

TRANSFORMERS AS MEASURE-THEORETIC ASSOCIATIVE MEMORY: A STATISTICAL PERSPECTIVE AND MINIMAX OPTIMALITY

Ryotaro Kawata^{1,2,*}, Taiji Suzuki^{1,2,§}

¹Department of Mathematical Informatics, University of Tokyo, Japan

²Center for Advanced Intelligence Project, RIKEN, Japan

*kawata-ryotaro725@g.ecc.u-tokyo.ac.jp

§taiji@mist.i.u-tokyo.ac.jp

ABSTRACT

Transformers excel through content-addressable retrieval and the ability to exploit contexts of, in principle, unbounded length. We recast associative memory at the level of probability measures, treating a context as a distribution over tokens and viewing attention as an integral operator on measures. Concretely, for mixture contexts $\nu = I^{-1} \sum_{i=1}^I \mu^{(i^*)}$ and a query $x_q(i^*)$, the task decomposes into (i) recall of the relevant component $\mu^{(i^*)}$ and (ii) prediction from (μ_{i^*}, x_q) . We study learned softmax attention (not a frozen kernel) trained by empirical risk minimization and show that a shallow measure-theoretic Transformer composed with an MLP learns the recall-and-predict map under a spectral assumption on the input densities. We further establish a matching minimax lower bound with the same rate exponent (up to multiplicative constants), proving sharpness of the convergence order. The framework offers a principled recipe for designing and analyzing Transformers that recall from arbitrarily long, distributional contexts with provable generalization guarantees.

1 INTRODUCTION

Transformers (Vaswani et al., 2017) have achieved strong empirical performance across natural language (Brown et al., 2020), vision (Dosovitskiy et al., 2021), and speech/audio (Dong et al., 2018). Two properties motivate our study: (i) *content-addressable retrieval* of associated information—an associative-memory view of attention—and (ii) the ability to leverage contexts of variable, in principle unbounded, length.

In this work, we cast associative memory at the *level of probability measures*, treating context as a distribution over tokens, and develop a rigorous statistical analysis of learned softmax-attention Transformers in this measure-theoretic setting.

Associative memory provides a unifying lens on how neural systems store and retrieve from partial cues: from early self-organizing and correlation memories to Hopfield attractors (Amari, 1972; Kohonen, 1972; Nakano, 1972; Hopfield, 1982; 1984). Transformers recast associative memory or recall via content-addressable attention, formally equivalent to Hopfield-style associative updates (Vaswani et al., 2017; Ramsauer et al., 2021). Recent studies quantify memory emergence and capacity (Bietti et al., 2023; Cabannes et al., 2024; Mahdavi et al., 2024; Kim et al., 2023; Jiang et al., 2024; Nichani et al., 2025).

As Transformers are engineered to ingest massive text corpora and long contexts, researchers have formalized this “context” as a *probability measure* over tokens, yielding a measure-theoretic handle on variable-size inputs. Summarizing the text data as one measure by the law of large numbers helps them to show results that are independent of the text length. A measure-theoretic view of Transformers formalizes attention as a map on distributions, enabling analysis of stability and emergent structure (Vuckovic et al., 2020; Sander et al., 2022; Geshkovski et al., 2025; Burger et al., 2025). On the expressivity side, Transformers can interpolate between input/output measures (Geshkovski

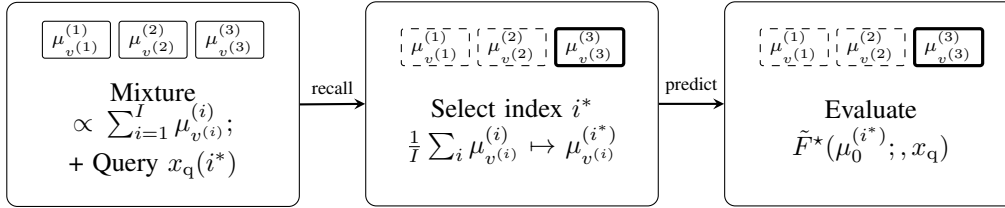


Figure 1: Associative recall at the level of measures (informal): the query $x_q(i^*)$ selects the relevant component measure $\mu_{v^{(i^*)}}$ from the mixture $\nu \propto \sum_{i=1}^I \mu_{v^{(i)}}$, followed by prediction from $(\mu_0^{(i^*)}, x_q)$. Note that each $\mu_{v^{(i)}}$ is constructed by $v^{(i)} \in \mathbb{S}^{d_1}$ and a measure $\mu_0^{(i)}$ on \mathbb{R}^{d_2} .

et al., 2024) and even uniformly approximate continuous in-context mappings where the context is itself a probability distribution (Furuya et al., 2025).

Recent work has developed statistical analyses of Transformers with infinite-dimensional inputs. Yet the link to *associative memory*—arguably a defining feature of attention—remains under-specified. Prior generalization results in distribution regression typically assume a *frozen* (non-learnable) attention kernel (Liu & Zhou, 2025), leaving unclear how *learned* attention retrieves the associated measure. Likewise, in a sequence-based in-context setting with infinite-dimensional inputs (Kim et al., 2024), the analysis was carried out under linear attention, whose limited expressiveness makes it difficult to realize the sharp, spiky weight distributions achievable by softmax attention (Han et al., 2024; Fan et al., 2025). These considerations motivate the central question:

Q. *Can a learned softmax-attention Transformer recall an infinite-dimensional (measure-valued) context and predict from it with provable generalization guarantees?*

“Associative memory” at the level of measures (informal). Consider a text corpus composed of I documents. We model each token as a vector $x = (v, z) \in \mathbb{R}^{d_1} \times \mathbb{R}^{d_2}$, where v encodes a document-level feature (e.g., topic) and z encodes token-level content. For document i , the document feature is a fixed vector $v^{(i)} \in \mathbb{S}^{d_1}$, while the content part z is sampled from some distribution $\mu_0^{(i)}$ on \mathbb{R}^{d_2} . In the limit of an infinitely long document, the empirical token distribution of document i converges to a probability measure $\mu_{v^{(i)}}$ on $\mathbb{R}^{d_1+d_2}$, namely the law of $x = (v^{(i)}, Z)$ with $Z \sim \mu_0^{(i)}$. The *context* seen by the model is then the mixture

$$\nu = \sum_{i=1}^I w_i \mu_{v^{(i)}}, \quad w_i \geq 0, \quad \sum_i w_i = 1,$$

representing a whole dataset containing many documents. Given such a mixture ν and a query $x_q \in \mathbb{R}^{d_1+d_2}$ whose first d_1 coordinates align with some document feature $v^{(i^*)}$, the desired “associative memory” behavior is:

- first, *recall* the component indexed by i^* from the mixture ν , and
- then *predict* a scalar quantity that depends only on the associated content distribution $\mu_0^{(i^*)}$ (and possibly on x_q).

We denote by $F^* : (\nu, x_q) \mapsto F^*(\nu, x_q) \in \mathbb{R}$ this ground-truth recall-and-predict map: by construction, its output depends on ν only through the single component $\mu_0^{(i^*)}$ selected by the query (Fig. 1). We study learned softmax-attention Transformers (with an integral/empirical measure view of attention) trained by empirical risk minimization (ERM) to implement this recall-and-predict pipeline. On the statistical side, we work in a very smooth regime: we endow the space of context measures with a reproducing kernel Hilbert space (RKHS) whose Mercer eigenvalues satisfy $\lambda_j \asymp \exp(-cj^\alpha)$ for some $\alpha > 0$ (as for Gaussian-type kernels (Schölkopf & Smola, 2002)). This spectral decay encodes strong smoothness of the underlying densities and induces an *effective dimension* that will govern our learning rates.

Contributions. We now outline the principal contributions of this work:

1. **Associative memory at the level of measures.** We formalize a general, mathematically rigorous framework for associative recall over measures: given a measure-valued context and a query, a *recall operator* selects the associated measure, and a *predictor* maps the recalled measure together with the query to an output. We formalize query-conditioned selection from arbitrarily long contexts: the model recalls the associated probability measure capturing the relevant content and predicts from its statistics.
2. **Generalization.** We show that a shallow (depth-2) measure-theoretic Transformer composed with an MLP can learn the recall-and-predict mapping at the level of measures (Theorem 1). In contrast to linear attentions (Kim et al., 2024) or frozen kernels (Zhou et al., 2024), softmax attention enables *sparse* and *adaptive* recall of the relevant measure. For empirical risk minimization over a bounded-parameter hypothesis class—provided the number of recall candidates is not excessively large—we establish the *sub-polynomial* population-risk bound $\exp\{-\Theta((\log n)^{\alpha/(\alpha+1)})\}$, showing that the statistical difficulty is governed by the kernel’s Mercer eigen-decay α .
3. **Minimax Optimality.** We prove a minimax lower bound with the same rate exponent $(\log n)^{\alpha/(\alpha+1)}$, establishing the sharpness of the convergence order (Theorem 2). Thus, under our spectral and mixture-growth assumptions, the proposed measure-theoretic Transformer is minimax-optimal *in the order of the exponent*, though multiplicative constants may differ.

The remainder of this paper is organized as follows. Section 2 reviews related work. Section 3 introduces the problem setting and the student model. Section 4 presents our main theoretical results. Section 5 concludes with discussion. Technical details and proofs are deferred to the appendices.

2 RELATED WORK

Associative Memory and Recall. Associative memory concepts originated in early neuroscience models (Hopfield, 1982; Amari, 1972; Kohonen, 1972; Nakano, 1972; Hopfield, 1982; 1984), followed by Graves et al. (2014); Weston et al. (2014); Ramsauer et al. (2021); Millidge et al. (2022). The Transformer architecture is closely related to associative memory by employing self-attention as a content-addressable mechanism (Vaswani et al., 2017). Recent work has increasingly focused on how associative memory emerges and scales within Transformer architectures (Bietti et al., 2023; Cabannes et al., 2024; Mahdavi et al., 2024; Kim et al., 2023; Jiang et al., 2024; Nichani et al., 2025).

Transformers for Infinite-Dimensional Inputs. A measure-theoretic perspective has enabled insightful analysis of Transformer architectures. Vuckovic et al. (2020); Sander et al. (2022) formalized self-attention as a map on probability measures. Its Lipschitzness is explored in Castin et al. (2024). Building on that framework, Geshkovski et al. (2023; 2025); Burger et al. (2025) modeled self-attention as an interacting particle system. Geshkovski et al. (2024) proved a universality result showing that Transformers can interpolate arbitrary input–output measure pairs, later strengthened by Furuya et al. (2025) to uniform approximation of continuous mappings over distributions and queries. On generalization, Liu & Zhou (2025) studied distribution regression, though restricted to a frozen attention kernel. In the context of sequential, infinite-dimensional inputs, Kim et al. (2024) studied in-context learning with linear attention, which essentially reduces to averaging behaviors; hence their analysis assumed *relaxed sparsity* and *orthonormality* of the recall candidates, reflecting the difficulty of achieving spiky one-hot recall in contrast to softmax attention (Han et al., 2024; Fan et al., 2025).

MLP Approximations of Functional Mappings. In statistical learning, Mhaskar & Hahm (1997) laid the groundwork by showing that multi-layer perceptrons (MLPs) can approximate continuous nonlinear functionals over function spaces in an optimal rate that was generalized by Stinchcombe (1999). Rossi et al. (2005) introduced a novel functional MLPs which is applicable to functional data, followed by variants (Yao et al., 2021; Song et al., 2023; Zhou et al., 2024). On the optimization front, Suzuki (2020); Nishikawa et al. (2022) established global optimization assurances for two-layer networks operating in an infinite-dimensional regime.

3 PROBLEM SETTING

Notations. For integers $N_1 \leq N_2$ and $\mathbf{v} \in \mathbb{R}^N$, we write $\mathbf{v}_{N_1:N_2} = (v_{N_1}, \dots, v_{N_2})^\top$. For a matrix A , $\|A\|_0$ denotes the number of nonzero entries and $\|A\|_\infty = \max_{i,j} |A_{i,j}|$. We write λ for the Lebesgue measure on \mathbb{R}^N , and $f_\# \mu$ for the pushforward of μ by f . For a measurable space \mathcal{X} , $\mathcal{P}(\mathcal{X})$ denotes the set of probability measures and $\mathcal{M}_+(\mathcal{X})$ the set of nonnegative measures. We use $(\Omega, \mathcal{F}, \mathbb{P})$ for a probability space, $\|f\|_{L^p}$ and $\|f\|_\infty$ for the usual L^p and essential sup norms, and \mathbb{P}_X for the law of a random variable X . Expectations are written $\mathbb{E}[\cdot]$ or $\mathbb{E}_X[\cdot]$ with the law of X .

Our Regression Problem. We now formalize the informal recall-and-predict scenario from the introduction.

Definition 1. Let $\mathcal{X}_0 \subset \mathbb{R}^{d_2}$ be a bounded token-content space and let $\mathcal{X} \subset \mathbb{R}^{d_1+d_2}$ denote the token space with the decomposition $x = (v, z)$, where $v \in \mathbb{R}^{d_1}$ encodes a document-level feature and $z \in \mathcal{X}_0$ encodes token-level content.

1. Mixture contexts and queries. Each document $i \in [I]$ is associated with a document feature $v^{(i)} \in \mathbb{S}^{d_1-1}$ and a content distribution $\mu_0^{(i)} \in \mathcal{M}_+(\mathcal{X}_0)$. Informally, $\mu_0^{(i)}$ represents the distribution of token contents (e.g., words or embeddings) appearing in document i . The corresponding token distribution on \mathcal{X} is the product measure

$$\mu_{v^{(i)}}^{(i)} := \delta_{v^{(i)}} \otimes \mu_0^{(i)} = (\text{Emb}_{v^{(i)}})_\# \mu_0^{(i)},$$

where $\text{Emb}_{v^{(i)}}(z) := (v^{(i)}, z)$ and $f_\# \mu$ denotes the pushforward of μ by f , so that $\mu_{v^{(i)}}^{(i)}$ is the joint distribution of the document feature $v^{(i)}$ and the token content in document i . A *context* is a mixture of these component measures,

$$\nu := \frac{1}{I} \sum_{i=1}^I \mu_{v^{(i)}}^{(i)} \in \mathcal{P}(\mathcal{X}),$$

which represents the token distribution of a whole dataset containing I documents. Concretely, ν is the law obtained by first sampling a document $i \in [I]$ at random and then a token $x = (v^{(i)}, z)$ from that document. Given such a mixture, a query $x_q \in \mathcal{X}$ is constructed so as to indicate a *distinguished* index $i^* \in [I]$. For concreteness, we take

$$x_q := \begin{bmatrix} v^{(i^*)} \\ 0_{d_2} \end{bmatrix} = \text{Emb}_{v^{(i^*)}}(0_{d_2}) \in \mathbb{R}^{d_1+d_2},$$

that is, the document feature $v^{(i^*)}$ padded with zeros in the last d_2 coordinates.¹

2. Ground-truth recall-and-predict map. The learning task is to predict a real-valued response

$$y = F^*(\nu, x_q) + \xi, \quad \xi \sim \mathcal{N}(0, \sigma^2),$$

from the pair (ν, x_q) . The key structural assumption is that F^* depends on the context ν only through the single component associated with the index i^* selected by the query. Equivalently, there exists a (hidden) functional \tilde{F}^* such that

$$F^*(\nu, x_q) = \tilde{F}^*(\mu_0^{(i^*)}, x_q), \quad \nu = \frac{1}{I} \sum_{i=1}^I \mu_{v^{(i)}}^{(i)},$$

so that the regression map $(\nu, x_q) \mapsto F^*(\nu, x_q)$ decomposes into the two conceptual stages

$$(\nu, x_q) \xrightarrow{\text{recall } i^*} \mu_0^{(i^*)} \xrightarrow{\text{predict}} \tilde{F}^*(\mu_0^{(i^*)}, x_q).$$

For instance, $F^*(\nu, x_q)$ could be a sentiment score of document i^* , or the probability that document i^* mentions an entity specified by the query.

In particular, F^* must first *associate* the query with the relevant component of the mixture context and then *predict* a scalar from the recalled component (e.g., in-context learning (Brown et al., 2020)). This formalizes an ‘‘associative memory’’ task at the level of probability measures.

¹More general queries, e.g. with a nonzero content part, can be treated as well; we fix the zero padding here for notational simplicity.

Statistical Estimation Problem. Motivated by the recall-and-predict regression task described above, We now formulate the associated statistical estimation problem. We observe n i.i.d. samples

$$\mathcal{S}_n := \{(\nu_t, x_{q_t}, y_t)\}_{t=1}^n$$

drawn from the joint distribution of (ν, x_q, y) . Let \mathcal{F} denote a hypothesis class of measurable functions $F : \mathcal{P}(\mathcal{X}_0) \times \mathbb{R}^{d_1+d_2} \rightarrow \mathbb{R}$. Given the training data, we define the empirical risk minimizer

$$\hat{F} := \arg \min_{F \in \mathcal{F}} \hat{\mathbb{E}}_n [(y - F(\nu, x_q))^2],$$

where $\hat{\mathbb{E}}_n$ denotes the empirical expectation over the n samples. Given a hypothesis (regressor) \hat{F} , our goal is to learn F^* so as to minimize the squared L^2 loss

$$R(F^*, \hat{F}) := \mathbb{E}_{\mathcal{S}_n} [\|F^*(\nu, x_q) - \hat{F}(\nu, x_q)\|_{L^2(\mathbb{P}_{\nu, x_q})}^2].$$

RKHS viewpoint on content measures. Before stating the assumptions, we briefly recall how the kernel K induces a function space for modeling token distributions. Given a positive definite kernel $K : \mathcal{X}_0 \times \mathcal{X}_0 \rightarrow \mathbb{R}$ on the bounded token-content domain $\mathcal{X}_0 \subset \mathbb{R}^{d_2}$, there exists a unique reproducing kernel Hilbert space (RKHS) (e.g., (Schölkopf & Smola, 2002)) \mathcal{H}_0 of functions $f : \mathcal{X}_0 \rightarrow \mathbb{R}$ such that

$$K(\cdot, x) \in \mathcal{H}_0 \quad \text{and} \quad f(x) = \langle f, K(\cdot, x) \rangle_{\mathcal{H}_0} \quad \text{for all } x \in \mathcal{X}_0.$$

Intuitively, \mathcal{H}_0 is the class of “smooth” functions compatible with K . In the Mercer basis

$$K(x, x') = \sum_{j \geq 1} \lambda_j e_j(x) e_j(x'),$$

every $f \in \mathcal{H}_0$ can be written as $f(x) = \sum_{j \geq 1} b_j e_j(x)$ with finite RKHS norm $\|f\|_{\mathcal{H}_0}^2 = \sum_{j \geq 1} b_j^2 / \lambda_j$. Large coefficients b_j in high-frequency directions (small λ_j) are penalized heavily, so \mathcal{H}_0 favours functions whose energy is concentrated on low-order eigen-components. In our setting we represent each content measure $\mu_0^{(i)}$ by its density $p_{\mu_0^{(i)}}$ on \mathcal{X}_0 , and we require these densities to lie in a fixed ball of \mathcal{H}_0 . The rapid eigenvalue decay $\lambda_j \asymp \exp(-cj^\alpha)$ corresponds to a very smooth (Gaussian Kernel-type (Schölkopf & Smola, 2002)) regime in which the effective dimension of this function class is small; this effective dimension will drive our statistical rates.

3.1 ASSUMPTIONS AND TECHNICAL SETTINGS

We now state the high-level assumptions used in our upper- and lower-bound analyses. Full technical versions are deferred to Appendix A.1. Throughout this subsection, contexts and queries are generated as described in the beginning of Section 3.

Assumption 1 (Smooth kernel and regular content measures). Let $\mathcal{X}_0 \subset \mathbb{R}^{d_2}$ be a bounded token-content domain and let $K : \mathcal{X}_0 \times \mathcal{X}_0 \rightarrow \mathbb{R}$ be a positive definite kernel with Mercer expansion (Schölkopf & Smola, 2002)

$$K(x, x') = \sum_{j \geq 1} \lambda_j e_j(x) e_j(x').$$

We assume (i) The eigenvalues decay exponentially, i.e., $\lambda_j \asymp \exp(-cj^\alpha)$ for some $c, \alpha > 0$ (a Gaussian-type smoothness regime), (ii) Each content distribution $\mu_0^{(i)}$ is a finite measure on \mathcal{X}_0 whose density lies in a fixed ball of the RKHS \mathcal{H}_0 induced by K . Informally, the token distributions are very smooth and their effective dimension is small, since most of the mass lies in low-order eigen-components. We focus on this exponentially decaying regime as a first step, to keep the analysis transparent.

Example 1 (Heat-kernel RKHS (Grigor’yan, 2006)). Consider $\mathcal{X}_0 = [0, 1]$ and the Laplace operator $\Delta = \frac{d^2}{dx^2}$ with Dirichlet boundary conditions. Its eigenfunctions and eigenvalues are $e_k(x) = \sqrt{2} \sin(k\pi x)$ and $\zeta_k = (k\pi)^2$ for $k \geq 1$. The heat kernel, the fundamental solution of the heat equation describing how heat placed at y at time 0 spreads to x by time t , is

$$K_t(x, y) = \sum_{k=1}^{\infty} e^{-\zeta_k t} e_k(x) e_k(y) = \sum_{k=1}^{\infty} e^{-\pi^2 k^2 t} e_k(x) e_k(y),$$

so the Mercer eigenvalues satisfy $\lambda_k = e^{-\pi^2 k^2 t} \asymp \exp(-ck^2)$ and Assumption 1 holds with $\alpha = 2$. More general constructions on compact manifolds are recalled in Appendix A.1.

Assumption 2 (Separated context vectors). The context vectors $v^{(i)} \in \mathbb{S}^{d_1-1}$ used to construct the mixture (1) are well separated: $\langle v^{(i)}, v^{(i')} \rangle \leq 0$, for all $1 \leq i < i' \leq I$, and we assume $I \leq d_1$. This guarantees that different documents are sufficiently distinguishable for the recall step².

Assumption 3 (Lipschitz ground-truth functional). There exists a metric d_{prod} on the space of pairs (μ_0, x) , induced by the RKHS structure above, such that the hidden functional $\tilde{F}^* : \{\mu_0^{(i)}\} \times \mathbb{R}^{d_1+d_2} \rightarrow \mathbb{R}$ is L -Lipschitz:

$$|\tilde{F}^*(\mu_0, x) - \tilde{F}^*(\mu'_0, x')| \leq L d_{\text{prod}}((\mu_0, x), (\mu'_0, x'))$$

for all admissible inputs (μ_0, x) and (μ'_0, x') . In the proofs, d_{prod} will be the sum of an RKHS-induced distance between densities and the Euclidean distance between queries; see Appendix A.1 for its precise form.

Assumption 4. Each content distribution $\mu_0^{(i)}$ is a probability measure on \mathcal{X}_0 .

Setting 1 (Probability setting for upper bound (Setting 1)). Contexts and queries are generated as in Section 3. Assumptions 1, 2, 3 and 4 hold. In other words, we consider smooth (Gaussian-type) content probability distributions in an RKHS ball, well-separated context vectors, and an L -Lipschitz ground-truth functional with respect to an RKHS-induced metric.

Structured model for lower bound. For the minimax lower bound, we work with a simplified random model for the content densities, following Lanthaler (2024): the density of μ_0 is generated by random coefficients in the Mercer expansion of K .

Assumption 5 (Informal structural model for densities). Let (λ_j, e_j) be the spectrum of K as in Assumption 1. We assume that $\frac{d\mu_0}{dX}(x) = \sum_{j \geq 1} \lambda_j^{\Theta(1)} Z_j e_j(x)$, where the coefficients Z_j are independent, bounded random variables with unit variance. The full set of structural conditions is given in Assumption 8 in Appendix A.1.

Setting 2 (Structured setting for lower bound (Setting 2)). Contexts and queries are generated as in Section 3. Assumption 1, 2 and 3 hold, and the densities of the content measures follow the structural model in Assumption 5. In this setting we derive minimax lower bounds under random Mercer coefficients.

3.2 STUDENT MODEL: MEASURE-THEORETIC TRANSFORMERS

We define our student model as a class of measure-theoretic Transformer architectures following Furuya et al. (2025).

Measure-Theoretic Attention. Given a set of tokens $X = (x_\ell)_{\ell=1}^w \in \mathbb{R}^{d_{\text{attn}} \times w}$ and a query $x \in \mathbb{R}^{d_{\text{attn}}}$ that encodes information about some of the tokens, a single unmasked attention head with parameters $\theta^{(h)} = (K^h, Q^h, V^h)$ in an ‘‘in-context’’ form (Furuya et al., 2025) computes

$$\text{SAttn}_{\theta^{(h)}}(X, x) = \sum_j \text{Softmax}(\langle Q^h x, K^h X \rangle) V^h x_j,$$

where $\text{Softmax}((z_1, \dots, z_N)) := (\exp(z_i) / \sum_j \exp(z_j))_{i=1}^N$. A standard multi-head attention with H heads is then $\text{MSAttn}_{\theta}(X) = \sum_{h=1}^H W^h \text{SAttn}_{\theta^{(h)}}(X, x)$ with $W^h \in \mathbb{R}^{d_{\text{attn}} \times d_{\text{attn}}}$. In the unmasked case, attention is permutation-equivariant in the token indices. This allows us to represent the input set by its empirical measure, in particular, in the form of a mixture measure

$$\nu_X = \frac{1}{I} \sum_{i=1}^I \left(\frac{1}{w_i} \sum_{\ell=1}^{w_i} \delta_{v^{(i)}} \otimes \delta_{u_\ell^{(i)}} \right) \in \mathcal{P}(\mathbb{R}^{d_{\text{attn}}}) \quad \text{in the limit as } w \rightarrow \infty,$$

where $v^{(i)} \in \mathbb{R}^{d_1}$ (group-shared, possibly indicated by the query), $u_\ell^{(i)} \in \mathbb{R}^{d_2}$ (token-specific) for $i \in [1 : I]$, $\ell \in [1 : w_i]$, $\sum_i w_i = w$ and rewrite attention in a measure-theoretic form, where the integral replaces the discrete sum in the limit.

²A similar separation/orthogonality structure is used in theoretical analyses of factual extraction in transformers, e.g., Ghosal et al. (2024)

Definition 2 (Measure-theoretic attention layer (Furuya et al., 2025)). Let d_{attn} be the embedding dimension for the attention. The *measure-theoretic attention layer* $\text{Attn}_\theta : \mathcal{P}(\mathbb{R}^{d_{\text{attn}}}) \times \mathbb{R}^{d_{\text{attn}}} \rightarrow \mathbb{R}^{d_{\text{attn}}}$ is defined by

$$\text{Attn}_\theta(\nu, x) = Ax + \sum_{h=1}^H W^h \int \text{Softmax}(\langle Q^h x, K^h y \rangle) V^h y \, d\nu(y),$$

where $\text{Softmax}(\langle Qx, Ky \rangle) := \exp(\langle Qx, Ky \rangle) / \int \exp(\langle Qx, Kz \rangle) \, d\nu(z)$. Here $A \in \mathbb{R}^{d_{\text{attn}} \times d_{\text{attn}}}$ applies a learned linear transformation to the skip connection.

When ν is an empirical mixture measure ν_X , this reduces to the standard discrete attention layer (we do not pursue the discrete case in this work). We expect that the query x indicates tokens in the i^* -th component in ν_X based on the group-shared vector $v^{(i^*)}$ and the model can recall them. We constrain these layers via the following bounded-parameter hypothesis class.

Definition 3 (Attention hypothesis class). For constants $B_a, B'_a, S_a, S'_a, F > 0$, define the class of H -head measure-theoretic attention layers

$$\begin{aligned} \mathcal{A}(d_{\text{attn}}, H, B_a, B'_a, S_a, S'_a) := & \{ \text{Attn}_\theta \mid \max_h \{ \|W^h\|_\infty, \|Q^h\|_\infty, \|K^h\|_\infty, \|V^h\|_\infty \} \leq B_a, \\ & \max_h \{ \|W^h\|_0, \|Q^h\|_0, \|K^h\|_0, \|V^h\|_0 \} \leq S_a, \|A\|_\infty \leq B'_a, \|A\|_0 \leq S'_a, \|\text{Attn}_\theta\|_\infty \leq F \}, \end{aligned}$$

where $\|M\|_\infty := \max_{i,j} |M_{ij}|$, $\|M\|_0$ is the number of non-zero entries, for a matrix M . We assume W^h, Q^h, K^h, V^h are square matrices for simplicity. We write $\mathcal{A}(d_{\text{attn}}, H, B_a, S_a)$ when $B'_a = B_a$ and $S'_a = S_a$.

MLP Layer. In addition to attention, a Transformer block also includes a feedforward component. Since this part does not depend on the underlying measure μ , we model it simply as a standard multilayer perceptron (MLP), defined below.

Definition 4 (MLP hypothesis class). Let $\ell \in \mathbb{N}$ be the depth and $\mathbf{p} = (p_0, p_1, \dots, p_{L+1}) \in \mathbb{N}^{L+2}$ be the layer widths. A neural network with architecture (ℓ, \mathbf{p}) is any function of the form

$$f : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{\ell+1}}, \quad \mathbf{x} \mapsto f(\mathbf{x}) = W_\ell \sigma_{\mathbf{v}_\ell} W_{\ell-1} \sigma_{\mathbf{v}_{\ell-1}} \cdots W_1 \sigma_{\mathbf{v}_1} W_0 \mathbf{x},$$

where $W_i \in \mathbb{R}^{p_{i+1} \times p_i}$ is the weight matrix of layer i and $\mathbf{v}_i \in \mathbb{R}^{p_i}$ is a shift (bias) vector applied through the activation $\sigma_{\mathbf{v}_i}(\mathbf{z}) := \sigma(\mathbf{z} - \mathbf{v}_i)$ with ReLU activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$. The hypothesis class of MLPs with architecture (ℓ, \mathbf{p}) is denoted

$$\begin{aligned} \mathcal{F}(\ell, \mathbf{p}, s, F) := & \{ f \text{ of the form above} \mid \\ & \max_j \{ \|W_j\|_\infty, |v_j|_\infty \} \leq 1, \sum_j \|W_j\|_0 + |v_j|_0 \leq s, \|f\|_\infty \leq F \}. \end{aligned}$$

Transformer Hypothesis Class via Composition. To formally describe a Transformer within our measure-theoretic framework, we introduce the notion of *composition* for measure-theoretic mappings, following Furuya et al. (2025). This will allow us to view a multi-layer Transformer as a successive composition of attention and feedforward layers, in exact analogy with the standard architecture.

The key idea is that such a mapping Γ acts simultaneously on a token x and on its generating distribution μ . Given $z \sim \mu$, applying Γ produces a transformed token $\Gamma(\mu, z)$ and induces a new distribution on transformed tokens $\Gamma(\mu, z)$, namely $(\Gamma(\mu, \cdot))_\# \mu$. Thus, composing two mappings means successively applying these joint transformations at both the sample and distribution levels.

Definition 5 (Composition of measure-theoretic mappings). Let $\Gamma_1 : \mathcal{P}(\mathbb{R}^{d_1^{(1)}}) \times \mathbb{R}^{d_1^{(1)}} \rightarrow \mathbb{R}^{d_1^{(2)}}$ and $\Gamma_2 : \mathcal{P}(\mathbb{R}^{d_2^{(1)}}) \times \mathbb{R}^{d_2^{(1)}} \rightarrow \mathbb{R}^{d_2^{(2)}}$ with $d_1^{(2)} = d_2^{(1)}$. Their composition is defined as

$$(\Gamma_2 \diamond \Gamma_1)(\nu, x) := \Gamma_2(\mu_1, \Gamma_1(\nu, x)), \quad \text{where} \quad \mu_1 := (\Gamma_1(\nu, \cdot))_\# \nu.$$

Remark 1. If we interpret ν as a limit of empirical measure $\nu_X = \lim_{w \rightarrow \infty} \frac{1}{w} \sum_{\ell=1}^w \delta_{x_\ell}$, then the composition $\Gamma_2 \diamond \Gamma_1$ acts by updating both the individual tokens and their empirical distribution. In this case, the construction is consistent with the standard layerwise composition in a Transformer: each layer maps the sequence of tokens (x_1, \dots, x_w) to a new sequence, while the corresponding empirical distribution is updated accordingly.

Building on the notion of measure-theoretic composition introduced above, we can now formally describe a Transformer as the successive composition of attention and feedforward layers. The following definition specifies the corresponding hypothesis class.

Definition 6 (Transformer hypothesis class). For parameters $(d_j, H_j, B_{a,j}, B'_{a,j}, S_{a,j}, S'_{a,j}, \ell_j, \mathbf{p}_j, s_j)_{j=1}^L$, define TF as the set of mappings of the form

$$\text{Attn}_{\theta_L} \diamond \text{MLP}_{\xi_L} \diamond \cdots \diamond \text{Attn}_{\theta_1} \diamond \text{MLP}_{\xi_1},$$

where $\text{Attn}_{\theta_j} \in \mathcal{A}(d_j, H_j, B_{a,j}, B'_{a,j}, S_{a,j}, S'_{a,j})$ and $\text{MLP}_{\xi_j} \in \mathcal{F}(\ell_j, \mathbf{p}_j, s_j)$. MLP layers are independent of μ (i.e. $\text{MLP}_{\xi_1}(\mu, x) := \text{MLP}_{\xi_1}(x)$), and all intermediate outputs are assumed uniformly bounded.

4 MAIN RESULTS

4.1 ESTIMATION ERROR OF MEASURE-THEORETIC TRANSFORMERS

We begin with the generalization performance of measure-theoretic Transformers in the recall-and-predict task. Throughout this subsection we work under the Probability Setting (Setting 1), where each content measure $\mu_0^{(i)}$ has a smooth RKHS density on \mathcal{X}_0 and the number of mixture components is not too large.

Theorem 1 (Sub-polynomial convergence; informal version of Theorem 3). *Let $F^*(\nu, x_q) = \tilde{F}^*(\mu_0^{(i^*)}, x_q)$ be a Lipschitz recall-and-predict map as in Section 3, and assume that the eigenvalues of the underlying kernel satisfy $\lambda_j \asymp \exp(-cj^\alpha)$ for some $\alpha > 0$. Suppose that either (i) the number of mixture components satisfies $I \leq d_1 \lesssim (\log n)^{1/(\alpha+1)}$, or (ii) \tilde{F}^* does not depend on x_q and $I \leq d_1 = n^{o(1)}$. Then, for a suitable choice of architecture parameters (defining a depth-2 measure-theoretic Transformer class TF_n), any empirical risk minimizer \hat{F}_n over TF_n satisfies*

$$R(F^*, \hat{F}_n) \lesssim \exp\{-\Omega((\log n)^{\alpha/(\alpha+1)})\}$$

under Setting 1.

Statistically unifying associative recall and infinite-token regimes. Prior work studied (i) *associative recall in Transformers* (Ramsauer et al., 2021) and (ii) *infinite-token / infinite-dimensional* inputs modeled as measures (Vuckovic et al., 2020). Theorem 1 integrates these threads by giving a *statistical* theory of measure-level associative recall: a measure-theoretic Transformer with learned softmax attention can *recall-and-predict at the level of measures*—sparsely isolating the query-relevant component of ν and basing the prediction on the recalled measure, in contrast to the universality or approximation results (Geshkovski et al., 2024; Furuya et al., 2025).

Informal interpretation: effective dimension. The spectral decay $\lambda_j \asymp \exp(-cj^\alpha)$ means that only the first few Mercer modes carry substantial signal. After the recall step, our Transformer effectively aggregates the first D Mercer coefficients

$$b_j \approx \int e_j d\mu_0^{(i^*)}, \quad j = 1, \dots, D,$$

so an infinite-dimensional measure is compressed into the D -dimensional vector $b = (b_1, \dots, b_D)$. Learning a Lipschitz function of b from n samples behaves like a D -dimensional problem (c.f., Schmidt-Hieber (2020)), with estimation error roughly

$$\text{Error of } D\text{-dim. problem} \approx n^{-\Theta(1/D)} \simeq \exp(-\Theta((\log n)/D)),$$

while truncating the Mercer expansion after D modes incurs a bias of order $\exp(-cD^\alpha)$. Balancing these terms yields an effective dimension $D_{\text{eff}}(n) \asymp (\log n)^{1/(\alpha+1)}$ and the sub-polynomial rate

$$R(F^*, \hat{F}_n) \approx \exp(-\Theta((\log n)^{\alpha/(\alpha+1)}))$$

stated in Theorem 1. In this sense, the estimator behaves as if it were fitting only $D_{\text{eff}}(n)$ degrees of freedom, despite each component being an infinite-dimensional measure. As a minimal sanity check, Appendix D presents a synthetic experiment whose convergence rate is consistent.

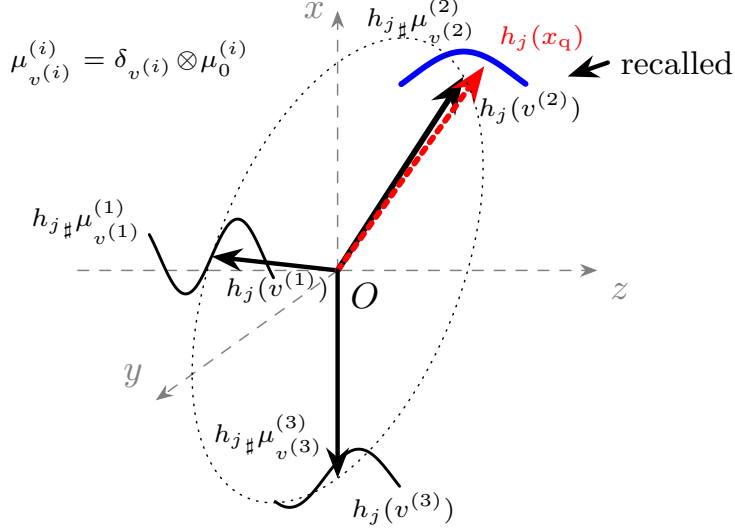


Figure 2: A geometric illustration of how query x_q and components $\mu_{v^{(i)}}^{(i)}$ are mapped by the (simplified) first layer $h_j((v, z)) = (v, e_j(z)) \in \mathbb{R}^{d_1+d_2}$, where e_j is the j th Mercer eigenfunction. The product of the first d_1 indices of $h_j(x_q)$ and $h_j(y) \sim h_{j\#}\mu_{v^{(i)}}^{(i)}$ tells whether $y = (v, z)$ is sampled from $\mu_{v^{(i^*)}}^{(i^*)}$ or not.

Mechanism: softmax attention as measure-valued associative memory. Given a mixture context $\nu = I^{-1} \sum_{i=1}^I \mu_{v^{(i)}}^{(i)}$ and a query $x_q = (v^{(i^*)}, 0)$, our depth-2 Transformer works as follows. An initial MLP embeds each token $y = (v^{(i)}, z)$ into a feature vector $h(y)$ that contains the first D Mercer features $(e_j(z))_{j=1}^D$. Softmax attention then computes scores $\langle Qh(x_q), Kh(y) \rangle$ and is parameterized so that these scores are large mainly when the document tag $v^{(i)}$ matches $v^{(i^*)}$ and small otherwise (Fig. 2).

Softmax attention then computes scores $\langle Qh(x_q), Kh(y) \rangle$ so that, after normalization, the weights concentrate on samples from $\mu_{v^{(i^*)}}^{(i^*)}$. Writing $w_{x_q}(y)$ for the resulting attention weight on token y , the value path computes, for each j ,

$$\hat{b}_j \approx \int e_j(z) w_{x_q}(y) d\nu(y) \approx \int e_j(z) d\mu_0^{(i^*)}(z), \quad j = 1, \dots, D,$$

yielding a D -dimensional descriptor $b = (\hat{b}_1, \dots, \hat{b}_D)$ of the recalled component measure. A final MLP maps (b, x_q) to the scalar prediction $\tilde{F}^*(\mu_0^{(i^*)}, x_q)$. In this way, softmax attention *filters* the mixture down to the relevant component and *integrates* its Mercer features, so that predicting from an infinite-dimensional measure reduces to learning a function of D summary statistics. This near one-hot, query-dependent filtering is beyond the limitations of frozen kernels (Zhou et al., 2024) and linear attentions (Kim et al., 2024).

4.2 MINIMAX OPTIMALITY (THE LOWER BOUND).

Next, we demonstrate the lower bound in Structured Setting (Setting 2). We show that the exponent of the obtained upper bound is essentially tight, albeit in Setting 2 whose technical assumptions are different from Setting 1.

Theorem 2 (Minimax Lower Bound). *In Setting 2, any $L^2(\mathbb{P}_{\nu, x_q}(\mathcal{H}_0))$ -estimator \hat{F} satisfies*

$$\sup_{\hat{F}^* \in \mathcal{F}^*} R(\hat{F}, F^*) \gtrsim \exp\left(-O((\ln n)^{\frac{\alpha}{\alpha+1}})\right)$$

where $F^*(\nu, x_q) := \tilde{F}^*(\mu_0^{(i^*)}, x_q)$, α is the decay rate of the eigenvalues of the kernel of \mathcal{H}_0 .

Optimality of Transformers. Taken together with Theorem 1, our information-theoretic lower bound in Theorem 2 shows that the statistical rate of empirical risk minimization over transformers achieves the minimax rate $(\ln n)^{\frac{\alpha}{\alpha+1}}$ up to multiplicative constants in the exponent. Equivalently, no method can improve the n -dependence beyond the stated exponent (up to universal constants), so the learned-softmax Transformer attains the best-possible sample complexity for this problem class. The optimality continues to hold for a fixed or slowly growing number of mixture components I , confirming that the learned softmax attention provides the right inductive bias for measure-level recall. This minimax optimality of softmax Transformers is consistent with statistical results for simple infinite-dimensional regression (Takakura & Suzuki, 2023) and with finite-dimensional in-context learning scenarios that require retrieval (Nishikawa et al., 2025).

Technical Contribution—Minimax Lower Bound. We first reduce associative recall to infinite-dimensional Lipschitz regression by observing that estimating from the mixed input ν is no easier than from the pure measure $\mu_0^{(i^*)}$. A truncation of Mercer coefficients plus anisotropic rescaling—modifying prior rescaling arguments (Lanthaler, 2024) to more general geometry with exponential decay—makes the induced geometry essentially isotropic, letting us embed the classical d -dimensional Lipschitz class and import standard packing bounds. Combining these bounds with the classical result (Yang & Barron, 1999) yields a rate matching our upper bound.

5 CONCLUSION AND DISCUSSION

We introduced the concept of measure-theoretic associative memory (recall) and established that learned softmax attention can realize sharp recall even in infinite-dimensional, measure-valued settings—something beyond the reach of frozen kernels and difficult for linear attention. Our analysis further shows that the statistical efficiency of Transformers extends beyond finite-dimensional contexts, offering a principled explanation of their recall ability. While the present results focus on exponentially decaying spectra under smooth eigenfunctions, they open the door to broader regimes: extending the rates to polynomial decay and incorporating eigenfunction smoothness into the analysis represent natural next steps toward a more complete theory.

LLM USAGE STATEMENT

Large language models are used for two purposes: to proofread and polish English writing, to help us find related works. We did not use any LLM assistant for designing the problem settings and constructing the proofs.

5.0.1 ACKNOWLEDGMENTS

RK and TS were partially supported by JSPS KAKENHI (24K02905) and JST CREST (JP-MJCR2115). This research is supported by the National Research Foundation, Singapore and the Ministry of Digital Development and Information under the AI Visiting Professorship Programme (award number AIVP-2024-004). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore and the Ministry of Digital Development and Information. RK was supported by the FY 2024 Self-directed Research Activity Grant of the University of Tokyo’s International Graduate Program “Innovation for Intelligent World” (IIW).

REFERENCES

- Shun-ichi Amari. Learning patterns and pattern sequences by self-organizing nets of threshold elements. *IEEE Transactions on Computers*, C-21(11):1197–1206, 1972.
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Hervé Jégou, and Léon Bottou. Birth of a transformer: A memory viewpoint. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- Emmanuel Boissard. Simple bounds for the convergence of empirical and occupation measures in 1-wasserstein distance. *Electronic Journal of Probability*, 16:2296–2333, 2011.

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Martin Burger, Samira Kabri, Yury Korolev, Tim Roith, and Lukas Weigand. Analysis of mean-field models arising from self-attention dynamics in transformer architectures with layer normalization. *Philosophical Transactions A*, 383(2298):20240233, 2025.
- Vivien Cabannes, Elvis Dohmatob, and Alberto Bietti. Scaling laws for associative memories. In *International Conference on Learning Representations (ICLR)*, 2024.
- Valérie Castin, Pierre Ablin, and Gabriel Peyré. How smooth is attention? In *International Conference on Machine Learning*, pp. 5817–5840. PMLR, 2024.
- Ronald R. Coifman and Stéphane Lafon. Diffusion maps. *Applied and Computational Harmonic Analysis*, 21(1):5–30, 2006.
- Linhao Dong, Shuang Xu, and Bo Xu. Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 5884–5888. IEEE, 2018.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Weinan E and Qingcan Wang. Exponential convergence of the deep neural network approximation for analytic functions. *Science China Mathematics*, 61(10):1733–1740, 2018. ISSN 1869-1862.
- Qihang Fan, Huaibo Huang, Yuang Ai, and Ran He. Rectifying magnitude neglect in linear attention, 2025.
- Takashi Furuya, Maarten V. de Hoop, and Gabriel Peyré. Transformers are universal in-context learners. In *The Thirteenth International Conference on Learning Representations*, 2025.
- Borjan Geshkovski, Cyril Letrouit, Yury Polyanskiy, and Philippe Rigollet. The emergence of clusters in self-attention dynamics. *Advances in Neural Information Processing Systems*, 36:57026–57037, 2023.
- Borjan Geshkovski, Philippe Rigollet, and Domènec Ruiz-Balet. Measure-to-measure interpolation using transformers. *arXiv preprint arXiv:2411.04551*, 2024.
- Borjan Geshkovski, Cyril Letrouit, Yury Polyanskiy, and Philippe Rigollet. A mathematical perspective on transformers. *Bulletin of the American Mathematical Society*, 62(3):427–479, 2025.
- Gaurav Ghosal, Tatsunori Hashimoto, and Aditi Raghunathan. Understanding finetuning for factual knowledge extraction, 2024.
- Alex Graves, Greg Wayne, and Ivo Danihelka. Neural Turing machines. *arXiv preprint arXiv:1410.5401*, 2014.
- Alexander Grigor’yan. Heat kernels on weighted manifolds and applications. *Cont. Math*, 398(2006):93–191, 2006.
- Dongchen Han, Yifan Pu, Zhuofan Xia, Yizeng Han, Xuran Pan, Xiu Li, Jiwen Lu, Shiji Song, and Gao Huang. Bridging the divide: Reconsidering softmax and linear attention. *Advances in Neural Information Processing Systems*, 37:79221–79245, 2024.
- 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.
- John J. Hopfield. Neurons with graded response have collective computational properties like those of two-state neurons. *Proceedings of the National Academy of Sciences*, 81(10):3088–3092, 1984.

- Y. Jiang, G. Rajendran, P. Ravikumar, and B. Aragam. Do llms dream of elephants (when told not to)? latent concept association and associative memory in transformers. *arXiv preprint arXiv:2406.18400*, 2024.
- J. Kim, M. Kim, and B. Mozafari. Provable memorization capacity of transformers. In *International Conference on Learning Representations (ICLR)*, 2023.
- Juno Kim, Tai Nakamaki, and Taiji Suzuki. Transformers are minimax optimal nonparametric in-context learners. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 106667–106713. Curran Associates, Inc., 2024.
- Teuvo Kohonen. Correlation matrix memories. *IEEE Transactions on Computers*, C-21(4):353–362, 1972.
- Samuel Lanthaler. Operator learning of lipschitz operators: An information-theoretic perspective, 2024.
- Peilin Liu and Ding-Xuan Zhou. Generalization analysis of transformers in distribution regression. *Neural Computation*, 37(2):260–293, 01 2025. ISSN 0899-7667.
- S. Mahdavi, R. Liao, and C. Thrampoulidis. Memorization capacity of multihead attention in transformers. In *International Conference on Learning Representations (ICLR)*, 2024.
- H. N. Mhaskar and Nahmwoo Hahm. Neural networks for functional approximation and system identification. *Neural Computation*, 9(1):143–159, 01 1997. ISSN 0899-7667.
- Beren Millidge, Tommaso Salvatori, Yuhang Song, Thomas Lukasiewicz, and Rafal Bogacz. Universal hopfield networks: A general framework for single-shot associative memory models. In *International Conference on Machine Learning*, pp. 15561–15583. PMLR, 2022.
- Boaz Nadler, Stephane Lafon, Ioannis Kevrekidis, and Ronald Coifman. Diffusion maps, spectral clustering and eigenfunctions of fokker-planck operators. In Y. Weiss, B. Schölkopf, and J. Platt (eds.), *Advances in Neural Information Processing Systems*, volume 18. MIT Press, 2005.
- Kenji Nakano. Associatron—a model of associative memory. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-2(3):380–388, 1972.
- Eshaan Nichani, Jason D. Lee, and Alberto Bietti. Understanding factual recall in transformers via associative memories. In *International Conference on Learning Representations (ICLR)*, 2025. Spotlight presentation, preprint arXiv:2412.06538.
- Naoki Nishikawa, Taiji Suzuki, Atsushi Nitanda, and Denny Wu. Two-layer neural network on infinite dimensional data: global optimization guarantee in the mean-field regime. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 32612–32623. Curran Associates, Inc., 2022.
- Naoki Nishikawa, Yujin Song, Kazusato Oko, Denny Wu, and Taiji Suzuki. Nonlinear transformers can perform inference-time feature learning. In *Forty-second International Conference on Machine Learning*, 2025.
- Hubert Ramsauer, Bernhard Schöfl, Anna Hopkins, Michael Widrich, Thomas Adler, Lukas Gruber, Markus Holzleitner, Zoltan Pavlović, Geir Kjetil Sandve, Vidar Greiff, Johannes Brandstetter, and Sepp Hochreiter. Hopfield networks is all you need. In *International Conference on Learning Representations (ICLR)*, 2021.
- Fabrice Rossi, Nicolas Delannay, Brieuc Conan-Guez, and Michel Verleysen. Representation of functional data in neural networks. *Neurocomputing*, 64:183–210, 2005. ISSN 0925-2312. Trends in Neurocomputing: 12th European Symposium on Artificial Neural Networks 2004.
- Michael E. Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. Sinkformers: Transformers with doubly stochastic attention. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 3515–3530. PMLR, 28–30 Mar 2022.

- Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation function. *The Annals of Statistics*, 48(4):1875–1897, 2020.
- Bernhard Schölkopf and Alexander J Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2002.
- Linhao Song, Jun Fan, Di-Rong Chen, and Ding-Xuan Zhou. Correction: Approximation of non-linear functionals using deep relu networks. *Journal of Fourier Analysis and Applications*, 29(5): 57, 2023. ISSN 1531-5851.
- M.B. Stinchcombe. Neural network approximation of continuous functionals and continuous functions on compactifications. *Neural Networks*, 12(3):467–477, 1999. ISSN 0893-6080.
- Taiji Suzuki. Generalization bound of globally optimal non-convex neural network training: Transportation map estimation by infinite dimensional langevin dynamics. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 19224–19237. Curran Associates, Inc., 2020.
- Shokichi Takakura and Taiji Suzuki. Approximation and estimation ability of transformers for sequence-to-sequence functions with infinite dimensional input. In *International Conference on Machine Learning*, pp. 33416–33447. PMLR, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- James Vuckovic, Aristide Baratin, and Remi Tachet des Combes. A mathematical theory of attention. *arXiv preprint arXiv:2007.02876*, 2020.
- Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2014.
- Weichun Xia and Lei Shi. Spectral algorithms on manifolds through diffusion. *arXiv preprint*, 2024.
- Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of convergence. *The Annals of Statistics*, 27(5):1564–1599, 1999. ISSN 00905364, 21688966.
- Junwen Yao, Jonas Mueller, and Jane-Ling Wang. Deep learning for functional data analysis with adaptive basis layers. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 11898–11908. PMLR, 18–24 Jul 2021.
- Tian-Yi Zhou, Namjoon Suh, Guang Cheng, and Xiaoming Huo. Approximation of rkhs functionals by neural networks. *CoRR*, abs/2403.12187, 2024.

A PRELIMINARIES AND NOTATIONAL REMARKS

We begin by fixing notation and conventions used throughout the paper. This includes basic vector and matrix operations, measure-theoretic notation, and standard identifications between measures and their densities.

For integers $N_1 \leq N_2$ and a vector $\mathbf{v} \in \mathbb{R}^N$, we define $\mathbf{v}_{N_1:N_2} := (v_{N_1}, \dots, v_{N_2})^\top$. For a matrix A , we define $\|A\|_0$ as the number of nonzero entries and $\|A\|_\infty := \max_{i,j} |A_{i,j}|$. We denote λ as the Lebesgue measure on \mathbb{R}^N , for a integer N . $f_\# \mu(\cdot) := \mu(f^{-1}(\cdot))$ denotes a pushforward of a measure μ by a mapping f . For a measurable space \mathcal{X} , we write $\mathcal{P}(\mathcal{X}) := \{\mu \mid \mu \text{ is a probability measure on } \mathcal{X}\}$, $\mathcal{M}_+(\mathcal{X}) := \{\mu \mid \mu \text{ is a nonnegative measure on } \mathcal{X}\}$.

Unless otherwise specified, we write $(\Omega, \mathcal{F}, \mathbb{P})$ for an underlying probability space. For a measurable function $f : \Omega \rightarrow \mathbb{R}$ and $1 \leq p < \infty$, the $L^p(\mathbb{P})$ norm is defined as $\|f\|_{L^p(\mathbb{P})} := \left(\int_\Omega |f|^p d\mathbb{P} \right)^{1/p}$, while the $L^\infty(\mathbb{P})$ norm is given by $\|f\|_{L^\infty(\mathbb{P})} := \text{ess sup}_{\omega \in \Omega} |f(\omega)|$. When the underlying measure is clear from context, we simply write $\|f\|_{L^p}$ and $\|f\|_\infty$.

For a random variable $X : \Omega \rightarrow \mathcal{X}$, we denote its distribution (the pushforward of \mathbb{P} under X) by $\mathbb{P}_X(\cdot) := \mathbb{P}(X \in \cdot)$. Expectations with respect to \mathbb{P} are denoted by $\mathbb{E}[\cdot]$, and if X is a random variable with law \mathbb{P}_X , we also write $\mathbb{E}_X[f(X)]$ for $\int f(x) d\mathbb{P}_X(x)$.

Remark 2 (Identification of measures and densities). Let λ be a reference measure on X (e.g., the Lebesgue measure). Given $f \in \mathcal{H}$ and a constant $c \in \mathbb{R}$, we define a probability measure μ by

$$\frac{d\mu}{d\lambda}(x) := f(x) + c,$$

where c is chosen so that $f + c \geq 0$ λ -a.e. and $\mu(X) = \int_X (f(x) + c) dx = 1$. In this case, we write $\mu \in \mathcal{H}$ by identifying μ with its density $f + c$. When sampling $\mu_0 \in B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$, please note that we choose c with no randomness and let $c = 0$ for simplicity, throughout this paper.

Definition 7 (Mercer expansion). Let \mathcal{H}_0 be a reproducing kernel Hilbert space (RKHS) on a domain \mathcal{X} with reproducing kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. By Mercer's theorem, K admits the decomposition

$$K(x, x') = \sum_{j=1}^{\infty} \lambda_j e_j(x) e_j(x'),$$

where $(\lambda_j)_{j \geq 1}$ are the (non-negative, non-increasing) Mercer eigenvalues and $(e_j)_{j \geq 1}$ are the corresponding L^2 -orthonormal eigenfunctions. For any $\mu \in \mathcal{H}_0$, its representation in the eigenbasis is

$$p_\mu = \sum_{j=1}^{\infty} b_j e_j, \quad \text{with coefficients } b_j = \langle p_\mu, e_j \rangle_{L^2}.$$

Definition 8 (Generalized RKHS norm). Let \mathcal{H}_0 be an RKHS with orthonormal basis $\{e_j\}_{j \geq 1}$ in L^2 and associated eigenvalues $\{\lambda_j\}_{j \geq 1}$. For a measure μ whose associated function in \mathcal{H}_0 admits the expansion

$$p_\mu = \sum_{j \geq 1} b_j e_j, \quad \text{where } \frac{d\mu}{d\lambda}(x) = p_\mu(x) + \text{constant}.$$

we define, for $a \in \mathbb{R}$, the *generalized norm*

$$\|\mu\|_{\mathcal{H}_0^a}^2 := \|p_\mu\|_{\mathcal{H}_0^a}^2 := \sum_{j \geq 1} \lambda_j^{-a} b_j^2.$$

Special cases include the case $a = 0$: $\|\cdot\|_{\mathcal{H}_0^0}$ coincides with the L^2 norm; $a = 1$: $\|\cdot\|_{\mathcal{H}_0^1}$ is the standard RKHS norm; $a = -1$: $\|\cdot\|_{\mathcal{H}_0^{-1}}$ is the MMD norm.

Definition 9 (Metric balls). Let (X, d) be a metric space. For $x \in X$ and $\epsilon > 0$, the (closed) metric ball of radius ϵ centered at x is

$$B(x, d, \epsilon) := \{y \in X \mid d(x, y) \leq \epsilon\}.$$

When $\epsilon = 1$, we simply write $B(x, d)$ and refer to it as the *unit ball*.

Definition 10 (Lipschitz functions). Let (X, d) be a metric space and $A \subset X$ be a fixed domain. A function $f : A \rightarrow \mathbb{R}$ is said to be *L-Lipschitz* on A with respect to d if

$$|f(x) - f(x')| \leq L d(x, x') \quad \forall x, x' \in A.$$

The set of all such functions is denoted by $\text{Lip}_L(A, d)$.

PRELIMINARIES ON METRIC ENTROPY

For the subsequent proofs, we will make repeated use of standard notions from metric entropy. In particular, coverings, packings, and their associated numbers provide a convenient way to quantify the complexity of hypothesis classes. We therefore collect the relevant definitions and basic lemmas here.

Definition 11 (ϵ -covering). Let (X, d) be a metric space, $A \subset X$, and $\epsilon > 0$. A finite set $\{x_1, \dots, x_N\} \subset X$ is called an ϵ -*covering* of A (with respect to d) if

$$A \subset \bigcup_{i=1}^N B(x_i, d, \epsilon),$$

where $B(x_i, d, \epsilon)$ denotes the closed metric ball of radius ϵ centered at x_i .

Definition 12 (ϵ -packing). Let (X, d) be a metric space, $A \subset X$, and $\epsilon > 0$. A finite set $\{x_1, \dots, x_N\} \subset A$ is called an ϵ -packing of A (with respect to d) if

$$d(x_i, x_j) > \epsilon \quad \text{for all } i \neq j.$$

Equivalently, the metric balls $B(x_i, d, \epsilon/2)$ are pairwise disjoint.

Definition 13 (Covering number). The *covering number* of A at scale ϵ with respect to d is

$$\mathcal{N}(A, \epsilon)_d := \min \{ N \mid \exists \epsilon\text{-covering of } A \text{ of size } N \}.$$

Definition 14 (Packing number). The *packing number* of A at scale ϵ with respect to d is

$$\mathcal{M}(A, \epsilon)_d := \max \{ M \mid \exists \epsilon\text{-packing of } A \text{ of size } M \}.$$

Lemma 1 (Covering–packing equivalence). *For any metric space (X, d) , any $A \subset X$, and $\epsilon > 0$, one has*

$$\mathcal{M}(A, 2\epsilon)_d \leq \mathcal{N}(A, \epsilon)_d \leq \mathcal{M}(A, \epsilon)_d.$$

In particular, the covering number and the packing number are equivalent up to constant factors in the scale parameter.

Lemma 2 (Monotonicity under metric domination). *Let d and d' be two metrics on X , and let $c > 0$ such that $d(x, y) \leq c d'(x, y)$ for all $x, y \in X$. Then, for any $A \subset X$ and $\epsilon > 0$,*

$$\mathcal{N}(A, \epsilon)_d \leq \mathcal{N}(A, \epsilon)_{cd'} = \mathcal{N}(A, c^{-1}\epsilon)_{d'}.$$

In words: if d is dominated by cd' , then ϵ -coverings with respect to cd' are also ϵ -coverings with respect to d , hence covering under d is no harder.

A.1 TECHNICAL VERSION OF SECTION 3.1

Now we will introduce two technical settings for the data generation. We begin with the assumptions required for establishing the estimation upper bound, and then turn to alternative structural assumptions that are used for deriving the lower bound.

A.1.1 PROBABILITY SETTING FOR UPPER BOUND.

Assumption 6 (Common Assumptions: The RKHS Structure, the Regularity, and the Lipschitzness: Technical Version of Assumption 1,2,3). Fix an integer $i^* \in \{1, \dots, I\}$. A query vector x_q , an input measure $\nu = \sum_i \delta_{v^{(i)}} \otimes \mu_0^{(i)}$, and an output $y = F^*(\nu, x_q) + \xi$ are generated as follows:

- *RKHS setting*: We assume the density of $\mu_0 \sim \mathbb{P}_{\mu_0}$ is in a space \mathcal{H}_0 ignoring a constant. \mathcal{H}_0 is an RKHS on a bounded domain $\mathcal{X}_0 \subset [-M, M]^{d_2}$ with Mercer decomposition (e.g., (Schölkopf & Smola, 2002))

$$K(x, x') = \sum_{j=1}^{\infty} \lambda_j e_j(x) e_j(x'), \quad \lambda_j \simeq \exp(-cj^\alpha),$$

where $c, \alpha > 0$ and eigenfunctions $(e_j)_{j \geq 1}$ are L^2 -orthonormal. μ_0 is nonnegative and $\frac{d\mu_0}{d\lambda} - c'$ lies in the metric ball $B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$ with $\gamma_b > 0$, where λ is the Lebesgue measure and c' is a global constant without any randomness. Let $c' = 0$ for simplicity, throughout this paper. The density is infinite-dimensional, but the informative content is sharply concentrated in low-order components, with high-frequency contributions decaying exponentially. For $a \in \mathbb{R}$, the \mathcal{H}_0^a -norm on \mathcal{H}_0 is defined as $\|\mu_0\|_{\mathcal{H}_0^a}^2 := \|\frac{d\mu_0}{d\lambda}\|_{\mathcal{H}_0^a}^2 := \sum_{j \geq 1} \lambda_j^{-a} b_j^2$ where $\frac{d\mu_0}{d\lambda} = \sum_j b_j e_j$. Similar RKHS structures can be found in Suzuki (2020); Nishikawa et al. (2022); Zhou et al. (2024); Liu & Zhou (2025).

- *Smoothness of eigenfunctions in \mathcal{H}_0* : (a) $(e_j)_{j \geq 1}$ are uniformly bounded, and analytic on $[-M, M]^{d_2}$; (b) $\mathcal{X}_0 \subset [-M + \delta, M - \delta]^{d_2}$ for some $\delta > 0$; (c) Each e_j admits an absolutely convergent power series $e_j(x) = \sum_{k \in \mathbb{N}^{d_2}} a_k x^k$ on $[-M, M]^{d_2}$. The three conditions are required to focus on the sample complexity regarding the decay rate α .

- *Distinguishability of sampled context vectors:* The vectors $(\mathbb{S}^{d_1-1})^{\otimes I} \ni (v^{(i)})_{i=1}^I \sim \mathbb{P}_v$ satisfy $\langle v^{(i)}, v^{(i')} \rangle \leq 0$, $1 \leq i < i' \leq I$. To satisfy this, we also require $I \leq d_1$. This ensures that contexts are sufficiently distinguishable for recall.
- *Lipschitzness of the target functional:* Remember that the output is generated as $y := \tilde{F}^*(\mu_0^{(i^*)}, x_q) + \xi$, $\xi \sim \mathcal{N}(0, \sigma^2)$. The hidden functional $\tilde{F}^* : B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}) \times \mathcal{X}_q \rightarrow \mathbb{R}$ is assumed to be Lipschitz: Let the metric on the product set $B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}) \times \mathcal{X}_q$ be

$$d_{\text{prod}}((\mu_0, x), (\mu'_0, x')) := \|\mu_0 - \mu'_0\|_{\mathcal{H}_0^{\gamma_f}} + \|x - x'\|_2, \quad \gamma_f < 0,$$

where $\mu_0, \mu'_0 \in B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$ and $x, x' \in \mathcal{X}_q$. We write $\mu_0 \in \mathcal{H}_0$ by identifying μ_0 with its density in \mathcal{H}_0 up to an additive constant (see Remark 2). We assume \tilde{F}^* is in $\text{Lip}_L(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}) \times \mathcal{X}_q, d_{\text{prod}})$, a set of L -Lipschitz functionals with respect to d_{prod} .

These assumptions specify the analytic and structural conditions of the RKHS and the distinguishability of contexts, which will be imposed throughout the analysis. We provide an example for our assumptions:

Example 2 (Rapid eigenvalue decay in Assumption 6). *It is known (Grigor'yan, 2006) that on a compact Riemannian manifold \mathcal{M} the Laplace-Beltrami operator has a discrete spectrum satisfying Weyl's law: the heat kernel expansion is $p_t(x, y) = \sum_{k=0}^{\infty} e^{-\Theta(k^{2/n})t} \varphi_k(x) \varphi_k(y)$, where $\{\varphi_k\}_k$ are eigenfunctions, $n = \dim(\mathcal{M})$, and $t > 0$. Settings with rapid eigenvalue decay have been investigated as a data structure (Nadler et al., 2005; Coifman & Lafon, 2006; Xia & Shi, 2024). Gaussian kernels, despite being non-compactly supported, are also widely used in ML tasks (e.g. Schölkopf & Smola (2002)).*

Next, we restrict the input measures to be probability measures:

Assumption 7 (Probability assumption for Setting 3: Technical Version of Assumption 4). Let $\mathcal{P}(X)$ denote the set of Borel probability measures on a measurable set X . For $\mu_0 \sim \mathbb{P}_{\mu_0}$, $\mu_0 \in \mathcal{P}(\mathcal{X}_0)$ almost surely (with probability 1).

Based on these assumptions, we can summarize the probabilistic setting for our upper bound analysis:

Setting 3 (Probability Setting for Upper Bound: Technical Version of Setting 1). Assumption 6 and Assumption 7 are satisfied. In short: $\mu_0 \sim \mathbb{P}_{\mu_0}$ is constrained as a probability measure whose density is in RKHS ball.

A.1.2 STRUCTURED SETTING FOR LOWER BOUND.

For the minimax lower bound, we relax the probability constraint in Assumption 7 and instead impose structural conditions (following Lanthaler (2024)) on the coefficients of the Mercer expansion:

Assumption 8 (Structural assumptions for Setting 4: Technical Version of Assumption 5). Let (Ω, \mathbb{P}) be a probability space. Instead of constraining $\mu_0 \sim \mathbb{P}_{\mu_0}$ to be a probability measure, μ_0 is generated in Definition 1 as follows:

- *Random RKHS element:* The “density” $p_{\mu_0} : \Omega \rightarrow \mathcal{H}_0$ associated with μ_0 has the expansion

$$\frac{d\mu_0}{d\lambda}(\omega)(\cdot) := p_{\mu_0}(\omega)(\cdot) = \sum_{j=1}^{\infty} \lambda_j^{\gamma_d} Z_j(\omega) e_j(\cdot), \quad \omega \in \Omega,$$

where $(e_j)_{j \geq 1}$ is an L^2 -orthonormal basis on \mathcal{X}_0 , $\gamma_d > 0$, and $\lambda_1^{\gamma_d} \geq \lambda_2^{\gamma_d} \geq \dots \geq 0$ are summable.

- *Random coefficients:* The variables $Z_j : \Omega \rightarrow \mathbb{R}$ are jointly independent, satisfy $\mathbb{E}[|Z_j|^2] = 1$, $Z_j \sim \rho_j(z) dz$, and obey the uniform bounds $\sup_j \|\rho_j\|_{\infty} \leq R$, $\lambda_1^{\gamma_d/2} \leq R$ for some $R > 0$.

Example 3. *Assumption 6 and Assumption 8 are compatible: Take $\rho_j = \frac{1}{2} \mathbb{1}_{\{-1,1\}}$ and $\lambda_j \leq A \exp(-cj^\alpha)$. Then, $\|\mu_0\|_{\mathcal{H}_0^{\gamma_b}} \leq \sum_j \lambda_j^{\gamma_d - \gamma_b}$ always converges if $\gamma_d > \gamma_b$ and $\|\mu_0\|_{\mathcal{H}_0^{\gamma_b}} = O(A^{\gamma_d - \gamma_b})$. Therefore, there exists some $A > 0$ such that $\|\mu_0\|_{\mathcal{H}_0^{\gamma_b}} \leq 1$ a.s.*

Finally, we summarize the corresponding setting:

Setting 4 (Structured Setting for Lower Bound: Technical Version of Setting 2). Assumption 6 and Assumption 8 are satisfied. In short: the density is sampled in the form of Mercer expansion $d\mu_0/d\lambda := \sum_j \lambda_j^{\gamma_d} Z_j e_j$ where Z_j are independent r.v.s. and μ_0 may not be the probability measure.

B ESTIMATION ERROR ANALYSIS (UPPER BOUND)

The overarching goal of this section is to derive a statistical upper bound for transformer-based estimators in Setting 3. Specifically, we establish that the empirical risk minimizer achieves a convergence rate of the form

$$R(\hat{F}, F^*) \lesssim \exp(-\Omega((\log n)^{\alpha/(\alpha+1)})).$$

This rate can be regarded as the infinite-dimensional analogue of the classical $n^{-\Theta(1/d)}$ risk bound for d -dimensional regression, where the effective dimension scales as $d \sim (\log n)^{1/(\alpha+1)}$. In particular, although the associative recall task requires handling measure-valued components, our analysis demonstrates that its statistical complexity coincides with that of a pure infinite-dimensional regression problem, whose minimax lower bound will be shown in Appendix C. The subsequent subsections establish this result step by step, through successive approximation bounds for the individual network layers.

B.1 PROOF SKETCH FOR THE ESTIMATION UPPER BOUND

In this part, we will explain how to prove the following theorem:

Theorem 3 (Sub-Polynomial Convergence Corresponding to Theorem 1, A Simplified Version of Theorems 5 and 6). *Let $\tilde{F}^* \in \text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}) \times \mathcal{X}_q, d_{\text{prod}})$ and assume one of the following cases:*

- (i) *the number of mixture components is bounded as $I \leq d_1 \lesssim (\ln n)^{\frac{1}{\alpha+1}}$;*
- (ii) *the hidden target function \tilde{F}^* is independent of x_q and $I \leq d_1 \simeq n^{o(1)}$.*

Let \hat{F} be the empirical risk minimizer with the transformer class $\text{TF}(\epsilon)$ as the set of mappings $\text{Attn}_{\theta_2} \diamond \text{MLP}_{\xi_2} \diamond \text{Attn}_{\theta_1} \diamond \text{MLP}_{\xi_1}$ such that, $H_1, B_{a,1}, \ell_2 = (\log \epsilon^{-1})^{O(1)}$, $\|\mathbf{p}_2\|_{\infty}, s_2 \lesssim \exp(O((\log \epsilon^{-1})^{1+\alpha^{-1}}))$, $d_{\text{attn}2}, S'_{a,2}, H_2, B'_{a,2} = 1$, $S_{a,2}, B_{a,2} = 0$ in both cases, and

in (i): $\ell_1, \|\mathbf{p}_1\|_{\infty}, s_1, d_{\text{attn}1}, S_{a,1} = (\log \epsilon^{-1})^{O(1)}$, ; in (ii): $\ell_1, \|\mathbf{p}_1\|_{\infty}, s_1, d_{\text{attn}1}, S_{a,1} \lesssim \exp(O((\log \epsilon^{-1})^{1+\alpha^{-1}}))$. Then in Setting 3,

$$R(F^*, \hat{F}) \lesssim \exp(-\Omega((\ln n)^{\frac{\alpha}{\alpha+1}})),$$

where α is the decay rate of the eigenvalues of the underlying kernel of \mathcal{H}_0 .

We will provide a proof sketch for the first case: (i) the number of mixture components is bounded as $I \leq d_1 \lesssim (\ln n)^{\frac{1}{\alpha+1}}$.

To derive a statistical rate we must calibrate the size of the measure-theoretic transformer hypothesis class. Concretely, we choose an architecture that grants an ϵ -approximation of F^* while keeping the covering entropy $V(\delta) := \log \mathcal{N}(\text{TF}(\epsilon); \delta)_{\infty}$ minimal; the risk bound then follows by balancing the approximation error ϵ with the estimation term governed by $V(\cdot)$. In short, the general theory in Schmidt-Hieber (2020) asserts that the excess risk can be bounded by the sum of approximation terms and a complexity term: for any $\delta > 0$, the L^2 -risk R is bounded by

$$R(F^*, \hat{F}) \lesssim \inf_{\hat{F}} \|\hat{F} - F^*\|_{L^2(\mathbb{P}_{\nu, x_q})}^2 + \delta + \frac{V(\delta)}{n}$$

Here the first term quantifies how well the architecture approximates F^* , while the second reflects the statistical price of searching over a class of size $V(\delta)$. Our proof thus first controls $V(\cdot)$ layer-wise and then selects ϵ to realize the optimal trade-off.

Our estimation bound relies on two main ingredients: (i) an approximation strategy for representing the target functional $F^*(\nu, x_q) = \tilde{F}^*(\mu_0^{(i^*)}, x_q)$ via a depth- $L = 2$ transformer, and (ii) covering entropy bounds for each component of the architecture, combined through a composition lemma for measure-theoretic mappings.

Step 1: Composition lemma (Appendix B.2). We first state a generic result for the covering number of compositions of measure-theoretic maps.

Lemma 3 (Composition lemma). *Let $\mathcal{G}_i, i = 1, 2$ be sets of maps $\Gamma_i : \mathcal{P}(\mathcal{X}_i^{(1)}) \times \mathcal{X}_i^{(1)} \rightarrow \mathcal{X}_i^{(2)}$ such that $\mathcal{X}_1^{(2)} \subset \mathcal{X}_2^{(1)}$, $\mathcal{N}(\mathcal{G}_i; \epsilon)_\infty \lesssim N_i$, and any $\Gamma_2 \in \mathcal{G}_2$ is $(L_{2,1}, L_{2,2})$ -Lipschitz with respect to the 1-Wasserstein and Euclidean metrics. Then,*

$$\mathcal{N}(\{\Gamma_2 \circ \Gamma_1; \Gamma_i \in \mathcal{G}_i, i = 1, 2\}; \epsilon)_\infty \lesssim \mathcal{N}(\mathcal{G}_2; \frac{\epsilon}{2})_\infty \cdot \mathcal{N}(\mathcal{G}_1; \frac{\epsilon}{2(L_{2,1} + L_{2,2})})_\infty.$$

The proof follows the standard finite-dimensional composition argument: approximate each map Γ_i by the nearest covering element, and bound the difference of the composed maps using the Lipschitz constants.

Step 2: Approximation by a depth-2 transformer. We approximate $F^*(\nu, x) = \tilde{F}^*(\mu_0^{(i^*)}, x_q)$ using the following architecture, focusing on the first $D \gtrsim (\log \epsilon^{-1})^{\alpha-1}$ Mercer coefficients:

- First MLP layer (Appendix B.3).** Construct MLP_{ξ_1} to augment the input (x_1, x_2) with evaluations $(e_i(x_2))_{i=1}^D$ of an analytic basis $\{e_i\}$ up to $O(\epsilon_1)$ error. Analyticity implies $\log \mathcal{N}(\mathcal{F}_1; \epsilon_1)_\infty \lesssim \text{poly} \log \epsilon_1^{-1} \cdot \text{poly} D$.
- First attention layer (Appendix B.4)** Apply Attn_{θ_1} to $\text{MLP}_{\xi_1}(\mu)$ to compute empirical means $\int e_j(y_2) d\mu_0(y_2)$ up to $O(\epsilon_1 + \epsilon_2)$ error. We approximate a one-hot selection of measures and compute as, informally,

$$\int \underbrace{\text{Softmax}(\text{MLP}_1(x_q)^\top Q^\top K \text{MLP}_1(y))}_{\mathbb{1}[y \sim \mu_{v^{(i^*)}}^{(i^*)}]} \text{MLP}_1(y) \underbrace{d\nu(y)}_{\propto d \sum_i \mu_{v^{(i)}}^{(i)}} \simeq \int \text{MLP}_1(y) d\mu_{v^{(i^*)}}^{(i^*)}(y).$$

By the construction of the first layer MLP_1 and the product decomposition of $\mu_{v^{(i^*)}}^{(i^*)}$, the RHS approximates $\int e_j(y_2) d\mu_0(y_2)$. Under sparsity constraints on the attention matrices,

$$\log \mathcal{N}(\mathcal{A}_1(d_{\text{attn}}, H, B_a, S_a); \epsilon_2)_\infty \lesssim \text{poly} D \cdot \text{poly} \log(I\epsilon_2^{-1}).$$

- Second MLP layer (Appendix B.5).** Approximate a Lipschitz map on $\mathbb{R}^{D+d_1+d_2}$ whose inputs are retained Mercer coefficients $(\int e_j(y_2) d\mu_0(y_2))_{j=1}^D$ and query x_q . We also show $D \gtrsim (\log \epsilon^{-1})^{\alpha-1}$ (this is ϵ , not ϵ_1) is sufficient to extract the features. This layer may have a large Lipschitz constant $L_{\text{MLP},2} \lesssim e^{-O(c_\epsilon (\log \epsilon^{-1})^{\frac{2+2\alpha}{\alpha}})}$ for $c_\epsilon \lesssim \text{poly} \log \log \epsilon^{-1}$, which constrains $\epsilon_1 \vee \epsilon_2$ to be super-polynomially small via the measure-theoretic composition discussed in Lemma 5.
- Second attention layer (Appendix B.6).** Implemented as a (one-dimensional) skip connection, with $O(1)$ Lipschitz constant, added to fit the formal hypothesis set definition.

Step 3: Bounding the covering entropy (Appendix B.7). Applying the composition lemma recursively over the four layers yields

$$\begin{aligned} \log \mathcal{N}(\text{TF}; \epsilon)_\infty &\lesssim \log \mathcal{N}(\mathcal{A}(d + D, 1, 0, 1, 0, d + D); \epsilon)_\infty + \log \mathcal{N}(\mathcal{F}(\ell_2, p_2, s_2); \tilde{\Omega}(\epsilon))_\infty \\ &\quad + \log \mathcal{N}(\mathcal{A}(d_{\text{attn}}, H, B_a, S_a); \tilde{\Omega}(L_{\text{MLP},2}^{-1}\epsilon))_\infty \\ &\quad + \log \mathcal{N}(\mathcal{F}(\ell_1, p_1, s_1); \tilde{\Omega}(L_{\text{MLP},2}^{-1}(L_{\text{Attn},1,W^1} + L_{\text{Attn},1,\|\cdot\|_2})^{-1}\epsilon))_\infty. \end{aligned}$$

This yields (also carefully bounding the term with respect to d_1 omitted above):

Lemma 4. *The covering entropy of the transformer class satisfies*

$$\log \mathcal{N}(\text{TF}; \epsilon)_\infty \lesssim \exp\left(O((\log \epsilon^{-1})^{(1+\min(\alpha,\beta))/\min(\alpha,\beta)})\right)$$

where α is the decay rate of RKHS \mathcal{H}_0 and $d_1 \simeq (\ln \epsilon^{-1})^{\beta-1}$.

Step 4: From covering entropy to risk bound (Appendix B.7). Applying the regression bound shown in Schmidt-Hieber (2020), with $V(\delta) = \log \mathcal{N}(\text{TF}; \delta)_\infty$, and choosing

$$\epsilon \simeq \exp\left(-\Theta(\log n)^{\frac{\min(\alpha, \beta)}{\min(\alpha, \beta)+1}}\right),$$

where $d_1 \simeq (\ln \epsilon^{-1})^{\beta-1}$, we obtain the sub-polynomial convergence rate:

Theorem 4 (Sub-polynomial convergence, a generalized version of Theorem 1). *In Probability Setting (Setting 3),*

$$\sup_{F^* \in \mathcal{F}^*} R(\hat{F}, F^*) \lesssim \exp\left(-\Omega\left((\log n)^{\frac{\min(\alpha, \beta)}{\min(\alpha, \beta)+1}}\right)\right),$$

where the eigenvalues are $\lambda_j \simeq \exp(-cj^\alpha)$ and the number of mixture components is $I \leq d_1 \lesssim (\ln n)^{\beta-1 \cdot \min(\alpha, \beta) / (\min(\alpha, \beta)+1)}$.

B.2 STEP 1: COMPOSITION LEMMA.

Lemma 5 (Composition Lemma. Restated). *Let \mathcal{G}_i , $i = 1, 2$ be sets of Γ_i , which are maps from $\mathcal{P}(\mathcal{X}_i^{(1)}) \times \mathcal{X}_i^{(1)}$ to $\mathcal{X}_i^{(2)}$ such that $\mathcal{X}_1^{(2)} \subset \mathcal{X}_2^{(1)}$, $\mathcal{N}(\mathcal{G}_i; \epsilon)_\infty \lesssim N_i$, and any $\Gamma_2 \in \mathcal{G}_2$ is $(L_{2,1}, L_{2,2})$ -Lipschitz with respect to 1-Wasserstein distance and Euclidean distance. Then, we have*

$$\mathcal{N}(\{\Gamma_2 \diamond \Gamma_1 \mid \Gamma_i \in \mathcal{G}_i\}; \epsilon)_\infty \lesssim \mathcal{N}(\mathcal{G}_2; \frac{\epsilon}{2})_\infty \cdot \mathcal{N}(\mathcal{G}_1; \frac{\epsilon}{2(L_{2,1} + L_{2,2})})_\infty.$$

Proof. First, remember that, for standard one-dimensional function classes \mathcal{F}, \mathcal{G} with covering numbers N_f, N_g ,

$$\mathcal{N}(\{f \circ g \mid f \in \mathcal{F}, g \in \mathcal{G}\}; \epsilon) \lesssim \mathcal{N}(\mathcal{F}; \epsilon) \cdot \mathcal{N}(\mathcal{G}; \epsilon/L_f)$$

where L_f is the upperbound of Lipschitz constants of $\forall f \in \mathcal{F}$. For measure-theoretic mappings $\Gamma_1 \in \mathcal{G}_1$ and $\Gamma_2 \in \mathcal{G}_2$, take Γ_1^i and Γ_2^j be the $(\epsilon_1$ and ϵ_2 nearest covering elements (i.e. $\sup_{\mu, x \in \mathcal{P}(\mathcal{X}_i^{(1)}) \times \mathcal{X}_i^{(1)}} |\Gamma_1(\mu, x) - \Gamma_1^i(\mu, x)| \leq \epsilon)$. Then we bound the difference of compositions as

$$\begin{aligned} & |(\Gamma_2 \diamond \Gamma_1)(\mu, x) - (\Gamma_2^j \diamond \Gamma_1^i)(\mu, x)| \\ &= |(\Gamma_2 \diamond \Gamma_1)(\mu, x) - (\Gamma_2 \diamond \Gamma_1^i)(\mu, x)| + |(\Gamma_2 \diamond \Gamma_1^i)(\mu, x) - (\Gamma_2^j \diamond \Gamma_1^i)(\mu, x)| \\ &\leq L_{2,1} W_1(\Gamma_1(\mu)_\# \mu, \Gamma_1^i(\mu)_\# \mu) + L_{2,2} \|\Gamma_1^i(\mu, x) - \Gamma_1(\mu, x)\|_2 + \epsilon_2 \\ &\leq (L_{2,1} + L_{2,2})\epsilon_1 + \epsilon_2. \end{aligned}$$

□

B.3 STEP 2-1: FIRST MLP LAYER

We begin by formalizing the approximation properties of the first MLP layer, which is responsible for embedding both the input tokens and auxiliary analytic features into a higher-dimensional representation. This layer plays a crucial role in ensuring that subsequent attention and MLP layers can operate on a sufficiently expressive feature space.

Lemma 6 (E & Wang (2018)). *Let f be an analytic function over $[-M, M]^{d_2}$ such that $\mathcal{X}_0 \subset [-M + \delta, M - \delta]^{d_2}$ for some $\delta > 0$ and the power series $f(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{N}^{d_2}} a_{\mathbf{k}} \mathbf{x}^{\mathbf{k}}$ is absolutely convergent over $[-M, M]^{d_2}$. Then, a deep ReLU network \hat{f} with depth $O((\log \epsilon^{-1})^{2d_2})$ and width $d_2 + 4$ (independent of ϵ) satisfies*

$$\sup_{x \in \mathcal{X}_0} |f(x) - \hat{f}(x)| \lesssim \epsilon.$$

Here, the notation $\mathbf{x}^{\mathbf{k}}$ denotes the multivariate monomial $\prod_{i=1}^{d_2} x_i^{k_i}$, and the absolute convergence condition ensures that the power series uniformly converges on $[-M, M]^{d_2}$, enabling uniform approximation on the interior domain \mathcal{X}_0 .

As a corollary of Lemma 6, we obtain the following result:

Corollary 1. Let $d = d_1 + d_2$. For a function f such that

$$f : [-M, M]^{d_1+d_2} \rightarrow \mathbb{R}^{d+D}, \quad f(x_1, x_2) = \begin{bmatrix} x_1 \\ x_2 \\ e_1(x_2) \\ \vdots \\ e_D(x_2) \end{bmatrix},$$

where e_j are defined in Assumption 6. That is, f preserves the first $d_1 + d_2$ coordinates (x_1, x_2) and augments them with D analytic feature functions e_j that depend only on x_2 . There exists a network $\hat{f} \in \mathcal{F}(\ell_1, \mathbf{p}_1, s_1, \infty)$, where $\ell_1 \lesssim (\log \epsilon_1^{-1})^{2d_2}$, $p_{1,j} \lesssim D + d$, $s_1 \lesssim d + (\log \epsilon_1^{-1})^{2d_2} \cdot d_2^2 \cdot D$, such that

$$\|f - \hat{f}\|_\infty \lesssim \epsilon_1.$$

Proof. For the first d_1 indices, we simulate $x_{1,i} = -\text{ReLU}(-\mathbf{e}_i \mathbf{e}_i^\top \cdot x_1) + \text{ReLU}(\mathbf{e}_i \mathbf{e}_i^\top \cdot x_1)$. This requires $O(1)$ depth, $O(d_1)$ width, and $O(d_1)$ parameters. The $i + d_1$ -th ($i = 1, \dots, d_2$) indices require $O(1)$ depth, $O(d_2)$ width, and $O(d_2)$ parameters. We refer to Lemma 6 for the rest of indices. \square

In other words, each analytic component e_j can be uniformly approximated by a ReLU network of logarithmic number of parameters, and the concatenated mapping f can be represented by a block-structured network with parameter bounds as stated.

To evaluate the covering number, we use the following lemma:

Lemma 7 (Schmidt-Hieber (2020)). Let $V = \prod_{i=0}^{\ell} (p_i + 1)$. Then,

$$\log \mathcal{N}(\mathcal{F}(\ell, p, s, F); \delta)_\infty \lesssim (s + 1) \log(\delta^{-1} \ell V).$$

The covering number is

$$\begin{aligned} \log \mathcal{N}(\mathcal{F}(\ell_1, p_1, s_1); \epsilon_1)_\infty &\lesssim (d + (\log \epsilon_1^{-1})^{2d_2} \cdot D) \cdot (\log \epsilon_1^{-1} + \log \log \epsilon_1^{-1} + (\log \epsilon_1^{-1})^{2d_2} \log(d + D)) \\ &\lesssim C_{D, \epsilon_1} (d + D (\log \epsilon_1^{-1})^{4d_2}), \end{aligned}$$

where $C_{D, \epsilon_1} \lesssim \text{poly log } D + \text{poly log } d + \text{poly log log } \epsilon_1^{-1}$.

B.4 STEP 2-2: FIRST ATTENTION LAYER

In the preceding section, we have constructed an MLP layer capable of approximating the basis functions e_j (Assumption 6) with high accuracy. The role of the first attention layer is now to process an input mixture measure ν_f , identify the *associated measure* corresponding to a given query component $y_{f,i}$, and then output the concatenation of the raw query coordinates with the integrals of e_j against that associated measure. Formally, this operation is realized by the mapping ϕ_2 defined below.

We now analyze the approximation properties and complexity of the first attention layer in our architecture. Recall that the attention operator Attn_θ has already been defined in the measure-theoretic form

$$\text{Attn}_\theta : \mathcal{P}(\mathbb{R}^{d_{\text{attn}}}) \times \mathbb{R}^{d_{\text{attn}}} \rightarrow \mathbb{R}^{d_{\text{attn}}},$$

where the first argument is a probability measure over token representations and the second argument is the query vector. The following lemmas show that, under appropriate structural assumptions on the input measures and functions:

1. the target mapping ϕ_2 can be realized to accuracy ϵ_2 by a member of the attention class $\mathcal{A}(d_{\text{attn}}, H, B_a, S_a)$ (Lemma 8);
2. such attention mappings are Lipschitz continuous with an explicit bound in terms of the model parameters (Lemma 9);
3. the ϵ_2 -covering number of \mathcal{A} admits an upper bound in the parameter regime above (Lemma 10).

We present these results in turn.

Before presenting Lemma 8, we clarify the role of the mapping ϕ_2 in the composition-of-maps view (cf. Definition 5). In our construction, the first MLP layer ϕ_1 approximates the Mercer features:

$$\phi_1(\mu, x) \simeq \begin{bmatrix} x_1 \\ \vdots \\ x_d \\ e_1(x_{d_1+1:d}) \\ \vdots \\ e_D(x_{d_1+1:d}) \end{bmatrix},$$

where $\{e_j\}_{j \geq 1}$ is the Mercer (RKHS) eigenbasis on $\mathcal{X}_0 \subset \mathbb{R}^{d_2}$. Accordingly, the push-forward measure after ϕ_1 is

$$\mu_1 := (\phi_1(\mu, \cdot))_{\#} \mu \in \mathcal{P}(\mathbb{R}^{d+D}),$$

and the composition rule yields

$$(\phi_2 \diamond \phi_1)(\mu, x) = \phi_2(\mu_1, \phi_1(\mu, x)).$$

In the present setting, ϕ_2 preserves the first d coordinates and replaces the last D coordinates by the (component-wise) integrals of the associated measure against the Mercer basis:

$$\phi_2(\mu_1, \phi_1(\mu, x)) = \begin{bmatrix} x_1 \\ \vdots \\ x_d \\ \int e_1 d\mu_0^{(i^*)} \\ \vdots \\ \int e_D d\mu_0^{(i^*)} \end{bmatrix},$$

where i^* denotes the index of the associated component selected by the query. Moreover, if $\mu_0^{(i^*)}$ admits a (Borel) density w.r.t. a reference measure λ , say $\frac{d\mu_0^{(i^*)}}{d\lambda} = f_{\mu_0^{(i^*)}}$ with $f_{\mu_0^{(i^*)}} \in L^2(\lambda)$ (cf. Remark 2), then, writing the Mercer expansion $f_{\mu_0^{(i^*)}} = \sum_{j \geq 1} b_j e_j$, we have, up to the immaterial constant term,

$$\int e_j d\mu_0^{(i^*)} = \int e_j f_{\mu_0^{(i^*)}} d\lambda = \langle e_j, f_{\mu_0^{(i^*)}} \rangle_{L^2(\lambda)} = b_j.$$

Hence ϕ_2 produces, in its last D coordinates, the (truncated) Mercer coefficients of the associated density.

Lemma 8 shows that, under our structural assumptions on the input measures and the transformation f , the target mapping ϕ_2 can be uniformly approximated to accuracy ϵ_2 by an attention mechanism with bounded parameters.

Intuitions of Lemma 8. The key point is in the structure of the first attention layer: For a fixed j , to extract the *associated* Mercer coefficient $b_j = \int e_j d\mu_0^{(i^*)}$, we construct QK-matrix W_{QK}^j such that, with tokens mapped by the (simplified) first MLP layer $(x_1, x_2) \mapsto \psi_j(x_1, x_2) = (x_1, e_j(x_2))$,

$$\psi_j(x_q)^\top W_{QK}^j \psi_j(y) = \begin{cases} \gg 1 & \text{if } i = i^*; \\ \leq 0 & \text{if } i \neq i^*, \end{cases} \quad \text{for } y \sim \mu_{v(i)}^{(i)} = \delta_{v(i)} \otimes \mu_0^{(i)}.$$

Then, the softmax value will be

$$\text{Softmax}(\psi_j(x_q)^\top W_{QK}^j \psi_j(y)) \simeq \mathbb{1}[i = i^*].$$

The construction is simple: take $W_{QK}^j \propto \sum_{k=1}^{d_1} e_k e_k^\top$ such that $\psi_j(x_q)^\top W_{QK}^j \psi_j(y) \propto \langle v^{(i^*)}, v^{(i)} \rangle$, and multiply a large scalar, where e_k is a k -th one-hot vector, since ψ preserves the

- $I \leq d_1$ and $\text{supp}(f_{\#}\mu_i) \subset [-B_y, B_y]^{d+D}$ for all i .

Moreover, for each i , define $x_i := [v_i^\top, 0^\top]^\top \in \mathbb{S}^{d_1} \times \{0_{d_2}\}$ and set $y_{f,i} := f(x_i)$.

Define the mapping $\phi_2 : \mathcal{P}_f \times \mathbb{R}^{d+D} \rightarrow \mathbb{R}^{d+D}$ by

$$\phi_2(\nu_f, y_{f,i}) := \begin{bmatrix} y_{f,i,1} \\ \vdots \\ y_{f,i,d} \\ \int d(\tilde{f}_1)_{\#} \tilde{\mu}_i \\ \vdots \\ \int d(\tilde{f}_D)_{\#} \tilde{\mu}_i \end{bmatrix},$$

where \tilde{f}_j denotes the j -th coordinate function of \tilde{f} .

Then, there exists an attention operator $\hat{\text{Attn}} \in \mathcal{A}(d_{\text{attn}}, H, B_a, S_a)$ such that

$$\sup_{\nu_f \in \mathcal{P}_f, y_{f,i}, v_i \in \mathbb{R}^d} \left\| \phi_2(\nu_f, y_{f,i}) - \hat{\text{Attn}}(\nu_f, y_{f,i}) \right\|_{\infty} \leq \epsilon_2,$$

where $d_{\text{attn}} = d + D$, $H = D$, $B_a \lesssim \sqrt{\log(I \cdot \epsilon_2^{-1})}$, and $S_a = d$.

Proof. Fix an arbitrary $h \in \{1, \dots, H = D\}$. We first specify the attention weight matrices as follows:

$$W^h = V^h = e_{d+h} e_{d+h}^\top,$$

where $d = d_1 + d_2$ and e_{d+h} denotes the $(d+h)$ -th standard basis vector in \mathbb{R}^{d+D} . Similarly, define

$$Q^h = K^h = c \begin{bmatrix} I_{d_1} & O \\ O & O_{d_2+D} \end{bmatrix},$$

for a sufficiently large constant $c \gtrsim \sqrt{\log(I^3 \cdot \epsilon_2^{-1})} \gtrsim \sqrt{\log(I \cdot \epsilon_2^{-1})}$.

The corresponding attention weight for a query x_i and a key $f(y)$ is given by

$$\begin{aligned} & \text{Softmax}(\langle Q^h x_i, K^h f(y) \rangle) \\ &= \frac{\exp(c^2 \langle x_{i,1:d_1}, y_{1:d_1} \rangle)}{\int \exp(c^2 \langle x_{i,1:d_1}, y_{1:d_1} \rangle) d\nu_f(y)} \\ &= \frac{\exp(c^2 \langle x_{i,1:d_1}, y_{1:d_1} \rangle)}{I^{-1} \sum_{i'=1}^I \int \exp(c^2 \langle x_{i,1:d_1}, y_{1:d_1} \rangle) d(\delta_{v_{i'}} \otimes \tilde{\mu}_{i'})(y)} \\ &= \frac{\exp(c^2 \langle x_{i,1:d_1}, y_{1:d_1} \rangle)}{I^{-1} \sum_{i'=1}^I \exp(c^2 \langle v_{i'}, y_{1:d_1} \rangle)} \\ &= \frac{\exp(c^2)}{\frac{1}{I} (\exp(c^2) + (I-1))} \mathbb{1}_{\{y_{1:d_1} = v_i\}} \quad (v_i \perp v_j, i \neq j) \\ &= I \cdot \mathbb{1}_{\{y_{1:d_1} = v_i\}} + O(I^{-1} \epsilon_2), \end{aligned}$$

where the indicator function $\mathbb{1}_{\{y_{1:d_1} = v_i\}}$ arises because the keys $y_{1:d_1}$ take values in the finite set $\{v_1, \dots, v_I\}$ with mutually nonpositive inner products. The last equality follows from the choice $c \gtrsim \sqrt{\log(I \cdot \epsilon_2^{-1})}$, which ensures exponential separation of the correct key from the others.

Next, applying the value and output projection matrices, we obtain

$$W^h \int \text{Softmax}(\langle Q^h x_i, K^h y \rangle) V^h y d\nu_f(y)$$

$$\begin{aligned}
&= \frac{1}{I} \sum_{i'=1}^I W^h \int \text{Softmax}(\langle Q^h x_i, K^h f(y) \rangle) V^h f(y) d\mu_{i'}(y) \\
&= \frac{1}{I} \sum_{i'=1}^I \int \left(I \cdot \mathbb{1}_{\{y_{1:d_1}=v_i\}} + O(I^{-1}\epsilon_2) \right) e_{d+h} \left(\tilde{f}_h(y_{d_1+1:d}) + O(\epsilon_1) \right) d(\delta_{v_{i'}} \otimes \tilde{\mu}_{i'})(y) \\
&\quad \left(\text{where } e_{d+h}^\top \tilde{f} = \tilde{f}_h \right) \\
&= \frac{1}{I} \sum_{i'=1}^I \left(\int \left(I \cdot \mathbb{1}_{\{y_{1:d_1}=v_i\}} + O(I^{-1}\epsilon_2) \right) d\delta_{v_{i'}}(y_{1:d_1}) \right) \left(\int \tilde{f}_h(y_{d_1+1:d}) d\tilde{\mu}_{i'}(y_{d_1+1:d}) \right) e_{d+h} \\
&= \left(\int \tilde{f}_h(y_{d_1+1:d}) d\tilde{\mu}_i(y_{d_1+1:d}) \right) e_{d+h} + O(\epsilon_2).
\end{aligned}$$

Finally, to incorporate the skip connection over the first d coordinates, let

$$A = \begin{bmatrix} I_d & O \\ O & O \end{bmatrix}.$$

Applying A to the input vector $y_{f,i}$ yields

$$Ay_{f,i} = \begin{bmatrix} y_{f,i,1} \\ \vdots \\ y_{f,i,d} \\ 0_D \end{bmatrix}.$$

Combining the attention output for each head h with this skip connection reproduces the target mapping ϕ_2 up to an error of order $O(\epsilon_2)$ in the ℓ_∞ norm. This establishes the desired approximation property. \square

Remark 3 (Why do we need a softmax attention?). We informally demonstrate how linear attentions struggle with one-hot selection of densities without orthogonality. The main problem is that the context vectors $(v^{(i)})_i$ may have a negative correlation. For example, we consider $I = 2$ and

$$v^{(1)} = -v^{(2)}.$$

If we only have access to a linear attention, with the same QK matrices in the lemma,

$$\text{LinAttn}(\langle Q^h f(x_i), K^h f(y) \rangle) \simeq \langle v^{(i)}, v^{(i^*)} \rangle = \begin{cases} 1 & \text{if } i = i^*; \\ -1 & \text{if } i \neq i^*. \end{cases}$$

This implies that it is hard for linear attentions to extract *only* the i^* -th measure through integration $\int \text{LinAttn}(\langle Q^h x_i, K^h y \rangle) V^h y d\nu_f(y)$. See Han et al. (2024); Fan et al. (2025) for empirical discussions; see also Kim et al. (2024), where strong assumptions such as relaxed sparsity and orthogonality of recall candidates were required to bypass this difficulty.

Thus, the first attention layer is expressive enough to implement the ‘‘association and extraction’’ operation: given a mixture, it can select the relevant component and compute the e_j -integrals needed for downstream processing. We next turn to the *stability* of such an operator with respect to perturbations in both the measure and the query vector.

The following lemma establishes a Lipschitz property of Attn_θ in both arguments. This lemma is inspired by Vuckovic et al. (2020). This quantitative stability will be essential for the subsequent covering number analysis.

The attention operator Attn_θ computes a weighted average of values using a softmax over inner products $\langle Q^h x, K^h y \rangle$. To bound its change when (μ, x) varies, we split the effect of μ and x .

For the measure part, Kantorovich–Rubinstein duality expresses the 1-Wasserstein distance W_1 (the standard Wasserstein metric) as the supremum of expectation differences over 1-Lipschitz functions, allowing us to control the change via the Lipschitz constant of the softmax kernel.

For the query part, we directly bound the kernel's Lipschitz dependence on x and apply sparsity of the output projection. Combining both yields the stated Lipschitz bound.

Lemma 9. *Let $\text{Attn}_\theta \in \mathcal{A}(d_{\text{attn}}, H, B_a, S_a)$ be the attention operator as defined in Section 3. Assume that the query inputs satisfy $\|x_1\|_\infty, \|x_2\|_\infty \leq B_x$, and that for each $i \in \{1, 2\}$, every $y \in \text{supp}(\mu_i)$ satisfies $\|y\|_\infty \leq B_y$. Then Attn_θ is Lipschitz in the joint variable (μ, x) in the sense that*

$$\begin{aligned} & \|\text{Attn}_\theta(\mu_1, x_1) - \text{Attn}_\theta(\mu_2, x_2)\|_\infty \\ & \lesssim H \exp(O(S_a^2 B_a^2 B_x B_y)) \cdot (W_1(\mu_1, \mu_2) + \|x_1 - x_2\|_2). \end{aligned}$$

Moreover, if $d_{\text{attn}} \lesssim D$, $H \lesssim D$, $B_a \lesssim \sqrt{\log(I\epsilon_2^{-1})}$, and $S_a \lesssim d$, $B_x, B_y \lesssim 1$, then the Lipschitz constant is bounded by $D \exp(O(d^2 \log(I\epsilon_2^{-1})))$.

Proof. Bounding the difference in μ . We first bound the difference in the μ -variable while keeping the query fixed. By (P-i) the matrix sparsity bound $\|Av\|_\infty \leq sb\|v\|_\infty$ when A has at most s nonzero entries per row and each entry bounded by b , we have

$$\begin{aligned} & \|\text{Attn}_\theta(\mu_1, x_1) - \text{Attn}_\theta(\mu_2, x_1)\|_\infty \\ & \leq \sum_h S_a^2 B_a^2 \left\| \int \frac{y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_1(z)} d\mu_1 - \int \frac{y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} d\mu_2 \right\|_\infty \end{aligned}$$

Next, Inserting intermediate terms to align denominators and numerators, we obtain,

$$\begin{aligned} & \|\text{Attn}_\theta(\mu_1, x_1) - \text{Attn}_\theta(\mu_2, x_1)\|_\infty \\ & \leq \sum_h S_a^2 B_a^2 \left\| \int \frac{y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_1(z)} d\mu_1 - \int \frac{y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} d\mu_1 \right. \\ & \quad \left. + \int \frac{y \exp(\langle Q^h x_1, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} d\mu_1 - \int \frac{y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} d\mu_2 \right\|_\infty \\ & \leq \underbrace{\sum_h S_a^2 B_a^2 \left\| \int \frac{y \exp(\langle Q^h x_1, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_1(z)} d\mu_1 - \int \frac{y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} d\mu_1 \right\|_\infty}_{(i)} \\ & \quad + \underbrace{\sum_h S_a^2 B_a^2 \left\| \int \frac{y \exp(\langle Q^h x_1, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} d\mu_1 - \int \frac{y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} d\mu_2 \right\|_\infty}_{(ii)}. \end{aligned}$$

Bounding the term (i). We have

$$\begin{aligned} (i) & \leq \sum_h S_a^2 B_a^2 \left\| \int y \exp(\langle Q^h x_1, K^h y \rangle) d\mu_1 \right\|_\infty \\ & \quad \times \left| \frac{1}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_1(z)} - \frac{1}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} \right| \\ & \leq \sum_h S_a^2 B_a^2 \left\| \int y \exp(\langle Q^h x_1, K^h y \rangle) d\mu_1 \right\|_\infty \\ & \quad \times \left(\min \left\{ \int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_1(z), \int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z) \right\} \right)^{-2} \\ & \quad \times \left| \int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_1(z) - \int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z) \right| \quad (\text{by (P-iii)}) \\ & \lesssim \sum_h S_a^2 B_a^2 B_y \exp(3S_a^2 B_a^2 B_x B_y) \left| \int \exp(\langle Q^h x_1, K^h z \rangle) d(\mu_1 - \mu_2)(z) \right| \quad (\text{by (P-ii)}), \end{aligned}$$

where we used:

(P-i) the matrix sparsity bound $\|Av\|_\infty, \|Av\|_1 \leq sb\|v\|_\infty$ when A has at most s nonzero entries per row and each entry bounded by b ;

(P-ii) $|\langle Q^h x_1, K^h y \rangle| \leq S_a B_a B_x \cdot S_a B_a B_y$. Indeed, since $\|Q^h\|_0, \|K^h\|_0 \leq S_a$ (total number of nonzero entries) and $|Q_{jk}^h|, |K_{jk}^h| \leq B_a$, while $\|x_1\|_\infty \leq B_x$ and $\|y\|_\infty \leq B_y$, we have

$$\begin{aligned} & |\langle Q^h x_1, K^h y \rangle| \\ & \leq \|Q^h x_1\|_\infty \|K^h y\|_1 \\ & \leq \left(\max_j \sum_k |Q_{jk}^h| |x_{1,k}| \right) \sum_{j,k} |K_{jk}^h| |y_k| \\ & \leq (S_a B_a B_x) (S_a B_a B_y). \end{aligned}$$

Here the bound $\|Q^h x_1\|_\infty \leq S_a B_a B_x$ follows because each coordinate of $Q^h x_1$ is a sum of at most S_a terms, each of magnitude at most $B_a B_x$; similarly, $\|K^h y\|_1 \leq \sum_{j,k} |K_{jk}^h| |y_k| \leq S_a B_a B_y$ since there are at most S_a nonzero matrix entries in total;

(P-iii) the bound

$$\left| \frac{1}{\alpha_1} - \frac{1}{\alpha_2} \right| \leq A^{-2} |\alpha_1 - \alpha_2|, \quad A < \min(\alpha_1, \alpha_2).$$

The RHS is bounded as

$$\begin{aligned} ((i) \lesssim) & \sum_h S_a^2 B_a^2 B_y \exp(3S_a^2 B_a^2 B_x B_y) \left| \int \exp(\langle Q^h x_1, K^h z \rangle) d(\mu_1 - \mu_2)(z) \right| \\ & \lesssim H S_a^4 B_a^4 B_x B_y \exp(4S_a^2 B_a^2 B_x B_y) W_1(\mu_1, \mu_2) \end{aligned}$$

using the Kantorovich–Rubinstein duality

$$W_1(\mu_1, \mu_2) = \sup_{\text{Lip}(\phi) \leq 1} \int \phi d(\mu_1 - \mu_2),$$

and the fact that $y \mapsto \exp(\langle Q^h x_1, K^h y \rangle)$ is $S_a^2 B_a^2 B_x \exp(S_a^2 B_a^2 B_x B_y)$ -Lipschitz on $[-B_y, B_y]^{d_{\text{attn}}}$ because

$$\begin{aligned} & |\exp(\langle Q^h x_1, K^h y_3 \rangle) - \exp(\langle Q^h x_1, K^h y_4 \rangle)| \\ & \leq \exp(S_a^2 B_a^2 B_x B_y) |\langle Q^h x_1, K^h (y_3 - y_4) \rangle| \\ & \leq S_a^2 B_a^2 B_x \exp(S_a^2 B_a^2 B_x B_y) \|y_3 - y_4\|_2. \end{aligned}$$

for $y_3, y_4 \in [-B_y, B_y]^{d_{\text{attn}}}$.

Bounding the term (ii). We have

$$\begin{aligned} (ii) & \lesssim \sum_h S_a^2 B_a^2 \left| \frac{1}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_2(z)} \right| \left\| \int y \exp(\langle Q^h x_1, K^h y \rangle) d(\mu_1 - \mu_2)(y) \right\|_\infty \\ & \lesssim \sum_h S_a^2 B_a^2 \exp(S_a^2 B_a^2 B_x B_y) \max_{i=1, \dots, d_{\text{attn}}} \left| \int y_i \exp(\langle Q^h x_1, K^h y \rangle) d(\mu_1 - \mu_2)(y) \right| \\ & \lesssim H(1 + S_a^2 B_a^2 B_x B_y) S_a^2 B_a^2 \exp(2S_a^2 B_a^2 B_x B_y) W_1(\mu_1, \mu_2) \end{aligned}$$

using the Kantorovich–Rubinstein duality

$$W_1(\mu_1, \mu_2) = \sup_{\text{Lip}(\phi) \leq 1} \int \phi d(\mu_1 - \mu_2),$$

and the fact that $y \mapsto y_i \exp(\langle Q^h x_1, K^h y \rangle)$ is $(1 + S_a^2 B_a^2 B_x B_y) \exp(S_a^2 B_a^2 B_x B_y)$ -Lipschitz on $[-B_y, B_y]^{d_{\text{attn}}}$.

Finally, we have

$$\|\text{Attn}_\theta(\mu_1, x_1) - \text{Attn}_\theta(\mu_2, x_1)\|_\infty$$

$$\lesssim HS_a^2 B_a^2 (S_a^2 B_a^2 B_x B_y \exp(4S_a^2 B_a^2 B_x B_y) + (1 + S_a^2 B_a^2 B_x B_y) \exp(2S_a^2 B_a^2 B_x B_y)) \\ \times W_1(\mu_1, \mu_2).$$

Bounding the difference in x . Next, we bound the difference in the query x . For fixed μ_1 , using similar interpolation and Lipschitz estimates in x ,

$$\begin{aligned} & \|\text{Attn}_\theta(\mu_1, x_1) - \text{Attn}_\theta(\mu_1, x_2)\|_\infty \\ & \leq \sum_h S_a^2 B_a^2 \left\| \int \frac{y \exp(\langle Q^h x_1, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_1(z)} - \frac{y \exp(\langle Q^h x_2, K^h y \rangle)}{\int \exp(\langle Q^h x_2, K^h z \rangle) d\mu_1(z)} d\mu_1(y) \right\|_\infty \\ & \quad + S_a B_a \|x_1 - x_2\|_\infty \\ & \leq \sum_h S_a^2 B_a^2 \left| \int \frac{y \exp(\langle Q^h x_1, K^h y \rangle)}{\int \exp(\langle Q^h x_1, K^h z \rangle) d\mu_1(z)} d\mu_1 - \int \frac{y \exp(\langle Q^h x_1, K^h y \rangle)}{\int \exp(\langle Q^h x_2, K^h z \rangle) d\mu_1(z)} d\mu_1(y) \right. \\ & \quad \left. + \int \frac{y \exp(\langle Q^h x_1, K^h y \rangle)}{\int \exp(\langle Q^h x_2, K^h z \rangle) d\mu_1(z)} d\mu_1 - \int \frac{y \exp(\langle Q^h x_2, K^h y \rangle)}{\int \exp(\langle Q^h x_2, K^h z \rangle) d\mu_1(z)} d\mu_1(y) \right| \\ & \quad + S_a B_a \|x_1 - x_2\|_2 \\ & \lesssim (S_a B_a \vee HS_a^4 B_a^4 B_x B_y \exp(4S_a^2 B_a^2 B_x B_y)) \|x_1 - x_2\|_2. \end{aligned}$$

□

By combining approximation and stability, we can control the complexity of the attention class via its covering number, as stated next.

We now bound the ϵ_2 -covering number of $\mathcal{A}(d_{\text{attn}}, H, B_a, S_a)$ in the parameter regime of interest.

Lemma 10. *The ϵ_2 -covering number of $\mathcal{A}(d_{\text{attn}}, H, B_a, S_a)$ is bounded by*

$$\begin{aligned} & \mathcal{N}(\mathcal{A}(d_{\text{attn}}, H, B_a, S_a); \epsilon_2)_\infty \\ & \lesssim (d_{\text{attn}}^2 \cdot \exp(O(\log(H) + S_a^2 B_a^2 B_x B_y)) (\epsilon_2^{-1} + 1))^{O(S_a H)}. \end{aligned}$$

Furthermore, if $d_{\text{attn}} \lesssim d + D$, $H = O(D)$, $B_a \lesssim \sqrt{\log(I \epsilon_2^{-1})}$, $S_a \lesssim d$, and $B_x, B_y = O(1)$, then the covering entropy is

$$\log \mathcal{N}(\mathcal{A}(d_{\text{attn}}, H, B_a, S_a), \epsilon_2) \lesssim C_D \cdot D d^3 (\log I + \log \epsilon_2^{-1})^2 \cdot \log \epsilon_2^{-1}$$

where $C_{d,D} \lesssim \text{poly}(\log D + \log d)$.

Proof. We define a $\Omega(\bar{\epsilon})$ -covering set of $\mathcal{A}(d_{\text{attn}}, H, B_a, S_a)$ as a set of mappings whose parameters can be constructed as follows:

- For each matrix W^h, Q^h, K^h, V^h in each $h \in \{1, \dots, H\}$,
 1. Choose S_a matrix entries among $O(d_{\text{attn}}^2)$ entries.
 2. For each matrix entry,
 - Set its value from $\{(j \cdot \tilde{\epsilon} - 1) \cdot B_a \mid j = 0, \dots, \lceil 2\tilde{\epsilon}^{-1} \rceil\}$ where $\tilde{\epsilon} \gtrsim \exp(-C_2(\log(H) + S_a^2 B_a^2 B_x B_y)) \bar{\epsilon}$ where C_2 is a sufficiently large constant.
- Set the value of chosen S_a entries in A from $\{(j \cdot \tilde{\epsilon} - 1) \cdot B_a \mid j = 0, \dots, \lceil 2\tilde{\epsilon}' \rceil\}$ where $\tilde{\epsilon}' \simeq (S_a B_a)^{-1} \bar{\epsilon}$.

Let us prove that the above set of mappings is a $\Omega(\bar{\epsilon})$ -covering. It is clear that for W^h $h = 1, \dots, H$, there exist matrices \hat{W}^h in the $\bar{\epsilon}$ -covering set such that

$$\|W^h - \hat{W}^h\|_\infty \lesssim H^{-1} \exp(-C_3(S_a^2 B_a^2 B_x B_y)) \bar{\epsilon}, \quad \|W^h - \hat{W}^h\|_0 \leq 2S_a$$

where C_3 is a sufficiently large constant. Similar inequalities hold true for Q^h, K^h, V^h , and A . Let $\hat{\theta} = (\hat{A}, (\hat{W}^h, \hat{Q}^h, \hat{K}^h, \hat{V}^h)_h)$. Then, for all $\mu \in \mathcal{P}([-B_y, B_y]^{d_{\text{attn}}})$ and $x \in [-B_x, B_x]^{d_{\text{attn}}}$,

$$\|\text{Attn}_\theta(\mu, x) - \text{Attn}_{\hat{\theta}}(\mu, x)\|_\infty$$

$$\begin{aligned}
&\leq \sum_h \left\| W^h \int \frac{V^h y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu - \hat{W}^h \int \frac{\hat{V}^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle)}{\int \exp(\langle \hat{Q}^h x, \hat{K}^h z \rangle) d\mu(z)} d\mu \right\|_\infty \\
&\quad + \|(A - \hat{A})x\|_\infty \\
&\lesssim \underbrace{\sum_h \left\| W^h \int \frac{V^h y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu - \hat{W}^h \int \frac{V^h y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu \right\|_\infty}_{(i)} \\
&\quad + \underbrace{\sum_h S_a B_a \left\| \int \frac{V^h y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu - \int \frac{V^h y \exp(\langle \hat{Q}^h x, K^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu \right\|_\infty}_{(ii)} \\
&\quad + \underbrace{\sum_h S_a B_a \left\| \int \frac{V^h y \exp(\langle \hat{Q}^h x, K^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu - \int \frac{V^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu \right\|_\infty}_{(iii)} \\
&\quad + \underbrace{\sum_h S_a B_a \left\| \int \frac{V^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu - \int \frac{V^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle)}{\int \exp(\langle \hat{Q}^h x, K^h z \rangle) d\mu(z)} d\mu \right\|_\infty}_{(iv)} \\
&\quad + \underbrace{\sum_h S_a B_a \left\| \int \frac{V^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle)}{\int \exp(\langle \hat{Q}^h x, K^h z \rangle) d\mu(z)} d\mu - \int \frac{V^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle)}{\int \exp(\langle \hat{Q}^h x, \hat{K}^h z \rangle) d\mu(z)} d\mu \right\|_\infty}_{(v)} \\
&\quad + \underbrace{\sum_h S_a B_a \left\| \int \frac{V^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle)}{\int \exp(\langle \hat{Q}^h x, \hat{K}^h z \rangle) d\mu(z)} d\mu - \int \frac{\hat{V}^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle)}{\int \exp(\langle \hat{Q}^h x, \hat{K}^h z \rangle) d\mu(z)} d\mu \right\|_\infty}_{(vi)} \\
&\quad + \underbrace{\|(A - \hat{A})x\|_\infty}_{(vii)} \\
&\lesssim \bar{\epsilon}.
\end{aligned}$$

Each term is bounded as follows:

(i). Let $\mathbf{w} := \int \frac{V^h y \exp(\langle Q^h x, K^h y \rangle)}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} d\mu$. Then, (i) $\leq \|(W^h - \hat{W}^h)\mathbf{w}\|_\infty \lesssim (2S_a) \cdot H^{-1} \exp(-C_3(S_a^2 B_a^2 B_x B_y)) \bar{\epsilon} \cdot \exp(O(S_a^2 B_a^2 B_x B_y)) \lesssim H^{-1} \bar{\epsilon}$ by (P-i,ii);

(ii). The second term is bounded by

$$\begin{aligned}
(ii) &\lesssim S_a B_a \left\| \sup_{\mathbf{y} \in \text{supp}(\mu)} \left| \frac{V^h y_i}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} \right. \right. \\
&\quad \times \left. \left. \left(\exp(\langle Q^h x, K^h y \rangle) - \exp(\langle \hat{Q}^h x, K^h y \rangle) \right) \right\|_\infty \\
&\lesssim S_a B_a \exp(O(S_a^2 B_a^2 B_x B_y)) \sup_{\mathbf{y} \in \text{supp}(\mu)} \left| \langle Q^h x, K^h y \rangle - \langle \hat{Q}^h x, K^h y \rangle \right| \\
&\lesssim S_a B_a \exp(O(S_a^2 B_a^2 B_x B_y)) \|(Q^h - \hat{Q}^h)x\|_\infty \sup_{\mathbf{y} \in \text{supp}(\mu)} \|K^h y\|_1 \\
&\lesssim \exp(O(S_a^2 B_a^2 B_x B_y)) \cdot H^{-1} \exp(-C_3(S_a^2 B_a^2 B_x B_y)) \bar{\epsilon}
\end{aligned}$$

$$\lesssim H^{-1}\bar{\epsilon}.$$

Note that the second inequality is derived by (P-ii,P-iv) and the fourth inequality is supported by (P-i).

(iii). The third term can be bounded in the same way as (ii).

(iv). The fourth term is bounded by

$$\begin{aligned} (iv) &\lesssim S_a B_a \left\| \int V^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle) d\mu \right\|_\infty \\ &\quad \times \left| \frac{1}{\int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)} - \frac{1}{\int \exp(\langle \hat{Q}^h x, K^h z \rangle) d\mu(z)} \right| \\ &\lesssim S_a B_a \left\| \int V^h y \exp(\langle \hat{Q}^h x, \hat{K}^h y \rangle) d\mu \right\|_\infty \\ &\quad \times \left(\min \left\{ \int \exp(\langle Q^h x, K^h z \rangle) d\mu(z)^{-2}, \int \exp(\langle \hat{Q}^h x, K^h z \rangle) d\mu(z) \right\} \right)^{-2} \\ &\quad \times \left| \int \exp(\langle \hat{Q}^h x, K^h z \rangle) - \exp(\langle Q^h x, K^h z \rangle) d\mu(z) \right| \quad (\text{by (P-iii)}) \\ &\lesssim \exp(O(S_a^2 B_a^2 B_x B_y)) \cdot \sup_{z \in \text{supp}(\mu)} \left| \langle Q^h x, K^h z \rangle - \langle \hat{Q}^h x, K^h z \rangle \right| \\ &\lesssim H^{-1}\bar{\epsilon}; \end{aligned}$$

(v). The fifth term is bounded in the similar way as (iv).

(vi). The sixth term is bounded in the same way as (i).

(vii). It is easily bounded by $H^{-1}\bar{\epsilon}$ using (P-i).

Please note that (P-i) $\|Ax\|_\infty, \|Ax\|_1 \leq sb\|x\|_\infty$ where the number of non-zero entries in a matrix A is bounded by s , and the absolute value of each entry is bounded by b , (P-ii) $|\langle Q^h x, K^h y \rangle| \leq S_a B_a B_x \cdot S_a B_a B_y$ because each coordinate of $Q^h x$ is a sum of at most S_a terms, each bounded by $B_a B_x$, and similarly each coordinate of $K^h y$ is bounded by $S_a B_a B_y$, (P-iii) $1/\alpha_1 - 1/\alpha_2 = (\alpha_1 - \alpha_2)/(\alpha_1 \alpha_2) \leq A^{-2}|\alpha_1 - \alpha_2|$ when $A < \alpha_1, \alpha_2$, and (P-iv) $y \mapsto \exp(y)$ and $y \mapsto y \exp(y)$ are $\exp(O(B))$ -Lipschitz over $[-B, B]$.

By the construction rule of the covering set, the covering number is bounded by

$$\begin{aligned} &\mathcal{N}(\mathcal{A}(d_{\text{attn}}, H, B_a, S_a); \bar{\epsilon})_\infty \\ &\lesssim \left(\binom{d_{\text{attn}}^2}{S_a} \bar{\epsilon}^{-S_a} \right)^{4H+1} \\ &\lesssim \left(\binom{d_{\text{attn}}^2}{S_a} \cdot \exp(O(S_a \log(H) + S_a^3 B_a^2 B_x B_y)) (\bar{\epsilon}^{-1} + 1)^{S_a} \right)^{4H+1}. \end{aligned}$$

□

Together, these results give a complete characterization of the first attention layer in the measure-theoretic setting: it can accurately realize ϕ_2 , does so in a stable manner, and has a covering number that scales favorably with D and ϵ_2 .

B.5 STEP 2-3: SECOND MLP LAYER

Having established in the previous subsections that the first MLP layer can approximate the Mercer basis functions e_j and that the attention mechanism can extract the corresponding coefficients

$\int e_j d\mu^{(i^*)}$ associated with the relevant component measure, we now turn to the next stage of the architecture.

In this step, the inputs to the model are effectively reduced to the finite collection of Mercer coefficients (b_1, \dots, b_D) together with the query vector x . The statistical problem is thus transformed into the approximation of a Lipschitz function defined over a $(D + O(1))$ -dimensional domain.

Our goal in this section is twofold: first, to determine the appropriate truncation dimension D that balances approximation error against complexity, and second, to establish approximation results for Lipschitz functions of D variables using neural networks.

B.5.1 DETERMINING THE DIMENSION D

As discussed above, after the first MLP and attention layers, the effective representation of the input measure μ_0 is reduced to its Mercer coefficients with respect to the kernel eigenbasis $\{e_i\}_{i \geq 1}$. In practice, however, only a finite number of coefficients can be retained. Thus, a key question is: how many terms D should be kept in the truncated expansion so that the approximation error remains negligible while the statistical complexity of the model is controlled? The following lemma quantifies the truncation error when approximating μ_0 by its projection onto the first D eigenfunctions.

Lemma 11. *Let $\mu_0 \in \mathcal{H}_0$ with Mercer expansion*

$$\frac{d\mu_0}{d\lambda} = \sum_{i=1}^{\infty} b_i e_i,$$

where $\{e_i\}$ are the Mercer eigenfunctions associated with kernel eigenvalues $\{\lambda_i\}$ and λ is the Lebesgue measure. Define the truncated approximation

$$\tilde{\mu}_0 = \sum_{i=1}^D b_i e_i.$$

If $\gamma_f < 0$ and $\gamma_b > 0$, then the truncation error in the γ_f -norm is bounded as

$$\|\mu_0 - \tilde{\mu}_0\|_{\gamma_f} \leq \lambda_{D+1}^{\frac{-\gamma_f + \gamma_b}{2}}.$$

Proof. The LHS is bounded by

$$\begin{aligned} & \|\mu_0 - \tilde{\mu}_0\|_{\gamma_f} \\ &= \sqrt{\sum_{i \geq D+1} \lambda_i^{-\gamma_f} b_i^2} \\ &= \sqrt{\lambda_{D+1}^{-\gamma_f + \gamma_b}} \sqrt{\sum_{i \geq D+1} \frac{\lambda_i^{-\gamma_f}}{\lambda_{D+1}^{-\gamma_f + \gamma_b}} b_i^2} \\ &\leq \sqrt{\lambda_{D+1}^{-\gamma_f + \gamma_b}} \sqrt{\sum_{i \geq D+1} \lambda_i^{-\gamma_b} b_i^2} \quad (\text{by } \lambda_{D+1} \geq \lambda_i, i \geq D+1 \text{ and } -\gamma_f > 0) \\ &\leq \sqrt{\lambda_{D+1}^{-\gamma_f + \gamma_b}} \quad (\text{by } \sqrt{\sum_i \lambda_i^{-\gamma_b} b_i^2} \leq 1) \end{aligned}$$

where we used $\frac{\lambda_i^{-\gamma_f}}{\lambda_{D+1}^{-\gamma_f}} \leq 1$ for $i \geq D+1$ in the first inequality, $\lambda_i^{-\gamma_b} \geq \lambda_{D+1}^{-\gamma_b}$ for $i \geq D+1$ and that μ_0 is in the ball in the last inequality. \square

Having controlled the truncation error of the Mercer expansion, we next turn to the regularity of the target regression function F^* . In particular, F^* is assumed to be Lipschitz continuous with respect to the product metric consisting of the γ_f -weighted RKHS distance on measures and the standard Euclidean distance on the query variable. Formally, there exists $L > 0$ such that

$$|F^*(\mu, x) - F^*(\nu, y)| \leq L \left(\|\mu - \nu\|_{\mathcal{H}_0^{\gamma_f}} + \|x - y\|_2 \right), \quad \forall \mu, \nu \in B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}), x, y \in \mathbb{R}^{d_1}.$$

This Lipschitz property ensures that once the infinite-dimensional measure μ is replaced by its truncated D -dimensional approximation, the induced error on F^* can be directly bounded. The following corollary makes this reduction explicit.

Corollary 2. *Define the truncated regression function*

$$\bar{F}_D(\mu_0, v) := F^* \left(\sum_{i=1}^D b_i e_i, x \right), \quad \frac{d\mu_0}{d\lambda} = \sum_{i=1}^{\infty} b_i e_i, \quad x = \begin{bmatrix} v \\ 0 \end{bmatrix}.$$

If $b_i = 0$ for $i \geq D + 1$, we simply write $\bar{F}_D(b_1, \dots, b_D, v)$. Then,

$$\sup_{\mu_0 \in B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0}^{\gamma_b})} |\bar{F}_D(\mu_0, v) - F^*(\mu_0, v)| \lesssim \lambda_{D+1}^{\frac{-\gamma_f + \gamma_b}{2}}.$$

Moreover, \bar{F}_D is Lipschitz with respect to the coefficients (b_1, \dots, b_D) and query v , satisfying

$$\begin{aligned} |\bar{F}_D(b, v) - \bar{F}_D(c, v')| &\leq L \sqrt{\max_{i \geq 1} \lambda_i^{-\gamma_f}} \|b - c\|_2 + L \|v - v'\|_2 \\ &\lesssim O(1) \cdot \left\| \begin{bmatrix} b - c \\ v - v' \end{bmatrix} \right\|_2 \end{aligned}$$

where $b = (b_1, \dots, b_D)$ and $c = (c_1, \dots, c_D)$. Note that $\gamma_f < 0$, so the multiplicative factor in front of $\|b - c\|_2$ is $\Theta(1)$.

B.5.2 APPROXIMATING FINITE-DIMENSIONAL LIPSCHITZ FUNCTIONS

Once the Mercer expansion has been truncated to $O(d_1 + D)$ coefficients, the infinite-dimensional regression problem reduces to approximating a Lipschitz function $f : [0, 1]^{O(d_1 + D)} \rightarrow \mathbb{R}$ with Lipschitz constant K . We now recall quantitative results on the approximation of such functions by deep ReLU networks.

Lemma 12 (Schmidt-Hieber (2020)). *For any function $f \in \text{Lip}_L([0, 1]^{\bar{D}})$ and any integers $m \geq 1$ and $N \geq \exp(\Omega(\bar{D}))$. There exists a network*

$$\tilde{f} \in \mathcal{F}(\ell, (\bar{D}, 6(D+1)N, \dots, 6(D+1)N, 1), s, \infty)$$

with depth

$$\ell \simeq (m+1)(1 + \log(\bar{D} + 1))$$

and number of parameters

$$s \lesssim (\bar{D} + 1)^{3+\bar{D}} N(m+6),$$

such that

$$\|\tilde{f} - f\|_{L^\infty} \lesssim (L+1)(1 + \bar{D}^2) 6^{\bar{D}} N 2^{-m} + LN^{-1/\bar{D}}.$$

This lemma shows that deep ReLU networks can approximate any Lipschitz function on $[0, 1]^{O(D)}$ with an explicit trade-off between network depth, width, and approximation error. The next remark connects this general result to our Mercer-RKHS setting.

Lemma 13 (Schmidt-Hieber (2020)). *Let $V = \prod_{i=0}^{\ell} (p_i + 1)$. Then,*

$$\log \mathcal{N}(\mathcal{F}(\ell, p, s, F); \delta)_\infty \lesssim (s+1) \log(\delta^{-1} \ell V).$$

This bound shows that the covering entropy grows at most logarithmically with the resolution δ^{-1} , once the architecture parameters (ℓ, p, s) are fixed. Applying our parameter selection yields the following implication.

Finally, for later use, we recall a useful estimate on the Lipschitz constant of a ReLU network in terms of its layer widths.

Lemma 14 (From the proof of lemma 5 in Schmidt-Hieber (2020)). *The Lipschitz constant of NN (w.r.t. infinity norm) is $\prod_{i=0}^{\ell} p_i$.*

From the above lemmas, we have a specialized approximation results for our Mercer-RKHS setting:

Corollary 3 (Specialization to Our Setting). *Under Setting 3 and assume $d_1 \simeq (\ln \epsilon^{-1})^{\beta^{-1}}$ so that*

$$\bar{D} \sim \left(\frac{-\gamma_f + \gamma_b}{2c}\right)^{-1/\alpha} (\ln \epsilon^{-1})^{1/\alpha} + d_1 \sim (\ln \epsilon^{-1})^{1/\min(\alpha, \beta)}.$$

Letting

$$m = O(\bar{D} \cdot \log \epsilon_3^{-1}), \quad N = \epsilon_3^{-\bar{D}}$$

in Lemma 12, there exists a ReLU network with depth

$$\ell_2 \lesssim (\text{poly log}(d + D) \cdot (d + D) \cdot \log \epsilon_3^{-1}),$$

width

$$\|\mathbf{p}_2\|_\infty \lesssim (d + D) \epsilon_3^{-d+D}$$

and the number of parameters

$$s_2 \lesssim \tilde{O}\left(\epsilon_3^{-O((\ln \epsilon^{-1})^{1/\min(\alpha, \beta)})}\right) \cdot ((\ln \epsilon)^{\min(\alpha, \beta)})^{O((\ln \epsilon)^{1/\min(\alpha, \beta)})} \cdot (\text{poly log}(d + D + \epsilon_3^{-1}))$$

that approximates \bar{F}_D , which was defined in Corollary 2, within sup-norm error $\lesssim \epsilon$.

Moreover, the covering entropy of the corresponding hypothesis class satisfies

$$\log \mathcal{N}(\mathcal{F}(\ell, p, s, F); \epsilon_3)_\infty \lesssim \epsilon_3^{-O((\ln \epsilon^{-1})^{1/\min(\alpha, \beta)})} \cdot \text{poly log}(\epsilon^{-1} \epsilon_3^{-1}),$$

and the Lipschitz constant of the network (with respect to the ℓ_∞ norm) is bounded as

$$\prod_{i=0}^{\ell} p_i \lesssim \epsilon_3^{-O(c_\epsilon (\ln \epsilon^{-1})^{\frac{2}{\min(\alpha, \beta)}} \cdot (\ln \epsilon_3^{-1}))}, \quad c_\epsilon \lesssim \text{poly log log } \epsilon^{-1}.$$

Instead of assuming $d_1 \lesssim (\log \epsilon^{-1})^{\beta^{-1}}$, if we assume that F^* is independent of x_q , then, by adding one layer $\mathbb{R}^{d+D} \ni x \mapsto x_{d+1:d+D} = -\text{ReLU}(-Ax) + \text{ReLU}(Ax) \in \mathbb{R}^D$ where $A = [O_{D \times d}; I_D]$, the ReLU network that approximates \bar{F}_D with sup-error $\lesssim \epsilon_3$ is constructed with depth

$$\ell_2 \lesssim (\text{poly log}(D) \cdot (D) \cdot \log \epsilon_3^{-1}),$$

width

$$\mathbf{p}_2 = (d + D, \underbrace{O(D\epsilon_3^{-D}), \dots, O(D\epsilon_3^{-D})}_{\ell - 1 \text{ times}})$$

and the number of parameters

$$s_2 \lesssim \tilde{O}\left(\epsilon_3^{-O((\ln \epsilon^{-1})^{1/\alpha})}\right) \cdot ((\ln \epsilon)^{1/\alpha})^{O((\ln \epsilon)^{1/\alpha})} \cdot (\text{poly log}(D + \epsilon_3^{-1}) + D).$$

The covering entropy is bounded as

$$\log \mathcal{N}(\mathcal{F}(\ell, p, s, F); \epsilon_3)_\infty \lesssim \epsilon_3^{-O((\ln \epsilon^{-1})^{1/\alpha})} \cdot \text{poly log}(\epsilon^{-1} d \epsilon_3^{-1}),$$

and the Lipschitz constant is

$$\prod_{i=0}^{\ell} p_i \lesssim d_1 \cdot \epsilon_3^{-O(c_\epsilon (\ln \epsilon^{-1})^{\frac{2}{\alpha}} \cdot (\ln \epsilon_3^{-1}))}, \quad c_\epsilon \lesssim \text{poly log log } \epsilon^{-1},$$

where $D \simeq (\ln \epsilon^{-1})^{\alpha^{-1}}$.

Remark 4. The original lemma of Schmidt-Hieber (2020) is stated for functions on $[0, 1]^{\bar{D}}$. In our setting, the domain is $[-O(1), O(1)]^{\bar{D}}$. A simple rescaling maps $[-O(1), O(1)]$ to $[0, 1]$, and this transformation only modifies the Lipschitz constant by a fixed multiplicative factor. Therefore, the approximation and covering results above remain valid up to universal constants.

This corollary consolidates the consequences of parameter selection in our setting: the effective input dimension \bar{D} grows like $(\ln \epsilon^{-1})^{1/\alpha}$, the network size scales sub-exponentially in $1/\epsilon$, the covering entropy is controlled by $\epsilon^{-O((\ln \epsilon^{-1})^{1/\alpha})}$, and the Lipschitz constant grows at most quasi-polynomially in ϵ^{-1} , when $\epsilon_3^{-1} \simeq \text{poly log } \epsilon^{-1} \cdot \epsilon^{-1}$

B.6 STEP 2-4: SECOND ATTENTION LAYER

Recall that the attention hypothesis class is parameterized as

$$\mathcal{A}(d_{\text{attn}}, H, B_a, B'_a, S_a, S'_a),$$

where d_{attn} is the embedding dimension, H is the number of heads, B_a, B'_a are bounds on the operator norms of the weight matrices, and S_a, S'_a are sparsity constraints.

In the present step, we only implement the skip connection of a scalar. We specialize to the case

$$d_{\text{attn}} = 1, \quad H = 1, \quad B_a = 0, \quad B'_a = 1, \quad S_a = 0, \quad S'_a = 1.$$

That is, the second attention layer belongs to the class

$$\mathcal{A}(1, 1, 0, 1, 0, 1).$$

This particular choice corresponds to a degenerate attention operator that is independent of the input measure and simply implements a skip connection acting as the identity on vectors, thereby ensuring consistency with the formal definition of the overall transformer class.

Lemma 15. *The ϵ_4 -covering number of $\mathcal{A}(d + D, 1, 0, 1, 0, d + D)$ satisfies*

$$\log \mathcal{N}(\mathcal{A}(1, 1, 0, 1, 0, 1); \delta)_\infty = O(\text{poly log } \epsilon_4^{-1}).$$

Proof. omitted. □

Lemma 16. *Every attention operator $\text{Attn} \in \mathcal{A}(1, 1, 0, 1, 0, 1)$ is $O(1)$ -Lipschitz with respect to the Euclidean norm.*

Proof. omitted. □

B.7 DERIVING AN ESTIMATION ERROR UPPER BOUND

We now combine the approximation bounds established in the previous subsections to derive an estimation error guarantee for transformer-type architectures. Let TF denote the hypothesis class consisting of transformer models with the architecture and parameter constraints described in Section 3.

Lemma 17 (Approximation by transformers). *In Setting 3, for all $\tilde{F}^* \in \text{Lip}_L(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^b}), \|\cdot\|_{\mathcal{H}_0^f})$, there exists $\hat{F} \in \text{TF}$ such that, for any input of the form*

$$(\nu, x_q) \quad \text{where} \quad \nu = \frac{1}{I} \sum_{i=1}^I \mu_{\nu^{(i)}}^{(i)}, \quad x_q = \text{Emb}_{\nu^{(i^*)}}(0_{d_2})$$

with $\mu_{\nu^{(i)}}^{(i)}$ generated according to Definition 1,

$$|F^*(\nu, x_q) - \hat{F}(\nu, x_q)| \lesssim \epsilon,$$

where the parameters $(d_j, H_j, B_{a,j}, B'_{a,j}, S_{a,j}, S'_{a,j}, \ell_j, \mathbf{p}_j, s_j)_{j=1}^2$ of the hypothesis set TF are defined as in Lemma 8 for $d_1, H_1, B_{a,1}, B'_{a,1}, S_{a,1}, S'_{a,1}$, Corollary 1 for $\ell_1, \mathbf{p}_1, s_1$, Lemma 15 for $d_2, H_2, B_{a,2}, B'_{a,2}, S_{a,2}, S'_{a,2}$, Corollary 3 for $\ell_2, \mathbf{p}_2, s_2$, respectively. The effective dimension D in them are determined in Corollary 2. Determination of $\epsilon_i, i = 1, \dots, 4$ are deferred to Lemma 18.

The above lemma shows that the transformer hypothesis class is sufficiently rich to approximate any Lipschitz target function \tilde{F}^* on the admissible input domain, with uniform accuracy ϵ . Please note that the output of each layer is uniformly bounded (we can add a clipping ReLU layer for each layer).

To analyze the statistical performance of ERM (empirical risk minimizer) within this class, we next require an upper bound on its covering entropy.

Lemma 18 (Covering entropy of transformers). *The covering entropy of the transformer hypothesis class satisfies*

$$\log \mathcal{N}(\text{TF}; \epsilon)_\infty \lesssim \epsilon^{-O((\ln \epsilon^{-1})^{1/\min(\alpha, \beta)})} = \exp\left(O(\ln \epsilon^{-1})^{(1+\min(\alpha, \beta))/\min(\alpha, \beta)}\right)$$

assuming that $I \leq d_1 \simeq (\ln \epsilon^{-1})^{\beta^{-1}}$ and $d_2 \simeq 1$.

Proof. The claim follows from applying the composition lemma (Lemma 5) for covering numbers. In particular,

$$\begin{aligned} & \log \mathcal{N}(\text{TF}; \epsilon)_\infty \\ & \lesssim \log \mathcal{N}(\mathcal{A}(1, 1, 0, 1, 0, 1); \epsilon)_\infty + \log \mathcal{N}\left(\mathcal{F}(\ell_2, p_2, s_2); \tilde{\Omega}(\epsilon)\right)_\infty \\ & \quad + \log \mathcal{N}\left(\{\text{Attn}_{\theta_1} \diamond \text{MLP}_{\xi_1}\}; \Omega(L_{\text{MLP}, 2}^{-1}\epsilon)\right) \\ & \lesssim \log \mathcal{N}\left(\mathcal{A}(1, 1, 0, 1, 0, 1); \underbrace{\epsilon}_{=: \epsilon_4}\right)_\infty \\ & \quad + \log \mathcal{N}\left(\mathcal{F}(\ell_2, p_2, s_2); \underbrace{\tilde{\Omega}(\epsilon)}_{=: \epsilon_3}\right)_\infty + \log \mathcal{N}\left(\mathcal{A}(d_{\text{attn}}, H, B_a, H_a); \underbrace{\tilde{\Omega}(L_{\text{MLP}, 2}^{-1}\epsilon)}_{=: \epsilon_2}\right)_\infty \\ & \quad + \log \mathcal{N}\left(\mathcal{F}(\ell_1, p_1, s_1); \underbrace{\tilde{\Omega}(L_{\text{MLP}, 2}^{-1}(L_{\text{Attn}, 1, W^1} + L_{\text{Attn}, 1, \|\cdot\|_2})^{-1}\epsilon)}_{=: \epsilon_1}\right)_\infty \\ & \lesssim \log \mathcal{N}(\mathcal{A}(d + D, 1, 0, 1, 0, d + D); \epsilon)_\infty \\ & \quad + \log \mathcal{N}\left(\mathcal{F}(\ell_2, p_2, s_2); \tilde{\Omega}(\epsilon)\right)_\infty + \log \mathcal{N}\left(\mathcal{A}(d_{\text{attn}}, H, B_a, H_a); \exp\left(-O(c_\epsilon \ln \epsilon^{-1})^{\frac{2+2\min(\alpha, \beta)}{\min(\alpha, \beta)}}\right)\right)_\infty \\ & \quad + \log \mathcal{N}\left(\mathcal{F}(\ell_1, p_1, s_1); I^{-1} \cdot \exp\left(-O(c_\epsilon \ln \epsilon^{-1})^{\frac{2+2\min(\alpha, \beta)}{\min(\alpha, \beta)}}\right)\right)_\infty, \end{aligned}$$

where $c_\epsilon \lesssim \text{poly log log } \epsilon^{-1}$. Here we used that

(i) Letting $\epsilon_4 = \epsilon$, the second attention layer is $\tilde{O}(1)$ -Lipschitz (Lemma 16);

(ii) Letting $\epsilon_3 \gtrsim \tilde{\Omega}(\epsilon)$, the Lipschitz constant of the second MLP layer is bounded as $L_{\text{MLP}, 2} \lesssim \exp(-O(c_\epsilon((\ln \epsilon^{-1})^{2/\min(\alpha, \beta)}) \cdot (\ln \epsilon_3^{-1})^2)) = \exp\left(c_\epsilon (\ln \epsilon^{-1})^{\frac{2+2\min(\alpha, \beta)}{\min(\alpha, \beta)}}\right)$ ($c_\epsilon \lesssim \text{poly log log } \epsilon^{-1}$) (Corollary 3);

(iii) Letting $\epsilon_2 \gtrsim \exp\left(-O(c_\epsilon \ln \epsilon^{-1})^{\frac{2+2\min(\alpha, \beta)}{\min(\alpha, \beta)}}\right)$, the Lipschitz constants of the first attention layer are bounded as $L_{\text{Attn}, 1, W^1} + L_{\text{Attn}, 1, \|\cdot\|_2} \lesssim D \exp(O(d^2 \log(I\epsilon_2^{-1}))) \lesssim I^{O(1)} \cdot \exp\left(O(c_\epsilon d^2 (\ln \epsilon^{-1})^{\frac{2+2\min(\alpha, \beta)}{\min(\alpha, \beta)}})\right)$ (Lemma 9);

(iv) We have $\epsilon_1 \gtrsim I^{-O(1)} \cdot \exp\left(-O(c_\epsilon d^2 (\ln \epsilon^{-1})^{\frac{2+2\min(\alpha, \beta)}{\min(\alpha, \beta)}})\right)$ for the first MLP layer.

By Lemma 15 and Corollary 3,

$$\log \mathcal{N}(\mathcal{A}(1, 1, 0, 1, 0, 1); \epsilon_4)_\infty \lesssim \text{poly log } \epsilon^{-1}$$

By Corollary 3,

$$\log \mathcal{N}(\mathcal{F}(\ell_2, p_2, s_2); \epsilon_3)_\infty \lesssim \epsilon^{-O((\ln \epsilon^{-1})^{1/\min(\alpha, \beta)})}.$$

By Lemma 10,

$$\log \mathcal{N}(\mathcal{A}(d_{\text{attn}}, H, B_a, H_a); \epsilon_2)_\infty \lesssim \text{poly}(\log \epsilon^{-1} + \log I) \cdot d^3.$$

By Corollary 1

$$\log \mathcal{N}(\mathcal{F}(\ell_1, p_1, s_1); \epsilon_1)_\infty \lesssim \text{poly log log } I \cdot \text{poly log } \epsilon^{-1} \cdot (((\log I) + \text{poly log } \epsilon^{-1} \cdot d)^{4d_2} \cdot D + d).$$

Assuming that $I \leq d_1 \simeq (\log \epsilon)^{\beta^{-1}}$ and $d_2 = O(1)$, we have

$$\begin{aligned} & \log \mathcal{N}(\mathcal{A}(d_{\text{attn}1}, H_1, B_{a,1}, H_{a,1}); \epsilon_2)_\infty + \log \mathcal{N}(\mathcal{F}(\ell_1, p_1, s_1); \epsilon_1)_\infty \\ & + \log \mathcal{N}(\mathcal{F}(\ell_2, p_2, s_2); \epsilon_3)_\infty + \log \mathcal{N}(\mathcal{A}(1, 1, 0, 1, 0, 1); \epsilon_4)_\infty \\ & \lesssim \epsilon^{-O((\ln \epsilon^{-1})^{1/\min(\alpha, \beta)})}. \end{aligned}$$

□

With these ingredients, we can invoke a general statistical learning bound for ERM.

Lemma 19 (Schmidt-Hieber (2020)). *Consider Gaussian regression, and let \hat{F} be the empirical risk minimizer over a hypothesis class $\mathcal{F} \subset L^2(\mathbb{P}_{\nu, x_q})$. Suppose $\|f\|_{L^\infty} \leq A$ for all $f \in \mathcal{F}$. Then, for any $\delta > 0$, if $V(\delta)$ denotes the covering entropy of \mathcal{F} , it holds that*

$$R(\bar{F}^*, \hat{F}) \lesssim \inf_{f \in \mathcal{F}} \|\hat{F} - \bar{F}^*\|_{L^2(\mathbb{P}_{\nu, x_q})}^2 + \frac{(A^2 + \sigma^2)V(\delta)}{n} + (A + \sigma)\delta.$$

We are now ready to state the statistical rate achieved by transformer ERM.

B.7.1 SUB-POLYNOMIAL CONVERGENCE RATE

Theorem 5 (Sub-polynomial convergence). *Let \hat{F} be the empirical risk minimizer whose hypothesis set is constructed in Lemma 17. In Setting 3, we have*

$$R(F^*, \hat{F}) \lesssim \exp\left(-\Omega\left((\ln n)^{\frac{\min(\alpha, \beta)}{\min(\alpha, \beta) + 1}}\right)\right)$$

assuming that the number of mixture components is bounded as $I \leq d_1 \lesssim (\ln n)^{\frac{\beta^{-1} \min(\alpha, \beta)}{\min(\alpha, \beta) + 1}}$ and $d_2 \simeq 1$.

Proof. Let $\epsilon \simeq \exp\left(-c'(\ln n)^{\frac{\min(\alpha, \beta)}{\min(\alpha, \beta) + 1}}\right)$ for sufficiently small constant $c' > 0$. Combining Lemmas 17 and 18 with Lemma 19, we obtain

$$\begin{aligned} R(F^*, \hat{F}) & \lesssim \epsilon + \frac{\exp\left(O\left((\ln \epsilon^{-1})^{1/\min(\alpha, \beta) + 1}\right)\right)}{n} \\ & \lesssim \epsilon + \frac{\exp\left(O\left(c'^{\min(\alpha, \beta)^{-1} + 1} (\ln n)\right)\right)}{n} \\ & \leq \epsilon + \frac{n^{c''}}{n} \quad \text{for some } 0 < c'' < 1, \\ & \lesssim \exp\left(-\Omega\left((\ln n)^{\frac{\min(\alpha, \beta)}{1 + \min(\alpha, \beta)}}\right)\right), \end{aligned}$$

assuming that the number of mixture components is bounded as $I \leq d_1 \lesssim (\ln \epsilon^{-1})^{\beta^{-1}}$ □

Remark 5 (Interpretation of Theorem 5). A common statistical learning bound for nonparametric regression takes the form

$$R(\hat{F}, \bar{F}^*) \lesssim n^{-\Theta(1/d)},$$

where d is the (effective) dimension of the problem. In our setting, however, the eigenvalue decay assumption $\lambda_j \simeq \exp(-cj^\alpha)$ implies that the effective dimension grows only as

$$d \sim (\ln n)^{1/(\alpha+1)}.$$

Consequently, the bound in Theorem 5 can be interpreted as a direct analogue of the classical $n^{-\Theta(1/d)}$ rate, but with d replaced by $(\ln n)^{1/(\alpha+1)}$. Importantly, this shows that the estimator bypasses the usual combinatorial difficulty of associative recall tasks. In our framework, each element to be recalled is not a finite symbol but rather a *probability measure*, i.e. an infinite-dimensional object. Despite this intrinsic complexity, the analysis reveals that the statistical behavior is governed purely by the eigenvalue decay of the underlying kernel, leading to the rate characteristic of infinite-dimensional regression.

B.7.2 BEYOND LOGARITHMIC CAPACITY

In Theorem 5, we discussed the case that the number of components (the ‘‘capacity’’ in terms of the associative memory) is bounded as

$$I \leq d_1 \lesssim (\ln n)^{\frac{\beta^{-1} \min(\alpha, \beta)}{\min(\alpha, \beta) + 1}},$$

which is logarithmic with respect to not only the sample size n , but also the number of ‘‘parameters’’, which we consider as the covering number, of our Transformer models. This is because the lipschitz functions over \mathbb{S}^{d_1-1} , not $B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$, becomes too complex when d_1 is large. On the other hand, in Appendix B.4, we observed that the number of the actual parameters attention matrix is linear in $d_1 (\geq I)$.

Here we consider the following additional assumption:

Assumption 9. The target function $\tilde{F}^*(\mu_0^{(i^*)}, x_q)$ is independent of x_q and only dependent of $\mu_0^{(i^*)}$. (i.e. we can write $\tilde{F}^*(\mu_0^{(i^*)}, x_q) = \tilde{F}^*(\mu_0^{(i^*)})$.)

Then, we have the polynomial ‘‘capacity’’ even in the associative recall task with the infinite-dimensional measure-valued components:

Lemma 20. Assume that

$$I \leq d_1 \simeq \exp\left(O((\log \epsilon^{-1})^{\frac{\alpha+1}{\alpha}})\right), \quad d_2 = O(1).$$

In Setting 3 and under the additional Assumption 9,

$$\log \mathcal{N}(\text{TF}; \epsilon)_\infty \lesssim \exp\left(O((\ln \epsilon^{-1})^{(1+\alpha)/\alpha})\right)$$

Proof. The main strategy of this lemma follows Lemma 18. We only mention the differences from the preceding lemma.

Let $\bar{D} = d + D \simeq d_1$ (consider the case $D \ll d_1$).

(i) Letting $\epsilon_4 = \epsilon$, the second attention layer is $O(1)$ -Lipschitz (Lemma 16);

(ii) Letting $\epsilon_3 \gtrsim \tilde{\Omega}(\epsilon)$, the Lipschitz constant of the second MLP layer is bounded as $L_{\text{MLP},2} \lesssim d_1 \exp(-O(c_\epsilon D^2 \cdot (\ln \epsilon^{-1})^2))$, ($c_\epsilon \lesssim \text{poly log log } \epsilon^{-1}$) (modifying Corollary 3 to ignore the first d indices);

(iii) Letting $\epsilon_2 \gtrsim d_1 \exp(-O(c_\epsilon D^2 \cdot (\ln \epsilon^{-1})^2))$, the Lipschitz constants of the first attention layer are bounded as $L_{\text{Attn},1,W^1} + L_{\text{Attn},1,\|\cdot\|_2} \lesssim D \exp(O(d^2 \log(I \epsilon_2^{-1}))) \lesssim \exp\left(O(c_\epsilon d_1^2 (\ln d_1)^{O(1)} (\ln \epsilon^{-1})^{2+2\alpha^{-1}})\right)$ (Lemma 9 and $I \leq d_1$);

(iv) We have $\epsilon_1 \gtrsim \exp\left(-O(c_\epsilon d_1^2 (\ln d_1)^{O(1)} (\ln \epsilon^{-1})^{2+2\alpha^{-1}})\right)$ for the first MLP layer.

By Lemma 15 and Corollary 3,

$$\log \mathcal{N}(\mathcal{A}(1, 1, 0, 1, 0, 1); \epsilon_4)_\infty \lesssim \text{poly log } \epsilon^{-1}$$

By modifying Corollary 3 to ignore the first d indices (corresponding to the query),

$$\log \mathcal{N}(\mathcal{F}(\ell_2, p_2, s_2); \epsilon_3)_\infty \lesssim \epsilon^{-O((\ln \epsilon^{-1})^{1/\alpha})}.$$

By Lemma 10,

$$\log \mathcal{N}(\mathcal{A}(d_{\text{attn}1}, H_1, B_{a,1}, H_{a,1}); \epsilon_2)_\infty \lesssim \text{poly}(\log(\epsilon^{-1} d_1)) \cdot d^3.$$

By Corollary 1

$$\log \mathcal{N}(\mathcal{F}(\ell_1, p_1, s_1); \epsilon_1)_\infty \lesssim \text{poly log log } I \cdot \text{poly log } \epsilon^{-1} \cdot (((\log I) + \text{poly log } \epsilon^{-1} \cdot d)^{4d_2} \cdot D + d).$$

Assuming that $I \leq d_1 \simeq \exp\left(O((\log \epsilon)^{\frac{\alpha+1}{\alpha}})\right)$ and $d_2 = O(1)$, we have

$$\log \mathcal{N}(\mathcal{A}(d_{\text{attn}1}, H_1, B_{a,1}, H_{a,1}); \epsilon_2)_\infty + \log \mathcal{N}(\mathcal{F}(\ell_1, p_1, s_1); \epsilon_1)_\infty$$

$$\begin{aligned} & + \log \mathcal{N}(\mathcal{F}(\ell_2, p_2, s_2); \epsilon_3)_\infty + \log \mathcal{N}(\mathcal{A}(1, 1, 0, 1, 0, 1); \epsilon_4)_\infty \\ & \lesssim \epsilon^{-O((\ln \epsilon^{-1})^{1/\alpha})}. \end{aligned}$$

□

In the same vein, we obtain the similar result as in Theorem 5:

Theorem 6. *Let \hat{F} be the empirical risk minimizer whose hypothesis set is constructed in Lemma 17. Assume that*

$$I \leq d_1 \simeq \exp(o(\log n)) = n^{o(1)}$$

and the target function $F^*(\mu_0^{(i^*)}, x_q)$ is independent of x_q (Assumption 9). In Setting 3, we have

$$R(F^*, \hat{F}) \lesssim \exp(-\Omega((\ln n)^{\frac{\alpha}{\alpha+1}})).$$

Proof. The proof is the same as in Theorem 5. □

C MINIMAX LOWER BOUND

C.1 PROOF SKETCH

The goal of this section is to establish an information-theoretic minimax lower bound for the associative recall problem in Setting 4. Our proof strategy consists of two main steps: a reduction to a pure infinite-dimensional regression task, and the derivation of covering/packing entropy bounds for the corresponding Lipschitz functionals.

Step 1: Reduction from Infinite-Dimensional Regression (Appendix C.2). We first reduce a regression problem on measures to the associative recall problem. Let

$$\mathcal{F}^* = \text{Lip}_L(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}) \times \mathcal{X}_q, d_{\text{prod}}), \quad \tilde{\mathcal{F}}^* = \text{Lip}_L(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}), \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}),$$

where \mathcal{F}^* is the full class of Lipschitz functions depending on both (μ, x) , and $\tilde{\mathcal{F}}^* \subset \mathcal{F}^*$ is the subclass depending only on μ .

Here, the key observation is that each ν can be written as an average of pushforward measures

$$\nu = \frac{1}{I} \sum_{i=1}^I \mu_{v^{(i)}}^{(i)} = \frac{1}{I} \sum_{i=1}^I (\text{Emb}_{v^{(i)}})_\# \mu_0^{(i)}.$$

Therefore, estimating $F(\mu_0^{(i^*)}, x)$ with the “noisy” input ν is at least as hard as estimating with the “pure” input $\mu_0^{(i^*)}$.

Formally, the two observation models differ as follows:

$$\begin{aligned} \mathcal{S}_n &= \{(\nu_t, (x_q)_t, y_t)\}_{t=1}^n, \quad y_t = \tilde{F}^*(\mu_0^{(i^*)}, (x_q)_t) + \xi_t, \quad \tilde{F}^* \in \mathcal{F}^*, \\ \mathcal{U}_n &= \{(\mu_0^{(i^*)}, y_t)\}_{t=1}^n, \quad y_t = \tilde{F}^*(\mu_0^{(i^*)}) + \xi_t, \quad \tilde{F}^* \in \tilde{\mathcal{F}}^* \subset \mathcal{F}^*. \end{aligned}$$

That is, under \mathcal{S}_n we observe outputs of a general Lipschitz function $\tilde{F}^* \in \mathcal{F}^*$, whereas under \mathcal{U}_n the outputs are restricted to the subclass $\tilde{\mathcal{F}}^*$.

Consequently, remembering that $F^*(\nu, x) := \tilde{F}^*(\mu_0^{(i^*)}, x)$, the minimax risk satisfies

$$\inf_{\hat{F}} \sup_{\tilde{F}^* \in \mathcal{F}^*} \mathbb{E}_{\mathcal{S}_n} [\|\hat{F}(\nu, x) - F^*(\nu, x)\|^2] \geq \inf_{\hat{F}} \sup_{\tilde{F}^* \in \tilde{\mathcal{F}}^*} \mathbb{E}_{\mathcal{U}_n} [\|\hat{F}(\mu_0) - \tilde{F}^*(\mu_0)\|^2].$$

In words: the associative recall problem with dataset \mathcal{S}_n and hypothesis class \mathcal{F}^* is at least as hard as the reduced regression problem with dataset \mathcal{U}_n and restricted class $\tilde{\mathcal{F}}^*$. This reduction allows us to focus on an *infinite-dimensional regression* setting.

Step 2: Entropy Bounds for Lipschitz Functionals. The minimax lower bound is based on the general Gaussian regression minimax bound in Yang & Barron (1999): in short,

$$\log \mathcal{M}(\mathcal{F}^*; \epsilon)_{L^2(\mathbb{P}_{\mu_0})} \simeq n\epsilon^2 \Rightarrow \inf_{\hat{F}} \sup_{F^*} R(F^*, \hat{F}) \gtrsim \epsilon^2,$$

where $\mathcal{M}(\mathcal{F}^*; \epsilon)_{L^2(\mathbb{P}_{\mu_0})}$ denotes the ϵ -packing number of the function class \mathcal{F}^* with respect to the $L^2(\mathbb{P}_{\mu_0})$ metric. Thus, to obtain a lower bound it suffices to evaluate the packing entropy (i.e., the metric entropy $\log \mathcal{M} (\simeq \log \mathcal{N})$) of the set of Lipschitz functionals G^* under $L^2(\mathbb{P}_{\mu_0})$. To apply the classical information-theoretic results, we derive both upper and lower bounds for the covering/packing entropy of the relevant Lipschitz functional class.

STEP 2-1: UPPER BOUND (APPENDIX C.3). It is known (see Boissard (2011)) that for a Lipschitz class,

$$\log \mathcal{N}(\text{Lip}(A, d); \epsilon)_{L^\infty} \lesssim \mathcal{N}(A; \epsilon)_d \cdot (\text{poly log } \epsilon^{-1}),$$

where A is the input set and d is the underlying metric. Applying this principle, we show that

$$\log \mathcal{N}(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}); \epsilon)_{\|\cdot\|_{\mathcal{H}_0^{\gamma_f}}} \lesssim \exp(c_u \log \epsilon^{-1})^{\frac{\alpha+1}{\alpha}}, \quad c_u < \infty.$$

The proof relies on the following isometric transformation: for $f = \sum_{i=1}^{\infty} b_i e_i$ and parameters a, b, c ,

$$f \mapsto \phi_{a,b,c}(f) := \sum_{i=1}^{\infty} \lambda_i^{\frac{c-b}{2}} b_i e_i,$$

which is an isometric bijection from $(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^a}), \|\cdot\|_{\mathcal{H}_0^b})$ to $(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{a-b+c}}), \|\cdot\|_{\mathcal{H}_0^c})$. We apply this with $a = \gamma_b$, $b = \gamma_f$, and $c = 0$. Then, we employ a standard argument of covering an infinite-dimensional ellipsoid endowed with ℓ^2 metric.

STEP 2-2: LOWER BOUND (APPENDIX C.4). The key difficulty is that the Lipschitz constant is *anisotropic*: differences along low-index directions (small j) are heavily penalized, while directions with larger j are effectively much smoother. Formally, for F in our class one has

$$|F(\sum_j b_j e_j) - F(\sum_j c_j e_j)| \leq L \sqrt{\sum_j \lambda_j^{-\gamma_f} (b_j - c_j)^2}, \quad \gamma_f < 0,$$

which clearly shows that directions with larger eigenvalues λ_j (small j) contribute far more to the Lipschitz bound than those with smaller eigenvalues (large j). In other words, the geometry of the function class is highly distorted across coordinates.

To make this structure explicit, we construct a rescaling map that embeds the standard cube $[0, 1]^d$ into our measure-input space. After rescaling each coordinate according to the eigenvalue decay $\{\lambda_j\}$, a Lipschitz function on $[0, 1]^d$ becomes a function on $B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$ with Lipschitz constant proportional to

$$\lambda_d^{(-\gamma_f + \gamma_a)/2}.$$

This shows that our class contains an embedded copy of the d -dimensional Lipschitz ball, up to a rescaling factor.

Consequently, the packing entropy of our class is at least as large as that of $\text{Lip}_1([0, 1]^d, \|\cdot\|_{\ell^\infty})$ at resolution $R\epsilon/\lambda_d^{(-\gamma_f + \gamma_a)/2}$. By combining this rescaling argument with known packing lower bounds for Lipschitz functions on $[0, 1]^d$ and the Yang–Barron information-theoretic inequality, we obtain the desired minimax lower bound for the associative recall problem.

Remark 6. In Setting 4, we assume that the density $\frac{d\mu_0}{d\lambda}$ is *nonnegative*, which can always be ensured by adding a sufficiently large constant shift. We then relax the additional constraint that μ_0 must be a probability measure, and instead only require that $\frac{d\mu_0}{d\lambda}$ belongs to a bounded ball in the ambient function space. This modification does not affect the minimax difficulty of the problem, since the essential hardness arises from the *infinite-dimensionality* of the domain, whereas the normalization constraint corresponds merely to a *finite-dimensional* restriction.

C.2 STEP 1: REDUCTION FROM INFINITE-DIMENSIONAL REGRESSION

By Lemma 21, we will show that the associative recall problem is at least as hard as a Gaussian regression problem where the input variables are measures. Remember that a standard Gaussian regression problem is defined as follows: we observe i.i.d. random variables $(X_t, Y_t)_{t=1}^n$ such that

$$Y_t = F^*(X_t) + \xi_t, \quad t = 1, \dots, n,$$

where ξ_t are i.i.d. Gaussian noise, independent of X_t . On the other hand, in our recall-and-predict problem, the observed input ν is noisy: $(\mu_{v^{(i)}}^{(i)})_{i \neq i^*}$ are mixed in ν , but they are irrelevant to the output y . We will show that we can obtain a better estimator when we “eliminate” the noises in the inputs. Formally, we obtain the following corollary:

Corollary 4. *Let*

$$\mathcal{F}^* := \text{Lip}_L(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}) \times \mathcal{X}_q, d_{\text{prod}}), \quad \tilde{\mathcal{F}}^* := \text{Lip}_L(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}), \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}).$$

Here, \mathcal{F}^* denotes the class of L -Lipschitz functions in both (μ, x) , while $\tilde{\mathcal{F}}^*$ denotes the subclass of L -Lipschitz functions depending only on μ . Note that $\tilde{\mathcal{F}}^* \subset \mathcal{F}^*$.

Consider datasets

$$\begin{aligned} \mathcal{S}_n &= \{(\nu_t, (x_q)_t, y_t)\}_{t=1}^n, \quad y_t = \tilde{F}^*(\mu_0^{(i^*)}, (x_q)_t) + \xi_t, \quad \tilde{F}^* \in \mathcal{F}^*, \\ \mathcal{U}_n &= \{(\mu_0^{(i^*)}, y_t)\}_{t=1}^n, \quad y_t = \tilde{F}^*(\mu_0^{(i^*)}) + \xi_t, \quad \tilde{F}^* \in \tilde{\mathcal{F}}^* \subset \mathcal{F}^*. \end{aligned}$$

sampled as in Setting 4. Then, remembering that $F^*(\nu, x_q) = \tilde{F}^*(\mu_0^{(i^*)}, x_q)$, we have

$$\begin{aligned} & \inf_{\mathcal{S}_n \mapsto \hat{F} \in L^2(\mathbb{P}_{\nu, x_q})} \sup_{\tilde{F}^* \in \mathcal{F}^*} \mathbb{E}_{\mathcal{S}_n} [\|\hat{F}(\nu, x) - F^*(\nu, x)\|_{L^2(\mathbb{P}_{\nu, x_q})}^2] \\ & \geq \inf_{\mathcal{U}_n \mapsto \hat{F} \in L^2(\mathbb{P}_{\mu_0})} \sup_{\tilde{F}^* \in \tilde{\mathcal{F}}^*} \mathbb{E}_{\mathcal{U}_n} [\|\hat{F}(\mu_0) - \tilde{F}^*(\mu_0)\|_{L^2(\mathbb{P}_{\mu_0})}^2]. \end{aligned}$$

In words: the estimation problem with query-dependent target functions is at least as hard as the restricted problem where the target depends only on μ . Hence, establishing a lower bound for the latter suffices.

To prove Corollary 4, we state the following lemma.

Lemma 21. *Let $\mathcal{S}_n = \{(\nu_t, x_t, y_t)\}_{t=1}^n$ and $\bar{\mathcal{S}}_n = \{(\mu_0^{(i^*)}, x_t, y_t)\}_{t=1}^n$ be datasets sampled as in Definition 1. Then, for any estimator $\hat{F} : \mathcal{S}_n \mapsto \hat{F}$, there exists an estimator $\tilde{F}_1 : \bar{\mathcal{S}}_n \mapsto \tilde{F}_1$ such that*

$$\mathbb{E}_{\mathcal{S}_n} [\|\hat{F}(\nu, x) - F^*(\nu, x)\|_{L^2(\mathbb{P}_{\nu, x_q})}^2] \geq \mathbb{E}_{\bar{\mathcal{S}}_n} [\|\tilde{F}_1(\mu_0, x) - \tilde{F}^*(\mu_0, x)\|_{L^2(\mathbb{P}_{\mu_0, x_q})}^2].$$

Moreover, if \tilde{F}^* is independent of the query x , i.e. $\tilde{F}^*(\mu, x) = \tilde{F}_2^*(\mu)$, then there exists an estimator \tilde{F}_2 depending only on $\{(\mu_0^{(i^*)}, y_t)\}_{t=1}^n$ such that

$$\mathbb{E}_{\mathcal{S}_n} [\|\hat{F}(\nu, x) - \tilde{F}_2^*(\mu_0)\|_{L^2}^2] \geq \mathbb{E}_{\{(\mu_0^{(i^*)}, y_t)\}_t} [\|\tilde{F}_2(\mu_0) - \tilde{F}_2^*(\mu_0)\|_{L^2}^2].$$

Proof. Remember that

$$\nu = \frac{1}{I} \sum_i \mu_{v^{(i)}}^{(i)} = \frac{1}{I} \sum_i (\text{Emb}_{v^{(i)}})_{\#} \mu_0^{(i)}, \quad \mu_0^{(i)} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{\mu_0}, \quad (v^{(i)})_{i=1}^I \sim \mathbb{P}_v.$$

We want to eliminate the dependency of $\mu_0^{(i)}$ and $v^{(i)}$ for $i \in [1 : I] \setminus \{i^*\}$ from the original estimator and make a better estimate. We will construct a Bayes-estimated mapping

$$\bar{\mathcal{S}}_n \mapsto \{(\mu_0^{(i^*)}, x_q) \mapsto \mathbb{E}_{\nu, \mathcal{T}_n} [\hat{F}_{\mathcal{S}_n}(\nu, x_q) \mid \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] =: \tilde{F}_1(\mu_0^{(i^*)}, x_q)\} \in L^2(\mathbb{P}_{\mu_0, x_q}),$$

where $\mathcal{T}_n := \{(\mu_0^{(i)})_t, (v^{(i)})_t \mid \text{for } \forall i \in [1 : I] \setminus \{i^*\} \text{ and } \forall t \in [1, n]\}$ and $\bar{\mathcal{S}}_n := ((\mu_0^{(i^*)})_t, (x_q)_t, y_t)_{t=1}^n$. This estimator $\bar{\mathcal{S}}_n \mapsto \tilde{F}_1$ is well-defined as the mapping from $\bar{\mathcal{S}}_n$ to a

function in $L^2(\mathbb{P}_{\mu_0, x_q})$ because (i) we can deterministically construct \mathcal{S}_n from $\bar{\mathcal{S}}_n$ and \mathcal{T}_n because $\nu = \frac{1}{I} \sum_i \mu_{v^{(i)}}^{(i)} = \frac{1}{I} \sum_i (\text{Emb}_{v^{(i)}})_\# \mu_0^{(i)}$ and (ii) ν only has the randomness of the noises $\{(\mu_0^{(i)}, v^{(i)})\}_{i \neq i^*}$, given $(\mu_0^{(i^*)}, x_q)$. Note that the above estimator does not use the oracle of sampling an input/output pair. In short, \hat{F}_1 is not cheating in the context of the standard Gaussian regression: To take the expectation with respect to $\nu | \mu_0^{(i^*)}$ and \mathcal{T}_n , we are additionally observing only the noises of the input ν , which do not exist in the standard Gaussian regression. From now, we will explicitly write $\hat{F} = \hat{F}_{\mathcal{S}_n}$. The loss is lower bounded as

$$\begin{aligned}
& \mathbb{E}_{\mathcal{S}_n} [\|\hat{F}_{\mathcal{S}_n}(\nu, x_q) - F^*(\nu, x_q)\|_{L^2(\mathbb{P}_{\nu, x_q})}^2] \\
&= \mathbb{E}_{\mathcal{S}_n} [\|(\hat{F}_{\mathcal{S}_n}(\nu, x_q) - \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n]) \\
&\quad + (\mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] - \tilde{F}^*(\mu_0^{(i^*)}, x_q))\|_{L^2(\mathbb{P}_{\nu, x_q})}^2] \\
&= \mathbb{E}_{\mathcal{S}_n} [\|\hat{F}_{\mathcal{S}_n}(\nu, x_q) - \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n]\|_{L^2(\mathbb{P}_{\nu, x_q})}^2] \quad \dots \text{(i)} \\
&\quad + 2\mathbb{E}_{\mathcal{S}_n} [\langle \hat{F}_{\mathcal{S}_n}(\nu, x_q) - \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n], \\
&\quad \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] - \tilde{F}^*(\mu_0^{(i^*)}, x_q) \rangle_{L^2(\mathbb{P}_{\nu, x_q})}] \quad \dots \text{(ii)} \\
&\quad + \mathbb{E}_{\mathcal{S}_n} [\|(\mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] - \tilde{F}^*(\mu_0^{(i^*)}, x_q))\|_{L^2(\mathbb{P}_{\nu, x_q})}^2] \\
&\geq \mathbb{E}_{\mathcal{S}_n} [\|(\tilde{F}_{1, \bar{\mathcal{S}}_n}(\mu_0, x_q) - \tilde{F}^*(\mu_0, x_q))\|_{L^2(\mathbb{P}_{\mu_0, x_q})}^2].
\end{aligned}$$

For the term (i), this is greater than zero. As for (ii), we have

$$\begin{aligned}
(ii) &= \mathbb{E}_{\mathcal{S}_n} [\langle \hat{F}_{\mathcal{S}_n}(\nu, x_q) - \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n], \\
&\quad \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] - \tilde{F}^*(\mu_0^{(i^*)}, x_q) \rangle_{L^2(\mathbb{P}_{\nu, x_q})}] \\
&= \mathbb{E}_{\bar{\mathcal{S}}_n} [\mathbb{E}_{\mathcal{T}_n} [\langle \hat{F}_{\mathcal{S}_n}(\nu, x_q) - \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n], \\
&\quad \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] - \tilde{F}^*(\mu_0^{(i^*)}, x_q) \rangle_{L^2(\mathbb{P}_{\nu, x_q})} | \bar{\mathcal{S}}_n]] \\
&= \mathbb{E}_{\bar{\mathcal{S}}_n} [\langle \mathbb{E}_{\mathcal{T}_n}[\hat{F}_{\mathcal{S}_n}(\nu, x_q) | \nu, x_q, \bar{\mathcal{S}}_n] - \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n], \\
&\quad \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] - \tilde{F}^*(\mu_0^{(i^*)}, x_q) \rangle_{L^2(\mathbb{P}_{\nu, x_q})}] \\
&\quad (\mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] \text{ and } \tilde{F}^*(\mu_0^{(i^*)}, x_q) \text{ are independent of } \mathcal{T}_n | \bar{\mathcal{S}}_n \text{ if } \bar{\mathcal{S}}_n \text{ is given}) \\
&= \mathbb{E}_{\bar{\mathcal{S}}_n} [\int \left\{ \left(\int (\mathbb{E}_{\mathcal{T}_n}[\hat{F}_{\mathcal{S}_n}(\nu, x_q) | \nu, x_q, \bar{\mathcal{S}}_n] - \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n]) d\mathbb{P}_\nu(\nu) \right) \right. \\
&\quad \left. \times (\mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] - \tilde{F}^*(\mu_0^{(i^*)}, x_q)) \right\} d\mathbb{P}_{x_q}(x_q)] \\
&= \mathbb{E}_{\bar{\mathcal{S}}_n} [\int \left\{ \left(\int (\mathbb{E}_{\mathcal{T}_n}[\hat{F}_{\mathcal{S}_n}(\nu, x_q) | \nu, x_q, \bar{\mathcal{S}}_n] - \mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n]) d\mathbb{P}_{\nu | \mu_0^{(i^*)}}(\nu) d\mathbb{P}_{\mu_0^{(i^*)}}(\mu_0^{(i^*)}) \right) \right. \\
&\quad \left. \times (\mathbb{E}_{\nu, \mathcal{T}_n}[\hat{F}_{\mathcal{S}_n} | \mu_0^{(i^*)}, x_q, \bar{\mathcal{S}}_n] - \tilde{F}^*(\mu_0^{(i^*)}, x_q)) \right\} d\mathbb{P}_{x_q}(x_q)] \\
&= 0.
\end{aligned}$$

In the same vein, we can omit the dependence of x_q if F^* is independent of x_q . We give an estimated mapping as

$$\mu_0 \mapsto \mathbb{E}_{x_q, ((x_q)_t)_{t=1}^n} [\tilde{F}_{1, \mathcal{S}_n} | \mu_0, ((\mu_0^{(i^*)})_t, y_t)_t]$$

that can be constructed only with the observation $((\mu_0^{(i^*)})_t, y_t)_t$. We omit the details for the second statement. \square

C.3 STEP 2-1: UPPER-BOUND OF THE ENTROPY

Thanks to Corollary 4, the lower-bound analysis reduces to a Gaussian regression problem in which the inputs $\mu_0^{(i^*)} \sim \mathbb{P}_{\mu_0}$ are generated according to Setting 4. The key step is to control the cov-

ring/packing numbers of the underlying function class so that we can apply the general minimax bound of Yang & Barron (1999).

Lemma 22 (Yang & Barron (1999)). *Let $\tilde{\mathcal{F}}^* := \text{Lip}_L(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}))$. Consider the dataset $\mathcal{U}_n = \{(\mu_0^{(i^*)}, y_i)\}_{i=1}^n$ generated as in Setting 4. Suppose there exist $\delta, \epsilon > 0$ such that*

$$\log_2 \mathcal{N}(\mathcal{F}^*; \epsilon)_{L^2} \leq \frac{n\epsilon^2}{2\sigma^2}, \quad \log_2 \mathcal{M}(\mathcal{F}^*; \delta)_{L^2} \geq \frac{2n\epsilon^2}{\sigma^2} + 2 \log 2.$$

Then

$$\inf_{\mathcal{U}_n \mapsto \hat{F}} \sup_{F^* \in \tilde{\mathcal{F}}^*} \mathbb{E}_{\mathcal{U}_n} \left[\|\hat{F} - F^*\|_{L^2(\mathbb{P}_{\mu_0})}^2 \right] \gtrsim \delta^2.$$

Remark 7. For lower-bound analysis, we adopt the L^2 norm when defining covering and packing numbers.

Goal of this subsection. To apply Lemma 22, we need tight control of the covering entropy of the Lipschitz function class. Specifically, we aim to establish an upper bound on

$$\log \mathcal{N} \left(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \epsilon) \right)_{L^2(\mathbb{P}_{\mu_0})}.$$

A lower bound is deferred to Appendix C.4.

By proposition B.2 in Boissard (2011), the covering entropy of a (1-)Lipschitz class can be controlled by the covering entropy of its input set:

$$\log \mathcal{N}(\text{Lip}(A, d); \epsilon)_{L^\infty} \lesssim \mathcal{N}(A; \epsilon)_d \cdot \text{polylog}(\epsilon^{-1}),$$

where A is the input domain (under some weak assumptions) and d is the underlying metric. Thus, our task reduces to bounding the entropy of the input set $A = B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$ equipped with the metric $\|\cdot\|_{\mathcal{H}_0^{\gamma_f}}$.

To this end, we establish an isometric correspondence between balls in weighted Hilbert spaces, which allows us to switch to the standard L^2 metric.

Lemma 23. *Let $a, b, c \in \mathbb{R}$. A mapping $\phi_{a,b,c}$:*

$$f = \sum_{i=1}^{\infty} b_i e_i \mapsto \phi(f) = \sum_{i=1}^{\infty} \lambda_i^{\frac{c-b}{2}} b_i e_i$$

is an isometric bijection from $(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^a}), \|\cdot\|_{\mathcal{H}_0^b})$ to $(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{a-b+c}}, \|\cdot\|_{\mathcal{H}_0^c})$.

Proof. First, for all $f \in B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^a})$,

$$\|\phi(f)\|_{\mathcal{H}_0^c} = \sum_i \lambda_i^{-c} \left(\lambda_i^{\frac{c-b}{2}} b_i \right)^2 = \sum_i \lambda_i^{-c+c-b} b_i^2 = \|f\|_{\mathcal{H}_0^b}.$$

This implies that ϕ is isometry and injective. Next, ϕ is also a surjection because

$$\begin{aligned} f &= \sum_{i=1}^{\infty} b_i e_i \in B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^a}) \\ \Leftrightarrow \sum_i \lambda_i^{-a} b_i^2 &\leq 1 \\ \Leftrightarrow \sum_i \lambda_i^{-a+b-c} \left(\lambda_i^{\frac{c-b}{2}} b_i \right)^2 &\leq 1 \\ \Leftrightarrow \phi(f) &= \sum_{i=1}^{\infty} \lambda_i^{\frac{c-b}{2}} b_i e_i \in B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{a-b+c}}) \end{aligned}$$

□

Using the above isometry and eigenvalue decay properties, we obtain the following upper bound on the entropy of the input set:

Lemma 24. *The covering entropy of $B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$ endowed with the distance $\|\cdot\|_{\mathcal{H}_0^{\gamma_f}}$ is upper bounded as*

$$\log \mathcal{N}(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}); \epsilon)_{\|\cdot\|_{\mathcal{H}_0^{\gamma_f}}} \lesssim \ln(\epsilon^{-1})^{\frac{\alpha+1}{\alpha}}.$$

Proof. By Lemma 23, letting $a = \gamma_b$, $b = \gamma_f$, $c = 0$, it is sufficient to show that

$$\log \mathcal{N}(B_\phi; \epsilon)_{L^2} \lesssim \ln(\epsilon^{-1})^{\frac{\alpha+1}{\alpha}}$$

where $B_\phi = B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b - \gamma_f}})$. We construct a J_ϵ -dimensional set $\mathcal{E}_{J_\epsilon} = \{(b_1, \dots, b_{J_\epsilon}) \mid f = \sum_i b_i e_i \in B_\phi\}$ such that, for all $f = \sum_i b_i e_i \in B_\phi$,

$$\sum_{j \geq J_\epsilon + 1} b_j^2 < \frac{1}{4} \epsilon^2.$$

This can be satisfied if $J_\epsilon \simeq (c^{-1}(\gamma_b - \gamma_f) \ln \epsilon^{-1})^{1/\alpha}$ because $\sum_{j \geq J_\epsilon + 1} b_j^2 \leq \sum_{j \geq J_\epsilon + 1} \lambda_j \lesssim \lambda_{J_\epsilon}$ with exponential decay and use $\lambda_j \simeq \exp(-cj^\alpha)$. To construct a $\frac{1}{2}\epsilon$ -covering on \mathcal{E}_{J_ϵ} , we need at most $O(1 \vee (\lambda_j^{\frac{\gamma_b - \gamma_f}{2}} / \epsilon))$ patterns for each dimension, so the covering entropy is bounded as

$$\begin{aligned} \log \mathcal{N}(\epsilon) &\lesssim \sum_{j=1}^{J_\epsilon} \log \frac{\lambda_j^{\frac{\gamma_b - \gamma_f}{2}}}{\epsilon} \\ &\lesssim \sum_{j=1}^{J_\epsilon} \left(-\left(\frac{\gamma_b - \gamma_f}{c}\right) j^\alpha + \left(\frac{\gamma_b - \gamma_f}{c}\right)^{-1} J_\epsilon^\alpha \right) \\ &\lesssim \left(\frac{\gamma_b - \gamma_f}{c}\right)^{-1} J_\epsilon^{\alpha+1} \\ &\simeq \left(\frac{\gamma_b - \gamma_f}{c}\right)^{-\alpha-1} (\ln \epsilon^{-1})^{(\alpha+1)/\alpha} \end{aligned}$$

□

Finally, results in Boissard (2011) and Lemma 24 yield the desired entropy bound for the Lipschitz class.

Lemma 25 (Based on Boissard (2011)). *The metric entropy of $\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})); \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}$ with respect to the Bochner $L^2(\mathbb{P}_{\mu_0})$ -norm satisfies*

$$\log \mathcal{N}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}), \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \epsilon)_{\|\cdot\|_{L^2(\mathbb{P}_{\mu_0})}} \lesssim \exp\left(c_u (\ln \epsilon^{-1})^{\frac{\alpha+1}{\alpha}}\right)$$

for some $c_u > 0$.

Proof. By Lemma 24, $\log \mathcal{N}(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}); \epsilon)_{\|\cdot\|_{\mathcal{H}_0^{\gamma_f}}} \lesssim \ln(\epsilon^{-1})^{\frac{\alpha+1}{\alpha}}$. Then, by Proposition B.2 in Boissard (2011) and $\epsilon^{-1} = \exp(\ln \epsilon^{-1}) = o(\exp((\ln \epsilon^{-1})^{\frac{\alpha+1}{\alpha}}))$ for any $\alpha > 0$, we have

$$\log \mathcal{N}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}), \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \epsilon)_{\|\cdot\|_{\infty}} \lesssim \exp\left(c_u (\ln \epsilon^{-1})^{\frac{\alpha+1}{\alpha}}\right).$$

By the inequality $\sqrt{\mathbb{P}_{\mu_0}(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}))} \|\cdot\|_{\infty} \geq \|\cdot\|_{L^2(\mathbb{P}_{\mu_0})}$ and $\sqrt{\mathbb{P}_{\mu_0}(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}))} = 1$, we obtain the desired bound. □

C.4 STEP 2-2: LOWER BOUND OF THE ENTROPY

To apply the general minimax bound in Lemma 22, We will provide the lower bound of the infinite-dimensional Lipschitz class.

The main difficulty lies in the anisotropic nature of the Lipschitz constant: for F in our class, one has

$$|F(\sum_j b_j e_j) - F(\sum_j c_j e_j)| \leq L \sqrt{\sum_j \lambda_j^{-\gamma_f} (b_j - c_j)^2}, \quad -\gamma_f > 0,$$

which implies that the functional is significantly smoother in directions corresponding to high-index coefficients e_j . To capture this effect, we construct the following embedding:

$$\iota_d : L^p([0, 1]^d) \hookrightarrow L^p(\mu_{\mathcal{H}_0}), \quad (\iota_d g)(f_\mu) := g\left(\Phi_1\left(\frac{f_{\mu,1}}{\lambda_1^{\gamma_d/2}}\right), \dots, \Phi_d\left(\frac{f_{\mu,d}}{\lambda_d^{\gamma_d/2}}\right)\right),$$

where

$$f_\mu = \sum_{j=1}^{\infty} \lambda_j^{\gamma_d/2} Z_j e_j, \quad d\mu(x) = f_\mu(x) dx,$$

with independent coefficients $Z_j \sim \rho_j$ and cumulative distribution functions $\Phi_j(z) = \int_{-\infty}^z \rho_j(u) du$.

We prove that

$$\iota_d(\text{Lip}_1([0, 1]^d)) \subset \text{Lip}_{R/\lambda_d^{(-\gamma_f+\gamma_d)/2}}(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}),$$

and hence obtain a packing lower bound

$$\log \mathcal{M}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \epsilon)_{L^p}) \geq \log \mathcal{M}\left(\text{Lip}_1([0, 1]^d, \|\cdot\|_{\ell^\infty}), \frac{R\epsilon}{\lambda_d^{(-\gamma_f+\gamma_d)/2}}\right)_{L^p}.$$

Combining this rescaling argument with standard Lipschitz-packing lower bounds and the information-theoretic results of Yang–Barron yields the desired minimax lower bound for the associative recall problem.

First, we construct an embedding ι_d . This construction suggests that, after a suitable rescaling of coordinates, functions on $[0, 1]^d$ can be embedded isometrically into our measure-input space. The next lemma formalizes this embedding and quantifies how the Lipschitz constant is rescaled.

Lemma 26 (An extension of Lanthaler (2024)). *Let \mathbb{P}_{μ_0} be a probability measure on $B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$ in Setting 4. For any $p \in [1, \infty)$ and $d \in \mathbb{N}$, there exists an isometric embedding*

$$\iota_d : L^p([0, 1]^d) \hookrightarrow L^p(\mathbb{P}_{\mu_0}),$$

such that $\iota_d(\text{Lip}_1([0, 1]^d; \|\cdot\|_\infty)) \subset \text{Lip}_{R/(\lambda_d^{(-\gamma_f+\gamma_d)/2})}(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}))$, where the Lipschitz norm on $[0, 1]^d$ is defined with respect to the ∞ -norm on $[0, 1]^d$.

Proof. By Assumption 8, $\mu \sim \mathbb{P}_{\mu_0}$ is sampled as

$$f_\mu = \sum_j \lambda_j^{\gamma_d/2} Z_j e_j, \quad \frac{d\mu}{d\lambda} = f_\mu,$$

where $Z_j, j \geq 1$ are independent and $Z_j \sim \rho_j(z) dz$. We define the cumulative distribution $\Phi_j(z) = \int_{-\infty}^z \rho_j dz$. We know that F_j is Lipschitz, whose Lipschitz constant is bounded by $\sup_j \|\rho_j\|_\infty \leq R$. We define $f_{\mu,i} = \langle f_\mu, e_i \rangle$. We will show that the mapping

$$\iota_d : L^p([0, 1]^d) \hookrightarrow L^p(\mathbb{P}_{\mu_0}), \quad (\iota_d g)(f_\mu) = g(\Phi_1(f_{\mu,1}/\lambda_1^{\gamma_d/2}), \dots, \Phi_d(f_{\mu,d}/\lambda_d^{\gamma_d/2}))$$

is the isometric embedding that we want. For $g \in L^p([0, 1]^d)$, the $L^p(\mathbb{P}_{\mu_0})$ -norm of $\iota_d g$ is equal to

$$\mathbb{E}_{f_\mu} |(\iota_d g)(f_\mu)|^p = \mathbb{E}_{f_\mu} |g(\Phi_1(f_{\mu,1}/\lambda_1^{\gamma_d/2}), \dots, \Phi_d(f_{\mu,d}/\lambda_d^{\gamma_d/2}))|^p$$

$$\begin{aligned}
&= \mathbb{E}_{f_\mu} |g(\Phi_1(Z_1), \dots, \Phi_d(Z_d))|^p \\
&= \int_{[0,1]^d} |g(x_1, \dots, x_d)|^p dx \quad (\Phi_j(Z_j) \sim \text{Unif}[0, 1]) \\
&= \|g\|_{L^p([0,1]^d)}^p,
\end{aligned}$$

which shows that ι_d is isometric embedding. Next, we evaluate the image of ι_d . A mapping

$$h_d : (B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}) \rightarrow ([0, 1]^d, \|\cdot\|_\infty), \quad f_\mu \mapsto (\Phi_1(f_{\mu,1}/\lambda_1^{\frac{\gamma_d}{2}}), \dots, \Phi_d(f_{\mu,d}/\lambda_d^{\frac{\gamma_d}{2}}))$$

is lipschitz because

$$\begin{aligned}
&\|h_d(f_{\mu_1}) - h_d(f_{\mu_2})\|_\infty \\
&\leq \max_j \left\{ \left| \Phi_j(f_{\mu_1,j}/\lambda_j^{\frac{\gamma_d}{2}}) - \Phi_j(f_{\mu_2,j}/\lambda_j^{\frac{\gamma_d}{2}}) \right| \right\} \\
&\leq \max_j \left\{ \text{Lip}(\Phi_j) \lambda_j^{-\frac{\gamma_d}{2}} |f_{\mu_1,j} - f_{\mu_2,j}| \right\} \\
&\leq \max_j \left\{ \text{Lip}(\Phi_j) \lambda_j^{\frac{\gamma_f - \gamma_d}{2}} \cdot \lambda_j^{-\frac{\gamma_f}{2}} |f_{\mu_1,j} - f_{\mu_2,j}| \right\} \\
&\leq \max_j \left\{ \text{Lip}(\Phi_j) \lambda_j^{\frac{\gamma_f - \gamma_d}{2}} \right\} \sum_j (\lambda_j^{-\frac{\gamma_f}{2}} |f_{\mu_1,j} - f_{\mu_2,j}|)
\end{aligned}$$

and thus

$$\text{Lip}(h_d) \leq \frac{R}{\lambda_d^{\frac{-\gamma_f + \gamma_d}{2}}}.$$

Therefore, $\forall g \in \text{Lip}_1([0, 1]^d, \infty)$, Lipschitz constant of $\iota_d g$ is bounded as

$$\text{Lip}(\iota_d g) = \text{Lip}(g \circ h_d) \leq \text{Lip}(g) \cdot \text{Lip}(h_d) \leq \frac{R}{\lambda_d^{\frac{-\gamma_f + \gamma_d}{2}}}.$$

Furthermore, we also have $\|\iota_d g\|_{C(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}))} \leq 1$. \square

Lemma 26 ensures that the Lipschitz class on $[0, 1]^d$ can be viewed as a subclass of our infinite-dimensional Lipschitz class, up to a scaling factor depending on λ_d . This immediately yields a lower bound on the packing numbers of our class in terms of the well-studied packing numbers of $\text{Lip}_1([0, 1]^d)$.

Corollary 5. *Under the assumptions of Lemma 26, we have*

$$\log \mathcal{M}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \epsilon)_{L^p(\mathbb{P}_{\mu_0})}) \gtrsim \log \mathcal{M} \left(\text{Lip}_1([0, 1]^d, \|\cdot\|_\infty), \frac{R\epsilon}{\lambda_d^{\frac{-\gamma_f + \gamma_d}{2}}} \right)_{L^p([0,1]^d)}$$

Proof. By rescaling the function, we have

$$\begin{aligned}
&\mathcal{M}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \epsilon)_{L^p}) \\
&= \mathcal{M} \left(\text{Lip}_{R/(\lambda_d^{(-\gamma_f + \gamma_d)/2})}(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \frac{R\epsilon}{(\lambda_d^{(-\gamma_f + \gamma_d)/2})} \right)_{L^p} \\
&\gtrsim \mathcal{M} \left(\iota_d(\text{Lip}_1([0, 1]^d; \ell^\infty)); \frac{R\epsilon}{(\lambda_d^{(-\gamma_f + \gamma_d)/2})} \right)_{L^p} \\
&\gtrsim \mathcal{M} \left(\text{Lip}_1([0, 1]^d; \ell^\infty); \frac{R\epsilon}{(\lambda_d^{(-\gamma_f + \gamma_d)/2})} \right)_{L^p([0,1]^d)}.
\end{aligned}$$

\square

Thus, the problem of estimating the packing entropy of the infinite-dimensional class reduces to that of estimating the entropy of finite-dimensional Lipschitz functions on the cube. Fortunately, sharp lower bounds for the latter are available in the literature. Note that packing and covering are almost equivalent when $\epsilon \rightarrow 0$.

Lemma 27 (Lanthaler (2024)). *For $p \in [1, \infty)$ and $d \in \mathbb{N}$, there exists a constant $c > 0$ independent of d such that*

$$\log \mathcal{M}(\text{Lip}_1([0, 1]^d, \|\cdot\|_\infty); \epsilon)_{\|\cdot\|_{L^p}} \gtrsim \left(\frac{c}{d\epsilon}\right)^d, \quad \forall \epsilon \in (0, c/d].$$

Combining the embedding argument with the finite-dimensional lower bounds above, we arrive at the following result, which provides the desired exponential lower bound on the entropy growth.

Lemma 28 (Parallel to Lanthaler (2024)). *Let \mathbb{P}_{μ_0} be a probability measure on $B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}})$ in Setting 4. The packing entropy of $\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}))$ with respect to the Bochner $L^p(\mathbb{P}_{\mu_0})$ -norm, satisfies*

$$\log \mathcal{M}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}), \epsilon)_{L^2(\mathbb{P}_{\mu_0})} \gtrsim \exp\left(c_l (\ln \epsilon^{-1})^{\frac{\alpha+1}{\alpha}}\right)$$

for some constant $c_l > 0$.

Proof. Combining Corollary 5 and Lemma 27,

$$\log \mathcal{M}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}), \epsilon)_{L^p} \gtrsim \left(\frac{c_1 \lambda_d^{-\frac{-\gamma_f + \gamma_d}{2}}}{8d\epsilon}\right)^d$$

where $c_1 > 0$ is a constant, provided $\epsilon \leq \frac{c_1 \lambda_d^{-\frac{-\gamma_f + \gamma_d}{2}}}{d}$. Then, by $\lambda_d \gtrsim \exp(-cd^\alpha)$, the RHS is lower bounded by

$$(\log \mathcal{M} \gtrsim) \left(\frac{c_1 \exp(-c'd^\alpha)}{d\epsilon}\right)^d \quad \text{if } \epsilon \leq c_2 d^{-1} \exp(-c'd^\alpha),$$

where $c' = \frac{c(-\gamma_f + \gamma_d)}{2}$ and $c_2 > 0$ is another constant. Assuming that ϵ is sufficiently small, let us take d as

$$d = \left(\frac{\log(c_2 \epsilon^{-1})}{c' + 1}\right)^{1/\alpha} \quad (\simeq (\ln \epsilon^{-1})^{1/\alpha}).$$

By rearranging the above inequality, we also obtain

$$\epsilon = c_2 \exp(-(c' + 1)d^\alpha),$$

which can satisfy $\epsilon \leq c_2 d^{-1} \exp(-c'd^\alpha)$ asymptotically, because $\exp(-d^\alpha) \ll d^{-1}$ where $d \simeq (\ln \epsilon^{-1})^{1/\alpha} \rightarrow \infty$ as $\epsilon \rightarrow 0$. Then, we have

$$\begin{aligned} & \log \mathcal{M}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}), \epsilon)_{L^p} \\ & \gtrsim \left(\frac{c_1 \exp(-c'd^\alpha)}{d\epsilon}\right)^d \\ & \gtrsim (\exp(-c'd^\alpha + (c' + 1)d^\alpha - \ln d + O(1)))^d \\ & \gtrsim \exp(\Omega(d^{\alpha+1})) \\ & \gtrsim \exp\left(\Omega((\ln \epsilon^{-1})^{\frac{\alpha+1}{\alpha}})\right), \end{aligned}$$

where we used $d^\alpha \simeq \ln \epsilon^{-1}$ in the fourth inequality. \square

C.5 PROOF OF MINIMAX LOWER BOUND

Theorem 7 (Minimax Lower Bound). *Under Setting 4, we have*

$$\inf_{S_n \rightarrow \hat{F}} \sup_{\hat{F}^* \in \mathcal{F}^*} \mathbb{E}_{S_n} \left[\|\hat{F} - F^*\|_{L^2(\mathbb{P}_{\nu, x_q})}^2 \right] \gtrsim \exp(-O((\ln n)^{\frac{\alpha}{\alpha+1}})).$$

Proof. By Lemma 21, it is sufficient to evaluate

$$\inf_{\mathcal{U}_n \mapsto \hat{F} \in L^2(\mathbb{P}_{\mu_0})} \sup_{F^* \in \mathcal{F}^*} \mathbb{E} \mathcal{U}_n [\|\hat{F}(\mu_0) - F^*(\mu_0)\|_{L^2(\mathbb{P}_{\mu_0})}^2]$$

where $\mathcal{U}_n = \{(\mu_0^{(i^*)}, y_t)\}_{t=1}^n$ sampled as in Setting 4, instead of the original problem. Let $V(\epsilon_n) := \log \mathcal{N}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \epsilon)_{L^2})$ and $M(\delta_n) := \log \mathcal{M}(\text{Lip}_1(B(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0^{\gamma_b}}, \|\cdot\|_{\mathcal{H}_0^{\gamma_f}}); \epsilon)_{L^2})$.

First, let $\epsilon_n = C_{1,\sigma} \exp(-c_E (\ln n)^{\frac{\alpha}{\alpha+1}})$ where $c_E = \left(\frac{1}{2c_u^2}\right)^{\frac{\alpha}{\alpha+1}}$. Then, by Lemma 25,

$$\begin{aligned} \frac{V(\epsilon_n)}{\epsilon_n^2} &\lesssim \exp\left(c_u (\log \epsilon_n^{-1})^{\frac{\alpha+1}{\alpha}} + 2c_E (\ln n)^{\frac{\alpha}{\alpha+1}}\right) \\ &\lesssim \exp\left(c_u (c_E (\ln n)^{\frac{\alpha}{\alpha+1}})^{\frac{\alpha+1}{\alpha}} + 2c_E (\ln n)^{\frac{\alpha}{\alpha+1}}\right) \\ &\lesssim \exp\left(\frac{1}{2} \ln n + O((\ln n)^{\frac{\alpha}{\alpha+1}})\right) \\ &\lesssim \frac{n}{\sigma^2}. \end{aligned}$$

Note that $\sigma = \Theta(1)$. Next, we lower bound the packing entropy. Taking $\delta_n = C_\sigma \exp(-c_d (\ln n)^{\frac{\alpha}{\alpha+1}})$ where $c_d = c_l^{-\frac{\alpha}{\alpha+1}}$ and $C_\sigma \gtrsim \frac{2}{\sigma^2} + \frac{2 \log 2}{1}$ is a sufficiently large constant only dependent of $\sigma = \Theta(1)$, using Lemma 28,

$$\frac{M(\delta_n)}{\epsilon_n^2} \gtrsim \exp\left(c_l (c_d^{\frac{\alpha+1}{\alpha}} \ln n) + O((\ln n)^{\frac{\alpha}{\alpha+1}})\right) \cdot \left(\frac{2}{\sigma^2} + \frac{2 \log 2}{1}\right) \gtrsim n \left(\frac{2}{\sigma^2} + \frac{2 \log 2}{n}\right).$$

Finally, applying Lemma 22, we have

$$\inf_{\mathcal{S}_n \mapsto \hat{F}} \sup_{F \in \mathcal{F}^\circ} \mathbb{E} \left[\|\hat{F} - F\|_{L^2}^2 \right] \gtrsim \delta_n^2 \gtrsim \exp(-O((\ln n)^{\frac{\alpha}{\alpha+1}})).$$

□

D SYNTHETIC EXPERIMENT ON MEASURE-VALUED ATTENTION

To provide a minimal empirical sanity check of our risk bounds, we design a simple synthetic experiment where the input is a *measure-valued context* on $[0, 1] \times \{-1, +1\}$ and the model is a single MLP→Attention→MLP block. The goal is to recover a scalar functional of an underlying “associative” measure from a mixture of associative and non-associative components.

Data-generating process. Fix a truncation level $M = 16 \in \mathbb{N}$ and an orthonormal trigonometric basis on $[0, 1]$,

$$\phi_0(x) = 1, \quad \phi_j(x) = \sqrt{2} \sin(\pi j x), \quad j = 1, \dots, M-1.$$

For a smoothness parameter $\alpha > 0$ we define eigenvalues

$$\lambda_j = \exp(-j^\alpha), \quad j = 1, \dots, M-1.$$

For each training example we sample two sets of coefficients $Z_1, Z_2 \sim \mathcal{N}(0, I_{M-1})$ independently, $Z_{1,0}, Z_{2,0} = 0$, and form the (unnormalized) densities on $[0, 1]$

$$\tilde{\mu}_k(x) = \sum_{j=0}^{M-1} \lambda_j Z_{k,j} \phi_j(x), \quad k \in \{1, 2\}.$$

We discretize $[0, 1]$ on a uniform grid $\{x_t\}_{t=1}^T$, $T = 32$, clamp the density to be nonnegative, and normalize to obtain a probability mass function $(p_k(t))_{t=1}^T$:

$$\mu_k^{\text{raw}}(x_t) = \tilde{\mu}_k(x_t), \quad \mu_k^+(x_t) = \max\{\mu_k^{\text{raw}}(x_t), \varepsilon\}, \quad p_k(t) = \frac{\mu_k^+(x_t)}{\sum_{s=1}^T \mu_k^+(x_s)},$$

with a small cutoff $\varepsilon > 0^3$

Independently, we sample a ‘‘query label’’ $v^{(1)} \in \{-1, +1\}$ uniformly and set $v^{(2)} := -v^{(1)}$. We then define a product-measure mixture on $[0, 1] \times \{-1, +1\}$ by

$$\nu = \frac{1}{2}(\mu_1 \otimes \delta_{v^{(1)}}) + \frac{1}{2}(\mu_2 \otimes \delta_{v^{(2)}}),$$

where μ_k is the discrete measure assigning mass $p_k(t)$ to x_t . To construct the input token sequence, we draw $n_{\text{tokens}} = 5000$ i.i.d. Monte Carlo samples $(X_i, Q_i) \sim \nu$:

$$(X_i, V_i) = \begin{cases} (x_{T_1}, v^{(1)}), & \text{with prob. } \frac{1}{2}, T_1 \sim p_1, \\ (x_{T_2}, v^{(2)}), & \text{with prob. } \frac{1}{2}, T_2 \sim p_2, \end{cases} \quad i = 1, \dots, n_{\text{tokens}}.$$

Finally, we append a single ‘‘query token’’ $(0, v^{(1)})$ at the end of the sequence, so that each input example is a sequence

$$\{(X_i, V_i)\}_{i=1}^{n_{\text{tokens}}} \cup \{(0, v^{(1)})\} \in ([0, 1] \times \{-1, +1\})^{n_{\text{tokens}}+1}.$$

The target output Y depends only on the *associative* measure μ_1 (and is independent of μ_2 and $v^{(2)}$):

$$Y \simeq \tilde{F}^*(\mu_1, \underbrace{X_{n_{\text{tokens}}+1}}_{\text{query token}}) := v^{(1)} \cdot \sum_{j=0}^{M-1} \lambda_j Z_{1,j}^2.$$

Intuitively, the model must use the final query token $(0, v^{(1)})$ to attend to tokens consistent with q_1 and recover information about the hidden coefficients Z_1 from Monte Carlo samples of μ_1 . Note that we add a small Gaussian noise with $\text{std} = 0.01$ in training.

Model and training. We use a minimal architecture that mirrors the theoretical measure-attention operator:

$$\text{context/query MLP} \rightarrow \text{measure attention} \rightarrow \text{MLP HEAD}.$$

For each example we construct a sequence of T_{ctx} context tokens $(x_t, v_t) \in \mathbb{R}^2$ together with a final query token $(0, v_{\text{query}})$. The context tokens and the query token are embedded by separate two-layer MLPs into $\mathbb{R}^{d_{\text{model}}}$ with $d_{\text{model}} = 8$ and hidden width $d_{\text{hidden}} = 8$. The resulting query embedding provides the Q vector, while the context embeddings provide the K and V vectors for a single 4-head softmax attention layer. This ‘‘measure-attention’’ layer outputs a single d_{model} -dimensional representation, which is fed into a final two-layer MLP head to produce the scalar prediction \hat{Y} . We train with the squared loss $\ell(\hat{Y}, Y) = (\hat{Y} - Y)^2$ using Adam with an exponentially decaying learning rate for 20 epochs.

For each $\alpha \in \{\alpha_1, \dots, \alpha_L\}$ we generate independent training sets of sizes $n \in \{n_{\text{min}}, \dots, n_{\text{max}}\}$ ($n = 2^k$ for $k = 2, \dots, 6$) and measure the empirical risk $L(n)$ on a held-out validation set ($n_{\text{val}} = 2000$).

Risk scaling. Theory predicts that in this setting the minimax risk decays as

$$L^*(n) \approx \exp(-c(\log n)^{\alpha/(\alpha+1)}),$$

up to multiplicative constants. To compare with this prediction, for each α we fit the parametric form

$$\log L(n) \approx A_\alpha - C_\alpha (\log n)^{\alpha/(\alpha+1)}$$

by least squares over (A_α, C_α) using the measured pairs $\{(\log n_i, \log L(n_i))\}_i$. Figure 4 shows $\log L(n)$ against $(\log n)^{\alpha/(\alpha+1)}$ together with the fitted curves.

As a minimal sanity check, this synthetic experiment (Fig. 4) in which varying the spectral decay parameter α systematically affects the convergence speed: heavier-tailed spectra (smaller α) lead to visibly slower decay of the empirical risk. This is qualitatively consistent with the theoretical prediction, although we do not attempt to match the precise asymptotic rate.

³In our theory, we did not explicitly investigated such an cutoff or the normalization for simplicity.

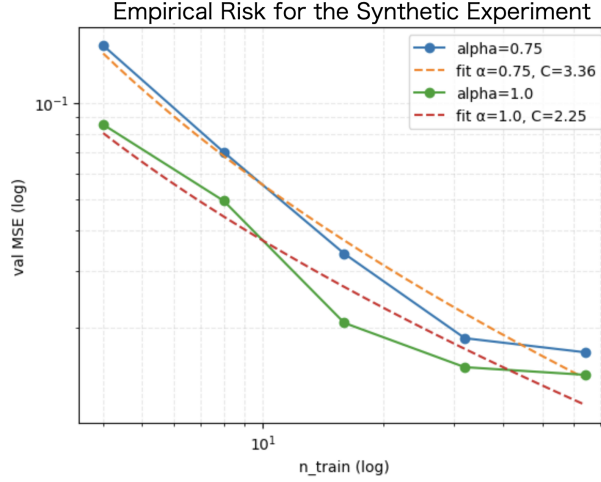


Figure 4: Empirical risk $L(n)$ for the synthetic measure-valued experiment, plotted on a transformed axis $(\log n)^{\alpha/(\alpha+1)}$ together with the fitted curves. Each risk was calculated with 2000 unknown samples.

Attention-weight analysis. To check whether the attention layer actually uses the query tag, we inspect the softmax attention weights of the trained model on the validation set. For a given example, let

$$\{(X_t, V_t)\}_{t=1}^T \in ([0, 1] \times \{-1, +1\})^T$$

denote the T context tokens, and let

$$(X_q, V_q) = (0, v^{(1)})$$

be the query token appended at the end of the sequence. For each attention head $h = 1, \dots, H$ we write

$$a^{(h)} \in [0, 1]^T$$

for the softmax attention weights from the query to the T context positions, so that

$$\sum_{t=1}^T a_t^{(h)} = 1.$$

We are interested in how the query token redistributes its attention mass over tokens whose tag matches the query versus those with the opposite tag. Accordingly, we define the index sets

$$S_{\text{same}} := \{1 \leq t \leq T : V_t = V_q\}, \quad S_{\text{diff}} := \{1 \leq t \leq T : V_t \neq V_q\},$$

that is, we only consider the context tokens and exclude the query token itself from both sets. For each head h we then compute the average per-token attention weight assigned by the query to same-tag and different-tag tokens,

$$\bar{w}_{\text{same}}^{(h)} := \frac{1}{|S_{\text{same}}|} \sum_{t \in S_{\text{same}}} a_t^{(h)}, \quad \bar{w}_{\text{diff}}^{(h)} := \frac{1}{|S_{\text{diff}}|} \sum_{t \in S_{\text{diff}}} a_t^{(h)},$$

as well as the total attention mass

$$m_{\text{same}}^{(h)} := \sum_{t \in S_{\text{same}}} a_t^{(h)}, \quad m_{\text{diff}}^{(h)} := \sum_{t \in S_{\text{diff}}} a_t^{(h)}.$$

In practice, we implement this by adding a flag to the attention module that, when enabled, stores the last softmax attention tensor $A \in \mathbb{R}^{B \times H \times 1 \times T}$ on the CPU after a forward pass (where B is the batch size), and we extract $a^{(h)}$ as the length- T vector $A_{b,h,1,:}$ corresponding to the query-to-context weights for each example b and head h .

Table 1: Average attention weight from the query token to same-tag vs. different-tag context tokens (mean and standard deviation over 1000 validation examples, with 64 training data and $\alpha = 1.0$). With $T_{\text{ctx}} = 5000$ and roughly half of the tokens sharing the query tag, the uniform baseline is $\bar{w} \approx 2 \times 10^{-4}$, while a head that focuses almost exclusively on same-tag tokens reaches $\bar{w}_{\text{same}} \approx 4 \times 10^{-4}$. We also report the validation MSE with the original queries and after shuffling queries within 1000 data for validation.

Head	\bar{w}_{same}	\bar{w}_{diff}	$\text{std}(w_{\text{same}})$	$\text{std}(w_{\text{diff}})$
0	1.96×10^{-4}	2.04×10^{-4}	2.80×10^{-4}	2.85×10^{-4}
1	1.44×10^{-19}	4.00×10^{-4}	2.10×10^{-19}	4.00×10^{-4}
2	4.00×10^{-4}	1.17×10^{-14}	4.00×10^{-4}	1.70×10^{-14}
3	4.00×10^{-4}	2.86×10^{-27}	4.00×10^{-4}	0 (too small)
mean	2.49×10^{-4}	1.51×10^{-4}	2.70×10^{-4}	1.71×10^{-4}

	original queries	shuffled queries
val MSE	1.44×10^{-2}	7.75×10^{-1}

We report in Table 1 the mean and standard deviation of $m_{\text{same}}^{(h)}$ and $m_{\text{diff}}^{(h)}$ over 1000 validation examples for each attention head h . On this synthetic task, two of the four heads concentrate essentially all of their attention mass on tokens whose tag matches the query tag, while another head exhibits the opposite preference and one head remains nearly symmetric. Averaged across heads, the query token assigns a larger total mass to tokens with the same tag than to those with the opposite tag, indicating a net bias toward tag-conditioned retrieval.

The absolute scale of the averaged weights \bar{w}_{same} and \bar{w}_{diff} is small (on the order of 10^{-4}) simply because the attention distribution is normalized over a long context of $T_{\text{ctx}} \approx n_{\text{tokens}} = 5000$ positions. Under an approximately uniform baseline, we would have

$$\bar{w}_{\text{unif}} \approx \frac{1}{T_{\text{ctx}}} \approx \frac{1}{5000} \approx 2 \times 10^{-4},$$

so the reported values should be interpreted relative to this $1/T_{\text{ctx}}$ scale rather than as absolute probabilities. In our construction, each context token independently comes from $\mu_1 \otimes \delta_{v(1)}$ or $\mu_2 \otimes \delta_{v(2)}$ with probability $1/2$, so typically $|S_{\text{same}}| \approx |S_{\text{diff}}| \approx T_{\text{ctx}}/2$. Consequently, values around $\bar{w} \approx 2 \times 10^{-4}$ correspond to almost-uniform attention over all context tokens, whereas values around $\bar{w}_{\text{same}} \approx 4 \times 10^{-4}$ and $\bar{w}_{\text{diff}} \approx 0$ indicate that a head places essentially all of its attention mass on the same-tag subset (and analogously for the opposite tag).

As a sanity check that the model genuinely uses the query input, we perform a ‘‘query shuffle’’ experiment at evaluation time: within each mini-batch we randomly permute the last (query) token across examples, while keeping the context tokens and targets fixed, and recompute the validation loss (the bottom of Table 1). On this synthetic task, shuffling the query tokens increases the validation MSE, confirming that the model relies nontrivially on the query input.

This minimal experiment is not intended as a thorough empirical study, but it provides a sanity check that the qualitative order of the risk predicted by our theory is reproducible in a simple measured attention setting.