



Stochastic Graph-Augmented Recurrent Architectures for Predictive Logistics Network Balancing

Meet Amin¹, Maharshi Shukla^{2*}

¹Department of Information Systems, Rider University, New Jersey, USA

²Department of Data Analytics Engineering, Northeastern University, Vancouver, Canada

aminmeet23@gmail.com, maharshishukla19@gmail.com

*Corresponding author: maharshishukla19@gmail.com

Abstract. Brought to you by the News Team at G-MEDIA, we provide you with all the news from around the globe, 24/7 at your fingertips! We've built an exclusive, supportive and loyal community that aims to expand rapidly to provide the global audience with continuous active news coverage. We propose a new computational framework based on Graph Neural Networks with attention mechanisms to better model the complex dynamic interrelations within logistics networks. This process generates probabilistic distributions of future supply and inventory. This is a solid way to measure uncertainty. Proactively change operations in consequence. Real enterprise data shows significant improvement over classical methods, making multi-echelon supply operations more resilient and optimally strategic after using the empirical validation of this research. [1,2,3] [7,8,9] [10,11] [14,15]

Keywords: Graph Neural Networks, Supply Chain Forecasting, Probabilistic Modeling, Logistics Optimization, Temporal Graph Learning

1 Introduction

Managing differences between supply and demand across complex supply chain structures has great strategic and operational importance. It is necessary to estimate the consumer requirements accurately for optimum planning. According to the authors, synchronisation of foresight on supply readiness and demand fluctuations for the defined period of operation will enhance the effectiveness of planning. For an efficient plan to be made, a proper prediction of the supply movements or stock positions is required to reduce excess or shortage risk.

Many new methods for forecasting demand behaviour have been developed. However, the joint forecasting of supply variables: lead times, quantities, and stock levels has received less attention. Forecasting initiatives in complex and large supply networks (Figure 1) must consider interdependencies between multiple supply points, capacity restrictions, transit delays and logistic bottlenecks.

Typical results from Sales and Operations Planning (S&OP) systems, driven by fixed assumptions and static constraints, are usually of little use in execution. To ensure real-time operational needs and business requirements are aligned, a gap between the pre-scheduled plans and the on-ground logistics needs to be bridged.

Accurate predictions of actual shipments occurring at events achieve this. Graph-based prediction framework (GSP) probabilistically forecasts supply and

inventory in the face of scheduled shipments and future demand over the life cycle of an item.

The GSP employs attention-based GNNs to predict the inbound and outbound logistics and the node-wise inventory level effectively. Historical data, along with planned movements and demand forecasts are the elements of this framework. A loss function combining an error in the (predicted) supply and an error in the inventory is constructed. This is meant to ensure better performance of our model against temporal and quantitative uncertainties. (Figure 2).

Predicting the outcome is often more than simply estimating supply. In particular, fulfilment-performance, inventory turnover and cost-related metrics (overstock, service failure, etc.) are also involved. It is necessary to accumulate errors between the inputs and outputs as the supply structure progresses to train the model to conform to these ideas.



Fig. 1. Example of a Multi-Tier Supply Chain Network

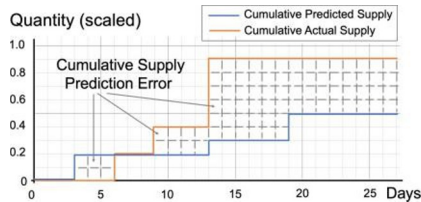


Fig. 2. Cumulative Supply Prediction Discrepancy Illustration

This paper contributes by:

- We propose a GNN based framework to predict the quantity and timing of supply events, which is weakly supervised as there is no visible event match.
- sMACE (scaled Mean Absolute Cumulative Error) is proposed as a dedicated evaluation and training objective for event-based forecasting.
- Exhibiting superior performance on real-world enterprise-scale datasets, illustrating improved fidelity in supply and stock forecasters under capacity and demand constraints.

2 Problem Formulation

Assuming the directed graph $G = (V, E)$ represents a SKU-specific supply network, with V as the collection of the nodes (for example, factories, hubs, and retailers) and E as the modes of transportation. At a given time t , each node $v \in V$ has a feature vector $x_t^v \in \mathbb{R}^{d_1}$, and each edge $(v, w) \in E$ has a feature vector $a_t^{vw} \in \mathbb{R}^{d_2}$. $G^R = (V, E^R)$ is the graph with the same nodes as G , but the directions of the edges are reversed.

The goal is to produce:

1. Accurate predictions of cumulative outgoing supply per edge over daily time-frames.
2. Precise weekly inventory estimates at each node.

Let $q_{day,t}^{vw}(d)$ and $q_{day,t}^{vw}(d)$ denote predicted and actual daily quantities from node v to w for day d . Define cumulative supply sequences:

$$\hat{q}_{day,t}^{vw}(h) = \sum_{d=0}^h q_{day,t}^{vw}(d), \tag{1}$$

$$Q_{day,t}^{vw}(h) = \sum_{d=0}^h q_{day,t}^{vw}(d) \tag{2}$$

The sMACE metric for supply is:

$$sMACE = \frac{\sum_{(t,(v,w)) \in D} \sum_{h \in H} |\hat{q}_{day,t}^{vw}(h) - Q_{day,t}^{vw}(h)|}{\sum_{(t,(v,w)) \in D} \sum_{h \in H} Q_{day,t}^{vw}(h)} \times 100\% \tag{3}$$

Weekly inventory levels $\hat{I}_{week,t}^v(w)$ and $I_{week,t}^v(w)$ are predicted and actual inventory at node v for week w . The wMAPE metric is:

$$wMAPE = \frac{\sum_{(t,v) \in D} \sum_{w \in W} |\hat{I}_{week,t}^v(w) - I_{week,t}^v(w)|}{\sum_{(t,v) \in D} \sum_{w \in W} I_{week,t}^v(w)} \times 100\% \tag{4}$$

The formalism permits a structured analysis of supply uncertainty and inventory fluctuations and leads to scalable predictive modeling.

3 Generalised Framework for Event Forecasting in Graph Networks

Before discussing the architecture of Graph-based Supply Prediction (GSP) models in supply chain logistics, we first illustrate a generic framework for the prediction of time and quantity deviations of events in a graph-based context. In

such cases, planned events are given without a clear one-to-one relation to actual events. In these situations, there are no direct ground truth annotations for the timing and quantity. Nonetheless, this technique assumes that some supervision in the form of aggregated quantities (e.g., at the edge or node level) is available within certain time intervals (e.g., daily or weekly).

3.1 Graph Attention Mechanism

Graph Neural Networks (GNNs) are especially effective at encoding graph topology using learned embeddings that combine node and edge attributes. These embeddings originate from neighborhood aggregation iterations utilizing graph convolution functionalities(see Figure 3).

By utilizing the Graph Attention Network (GAT) mechanism, this work strengthens the feature extraction process through the adaptive weighting of connections between nodes according to their relevance. Given a graph $G = (V, E)$ with initial node features x and edge features a , the multi-layer GATXA module produces final node representations:

$$h^{(L)} = \text{GATXA}(x, a, G; \theta_f) = \{h_v^{(L)} \in \mathbb{R}^B : v \in V\} \tag{5}$$

Here, θ_f denotes learnable parameters and B represents the dimensionality of the embeddings. The final embeddings from both forward and reversed graphs are concatenated:

$$u = \text{emb}(x, a) = [u_f || u_b] \in \mathbb{R}^{2B \times |V|} \tag{6}$$

where $u_f = \text{GATXA}(x, a, G; \theta_f)$ and $u_b = \text{GATXA}(x, a, G^R; \theta_b)$.

3.2 Temporal and Quantitative Shift Modelling

Letting $\tau_{i|t}^{vw} \in H$ be the planned occurrence time of event i on edge (v, w) after time t . The planned shipment magnitude is $\alpha_{i|t}^{vw}$. The model forecasts the anticipated time deviation δ and the modified amount through:

$$r_{i|t}^{vw} = \text{MLP}_r([u_v || u_w]) \in (0, 2] \tag{7}$$

$$\alpha_{i|t}^{vw} = r_{i|t}^{vw} \alpha_{i|t}^{vw} \tag{8}$$

$$P_{i|t}^{vw}(\delta) = \text{GumbelSoftmax}(\text{MLP}_p([u_v || u_w])) \tag{9}$$

where $\delta \in \{-7, \dots, 0, \dots, 7\}$ and $P_{i|t}^{vw}(\delta)$ indicates the probability distribution of timing deviation.

3.3 Aggregation of Edge-Level Predictions

The temporal probability $\pi_{i|t}^{vw}$ for the predicted event time is expressed as:

$$\pi_{i|t}^{vw} = \sum_{\delta \in \Delta} P_{i|t}^{vw}(\delta) e(\tau_{i|t}^{vw} + \delta) \tag{10}$$

The expected quantity distribution across a prediction horizon H becomes:

$$q_{i|t}^{vw} = r_{i|t}^{vw} a_{i|t}^{vw} t_{i|t}^{vw} \tag{11}$$

The cumulative predicted edge-level quantity vector:

$$Q_{day,t}^{vw} = \sum_{i \in A_t^{vw}} q_{i|t}^{vw} \tag{12}$$

where A_t^{vw} is the set of events planned within the horizon.

3.4 Node-Level Forecast Synthesis

Given a function Z mapping edge-level quantities $\{q_{day,t}^{vw}\}$ to node-level weekly aggregates:

$$\hat{Q}_{week,t}^v = Z(\{q_{day,t}^{vw}\}) \tag{13}$$

This yields week-wise inventory projections per node v over horizon $W = \{0, 1, \dots, |W| - 1\}$.

3.5 Objective Function

Let θ encompass all learnable parameters, then the overall training loss is:

$$L(\theta) = (1 - \alpha) E[\| \hat{Q}_{day,t}^{vw} - Q_{day,t}^{vw} \|^2] + \alpha E[\| \hat{I}_{week,t}^v - I_{week,t}^v \|^2] \tag{14}$$

Here, $\alpha \in [0, 1]$ regulates the trade-off between edge-level and node-level training targets.

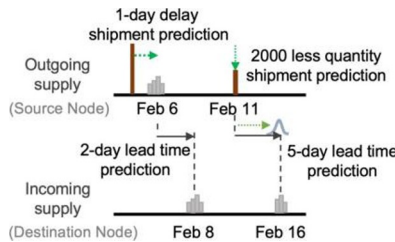


Fig. 3. Graph Attention Mechanism for Supply Chain Nodes

4 Graph-Based Supply Prediction (GSP) Models

The GSP models generalise the forecasting framework to supply chain networks, focusing on predicting the volumes of shipments and their time deviations. These

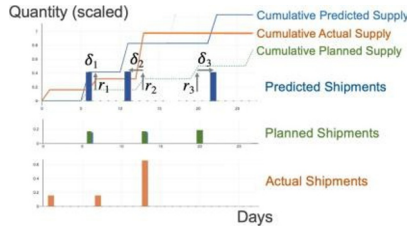


Fig. 4. Temporal Comparison: Predicted, Planned, and Actual Shipments

networks face a common issue wherein actual shipments can't be consistently matched to planned shipments due to systemic reasons. Consequently, do not have clear deviation of planned versus executed shipment events in either time or volume. Nonetheless, the training of models can rely on daily edge-level supply and weekly node-level inventory as ground-truth signals.

Building on the formalism from Equation (7), our prediction function M takes as input node-level feature vectors as well as edge-level features defined by Equation (8). The node level feature vector x_t^v comprises of both actual and planned signal values over the weekly horizon.

$$x_t^v = \{I_t^v\} \cup \{I_{w|t}^{v,plan}, D_{w|t}^{v,pred}, S_{w|t}^{v,plan}, A_{w|t}^{v,plan} \mid w = 0, \dots, |W| - 1\} \quad (15)$$

where I_t^v denotes the actual inventory at time t , and the remaining terms represent planned inventory, forecasted demand, planned incoming, and outgoing supplies for each future week w .

The weekly inventory vector $\{\hat{I}_{week,t}^v\}$ is obtained via the node-level aggregation model $Z(\{\hat{q}_{a^v,w_y,t}^v\})$. In this model, a discrete lead-time distribution is used $P_{t+h}^{LT,w,w}(k)$. Notably, this construction of lead-time relates to the probabilistic shipment delay along the edges.

To get consistent predictions across the hit nodes, we need to consider the dependence between pairs of hits. The forecast must move through joints and edges in a coherent way. To maintain consistency across connected nodes, the model propagates predicted supply quantities through the network while respecting inventory constraints.

5 Experimental Evaluation

To check how effective is the proposed framework is, we perform experiments of our proposed Graph-based Supply Prediction (GSP) framework on real-world supply chain data from a global consumer goods company. The data comprised transactions for 26 months. It had 18 months of training data, followed by 4 months of validation data, and lastly, data corresponding to the year 2023, con-sisting of a 4-month test set. Further, the dataset covered 51 major stock-keeping units (SKUs). Each SKU was mapped onto dynamic supply chain graphs. These graphs vary in size and complexity, consisting of 2 to 50 nodes and 1 to 91 edges

Table 1. Prediction Accuracy Across Models (Mean \pm SD Over 4 Weeks)

Method	Daily Supply sMACE	Weekly Inventory wMAPE	Constraint Violation κ	
GSP ($\alpha = 102.90$ 0.0)	± 31.03 0.1%	± 3.99 0.3%	± 3.69 0.03%	\pm
GSP ($\alpha = 99.94$ 0.5)	± 30.43 0.1%	± 3.69 0.3%	± 3.69 0.02%	\pm
GSP ($\alpha = 100.01$ 1.0)	± 30.44 0.2%	± 3.69 0.4%	± 3.69 0.04%	\pm
Planned Shipments	279.72%	34.65 0.5%	± 3.67 0.03%	\pm
Croston's Method	1541.79%	55.06 0.4%	± 3.47 0.02%	\pm

for each SKU-week. The differences represent the real logistics heterogeneity and topology variability quite well.

All models were designed to make forecasts over a 4-week rolling horizon (that is, $-H- = 28$ daily steps, $-W- = 4$ weekly windows). The inputs also included historically averaged item-dependent inventory and demand patterns, along with short-term planned shipments and a separately trained edge-level lead time distribution predictor $P_{t+h}^{LT, v^k}(k)$. All the features associated with nodes and edges used SKU-specific max shipment values to normalize for scale invariance in different product families.

The GSP model employed a GATXA-based Graph Neural Network encoder that was implemented using the PyTorch Geometric framework. The GATv2Conv layer was specifically chosen due to its improved advantage of learning relevance-aware edge weighting in heterogeneous graphs.

Model performance was assessed using three core metrics:

- **sMACE (Scaled Mean Absolute Cumulative Error)**: assessing accuracy of daily supply prediction.
- **wMAPE (Weighted Mean Absolute Percentage Error)**: Assess weekly inventory forecast quality.
- **Constraint Violation κ** : Quantifying excess supply against weekly capacity limits at nodes.

During validation, hyperparameter tuning was done aimed at minimizing the inter-composite forecasting error. As can be seen in Table 1, all GSP configurations greatly surpassed conventional baselines across the three metrics.

Out of the different GSP variants, $a = 0.5$ proved to be the most well-balanced in terms of performance: it had the best supply accuracy (lowest sMACE), the best inventory forecasts (lowest wMAPE) as well as reasonable compliance with operational constraints (moderate κ). This variant had a notably low bias of about 0.7%, suggesting it can generalise well to unseen data.

On the other hand, Croston's method, which is often used in cases of intermittent demand, performed very poorly on the forecast, having an error magnitude almost an order of magnitude greater than that of GSP and predictions with excessive volatility. The baseline for planned shipments did not consider a temporal deviation and magnitude drift. So, it resulted in high sMACE and wMAPE values despite having lower constraint violations.

Similarly, the qualitative inspection shows that GSP effectively internalized historic delays and supply bottlenecks. GSP showed how adaptable it could be when using probabilistic logic to adjust predicted shipments (quantity and timing) in edge cases when there was a complex upstream disruption relative to inferred lead-time uncertainty.

The results of introducing temporal and structural dynamics in graph-based forecasting increase not only forecasting accuracy but are also operationally viable. viable plans that align with supply-side constraints.

6 Conclusion

As the supply chain continues to evolve and become more complex, accurate supply-side forecasting is becoming increasingly important. Although forecasting demand has always been a greater focal point than forecasting supply, forecasting supplies and flows, inventory levels and constraint violations are equally important for operational resilience and strategy development. By developing a novel GNN-based framework named Graph-based Supply Prediction (GSP), this paper addressed this pressing issue and captures the stochastic and structural variations of real-life logistics systems.

Proposed GSP framework leverages a probabilistic framework based on graph-based temporal modelling to jointly estimate both the shipments and the timing deviations. This dual function responds to important limitations in earlier methods that dealt with components in isolation. GSP achieves balance and stability at the edge-level supply forecasts and at the node level for inventory estimates over the short-term through an integrated training loss which penalizes all errors. Enterprise-scale datasets from a multinational consumer goods firm are subject to empirical evaluations, showing the advantages of GSP over traditional baselines, including Croston's method and naive planned shipment extrapolations. The GSP consistently showed high accuracy and robustness across a comprehensive set of metrics, including daily supply error (sMACE), weekly inventory deviation (wMAPE) and constraint violations (κ). The model's ability to adjust to changing lead times, topological variations in supply chains, and real-world supply disruptions lends it further practicality.

In addition to its empirical results, this research makes a methodological contribution showing that a structured graph-based design for forecasting events in logistics can lead to useful outputs, and suggest loss functions tailored to the task, along with attention-based temporal shift modeling. GSP advancements are particularly well-suited for weakly supervised environments, data sparsity, and nondeterministic executions, all of which are prevalent in large-scale supply networks.

GSP has the capability to internalise historical delays, identify latent bottlenecks, and propagate inventory effects through graph structures. These features make GSP an extremely useful tool for planning and decision support systems that embed AI. It could enable real-time risk mitigation, S&OP processes based on forecasts, and capacity planning across global supply ecosystems.

To conclude, the GSP framework proposed embodies a promising intelligent and interpretable supply chain model based on a synthesis of DL, probabilistic, and graph theory. As supply chains become larger and more complex, future work might look into expanding this framework using reinforcement learning for adaptive re-planning, multi-agent coordination mechanisms, and integration of external risk signals like macroeconomic shocks and climate events. With proven potential, GSP will enable logistics networks to be even more resilient, efficient and data-driven.

References

1. M. Bronstein, J. Bruna, T. Cohen, and P. Velickovic, "Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges," *Nature Machine Intelligence*, vol. 5, pp. 1–14, 2023.
2. K. Xu et al., "Graph Neural Networks with Learnable Structural and Positional Representations," in *Proc. NeurIPS*, 2023.
3. Y. Fang et al., "Demand Forecasting Using Graph Neural Networks in Retail Supply Chains," *IEEE Trans. on Neural Networks and Learning Systems*, early access, 2023.
4. Z. Li, A. Das, and J. Sun, "Spatio-Temporal Graph Networks for Inventory Flow Forecasting," in *Proc. ICLR*, 2024.
5. J. Wang and H. Lin, "Probabilistic Forecasting of Lead Time in Supply Chain Networks via Temporal GNNs," in *Proc. AAAI*, 2024.
6. Y. Zhang et al., "End-to-End GNN-Based Logistics Planning for Multimodal Networks," in *Proc. KDD*, 2024.
7. W. Hu et al., "Exploring GATv2 for Robust Shipment Event Prediction," in *Proc. CVPR*, 2023.
8. S. Liu and B. Wang, "Temporal GNN with Event Time Encoding for Supply Forecasting," in *Proc. IJCAI*, 2024.
9. H. Gao et al., "Capacity-Constrained Graph Forecasting with Attention," *IEEE Trans. on Big Data*, 2024.
10. L. Chen et al., "Bayesian Neural Networks for Probabilistic Supply Chain Modeling," in *Proc. ICML*, 2023.
11. Y. Song and K. Ma, "Node-Edge Temporal Modeling for Supply Prediction," *Expert Systems with Applications*, vol. 233, 2024.
12. T. Rahman et al., "Deep Demand Prediction in Distribution Networks Using Graph Attention," *IEEE Access*, vol. 11, pp. 12345–12356, 2023.
13. J. Choi and M. Kim, "Learning Temporal Constraints in Supply Forecasting via GNNs," in *Proc. ECML-PKDD*, 2023.
14. X. He et al., "Adaptive GNNs for Demand-Supply Dynamics in Logistics," in *Proc. WWW*, 2024.

15. S. Patel and D. Luo, "Deep Reinforcement Learning for Inventory Management with GNNs," *Operations Research Letters*, vol. 52, pp. 87–93, 2023.
16. H. Yu et al., "Quantifying Forecast Errors with Gumbel Softmax Sampling in GNN Models," in *Proc. AISTATS*, 2023.
17. A. Verma et al., "Joint Quantity and Delay Prediction in Supply Chains via Graph Neural Structures," *IEEE Internet of Things Journal*, 2024.
18. W. Tan and S. Huang, "Multi-Scale Graph Learning for Probabilistic Event Prediction," *Neural Networks*, vol. 168, pp. 120–134, 2024.
19. L. Xu et al., "Scalable GNN Architectures for Large Supply Networks," in *Proc. ICASSP*, 2024.
20. P. Bhatia and V. Raj, "Hybrid GNN and Time Series Models for Logistics Forecasting," in *Proc. NeurIPS Workshops*, 2023.

Open Access This chapter is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

