



QUANTS LEARN TO READ (FINALLY)

Harindra de Silva

The history of quantitative finance is fundamentally a history of information processing. For decades, the dominant paradigm of systematic investing was entirely numerical. Quantitative researchers built their empires on the statistical analysis of price, volume, volatility, and structured accounting data. Yet, while mathematical models grew increasingly sophisticated, a vast ocean of market-moving information remained largely inaccessible to algorithms: the written word.

Financial markets are driven as much by narrative as they are by numbers. The subtle linguistic shifts in a central bank press release, the defensive tone of a chief executive officer during a quarterly earnings call, and the heavily lawyered disclosures buried deep within the footnotes of an annual report contain critical signals regarding future asset performance. For early quantitative funds, this unstructured textual data represented the ultimate frontier — a reservoir of untapped alpha. However, teaching machines to “read” financial text and extract actionable, statistically significant insights has required a multi-decade technological odyssey.

This report documents the evolution of textual analysis in quantitative finance. It traces the trajectory from the dark ages of physical paper and manual data entry, through the heuristic “bag-of-words” approaches of the early 2000s, to the neural network revolution of domain-specific transformers like FinBERT. Finally, it explores the current frontier: the deployment of reasoning-based Large Language Models (LLMs) capable of executing complex, multi-step financial logic. Throughout this evolution, **a critical finding has emerged: generic linguistic models generally struggle in the financial domain.** The ultimate quantitative edge is, and often has been, derived from custom, domain-specific training. For investment consultants, finance professionals, and sophisticated investment teams now being exposed to modern artificial intelligence tools, understanding this historical perspective is not merely an academic exercise; it is a prerequisite for evaluating which technologies will likely deliver an informational advantage in highly competitive capital markets.

PART I: THE DARK AGES AND THE DAWN OF DIGITIZATION



THE PRE-DIGITAL FOUNDATIONS OF QUANTITATIVE INVESTING

To appreciate the sophistication of modern artificial intelligence in finance, one must first understand the severe limitations of the era that preceded it. The intellectual foundations of quantitative finance stretch back over a century, beginning with Louis Bachelier's 1900 dissertation on the theory of speculation, which first applied mathematical principles to financial markets.¹ As computational power slowly increased, the theoretical gave way to the practical. The 1970s and 1980s saw the emergence of foundational models like the Black-Scholes options pricing model and the proliferation of computerized trading on the New York Stock Exchange.¹

1900

Bachelier's dissertation applies mathematics to financial markets for the first time

1970s–80s

Black-Scholes model and computerized trading on the NYSE emerge

1982

Renaissance Technologies founded by James Simons

1988

D.E. Shaw & Co. founded; 18% annualized returns via statistical arbitrage

By the late 1980s and early 1990s, pioneering quantitative firms like Renaissance Technologies (founded in 1982 by James Simons) and D.E. Shaw & Co. (founded in 1988 by David E. Shaw) began dominating the markets.³ David Shaw famously described finance as a “wonderfully pure information-processing business”.⁴ These firms achieved unprecedented returns — with D.E. Shaw reportedly delivering 18 percent annualized returns in its early years — by assembling teams of computer scientists and physicists to identify microscopic statistical arbitrages.³ Yet, these early strategies were entirely constrained by the data formats available to them. They were parsing price ticks and standardized numerical feeds; the concept of automated, systematic textual analysis was practically science fiction.

THE ERA OF PHYSICAL PAPER AND THE MANUAL DATA BOTTLENECK

During the formative years of quantitative investing, corporate disclosures, annual reports (10-Ks), and earnings transcripts were not available in machine-readable formats. They were physical documents printed on paper and distributed via mail or accessible only in centralized reading rooms. For a quantitative researcher seeking to systematically analyze corporate fundamentals across a universe of thousands of equities, the primary bottleneck was not computational power or algorithmic sophistication, but raw data acquisition.²

The early solution to this bottleneck was brute-force human labor. The financial data industry, most notably Bloomberg L.P., built its initial dominance by solving this exact physical constraint. Founded in 1981, Bloomberg launched its flagship terminal in 1982, rapidly expanding to over 10,000 installed units within a decade.⁷ But the data powering these terminals did not magically appear from corporate mainframes. To provide historical financial statements and balance sheet data, Bloomberg and its competitors employed vast armies of data entry clerks.¹⁰

These teams, primarily located in data centers such as Bloomberg's Princeton facility, were tasked with physically reading corporate financial statements and manually keying the data into proprietary databases, subjecting it to vigorous quality assurance processes.¹⁰ This process was excruciatingly slow, highly expensive, and inherently prone to human error.⁶ A typographical error in a manually entered cash flow statement could drastically skew a quantitative model's output, and the cost of poor data quality in financial systems remains a multi-trillion-dollar issue globally.⁶ Consequently, the earliest systematic text analysis was practically non-existent; quants struggled merely to digitize the numerical tables, leaving the rich, contextual narratives of Management Discussion and Analysis (MD&A) sections and the nuanced tone of executive transcripts entirely untouched and unquantified.

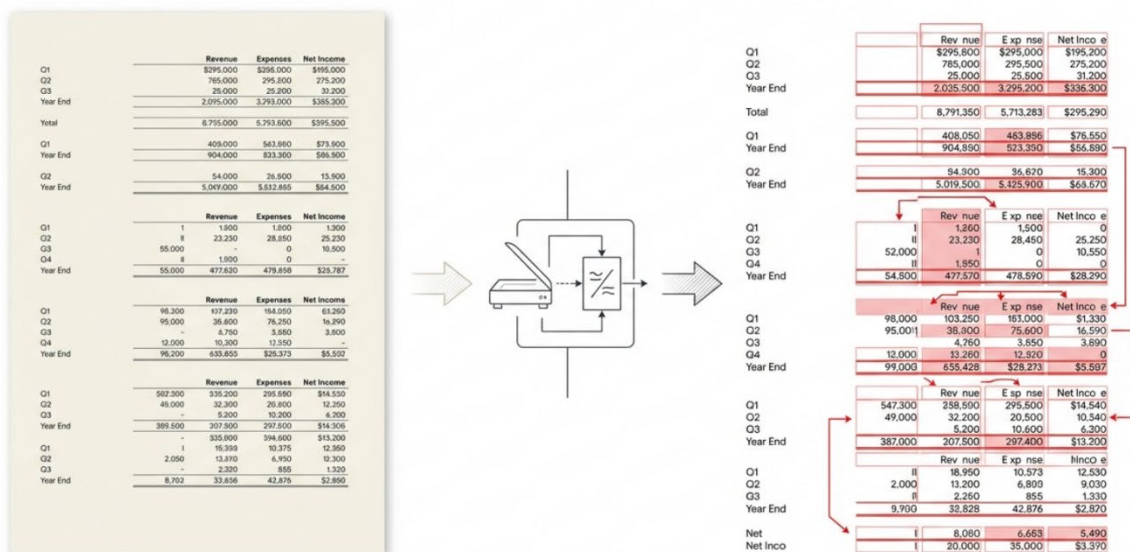
THE PROMISE AND PERIL OF OPTICAL CHARACTER RECOGNITION (OCR)

As computational power expanded in the 1990s and early 2000s, the financial industry turned to Optical Character Recognition (OCR) technology in a desperate bid to automate the digitization process.¹³ The roots of OCR stretch back significantly further than the digital age; Emanuel Goldberg invented a machine utilizing photoelectric cells for pattern recognition as early as the 1910s and 1920s to retrieve microfilmed financial records, and Ray Kurzweil commercialized omni-font OCR in the 1970s.¹³

By the late 1990s, researchers were attempting to leverage commercial OCR software to transcribe the vast archives of historical paper documents.¹⁵ While OCR represented a conceptual leap forward from manual data entry, it introduced a new paradigm of severe technical challenges for quantitative researchers. Financial documents are notoriously complex in their layout. They are not simple, continuous streams of prose; they feature multi-column text, intricate hierarchical tables, dense footnotes, and varying typography.

Early OCR algorithms frequently failed to accurately parse these complex spatial topographies.¹⁸ Techniques designed for standard documents, such as adaptive run-length smoothing and skeleton segmentation paths, routinely broke down when confronted with dense financial tables.¹⁹ The resulting data was often plagued by severe "noise" — typographical errors, punctuation inconsistencies, and the catastrophic breakdown of tabular structures where neat columns of numbers were flattened into meaningless, continuous strings of text.¹⁸

THE INFORMATION BOTTLENECK: THE REALITY OF EARLY OCR DIGITIZATION



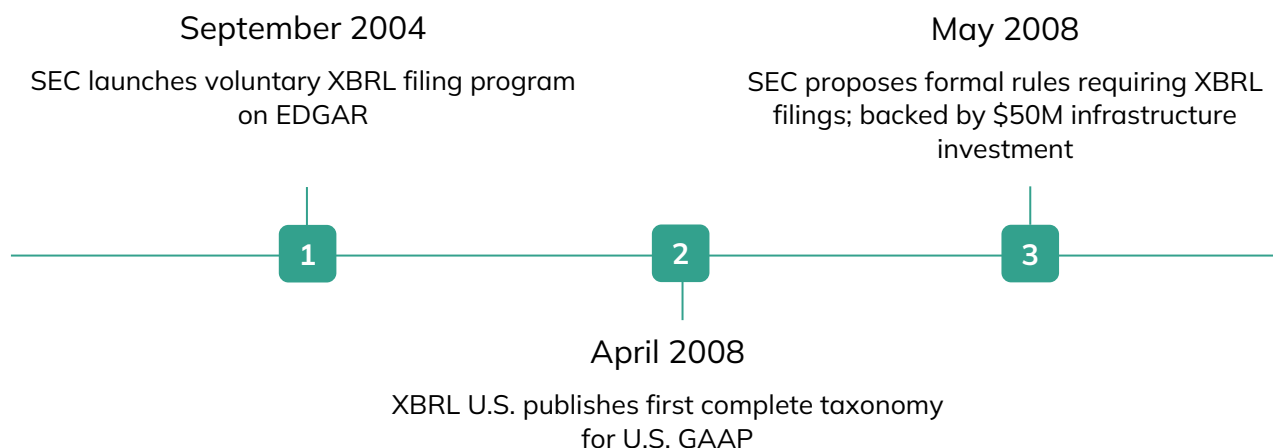
Prior to standardized electronic filing, quantitative researchers relied on Optical Character Recognition (OCR) to parse physical documents. The technology frequently failed to preserve the complex spatial relationships of financial tables, resulting in noisy, unstructured data that hindered systematic analysis.

Quantitative teams found themselves dedicating immense engineering resources to the mundane task of data cleaning rather than alpha generation. Researchers had to write elaborate, highly brittle regular expression (Regex) scripts and heuristic parsing rules to comb through OCR outputs, attempting to reconstruct the original document's meaning.²⁰ The sheer friction of transforming non-machine-readable data into a usable format meant that textual analysis remained an esoteric, resource-intensive niche within quantitative finance.

THE EDGAR REVOLUTION AND STRUCTURED DATA

The paradigm shifted fundamentally in the mid-2000s, driven not by a sudden technological breakthrough in artificial intelligence, but by aggressive regulatory mandate. The U.S. Securities and Exchange Commission (SEC) initiated a sweeping modernization of its Electronic Data Gathering, Analysis, and Retrieval (EDGAR) system.²¹

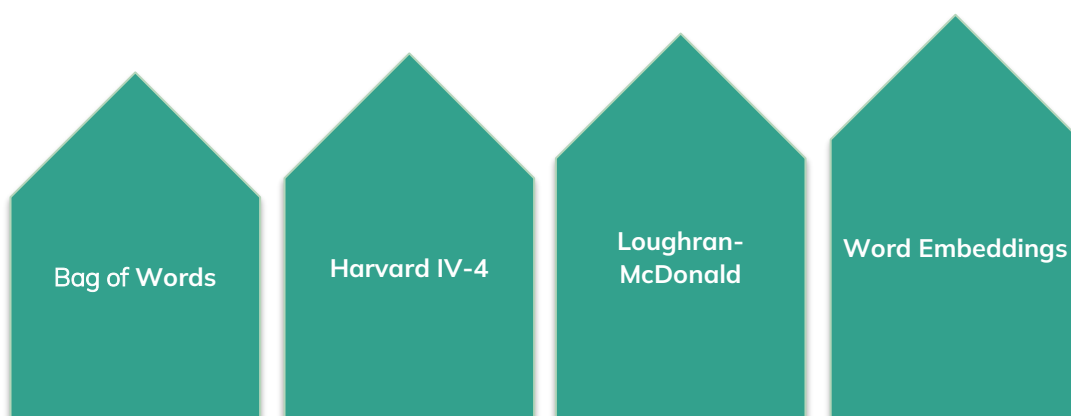
Recognizing that the capital markets required faster, more reliable information flow to maintain efficiency, the SEC began advocating for the use of eXtensible Business Reporting Language (XBRL).²¹ XBRL was specifically developed to provide a standardized, machine-readable format for financial reporting, effectively tagging every individual data point within a financial statement with a specific, universally understood digital identifier.²¹ In September 2004, the SEC launched a voluntary filing program, allowing registrants to supplement their traditional filings with XBRL data.²¹



The momentum accelerated dramatically when XBRL U.S. published the first complete taxonomy for U.S. GAAP in April 2008.²³ Shortly thereafter, in May 2008, the SEC proposed formal rules requiring public companies and mutual funds to file their financial statements in XBRL format.²³ Backed by a \$50 million investment in infrastructure by the SEC, this mandate forced the corporate world into the machine-readable era.²³

This regulatory shift was a profound watershed moment for quantitative finance. Suddenly, financial statements were uniformly structured.²³ A quantitative model could algorithmically query a 10-K filing and instantly extract the exact value for “Net Income” or “Operating Cash Flow” without relying on manual data entry clerks or error-prone OCR systems.²⁴ However, while XBRL brilliantly solved the problem of extracting numerical accounting data, the narrative text within these filings — the qualitative explanations of performance, the risk factors, the forward-looking guidance — remained unstructured. Quants now had access to vast oceans of clean, digital text, but they still lacked the computational tools to genuinely understand it. The race to teach algorithms to “read” had officially begun.

PART II: THE LEXICON ERA AND THE DISCOVERY OF CONTEXT



THE “BAG OF WORDS” AND THE HARVARD IV-4 FALLACY

With vast corpora of clean, machine-readable text now readily available via the EDGAR database, early quantitative researchers adopted rudimentary techniques from the nascent field of computational linguistics. The initial, dominant approach was highly heuristic, relying on the “bag-of-words” model.²⁵ In this mathematical framework, a document is stripped of its grammar, syntax, and word order. It is simply treated as an unstructured collection (a “bag”) of individual tokens.²⁵

To extract sentiment — to determine if a corporate document was broadly optimistic or pessimistic regarding the firm’s prospects — researchers relied on pre-existing psychological and sociological dictionaries. The most prominent and widely adopted of these was the General Inquirer program, specifically utilizing the Harvard-IV-4 psychosocial dictionary, and its “TagNeg” (H4N) file for identifying negative sentiment.²⁵ A quantitative algorithm would simply scan a corporate 10-K, count the frequency of words that appeared on the Harvard negative list, scale that count by the document’s total length (term frequency), and output a static sentiment score.²⁵

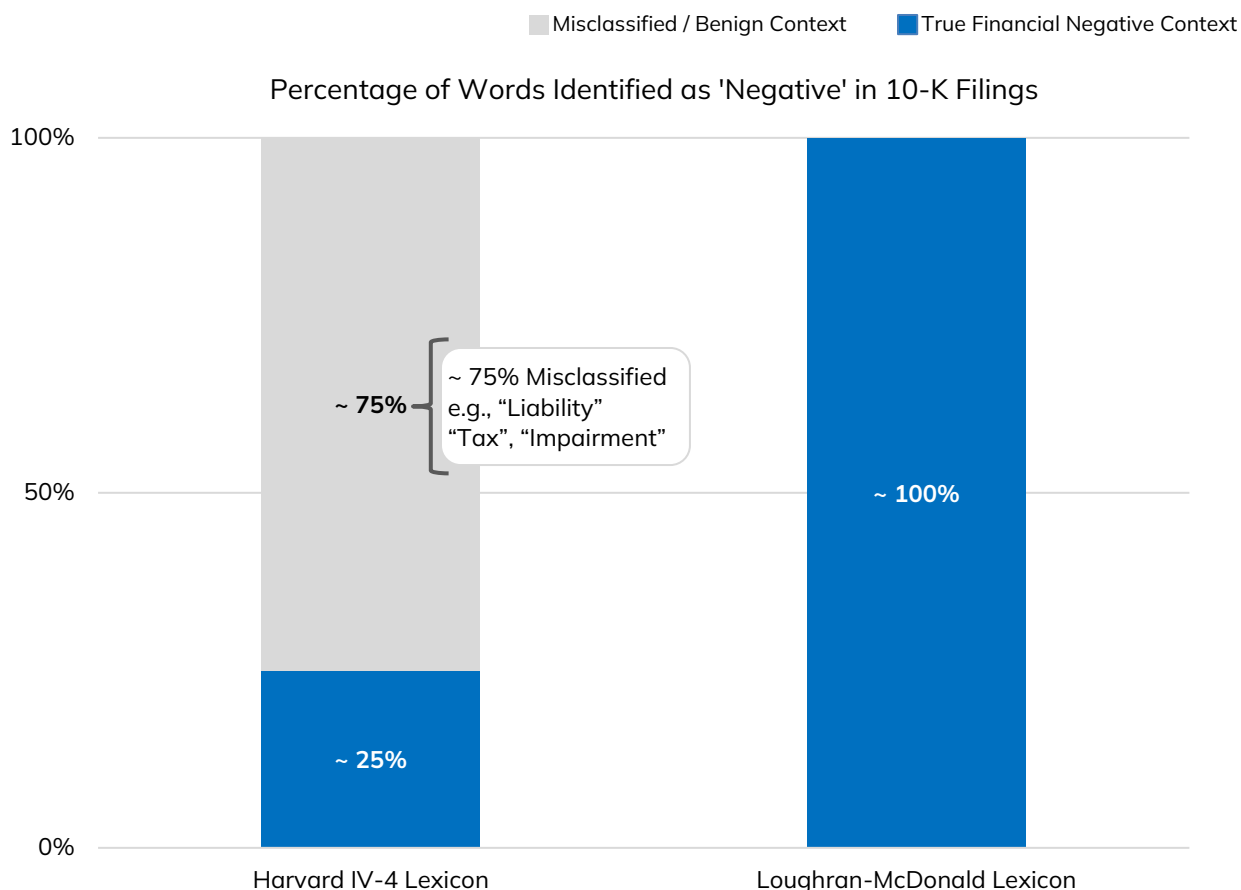
This rudimentary text processing methodology immediately revealed glaring flaws when applied to the capital markets. The Harvard-IV-4 dictionary, containing 2,005 negative words, was designed for general English and psychosocial analysis, not the highly specialized dialect of corporate finance.²⁷ When applied to SEC filings, the dictionary systematically and severely misclassified the tone of the documents.

THE LOUGHRAN-McDONALD BREAKTHROUGH: CUSTOMIZATION AS AN EDGE

The critical inflection point in early financial text analysis occurred in 2011 with the publication of a seminal paper by Tim Loughran and Bill McDonald. Their exhaustive research demonstrated a devastating flaw in the prevailing quantitative methodology: almost three-fourths (nearly 75%) of the words identified as negative by the widely used Harvard Dictionary were not actually negative when deployed in a financial context.²⁷

Consider the standard vocabulary of a corporate financial report. Words like “tax,” “liability,” “board,” “mine,” “requirements,” “impairment,” and “cancer” are frequently used.²⁷ To a general psychological dictionary, “liability” implies a personal burden, “mine” might imply an explosive device, and “cancer” is inherently and universally negative.²⁷ However, in a 10-K filing, a “liability” is simply the standard accounting term for the right side of a balance sheet. A “mine” is a core operational asset for a commodities or extraction firm. “Cancer” is a target disease for a pharmaceutical company’s research pipeline. A company plainly stating, “We paid our tax liability,” or an oncology firm discussing its “cancer trials” would be erroneously flagged as overwhelmingly negative by the Harvard dictionary.²⁷

THE DANGER OF GENERIC LEXICONS IN FINANCIAL TEXT ANALYSIS



Analysis of 10-K filings from 1994 to 2008 revealed that generic psychosocial dictionaries (Harvard IV-4) misclassified nearly 75% of words identified as 'negative' because they failed to account for financial context (e.g., words like 'liability' or 'tax').

Data sources: arXiv (S12)), Annual Reviews (S42), Journal of Finance (Loughran & McDonald) (S43)

CONSTRUCTION AND IMPACT

To solve this, Loughran and McDonald constructed a custom, domain-specific financial sentiment dictionary.²⁷ They manually parsed through the full EDGAR 10-K archive to classify words based specifically on their implications for corporate performance and stock returns.²⁷ Their refined lexicon included 2,345 strictly negative terms tailored for finance (along with 347 positive terms).²⁹ They also expanded their analysis to include vital categories beyond mere positive and negative sentiment, introducing specific classifications for "uncertainty," "litigiousness," "weak modal," "strong modal," and "constraining" language.³²

The introduction of the Loughran-McDonald (LM) lexicon was a revelation for the quantitative community. When backtested against historical market data, the LM dictionary significantly outperformed the Harvard IV-4 dictionary in predicting subsequent financial outcomes. It demonstrated statistically significant correlations with 10-K filing returns, trading volume, return volatility, unexpected earnings, and even instances of material weakness and fraud.²⁷

This era established a core principle for quantitative researchers utilizing Natural Language Processing (NLP): **generic models are often limited when applied to specialized domains**. The investment edge does not come merely from accessing the text, but from customizing the analytical tool to understand the unique structural and semantic rules of the financial ecosystem. This principle would echo significantly in the decades to follow as the industry transitioned from simple word counts to complex neural networks.

BEYOND THE BAG: SYNTAX, SEMANTICS, AND WORD EMBEDDINGS

While the Loughran-McDonald dictionary provided a massive leap in analytical accuracy, the underlying “bag-of-words” methodology remained fundamentally limited. By stripping text of its structural order, the methodology was completely blind to context, syntax, and negation.²⁵ For instance, a dictionary-based algorithm would parse the sentence “The company’s core profits did not increase” and potentially flag it as positive due to the presence of the word “profits” and “increase,” entirely missing the negating phrase “did not”.²⁶

Furthermore, simple term-frequency counts failed to capture semantic relationships. A dictionary model does not inherently understand that “revenue,” “sales,” and “top-line growth” are conceptually related unless they are manually hardcoded into the same exact category.

The first major paradigm shift away from rigid lexicons occurred with the advent of Neural Language Models (NLMs) and word embeddings, most notably algorithms like Word2Vec and GloVe introduced around 2013.³³ Instead of relying on human-curated lists, Word2Vec utilized shallow neural networks to process vast amounts of text and map words into dense, multi-dimensional vector spaces.³³

In this mathematical vector space, words that appeared in similar contexts were positioned closer together. This allowed quantitative models to capture complex semantic relationships algebraically. The famous linguistic equation “king - man + woman = queen” could now be replicated in finance: “equity - stock + debt = bond.” To power these models, researchers curated massive domain-specific datasets, such as the EDGAR-CORPUS, which comprised annual reports from all publicly traded US companies spanning over 25 years and containing billions of tokens.²⁴ This spatial representation of language allowed quants to build predictive models that understood synonyms, nuances, and thematic clusters without explicit, manual programming.²⁴

THE TRANSFORMER EPOCH AND THE FINBERT TRIUMPH

The true revolution in financial NLP, however, arrived with the development of the Transformer architecture, and specifically the release of BERT (Bidirectional Encoder Representations from Transformers) by Google.

Prior to Transformers, models processed text sequentially (using Recurrent Neural Networks or LSTMs). They struggled with “long-term dependencies” — remembering a critical modifying clause at the beginning of a long paragraph when analyzing a noun at the end. Financial documents, such as annual reports and prospectuses, are notoriously lengthy and structurally convoluted, making sequential processing highly inefficient and prone to information loss.¹⁸

The Transformer architecture introduced the “self-attention” mechanism.³⁵ Instead of reading a sentence linearly from left to right, a self-attention layer looks at every word in a sequence simultaneously and mathematically weighs the importance (the “attention”) each word should pay to every other word.³⁵ This allowed the model to maintain deep, bidirectional context. If a 10-K stated, “Despite the severe supply chain disruptions experienced in the third quarter, our domestic profit margins expanded,” the attention mechanism could seamlessly link the expansion of margins specifically to the adversarial context of the disruptions, generating a highly nuanced interpretation of management’s operational resilience.

When generic BERT was released, quantitative researchers eagerly applied it to financial sentiment analysis. However, history quickly repeated itself. Just as the generic Harvard IV-4 dictionary failed because it lacked financial context, the generic BERT model — which was pre-trained on a massive corpus of standard English text, primarily Wikipedia and BookCorpus — struggled to interpret the nuances of corporate finance.³⁶ To a neural network trained on Wikipedia, the phrase “the stock was heavily shorted” might not register as a critical indicator of market sentiment, as the grammatical structure and vocabulary differ wildly from standard encyclopedic text. The solution to this deficiency was the creation of **FinBERT**.³⁶

HOW FINBERT WAS BUILT

FinBERT’s architecture mirrors the standard BERT model, but its foundational knowledge base is radically different. To create FinBERT, researchers subjected the architecture to a rigorous pre-training phase utilizing a massive, specialized financial corpus. This corpus included billions of tokens scraped directly from corporate 10-Ks, 10-Qs, earnings call transcripts, and professional analyst reports.³⁶

Specialized Pre-Training Corpus

Billions of tokens from corporate 10-Ks, 10-Qs, earnings call transcripts, and professional analyst reports

Masked Language Modeling (MLM)

Random words in financial text are hidden; the model must predict them based on context, absorbing the dialect of Wall Street.

Next Sentence Prediction (NSP)

Model learns the structural flow of financial discourse from CFOs, central bankers, and equity analysts

Fine-Tuning on Annotated Data

10,000 manually annotated financial sentences labeled for positive, negative or neutral sentiment

During this pre-training process, FinBERT engaged in Masked Language Modeling (MLM) — where random words in a financial text are hidden, and the model must mathematically predict them based on context — and Next Sentence Prediction (NSP).³⁷ By forcing the neural network to repeatedly predict the linguistic patterns of chief financial officers, central bankers, and equity analysts, FinBERT fundamentally absorbed the dialect of Wall Street.³⁶ Following pre-training, the model underwent fine-tuning using 10,000 manually annotated financial sentences labeled for positive, negative, or neutral sentiment.³⁶

THE EVIDENCE FOR DOMAIN-SPECIFIC TRAINING

The empirical results were striking. In head-to-head benchmarking for financial sentiment classification, FinBERT consistently outperformed both the legacy Loughran-McDonald dictionary approach and the generic, vanilla BERT model.³⁸

TABLE 1: EVOLUTION OF FINANCIAL SENTIMENT ACCURACY

Model Architecture	Approach Type	Primary Training Corpus	Contextual Awareness	Financial Sentiment F1 Score
Harvard IV-4	Lexicon / Dictionary	General Psychology	None (Bag of Words)	Poor (Systematic Misclassification)
Loughran-McDonald	Lexicon / Dictionary	Financial (10-Ks)	None (Bag of Words)	Baseline standard for Finance
Generic BERT	Transformer (Deep Learning)	Wikipedia / Books	High (Bidirectional)	~0.84
FinBERT	Transformer (Deep Learning)	SEC Filings / Analyst Reports	High (Bidirectional)	~0.85 - 0.86+

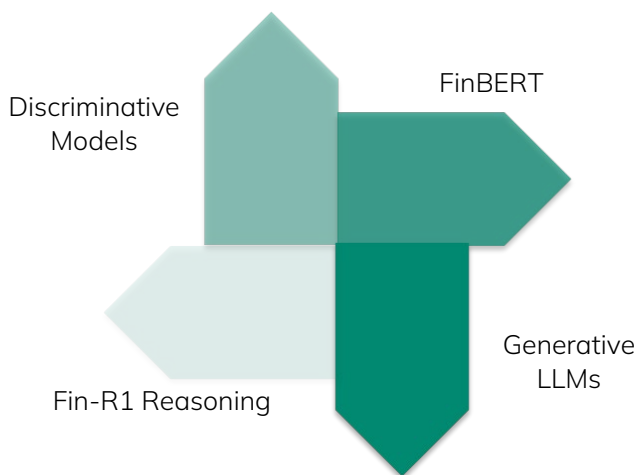
Data represents generalized performance metrics from comparative studies on financial text classification tasks, illustrating the superiority of domain-specific training.

SECTOR-LEVEL NUANCE AND THE LIMITS OF CLASSIFICATION

The superiority of FinBERT solidified a critical paradigm in modern quantitative finance: the deep integration of advanced NLP techniques and domain-specific Machine Learning models is an indispensable tool for extracting alpha.³⁶ **However, studies revealed that the efficacy of these models varied significantly across different market sectors.** For instance, advanced sentiment analysis yielded dramatic improvements in forecast accuracy within the Energy (+14.31%) and Aviation (+3.81%) sectors, while showing less impact or even deterioration in heavily structured, quantitative sectors like Banking (-1.50%) and Steel (-10.03%).⁴¹ This highlighted that textual sentiment is not a blunt instrument, but a highly contextual signal that must be calibrated to specific market microstructures.⁴¹

By the early 2020s, the use of custom-trained, transformer-based language models had become table stakes for sophisticated hedge funds. The industry had successfully moved from merely counting domain-specific words to deeply understanding domain-specific context. Yet, the capability of these models was still largely confined to classification tasks — labeling a paragraph as bullish or bearish. The next evolutionary leap would require models not just to classify text, but to reason through it.

PART III: THE LLM FRONTIER AND DEEP FINANCIAL REASONING



THE GENERATIVE SHIFT AND THE ALLURE OF ZERO-SHOT ANALYSIS

The public release of highly scaled, generative Large Language Models (LLMs) — most notably the GPT-4 architecture from OpenAI and subsequent models from Meta and Google — initiated a tectonic shift across all knowledge industries, and quantitative finance was no exception.⁴² While previous models like FinBERT were discriminative (specifically designed to classify input data into predefined categories), LLMs were generative. They were capable of synthesizing vast amounts of information, generating executive summaries, and purportedly executing complex analytical workflows natively.⁴³ The prospect of an artificial intelligence agent reading an entire S-1 prospectus and outputting a comprehensive investment thesis seemingly overnight captivated the industry.

The initial hype surrounding the application of generic LLMs to finance reached a crescendo with a highly publicized working paper originating from researchers at the University of Chicago Booth School of Business.⁴⁴ The researchers — Alex Kim, Maximilian Muhn, and Valeri Nikolaev — sought to determine if a general-purpose LLM (specifically GPT-4 Turbo) could perform rigorous financial statement analysis and predict future earnings changes with the same acumen as professional human analysts.⁴⁴

THE CHICAGO CONTROVERSY: THE DANGER OF DATA CONTAMINATION

To conduct their study, the researchers collected and anonymized financial statements from over 15,000 corporations spanning from 1968 to 2021.⁴⁴ They stripped away company names and specific years, replacing them with generic labels (e.g., Year t , Year $t - 1$), and standardized the format to match Compustat’s balancing model.⁴⁸ They then fed this stripped data into the LLM.

The results initially reported were staggering. Utilizing a prompting strategy that instructed the model to analyze the statements step-by-step, the LLM purportedly achieved a 60% accuracy rate in predicting the direction of future earnings.⁴⁴ This significantly outperformed the consensus of human financial analysts, who averaged roughly 53% accuracy on the same cohort, and bested naive models which hovered at 49%.⁴⁴ Furthermore, the researchers claimed that trading strategies built upon the LLM’s predictions yielded higher Sharpe ratios and significant alpha compared to traditional machine learning baselines.⁴⁴

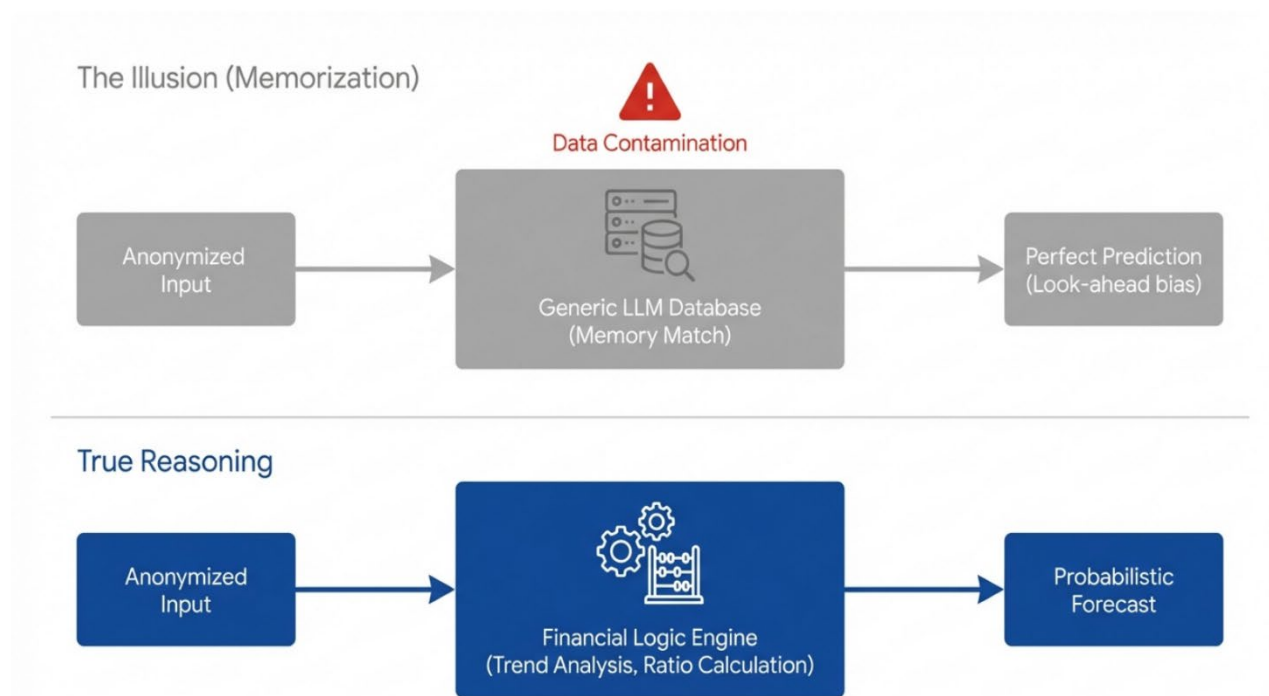
ACCURACY

49%	53%	60%
Naive Models	Human Analysts	LLM Accuracy
Average of naive baseline models on the same task	Average accuracy of professional human financial analysts on the same cohort	GPT-4 Turbo’s purported accuracy in predicting direction of future earnings
15K+ Corporations Financials 1968 – 2021		

The narrative that a generic AI could instantly outperform highly trained equity analysts without any industry-specific narrative context sent shockwaves through the investment community.⁴⁸ However, the triumph was short-lived, and **the paper soon became the center of significant academic and industry controversy.**⁴⁶

Upon closer inspection and attempts at replication, severe methodological flaws were identified by the quantitative community and the authors themselves. The primary issue was **massive data contamination and look-ahead bias**.⁴⁹ Despite the researchers' attempts to "anonymize" the financial statements by removing text identifiers, the generic LLM — having been trained on vast swaths of the internet, including historical financial databases, SEC filings, and Wikipedia — was essentially able to memorize and recognize the exact numerical fingerprints of the companies.⁴⁴ The model was **not necessarily "reasoning"** about the future fundamentals of an unknown company; it was **simply recalling the future** from its extensive training memory. Recognizing these inconsistencies and the failure to replicate the findings cleanly out-of-sample, the authors withdrew the working paper from circulation.⁴⁷

THE ILLUSION OF ALPHA: DATA CONTAMINATION IN GENERIC LLMs



When generic LLMs analyze historical financial statements, high predictive accuracy often stems from data contamination rather than true analytical capability. The model recognizes the 'anonymized' numerical patterns from its vast training data, retrieving the known future outcome instead of executing fundamental financial reasoning.

This episode forcefully reiterated the central lesson of the Loughran-McDonald era: **generic tools are fraught with peril in finance**.⁵¹ Generic LLMs hallucinate, they struggle with strict numerical reasoning, they fail to grasp the rigid logic of accounting identities, and their broad training data makes them entirely unsuitable for backtesting due to data leakage.⁴² When evaluating complex financial documents, standard models often misread charts, break tabular cell mappings, and fail to understand the dependencies and mathematical formulas embedded within unstructured data like spreadsheets.⁴²

ENGINEERING SYSTEM 2 THINKING: CHAIN OF THOUGHT

Despite the setbacks regarding generic, zero-shot prediction, the core underlying technology of LLMs presented undeniable utility, provided their output could be strictly controlled. To harness this power, quantitative researchers began adapting advanced prompting architectures, most notably Chain-of-Thought (CoT) reasoning.⁵⁴

Traditional machine learning and generic LLM interactions operate predominantly on “System 1” thinking — fast, intuitive, pattern-matching responses.⁵⁶ If asked a complex financial valuation question, a standard model attempts to generate the answer in a single, continuous forward pass, frequently resulting in catastrophic mathematical or logical errors because it is merely guessing the most probable next word rather than executing an equation.

Chain-of-Thought forces the model into “System 2” thinking — slow, deliberate, step-by-step logical progression.⁵⁴ Instead of passively generating a summary, a CoT prompt explicitly forces the LLM to break down a complex problem into intermediate steps that mimic the rigorous analytical workflow of a Chartered Financial Analyst (CFA).⁴⁴

For example, when parsing a complex narrative disclosure about a corporate acquisition, a CoT framework mandates that the model first explicitly identify the target entities, then extract the purchase price, then identify the financing mechanism (debt vs. equity ratios), then calculate the implied valuation multiples, and only then generate an assessment of the transaction’s impact.⁵⁴ Researchers have developed sophisticated zero-shot prompts, such as FinCoT, that inject expert financial workflows — encoded as structural blueprints — directly into the model’s context window.⁵⁵ By enforcing this explicit intermediate structure, researchers can significantly reduce hallucinations, verify the model’s adherence to financial formulas (verifying units, basis points, and boundary conditions), and create an auditable trail of reasoning.⁵⁵

THE ULTIMATE EDGE: DOMAIN-SPECIFIC REASONING MODELS (FIN-R1)

While CoT prompting improves the performance of generic models, it is essentially a software patch applied to a fundamentally generalized engine. The true frontier of quantitative NLP — the modern equivalent of the FinBERT leap — is the development of explicitly trained, domain-specific financial reasoning models. The current paradigm shift involves moving away from models that simply predict the next word, to models that are trained via Reinforcement Learning (RL) to search for correct logical pathways.⁵⁷

This shift is perfectly exemplified by the emergence of specialized reasoning architectures like Fin-R1. Developed by researchers at the Shanghai University of Finance and Economics, Fin-R1 was designed specifically to tackle the fragmented data, weak business generalization, and uncontrollable reasoning logic endemic to the financial sector.⁵⁸

Fin-R1 represents a masterclass in financial AI architecture. Despite possessing a relatively compact parameter scale (built upon the Qwen2.5-7B-Instruct architecture, meaning it operates with 7 billion parameters, compared to generic behemoths like GPT-4 which exceed a trillion parameters), it achieves state-of-the-art performance on complex financial tasks.⁵⁸ It accomplishes this by utilizing a sophisticated two-stage training pipeline.⁵⁷

Stage 1: Supervised Fine-Tuning (SFT) on High-Quality CoT Data The foundation of a reasoning model is its data. To cure the model of generic, generalized responses, researchers constructed an expansive, proprietary dataset (Fin-R1-Data) containing over 60,000 highly curated, verified Chain-of-Thought entries.⁵⁸ This dataset encompasses multidimensional financial expertise in both Chinese and English, including core banking, securities scenarios, and quantitative code generation.⁵⁸ Crucially, it incorporates data like the Financial Postgraduate Entrance Exam (FinPEE) dataset, consisting of rigorous calculation problems.⁵⁸ By fine-tuning the base model strictly on expert-level logical pathways, the model learns the structural grammar and mathematical rigor of financial analysis.⁵⁷

Stage 2: Reinforcement Learning via GRPO The critical differentiation occurs in the second stage. Fin-R1 employs Reinforcement Learning, specifically utilizing an algorithm known as Group Relative Policy Optimization (GRPO).⁵⁸ In standard language modeling, the model is rewarded simply for producing text that sounds human and plausible. In financial reinforcement learning, the model is explicitly penalized for logical inconsistencies, mathematical calculation errors, or non-compliance with accounting rules.

Through GRPO, the model learns to independently construct deep, internal reasoning chains before emitting a final answer. It evaluates multiple potential logical paths to solve a financial problem and updates its internal policy to favor the path that adheres strictly to verifiable financial truth.⁵⁸ This dual reward mechanism enhances both the formatting accuracy and the profound correctness of the content.⁵⁸

The performance metrics of this architecture provide strong evidence for the benefits of custom training. When evaluated on rigorous industry benchmarks like FinQA (numerical reasoning over complex financial reports) and ConvFinQA (multi-turn conversational financial reasoning), the 7B parameter Fin-R1 model achieves significant performance advantages.⁵⁷

TABLE 2: PERFORMANCE EVALUATION ON COMPLEX FINANCIAL REASONING BENCHMARKS

Model	Parameters	FinQA Score	ConvFinQA Score	Average Score
Qwen2.5-7B-Instruct	7B (Base)	60.0	66.0	65.6
DeepSeek-R1-Distill-Qwen	7B	55.0	62.0	58.0
Fin-R1	7B (Custom Reasoning)	76.0	85.0	75.2

The specialized Fin-R1 model, trained explicitly with Reinforcement Learning for financial logic, achieves state-of-the-art results, dramatically outperforming both its base instruction model and generic distilled reasoning models of the same size.⁵⁷

The ablation studies surrounding these models reveal a fascinating truth regarding the deployment of AI in finance: applying Reinforcement Learning to a base model yields only modest gains.⁵⁷ The massive leap in capability — the true alpha — is only unlocked when Reinforcement Learning is applied *on top of* the domain-specific Supervised Fine-Tuning. The model must first be explicitly taught the vocabulary and structure of finance, and only then can it be successfully taught to reason within that structure.⁵⁷

These reasoning models represent **a profound paradigm shift**. They are not merely reading text and extracting positive or negative sentiment; they are interpreting complex multi-tiered tables within 10-Ks, tracking the spatial relationships of hierarchical numerical data across XBRL tags, evaluating the regulatory compliance of unstructured text, and generating rigorous, auditable financial code.⁴² They solve the very problems of complex document layout and spreadsheet logic preservation that have plagued quantitative researchers since the earliest days of Optical Character Recognition.¹⁹

CUSTOM ALPHA: WHEN DOMAIN TRAINING PAYS OFF

To validate the theoretical trajectory of financial NLP, a rigorous empirical simulation benchmarked three sentiment methodologies against a custom-tuned LLM architecture using a uniform corpus of corporate earnings call transcripts across the US Equity Universe from

December 31, 2014, to December 31, 2024. All signals were sector-demeaned, subject to liquidity filters ($ADVT \geq \$250K$, $MKTCAP \geq \$100M$), and lagged by 7 days. Results reflect long-short quintile portfolios with equal weighting.

The three benchmarks were: Loughran-McDonald (LM), a financial sentiment dictionary; FinBERT, a BERT-based model pre-trained on financial data; and ChatGPT, a generic LLM prompted to classify sentiment statement-by-statement. All three share a critical structural limitation: they analyze individual statements in isolation, discarding the narrative arc and conversational dynamics of the full transcript. The custom model was designed to overcome this by processing each transcript end-to-end with full global context, using the multi-step reasoning and domain-specific reinforcement learning frameworks discussed in the prior section.

The distributional properties of the benchmark signals reveal an immediate structural weakness. Both generic models classify over 70% of earnings call statements as neutral. More tellingly, ChatGPT’s negative classification rate (2.90%) is less than half that of FinBERT (7.13%) — a systematic optimism bias that materially impairs short-side signal generation in a long-short portfolio.

TABLE 3: DISTRIBUTION OF SENTIMENT CLASSIFICATIONS

Sentiment	FinBERT	ChatGPT
Neutral	73.90%	71.36%
Positive	18.97%	25.74%
Negative	7.13%	2.90%

Despite being applied to identical transcripts, pairwise signal correlations are moderate. FinBERT and ChatGPT are most similar at 40%. The custom model shows the lowest correlation to all three benchmarks (28%–34%), confirming it extracts fundamentally different information from the same underlying documents.

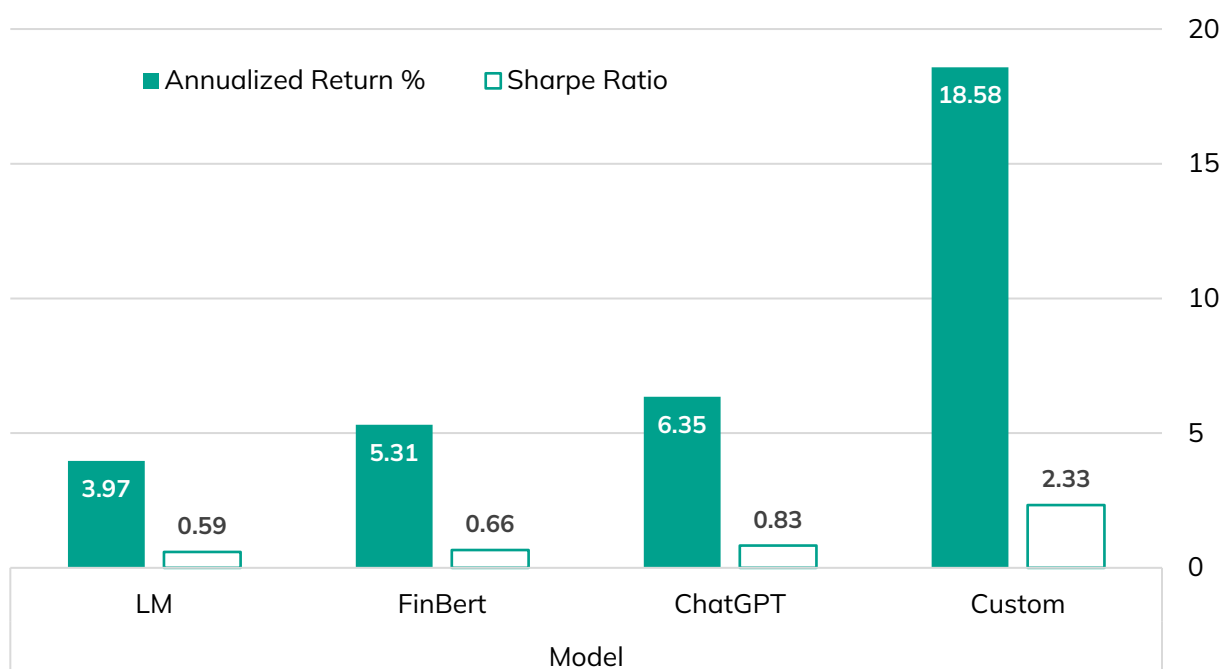
TABLE 4: CORRELATION MATRIX OF INVESTMENT SIGNALS

Correlation	LM	FinBERT	ChatGPT	Custom
LM	1.00	38%	37%	34%
FinBERT	38%	1.00	40%	30%
ChatGPT	37%	40%	1.00	28%
Custom	34%	30%	28%	1.00

PERFORMANCE RESULTS: A CLEAR STAIRCASE

The performance results reflect a clear staircase: each step up in model sophistication — from dictionary to domain-trained transformer to generic LLM — delivers progressively higher risk-adjusted returns, with Sharpe Ratios rising from 0.59 (LM) to 0.66 (FinBERT) to 0.83 (ChatGPT). The custom model breaks from this progression entirely, delivering an 18.58% annualized return and a Sharpe of 2.33 — a 3–5x improvement over the next-best alternative. Notably, this outperformance held across divergent market regimes, including years where the benchmark models each generated negative returns.

TABLE 5: SUMMARY PERFORMANCE — LONG-SHORT QUINTILE PORTFOLIOS, EQUAL WEIGHTED, US_MC UNIVERSE, 2015–2024



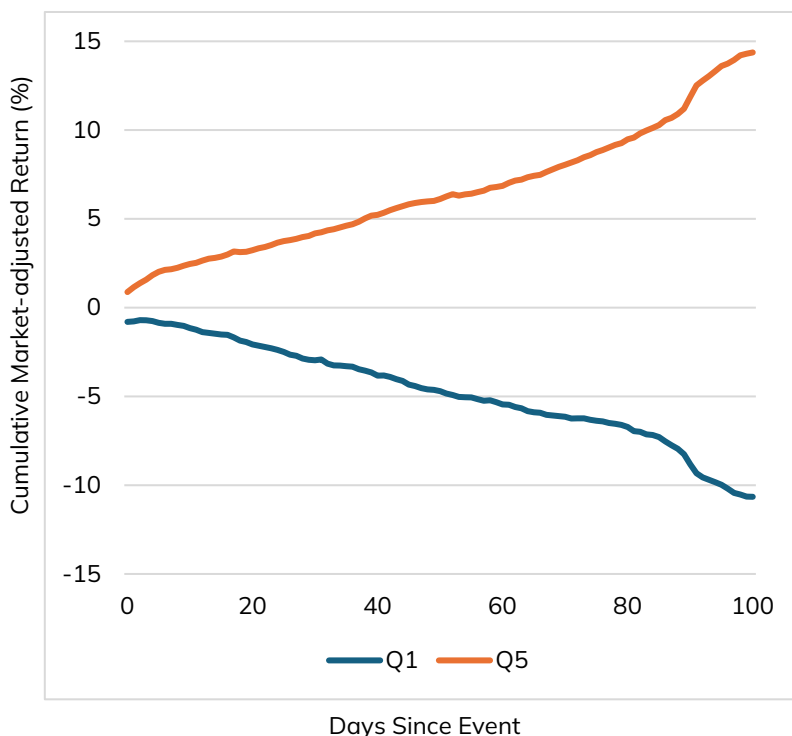
	LM	FinBERT	ChatGPT	Custom
Annualized Return	3.97%	5.31%	6.35%	18.58%
Sharpe Ratio	0.59	0.66	0.83	2.33

THE MECHANISM OF PERSISTENCE

Risk factor analysis confirms **the return is not disguised beta**. All four models show minimal or negative exposure to value, dividend yield, and volatility. The custom model's most notable differentiation is a higher momentum loading (0.28 vs. 0.16–0.17 for the benchmarks), consistent with a signal that identifies earnings quality momentum the market is slow to price, and a higher size exposure (0.20 vs. 0.09–0.13), indicating a practical tilt toward larger, more liquid names suited to institutional deployment.

Event-study analysis reveals **the mechanism behind this persistence: a slow decay driven by subtle signals**. Unlike raw sentiment — typically priced in on the announcement date — the custom model's alpha accrues gradually over the full 100-day post-earnings window. Management tone, Q&A dynamics, and qualitative disclosures slowly diffuse into markets through analyst revisions and guidance updates converting a statistical finding into an implementable systematic edge.

CUMULATIVE PERFORMANCE OF SENTIMENT MODELS (2014–2024 LONG-SHORT QUINTILE PORTFOLIO)



SYNTHESIS AND CONCLUSION

The journey of textual data in quantitative finance is a testament to the industry's relentless pursuit of informational advantage. The evolution from the manual data entry clerks of the 1980s, hand-keying balance sheets into early Bloomberg terminals, to autonomous, reasoning-based AI agents today represents a fundamental transformation in how market intelligence is synthesized and weaponized.

Early quantitative strategies viewed **unstructured text as an impenetrable wall of noise**, relying on rudimentary OCR that stripped financial documents of their vital structural integrity.¹³ The mandate of XBRL by the SEC provided the necessary numerical structure, democratizing access to hard data but leaving the rich narrative context unmined and unresolved.²²

As quants sought to extract alpha from this narrative, the limitations of generic approaches became starkly, and sometimes painfully, apparent. The failure of the Harvard IV-4 dictionary proved that words like “liability” and “impairment” carry unique, domain-specific weight in the capital markets, leading to the creation of the Loughran-McDonald lexicon — the first definitive, empirical proof that customization provides an analytical edge over generic heuristic tools.²⁷ This principle held true through the deep learning revolution; generic BERT models failed to grasp the semantic nuances of SEC filings, necessitating the computationally expensive pre-training of FinBERT on specialized financial corpora to achieve acceptable classification accuracy.³⁶

Today, the industry stands at the precipice of a new era defined by complex reasoning capabilities. The initial enthusiasm for utilizing generic Large Language Models to forecast stock movements was tempered by the harsh realities of data contamination, look-ahead bias, and the models' inability to perform rigorous, multi-step financial logic without hallucinating.⁴⁹

To remain competitive, sophisticated investment teams are turning away from monolithic generic AI and toward highly specialized architectures. By combining Vision-Language Models (VLMs) to parse non-searchable PDFs and complex tabular structures⁴², utilizing Chain-of-Thought prompting to enforce analytical rigor⁵⁵, and deploying Reinforcement Learning (GRPO) to mandate strict adherence to financial axioms, tools like Fin-R1 are bridging the gap between basic natural language processing and genuine, systematic financial acumen.⁵⁸

For investment consultants, fund allocators, and finance professionals navigating this rapidly shifting landscape, the historical perspective is vital. The application of artificial intelligence in finance is not a monolithic plug-and-play solution. The market is semi-strong efficient; any signal that can be easily extracted by a generic, off-the-shelf model is often rapidly commoditized and priced in by the broader market. Alpha is typically generated at the absolute margins of complexity.

The enduring lesson of the past three decades of quantitative research is significant: technological advantages in finance often belong to those who painstakingly tailor their algorithms and models to the unique structural, semantic, and logical realities of the capital markets. Custom training is a key mechanism by which an edge can be created and sustained.

WORKS CITED

1. A history of quant - Hermes Investment, May 2025, <https://www.hermes-investment.com/uploads/2025/06/dfc8397a2ca3f72013414f3fc4cda12c/fhl-a-history-of-quant-05-2025.pdf>
2. A history of quant | Federated Hermes Limited, June 18, 2025, <https://www.hermes-investment.com/uk/en/institutions/insights/macro/a-history-of-quant/>
3. D. E. Shaw & Co: Inside the Quiet Giant of Quant Finance - Quartr, September 8, 2025, <https://quartr.com/insights/company-research/de-shaw-and-co-inside-the-quiet-giant-of-quant-finance>
4. Cyber profits add up to secretive success story - Center for Public Integrity, August 15, 2000, <https://publicintegrity.org/politics/cyber-profits-add-up-to-secretive-success-story/>
5. Renaissance Technologies - Wikipedia, https://en.wikipedia.org/wiki/Renaissance_Technologies
6. From Manual to Automated Data Entry: A Guide for Finance Teams - Copia Wealth Studios, August 7, 2025, <https://copiawealthstudios.com/blog/from-manual-to-automated-data-entry-a-guide-for-finance-teams>
7. Bloomberg L.P. - Encyclopedia.com, <https://www.encyclopedia.com/books/politics-and-business-magazines/bloomberg-lp>
8. How The Bloomberg Terminal Made History — And Stays Ever Relevant - 3BL Media, March 24, 2017, <https://www.3blmedia.com/news/how-bloomberg-terminal-made-history-and-stays-ever-relevant>
9. Bloomberg: Overview and History of the Financial News Company - YouTube, <https://www.youtube.com/watch?v=XASwAPAQvr0>
10. Where does Bloomberg source its data from? - Quora, <https://www.quora.com/Where-does-Bloomberg-source-its-data-from>
11. How far back does Bloomberg data go? - Quora, <https://www.quora.com/How-far-back-does-Bloomberg-data-go>
12. OCR in Finance: Turning Receipts, Invoices, and Bank Statements into Actionable Data, <https://medium.com/intelligent-document-insights/ocr-in-finance-and-accounting-3ed470292224>
13. A brief history of Optical Character Recognition (OCR) - Pitney Bowes, <https://www.pitneybowes.com/content/dam/pitneybowes/uk/en/business-automation/dip-optical-character-recognition-blog.pdf>
14. OCR and People: Your Dynamic Data-Entry Duo - CloudFactory, <https://www.cloudfactory.com/blog/ocr-and-people-your-dynamic-data-entry-duo>
15. (PDF) Open source optical character recognition for historical research - ResearchGate, https://www.researchgate.net/publication/242337212_Open_source_optical_character_recognition_for_historical_research
16. A Quantitative/Qualitative Approach to OCR Error Detection and Correction in Old Newspapers for Corpus-assisted Discourse Studie - CEUR-WS.org, <https://ceur-ws.org/Vol-2816/paper5.pdf>
17. History of Quantitative Text Analysis - Summer Institute in Computational Social Science, https://sicss.io/2019/materials/day3-text-analysis/history-text-analysis/Introduction_to_Text_as_Data.html
18. Top 7 Challenges in Financial Text Preprocessing - Phoenix Strategy Group, <https://www.phoenixstrategy.group/blog/top-7-challenges-in-financial-text-preprocessing>
19. Document Parsing Unveiled: Techniques, Challenges, and Prospects for Structured Information Extraction - arXiv.org, October 28, 2024, <https://arxiv.org/html/2410.21169v1>
20. Quants who parse SEC filings — where are the biggest bottlenecks? - Reddit, https://www.reddit.com/r/quant/comments/1jpsdeg/quants_who_parse_sec_filings_where_are_the/
21. The SEC's XBRL Voluntary Filing Program on EDGAR: A Case for Quality Assurance, <https://scholarship.libraries.rutgers.edu/esploro/outputs/journalArticle/The-SECs-XBRL-Voluntary-Filing-Program/991031550015104646>

22. XBRL Voluntary Financial Reporting Program on the EDGAR System - SEC.gov, February 3, 2005, <https://www.sec.gov/rules-regulations/2005/02/xbrl-voluntary-financial-reporting-program-edgar-system>
23. 10 years of SEC Compliance Filing - A retrospective Glance, November 25, 2019, <https://irisarbon.com/10-years-of-sec-compliance-filing-a-retrospective-glance/>
24. [2109.14394] EDGAR-CORPUS: Billions of Tokens Make The World Go Round - arXiv, September 29, 2021, <https://arxiv.org/abs/2109.14394>
25. Textual sentiment in finance: A survey of methods and models - SSRN, https://papers.ssrn.com/sol3/Delivery.cfm/SSRN_ID2759541_code1450135.pdf?abstractid=2213801
26. Understanding Sentiment Through Context, <https://www.iimb.ac.in/sites/default/files/inline-files/Crowley-Wong.pdf>
27. When Is a Liability Not a Liability Textual Analysis, Dictionaries, and 10-Ks - UTS, February 2011, https://www.uts.edu.au/globalassets/sites/default/files/adg_cons2015_loughran-mcdonald-je-2011.pdf
28. Sentiment analysis based on a social media customised dictionary - PMC, <https://pmc.ncbi.nlm.nih.gov/articles/PMC8374646/>
29. The colour of finance words* - University of Colorado Boulder, November 8, 2022, <https://leeds-faculty.colorado.edu/garcia/finsentiment20221108.pdf>
30. Textual Analysis in Finance - Annual Reviews, <https://www.annualreviews.org/doi/pdf/10.1146/annurev-financial-012820-032249>
31. Sentiment Analysis of Economic Text: A Lexicon-Based Approach - arXiv.org, November 21, 2024, <https://arxiv.org/html/2411.13958v1>
32. Loughran-McDonald Master Dictionary w/ Sentiment Word Lists | Software Repository for Accounting and Finance, <https://sraf.nd.edu/loughranmcdonald-master-dictionary/>
33. Evolution from Traditional NLP Models to State-of-the-Art LLMs | by Noor Fatima | Medium, <https://medium.com/@noorfatimaafzalbutt/evolution-from-traditional-nlp-models-to-state-of-the-art-llms-00cef9829a65>
34. Full article: Narrative and computational text analysis in business and economic history, <https://www.tandfonline.com/doi/full/10.1080/03585522.2021.1984299>
35. Transformers in Sentiment Analysis: A Paradigm Shift in Deep Learning Research - Journal of Information Systems Engineering and Management, <https://jisem-journal.com/index.php/journal/article/download/612/210>
36. FINANCIAL SENTIMENT ANALYSIS IN BIST100 COMPANIES' ANNUAL REPORTS - Middle East Technical University, <https://open.metu.edu.tr/bitstream/handle/11511/113400/10565277.pdf>
37. FinBERT: Financial Sentiment Analysis with Pre-trained Language Models - ar5iv - arXiv, August 27, 2019, <https://ar5iv.labs.arxiv.org/html/1908.10063>
38. Pre-trained Large Language Models for Financial Sentiment Analysis - arXiv.org, January 10, 2024, <https://arxiv.org/html/2401.05215v1>
39. Does sentiment help in asset pricing? A novel approach using large language models and market-based labels - SSRN, <https://papers.ssrn.com/sol3/Delivery.cfm/4905533.pdf?abstractid=4905533&mirid=1>
40. Financial report sentiment analysis using Loughran-McDonald dictionary and BERT - WSEAS, 2024, [https://wseas.com/journals/fe/2024/a30fe-017\(2024\).pdf](https://wseas.com/journals/fe/2024/a30fe-017(2024).pdf)
41. Comparison of BERT, FinBERT and traditional model predictions with... - ResearchGate, https://www.researchgate.net/figure/Comparison-of-BERT-FinBERT-and-traditional-model-predictions-with-actual-categories_fig2_391446592
42. Why LLMs fail on financial documents and how AI extraction can finally fix it - Desia, <https://www.desia.ai/news/why-llms-fail-financial-documents-ai-extraction>
43. Natural language processing in finance: A survey - SenticNet, <http://www.sentic.net/nlp-in-finance.pdf>

44. Financial Statement Analysis with Large Language Models, May 2024, <https://bfi.uchicago.edu/wp-content/uploads/2024/05/Financial-Statement-Analysis-with-Large-Language-Models.pdf>
45. Financial Statement Analysis with Large Language Models - Fama Miller Center, <https://www.chicagobooth.edu/research/fama-miller/finance-research/funding/a-demand-system-approach-for-fixed-income/financial-statement-analysis-with-large-language-models>
46. The Untold Story of Alex Kim – Financial Statement Analysis with Large Language Models, https://www.reddit.com/r/ArtificialIntelligence/comments/1fb48pk/the_untold_story_of_alex_kim_financial_statement/
47. [2407.17866] Financial Statement Analysis with Large Language Models - arXiv.org, July 25, 2024, <https://arxiv.org/abs/2407.17866>
48. Financial Statement Analysis with Large Language Models - arXiv, July 25, 2024, <https://arxiv.org/html/2407.17866v1>
49. UChicago: GPT better than humans at predicting earnings : r/quant - Reddit, https://www.reddit.com/r/quant/comments/1d2mxg3/uchicago_gpt_better_than_humans_at_predicting/
50. Ongoing Research Projects | The University of Chicago Booth School of Business, <https://faculty.chicagobooth.edu/valeri-nikolaev/ongoing-research-projects>
51. Generic LLMs vs. Domain-Specific LLMs: What's the Difference? - Dataversity, <https://www.dataversity.net/articles/generic-llms-vs-domain-specific-llms-whats-the-difference/>
52. Smarter, smaller, safer: The case for small language models in financial services - Infosys, <https://www.infosys.com/iki/perspectives/small-language-models-financial-services.html>
53. Beyond Classification: Financial Reasoning in State-of-the-Art Language Models - arXiv, May 2, 2023, <https://arxiv.org/html/2305.01505v2>
54. Chain-of-Thought (CoT) Prompting in AI-Powered Financial Analysis, <https://corporatefinanceinstitute.com/resources/financial-modeling/chain-of-thought-prompting-financial-analysis/>
55. FinCoT: Grounding Chain-of-Thought in Expert Financial Reasoning - arXiv, June 19, 2025, <https://arxiv.org/html/2506.16123v1>
56. Reasoning or Overthinking: Evaluating Large Language Models on Financial Sentiment Analysis - arXiv, June 5, 2025, <https://arxiv.org/html/2506.04574v1>
57. Fin-R1: A Large Language Model for Financial Reasoning through Reinforcement Learning - arXiv, March 20, 2025, <https://arxiv.org/html/2503.16252v3>
58. Fin-R1: A Large Language Model for Financial Reasoning through Reinforcement Learning, March 20, 2025, <https://arxiv.org/html/2503.16252v1>
59. Augmenting large language models for financial sentiment analysis: a heuristic sparse mixture-of-experts framework - PeerJ, <https://peerj.com/articles/cs-3607/>
60. README_en.md · SUFE-AIFLM-Lab/Fin-R1 at main - Hugging Face, https://huggingface.co/SUFE-AIFLM-Lab/Fin-R1/blob/main/README_en.md
61. Fin-R1:A Specialized Large Language Model for Financial Reasoning and Decision-Making : r/LocalLLaMA - Reddit, https://www.reddit.com/r/LocalLLaMA/comments/1jk97sp/finr1a_specialized_large_language_model_for/
62. Fin-R1's Financial Reasoning: Excels in Financial Table & Conversation AI - Medium, <https://medium.com/aimonks/fin-r1s-financial-reasoning-excels-in-financial-table-conversation-ai-7625085ef057>
63. The-FinAI/Fino1: This is the repo of developing reasoning models in the specific domain of financial, aim to enhance models capabilities in handling financial reasoning tasks. - GitHub, <https://github.com/The-FinAI/Fino1>