

PLOTS UNLOCK TIME-SERIES UNDERSTANDING IN MULTIMODAL MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

While multimodal foundation models can now natively work with data beyond text, they remain underutilized in analyzing the considerable amounts of multi-dimensional time-series data in fields like healthcare, finance, and social sciences, representing a missed opportunity for richer, data-driven insights. This paper proposes a simple but effective method that leverages the existing vision encoders of these models to “see” time-series data via plots, avoiding the need for additional, potentially costly, model training. Our empirical evaluations show that this approach outperforms providing the raw time-series data as text, with the additional benefit that visual time-series representations demonstrate up to a 90% reduction in model API costs. We validate our hypothesis through synthetic data tasks of increasing complexity, progressing from simple functional form identification on clean data, to extracting trends from noisy scatter plots. To demonstrate generalizability from synthetic tasks with clear reasoning steps to more complex, real-world scenarios, we apply our approach to consumer health tasks – specifically fall detection, activity recognition, and readiness assessment – which involve heterogeneous, noisy data and multi-step reasoning. The overall success in plot performance over text performance (up to an 120% performance increase on zero-shot synthetic tasks, and up to 150% performance increase on real-world tasks), across both GPT and Gemini model families, highlights our approach’s potential for making the best use of the native capabilities of foundation models.

1 INTRODUCTION

Multimodal models like GPT4 (Achiam et al., 2023) and Gemini (Gemini Team et al., 2023) are trained to understand visual information natively. However, they are not specifically trained to understand time-series data – in particular, the tokenizers for large language models (LLMs) are not well-suited for representing large sequences of floating point numbers (Spathis & Kawsar, 2024). This mirrors the human approach (Card et al., 1999; Yalçin et al., 2016); we cannot easily make sense of a long array of floating point numbers - instead our first instinct is often to visualize the data through plotting, followed by extracting insights through statistical analysis.

We investigate the hypothesis that multimodal models understand time-series data better through their vision encoders than through the textual representation of the sequences using synthetic and real-world data experiments. Our synthetic data experiments allow us to closely control the difficulty of tasks through the addition of noise and by changing the number of points in each function. We also use a mix of tasks that require a differing number of reasoning steps, as well as different kinds of reasoning, to get a correct answer.

Fall detection and activity recognition are both real-world tasks that make use of inertial measurement units (IMUs) from mobile phones or wearable devices. IMUs are 6-dimensional waveforms consisting of 3 axes of acceleration data and 3 axes of angular velocity data. The fall detection task consists of classifying an IMU waveform segment into one of three classes: Fall, Active Daily Living (ADL) or Near Fall (a hard negative class). The activity recognition task consists of classifying a waveform segment into one of five classes: Sitting, Standing, Walking, Cycling or Stairs.

By contrast, the readiness assessment task is a binary classification of 28 days of training load data from a single user into a state of undertraining or overtraining. Because of the tabular nature of the

054 data, the plot version is presented as a bar plot. This is not the ideal setting for our method - we
055 believe it's best used when the amount of data exceeds what's reasonably presentable in a text table.

056
057 Our findings show that when using our plot-based approach multimodal models perform much better
058 on tasks where the result is dependent on understanding the overall trend. We find specific examples
059 of this when identifying the functional form, the number of clusters, the correlation between two
060 functions, and on the real-world pattern-recognition tasks of activity recognition and fall detection.
061 For example, GPT4o using plots on the functional form identification task shows a performance
062 improvement of 122% over using the text representation. On other tasks that require more advanced
063 reasoning such as multi-step or connecting trend shapes with sequence magnitudes (e.g. identifying
064 derivatives), and on tasks with tabular data (e.g. readiness assessment), the performance is equiv-
065 alent. However, there is a substantial cost difference between vision and text prompts, which is
066 particularly pronounced on very long-context tasks, as the same information in a long sequence that
067 requires many (10,000's to 100,000's) text tokens can be represented in one plot with many fewer
068 (100's to 1000's) vision tokens. While vision tokens are more expensive than text tokens, the differ-
069 ence in unit cost is much lower than the orders of magnitude difference in overall prompt length, so
070 that the total cost is still much lower using the vision approach. This difference is especially rele-
071 vant on tasks where extensive few-shots are required to achieve good performance, and optimizing
072 token efficiency translates to significant resource savings. Not only does this plot-based approach
073 achieve better performance while being more efficient, it is also completely generalizable across any
074 task that involves reasoning about a long, complex time-series as it requires zero additional model
075 training.

076
077 Our work empirically evaluates the relative performance of the native capabilities of existing mul-
078 timodal foundation models on visual versus textual representations of time-series data. This con-
079 tribution furthers the understanding of modern foundation models, which is important in real-world
080 contexts as user-facing products continue to develop multimodal sophistication and users interrogate
081 data types that are more complex than can be easily represented with text only. While we do not
082 claim to achieve the same absolute performance as models trained for specific tasks, and likely our
083 approach will not match such models, our results presented here nonetheless show that in contexts
084 when one relies on a foundation model to ingest any general time-series data with *a priori* unknown
085 characteristics, a visual representation is on balance likely to yield better, and cheaper, results.

086 2 RELATED WORK

087
088 **Forecasting** We use the term “time-series understanding” in this paper to distinguish from time-
089 series forecasting. Time-series forecasting predicts future data points based on points seen so far,
090 whereas we are primarily interested in the setting where we connect the time-series data to a multi-
091 modal model for further analysis. In particular, we want to show that multimodal models can reason
092 about overall trends, the relationship between multiple time-series, overall clustering of data, and
093 other time-series understanding tasks. Forecasting has been a very productive area of the field –
094 for a closer look at the literature, the survey paper by Zhang et al. (2024) tracks a wide range of
095 time-series understanding and forecasting work.

096
097 **Time-series models** There are several existing approaches that train time-series encoders for specific
098 tasks or domains. For example, Chan et al. (2024) trained a domain-specific encoder for multimodal
099 medical time-series analysis. Similarly, Cosentino et al. (2024) trained an encoder for the sleep
100 data in their Digital Well-being task. In this paper, we are not claiming that our approach would
101 outperform a task-specific model on a specific task – we claim that one can achieve much better
102 performance from a foundation model by exploiting its native multimodal capabilities compared to
103 using only text.

104
105 Others have also shown that training foundation models with Transformer architectures specifically
106 to work in time-series contexts can lead to good results across tasks including mostly forecasting
107 but also classification, anomaly detection and imputation. These trained models do especially well
108 when the input time-series are carefully pre-processed and tokenized, including patching, scaling
109 and quantization Das et al. (2023); Woo et al. (2024); Goswami et al. (2024); Ansari et al. (2024);
110 Cai et al. (2023). While they did not train a new model, in LLMTime (Gruver et al., 2024) the
111 authors showed that with careful tokenization, text-only LLMs can perform well at forecasting tasks;

108 we perform ablations based on their methods and select the best tokenizations accordingly for our
109 text baselines.

110 In this work, our goal is to show that simply plotting the data without additional data preprocessing
111 or model training is at the very least an easy first step, and might be a helpful approach when training
112 a task-specific encoder from scratch may not be feasible due to the requirements on having additional
113 paired data, compute and expertise.

114 **Vision models and visual representations** While studying multimodal models’ abilities to reason
115 about visual inputs, Rahmanzadehgervi et al. (2024) found that multimodal models are unable to
116 reason effectively, although some follow-up work by Corin (2024) indicated that prompt engineering
117 can fix losses. In any case, our results do not necessarily contradict this - for many of our tasks
118 humans may be able to get perfect scores, and indeed the multimodal models do not. Regardless,
119 our main claim that plots are better suited than text as input to a multimodal model for time-series
120 understanding holds true.

121 Perhaps an inversion of our approach, DePlot (Liu et al., 2023) translated visual plots to numeric
122 tables and operated on the tabular data. This approach may be sound for discrete data where the
123 number of points remains small – cases where a human would be expected to understand a table of
124 data well.

125 Closely related to our work are those methods that use vision-embedding models like Contrastive
126 Language–Image Pre-training (CLIP) by Radford et al. (2021). Wimmer & ReKabsaz (2023) used
127 CLIP to embed plots of financial time-series data, from which features are extracted for use by down-
128 stream classifiers. Since we use the vision encoders of the multimodal foundation models directly,
129 there is no need for further feature extraction or downstream classifiers in our approach. IMU2CLIP
130 (Moon et al., 2023) and ImageBind (Girdhar et al., 2023) used video and image data paired with
131 waveforms to learn joint embeddings. Both of these works rely on existing paired waveform and
132 video data to “bind” the modalities together, whereas we can simply plot the waveforms and use the
133 existing multimodal vision encoder to derive our time-series embeddings.

134 **Measuring understanding** Past work has also investigated various approaches to measuring the
135 degree to which models can understand and reason about various modalities of inputs, such as charts
136 (CharXiv (Wang et al., 2024)), tables and figures (SPIQA (Pramanick et al., 2024)) and time-series
137 themselves (TimeSeriesExam (Cai et al., 2024)). These works generally involve generating novel
138 evaluation datasets, and in some cases (e.g. (Wang et al., 2024), (Pramanick et al., 2024)) rely
139 on language models to generate the questions themselves. In our work we deliberately avoid this
140 approach as it can introduce biases during evaluation that are hard to account for (e.g. favoring their
141 own output as in (Panickssery et al., 2024)). In TimeSeriesExam (Cai et al., 2024), the authors also
142 found, as we do, that models perform better on plot-based representations of time-series than the
143 analogues text-based representations, though only demonstrate this on synthetic data as part of a
144 carefully optimised exam generation algorithm.

145 146 147 3 METHODOLOGY 148

149 We evaluate our visual prompting method on both synthetic data and real-world use-cases. Syn-
150 thetic data allows us to control the difficulty of the task by adding noise and altering the number
151 of data points, and to investigate specific kinds of reasoning in isolation. We chose the synthetic
152 tasks to align with the different steps of reasoning we hypothesize are required for the representa-
153 tive real-world use cases we test on. Note that in this context, “reasoning” refers to the high-level
154 steps we believe humans would take to get to the right answer, rather than any formal modelling ap-
155 proach such as chain of thought. These tasks include understanding the local and global longitudinal
156 signatures (trend and magnitude) of a time-series, and potentially comparing it with several other
157 time-series (as in the case of multidimensional sensing). We summarize in Section 4.1 the different
158 tasks in our experiments, along with the type of reasoning we are probing.

159 The goal of our work is to study specifically the difference in performance achieved by models
160 when ingesting visual versus textual representations rather than the absolute performance on any
161 one task with either modality. As such the appropriate baseline, and the one we use, is the models’
performances on textual representations. We nonetheless include random choice baselines for con-

162 text to show there is utility in leveraging these models at all, and compare against state-of-the-art
163 task-specific models (for two of the real-world tasks) for context.

164 3.1 STRUCTURED PROMPTING

165 We used the open-source structured prompting library Langfun (Peng, 2023) for all tasks in this paper,
166 with the exception of the Readiness task (Section 4.2) which is processed in a privacy-preserving
167 sandbox environment. The prompts and Langfun code snippets for all tasks (except Readiness) are
168 provided in Appendix A.4 for reproducibility. The structured prompting approach in Langfun allows
169 us to use target schemas for outputs, though we do not use the controlled generation feature (Gemini
170 models) or structured output (GPT4o models), simply relying on the native formatting of the model
171 to the correct schema.
172
173

174 3.2 MODELS

175 We tested all synthetic data tasks on two frontier models: Gemini Pro 1.5 (gemini-pro-001) and
176 GPT4o (gpt4o-2024-08-06) and two smaller models Gemini Flash 1.5 (gemini-flash-001) and
177 GPT4o-mini (gpt4o-mini-2024-07-18). We use a temperature of 0.1 for all our experiments, Sup-
178plementary Tables S33-S35 includes our ablations on temperature. All other sampling parameters
179 remain at API defaults.
180
181

182 3.3 FLOATING POINT REPRESENTATION

183 In order to find the best textual representations, we ran ablations (Appendix A.3) inspired by LLM-
184 Time (Gruber et al., 2024) on which floating point precision (2, 4, 8, 16) and separator (space or
185 comma and space) to use. We also tested the scaling approach suggested by LLMTime. We found
186 that the lower precision led to better performance. On synthetic tasks, the best performing separator
187 differs per model, on Gemini we make use of the space separator whereas on the GPT4o family we
188 use comma and space. On real-world tasks we make use of the space separator for all models as
189 we did not observe a difference in performance on these tasks and the space separator uses fewer
190 tokens.
191

192 3.4 STATISTICAL METHODS

193 For our aggregate results, we aggregate individual model responses to an overall performance quality
194 metric (accuracy or mean absolute error (MAE)) over the task dimensions as described in Supple-
195 mental Section A.1.1. This produces multiple points from which we extract a distribution presented
196 as a box-plot where the central line is the median, the edges of the boxes are the inter-quartile range
197 (IQR), the whisker lengths extend to 1.5 times the IQR and outliers are presented as individual
198 points. For our more detailed plots in Appendix A.2 we show 95% confidence intervals constructed
199 from 1.96 times the standard error of the mean of the metric.
200

201 For real-world tasks, since we don't have the ability to regenerate the same problem, we instead make
202 use of bootstrapping (with 1,000 replicates) to produce distributions of the macro-averaged F_1 scores
203 from which we construct similar box-plots as for the synthetic tasks. Note that the distributions
204 plotted in the real-world box-plots are thus expected to be tighter than the synthetic task plots, as
205 they don't reflect independent replicates.

206 In Table 2, we present the median and IQRs of the differences between plot and text performances,
207 with the difference taken such that a positive value always means the plot method performs better
208 than the text method. For synthetic tasks, the median and IQR are calculated by directly creating
209 the distribution of differences between the plot and text performances of different instances of the
210 experiment, while for real-world tasks we create a distribution of differences by randomly sampling
211 1,000 random pairs of the bootstrapped metric distributions described immediately earlier.

212 For synthetic tasks only, we test for significant differences between the plot and text performances
213 with a two-sided Wilcoxon signed-rank test (Wilcoxon, 1945). We apply a Bonferroni (Bland &
214 Altman, 1995) correction for multiple comparisons within a task block. We could not apply the
215 same hypothesis testing framework to the real-world tasks as the performance distributions were
bootstrapped and thus not independent, violating the assumptions of the Wilcoxon test.

4 EXPERIMENTS AND RESULTS

Type of data	Task	Reasoning	Time-series length
Synthetic	Functional form id.	Understanding of one overall trend	10's - 1,000's
	Correlation of two lines	Understanding of two overall trends	10's - 1,000's
	2D cluster counting	Understanding and counting N overall trends	10's - 1,000's
	Derivative id.	Multi-step reasoning connecting two overall trends	10's - 1,000's
	Quadratic derivative id.	Multi-step reasoning connecting two trends and magnitudes	10's - 1,000's
Real-world	Fall detection from IMU data	Classify a pattern based on local spikes in multiple signals	10,000's
	Activity recognition from IMU data	Classify a pattern based on global trends in multiple signals	10,000's
	Readiness from wearable measures of training intensity	Compare a local trend with a global trend	10's

Table 1: Summary of the tasks we study in this paper, including the reasoning each requires and the length of the input time-series (based on the number of points). The “reasoning” column is a high-level summary of the steps needed to answer the tasks’ questions correctly; tasks are detailed in Sections 4.1 and 4.2 and Supplementary Information Section A.1.

Task (Metric)	Few-Shots	GPT4o-mini	Gemini Flash 1.5	GPT4o	Gemini Pro 1.5
Functional form id. (Accuracy)	0	0.32 (0.18, 0.41)*	0.22 (0.11, 0.40)*	0.46 (0.31, 0.52)*	0.04 (-0.08, 0.25)
Correlation of two lines (Accuracy)	0	0.33 (0.17, 0.33)*	0.25 (0.17, 0.33)*	0.17 (0.00, 0.17)*	0.33 (0.17, 0.50)*
2D cluster counting (MAE)	0	1.02 (0.53, 1.09)*	1.67 (1.49, 1.80)*	2.29 (1.84, 2.44)*	1.82 (1.36, 2.06)*
Derivative id. (Accuracy)	0	0.16 (0.12, 0.24)*	0.08 (-0.04, 0.20)	0.00 (-0.18, 0.12)	-0.02 (-0.16, 0.08)
Quadratic derivative id. (Accuracy)	0	0.27 (0.23, 0.38)*	0.15 (0.02, 0.30)*	-0.17 (-0.30, -0.10)*	0.17 (-0.01, 0.24)
	3	0.17 (0.12, 0.33)*	0.32 (0.13, 0.34)*	0.10 (-0.07, 0.17)	0.28 (0.19, 0.43)*
Fall detection (F_1 score)	1	0.03 (0.02, 0.05)	0.11 (0.09, 0.12)	0.32 (0.31, 0.34)	0.13 (0.10, 0.15)
	10	0.21 (0.19, 0.22)	0.17 (0.15, 0.19)	0.50 (0.49, 0.52)	0.40 (0.38, 0.42)
Activity detection (F_1 score)	1	0.09 (0.07, 0.11)	0.12 (0.10, 0.14)	0.03 (0.01, 0.06)	0.20 (0.18, 0.22)
	5	0.11 (0.09, 0.13)	0.23 (0.21, 0.25)	0.18 (0.15, 0.20)	0.23 (0.21, 0.25)
Readiness (F_1 score)	0	–	-0.08 (-0.11, -0.06)	–	0.07 (0.05, 0.09)

Table 2: Summary of the experiments with 38 out of 42 results showing better performance on plots (bold numbers). Cells contain metric medians and IQRs. Stars in synthetic tasks only indicate statistically significant differences between plot and text metrics at 95% confidence corrected for multiple comparisons; we could not perform the same hypothesis testing on the real-world tasks. See Section 3.4 for statistical details and Supplementary Table S2 for relative differences between approaches.

4.1 SYNTHETIC DATA TASKS

Figure 1 summarizes all the zero-shot versions of the synthetic data tasks showing that plot-based methods outperform the text-based methods across GPT and Gemini model families, with few exceptions. Detailed task descriptions are provided in Supplementary Information Section A.1.1, and non-aggregated results over dataset parameters such as number of points and noise level are available in Supplementary Information Section A.2.

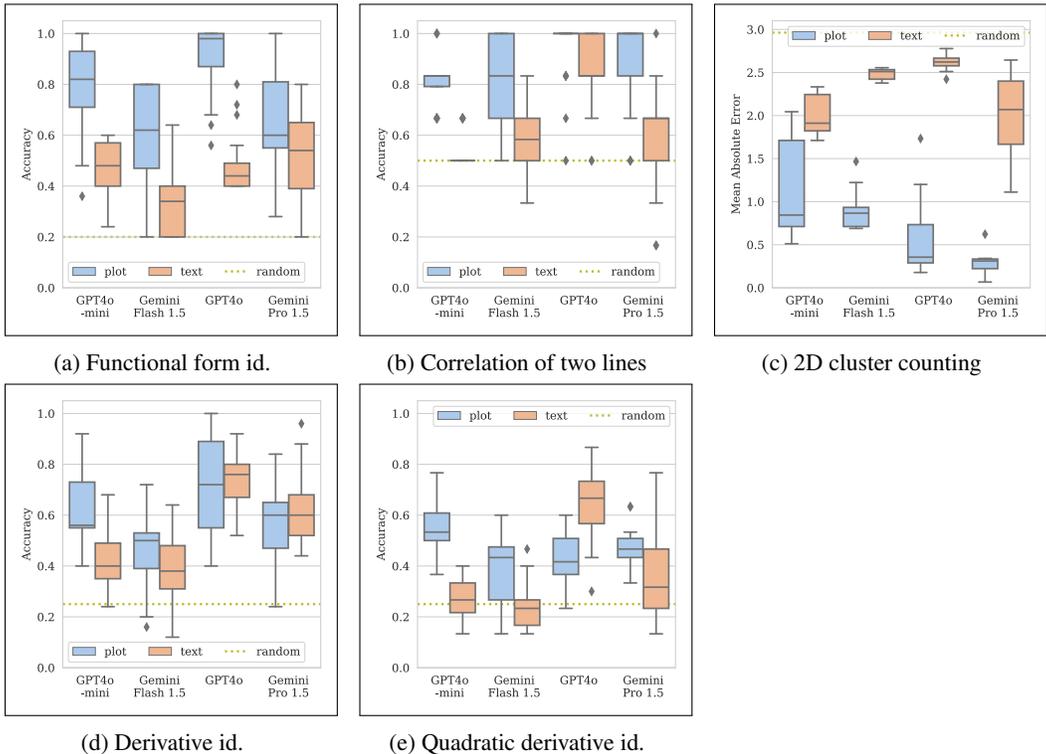


Figure 1: Zero-shot synthetic data results showing plot- and text-based accuracy (MAE for the cluster counting task) distributions for all models, with horizontal lines representing random performance. The results generally show better performance for plots compared to text across models.

Functional form identification (id.) This is the simplest task that requires only identifying one overall trend and correctly classifying it into one of five functional tasks (linear, quadratic, cubic, exponential or periodic). We generate a set number of points with a controlled level of noise according to one of the five function classes, and test the model’s ability to label the global trend into the correct class.

Correlation of two lines This task now requires understanding the trends in two lines and comparing them against each other to correct identify whether the lines are positively or negatively correlated. Here we generate two linear series with controlled number of points and noise and predetermined slopes, measure the sign of the correlation analytically using the Pearson correlation coefficient, and probe the model’s ability to identify the directional correlation.

2D cluster counting In this task, the model needs to correctly identify the N number of clusters present in a set of points. To test this, we generate random points on a 2D grid with a set number of clusters and controlled minimum distance between cluster centers and standard deviation of the points about the centres. The model is then instructed to count the number of clusters.

Derivative identification This is a harder task: the model must now identify the correct first derivative (out of four choices) of the function provided in the question. The function and choices are presented as either plots or text. We pass various known functions to the model alongside four synthetic first derivatives, each corresponding to different functional classes, with controlled levels of

noise and number of points. The models are then asked to identify which of the multiple choices corresponds to the true derivative.

Quadratic derivative identification As a hard variant of derivative identification, we always pass a quadratic function and present four different linear functions as multiple choice answers, with controlled levels of noise and number of points. The model must now identify the correct linear function (out of four choices) that corresponds to the quadratic function’s derivative, so it must pay attention to both the sign and magnitude of the linear slopes.

In order to investigate the quadratic derivative task further, we also ran experiments providing few-shot examples with reasoning traces with results shown in Figure 2. Here we find that GPT4o text zero-shot remains an outlier in its strong performance, but for the other models plot outperforms text, with few shots improving performance in the Gemini family for both plot and text, but reducing performance in the GPT family of models for both plot and text.

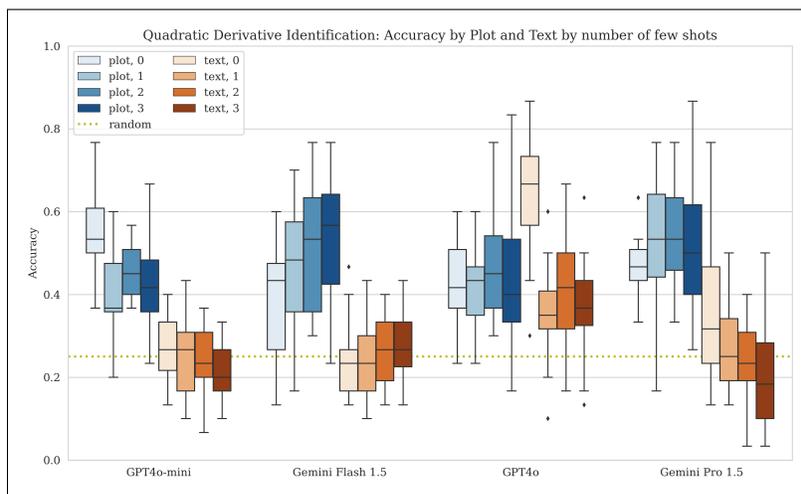


Figure 2: Quadratic derivative identification results show zero-shot plots outperform text, except for the outlier GPT4o model. When using few-shots, more examples generally improves the gain.

4.2 REAL-WORLD TASKS

Our synthetic tasks built up in complexity from understanding one trend globally (functional form identification) to understanding several trends simultaneously using global and local information (correlation and cluster counting) and complex multi-step reasoning (derivative identifications). We now probe tasks on real-world data that require a mix of these reasoning abilities, including simultaneous understanding of multiple sensor signals to uncover either local or global trends in the first two tasks (fall detection and activity recognition), and extracting two trends of different timescales in the last task (readiness).

Fall detection from inertial measurement units (IMUs) An IMU segment is a 6D-vector composed of 3-axes accelerometer signals and 3-axes gyroscope signals. The first real-world task we evaluate (results in Figure 3) is to classify whether a 15-second IMU segment recorded at 128hz contains a fall, a “near” fall or showed “active daily living” (ADL). The dataset used in the open-source IMU Fall Detection dataset (IMUFD, Aziz et al. (2017)).

Few-shot fall detection is a pattern-recognition task - typically a fall shows up on the IMU as a big spike in magnitude on multiple axes. What makes the task hard is the inclusion of the hard negative class of “Near” falls, where the participants of the study pretend to trip but recover before actually falling, creating similarly large changes in magnitude on the IMU.

Activity recognition from IMUs A further real-world IMU task we evaluate is to classify whether a 15-second IMU segment is one of five activity classes: “sit”, “stand”, “stairs”, “walk” or “bike” (results in Figure 4). The dataset used in the open-source Heterogeneity Human Activity Recogni-

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

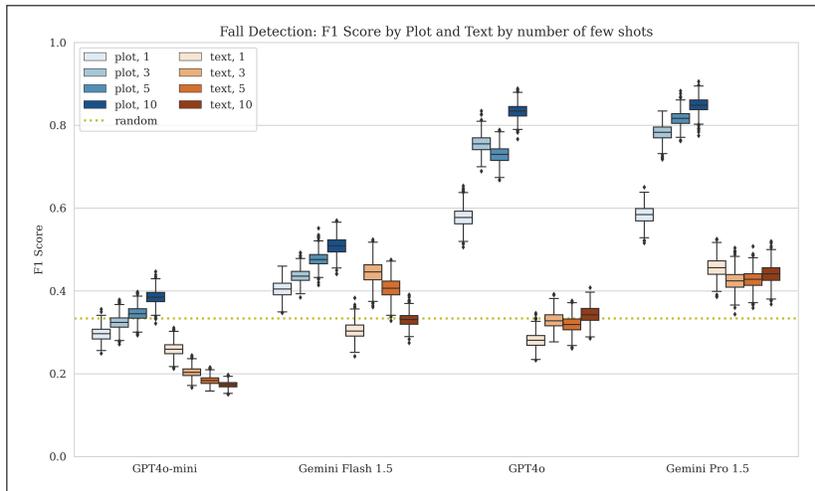


Figure 3: Results of fall detection task show consistently better plot performances across models and number of few-shots, with plot performance generally increasing with number of shots. The top plot models have 10-shot (sensitivity, specificity) as follows: Gemini Pro 1.5 - (0.84, 0.95) and GPT4o - (0.92, 0.81), compared to the state-of-the-art task-specific support-vector machine model reported by Aziz et al. (2017) which achieves (0.96, 0.96) (see Supplementary Table S1 for more details).

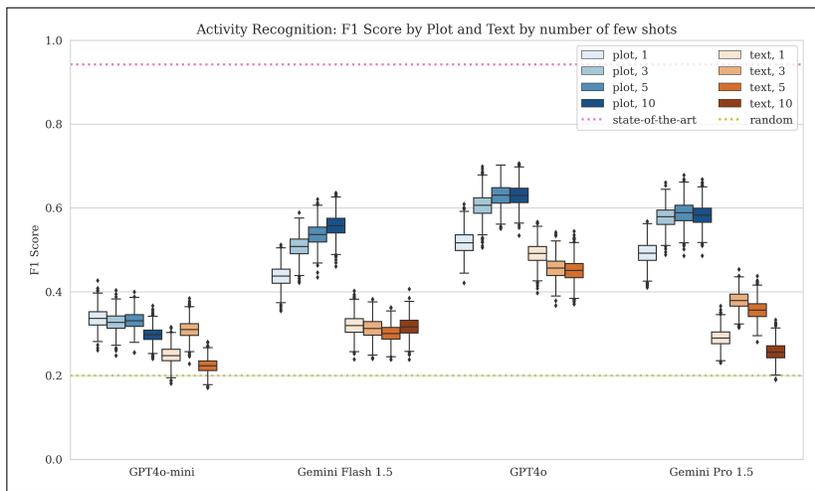


Figure 4: Results of activity detection task for all models across few-shot numbers (where context length allowed), showing overall improved performance for plots. The performance of the state-of-the-art deep-learning model reported by Kumar & Selvam (2022) is included for reference.

tion dataset (HHAR, Stisen et al. (2015)). As with Fall Detection we test performance with 1, 3, 5 and 10 few-shot examples, with 383 samples per model and number of few-shots. For this task the text representation of the 10 few-shot examples exceeded the GPT4o and GPT4o-mini 128k context windows, so these results are only shown for the Gemini models.

Activity recognition requires evaluating the entire IMU segment and correlating signals between different axes and sensors to determine the likely activity, as the noisy signals may only subtly change between “sit” and “stand” or “walk” and “stairs”. The HHAR dataset was deliberately collected to be heterogeneous containing data collected from four different types of smartphone and two different types of smartwatch. The classes in the dataset aren’t balanced so performance is reported using an F_1 score.

Readiness Estimating fitness readiness for a workout is a multicomponent task that involves assessment of health metrics, sleep, training load, and subjective feedback (Cosentino et al., 2024). Among those, training load analysis can be evaluated quantitatively and involves plot interpretation. Therefore, we framed the task as a binary classification problem (training load trending upwards or downwards) and use the calculated acute-chronic workload ratio (ACWR) to obtain ground truth labels.

ACWR is a ratio of acute training load (total training impulse, or TRIMP, over the past 7 days) divided by chronic training load (28-day average of acute load). An ACWR equal to 1.0 means that the user has exercised at the same intensity continuously over the past week compared to the month, below 1.0 means that they are trending downward, and above 1.0 means they are trending upwards. Precise ACWR calculation involves multiple mathematical operations, so we assess the model’s ability to understand the trend from monthly TRIMP values without explicit calculation.

We use training load data from 350 fitness case studies (Cosentino et al., 2024) and present it as tables or TRIMP bar plots (results shown in Figure 5). Each case study contains data from 30 consecutive days. We use a simplified version of the textual prompt and visualization from Cosentino et al. (2024) and do not split TRIMP in different heart rate zones. We tested Gemini 1.5 Pro and Gemini 1.5 Flash for both plot and text approaches as zero-shot tasks. Since this task involved analyzing just 30 data points we did not expect the plot prompt to excel here. Interestingly, models of different sizes showed opposite trends: Gemini 1.5 Pro had a slight increase in performance when using plots, while Gemini 1.5 Flash had a slight increase in performance when using text, though the magnitudes of the gains were small.

4.3 ABLATIONS

We considered a variety of text and plot ablations to confirm if there were any large gains. All ablations were performed on Gemini Pro 1.5. We used the function identification task to test for any performance differences; details, results and visualizations are reported in Section A.3.

4.4 COST

Using plots for time-series data can often be more cost-efficient and token-efficient. Token efficiency matters when your context is large and the context-window is limited; for example we needed to downsample our raw signals to fit them into the 128k context window for GPT4o(-mini), particularly with the large few-shot experiments (Section 4.2). For example, when using the Gemini API (Google, 2024), images account for 258 tokens if both dimensions are less than 384 x 384 pixels, after which 4 additional crops are added for a total of 1290 tokens. Text tasks can easily be 10x larger (e.g. 10-shot activity recognition Section 4.2) requiring more than the entire 128k context available. Depending on the task it may be possible to reduce the number of text tokens by further sub-sampling of the data, but this may result in reduced task performance. The optimal sampling rate may also be task- and dataset-specific.

Plot experiments also end up being cheaper. As an example, for our most expensive experiment on few-shot activity recognition, we can estimate the input token cost of our 5-shot experiment for both plot and text on GPT4o (OpenAI, 2024). The text version of this task required nearly 128k text tokens (costing \$0.32 per 5-shot question); by contrast, the plot version required 50 images (costing \$0.032), a 10x difference in overall costs for input tokens. In addition to being cheaper, the plotting approach also scales better for longer time-series, as the number of tokens required for the textual

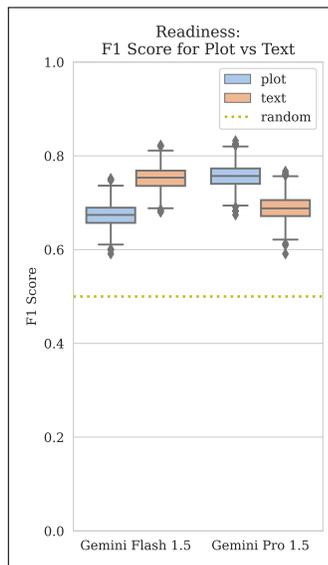


Figure 5: Results of readiness task for Gemini models only (as the dataset cannot be sent to other models), demonstrating approximate parity between the text and plot approaches.

486 approach grows with the number of data points while a plot of the same size will generate the same
487 number of tokens independent of the number of data points being plotted.

489 5 CONCLUSIONS AND FUTURE WORK

491 The key finding from our rigorous empirical evaluation is that engaging the vision encoder of a mul-
492 timodal foundation model through the use of plot representations leads to significant performance
493 and efficiency gains on time-series understanding tasks, compared with relying on the text encoder.
494 By processing data visually instead of textually, these models can better capture temporal patterns
495 and relationships. We established our results on synthetic data with well-controlled characteristics
496 and reasoning types, and also showed that this approach holds on real, noisy and complex tasks re-
497 lated to making sense of consumer health signals. This is analogous to the gains that humans benefit
498 from when visualising data (Card et al., 1999; Yalçin et al., 2016), though we do not claim that the
499 mechanisms by which visual data are processed by the models we study here are the same as those
500 with which humans process visual stimuli.

501 The method presented here is powerful in its simplicity and generalizability and relies on the native
502 capabilities of multimodal models requiring no additional training – we believe that it is particularly
503 useful when the following conditions are met:

- 505 • You want to use an off-the-shelf multimodal model to interpret your time-series data, as
506 might be the case in user-facing applications that rely on general models to understand a
507 broad range of potential user inputs including natural language.
- 508 • Your use-case is not restricted to a specific task or modality, and generalizability across
509 tasks is more important than accuracy on a single task. We showed that plots act as a gen-
510 eralizable time-series encoder across many tasks, even though they may not be better than
511 a task-specific encoder trained for one task. Training task-specific encoders for multimodal
512 models can be limited by availability of paired training data, compute and expertise.
- 513 • You don’t want to downsample your data. In many cases the textual representation of real
514 time-series outstrips the maximum context length, and so the plot-based approach is the
515 only way to present the data without downsampling.

516 Our focus in this work is specific to time-series understanding (i.e., reasoning about known data).
517 Forecasting is also important to time-series analysis and in the future we suspect that leveraging the
518 vision components of multimodal models might yield positive results in this area too.

519 In this work, all plots were generated by human-written code in order to avoid any bias. As such we
520 do not rigorously study what the optimal plotted form of a certain time-series might be for visual
521 understanding; this is likely to be a function of the exact downstream task or user request, but could
522 in theory be automated and forms the basis for future work. Looking forward, in real applications
523 we anticipate that plotting could be part of a tool-use framework, where the model is prompted to
524 choose how and when to plot the data, after which it uses the plot representation it created.

525 Lastly, further work remains in the explainability context – while we demonstrate empirically that
526 visual representations generally outperform textual representations of time-series data, we have not
527 probed why mechanistically this is so.

530 6 REPRODUCIBILITY STATEMENT

531 Our evaluations are run on publicly available models that have publicly available APIs. The exact
532 model versions used are detailed in Section 3.2.

534 Our structured prompting methods are detailed in 3.1 and we include the actual prompts and target
535 dataclasses used in Supplementary Section A.4.

536 All the data for the synthetic tasks is by definition synthesized and the detailed task descriptions in
537 Supplementary Section A.1.1 provide the necessary details to recreate these synthetic datasets.

538 The IMU Fall Detection Dataset (IMUFD, Aziz et al. (2017)) and Heterogeneity Human Activity
539 Recognition dataset (HHAR, Stisen et al. (2015)) used respectively for the fall detection and ac-

540 tivity recognition tasks are both publicly available. Pre- and post-processing steps are detailed in
541 Supplementary Section A.1.2.

542 The dataset for the Readiness task is not currently publicly available. However the task details in
543 Section 4.2 would enable reproduction of the results with access to a comparable dataset.
544

545 REFERENCES

546 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
547 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
548 report. *arXiv preprint arXiv:2303.08774*, 2023.
549

550 Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen,
551 Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al.
552 Chronos: Learning the language of time series. *arXiv preprint arXiv:2403.07815*, 2024.
553

554 Omar Aziz, Magnus Musngi, Edward J Park, Greg Mori, and Stephen N Robinovitch. A com-
555 parison of accuracy of fall detection algorithms (threshold-based vs. machine learning) using
556 waist-mounted tri-axial accelerometer signals from a comprehensive set of falls and non-fall tri-
557 als. *Medical & biological engineering & computing*, 55:45–55, 2017.
558

559 J Martin Bland and Douglas G Altman. Multiple significance tests: the bonferroni method. *Bmj*,
560 310(6973):170, 1995.

561 Yifu Cai, Mononito Goswami, Arjun Choudhry, Arvind Srinivasan, and Artur Dubrawski. Jolt:
562 Jointly learned representations of language and time-series. In *Deep Generative Models for*
563 *Health Workshop NeurIPS 2023*, 2023.
564

565 Yifu Cai, Arjun Choudhry, Mononito Goswami, and Artur Dubrawski. Timeseriesexam: A time
566 series understanding exam. *arXiv preprint arXiv:2410.14752*, 2024.

567 Stuart K Card, Jock Mackinlay, and Ben Shneiderman. *Readings in information visualization: using*
568 *vision to think*. Morgan Kaufmann, 1999.
569

570 Nimeesha Chan, Felix Parker, William Bennett, Tianyi Wu, Mung Yao Jia, James Fackler, and Kimia
571 Ghobadi. Medtsllm: Leveraging llms for multimodal medical time series analysis. *arXiv preprint*
572 *arXiv:2408.07773*, 2024.

573 Daniel Corin. VLMs Aren’t Blind. [https://www.danielcorin.com/posts/2024/
574 vlms-arent-blind/](https://www.danielcorin.com/posts/2024/vlms-arent-blind/), 2024. Accessed September 24, 2024.
575

576 Justin Cosentino, Anastasiya Belyaeva, Xin Liu, Nicholas A Furlotte, Zhun Yang, Chace Lee, Erik
577 Schenck, Yojan Patel, Jian Cui, Logan Douglas Schneider, et al. Towards a personal health large
578 language model. *arXiv preprint arXiv:2406.06474*, 2024.

579 Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
580 time-series forecasting. *arXiv preprint arXiv:2310.10688*, 2023.
581

582 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,
583 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly
584 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

585 Rohit Girdhar, Alaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand
586 Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of*
587 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15180–15190, 2023.

588 Google. Gemini API Cost. [https://cloud.google.com/vertex-ai/
589 generative-ai/docs/multimodal/image-understanding](https://cloud.google.com/vertex-ai/generative-ai/docs/multimodal/image-understanding), 2024. Accessed
590 September 26, 2024.
591

592 Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski.
593 Moment: A family of open time-series foundation models. *arXiv preprint arXiv:2402.03885*,
2024.

- 594 Nate Gruver, Marc Finzi, Shikai Qiu, and Andrew G Wilson. Large language models are zero-shot
595 time series forecasters. *Advances in Neural Information Processing Systems*, 36, 2024.
- 596
- 597 Prabhat Kumar and Suresh Selvam. Deeptranshhar: Inter-subjects heterogeneous activity recog-
598 nition approach in the non-identical environment using wearable sensors. *National Academy*
599 *Science Letters*, 2022.
- 600 Fangyu Liu, Julian Eisenschlos, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee,
601 Mandar Joshi, Wenhui Chen, Nigel Collier, and Yasemin Altun. DePlot: One-shot visual lan-
602 guage reasoning by plot-to-table translation. In *Findings of the Association for Computational*
603 *Linguistics: ACL 2023*, pp. 10381–10399, 2023.
- 604 Seungwhan Moon, Andrea Madotto, Zhaojiang Lin, Aparajita Saraf, Amy Bearman, and Babak
605 Damavandi. IMU2CLIP: Language-grounded Motion Sensor Translation with Multimodal Con-
606 trastive Learning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*,
607 pp. 13246–13253, 2023.
- 608
- 609 OpenAI. OpenAI API Pricing. <https://openai.com/api/pricing/>, 2024. Accessed
610 September 26, 2024.
- 611 Arjun Panickssery, Samuel R Bowman, and Shi Feng. Llm evaluators recognize and favor their own
612 generations. *arXiv preprint arXiv:2404.13076*, 2024.
- 613
- 614 Daiyi Peng. Langfun. <https://github.com/google/langfun> (Version 0.0.1), September
615 2023.
- 616 Shraman Pramanick, Rama Chellappa, and Subhashini Venugopalan. Spiga: A dataset for multi-
617 modal question answering on scientific papers. *arXiv preprint arXiv:2407.09413*, 2024.
- 618
- 619 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
620 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
621 models from natural language supervision. In *International conference on machine learning*, pp.
622 8748–8763. PMLR, 2021.
- 623 Pooyan Rahmazadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision
624 language models are blind. *arXiv preprint arXiv:2407.06581*, 2024.
- 625
- 626 Dimitris Spathis and Fahim Kawsar. The first step is the hardest: pitfalls of representing and tok-
627 enizing temporal data for large language models. *Journal of the American Medical Informatics*
628 *Association*, 31(9):2151–2158, 07 2024. ISSN 1527-974X. doi: 10.1093/jamia/ocae090. URL
629 <https://doi.org/10.1093/jamia/ocae090>.
- 630 Allan Stisen, Henrik Blunck, Sourav Bhattacharya, Thor Siiger Prentow, Mikkel Baun Kjærgaard,
631 Anind Dey, Tobias Sonne, and Mads Møller Jensen. Smart devices are different: Assessing and
632 mitigating mobile sensing heterogeneities for activity recognition. In *Proceedings of the 13th ACM*
633 *conference on embedded networked sensor systems*, pp. 127–140, 2015.
- 634 Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi
635 Wu, Haotian Liu, Sathika Malladi, et al. Charxiv: Charting gaps in realistic chart understanding
636 in multimodal llms. *arXiv preprint arXiv:2406.18521*, 2024.
- 637
- 638 Frank Wilcoxon. Individual comparisons by ranking methods. In *Breakthroughs in statistics:*
639 *Methodology and distribution*, pp. 196–202. Springer, 1945.
- 640 Christopher Wimmer and Navid Rekabsaz. Leveraging vision-language models for granular market
641 change prediction. *arXiv preprint arXiv:2301.10166*, 2023.
- 642
- 643 Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sa-
644 hoo. Unified training of universal time series forecasting transformers. *arXiv preprint*
645 *arXiv:2402.02592*, 2024.
- 646 M Adil Yalçın, Niklas Elmquist, and Benjamin B Bederson. Cognitive stages in visual data explo-
647 ration. In *Proceedings of the Sixth Workshop on Beyond Time and Errors on Novel Evaluation*
Methods for Visualization, pp. 86–95, 2016.

648 Xiyuan Zhang, Ranak Roy Chowdhury, Rajesh K. Gupta, and Jingbo Shang. Large language models
 649 for time series: A survey. In Kate Larson (ed.), *Proceedings of the Thirty-Third International Joint*
 650 *Conference on Artificial Intelligence, IJCAI-24*, pp. 8335–8343. International Joint Conferences
 651 on Artificial Intelligence Organization, 8 2024. doi: 10.24963/ijcai.2024/921. URL <https://doi.org/10.24963/ijcai.2024/921>. Survey Track.

654 A SUPPLEMENTARY INFORMATION

655 A.1 DETAILED TASK DESCRIPTIONS

656 In this section we provide further details of each task implementation to assist with reproducibility
 657 of results. We also provide Python code snippets for synthetic data generation and plotting.

658 A.1.1 SYNTHETIC TASKS

659 **Functional form identification**

- 660 • We generate (x, y) series of linear, quadratic, cubic, exponential and periodic functions
 661 with variable number of points and noise (injected into the function domains) over the
 662 range $x \in [-10, 10]$.
- 663 • We perform five repeats across different numbers of points (50, 500, 1000 and 2500) and
 664 noise levels (0.0, 0.5, 1.0, 2.0 and 5.0), giving 500 samples per model across number of
 665 points, noise level, function type and random replica dimensions.
- 666 • These results are then passed either as a stringified series of x and y vectors (“text” task),
 667 or as a matplotlib figure (the “plot” task) to the model.
- 668 • This task only requires that the model is able to understand and label the global longitudinal
 669 trend of the function, without overly needing to reason about magnitudes.

```

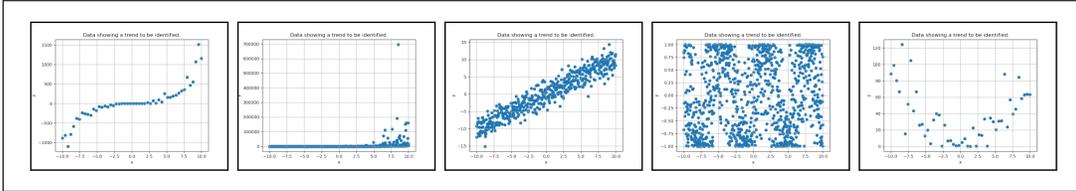
670 def generate(rng: np.random.Generator, func_type: str,
671             x_range: tuple[int, int], num_points: int,
672             noise_level: float):
673     x = np.linspace(x_range[0], x_range[1], num_points)
674     noise = rng.normal(0, noise_level, num_points)
675     if func_type == "exponential":
676         y = np.exp(x + noise)
677     elif func_type == "periodic":
678         y = np.sin(x + noise)
679     elif func_type == "quadratic":
680         y = (x + noise) ** 2
681     elif func_type == "linear":
682         y = x + noise
683     elif func_type == "cubic":
684         y = (x + noise) ** 3
685     else:
686         raise ValueError("Invalid function type %s" % func_type)
687     return x, y
  
```

688 Functional form identification - Data generation code

```

689 fig = plt.figure(figsize=(6.4, 4.8)) # Rendered at 100 dpi
690 ax = fig.gca()
691 ax.scatter(xs, ys)
692 ax.set_title("Data showing a trend to be identified.")
693 ax.set_xlabel("x")
694 ax.set_ylabel("y")
695 ax.grid(True)
  
```

Functional form identification - Matplotlib plotting code



Supplementary Figure S1: Plots used for functional form identification task examples chosen at random representing at various noise levels (from left to right): cubic, exponential, linear, periodic, quadratic.

Correlation of two lines

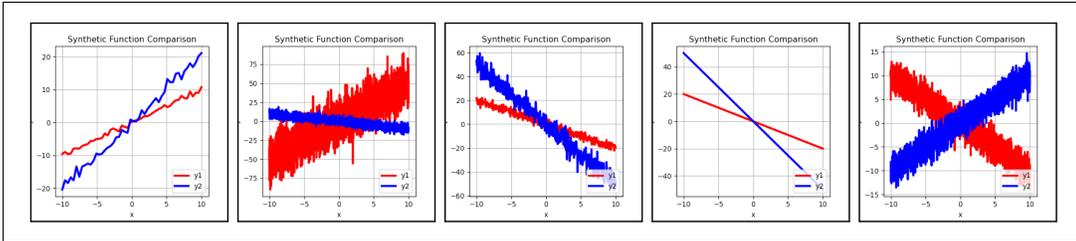
- We generate (x, y_1) and (x, y_2) series that represent linear functions $y = m \cdot x$ with variable pairs of slopes $(m_1, m_2) \in \{(1, 2), (-1, 1), (5, -1), (-2, -5), (-3, 2), (2, 3)\}$.
- We perform trials across different numbers of points (50, 500, 1000 and 2500) and noise levels (0, 0.25, 0.5, 1, 1.5, 2.0, 3.0 and 5.0) over the range $x \in [-10, 10]$, with 192 samples per model across number of points, noise level and random replica dimensions.
- These results are then passed either as a stringified series of x , y_1 and y_2 vectors (“text” task), or as a matplotlib figure (the “plot” task) to the model. The model is then asked to classify whether the two lines $y_1(x)$ and $y_2(x)$ are positively or negatively correlated.
- This task requires first understanding two global trends, and then comparing them with each other.

```
def generate(rng: np.random.Generator, x_range: tuple[int, int],
            slope_coeffs: tuple[int, int], num_points: int,
            noise_level: float):
    x = np.linspace(x_range[0], x_range[1], num_points)
    y1_noise = rng.normal(0, noise_level, num_points)
    y2_noise = rng.normal(0, noise_level, num_points)
    y1 = slope_coeffs[0] * (x + y1_noise)
    y2 = slope_coeffs[1] * (x + y2_noise)
    return x, y1, y2
```

Correlation of two lines - Data generation code

```
fig = plt.figure(figsize=(4, 4)) # Rendered at 96 dpi
ax = fig.gca()
ax.set_title("Synthetic Function Comparison")
ax.plot(xs, y1s, color="red", linewidth=3, label="y1")
ax.plot(xs, y2s, color="blue", linewidth=3, label="y2")
ax.legend(loc="lower right")
ax.set_xlabel("x")
ax.set_ylabel("y")
ax.grid(True)
```

Correlation of two lines - Matplotlib plotting code



Supplementary Figure S2: Plots used for correlation of two lines task examples chosen at random representing positive and negative correlations at various noise levels.

2D cluster counting

- We generate a series of points corresponding to N distinct clusters, parameterised by the standard deviation from the cluster center which also controls the difficulty of the task.
- The cluster centers are chosen randomly, and we enforce a minimum distance between clusters.
- We perform five repeats across different levels of standard deviation (0.025, 0.05 and 0.075), different levels of number of points per clusters (5, 50 and 100) and with the number of clusters from 1 to 9, giving 405 samples per model across standard deviation, number of clusters, number of points per clusters and random replica dimensions.
- Extending the correlation task, this task now requires that the model is able to simultaneously identify and keep separate track of N different patterns.

```

def _generate_centers(num, rng, radius = 1.0, margin = 0.1):
    def _has_close_centers(center_coordinates, distance_threshold):
        distances = scipy.spatial.distance.cdist(
            center_coordinates, center_coordinates, "euclidean")
        mask = np.triu(np.ones_like(distances), k=1).astype(bool)
        return np.any(distances[mask] < distance_threshold)
    for _ in range(10000):
        coordinates = [
            (radius * x, radius * y)
            for x, y in zip(
                rng.uniform(-radius+margin, radius-margin, num),
                rng.uniform(-radius+margin, radius-margin, num),
            )
        ]
        if not _has_close_centers(coordinates, radius * 0.3):
            return coordinates
        raise ValueError("Could not find well separated centers")

def generate(rng: np.random.Generator, seed: int,
            cluster_points: int, cluster_count: int, cluster_std: float):
    generated_samples, _, _ = sklearn.datasets.make_blobs(
        n_samples=cluster_points * cluster_count,
        centers=_generate_centers(cluster_count, rng),
        cluster_std=cluster_std,
        random_state=seed,
        return_centers=True,
        center_box=(-1, 1),
    )
    return generated_samples[:, 0], generated_samples[:, 1]

```

2D cluster counting - Data generation code

810

```

811 fig = plt.figure(figsize=(8, 8)) # Rendered at 96 dpi
812 ax = fig.gca()
813 ax.scatter(xs, ys, c="black", s=50)
814 ax.set_title("Synthetic Clustered Data")
815 ax.set_xlabel("x")
816 ax.set_ylabel("y")
817 ax.set_xlim(-1, 1)
818 ax.set_ylim(-1, 1)
819 ax.grid(True)

```

820

2D cluster counting - Matplotlib plotting code

821

822

823

824

825

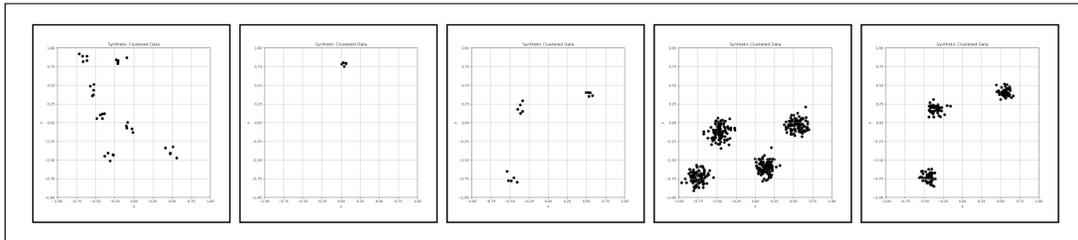
826

827

828

829

830



831

Supplementary Figure S3: Plots used for 2D cluster counting task examples chosen at random representing varying cluster parameterizations.

832

833

834

Derivative identification

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

```

def generate(rng: np.random.Generator, func_type: str,
            x_range: tuple[int, int], num_points: int,
            noise_level: float):
    x = np.linspace(x_range[0], x_range[1], num_points)
    noise = rng.normal(0, noise_level, num_points)
    if func_type == "exponential":
        y = np.exp(x + noise)
        dy = np.exp(x + noise)
    elif func_type == "periodic":
        y = np.sin(x + noise)

```

```

864     dy = np.cos(x + noise)
865     elif func_type == "quadratic":
866         y = (x + noise) ** 2
867         dy = 2.0 * (x + noise)
868     elif func_type == "linear":
869         y = x + noise
870         dy = np.ones(len(x)) + noise
871     elif func_type == "cubic":
872         y = (x + noise) ** 3
873         dy = 3.0 * (x + noise) ** 2
874     else:
875         raise ValueError("Invalid function type %s" % func_type)
876     return x, y, dy

```

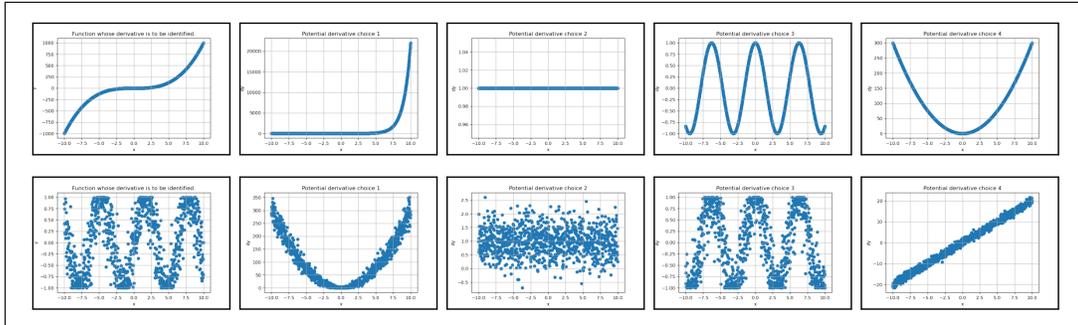
Derivative identification - Data generation code

```

879 fig = plt.figure(figsize=(6, 4)) # Rendered at 100 dpi
880 ax = fig.gca()
881 ax.scatter(x, y)
882 ax.set_title(
883     f"Potential derivative choice {choice_idx}" if is_derivative
884     else "Function whose derivative is to be identified.")
885 ax.set_xlabel("x")
886 ax.set_ylabel("dy" if is_derivative else "y")
887 ax.grid(True)

```

Derivative identification - Matplotlib plotting code



Supplementary Figure S4: Each row shows plots used for a randomly selected derivative identification task example. The leftmost plot in the row is the function to identify the derivative of, and the remaining plots are the four multiple choices.

Quadratic derivative identification

- We generate (x, y) series of quadratics (of form $y(x) = A \cdot x^2$) and their derivatives $y'(x)$ over a range of scales $A \in \{-10, -5, -1, 1, 5, 10\}$ over the range $x \in [-10, 10]$.
- We perform five repeats across different numbers of points (50, 500, 1000 and 2000) and noise levels (0.0, 0.5, 1.0, 2.0 and 5.0), with 600 samples per model and number of few-shots across number of points, noise level, function type and random replica dimensions.
- There are four multiple choices per sample, and each choice is a random selection of derivatives of quadratic functions with a range of scales $A \in \{-20, -15, -10, -5, -1, 1, 5, 10, 15, 20\}$, with the same noise level and number of points as the quadratic function in question.

- We create few-shot examples in the same manner as the main task, with the quadratic functions being randomly sampled from a range of scales $A \in \{-20, -15, -3, 3, 15, 20\}$ with 0 noise and 50 data points.
- These results are then passed either as a stringified series of x and y vectors (“text” task), or as a matplotlib figure (the “plot” task) to the model.
- This hard variant of the derivatives task introduces a new requirement for the model to reason about function magnitudes as all the possible choices are the correct functional form (i.e., linear).

```

def generate(rng: np.random.Generator, quadratic_scale: float,
            x_range: tuple[int, int], num_points: int,
            noise_level: float):
    noise = rng.normal(0, noise_level, num_points)
    x = np.linspace(x_range[0], x_range[1], num_points)
    y = quadratic_scale * (x + noise) ** 2
    dy = 2.0 * quadratic_scale * (x + noise)
    return x, y, dy

```

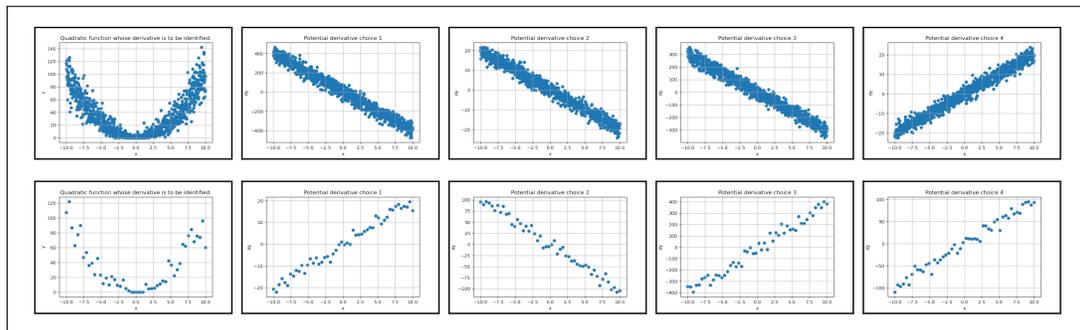
Quadratic derivative identification - Data generation code

```

fig = plt.figure(figsize=(6, 4)) # Rendered at 100 dpi
ax = fig.gca()
ax.scatter(x, y)
ax.set_title(
    f"Potential derivative choice {choice_idx}" if is_derivative
    else "Quadratic function whose derivative is to be identified.")
ax.set_xlabel("x")
ax.set_ylabel("dy" if is_derivative else "y")
ax.grid(True)

```

Quadratic derivative identification - Matplotlib plotting code



Supplementary Figure S5: Each row shows plots used for a randomly selected quadratic derivative identification task example. The leftmost plot in the row is the function to identify the derivative of, and the remaining plots are the four multiple choices.

A.1.2 REAL-WORLD TASKS

Fall detection

The task is hard to define as zero-shot, since there isn’t a natural way of explaining the plots and text, so we frame this as a few-shot task. We test few-shot examples with 1, 3, 5 and 10 examples, with 480 samples for each body location per model and number of few-shots. Model API errors were ignored as long as at least one body location succeeded for that example.

Pre-processing and post-processing

972 Due to context-window limitations for the text-only task, we apply a 1D average pool with a kernel
 973 size of 10 and stride of 10 to the data. This allows us to use up to 10 few-shot examples in the
 974 text-only tasks and still fit into the GPT4o and GPT4o-mini 128k context window.

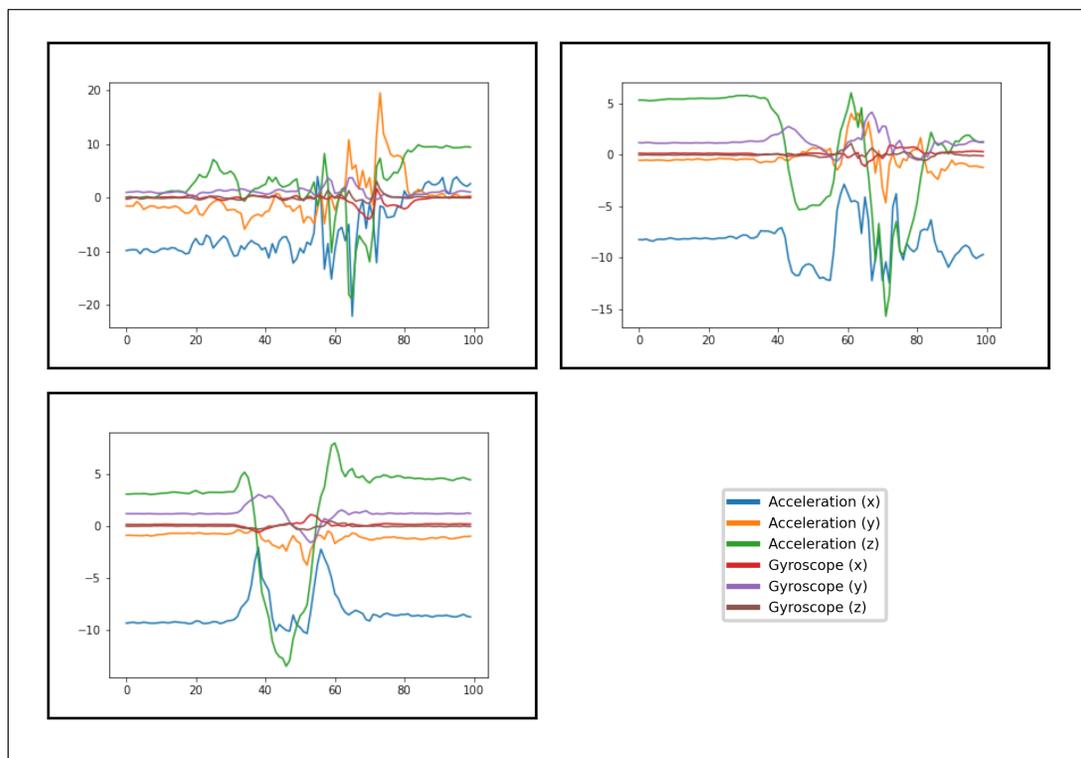
975 The IMUFD dataset provides multiple IMUs from 7 body locations. We consider a subset of head,
 976 waist, left and right thigh. We sample IMUs from each as separate examples - when evaluating either
 977 text or plot approach, the final prediction is a majority vote from the predictions from each of these
 978 body parts.

979 The dataset is stratified into train (20%) and test (80%) based on participant ID. Few-shot examples
 980 are chosen from the “train” set while eval data is chosen from the “test” set. This ensures the model
 981 never sees any data from the same participant.
 982

983 *Plotting*

```
984 fig = plt.figure(figsize=(6, 4)) # Rendered at 100 dpi
985 ax = fig.gca()
986 for i in range(6):
987     ax.plot(example[i])
988
```

989 Fall detection - Matplotlib plotting code



1016 Supplementary Figure S6: Illustrative plots of example entries from the IMU Fall Detection Dataset
 1017 (IMUFD, Aziz et al. (2017)) that were correctly categorized by Gemini Pro (10-shot). Top-Left:
 1018 Fall. Top-Right: Near Fall. Bottom-Left: active daily living (ADL). Examples were collected with
 1019 an IMU sensor at the waist and each of the three acceleration and gyroscope axes were plotted as a
 1020 different series on the same plot. The legend was not provided to the model, but is shown here to
 1021 aid visualisation.

1022 *Comparison with task-specific model*

1023 In Table S1 we compare the results from most performant foundation models studied here (GPT4o
 1024 and Gemini Pro 1.5) with a task-specific model reported in (Aziz et al., 2017). While we don't
 1025

1026 achieve the same sensitivity and specificity as expected for a general model, our plot results are of
 1027 the same order of magnitude.

1029 Model	Sensitivity	Specificity
1030 GPT4o, plot, 10-shot	0.92	0.81
1031 GPT4o, text, 10-shot	0.70	0.49
1032 Gemini Pro 1.5, plot, 10-shot	0.84	0.95
1033 Gemini Pro 1.5, text, 10-shot	0.61	0.70
1034 Task-specific SVM	0.96	0.96

1035
 1036 Supplementary Table S1: Comparison of most performant foundation models with task-specific fall
 1037 detection support vector machine (SVM) as reported in (Aziz et al., 2017).

1038 **Activity recognition**

1039 *Pre-processing*

1040
 1041 As with the fall detection data, we apply a 1D average pool over the raw data to limit the number of
 1042 text tokens. As the raw IMU sample rates varied between 70-200Hz the kernel size, and matching
 1043 stride, was chosen to target a downsampled frequency of 10Hz for each example. We select few-shot
 1044 examples using leave-one-out cross-validation at the dataset user level to maximise the number of
 1045 examples used for validation while ensuring we exclude any few-shot examples from the same user,
 1046 even those from different device types.

1047
 1048 The raw HHAR dataset consists of examples with varying durations and longer examples typically
 1049 contain multi-second gaps with no samples from the IMU. We chunk each raw example by splitting
 1050 on any gaps longer than 2 seconds and take a central 15 second crop of the longest chunk to create
 1051 the examples used in this study. Any examples with mismatching sample rates for the accelerometer
 1052 and gyroscope are filtered out. The raw labels “stairsup” and “stairsdown” are coalesced into a
 1053 single “stairs” class.

1054 *Plotting*

1055
 1056 For the activity recognition task, we plot accelerometer and gyroscope signals separately – we find
 1057 empirically that plotting the signals separately has a marginally beneficial effect on performance
 1058 over plotting together (see Supplementary Figure S8).

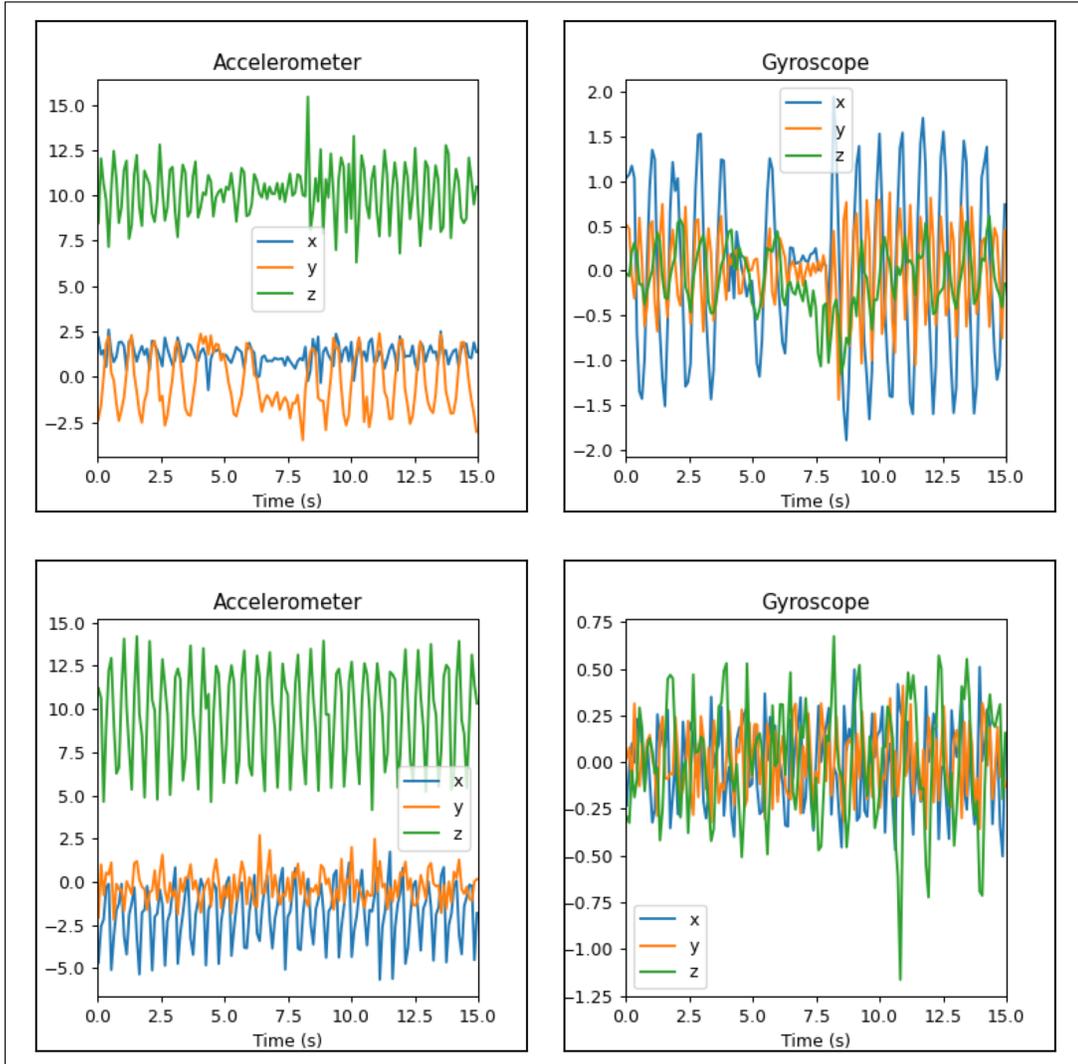
```

1059 figs = []
1060 for sensor_label, idxs in [
1061     ("Accelerometer", [1, 2, 3]),
1062     ("Gyroscope", [4, 5, 6])]:
1063     fig = plt.figure(figsize=(4, 4)) # Rendered at 90 dpi
1064     ax = fig.gca()
1065     ax.set_title(sensor_label)
1066     for axis, idx in zip(["x", "y", "z"], idxs):
1067         ax.plot(example[0], example[idx], label=axis)
1068     ax.legend()
1069     ax.set_xlabel("Time (s)")
1070     ax.set_xlim(0, 15)
1071     figs.append(fig)
  
```

1072 Activity recognition - Matplotlib plotting code

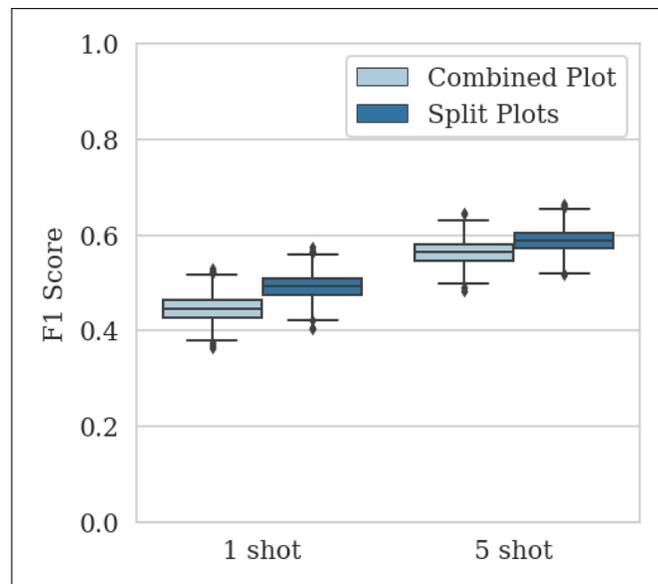
1073
 1074
 1075
 1076
 1077
 1078
 1079

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133



Supplementary Figure S7: Illustrative plots of example entries from the Heterogeneity Human Activity Recognition dataset (HHAR, Stisen et al. (2015)) that were correctly categorized by Gemini Pro (10-shot). Top: Bike. Bottom: Walk. These examples were collected from a Nexus4 device. The signals for each sensor type are plotted separately.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187



Supplementary Figure S8: Performance comparison between combining all 6 IMU axes in a single plot vs splitting the accelerometer and gyroscope into distinct plots for the Activity Recognition task. Results shown for 1-shot and 5-shot on Gemini Pro indicate a marginal performance increase when the data is split into 2 plots, which is reduced when moving from 1-shot to 5-shot.

A.2 FURTHER RESULTS

We first present the relative differences in task performance between plot and text approaches in Supplementary Table S2.

Task (Metric)	Shots	GPT4o-mini	Gemini Flash 1.5	GPT4o	Gemini Pro 1.5
Functional form id. (Accuracy)	0	70.8	82.4	122.7	11.1
Correlation of two lines (Accuracy)	0	66.7	42.9	20.0	50.0
2D cluster counting (MAE)	0	55.8	65.5	86.4	85.0
Derivative id. (Accuracy)	0	40.0	31.6	-5.3	-0.0
Quadratic derivative id. (Accuracy)	0	100.0	85.7	-37.5	47.4
	1	37.5	107.1	23.8	113.3
	2	92.9	100.0	8.0	128.6
	3	108.3	112.5	9.1	172.7
Fall detection (F_1)	1	11.6	36.5	114.4	27.3
	10	120.6	49.1	153.0	92.4
Activity detection (F_1)	1	36.0	36.5	6.3	69.3
	5	48.0	77.2	39.7	64.6
Readiness (F_1)	0	–	-10.8	–	10.0

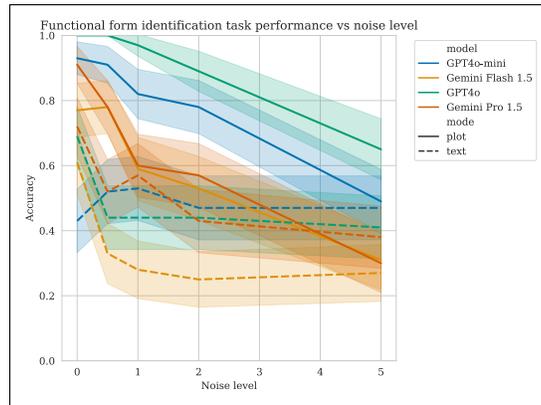
Supplementary Table S2: Percent differences in median plot performances relative to median text performances.

Next, we present results of the synthetic task performances as functions of the various dataset parameters that control the difficulty of the task example. The variable model responses over different combinations of the dataset parameters described in this section create the distributions shown in Figure 1. Depending on the task, these dataset parameters include:

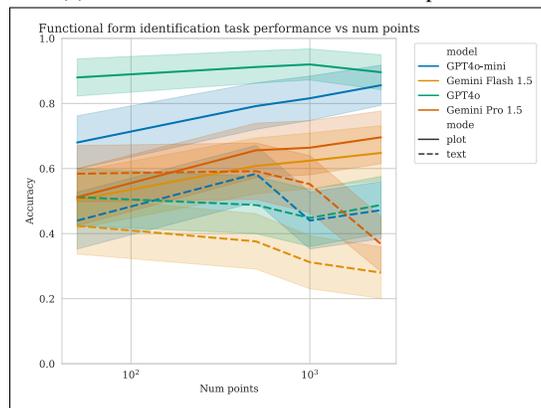
- **Number of points:** the number of points per series. For all tasks except for 2D cluster counting, this means the number of function samples between $x = -10$ and $x = 10$. For the cluster counting task, this is simply the number of points per cluster.
- **Noise level:** the amount of noise injected into the functions, i.e. $y = y(x + noise.level)$. This parameter is relevant for all synthetic tasks except the 2D cluster counting.
- **Standard deviation:** for the 2D cluster counting task, this controls the tightness of each cluster.

We also show the performance of each task over the set of possible correct answers specific to that task.

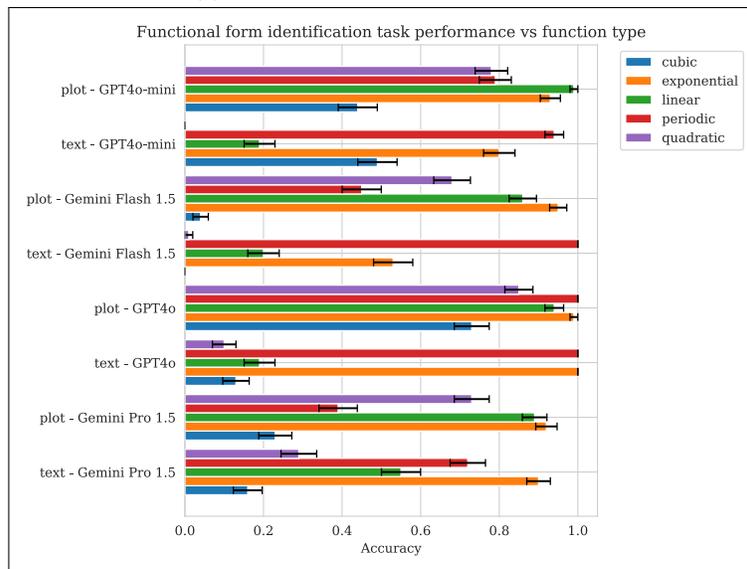
A.2.1 FUNCTIONAL FORM IDENTIFICATION



(a) Results as a function of number of points



(b) Results as a function of noise level

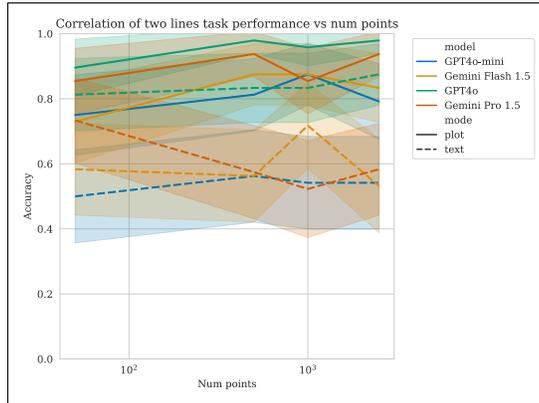


(c) Results as a function of true function type

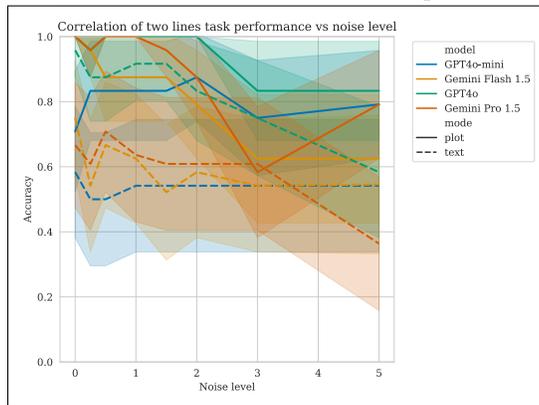
Supplementary Figure S9: Functional form identification performance over various dataset parameters

A.2.2 CORRELATION OF TWO LINES

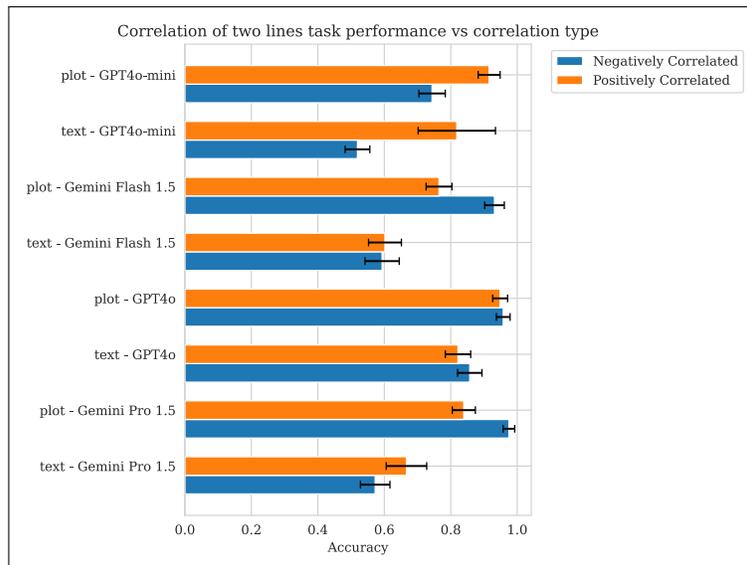
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349



(a) Results as a function of number of points



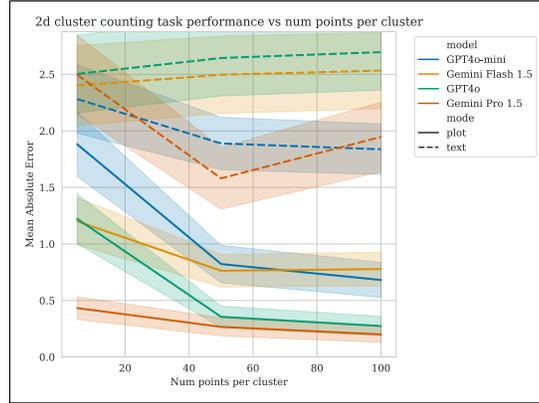
(b) Results as a function of noise level



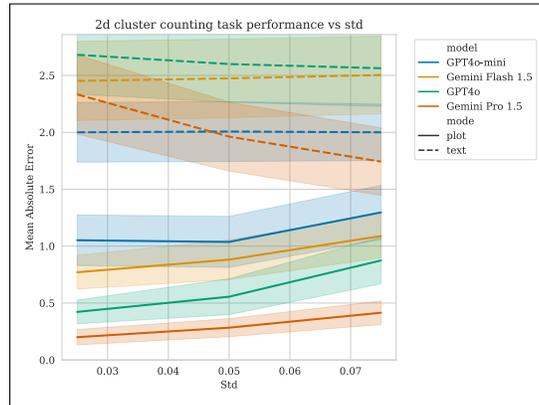
(c) Results as a function of true correlation type

Supplementary Figure S10: Correlation of two lines performance over various dataset parameters

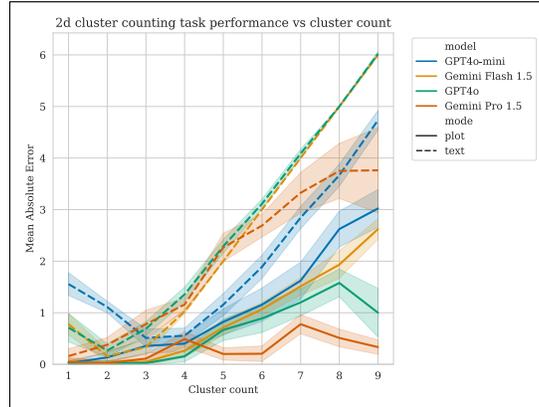
A.2.3 2D CLUSTER COUNTING



(a) Results as a function of number of points



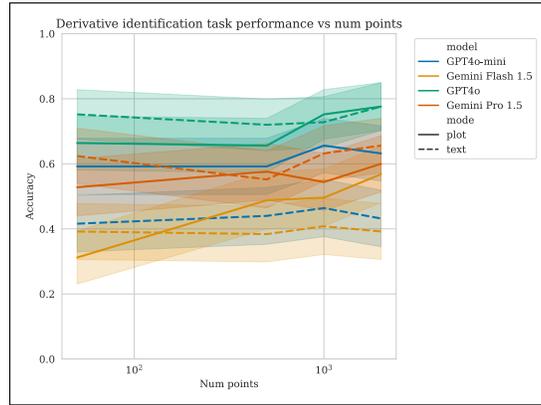
(b) Results as a function of cluster standard deviation



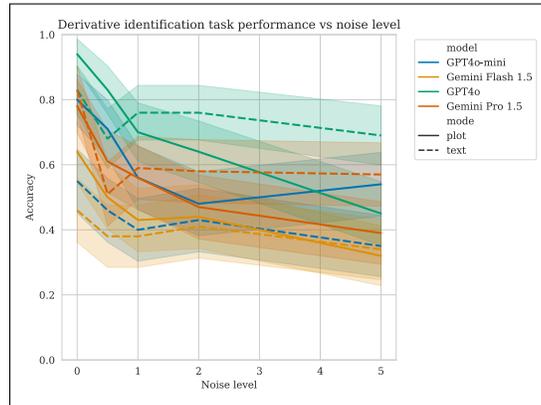
(c) Results as a function of true number of clusters

Supplementary Figure S11: 2D cluster counting performance over various dataset parameters (lower is better)

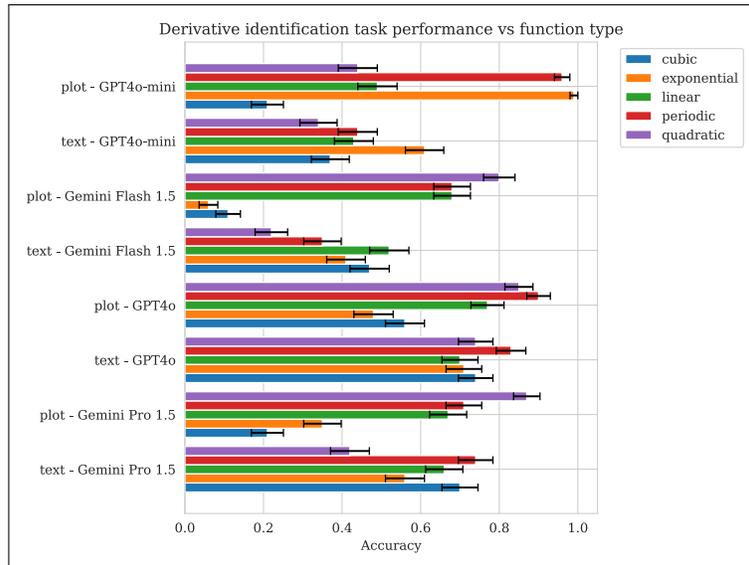
A.2.4 DERIVATIVE IDENTIFICATION



(a) Results as a function of number of points



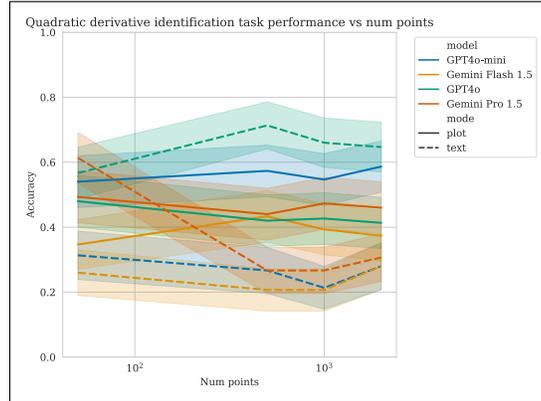
(b) Results as a function of noise level



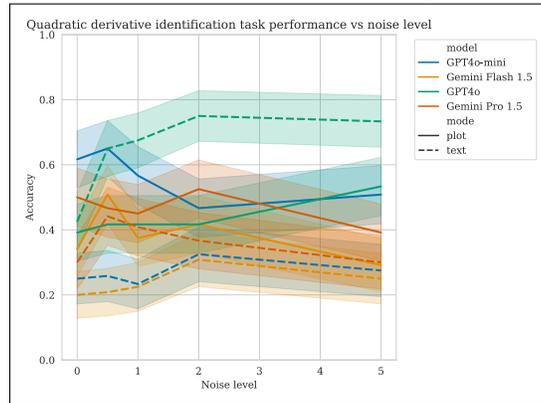
(c) Results as a function of input function class

Supplementary Figure S12: Derivative identification performance over various dataset parameters

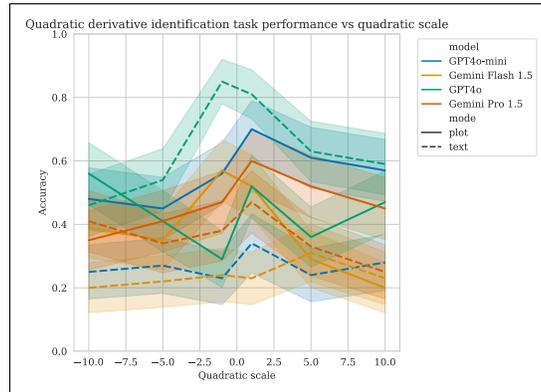
A.2.5 QUADRATIC DERIVATIVE IDENTIFICATION



(a) Results as a function of number of points



(b) Results as a function of noise level



(c) Results as a function of input quadratic scale (A in $y = A \cdot x^2$)

Supplementary Figure S13: Zero-shot quadratic derivative identification performance over various dataset parameters

A.3 ABLATIONS

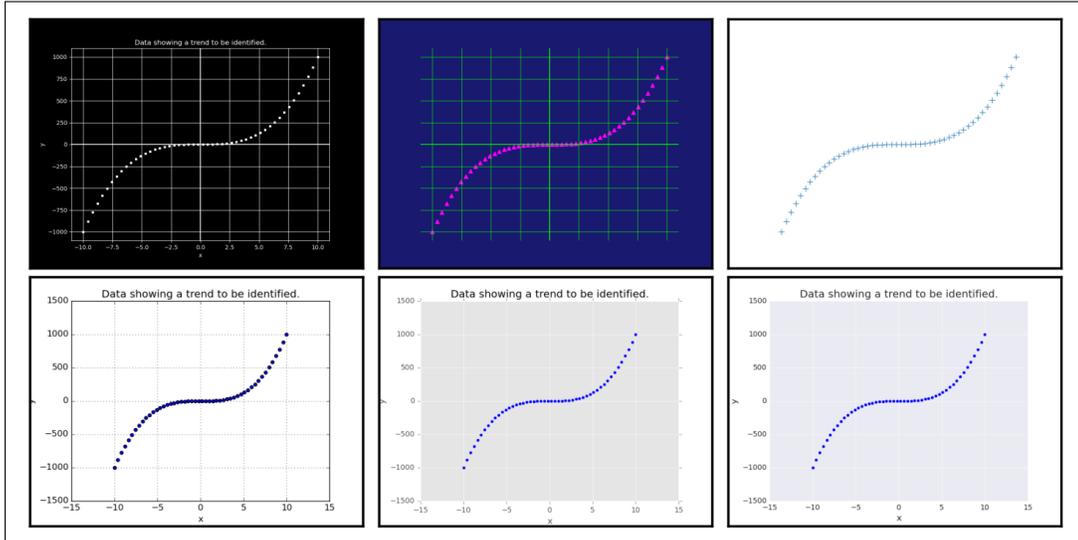
A.3.1 METHODOLOGY

We perform the ablations described here on the functional form identification task using Gemini Pro 1.5.

For the text-only task we tested the effect of comma versus space separation of numbers. We were also guided by the approach in LLMTime (Gruver et al., 2024) and tested the effect of fixed precision of floating point numbers (2, 4, 8, 16) and rescaling the input numbers.

For the plot tasks, we considered the following ablations, with Supplementary Figure S14 showing examples:

- Resolution (in dpi): 25, 50, 100, 200, 400
- Figure size: (3.5, 3.5), (4, 3), (7, 7), (8, 6), (12, 12)
- Plot style: Default, classic, ggplot, seaborn-whitegrid, seaborn-darkgrid
- Different color palettes for text, background and scatter
- Marker types: circle, square, triangle, x-mark, plus
- Marker sizes: small (10), medium (50), large (100)
- Plot components: all (title, axis labels, spines, ticks, grid, axes), minimal (grid and axes only) or none
- Temperature: 0.0, 0.1, 0.3, 0.55, 1.0



Supplementary Figure S14: A sample of the different ablations of plots used for a cubic function identification task.

We also test the combined plot and text version of the task, across both possible prompt orderings (i.e., text followed by plot, and plot followed by text).

Table cells contain mean accuracies with 95% confidence interval derived from 1,000 bootstrap repeats in brackets; rows and columns refer to different combinations of ablation and dataset parameters. We note that generally speaking, nothing had a strong positive effect on the results.

A.3.2 TEXT-BASED RESULTS

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598

Separator	50	500	1000	2500
Comma and space	0.63 (0.59, 0.67)	0.50 (0.46, 0.54)	0.50 (0.46, 0.54)	0.36 (0.32, 0.40)
Space	0.62 (0.58, 0.66)	0.54 (0.50, 0.59)	0.48 (0.44, 0.53)	0.36 (0.32, 0.40)

Supplementary Table S3: Text-based functional form identification ablation results: modifying value separator for different numbers of points.

Separator	0.0	0.5	1.0	2.0	5.0
Comma and space	0.71 (0.66, 0.75)	0.49 (0.44, 0.54)	0.44 (0.39, 0.49)	0.46 (0.41, 0.51)	0.39 (0.34, 0.44)
Space	0.73 (0.69, 0.78)	0.49 (0.44, 0.54)	0.46 (0.41, 0.51)	0.42 (0.37, 0.48)	0.40 (0.34, 0.44)

Supplementary Table S4: Text-based functional form identification ablation results: modifying value separator for different noise levels.

Separator	cubic	exponential	linear	periodic	quadratic
Comma and space	0.12 (0.08, 0.14)	0.91 (0.88, 0.94)	0.48 (0.43, 0.53)	0.73 (0.69, 0.78)	0.26 (0.21, 0.30)
Space	0.10 (0.07, 0.12)	0.83 (0.80, 0.87)	0.45 (0.40, 0.50)	0.87 (0.84, 0.91)	0.26 (0.22, 0.30)

Supplementary Table S5: Text-based functional form identification ablation results: modifying value separator for different function classes.

Fixed precision	50	500	1000	2500
2	0.70 (0.64, 0.75)	0.53 (0.47, 0.59)	0.51 (0.45, 0.57)	0.31 (0.25, 0.37)
4	0.62 (0.55, 0.68)	0.53 (0.47, 0.59)	0.50 (0.44, 0.56)	0.40 (0.34, 0.46)
8	0.60 (0.54, 0.66)	0.51 (0.45, 0.57)	0.52 (0.46, 0.58)	0.36 (0.30, 0.42)
16	0.59 (0.53, 0.64)	0.52 (0.46, 0.59)	0.44 (0.38, 0.51)	0.36 (0.31, 0.43)

Supplementary Table S6: Text-based functional form identification ablation results: modifying floating point fixed precision for different numbers of points.

1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631

Fixed precision	0.0	0.5	1.0	2.0	5.0
2	0.77 (0.71, 0.82)	0.51 (0.43, 0.57)	0.45 (0.38, 0.52)	0.43 (0.36, 0.51)	0.40 (0.33, 0.47)
4	0.75 (0.68, 0.81)	0.46 (0.40, 0.53)	0.46 (0.40, 0.53)	0.46 (0.39, 0.53)	0.43 (0.36, 0.50)
8	0.69 (0.62, 0.75)	0.45 (0.38, 0.52)	0.48 (0.41, 0.55)	0.47 (0.40, 0.53)	0.41 (0.34, 0.48)
16	0.68 (0.61, 0.74)	0.56 (0.49, 0.63)	0.42 (0.35, 0.49)	0.41 (0.34, 0.47)	0.33 (0.27, 0.40)

Supplementary Table S7: Text-based functional form identification ablation results: modifying floating point fixed precision for different noise levels.

Fixed precision	cubic	exponential	linear	periodic	quadratic
2	0.12 (0.07, 0.16)	0.95 (0.92, 0.98)	0.58 (0.52, 0.65)	0.55 (0.48, 0.62)	0.35 (0.29, 0.42)
4	0.10 (0.07, 0.15)	0.95 (0.93, 0.98)	0.41 (0.35, 0.48)	0.81 (0.76, 0.86)	0.27 (0.20, 0.32)
8	0.07 (0.04, 0.12)	0.84 (0.79, 0.89)	0.45 (0.38, 0.51)	0.89 (0.84, 0.93)	0.24 (0.18, 0.30)
16	0.12 (0.08, 0.17)	0.74 (0.68, 0.80)	0.40 (0.33, 0.47)	0.96 (0.94, 0.98)	0.17 (0.12, 0.23)

Supplementary Table S8: Text-based functional form identification ablation results: modifying floating point fixed precision for different function classes.

1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664

rescaling	50	500	1000	2500
llmtime ($\alpha = 0.5, \beta = 0.0$)	0.61 (0.52, 0.69)	0.31 (0.23, 0.40)	0.32 (0.24, 0.40)	0.40 (0.32, 0.48)
llmtime ($\alpha = 0.5, \beta = 0.15$)	0.55 (0.46, 0.65)	0.38 (0.30, 0.46)	0.35 (0.27, 0.43)	0.30 (0.22, 0.38)
llmtime ($\alpha = 0.5, \beta = 0.3$)	0.50 (0.42, 0.59)	0.37 (0.28, 0.46)	0.33 (0.26, 0.41)	0.27 (0.19, 0.35)
llmtime ($\alpha = 0.5, \beta = 0.5$)	0.57 (0.48, 0.66)	0.38 (0.29, 0.46)	0.34 (0.26, 0.43)	0.30 (0.22, 0.38)
llmtime ($\alpha = 0.7, \beta = 0.0$)	0.57 (0.48, 0.66)	0.32 (0.24, 0.41)	0.29 (0.21, 0.38)	0.38 (0.30, 0.47)
llmtime ($\alpha = 0.7, \beta = 0.15$)	0.54 (0.46, 0.62)	0.35 (0.27, 0.44)	0.30 (0.22, 0.38)	0.30 (0.22, 0.38)
llmtime ($\alpha = 0.7, \beta = 0.3$)	0.57 (0.48, 0.66)	0.34 (0.26, 0.43)	0.30 (0.22, 0.38)	0.30 (0.22, 0.38)
llmtime ($\alpha = 0.7, \beta = 0.5$)	0.54 (0.45, 0.62)	0.40 (0.32, 0.49)	0.32 (0.24, 0.40)	0.28 (0.21, 0.36)
llmtime ($\alpha = 0.9, \beta = 0.0$)	0.52 (0.43, 0.61)	0.33 (0.24, 0.41)	0.29 (0.20, 0.37)	0.38 (0.30, 0.47)
llmtime ($\alpha = 0.9, \beta = 0.15$)	0.52 (0.43, 0.61)	0.38 (0.30, 0.46)	0.32 (0.25, 0.40)	0.27 (0.19, 0.35)
llmtime ($\alpha = 0.9, \beta = 0.3$)	0.50 (0.41, 0.58)	0.31 (0.23, 0.39)	0.28 (0.21, 0.36)	0.30 (0.23, 0.38)
llmtime ($\alpha = 0.9, \beta = 0.5$)	0.53 (0.44, 0.62)	0.32 (0.24, 0.39)	0.32 (0.24, 0.40)	0.31 (0.23, 0.39)
llmtime ($\alpha = 0.99, \beta = 0.0$)	0.56 (0.47, 0.65)	0.31 (0.23, 0.40)	0.31 (0.24, 0.39)	0.41 (0.33, 0.50)
llmtime ($\alpha = 0.99, \beta = 0.15$)	0.54 (0.46, 0.62)	0.34 (0.25, 0.42)	0.38 (0.29, 0.46)	0.32 (0.24, 0.40)
llmtime ($\alpha = 0.99, \beta = 0.3$)	0.54 (0.46, 0.63)	0.33 (0.24, 0.41)	0.32 (0.24, 0.40)	0.29 (0.22, 0.38)
llmtime ($\alpha = 0.99, \beta = 0.5$)	0.51 (0.43, 0.60)	0.34 (0.26, 0.42)	0.33 (0.25, 0.41)	0.30 (0.22, 0.37)
minmax	0.61 (0.53, 0.70)	0.42 (0.34, 0.51)	0.35 (0.27, 0.43)	0.24 (0.17, 0.32)
none	0.62 (0.54, 0.70)	0.58 (0.49, 0.66)	0.57 (0.47, 0.66)	0.35 (0.27, 0.43)

Supplementary Table S9: Text-based functional form identification ablation results: rescaling values as suggested in (Gruver et al., 2024) for different numbers of points.

1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697

rescaling	0.0	0.5	1.0	2.0	5.0
lmltime ($\alpha = 0.5, \beta = 0.0$)	0.41 (0.31, 0.50)	0.43 (0.33, 0.53)	0.43 (0.33, 0.52)	0.39 (0.29, 0.48)	0.39 (0.30, 0.49)
lmltime ($\alpha = 0.5, \beta = 0.15$)	0.48 (0.38, 0.58)	0.42 (0.32, 0.51)	0.40 (0.31, 0.50)	0.36 (0.26, 0.45)	0.31 (0.22, 0.40)
lmltime ($\alpha = 0.5, \beta = 0.3$)	0.48 (0.38, 0.57)	0.42 (0.33, 0.52)	0.42 (0.32, 0.52)	0.28 (0.19, 0.37)	0.24 (0.16, 0.33)
lmltime ($\alpha = 0.5, \beta = 0.5$)	0.49 (0.39, 0.59)	0.44 (0.34, 0.54)	0.42 (0.33, 0.51)	0.36 (0.27, 0.46)	0.27 (0.18, 0.36)
lmltime ($\alpha = 0.7, \beta = 0.0$)	0.33 (0.24, 0.42)	0.44 (0.34, 0.54)	0.43 (0.33, 0.53)	0.35 (0.26, 0.44)	0.40 (0.31, 0.49)
lmltime ($\alpha = 0.7, \beta = 0.15$)	0.45 (0.35, 0.54)	0.47 (0.37, 0.56)	0.34 (0.25, 0.44)	0.34 (0.25, 0.43)	0.26 (0.18, 0.34)
lmltime ($\alpha = 0.7, \beta = 0.3$)	0.43 (0.34, 0.53)	0.39 (0.30, 0.48)	0.37 (0.27, 0.47)	0.38 (0.29, 0.47)	0.31 (0.22, 0.40)
lmltime ($\alpha = 0.7, \beta = 0.5$)	0.49 (0.38, 0.58)	0.47 (0.36, 0.57)	0.41 (0.31, 0.50)	0.30 (0.20, 0.39)	0.25 (0.17, 0.33)
lmltime ($\alpha = 0.9, \beta = 0.0$)	0.39 (0.30, 0.49)	0.41 (0.31, 0.50)	0.42 (0.32, 0.52)	0.38 (0.29, 0.48)	0.30 (0.21, 0.39)
lmltime ($\alpha = 0.9, \beta = 0.15$)	0.45 (0.35, 0.55)	0.40 (0.32, 0.49)	0.39 (0.30, 0.49)	0.35 (0.26, 0.44)	0.27 (0.19, 0.36)
lmltime ($\alpha = 0.9, \beta = 0.3$)	0.39 (0.29, 0.49)	0.41 (0.31, 0.51)	0.35 (0.26, 0.44)	0.34 (0.25, 0.43)	0.25 (0.17, 0.33)
lmltime ($\alpha = 0.9, \beta = 0.5$)	0.44 (0.34, 0.54)	0.45 (0.35, 0.54)	0.44 (0.34, 0.53)	0.33 (0.24, 0.43)	0.19 (0.12, 0.27)
lmltime ($\alpha = 0.99, \beta = 0.0$)	0.41 (0.32, 0.50)	0.43 (0.33, 0.53)	0.38 (0.29, 0.48)	0.38 (0.29, 0.47)	0.39 (0.29, 0.48)
lmltime ($\alpha = 0.99, \beta = 0.15$)	0.50 (0.40, 0.60)	0.43 (0.34, 0.53)	0.44 (0.35, 0.53)	0.30 (0.21, 0.40)	0.30 (0.21, 0.40)
lmltime ($\alpha = 0.99, \beta = 0.3$)	0.43 (0.34, 0.54)	0.44 (0.35, 0.53)	0.40 (0.31, 0.49)	0.27 (0.19, 0.36)	0.31 (0.22, 0.40)
lmltime ($\alpha = 0.99, \beta = 0.5$)	0.47 (0.38, 0.57)	0.45 (0.35, 0.55)	0.40 (0.31, 0.50)	0.27 (0.18, 0.36)	0.25 (0.17, 0.34)
minmax	0.39 (0.30, 0.49)	0.53 (0.43, 0.62)	0.48 (0.38, 0.58)	0.35 (0.26, 0.44)	0.28 (0.20, 0.36)
none	0.77 (0.68, 0.85)	0.54 (0.44, 0.64)	0.55 (0.45, 0.64)	0.45 (0.35, 0.55)	0.34 (0.25, 0.44)

Supplementary Table S10: Text-based functional form identification ablation results: rescaling values as suggested in (Gruver et al., 2024) for different noise levels.

1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730

rescaling	cubic	exponential	linear	periodic	quadratic
llmtime ($\alpha=0.5, \beta=0.0$)	0.01 (0.00, 0.03)	0.68 (0.58, 0.76)	0.50 (0.40, 0.60)	0.62 (0.53, 0.72)	0.24 (0.15, 0.33)
llmtime ($\alpha=0.5, \beta=0.15$)	0.01 (0.00, 0.03)	0.58 (0.49, 0.68)	0.55 (0.45, 0.64)	0.58 (0.49, 0.67)	0.25 (0.17, 0.34)
llmtime ($\alpha=0.5, \beta=0.3$)	0.01 (0.00, 0.03)	0.59 (0.50, 0.69)	0.42 (0.33, 0.51)	0.57 (0.47, 0.67)	0.25 (0.17, 0.33)
llmtime ($\alpha=0.5, \beta=0.5$)	0.01 (0.00, 0.03)	0.56 (0.46, 0.65)	0.52 (0.42, 0.62)	0.60 (0.50, 0.69)	0.29 (0.21, 0.38)
llmtime ($\alpha=0.7, \beta=0.0$)	0.02 (0.00, 0.05)	0.71 (0.62, 0.79)	0.44 (0.34, 0.54)	0.62 (0.52, 0.72)	0.16 (0.09, 0.24)
llmtime ($\alpha=0.7, \beta=0.15$)	0.00 (0.00, 0.00)	0.54 (0.44, 0.64)	0.47 (0.38, 0.56)	0.61 (0.51, 0.70)	0.24 (0.16, 0.32)
llmtime ($\alpha=0.7, \beta=0.3$)	0.00 (0.00, 0.00)	0.62 (0.52, 0.72)	0.46 (0.36, 0.56)	0.57 (0.47, 0.67)	0.23 (0.15, 0.31)
llmtime ($\alpha=0.7, \beta=0.5$)	0.03 (0.00, 0.07)	0.57 (0.47, 0.67)	0.49 (0.40, 0.59)	0.59 (0.49, 0.69)	0.24 (0.16, 0.33)
llmtime ($\alpha=0.9, \beta=0.0$)	0.01 (0.00, 0.03)	0.68 (0.59, 0.77)	0.36 (0.27, 0.45)	0.61 (0.52, 0.70)	0.24 (0.16, 0.32)
llmtime ($\alpha=0.9, \beta=0.15$)	0.00 (0.00, 0.00)	0.61 (0.52, 0.70)	0.48 (0.39, 0.58)	0.57 (0.47, 0.67)	0.20 (0.12, 0.28)
llmtime ($\alpha=0.9, \beta=0.3$)	0.02 (0.00, 0.05)	0.48 (0.38, 0.57)	0.40 (0.31, 0.50)	0.65 (0.55, 0.74)	0.19 (0.12, 0.27)
llmtime ($\alpha=0.9, \beta=0.5$)	0.03 (0.00, 0.07)	0.53 (0.43, 0.63)	0.53 (0.44, 0.63)	0.53 (0.43, 0.63)	0.23 (0.15, 0.32)
llmtime ($\alpha=0.99, \beta=0.0$)	0.01 (0.00, 0.03)	0.71 (0.62, 0.80)	0.40 (0.30, 0.50)	0.62 (0.53, 0.71)	0.25 (0.17, 0.34)
llmtime ($\alpha=0.99, \beta=0.15$)	0.00 (0.00, 0.00)	0.59 (0.48, 0.69)	0.48 (0.38, 0.58)	0.63 (0.54, 0.72)	0.27 (0.18, 0.36)
llmtime ($\alpha=0.99, \beta=0.3$)	0.01 (0.00, 0.03)	0.61 (0.51, 0.71)	0.45 (0.35, 0.55)	0.58 (0.48, 0.68)	0.20 (0.12, 0.28)
llmtime ($\alpha=0.99, \beta=0.5$)	0.01 (0.00, 0.03)	0.51 (0.42, 0.60)	0.47 (0.37, 0.57)	0.61 (0.51, 0.70)	0.24 (0.16, 0.33)
minmax	0.00 (0.00, 0.00)	0.57 (0.47, 0.66)	0.59 (0.49, 0.67)	0.67 (0.58, 0.76)	0.20 (0.13, 0.28)
none	0.17 (0.10, 0.25)	0.92 (0.87, 0.97)	0.48 (0.38, 0.58)	0.79 (0.71, 0.87)	0.29 (0.20, 0.37)

Supplementary Table S11: Text-based functional form identification ablation results: rescaling values as suggested in (Gruver et al., 2024) for different function classes.

1731 A.3.3 PLOT-BASED RESULTS
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784

1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817

dpi	50	500	1000	2500
25	0.66 (0.57, 0.74)	0.67 (0.60, 0.75)	0.70 (0.62, 0.78)	0.66 (0.58, 0.74)
50	0.59 (0.51, 0.67)	0.70 (0.62, 0.78)	0.70 (0.62, 0.78)	0.74 (0.66, 0.82)
100	0.67 (0.58, 0.75)	0.70 (0.62, 0.78)	0.68 (0.60, 0.76)	0.74 (0.66, 0.81)
200	0.72 (0.64, 0.80)	0.70 (0.62, 0.78)	0.72 (0.64, 0.80)	0.73 (0.65, 0.81)
400	0.63 (0.55, 0.71)	0.71 (0.63, 0.79)	0.71 (0.63, 0.79)	0.73 (0.65, 0.81)

Supplementary Table S12: Plot-based functional form identification ablation results: modifying figure dpi for different numbers of points.

dpi	0.0	0.5	1.0	2.0	5.0
25	0.91 (0.85, 0.96)	0.83 (0.76, 0.90)	0.61 (0.52, 0.71)	0.61 (0.51, 0.70)	0.41 (0.30, 0.51)
50	0.92 (0.86, 0.97)	0.87 (0.80, 0.93)	0.68 (0.59, 0.77)	0.59 (0.49, 0.69)	0.35 (0.26, 0.44)
100	0.95 (0.90, 0.99)	0.89 (0.82, 0.95)	0.65 (0.56, 0.74)	0.60 (0.50, 0.70)	0.40 (0.31, 0.50)
200	0.98 (0.95, 1.00)	0.89 (0.83, 0.95)	0.69 (0.60, 0.78)	0.61 (0.52, 0.70)	0.42 (0.33, 0.52)
400	0.94 (0.89, 0.98)	0.89 (0.82, 0.95)	0.68 (0.59, 0.77)	0.58 (0.49, 0.68)	0.39 (0.30, 0.49)

Supplementary Table S13: Plot-based functional form identification ablation results: modifying figure dpi for different noise levels.

dpi	cubic	exponential	linear	periodic	quadratic
25	0.22 (0.14, 0.30)	0.95 (0.91, 0.99)	0.99 (0.97, 1.00)	0.37 (0.28, 0.47)	0.84 (0.77, 0.91)
50	0.38 (0.29, 0.47)	0.95 (0.91, 0.99)	0.98 (0.95, 1.00)	0.32 (0.23, 0.41)	0.78 (0.70, 0.85)
100	0.37 (0.28, 0.47)	0.96 (0.92, 0.99)	0.97 (0.93, 1.00)	0.45 (0.35, 0.55)	0.74 (0.65, 0.82)
200	0.42 (0.32, 0.51)	0.97 (0.93, 1.00)	0.98 (0.95, 1.00)	0.45 (0.36, 0.55)	0.77 (0.68, 0.84)
400	0.43 (0.33, 0.53)	0.94 (0.89, 0.98)	0.98 (0.95, 1.00)	0.39 (0.29, 0.49)	0.74 (0.65, 0.83)

Supplementary Table S14: Plot-based functional form identification ablation results: modifying figure dpi for different function classes.

1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850

figsize	50	500	1000	2500
(4, 3)	0.67 (0.59, 0.75)	0.70 (0.62, 0.78)	0.71 (0.63, 0.79)	0.71 (0.63, 0.79)
(3.5, 3.5)	0.71 (0.63, 0.79)	0.73 (0.65, 0.81)	0.76 (0.68, 0.83)	0.77 (0.70, 0.83)
(6.4, 4.8)	0.69 (0.60, 0.76)	0.72 (0.65, 0.80)	0.70 (0.63, 0.78)	0.72 (0.65, 0.80)
(8, 6)	0.65 (0.57, 0.74)	0.70 (0.62, 0.78)	0.69 (0.60, 0.77)	0.72 (0.65, 0.79)
(7, 7)	0.63 (0.55, 0.72)	0.74 (0.66, 0.81)	0.71 (0.63, 0.79)	0.75 (0.67, 0.82)
(12, 12)	0.70 (0.62, 0.78)	0.66 (0.58, 0.74)	0.68 (0.60, 0.75)	0.71 (0.63, 0.79)

Supplementary Table S15: Plot-based functional form identification ablation results: modifying figure size for different numbers of points.

figsize	0.0	0.5	1.0	2.0	5.0
(4, 3)	0.98 (0.95, 1.00)	0.86 (0.79, 0.93)	0.65 (0.55, 0.75)	0.63 (0.54, 0.72)	0.37 (0.28, 0.47)
(3.5, 3.5)	0.98 (0.95, 1.00)	0.90 (0.83, 0.95)	0.76 (0.68, 0.84)	0.69 (0.60, 0.77)	0.38 (0.29, 0.48)
(6.4, 4.8)	0.97 (0.93, 1.00)	0.89 (0.83, 0.95)	0.66 (0.57, 0.75)	0.59 (0.49, 0.69)	0.43 (0.34, 0.53)
(8, 6)	0.95 (0.90, 0.99)	0.88 (0.82, 0.94)	0.67 (0.59, 0.75)	0.59 (0.49, 0.68)	0.35 (0.26, 0.44)
(7, 7)	0.95 (0.90, 0.99)	0.88 (0.81, 0.94)	0.76 (0.68, 0.84)	0.60 (0.51, 0.70)	0.35 (0.26, 0.44)
(12, 12)	0.94 (0.89, 0.98)	0.79 (0.70, 0.87)	0.69 (0.60, 0.78)	0.63 (0.54, 0.72)	0.38 (0.29, 0.48)

Supplementary Table S16: Plot-based functional form identification ablation results: modifying figure size for different noise levels.

figsize	cubic	exponential	linear	periodic	quadratic
(4, 3)	0.35 (0.26, 0.44)	0.95 (0.90, 0.99)	0.99 (0.97, 1.00)	0.39 (0.29, 0.48)	0.81 (0.73, 0.89)
(3.5, 3.5)	0.57 (0.47, 0.66)	0.96 (0.92, 0.99)	0.99 (0.97, 1.00)	0.36 (0.27, 0.46)	0.83 (0.76, 0.90)
(6.4, 4.8)	0.37 (0.28, 0.46)	0.96 (0.92, 0.99)	0.99 (0.96, 1.00)	0.48 (0.38, 0.57)	0.74 (0.65, 0.83)
(8, 6)	0.40 (0.30, 0.50)	0.94 (0.89, 0.98)	0.97 (0.93, 1.00)	0.39 (0.30, 0.49)	0.74 (0.66, 0.83)
(7, 7)	0.51 (0.41, 0.61)	0.93 (0.88, 0.98)	0.95 (0.90, 0.99)	0.35 (0.26, 0.44)	0.80 (0.72, 0.88)
(12, 12)	0.31 (0.22, 0.40)	0.91 (0.85, 0.96)	1.00 (1.00, 1.00)	0.43 (0.33, 0.53)	0.78 (0.70, 0.85)

Supplementary Table S17: Plot-based functional form identification ablation results: modifying figure size for different function classes.

1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883

Plot style	50	500	1000	2500
None	0.66 (0.57, 0.74)	0.72 (0.64, 0.79)	0.70 (0.62, 0.78)	0.72 (0.64, 0.79)
classic	0.70 (0.62, 0.78)	0.72 (0.64, 0.79)	0.70 (0.62, 0.78)	0.73 (0.65, 0.80)
ggplot	0.65 (0.56, 0.73)	0.70 (0.62, 0.79)	0.70 (0.62, 0.78)	0.74 (0.67, 0.82)
seaborn-v0.8-darkgrid	0.68 (0.60, 0.76)	0.71 (0.63, 0.78)	0.69 (0.60, 0.77)	0.74 (0.66, 0.82)
seaborn-v0.8-whitegrid	0.66 (0.58, 0.74)	0.70 (0.62, 0.78)	0.70 (0.62, 0.78)	0.72 (0.63, 0.80)

Supplementary Table S18: Plot-based functional form identification ablation results: modifying plotting style for different numbers of points.

Plot style	0.0	0.5	1.0	2.0	5.0
None	0.95 (0.91, 0.99)	0.92 (0.86, 0.97)	0.65 (0.55, 0.74)	0.56 (0.47, 0.65)	0.41 (0.31, 0.50)
classic	0.96 (0.92, 0.99)	0.91 (0.85, 0.96)	0.67 (0.57, 0.76)	0.59 (0.50, 0.68)	0.43 (0.33, 0.52)
ggplot	0.94 (0.89, 0.98)	0.91 (0.85, 0.96)	0.68 (0.59, 0.77)	0.60 (0.50, 0.70)	0.37 (0.28, 0.46)
seaborn-v0.8-darkgrid	0.96 (0.92, 0.99)	0.87 (0.80, 0.93)	0.67 (0.57, 0.76)	0.61 (0.51, 0.70)	0.41 (0.32, 0.51)
seaborn-v0.8-whitegrid	0.96 (0.92, 0.99)	0.87 (0.81, 0.94)	0.67 (0.58, 0.75)	0.58 (0.48, 0.68)	0.41 (0.32, 0.51)

Supplementary Table S19: Plot-based functional form identification ablation results: modifying plotting style for different noise levels.

Plot style	cubic	exponential	linear	periodic	quadratic
None	0.39 (0.29, 0.48)	0.97 (0.93, 1.00)	0.97 (0.93, 1.00)	0.43 (0.34, 0.53)	0.73 (0.64, 0.82)
classic	0.41 (0.32, 0.51)	0.96 (0.92, 0.99)	0.98 (0.95, 1.00)	0.46 (0.37, 0.56)	0.75 (0.67, 0.83)
ggplot	0.42 (0.32, 0.51)	0.94 (0.89, 0.98)	0.96 (0.92, 0.99)	0.43 (0.34, 0.53)	0.75 (0.67, 0.84)
seaborn-v0.8-darkgrid	0.36 (0.26, 0.45)	0.95 (0.91, 0.99)	1.00 (1.00, 1.00)	0.46 (0.36, 0.55)	0.75 (0.65, 0.84)
seaborn-v0.8-whitegrid	0.35 (0.26, 0.45)	0.97 (0.93, 1.00)	0.97 (0.93, 1.00)	0.47 (0.38, 0.57)	0.73 (0.64, 0.82)

Supplementary Table S20: Plot-based functional form identification ablation results: modifying plotting styles for different function classes.

Color palette	50	500	2500
black and white	0.50 (0.48, 0.51)	0.64 (0.63, 0.65)	0.67 (0.66, 0.69)
default	0.50 (0.49, 0.52)	0.63 (0.61, 0.64)	0.65 (0.64, 0.67)
high contrast	0.49 (0.48, 0.51)	0.60 (0.58, 0.61)	0.63 (0.61, 0.64)
invert	0.51 (0.50, 0.53)	0.64 (0.62, 0.65)	0.66 (0.65, 0.68)
low contrast	0.48 (0.47, 0.50)	0.61 (0.60, 0.63)	0.65 (0.63, 0.66)

Supplementary Table S21: Plot-based functional form identification ablation results: modifying color palette for different numbers of points.

1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916

Color palette	0.0	1.0	5.0
black and white	0.88 (0.87, 0.89)	0.63 (0.62, 0.65)	0.30 (0.29, 0.32)
default	0.88 (0.87, 0.89)	0.63 (0.62, 0.65)	0.27 (0.26, 0.29)
high contrast	0.85 (0.84, 0.86)	0.62 (0.60, 0.63)	0.25 (0.24, 0.26)
invert	0.88 (0.87, 0.89)	0.64 (0.62, 0.65)	0.30 (0.28, 0.31)
low contrast	0.86 (0.85, 0.87)	0.61 (0.60, 0.63)	0.27 (0.25, 0.28)

Supplementary Table S22: Plot-based functional form identification ablation results: modifying color palette for different noise levels.

Color palette	cubic	exponential	linear	periodic	quadratic
black and white	0.35 (0.33, 0.37)	0.78 (0.76, 0.79)	0.90 (0.88, 0.91)	0.30 (0.28, 0.32)	0.71 (0.69, 0.72)
default	0.32 (0.30, 0.34)	0.76 (0.74, 0.77)	0.90 (0.89, 0.91)	0.31 (0.29, 0.33)	0.69 (0.67, 0.71)
high contrast	0.27 (0.25, 0.28)	0.73 (0.71, 0.74)	0.90 (0.89, 0.91)	0.28 (0.26, 0.30)	0.69 (0.68, 0.71)
invert	0.32 (0.30, 0.34)	0.79 (0.77, 0.80)	0.90 (0.89, 0.91)	0.31 (0.29, 0.33)	0.70 (0.69, 0.72)
low contrast	0.29 (0.28, 0.31)	0.75 (0.74, 0.77)	0.88 (0.87, 0.89)	0.29 (0.27, 0.30)	0.69 (0.67, 0.71)

Supplementary Table S23: Plot-based functional form identification ablation results: modifying color palette for different function classes.

Plot marker	50	500	2500
+	0.48 (0.47, 0.49)	0.61 (0.60, 0.63)	0.64 (0.63, 0.66)
^	0.50 (0.49, 0.52)	0.64 (0.63, 0.66)	0.67 (0.66, 0.69)
o	0.52 (0.51, 0.54)	0.63 (0.62, 0.65)	0.65 (0.64, 0.67)
s	0.50 (0.48, 0.52)	0.62 (0.60, 0.64)	0.64 (0.63, 0.66)
x	0.49 (0.47, 0.50)	0.61 (0.59, 0.62)	0.65 (0.64, 0.67)

Supplementary Table S24: Plot-based functional form identification ablation results: modifying plot markers for different numbers of points.

Plot marker	0.0	1.0	5.0
+	0.86 (0.85, 0.87)	0.62 (0.60, 0.63)	0.26 (0.25, 0.27)
^	0.88 (0.87, 0.89)	0.64 (0.63, 0.66)	0.30 (0.29, 0.32)
o	0.88 (0.87, 0.89)	0.63 (0.62, 0.65)	0.29 (0.28, 0.31)
s	0.88 (0.87, 0.89)	0.62 (0.60, 0.63)	0.27 (0.26, 0.28)
x	0.86 (0.85, 0.87)	0.62 (0.61, 0.64)	0.27 (0.25, 0.28)

Supplementary Table S25: Plot-based functional form identification ablation results: modifying plot markers for different noise levels.

1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949

Plot marker	cubic	exponential	linear	periodic	quadratic
+	0.29 (0.27, 0.31)	0.75 (0.74, 0.77)	0.89 (0.88, 0.90)	0.28 (0.26, 0.30)	0.68 (0.66, 0.69)
^	0.31 (0.29, 0.33)	0.77 (0.75, 0.78)	0.89 (0.87, 0.90)	0.34 (0.33, 0.36)	0.73 (0.71, 0.74)
o	0.32 (0.30, 0.34)	0.78 (0.77, 0.80)	0.90 (0.89, 0.91)	0.30 (0.28, 0.32)	0.71 (0.69, 0.73)
s	0.33 (0.31, 0.35)	0.75 (0.74, 0.77)	0.91 (0.90, 0.92)	0.28 (0.26, 0.29)	0.68 (0.66, 0.70)
x	0.30 (0.29, 0.32)	0.75 (0.73, 0.76)	0.89 (0.87, 0.90)	0.28 (0.26, 0.30)	0.70 (0.68, 0.72)

Supplementary Table S26: Plot-based functional form identification ablation results: modifying plot markers for different function classes.

Markers size	50	500	2500
large	0.51 (0.49, 0.52)	0.62 (0.61, 0.63)	0.65 (0.64, 0.66)
medium	0.50 (0.49, 0.51)	0.62 (0.61, 0.63)	0.65 (0.64, 0.66)
small	0.49 (0.48, 0.50)	0.63 (0.62, 0.64)	0.65 (0.64, 0.66)

Supplementary Table S27: Plot-based functional form identification ablation results: modifying plot maker sizes for different numbers of points.

Marker size	0.0	1.0	5.0
large	0.87 (0.86, 0.88)	0.63 (0.62, 0.64)	0.27 (0.26, 0.28)
medium	0.87 (0.86, 0.88)	0.63 (0.61, 0.64)	0.28 (0.27, 0.29)
small	0.87 (0.86, 0.87)	0.62 (0.61, 0.63)	0.28 (0.27, 0.29)

Supplementary Table S28: Plot-based functional form identification ablation results: modifying plot maker sizes for different noise levels.

Marker size	cubic	exponential	linear	periodic	quadratic
large	0.32 (0.30, 0.33)	0.75 (0.73, 0.76)	0.91 (0.90, 0.92)	0.30 (0.28, 0.31)	0.70 (0.68, 0.71)
medium	0.31 (0.30, 0.33)	0.76 (0.75, 0.77)	0.90 (0.89, 0.91)	0.29 (0.27, 0.30)	0.70 (0.69, 0.72)
small	0.30 (0.29, 0.32)	0.78 (0.76, 0.79)	0.88 (0.87, 0.89)	0.30 (0.29, 0.32)	0.69 (0.68, 0.71)

Supplementary Table S29: Plot-based functional form identification ablation results: modifying plot maker sizes for different function classes.

Plot components	50	500	2500
all	0.53 (0.51, 0.54)	0.65 (0.64, 0.66)	0.68 (0.67, 0.69)
minimal	0.47 (0.46, 0.48)	0.61 (0.60, 0.62)	0.63 (0.62, 0.64)
none	0.50 (0.49, 0.51)	0.61 (0.60, 0.62)	0.66 (0.64, 0.67)

Supplementary Table S30: Plot-based functional form identification ablation results: modifying plot components for different numbers of points.

1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982

Plot components	0.0	1.0	5.0
all	0.89 (0.88, 0.90)	0.65 (0.63, 0.66)	0.32 (0.31, 0.33)
minimal	0.83 (0.82, 0.84)	0.60 (0.59, 0.61)	0.28 (0.27, 0.29)
none	0.89 (0.88, 0.90)	0.63 (0.62, 0.65)	0.24 (0.23, 0.25)

Supplementary Table S31: Plot-based functional form identification ablation results: modifying plot components for different noise levels.

Plot components	cubic	exponential	linear	periodic	quadratic
all	0.37 (0.36, 0.39)	0.82 (0.81, 0.83)	0.94 (0.94, 0.95)	0.29 (0.28, 0.31)	0.66 (0.65, 0.68)
minimal	0.24 (0.23, 0.25)	0.75 (0.74, 0.76)	0.89 (0.88, 0.90)	0.25 (0.24, 0.27)	0.71 (0.69, 0.72)
none	0.32 (0.30, 0.33)	0.71 (0.70, 0.72)	0.85 (0.84, 0.86)	0.34 (0.33, 0.35)	0.72 (0.71, 0.74)

Supplementary Table S32: Plot-based functional form identification ablation results: modifying plot components for different function classes.

Temperature	50	500	1000	2500
0.00	0.70 (0.62, 0.77)	0.71 (0.63, 0.79)	0.68 (0.60, 0.75)	0.72 (0.64, 0.80)
0.10	0.66 (0.58, 0.74)	0.71 (0.63, 0.78)	0.70 (0.62, 0.78)	0.73 (0.65, 0.80)
0.30	0.72 (0.64, 0.80)	0.72 (0.64, 0.80)	0.68 (0.60, 0.77)	0.73 (0.65, 0.80)
0.55	0.69 (0.60, 0.78)	0.70 (0.62, 0.78)	0.70 (0.62, 0.77)	0.74 (0.66, 0.82)
1.00	0.71 (0.64, 0.78)	0.70 (0.62, 0.78)	0.68 (0.59, 0.76)	0.74 (0.66, 0.81)

Supplementary Table S33: Plot-based functional form identification ablation results: modifying temperature for different numbers of points.

Temperature	0.0	0.5	1.0	2.0	5.0
0.00	0.96 (0.92, 0.99)	0.90 (0.84, 0.96)	0.67 (0.58, 0.75)	0.58 (0.47, 0.67)	0.40 (0.31, 0.50)
0.10	0.96 (0.92, 0.99)	0.90 (0.84, 0.96)	0.67 (0.58, 0.76)	0.57 (0.48, 0.66)	0.40 (0.31, 0.51)
0.30	0.96 (0.91, 0.99)	0.93 (0.87, 0.98)	0.67 (0.58, 0.76)	0.62 (0.53, 0.71)	0.38 (0.29, 0.48)
0.55	0.95 (0.90, 0.99)	0.90 (0.84, 0.95)	0.69 (0.60, 0.78)	0.59 (0.49, 0.69)	0.41 (0.32, 0.51)
1.00	0.96 (0.92, 0.99)	0.91 (0.85, 0.96)	0.68 (0.58, 0.77)	0.61 (0.52, 0.70)	0.38 (0.28, 0.47)

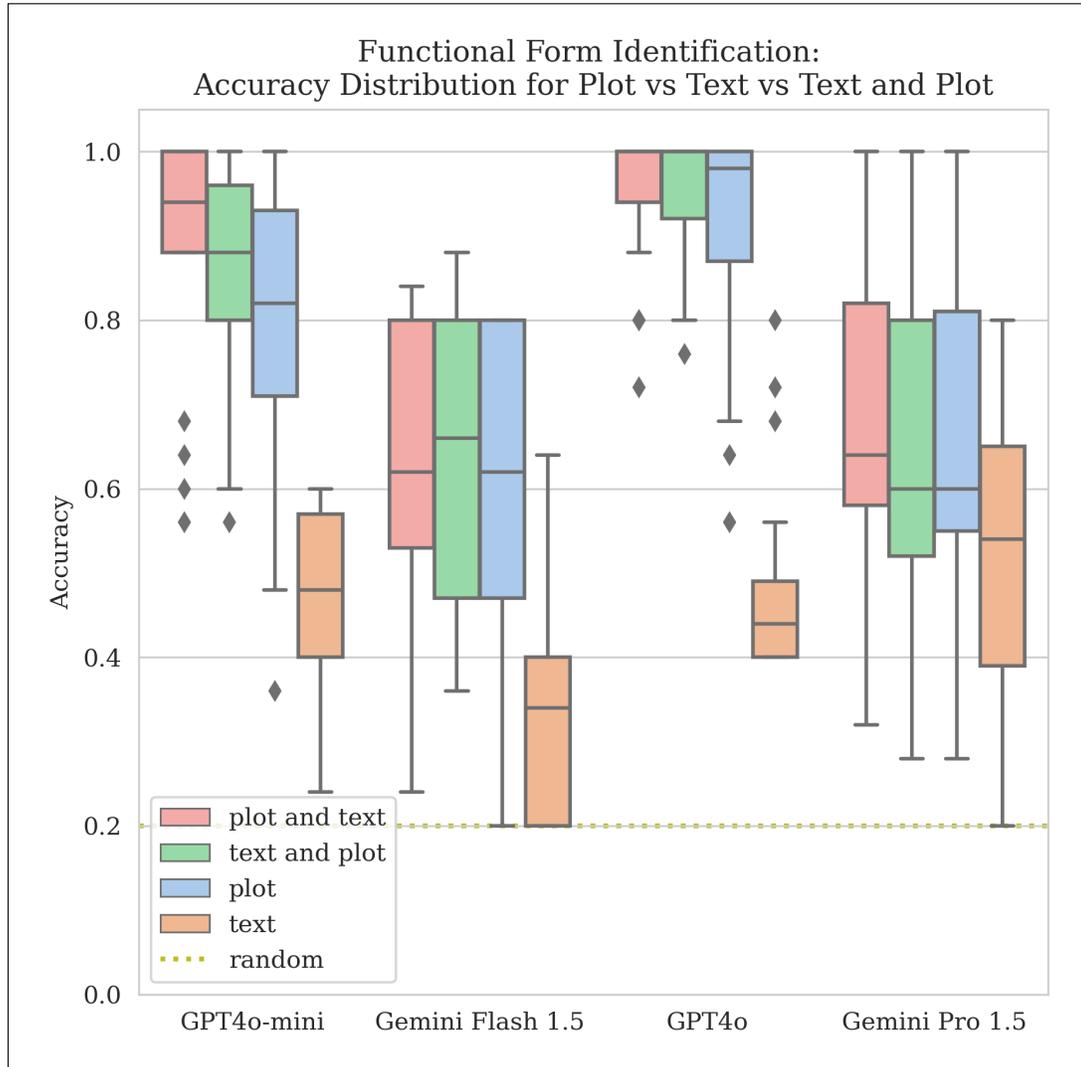
Supplementary Table S34: Plot-based functional form identification ablation results: modifying temperature for different noise levels.

1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015

Temperature	cubic	exponential	linear	periodic	quadratic
0.0	0.37 (0.27, 0.46)	0.95 (0.91, 0.99)	0.97 (0.93, 1.00)	0.48 (0.38, 0.58)	0.74 (0.65, 0.82)
0.1	0.40 (0.30, 0.50)	0.94 (0.89, 0.98)	0.97 (0.93, 1.00)	0.45 (0.35, 0.54)	0.74 (0.66, 0.83)
0.3	0.42 (0.33, 0.52)	0.92 (0.86, 0.97)	0.99 (0.97, 1.00)	0.47 (0.38, 0.57)	0.76 (0.67, 0.84)
0.55	0.39 (0.30, 0.49)	0.96 (0.92, 0.99)	0.99 (0.97, 1.00)	0.46 (0.37, 0.56)	0.74 (0.65, 0.82)
1.0	0.37 (0.28, 0.46)	0.96 (0.91, 0.99)	0.96 (0.92, 0.99)	0.48 (0.38, 0.58)	0.77 (0.69, 0.85)

Supplementary Table S35: Plot-based functional form identification ablation results: modifying temperature for different function classes.

A.3.4 COMBINED MODALITY RESULTS



Supplementary Figure S15: Results of functional form identification task across all modalities, including combined plot and text.

A.4 PROMPTS AND TARGET DATACLASSES

For reproducibility, we provide here the prompt templates and target dataclasses we provided to Langfun for each task, except for readiness as that was performed on a proprietary dataset.

A.4.1 FUNCTIONAL FORM IDENTIFICATION

```

FunctionType = typing.Literal[
    'exponential', 'periodic', 'quadratic', 'linear', 'cubic'
]

@dataclasses.dataclass
class FunctionClassification:
    reason: str

```

2070 function_type: FunctionType
2071

2072 **Target Dataclass**

2073
2074 You are a professional data scientist with a great intuition for
2075 looking for trends in data. You are taking your next exam
2076 which has a choose-the-correct-answer format. Answer the
2077 question below to the best of your ability.
2078 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS
2079 , even if you are not sure.

2080 Q: Classify the trend in the plot shown into one of the
2081 following types. Note that a periodic function can be sine,
2082 cosine, etc.

```
2083     {%- for f in function_type %}
2084         {{ f }}
2085     {%- endfor %}
```

```
2086
2087     {{ plot }}
```

2088
2089 The data may be noisy, so do your best to see the
2090 underlying trend. Provide a reason for your answer.

2091 **Plot Prompt**

2092
2093 You are a professional data scientist with a great intuition for
2094 looking for trends in data. You are taking your next exam
2095 which has a choose-the-correct-answer format. Answer the
2096 question below to the best of your ability.
2097 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS
2098 , even if you are not sure.

2099 Q: Classify the trend in the data shown into one of the
2100 following types. Note that a periodic function can be sine,
2101 cosine, etc.

```
2102     {%- for f in function_type %}
2103         {{ f }}
2104     {%- endfor %}
```

```
2105
2106     x: {{ x }}
2107     y: {{ y }}
```

2108
2109 The data may be noisy, so do your best to see the
2110 underlying trend. Provide a reason for your answer.

2111 **Text Prompt**

2112
2113 You are a professional data scientist with a great intuition for
2114 looking for trends in data. You are taking your next exam
2115 which has a choose-the-correct-answer format. Answer the
2116 question below to the best of your ability.
2117 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS
2118 , even if you are not sure.

2119 Q: Classify the trend in the data shown into one of the
2120 following types. Note that a periodic function can be sine,
2121 cosine, etc.

```
2122     {%- for f in function_type %}
2123         {{ f }}
```

```

2124         {%- endfor %}
2125
2126     Here are the data presented as lists of values:
2127     x: {{ x }}
2128     y: {{ y }}
2129
2130     Here are the same data presented as a plot:
2131     {{ plot }}
2132
2133     The data may be noisy, so do your best to see the
2134     underlying trend. Provide a reason for your answer.

```

Text and Plot Prompt

```

2137 You are a professional data scientist with a great intuition for
2138 looking for trends in data. You are taking your next exam
2139 which has a choose-the-correct-answer format. Answer the
2140 question below to the best of your ability.
2141 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS
2142 , even if you are not sure.
2143
2144     Q: Classify the trend in the data shown into one of the
2145     following types. Note that a periodic function can be sine,
2146     cosine, etc.
2147     {%- for f in function_type %}
2148     {{ f }}
2149     {%- endfor %}
2150
2151     Here are the same data presented as a plot:
2152     {{ plot }}
2153
2154     Here are the data presented as lists of values:
2155     x: {{ x }}
2156     y: {{ y }}
2157
2158     The data may be noisy, so do your best to see the
2159     underlying trend. Provide a reason for your answer.

```

Plot and Text Prompt

A.4.2 CORRELATION OF TWO LINES

```

2163 CorrelationType = typing.Literal[
2164     'Positively_Correlated', 'Negatively_Correlated'
2165 ]
2166
2167 @dataclasses.dataclass
2168 class FunctionCorrelationResult:
2169     reason: str
2170     correlation_type: CorrelationType
2171

```

Target Dataclass

```

2174
2175 #####
2176 #####
2177 You are a professional data scientist with a great intuition for
    looking for trends in data.

```

2178 Answer the question below to the best of your ability.
 2179 The data might be noisy.
 2180 Provide reasoning.
 2181 #####
 2182 #####
 2183
 2184 Q: Observe two plots $y_1=f(x)$ and $y_2=g(x)$ over the same range of x .
 2185
 2186 plot: {{ plot }}
 2187
 2188 Find out if they are positively or negatively correlated. Provide
 2189 your reasoning.
 2190
 2191 You may want to follow the following steps to solve this problem:
 2192 Analyze the two functions and find out when one is increasing
 2193 whether the other one is decreasing.
 2194 If both functions tend to increase and decrease together, they are
 2195 positively correlated.
 2196 If one function tends to increase and the other one tends to
 2197 decrease, they are negatively correlated.

2197 Plot Prompt

2198
 2199
 2200 #####
 2201 #####
 2202 You are a professional data scientist with a great intuition for
 2203 looking for trends in data.
 2204 Answer the question below to the best of your ability.
 2205 The data might be noisy.
 2206 Provide reasoning.
 2207 #####
 2208 #####
 2209
 2210 Q: Analyze two functions $y_1=f(x)$ and $y_2=g(x)$ defined over the same
 2211 range of x .
 2212
 2213 x: {{ x }}
 2214 y1: {{ y1 }}
 2215 y2: {{ y2 }}
 2216
 2217 Find out if they are positively or negatively correlated. Provide
 2218 your reasoning.
 2219
 2220 You may want to follow the following steps to solve this problem:
 2221 Analyze the two functions and find out when one is increasing
 2222 whether the other one is decreasing.
 2223 If both functions tend to increase and decrease together, they are
 2224 positively correlated.
 2225 If one function tends to increase and the other one tends to
 2226 decrease, they are negatively correlated.

2226 Text Prompt

2228 A.4.3 2D CLUSTER COUNTING

2229
 2230 @dataclasses.dataclass
 2231 **class** ClusterCount:
 cluster_count: **int**

2232 reason: **str**

2233
2234 Target Dataclass

2235
2236
2237 Q: Here is a scatter plot of data points. The points form radial
2238 clusters. The cluster count can be 1, 2, 3, 4, 5, 6, 7, 8, or 9.
2239 Your task is to count the number of clusters from this plot.
2240 As an intermediate step, estimate the positions of the cluster
2241 centers, followed by the number of clusters.

2242 {{ plot }}

2243
2244 A:

2245
2246 Plot Prompt

2247
2248
2249 Q: Count the number of clusters. The possible number of clusters
2250 is 1, 2, 3, 4, 5, 6, 7, 8, or 9. The clusters are radial in shape.
2251 As an intermediate step, estimate the positions of the cluster
2252 centers.

2253
2254 Here are the data points. Only provide the cluster count, a rough
2255 estimate is fine too. The data may be noisy, just make your best
2256 guess, no explanations needed.

2257 x: {{ x }}

2258 y: {{ y }}

2259
2260 A:

2261 Text Prompt

2262
2263 A.4.4 DERIVATIVE IDENTIFICATION

2264
2265 MCQChoice = typing.Literal['1', '2', '3', '4']

2266
2267
2268 @dataclasses.dataclass
2269 **class** DerivativeMCQChoice:
2270 reason: **str**
2271 mcq_choice: MCQChoice

2272 Target Dataclass

2273
2274
2275 #####
2276 #####
2277 You are a professional data scientist with a great intuition for
2278 looking for trends in data.
2279 You are taking your next exam which has a choose-the-correct-
2280 answer format.
2281 Answer the question below to the best of your ability.
2282 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS,
2283 even if you are not sure.
2284 #####
2285 #####

2286 Q: Observe the trend of the first plot below. Now, based on the
 2287 trend, select one of the
 2288 following plots that corresponds to the derivative of the original
 2289 plot.
 2290
 2291 Provide your reasoning, and then return an answer. Your reasoning
 2292 should
 2293 include a description of the trend of the original plot, the
 2294 description of the trends of all the choices (1-4), and then
 2295 careful reasoning to select the correct answer. Please find the
 2296 plots
 2297 below.

2298 Original plot:
 2299 `{{ plots[0] }}`
 2300
 2301 `{%- for plot in plots[1:] %}`
 2302 `Choice: {{ loop.index }}`
 2303 `{{ plot }}`
 2304 `{%- endfor %}`

Plot Prompt

2308 #####
 2309 #####
 2310 You are a professional data scientist with a great intuition for
 2311 looking for trends in data.
 2312 You are taking your next exam which has a choose-the-correct-
 2313 answer format.
 2314 Answer the question below to the best of your ability.
 2315 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS,
 2316 even if you are not sure.
 2317 #####
 2318 #####

2319 Q: Consider the first set of data provided as lists of x and y
 2320 points, and
 2321 consider its trend.
 2322
 2323 The next four sets of data are labelled 1, 2, 3, 4 based on the
 2324 order in
 2325 which I give them to you, and represent potential derivatives as
 2326 lists of x
 2327 and dy points. Of these, choose the dataset number (1, 2, 3 or 4)
 2328 that
 2329 corresponds to the derivative of the original data. Provide your
 2330 reasoning
 2331 before your answer.

2332 Original data:
 2333 `x: {{ x }}`
 2334 `y: {{ y }}`
 2335
 2336 `{%- for d in derivatives %}`
 2337 Dataset `{{ loop.index }}`:
 2338 `x: {{ d[0] }}`
 2339 `dy: {{ d[1] }}`
`{%- endfor %}`

2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393

Text Prompt

A.4.5 QUADRATIC DERIVATIVE IDENTIFICATION

```
MCQChoice = typing.Literal['1', '2', '3', '4']

@dataclasses.dataclass
class QuadraticDerivativeMCQChoice:
    reason: str
    mcq_choice: MCQChoice
```

Target Dataclass

```
#####
#####
You are a professional data scientist with a great intuition for
looking for trends in data.
You are taking your next exam which has a choose-the-correct-
answer format.
Answer the question below to the best of your ability.
To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS,
even if you are not sure.
#####
#####

Q: Observe the slope and magnitude of the first plot of a
quadratic function
below. Now, based on both slope and magnitude, select one of the
following
plots that corresponds to the derivative of the original plot.
Note that
many of the choices might have the same slope sign, so you will
have to also
consider the magnitude of the y-values to get a correct answer.

The following reasoning will help you choose the right answer:
- From the shape of the original data, determine the function
class, and
thus the function class of the derivative.
- From the shape of the original data, determine the trend of the
expected
derivative, and thus the sign of its parameters.
- Use a few points from the original data to determine the
mangitude of the
original function's parameters, and thus the magnitude of the
expected
derivative's parameters.
- For each possible derivative choice, consider the shape of the
derivative and thereform the sign of its parameters. Also use a
few points
to determine the magnitude of the derivative choice's parameters.
- Compare the trend, sign and magnitudes of each derivative choice
with the
expected result from observing the original data, and choose the
answer that
```

2394 best matches.
 2395
 2396 Provide your reasoning, and then return an answer. Your reasoning
 2397 should
 2398 include a description of the slope and magnitude of the original
 2399 plot, the
 2400 description of the slope and magnitudes of all the choices (1-4),
 2401 and then
 2402 careful reasoning to select the correct answer. Please find the
 2403 plots
 2404 below.

2405 Original plot:
 2406 {{ plots[0] }}
 2407

2408 Choices below

```
2409
2410 {%- for plot in plots[1:] %}
2411 Choice: {{ loop.index }}
2412     {{ plot }}
2413 {%- endfor %}
```

2414 Plot Prompt - zero-shot

```
2415
2416
2417 #####
2418 #####
2419 You are a professional data scientist with a great intuition for
2420 looking for trends in data.
2421 You are taking your next exam which has a choose-the-correct-
2422 answer format.
2423 Answer the question below to the best of your ability.
2424 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS,
2425 even if you are not sure.
2426 #####
2427 #####
```

2428 Q: Observe the slope and magnitude of the first set of x and y
 2429 points
 2430 sampled from a potentially noisy quadratic function below. Now,
 2431 based on
 2432 both slope and magnitude, select one of the following choices of x
 2433 and y
 2434 points that corresponds to the derivative of the original data.
 2435 Note that
 2436 many of the choices might have the same slope sign, so you will
 2437 have to also
 2438 consider the magnitude of the y-values to get a correct answer.
 2439 The following reasoning will help you choose the right answer:
 2440 - From the shape of the original data, determine the function
 2441 class, and
 2442 thus the function class of the derivative.
 2443 - From the shape of the original data, determine the trend of the
 2444 expected
 2445 derivative, and thus the sign of its parameters.
 2446 - Use a few points from the original data to determine the
 2447 magnitude of the
 original function's parameters, and thus the magnitude of the
 expected

2448 derivative's parameters.
 2449 - For each possible derivative choice, consider the shape of the
 2450 derivative and thereform the sign of its parameters. Also use a
 2451 few points
 2452 to determine the magnitude of the derivative choice's parameters.
 2453 - Compare the trend, sign and magnitudes of each derivative choice
 2454 with the
 2455 expected result from observing the original data, and choose the
 2456 answer that
 2457 best matches.

2458 Provide your reasoning, and then
 2459 return an answer. Your reasoning should include a description of
 2460 the slope
 2461 and magnitude of the original data, the description of the slope
 2462 and
 2463 magnitudes of all the choices (1-4), and then careful reasoning to
 2464 select
 2465 the correct answer. Please find the data below.

2466 Original data:
 2467 x: {{ x }}
 2468 y: {{ y }}

2470 Choices below

```
2472 {%- for d in derivatives %}
2473 Choice: {{ loop.index }}
2474     x: {{ d[0] }}
2475     dy: {{ d[1] }}
2476 {%- endfor %}
```

Text Prompt - zero-shot

```
2480 #####
2481 #####
2482 You are a professional data scientist with a great intuition for
2483 looking for trends in data.
2484 You are taking your next exam which has a choose-the-correct-
2485 answer format.
2486 Answer the question below to the best of your ability.
2487 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS,
2488 even if you are not sure.
2489 #####
2490 #####
```

2491 Q: Observe the slope and magnitude of the first plot of a
 2492 quadratic function
 2493 below. Now, based on both slope and magnitude, select one of the
 2494 following
 2495 plots that corresponds to the derivative of the original plot.
 2496 Note that
 2497 many of the choices might have the same slope sign, so you will
 2498 have to also
 2499 consider the magnitude of the y-values to get a correct answer.
 2500
 2501 Here are some examples of how to approach this task:

```

2502 {% - for example in fewshots %}
2503
2504 ***** Example *****
2505 Original plot:
2506   {{ example.plots[0] }}
2507
2508 Choices below:
2509
2510 {% - for plot in example.plots[1:] %}
2511 Choice: {{ loop.index }}
2512   {{ plot }}
2513 {% - endfor %}
2514
2515 Reasoning:
2516 I know that the original function is quadratic, therefore the
2517 derivative
2518 must be linear.
2519
2520 I see that the original quadratic function opens
2521 {{ "up" if example.quadratic_scale > 0 else "down" }}, therefore
2522 the
2523 derivative must have a
2524 {{ "positive" if example.quadratic_scale > 0 else "negative" }}
2525 slope.
2526
2527 Furthermore I see that the original function goes from a value
2528 of y=0
2529 around x=0 to a value of y={{ example.quadratic_scale }} around
2530 x=1,
2531 therefore the original function must be of the form
2532 y={{ example.quadratic_scale }} * x^2, therefore the derivative's
2533 slope
2534 must have a magnitude of {{ 2 * example.quadratic_scale }}.
2535
2536 Of the choices I've been given:
2537 {% - for scale in example.mcq_scales %}
2538 Choice {{ loop.index }}:
2539 - the line is {{ "increasing" if scale > 0 else "decreasing"
2540 }}}, so the
2541 slope must be {{ "positive" if scale > 0 else "negative" }}
2542 - the line goes from having a value of y=0 around x=0 to a
2543 value of
2544 y={{ 2 * scale }} around x=1, so the value of the slope must
2545 be
2546 {{ 2 * scale }}
2547 - hence the derivative is a line with slope {{ 2 * scale }} and
2548 corresponds to an original quadratic function that opens
2549 {{ "up" if scale > 0 else "down" }}
2550 {% - endfor %}
2551
2552 Only choice number {{ example.mcq_correct_idx_one_indexed }} has
2553 the
2554 correct direction and slope magnitude.
2555
2556 Therefore the correct answer must be *choice number {{ example.
2557 mcq_correct_idx_one_indexed }}*.
2558
2559 {% - endfor %}

```

```

2556 ***** Your turn *****
2557
2558 Provide your reasoning, and then return an answer. Your reasoning
2559 should
2560 include a description of the slope and magnitude of the original
2561 plot, the
2562 description of the slope and magnitudes of all the choices (1-4),
2563 and then
2564 careful reasoning to select the correct answer. Please find the
2565 plots
2566 below.
2567
2568 Original plot:
2569 {{ plots[0] }}
2570
2571 Choices below
2572
2573 {%- for plot in plots[1:] %}
2574 Choice: {{ loop.index }}
2575         {{ plot }}
2576 {%- endfor %}

```

Plot Prompt - few-shot

```

2577
2578
2579 #####
2580 #####
2581 You are a professional data scientist with a great intuition for
2582 looking for trends in data.
2583 You are taking your next exam which has a choose-the-correct-
2584 answer format.
2585 Answer the question below to the best of your ability.
2586 To maximise your score on the exam ALWAYS PROVIDE A BEST GUESS,
2587 even if you are not sure.
2588 #####
2589 #####
2590
2591 Q: Observe the slope and magnitude of the first set of x and y
2592 points
2593 sampled from a potentially noisy quadratic function below. Now,
2594 based on
2595 both slope and magnitude, select one of the following choices of x
2596 and y
2597 points that corresponds to the derivative of the original data.
2598 Note that
2599 many of the choices might have the same slope sign, so you will
2600 have to also
2601 consider the magnitude of the y-values to get a correct answer.
2602
2603 Here are some examples of how to approach this task:
2604
2605 {%- for example in fewshots %}
2606
2607 ***** Example *****
2608 Original data:
2609     x: {{ example.x }}
2610     y: {{ example.y }}
2611
2612 Choices below:

```

```

2610
2611     {%- for d in example.derivatives %}
2612     Choice: {{ loop.index }}
2613         x: {{ d[0] }}
2614         dy: {{ d[1] }}
2615     {%- endfor %}
2616
2617     Reasoning:
2618     I know that the original function is quadratic, therefore the
2619     derivative
2620     must be linear.
2621
2622     I see that the original quadratic function opens
2623     {{ "up" if example.quadratic_scale > 0 else "down" }}, therefore
2624     the
2625     derivative must have a {{ "positive" if example.quadratic_scale
2626     > 0 else "negative" }} slope.
2627
2628     Furthermore I see that the original function goes from a value
2629     of y=0
2630     around x=0 to a value of y={{ example.quadratic_scale }} around
2631     x=1,
2632     therefore the original function must be of the form
2633     y={{ example.quadratic_scale }} * x^2, therefore the derivative's
2634     slope
2635     must have a magnitude of {{ 2 * example.quadratic_scale }}.
2636
2637     Of the choices I've been given:
2638     {%- for scale in example.mcq_scales %}
2639     Choice {{ loop.index }}:
2640     - the line is {{ "increasing" if scale > 0 else "decreasing"
2641     }}), so the
2642     slope must be {{ "positive" if scale > 0 else "negative" }}
2643     - the line goes from having a value of y=0 around x=0 to a
2644     value of
2645     y={{ 2 * scale }} around x=1, so the value of the slope must
2646     be {{ 2 * scale}}
2647     - hence the derivative is a line with slope {{ 2 * scale}} and
2648     corresponds to an original quadratic function that opens {{ "
2649     up" if scale > 0 else "down" }}
2650     {%- endfor %}
2651
2652     Only choice number {{ example.mcq_correct_idx_one_indexed }} has
2653     the
2654     correct direction and slope magnitude.
2655
2656     Therefore the correct answer must be *choice number {{ example.
2657     mcq_correct_idx_one_indexed }}*.
2658
2659     {%- endfor %}
2660
2661     ***** Your turn *****
2662
2663     Provide your reasoning, and then return an answer. Your reasoning
2664     should
2665     include a description of the slope and magnitude of the original
2666     data, the
2667     description of the slope and magnitudes of all the choices (1-4),
2668     and then

```

```

2664 careful reasoning to select the correct answer.
2665
2666 Please find the data below.
2667
2668 Original data:
2669     x: {{ x }}
2670     y: {{ y }}
2671
2672 Choices below
2673
2674 {%- for d in derivatives %}
2675 Choice: {{ loop.index }}
2676     x: {{ d[0] }}
2677     dy: {{ d[1] }}
2678 {%- endfor %}

```

Text Prompt - few-shot

A.4.6 FALL DETECTION FROM IMU DATA

```

2683 @dataclasses.dataclass
2684 class Fall:
2685     fall_type: Literal["ADLs", "Falls", "Near"]

```

Target Dataclass

```

2688
2689 {%- for img, label in zip(few_shot_series, few_shot_labels) %}
2690 Given that the following plot was classified as {{ label }}:
2691 {{ img }}
2692 {%- endfor %}
2693 Classify the following plot in one of the following classes: ADLs,
2694 Falls, Near.
2695 {{ sample }}
2696 ALWAYS provide a best guess, since you will be graded on your
2697 response.

```

Plot Prompt

```

2700
2701 {%- for data, label in zip(few_shot_series, few_shot_labels) %}
2702 Given that the following time-series data was classified as {{
2703 label }}:
2704 {{ data }}
2705 {%- endfor %}
2706 Classify the following time-series data as 'ADLs', 'Falls', or '
2707 Near':
2708 {{ sample }}
2709 ALWAYS provide a best guess, since you will be graded on your
2710 response.

```

Text Prompt

A.4.7 ACTIVITY RECOGNITION FROM IMU DATA

```

2714 @dataclasses.dataclass
2715 class ActivityNoUnknown:
2716     activity_type: Literal["bike", "sit", "stand", "walk", "stairs"]

```

2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753
2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771

Target Dataclass

```
{%- for imgs, label in zip(few_shot_series, few_shot_labels) %}
Given that the following plots were classified as {{ label }}:
{%- for img in imgs %}
{{ img }}
{%- endfor %}
{%- endfor %}
Classify the following plots in one of the following classes: {{
classes }}.
{%- for img in sample %}
{{ img }}
{%- endfor %}
ALWAYS provide a best guess, since you will be graded on your
response.
```

Plot Prompt

```
{%- for data, label in zip(few_shot_series, few_shot_labels) %}
Given that the following time-series data was classified as {{
label }}:
{{ data }}
{%- endfor %}
Classify the following time-series data in one of the following
classes: {{ classes }}.
{{ sample }}
ALWAYS provide a best guess, since you will be graded on your
response.
```

Text Prompt