

FinCall-Surprise: A Large Scale Multi-modal Benchmark for Earning Surprise Prediction

Anonymous ACL submission

Abstract

Predicting corporate earnings surprises is a profitable yet challenging task, as accurate forecasts can inform significant investment decisions. However, progress in this domain has been constrained by a reliance on expensive, proprietary, and text-only data, limiting the development of advanced models. To address this gap, we introduce **FinCall-Surprise** (Financial Conference Call for Earning Surprise Prediction), the first large-scale, open-source, and multi-modal dataset for earnings surprise prediction. Comprising 2,688 unique corporate conference calls from 2019 to 2021, our dataset features word-to-word conference call textual transcripts, full audio recordings, and corresponding presentation slides. We establish a comprehensive benchmark by evaluating 26 state-of-the-art unimodal and multi-modal LLMs. Our findings reveal that (1) while many models achieve high accuracy, this performance is often an illusion caused by significant class imbalance in the real-world data. (2) Some specialized financial models demonstrate unexpected weaknesses in instruction-following and language generation. (3) Although incorporating audio and visual modalities provides some performance gains, current models still struggle to leverage these signals effectively. These results highlight critical limitations in the financial reasoning capabilities of existing LLMs and establish a challenging new baseline for future research. The FinCall-Surprise dataset is available at <https://anonymous.4open.science/r/FinCall-Surprise-F212>.

1 Introduction

In the finance domain, a conference call, often referred to as an earnings call, serves as a critical communication channel between the management of a public company and its stakeholders, including analysts, investors, and the media (Kimbrough, 2005). During these calls, executives present the

firm’s financial results, discuss recent performance, and provide forward-looking guidance. Analysts and investors closely analyze this information to predict their expectations of the company’s earnings per share (EPS) (Patell, 1976). An *earnings surprise* occurs when the company’s reported EPS differs from market expectations (Latane and Jones, 1979). Historically, a positive earnings surprise, where actual earnings exceed the reported EPS, often correlates with a subsequent increase in the company’s stock price, while a negative surprise typically leads to a decline. Therefore, the ability to accurately predict an earnings surprise ahead of the official announcement is a significant challenge and an opportunity for investors seeking to inform their trading decisions (Skinner and Sloan, 2002).

Given the potential financial rewards, researchers and investors have long sought to systematically analyze conference calls to predict earnings surprises. This pursuit began with traditional machine learning models (Dhar and Chou, 2001; Doyle et al., 2006; Bissessur and Veenman, 2016), and has recently shifted attention toward large language models (LLMs) (Zhang et al., 2025; Zhang and He, 2025; Zhu et al., 2025; Liang and Carasco Kind, 2025). However, while the performance of these models is often impressive, their capabilities are inherently constrained by a fundamental limitation, which is the unimodal and text-only nature of the datasets they are trained on. In reality, human analysts do not just analyze what executives say, but also how they say it. The vocal tone and rhythm of a speaker can reveal confidence or uncertainty that is lost in a plain transcript. Similarly, visual information, such as presentation slides, provides an essential context that is often referenced during the discussion.

To address this critical gap, we introduce **FinCall-Surprise**, a novel, large-scale, multi-modal dataset specifically designed for the task of earnings surprise prediction. The dataset is com-

posed of 2,688 unique corporate conference calls spanning from 2019 to 2021, with 919 calls from 2019, 704 from 2020, and 1,065 from 2021. Most importantly, each call in the dataset is complete, containing three synchronized modalities: 1) the word-to-word textual transcript, 2) the full audio recording of the call, and 3) the corresponding presentation slides referenced by executives. With the introduction of this dataset, our primary contributions are as follows:

- We present FinCall-Surprise, the first large-scale, multi-modality dataset focused on real-world corporate earnings calls, providing a rich resource for developing and evaluating earning surprise prediction models.
- We establish a comprehensive benchmark by evaluating more than 20 state-of-the-art unimodal and multi-modal LLMs.
- Our benchmarking findings reveal that current models struggle to effectively leverage the multi-modal signals for this task, establishing a challenging baseline and highlighting the need for more sophisticated architectures.

2 Related Work

2.1 Earning Surprise Prediction Models

The use of AI in predicting earnings surprises has a long history that began with traditional statistical and econometric models. In early foundational research, researchers like [Numbers \(1968\)](#) relied on simple models such as linear regressions. A pivotal shift occurred when [Ou and Penman \(1989\)](#) used a broader set of financial statement data with a logit model to predict the direction of earnings, foreshadowing the machine learning era. From approximately 2000 to 2015, classic machine learning models like SVMs, Random Forests, and Gradient Boosting took center stage ([Dhar and Chou, 2001](#); [Skinner and Sloan, 2002](#); [Zolotoy, 2012](#)). Their key advantage was the ability to analyze hundreds of financial variables simultaneously, capturing complex, non-linear relationships that eluded older methods. The period after 2015 saw the rise of deep learning, particularly RNNs and LSTMs, which introduced unstructured text as a new and powerful data source ([Vargas et al., 2017](#); [Meursault et al., 2023](#); [Zhu, 2020](#); [Sawhney et al., 2020](#)).

The contemporary approach to earnings surprise prediction is dominated by LLMs, which represent

a revolutionary leap in analytical capability. Researchers began feeding models text from conference call transcripts and financial reports to extract sentiment and forward-looking statements ([Koval et al., 2023](#)). Since LLMs can possess a deep and nuanced understanding of financial jargon and context, they can dissect conference call transcripts, press releases, and news with unprecedented accuracy ([Araci, 2019](#); [Kim et al., 2024](#); [Lee et al., 2025](#)). However, due to the lack of high-quality multimodal conference call datasets, models for earnings surprise prediction have not yet reached the multimodal level.

2.2 Earning Surprise Prediction Datasets

Due to the high value and sensitivity of corporate financial data, most resources are stored within major commercial platforms such as WRDS, Bloomberg, and Refinitiv, which require paid subscriptions or API access. Despite this barrier, there has been steady progress in developing open-source financial datasets, each typically targeting narrower and more specialized tasks in the finance domain to advance AI research in this field. For instance, datasets such as ([Chen et al., 2021, 2024](#); [Zhu et al., 2021](#); [Lai et al., 2024](#); [Reddy et al., 2024](#); [Yuan et al., 2024](#)) focus on question answering over financial reports. Other datasets have been developed for stock market prediction ([Dong et al., 2024](#); [Rao, 2021](#); [Patel, 2021](#); [Qin and Yang, 2019](#)), fraud detection ([Feng et al., 2023](#)), sentiment analysis ([Borhani, 2024](#); [Cortis et al., 2017](#)), and misinformation classification ([Rangapur et al., 2025](#)). More recent efforts have extended financial datasets into the multimodal space, motivated by both the limitations of single-modality data and the growing capabilities of LLMs to understand information beyond text ([Li et al., 2020](#); [Shu et al., 2025](#); [Luo et al., 2025](#)).

Earnings surprise prediction has recently gained increasing attention in the context of LLMs, with several studies proposing new methods to improve performance on this task ([Zhu et al., 2025](#); [Zhang et al., 2025](#)). However, relatively few works have focused on dataset construction. Early efforts such as [Koval et al. \(2023, 2024\)](#) introduced conference call transcript datasets for this purpose, but these resources are not fully open source, limiting their accessibility and impact. Consequently, many current studies continue to depend on expensive commercial data providers, often accessing conference call content through paid APIs ([Heater et al., 2025](#);

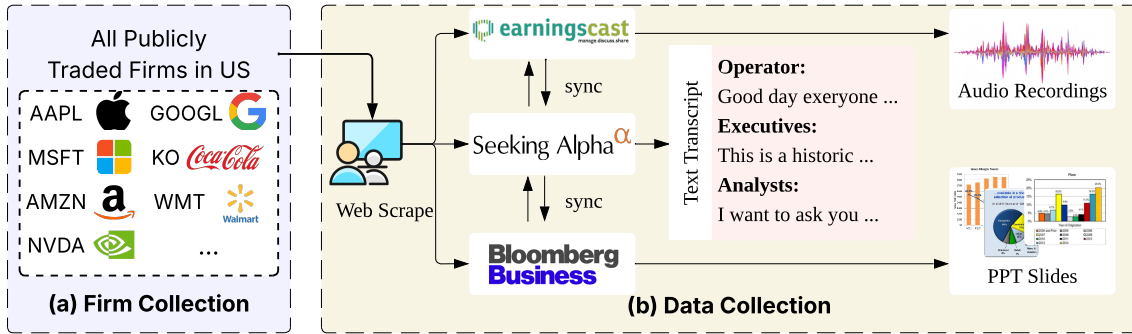


Figure 1: Overview of our data construction pipeline, which consists of two stages: **(a) Firm Collection (left)**: We select large, publicly traded US firms based on market capitalization ($> \$1B$) and daily trading volume ($> \$50M$). **(b) Data Collection (right)**: For each firm, we gather and synchronize three modalities for each quarterly earnings call: textual transcripts from Seeking Alpha, audio recordings from EarningsCast, and presentation slides from sources like Bloomberg Business.

Alsabah, 2025; Zhang and He, 2025).

Our proposed FinCall-Surprise dataset addresses this challenge by introducing the first fully open-source, multi-modal conference call dataset specifically designed for earnings surprise prediction, covering the period from 2019 to 2021. Each data in our dataset provides synchronized modalities: the complete text transcript of the conference call, the full audio recording, and the accompanying presentation slides referenced by executives. By releasing this dataset, we aim to remove the reliance on commercial APIs and enable the research community to explore earnings surprise prediction at a multimodal level, opening new directions for both financial NLP and multimodal learning.

3 FinCall-Surprise Construction

In this section, we detail the creation process for the FinCall-Surprise dataset.

3.1 Data Construction Pipeline

The construction of FinCall-Surprise follows a comprehensive two-stage pipeline designed to collect, synchronize, and annotate data from multiple sources. Our objective is to build a robust benchmark that integrates the textual, auditory, and visual dimensions of corporate conference calls. As illustrated in Figure 1, the first stage involves assembling the universe of all publicly traded firms in the United States with substantial size and liquidity. To focus on economically significant firms, we restrict the sample to companies with market capitalizations exceeding \$1 billion and average daily trading volumes above \$50 million. This initial firm collection yielded an initial pool of more

than 4,000 unique companies.

In the second stage, we acquire the three distinct modalities for each firm’s earnings conference calls. For textual transcripts, we systematically web scrape content from established financial platforms, primarily Seeking Alpha¹. These transcripts are well structured, containing both speaker identifiers (operators, executives and analysts) and the corresponding speech content. For the auditory component, we web scrape the associated audio recordings from EarningsCast². Finally, we construct a dataset of corporate presentation slides by collecting files from multiple sources, including Bloomberg News and company websites. Most publicly traded firms host one conference call per fiscal quarter. To ensure temporal consistency across modalities, we align all data sources by quarterly reporting periods. In addition, we use the conference call titles to cross-check and verify that the transcript, audio recording, and presentation slides correspond to the same event, ensuring that our multi-modal dataset is accurately matched at the event level.

3.2 Earning Surprise Label Preliminaries

Before classifying each conference call label as positive or negative, we first collected the reported Earnings Per Share (EPS) and the analyst consensus EPS forecasts from the IBES database. Following Latane and Jones (1979), we measure the earnings surprise (ES) using the Standardized Unexpected Earnings (SUE), defined as the difference between the reported EPS and the analyst consen-

¹<https://seekingalpha.com/>

²<https://earningscast.com/>

Year	Transcript (words)			Slide (pages)			Audio (sec)			Label (Percentage)		Total
	Mini	Max	Avg.	Mini	Max	Avg.	Mini	Max	Avg.	Positive	Negative	
2019	905	57,485	8,634.66	1	156	26.05	1096.07	8870.97	3642.66	0.79	0.21	919
2020	865	46,623	8,925.10	1	170	27.80	1404.06	7585.31	3859.25	0.86	0.14	704
2021	752	45,632	8,646.77	1	113	27.46	1197.95	11938.90	3720.29	0.89	0.11	1065

Table 1: Data statistics of the FinCall-Surprise. The dataset spans three years (2019–2021), with each conference call containing three synchronized modalities: text transcripts, presentation slides, and audio recordings. We report the minimum, maximum, and average values for each modality. Transcript length is measured in words, slides in pages, and audio in seconds. For labels, we present the percentage distribution of positive and negative classes. The last column shows how much data we have for each year.

249 sus estimate, scaled by the standard deviation of
250 analyst forecasts. The consensus estimate is calcu-
251 lated as the mean of the most recent valid analyst
252 forecasts issued within one month after the confer-
253 ence call, allowing analysts to revise their expect-
254 ations based on the call content and recent finan-
255 cial disclosures. This design provides a forward-
256 looking measure of market expectations and yields
257 a more realistic, yet challenging, prediction task.
258 The average time span between the input transcript
259 and the target earnings event is about three months,
260 further highlighting the difficulty of the task.

$$261 \quad ES = \frac{EPS_{\text{reported}} - \text{Avg}(EPS_{\text{estimated}})}{\text{Std}(EPS_{\text{estimated}})} \quad (1)$$

$$262 \quad y = \begin{cases} 0, & ES \leq -\delta \\ 1, & ES \geq \delta \end{cases} \quad (2)$$

263 We convert the continuous earnings surprise (ES)
264 into a binary classification task by assigning a label
265 of Positive (+1) when $ES > \delta$ and Negative (0)
266 when $ES < -\delta$, where $\delta = 0.50$. This thresh-
267 old follows prior studies on standardized unex-
268 pected earnings (Eli Bartov, 1992) and price mo-
269 mentum (Luo et al., 2022), which classify earnings
270 surprises as large when $|SUE| \geq 0.5$. The cho-
271 sen cutoff balances sample size with event signifi-
272 cance. Observations with immaterial surprises (i.e.,
273 $ES \in [-0.50, 0.50]$) are excluded, as these near-
274 zero values typically elicit weak market responses
275 and may reflect earnings management. Although
276 ES is continuous, market reactions are largely bi-
277 nary, responding more to the direction than to the
278 magnitude of the surprise. We therefore focus on
279 material surprises that are more likely to influence
280 investor behavior and asset pricing.
281

282 3.3 Data Analysis

283 A detailed statistical analysis of our dataset is pre-
284 sented in Table 1. The statistics highlight the sub-
285 stantial scale and complexity of the data across

286 all three modalities. On average, the textual tran-
287 scriptions for each conference call contain approxi-
288 mately 8,600 words. The accompanying presen-
289 tation slides average 27 pages in length. Finally,
290 the audio recordings have an average duration of
291 approximately 3,700 seconds. Each data point in
292 our dataset is annotated with a binary label indicat-
293 ing either a positive or negative earnings surprise.
294 As shown in the last column, the distribution of
295 labels is highly imbalanced, with the proportion
296 of positive labels significantly exceeding that of
297 negative labels. This imbalance is unavoidable,
298 since our dataset is constructed from real-world
299 financial data, where positive earnings surprises
300 occur more frequently. The content of the confer-
301 ence calls is consistently structured around three
302 primary speaker roles: Operators, Executives, and
303 Analysts. An example is provided in Figure 2. The
304 Operator typically opens and closes the call. The
305 Executives, including chief members such as the
306 CEO and CFO, present the company’s financial
307 results and strategic outlook. Finally, the Analysts,
308 representing investment firms and financial insti-
309 tutions, pose questions to the executives to gain
310 deeper insights into the company’s performance.

311 4 Benchmark Setting

312 4.1 Baseline LLMs

313 To establish a comprehensive performance bench-
314 mark for FinCall-Surprise, we selected a diverse
315 set of 26 state-of-the-art, open- and close-source
316 models. For a structured analysis, these models
317 were categorized into four distinct groups.

318 The first group consists of general purpose uni-
319 modal models, which process text-only inputs, and
320 includes GPT-oss-20B (OpenAI, 2025), Qwen-2.5-
321 14B (Team, 2024), Mistral-7B (Jiang et al., 2023),
322 LLaMa-3.2-3B (Grattafiori et al., 2024), Gemma-
323 3-1B (Team, 2025a). The second group con-

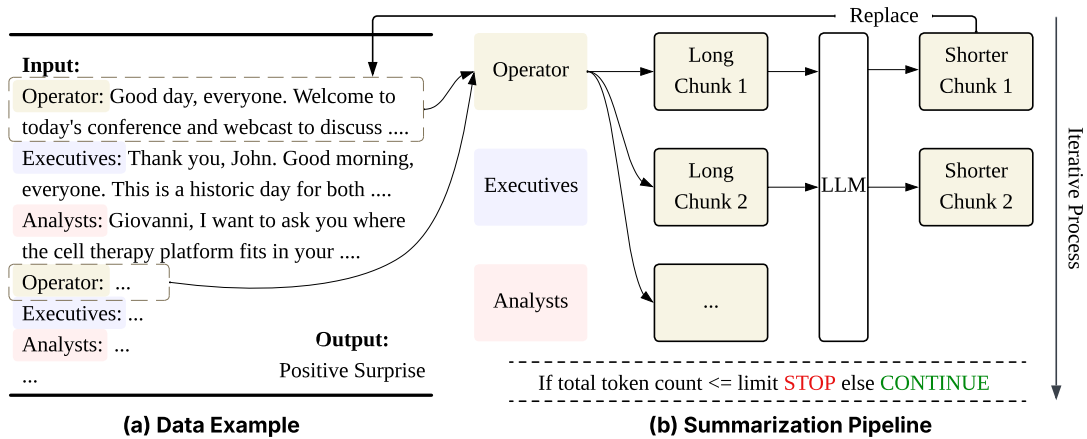


Figure 2: Illustration of our summarization pipeline. (a) A data example showing a conference call transcript with speaker turns (Operator, Executives, Analysts) and its corresponding earning surprise label. (b) The summarization pipeline, where transcripts are grouped by speaker and iteratively summarized by the LLM until the total token count falls below the predefined limit.

324 consists of finance-specialized unimodal models fine-tuned for the financial domain, including Finance-LLaMa3-8B (Cheng et al., 2024), Finance-LLaMa-8B (tarun7r, 2025), Finance-R1 (KhazarAI, 2024), LLaMa-RAG (Liu et al., 2024). The third category features Vision-Language Models (VLMs), designed to process both images and text. For this group, we evaluated GPT-5, GPT-5-mini, GPT-5-nano, GPT-4.1, Mistral-3.1-24B (Jiang et al., 2023), Gemma-3-12B (Team, 2025a), LLaMa-3.2-11B (Grattafiori et al., 2024), Sa2VA-8B (Yuan et al., 2025), Qwen-2.5-VL-7B (Team, 2025b), Qwen-2-VL-7B (Wang et al., 2024), Cosmos-7B (Azzolini et al., 2025), and LLaVa-1.6-7B (Liu et al., 2023). Finally, we assessed audio-language models, which handle audio and text inputs, including Voxtral-24B (Liu et al., 2025), DeSTA-2.5-Audio-8B (Lu et al., 2025), Qwen-2.5-Omni-7B (Xu et al., 2025), Qwen-2-Audio-7B (Chu et al., 2024), Gemma-3n-E4B (Team, 2025a).

344 4.2 Benchmark Input Design

345 To systematically evaluate the performance of different models, we designed a customized input and prompting strategy for each modality.

348 **Text-Only Modality.** To ensure a fair evaluation across all baseline models, we standardized the input length to accommodate the smallest context window of 32K tokens. We set the threshold to 31K, because we need to reserve some space for the instruction prompt. As approximately 20% of the transcripts in our dataset exceed this limit, we employed a targeted summarization strategy

356 for these longer texts. Transcripts already under the 31K token threshold were left unmodified. To minimize information loss, we adopted a conservative, iterative summarization process as shown in Figure 2. First, each transcript is segmented into chunks based on speaker type (e.g., Operator, Executive, Analyst). We then incrementally summarize the chunks that are least critical to the financial context, beginning with those from the “Operator”. These sections typically contain procedural dialogue, such as call introductions and closings, rather than useful financial discussion. Every time a chunk is summarized, we recalculate the total token count, and the iterative process stops as soon as the transcript length falls below the 31K token threshold. If summarizing all “Operator” chunks is insufficient, the iterative process continues with chunks from “Executives” and then “Analysts”. In practice, summarization usually stops midway through this process, leaving the majority of the transcript untouched and thus preserving as much original information as possible. Finally, the processed text was formatted into a single concatenated sequence for model input: “{Instruction} {Text Transcript} {Answer Format}”. This structure provides the model with the complete verbal context of the conference call.

383 For this summarization task, we utilized the BART-Large-CNN model (Lewis et al., 2019), because it is lightweight, reliable, and does not require a specific instruction prompt to function effectively. For token number checking, we used “tiktoken” python library. While summarization was a necessary step for our experiments, we will release the original,

390 full-length transcripts to support future research
391 that may leverage models with larger context win-
392 dows or more advanced summarization techniques.

393 **Image-Text Modality.** Evaluating multi-modal
394 models that accept text and images required a pre-
395 processing step for the presentation slides, which
396 were originally in PDF format. As most models do
397 not support PDF inputs, we converted each page of
398 a presentation into a separate image. To manage
399 the image input limits common to these models, we
400 developed a random sampling strategy. For each
401 presentation, we randomly selected three images
402 from the slide range, excluding the first and last
403 pages, which typically contain non-substantive con-
404 tent like title or closing. In cases where a presenta-
405 tion had three or fewer content slides, all available
406 slides were used. The final input for this modality
407 was structured as: “{Images} {Instruction} {Text
408 Transcript} {Answer Format}”. Note that although
409 our benchmark experiment uses only three images,
410 we still release the full PDF in our dataset.

411 **Audio-Text Modality.** For models capable of
412 processing audio, we intentionally omitted the text
413 transcript from the input. Since the audio record-
414 ing contains the same verbal information as the
415 transcript, this approach isolates the model’s abil-
416 ity to comprehend and reason based on auditory
417 signals alone, such as speaker tone and inflection.
418 The input was therefore constructed as: “{Audio}
419 {Instruction} {Answer Format}”.

420 Across all modalities, the prompt structure was
421 designed to test both the model’s reasoning capa-
422 bilities and its ability to follow instructions. We
423 required the final answer to be placed in a specific,
424 designated location within the output, which en-
425 abled reliable parsing for our automated evaluation
426 pipeline. Detailed examples of the prompts for
427 each modality are provided in the Appendix A.

428 4.3 Benchmark Metrics

429 Since our benchmark involves only two possible
430 outputs, accuracy naturally serves as a primary
431 evaluation metric. However, accuracy alone can
432 be misleading, particularly under class imbalance.
433 To provide a more balanced evaluation, we report
434 macro-averaged precision, macro-averaged recall,
435 and macro-averaged F1-score. These metrics treat
436 both positive and negative classes equally, regard-
437 less of their frequency in the dataset. *By incorpo-*
438 *rating macro-averaged metrics alongside accuracy,*
439 *we account for the dataset’s class imbalance and*

440 *obtain a more reliable assessment of model perfor-*
441 *mance across both positive and negative cases.*

442 4.4 Implementation Details

443 All experiments were conducted on a single
444 NVIDIA A100 SXM4 GPU with 80GB of mem-
445 ory. For all baseline models, we maintained their
446 official repository configurations and loaded them
447 with bfloat16 precision to optimize computational
448 efficiency. Beyond this precision setting, no other
449 modifications were made to the models’ default
450 parameters, ensuring a fair and reproducible com-
451 parison across all baseline evaluations.

452 5 Benchmark Results

453 We have listed all 26 models’ result in Table 2 in-
454 cluding 5 general purpose text modality models, 4
455 financial finetuned models, 12 image-text modal-
456 ities models, 5 audio-text modalities models. We
457 have the following observations.

458 5.1 High Performance Illusion

459 As shown in Table 2, many models appear to
460 achieve strong results when evaluated only by ac-
461 curacy. In fact, 15 out of 26 models surpass 70%
462 accuracy. However, this impression of high perfor-
463 mance is misleading and illusional. Our dataset is
464 heavily imbalanced toward positive samples, and
465 when we examine precision, recall, and F1 score,
466 most of the models fall below 55%. This indi-
467 cates that most models are biased toward predict-
468 ing the majority (positive) class and fail to gener-
469 alize to negative cases. A plausible interpretation
470 for this biased behavior is that the models have
471 simply learned the statistical distribution of the fi-
472 nancial domain, where positive outcomes are more
473 frequent. This aligns with our own findings during
474 the dataset construction phase, which confirmed
475 that real-world financial data often exhibits signifi-
476 cant imbalance.

477 A closer look also reveals that model accuracy
478 tends to increase on the 2021 subset compared
479 to earlier years. For instance, Qwen-2.5 in the
480 General-Text category achieves accuracies of 0.78
481 and 0.77 in 2019 and 2020, respectively, but rises
482 to 0.84 in 2021. However, this gain does not re-
483 flect genuine model improvement. Rather, it stems
484 from the fact that the 2021 data is more imbalanced,
485 with 89% positive and only 11% negative samples.
486 When models default to predicting the majority
487 class, their accuracy naturally rises under such con-
488 ditions. This improvement is therefore illusory, as

	2019				2020				2021				Overall			
	Acc	Prec	Rec	F1	Acc	Prec	Rec	F1	Acc	Prec	Rec	F1	Acc	Prec	Rec	F1
<i>General-Text</i>																
GPT-oss (20B)	0.73	0.57	0.59	0.58	0.63	0.49	0.49	0.48	0.75	0.52	0.52	0.52	0.71	0.53	0.54	0.53
Qwen-2.5 (14B)	0.78	0.51	0.50	0.49	0.77	0.52	0.52	0.52	0.84	0.48	0.49	0.48	0.80	0.51	0.50	0.50
Mistral (7B)	0.73	0.50	0.50	0.50	0.68	0.52	0.53	0.51	0.80	0.51	0.51	0.51	0.74	0.51	0.51	0.51
LLaMa-3.2 (3B)	0.74	0.50	0.50	0.50	0.75	0.51	0.51	0.51	0.83	0.52	0.52	0.52	0.78	0.51	0.51	0.51
Gemma-3 (1B)	0.70	0.43	0.45	0.44	0.73	0.45	0.45	0.45	0.77	0.45	0.45	0.45	0.74	0.44	0.45	0.45
<i>Finance-FT</i>																
Fin-LLaMa3 (8B)	0.10	0.17	0.08	0.09	0.12	0.23	0.07	0.10	0.14	0.27	0.08	0.12	0.12	0.22	0.08	0.11
Fin-LLaMa (8B)	0.70	0.46	0.47	0.46	0.77	0.47	0.47	0.47	0.77	0.46	0.46	0.46	0.74	0.47	0.47	0.47
Finance-R1 (1.7B)	0.49	0.37	0.32	0.34	0.48	0.42	0.32	0.35	0.57	0.43	0.35	0.38	0.52	0.40	0.33	0.36
LLaMa-RAG (8B)	0.08	0.14	0.05	0.07	0.10	0.21	0.06	0.09	0.12	0.26	0.08	0.11	0.10	0.20	0.06	0.09
<i>Image-Text</i>																
GPT-5	0.57	0.54	0.56	0.51	0.52	0.54	0.59	0.46	0.71	0.54	0.58	0.53	0.62	0.55	0.58	0.52
GPT-5-mini	0.60	0.53	0.54	0.52	0.53	0.52	0.54	0.45	0.73	0.53	0.55	0.52	0.64	0.53	0.55	0.51
GPT-5-nano	0.61	0.53	0.54	0.52	0.48	0.52	0.54	0.43	0.73	0.52	0.53	0.52	0.62	0.52	0.54	0.50
GPT-4.1	0.66	0.53	0.54	0.53	0.57	0.53	0.56	0.48	0.79	0.54	0.55	0.54	0.69	0.53	0.55	0.53
Mistral-3.1 (24B)	0.72	0.49	0.50	0.49	0.72	0.51	0.52	0.51	0.86	0.52	0.51	0.50	0.77	0.51	0.51	0.51
Gemma-3 (12B)	0.66	0.50	0.50	0.50	0.65	0.52	0.54	0.50	0.83	0.55	0.54	0.54	0.73	0.52	0.53	0.52
LLaMa-3.2 (11B)	0.68	0.52	0.54	0.50	0.49	0.52	0.54	0.43	0.52	0.51	0.51	0.47	0.58	0.52	0.54	0.48
Sa2VA (8B)	0.72	0.43	0.46	0.44	0.77	0.45	0.46	0.45	0.80	0.44	0.45	0.45	0.76	0.44	0.46	0.45
Qwen-2.5 (7B)	0.78	0.52	0.51	0.48	0.82	0.47	0.49	0.47	0.88	0.45	0.49	0.47	0.82	0.49	0.50	0.48
Qwen-2 (7B)	0.70	0.43	0.45	0.44	0.71	0.49	0.48	0.48	0.81	0.49	0.49	0.49	0.74	0.47	0.47	0.47
Cosmos (7B)	0.75	0.44	0.48	0.45	0.79	0.49	0.49	0.49	0.86	0.47	0.49	0.48	0.81	0.47	0.49	0.47
LLaVa-1.6 (7B)	0.68	0.51	0.51	0.51	0.66	0.49	0.48	0.47	0.84	0.54	0.53	0.53	0.74	0.52	0.52	0.52
<i>Audio-Text</i>																
Voxtral (24B)	0.71	0.51	0.51	0.51	0.69	0.51	0.52	0.50	0.80	0.53	0.53	0.53	0.74	0.52	0.52	0.52
DeSTA-2.5 (8B)	0.54	0.50	0.49	0.47	0.52	0.50	0.49	0.43	0.55	0.50	0.49	0.43	0.54	0.50	0.49	0.45
Qwen-2.5 (7B)	0.78	0.46	0.50	0.45	0.84	0.53	0.51	0.50	0.88	0.51	0.50	0.48	0.84	0.51	0.50	0.47
Qwen-2 (7B)	0.62	0.43	0.42	0.43	0.69	0.46	0.44	0.45	0.67	0.45	0.42	0.43	0.66	0.44	0.43	0.43
Gemma-3n (8B)	0.35	0.53	0.53	0.35	0.32	0.50	0.50	0.31	0.33	0.48	0.46	0.31	0.34	0.50	0.51	0.33

Table 2: Model comparison on the FinCall-Surprise benchmark. The table is divided into four sections: General-Text refers to unimodal LLMs trained for general purposes; Finance-FT refers to unimodal LLMs finetuned on financial datasets; Image-Text refers to vision-language models that accept both images and text; and Audio-Text refers to audio-language models that accept both audio and text. For each metric and year, the highest score within each category is highlighted in **bold**. (Prec, Rec, and F1 represent Macro Average Precision, Macro Average Recall, and Macro Average F1-score, respectively.)

confirmed by its precision, recall, and F1 scores. For example, although Qwen-2.5 achieves 0.84 accuracy in 2021, its precision, recall, and F1 scores show no corresponding increase.

5.2 Balanced Predictions and Robustness of Closed-Source Models

Surprisingly, the four closed-source models evaluated, GPT-5, GPT-5-mini, GPT-5-nano, and GPT-4.1, achieved lower overall accuracy scores in the Image-Text category compared to smaller, open-source models. For instance, the highest accuracy among these closed-source models was 0.69, considerably lower than the 0.82 achieved by Qwen-2.5. However, we argue that this accuracy metric is misleading (discussed in section 5.1). A deeper analysis reveals that the GPT family, particularly GPT-5, consistently achieved the highest precision and recall scores nearly every year. This

discrepancy arises from the models’ prediction behavior. Manual inspection of the outputs showed that unlike many open-source models which predominantly default to a “Positive” label, the closed-source models generated a much more balanced distribution of “Positive” and “Negative” predictions. This avoidance of majority-class bias is what lowers their raw accuracy score. We interpret this as a positive sign, as it suggests the GPT models possess a more robust financial reasoning capability, allowing them to avoid the majority-class bias inherent in the dataset.

5.3 Failure Analysis of Finetuned Financial Models

While most models demonstrated reasonable performance on accuracy, a subset of models in the Finance-FT category performed exceptionally poorly. Notably, Fin-LLaMa3 and LLaMa-RAG

525 achieved overall accuracy scores of only 12% and
526 10%, respectively. A qualitative analysis of their
527 outputs revealed three primary failure reasons: re-
528 sponse bias, poor instruction following, and de-
529 graded language generation.

530 First, unlike most other models, these two mod-
531 els show a strong tendency to predict “negative” la-
532 bels. The reason for this bias is unclear, but it may
533 stem from the nature of their finetuning data. Sec-
534 ond, the models exhibit difficulty following instruc-
535 tions. To evaluate model performance systemati-
536 cally, our benchmark requires answers in a specific
537 format: `Result = [[POSITIVE or NEGATIVE]]`.
538 However, the finetuned models often failed to com-
539 pply, producing outputs in inconsistent formats. Al-
540 though we attempted to accommodate these varia-
541 tions, some responses could not be parsed reliably.
542 We treated such cases as incorrect, since the ability
543 to follow task instructions should be considered
544 an essential component of performance. Third,
545 in several instances the models failed to generate
546 grammatically correct or meaningful sentences. We
547 suspect that both the instruction-following failures
548 and the degradation in basic language generation
549 stem from over-finetuning on highly specialized
550 financial datasets. In pursuit of domain-specific
551 performance, finetuning may inadvertently over-
552 penalize general language capabilities, diminishing
553 the model’s broader usability. These findings high-
554 light that when developing task-specialized models,
555 it is important to ensure that this does not come at
556 the expense of fundamental language competence
557 and instruction-following ability.

558 **5.4 Do Image and Audio Inputs Improve** 559 **Model Performance?**

560 As shown in Table 2, model performance varies sig-
561 nificantly across categories. Interestingly, we ob-
562 serve both modest improvements and notable degra-
563 dations in accuracy when additional modalities are
564 introduced. For instance, within the Qwen-2.5 fam-
565 ily, the General-Text model with 14B parameters
566 achieves an overall accuracy of 0.80. When image
567 data is incorporated, the 7B Image-Text model im-
568 proves to 0.82. With audio data, the 7B Audio-Text
569 model increases the accuracy to 0.84. While these
570 gains may appear modest, the fact that smaller mul-
571 timodal models outperform their larger text-only
572 counterpart suggests that visual and auditory inputs
573 can provide complementary signals beyond the raw
574 conference call transcripts.

575 However, in other cases, the addition of new

576 modalities leads to performance degradation. For
577 example, the text-only LLaMa-3.2 model with 3B
578 parameters achieves a respectable accuracy of 0.78,
579 yet its larger Image-Text variant with 11B param-
580 eters performs significantly worse, dropping to 0.58.
581 Audio-Text models also exhibit highly inconsistent
582 performance. While Voxtral achieves 0.74 accu-
583 racy and Qwen-2.5 achieves 0.84 accuracy, other
584 models like DeSTA-2.5 and Gemma-3n reach only
585 0.54 and 0.34, respectively. These negative results
586 suggest that, for many state-of-the-art models, sup-
587plementary signals from presentation slides and au-
588dio may introduce noise that cannot be effectively
589 understood with textual content, ultimately hinder-
590 ing predictive capabilities. More broadly, current
591 multimodal models often lack the robustness and
592 consistency that text-only LLMs demonstrate when
593 applied to single-modal financial data.

594 Taken together, our results highlight both the po-
595 tential and the limitations of multimodal learning
596 in this domain. On the one hand, certain models
597 clearly benefit from the additional information that
598 pure text alone cannot provide. On the other hand,
599 many existing multimodal models cannot interpret
600 complex financial signals across modalities. We
601 anticipate that future research leveraging our multi-
602 modal benchmark will unlock deeper insights and
603 enable models to capture information that plain
604 transcripts alone cannot convey.

605 **6 Conclusions**

606 In this work, we introduced FinCall-Surprise, the
607 first large-scale, open-source, multimodal bench-
608 mark for earnings surprise prediction, compris-
609 ing synchronized word-level transcripts, full audio
610 recordings, and presentation slides from more than
611 2,600 corporate conference calls between 2019 and
612 2021. Our evaluation of 26 state-of-the-art models
613 highlights both the potential and the limitations of
614 current approaches. We find that many open-source
615 models achieve deceptively high accuracy due to
616 class imbalance, while closed-source models gener-
617 ally provide more balanced predictions. Financial
618 models fine-tuned on narrow datasets often lose ba-
619 sic instruction-following and language generation
620 abilities, and existing multimodal models continue
621 to struggle with effectively integrating image and
622 audio in the financial domain. We release FinCall-
623 Surprise to reduce reliance on proprietary data and
624 provide an open, challenging benchmark for ad-
625 vancing models in earnings surprise prediction.

626 Limitations

627 A potential limitation of FinCall-Surprise lies in the
628 natural class imbalance between positive and nega-
629 tive earnings surprises. Because positive surprises
630 occur more frequently in real-world financial mar-
631 kets, our dataset inevitably reflects this biased distri-
632 bution. While this imbalance makes the benchmark
633 realistic, it also poses challenges for model evalua-
634 tion. Future research that employs FinCall-Surprise
635 for training could consider techniques such as up-
636 sampling, downsampling, or other modern data-
637 balancing strategies to mitigate this limitation and
638 better assess model robustness across both positive
639 and negative cases.

640 References

641 Khaled Alsabab. 2025. Love me do: Twitter likes and
642 earnings surprise. *Journal of Behavioral Finance*,
643 26(3):283–302.

644 Dogu Araci. 2019. Finbert: Financial sentiment analy-
645 sis with pre-trained language models. *arXiv preprint*
646 *arXiv:1908.10063*.

647 Alisson Azzolini, Junjie Bai, Hannah Brandon, Jiaxin
648 Cao, Prithvijit Chattopadhyay, Huayu Chen, Jinju
649 Chu, Yin Cui, Jenna Diamond, Yifan Ding, and 1
650 others. 2025. Cosmos-reason1: From physical com-
651 mon sense to embodied reasoning. *arXiv preprint*
652 *arXiv:2503.15558*.

653 Sanjay W Bissessur and David Veenman. 2016. Analyst
654 information precision and small earnings surprises.
655 *Review of Accounting Studies*, 21(4):1327–1360.

656 Taha Borhani. 2024. [Twitter financial news sentiment](#)
657 [dataset](#). Kaggle. Accessed: 2025-09-13.

658 Jian Chen, Peilin Zhou, Yining Hua, Yingxin Loh,
659 Kehui Chen, Ziyuan Li, Bing Zhu, and Junwei
660 Liang. 2024. Fintextqa: A dataset for long-
661 form financial question answering. *arXiv preprint*
662 *arXiv:2405.09980*.

663 Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena
664 Shah, Iana Borova, Dylan Langdon, Reema Moussa,
665 Matt Beane, Ting-Hao Huang, Bryan Routledge,
666 and 1 others. 2021. Finqa: A dataset of numerical
667 reasoning over financial data. *arXiv preprint*
668 *arXiv:2109.00122*.

669 Daixuan Cheng, Yuxian Gu, Shaohan Huang, Junyu Bi,
670 Minlie Huang, and Furu Wei. 2024. [Instruction pre-](#)
671 [training: Language models are supervised multitask](#)
672 [learners](#). In *Proceedings of the 2024 Conference on*
673 *Empirical Methods in Natural Language Processing*,
674 pages 2529–2550, Miami, Florida, USA. Association
675 for Computational Linguistics.

Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei,
Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng
He, Junyang Lin, Chang Zhou, and Jingren Zhou.
2024. Qwen2-audio technical report. *arXiv preprint*
arXiv:2407.10759.

Keith Cortis, André Freitas, Tobias Daudert, Manuela
Huerlimann, Manel Zarrouk, Siegfried Handschuh,
and Brian Davis. 2017. [SemEval-2017 task 5: Fine-](#)
[grained sentiment analysis on financial microblogs](#)
[and news](#). In *Proceedings of the 11th International*
Workshop on Semantic Evaluation (SemEval-2017),
pages 519–535, Vancouver, Canada. Association for
Computational Linguistics.

Vasant Dhar and Dashin Chou. 2001. A comparison of
nonlinear methods for predicting earnings surprises
and returns. *IEEE Transactions on Neural networks*,
12(4):907–921.

Zihan Dong, Xinyu Fan, and Zhiyuan Peng. 2024. Fn-
spid: A comprehensive financial news dataset in time
series. In *Proceedings of the 30th ACM SIGKDD*
Conference on Knowledge Discovery and Data Min-
ing, pages 4918–4927.

Jeffrey T Doyle, Russell J Lundholm, and Mark T Soli-
man. 2006. The extreme future stock returns follow-
ing *i/b/e/s* earnings surprises. *Journal of Accounting*
Research, 44(5):849–887.

Eli Bartov. 1992. [Patterns in Unexpected Earnings as an](#)
[Explanation for Post-Announcement Drift](#). *The Ac-*
counting Review, 67(3):610–622. Publisher: Ameri-
can Accounting Association.

Duanyu Feng, Yongfu Dai, Jimin Huang, Yifang Zhang,
Qianqian Xie, Weiguang Han, Zhengyu Chen, Ale-
jandro Lopez-Lira, and Hao Wang. 2023. Empow-
ering many, biasing a few: Generalist credit scor-
ing through large language models. *arXiv preprint*
arXiv:2310.00566.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
Abhinav Pandey, Abhishek Kadian, Ahmad Al-
Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,
Alex Vaughan, and 1 others. 2024. The llama 3 herd
of models. *arXiv preprint arXiv:2407.21783*.

John C Heater, Ye Liu, Qin Tan, and Frank Zhang. 2025.
Winning is not enough: Changing landscapes of earn-
ings surprises and the market reaction. *Contempo-*
rary Accounting Research.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-
sch, Chris Bamford, and Devendra Singh Chaplot.
2023. Diego de las casas. *Florian Bressand, Gianna*
Lengyel, Guillaume Lample, Lucile Saulnier, L elio
Renard Lavaud, Marie-Anne Lachaux, Pierre Stock,
Teven Le Scao, Thibaut Lavril, Thomas Wang, Timo-
th ee Lacroix, and William El Sayed, pages 50–72.

KhazarAI. 2024. Personal-finance-r1.
[https://huggingface.co/khazarai/](https://huggingface.co/khazarai/Personal-Finance-R1)
[Personal-Finance-R1](#). Accessed: 2025-09-
14.

732	Alex Kim, Maximilian Muhn, and Valeri Nikolaev.	Wanlong Liu, Junying Chen, Ke Ji, Li Zhou, Wenyu	784
733	2024. Financial statement analysis with large lan-	Chen, and Benyou Wang. 2024. Rag-instruct: Boost-	785
734	guage models. <i>arXiv preprint arXiv:2407.17866</i> .	ing llms with diverse retrieval-augmented instruc-	786
		tions. <i>Preprint</i> , arXiv:2501.00353.	787
735	Michael D Kimbrough. 2005. The effect of conference	Ke-Han Lu, Zhehuai Chen, Szu-Wei Fu, Chao-	788
736	calls on analyst and market underreaction to earnings	Han Huck Yang, Sung-Feng Huang, Chih-Kai Yang,	789
737	announcements. <i>The Accounting Review</i> , 80(1):189–	Chee-En Yu, Chun-Wei Chen, Wei-Chih Chen,	790
738	219.	Chien-yu Huang, and 1 others. 2025. Desta2.5-audio:	791
		Toward general-purpose large audio language model	792
739	Ross Koval, Nicholas Andrews, and Xifeng Yan. 2023.	with self-generated cross-modal alignment. <i>arXiv</i>	793
740	Forecasting earnings surprises from conference call	<i>preprint arXiv:2507.02768</i> .	794
741	transcripts. In <i>Findings of the association for compu-</i>		
742	<i>tational linguistics: ACL 2023</i> , pages 8197–8209.	Cheng Luo, Enrichetta Ravina, Marco Sammon, and	795
		Luis M. Viceira. 2022. Retail investors’ contrarian	796
743	Ross Koval, Nicholas Andrews, and Xifeng Yan. 2024.	behavior around news, attention, and the momentum	797
744	Financial forecasting from textual and tabular time	effect . Technical report, Social Science Research	798
745	series. In <i>Findings of the Association for Computa-</i>	Network. Posted: 5 Apr 2020; Last revised: 26 May	799
746	<i>tional Linguistics: EMNLP 2024</i> , pages 8289–8300.	2022.	800
747	Viet Dac Lai, Michael Krumdick, Charles Lovering,	Junyu Luo, Zhizhuo Kou, Liming Yang, Xiao Luo,	801
748	Varshini Reddy, Craig Schmidt, and Chris Tanner.	Jinsheng Huang, Zhiping Xiao, Jingshu Peng,	802
749	2024. Sec-qa: A systematic evaluation corpus for	Chengzhong Liu, Jiaming Ji, Xuanzhe Liu, and 1 oth-	803
750	financial qa. <i>arXiv preprint arXiv:2406.14394</i> .	ers. 2025. Finmme: Benchmark dataset for financial	804
		multi-modal reasoning evaluation. <i>arXiv preprint</i>	805
		<i>arXiv:2505.24714</i> .	806
751	Henry A Latane and Charles P Jones. 1979. Standard-	Vitaly Meursault, Pierre Jinghong Liang, Bryan R.	807
752	ized unexpected earnings–1971-77. <i>The Journal of</i>	Routledge, and Madeline Marco Scanlon. 2023.	808
753	<i>Finance</i> , 34(3):717–724.	Pead.txt: Post-earnings-announcement drift using	809
		text . <i>Journal of Financial and Quantitative Analysis</i> ,	810
754	Jean Lee, Nicholas Stevens, and Soyeon Caren Han.	58(6):2299–2326.	811
755	2025. Large language models in finance (finllms).		
756	<i>Neural Computing and Applications</i> , pages 1–15.	Income Numbers. 1968. An empirical evaluation of	812
		accounting. <i>Journal of Accounting Research</i> .	813
757	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	OpenAI. 2025. gpt-oss-120b & gpt-oss-20b model card .	814
758	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	<i>Preprint</i> , arXiv:2508.10925.	815
759	Veselin Stoyanov, and Luke Zettlemoyer. 2019.		
760	BART: denoising sequence-to-sequence pre-training	Jane A Ou and Stephen H Penman. 1989. Financial	816
761	for natural language generation, translation, and com-	statement analysis and the prediction of stock returns.	817
762	prehension . <i>CoRR</i> , abs/1910.13461.	<i>Journal of accounting and economics</i> , 11(4):295–	818
		329.	819
763	Jiazheng Li, Linyi Yang, Barry Smyth, and Ruihai Dong.	Varpit Patel. 2021. Google stock data . Kaggle. Ac-	820
764	2020. Maec: A multimodal aligned earnings confer-	cessed: 2025-09-13.	821
765	ence call dataset for financial risk prediction. In <i>Pro-</i>		
766	<i>ceedings of the 29th ACM International Conference</i>	James M Patell. 1976. Corporate forecasts of earnings	822
767	<i>on Information & Knowledge Management</i> , pages	per share and stock price behavior: Empirical test.	823
768	3063–3070.	<i>Journal of accounting research</i> , pages 246–276.	824
769	Qingwen Liang and Matias Carrasco Kind. 2025. How	Yu Qin and Yi Yang. 2019. What you say and how you	825
770	does managers’ willingness to disclose affect anal-	say it matters: Predicting stock volatility using verbal	826
771	ysts’ earning forecasts—a measurement by llms.	and vocal cues . In <i>Proceedings of the 57th Annual</i>	827
772	<i>Available at SSRN 5199752</i> .	<i>Meeting of the Association for Computational Lin-</i>	828
		<i>guistics</i> , pages 390–401, Florence, Italy. Association	829
773	Alexander H. Liu, Andy Ehrenberg, Andy Lo, Clément	for Computational Linguistics.	830
774	Denoix, Corentin Barreau, Guillaume Lample, Jean-		
775	Malo Delignon, Khyathi Raghavi Chandu, Patrick	Aman Rangapur, Haoran Wang, Ling Jian, and Kai Shu.	831
776	von Platen, Pavankumar Reddy Muddireddy, Sanchit	2025. Fin-fact: A benchmark dataset for multimodal	832
777	Gandhi, Soham Ghosh, Srijan Mishra, Thomas Fou-	financial fact-checking and explanation generation .	833
778	bert, Abhinav Rastogi, Adam Yang, Albert Q. Jiang,	In <i>Companion Proceedings of the ACM on Web Con-</i>	834
779	Alexandre Sablayrolles, Amélie Héliou, and 87 oth-	<i>ference 2025</i> , WWW ’25, page 785–788, New York,	835
780	ers. 2025. Voxtral . <i>Preprint</i> , arXiv:2507.13264.	NY, USA. Association for Computing Machinery.	836
781	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae	Rohan Rao. 2021. Nifty 50 stock market data (2000 -	837
782	Lee. 2023. Improved baselines with visual instruc-	2021) . Kaggle. Accessed: 2025-09-13.	838
783	tion tuning . <i>Preprint</i> , arXiv:2310.03744.		

839	Varshini Reddy, Rik Koncel-Kedziorski, Viet Dac Lai,	Ziqiang Yuan, Kaiyuan Wang, Shoutai Zhu, Ye Yuan,	892
840	Michael Krumdick, Charles Lovering, and Chris Tan-	Jingya Zhou, Yanlin Zhu, and Wenqi Wei. 2024. Fin-	893
841	ner. 2024. Docfinqa: A long-context financial reason-	llms: A framework for financial reasoning dataset	894
842	ing dataset. <i>arXiv preprint arXiv:2401.06915</i> .	generation with large language models. <i>IEEE Trans-</i>	895
		<i>actions on Big Data</i> .	896
843	Ramit Sawhney, Piyush Khanna, Arshiya Aggarwal,	Cong Zhang and Zhenzhi He. 2025. Cross-sectional	897
844	Taru Jain, Puneet Mathur, and Rajiv Ratn Shah. 2020.	spillovers of earnings surprises and asset price anom-	898
845	VoLTAGE: Volatility forecasting via text audio fusion	lies. Available at SSRN 5415255.	899
846	with graph convolution networks for earnings calls .		
847	In <i>Proceedings of the 2020 Conference on Empirical</i>	Huopu Zhang, Yanguang Liu, and Mengnan Du. 2025.	900
848	<i>Methods in Natural Language Processing (EMNLP)</i> ,	Sae-fire: Enhancing earnings surprise predictions	901
849	pages 8001–8013, Online. Association for Computa-	through sparse autoencoder feature selection. <i>arXiv</i>	902
850	tional Linguistics.	<i>preprint arXiv:2505.14420</i> .	903
851	Dong Shu, Haoyang Yuan, Yuchen Wang, Yanguang	Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao	904
852	Liu, Huopu Zhang, Haiyan Zhao, and Mengnan Du.	Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and	905
853	2025. Finchart-bench: Benchmarking financial chart	Tat-Seng Chua. 2021. Tat-qa: A question answering	906
854	comprehension in vision-language models. <i>arXiv</i>	benchmark on a hybrid of tabular and textual content	907
855	<i>preprint arXiv:2507.14823</i> .	in finance. <i>arXiv preprint arXiv:2105.07624</i> .	908
856	Douglas J Skinner and Richard G Sloan. 2002. Earn-	Yongqiong Zhu. 2020. Stock price prediction using the	909
857	ings surprises, growth expectations, and stock returns	rnn model . <i>Journal of Physics: Conference Series</i> ,	910
858	or don't let an earnings torpedo sink your portfolio.	1650:032103.	911
859	<i>Review of accounting studies</i> , 7(2):289–312.		
860	tarun7r. 2025. tarun7r/finance-llama-8b: A llama 3.1 8b	Yu Zhu, Xiao Liu, and Olivia R Liu Sheng. 2025. Post-	912
861	model fine-tuned on josephflowers/finance-instruct-	earnings-announcement drift prediction: Leveraging	913
862	500k. https://huggingface.co/tarun7r/	postevent investor responses with multitask learning.	914
863	Finance-Llama-8B .	<i>Information Systems Research</i> .	915
864	Gemma Team. 2025a. Gemma 3 .	Leon Zolotoy. 2012. Earnings surprise implicit in	916
865	Qwen Team. 2024. Qwen2.5: A party of foundation	stock prices: which earnings forecasting models are	917
866	models .	investors using and what determines their choice?	918
867	Qwen Team. 2025b. Qwen2.5-vl .	<i>Journal of Business Finance & Accounting</i> , 39(9-	919
868	Manuel R Vargas, Beatriz SLP De Lima, and Alexan-	10):1161–1179.	920
869	dre G Evsukoff. 2017. Deep learning for stock mar-		
870	ket prediction from financial news articles. In <i>2017</i>		
871	<i>IEEE international conference on computational in-</i>		
872	<i>telligence and virtual environments for measurement</i>		
873	<i>systems and applications (CIVEMSA)</i> , pages 60–65.		
874	IEEE.		
875	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-		
876	hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin		
877	Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei		
878	Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang		
879	Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-		
880	vl: Enhancing vision-language model's perception		
881	of the world at any resolution. <i>arXiv preprint</i>		
882	<i>arXiv:2409.12191</i> .		
883	Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting		
884	He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan,		
885	Kai Dang, and 1 others. 2025. Qwen2. 5-omni tech-		
886	nical report. <i>arXiv preprint arXiv:2503.20215</i> .		
887	Haobo Yuan, Xiangtai Li, Tao Zhang, Zilong Huang		
888	Huang, Shilin Xu, Shunping Ji, Yunhai Tong, Lu Qi,		
889	Jiashi Feng, and Ming-Hsuan Yang. 2025. Sa2va:		
890	Marrying sam2 with llava for dense grounded under-		
891	standing of images and videos. <i>arXiv preprint</i> .		

A Prompt Used in Our Experiments

A.1 Prompt Used in Text-Only Modality

Text-only Modality

You are an expert equity research analyst. Your task is to read the following earnings-conference call transcript and decide whether next quarter's earnings per share (EPS) will **beat** or **miss** consensus estimates.

Instructions

- **Label only:** respond with exactly one of:
 - POSITIVE (indicating an expected positive EPS surprise)
 - NEGATIVE (indicating an expected negative EPS surprise)
- **Brief rationale** (1–2 sentences) explaining the key driver
- **No other text**

Definitions

- **POSITIVE:** indications of stronger-than-expected earnings or revenue growth, upbeat guidance, new partnerships, successful product launches, cost savings, market expansion, etc.
- **NEGATIVE:** indications of missed earnings or revenue declines, weak or withdrawn guidance, regulatory/legal setbacks, unexpected costs, competitive pressures, etc.

Transcript:

{transcript}

Answer format

Rationale = [[your rationale]]

Result = [[POSITIVE or NEGATIVE]]

A.2 Prompt Used in Image-Text Modality

Multi-modal Modality

[image_1, image_2, image_3]

You are an expert equity research analyst. Your task is to read the following earnings-conference call transcript and the accompanying slide deck images to decide whether next quarter's earnings per share (EPS) will **beat** or **miss** consensus estimates.

Instructions

1. **Label only:** respond with exactly one of:
 - POSITIVE (indicating an expected positive EPS surprise)
 - NEGATIVE (indicating an expected negative EPS surprise)
2. **Brief rationale** (1–2 sentences) explaining the key driver
3. **No other text**

Definitions

- **POSITIVE:** indications of stronger-than-expected earnings or revenue growth, upbeat guidance, new partnerships, successful product launches, cost savings, market expansion, etc.
- **NEGATIVE:** indications of missed earnings or revenue declines, weak or withdrawn guidance, regulatory/legal setbacks, unexpected costs, competitive pressures, etc.

Transcript:

{transcript}

Answer format

Rationale = [[your rationale]]

Result = [[POSITIVE or NEGATIVE]]

A.3 Prompt Used in Audio-Text Modality

Audio-only Modality

{audio}

You are an expert equity research analyst. Your task is to analyze the provided audio from an earnings conference call and predict whether next quarter's earnings per share (EPS) will **beat** or **miss** consensus estimates.

Instructions

1. **Label only:** respond with exactly one of:
 - POSITIVE (indicating an expected positive EPS surprise)
 - NEGATIVE (indicating an expected negative EPS surprise)
2. **Brief rationale** (1–2 sentences) explaining the key driver
3. **No other text**

Definitions

- **POSITIVE:** indications of stronger-than-expected earnings or revenue growth, upbeat guidance, new partnerships, successful product launches, cost savings, market expansion, etc.
- **NEGATIVE:** indications of missed earnings or revenue declines, weak or withdrawn guidance, regulatory/legal setbacks, unexpected costs, competitive pressures, etc.

Answer format

Rationale = [[your rationale]]

Result = [[POSITIVE or NEGATIVE]]

product launches and a robust early-stage pipeline. The combination of Bristol-Myers

Finance-R1

<think> Okay, let's tackle this. The user wants to know if next quarter's EPS will beat or miss consensus estimates based on the earnings conference call transcript provided. First, I need to scan through the transcript for any mentions of EPS, revenue growth, cost savings, market expansion, or other factors that would indicate a positive or negative surprise. Looking at the transcript, the main points seem to be about the merger between Bristol-Myers Squibb and Celgene. They mention creating a leading biopharma company, combining their assets, and the strategic benefits. There are talks about expanding their portfolio, new product launches (like Revlimid, Pomalyst, and others), and the potential for growth in oncology and immunology areas. The executives talk about the financial implications of the deal, stating that they'll deliver strong returns with immediate EPS accretion. They also mention the combined pipeline with potential revenue from 6 near-term product launches. There's a lot of emphasis on the synergy and the long-term growth opportunities. However, the transcript doesn't mention any actual EPS numbers or revenue figures for the current quarter. The focus is more on the strategic moves and future outlook rather than past or current EPS performance. The key drivers here are the strategic alignment, combined pipeline potential, and the expected revenue growth from new products. The definition of POSITIVE is when earnings or revenue growth is stronger than expected, or there's new partnerships, successful product launches, cost savings, etc. The transcript mentions new product launches, expanded portfolio, and the potential for growth in key areas. Since these are positive factors, it's likely that the EPS will beat expectations. The NEGATIVE would be if earnings or revenue declines are expected, but the transcript doesn't mention any declines. The focus is on growth and positive outcomes. So the key driver here is the strategic expansion and the potential for strong revenue from the new product launches and expanded portfolio, which would lead to a positive EPS surprise. </think> Rationale = [The transcript emphasizes strategic alignment, expanded portfolio, and high-growth potential from combined pipelines, including 6 near-term product launches. These factors suggest strong revenue and EPS growth, indicating a positive surprise.] Result = [POSITIVE]

LLaMa-RAG

Good a good an the the the and **Fees to have a good a positive to and to ** to the (and to come the positive (best (1 the negative The. We the ** to have an overall to come and **to the positive the 2018 to come the **the same to come the ** to come the best of the 2018 and ** to come the best of the ** (Positive the positive the next year we the to a good a good (the negative and the negative the 1 ** the positive the **to the **the. We **the negative the higher the positive the positive (the (negative. We the 1. **the positive (Best of a positive the we, **the positive (1 **the positive and **the same the **the negative to come the positive the positive (1 **the negative the negative the positive (negative Following the The The We the negative the best (the positive (The positive, positive (negative the positive, the negative The positive to come the Positive ** to come the positive the positive, and Positive the positive, positive (negative the negative the following an overall the

Example of Response Bias

Transcript:

We expect the transaction will bring significant financial benefits to shareholders of both companies from day one. We believe we'll generate returns in excess of both companies cost of capital and deliver more than 40% accretion in the first full year. From a balance sheet perspective, we will remain in a very strong position. We project substantial free cash flow in excess of \$45 billion in the first 3 years...

Finance-R1 Output:

... The company faces potential risks in the global market and slight uncertainty in supply chains, which could impact future margins ... Result = [NEGATIVE]

Example of Poor Instruction Following

Finance-R1 Output:

... The executives seem very confident about the new product launch. The sentiment is bullish and I expect them to beat the estimates. Therefore, the next quarter earning surprise is positive.

938

939

940

941

942