while intermediate layers are patched with counterfactual values.

Given a prompt $P_{(e,r) \rightarrow a}$ and ℓ_{attr} we say that a path prefix $[\alpha_1, \dots, \alpha_k]$, $\alpha_k \leq \ell_{\text{attr}}$ is a *valid* prefix of a *computation path* iff:

$$M\left(P \mid isolate([\alpha_1, \dots, \alpha_k]), lock(O^{(e,\ell_{\text{attr}})})\right) = M(P) = a$$

We construct a minimal computational path via an iterative greedy search starting from the entity embedding³. If the embedding layer itself is sufficient (i.e., locking $E^{(e)}$ yields a), the path is empty. Otherwise, We start with an empty prefix (feeding the clean embedding through all the layers) and at each step, we extend the prefix by skipping as many layers as possible after its last layer while remaining a valid prefix. This strategy maximizes the number of skipped layers at every iteration, ensuring the final sequence is minimal. Formally,

³In our analysis (e.g., when calculating path length), we consider the entity embedding E as the implicit start of all paths.

we define the extension operation:

$$\begin{aligned} \text{extend}(\text{path}) &:= \text{path} \oplus \text{next}(\text{path}) \\ \text{next}([\alpha_1, \dots, \alpha_{k-1}]) &:= \\ &\max_{\alpha_{k-1} < \alpha_k \leq \ell_{\text{attr}}} \text{ s.t. } ([\alpha_1, \dots, \alpha_{k-1}, \alpha_k] \text{ is valid}) \end{aligned}$$

and repeatedly apply *extend*, starting from an empty path, until either the last path layer $\alpha_k = \ell_{\text{attr}}$ or $M(P|lock(O^{(e,\alpha_k)})) = MP(P) = a$.⁴ The process is illustrated in Figure 2.

5.3 Alternative Minimal Path

The procedure in subsection 5.2 searches for a minimal path that reaches the sufficient representation ℓ_{attr} . What if we relax this restriction, and allow the search procedure to construct paths that end after layer ℓ_{attr} ? Under this definition a path prefix is valid if $M(P \mid isolate([\alpha_1, \dots, \alpha_k])) = M(P) = a$ and we stop the prefix extension process when locking the last prefix layer yields the correct attribute.

This procedure reveals layers capable of executing the attribute computation even when the information integration occurs at a different depth than in the original forward pass.

6 Experiments and Results

Models. we experiment with two decoder-only transformer models: LLaMA 3.1 8B (32 layers) (Dubey et al., 2024) and Qwen 3 8B (36 layers) (Yang et al., 2025).

Data We use the CounterFact dataset (Meng et al., 2022), which provides prompts in the format described above, each paired with the correct target attribute. For each model, we select 2,000 prompts for which the target attribute is a single token and this token is correctly predicted by the model.

6.1 Main Experiments

We applied our path identification method to the dataset to extract the causal computation paths for all entity–attribute pairs. We implemented the interventions using the **NNsight** package (Fiotto-Kaufman et al., 2025). The code will be available with camera ready.

⁴We allow early stopping when $\alpha_k < \ell_{\text{attr}}$ as skipping intermediate layers may allow the path to reach sufficiency earlier than the initial upper bound. In such cases, we redefine $\ell_{\text{attr}} = \alpha_k$ for later analyses.

When is attribute knowledge available?

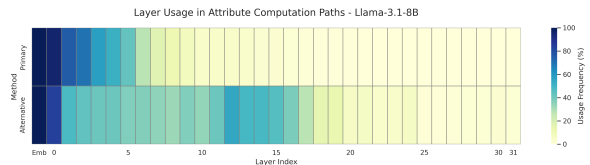
For both models, the attribute computations complete in the early-to-mid layers (average ℓ_{attr} of 4.61 for LLaMA and 7.97 for Qwen, although for some cases we reach up layer 15 and even 20 for both models (See Appendix B Figure 5(a) for the complete distribution). This aligns with previous studies on factual recall (Geva et al., 2023; Yu et al., 2025) that also place the conclusion of entity enrichment (for all attributes) at these early-to-mid layers.

How is attribute knowledge constructed?

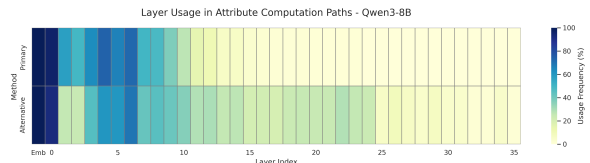
We now turn to examine the paths through which we arrive at the final representations ℓ_{attr} . The overwhelming majority of paths lengths (the number of layers in a path, including the embedding layer) are > 2 , with an average of 5.91 (LLaMA) and 7.97 (Qwen). However, many paths (33.1% of the LLaMA cases and 78.6% of Qwen’s) skip at least one layer, with an average skip size of 0.7 for LLaMA and 2.0 for Qwen (Appendix B Figure 6). For both models, the entities need to be processed by several (but not all) model layers in order for the knowledge to be available to use. The overall number of needed layers is similar for both models, but the Qwen paths are longer and sparser on average.

Are the constructed representations sufficient? To characterize the functional role of ℓ_{attr} output ($O(e, \ell_{attr})$), we performed a suite of four targeted patching interventions, ranging from overriding this specific representation to injecting it into all model layers (see Appendix E for comprehensive methodology and results).

These experiments yield two insights. First, the representation at ℓ_{attr} is strictly necessary for factual recall; overriding it and following layers with a counterfactual state disrupts the correct prediction in $> 94\%$ of cases across both models. Second, we find that the representation constructed at ℓ_{attr} on its own is often not sufficient: while injecting this representation into all layers from ℓ_{attr} in a counterfactual run in some configurations restores the desired attribute in almost 56% of the cases, it fails to do so in the remaining 44%, for which activations at intermediate path layers are also required. This indicates that the computation path often acts as a coherent functional unit, where earlier entity representations are used for the final



(a) LLaMA 3.1 8B



(b) Qwen 3 8B

Figure 3: **Layer Usage Patterns.** Aggregated layer utilization frequency for Primary (top row in each panel) vs. Alternative (bottom row) computation paths.

prediction.

Existence of alternate paths. The alternative path search confirms that models contain multiple, functionally equivalent circuits for factual computation, identifying non-identical alternative paths in 80.1% (LLaMA) and 82.7% (Qwen) of cases. Quantitative analysis (Appendix B, Figure 7) reveals that alternative paths are universally longer (mean length: 9.47 for LLaMA, 10.47 for Qwen), relying on massive non-sequential jumps (mean skip size: 6.46 layers for LLaMA, 10.4 for Qwen).

In addition, ℓ_{attr} shifts significantly deeper compared to the primary path—averaging 13.93 for LLaMA and 18.06 for Qwen (Appendix B, Figure 5). Layer usage heatmaps aggregated across paths (Figure 3) further confirm that these paths recruit a distinct, deeper set of modules, effectively bypassing the primary processing centers. Thus, while redundant paths exist, the primary path consistently offers the most compact route to attribute sufficiency.

Constructing additional paths via recombination. We test the viability of “hybrid” paths composed of early steps from a primary path and later steps from an alternative path, by concatenating a prefix from a primary path with a valid suffix from the corresponding alternative path, enforcing that the resulting path is distinct from the originals and non-trivial (the chosen prefix from the primary does not contain the corresponding initial segment of the

alternative). Each hybrid is evaluated using our counterfactual patching method (subsection 5.2); success is defined as restoring the target prediction and achieving attribute sufficiency at the final layer.

The results reveal significant modular flexibility. In LLaMA-3.1-8B, 35.4% of prompts admit at least one successful hybrid path (avg. 2.75 successful combinations per prompt), while Qwen3-8B shows a 32.8% success rate with higher redundancy (avg. 4.13 combinations). This proves the existence of additional computation paths beyond those identified in previous experiments, demonstrating that different layer sequences can perform functionally equivalent transformations. It also indicates that layer roles in the primary and alternate paths do not necessarily map one-to-one: in some cases a single layer in the primary path is accounted for by several layers in the alternate path, or vice versa.

6.2 Additional Experiments and Analyses

Effect of order. While in most prompts the relation appears before the entity (“*The mother tongue of X is*”), in few of them the order is reversed (“*X’s mother tongue is*”). Does this affect the retrieval process? We could not identify a strong trend, beyond the paths concluding on later layers in LLaMa (but not in Qwen). However, this could be due to small sample size, and could be the topic of a future study. Further details are available in Appendix C.

Paths in Individual Relations. We analyzed relation types with sufficient coverage ($N \geq 50$). We observe a general correlation between semantic specificity and computational cost (Appendix D). Broad categorical relations (e.g., *Continent*) typically require fewer layers than more specific factual retrievals (e.g., *Headquarters*), suggesting that the model often recruits longer circuits to resolve specific attributes. Furthermore, per-relation layer-usage heatmaps in Appendix D confirm that distinct relations recruit specific sets of layers beyond simple path length differences.

However, this length hierarchy is not universal: we observe significant rank discordance between models and between primary and alternate paths. The fact that relative difficulty

shifts in the alternative setting implies that “backup” mechanisms utilize distinct processing logic rather than simply mirroring primary circuits.

Relation to Entity Resolution. To examine the relationship between entity resolution (ER) (Ghandeharioun et al., 2024) and attribute computation, we apply their Patchscopes method to *multi-token* entities along our identified paths ($N_{LLaMA} = 1,972$, $N_{Qwen} = 1,945$). We use the same constant prompt used by Ghandeharioun et al., 2024 of the form “ $e_1 : d_1, \dots, e_k : d_k, x$ ” where each description d_i begins with the entity name e_i . We extract the entity’s final token representation from every layer along its computation path (for both methods) and inject it into the x position across all target layers. We then identify the first source–target layer pair, if any, that decodes the full multi-token entity name.

The results (Appendix F) reveal a strong dissociation between ER and attribute sufficiency. Explicit ER steps are detected in a minority of primary paths (12.9% for LLaMA, 9.1% for Qwen) and even fewer alternative paths (9.1% and 6.6%). This implies that detectable ER is not a universal prerequisite for factual recall; while primary paths retain slightly stronger identity traces, attribute retrieval often occurs implicitly without a distinguishable “identity state” accessible to Patchscopes.

7 Discussion and Conclusion

We introduced the concept of an Attribute-Computation Path and a novel patching protocol to identify such paths in LLMs. Our analysis confirmed that the computation of an attribute is a multi-step, non-contiguous process that frequently skips layers and requires multiple necessary steps. Furthermore, we show that models possess multiple, functionally equivalent computation paths for the same fact, emphasizing a high degree of computational redundancy.

Our findings suggests that the dominant folk view of factual storage and retrieval, in which facts are stored in specific MLP layers, is simplistic and misleading: the reality appears to be significantly more complex, with knowledge being stored in a distributed way across multiple layers, that needs to be processed together for an effective recall of an individual fact.

662 Limitations

663 While we showed redundancy in minimal at-
664 tribute computation paths exists, we did not
665 fully quantify it. Our estimation of model re-
666 dundancy is constrained by our search strategy:
667 we explicitly targeted two distinct path types
668 by defining specific upper bounds for the iter-
669 ative greedy search (ℓ_{attr} and the final model
670 layer). Consequently, our findings likely rep-
671 resent only a *lower bound* on the true degree
672 of redundancy; it is plausible that a granular
673 sweep of search bounds across all intermediate
674 layers would reveal a much larger spectrum of
675 functionally equivalent paths that our current
676 analysis did not explore.

677 Second, while our interventions demonstrate
678 that these redundant paths can be mechanis-
679 tically accessed in an intervention, it remains
680 unclear whether the model utilizes these deeper
681 layers for entity enrichment during a standard,
682 unperturbed forward pass, or if they represent
683 degenerate remnants of the training process.

684 References

685 Emmanuel Ameisen, Jack Lindsey, Adam Pearce,
686 Wes Gurnee, Nicholas L. Turner, Brian Chen,
687 Craig Citro, David Abrahams, Shan Carter, Basil
688 Hosmer, Jonathan Marcus, Michael Sklar, Adly
689 Templeton, Trenton Bricken, Callum McDougall,
690 Hoagy Cunningham, Thomas Henighan, Adam
691 Jermyn, Andy Jones, and 8 others. 2025. *Cir-
692 cuit tracing: Revealing computational graphs in
693 language models*. *Transformer Circuits Thread*.

694 Bilal Chughtai, Alan Cooney, and Neel Nanda. 2024.
695 Summing up the facts: Additive mechanisms
696 behind factual recall in llms. *arXiv preprint
697 arXiv:2402.07321*.

698 Roi Cohen, Mor Geva, Jonathan Berant, and
699 Amir Globerson. 2023. Crawling the inter-
700 nal knowledge-base of language models. *arXiv
701 preprint arXiv:2301.12810*.

702 Damai Dai, Li Dong, Yaru Hao, Zhifang Sui,
703 Baobao Chang, and Furu Wei. 2021. Knowl-
704 edge neurons in pretrained transformers. *arXiv
705 preprint arXiv:2104.08696*.

706 Abhimanyu Dubey, Abhinav Jauhri, Abhinav
707 Pandey, Abhishek Kadian, Ahmad Al-Dahle,
708 Aiesha Letman, Akhil Mathur, Alan Schelten,
709 Amy Yang, Angela Fan, and 1 others. 2024. The
710 llama 3 herd of models. *arXiv e-prints*, pages
711 arXiv–2407.

712 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom
713 Henighan, Nicholas Joseph, Ben Mann, Amanda

Askill, Yuntao Bai, Anna Chen, Tom Conerly, 714
Nova DasSarma, Dawn Drain, Deep Ganguli, 715
Zac Hatfield-Dodds, Danny Hernandez, Andy 716
Jones, Jackson Kernion, Liane Lovitt, Kamal 717
Ndousse, and 6 others. 2021. A mathemati- 718
cal framework for transformer circuits. *Trans- 719
former Circuits Thread*. [https://transformer- 720
circuits.pub/2021/framework/index.html](https://transformer-circuits.pub/2021/framework/index.html). 721

Jaden Fried Fiotto-Kaufman, Alexander Russell 722
Loftus, Eric Todd, Jannik Brinkmann, Koyena 723
Pal, Dmitrii Troitskii, Michael Ripa, Adam 724
Belfki, Can Rager, Caden Juang, Aaron Mueller, 725
Samuel Marks, Arnab Sen Sharma, Francesca 726
Lucchetti, Nikhil Prakash, Carla E. Brodley, Ar- 727
jun Guha, Jonathan Bell, Byron C Wallace, and 728
David Bau. 2025. *NNsight and NDIF: Democ- 729
ratizing access to foundation model internals*.
In *The Thirteenth International Conference on 730
Learning Representations*. 731

Mor Geva, Jasmijn Bastings, Katja Filippova, and 732
Amir Globerson. 2023. Dissecting recall of fac- 733
tual associations in auto-regressive language mod- 734
els. *arXiv preprint arXiv:2304.14767*. 735

Mor Geva, Avi Caciularu, Kevin Wang, and Yoav 736
Goldberg. 2022. Transformer feed-forward layers
build predictions by promoting concepts in the 737
vocabulary space. In *Proceedings of the 2022 con- 738
ference on empirical methods in natural language 739
processing*, pages 30–45. 740

Mor Geva, Roei Schuster, Jonathan Berant, and 741
Omer Levy. 2021. Transformer feed-forward lay- 742
ers are key-value memories. In *Proceedings of 743
the 2021 Conference on Empirical Methods in 744
Natural Language Processing*, pages 5484–5495. 745

Asma Ghandeharioun, Avi Caciularu, Adam Pearce, 746
Lucas Dixon, and Mor Geva. 2024. Patchscopes:
A unifying framework for inspecting hidden rep- 747
resentations of language models. *arXiv preprint 748
arXiv:2401.06102*. 749

Nicholas Goldowsky-Dill, Chris MacLeod, Lucas 750
Sato, and Aryaman Arora. 2023. Localizing
model behavior with path patching. *arXiv 751
preprint arXiv:2304.05969*. 752

Daniela Gottesman and Mor Geva. 2024. Estim- 753
ating knowledge in large language models with- 754
out generating a single token. *arXiv preprint 755
arXiv:2406.12673*. 756

Wes Gurnee and Max Tegmark. 2023. Language 757
models represent space and time. *arXiv preprint 758
arXiv:2310.02207*. 759

Peter Hase, Mohit Bansal, Been Kim, and Asma 760
Ghandeharioun. 2023. Does localization inform 761
editing? surprising differences in causality-based 762
localization vs. knowledge editing in language 763
models. *Advances in Neural Information Pro- 764
cessing Systems*, 36:17643–17668. 765

770	Evan Hernandez, Arnab Sen Sharma, Tal Haklay,	Yunzhi Yao, Ningyu Zhang, Zekun Xi, Mengru	824
771	Kevin Meng, Martin Wattenberg, Jacob Andreas,	Wang, Ziwen Xu, Shumin Deng, and Huajun	825
772	Yonatan Belinkov, and David Bau. 2023. Linear-	Chen. 2024. Knowledge circuits in pretrained	826
773	ity of relation decoding in transformer language	transformers. <i>Advances in Neural Information</i>	827
774	models. <i>arXiv preprint arXiv:2308.09124</i> .	<i>Processing Systems</i> , 37:118571–118602.	828
775	Zhengbao Jiang, Frank F Xu, Jun Araki, and Gra-	Zeping Yu and Sophia Ananiadou. 2024. Neuron-	829
776	ham Neubig. 2020. How can we know what lan-	level knowledge attribution in large language	830
777	guage models know? <i>Transactions of the Associ-</i>	models. In <i>Proceedings of the 2024 Conference</i>	831
778	<i>ation for Computational Linguistics</i> , 8:423–438.	<i>on Empirical Methods in Natural Language Pro-</i>	832
779	Shahar Katz, Yonatan Belinkov, Mor Geva, and	<i>cessing</i> , pages 3267–3280.	833
780	Lior Wolf. 2024. Backward lens: Projecting lan-	Zeping Yu, Yonatan Belinkov, and Sophia Anani-	834
781	guage model gradients into the vocabulary space.	adou. 2025. Back attention: Understanding and	835
782	<i>arXiv preprint arXiv:2402.12865</i> .	enhancing multi-hop reasoning in large language	836
783	Thomas McGrath, Matthew Rahtz, Janos Kra-	models. <i>arXiv preprint arXiv:2502.10835</i> .	837
784	mar, Vladimir Mikulik, and Shane Legg. 2023.		
785	The hydra effect: Emergent self-repair in lan-		
786	guage model computations. <i>arXiv preprint</i>		
787	<i>arXiv:2307.15771</i> .		
788	Kevin Meng, David Bau, Alex Andonian, and		
789	Yonatan Belinkov. 2022. Locating and editing		
790	factual associations in gpt. <i>Advances in neural</i>		
791	<i>information processing systems</i> , 35:17359–17372.		
792	Neel Nanda, Senthooran Rajamanoharan, János		
793	Kramár, and Rohin Shah. 2023. Fact finding:		
794	Attempting to reverse-engineer factual recall on		
795	the neuron level . AI Alignment Forum.		
796	Chris Olah, Nick Cammarata, Ludwig Schubert,		
797	Gabriel Goh, Michael Petrov, and Shan Carter.		
798	2020. Zoom in: An introduction to circuits . <i>Dis-</i>		
799	<i>till</i> . https://distill.pub/2020/circuits/zoom-in .		
800	Fabio Petroni, Tim Rocktäschel, Patrick Lewis,		
801	Anton Bakhtin, Yuxiang Wu, Alexander H		
802	Miller, and Sebastian Riedel. 2019. Language		
803	models as knowledge bases? <i>arXiv preprint</i>		
804	<i>arXiv:1909.01066</i> .		
805	Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov,		
806	Sharon Qian, Daniel Nevo, Yaron Singer, and		
807	Stuart Shieber. 2020. Investigating gender bias		
808	in language models using causal mediation anal-		
809	ysis. <i>Advances in neural information processing</i>		
810	<i>systems</i> , 33:12388–12401.		
811	Kevin Wang, Alexandre Variengien, Arthur Conmy,		
812	Buck Shlegeris, and Jacob Steinhardt. 2022. In-		
813	terpretability in the wild: a circuit for indirect ob-		
814	ject identification in gpt-2 small. <i>arXiv preprint</i>		
815	<i>arXiv:2211.00593</i> .		
816	Zijian Wang and Chang Xu. 2025. Functional ab-		
817	straction of knowledge recall in large language		
818	models. <i>arXiv preprint arXiv:2504.14496</i> .		
819	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,		
820	Binyuan Hui, Bo Zheng, Bowen Yu, Chang		
821	Gao, Chengen Huang, Chenxu Lv, and 1 others.		
822	2025. Qwen3 technical report. <i>arXiv preprint</i>		
823	<i>arXiv:2505.09388</i> .		

838	A Visualization of ℓ_{attr} Identification	D Relation-Specific Path Analysis	884
839	Figure 4 provides a detailed visualization of the	In Section 6.2, we discussed how path length	885
840	ℓ_{attr} identification process, using the prompt	varies by relation. This section provides a com-	886
841	“The mother tongue of Angela Merkel is” (Tar-	prehensive breakdown of these variations. First,	887
842	get: “German”). The figure illustrates this	Figure 10 illustrates the mean path length	888
843	step-by-step procedure in which we lock the en-	across all relations, sorted by the primary	889
844	tity representation at each layer—starting from	path length in LLaMA-3.1-8B. Following this	890
845	the embeddings and moving upwards—until the	overview, Figures 11 and 12 present the granu-	891
846	operation successfully reproduces the correct	lar layer-usage distributions for these relations,	892
847	target token.	comparing LLaMA-3.1-8B and Qwen3-8B.	893
848	B Extended Quantitative Analysis	E Sufficiency of constructed	894
849	In this section, we provide the detailed statisti-	representations	895
850	cal breakdown of the identified paths. Figure	To validate that the identified computational	896
851	5 illustrates the distribution of the ℓ_{attr} for	path acts as the causal driver of factual recall,	897
852	both the primary and alternative strategies, as	we analyze the functional role of the entity rep-	898
853	well as the layers skipping ratio in the paths.	resentation at ℓ_{attr} . We conduct four targeted	899
854	Figure 7 offers a direct comparison between	patching experiments to test the representa-	900
855	the primary and alternative paths, quantifying	tion’s necessity and sufficiency:	901
856	the increase in depth and path length when		
857	the model is forced to utilize later layers for	1. State Knockout: We execute the min-	902
858	attribute computation.	imal path up to ℓ_{attr} , then immediately	903
		overwrite the entity representation with a	904
859	C Analysis by Prompt Structure	counterfactual one. This tests if the con-	905
860	(Entity Position)	structed state is necessary for prediction.	906
861	Here, we provide detailed results of the experi-	2. Forward Propagation: We inject the	907
862	ments. First, we classified each prompt as one	representation from ℓ_{attr} into all subse-	908
863	of the following types:	quent layers inputs ($l > \ell_{attr}$) while keep-	909
		ing previous layers counterfactual. This	910
864	• entity_first: The entity appears before	tests if the final state alone, without the	911
865	the relation (e.g., “Albert Einstein was	path’s history, is sufficient to drive the	912
866	born in...”).	upper model layers.	913
867	• relation_first: The relation appears be-	3. Path + Continuation: We inject the	914
868	fore the entity (e.g., “The birthplace of	representation from ℓ_{attr} into all path lay-	915
869	Albert Einstein is...”).	ers and to all higher layers ($l > \ell_{attr}$),	916
		leaving only off-path intermediate layers	917
870	To perform this classification, we used GPT-	as counterfactual.	918
871	4o. We provided the model with the formatted	4. Global Broadcast: We inject the repre-	919
872	prompt (e.g., “The official religion of Edwin of	sentation from ℓ_{attr} into <i>all</i> model layers.	920
873	Northumbria is”) and requested a classification		
874	using the few-shot prompt structure in Figure	This procedure is applied to both the primary	921
875	8.	and alternative paths.	922
876	Figure 9 shows the comparison of ℓ_{attr} and	Results. Both models show a strong reliance	923
877	path length between these two prompt types	on the output of ℓ_{attr} identified by our methods.	924
878	for the primary and the alternative paths. We	For LLaMA 3.1 8B, replacing the representa-	925
879	restricted our comparison to relations that had	tion with that from the counterfactual prompt	926
880	at least 50 samples for both structure types to-	disrupted the correct prediction in 99.50% of	927
881	gether across both models. The number above	cases in the primary method and in 94.95% of	928
882	each bar indicates the sample size used to cal-	cases in the alternative method. For Qwen 3	929
883	culate the average.	8B, we observed a similar trend, with prediction	930

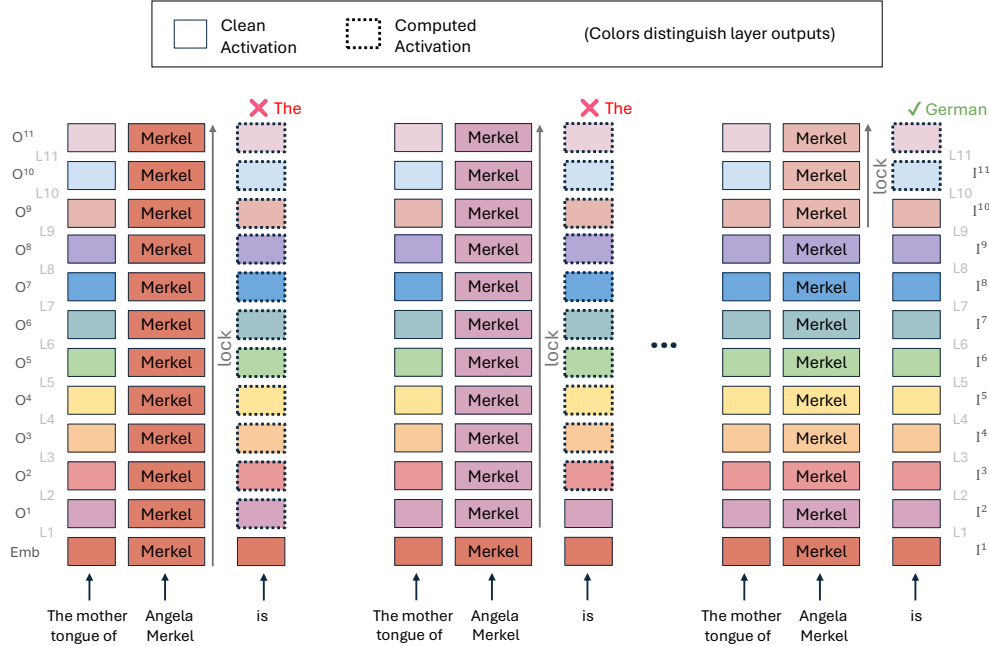


Figure 4: ℓ_{attr} **Identification**. This process identifies ℓ_{attr} , the first layer where the entity representation is robust enough to elicit the target attribute without further processing. The labeled squares represent the entity representation at different layers. **Left & Middle (Testing an Insufficient Layer)**: We lock the “Angela Merkel” representation at an early layer. The model fails to retrieve the correct attribute, predicting a generic or nonsensical token (e.g., “The”), indicating the representation is not yet sufficient. **Right (Testing a Sufficient Layer)**: We lock the representation at a deeper layer. The model successfully predicts “German”, identifying this layer as ℓ_{attr} .

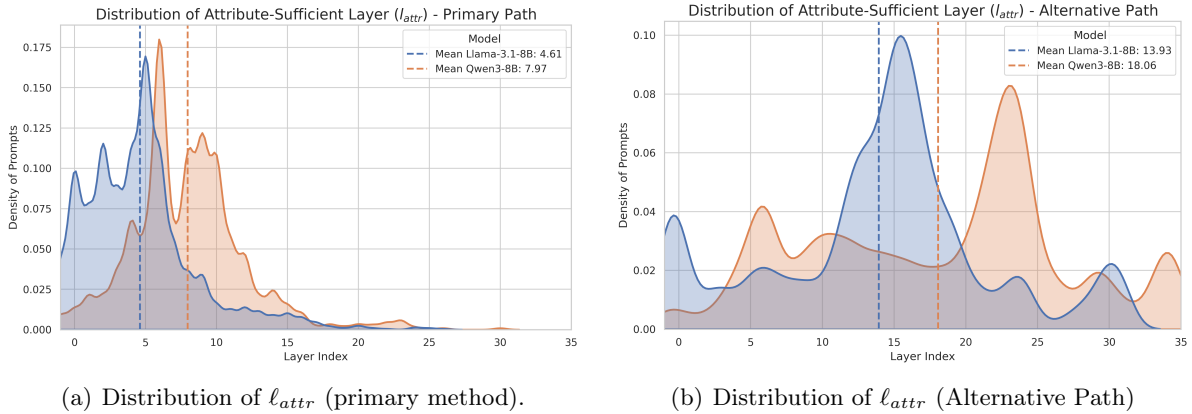
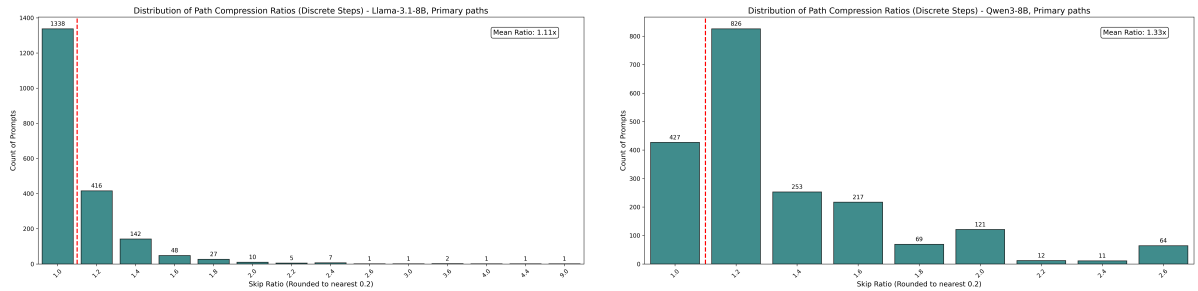


Figure 5: Analysis of minimal computation paths for the primary method. (a) The distribution of ℓ_{attr} , showing both models complete attribute computation in the early-to-mid layers. (b) The distribution of the final ℓ_{attr} for the alternative path, which is significantly deeper and wider than the primary path’s distribution.

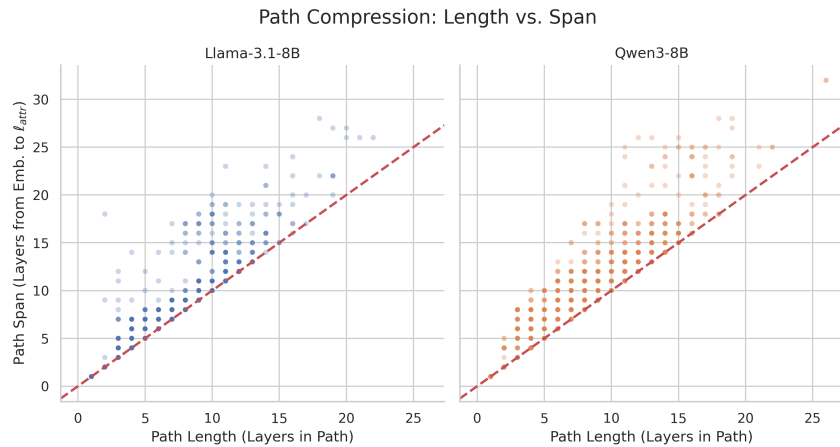
931 disruption in 99.40% and 94.35% of cases for
 932 the primary and alternative methods, respec-
 933 tively. These findings demonstrate that in both
 934 models, $O^{(e, \ell_{attr})}$ is not a passive byproduct of
 935 earlier processing, but a meaningful functional
 936 state that the model actively uses for factual
 937 recall.

Table 1 details the performance impact of
 each intervention. Both models show a strong
 reliance on the output of ℓ_{attr} identified by our
 methods. In the **State Knockout** setting, re-
 placing the representation disrupted the correct
 prediction in $> 94\%$ of cases across both mod-
 els and methods, confirming that $O^{(e, \ell_{attr})}$ is a

938
 939
 940
 941
 942
 943
 944



(a) Distribution of Path Compression Ratios (LLaMA). (b) Distribution of Path Compression Ratios (Qwen)



(c) Path length vs. path span (primary method).

Figure 6: Analysis of minimal computation paths for the primary method. (a) Distribution of Path Compression Ratios for LLaMA 3.1 8B. The ratio is computed as the path depth divided by the number of steps. The first bar (ratio=1.0) represents purely sequential paths. The red dashed line acts as a boundary, separating these sequential paths from the compressed paths (ratio > 1.0) to the right. Ratios > 1.0 are rounded to the nearest 0.2 for visualization. (b) The corresponding distribution for Qwen 3 8B. (c) A comparison of actual path length vs. path span for the primary paths. Dots below the $y = x$ line represent “compressed” paths that skip layers.

necessary functional state.

However, the other experiments reveal that the final state alone is often insufficient; the model requires the intermediate path computations to fully restore the prediction (see Table 1).

F Entity Resolution Detection Results

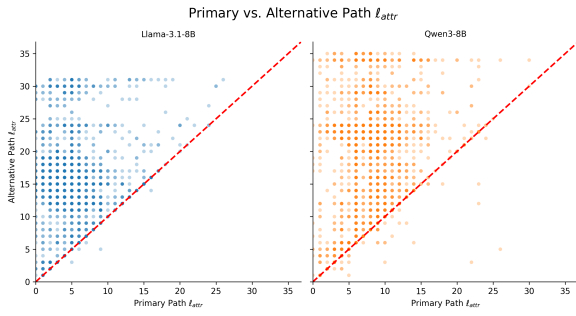
In Section 6.2, we analyzed the relationship between the attribute computation path and the explicit resolution of the entity’s identity.

To perform this analysis, we employed the Patchscopes method using the identical constant prompt configuration from Ghandeharioun et al., 2024. The specific few-shot prompt is:

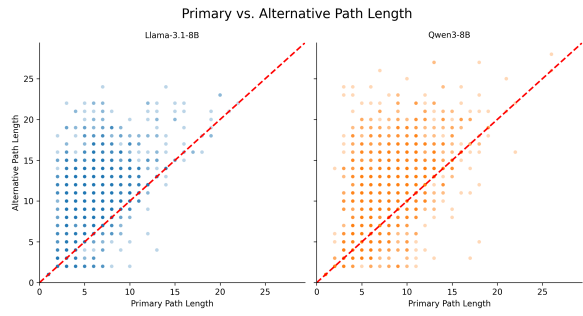
“Syria: Syria is a country in the Middle East, Leonardo DiCaprio: Leonardo DiCaprio is an American actor, Samsung: Samsung is a South Korean multinational corporation, x ”

where x serves as the placeholder for the injected representation.

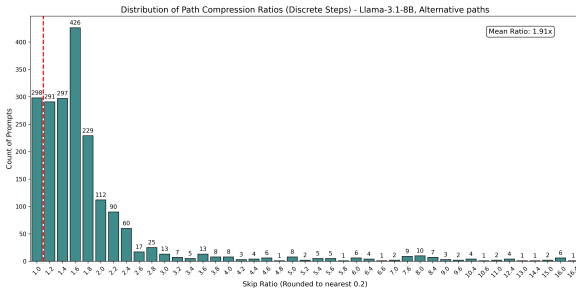
Figure 13 presents the quantitative results of this analysis. It details the success rate of detecting Entity Resolution (ER) along both the primary and alternative computation paths for LLaMA-3.1-8B and Qwen3-8B.



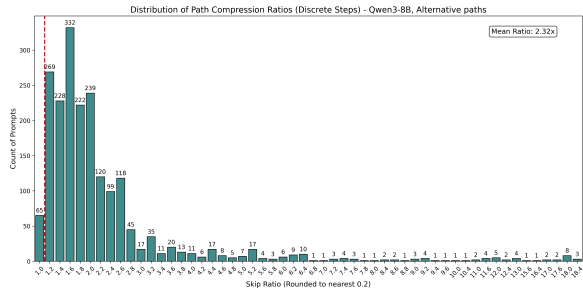
(a) ℓ_{attr} : Primary vs. Alternative Path



(b) Path Length: Primary vs. Alternative Path



(c) Distribution of Path Compression Ratios (LLaMA).



(d) Distribution of Path Compression Ratios (Qwen)

Figure 7: Analysis of the alternative computation path. (a) Comparison of ℓ_{attr} for the primary path (x-axis) vs. the alternative path (y-axis). Most points lie above the $y = x$ line, showing the alternative path is deeper. (b) Comparison of path length. Points are again mostly above the $y = x$ line, indicating the alternative path is also longer (more computational steps). (c) Distribution of Path Compression Ratios for LLaMA 3.1 8B. The ratio is computed as the path depth divided by the number of steps. The first bar (ratio=1.0) represents purely sequential paths. Ratios > 1.0 are rounded to the nearest 0.2 for visualization. (d) The corresponding distribution for Qwen 3 8B.

Example 1:

Sentence: The official religion of Edwin of Northumbria is

Answer: relation

Example 2:

Sentence: In Nykarleby, the language spoken is

Answer: entity

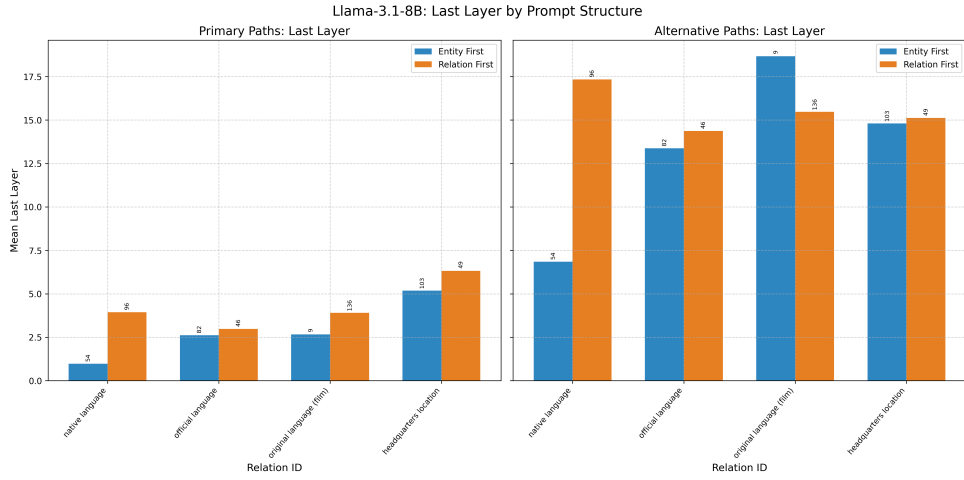
Now, analyze the following sentence.

What appears first: the relation or the entity?

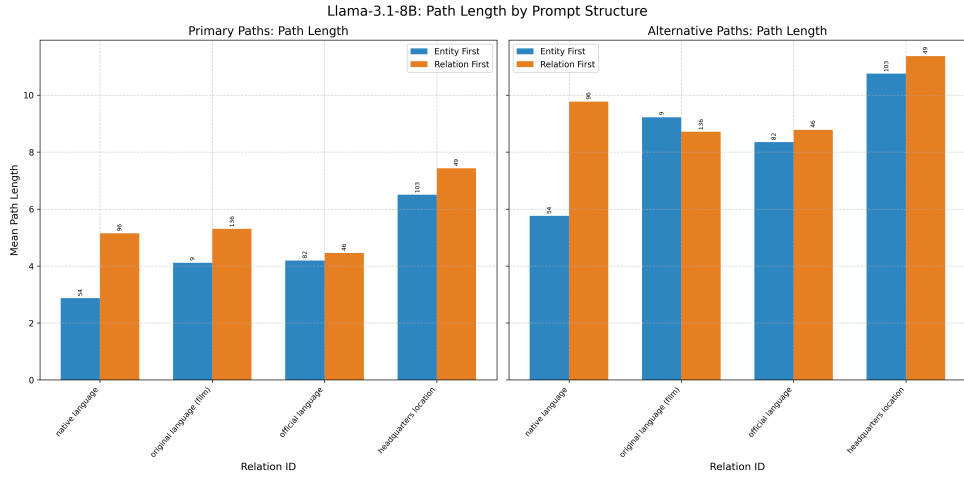
Sentence: {prompt}

Answer with only 'relation' or 'entity'.

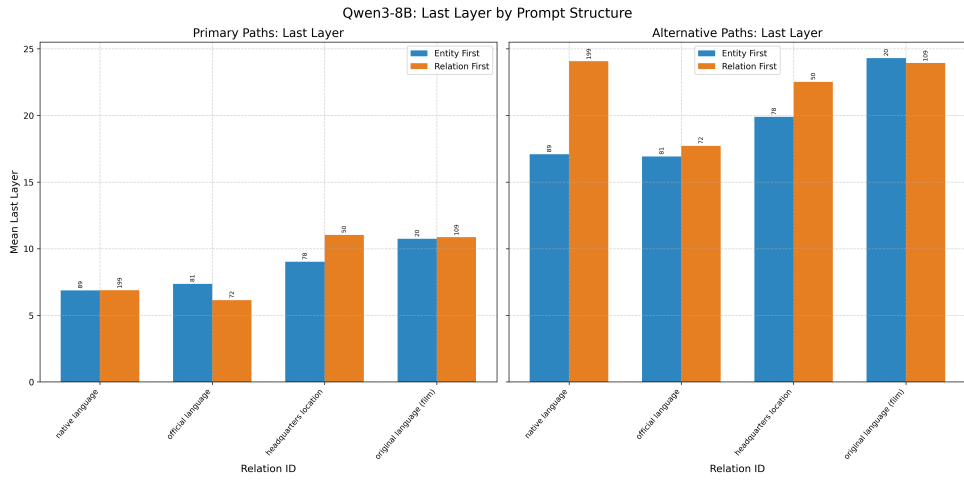
Figure 8: The few-shot prompt used to classify the ordering of entity and relation in the query.



(a) ℓ_{attr} vs. Entity Position - LLaMA 3.1 8B.

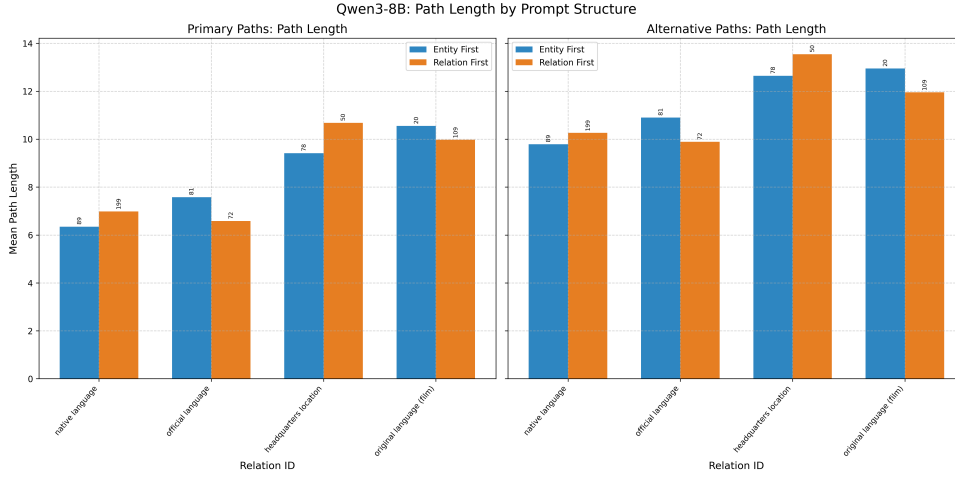


(b) Path Length vs. Entity Position - LLaMA 3.1 8B.



(c) ℓ_{attr} vs. Entity Position - Qwen 3 8B.

Figure 9: Analysis of ℓ_{attr} and Path Lengths by Prompt Structure.



(a) Path Length vs. Entity Position - Qwen 3 8B.

Figure 9: Analysis of ℓ_{attr} and Path Lengths by Prompt Structure. (continued)

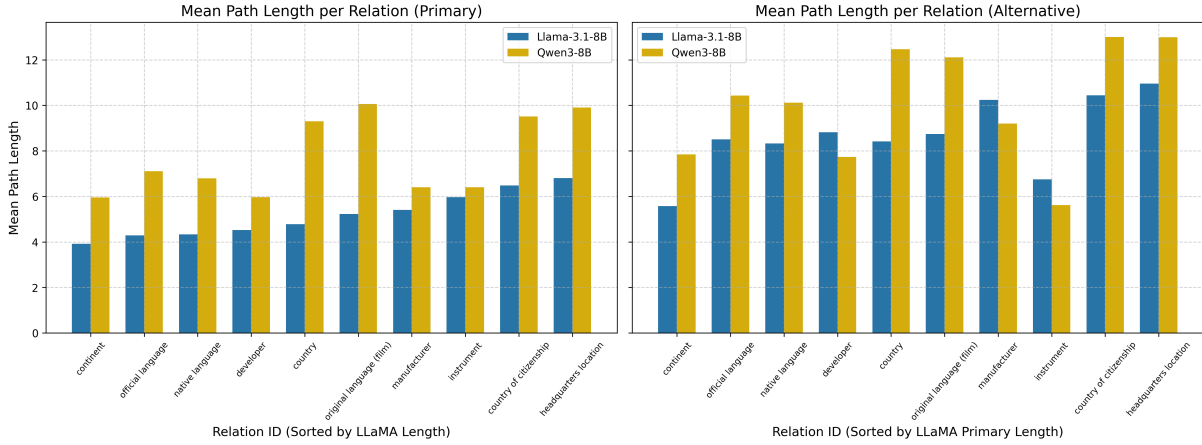


Figure 10: Mean path length per relation. Relations are sorted by LLaMA primary path length.

Intervention Type	LLaMA 3.1 8B		Qwen 3 8B	
	Primary	Alternative	Primary	Alternative
State Knockout	99.50%	94.95%	99.40%	94.35%
Forward Prop	44.50%	76.60%	56.05%	76.85%
Path + Cont.	34.45%	48.25%	53.20%	64.55%
Global Broadcast	37.00%	61.95%	62.65%	78.65%

Table 1: Prediction failure rates across functional analysis experiments. All values indicate the percentage of cases where the model failed to output the correct attribute.

Layer Usage Distribution

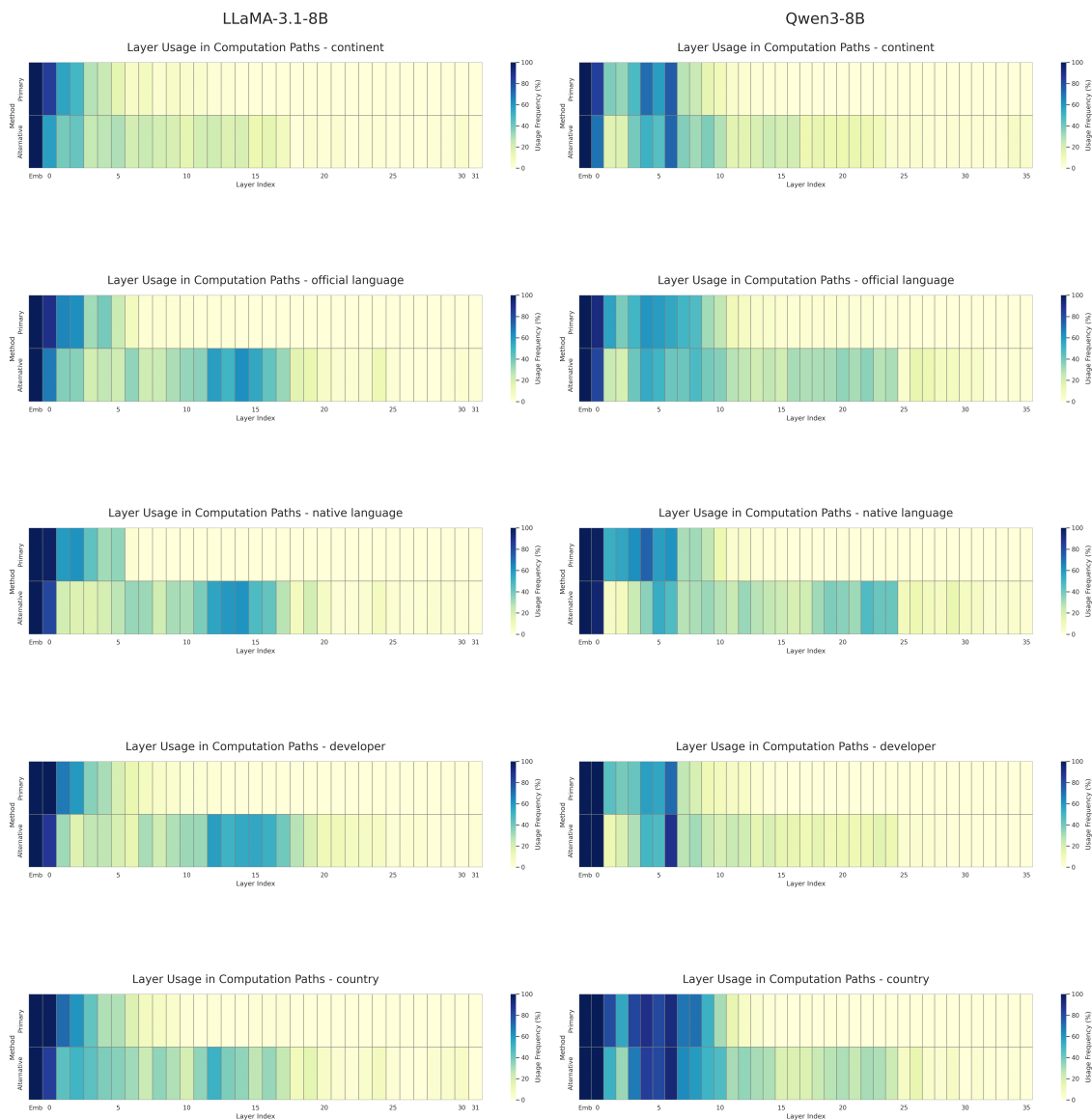


Figure 11: **Layer Usage per Relation (Part 1)**. Comparison of Relations 1–5. LLaMA-3.1-8B (Left) vs. Qwen3-8B (Right). Warmer colors indicate higher usage probability.

Layer Usage Distribution

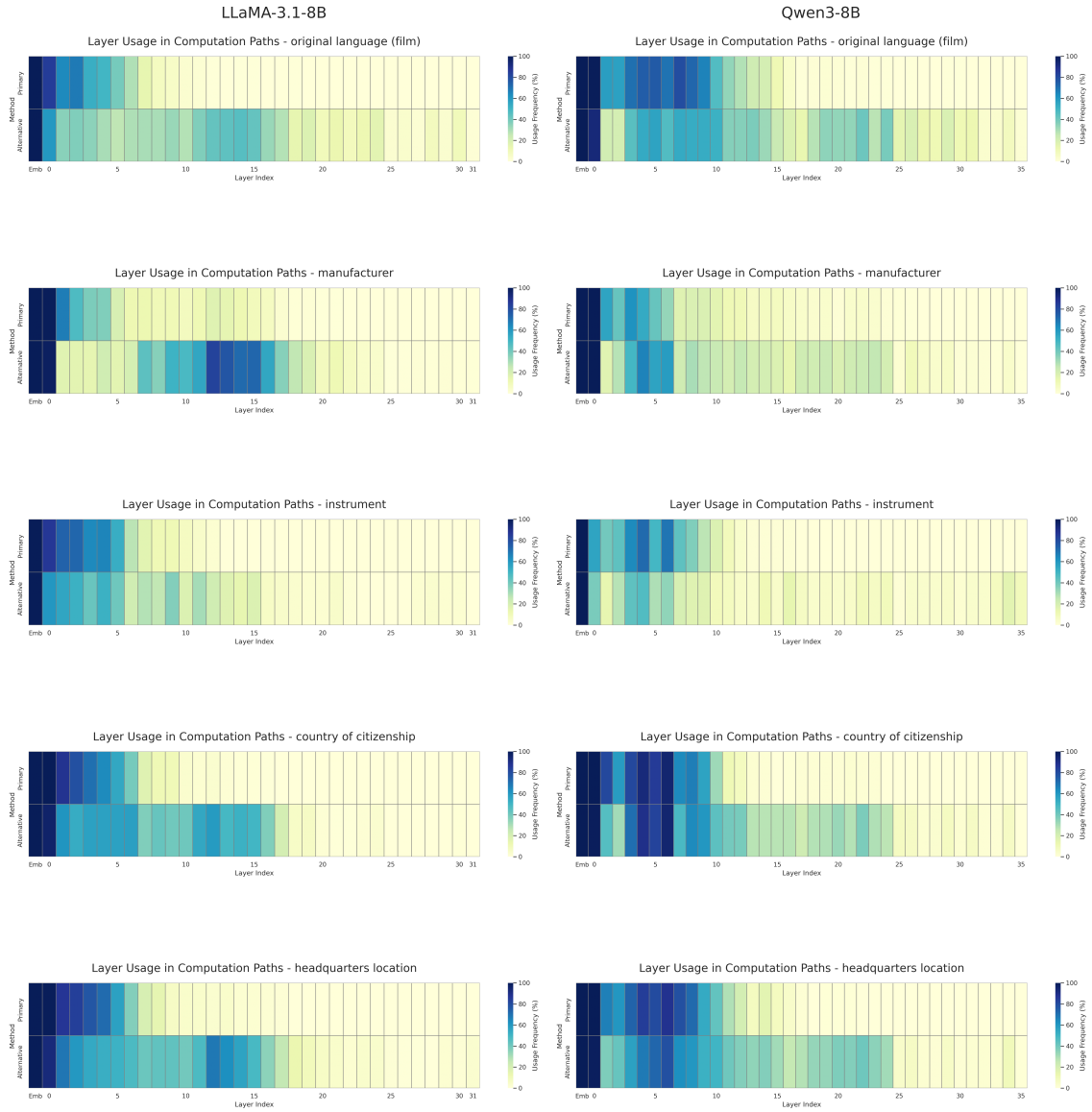


Figure 12: **Layer Usage per Relation (Part 2).** Comparison of Relations 6–10.

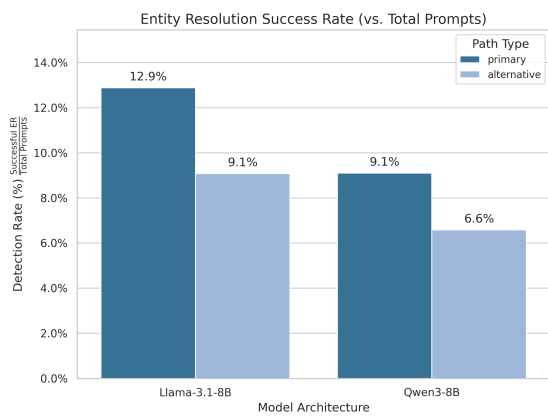


Figure 13: **Entity Resolution (ER) Detection Success Rate.** The figure shows the percentage of prompts where entity resolution was successfully detected along the minimal computation path (relative to the prompts analyzed for each method and model).