

Understanding Factual Recall in Transformers via Associative Memories

Eshaan Nichani*
Princeton University

Jason D. Lee
Princeton University

Alberto Bietti
Flatiron Institute

December 10, 2024

Abstract

Large language models have demonstrated an impressive ability to perform factual recall. Prior work has found that transformers trained on factual recall tasks can store information at a rate proportional to their parameter count. In our work, we show that shallow transformers can use a combination of associative memories to obtain such near optimal storage capacity. We begin by proving that the storage capacities of both linear and MLP associative memories scale linearly with parameter count. We next introduce a synthetic factual recall task, and prove that a transformer with a single layer of self-attention followed by an MLP can obtain 100% accuracy on the task whenever either the total number of self-attention parameters or MLP parameters scales (up to log factors) linearly with the number of facts. In particular, the transformer can trade off between using the value matrices or the MLP as an associative memory to store the dataset of facts. We complement these expressivity results with an analysis of the gradient flow trajectory of a simplified linear attention model trained on our factual recall task, where we show that the model exhibits sequential learning behavior.

1 Introduction

One hallmark capability of transformer-based large language models (LLMs) is factual recall [40, 21, 43]. Given a prompt of the form “In what year was George Washington born?” an LLM will correctly respond with “1732.” Language models thus act as databases, storing somewhere in their parameters mappings of the form (George Washington, birth year) \mapsto (1732) which can be easily accessed during inference time.

Prior work [2] has observed that transformers trained on factual recall tasks can store information at a rate proportional to their parameter count. Other studies [e.g., 35, 12, 38, 29] have sought to understand the specific mechanism by which transformers implement factual recall, probing models to understand specifically which transformer blocks “contain” certain facts. However, these studies do not consider the memorization capacity of such constructions, and it is thus an

*eshnich@princeton.edu. This work was conducted while EN was an intern at the Flatiron Institute.

open question to understand how transformers optimally encode such factual information within their weights.

In this work, we show that shallow transformers can use a combination of *associative memories* to obtain near-optimal storage capacity for factual recall tasks. Associative memories store pairs of input-output embeddings through their outer products, and are thus well-suited for modeling the weight matrices of a transformer. Prior work [5] has shown that this associative memory model is a key primitive towards understanding both the representational capacity and optimization dynamics of transformers on synthetic tasks.

Our specific contributions are as follows:

- In Section 3 we begin by studying the ability of linear and MLP associative memory models to store associations between discrete vocabularies. We prove that when the embeddings are sampled randomly over the sphere, these models can store a number of associations proportional to their parameter count, significantly improving over the case where the embeddings are orthogonal.
- In Section 4, we introduce a synthetic next-token prediction task which models factual recall. The data distribution consists of prompts containing a subject token s and relation token r hidden amongst a set of noise tokens, which the learner must map to a ground truth answer $a^*(s, r)$. Our main theorem is that a transformer consisting of a single multi-head self-attention layer followed by an MLP can obtain 100% accuracy when *either* the number of self-attention parameters or MLP parameters scales (up to logs) proportionally with the dataset size.
- In Section 5, we study the gradient descent dynamics of a single linear self-attention head trained on the synthetic task. We prove that the model undergoes a sequential learning dynamics, consisting of a “hallucination” stage where the model outputs the conditional distribution for the answer based on only the relation.
- Finally, in Section 6 we complement our constructions with lower bounds, showing that they are optimal up to logarithmic factors.

Overall, our work makes progress towards understanding the mechanism by which transformers learn and store factual information.

2 Related Work

Associative memories. Associative memories have a long history in the neural computation literature [16, 24, 49]. More recently there has been renewed interest in extensions of such models with larger capacity [25, 8, 28]. These have been linked to the attention blocks in Transformers [42, 44], with [26, 15] in particular using the connection between self-attention and associative memories to design new variants of the attention module. [41] show that overparameterized autoencoders can also behave as associative memories. However, these connections differs from our work, where we consider instead the role of both self-attention and MLP weights as associative memories, in a similar vein to [5, 7].

Memorization and factual recall. Large language models are known to store vast amounts of factual knowledge in their weights [21, 43, 11]. Several recent works in the mechanistic interpretability literature have attempted to understand how transformers store facts [35, 12, 38, 29]. Allen-Zhu and Li [2] empirically studied the memorization capacity for Transformer language models of different sizes trained on synthetic factual recall tasks, and observed near-linear scaling with the number of parameters. Jiang et al. [20] demonstrate how shallow transformers can solve a related latent concept association task by viewing the weight matrices as associative memories. At a more basic level, several works have studied the memorization capacity of neural networks, using constructions that differ from our associative memory approach, both in the context of regression [6, 47, 30] and (next) token prediction [33, 22, 23, 31].

Gradient dynamics. Training dynamics of transformer models on various tasks has been a popular recent line of research [18, 45, 27, 5, 46, 39]. Zhang et al. [50], Mahankali et al. [32] studied training dynamics of transformers with linear attention on in-context learning tasks. Ghosal et al. [13] studied the fine-tuning dynamics on a similar factual recall task, showing how training on lesser-known facts may hurt performance. Our emphasis differs in that we consider non-orthogonal embeddings, and require the model to additionally filter out the relevant subject and relation tokens from the noise tokens, which requires learning of the key and query matrices.

3 Associative Memories

In this section, we show that associative memories have a storage capacity on the order of the number of parameters (up to logarithmic factors), which is near-optimal (as we show in Section 6).

Setup. Our setting follows that of Cabannes et al. [7]. Let $[N]$ be the set of input tokens, and $[M]$ be the set of output tokens. Our goal is to store a set of associations given by the function $f^* : [N] \rightarrow [M]$. For each input token $x \in [N]$ we assign a corresponding embedding vector $\mathbf{e}_x \in \mathbb{R}^d$, and likewise for each output token $y \in [M]$ we associate an unembedding vector $\mathbf{u}_y \in \mathbb{R}^d$. We primarily focus on the setting where the embeddings $\{\mathbf{e}_x\}_{x \in [N]}$ and $\{\mathbf{u}_y\}_{y \in [M]}$ are drawn i.i.d uniformly from the sphere of radius 1. Let $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be our model which “stores” the associations f^* . Given such an F , the prediction $\hat{f}(x)$ for $f^*(x)$ is given by the arg-max decoding $\hat{f}(x) := \arg \max_{y \in [M]} \mathbf{u}_y^\top F(\mathbf{e}_x)$.

3.1 Linear Associative Memories

We first consider the case where F is a linear map $F(\mathbf{e}_x) = \mathbf{W} \mathbf{e}_x$.

Theorem 1. *Assume that f^* is injective. If $d^2 \gtrsim N \text{ poly log } N$, then with high probability over the draw of the embeddings, there exists a \mathbf{W} such that*

$$\arg \max_{y \in [M]} \mathbf{u}_y^\top \mathbf{W} \mathbf{e}_x = f^*(x) \quad \text{for all } x \in [N]. \quad (1)$$

This capacity is obtained by the construction $\mathbf{W} = \sum_{x \in [N]} \mathbf{u}_{f^(x)} \mathbf{e}_x^\top$. Furthermore, if \mathbf{W} is restricted to be a rank m matrix, then such a \mathbf{W} exists when $md \gtrsim N \text{ poly log } N$; this construction is $\mathbf{W} = \sum_{x \in [N]} \mathbf{u}_{f^*(x)} \mathbf{e}_x^\top \sum_{i=1}^m \mathbf{v}_i \mathbf{v}_i^\top$, where $\mathbf{v}_i \in \mathbb{R}^d$ are drawn i.i.d from the standard Gaussian.*

Since \mathbf{W} has d^2 parameters, Theorem 1 shows that the number of associations that can be stored scales (up to log factors) linearly in the number of parameters. We note that in this linear case, the injectivity assumption on f^* is important, as otherwise the capacity may be as low as d , as in [7]. Additionally, we remark that these constructions are easily obtained by gradient descent; the general \mathbf{W} construction corresponds to one-step of GD on the correlation loss $L(\mathbf{W}) = -\sum_x \mathbf{u}_{f^*(x)}^\top \mathbf{W} e_x$, while the low-rank construction corresponds to parameterizing $\mathbf{W} = \mathbf{U}\mathbf{V}^\top$ for $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times m}$, and taking one step of GD on \mathbf{U} while \mathbf{V} is fixed to random Gaussian initialization. The proof of Theorem 1 is deferred to Appendix B.

Remarks. Our setting bears similarity to associative Hopfield networks [16], yet differs in that we decode to a fixed discrete set of output tokens $[M]$ rather than exactly matching the target output. This more closely resembles the language modeling framework, and allows us to improve the memorization capacity from d to d^2 [34]. Next, we note that non-orthogonality of the embeddings is necessary for Theorem 1, as the optimal storage capacity for one-hot embeddings is only $N = d$. Since our constructions are in the regime $N \gg d$, the associative memory \mathbf{W} is a *superposition* [10] of the outer products $\mathbf{u}_{f^*(x)} e_x^\top$. Finally, we remark that the random, rather than trainable, embeddings setting was also studied in Cabannes et al. [7]. The embeddings can be viewed as global quantities in a larger network, of which the associative memory is implementing some subtask, and is thus not able to optimize these embeddings in order to solve its specific task.

3.2 MLP Associative Memories

Next, we consider the case where F is a two-layer neural network with hidden width m ; that is $F(e_x) = \mathbf{V}^\top \sigma(\mathbf{W} e_x)$ for $\mathbf{V}, \mathbf{W} \in \mathbb{R}^{m \times d}$.

Theorem 2 (Informal). *If $md \gtrsim N \text{ poly } \log N$, then with high probability over the draw of the embeddings, there exists \mathbf{V}, \mathbf{W} such that*

$$\arg \max_{y \in [M]} \mathbf{u}_y^\top \mathbf{V}^\top \sigma(\mathbf{W} e_x) = f^*(x) \quad \text{for all } x \in [N]. \quad (2)$$

Since the MLP has $2md$ parameters, Theorem 2 shows that the MLP associative memory scheme has storage capacity which is (nearly) linear in the parameter count.

Proof Sketch. The construction for Theorem 2 mimics that of Theorem 1, after an appropriate random feature transformation. First, sample the rows of \mathbf{W} from the standard Gaussian. Then, set each $v_i = m^{-1} \sum_x \mathbf{u}_{f^*(x)} h_k(\langle \mathbf{w}_i, e_x \rangle)$, where h_k is the k th Hermite polynomial (see Appendix F.1). We then see that

$$F(e_x) = \frac{1}{m} \sum_{i=1}^m \sum_{x' \in [N]} \mathbf{u}_{f^*(x')} h_k(\langle \mathbf{w}_i, e_{x'} \rangle) \sigma(\langle \mathbf{w}_i, e_x \rangle) \approx \sum_{x' \in [N]} \mathbf{u}_{f^*(x')} \langle e_x, e_{x'} \rangle^k. \quad (3)$$

for sufficiently large m . Such *polynomial* associative memory is reminiscent of that in Krotov and Hopfield [25], and can store many more associations for large k . By choosing $k \approx \log_d m$ and appropriately dealing with concentration, one can obtain the $\tilde{O}(md)$ storage capacity (for technical reasons, we must also use the Neural Tangent Kernel [17] rather the random feature model). The full proof of Theorem 2 is deferred to Appendix B.

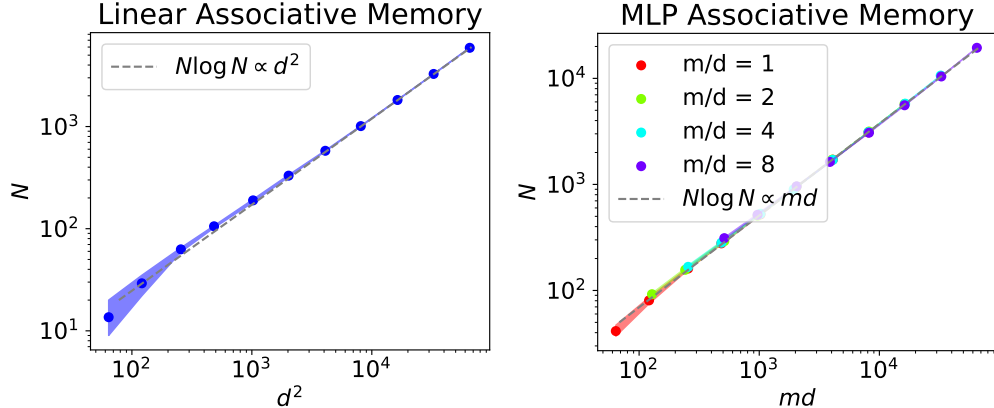


Figure 1: We train linear and MLP associative memories to store the association $f^*(x) = x$. (Left) A linear associative memory requires $d^2 \propto N \log N$ parameters to store N associations. (Right) The MLP associative memory requires $md \propto N \log N$ parameters to store N associations, as predicted by Theorem 2.

On Optimal Storage Capacity. Prior works [6, 30, 31] studying the memorization capacity of neural networks focus on the regression setting, and thus do not directly apply to our setup with multiple outputs and a discrete set of output tokens. Other works [47, 22, 23] show that one can memorize N arbitrary labels with $\tilde{O}(\sqrt{N})$ parameters, at the expense of using a bit complexity of $\tilde{\Omega}(\sqrt{N})$. Such networks still require $\Omega(N)$ bits, which matches our lower bounds in Section 6. These constructions, however, are unwieldy, and are not learnable if we restrict the precision to be poly log N . Instead, our constructions are learnable – the linear construction results from one step of GD, while the ReLU construction uses the NTK and can thus be learned via GD on a convex loss. In Corollary 2, we show that a quantized version of the construction from Theorem 1 indeed succeeds with bit precision $\tilde{O}(1)$, and thus more accurately captures realistic training regimes where models do seem to succeed with low precision [9, 2].

Empirical Validation. In Figure 1, we train both linear and MLP associative memories to store the association $f^*(x) = x$. Given a fixed model size (d, m) , we fit datasets with increasing values of N using the cross entropy loss, and plot the largest value of N for which we can obtain at least 99% accuracy. We observe that the linear associative memory can store $\tilde{\Theta}(d^2)$ associations, while the MLP associative memory can store $\tilde{\Theta}(md)$ associations. See Appendix A for additional experiments where the number of output tokens M does not scale with N .

4 A Synthetic Task for Factual Recall

In this section we introduce a synthetic factual recall task, and show that one-layer transformers constructed via associative memories can store a number of facts proportional to parameter count.

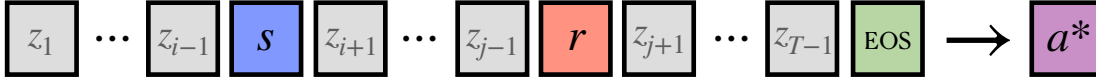


Figure 2: A diagram of the synthetic factual recall task.

4.1 The Task

We first define a global dictionary of facts. Let \mathcal{S} be a set of subject tokens and \mathcal{R} be a set of relation tokens, where $S = |\mathcal{S}|$, $R = |\mathcal{R}|$. Let \mathcal{A} be the set of answer tokens. We let $a^* : \mathcal{S} \times \mathcal{R} \rightarrow \mathcal{A}$ be the ground truth association function, which maps subject-relation tuples (s, r) to their corresponding answer $a^*(s, r)$ ¹. A similar such task was considered in Petroni et al. [40], Ghosal et al. [13].

Define \mathcal{A}_r to be the set of answers corresponding to a relation r , i.e $\mathcal{A}_r := \{a^*(s, r) : s \in \mathcal{S}\}$. Define $\mathcal{A}_s := \{a^*(s, r) : r \in \mathcal{R}\}$ analogously. We assume that each relation corresponds to a disjoint set of answers:

Assumption 1. $\mathcal{A}_r \cap \mathcal{A}_{r'} = \emptyset$ for $r, r' \in \mathcal{R}$ with $r \neq r'$. Furthermore, define $D := \max_{r \in \mathcal{R}} |\mathcal{A}_r|$.

For example, \mathcal{S} could be the set of all countries, while \mathcal{R} could be {president, capital}; in this case, the set of all presidents and set of all capitals are disjoint.

We next define our data distribution \mathcal{D} over sequences. Let $T > 0$ be the context length. Let \mathcal{N} be a set of noise tokens, and define the vocabulary to be $\mathcal{V} := \mathcal{S} \cup \mathcal{R} \cup \mathcal{A} \cup \mathcal{N} \cup \{\text{EOS}\}$, where EOS is a special “end-of-sequence” token. The data distribution is over length $T + 1$ sequences $z_{1:T+1} := (z_1, z_2, \dots, z_T, z_{T+1}) \in \mathcal{V}^{T+1}$, generated via the following procedure:

1. First, sample a subject and relation tuple (s, r) from some distribution p over $\mathcal{S} \times \mathcal{R}$.
2. Next, sample two distinct indices $i, j \in [T - 1]$. Set $z_i = s$ and $z_j = r$.
3. For the remainder of tokens z_k where $k \in [T - 1] \setminus \{i, j\}$, draw z_k uniformly at random from the noise tokens \mathcal{N} .
4. Set $z_T = \text{EOS}$.
5. Finally, set $z_{T+1} = a^*(s, r)$.

The goal of this task is to predict z_{T+1} from (z_1, \dots, z_T) . A model which can successfully do so must (1) be able to isolate the relevant subject and relation tokens from the noise tokens and (2) store all of the associations $(s, r) \mapsto a^*(s, r)$. See Figure 2 for a diagram of the task.

¹We focus on one-to-one relations, where each (s, r) pair corresponds to a unique a^* . This is in contrast to the one-to-many setting, where each (s, r) maps to many possible answers (for example, $s = \text{“France,”}$ $r = \text{“city,”}$ $a^* \in \{\text{“Paris”}, \text{“Toulouse”}, \dots\}$)

4.2 The Model: One-Layer Transformer

Our learner for the task is a single layer of multi-head self attention followed by an MLP. Define d to be the embedding dimension. The input to the transformer is a sequence of vectors $\mathbf{X} := (\mathbf{x}_1, \dots, \mathbf{x}_T)^\top \in \mathbb{R}^{T \times d}$. Each self attention head is parameterized by the key, query, and value matrices $\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V \in \mathbb{R}^{d_h \times d}$, where d_h is the *head dimension*. The self attention head is then a map $\text{attn}(\cdot; \mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V) : \mathbb{R}^{T \times d} \rightarrow \mathbb{R}^{d_h}$, which operates as

$$\text{attn}(\mathbf{X}; \mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V) = \mathbf{W}_V \mathbf{X}^\top \mathcal{S}(\mathbf{X} \mathbf{W}_K^\top \mathbf{W}_Q \mathbf{x}_T), \quad (4)$$

where $\mathcal{S}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$ is the softmax operator.

A multi-head self-attention layer with H heads is parameterized by H different key, query, and value matrices, along with H output matrices. Let $\boldsymbol{\theta} := \{(\mathbf{W}_K^{(h)}, \mathbf{W}_Q^{(h)}, \mathbf{W}_V^{(h)}, \mathbf{W}_O^{(h)})\}_{h \in [H]}$, where $\mathbf{W}_K^{(h)}, \mathbf{W}_Q^{(h)}, \mathbf{W}_V^{(h)}, \mathbf{W}_O^{(h)} \in \mathbb{R}^{d_h \times d}$. A multi-head self-attention layer is then a map $F_{\text{MHSA}}(\cdot; \boldsymbol{\theta}) : \mathbb{R}^{T \times d} \rightarrow \mathbb{R}^d$ given by

$$F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta}) = \sum_{h \in [H]} \mathbf{W}_O^{(h)\top} \text{attn}(\mathbf{X}; \mathbf{W}_K^{(h)}, \mathbf{W}_Q^{(h)}, \mathbf{W}_V^{(h)}). \quad (5)$$

Finally, a single-layer transformer combines a multi-head self-attention layer with an MLP. Let m be the MLP width. Let $\mathbf{V}, \mathbf{W} \in \mathbb{R}^{m \times d}$ be the MLP parameters, and define $\boldsymbol{\theta}_{\text{TF}} := \boldsymbol{\theta} \cup \{\mathbf{V}, \mathbf{W}\}$. Then, a single-layer transformer is the map $F_{\text{TF}}(\cdot; \boldsymbol{\theta}_{\text{TF}}) : \mathbb{R}^{T \times d} \rightarrow \mathbb{R}^d$ given by

$$F_{\text{TF}}(\mathbf{X}; \boldsymbol{\theta}_{\text{TF}}) = F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta}) + \mathbf{V}^\top \sigma(\mathbf{W} F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta})). \quad (6)$$

A single-layer transformer is parameterized by the tuple of hyperparameters (d, H, d_h, m) . The model has $4Hdd_h$ self-attention parameters, and $2md$ MLP parameters.

4.3 One-Layer Transformers have (Almost) Linear Storage Capacity

We next characterize how large a single-layer transformer must be in order to obtain 100% accuracy on the synthetic task. For each token $z \in \mathcal{V}$, sample its embedding vectors $\boldsymbol{\varphi}(z) \in \mathbb{R}^d$ i.i.d uniformly over the sphere of radius 1. An input sequence (z_1, \dots, z_T) gets embedded as $\mathbf{X} = (\boldsymbol{\varphi}(z_1), \dots, \boldsymbol{\varphi}(z_T))^\top$. We use argmax decoding to predict the next token; that is,

$$\hat{f}(z_{1:T}) = \arg \max_{z \in \mathcal{V}} \boldsymbol{\varphi}(z)^\top F_{\text{TF}}(\mathbf{X}; \boldsymbol{\theta}_{\text{TF}}). \quad (7)$$

Our first result is that there exists an attention-only single-layer transformer that obtain 100% accuracy on the factual recall task, as long as the total number of self-attention parameters $4Hdd_h$ scales (up to logarithmic factors) linearly with the dataset size SR .

Theorem 3 (Attention-only, informal). *Assume that $d \geq \tilde{\Omega}(\max(R, D))$ and $Hd_h \geq \tilde{\Omega}(S + R)$. With high probability over the embeddings, there exists a single-layer attention-only transformer $F_{\text{TF}}(\cdot; \boldsymbol{\theta}_{\text{TF}})$ with embedding dimension d , number of heads H and head dimension d_h such that*

$$\mathbb{P}_{z_{1:T+1} \sim \mathcal{D}} \left[\arg \max_{z \in \mathcal{V}} \boldsymbol{\varphi}(z)^\top F_{\text{TF}}(\mathbf{X}; \boldsymbol{\theta}_{\text{TF}}) = z_{T+1} \right] = 1. \quad (8)$$

We next show that a single-layer transformer with an MLP can obtain 100% accuracy on the factual recall task, if the number of MLP parameters md scales linearly with the dataset size:

Theorem 4 (Attention + MLP, informal). *Assume that σ is a polynomial of sufficiently large degree. Define $C(a) = |\{(s, r) : a^*(s, r) = a\}|$. Let (d, H, d_h, m) satisfy*

$$d \geq \tilde{\Omega}(1) \quad Hd_h \geq \tilde{\Omega}(S + R) \quad m \geq \tilde{\Omega}(\max_a C(a)) \quad md \geq \tilde{\Omega}(SR). \quad (9)$$

Then with high probability over the embeddings there exists a single-layer transformer $F_{\text{TF}}(\cdot; \theta_{\text{TF}})$ with embedding dimension d , number of heads H , head dimension d_h , and MLP width m such that

$$\mathbb{P}_{z_{1:T+1} \sim \mathcal{D}} \left[\arg \max_{z \in \mathcal{V}} \varphi(z)^\top F_{\text{TF}}(\mathbf{X}; \theta_{\text{TF}}) = z_{T+1} \right] = 1. \quad (10)$$

The proofs of Theorem 3 and Theorem 4 are deferred to Appendix C.

Remarks. Theorems 3 and 4 each have two main constraints on the size of the architecture needed to obtain 100% accuracy. First, the quantity Hd_h must be larger than $S + R$. This corresponds to self-attention having sufficient capacity to filter out the tokens in $\mathcal{S} \cup \mathcal{R}$ from the noise tokens \mathcal{N} . For the attention-only architecture, we additionally require $d = \tilde{\Omega}(\max(R, D))$. When $R \geq D$, the total number of parameters Hdd_h is (up to logs) at least the total number of facts SR . For the MLP construction, the second condition is that the number of MLP parameters, md , scales nearly linearly with the number of facts SR . As such, as long as either the total number of self-attention parameters or the total number of MLP parameters is large enough, 100% accuracy can be obtained. The single-layer transformer can thus trade off MLP and self-attention parameters while still maintaining perfect accuracy. This phenomenon is reflected in the experiments in Section 4.4. We remark that it is straightforward to extend our construction to the case where we only need to store a size M subset of $\mathcal{S} \times \mathcal{R}$, where the constraints now become $Hdd_h, md = \tilde{\Omega}(M)$.

Proof Sketch. Theorems 3 and 4 both utilize the associative memory framework of Section 3. First, the key and query matrices of each self-attention head act as a *denoiser*, selecting the relevant subject and relation tokens in $z_{1:T}$ while ignoring the noise tokens. To the h th attention head, we associate a subset $\mathcal{S}^{(h)} \subset \mathcal{S} \cup \mathcal{R}$ of subject and relation tokens. Then, setting

$$\mathbf{W}_K^{(h)\top} \mathbf{W}_Q^{(h)} \approx \beta \sum_{z \in \mathcal{S}^{(h)}} \varphi(z) \varphi(\text{EOS})^\top \quad (11)$$

for a large constant β , we see that the h th head will only attend to the tokens in the subset $\mathcal{S}^{(h)}$. We remark that since the embeddings are d -dimensional, at most $d / \text{poly} \log(d)$ embeddings can be in superposition, and thus we must have $|\mathcal{S}^{(h)}| \leq d / \text{poly} \log(d)$.

For the attention-only construction, the output-value matrix $\mathbf{W}_O^{(h)\top} \mathbf{W}_V^{(h)}$ acts as a linear associative memory, mapping each z in $\mathcal{S}^{(h)}$ to a superposition of all possible answers associated with the

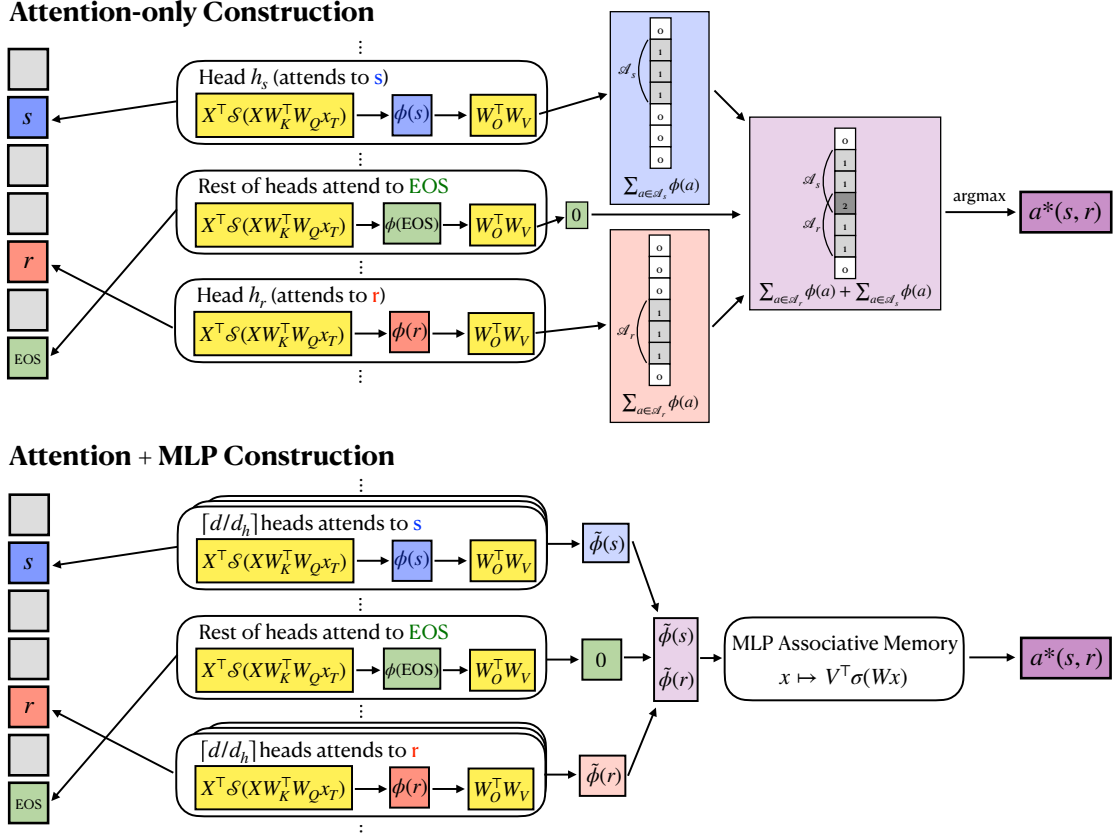


Figure 3: Both the Attention-only and Attention+MLP constructions for the factual recall task.

subject/relation z . Letting P_h be a projection onto a random d_h -dimensional subspace of \mathbb{R}^d , we set

$$\mathbf{W}_O^{(h)\top} \mathbf{W}_V^{(h)} \propto \sum_{z \in \mathcal{S}^{(h)}} \sum_{a \in \mathcal{A}_z} \varphi(a) \varphi(z)^\top P_h. \quad (12)$$

In Lemma 2, we show that this construction stores at most d_h tokens per head (i.e. $|\mathcal{S}^{(h)}| \lesssim d_h$), and requires the dimension to scale with the number of elements in superposition (i.e. $|\mathcal{A}^z| \lesssim d$). Since $|\mathcal{A}^z| \leq R + D$, and the $\mathcal{S}^{(h)}$ partition $\mathcal{S} \cup \mathcal{R}$, it suffices to take $d \gtrsim R + D$ and $Hd_h \gtrsim S + R$.

For the MLP construction, we instead associate the subset $\mathcal{S}^{(h)}$ with $\lceil d/d_h \rceil$ attention heads. This is equivalent to having a single full-rank attention head per subset. We set the aggregate output-value matrix to the identity, so that the output of the self-attention layer is $F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta}) = \varphi(s) + \varphi(r)$. Finally, the MLP layer acts as an MLP associative memory, mapping $\varphi(s) + \varphi(r)$ to $\varphi(a^*(s, r))$ for each (s, r) pair. Via a similar computation to Theorem 2, it suffices to make the total number of parameters md be $md = \tilde{\Omega}(SR)$. Since the $\mathcal{S}^{(h)}$ partition $\mathcal{S} \cup \mathcal{R}$, it suffices to take $Hd_h \gtrsim S + R$ as well. See Figure 3 for a diagram describing both constructions.

4.4 Empirical Validation

We next empirically validate the claims of Theorems 3 and 4 that 100% accuracy can be obtained as long as either the total number of self-attention or MLP parameters scales with SR . We further

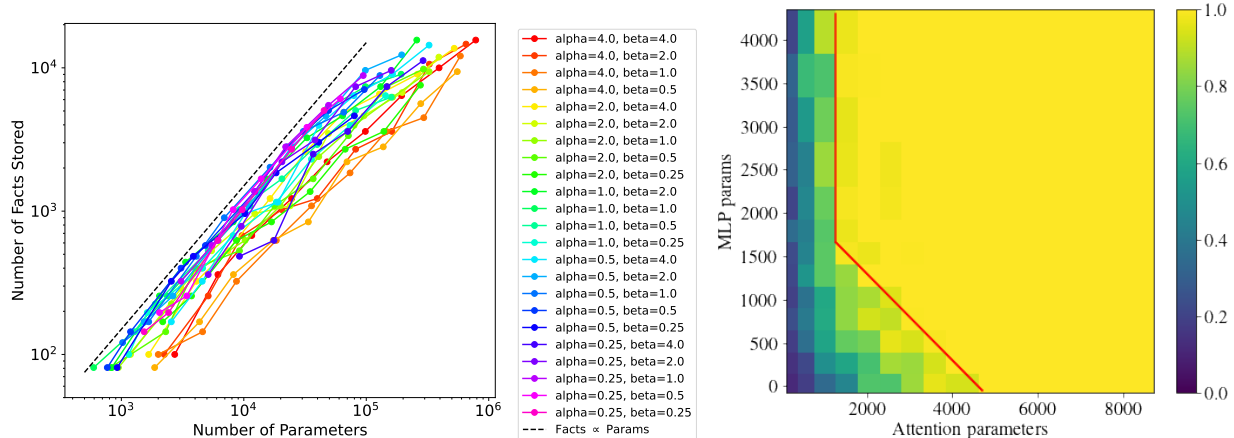


Figure 4: (Left) The number of facts stored scales linearly with the total number of parameters, for a wide range of model sizes. (Right) For a fixed dataset, the model can trade off MLP parameters for attention parameters to obtain 100% accuracy. The heatmap color corresponds to model accuracy.

observe that 100% accuracy can be achieved as long as the *total* number of parameters scales with SR , providing evidence that the model can simultaneously use attention and the MLP to store facts.

In Figure 4, we train a wide range of models of various “shapes” on datasets of varying sizes. A model shape is defined by the tuple $(\alpha, \beta, H)^2$, and corresponds to the family of models satisfying $Hd_h = \alpha d$ and $m = \beta d$. The total number of model parameters is $4Hd_h d + 2md = (4\alpha + 2\beta)d^2$, which can thus be varied by increasing d . For a fixed model size (d, H, d_h, m) , we binary search on the largest dataset size that can be memorized. Specifically, we fix $D = 8$ and vary S, R jointly as $S = R$. Experiments with different scalings are considered in Appendix A. For each (S, R, D) , the fact dataset is generated at random by selecting $|\mathcal{A}| = RD$, $|\mathcal{N}| = S + R$, and for each s sampling $a^*(s, r)$ uniformly at random from $\{(r - 1)D + 1, \dots, rD\}$. We say the dataset was successfully memorized, and as such SR facts were stored, if the model can obtain an accuracy of at least 99%.

On the left panel of Figure 4 we observe that, across different model shapes, the maximum number of facts stored scales linearly with the total number of parameters. On the right panel, we consider a specific dataset with $S = 32, R = 32, D = 8$, and plot the accuracy as the number of parameters vary. We observe that the model can trade off MLP parameters for self-attention parameters, while still maintaining an accuracy of near 1. However, we do still require the total number of attention parameters to be large enough; this corresponds to the $Hd_h = \tilde{\Omega}(S + R)$ constraint.

5 Optimization Dynamics

We next study the optimization dynamics of the factual recall task. To simplify the model, we consider a linear attention transformer (i.e., the softmax is replaced with the identity map) with orthogonal embeddings. We set $d = |\mathcal{V}|$, and let the embedding vectors $\{\varphi(z)\}_{z \in \mathcal{V}}$ satisfy $\langle \varphi(z), \varphi(z') \rangle = \delta_{z=z'}$. Such linear attention or orthogonal embeddings assumptions are common

²The “standard” transformer scaling takes $\alpha = 1, \beta = 4$.

in prior works studying the gradient descent dynamics of transformers [27, 48, 1, 32, 50, 39].

The linear attention model is given by

$$F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) := \mathbf{W}_{OV} \mathbf{X}^\top \mathbf{X} \mathbf{W}_{KQ} \mathbf{x}_T, \quad (13)$$

where we set $d_h = d$ and let $\mathbf{W}_{OV} := \mathbf{W}_O^\top \mathbf{W}_V$, $\mathbf{W}_{KQ} := \mathbf{W}_K^\top \mathbf{W}_Q$ denote the non-factorized output-value and key-query matrices. Let $\hat{p}(\cdot | z_{1:T}) \in \Delta_{\mathcal{A}}$ be the predicted next token distribution on an input sequence $z_{1:T}$, i.e

$$\hat{p}(a | z_{1:T}) := \frac{\exp(\langle \boldsymbol{\varphi}(a), F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) \rangle)}{\sum_{a' \in \mathcal{A}} \exp(\langle \boldsymbol{\varphi}(a'), F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) \rangle)}. \quad (14)$$

One can then rewrite the cross entropy loss as

$$L(\boldsymbol{\theta}) = \mathbb{E}_{z_{1:T+1}}[-\log \hat{p}(z_{T+1} | z_{1:T})]. \quad (15)$$

We would like to characterize the output of running gradient flow, (i.e $\dot{\boldsymbol{\theta}} = -\nabla L(\boldsymbol{\theta})$) with respect to the non-factorized parameters $\boldsymbol{\theta} := \{\mathbf{W}_{OV}, \mathbf{W}_{KQ}\}$ on the cross-entropy loss (15). For notational convenience, we denote $\mathbf{W}_{OV}(a, z) := \boldsymbol{\varphi}(a)^\top \mathbf{W}_{OV} \boldsymbol{\varphi}(z)$, $\mathbf{W}_{KQ}(z) := \boldsymbol{\varphi}(z)^\top \mathbf{W}_{KQ} \boldsymbol{\varphi}(\text{EOS})$, and note that by isometry gradient flow on $\boldsymbol{\theta}$ is equivalent to gradient flow on these quantities.

Let us assume that we start from the following ‘‘balanced’’ initialization.

Assumption 2. *Given an initialization scale $\alpha > 0$, set $\mathbf{W}_{OV}(a, z) = \alpha$ and $\mathbf{W}_{KQ}(z) = \alpha \sqrt{|\mathcal{A}| + 1}$ for each $a \in \mathcal{A}$, $z \in \mathcal{V}$.*

Our first result is that the gradient flow indeed converges to zero loss. As a consequence, the predicted next token probabilities $\hat{p}(z_{T+1} | z_{1:T})$ converge to $\mathbf{1}(z_{T+1} = a^*(s, r))$, where s, r are the subject and relation contained in the sequence $z_{1:T}$.

Theorem 5 (Global Convergence). *For $t \geq 0$, let $\boldsymbol{\theta}(t)$ be the output of running gradient flow for t time. For any $\delta > 0$, there exists a time t_δ such that for $t \geq t_\delta$, $L(\boldsymbol{\theta}(t)) \leq \delta$.*

We next show that the model undergoes a sequential learning dynamics. Let us assume that the number of subjects S is much greater than the number of facts R . We show that during the first stage of training only the $\mathbf{W}_{OV}(a, r)$ and $\mathbf{W}_{KQ}(r)$ components grow for relations $r \in \mathcal{R}$, while the remainder of the parameters stay close to zero. As such, the model gets close to outputting the best predictor based on just the relation token r . Define $p^*(\cdot | r)$ to be the conditional distribution of the answer, given the relation r , i.e $p^*(a | r) := \sum_{s \in \mathcal{S}} p(s | r) \mathbf{1}(a = a^*(s, r))$.

Theorem 6 (Sequential Learning). *Assume that $S \geq 8R\sqrt{2D}$, and $|\mathcal{N}| \geq 4R\sqrt{2DT}$. Let $p(s, r) = \frac{1}{SR}$. Pick $\epsilon > 0$. There exists runtime T^* and initialization scale α (both depending on ϵ) such that:*

1. For all $t \leq T^*$ and $z \in \mathcal{S} \cup \mathcal{N}$, $a \in \mathcal{A}$, we have $|\mathbf{W}_{OV}(a, z)|, |\mathbf{W}_{KQ}(z)| \leq \alpha^{1/2}$
2. There exists $t \leq T^*$ such that, for any input sequence $z_{1:T}$ containing a relation r ,

$$\sum_{a \in \mathcal{A}} (p^*(a | r) - \hat{p}(a | z_{1:T}))^2 \leq \epsilon^2. \quad (16)$$

Proofs of Theorems 5 and 6 are deferred to Appendix D.

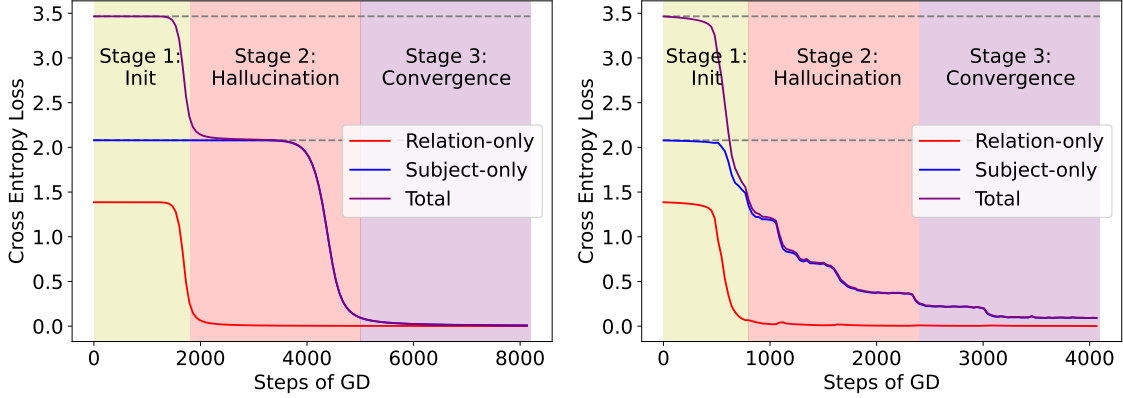


Figure 5: (Left) Loss of the linear attention model with orthogonal embeddings. There is an intermediate *hallucination* stage where the loss plateaus and the model predicts based on only the relation. (Right) Loss of the softmax attention model with random embeddings. We again observe an intermediate hallucination stage, where the relation-only loss is zero but the total loss is still large.

Remarks. Theorem 6 tells us that at some intermediate time, the prediction of the model $\hat{p}(\cdot | z_{1:T})$ is approximately equal to $p^*(\cdot | r)$, the conditional distribution of the answer given the relation r . At this stage, the model ignores all other tokens in the sequence $z_{1:T}$ – including the useful subject token s – and predicts based only on the relation r . For example, if S is the set of all countries and r is the relation “capital,” then on the prompt “What is the capital of France?” the model will output a random countries’ capital. We view this as an instance of *hallucination*: the model is outputting a plausible, yet ultimately incorrect, answer to the prompt. We remark that without the assumption that $S \gg R$, it is possible for this intermediate hallucination stage to exhibit different behavior.

Empirical Validation. We next empirically verify Theorems 5 and 6. We first train the linear attention model with orthogonal embeddings (15) with $S = 16$, $R = 4$ and $D = 8$, and plot the loss over time. In the left pane of Figure 5, we observe three distinct stages. At the start of training, the prediction is close to uniform over all possible answers, and the model obtains a loss of $\log |\mathcal{A}|$. Next, the loss plateaus at $\log D$, and the model outputs the conditional distribution of a given the relation r . Finally, as training continues, the model escapes the plateau and converges to zero loss. We include the “relation-only loss” in the plot, defined as $\mathbb{E}_{z_{1:T+1}}[-\log(\sum_{a \in \mathcal{A}_r} p(a | z_{1:T}))]$, where any probability mass assigned to an answer which is valid for the relation r is considered to be correct; the subject-only loss is defined analogously.

In the right pane of Figure 5, we plot the loss of a single softmax attention head with random embeddings trained on the same factual recall task. We observe similar phenomenology as for linear attention, and identify an intermediate “hallucination” stage where the relation-only loss drops to zero, but the subject-only loss is still far from zero.

6 Lower Bounds

In this section, we argue via information-theoretic arguments that the results from Sections 3 and 4 are optimal up to logarithmic factors. Proofs are deferred to Appendix E.

Associative Memories. Let $[N]$ and $[M]$ be the input and output vocabularies, respectively. To establish a lower bound, we must consider a *distribution* over association functions f^* . For each $x \in [N]$, assume that the output $f^*(x)$ is sampled independently from the uniform distribution over $[M]$. We model the learning protocol as follows. At train time, the learner observes the randomly sampled ground truth f^* , and writes down a B bit model \mathbf{F} . At test time, the learner generates a set of predictions $\hat{f} \in \mathbb{R}^{N \times M}$ from \mathbf{F} , where $\hat{f}(x) \in \Delta_M$ is the prediction for $f^*(x)$. Both the mappings $f^* \rightarrow \mathbf{F}$ and $\mathbf{F} \rightarrow \hat{f}$ can be randomized. Let p be a probability distribution over the input space $[N]$; assume WLOG that $p(1) \geq p(2) \geq \dots \geq p(N)$. The goal of the learner is to minimize the cross entropy loss $L(\hat{f}) = -\sum_{x \in [N]} p(x) \log \hat{f}(x)_{f^*(x)}$.

Theorem 7. *The expected loss of the learner can be lower bounded by*

$$\mathbb{E}_{f^*, \hat{f}} [L(\hat{f})] \geq \log M \cdot \sum_{x \geq \lceil \frac{B}{\log M} \rceil} p(x). \quad (17)$$

We thus see that in order to obtain zero loss, the learner must use $B \geq N \log M$ bits; this matches the construction from Theorem 2 up to \log factors. As a corollary of Theorem 7, we can obtain scaling law lower bounds with respect to model size.

Corollary 1. *Assume that p is a power law, i.e $p(x) \propto x^{-\alpha}$ for $\alpha > 1$. Then $\mathbb{E}_{f^*, \hat{f}} [L(\hat{f})] \gtrsim B^{1-\alpha}$.*

This lower bound is obtained by the MLP associative memory by storing the most probable $\tilde{O}(B)$ associations. This matches the scaling law with respect to model size considered in Michaud et al. [36], Cabannes et al. [7], which also considered storing the $\tilde{O}(B)$ most frequent associations.

Remark. The constructions in Section 3 require storing $\tilde{O}(N)$ network parameters, along with input and output embeddings. We view \mathbf{F} in Theorem 7 as containing only the network parameters, while the embeddings are “global” quantities, independent of the ground truth f^* , used to compute the predictions \hat{f} . This matches our interpretation of the embeddings as fixed global quantities which cannot be modified by the associative memory. We remark that the associative memory constructions from Section 3 match the lower bound, since they hold for $\tilde{O}(1)$ -bit precision (Corollary 2).

Factual Recall. We next prove a lower bound for the factual recall task; a similar bound was proven in Allen-Zhu and Li [2]. Let \mathcal{S} and \mathcal{R} be the fixed set of subjects and relations and \mathcal{V} be the full vocabulary, where $|\mathcal{V}| \gg |\mathcal{S}|, |\mathcal{R}|$. The association function $a^* : \mathcal{S} \times \mathcal{R} \rightarrow \mathcal{V}$ is sampled randomly as follows. First, for each relation $r \in \mathcal{R}$, the answer set \mathcal{A}_r is chosen to be a uniformly random size D subset of \mathcal{V} , conditional on all subsets \mathcal{A}_r being disjoint. For each $s \in \mathcal{S}$, the answer $a^*(s, r)$ is sampled uniformly at random from \mathcal{A}_r . The learner sees the association a^* ,

writes down a B bit model \mathbf{F} , and from \mathbf{F} generates a set of predictions $\hat{f} \in \mathbb{R}^{S \times R \times |\mathcal{V}|}$, where $\hat{f}(s, r) \in \Delta_{\mathcal{V}}$ is the prediction for $a^*(s, r)$. We lower bound \mathbf{L} , the expected cross entropy loss with respect to a distribution $p(s, r)$ over $\mathcal{S} \times \mathcal{R}$, defined as follows:

$$\mathbf{L} := \mathbb{E}_{a^*, \hat{f}} [L(\hat{f})] = \mathbb{E}_{a^*, \hat{f}} \left[- \sum_{s, r} p(s, r) \log \hat{f}(s, r)_{a^*(s, r)} \right]. \quad (18)$$

Theorem 8. Assume that $|\mathcal{V}| \geq 2RD$ and $S \geq CD \log(2D^2 \log |\mathcal{V}|)$ for sufficiently large constant C . There exists a constant $c \in (0, 1)$ such that, if $\mathbf{L} = 0$, the number of bits B must satisfy

$$B \geq SR \log D + (1 - c)RD \log (|\mathcal{V}|/D) \quad (19)$$

We thus see that $\tilde{O}(SR)$ parameters are needed to achieve a loss of zero. For this lower bound, the learner knows the sets \mathcal{S} and \mathcal{R} and does not have to distinguish them from the noise tokens \mathcal{N} , making it a strictly easier problem than the factual recall task in Section 4.

7 Discussion

In this work, we showed that shallow transformers can use associative memories to obtain near optimal storage capacity for factual recall tasks. Furthermore, by studying the optimization dynamics of a simplified model, we also showed that transformers undergo an intermediate hallucination stage. One interesting direction of future work is to better understand the role of the embeddings, and whether there exists an optimal choice of (non-random) embeddings leading to more efficient constructions. Another important direction is to understand the implication of our results towards understanding empirical LLM scaling laws [14]. In particular, does there exist a scaling law lower bound for the factual recall task? Finally, it would be very interesting to understand the extent to which larger models utilize similar associative memory constructions, and if one can probe whether specific “facts” are stored in either the self-attention matrices or the MLP.

Acknowledgements

JDL acknowledges support of NSF CCF 2002272, NSF IIS 2107304, NSF CIF 2212262, ONR Young Investigator Award, and NSF CAREER Award 2144994. This work was conducted in part at the Simons Institute.

References

- [1] Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [2] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity scaling laws. *arXiv preprint arXiv:2404.05405*, 2024.
- [3] William Beckner. Inequalities in fourier analysis. *Annals of Mathematics*, 102(1):159–182, 1975. ISSN 0003486X, 19398980. URL <http://www.jstor.org/stable/1970980>.

- [4] William Beckner. Sobolev inequalities, the poisson semigroup, and analysis on the sphere sn. *Proceedings of the National Academy of Sciences of the United States of America*, 89 (11):4816–4819, 1992. ISSN 00278424, 10916490. URL <http://www.jstor.org/stable/2359537>.
- [5] Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [6] Sébastien Bubeck, Ronen Eldan, Yin Tat Lee, and Dan Mikulincer. Network size and size of the weights in memorization with two-layers neural networks. In *Advances in Neural Information Processing Systems*, 2020.
- [7] Vivien Cabannes, Elvis Dohmatob, and Alberto Bietti. Scaling laws for associative memories. In *International Conference on Learning Representations (ICLR)*, 2024.
- [8] Mete Demircigil, Judith Heusel, Matthias Löwe, Sven Uppgang, and Franck Vermet. On a model of associative memory with huge storage capacity. *Journal of Statistical Physics*, 168: 288–299, 2017.
- [9] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35:30318–30332, 2022.
- [10] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. URL https://transformer-circuits.pub/2022/toy_model/index.html.
- [11] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021.
- [12] Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual associations in auto-regressive language models. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023.
- [13] Gaurav Ghosal, Tatsunori Hashimoto, and Aditi Raghunathan. Understanding finetuning for factual knowledge extraction. *arXiv preprint arXiv:2406.14785*, 2024.
- [14] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [15] Benjamin Hoover, Yuchen Liang, Bao Pham, Rameswar Panda, Hendrik Strobelt, Duen Horng Chau, Mohammed Zaki, and Dmitry Krotov. Energy transformer. *Advances in Neural Information Processing Systems*, 36, 2023.

- [16] John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- [17] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- [18] Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [19] Samy Jelassi, Clara Mohri, David Brandfonbrener, Alex Gu, Nikhil Vyas, Nikhil Anand, David Alvarez-Melis, Yuanzhi Li, Sham M Kakade, and Eran Malach. Mixture of parrots: Experts improve memorization more than reasoning. *arXiv preprint arXiv:2410.19034*, 2024.
- [20] Yibo Jiang, Goutham Rajendran, Pradeep Ravikumar, and Bryon Aragam. Do llms dream of elephants (when told not to)? latent concept association and associative memory in transformers. *arXiv preprint arXiv:2406.18400*, 2024.
- [21] Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8: 423–438, 2020.
- [22] Tokio Kajitsuka and Issei Sato. Are transformers with one layer self-attention using low-rank weight matrices universal approximators? *arXiv preprint arXiv:2307.14023*, 2023.
- [23] Tokio Kajitsuka and Issei Sato. Optimal memorization capacity of transformers. *arXiv preprint arXiv:2409.17677*, 2024.
- [24] Teuvo Kohonen. Correlation matrix memories. *IEEE Transactions on Computers*, 1972.
- [25] Dmitry Krotov and John J Hopfield. Dense associative memory for pattern recognition. *Advances in neural information processing systems*, 29, 2016.
- [26] Hung Le, Truyen Tran, and Svetha Venkatesh. Self-attentive associative memory. In *International conference on machine learning*, pages 5682–5691. PMLR, 2020.
- [27] Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do transformers learn topic structure: Towards a mechanistic understanding. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2023.
- [28] Carlo Lucibello and Marc Mézard. Exponential capacity of dense associative memories. *Physical Review Letters*, 132(7):077301, 2024.
- [29] Ang Lv, Kaiyi Zhang, Yuhan Chen, Yulong Wang, Lifeng Liu, Ji-Rong Wen, Jian Xie, and Rui Yan. Interpreting key mechanisms of factual recall in transformer-based language models. *arXiv preprint arXiv:2403.19521*, 2024.
- [30] Liam Madden and Christos Thrampoulidis. Memory capacity of two layer neural networks with smooth activations. *SIAM Journal on Mathematics of Data Science*, 6(3):679–702, 2024.

- [31] Liam Madden, Curtis Fox, and Christos Thrampoulidis. Upper and lower memory capacity bounds of transformers for next-token prediction. *arXiv preprint arXiv:2405.13718*, 2024.
- [32] Arvind Mahankali, Tatsunori B Hashimoto, and Tengyu Ma. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
- [33] Sadegh Mahdavi, Renjie Liao, and Christos Thrampoulidis. Memorization capacity of multi-head attention in transformers. *arXiv preprint arXiv:2306.02010*, 2023.
- [34] R. McEliece, E. Posner, E. Rodemich, and S. Venkatesh. The capacity of the hopfield associative memory. *IEEE Transactions on Information Theory*, 33(4):461–482, 1987. doi: 10.1109/TIT.1987.1057328.
- [35] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [36] Eric Michaud, Ziming Liu, Uzey Girit, and Max Tegmark. The quantization model of neural scaling. *Advances in Neural Information Processing Systems*, 36, 2023.
- [37] Ashley Montanaro. Some applications of hypercontractive inequalities in quantum information theory. *Journal of Mathematical Physics*, 53(12), 2012.
- [38] Neel Nanda, S Rajamanoharan, J Kramár, and R Shah. Fact finding: Attempting to reverse-engineer factual recall on the neuron level. *AI Alignment Forum*, 2023. URL <https://www.alignmentforum.org/posts/iGuwZTHWb6DFY3sKB/fact-finding-attempting-to-reverse-engineer-factual-recall>.
- [39] Eshaan Nichani, Alex Damian, and Jason D Lee. How transformers learn causal structure with gradient descent. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2024.
- [40] Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. Language models as knowledge bases? In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2019.
- [41] Adityanarayanan Radhakrishnan, Mikhail Belkin, and Caroline Uhler. Overparameterized neural networks implement associative memory. *Proceedings of the National Academy of Sciences*, 117(44):27162–27170, 2020.
- [42] Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Thomas Adler, Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Sandve, et al. Hopfield networks is all you need. *arXiv preprint arXiv:2008.02217*, 2020.
- [43] Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model? In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020.

- [44] Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight programmers. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [45] Charlie Snell, Ruiqi Zhong, Dan Klein, and Jacob Steinhardt. Approximating how single head attention learns. *arXiv preprint arXiv:2103.07601*, 2021.
- [46] Yuandong Tian, Yiping Wang, Beidi Chen, and Simon Du. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [47] Gal Vardi, Gilad Yehudai, and Ohad Shamir. On the optimal memorization power of relu neural networks. *arXiv preprint arXiv:2110.03187*, 2021.
- [48] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023.
- [49] David J Willshaw, O Peter Buneman, and Hugh Christopher Longuet-Higgins. Non-holographic associative memory. *Nature*, 222(5197):960–962, 1969.
- [50] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024.

A Additional Experiments

A.1 MLP Associative Memory

In Figure 6, we train MLP associative memories to store the association $f^*(x) = x \bmod M$. We fix $M = 32$ throughout. In this case, we see that the number of associations N which can be stored by a model with md parameters scales as $N \propto md$. This linear scaling in the absence of logarithmic factors is due to the fact that the number of output tokens M is a constant, and does not scale with the number of input tokens N .

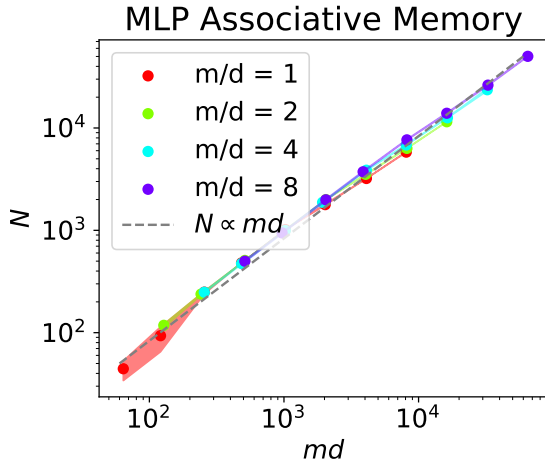


Figure 6: We train an MLP associative memory to store the association $f^*(x) = x \bmod 32$. Empirically, $md \propto N$ parameters are required to store N associations.

A.2 Factual Recall

In Figures 7 to 9, we repeat the experiment in the left pane of Figure 4, for different choices of H and scalings of (S, R, D) . In all plots, we observe the general trend that the number of facts stored scales proportionally to the number of model parameters.

A.3 Experimental Details

Figures 1, 6: For the MLP associative memory experiments, for each choice of m, d, N , we first sample random embeddings $\{e_x\}_{x \in [N]}$, $\{u_y\}_{y \in [M]}$ i.i.d uniformly over the sphere. We train a two-layer neural network on the cross entropy loss to predict the association $f^*(x) = x$. We use standard parameterization and initialization, and the activation $\sigma = \text{ReLU}$. The network is trained using ADAM with a learning rate of 10^{-2} for 2^{14} steps. We compute the maximum accuracy the network achieves over the training run, and say that the network has “stored” the dataset if the highest accuracy is at least 99%. We repeat this procedure to binary search over N , to find the largest value of N such that the network achieves an accuracy of at least 99%. Error bars are shown for 5 random seeds.

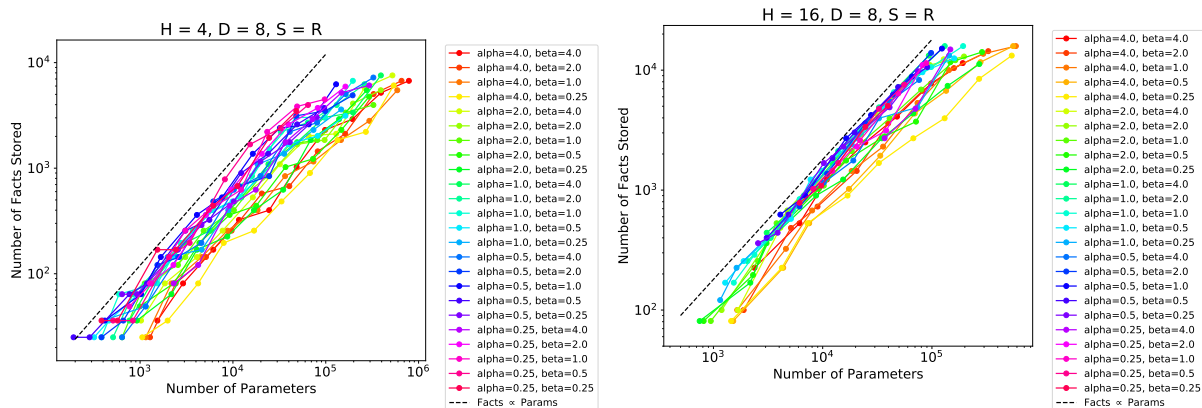


Figure 7: We repeat the experiment in Figure 4, varying the number of head to be 4 (Left) or 16 (Right). In both cases, we observe that the number of facts stored scales linear with the parameter count.

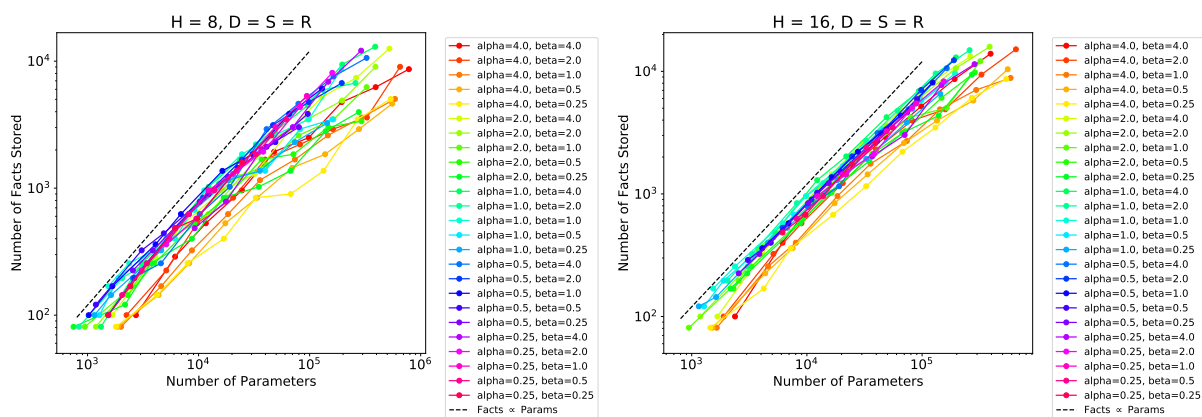


Figure 8: We repeat Figure 4 on factual recall tasks where each subject and relation map to a distinct answer (i.e $D = S$).

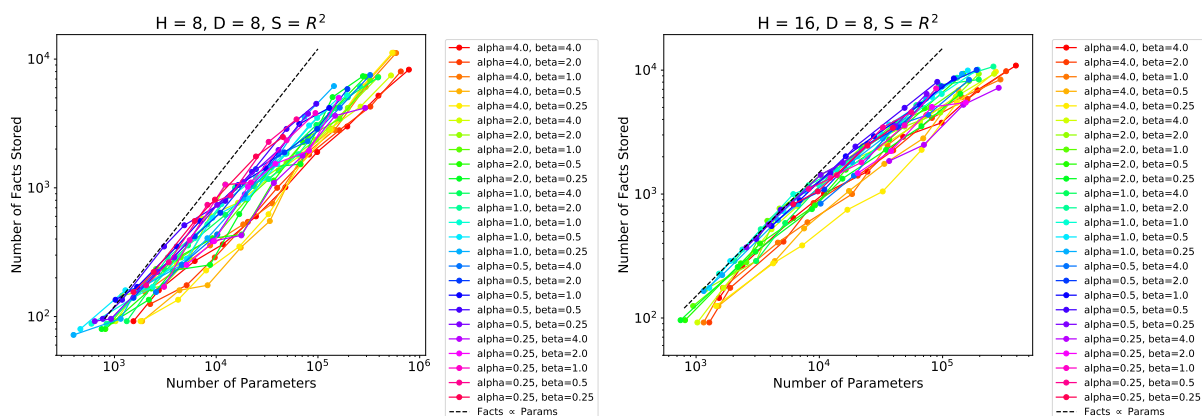


Figure 9: We repeat Figure 4 on factual recall tasks where the number of subjects is much larger than the number of relations; specifically, we take $S = R^2$.

Figures 4, 7, 8, 9: We consider a fixed prompt length of $T = 32$, and train the models via online batch gradient descent with batch size 1024 on the population loss (i.e we sample an independent batch at each timestep). We use standard parameterization and initialization for both self-attention and the MLP. For a fixed model size, we binary search over the maximum value of SR such that the model achieves an accuracy of at least 99%. All models were trained using ADAM for 2^{14} steps, with a sweep over learning rates in $\{.001, .003, .01\}$ (where we consider the best performing model over all learning rates).

Figure 5: In the left pane we train a linear attention head with orthogonal embeddings. The weights are all initialized to be equal to 10^{-5} . In the right plot, we train a softmax attention head with random embeddings, which are fixed throughout training.

B Proofs for Section 3

Proof of Theorem 1. Let us set $\mathbf{W} = \sum_{z \in [N]} \mathbf{u}_{f^*(z)} \mathbf{e}_z^\top$. For $y \neq f^*(x)$, define the quantity γ_{xy} by

$$\gamma_{xy} = (\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top \mathbf{W} \mathbf{e}_x.$$

We first see that (where the expectation is taken over the randomness of the embedding vectors)

$$\begin{aligned} \mathbb{E}[\gamma_{xy}] &= \sum_{z \in [N]} \mathbb{E}[(\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top \mathbf{u}_{f^*(z)}] \mathbb{E}[\mathbf{e}_z^\top \mathbf{e}_x] \\ &= \mathbb{E}[(\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top \mathbf{u}_{f^*(x)}] \\ &= 1. \end{aligned}$$

We can next compute the second moment of γ_{xy} . Since the \mathbf{e}_z are drawn uniformly on the sphere, the $\mathbf{e}_z^\top \mathbf{e}_x$ terms for $z \neq x$ are independent and mean zero. Therefore

$$\begin{aligned} \mathbb{E}[\gamma_{xy}^2] &= \sum_{z \in [N]} \mathbb{E}[(\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top \mathbf{u}_{f^*(z)}]^2 \mathbb{E}[(\mathbf{e}_z^\top \mathbf{e}_x)^2] \\ &= \mathbb{E}[(1 - \mathbf{u}_{f^*(x)}^\top \mathbf{u}_y)^2] + \frac{1}{d} \sum_{z \neq x} \mathbb{E}[(\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top \mathbf{u}_{f^*(z)}]^2 \\ &= 1 + \frac{1}{d} + \frac{1}{d} \sum_{z \neq x} \left(\frac{2}{d} \cdot \mathbf{1}(f^*(z) \neq y) + \left(1 + \frac{1}{d}\right) \cdot \mathbf{1}(f^*(z) = y) \right) \\ &= \left(1 + \frac{1}{d}\right)^2 + \frac{2(N-2)}{d^2}. \end{aligned}$$

Therefore

$$\text{Var}(\gamma_{xy}) = \frac{2}{d} + \frac{1}{d^2} + \frac{2(N-2)}{d^2} \lesssim \frac{1}{d} + \frac{N}{d^2}.$$

Let δ be a fixed failure probability, and let $\delta' = \frac{\delta}{NM}$. Observe that γ_{xy} is a degree 4 polynomial. By Lemma 14, by choosing $d \geq C \log^4(1/\delta')$ and $d^2 \geq CN \log^4(1/\delta')$ for a sufficiently large constant C , we have that $\frac{16e^{-1} \log^4(1/\delta) \text{Var}(\gamma_{xy})}{(\mathbb{E}\gamma_{xy})^2} \leq 1$, and thus $\mathbb{P}(\gamma_{xy} \leq 0) \leq \delta'$.

Therefore union bounding over all (x, y) pairs with $y \neq f^*(x)$, we have that

$$\mathbb{P}(\exists \gamma_{xy} \leq 0) \leq N(M-1)\delta' \leq \delta.$$

Thus with probability $1-\delta$, $\gamma_{xy} > 0$ for all x and $y \neq f^*(x)$, and on this event $\arg \max_{y \in [M]} \mathbf{u}_y^\top \mathbf{W} \mathbf{e}_x = f^*(x)$ for all $x \in [N]$. □

For Theorem 2, we need the following assumption on the activation σ .

Assumption 3. σ is a polynomial of degree q . Furthermore, if $\sigma(z) = \sum_{k=0}^q c_k h_k(z)$ is the Hermite decomposition of σ , then $c_k \neq 0$ for all $0 \leq k \leq q$.

We prove the following formal version of Theorem 2.

Theorem 9. Let $\epsilon \in (0, 1)$ be a fixed constant. Assume that $d \geq N^\epsilon$ and $N \geq C_1(\epsilon)$, where $C_1(\epsilon)$ is a constant depending only on ϵ . Assume that q in Assumption 3 satisfies $q = \frac{C_2}{\epsilon}$ for some $C_2 > 2$. Then, if $md \gtrsim N(C_3 \log(MN/\delta))^{C_4/\epsilon}$, with probability $1 - \delta$ over the draw of the embeddings, there exists \mathbf{V}, \mathbf{W} such that

$$\arg \max_{y \in [M]} \mathbf{u}_y^\top \mathbf{V}^\top \sigma(\mathbf{W} \mathbf{e}_x) = f^*(x) \quad (20)$$

for all $x \in [N]$.

Proof of Theorem 2. Let us consider the linearization, or Neural Tangent Kernel, of F :

$$F_{\text{NTK}}(\mathbf{z}) = \mathbf{V}^\top (\sigma'(\mathbf{W}^0 \mathbf{z}) \odot (\mathbf{W} - \mathbf{W}^0) \mathbf{z}) = \sum_{i \in [m]} \mathbf{v}_i \sigma'(\langle \mathbf{w}_i^0, \mathbf{z} \rangle) \langle \mathbf{w}_i - \mathbf{w}_i^0, \mathbf{z} \rangle.$$

where \mathbf{W}^0 is the initialization which we are linearizing with respect to. By rescaling the parameters as $\mathbf{W} - \mathbf{W}^0 \leftarrow \epsilon(\mathbf{W} - \mathbf{W}^0)$ and $\mathbf{V} \leftarrow \epsilon^{-1} \mathbf{V}$, we see that $F \rightarrow F_{\text{NTK}}$ as $\epsilon \rightarrow 0$. It thus suffices to work with F_{NTK} instead of F . For ease of notation, we redefine \mathbf{W}^0 as \mathbf{W} , and $\mathbf{W} - \mathbf{W}^0$ as \mathbf{Q} , so that

$$F(\mathbf{z}) = \sum_{i \in [m]} \mathbf{v}_i \sigma'(\langle \mathbf{w}_i, \mathbf{z} \rangle) \langle \mathbf{q}_i, \mathbf{z} \rangle$$

Let k be a even integer, to be chosen later. Assume without loss of generality that $c_{k+1} > 0$ (if it is negative, we can simply negate all the \mathbf{q}_i in the construction below). Set

$$\mathbf{q}_i = \frac{1}{m} \sum_{z \in [N]} h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle) \langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle \mathbf{e}_z,$$

where h_k is the k th Hermite polynomial. Then

$$F(\mathbf{e}_x) = \frac{1}{m} \sum_{i \in [m], z \in [N]} \mathbf{v}_i \langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle \sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle) h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle) \langle \mathbf{e}_x, \mathbf{e}_z \rangle$$

As in the proof of Theorem 1, define the margin between x and some $y \neq f^*(x)$ as

$$\begin{aligned}\gamma_{xy} &= (\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top F(\mathbf{e}_x) \\ &= \frac{1}{m} \sum_{i \in [m], z \in [N]} \langle \mathbf{v}_i, \mathbf{u}_{f^*(x)} - \mathbf{u}_y \rangle \langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle \sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle) h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle) \langle \mathbf{e}_x, \mathbf{e}_z \rangle.\end{aligned}$$

We will show that, with high probability over the draw of the embeddings over the sphere, *and* the $\mathbf{v}_i, \mathbf{w}_i$ independently from the standard Gaussian, that $\gamma_{xy} > 0$ for all $y \neq f^*(x)$.

The expectation of the margin is

$$\begin{aligned}\mathbb{E}[\gamma_{xy}] &= \sum_z \mathbb{E}[\langle \mathbf{v}_i, \mathbf{u}_{f^*(x)} - \mathbf{u}_y \rangle \langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle \sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle) h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle) \langle \mathbf{e}_x, \mathbf{e}_z \rangle] \\ &= c_{k+1} \sum_z \mathbb{E}[\langle \mathbf{u}_{f^*(x)} - \mathbf{u}_y, \mathbf{u}_{f^*(z)} \rangle \langle \mathbf{e}_z, \mathbf{e}_x \rangle^{k+1}] \\ &= c_{k+1}.\end{aligned}$$

We next compute the variance. Define ω_{iz}^{xy} as

$$\omega_{iz}^{xy} = \langle \mathbf{v}_i, \mathbf{u}_{f^*(x)} - \mathbf{u}_y \rangle \langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle \sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle) h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle) \langle \mathbf{e}_x, \mathbf{e}_z \rangle,$$

so that $\gamma_{xy} = \frac{1}{m} \sum_{i,z} \omega_{iz}^{xy}$. First, observe that when $z \neq z'$, we have $\mathbb{E}[\omega_{iz}^{xy} \omega_{jz'}^{xy}] = 0$, since k is even. For $i \neq j$, we have that

$$\mathbb{E}[\omega_{iz}^{xy} \omega_{jz}^{xy}] = c_{k+1}^2 \mathbb{E}[\langle \mathbf{u}_{f^*(x)} - \mathbf{u}_y, \mathbf{u}_{f^*(z)} \rangle^2] \mathbb{E}[\langle \mathbf{e}_z, \mathbf{e}_x \rangle^{2(k+1)}]$$

First, by Lemma 12 we have that

$$\mathbb{E}[\langle \mathbf{e}_z, \mathbf{e}_x \rangle^{2(k+1)}] \leq \begin{cases} 1 & x = z \\ (2k+2)^{k+1} d^{-(k+1)} & x \neq z \end{cases}.$$

Next, we see that

$$\mathbb{E}[\langle \mathbf{u}_{f^*(x)} - \mathbf{u}_y, \mathbf{u}_{f^*(z)} \rangle^2] = \begin{cases} 1 + \frac{1}{d} & f^*(x) = f^*(z) \text{ or } y = f^*(z) \\ \frac{2}{d} & \text{otherwise} \end{cases}.$$

Finally, we have that

$$\begin{aligned}\mathbb{E}[\omega_{iz}^{xy} \omega_{iz}^{xy}] &= \mathbb{E}[\langle \mathbf{v}_i, \mathbf{u}_{f^*(x)} - \mathbf{u}_y \rangle^2 \langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle^2 \sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle)^2 h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle)^2 \langle \mathbf{e}_x, \mathbf{e}_z \rangle^2] \\ &= \mathbb{E}[\langle \mathbf{v}_i, \mathbf{u}_{f^*(x)} - \mathbf{u}_y \rangle^2 \langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle^2] \mathbb{E}[\sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle)^2 h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle)^2 \langle \mathbf{e}_x, \mathbf{e}_z \rangle^2].\end{aligned}$$

The first quantity can be bounded as

$$\begin{aligned}\mathbb{E}[\langle \mathbf{v}_i, \mathbf{u}_{f^*(x)} - \mathbf{u}_y \rangle^2 \langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle^2] &\leq \mathbb{E}[\langle \mathbf{v}_i, \mathbf{u}_{f^*(x)} - \mathbf{u}_y \rangle^4]^{1/2} \mathbb{E}[\langle \mathbf{v}_i, \mathbf{u}_{f^*(z)} \rangle^4]^{1/2} \\ &\leq 2 \cdot \sqrt{3} \cdot \sqrt{3} = 6.\end{aligned}$$

The second term is bounded as

$$\mathbb{E}[\sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle)^2 h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle)^2 \langle \mathbf{e}_x, \mathbf{e}_z \rangle^2] \leq \mathbb{E}[\sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle)^8]^{1/4} \mathbb{E}[h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle)^8]^{1/4} \mathbb{E}[\langle \mathbf{e}_x, \mathbf{e}_z \rangle^4]^{1/2}$$

By Gaussian hypercontractivity (Lemma 13),

$$\mathbb{E}[\sigma'(\langle \mathbf{w}_i, \mathbf{e}_x \rangle)^8]^{1/4} = \|\sigma'\|_{L^8}^2 \leq 8^q \|\sigma'\|_{L^2}^2 \lesssim 8^q,$$

and likewise

$$\mathbb{E}[h_k(\langle \mathbf{e}_z, \mathbf{w}_i \rangle)^8]^{1/4} \leq 8^k.$$

Finally, $\mathbb{E}[\langle \mathbf{e}_x, \mathbf{e}_z \rangle^4]^{1/2} \leq 4d^{-1}$ if $x \neq z$, and 1 otherwise. Altogether,

$$\mathbb{E}[\omega_{iz}^{xy} \omega_{iz}^{xy}] \lesssim 2^{3q+3k} (d^{-1} + \mathbf{1}(x = z)).$$

Altogether, we get that

$$\mathbb{E}[\gamma_{xy}^2] = \frac{m-1}{m} \sum_z \mathbb{E}[\omega_{iz}^{xy} \omega_{jz}^{xy}] + \frac{1}{m} \sum_z \mathbb{E}[\omega_{iz}^{xy} \omega_{iz}^{xy}].$$

The first quantity is

$$\sum_z \mathbb{E}[\omega_{iz}^{xy} \omega_{jz}^{xy}] \leq \left(1 + \frac{1}{d}\right) c_{k+1}^2 + c_{k+1}^2 (2k+2)^{k+1} d^{-(k+1)} \cdot 2N.$$

The second quantity is

$$\sum_z \mathbb{E}[\omega_{iz}^{xy} \omega_{iz}^{xy}] \leq 2^{3q+3k} \left(1 + \frac{N}{d}\right).$$

Therefore

$$\text{Var}(\gamma_{xy}) \lesssim \frac{1}{d} + \frac{(2k)^{k+1} N}{d^{k+1}} + \frac{2^{3q+3k} N}{md}.$$

Choose $k = 2\lceil \frac{1}{\epsilon} \rceil$; then

$$\frac{d^k}{(2k)^{k+1} N} \geq \frac{N}{(4/\epsilon)^{1+2/\epsilon}} \geq 1$$

for $N \geq C_1(\epsilon)$, and so

$$\text{Var}(\gamma_{xy}) \lesssim \frac{1}{d} + \frac{2^{3q+3k} N}{md}.$$

Observe that γ_{xy} is a degree $2q + 2k + 4 \leq 4q + 4$ polynomial. If $md \gtrsim N \cdot C_3^q \log^{4q+4}(1/\delta')$ for unspecified constant C_3 , then $\frac{2^{4q+4} e^{-1} \log^{4q+4}(1/\delta') \text{Var}(\gamma_{xy})}{(\mathbb{E}\gamma_{xy})^2} \leq 1$, and thus by Lemma 14, we have that $\mathbb{P}(\gamma_{xy} \leq 0) \leq \delta'$. Choosing $\delta' = \frac{\delta}{MN}$ and union bounding over all (x, y) pairs with $y \neq f^*(x)$ yields the desired result. \square

B.1 Bounded Bit Complexity

Corollary 2. *Under the setting of Theorem 1, if $d^2 \gtrsim N \text{ poly log } N$, then with high probability there exists a quantized weight matrix $\tilde{\mathbf{W}}$, where each weight requires $O(\log d)$ bits to store, such that*

$$\arg \max_{y \in [M]} \mathbf{u}_y^\top \tilde{\mathbf{W}} \mathbf{e}_x = f^*(x) \quad \text{for all } x \in [N]. \quad (21)$$

Proof. One sees from the proof of Theorem 1 that, with high probability over the embeddings, the weight matrix $\mathbf{W} = \sum_{z \in [N]} \mathbf{u}_{f^*(z)} \mathbf{e}_z^\top$ has a margin γ_{xy} satisfies $\gamma_{xy} \geq \frac{1}{2}$ for all $y \neq f^*(x)$. Each entry of \mathbf{W} lies in the interval $[-N, N]$. For some $\epsilon > 0$, define $\tilde{\mathbf{W}}$ by rounding each entry of \mathbf{W} to the nearest multiple of ϵ . By definition, $\|\mathbf{W} - \tilde{\mathbf{W}}\|_\infty \leq \epsilon$. We also see that

$$\left| \mathbf{u}_y^\top (\mathbf{W} - \tilde{\mathbf{W}}) \mathbf{e}_x \right| \leq \|\mathbf{W} - \tilde{\mathbf{W}}\|_\infty \|\mathbf{u}_y \mathbf{e}_x^\top\|_1 \leq d\epsilon.$$

Thus choosing $\epsilon < \frac{1}{8d}$, the margin of the quantized network satisfies

$$\begin{aligned} \tilde{\gamma}_{xy} &:= (\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top \tilde{\mathbf{W}} \mathbf{e}_x \\ &\geq (\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top \mathbf{W} \mathbf{e}_x - \left| (\mathbf{u}_{f^*(x)} - \mathbf{u}_y)^\top (\mathbf{W} - \tilde{\mathbf{W}}) \mathbf{e}_x \right| \\ &\geq \frac{1}{2} - 2d\epsilon \\ &> 0. \end{aligned}$$

Finally, the number of bits required to store each weight is $\log(2N/\epsilon) = \log(16Nd) = O(\log d)$. \square

We remark that a similar quantization argument was proven in Jelassi et al. [19].

C Proofs for Section 4

Lemma 1. *Let $\mathcal{V}^{(h)} \subset \mathcal{S} \cup \mathcal{R}$. Assume that $d \gtrsim |\mathcal{V}^{(h)}| \log(|\mathcal{V}|/\delta)$. Define $\mathbf{v} := \sum_{z \in \mathcal{V}^{(h)}} \varphi(z) + \frac{1}{2}\varphi(\text{EOS})$. Then, with probability $1 - \delta$ over the draw of the embeddings,*

$$\langle \mathbf{v}, \varphi(z) \rangle > \langle \mathbf{v}, \varphi(\text{EOS}) \rangle + \frac{1}{4} > \langle \mathbf{v}, \varphi(z') \rangle + \frac{1}{2}$$

for any $z \in \mathcal{V}^{(h)}$ and $z' \notin \mathcal{V}^{(h)}$.

Proof. Define γ_z as

$$\gamma_z := \begin{cases} \langle \mathbf{v}, \varphi(z) \rangle - 1 & z \in \mathcal{V}^{(h)} \\ \langle \mathbf{v}, \varphi(\text{EOS}) \rangle - \frac{1}{2} & z = \text{EOS} \\ \langle \mathbf{v}, \varphi(z) \rangle & z \notin \mathcal{V}^{(h)} \end{cases}$$

We first see that $\mathbb{E}[\gamma_z] = 0$.

Next, observe that

$$\gamma_z = \begin{cases} \sum_{z' \in \mathcal{V}^{(h)}} \langle \varphi(z), \varphi(z') \rangle & z \notin \mathcal{V}^{(h)} \\ \sum_{z' \in \mathcal{V}^{(h)} \setminus \{z\}} \langle \varphi(z), \varphi(z') \rangle & z \in \mathcal{V}^{(h)} \end{cases}$$

Since each of the $\langle \varphi(z), \varphi(z') \rangle$ are independent subGaussian variables with variance proxy $1/d$, by Hoeffding's inequality we have that, with probability $1 - \delta'$,

$$|\gamma_z| \lesssim \sqrt{\frac{|\mathcal{V}^{(h)}| \cdot \log(1/\delta')}{d}}.$$

Setting $\delta' = \delta/|\mathcal{V}|$ and union bounding over all $z \in \mathcal{V}$ yields the desired result. \square

C.1 Construction via Self-Attention

Lemma 2. Let $\mathcal{V}^{(h)} \subset \mathcal{V}$, and for each $z \in \mathcal{V}^{(h)}$, let $\mathcal{A}^z \subset \mathcal{V}$. Assume that $d \gtrsim \max_{z \in \mathcal{V}^{(h)}} |\mathcal{A}^z| \log^6(|\mathcal{V}|/\delta)$ and $d_h \gtrsim |\mathcal{V}^{(h)}| \log^6(|\mathcal{V}|/\delta)$. Define

$$\mathbf{W} := \frac{d}{d_h} \sum_{z \in \mathcal{V}^{(h)}} \sum_{a \in \mathcal{A}^z} \sum_{i=1}^{d_h} \varphi(a) \varphi(z)^\top \mathbf{w}_i \mathbf{w}_i^\top,$$

where \mathbf{w}_i are chosen uniformly on the sphere of radius 1, conditioned on being orthogonal to $\varphi(\text{EOS})$. Then, with probability $1 - \delta$ over the draw of the embeddings and the \mathbf{w}_i ,

$$|\varphi(a)^\top \mathbf{W} \varphi(z) - \mathbf{1}(a \in \mathcal{A}^z)| \leq \frac{1}{5}$$

for all $z \in \mathcal{V}^{(h)}$, $a \in \mathcal{V}$.

Proof. Define

$$\gamma_{az} := \varphi(a)^\top \mathbf{W} \varphi(z).$$

We first see that

$$\mathbb{E}[\gamma_{az}] = \mathbb{E} \left[\sum_{z' \in \mathcal{V}^{(h)}} \sum_{a' \in \mathcal{A}^{z'}} \langle \varphi(a), \varphi(a') \rangle \langle \varphi(z), P_{\varphi(\text{EOS})}^\perp \varphi(z') \rangle \right] = \frac{d-1}{d} \cdot \mathbf{1}(a \in \mathcal{A}^z).$$

We next compute the variance. For $a \notin \mathcal{A}^z$,

$$\mathbb{E}[\gamma_{az}^2] = \frac{d^2}{d_h^2} \mathbb{E} \left[\left(\sum_{z' \in \mathcal{V}^{(h)}} \sum_{a' \in \mathcal{A}^{z'}} \sum_{i=1}^{d_h} \langle \varphi(a), \varphi(a') \rangle \langle \varphi(z), \mathbf{w}_i \rangle \langle \mathbf{w}_i, \varphi(z') \rangle \right)^2 \right]$$

Define $\omega_{a'z'i} = \langle \varphi(a), \varphi(a') \rangle \langle \varphi(z), \mathbf{w}_i \rangle \langle \mathbf{w}_i, \varphi(z') \rangle$. We see that $\mathbb{E}[\omega_{a_1 z_1 i} \omega_{a_2 z_2 j}]$ is nonzero only if $a_1 = a_2$ and $z_1 = z_2$. For $i \neq j$, we have that

$$\begin{aligned} \mathbb{E}[\omega_{a'z'i} \omega_{a'z'j}] &= \mathbb{E}[\langle \varphi(a), \varphi(a') \rangle^2 \langle \varphi(z), \mathbf{w}_i \rangle \langle \mathbf{w}_i, \varphi(z') \rangle \langle \varphi(z), \mathbf{w}_j \rangle \langle \mathbf{w}_j, \varphi(z') \rangle] \\ &= d^{-2} \mathbb{E}[\langle \varphi(a), \varphi(a') \rangle^2] \mathbb{E}[\langle \varphi(z), P_{\varphi(\text{EOS})}^\perp \varphi(z') \rangle^2] \\ &\leq d^{-2} \rho_{aa'} \rho_{zz'}, \end{aligned}$$

where $\rho_{ij} = \begin{cases} 1 & i = j \\ d^{-1} & i \neq j \end{cases}$. Also,

$$\mathbb{E}[\omega_{a'z'i}^2] = \mathbb{E}[\langle \varphi(a), \varphi(a') \rangle^2 \langle \varphi(z), \mathbf{w}_i \rangle^2 \langle \mathbf{w}_i, \varphi(z') \rangle^2].$$

Since $\mathbb{E}[\mathbf{w}_i^{\otimes 4}] = \frac{3}{(d-1)(d+1)} \text{Sym}\left((P_{\varphi(\text{EOS})}^\perp)^{\otimes 2}\right)$,

$$\begin{aligned} \mathbb{E}[\omega_{a'z'i}^2] &\leq (d^2 - 1)^{-1} \mathbb{E}[\langle \varphi(a), \varphi(a') \rangle^2] \left(\mathbb{E}[\|P_{\varphi(\text{EOS})}^\perp \varphi(z)\|^2]^2 + 2\mathbb{E}[\langle \varphi(z), P_{\varphi(\text{EOS})}^\perp \varphi(z') \rangle^2] \right) \\ &\leq d^{-2} \rho_{aa'} (1 + 2\rho_{zz'}). \end{aligned}$$

Altogether,

$$\begin{aligned} \mathbb{E}[\gamma_{az}^2] &= \frac{d^2}{d_h^2} \sum_{z' \in \mathcal{V}^{(h)}} \sum_{a' \in \mathcal{A}^{z'}} \sum_{i,j=1}^{d_h} \mathbb{E}[\omega_{a'z'i} \omega_{a'z'j}] \\ &\leq \sum_{z' \in \mathcal{V}^{(h)}} \sum_{a' \in \mathcal{A}^{z'}} \left(\frac{d_h - 1}{d_h} \rho_{aa'} \rho_{zz'} + \frac{1}{d_h} \rho_{aa'} (1 + 2\rho_{zz'}) \right) \\ &= \frac{d_h + 2}{d_h} \sum_{a' \in \mathcal{A}^z} \rho_{aa'} + \frac{d + d_h + 1}{dd_h} \sum_{z' \in \mathcal{V}^{(h)} \setminus \{z\}} \sum_{a' \in \mathcal{A}^{z'}} \rho_{aa'} \\ &\leq \left(\frac{d_h + 2}{d_h} \right) \left(\mathbf{1}(a \in \mathcal{A}^z) + \frac{|\mathcal{A}^z|}{d} \right) + \frac{d + d_h + 1}{dd_h} \cdot \sum_{z' \in \mathcal{V}^{(h)} \setminus \{z\}} \left(\mathbf{1}(a \in \mathcal{A}^{z'}) + \frac{|\mathcal{A}^{z'}|}{d} \right), \end{aligned}$$

and thus

$$\begin{aligned} \text{Var}(\gamma_{az}) &= \mathbb{E}[\gamma_{az}^2] - \frac{d-1}{d} \cdot \mathbf{1}(a \in \mathcal{A}^z) \\ &\lesssim \frac{|\mathcal{A}^z|}{d} + \frac{1}{d_h} \sum_{z' \in \mathcal{V}^{(h)}} \left(\mathbf{1}(a \in \mathcal{A}^{z'}) + \frac{|\mathcal{A}^{z'}|}{d} \right) \\ &\lesssim \frac{\max_{z \in \mathcal{V}^{(h)}} |\mathcal{A}^z|}{d} + \frac{|\mathcal{V}^{(h)}|}{d_h}, \end{aligned}$$

since $d \geq |\mathcal{A}^z|$, $d_h \geq |\mathcal{V}^{(h)}|$. Next, since γ_{az} is a degree 6 polynomial, with probability $1 - \frac{\delta}{|\mathcal{V}||\mathcal{V}^{(h)}|}$

we have that

$$\begin{aligned}
|\gamma_{az} - \mathbf{1}(a \in \mathcal{A}^z)| &\lesssim \sqrt{\text{Var}(\gamma_{az}) \log^6(|\mathcal{V}|/\delta)} \\
&\lesssim \sqrt{\frac{\max_{z \in \mathcal{V}^{(h)}} |\mathcal{A}^z| \log^6(|\mathcal{V}|/\delta)}{d} + \frac{|\mathcal{V}^{(h)}| \log^6(|\mathcal{V}|/\delta)}{d_h}} \\
&\leq \frac{1}{5}.
\end{aligned}$$

Union bounding over all $z \in \mathcal{V}^{(h)}$, $a \in \mathcal{V}$ yields the desired result. \square

Let us state the formal version of Theorem 3 which we aim to prove:

Theorem 10. *Assume that $d \gtrsim \max(R, D) \cdot \log^6(|\mathcal{V}|SR/\delta)$ and $Hd_h \gtrsim S \log^6(|\mathcal{V}|SR/\delta)$. Then, with probability $1 - \delta$, there exists a single-layer attention-only transformer $F_{\text{TF}}(\cdot; \boldsymbol{\theta}_{\text{TF}})$ with embedding dimension d , number of heads H and head dimension d_h such that*

$$\mathbb{P}_{z_{1:T+1} \sim \mathcal{D}} \left[\arg \max_{z \in \mathcal{V}} \boldsymbol{\varphi}(z)^\top F_{\text{TF}}(\mathbf{X}; \boldsymbol{\theta}_{\text{TF}}) = z_{T+1} \right] = 1.$$

Remark. When $R \geq D$, one can obtain an accuracy of 100% whenever the total parameter count is

$$Hdd_h \gtrsim SR \text{poly} \log(|\mathcal{V}|SR/\delta).$$

Proof of Theorem 10. Partition \mathcal{S} into the sets $\mathcal{S}^{(1)}, \dots, \mathcal{S}^{(N_S)}$ and \mathcal{R} into the sets $\mathcal{R}^{(1)}, \dots, \mathcal{R}^{(N_R)}$, such that $|\mathcal{S}^{(i)}|, |\mathcal{R}^{(j)}| \leq M$ and $N_S = \lceil \frac{S}{M} \rceil, N_R = \lceil \frac{R}{M} \rceil$.

Let us choose M so that $d \geq d_h \gtrsim M \log^6(|\mathcal{V}|/\delta)$. The total number of attention heads is then

$$H = N_S + N_R \gtrsim \frac{S \log^6(|\mathcal{V}|/\delta)}{d_h}.$$

For each $i \in [N_S]$, we construct the attention head i as follows. First, let

$$\mathbf{W}_K^{(i)\top} \mathbf{W}_Q^{(i)} = \beta \sum_{z \in \mathcal{S}^{(i)}} \boldsymbol{\varphi}(z) \boldsymbol{\varphi}(\text{EOS})^\top + \frac{\beta}{2} \boldsymbol{\varphi}(\text{EOS}) \boldsymbol{\varphi}(\text{EOS})^\top$$

for a large constant β . Next, set

$$\mathbf{W}_O^{(i)\top} \mathbf{W}_V^{(i)} = \frac{d}{d_h} \sum_{z \in \mathcal{S}^{(i)}} \sum_{a \in \mathcal{A}_z} \sum_{i=1}^{d_h} \boldsymbol{\varphi}(a) \boldsymbol{\varphi}(z)^\top \mathbf{w}_i \mathbf{w}_i^\top,$$

for \mathbf{w}_i sampled uniformly on the sphere, orthogonal to $\boldsymbol{\varphi}(\text{EOS})$.

Consider an input sequence (z_1, \dots, z_T) , and let s be the subject token in this sequence. On the event that Lemma 1 holds, if $s \notin \mathcal{S}^{(i)}$, then $z_t \notin \mathcal{S}^{(i)}$, and thus

$$\varphi(z_t)^\top \mathbf{W}_K^{(i)\top} \mathbf{W}_Q^{(i)} \varphi(\text{EOS}) < \varphi(\text{EOS})^\top \mathbf{W}_K^{(i)\top} \mathbf{W}_Q^{(i)} \varphi(\text{EOS}) - \frac{\beta}{2}$$

for all $t < T$. As $\beta \rightarrow \infty$, the self-attention module fully attends to the EOS token. On the other hand, if $s \in \mathcal{S}^{(i)}$, then if $z_{t^*} = s$ we have

$$\varphi(z_t)^\top \mathbf{W}_K^{(i)\top} \mathbf{W}_Q^{(i)} \varphi(\text{EOS}) < \varphi(z_{t^*})^\top \mathbf{W}_K^{(i)\top} \mathbf{W}_Q^{(i)} \varphi(\text{EOS}) - \frac{\beta}{2}$$

for all $t \neq t^*$. Likewise, as $\beta \rightarrow \infty$, the softmax converges to a hardmax on the z_{t^*} token. Altogether, we get that

$$\mathbf{X}^\top \mathcal{S} \left(\mathbf{X} \mathbf{W}_K^{(i)\top} \mathbf{W}_Q^{(i)} \mathbf{x}_T \right) = \begin{cases} \varphi(\text{EOS}) & s \notin \mathcal{S}^{(i)} \\ \varphi(s) & s \in \mathcal{S}^{(i)} \end{cases}.$$

Next, on the event that Lemma 2 holds, since $d \gtrsim R \log^6(|\mathcal{V}|/\delta) \geq |\mathcal{A}_s| \log^6(|\mathcal{V}|/\delta)$, we have that

$$\left| \varphi(a)^\top \mathbf{W}_O^{(i)\top} \mathbf{W}_V^{(i)} \varphi(s) - \mathbf{1}(a \in \mathcal{A}_s) \right| \leq \frac{1}{6}$$

for $s \in \mathcal{S}^{(i)}$. Defining $\text{attn}_i := \mathbf{W}_O^{(i)\top} \mathbf{W}_V^{(i)} \mathcal{S} \left(\mathbf{X} \mathbf{W}_K^{(i)\top} \mathbf{W}_Q^{(i)} \mathbf{x}_T \right)$, we have that

$$\varphi(a)^\top \text{attn}_i \in \begin{cases} \{0\} & s \notin \mathcal{S}^{(i)} \\ [-\frac{1}{5}, \frac{1}{5}] & s \in \mathcal{S}^{(i)}, a \notin \mathcal{A}_s \\ [\frac{4}{5}, \frac{6}{5}] & s \in \mathcal{S}^{(i)}, a \in \mathcal{A}_s \end{cases}.$$

By an identical construction, for each $j \in [N_R]$, with probability $1 - 2\delta$ we can construct the attention head $N_S + j$ such that

$$\varphi(a)^\top \text{attn}_{N_S+j} \in \begin{cases} \{0\} & r \notin \mathcal{R}^{(i)} \\ [-\frac{1}{5}, \frac{1}{5}] & r \in \mathcal{R}^{(i)}, a \notin \mathcal{A}_r \\ [\frac{4}{5}, \frac{6}{5}] & r \in \mathcal{R}^{(i)}, a \in \mathcal{A}_r \end{cases}$$

as long as $d \gtrsim D \log^6(|\mathcal{V}|/\delta) \geq |\mathcal{A}_r| \log^6(|\mathcal{V}|/\delta)$. Therefore by a union bound, with probability $1 - 2SR\delta$ we have that (where $s \in \mathcal{S}^{(i)}$ and $r \in \mathcal{S}^{(j)}$)

$$\begin{aligned} \varphi(a)^\top F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta}) &= \sum_{h=1}^{N_S+N_R} \varphi(a)^\top \text{attn}_h \\ &= \varphi(a)^\top \text{attn}_i + \varphi(a)^\top \text{attn}_{N_S+j} \end{aligned}$$

If $a = a^*(s, r)$, then $a \in \mathcal{A}_s \cap \mathcal{A}_r$, and thus

$$\varphi(a)^\top F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta}) \geq \frac{4}{5} + \frac{4}{5} = \frac{8}{5}.$$

Otherwise, either $\varphi(a)^\top \text{attn}_i$ or $\varphi(a)^\top \text{attn}_{N_S+j}$ is $\leq \frac{1}{5}$ and thus

$$\varphi(a)^\top F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta}) \geq \frac{6}{5} + \frac{1}{5} = \frac{7}{5}.$$

Therefore $\arg \max_{a \in \mathcal{V}} \varphi(a)^\top F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta}) = a^*(s, r)$. Replacing $2SR\delta$ with δ yields the desired result. \square

C.2 Construction via MLP

Lemma 3. Let ϵ be a fixed constant. Assume that q in Assumption 3 satisfies $q = \frac{C_2}{\epsilon}$ for some $C_2 > 2$. Assume that $d \geq S^\epsilon, R^\epsilon$. Define $C(a) = |\{(s, r) : a^*(s, r) = a\}|$.

Let d be odd, and let P, Q be orthogonal $\lfloor d/2 \rfloor$ dimensional subspaces of \mathbb{R}^d . Define $\tilde{\varphi}(s) = \Pi_P \varphi(s), \tilde{\varphi}(r) = \Pi_Q \varphi(r)$.

There exists universal constants C_3, C_4 such that if

$$\begin{aligned} d &\gtrsim (C_3 \log(|\mathcal{V}|/\delta)/\epsilon)^{C_4/\epsilon} \\ m &\gtrsim (C_3 \log(|\mathcal{V}|/\delta))^{C_4/\epsilon} \cdot \max_a C(a) \\ md &\gtrsim (C_3 \log(|\mathcal{V}|/\delta))^{C_4/\epsilon} \cdot SR, \end{aligned}$$

then with probability $1 - \delta$ over the draw of the embeddings there exists a two-layer neural network $F(\mathbf{z}) = \sum_{i \in [m]} \mathbf{v}_i \sigma(\mathbf{w}_i^\top \mathbf{z})$ of width m satisfying

$$\arg \max_{a \in \mathcal{V}} \varphi(a)^\top F(\tilde{\varphi}(s) + \tilde{\varphi}(r)) = a^*(s, r)$$

for all $s \in \mathcal{S}, r \in \mathcal{R}$.

Proof. For odd integers p, k to be determined later, let us set

$$\mathbf{v}_i = \frac{1}{m} \sum_{s, r} \langle \mathbf{He}_{p+k}(\mathbf{w}_i), \tilde{\varphi}(s)^{\otimes p} \otimes \tilde{\varphi}(r)^{\otimes k} \rangle \cdot \varphi(a^*(s, r)),$$

where $\mathbf{He}_{p+k} : \mathbb{R}^d \rightarrow (\mathbb{R}^d)^{\otimes(p+k)}$ is the Hermite tensor of degree $p+k$ (see Appendix F.1). Assume without loss of generality that $c_{p+k} := \mathbb{E}[\sigma^{(p+k)}(z)]$, the $(p+k)$ th Hermite coefficient of σ is positive (the negative case can be handled by negating all the \mathbf{v}_i in the construction)

For some (s, r) , the margin for some $a \neq a^*(s, r)$ is

$$\begin{aligned} \gamma_{sra} &= \langle \varphi(a^*(s, r)) - \varphi(a), F(\tilde{\varphi}(s) + \tilde{\varphi}(r)) \rangle \\ &= \frac{1}{m} \sum_{i \in [m]} \sum_{s', r'} \sigma(\langle \mathbf{w}_i, \tilde{\varphi}(s) + \tilde{\varphi}(r) \rangle) \langle \mathbf{He}_{p+k}(\mathbf{w}_i), \varphi(s')^{\otimes p} \otimes \tilde{\varphi}(r')^{\otimes k} \rangle \\ &\quad \cdot \langle \varphi(a^*(s', r')), \varphi(a^*(s, r)) - \varphi(a) \rangle. \end{aligned}$$

We first see that

$$\begin{aligned} &\mathbb{E}[\gamma_{sra}] \\ &= \sum_{s', r'} \mathbb{E}[\sigma(\langle \mathbf{w}_i, \tilde{\varphi}(s) + \tilde{\varphi}(r) \rangle) \langle \mathbf{He}_{p+k}(\mathbf{w}_i), \varphi(s')^{\otimes p} \otimes \tilde{\varphi}(r')^{\otimes k} \rangle] \\ &\quad \cdot (\mathbf{1}(a^*(s, r) = a^*(s', r')) - \mathbf{1}(a = a^*(s', r'))) \\ &= \sum_{s', r'} \mathbb{E}[\sigma^{(p+k)}(\langle \mathbf{w}_i, \tilde{\varphi}(s) + \tilde{\varphi}(r) \rangle) \langle (\tilde{\varphi}(s) + \tilde{\varphi}(r))^{\otimes(p+k)}, \varphi(s')^{\otimes p} \otimes \tilde{\varphi}(r')^{\otimes k} \rangle] \\ &\quad \cdot (\mathbf{1}(a^*(s, r) = a^*(s', r')) - \mathbf{1}(a = a^*(s', r'))) \end{aligned}$$

If either $s' \neq s$ or $r \neq r'$, we see that conditioned on $\tilde{\varphi}(s), \tilde{\varphi}(r)$, the quantity $\varphi(s')^{\otimes p} \otimes \tilde{\varphi}(r')^{\otimes k}$ is mean zero. Therefore the only nonzero term in the sum is when $(s, r) = (s', r')$, and so

$$\mathbb{E}[\gamma_{sra}] = \mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\sigma^{(p+k)} \left(Z \sqrt{\|\tilde{\varphi}(s)\|^2 + \|\tilde{\varphi}(r)\|^2} \right) \cdot \|\tilde{\varphi}(s)\|^{2p} \|\tilde{\varphi}(r)\|^{2k} \right]$$

The quantities $\|\tilde{\varphi}(s)\|^2 - \frac{1}{2}, \|\tilde{\varphi}(r)\|^2 - \frac{1}{2}$ are subexponential random variables with Orlicz norm $1/d$, and therefore we can bound

$$|\mathbb{E}[\gamma_{sra}] - c_{p+k} 2^{-p-k}| \lesssim \frac{1}{d}$$

We next compute the variance. Define $\omega_{is'r'}$ by

$$\omega_{is'r'} = \sigma(\langle \mathbf{w}_i, \tilde{\varphi}(s) + \tilde{\varphi}(r) \rangle) \langle \mathbf{H}_{e_{p+k}}(\mathbf{w}_i), \varphi(s')^{\otimes p} \otimes \tilde{\varphi}(r')^{\otimes k} \rangle \cdot \langle \varphi(a^*(s', r')), \varphi(a^*(s, r)) - \varphi(a) \rangle$$

We first observe that $\mathbb{E}[\omega_{is_1 r_1} \omega_{js_2 r_2}]$ is zero, unless $s_1 = s_2$ and $r_1 = r_2$. Next, we compute the expectation of $\omega_{is'r'}$, conditioned on the embeddings (i.e with respect to the randomness \mathbf{w}_i):

$$\begin{aligned} \mathbb{E}[\omega_{is'r'} \mid \varphi] &= \mathbb{E}_{\mathbf{w}_i} \left[\sigma^{(p+k)}(\langle \mathbf{w}_i, \tilde{\varphi}(s) + \tilde{\varphi}(r) \rangle) \cdot \langle (\tilde{\varphi}(s) + \tilde{\varphi}(r))^{\otimes (p+k)}, \tilde{\varphi}(s')^{\otimes p} \otimes \tilde{\varphi}(r')^{\otimes k} \rangle \right. \\ &\quad \left. \cdot \langle \varphi(a^*(s', r')), \varphi(a^*(s, r)) - \varphi(a) \rangle \right] \\ &= \mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\sigma^{(p+k)} \left(Z \sqrt{\|\tilde{\varphi}(s)\|^2 + \|\tilde{\varphi}(r)\|^2} \right) \right] \cdot \langle \tilde{\varphi}(s), \tilde{\varphi}(s') \rangle^p \cdot \langle \tilde{\varphi}(r), \tilde{\varphi}(r') \rangle^k \cdot \langle \varphi(a^*(s', r')), \varphi(a^*(s, r)) - \varphi(a) \rangle \end{aligned}$$

Therefore for $i \neq j$,

$$\begin{aligned} \mathbb{E}[\omega_{is'r'} \omega_{js'r'}] &= \mathbb{E}[\mathbb{E}[\omega_{is'r'} \mid \varphi]^2] \\ &\lesssim c_{p+k}^2 \mathbb{E}[\langle \tilde{\varphi}(s), \varphi(s') \rangle^{2p}] \mathbb{E}[\langle \tilde{\varphi}(r), \tilde{\varphi}(r') \rangle^{2k}] \mathbb{E}[\langle \varphi(a^*(s', r')), \varphi(a^*(s, r)) - \varphi(a) \rangle^2]. \end{aligned}$$

When $(s, r) = (s', r')$, then

$$\begin{aligned} \mathbb{E}[\omega_{isr} \omega_{jsr}] &= \mathbb{E} \left[\mathbb{E}_{Z \sim \mathcal{N}(0,1)} \left[\sigma^{(p+k)} \left(Z \sqrt{\|\tilde{\varphi}(s)\|^2 + \|\tilde{\varphi}(r)\|^2} \right) \right]^2 \|\tilde{\varphi}(s)\|^{4p} \|\tilde{\varphi}(r)\|^{4k} \right] \cdot (1 + d^{-1}) \\ &= c_{p+k}^2 2^{-2p-2k} + O(1/d) \end{aligned}$$

Next, define the quantities

$$\begin{aligned} \rho_{ss'} &= \mathbb{E}[\langle \tilde{\varphi}(s), \tilde{\varphi}(s') \rangle^{2p}] \\ \rho_{rr'} &= \mathbb{E}[\langle \tilde{\varphi}(r), \tilde{\varphi}(r') \rangle^{2k}] \\ \rho_{aa'} &= \mathbb{E}[\langle \varphi(a), \varphi(a') \rangle^2], \end{aligned}$$

so that

$$\mathbb{E}[\omega_{is'r'} \omega_{js'r'}] \lesssim c_{p+k}^2 \rho_{ss'} \rho_{rr'} (\rho_{aa^*(s', r')} + \rho_{a^*(s, r) a^*(s', r')}).$$

We see that for $s \neq s', r \neq r', a \neq a'$,

$$\begin{aligned}\rho_{ss'} &\leq (2p)^p (d/2)^{-p} = (4p)^p d^{-p} \\ \rho_{rr'} &\leq (2k)^k (d/2)^{-k} = (4k)^k d^{-k} \\ \rho_{aa'} &\leq d^{-1}\end{aligned}$$

Next, see that

$$\begin{aligned}\mathbb{E}[\omega_{is'r'}^2] &= \mathbb{E}[\sigma(\langle \mathbf{w}_i, \tilde{\varphi}(s) + \tilde{\varphi}(r) \rangle)^2 \langle \mathbf{He}_{p+k}(\mathbf{w}_i), \tilde{\varphi}(s')^{\otimes p} \otimes \tilde{\varphi}(r')^{\otimes k} \rangle^2] \mathbb{E}[\langle \varphi(a^*(s', r')), \varphi(a^*(s, r)) - \varphi(a) \rangle^2] \\ &\leq \mathbb{E}[\sigma(\langle \mathbf{w}_i, \tilde{\varphi}(s) + \tilde{\varphi}(r) \rangle)^4]^{1/2} \mathbb{E}[\langle \mathbf{He}_{p+k}(\mathbf{w}_i), \tilde{\varphi}(s')^{\otimes p} \otimes \tilde{\varphi}(r')^{\otimes k} \rangle^4]^{1/2} (\rho_{aa^*(s', r')} + \rho_{a^*(s, r)a^*(s' r')}) \\ &\leq 2^{4q} (\rho_{aa^*(s', r')} + \rho_{a^*(s, r)a^*(s' r')}),\end{aligned}$$

where we have applied Lemma 13 to the first two expectations. Altogether, we have that

$$\begin{aligned}\mathbb{E}[\gamma_{sra}^2] &= \sum_{s', r'} \left(\frac{m-1}{m} \mathbb{E}[\omega_{is'r'} \omega_{js'r'}] + \frac{1}{m} \mathbb{E}[\omega_{is'r'}^2] \right) \\ &= \frac{m-1}{m} \mathbb{E}[\omega_{isr} \omega_{jsr}] + \frac{m-1}{m} \sum_{(s', r') \neq (s, r)} \mathbb{E}[\omega_{is'r'} \omega_{js'r'}] + \frac{1}{m} \sum_{s', r'} \mathbb{E}[\omega_{is'r'}^2],\end{aligned}$$

and thus

$$\begin{aligned}\text{Var}(\gamma_{sra}) &= \mathbb{E}[\gamma_{sra}^2] - c_{p+k}^2 \\ &\lesssim c_{p+k}^2 \sum_{(s', r') \neq (s, r)} \rho_{ss'} \rho_{rr'} (\rho_{aa^*(s', r')} + \rho_{a^*(s, r)a^*(s' r')}) + \frac{2^{4q}}{m} \sum_{s', r'} (\rho_{aa^*(s', r')} + \rho_{a^*(s, r)a^*(s' r')})\end{aligned}$$

For the first sum, we can bound

$$\begin{aligned}\sum_{(s', r') \neq (s, r)} \rho_{ss'} \rho_{rr'} (\rho_{aa^*(s', r')} + \rho_{a^*(s, r)a^*(s' r')}) &\leq \sum_{(s', r') \neq (s, r)} \rho_{ss'} \rho_{rr'} \\ &\leq SR \cdot (4p)^p (4k)^k d^{-p-k} + S(4p)^p d^{-p} + R(4k)^k d^{-k}\end{aligned}$$

For the second sum, we get that

$$\sum_{s', r'} (\rho_{aa^*(s', r')} + \rho_{a^*(s, r)a^*(s' r')}) \leq C(a) + C(a^*(s, r)) + \frac{2SR}{d}$$

Altogether,

$$\frac{\text{Var}(\gamma_{sra})}{\mathbb{E}[\gamma_{sra}]^2} \lesssim \frac{S(4p)^p}{d^p} + \frac{R(4k)^k}{d^k} + \frac{S(4p)^p}{d^p} \cdot \frac{R(4k)^k}{d^k} + \frac{2^{4q}(C(a) + C(a^*(s, r)))}{mc_{p+k}^2} + \frac{2^{4q}SR}{mdc_{p+k}^2}.$$

Let δ' be a fixed failure probability. We see that, for $p = 2\lceil \frac{1}{\epsilon} \rceil + 1$

$$\frac{2^{4q} \log^{4q+2}(1/\delta') S(4p)^p}{d^p} \lesssim \frac{2^{\frac{4C_2}{\epsilon}} \log^{\frac{4C_2}{\epsilon}+2}(1/\delta') (\frac{8}{\epsilon})^{2/\epsilon}}{d^{\frac{1}{\epsilon}}} \lesssim 1,$$

whenever $d \gtrsim (C_3 \log(1/\delta')/\epsilon)^{C_4/\epsilon}$ for appropriately chosen constants C_3, C_4 . Likewise, setting $k = 2\lceil \frac{1}{\epsilon} \rceil + 1$, we get that

$$\frac{2^{4q} \log^{4q+2}(1/\delta') R(4k)^k}{d^k} \lesssim 1.$$

Next, setting $m \gtrsim 2^{8q+2} \log^{4q+2}(1/\delta') c_{p+k}^{-2} \cdot \max_a C(a)$ and $md \gtrsim 2^{8q+2} \log^{4q+2}(1/\delta') c_{p+k}^{-2} SR$ yields

$$2^{4q+2} \log^{4q+2}(1/\delta') \cdot \frac{2^{4q}(C(a) + C(a^*(s, r)))}{m c_{p+k}^2} \lesssim 1$$

$$2^{4q+2} \log^{4q+2}(1/\delta') \cdot \frac{2^{4q} SR}{m d c_{p+k}^2} \lesssim 1.$$

Altogether, by choosing constants appropriately, we get that

$$2^{4q+2} \log^{4q+2}(1/\delta') e^{-1} \cdot \frac{\text{Var}(\gamma_{sra})}{\mathbb{E}[\gamma_{sra}]^2} \leq 1.$$

Therefore by Lemma 14, with probability $1 - \delta'$ we have that $\gamma_{sra} > 0$. Union bounding over all s, r, a and setting $\delta' = \frac{\delta}{SR|\mathcal{V}|}$ yields the desired result. \square

We next state the formal version of Theorem 4, which we wish to prove:

Theorem 11. *Let ϵ be a fixed constant. Assume that σ is a degree q polynomial, where $q = C_1/\epsilon$ for some $C_1 > 2$. Assume that $d \geq S^\epsilon, R^\epsilon$. Define $C(a) = |\{(s, r) : a^*(s, r) = a\}|$.*

Let (d, H, d_h, m) satisfy

$$d \gtrsim (C_2 \log(|\mathcal{V}|/\delta)/\epsilon)^{C_3/\epsilon}$$

$$H d_h \gtrsim (S + R) \log(|\mathcal{V}|/\delta)$$

$$m \gtrsim (C_2 \log(|\mathcal{V}|/\delta))^{C_3/\epsilon} \cdot \max_a C(a)$$

$$md \gtrsim (C_2 \log(|\mathcal{V}|/\delta))^{C_3/\epsilon} \cdot SR,$$

Then, with probability $1 - \delta$, there exists a single-layer transformer $F_{\text{TF}}(\cdot; \boldsymbol{\theta}_{\text{TF}})$ with embedding dimension d , number of heads H , head dimension d_h , and MLP width m such that

$$\mathbb{P}_{z_1:T+1 \sim \mathcal{D}} \left[\arg \max_{z \in \mathcal{V}} \boldsymbol{\varphi}(z)^\top F_{\text{TF}}(\mathbf{X}; \boldsymbol{\theta}_{\text{TF}}) = z_{T+1} \right] = 1.$$

Remark. Ignoring polylog factors, and treating ϵ as a constant, the constraints on the architecture size become

$$Hd_h \gtrsim S + R \quad \text{and} \quad m \gtrsim C(a) \quad \text{and} \quad md \gtrsim SR.$$

We first note that $C(a) \leq S$, and so $m \gtrsim S$ is sufficient. It is possible for $C(a)$ to be much smaller; on average we expect $C(a) \approx S/D$, and we also note that it is possible for $C(a) = 1$. The main constraint is that $md \gtrsim SR$, i.e that the number of MLP parameters scales linearly with the number of facts that need to be stored.

Proof of Theorem 11. Partition \mathcal{S} into the sets $\mathcal{S}^{(1)}, \dots, \mathcal{S}^{(N_S)}$ and \mathcal{R} into the sets $\mathcal{R}^{(1)}, \dots, \mathcal{R}^{(N_R)}$, such that $|\mathcal{S}^{(i)}|, |\mathcal{R}^{(j)}| \leq M$ and $N_S = \lceil \frac{S}{M} \rceil, N_R = \lceil \frac{R}{M} \rceil$. Assume that $d = \Theta(M \log(|\mathcal{V}|/\delta'))$

Let $H = \lceil d/d_h \rceil$. For each $i \in [N_S]$, we construct the H' attention heads corresponding to $h \in \{(i-1)H'+1, \dots, iH'\}$ as follows. First, for all such h , let

$$\mathbf{W}_K^{(h)\top} \mathbf{W}_Q^{(h)} = \beta \sum_{z \in \mathcal{S}^{(i)}} \varphi(z) \varphi(\text{EOS})^\top + \frac{\beta}{2} \varphi(\text{EOS}) \varphi(\text{EOS})^\top$$

for a large constant β . By an identical argument to as in Theorem 3, on the event that Lemma 1 holds we have that

$$\mathbf{X}^\top \mathcal{S} \left(\mathbf{X} \mathbf{W}_K^{(h)\top} \mathbf{W}_Q^{(h)} \mathbf{x}_T \right) = \begin{cases} \varphi(\text{EOS}) & s \notin \mathcal{S}^{(i)} \\ \varphi(s) & s \in \mathcal{S}^{(i)}. \end{cases}$$

The total contribution from these attention heads is then

$$\sum_{h=(i-1)H'+1}^{iH'} \mathbf{W}_O^{(h)\top} \text{attn}(\mathbf{X}; \mathbf{W}_K^{(h)}, \mathbf{W}_Q^{(h)}, \mathbf{W}_V^{(h)}) = \left(\sum_{h=(i-1)H'+1}^{iH'} \mathbf{W}_O^{(h)\top} \mathbf{W}_V^{(h)} \right) \cdot \begin{cases} \varphi(\text{EOS}) & s \notin \mathcal{S}^{(i)} \\ \varphi(s) & s \in \mathcal{S}^{(i)} \end{cases}$$

Since $H'd_h \geq d$, we can let $\sum_{h=(i-1)H'+1}^{iH'} \mathbf{W}_O^{(h)\top} \mathbf{W}_V^{(h)}$ be a projection onto a $\lceil d/2 \rceil$ dimensional subspace P , orthogonal to $\varphi(\text{EOS})$, and thus

$$\sum_{h=(i-1)H'+1}^{iH'} \mathbf{W}_O^{(h)\top} \text{attn}(\mathbf{X}; \mathbf{W}_K^{(h)}, \mathbf{W}_Q^{(h)}, \mathbf{W}_V^{(h)}) = \begin{cases} 0 & s \notin \mathcal{S}^{(i)} \\ \Pi_P \varphi(s) & s \in \mathcal{S}^{(i)} \end{cases}$$

Altogether, if the sequence (z_1, \dots, z_T) contains the subject s , then

$$\sum_{h=1}^{H'N_S} \mathbf{W}_O^{(h)\top} \text{attn}(\mathbf{X}; \mathbf{W}_K^{(h)}, \mathbf{W}_Q^{(h)}, \mathbf{W}_V^{(h)}) = \Pi_P \varphi(s)$$

Similarly, if we let Q be a $\lceil d/2 \rceil$ dimensional subspace orthogonal to P and $\varphi(\text{EOS})$, then we can construct the attention heads $h \in \{H'N_S + 1, \dots, H'N_S + H'N_R\}$ such that

$$\sum_{h=H'N_S+1}^{H'N_S+H'N_R} \mathbf{W}_O^{(h)\top} \text{attn}(\mathbf{X}; \mathbf{W}_K^{(h)}, \mathbf{W}_Q^{(h)}, \mathbf{W}_V^{(h)}) = \Pi_Q \varphi(r),$$

where r is the relation in the sequence (z_1, \dots, z_T) . Such a construction exists with probability $1 - (N_S + N_R)\delta'$. The total number of heads is

$$H = H'N_S + H'N_R \propto \frac{d(S+R)}{d_h M} \propto \frac{(S+R) \log(|\mathcal{V}|/\delta')}{d_h}.$$

The output of the self-attention component is then

$$F_{\text{MHSA}}(\mathbf{X}; \boldsymbol{\theta}) = \Pi_P \boldsymbol{\varphi}(s) + \Pi_Q \boldsymbol{\varphi}(r) = \tilde{\boldsymbol{\varphi}}(s) + \tilde{\boldsymbol{\varphi}}(r).$$

On the event that Lemma 3 holds, we have that there exists a two-layer neural network $F(z) = \sum_{i \in [m]} \mathbf{v}_i \sigma(\mathbf{w}_i^\top z)$ of width m such that

$$\arg \max_a \boldsymbol{\varphi}(a)^\top F(\boldsymbol{\varphi}(s) + \tilde{\boldsymbol{\varphi}}(r)) = a^*(s, r).$$

Scaling \mathbf{V} by a large enough constant ensures that

$$\arg \max_{z \in \mathcal{V}} \boldsymbol{\varphi}(z)^\top F_{\text{TF}}(\mathbf{X}; \boldsymbol{\theta}_{\text{TF}}) = a^*(s, r).$$

Union bounding over all the high probability events and setting $\delta = \delta'/(N_S + N_R + 1)$ yields the desired result. \square

D Proofs for Section 5

D.1 Preliminaries

Recall that the parameters are $\boldsymbol{\theta} := \{\mathbf{W}_{OV}(a, z)\}_{a \in \mathcal{A}, z \in \mathcal{V}} \cup \{\mathbf{W}_{KQ}(z)\}_{z \in \mathcal{V}}$, and that the cross entropy loss is

$$L(\boldsymbol{\theta}) := \mathbb{E}_{z_1:T+1} \left[-\langle \boldsymbol{\varphi}(z_{T+1}), F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) \rangle + \log \left(\sum_{a \in \mathcal{A}} \exp(\langle \boldsymbol{\varphi}(a), F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) \rangle) \right) \right]$$

where

$$\boldsymbol{\varphi}(a)^\top F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) = \sum_{t=1}^T \mathbf{W}_{OV}(a, z_t) \mathbf{W}_{KQ}(z_t).$$

We consider running gradient flow:

$$\dot{\boldsymbol{\theta}} = -\nabla L(\boldsymbol{\theta})$$

from the initialization $\mathbf{W}_{OV}(a, z) = \alpha$, $\mathbf{W}_{KQ}(z) = \alpha \sqrt{|\mathcal{A}| + 1}$ for some $\alpha > 0$.

We also define Θ by

$$\Theta(a, z) = \mathbf{W}_{KQ}(z) \mathbf{W}_{OV}(a, z),$$

and remark that the loss L is convex in Θ .

Lemma 4 (Balancedness). Let $C(z_{1:T}, z)$ denote the number of tokens in $z_{1:T}$ equal to z . The loss gradients are given by

$$\begin{aligned}\partial_{\mathbf{W}_{VO}(a,z)}L(\boldsymbol{\theta}) &= -\mathbf{W}_{KQ}(z) \cdot \mathbb{E}_{z_{1:T}}[C(z_{1:T}, z) \cdot (\mathbf{1}(a = a^*(z_{1:T})) - \hat{p}(a | z_{1:T}))] \\ \partial_{\mathbf{W}_{KQ}(z)}L(\boldsymbol{\theta}) &= -\sum_a \mathbf{W}_{OV}(a, z) \cdot \mathbb{E}_{z_{1:T}}[C(z_{1:T}, z) \cdot (\mathbf{1}(a = a^*(z_{1:T})) - \hat{p}(a | z_{1:T}))]\end{aligned}$$

As such, the quantity

$$\mathbf{W}_{KQ}(z)^2 - \sum_{a \in \mathcal{A}} \mathbf{W}_{VO}(a, z)^2$$

is constant throughout the gradient flow trajectory.

Proof. We first see that

$$\partial_{\mathbf{W}_{VO}(a,z)}(\boldsymbol{\varphi}(a')^\top F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta})) = \mathbf{1}(a = a') \cdot C(z_{1:T}, z) \cdot \mathbf{W}_{KQ}(z),$$

Similarly,

$$\partial_{\mathbf{W}_{KQ}(z)}(\boldsymbol{\varphi}(a')^\top F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta})) = C(z_{1:T}, z) \cdot \mathbf{W}_{OV}(a', z).$$

Therefore

$$\begin{aligned}\partial_{\mathbf{W}_{VO}(a,z)}L(\boldsymbol{\theta}) &= \mathbf{W}_{KQ}(z) \cdot \mathbb{E} \left[-\mathbf{1}(z_{T+1} = a) \cdot C(z_{1:T}, z) + \frac{\sum_{a' \in \mathcal{A}} \exp(\langle \boldsymbol{\varphi}(a'), F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) \rangle) \cdot \mathbf{1}(a = a') \cdot C(z_{1:T}, z)}{\sum_{a' \in \mathcal{A}} \exp(\langle \boldsymbol{\varphi}(a'), F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) \rangle)} \right] \\ &= -\mathbf{W}_{KQ}(z) \cdot \mathbb{E}_{z_{1:T}}[C(z_{1:T}, z) \cdot (\mathbf{1}(a = a^*(z_{1:T})) - \hat{p}(a | z_{1:T}))].\end{aligned}$$

By a similar computation,

$$\begin{aligned}\partial_{\mathbf{W}_{KQ}(z)}L(\boldsymbol{\theta}) &= \mathbb{E}_{z_{1:T}} \left[-\mathbf{W}_{OV}(z_{T+1}, z) \cdot C(z_{1:T}, z) + \sum_a \hat{p}(a | z_{1:T}) \mathbf{W}_{OV}(a, z) \cdot C(z_{1:T}, z) \right] \\ &= \mathbb{E}_{z_{1:T}} \left[C(z_{1:T}, z) \cdot \left(-\mathbf{W}_{OV}(a^*(z_{1:T}), z) + \sum_a \hat{p}(a | z_{1:T}) \mathbf{W}_{OV}(a, z) \right) \right] \\ &= -\sum_a \mathbf{W}_{OV}(a, z) \cdot \mathbb{E}_{z_{1:T}}[C(z_{1:T}, z) \cdot (\mathbf{1}(a = a^*(z_{1:T})) - \hat{p}(a | z_{1:T}))].\end{aligned}$$

Under gradient flow, we see that

$$\begin{aligned}
& \frac{1}{2} \frac{d}{dt} \left(\mathbf{W}_{KQ}(z)^2 - \sum_{a \in \mathcal{A}} \mathbf{W}_{VO}(a, z)^2 \right) \\
&= \mathbf{W}_{KQ}(z) \cdot \frac{d}{dt} \mathbf{W}_{KQ}(z) - \sum_{a \in \mathcal{A}} \mathbf{W}_{VO}(a, z) \cdot \frac{d}{dt} \mathbf{W}_{VO}(a, z) \\
&= -\mathbf{W}_{KQ}(z) \cdot \partial_{\mathbf{W}_{KQ}(z)} L(\boldsymbol{\theta}) + \sum_{a \in \mathcal{A}} \mathbf{W}_{VO}(a, z) \cdot \partial_{\mathbf{W}_{VO}(a, z)} L(\boldsymbol{\theta}) \\
&= \mathbf{W}_{KQ}(z) \sum_a \mathbf{W}_{OV}(a, z) \cdot \mathbb{E}_{z_{1:T}} [C(z_{1:T}, z) \cdot (\mathbf{1}(a = a^*(z_{1:T})) - \hat{p}(a | z_{1:T}))] \\
&\quad - \sum_{a \in \mathcal{A}} \mathbf{W}_{OV}(a, z) \mathbf{W}_{KQ}(z) \cdot \mathbb{E}_{z_{1:T}} [C(z_{1:T}, z) \cdot (\mathbf{1}(a = a^*(z_{1:T})) - \hat{p}(a | z_{1:T}))] \\
&= 0.
\end{aligned}$$

□

Corollary 3. Throughout the gradient flow trajectory, $\mathbf{W}_{KQ}(z) \geq \alpha$.

Proof. At initialization, $\mathbf{W}_{KQ}(z)^2 - \sum_{a \in \mathcal{A}} \mathbf{W}_{VO}(a, z)^2 = \alpha^2$. Since this quantity is an invariant of gradient flow, it is impossible for $\mathbf{W}_{KQ}(z) = 0$, and thus $\mathbf{W}_{KQ}(z) > 0$ throughout the entire trajectory. Furthermore,

$$\mathbf{W}_{KQ}(z)^2 = \sum_{a \in \mathcal{A}} \mathbf{W}_{VO}(a, z)^2 + \alpha^2 \geq \alpha^2,$$

and thus $\mathbf{W}_{KQ}(z) \geq \alpha$.

□

D.2 Proof of Theorem 5

Proof of Theorem 5. Let us select

$$\epsilon \leq \min \left(\frac{1}{2} \alpha p(s, r) |\mathcal{A}|^{-1} T^{-2} |\mathcal{N}|^{-(T-3)}, \frac{1}{2} \alpha |\mathcal{A}|^{-1} S^{-1} R^{-1} \delta \right).$$

There exists a time T_ϵ such that for all $t \geq T_\epsilon$, $\|\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}(t))\| \leq \epsilon$. Let us set $t_\delta = T_\epsilon$. Now, consider some iterate $\boldsymbol{\theta} := \boldsymbol{\theta}(t)$ for $t \geq t_\delta$.

First, see that for $s \in \mathcal{S}$,

$$\begin{aligned}
\partial_{\mathbf{W}_{OV}(a, s)} L(\boldsymbol{\theta}) &= -\mathbf{W}_{KQ}(s) \cdot \mathbb{E}_{z_{1:T}} [C(z_{1:T}, z) \cdot (\mathbf{1}(a = a^*(z_{1:T})) - \hat{p}(a | z_{1:T}))] \\
&= -\mathbf{W}_{KQ}(s) \cdot p(s) \cdot \mathbb{E}_{z_{1:T}} [\mathbf{1}(a = a^*(z_{1:T})) - \hat{p}(a | z_{1:T}) \mid s \in z_{1:T}].
\end{aligned}$$

Consider some $a \notin \mathcal{A}_s$. Then $\mathbb{E}_{z_{1:T}} [\mathbf{1}(a = a^*(z_{1:T})) \mid s \in z_{1:T}] = 0$, and thus

$$\begin{aligned}
\partial_{\mathbf{W}_{OV}(a, s)} L(\boldsymbol{\theta}) &= \mathbf{W}_{KQ}(s) \cdot p(s) \cdot \mathbb{E}_{z_{1:T}} [\hat{p}(a | z_{1:T}) \mid s \in z_{1:T}] \\
&= \mathbf{W}_{KQ}(s) \sum_{r \in \mathcal{R}} p(s, r) \cdot \mathbb{E}_{z_{1:T}} [\hat{p}(a | z_{1:T}) \mid s, r \in z_{1:T}]
\end{aligned}$$

As such, since $|\partial_{\mathbf{w}_{OV(a,s)}} L(\boldsymbol{\theta})| \leq \epsilon$,

$$\mathbb{E}_{z_{1:T}}[\hat{p}(a | z_{1:T}) | s, r \in z_{1:T}] \leq \epsilon \alpha^{-1} p(s, r)^{-1}.$$

By an identical argument, since $|\partial_{\mathbf{w}_{OV(a,r)}} L(\boldsymbol{\theta})| \leq \epsilon$, then for $a \notin \mathcal{A}_r$

$$\mathbb{E}_{z_{1:T}}[\hat{p}(a | z_{1:T}) | s, r \in z_{1:T}] \leq \epsilon \alpha^{-1} p(s, r)^{-1}.$$

For any $a \neq a^*(s, r)$, either $a \notin \mathcal{A}_s$ or $a \notin \mathcal{A}_r$. Therefore $\mathbb{E}_{z_{1:T}}[\hat{p}(a | z_{1:T}) | s, r \in z_{1:T}] \leq \epsilon \alpha^{-1} p(s, r)^{-1}$ for all $a \neq a^*(s, r)$, and thus

$$\mathbb{E}_{z_{1:T}}[\hat{p}(a^*(s, r) | z_{1:T}) | s, r \in z_{1:T}] \geq 1 - \epsilon \alpha^{-1} p(s, r)^{-1} |\mathcal{A}|.$$

There are at most $T^2 |\mathcal{N}|^{T-3}$ sequences $z_{1:T}$ containing (s, r) , each of which occurs with equal probability. Therefore

$$\hat{p}(a^*(s, r) | z_{1:T}) \geq 1 - T^2 |\mathcal{N}|^{T-3} \cdot \epsilon \alpha^{-1} p(s, r)^{-1} |\mathcal{A}|$$

for all such $z_{1:T}$. Then, bounding $-\log(1 - z) \leq 2z$ for $z \in [0, \frac{1}{2}]$,

$$\begin{aligned} \mathbb{E}[-\log \hat{p}(a^*(s, r) | z_{1:T}) | s, r \in z_{1:T}] &\leq 2\mathbb{E}[1 - \hat{p}(a^*(s, r) | z_{1:T}) | s, r \in z_{1:T}] \\ &\leq 2\epsilon \alpha^{-1} p(s, r)^{-1} |\mathcal{A}|. \end{aligned}$$

Altogether, the loss is

$$\begin{aligned} \mathbb{E}[-\log \hat{p}(z_{T+1} | z_{1:T})] &= \sum_{s,r} p(s, r) \cdot \mathbb{E}[-\log \hat{p}(a^*(s, r) | z_{1:T}) | s, r \in z_{1:T}] \\ &\leq 2\epsilon \alpha^{-1} |\mathcal{A}| SR \\ &\leq \delta, \end{aligned}$$

as desired. □

D.3 Sequential Learning

The goal of this section is to show that the model learns *sequentially*; first, the relation components grow, then the subject components grow. This is given formally by Theorem 6

We first prove that weights corresponding to the subject and noise tokens stay bounded during the beginning of the trajectory.

Lemma 5. For $s \in \mathcal{S}$,

$$\mathbf{W}_{KQ}(z) \leq \exp(2p(s)t) \cdot \alpha \sqrt{|\mathcal{A}| + 1}.$$

Likewise, for $z \in \mathcal{N}$,

$$\mathbf{W}_{KQ}(z) \leq \exp(2Tt/|\mathcal{N}|) \cdot \alpha \sqrt{|\mathcal{A}| + 1}.$$

Proof. Recall that the update for $\mathbf{W}_{KQ}(s)$ is

$$\begin{aligned}\dot{\mathbf{W}}_{KQ}(s) &= p(s) \langle \mathbf{W}_{OV}(\cdot, s), p^*(\cdot | s) - \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | s \in z_{1:T}] \rangle \\ &\leq p(s) \|\mathbf{W}_{OV}(\cdot, s)\| \|p^*(\cdot | s) - \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | s \in z_{1:T}]\| \\ &\leq 2p(s) \|\mathbf{W}_{OV}(\cdot, s)\| \\ &\leq 2p(s) \mathbf{W}_{KQ}(s)\end{aligned}$$

Therefore by Gronwall's inequality,

$$\mathbf{W}_{KQ}(s) \leq \exp(2p(s)t) \cdot \alpha \sqrt{|\mathcal{A}| + 1}.$$

Similarly, the update for $\mathbf{W}_{KQ}(z)$ for $z \in \mathcal{N}$ is

$$\begin{aligned}\dot{\mathbf{W}}_{KQ}(z) &= \langle \mathbf{W}_{OV}(\cdot, z), \mathbb{E}_{z_{1:T}}[C(z_{1:T}, z) \cdot (\mathbf{1}(\cdot = a^*(z_{1:T})) - \hat{p}(\cdot | z_{1:T}))] \rangle \\ &\leq \|\mathbf{W}_{OV}(\cdot, z)\| \cdot \mathbb{E}[C(z_{1:T}, z) \|\mathbf{1}(\cdot = a^*(z_{1:T})) - \hat{p}(\cdot | z_{1:T})\|] \\ &\leq 2\mathbf{W}_{OV}(\cdot, z) \mathbb{E}[C(z_{1:T}, z)] \\ &\leq \frac{2T}{|\mathcal{N}|} \mathbf{W}_{KQ}(z).\end{aligned}$$

Again by Gronwall's inequality,

$$\mathbf{W}_{KQ}(z) \leq \exp(2Tt/|\mathcal{N}|) \cdot \alpha \sqrt{|\mathcal{A}| + 1}.$$

□

The following lemma is our key result, and shows that, assuming that the subject and noise weights stay bounded, the relation weights grow until the output of the model approximates the best relation-only prediction.

Lemma 6. *Let $\alpha_{sm}, \epsilon > 0$ be arbitrary parameters satisfying*

$$\begin{aligned}\alpha_{sm}^2 T &\leq \frac{1}{150} \log \left(\frac{\epsilon^2}{\alpha^2(|\mathcal{A}| + 1)} \right)^{-1} \cdot \min_r \|p^*(\cdot | r) - p_0\| \\ \epsilon^2 &\leq \frac{1}{50(|\mathcal{A}| + 1)} \cdot \min_r \|p^*(\cdot | r) - p_0\|\end{aligned}$$

For a target accuracy $\epsilon_{min} > 0$, define T^ by*

$$T^* = \max_r p(r)^{-1} \|p^*(\cdot | r) - p_0\|^{-1} \log \left(\frac{\epsilon}{\alpha \sqrt{|\mathcal{A}| + 1}} \right) + 100(|\mathcal{A}| + 1) \log |\mathcal{A}| \epsilon^{-2} \epsilon_{min}^{-2}$$

Assume that for $z \in \mathcal{S} \cup \mathcal{N}$ that $\mathbf{W}_{KQ}(z) \leq \alpha_{sm}$. Then, there exists $t \leq T^$ such that*

$$\sum_r p(r)^2 \|p^*(\cdot | r) - \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}]\|^2 \leq \epsilon_{min}^2.$$

Proof. The proof proceeds in three stages. First, we bound the time required for the relation weights to escape the origin. Next, we prove that the relation weights stay large. Finally, we show convergence.

Stage 1: Escaping the origin. The gradient flow update on $\mathbf{W}_{OV}(a, r)$ is

$$\dot{\mathbf{W}}_{OV}(a, r) = \mathbf{W}_{KQ}(r) \cdot p(r)(p^*(a | r) - \mathbb{E}_{z_{1:T}}[\hat{p}(a | z_{1:T}) | r \in z_{1:T}])$$

We thus have

$$\left\| \dot{\mathbf{W}}_{OV}(\cdot, r) - \mathbf{W}_{KQ}(r) \cdot p(r)(p^*(\cdot | r) - p_0(a)) \right\| \leq \mathbf{W}_{KQ}(r) \cdot p(r) \|\mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}] - p_0\|$$

Define $p_0 = \frac{1}{|\mathcal{A}|} \mathbf{1}_{\mathcal{A}}$. Observe that

$$\begin{aligned} \|\mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}] - p_0\| &\leq \mathbb{E}_{z_{1:T}}[\|\hat{p}(a | z_{1:T}) - p_0\| | r \in z_{1:T}] \\ &\leq \mathbb{E}_{z_{1:T}} \left[\sum_t \mathbf{W}_{KQ}(z_t) \|\mathbf{W}_{OV}(\cdot, z_t)\| | r \in z_{1:T} \right] \\ &\leq \mathbf{W}_{KQ}(r) \|\mathbf{W}_{OV}(\cdot, r)\| + T\alpha_{sm}^2 \\ &\leq \mathbf{W}_{KQ}(r)^2 + T\alpha_{sm}^2. \end{aligned}$$

Thus

$$\left\| \dot{\mathbf{W}}_{OV}(\cdot, r) - \mathbf{W}_{KQ}(r) \cdot p(r)(p^*(\cdot | r) - p_0(a)) \right\| \leq p(r) \mathbf{W}_{KQ}(r) (\mathbf{W}_{KQ}(r)^2 + T\alpha_{sm}^2)$$

Likewise,

$$\dot{\mathbf{W}}_{KQ}(r) = p(r) \langle \mathbf{W}_{OV}(\cdot, r), (p^*(\cdot | r) - \mathbb{E}_{z_{1:T}}[\hat{p}(a | z_{1:T}) | r \in z_{1:T}]) \rangle,$$

and thus

$$\begin{aligned} \left| \dot{\mathbf{W}}_{KQ}(r) - p(r) \langle \mathbf{W}_{OV}(\cdot, r), (p^*(\cdot | r) - p_0) \rangle \right| &\leq p(r) \|\mathbf{W}_{OV}(\cdot, r)\| \|\mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}] - p_0\| \\ &\leq p(r) \mathbf{W}_{KQ}(r) (\mathbf{W}_{KQ}(r)^2 + T\alpha_{sm}^2) \end{aligned}$$

Define the vector $\mathbf{u} \in \mathbb{R}^2$ by

$$\mathbf{u} = \begin{bmatrix} \mathbf{W}_{KQ}(r) \\ \langle \mathbf{W}_{OV}(\cdot, r), \frac{p^*(\cdot | r) - p_0}{\|p^*(\cdot | r) - p_0\|} \rangle \end{bmatrix}$$

We see that

$$\left\| \dot{\mathbf{u}} - p(r) \|p^*(\cdot | r) - p_0\| \cdot \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \mathbf{u} \right\| \leq 2p(r) \mathbf{W}_{KQ}(r) (\mathbf{W}_{KQ}(r)^2 + T\alpha_{sm}^2)$$

Therefore

$$\begin{aligned} \frac{d}{dt} (\|\mathbf{u}\|^2) &\leq 2 \langle \dot{\mathbf{u}}, \mathbf{u} \rangle \\ &\leq 2p(r) \|p^*(\cdot | r) - p_0\| \|\mathbf{u}\|^2 + 4p(r) \|\mathbf{u}\| \mathbf{W}_{KQ}(r) (\mathbf{W}_{KQ}(r)^2 + T\alpha_{sm}^2) \\ &\leq 2p(r) \|p^*(\cdot | r) - p_0\| \|\mathbf{u}\|^2 + 4p(r) \|\mathbf{u}\|^2 (\|\mathbf{u}\|^2 + T\alpha_{sm}^2) \\ &\leq 2p(r) (\|p^*(\cdot | r) - p_0\| + 2T\alpha_{sm}^2) \|\mathbf{u}\|^2 + 4p(r) \|\mathbf{u}\|^4. \end{aligned}$$

where the last inequality bounds $\mathbf{W}_{KQ}^2(r) \leq \|\mathbf{u}\|^2$.

Define $\gamma_r := 2p(r)(\|p^*(\cdot | r) - p_0\| + 2T\alpha_{sm}^2)$. By Lemma 7 we have that for

$$t < \gamma_r^{-1} \log \left(\frac{\gamma_r}{4p(r)\|\mathbf{u}_0\|^2} + 1 \right),$$

$$\|\mathbf{u}\|^2 \leq \frac{\gamma_r \|\mathbf{u}_0\|^2 \exp(\gamma_r t)}{\gamma_r + 4p(r)^2(1 - \exp(\gamma_r t))}$$

Let T_ϵ be the first time that $\|\mathbf{u}\| \geq \epsilon$. If $T_\epsilon < \gamma_r^{-1} \log \left(\frac{\gamma_r}{4p(r)\|\mathbf{u}_0\|^2} + 1 \right)$, then

$$\epsilon^2 \leq \|\mathbf{u}\|^2 \leq \frac{\gamma_r \|\mathbf{u}_0\|^2 \exp(\gamma_r T_\epsilon)}{\gamma_r + 4p(r)^2(1 - \exp(\gamma_r T_\epsilon))} \leq \frac{\gamma_r \alpha^2(|\mathcal{A}| + 1) \exp(\gamma_r T_\epsilon)}{\gamma_r + 4p(r)^2(1 - \exp(\gamma_r T_\epsilon))}.$$

Therefore

$$T_\epsilon \geq \gamma_r^{-1} \log \left(\frac{\epsilon^2 \gamma_r + 4p(r)\epsilon^2 \alpha^2(|\mathcal{A}| + 1)}{\alpha^2(|\mathcal{A}| + 1) \gamma_r + 4p(r)\epsilon^2 \alpha^2(|\mathcal{A}| + 1)} \right) \geq \gamma_r^{-1} \log \left(\frac{\epsilon^2}{2\alpha^2(|\mathcal{A}| + 1)} \right)$$

for $\epsilon^2 \leq \frac{\gamma_r}{4p(r)}$. On this assumption, $\frac{\epsilon^2}{2\alpha^2(|\mathcal{A}| + 1)} \leq \frac{\gamma_r}{4p(r)\|\mathbf{u}_0\|^2}$, and thus we always have $T_\epsilon \geq \gamma_r^{-1} \log \left(\frac{\epsilon^2}{2\alpha^2(|\mathcal{A}| + 1)} \right)$.

Define L_r by

$$L_r(\boldsymbol{\theta}) := p(r) \mathbb{E}_{z_{1:T+1}} \left[-\langle \boldsymbol{\varphi}(z_{T+1}), F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) \rangle + \log \left(\sum_{a \in \mathcal{A}} \exp(\langle \boldsymbol{\varphi}(a), F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta}) \rangle) \right) \mid r \in z_{1:T} \right]$$

Let us define the relation-only model as

$$\boldsymbol{\varphi}(a)^\top F_{\text{rel}}(\mathbf{X}; \boldsymbol{\theta}) = \mathbf{W}_{OV}(a, r) \mathbf{W}_{KQ}(r)$$

where $r \in z_{1:T}$. We see that

$$|\boldsymbol{\varphi}(a)^\top F_{\text{rel}}(\mathbf{X}; \boldsymbol{\theta}) - \boldsymbol{\varphi}(a)^\top F_{\text{lin}}(\mathbf{X}; \boldsymbol{\theta})| \leq (T-1)\alpha_{sm}^2.$$

Define $g : \mathbb{R}^{|\mathcal{A}|} \rightarrow \mathbb{R}$ by $g(\mathbf{z}) = \log(\sum_a \exp(z_a))$. We see that $\nabla_{\mathbf{z}} g(\mathbf{z}) = \mathcal{S}(\mathbf{z})$, where \mathcal{S} is the softmax, and thus

$$|g(\mathbf{z}_1) - g(\mathbf{z}_2)| \leq \sup_{\mathbf{z}} \|\nabla_{\mathbf{z}} g(\mathbf{z})\|_1 \cdot \|g(\mathbf{z}_1) - g(\mathbf{z}_2)\|_\infty \leq \|g(\mathbf{z}_1) - g(\mathbf{z}_2)\|_\infty.$$

Therefore defining the relation-only loss \bar{L}_r as

$$\begin{aligned} \bar{L}_r(\boldsymbol{\theta}) &:= p(r) \mathbb{E}_s \left[-\langle \boldsymbol{\varphi}(a^*(s, r)), F_{\text{rel}}(r; \boldsymbol{\theta}) \rangle + \log \left(\sum_{a \in \mathcal{A}} \exp(\langle \boldsymbol{\varphi}(a), F_{\text{lin}}(r; \boldsymbol{\theta}) \rangle) \right) \right] \\ &= -p(r) \sum_a p(a | r) \mathbf{W}_{OV}(a^*(s, r), r) \mathbf{W}_{KQ}(r) + p(r) \log \left(\sum_{a \in \mathcal{A}} \exp(\mathbf{W}_{OV}(a, r) \mathbf{W}_{KQ}(r)) \right), \end{aligned}$$

we see that

$$|L_r(\boldsymbol{\theta}) - \bar{L}_r(\boldsymbol{\theta})| \leq 2(T-1)\alpha_{sm}^2.$$

Since log-sum-exp is 1-strongly-convex, recalling that $\boldsymbol{\Theta}(a, r) := \mathbf{W}_{OV}(a, r)\mathbf{W}_{KQ}(r)$,

$$\log \left(\sum_a \exp(\boldsymbol{\Theta}(a, r)) \right) \leq \log(|\mathcal{A}|) + \sum_a \frac{1}{|\mathcal{A}|} \boldsymbol{\Theta}(a, r) + \frac{1}{2} \|\boldsymbol{\Theta}(\cdot, r)\|^2.$$

Therefore

$$\begin{aligned} \bar{L}_r &\leq p(r) \log |\mathcal{A}| - p(r) \langle \boldsymbol{\Theta}(\cdot, r), p^*(\cdot | r) - p_0 \rangle + \frac{1}{2} p(r) \|\boldsymbol{\Theta}(\cdot, r)\|^2 \\ &= L_{r,0} - p(r) \|p^*(\cdot | r) - p_0\| \cdot \mathbf{u}_1 \mathbf{u}_2 + \frac{1}{2} p(r) \|\boldsymbol{\Theta}(\cdot, r)\|^2. \end{aligned}$$

We next track the evolution of $\mathbf{u}_1 \mathbf{u}_2$:

$$\begin{aligned} \frac{d}{dt}(\mathbf{u}_1 \mathbf{u}_2) &= \dot{\mathbf{u}}_1 \mathbf{u}_2 + \mathbf{u}_1 \dot{\mathbf{u}}_2 \\ &\geq p(r) \|p^*(\cdot | r) - p_0\| \|\mathbf{u}\|^2 - 4p(r) \|\mathbf{u}\|^2 (\|\mathbf{u}\|^2 + (T-1)\alpha_{sm}^2) \\ &\geq p(r) (\|p^*(\cdot | r) - p_0\| - 4\|\mathbf{u}\|^2 - 4(T-1)\alpha_{sm}^2) \|\mathbf{u}\|^2 \\ &\geq p(r) (\|p^*(\cdot | r) - p_0\| - 4T\alpha_{sm}^2) \|\mathbf{u}\|^2 - 4p(r) \|\mathbf{u}\|^4. \end{aligned}$$

for $t \leq T_\epsilon$. Since $\|\mathbf{u}\| \leq \epsilon$, this is increasing in $\|\mathbf{u}\|$.

We first have the bound $\|\mathbf{u}\|^2 \geq \mathbf{W}_{KQ}(r)^2 \geq \alpha^2$. Next, we have the bound $\|\mathbf{u}\|^2 \geq 2\mathbf{u}_1 \mathbf{u}_2$. Pick some time $\tau \leq T_\epsilon$. Define $\gamma_r^- := 2p(r) (\|p^*(\cdot | r) - p_0\| - 4T\alpha_{sm}^2)$. We see that

$$(\mathbf{u}_1 \mathbf{u}_2)(\tau) \geq \left(\frac{1}{2} \gamma_r^- \alpha^2 - 4p(r) \alpha^4 \right) \tau \geq \frac{1}{4} \gamma_r^- \alpha^2 \tau$$

Next, by Lemma 7, for $t \leq T_\epsilon$ we have

$$(\mathbf{u}_1 \mathbf{u}_2)(t) \geq \frac{\gamma_r^- (\mathbf{u}_1 \mathbf{u}_2)(\tau) \exp(\gamma_r^- (t - \tau))}{\gamma_r^- + 8p(r) (\mathbf{u}_1 \mathbf{u}_2)(\tau) \exp(\gamma_r^- (t - \tau))}.$$

Plugging in $t = \gamma_r^{-1} \log \left(\frac{\epsilon^2}{2\alpha^2(|\mathcal{A}|+1)} \right)$,

$$(\mathbf{u}_1 \mathbf{u}_2)(\tau) \exp(\gamma_r^- (t - \tau)) \geq \frac{1}{4} \gamma_r^- \alpha^2 \tau \exp(-\gamma_r^- \tau) \cdot \exp \left(\frac{\gamma_r^-}{\gamma_r} \log \left(\frac{\epsilon^2}{2\alpha^2(|\mathcal{A}|+1)} \right) \right)$$

Selecting $\tau = 1/\gamma_r^-$, we get

$$\begin{aligned}
(\mathbf{u}_1 \mathbf{u}_2)(\tau) \exp(\gamma_r^-(t - \tau)) &\geq \frac{\alpha^2}{4e} \cdot \exp\left(\frac{\|p^*(\cdot | r) - p_0\| - 4T\alpha_{sm}^2 \log\left(\frac{\epsilon^2}{2\alpha^2(|\mathcal{A}| + 1)}\right)}{\|p^*(\cdot | r) - p_0\| + 2T\alpha_{sm}^2}\right) \\
&\geq \frac{\alpha^2}{4e} \cdot \frac{\epsilon^2}{2\alpha^2(|\mathcal{A}| + 1)} \cdot \exp\left(\frac{-6T\alpha_{sm}^2}{\|p^*(\cdot | r) - p_0\| + 2T\alpha_{sm}^2} \log\left(\frac{\epsilon^2}{2\alpha^2(|\mathcal{A}| + 1)}\right)\right) \\
&\geq \frac{\alpha^2}{4e} \cdot \frac{\epsilon^2}{2\alpha^2(|\mathcal{A}| + 1)} \cdot \left(1 - \frac{6T\alpha_{sm}^2}{\|p^*(\cdot | r) - p_0\| + 2T\alpha_{sm}^2} \log\left(\frac{\epsilon^2}{2\alpha^2(|\mathcal{A}| + 1)}\right)\right) \\
&\geq \frac{\epsilon^2}{50(|\mathcal{A}| + 1)}.
\end{aligned}$$

whenever $\frac{T\alpha_{sm}^2}{\|p^*(\cdot | r) - p_0\| + 2T\alpha_{sm}^2} \log\left(\frac{\epsilon^2}{\alpha^2(|\mathcal{A}| + 1)}\right) \leq \frac{1}{150}$.

Therefore

$$\begin{aligned}
(\mathbf{u}_1 \mathbf{u}_2)(t) &\geq \frac{\epsilon^2}{50(|\mathcal{A}| + 1)} \cdot \frac{1}{1 + 8\frac{p(r)}{\gamma_r^-} \frac{\epsilon^2}{50(|\mathcal{A}| + 1)}} \\
&\geq \frac{\epsilon^2}{50(|\mathcal{A}| + 1)} \cdot \frac{1}{1 + \frac{4\epsilon^2}{\|p^*(\cdot | r) - p_0\| 50(|\mathcal{A}| + 1)}} \\
&\geq \frac{\epsilon^2}{100(|\mathcal{A}| + 1)},
\end{aligned}$$

Altogether, we get that the loss is

$$\begin{aligned}
\bar{L}_r &\leq L_{r,0} - p(r) \|p^*(\cdot | r) - p_0\| \cdot \frac{\epsilon^2}{100(|\mathcal{A}| + 1)} + \frac{1}{2} p(r) \epsilon^4 \\
&\leq L_{r,0} - p(r) \|p^*(\cdot | r) - p_0\| \cdot \frac{\epsilon^2}{200(|\mathcal{A}| + 1)}
\end{aligned}$$

whenever $\epsilon^2 \leq \frac{\|p^*(\cdot | r) - p_0\|}{100(|\mathcal{A}| + 1)}$.

Stage 2: Norm stays large Next, we want to show that $\mathbf{W}_{KQ}(r)$ stays large. We first show that the relation-only loss \bar{L} is decreasing. We can compute that

$$\frac{d}{dt} \bar{L}_r(\boldsymbol{\theta}) = \langle \nabla_{\boldsymbol{\theta}} \bar{L}_r, \nabla_{\boldsymbol{\theta}} L_r \rangle$$

Define $\hat{p}(\cdot | r)$ by

$$\hat{p}(a | r) = \frac{\exp(\mathbf{W}_{KQ}(r) \mathbf{W}_{OV}(a, r))}{\sum_{a'} \exp(\mathbf{W}_{KQ}(r) \mathbf{W}_{OV}(a', r))}.$$

We observe that

$$\begin{aligned}
\partial_{\mathbf{W}_{OV}(\cdot, r)} \bar{L}_r &= \mathbf{W}_{KQ}(r) p(r) (p^*(\cdot | r) - \hat{p}(\cdot | r)) \\
\partial_{\mathbf{W}_{OV}(\cdot, r)} L_r &= \mathbf{W}_{KQ}(r) p(r) (p^*(\cdot | r) - \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}])
\end{aligned}$$

and thus

$$\begin{aligned} \|\partial_{\mathbf{W}_{OV}(\cdot, r)} \bar{L}_r - \partial_{\mathbf{W}_{OV}(\cdot, r)} L_r\| &\leq \mathbf{W}_{KQ}(r)p(r) \|\mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}] - \hat{p}(\cdot | r)\| \\ &\leq \mathbf{W}_{KQ}(r)p(r)T\alpha_{sm}^2. \end{aligned}$$

Likewise,

$$\begin{aligned} |\partial_{\mathbf{W}_{KQ}(r)} \bar{L}_r - \partial_{\mathbf{W}_{KQ}(r)} L_r| &= p(r) |\langle \mathbf{W}_{OV}(\cdot, r), \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}] - \hat{p}(\cdot | r) \rangle| \\ &\leq p(r) \|\mathbf{W}_{OV}(\cdot, r)\| \|\mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}] - \hat{p}(\cdot | r)\| \\ &\leq \mathbf{W}_{KQ}(r)p(r)T\alpha_{sm}^2. \end{aligned}$$

Therefore

$$\begin{aligned} \frac{d}{dt} L_r(\boldsymbol{\theta}) &\geq \|\nabla_{\boldsymbol{\theta}} \bar{L}_r\|^2 - \|\nabla_{\boldsymbol{\theta}} \bar{L}_r\| \|\nabla_{\boldsymbol{\theta}} \bar{L}_r - \nabla_{\boldsymbol{\theta}} L_r\| \\ &\geq \|\nabla_{\boldsymbol{\theta}} \bar{L}_r\|^2 - \|\nabla_{\boldsymbol{\theta}} \bar{L}_r\| \sqrt{2} \mathbf{W}_{KQ}(r)p(r)T\alpha_{sm}^2. \end{aligned}$$

Assume that $\frac{d}{dt} L_r(\boldsymbol{\theta}) < 0$. Then

$$\sqrt{2} \mathbf{W}_{KQ}(r)p(r)T\alpha_{sm}^2 \geq \|\nabla_{\boldsymbol{\theta}} \bar{L}_r\| \geq \mathbf{W}_{KQ}(r)p(r) \|p^*(\cdot | r) - \hat{p}(\cdot | r)\|,$$

i.e

$$\|p^*(\cdot | r) - \hat{p}(\cdot | r)\| \leq \sqrt{2}T\alpha_{sm}^2.$$

Assuming that $\sqrt{2}T\alpha_{sm}^2 < \frac{1}{2S}$, since $p^*(a | r) > \frac{1}{S}$ for $p^*(a | r) > 0$ we have that

$$\left| p^*(a | r) \log \frac{p^*(a | r)}{\hat{p}(a | r)} \right| = p^*(a | r) \left| \log \left(1 + \frac{\hat{p}(a | r) - p^*(a | r)}{p^*(a | r)} \right) \right| \leq 2|\hat{p}(a | r) - p^*(a | r)|$$

Therefore $\bar{L}_r - p(r)H(p^*(\cdot | r)) \leq 2p(r)\|\hat{p}(\cdot | r) - p^*(\cdot | r)\|_1 \leq \sqrt{2}T\alpha_{sm}^2 D$. As such, we have that \bar{L}_r stays below $L_{r,0} - p(r)\|p^*(\cdot | r) - p_0\| \cdot \frac{\epsilon^2}{200(|\mathcal{A}|+1)}$ for the remainder of the gradient flow trajectory.

By convexity in Θ space,

$$\begin{aligned} L_{r,0} - \bar{L}_r(\Theta) &\leq -\langle \nabla_{\Theta(\cdot, r)} \bar{L}(\mathbf{0}), \Theta(\cdot, r) \rangle \\ &\leq \|\nabla_{\Theta(\cdot, r)} \bar{L}(\mathbf{0})\| \|\Theta(\cdot, r)\| \\ &\leq p(r) \cdot \|p^*(\cdot | r) - p_0\| \cdot \|\Theta(\cdot, r)\|. \end{aligned}$$

Therefore

$$\|\Theta(\cdot, r)\| \geq \frac{\epsilon^2}{100(|\mathcal{A}|+1)}$$

Stage 3: Convergence. Next, we can bound the loss decrease by

$$\begin{aligned}
\frac{d}{dt}L(\boldsymbol{\theta}) &= -\|\nabla_{\boldsymbol{\theta}}L\|^2 \\
&\leq -\sum_r \|\partial_{\mathbf{w}_{OV}(\cdot, r)}L\|^2 \\
&= -\sum_r \|\partial_{\mathbf{w}_{OV}(\cdot, r)}L_r\|^2 \\
&= -\sum_r \mathbf{W}_{KQ}(r)^2 p(r)^2 \|p^*(\cdot | r) - \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}]\| \\
&\leq -\sum_r \|\boldsymbol{\Theta}(\cdot, r)\| p(r)^2 \|p^*(\cdot | r) - \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}]\|^2 \\
&\leq -\frac{\epsilon^2}{200(|\mathcal{A}| + 1)} \sum_r p(r)^2 \|p^*(\cdot | r) - \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}]\|^2
\end{aligned}$$

Since $L(\boldsymbol{\Theta}) \leq L(\boldsymbol{\Theta}_0) = \log |\mathcal{A}|$, and $L(\boldsymbol{\Theta}) \geq 0$, there exists some $t \leq \max_r t^* + 100(|\mathcal{A}| + 1) \log |\mathcal{A}| \epsilon^{-2} \epsilon_{min}^{-2}$ such that

$$\sum_r p(r)^2 \|p^*(\cdot | r) - \mathbb{E}_{z_{1:T}}[\hat{p}(\cdot | z_{1:T}) | r \in z_{1:T}]\|^2 \leq \epsilon_{min}^2,$$

as desired. □

To conclude, we must set α, α_{sm}, T^* appropriately in terms of ϵ in order to apply Lemma 6.

Proof of Theorem 6. Let $\epsilon = \epsilon_{min} \leq \frac{1}{100(|\mathcal{A}|+1)} \cdot \min_r \|p^*(\cdot | r) - p_0\|$ be the target accuracy. Let us choose the initialization α so that $\log\left(\frac{\epsilon}{\alpha\sqrt{|\mathcal{A}|+1}}\right) = \iota$, where ι is chosen so that

$$\iota \geq 100(|\mathcal{A}| + 1) \log |\mathcal{A}| \epsilon^{-4} p(r) \|p^*(\cdot | r) - p_0\|.$$

In this case, we see that

$$T^* \leq 2\iota \max_r (p(r)^{-1} \|p^*(\cdot | r) - p_0\|^{-1}).$$

Since $p^*(\cdot | r)$ is supported on at most D elements, and $|\mathcal{A}| \geq 2D$, we have $\|p^*(\cdot | r) - p_0\|^{-1} \leq \sqrt{2D}$. Therefore $T^* \leq 2R\sqrt{D} \cdot \iota$. Let us compute α_{sm} . We see that, since $S \geq 8R\sqrt{2D}$,

$$\begin{aligned}
\mathbf{W}_{KQ}(s) &\leq \exp\left(\frac{4R\sqrt{D}}{S}\iota\right) \cdot \alpha\sqrt{|\mathcal{A}| + 1} \\
&\leq \exp\left(\frac{1}{2}\iota\right) \alpha\sqrt{|\mathcal{A}| + 1} \\
&= \sqrt{\epsilon\alpha} \cdot (|\mathcal{A}| + 1)^{1/4} \\
&\leq \sqrt{\alpha}
\end{aligned}$$

Similarly, since $\mathcal{N} \geq 4R\sqrt{2DT}$,

$$\begin{aligned} \mathbf{W}_{KQ}(s) &\leq \exp\left(\frac{4R\sqrt{DT}}{|\mathcal{N}|}\iota\right) \cdot \alpha\sqrt{|A|+1} \\ &\leq \exp\left(\frac{1}{2}\iota\right) \alpha\sqrt{|A|+1} \\ &\leq \sqrt{\alpha} \end{aligned}$$

Therefore the assumption holds for $\alpha_{sm} = \sqrt{\alpha}$. To conclude, we must verify that

$$T\alpha \leq \frac{1}{300}\iota^{-1} \cdot \min_r \|p^*(\cdot | r) - p_0\|.$$

But since $\alpha = \frac{\epsilon}{\sqrt{|A|+1}}e^{-\iota}$, the RHS scales with $e^{-\iota}$ and the RHS scales with ι , and thus the condition can be obtained for choosing ι sufficiently large.

Under the setting of parameters the conditions of Lemma 6 are satisfied, and thus the claim holds. \square

D.4 Helper Lemma

Lemma 7. *Let $z(t) \geq 0$ satisfy*

$$\dot{z} \leq Az + Bz^2.$$

for positive constants A, B . Then

$$z(t) \leq \frac{Az(0)e^{At}}{A + Bz(0)(1 - e^{At})}.$$

Furthermore, if

$$\dot{z} \geq Az - Bz^2,$$

and $z \in [0, \frac{A}{2B}]$ on the interval $[0, T]$, then

$$z(t) \geq \frac{Az(0)e^{At}}{A + Bz(0)(e^{At} - 1)}.$$

for all $t \in [0, T]$.

Both claims follow from the Bihari-LaSalle inequality.

E Proofs from Section 6

E.1 Associative Memories

Proof of Theorem 7. $f^* \rightarrow \mathbf{F} \rightarrow \hat{f}$ is a Markov chain, so by the data processing inequality,

$$I(f^*; \hat{f}) \leq I(f^*; \mathbf{F}).$$

Also, by definition of mutual information

$$I(f^*; \mathbf{F}) \leq H(\mathbf{F}) \leq B,$$

where the last inequality follows since \mathbf{F} is an B -bit message. Thus $I(f^*; \hat{f}) \leq B$.

Let $q_x(\cdot | \hat{f})$ be the conditional distribution of $f^*(x)$ given \hat{f} . Consider some fixed \hat{f} . $\hat{f}(x)$ is also a probability distribution over $[M]$, and thus by Gibbs' inequality

$$\mathbb{E}_{y \sim q_x(\cdot | \hat{f})} \left[-\log \hat{f}(x)_y \right] \geq \mathbb{E}_{y \sim q_x(\cdot | \hat{f})} \left[-\log q_x(y | \hat{f}) \right].$$

Therefore, letting q be the marginal distribution over \hat{f} and q_x the marginal over $f^*(x)$,

$$\begin{aligned} \mathbb{E}_{f^*, \hat{f}} \left[-\log \hat{f}(x)_{f^*(x)} \right] &= \mathbb{E}_{\hat{f}} \left[\mathbb{E}_{y \sim q_x(\cdot | \hat{f})} \left[-\log \hat{f}(x)_y \right] \right] \\ &\geq \mathbb{E}_{\hat{f}} \left[\mathbb{E}_{y \sim q_x(\cdot | \hat{f})} \left[-\log q_x(y | \hat{f}) \right] \right] \\ &= \mathbb{E}_{f^*, \hat{f}} \left[-\log q_x(f^*(x) | \hat{f}) \right] \\ &= \mathbb{E}_{f^*, \hat{f}} \left[-\log \frac{q_x(f^*(x), \hat{f})}{q(\hat{f})q_x(f^*(x))} - \log q_x(f^*(x)) \right] \\ &= -I(f^*(x); \hat{f}) + \log M. \end{aligned}$$

where in the last step we use the fact that q_x is uniform over $[M]$, and plug in the definition of mutual information. The total loss is thus

$$\mathbb{E}_{f^*, \hat{f}} \left[L(\hat{f}) \right] \geq \sum_{x \in [N]} p(x) \left(-I(f^*(x); \hat{f}) + \log M \right). \quad (22)$$

Since the y_i are independent,

$$B \geq I(f^*; \hat{f}) \geq \sum_{x \in [N]} I(f^*(x); \hat{f}).$$

Also, $0 \leq I(f^*(x); \hat{f}) \leq H(f^*(x)) = \log M$. Therefore equation 22 is minimized when $I(f^*(x); \hat{f}) = \log M$ for the $B/\log M$ most frequent tokens. Altogether,

$$\mathbb{E}_{f^*, \hat{f}} \left[L(\hat{f}) \right] \geq \log M \cdot \sum_{x > \lceil \frac{B}{\log M} \rceil} p(x).$$

□

Proof of Corollary 1. Let $p(x) = Z_\alpha x^{-\alpha}$, where $Z_\alpha = \sum_{x \in [N]} x^{-\alpha}$. We can bound

$$\sum_{x=k}^n p(x) = \frac{\sum_{x=k}^n x^{-\alpha}}{\sum_{x=1}^n x^{-\alpha}} \asymp k^{1-\alpha}.$$

Therefore

$$\mathbb{E}_{f^*, \hat{f}} \left[L(\hat{f}) \right] \geq \log M \cdot \sum_{x > \lceil \frac{B}{\log M} \rceil} p(x) \gtrsim \log M \left(\frac{B}{\log M} \right)^{1-\alpha} \gtrsim B^{1-\alpha}.$$

□

E.2 Factual Recall

Proof. Define $\ell(s, r) := \mathbb{E}_{\mathcal{D}, \hat{f}} \left[-\log \hat{f}(s, r)_{a^*(s, r)} \right]$ so that

$$\mathbf{L} = p(s, r) \cdot \ell(s, r).$$

Let us define the expanded dataset $\mathcal{D} := \{\mathcal{A}_r\}_{r \in \mathcal{R}} \cup \{a^*(s, r)\}_{s \in \mathcal{S}, r \in \mathcal{R}}$. We observe that $\mathcal{D} \rightarrow a^* \rightarrow \mathbf{F} \rightarrow \hat{f}$ is a Markov chain, and thus by the data processing inequality

$$B \geq I(\mathcal{D}; \hat{f}).$$

Next, by the chain rule, we can decompose

$$\begin{aligned} I(\mathcal{D}; \hat{f}) &= I(\mathcal{A}_1, \dots, \mathcal{A}_R; \hat{f}) + I(a^*; \hat{f} \mid \mathcal{A}_1, \dots, \mathcal{A}_R) \\ &\geq I(\mathcal{A}_1, \dots, \mathcal{A}_R; \hat{f}) + \sum_{s, r} I(a^*(s, r); \hat{f} \mid \mathcal{A}_1, \dots, \mathcal{A}_R) \\ &= I(\mathcal{A}_1, \dots, \mathcal{A}_R; \hat{f}) + \sum_{s, r} I(a^*(s, r); \hat{f} \mid \mathcal{A}_r), \end{aligned}$$

where the first inequality uses the fact that the $a^*(s, r)$ are conditionally independent given the \mathcal{A}_r , and the second uses that $a^*(s, r)$ is independent of $\mathcal{A}_{r'}$ given \mathcal{A}_r , for $r \neq r'$.

We can decompose the first mutual information term, using the fact that the \mathcal{A}_r are nearly independent:

Lemma 8. *Assume that $|\mathcal{V}| \geq 2RD$. Then*

$$I(\mathcal{A}_1, \dots, \mathcal{A}_R; \hat{f}) \geq \sum_r I(\mathcal{A}_r; \hat{f}) - \frac{2R^2D^2}{|\mathcal{V}|}.$$

We next relate $I(\mathcal{A}_r; \hat{f})$ to the loss. The intuition for this lemma is that for a fixed r the quantity $\sum_s \ell(s, r)$ is small, then the predictor \hat{f} must contain information about the answer set \mathcal{A}_r .

Lemma 9. *Assume that $|\mathcal{V}| \geq 2D$. Define $\eta := C\sqrt{\frac{D}{S} \log(2D^2 \log |\mathcal{V}|)}$ for a sufficiently large constant C , and assume that $\eta \leq 1$. Then*

$$I(\mathcal{A}_r; \hat{f}) \geq -(1 + \eta) \frac{D}{S} \cdot \sum_{s \in [S]} \ell(s, r) + D \log \frac{|\mathcal{V}|}{D} - \underbrace{\frac{2D \log |\mathcal{V}|}{|\mathcal{V}|} - \frac{2D^2}{|\mathcal{V}|}}_{\text{lower order term}} - \eta D - 1$$

Finally, we relate $I(a^*(s, r); \hat{f} \mid \mathcal{A}_r)$ to the loss. Similarly, the intuition for this lemma is that if the loss $\ell(s, r)$ is small, then \hat{f} must contain information about the true association $a^*(s, r)$.

Lemma 10. *For all s, r ,*

$$I(a^*(s, r); \hat{f} \mid \mathcal{A}_r) \geq \log D - \ell(s, r).$$

The proofs for Lemmas 8 to 10 are deferred to Appendix E.3.

Combining Lemmas 8 to 10, we get

$$\begin{aligned}
B &\geq I(\mathcal{A}_1, \dots, \mathcal{A}_R; \hat{f}) + \sum_{s,r} I(a^*(s, r); \hat{f} \mid \mathcal{A}_r) \\
&= -(1 + \eta) \frac{D}{S} \sum_{s,r} \ell(s, r) + RD \log \frac{|\mathcal{V}|}{D} - \underbrace{\frac{2RD \log |\mathcal{V}|}{|\mathcal{V}|} - \frac{2RD^2}{|\mathcal{V}|} - \eta RD - R - \frac{2R^2 D^2}{|\mathcal{V}|}}_{\text{lower order term}} \\
&\quad + SR \log D - \sum_{s,r} \ell(s, r) \\
&= - \left((1 + \eta) \frac{D}{S} + 1 \right) \sum_{s,r} \ell(s, r) + SR \log D + RD \log \frac{|\mathcal{V}|}{D} - \varepsilon_{lot},
\end{aligned}$$

where $\varepsilon_{lot} := \frac{2RD \log |\mathcal{V}|}{|\mathcal{V}|} + \frac{2RD^2}{|\mathcal{V}|} + \eta RD + R + \frac{2R^2 D^2}{|\mathcal{V}|} \ll RD \log \frac{|\mathcal{V}|}{D}$ is a lower order term.

Altogether, we see that in order for all the losses $\ell(s, r)$ to equal zero, we require

$$B \geq SR \log D + RD \log \frac{|\mathcal{V}|}{D} - \varepsilon_{lot} \geq SR \log D + (1 - c)RD \log \frac{|\mathcal{V}|}{D}$$

Furthermore, when $p(s, r) = \frac{1}{RS}$, then $\mathbf{L} = \frac{1}{SR} \sum_{s,r} \ell(s, r)$, and the bound becomes

$$\begin{aligned}
B &\geq -((1 + \eta)RD + RS) \cdot \mathbf{L} + SR \log D + RD \log \frac{|\mathcal{V}|}{D} - \varepsilon_{lot} \\
&\quad - ((1 + c)RD + RS) \cdot \mathbf{L} + SR \log D + (1 - c)RD \log \frac{|\mathcal{V}|}{D}.
\end{aligned}$$

□

E.3 Auxiliary Lemmas

Lemma 11. For random variables X, Y, Z ,

$$I(X, Y; Z) \geq I(X; Z) + I(Y; Z) - I(X; Y)$$

Proof. By standard properties of mutual information:

$$\begin{aligned}
&I(X, Y; Z) - I(X; Z) - I(Y; Z) \\
&= H(X, Y) - H(X, Y \mid Z) - H(X) + H(X \mid Z) - H(Y) + H(Y \mid Z) \\
&= I(X; Y \mid Z) - I(X; Y) \\
&\geq -I(X; Y).
\end{aligned}$$

□

Proof of Lemma 8. By Lemma 11,

$$\begin{aligned}
I(\mathcal{A}_1, \dots, \mathcal{A}_R; \hat{f}) &\geq \sum_r I(\mathcal{A}_r; \hat{f}) - \sum_r I(\mathcal{A}_r; \mathcal{A}_1, \dots, \mathcal{A}_{r-1}) \\
&= \sum_r I(\mathcal{A}_r; \hat{f}) - \left(\sum_r H(\mathcal{A}_r) + H(\mathcal{A}_1, \dots, \mathcal{A}_{r-1}) - H(\mathcal{A}_1, \dots, \mathcal{A}_r) \right) \\
&= \sum_r I(\mathcal{A}_r; \hat{f}) - \sum_r H(\mathcal{A}_r) + H(\mathcal{A}_1, \dots, \mathcal{A}_R).
\end{aligned}$$

Since each \mathcal{A}_r is a uniformly random subset of \mathcal{V} , we have $H(\mathcal{A}_r) = \log \binom{|\mathcal{V}|}{D}$. Also, we can bound

$$H(\mathcal{A}_1, \dots, \mathcal{A}_R) = \log \left(\binom{|\mathcal{V}|}{D} \binom{|\mathcal{V}| - D}{D} \dots \binom{|\mathcal{V}| - (R-1)D}{D} \right).$$

Thus

$$\begin{aligned}
\sum_r H(\mathcal{A}_r) - H(\mathcal{A}_1, \dots, \mathcal{A}_R) &\leq R \log \frac{\binom{|\mathcal{V}|}{D}}{\binom{|\mathcal{V}| - (R-1)D}{D}} \\
&= R \log \frac{|\mathcal{V}|!(|\mathcal{V}| - RD)!}{(|\mathcal{V}| - D)!(|\mathcal{V}| - (R-1)D)!} \\
&\leq RD \log \frac{|\mathcal{V}|}{|\mathcal{V}| - RD} \\
&\leq \frac{2R^2 D^2}{|\mathcal{V}|},
\end{aligned}$$

where we used the bound $\log \frac{1}{1-x} \leq 2x$ on $(0, \frac{1}{2})$. Plugging in yields the desired bound. \square

Proof of Lemma 9. Let (z_1, \dots, z_D) be a random permutation of \mathcal{A}_r . We first aim to relate $I(\mathcal{A}; \hat{f})$ to $I(z_i; \hat{f})$. By the data processing inequality,

$$I(\mathcal{A}_r; \hat{f}) \geq I(z_1, \dots, z_D; \hat{f}).$$

By Lemma 11,

$$\begin{aligned}
I(z_1, \dots, z_D; \hat{f}) &\geq \sum_i I(z_i; \hat{f}) - \sum_i I(z_i; z_1, \dots, z_{i-1}) \\
&= \sum_i I(z_i; \hat{f}) - \sum_i H(z_i) + H(z_1, \dots, z_D).
\end{aligned}$$

The tuple (z_1, \dots, z_D) is chosen uniformly at random from \mathcal{V}^D , conditioned on all the z_i being distinct. Therefore $H(z_i) = \log |\mathcal{V}|$, and $H(z_1, \dots, z_D) = \log (|\mathcal{V}| \dots (|\mathcal{V}| - D + 1))$. Thus

$$\begin{aligned}
\sum_i H(z_i) - H(z_1, \dots, z_D) &= \log \left(\frac{|\mathcal{V}|^D}{|\mathcal{V}| \dots (|\mathcal{V}| - D + 1)} \right) \\
&\leq D \log \frac{|\mathcal{V}|}{|\mathcal{V}| - D} \\
&\leq \frac{2D^2}{|\mathcal{V}|}.
\end{aligned}$$

Altogether,

$$I(\mathcal{A}_r; \hat{f}) \geq \sum_i I(z_i; \hat{f}) - \frac{2D^2}{|\mathcal{V}|}.$$

Next, using the definition of mutual information and Gibbs' inequality,

$$I(z_i; \hat{f}) = \mathbb{E}_{z_i, \hat{f}} \left[\log \frac{\mathbb{P}(z_i | \hat{f})}{\mathbb{P}(z_i)} \right] \geq \mathbb{E}_{z_i, \hat{f}} \left[\log \frac{q(z_i | \hat{f})}{\mathbb{P}(z_i)} \right]$$

for any probability distribution q . Let us define q as follows. First, define $\tilde{f}(s, r) := (1-\epsilon)\hat{f}(s, r) + \frac{\epsilon}{|\mathcal{V}|} \mathbf{1} \in \Delta_{\mathcal{V}}$, for a small constant ϵ to be chosen later. Next, define

$$q(z | \hat{f}) := \frac{1}{S} \sum_s \tilde{f}(s, r)_z.$$

Plugging in, and observing that $\mathbb{P}(z_i) = \frac{1}{|\mathcal{V}|}$, we get that

$$I(z_i; \hat{f}) \geq \mathbb{E}_{z_i, \hat{f}} \left[\log \left(\frac{1}{S} \sum_s \tilde{f}(s, r)_{z_i} \right) \right] + \log |\mathcal{V}|.$$

Define $\mathcal{N}_z := \{s : a^*(s, r) = z\}$. Let \mathcal{E} be the event that $|\mathcal{N}_z| \geq M$ for all $z \in \mathcal{A}$. On the event \mathcal{E} , we can bound

$$\begin{aligned} \log \left(\frac{1}{S} \sum_s \tilde{f}(s, r)_{z_i} \right) &\geq \log \left(\frac{1}{S} \sum_{s \in \mathcal{N}_{z_i}} \tilde{f}(s, r)_{a^*(s, r)} \right) \\ &= \log \left(\frac{1}{|\mathcal{N}_{z_i}|} \sum_{s \in \mathcal{N}_{z_i}} \tilde{f}(s, r)_{a^*(s, r)} \right) + \log \frac{|\mathcal{N}_{z_i}|}{S} \\ &\geq \frac{1}{|\mathcal{N}_{z_i}|} \sum_{s \in \mathcal{N}_{z_i}} \log \tilde{f}(s, r)_{a^*(s, r)} + \log \frac{|\mathcal{N}_{z_i}|}{S} \\ &\geq \frac{1}{M} \sum_{s \in \mathcal{N}_{z_i}} \log \tilde{f}(s, r)_{a^*(s, r)} + \log \frac{M}{S}. \end{aligned}$$

Thus

$$\begin{aligned} \sum_{i \in [D]} \log \left(\frac{1}{S} \sum_s \tilde{f}(s, r)_{z_i} \right) &\geq \frac{1}{M} \sum_{i \in [D]} \sum_{s \in \mathcal{N}_{z_i}} \log \tilde{f}(s, r)_{a^*(s, r)} + D \log \frac{M}{S} \\ &= \frac{1}{M} \sum_{s \in [S]} \log \tilde{f}(s, r)_{a^*(s, r)} + D \log \frac{M}{S} \\ &\geq \frac{1-\epsilon}{M} \sum_{s \in [S]} \log \hat{f}(s, r)_{a^*(s, r)} - \frac{S\epsilon}{M} \log |\mathcal{V}| + D \log \frac{M}{S}. \end{aligned}$$

On $\bar{\mathcal{E}}$, we have the naive bound

$$\sum_{i \in [D]} \log \left(\frac{1}{S} \sum_s \tilde{f}(s, r)_{z_i} \right) \geq D \log \frac{\epsilon}{|\mathcal{V}|}.$$

Altogether, we have

$$\begin{aligned} & \sum_{i \in [D]} I(z_i; \hat{f}) \\ & \geq \mathbb{E}_{z_i, \hat{f}} \left[\sum_{i \in [D]} \log \left(\frac{1}{S} \sum_s \tilde{f}(s, r)_{z_i} \right) \right] + D \log |\mathcal{V}| \\ & \geq \mathbb{E}_{z_i, \hat{f}} \left[\mathbf{1}(\mathcal{E}) \cdot \left(\frac{1-\epsilon}{M} \sum_{s \in [S]} \log \hat{f}(s, r)_{a^*(s, r)} - \frac{S\epsilon}{M} \log |\mathcal{V}| + D \log \frac{M}{S} \right) \right] \\ & \quad + \mathbb{P}(\bar{\mathcal{E}}) \cdot D \log \frac{\epsilon}{|\mathcal{V}|} + D \log |\mathcal{V}| \\ & \geq \frac{1}{M} \mathbb{E}_{z_i, \hat{f}} \left[\sum_{s \in [S]} \log \hat{f}(s, r)_{a^*(s, r)} \right] - \frac{S\epsilon}{M} \log |\mathcal{V}| + D \log \frac{M}{S} + \mathbb{P}(\bar{\mathcal{E}}) \cdot D \log \frac{\epsilon}{|\mathcal{V}|} + D \log |\mathcal{V}| \\ & = -\frac{1}{M} \cdot \sum_{s \in [S]} \ell(s, r) - \frac{S\epsilon}{M} \log |\mathcal{V}| + D \log \frac{M}{S} + \mathbb{P}(\bar{\mathcal{E}}) \cdot D \log \frac{\epsilon}{|\mathcal{V}|} + D \log |\mathcal{V}| \end{aligned}$$

By Bernstein's inequality and a union bound (a similar such concentration argument was used in the lower bound proof in Allen-Zhu and Li [2]), there exists a constant C such that $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$ for

$$M = \frac{S}{D} - C \sqrt{\frac{S}{D} \log(D/\delta)},$$

as long as $S \geq D \log(D/\delta)$. Set $\epsilon = \frac{1}{|\mathcal{V}|}$, $\delta = \frac{1}{2D \log |\mathcal{V}|}$, and define $\eta := 2C \sqrt{\frac{D}{S} \log(2D^2 \log |\mathcal{V}|)} \leq 1$. We have that

$$\frac{M}{S} = \frac{1}{D} \left(1 - C \sqrt{\frac{D}{S} \log(D/\delta)} \right) = \frac{1}{D} \left(1 - \frac{\eta}{2} \right),$$

and thus

$$\frac{S}{M} = \frac{D}{1 - \eta/2} \leq D(1 + \eta).$$

Therefore

$$\begin{aligned} \sum_{i \in [D]} I(z_i; \hat{f}) & \geq -(1 + \eta) \frac{D}{S} \cdot \sum_{s \in [S]} \ell(s, r) - (1 + \eta) \frac{D \log |\mathcal{V}|}{|\mathcal{V}|} + D \log \frac{|\mathcal{V}|}{D} + D \log(1 - \eta/2) - 1 \\ & \geq -(1 + \eta) \frac{D}{S} \cdot \sum_{s \in [S]} \ell(s, r) + D \log \frac{|\mathcal{V}|}{D} - \frac{2D \log |\mathcal{V}|}{|\mathcal{V}|} - \eta D - 1. \end{aligned}$$

Altogether, we have

$$\begin{aligned} I(\mathcal{A}_r; \hat{f}) &\geq \sum_i I(z_i; \hat{f}) - \frac{2D^2}{|\mathcal{V}|} \\ &\geq -(1 + \eta) \frac{D}{S} \cdot \sum_{s \in [S]} \ell(s, r) + D \log \frac{|\mathcal{V}|}{D} - \frac{2D \log |\mathcal{V}|}{|\mathcal{V}|} - \frac{2D^2}{|\mathcal{V}|} - \eta D - 1, \end{aligned}$$

as desired. \square

Proof of Lemma 10. By the definition of mutual information and Gibbs' inequality,

$$\begin{aligned} I(a^*(s, r); \hat{f} \mid \mathcal{A}_r) &= \mathbb{E}_{\mathcal{A}_r} \left[\mathbb{E}_{a^*(s, r), \hat{f} \mid \mathcal{A}_r} \left[\log \frac{\mathbb{P}(a^*(s, r) \mid \hat{f}, \mathcal{A}_r)}{\mathbb{P}(a^*(s, r) \mid \mathcal{A}_r)} \right] \right] \\ &= \mathbb{E}_{\mathcal{A}_r} \left[\mathbb{E}_{a^*(s, r), \hat{f} \mid \mathcal{A}_r} \left[\log \mathbb{P}(a^*(s, r) \mid \hat{f}, \mathcal{A}_r) \right] \right] + \log D \\ &\geq \mathbb{E}_{\mathcal{A}_r} \left[\mathbb{E}_{a^*(s, r), \hat{f} \mid \mathcal{A}_r} \left[\log q(a^*(s, r) \mid \hat{f}, \mathcal{A}_r) \right] \right] + \log D \end{aligned}$$

where $q(\cdot \mid \hat{f}, \mathcal{A}_r)$ is any distribution over \mathcal{V} . Let us define q to be

$$q(a \mid \hat{f}, \mathcal{A}_r) \propto \hat{f}(s, r)_a \cdot \mathbf{1}(a \in \mathcal{A}_r)$$

Since $a^*(s, r) \in \mathcal{A}_r$ always, we have that $q(a^*(s, r) \mid \hat{f}, \mathcal{A}_r) \geq \hat{f}(s, r)_{a^*(s, r)}$, and thus

$$\begin{aligned} I(a^*(s, r); \hat{f} \mid \mathcal{A}_r) &\geq \mathbb{E}_{\mathcal{A}_r} \left[\mathbb{E}_{a^*(s, r), \hat{f} \mid \mathcal{A}_r} \left[\log \hat{f}(s, r)_{a^*(s, r)} \right] \right] + \log D \\ &= \mathbb{E}_{\mathcal{D}} \left[\log \hat{f}(s, r)_{a^*(s, r)} \right] + \log D \\ &= \log D - \ell(s, r). \end{aligned}$$

\square

F Technical Lemmas

Lemma 12. *Let u, v be drawn uniformly over the d -dimensional sphere of radius 1. Then*

$$\mathbb{E}[\langle u, v \rangle^{2p}] \leq (2p)^p d^{-p}$$

Lemma 13 (Hypercontractivity for product distributions). *Let $f : (\mathbb{S}^{d-1})^k \times \mathbb{R}^m \rightarrow \mathbb{R}$ be a polynomial of total degree at most p . Then*

$$\|f\|_{L^q(\nu_d^{\otimes k} \otimes \mu_m)} \leq (q-1)^{p/2} \|f\|_{L^2(\nu_d^{\otimes k} \otimes \mu_m)},$$

where ν_d is the uniform distribution over the sphere \mathbb{S}^{d-1} , and μ_m is the standard Gaussian in m dimensions.

Hypercontractivity for the Boolean hypercube (which implies hypercontractivity for Gaussian space) and for the sphere are consequences of Beckner [3, 4]. To show Lemma 13, one can use similar techniques to the proof of Corollary 12 in Montanaro [37].

Lemma 14. *Let $f : (\mathbb{S}^{d-1})^k \times \mathbb{R}^m \rightarrow \mathbb{R}$ be a polynomial of total degree at most p . Assume that $\mathbb{E}f \geq 0$, where the expectation is taken with respect to $\nu_d^{\otimes k} \otimes \mu_m$. Then if*

$$\frac{2^p e^{-1} \log^p(1/\delta) \text{Var}(f)}{(\mathbb{E}f)^2} \leq 1,$$

$f \leq 0$ with probability at most δ .

Proof. By Markov's inequality,

$$\begin{aligned} \mathbb{P}(f \leq 0) &\leq \mathbb{P}(|f - \mathbb{E}f| \geq \mathbb{E}f) \\ &\leq \mathbb{P}(|f - \mathbb{E}f|^q \geq (\mathbb{E}f)^q) \\ &\leq \frac{\mathbb{E}[|f - \mathbb{E}f|^q]}{(\mathbb{E}f)^q}. \end{aligned}$$

Since f is a degree p polynomial, by Lemma 13 we have that

$$\mathbb{E}[|f - \mathbb{E}f|^q]^{1/q} \leq q^{p/2} \text{Var}(f)^{1/2}.$$

Therefore

$$\mathbb{P}(f \leq 0) \leq \left(\frac{q^p \text{Var}(f)}{(\mathbb{E}f)^2} \right)^{q/2}.$$

Setting $q = 2 \log(1/\delta)$, we see that whenever

$$\frac{2^p e^{-1} \log^p(1/\delta) \text{Var}(f)}{(\mathbb{E}f)^2} \leq 1,$$

we have

$$\mathbb{P}(f \leq 0) \leq \left(\frac{q^p \text{Var}(f)}{(\mathbb{E}f)^2} \right)^{q/2} \leq \delta,$$

as desired. □

F.1 Hermite Polynomials

Let μ be the standard Gaussian in 1 dimension, and let $L^2(\mu)$ be the function space of square-integrable functions with respect to this Gaussian measure. The Hermite polynomials $\{h_k\}_{k \geq 0}$ form an orthonormal basis of $L^2(\mu)$. In particular, h_k is a degree k polynomial, satisfying

$$\langle h_i, h_k \rangle_{L^2(\mu)} = \delta_{ij}.$$

One useful property of Hermite polynomials is the following:

Lemma 15. Let $\mathbf{u}, \mathbf{w} \in \mathbb{R}^d$ with $\|\mathbf{u}\| = \|\mathbf{w}\| = 1$, and let $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$. Then

$$\mathbb{E}_{\mathbf{x}}[h_k(\langle \mathbf{u}, \mathbf{x} \rangle)h_k(\langle \mathbf{w}, \mathbf{x} \rangle)] = \langle \mathbf{u}, \mathbf{w} \rangle^k.$$

Next, let μ_d be the standard Gaussian in d dimensions. The function space $L^2(\mu_d)$ has an orthonormal basis of Hermite tensors $\{\mathbf{He}_k\}_{k \geq 0}$, where $\mathbf{He}_k : \mathbb{R}^d \rightarrow (\mathbb{R}^d)^{\otimes k}$:

Definition 1. Let the k th Hermite tensor $\mathbf{He}_k : \mathbb{R}^d \rightarrow (\mathbb{R}^d)^{\otimes k}$ be defined as

$$\mathbf{He}_k(\mathbf{x}) = (-1)^k \frac{\nabla^k \mu_d(\mathbf{x})}{\mu_d(\mathbf{x})},$$

where $\mu_d(\mathbf{x}) = (2\pi)^{-d/2} \exp(-\frac{1}{2}\|\mathbf{x}\|^2)$ is the Gaussian density. We remark that each entry of $\mathbf{He}_k(\mathbf{x})$ is a degree k polynomial in \mathbf{x} , and

The Hermite tensors satisfy the following useful properties:

Lemma 16 (Properties of Hermite Tensors).

- (Connection to Hermite Polynomials) If $\mathbf{w} \in \mathbb{R}^d$, $\|\mathbf{w}\| = 1$, then

$$h_k(\langle \mathbf{w}, \mathbf{x} \rangle) = \langle \mathbf{He}_k(\mathbf{x}), \mathbf{w}^{\otimes k} \rangle$$

- (Stein's Lemma) For $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$, $f \in L^2(\mu_d)$,

$$\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})\mathbf{He}_k(\mathbf{x})] = \mathbb{E}_{\mathbf{x}}[\nabla^k f(\mathbf{x})].$$