

FinReflectKG: Agentic Construction and Evaluation of Financial Knowledge Graphs

Abhinav Arun
Domyn
New York, US
abhinav.arun@domyn.com

Fabrizio Dimino
Domyn
New York, US
fabrizio.dimino@domyn.com

Tejas Prakash Agarwal
Domyn
New York, US
tejas.p.agrawal@domyn.com

Bhaskarjit Sarmah
Domyn
Gurgaon, India
bhaskarjit.sarmah@domyn.com

Stefano Pasquali
Domyn
New York, US
stefano.pasquali@domyn.com

Abstract

The financial domain poses unique challenges for knowledge graph (KG) construction at scale due to the complexity and regulatory nature of financial documents. Despite the critical importance of structured financial knowledge, the field lacks large-scale, open-source datasets capturing rich semantic relationships from corporate disclosures. We introduce an **open-source, large-scale financial knowledge graph dataset built from the latest annual SEC 10-K filings of all S&P 100 companies** - a comprehensive resource designed to catalyze research in financial AI.

We propose a robust and generalizable knowledge graph (KG) construction framework that integrates intelligent document parsing, table-aware chunking, and schema-guided iterative extraction with a reflection-driven feedback loop. Our system incorporates a comprehensive evaluation pipeline, combining rule-based checks, statistical validation, and LLM-as-a-Judge assessments to holistically measure extraction quality. We support three extraction modes—single-pass, multi-pass, and reflection-agent-based allowing flexible trade-offs between efficiency, accuracy, and reliability based on user requirements. Empirical evaluations demonstrate that the reflection-agent-based mode consistently achieves the best balance, attaining a 64.8% compliance score against all rule-based policies (CheckRules) and outperforming baseline methods (single-pass & multi-pass) across key metrics such as precision, comprehensiveness, and relevance in LLM-guided evaluations.

The utility of our KG pipeline is demonstrated through its flexible extraction modes, coupled with a multi-faceted evaluation methodology. By releasing a high-quality, thoroughly evaluated dataset along with a comprehensive KG construction & evaluation framework, we aim to advance transparency, reproducibility, and innovation in financial KG research. **The dataset is publicly available at:** <https://anonymous.4open.science/r/KG-Financial-Datasets-SP-100-529B/README.md>

Keywords: Knowledge Graphs, Financial Data, SEC Filings, Natural Language Processing, Information Extraction, LLMs

ACM Reference Format:

Abhinav Arun, Fabrizio Dimino, Tejas Prakash Agarwal, Bhaskarjit Sarmah, and Stefano Pasquali. 2025. FinReflectKG: Agentic Construction and Evaluation of Financial Knowledge Graphs. In *6th ACM International Conference on AI in Finance (ICAIF '25)*, November 15–18, 2025, Singapore, Singapore. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3768292.3770363>

1 Introduction

Knowledge graphs (KGs) have become foundational for organizing and reasoning over complex, interconnected information in domains ranging from biomedicine to finance [2, 6, 12]. In the financial sector, the construction of high-quality KGs is particularly challenging due to the heterogeneous and highly interconnected nature of documents such as SEC 10-K filings. Although recent advances in large language models (LLMs) have enabled significant progress in information extraction [10, 27], there is a general lack of reliable KG benchmark datasets for the financial domain [9, 15], which hinders the widespread adoption of KGs for downstream financial applications. Additionally, most existing financial KGs are either limited in scope, relying primarily on news feeds [6, 12], or lack the rigorous evaluation required for building trust for wide scale deployment & adoption within the financial industry [8].

In this work, we address the gaps by introducing a large-scale financial KG dataset built exclusively using annual SEC 10-K filings of all the S&P 100 companies for the year 2024. Our approach is inspired by recent developments in prompt-driven and iterative extraction [3, 7, 23], schema canonicalization [27], and self-reflective LLM agents [12]. We have built a robust pipeline that combines intelligent document parsing, table-aware chunking, schema-guided iterative extraction, and reflection-driven feedback, culminating in a thoroughly evaluated resource for the financial AI community.

Our framework supports three extraction modes—single-pass, multi-pass, and reflection-agent-based and we further introduce a dynamic schema configuration process, leveraging both LLMs and domain experts to ensure business relevance and adaptability. The resulting dataset and methodology empower a range of downstream applications, including KG-powered entity search, multi-hop question answering, signal generation via news integration, and advanced graph powered predictive models & network analytics.

The main contributions of our paper are as follows:

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICAIF '25, Singapore, Singapore

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-2220-2/2025/11
<https://doi.org/10.1145/3768292.3770363>

- **Large-Scale Open Source Financial KG Dataset:** We release an open-source, comprehensive financial knowledge graph constructed from SEC 10-K filings of all S&P 100 companies, providing a massive, high-quality resource for financial AI research and applications.
- **Novel Reflection Driven Extraction Framework:** We introduce a three-mode extraction pipeline (single-pass, multi-pass, reflection-agent-based) with reflection driven feedback systematically improving extraction quality through iterative refinement, achieving 64.8% compliance across all Check-Rules and best voted mode in the LLM-as-a-Judge evaluation.
- **Generalizable Evaluation Methodology:** We develop a holistic evaluation framework encompassing rule-based compliance checks, coverage analysis, semantic diversity metrics, and LLM-as-a-Judge comparative assessment, establishing new benchmarks for financial KG evaluation.

By releasing this high-quality, rigorously evaluated dataset and pipeline, we aim to advance transparency, reproducibility, and innovation in financial knowledge graph research and its applications.

2 Related Work

In this section, we provide details on existing research that has been happening in the realm of Knowledge Graph Construction with a focus on the financial domain.

Knowledge Graph Construction: Knowledge graph construction (KGC) has traditionally followed a pipeline architecture composed of discrete subtasks such as entity recognition [14, 28] and relation classification [25, 26], often relying on supervised learning with domain-specific annotations. However, the recent success of LLMs in information extraction tasks [10] highlights a paradigm shift toward prompt-based strategies for robust knowledge graph construction.

Recent work has leveraged large language models (LLMs) for zero and few-shot extraction using iterative prompting strategies. Carta [3] propose a zero-shot pipeline using GPT-3.5, where a sequence of prompts progressively identifies entities, types them, extracts relations, and resolves co-references without supervision or external knowledge bases. Jiang et al. [7] introduce the Retrieval-And-Structuring (RAS) framework, which alternates between query planning, retrieval, and triple extraction to incrementally construct task-specific mini-KGs, outperforming standard RAG methods. These approaches complement interactive designs like ChatIE [23], which frames triple extraction as a multi-turn QA task. Collectively, they underscore the emergence of modular, prompt-driven workflows for scalable KG construction with LLMs.

Normalization and Schema Canonicalization : While open information extraction (OIE) methods enable large-scale triples extraction, they often produce unstandardized and semantically redundant outputs. Without canonicalization, multiple surface-level variations of the same relation (e.g., *supplies*, *is supplier of*) can co-exist in the knowledge graph, introducing ambiguity and reducing its utility for downstream tasks. Traditional approaches depend on whether a target schema is available: alignment-based methods use lexical resources like WordNet, while schema-free methods such as CESI [21] cluster relations using embeddings and external signals.

However, clustering often over-generalizes, merging distinct semantics. The recent Extract-Define-Canonicalize (EDC) framework [27] offers a more robust, LLM-native alternative. It extracts triples, defines schema candidates in context, and canonicalizes them using LLM-generated definitions and reasoning-avoiding brittle heuristics or external ontologies.

Financial Knowledge Graphs: Early work on financial knowledge graphs (KGs) focused on extracting structured event representations from unstructured news. Benetka et al. [2] jointly extract all attributes of economic transactions as quintuples, aggregating information across multiple mentions to build unified event-centric graphs. Subsequent efforts, such as Elhammadi et al. [6], combined semantic role labeling, dependency parsing, and domain-specific dictionaries to construct high-precision financial KGs from news. More recently, Li and Passino [12] introduced FinDKG, a dynamic, time-varying financial KG generated from news using a fine-tuned LLM pipeline, systematically extracting entities and relations as event quadruples. Collectively, these works highlight the evolution from rule-based and supervised pipelines to LLM-driven, schema-flexible approaches for building robust financial knowledge graphs.

Building upon these foundations, we present an open-source, large-scale financial knowledge graph constructed exclusively from SEC 10-K filings of all S&P 100 companies. Our work addresses key limitations in existing approaches: while previous financial KGs focused on news-based event extraction, we target the most authoritative financial documents (SEC filings) with comprehensive schema-guided extraction, establishing new benchmarks for financial KG construction and evaluation.

3 Problem Formulation and Schema Design

Building on recent advances in large language models and prompt engineering, we design a pipeline to extract high-quality financial knowledge graph triples of the form (Head Entity, Head Type, Relationship, Tail Entity, Tail Type) from the SEC 10-K filings.

Definition 3.1. A *triple* in our context is a 5-tuple of the form (Head Entity, Head Type, Relationship, Tail Entity, Tail Type), where the Head and Tail Entities are linked by a semantic Relationship, and each entity is annotated with its type. Throughout this paper, we refer to these 5-tuples as triples.

Our pipeline begins with a business-driven schema configuration approach (closed information extraction), where entity types and relationships are primarily defined by business subject matter experts (SMEs), specific to input financial data feeds and downstream applications [8, 12]. For instance, the same SEC 10-K filing can be utilized differently by portfolio managers and risk managers. We employ an AI-assisted reconciliation approach where we prompt the underlying LLM to propose a schema given sample documents, but the schemas are ultimately approved and defined by SMEs. This approach ensures domain relevance and business alignment.

We focus on extracting KG triples in a closed information extraction setting, which leads to less noisy knowledge graphs and enables the construction of reliable KGs that can be directly integrated into downstream applications. For SEC filings, we have created a comprehensive schema through collaboration between LLMs and financial SMEs, a subset of which is presented in Tables 1

and 2. This schema captures the complex relationships and entities specific to financial reporting and regulatory compliance.

Table 1: Subset of Pre-Configured Entity Types and their Definitions for Financial KG Construction

Entity Type	Definition
ORG	Filing Company (Issuer: The public company that is the subject of the 10-K filing)
PERSON	Key individuals (e.g., CEO, CFO, Board members)
COMP	External companies referenced in the filing, including competitors, suppliers, customers, or partners
PRODUCT	Products or services offered by the company or competitors (e.g., iPhone, AWS)
SEGMENT	Internal divisions or business segments of the filer ORG (e.g., Cloud segment, North America retail)
FIN_METRIC	Financial metrics or values (e.g., Net Income, EBITDA, CapEx, Revenue)
RISK_FACTOR	Documented risks (e.g., market risk, supply chain risk, regulatory risk)
EVENT	Material events such as pandemics, natural disasters, M&A events, regulatory changes
REGULATORY_REQUIREMENT	Specific regulations or legal frameworks (e.g., Basel III, GDPR, SEC requirements)
ESG_TOPIC	Environmental, Social, and Governance themes (e.g., Carbon Emissions, DEI, Climate Risk)

Table 2: Subset of Pre-Configured Relationship Types and their Definitions for Financial KG Construction

Rel. Type	Definition
Has_Stake_In	Indicates full or partial ownership or equity interest
Operates_In	Indicates operational geography or market presence
Produces	Manufactures or develops a product or service
Impacts	Specifies the broad influence or effect an entity or event has on financial performance, market trends, or other key outcomes
Involved_In	Specifies direct involvement in an event such as a merger, acquisition, or litigation
Impacted_By	Indicates that the entity was materially affected by a major event
Discloses	Reveals or reports information, metrics, or developments
Complies_With	Meets regulatory or policy requirements
Supplies	Indicates vendor or supplier relationship
Partners_With	Indicates formal or strategic collaboration

4 Methodology

Once we have configured the schema & defined the ontology, our KG construction pipeline comprises of four interlocking components: (1) Intelligent Document Parsing Layer, (2) Table-Aware

Semantic Chunking Layer, (3) Iterative Prompt & Agent Driven Triples Extraction Layer, and (4) Robust Evaluation Layer. By integrating best practices from AI-driven information extraction and financial KG construction, we achieve both high precision and domain relevance. The overall pipeline is showcased in the Figure 1

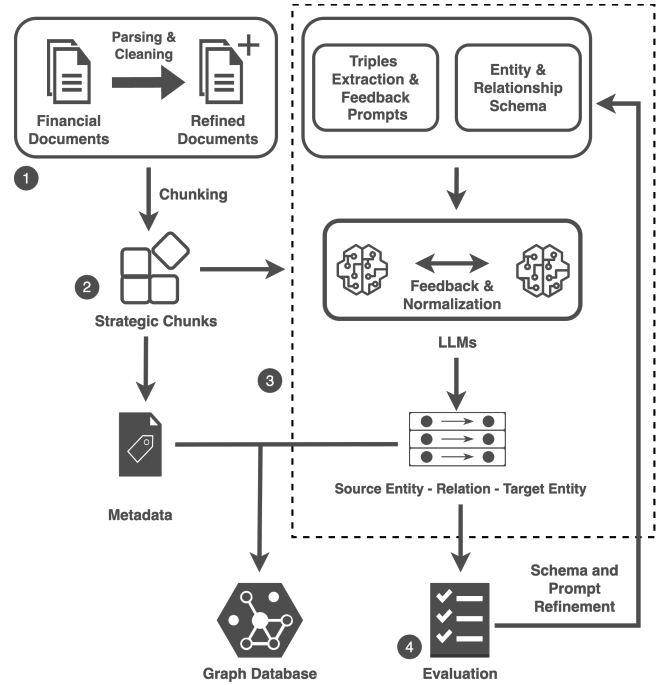


Figure 1: Overview of the Graph RAG microservice and its integration with the KG construction pipeline.

4.1 Intelligent Document Parsing Layer

Our pipeline leverages an advanced document parsing layer built on `docling` [1], enabling robust extraction and retention of the diverse formats present in SEC 10-K filings. This layer preserves the multiformal architecture of the source documents, including narrative text, tables, and images. Tables are especially critical in financial filings, as they often encapsulate key quantitative disclosures (e.g., revenue by segment, risk breakdowns, or financial metrics) essential for downstream knowledge graph construction.

The parser operates in two modes: a **multimodal mode**, which combines OCR and image annotation to capture both textual and visual information, and a **text-only mode** for efficient processing when only narrative content is required. For the scope of this paper, we used only the text mode. In both modes, tables structures are retained as markdown context, maintaining row-column associations and semantic context. Narrative sections are tagged with section headers (e.g., “Risk Factors,” “Management’s Discussion”) to facilitate schema alignment in subsequent pipeline stages.

This intelligent parsing approach ensures that no critical information is lost during preprocessing, and provides a high-fidelity, semantically annotated representation of the original document for downstream chunking and extraction.

4.2 Table-Aware Semantic Chunking Layer

To accommodate the token limits of modern large language models, our pipeline incorporates a custom chunking algorithm that strategically segments documents while preserving semantic and structural integrity. Unlike naive sliding window or fixed-length approaches, our chunker is **table-aware**: any detected table in the markdown is retained as a single atomic chunk, ensuring that row and column relationships—and thus the full context of financial data—are never fragmented. This is particularly important for SEC 10-K filings, where tables often encode critical quantitative disclosures.

We employ a **section-aware text segmentation**, splitting text at logical boundaries such as paragraphs or subsection headings to maintain topical coherence. Each chunk is constrained to a maximum size of 2048 tokens (`CHUNK_SIZE = 2048`), balancing context preservation with LLM input requirements.

This approach ensures that both tabular and textual information are optimally prepared for downstream extraction, maximizing the fidelity and relevance of the knowledge graph construction process.

4.3 Iterative Prompt & Agent Driven Triples Extraction Layer

Leveraging the predefined schema as highlighted in Tables 1 and 2, we adopt an iterative empirical approach to identify the paradigm that gives us the most reliable and grounded KG triples. We have used Qwen2.5-72B-Instruct as the LLM for KG construction. The prompts explicitly enumerate these categories to constrain LLM outputs.

4.3.1 Single-Pass Workflow. In the single-pass mode, we employ a single, comprehensive prompt that instructs the language model to extract all valid knowledge graph triples from each document chunk in one step. The prompt enforces the use of only pre-defined entity and relation types, requires normalization of entity names (e.g., mapping all company references to the ticker), and outputs the results in a strict JSON format for downstream processing. While efficient, this approach may still yield occasional inconsistencies in entity normalization or relation assignment due to the inherent limitations of single-turn LLM prompting. Mathematically, for each chunk c and schema S :

$$T_c^{(1)} = \text{Extract}(c, S, \phi_{\text{sp}}) \quad (1)$$

where $T_c^{(1)}$ is the set of extracted and normalized triples for chunk c , and ϕ_{sp} denotes the single pass (sp) mode (prompts).

4.3.2 Multi-Pass Workflow. To improve extraction quality and consistency, we adopt a multi-pass prompting strategy. In this approach, the language model first extracts candidate triples from each chunk using the pre-defined schema. In a second pass, the model re-ingests its own output alongside the original chunk and applies a dedicated normalization prompt to:

- Enforce canonical naming (e.g., ticker substitution for company references).
- Filter to schema-compliant entity and relation types.
- Merge duplicate or redundant entities and relationships.
- Validate directionality and ordering for all relations.

- Remove or correct invalid or ambiguous triples.

This two-step process leverages the LLM’s reasoning capabilities for both extraction and refinement, resulting in higher precision and more consistent knowledge graph triples. Mathematically, for each chunk c and schema S :

$$T_c^{(1)} = \text{Extract}(c, S, \phi_{\text{mp}}) \quad (2)$$

$$T_c^{(2)} = \text{Normalize}(c, T_c^{(1)}, S, \phi_{\text{mp}}) \quad (3)$$

where $T_c^{(1)}$ is the initial set of extracted triples, $T_c^{(2)}$ is the refined, schema-compliant set after normalization, and ϕ_{mp} denotes the multi pass (mp) mode (prompts).

4.3.3 Reflection-Driven Agentic Workflow and Meta-Analysis. We deploy a dedicated reflection agent that iteratively refines the initial set of triples by simulating a multi turn interaction between the feedback (critic) & correction LLM. The critic LLM specifically:

- Verifies the entity labels and relation assignments against the domain schema.
- Assesses business relevance and flags low-value or contradictory triples.

Feedback is returned in a structured JSON schema enabling automated ingestion as shown in box 4.1. All critique instances are logged for meta-analysis, revealing recurrent error patterns and informing prompt redesign. Our reflection approach is inspired by recent advances in self-reflective and memory-augmented language agents [11, 18].

Box 4.1: Sample Response from the Feedback LLM

```
{
  "triple_number": "Triple N",
  "triple": ["We", "ORG", "Impacted_By", "supply chain disruptions", "RISK_TYPE"],
  "issue": "Indirect reference to an entity in the triple. RISK_TYPE is not a valid preconfigured category",
  "suggestion": "replace We with NVDA as this information comes from Nvidia's 10-K file; substitute RISK_TYPE with RISK_FACTOR from the configured entity types"
}
```

Let $T_c^{(1)}$ be the initial set of triples for chunk c (from the extraction LLM), S is the preconfigured schema and ϕ_{re} denotes the reflection (re) extraction mode. For each reflection step $t = 1, 2, \dots, n$:

- **Extraction LLM** generates an initial set of triples $T_c^{(1)}$ which are further validated and refined in a cyclic loop of feedback & correction LLMs.
- **Feedback LLM** analyzes $T_c^{(t-1)}$ and produces a set of issues $F_c^{(t)}$ and suggestions.
- **Correction LLM** updates problematic triples (or drops them) to produce $T_c^{(t)}$.

$$T_c^{(1)} = \text{Extract}(c, S, \phi_{\text{re}}) \quad (4)$$

$$F_c^{(t)} = \text{Feedback}(c, T_c^{(t-1)}, S, \phi_{\text{re}}) \quad (5)$$

$$T_c^{(t)} = \text{Correct}(c, T_c^{(t-1)}, F_c^{(t)}, S, \phi_{\text{re}}) \quad (6)$$

where ϕ_{re} denotes the reflection (re) mode.

Stopping criteria:

Stop at step t^* if $F_c^{(t^*)} = \emptyset$ (no issues found) or $t^* = n_{\max}$ (max steps reached which can be empirically determined for different LLMs). This is intelligently determined using the agentic prowess of the underlying LLMs.

Final output:

$$T_c^{(*)} = T_c^{(t^*)} \quad (7)$$

For all chunks in document D :

$$T_D^{(*)} = \bigcup_{i=1}^N T_{c_i}^{(*)} \quad (8)$$

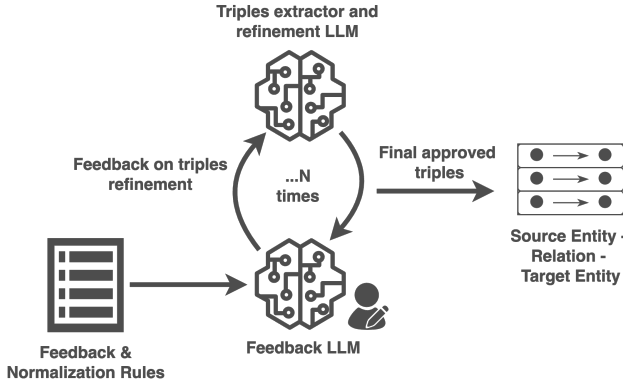


Figure 2: Overview of Reflection Agent where the triples are refined using an iterative feedback loop

This iterative prompt engineering & agentic approach advances the current state of domain-specific KG extraction and offers a reusable framework for other regulated domains.

5 Evaluation

Evaluating knowledge graph construction requires overcoming limitations of traditional metrics that rely on simplistic or single-dimensional measures, especially when comprehensive ground-truth annotations are unavailable [5]. To address this, we propose a holistic evaluation framework integrating complementary methodologies:

5.1 CheckRules

Knowledge graph construction systems often suffer from inconsistent entity normalization and schema violations, leading to redundant representations and semantic ambiguity. A critical issue arises from abstract entity references that lack semantic clarity: pronouns like "the company", "we", "our", or "it" create ambiguous nodes that cannot be properly linked or reasoned upon. For instance, in a financial document about Apple Inc., the same corporate entity may be represented as "Apple Inc.", "the company", "we", or the standardized ticker "AAPL", creating artificial multiplicity and semantic confusion that degrades graph quality and downstream reasoning capabilities.

To address these systematic issues, **CheckRules** evaluates each extracted triple against a set of rules:

- **Subject Reference:** Identifies and flags abstract entity references (e.g., "the company", "we", "our", "it") that lack semantic specificity (e.g., "AAPL" for Apple Inc.).
- **Entity Length Constraint:** Limits entity names to maximum 5 words to prevent verbose representations.
- **Entity Schema Compliance:** Ensures extracted entities types adhere to a predefined schema.
- **Relationship Schema Compliance:** Ensures extracted relationship adhere to a predefined schema.

Each extracted triple is individually evaluated against these rules. For a triple t with R rules, the CheckRules score is:

$$CR(t) = \frac{1}{R} \sum_{i=1}^R \phi_i(t) \quad (9)$$

where $\phi_i(t) \in \{0, 1\}$ indicates compliance.

Table 3: CheckRules scores for different extraction modes

Rule(↑)	Single Pass (%)	Multi Pass (%)	Reflection (%)
Subject Reference	99.9	100.0	100.0
Entity Length Constraint	68.2	79.5	78.0
Entity Schema Compliance	95.9	96.6	98.1
Rel. Schema Compliance	64.6	62.0	84.2

Table 4: Proportion of valid triples with increasingly strict criteria

Metric (↑)	Single Pass (%)	Multi Pass (%)	Reflection (%)
At least 1 rule	100.0	100.0	100.0
At least 2 rules	99.1	99.4	99.8
At least 3 rules	87.3	90.8	95.6
At least 4 rules	42.3	47.3	64.8

While Table 3 identifies specific compliance bottlenecks, Table 4 quantifies the proportion of triples considered valid under increasingly strict criteria. The results show that company references and entity-type compliance are consistently handled across all modes. However, relationship-schema compliance lags behind, indicating the need for further refinements of the predefined schema[22]. Notably, the Reflection paradigm substantially improves compliance, particularly concerning entity naming length and relationship schema rules.

The second table indeed confirms this by highlighting Reflection's notable performance advantage at higher compliance thresholds.

5.2 Local Extraction Efficiency

To assess extraction effectiveness and comprehensiveness, we compute Coverage Ratios, quantifying diversity and completeness across entities, entity types, and relationships. Specifically, we measure the proportion of unique entities and entity types relative to total extracted elements (Entity Coverage Ratio, ECR; Type Coverage Ratio,

Table 5: Local Extraction Efficiency scores across different extraction modes

Metric (\uparrow)	Single Pass (%)	Multi Pass (%)	Reflection (%)
Triples (per chunk)	13.3	12.4	15.8
ECR	0.30	0.31	0.53
TCR	0.09	0.10	0.27
TCR-N	0.13	0.13	0.18
RCR	0.21	0.22	0.38
RCR-N	0.13	0.13	0.14

TCR), and evaluate schema utilization through normalized coverage ratios (TCR-N, RCR-N). Similar computations for relationships yield Relationship Coverage Ratios (RCR, RCR-N).

Table 5 highlights a clear performance hierarchy, with Reflection consistently outperforming Single Pass and Multi Pass across all coverage metrics. Reflection generates more triples per chunk and achieves substantially higher entity, type, and relationship coverage, indicating a richer and more diverse extraction of semantic content. Normalized ratios (TCR-N, RCR-N), although improved by Reflection, it still suggests the underutilization of schema and highlights possible avenue for schema-design improvements and better alignment with extracted data.

5.3 Global Semantic Diversity

To quantify semantic diversity within the dataset, we analyze the distribution of extracted entities, types, and relationships using information-theoretic measures. Given frequency distributions $\mathbf{p} = \{p_1, p_2, \dots, p_n\}$ where p_i represents the normalized frequency of element i in the extracted set, we compute Shannon Entropy [17], which provides insights into overall diversity by measuring the balance or skewness of distributions across extracted elements:

$$H(X) = - \sum_{i=1}^n p_i \log_2 p_i \quad (10)$$

In addition, Schema-Normalized Entropy contextualizes diversity within the predefined schema:

$$H_{\text{norm}}(X) = \frac{H(X)}{\log_2 |S|} \quad (11)$$

where $|S|$ is the total number of elements defined in the schema.

Furthermore, Rényi Entropy [16], with parameter $\alpha = 2$, emphasizes concentration on frequent or rare elements:

$$H_2(X) = - \log_2 \left(\sum_{i=1}^n p_i^2 \right) \quad (12)$$

For each extraction method, we compute these measures on three distributions: (1) entity frequencies \mathbf{p}_E , (2) entity type frequencies \mathbf{p}_T , and (3) relationship frequencies \mathbf{p}_R , aggregated across all document chunks.

Table 6 reveals a clear trade-off between extraction completeness and variety. The Reflection method achieves the highest coverage but exhibits the lowest entropy across all dimensions. This result demonstrates that Reflection intentionally reduces diversity to yield a more compact, connected, and navigable graph consistent with

Table 6: Entropy scores across different extraction modes

Metric	Single Pass	Multi Pass	Reflection
<i>Shannon Entropy</i>			
Entity	7.5383	7.2845	7.1779
Entity Type	3.0290	2.9835	2.8665
Relationship	5.5438	5.6116	4.3164
<i>Normalized Schema Entropy</i>			
Entity Type	0.6607	0.6507	0.6252
Relationship	1.1412	1.1552	0.8885
<i>Rényi Entropy ($\alpha = 2$)</i>			
Entity	2.8619	2.5883	2.4574
Entity Type	2.1312	2.0851	1.9877
Relationship	3.6355	3.5691	2.6814
<i>Normalized Schema Rényi Entropy ($\alpha = 2$)</i>			
Entity Type	0.4648	0.4548	0.4335
Relationship	0.7484	0.7347	0.5520

predefined extraction rules. Given the significant improvement of Reflection on complementary metrics, this reduction in entropy is well within the acceptable range. Future work will monitor entropy drift and adapt the Reflection rules when diversity falls below the predefined threshold. All extraction methods exhibit moderate schema-normalized entropy, suggesting that the current ontology effectively captures core semantic categories in the corpus, while retaining room for schema enhancement to accommodate greater semantic detail.

5.4 Comparative Evaluation: LLM-as-a-Judge

The absence of benchmark ground truth triples for knowledge graph evaluation necessitates an alternative evaluation approach. We leverage LLMs as comparative judges to assess:

- **Precision:** Assesses the clarity, specificity, and uniqueness of the extracted triples.
- **Faithfulness:** Measures the factual accuracy and grounding of the triples within the source text.
- **Comprehensiveness:** Evaluates how completely the generated triples capture the core informational content of the source text.
- **Relevance:** Determines the contextual alignment of triples with the main topics and themes of the source text.

This methodology enables ground truth-agnostic metrics that provide relative comparative evaluations rather than absolute measurements, circumventing the need for predefined reference datasets. According to recent findings [4, 19], Chain-of-Thought (CoT) reasoning [24] can underperform in linear, fast, and intuitive tasks, potentially introducing unnecessary complexity and overthinking. Consequently, we adopt a prompting strategy similar to that used by Lopez et al. [13], employing direct instructions without intermediate reasoning steps. In this approach, the model initially commits to a judgment and subsequently provides an explanation.

This aligns with cognitive insights indicating that post-hoc rationalization can offer improved human interpretability compared to reasoning that precedes commitment [20]. In our evaluation, we utilized the Qwen3-32B model without reasoning prompts and set the temperature parameter to 0.1.

Table 7: LLM as a Judge Score for various extraction modes

Metric (\uparrow)	Single Pass (%)	Multi Pass (%)	Reflection (%)
Precision	22.3	38.6	39.1
Faithfulness	40.1	24.4	35.5
Comprehensiveness	36.3	15.6	48.1
Relevance	34.6	28.1	37.3

To ensure the reliability of our comparative evaluations, we implement a robustness assessment protocol. For each three-way comparison on chunk c_{ij} and metric μ , we conduct $n_{\text{votes}} = 3$ independent evaluations:

$$\mathcal{V}_{\mu}(c_{ij}) = \{J_{\mu}^{(k)}(c_{ij}, \mathcal{T}_{ij}^{(m_1)}, \mathcal{T}_{ij}^{(m_2)}, \mathcal{T}_{ij}^{(m_3)})\} k = 1^3 \quad (13)$$

where k represents the k -th independent vote. In cases where the votes do not reach consensus, we request a fourth decisive vote.

Table 8: Consistency across the three LLM as a Judge runs

Metric (\uparrow)	Agreement (%)
Precision	82.1
Faithfulness	81.3
Comprehensiveness	86.7
Relevance	83.7

Table 7 reveals nuanced performance patterns across extraction modes. The reflection mode outperforms other modes in precision, comprehensiveness, and relevance, whereas the single-pass mode excels in faithfulness. This dichotomy underscores a fundamental trade-off between comprehensiveness of extraction and factual grounding. The reflection mode’s relatively lower faithfulness score, despite greater comprehensiveness, suggests increased triple generation potentially exceeds source-constrained accuracy boundaries. Moreover, Table 8 suggests an inverse correlation, where lower agreement is indicative of higher contextual uncertainty.

6 Discussion

Our comprehensive evaluation framework reveals that the reflection mode delivers the optimal reliability-coverage balance among the three extraction modes. The reflection paradigm demonstrates superior performance across multiple complementary metrics: it achieves the highest CheckRules compliance (64.8% for all four rules), generates the most triples per chunk (15.8), and wins a clear plurality in LLM-as-a-Judge evaluations across precision, comprehensiveness, and relevance dimensions.

The iterative critic-corrector loop systematically addresses schema violations while expanding triple coverage, yielding a denser yet

cleaner knowledge graph. This improvement stems from the reflection mechanism’s ability to identify and correct extraction errors while discovering previously missed but valid triples. The reflection mode’s superior entity coverage ratio (ECR: 0.53 vs. 0.30-0.31) and relationship coverage ratio (RCR: 0.38 vs. 0.21-0.22) demonstrate its effectiveness in capturing diverse semantic content. The diversity analysis further illuminates the reflection mode’s approach: while achieving the highest coverage, it exhibits lower entropy across all dimensions, indicating a deliberate reduction in uncertainty to yield a more compact, connected, and navigable graph.

However, these gains come with some trade-offs. The reflection agent requires additional inference rounds, potentially limiting its suitability for real-time applications requiring fast turnaround (e.g., intraday news feeds). In such scenarios, the single pass strategy may offer a viable alternative, recovering most normalization benefits with reduced computational overhead.

Our evaluation results reveal some limitations. First, cross document co-reference resolution is only partially addressed as the reflection loop operates on isolated filings. Second, our evaluation methodology depends on LLM-voting surrogate ground truth, risking propagation of biases inherent in the underlying judge models.

7 Conclusion and Future Work

In this work, we democratize the access to a reliable financial knowledge graph (KG) triples by releasing an open-source dataset constructed from SEC 10-K filings of all the S&P 100 companies. Given the universal and utilitarian nature of SEC 10-K filings, we believe that this open-source dataset will empower developers and researchers to build robust, finance-specific applications, advancing transparency and innovation in the financial domain.

Our pipeline demonstrates that reliable KG triples can be generated using relatively compact language models (<100B parameters) without costly fine-tuning, provided a robust evaluation and normalization framework is in place—an important consideration for regulated industries where trust and auditability are paramount. This motivated our development of a holistic novel framework for evaluating our KG construction pipeline and its results.

We are currently pursuing the below enhancements to further broaden the impact and flexibility of our approach.

Schema-Free KG Construction and Self-Improvement: Inspired by the Extract-Define-Canonicalize (EDC) paradigm [27], we are developing a schema-free pipeline capable of creating schemas from scratch and iteratively refining existing ones. This is particularly valuable for private financial data sources where schema requirements are unknown. Our evaluation results demonstrate that pre-configured schemas may not capture all nuances within financial data, necessitating a schema discovery and improvement pipeline. It would be interesting to compare the results of closed vs. open schema frameworks in terms of coverage and performance degradation. Careful processing & retrieval of expanded schemas is critical for efficient instruction following during inference.

Temporal Knowledge Graphs for Thematic Investing: Building on recent advances in financial temporal KGs [12], we plan to enrich our dataset with time-aware features to enable structured event timelines and causal reasoning. Our approach will incorporate temporal annotations capturing the evolution of financial

relationships—market events, regulatory changes, and corporate developments—across various time horizons. Temporal KGs can enhance explainability for signal generation for market movements and offer actionable insights for traders and risk managers.

Enhanced LLM-as-a-Judge Evaluation Methodology: Several factors warrant consideration when interpreting the results of LLM-as-a-judge approach. First, our use of Qwen3-32B without reasoning capabilities may not capture the full potential of more sophisticated reasoning models. Second, the extrinsic validation by using diverse LLM families is necessary to assess cross-model reliability and can help normalize the biases inherent in using a specific model family.

Table-Aware Serialization for Improved Extraction:

We observed that smaller language models exhibit conservative extraction behavior from markdown-formatted tables missing semantic relationships despite structured data availability. A dedicated table-aware serialization module could significantly enhance knowledge graph coverage and completeness, particularly for quantitative financial metrics and their inter relationships.

We have used the underlying KG for essential downstream applications, such as question answering, proactive recommendations by integrating real-time news feeds with the structured KG. The promising results highlight the broad applicability of our KG for real-world financial use cases and can be explored at scale.

Looking ahead, we are making a concerted effort to not only augment our pipeline with above enhancements but also expand our KG triple dataset to include all the S&P 500 companies with data from the last 10 years of annual SEC 10-K filings, significantly broadening the scope and temporal coverage of our financial KG. We hope to catalyze further research and practical adoption of knowledge graphs in the financial sector and beyond.

References

- [1] Christoph Auer, Maksym Lysak, Ahmed Nassar, Michele Dolfi, Nikolaos Livathinos, Panos Vagenas, Cesar Ramis Berrospi, Matteo Omenetti, Fabian Lindlbauer, Kasper Dinkla, Lokesh Mishra, Yusik Kim, Shubham Gupta, Rafael Teixeira de Lima, Valery Weber, Lucas Morin, Ingmar Meijer, Viktor Kuropiatnyk, and Peter W. J. Staar. 2024. *Docling: An Open-Source, Extensible and Efficient PDF Document Converter*. Technical Report. IBM Research. <https://arxiv.org/pdf/2408.09869> arXiv preprint arXiv:2408.09869.
- [2] Robert Benetka, Pedro Szekely, and Craig A Knoblock. 2017. Financial news event extraction using semi-supervised learning. *arXiv preprint arXiv:1708.05663* (2017). <https://arxiv.org/abs/1708.05663>
- [3] Luca Carta. 2023. Iterative Zero-Shot Prompting for Knowledge Graph Construction with Large Language Models. *arXiv preprint arXiv:2310.08455* (2023).
- [4] Xingyu Chen, Jiahao Xu, Tian Liang, Zhiwei He, Jianhui Pang, Dian Yu, Linfeng Song, Qiuzhi Liu, Mengfei Zhou, Zhuosheng Zhang, et al. 2024. Do not think that much for $2+3=?$ on the overthinking of o1-like llms. *arXiv preprint arXiv:2412.21187* (2024).
- [5] S. Choi and Y. Jung. 2025. Knowledge Graph Construction: Extraction, Learning, and Evaluation. *Applied Sciences* 15, 7 (2025), 3727. <https://doi.org/10.3390/app15073727>
- [6] Ahmed Elhammedi, Saptarshi Ghosh, et al. 2020. A pipeline for financial knowledge graph construction from news. *arXiv preprint arXiv:2010.13604* (2020). <https://aclanthology.org/2020.coling-main.550>
- [7] Yongqi Jiang, Deming Ye, Minghan Bi, and Danqi Chen. 2025. RAS: Retrieval-And-Structuring for Knowledge-Intensive Generation. *arXiv preprint arXiv:2504.02483* (2025).
- [8] Natthawut Kertkeidkachorn, Rungsiman Nararatwong, Ziwei Xu, and Ryutaro Ichise. 2023. FinKG: A Core Financial Knowledge Graph for Financial Analysis. In *2023 IEEE 17th International Conference on Semantic Computing (ICSC)*. IEEE, 90–93. <https://doi.org/10.1109/ICSC56153.2023.00020>
- [9] Rik Koncel-Kedziorski, Michael Krumdieck, Viet Lai, Varshini Reddy, Charles Lovering, and Chris Tanner. 2023. BizBench: A Quantitative Reasoning Benchmark for Business and Finance. *arXiv preprint arXiv:2311.06602* (2023). <https://arxiv.org/pdf/2311.06602>
- [10] Bo Li, Gexiang Fang, Yang Yang, Quansen Wang, Wei Ye, Wen Zhao, and Shikun Zhang. 2023. Evaluating ChatGPT’s Information Extraction Capabilities: An Assessment of Performance, Explainability, Calibration, and Faithfulness. *arXiv preprint arXiv:2304.11633* (2023). <https://arxiv.org/abs/2304.11633>
- [11] Tianjun Li, Bowen Zhang, Xiangning Li, et al. 2024. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2405.06682* (2024). <https://arxiv.org/pdf/2405.06682>
- [12] Yifan Li and Kevin M Passino. 2024. FinDKG: Dynamic Financial Knowledge Graph Construction from News with Large Language Models. *arXiv preprint arXiv:2402.02413* (2024). <https://arxiv.org/abs/2402.02413>
- [13] Alejandro Lopez-Lira and Yuehua Tang. 2023. Can chatgpt forecast stock price movements? return predictability and large language models. *arXiv preprint arXiv:2304.07619* (2023).
- [14] Pedro Henrique Martins et al. 2019. Joint learning of named entity recognition and disambiguation. *Proceedings of the 2019 Conference of the European Chapter of the Association for Computational Linguistics (EACL)* (2019).
- [15] Varshini Reddy, Rik Koncel-Kedziorski, Viet Dac Lai, Michael Krumdieck, Charles Lovering, and Chris Tanner. 2024. DocFinQA: A Long-Context Financial Reasoning Dataset. *arXiv preprint arXiv:2401.06915* (2024). <https://arxiv.org/pdf/2401.06915>
- [16] Alfréd Rényi. 1961. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, Vol. 1. University of California Press, 547–561.
- [17] Claude E. Shannon. 1948. A Mathematical Theory of Communication. *Bell System Technical Journal* 27, 3 (1948), 379–423.
- [18] Maxwell Shinn, Andrew Labash, et al. 2023. Reflexion: Language Agents with Verbal Reinforcement Learning. *arXiv preprint arXiv:2303.11366* (2023). <https://arxiv.org/pdf/2303.11366>
- [19] Yang Sui, Yu-Neng Chuang, Guanchu Wang, Jiamu Zhang, Tianyi Zhang, Jiayi Yuan, Hongyi Liu, Andrew Wen, Shaochen Zhong, Hanjie Chen, et al. 2025. Stop overthinking: A survey on efficient reasoning for large language models. *arXiv preprint arXiv:2503.16419* (2025).
- [20] Dimitris Vamvourellis and Dhagash Mehta. 2025. Reasoning or Overthinking: Evaluating Large Language Models on Financial Sentiment Analysis. *arXiv preprint arXiv:2506.04574* (2025).
- [21] Shikhar Vashishth, Prince Jain, and Partha Talukdar. 2018. CESI: Canonicalizing Open Knowledge Bases using Embeddings and Side Information. *Proceedings of the 2018 World Wide Web Conference (WWW)* (2018), 1317–1327.
- [22] Somn Wadhwa, Silvio Amir, and Byron C. Wallace. 2023. Revisiting Relation Extraction in the era of Large Language Models. *arXiv preprint arXiv:2305.05003* (2023). <https://arxiv.org/pdf/2305.05003>
- [23] Jindong Wei, Lema Li, Zuchao Li, and Hai Zhao. 2023. ChatIE: Language Model as an Interactive Information Extractor. *arXiv preprint arXiv:2305.15452* (2023).
- [24] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *arXiv preprint arXiv:2201.11903* (2022). <https://doi.org/10.48550/arXiv.2201.11903>
- [25] Daojian Zeng et al. 2014. Relation classification via convolutional deep neural network. *Proceedings of the 25th International Conference on Computational Linguistics (COLING)* (2014).
- [26] Daojian Zeng et al. 2015. Distant supervision for relation extraction via piecewise convolutional neural networks. *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (2015).
- [27] Bowen Zhang and Harold Soh. 2024. Extract, Define, Canonicalize: An LLM-based Framework for Knowledge Graph Construction. *arXiv preprint arXiv:2404.03868* (2024). <https://arxiv.org/pdf/2404.03868>
- [28] Peter Žukov Gregorič et al. 2018. A neural multi-task learning framework to jointly model medical named entity recognition and normalization. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (2018).