
Localization and Steering of Economic Uncertainty in Large Language Models

Rice (Xilin) Wang
Northeastern University
wang.xil@northeastern.edu

Verónica C. Pérez
Boston University
vcperez@bu.edu

Claire Schlesinger
Northeastern University
schlesinger.e@northeastern.edu

Luze Sun
Northeastern University
sun.luz@northeastern.edu

Abstract

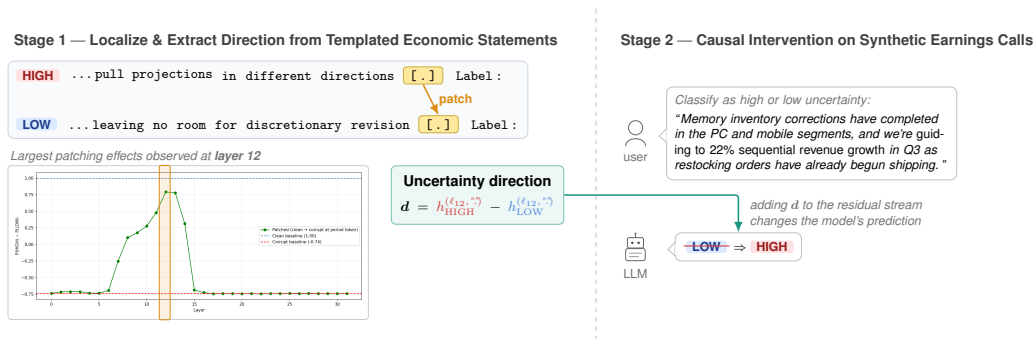
Economic uncertainty shapes investment, hiring, and asset prices, making its measurement from text a central task in economics and finance. While LLMs accurately measure uncertainty in financial texts such as earnings call transcripts, it remains unclear whether they develop coherent internal representations or merely pattern-match on surface lexical cues. In this study, we discover that LLMs *linearly represent economic uncertainty* with a single direction in the residual stream. Using activation patching on two synthetic datasets of contrastive earnings-call statements with varying linguistic styles, we localize this direction and find that models aggregate the uncertainty signal at the final token regardless of lexical patterns. The extracted direction perfectly separates held-out high- and low-uncertainty statements both within- and cross-dataset. Causal interventions show that adding or subtracting this direction monotonically flips uncertainty predictions, including cross-dataset transfer from templated to naturalistic text. Finally, in a downstream portfolio allocation task using real earnings-call excerpts, steering along the uncertainty direction shifts model investment toward safe assets, consistent with economic theory. Together, our results establish that LLMs encode economic uncertainty as a structured, causally active, and transferable representation, offering a foundation for interpretability-based auditing and control of LLMs deployed in financial analysis.

1 Introduction

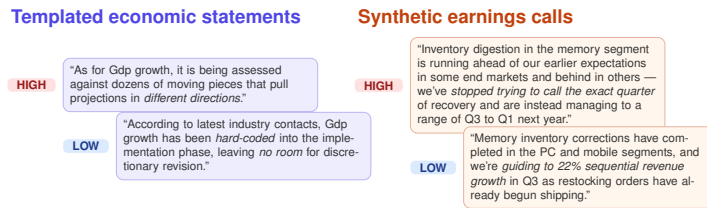
Economic uncertainty—the degree to which the future path of key economic variables is unknown—shapes investment, hiring, and asset prices [Bloom et al., 2018, Baker et al., 2016, Romer, 1990, Pástor and Veronesi, 2013]. Measuring it accurately from text has therefore attracted sustained attention across economics and finance. Earnings calls, quarterly conferences in which executive and analysts discuss financial results and future outlook into the economy, have emerged as a particularly valuable source of information for this work. Unlike news articles, earnings calls combine forward-looking firm perspectives with real-time analyst questions, making them a rich, unscripted, source of information [Hassan et al., 2019, 2023, 2024, 2025]. Traditional measurement methods for uncertainty, such as Bag-of-Words (BoW), have recently given way to LLM-based approaches that offer richer contextual sensitivity [Clayton et al., 2025, Audrino et al., 2024], and practitioners have begun deploying “agentic financial analysis” models, to process, and synthesize earnings-call content at scale [Godsall et al., 2025, Gupta et al., 2025]. Yet these LLM results remain fundamentally opaque: when a model classifies an earnings-call statement as expressing *high* uncertainty, it is unclear what internal computation it performs. Prior work has focused on applying pre-trained Sparse

Autoencoders (SAEs) to large language models, discovering several financial concepts including risk-related ones from the internal activations Chen et al. [2026]. However, whether a similarly coherent representation exists for uncertainty or not in the earnings-call domain remains unexamined. Do models activate a structured, stable representation of the economic uncertainty, or does it rely on diffuse lexical cues that happen to correlate with the label? For policy-relevant analysis, this distinction has practical stakes: a model with a coherent internal uncertainty feature can be audited and steered reliably; one that merely pattern-matches on surface tokens cannot.

In this work, we take a mechanistic approach to discover internal representations of economic uncertainty in LLMs on earnings-call-style text (Figure 1). First we aim to localize a model representation that linearly encodes uncertainty. On two synthetically generated datasets of contrastive pairs of earnings call statements, we use activation patching to identify a localized uncertainty direction from the residual stream activations. Our experiments show that this direction generalizes to held-out samples with high classification accuracy. We then ask: is this direction causally implicated in the model’s predictions across settings? Through causal intervention, we find that this direction actively governs predictions across settings: adding it to the residual stream moves predictions towards high-uncertainty labels and vice versa. Crucially, the uncertainty direction is transferrable: extracted from the programmatically generated, templated dataset, it exhibits high intervention effectiveness on more natural, lexically diverse texts. Finally, we study the usefulness of the uncertainty representation in a realistic investment setting – we ask the LLM to behave as a strategic financial analyst, analyzing a company’s financial data history and real earnings call statements and allocating between the company’s stock and low-risk U.S. treasuries. Steering the uncertainty direction consistently controls the LLM’s allocation, where a positive steering produces much more conservative investment strategies. Together, these results offer a detailed insight into how LLMs internally represent economic uncertainty, shedding light on practical usages of interpretability techniques on promoting reliability and controllability of LLM-based economic measurement.



(a) Pipeline for uncertainty direction extraction and causal intervention.



(b) Example contrastive statement pairs from each dataset.

Figure 1: Discovering economic uncertainty in Large Language Models. **Stage 1:** We localize uncertainty information via activation patching, with layer 12 yielding the largest patching effect (llama 3.1 8b instruct); we then compute mean activation differences between high- and low-uncertainty statements to extract the uncertainty direction. **Stage 2:** Causal intervention reliably changes the model’s predictions on synthetic earnings call statements.

2 Related Work

The measurement of economic uncertainty from text has evolved in two phases. The first relies on Bag-of-Words (BoW) methods: Baker et al. [2016] construct an Economic Policy Uncertainty (EPU) index by counting co-occurrences of economy-, policy-, and uncertainty-related terms in newspaper archives, an approach extended to geopolitical risk [Caldara and Iacoviello, 2022], social media [Baker et al., 2021], and corporate earnings calls [Hassan et al., 2019, 2023]. Loughran and McDonald [2011] refined this paradigm for finance, demonstrating that domain-specific word lists outperform general-purpose sentiment lexicons on 10-K filings. The second phase leverages LLMs as richer, context-sensitive classifiers: Clayton et al. [2025] prompt LLMs to extract geopolitical pressure signals from full earnings call transcripts, while Audrino et al. [2024] quantify monetary policy, financial-market, and geopolitical uncertainty from newspaper text, obtaining stronger correlations with macroeconomic variables than BoW baselines. Beyond uncertainty, LLMs have shown strong performance across economic text analysis tasks including consumption forecasting [Chen et al., 2022], business-cycle monitoring [Bybee, 2023], and central bank communication analysis [Fang et al., 2025]. Whether these gains reflect genuine conceptual understanding or sophisticated pattern-matching remains an open question.

The earliest form of feature interpretability in economic text analysis was the BoW model itself: Tetlock [2007] showed that word-frequency features from Wall Street Journal columns predict next-day excess stock returns, establishing that textual features carry genuine predictive content for economic outcomes. On the other hand, mechanistic interpretability approaches, focusing mainly on LLMs, provides techniques for understanding specific computations within transformer networks. Probing classifiers, linear models trained on frozen intermediate-layer activations, reveal that representations encode rich linguistic properties in a layered, structured fashion [Tenney et al., 2019, Belinkov, 2022]. The *linear representation hypothesis* formalizes this by positing that high-level concepts correspond to directions in the model’s residual stream [Park et al., 2024]. Activation patching operationalizes this idea causally: by replacing activations from one forward pass with those of a contrastive run, one can precisely localize which layer carries information critical to a given behavior [Vig et al., 2020a, Meng et al., 2022a]; this technique has been used to reverse-engineer functional circuits for syntactic agreement [Vig et al., 2020a] and indirect-object identification [Wang et al., 2023]. Applied to finance, Chen et al. [2026] use sparse autoencoders over learned neural features to decompose LLM activations from financial news into interpretable clusters (sentiment, risk, timing) that predict asset returns and can be steered to correct optimism bias. Our work differs from theirs along two dimensions. First, their approach discovers structure post-hoc and establishes correlational links between features and outputs, while we start from a theoretically grounded economic concept and use activation patching to establish that the extracted representation *causally* governs model predictions. Second, we operate on earnings calls rather than financial news; crucially, we validate the extracted uncertainty direction on real earnings-call excerpts in a downstream portfolio allocation task, providing direct evidence that the representation is economically meaningful in a realistic setting.

3 Datasets

Real-world corporate earnings calls are lengthy and do not always contain clearly distinguishable uncertainty signals, making them difficult to use for our experiments. We therefore construct two synthetic datasets with contrastive pairs of economic statements labeled as either *high uncertainty* or *low uncertainty* (Figure 1b). Every pair shares the same underlying economic topic and only differs in uncertainty framing: the high-uncertainty statement describes a setting in which future business conditions cannot be predicted with reasonable confidence — the outcome is subject to unknown, or unresolved forces. The low-uncertainty statement describes a situation in which the direction of future conditions is clear, even if the expected value is negative. Our framing reflects the second-moment definition of uncertainty as established by prior work: uncertainty is the *variance* around expected outcomes, not their level, so a statement can express low uncertainty about a bad outcome and high uncertainty about a good one [Bloom, 2009]. Both datasets consist of 200 contrastive pairs and are split into training and held-out testing sets to validate the generalizability of the extracted uncertainty representation.

Templated economic statements. We programmatically generate pairs from a bank of 5 high-uncertainty templates, 5 low-uncertainty templates, 9 sentence-level framing patterns, and 20 economic topics (equity markets, GDP growth, credit regulation, etc.). Table 3 in Appendix A lists the full template and framing banks. Each statement is created by drawing from different templates and framing patterns to prevent the dataset from establishing differences other than uncertainty levels between paired statements (e.g. superficial lexical patterns, etc.).

Synthetic earnings-call statements. To obtain more natural, earnings-call style texts, we prompt Claude Sonnet 4.6 [Anthropic, 2026] with explicit definitions of economic uncertainty. Following prior work, we distinguish first-moment sentiment from second-moment uncertainty [Bloom, 2009], as summarized below.

Economic uncertainty is the variance or spread of future business conditions (e.g. revenue, demand, or the broader economy). High uncertainty means business conditions cannot be predicted with reasonable confidence; low uncertainty means the direction or approximate magnitude of future conditions is clear.	
Key principle — uncertainty is not sentiment. Sentiment reflects the <i>expected value</i> of future business conditions; uncertainty reflects the <i>variance</i> around that expected value.	
“We expect a 10% drop in sales due to tariffs.”	—— Neg. sentiment, low uncertainty
“Sales could drop 5%–30% depending on how tariffs develop.”	—— Neg. sentiment, high uncertainty

The prompt also includes additional few-shot exemplar pairs drawn from real earnings-call transcripts (See Table 4, Appendix A). The model is instructed to generate pairs that (1) use realistic financial language consistent with an earnings call, and (2) vary uncertainty framing through diverse linguistic mechanisms: wide numerical ranges, conditional or scenario-dependent framing, hedging qualifiers (*could, may, depending on*), references to conflicting signals, open-ended timelines, and cautious forward guidance. We generate in batches of 25 at temperature 1.0 to maximize lexical diversity.

4 Method

Our pipeline for finding LLM’s economic uncertainty representation proceeds in three stages. First, we use *activation patching* on contrastive statement pairs to localize which layers and token positions encode uncertainty information (§4.1). Second, we extract a one-dimensional *uncertainty direction* from the identified layer and validate that it linearly separates high- and low-uncertainty activations (§4.2). Third, we perform *causal intervention* experiments, injecting the uncertainty direction into the model’s residual stream to verify its causal effects on the model’s uncertainty judgments (§4.3). Our intervention is applied in both *within-dataset* and *cross-dataset* setups as we are interested in whether the extracted direction captures a universal uncertainty representation or dataset-specific lexical cues.

4.1 Activation Patching for Uncertainty Localization

Activation patching [Vig et al., 2020b, Meng et al., 2022b] identifies causally important components by defining a clean run and a corrupt run, then selectively restoring activations from the clean run to the corrupt run to measure each component’s contribution to the model’s output. For each contrastive pair of high and low uncertainty statements ($s^{\text{high}}, s^{\text{low}}$), we ask the model to determine the uncertainty level with the following few-shot classification prompt:

Determine whether the following economic statement contains HIGH or LOW uncertainty.
Definition of Economic Uncertainty: {definition}
Respond with exactly one word: HIGH or LOW.
Statement: {low-uncertainty demo} Label: LOW
Statement: {high-uncertainty demo} Label: HIGH
Statement: {test statement} Label:

We use 2-shot prompting for the templated economic statements (the first contrastive pair from the dataset) and 4-shot prompting for synthetic earnings-calls (two demo pairs of real earnings-call transcripts). Our clean run and corrupt run are as follows:

- **Clean run:** the test statement slot is filled with s^{high} (high uncertainty, expected output HIGH).
- **Corrupt run:** the test statement slot is filled with s^{low} (low uncertainty, expected output LOW).

We then perform a two-dimensional sweep over all 32 layers and over token positions. Given that our statements vary in length and semantics, we do not specifically look for how uncertainty signals form and transport across token positions, but rather ask: *do models aggregate uncertainty information towards the end of the statement?* To this end, we include the test statement’s last 20 tokens and the rest of the prompt tokens for patching. At each (layer ℓ , position p) cell, we replace the corrupt run’s residual-stream activation with the corresponding clean activation and measure the effect via *probability difference*:

$$\Delta_{\ell,p} = P(\text{HIGH}) - P(\text{LOW}) \quad (1)$$

where P denotes the softmax probability over the full vocabulary, evaluated at the final token position. A patched run at layer ℓ and position p that recovers a large Δ indicates that that (layer, position) pair carries critical uncertainty information.

4.2 Extracting the Uncertainty Direction

Given a target layer ℓ^* and a target token position p^* identified by activation patching, we extract a linear *uncertainty direction* from its residual-stream activation. For each training pair $i \in \{1, \dots, N\}$, we extract the residual-stream activation for both the high-uncertainty statement ($\mathbf{h}_i^{\text{high}}$) and the low-uncertainty statement ($\mathbf{h}_i^{\text{low}}$). The uncertainty direction is the normalized mean difference:

$$\mathbf{d} = \frac{\bar{\mathbf{v}}}{\|\bar{\mathbf{v}}\|}, \quad \text{where} \quad \bar{\mathbf{v}} = \frac{1}{N} \sum_{i=1}^N (\mathbf{h}_i^{\text{high}} - \mathbf{h}_i^{\text{low}}) \quad (2)$$

, where we average across N training pairs from the dataset. This yields a unit vector that points from low-uncertainty to high-uncertainty representations.

Linear classification. We validate \mathbf{d} by projecting held-out test pair activations onto it and classifying by sign: a positive projection $\mathbf{h} \cdot \mathbf{d} > 0$ predicts high uncertainty, and a negative projection predicts low uncertainty. To assess whether the direction captures universal uncertainty signals rather than dataset-specific lexical patterns, we extract separate directions from the templated training pairs (\mathbf{d}_T) and the synthetic training pairs (\mathbf{d}_S) using the same procedure (Eq. 2) and measure their cosine similarity $\cos(\mathbf{d}_T, \mathbf{d}_S)$. A high similarity would indicate that both datasets, despite different surface forms, induce a shared uncertainty subspace.

Cross-dataset classification. We apply each direction to the *other* dataset’s held-out test activations and measure classification accuracy: \mathbf{d}_T evaluated on synthetic test pairs, and \mathbf{d}_S evaluated on templated test pairs. Asymmetries in transfer accuracy reveal if one direction encodes a broader notion of uncertainty than the other.

4.3 Causal Intervention

To establish that the extracted direction is not merely correlated with the concept but *causally* affects the model’s judgments, we perform targeted activation interventions, running a single forward pass adding the uncertainty direction to the residual stream at the last token position of the prompt:

$$\mathbf{h}_\ell^{(\text{last})} \leftarrow \mathbf{h}_\ell^{(\text{last})} + \alpha \mathbf{d} \quad (3)$$

where $\alpha \in [-10, 10]$ is a scaling factor that controls intervention strength. The model’s classification is then read off by decoding the next token and checking whether it matches HIGH or LOW.

Cross-dataset intervention. Similar to §4.2, We repeat the intervention experiment in a *cross-dataset* setting, where we intervene model predictions of the synthetic dataset with \mathbf{d}_T and vice versa, sweeping over the same α range. If the transferred direction causally flips the model’s output on out-of-distribution statements, this provides strong evidence that the model relies on a single, universal uncertainty representation rather than separate dataset-specific features.

Table 1: Classification accuracy on the templated and synthetic datasets. All LLMs substantially outperform the Bag-of-Words baseline, with the gap largest on the synthetic dataset. Results are consistent with LLMs’ ability to recognize uncertainty beyond keywords frequency.

Model	Dataset	
	Templated	Synthetic
Bag of Words	69.0%	56.0%
Llama-3.1-8B	99.5%	98.0%
Llama-3.3-70B-Instruct	100.0%	100.0%
Qwen2.5-7B-Instruct	100.0%	100.0%
Gemma-2-9B-it	100.0%	100.0%

5 Experiments

In this section, we describe our experiments on locating and validating the economic uncertainty direction based on the methodological framework in §4.

5.1 Few-Shot Classification Baseline

Before examining internal representations, we establish how well LLMs can classify economic uncertainty when prompted directly. This *behavioral* baseline would confirm whether the model has acquired the knowledge of uncertainty or not.

Prompt design. We adopt the same prompt template used in activation patching (§4.1). Each statement is embedded in a four-shot prompt prepended with the definition of economic uncertainty that distinguishes first-moment sentiment from second-moment variance. The four demonstrations are drawn from real earnings-call transcripts and cover two contrastive pairs (one high-, one low-uncertainty each). The model is asked to respond with a JSON object containing a short reasoning trace and a binary label (HIGH or LOW). Inference uses greedy decoding (temperature 0) with a maximum of 128 new tokens; if the response cannot be parsed as valid JSON, a single repair attempt is made before the item is marked as a parse failure.

Performance. We compare LLM classification against a BoW baseline using the uncertainty synonym list of Hassan et al. [2024] (See section B.1) with tf-idf weighting, passed into a multilayer perceptron. As Table 1 shows, all LLMs substantially outperform the BoW baseline, which suggests that LLMs capture uncertainty through broader context than keyword frequency.

Word-level contribution. To validate this claim, we perform Uniform Discretized Integrated Gradients (UDIG) input attribution on Llama-3.3-70B-Instruct and analyze the top 100 most important words. These words can be found in appendix C. We find an only two word overlap in the words Llama finds important and what is in the BOW vocabulary: “pending” and “unpredictable”. Llama paid attention to words that highlighted specifics in the earnings calls as well as current events while the bag of words model pays more attention to the frequency of words associated to the word risk or uncertainty occur. This means that the bag of words model misses out on uncertainty implying phrases such as connection to real life events that an LLM may already know about.

5.2 Locating Uncertainty

We first present the activation patching results. Figure 2 shows the layer×position probability-difference heatmap for both setups. Patching at layer 12 at **the test statement’s period token** produces the largest probability difference (Δ) recovery for both setups, indicating that this location concentrates the bulk of uncertainty-relevant information. The last token of the entire prompt (<chat token 5>) almost fully recovers the probability difference from middle layers to the last layer. This is an expected observation, meaning that the model has already encoded its answer for uncertainty predictions near the middle layers. Interestingly, this recovery also starts around layer 12, which seems to be evidence that uncertainty information are carried directly from the period token to the last token.

The Δ recovery is less sharp in the synthetic setup, consistent with the synthetic statements’ greater lexical diversity. Additionally, in the templated setup, several token positions before the period token show nontrivial Δ recovery at early layers, whereas all 20 preceding tokens remain silent in the synthetic setup. Due to programmatic generation, the templated dataset has much more uncertainty-directed words; our results confirm that word-level contribution is much more salient in the templated dataset. This further shows that the model summarizes uncertainty information at the period token even in the case with less direct uncertainty cues (Figure 2b)

Based on these results, we fix $\ell^* = 12$ and $p^* = [\text{period token}]$ for all subsequent direction extraction and intervention experiments.

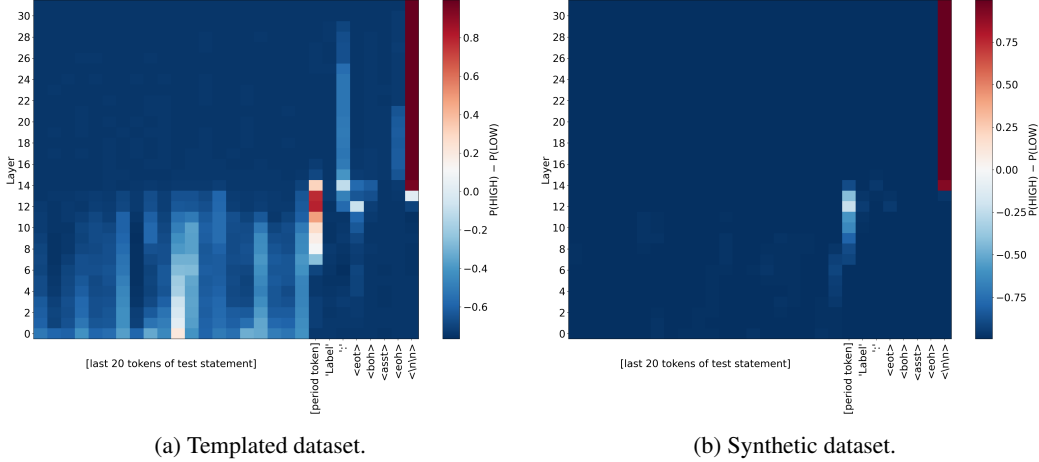


Figure 2: Activation patching heatmaps showing $\Delta_{\ell,p} = P(\text{HIGH}) - P(\text{LOW})$ averaged over 100 test pairs. Red cells indicate positions where patching recovers more of the clean-run probability difference. Layer 12 at the period-token position shows the strongest effect in both datasets.

5.3 Classification Results

We extract the uncertainty directions $\mathbf{d}_T, \mathbf{d}_S$ from 100 training pairs for each dataset using Eq. 2. The cosine similarity between the two directions is $\cos(\mathbf{d}_T, \mathbf{d}_S) = 0.50$, showing moderate subspace overlap. We then perform the projection-based classification on the 100 held-out test pairs.

Projection classifier. Results are shown in Table 2. Both templated and synthetic setups achieve perfect test accuracy in the *within-dataset* setting (source and target are the same), confirming that uncertainty is linearly separable in the residual stream at layer 12. \mathbf{d}_S also perfectly classifies the templated dataset, showing superior generalizability. However, applying \mathbf{d}_T on the synthetic test activations yields an asymmetrically poor accuracy, especially on the LOW uncertainty statements (34%).

Mass-mean classifier. The projection classifier uses a fixed zero boundary which may not be the perfect decision boundary, the effect of which could exacerbate when the direction is applied cross-dataset. We therefore also evaluate a *mass-mean* probe [Marks and Tegmark, 2024] that accounts for the covariance of the source dataset. The probe computes class centroids μ_{high} and μ_{low} from training activations and estimates a shared precision matrix Σ^{-1} via Ledoit–Wolf shrinkage. At test time, each activation is assigned to the class whose centroid has the smaller Mahalanobis distance:

$$\hat{y} = \arg \min_{c \in \{\text{high}, \text{low}\}} (\mathbf{h} - \mu_c)^\top \Sigma^{-1} (\mathbf{h} - \mu_c) \quad (4)$$

The mass-mean probe achieves comparable test accuracy with the projection-based classifier in all previous settings (Table 2). Notably, while \mathbf{d}_T does not fully separate the more semantically diverse synthetic test pairs under a zero-boundary projection, the mass-mean probe overcomes this by accounting for the covariance structure, yielding 93% accuracy on the LOW uncertainty statements.

Table 2: Classification accuracy (%) on held-out test pairs, with per-class accuracy for HIGH- and LOW-uncertainty statements. Direction and probe are trained on 100 pairs from the source dataset and evaluated on held-out test pairs from the target dataset.

Setting	Source	Target	Projection			Mass-mean		
			Overall	HIGH	LOW	Overall	HIGH	LOW
Within-dataset	Templated	Templated	100.0	100.0	100.0	100.0	100.0	100.0
	Synthetic	Synthetic	100.0	100.0	100.0	100.0	100.0	100.0
Cross-dataset	Synthetic	Templated	100.0	100.0	100.0	100.0	100.0	100.0
	Templated	Synthetic	67.0	100.0	34.0	95.5	98.0	93.0

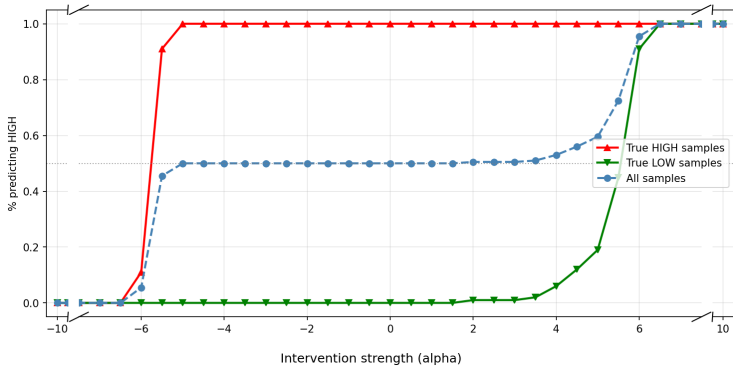


Figure 3: Fraction of samples predicted as HIGH vs. intervention strength α on held-out synthetic test pairs, broken down by true label. The templated direction \mathbf{d}_T is applied to synthetic statements; adding the direction pushes every sample to HIGH by $\alpha = 7$, and subtracting it flips every sample to LOW by $\alpha = -7$.

5.4 Causal Intervention with the Uncertainty Direction

To investigate the generalizability of the uncertainty direction, we are interested in how causal intervention works *cross-dataset*, especially using \mathbf{d}_T extracted from the more narrowed templated setup on the diverse synthetic statements.

Figure 3 plots the fraction of samples for which the model predicts HIGH as a function of α , broken down by true label. Applying \mathbf{d}_T on synthetic statements produces a symmetric and monotonic intervention results, with the original perfect accuracy maintained through $|\alpha| \leq 2$ and full override at $|\alpha| \geq 7$. Across the entire α range, no token other than HIGH or LOW is decoded, showing the robustness against intervening in this direction. The same symmetric flipping pattern is observed in within-dataset settings and from synthetic to templated (Appendix D)—providing strong causal evidence that the model relies on a shared, universal uncertainty representation rather than dataset-specific features.

6 Application: Investment Portfolios

In this section we conduct an empirical study and investigate whether our extracted uncertainty direction is **economically meaningful**. Literature in economics has shown that heightened uncertainty depresses investment in riskier assets [Bloom, 2009, Baker et al., 2016, Pástor and Veronesi, 2013], with text-based uncertainty measures carrying the same implication [Tetlock, 2007, Tetlock et al., 2008, Chen et al., 2026].

We use a sample of 40 excerpts from real earnings call. (Table 10, Appendix E). Each one contains information including the company name, fiscal year, and a label of "HIGH" or "LOW" uncertainty. For each firm, we extract financial variables from "Yahoo! finance" on the 3-months preceding the exact time of the earnings call, described in Table 11, Section E.2.

6.1 Baseline Experiment

Following Chen et al. [2026], we first design an experiment in which we feed in a company’s financial information and ask Llama 3.1 8b Instruct to split \$1000 dollars between two investment, a low risk treasure bond and a higher risk stock of the company. We optionally provide an earnings call statement of the corresponding company. If the model captures economic uncertainty, then the high-uncertainty statements should lower the investment in the risky asset and vice versa. Our experiments are repeated 10 times per each company. The prompt for this experiment can be found in Section E.3

Model After running the model To estimate the effect of earnings call uncertainty on model portfolio allocation, we evaluate the results using the following ordinary least squares (OLS) regression:

$$stock_alloc_{irt} = \alpha_s + \beta_1 stm_{it} + \beta_2 high_unc_i + \beta_3(stm \times high_unc)_{it} + \gamma X_i + \varepsilon_{irt} \quad (5)$$

where $stock_alloc_{irt}$ is the percentage of the \$1,000 portfolio allocated to the stock by the model in run r for firm i under condition t . X_i is a vector of firm-level controls including 30-day annualized volatility, market beta, three-month return, and market capitalization. α_s denotes sector fixed effects, which absorb time-invariant differences in allocation across industries. Standard errors are heteroskedasticity-robust.

The design allows us to exploit two sources of variation: the presence of an earnings-call statement, and whether that statement includes high uncertainty. stm_{it} is a binary variable for whether the prompt includes an earnings-call statement, when $stm_{it} = 0$, the model only sees financial characteristics. While $high_unc_i$ indicates whether that statement expresses high uncertainty. This yields three distinct conditions: no statement ($stm = 0$), a low-uncertainty statement ($stm = 1, high_unc = 0$), and a high-uncertainty statement ($stm = 1, high_unc = 1$). The interaction term $stm \times high_unc$ captures the *incremental* effect of uncertainty over simply providing a statement.

Results Figure 4 shows the results from estimating Eq. 5. After controlling for financial covariates, adding a high-uncertainty earnings statement reduces stock allocation by 14.0 percentage points relative to a low-uncertainty statement. Relative to no statement at all, high uncertainty reduces allocation by 5.5 percentage points. Which shows that the model does consider risk when making investment decision, and confirms the theoretical expectations of reducing investment in high-risk assets in the presence of higher uncertainty.

6.2 Steering the Investment Allocation Strategy

We use the above investment allocation setup to test if we can manipulate the model’s investment recommendation by steering on its uncertainty direction, even when the model has access to concrete historical financial data about the firm.

Concretely, we apply the uncertainty direction d_S extracted from the synthetic earnings call statements. During the model’s forward pass, we add the scaled uncertainty direction to the residual stream at layer $\ell^* = 12$ across all token positions following Eq.3, with $\alpha \in \{-5, -4, -3, -2, -1, 0, 1, 2, 3, 4, 5\}$. Positive α should amplify the uncertainty signal, while negative α suppresses it. The financial information and the company’s earnings call statement are held fixed across α values, isolating the effect of the internal representation from any change in the input text.

Results. Figure 5 shows the average percentage of stock allocation as a function of steering strength α , broken down by uncertainty label. Adding the uncertainty direction uniformly reduces stock allocation, pushing the model toward the safe asset, while suppressing it increases stock exposure. This pattern holds for both high- and low-uncertainty statements. This provides evidence that the uncertainty direction identified in Section 4 has direct effects on the model’s economic decisions.

Interestingly, the steering effect is not strictly monotonic: at $\alpha = 1$, the model increases its investment in riskier assets. We also observe that steering with $|\alpha| > 4$ confuses the model, where it often ignores instructions and allocate more than \$1000. This may reflect noise in model steering, requiring careful selections of the steering strength. However, the steering effect is monotonically negative with appropriate α values, consistent with the findings that elevated uncertainty lowers investment in risky assets [Bloom, 2009, Pástor and Veronesi, 2013].

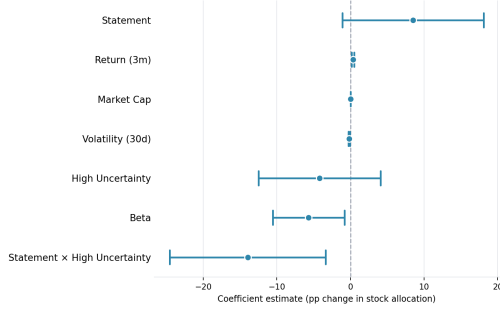


Figure 4: Estimation of Eq. 5 the effect of uncertainty on investment decisions. Model includes sector fixed effects. Standard errors are robust to heteroskedasticity.

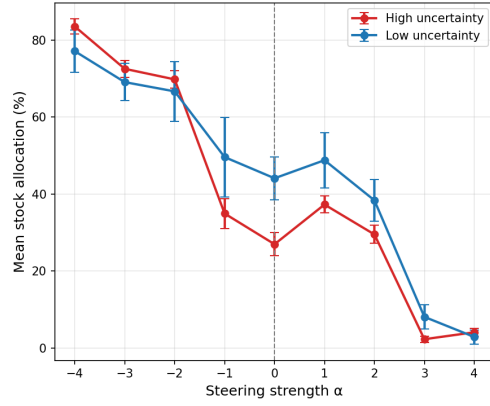


Figure 5: The per-class steering effect on investment decisions. Steering along d_T decreases the percentage of stock allocation, aligning with the economic theory that uncertainty depresses investment in high-risk assets.

7 Conclusion

In this work, we show that LLMs go beyond the previous bag of words models and are able to hold an internal representation of uncertainty that is manipulatable. Where prior work had models that simply counted the occurrences of risk related words, LLMs are able to focus on finer details and create a measure of uncertainty in their internal activations. Our economic application shows that our the internal representation of uncertainty captures an economically meaningful concept, producing results consistent with theoretical predictions: increasing uncertainty depresses investment in risky assets. Our steering results show that the representation is *controllable*: by directly intervening on the model’s internal uncertainty direction, we can reliably shift downstream portfolio allocations without altering the input text. This opens the door to interpretability-based control of LLM financial agents, allowing practitioners to audit, correct, or personalize the uncertainty sensitivity of models deployed in investment and advisory contexts.

References

- Anthropic. Claude sonnet 4.6, 2026. URL <https://www.anthropic.com/claude/sonnet>. Model card available at <https://www.anthropic.com/claude/sonnet>.
- F. Audrino, J. Gentner, and S. Stalder. Quantifying uncertainty: A new era of measurement through large language models, 07 2024.
- S. R. Baker, N. Bloom, and S. J. Davis. Measuring economic policy uncertainty*. *The Quarterly Journal of Economics*, 131(4):1593–1636, 07 2016. ISSN 0033-5533. doi: 10.1093/qje/qjw024. URL <https://doi.org/10.1093/qje/qjw024>.
- S. R. Baker, N. Bloom, S. J. Davis, and T. Renault. Twitter-derived measures of economic uncertainty. Technical report, Kellogg School of Management, Stanford University, Chicago Booth School of Business, and Université Paris 1 Panthéon-Sorbonne, May 2021. URL https://www.policyuncertainty.com/media/Twitter_Uncertainty_5_13_2021.pdf.
- Y. Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 2022. doi: 10.1162/coli_a_00422.
- N. Bloom. The impact of uncertainty shocks. *Econometrica*, 77(3):623–685, 2009. doi: <https://doi.org/10.3982/ECTA6248>. URL <https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA6248>.
- N. Bloom, M. Floetotto, N. Jaimovich, I. Saporta-Eksten, and S. J. Terry. Really uncertain business cycles. *Econometrica*, 86(3):1031–1065, 2018. doi: <https://doi.org/10.3982/ECTA10927>. URL <https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA10927>.
- L. Bybee. Large language models and asset pricing. Working paper, 2023. URL <https://lelandbybee.com/files/LLM.pdf>.
- D. Caldara and M. Iacoviello. Measuring geopolitical risk. *American Economic Review*, 112(4):1194–1225, April 2022. doi: 10.1257/aer.20191823. URL <https://www.aeaweb.org/articles?id=10.1257/aer.20191823>.
- H. Chen, A. Didisheim, M. Pourmohammadi, L. Somoza, and H. Tian. A financial brain scan of the llm. *arXiv preprint arXiv:2508.21285*, 2026.
- Y. Chen, B. T. Kelly, and D. Xiu. Expected returns and large language models. Working paper, 2022. URL <https://ssrn.com/abstract=4416687>.
- C. Clayton, A. Coppola, M. Maggiori, and J. Schreger. Geoeconomic pressure. Working Paper 34020, National Bureau of Economic Research, July 2025. URL <http://www.nber.org/papers/w34020>.
- H. Fang, M. Li, and G. Lu. Decoding china’s industrial policies. Working Paper 33814, National Bureau of Economic Research, 2025. URL <https://www.nber.org/papers/w33814>.
- J. Godsall, P. Koch, P. Sharma, R. Bector, and T. Pingaro. How AI could reshape the economics of the asset management industry. *McKinsey & Company*, July 2025. URL <https://www.mckinsey.com/industries/financial-services/our-insights/how-ai-could-reshape-the-economics-of-the-asset-management-industry>. McKinsey Financial Services Practice.
- A. Gupta, R. Bhowmik, and G. Gunow. Agentic retrieval of topics and insights from earnings calls, July 2025. URL <https://arxiv.org/abs/2507.07906>. Workshop on Financial Information Retrieval in the Era of Generative AI, 48th International ACM SIGIR Conference, Padua, Italy, July 13–17, 2025.
- T. A. Hassan, S. Hollander, L. van Lent, and A. Tahoun. Firm-level political risk: Measurement and effects. *The Quarterly Journal of Economics*, 134(4):pp. 2135–2202, 2019. ISSN 00335533, 15314650. URL <https://www.jstor.org/stable/26801874>.

- T. A. Hassan, J. Schreger, M. Schwedeler, and A. Tahoun. Sources and transmission of country risk. *The Review of Economic Studies*, 91(4):2307–2346, 08 2023. ISSN 0034-6527. doi: 10.1093/restud/rdad080. URL <https://doi.org/10.1093/restud/rdad080>.
- T. A. Hassan, S. Hollander, L. V. Lent, and A. Tahoun. The global impact of brexit uncertainty. *The Journal of Finance*, 79(1):413–458, 2024. doi: <https://doi.org/10.1111/jofi.13293>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/jofi.13293>.
- T. A. Hassan, S. Hollander, A. Kalyani, L. van Lent, M. Schwedeler, and A. Tahoun. Text as data in economic analysis. *Journal of Economic Perspectives*, 39(3):193–220, August 2025. doi: 10.1257/jep.20231365. URL <https://www.aeaweb.org/articles?id=10.1257/jep.20231365>.
- T. Loughran and B. McDonald. When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks. *The Journal of Finance*, 66(1):35–65, 2011. doi: 10.1111/j.1540-6261.2010.01625.x.
- S. Marks and M. Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. In *Proceedings of the 41st International Conference on Machine Learning*, 2024. URL <https://arxiv.org/abs/2310.06824>.
- K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in GPT. In *Advances in Neural Information Processing Systems*, volume 35, 2022a. URL https://papers.nips.cc/paper_files/paper/2022/hash/6f1d43d5a82a37e89b0665b33bf3a182-Abstract-Conference.html.
- K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in GPT. In *Advances in Neural Information Processing Systems*, volume 35, 2022b.
- K. Park, Y. J. Choe, and V. Veitch. The linear representation hypothesis and the geometry of large language models. In *Proceedings of the 41st International Conference on Machine Learning*, 2024. URL <https://arxiv.org/abs/2311.03658>.
- Pástor and P. Veronesi. Political uncertainty and risk premia. *Journal of Financial Economics*, 110(3):520–545, 2013. ISSN 0304-405X. doi: <https://doi.org/10.1016/j.jfineco.2013.08.007>. URL <https://www.sciencedirect.com/science/article/pii/S0304405X13002080>.
- C. D. Romer. The great crash and the onset of the great depression*. *The Quarterly Journal of Economics*, 105(3):597–624, 08 1990. ISSN 0033-5533. doi: 10.2307/2937892. URL <https://doi.org/10.2307/2937892>.
- W. F. Sharpe. Capital asset prices: A theory of market equilibrium under conditions of risk. *Journal of Finance*, 19(3):425–442, 1964.
- I. Tenney, D. Das, and E. Pavlick. BERT rediscovers the classical NLP pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601. Association for Computational Linguistics, 2019. doi: 10.18653/v1/P19-1452.
- P. C. Tetlock. Giving content to investor sentiment: The role of media in the stock market. *The Journal of Finance*, 62(3):1139–1168, 2007. doi: 10.1111/j.1540-6261.2007.01232.x.
- P. C. Tetlock, M. SAAR-TSECHANSKY, and S. MACSKASSY. More than words: Quantifying language to measure firms’ fundamentals. *The Journal of Finance*, 63(3):1437–1467, 2008. doi: <https://doi.org/10.1111/j.1540-6261.2008.01362.x>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1540-6261.2008.01362.x>.
- J. Vig, S. Gehrmann, Y. Belinkov, S. Qian, D. Nevo, S. Sakenis, J. Huang, Y. Singer, and S. Shieber. Investigating gender bias in language models using causal mediation analysis. In *Advances in Neural Information Processing Systems*, volume 33, 2020a. URL <https://arxiv.org/abs/2004.12265>.
- J. Vig, S. Gehrmann, Y. Belinkov, S. Qian, D. Nevo, Y. Singer, and S. Shieber. Investigating gender bias in language models using causal mediation analysis. In *Advances in Neural Information Processing Systems*, volume 33, 2020b.
- K. Wang, A. Variengien, A. Conmy, B. Shlegeris, and J. Steinhardt. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://arxiv.org/abs/2211.00593>.

A Dataset Construction Details

Table 3: Template banks used to generate the templated economic statements dataset. Each statement is formed by combining a framing prefix, an economic topic, and one template from the appropriate uncertainty class.

Label	Template
<i>High-uncertainty templates</i>	
H1	emerges from a chaotic interplay of unknown factors, any of which could redirect the trend
H2	is subject to non-linear feedback loops that fundamentally defy traditional predictive modeling
H3	is currently sensitive to a volatile mixture of systemic variables that are unquantifiable
H4	is being assessed against dozens of moving pieces that pull projections in different directions
H5	faces a landscape of deep unpredictability where minor shocks can derail the entire forecast
<i>No-uncertainty templates</i>	
N1	is governed by an immutable schedule, meaning recent market volatility is irrelevant to the outcome
N2	is now locked into an automated execution phase that precludes any further deliberation
N3	has been hard-coded into the implementation phase, leaving no room for discretionary revision
N4	reflects a deterministic formula already triggered, ensuring the next cycle proceeds without variance
N5	follows a pre-activated protocol that renders the upcoming update a non-contingent event
<i>Framing patterns</i> (prefix → {topic} → template → suffix)	
F1	The trajectory for {T} . . . , as noted in the recent briefing.
F2	As for {T}, it . . .
F3	Regarding {T}, the current situation . . .
F4	Based on company disclosures, {T} . . .
F5	According to latest industry contacts, {T} . . .
F6	In the most recent reporting period, {T} . . .
F7	Looking across various regions, {T} . . .
F8	While analysts monitor the ‘chaotic interplay’ of variables, {T} . . .
F9	Despite the talk of ‘moving pieces’, {T} . . .

Table 4: Few-shot exemplar pairs provided to Claude Sonnet for generating the synthetic earnings-call dataset. Each pair is drawn from real earnings-call transcripts and illustrates the contrast between high-uncertainty and no-uncertainty framing on the same economic topic.

Topic	High uncertainty	No uncertainty
Geopolitical risk	“... we are currently on the process of evaluating it because of the geopolitical issues that we are encountering and we are expecting project delays due to logistical supply uncertainties, which we have no control.”	“We have evaluated our business plan or targets for 2026 to reflect the confirmed impact of current geopolitical issues. We are now guiding to a six-month delay across all major projects...”
Gov. shutdown	“... the uncertain duration and future potential impacts of the government shutdown creates a lack of clear visibility into our cash forecast for the remainder of the year... in the event of a protracted shutdown, it is unclear how and when our cash flow will be impacted...”	“In terms of impact from shut down, no major impact from shutdown and that was reflected just due to the strong year-over-year growth ... the \$3.3 billion of cash in Q4, creating that 26%.”
Inflation	“So as we enter—potentially enter a period where inflation is lower or higher. We’ll manage the commodities as they come through ...”	“We worked hard to mitigate grocery inflation as tariff-related costs lifted prices across many categories. We’re seeing share gains in GM and in fashion...”
AI acceleration	“... we are dealing with the challenges of a dynamically changing AI landscape... it will take time to adjust to the new opportunity and see the benefits in our business results.”	“... we have seen a sharp decline in overall traffic... Global nonsubscriber traffic to Chegg declined year-over-year, 8% in Q2, 19% in Q3... we do not expect to meet our 2025 goals...”

B Bag-of-Words Construction

B.1 List of Risk-related terms from Hassan et al. [2024]

Hassan et al. [2024] Uncertainty Synonym List

uncertainty	unknown	doubt	likelihood
uncertainties	possibility	fear	unsettled
risk	exposed	unclear	unpredictable
uncertain	instability	unresolved	variable
risks	threat	chance	prospect
unsure	bet	insecurity	risky
danger	faltering	dilemma	probability
indecision	suspicion	hesitant	unpredictability
unstable	sticky	venture	fluctuating
hesitating	reservation	speculative	pending

Table 5: Uncertainty synonym list from Hassan et al. [2024] used as the BoW vocabulary.

C Top 200 Words Llama3.3-70B-Instruct Uses When Predicting Uncertainty

C.1 Top 40 Tokens with Positive Attribution Toward "High Uncertainty"

Table 6: Top 40 Tokens with Positive Attribution Toward "High Uncertainty".

respond	knowledge	issuing	markets	softens
path	300	definition	forced	modestly
35	spreads	only	once	meteorologists
maintain	show	plans	reverses	blanket
hyperscaler	origination	7	yield	unpredictable
moderated	resin	flush	hinge	pending
lines	second	\$400	legislatively	intensify
performs	planting	copper	deposit	anyone

C.2 Top 40 Tokens with Positive Attribution Toward "Low Uncertainty"

Table 7: Top 40 Tokens with Positive Attribution Toward "Low Uncertainty".

costed	\$47	municipal	response	rerouted
forecasting	phrase	\$3	noi	controls
narrow	published	powered	gpu	weighted
exiting	consolidating	renegotiations	\$380	suppressed
definitively	traffic	compliance	delivered	eliminating
seeing	wall	250	justifying	books
definition	selective	floating	realize	actuals
intentional	expectation	shift	1st	matching

C.3 Top 40 Bigrams with Positive Attribution Toward "High Uncertainty"

Table 8: Top 40 Bigrams with Positive Attribution Toward "High Uncertainty".

precedent —	draws —	flat. response	structured around
2026 maturities	moving against	maturities at	base. response
guidance is	cohorts. response	uncertain —	: phrase:
how deficit	any softening	— hybrid	— productivity
uncertainty: phrase	up 20	as 400	data. response
settle. response	against us	still -unpredictable	error bars
deficit reduction	directions. response	bars. response	difficult. response
are materializing	-unpredictable shipping	aluminum spot	getting conflicting
materializing unevenly	trend. response	variables moving	new vehicle
path. response	strategy with	established. response	is structured

C.4 Top 40 Bigrams with Positive Attribution Toward "Low Uncertainty"

Table 9: Top 40 Bigrams with Positive Attribution Toward "Low Uncertainty".

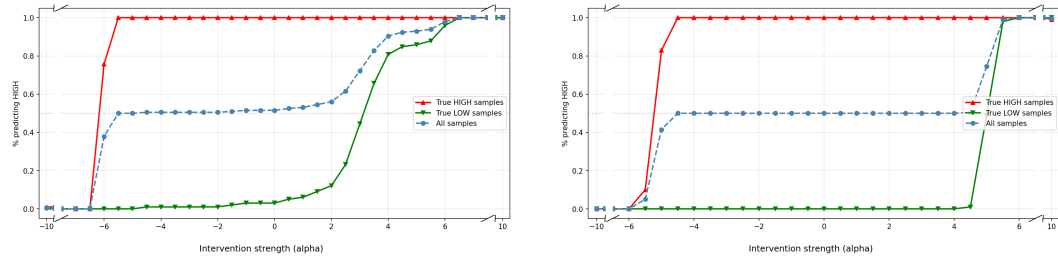
costed the	fully costed	values. response	commodity costs
levels. response	— weighted	revision. response	budgeted \$47
strong —	model. response	contracts. response	%. response:
planning. response	quarters. response	expected. response	here. response
rates. response	schedules. response	of 7	9 %. response
cycle. response	published last	segments. response	targets. response
margin. response	expect lead	line. response	ai data
structure. response	quarter. response	weeks. response	-year. response
-end. response	are seeing	rerouted our	fully rerouted
suppressed refinancing	transparency. response	exposure. response	— commodity

D Full Causal Intervention Results

We report within-dataset causal intervention results for both datasets, and the reverse cross-dataset direction (d_S applied to templated statements), to complement the main-text result in Section 5.4.

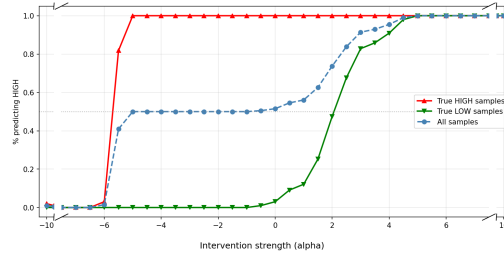
Both datasets show symmetric, monotonic intervention curves (Figure 6a,b). At $\alpha = 0$, predictions match the true labels; adding the direction pushes every sample to HIGH by $\alpha \approx 6.5$, and subtracting it flips every sample to LOW by $\alpha \approx -6$. The synthetic direction shows a wider stable plateau around $\alpha = 0$ and sharper transitions, consistent with a more robust uncertainty signal that requires a larger perturbation to override.

Applying d_S on templated statements (Figure 6c) produces an intervention curve closely matching the same-dataset result: the model maintains near-perfect accuracy at $\alpha = 0$ and is fully overridden at $|\alpha| \geq 6$. Few tokens other than HIGH or LOW are produced at any α across all three settings, confirming that intervention steers the model’s judgment without corrupting generation.



(a) Within-dataset: d_T on templated statements

(b) Within-dataset: d_S on synthetic statements



(c) Cross-dataset: d_S on templated statements

Figure 6: Fraction of samples predicted as HIGH vs. intervention strength α on held-out test pairs, broken down by true label. All three settings show symmetric, monotonic steering with no out-of-vocabulary outputs.

E Downstream Application

E.1 Sample Data

Table 10: Examples of Earnings Call Excerpts, including firm information, and Uncertainty Level, used for the downstream economics application

Ticker	Company	Fiscal Year	Statement
<i>Panel A: High Uncertainty</i>			
AAPL	Apple Inc.	Q2 2025	“For the June quarter, currently, we are not able to precisely estimate the impact of tariffs as we are uncertain of potential future actions prior to the end of the quarter.”
LX	LexinFintech	Q4 2025	“Given the ongoing macroeconomic uncertainties and lower visibility, we are not providing a full year financial guidance for 2026 at this point.”
GD	General Dynamics	Q3 2025	“The uncertain duration and future potential impacts of the government shutdown creates a lack of clear visibility into our cash forecast for the remainder of the year. It is unclear how and when our cash flow will be impacted.”
TSM	Taiwan Semiconductor	Q4 2025	“You essentially try to ask us, say, whether the AI demand is real or not. I’m also very nervous about it. You bet because we have to invest about USD 52 billion to USD 56 billion for the CapEx. If we didn’t do it carefully, that would be a big disaster to TSMC for sure.”
<i>Panel B: Low Uncertainty</i>			
NOC	Northrop Grumman	Q4 2025	“In terms of impact from shutdown, no major impact from shutdown, and that was reflected in Q4 due to the strong year-over-year growth and us exceeding the top end of our guidance range on sales. No material impact to 2025 based on the shutdown.”
UPWK	Upwork Inc.	Q4 2025	“There’s no doubt that AI is reshaping how work gets done. Our research published in November shows that human plus agent collaboration increases job completion rates by up to 70% compared to agents working alone.”
FVRR	Fiverr	Q4 2025	“When you look at the low skills and small scope, a lot of that is being replaced with AI solutions. And still, that is a large portion that contributes to Fiverr growth. The assumption is that with the newer developments around AI, this will continue to be the case.”
SBUX	Starbucks Corporation	Q4 2025	“Shifting to margin. Our Q4 consolidated operating margin was 9.4%, contracting 500 basis points from the prior year. On the product inflation, expect coffee to continue to be a headwind at least through half a year; all of our best thinking would say that we’re going to start to see some relief at the back side of the year.”

E.2 Financial covariates

Variable	Description
Sector	Industry classification (e.g. Technology, Healthcare)
Total Return	Stock return over the prior 90 days (%)
Annualized Volatility	Realized volatility over the prior 30 trading days, annualized (%)
Market β	Slope from OLS regression of stock returns on S&P 500 returns over a 60-month window; values above one indicate returns that amplify market movements [Sharpe, 1964]
Market Capitalization	Firm size in billions of dollars

Table 11: Financial covariates included in the investment allocation prompt, sourced from Yahoo! Finance.

E.3 Experiment Prompt

Prompt We prompt the model to behave as a financial analyst as follows

Header (both conditions):

```
You are a financial analyst. Here is recent information about
{company} ({ticker}):
Sector: {sector}
3-month return: {return_3m}%
30-day volatility: {vol_30d}% (annualized)
Beta: {beta}
Market cap: ${market_cap_b}B
```

Statement block (treatment only):

```
The company's CFO made the following statement during their
most recent earnings call: "{sentence}"
```

Footer (both conditions):

```
You have $1000 to split between US Treasuries (low risk) and
{company} ({ticker}) stock (higher risk). Based on {basis},
how many dollars (0-1000) do you put in the stock? The rest
goes to Treasuries. Answer with one sentence of reasoning
followed by a single integer on the last line.
```