

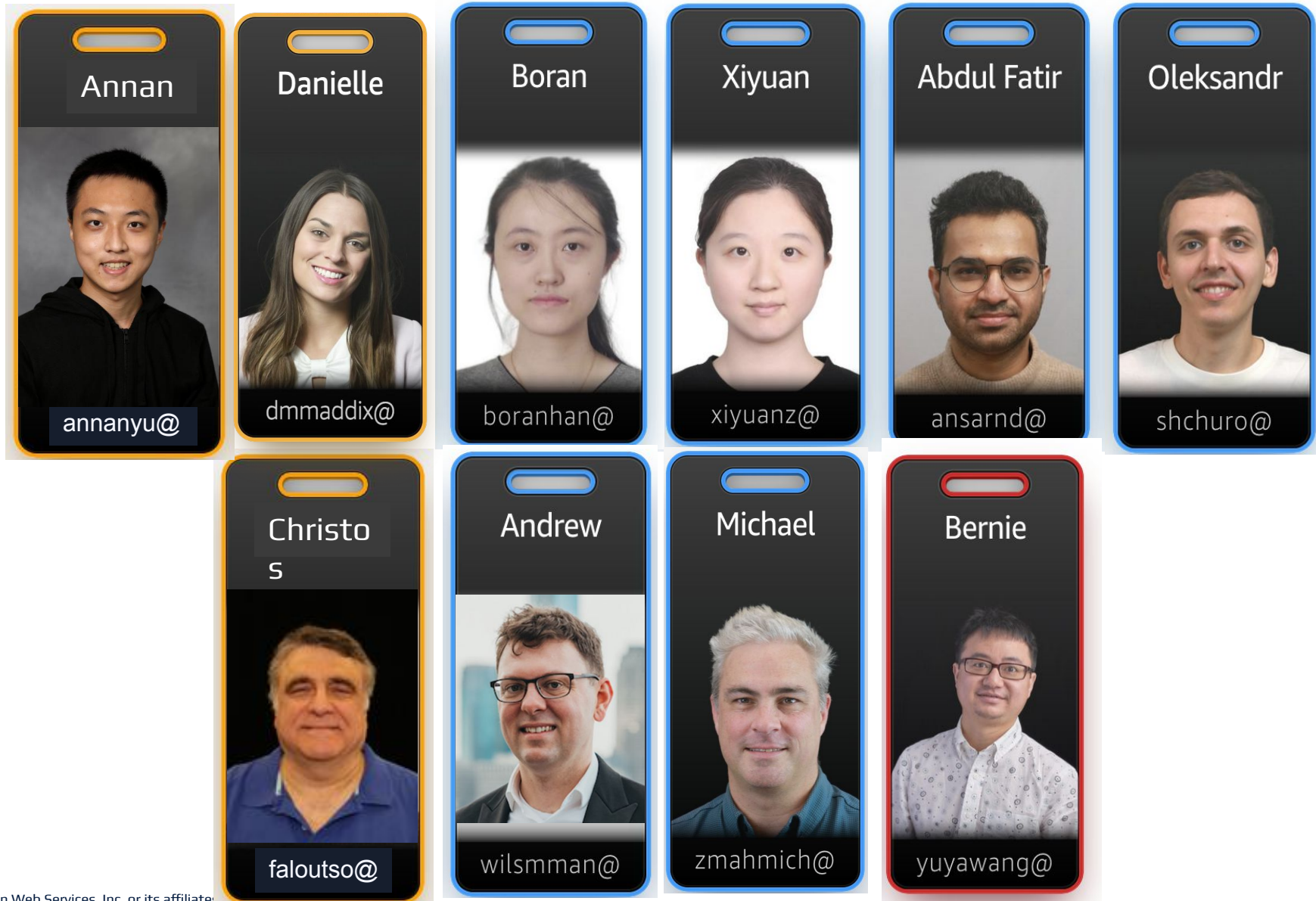
Understanding the Bitter Lesson in Time Series Foundation Models

Danielle Maddix Robinson

Senior Applied Scientist, AWS AI



Collaborators

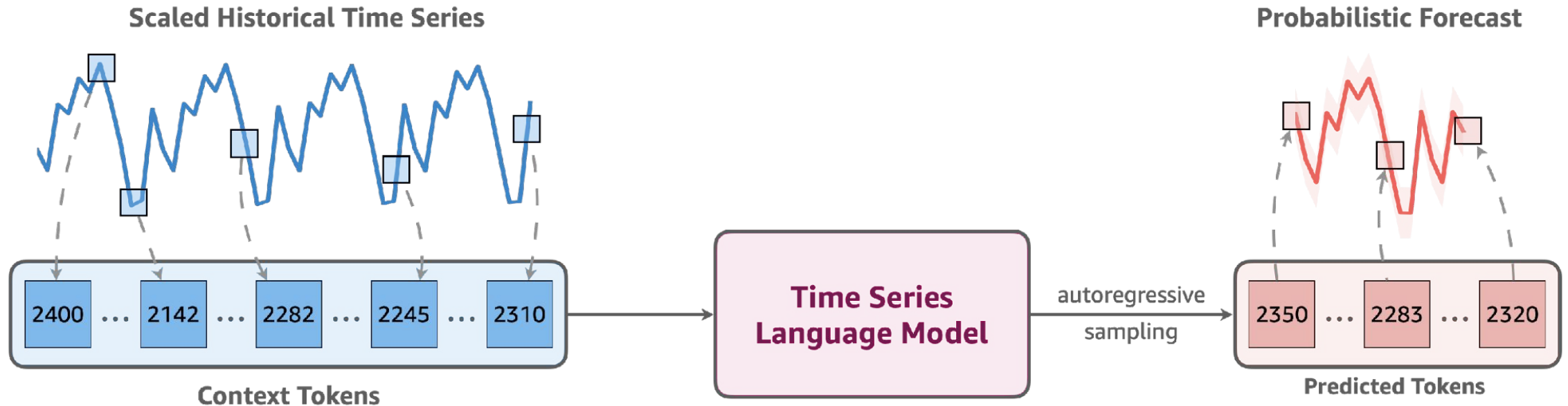


Chronos

A language modeling framework for time series data

that encodes time series into discrete tokens
and trains a language model on them

- ✓ probabilistic by design
- ✓ requires no changes to the language model architecture or training procedure



Ansari, A.F., et al., "Chronos: Learning the Language of Time Series", TMLR, 2024.

Baselines

Pretrained Models

single model used across all tasks

- LLMTime
- ForecastPFN
- LagLlama
- Moirai

CHRONOS
goes here

Task-specific Models

separate model trained/fine-tuned for each task

- PatchTST
- DeepAR
- WaveNet
- TFT
- DLinear
- NBEATS
- NHiTS
- GPT4TS

Local Models

separate model for each time series

- Naive
- SeasonalNaive
- AutoETS
- AutoTheta
- AutoARIMA

Benchmarks

Benchmark I

15 in-domain datasets for **CHRONOS**

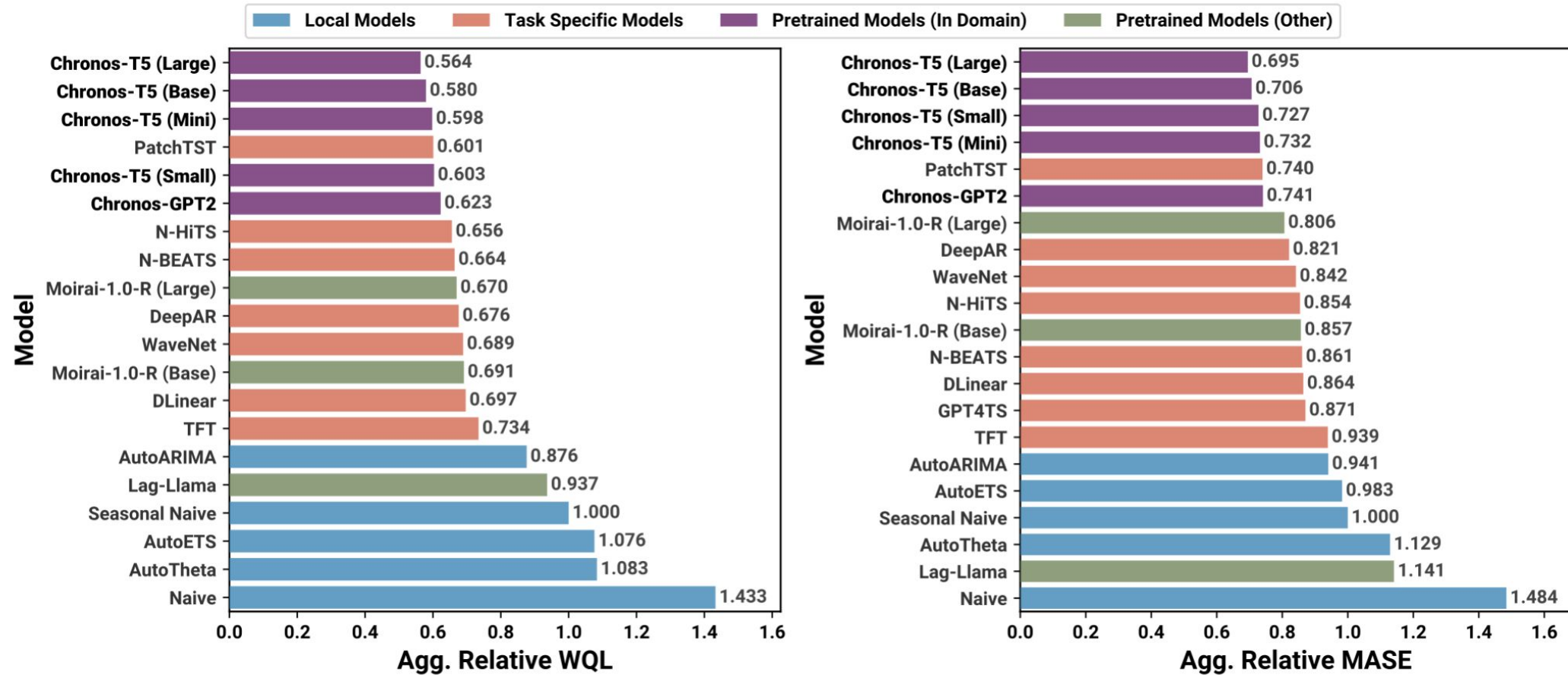
- Electricity (15 Min.)
- Electricity (Hourly)
- Electricity (Weekly)
- KDD Cup 2018
- M4 (Daily)
- M4 (Hourly)
- M4 (Monthly)
- M4 (Weekly)
- Pedestrian Counts
- Taxi (30 Min.)
- Uber TLC (Hourly)
- Uber TLC (Daily)
- Rideshare
- Temperature-Rain
- London Smart Meters

Benchmark II

27 zero-shot datasets for **CHRONOS**

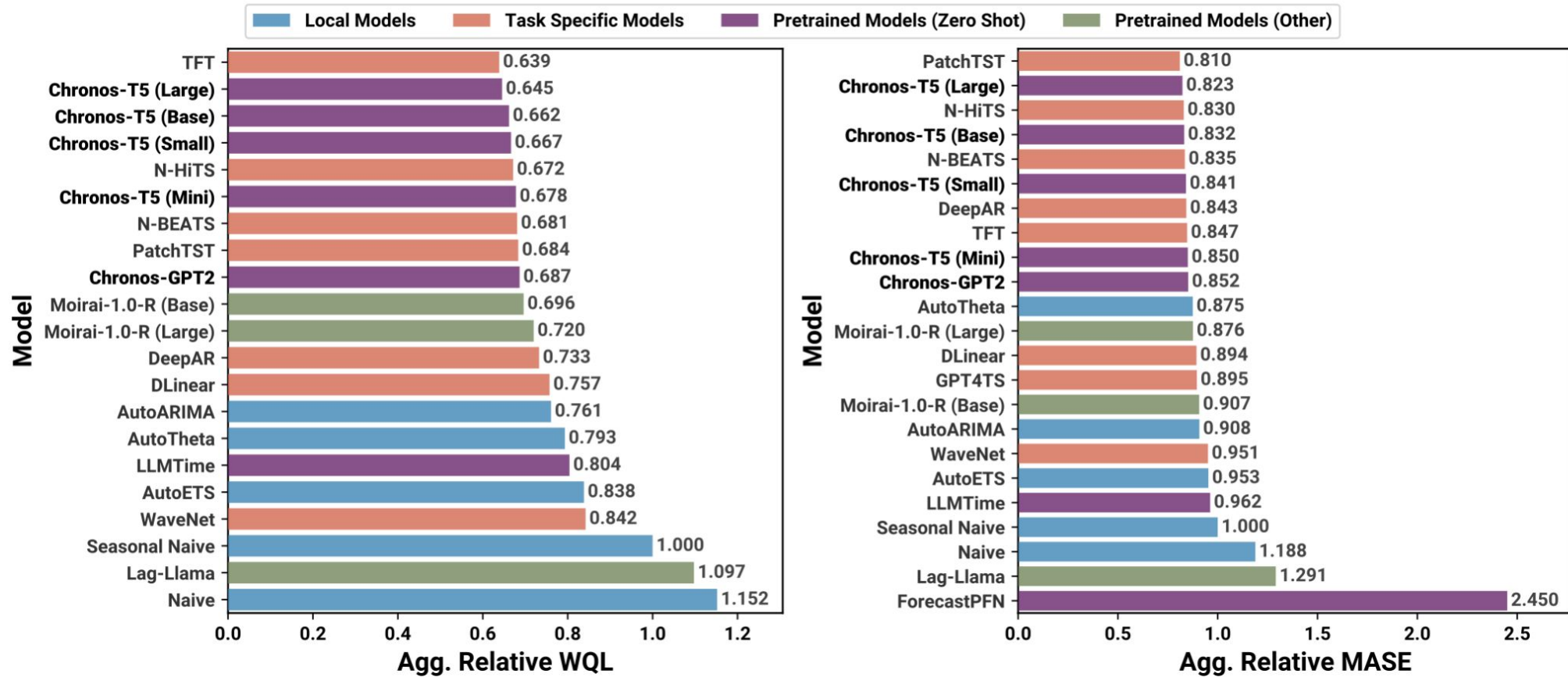
- Australian Electricity
- Car Parts
- CIF 2016
- Covid Deaths
- Dominick
- ERCOT Load
- ETT (15 Min.)
- ETT (Hourly)
- Exchange Rate
- FRED-MD
- Hospital
- M1 (Quarterly)
- M1 (Monthly)
- M1 (Yearly)
- M3 (Monthly)
- M3 (Quarterly)
- M3 (Yearly)
- M4 (Quarterly)
- M4 (Yearly)
- M5
- NN5 (Daily)
- NN5 (Weekly)
- Tourism (Monthly)
- Tourism (Quarterly)
- Tourism (Yearly)
- Traffic
- Weather

Chronos: In-domain Results



In-domain: 15 datasets that were part of the training corpus of **CHRONOS**

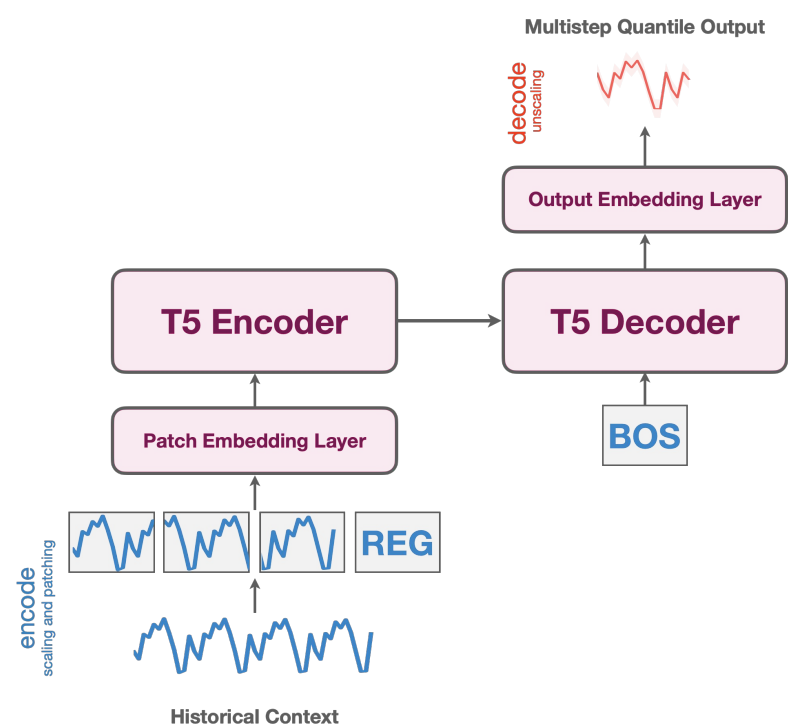
Chronos: Zero-shot Results



Zero-shot: 27 datasets not seen by CHRONOS during training

Chronos-Bolt ⚡

More accurate and 250x faster than the original Chronos models



	Chronos-Bolt	Chronos
Input (tokens)	Patches	Individual observations
Output (forecast)	Multi-step quantile forecast	Autoregressive sampling
Loss function	Quantile loss	Cross-entropy loss
Context length	2048	512
Inference device	CPU or GPU	GPU

Ansari, A.F., et al., “Fast and accurate zero-shot forecasting with Chronos-Bolt and AutoGluon”, AWS Technical Report, 2025.



Chronos-Bolt ⚡ : 250x faster than Chronos



Chronos: zero-shot results

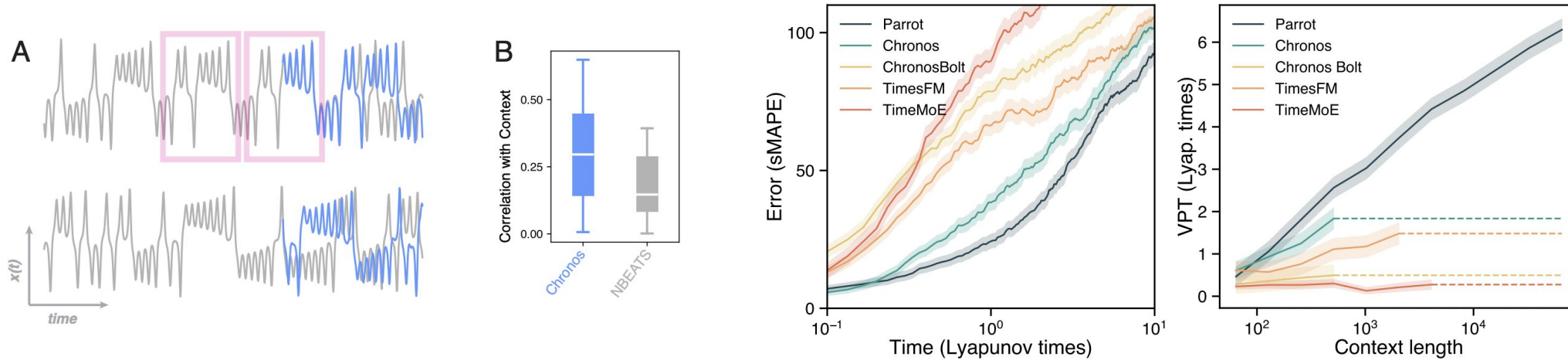
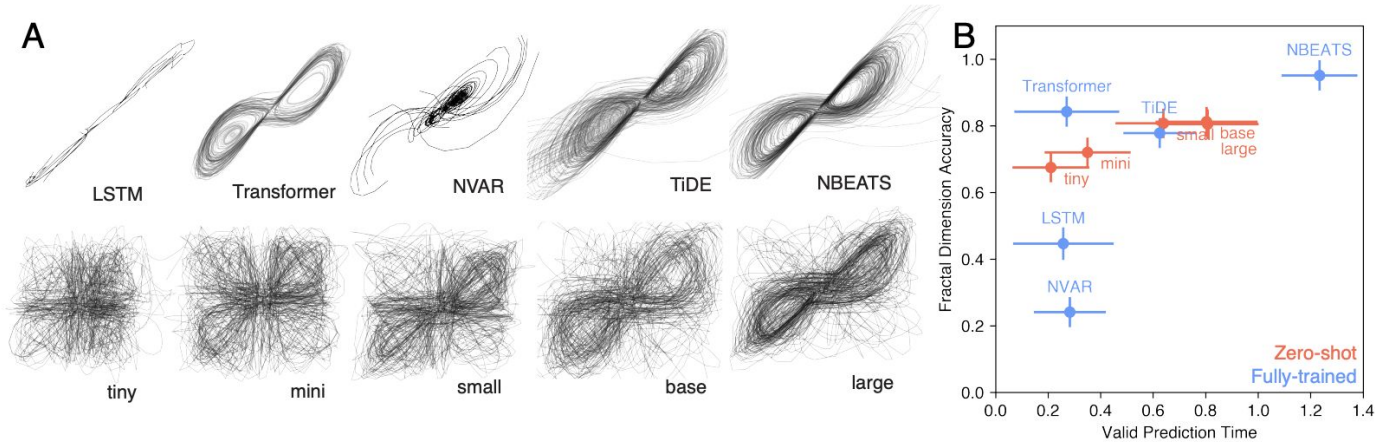
Forecast Error

Chronos-Bolt: zero-shot results



Forecast Error

Long-term Behavior on Chaotic Systems



Zhang, Y. et al. "Zero-shot Forecasting of Chaotic Systems," ICLR, 2025.

Zhang, Y. et al., "Context parroting: A simple but tough-to-beat baseline for foundation models in scientific machine learning", arXiv preprint arXiv:2505.11349, 2025.

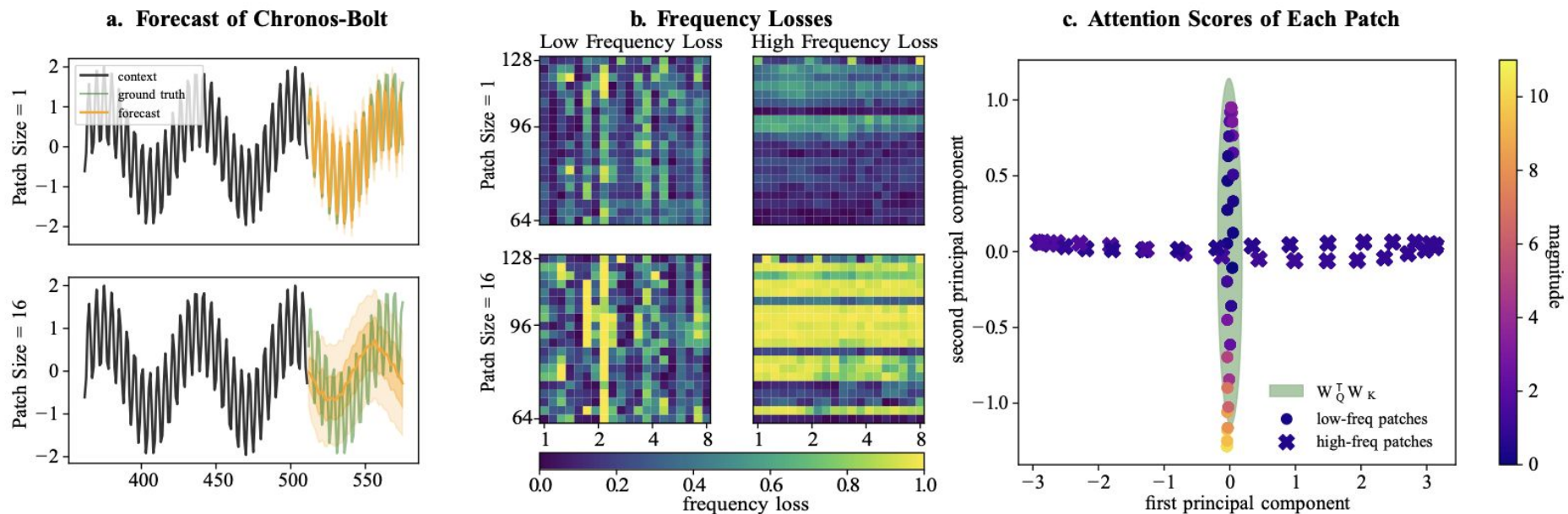
Design Choices of TSFMs

Yu, A. et al., “Understanding the Implicit Biases of Design Choices for Time Series Foundation Models”, arXiv preprint arXiv:2510.19236, Under Review, 2025.

Inductive Biases Overview

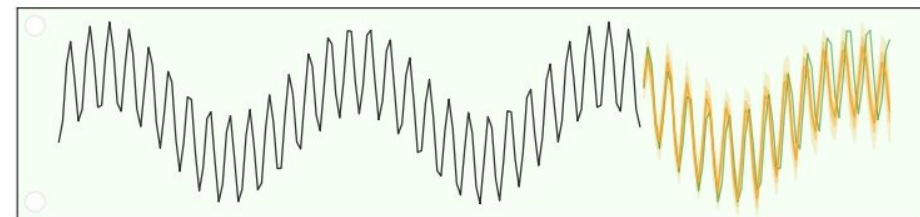
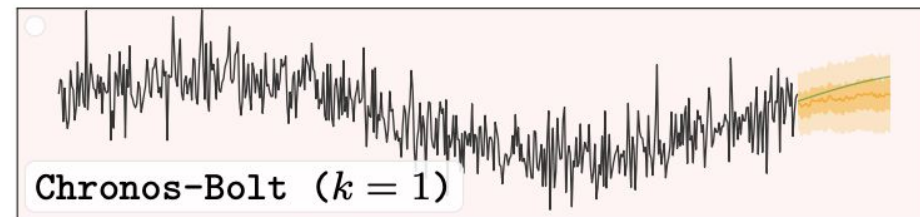
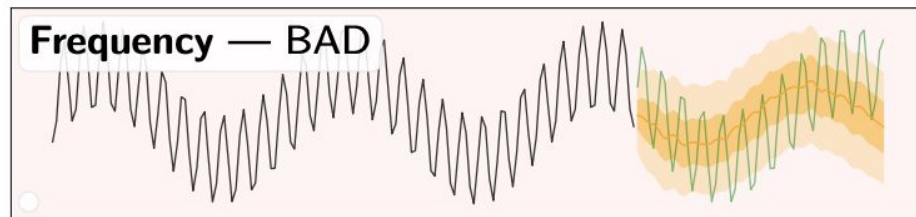
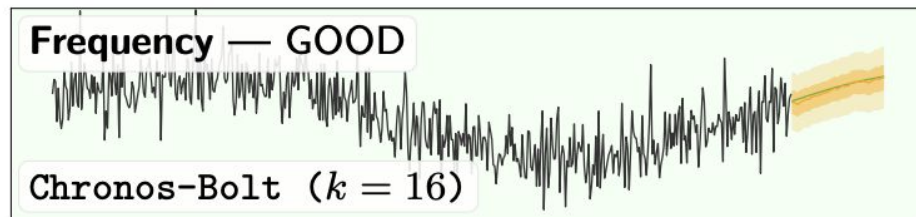
How Do TSFMs Learn Time?

Temporal Frequency Bias

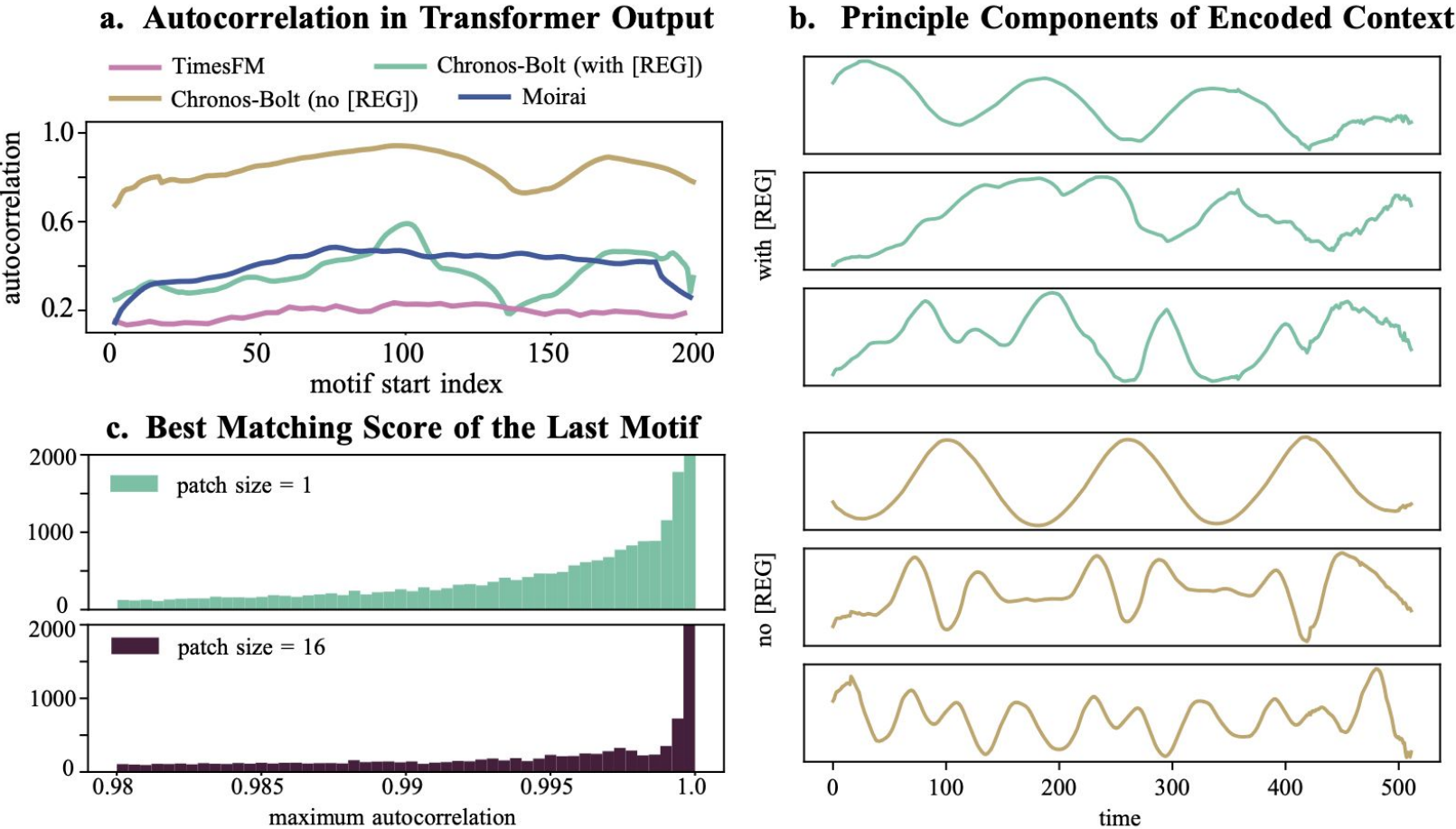


- Models with larger patch size k captures only the low frequency mode
- Chronos-Bolt $k = 16$ fails at capturing high-frequency information
- Attention scores are heavily dominated by the low-frequency patches

Frequency Bias: Good or Bad?

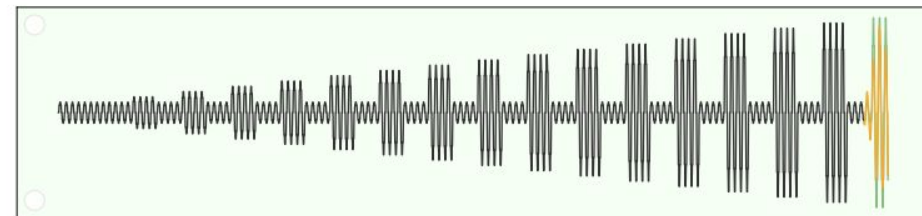
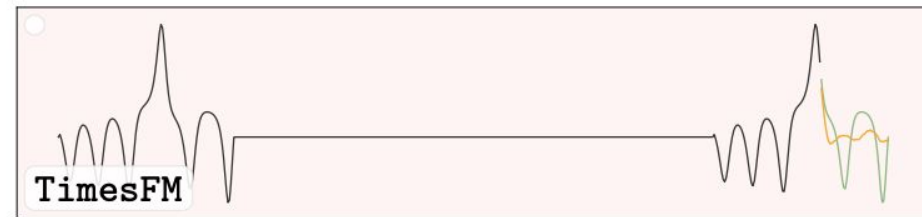
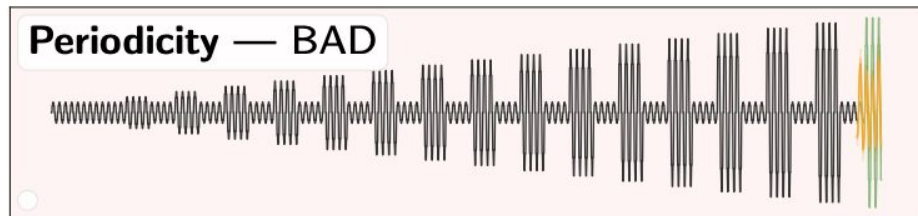


Temporal Periodicity Bias



Controlled by alignment of patch size k with underlying recurrent motifs

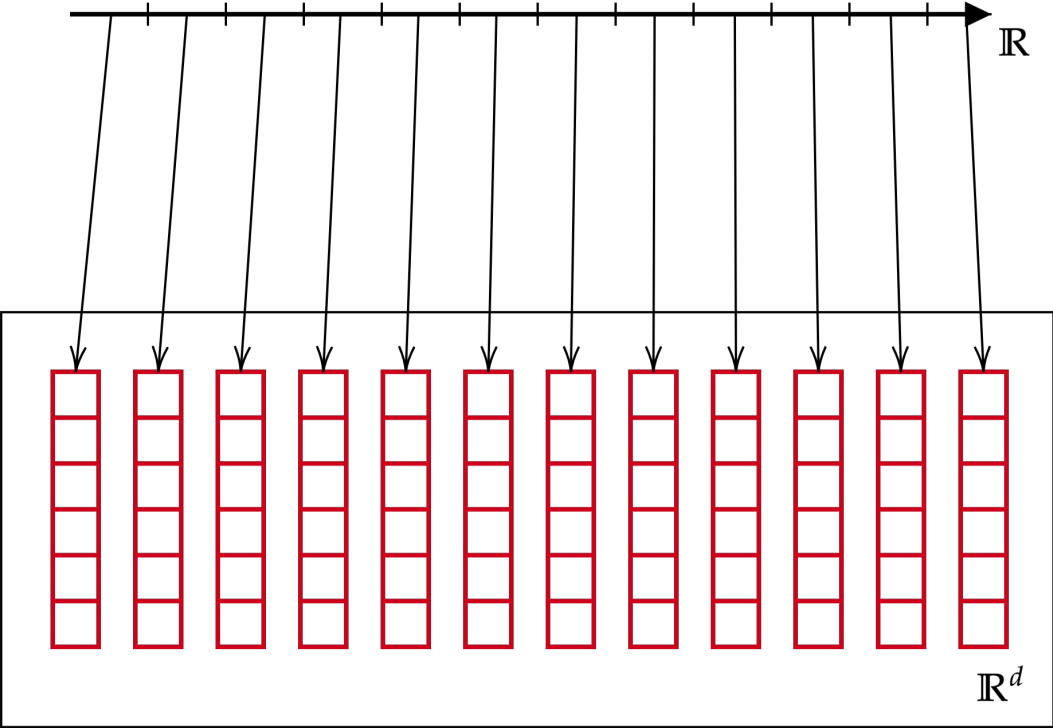
Periodicity Bias: Good or Bad?



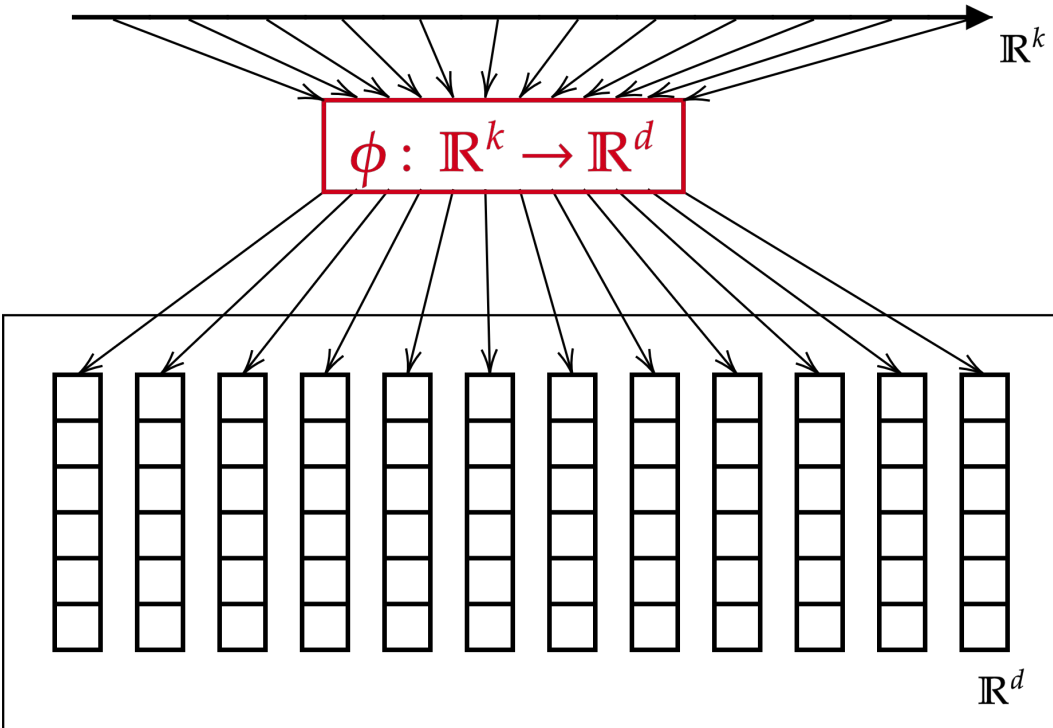
How Do TSFMs Learn Geometry?

Design Choice: Embedding Type

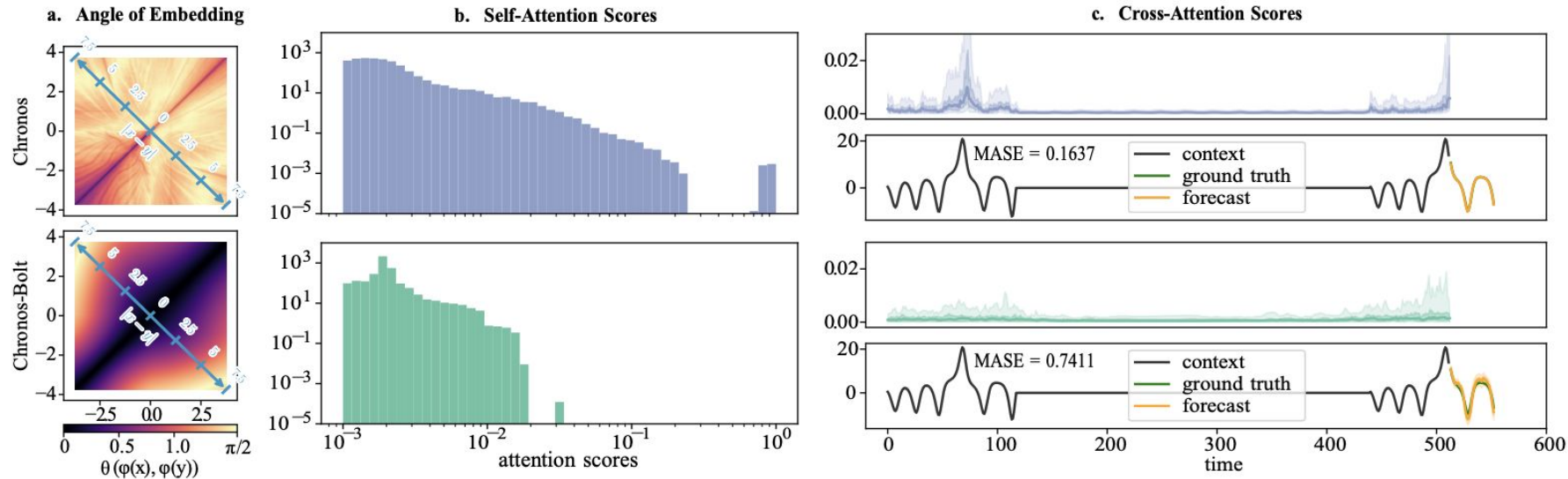
Quantization



Continuous Embedding

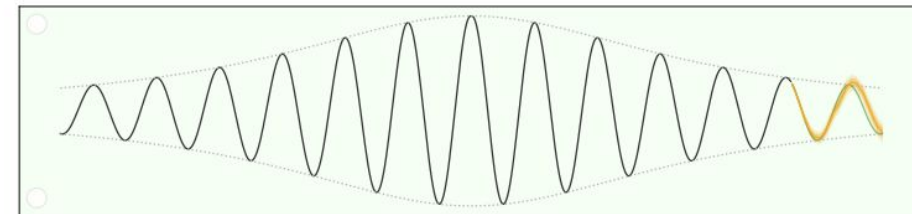
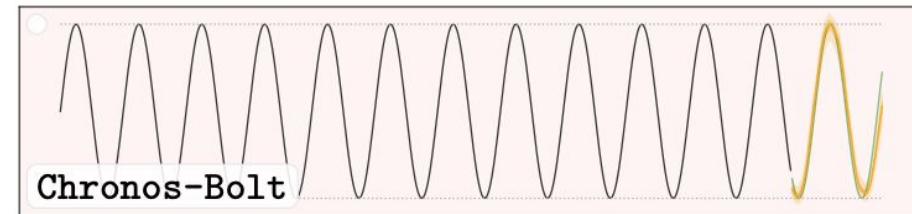
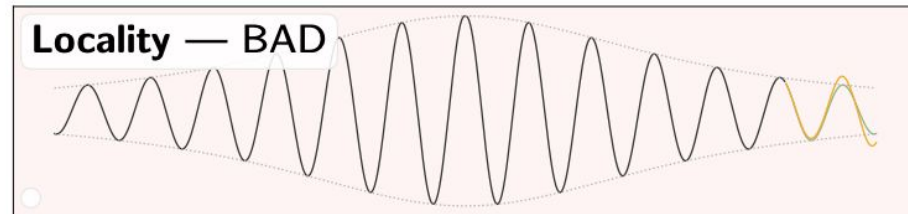
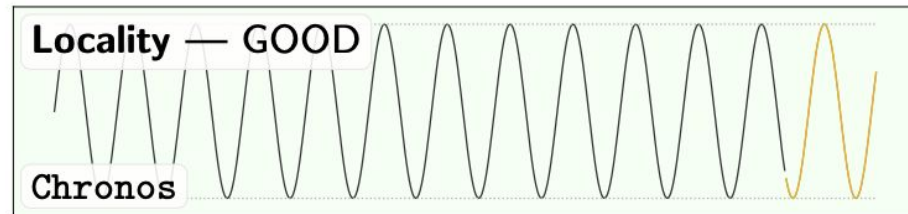


Geometric Angular Bias

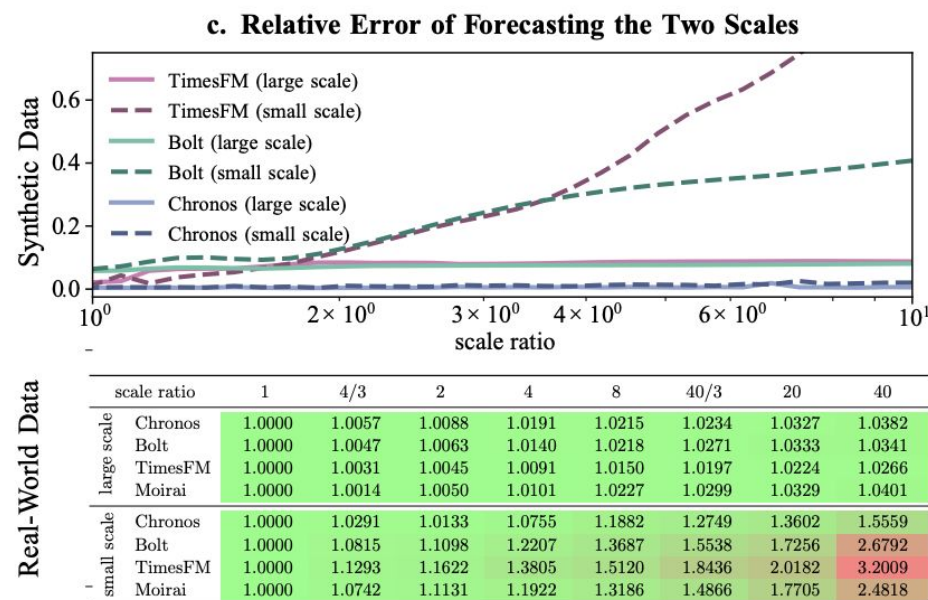
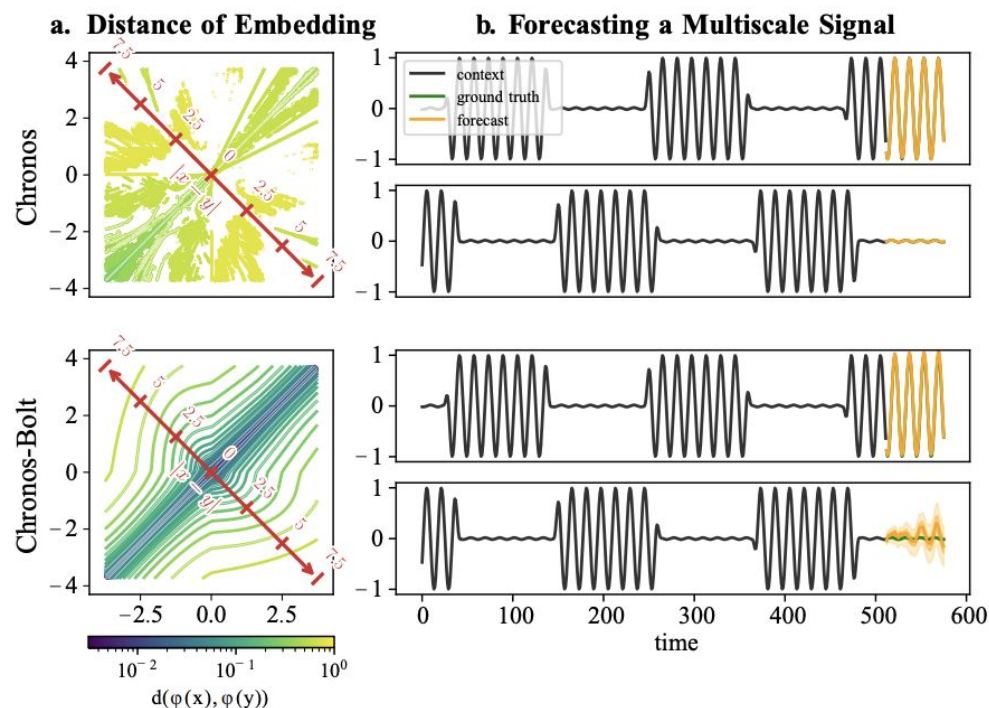


- Attention scores of Chronos are bimodal - tokens place high attention weights on their neighbors (locality)
- Attention scores of Chronos-Bolt are more evenly distributed (global)
- Chronos achieves a significantly lower MASE than Chronos-Bolt on a context formed by repeating a motif of a chaotic system by its "parroting"

Angular Bias: Good or Bad?

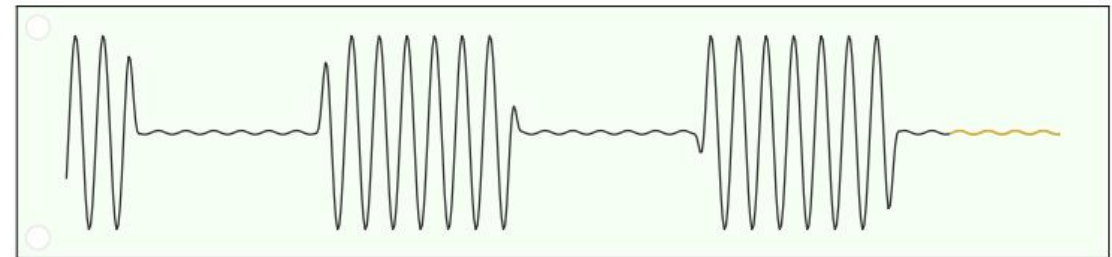
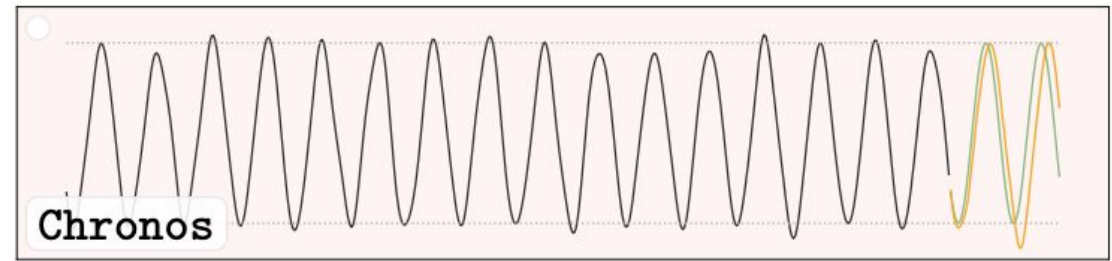
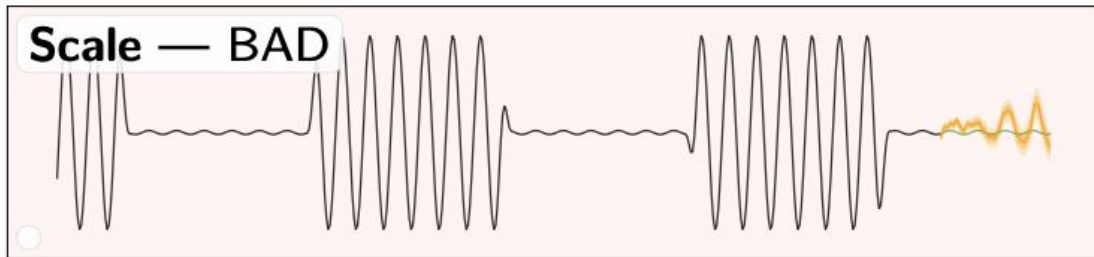
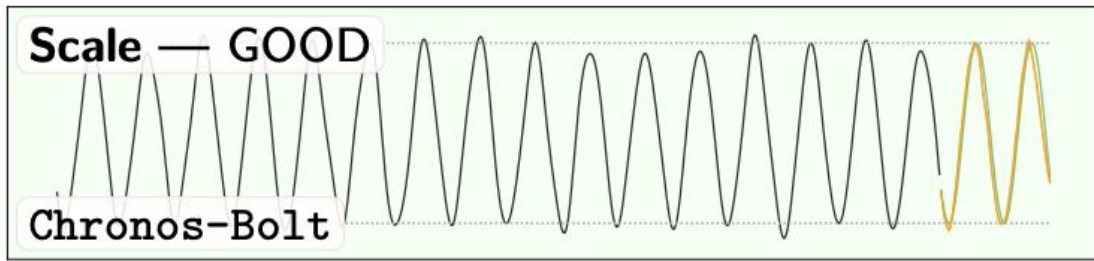


Geometric Distance Bias

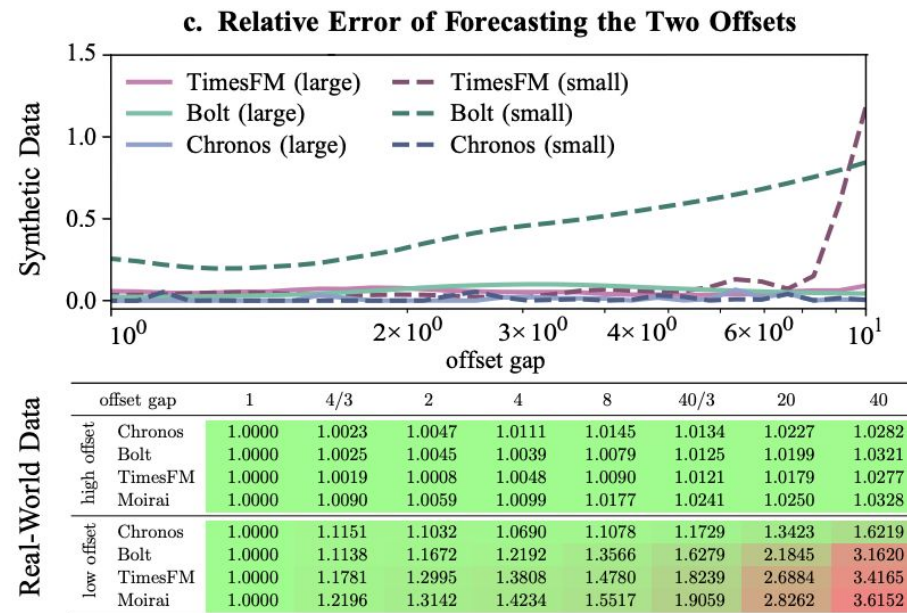
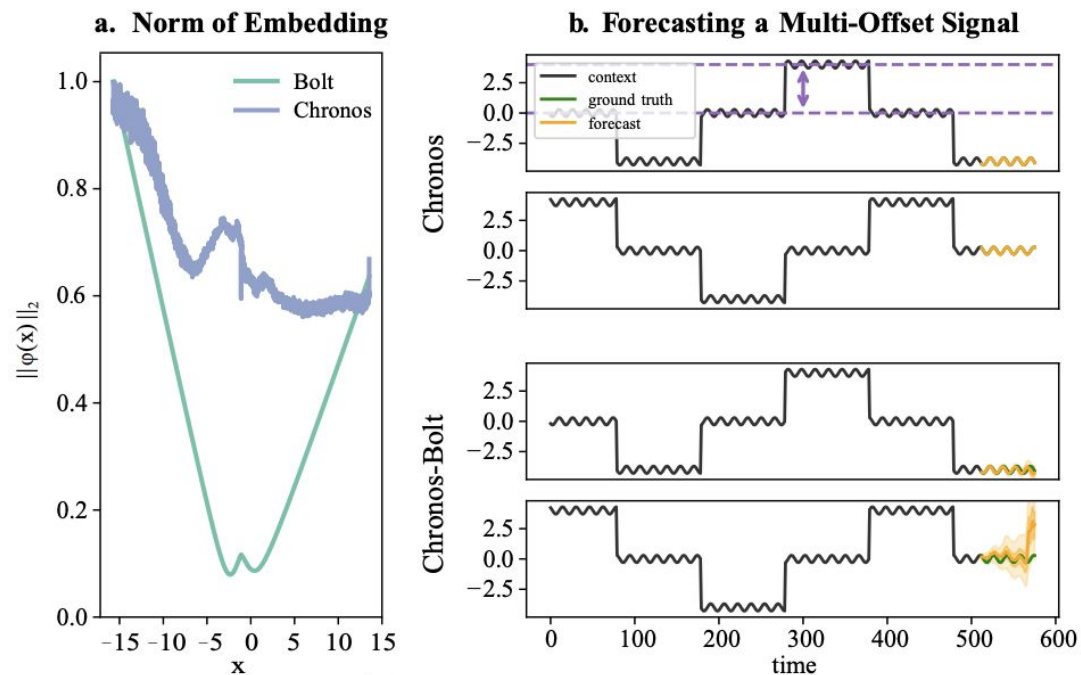


- Quantization-based embedding in Chronos magnifies smaller scales, which makes nearby numbers appear more distinct in hidden space
- Continuous-embedding in Chronos-Bolt maps nearby numbers to nearby vectors in hidden space
- For multi-scale structure, Chronos is more sensitive and better at learning fine-scale patterns than Chronos-Bolt

Distance Bias: Good or Bad?

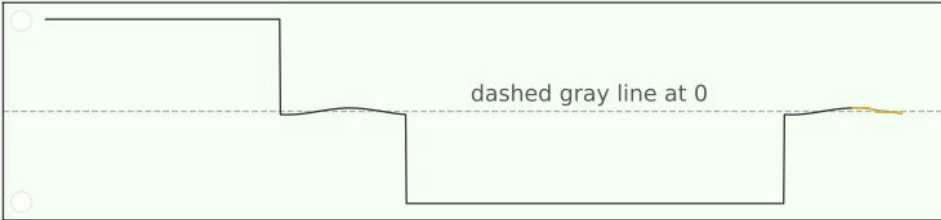
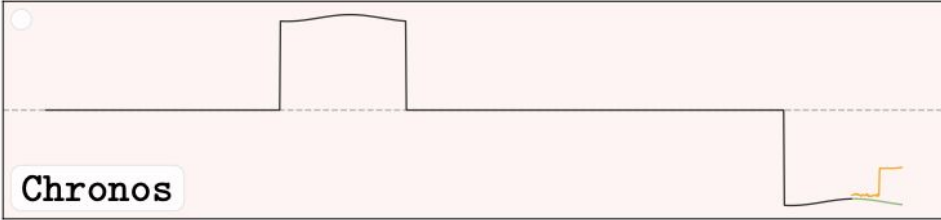
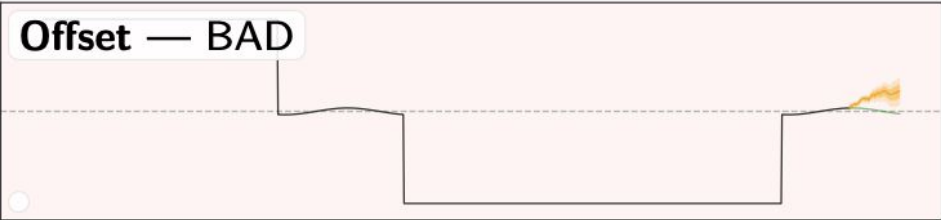
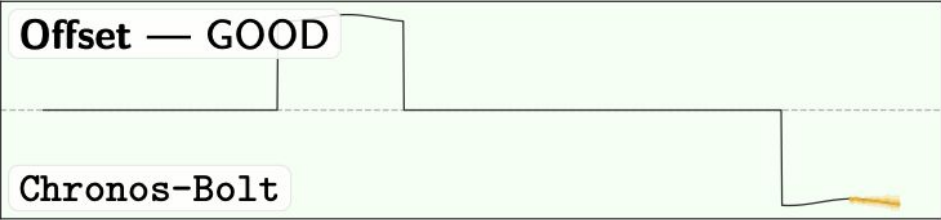


Geometric Norm Bias



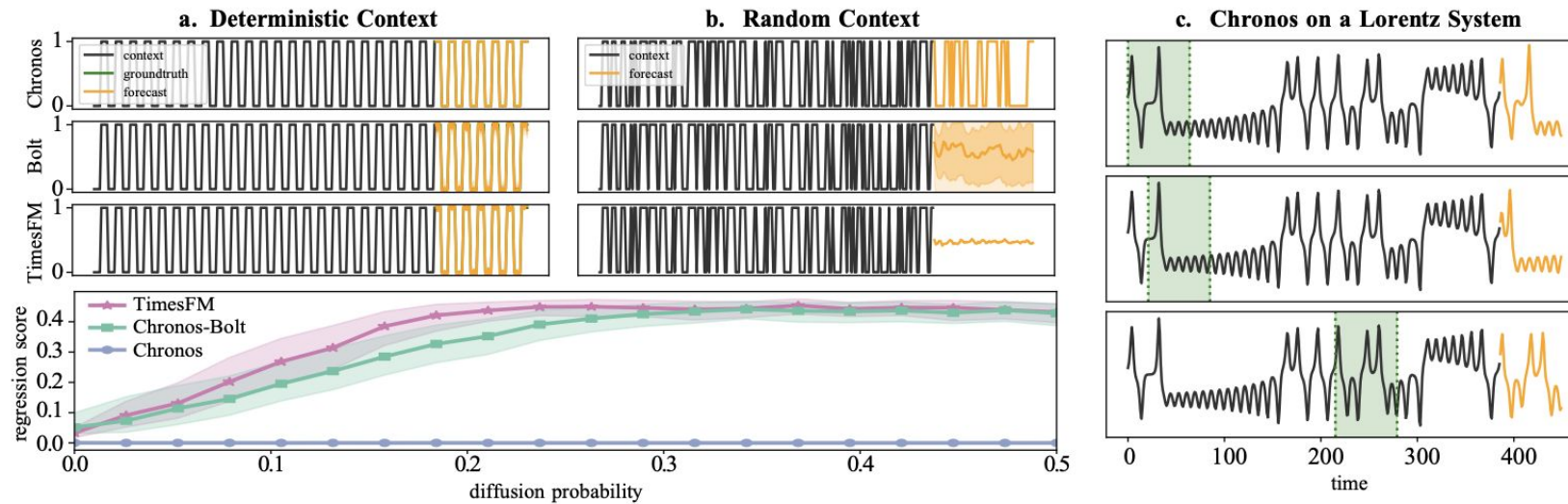
Chronos-Bolt struggles to forecast the near-zero period in the multi-offset time series

Norm Bias: Good or Bad?



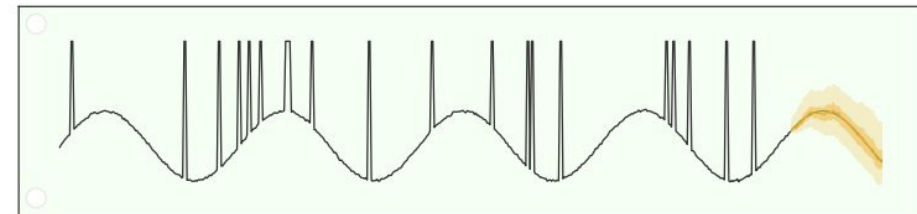
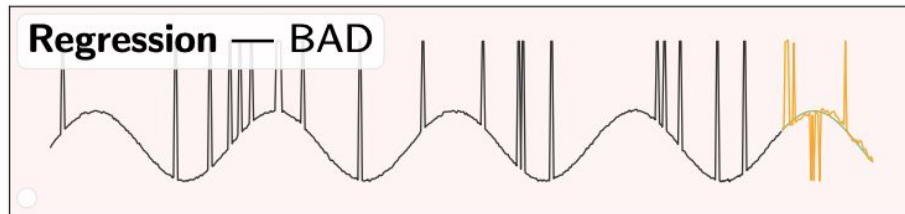
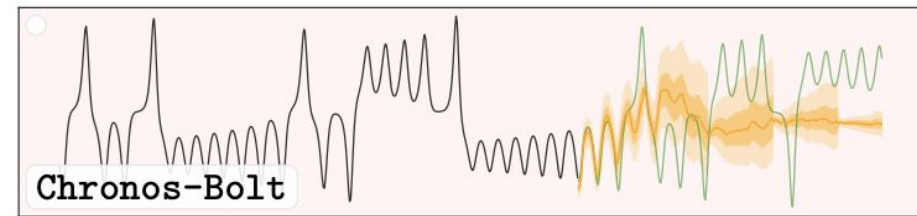
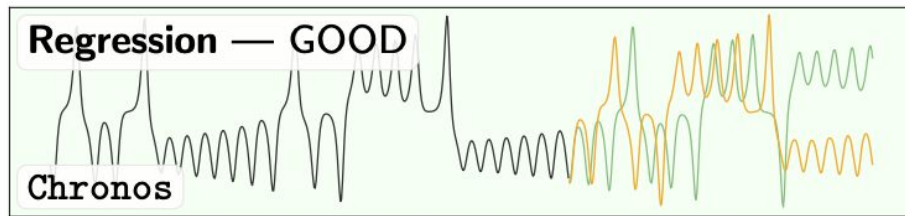
How Do TSFMs Regress to the Mean?

Regression-to-the-Mean Bias



- Models trained with L^2 or L^1 regression losses regress to mean or median
- Chronos with cross-entropy loss models full probability distribution and settles on a mode
- In chaotic systems, fractal dimension is measurement of long-term geometry of the trajectories and regressing to mean/median can severely damage it
- Chronos “parrots” three distinct outcome branches from the context of a Lorentz chaotic system

Regression-to-the-Mean Bias: Good or Bad?



Understanding Transformers for Time Series

Yu, A., **Maddix, D.C.**, et al., "Understanding Transformers for Time Series: Rank Structure, Flow-of-ranks, and Compressibility", arXiv preprint arXiv:2510.03358, Under Review, 2025.

Understanding Transformers for Time Series

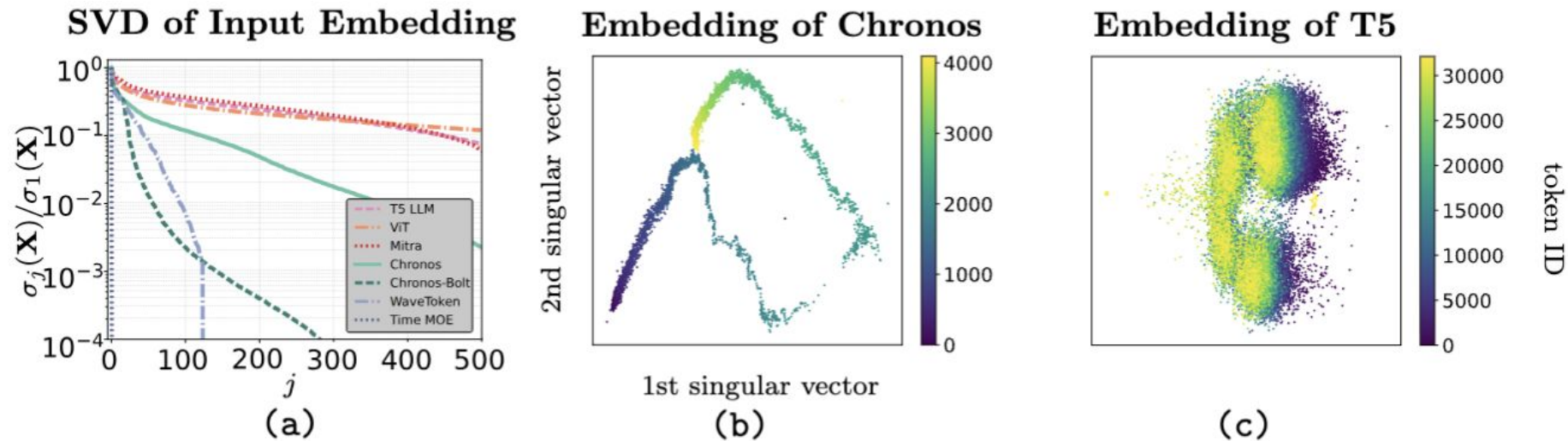
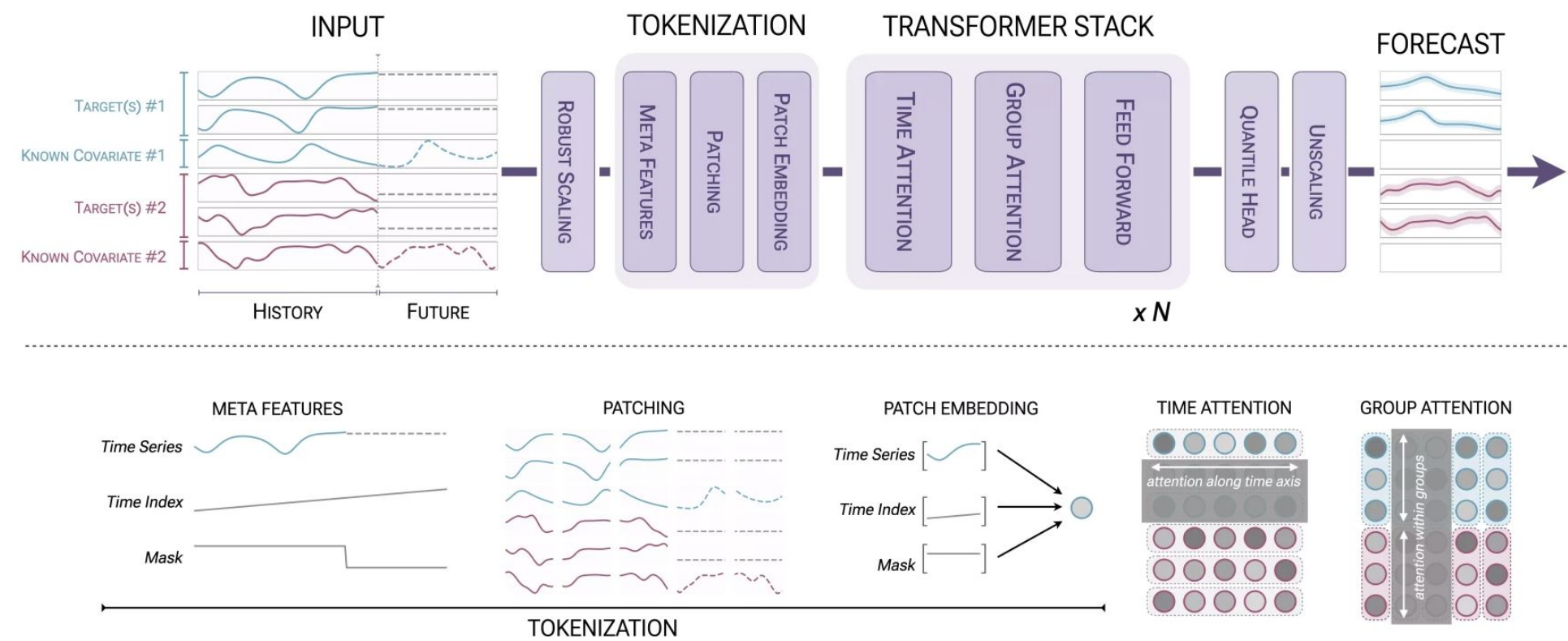


Figure 2: (a): Singular values of the embedded input matrices from many different TSFMs, a TFM, a ViT, and an LLM. (b, c): Embedding space of Chronos and a T5 LLM, respectively, visualized by projecting them onto the leading two singular vectors of the embedding matrix.

Conclusion: Have we achieved a BERT moment?

- Identify design choices in TSFMs that cause 3 inductive biases:
 - Temporal, geometric and regression-to-mean
- Careful numerical analysis and design of TSFMs is required
- Temporal data is a different data-modality, e.g., more compressible, has frequency parameters, continuity in time
- Bitter Lesson
 - Adding traditional forecasting inductive biases can help improve performance on classical benchmarks
 - But it can hurt generalization on unseen domains and tasks, e.g., chaotic systems

Chronos-2: From Univariate to Multivariate



Ansari, A.F., et al., "Chronos-2: From Univariate to Universal Forecasting", arXiv preprint arXiv:2510.15821, 2025.

fev-bench : Realistic benchmark for time series forecasting

- Large-scale evaluation on real-world forecasting tasks
 - 100 univariate & multivariate tasks (incl. 46 with covariates)
- Statistically sound aggregation methods
 - Reliable model comparisons using bootstrap confidence intervals
- Extensible infrastructure for reproducible evaluation
 - Lightweight Python wrapper on top of 🤗 datasets library
- Paper: arxiv.org/abs/2509.26468
- Code: github.com/autogluon/fev
- Leaderboard: huggingface.co/spaces/autogluon/fev-bench

🏆 Leaderboard

Results for various forecasting models on 100 tasks of the fev-bench benchmark, as described in the paper [fev-bench: A Realistic Benchmark for Time Series Forecasting](https://arxiv.org/abs/2509.26468).

Model Name	Avg. win rate (%)	Skill score (%)	Median runtime (s)	Leakage (%)	Failed tasks (%)	Zero-shot	Organization	Link
Chronos-2	91.4	47.3	3.6	0	0	✓	AWS	Link
TiRex	82.7	42.6	1.4	1	0	✓	NX-AI	Link
TimesFM-2.5	77.6	42.2	16.9	10	0	✓	Google	Link
Toto-1.0	70.2	40.7	90.7	8	0	✓	Datadog	Link
Chronos-Bolt	64.4	38.9	1.0	0	0	✓	AWS	Link
Moirai-2.0	64.4	39.3	2.5	28	0	✓	Salesforce	Link





Pairwise Win Rate (SQL) with 95% CIs

Model 2

	Chronos-2	TiRex	TimesFM-2.5	Toto-1.0	COSMIC	Moirai-2.0	Chronos-Bolt	TabPFN-TS	Sundial	Stat. Ensemble	AutoARIMA	AutoETS	AutoTheta	SeasonalNaive	Naive
Model 1															
Chronos-2	50 (50.0, 50.0)	72 (64.0, 81.0)	74 (65.0, 82.0)	78 (70.0, 86.0)	91 (85.0, 96.0)	93 (88.0, 98.0)	94 (89.0, 98.0)	88 (81.0, 94.0)	96 (92.0, 99.0)	95 (90.0, 99.0)	96 (92.0, 99.0)	94 (89.0, 98.0)	99 (97.0, 100)	100 (100, 100)	100 (100, 100)
TiRex	28 (19.0, 36.0)	50 (50.0, 50.0)	54.5 (44.0, 64.0)	68.5 (59.0, 77.5)	74 (65.0, 82.0)	82.5 (75.0, 89.5)	83.5 (76.0, 90.5)	72 (63.0, 80.0)	90.5 (84.5, 95.5)	92 (86.0, 97.0)	95 (90.0, 99.0)	92 (87.0, 96.0)	99 (97.0, 100)	100 (100, 100)	99 (97.0, 100)
TimesFM-2.5	26 (18.0, 35.0)	45.5 (36.0, 56.0)	50 (50.0, 50.0)	57 (48.0, 66.0)	66 (57.0, 75.0)	75 (67.5, 83.0)	73 (65.0, 81.0)	70 (60.0, 79.0)	92.5 (87.0, 97.0)	84 (76.0, 91.0)	92 (86.0, 96.0)	87 (80.0, 93.0)	96 (91.0, 99.0)	99 (97.0, 100)	99 (97.0, 100)
Toto-1.0	22 (14.0, 30.0)	31.5 (22.5, 41.0)	43 (34.0, 52.0)	50 (50.0, 50.0)	47 (37.0, 58.0)	60 (51.5, 69.0)	57 (48.0, 66.0)	57 (48.0, 66.0)	82.5 (74.0, 89.5)	80 (72.0, 87.0)	86 (79.0, 92.0)	83 (75.0, 90.0)	90 (84.0, 95.0)	96 (92.0, 99.0)	98 (95.0, 100)

Shchur, O., et al., "fev-bench: A Realistic Benchmark for Time Series Forecasting", arXiv preprint arXiv:2509.26468, 2025.

Chronos in the Open Source

- Inference code available on [GitHub](#) 
- Model weights available on [Hugging Face](#) 
- Deploy Chronos-2 on AWS using [SageMaker JumpStart](#) 
- Run Chronos with 1 line of code using [AutoGluon](#)  (Chronos-2 coming soon!)

Thank you!

Danielle Maddix Robinson

dmmaddix@amazon.com

<https://dcmaddix.github.io>

