

---

# MODELING THE EVOLUTION OF AI LANGUAGE

---

A Single-Firm Case Study

Marco Siliezar<sup>1, 2,†</sup>

<sup>1</sup> Northwestern University

<sup>2</sup> <https://meatloaf02.github.io>

† [marco.siliezar@gmail.com](mailto:marco.siliezar@gmail.com)

## Abstract

With the popularity and fast evolution of artificial intelligence (AI) over the last several years, enterprise software firms have increasingly redistributed internal resources toward the development and integration of AI technologies (McKinsey & Company 2025). This includes R&D spending, product roadmaps, and strategic priorities. As noted in the McKinsey article, executives of these companies view AI as likely to drive “potentially higher margins and earnings” which translate to increases in shareholder value. This research paper attempts to answer the question, has this shareholder value been recognized?

This research aims to review the evolution of AI related language, product capabilities, market positioning, and risk disclosures from 2015 to present for Workday, Inc., a business-to-business enterprise software company. Information about this company will be extracted from various sources online including regulatory filings, official company press releases, and events. This data will be ingested and stored in a knowledge base. This will aim to contribute to a predictive model for predicting the next month return direction for this company’s stock as a proxy for shareholder value. The work completed so far: developing a structured plan, creating seed URLs from which this information will be extracted, drafting of a knowledge graph schema blueprint, and a ground-truth approach.

# Table of Contents

Abstract .....	i
Introduction and Problem Statement .....	1
Literature Review .....	1
Textual Analysis of SEC Filings as Financial Signal .....	1
Corporate Disclosures and Knowledge Graphs .....	2
AI-Related Language in Disclosures and Risks .....	2
Research Gap and Contribution .....	2
Methods .....	2
Results .....	3
Conclusions.....	4
References.....	4
Appendix.....	5

## Introduction and Problem Statement

Core enterprise software-as-a-service enterprise platform companies sell to other businesses products that aim to increase productivity. Workday is a publicly traded company that builds these products for managing human resources, finances, payroll, and other core aspects of corporate and government operations. These aspects are essential to corporate operations thus there is an inherent efficiency maximization problem that is to be solved. Therefore, Workday is an interesting single-company case study to research whether its use of AI addresses this efficiency problem for its customers.

This research aims to compare the increase in AI-related language over time against the stock price. A predictive model will be powered by a knowledge graph containing information about this company. The predictive model will attempt to predict next month's stock direction. The model will not answer whether the stock should be bought or sold, rather it will output the probability the direction next month is up based on signals and measures derived from the knowledge base.

Likely users of this application will be analysts and investors conducting research who will be looking for support for a trading-style direction rule using this firm's AI strategy as input. The application will not be for users who want to answer the question, "should I buy or sell WDAY?"

## Literature Review

### Textual Analysis of SEC Filings as Financial Signal

Finance researchers have previously studied the link between financial tone of text in 10-K filings to financial performance and outcomes. By creating finance-domain specific word lists, Loughran and McDonald (2011) showed that these lists provide a better understanding of the impact of filing text on stock returns more than generic nonbusiness word lists. The text in filings is an important aspect but more so is its readability. When analyzing the relationship between annual report readability and company performance, Li (2008) found that firms with poor performance have more difficult to read annual reports. Thus, financial tone and readability are signals and are features worth comparing against embeddings.

## Corporate Disclosures and Knowledge Graphs

Financial narrative researchers have previously explored how to process unstructured corporate disclosure documents using natural language processing techniques to store them in knowledge graphs to discover underlying networks. An approach by Cavar (2018) creates a pipeline to process 10-K filings, analyze their semantics using publicly available Python packages for NLP processing, creates a uniform data structure, and implements a knowledge graph for storing raw text as tuples. This demonstrates the feasibility of domain-specific knowledge graphs for storing and analyzing corporate disclosure data.

## AI-Related Language in Disclosures and Risks

The rise of AI-related language has been quantified by industry analysis. In a report by Arize (2024), an AI observability and LLM evaluation company, quantified how often companies mention AI and how often this language is associated with risk. Their report shows that in 2024 the Fortune 500 companies mention AI 250.1% more than in 2022. These companies cited AI as a risk factor more than 473.5% more in 2024 than in 2022. This analysis supports the hypothesis about AI as an opportunity versus a risk.

## Research Gap and Contribution

Prior work demonstrates that language in 10-K filings is informative, provides signals, and representing their information in a knowledge graph, provides relationships that would otherwise be nonobvious. However, there is limited work integrating longitudinal NLP, knowledge graph modeling, and firm-specific context within a single framework. This research contributes by modeling Workday, Inc. as a longitudinal single-firm case study, constructing a knowledge graph that captures the evolution of AI-related language from 2015 to present. This approach enables exploratory signals relevant to investment decision-making.

## Methods

The initial blueprint of the database schema includes entities and edges that are Resource Description Framework (RDF) compatible. Table 1 in the Appendix is a draft entity overview and Table 2 in the Appendix is a draft nodes overview. As the goal for this KG is to deliver buy/sell signal features which originate from information in documents created at a point in time, the database is designed to store documents which contain excerpts about mentions and claims. These mentions will become the evidence for the signal. Events will also become important. Events include filings, earnings calls, acquisition announcements, among

other activity. These events will be time anchors. If analyst reports become obtainable such that they can be stored in the database, they are also documents and measures.

Information sources for the database will come from SEC filings, official company press releases, blog posts, event pages, partner pages, and analyst report landing pages for Workday, its direct competitors, and customers. Information on SEC filings will be extracted via the SEC's EDGAR system using APIs. Information on website pages will be extracted via a focused web crawler developed specifically for this research. The GDELT database will be used for extracting events. Really Simple Syndication (RSS) feeds will be used for incremental updates to the database.

## Results

These information sources will contain only publicly available information outside any paywalls. To minimize noise and poor data, seed pages will be manually tagged on AI-related pages for the prediction model. Additionally, the data will be normalized and deduped since duplicate data is surely to exist as information is sourced from various locations on the web. Although rich information is available from analyst reports, almost all require a paid subscription, credentialed or authenticated access, or are licensed data. As a fallback strategy, public analyst quotes in news articles will be fetched and analyst ratings will be fetched from a publicly available source such as Yahoo Finance. Additionally, to respect websites' Terms of Service, "robots.txt" will be checked before crawling any website domain and rate limiting will be implemented to maintain compliance of websites' terms and conditions.

The MVP scope of this research is to ingest 300 – 500 documents to balance comprehensive coverage with project feasibility. Ideally, these documents are distributed across the analysis period (2015 to present) with the priority of documents sourced from SEC filings. The success criteria in as far as document ingestion is that MVP is achieved when at least one document is ingested for every quarter in the analysis period, all SEC filings are included, and there are at least 30 AI-relevant documents per year at least starting in 2020. From an exploratory analysis of GDELT reveals that this may be viable.

## Conclusions

The knowledge base that will be built for this research is designed to power a predictive model for the analyst and investor with the ability to answer questions such as: “What AI capabilities does Workday claim, where, and how often?”, “Which products are emphasized over time?”, “What risk topics are emerging?”. Additionally, it will power buy and sell signal features from media measures, impacts of events, and shifts in the language of filings over time. This depends on the quality and quantity of data which can be gathered from public sources. At the same time, the timespan of 2015 to present should provide sufficient data for the model. Lastly, this application will not claim certainty for a buy or sell position of a stock in the company, and it will not claim it has solved stock prediction. The success of this application in production for its users paves the way to implement future applications to additional B2B SaaS companies in the information technology sector.

## References

- Arize AI, Inc. 2024. “The Rise of Generative AI in SEC Filings”. <https://arize.com/wp-content/uploads/2024/07/The-Rise-of-Generative-AI-In-SEC-Filings-Arize-AI-Report-2024.pdf>
- Cavar, Damir and Matthew Josefy. 2018. “Mapping Deep NLP to Knowledge Graphs: An Enhanced Approach to Analyzing Corporate Filings with Regulators.” *Proceedings of The first financial narrative processing workshop (FNP 2018)*.
- Li, Feng. 2008. “Annual report readability, current earnings, and earnings persistence”. *Journal of Accounting & Economics* 45, no. 2 – 3 (August 2008): 221 – 247.  
<https://doi.org/10.1016/j.jacceco.2008.02.003>
- Loughran, Tim and Bill McDonald. 2011. “When is a Liability Not a Liability? Textual Analysis, Dictionaries, and 10-Ks”. *The Journal of Finance* 66, no. 1 (February 2011): 35 – 65.  
<https://doi.org/10.1111/j.1540-6261.2010.01625>.
- McKinsey & Company. 2025. “The AI-centric imperative: Navigating the next software frontier. October 16, 2025. <https://www.mckinsey.com/industries/technology-media-and-telecommunications/our-insights/the-ai-centric-imperative-navigating-the-next-software-frontier>

## Appendix

Edges	Examples
document_entity_mention	entity mentioned in excerpt
claim	statement from excerpt
document_event_link	document references an event
event_entity_role	entities involved in event
company_product	product owner
product_capability	capability enabled by product
company_market_segment	where company competes
company_risk_topic	topic of risk exposure
analyst_report_output	recommendation and confidence
media_measures	launch new AI product at earnings

Table 1. Draft Entity Schema

Nodes	Examples
company	Workday; partner; competitor
product	HCM; Financials
technology_capability	AI; platform
market_segment	Higher Ed; Financials; Healthcare; Government
risk_topic	regulatory; competition; security
person	exec; analyst
event	earnings; investor day; product launch
source	SEC/EDGAR; Investor Relations
document	filing; news; press release
excerpt	text within a document

Table 2. Draft Nodes Schema