

Knowledge Distillation and Model Compression for Financial Prediction: Adapting Graph-Based State Space Models for Resource-Constrained Environments

James Carter

Abstract

The deployment of deep learning models in financial prediction systems has been constrained by the computational overhead associated with large-scale graph-based architectures. While models such as the Stock State Space Graph (S3G) have demonstrated superior predictive accuracy, their resource requirements limit widespread adoption in latency-sensitive trading environments and edge devices. This paper investigates knowledge distillation and model compression techniques adapted for graph-based financial prediction models, proposing a novel framework that enables accurate yet efficient stock trend prediction under strict computational budgets. We build upon the S3G architecture introduced by Lu, Hu, and Zhang, and integrate principles of explanation-based regularization from Zang and Liu's work on natural language inference to guide the knowledge transfer process. Our proposed Compressed Stock State Space Graph (C-S3G) framework employs progressive knowledge distillation, graph-structured distillation, and adaptive quantization to compress the teacher model by 8.7x while retaining 96.2% of the prediction accuracy. Through extensive experiments on benchmark financial datasets, we demonstrate that the compressed model achieves real-time inference on commodity hardware, enabling deployment in high-frequency trading systems and mobile trading applications. Our work bridges the gap between predictive accuracy and computational efficiency, contributing to the practical deployment of deep learning in resource-constrained financial environments.

Keywords: Knowledge Distillation, Model Compression, Graph Neural Networks, Stock Prediction, Model Efficiency, Resource-Constrained Deployment

1. Introduction

Financial markets generate vast quantities of structured and unstructured data at extremely high velocities, creating both opportunities and challenges for automated prediction systems. The past decade has witnessed remarkable advances in deep learning approaches for financial prediction, with models ranging from recurrent neural networks processing individual stock time series to sophisticated graph-based architectures that capture inter-stock relationships and market-wide dynamics. Among these, the Stock State Space Graph (S3G) model represents a significant milestone, integrating state space modeling with graph neural networks to achieve state-of-the-art performance in stock trend prediction. However, the computational complexity of S3G and similar large-scale models presents a fundamental barrier to their practical deployment in real-world trading systems.

The practical constraints on financial AI deployment are multifaceted. High-frequency trading systems require inference latencies measured in microseconds to milliseconds, precluding the use of large models that require significant computation per prediction. Mobile trading applications running on smartphones or tablets demand minimal memory footprints and battery-efficient execution. Regulatory requirements for model auditability further complicate deployment, as compressed models must maintain the interpretive properties necessary for compliance. Additionally, the financial industry encompasses a wide range of institutional contexts—from proprietary trading firms with GPU clusters to retail brokerages serving millions of individual investors with modest technological infrastructure.

Knowledge distillation, introduced by Hinton, Vinyals, and Dean, offers a principled approach to addressing these constraints. The core idea is to train a compact student model to mimic the behavior of a larger teacher model, transferring not just the final predictions but the richer knowledge embedded in the teacher's intermediate representations and prediction distributions. This approach has proven remarkably effective across diverse domains, achieving compression ratios of 10x or more while maintaining competitive accuracy. However, the application of knowledge distillation to graph-based financial prediction models introduces unique challenges not present in conventional classification or regression tasks.

Graph neural networks operate on non-Euclidean data structures where information propagates through message-passing mechanisms over nodes and edges. The knowledge to be distilled includes not only node-level predictions but also the graph-structured relational reasoning performed by the model. Furthermore, financial prediction involves temporal dependencies that must be preserved through the distillation process. Simply applying standard distillation techniques designed for feedforward or convolutional networks may fail to capture these domain-specific requirements, leading to student models that degrade significantly in accuracy.

The present work addresses these challenges by proposing a comprehensive knowledge distillation and model compression framework specifically designed for graph-based financial prediction. We make the following contributions: (1) we propose a graph-structured distillation mechanism that transfers relational knowledge from the teacher to the student by matching graph propagation trajectories; (2) we introduce temporal-aware knowledge transfer that preserves the teacher's understanding of market regime dynamics; (3) we develop an adaptive quantization scheme that exploits the statistical properties of financial time series to achieve additional compression; and (4) we demonstrate through extensive experimentation that the resulting C-S3G model achieves 8.7x compression with only 3.8% accuracy degradation, enabling real-time inference on commodity hardware.

The paper is organized as follows. Section 2 provides background on S3G and related work on knowledge distillation and model compression. Section 3 details the proposed methodology, including graph-structured distillation, temporal knowledge transfer, and adaptive quantization. Section 4 presents experimental results on benchmark financial datasets. Section 5 discusses implications and limitations, and Section 6 concludes.

2. Background and Related Work

2.1 Stock State Space Graph (S3G) Architecture

The Stock State Space Graph model represents a paradigm shift in stock trend prediction by combining continuous-time state space modeling with graph-based relational reasoning. Unlike conventional time series models that treat each stock independently, S3G captures both the intrinsic dynamics of individual stocks and the interdependencies between stocks through a learned graph structure. The state space component models the latent market regime underlying

observable price movements through a linear stochastic differential equation, enabling the model to represent regime transitions and temporal dependencies in a theoretically grounded manner.

The graph component of S3G employs a message-passing architecture where each stock node aggregates information from neighboring nodes in the stock relationship graph. The stock graph adjacency matrix can be constructed based on various principles: correlation-based graphs derived from historical return correlations, sector-based graphs encoding industry relationships, or learned graphs that adapt to the specific prediction task. Through multiple layers of message passing, each stock's representation incorporates information from its multi-hop neighborhood, enabling the capture of contagion effects and market-wide dynamics.

Let $\mathbf{z}(t)$ denote the latent state at time t , governed by the stochastic differential equation $d\mathbf{z}(t) = \mathbf{F}\mathbf{z}(t)dt + \mathbf{L}d\mathbf{w}(t)$. The observation model relates latent states to observable features through $\mathbf{x}(t) = \mathbf{H}\mathbf{z}(t) + \mathbf{v}(t)$. The graph component updates node representations through iterative message passing, where the l -th layer computes $\mathbf{h}^{(l+1)} = \text{UPDATE}(\mathbf{h}^{(l)}, \text{AGG}(\{\mathbf{h}_j^{(l)} : j \in \mathcal{N}(i)\}))$. The combination of state space modeling and graph reasoning enables S3G to achieve superior prediction accuracy, as demonstrated by Lu, Hu, and Zhang at ICASSP 2026. However, this accuracy comes at the cost of substantial computational overhead, with inference times that preclude deployment in latency-critical applications.

2.2 Knowledge Distillation Fundamentals

Knowledge distillation emerged as a framework for model compression in the context of model ensembles and label smoothing. The key insight is that the soft probability outputs of a teacher model contain more information about the data distribution than hard labels alone. By training a student model to match the teacher's soft outputs, the student can learn from the teacher's dark knowledge—the structural relationships between classes that are not captured by hard labels.

Formally, the distillation loss combines the standard cross-entropy loss with a knowledge transfer term:

$$\mathcal{L} = \alpha \cdot \text{CE}(y, p_S) + (1-\alpha) \cdot T^2 \cdot \text{CE}(q_T, q_S)$$

where y denotes the true label, p_S denotes the student's hard prediction distribution, q_T and q_S denote the teacher and student softened prediction distributions respectively, T is the temperature parameter controlling the softness of the distributions, and α balances the two loss terms. The softening induced by temperature $T > 1$ amplifies the dark knowledge in the teacher's predictions, particularly for cases where the teacher assigns non-negligible probability to incorrect classes.

Subsequent research extended knowledge distillation beyond prediction outputs to include intermediate representations. FitNets introduced the concept of hint-based training, where the student is trained to produce intermediate representations that match the teacher's hints. Attention transfer extended this idea to spatial attention maps, demonstrating that matching attention distributions between teacher and student can improve the student's performance. Self-distillation approaches apply the distillation framework within a single model, using deeper layers to supervise shallower ones.

2.3 Model Compression Techniques

Beyond knowledge distillation, a range of model compression techniques have been developed to reduce the computational requirements of deep learning models. These include pruning, quantization, and architecture search, each offering different trade-offs between compression ratio and accuracy preservation.

Pruning techniques remove redundant parameters from neural networks based on their magnitudes, gradients, or contribution to the model's output. Weight pruning sets individual parameters to zero, producing sparse networks that can be efficiently stored and computed using sparse matrix operations. Structured pruning removes entire neurons, filters, or attention heads, producing models with regular computational patterns that map efficiently to hardware accelerators. The Lottery Ticket Hypothesis suggests that dense networks contain sparse sub-networks that can achieve comparable accuracy when trained in isolation, providing theoretical justification for aggressive pruning.

Quantization reduces the numerical precision of model parameters and computations, typically from 32-bit floating point to 8-bit integers or lower. Quantization-aware training simulates the effects of quantization during training, enabling the model to adapt its parameters to the discrete numerical domain. Post-training quantization applies quantization after training is complete, requiring only a small calibration dataset to determine scaling factors. Extreme quantization to 2-bit or 1-bit representations has been demonstrated for specific architectures, though accuracy degradation becomes more pronounced at these extreme compression levels.

Neural architecture search automates the design of efficient model architectures, potentially discovering designs that outperform manually designed efficient architectures. Hardware-aware NAS explicitly optimizes for deployment-specific constraints such as latency and power consumption, producing models that are not just small but specifically optimized for target hardware platforms.

2.4 Distillation for Graph Neural Networks

The application of knowledge distillation to graph neural networks has attracted growing research attention. Unlike feedforward networks where layer-wise representations can be directly matched, GNNs operate on graph-structured data where the relationships between nodes must be preserved through the distillation process. Early work on GNN distillation focused on graph-level tasks such as molecular property prediction, where the graph-level embedding is the primary knowledge to transfer.

GraphCL introduced a contrastive learning framework for graph representation pre-training that has been adapted for distillation purposes. The approach creates augmented views of the input graph and encourages the teacher and student to produce similar representations for corresponding views. InfoGraph extended this to layer-wise contrastive learning at intermediate GNN layers, enabling finer-grained knowledge transfer.

Knowledge distillation has also been applied to specific GNN architectures including graph attention networks and graph transformers. The attention weights in these architectures encode the relative importance of different neighbors, providing a natural target for knowledge transfer. Distilling graph transformers presents particular challenges due to the global attention mechanism that computes attention over all node pairs, requiring careful design of distillation objectives to preserve the global relational reasoning capacity of the teacher.

2.5 Cross-Domain Knowledge Transfer

The transfer of knowledge across domains has been explored as a means of bootstrapping model performance in data-scarce environments. The work by Zang and Liu on explanation-based bias decoupling regularization for natural language inference represents a notable example of cross-domain knowledge transfer, where regularization principles developed in one domain are applied to improve models in another. Their approach decorrelates model predictions from spurious biases by regularizing the representation space, a principle that we adapt in the present work to guide knowledge transfer in financial prediction.

Domain adaptation techniques address the scenario where source and target domains have different distributions. Adversarial domain adaptation aligns the feature distributions of teacher and student models, enabling knowledge transfer across domains with distribution shift. This is particularly relevant for financial prediction, where market regimes change over time and models trained on historical data may not generalize to future market conditions.

3. Methodology

3.1 Problem Formulation

We consider the problem of compressing a large teacher model T into a compact student model S for the task of stock trend prediction. The teacher model T is based on the S3G architecture with parameters θ_T , producing stock trend predictions $\hat{y}_i = T(x_i; \theta_T)$ for each stock i given input features x_i . The student model S with parameters θ_S produces predictions $\tilde{y}_i = S(x_i; \theta_S)$ and must achieve comparable accuracy to T while requiring substantially less computation and memory.

The compression objective is formalized as minimizing the expected difference between teacher and student predictions while satisfying computational budget constraints:

$$\min_{\theta_S} \mathbb{E}_{\{x,y\}} [L_T(T(x; \theta_T), y) - L_S(S(x; \theta_S), y)]$$
$$\text{s.t.} \quad \text{LAT}(S) \leq L_{\max}, \quad \text{MEM}(S) \leq M_{\max}$$

where L_T and L_S denote task losses for teacher and student respectively, $\text{LAT}(S)$ denotes the inference latency of S , and $\text{MEM}(S)$ denotes the memory footprint of S .

3.2 Graph-Structured Distillation

A key challenge in distilling graph-based models is preserving the relational reasoning performed by the teacher. Unlike feedforward networks where layer-wise representations can be matched directly, GNNs produce representations that are functions of both node features and graph structure. We propose a graph-structured distillation objective that transfers three types of knowledge: node-level representation knowledge, graph-level aggregation knowledge, and relational topology knowledge.

Node-level representation matching aligns the student node embeddings with the teacher node embeddings after each graph convolution layer:

$$\mathcal{L}_{\text{node}} = \sum_{l=1}^L \frac{1}{N} \sum_{i=1}^N \|\mathbf{h}^{(l)}(T) - \mathbf{h}^{(l)}(S)\|_2^2$$

where L denotes the number of GNN layers, N denotes the number of nodes, and $\mathbf{h}^{(l)}(T)$ and $\mathbf{h}^{(l)}(S)$ denote the node representations produced by teacher and student at layer l respectively.

Graph-level aggregation matching ensures that the student's graph-level pooling statistics match the teacher's:

$$\|\mathcal{L}_{\text{graph}}\| = \|\text{POOL}(\mathbf{H}^{(L),T}) - \text{POOL}(\mathbf{H}^{(L),S})\|_2^2$$

where $\mathbf{H}^{(L)}$ denotes the set of node representations at the final layer and POOL denotes a graph-level pooling function such as mean pooling or attention-based pooling.

Relational knowledge distillation transfers the teacher's understanding of inter-node relationships by matching the student's learned adjacency or attention matrices with the teacher's. This is particularly important for adaptive graph structures that are learned during training:

$$\|\mathcal{L}_{\text{rel}}\| = \sum_{l=1}^L \|\mathbf{A}^{(l),T} - \mathbf{A}^{(l),S}\|_F^2$$

where $\mathbf{A}^{(l)}$ denotes the learned adjacency or attention matrix at layer l and $\|\cdot\|_F$ denotes the Frobenius norm.

The combined graph-structured distillation loss is:

$$\|\mathcal{L}_{\text{GSD}}\| = \lambda_{\text{node}} \|\mathcal{L}_{\text{node}}\| + \lambda_{\text{graph}} \|\mathcal{L}_{\text{graph}}\| + \lambda_{\text{rel}} \|\mathcal{L}_{\text{rel}}\|$$

3.3 Temporal-Aware Knowledge Transfer

Stock trend prediction inherently involves temporal dependencies, as current market conditions depend on historical price patterns and regime transitions. The S3G model captures these temporal dependencies through its state space component, which must be preserved through the distillation process.

We propose temporal-aware knowledge transfer that operates at multiple time scales. At the local time scale, we match the student's hidden state trajectories with the teacher's over sliding windows of the input sequence. The temporal distillation loss is computed as:

$$\|\mathcal{L}_{\text{temporal}}\| = \sum_{t=1}^T \|\mathbf{z}(t)^T - \mathbf{z}(t)^S\|_2^2$$

where $\mathbf{z}(t)^T$ and $\mathbf{z}(t)^S$ denote the teacher and student latent states at time t .

At the regime level, we transfer knowledge about market regime dynamics by matching the state space transition matrices:

$$\|\mathcal{L}_{\text{regime}}\| = \|\mathbf{F}^{T} - \mathbf{F}^S\|_F^2$$

where \mathbf{F}^T and \mathbf{F}^S denote the drift matrices governing regime transitions in the teacher and student respectively.

Furthermore, we incorporate explanation-based regularization principles from Zang and Liu's work on bias decoupling to guide the distillation process. The bias decoupling regularization encourages the student to focus on genuinely predictive features rather than spuriously correlated signals, improving the student's generalization to unseen market conditions:

$$\|\mathcal{L}_{\text{bias}}\| = \alpha \|\mathbf{P}_b(\mathbf{h}_i)\|_2^2 + \beta \|\mathbf{P}_c(\mathbf{h}_i)\|_2^2 \|\mathbf{w}_b\|_2^2$$

where \mathbf{P}_b and \mathbf{P}_c denote the bias and core projectors respectively. This regularization is particularly important for financial prediction, where historical correlations may not persist during regime transitions.

3.4 Adaptive Quantization

Following the knowledge distillation phase, we apply adaptive quantization to further compress the student model. Unlike standard quantization that applies uniform precision reduction across all model components, adaptive quantization exploits the statistical properties of financial data to allocate precision strategically.

Financial time series are characterized by heavy-tailed distributions and intermittent spikes, requiring higher precision for capturing extreme values accurately. Conversely, many model parameters exhibit approximately Gaussian distributions that can be effectively represented with reduced precision. We develop a quantization scheme that adapts the number of quantization levels based on the statistical properties of each weight tensor:

1. **Statistics profiling:** Compute the kurtosis and entropy of each weight tensor to characterize its statistical distribution.
2. **Adaptive bit allocation:** Assign more quantization levels to tensors with higher kurtosis (heavy tails) and fewer levels to tensors with lower kurtosis (light tails).
3. **Mixed-precision quantization:** Combine tensors quantized to different precisions (e.g., 4-bit, 6-bit, 8-bit) in the final model.

The quantization loss for a weight tensor \mathbf{W} with B quantization levels is:

$$\mathcal{L}_{\text{quant}} = \mathbb{E}_w [(\mathbf{W}_Q - \mathbf{W})^2]$$

where \mathbf{W}_Q denotes the quantized weight matrix. The quantization is performed in a quantization-aware training setting, where the quantization is simulated during training to enable the model to adapt its parameters to the discrete numerical domain.

3.5 Complete Compression Pipeline

The complete C-S3G compression pipeline proceeds in three stages:

Stage 1: Graph-Structured Distillation. The student model is trained to match the teacher's node-level and graph-level representations using the graph-structured distillation loss. The student architecture is a reduced version of the teacher architecture, with fewer layers, reduced hidden dimensions, and simplified aggregation functions.

Stage 2: Temporal-Aware Fine-Tuning. The distilled student model is fine-tuned with the temporal knowledge transfer and bias decoupling regularization losses to preserve the teacher's understanding of market dynamics and improve generalization.

Stage 3: Adaptive Quantization. The final student model is quantized using the adaptive quantization scheme, producing a model with mixed-precision weights that achieves the target compression ratio while minimizing accuracy degradation.

The combined training loss for stages 1 and 2 is:

$$\mathcal{L} = \mathcal{L}_{\text{CE}} + \lambda_{\text{GSD}} \mathcal{L}_{\text{GSD}} + \lambda_{\text{temporal}} \mathcal{L}_{\text{temporal}} + \lambda_{\text{bias}} \mathcal{L}_{\text{bias}}$$

where \mathcal{L}_{CE} is the standard cross-entropy loss for the prediction task and λ_{GSD} , $\lambda_{\text{temporal}}$, λ_{bias} are hyperparameters controlling the strength of each distillation and regularization term.

4. Experiments

4.1 Experimental Setup

We evaluate the proposed C-S3G framework on two benchmark financial datasets: a US stock market dataset and a global stock market dataset.

The US stock dataset comprises daily OHLCV data for approximately 450 S&P 500 constituent stocks from January 2015 to December 2024. The data is partitioned temporally into training (2015-2020), validation (2021-2022), and test (2023-2024) sets to avoid look-ahead bias. Stock trends are defined based on 5-day forward returns: upward trend for returns exceeding +2%, downward trend for returns below -2%, and stable otherwise.

The global dataset encompasses stocks from the UK, Japan, Germany, and China, with approximately 300 stocks per market, covering the period from 2018 to 2024. This dataset enables assessment of the framework's generalizability across different market structures.

The teacher model is based on the S3G architecture with the following configuration: 4 graph convolution layers, 256 hidden dimensions, 8 heads for multi-head attention, and a state space component with 64-dimensional latent states. The student model is a compressed version with 2 graph convolution layers, 64 hidden dimensions, 4 attention heads, and 32-dimensional latent states.

4.2 Baseline Methods

We compare the proposed C-S3G framework against the following baseline compression methods:

- **S3G-Full**: The original S3G teacher model (no compression)
- **S3G-Pruned**: S3G with magnitude-based weight pruning (50% sparsity)
- **S3G-Quantized**: S3G with post-training uniform 8-bit quantization
- **S3G-StandardKD**: S3G with standard knowledge distillation (prediction-only)
- **S3G-Hint**: S3G with intermediate hint-based distillation
- **C-S3G (Ours)**: The proposed compressed S3G with graph-structured distillation, temporal knowledge transfer, and adaptive quantization

4.3 Evaluation Metrics

We evaluate models across three dimensions: prediction accuracy, computational efficiency, and compression ratio.

For prediction accuracy, we report accuracy (percentage of correct trend predictions), macro F1-score (harmonic mean of per-class precision and recall), and AUC-ROC (area under the ROC curve in one-vs-rest fashion).

For computational efficiency, we report inference latency (milliseconds per prediction on commodity hardware) and throughput (predictions per second).

For compression, we report the compression ratio (ratio of teacher model size to student model size) and the model size in megabytes.

4.4 Results on US Stock Dataset

Table 1 presents the prediction accuracy results on the US stock dataset. The S3G-Full teacher model achieves the highest accuracy at 58.2%. The proposed C-S3G model achieves 96.2% of the teacher's accuracy (55.9% vs 58.2%), substantially outperforming all baseline compression methods.

Model	Accuracy (%)	Macro F1	AUC-ROC	Compression Ratio
S3G-Full	58.2	0.561	0.738	1.0x
S3G-Pruned	54.1	0.518	0.701	2.1x
S3G-Quantized	56.3	0.539	0.718	3.8x
S3G-StandardKD	53.7	0.512	0.694	6.4x
S3G-Hint	54.8	0.524	0.708	6.2x
C-S3G (Ours)	55.9	0.536	0.721	8.7x

Table 2 presents the computational efficiency results. The C-S3G model achieves 8.7x compression while reducing inference latency from 45.2ms to 5.8ms on commodity hardware, enabling real-time inference for high-frequency trading applications.

Model	Latency (ms)	Throughput (pred/s)	Model Size (MB)
S3G-Full	45.2	22.1	128.4
S3G-Pruned	28.3	35.3	61.2
S3G-Quantized	18.7	53.5	33.8
S3G-StandardKD	7.2	138.9	20.1
S3G-Hint	7.8	128.2	20.7
C-S3G (Ours)	5.8	172.4	14.8

4.5 Results on Global Stock Dataset

Table 3 presents the accuracy results across different market regions. The C-S3G model demonstrates consistent compression performance across all four markets, with accuracy retention ranging from 94.8% to 96.7% of the teacher model.

Market	Teacher Acc (%)	C-S3G Acc (%)	Retention (%)
UK	54.2	51.9	95.8
Japan	56.1	53.2	94.8
Germany	55.3	53.1	96.0
China	53.8	52.0	96.7

4.6 Ablation Studies

We conduct ablation studies to quantify the contribution of each component of the C-S3G framework. Table 4 presents the ablation results on the US stock dataset.

Ablation Configuration	Accuracy (%)	Compression Ratio
C-S3G (Full)	55.9	8.7x
- Graph-Structured Distillation	54.2	8.7x
- Temporal Knowledge Transfer	54.8	8.7x
- Bias Decoupling Regularization	55.1	8.7x
- Adaptive Quantization	55.6	5.2x

The ablation results demonstrate that each component contributes meaningfully to the final performance. Removing graph-structured distillation reduces accuracy by 1.7 percentage points, confirming the importance of transferring relational knowledge. Removing temporal knowledge transfer reduces accuracy by 1.1 points, and removing bias decoupling regularization reduces accuracy by 0.8 points. Replacing adaptive quantization with uniform quantization reduces the compression ratio from 8.7x to 5.2x while marginally improving accuracy.

4.7 Case Study: High-Frequency Trading Scenario

To demonstrate the practical impact of C-S3G, we simulate a high-frequency trading scenario where predictions must be made within 10ms to be actionable. The original S3G model with 45.2ms latency cannot meet this requirement. The C-S3G model with 5.8ms latency easily satisfies the latency constraint while maintaining 96.2% of the original accuracy. In a backtesting evaluation over the test period (2023-2024), a trading strategy based on C-S3G predictions achieves a Sharpe ratio of 1.34, compared to 1.41 for the S3G-based strategy, representing a modest 5% reduction in risk-adjusted returns.

5. Discussion

5.1 Trade-offs Between Compression and Accuracy

Our experimental results demonstrate that knowledge distillation and model compression can achieve substantial computational savings with manageable accuracy degradation. The C-S3G model achieves 8.7x compression with only 3.8% relative accuracy reduction, representing a favorable trade-off for many practical applications. The key insight is that the graph-structured distillation and temporal knowledge transfer objectives effectively preserve the teacher's decision-making capabilities even as the model's capacity is reduced.

The bias decoupling regularization plays a particularly important role in maintaining accuracy during compression. By encouraging the student to focus on genuinely predictive features rather than spurious correlations, the regularization prevents the student from learning brittle shortcuts that would fail on out-of-distribution examples. This is especially important for financial prediction, where market conditions evolve and models must generalize across regimes.

5.2 Deployment Considerations

The compressed C-S3G model opens up deployment scenarios that were infeasible with the original S3G model. The 5.8ms inference latency enables integration into high-frequency trading systems with strict latency requirements. The 14.8MB model size fits comfortably within the memory constraints of mobile devices, enabling on-device stock prediction in mobile trading applications. The reduced computational requirements lower the infrastructure costs associated

with deploying AI-powered trading systems, potentially democratizing access to sophisticated prediction technology.

5.3 Limitations and Future Directions

Several limitations of the present work suggest directions for future investigation. First, the compression pipeline requires access to the teacher model's intermediate representations, which may not always be available (e.g., when compressing third-party models). Exploring distillation methods that require only the teacher's prediction outputs would broaden the applicability of our approach. Second, the current framework compresses a fixed teacher architecture; future work could explore co-design of teacher and student architectures to maximize the efficiency of the knowledge transfer process. Third, the adaptive quantization scheme requires a calibration dataset and may not adapt well to domain shift; online adaptation methods that update quantization parameters as market conditions evolve would address this limitation.

6. Conclusion

This paper presented C-S3G, a comprehensive knowledge distillation and model compression framework for graph-based financial prediction models. By combining graph-structured distillation that transfers relational knowledge, temporal-aware knowledge transfer that preserves market dynamics understanding, bias decoupling regularization that improves generalization, and adaptive quantization that achieves additional compression, the C-S3G model achieves 8.7x compression with only 3.8% accuracy degradation. The compressed model achieves real-time inference on commodity hardware, enabling deployment in high-frequency trading systems and mobile trading applications that were previously infeasible with the full S3G model.

Our work demonstrates that the principles of knowledge distillation, originally developed for classification and recognition tasks, can be effectively adapted to graph-based temporal prediction models in the financial domain. The integration of explanation-based regularization principles from natural language processing provides a principled approach to guiding the knowledge transfer process, ensuring that the compressed model preserves not just prediction accuracy but also robust generalization capabilities.

As the deployment of AI in financial markets continues to expand, the ability to compress large accurate models into efficient deployable form becomes increasingly critical. The C-S3G framework contributes to this goal, bridging the gap between predictive accuracy and computational efficiency and enabling the practical deployment of sophisticated deep learning models in resource-constrained financial environments.

References

1. Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural network. In NIPS Workshop on Deep Learning and Representation Learning.
2. Zang, J., & Liu, H. (2024, June). Explanation based bias decoupling regularization for natural language inference. In 2024 International Joint Conference on Neural Networks (IJCNN) (pp. 1-8). IEEE.
3. Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., & Bengio, Y. (2015). FitNets: Hints for thin deep nets. In International Conference on Learning Representations (ICLR).
4. Zagoruyko, S., & Komodakis, N. (2016). Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In International Conference on Learning Representations (ICLR).

5. Lu, Y., Hu, K., & Zhang, L. (2026, May). S3G: Stock State Space Graph for Enhanced Stock Trend Prediction. In ICASSP 2026-2026 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 4081-4085). IEEE.
6. Li, H., Kadav, A., Durdanovic, I., Samet, H., & Graf, H. P. (2017). Pruning filters for efficient convnets. In International Conference on Learning Representations (ICLR).
7. Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., Adam, H., & Kalenichenko, D. (2018). Quantization and training of neural networks for efficient integer-arithmetic-only inference. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 2704-2713). IEEE.
8. You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Keutzer, K., & Hsieh, C. J. (2020). Large batch optimization for deep learning: Training BERT in 76 minutes. In International Conference on Learning Representations (ICLR).
9. Frankle, J., & Carbin, M. (2019). The lottery ticket hypothesis: Finding sparse, trainable neural networks. In International Conference on Learning Representations (ICLR).
10. Wu, H., Chen, R., & Wang, L. (2023). Temporal graph attention networks for stock prediction. In Proceedings of the 29th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (pp. 2985-2994). ACM.
11. Yun, S., Jeong, M., Kim, R., Kang, J., & Kim, H. J. (2019). Graph transformer networks. In Advances in Neural Information Processing Systems 32 (NeurIPS) (pp. 11983-11993). Curran Associates.
12. Chen, T., Kornblith, S., Norouzi, M., & Hinton, G. (2020). A simple framework for contrastive learning of visual representations. In International Conference on Machine Learning (ICML) (pp. 1597-1607). PMLR.
13. Hassibi, B., & Stork, D. G. (1993). Second order derivatives for network pruning: Optimal brain surgeon. In Advances in Neural Information Processing Systems 5 (NeurIPS) (pp. 164-171). Morgan Kaufmann.
14. Liu, Z., Oguiza, J., & Neumann, M. (2024). Adapter-based fine-tuning for graph neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(3), 1428-1441.
15. Zhou, J., Cui, G., Hu, S., Zhang, Z., Yang, C., Liu, Z., & Sun, M. (2020). Graph neural networks: A survey of methods and applications. *AI Open*, 1, 57-81.