# Financial Time Series Representation Learning

**Philippe Chatigny**[1*] , **Jean-Marc Patenaude**[2] , **Shengrui Wang**[1]

[1]Department of Computer Science, Université de Sherbrooke, QC, Canada
[2]Laplace Insights, Sherbrooke, QC, Canada

{Philippe.Chatigny, Shengrui.Wang}@usherbrooke.ca, JeanMarc@LaplaceInsights.com

## Abstract

This paper addresses the difficulty of forecasting multiple financial time series (TS) conjointly using deep neural networks (DNN). We investigate whether DNN-based models could forecast these TS more efficiently by learning their representation directly. To this end, we make use of the dynamic factor graph (DFG) from that we enhance by proposing a novel variable-length attention-based mechanism to render it memory-augmented. Using this mechanism, we propose an unsupervised DNN architecture for multivariate TS forecasting that allows to learn and take advantage of the relationships between these TS. We test our model on two datasets covering 19 years of investment funds activities. Our experimental results show that our proposed approach outperforms significantly typical DNN-based and statistical models at forecasting their 21-day price trajectory.

## 1 Introduction

In recent decades, DNN has helped improved TS forecast accuracy in various social settings [Makridakis *et al.*, 2019]. Besides their ability to handle non-linear processes, they provide a cost-effective approach to uncovering relations between TS automatically. They do so by enforcing a hierarchical structure for pattern detectors throughout its hidden layers, from which, by using sensitivity analysis methods [Ribeiro *et al.*, 2016], we can determine the factors that influence the forecasted trajectory. DNN has permitted increasing accuracy on mostly homogeneous datasets with multiple measurements as well as in applications where there exist exogenous variables that are strongly related to the variable(s) of interest; e.g.: traffic or electricity load forecasting [Hyndman and Athanasopoulos, 2018]. However, training a DNN remains difficult for settings in which the TS are non-ergodic, heteroskedastic, non-stationary, or with high noise to signal ratios. Such cases are often found in financial TS. Few DNN-based models have demonstrated consistent accuracy on datasets spanning over multiple years for different asset classes [Sezer *et al.*, 2019].

A major reason for the difficulties in forecasting financial TS is that most DNN learning frameworks do not appear to be adapted for this setting despite the large quantity of TS available. Training a DNN needs a sufficiently large dataset of independent training samples that are representative of the data to infer. At the exception of applications such as intra-day forecasting [Qin *et al.*, 2017], most financial applications rely on TS that have a relatively limited number of measurements [Makridakis *et al.*, 2018b]. Additionally, historical price trajectories can be very noisy and their behaviors follow a more complex cyclical effect than that found in intra-day data: the market cycle [Hamilton and Lin, 1996][1]. As it is not possible to obtain multiple independent realizations of a specific asset's price fluctuation under different circumstances for the same time period [Marshall, 2009], the nature of financial TS enforce both a deprivation on the amount of training data and the well-known difficulty of modeling these long-term effects [Bengio *et al.*, 1994].

The aim of this paper is to propose a more efficient DNN framework for forecasting multiple financial assets conjointly. The key contributions of this paper are as follows:

(1) We propose a novel attention mechanism for the Dynamic Factor Graph (DFG) framework. This mechanism offers the capability of considering a variable number of past latent states. We make use of this mechanism to vary the order of an autoregressive (AR) generative function over time.

(2) When integrated into a spatiotemporal neural network (NN) model, we provide an unsupervised deep generative approach, to model interactions between multiple TS for multivariate TS forecasting. Our spatiotemporal adaptive neural network (STANN) has the particularity to allows the discovery of the interrelations between each TS if they are not given as prior.

(3) Our experimental evaluation shows that the proposed model provides a more effective learning framework for the addressed setting of forecasting 21 daily return trajectories of exchanged traded (ETFs) funds and mutual funds (MFs).

The remainder of this paper is organized as follows: Section 2

---

*Contact Author

[1]The market cycle can be expressed by irregular periodic fluctuations of assets prices observed in different market conditions.

reviews major existing work in modeling TS in social settings and relevant models similar to our work. In Section 3, we present our model and describe its training procedure. In Section 4, we present the setup of our empirical evaluation, which extends around 19 years of financial market activities, and describe our results. Section 5 presents our conclusion.

## 2 Related Work

DNN-based approaches, such as the recurrent neural network (RNN), have been extensively used to model TS, given their ability to represent and to forecast sequential data. Various approaches have been introduced to facilitate their application on TS [Rubanova *et al.*, 2019]. While promising results have been achieved recently for financial TS prediction such as [Borovykh *et al.*, 2018], it has been pointed out that much published ML work in the TS literature claims satisfactory accuracy without making an adequate comparison of their proposed against statistical approaches [Makridakis *et al.*, 2018b].

In fact, only a few authors, such as [Rangapuram *et al.*, 2018; Smyl, 2020], have been able to explicitly evaluate that their models yield better performance on multiple TS over simple statistical models like ARIMA or even a naive forecast. Misleading results are often concluded by using non-scaled error metrics that are known to be ambiguous when comparing TS forecasts [Hyndman and Koehler, 2006]. The myriad of proposed DNN-based models [Sezer *et al.*, 2019] applied to financial settings and the results presented around them have raised undue expectations that such methodologies provide accurate predictions on various financial TS applications, while there is clearly a lack of experimental proof that they outperform simple baselines for the majority of them.

Nonetheless, large gains can still be achieved using DNN and ML approaches. Recently, state-of-the-art accuracy was achieved at the M4 competition [Makridakis *et al.*, 2018a] from which the top 2 entries used DNN-based or ML techniques along with statistical models. Subsequently to these findings, [Oreshkin *et al.*, 2020] were the first to show that it was possible to build a pure DNN-based model for this task and achieved greater gains than the best competition entry [Smyl, 2020]. Given the wide range of TS to forecast[2], the top performing models submited relied on ensemble techniques to be robust over the different types of series. The direct comparison between single and ensemble models is generally unfair but the findings from these models can be investigated for building better individual models. For instance, the DNN-based models which performed well on this dataset [Smyl, 2020; Oreshkin *et al.*, 2020] provide insights into which techniques to use for improving the performance of individual models: Residual connections between hidden layers, adaptive learning rate scheduling, input prepossessing and both seasonal and trend decomposition embedded directly in the model.

Most of these techniques are *"tricks"* to facilitate learning DNN. However, the idea of applying a TS decomposition within a neural network is promising and several au-

---

[2]The M4 dataset contained 100'000 individual TS from which approximately 25% were financial TS of different types.
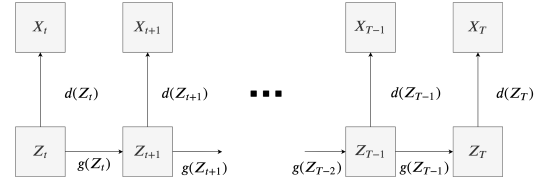


Figure 1: Illustration of DFG's architecture. Both decoder and dynamic modules can be implemented as parametric functions and be trained using gradient descent.

thors [Oreshkin *et al.*, 2020; Godfrey and Gashler, 2017; Hansen and Nelson, 2003] have shown its effectiveness on real-world datasets. Given the well-known difficulty of dealing with the raw signal of financial TS, we raise the question whether we could learn a better representation of these TS directly and apply such decomposition. Fortunately, there exist a DNN learning framework well suited for this task: the Dynamic Factor Graph (DFG) [Piotr and LeCun, 2009]. Illustrated in Fig. 1, DFG consists of a directed acyclic state-space model with continuous latent factors $Z$ that permits modeling the joint probability $P(Z, X)$. Contrary to typical RNNs, DFG learns the representation of a TS ($X$) directly within $Z$ instead of computing hidden states in an iterative fashion using a sequence of past inputs. It does so by using two modules: a decoder module and a dynamic module. The decoder module $d()$ infers values over the latent space at time $t$, $\tilde{X}_t = d(Z_t)$, and the dynamic module $g()$ captures the dynamic process of a TS through the AR function $Z_{t+1} = g(Z_t)$.

The main difference between a RNN and a DFG in its simplest form is that the state space component is used differently. a RNN with self loops on the hidden cells is usually formulated as $Z_{t+1} = g(Z_t, X_t)$ with $X_t$ being the measurement obtained at time $t$. DFG does not require $X_t$ as the dynamics of the series is stored on external memory and captured entirely in the latent space: $Z_{t+1} = g(Z_t)$. Thus, DFG is a particular case of a RNN where the hidden states are directly learned instead of being computed explicitly by a function of past inputs. Like most AR functions, we could specify the AR order of $g()$ by changing its configuration: $Z_{t+1} = g(Z_t, ..., Z_{t-v})$. Note that doing so makes the assumption that the AR parameters $(\varphi_1, ..., \varphi_n)$ are constant. However, this AR order is an additional hyperparameter (HP) to tune and assuming that the AR parameters are constant often impair training as the resulting AR weights are optimized to reduce the average error. This limitation is problematic if the training data contains multiple TS dynamics. In this work, we address these limitations by proposing an attention mechanism that permits DFG to select its AR order automatically and adjust its AR parameters over time.

## 3 The STANN model

### 3.1 Notation and task

Given $X : \mathbb{R}^{T \times n \times m}$, a 3-dimensional tensor representing a set of $n$ TS of length $T$ and dimensionality $m$, we define $X_{t,i,j}$ as the value of dimension $j$ for TS $i$ at time $t$. The task

of interest is to predict $n$ multivariate TS $\tau$ time steps ahead $\tilde{X} : \mathbb{R}^{\tau \times n \times m}$. We represent the spatial relationship between series within a 3-dimensional tensor $W : \mathbb{R}^{n \times R \times n}$ where $R$ is the number of relations considered. Thus, our aim is to train a model $f : \mathbb{R}^{T \times n \times m} + [\mathbb{R}^{n \times R \times n}] \rightarrow \mathbb{R}^{\tau \times n \times m}$.

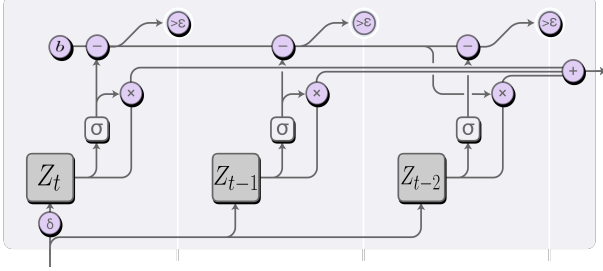## 3.2 Model definition without spatial dependencies



Figure 2: Illustration of the proposed attention mechanism. For illustration purposes, $b$ includes both $b_{\texttt{time}}$ and $b_{\texttt{cost}}$, and $\sigma$ denotes the AR weight produced by $actm(Z_{t-i})$. The drawing was adapted from [Olah and Carter, 2016].

In STANN, we use a deterministic variant of the DFG. The decoder module $d()$ decodes instead the expected variation between $X_{t-1}$ and $X_t$ from the latent factor $Z_t$, i.e. $X_t : \tilde{\Delta}_{X_{t-1}}^{X_t} \leftarrow Z_t$. The decoder is defined as in Eq. 1, where $\tilde{X}_t$ is the prediction computed at time $t$. The dynamical module $g$ is defined by Eq. 2 and considers the past $\upsilon$ relevant latent factors $Z_{t-\upsilon:t}$, i.e. $Z_{t-\upsilon}$ to $Z_t$. $d()$ and $g()$ are implemented as doubly residual stacking NNs like that used in N-BEATS [Oreshkin *et al.*, 2020] in which we use Convolutional Neural Network in its basic block. Contrary to N-BEATS, we apply the TS decomposition on the latent factors instead of the raw signals. As mentioned in the previous section, assuming that the forecast depends on a fixed AR order covering past $\upsilon$ observations is a strong assumption that can impair model training if not selected correctly. One can view that RNNs, like the LSTM, allow $\upsilon$ to vary over time by maintaining in memory a state vector that allows them to retain information as long as required and forget it when it is no longer relevant [Hochreiter and Schmidhuber, 1997].

$$\tilde{X}_t = X_{t-1} + d(Z_t) \quad (1) \quad \tilde{Z}_{t+1} = g(Z_{t-\upsilon:t}) \quad (2)$$

Alike LSTM, we improve upon DFG by updating $\upsilon$ adaptively such that Eq. 2 is done using a variable $\upsilon$ of past steps. To do so, we propose an adaptive attention-based mechanism to enable DFG to be *memory-augmented*. Our attention mechanism is inspired by the *Adaptive Computation Time* (ACT) algorithm proposed in [Graves, 2016], implemented as $actm()$ in our model. $actm()$ generates a probability distribution on $Z$ from which we estimate the AR weights for each latent factor: $\varphi_{Z_t} = \sigma(W_\varphi Z_t + b_\varphi)$ with associated weight matrix $W_\varphi$ and bias $b_\varphi$. $Z_{t-\upsilon:t}$ is computed using the sums of all past latent factors weighted by their autoregressive weights as in Eq. 3. $actm()$ uses two budgets $b(t) = \{b_{\texttt{time}} = t, b_{\texttt{cost}} = 1\}$: one to keep account of

available past time steps and one to track the cost of considering a latent factor. Each time we consider a latent factor $Z_t$, we reduce our budget $b_{\texttt{time}}$ by 1 and $b_{\texttt{cost}}$ by $\varphi_{Z_t}$ which is bounded within $]0, 1[$. If a budget goes below $\epsilon$ (either $b_{\texttt{cost}} < \kappa$ or $b_{\texttt{time}} = 0$), we stop considering any more latent factor and attribute the remaining cost budget to the last factor considered. $\kappa$ is a small constant (0.01 for the experiments in this paper), whose purpose is to allow the selection of an AR(1) process.

$$Z_{t-\upsilon:t} = actm(Z_t) = \sum_{0 \leq k < \upsilon : [b(t) > \epsilon]} \varphi_{Z_{t-k}} Z_{t-k} \quad (3)$$

We can interpret $actm()$'s objective as evaluating the quality of each past latent factor and assigning the appropriate autoregressive weight at times $t$ that maximizes the log likelihood of the generative process modeled by Eq. 2. Since $actm()$ uses $b_{\texttt{time}}$ to determine how many past steps are available, $actm()$ can theoretically account all previous learned factors if $\sum_{k=0}^{T} P(\sigma_{t-k}|Z_{t-k}) > \epsilon$. Note that the imposed budget restricts each autoregressive weight to be between 0 and 1 with the sums of all the weights being equal to 1. We apply this mechanism solely within $g()$ to facilitate the training of both components. The attention mechanism is summarized in Figure 2 and can be designed as any configuration of a feedforward network with a sigmoid activation function.

The training procedure consists of minimizing this bi-objectives loss function (4):

$$\mathcal{L}(d, g, Z) = \frac{1}{T} \Delta(X_{t-1} + d(Z_t), X_t) \quad (4a)$$

$$+ \frac{1}{T} \sum_{t=1}^{T-1} ||Z_{t+1} - g(Z_{t-\upsilon:t})||^2 \quad (4b)$$

The first term (4a) measures the ability of the model to reconstruct $X_t$ from $Z_t$. The second term (4b) measures the ability of the system to capture the dynamicity of the equation by its ability to link states of $Z$ in sequential order. $\Delta$ is a loss function that measures the difference between the prediction $\tilde{X}_t$ and the ground truth $X_t$. In this work, we choose to train our model from end-to-end [Tesauro, 1995] instead of using the proposed expectation-minimisation based approach [Piotr and LeCun, 2009] since we have observed that it produces better results.

## 3.3 Model Definition Including Spatial Dependencies

Let us now introduce the way interrelations between TS are captured. As pointed out in [Schwendener, 2010], multiple types of relations between financial time TS has been uncovered. One must be able to test whether this prior knowledge has predictive capability. Thus, we propose that the relationships between the dynamic processes of multiple TS be given as additional prior inputs $W \in \mathbb{R}^{n \times R \times n}$ to the model like in [Ziat *et al.*, 2017]. We will first formalize how relationships between series are incorporated into the model and how this

allows us to "virtually" have a high number of training samples. Then, we will describe two extensions of this approach. The first extension allows to weight the strength of these relations and the second allows the model to learn these relations directly without any prior information.

Relationships between the dynamic processes of $n$ TS are incorporated via a tensor $W \in \mathbb{R}_+^{n \times R \times n}$, where $R$ is the number of relation types given as prior. We formulate that $Z_{t+1,i}$ depends on its own latent representation at time $t$ (intradependency) and on the representations of other series at time $t$ (interdependency). Intradependency is modeled through a linear mapping $\Theta^{(0)} \in \mathbb{R}^{n \times n}$ and interdependency is modeled by $R$ transition matrices $\Theta^{(R)} \in \mathbb{R}^{R \times n \times n}$. Thus, to evaluate $Z_{t+1}$, we compute the matrix product between the latent space $Z_t$ and its dependencies $(\Theta^{(0)}, \Theta^{(R)})$ as in Eq. 5. The decoder, follows along by using $Z_{t,i}$, as inputs, and computes the expected variation as in Eq. 6. $h_g, h_d$ are the respective activation functions of $g()$ and $d()$.

$$Z_{t+1,i} = g(Z_{t-v:t}, i) =$$
$$h_g(actm_g(Z_t \theta_i^{(0)} + \sum_{r \in R} W_i^{(r)} Z_t \theta_i^{(r)})) \quad (5)$$

$$\tilde{\Delta}_{X_{t+1,i}}^{X_{t,i}} = d(Z_{t,i}) = h_d(actm_d(Z_{t,i})) \quad (6)$$

Note that $Z_t$ is shared between all series with respect to $g()$, but the representation of each series is disentangled explicitly by means of $W$, i.e.: $d()$ takes as input $Z_{t,i}$, the hidden factor of the $i^{\text{th}}$ TS. Doing so has two advantages: (1) $g()$ can forecast $\tilde{Z}_{t+1}$ with fewer regressors. (2) It "virtually" increases the amount of training samples as we can use time and positional coordinates to make $T \times n$ training samples.

### 3.4 Model extensions

The two possible extensions proposed in [Ziat *et al.*, 2017] can also be applied to our model. We summarize the extensions here; readers are invited to refer to the original paper [Ziat *et al.*, 2017] for a more detailed explanation. The first extension, denoted by STANN-R, consists of adding a learned matrix of weights $\Gamma^r \in \mathbb{R}_+^{n \times n}$ that can reduce the strength of relations given as prior. The second extension, denoted by STANN-D, consists of replacing $W$ with $\Gamma$ such that the model learns both the relational structure and the relation weights via $\Gamma$. Applying the STANN-R or STANN-D extension formalizes Eq. 5 as in Eq. 7 or Eq. 8, respectively, where $\odot$ signifies element-wise multiplication between two matrices:

$$Z_{t+1,i} = g(Z_{t-v:t,i}) = g(actm_g(Z_t \theta_i^{(0)}$$
$$+ \sum_{r \in R} (\Gamma_i^{(r)} \odot W_i^{(r)}) Z_t \theta_i^{(r)}) \quad (7)$$

$$Z_{t+1,i} = g(Z_{t-v:t,i}) = g(actm_g(Z_t \theta_i^{(0)} + \sum_{r \in R} \Gamma_i^{(r)} \theta^{(r)}))$$
$$\quad (8)$$

The optimization problem can thus be adjusted for $\Gamma$, depending on whether the dynamic function is specified by Eq. 7 or Eq. 8, and can be written as Eq. 9. $|\Gamma|$ is a $l_1$ regularizing term intended to sparsify $\Gamma^{(r)}$, and $\gamma$ is a hyperparameter set to tune this term and $\lambda$ is a factor set to balance the importance between $g()$ and $d()$.

$$d^* g^*, actm_g^*, \Gamma^* =$$
$$\operatorname*{argmin}_{d,Z,\Gamma,} \frac{1}{T} \sum_t \Delta(d(Z_t) + X_{t-1}, X_t)$$
$$+ \gamma |\Gamma| + \lambda \frac{1}{T} \sum_{t=1}^{T-1} ||Z_{t+1} - g(Z_{t-v:t})||^2 \quad (9)$$

## 4 Experiments

### 4.1 Datasets and experimentation procedure

We report here the results of experimental evaluation of our forecasting methods on two datasets: Fasttrack and Fasttrack extended. Both datasets contain daily closing prices of US MFs and ETFs traded on US financial markets each covering different types of asset classes including stocks, bonds, commodities, currencies and market indexes, or a proxy for a market index. When both are combined, they cover 19 years of financial market activities and provide an overall view of the whole financial ecosystem. Each TS of these datasets represents the aggregation of multiple individual financial assets. In some of these TS, like VFICX, the aggregation of these individual TS are subject to vary over time with respect to management activities associated with these funds.

The two datasets used are summarized in Table 1 and were obtained through FastTrack[3], a professional grade investment data provider. We used the adjusted closing price to simulate how each model would have performed in a real-case scenario. Fasttrack includes the following funds: SPY, EWJ, VSMGX, FNMIX, VEXMX, VFITX, VFICX, DXY-Z, VBISX, VUSTX. Fasttrack extended included the following funds: DBC, DIA, EFA, EWA, EWC, EWD, EWG, EWH, EWI, EWJ, EWK, EWL, EWM, EWN, EWO, EWP, EWQ, EWS, EWU, EWW, FCYIX, FNARX, FNMIX, FSCPX, FSDCX, FSLEX, FSNGX, GLD, IEF, LQD, MDY, PCY, QQQ, SHY, SPY, TIP, TLT, USO, UUP, VASVX, VAW, VBISX, VCR, VDC, VDE, VEIEX, VFH, VFICX, VFISX, VFITX, VGSIX, VGT, VHT, VINEX, VIS, VNQ, VOX, VPU, VTSMX, VWO, VXF, XLB, XLE, XLF, XLI, XLK, XLP, XLU, XLV. Both datasets are proprietary, which we do not have the permission to share publicly. However, for the sake of reproducibility, we enumerated the tickers used for our experiments to help interested readers reconstruct the datasets from public data sources.

To evaluate the average forecasting performance of our model, we use the following metrics: The Mean Absolute Scaled Error (MASE) [Hyndman and Koehler, 2006], the Theil U2 score (TheilU) [Theil, 1971], the sDILATE presented in [Vincent and Thome, 2019] and the Mean directional accuracy (MDA) score [Schnader and Stekler, 1990].

---

[3]https://investorsfasttrack.com

Table 1: Datasets for experimental evaluation

| Dataset | $T$ | $n$ | Data type | Time horizon | $\tau$ | # Runs per model |
|---|---|---|---|---|---|---|
| Fasttrack | 2186 | 10 | daily adj. close | 1996/07/08-2007/08/22 | 21 | 100 |
| Fasttrack extended | 2000 | 69 | daily adj. close | 2011/05/31 -2019/05/10 | 21 | 54 |

$T$ is the total number of time points, $n$ the number of series, $\tau$ the number of steps ahead to forecast and # Runs the total number of evaluation runs made. For all datasets, we considered only the closing price $m = 1$.

Table 2: Average forecasting performance of tested models on Fasttrack and Fasttrack extended datasets

| | FAST TRACK | | | | FAST TRACK EXTENDED | | | |
|---|---|---|---|---|---|---|---|---|
| **Model:** | *MASE* | *THEILU* | *sDILATE* | *MDA* | *MASE* | *THEILU* | *sDILATE* | *MDA* |
| Naive | $1.0000 \pm 0.0000$ | $1.0000 \pm 0.0000$ | $1.0000 \pm 0.0000$ | $0.0180 \pm 0.0149^{****}$ | $1.0000 \pm 0.0000$ | $1.0000 \pm 0.0000$ | $1.0000 \pm 0.0000$ | $0.0128 \pm 0.0082^{****}$ |
| AR | $1.0707 \pm 0.1517^{****}$ | $1.0757 \pm 0.1577^{****}$ | $1.1819 \pm 0.3560^{****}$ | $0.5085 \pm 0.1469^{****}$ | $1.0337 \pm 0.0844^{****}$ | $1.0306 \pm 0.1033^{***}$ | $1.0723 \pm 0.2109^{***}$ | $0.4788 \pm 0.0955^{*}$ |
| ARIMA | $1.0030 \pm 0.1205^{t}$ | $1.0133 \pm 0.1204^{t}$ | $1.0412 \pm 0.2457^{*}$ | $\mathbf{0.5817 \pm 0.1834}$ | $1.0011 \pm 0.0945^{*}$ | $1.0008 \pm 0.1193^{t}$ | $1.0156 \pm 0.2373^{t}$ | $0.2748 \pm 0.1222^{****}$ |
| LSTM | $1.3399 \pm 0.6020^{****}$ | $1.3405 \pm 0.6332^{****}$ | $2.1941 \pm 2.5503^{****}$ | $0.4861 \pm 0.1624^{****}$ | $1.2543 \pm 0.3020^{****}$ | $1.2311 \pm 0.3002^{****}$ | $1.6041 \pm 0.8031^{****}$ | $0.4821 \pm 0.1426^{t}$ |
| WaveNet | $1.5936 \pm 0.7655^{****}$ | $1.6093 \pm 0.8133^{****}$ | $3.2449 \pm 3.7111^{****}$ | $0.4844 \pm 0.1606^{****}$ | $1.3988 \pm 0.5445^{****}$ | $1.4071 \pm 0.5930^{****}$ | $2.3042 \pm 2.4721^{****}$ | $0.4864 \pm 0.1472^{*}$ |
| STNN | $0.9852 \pm 0.0693$ | $0.9920 \pm 0.0756$ | $0.9897 \pm 0.1484$ | $\underline{0.5942 \pm 0.1816}$ | $1.0020 \pm 0.1536$ | $0.9959 \pm 0.1591$ | $1.0165 \pm 0.3354$ | $0.5259 \pm 0.1822$ |
| STNN-R | $0.9860 \pm 0.0785$ | $0.9900 \pm 0.0791^{*}$ | $0.9863 \pm 0.1431^{*}$ | $0.5450 \pm 0.1965^{*}$ | $1.0122 \pm 0.1707$ | $1.0047 \pm 0.1698$ | $1.0369 \pm 0.3687$ | $0.5241 \pm 0.1693$ |
| STNN-D | $1.0812 \pm 0.2957^{***}$ | $1.0808 \pm 0.2765^{**}$ | $1.2439 \pm 0.8131^{**}$ | $0.5585 \pm 0.1533^{t}$ | $0.9814 \pm 0.1147$ | $0.9791 \pm 0.1255$ | $0.9743 \pm 0.2495$ | $0.5401 \pm 0.2052$ |
| STANN | $\underline{\mathbf{0.9792 \pm 0.1045}}$ | $\underline{\mathbf{0.9828 \pm 0.1114}}$ | $0.9783 \pm 0.2174$ | $0.5363 \pm 0.1914$ | $0.9832 \pm 0.1023$ | $0.9814 \pm 0.1084$ | $0.9750 \pm 0.2148$ | $0.5360 \pm 0.2030$ |
| STANN-R | $\mathbf{0.9806 \pm 0.0784}$ | $\mathbf{0.9863 \pm 0.0804}$ | $0.9793 \pm 0.1562$ | $0.5864 \pm 0.1873$ | $0.9836 \pm 0.1026$ | $0.9816 \pm 0.1098$ | $0.9755 \pm 0.2189$ | $0.5401 \pm 0.2051$ |
| STANN-D | $\mathbf{0.9864 \pm 0.0381}$ | $\mathbf{0.9870 \pm 0.0374}$ | $\underline{\mathbf{0.9756 \pm 0.0707}}$ | $0.5642 \pm 0.1956^{t}$ | $\underline{\mathbf{0.9795 \pm 0.1016}}$ | $\underline{\mathbf{0.9785 \pm 0.1096}}$ | $\underline{\mathbf{0.9694 \pm 0.2176}}$ | $\mathbf{0.5406 \pm 0.2055}$ |

Averaged forecasting results of the 21 days multivariate trajectory forecasts for both datasets. We highlight the bests methods in bold by using the Wilcoxon signed-rank test with significance level of $p - value < 0.10$. We also indicate the statistical significance between the best performing model in regards to their associate metrics (t: P$\leq$ 0.10; *: P$\leq$ 0.05; **: P$\leq$0.01; ***: P$\leq$0.001; ****: P$\leq$0.0001). We underline the best performing model used for comparing the significance level on all metrics.

TheilU is the equivalent of MASE that use the RMSE instead for comparing the difference between predicted value and the ground truth. Both MASE and TheilU account for the bias and variance of error residual with TheilU penalizing more large errors than MASE. [Vincent and Thome, 2019] argued that such error measures have limitations since we can obtain the same loss value for completely different forecasts. Therefore, we use a scaled version of the DILATE loss. sDILATE gives more importance to the shape of the forecast than the distribution of the error residual. Finally, we use a modified version of the MDA for measuring whether the direction of the forecast is the same as the price trajectory in comparison to the last known value of each TS, i.e.:
$$MDA(\tilde{X}_{t:t+\tau}, X_{t:t+\tau}) = \frac{1}{N} \sum_{i=0}^{\tau} sign(\tilde{X}_{t:t+i} - X_{t-1}) = sign(X_{t:t+i} - X_{t-1}).$$

To train each model, we carried out an evaluation on a rolling forecasting origin and by setting $\tau$, i.e. the number of time steps, to 21 days for simulating forecasting on a monthly basis. All models were trained on normalized TS using the interquartile range method. Produced forecasts were unscaled back to the original TS scales to measure forecast's error. All DNN-based models were trained using stochastic gradient descent SGD with Adam [Kingma and Ba, 2015] and a learning rate scheduler [Loshchilov and Hutter, 2017]. The number of steps for SGD and other model parameters, like the optimal training window, were determined by a hyperparameter search.

### 4.2 Baseline models

As a limit on the scope of our evaluation, we considered only models that can forecast multivariate TS directly, with the exception of two baseline models. We performed experiments with the following models:

1. **Naive:** A simple heuristic that assumes that the $\tau$ future steps will be the same as last previously observed.

2. **AR:** A classical univariate autoregressive process in which each TS is forecasted individually. The prediction is a linear function of past lags $l$.

3. **ARIMA:** An autoregressive integrated moving average model that forecasts each TS individually. Implementation of ARIMA was done with [Taylor G Smith, 2019] to automatize the selection of the best parameters over the training set.

4. **LSTM:** A long short-term memory model which forecasts $\tau$ steps ahead in an iterative fashion.

5. **WaveNet:** A convolutional neural network using dilated causal convolutions [Van Den Oord et al., 2016].

6. **STNN:** The closest model to ours. STNN can be considered as a particular case of our model, i.e. model with $v = 0$. The two extensions of STNN (STNN-R and STNN-D) [Ziat et al., 2017] are also considered. The Pearson correlations between TS was computed over the training set to define $W$. We use the same training strategy as STANN, i.e. modeling the variation only and trainning the model from end-to-end to establish a fair comparison between model architectures.

7. **STANN:** The model proposed in this paper. The two extensions presented in Section 3.4 are also considered. The extensions expressed in Eq. 7 and Eq. 8 are denoted by STANN-R and STANN-D, respectively. Pearson correlation was used to define $W$.

Comparison between STANN-D and STNN-D forecast error
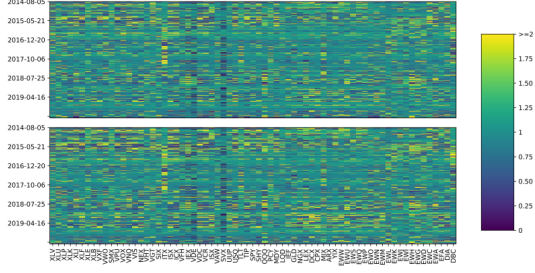for `Fasttrack extended` dataset



Figure 3: Concatenation of the 21 daily return forecasts of STANN-D (top) and STNN-D (bottom). The absolute scaled error per series is presented.

### 4.3 Results

Our experimental results are summarized in Table 2. First, we analyze the average performance of all models and the statistical significance of the results obtained. Our model outperform all DNN-based and statistical baselines considered in regards to all metrics on both datasets. This result suggests the superiority of the unsupervised training framework of the DFG for modeling these TS conjointly. The addition of the proposed attention mechanism and the TS decomposition appear to slightly improve the performance over its based model (STNN) but we can not claim a significant statistical difference for the accuracy improvement of forecasting these trajectories. Nonetheless, these results are very promising when considering that (1) our approach achieved such results by using a relatively few amount of TS and that (2) it was trained solely by using historical prices.

We can qualitatively compare our models by plotting the absolute scaled error of the individual point forecast (IPF) for all the TS forecasted and comparing where our model fails. We observe that our approach is relatively consistent at forecasting the trajectory of each assets (Fig. 3), but the majority of the residuals appears to occurs in epistemic fashion, i.e.: the TS' forecasting difficulty varies over time. Our proposed attention mechanism increases slightly the forecast accuracy in these episodes of forecast instability over its based model, which explains the majority of the additional average gain in accuracy. Next, we performed an ablation study by comparing what effect the TS decomposition technique and the attention mechanism have to help obtain better forecast. We do so by plotting the probability distribution of each step-ahead forecast of STNN-D and STANN-D along with versions of our model where one of the two components is missing. Interestingly, adding soley the attention mechanism increase the overall error but reduce the error propagation often found in recursive approach. When combined with the TS decomposition architecture presented in [Oreshkin *et al.*, 2020], we observe a significant error reduction for the last 14 days of the forecast trajectory. Independently from the step-ahead forecasted, our approach is significantly better than its based model at reducing IPF ****.

Individual step-ahea MASE distribution of our approach on the `Fasttrack extended` dataset.
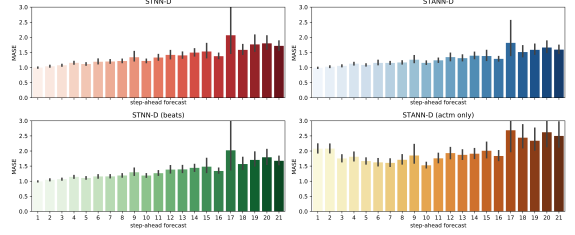


Figure 4: An ablation studies of STANN-D is presented. The median absolute scaled error for STANN-D is equal to 0.9671 and for STNN-D is equal to 0.9810

## 5 Discussion

This work proposed a new unsupervised deep generative model (STANN) for forecasting multivariate TS conjointly which explicitly models the interactions between TS. Experiments were performed on two financial datasets covering over 19 years of market history. Our experiments indicate that STANN provides a more effective learning framework than both DNN-based approaches and statistical baselines. We showed that this class of models perform wells in a low-data setting and that our proposed attention mechanism helps improve forecasting performances over its based model. Finally, we showed that having few training TS samples is not an issue for training this class of models as they permits to "virtualy" increase its amount of training sample.

We would like to emphasize on the limited understanding of its effectiveness in relation to the selection of HPs. Indeed, a mis-selection of HPs can have a large impact on the model's performance which can render difficult its application at large scale. Hence, we advocate the pursue of future works to enlarge our theoretical understanding on this class of models as well as testing if similar results can be achieved at larger scale.

## References

[Bengio *et al.*, 1994] Yoshua Bengio, Patrice Simard, Paolo Frasconi, et al. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166, 1994.

[Borovykh *et al.*, 2018] Anastasia Borovykh, Sander Bohte, and Cornelis W Oosterlee. Dilated convolutional neural networks for time series forecasting. *Journal of Computational Finance, Forthcoming*, 2018.

[Godfrey and Gashler, 2017] Luke B Godfrey and Michael S Gashler. Neural decomposition of time-series data for effective generalization. *IEEE transactions on neural networks and learning systems*, 29(7):2973–2985, 2017.

[Graves, 2016] Alex Graves. Adaptive computation time for recurrent neural networks. *arXiv preprint arXiv:1603.08983*, 2016.

[Hamilton and Lin, 1996] James D Hamilton and Gang Lin. Stock market volatility and the business cycle. *Journal of applied econometrics*, 11(5):573–593, 1996.

[Hansen and Nelson, 2003] JV Hansen and RD Nelson. Forecasting and recombining time-series components by using neural networks. *Journal of the Operational Research Society*, 54(3):307–317, 2003.

[Hochreiter and Schmidhuber, 1997] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[Hyndman and Athanasopoulos, 2018] Rob J Hyndman and George Athanasopoulos. *Forecasting: principles and practice*. OTexts, 2018.

[Hyndman and Koehler, 2006] Rob J Hyndman and Anne B Koehler. Another look at measures of forecast accuracy. *International journal of forecasting*, 22(4):679–688, 2006.

[Kingma and Ba, 2015] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.

[Loshchilov and Hutter, 2017] I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations (ICLR) 2017 Conference Track*, April 2017.

[Makridakis *et al.*, 2018a] Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. The m4 competition: Results, findings, conclusion and way forward. *International Journal of Forecasting*, 34(4):802–808, 2018.

[Makridakis *et al.*, 2018b] Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. Statistical and machine learning forecasting methods: Concerns and ways forward. *PloS one*, 13(3):e0194889, 2018.

[Makridakis *et al.*, 2019] Spyros Makridakis, Rob J Hyndman, and Fotios Petropoulos. Forecasting in social settings: the state of the art. *International Journal of Forecasting*, 2019.

[Marshall, 2009] Alfred Marshall. *Principles of economics: unabridged eighth edition*. Cosimo, Inc., 2009.

[Olah and Carter, 2016] Chris Olah and Shan Carter. Attention and augmented recurrent neural networks. *Distill*, 2016.

[Oreshkin *et al.*, 2020] Boris N Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis expansion analysis for interpretable time series forecasting. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 27-30*, 2020.

[Piotr and LeCun, 2009] Mirowski Piotr and Yann LeCun. Dynamic factor graphs for time series modeling. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 128–143. Springer, 2009.

[Qin *et al.*, 2017] Yao Qin, Dongjin Song, Haifeng Chen, Wei Cheng, Guofei Jiang, and Garrison Cottrell. A dual-stage attention-based recurrent neural network for time series prediction. *International Joint Conference on Artificial Intelligence*, 2017.

[Rangapuram *et al.*, 2018] Syama Sundar Rangapuram, Matthias W Seeger, Jan Gasthaus, Lorenzo Stella, Yuyang Wang, and Tim Januschowski. Deep state space models for time series forecasting. In *Advances in NIPS*, pages 7785–7794, 2018.

[Ribeiro *et al.*, 2016] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144. ACM, 2016.

[Rubanova *et al.*, 2019] Yulia Rubanova, Tian Qi Chen, and David K Duvenaud. Latent ordinary differential equations for irregularly-sampled time series. In *Advances in Neural Information Processing Systems*, pages 5321–5331, 2019.

[Schnader and Stekler, 1990] Marjorie H Schnader and Herman O Stekler. Evaluating predictions of change. *Journal of Business*, pages 99–107, 1990.

[Schwendener, 2010] Alvin Schwendener. *The estimation of financial markets by means of a regime-switching model*. PhD thesis, University of St. Gallen, 2010.

[Sezer *et al.*, 2019] Omer Berat Sezer, Mehmet Ugur Gudelek, and Ahmet Murat Ozbayoglu. Financial time series forecasting with deep learning: A systematic literature review: 2005-2019. *arXiv preprint arXiv:1911.13288*, 2019.

[Smyl, 2020] Slawek Smyl. A hybrid method of exponential smoothing and recurrent neural networks for time series forecasting. *International Journal of Forecasting*, 36(1):75–85, January 2020.

[Taylor G Smith, 2019] Aaron Smith Steven Hoelscher Taylor G Smith, Charles Drotar. pmdarima. https://github.com/tgsmith61591/pmdarima, 2019.

[Tesauro, 1995] Gerald Tesauro. Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):58–68, 1995.

[Theil, 1971] Henri Theil. Applied economic forecasting. 1971.

[Van Den Oord *et al.*, 2016] Aäron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *SSW*, 125, 2016.

[Vincent and Thome, 2019] LE Vincent and Nicolas Thome. Shape and time distortion loss for training deep time series forecasting models. In *Advances in NIPS*, pages 4191–4203, 2019.

[Ziat *et al.*, 2017] Ali Ziat, Edouard Delasalles, Ludovic Denoyer, and Patrick Gallinari. Spatio-temporal neural networks for space-time series forecasting and relations discovery. In *2017 IEEE ICDM*, pages 705–714. IEEE, 2017.