
A Study of In-Context Learning for Periodic Functions

Anonymous Author(s)

Affiliation

Address

email

Anonymous Author(s)

Abstract

1 We present a comprehensive investigation into the capacity of GPT-style transform-
2 ers to learn and generalize periodic functions via in-context learning. Building on
3 the findings of Garg et al. (2023) [1], we evaluate the the ability of small GPT-style
4 transformers to fit sinusoidal superpositions within and beyond its training domain.

5 1 Introduction

6 Recent work has demonstrated remarkable in-context learning capabilities of transformers across
7 diverse function classes [1]. However, the specific challenge of modeling periodic functions remains
8 unaddressed. Periodic signals — ubiquitous in scientific and engineering applications — pose
9 unique learning requirements, motivating our study. We claim that transformers are capable of
10 interpolating between points and slightly outside of a context given to them, but struggle to
11 learn true periodic functions. To check our claim we aim to (1) assess whether standard GPT-derived
12 architectures can implicitly learn periodicity via context alone and (2) evaluate whether augmenting
13 these models with sinusoidal activations enables robust periodic learning.

14 2 Related Work

15 **In-Context Learning on Synthetic Functions** Garg et al. (2023) provide an extensive empirical
16 study of Transformers’ ability to perform in-context learning on synthetic function families. They
17 train standard Transformer models from scratch so that, given a small set of input–output examples
18 from an unseen function, the model can predict the function’s value at new query points purely
19 via its forward pass, without any weight updates. Notably, these Transformers match the analytic
20 least-squares estimator on linear regression tasks and generalize to more complex families such as
21 sparse linear models, two-layer MLPs, and decision trees [1]. This work establishes that Transformers
22 can internalize learning algorithms; our study extends this paradigm to the challenging domain of
23 periodic functions.

24 **Neural Networks and Periodic Inductive Biases** Ziyin et al. (2020) [5] demonstrate that standard
25 feed-forward networks with ReLU, tanh, or sigmoid activations fail to extrapolate even simple
26 periodic signals beyond their training range, owing to a lack of any built-in periodic bias. They
27 introduce the activation $\phi(x) = x + \sin^2(x)$, which embeds a sinusoidal component into each layer
28 and empirically recovers and extrapolates periodic patterns (e.g. temperature cycles, financial time
29 series) that standard networks cannot . Their work highlights the necessity of architectural periodicity.

30 **Fourier Feature Representations** Tancik et al. (2020) [3] show that neural networks suffer from
31 a spectral bias toward low-frequency functions , making it hard to learn fine-scale or high-frequency
32 variation. To address this, they prepend inputs with *random Fourier features*—i.e. mappings of the
33 form $\sin(\mathbf{B}x)$ and $\cos(\mathbf{B}x)$ —thereby furnishing the network with a rich basis of periodic components.

34 This transformation enables simple MLPs to fit and extrapolate high-frequency functions, proving
35 particularly effective for coordinate-based neural representations of signals. These findings suggest
36 that explicit sinusoidal embeddings can compensate for a network’s native deficiencies in learning
37 periodic structure.

38 **Transformers for Continuous Time-Series** Zhou. (2022) [4] propose the Frequency Enhanced
39 Decomposed Transformer (FEDformer) for long-term time-series forecasting. FEDformer integrates
40 a seasonal–trend decomposition step with sparse Fourier blocks that extract dominant frequency
41 components before applying self-attention. This design captures periodic (seasonal) patterns more
42 efficiently than vanilla Transformers, achieving state-of-the-art forecasting accuracy on benchmarks
43 with clear periodic structure. The success of FEDformer underscores that combining Transformer
44 architectures with frequency-domain inductive biases can markedly improve periodic signal modeling.
45 In contrast, our work examines whether unmodified Transformers can learn periodic functions
46 in-context, or whether similar frequency-aware mechanisms are essential.

47 **3 Architecture and Codebase**

48 We began our experiments extending the codebase from Garg et al. [1]. However the code was poorly
49 optimized for execution on GPUs, requiring extensive data transfer between CPU and GPU and thus
50 resulting in excessive training times (over 8 hours per train for the standard-sized model—even using
51 an A100). We conducted initial experiments based on the ‘tiny’ model (3 layers, 2 attention heads,
52 64-dimensional embeddings, 0.2M parameters). However, since it yielded no results, we ended up
53 rewriting the core data generation and training loops from scratch to minimize training time and allow
54 for the training of larger models. Our GPU-optimized solution is contained in the Github [6] at the
55 end of this file. The architecture used by our final version is discussed below.

56 We developed a regression model utilizing the GPT-2 architecture as the backbone. Our model
57 architecture begins with an initial linear layer mapping input dimensions from 1 to 256, followed by
58 a GPT-2-based transformer backbone composed of 12 transformer layers, each containing 8 attention
59 heads with an embedding dimension of 256. The architecture concludes with a final linear layer
60 projecting from 256 dimensions down to a single scalar output. We trained two distinct versions of
61 this architecture: one variant employing the default activation functions (GeLU) and a second variant
62 where each activation function was replaced by a custom activation defined as $\phi(x) = x + \sin^2(x)$.

63 **4 Approach**

64 **4.1 Data Generation**

65 Training and evaluation datasets consist of superpositions of sinusoidal functions of the form

$$f(x) = \sum_{n=1}^k a_n \sin(2\pi n x + c_n).$$

66 All the basis functions are periodic with period 1, and so the entire function is also periodic with
67 period 1. The coefficients a_n are sampled from a standard normal distribution, and c_n is chosen to
68 be either 0 or $\frac{\pi}{2}$ with equal probability. First, m points are sampled at uniform intervals sufficient
69 for perfect reconstruction according to the Shannon-Nyquist theorem. Then, another set of m input
70 values is randomly drawn from a standard normal distribution and appended to the grid values.
71 Finally, the function f is evaluated at each of these points and the evaluations are interleaved with the
72 t-values (see (1)).

73 **4.2 Training Methodology**

74 Our regression model was trained to approximate the value of a target function f at a specific test
75 point t_{test} , given a contextual input comprising m pairs of points. The value of m is selected based
76 on criteria derived from the Shannon-Nyquist reconstruction theorem. Specifically, the input to our
77 model is structured as follows:

$$\left[\underbrace{t_0, f(t_0), t_1, f(t_1), \dots, t_m, f(t_m)}_{\text{context points}}, \underbrace{t_{\text{random}_1}, f(t_{\text{random}_1}), \dots, t_{\text{random}_m}, f(t_{\text{random}_m})}_{\text{randomly generated}} \right]. \quad (1)$$

78 In subsequent discussions, the initial set of $(t, f(t))$ pairs serve as contextual information, whereas
 79 the latter set of randomly sampled t values constitute the prediction targets.

80 During training, the model predicts the function values $f(t)$ for these randomly sampled points.
 81 Ground truth $f(t)$ values at these prediction points are masked (set to zero) to enforce predictive
 82 learning, and the loss is computed exclusively on these masked $f(t)$ values. Importantly, loss
 83 computations exclude all contextual information and input t -values.

84 We evaluated two distinct training methodologies: standard training and curriculum learning. Under
 85 standard training conditions, each training batch comprises the superposition of n sine waves,
 86 with n varying cyclically ($1 \rightarrow 2 \rightarrow \dots \rightarrow n \rightarrow 1 \rightarrow 2 \rightarrow \dots$). In contrast, our curriculum
 87 learning strategy incrementally increases the complexity of training functions by progressively
 88 introducing additional sine-wave components. Specifically, the model undergoes extensive training
 89 (approximately 3,000 to 11,000 epochs) at each incremental step before advancing to functions
 90 composed of more sine waves.

91 As a scoring metric, we use MSE (mean squared error) with a perfect reconstruction as a baseline.
 92 The choice of a perfect reconstruction as a baseline is reasonable, since it is easy to perfectly recover
 93 the coefficients of each component wave of $f(t)$ [2]. We demonstrate this process in a separate
 94 notebook on our GitHub [6].

95 4.3 Resource Budget

96 Our resource budget for this project is \$400. We have used \$100 of it so far (various Google Colab
 97 subscriptions for each member of the team). All training was done using Google Colab on 40GB
 98 NVIDIA A100s.

99 4.4 Evaluation Methodology

100 To evaluate our models, we constructed several periodic functions and evaluated our model on a
 101 dense grids of points over (1) the same period as the context, (2) one period beyond the context, and
 102 (3) a large range that includes the range of the context.

103 5 Experiments & Results

104 5.1 Training

105 Across multiple runs, we found that, controlling for the number of epochs, curriculum learning
 106 converged to a much lower training loss than standard training, similar to the findings of Garg et al.
 107 [1]. See below for a comparison of loss curves for curriculum and non-curriculum learning.

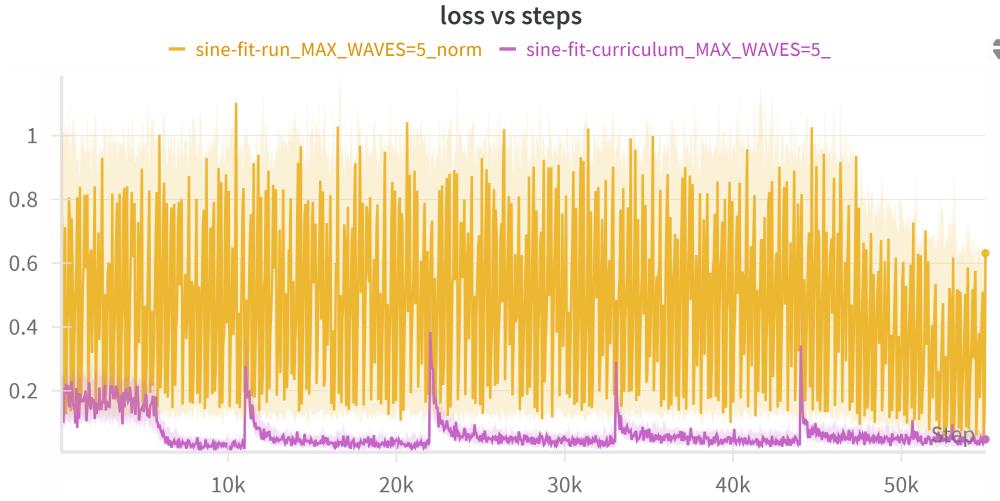


Figure 1: Loss vs. steps for curriculum learning (purple) and standard training (yellow)

108 **5.2 Model Predictions**

109 Upon evaluation, all of our models reconstruct signals quite well in the neighborhood that they were
110 evaluated on during training (the period after the context).

111 After doing local runs we rewrote the code base so that it was optimized for GPU parallelization. With
112 this optimization we were able to run the code base with the standard model. With these parameters.
113 Standard 256 embedding size, 12 layers, 8 heads, and 9.5M total parameters. After these changes, we
114 observed a massive difference in our experiments:

115 For our evaluation we ran 150 experiments where we predict 200 points between $t = 0$ and $t = 3$.
116 We did this for a standard model with GeLU activation and a model with sin activation. We did this
117 for each with 4, 5, and 6 superpositions for the sinusoidal waves. Table 5.2 contains a summary of
118 the data.

Activation Type	# Of Sines	MSE Over 1st Period (0,1)	MSE Over 2nd Period (1,2)	MSE Over 3rd Period (2,3)	Overall
Standard	4	0.44893842	0.31079719	0.2712146	0.34373317
Sinusoidal	4	0.34100499	0.4239749	0.32062751	0.36196518
Standard	5	0.03727583	0.04325096	0.06601895	0.04886744
Sinusoidal	5	0.03010587	0.0315977	0.04238979	0.03458119
Standard	6	0.64249314	0.56375453	0.44534295	0.54915387
Sinusoidal	6	0.32312029	0.33671242	0.36268511	0.33958891

Table 1: Comparison of average MSE Between Standard and Sinusoidal Activations Across Periodic Intervals with Varying Numbers of Sine Components (lower is better).

119 Note that the sinusoidal activation model generally performs better than the standard activation
120 model. The mechanism behind the dominance of the standard activation model under four waves is
121 not fully understood. Also note that we have not evaluated the effectiveness of sinusoidal activations
122 on other transformer tasks. The findings of Ziyin et al. [5] suggest that severe performance hits
123 might not be expected.

124
125 Representative predictions are illustrated in Figures 2–7.

126 **6 Figures**

127 Model predictions and ground truth values: Standard Activations meaning GeLU.

128 An insight that we see is that the model actually improves the further it gets from $t(0)$ until $t(3)$ where
129 the prediction flat lines and error shoots up. This is also the only tested superposition $k = 4$ where
130 GeLU outperforms sine activation.

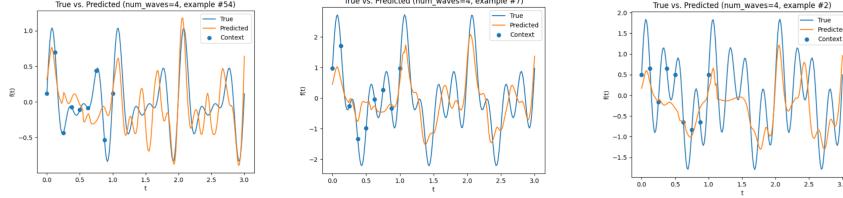


Figure 2: Standard Activations with 4 Superpositions

131 Unsurprisingly the model being trained on 5 superpositions does quite well. Sine activation outper-
132 forms this GeLU activation. Its mean squared error is a factor of ten less than the other superpositions.

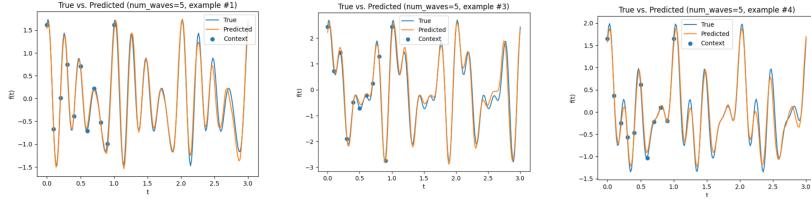


Figure 3: Standard Activations with 5 Superpositions

133 A similar insight to 4 superpositions is seen here. We observe that the model actually improves the
134 further it gets from $t = 0$ until $t = 3$ where the prediction flat lines and error shoots up.

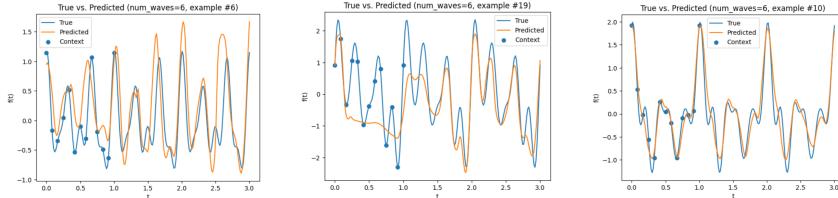


Figure 4: Standard Activations with 6 Superpositions

135 An interesting observation is that the MSE increases from $t = 1$ to $t = 2$ and then decreases from
136 $t = 2$ to $t = 3$ to even beyond what it predicts within the context $t = 0$ to $t = 1$. Furthermore this is
137 the only experimental setup where sine-based activation is outperformed by GeLU activation.

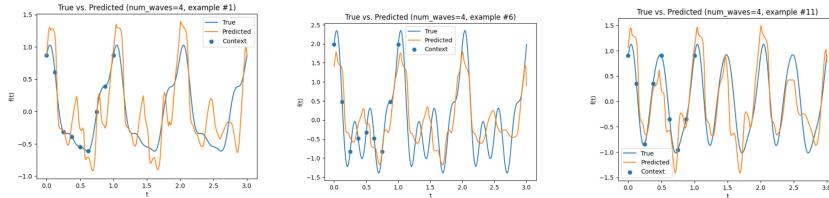


Figure 5: Sinusoidal Activations with 4 Superpositions

138 This is the best current model.

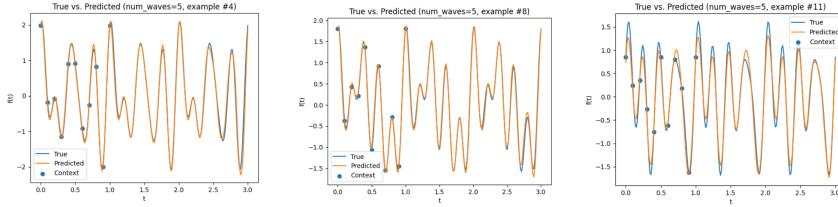


Figure 6: Sinusoidal Activations with 5 Superpositions (max number trained on)

139 Here we see what we expect with MSE increasing as predictions get further from the model's given
140 context. It outperforms its GeLU counterpart.

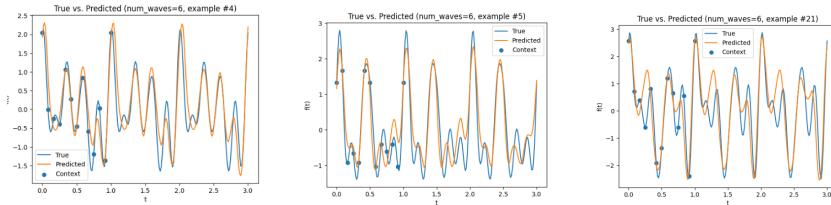


Figure 7: Sinusoidal Activations with 6 Superpositions

141 7 Next Steps, Current Problems

- 142 1. Evaluating the model's ability to generalize over a large interval by training on multiple
143 smaller subintervals that collectively span the full range (e.g., training on intervals of size
144 10 covering the range [0, 100]).
- 145 2. Training the model on disjoint intervals and evaluating if it is able to interpolate between
146 those intervals. (e.g. training on points from [0, 1] & [2, 3]; evaluate its on points from
147 [1, 2]).
- 148 3. Train using more robust curriculum schedules that incorporate small amounts of data using
149 all previous numbers of sine waves during each step to ensure model generalization.
- 150 4. Try autoregressive prediction to see what effects this has on model generation at far away
151 points.
- 152 5. Train on noisy data.
- 153 6. Experiment with real-life data, e.g., crop yield and weather.
- 154 7. Normalize the training data (divide sum of sine waves by total number of waves).
- 155 8. Shuffling and testing adding extra random points from the sine wave into the context.
- 156 9. Apply alternative regularization schemes (e.g., weight decay, dropout variations).
- 157 10. Selectively apply the sinusoidal activation function to certain parts of the transformer
158 architecture instead of swapping all activations.

159 References

- 160 [1] S. Garg et al. What Can Transformers Learn in Context? (2023). <https://arxiv.org/pdf/2208.01066.pdf>.
- 162 [2] Shannon-Nyquist sampling theorem. <https://mathworld.wolfram.com/SamplingTheorem.html>
- 164 [3] Tancik et al. Spectral Bias and Fourier Features. <https://mathworld.wolfram.com/SamplingTheorem.html>

- 166 [4] Zhou et al. Frequency Enhanced Decomposed Transformer //arxiv.org/abs/2201.12740
167 [5] L. Ziyin et al. Neural Networks Fail to Learn Periodic Functions and How to Fix It. (2020).
168 <https://arxiv.org/abs/2006.08195>.
169 [6] Our Code Base https://github.com/drnestelloop/periodic_icl
170 [7] RunPod / Lambda GPU Pricing. <https://www.runpod.io/pricing>; <https://lambda.ai/service/gpu-cloud?matchtype=...>

172 8 Self Review

173 8.1 Main Goal and Claims

174 The main goal of the project is to evaluate the ability of GPT-style transformers to learn and generalize
175 periodic functions via in-context learning, and to determine whether architectural modifications, such
176 as sinusoidal activations, improve their ability to do so.
177 We claim that transformers are capable of interpolating between points and slightly outside
178 of a context given to them, but struggle to learn true periodic functions (i.e. accurately model the
179 periodic nature of a function over arbitrarily many periods).

180 8.2 Experimental Design and Evaluation Protocol

181 We trained two GPT-style transformers on synthetic datasets composed of sinusoidal superpositions.
182 We began with a 'tiny' GPT-2 style model (not shown) (3 layers, 2 heads, 64-dim embeddings) and
183 then progressed to a 'standard' model (12 layers, 8 heads, 256-dim embeddings). We trained one
184 model with standard GeLu activations and the other with a sine-based activation function.
185 Models are evaluated using mean-squared error (MSE) measured on freshly-generated periodic func-
186 tion samples for both interpolation (within-domain) and extrapolation (outside-domain). Performance
187 is compared against ground truth. We compared MSE loss across intervals and further visually
188 evaluated the performance by plotting the predictions versus the ground truth.

189 8.3 Data and Task

190 The data used for training and evaluation are synthetic periodic functions generated as superpositions
191 of sinusoids with randomly sampled parameters and Gaussian noise added to ensure robust training
192 signals. Each training example presents a sequence of input/output pairs as context; the task is to
193 predict the output of a function f at a new query t_{new} .

194 8.4 Support for Claims

195 The experiments support our claims. For the case of five superimposed sinusoids, the standard
196 Transformer records an average mean squared error (MSE) of 0.0489 on held-out query points,
197 showing that it can interpolate between examples with decent accuracy. The variant with sinusoidal
198 activations reduces its MSE to 0.0346—a relative improvement of nearly 30% demonstrating that
199 embedding a periodic bias enhances the model's ability to learn and generalize periodic patterns. The
200 predictions gradually get worse as they move further from context thus not showing true growth.

201 8.5 Limitations, Strengths, and Weaknesses

202 Limitations are acknowledged, including the transformers' struggles with extrapolation tasks and the
203 inefficiency of the initial code base. However, deeper theoretical reasons behind the transformer's
204 difficulty with periodic extrapolation are not extensively discussed. Additionally, while our synthetic
205 data regime enables precise control over signal complexity and noise, we do not evaluate real-world
206 periodic signals (e.g. sensor or climate data), leaving questions of practical applicability open to
207 future work.

208 Strengths of the paper include clearly defined objectives, well-structured experimentation, and
209 careful step-by-step troubleshooting and optimization. We have transparently documented our
210 implementation issues and systematically explored potential remedies.

211 Weaknesses include the absence of quantifying statistical significance in reported improvements,
212 extensive justification for hyperparameter choices, and the limited interpretability of initial experiment
213 outcomes. Additionally, a broader set of comparisons with other advanced regression techniques
214 beyond 3-NN would strengthen the findings.

215 **8.6 Suggestions for Improvement**

216 To strengthen our project’s impact, we believe we could (1) include experiments on real-world
217 periodic datasets, (2) report variance estimates over multiple random seeds or cross-validation folds,
218 and (3) provide a more systematic ablation of the sinusoidal activation parameters (e.g. frequency
219 scaling factors). Additionally, we could implement more rigorous theoretical analysis or experimental
220 comparisons with other function approximation methods, such as Fourier-based neural networks or
221 state-of-the-art recurrent architectures, to better contextualize the transformers’ performance.

222 **8.7 Related Work and Reproducibility**

223 Relevant related work is clearly identified, notably in-context learning of functions in Garg et al.
224 (2023) [1] and periodic induction via custom activations in Ziyin et al. (2020) [5].

225 Reproducibility of our work is enabled by the clear description of model architectures, data generation
226 procedures, and optimizer settings. We share the optimized code base via GitHub [6] with descriptions
227 for each code cell. Rerunning the experiments and reproducing results is very feasible, assuming
228 access to comparable computational resources.

229 **8.8 Presentation and Readability**

230 Most figures in the paper have well-labeled axes and clear legends, though some might benefit from
231 axis explanations or error bars to convey variance. The methodology is explained in sufficient detail.

232 **8.9 Quality of Writing**

233 We have written this paper to the best of our abilities. Our process involved one contributor authoring
234 a section, followed by a peer review from the other 3 contributors and the use of online tools such
235 as Grammarly. We are therefore confident that the English prose is generally clear, with only minor
236 stylistic improvements possible in terms of conciseness, brevity, and succinctness.

237 **8.10 TODO**

238 We have a strong understading of what we should do moving forward. I think robust testing for
239 if our model can predict between periods, experimenting with real life data, and normalizing the
240 training data are high priority. After that trying out curriculum learning while bringning in training
241 from priorThe rest seem to be lower priority. Perhaps they can display their predictions against
242 nearest-neighbors to display the improvement.