

Visual Time Series Forecasting: An Image-driven Approach

Naftali Cohen
J. P. Morgan AI Research
New York, NY
naftali.cohen@jpmchase.com

Srijan Sood
J. P. Morgan AI Research
New York, NY
srijan.sood@jpmchase.com

Zhen Zeng
J. P. Morgan AI Research
New York, NY
zhen.zeng@jpmchase.com

Tucker Balch
J. P. Morgan AI Research
New York, NY
tucker.balch@jpmchase.com

Manuela Veloso
J. P. Morgan AI Research
New York, NY
manuela.veloso@jpmchase.com

ABSTRACT

In this work, we address time-series forecasting as a computer vision task. We capture input data as an image and train a model to produce the subsequent image. This approach results in predicting distributions as opposed to pointwise values. To assess the robustness and quality of our approach, we examine various datasets and multiple evaluation metrics. Our experiments show that our forecasting tool is effective for cyclic data but somewhat less for irregular data such as stock prices. Importantly, when using image-based evaluation metrics, we find our method to outperform various baselines, including ARIMA, and a numerical variation of our deep learning approach.

CCS CONCEPTS

• Computing methodologies → Image representations; • Mathematics of computing → Time series analysis.

KEYWORDS

time-series forecasting, Image representations, ARIMA, visualizations

ACM Reference Format:

Naftali Cohen, Srijan Sood, Zhen Zeng, Tucker Balch, and Manuela Veloso. 2021. Visual Time Series Forecasting: An Image-driven Approach. In *MileTS '21: 7th SIGKDD Workshop on Mining and Learning from Time Series, August 14th, 2021, Virtual Conference*. ACM, New York, NY, USA, 5 pages.

1 INTRODUCTION AND RELATED WORK

Time series forecasting is a standard statistical task that concerns predicting future values given historical information. Conventional forecasting tasks range from uncovering simple periodic patterns to forecasting intricate nonlinear patterns. The prevailing and most widely used forecasting techniques include linear regression, exponential smoothing, and ARIMA (e.g., [10, 17, 22]). In recent years, modern approaches emerge as tree-based algorithms, ensemble

methods, neural network autoregression, and recurrent neural networks (e.g., [10]). These methods are useful for highly nonlinear and inseparable data but are often considered less stable than the more traditional approaches (e.g., [17, 18]).

In the last few years, deep learning approaches have been applied in the domain of time series analysis, for forecasting [4, 11, 27, 28], as well as unsupervised approaches for pre-training, clustering, and distance calculation [1, 24, 29, 31]. The common theme across these works is their use of stacked autoencoders (with different variations – vanilla, convolutional, recurrent, etc.) on numeric time series data. Autoencoders have also shown promise in the computer vision domain across tasks as image denoising [3, 14], image compression [2], and image completion and in-painting [20, 23].

This paper follows these studies and presents a new perspective on numerical time series forecasting by transforming the problem completely into the computer-vision domain. We capture input data as images and build a network that outputs corresponding subsequent images. To the best of our knowledge, this is the first study that aims at explicit visual forecasting of time series data as plots. Previous researches leveraged computer vision for time-series data but focused on classifying trade patterns [7, 8], numeric forecast [6], learning weights to combine multiple statistical forecasting methods [19], and video prediction for multivariate economic forecasting [32]. We follow up on these approaches but focus on an explicit regression-like image prediction task.

This work presents a few advantages. Visual time series forecasting is a data-driven non-parametric method, not constrained to a predetermined set of parameters. Thus, the approach is flexible and adaptable to many data forms, as shown by application across various datasets. This bears a stark contrast with classical time series forecasting approaches that are often tailored to the particularity of the data in hand. The main advantage of this method is that its prediction is independent of other techniques. This is important as it was repeatedly shown that an aggregate of independent techniques outperforms the best-in-class method (e.g., [10, 12, 15]). Secondly, visual predictions result in inherent uncertainty estimates as opposed to pointwise estimates, as they represent distributions over pixels as opposed to explicit value prediction. In addition, financial time series data are often presented and act upon without having access to the underlying numeric information (e.g., financial trading using the smartphone applications). Thus, it seems viable to examine the value in inferring using visualizations alone. Lastly, as

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MileTS '21, August 14th, 2021, Virtual Conference

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.

will be discussed later on, we show that transforming the continuous numeric data to a discrete bounded space using visualization results in robust and stable predictions. We evaluate predictions using multiple metrics. When considering object-detection metrics such as Intersection-over-Union (IoU), visual forecasting outperforms the corresponding numeric baseline. However, when utilizing more traditional time-series evaluation metrics as the symmetric mean absolute percentage error (SMAPE), we find the visual view to perform similarly to its numerical baselines.

2 DATASETS

This paper uses four datasets - two synthetic and two real - with varying degrees of periodicity and complexity to examine the utility of forecasting using images. Each dataset consisted of approximately 40-50k training samples, 4-5k validation samples, and 15k testing samples. Figure 1 shows examples of the data, and the supplementary material contains a detailed description of each of the datasets and how they were curated.

Harmonic: multi-periodic data sampled from harmonic functions. It is derived synthetically with a linearly additive two-timescale harmonic generating function consisting of sine waves on two time-scales: short oscillations that are composed on a much longer wave trains.

OU: synthesized using mean-reverting time series based on Ornstein-Uhlenbeck (OU) processes [5]. These often resemble characteristics of financial interest rates or volatility: noisy on finer scales but predictable on the larger scale.

ECG: signals measured from 17 different people adopted from MIT-BIH Normal Sinus Rhythm Database [13]. The data has prominent spikes about every second, which makes the data predictable. However, there is noticeable noise between spikes that is much harder to predict.

Financial: stock data from Yahoo! Finance consisting of daily Adjusted Close values of stocks that contributed to the S&P-500 index from 2000-2019. Each time series segment consists of 80 days and is generally considered notoriously hard to predict [26].

3 PROBLEM STATEMENT

Given a time series signal, our goal is to produce a visual forecast of its future. We approach this problem by first converting the numeric time series into an image (detailed procedure described in supplementary material), and then producing a corresponding forecast image using deep-learning techniques. By doing so, we obtain an image in which the pixel values in each column sum to 1; each column can be perceived as a discrete probability distribution (see Figure 2). Columns represent the independent variable time, while rows capture the dependent variable: pixel intensity. The value of the time series S at time t is now simply the pixel index r (row) at that time (column) with the highest intensity.

Let X be the set of images of input time series signals, and Y be the set of corresponding forecast output images. The overlap constant c defines the overlap fraction between the input image $x \in X$ and the forecast $y \in Y$, where $c = 1$ implies $x = y, \forall x \in X$, and $c = 0$ implies that $x \cap y = \emptyset, \forall x \in X$, i.e., x and y are distinct. In our experiments, we use $c = 0.75$ which means the first 75% of the forecast image y is simply a reconstruction of the later

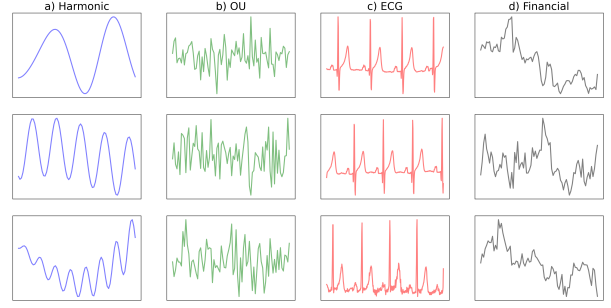


Figure 1: Sampled examples of the four datasets: Harmonic, OU, ECG, and Financial.

75% of the input image x , and the rest 25% of y corresponds to visual forecasting of the future. We chose $c = 0.75$ such that the reconstructed overlap region (first 75% in y) serves as a sanity check on the effectiveness of a forecasting method, and the prediction region (later 25% in y) provides forecasting into the near future. Please refer to the supplementary material for an illustration.

4 METHOD

4.1 Image-to-Image Regression

As mentioned in Section 1, recent work has seen the extensive use of autoencoders in both the time series and computer vision domains. Following these, we extend the use of autoencoders to our image-to-image time series forecasting setting. We use a simplistic convolutional autoencoder to produce a visual forecast image with the continuation of an input time series image, by learning an undercomplete mapping $g \circ f$,

$$\hat{y} = g(f(x)), \forall x \in X,$$

where the encoder network $f(\cdot)$ learns meaningful patterns and projects the input image x into an embedding vector, and the decoder network $g(\cdot)$ reconstructs the forecast image from the embedding vector. We purposely do not use sequential information or LSTM cells as we wish to examine the benefits of framing the regression problem in an image setting. This can later be extended to more complex architectures.

We call this method **VisualAE**. We used 2D convolutional layers with a kernel size of 5×5 , stride 2, and padding 2. All layers are followed by ReLU activation and batch normalization. The encoder network consists of 3 convolutional layers which transform a $80 \times 80 \times 1$ input image to $10 \times 10 \times 512$, after which we obtain an embedding vector of length 512 using a fully connected layer. This process is then mirrored for the decoder network, resulting in a forecast image of dimension 80×80 . We include a diagram illustrating this architecture in the supplementary material.

4.2 Loss Functions

We care about the likelihood of pixel intensity in a particular location (row) in each column of the forecast image. This can be achieved by leveraging metrics that compare two probability distributions. We do so in a column-wise manner: the loss L to compare target ground-truth (GT) image y with prediction image \hat{y} is the

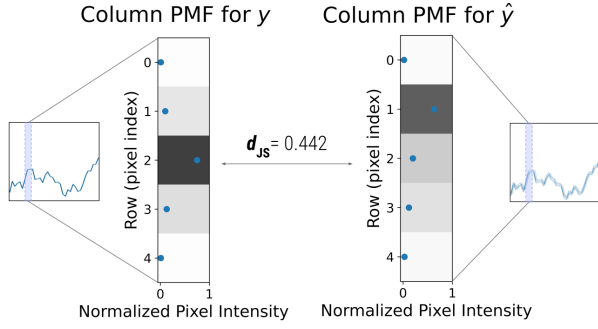


Figure 2: A depiction of comparison of two sample column probability distributions $y = [0.01, 0.1, 0.75, 0.13, 0.01]$ and $\hat{y} = [0.02, 0.63, 0.2, 0.12, 0.03]$.

sum of column-wise distances between the two,

$$L(y, \hat{y}) = \sum_{i=1}^w d(y_i, \hat{y}_i),$$

where y_i, \hat{y}_i are the i th column in the ground truth and forecast images, d is any distance measure between two distributions (y_i and \hat{y}_i in this case), and w is the width of images. This process is depicted in Figure 2.

Measures such as the Kullback-Leibler Divergence have been extensively used as loss functions ([15]), as they provide a way of computing the distance from an *approximate distribution* Q to a *true distribution* P . In this study, following [16], we choose d to be the **Jensen-Shannon Divergence** (JSD), which is a symmetric, more stable version of the Kullback-Leibler Divergence having the property that $D_{JS}(P\|Q) = D_{JS}(Q\|P)$. Here, JSD is computed as

$$D_{JS}(P\|Q) = \frac{1}{2}D_{KL}(P\|M) + \frac{1}{2}D_{KL}(Q\|M)$$

where $M = \frac{1}{2}(P + Q)$.

5 EXPERIMENTS

We experimented with four datasets: Harmonic, OU, ECG, and Financial, as they cover a wide range of complexity and predictability in time series data (illustrated in the dataset analysis in supplementary). In this study we used PyTorch Lightning [9, 25] for implementation and Nvidia Tesla T4 GPUs in our experiments. As described in Section 3, there is a 75% overlap between input and output. In our experiments, each sample contains 80 datapoints; we aim to forecast the last 20 datapoints (last 25%) of the output image. We benchmark the proposed method against three baseline methods.

5.1 Methods

We summarize the benchmarked methods as following. Please refer to a thorough description (including training details, and data preprocessing) of each method in the supplementary material.

- **VisualAE**: This is the proposed method as discussed in Section 4.1. We train on images with size 80×80 .
- **NumAE** (Numeric AE): We also train an autoencoder network to produce numeric forecasts of the original numerical time series signal.

- **ARIMA**: Autoregressive Integrated Moving Average (ARIMA) models are a class of methods that are designed to capture autocorrelations in the data.
- **RandomWalk**: We used the random walk without drift model as a naive numeric forecasting baseline for comparison ([30]).

5.2 Forecast Accuracy Metrics

We use a variety of measures to assess the accuracy of forecast predictions from each method. Some of these metrics are extensively used in the time series forecasting domain, whereas the others we extend from the overarching machine learning field to this task.

The baseline methods **ARIMA**, **NumAE** and **RandomWalk** produce continuous numeric forecasts, whereas our method **VisualAE** produces an image. Accordingly, we convert this image back to a numeric forecast which we can use to assess predictions using the metrics described in Section 5.2.1. Similarly, to leverage the image based metrics described in Section 5.2.2, we transform the numeric predictions of the baseline methods into images. We discuss the interplay between these metrics across Section 5.3, along with in-depth discussions and insights regarding comparisons between numeric and image based metrics in the supplementary material.

5.2.1 Numeric Measures. We use the Symmetric Mean Absolute Percentage Error (**SMAPE**), and the Mean Absolute Scaled Error (**MASE**) to evaluate the numeric forecasts. These two metrics are widely used in the literature for forecast accuracy evaluation [21]. Please refer to supplementary materials for equations of these metrics.

5.2.2 Image based Measures. In addition to utilizing traditional forecasting error metrics, we can measure the similarity between the predicted image and the ground-truth image in our setting to evaluate forecast accuracy. We use Jensen-Shannon Divergence (**JSD**), which is the same as the loss described in Section 4.2. In addition, we use an extended version of Intersection-over-Union (**IoU**) to measure image similarity columnwise. We first obtain the 1D bounding box of non-zero pixels for each column, then compute the IoU between bounding boxes of each corresponding column in the ground-truth and predicted images. This ranges from 0.0 to 1.0, with higher values indicating better forecasts.

5.3 Results

All reported metrics mentioned in Section 5.2 are over the unseen future prediction region. For both **VisualAE** and **NumAE**, we averaged these metrics over five independently trained models with different random weight initializations. We demonstrate that the proposed method **VisualAE** outperforms baseline methods **NumAE**, **RandomWalk**, and **ARIMA** across all four datasets when evaluated using image-based metrics (such as IoU). However, as we will discuss in this section, traditional numeric metrics are inconsistent with this finding. We demonstrate the value of using a visual approach to time-series forecasting, and how image-based evaluation metrics can help address some of the caveats of traditional numeric metrics.

We report the mean and standard deviation of various prediction accuracy metrics in Table 1. **VisualAE** achieves higher IoU

	Method	SMAPE $\mu \pm \sigma$	MASE $\mu \pm \sigma$	IoU $\mu \pm \sigma$	JSD $\mu \pm \sigma$
Harmonic	RandomWalk	1.239 \pm 0.440	5.106 \pm 3.405	0.179 \pm 0.060	0.501 \pm 0.043
	NumAE	0.480 \pm 0.297	1.258 \pm 1.081	0.423 \pm 0.107	0.334 \pm 0.103
	ARIMA	0.580 \pm 0.398	2.694 \pm 3.350	0.447 \pm 0.238	0.343 \pm 0.186
	VisualAE	0.527 \pm 0.303	1.518 \pm 1.482	0.460 \pm 0.088	0.271 \pm 0.115
OU	RandomWalk	0.018 \pm 0.069	1.007 \pm 0.385	0.257 \pm 0.021	0.381 \pm 0.019
	NumAE	0.014 \pm 0.056	471.411 \pm 8486.706	0.165 \pm 0.076	0.543 \pm 0.052
	ARIMA	0.014 \pm 0.056	0.736 \pm 0.133	0.141 \pm 0.014	0.556 \pm 0.011
	VisualAE	0.014 \pm 0.060	0.748 \pm 0.119	0.469 \pm 0.017	0.257 \pm 0.010
ECG	RandomWalk	1.173 \pm 0.463	1.551 \pm 1.384	0.164 \pm 0.014	0.501 \pm 0.021
	NumAE	1.097 \pm 0.200	0.979 \pm 0.280	0.278 \pm 0.047	0.463 \pm 0.051
	ARIMA	1.409 \pm 0.305	1.535 \pm 1.688	0.160 \pm 0.011	0.576 \pm 0.009
	VisualAE	0.596 \pm 0.254	1.658 \pm 0.321	0.485 \pm 0.022	0.230 \pm 0.041
Financial	RandomWalk	0.036 \pm 0.028	3.364 \pm 2.217	0.186 \pm 0.054	0.475 \pm 0.050
	NumAE	0.036 \pm 0.028	3.364 \pm 2.205	0.132 \pm 0.069	0.598 \pm 0.059
	ARIMA	0.042 \pm 0.035	4.034 \pm 14.697	0.119 \pm 0.072	0.606 \pm 0.063
	VisualAE	0.043 \pm 0.028	4.007 \pm 2.084	0.212 \pm 0.080	0.511 \pm 0.070

Table 1: Summary of various metrics on out-of-sample data with mean \pm standard deviation for the forecast region. Lower SMAPE/MASE/JSD error (or higher IoU score) implies better prediction accuracy.

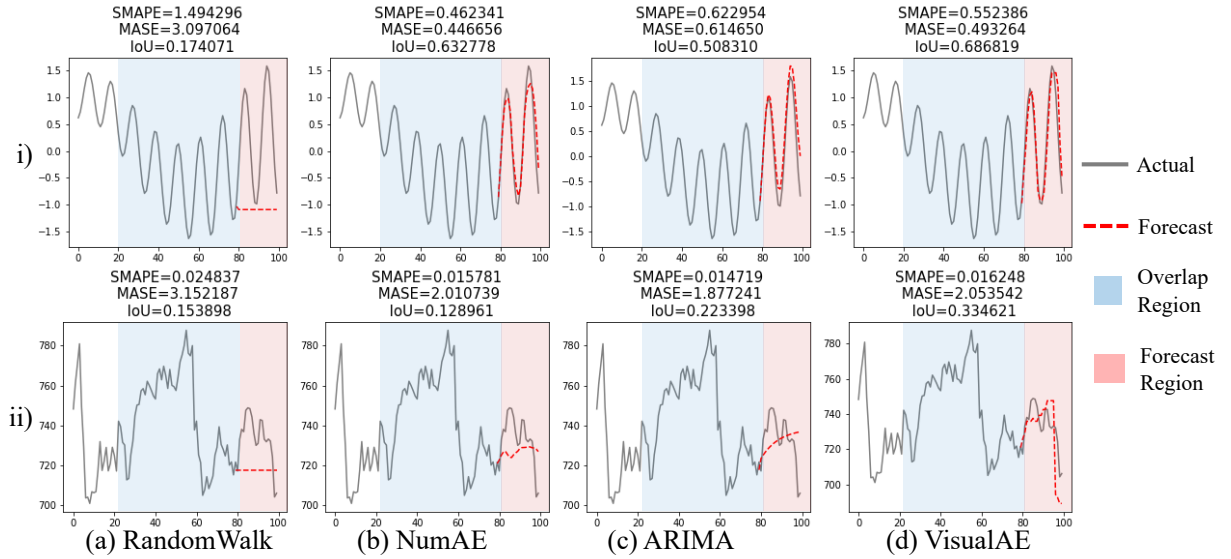


Figure 3: IoU metric better captures visual forecast accuracy compared to traditional numeric metrics SMAPE and MASE.

scores than all baselines across the four datasets. The same holds true for JSD (with the exception of **RandomWalk** scoring better in the Financial dataset). The numeric metrics are often inconsistent – within themselves (SMAPE and MASE) – as well as across the four datasets. According to the numeric metrics, **VisualAE** is a close second (if not similar) to **NumAE**, with the exception of the ECG dataset, where **VisualAE** performs the best, and the OU dataset, where **ARIMA** and **VisualAE** perform similarly to **NumAE**. Please refer to our supplementary material for a more detailed discussion on the characteristics of benchmarked methods, along with qualitative examples.

5.4 Numeric vs. Image based Metrics

In Table 1, numeric metrics are often inconsistent with the image-based ones, and sometimes do not agree amongst each other (e.g., SMAPE & MASE values for OU dataset). They are sensitive and often fail to recognize good quality forecasts (e.g., **RandomWalk** reportedly performing the best for the Financial dataset). Picking a percentage error such as SMAPE also carries the inability to compare forecast method quality across series (e.g., low errors in the Financial dataset do not capture that it is the most challenging to predict).

The IoU metric is able to capture this information across the datasets, along with preserving rank-ordering of forecast quality

amongst the four methods. As shown in Figure 3, the IoU metric is better at discerning which forecast better captures ground-truth trends. This is evident with higher IoU values when the visual shape of predictions matches the ground truth well. We believe using a two-pronged approach of utilizing both numeric and visual approaches holds immense value for the field of time series forecasting.

6 SUMMARY AND CONCLUSION

To the best of our knowledge, this study is the first to explicitly forecast time series using visual representations of numeric data. We show that image-based measures can capture prediction quality more consistently than traditional numeric metrics. The proposed visual forecasting approach, albeit simplistic, performs well across datasets. Our findings show promising results for both periodic time series (including abrupt spikes in ECG) and irregular financial data. We believe that leveraging visual approaches holds immense promise for the field of time series forecasting in the future, especially when used in conjunction with traditional methods.

Disclaimer: This paper was prepared for information purposes by the Artificial Intelligence Research group of J. P. Morgan Chase & Co. and its affiliates (“J. P. Morgan”), and is not a product of the Research Department of J. P. Morgan. J. P. Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

©2021 J. P. Morgan Chase & Co. All rights reserved.

REFERENCES

- [1] Abubakar Abid and James Y Zou. 2018. Learning a warping distance from unlabeled time series using sequence autoencoders. In *Advances in Neural Information Processing Systems*. 10547–10555.
- [2] Pinar Akyazi and Touradj Ebrahimi. 2019. Learning-Based Image Compression using Convolutional Autoencoder and Wavelet Decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.
- [3] Guillaume Alain and Yoshua Bengio. 2014. What regularized auto-encoders learn from the data-generating distribution. *The Journal of Machine Learning Research* 15, 1 (2014), 3563–3593.
- [4] Wei Bao, Jun Yue, and Yulei Rao. 2017. A deep learning framework for financial time series using stacked autoencoders and long-short term memory. *PloS one* 12, 7 (2017), e0180944.
- [5] David Byrd. 2019. Explaining agent-based financial market simulation. *arXiv preprint arXiv:1909.11650* (2019).
- [6] Naftali Cohen, Tucker Balch, and Manuela Veloso. 2019. The effect of visual design in image classification. *arXiv preprint arXiv:1907.09567* (2019).
- [7] Naftali Cohen, Tucker Balch, and Manuela Veloso. 2019. Trading via image classification. *arXiv preprint arXiv:1907.10046* (2019).
- [8] Bairui Du and Paolo Barucca. 2020. Image Processing Tools for Financial Time Series Classification. *arXiv preprint arXiv:2008.06042* (2020).
- [9] WA Falcon. 2019. PyTorch Lightning. *GitHub*. Note: <https://github.com/PyTorchLightning/pytorch-lightning> 3 (2019).
- [10] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. 2001. *The elements of statistical learning*. Springer series in statistics New York.
- [11] Andre Gensler, Janosch Henze, Bernhard Sick, and Nils Raabe. 2016. Deep Learning for solar power forecasting—An approach using AutoEncoder and LSTM Neural Networks. In *2016 IEEE international conference on systems, man, and cybernetics (SMC)*. IEEE, 002858–002865.
- [12] Aurélien Geron. 2019. *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems*. O'Reilly Media.
- [13] Ary I Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. 2000. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *circulation* 101, 23 (2000), e215–e220.
- [14] Lovedeep Gondara. 2016. Medical image denoising using convolutional denoising autoencoders. In *2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW)*. IEEE, 241–246.
- [15] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. 2016. *Deep learning*. MIT press Cambridge.
- [16] Ferenc Huszar. 2015. How (not) to train your generative model: Scheduled sampling, likelihood, adversary? *arXiv preprint arXiv:1511.05101* (2015).
- [17] Rob J Hyndman and George Athanasopoulos. 2018. *Forecasting: principles and practice*. OTexts.
- [18] R Krispin. 2019. *Hands-On Time Series Analysis with R: Perform time series analysis and forecasting using R*.
- [19] Xixi Li, Yanfei Kang, and Feng Li. 2020. Forecasting with time series imaging. *Expert Systems with Applications* 160 (2020), 113680.
- [20] Yijun Li, Sifei Liu, Jimei Yang, and Ming-Hsuan Yang. 2017. Generative face completion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3911–3919.
- [21] Spyros Makridakis and Michele Hibon. 2000. The M3-Competition: results, conclusions and implications. *International journal of forecasting* 16, 4 (2000), 451–476.
- [22] Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. 2020. The M4 Competition: 100,000 time series and 61 forecasting methods. *International Journal of Forecasting* 36, 1 (2020), 54–74.
- [23] Xiao-Jiao Mao, Chunhua Shen, and Yu-Bin Yang. 2016. Image restoration using convolutional auto-encoders with symmetric skip connections. *arXiv preprint arXiv:1606.08921* (2016).
- [24] S Mostafa Mousavi, Weiqiang Zhu, William Ellsworth, and Gregory Beroza. 2019. Unsupervised clustering of seismic signals using deep convolutional autoencoders. *IEEE Geoscience and Remote Sensing Letters* 16, 11 (2019), 1693–1697.
- [25] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems*. 8026–8037.
- [26] Lasse Heje Pedersen. 2019. *Efficiently inefficient: how smart money invests and market prices are determined*. Princeton University Press.
- [27] Pablo Romeu, Francisco Zamora-Martinez, Paloma Botella-Rocamora, and Juan Pardo. 2015. Stacked denoising auto-encoders for short-term time series forecasting. In *Artificial Neural Networks*. Springer, 463–486.
- [28] Alaa Sagheer and Mostafa Kotb. 2019. Time series forecasting of petroleum production using deep LSTM recurrent networks. *Neurocomputing* 323 (2019), 203–213.
- [29] Alaa Sagheer and Mostafa Kotb. 2019. Unsupervised pre-training of a Deep LSTM-based Stacked Autoencoder for Multivariate time series forecasting problems. *Scientific Reports* 9, 1 (2019), 1–16.
- [30] Steven E Shreve. 2004. *Stochastic calculus for finance II: Continuous-time models*. Vol. 11. Springer Science & Business Media.
- [31] Neda Tavakoli, Sima Siame-Namini, Mahdi Adl Khanghah, Fahimeh Mirza Soltani, and Akbar Siame Namin. 2020. Clustering Time Series Data through Autoencoder-based Deep Learning Models. *arXiv preprint arXiv:2004.07296* (2020).
- [32] Zeng Zhen, Tucker Balch, and Manuela Veloso. 2021. Deep Video Prediction for Economic Forecasting. *arXiv preprint arXiv:2102.12061* (2021).

Supplementary: Visual Time Series Forecasting: An Image-driven Approach

Naftali Cohen
J. P. Morgan AI Research
New York, NY
naftali.cohen@jpmchase.com

Srijan Sood
J. P. Morgan AI Research
New York, NY
srijan.sood@jpmchase.com

Zhen Zeng
J. P. Morgan AI Research
New York, NY
zhen.zeng@jpmchase.com

Tucker Balch
J. P. Morgan AI Research
New York, NY
tucker.balch@jpmchase.com

Manuela Veloso
J. P. Morgan AI Research
New York, NY
manuela.veloso@jpmchase.com

ACM Reference Format:

Naftali Cohen, Srijan Sood, Zhen Zeng, Tucker Balch, and Manuela Veloso. 2021. Supplementary: Visual Time Series Forecasting: An Image-driven Approach. In *MileTS '21: 7th SIGKDD Workshop on Mining and Learning from Time Series, August 14th, 2021, Virtual Conference*. ACM, New York, NY, USA, 5 pages.

1 DATASETS

This paper uses four datasets, two synthetic and two real, with varying degrees of periodicity and complexity to examine the utility of forecasting using images.

1.1 Synthetic Data

We generate two different series for the synthetic datasets: multiperiodic data sampled from harmonic functions and mean-reverting data generated following the Ornstein–Uhlenbeck process.

1.1.1 Harmonic data. The first dataset is derived synthetically and is designed to be involved but still with a prominent, repeated signal. We synthesized the time series s_t with a linearly additive two-timescale harmonic generating function,

$$s_t = (A_1 + B_1 t) \sin(2\pi t/T_1 + \phi_1) + (A_2 + B_2 t) \sin(2\pi t/T_2 + \phi_2),$$

where the time t varies from $t = 1$ to $t = T$, and T denotes the total length of the time series. The multiplicative amplitudes A_1 and A_2 are randomly sampled from a Gaussian distribution $\mathcal{N}(1, 0.5)$, while the amplitude of the linear trends B_1 and B_2 are sampled from a uniform distribution $\mathcal{U}(-1/T, 1/T)$. The driving time scales are short (T_1) and long (T_2) relative to the total length of T . Thus, $T_1 \sim \mathcal{N}(T/5, T/10)$, while $T_2 \sim \mathcal{N}(T, T/2)$. Lastly, the phase shifts ϕ_1 and ϕ_2 are sampled from a uniform distribution $\mathcal{U}(0, 2\pi)$. We generated and used 42,188 examples as a train set, 4,687 for the validation set, and 15,625 for the test set. Each time series differ concerning the possible combination of tuning parameters. Panel (a) (Figure 1 in the main paper) shows three sampled examples of

the harmonic data and it is easy to see that the synthetic time series consist of two time-scales: short oscillations that are composed on a much longer wave trains.

1.1.2 OU Data. We synthesized mean-reverting time series based on Ornstein–Uhlenbeck (OU) process as described in [1]. A mean-reverting time series tends to drift towards a fundamental mean value. We chose to synthesize the mean-reverting time series to resemble the characteristics of financial interest rates or volatility. OU’s stochastic nature makes it noisy on fine scales but predictable on the larger scale, which is the focus of this study. Specifically, we generated the OU dataset following the equation adopted from [1] with,

$$s_t \sim \mathcal{N}(\mu + (s_{t-1} - \mu)e^{-\gamma t}, \frac{\sigma^2}{2\gamma}(1 - e^{-2\gamma t})),$$

where μ is the mean value that the time series reverts back to, and s_0 starts at μ . We used mean reversion rate $\gamma \sim \mathcal{N}(8e^{-8}, 4e^{-8})$ with units ns^{-1} , and a volatility value $\sigma \sim \mathcal{N}(1e^{-2}, 5e^{-3})$. Overall, we generated the time series by sampling s_t at every minute. We generated and used 45,000 examples as a train set, 5,000 for the validation set, and 15,000 for the test set. Similar to the Harmonic data, each time series differ concerning the possible combination of tuning parameters. Panel (b) (Figure 1 in the main paper) shows three samples of the OU data. One can see that the OU data tend to be noisy with uncorrelated ups and downs, but on larger scales, the data is concentrated in the middle of the image as values drift toward the mean due to its reversion constraint.

1.2 Real Data

Along with the synthetically generated data, we use two real-world time series datasets.

1.2.1 ECG data. The ECG data is measured information from 17 different people adopted from MIT-BIH Normal Sinus Rhythm Database [3]. We curated 18 hours of data for each subject after manually examining the data’s consistency and validity by analyzing the mean and standard deviation of the time series data for each subject (not shown). For each subject, we consider segments of 2.56 seconds (corresponding to 320 data points) sampled randomly from the data. These are then downsampled to 80 data points to be on-par with the other datasets. 13 out of the 17 subjects are used as training data while the other 4 are used as out-of-sample

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MileTS '21, August 14th, 2021, Virtual Conference

© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM.

testing data. Overall, we sampled 42,188 examples for the training set, while from the test data, we sampled 4,687 as a validation set and 15,625 as a test set. Panel (c) (Figure 1 in the main paper) shows three sampled examples of the ECG data. One can see that the data has prominent spikes about every second, which makes the data predictable. However, there is noticeable noise between spikes that is much harder to predict.

1.2.2 Financial data. The last dataset is financial stock data from Yahoo! Finance. The data consists of daily Adjusted Close values of stocks that contributed to the S&P-500 index since 2000. Each time series segment consists of 80 days and is standardized by subtracting the mean and dividing by the standard deviation for each segment separately. For the train data, we sampled information randomly from the year 2000 to 2014, while for the test, we sampled information from 2016 to 2019. Overall, we sampled 38,746 examples as a training set, while from the test data, we sampled 4,306 as a validation set and 15,625 as a test set.

Panel (d) (Figure 1 in the main paper) shows three sampled examples of the financial data. Here, one can see that the data is much less predictable than the previous three. Although financial data is persistent with sequentially related information, it is hard to spot repeated signals that will make the data predictable. Indeed, the prevailing theory of financial markets argues that markets are very efficient, and their future movements are notoriously hard to predict, especially given price information alone (e.g., [6]).

1.3 Complexity of Time-Series Data

To provide a reference for how the time series across our datasets vary, we measured the complexity of each with Weighted Permutation Entropy [2] (WPE). WPE measures the entropy of the ordinal patterns in time series data. In our experiments, we consider the ordinal patterns of each triplet (s_{t-1}, s_t, s_{t+1}) along a given 1-d time series s where t stands for time. Empirically using a triplet (other than segment of length at 2, or 4 or larger numbers) to calculate the complexity of time-series is informative and also easily tractable. We normalize WPE into $[0, 1]$ with normalized entropy. The larger WPE, the more complex the data is.

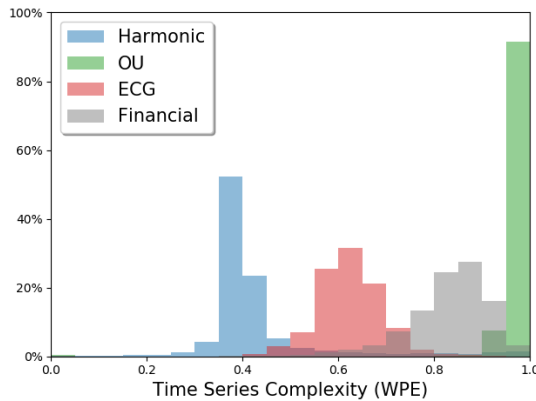


Figure 1: Distribution of dataset's complexity measured using Weighted Permutation Entropy.

As shown in Figure 1, we can see that the datasets cover a broad range of complexity. As expected, the simplest data is Harmonic with its deterministic periodicity. The ECG data is also periodic but more complex due to its irregularities between the spikes. The Financial data is not cyclic at all with almost random movement of fine scales, and, therefore, more complex than both the Harmonic and ECG. The OU data exhibits even more random oscillations and abrupt changes compared to other datasets, thus it is measured as the most complex dataset. However, on the larger scale, the OU data bounces around a hidden mean value with bounded noise, making it possible to predict future value mean and ranges.

2 ILLUSTRATION OF PREDICTION FRAMEWORK

Figure 2 illustrates the overlap between the input and output image. The network is tasked to learn to reconstruct the overlap region, as well as to predict the future region in the output image.

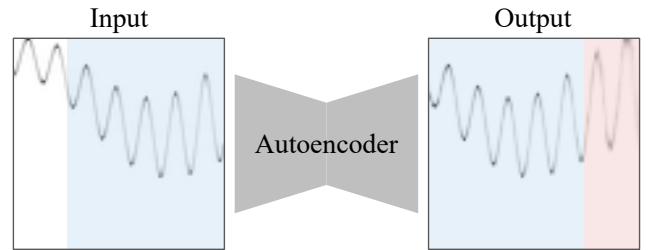


Figure 2: An overview of our problem setup. The blue region shows a 75% overlap between the input x and output y ; the forecast region \hat{y} is denoted in red.

3 METHOD DETAILS

3.1 Covert Numeric Time Series to Images

Given a 1-d numeric time series $S = [s_0, \dots, s_T]$ with $s_t \in \mathbb{R}$, we convert S into a 2-d image x by plotting it out, with t being the horizontal axis and s_t being the vertical axis¹. We standardize each converted image x through following pre-processing steps. First, pixels in x are scaled to $[0, 1]$ and negated (i.e., $x = 1 - x/255$) so that the pixels corresponding to the plotted time series signal are bright (values close to 1), whereas the rest of the background pixels become dark (values close to 0). Note that there can be multiple bright (non-zero) pixels in each column due to anti-aliasing while plotting the images.

Upon normalizing each column in x such that the pixel values in each column sum to 1, each column can be perceived as a discrete probability distribution (see Figure 2 from the main paper). Columns represent the independent variable time, while rows capture the dependent variable: pixel intensity. The value of the time series S at time t is now simply the pixel index r (row) at that time (column) with the highest intensity.

¹We plotted each time series S with bounded intervals. The interval for x -axis is $[0 - \epsilon, T + \epsilon]$, whereas the interval for y -axis is $[\min(s_t) - \epsilon, \max(s_t) + \epsilon]$, where $\epsilon = 10^{-6}$.

Predictions are made over normalized data. To preserve the ability to forecast in physical units, we utilize the span of the input raw data values to transform forecasts to the corresponding physical scales.

3.2 VisualAE Architecture

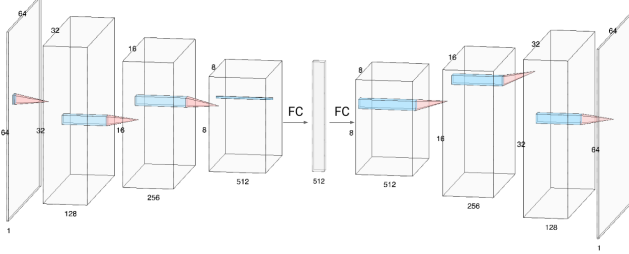


Figure 3: The architecture of undercomplete convolutional autoencoder network used in this study.

4 DETAILED DESCRIPTIONS ON BENCHMARKED METHODS

4.1 VisualAE

This is the proposed method as discussed in Section 4.1 from the main paper. We train on images with size 80×80 . We use a batch size of 128 and early stopping after 15 consecutive non-improving validation epochs to avoid overfitting during training. We start with a learning rate of 0.1, which is decayed by a factor of 0.1 (till $1e-4$) after every 5 non-improving validation epochs.

4.2 NumAE (Numeric AE)

We also train an autoencoder network to produce numeric forecasts of the original numerical time series signal. The numeric input and output time series are standardized using min-max normalization (with bounds obtained from the input to avoid leakage to future). The autoencoder is trained to predict the output time series by minimizing the Huber loss [4].

The architecture, though similar to Figure 3, is shallower (as the dimension of numeric input is much smaller than the images), and uses 1D convolutional layers of kernel size of 5×5 , stride 2 and padding 2. All layers are followed by ReLU activation and batch normalization. The encoder part consists of a series of 2 convolutional layers (of $T/2$ and $T/4$ filters, where T is the length of the signal) and a fully connected layer, which gives us a latent representation of embedding length $T/4$. The decoder is a mirrored encoder.

We use a batch size of 128, along with a learning rate of 0.01 which is decayed by a factor of 0.1 after every 5 consecutive non-improving validation epochs. We also utilize early stopping, as described earlier.

4.3 ARIMA

Autoregressive Integrated Moving Average (ARIMA) models are a class of methods that are designed to capture autocorrelations

in the data using a combination approach of autoregressive model, moving average model, and differencing (e.g., [8]). The purpose of each of these three features is to make the model fit the data as well as possible. We used auto arima from pmdarima² library in our experiments.

4.4 RandomWalk

We used the random walk without drift model as a naive numeric forecasting baseline for comparison (e.g., [7]). Specifically, this model assumes that the first difference of the time series data is not time-dependent, and follows a Gaussian distribution $\mathcal{N}(0, \sigma)$. Given a numeric input time series $\{s_0, \dots, s_{t-1}, s_t\}$, in order to predict $\{s_{t+1}, \dots, s_{t+n}\}$, we first estimates σ as

$$\sigma = \sqrt{\frac{t}{t-1} \mathbb{E} [(s_i - s_{i-1})^2]}$$

and the prediction at future time $t + k$ follows

$$s_{t+k} \sim \mathcal{N}(s_t, \sqrt{k}\sigma).$$

This results in a naive numeric forecast that simply extrapolates the last observed value into the future. If we wish to obtain the corresponding image, this forecast is accompanied with a growing uncertainty cone obtained through the equation above.

5 DETAILED EQUATIONS OF NUMERIC MEASURES

SMAPE. The Symmetric Mean Absolute Percentage Error (or SMAPE) is a widely used measure of forecast accuracy [5]. It is calculated as:

$$\frac{1}{T} \sum_{i=1}^T \frac{|y_i - \hat{y}_i|}{(|y_i| + |\hat{y}_i|)/2}$$

where \hat{y}_i is the forecast, y_i the corresponding observed ground-truth, and t is the length of the time series. It ranges from 0.0 to 2.0, with lower values indicating better forecasts.

MASE. The Mean Absolute Scaled Error (or MASE) is another commonly used measure of forecast accuracy [5]. It is the mean absolute error of a forecast divided by the mean absolute first-order difference of actuals, calculated as

$$\frac{|e_j|}{\frac{1}{T-1} \sum_{i=2}^T |y_i - y_{i-1}|}$$

where the numerator $|e_j|$ is the mean absolute error of the forecast, and the denominator is the mean absolute first-order difference over the ground-truth data period T . Errors less than 1 imply that the forecast performs better than the naive one-step method, with lower values indicating better predictions.

6 DETAILED DATASET-BREAKDOWN OF FORECASTING RESULTS

Harmonic Data: The Harmonic dataset are dominated by cyclic patterns. As explained in section 1, each time series in the Harmonic dataset is a mixture (superposition) of two randomly generated individual sinusoids, and each sinusoid exhibits short cyclic patterns

²<http://alkaline-ml.com/pmdarima/>

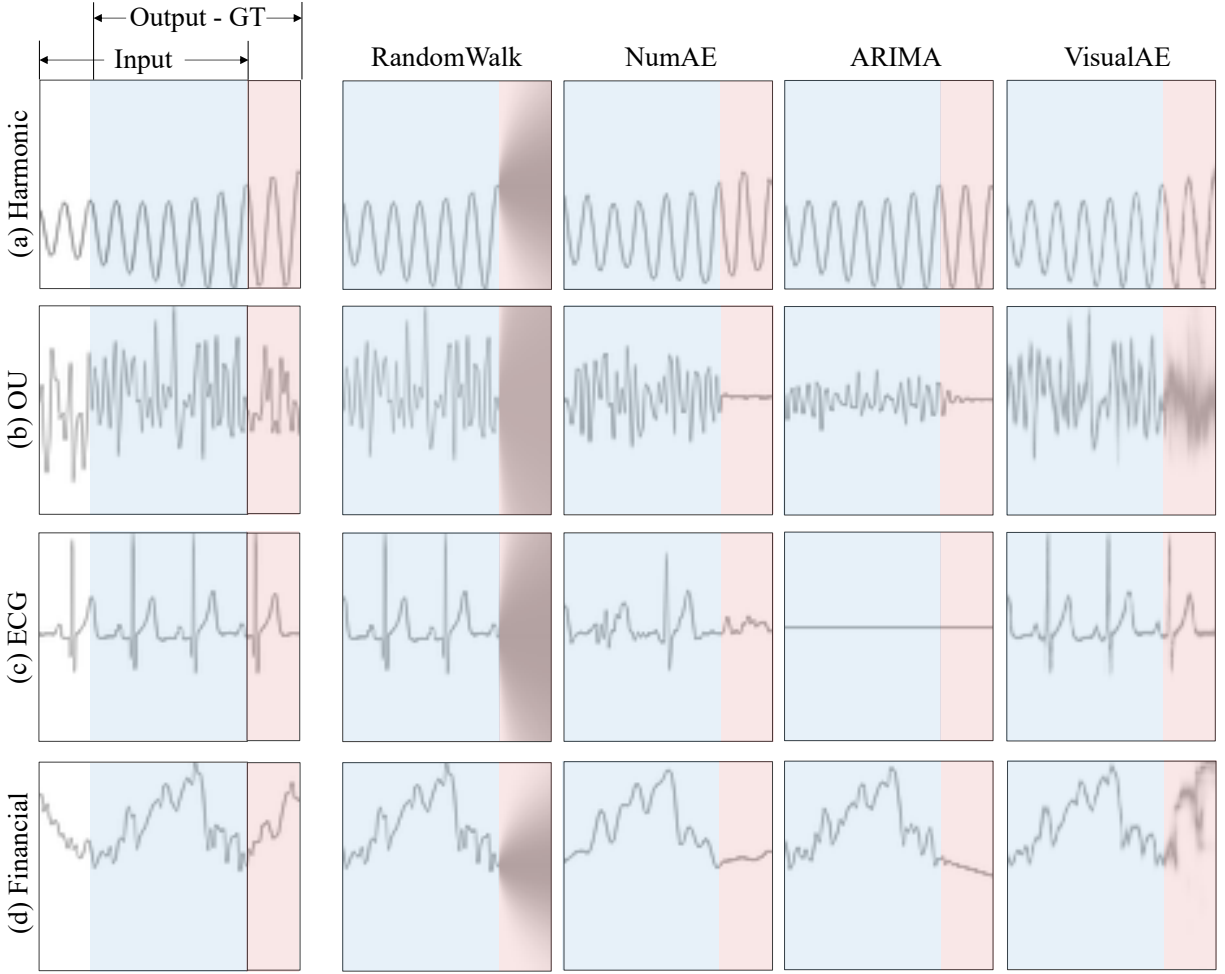


Figure 4: Example out-of-sample forecast predictions using the baseline methods RandomWalk, NumAE, ARIMA, and the proposed method VisualAE. The blue region indicates overlap between the input and output, whereas the red area denotes the future forecast. We show the reconstructed (or fitted) and forecast time series in blue and red region respectively.

along with long term damping or magnifying trends. **ARIMA** performs well but not best because on many occasions, individual sequences don't span the full range of variability needed to tune the model's parameters. This is a barrier for **ARIMA** as it doesn't have the capability to cross-learn between multiple independent time-series. The naive forecasting baseline **RandomWalk** can only learn time-independent stepwise value changes, thus cannot model cyclic patterns. **NumAE** and **VisualAE** capture these patterns well, as observed in Figure 4(a) and Table 1 in the main paper, with **NumAE** performing slightly better according to numeric metrics, and **VisualAE** taking the lead in image-based metrics.

OU Data: It is hard to predict the exact daily changes in the OU data owing to the random process' nature. However, over a larger scale, the OU data is predictable as a mean-reverting process. Visually, we expect majority of the values to concentrate around the mean value of the time series with some noise. **RandomWalk** extrapolates the last observed value, whereas **ARIMA** extrapolates

the in-sample mean value as a steady line for each sample; **NumAE** predicts a slightly jagged version of the same. This becomes evident in the abnormally large MASE error, which is sensitive to division by a small term $\epsilon = 1e - 4$ when the naive one-step denominator approaches 0. The SMAPE metric appears to be similarly non-informative, as it cannot disambiguate the performance of the four methods. The intricacies of these forecasts are captured in the IoU and JSD metrics, according to which **VisualAE** performs the best. This is evident in Figure 4(b), where we show that **VisualAE** concentrates on the hidden mean value, and was also able to partially recover the range of the noise – unlike the other baselines.

ECG Data: ECG time series are periodical with intermittent spikes, and hence inherently predictable. They have relatively constant frequency and do not possess much time dependent uncertainty. Figure 4(c) shows that **VisualAE** is able to capture these cyclic patterns well, as evidenced by all metrics – image-based and numeric. **VisualAE** is able to handle data with sharp and abrupt

changes, and better recovers the heart beat spikes as compared to **NumAE**. Similar to the Harmonic dataset, **ARIMA** is unable to capture the spiky patterns in ECG dataset, and **RandomWalk** simply extrapolates the last value.

Financial Data: Financial time series are the most challenging to forecast amongst the four datasets. According to the prevailing literature (e.g., [6]), financial data is close to random on short scales and shows no apparent periodicity on large scales. Figure 4(d) shows that similar to the OU predictions, **NumAE** and **ARIMA** predicted the future with a weak linear trend, while **VisualAE** outperformed with a predicted curve that captures some of the finer details along with the overall nonlinear trend. This is captured by the IoU metric, but if judged according to SMAPE and MASE, **RandomWalk** would be the best-performer, tied with **NumAE**. This is rather concerning, as solely using numeric metrics would lead us to misleading conclusions, further demonstrating the benefit of using a visual approach in conjunction with traditional numeric methods.

7 FURTHER ANALYSIS ON ADVANTAGE OF IOU METRIC

We can see from Figure 3 in the main paper that the IoU metric better captures how well the visual shape of the predicted time series matches the ground truth. To carry the discussion from the main paper to more depth, Figure 3 in the main paper (i) shows an example for the Harmonic dataset, where according to MASE and SMAPE, the **NumAE** forecast (column b) is the best performer.

This is disputed by IoU, according to which the **VisualAE** forecast (column d) is better, and a qualitative visual inspection also suggests the same. Similarly, (ii) in the same figure shows a hard-to-predict example of the Financial dataset. Just looking at MASE and SMAPE metrics would suggest that both **NumAE** and **VisualAE** forecasts are of similar quality, whereas a visual inspection shows that **VisualAE** captures that long-term trend whereas **NumAE** absolutely does not. Once again, the IoU measure captures this difference, reinforcing our belief that a two-pronged approach of utilizing both numeric and visual approaches holds immense value for the field of time series forecasting.

REFERENCES

- [1] David Byrd. 2019. Explaining agent-based financial market simulation. *arXiv preprint arXiv:1909.11650* (2019).
- [2] Bilal Fadlallah, Badong Chen, Andreas Keil, and José Principe. 2013. Weighted-permutation entropy: A complexity measure for time series incorporating amplitude information. *Physical Review E* 87, 2 (2013), 022911.
- [3] Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. 2000. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *circulation* 101, 23 (2000), e215–e220.
- [4] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. 2001. *The Elements of Statistical Learning*. Springer New York Inc., New York, NY, USA.
- [5] Spyros Makridakis and Michele Hibon. 2000. The M3-Competition: results, conclusions and implications. *International journal of forecasting* 16, 4 (2000), 451–476.
- [6] Lasse Heje Pedersen. 2019. *Efficiently inefficient: how smart money invests and market prices are determined*. Princeton University Press.
- [7] Steven E Shreve. 2004. *Stochastic calculus for finance II: Continuous-time models*. Vol. 11. Springer Science & Business Media.
- [8] Daniel S Wilks. 2011. *Statistical methods in the atmospheric sciences*. Vol. 100. Academic press.