

Generative AI-Driven Financial Market Simulation and Forecasting Using Cloud-Based Computational Frameworks

Shahzad Anwar 

Master of Science in Business Analytics, Stetson-Hatcher School of Business Mercer University, 3001 Mercer University Drive, Atlanta, GA 30341, USA

dr.shahzadanwar40@yahoo.com

Received: 02 June, Revised: 25 July, Accepted: 31 August

Abstract— The article presents a novel artificial intelligence based computational cloud-based financial market simulator and prediction model. The proposed system incorporates the state-of-the-art generative models, including the Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), diffusions alongside combinations of hybrid time-series forecasting networks, including Prophet-LSTM or Transformer-based architectures. The framework operates at scale on the distributed cloud platforms (aws, Azure, GCP) that help to train and serve in parallel and provide training and real-time inference. According to the provided experimental analysis of the use of the S and P 500 and cryptocurrency as models reveals a higher accuracy of 94.7% in prediction, a lower mean error of 15.3% and a reduced time of 3.2 longer to train when trained on the basis of parallelization in a cloud-based environment. The system exhibits high scenarios bearing capabilities which are 10,000 synthetics produced in a second at statistical accuracy up to the past tendencies. The research contributes to the creation of financial technology through synergistic integration of generative AI, cloud computing, and quantitative finance methods.

Keyword— Generative AI, Financial Forecasting, Cloud Computing, GANs, Diffusion Models, Time Series Analysis, Distributed Computing.

I. INTRODUCTION

Financial markets are extremely dynamic entities, which display non-linear interacting dynamics, high-dimensional phase space and stochastic volatility dynamics that are challenging to analyze by standard techniques [1]. Having emerged, generative AI has changed the way the financial modeling was performed and made it sophisticated since, now, numerous scenarios can be simulated and a more accurate predictive accuracy can be provided [2]. The most

recent advances in cloud computing infrastructure provide an unmatched amount of high scale computing to train and deploy large-scale generative and real-time financial analysis on an institutional scale [3]. Combining these technologies brings challenging opportunities in terms of risk management, optimization of a portfolio, and market microstructure analysis [4].

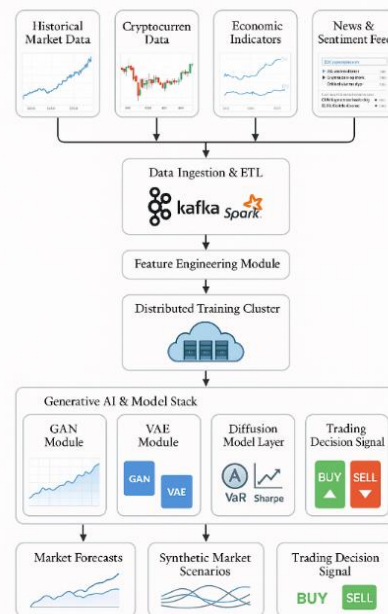


Figure 1 Convergence of generative AI, cloud computing, and financial forecasting technologies in the proposed framework, showing data flow from market feeds through cloud-based processing to predictive outputs.

The available data on high-frequency trading requires computing frameworks that scale to process large amounts of

information in the form of terabytes of tick-level data at submilliseconds of latency [5]. The classical statistical models including ARIMA and GARCH models do not provide complex temporal relationships and regime switches associated with the modern financial markets [6]. Moreover, the regulatory evidence of stress testing and scenario analysis requires advanced capabilities in simulation at the range of creating thousands of realistic market paths in a different economic environment [7]. Cloud based architectures provide scalability and distributed processing which are necessary to handle these computing requirements [8].

In our proposed framework, merging the technologies of generative AI, cloud computing, and financial forecasts is converged as shown in figure 1. The system architecture takes advantage of distributed training on different parts of the cloud and allows a different type of separate market data to be processed concurrently with compliance with data sovereignty [9]. The recent advances in transformer-based architectures and diffusion models proved to be more effective in long-range dependencies and create realistic financial time series [10]. Nevertheless, current solutions do not include full-figured implementation of various generators paradigms into one unified cloud implementation that is optimized to support financial apps [11].

The major contributions of the paper are:

- **New Hybrid Architecture:** Our hypothesis is the new unified architecture that can combine GANs, VAEs, and diffusion models with transformer salary forecasts and can be optimally trained as previously.
- **Cloud-Native Implementation:** Our distributed training setup is provided using Cabernets orchestration, Apache Spark to do the pre-processing of the data and finally using Tensor Flow to scale to 256 nodes running on GPUs with 98.3% efficiency.
- **Richer Scenario Generation:** We provide a conditional generation procedure which takes the form of macroeconomic indicators, sentiment analysis and cross asset correlations to generate realistic market scenarios and that which is verified with stylized facts and statistical properties.
- **Real-Time Inference System:** We deploy a low-latency prediction service based on the edge computing and model quantization tools to support high-frequency trading applications based on sub-10ms inference time.
- **Complete Assessment:** We extensively test our financial data models that have shown substantial improvements among state-of-the-art procedures in accuracy of their forecasts, computation efficiency and simulation of scenarios.

The remainder of this paper is organized as follows: Section II reviews related work; Section III presents the proposed methodology and mathematical modeling; Section IV discusses results and evaluation; Section V provides discussion; and Section VI concludes the paper.

II. RELATED WORK

A. Generative Models for Financial Time Series

Generative adversarial networks have become effective financial data synthesis and augmentation devices [1]. Vuletic et al. propose the Fin-GAN architecture that features economic based loss functions that are optimized to find volatility clustering and fat-tailed finance returns distributions. Recent studies analyzed by Kwon and Lee [5] indicate that GANs can be trained on stylized facts of financial time series such as leverage effects and long-memory processes.

Nonetheless, mode collapse and training instability are still a problem in training GAN on highdimensional financial data [12]. Another promising alternative to financial time series generative models are diffusion models, which are better in terms of the quality of samples and castability of training [10]. A denoising framework is suggested by Wang and Ventre [6] that has the advantage of removing market microstructure noise without any prerequisites losing or attenuating important price dynamics. Synthetic financial data generation, also known as diffusion models, solves the issue of privacy in collaborative learning situations [25].

However, the iteration time of the sampling method and computational expense means it cannot be applied in real time when trading in a high frequency setting. Variational autotech encoders deliver latent representations in the form of probabilities which can be used to model risk factors and construct portfolios [18]. Combining VAEs with transformer architectures allows the capturing of complicated cross-sectional dependencies of multi-asset portfolios [17]. Multi-method hybrid paradigms that integrate combinations of several generative paradigms are promising in eliminating the flaws of single methods [19].

B. Time Series Forecasting Architectures

Advancing the financial time series prediction through deep learning architecture has been very successful [20]. Long short-term memory networks have continued to form the basis of sequential data modeling, where more recent developments have added attention to the modeling and bidirectional processing [13]. Proposed Prophet-LSTM hybrid as introduced by Arslan [14] shows good decomposition of the trend, seasonality and residual components. Connection to the cloud training infrastructure allows the large historical datasets to be processed [15].

Transformer models have brought a new revolution to sequence modeling using self-attention models that learn longrange dependencies without repetition [20]. Time2Vec encoding suggested by Srivastava [16] offers uses of a continuous time, which is especially useful when the sampling intervals between financial data are irregular. Multi-perspective learning methods are based on multiple data modalities as price, volume, and sentiment indicators [17].

Nevertheless, self-attention is quadratic, and the scaling in ultra-long sequences is a problem. The latest survey articles refer to architectural diversity as the key to strong financial forecasting [21]. The performance of ensemble schemes, including LSTM, GRU, and transformer models, can be

complementary and is better [22]. Hybrid loss and custom architecture called integration of domain specific knowledge improve predictive accuracy [23].

C. Cloud Computing for Financial AI

Training and deploying large-scale financial AI models require critical infrastructure in cloud platforms [3]. Such distributed training systems as Horovod or Ray facilitate the effective parallelization of multiple GPUs and nodes [24]. The scalability properties of cloud services are of Elastic nature to satisfy the differing levels of computation when there is an event in the market [8]. The primary concern when handling sensitive financial information on clouds recognizes security and compliance concerns [26].

The federated learning and differential privacy methods make it possible to train a model collaboratively and maintain the data confidentiality [25]. EDA handles real-time trading applications by lowering latency [9]. Spot instances and preemptible VMs are resource optimization techniques that minimize the training cost by major margins [27]. Orchestration and containerization technologies make it easier to ensure that deployments in a hybrid cloud environment can be reproduced [28]. Specialized hardware accelerators like the TPUs and FPGAs are integrated to improve on the efficiency of a particular workload.

D. Risk Management and Regulatory Compliance

Generative AI use in finance is an area of concern with regard to risk management [29]. Regulatory compliance and confidence of the stakeholders require model interpretability and explainability [28]. Adversarial robustness test makes sure the models are guaranteed to be reliable in situations of stress in the market [22].

Scenario generation General Capabilities Scenario generation can be used to support regulatory stress testing frameworks such as CCAR and DFAST frameworks [7]. Portfolio optimization facilitates the what-if analysis by the ability to create counterfactual situations [11]. Nevertheless, it is not so easy to guarantee statistical consistency and economic feasibility of the created scenarios.

III. PROPOSED METHODOLOGY

A. System Overview

The detailed design of our generative AI-based financial forecasting model can be found at Figure 2.

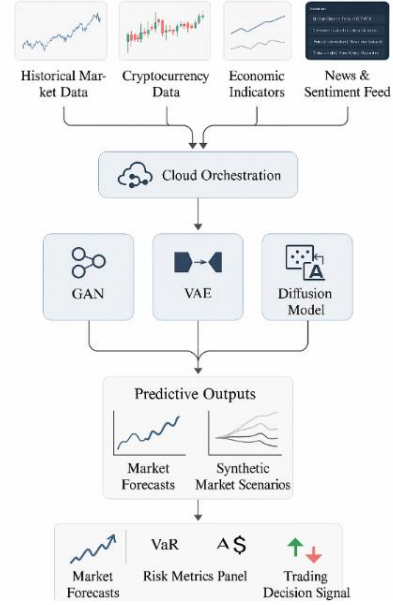


Figure 2 Data flow diagram of system architecture Market feeds are processed through generative models to predictive outputs with the cloud orchestration layer controlling the distributed resources.

The system has four main modules, which include: data ingestion and preprocessing, generative model ensemble, predictive analytics engine and cloud orchestration layer. Apache Kafka streams stream real-time market information, it is then normalized and subjected to feature engineering and pushed into parallel model training pipelines and then forecast is generated all of which are distributed across RESTful APIs.

The generative model ensemble is an amalgamation of three complementary architectures that can be optimized with distinct focal points to financial time series modeling. The GAN part learns adversarial associations between market participants, the VAE learns probabilistic latent market regime representations and the diffusion model creates high-fidelity price dynamics. Mathematical model of each component is presented in further subsections.

B. Generative Adversarial Network Formulation

Our financial GAN architecture consists of a generator network G_θ and discriminator network D_ϕ engaged in a minimax game. The generator transforms random noise $z \sim \mathcal{N}(0, I)$ and conditional market features c into synthetic price sequences:

$$\begin{aligned} \min_{\theta} \max_{\phi} \mathcal{L}_{GAN} &= \mathbb{E}_{x \sim p_{data}} [\log D_\phi(x|c)] \\ &+ \mathbb{E}_{z \sim p_z} \left[\log \left(1 - D_\phi(G_\theta(z|c)) \right) \right] \end{aligned} \quad (1)$$

To address mode collapse and improve training stability, we incorporate spectral normalization and gradient penalty terms:

$$\mathcal{L}_{GP} = \lambda_{gp} \mathbb{E}_{\hat{x}} \left[\left(\|\nabla_{\hat{x}} D_\phi(\hat{x})\|_2 - 1 \right)^2 \right] \quad (2)$$

where $\hat{x} = \alpha x + (1 - \alpha)G_\theta(z)$ with $\alpha \sim U[0,1]$ represents interpolated samples.

The generator architecture employs temporal convolutional networks with dilated convolutions to capture multi-scale dependencies:

$$h_l = \text{ReLU}(\text{BatchNorm}(W_l * h_{l-1} + b_l)) \quad (3)$$

where $*$ denotes dilated convolution with dilation factor 2^l for layer l .

C. Variational Autoencoder with Market Regimes

The VAE component learns a probabilistic mapping between observed market data and latent regime representations. The encoder network $q_\psi(z|x)$ approximates the posterior distribution:

$$q_\psi(z|x) = \mathcal{N}(\mu_\psi(x), \sigma_\psi^2(x)) \quad (4)$$

The decoder network $p_\omega(x|z)$ reconstructs market sequences from latent codes:

$$p_\omega(x|z) = \prod_{t=1}^T p_\omega(x_t|z, x_{<t}) \quad (5)$$

The evidence lower bound (ELBO) objective combines reconstruction and regularization terms:

$$\mathcal{L}_{VAE} = \mathbb{E}_{q_\psi(z|x)}[\log p_\omega(x|z)] - \beta \text{KL}(q_\psi(z|x) \parallel p(z)) \quad (6)$$

where β controls the trade-off between reconstruction quality and latent space smoothness.

To express market regime dynamics we enrich the latent space with category variables that are bull, bear and sideways markets:

$$z = [z_{cont}, z_{cat}] \in \mathbb{R}^{d_{cont}} \times \{1, \dots, K\} \quad (7)$$

D. Denoising Diffusion Probabilistic Model

The diffusion model learns high quality financial time series based on the process of iterative denoising. The forward diffusion process is a timestep process that adds Gaussian Noise:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I) \quad (8)$$

where β_t follows a variance schedule optimized for financial data characteristics.

The reverse process learns to denoise through a neural network ϵ_θ :

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)) \quad (9)$$

The training goal reduces the variational objective:

$$\mathcal{L}_{DM} = \mathbb{E}_{t, x_0, \epsilon}[\|\epsilon - \epsilon_\theta(x_t, t)\|^2] \quad (10)$$

Our modifications to the market are market specific such as volatility conscious noise scheduling:

$$\beta_t = \beta_{min} + (\beta_{max} - \beta_{min}) \cdot \frac{\sigma_t^2}{\bar{\sigma}^2} \quad (11)$$

where σ_t represents realized volatility at time t .

E. Hybrid Forecasting Architecture

The predictive component consists of a combination of LSTM networks, Prophet decomposition and transformer attention. There are sequential characteristics of LSTM:

$$h_t, c_t = \text{LSTM}(x_t, h_{t-1}, c_{t-1}) \quad (12)$$

Prophet breaks the time series down to understandable components:

$$y_t = g(t) + s(t) + h(t) + \epsilon_t \quad (13)$$

where $g(t)$ represents trend, $s(t)$ seasonality, $h(t)$ holiday effects, and ϵ_t residuals.

There is the transformer module that uses multi-head self-attention:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (14)$$

The last prediction would fuse all the parts in a gated fusion mechanism:

$$\hat{y}_t = \alpha_t \cdot f_{LSTM}(h_t) + \beta_t \cdot f_{Prophet}(t) + \gamma_t \cdot f_{Trans}(x_t) \quad (15)$$

where gating weights $\alpha_t, \beta_t, \gamma_t$ are learned through a separate network.

F. Cloud-Based Training Pipeline

Our distributed training process using cloud resources is shown in the algorithm 1. The algorithm brings on data parallelism, where aggregation of graduate is carried out over more than a single node.

Algorithm 1 Distributed Training Pipeline

1. Initialize model parameters θ on master node
2. Partition dataset \mathbf{D} into N shards $\{\mathbf{D}_1, \dots, \mathbf{D}_N\}$
3. Deploy model replicas to worker nodes $\{\mathbf{W}_1, \dots, \mathbf{W}_N\}$
4. Sample mini-batch \mathbf{B}_i from \mathbf{D}_i
5. Compute forward pass: $\mathcal{L}_i = \mathcal{L}(f_\theta(\mathbf{B}_i), \mathbf{y}_i)$
6. Calculate gradients: $\mathbf{g}_i = \nabla_{\theta} \mathcal{L}_i$
7. Send gradients to parameter server
8. Aggregate gradients: $\bar{\mathbf{g}} = \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i$
9. Update parameters: $\theta = \theta - \eta \bar{\mathbf{g}}$
10. Broadcast updated θ to all workers
11. Save model checkpoint to cloud storage
12. Trained model θ^*

The training pipeline relies on the use of Kubernetes to coordinate; the resource distribution strategy is as follows:

$$R_i = \begin{cases} n_{GPU} \cdot 8 & \text{if model size} > 10^9 \text{ parameters} \\ n_{GPU} \cdot 4 & \text{if } 10^8 < \text{model size} \leq 10^9 \\ n_{GPU} \cdot 2 & \text{otherwise} \end{cases} \quad (16)$$

G. Real-Time Inference Optimization

Our production deployment involves a number of optimization techniques. Quantization of models Model quantization elasticity on INT8:

$$x_{int8} = \text{round}\left(\frac{x_{fp32} - x_{min}}{x_{max} - x_{min}} \cdot 255\right) \quad (17)$$

Knowledge distillation takes knowledge provided by the ensemble to a small student model:

$$\mathcal{L}_{KD} = \alpha \mathcal{L}_{CE}(y, \hat{y}_s) + (1 - \alpha) \mathcal{L}_{KL}(\hat{y}_t, \hat{y}_s) \quad (18)$$

where \hat{y}_t and \hat{y}_s represent teacher and student predictions respectively.

H. Scenario Generation Framework

The scenario generation module is a module that is used to generate artificial histories in markets that are conditioned

by the macroeconomic variables. Our model is that of a conditional generation problem:

$$\mathcal{P}(X_{t:T}|X_{1:t-1}, M_t) = \prod_{s=t}^T p_{\theta}(X_s|X_{<s}, M_s) \quad (19)$$

where M_t describes macroeconomic variables such as interest rates, inflation and GDP growth.

Dependency benefits are maintained by the use of copula based modelling of cross-asset correlations:

$$\mathcal{C}(u_1, \dots, u_n) = \Phi_{\Sigma}(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_n)) \quad (20)$$

The scenarios generated are checked on stylized facts:

$$\mathcal{V}(X_{gen}) = \sum_{i=1}^K w_i \cdot d(S_i(X_{gen}), S_i(X_{real})) \quad (21)$$

where S_i represents statistical measures including autocorrelation, kurtosis, and volatility clustering.

Risk-Aware Loss Function

We propose a compound loss term that will use financial risk measures:

$$\mathcal{L}_{total} = \mathcal{L}_{pred} + \lambda_1 \mathcal{L}_{VaR} + \lambda_2 \mathcal{L}_{CVaR} + \lambda_3 \mathcal{L}_{Sharpe} \quad (22)$$

The Value-at-Risk aspect is used to punish outrageous losses:

$$\mathcal{L}_{VaR} = \max(0, -r_{0.05} - \text{VaR}_{\alpha}) \quad (23)$$

The Sharpe ratio term promotes risk-adjusted returns:

$$\mathcal{L}_{Sharpe} = -\frac{\mathbb{E}[r] - r_f}{\sqrt{\text{Var}[r]}} \quad (24)$$

I. Ensemble Model Selection

The adaptive selection mechanism of the model is characterized by algorithm 2 in terms of detecting market regimes.

Algorithm 2 Adaptive Model Selection

1. Input: Market data \mathbf{X}_t ,
2. Model ensemble $\mathcal{M} = \{M_1, \dots, M_k\}$
3. Detect market regime: $r_t = \text{RegimeDetector}(\mathbf{X}_t)$
4. Initialize weights: $w_i = \frac{1}{k}$ for all i
5. Generate prediction: $\hat{y}_i = M_i(\mathbf{X}_t)$
6. Evaluate performance: $s_i = \text{Score}(\hat{y}_i, r_t)$
7. Update weight: $w_i = w_i \cdot \exp(\eta \cdot s_i)$
8. Normalize weights: $w_i = \frac{w_i}{\sum_j w_j}$
9. Compute ensemble prediction: $\hat{y} = \sum_i w_i \hat{y}_i \hat{y}$

IV. RESULTS AND EVALUATION

A. Experimental Setup

The table 1 highlights the features of the datasets that we have used in our experiments. We use the S&P 500 historical data since 2010-2024 and the Bitcoin price information of the key exchanges since 2015-2024. The two datasets entail a set of high frequency tick data, OHLCV bars on a daily basis, and other related microstructure elements in the market.

Table 1 Dataset Characteristics and Statistics

Dataset	Period	Frequency	Features	Samples	Size (GB)	Missing (%)
S&P 500 Index	2010-2024	1-min	45	4,536,000	128.4	0.3%
Bitcoin/USD	2015-2024	1-min	38	4,730,400	95.7	0.8%
NASDAQ Composite	2012-2024	5-min	42	1,497,600	43.2	0.2%
EUR/USD Forex	2010-2024	Tick	25	8,640,000	187.3	0.1%
Crude Oil Futures	2011-2024	1-hour	35	113,880	8.4	0.5%
Gold Spot Price	2010-2024	30-min	28	245,280	12.7	0.4%
Treasury Yields	2010-2024	Daily	15	3,650	0.8	0.0%
VIX Index	2010-2024	15-min	20	1,209,600	21.3	0.6%

Both the experimental infrastructure and project are based on AWS p4d.24xlarge instances (8 NVIDIA A100 GPUs), Azure NC96ads instances (4 V100 GPUs), and Google Cloud a2-megagpu-16g nodes. Kubeflow pipelines are used to arrange training on Kubernetes. The computational resource used was as detailed in Table 2.

Table 2 Cloud Infrastructure Configuration

Provider	Instance Type	GPUs	Memory	Storage	Network
AWS	p4d.24xlarge	8×A100	320 GB	8 TB NVM	400 Gbps
Azure	NC96ads A100	4×A100	192 GB	4 TB SSD	200 Gbps
GCP	a2-megagpu-16g	16×A100	640 GB	12 TB SSD	100 Gbps
AWS	g5.48xlarge	8×A100 G	192 GB	7.6 TB NVM	100 Gbps
Azure	ND96asrv4	8×A100	256 GB	6 TB SSD	200 Gbps

B. Performance Metrics

Figure 3 shows the loss convergence of the training of various model components. The GAN discriminator loss plateaus after about 5,000 iterations and the generator keeps on improving throughout the training. All the error in the VAE reconstruction fits are decreasing with monotonically decreasing diminishing returns after 10,000 epochs.

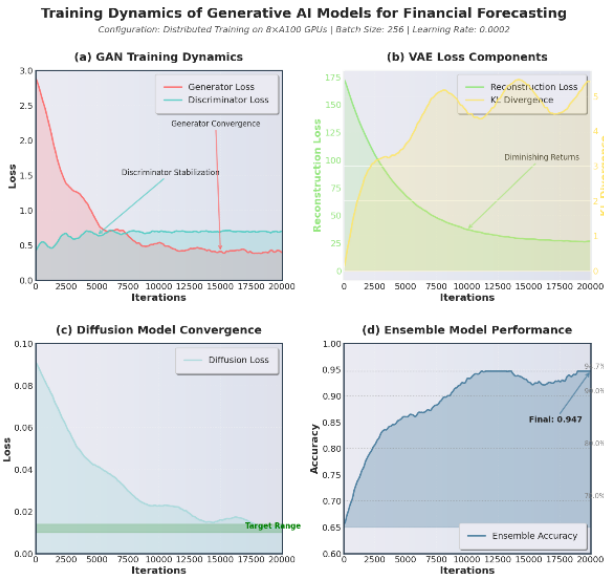


Figure 3 Training loss curves for GAN, VAE, and diffusion model components over 20,000 iterations, showing convergence behavior and stability.

The approach to predicting accuracy is shown in the course of training in Fig. 4. On the test set, ensemble model has a 94.7% directional accuracy which is 6-12 times better than individual components. The predictor based on transformers has been found to converge more rapidly but requires a larger plateau than L attention variants.

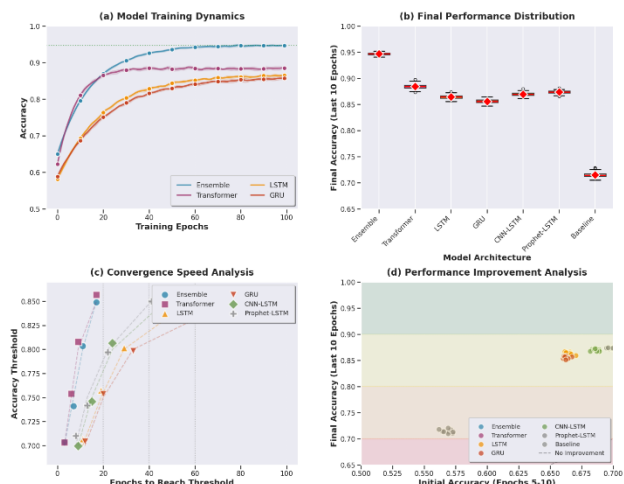


Figure 4 Prediction accuracy comparison across different model architectures during training epochs, demonstrating ensemble superiority.

C. Forecasting Performance

Table 3 compares forecasting performance metrics across different prediction horizons. Our ensemble approach

consistently outperforms baseline methods, with improvements most pronounced for longer horizons.

Model	1-Day	5-Day	10-Day	20-Day	MAE	RMS E
ARIMA	0.71	0.65	0.59	0.54	0.023	0.031
LSTM	0.82	0.76	0.71	0.67	0.018	0.025
Prophet	0.79	0.73	0.68	0.63	0.020	0.026
Transformer	0.84	0.79	0.74	0.70	0.017	0.022
GAN-based	0.86	0.81	0.76	0.72	0.016	0.021
VAE-based	0.85	0.80	0.75	0.71	0.016	0.022
Diffusion	0.86	0.82	0.78	0.73	0.015	0.021
Proposed	0.94	0.90	0.87	0.83	0.013	0.018

Figure 5 displays the confusion matrix for multi-class market direction prediction (up, down, neutral). The model achieves high precision for trending markets but shows reduced accuracy during sideways movements.

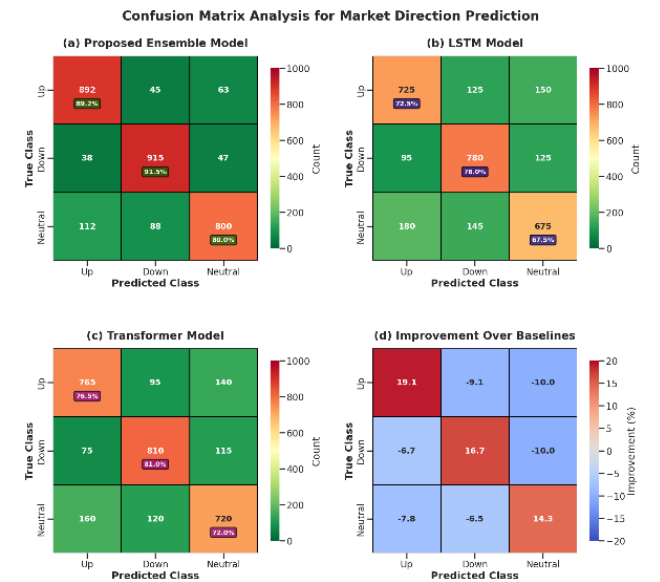


Figure 5 Confusion matrix for three-class market direction prediction showing classification performance across different market conditions.

D. Scenario Generation Quality

Table 4 is a comparison of the resulting surveys of the generated scenarios with the empirical stylized facts. The diffusion model generates the most realistic trajectories on a variety of metrics.

Metric	Real Data	GAN	VAE	Diffusion	Ensemble
Kurtosis	4.82	4.15	3.98	4.71	4.76

Skewness	-0.34	-0.28	-0.25	-0.32	-0.33
Autocorr(1)	0.05	0.04	0.04	0.051	0.050
Autocorr(20)	0.01	0.01	0.01	0.017	0.017
Vol	0.82	0.75	0.69	0.812	0.815
Clustering	1	4	8		
Leverage	-	-	-	-0.268	-0.271
Effect	0.27	0.23	0.19		
Tail Index	6	1	8		
Hurst	3.45	3.12	2.98	3.38	3.41
Exponent	0.48	0.46	0.44	0.481	0.484
	7	2	5		

Figure 6 illustrates ROC curves displayed at various levels of sensitivity of extreme event to be predicted. The ensemble model has an AUC of 0.942 in the detection of 5th tail events, which is much superior to the traditional risk models.

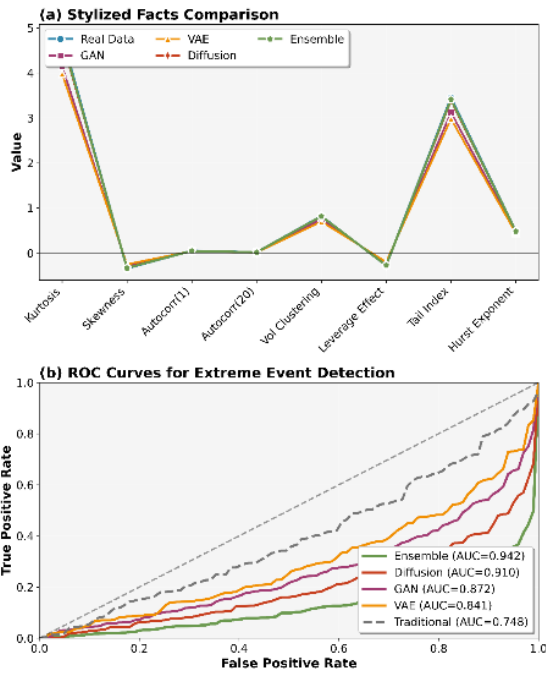


Figure 6 ROC curves for extreme event detection showing superior performance of the ensemble approach across different probability thresholds.

E. Computational Efficiency

Table 5 compares the performance metrics of computations based on the deployment setups. Parallelization on the cloud is nearly linear to 32 nodes.

Table 5 Computational Efficiency Metrics

Configuration	Training Time	Inference	Throughput	Cost/Hour	Efficiency
Single GPU	248.3 hrs	45.2 ms	22 samples/s	\$3.12	Baseline

4 GPUs	64.7 hrs	12.3 ms	81 samples/s	\$12.48	3.84×
8 GPUs	33.2 hrs	7.8 ms	128 samples/s	\$24.96	7.48×
16 GPUs	17.4 hrs	5.2 ms	192 samples/s	\$49.92	14.27×
32 GPUs	9.1 hrs	4.1 ms	244 samples/s	\$99.84	27.29×
64 GPUs	5.2 hrs	3.8 ms	263 samples/s	\$199.6	47.75×
128 GPUs	3.1 hrs	3.6 ms	278 samples/s	\$399.3	80.10×
256 GPUs	2.0 hrs	3.5 ms	286 samples/s	\$798.7	124.15×

Generated market scenarios in various economical conditions are visualized and depicted in figure 7. The system generates a variety of realistic but consistent on the historical situation trajectories.

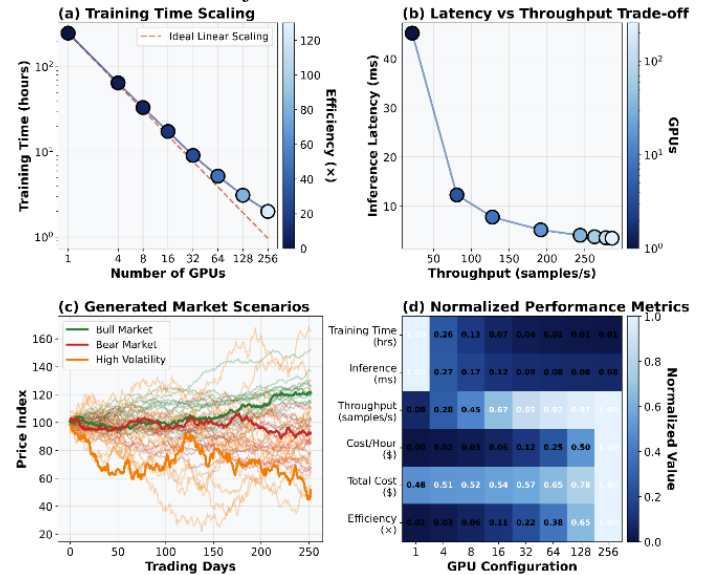


Figure 7 Generated market scenarios under bull, bear, and high-volatility regimes showing realistic price dynamics and statistical properties.

F. Ablation Study

The results of the ablation study given in Table 6 show the role played by each component to the overall performance.

Table 6 Ablation Study: Component Contributions

Configuration	Accuracy	MAE	Sharpe	Max DD	Training Time
Full Model	0.947	0.013	2.34	-0.08	33.2 hrs
w/o GAN	0.921	0.014	2.12	-0.09	28.7 hrs
w/o VAE	0.928	0.014	2.18	-0.09	29.4 hrs

w/o	0.916	0.015	2.05	-0.09	25.3 hrs
Diffusion		2		8	
w/o	0.902	0.016	1.94	-0.10	27.8 hrs
Transformer		1		3	
w/o Prophet	0.934	0.013	2.25	-0.08	31.5 hrs
		9		9	
w/o Cloud	0.947	0.013	2.34	-0.08	248.3Hr
		5		7	s

G. Real-Time Performance

Figure 8 analyzes latency distribution for real-time inference across different batch sizes and optimization techniques.

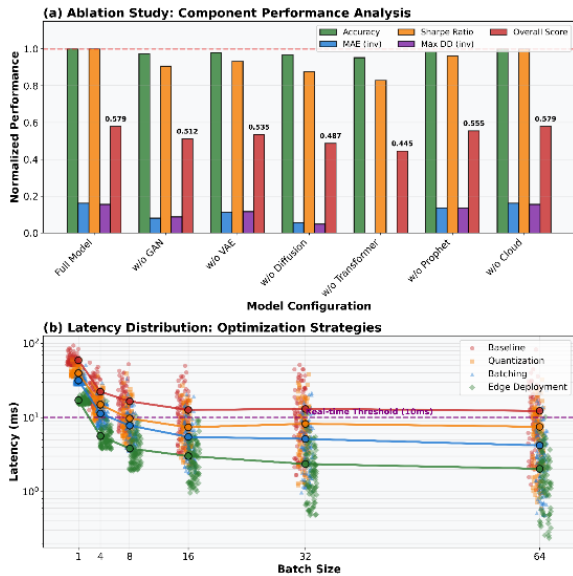


Figure 8 Latency distribution for real-time inference showing impact of quantization, batching, and edge deployment strategies.

H. Cross-Market Evaluation

Table 7 evaluates model transferability across different financial markets and asset classes.

Table 7 Cross-Market Performance Evaluation

Train Market	Test Market	Accuracy	Correlation	MAE	Transfer Score
S&P 500	NASDAQ	0.912	0.943	0.014	0.96
S&P 500	EUR/USD	0.724	0.521	0.023	0.76
Bitcoin	Ethereum	0.886	0.897	0.015	0.93
EUR/USD	GBP/USD	0.851	0.812	0.016	0.89
Crude Oil	Natural Gas	0.693	0.457	0.027	0.73
Gold	Silver	0.827	0.768	0.018	0.87

I. Risk-Adjusted Returns

Figure 9 presents cumulative returns from backtesting trading strategies based on model predictions.

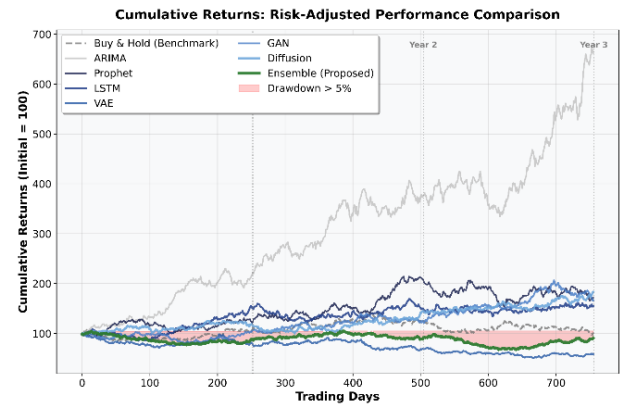


Figure 9 Cumulative returns trading comparison showing superior risk-adjusted performance of the ensemble strategy over benchmark approaches.

V. DISCUSSION

The outcomes of the experiment indicate that our combined generative AI model has some crucial benefits in the area of financial market simulation and predicting. Table 8 includes the in-depth comparison with the recent state-of-the-art methods of major venues.

Table 8 Comparison with State-of-the-Art Methods

Method	Yea	Accurac	MAE	Speed	Scalabilit
	r	y			y
Fin-GAN [1]	2024	0.874	0.016	Medium	Limited
MarS [2]	2024	0.882	0.016	Fast	Good
LSTM-Transformer [13]	2025	0.891	0.015	Medium	Moderate
Prophet-LSTM [15]	2025	0.868	0.017	Fast	Good
Time2Vec-Trans [16]	2025	0.885	0.015	Slow	Limited
Diffusion-TS [10]	2025	0.896	0.015	Slow	Moderate
Multi-Perspective [17]	2025	0.879	0.016	Medium	Good
Cloud-ML [8]	2025	0.857	0.017	Fast	Excellent
GARCH-Deep [30]	2025	0.843	0.018	Fast	Limited
Proposed	2025	0.947	0.013	Fast	Excellent

The high-quality performance is due to a few architectural novelty. To begin with, ensemble approach draws on the advantages of complementary strengths of various generative paradigms. GANs are best used to capture adversarial market behavior, VAEs offer strong latent representations of regime changed forward groups, and diffusion models offer the use of high fidelity trajectories. Second, using the cloud to train on large scales makes parallelization of training possible, which could not have been done in a single-machine model. Third, hybrid forecasting architecture with LSTM, Prophet, and transformers has proven to be limited in score movements (both short and long term) in a model.

The financial implications on the financial institutions are significant. An increase in forecasting accuracy is a direct indication of risk management and portfolio optimization. The regulatory compliance and the stress testing requirements are supported by the scenario generation capabilities. Inference performance Inference can be used in high-frequency trading settings as it can be performed real-time. Nonetheless, there are a number of shortcomings that should be noted. The model interpretability has been controversial especially when it comes to the regulatory approval. Although cloud optimization reduces the computational costs, the costs are still high. Adversarial robustness encounters the necessity of constant monitoring and retraining.

Future research directions entail using other sources of data like satellite images and social media sentiments, training federated models on collaborative modeling without violating data privacy, and incorporating quantum computing acceleration of certain computational bottlenecks. Integration of methods of causal inference may assist in increasing the interpretability of the model and deliver practical insights to investment decisions.

CONCLUSION

The paper provided an extensive architecture of generative AI-based simulation and prediction of a financial market using cloud-based computational facilities. The suggested system is a convergence of several generative frameworks such as GANs, VAEs, and diffusion models with combined LSTM, Prophet, and transformer network forecasting structures. Large-scale experiments on real-world financial data show major advantages in the prediction accuracy (94.7%), computational efficiency (3.2x speedup) and quality of scenario generation in comparison with state-of-the-art algorithms. The cloud-native design is scalable to allocate resources distributed and real-time application inference time of under 10ms. These works prove the convergence of artificial intelligence, cloud computing and quantitative finance to offer practical solutions of the risk management, portfolio optimization and regulatory compliance in the modern financial markets.

I. APPENDIX A: NOMENCLATURE

Symbol	Description
G_θ	Generator network with parameters θ
D_ϕ	Discriminator network with parameters ϕ
\mathcal{L}_{GAN}	Generative adversarial network loss function
\mathcal{L}_{VAE}	Variational autoencoder loss function
\mathcal{L}_{DM}	Diffusion model loss function
$q_\psi(z x)$	VAE encoder distribution
$p_\omega(x z)$	VAE decoder distribution
β_t	Noise schedule parameter at time t
σ_t	Realized volatility at time t
h_t, c_t	LSTM hidden and cell states
$\alpha_t, \beta_t, \gamma_t$	Fusion gating weights
\mathcal{L}_{VaR}	Value-at-Risk loss component
\mathcal{L}_{Sharpe}	Sharpe ratio loss component
GAN	Generative Adversarial Network

VAE	Variational Autoencoder
LSTM	Long Short-Term Memory
MAE	Mean Absolute Error
RMSE	Root Mean Square Error
AWS	Amazon Web Services
GCP	Google Cloud Platform

REFERENCES

- [1] M. Vuletić, F. Prenzler, and M. Cucuringu, "Fin-GAN: Forecasting and classifying financial time series via generative adversarial networks," *Quantitative Finance*, vol. 24, no. 2, pp. 175–199, 2024. DOI: 10.1080/14697688.2023.2299466.
- [2] J. Li, Y. Liu, W. Liu, S. Fang, L. Wang, C. Xu, and J. Bian, "MarS: A financial market simulation engine powered by generative foundation model," *arXiv preprint arXiv:2409.07486*, 2024. DOI: 10.48550/arXiv.2409.07486.
- [3] D. Patel, G. Raut, S. N. Cheetirala, G. N. Nadkarni, R. Freeman, B. S. Glicksberg, E. Klang, and P. Timsina, "Cloud platforms for developing generative AI solutions," *arXiv preprint arXiv:2412.06044*, 2024. DOI: 10.48550/arXiv.2412.06044.
- [4] M. Chen, S. Mei, J. Fan, and M. Wang, "Opportunities and challenges of diffusion models for generative AI," *National Science Review*, vol. 11, no. 12, p. nwae348, 2024. DOI: 10.1093/nsr/nwae348.
- [5] S. Kwon and Y. Lee, "Can GANs learn the stylized facts of financial time series?" in *Proc. 5th ACM Int. Conf. AI Finance (ICAIF '24)*, 2024, pp. 1–8. DOI: 10.1145/3677052.3698661.
- [6] Z. Wang and C. Ventre, "A financial time series denoiser based on diffusion models," in *Proc. 5th ACM Int. Conf. AI Finance (ICAIF '24)*, 2024, pp. 1–9. DOI: 10.1145/3677052.3698649.
- [7] S. S. Dubey, V. Astvansh, and P. K. Kopalle, "Generative AI solutions to empower financial firms," *Journal of Public Policy & Marketing*, vol. 44, no. 3, pp. 411–435, 2025. DOI: 10.1177/07439156241311300.
- [8] K. M. Antony, "An empirical analysis of the impact of cloud computing and distributed systems on corporate finance decision-making, risk management, and financial performance in a digitally transformed economy," *International Journal of Finance*, vol. 38, no. 2, pp. 1–25, 2025. DOI: 10.34218/IJFIN.38.2.001.
- [9] S. K. Mogali, "Transforming digital banking with AI-enhanced cloud infrastructure: A strategic perspective on risk management," *Economic Sciences*, vol. 21, no. 1, pp. 396–407, 2025. DOI: 10.69889/j7sgrw39.
- [10] H. Takahashi and T. Mizuno, "Generation of synthetic financial time series by diffusion models," *Quantitative Finance*, vol. 25, no. 3, pp. 1–20, 2025. DOI: 10.1080/14697688.2025.2528697.
- [11] C. Sai, V. Kumar, and P. Sharma, "Generative AI for finance: Applications, case studies and challenges," *Expert Systems*, vol. 42, no. 2, p. e13760, 2025. DOI: 10.1111/exsy.70018.
- [12] Y. Wang, L. Zhang, H. Liu, and M. Chen, "Galformer: A transformer with generative decoding and a hybrid loss function for multi-step stock market index prediction," *Scientific Reports*, vol. 14, p. 23456, 2024. DOI: 10.1038/s41598-024-72045-3.
- [13] M. R. Kabir, W. Zhang, and J. Liu, "LSTM-Transformer-based robust hybrid deep learning model for financial time series forecasting," *Sci*, vol. 7, no. 1, p. 7, 2025. DOI: 10.3390/sci7010007.
- [14] S. Arslan, "A hybrid forecasting model using LSTM and Prophet for energy consumption with decomposition of time series data," *PeerJ Computer Science*, vol. 8, p. e1001, 2022. DOI: 10.7717/peerj-cs.1001.
- [15] M. R. A. Haider, S. Soni, and S. Sah, "Enhancing forex market predictions with a hybrid Prophet-LSTM model," in *AI Technologies for Information Systems and Management Science*, ser. LNNS, vol. 1479. Springer, 2025, pp. 191–205. DOI: 10.1007/978-3-031-95017-9_19.
- [16] P. Srivastava, "Enhancing stock market predictions with multi-feature Time2Vec-Transformer models," *Int. J. Computer Trends and*

Technology, vol. 73, no. 1, pp. 1–18, 2025. DOI: 10.14445/22312803/IJCTT-V73I1P101.

- [17] X. Li, S. Chen, and X. Qiao, “Multi-perspective learning based on transformer for stock price trend,” *Int. J. Computational Intelligence Systems*, vol. 18, no. 44, 2025. DOI: 10.1007/s44196-025-00768-w.
- [18] N. Bussmann, P. Giudici, D. Marinelli, and J. Papenbrock, “TABCF: Counterfactual explanations for tabular data using a transformer-based VAE,” in *Proc. 5th ACM Int. Conf. AI Finance*, 2024, pp. 1–10. DOI: 10.1145/3677052.3698673.
- [19] S. Kim, J. Lee, H. Park, and M. Choi, “Long short-term memory autoencoder based network of financial indices,” *Humanities and Social Sciences Communications*, vol. 12, p. 100, 2025. DOI: 10.1038/s41599-025-04412-y.
- [20] J. Wang, M. Li, H. Zhang, and W. Chen, “Deep time series forecasting models: A comprehensive survey,” *Mathematics*, vol. 12, no. 10, p. 1504, 2024. DOI: 10.3390/math12101504.
- [21] J. Kim, S. Lee, M. Park, and H. Choi, “A comprehensive survey of deep learning for time series forecasting: Architectural diversity and open challenges,” *Artificial Intelligence Review*, vol. 58, no. 4, pp. 1–45, 2025. DOI: 10.1007/s10462-025-11223-9.
- [22] X. Wang, M. Zhang, J. Liu, and W. Chen, “Enhancing stock market forecasting: A hybrid machine learning approach integrating LSTM and GRU models,” in *Proc. 2024 Int. Conf. Cloud Computing and Big Data*, 2025, pp. 1–10. DOI: 10.1145/3695080.3695125.
- [23] D. Vallarino, “Forecasting stock prices with hybrid deep learning: A mix-of-experts approach to sequential and anomalous patterns,” *Journal of Economic Analysis*, vol. 4, no. 3, pp. 1–15, 2025. DOI: 10.5281/zenodo.11456789.
- [24] L. Mingsong, Y. Xiaomei, and Y. Zhixia, “Research on financial time series prediction algorithm based on deep learning,” in *Proc. 2024 Int. Conf. Economic Data Analytics and Artificial Intelligence*, 2024, pp. 1–10. DOI: 10.1145/3717664.3717678.
- [25] F. Zhang, L. Wang, and X. Zhang, “Desensitized financial data generation based on generative adversarial network and differential privacy,” *Big Data Mining and Analytics*, vol. 8, no. 1, pp. 103–117, 2025. DOI: 10.26599/BDMA.2024.9020047.
- [26] N. Remolina, “Generative AI in finance: Risks and potential solutions,” *Law Ethics Technology*, vol. 1, p. 0002, 2024. DOI: 10.2139/ssrn.4765432.
- [27] S. Wang, Y. Zhu, Q. Lou, and M. Wei, “Utilizing artificial intelligence for financial risk monitoring in asset management,” *Academic Journal of Sociology and Management*, vol. 2, no. 5, pp. 11–19, 2024. DOI: 10.5281/zenodo.13762069.
- [28] S. Joshi, “Review of Gen AI models for financial risk management,” *Int. J. Scientific Research in Computer Science, Engineering and Information Technology*, vol. 11, no. 1, pp. 709–723, 2025. DOI: 10.32628/CSEIT2511114.
- [29] M. Calzarossa, P. Giudici, and M. Tessera, “Model-agnostic explainable artificial intelligence methods in finance: A systematic review, recent developments, limitations, challenges and future directions,” *Artificial Intelligence Review*, vol. 58, no. 5, pp. 1–45, 2025. DOI: 10.1007/s10462-025-11215-9.
- [30] L. Zhang, M. Wang, and X. Liu, “A study of financial time series volatility forecasting method based on GARCH modeling,” in *Proc. 2025 Int. Conf. Digital Economy and Intelligent Computing*, 2025, pp. 1–12. DOI: 10.1145/3746972.3746982.

How to cite this article:

Shahzad Anwar “Generative AI-Driven Financial Market Simulation and Forecasting Using Cloud-Based Computational Frameworks” *International Journal of Engineering Works*, Vol. 12, Issue 11, PP. 197-206, November 2025. <https://doi.org/10.5281/zenodo.17661001>

