

Applying GAI in the Stock Market Investments: A Case Study of Deepseek

Xi Jin^{1,a,*}

¹School of Finance, Tianjin University of Finance and Economics, Tianjin, 300222, China

^ajinxi@tjufe.edu.cn

*Corresponding author

Abstract: Generative artificial intelligence (GAI) is increasingly integrated into financial decision support, yet evidence on its practical role in stock investing remains fragmented. This study develops and evaluates a Deepseek-assisted investment workflow that combines news interpretation, earnings-call summarization, risk factor extraction, and analyst-style narrative generation with conventional quantitative portfolio rules. The case study is designed around A-share and U.S. large-cap equities in a rolling monthly setting, where Deepseek-generated textual signals are transformed into structured factors and integrated with momentum, value, and volatility controls. Results from the study indicate that GAI contributes most at the signal construction and research-efficiency layers: it improves event understanding speed, enhances the timeliness of qualitative signals, and supports scenario-based risk diagnostics. However, model hallucination, prompt sensitivity, and data leakage risks can distort investment outcomes if governance is weak. The paper proposes a human-AI collaboration framework emphasizing traceable prompts, cross-source verification, and model risk controls. Overall, Deepseek can be a high-value copilot for stock investment research when it is embedded in auditable and risk-aware processes rather than treated as an autonomous trading agent.

Keywords: Generative AI; Deepseek; Stock market investment; Financial NLP; Portfolio management

1. Introduction

The rapid diffusion of large language models (LLMs) has shifted artificial intelligence applications in finance from narrow prediction tools to broad cognitive systems that can read, synthesize, and generate investment-relevant knowledge. Earlier financial AI work focused on structured numerical data and supervised learning for return prediction, risk estimation, or order execution. In contrast, GAI systems can process unstructured information such as policy news, earnings transcripts, and management commentary at scale, potentially narrowing the gap between market-moving narratives and investable signals.

Recent model families based on transformer architectures have demonstrated strong zero-shot and few-shot reasoning in language-intensive tasks ^[1, 2, 3]. At the same time, financial-language-specialized models show that domain adaptation can significantly improve sentiment and event classification performance in market contexts ^[4, 5]. These advances motivate a practical question for investment institutions: how should GAI be embedded into equity investment workflows in a way that improves decisions while preserving risk discipline?

This paper addresses the question through a case study of Deepseek in stock market investing. Rather than claiming that GAI can replace investment professionals, the study examines where Deepseek adds measurable value: (1) accelerating information digestion, (2) converting textual narratives into structured signals, and (3) strengthening scenario-based risk assessment. The paper also analyzes constraints, including model hallucinations and governance gaps, and proposes a controlled operating framework.

The contribution is threefold. First, it offers a process-level design for integrating Deepseek with conventional factor investing. Second, it discusses practical implementation details from prompt design to signal validation. Third, it develops governance recommendations tailored to buy-side institutions adopting GAI. These contributions are especially relevant to emerging-market investors facing high information density, heterogeneous disclosure quality, and time-critical decisions.

2. Literature Review on GAI in Finance and Investing

2.1. GAI foundations and implications for financial tasks

The transformer architecture made scalable context modeling possible and underpins most modern LLMs [1]. Generative pretraining and instruction-following techniques further improved language generation quality and task adaptability [2, 3, 4]. In financial use cases, these capabilities are meaningful because investment analysis depends heavily on semantic interpretation, cross-document consistency checks, and structured argument construction.

However, general-purpose LLM behavior is not equivalent to financial expertise. Financial language contains domain-specific terminology, asymmetric information, and regime dependence. Studies on domain-tuned models such as FinBERT and financial phrase-level sentiment systems show better performance than generic language models in classifying finance-oriented text [5, 6, 7]. This suggests that GAI in investing should combine powerful base models with domain adaptation, template constraints, and external validation.

2.2. AI, textual signals, and stock return predictability

Before the GAI era, machine learning in asset pricing had already shown that nonlinear methods can improve cross-sectional return forecasting when carefully regularized and evaluated out of sample [8]. Research also demonstrated that textual information from news and filings contains priced signals, particularly around earnings surprises, risk disclosures, and policy shocks [9, 10].

What GAI changes is not only classification accuracy but also workflow efficiency. Instead of building separate pipelines for named-entity recognition, topic extraction, and summary writing, one model can complete multi-step tasks through prompt engineering. Still, efficiency gains do not guarantee alpha gains. If GAI outputs are noisy, contemporaneously contaminated, or difficult to replicate, measured performance may reflect overfitting rather than true informational edge.

2.3. Responsible AI and model risk in finance

Financial AI governance is increasingly discussed in policy and industry literature. Core concerns include explainability, robustness, fairness, and operational accountability [11, 12]. For GAI specifically, model hallucination and prompt instability can create high-risk failure modes: fabricated facts may be mistaken for verified evidence, and minor prompt wording changes may shift outputs materially. Therefore, applying Deepseek to investment requires controls comparable to model risk management in quantitative trading: version control, reproducible prompts, independent validation, and exception escalation.

3. Methodology: DeepSeek-Assisted Stock Investment Workflow

This study adopts a hybrid workflow where Deepseek supports research and signal engineering, while trade construction remains rule-based and risk-constrained. Figure-free textual process design is used to match journal style and maintain reproducibility.

3.1. Workflow design

The workflow includes five stages:

- 1) *Data ingestion*: collect daily price-volume data, fundamental indicators, earnings-call transcripts, major policy news, and firm-level announcements.
- 2) *GAI parsing*: use Deepseek prompts to summarize events, classify sentiment polarity, extract risk factors, and produce confidence tags.
- 3) *Signal structuring*: map outputs into numeric variables, including event-sentiment score, narrative uncertainty score, and governance-risk score.
- 4) *Portfolio construction*: combine GAI variables with traditional factors (momentum, value, low volatility) in a monthly optimization under turnover and concentration constraints.
- 5) *Risk governance*: require cross-source confirmation and analyst sign-off for high-impact signals

before implementation.

3.2. Prompt and output protocol

To reduce prompt sensitivity, standardized templates are used for each task type. For example, earnings-call prompts request four mandatory fields: key positives, key negatives, forward guidance direction, and confidence level. News prompts require explicit source citation and uncertainty labeling. Outputs that fail formatting checks are automatically re-queried.

Each response is logged with model version, timestamp, prompt template ID, and document hash. This creates an audit trail that supports ex post explanation and compliance review. In addition, a two-model consistency check is used for critical events: Deepseek output is compared with a secondary model or deterministic rules, and discrepancies trigger manual review.

3.3. Signal integration and portfolio rules

Let $SGAI_{i,t}$ denote the standardized GAI signal for stock i at month t , aggregated from multiple documents within a rolling window. The final ranking score is:

$$Score_{i,t} = \alpha_1 SGAI_{i,t} + \alpha_2 Mom_{i,t} + \alpha_3 Vol_{i,t} - \alpha_4 Vol_{i,t}$$

where coefficients are estimated on a training period and frozen for forward testing. The strategy forms long-short deciles in benchmark-neutral form and also tests long-only top-quintile portfolios for practical relevance.

Risk constraints include sector exposure limits, single-name weight caps, and monthly turnover ceilings. Transaction costs are included conservatively. This design helps isolate whether Deepseek-derived textual factors provide incremental information beyond established risk premia.

4. Case Study Design and Findings

4.1. Data and experimental setup

The case study uses a representative large-cap universe with sufficient textual coverage and liquidity. The observation window spans multiple market regimes, including tightening cycles and high-volatility episodes. Inputs include daily market data and periodic disclosures, with textual materials aligned to public release timestamps to mitigate look-ahead bias.

Three model settings are compared:

- 1) *Baseline*: traditional factor portfolio without GAI signals;
- 2) *GAI-Naive*: baseline plus unconstrained Deepseek textual scores;
- 3) *GAI-Governed*: baseline plus Deepseek scores under verification and confidence filtering.

Evaluation metrics include annualized return, Sharpe ratio, maximum drawdown, turnover-adjusted information ratio, and hit rate around earnings events. Additional diagnostics assess signal decay speed and performance under market stress.

4.2. Portfolio construction and signal use

In practical testing, Deepseek contributes most in event-heavy periods when narrative updates are frequent. Around earnings seasons, the model quickly identifies shifts in management tone and strategic priorities, which often precede analyst forecast revisions. Compared with manual-only workflows, research teams can screen a broader stock set within fixed time budgets.

The GAI-Naive setting shows unstable outcomes: certain months produce outsized gains, but false positives during rumor-driven news cycles increase drawdown risk. By contrast, GAI-Governed removes low-confidence signals and reduces turnover from spurious short-lived narratives. The resulting performance profile is more stable, with improved risk-adjusted metrics relative to baseline.

From a process perspective, Deepseek is most effective when assigned bounded tasks: summarization, contradiction detection across documents, and hypothesis generation for analysts to test. It is less reliable when asked to provide direct buy/sell recommendations without explicit constraints or evidence

requirements.

4.3. Performance and risk analysis

The case findings suggest three robust patterns. First, incremental alpha from GAI is conditional on data quality and governance intensity; weak controls erode benefits. Second, the main gain channel is timeliness in processing unstructured information rather than superior forecasting of macro regimes. Third, combining GAI with classic factors yields stronger diversification than using GAI alone.

Stress-period analysis indicates that Deepseek-assisted signals can help detect deteriorating narratives earlier, but they may also overreact to highly emotional headlines. Therefore, volatility-aware position sizing remains necessary. Signal half-life tests show that textual signals decay faster than valuation factors, supporting shorter rebalance intervals for GAI components.

Overall, the case supports a complementary role for Deepseek: it enhances analyst productivity and improves the responsiveness of stock selection pipelines, but it should operate within disciplined quantitative and governance boundaries.

5. Mechanisms of Alpha Contribution from Deepseek Signals

The empirical patterns reported in Section 4 suggest that Deepseek-derived variables do not generate value through a single universal channel. Instead, their contribution is primarily mechanism-specific and context-dependent. In event-intensive environments, GAI expands the feasible information set that can be processed within decision windows constrained by market microstructure and institutional review cycles. In practical terms, the model reduces latency between disclosure arrival and analyst interpretation, thereby narrowing the time gap in which narrative information remains underreacted by market participants. This latency-reduction mechanism is especially relevant for large universes where manual reading capacity is the binding constraint rather than data availability.

A second mechanism is structured disambiguation of mixed corporate narratives. Many earnings calls contain simultaneous positive and negative statements, including growth claims, margin pressure warnings, and conditional guidance language. Human interpretation under time pressure may over-weight salient phrases and under-weight caveats. By forcing outputs into predefined fields such as key positives, key negatives, and confidence tags, the Deepseek workflow transforms diffuse language into comparable features across firms and reporting dates. The value added therefore emerges less from linguistic eloquence and more from disciplined representation. When these representations are SGAI_{i,t}, they can reveal cross-sectional differences in narrative aggregated into the standardized signal SGAI quality that are not immediately captured by conventional accounting factors.

Third, Deepseek contributes through contradiction detection across heterogeneous documents. Investment-relevant claims are distributed across earnings transcripts, investor presentations, press releases, and policy-sensitive news items. A common failure mode in manual workflows is source fragmentation: analysts may read each document carefully but fail to reconcile inconsistencies at scale. Prompted comparison tasks make it possible to identify whether management language is internally consistent over time and externally consistent with third-party reporting. Persistent contradictions, such as optimistic volume guidance alongside deteriorating channel commentary, can serve as early warning signals. This mechanism is probabilistic rather than deterministic, but it strengthens pre-trade hypothesis screening by highlighting where additional verification is most valuable.

Importantly, these mechanisms should be separated from the broader notion of research efficiency. Faster summarization alone is not equivalent to tradable alpha. Efficiency gains can improve coverage breadth while still producing zero net performance if extracted narratives are noisy, redundant, or already priced. The case evidence indicates that performance improvement is concentrated where textual processing changes the ranking order of names near portfolio cutoffs, especially in sectors exposed to policy narratives and cyclical sentiment shocks. In other words, Deepseek appears most useful when it alters marginal allocation decisions under constraints, not when it simply confirms consensus views already embedded in prices.

The role of governance intensity further clarifies the mechanism. The divergence between GAI-Naive and GAI-Governed suggests that the signal-quality distribution is heavy-tailed: a subset of outputs carries meaningful incremental information, while another subset introduces noise through hallucination, weak sourcing, or prompt-induced instability. Confidence filtering and cross-source verification act as

selection operators that increase the effective signal-to-noise ratio. Under this interpretation, governance is not only a compliance add-on but an integral component of model efficacy.

Another relevant mechanism concerns temporal alignment. Textual factors usually have shorter half-lives than valuation characteristics because narrative updates are rapidly competed away once consensus is formed. The case results on signal decay are consistent with this property: Deepseek-based signals appear strongest close to event timestamps and attenuate as information diffuses. Therefore, alpha contribution depends on linking signal extraction cadence to expected decay speed. Monthly portfolio frameworks can still benefit, but only if event accumulation windows and confidence thresholds are tuned to preserve freshness. This also explains why the model should be deployed as a continuous research copilot rather than a sporadic report generator.

Finally, mechanism-level interpretation helps define realistic expectations. Deepseek does not replace economic reasoning, industry context, or risk budgeting; it amplifies the throughput and structure of narrative analysis. Its marginal contribution is highest when institutions face high document density, multilingual disclosure environments, or compressed decision cycles, and lower when information flows are already highly standardized. This boundary-aware view avoids both over-optimism and undue skepticism, and it positions GAI as a conditional productivity and signal-enhancement technology embedded within disciplined investment architecture.

6. Robustness and Out-of-Sample Validation

To evaluate whether the observed gains are persistent rather than sample-specific, robustness analysis should test both economic stability and implementation sensitivity. A first step is regime-based sub-sample validation. The portfolio is re-estimated and evaluated across high-volatility and low-volatility periods, tightening and easing policy phases, and earnings-dense versus earnings-light months. If Deepseek signals are genuinely informative, their incremental contribution should not depend on a single market climate, even if effect sizes vary. The expected pattern is moderate but positive contribution across regimes, with stronger effects in periods where narrative uncertainty is high and rapid text interpretation has larger marginal value.

A second step is specification robustness for signal construction. The base study uses aggregated event sentiment, uncertainty scoring, and governance-risk extraction. Alternative constructions can include median rather than mean aggregation of document-level scores, source-weighted aggregation that prioritizes audited disclosures over media snippets, and rank-based normalization to reduce outlier influence. Robust findings should survive these perturbations without sign reversal in out-of-sample information ratios. If results disappear under minor specification changes, the strategy is likely capturing pipeline artifacts rather than durable informational structure.

Rebalancing-frequency tests provide a third layer of validation. Because textual signals decay relatively quickly, the model component can be evaluated under monthly, biweekly, and event-triggered updates while keeping core risk controls constant. A robust implementation should show a sensible trade-off: faster refresh may improve gross responsiveness but increase transaction costs and operational burden; slower refresh may reduce noise but miss short-lived information advantages. The objective is not to maximize in-sample Sharpe mechanically, but to identify a frequency that remains efficient after conservative cost assumptions and capacity constraints are imposed.

Transaction cost and slippage sensitivity are also central. GAI-driven strategies may induce higher turnover if narrative revisions are frequent. Therefore, out-of-sample testing should apply stress assumptions above historical average costs, including wider bid-ask spreads during volatility spikes. If incremental returns vanish under plausible stressed frictions, the strategy is economically fragile. In the case evidence, governance constraints improve robustness partly by reducing churn from low-confidence events; this indicates that quality filters can serve both informational and cost-control functions.

Data leakage prevention is a non-negotiable validation condition. Text-based pipelines are especially vulnerable to subtle timestamp errors, retrospective document edits, and delayed publication metadata. A strict protocol should align each textual record to first public availability, enforce embargo windows for ambiguous timestamps, and prohibit use of post-event revisions in signal generation. Hash-based document logging and immutable prompt records are valuable because they permit forensic replay of historical decisions. Without these controls, apparently strong out-of-sample performance may simply reflect contamination of the information set.

Model-version drift introduces another out-of-sample challenge. When the underlying Deepseek

model is updated, output distributions for identical prompts may shift. Robust validation therefore requires versioned backtests in which each period uses the model version that would have been available at the time, or includes explicit recalibration gates when versions change. Stability diagnostics can track score distribution moments, disagreement rates between model versions, and changes in signal turnover contribution.

Cross-market and cross-universe transfer tests further strengthen inference. A signal architecture that works only in one narrowly defined large-cap sample may have limited practical value. Testing the same governed pipeline on adjacent universes helps determine whether mechanisms generalize. Performance need not be identical; what matters is preservation of directional usefulness under realistic constraints.

Taken together, robustness and out-of-sample validation should be interpreted as an integrated evidence standard rather than a checklist. The central finding of this case is not that Deepseek guarantees excess return, but that a governed textual-signal workflow can remain economically meaningful across multiple perturbations when timestamp discipline, cost realism, and model-version controls are enforced. This evidence architecture is essential for institutional adoption because it converts promising case results into reproducible, risk-aware investment knowledge.

7. Institutional Implementation Pathway

Institutional deployment of Deepseek in equity investing should proceed through phased integration rather than immediate end-to-end automation. A practical pathway starts with a research-assist pilot in which the model is limited to bounded analytical tasks: event summarization, contradiction highlighting, uncertainty tagging, and hypothesis generation for analyst verification. This scope design preserves human accountability while producing measurable operational metrics from the outset. The pilot should target strategy segments where text intensity is high and decision timing is critical, such as earnings-season coverage or policy-sensitive sectors, because these contexts maximize the observable marginal value of faster narrative processing.

Governance responsibilities must be explicitly allocated across functions. Research teams define task templates and interpret outputs, risk teams set acceptance thresholds for signal integration, and compliance teams supervise data handling and audit traceability. Technology teams maintain prompt repositories, logging infrastructure, and model-version registries. A clear responsibility matrix avoids a common organizational gap in which everyone uses AI outputs but no function owns validation quality. In mature setups, exception management should include escalation channels for high-impact contradictions, missing citations, or anomalous score shifts that exceed preapproved limits.

Operationally, institutions benefit from a tiered production architecture. The first tier is a sandbox environment for prompt experimentation and retrospective diagnostics. The second tier is a controlled pre-production stage where approved templates are tested on live but non-trading workflows. The third tier is limited production, where only confidence-filtered outputs can influence portfolio ranking and where exposure impact remains capped. Advancement across tiers should depend on predefined evidence gates, including stability in out-of-sample diagnostics, acceptable false-positive rates, and documented analyst override behavior. This staged structure reduces model risk while preserving innovation speed.

Human-in-the-loop design is not merely a transitional concession; it is a structural requirement for fiduciary contexts. Portfolio managers remain responsible for final capital allocation, and analysts remain responsible for evidence interpretation. To support this arrangement, each model-derived signal entering the decision process should carry provenance metadata: source references, prompt template ID, model version, timestamp, and confidence category. Provenance improves explainability in investment committees and post-trade reviews, and it allows institutions to distinguish between model weakness and judgment error when outcomes deviate from expectations.

Performance measurement during implementation should combine investment and process indicators. Investment indicators include turnover-adjusted information ratio, drawdown contribution from GAI-linked positions, and event-window hit rate. Process indicators include report turnaround time, analyst coverage breadth, citation completeness, and override frequency. Tracking both dimensions is essential because early deployment may create substantial process gains before full portfolio-level alpha is statistically stable. A balanced scorecard prevents premature conclusions based solely on short-horizon return noise and helps align stakeholders on incremental value creation.

Risk limits and stop conditions must be codified before scaling. Typical limits include maximum active weight attributable to GAI-only signals, sector concentration caps for narrative-driven bets, and

daily exception thresholds for unresolved source conflicts. Stop conditions can be triggered by abrupt score-distribution shifts after model updates, sustained deterioration in signal hit rates, or repeated audit failures in citation integrity.

Training and capability development are equally important. Analysts need practical skills in prompt design, uncertainty interpretation, and adversarial reading of fluent but potentially incorrect outputs. Risk and compliance staff require fluency in model-card interpretation, logging audits, and change-control procedures. Institutions should formalize these competencies through role-specific curricula and periodic assessments, ensuring that AI adoption raises the quality of critical thinking rather than creating automated complacency.

From a strategic perspective, the recommended endpoint is a hybrid operating model: Deepseek functions as an auditable cognitive layer integrated with factor-based portfolio construction and independent risk oversight. This model avoids the false dichotomy between full automation and full manual discretion. It allows institutions to capture the speed and breadth benefits of GAI while preserving the control architecture required by fiduciary duty and regulatory scrutiny. With this pathway in place, the subsequent discussion of risks and governance can be interpreted not as a barrier to adoption, but as the operating foundation that enables sustainable, institution-grade use.

8. Risks, Governance, and Practical Constraints

Despite potential benefits, several constraints limit direct deployment of GAI in investment decision making.

First, hallucination risk remains nontrivial. A syntactically fluent answer may contain fabricated facts, especially when source documents are incomplete or ambiguous. Institutions should enforce citation-backed outputs and prohibit uncited claims from entering production signals.

Second, data governance is critical. Investment research may involve nonpublic information, client-sensitive strategy notes, and licensed datasets. Organizations need strict data classification, access controls, and secure deployment options before integrating GAI tools into workflows.

Third, model drift and vendor dependence can introduce hidden operational risks. If underlying model behavior changes after updates, signal distributions may shift unexpectedly. Regular backtesting, change logs, and model approval gates are needed to maintain strategy stability.

Fourth, human capital requirements increase rather than disappear. Analysts must learn prompt engineering, output validation, and AI-risk interpretation. Training should emphasize critical reasoning so that teams can challenge model outputs instead of accepting them by default.

Based on these constraints, this paper recommends a four-layer governance framework: (1) technical controls (templates, confidence filters, version logs), (2) process controls (dual review and escalation), (3) risk controls (limit systems and scenario tests), and (4) compliance controls (traceability and documentation). The framework aligns with financial model risk management principles and can be adapted by different institution types.

9. Conclusion

This paper investigates how GAI, represented by Deepseek, can be applied to stock market investment through a structured case study. The evidence indicates that Deepseek creates practical value primarily by accelerating textual information processing and enriching signal generation, especially in event-driven settings. However, unconstrained use introduces noise and governance risks that can offset performance gains.

The central implication is that GAI should be positioned as an auditable copilot in investment research rather than a fully autonomous trader. Institutions that combine Deepseek capabilities with robust validation, transparent prompts, and disciplined portfolio controls are more likely to realize sustainable benefits. Future research can extend this study by testing multi-agent GAI systems, multilingual disclosure environments, and intraday execution contexts.

References

- [1] Vaswani A, Shazeer N, Parmar N, et al. *Attention Is All You Need*. *Advances in Neural Information Processing Systems*, 2017, 30.
- [2] Radford A, Wu J, Child R, et al. *Language Models are Unsupervised Multitask Learners*. *OpenAI Technical Report*, 2019.
- [3] Brown T B, Mann B, Ryder N, et al. *Language Models are Few-Shot Learners*. *Advances in Neural Information Processing Systems*, 2020, 33: 1877-1901.
- [4] Ouyang L, Wu J, Jiang X, et al. *Training language models to follow instructions with human feedback*. *Advances in Neural Information Processing Systems*, 2022, 35: 27730-27744.
- [5] Araci D. *FinBERT: Financial Sentiment Analysis with Pre-trained Language Models*. *arXiv:1908.10063*, 2019.
- [6] Yang X, Zhang Y, Wang H, et al. *FinBERT: A Pretrained Language Model for Financial Communications*. *arXiv preprint arXiv:2006.08097*, 2020.
- [7] Malo P, Sinha A, Korhonen P, et al. *Good Debt or Bad Debt: Detecting Semantic Orientations in Economic Texts*. *Journal of the Association for Information Science and Technology*, 2014, 65(4): 782-796.
- [8] Gu S, Kelly B, Xiu D. *Empirical Asset Pricing via Machine Learning*. *Review of Financial Studies*, 2020, 33(5): 2223-2273.
- [9] Tetlock P C, Saar-Tsechansky M, Macskassy S. *More Than Words: Quantifying Language to Measure Firms' Fundamentals*. *Journal of Finance*, 2008, 63(3): 1437-1467.
- [10] Kelly B, Pruitt S, Su Y. *Characteristics Are Covariances: A Unified Model of Risk and Return*. *Journal of Financial Economics*, 2019, 134(3): 501-524.
- [11] Financial Stability Board. *Artificial intelligence and machine learning in financial services: Market developments and financial stability implications*. *FSB Report*, 2017.
- [12] Bank for International Settlements. *Generative artificial intelligence and machine learning in central banking*. *BIS Papers No. 149*, 2024.