

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 UNDERSTANDING TRANSFORMERS FOR TIME SERIES: RANK STRUCTURE, FLOW-OF-RANKS, AND COMPRESSIBILITY

Anonymous authors

Paper under double-blind review

ABSTRACT

Transformers are widely used across data modalities, and yet the principles distilled from text models often transfer imperfectly to models trained to other modalities. In this paper, we analyze Transformers through the lens of rank structure. Our focus is on the time series setting, where the structural properties of the data differ remarkably from those of text or vision. We show that time-series embeddings, unlike text or vision, exhibit sharply decaying singular value spectra: small patch sizes and smooth continuous mappings concentrate the data into low-rank subspaces. From this, we prove that the associated $Q/K/V$ projections admit accurate low-rank approximations, and that attention layers become compressible in proportion to the decay of the embedding spectrum. We introduce the concept of *flow-of-ranks*, a phenomenon by which nonlinear mixing across depth inflates the rank, explaining why early layers are most amenable to compression and why ranks grow with depth. Guided by these theoretical and empirical results, we use these insights to compress Chronos, a large time series foundation model, achieving a reduction of 65% in inference time and 81% in memory, without loss of accuracy. Our findings provide principled guidance for allocating width, depth, and heads in time series foundation models, and for exploiting their inherent compressibility.

1 INTRODUCTION

Transformers, originally designed for language (Lewis et al., 2020; Achiam et al., 2023), are now widely deployed, e.g., time series (Ansari et al., 2024a; Das et al., 2024; Shi et al., 2025; Wolff et al., 2025), images (Liu et al., 2021; Dosovitskiy et al., 2020), molecules (Maziarka et al., 2020; Leon et al., 2024), and DNA sequences (Ji et al., 2021; Le et al., 2021; Nguyen et al., 2024). A common approach to apply Transformers to these other data modalities is to directly transfer architectural parameters (e.g., width, heads, depth) from text-based models, on the assumption that what works for text should generalize. However, this assumption is fragile. As an example, we show that time series differ fundamentally from language in how signals are tokenized and embedded. This leads to the more general question: how well do community insights on pretraining and hyperparameter tuning, based largely on Transformers applied to text data, port to Transformers applied to other data modalities? Understanding the answer to this question is particularly important when Transformers are applied in domains where data are less abundant than in the text domain.

Here, we address this question in the context of time-series data and time-series forecasting. A priori, the answer to this question is not obvious: time series data have similarities with text data (e.g., they have obvious sequential properties), but they also have many differences (e.g., it is not clear that they are well-modelable by a discrete set of tokens). The question, however, is timely: time series data are ubiquitous in many domains, including scientific (Abhishek et al., 2012; Zhang & Gilpin, 2025; Lai et al., 2025), industrial (Hong et al., 2016), and financial applications (Zhang et al., 2001), where they facilitate critical tasks, e.g., forecasting (Hyndman & Athanasopoulos, 2018), imputation (Yoon et al., 2018), and anomaly detection (Blázquez-García et al., 2021). In addition, there has been recent growing interest in developing so-called time series foundation models (TSFMs) (Brown et al., 2020; Radford et al., 2021; Ansari et al., 2024a). These models are large pretrained models, designed with the hope of providing a foundation (Bommasani et al., 2021), to be adaptable to a

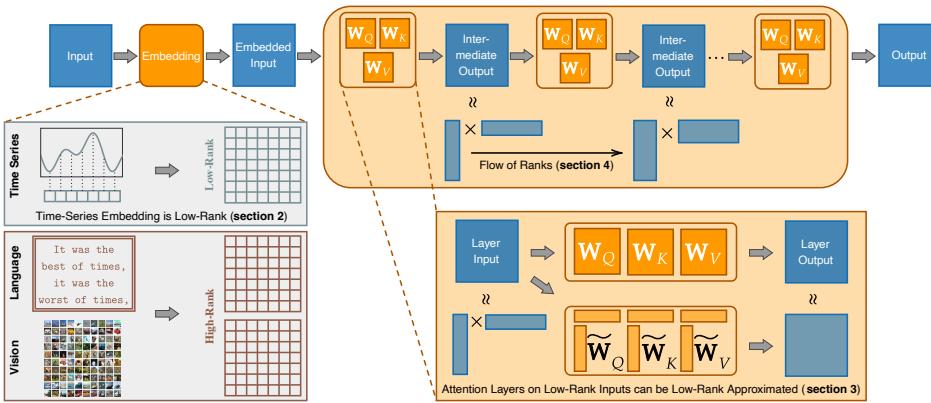


Figure 1: Overview of our results. We show that the embedded inputs of Transformers trained with time-series data have much lower ranks than those of other modalities, including Vision Transformers and Transformers trained with language data (see section 2); we prove that attention matrices on low-rank inputs are well-approximated by low-rank matrices (see section 3); and we introduce and demonstrate a concept called flow-of-ranks, describing how attention matrices in earlier layers are more compressible than those in later layers (see section 4).

wide range of domains and time series tasks. As with other models that aim to provide such a “foundation,” TSFMs are appealing because they reduce the need for task-specific architectures and parameter computation, thereby enabling the transfer to new settings with relatively little effort.

In this paper, we develop a framework for analyzing design decisions in Transformers; see Figure 1 for an overview and Appendix A for related work. Our view is that these decisions are guided by an understanding of the structural properties of the data modality, which are then inherited by properties of the model. Our framework enables a detailed linear algebraic analysis of the attention layers within a Transformer, which consist of three linear transformations to form the queries, keys, and values. Our approach is general, and we apply it to the time series domain. Lastly, we show that we can use our approach as a practical tool for compressing TSFMs. We summarize the flow of our paper and our main contributions as follows:

1. **Data modality and rank structure.** We compare Transformers trained to different data modalities through the lens of numerical rank, and we demonstrate that time-series data lead to a particularly low-rank structure, at the level of inputs. To the best of our knowledge, our work is the first to directly study which data modality features make Transformer models well-approximated via a truncated singular value decomposition (SVD), and why. We also show how standard time-series embeddings preserve low-rank structure in the hidden space, which differs from large-vocabulary text and other modalities. (See section 2.)
2. **From low-rank inputs to low-rank attention.** Next, we provide the first general theoretical results that connect low-rank embeddings to low-rank attention matrices, and we make clear how width and the number of heads control the quality of low-rank approximations. **While these results apply generally to any low-rank embeddings, we illustrate them in the case of time series.** In addition, our results are sharp: we show that high-rank embedded inputs, which appear in data modalities such as text and vision, lead to incompressible attention matrices. (See section 3.)
3. **Flow-of-ranks.** We introduce the “flow-of-ranks” concept to describe how the numerical rank changes across layers in a deep Transformer. This extends our earlier analysis from a single attention layer to the setting of a deep Transformer, where nonlinear activations, residual mixing, and normalization gradually increase the rank of a representation. We note that “flow-of-ranks” explains why early layers are often better approximated by SVD than later ones. (See section 4.)
4. **Compressibility of real-world TSFMs.** Finally, we illustrate one application of our insights: compressing a real-world TSFM. We demonstrate that the same set of hyperparameters (i.e., width, depth, and number of heads) more severely over-parameterizes TSFMs than it does LLMs. This leads us to develop two complementary compression strategies: compressing a pretrained model and pretraining a model that is compressed by design. In particular, we show that by compressing a Chronos model, we can reduce the inference time by 65.4% and memory usage by 81.4%, at no cost of predictive performance. Overall, our results demonstrate and explain

Table 1: Comparison of embedding strategies and patch sizes across TSFMs.

Model	Chronos	WaveToken	TOTEM	Time-MOE	Chronos-Bolt	TimesFM
Strategy	Quantization	Quantization	Quantization	Continuous	Continuous	Continuous
Patch Size	1	1	1	1	16	32

why state-of-the-art TSFMs are highly compressible, in practice, compared to (language-trained) LLMs of the same size. (See section 5.)

Additional supporting material may be found in the appendices.

2 DATA MODALITY AND RANK STRUCTURE OF EMBEDDING

In this section, we investigate the structure of time-series embeddings and compare them to embeddings from other data modalities. Our goal is to understand, from both a theoretical and empirical perspective, why time-series inputs often look low-rank after embedding.

We consider univariate time series. In particular, let $\mathbf{x} = (x_1, \dots, x_T) \in \mathbb{R}^{1 \times T}$ be a univariate input of length T . Note that, unlike other data modalities, the input \mathbf{x} is a rank-1 matrix.¹ The first step in a TSFM is to map \mathbf{x} into a high-dimensional sequence via an embedding function $\Phi : \mathbb{R}^{1 \times T} \rightarrow \mathbb{R}^{d \times L}$, where d denotes the hidden dimension of the model and L denotes the new sequence length, possibly different from T due to patching. This embedding is typically constructed by applying a trainable function $\phi : \mathbb{R}^k \rightarrow \mathbb{R}^d$ to disjoint patches of size k . Assuming $L = T/k$ is an integer, this gives:

$$\Phi(\mathbf{x}) = (\phi(x_1, \dots, x_k), \phi(x_{k+1}, \dots, x_{2k}), \dots, \phi(x_{(L-1)k+1}, \dots, x_{Lk} = x_T)).$$

These patch embedding functions ϕ fall into two main categories (see Figure 8 in Appendix B):

- A **quantization-based embedding** partitions the input space \mathbb{R}^k into V disjoint regions, $\mathbb{R}^k = \biguplus_{i=1}^V R_i$. Each region R_i is mapped to a unique trainable vector $\mathbf{u}_i \in \mathbb{R}^d$, so that $\phi(\mathbf{x}) = \sum_{i=1}^V \mathbb{1}_{\{\mathbf{x} \in R_i\}} \mathbf{u}_i$. This approach is used in several TSFMs, e.g., Chronos (Ansari et al., 2024a), WaveToken (Masserano et al., 2025), and TOTEM (Talukder et al., 2024).
- A **continuous embedding** uses a parameterized function, typically a neural network $\phi(\cdot; \theta)$, to map patches directly. This strategy is also used in many TSFMs, e.g., Chronos-Bolt (Ansari et al., 2024b), Moirai (Woo et al., 2024), TimesFM (Das et al., 2024), and Time-MoE (Shi et al., 2025).

Table 1 summarizes the design choices for these prominent TSFMs.

In most TSFMs, the patch size is significantly smaller than the hidden dimension (i.e., $k \ll d$), meaning ϕ maps from a low-dimensional space to a higher-dimensional one. Intuitively, if ϕ is well-behaved, it should embed the low-dimensional space \mathbb{R}^k into a corresponding low-dimensional submanifold $\phi(\mathbb{R}^k) \subset \mathbb{R}^d$. To formalize this notion of dimensionality, we use singular values. Let \mathbf{U} be a linear operator between some Hilbert spaces and $\sigma_1 \geq \dots \geq \sigma_n \geq 0$ be its singular values, where $1 \leq n \leq \infty$. While the algebraic rank of an object \mathbf{U} , i.e., the number of its non-zero singular values, is a strict measure, real-world data is noisy; therefore, we use the more practical concept of numerical rank. For a tolerance $\varepsilon > 0$, the ε -rank of \mathbf{U} denotes the number of its singular values that are significant relative to the largest one:

$$\text{rank}_\varepsilon(\mathbf{U}) = |\{j \mid \sigma_j(\mathbf{U})/\sigma_1(\mathbf{U}) > \varepsilon\}|. \quad (1)$$

A low numerical rank implies \mathbf{U} is well-approximated by an operator $\tilde{\mathbf{U}}$ with $\text{rank}(\tilde{\mathbf{U}}) = \text{rank}_\varepsilon(\mathbf{U})$.

Our central hypothesis is that for a large corpus of input patches $\{\mathbf{x}^{(i)}\}_{i=1}^N$, the resulting embedded matrix, via quantization or a continuous embedding, $\mathbf{U} = [\phi(\mathbf{x}^{(1)}) \ \dots \ \phi(\mathbf{x}^{(N)})]$ has a low numerical rank, which is smaller than the ambient dimension d . To test this hypothesis, we sample

¹In the case of a few-variate time series, \mathbf{x} becomes a few-rank matrix, making the analysis in the paper directly generalizable by viewing the number of variates as another patch dimension. We note that in some time-series applications with a large number of input variables, the embedded inputs may not be low-rank, and our analysis may not directly apply in such cases.

162 thousands of patches from diverse signals (e.g., sinusoids, exponential functions, and white noise),
 163 and we compute the singular value decay of their embeddings. For contrast, we perform the same
 164 analysis on a tabular foundation model (TFM) (Mitra) (Zhang & Robinson, 2025) using 1000 syn-
 165 synthetic tables, a T5 LLM processing text from Dickens’ *A Tale of Two Cities*, and a ViT processing
 166 1000 randomly sampled images from CIFAR-10.

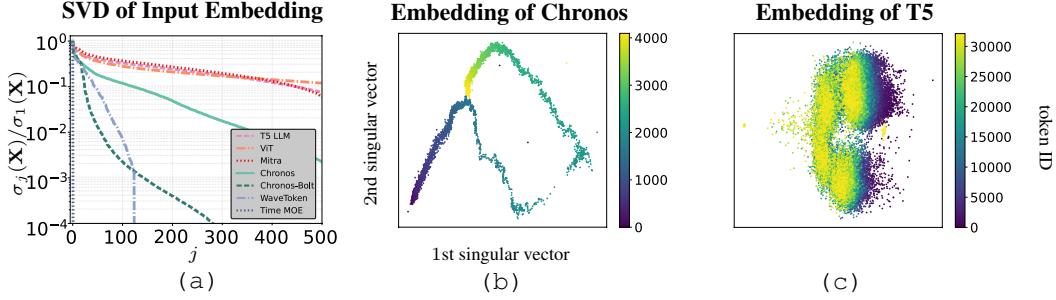


Figure 2: (a): Singular values of the embedded input matrices from many different TSFMs, a TFM, a ViT, and an LLM. (b, c): Embedding space of Chronos and a T5 LLM, respectively, visualized by projecting them onto the leading two singular vectors of the embedding matrix.

As shown in Figure 2 (a), the singular values of TSFM embeddings decay dramatically faster than those from the tabular and language models, which confirms their significantly lower numerical rank. Figure 2 (b, c) visualize the embedding spaces of Chronos (quantized) and T5 (language tokens) by projecting them onto their top two singular vectors. The Chronos embedding, mapping a quantized real line, reveals a clear low-dimensional structure, whereas the T5 embedding, likely due to its vocabulary properties, appears far less well-structured.

One may wonder: why is it that Chronos-Bolt ($k = 16$) produces a lower-rank embedding than Chronos ($k = 1$)? This seeming surprise arises from their different embedding mechanisms. A quantization-based model like Chronos initially maps adjacent values to random, unstructured vectors; it must learn the geometry of the real line during training. In contrast, a continuous embedding like Chronos-Bolt uses a smooth neural network $\phi(\cdot; \theta)$. This architectural choice imposes smoothness from the start, ensuring that even a randomly initialized model maps the low-dimensional patch space \mathbb{R}^k to a low-dimensional submanifold in \mathbb{R}^d (see Appendix H).

Although understanding a quantization-based embedding requires looking into the training dynamics, which is beyond the scope of this paper, we theoretically analyze how a continuous embedding preserves low-rank structures. The following theorem formalizes this intuition.

Theorem 1. Given any hidden dimension $d > 1$, let $\phi : [-1, 1] \rightarrow \mathbb{R}^d$ be a function that embeds $[-1, 1]$ into \mathbb{R}^d . Given L arbitrary points x_1, \dots, x_L sampled from $[-1, 1]$, define

$$\Xi = \begin{bmatrix} \phi_1(x_1) & \cdots \\ \vdots & \ddots \\ \phi_d(x_1) & \cdots \end{bmatrix} \in \mathbb{R}^{d \times [-1, 1]}, \quad \Psi = \begin{bmatrix} \phi_1(x_1) & \cdots & \phi_1(x_L) \\ \vdots & \ddots & \vdots \\ \phi_d(x_1) & \cdots & \phi_d(x_L) \end{bmatrix} \in \mathbb{R}^{d \times L}.$$

Let $s_1 \geq \dots \geq s_d \geq 0$ and $\sigma_1 \geq \dots \geq \sigma_d \geq 0$ be the singular values of the quasimatrix Ξ and matrix Ψ , respectively. Then, the following statement holds:

1. If, for some $V > 0$, $\nu \geq 1$, and every $1 \leq i \leq d$, we have ϕ_i and its derivative through $\phi_i^{(\nu)}$ are absolutely continuous on $[-1, 1]$ and $\phi_i^{(\nu)}$ is of bounded variation V , then we have

$$s_{j+1} \leq \frac{4V\sqrt{d}}{\pi\nu(j-1-\nu)^\nu} = \mathcal{O}(j^{-\nu}\sqrt{d}), \quad \sigma_{j+1} \leq \frac{2V\sqrt{dL}}{\pi\nu(j-1-\nu)^\nu} = \mathcal{O}(j^{-\nu}\sqrt{dL}), \quad \nu+1 < j \leq d-1.$$

2. If, for some $M > 0$ and every $1 \leq i \leq d$, ϕ_i has an analytic continuation to the Bernstein ellipse of radius $\rho > 1$ (see Trefethen (2019)), whose infinity norm is no greater than M , then we have

$$s_{j+1} \leq \frac{4M\sqrt{d}\rho^{-j+1}}{\rho-1} = \mathcal{O}(\rho^{-j}\sqrt{d}), \quad \sigma_{j+1} \leq \frac{2M\sqrt{dL}\rho^{-j+1}}{\rho-1} = \mathcal{O}(\rho^{-j}\sqrt{dL}), \quad 0 \leq j \leq d-1.$$

We refer interested readers to [Townsend & Trefethen \(2015\)](#) for the precise definition of a quasimatrix (informally, it is a “matrix” in which one of the dimensions is discrete as usual but the other is continuous). Theorem 1 guarantees that for univariate patches ($k = 1$), a smooth embedding function yields singular values with guaranteed decay rates: polynomial decay of order ν for functions with ν continuous derivatives, and exponential decay for analytic functions.

See Appendix B for a proof of Theorem 1 using classic univariate polynomial approximation techniques. Using this result, we can directly explain the low-rank structure observed in models like Time-MoE (see Figure 2 and Corollary 2). While multivariate polynomial approximation results enable us to extend Theorem 1, the size of the polynomial basis up to a fixed degree increases exponentially with k , which makes it less practically relevant for a larger patch size k . Instead, for Chronos-Bolt, where $k = 16$, we seek an ad hoc result that works when the embedding $\phi(\cdot; \theta)$ is an MLP (see Appendix C for a proof of Theorem 2 below).

Theorem 2. Consider the input embedding defined by a two-layer residual MLP:

$$\Phi(\mathbf{X}) = \mathbf{W}_3\mathbf{X} + \mathbf{W}_2\omega(\mathbf{W}_1\mathbf{X}), \quad \mathbf{X} \in \mathbb{R}^{k \times L}, \quad \mathbf{W}_1 \in \mathbb{R}^{d_f \times k}, \quad \mathbf{W}_2 \in \mathbb{R}^{d \times d_f}, \quad \mathbf{W}_3 \in \mathbb{R}^{d \times k},$$

where k denotes the patch size, L denotes the number of patches, d_f denotes the hidden-layer dimension in an MLP, $d > k$ denotes the hidden dimension of the Transformer, and ω denotes any activation function satisfying that $|\omega(x)| \leq |x|$ for every $x \in \mathbb{R}$. Then, for any $\varepsilon > 0$, we have

$$|\{j \mid \sigma_j(\Phi(\mathbf{X})) > \varepsilon \|\mathbf{W}_2\|_2 \|\mathbf{W}_1\mathbf{X}\|_2\}| \leq \min\{d, (1 + \varepsilon^{-2})k\}.$$

Theorem 2 states that the numerical rank of the continuous embedding is bounded by a quantity dependent linearly on the patch size k , not the much larger ambient dimension d . The term $\|\mathbf{W}_2\|_2 \|\mathbf{W}_1\mathbf{X}\|_2$ reflects the natural scaling of $\sigma_1(\Phi(\mathbf{X}))$ in the definition of rank_ε . In practice, we have $k \ll d$, meaning that MLP embeds an input patch into a low-dimensional subspace in \mathbb{R}^d .

To illustrate this theorem, we pretrain 6 Chronos-Bolt models of different patch sizes, and we compute the singular values of their embeddings. Figure 3 (a) shows that the numerical rank of the embedding increases with the patch size k . We also perform an in-depth analysis, where for each pair of embedded inputs $\phi(\mathbf{x}^{(i)})$ and $\phi(\mathbf{x}^{(j)}) \in \mathbb{R}^d$, we compute the angle between them as follows:

$$\theta(\phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)})) = |\phi(\mathbf{x}^{(i)})^\top \phi(\mathbf{x}^{(j)})| / (\|\phi(\mathbf{x}^{(i)})\|_2 \|\phi(\mathbf{x}^{(j)})\|_2) \in [0, \pi/2]. \quad (2)$$

The larger the θ , the more linearly independent the two vectors are. We illustrate this in Figure 3 (b), where brighter heatmaps correspond to higher rank matrices. We see that the embedded input, $\Phi(\mathbf{X})$ from Chronos-Bolt, which is a subset of $\mathbb{R}^d = \mathbb{R}^{768}$, spans a subspace of significantly smaller dimension, i.e., the image of $\mathbb{R}^k = \mathbb{R}^{16}$ under ϕ .

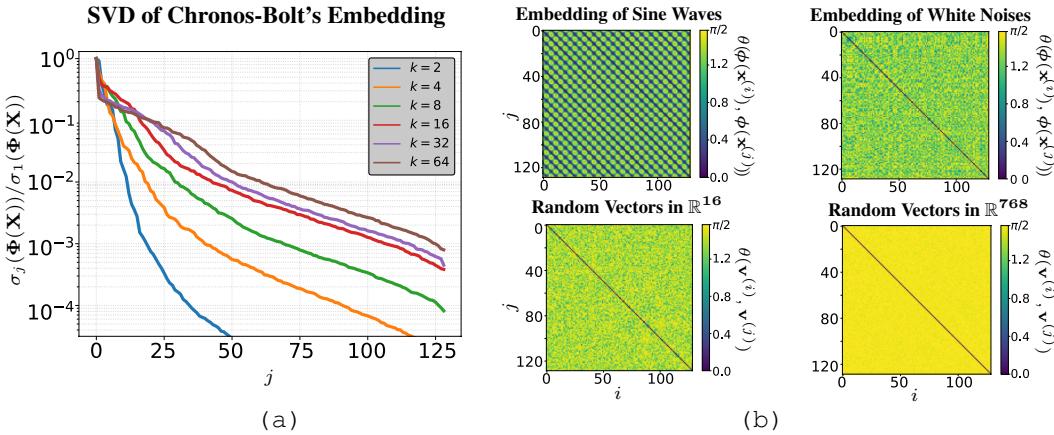


Figure 3: (a) Singular values of the embedded input matrices from Chronos-Bolt models pre-trained with different patch sizes k . (b) Angles between Chronos-Bolt’s embedded vectors in $\mathbb{R}^d = \mathbb{R}^{768}$ defined in eq. (2), where the patches $\mathbf{x}^{(i)}$ are from a sinusoidal wave and Gaussian white noises, respectively. We also plot the angles between i.i.d. random Gaussian vectors in $\mathbb{R}^k = \mathbb{R}^{16}$ and $\mathbb{R}^d = \mathbb{R}^{768}$ for comparison.

270 **3 FROM LOW-RANK INPUTS TO LOW-RANK ATTENTION MATRICES**
 271

272 Let $\mathbf{U} \in \mathbb{R}^{d \times L}$ be an input embedded in the hidden space. Recall that in section 2 we showed that for
 273 TSFMs, \mathbf{U} often has a low numerical rank. This immediately implies \mathbf{U} can be expressed in a low-
 274 rank format: $\mathbf{U} \approx \mathbf{U}_1 \mathbf{U}_2$, where $\mathbf{U}_1 \in \mathbb{R}^{d \times \tilde{d}}$ and $\mathbf{U}_2 \in \mathbb{R}^{\tilde{d} \times L}$ for some $\tilde{d} \ll d$. This yields faster
 275 matrix-matrix products with \mathbf{U} , but a limitation is that this representation requires an expensive
 276 rank-revealing matrix factorization (Damle et al., 2024), which adds overhead, particularly during
 277 backpropagation. If so, how do we leverage the low-rank structure of TSFM embeddings?

278 From basic linear algebra, it is known that for a linear operator $\mathbf{T} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ to act on an r -
 279 dimensional subspace, one only needs to specify the operator in r directions, in which case the
 280 operator can then be well-approximated by a low-rank matrix $\tilde{\mathbf{T}}$ whose rank scales with r instead of
 281 the full width d (Damle et al., 2024; Ipsen & Saibaba, 2024). An attention layer, defined by
 282

$$283 \quad \text{Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) = \mathbf{W}_V \mathbf{U} \text{softmax}\left(\frac{\mathbf{U}^\top \mathbf{W}_Q^\top \mathbf{W}_K^\top \mathbf{U}}{\sqrt{d}}\right), \quad \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d},$$

284

285 while nonlinear, contains three linear transformations: namely, the queries, the keys, and the values.
 286 In this section, we establish Theorem 3, which supports the low-rank representations of \mathbf{W}_Q , \mathbf{W}_K ,
 287 and \mathbf{W}_V given a low-rank \mathbf{U} (see Appendix D for the proof).

288 **Theorem 3.** Let $C > 0$ be a constant. Let $\Xi = [\mathbf{x}_1 \cdots \mathbf{x}_N] \in \mathbb{R}^{d \times N}$ be given for some
 289 $d, N \geq 1$, $\mathbf{x}_j \in \mathbb{R}^d$, and $\|\mathbf{x}_j\|_2 \leq C$ for all $1 \leq j \leq N$. Let $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$ be matrices
 290 such that $\|\mathbf{W}_Q^\top \mathbf{W}_K\|_2 \leq C\sqrt{d}$, and $\|\mathbf{W}_V\|_2 \leq C\sqrt{d}$. The following two statements hold:

- 291 1. **(Attention matrices are compressible on low-rank inputs.)** For any $\tilde{d} < d$ such that
 292 $\sigma_{\tilde{d}+1} := \sigma_{\tilde{d}+1}(\Xi) \leq 1$, there exist $\tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V \in \mathbb{R}^{d \times d}$ with $\text{rank}(\tilde{\mathbf{W}}_Q) = \text{rank}(\tilde{\mathbf{W}}_K) =$
 293 $\text{rank}(\tilde{\mathbf{W}}_V) = \tilde{d}$, $\|\tilde{\mathbf{W}}_Q^\top \tilde{\mathbf{W}}_K\|_2 \leq \|\mathbf{W}_Q^\top \mathbf{W}_K\|_2$, and $\|\tilde{\mathbf{W}}_V\|_2 \leq \|\mathbf{W}_V\|_2$, such that given any
 294 matrix $\mathbf{U} \in \mathbb{R}^{d \times L}$ for any $L \geq 1$, where each column of \mathbf{U} is a column of Ξ , we have that
 295

$$296 \quad \left\| \text{Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) - \text{Attention}(\mathbf{U}; \tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V) \right\|_F \leq \mathcal{O}(\sqrt{d} \sigma_{\tilde{d}+1}), \quad (3)$$

297 where the constant in the \mathcal{O} -notation only depends on C .

- 298 2. **(Attention matrices are incompressible on high-rank inputs.)** The upper bound in eq. (3) is
 299 tight up to a factor of \sqrt{d} . That is, fix some $d \geq 1$, $L \geq d$ and $1 \geq \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d > 0$.
 300 There exist $\mathbf{U} \in \mathbb{R}^{d \times L}$ with $\sigma_j(\mathbf{U}) = \sigma_j$ for all $1 \leq j \leq d$, and $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d}$ such that for
 301 any $\tilde{d} < d$, any orthogonal matrix $\mathbf{W}_V \in \mathbb{R}^{d \times d}$, and any rank- \tilde{d} matrix $\tilde{\mathbf{W}}_V \in \mathbb{R}^{d \times d}$, we have
 302

$$303 \quad \left\| \text{Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V) - \text{Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \tilde{\mathbf{W}}_V) \right\|_F \geq \frac{1}{4} \sigma_{\tilde{d}+1}. \quad (4)$$

304 Note that Theorem 3 is a purely numerical statement that holds for any low-rank embedding Ξ .
 305 It applies to time series because our analysis in section 2 reveals the low-rank nature of Ξ , but it
 306 nonetheless works for other modalities provided that Ξ is low-rank. The first statement of Theorem 3
 307 says, at a high level, that if the inputs \mathbf{U} come from a low-rank vocabulary Ξ , i.e., $\sigma_{\tilde{d}+1}$ is small,
 308 then one only needs low-rank attention matrices $\tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V$ on \mathbf{U} . Here, it is important that
 309 we fix the low-rank embedded space Ξ and prove a uniform bound eq. (3) that holds for all input \mathbf{U} ,
 310 so that our low-rank approximation $\tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V$ is not input-dependent. Since TSFMs have a
 311 low-rank vocabulary Ξ (see section 2), high-rank attention matrices in TSFMs can be approximated
 312 by low-rank ones. See Appendix J for a numerical experiment. We emphasize that this low-rank
 313 property does not depend on the temporal simplicity of a particular time series (e.g., being constant),
 314 but rather on the intrinsic low-dimensionality of the input embedding space Ξ itself, which makes \mathbf{U}
 315 low-rank-sufficient for any time-series input. For other data modalities than time-series, e.g., TFMNs,
 316 ViTs, and LLMs (see Figure 2), where \mathbf{U} has a higher rank, the second statement of Theorem 3
 317 suggests that the attention matrices are less compressible, as we will show in section 5.

318 To illustrate the concepts in Theorem 3 on TSFMs, we take pretrained Chronos models and T5
 319 LLMs of different hidden dimensions d . These two models have the same Transformer size and
 320 architecture, and differ only in the pretraining data modality. For each model and a fixed $\varepsilon > 0$, we
 321

compute the averaged ε -rank of all query projection matrices \mathbf{W}_Q from the model. Figure 4 (a, b) show that as the size of \mathbf{W}_Q increases, the matrix \mathbf{W}_Q stays low-rank in a Chronos model, and its rank scales almost linearly with d in a T5 LLM. For Chronos models, where the vocabulary is embedded in a low-rank subspace, low-rank attention matrices suffice to capture most information in the inputs. This empirical observation goes beyond Theorem 3, which proves the sufficiency of low-rank $\mathbf{W}_{Q/K/V}$, but it does not guarantee their appearance via training. Figure 4 (a) shows that training, while independent from expressiveness, gives rise to low-rank weights (see Appendix G for more insights). This provides motivation for pretraining a compressed model and for compressing a pretrained one (see section 5). In Appendix G, we show an analysis of the distribution of the singular values of attention matrices, in pretrained foundation models and also during training.

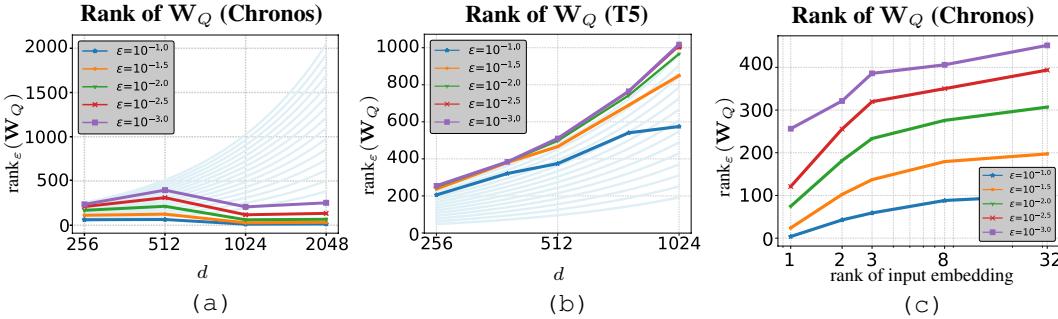


Figure 4: The averaged ε -rank of query projection matrices \mathbf{W}_Q in pretrained Chronos models and T5 LLMs. In (a, b), we vary the hidden dimension d . The light blue curves are contours of the ratio between the horizontal and the vertical axes in the semilog-x scale. In (c), we fix the hidden dimension $d = 512$ and change the rank of the fixed input embedding Ξ (see Appendix I).

To further corroborate the role of the rank \tilde{d} in Theorem 3, we pretrain Chronos models with a fixed vocabulary. We increase the (algebraic) rank of Ξ from 1 to 32 (see Appendix I), where we observe that the numerical ranks of attention matrices increase with $\text{rank}(\Xi)$ (see Figure 4 (c)).

While our discussion in this section assumes a single head in the attention layer, in Appendix E we extend the result to the multi-head case. From the standpoint of representation complexity, Theorem 3 applies equally to multi-head attention, and it is independent of the number of heads (see Theorem 5 in Appendix D). In practice, however, we observe that the $\mathbf{W}_{Q/K/V}$ matrices in pretrained multi-head layers tend to have higher numerical ranks (see Figure 9). This effect can be understood through the lens of numerical linear algebra, via a concept called “sketching” (see Appendix E). It helps explain why additional heads can improve robustness and training stability, even though the underlying complexity result is head-agnostic (see Theorem 6 in Appendix E).

4 FLOW-OF-RANKS: MOVING THROUGH A TRANSFORMER

Our prior discussions in section 3 focused on compressing a single attention layer, but a Transformer in a TSFM has many layers. While we showed that the input into the first attention layer is low-rank (see section 2), one may wonder: what about the inputs into a later layer? If you apply a linear transformation to a low-rank input, then it at most preserves, if not decreases, the rank of the input; however, each layer in a Transformer is generally nonlinear. Here, we demonstrate that these nonlinear layers can increase the rank of the input.

We refer to the phenomenon of the increasing rank of the input as it goes deeper into the layers of a model as the *flow-of-ranks*. Our contribution here is twofold: we quantify this rank growth across layers in a deep Transformer by proving Theorem 4 and then show what that means to Transformers applied to time-series data. To see the “flow-of-ranks” in practice, we show the numerical ranks of attention matrices per layer in a Chronos model and a T5 LLM in Figure 5 (a, b). The numerical ranks of attention matrices become higher for deeper layers because of the flow-of-ranks (see Figure 5 (c)). For a T5 LLM, only with a large ε do we observe an increase in the ε -rank, because the ε -rank with a small ε is already saturated (i.e., close to the matrix’s dimension) to capture the high-rank vocabulary Ξ . To formalize this discussion, we prove how an attention layer increases the small singular values of a low-rank input matrix. Conceptually, this “flow-of-ranks” can be viewed as a mechanism by which the model lifts a low-dimensional but complex signal into

378
 379
 380
 381
 382
 383
 384
 385
 386
 387
 388
 389
 390
 391
 392
 393
 394
 395
 396
 397
 398
 399
 400
 401
 402
 403
 404
 405
 406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 a higher-dimensional space where the underlying autocorrelation becomes simpler. This is analogous to the Koopman operator framework, which similarly transforms nonlinear dynamics into a higher-dimensional representation with approximately linear evolution.

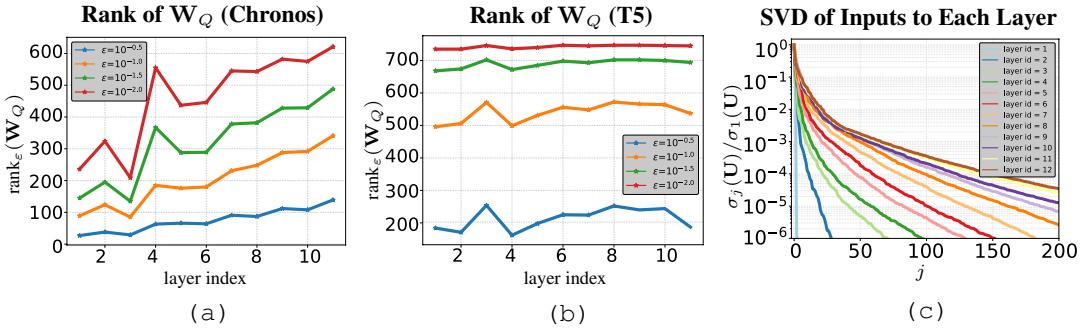


Figure 5: (a, b): The ϵ -rank of every query projection \mathbf{W}_Q in the encoder of a Chronos model and a T5 model of the same size, respectively. (c): Singular values of the input matrix to each Chronos' encoder layer, starting with a constant signal $(x_1, \dots, x_T) = (0, \dots, 0)$.

Theorem 4. Given positive integers $d \leq L$, let $\mathbf{U} \in \mathbb{R}^{d \times L}$ be an input matrix with singular values $1 = \sigma_1 \geq \dots \geq \sigma_d > 0$. Let $D \geq 1$ be the number of layers of the model. Let h be the number of heads and d_h be the per-head dimension so that $d = h \times d_h$. The following statements hold:

1. Suppose $\sqrt{D} \geq 2e^2h$. For every $1 \leq i \leq h$, let $\|\mathbf{W}_Q^{(i)^\top} \mathbf{W}_K^{(i)}\|_2 \leq \sqrt{d_h}$, and $\|\mathbf{W}_V^{(i)}\|_2 \leq 1$ be any attention matrices (see eq. (18) for the notation). For any $1 \leq k \leq d$, we have

$$\sigma_k(\mathbf{Z})/\sigma_1(\mathbf{Z}) \leq 2 \min_{1 \leq j \leq k} \left(\sigma_{k-j+1} + \frac{e^2h}{\sqrt{D}} \sigma_{\lfloor (j-1)/h \rfloor + 1} \right), \quad (5)$$

where $\mathbf{Z} = \mathbf{U} + \mathbf{Y}/\sqrt{D}$ and $\mathbf{Y} = \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h)$.

2. The upper bound above is tight. More precisely, given any d, L, N, h , and singular values $1 = \sigma_1 \geq \dots \geq \sigma_d > 0$, there exist an input $\mathbf{U} \in \mathbb{R}^{d \times L}$ with singular values $\sigma_1, \dots, \sigma_d$, and matrices $\mathbf{W}_Q, \mathbf{W}_K$, and \mathbf{W}_V with $\|\mathbf{W}_V^{(i)}\|_2 \leq 1$, such that for any $1 \leq k \leq d$, we have

$$\sigma_k(\mathbf{Z})/\sigma_1(\mathbf{Z}) \geq \frac{1}{4} \left(\sigma_{(i-1)d_h + \lceil k/i \rceil} + \frac{1}{\sqrt{D}} \sigma_{\lceil k/i \rceil} \right), \quad (6)$$

for every $i \leq h$ such that $\lceil k/i \rceil < d_h$.

See Appendix F for the proof of Theorem 4 and Appendix J for a numerical experiment. The matrix \mathbf{Z} can be interpreted as the output of a residual attention layer given input \mathbf{U} , where the scaling factor $1/\sqrt{D}$ is commonly used in a Transformer to stabilize and accelerate training (De & Smith, 2020). The upper bound in eq. (5) provides a guarantee that the rank of the output cannot be arbitrarily large if the input is low-rank. In this bound, there are two terms balancing each other: σ_{k-j+1} , which increases as j increases, and $\sigma_{\lfloor (j-1)/h \rfloor + 1}$, which decreases as j increases. Without knowing the precise distribution of the singular values, the j which minimizes eq. (5) cannot be determined. To understand why the upper bound is tight, we prove a corollary in a simplified one-layer setting.

Corollary 1. Using the notations in Theorem 4 and assuming $D = 1$, we have $\sigma_k(\mathbf{Z}) \leq \mathcal{O}(h)\sigma_{\lfloor (k+1)/(h+1) \rfloor + 1}$. In addition, the lower bound eq. (6) satisfies that $\sigma_k(\mathbf{Z}) \geq \mathcal{O}(1)\sigma_{\lceil k/h \rceil}$ for every $k \leq d - d_h$. The constants in both \mathcal{O} -notations are universal.

The corollary states that when \mathbf{U} goes through an attention layer, the k^{th} singular value of the output \mathbf{Z} can be increased to at most the order of magnitude of the $\lceil k/h \rceil^{\text{th}}$ singular value of \mathbf{U} . The lower bound suggests that this can be achieved by some inputs, which proves the sharpness of the upper bound. Interestingly, this result says that the number of heads h plays an important role in increasing the ranks of an input matrix. The term to note here is not the $\mathcal{O}(h)$ factor multiplied to the upper bound but the division by $h+1$ in the subindex $\lfloor (k+1)/(h+1) \rfloor + 1$. When h is large and the singular values decay fast, $\sigma_{\lfloor (k+1)/(h+1) \rfloor + 1}(\mathbf{U})$ can be significantly larger than $\sigma_k(\mathbf{U})$ (see Figure 17).

In Figure 6, we also verify the flow-of-ranks in many TSFMs beyond Chronos. In general, we observe that the rank of an attention weight matrix in many TSFMs grows as a function of the layer

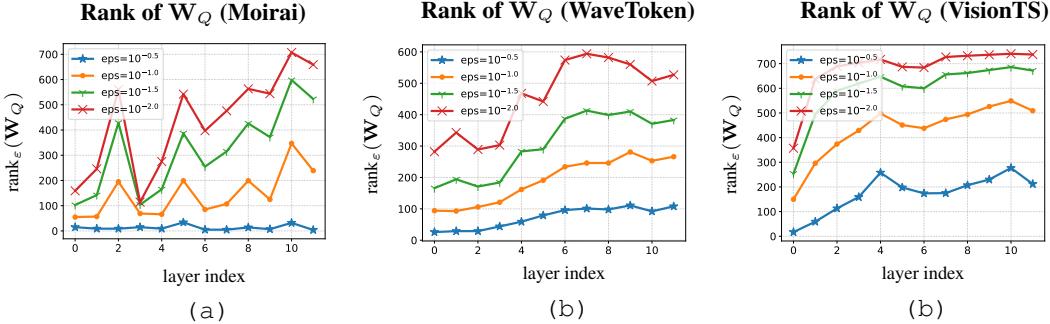


Figure 6: Flow-of-ranks is also observed for many TSFMs other than Chronos, including Moirai-1.0-R-base (Woo et al., 2024), WaveToken (base) (Masserano et al., 2025), and VisionTS (Chen et al., 2025). The hidden dimension d in all three models are 768.

index. The three models compared in Figure 6 also have distinct input space dimensions, with that of VisionTS significantly larger than those of Moirai and WaveToken, which is also reflected by the overall larger numerical ranks of the weight matrices of VisionTS.

5 USING THESE INSIGHTS: HOW TO COMPRESS A TSFM?

In this section, we apply our theoretical framework to compress large TSFMs. We pursue two approaches: first, motivated by the observation that attention matrices in many pretrained TSFMs are already low rank, we apply truncated SVD to each attention matrix of a pretrained model. Second, to achieve stronger compression, we design architectures that parameterize attention matrices in low rank from the start and pretrain them from scratch, using a layer-dependent hyperparameter schedule that accounts for the “flow-of-ranks.” We present results on compressing Chronos, and we provide more results on another TSM in Appendix K to support the broad applicability of our methods. **Unlike LoRA’s low-rank fine-tuning updates, we factorize our attention weights themselves, yielding inference-time efficiency and enabling pretraining with a compressed backbone.**

Compressing a Pretrained Model. In Figure 4, we see that attention matrices in a large pretrained model usually do not have a full numerical rank. That means we can well-approximate a large matrix with a low-rank one. More precisely, let $\mathbf{W} \in \mathbb{R}^{d \times d}$ be an attention matrix and let $\mathbf{W} = \sum_{j=1}^d \sigma_j \mathbf{u}_j \mathbf{v}_j^\top$ be the singular value decomposition (SVD) of \mathbf{W} . Then, for a fixed $\varepsilon > 0$, the truncated SVD satisfies a relative error bound: $\left\| \mathbf{W} - \sum_{j=1}^{\text{rank}_\varepsilon(\mathbf{W})} \sigma_j \mathbf{u}_j \mathbf{v}_j^\top \right\|_2 \leq \varepsilon \|\mathbf{W}\|_2$ (see eq. (1)). For a fixed $\varepsilon > 0$, we apply the truncated SVD to every attention matrix of the pretrained Chronos model without fine-tuning, which reduces the number of parameters. We compute the WQL and MASE losses (see Ansari et al. (2024a)) and their geometric means relative to the original model. A score below 1 means the reduced model performs better while it performs worse otherwise.

Table 2: Results of compressing a pretrained Chronos or T5 model. We apply truncated SVD with a specific ε , which results in a model whose attention matrices are compressed to the “Ratio” of the original size. We compare the performance scores relative to the original pretrained model, where “LPPL” stands for the logarithm of the perplexity. Overlap is obtained by selecting the top-10 tokens from the original model and the compressed one, and computing their Jaccard overlap.

Ratio	Chronos				T5		
	In-Domain ↓ WQL	MASE	Zero-Shot ↓ WQL	MASE	Overlap ↑	LPPL ↓	Overlap ↑
0.073	4.409	4.095	3.562	3.435	0.308	3.313	0.227
0.151	1.991	2.412	1.566	1.576	0.717	2.530	0.301
0.237	1.053	1.005	1.030	1.011	0.883	1.652	0.345
0.393	1.009	1.024	0.990	0.994	0.979	1.544	0.568
0.569	1.003	0.952	0.945	0.995	0.997	1.290	0.730
0.755	1.007	1.021	1.027	1.000	0.999	1.085	0.871
0.889	1.000	0.960	1.016	1.001	0.999	1.028	0.954
1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000

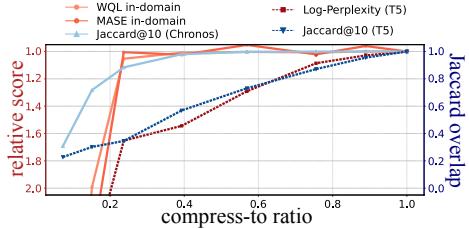


Table 2 shows that we can compress the attention matrices up to around 23.7% of the original size without any loss of performance. Reducing the size further takes away key information in attention matrices and results in a rapid performance deterioration, leading us to the next compression method.

Table 3: Results of pretraining a compressed Moirai-1.0-R-base model. We compare the performance scores *relative to the original pretrained model*. We show prediction losses when the flow-of-ranks is or is not used. The last row is the baseline.

Size Ratio	With Flow-of-ranks				Without Flow-of-ranks			
	\tilde{d}_0	α	WQL ↓	MASE ↓	\tilde{d}_0	α	WQL ↓	MASE ↓
0.250	8	0.34	1.001	1.014	16	0.00	1.069	1.050
0.500	10	0.58	0.996	1.007	32	0.00	1.038	1.036
1.000	-	-	-	-	64	0.00	1.000	1.000

Pretraining a Compressed Model. Table 2 shows a hard limit of compression: compressing the size to below about 20% makes the performance significantly worse. If one is given enough budget, a more robust method to obtain a smaller TSFM is by pretraining a compressed model. We parameterize a $d \times d$ attention matrix by a rank- \tilde{d} representation. Driven by “flow-of-ranks” from section 4, we choose to let \tilde{d} be layer-dependent. In the i th layer of the model, we use a $\tilde{d}_i = \tilde{d}(i) = \lceil \tilde{d}_0(1+i)^\alpha \rceil$, where $\tilde{d}_0 > 0$ and $\alpha \geq 0$ are two hyperparameters. This design is motivated by the fact that the numerical rank of an input to a layer increases as the layer index, and we need attention matrices of higher ranks to capture them (see Theorem 3). In Figure 7, we see how pretraining a compressed model allows us to expand the time–accuracy Pareto frontier of TSFMs, which show that pretraining a compressed model is more robust than compressing a pretrained one. In fact, in Table 4 of Appendix K, we show that pretraining a compressed Chronos-Bolt allows us to outperform even traditional local methods on *both* time and accuracy.

\tilde{d}_0	α	Size Ratio	Inference		Embedding From Scratch		Reuse Embedding	
			Time	Space	In-Domain	Zero-Shot	In-Domain	Zero-Shot
					WQL ↓	MASE ↓		
3	0.27	0.075	0.346	0.186	1.034	0.988	0.966	0.982
5	0.35	0.150	0.398	0.312	1.048	0.982	1.080	1.055
7	0.40	0.250	0.494	0.440	1.021	0.949	0.996	1.019
64	0.00	1.000	1.000	1.000	1.000	1.000	1.000	1.000

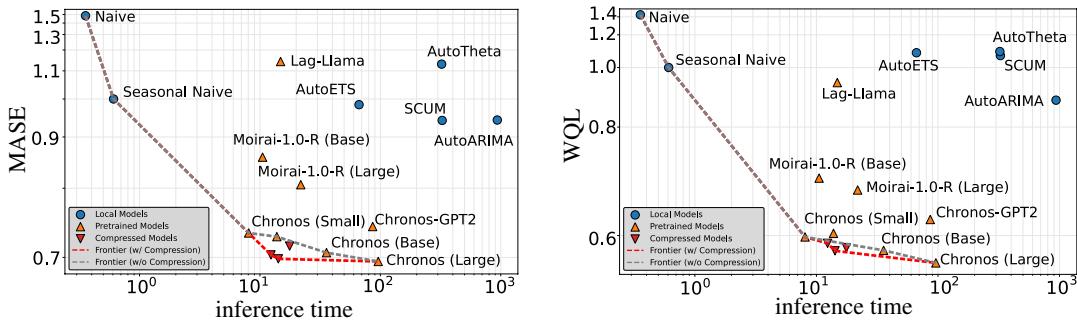


Figure 7: Results of pretraining a compressed Chronos model. We compare the performance scores *relative to the original pretrained model*. We show prediction losses for both models whose embedding matrix is randomly initialized and models whose embedding matrix is inherited from the original pretrained model. The last row is the baseline.

To show the generality of our analysis in this paper, we also pretrain compressed Moirai-1.0-R-base models. We show the results in Table 3. It clearly shows the benefit of the flow-of-ranks design: given a fixed compression ratio, a model with a different reduced rank per layer performs significantly better than one whose rank remains layer-independent.

6 CONCLUSION

We have developed a data/modality-dependent framework via the lens of rank structure for analyzing the structure of and design decisions for Transformers, and we have applied it to Chronos and Chronos-Bolt, two popular TSFMs. Our results highlight how properties of the model depend on and interact with properties of the input data, and they lead to concrete principles for the design of models, parameters, and hyperparameters. We illustrate our proof-of-principle results by showing the compressibility of TSFMs in comparison to Transformers trained on text data.

540 REFERENCES
541

542 Kumar Abhishek, Maheshwari Prasad Singh, Saswata Ghosh, and Abhishek Anand. Weather fore-
543 casting model using artificial neural network. *Procedia Technology*, 4:311–318, 2012.

544 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
545 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
546 report. *arXiv preprint arXiv:2303.08774*, 2023.
547

548 Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen,
549 Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al.
550 Chronos: Learning the language of time series. *arXiv preprint arXiv:2403.07815*, 2024a.

551 Abdul Fatir Ansari, Caner Turkmen, Oleksandr Shchur, and Lorenzo Stella.
552 Fast and accurate zero-shot forecasting with Chronos-Bolt and Autogluon,
553 2024b. URL [https://aws.amazon.com/blogs/machine-learning/
554 fast-and-accurate-zero-shot-forecasting-with-chronos-bolt-and-autogluon/](https://aws.amazon.com/blogs/machine-learning/fast-and-accurate-zero-shot-forecasting-with-chronos-bolt-and-autogluon/).
555

556 Ane Blázquez-García, Angel Conde, Usue Mori, and Jose A Lozano. A review on outlier/anomaly
557 detection in time series data. *ACM computing surveys (CSUR)*, 54(3):1–33, 2021.

558 Rishi Bommasani et al. On the opportunities and risks of foundation models. *arXiv preprint
559 arXiv:2108.07258*, 2021.
560

561 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
562 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
563 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

564 Cristian Challu, Kin G Olivares, Boris N Oreshkin, Federico Garza Ramirez, Max Mergenthaler
565 Canseco, and Artur Dubrawski. Nhits: Neural hierarchical interpolation for time series forecast-
566 ing. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 6989–6997,
567 2023.
568

569 Beidi Chen, Tri Dao, Eric Winsor, Zhao Song, Atri Rudra, and Christopher Ré. Scatterbrain: Uni-
570 fying sparse and low-rank attention. *Advances in Neural Information Processing Systems*, 34:
571 17413–17426, 2021.

572 Mouxiang Chen, Lefei Shen, Zhuo Li, Xiaoyun Joy Wang, Jianling Sun, and Chenghao Liu. Vi-
573 sionTS: Visual masked autoencoders are free-lunch zero-shot time series forecasters. *PMLR*,
574 2025.
575

576 Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas
577 Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention
578 with performers. *International Conference on Learning Representations*, 2021.

579 Anil Damle, Silke Glas, Alex Townsend, and Annan Yu. How to reveal the rank of a matrix? *arXiv
580 preprint arXiv:2405.04330*, 2024.
581

582 Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
583 time-series forecasting. In *Forty-first International Conference on Machine Learning*, 2024.
584

585 Soham De and Samuel L. Smith. Batch normalization biases residual blocks towards the identity
586 function in deep networks. In *Neural Information Processing Systems Foundation*, 2020.

587 Emily L Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting linear
588 structure within convolutional networks for efficient evaluation. *Advances in neural information
589 processing systems*, 27, 2014.
590

591 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
592 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An
593 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint
arXiv:2010.11929*, 2020.

- 594 P. Drineas and M. W. Mahoney. RandNLA: Randomized numerical linear algebra. *Communications*
595 *of the ACM*, 59:80–90, 2016.
- 596
- 597 Shibo Feng, Peilin Zhao, Liu Liu, Pengcheng Wu, and Zhiqi Shen. Hdt: Hierarchical discrete
598 transformer for multivariate time series forecasting. *arXiv preprint arXiv:2502.08302*, 2025.
- 599
- 600 Everett S Gardner Jr. Exponential smoothing: The state of the art. *Journal of forecasting*, 4(1):
601 1–28, 1985.
- 602
- 603 Leon Götz, Marcel Kollovieh, Stephan Günnemann, and Leo Schwinn. Byte pair encoding for
604 efficient time series forecasting. *arXiv preprint arXiv:2505.14411*, 2025.
- 605
- 606 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
607 *preprint arXiv:2312.00752*, 2023.
- 608
- 609 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
610 state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- 611
- 612 Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization
613 of diagonal state space models. *Advances in Neural Information Processing Systems*, 35:35971–
614 35983, 2022.
- 615
- 616 Haokun Gui, Xiucheng Li, and Xinyang Chen. Vector quantization pretraining for eeg time series
617 with random projection and phase alignment. In *International Conference on Machine Learning*,
618 pp. 16731–16750. PMLR, 2024.
- 619
- 620 Tao Hong, Pierre Pinson, Shu Fan, Hamidreza Zareipour, Alberto Troccoli, and Rob J Hyndman.
621 Probabilistic energy forecasting: Global energy forecasting competition 2014 and beyond, 2016.
- 622
- 623 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
624 Weizhu Chen, et al. LoRA: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- 625
- 626 Rob J Hyndman and George Athanasopoulos. *Forecasting: principles and practice*. OTexts, 2018.
- 627
- 628 Ilse CF Ipsen and Arvind K Saibaba. Stable rank and intrinsic dimension of real and complex
629 matrices. *arXiv preprint arXiv:2407.21594*, 2024.
- 630
- 631 Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks
632 with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.
- 633
- 634 Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V Davuluri. Dnabert: pre-trained bidirectional
635 encoder representations from transformers model for dna-language in genome. *Bioinformatics*,
636 37(15):2112–2120, 2021.
- 637
- 638 Jeffrey Lai, Anthony Bao, and William Gilpin. Panda: A pretrained forecast model for universal
639 representation of chaotic dynamics. *arXiv preprint arXiv:2505.13755*, 2025.
- 640
- 641 Nguyen Quoc Khanh Le, Quang-Thai Ho, Trinh-Trung-Duong Nguyen, and Yu-Yen Ou. A trans-
642 former architecture based on bert and 2d convolutional neural network to identify dna enhancers
643 from sequence information. *Briefings in bioinformatics*, 22(5):bbab005, 2021.
- 644
- 645 Vadim Lebedev, Yaroslav Ganin, Maksim Rakhuba, Ivan Oseledets, and Victor Lempitsky.
646 Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *arXiv preprint*
647 *arXiv:1412.6553*, 2014.
- 648
- 649 Miguelangel Leon, Yuriy Perezhohin, Fernando Peres, Aleš Popovič, and Mauro Castelli. Compar-
650 ing smiles and selfies tokenization for enhanced chemical language modeling. *Scientific Reports*,
651 14(1):25016, 2024.
- 652
- 653 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer
654 Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training
655 for natural language generation, translation, and comprehension. *Proceedings of the 58th Annual*
656 *Meeting of the Association for Computational Linguistics*, 2020.

- 648 Bryan Lim, Sercan Ö Arik, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for
 649 interpretable multi-horizon time series forecasting. *International journal of forecasting*, 37(4):
 650 1748–1764, 2021.
- 651
- 652 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
 653 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the
 654 IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 655 Spyros Makridakis and Michele Hibon. Arma models and the box-jenkins methodology. *Journal
 656 of forecasting*, 16(3):147–163, 1997.
- 657
- 658 Charles H Martin and Michael W Mahoney. Implicit self-regularization in deep neural networks:
 659 Evidence from random matrix theory and implications for learning. *Journal of Machine Learning
 660 Research*, 22(165):1–73, 2021.
- 661 Luca Masserano, Abdul Fatir Ansari, Boran Han, Xiyuan Zhang, Christos Faloutsos, Michael W
 662 Mahoney, Andrew Gordon Wilson, Youngsuk Park, Syama Rangapuram, Danielle C Maddix,
 663 et al. Enhancing foundation models for time series forecasting via wavelet-based tokenization. In
 664 *International Conference on Machine Learning*. PMLR, 2025.
- 665 Łukasz Maziarka, Tomasz Danel, Sławomir Mucha, Krzysztof Rataj, Jacek Tabor, and Stanisław
 666 Jastrzebski. Molecule attention transformer. *arXiv preprint arXiv:2002.08264*, 2020.
- 667
- 668 Eric Nguyen, Michael Poli, Matthew G Durrant, Brian Kang, Dhruva Katrekar, David B Li, Liam J
 669 Bartie, Armin W Thomas, Samuel H King, Garyk Brix, et al. Sequence modeling and design
 670 from molecular to genome scale with evo. *Science*, 386(6723):eado9336, 2024.
- 671 Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64
 672 words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*, 2022.
- 673
- 674 Boris N Oreshkin, Dmitri Carpow, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis
 675 expansion analysis for interpretable time series forecasting. *arXiv preprint arXiv:1905.10437*,
 676 2019.
- 677 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
 678 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
 679 models from natural language supervision. In *International conference on machine learning*, pp.
 680 8748–8763. PMLR, 2021.
- 681 David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic fore-
 682 casting with autoregressive recurrent networks. *International journal of forecasting*, 36(3):1181–
 683 1191, 2020.
- 684
- 685 Pratyusha Sharma, Jordan T Ash, and Dipendra Misra. The truth is in there: Improving reasoning
 686 in language models with layer-selective rank reduction. *arXiv preprint arXiv:2312.13558*, 2023.
- 687
- 688 Xiaoming Shi, Shiyu Wang, Yuqi Nie, Dianqi Li, Zhou Ye, Qingsong Wen, and Ming Jin. Time-moe:
 689 Billion-scale time series foundation models with mixture of experts. *International Conference on
 690 Learning Representations*, 2025.
- 691
- 692 Sabera Talukder, Yisong Yue, and Georgia Gkioxari. Totem: Tokenized time series embeddings for
 693 general time series analysis. *arXiv preprint arXiv:2402.16412*, 2024.
- 694
- 695 Alex Townsend and Lloyd N Trefethen. Continuous analogues of matrix factorizations. *Proceedings
 696 of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 471(2173):20140585,
 697 2015.
- 698
- 699 Lloyd N Trefethen. *Approximation theory and approximation practice, extended edition*. SIAM,
 700 2019.
- 701
- 702 Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention
 703 with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020.
- 704
- 705 Andrew Gordon Wilson. Deep learning is not so mysterious or different. *PMLR*, 2025.

- 702 Malcolm L Wolff, Shenghao Yang, Kari Torkkola, and Michael W Mahoney. Using pre-trained
703 LLMs for multivariate time series forecasting. *arXiv preprint arXiv:2501.06386*, 2025.
704
- 705 Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo.
706 Unified training of universal time series forecasting transformers. 2024.
- 707 Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and
708 Vikas Singh. Nyströmformer: A nyström-based algorithm for approximating self-attention. In
709 *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 14138–14148, 2021.
710
- 711 Jinsung Yoon, William R Zame, and Mihaela Van Der Schaar. Estimating missing data in temporal
712 data streams using multi-directional recurrent neural networks. *IEEE Transactions on Biomedical
Engineering*, 66(5):1477–1490, 2018.
- 713
- 714 Annan Yu, Arnur Nigmatov, Dmitriy Morozov, Michael W Mahoney, and N Benjamin Erichson. Ro-
715 bustifying state-space models for long sequences via approximate diagonalization. *arXiv preprint
arXiv:2310.01698*, 2023.
- 716
- 717 Annan Yu, Michael W Mahoney, and N Benjamin Erichson. Hope for a robust parameterization of
718 long-memory state space models. *arXiv preprint arXiv:2405.13975*, 2024.
- 719
- 720 Michael Yuanjie Zhang, Jeffrey R Russell, and Ruey S Tsay. A nonlinear autoregressive conditional
721 duration model with applications to financial transaction data. *Journal of Econometrics*, 104(1):
722 179–207, 2001.
- 723
- 724 Xiyuan Zhang and Danielle Maddix Robinson. Mitra: Mixed synthetic priors for enhanc-
725 ing tabular foundation models, 2025. URL [https://www.amazon.science/blog/
mitra-mixed-synthetic-priors-for-enhancing-tabular-foundation-models](https://www.amazon.science/blog/mitra-mixed-synthetic-priors-for-enhancing-tabular-foundation-models).
- 726
- 727 Yuanzhao Zhang and William Gilpin. Zero-shot forecasting of chaotic systems. *International Con-
ference on Learning Representations*, 2025.
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

756 **A RELATED WORK**

757 **A.1 TIME-SERIES MODELING**

760 Time series data are ubiquitous in many domains, including scientific (Abhishek et al., 2012; Zhang
 761 & Gilpin, 2025; Lai et al., 2025), industrial (Hong et al., 2016), and financial applications (Zhang
 762 et al., 2001), where they facilitate critical tasks, such as forecasting (Hyndman & Athanasopoulos,
 763 2018), imputation (Yoon et al., 2018), and anomaly detection (Blázquez-García et al., 2021).
 764 Classical time-series forecasting (Hyndman & Athanasopoulos, 2018) has deep roots in traditional
 765 methods such as ARIMA (Makridakis & Hibon, 1997) and exponential smoothing (Gardner Jr,
 766 1985). Many modern approaches involve deep learning; for example, sequence models such as
 767 DeepAR (Salinas et al., 2020) popularized probabilistic forecasting at scale, while Lim et al. (2021)
 768 combined the attention mechanism with interpretability for multi-horizon tasks. There are other neu-
 769 ral forecasters without attention, such as N-BEATS (Oreshkin et al., 2019) and N-HiTS (Challu et al.,
 770 2023). An evolving modern alternative that targets long contexts efficiently is the recently developed
 771 state-space models (Gu et al., 2021) and Mamba (Gu & Dao, 2023). These neural-network-based
 772 methods are usually considered task-specific — that is, given a task, one trains a model to obtain a
 773 specific set of weights tailored to that task.

774 **A.2 TIME-SERIES FOUNDATION MODELS**

776 Recent work pretrains general-purpose time series models across domains and tasks. Examples
 777 include TimesFM (Das et al., 2024), Chronos (Ansari et al., 2024a), Moirai (Woo et al., 2024), and
 778 Time MOE (Shi et al., 2025). These models differ in tokenization choices, training corpora, and
 779 zero-shot protocols, but they share the goal of one model that transfers across datasets and horizons.

780 **A.3 TIME-SERIES TOKENIZATION AND EMBEDDING**

783 Designing the input representation is the key to leveraging the high flexibility and expressive power
 784 of large Transformers. Patching turns small motifs into tokens, as in Nie et al. (2022) and Das
 785 et al. (2024). Discrete tokenization via quantization has been explored by Ansari et al. (2024a)
 786 and Talukder et al. (2024), and other discrete designs include HDT (Feng et al., 2025) and vector-
 787 quantized methods (Gui et al., 2024). Another line of research is by using frequency-based tokeniz-
 788 ers, such as WaveToken (Masserano et al., 2025). Lately, traditional language embedding strategies
 789 such as Byte-Pair Encoding are also considered (Götz et al., 2025).

790 **A.4 LOW-RANK STRUCTURES IN DEEP LEARNING**

792 Low-rank structure shows up in the deep learning community as both an inductive bias (Martin
 793 & Mahoney, 2021; Wilson, 2025; Yu et al., 2024) and a compression tool (Sharma et al., 2023;
 794 Gu et al., 2022; Yu et al., 2023). There is also a line of research that looks into parameter-efficient
 795 finetuning called LoRA (Hu et al., 2022), which learns low-rank updates to weight matrices. We note
 796 that LoRA is different from our low-rank model compression strategy because LoRA essentially
 797 learns a high-rank-plus-low-rank expression of the weight matrices, and does not facilitate storage
 798 or inference of a model. Earlier work reduced cost via matrix or tensor factorizations, especially in
 799 the CNN community (Denton et al., 2014; Jaderberg et al., 2014; Lebedev et al., 2014). These ideas
 800 motivate studying when low-rank operators suffice for sequence models, while our paper provides
 801 clean justifications for why low-rank operators suffice for a time-series foundation model.

802 **A.5 LOW-RANK STRUCTURES IN TRANSFORMERS**

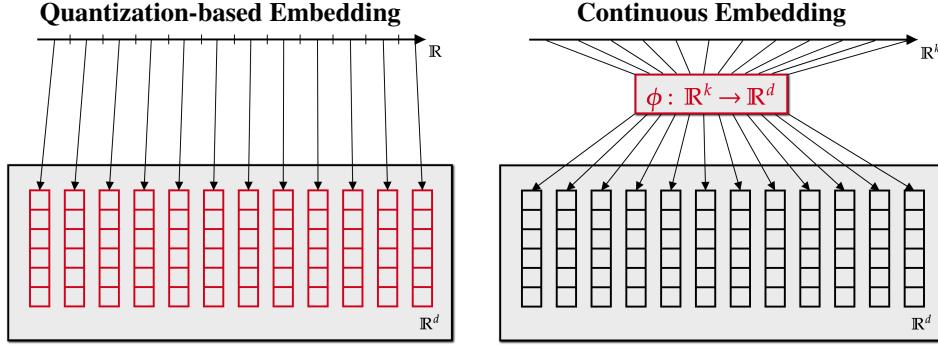
804 Low-rank structure in Transformers and attention has been extensively exploited for both effi-
 805 ciency and interpretability. A first line of work directly approximates the quadratic self-attention
 806 matrix with low-rank kernels: Linformer (Wang et al., 2020) projects keys and values to a low-
 807 dimensional subspace, yielding self-attention with linear complexity in sequence length, while
 808 Nyströmformer (Xiong et al., 2021) uses a Nyström approximation of the softmax kernel based
 809 on a small set of landmark tokens. Random-feature methods such as Performer (Choromanski et al.,
 2021) further view softmax attention as a kernel and approximate it with low-rank random feature

maps, obtaining linear-time and memory. Scatterbrain (Chen et al., 2021) analyzes when sparse versus low-rank attention yields better approximations and proposes a unified sparse+low-rank estimator that attains lower error than either component alone. Complementary approaches parameterize the weight matrices themselves in low rank: LoRA (Hu et al., 2022) learns low-rank updates on top of frozen full-rank weights for parameter-efficient fine-tuning of large language models, and Sharma et al. (2023) perform layer-selective rank reduction to probe and improve the reasoning behavior of Transformers. These works treat low rank primarily as an architectural or algorithmic assumption used to accelerate or regularize attention. By contrast, our analysis starts from the observed low numerical rank of time-series embeddings, proves when attention matrices over such embeddings are provably compressible, and then uses this data-driven perspective to design TSFMs whose attention layers are low rank by construction and compressible, which accelerates both pretraining and inference.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864 **B PROOF OF THEOREM 1 AND COROLLARY 2**
 865

866
 867 Before we go into the proofs of Theorem 1 and Corollary 2, we present a figure to illustrate the differ-
 868 ences between the quantization-based embedding and continuous embedding discussed in Section 2.



881 **Figure 8:** A comparison of the quantization-based embeddings and continuous embeddings. In a
 882 quantization-based one, we discretize the input domain into a number of regions and map each of
 883 them to a trainable vector. In a continuous embedding, we use a trainable function (usually an MLP)
 884 to map each element into the hidden space \mathbb{R}^d . The red parts in the figure highlight the trainable
 885 components.
 886

887 The proof of Theorem 1 relies on classic polynomial approximation results. Intuitively, if a function
 888 ϕ is well-approximated by a low-degree polynomial, then the low-degree polynomial can be repre-
 889 sented with a small basis, spanning a low-dimensional subspace. We now formalize this notion by
 890 providing the proof.

891
 892
 893 *Proof of Theorem 1.* Fix some $1 \leq j \leq d$. Let $P_{i,j}$ be the best degree- $(j-1)$ polynomial approxi-
 894 mation of ϕ_i in the infinity norm, and let
 895

$$\delta_j = \max_{1 \leq i \leq d} \|\phi_i - P_{i,j}\|_{L^\infty([-1,1])}.$$

896 Since $P_{i,j}$ is a polynomial of degree $\leq j-1$, we can write it into
 897

$$P_{i,j}(x) = \sum_{k=0}^{j-1} a_{i,k} T_k(x),$$

901 where T_k is the degree- k Chebyshev polynomial. Then, we can write Ξ into
 902

$$\Xi = \underbrace{\begin{bmatrix} -P_{1,j}(x) \\ \vdots \\ -P_{d,j}(x) \end{bmatrix}}_{\mathbf{E}_j} + \underbrace{\begin{bmatrix} -(\phi_1 - P_{1,j})(x) \\ \vdots \\ -(\phi_d - P_{d,j})(x) \end{bmatrix}}_{\Xi_j} = \underbrace{\begin{bmatrix} a_{1,0} & \cdots & a_{1,j-1} \\ \vdots & \ddots & \vdots \\ a_{d,0} & \cdots & a_{d,j-1} \end{bmatrix}}_{\mathbf{E}_j} \underbrace{\begin{bmatrix} -T_0(x) \\ \vdots \\ -T_{j-1}(x) \end{bmatrix}}_{\Xi_j} + \mathbf{E}_j.$$

913 Since Ξ_j is a quasimatrix of rank at most j , by (Townsend & Trefethen, 2015, Thm. 4.2), we have
 914 that

$$s_{j+1} \leq \|\mathbf{E}_j\|_F = \sqrt{\sum_{i=1}^d \|\phi_i - P_{i,j}\|_{L^2([-1,1])}^2} \leq 2 \sqrt{\sum_{i=1}^d \|\phi_i - P_{i,j}\|_{L^\infty([-1,1])}^2} \leq 2\sqrt{d} \delta_j. \quad (7)$$

918 Similarly, note that
 919

$$\begin{aligned}
 920 \quad \Psi &= \begin{bmatrix} P_{1,j}(x_1) & \cdots & P_{1,j}(x_L) \\ \vdots & \ddots & \vdots \\ P_{d,j}(x_1) & \cdots & P_{d,j}(x_L) \end{bmatrix} + \underbrace{\begin{bmatrix} (\phi_1 - P_{1,j})(x_1) & \cdots & (\phi_1 - P_{1,j})(x_L) \\ \vdots & \ddots & \vdots \\ (\phi_d - P_{d,j})(x_1) & \cdots & (\phi_d - P_{d,j})(x_L) \end{bmatrix}}_{\mathbf{F}_j} \\
 924 \\
 925 \quad &= \underbrace{\begin{bmatrix} a_{1,0} & \cdots & a_{1,j-1} \\ \vdots & \ddots & \vdots \\ a_{d,0} & \cdots & a_{d,j-1} \end{bmatrix}}_{\Psi_j} \begin{bmatrix} T_0(x_1) & \cdots & T_0(x_L) \\ \vdots & \ddots & \vdots \\ T_{j-1}(x_1) & \cdots & T_{j-1}(x_1) \end{bmatrix} + \mathbf{F}_j.
 \end{aligned}$$

930 Since the rank of Ψ_j is at most j , by the Eckart–Young inequality, we have that
 931

$$\sigma_{j+1} \leq \|\mathbf{F}_j\|_2 \leq \|\mathbf{F}_j\|_F \leq \sqrt{dL}\delta_j. \quad (8)$$

932 From (Trefethen, 2019, Thm. 7.2 & 8.2), we have that
 933

$$\begin{aligned}
 935 \quad \delta_j &\leq \frac{2V}{\pi\nu(j-1-\nu)^\nu} \quad \text{and} \quad \delta_j \leq \frac{2M\rho^{-j+1}}{\rho-1} \\
 937 \\
 938 \quad \text{when } \phi_i \text{ satisfies the condition in the first and the second statement of the theorem, respectively.} \\
 939 \quad \text{Hence, the two statement are proved by combining eq. (7), (8), and (9).} \quad \square
 \end{aligned} \quad (9)$$

940 The activation function swish_β used in Time MOE is analytic, and its domain of analyticity increases
 941 as $\beta \rightarrow 0^+$. Hence, we can use Theorem 1 to prove that Time MOE’s embedding is low-rank.
 942

943 **Corollary 2.** The embedding used in Time-MoE, defined by $\phi_i(x) = \text{swish}_\beta(w_i x) \cdot (v_i x) =$
 944 $(w_i v_i x^2)/(1 + e^{-\beta x})$, where we assume $|w_i v_i| \leq 1$ and $\beta > 0$, satisfies that

$$\sigma_{j+1} = \mathcal{O}\left(\sqrt{dL}(\beta + \beta^{-1})(1 + \pi/(2\beta))^{-j+1}\right),$$

945 where σ_{j+1} is defined in Theorem 1 for any Ψ and the constant in the \mathcal{O} -notation is universal.
 946

947 *Proof.* The function $\phi_i(z)$ is a meromorphic function with poles at $\{z \mid \exp(-\beta z) = -1\} = \{(2k+1)\pi i/\beta \mid k \in \mathbb{Z}\}$. Set $\rho = 1 + \pi/(2\beta)$. Then, within the Bernstein ellipse E_ρ of radius ρ , we have
 948 that

$$\text{Re}(e^{-\beta x}) > -c \Rightarrow |\phi_i(x)| \leq \frac{|x|^2}{1-c} = \mathcal{O}(1 + \beta^{-2}), \quad x \in E_\rho,$$

949 where $0 < c < 1$ is a universal constant and the constant in the \mathcal{O} -notation is also universal. The
 950 corollary follows from Theorem 1. \square
 951

952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

972 **C PROOF OF THEOREM 2**
 973

974 The intuition in Theorem 2 is that if the MLP does not have a nonlinear activation function, then
 975 the linear MLP clearly maps a low-rank matrix to a low-rank matrix. While a nonlinear activation
 976 function complicates things, a contractive function does not increase the Frobenius norm of a matrix,
 977 which is equivalent to the ℓ^2 -norm of all singular values.
 978

979 *Proof of Theorem 2.* Since \mathbf{X} is a rank- k matrix, the rank of $\mathbf{W}_1 \mathbf{X}$ is at most k . Hence, the Frobe-
 980 nius norm of the matrix satisfies that

$$981 \quad \| \mathbf{W}_1 \mathbf{X} \|_F^2 = \sum_{j=1}^k \sigma_j(\mathbf{W}_1 \mathbf{X})^2 \leq k \sigma_1(\mathbf{W}_1 \mathbf{X})^2.$$

982 Given that $|\omega(x)| \leq |x|$ for every $x \in \mathbb{R}$, we have

$$983 \quad k \sigma_1(\mathbf{W}_1 \mathbf{X})^2 \geq \| \mathbf{W}_1 \mathbf{X} \|_F^2 \geq \| \omega(\mathbf{W}_1 \mathbf{X}) \|_F^2 = \sum_{j=1}^{\min(d_f, L)} \sigma_j(\omega(\mathbf{W}_1 \mathbf{X}))^2.$$

984 Using the singular value inequalities, we have that

$$985 \quad \begin{aligned} \sum_{j=1}^{\min(d_f, L)} \sigma_j(\mathbf{W}_2 \omega(\mathbf{W}_1 \mathbf{X}))^2 &\leq \sigma_1(\mathbf{W}_2)^2 \sum_{j=1}^{\min(d_f, L)} \sigma_j(\omega(\mathbf{W}_1 \mathbf{X}))^2 \\ &\leq k \sigma_1(\mathbf{W}_2)^2 \sigma_1(\mathbf{W}_1 \mathbf{X})^2 = k \| \mathbf{W}_2 \|_2^2 \| \mathbf{W}_1 \mathbf{X} \|_2^2. \end{aligned}$$

986 This is, we can control the number of large singular values of $\mathbf{W}_2 \omega(\mathbf{W}_1 \mathbf{X})$ by

$$987 \quad |\{j \mid \sigma_j(\mathbf{W}_2 \omega(\mathbf{W}_1 \mathbf{X})) > \varepsilon \| \mathbf{W}_2 \|_2 \| \mathbf{W}_1 \mathbf{X} \|_2\}| \leq \varepsilon^{-2} k.$$

988 Moreover, the rank of the matrix $\mathbf{W}_3 \mathbf{X}$ is at most k . Using the singular value inequality that

$$989 \quad \sigma_{i+j-1}(\mathbf{A} + \mathbf{B}) \leq \sigma_i(\mathbf{A}) + \sigma_j(\mathbf{B}), \quad \mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times L}, \quad i + j - 1 \leq \min(d, L),$$

990 we have

$$991 \quad \sigma_{k+j}(\phi(\mathbf{X})) \leq \sigma_{k+1}(\mathbf{W}_3 \mathbf{X}) + \sigma_j(\mathbf{W}_2 \omega(\mathbf{W}_1 \mathbf{X})) = \sigma_j(\mathbf{W}_2 \omega(\mathbf{W}_1 \mathbf{X})), \quad j \leq \min(d, L) - k.$$

992 This gives us

$$993 \quad \begin{aligned} |\{j \mid \sigma_j(\phi(\mathbf{X})) > \varepsilon \| \mathbf{W}_2 \|_2 \| \mathbf{W}_1 \mathbf{X} \|_2\}| &\leq |\{j \mid \sigma_j(\mathbf{W}_2 \omega(\mathbf{W}_1 \mathbf{X})) > \varepsilon \| \mathbf{W}_2 \|_2 \| \mathbf{W}_1 \mathbf{X} \|_2\}| + k \\ &\leq (1 + \varepsilon^{-2}) k, \end{aligned}$$

994 which proves the result. □

1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

1026 **D PROOF OF THEOREM 3**
 1027

1028 To prove the upper bound in Theorem 3, we will prove a stronger result concerning not only a
 1029 single-head attention but also a multi-head one.

1030 **Theorem 5.** Let $C > 0$ be a constant. Let $\Xi = [\mathbf{x}_1 \cdots \mathbf{x}_N] \in \mathbb{R}^{d \times N}$ be an embedding matrix,
 1031 where d is the hidden dimension, N is the vocabulary size, $\mathbf{x}_j \in \mathbb{R}^d$, and $\|\mathbf{x}_j\|_2 \leq C$ for all
 1032 $1 \leq j \leq N$. Let $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$ be multi-head attention matrices with h heads, such that
 1033 $\|(\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)}\|_2 \leq C\sqrt{d_h}$, and $\|\mathbf{W}_V^{(i)}\|_2 \leq C\sqrt{d_h}$ for every $1 \leq i \leq h$. For any $\tilde{d} < d$ such
 1034 that $\sigma_{\tilde{d}+1} := \sigma_{\tilde{d}+1}(\Xi) \leq 1$, there exists a stable low-rank approximation $\tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V \in \mathbb{R}^{d \times d}$
 1035 with $\text{rank}(\tilde{\mathbf{W}}_Q) = \text{rank}(\tilde{\mathbf{W}}_K) = \text{rank}(\tilde{\mathbf{W}}_V) = \tilde{d}$, $\|\tilde{\mathbf{W}}_Q^\top \tilde{\mathbf{W}}_K\|_2 \leq \|\mathbf{W}_Q^\top \mathbf{W}_K\|_2$, and $\|\tilde{\mathbf{W}}_V\|_2 \leq$
 1036 $\|\mathbf{W}_V\|_2$, such that given any input matrix $\mathbf{U} \in \mathbb{R}^{d \times L}$ for any $L \geq 1$, where each column of \mathbf{U} is a
 1037 column of Ξ , we have that the low-rank attention matrices uniformly approximate the original one:
 1038

$$1040 \quad \left\| \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) - \text{MH-Attention}(\mathbf{U}; \tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V) \right\|_F \leq \mathcal{O}(\sqrt{d} \sigma_{\tilde{d}+1}), \quad (10)$$

1042 where the constant in the \mathcal{O} -notation only depends on C .
 1043

1044 *Proof.* Fix a $\tilde{d} < d$. For the sake of simplicity, we assume, without loss of generality, that $C = 1$.
 1045 Let

$$1046 \quad \Xi = \mathbf{U}_{\tilde{d}} \mathbf{S}_{\tilde{d}} \mathbf{V}_{\tilde{d}}^\top, \quad \mathbf{U}_{\tilde{d}} \in \mathbb{R}^{d \times \tilde{d}}, \quad \mathbf{S}_{\tilde{d}} \in \mathbb{R}^{\tilde{d} \times \tilde{d}}, \quad \mathbf{V}_{\tilde{d}} \in \mathbb{R}^{N \times \tilde{d}}$$

1047 be the truncated SVD of Ξ . Let $\mathbf{Q}_{\tilde{d}} = \mathbf{S}_{\tilde{d}} \mathbf{V}_{\tilde{d}}^\top \in \mathbb{R}^{\tilde{d} \times N}$. Therefore, we have

$$1050 \quad \|\Xi - \mathbf{U}_{\tilde{d}} \mathbf{Q}_{\tilde{d}}\|_2 \leq \sigma_{\tilde{d}+1}(\Xi).$$

1052 Define $\tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K$, and $\tilde{\mathbf{W}}_V$ by

$$1054 \quad \tilde{\mathbf{W}}_Q = \mathbf{W}_Q \mathbf{U}_{\tilde{d}} \mathbf{U}_{\tilde{d}}^\top \in \mathbb{R}^{d \times d}, \quad \tilde{\mathbf{W}}_K = \mathbf{W}_K \mathbf{U}_{\tilde{d}} \mathbf{U}_{\tilde{d}}^\top \in \mathbb{R}^{d \times d}, \quad \tilde{\mathbf{W}}_V = \mathbf{W}_V \mathbf{U}_{\tilde{d}} \mathbf{U}_{\tilde{d}}^\top \in \mathbb{R}^{d \times d}.$$

1055 Since $\mathbf{U}_{\tilde{d}}$ has orthonormal columns, we have that $\|\mathbf{U}_{\tilde{d}}\|_2 = \|\mathbf{U}_{\tilde{d}}^\top\|_2 = 1$. By the sub-multiplicity of
 1056 the spectral norm, we have that

$$1058 \quad \|\tilde{\mathbf{W}}_Q^\top \tilde{\mathbf{W}}_K\|_2 \leq \|\mathbf{W}_Q^\top \mathbf{W}_K\|_2, \quad \|\tilde{\mathbf{W}}_V\|_2 \leq \|\mathbf{W}_V\|_2,$$

1060 which proves the stability of the low-rank representation. Let $\mathbf{Q}_U \in \mathbb{R}^{\tilde{d} \times L}$ be the matrix defined by
 1061 the condition that the i th column of \mathbf{Q}_U is the j th column of $\mathbf{Q}_{\tilde{d}}$ if and only if the i th column of \mathbf{U}
 1062 is the j th column of Ξ , and let

$$1063 \quad \Delta_U = \mathbf{U}_{\tilde{d}} \mathbf{Q}_U - \mathbf{U}.$$

1064 Since every column of Δ_U is a column of $\Xi - \mathbf{U}_{\tilde{d}} \mathbf{Q}_{\tilde{d}}$, its norm is no greater than $\sigma_{\tilde{d}+1}(\Xi)$, i.e.,
 1065

$$1066 \quad \begin{aligned} \|\Delta_U(:, j)\|_2 &\leq \sigma_{\tilde{d}+1}(\Xi), \quad 1 \leq j \leq L, \\ 1067 \quad \|\Delta_U\|_2 &\leq \|\Delta_U\|_F \leq \sqrt{L} \sigma_{\tilde{d}+1}(\Xi). \end{aligned} \quad (11)$$

1069 Next, we have

$$\begin{aligned} 1071 \quad \|\Delta_V^{(i)}\|_2 &\leq \|\Delta_U^{(i)}\|_2 \leq \sqrt{d_h L} \sigma_{\tilde{d}+1}(\Xi), \\ 1072 \quad \Delta_V^{(i)} &:= \mathbf{W}_V^{(i)} \mathbf{U} - \tilde{\mathbf{W}}_V^{(i)} \mathbf{U} = \mathbf{W}_V^{(i)} \mathbf{U} - \mathbf{W}_V^{(i)} \mathbf{U}_{\tilde{d}} \mathbf{U}_{\tilde{d}}^\top \mathbf{U} \\ 1073 &= \mathbf{W}_V^{(i)} \mathbf{U} - \mathbf{W}_V^{(i)} \mathbf{U}_{\tilde{d}} \mathbf{U}_{\tilde{d}}^\top \mathbf{U}_{\tilde{d}} \mathbf{Q}_U \\ 1074 &= \mathbf{W}_V^{(i)} \mathbf{U} - \mathbf{W}_V^{(i)} \mathbf{U}_{\tilde{d}} \mathbf{Q}_U = -\mathbf{W}_V^{(i)} \Delta_U, \\ 1075 & \end{aligned}$$

1077 where the first inequality comes from the sub-multiplicity of the spectral norm. Moreover, every
 1078 column of $\Delta_V^{(i)}$ is a column of $-\mathbf{W}_V^{(i)} \Delta_U$ and must satisfy that

$$1079 \quad \|\Delta_V^{(i)}(:, j)\|_2 \leq \sqrt{d_h} \sigma_{\tilde{d}+1}(\Xi), \quad 1 \leq j \leq L. \quad (12)$$

1080 Similarly, we have that
 1081

$$\begin{aligned}
 \| \Delta_Q^{(i)} \|_2 &\leq \sqrt{d_h} L (2\sigma_{\tilde{d}+1}(\Xi) + \sigma_{\tilde{d}+1}(\Xi)^2), \\
 \Delta_Q^{(i)} &:= \mathbf{U}^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \mathbf{U} - \mathbf{U}^\top (\tilde{\mathbf{W}}_Q^{(i)})^\top \tilde{\mathbf{W}}_K^{(i)} \mathbf{U} \\
 &= \mathbf{U}^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \mathbf{U} - \mathbf{U}^\top \mathbf{U}_{\tilde{d}} \mathbf{U}_{\tilde{d}}^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \mathbf{U}_{\tilde{d}} \mathbf{U}_{\tilde{d}}^\top \mathbf{U} \\
 &= \mathbf{U}^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \mathbf{U} - \mathbf{Q}_X^\top \mathbf{U}_{\tilde{d}}^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \mathbf{U}_{\tilde{d}} \mathbf{Q}_U \\
 &= -\Delta_U^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \mathbf{U} - \mathbf{U}^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \Delta_U - \Delta_U^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \Delta_U,
 \end{aligned}$$

1089 where the inequality is obtained by recalling that $\|\mathbf{U}\|_2 \leq \sqrt{L}$ and $\|\Delta_U\|_2 \leq \sqrt{L} \sigma_{\tilde{d}+1}(\Xi)$.
 1090 Moreover, since we have $\|\mathbf{U}(:, j)\|_2 \leq 1$ and $\|\Delta_U(:, j)\|_2 \leq \sigma_{\tilde{d}+1}(\Xi)$ for all $1 \leq j \leq N$, every
 1091 entry of $\Delta_Q^{(i)}$ satisfies that
 1092

$$\| \Delta_Q^{(i)}(t, j) \|_2 \leq \sqrt{d_h} (2\sigma_{\tilde{d}+1}(\Xi) + \sigma_{\tilde{d}+1}(\Xi)^2), \quad 1 \leq t \leq L, \quad 1 \leq j \leq L. \quad (13)$$

1093 For each fixed $1 \leq i \leq h$, define the notations
 1094

$$\begin{aligned}
 \mathbf{G}^{(i)} &= \frac{\mathbf{U}^\top (\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)} \mathbf{U}}{\sqrt{d_h}}, & \tilde{\mathbf{G}}^{(i)} &= \frac{\mathbf{U}^\top (\tilde{\mathbf{W}}_Q^{(i)})^\top \tilde{\mathbf{W}}_K^{(i)} \mathbf{U}}{\sqrt{d_h}}, \\
 \mathbf{g}_j^{(i)} &= \mathbf{G}(:, j), & \tilde{\mathbf{g}}_j^{(i)} &= \tilde{\mathbf{G}}(:, j).
 \end{aligned}$$

1100 For simplicity, we drop the superscripts (i) on \mathbf{G} , $\tilde{\mathbf{G}}$, \mathbf{g} , and $\tilde{\mathbf{g}}$. Since we assumed $\|\mathbf{u}_j\|_2 \leq 1$ for all
 1101 $1 \leq j \leq N$ and $\|(\mathbf{W}_Q^{(i)})^\top \mathbf{W}_K^{(i)}\|_2 \leq \sqrt{d_h}$, we have that
 1102

$$\| \mathbf{g}_j \|_{\max} \leq 1, \quad \| \tilde{\mathbf{g}}_j \|_{\max} \leq 1, \quad 1 \leq j \leq L.$$

1103 Denote by g_j^t and \tilde{g}_j^t the t th entry of \mathbf{g}_j and $\tilde{\mathbf{g}}_j$, respectively. For every fixed $1 \leq j \leq L$ and
 1104 $1 \leq t \leq L$, we have
 1105

$$\begin{aligned}
 & \left| \underbrace{\text{softmax}(\mathbf{g}_j)^t - \text{softmax}(\tilde{\mathbf{g}}_j)^t}_{d_j^t} \right| = \left| \frac{\exp(g_j^t)}{\sum_{k=1}^L \exp(g_j^k)} - \frac{\exp(\tilde{g}_j^t)}{\sum_{k=1}^L \exp(\tilde{g}_j^k)} \right| \\
 &= \frac{\left| \sum_{k=1}^L (\exp(g_j^t) \exp(\tilde{g}_j^k) - \exp(\tilde{g}_j^t) \exp(g_j^k)) \right|}{\left| \left(\sum_{k=1}^L \exp(g_j^k) \right) \left(\sum_{k=1}^L \exp(\tilde{g}_j^k) \right) \right|} \leq \frac{L \max_k |\exp(g_j^t) \exp(\tilde{g}_j^k) - \exp(\tilde{g}_j^t) \exp(g_j^k)|}{L^2 \exp(-2)} \\
 &= \frac{L \max_k |\exp(g_j^t)(\exp(g_j^k) + (\exp(\tilde{g}_j^k) - \exp(g_j^k))) - (\exp(g_j^t) + (\exp(\tilde{g}_j^t) - \exp(g_j^t))) \exp(g_j^k)|}{L^2 \exp(-2)} \\
 &\leq \frac{\max_k (\exp(1)(|\exp(g_j^k) - \exp(\tilde{g}_j^k)| + |\exp(g_j^t) - \exp(\tilde{g}_j^t)|) + |\exp(g_j^k) - \exp(\tilde{g}_j^k)| |\exp(g_j^t) - \exp(\tilde{g}_j^t)|)}{L \exp(-2)} \\
 &\leq \frac{2 \exp(1)(2\sigma_{\tilde{d}+1}(\Xi) + \sigma_{\tilde{d}+1}(\Xi)^2) + (2\sigma_{\tilde{d}+1}(\Xi) + \sigma_{\tilde{d}+1}(\Xi)^2)^2}{L \exp(-2)}, \\
 &\underbrace{\quad \quad \quad D_Q}_{(14)}
 \end{aligned}$$

1122 where the last inequality follows from eq. (13). Hence, for each fixed $1 \leq j \leq L$ and $1 \leq i \leq h$, we
 1123 have that
 1124

$$\begin{aligned}
 & \left\| \mathbf{W}_V^{(i)} \mathbf{U} \text{softmax}(\mathbf{g}_j^{(i)}) - \tilde{\mathbf{W}}_V^{(i)} \mathbf{U} \text{softmax}(\tilde{\mathbf{g}}_j^{(i)}) \right\|_2 \\
 &= \left\| \sum_{i=1}^L \left(\mathbf{W}_V^{(i)} \mathbf{U}(:, i) \text{softmax}(g_j^t) - \tilde{\mathbf{W}}_V^{(i)} \mathbf{U}(:, i) \text{softmax}(\tilde{g}_j^t) \right) \right\|_2 \\
 &\leq \sqrt{L} \max_{1 \leq t \leq L} \left\| \left(\mathbf{W}_V^{(i)} \mathbf{U}(:, t) \text{softmax}(g_j^t) - (\mathbf{W}_V^{(i)} \mathbf{U}(:, t) - \Delta_V^{(i)}(:, t)) (\text{softmax}(g_j^t) - d_j^t) \right) \right\|_2 \\
 &\leq \sqrt{L} \max_{1 \leq t \leq L} \left(\|\Delta_V^{(i)}(:, t)\|_2 |\text{softmax}(\tilde{g}_j^t)| + \|\mathbf{W}_V^{(i)} \mathbf{U}(:, t)\|_2 D_Q \right) \\
 &\leq \sqrt{L} \left(\sqrt{d_h} \sigma_{\tilde{d}+1} \Theta(1/L) + \sqrt{d_h} \mathcal{O} \left((\sigma_{\tilde{d}+1} + \sigma_{\tilde{d}+1}^2) / L \right) \right) = \mathcal{O} \left(\sqrt{d_h} \sigma_{\tilde{d}+1} / \sqrt{L} \right),
 \end{aligned}$$

1134 where the last inequality follows from eq. (12) and (14). Since $\mathbf{W}_V^{(i)} \mathbf{U} \text{ softmax}(\mathbf{g}_j^{(i)}) -$
 1135 $\tilde{\mathbf{W}}_V^{(i)} \mathbf{U} \text{ softmax}(\tilde{\mathbf{g}}_j^{(i)})$ is the j th column of
 1136

$$1137 \quad \text{Attention}(\mathbf{U}; \mathbf{W}_Q^{(i)}, \mathbf{W}_K^{(i)}, \mathbf{W}_V^{(i)}) - \text{Attention}(\mathbf{U}; \tilde{\mathbf{W}}_Q^{(i)}, \tilde{\mathbf{W}}_K^{(i)}, \tilde{\mathbf{W}}_V^{(i)}),$$

1139 we have that

$$\begin{aligned} 1141 \quad & \| \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) - \text{Attention}(\mathbf{U}; \tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V, h) \|_F \\ 1142 \quad &= \sqrt{\sum_{i=1}^h \| \text{Attention}(\mathbf{U}; \mathbf{W}_Q^{(i)}, \mathbf{W}_K^{(i)}, \mathbf{W}_V^{(i)}) - \text{Attention}(\mathbf{U}; \tilde{\mathbf{W}}_Q^{(i)}, \tilde{\mathbf{W}}_K^{(i)}, \tilde{\mathbf{W}}_V^{(i)}) \|_F^2} \\ 1143 \quad &\leq \sqrt{\sum_{i=1}^h \sum_{j=1}^L \| \mathbf{W}_V^{(i)} \mathbf{U} \text{ softmax}(\mathbf{g}_j^{(i)}) - \tilde{\mathbf{W}}_V^{(i)} \mathbf{U} \text{ softmax}(\tilde{\mathbf{g}}_j^{(i)}) \|_2^2} \\ 1144 \quad &= \sqrt{hL} \mathcal{O}\left(\sqrt{d_h} \sigma_{\tilde{d}+1} / \sqrt{L}\right) = \mathcal{O}\left(\sqrt{d} \sigma_{\tilde{d}+1}\right). \end{aligned}$$

1150 The proof is complete. \square

1153 Theorem 5 immediately proves the upper bound in Theorem 3. The lower bound needs a separate
 1154 argument.

1155 *Proof of Theorem 3.* The upper bound follows immediately from Theorem 5 by setting $h = 1$. To
 1156 prove the lower bound, let $\mathbf{U}' \in \mathbb{R}^{d \times \tilde{d}}$ and $\mathbf{U} \in \mathbb{R}^{d \times L}$ be such that

$$1159 \quad \mathbf{U} = [\mathbf{U}' \mid \mathbf{0}], \quad \mathbf{U}' = \text{diag}(\sigma_1, \dots, \sigma_d) \in \mathbb{R}^{d \times d}, \quad \mathbf{0} \in \mathbb{R}^{d \times (L-d)}.$$

1160 Clearly, the singular values of \mathbf{U} are $\sigma_1, \dots, \sigma_d$. Define \mathbf{W}_Q and \mathbf{W}_K such that

$$1162 \quad \mathbf{W}_Q^\top \mathbf{W}_K = \log(4d) \sigma_d^{-2} \sqrt{d} \mathbf{I}_d.$$

1163 Then, we have

$$1165 \quad \mathbf{T} := \frac{\mathbf{U}^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{U}}{\sqrt{d}} = \log(4d) \sigma_d^{-2} \left[\begin{array}{c|c} \text{diag}(\sigma_1^2, \dots, \sigma_d^2) & \mathbf{0}_{d \times (L-d)} \\ \hline \mathbf{0}_{(L-d) \times d} & \mathbf{0}_{(L-d) \times (L-d)} \end{array} \right].$$

1168 Let $\mathbf{G} = \text{softmax}(\mathbf{T})$ and let \mathbf{g}_j be the j th column of \mathbf{G} . Then, for every $1 \leq j \leq d$, we have

$$1170 \quad \mathbf{g}_j = [g_j^1, \dots, g_j^L]^\top, \quad g_j^i = \begin{cases} \frac{1}{(L-1) + \exp(\log(4d) \sigma_d^{-2} \sigma_j^2)} \leq \frac{1}{4d}, & i \neq j, \\ \frac{\exp(\log(4d) \sigma_d^{-2} \sigma_j^2)}{(L-1) + \exp(\log(4d) \sigma_d^{-2} \sigma_j^2)} \geq \frac{1}{2}, & i = j. \end{cases} \quad (15)$$

1175 Write \mathbf{G} into

$$1176 \quad \mathbf{G} = \left[\begin{array}{c|c} \mathbf{G}_{11} & \mathbf{G}_{12} \\ \hline \mathbf{G}_{21} & \mathbf{G}_{22} \end{array} \right],$$

1178 where

$$1180 \quad \mathbf{G}_{11} = \mathbf{G}' \in \mathbb{R}^{d \times d}, \quad \mathbf{G}_{12} \in \mathbb{R}^{d \times (L-d)}, \quad \mathbf{G}_{21} \in \mathbb{R}^{(L-d) \times d}, \quad \mathbf{G}_{22} \in \mathbb{R}^{(L-d) \times (L-d)},$$

1181 and write \mathbf{G}_{11} into the sum of its diagonal part and off-diagonal part:

$$1183 \quad \mathbf{G}_{11} = \mathbf{G}_{\text{diag}} + \mathbf{G}_{\text{off-diag}},$$

1184 where \mathbf{G}_{diag} is a diagonal matrix and $\mathbf{G}_{\text{off-diag}}$ is a matrix with zero diagonal entries. Then,
 1185 by eq. (15), we have

$$1187 \quad \sigma_d(\mathbf{G}_{\text{diag}}) = \min_{1 \leq j \leq d} g_j^j \geq \frac{1}{2}, \quad \|\mathbf{G}_{\text{off-diag}}\|_2 \leq \|\mathbf{G}_{\text{off-diag}}\|_F \leq d \frac{1}{4d} = \frac{1}{4}.$$

1188 Hence, by Weyl's inequality, we have
 1189

$$1190 \sigma_d(\mathbf{G}_{11}) \geq \sigma_d(\mathbf{G}_{\text{diag}}) - \|\mathbf{G}_{\text{off-diag}}\|_2 \geq \frac{1}{4}. \quad (16)$$

1192 Using the results above, we have
 1193

$$\begin{aligned} 1194 \mathbf{W}_V \mathbf{U} \text{ softmax}\left(\frac{\mathbf{U}^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{U}}{\sqrt{d}}\right) - \tilde{\mathbf{W}}_V \mathbf{U} \text{ softmax}\left(\frac{\mathbf{U}^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{U}}{\sqrt{d}}\right) \\ 1195 \\ 1196 \\ 1197 = (\mathbf{W}_V \mathbf{U} - \tilde{\mathbf{W}}_V \mathbf{U}) \mathbf{G} = \left[\begin{array}{c|c} \mathbf{W}_V \mathbf{U}' - \tilde{\mathbf{W}}_V \mathbf{U}' & \mathbf{0}_{d \times (L-d)} \end{array} \right] \left[\begin{array}{c|c} \mathbf{G}_{11} & \mathbf{G}_{12} \\ \hline \mathbf{G}_{21} & \mathbf{G}_{22} \end{array} \right] \\ 1198 \\ 1199 = \left[\begin{array}{c|c} (\mathbf{W}_V \mathbf{U}' - \tilde{\mathbf{W}}_V \mathbf{U}') \mathbf{G}_{11} & (\mathbf{W}_V \mathbf{U}' - \tilde{\mathbf{W}}_V \mathbf{U}') \mathbf{G}_{12} \end{array} \right], \\ 1200 \end{aligned}$$

1201 but $\tilde{\mathbf{W}}_V \mathbf{U}'$ is a rank- \tilde{d} matrix so we must have
 1202

$$1203 \|\mathbf{W}_V \mathbf{U}' - \tilde{\mathbf{W}}_V \mathbf{U}'\|_2 \geq \sigma_{\tilde{d}+1}(\mathbf{W}_V \mathbf{U}') = \sigma_{\tilde{d}+1}. \quad (17)$$

1204 Hence, by the singular value inequalities, we have
 1205

$$\begin{aligned} 1206 \left\| \mathbf{W}_V \mathbf{U} \text{ softmax}\left(\frac{\mathbf{U}^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{U}}{\sqrt{d}}\right) - \tilde{\mathbf{W}}_V \mathbf{U} \text{ softmax}\left(\frac{\mathbf{U}^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{U}}{\sqrt{d}}\right) \right\|_2 \\ 1207 \\ 1208 \\ 1209 \geq \|(\mathbf{W}_V \mathbf{U}' - \tilde{\mathbf{W}}_V \mathbf{U}') \mathbf{G}_{11}\|_2 \geq \sigma_1(\mathbf{W}_V \mathbf{U}' - \tilde{\mathbf{W}}_V \mathbf{U}') \sigma_d(\mathbf{G}_{11}) \geq \frac{1}{4} \sigma_{\tilde{d}+1}, \\ 1210 \end{aligned}$$

1211 where the last inequality is obtained by combining eq. (16) and (17). The proof is complete. \square
 1212

1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

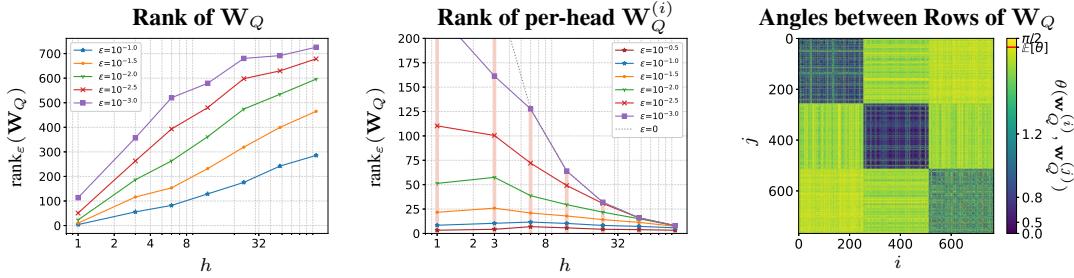
1242 **E COMPRESSING A MULTI-HEAD ATTENTION LAYER**
 1243

1244 Theorem 3 concerns the compressibility of a single-head attention layer. In practice, most TSFMs
 1245 use multi-head attention instead:²

$$\text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) = \begin{bmatrix} \text{Attention}(\mathbf{U}; \mathbf{W}_Q^{(1)}, \mathbf{W}_K^{(1)}, \mathbf{W}_V^{(1)}) \\ \vdots \\ \text{Attention}(\mathbf{U}; \mathbf{W}_Q^{(h)}, \mathbf{W}_K^{(h)}, \mathbf{W}_V^{(h)}) \end{bmatrix}, \quad (18)$$

$$\mathbf{W}_{Q/K/V} = \left[(\mathbf{W}_{Q/K/V}^{(1)})^\top \cdots (\mathbf{W}_{Q/K/V}^{(h)})^\top \right]^\top, \quad \mathbf{W}_{Q/K/V}^{(i)} \in \mathbb{R}^{d_h \times d}, \quad d_h = d/h.$$

1250 This leads to the following question: does the number of heads have an effect on the numerical ranks
 1251 of the attention matrices in trained Chronos models? The answer to this is positive, but there are
 1252 some important subtleties. If we look at the left panel of Figure 9, then we see that if we pretrain
 1253 Chronos models with fixed hidden dimension $d = 1024$ and only increase the number of heads,
 1254 the numerical rank of attention matrices increases. At first, it may seem that this happens because
 1255 a multi-head attention is less compressible than a single-head one, but this is not the case. It is
 1256 straightforward to show that Theorem 3 holds as well for multi-head attentions (see Theorem 5
 1257 in Appendix D). In particular, if Ξ has an algebraic rank of \tilde{d} , then all multi-head attention matrices
 1258 can be compressed to rank- \tilde{d} without affecting the output.



1270 **Figure 9:** The left panel shows the averaged ε -rank of query projection matrices \mathbf{W}_Q in pretrained
 1271 Chronos models with varying number of heads. The middle panel shows the averaged ε -rank of
 1272 query projection submatrices $\mathbf{W}_Q^{(i)}$ for every head. The right panel shows the angle between every
 1273 pair of rows of \mathbf{W}_Q in the first layer of a 3-head pretrained Chronos model.
 1274

1275 If Theorem 5 holds for multi-head attention, then why does a multi-head attention exhibit a higher
 1276 numerical rank? To understand this, we need to understand the mechanism of a low-rank weight
 1277 matrix. We use \mathbf{W}_Q for illustration, and the same analogous applies to \mathbf{W}_K and \mathbf{W}_V . Let $\mathbf{W}_Q \in \mathbb{R}^{d \times d}$ be approximated by a rank- \tilde{d} matrix, i.e., $\mathbf{W}_Q \approx \mathbf{W}_1 \mathbf{W}_2$, where $\mathbf{W}_1 \in \mathbb{R}^{d \times \tilde{d}}$ and $\mathbf{W}_2 \in \mathbb{R}^{\tilde{d} \times d}$. We view the action $\mathbf{W}_Q \mathbf{U}$, which maps every column of \mathbf{U} from \mathbb{R}^d into \mathbb{R}^d , as a two-step process. First, we multiply \mathbf{W}_2 to \mathbf{U} to drop the columns of \mathbf{U} from \mathbb{R}^d to $\mathbb{R}^{\tilde{d}}$. In numerical linear algebra, this is known as ‘‘sketching,’’ (Drineas & Mahoney, 2016) i.e., \mathbf{W}_2 ‘‘sketches’’ a \tilde{d} -dimensional subspace $R(\mathbf{W}_2 \mathbf{U})$ in \mathbb{R}^L , the row space of $\mathbf{W}_2 \mathbf{U}$, from a d -dimensional subspace $R(\mathbf{U})$ in \mathbb{R}^L .³ If $\sigma_{\tilde{d}+1}(\mathbf{U})$ is small, then there exists a sketching matrix \mathbf{W}_2 so that $R(\mathbf{W}_2 \mathbf{U}) \approx R(\mathbf{U}) \supset R(\mathbf{W}_Q \mathbf{U})$; hence, we can apply a matrix \mathbf{W}_1 to lift the columns of $\mathbf{W}_2 \mathbf{U}$ from $\mathbb{R}^{\tilde{d}}$ back to \mathbb{R}^d and have $\mathbf{W}_1(\mathbf{W}_2 \mathbf{U}) \approx \mathbf{W}_Q \mathbf{U}$.

1287 If we apply a low-rank approximation to a multi-head attention matrix \mathbf{W}_Q , then we obtain:

$$\begin{bmatrix} \mathbf{W}_Q^{(1)} \\ \vdots \\ \mathbf{W}_Q^{(h)} \end{bmatrix} = \mathbf{W}_Q \approx \mathbf{W}_1 \mathbf{W}_2 = \begin{bmatrix} \mathbf{W}_1^{(1)} \\ \vdots \\ \mathbf{W}_1^{(h)} \end{bmatrix} \mathbf{W}_2 = \begin{bmatrix} \mathbf{W}_1^{(1)} \mathbf{W}_2 \\ \vdots \\ \mathbf{W}_1^{(h)} \mathbf{W}_2 \end{bmatrix}, \quad \mathbf{W}_1^{(i)} \in \mathbb{R}^{d_h \times \tilde{d}}. \quad (19)$$

1293 ²In eq. (18), the $1/\sqrt{d}$ normalization factor in an attention is changed into $1/\sqrt{d_h}$.

1294 ³In numerical linear algebra (NLA), sketching is usually considered as an operation on the column space.
 1295 To adapt it to our framework, we sketch the row space instead. This is achieved by taking the transpose of the
 matrix \mathbf{U} and sketching its column space as usually done in NLA.

When we apply \mathbf{W}_Q to \mathbf{U} , for each head $1 \leq i \leq h$, \mathbf{W}_2 sketches a \tilde{d} -dimensional subspace from the row space of \mathbf{U} while $\mathbf{W}_1^{(i)}$ forms d_h rows from this \tilde{d} -dimensional space. This makes the problem clear: in an “optimal” low-rank approximation of \mathbf{W}_Q , the sketching matrix \mathbf{W}_2 is shared across all heads. Since each head is independent from the others in a multi-head attention, there is no guarantee that, in a pretrained model, the sketching matrices will be shared. In the right panel of Figure 9, we visualize the angles (see eq. (2)) between every pair of rows of \mathbf{W}_Q . If \mathbf{W}_2 is shared across all heads, then all rows of \mathbf{W}_Q should exhibit some linear dependency. We only see low-rank structures within each head $\mathbf{W}_Q^{(i)}$ (i.e., the dark $d_h \times d_h$ blocks along the diagonal). This shows that instead of eq. (19), \mathbf{W}_Q from a pretrained Chronos model looks more like

$$\begin{bmatrix} \mathbf{W}_Q^{(1)} \\ \vdots \\ \mathbf{W}_Q^{(h)} \end{bmatrix} = \mathbf{W}_Q \approx \begin{bmatrix} \mathbf{W}_1^{(1)} \mathbf{W}_2^{(1)} \\ \vdots \\ \mathbf{W}_1^{(h)} \mathbf{W}_2^{(h)} \end{bmatrix}, \quad \mathbf{W}_1^{(i)} \in \mathbb{R}^{d_h \times \tilde{d}}, \quad \mathbf{W}_2^{(i)} \in \mathbb{R}^{\tilde{d} \times d}, \quad (20)$$

which leverages a head-dependent sketching $\mathbf{W}_2^{(i)}$. The rank of the right-hand side matrix in eq. (20), which consists of h rank- \tilde{d} submatrices, can be as large as $hd\tilde{d}$. This finding explains the phenomenon we observed at the beginning of this section, i.e., while the numerical rank \tilde{d} of the input stays the same, the rank of the attention matrices increases with the number of heads h . The middle panel of Figure 9 confirms this numerically. We see that the numerical rank of the projection matrix $\mathbf{W}_Q^{(i)}$ in each head remains relatively constant for reasonably large ε as we increase the number of heads.

Our observation in Figure 9 is empirical and explanatory: it says in a pretrained TSFM, the attention matrices $\mathbf{W}_{Q/K/V}$ are essentially doing the head-dependent sketching (see eq. (20)). From a methodological point of view, if we train a compressed model (i.e., the one that is parameterized by \mathbf{W}_1 and \mathbf{W}_2 instead of \mathbf{W}_Q) from scratch, should we adopt a parameter-efficient head-independent sketching in eq. (19) or a head-dependent one, which requires more parameters? If we know the numerical rank \tilde{d} of Ξ so that $\sigma_{\tilde{d}+1}(\Xi)$ is tiny, then from an approximation theory perspective, there is little loss using eq. (19). The only problem is: we do not know \tilde{d} a priori. If we choose \tilde{d} small enough that $\sigma_{\tilde{d}+1}(\Xi)$ is still relatively large, then the following theorem (see Appendix E for the proof) shows the benefit of using a head-dependent sketching $\mathbf{W}_2^{(h)}$ in eq. (20).

Theorem 6. Fix some $d = h \times d_h$ for two positive integers h and d_h and $L \geq d$. Let $C \geq 1$ be any stability bound. Let $1 \geq \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d > 0$ be any fixed sequence. There exist an input $\mathbf{U} \in \mathbb{R}^{d \times L}$ with $\sigma_j(\mathbf{U}) = \sigma_j$ for all $1 \leq j \leq d$ and three matrices $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$, such that the following statements hold for any $\tilde{d} < d_h$:

- Given any rank- \tilde{d} matrix $\tilde{\mathbf{W}}_V \in \mathbb{R}^{d \times d}$ with $\|\tilde{\mathbf{W}}_V\|_F \leq C\|\mathbf{W}_V\|_F$, we have

$$\left\| \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) - \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \tilde{\mathbf{W}}_V, h) \right\|_2 \geq \frac{1}{2}\sigma_{\tilde{d}+1}.$$

- There exist low-rank matrices $\tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V \in \mathbb{R}^{d \times d}$ with $\|\tilde{\mathbf{W}}_Q\|_F \leq \|\mathbf{W}_Q\|_F, \|\tilde{\mathbf{W}}_K\|_F \leq \|\mathbf{W}_K\|_F$, and $\|\tilde{\mathbf{W}}_V\|_F \leq \|\mathbf{W}_V\|_F$, such that

$$\tilde{\mathbf{W}}_{Q/K/V} = \begin{bmatrix} \mathbf{W}_{Q/K/V,1}^{(1)} \mathbf{W}_{Q/K/V,2}^{(1)} \\ \vdots \\ \mathbf{W}_{Q/K/V,1}^{(h)} \mathbf{W}_{Q/K/V,2}^{(h)} \end{bmatrix}, \quad \mathbf{W}_{Q/K/V,1}^{(i)} \in \mathbb{R}^{d_h \times (\tilde{d}+1)}, \quad \mathbf{W}_{Q/K/V,2}^{(i)} \in \mathbb{R}^{(\tilde{d}+1) \times d},$$

where every row of $\mathbf{W}_{Q/K/V,2}^{(i)}$, except the last row of $\mathbf{W}_{Q/K/V,2}^{(h)}$, contains at most d_h non-zero entries, and satisfies that

$$\left\| \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) - \text{MH-Attention}(\mathbf{U}; \tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V, h) \right\|_2 \leq 4\sqrt{h}\sigma_{h\tilde{d}+1}.$$

Theorem 6 highlights an “error-correcting” mechanism of the head-dependent low-rank design in eq. (20) in the following sense: if we choose a \tilde{d} where $\sigma_{\tilde{d}+1}$ is large, then the first statement

shows that the shared sketching in eq. (19) inevitably limits the expressiveness of the reduced multi-head attention layer. Using a head-dependent low-rank representation relaxes the error to the order of $\sigma_{h\tilde{d}+1}$, which can be significantly smaller than $\sigma_{\tilde{d}+1}$ when the number of heads h is large and the singular values decay fast. While using eq. (20) requires a larger number of parameters, Theorem 6 suggests the potential of a sparse parameterization of the sketching matrices $\mathbf{W}_{Q/K/V,2}^{(i)}$. We leave this as a promising future direction.

E.1 PROOF OF THEOREM 6

We prove that a sparse but head-dependent sketching performs better in theory than a head-independent sketching, especially when \tilde{d} is lower than the numerical rank of the input matrix. The intuition is that if \tilde{d} is too small, then the sketching performed by a single sketching matrix must be bad, but if we use head-dependent sketchings, then we can leverage multiple of them to still obtain a lot from the input matrix.

Proof of Theorem 6. Set $\mathbf{W}_V = \mathbf{I}_d$. Interleave the singular values $\sigma_1, \dots, \sigma_d$ as follows:

$$\begin{aligned}\boldsymbol{\sigma}^{(1)} &= (\sigma_1^{(1)}, \sigma_2^{(1)}, \dots, \sigma_{d_h}^{(1)}) = (\sigma_1, \sigma_{h+1}, \dots, \sigma_{(d_h-1)h+1}), \\ \boldsymbol{\sigma}^{(2)} &= (\sigma_1^{(2)}, \sigma_2^{(2)}, \dots, \sigma_{d_h}^{(2)}) = (\sigma_2, \sigma_{h+2}, \dots, \sigma_{(d_h-1)h+2}), \\ &\vdots \\ \boldsymbol{\sigma}^{(h)} &= (\sigma_1^{(h)}, \sigma_2^{(h)}, \dots, \sigma_{d_h}^{(h)}) = (\sigma_h, \sigma_{2h}, \dots, \sigma_{d_h \cdot h}).\end{aligned}$$

Define the input matrix by

$$\mathbf{U} = \left[\underbrace{\text{diag}(\text{diag}(\boldsymbol{\sigma}^{(1)}), \dots, \text{diag}(\boldsymbol{\sigma}^{(h)}))}_{\mathbf{U}_D} \mid \mathbf{0}_{d \times (L-d)} \right].$$

For every $1 \leq i \leq h$ and $1 \leq j \leq d_h$, let $\mathbf{v}_j^{(i)} \in \mathbb{R}^{d_h}$ be a unit vector satisfying that

$$\left\| \mathbf{v}_j^{(i)\top} \mathbf{v}_{j'}^{(i')} \right\|_2 < 1, \quad (i, j) \neq (i', j').$$

Then, consider the following matrix:

$$\overline{\mathbf{W}}_Q = \overline{\mathbf{W}}_K = \left[\mathbf{v}_1^{(1)}/\sigma_1^{(1)} \quad \cdots \quad \mathbf{v}_{d_h}^{(1)}/\sigma_{d_h}^{(1)} \mid \cdots \mid \mathbf{v}_1^{(h)}/\sigma_1^{(h)} \quad \cdots \quad \mathbf{v}_{d_h}^{(h)}/\sigma_{d_h}^{(h)} \right] \in \mathbb{R}^{d_h \times d}.$$

Then, by our construction, the matrix $\mathbf{U}^\top \overline{\mathbf{W}}_Q^\top \overline{\mathbf{W}}_K \mathbf{U}$ satisfies that

$$(\mathbf{U}^\top \overline{\mathbf{W}}_Q^\top \overline{\mathbf{W}}_K \mathbf{U})(j, j) = 1, \quad |(\mathbf{U}^\top \overline{\mathbf{W}}_Q^\top \overline{\mathbf{W}}_K \mathbf{U})(i, j)| < 1, \quad i \neq j.$$

Set

$$\varepsilon = \min \left\{ (1 + C)^{-1} \sigma_{\tilde{d}+1} / 2, \min_{1 \leq i \leq h} \sigma_{d+1}^{(i)} / (2 \|\mathbf{W}_V \mathbf{U}\|_2 + 2\sigma_1), 1 \right\}.$$

By choosing a sufficiently large $\alpha > 0$, we guarantee that

$$\left\| \underbrace{\text{softmax} \left(\frac{\mathbf{U}^\top (\alpha \overline{\mathbf{W}}_Q^\top)(\alpha \overline{\mathbf{W}}_K) \mathbf{U}}{\sqrt{d_h}} \right)}_{\mathbf{G}} - \underbrace{\left[\begin{array}{c|c} \mathbf{I}_d & L^{-1} \mathbf{1}_{d \times (L-d)} \\ \hline \mathbf{0}_{(L-d) \times d} & L^{-1} \mathbf{1}_{(L-d) \times (L-d)} \end{array} \right]}_{\mathbf{G}} \right\|_F \leq \varepsilon. \quad (21)$$

Let query and key matrices \mathbf{W}_Q and \mathbf{W}_K be defined in eq. (18) using submatrices $\mathbf{W}_Q^{(i)} = \alpha \overline{\mathbf{W}}_Q$ and $\mathbf{W}_K^{(i)} = \alpha \overline{\mathbf{W}}_K$, respectively, for all $1 \leq i \leq h$. We use these matrices \mathbf{W}_Q , \mathbf{W}_K , \mathbf{W}_V , and \mathbf{U} to prove the two claims.

Claim 1: No low-rank parameterization gets to an accuracy below $\sigma_{\tilde{d}+1}$. Let $\tilde{\mathbf{W}}_V \in \mathbb{R}^{d \times d}$ be any rank- \tilde{d} matrix. We can write the approximation error as

$$\begin{aligned} & \left\| \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) - \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \tilde{\mathbf{W}}_V, h) \right\|_2 \\ &= \left\| \begin{bmatrix} \mathbf{W}_V^{(1)} \mathbf{U} \mathbf{G} \\ \vdots \\ \mathbf{W}_V^{(h)} \mathbf{U} \mathbf{G} \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{W}}_V^{(1)} \mathbf{U} \mathbf{G} \\ \vdots \\ \tilde{\mathbf{W}}_V^{(h)} \mathbf{U} \mathbf{G} \end{bmatrix} \right\|_2 \\ &\geq \underbrace{\left\| \begin{bmatrix} \mathbf{W}_V^{(1)} \mathbf{U} \tilde{\mathbf{G}} \\ \vdots \\ \mathbf{W}_V^{(h)} \mathbf{U} \tilde{\mathbf{G}} \end{bmatrix} - \begin{bmatrix} \tilde{\mathbf{W}}_V^{(1)} \mathbf{U} \tilde{\mathbf{G}} \\ \vdots \\ \tilde{\mathbf{W}}_V^{(h)} \mathbf{U} \tilde{\mathbf{G}} \end{bmatrix} \right\|_2}_{E_{\text{lead}}} - \underbrace{\left\| \begin{bmatrix} (\mathbf{W}_V^{(1)} - \tilde{\mathbf{W}}_V^{(1)}) \mathbf{U} (\mathbf{G} - \tilde{\mathbf{G}}) \\ \vdots \\ (\mathbf{W}_V^{(h)} - \tilde{\mathbf{W}}_V^{(h)}) \mathbf{U} (\mathbf{G} - \tilde{\mathbf{G}}) \end{bmatrix} \right\|_2}_{E_{\text{trail}}}. \end{aligned}$$

To control E_{trail} , we note that

$$\begin{aligned} E_{\text{trail}} &\leq \max_{1 \leq i \leq h} \left\| (\mathbf{W}_V^{(i)} - \tilde{\mathbf{W}}_V^{(i)}) \mathbf{U} (\mathbf{G} - \tilde{\mathbf{G}}) \right\|_2 \leq \max_{1 \leq i \leq h} \left(\left\| \mathbf{W}_V^{(i)} \right\|_2 + \left\| \tilde{\mathbf{W}}_V^{(i)} \right\|_2 \right) \left\| \mathbf{G} - \tilde{\mathbf{G}} \right\|_2 \\ &\leq (1 + C)\varepsilon = \frac{1}{2}\sigma_{\tilde{d}+1}. \end{aligned} \quad (22)$$

To control E_{lead} , we note that

$$\begin{bmatrix} (\mathbf{W}_V^{(1)} - \tilde{\mathbf{W}}_V^{(1)}) \mathbf{U} \tilde{\mathbf{G}} \\ \vdots \\ (\mathbf{W}_V^{(h)} - \tilde{\mathbf{W}}_V^{(h)}) \mathbf{U} \tilde{\mathbf{G}} \end{bmatrix} = \begin{bmatrix} (\mathbf{W}_V^{(1)} - \tilde{\mathbf{W}}_V^{(1)}) [\mathbf{U}_D \mid L^{-1} \mathbf{U}_D \mathbf{1}_{d \times (L-d)}] \\ \vdots \\ (\mathbf{W}_V^{(h)} - \tilde{\mathbf{W}}_V^{(h)}) [\mathbf{U}_D \mid L^{-1} \mathbf{U}_D \mathbf{1}_{d \times (L-d)}] \end{bmatrix}.$$

Hence, we have

$$\begin{aligned} E_{\text{lead}} &= \left\| \begin{bmatrix} (\mathbf{W}_V^{(1)} - \tilde{\mathbf{W}}_V^{(1)}) \mathbf{U} \tilde{\mathbf{G}} \\ \vdots \\ (\mathbf{W}_V^{(h)} - \tilde{\mathbf{W}}_V^{(h)}) \mathbf{U} \tilde{\mathbf{G}} \end{bmatrix} \right\|_2 \geq \left\| \begin{bmatrix} (\mathbf{W}_V^{(1)} - \tilde{\mathbf{W}}_V^{(1)}) \mathbf{U}_D \\ \vdots \\ (\mathbf{W}_V^{(h)} - \tilde{\mathbf{W}}_V^{(h)}) \mathbf{U}_D \end{bmatrix} \right\|_2 \\ &= \left\| \mathbf{W}_V \mathbf{U}_D - \tilde{\mathbf{W}}_V \mathbf{U}_D \right\|_2 \geq \sigma_{\tilde{d}+1}, \end{aligned} \quad (23)$$

where the last inequality follows from the fact that the singular values of $\mathbf{W}_V \mathbf{U}_D$ are exactly $\sigma_1, \dots, \sigma_d$ and $\tilde{\mathbf{W}}_V \mathbf{U}_D$ is a matrix whose rank is at most d . Combining eq. (22) and (23), we have

$$\begin{aligned} & \left\| \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) - \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \tilde{\mathbf{W}}_V, h) \right\|_2 \\ &\geq E_{\text{lead}} - E_{\text{trail}} \geq \frac{1}{2}\sigma_{\tilde{d}+1}. \end{aligned}$$

Claim 2: A sparse parameterization gets to an accuracy below $\sigma_{h\tilde{d}+1}$. For every $1 \leq i \leq h$, define $\mathbf{W}_{V,L}^{(i)} \in \mathbb{R}^{d_h \times \tilde{d}}$ and $\mathbf{W}_{V,R}^{(i)} \in \mathbb{R}^{\tilde{d} \times d}$ as follows:

$$\mathbf{W}_{V,L}^{(i)} = \begin{bmatrix} \mathbf{I}_{\tilde{d}} \\ \mathbf{0}_{(d_h - \tilde{d}) \times \tilde{d}} \end{bmatrix}, \quad \mathbf{W}_{V,R}^{(i)} = \begin{bmatrix} \mathbf{0}_{\tilde{d} \times d_h(i-1)} \mid \mathbf{I}_{\tilde{d}} \mid \mathbf{0}_{\tilde{d} \times ((d_h - \tilde{d}) + d_h(h-i))} \end{bmatrix}.$$

Then, we have

$$\mathbf{W}_{V,L}^{(i)} \mathbf{W}_{V,R}^{(i)} \mathbf{U} = \begin{bmatrix} \mathbf{0}_{\tilde{d} \times d_h(i-1)} & \text{diag}(\sigma_1^{(i)}, \dots, \sigma_{\tilde{d}}^{(i)}) & \mathbf{0}_{\tilde{d} \times ((d_h - \tilde{d}) + d_h(h-i))} \\ \mathbf{0}_{(d_h - \tilde{d}) \times d_h(i-1)} & \mathbf{0}_{(d_h - \tilde{d}) \times \tilde{d}} & \mathbf{0}_{(d_h - \tilde{d}) \times ((d_h - \tilde{d}) + d_h(h-i))} \end{bmatrix},$$

while

$$\mathbf{W}_V^{(i)} \mathbf{U} = \begin{bmatrix} \mathbf{0}_{\tilde{d} \times d_h(i-1)} \mid \text{diag}(\sigma_1^{(i)}, \dots, \sigma_{d_h}^{(i)}) \mid \mathbf{0}_{\tilde{d} \times d_h(h-i)} \end{bmatrix}.$$

1458 Hence, we have that

$$\| \mathbf{W}_V^{(i)} \mathbf{U} - \mathbf{W}_{V,L}^{(i)} \mathbf{W}_{V,R}^{(i)} \mathbf{U} \|_2 = \sigma_{\tilde{d}+1}^{(i)}. \quad (24)$$

1460
1461 Similarly, let $\mathbf{W}_{Q,L}^{(i)}, \mathbf{W}_{K,L}^{(i)} \in \mathbb{R}^{d_h \times (\tilde{d}+1)}$ and $\mathbf{W}_{Q,R}^{(i)}, \mathbf{W}_{K,R}^{(i)} \in \mathbb{R}^{(\tilde{d}+1) \times d}$ be given by

$$\begin{aligned} \mathbf{W}_{Q,L}^{(i)} &= \mathbf{W}_{K,L}^{(i)} = \alpha \begin{bmatrix} \mathbf{v}_1^{(i)} / \sigma_1^{(i)} & \cdots & \mathbf{v}_{\tilde{d}}^{(i)} / \sigma_{\tilde{d}}^{(i)} & \mathbf{v}_1^{(i')} \end{bmatrix}, \quad i' = \text{mod}(i, h) + 1, \\ \mathbf{W}_{Q,R}^{(i)} &= \begin{cases} \left[\begin{array}{c|c} \mathbf{I}_{\tilde{d}} & \mathbf{0}_{d-\tilde{d}} \\ \hline \mathbf{0}_{1 \times \tilde{d}} & [1/\sigma_{\tilde{d}+1}, 0, \dots, 0] \end{array} \right], & i = 1, \\ \left[\begin{array}{c|c|c} \mathbf{0}_{\tilde{d} \times d_h(i-1)} & \mathbf{I}_{\tilde{d}} & \mathbf{0}_{\tilde{d} \times ((d_h-\tilde{d})+d_h(h-i))} \\ \hline [1/\sigma_1, 0, \dots, 0] & \mathbf{0}_{1 \times \tilde{d}} & \mathbf{0}_{1 \times ((d_h-\tilde{d})+d_h(h-i))} \end{array} \right], & i > 1, \end{cases} \\ \mathbf{W}_{K,R}^{(i)} &= \left[\begin{array}{c|c|c} \mathbf{0}_{\tilde{d} \times d_h(i-1)} & \mathbf{I}_{\tilde{d}} & \mathbf{0}_{\tilde{d} \times ((d_h-\tilde{d})+d_h(h-i))} \\ \hline [1/\sigma_1, \dots, 1/\sigma_{d_h(i-1)}] & \mathbf{0}_{1 \times \tilde{d}} & [1/\sigma_{d_h(i-1)+\tilde{d}+1}, \dots, 1/\sigma_d] \end{array} \right]. \end{aligned}$$

1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477 Note that since we have $\|\mathbf{v}_j^{(i)}\|_2$ for every $1 \leq i \leq h$ and $1 \leq j \leq d_h$, we clearly have that $\|\mathbf{W}_{Q,L}^{(i)} \mathbf{W}_{Q,R}^{(i)}\|_F \leq \|\alpha \bar{\mathbf{W}}_Q\|_F$ and $\|\mathbf{W}_{K,L}^{(i)} \mathbf{W}_{K,R}^{(i)}\|_F = \|\alpha \bar{\mathbf{W}}_K\|_F$ for every $1 \leq i \leq h$. Hence, we have $\|\bar{\mathbf{W}}_Q\|_F \leq \|\mathbf{W}_Q\|_F$ and $\|\bar{\mathbf{W}}_K\|_F = \|\mathbf{W}_K\|_F$. From now on, we only argue for the first head, i.e., $i = 1$. The rest follows easily from symmetry. Set $i = 1$, we have

$$\underbrace{\left(\mathbf{U}^\top \mathbf{W}_{Q,R}^{(1)}{}^\top \mathbf{W}_{Q,L}^{(1)}{}^\top \mathbf{W}_{K,L}^{(1)} \mathbf{W}_{K,R}^{(1)} \mathbf{U} \right)}_{\tilde{\mathbf{T}}^{(1)}} (j, k) = \alpha^2 \begin{cases} \mathbf{v}_j^{(1)\top} \mathbf{v}_k^{(1)}, & 1 \leq j \leq \tilde{d}, \quad 1 \leq k \leq \tilde{d}, \\ \mathbf{v}_j^{(1)\top} \mathbf{v}_1^{(2)}, & 1 \leq j \leq \tilde{d}, \quad \tilde{d} < k \leq d, \\ \mathbf{v}_1^{(2)\top} \mathbf{v}_k^{(1)}, & j = \tilde{d} + 1, \quad 1 \leq k \leq \tilde{d}, \\ \mathbf{v}_1^{(2)\top} \mathbf{v}_1^{(2)}, & j = \tilde{d} + 1, \quad \tilde{d} < k \leq d, \\ 0, & \text{otherwise}. \end{cases}$$

1478
1479
1480
1481
1482
1483
1484 That is, for $1 \leq j \leq d$, the j th column of $\tilde{\mathbf{T}}^{(1)}$ has exactly one entry that equals 1, which is the j th element for $1 \leq j \leq \tilde{d}$ and the $(j+1)$ st element for $j > \tilde{d}$. Moreover, for $j > d$, the j th column of $\tilde{\mathbf{T}}^{(1)}$ is zero. Hence, from the definition of α , we have

$$\left\| \text{softmax} \left(\tilde{\mathbf{T}}^{(1)} \right) (1 : \tilde{d}, :) - \mathbf{G}(1 : \tilde{d}, :) \right\|_2 \leq \left\| \text{softmax} \left(\tilde{\mathbf{T}}^{(1)} \right) (1 : \tilde{d}, :) - \tilde{\mathbf{G}}(1 : \tilde{d}, :) \right\|_2 + \left\| \tilde{\mathbf{G}} - \mathbf{G} \right\|_2 \leq 2\varepsilon. \quad (25)$$

1491 Combine eq. (24) and (25). We have

$$\begin{aligned} &\left\| \mathbf{W}_V^{(1)} \mathbf{U} \mathbf{G} - \mathbf{W}_{V,L}^{(1)} \mathbf{W}_{V,R}^{(1)} \mathbf{U} \text{softmax} \left(\tilde{\mathbf{T}}^{(1)} \right) \right\|_2 \\ &\leq \left\| \mathbf{W}_V^{(1)} \mathbf{U} \mathbf{G} - \mathbf{W}_{V,L}^{(1)} \mathbf{W}_{V,R}^{(1)} \mathbf{U} \mathbf{G} \right\|_2 + \left\| \mathbf{W}_{V,L}^{(1)} \mathbf{W}_{V,R}^{(1)} \mathbf{U} \mathbf{G} - \mathbf{W}_{V,L}^{(1)} \mathbf{W}_{V,R}^{(1)} \mathbf{U} \text{softmax} \left(\tilde{\mathbf{T}}^{(1)} \right) \right\|_2 \\ &= \left\| \mathbf{W}_V^{(1)} \mathbf{U} \mathbf{G} - \mathbf{W}_{V,L}^{(1)} \mathbf{W}_{V,R}^{(1)} \mathbf{U} \mathbf{G} \right\|_2 + \left\| \mathbf{W}_{V,L}^{(1)} \mathbf{W}_{V,R}^{(1)} \mathbf{U} \left(\mathbf{G} - \text{softmax} \left(\tilde{\mathbf{T}}^{(1)} \right) \right) (1 : \tilde{d}, :) \right\|_2 \\ &\leq \left\| \mathbf{W}_V^{(1)} \mathbf{U} - \mathbf{W}_{V,L}^{(1)} \mathbf{W}_{V,R}^{(1)} \mathbf{U} \right\|_2 \|\mathbf{G}\|_2 + \left(\left\| \mathbf{W}_V^{(1)} \mathbf{U} \right\|_2 + \sigma_{\tilde{d}+1}^{(1)} \right) \left\| \left(\mathbf{G} - \text{softmax} \left(\tilde{\mathbf{T}}^{(1)} \right) \right) (1 : \tilde{d}, :) \right\|_2 \\ &\leq \sigma_{\tilde{d}+1}^{(1)} (2 + \varepsilon) + 2\varepsilon \left\| \mathbf{W}_{V,L}^{(1)} \mathbf{W}_{V,R}^{(1)} \mathbf{U} \right\|_2 \leq 4\sigma_{\tilde{d}+1}^{(1)} \leq 4\sigma_{h\tilde{d}+1}, \end{aligned}$$

1501
1502
1503
1504
1505
1506 where the first inequality follows from the triangle inequality, the first equation follows from the sparsity of $\mathbf{W}_{V,R}^{(1)}$, the second inequality follows from the sub-multiplicity of the spectral norm and eq. (24), the third inequality follows from eq. (24) and eq. (25), the fifth inequality follows from the definition of ε , and the last inequality follows from the definition of $\sigma_{\tilde{d}+1}^{(1)}$. Notably, this inequality holds for every head $1 \leq i \leq h$. Hence, we have

$$\begin{aligned} &\left\| \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) - \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \tilde{\mathbf{W}}_V, h) \right\|_2 \\ &\leq \sqrt{h} \max_{1 \leq i \leq h} \left\| \mathbf{W}_V^{(i)} \mathbf{U} \mathbf{G} - \mathbf{W}_{V,L}^{(i)} \mathbf{W}_{V,R}^{(i)} \mathbf{U} \text{softmax} \left(\tilde{\mathbf{T}}^{(i)} \right) \right\|_2 \leq 4\sqrt{h} \sigma_{h\tilde{d}+1}. \end{aligned}$$

1511 The proof is complete. \square

1512 **F PROOF OF THEOREM 4 AND COROLLARY 1**
 1513

1514 In this section, we prove the idea of flow-of-ranks. The upper bound is proved via a straightforward
 1515 application of singular value inequalities while the lower bound requires a careful construction.
 1516

1517 *Proof of Theorem 4.* We prove the two statements separately.
 1518

1519 **The upper bound.** Fix a head index $1 \leq i \leq h$ and let
 1520

$$\mathbf{Y}^{(i)} = \mathbf{W}_V^{(i)} \mathbf{U} \text{ softmax}\left(\mathbf{T}^{(i)}\right), \quad \mathbf{T}^{(i)} = \frac{\mathbf{U}^\top \mathbf{W}_Q^{(i)\top} \mathbf{W}_K^{(i)} \mathbf{U}}{\sqrt{d_h}}.$$

1523 Then, every entry of $\mathbf{T}^{(i)}$ has a magnitude no greater than 1 because every column of \mathbf{U} has a 2-
 1524 norm no greater than 1 and $\|\mathbf{W}_Q^{(i)\top} \mathbf{W}_K^{(i)}\|_2 \leq \sqrt{d_h}$. Hence, every entry of $\text{softmax}(\mathbf{T}^{(i)})$ has a
 1525 magnitude between $L^{-1}e^{-2}$ and $L^{-1}e^2$. This means
 1526

$$\sigma_1\left(\text{softmax}\left(\mathbf{T}^{(i)}\right)\right) = \left\| \text{softmax}\left(\mathbf{T}^{(i)}\right) \right\|_2 \leq \left\| \text{softmax}\left(\mathbf{T}^{(i)}\right) \right\|_F \leq e^2.$$

1529 Moreover, since $\left\| \mathbf{W}_V^{(i)} \right\|_2 \leq 1$, we have that
 1530

$$\sigma_j\left(\mathbf{Y}^{(i)}\right) \leq e^2 \sigma_j(\mathbf{U}), \quad 1 \leq j \leq d_h.$$

1533 Fix some $1 \leq j \leq d$ and let $j_h = \lfloor (j-1)/h \rfloor + 1$. When we concatenate $\mathbf{Y}^{(i)}$ to form \mathbf{Y} , using
 1534 the singular value inequality that $\sigma_{i+j-1}(\mathbf{A} + \mathbf{B}) \leq \sigma_i(\mathbf{A}) + \sigma_j(\mathbf{B})$ and defining $\mathbf{Y}^{\mathcal{I}}$ to be the
 1535 concatenation of $\mathbf{Y}^{(i)}$ for all $i \in \mathcal{I}$, we have
 1536

$$\begin{aligned} \sigma_j(\mathbf{Y}) &\leq \sigma_{j_h}\left(\mathbf{Y}^{(h)}\right) + \sigma_{j-j_h+1}\left(\mathbf{Y}^{[h-1]}\right) \leq \sigma_{j_h}\left(\mathbf{Y}^{(h)}\right) + \sigma_{j_h}\left(\mathbf{Y}^{(h-1)}\right) + \sigma_{j-2j_h+2}\left(\mathbf{Y}^{[h-2]}\right) \\ &\leq \dots \leq \left(\sum_{i=2}^h \sigma_{j_h}\left(\mathbf{Y}^{(i)}\right) \right) + \sigma_{j-(h-1)(j_h-1)}\left(\mathbf{Y}^{(1)}\right) \leq \sum_{i=1}^h \sigma_{j_h}\left(\mathbf{Y}^{(i)}\right) \leq e^2 h \sigma_{j_h}(\mathbf{U}). \end{aligned}$$

1542 Hence, applying the same inequality again, we have that
 1543

$$\sigma_k(\mathbf{Z}) = \sigma_k\left(\mathbf{U} + \frac{1}{\sqrt{N}} \mathbf{Y}\right) \leq \sigma_{k-j+1}(\mathbf{U}) + \sigma_j\left(\frac{1}{\sqrt{N}} \mathbf{Y}\right) \leq \sigma_{k-j+1} + \frac{e^2 h}{\sqrt{N}} \sigma_{\lfloor (j-1)/h \rfloor + 1}. \quad (26)$$

1546 Moreover, from the triangle inequality, we have
 1547

$$\sigma_1(\mathbf{Z}) \geq \sigma_1(\mathbf{U}) - \sigma_1\left(\frac{1}{\sqrt{N}} \mathbf{Y}\right) \geq \sigma_1(\mathbf{U}) - \frac{1}{\sqrt{N}} e^2 h \sigma_1(\mathbf{U}) \geq \frac{1}{2} \sigma_1(\mathbf{U}). \quad (27)$$

1550 Combining eq. (26) and (27), we prove the theorem.
 1551

1552 **The lower bound.** Let $1 = \sigma_1 \geq \dots \geq \sigma_d > 0$ be given. Define the input matrix to be
 1553

$$\mathbf{U} = [\text{diag}(\sigma_1, \dots, \sigma_d) \quad \mathbf{0}_{d \times (L-d)}].$$

1554 For every $1 \leq i \leq h$, let $\mathbf{W}_V^{(i)}$ be
 1555

$$\mathbf{W}_V^{(i)} = \begin{bmatrix} \mathbf{I}_{d_h-1} & \mathbf{0}_{(d_h-1) \times (d-d_h+1)} \\ \mathbf{0}_{1 \times (d_h-1)} & \mathbf{0}_{1 \times (d-d_h+1)} \end{bmatrix}.$$

1556 Set $\varepsilon = \sqrt{N} \sigma_d / (5h)$. For every $1 \leq i \leq h$, let $\mathbf{P}^{(i)} \in \mathbb{R}^{L \times L}$ be the matrix so that $\mathbf{P}^{(i)}(1 : d_h, ((i-1)d_h + 1) : id_h) = \mathbf{I}_{d_h}$, $\mathbf{P}^{(i)}(d_h, k) = 1$ for every $1 \leq k < (i-1)d_h$ and every $k > id_h$,
 1557 and is zero elsewhere. For a sufficiently large $\alpha > 0$ determined later on, let
 1558

$$\mathbf{W}_Q^{(i)} = \alpha [\text{diag}(\sigma_1^{-1}, \dots, \sigma_{d_h}^{-1}) \quad \mathbf{0}_{d_h \times (d-d_h)}],$$

$$\mathbf{W}_K^{(i)} = \left[\begin{array}{c|c|c} \mathbf{0}_{(d_h-1) \times (i-1)d_h} & \text{diag}(\sigma_{(i-1)d_h+1}^{-1}, \dots, \sigma_{id_h}^{-1}) & \mathbf{0}_{(d_h-1) \times (h-i)d_h} \\ \hline [\sigma_1^{-1} \quad \dots \quad \sigma_{(i-1)d_h}^{-1}] & & [\sigma_{id_h+1}^{-1} \quad \dots \quad \sigma_d^{-1}] \end{array} \right].$$

1566 Then, we have
 1567

$$1568 \mathbf{T}^{(i)} := \mathbf{U}^\top \mathbf{W}_Q^{(i)\top} \mathbf{W}_K^{(i)} \mathbf{U} \\ 1569 \\ 1570 = \alpha \begin{bmatrix} \mathbf{I}_{d_h} & \mathbf{0}_{d_h \times (L-d_h)} \end{bmatrix}^\top \begin{bmatrix} \mathbf{0}_{(d_h-1) \times (i-1)d_h} & \mathbf{I}_{d_h} & \mathbf{0}_{(d_h-1) \times (h-i)d_h} & \mathbf{0}_{d_h \times (L-d)} \\ \mathbf{1}_{1 \times (i-1)d_h} & & \mathbf{1}_{1 \times (h-i)d_h} & \end{bmatrix} = \alpha \mathbf{P}^{(i)}. \\ 1571 \\ 1572$$

1573 Therefore, by picking α sufficiently large, we have
 1574

$$1575 \|\mathbf{G}^{(i)} - \tilde{\mathbf{G}}^{(i)}\|_2 \leq \epsilon, \quad \mathbf{G}^{(i)} := \text{softmax}(\mathbf{T}^{(i)}), \quad \tilde{\mathbf{G}}^{(i)}(j, k) = \begin{cases} 1, & k = j + (i-1)d_h, \\ 1, & j = d_h \text{ and } k \leq (i-1)d_h, \\ 1, & j = d_h \text{ and } k > id_h, \\ L^{-1}, & k > d, \\ 0, & \text{otherwise.} \end{cases} \\ 1576 \\ 1577 \\ 1578 \\ 1579 \\ 1580 \\ 1581 \quad \text{Note that we have} \\ 1582 \\ 1583 \mathbf{W}_V^{(i)} \mathbf{U} \tilde{\mathbf{G}}^{(i)} = \begin{bmatrix} \text{diag}(\sigma_1, \dots, \sigma_{h_d-1}) & 0 & \mathbf{0}_{(d_h-1) \times (L-d_h)} \\ \mathbf{0}_{1 \times (d_h-1)} & 0 & \mathbf{0}_{1 \times (L-d_h)} \end{bmatrix} \\ 1584 \\ 1585 \\ 1586 \\ 1587 \\ 1588 \\ 1589 \\ 1590 \\ 1591 \times \begin{bmatrix} \mathbf{0}_{(d_h-1) \times (i-1)d_h} & \mathbf{I}_{d_h} & \mathbf{0}_{(d_h-1) \times (h-i)d_h} & L^{-1} \mathbf{1}_{d_h \times (L-d)} \\ \mathbf{1}_{1 \times (i-1)d_h} & & \mathbf{1}_{1 \times (h-i)d_h} & \\ \mathbf{0}_{(L-d_h) \times (i-1)d_h} & \mathbf{0}_{(L-d_h) \times d_h} & \mathbf{0}_{(L-d_h) \times (h-i)d_h} & L^{-1} \mathbf{1}_{(L-d_h) \times (L-d)} \end{bmatrix} \\ = \underbrace{\begin{bmatrix} \mathbf{0}_{d_h \times (i-1)d_h} & \text{diag}(\sigma_1, \dots, \sigma_{h_d-1}, 0) & \mathbf{0}_{d_h \times (h-i)d_h} \\ \mathbf{Y}_{\text{lead}}^{(i)} & & \end{bmatrix}}_{\mathbf{Y}_{\text{lead}}} \underbrace{\begin{bmatrix} L^{-1} \text{diag}(\sigma_1, \dots, \sigma_{h_d-1}) \mathbf{1}_{(d_h-1) \times (L-d)} \\ \mathbf{0}_{1 \times (L-d)} \\ \mathbf{Y}_{\text{trail}}^{(i)} \end{bmatrix}}_{\mathbf{Y}_{\text{trail}}}.$$

1592 Note that the mulithead attention output is given by
 1593

$$\mathbf{Y} = \text{MH-Attention}(\mathbf{U}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) = [\mathbf{Y}_{\text{lead}} \quad \mathbf{Y}_{\text{trail}}] + \mathbf{E}, \\ \mathbf{Y}_{\text{lead}} = \left[\mathbf{Y}_{\text{lead}}^{(1)\top} \quad \dots \quad \mathbf{Y}_{\text{lead}}^{(h)\top} \right]^\top, \quad \mathbf{Y}_{\text{trail}} = \left[\mathbf{Y}_{\text{trail}}^{(1)\top} \quad \dots \quad \mathbf{Y}_{\text{trail}}^{(h)\top} \right]^\top,$$

1594 where
 1595

$$1596 \|\mathbf{E}\|_2 \leq \sum_{i=1}^h \|\mathbf{W}_V^{(i)} \mathbf{U} \tilde{\mathbf{G}}^{(i)} - \mathbf{W}_V^{(i)} \mathbf{U} \mathbf{G}^{(i)}\|_2 \leq \sum_{i=1}^h \|\mathbf{W}_V^{(i)} \mathbf{U}\|_2 \|\mathbf{G}^{(i)} - \tilde{\mathbf{G}}^{(i)}\|_2 \leq h\varepsilon,$$

1597 and the layer output is given by
 1598

$$\mathbf{Z} = \mathbf{U} + \frac{1}{\sqrt{N}} \mathbf{Y} = [\mathbf{Z}_{\text{lead}} \quad \mathbf{Z}_{\text{trail}}] + \frac{1}{\sqrt{N}} \mathbf{E}, \\ \mathbf{Z}_{\text{lead}} = \text{diag}(\sigma_1, \dots, \sigma_d) + \frac{1}{\sqrt{N}} \mathbf{Y}_{\text{lead}}, \quad \mathbf{Z}_{\text{trail}} = \frac{1}{\sqrt{N}} \mathbf{Y}_{\text{trail}}.$$

1600 Since \mathbf{Z}_{lead} is a diagonal matrix with nonnegative entries, its singular values equal its diagonal entries, which are
 1601

$$\begin{array}{cccccc} \sigma_1 + \frac{1}{\sqrt{N}} \sigma_1 & \sigma_2 + \frac{1}{\sqrt{N}} \sigma_2 & \cdots & \sigma_{d_h-1} + \frac{1}{\sqrt{N}} \sigma_{d_h-1} & \sigma_{d_h} \\ \sigma_{d_h+1} + \frac{1}{\sqrt{N}} \sigma_1 & \sigma_{d_h+2} + \frac{1}{\sqrt{N}} \sigma_2 & \cdots & \sigma_{d_h+(d_h-1)} + \frac{1}{\sqrt{N}} \sigma_{d_h-1} & \sigma_{d_h+d_h} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \sigma_{(h-1)d_h+1} + \frac{1}{\sqrt{N}} \sigma_1 & \sigma_{(h-1)d_h+2} + \frac{1}{\sqrt{N}} \sigma_2 & \cdots & \sigma_{(h-1)d_h+(d_h-1)} + \frac{1}{\sqrt{N}} \sigma_{d_h-1} & \sigma_{(h-1)d_h+d_h} \end{array} \quad (29)$$

1616 Moreover, since \mathbf{Z}_{lead} is a submatrix of \mathbf{Z} , the singular values of \mathbf{Z} are no smaller than the singular
 1617 values of \mathbf{Z}_{lead} . Given $1 \leq i \leq h$ and $1 \leq j \leq d_h$, let $\xi_{i,j} = \sigma_{(i-1)d_h+j} + \mathbb{1}_{\{j \neq d_h\}} \sigma_j / \sqrt{N}$ be the
 1618 element in row i and column j in the array in eq. (29). Then, we have that
 1619

$$\xi_{i,j} \leq \xi_{i',j'}, \quad \text{for any } i \geq i' \text{ and } j \geq j'.$$

1620 Hence, $\xi_{i,j}$ is no greater than the $(i \times j)$ th singular value of \mathbf{Z}_{lead} . That means for every $1 \leq i \leq h$,
 1621 the k th singular value of \mathbf{Z}_{lead} satisfies
 1622

$$1623 \sigma_k(\mathbf{Z}_{\text{lead}}) \geq \xi_{i,\lceil k/i \rceil} = \sigma_{(i-1)d_h + \lceil k/i \rceil} + \frac{\mathbb{1}_{\lceil k/i \rceil \neq d_h}}{\sqrt{N}} \sigma_{\lceil k/i \rceil}. \\ 1624$$

1625 Using Weyl's inequality, we have
 1626

$$1627 \sigma_k(\mathbf{Z}) \geq \sigma_k(\mathbf{Z}_{\text{lead}}) - \frac{1}{\sqrt{N}} \sigma_1(\mathbf{E}) \geq \sigma_k(\mathbf{Z}_{\text{lead}}) - \frac{h\varepsilon}{\sqrt{N}} \\ 1628 \\ 1629 \geq \left(1 - \frac{1}{5}\right) \left(\sigma_{(i-1)d_h + \lceil k/i \rceil} + \frac{\mathbb{1}_{\lceil k/i \rceil \neq d_h}}{\sqrt{N}} \sigma_{\lceil k/i \rceil}\right). \\ 1630$$

1631 Moreover, we have
 1632

$$1633 \sigma_1(\mathbf{Z}) \leq \sigma_1(\mathbf{Z}_{\text{lead}}) + \sigma_1(\mathbf{Z}_{\text{trail}}) + \frac{1}{\sqrt{N}} \sigma_1(\mathbf{E}) \leq \left(1 + \frac{1}{\sqrt{N}}\right) + \frac{1}{\sqrt{N}} \|\mathbf{Y}_{\text{trail}}\|_F + \frac{1}{\sqrt{N}} h\varepsilon \\ 1634 \\ 1635 \leq 2 + \frac{1}{\sqrt{N}} \sqrt{L \times d \times L^{-2}} + \frac{1}{5} \leq \frac{16}{5}. \\ 1636$$

1638 Combining the two inequalities above, we obtain the theorem. \square
 1639

1640 The proof of Corollary 1 is a straightforward manipulation of the ceiling and floor operators.
 1641

1642 *Proof of Corollary 1.* Set $j = \lfloor (hk+1)/(h+1) \rfloor$, we have

$$1643 \sigma_{k-j+1} = \sigma_{k-\lfloor (hk+1)/(h+1) \rfloor + 1} \leq \sigma_{\lceil k-(hk+1)/(h+1)+1 \rceil} = \sigma_{\lceil (k+h)/(h+1) \rceil}, \\ 1644$$

1645 and

$$1646 \sigma_{\lfloor (j-1)/h \rfloor + 1} = \sigma_{\lfloor (\lfloor (hk+1)/(h+1) \rfloor - 1)/h \rfloor + 1} \leq \sigma_{\lfloor ((hk+1)/(h+1)-1)/h \rfloor + 1} = \sigma_{\lfloor (k+1)/(h+1) \rfloor + 1}. \\ 1647$$

1648 The corollary follows from Theorem 4. \square

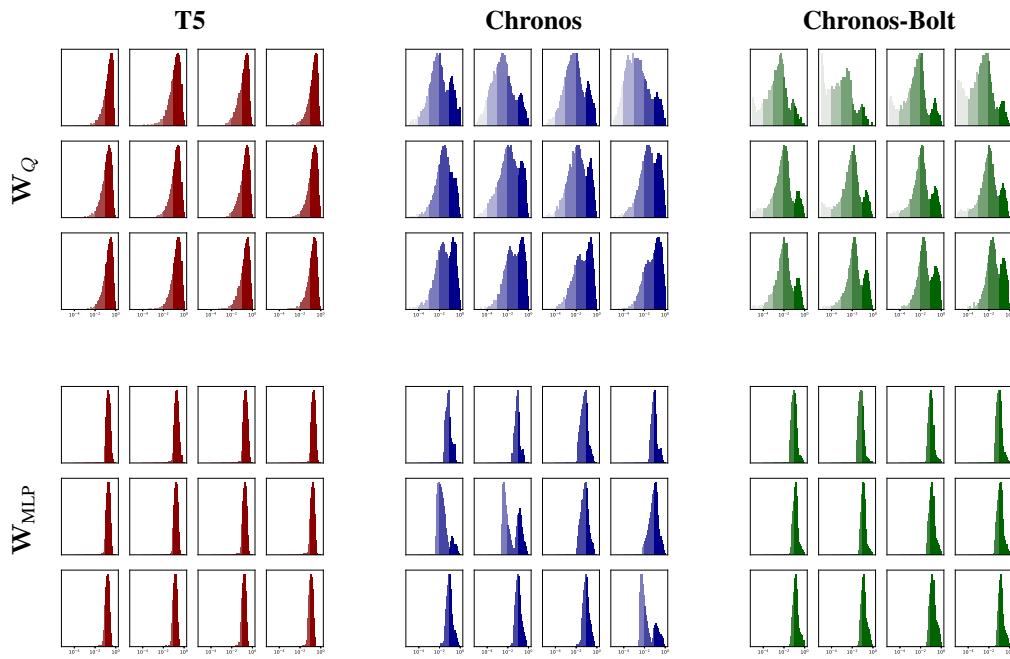
1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

1674 **G WATCHING THE WEIGHTS OF LARGE MODELS**
 1675

1676 The central claims made in this paper are heavily based on the low-rank structures of weight matrices
 1677 in large-scale time series foundation models. In this section, we provide more empirical analysis of
 1678 the singular values of these matrices.
 1679

1680 **G.1 COMPARING THE WEIGHTS OF TSFMS AND LLMs**
 1681

1682 For the first set of comparisons, we consider three models: T5, Chronos, and Chronos-Bolt. While
 1683 T5 is an LLM, Chronos and Chronos-Bolt are TSFMs. The three models we compare have the same
 1684 base size, and each model contains 12 encoder layers and 12 decoder layers. From each layer, we
 1685 take out a matrix \mathbf{W} , which is either the query projection matrix $\mathbf{W}_Q \in \mathbb{R}^{768 \times 768}$ or the first matrix
 1686 of the MLP layer $\mathbf{W}_{\text{MLP}} \in \mathbb{R}^{3072 \times 768}$. We apply an SVD to the weight matrix $\mathbf{W} = \mathbf{U}\Sigma\mathbf{V}^\top$ and
 1687 construct a histogram out of all relative singular values in $\text{diag}(\Sigma)/\Sigma_{1,1}$.
 1688



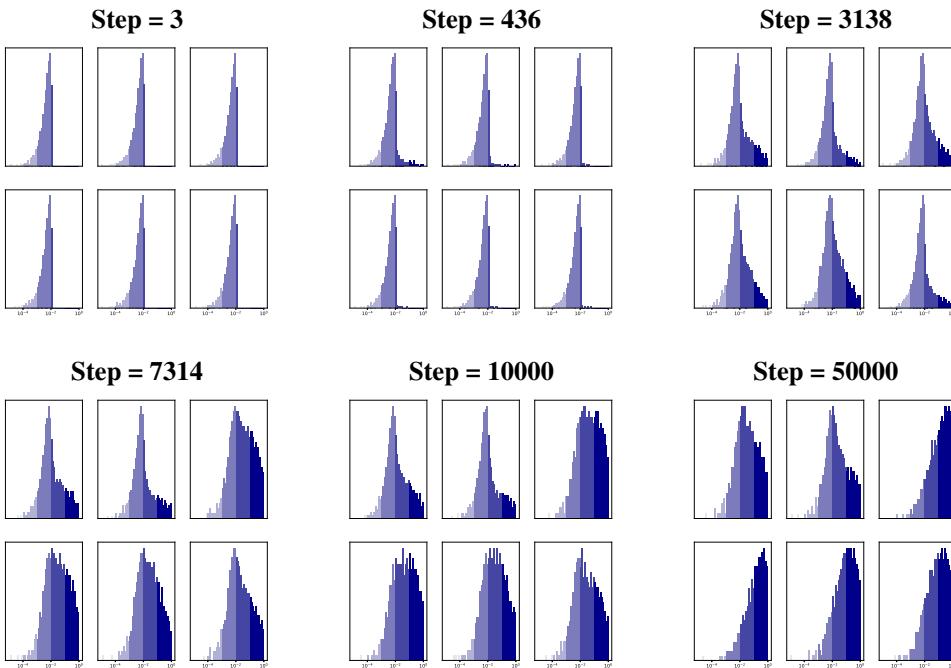
1711 **Figure 10:** The distribution of all relative singular values in a weight matrix \mathbf{W} in a T5, Chronos,
 1712 or Chronos-Bolt model. We use five progressively fainter face colors to indicate the range where the
 1713 relative singular values are in the $(1e-5, 1e-4]$, $(1e-4, 1e-3]$, $(1e-3, 1e-2]$, $(1e-2, 1e-1]$,
 1714 and $(1e-1, 1e0]$ ranges, respectively. Each model contains 12 encoder layers, and we arrange
 1715 the corresponding 12 weight matrices in row-major order. Note that all histograms are made on a
 1716 semilog-x scale.
 1717

1718 From Figure 10, we make some observations that align with the major claims made in the main
 1719 manuscript:
 1720

1. The relative singular values decay faster for Chronos and Chronos-Bolt and they decay much
 1721 slower for T5. The reason is that TSFMs leverage low-rank embeddings and do not require
 1722 high-rank attention matrices, which is not the case for T5, which inevitably needs a high-rank
 1723 embedding.
2. In Chronos and Chronos-Bolt, when the layers get deeper, we generally have a higher-rank struc-
 1724 ture. This aligns with the flow-of-ranks idea: as an input is pushed through more nonlinear
 1725 attention and MLP layers, its rank gets higher, and we need higher-rank attention matrices.
3. The attention matrices are generally more compressible than MLP matrices. This is not sur-
 1726 prising, because the attention matrices in these models are square matrices, making low-rank
 1727 structures much easier to emerge than in a rectangular MLP matrix.

1728 **G.2 WATCHING THE WEIGHTS OF TSFMS DURING TRAINING**
 1729

1730 Another mysterious finding that we used in the main manuscript is that the attention matrices in a
 1731 pretrained Chronos or Chronos-Bolt model all demonstrate low-rank structure; however, while The-
 1732 orem 3 only suggests that the attention matrices *can* be written in a low-rank form, it does not
 1733 preclude these matrices from having a high-rank form. To understand why the attention matrices
 1734 happen to exhibit a low-rank structure, we watch the training dynamics of these attention matrices,
 1735 where we pretrain a small Chronos model and record the six weight matrices in its six encoder layers
 1736 at a few different training steps. Unlike Figure 10, in Figure 11, we do not show the relative singular
 1737 values σ_j/σ_1 . The reason will be clear later.



1761 **Figure 11:** The distribution of all absolute singular values in a weight matrix \mathbf{W}_Q in a Chronos
 1762 model over training. The model contains 6 encoder layers, and we arrange the corresponding 6
 1763 weight matrices in row-major order. Note that all histograms are made on a semilog-x scale.

1764 We note that the initialization of Chronos typically relies on a very small scaling factor. That is, the
 1765 weights are all initialized to be small. From Figure 11, we see that the weight matrices are eventually
 1766 learned to be larger: that is, by looking at the absolute singular values instead of the relative ones,
 1767 we see that the leading singular value, i.e., the norm of the weight matrix \mathbf{W}_Q , eventually gets
 1768 larger. When they get larger, we note that the low-rank structure evolves. This “learning the leading
 1769 singular direction” interpretation is more plausible than if the weight matrix \mathbf{W}_Q is initialized large,
 1770 because it has no incentive to “forget the residual ranks” when it does not need to.
 1771

1772 **G.3 WATCHING THE WEIGHTS IN A MULTIHEAD ATTENTION**
 1773

1774 In Figure 9, we showed the “correlation matrix” of a three-head attention matrix \mathbf{W}_Q . In Figure 12,
 1775 we show more of these matrices when a different number of heads h . For each number of heads
 1776 h , we can see clearly h low-rank blocks on the diagonal, corresponding to the in-head low-rank
 1777 attention.
 1778

1779 One interesting observation we can make is that the off-diagonal parts of the heatmaps, while much
 1780 more yellowish than the diagonal parts, also have block structures. This is expected: if the rows in
 1781 each head are very colinear, then when you compare rows between two different heads, they should
 share a similar angle. That is, we have a duality that not only holds for attention matrices but also

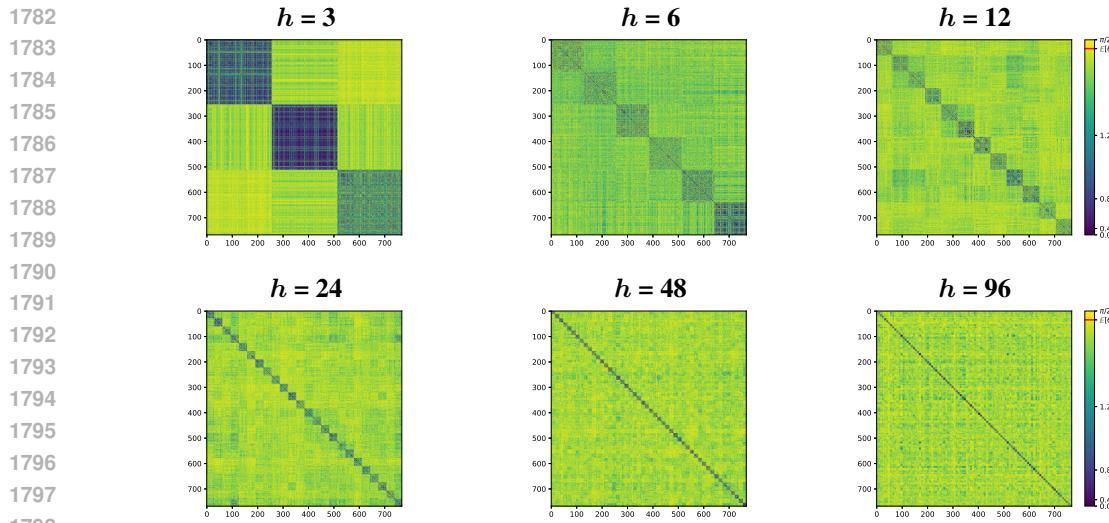


Figure 12: We pretrain a Chronos model with a different number of heads h . As h changes, we show the angle between every pair of rows of \mathbf{W}_Q in the first encoder layer.

holds in general: the darker a diagonal block is, the clearer the edges of off-diagonal blocks in the same row of blocks will be.

The off-diagonal parts of attention matrices, while much more orthogonal than the rows within each head, are still far from random. This means the heads themselves are also somewhat correlated, which is not surprising given that the row spaces of query, key, and value matrices are subspaces of the potentially low-dimensional row space of the input and that the models are trained with a single objective. There are two forces pulling against each other: a head-dependent random initialization, which “orthogonalizes” different heads, and a head-independent training objective that tries to align these heads.

There is a third trend (and the second duality) that is also very interesting: darker diagonal blocks seem to correspond to lighter off-diagonal blocks. There is one potential explanation for that: if the diagonal blocks are very dark, that means each head only attends to a tiny bit of information of the input. In order to obtain good results, other heads must incorporate the remaining bits of the input, resulting in very different sketchings and larger angles. On the other hand, if the diagonal blocks are brighter, that means it already contains a fair amount of information about the input, and lots of the input information is shared across heads. Hence, they should be much more correlated during training.

1836 H VISUALIZATION OF TOKENIZERS

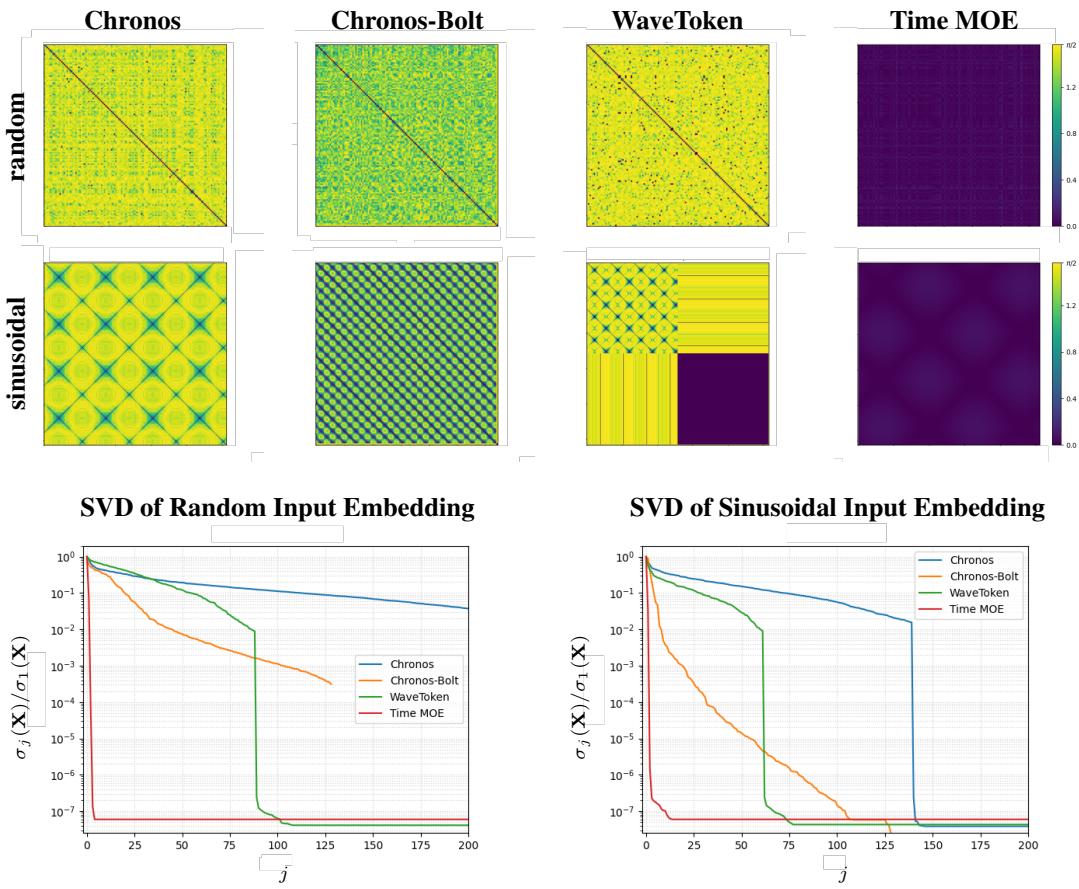
1838 In section 2, we consider a large corpus of inputs. Here, we perform two case studies to look into
 1839 how each of the embeddings works. In particular, we select two time-series input data:

- 1841 • A sinusoidal wave $\sin(t)$.
- 1842 • A random Gaussian input, where each entry is i.i.d. $\mathcal{N}(0, 1)$.

1844 Let $\mathbf{X} \in \mathbb{R}^{d \times L}$ be the embedded input. For each pair of vectors \mathbf{x}_i and \mathbf{x}_j of \mathbf{X} , we compute the
 1845 their correlation:

$$1846 \theta_{i,j} = \arccos \left(\frac{|\mathbf{x}_i^\top \mathbf{x}_j|}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_j\|_2} \right).$$

1848 These correlation matrices are shown in Figure 13. There are two interesting observations to make:
 1849 first, for WaveToken, the correlation matrix given a sinusoidal input has clearly four blocks —
 1850 corresponding to the low-frequency wavelets and high-frequency ones. For Time MOE, the angles
 1851 are so much darker because we apply a continuous embedding on a one-dimensional input space,
 1852 leading to an ultimately low-rank structure.



1881 **Figure 13:** The heat maps show the correlation matrix of the embedded input with a random context
 1882 or a sinusoidal one, using four different embedding strategies. The two line plots show the relative
 1883 singular values of the embedded matrix \mathbf{X} .

1884 We also observe that the numerical rank of the embedded input \mathbf{X} is generally higher when the
 1885 context is random than sinusoidal. This is not surprising either, because the temporal relationship
 1886 within a random context is much more complex than that in a sinusoidal one.

1890 I MORE ON THE CHEBYSHEV EMBEDDING

1891
 1892 In Figure 4, we show an experiment where we increase the rank of a fixed input embedding and
 1893 watch the numerical ranks of pretrained Chronos models with that embedding. Here, we further ex-
 1894 plain how we can control the rank of this fixed input embedding function using Chebyshev poly-
 1895 nomials. To motivate our design, we first consider how we can compute a rank-1 embedding function
 1896 $\Phi_1 : \mathbb{R} \rightarrow \mathbb{R}^d$. In this case, what Φ_1 needs to do is to map the real line linearly onto a one-
 1897 dimensional subspace of \mathbb{R}^d . That is, we can set $\Phi_1(x) = x\mathbf{u}$ for some fixed unit vector $\mathbf{u} \in \mathbb{R}^d$.

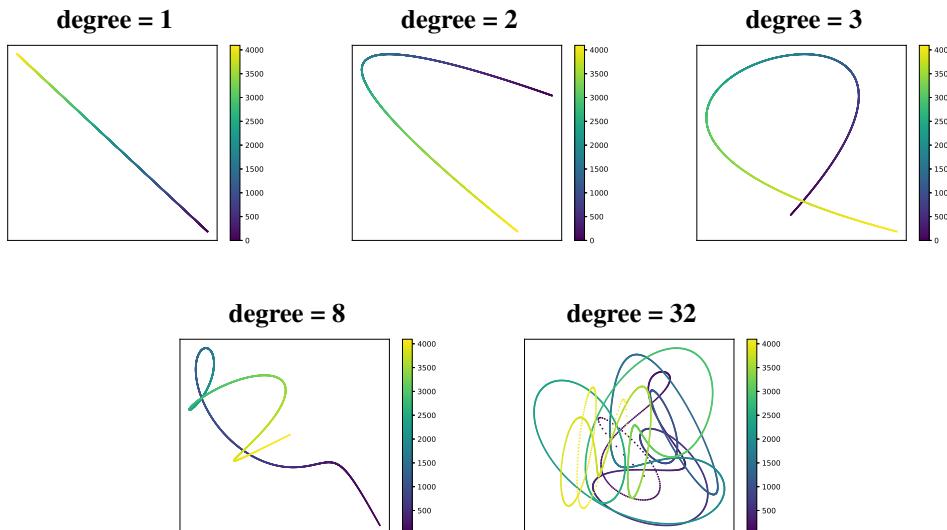
1898 Now, how can we use this idea to design a rank- k embedding $\Phi_k(x)$? One way to do that is by
 1899 considering the following extension of Φ_k :

$$1900 \quad \Phi_k : \mathbb{R} \rightarrow \mathbb{R}^d, \quad x \mapsto f_1(x)\mathbf{u}_1 + \cdots + f_k(x)\mathbf{u}_k.$$

1901 As long as the functions f_1, \dots, f_k are linearly independent, the image of the real line, $\Phi_k(\mathbb{R})$,
 1902 is a subset of a k -dimensional subspace of \mathbb{R}^d . The only question that remains is: how to choose
 1903 $f_1(x), \dots, f_k(x)$. Perhaps the easiest way to choose such a basis is by making them monomials:
 1904 $f_j(x) = x^j$. However, the monomial basis often suffers from many numerical stability issues and is
 1905 not ideal for the embedding, e.g., for a large j , the set $\{x, \dots, x^j\}$ is very ill-conditioned. To choose
 1906 a well-conditioned basis, we use Chebyshev polynomials, which are orthogonal on $[-1, 1]$. That
 1907 is, we set $f_j(x) = T_j(x/x_{\max})$, where x_{\max} is the maximum number considered in quantization,
 1908 which equals 15 in the case of Chronos. Given a set of points x_1, \dots, x_L to embed, we assemble
 1909 the embedded matrix $\mathbf{X} \in \mathbb{R}^{d \times L}$ as follows:

- 1910 1. Sample orthonormal random column vectors $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{R}^{d \times 1}$.
- 1911 2. Compute row vectors $\mathbf{f}_1, \dots, \mathbf{f}_k \in \mathbb{R}^{1 \times L}$, where the i th entry of \mathbf{f}_j is $T_j(x_i/x_{\max})$.
- 1912 3. Compute the outer product $\mathbf{X} = \sum_{j=1}^k \mathbf{u}_j \mathbf{f}_j$.

1913 We show the visualization of our Chebyshev embeddings in Figure 14, where as see that as we
 1914 increase the rank k , the embedding becomes more “complicated.”



1915
 1916 **Figure 14:** Visualization of the Chebyshev embeddings. For each embedding, we visualize the first
 1917 two dimensions out of the 768 in the hidden space. There are 4096 points embedded in the hidden
 1918 space, whose corresponding values in the input domain range from -15 to 15 .

1944
1945

J NUMERICAL EXPERIMENTS TO CORROBORATE OUR THEOREM

1946
1947
1948

In this appendix, we provide numerical experiments, simulated in MATLAB R2024b, that verify the theoretical statements we made in the main manuscript.

1949
1950

J.1 TWO NUMERICAL EXPERIMENTS ON THEOREM 3

1951
1952
1953
1954
1955

To verify Theorem 3, we fix the hidden dimension to be $d = 512$ and the vocabulary size $N = 4096$. We create a random embedded matrix $\mathbf{X} \in \mathbb{R}^{d \times N}$ with singular values $\mathbf{s} \in \mathbb{R}^d$ ranging from e^{-5} to e^0 and uniformly distributed on a logarithmic scale. This matrix is computed by randomly sampling an orthogonal matrix $\mathbf{U} \in \mathbb{R}^{d \times d}$ and a matrix $\mathbf{V} \in \mathbb{R}^{N \times d}$ with orthonormal columns, via QR-decomposing random matrices, and setting $\mathbf{X} = \mathbf{U} \text{diag}(\mathbf{s}) \mathbf{V}^\top$.

1956
1957
1958
1959
1960
1961
1962

Next, we randomly sample three attention matrices $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$, and for each reduced-order $\tilde{d} = 1, \dots, d - 1$, we use the constructive formulas given in the proof of Theorem 3 to compute the reduced matrices $\tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V$. We set the input to be the entire vocabulary matrix \mathbf{X} and compute the output of the original attention layer, defined by $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$, as well as that of the reduced one, defined by $\tilde{\mathbf{W}}_Q, \tilde{\mathbf{W}}_K, \tilde{\mathbf{W}}_V$. We evaluate the Frobenius norm of the difference between the two outputs.

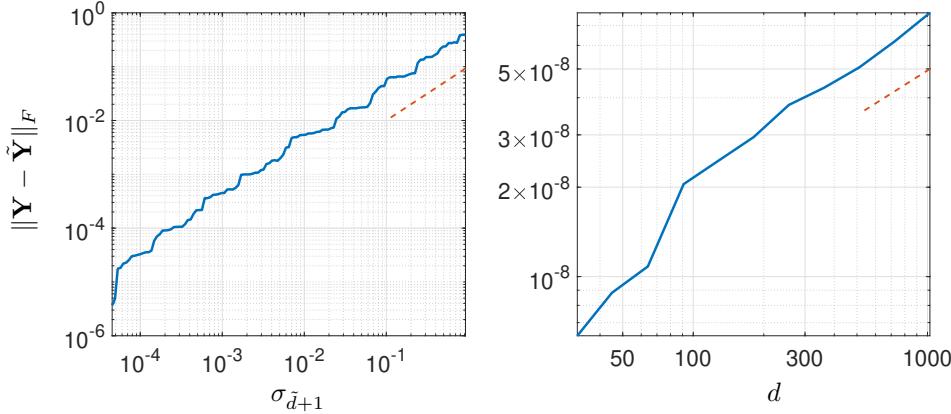
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
19771978
1979
1980
1981
1982
1983
1984

Figure 15: The left panel shows the relationship between the reduced-order \tilde{d} and the approximation error of a randomly sampled attention layer applied to a randomly sampled input matrix \mathbf{X} with controlled singular values $\sigma_1, \dots, \sigma_d$. The reference line has a slope of 1 in the log-log plot. The right panel shows the relationship between the dimension of the hidden space d and the approximation error of a fixed-degree reduced-order model. For each hidden space d , we randomly resample the attention matrices and embedded matrix \mathbf{X} , holding its leading singular values unchanged. The reference line has a slope of $1/2$ in the log-log plot.

1985
1986
1987
1988
1989
1990
1991

Figure 15 shows two controlled experiments by changing a different variable. On the left, we change the truncation degree \tilde{d} , which in turn controls the singular value $\sigma_{\tilde{d}+1}$. The reference line has a slope of 1, revealing a linear relationship between $\|\mathbf{Y} - \tilde{\mathbf{Y}}\|_F$ and $\sigma_{\tilde{d}+1}$, as indicated in Theorem 3. On the right, we change the size of the hidden space d , and the reference line has a slope of $1/2$, which verifies the \sqrt{d} factor in the statement of Theorem 3.

1992
1993

J.2 A NUMERICAL EXPERIMENT ON THEOREM 6

1994
1995
1996
1997

The essence of Theorem 6 is that a sparse multi-head sketching is more effective than a dense single-head sketching. We use an experiment to verify that. In our setting, we set $d = 2048$ and $L = 4096$, and we randomly sample an input matrix $\mathbf{X} \in \mathbb{R}^{d \times L}$ with exponentially decaying singular values, as shown in the left panel of Figure 16. Our sampling method is the same as the one outlined in the previous experiment.

We set our reduced rank to $\tilde{d} = 4$. If we just use a single head to sketch a rank- \tilde{d} space from the row space of \mathbf{X} , then its accuracy is lower-bounded by $\sigma_{\tilde{d}+1}$, which is still a large number. To explore the potential of sparse head-dependent sketching, we increase the number of heads from $h = 1$ to 128 and keep it an integral divisor of d . For each h , we use the sparse sketching; that is, we assemble a random matrix

$$\mathbf{W}_{2,h} = \begin{bmatrix} \mathbf{W}_{2,h}^{(1),\top} & \dots & \mathbf{W}_{2,h}^{(1),\top} \end{bmatrix}^\top, \quad \mathbf{W}_{2,h}^{(i)} \in \mathbb{R}^{\tilde{d} \times d}.$$

The matrix $\mathbf{W}_{2,h}$ is sparse in the sense that it has exactly $d_h = d/h$ non-zero entries, whose positions are randomly chosen, and for each non-zero position, we sample its value i.i.d. from $\mathcal{N}(0, 1)$. Note that given this sparse design, $\mathbf{W}_{2,h}$ has the same number of non-zero entries for any h .

Now, our question is: what is the difference between the row space of $\mathbf{W}_{2,h}\mathbf{X}$ and that of \mathbf{X} ? In other words, how good is the sketching using $\mathbf{W}_{2,h}$? To this end, we clearly have that $\mathcal{R}(\mathbf{W}_{2,h}\mathbf{X}) \subset \mathcal{R}(\mathbf{X})$, where we use the notation \mathcal{R} for the row space of a matrix. Hence, the remaining question is how much in $\mathcal{R}(\mathbf{X})$ is not filled by $\mathcal{R}(\mathbf{W}_{2,h}\mathbf{X})$. This can be measured by projecting $\mathcal{R}(\mathbf{X})$ onto $\mathcal{R}(\mathbf{W}_{2,h}\mathbf{X})$ and measure the loss by the projection:

$$d(\mathcal{R}(\mathbf{X}), \mathcal{R}(\mathbf{W}_{2,h}\mathbf{X})) = \|\text{orth}(\mathbf{X}^\top \mathbf{W}_{R,h}^\top) \text{orth}(\mathbf{X}^\top \mathbf{W}_{R,h}^\top)^\top \mathbf{X}^\top - \mathbf{X}^\top\|_2. \quad (30)$$

We show this distance in the right panel of Figure 16, and we relate this to the singular values of \mathbf{X} . In that sense, since the vertical rules in the left panel are almost evenly spaced, it shows that the “effective rank” of $\mathcal{R}(\mathbf{W}_{2,h}\mathbf{X})$ grows proportionally with respect to h , indicating that the quality of a sparse multi-head sketching is comparable to the quality of a dense multi-head sketching.

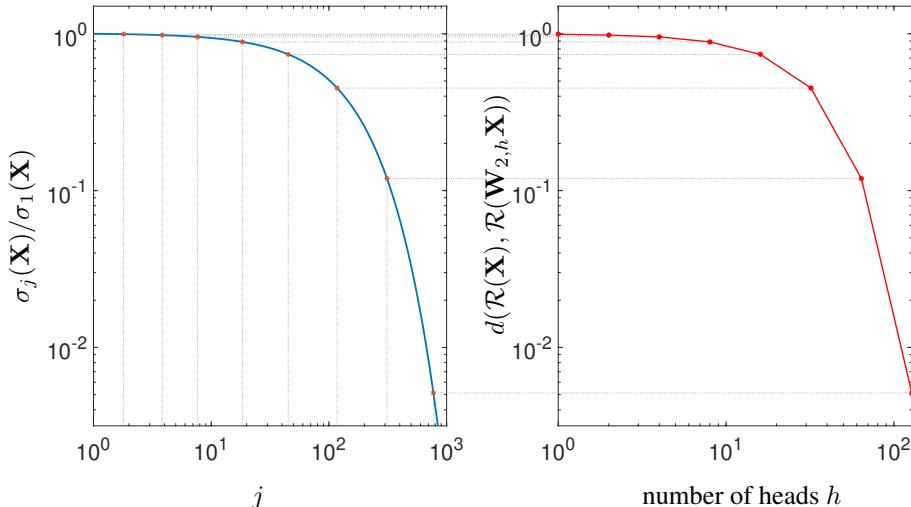


Figure 16: On the left, we show the singular values of the input matrices. On the right, we compute the “sketching error” defined in eq. (30), as we change the number of heads in a sparse sketching. We map the sketching errors back to the singular value plot to indicate the index j such that our sketching error achieves an error below σ_{j+1} .

J.3 A NUMERICAL EXPERIMENT ON THEOREM 4

Our final numerical experiment considers the flow-of-ranks. Theorem 4 suggests that a larger number of heads also facilitates the flow-of-rank, and this is hard to validate empirically with pretrained Chronos models. In our targeted numerical experiment, we fix an input matrix \mathbf{X} with predefined exponentially decaying singular values, which is, again, randomly sampled from the product of a random orthogonal matrix, a predefined diagonal matrix, and the transpose of a random matrix with orthonormal columns (see the previous two subsections). Then, we randomly select $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ and compute the output of the attention mechanism defined by these three matrices together with a number of heads parameter h . Figure 17 gives us that the output matrix

$$\mathbf{Y} = \text{MH-Attention}(\mathbf{X}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) + \mathbf{X}$$

2052 is higher-rank as h increases.
 2053

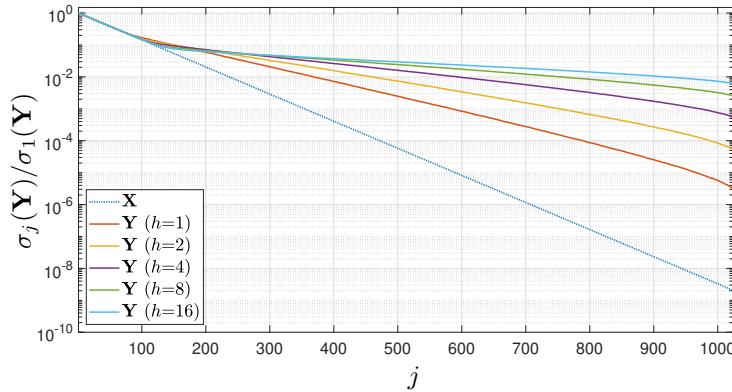


Figure 17: Relative singular values of the input matrix \mathbf{X} and output matrices $\mathbf{Y} = \text{MH-Attention}(\mathbf{X}; \mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, h) + \mathbf{X}$ as we change the number of heads h . We see that the numerical rank of \mathbf{Y} increases with h .

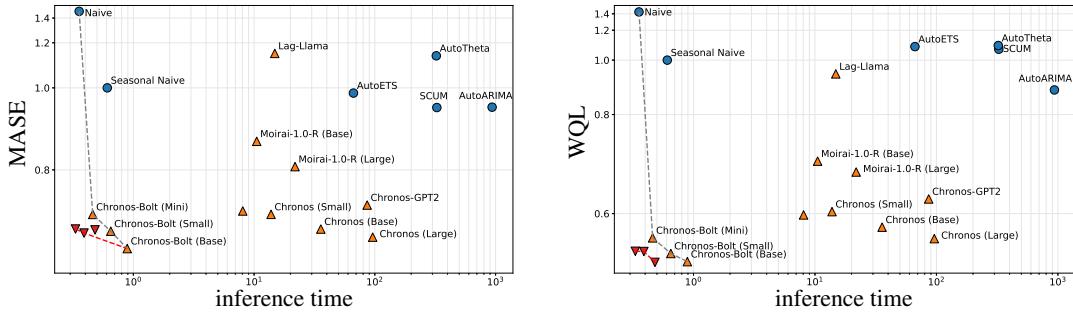
2054
 2055
 2056
 2057
 2058
 2059
 2060
 2061
 2062
 2063
 2064
 2065
 2066
 2067
 2068
 2069
 2070
 2071
 2072
 2073
 2074
 2075
 2076
 2077
 2078
 2079
 2080
 2081
 2082
 2083
 2084
 2085
 2086
 2087
 2088
 2089
 2090
 2091
 2092
 2093
 2094
 2095
 2096
 2097
 2098
 2099
 2100
 2101
 2102
 2103
 2104
 2105

2106 K ADDITIONAL EXPERIMENTS AND DISCUSSIONS

2108 **Pretraining a Compressed Chronos-Bolt.** In addition to Chronos, we also pretrain a compressed
 2109 Chronos-Bolt model. To show the promise of compression, we start with an already-small Chronos-
 2110 Bolt (small) model, based on the T5 (small) architecture. The table and figures in Table 4 can be
 2111 read in the same way as those in Figure 7. In particular, we see that for both MASE and WQL, the
 2112 compressed Chronos-Bolt models completely form the Pareto frontier on our evaluation benchmark.
 2113 That is, given any local method or pretrained foundation model, there exists a compressed Chronos-
 2114 Bolt model that is simultaneously faster and more accurate.

2115 **Table 4:** Results of pretraining a compressed Chronos-Bolt (small) model. We compare the performance scores *relative to the original pretrained model*. The last row is the baseline.

\tilde{d}_0	α	Inference Time		In-Domain		Zero-Shot	
		Space	WQL ↓	MASE ↓	WQL ↓	MASE ↓	
2	0.25	0.517	0.802	1.005	1.004	1.013	0.999
3	0.50	0.601	0.821	1.005	0.996	1.009	1.001
5	0.70	0.740	0.840	0.980	1.003	0.991	1.010
64	0.00	1.000	1.000	1.000	1.000	1.000	1.000



2135 While it is hard to find a rule of thumb that works for any task, we propose a guideline that should
 2136 work fine in most cases. Let D be the number of layers and d the hidden size. Set the per-layer
 2137 target rank

$$2138 \quad r_\ell = r_1 + (r_D - r_1) \left(\frac{\ell-1}{D-1} \right)^\alpha, \quad \ell = 1, \dots, D,$$

2139 where we can set r_{\min} = the median numerical rank of a small sample of input embeddings (or
 2140 16 if unknown), $r_{\max} = d/2$, and $\gamma = 0.5$. This design guarantees that r_ℓ grows smoothly and
 2141 monotonically from r_{\min} to r_{\max} (concave in depth), aligning capacity with the observed flow-of-
 2142 ranks while keeping early layers compact.

2144
 2145
 2146
 2147
 2148
 2149
 2150
 2151
 2152
 2153
 2154
 2155
 2156
 2157
 2158
 2159

2160 **L THE USE OF LARGE LANGUAGE MODELS (LLMs)**

2161
2162 LLMs are used for polishing the writing and word choices of a few sections in the main text. They
2163 are not used in the conceptualization and implementation of research.

2164

2165

2166

2167

2168

2169

2170

2171

2172

2173

2174

2175

2176

2177

2178

2179

2180

2181

2182

2183

2184

2185

2186

2187

2188

2189

2190

2191

2192

2193

2194

2195

2196

2197

2198

2199

2200

2201

2202

2203

2204

2205

2206

2207

2208

2209

2210

2211

2212

2213