## ChatGPT

# NLP Pipelines in Investment Research: Applications, Data, and Database Design

## NLP Applications in Investment Research and Analysis

Financial firms increasingly use **Natural Language Processing (NLP)** to aid investment research and analysis. These NLP pipelines help analysts sift through **unstructured text data** (like news articles, reports, and transcripts) and convert it into **structured insights** that can drive decisions [1] . In simple terms, NLP techniques enable computers to read and understand financial texts (written or spoken language) and turn them into useful data points or signals. This addresses the modern "information overload" – with countless news stories, filings, and transcripts daily, **AI/NLP can automatically identify relevant information, trends, or red flags** much faster than a human could [2] [3] . The end goal is to not miss any critical insight that could affect investment decisions.

**Common NLP functions** in finance include *named entity recognition* (identifying key entities like company names, people, figures), *topic or theme extraction* (finding what themes or issues are being discussed), *sentiment analysis* (detecting positive or negative tone), and *information extraction* (pulling out facts like financial metrics or events) [4] . Often these are combined to produce higher-level insights – for example, tracking how the sentiment on a company changes over time on social media, or how the wording in its annual report differs from last year [4] . In practice, NLP "reads" the text and produces structured outputs (like a sentiment score, a list of mentioned companies, or a summary) that an analyst or a model can use [1] .

To illustrate the nature of these applications, here are several **key NLP use-cases in investment research**:

### Sentiment Analysis of News and Social Media

**Sentiment analysis** is widely used to gauge the market's mood. This involves analyzing text to determine if the sentiment is positive, negative, or neutral. Financial news articles and even **social media posts (like tweets or Reddit discussions)** are analyzed in real-time to see how investors feel about a company or market event [5] . For example, a hedge fund might analyze the sentiment of thousands of tweets about a company right before its earnings release. In one case, *analysts detected a surge of positive sentiment on social media about a pharma company before earnings*, indicating high market optimism [6] . At the same time, they used topic modeling on news articles which highlighted the company's "innovative drug pipeline," reinforcing a positive outlook [6] . Acting on these NLP insights, the fund took a bullish position – and indeed the company's earnings beat expectations, so the stock jumped, validating the NLP analysis [7] . This example shows how **combining social-media sentiment with news topic analysis** can give an investment edge.

Industry practitioners confirm that **sentiment from text is a crucial signal**. Buy-side investment teams routinely monitor sentiment and topical trends in news articles or tweets to inform their decisions [8] . Many firms purchase **news analytics feeds** from vendors that score news stories in real time for sentiment

and relevance. For instance, Refinitiv (formerly Thomson Reuters) and Bloomberg offer such feeds, and companies like RavenPack specialize in tagging news with sentiment and entities for quantitative trading. The appeal is that **NLP can flag bullish or bearish news faster** than humans – helping portfolio managers not miss market-moving information. Sentiment analysis is also applied to **financial discussion forums** and blogs to gauge retail investor mood (e.g. detecting if a stock is "meme-worthy"). In fact, industry research notes that investors now *"systematically track social media, news and blogs"* as part of their market monitoring strategy [9] . All of these sources are unstructured text streams that NLP can transform into sentiment indices or alerts.

Importantly, sentiment analysis isn't just about labeling text positive/negative – it also helps detect **anomalies or inconsistency** in narratives. For example, if a company's management commentary is normally upbeat but suddenly becomes cautious, an NLP sentiment tool would flag that change. One financial research platform notes that NLP-based sentiment features can *"spot anomalies or inconsistencies in company documents,"* alerting analysts to possible issues that merit deeper investigation [10] . Overall, sentiment NLP pipelines help researchers keep a pulse on market emotion and spot turning points in sentiment that could precede price movements.

## Analyzing Earnings Call Transcripts

**Earnings call transcripts** (the text of quarterly conference calls between company executives and analysts) are a goldmine of qualitative information. NLP pipelines are used to **automatically analyze these transcripts** for tone, content, and clues about the company's outlook. Early applications simply tried to summarize or skim transcripts to save analysts' time. Now, more advanced NLP looks for linguistic patterns that might reveal management's true feelings. For instance, NLP models can pick up on **hesitation, evasive phrasing, or overly optimistic language** by executives. One fintech firm explained that NLP can analyze the language of an earnings call and even apply *psycholinguistics* – essentially gauging executives' emotional state or confidence level from their words [11] . By doing so, analysts can infer if management is hiding concern or is genuinely confident. The **qualitative "value signals" from how things are said** can be quantified with NLP [12] .

In practice, an NLP pipeline might flag if a CEO's tone is significantly more negative this quarter compared to last, or if they are using more uncertain words (like "we hope" or "maybe"). These subtleties can indicate future performance issues. **Financial firms use such transcript analysis to augment their research** – for example, constructing a "behavioral profile" of executives. AlphaSense (a research platform) noted that beyond just summarizing transcripts, NLP can *"establish inferences about executives' character and state of mind"* from their spoken words [11] . If NLP quantifies a drop in confidence or a spike in anxious language, analysts might treat that as a warning signal. Conversely, consistently upbeat language could reinforce a positive thesis. Some funds even integrate these **transcript sentiment scores** into trading algorithms. Overall, NLP turns the once-hard-to-quantify aspects of an earnings call (tone, sentiment, topic frequency) into data that can be tracked and compared.

Beyond sentiment, NLP also helps with **information retrieval from transcripts**. For example, theme extraction tools can pull out how often certain topics or keywords were mentioned on calls. Analysts use this to see what management teams are focusing on or worrying about. An industry use-case described *tracking the mention counts* of key topics in calls over time [13] – e.g. how often "supply chain" is mentioned each quarter. If mentions of a topic are increasing across many companies, it might signal an emerging trend or concern in that industry. One real-world example: SAP's strategy team used NLP on transcripts to

confirm that interest in "digital transformation" topics was rising among company executives, validating a hypothesis that digitally mature companies perform better [14] . This shows how **counting and analyzing the language in transcripts** can reveal trend lines that inform investment theses.

## Parsing Financial Filings and Reports

**Financial reports and regulatory filings** (like annual 10-K reports, quarterly 10-Qs, or other disclosures) contain a mixture of structured numbers and lengthy text discussions. NLP is used to **extract and analyze the text portions** of these documents. This can include the Management Discussion & Analysis (MD&A), risk factor sections, or footnotes in filings. The goal is often to detect changes and signals in a company's own disclosures. For example, NLP can do a *"change analysis"* on annual reports – highlighting what changed in the language year-over-year (e.g., a new risk factor added, or tone of outlook section turning more cautious) [4] . These changes can be early warnings of issues. NLP can also classify and tag sections of reports (like identify all accounting policy changes mentioned, or all forward-looking statements). This saves analysts from manually combing through hundreds of pages per filing.

Another important application is **key metric extraction**. Rather than reading tables, NLP can be trained to find specific figures or ratios in the text. For instance, an NLP system might automatically pull out *revenue, profit, or cash flow numbers* from a 10-K and populate a database. One source explains that NLP can *"quickly pull data like revenue, cash flow, and debt ratios from financial statements"*, automating what used to be a manual data entry task [15] . This not only speeds up research but also reduces errors. For example, if an analyst wants to compare **debt levels** of companies, an NLP pipeline could extract the latest debt figures from each company's filings and store them for easy querying, instead of the analyst opening each report.

NLP on filings is also used for **sentiment and risk analysis** similar to news. Researchers have created **"tone" metrics for annual reports**, where the text is scored for positivity/negativity (an early famous example is counting occurrences of negative words in 10-Ks to predict stock volatility). Today's NLP goes further by understanding context – e.g., distinguishing negative words in a benign context versus actual risk discussion. Some asset managers incorporate an *"optimism score"* or a *"readability score"* of a company's filings into their evaluation. For instance, a very complicated, jargon-heavy report might score poorly on readability (which studies have linked to potential red flags), and NLP can quantify that complexity. Another modern trend is **ESG (Environmental, Social, Governance) analysis**: NLP tools scan filings and reports for mentions of ESG topics (like "climate change" or "diversity") and even gauge how substantively the company addresses them. This helps investors compare companies on qualitative factors that are hard to measure otherwise.

## Thematic Search and Market Intelligence

Investment research often boils down to **finding relevant information quickly**: "What has been said about X topic or Y company lately?" NLP-powered search and discovery tools are therefore crucial. Unlike a simple keyword search, modern NLP-based search engines can understand synonyms, context, and even provide **summaries** of the results. Financial firms deploy internal search platforms (or use vendors like AlphaSense, Sentieo, etc.) where an analyst can query a topic and the system will search a vast document database (news, filings, transcripts, research reports) for relevant mentions. These platforms use NLP for features like **synonym recognition** (so if you search "inflation", it also finds "price increases"), **entity linking** (knowing that "Apple" refers to Apple Inc. and not just the fruit), and **intelligent filtering** (e.g. filter results by company or industry automatically).

One advanced feature described is a **"Search Summary"**, which aggregates how a topic is discussed across many documents [16] [17]. For example, analysts looked at **European companies' earnings calls discussing share buybacks**. Using an NLP search tool, they filtered to European companies and got over 800 mentions of "buybacks" in Q2 of 2021 [18]. The tool then broke down the results by country, industry, and even sentiment, giving a macro view of how prevalent and in what tone the topic appeared [19]. This kind of **NLP-driven thematic analysis** helps analysts spot trends ("Are more companies talking about supply chain issues this quarter?") and identify which companies or sectors are most involved.

Another example: M&A bankers might use **topic extraction** to find information not obvious in standard financials. AlphaSense noted that bankers used NLP *"theme extraction to find financial metrics not found in standard 10-Ks and parse alternative sources for companies in emerging industries"* [20]. In other words, if a company is new or not well-covered, bankers can use NLP search to gather intelligence from news, patents, or niche reports about it. NLP can also generate alerts – e.g., if a key topic or keyword appears in any source (say a specific lawsuit or a regulatory change), the system can automatically notify the users. This greatly **enhances situational awareness**.

In essence, NLP turns a big text repository into a **financial knowledge base** that can be queried. The search engine doesn't just retrieve raw documents; it often provides a snippet or summary of why that document is relevant. According to one source, an NLP-powered financial search engine will *"extract key components, concepts, and notions"* from documents and then present **a summary of the most relevant information** for the query [21]. This saves the analyst from reading dozens of full documents. For example, an analyst could ask, "What are the main points from all analyst reports about Tesla's latest earnings?" and the system could return a synthesized summary of common themes. Such **automated summarization** is another NLP capability used in investment research, especially now with advanced language models.

## Risk Monitoring and Compliance (NLP for Red Flags)

Although slightly beyond pure "investment research," many firms use NLP for **risk analysis** which ultimately feeds into investment decisions (avoiding losses, compliance issues, etc.). NLP can monitor news and reports for any mention of potential risks related to portfolio companies – for instance, scanning for phrases that indicate accounting irregularities, legal troubles, or regulatory actions. An example is **monitoring regulatory updates**: NLP systems ingest notices from regulators (like the SEC or Federal Reserve) and flag anything relevant to the firm's investments [22]. If the Federal Reserve releases meeting minutes, NLP can analyze the language to gauge how hawkish or dovish it is, which is valuable insight for bond traders.

In compliance, NLP automates reading of things like **communications or filings for violations**. For instance, banks use NLP to scan emails or chat transcripts of traders for compliance triggers (certain words that indicate policy breaches). While this is more on the operational risk side, it employs the same technology – reading unstructured text (emails) and flagging meaning.

From an investment perspective, **counterparty risk** can be monitored by NLP: scanning news about key partners, suppliers, or borrowers for any hints of trouble (say, "bankruptcy", "downgrade", or "fraud" mentions) [23]. If an important supplier to several companies is getting negative news coverage about financial stress, an NLP alert could prompt analysts to investigate impact on their portfolio. Similarly, **fraud detection** uses NLP to find anomalies in textual data (like unusual language in invoices or reports) that might indicate something is off [24].

In summary, the nature of NLP applications in investment research is that they **automate the reading and analysis of text-based information** that investors have always cared about, but at a much larger scale and in real-time. Whether it's gauging the sentiment of the entire market, extracting facts from a single filing, or summarizing a week's worth of news on a sector, NLP pipelines aim to provide **comprehensive, up-to-date insights** without requiring humans to read every word. This helps analysts generate ideas, validate hypotheses, manage risks, and ultimately make more informed investment decisions [25].

## Key Data Types and Sources for Financial NLP Pipelines

To build NLP pipelines like those above, one needs to gather a variety of **financial textual data**. The nature of each data source and how it's used can differ greatly. Below is an overview of the main types of data used in investment-focused NLP, along with details on their characteristics, availability, and whether they are structured or unstructured.

- **Financial News Articles:** News is a primary input for many NLP models in finance. News articles are **unstructured text** written by journalists or newswires about companies, markets, economics, etc. They tend to be relatively short-form (a few hundred words) and published in real-time. The content can range from breaking news (e.g., a CEO resignation) to analysis pieces. **Nature:** News text often contains facts, quotes, and a narrative of an event. It can be parsed for sentiment (positive/negative tone), for specific events (mergers, product launches), or for topics mentioned. **Data examples:** This includes feeds from Reuters, Bloomberg, CNBC, Financial Times, Wall Street Journal, etc. **Sources and availability:** Many news sites offer some free articles or RSS feeds, but comprehensive real-time data is usually via subscriptions (e.g., a Bloomberg Terminal or Refinitiv feed). For a class project, one can use *publicly available news* from sources like Yahoo Finance (which aggregates news stories on companies) or free APIs like NewsAPI (which provides news article JSON for given queries, with limitations). Additionally, some datasets of historical news headlines are available (for example, Kaggle has had datasets of news headlines for stock market prediction). Overall, news articles are **publicly available to an extent** – headlines and summaries are often free, full text sometimes behind paywalls – but at least some open data can be gathered for academic use. In all cases, news articles are **unstructured** (plain text), though often with metadata like date, source, and tickers.

- **Social Media and Online Forums:** Social platforms have become important data sources for sentiment and trend analysis. These include **Twitter posts (tweets)**, Reddit posts (e.g., the WallStreetBets subreddit), stock forums, and blogs. **Nature:** Social media text is *highly unstructured*, very informal, and often noisy. Posts may include slang, emojis, abbreviations (e.g., "$APPL to the moon "). They are usually short (280 characters for tweets, longer for Reddit comments) but extremely high-volume. Despite the noise, social media can reveal retail investor sentiment, rumors, and emerging narratives (for example, viral discussions that can drive stock momentum, as seen in the GameStop saga). NLP techniques like sentiment analysis and topic clustering are applied here to gauge the crowd's mood. **Data examples:** Tweets mentioning certain stock tickers, Reddit threads about a sector, blog articles analyzing a company. **Sources and availability:** Social media data is **publicly accessible but with some restrictions**. Twitter provides a public API for researchers (however, as of 2023–2025, API access has become more limited or paid for large volumes). Reddit data can often be accessed via APIs or public data dumps (the Pushshift.io dataset archives Reddit comments, for instance). There have been public examples of using Twitter for finance – *studies show incorporating social media and news data can improve market monitoring* [9]. For a class project, one might collect tweets using a developer API key (limited amounts) or use an existing public dataset

(some researchers have shared collections of tweets or Reddit posts related to finance). Social/blog data is **unstructured text**. It often requires extra cleaning (to remove spam, handle misspellings, etc.), but it's very useful once processed, since it provides a gauge of *public/retail opinion and chatter*, which traditional data lacks.

- **Company Filings and Financial Reports:** These are official documents that companies are required to file (in the U.S., with the SEC) or release periodically. Key examples are the **10-K (annual report)**, **10-Q (quarterly report)**, 8-K (current report for major events), annual shareholder reports, prospectuses, and other regulatory filings. **Nature:** Filings are **semi-structured**. They contain a lot of text in a structured format – sections with headings (business overview, risk factors, MD&A, financial statements, etc.). They also contain tables of numbers (the financial statements themselves), but from an NLP perspective, we focus on the textual sections. These texts are formal, often legalistic, and can be very lengthy (hundreds of pages). They're a rich source of information about a company's performance, strategies, and risks in its own words. NLP can extract specific items (like risk factor keywords, or numerical values embedded in text), analyze sentiment or tone (e.g., is the language more negative this year?), and compare changes over time. **Data examples:** A 10-K filing's risk factor section to see if new risks were added, or the CEO's letter to shareholders for tone. **Sources and availability:** In the US, these filings are **publicly available via the SEC's EDGAR database**. EDGAR provides free access to millions of filings [26] . One can search and download filings by company name or ticker. This is a great resource for a class project because you can get actual 10-K/10-Q texts for many companies without cost. Outside the US, many countries have similar public filing systems. Some companies also post these reports on their investor relations websites. Filings often come in HTML or PDF formats; extracting text from them might require parsing. Overall, filings are **public and structured in format but contain a lot of unstructured text** within them. They are a prime candidate for building an NLP database (e.g., a database of all 10-K risk factor texts for S&P 500 companies).

- **Earnings Call Transcripts:** These are transcripts of conference calls (or webcast presentations) where company management discusses financial results with analysts. Many public companies hold these calls each quarter. **Nature:** Transcripts are unstructured text, essentially dialogue. They have a Q&A format: prepared remarks by executives followed by analysts' questions and answers. The language is spoken (and then transcribed), which means it can be more informal or fragmented than written reports. They often contain valuable forward-looking statements, clarifications of results, and the tone or attitude of management under questioning. NLP uses transcripts for sentiment analysis (how positive or negative does management sound), emotion or **tone analysis** (e.g., measuring stress or confidence), and content extraction (what topics are discussed, how frequently). **Data examples:** For instance, extracting all statements about "pricing pressure" from the last call, or computing a sentiment score for each answer the CEO gives. **Sources and availability: Transcripts can be a bit tricky to source publicly.** Many are provided by services like Thomson Reuters, FactSet, or Seeking Alpha. SeekingAlpha provides free transcripts for a short time after the call (they are crowdsourced or provided as a courtesy; older ones may require a subscription). Sometimes the company's own website will post a transcript or at least an audio recording (larger companies often post transcripts in their IR section). There are also third-party repositories – for example, Nasdaq's website or Motley Fool occasionally publishes transcripts as articles. For a project, one could gather a few transcripts by manually copying them from such sources, or use any publicly released transcripts. Additionally, some academic datasets or libraries have collected a bunch of transcripts (you might find on Kaggle or other forums datasets like "Earnings calls transcripts dataset"). The

transcripts are **unstructured text**. They may need cleaning (speaker names, timestamps, etc., sometimes need to be removed). Despite not being as readily API-available as filings, they are considered public information (since anyone can listen to the call, and transcripts often circulate in the public domain shortly after). Analysts and NLP models frequently leverage transcripts – *they are cited as a frequently leveraged text source alongside filings and social media* [27] .

- **Analyst Research Reports:** These are reports written by financial analysts (e.g., equity research reports from banks or independent research firms) giving their analysis and opinions on a company or sector. **Nature:** They are unstructured text documents, often in PDF, with a mix of narrative, charts, and tables. They contain analysts' insights, forecasts, and ratings. For NLP, these can be used to gauge consensus opinions or to extract data points (like price targets, or key arguments pro/con about a stock). They might also be used as a knowledge base for question-answering systems (e.g., answering, "What is the expected growth for Company X next year?" by finding it in a report). **Data examples:** An example is extracting all analyst price upgrade/downgrade rationales. **Sources and availability: Generally, these are *not publicly available*** for free – they are proprietary content from brokerages and usually require subscription (like via Bloomberg or Refinitiv or the brokerage itself). Because this is a class project with only public data, you likely *cannot get a large corpus of real analyst reports* easily. Sometimes, companies will include excerpts of analyst comments in press releases or news articles ("Analyst at XYZ Bank says …"), which are public. Alternatively, older reports might be found on forums or via library databases if your university has access. But for designing a database, you might skip this due to access issues. If included, treat them as **unstructured, private text**. In industry, some hedge funds do NLP on these reports (scanning for sentiment or extracting mentioned companies, etc.), but for academic purposes, you'd focus on more accessible data like news and filings.

- **Press Releases and Corporate Announcements:** These are official statements issued by companies (or other entities like regulators). **Nature:** They are structured as formal announcements – often template-like (headline, date, location, body text, quotes from executives, etc.). Press releases about earnings, product launches, M&A deals, leadership changes, etc., are valuable to parse. They often contain the first news of an event (e.g., "Company X to acquire Company Y"). NLP can be used to immediately interpret a press release – for instance, categorize it (merger announcement vs. earnings vs. partnership), extract key details (who, what, when, financial terms), and even assess market sentiment (though press releases are usually upbeat by nature, since they're corporate PR). **Data examples:** An example is automatically extracting from an earnings press release the reported revenue and comparing it to consensus estimates (some NLP/finance systems do this to generate instant alerts on earnings beats or misses). **Sources and availability:** Press releases are **publicly available**. Companies post them on their websites (often under "Media" or "Investor" sections). They are also distributed via newswire services like PR Newswire or Business Wire (these often show up on finance portals). Many financial data platforms aggregate press releases (Yahoo Finance often lists recent press releases of a company). They are free to use, generally, since they are public communications. The text is unstructured but follows a predictable format, and key data may be embedded in the text.

- **Economic and Market Data (Textual):** Not all data is company-specific. There are textual data like **economic reports (e.g., central bank statements, economic outlook reports)** and **market commentary**. **Nature:** For example, the U.S. Federal Reserve's meeting minutes or statements are text that can be parsed for hints about policy changes (NLP can gauge how the language changed –

a well-known practice is to compare Fed statements word-for-word). Similarly, articles or reports about macroeconomic conditions can be analyzed to create sentiment or topic indicators for the economy. **Sources:** Central bank websites (Federal Reserve, ECB, etc.) provide statements and minutes publicly. International organizations (IMF, etc.) publish reports. These are public and structured as documents/PDFs (unstructured text inside).

In summary, most **data feeding financial NLP pipelines is unstructured text**, ranging from news and social media to formal corporate filings [28] . According to industry estimates, about *80-90% of useful data in finance is unstructured text* [29] – which is why NLP is so crucial. A report by the London Stock Exchange Group noted that **news, documents, earnings calls, and even alternative textual sources are being leveraged heavily** as firms look for new insights [28] [27] . For a class project, **publicly available sources** like SEC's EDGAR (for filings) [26] , company websites (for press releases and maybe some reports), financial news APIs or open datasets, and social media (via Twitter/Reddit APIs or datasets) are your go-to. All these texts will typically be ingested in raw form (perhaps HTML/PDF to text) and then processed by your NLP pipeline.

It's worth noting that along with the text itself, you will often collect **metadata**: e.g., for a news article, metadata is the date/time, source, tickers/companies involved; for a tweet, the timestamp and author; for a filing, the company and filing date, etc. This metadata is important for organizing the database and enabling queries (like "news in the last 7 days" or "filings for company X in 2021").

## Designing a Database for Financial NLP Data

Building a database to support these NLP applications involves deciding how to **store and organize both the raw text and the structured insights extracted from it**. The design must accommodate large volumes of unstructured data (documents, transcripts, etc.) as well as structured data (like numbers, dates, categories, and relationships that the NLP pipeline produces). Below is a general architectural overview and some design considerations for such a database.

**1. Overall Architecture (Pipeline -> Database):** Typically, the system will have an **ingestion pipeline** that collects data from various sources (news API, EDGAR filings, etc.), applies NLP processing, and then stores results in a database or search index. In industry, this often looks like: *Text Data Acquisition → Preprocessing/ Cleaning → NLP Analysis/Tagging → Storage/Indexing* [30] [31] . For example, you might gather a bunch of articles and transcripts, convert PDFs to text, clean the text, then run NLP models (for sentiment, entity extraction, etc.), and finally store both the original text and the extracted information. The **database design** needs to handle two things: storing the *documents* in a way that they can be searched (by keywords, by date, by company, etc.), and storing the *extracted structured data* (like a sentiment score, or a list of entities in each document).

**2. Choosing a Storage Paradigm:** Traditional relational databases (SQL databases) are good for structured data, but not ideal for full-text search. Often, financial firms use a combination of a **search engine or NoSQL store for documents** and a **relational or analytical database for structured facts**. For instance, a common choice is to use **Elasticsearch** (a search/indexing engine) to store text documents with their metadata, since it allows very fast keyword and semantic search over the text [32] . Elasticsearch can index every word and also store fields like date, source, sentiment score, etc., enabling queries like "find news articles about *inflation* in 2023 with negative sentiment". Meanwhile, for structured data (like numerical financial data or aggregated indicators), something like a SQL database or a time-series database might be used. In modern data architecture, these could both be parts of a data lake or warehouse. Indeed, an

industry report notes that many are using **Elasticsearch or cloud databases (AWS Athena, Google BigQuery)** in their NLP pipelines [32] – Elasticsearch for text search and BigQuery/Athena for scalable querying of structured data in the cloud.

For a class project, you might simplify and use one database that can handle text and structure. Some relational databases (PostgreSQL, MySQL) have full-text search capabilities and JSON fields which can store unstructured data. Or you could use an open-source search engine (like whoosh or Elasticsearch) for the text and a small relational DB for metadata. The design choice should ensure you can retrieve documents by content and also do analytical queries on the extracted data.

**3. Schema Design (Data Model):** In a relational sense, you would likely have tables (or collections, if using a NoSQL document store) that represent the main entities:

- A **Documents table** (or collection) that stores each text document (news article, transcript, filing, etc.). This would have fields like: *document_id*, *type* (news, filing, tweet, etc.), *date*, *source*, *title/ headline*, and either the full text or a pointer to where the text is stored (sometimes large texts might be stored as files, but often you can store them in a TEXT column). You would also include fields for any readily known structured info – for example, a news article might have a *tickers* field listing the companies mentioned (if known), or a transcript might have a *company_id* and *quarter*. This is your primary storage of raw content. It can be quite large, so indexing is important. In a search engine context, the document store *is* the index.

- **Entities/Reference tables:** If your NLP pipeline does entity recognition (e.g., finds that an article mentions *Apple Inc.* or *CEO John Doe*), you might have a table for *Companies* (with a unique ID, ticker, name, etc.), a table for *People* (if tracking executives), etc. Then you might have a linking table like *Document_Company* to link documents to the companies they mention. For instance, all news articles that mention Apple would have entries linking to Apple's ID. This would allow queries like "give me all documents related to Apple". If you only care about companies, you might not need a separate table and just store a list of tickers in the document record itself. But having a company master table is useful if integrating with stock price data later.

- **NLP Output tables:** Depending on what analyses you do, you may store results in separate tables or fields. For example, if you do **sentiment analysis** on each document, you could add a column *sentiment_score* (or category) to the Documents table. That way you can query `WHERE sentiment_score < 0` to find negative sentiment docs. If you do **topic classification** (say you assign each document a topic like "Market" vs "Product News" vs "Earnings"), that can also be a field or a separate classification table. For **numeric data extracted**, like if you parse a 10-K and pull out revenue, you might have a table of FinancialMetrics (doc_id, metric_name, value) or you could store key metrics as columns for certain doc types (e.g., revenue, EPS for earnings releases). The design can be adjusted based on how you'll query it. For flexible analysis, storing each extracted fact as a row (with doc reference) in a separate table is more normalized. But for easier grabbing of a few things, adding columns to the document is fine.

- **Time-series Data integration:** Often, after extracting signals from text, you want to compare with market data (prices, returns). It could be beyond scope, but you might have a Price table for stock prices (date, ticker, closing price, etc.) if needed. This would let you do queries like "show me sentiment vs stock price over time." Since the question is more about NLP and text, you may not

delve deep here, but be aware that a complete system often merges text-derived data with traditional financial data. In fact, practitioners backtest NLP-generated signals by combining them with structured historical price data [33] .

- **Example of a simple schema:** For instance, a **NewsArticles table** and a **Transcripts table** could both feed into a unified view via a Documents-like abstraction, or they could be separate if their fields differ a lot. Let's say we keep them together with a type flag. The table might look like: `Documents(doc_id, doc_type, date, company_id, source, title, text, sentiment_score, language_complexity, etc.)`.
  Then a `Companies(company_id, ticker, company_name, sector, etc.)`.
  Then a linking like `Document_Company(doc_id, company_id)` if a document can involve multiple companies (news often does mention several companies).
  If using a NoSQL/document store approach, you might store each document as a JSON with these fields and an array of mentioned companies, etc. The key is to allow flexible queries.

**4. Indexing and Query Capability:** In designing the database, think about **what questions you need to answer (see next section)**. If you need full-text search (like find all documents where text contains "inflation"), you must have an index on the text. Dedicated search engines automatically index text for this. In a SQL database, you'd use full-text indexes. Also, index fields like date and company for filtering. For example, you might index the combination of (company_id, date) on the Documents table so that queries for documents about a certain company in a date range are fast. The database should also be optimized for updating with new documents (news comes daily). Modern pipelines often continuously load new data and update the indexes (for real-time dashboards).

**5. Handling Unstructured Data in the DB:** As mentioned, storing large raw texts can be done directly (many databases can handle text blobs) but search is the bigger issue. Often the solution is to use a search engine in tandem. An alternative approach is to generate **vector embeddings** of documents (via NLP models) and store those for semantic search – but that might be beyond the scope here. However, note that some financial firms do convert text to embeddings and store them to power semantic similarity searches (especially with the rise of language models like BERT). If one were to incorporate that, the database might have an Embeddings table mapping doc_id to an embedding vector, and a specialized vector search service.

**6. Data Lake vs Database:** In industry, it's common to keep a **data lake** of raw text files (just to have all source data stored), then have processed outputs in databases. For your project, you can treat the database as both the repository of content and results. Just be mindful of storage considerations.

**7. Example – Combining everything:** Imagine you have collected 10,000 news articles and 500 transcripts. You run an NLP pipeline that tags each document with the companies mentioned, calculates a sentiment score, and extracts any financial metrics. In your database, you'd have those 10,500 records in Documents. Each has a type (news or transcript), a date, etc. Suppose 1,000 of those mention Apple – in Document_Company you'll have ~1,000 rows linking those docs to Apple's company_id. Suppose you extracted 200 metric values (like revenue from some filings or calls) – those go into Metrics table (with references to doc and maybe to a standardized metric name). Now an analyst can query the database: for example, "Give me all negative-sentiment news articles about Apple in 2023" – this query would look up Apple's id, filter Documents where company_id = Apple and sentiment_score < 0 and date between 2023-01-01 and 2023-12-31 (for instance), and return the titles and dates. Or "find transcripts where the CFO

mentioned 'supply chain' more than 5 times" – this might require having stored the count of a keyword per document or running a text search query.

**8. Ensure flexibility:** The design should be **flexible to add new NLP features**. For example, if later you decide to add a "risk flag" field (maybe NLP detects if a document contains some risk keywords), you should be able to add that without redesigning everything – ideally just another column or table. This is why a good design often separates the core document info and the various analysis outputs (one-to-one or one-to-many relationships).

In practice, many financial institutions build central data platforms where all unstructured data is aggregated, and they often rely on specialized tools to query it. They might have a **central NLP pipeline** that feeds many use cases [34] [35]. For a project, a well-structured relational database or a combination of a SQL DB + search library will do the job.

To summarize the database design: **store each document with metadata (date, type, source, etc.), link documents to key entities like companies, and store NLP-generated fields (sentiments, topics, etc.) for each document.** Use indexes or specialized stores to enable fast text queries. The result is a dataset where you can ask complex questions across thousands of documents in seconds, rather than a human spending days reading them.

*(Note: If using cloud services, one could use something like a BigQuery table for all documents, since it can do SQL queries on text and has some full-text search functions. But on a local setup, a combination of PostgreSQL and an Elasticsearch might be a good architecture.)*

## Example Queries on the NLP Data (in English)

Finally, to illustrate how investment analysts or systems would query this database, here are some **important queries you might want to run**, expressed in plain English. These queries demonstrate the kinds of questions a well-designed financial NLP database can answer:

- **Query 1:** *"What is the recent sentiment on Company X in the news?"* – (This would retrieve news articles about Company X from the last, say, one month and compute or list their sentiment scores, to show if coverage is mostly positive or negative recently.)

- **Query 2:** *"Show me all negative news headlines for Tech sector companies in the past week."* – (This would filter the Documents for type=news, date within 7 days, sector = Technology for the company mentioned, and sentiment_score < 0, then output the headlines and company names. This could help identify potential problem news in the tech sector quickly.)

- **Query 3:** *"Find mentions of 'supply chain issues' in any S&P 500 company filings or transcripts in 2023."* – (This involves a full-text search across 10-Ks, 10-Qs, and earnings call transcripts for the phrase "supply chain". The result might list which companies and documents discussed supply chain problems, indicating broader trends or specific companies impacted.)

- **Query 4:** *"Which companies had a significant drop in management sentiment in their last earnings call?"* – (Here the system would compare sentiment or tone metrics from the last earnings call to previous

calls. It would output companies where the sentiment score fell sharply. For example, if last quarter's call was very positive and this quarter's is very cautious, the query would flag that. This helps pinpoint companies that might be facing new challenges.)

• **Query 5:** *"List the top 5 topics mentioned in Banking industry earnings calls this quarter."* – (This would aggregate transcripts of all banks for the current quarter, run topic extraction or simply count keywords. It might output topics like "interest rates", "loan loss reserves", "regulatory capital", etc., possibly with the frequency. Essentially, it tells an analyst what themes are most discussed by bank executives recently.)

• **Query 6:** *"Give me a summary of all news and filings related to Company Y's merger announcement."* – (Suppose Company Y announced a merger. This query would gather the press release, any news stories, and any 8-K filings about that event, and ideally summarize them. While automated summarization might be advanced, the system could at least compile the relevant documents. An analyst could then quickly read the key points – or an NLP summary if available.)

• **Query 7:** *"What are the common risk factors mentioned by semiconductor companies in their annual reports?"* – (This would search the Risk Factors sections of 10-K filings for all companies in the semiconductor industry and find recurring phrases or topics. The answer might be, for example, "supply chain disruptions, foreign export regulations, cyclical demand, competition from XYZ…" along with which companies mention them. It gives a sector-level view of concerns.)

• **Query 8:** *"Which stocks showed a spike in negative social media sentiment right before a drop in stock price?"* – (This is a more complex analytical query. It would require linking sentiment data from social media with stock price data. Essentially, for each company, look for instances where Twitter/Reddit sentiment turned sharply negative and then the stock fell soon after. This can validate that your sentiment metric has predictive power or identify cases like scandal or bad news that first broke on social media. Running this query assumes you have both sentiment time-series and price time-series in your database and can correlate them.)

• **Query 9:** *"Retrieve all press releases of Company Z in the last year and highlight any that contain the word 'partnership'."* – (This fetches Company Z's press releases (which would be stored as documents with type "press_release") over the last year, and then filters or marks those mentioning "partnership". This could be useful if, say, an analyst wants to quickly see all partnership announcements by the company to evaluate its strategy.)

• **Query 10:** *"Has the tone of Federal Reserve statements changed regarding inflation in the past 6 months?"* – (This query targets a specific textual source – the Fed's statements – and would use NLP (perhaps sentiment or a custom dictionary search for hawkish/dovish terms) to compare statements over time. The result might be an analysis showing that the language has become more aggressive on inflation, which is useful for macro investment decisions.)

Each of these English queries would translate to a combination of text search, filters, and possibly joins on the database we designed. For instance, Query 1 might be implemented by searching the Documents for company_id = X and type=news, then looking at sentiment_score. Query 5 would involve filtering by industry and date, then doing a frequency count of keywords in transcripts. The key point is that **with a**

**well-structured database of financial text and metadata, you can ask complex questions that combine content and context**.

In real-world use, these queries could be done via a user interface (like a dashboard or search bar where the back-end runs the SQL/full-text search) or by data analysts writing queries. The examples above reflect common analyst tasks: monitoring sentiment, finding relevant discussions on a topic, comparing language across time or peers, and connecting textual signals to numeric outcomes.

By focusing on publicly available data and using NLP to structure it, even a class project database can demonstrate these capabilities. The end result is akin to a mini "financial intelligence platform" where one can query the collective knowledge hidden in unstructured financial texts and gain actionable insights, much like professional platforms do [36] .

**Sources:**

1. AlphaSense – *"3 Ways to Apply NLP in Financial Research"* (examples of earnings call analysis and trend detection) [37]  [20]

2. Needl.ai – *"Top 9 Applications of NLP in Finance"* (use cases such as sentiment-driven investing and NLP-powered search engines) [6]  [21]

3. Phoenix Strategy Group – *"How NLP Improves Financial Data Analysis"* (overview of NLP tasks like entity recognition, sentiment, and their value in finance) [38]  [39]

4. LSEG (London Stock Exchange Group) Report – *"NLP in Financial Services"* (industry perspective on data sources and adoption of NLP, noting that **news, documents, earnings calls, and social media** are heavily leveraged data sources) [28]  [9]

5. SEC Investor.gov – Definition of EDGAR (confirmation that SEC filings are freely available public data) [26]

6. Refinitiv/Financial Industry Insights – (Not directly quoted above, but referenced in context: combining unstructured and structured data for signals [33]  and typical tools like Elasticsearch, BigQuery for text data storage [32] ).

---

1  4  8  9  27  28  29  30  31  32  33  34  35  NLP in Financial Services
https://www.lseg.com/content/dam/lseg/en_us/documents/research-findings/nlp-in-financial-services.pdf

2  3  10  11  12  13  14  16  17  18  19  20  36  37  3 Ways to Apply Natural Language Processing (NLP) in Financial Research
https://www.alpha-sense.com/blog/product/natural-language-processing-financial-research/

5  15  22  23  24  38  39  How NLP Improves Financial Data Analysis - Phoenix Strategy Group
https://www.phoenixstrategy.group/blog/how-nlp-improves-financial-data-analysis

6  7  21  25  NLP in Finance: Empowering Data-Driven Insights and Decision-Making | Needl.ai
https://www.needl.ai/blog/top-8-applications-of-nlp-in-finance

26  EDGAR | Investor.gov
https://www.investor.gov/introduction-investing/investing-basics/glossary/edgar