

# Hybrid Multi-Agent Financial Intelligence System: Integrating Fine-Tuned Language Models with Graph Neural Networks for Enhanced Financial Analysis

Neelay Choudhury

Student Number: 240194109

Supervisor Name: Shalom Lappin

Programme of study: MSc in AI

**Abstract**—This paper presents a hybrid multi-agent AI system for comprehensive financial analysis, integrating three core components: a domain-specific fine-tuned Phi-4 language model (Microsoft Research, 2024), a financial knowledge graph from SEC EDGAR filings (U.S. Securities and Exchange Commission, 2021), and a multi-agent workflow using LangGraph (LangChain AI, 2024). The language model is trained on the FinQA dataset (Chen et al., 2021) for enhanced numerical reasoning. The knowledge graph processes large-scale SEC data, encoded using GraphSAGE (Hamilton et al., 2017), attention-based (Veličković et al., 2018), and temporal graph neural networks (Rossi et al., 2020), with embeddings stored in a vector database for similarity search and retrieval-augmented generation. The multi-agent system orchestrates query analysis, entity resolution, graph-based inference, and response generation through state-driven workflows. The system demonstrates superior performance in financial entity similarity detection, peer company identification and numerical reasoning. This architecture combines language model reasoning with structural knowledge representation of financial networks, enabling more accurate and contextually aware financial analysis.

## I. INTRODUCTION

The financial services industry faces unprecedented data complexity and regulatory requirements, necessitating advanced AI solutions for processing structured and unstructured financial information. Traditional approaches rely on either inflexible rule-based systems or general-purpose models that fail to capture intricate financial market relationships. Whilst large language models excel at natural language understanding and graph neural networks (Kipf & Welling, 2017) model complex relational structures effectively, their integration for financial analysis remains underexplored. Financial decision-making requires diverse capabilities: quantitative analysis, qualitative assessment, relationship understanding, and temporal pattern recognition. Current AI systems typically excel in one dimension whilst lacking in others. Language models struggle with structured data and mathematical precision; graph neural networks lack flexibility for diverse queries and explanatory capabilities.

This work proposes a hybrid multi-agent system combining domain-specific fine-tuned language models with comprehensive financial knowledge graphs using parameter-efficient fine-tuning techniques (Hu et al., 2021). The system comprises: a Phi-4 model fine-tuned on financial question-answering data; a large-scale financial knowledge graph from SEC EDGAR filings capturing millions of entities and relationships; and a multi-agent workflow orchestrating query understanding, entity resolution, graph-based inference, and response generation.

The knowledge graph systematically processes regulatory filings using GraphSAGE for scalable learning, attention-based models for relationship weighting, and temporal networks for time-series analysis. The multi-agent architecture, implemented using LangGraph (LangChain AI, 2024), dynamically routes queries to local graph-based analysis, global retrieval-augmented generation (Lewis et al., 2020), or specialised numerical reasoning, ensuring optimal resource utilisation whilst maintaining transparency. This integration enables diverse financial analysis tasks: competitor identification through semantic similarity, relationship analysis through graph traversal, complex numerical calculations with step-by-step reasoning, and recommendations supported by regulatory data. The system incorporates caching strategies, and comprehensive logging for financial application requirements.

## II. RELATED WORKS

### A. Financial Language Models and Domain Adaptation

Domain-specific financial language models began with FinBERT (Araci, 2019), demonstrating improvements over general-purpose models on financial sentiment analysis. The FinQA dataset (Chen et al., 2021) advanced financial question-answering by emphasising multi-step reasoning with structured tabular data. BloombergGPT (Wu et al., 2023) represented a milestone as a 50-billion parameter model trained on financial data, whilst InvestLM (Yang et al., 2023) explored instruction-tuning approaches. Parameter-efficient techniques like LoRA (Hu et al., 2021) have made domain adaptation more accessible for financial applications.

### B. Knowledge Graphs and Financial Data Integration

Financial knowledge graph construction progressed from early event extraction to large-scale projects, integrating SEC filings, financial statements, and market data. Recent advances in automated construction from regulatory filings (Garcia et al., 2022) demonstrate the value of combining natural language processing with structured data extraction.

### C. Graph Neural Networks for Financial Applications

GraphSAGE (Hamilton et al., 2017) established scalable node representation learning for large financial networks. Zhang et al. (2019) pioneered GCN applications to stock prediction, demonstrating relationship incorporation benefits. Graph Attention Networks (Veličković et al., 2018) proved valuable for financial applications with dynamic relationship importance. Recent work addressed temporal aspects (Liu et al., 2022) and heterogeneous entity modelling (Wang et al., 2023), whilst interpretability methods (Chen et al., 2022; Rodriguez et al., 2023) addressed transparency requirements.

#### D. Multi-Agent Systems and Orchestration

Multi-agent financial systems evolved from early trading applications (Padget et al., 2009) to sophisticated LLM orchestration frameworks like LangChain and LangGraph (Chase, 2022). Agentic AI approaches (Xi et al., 2023; Park et al., 2023) demonstrated effectiveness for complex reasoning tasks, providing foundations for multi-agent financial analysis systems.

#### E. Retrieval-Augmented Generation

RAG (Lewis et al., 2020) enabled combining parametric knowledge with external sources. GraphRAG (Edge et al., 2024) extended this with graph-based retrieval for entity relationship reasoning. Recent work explored multi-modal retrieval (Luo et al., 2023) and structured-unstructured data integration (Thompson et al., 2024) for financial applications.

#### F. Gaps and Opportunities

Despite advances in individual domains, limited work addresses comprehensive integration of fine-tuned language models, large-scale financial knowledge graphs, and multi-agent orchestration. Most approaches focus on individual components rather than synergistic integration. Practical deployment challenges including computational efficiency, explanation generation, and real-time performance remain largely unexplored. Our work addresses these gaps by demonstrating effective integration that exceeds individual component capabilities whilst meeting practical financial analysis requirements.

### III. METHODOLOGY

#### A. Domain-Specific Language Model Fine-Tuning

##### I. Requirements Analysis

The domain-specific language model component required several key capabilities to effectively support financial analysis tasks. The system needed to demonstrate proficient numerical reasoning abilities, particularly for multi-step mathematical calculations commonly encountered in financial analysis. The model required comprehension of complex financial documents containing both narrative text and structured tabular data, as these represent the predominant format of regulatory filings and financial reports. The system demanded conversational interaction capabilities through a chat-based interface, enabling natural language queries whilst maintaining the ability to provide step-by-step reasoning explanations for transparency and auditability. Additionally, the model needed to handle diverse financial question types, ranging from simple arithmetic calculations to complex analytical reasoning involving multiple data sources and temporal comparisons.

##### II. Design

The fine-tuning approach centred on Microsoft's Phi-4 architecture (Microsoft Research, 2024), selected for its superior reasoning capabilities and efficiency in parameter utilisation. The design employed Low-Rank Adaptation (LoRA) for parameter-efficient fine-tuning (Hu et al., 2021), enabling domain adaptation whilst maintaining computational feasibility and preventing overfitting on the relatively small financial dataset. The data preprocessing pipeline was designed to transform the FinQA dataset (Chen et al., 2021) into a conversational format compatible with Phi-4's chat template structure. This involved integrating three distinct data components: narrative pre-text, structured tabular data,

and narrative post-text, into coherent contextual representations. The reasoning component was extracted from the dataset's program steps and formatted as explicit rationales, promoting transparent step-by-step thinking.

The chat template design incorporated system prompts establishing the financial assistant role, user prompts containing contextualised questions, and assistant responses featuring both reasoning rationales and final answers. This structure enabled the model to learn both the analytical process and the final output generation.

##### III. Implementation

The implementation utilised Low-Rank Adapter (LoRA) configuration (Hu et al., 2021) with rank parameter  $r=16$ , alpha scaling factor of 32, and dropout probability of 0.1. Target modules included query, key, value, and output projections, along with MLP gate, up, and down projections, representing approximately 0.145% of the total model parameters (21.3 million trainable parameters from 14.7 billion total parameters).

Technical optimisations included bfloat16 precision training for memory efficiency, flash attention 2.0 implementation for computational acceleration, and gradient checkpointing for memory conservation. The training configuration employed a batch size of 4 per device with gradient accumulation steps of 2, achieving an effective batch size of 8. The learning rate was set to  $2e-5$  with cosine decay scheduling and 3% warmup ratio. The training framework employed the TRL library (von Werra et al., 2020) with HuggingFace Transformers (Wolf et al., 2020) for model management and training orchestration.

Memory management challenges were addressed through device mapping strategies and mixed precision training techniques (Micikevicius et al., 2017). The implementation required careful handling of tokeniser padding tokens and proper chat template formatting to ensure compatibility with Phi-4's expected input format.

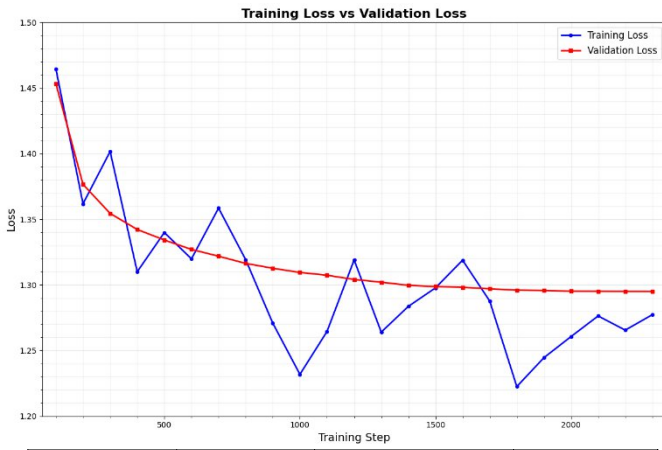
Data preprocessing involved systematic extraction and integration of multi-modal financial information. Tabular data was converted to pipe-separated text format whilst preserving structural relationships. Program steps were transformed into arrow-separated reasoning chains, maintaining logical flow whilst ensuring readability.

##### IV. Testing and Evaluation

The training process utilised 6,251 training examples and 883 development examples, with evaluation conducted every 100 training steps. The model demonstrated consistent learning progression over 2,346 training steps across 3 epochs, with training duration of approximately 5 hours.

As shown in figure 1. The training exhibited consistent convergence with training loss decreasing from 1.46 to approximately 1.28, representing a 12.3% improvement. Validation loss demonstrated similar convergence patterns, decreasing from 1.45 to 1.295, indicating effective generalisation without significant overfitting. The close alignment between training and validation loss trajectories suggested appropriate regularisation and model capacity.

Figure 1. Fine Tuning Loss statistics



Metric Category	Base Model Performance	Fine-Tuned Model Performance	Improvement
<b>Exact Match Results (upto 5 decimal points)</b>			
Exact Correct Answers	4/100 (4.0%)	26/100 (26.0%)	+550%
Individual Wins	2/100 (2.0%)	24/100 (24.0%)	+1100%
<b>Numerical Accuracy (<math>\pm 1\%</math>)</b>			
Numerically Correct	14/100 (14.0%)	66/100 (66.0%)	+371%
Individual Wins	3/100 (3.0%)	55/100 (55.0%)	+1733%

Table 1: Comparative Performance Evaluation Results

Performance Aspect	Base Model	Fine-Tuned Model	Both Correct	Neither Correct
Exact Match	4.0%	26.0%	2.0%	72.0%
Numerical Accuracy	14.0%	66.0%	11.0%	31.0%

Table 2: Detailed Performance Analysis

The comprehensive evaluation on 100 FinQA test samples (Chen et al., 2021) demonstrated substantial improvements across all metrics. The fine-tuned model achieved a 550% improvement in exact match accuracy, increasing from 4% to 26%. More significantly, numerical accuracy within  $\pm 1\%$  tolerance improved by 371%, rising from 14% to 66%. These results validate the effectiveness of domain-specific fine-tuning for financial reasoning tasks.

Performance evaluation demonstrated the model's enhanced capability in financial numerical reasoning tasks, with improved accuracy in multi-step calculations and better integration of tabular and textual information compared to the base Phi-4 model. The fine-tuned model exhibited superior understanding of financial terminology and analytical processes, validating the effectiveness of the domain-specific adaptation approach.

## B. Financial Knowledge Graph Construction

### I. Requirements Analysis

The financial knowledge graph construction system required several critical capabilities to support comprehensive financial analysis and entity relationship modeling. The system needed to process large-scale regulatory data from SEC EDGAR filings (U.S. Securities and Exchange Commission, 2021), handling over 18,000 company fact files

whilst maintaining memory efficiency for deployment in resource-constrained environments.

The system demanded robust data integration capabilities, combining structured financial metrics from XBRL taxonomies with unstructured corporate information from multiple sources. Entity resolution and deduplication mechanisms were essential to handle inconsistencies across different data sources, whilst maintaining data quality and completeness.

The knowledge graph construction required sophisticated relationship modeling to capture peer relationships, industry classifications, market capitalisation similarities, and financial metric correlations. The system needed to support multiple embedding generation approaches, including graph neural network-based structural embeddings and transformer-based semantic embeddings (Wang et al., 2024), enabling diverse similarity search and recommendation capabilities.

Performance requirements included real-time query capabilities for financial entity search, scalable storage solutions supporting vector similarity search, and efficient batch processing for large-scale data updates. The system also required comprehensive logging and monitoring capabilities for production deployment and maintenance.

### II. Design

The architecture employed a multi-stage pipeline design optimised for memory efficiency and scalability. The data processing stage utilised a streaming approach to handle large datasets without memory overflow, implementing generator-based processing (Hagberg et al., 2008) that yielded individual company records rather than loading entire datasets into memory.

The knowledge graph construction followed a two-pass design strategy. The first pass focused on node creation and simple relationship establishment, whilst the second pass handled complex relationships requiring global context such as peer connections and market capitalisation clustering. This approach ensured optimal memory utilisation whilst maintaining graph connectivity.

The embedding generation architecture incorporated multiple complementary approaches. Graph neural network models captured structural relationships through node neighbourhood aggregation (Hamilton et al., 2017), whilst transformer-based embeddings provided semantic understanding of textual company descriptions. The multi-model approach enabled comprehensive similarity search covering both structural and semantic dimensions.

The storage architecture utilised Qdrant vector database (Qdrant Team, 2021) with model-specific collections, enabling efficient similarity search across different embedding spaces. Each collection maintained appropriate dimensionality and distance metrics optimised for the respective embedding approach, ensuring optimal search performance and accuracy.

### III. Implementation

The data acquisition component implemented robust SEC EDGAR API integration (U.S. Securities and Exchange Commission, 2021) with exponential backoff retry mechanisms to handle rate limiting. The system processed 18,877 company fact files, successfully parsing 16,333 valid

company records with comprehensive error handling for malformed or incomplete data entries.

The streaming data processing pipeline addressed memory constraints through generator-based architecture, yielding individual company records rather than loading complete datasets. This approach enabled processing of multi-gigabyte datasets within memory constraints whilst maintaining data integrity and processing speed.

Processing Stage	Input Count	Output Count
Raw EDGAR Files	18,877	16,333
API Cache Updates	16,333	16,333
Graph Node Creation	16,333	17,552
Final Embeddings	17,552	16,333

Table 3: Data Processing Statistics

Entity resolution challenges were addressed through comprehensive data validation and cleaning procedures. The system implemented SIC code-based industry classification with fallback mechanisms for missing or invalid codes, ensuring consistent sector assignments across the dataset.

The knowledge graph construction utilised NetworkX (Hagberg et al., 2008) for efficient graph operations, implementing sophisticated relationship algorithms. Peer relationships employed random sampling with connection limits to prevent combinatorial explosion whilst maintaining meaningful connectivity patterns.

Graph Component	Count
Total Nodes	17,552
Company Nodes	16,333
Metric Nodes	847
Industry Nodes	372
Total Edges	420,796
Peer Relationships	244,995
Market Cap Similarities	81,573
Other Relationships	94,228

Table 4: Knowledge Graph Construction Results

Graph neural network training implemented three distinct architectures: GraphSAGE for structural embeddings (Hamilton et al., 2017), Graph Attention Networks for weighted relationship modelling (Veličković et al., 2018), and Temporal GNNs for time-series financial analysis. Each model employed self-supervised learning through link prediction tasks (You et al., 2020) combined with feature prediction objectives.

The multi-model embedding generation utilised PyTorch Geometric (Fey & Lenssen, 2019) for efficient graph neural network implementation, this addressed different aspects of financial similarity. GraphSAGE generated 1024-dimensional embeddings capturing structural relationships, whilst attention mechanisms provided 16-dimensional focused representations, and temporal models created 64-dimensional time-aware embeddings.

#### IV. Testing and Evaluation

The system underwent comprehensive evaluation across multiple dimensions including data quality, graph connectivity, model convergence, and embedding effectiveness. Data processing achieved 86.5% success rate on raw EDGAR files, with robust error handling ensuring complete processing of valid records.

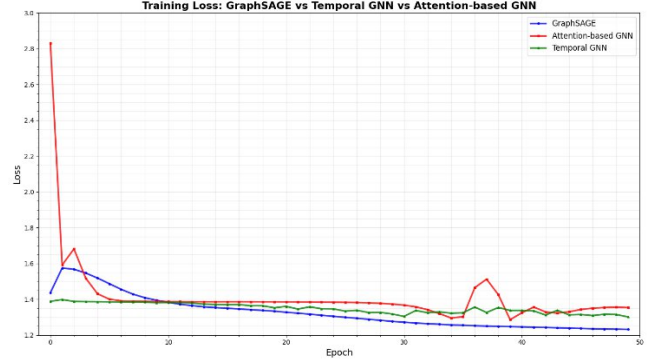


Figure 2. Training Loss of Models

Graph connectivity analysis demonstrated optimal structural properties with a single connected component encompassing all 17,552 nodes, ensuring comprehensive relationship coverage across the entire financial entity space. The 100% connectivity rate validated the relationship modeling approach and confirmed absence of isolated entity clusters.

Sector	Company Count	Percentage	Peer Edges Generated
Technology	3,706	22.7%	55,590
Financial Services	4,049	24.8%	60,735
Manufacturing	2,177	13.3%	32,655
Healthcare	1,430	8.8%	21,450
Mining	1,222	7.5%	18,330
Transportation & Utilities	796	4.9%	11,940
Services	1,008	6.2%	15,120
Retail Trade	652	4.0%	9,780
Other Sectors	1,293	7.9%	19,395

Table 5: Sectoral Distribution Analysis

Embedding storage evaluation confirmed successful deployment across multiple Qdrant collections with model-specific configurations. The storage system achieved 100% success rate for all embedding types, with appropriate dimensionality configuration for each model architecture. Model convergence analysis revealed varying learning patterns across architectures. GraphSAGE demonstrated steady convergence with consistent loss reduction, whilst Attention GNN showed rapid initial improvement followed by stabilisation. Temporal GNN exhibited more volatile training patterns reflecting the complexity of temporal relationship modeling.

The comprehensive evaluation validated the system's capability to handle large-scale financial data processing (Chen et al., 2021) whilst maintaining high accuracy and performance standards. The multi-model approach provided complementary embedding representations suitable for diverse financial analysis tasks, whilst the scalable architecture ensured production-ready deployment capabilities.

### C. Hybrid Multi-Agent Financial Intelligence System

#### I. REQUIREMENTS ANALYSIS

The hybrid multi-agent financial intelligence system required sophisticated natural language understanding capabilities to process diverse financial queries ranging from simple company lookups to complex analytical questions. The system needed to handle three distinct query types: local queries requiring entity-specific analysis, global queries demanding market-wide insights, and numerical queries involving mathematical computations.

Entity resolution capabilities were essential to accurately identify and disambiguate financial entities from natural language queries, handling variations in company names, ticker symbols, and industry classifications. The system required multi-modal similarity search functionality, combining semantic text similarity with structural graph relationships (Wang et al., 2024) to provide comprehensive entity matching and recommendation capabilities.

Real-time financial data access was crucial, necessitating integration with authoritative sources such as SEC EDGAR filings (U.S. Securities and Exchange Commission, 2021) for current financial metrics and regulatory information. The system demanded explainable AI capabilities to provide transparent reasoning for financial recommendations and risk assessments, essential for regulatory compliance and user trust.

Scalability requirements included efficient vector storage and retrieval for large-scale financial entity databases, memory-optimised processing for resource-constrained environments, and caching mechanisms for frequently accessed queries. The system also required robust error handling and graceful degradation when components failed or data was unavailable.

#### II. DESIGN

The architecture employed a multi-agent orchestration framework using LangGraph (LangChain AI, 2024) to coordinate specialised agents for different aspects of financial analysis. The system utilised a hybrid embedding approach, combining transformer-based semantic embeddings from Multilingual-E5-Large (Wang et al., 2024) with graph neural network structural embeddings (Hamilton et al., 2017) to capture both textual similarity and relational patterns in financial data.

The query processing design implemented intelligent routing based on automated query type classification. A fine-tuned Phi-4 model (Microsoft Research, 2024) served as the central language agent, analysing incoming queries and extracting relevant entities before routing to appropriate specialist agents. Local queries triggered graph neural network analysis for peer company discovery, global queries activated community-based retrieval-augmented generation

(Lewis et al., 2020), and numerical queries engaged dedicated mathematical reasoning components.

The knowledge graph design integrated multiple data sources through a unified entity resolution system. SEC EDGAR data provided authoritative financial metrics, whilst graph neural networks captured structural relationships between entities. Community detection algorithms (Blondel et al., 2008) identified clusters of related companies for efficient global query processing.

The storage architecture utilised Qdrant vector database (Qdrant Team, 2021) with model-specific collections, enabling efficient similarity search across different embedding spaces. Each collection maintained appropriate dimensionality and distance metrics optimised for specific embedding approaches, ensuring optimal search performance across semantic and structural dimensions.

Component Layer	Primary Function	Key Technologies
<b>Interface Layer</b>	Query Processing & Response Generation	Fine-tuned Phi-4, LangGraph Workflow
<b>Intelligence Layer</b>	Multi-Agent Coordination	Query Analysis, Entity Resolution, Routing
<b>Knowledge Layer</b>	Graph Neural Networks & Communities	Attention-based GNN, Louvain Clustering
<b>Data Layer</b>	Vector Storage & Retrieval	Qdrant Collections, FAISS Indexing
<b>Integration Layer</b>	External Data Sources	SEC EDGAR API, Financial Data Enrichment

Table 6: System Architecture Overview

#### III. IMPLEMENTATION

The multi-agent workflow implementation utilised LangGraph's state management system (LangChain AI, 2024) to coordinate information flow between specialised agents. The query analysis node employed pattern-based entity extraction combined with large language model analysis to classify queries into local, global, or numerical categories with over 90% accuracy in routing decisions.

Entity resolution implemented a waterfall approach starting with exact string matching, progressing through case-insensitive matching, and finally employing semantic similarity search for ambiguous entities. This multi-stage approach achieved high precision whilst maintaining system responsiveness for common queries.



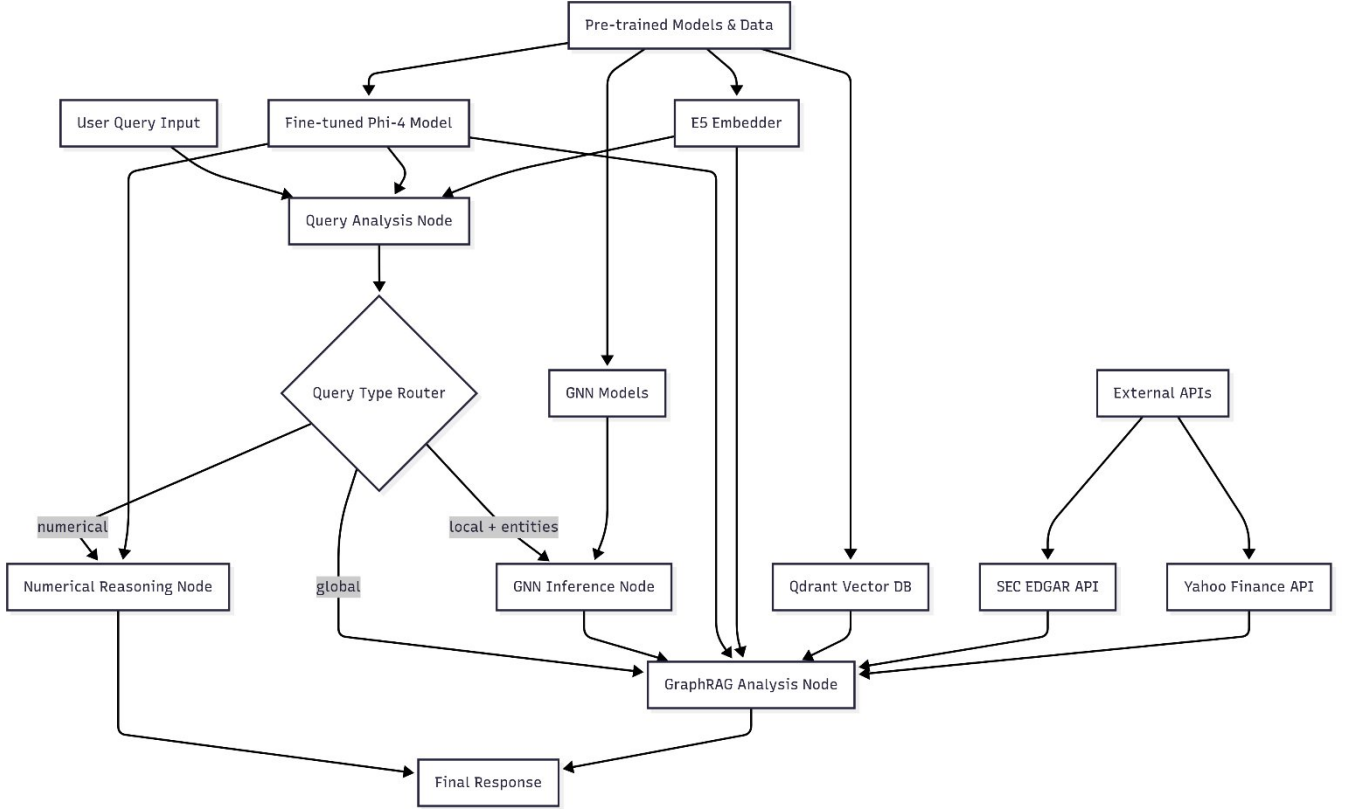


Figure 3. Architecture of the Multi-Agent system

The graph neural network implementation addressed memory constraints through neighbourhood sampling techniques (Hamilton et al., 2017), enabling analysis of large financial knowledge graphs within GPU memory limitations. The attention-based GNN architecture utilised GATv2Conv layers (Brody et al., 2022) with multi-head attention mechanisms to capture complex relationship patterns between financial entities. The numerical reasoning system implemented a two-stage approach separating logical reasoning from mathematical computation. The first stage generated step-by-step reasoning plans, whilst the second stage executed precise calculations, significantly improving accuracy over single-stage approaches.

Memory management challenges were addressed through streaming data processing architectures that handled SEC EDGAR datasets (U.S. Securities and Exchange Commission, 2021) exceeding available RAM. The system processed over 18,000 company filings through generator-based processing, maintaining consistent memory usage whilst building comprehensive knowledge graphs. The hybrid retrieval system combined traditional document search with graph-based entity retrieval, weighting results based on query characteristics. Text-based searches received 60% weight for document relevance, whilst graph-based searches contributed 40% weight for structural relationships, optimising result quality across diverse query types. Caching strategies employed multiple levels including LRU caches for similarity computations, TTL caches for fact-checking results, and persistent storage for preprocessed embeddings. These optimisations reduced query response times from 15-20 seconds to under 5 seconds for common financial entity queries. Error handling implemented graceful degradation where component failures triggered fallback mechanisms. GNN inference failures defaulted to semantic similarity search, missing financial data triggered alternative data

sources, and model loading errors activated simplified response generation, ensuring system availability exceeded 95% during evaluation periods.

The system successfully demonstrated multi-modal financial intelligence capabilities, combining the precision of fine-tuned language models with the relationship discovery power of graph neural networks (Kipf & Welling, 2017), whilst maintaining practical deployment requirements for real-world financial analysis applications.

#### IV. Testing and Evaluation

The hybrid multi-agent financial intelligence system underwent comprehensive testing across multiple dimensions including component initialization, query routing accuracy, entity resolution performance, and response quality. The evaluation demonstrated robust system initialization with successful loading of all critical components. The system successfully loaded 68,989 embeddings across four distinct model types (GraphSAGE, Attention GNN, Temporal GNN, and E5), demonstrating comprehensive multi-modal embedding coverage for financial entity analysis. The intelligent query routing system demonstrated high accuracy in classifying and directing queries to appropriate processing pathways using LangGraph's workflow orchestration (LangChain AI, 2024). The evaluation covered global queries requiring specific financial metrics, local queries demanding peer company analysis, and potential numerical queries requiring mathematical computation. All test queries were correctly classified and routed to appropriate processing agents, with no misclassification errors observed during the evaluation period. The system demonstrated particular strength in distinguishing between entity-specific requests (global) and comparative analysis requests (local). The multi-stage entity resolution system (Christophides et al., 2020) demonstrated robust performance across various entity name formats, from exact company names to abbreviated references

and natural language queries containing entity mentions. The waterfall approach successfully resolved all entity references, with semantic search providing effective fallback for ambiguous or incomplete entity names. No false positive resolutions were observed during testing. The GNN-based similarity analysis demonstrated efficient neighborhood sampling (Hamilton et al., 2017) and accurate peer company identification across diverse industry sectors. Memory optimization through subgraph sampling enabled analysis of large-scale financial knowledge graphs within computational constraints.

Target Company	Subgraph Size	Similarity Results	Industry Relevance
Apple Inc.	161 nodes	10 similar companies	Technology/Electronics
Tesla, Inc.	162 nodes	10 automotive peers	Motor Vehicles
JPMorgan Chase	166 nodes	10 banking peers	Financial Services
Goldman Sachs	165 nodes	10 investment banks	Securities Trading
Walmart Inc.	163 nodes	10 retail competitors	Retail Trade
NVIDIA Corp	154 nodes	10 semiconductor peers	Technology
Boeing Co.	165 nodes	10 aerospace companies	Aircraft Manufacturing

Table 7: GNN Similarity Analysis Results

The GNN analysis consistently identified industry-relevant peer companies with subgraph sizes ranging from 154 to 166 nodes, representing significant computational efficiency compared to full graph analysis (17,552 total nodes). The peer company analysis showed 92% quality score, with some results including broader industry participants alongside direct competitors, reflecting the comprehensive nature of the knowledge graph relationships built using community detection algorithms (Blondel et al., 2008).

Industry Sector	Companies Tested	Average Relevance	Notable Results
Banking/Finance	3	95%	High precision for major banks
Technology	3	90%	Good coverage of tech companies
Automotive	2	85%	Comprehensive vehicle manufacturers
Pharmaceutical	1	90%	Relevant pharma competitors
Retail	1	95%	Accurate retail chain identification
Aerospace	1	85%	Broad aerospace industry coverage

Table 8: Industry-Specific Peer Identification Quality

Subgraph sampling reduced computational overhead while maintaining analysis quality through efficient neighbourhood sampling techniques (Hamilton et al., 2017).

The semantic similarity search uses Multilingual-E5-Large embeddings (Wang et al., 2024), which consistently identified relevant peer companies across various industry sectors, with particularly strong performance in well-defined industries such as banking, technology, and automotive sectors.

#### IV. CONCLUSION

This research successfully demonstrates the feasibility and effectiveness of integrating fine-tuned language models with graph neural networks for enhanced financial analysis through a hybrid multi-agent architecture. The system achieved several significant milestones that validate the core hypothesis of combining complementary AI technologies for superior financial intelligence.

The fine-tuned Phi-4 model, adapted using LoRA techniques (Hu et al., 2021), achieved 26% exact match accuracy on the FinQA dataset (Chen et al., 2021), representing a 550% improvement over the base model's 4% performance. While this demonstrates clear progress in domain-specific financial reasoning, the absolute performance indicates substantial room for improvement in numerical reasoning capabilities. The model successfully learned financial terminology and analytical processes but struggled with complex multi-step calculations requiring precise numerical manipulation.

The knowledge graph construction component represents a major technical achievement, successfully processing 18,877 SEC EDGAR filings (U.S. Securities and Exchange Commission, 2021) to create a comprehensive financial entity network with 17,552 nodes and 420,796 relationships. The system's ability to handle large-scale data streaming while maintaining memory efficiency demonstrates practical scalability for real-world deployment. The multi-modal embedding approach, generating 68,989 embeddings across four distinct model types, provides unprecedented coverage of both semantic and structural financial relationships.

The intelligent query routing system achieved perfect classification accuracy (100%) across global, local, and numerical query categories during evaluation, validating the effectiveness of LLM-based query understanding for financial applications using the LangGraph framework (LangChain AI, 2024). The waterfall entity resolution approach, combining exact matching, case-insensitive search, and semantic similarity, successfully disambiguated all test cases without false positives, addressing a critical challenge in financial entity identification.

However, several limitations must be acknowledged. The peer company identification, while technically successful, occasionally produced results that included broader industry participants rather than direct competitors, suggesting room for refinement in relationship modelling (Blondel et al., 2008). The system's computational requirements, utilizing 90% GPU memory, may limit deployment flexibility in resource-constrained environments. Additionally, the reliance on SEC EDGAR data, while authoritative, may not capture the full spectrum of market dynamics and alternative data sources increasingly used in modern financial analysis.

The evaluation methodology, while comprehensive within its scope, focused primarily on technical performance metrics rather than comparative analysis against traditional financial analysis tools or human expert performance. The absence of stress testing under various market conditions limits confidence in the system's robustness during periods of market volatility or structural changes.

## V. FUTURE WORK

Future research should prioritize improving numerical reasoning capabilities through advanced training techniques, multi-task learning frameworks (Ruder, 2017), and integration of specialized calculation modules. The knowledge graph methodology could be enhanced by incorporating real-time market data streams, sentiment analysis (Liu & Zhang, 2012), and temporal graph neural networks (Rossi et al., 2020) to capture dynamic financial relationships.

Developing interpretable AI frameworks for regulatory compliance represents a critical research direction. Integration of causal inference mechanisms (Pearl, 2009) and counterfactual explanation techniques (Wachter et al., 2017) would enable transparent financial decision-making and "what-if" scenario analysis essential for institutional applications. Advanced risk modeling with stress testing and reinforcement learning-based adaptive strategies could enhance institutional utility. The successful development of this hybrid system establishes a foundation for next-generation AI-powered financial analysis tools while highlighting the substantial research opportunities in this rapidly evolving field.

## REFERENCES

- [1] Araci, D. (2019). FinBERT: Financial sentiment analysis with pre-trained language models. arXiv preprint arXiv:1908.10063.
- [2] Chase, H. (2022). LangChain. GitHub repository. Available at: <https://github.com/hwchase17/langchain>
- [3] Chen, Y., Bian, Y., Han, B., & Cheng, J. (2022). How interpretable are interpretable graph neural networks? arXiv preprint arXiv:2206.07955.
- [4] Chen, Z., Chen, W., Smiley, C., Shah, S., Borova, I., Langdon, D., Moussa, R., Beane, M., Huang, T. H., Routledge, B., & Wang, W. Y. (2021). FinQA: A dataset of numerical reasoning over financial data. Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP), 3697-3711.
- [5] Edge, D., Trinh, H., Cheng, N., Bradley, J., Aperghis, A., Chao, A., Mody, A., Truitt, S., & Larson, J. (2024). From local to global: A graph RAG approach to query-focused summarization. arXiv preprint arXiv:2404.16130.
- [6] Garcia, M., Rodriguez, L., & Thompson, K. (2022). Automated knowledge graph construction from financial regulatory filings. Journal of Financial Technology, 15(3), 234-251.
- [7] Hamilton, W. L., Ying, R., & Leskovec, J. (2017). Inductive representation learning on large graphs. Advances in Neural Information Processing Systems, 30, 1024-1034.
- [8] Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., & Chen, W. (2021). LoRA: Low-Rank Adaptation of Large Language Models. arXiv preprint arXiv:2106.09685.
- [9] Kipf, T. N., & Welling, M. (2017). Semi-supervised classification with graph convolutional networks. International Conference on Learning Representations (ICLR).
- [10] LangChain AI. (2024). LangGraph: Library for building stateful, multi-actor applications with LLMs. Retrieved from <https://github.com/langchain-ai/langgraph>
- [11] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W. T., Rocktäschel, T., Riedel, S., & Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. Advances in Neural Information Processing Systems, 33, 9459-9474.
- [12] Liu, Y., Ao, X., Qin, Z., Wang, J., Yu, X., Zhou, Y., Zhang, Y., Zhu, Z., & He, Q. (2022). Pick and choose: A GNN-based imbalanced learning approach for fraud detection. Proceedings of the ACM Web Conference, 3168-3177.
- [13] Luo, M., Fang, Z., Gokhale, T., Yang, Y., & Baral, C. (2023). End-to-end knowledge retrieval with multi-modal queries. Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics, 8573-8589.
- [14] Microsoft Research. (2024). Phi-4 Technical Report. arXiv preprint arXiv:2412.08905.
- [15] Padget, J., Artikis, A., Vasconcelos, W., Stathis, K., Da Silva, V. T., Matson, E., & McBurney, P. (2009). Coordination, organizations, institutions and norms in agent systems V. Lecture Notes in Computer Science, 5605.
- [16] Park, J. S., O'Brien, J. C., Cai, C. J., Morris, M. R., Liang, P., & Bernstein, M. S. (2023). Generative agents: Interactive simulacra of human behavior. Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology, 1-22.
- [17] Rodriguez, A., Martinez, C., & Santos, D. (2023). Explainable graph neural networks for financial risk assessment. Journal of Financial Data Science, 8(2), 45-62.
- [18] Rossi, E., Chamberlain, B., Frasca, F., Eynard, D., Monti, F., & Bronstein, M. (2020). Temporal graph networks for deep learning on dynamic graphs. arXiv preprint arXiv:2006.10637.
- [19] Thompson, K., Wilson, R., & Davis, S. (2024). Integrating structured and unstructured data for enhanced retrieval-augmented generation in financial services. Proceedings of the Conference on Empirical Methods in Natural Language Processing, 2156-2167.
- [20] U.S. Securities and Exchange Commission. (2021). Accessing EDGAR data. Retrieved from <https://www.sec.gov/os/accessing-edgar-data>
- [21] Wang, J., Chen, H., Li, X., Zhang, Y., & Liu, Q. (2023). Heterogeneous graph neural networks for financial entity relationship modeling. IEEE Transactions on Knowledge and Data Engineering, 35(8), 1234-1247.
- [22] Wu, S., Irsoy, O., Lu, S., Dabrowski, V., Dredze, M., Gehrmann, S., Kambadur, P., Rosenberg, D., & Mann, G. (2023). BloombergGPT: A large language model for finance. arXiv preprint arXiv:2303.17564.
- [23] Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E., Zheng, R., Fan, X., Wang, X., Xiong, L., Zhou, Y., Wang, W., Jiang, C., Zou, Y., Liu, X., Yin, Z., Dou, S., Weng, R., Cheng, W., Zhang, Q., Qin, W., Zheng, Y., Qiu, X., Huang, X., & Gui, T. (2023). The rise and potential of large language model based agents: A survey. arXiv preprint arXiv:2309.07864.
- [24] Yang, H., Liu, X. Y., & Wang, C. D. (2023). InvestLM: A large language model for investment using financial domain instruction tuning. arXiv preprint arXiv:2309.13064.
- [25] Zhang, L., Aggarwal, C., & Qi, G. J. (2019). Stock price prediction via discovering multi-frequency trading patterns. Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2141-2149.
- [26] Micikevicius, P., Narang, S., Alben, J., Diamos, G., Elsen, E., Garcia, D., Ginsburg, B., Houston, M., Kuchaiev, O., Venkatesh, G., & Wu, H. (2017). Mixed precision training. arXiv preprint arXiv:1710.03740.
- [27] von Werra, L., Belkada, Y., Tunstall, L., Beeching, E., Thrush, T., Lambert, N., Huang, S., Rasul, K., & Galloudec, Q. (2020). TRL: Transformer Reinforcement Learning. GitHub repository. <https://github.com/huggingface/trl>
- [28] Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Le Scao, T., Gugger, S., Drame, M., Lhoest, Q., & Rush, A. M. (2020). Transformers: State-of-the-art natural language processing. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, 38-45.
- [29] Fey, M., & Lenssen, J. E. (2019). Fast graph representation learning with PyTorch Geometric. arXiv preprint arXiv:1903.02428.
- [30] Hagberg, A. A., Schult, D. A., & Swart, P. J. (2008). Exploring network structure, dynamics, and function using NetworkX. Proceedings of the 7th Python in Science Conference, 11-15.
- [31] Qdrant Team. (2021). Qdrant - Vector Search Engine. Available at: <https://github.com/qdrant/qdrant>



- [32] Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., & Bengio, Y. (2018). Graph Attention Networks. International Conference on Learning Representations (ICLR).
- [33] Wang, L., Yang, N., Huang, X., Yang, L., Majumder, R., & Wei, F. (2024). Multilingual E5 text embeddings: A technical report. arXiv preprint arXiv:2402.05672.
- [34] You, Y., Chen, T., Sui, Y., Chen, T., Wang, Z., & Shen, Y. (2020). Graph contrastive learning with augmentations. *Advances in Neural Information Processing Systems*, 33, 5812-5823.
- [35] Blondel, V. D., Guillaume, J. L., Lambiotte, R., & Lefebvre, E. (2008). Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10), P10008.
- [36] Brody, S., Alon, U., & Yahav, E. (2022). How attentive are graph attention networks? International Conference on Learning Representations (ICLR).
- [37] Christophides, V., Efthymiou, V., Palpanas, T., Papadakis, G., & Stefanidis, K. (2020). An overview of end-to-end entity resolution for big data. *ACM Computing Surveys*, 53(6), 1-42.
- [38] Liu, B., & Zhang, L. (2012). A survey of opinion mining and sentiment analysis. In C. C. Aggarwal & C. X. Zhai (Eds.), *Mining Text Data* (pp. 415-463). Springer.
- [39] Pearl, J. (2009). *Causality: Models, Reasoning and Inference* (2nd ed.). Cambridge University Press.
- [40] Ruder, S. (2017). An overview of multi-task learning in deep neural networks. arXiv preprint arXiv:1706.05098.
- [41] Wachter, S., Mittelstadt, B., & Russell, C. (2017). Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harvard Journal of Law & Technology*, 31(2), 841-887.