
Do LLMs dream of elephants (when told not to)?

Latent concept association and associative memory in transformers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large Language Models (LLMs) have the capacity to store and recall facts.
2 Through experimentation with open-source models, we observe that this ability
3 to retrieve facts can be easily manipulated by changing contexts, even without
4 altering their factual meanings. These findings highlight that LLMs might behave
5 like an associative memory model where certain tokens in the contexts serve as
6 clues to retrieving facts. We mathematically explore this property by studying how
7 transformers, the building blocks of LLMs, can complete such memory tasks. We
8 study a simple latent concept association problem with a one-layer transformer
9 and we show theoretically and empirically that the transformer gathers information
10 using self-attention and uses the value matrix for associative memory.

11 1 Introduction

12 What is the first thing that would come to mind if you were asked *not* to think of an elephant? Chances
13 are, you would be thinking about elephants. What if we ask the same thing to Large Language Models
14 (LLMs)? Obviously, one would expect the outputs of LLMs to be heavily influenced by tokens in the
15 context [Bro+20]. Could such influence potentially prime LLMs into changing outputs in a nontrivial
16 way? To gain a deeper understanding, we focus on one specific task called fact retrieval [Men+22;
17 Men+23] where expected output answers are given. LLMs, which are trained on vast amounts of
18 data, are known to have the capability to store and recall facts [Men+22; Men+23; DCAT21; Mit+21;
19 Mit+22; Dai+21]. This ability raises natural questions: *How robust is fact retrieval, and to what extent*
20 *does it depend on semantic meanings within contexts? What does it reveal about memory in LLMs?*

21 In this paper, we first demonstrate that fact retrieval is not robust and LLMs can be easily fooled by
22 varying contexts. For example, when asked to complete “The Eiffel Tower is in the city of”, GPT-2
23 [Rad+19] answers with “Paris”. However, when prompted with “The Eiffel Tower is not in Chicago.
24 The Eiffel Tower is in the city of”, GPT-2 responds with “Chicago”. See Figure 1 for more examples,
25 including Gemma and LLaMA. On the other hand, humans do not find the two sentences factually
26 confusing and would answer “Paris” in both cases. We call this phenomenon *context hijacking*.
27 Importantly, these findings suggest that LLMs might behave like an associative memory model. In
28 which, tokens in contexts guide the retrieval of memories, even if such associations formed are not
29 inherently semantically meaningful.

30 This associative memory perspective raises further interpretability questions about how LLMs form
31 such associations. Answering these questions can facilitate the development of more robust LLMs.
32 Unlike classical models of associative memory in which distance between memory patterns are
33 measured directly and the associations between inputs and outputs are well-specified, fact retrieval
34 relies on a more nuanced notion of similarity measured by latent (unobserved) semantic concepts.

Context Hijacking		
MODEL	CONTEXT	NEXT TOKEN
All models	The Eiffel Tower is in the city of	Paris
GPT-2 / Gemma-2B	The Eiffel Tower is not in Chicago. Therefore, the Eiffel Tower is in the city of	Chicago
Gemma-2B-IT	The Eiffel Tower is not in Chicago. However, the Chicago river is in Chicago. Therefore, the Eiffel Tower is in the city of	Chicago
LLaMA-7B	The Eiffel Tower is not in Chicago. The Eiffel Tower is not in Chicago. The Eiffel Tower is not in Chicago. The Eiffel Tower is not in Chicago. The Eiffel Tower is not in Chicago. The Eiffel Tower is not in Chicago. Therefore, the Eiffel Tower is in the city of	Chicago

Figure 1: Examples of context hijacking for various LLMs, showcasing that fact retrieval is not robust.

35 To model this, we propose a synthetic task called *latent concept association* where the output token is
 36 closely related to sampled tokens in the context but wherein similarity is measured via a latent space
 37 of semantic concepts. We then investigate how a one-layer transformer [Vas+17], a fundamental
 38 component of LLMs, can tackle this memory retrieval task in which various context distributions
 39 correspond to distinct memory patterns. We demonstrate that the transformer accomplishes the
 40 task in two stages: The self-attention layer gathers information, while the value matrix functions
 41 as associative memory. Moreover, low-rank structure also emerges in the embedding space of trained
 42 transformers. These findings provide additional theoretical validation for numerous existing low-rank
 43 editing and fine-tuning techniques [Men+22; Hu+21].

44 **Contributions** Specifically, we make the following contributions:

- 45 1. We systematically demonstrate context hijacking for various open source LLM models
 46 including GPT-2 [Rad+19], LLaMA-2 [Tou+23] and Gemma [Tea+24], which show
 47 that fact retrieval can be misled by contexts (Appendix B), reaffirming that LLMs lack
 48 robustness to context changes [Shi+23; Pet+20; CSH22; Yor+23; PE21].
- 49 2. We propose a synthetic memory retrieval task termed latent concept association, allowing
 50 us to analyze how transformers can accomplish memory recall (Section 3). Unlike
 51 classical models of associative memory, our task creates associations in a latent, semantic
 52 concept space as opposed to directly between observed tokens. This perspective is crucial
 53 to understanding how transformers can solve fact retrieval problems by implementing
 54 associative memory based on similarity in the latent space.
- 55 3. We theoretically (Section 4) and empirically (Appendix D) study trained transformers on
 56 this latent concept association problem, showing that self-attention is used to aggregate
 57 information while the value matrix serves as associative memory. And moreover, we
 58 discover that the embedding space can exhibit a low-rank structure, offering additional
 59 support for existing editing and fine-tuning methods [Men+22; Hu+21].

60 2 Context hijacking in LLMs

61 We systematically examine the phenomenon of context hijacking with the COUNTERFACT dataset
 62 [Men+22]. Due to the page limit, more details can be found in Appendix B. Overall, the experimental
 63 results show that even prepending contexts with factually correct sentences can cause LLMs to output
 64 incorrect tokens.

65 Context hijacking indicates that fact retrieval in LLMs is not robust and that accurate fact recall
 66 does not necessarily depend on the semantics of the context. As a result, one hypothesis is to view
 67 LLMs as an associative memory model where special tokens in contexts, associated with the fact,
 68 provide partial information or clues to facilitate memory retrieval [Zha23]. To better understand
 69 this perspective, we design a synthetic memory retrieval task to evaluate how the building blocks of
 70 LLMs, transformers, can solve it.

71 3 Problem setup

72 In the context of LLMs, fact or memory retrieval, can be modeled as a next token prediction problem.
73 Given a context (e.g., “The capital of France is”), the objective is to accurately predict the next token
74 (e.g., “Paris”) based on the factual relation between context and the following token.

75 Previous papers [Ram+20; Mil+22; BP21; Zha23] have studied the connection between attention and
76 autoassociative and heteroassociative memory. For autoassociative memory, contexts are modeled as
77 a set of existing memories and the goal of self-attention is to select the closest one or approximations
78 to it. On top of this, heteroassociative memory [Mil+22; BP21] has an additional projection to remap
79 each output to a different one, whether within the same space or otherwise. In both scenarios, the
80 goal is to locate the closest pattern within the context when provided with a query (up to a remapping
81 if it’s heteroassociative).

82 Fact retrieval, on the other hand, does not strictly follow this framework. The crux of the issue
83 is that the output token is not necessarily close to any particular token in the context but rather a
84 combination of them and the “closeness” is intuitively measured by latent semantic concepts. For
85 example, consider context sentence “The capital of France is” with the output “Paris”. Here, none of
86 the tokens in the context directly corresponds to the word “Paris”. Yet some tokens contain partial
87 information about “Paris”. Intuitively, “capital” aligns with the “isCapital” concept of “Paris”, while
88 “France” corresponds to the “isFrench” concept linked to “Paris” where all the concepts are latent. To
89 model such phenomenon, we propose a synthetic task called *latent concept association* where the
90 output token is closely related to tokens in the context and similarity is measured via the latent space.

91 3.1 Latent concept association

92 We propose a synthetic prediction task where for each output token y , tokens in the context (denoted
93 by x) are sampled from a conditional distribution given y . Tokens that are similar to y will be
94 favored to appear more in the context, except for y itself. The task of latent concept association is to
95 successfully retrieve the token y given samples from $p(x|y)$. The synthetic setup simplifies by not
96 accounting for the sequential nature of language, a choice supported by previous experiments on
97 context hijacking (Appendix B). We formalize this task below.

98 To measure similarity, we define a latent space. Here, the latent space is a collection of m binary
99 latent variables Z_i . These could be viewed as semantic concept variables. Let $Z = (Z_1, \dots, Z_m)$ be
100 the corresponding random vector, z be its realization, and \mathcal{Z} be the collection of all latent binary
101 vectors. For each latent vector z , there’s one associated token $t \in [V] = \{0, \dots, V - 1\}$ where V is
102 the total number of tokens. Here we represent the tokenizer as ι where $\iota(z) = t$. In this paper, we
103 assume that ι is the standard tokenizer where each binary vector is mapped to its decimal number. In
104 other words, there’s a one to one map between latent vectors and tokens. Because the map is one to
105 one, we sometimes use latent vectors and tokens interchangeably. We also assume that every latent
106 binary vector has a unique corresponding token, therefore $V = 2^m$.

107 Under the latent concept association model, the goal is to retrieve specific output tokens given partial
108 information in the contexts. This is modeled by the latent conditional distribution:

$$p(z|z^*) = \omega\pi(z|z^*) + (1 - \omega)\text{Unif}(\mathcal{Z})$$

109 where

$$\pi(z|z^*) \propto \begin{cases} \exp(-D_H(z, z^*)/\beta) & z \in \mathcal{N}(z^*), \\ 0 & z \notin \mathcal{N}(z^*). \end{cases}$$

110 Here D_H is the Hamming distance, $\mathcal{N}(z^*)$ is a subset of $\mathcal{Z} \setminus \{z^*\}$ and $\beta > 0$ is the temperature param-
111 eter. The use of Hamming distance draws a parallel with the notion of distributional semantics in natural
112 language: “a word is characterized by the company it keeps” [Fir57]. In words, $p(z|z^*)$ says that with
113 probability $1 - \omega$, the conditional distribution uniformly generate random latent vectors and with prob-
114 ability ω , the latent vector is generated from the *informative conditional distribution* $\pi(z|z^*)$ where
115 the support of the conditional distribution is $\mathcal{N}(z^*)$. Here, π represents the informative conditional dis-
116 tribution that depends on z^* whereas the uniform distribution is uninformative and can be considered
117 as noise. The mixture model parameter ω determines the signal to noise ratio of the contexts.

118 Therefore, for any latent vector z^* and its associated token, one can generate L context token words
119 with the aforementioned latent conditional distribution:

- 120 • Uniformly sample a latent vector z^*
- 121 • For $l = 1, \dots, L - 1$, sample $z_l \sim p(z|z^*)$ and $t_l = \iota(z_l)$.
- 122 • For $l = L$, sample $z \sim \pi(z|z^*)$ and $t_L = \iota(z)$.

123 Consequently, we have $x = (t_1, \dots, t_L)$ and $y = \iota(z^*)$. The last token in the context is generated
 124 specifically to make sure that it is not from the uniform distribution. This ensures that the last token
 125 can use attention to look for clues, relevant to the output, in the context. Let \mathcal{D}^L be the sampling
 126 distribution to generate (x, y) pairs. The conditional probability of y given x is given by $p(y|x)$.
 127 With slight abuse of notation, given a token $t \in [V]$, we define $\mathcal{N}(t) = \mathcal{N}(\iota^{-1}(t))$. We also define
 128 $D_H(t, t') = D_H(\iota^{-1}(t), \iota^{-1}(t'))$ for any pair of tokens t and t' .

129 For any function f that maps the context to estimated logits of output labels, the training objective
 130 is to minimize this loss of the last position: $\mathbb{E}_{(x,y) \in \mathcal{D}^L} [\ell(f(x), y)]$ where ℓ is the cross entropy loss
 131 with softmax. The error rate of latent concept association is defined by the following: $R_{\mathcal{D}^L}(f) =$
 132 $\mathbb{P}_{(x,y) \sim \mathcal{D}^L} [\text{argmax } f(x) \neq y]$ And the accuracy is $1 - R_{\mathcal{D}^L}(f)$.

133 3.2 Transformer network architecture

134 Given a context $x = (t_1, \dots, t_L)$ which consists of L tokens, we define $X \in \{0, 1\}^{V \times L}$ to be its
 135 one-hot encoding where V is the vocabulary size. Here we use χ to represent the one-hot encoding
 136 function (i.e., $\chi(x) = X$). Similar to [LLR23; Tar+23a; Li+24], we also consider a simplified
 137 one-layer transformer model without residual connections and normalization:

$$f^L(x) = \left[W_E^T W_V \text{attn}(W_E \chi(x)) \right]_{:L} \quad (3.1)$$

138 where

$$\text{attn}(U) = U \sigma \left(\frac{(W_K U)^T (W_Q U)}{\sqrt{d_a}} \right),$$

139 $W_K \in \mathbb{R}^{d_a \times d}$ is the key matrix, and $W_Q \in \mathbb{R}^{d_a \times d}$ is the query matrix and d_a is the attention head
 140 size. $\sigma : \mathbb{R}^{L \times L} \rightarrow (0, 1)^{L \times L}$ is the column-wise softmax operation. $W_V \in \mathbb{R}^{d \times d}$ is the value
 141 matrix and $W_E \in \mathbb{R}^{d \times V}$ is the embedding matrix. Here, we adopt the weight tie-in implementation
 142 which is used for Gemma [Tea+24]. We focus solely on the prediction of the last position, as it is
 143 the only one relevant for latent concept association. For convenience, we also use $h(x)$ to mean
 144 $[\text{attn}(W_E \chi(x))]_{:L}$, which is the hidden representation after attention for the last position, and $f_t^L(x)$
 145 to represent the logit for output token t .

146 4 Theoretical analysis

147 In this section, we theoretically investigate how a single-layer transformer can solve the latent
 148 concept association problem. We first introduce a hypothetical associative memory model that utilizes
 149 self-attention for information aggregation and employs the value matrix for memory retrieval. This
 150 hypothetical model turns out to mirror trained transformers in experiments. We also examine the
 151 role of each individual component of the network: the value matrix, embeddings, and the attention
 152 mechanism. We validate our theoretical claims in Appendix D.

153 4.1 Hypothetical associative memory model

154 In this section, we show that a simple single-layer transformer network can solve the latent concept
 155 association problem. The formal result is presented below in Theorem 1; first we require a few more
 156 definitions. Let $W_E(t)$ be the t -th column of the embedding matrix W_E . In other words, this is the
 157 embedding for token t . Given a token t , define $\mathcal{N}_1(t)$ to be the subset of tokens whose latent vectors
 158 are only 1 Hamming distance away from t 's latent vector: $\mathcal{N}_1(t) = \{t' : D_H(t', t) = 1\} \cap \mathcal{N}(t)$.
 159 For any output token t , $\mathcal{N}_1(t)$ contains tokens with the highest probabilities to appear in the context.

160 The following theorem formalizes the intuition that a one-layer transformer that uses self-attention
 161 to summarize statistics about the context distributions and whose value matrix uses aggregated
 162 representations to retrieve output tokens can solve the latent concept association problem defined in
 163 Section 3.1.

164 **Theorem 1** (informal). *Suppose the data generating process follows Section 3.1 where $m \geq 3$,
 165 $\omega = 1$, and $\mathcal{N}(t) = V \setminus \{t\}$. Then for any $\varepsilon > 0$, there exists a transformer model given by (3.1)
 166 that achieves error ε , i.e. $R_{\mathcal{D}^L}(f^L) < \varepsilon$ given sufficiently large context length L .*

167 More precisely, for the transformer in Theorem 1, we will have $W_K = 0$ and $W_Q = 0$. Each row of
 168 W_E is orthogonal to each other and normalized. And W_V is given by

$$W_V = \sum_{t \in [V]} W_E(t) \left(\sum_{t' \in \mathcal{N}_1(t)} W_E(t')^T \right) \quad (4.1)$$

169 A more formal statement of the theorem and its proof is given in Appendix E (Theorem 7).

170 Intuitively, Theorem 1 suggests having more samples from $p(x|y)$ can lead to a better recall rate. On
 171 the other hand, if contexts are modified to contain more samples from $p(x|\tilde{y})$ where $\tilde{y} \neq y$, then it is
 172 likely for transformer to output the wrong token. This is similar to context hijacking (see Section 4.4).
 173 The construction of the value matrix is similar to the associative memory model used in [Bie+24;
 174 CSB24], but in our case, there is no explicit one-to-one input and output pairs stored as memories.
 175 Rather, a combination of inputs are mapped to a single output.

176 While the construction in Theorem 1 is just one way that a single-layer transformer can tackle this task,
 177 it turns out empirically this construction of W_V is close to the trained W_V , even in the noisy case ($\omega \neq$
 178 1). In Appendix D.1, we will demonstrate that substituting trained value matrices with constructed
 179 ones can retain accuracy, and the constructed and trained value matrices even share close low-rank
 180 approximations. Moreover, in this hypothetical model, a simple uniform attention mechanism is
 181 deployed to allow self-attention to count occurrences of each individual tokens. Since the embeddings
 182 are orthonormal vectors, there is no interference. Hence, the self-attention layer can be viewed as
 183 aggregating information of contexts. It is worth noting that, in different settings, more sophisticated
 184 embedding structures and attention patterns are needed. This is discussed in the following sections.

185 4.2 On the role of the value matrix

186 The construction in Theorem 1 relies on the value matrix acting as associative memory. But is it
 187 necessary? Could we integrate the functionality of the value matrix into the self-attention module to
 188 solve the latent concept association problem? Empirically, the answer seems to be negative as will be
 189 shown in Appendix D.1. In particular, when the context length is small, setting the value matrix to be
 190 the identity would lead to subpar memory recall accuracy.

191 This is because if the value matrix is the identity, the transformer would be more susceptible to the
 192 noise in the context. To see this, notice that given any pair of context and output token (x, y) , the
 193 latent representation after self-attention $h(x)$ must live in the polyhedron S_y to be classified correctly
 194 where S_y is defined as:

$$S_y = \{v : (W_E(y) - W_E(t))^T v > 0 \text{ where } t \notin [V] \setminus \{y\}\}$$

195 Note that, by definition, for any two tokens y and \tilde{y} , $S_y \cap S_{\tilde{y}} = \emptyset$. On the other hand, because of the
 196 self-attention mechanism, $h(x)$ must also live in the convex hull of all the embedding vectors:

$$CV = \text{Conv}(W^E(0), \dots, W^E(|V| - 1))$$

197 In other words, for any pair (x, y) to be classified correctly, $h(x)$ must live in the intersection of S_y
 198 and CV . Due to the stochastic nature of x , it is likely for $h(x)$ to be outside of this intersection. The
 199 remapping effect of the value matrix can help with this problem. The following lemma explains this
 200 intuition.

201 **Lemma 2.** *Suppose the data generating process follows Section 3.1 where $m \geq 3$, $\omega = 1$ and
 202 $\mathcal{N}(t) = \{t' : D_H(t, t') = 1\}$. For any single layer transformer given by (3.1) where each row of
 203 W_E is orthogonal to each other and normalized, if W_V is constructed as in (4.1), then the error rate
 204 is 0. If W_V is the identity matrix, then the error rate is strictly larger than 0.*

205 Another intriguing phenomenon occurs when the value matrix is the identity matrix. In this case, the
 206 inner product between embeddings and their corresponding Hamming distance varies linearly. This
 207 relationship can be formalized by the following theorem.

208 **Theorem 3.** *Suppose the data generating process follows Section 3.1 where $m \geq 3$, $\omega = 1$ and*
 209 *$\mathcal{N}(t) = V \setminus \{t\}$. For any single layer transformer given by (3.1) with W_V being the identity matrix,*
 210 *if the cross entropy loss is minimized so that for any sampled pair (x, y) ,*

$$p(y|x) = \hat{p}(y|x) = \text{softmax}(f_y^L(x))$$

211 *there exists $a > 0$ and b such that for two tokens $t \neq t'$,*

$$\langle W_E(t), W_E(t') \rangle = -aD_H(t, t') + b$$

212 4.3 Embedding training and geometry

213 The hypothetical model in Section 4.1 requires embeddings to form an orthonormal basis. In
 214 the overparameterization regime where the embedding dimension d is larger than the number of
 215 tokens V , this can be approximately achieved by Gaussian initialization. However, in practice, the
 216 embedding dimension is typically smaller than the vocabulary size, in which case it is impossible
 217 for the embeddings to constitute such a basis. Empirically, in Appendix D.2, we observe that with
 218 overparameterization ($d > V$), embeddings can be frozen at their Gaussian initialization, whereas in
 219 the underparameterized regime, embedding training is required to achieve better recall accuracy.

220 This raises the question: What kind of embedding geometry is learned in the underparameterized
 221 regime? Experiments reveal a close relationship between the inner product of embeddings for two
 222 tokens and the Hamming distance of these tokens (see Figure 3b and Figure G.5 in Appendix G.2).
 223 Approximately, we have the following relationship:

$$\langle W_E(t), W_E(t') \rangle = \begin{cases} b_0 & t = t' \\ -aD_H(t, t') + b & t \neq t' \end{cases} \quad (4.2)$$

224 for any two tokens t and t' where $b_0 > b$ and $a > 0$. One can view this as a combination of the
 225 embedding geometry under Gaussian initialization and the geometry when W_V is the identity matrix
 226 (Theorem 3). Importantly, this structure demonstrates that trained embeddings inherently capture
 227 similarity within the latent space. Theoretically, this embedding structure (4.2) can also lead to low
 228 error rate under specific conditions on b_0 , b and a , which is articulated by the following theorem.

229 **Theorem 4 (Informal).** *Following the same setup as in Theorem 1, but embeddings obey (4.2), then*
 230 *under certain conditions on a , b and if b_0 and context length L are sufficiently large, the error rate*
 231 *can be arbitrarily small, i.e. $R_{\mathcal{D}^L}(f^L) < \varepsilon$ for any $0 < \varepsilon < 1$.*

232 The formal statement of the theorem and its proof is given in Appendix E (Theorem 8).

233 Notably, this embedding geometry also implies a low-rank structure. Let’s first consider the special
 234 case when $b_0 = b$. In other words, the inner product between embeddings and their corresponding
 235 Hamming distance varies linearly.

236 **Lemma 5.** *If embeddings follow (4.2) and $b = b_0$ and $\mathcal{N}(t) = V \setminus \{t\}$, then $\text{rank}(W_E) \leq m + 2$.*

237 When $b_0 > b$, the embedding matrix will not be strictly low rank. However, it can still exhibit
 238 approximate low-rank behavior, characterized by an eigengap between the top and bottom singular
 239 values. This is verified empirically (see Figure G.9-G.12 in Appendix G.4).

240 4.4 Context hijacking and the misclassification of memory recall

241 In light of the theoretical results on latent concept association, a natural question arises: How do these
 242 results connect to context hijacking in LLMs? In essence, for the latent concept association problem,
 243 the differentiation of output tokens is achieved by distinguishing between the various conditional
 244 distributions $p(x|y)$. Thus, adding or changing tokens in the context x so that it resembles a different
 245 conditional distribution can result in misclassification. In Appendix G.5, we present experiments
 246 showing that mixing different contexts can cause transformers to misclassify. This partially explains
 247 context hijacking in LLMs (Appendix B). On the other hand, it is well-known that the error rate
 248 is related to the KL divergence between conditional distributions of contexts [Cov99]. The closer
 249 the distributions are, the easier it is for the model to misclassify. Here, longer contexts, primarily
 250 composed of i.i.d samples, suggest larger divergences, thus higher memory recall rate. This is
 251 theoretically implied by Theorem 1 and Theorem 4 and empirically verified in Appendix G.6. Such
 252 result is also related to reverse context hijacking (Appendix F) where prepending sentences including
 253 true target words can improve fact recall rate.

254 References

- 255 [AZL23a] Z. Allen-Zhu and Y. Li. *Physics of language models: part 3.1, knowledge storage and*
256 *extraction*. 2023. arXiv: [2309.14316](#) [cs.CL] (cit. on p. 12).
- 257 [AZL23b] Z. Allen-Zhu and Y. Li. *Physics of language models: part 3.2, knowledge manipulation*.
258 2023. arXiv: [2309.14402](#) [cs.CL] (cit. on p. 12).
- 259 [AZL24] Z. Allen-Zhu and Y. Li. *Physics of language models: part 3.3, knowledge capacity*
260 *scaling laws*. 2024. arXiv: [2404.05405](#) [cs.CL] (cit. on p. 12).
- 261 [Apr+22] G. Apruzzese, H. S. Anderson, S. Dambra, D. Freeman, F. Pierazzi, and K. A. Roundy.
262 "real attackers don't compute gradients": bridging the gap between adversarial ml
263 *research and practice*. 2022. arXiv: [2212.14315](#) [cs.CR] (cit. on p. 12).
- 264 [Bai+24] Y. Bai, F. Chen, H. Wang, C. Xiong, and S. Mei. "Transformers as statisticians:
265 provable in-context learning with in-context algorithm selection". *Advances in neural*
266 *information processing systems* (2024) (cit. on p. 12).
- 267 [BYBOS95] R. Ben-Yishai, R. L. Bar-Or, and H. Sompolinsky. "Theory of orientation tuning in
268 visual cortex." *Proceedings of the National Academy of Sciences* 9 (1995) (cit. on
269 p. 12).
- 270 [Bie+24] A. Bietti, V. Cabannes, D. Bouchacourt, H. Jegou, and L. Bottou. "Birth of a trans-
271 former: a memory viewpoint". *Advances in Neural Information Processing Systems*
272 (2024) (cit. on pp. 5, 12).
- 273 [BP21] T. Bricken and C. Pehlevan. "Attention approximates sparse distributed memory".
274 *Advances in Neural Information Processing Systems* (2021) (cit. on pp. 3, 12).
- 275 [Bro+20] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan,
276 P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan,
277 R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler,
278 M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I.
279 Sutskever, and D. Amodei. *Language models are few-shot learners*. 2020. arXiv:
280 [2005.14165](#) [cs.CL] (cit. on p. 1).
- 281 [CDB23] V. Cabannes, E. Dohmatob, and A. Bietti. "Scaling laws for associative memories".
282 *arXiv preprint arXiv:2310.02984* (2023) (cit. on p. 12).
- 283 [CSB24] V. Cabannes, B. Simsek, and A. Bietti. "Learning associative memories with gradient
284 descent". *arXiv preprint arXiv:2402.18724* (2024) (cit. on pp. 5, 12).
- 285 [Cha22] F. Charton. "What is my math transformer doing?—three results on interpretability and
286 generalization". *arXiv preprint arXiv:2211.00170* (2022) (cit. on p. 12).
- 287 [Cho+24] A. G. Chowdhury, M. M. Islam, V. Kumar, F. H. Shezan, V. Jain, and A. Chadha.
288 "Breaking down the defenses: a comparative survey of attacks on large language
289 models". *arXiv preprint arXiv:2403.04786* (2024) (cit. on p. 12).
- 290 [Cov99] T. M. Cover. *Elements of information theory*. 1999 (cit. on p. 6).
- 291 [CSH22] A. Creswell, M. Shanahan, and I. Higgins. "Selection-inference: exploiting large
292 language models for interpretable logical reasoning". *arXiv preprint arXiv:2205.09712*
293 (2022) (cit. on pp. 2, 12, 13).
- 294 [Dai+21] D. Dai, L. Dong, Y. Hao, Z. Sui, B. Chang, and F. Wei. "Knowledge neurons in
295 pretrained transformers". *arXiv preprint arXiv:2104.08696* (2021) (cit. on pp. 1, 12).
- 296 [DCAT21] N. De Cao, W. Aziz, and I. Titov. "Editing factual knowledge in language models".
297 *arXiv preprint arXiv:2104.08164* (2021) (cit. on pp. 1, 12).
- 298 [Dev83] L. Devroye. "The equivalence of weak, strong and complete convergence in l1 for
299 kernel density estimates". *The Annals of Statistics* 3 (1983) (cit. on p. 16).
- 300 [Ede+24] B. L. Edelman, E. Edelman, S. Goel, E. Malach, and N. Tsilivis. "The evolution
301 of statistical induction heads: in-context learning markov chains". *arXiv preprint*
302 *arXiv:2402.11004* (2024) (cit. on p. 12).
- 303 [Elh+21] N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai,
304 A. Chen, T. Conerly, et al. "A mathematical framework for transformer circuits".
305 *Transformer Circuits Thread* (2021) (cit. on p. 12).
- 306 [EI+24] M. Emrullah Ildiz, Y. Huang, Y. Li, A. Singh Rawat, and S. Oymak. "From self-
307 attention to markov models: unveiling the dynamics of generative transformers". *arXiv*
308 *e-prints* (2024) (cit. on p. 12).

- 309 [Fel20] V. Feldman. “Does learning require memorization? a short tale about a long tail”. In:
310 *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*.
311 2020 (cit. on p. 12).
- 312 [FZ20] V. Feldman and C. Zhang. “What neural networks memorize and why: discovering
313 the long tail via influence estimation”. *Advances in Neural Information Processing*
314 *Systems* (2020) (cit. on p. 12).
- 315 [Fir57] J. Firth. “A synopsis of linguistic theory, 1930-1955”. *Studies in linguistic analysis*
316 (1957) (cit. on p. 3).
- 317 [Gar+22] S. Garg, D. Tsipras, P. S. Liang, and G. Valiant. “What can transformers learn in-
318 context? a case study of simple function classes”. *Advances in Neural Information*
319 *Processing Systems* (2022) (cit. on p. 12).
- 320 [Gei+21] A. Geiger, H. Lu, T. Icard, and C. Potts. “Causal abstractions of neural networks”.
321 *Advances in Neural Information Processing Systems* (2021) (cit. on p. 12).
- 322 [Gei+22] A. Geiger, Z. Wu, H. Lu, J. Rozner, E. Kreiss, T. Icard, N. Goodman, and C. Potts.
323 “Inducing causal structure for interpretable neural networks”. In: *International Con-*
324 *ference on Machine Learning*. PMLR. 2022 (cit. on p. 12).
- 325 [Gei+24] A. Geiger, Z. Wu, C. Potts, T. Icard, and N. Goodman. “Finding alignments between
326 interpretable causal variables and distributed neural representations”. In: *Causal*
327 *Learning and Reasoning*. PMLR. 2024 (cit. on p. 12).
- 328 [Gev+23] M. Geva, J. Bastings, K. Filippova, and A. Globerson. “Dissecting recall of factual
329 associations in auto-regressive language models”. *arXiv preprint arXiv:2304.14767*
330 (2023) (cit. on p. 12).
- 331 [Has+24] P. Hase, M. Bansal, B. Kim, and A. Ghandeharioun. “Does localization inform edit-
332 ing? surprising differences in causality-based localization vs. knowledge editing in
333 language models”. *Advances in Neural Information Processing Systems* (2024) (cit. on
334 p. 12).
- 335 [Hen+23] T. Henighan, S. Carter, T. Hume, N. Elhage, R. Lasenby, S. Fort, N. Schiefer, and
336 C. Olah. “Superposition, memorization, and double descent”. *Transformer Circuits*
337 *Thread* (2023) (cit. on p. 12).
- 338 [Hop82] J. J. Hopfield. “Neural networks and physical systems with emergent collective com-
339 putational abilities.” *Proceedings of the national academy of sciences* 8 (1982) (cit. on
340 p. 12).
- 341 [Hu+21] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. *Lora:*
342 *low-rank adaptation of large language models*. 2021. arXiv: 2106.09685 [cs.GL]
343 (cit. on p. 2).
- 344 [JP20] Y. Jiang and C. Pehlevan. “Associative memory in iterated overparameterized sigmoid
345 autoencoders”. In: *International conference on machine learning*. PMLR. 2020 (cit. on
346 p. 12).
- 347 [Jia+24] Y. Jiang, G. Rajendran, P. Ravikumar, B. Aragam, and V. Veitch. “On the origins of
348 linear representations in large language models”. *arXiv preprint arXiv:2403.03867*
349 (2024) (cit. on p. 12).
- 350 [Jin+23] T. Jin, N. Clement, X. Dong, V. Nagarajan, M. Carbin, J. Ragan-Kelley, and G. K.
351 Dziugaite. “The cost of down-scaling language models: fact recall deteriorates before
352 in-context learning”. *arXiv preprint arXiv:2310.04680* (2023) (cit. on p. 12).
- 353 [KKM22] J. Kim, M. Kim, and B. Mozafari. “Provable memorization capacity of transformers”.
354 In: *The Eleventh International Conference on Learning Representations*. 2022 (cit. on
355 p. 12).
- 356 [Li+24] Y. Li, Y. Huang, M. E. Ildiz, A. S. Rawat, and S. Oymak. “Mechanics of next token
357 prediction with self-attention”. In: *International Conference on Artificial Intelligence*
358 *and Statistics*. PMLR. 2024 (cit. on pp. 4, 12).
- 359 [LLR23] Y. Li, Y. Li, and A. Risteski. “How do transformers learn topic structure: towards
360 a mechanistic understanding”. In: *International Conference on Machine Learning*.
361 PMLR. 2023 (cit. on pp. 4, 12).
- 362 [Liu+23a] Y. Liu, G. Deng, Y. Li, K. Wang, T. Zhang, Y. Liu, H. Wang, Y. Zheng, and Y.
363 Liu. “Prompt injection attack against llm-integrated applications”. *arXiv preprint*
364 *arXiv:2306.05499* (2023) (cit. on p. 12).

- 365 [Liu+23b] Y. Liu, G. Deng, Z. Xu, Y. Li, Y. Zheng, Y. Zhang, L. Zhao, T. Zhang, and Y. Liu.
366 “Jailbreaking chatgpt via prompt engineering: an empirical study”. *arXiv preprint*
367 *arXiv:2305.13860* (2023) (cit. on p. 12).
- 368 [Liu+22] Z. Liu, O. Kitouni, N. S. Nolte, E. Michaud, M. Tegmark, and M. Williams. “Towards
369 understanding grokking: an effective theory of representation learning”. *Advances in*
370 *Neural Information Processing Systems* (2022) (cit. on p. 12).
- 371 [LH17] I. Loshchilov and F. Hutter. “Decoupled weight decay regularization”. *arXiv preprint*
372 *arXiv:1711.05101* (2017) (cit. on p. 26).
- 373 [MLT23] S. Mahdavi, R. Liao, and C. Thrampoulidis. “Memorization capacity of multi-head
374 attention in transformers”. *arXiv preprint arXiv:2306.02010* (2023) (cit. on p. 12).
- 375 [Mak+24] A. V. Makkuva, M. Bondaschi, A. Girish, A. Nagle, M. Jaggi, H. Kim, and M. Gastpar.
376 *Attention with markov: a framework for principled analysis of transformers via markov*
377 *chains*. 2024. arXiv: 2402.04161 [cs.LG] (cit. on p. 12).
- 378 [McG+23] T. McGrath, M. Rahtz, J. Kramar, V. Mikulik, and S. Legg. “The hydra effect: emergent
379 self-repair in language model computations”. *arXiv preprint arXiv:2307.15771* (2023)
380 (cit. on p. 12).
- 381 [Men+22] K. Meng, D. Bau, A. Andonian, and Y. Belinkov. “Locating and editing factual
382 associations in gpt”. *Advances in Neural Information Processing Systems* (2022) (cit.
383 on pp. 1, 2, 12, 13, 24).
- 384 [Men+23] K. Meng, A. S. Sharma, A. Andonian, Y. Belinkov, and D. Bau. *Mass-editing memory*
385 *in a transformer*. 2023. arXiv: 2210.07229 [cs.CL] (cit. on pp. 1, 12).
- 386 [Mil+22] B. Millidge, T. Salvatori, Y. Song, T. Lukasiewicz, and R. Bogacz. “Universal hop-
387 field networks: a general framework for single-shot associative memory models”. In:
388 *International Conference on Machine Learning*. PMLR. 2022 (cit. on pp. 3, 12).
- 389 [Mit+21] E. Mitchell, C. Lin, A. Bosselut, C. Finn, and C. D. Manning. “Fast model editing at
390 scale”. *arXiv preprint arXiv:2110.11309* (2021) (cit. on pp. 1, 12).
- 391 [Mit+22] E. Mitchell, C. Lin, A. Bosselut, C. D. Manning, and C. Finn. “Memory-based model
392 editing at scale”. In: *International Conference on Machine Learning*. PMLR. 2022
393 (cit. on pp. 1, 12).
- 394 [Nan+23] N. Nanda, L. Chan, T. Lieberum, J. Smith, and J. Steinhardt. “Progress measures for
395 grokking via mechanistic interpretability”. *arXiv preprint arXiv:2301.05217* (2023)
396 (cit. on p. 12).
- 397 [NDL24] E. Nichani, A. Damian, and J. D. Lee. *How transformers learn causal structure with*
398 *gradient descent*. 2024. arXiv: 2402.14735 [cs.LG] (cit. on p. 12).
- 399 [Ols+22a] C. Olsson, N. Elhage, N. Nanda, N. Joseph, N. DasSarma, T. Henighan, B. Mann,
400 A. Askell, Y. Bai, A. Chen, T. Conerly, D. Drain, D. Ganguli, Z. Hatfield-Dodds,
401 D. Hernandez, S. Johnston, A. Jones, J. Kernion, L. Lovitt, K. Ndousse, D. Amodei,
402 T. Brown, J. Clark, J. Kaplan, S. McCandlish, and C. Olah. *In-context learning and*
403 *induction heads*. 2022. arXiv: 2209.11895 [cs.LG] (cit. on p. 12).
- 404 [Ols+22b] C. Olsson, N. Elhage, N. Nanda, N. Joseph, N. DasSarma, T. Henighan, B. Mann,
405 A. Askell, Y. Bai, A. Chen, et al. “In-context learning and induction heads”. *arXiv*
406 *preprint arXiv:2209.11895* (2022) (cit. on p. 12).
- 407 [PE21] L. Pandia and A. Ettinger. “Sorting through the noise: testing robustness of information
408 processing in pre-trained language models”. *arXiv preprint arXiv:2109.12393* (2021)
409 (cit. on pp. 2, 12, 13).
- 410 [PR22] F. Perez and I. Ribeiro. “Ignore previous prompt: attack techniques for language
411 models”. *arXiv preprint arXiv:2211.09527* (2022) (cit. on p. 12).
- 412 [Pet+20] F. Petroni, P. Lewis, A. Piktus, T. Rocktäschel, Y. Wu, A. H. Miller, and S.
413 Riedel. “How context affects language models’ factual predictions”. *arXiv preprint*
414 *arXiv:2005.04611* (2020) (cit. on pp. 2, 12, 13).
- 415 [Rad+19] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. “Language
416 models are unsupervised multitask learners”. *OpenAI blog* 8 (2019) (cit. on pp. 1, 2,
417 12).
- 418 [RBU20] A. Radhakrishnan, M. Belkin, and C. Uhler. “Overparameterized neural networks
419 implement associative memory”. *Proceedings of the National Academy of Sciences*
420 44 (2020) (cit. on p. 12).

- 421 [Raj+24] G. Rajendran, S. Buchholz, B. Aragam, B. Schölkopf, and P. Ravikumar. “Learning
422 interpretable concepts: unifying causal representation learning and foundation models”.
423 *arXiv preprint arXiv:2402.09236* (2024) (cit. on p. 12).
- 424 [Ram+20] H. Ramsauer, B. Schöfl, J. Lehner, P. Seidl, M. Widrich, T. Adler, L. Gruber, M.
425 Holzleitner, M. Pavlović, G. K. Sandve, et al. “Hopfield networks is all you need”.
426 *arXiv preprint arXiv:2008.02217* (2020) (cit. on pp. 3, 12).
- 427 [Rao+23] A. Rao, S. Vashistha, A. Naik, S. Aditya, and M. Choudhury. “Tricking llms into
428 disobedience: understanding, analyzing, and preventing jailbreaks”. *arXiv preprint*
429 *arXiv:2305.14965* (2023) (cit. on p. 12).
- 430 [Sak+23] M. Sakarvadia, A. Ajith, A. Khan, D. Grzenda, N. Hudson, A. Bauer, K. Chard,
431 and I. Foster. “Memory injections: correcting multi-hop reasoning failures during
432 inference in transformer-based language models”. *arXiv preprint arXiv:2309.05605*
433 (2023) (cit. on p. 12).
- 434 [Seu96] H. S. Seung. “How the brain keeps the eyes still”. *Proceedings of the National*
435 *Academy of Sciences* 23 (1996) (cit. on p. 12).
- 436 [SMR23] M. Shanahan, K. McDonell, and L. Reynolds. “Role play with large language models”.
437 *Nature* 7987 (2023) (cit. on p. 12).
- 438 [She+24] H. Sheen, S. Chen, T. Wang, and H. H. Zhou. “Implicit regularization of gradient
439 flow on one-layer softmax attention”. *arXiv preprint arXiv:2403.08699* (2024) (cit. on
440 p. 12).
- 441 [Shi+23] F. Shi, X. Chen, K. Misra, N. Scales, D. Dohan, E. H. Chi, N. Schärli, and D. Zhou.
442 “Large language models can be easily distracted by irrelevant context”. In: *International*
443 *Conference on Machine Learning*. PMLR. 2023 (cit. on pp. 2, 12, 13).
- 444 [Si+22] W. M. Si, M. Backes, J. Blackburn, E. De Cristofaro, G. Stringhini, S. Zannettou, and
445 Y. Zhang. “Why so toxic? measuring and triggering toxic behavior in open-domain
446 chatbots”. In: *Proceedings of the 2022 ACM SIGSAC Conference on Computer and*
447 *Communications Security*. 2022 (cit. on p. 12).
- 448 [Ska+94] W. Skaggs, J. Knierim, H. Kudrimoti, and B. McNaughton. “A model of the neural
449 basis of the rat’s sense of direction”. *Advances in neural information processing*
450 *systems* (1994) (cit. on p. 12).
- 451 [SS22] J. Steinberg and H. Sompolinsky. “Associative memory of structured knowledge”.
452 *Scientific Reports* 1 (2022) (cit. on p. 12).
- 453 [Tar+23a] D. A. Tarzanagh, Y. Li, C. Thrampoulidis, and S. Oymak. “Transformers as support
454 vector machines”. *arXiv preprint arXiv:2308.16898* (2023) (cit. on pp. 4, 12).
- 455 [Tar+23b] D. A. Tarzanagh, Y. Li, X. Zhang, and S. Oymak. “Margin maximization in attention
456 mechanism”. *arXiv preprint arXiv:2306.13596* (2023) (cit. on p. 12).
- 457 [Tea+24] G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M.
458 Rivière, M. S. Kale, J. Love, et al. “Gemma: open models based on gemini research
459 and technology”. *arXiv preprint arXiv:2403.08295* (2024) (cit. on pp. 2, 4, 12).
- 460 [Tia+24] B. Tian, S. Cheng, X. Liang, N. Zhang, Y. Hu, K. Xue, Y. Gou, X. Chen, and H. Chen.
461 “Instructedit: instruction-based knowledge editing for large language models”. *arXiv*
462 *preprint arXiv:2402.16123* (2024) (cit. on p. 12).
- 463 [Tia+23a] Y. Tian, Y. Wang, B. Chen, and S. S. Du. “Scan and snap: understanding training
464 dynamics and token composition in 1-layer transformer”. *Advances in Neural Inform-*
465 *ation Processing Systems* (2023) (cit. on p. 12).
- 466 [Tia+23b] Y. Tian, Y. Wang, Z. Zhang, B. Chen, and S. Du. “Joma: demystifying multilayer trans-
467 formers via joint dynamics of mlp and attention”. *arXiv preprint arXiv:2310.00535*
468 (2023) (cit. on p. 12).
- 469 [Tou+23] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov,
470 S. Batra, P. Bhargava, S. Bhosale, et al. “Llama 2: open foundation and fine-tuned
471 chat models”. *arXiv preprint arXiv:2307.09288* (2023) (cit. on pp. 2, 12).
- 472 [Vas+17] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
473 I. Polosukhin. “Attention is all you need”. *Advances in neural information processing*
474 *systems* (2017) (cit. on pp. 2, 12).

- 475 [Wan+23a] B. Wang, W. Chen, H. Pei, C. Xie, M. Kang, C. Zhang, C. Xu, Z. Xiong, R. Dutta,
476 R. Schaeffer, et al. “Decodingtrust: a comprehensive assessment of trustworthiness in
477 gpt models”. *arXiv preprint arXiv:2306.11698* (2023) (cit. on p. 12).
- 478 [Wan+23b] J. Wang, X. Hu, W. Hou, H. Chen, R. Zheng, Y. Wang, L. Yang, H. Huang, W. Ye,
479 X. Geng, et al. “On the robustness of chatgpt: an adversarial and out-of-distribution
480 perspective”. *arXiv preprint arXiv:2302.12095* (2023) (cit. on p. 12).
- 481 [Wan+22] K. Wang, A. Variengien, A. Conmy, B. Shlegeris, and J. Steinhardt. “Interpretability
482 in the wild: a circuit for indirect object identification in gpt-2 small”. *arXiv preprint*
483 *arXiv:2211.00593* (2022) (cit. on p. 12).
- 484 [Wu+24] Z. Wu, A. Geiger, T. Icard, C. Potts, and N. Goodman. “Interpretability at scale:
485 identifying causal mechanisms in alpaca”. *Advances in Neural Information Processing*
486 *Systems* (2024) (cit. on p. 12).
- 487 [Xie+21] S. M. Xie, A. Raghunathan, P. Liang, and T. Ma. “An explanation of in-context
488 learning as implicit bayesian inference”. *arXiv preprint arXiv:2111.02080* (2021)
489 (cit. on p. 12).
- 490 [Xu+23] X. Xu, K. Kong, N. Liu, L. Cui, D. Wang, J. Zhang, and M. Kankanhalli. “An llm
491 can fool itself: a prompt-based adversarial attack”. *arXiv preprint arXiv:2310.13345*
492 (2023) (cit. on p. 12).
- 493 [Yor+23] O. Yoran, T. Wolfson, O. Ram, and J. Berant. “Making retrieval-augmented language
494 models robust to irrelevant context”. *arXiv preprint arXiv:2310.01558* (2023) (cit. on
495 pp. 2, 12, 13).
- 496 [Zha+22] Y. Zhang, A. Backurs, S. Bubeck, R. Eldan, S. Gunasekar, and T. Wagner. “Unveiling
497 transformers with lego: a synthetic reasoning task”. *arXiv preprint arXiv:2206.04301*
498 (2022) (cit. on p. 12).
- 499 [Zha+23] Z. Zhang, M. Fang, L. Chen, M.-R. Namazi-Rad, and J. Wang. “How do large language
500 models capture the ever-changing world knowledge? a review of recent advances”.
501 *arXiv preprint arXiv:2310.07343* (2023) (cit. on p. 12).
- 502 [Zha23] J. Zhao. “In-context exemplars as clues to retrieving from large associative memory”.
503 *arXiv preprint arXiv:2311.03498* (2023) (cit. on pp. 2, 3, 12, 13).
- 504 [Zhe+23] C. Zheng, L. Li, Q. Dong, Y. Fan, Z. Wu, J. Xu, and B. Chang. “Can we edit factual
505 knowledge by in-context learning?” *arXiv preprint arXiv:2305.12740* (2023) (cit. on
506 p. 12).
- 507 [Zho+24] Z. Zhong, Z. Liu, M. Tegmark, and J. Andreas. “The clock and the pizza: two stories
508 in mechanistic explanation of neural networks”. *Advances in Neural Information*
509 *Processing Systems* (2024) (cit. on p. 12).
- 510 [Zhu+23] K. Zhu, J. Wang, J. Zhou, Z. Wang, H. Chen, Y. Wang, L. Yang, W. Ye, N. Z. Gong,
511 Y. Zhang, et al. “Promptbench: towards evaluating the robustness of large language
512 models on adversarial prompts”. *arXiv preprint arXiv:2306.04528* (2023) (cit. on
513 p. 12).
- 514 [Zou+23] A. Zou, Z. Wang, N. Carlini, M. Nasr, J. Z. Kolter, and M. Fredrikson. *Universal*
515 *and transferable adversarial attacks on aligned language models*. 2023. arXiv: 2307.
516 15043 [cs.CL] (cit. on p. 12).

517 A Literature review

518 **Associative memory** Associative memory has been explored within the field of neuroscience
519 [Hop82; Seu96; BYBOS95; Ska+94; SS22]. The most popular models among them is the Hopfield
520 network [Hop82] and its modern successors [Ram+20; Mil+22; Zha23] are closely related to the
521 attention layer used in transformers [Vas+17]. In addition, the attention mechanism has also been
522 shown to approximate another associative memory model known as sparse distributed memory
523 [BP21]. Beyond attention, Radhakrishnan et al. [RBU20] and Jiang and Pehlevan [JP20] show
524 that overparameterized autoencoders can implement associative memory as well. This paper studies
525 fact retrieval as a form of associative memory. Another closely related area of research focuses on
526 memorization in deep neural networks. Henighan et al. [Hen+23] shows that a simple neural network
527 trained on toy model will store data points in the overfitting regime while storing features in the
528 underfitting regime. Feldman [Fel20] and Feldman and Zhang [FZ20] study the interplay between
529 memorization and long tail distributions while Kim et al. [KKM22] and Mahdavi et al. [MLT23]
530 study the memorization capacity of transformers.

531 **Interpreting transformers and LLMs** There’s a growing body of work on understanding how
532 transformers and LLMs work [LLR23; AZL23a; AZL23b; AZL24; EI+24; Tar+23b; Tar+23a; Li+24],
533 including training dynamics [Tia+23a; Tia+23b; She+24] and in-context learning [Xie+21; Gar+22;
534 Bai+24; Bai+24]. Recent papers have introduced synthetic tasks to better understand the mechanisms
535 of transformers [Cha22; Liu+22; Nan+23; Zha+22; Zho+24], such as those focused on Markov
536 chains [Bie+24; Ede+24; NDL24; Mak+24]. Most notably, Bietti et al. [Bie+24] and subsequent
537 works [CDB23; CSB24] study weights in transformers as associative memory but their focus is
538 on understanding induction head [Ols+22b] and one-to-one map between input query and output
539 memory. An increasing amount of research is dedicated to understanding the internals of pre-trained
540 LLMs, broadly categorized under the term “mechanistic interpretability” [Elh+21; Ols+22a; Gev+23;
541 Men+22; Men+23; Jia+24; Raj+24; Has+24; Wan+22; McG+23; Gei+21; Gei+22; Gei+24; Wu+24].

542 **Knowledge editing and adversarial attacks on LLMs** Fact recall and knowledge editing have
543 been extensively studied [Men+22; Men+23; Has+24; Sak+23; DCAT21; Mit+21; Mit+22; Dai+21;
544 Zha+23; Tia+24; Jin+23], including the use of in-context learning to edit facts [Zhe+23]. This
545 paper aims to explore a different aspect by examining the robustness of fact recall to variation in
546 prompts. A closely related line of work focuses on adversarial attacks on LLMs [see Cho+24, for a
547 review]. Specifically, prompt-based adversarial attacks [Xu+23; Zhu+23; Wan+23b] focus on the
548 manipulation of answers within specific classification tasks while other works concentrate on safety
549 issues [Liu+23a; PR22; Zou+23; Apr+22; Wan+23a; Si+22; Rao+23; SMR23; Liu+23b]. There
550 are also works showing LLMs can be distracted by irrelevant contexts in problem solving [Shi+23],
551 question answering [Pet+20; CSH22; Yor+23] and factual reasoning [PE21]. Although phenomena
552 akin to context hijacking have been reported in different instances, the goals of this paper are to give
553 a systematic robustness study for fact retrieval, offer a framework for interpreting it in the context of
554 associative memory, and deepen our understanding of LLMs.

555 B Context hijacking in LLMs

556 In this section, we run experiments on LLMs including GPT-2 [Rad+19], Gemma [Tea+24] (both
557 base and instruct models) and LLaMA-2-7B [Tou+23] to explore the effects of context hijacking
558 on manipulating LLM outputs. As an example, consider Figure 1. When we prompt the LLMs
559 with the context “The Eiffel Tower is in the city of”, all 4 LLMs output the correct answer (“Paris”).
560 However, as we see in the example, we can actually manipulate the output of the LLMs simply by
561 modifying the context with additional *factual* information that would not confuse a human. We call
562 this *context-hijacking*. Due to the different capacities and capabilities of each model, the examples in
563 Figure 1 use different hijacking techniques. This is most notable on LLaMA-2-7B, which is a much
564 larger model than the others. Of course, as expected, the more sophisticated attack on LLaMA also
565 works on GPT-2 and Gemma. Additionally, the instruction-tuned version of Gemma can understand
566 special words like “not” to some extent. Nevertheless, it is still possible to systematically hijack
567 these LLMs, as demonstrated below.

568 We explore this phenomenon at scale with the COUNTERFACT dataset introduced in [Men+22], a
569 dataset of difficult counterfactual assertions containing a diverse set of subjects, relations, and linguistics.

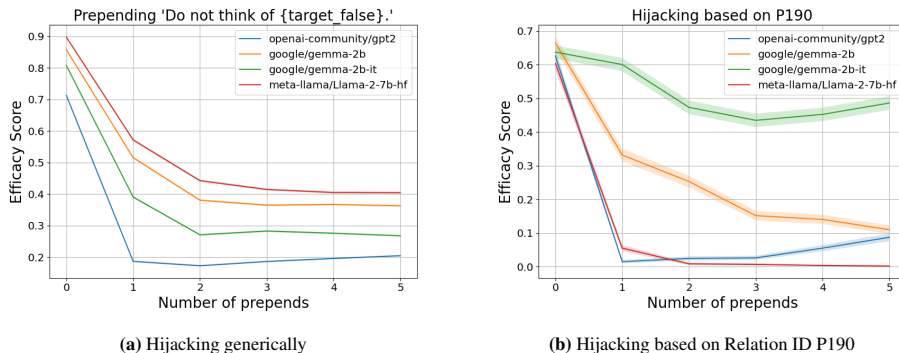


Figure 2: Context hijacking can cause LLMs to output false target. The figure shows efficacy score versus the number of preponds for various LLMs on the COUNTERFACT dataset under two hijacking schemes.

570 tic variations. COUNTERFACT has 21, 919 samples, each of which are given by a tuple (p, o_*, o_-, s, r) .
 571 From each sample, we have a context prompt p with a true target answer o_* (target_true) and a
 572 false target answer o_- (target_false), e.g. the prompt $p = \text{“Eiffel Tower can be found in”}$ has true
 573 target $o_* = \text{“Paris”}$ and false target $o_- = \text{“Guam”}$. Additionally, the main entity in p is the subject
 574 s ($s = \text{“Eiffel Tower”}$) and the prompt is categorized into relations r (for instance, other samples
 575 with the same relation ID as the example above could be of the form $\text{“The location of \{subject\} is”}$,
 576 $\text{“\{subject\} can be found in”}$, $\text{“Where is \{subject\}? It is in”}$). For additional details on how the dataset
 577 was collected, see [Men+22].

578 For a hijacking scheme, we report the Efficacy Score (ES) [Men+22], which is the proportion of
 579 samples for which the token probabilities satisfy $Pr[o_-] > Pr[o_*]$ after modifying the context,
 580 that is, the proportion of the dataset that has been successfully manipulated. We experiment with
 581 two hijacking schemes for this dataset. We first hijack by prepending the text $\text{“Do not think of$
 582 $\{target_false\}”}$ to each context. For instance, the prompt $\text{“The Eiffel Tower is in”}$ gets changed to
 583 $\text{“Do not think of Guam. The Eiffel Tower is in”}$. In Figure 2a, we see that the efficacy score drops
 584 significantly after hijacking. Here, we prepend the hijacking sentence k times for $k = 0, \dots, 5$ where
 585 $k = 0$ yields the original prompt. We see that additional preponds decrease the score further.

586 In the second scheme, we make use of the relation ID r to prepend factually correct sentences. For
 587 instance, one can hijack the example above to $\text{“The Eiffel Tower is not located in Guam. The Eiffel$
 588 Tower is in” . We test this hijacking philosophy on different relation IDs. In particular, Figure 2b
 589 reports hijacking based on relation ID $P190$ (“twin city”). And we see similar patterns that with
 590 more preponds, the ES score gets lower. It is also worth noting that one can even hijack by only
 591 including words that are semantically close to the false target (e.g., “France” for false target “French”).
 592 This suggests that context hijacking is more than simply the LLM copying tokens from contexts.
 593 Additional details and experiments for both hijacking schemes and for other relation IDs are in
 594 Appendix F.

595 These experiments show that context hijacking changes the behavior of LLMs, leading them to
 596 output incorrect tokens, without altering the factual meaning of the context. It is worth noting that
 597 similar fragile behaviors of LLMs have been observed in the literature in different contexts [Shi+23;
 598 Pet+20; CSH22; Yor+23; PE21]. See Appendix A for more details.

599 Context hijacking indicates that fact retrieval in LLMs is not robust and that accurate fact recall
 600 does not necessarily depend on the semantics of the context. As a result, one hypothesis is to view
 601 LLMs as an associative memory model where special tokens in contexts, associated with the fact,
 602 provide partial information or clues to facilitate memory retrieval [Zha23]. To better understand
 603 this perspective, we design a synthetic memory retrieval task to evaluate how the building blocks of
 604 LLMs, transformers, can solve it.

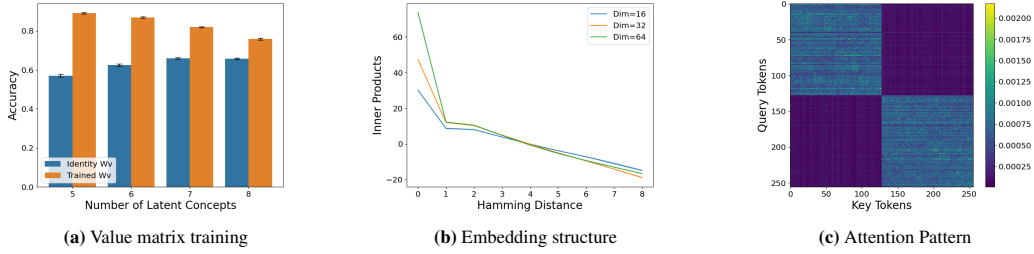


Figure 3: Key components of the single-layer transformer working together on the latent concept association problem. (a) Fixing the value matrix W_V as the identity matrix results in lower accuracy compared to training W_V . The figure reports average accuracy for both fixed and trained W_V with $L = 64$. (b) When training in the underparameterized regime, the embedding structure is approximated by (4.2). The graph displays the average inner product between embeddings of two tokens against the corresponding Hamming distance between these tokens when $m = 8$. (c) The self-attention layer can select tokens within the same cluster. The figure shows average attention score heat map with $m = 8$ and the cluster structure from Appendix C.1.

605 C Additional theoretical results

606 C.1 The role of attention selection

607 As of now, attention does not play a significant role in the analysis. But perhaps unsurprisingly, the
 608 attention mechanism is useful in selecting relevant information. To see this, let's consider a specific
 609 setting where for any latent vector z^* , $\mathcal{N}(z^*) = \{z : z_1^* = z_1\} \setminus \{z^*\}$.

610 Essentially, latent vectors are partitioned into two clusters based on the value of the first latent variable,
 611 and the informative conditional distribution π only samples latent vectors that are in the same cluster
 612 as the output latent vector. Empirically, when trained under this setting, the attention mechanism
 613 will pay more attention to tokens within the same cluster (Appendix D.3). This implies that the
 614 self-attention layer can mitigate noise and concentrate on the informative conditional distribution π .

615 To understand this more intuitively, we will study the gradient of unnormalized attention scores. In
 616 particular, the unnormalized attention score is defined as:

$$u_{t,t'} = (W_K W_E(t))^T (W_Q W_E(t')) / \sqrt{d_a}.$$

617 **Lemma 6.** Suppose the data generating process follows Section 3.1 and $\mathcal{N}(z^*) = \{z : z_1^* =$
 618 $z_1\} \setminus \{z^*\}$. Given the last token in the sequence t_L , then

$$\nabla_{u_{t,t_L}} \ell(f^L) = \nabla \ell(f^L)^T (W_E)^T W^V (\alpha_t \hat{p}_t W_E(t) - \hat{p}_t \sum_{l=1}^L \hat{p}_l W_E(t_l))$$

619 where for token t , $\alpha_t = \sum_{l=1}^L \mathbf{1}[t_l = t]$ and \hat{p}_t is the normalized attention score for token t .

620 Typically, α_t is larger when token t and t_L belong to the same cluster because tokens within the
 621 same cluster tend to co-occur frequently. As a result, the gradient contribution to the unnormalized
 622 attention score is usually larger for tokens within the same cluster.

623 D Experiments

624 The main implications of the theoretical results in the previous section are:

- 625 1. The value matrix is important and has associative memory structure as in (4.1).
- 626 2. Training embeddings is crucial in the underparameterized regime, where embeddings exhibit
 627 certain geometric structures.
- 628 3. Attention mechanism is used to select the most relevant tokens.

629 To evaluate these claims, we conduct several experiments on synthetic datasets. Additional experi-
 630 mental details and results can be found in Appendix G.

631 **D.1 On the value matrix W_V**

632 In this section, we study the necessity of the value matrix W_V and its structure. First, we conduct ex-
633 periments to compare the effects of training versus freezing W_V as the identity matrix, with the context
634 lengths L set to 64 and 128. Figure 3a and Figure G.1 show that when the context length is small, freez-
635 ing W_V can lead to a significant decline in accuracy. This is inline with Lemma 2 and validates it in a
636 general setting, implying the significance of the value matrix in maintaining a high memory recall rate.

637 Next, we investigate the degree of alignment between the trained value matrix W_V and the con-
638 struction in (4.1). The first set of experiments examines the similarity in functionality between the
639 two matrices. We replace value matrices in trained transformers with the constructed ones like in
640 (4.1) and then report accuracy with the new value matrix. As a baseline, we also consider randomly
641 constructed value matrix, where the outer product pairs are chosen randomly (detailed construction
642 can be found in Appendix G.1). Figure G.2 indicates that the accuracy does not significantly decrease
643 when the value matrix is replaced with the constructed ones. Furthermore, not only are the constructed
644 value matrix and the trained value matrix functionally alike, but they also share similar low-rank
645 approximations. We use singular value decomposition to get the best low rank approximations of
646 various value matrices where the rank is set to be the same as the number of latent variables (m). We
647 then compute smallest principal angles between low-rank approximations of trained value matrices
648 and those of constructed, randomly constructed, and Gaussian-initialized value matrices. Figure G.3
649 shows that the constructed ones have, on average, smallest principal angles with the trained ones.

650 **D.2 On the embeddings**

651 In this section, we explore the significance of embedding training in the underparameterized regime
652 and embedding structures. We conduct experiments to compare the effects of training versus freezing
653 embeddings with different embedding dimensions. The learning rate is selected as the best option
654 from $\{0.01, 0.001\}$ depending on the dimensions. Figure G.4 clearly shows that when the dimension
655 is smaller than the vocabulary size ($d < V$), embedding training is required. It is not necessary in
656 the overparameterized regime ($d > V$), partially confirming Theorem 1 because if embeddings are
657 initialized from a high-dimensional multi-variate Gaussian, they are approximately orthogonal to
658 each other and have the same norms.

659 The next question is what kind of embedding structures are formed for trained transformers in the
660 underparameterized regime. From Figure 3b and Figure G.5, it is evident that the relationship between
661 the average inner product of embeddings for two tokens and their corresponding Hamming distance
662 roughly aligns with (4.2). Perhaps surprisingly, if we plot the same graph for trained transformers
663 with a fixed identity value matrix, the relationship is mostly linear as shown in Figure G.6, confirming
664 our theory (Theorem 3).

665 As suggested in Section 4.3, such embedding geometry (4.2) can lead to low rank structures. We verify
666 this claim by studying the spectrum of the embedding matrix W_E . As illustrated in Appendix G.4,
667 Figure G.9-G.12 demonstrate that there are eigengaps between top and bottom singular values,
668 suggesting low-rank structures.

669 **D.3 On the attention selection mechanism**

670 In this section, we examine the role of attention pattern by considering a special class of latent
671 concept association model as defined in Appendix C.1. Figure 3c and Figure G.7 clearly show
672 that the self-attention select tokens in the same clusters. This suggests that attention can filter out
673 noise and focus on the informative conditional distribution π . We extend experiments to consider
674 cluster structures that depend on the first two latent variables (detailed construction can be found in
675 Appendix G.3) and Figure G.8 shows attention pattern as expected.

676 **E Additional Theoretical Results and Proofs**

677 **E.1 Proofs for Section 4.1**

678 Theorem 1 can be stated more formally as follows:

679 **Theorem 7.** Suppose the data generating process follows Section 3.1 where $m \geq 3$, $\omega = 1$, and
680 $\mathcal{N}(t) = V \setminus \{t\}$. Assume there exists a single layer transformer given by (3.1) such that a) $W_K = 0$
681 and $W_Q = 0$, b) Each row of W_E is orthogonal to each other and normalized, and c) W_V is given by

$$W_V = \sum_{i \in [V]} W_E(i) \left(\sum_{j \in \mathcal{N}_1(i)} W_E(j)^T \right).$$

682 Then if $L > \max\left\{ \frac{100m^2 \log(3/\varepsilon)}{(\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}))^2}, \frac{80m^2 |\mathcal{N}(y)|}{(\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}))^2} \right\}$ for any y , then

$$R_{\mathcal{D}^L}(f^L) \leq \varepsilon,$$

683 where $0 < \varepsilon < 1$.

684 *Proof.* First of all, the error is defined to be:

$$\begin{aligned} R_{\mathcal{D}^L}(f^L) &= \mathbb{P}_{(x,y) \sim \mathcal{D}^L} [\operatorname{argmax} f^L(x) \neq y] \\ &= \mathbb{P}_y \mathbb{P}_{x|y} [\operatorname{argmax} f^L(x) \neq y] \end{aligned}$$

685 Let's focus on the conditional probability $\mathbb{P}_{x|y} [\operatorname{argmax} f^L(x) \neq y]$.

686 By construction, the single layer transformer model has uniform attention. Therefore,

$$h(x) = \sum_{i \in \mathcal{N}(y)} \alpha_i W_E(i)$$

687 where $\alpha_i = \frac{1}{L} \sum_{k=1}^L \mathbf{1}\{t_k = i\}$ which is the number of occurrence of token i in the sequence.

688 By the latent concept association model, we know that

$$p(i|y) = \frac{\exp(-D_H(i, y)/\beta)}{Z}$$

689 where $Z = \sum_{i \in \mathcal{N}(y)} \exp(-D_H(i, y)/\beta)$.

690 Thus, the logit for token y is

$$f_y^L(x) = \sum_{i \in \mathcal{N}_1(y)} \alpha_i$$

691 And the logit for any other token \tilde{y} is

$$f_{\tilde{y}}^L(x) = \sum_{i \in \mathcal{N}_1(\tilde{y})} \alpha_i$$

692 For the prediction to be correct, we need

$$\max_{\tilde{y}} f_y^L(x) - f_{\tilde{y}}^L(x) > 0$$

693 By Lemma 3 of [Dev83], we know that for all $\Delta \in (0, 1)$, if $\frac{|\mathcal{N}(y)|}{L} \leq \frac{\Delta^2}{20}$, we have

$$\mathbb{P}\left(\max_{i \in \mathcal{N}(y)} |\alpha_i - p(i|y)| > \Delta \right) \leq \mathbb{P}\left(\sum_{i \in \mathcal{N}(y)} |\alpha_i - p(i|y)| > \Delta \right) \leq 3 \exp(-L\Delta^2/25)$$

694 Therefore, if $L \geq \max\left\{ \frac{25 \log(3/\varepsilon)}{\Delta^2}, \frac{20|\mathcal{N}(y)|}{\Delta^2} \right\}$, then with probability at least $1 - \varepsilon$, we have,

$$\max_{i \in \mathcal{N}(y)} |\alpha_i - p(i|y)| \leq \Delta$$

$$\begin{aligned}
f_y^L(x) - f_{\tilde{y}}^L(x) &= \sum_{i \in \mathcal{N}_1(y)} \alpha_i - \sum_{j \in \mathcal{N}_1(\tilde{y})} \alpha_j \\
&= \sum_{i \in \mathcal{N}_1(y)} \alpha_i - \sum_{i \in \mathcal{N}_1(y)} p(i|y) + \sum_{i \in \mathcal{N}_1(y)} p(i|y) \\
&\quad - \sum_{j \in \mathcal{N}_1(\tilde{y})} p(j|y) + \sum_{j \in \mathcal{N}_1(\tilde{y})} p(j|y) - \sum_{j \in \mathcal{N}_1(\tilde{y})} \alpha_j \\
&\geq \sum_{i \in \mathcal{N}_1(y)} p(i|y) - \sum_{j \in \mathcal{N}_1(\tilde{y})} p(j|y) - 2m\Delta \\
&\geq \exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}) - 2m\Delta
\end{aligned}$$

695 Note that because of Lemma 10, there's no neighboring set that is the superset of another.

696 Therefore as long as $\Delta < \frac{\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta})}{2m}$,

$$f_y^L(x) - f_{\tilde{y}}^L(x) > 0$$

697 for any \tilde{y} .

698 Finally, if $L > \max\left\{\frac{100m^2 \log(3/\varepsilon)}{(\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}))^2}, \frac{80m^2 |\mathcal{N}(y)|}{(\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}))^2}\right\}$ for any y , then

$$\mathbb{P}_{x|y}[\operatorname{argmax} f^L(x) \neq y] \leq \varepsilon$$

699 And

$$\begin{aligned}
R_{\mathcal{D}^L}(f^L) &= \mathbb{P}_{(x,y) \sim \mathcal{D}^L}[\operatorname{argmax} f^L(x) \neq y] \\
&= \mathbb{P}_y \mathbb{P}_{x|y}[\operatorname{argmax} f^L(x) \neq y] \leq \varepsilon
\end{aligned}$$

700

□

701 E.2 Proofs for Section 4.2

702 **Lemma 2.** Suppose the data generating process follows Section 3.1 where $m \geq 3$, $\omega = 1$ and
703 $\mathcal{N}(t) = \{t' : D_H(t, t') = 1\}$. For any single layer transformer given by (3.1) where each row of
704 W_E is orthogonal to each other and normalized, if W_V is constructed as in (4.1), then the error rate
705 is 0. If W_V is the identity matrix, then the error rate is strictly larger than 0.

706 *Proof.* Following the proof for Theorem 7, let's focus on the conditional probability:

$$\mathbb{P}_{x|y}[\operatorname{argmax} f^L(x) \neq y]$$

707 By construction, we have

$$h(x) = \sum_{i \in \mathcal{N}_1(y)} \alpha_i W_E(i)$$

708 where $\alpha_i = \frac{1}{L} \sum_{k=1}^L \mathbf{1}\{t_k = i\}$ which is the number of occurrence of token i in the sequence.

709 Let's consider the first case where W_V is constructed as in (4.1). Then we know that for some other
710 token $\tilde{y} \neq y$,

$$f_y^L(x) - f_{\tilde{y}}^L(x) = \sum_{i \in \mathcal{N}_1(y)} \alpha_i - \sum_{i \in \mathcal{N}_1(\tilde{y})} \alpha_i = 1 - \sum_{i \in \mathcal{N}_1(\tilde{y})} \alpha_i$$

711 By Lemma 10, we have that for any token $\tilde{y} \neq y$,

$$f_y^L(x) - f_{\tilde{y}}^L(x) > 0$$

712 Therefore, the error rate is always 0.

713 Now let's consider the second case where W_V is the identity matrix. Let j be a token in the set $\mathcal{N}_1(y)$.
 714 Then there is a non-zero probability that context x contains only j . In that case,

$$h(x) = W_E(j)$$

715 However, we know that by the assumption on the embedding matrix,

$$f_y^L(x) - f_j^L(x) = (W_E(y) - W_E(j))^T h(x) = -\|W_E(j)\|^2 < 0$$

716 This implies that there's non zero probability that y is misclassified. Therefore, when W_V is the
 717 identity matrix, the error rate is strictly larger than 0. \square

718 **Theorem 3.** Suppose the data generating process follows Section 3.1 where $m \geq 3$, $\omega = 1$ and
 719 $\mathcal{N}(t) = V \setminus \{t\}$. For any single layer transformer given by (3.1) with W_V being the identity matrix,
 720 if the cross entropy loss is minimized so that for any sampled pair (x, y) ,

$$p(y|x) = \hat{p}(y|x) = \text{softmax}(f_y^L(x))$$

721 there exists $a > 0$ and b such that for two tokens $t \neq t'$,

$$\langle W_E(t), W_E(t') \rangle = -aD_H(t, t') + b$$

722 *Proof.* Because for any pair of (x, y) , the estimated conditional probability matches the true condi-
 723 tional probability. In particular, let's consider two target tokens y_1, y_2 and context $x = (t_i, \dots, t_i)$ for
 724 some token t_i such that $p(x|y_1) > 0$ and $p(x|y_2) > 0$, then

$$\frac{p(y_1|x)}{p(y_2|x)} = \frac{p(x|y_1)p(y_1)}{p(x|y_2)p(y_2)} = \frac{p(x|y_1)}{p(x|y_2)} = \frac{\hat{p}(x|y_1)}{\hat{p}(x|y_2)} = \exp((W_E(y_1) - W_E(y_2))^T h(x))$$

725 The second equality is because $p(y)$ is the uniform distribution. By our construction,

$$\frac{p(x|y_1)}{p(x|y_2)} = \frac{p(t_i|y_1)^L}{p(t_i|y_2)^L} = \exp((W_E(y_2) - W_E(y_1))^T h(x)) = \exp((W_E(y_1) - W_E(y_2))^T W_E(t_i))$$

726 By the data generating process, we have that

$$\frac{L}{\beta}(D_H(t_i, y_2) - D_H(t_i, y_1)) = (W_E(y_1) - W_E(y_2))^T W_E(t_i)$$

727 Let $t_i = y_3$ such that $y_3 \neq y_1, y_3 \neq y_2$, then

$$\frac{L}{\beta}D_H(y_3, y_1) - W_E(y_1)^T W_E(y_3) = \frac{L}{\beta}D_H(y_3, y_2) - W_E(y_2)^T W_E(y_3)$$

728 For simplicity, let's define

$$\Psi(y_1, y_2) = \frac{L}{\beta}D_H(y_1, y_2) - W_E(y_1)^T W_E(y_2)$$

729 Therefore,

$$\Psi(y_3, y_1) = \Psi(y_3, y_2)$$

730 Now consider five distinct labels: y_1, y_2, y_3, y_4, y_5 . We have,

$$\Psi(y_3, y_1) = \Psi(y_3, y_2) = \Psi(y_4, y_2) = \Psi(y_4, y_5)$$

731 In other words, $\Psi(y_3, y_1) = \Psi(y_4, y_5)$ for arbitrarily chosen distinct labels y_1, y_3, y_4, y_5 . Therefore,
 732 $\Psi(t, t')$ is a constant for $t \neq t'$.

733 For any two tokens $t \neq t'$,

$$\frac{L}{\beta}D_H(t, t') - W_E(t)^T W_E(t') = C$$

734 Thus,

$$W_E(t)^T W_E(t') = -\frac{L}{\beta}D_H(t, t') + C$$

735 \square

736 **E.3 Proofs for Section 4.3**

737 Theorem 4 can be formalized as the following theorem.

738 **Theorem 8.** *Following the same setup as in Theorem 7, but embeddings follow (4.2) then if $b > 0$,*
 739 $\Delta_1 > 0$, $0 < \Delta < \frac{\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta})}{2m}$, $L \geq \max\{\frac{25 \log(3/\varepsilon)}{\Delta^2}, \frac{20|\mathcal{N}(y)|}{\Delta^2}\}$ for any y , and

$$0 < a < \frac{2 \exp(\frac{1}{\beta})}{(|V| - 2)m^2}$$

740 and

$$b_0 > \max\left\{\frac{a(m-2)m + \Delta_1}{\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}) - 2m\Delta} + b, \frac{(b-a)\Delta_1 - \frac{|V|-2}{2}abm^2 \exp(-\frac{1}{\beta}) + \frac{|V|-2}{2}a^2(m-2)m^2}{1 - \frac{|V|-2}{2}am^2 \exp(-\frac{1}{\beta})}\right\}$$

741 we have

$$R_{\mathcal{D}^L}(f^L) \leq \varepsilon$$

742 where $0 < \varepsilon < 1$.

743 *Proof.* Following the proof of Theorem 7, let's also focus on the conditional probability

$$\mathbb{P}_{x|y}[\operatorname{argmax} f^L(x) \neq y]$$

744 By construction, the single layer transformer model has uniform attention. Therefore,

$$h(x) = \sum_{i \in \mathcal{N}(y)} \alpha_i W_E(i)$$

745 where $\alpha_i = \frac{1}{L} \sum_{k=1}^L \mathbf{1}\{t_k = i\}$ which is the number of occurrence of token i in the sequence. For
 746 simplicity, let's define $\alpha_y = 0$ such that

$$h(x) = \sum_{i \in [V]} \alpha_i W_E(i)$$

747 Similarly, we also have that if $L \geq \max\{\frac{25 \log(3/\varepsilon)}{\Delta^2}, \frac{20|\mathcal{N}(y)|}{\Delta^2}\}$, then with probability at least $1 - \varepsilon$,
 748 we have,

$$\max_{i \in [V]} |\alpha_i - p(i|y)| \leq \Delta$$

749 Also define the following:

$$\begin{aligned} \phi_k(x) &= \sum_{j \in \mathcal{N}_1(k)} W_E(j)^T \left(\sum_{i \in [V]} \alpha_i W_E(i) \right) \\ v_k(y) &= W_E(y)^T W_E(k) \end{aligned}$$

750 Thus, the logit for token y is

$$f_y^L(x) = \sum_{k=0}^{|V|-1} v_k(y) \phi_k(x)$$

751 Let's investigate $\phi_k(x)$. By Lemma 9,

$$\begin{aligned} \phi_k(x) &= \sum_{i \in [V]} \alpha_i \left(\sum_{j \in \mathcal{N}_1(k)} W_E(j)^T W_E(i) \right) \\ &= (b_0 - b) \sum_{j \in \mathcal{N}_1(k)} \alpha_j + \sum_{i \in [V]} \alpha_i (-a(m-2)D_H(k, i) + (b-a)m) \end{aligned}$$

752 Thus, for any $k_1, k_2 \in [V]$,

$$\begin{aligned} \phi_{k_1}(x) - \phi_{k_2}(x) &= (b_0 - b) \left(\sum_{j_1 \in \mathcal{N}_1(k_1)} \alpha_{j_1} - \sum_{j_2 \in \mathcal{N}_1(k_2)} \alpha_{j_2} \right) \\ &\quad + \sum_{i \in [V]} \alpha_i a(m-2) (D_H(k_2, i) - D_H(k_1, i)) \end{aligned}$$

753 Because $-m \leq D_H(k_2, i) - D_H(k_1, i) \leq m$, we have

$$\begin{aligned} (b_0 - b) \left(\sum_{j_1 \in \mathcal{N}_1(k_1)} \alpha_{j_1} - \sum_{j_2 \in \mathcal{N}_1(k_2)} \alpha_{j_2} \right) - a(m-2)m \\ \leq \phi_{k_1}(x) - \phi_{k_2}(x) \leq \\ (b_0 - b) \left(\sum_{j_1 \in \mathcal{N}_1(k_1)} \alpha_{j_1} - \sum_{j_2 \in \mathcal{N}_1(k_2)} \alpha_{j_2} \right) + a(m-2)m \end{aligned}$$

754 For prediction to be correct, we need

$$\max_{\tilde{y}} f_{\tilde{y}}^L(x) - f_{\tilde{y}}^L(x) > 0$$

755 This also means that

$$\max_{\tilde{y}} \sum_{k=0}^{|V|-1} (v_k(y) - v_k(\tilde{y})) \phi_k(x) > 0$$

756 One can show that for any k , if $\iota^{-1}(\tilde{k}) = \iota^{-1}(y) \otimes \iota^{-1}(\tilde{y}) \otimes \iota^{-1}(k)$ where \otimes means bitwise XOR,
757 then

$$v_k(y) - v_k(\tilde{y}) = v_{\tilde{k}}(\tilde{y}) - v_{\tilde{k}}(y) \tag{E.1}$$

758 First of all, if $k = y$, then $\tilde{k} = \tilde{y}$, which means

$$v_k(y) - v_k(\tilde{y}) = v_{\tilde{k}}(\tilde{y}) - v_{\tilde{k}}(y) = b_0 + aD_H(y, \tilde{y}) - b$$

759 If $k \neq y, \tilde{y}$, then (E.1) implies that

$$D_H(k, y) - D_H(k, \tilde{y}) = D_H(\tilde{k}, \tilde{y}) - D_H(\tilde{k}, y)$$

760 We know that $D_H(k, y)$ is the number of 1s in $\iota^{-1}(k) \otimes \iota^{-1}(y)$ and,

$$\iota^{-1}(\tilde{k}) \otimes \iota^{-1}(y) = \iota^{-1}(y) \otimes \iota^{-1}(\tilde{y}) \otimes \iota^{-1}(k) \otimes \iota^{-1}(y) = \iota^{-1}(\tilde{y}) \otimes \iota^{-1}(k)$$

761 Similarly,

$$\iota^{-1}(\tilde{k}) \otimes \iota^{-1}(\tilde{y}) = \iota^{-1}(y) \otimes \iota^{-1}(k)$$

762 Therefore, (E.1) holds and we can rewrite $f_{\tilde{y}}^L(x) - f_{\tilde{y}}^L(x)$ as

$$\begin{aligned} f_{\tilde{y}}^L(x) - f_{\tilde{y}}^L(x) &= \sum_{k=0}^{|V|-1} (v_k(y) - v_k(\tilde{y})) \phi_k(x) \\ &= (b_0 - b + aD_H(y, \tilde{y})) (\phi_y(x) - \phi_{\tilde{y}}(x)) \\ &\quad + \sum_{k \neq y, \tilde{y}, D_H(k, y) \geq D_H(k, \tilde{y})} a(D_H(k, y) - D_H(k, \tilde{y})) (\phi_k(x) - \phi_{\tilde{k}}(x)) \end{aligned}$$

763 We already know that $b_0 > b > 0$ and $a > 0$, thus, $b_0 - b + aD_H(y, \tilde{y}) > 0$ for any pair y, \tilde{y} .

764 We also want $\phi_y(x) - \phi_{\tilde{y}}(x)$ to be positive. Note that

$$\phi_y(x) - \phi_{\tilde{y}}(x) \geq (b_0 - b) \left(\exp\left(-\frac{1}{\beta}\right) - \exp\left(-\frac{2}{\beta}\right) - 2m\Delta \right) - a(m-2)m$$

765 We need $\Delta < \frac{\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta})}{2m}$ and for some positive $\Delta_1 > 0$, b_0 needs to be large enough such that

$$\phi_y(x) - \phi_{\tilde{y}}(x) > \Delta_1$$

766 which implies that

$$b_0 > \frac{a(m-2)m + \Delta_1}{\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}) - 2m\Delta} + b \quad (\text{E.2})$$

767 On the other hand, for $k \neq y, \tilde{y}$, we have

$$\begin{aligned} \phi_k(x) - \phi_{\tilde{k}}(x) &\geq (b_0 - b) \left(\sum_{j_1 \in \mathcal{N}_1(k)} \alpha_{j_1} - \sum_{j_2 \in \mathcal{N}_1(\tilde{k})} \alpha_{j_2} \right) - a(m-2)m \\ &\geq (b_0 - b) \left(-(m-1) \exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}) - 2m\Delta \right) - a(m-2)m \\ &\geq (b_0 - b) \left(-(m-1) \exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}) + \exp(-\frac{2}{\beta}) - \exp(-\frac{1}{\beta}) \right) - a(m-2)m \\ &\geq -(b_0 - b)m \exp(-\frac{1}{\beta}) - a(m-2)m \end{aligned}$$

768 Then, we have

$$\begin{aligned} f_y^L(x) - f_{\tilde{y}}^L(x) &\geq (b_0 - b + a)\Delta_1 - \frac{|V|-2}{2} \left((b_0 - b)am^2 \exp(-\frac{1}{\beta}) + a^2(m-2)m^2 \right) \\ &\geq \left(1 - \frac{|V|-2}{2} am^2 \exp(-\frac{1}{\beta}) \right) b_0 - (b-a)\Delta_1 + \frac{|V|-2}{2} abm^2 \exp(-\frac{1}{\beta}) - \frac{|V|-2}{2} a^2(m-2)m^2 \end{aligned}$$

769 The lower bound is independent of \tilde{y} , therefore, we need it to be positive to ensure the prediction is
770 correct. To achieve this, we want

$$1 - \frac{|V|-2}{2} am^2 \exp(-\frac{1}{\beta}) > 0$$

771 which implies that

$$a < \frac{2 \exp(\frac{1}{\beta})}{(|V|-2)m^2} \quad (\text{E.3})$$

772 And finally we need

$$b_0 > \frac{(b-a)\Delta_1 - \frac{|V|-2}{2} abm^2 \exp(-\frac{1}{\beta}) + \frac{|V|-2}{2} a^2(m-2)m^2}{1 - \frac{|V|-2}{2} am^2 \exp(-\frac{1}{\beta})} \quad (\text{E.4})$$

773 To summarize, if $b > 0$, $\Delta_1 > 0$, $0 < \Delta < \frac{\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta})}{2m}$, $L \geq \max\{\frac{25 \log(3/\varepsilon)}{\Delta^2}, \frac{20|\mathcal{N}(y)|}{\Delta^2}\}$ for
774 any y , and

$$0 < a < \frac{2 \exp(\frac{1}{\beta})}{(|V|-2)m^2}$$

775 and

$$b_0 > \max\left\{ \frac{a(m-2)m + \Delta_1}{\exp(-\frac{1}{\beta}) - \exp(-\frac{2}{\beta}) - 2m\Delta} + b, \frac{(b-a)\Delta_1 - \frac{|V|-2}{2} abm^2 \exp(-\frac{1}{\beta}) + \frac{|V|-2}{2} a^2(m-2)m^2}{1 - \frac{|V|-2}{2} am^2 \exp(-\frac{1}{\beta})} \right\}$$

776 we have

$$R_{\mathcal{D}^L}(f^L) \leq \varepsilon$$

777 where $0 < \varepsilon < 1$.

778

□

779 **Lemma 5.** If embeddings follow (4.2) and $b = b_0$ and $\mathcal{N}(t) = V \setminus \{t\}$, then $\text{rank}(W_E) \leq m + 2$.

780 *Proof.* By (4.2), we have that

$$\langle W_E(i), W_E(j) \rangle = -aD_H(i, j) + b$$

781 Therefore,

$$(W_E)^T W_E = -aD_H + b\mathbf{1}\mathbf{1}^T$$

782 Let's first look at D_H which has rank at most $m + 1$. To see this, let's consider a set of $m + 1$ tokens:

783 $\{e_0, e_1, \dots, e_m\} \subseteq V$ where $e_k = 2^k$. Here e_0 is associated with the latent vector of all zeroes and
784 the latent vector associated with e_k has only the k -th latent variable being 1.

785 On the other hand, for any token i , we have that,

$$i = \sum_{k:\iota^{-1}(i)_k=1} e_k$$

786 In fact,

$$D_H(i) = \sum_{k:\iota^{-1}(i)_k=1} \left(D_H(e_k) - D_H(e_0) \right) + D_H(e_0)$$

787 where $D_H(i)$ is the i -th row of D_H , and for each entry j of $D_H(i)$, we have that

$$D_H(i, j) = \sum_{k:\iota^{-1}(i)_k=1} \left(D_H(e_k, j) - D_H(e_0, j) \right) + D_H(e_0, j)$$

788 This is because

$$D_H(e_k, j) - D_H(e_0, j) = \begin{cases} +1 & \text{if } \iota^{-1}(j)_k = 0 \\ -1 & \text{if } \iota^{-1}(j)_k = 1 \end{cases}$$

789 Thus, we can rewrite $D_H(i, j)$ as

$$\begin{aligned} D_H(i, j) &= \sum_{k:\iota^{-1}(i)_k=1} \left(\mathbf{1}[\iota^{-1}(i)_k = 1, \iota^{-1}(j)_k = 0] - \mathbf{1}[\iota^{-1}(i)_k = 1, \iota^{-1}(j)_k = 1] \right) + D_H(e_0, j) \\ &= \sum_{k=1}^m \left(\mathbf{1}[\iota^{-1}(i)_k = 1, \iota^{-1}(j)_k = 0] - \mathbf{1}[\iota^{-1}(i)_k = 1, \iota^{-1}(j)_k = 1] \right) \\ &\quad + \sum_{k=1}^m \left(\mathbf{1}[\iota^{-1}(i)_k = 0, \iota^{-1}(j)_k = 1] + \mathbf{1}[\iota^{-1}(i)_k = 1, \iota^{-1}(j)_k = 1] \right) \\ &= \sum_{k=1}^m \mathbf{1}[\iota^{-1}(i)_k = 1, \iota^{-1}(j)_k = 0] + \mathbf{1}[\iota^{-1}(i)_k = 0, \iota^{-1}(j)_k = 1] \\ &= D_H(i, j) \end{aligned}$$

790 Therefore, every row of D_H can be written as a linear combination of
791 $\{D_H(e_0), D_H(e_1), \dots, D_H(e_m)\}$. In other words, D_H has rank at most $m + 1$.

792 Therefore,

$$\text{rank}((W_E)^T W_E) = \text{rank}(W_E) \leq m + 2.$$

793

□

794 **Lemma 9.** Let $z^{(0)}$ and $z^{(1)}$ be two binary vectors of size m where $m \geq 2$. Then,

$$\sum_{z:D_H(z^{(0)}, z)=1} D_H(z, z^{(1)}) = (m - 2)D_H(z^{(0)}, z^{(1)}) + m$$

795 *Proof.* For z such that $D_H(z, z^{(0)}) = 1$, we know that there are two cases. Either z differs with $z^{(0)}$
796 on a entry but agrees with $z^{(1)}$ on that entry or z differs with both $z^{(0)}$ and $z^{(1)}$.

797 For the first case, we know that there are $D_H(z^{(0)}, z^{(1)})$ such entries. In this case, $D_H(z, z^{(1)}) =$
798 $D_H(z^{(0)}, z^{(1)}) - 1$. For the second case, $D_H(z, z^{(1)}) = D_H(z^{(0)}, z^{(1)}) + 1$.

799 Therefore,

$$\begin{aligned}
& \sum_{z: D_H(z, z^{(0)})=1} D_H(z, z^{(1)}) \\
&= D_H(z^{(0)}, z^{(1)})(D_H(z^{(0)}, z^{(1)}) - 1) + (m - D_H(z^{(0)}, z^{(1)}))(D_H(z^{(0)}, z^{(1)}) + 1) \\
&= (m - 2)D_H(z^{(0)}, z^{(1)}) + m
\end{aligned}$$

800

□

801 **Lemma 10.** *If $m \geq 3$ and $\mathcal{N}(t) = V \setminus \{t\}$, then $\mathcal{N}_1(t) \not\subseteq \mathcal{N}_1(t')$ for any $t, t' \in [V]$.*

802 *Proof.* For any token t , $\mathcal{N}_1(t)$ contains any token t' such that $D_H(t, t') = 1$ by the conditions. Then
803 given a set $\mathcal{N}_1(t)$, one can uniquely determine token t . This is because for the set of latent vectors
804 associated with $\mathcal{N}_1(t)$, at each index, there could only be one possible change. □

805 E.4 Proofs for Appendix C.1

806 **Lemma 6.** *Suppose the data generating process follows Section 3.1 and $\mathcal{N}(z^*) = \{z : z_1^* =$
807 $z_1\} \setminus \{z^*\}$. Given the last token in the sequence t_L , then*

$$\nabla_{u_{t,t_L}} \ell(f^L) = \nabla \ell(f^L)^T (W_E)^T W^V (\alpha_t \hat{p}_t W_E(t) - \hat{p}_t \sum_{l=1}^L \hat{p}_{t_l} W_E(t_l))$$

808 where for token t , $\alpha_t = \sum_{l=1}^L \mathbf{1}[t_l = t]$ and \hat{p}_t is the normalized attention score for token t .

809 *Proof.* Recall that,

$$\begin{aligned}
f^L(x) &= \left[W_E^T W_V \text{attn}(W_E \chi(x)) \right]_{:L} \\
&= W_E^T W_V \sum_{l=1}^L \frac{\exp(u_{t_l, t_L})}{Z} W_E(t_l)
\end{aligned}$$

810 where Z is a normalizing constant.

811 Define $\hat{p}_{t_l} = \frac{\exp(u_{t_l, t_L})}{Z}$. Then we have

$$f^L(x) = W_E^T W_V \sum_{l=1}^L \hat{p}_{t_l} W_E(t_l)$$

812 Note that if $t_l = t$ then,

$$\frac{\partial \hat{p}_{t_l}}{\partial u_{t,t_L}} = \hat{p}_{t_l} (1 - \hat{p}_{t_l})$$

813 Otherwise,

$$\frac{\partial \hat{p}_{t_l}}{\partial u_{t,t_L}} = -\hat{p}_{t_l} \hat{p}_t$$

814 By the chain rule, we know that

$$\nabla_{u_{t,t_L}} \ell(f^L) = \nabla \ell(f^L)^T (W_E)^T W^V \left(\sum_{l=1}^L \mathbf{1}[t_l = t] \hat{p}_{t_l} W_E(t) - \sum_{l=1}^L \hat{p}_{t_l} \hat{p}_t W_E(t_l) \right)$$

815 Therefore,

$$\nabla_{u_{t,t_L}} \ell(f^L) = \nabla \ell(f^L)^T (W_E)^T W^V (\alpha_t \hat{p}_t W_E(t) - \hat{p}_t \sum_{l=1}^L \hat{p}_{t_l} W_E(t_l))$$

816 where $\alpha_t = \sum_{l=1}^L \mathbf{1}[t_l = t]$. □

817 **F Additional experiments – context hijacking**

818 In this section, we show the results of additional context hijacking experiments on the COUNTERFACT
 819 dataset [Men+22].

820 **Reverse context hijacking** In Figure 2a, we saw the effects of hijacking by adding in “Do not think
 821 of {target_false}.” to each context. Now, we measure the effect of the reverse: What if we prepend
 822 “Do not think of {target_true}.” ?

823 Based on the study in this paper on how associative memory works in LLMs, we should expect the
 824 efficacy score to increase. Indeed, this is what happens, as we see in Figure F.1.

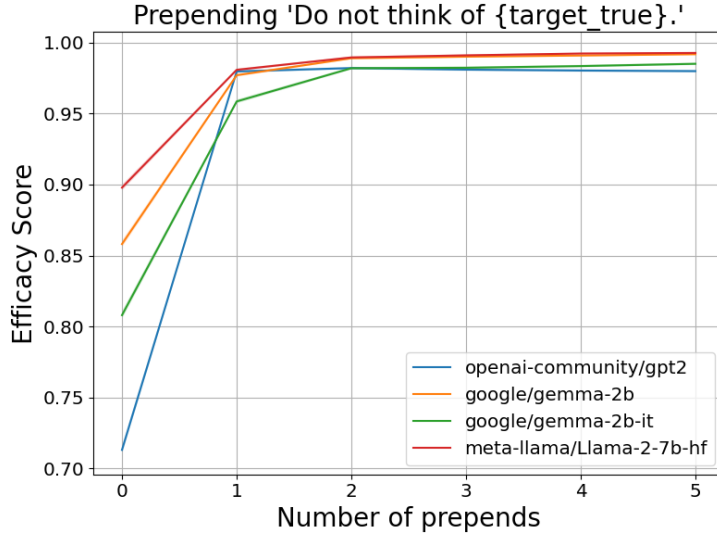


Figure F.1: Prepending ‘Do not think of {target_true}.’ can increase the chance of LLMs to output correct tokens. This figure shows efficacy score versus the number of preprends for various LLMs on the COUNTERFACT dataset with the reverse context hijacking scheme.

825 **Hijacking based on relation IDs** We first give an example of each of the 4 relation IDs we hijack
 826 in Table 1.

Table 1: Examples of contexts in Relation IDs from COUNTERFACT

RELATION ID r	CONTEXT p	TRUE TARGET o_*	FALSE TARGET o_-
P190	Kharkiv is a twin city of	Warsaw	Athens
P103	The native language of Anatole France is	French	English
P641	Hank Aaron professionally plays the sport	baseball	basketball
P131	Kalamazoo County can be found in	Michigan	Indiana

Table 2: Examples of hijack and reverse hijack formats based on Relation IDs

RELATION ID r	CONTEXT HIJACK SENTENCE	REVERSE CONTEXT HIJACK SENTENCE
P190	The twin city of {subject} is not {target_false}	The twin city of {subject} is {target_true}
P103	{subject} cannot speak {target_false}	{subject} can speak {target_true}
P641	{subject} does not play {target_false}	{subject} plays {target_true}
P131	{subject} is not located in {target_false}	{subject} is located in {target_true}

827 Similar to Figure 2b, we repeat the hijacking experiments where we prepend factual sentences
 828 generated from the relation ID. We use the format illustrated in Table 2 for the prepended sentences.

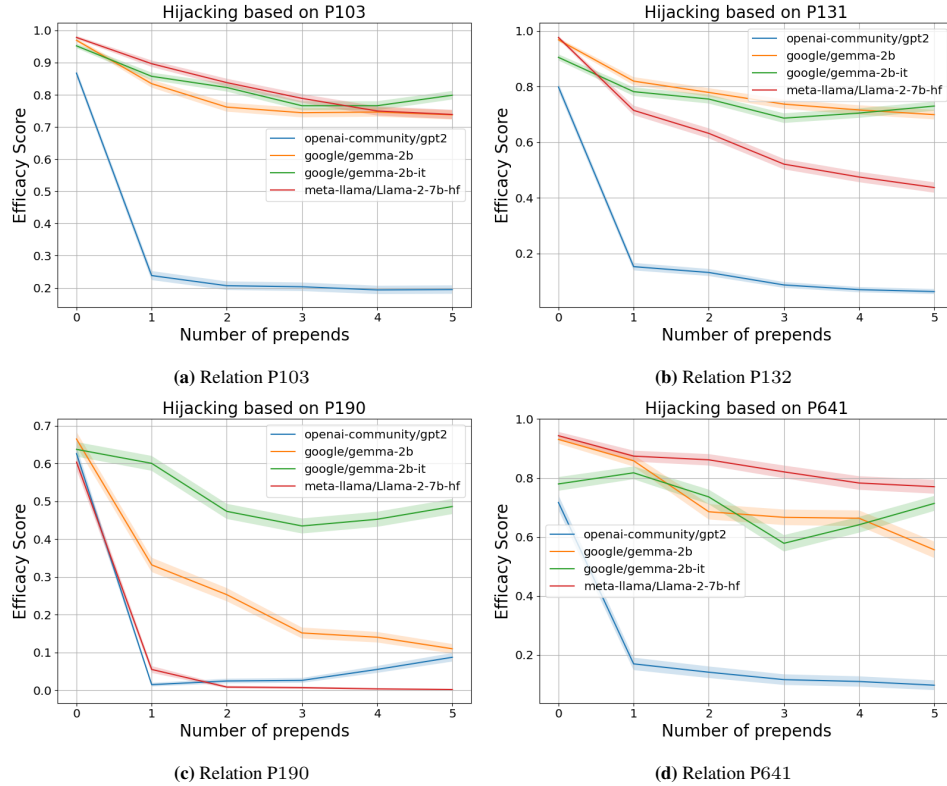


Figure F.2: Context hijacking based on relation IDs can result in LLMs output incorrect tokens. This figure shows efficacy score versus the number of preprends for various LLMs on the COUNTERFACT dataset with hijacking scheme presented in Table 2.

829 We experiment with 3 other relation IDs and we see similar trends for all the LLMs in Figure F.2a,
 830 F.2b, and F.2d. That is, the efficacy score drops for the first preprend and as we increase the number of
 831 preprends, the trend of ES dropping continues. Therefore, this confirms our intuition that LLMs can
 832 be hijacked by contexts without changing the factual meaning.

833 Similar to Figure F.1, we experiment with reverse context hijacking where we give the answers based
 834 on relation IDs, as shown in Table 2. We again experiment with the same 4 relation IDs and the
 835 results are in Figure F.3a - F.3d. We see that the efficacy score increases when we preprend the answer
 836 sentence, thereby verifying the observations of this study.

837 **Hijacking without exact target words** So far, the experiments use prompts that either contain
 838 true or false target words. It turns out, the inclusion of exact target words are not necessary. To see
 839 this, we experiment a variant of the generic hijacking and reverse hijacking experiments. But instead
 840 of saying “Do not think of {target_false}” or “Do not think of {target_true}”. We replace target
 841 words with words that are semantically close. In particular, for relation P1412, we replace words
 842 representing language (e.g., “French”) with their associated country name (e.g., “France”). As shown
 843 in Figure F.4, context hijacking and reverse hijacking still work in this case.

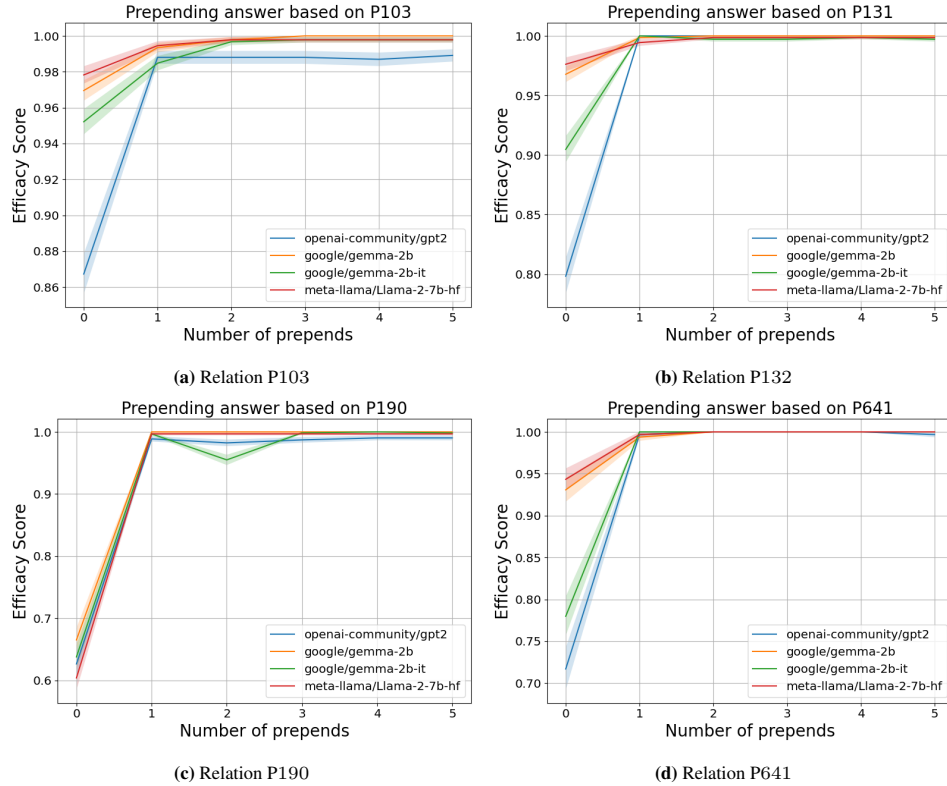


Figure F.3: Reverse context hijacking based on relation IDs can result in LLMs to be more likely to be correct. This figure shows efficacy score versus the number of preprends for various LLMs on the COUNTERFACT dataset with the reverse hijacking scheme presented in Table 2.

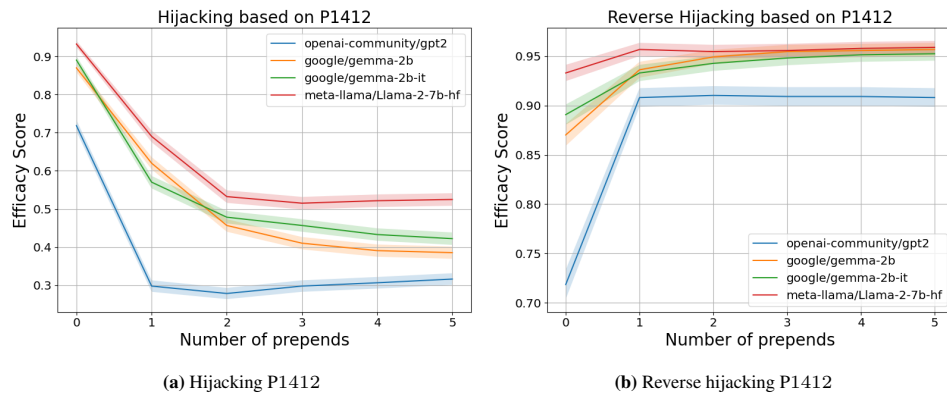


Figure F.4: Hijacking and reverse hijacking experiments on relation P1412 show that context hijacking does not require exact target word to appear in the context. This figure shows efficacy score versus the number of preprends for various LLMs on the COUNTERFACT dataset.

844 G Additional experiments and figures – latent concept association

845 In this appendix section, we present additional experimental details and results from the synthetic
846 experiments on latent concept association.

847 **Experimental setup** Synthetic data are generated following the model in Section 3.1. Unless
848 otherwise stated, the default setup has $\omega = 0.5$, $\beta = 1$ and $\mathcal{N}(i) = V \setminus \{i\}$ and $L = 256$. The
849 default hidden dimension of the one-layer transformer is also set to be 256. The model is optimized
850 using AdamW [LH17] where the learning rate is chosen from $\{0.01, 0.001\}$. The evaluation dataset

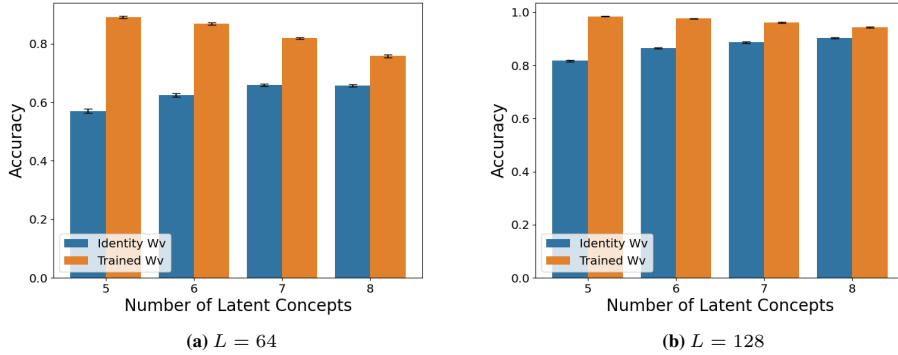


Figure G.1: Fixing the value matrix W_V as the identity matrix results in lower accuracy compared to training W_V , especially for smaller context length L . The figure reports accuracy for both fixed and trained W_V settings, with standard errors calculated over 10 runs.

851 is drawn from the same distribution as the training dataset and consists of 1024 (x, y) pairs. Although
 852 theoretical results in Section 4 may freeze certain parts of the network for simplicity, in this section,
 853 unless otherwise specified, all layers of the transformers are trained jointly. Also, in this section, we
 854 typically report accuracy which is $1 - \text{error}$.

855 G.1 On the value matrix W_V

856 In this section, we provide additional figures of Appendix D.1. Specifically, Figure G.1 shows that
 857 fixing the value matrix to be the identity will negatively impact accuracy. Figure G.2 indicates that re-
 858 placing trained value matrices with constructed ones can preserve accuracy to some extent. Figure G.3
 859 suggests that trained value matrices and constructed ones share similar low-rank approximations. For
 860 the last two sets of experiments, we consider randomly constructed value matrix, where the outer
 861 product pairs are chosen randomly, defined formally as follows:

$$W_V = \sum_{i \in [V]} W_E(i) \left(\sum_{\{j\} \sim \text{Unif}([V])^{|N_1(i)|}} W_E(j)^T \right)$$

862 G.2 On the embeddings

863 This section provides additional figures from Appendix D.2. Figure G.4 shows that in the under-
 864 parameterized regime, embedding training is required. Figure G.5 indicates that the embedding
 865 structure in the underparameterized regime roughly follows (4.2). Finally Figure G.6 shows that,
 866 when the value matrix is fixed to the identity, the relationship between inner product of embeddings
 867 and their corresponding Hamming distance is mostly linear.

868 G.3 On the attention selection mechanism

869 This section provides additional figures from Appendix D.3. Figure G.7-G.8 show that attention
 870 mechanism selects tokens in the same cluster as the last token. In particular, for Figure G.8, we
 871 extend experiments to consider cluster structures that depend on the first two latent variables. In other
 872 words, for any latent vector z^* , we have

$$\mathcal{N}(z^*) = \{z : z_1^* = z_1 \text{ and } z_2^* = z_2\} \setminus \{z^*\}$$

873 G.4 Spectrum of embeddings

874 We display several plots of embedding spectra (Figure G.9, Figure G.10, Figure G.11, Figure G.12)
 875 that exhibit eigengaps between the top and bottom eigenvalues, suggesting low-rank structures.

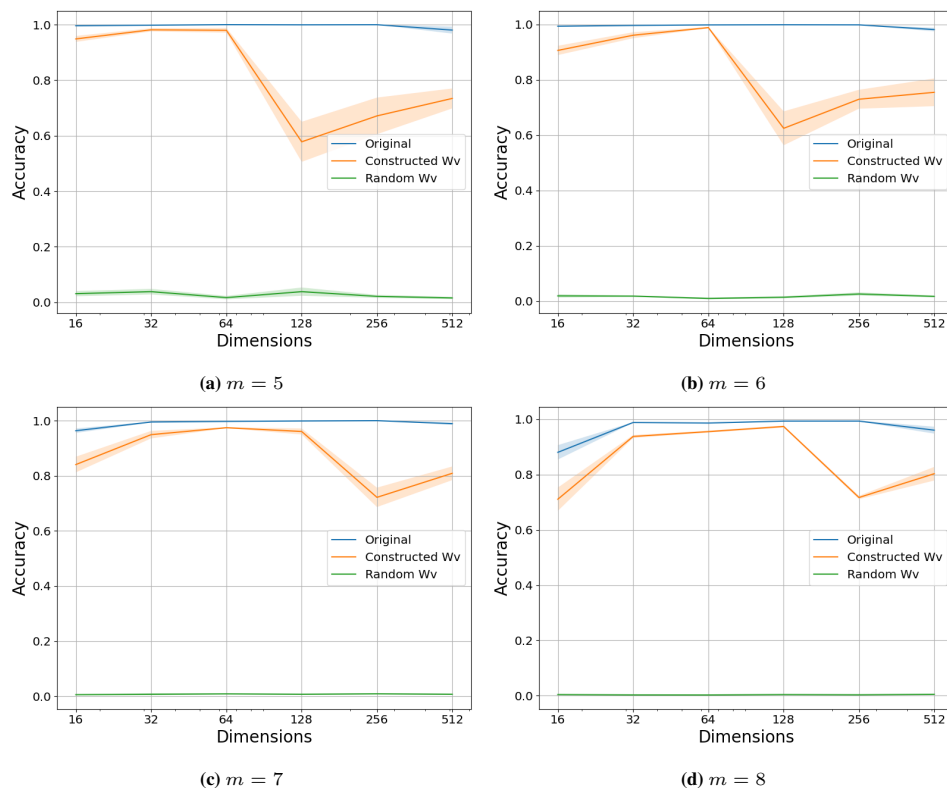


Figure G.2: When the value matrix is replaced with the constructed one in trained transformers, the accuracy does not significantly decrease compared to replacing the value matrix with randomly constructed ones. The graph reports accuracy under different embedding dimensions and standard errors are over 5 runs.

876 G.5 Context hijacking in latent concept association

877 In this section, we want to simulate context hijacking in the latent concept association model. To
 878 achieve that, we first sample two output tokens y^1 (true target) and y^2 (false target) and then generate
 879 contexts $x^1 = (t_1^1, \dots, t_L^1)$ and $x^2 = (t_1^2, \dots, t_L^2)$ from $p(x^1|y^1)$ and $p(x^2|y^2)$. Then we mix the two
 880 contexts with rate p_m . In other words, for the final mixed context $x = (t_1, \dots, t_L)$, t_l has probability
 881 $1 - p_m$ to be t_l^1 and p_m probability to be t_l^2 . Figure G.13 shows that, as the mixing rate increases
 882 from 0.0 to 1.0, the trained transformer tends to favor predicting false targets. This mirrors the
 883 phenomenon of context hijacking in LLMs.

884 G.6 On the context lengths

885 As alluded in Section 4.4, the memory recall rate is closely related to the KL divergences between
 886 context conditional distributions. Because contexts contain mostly i.i.d samples, longer contexts
 887 imply larger divergences. This is empirically verified in Figure G.14 which demonstrates that longer
 888 context lengths can lead to higher accuracy.

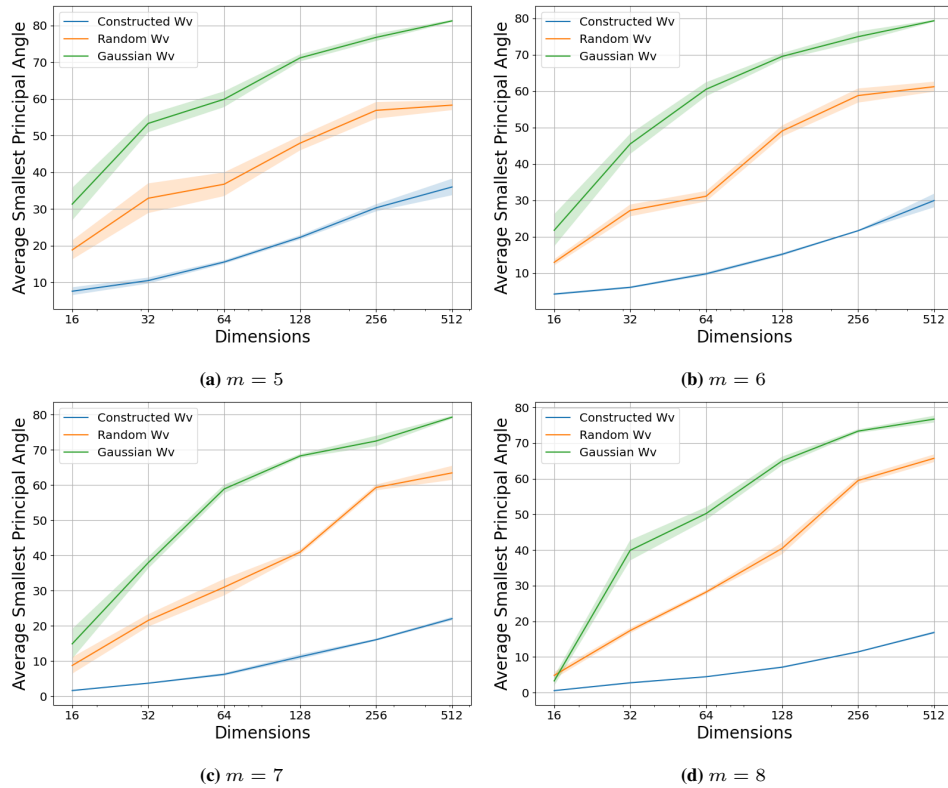


Figure G.3: The constructed value matrix W_V has similar low rank approximation with the trained value matrix. The figure displays average smallest principal angles between low-rank approximations of trained value matrices and those of constructed, randomly constructed, and Gaussian-initialized value matrices. Standard errors are over 5 runs.

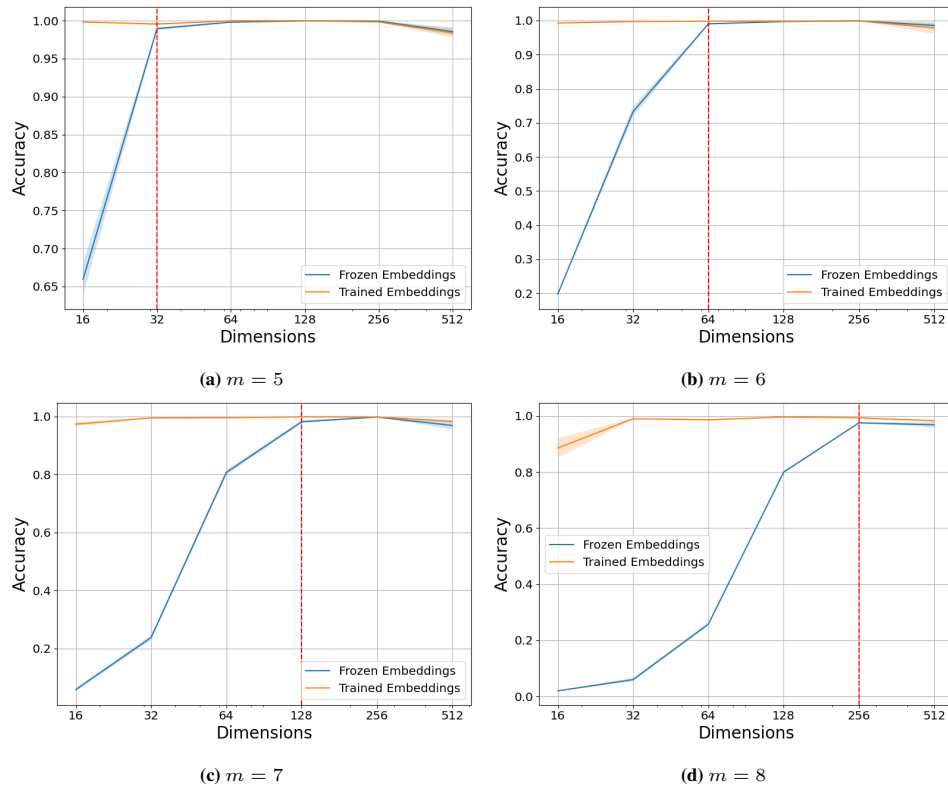


Figure G.4: In the underparameterized regime ($d < V$), freezing embeddings to initializations causes a significant decrease in performance. The graph reports accuracy with different embedding dimensions and the standard errors are over 5 runs. Red lines indicate when $d = V$.

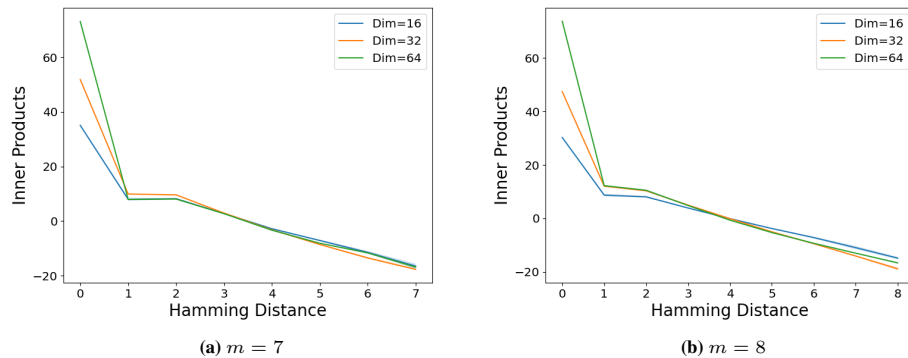


Figure G.5: The relationship between inner products of embeddings and corresponding Hamming distances of tokens can be approximated by (4.2). The graph displays the average inner product between embeddings of two tokens against the corresponding Hamming distance between these tokens. Standard errors are over 5 runs.

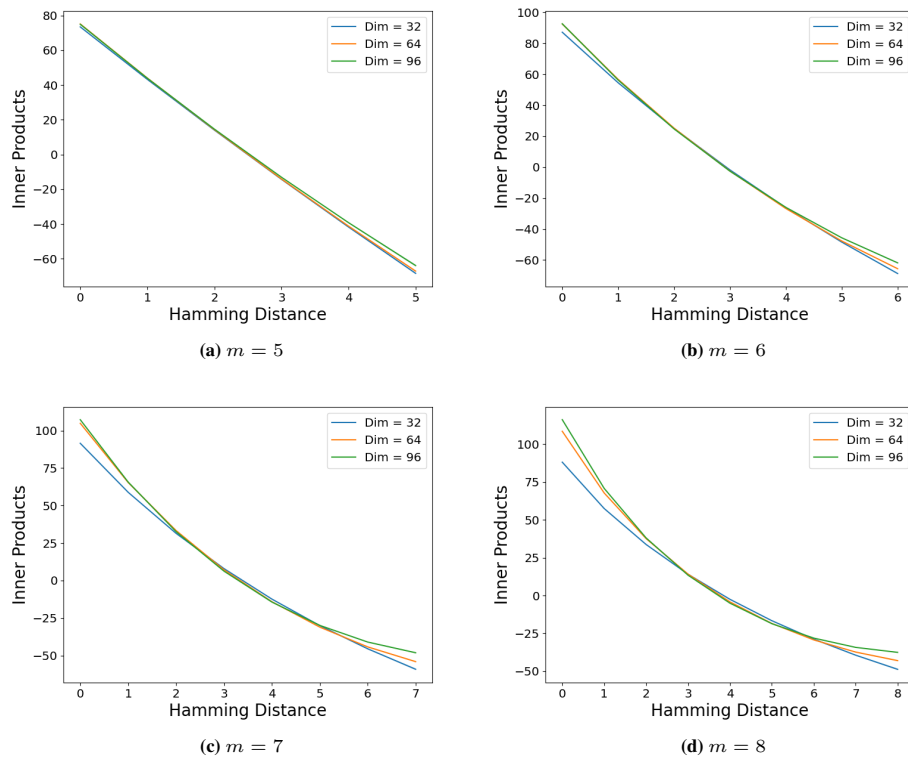


Figure G.6: The relationship between inner products of embeddings and corresponding Hamming distances of tokens is mostly linear when the value matrix W_V is fixed to be the identity. The graph displays the average inner product between embeddings of two tokens against the corresponding Hamming distance between these tokens. Standard errors are over 10 runs.

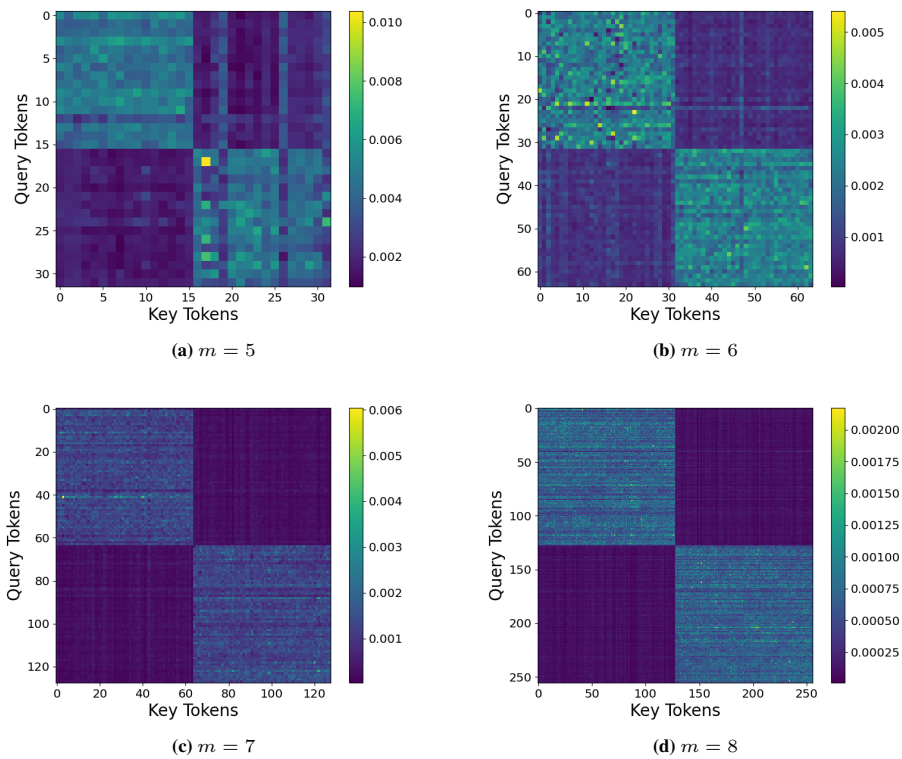


Figure G.7: The attention patterns show the underlying cluster structure of the data generating process. Here, for any latent vector, we have $\mathcal{N}(z^*) = \{z : z_1^* = z_1\} \setminus \{z^*\}$. The figure shows attention score heat maps that are averaged over 10 runs.

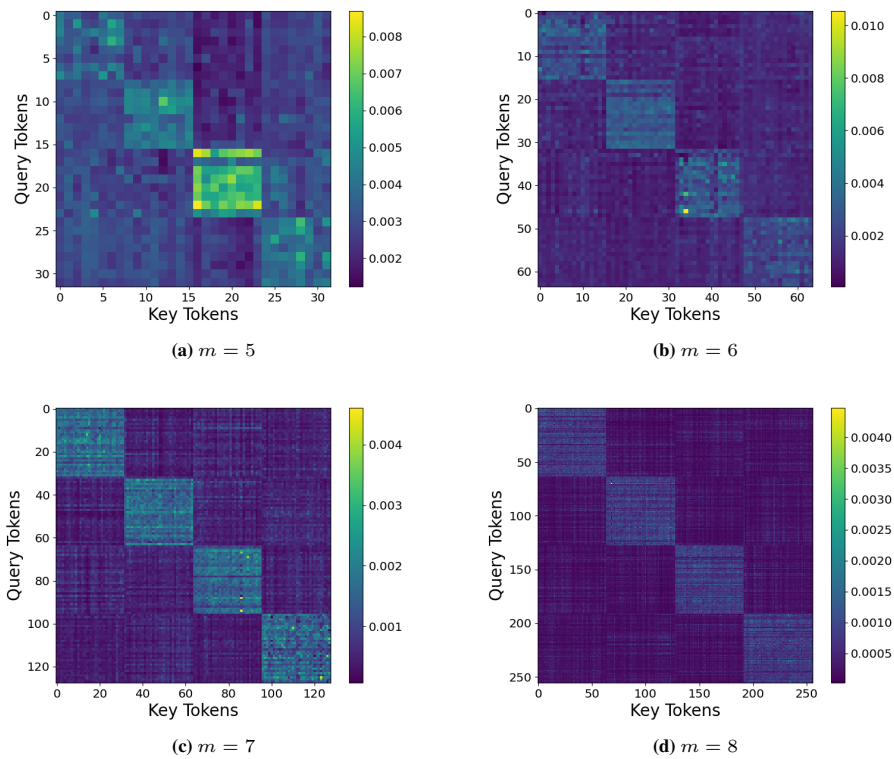
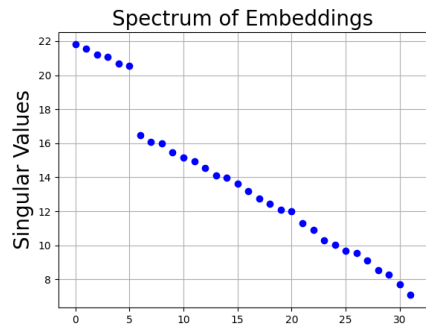
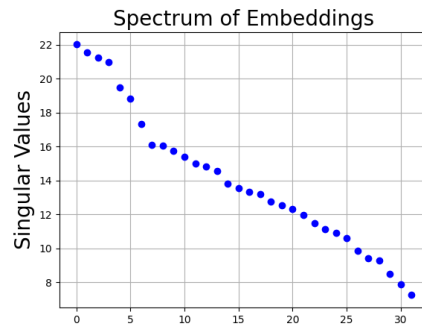


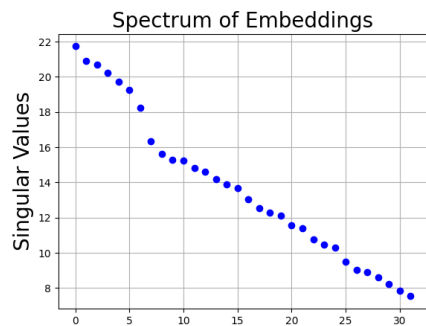
Figure G.8: The attention patterns show the underlying cluster structure of the data generating process. Here, for any latent vector, we have $\mathcal{N}(z^*) = \{z : z_1^* = z_1 \text{ and } z_2^* = z_2\} \setminus \{z^*\}$. The figure shows attention score heat maps that are averaged over 10 runs.



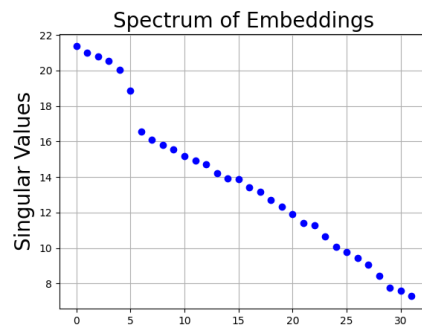
(a) Sample 1



(b) Sample 2

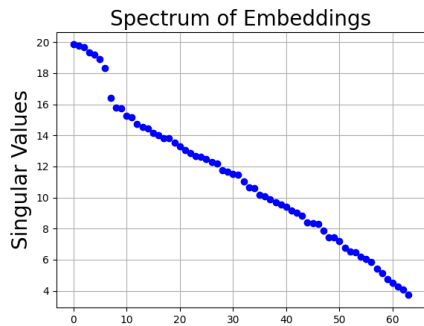


(c) Sample 3

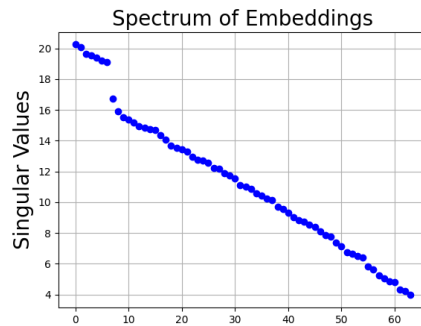


(d) Sample 4

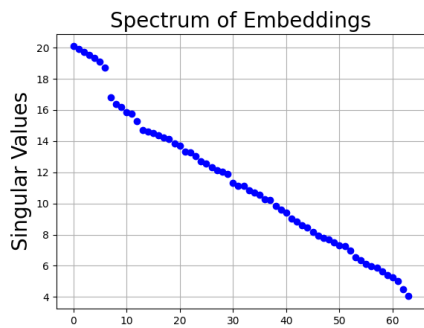
Figure G.9: The spectrum of embedding matrix W_E has eigengaps between the top and bottom eigenvalues, indicating low rank structures. The figure shows results from 4 experimental runs. Number of latent variable m is 7 and the embedding dimension is 32.



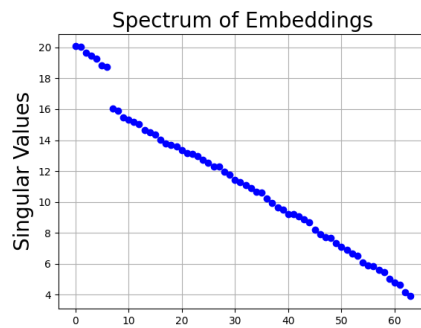
(a) Sample 1



(b) Sample 2

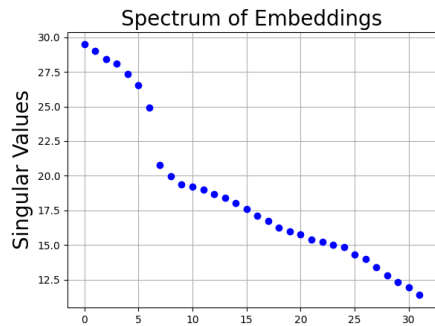


(c) Sample 3

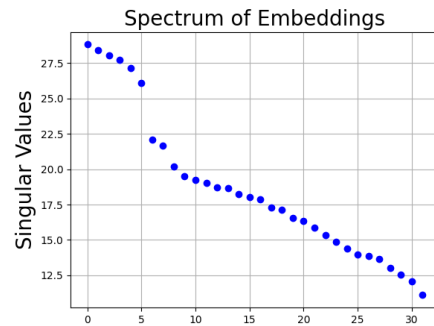


(d) Sample 4

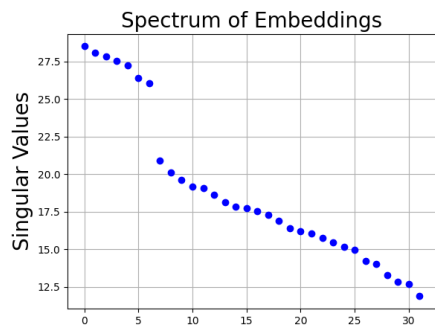
Figure G.10: The spectrum of embedding matrix W_E has eigengaps between the top and bottom eigenvalues, indicating low rank structures. The figure shows results from 4 experimental runs. Number of latent variable m is 7 and the embedding dimension is 64.



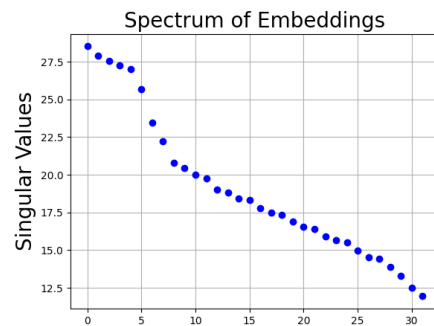
(a) Sample 1



(b) Sample 2

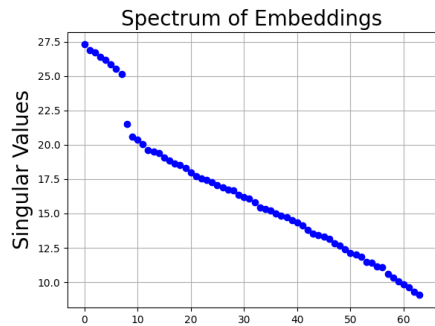


(c) Sample 3

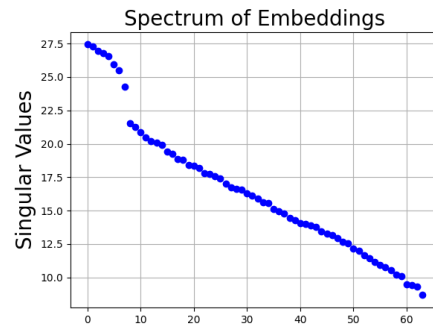


(d) Sample 4

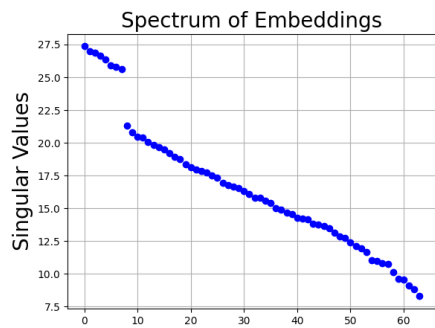
Figure G.11: The spectrum of embedding matrix W_E has eigengaps between the top and bottom eigenvalues, indicating low rank structures. The figure shows results from 4 experimental runs. Number of latent variable m is 8 and the embedding dimension is 32.



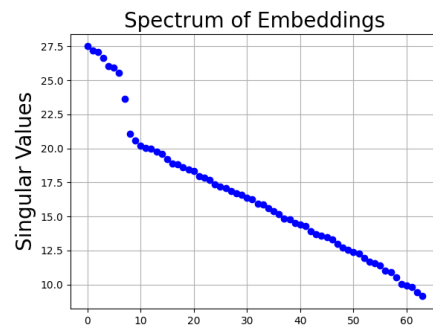
(a) Sample 1



(b) Sample 2

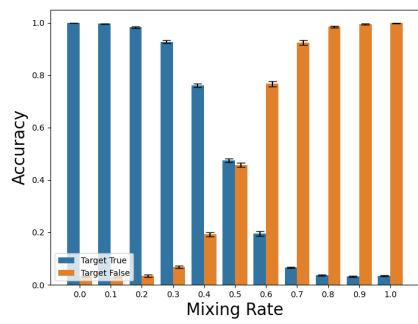


(c) Sample 3

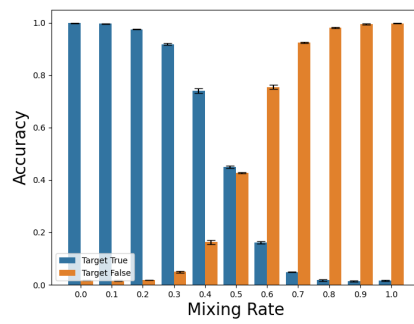


(d) Sample 4

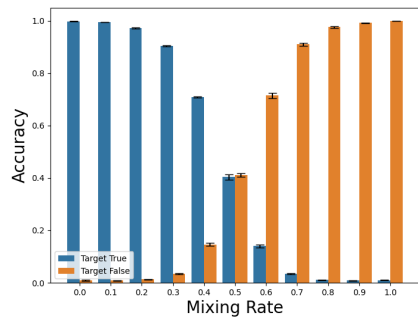
Figure G.12: The spectrum of embedding matrix W_E has eigengaps between the top and bottom eigenvalues, indicating low rank structures. The figure shows results from 4 experimental runs. Number of latent variable m is 8 and the embedding dimension is 64.



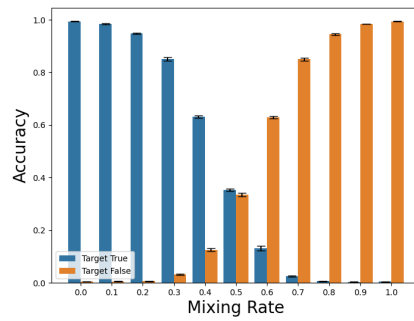
(a) $m = 5$



(b) $m = 6$



(c) $m = 7$



(d) $m = 8$

Figure G.13: Mixing contexts can cause misclassification. The figure reports accuracy for true target and false target under various context mixing rate. Standard errors are over 5 runs.

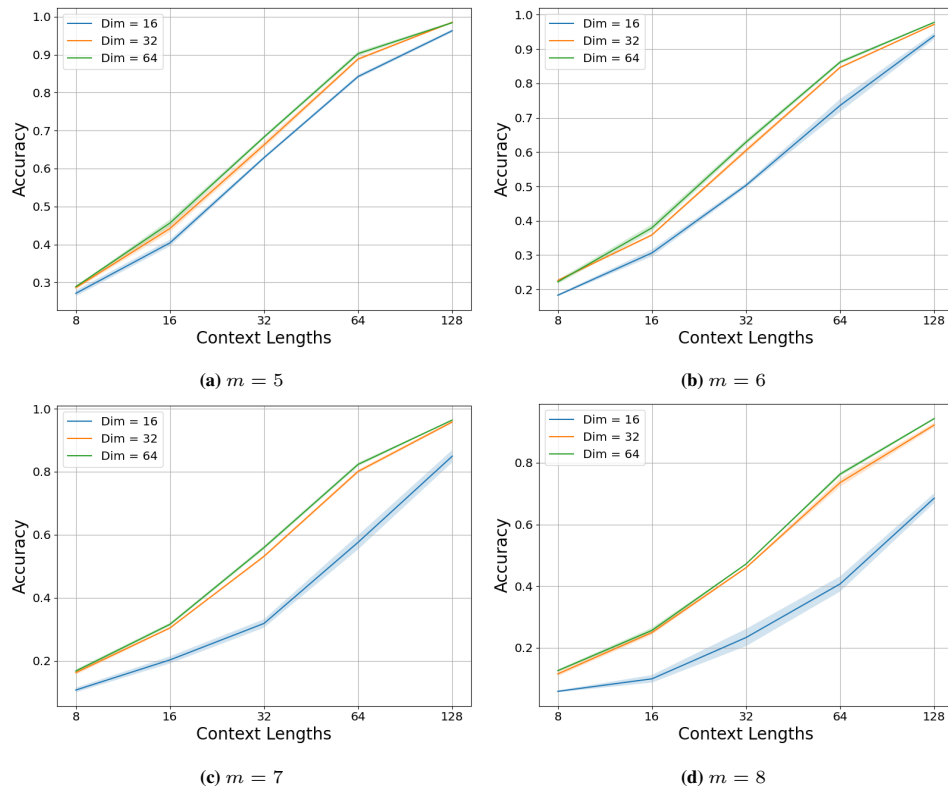


Figure G.14: Increasing context lengths can improve accuracy. The figure reports accuracy across various context lengths and dimensions. Standard errors are over 5 runs.