

Large Language Models in Business Data Analysis and Forecasting

Introduction

Since the debut of GPT-3 in 2020, large language models (LLMs) have been increasingly explored for business analytics tasks – from interpreting numerical data (sales figures, financial metrics, market trends) to generating forecasts and narrative reports. Recent studies and industry reports suggest that LLMs can **mimic certain analytical capabilities** of humans and traditional models, sometimes with surprising accuracy, while also revealing clear limitations. Researchers have begun systematically evaluating LLM performance on tasks like time-series forecasting, financial statement analysis, and data-driven report generation ¹. This report synthesizes findings since 2020 on how LLMs such as GPT-3/4, Claude, LLaMA, and others perform in analyzing numerical business data and preparing forecasts or business reports, highlighting successful applications, failure modes, comparative benchmarks, and real-world case studies.

LLMs for Time-Series Forecasting and Trend Analysis

One line of research treats **time-series forecasting** as a sequence-prediction problem for LLMs. By encoding a time series as a string of digits, forecasting can be cast as next-token prediction – essentially asking an LLM to continue the sequence. Gruver *et al.* (2023) introduced “LLMTIME,” demonstrating that general LLMs (GPT-3 and LLaMA-2) can *zero-shot extrapolate* time series with performance **comparable to or exceeding** classic forecasting models ². In zero-shot tests on standard benchmarks, the LLM approach achieved **lower error (MAE)** than dedicated models like ARIMA, temporal CNNs, and N-HiTS in many cases ³. This suggests LLMs can capture patterns (trends, seasonality) in numeric sequences without task-specific training. The authors attribute this success to LLMs’ ability to represent multimodal predictive distributions and their inherent biases toward simplicity and repetition (which align with seasonal patterns) ⁴. Notably, with careful prompt design (e.g. ensuring each digit is a separate token), GPT-based models produced plausible forecasts and well-calibrated uncertainty estimates ³ ⁵.

However, **not all LLMs or settings excel**. The same study found GPT-4 actually underperformed GPT-3 on some time-series tasks due to technical reasons: GPT-4’s tokenization of numbers and post-training alignment (RLHF) hurt its numeric forecasting ability ⁶. For instance, GPT-4’s tendency to avoid uncertain answers led to poorer calibrated forecasts than GPT-3 ⁶. This highlights that newer or larger LLMs are not universally better for numeric tasks – details like how numbers are encoded can significantly impact performance.

Other researchers have evaluated LLMs on **understanding and summarizing time-series data**. A 2024 JP Morgan study created a taxonomy of time-series features (trends, seasonality, outliers, structural breaks, etc.) and tested state-of-the-art LLMs on identifying these patterns and generating natural-language summaries ⁷ ⁸. Their findings show that **GPT-4 is notably strong** at detecting clear patterns: in zero-shot trials, GPT-4 correctly recognized trend directions and seasonal cycles far more often than LLaMA-2, Vicuna, or GPT-3.5, substantially **outperforming those smaller models** ⁹. With chain-of-thought

prompting (asking the model to reason stepwise), GPT-4's accuracy improved further on these tasks ¹⁰. For example, GPT-4 was highly robust in identifying upward or downward trends and periodic seasonality, whereas models like Vicuna missed these more frequently ⁹. GPT-4 also excelled in **time-series Q&A** – it could retrieve exact values from a given data series with perfect accuracy and compute derived metrics like minima or maxima with low error ¹¹. These are encouraging results, suggesting GPT-4 can serve as a capable “analyst” to summarize datasets or answer business queries about the data.

That said, even GPT-4 struggled with more complex statistical features. The same benchmark noted **all models had trouble** recognizing subtle structural breaks or volatility changes, with low accuracy across the board ¹². GPT-4 in particular tended to give a default negative or uncertain answer when asked about stationarity or fat-tailed distributions, sometimes even refusing to answer with the explanation that “it is only an AI model and cannot perform the necessary statistical tests” ¹³. This points to current LLMs’ limitations in higher-order statistical reasoning – they may know common patterns but lack true statistical inference skills. In multivariate time-series scenarios (multiple related metrics), performance was only moderate for all models, indicating room for improvement in handling more intricate business datasets ¹⁴.

LLMs for Financial Statement Analysis and Forecasting

Another promising application is using LLMs to analyze **financial reports and metrics** to forecast business outcomes. A prominent example is the work by Kim *et al.* (2024), who investigated GPT-4's ability to act as a financial analyst. In their study, they fed GPT-4 with **standardized, anonymized financial statements** (balance sheets, income statements, etc.) for a set of companies and asked the model to determine whether each company's future earnings would increase or decrease ¹⁵. Remarkably, even without any industry-specific context or fine-tuning, **GPT-4 (Turbo)** was able to **predict earnings direction more accurately than human financial analysts** in their sample ¹⁵. The model's directional predictions of earnings changes not only outperformed the average human equity analyst, but also matched the accuracy of a state-of-the-art machine learning model trained specifically for this prediction task ¹⁵. GPT-4 showed particular advantage on firms that humans find difficult to analyze (e.g. smaller or loss-making companies), suggesting it picked up on subtler clues in the financial data ¹⁶. Crucially, the researchers concluded the LLM was not simply regurgitating its training data (the financials were anonymized to prevent memorization); rather, **GPT-4 generated useful analytical insights** about each company's prospects to justify its forecasts ¹⁷. In essence, the model performed a form of *financial statement analysis* akin to an equity research report – identifying patterns in ratios and trends that signal future growth or decline. As a real-world impact test, the authors built simple trading strategies based on GPT-4's stock forecasts. The **LLM-driven strategy yielded higher risk-adjusted returns (Sharpe ratios and alpha)** than those based on either human analysts' predictions or other models ¹⁸. This suggests that GPT-4's analytic edge translated into meaningful financial performance, a strong endorsement of LLM potential in this domain.

It is worth noting that the above study employed **chain-of-thought prompting** to boost GPT-4's reasoning ¹⁹. By asking the model to explain its analysis step by step, they obtained more accurate and contextually grounded outputs. This aligns with other findings that prompt techniques and few-shot examples can improve LLMs' numerical reasoning ²⁰. Indeed, domain experts are beginning to treat LLMs as collaborative analytical tools: the model might draft an analysis or forecast, accompanied by a narrative explanation, which a human can then review.

Despite these successes, **not all forecasting tasks favor LLMs over humans or traditional methods**. A comprehensive 2025 evaluation examined LLMs' abilities on classic financial analysis exercises – computing

ratios, performing *DuPont analysis* (breaking down return on equity), predicting bankruptcy via Altman Z-scores, and explicitly forecasting metrics like next-year sales or EBITDA ²¹ ²². In this study, various open-source LLMs (fine-tuned LLaMA variants and others) were benchmarked against human experts and statistical models. The results were mixed. On the one hand, LLMs did fairly well in **automation of financial calculations**: for example, models could calculate common financial ratios from input data with reasonable accuracy, and even plug those into formulas like DuPont or Z-score models with only moderate errors ²³ ²⁴. This indicates that LLMs can replicate the mechanical aspects of financial analysis (extracting numbers, doing basic arithmetic or classification of financial health). On the other hand, when it came to **forecasting complex metrics** like a company's revenue or profit for the next year, the gap between LLMs and humans became evident. In a head-to-head comparison, a seasoned human financial analyst achieved much lower forecasting error (sMAPE of ~25 for sales and ~45 for EBITDA) than any LLM configuration the researchers tested ²⁵. The **best LLM-based approach underperformed the human by a large margin**, highlighting that experienced analysts still add significant value, likely through domain knowledge and nuanced judgment ²⁵. In fact, an ensemble of traditional predictive models in that study beat all standalone LLMs, underscoring that current LLMs are **not yet matching the precision of specialized statistical methods or expert intuition for detailed financial forecasting** ²⁵.

These findings illustrate a pattern: **LLMs can excel at qualitative analysis of financial data (identifying whether a situation looks better or worse, generating commentary, etc.) but may falter at precise quantitative predictions**. Predicting the *direction* of change (up or down) – essentially a classification task – appears easier for LLMs than estimating exact magnitudes. GPT-4 could tell which companies are likely to improve earnings (often beating humans on that), yet predicting *how much* earnings will be next quarter or next year is a harder regression task where humans or specialized models still hold an edge. Moreover, LLMs' known issues with arithmetic and multi-step calculations contribute to error in precise forecasting ²³. Even GPT-4, despite improvements, can misinterpret numbers or make calculation mistakes without special prompting (e.g. it may confuse units, drop or add zeros, or misapply growth rates) ²⁶. This means for scenarios like budgeting or financial modeling where accuracy is paramount, LLMs should be used with caution or in tandem with tools that ensure correctness (such as letting the LLM generate a formula which is then executed by a spreadsheet or code).

Comparisons: LLMs vs. Traditional Methods and vs. Each Other

LLMs vs. statistical models: Early evidence shows LLMs can compete surprisingly well with traditional forecasting techniques in zero-shot settings. For instance, the LLMTIME study found GPT-based models producing forecasts with accuracy on par with ARIMA and deep learning time-series models without any task-specific training ³. This is notable because time-series models usually require training on large historical datasets for each task, whereas an LLM can be applied directly with an appropriate prompt. The advantage of LLMs becomes more pronounced when the forecasting task involves **multimodal reasoning or external knowledge**. Traditional models purely extrapolate numbers, but an LLM could, in principle, incorporate contextual knowledge (e.g. "holiday sales tend to spike in Q4" or "a new competitor might affect growth") if provided in the prompt. That said, careful evaluations warn against being too optimistic: some LLM "forecasts" may simply echo patterns seen in their training data (for example, parroting the fact that sales usually go up in December) rather than truly discovering patterns in a new dataset. Researchers from ETH Zürich point out that LLMs might even *cheat* by piggybacking on human forecasts embedded in their training corpus ²⁷. For example, if an LLM was trained on news that *analysts expect a recession next year*, it might "forecast" a recession when prompted – not due to its own analysis, but because it recalls the consensus prediction. This blurs the line between genuine forecasting ability and mere retrieval of known

outcomes. To mitigate this, evaluation methods are being refined to ensure we test LLMs' independent reasoning (e.g. only asking about truly novel future events or holding out data the model couldn't have seen) ²⁸ ²⁹. When those controls are in place, current LLMs have yet to conclusively prove they can outperform robust statistical methods on **quantitative accuracy** of forecasts – but they often match or exceed them on **qualitative insights and narrative utility**.

LLMs vs. human experts: Studies so far paint a nuanced picture. On tasks that require scanning **large volumes of data quickly or maintaining consistency**, LLMs have an edge. For example, GPT-4 can instantly read through years of financial statements or thousands of rows of metrics and provide a coherent summary or classification (positive/negative outlook), a feat hard for a human to do with equal speed and breadth. This is why GPT-4 was able to outperform human analysts in *directional* predictions of earnings changes ³⁰. The model systematically detected patterns that a human might overlook or take much longer to discern, especially when dealing with many companies or periods. LLMs also don't suffer from fatigue or cognitive bias, which can sometimes lead human forecasters astray. On the flip side, human experts bring **contextual understanding and skepticism** that LLMs lack. In the FinNLP 2025 study, the human expert's forecasts for sales/EBITDA were far more accurate than any model's ²⁵, likely because the human could factor in one-off events, industry trends, or data quirks that an LLM (which only saw the raw numbers) did not understand. Humans can also question data quality – an analyst noticing a weird spike might investigate if it's a data error, whereas an LLM will dutifully incorporate it into its output unless instructed otherwise. Moreover, current LLMs have **no true understanding of causality** in business processes; they cannot easily distinguish correlation from causation or foresee regime shifts that differ from historical patterns. Therefore, while LLMs can **augment** human analysts by handling routine analysis and providing preliminary assessments, expert oversight remains crucial for high-stakes forecasting and strategy.

Model-to-model comparisons: Among LLMs themselves, **GPT-4 consistently emerges as a top performer** on complex business analytics tasks in studies that include it. As noted, GPT-4 outshone smaller models like GPT-3.5, LLaMA-2 (70B), Vicuna, etc., in tasks such as identifying time-series features ⁹ and performing arithmetic or retrieval on financial data ¹¹. Its larger training corpus and more advanced architecture likely give it a better grasp of numerical relationships described in text (e.g. it might *remember* or *understand* accounting principles or common growth patterns more deeply). Claude (Anthropic's model) has been less frequently evaluated in published literature for quantitative tasks, but one Federal Reserve study in 2025 used **Claude 3.5 (codename "Claude-Sonnet")** to test macroeconomic knowledge. Claude showed impressive *memory* for certain facts – it could recall historical unemployment rates and inflation figures with high accuracy decades back ³¹. This suggests that domain-specific data present in training (e.g. major economic indicators) can be retained by LLMs to a surprising extent. However, Claude struggled with volatile series like GDP growth, indicating that it likely returned an averaged view (missing the extremes) ³². Interestingly, Claude's answers for GDP seemed to **blend knowledge from different time periods (first-release data vs. later revisions)**, effectively smoothing the data ³³. This kind of "blurred" recall exemplifies how a model's vast but imprecise knowledge can be a double-edged sword – it knows the general shape of the data but not the fine details, and it may inadvertently use future information when asked about past forecasts (a form of *lookahead bias*).

Open-source models like **LLaMA-2** (and derivatives such as Vicuna or fine-tuned versions) have shown decent capabilities on some benchmarks, but generally lag behind the frontier models on business tasks. In the JP Morgan time-series understanding benchmark, LLaMA-2 and Vicuna had lower accuracy in detecting patterns, and they struggled more with providing complete, correct answers ³⁴ ³⁵. Fine-tuning can narrow this gap – e.g. a tailored LLaMA-based model in the FinNLP study (referred to as "Llama 3.1/3.2" in

results) achieved reasonably high F1 scores in bankruptcy prediction after fine-tuning, even surpassing GPT-3.5 in certain classifications ³⁶. This implies that with domain-specific training data or optimization (e.g. incorporating financial calculation steps), smaller LLMs can become competent for specific tasks. Nonetheless, **the largest, general models like GPT-4 still hold an advantage in versatility and out-of-the-box accuracy**, especially in zero-shot or instruction-following scenarios. Table 1 summarizes some key comparative findings across different studies and models.

Study (Year)	LLM(s) Evaluated	Task / Domain	Key Findings
<i>Gruver et al., NeurIPS 2023</i> ² ³	GPT-3 (davinci), LLaMA-2 70B	Time-series forecasting (multiple datasets)	With proper prompting (“LLMTIME”), GPT-3 and LLaMA-2 achieved forecast accuracy comparable to or better than classical models (ARIMA, N-HiTS, etc.) in zero-shot mode ² . LLMs obtained the best or second-best MAE on each benchmark, indicating competitive performance without specialized training ³ .
<i>JP Morgan (Fons et al., EMNLP 2024)</i> ⁹ ¹¹	GPT-4, GPT-3.5, LLaMA-2, Vicuna	Time-series feature detection & summarization	GPT-4 outperformed smaller LLMs in identifying trends and seasonality (zero-shot), and answered numeric queries about the data with perfect accuracy ⁹ ¹¹ . All models struggled on advanced features (volatility shifts, stationarity), and even GPT-4 sometimes refused such queries due to its alignment limits ¹² .
<i>Kim et al. 2024 (UChicago)</i> ³⁰ ¹⁶	GPT-4 Turbo	Financial statement analysis (earnings prediction)	GPT-4, using chain-of-thought prompting, beat human analysts at predicting the direction of earnings changes from anonymized financial statements ³⁰ . It also matched a bespoke ML model’s accuracy. GPT-4’s advantage was strongest on hard-to-analyze firms (small or negative earnings) and its forecasts enabled trading strategies with higher Sharpe ratios than human-based strategies ¹⁶ .
<i>FinNLP 2025 (Wu et al.)</i> ²⁵	Fine-tuned LLaMA variants, Mistral, GPT-3.5 (tasks vary)	Financial analysis (ratios, bankruptcy, forecasting)	LLMs were able to calculate financial ratios and plug them into models (DuPont ROE, Altman Z) with moderate success. However, in forecasting next-year sales/EBITDA , even the best LLM model had much higher error than an expert human – the human’s sMAPE was ~25 vs LLMs’ ~80+ on Sales, showing a large gap in precision ²⁵ . Ensemble ML models outperformed all individual LLMs on these forecasts.

Study (Year)	LLM(s) Evaluated	Task / Domain	Key Findings
<i>Federal Reserve (Ma et al., 2025)</i> ³¹ ³³	Claude 3.5 (Anthropic)	Macroeconomic data recall & “nowcasting”	Claude exhibited remarkable recall of certain historical data (e.g. quarterly unemployment and CPI since the 1940s, with high accuracy) ³¹ . But it performed poorly on volatile series (GDP or industrial production), missing rapid swings ³² . The model’s estimates tended to mix initial reported figures with later revisions, implying it leaked future information and produced overly smooth “forecasts” ³³ .
<i>ETH Zürich (SPY Lab) 2025</i> ²⁷	GPT-4, others (conceptual evaluation)	Forecasting evaluation methodology	Identified pitfalls in current LLM forecasting benchmarks: models can exploit leaks or piggyback on known human forecasts in training data ²⁷ . Also noted that knowledge cut-off dates are unreliable (models often know events beyond stated training cutoff) and that evaluating forecast skill requires avoiding temporal leakage and trivial clues ²⁹ ²⁷ . Urges more rigorous, future-facing evaluations before concluding LLMs match human forecasters ³⁷ .

Table 1: Select evaluations of LLM performance on business data analysis and forecasting tasks (2023–2025). GPT-4 generally leads on interpretive tasks, while precise forecasting of numeric values remains challenging for LLMs in comparison to humans and specialized models.

Documented Limitations and Failure Modes

While LLMs have demonstrated impressive capabilities, researchers have also catalogued **systematic weaknesses** when it comes to numerical business analytics:

- **Arithmetic and Numerical Precision:** Even advanced LLMs are prone to basic calculation errors. For instance, models frequently miscompute financial ratios or growth percentages if not explicitly guided through each step ²³ . GPT-3.5/4 can err on multi-step math like compounding growth or summing across categories, because they rely on learned patterns rather than deterministic calculation. This has led to **inaccurate financial projections** in some cases when the model’s reasoning was flawed ²⁶ . Techniques like forcing the model to show its work (chain-of-thought) or using external tools (code execution, retrieval of formulas) are often needed to improve accuracy ²⁰ .
- **Hallucinated or Inconsistent Analysis:** LLMs may produce confident-sounding business reports that contain **fabricated or inconsistent numbers**. If asked to explain a trend, a model might incorrectly cite a percentage growth that doesn’t actually match the data given, especially if the input

data is long or if the model drifts from it. Ensuring the model sticks to provided data (via careful prompting or fine-tuning on “grounded” responses) is an active area of development.

- **Lack of True Causality or Domain Understanding:** LLMs do not truly *understand* economics or business causation; they recognize patterns in text. Thus, they might attribute a revenue increase to the wrong factor or fail to foresee a turning point that isn't reflected in past patterns. For example, a model wouldn't inherently know that a semiconductor shortage could cap a car company's sales next year unless that scenario was described in training data. This limitation means forecasts can be brittle – accurate only so long as conditions resemble historical correlations.
- **Dependence on Input Formatting and Context:** Seemingly minor changes in how data is presented can affect results. The JP Morgan study found that the **format of time-series data (plain CSV vs. JSON vs. other separators)** influenced LLM performance on certain tasks ³⁸. Likewise, the LLMTIME authors showed that inserting spaces between digits dramatically changed GPT-3's forecasting quality ⁵ ³⁹. LLMs can be sensitive to prompt phrasing; an ambiguous query might yield irrelevant answers. This brittleness requires careful prompt engineering when applying LLMs to business data.
- **Alignment and Refusal Behaviors:** As noted, GPT-4 sometimes refuses to engage in valid analytical tasks (like statistical tests or making speculative forecasts) because its alignment tuning discourages giving uncertain answers ¹³. While caution can be good (we don't want the model to bluff about important decisions), excessive refusal or safe responses can render it less useful as an analytic aide. Tuning a balance between **helpfulness and correctness** is still ongoing – too much freedom and the model hallucinates; too much caution and it declines legitimate requests.
- **Temporal Generalization and Data Cutoffs:** LLMs have a fixed training cutoff (e.g. late 2021 for GPT-3, September 2021 for early GPT-4, etc. – though newer versions and Claude have more recent training). This means their knowledge of business events or economic conditions might be outdated. They also lack *real-time data access* by default. An LLM might confidently analyze last quarter's sales with outdated assumptions if not updated. Solutions include retrieval augmentation (feeding the model up-to-date data) or fine-tuning on recent info, but without these, the model could make **anachronistic errors**. Additionally, as the SPY Lab analysis showed, even models with a stated cutoff can sometimes leak post-cutoff knowledge ⁴⁰ ⁴¹, which complicates trust in what the model truly “knows” when simulating forecasts from a past perspective.
- **Copying Training Data Predictions:** Perhaps the most insidious failure mode for forecasting is that an LLM might output a forecast it saw in its training set, rather than deducing it. For example, if many sources predicted “Market X will grow 10% next year,” a model might parrot that 10% figure when asked, instead of analyzing fresh data. This **piggybacking on human forecasts** means the model isn't adding new value – and if those forecasts were wrong, the model will be wrong too. Researchers caution that without robust evaluation, an LLM could appear accurate in hindsight simply because it regurgitated the consensus that was available to it ²⁷.

In summary, current LLMs **lack reliability in high-stakes numeric tasks**. They need oversight and often integration with tools (like calculators, databases, or factual references) to ensure the outputs are correct. Many failures have been documented in academic settings to help identify where improvements are needed. These include models misreading tables, generating infeasible projections (like negative expenses

when that's impossible), or failing to adjust their answers when obvious contextual clues change (e.g. not realizing a one-time event invalidates a trend). Such limitations underscore that **LLMs are not a wholesale replacement for traditional analytical methods or expert judgment**, though they can automate and enhance parts of the analysis process.

Real-World Deployments and Case Studies

Despite limitations, industry has quickly embraced LLMs for business intelligence and forecasting tasks, often in a **human-in-the-loop** fashion. Several notable deployments illustrate how LLMs are being used in practice:

- **HARMAN's ForecastGPT (2024):** Global tech company HARMAN introduced "ForecastGPT," a platform that leverages GPT-based models for enterprise forecasting across various domains ⁴². It is designed to generate **forecasts for sales, demand, inventory, financial KPIs, and marketing metrics** with accompanying natural-language commentary explaining the trends ⁴³. For example, an analyst can input historical revenue streams and ForecastGPT will output not only the projected revenue for upcoming quarters but also a narrative highlighting patterns (seasonal upticks, anomalies, etc.) ⁴³. The goal is to enable decision-makers to obtain **quick, interpretable forecasts** without manually building models. HARMAN claims their system can integrate with various data sources (Excel, SQL databases, etc.) and produce accurate predictions in dynamic, fast-changing environments ⁴² ⁴³. This reflects a trend of wrapping LLMs in user-friendly tools that address business-specific needs – essentially bringing the power of GPT to non-technical forecasters.
- **BloombergGPT (2023):** Bloomberg L.P. created a custom 50-billion parameter LLM trained on a vast corpus of financial data (news, filings, market data) alongside general text. While *BloombergGPT* is primarily aimed at NLP tasks in finance (like question answering, classification of news sentiment, etc.), it demonstrated the value of domain-specific LLMs. According to its creators, BloombergGPT significantly outperforms comparably sized general models on financial tasks, *without* sacrificing performance on general language tasks ⁴⁴. In effect, this model can better understand financial jargon, numerical expressions in financial text, and presumably answer business questions more accurately. For instance, it might better handle a query like "What was Company X's EBITDA in 2022 and how did it compare to 2021?" by extracting the correct numbers from its training data. BloombergGPT showcases an industry investment in tailoring LLMs to **business and finance verticals**, which likely improves their reliability when analyzing those specific types of numerical data. (Its performance relative to GPT-4 on pure forecasting isn't public, but specialization seems to yield gains in relevant sub-tasks.)
- **PwC and others adopting GPT-4:** Large consulting and financial firms have started integrating GPT-4 into their workflows. PwC announced a partnership with OpenAI to **embed GPT-4 in its tax, legal, and finance advisory services**, citing the need to boost analysts' efficiency in reviewing documents and performing analyses ⁴⁵. In practical terms, this means an analyst might use ChatGPT to summarize a 100-page financial report or to draft an initial financial risk assessment, which the human then fine-tunes. JPMorgan Chase's analysts, as another example, have experimented with ChatGPT for tasks like *drafting portions of equity research reports* or generating spreadsheet commentary (with compliance oversight). These deployments illustrate that **LLMs are being used as copilots** – generating first drafts of analyses, answering ad-hoc questions about data, or suggesting scenarios, thereby freeing human experts to focus on judgment-intensive parts of the

job. Importantly, most firms stress a *human review* step, given the known issues around accuracy and accountability.

- **Automated Insights/Narrative Science 2.0:** Even before GPTs, companies like Narrative Science specialized in “data-to-text” generation for business intelligence (turning charts and tables into written summaries). Now, GPT-based solutions have taken this further. For instance, Microsoft’s Power BI now includes generative AI features that allow users to ask questions about their data in natural language and get explanations or visuals. This is essentially an LLM parsing the query, performing an internal analysis (sometimes via an integrated calculation engine), and producing a human-readable answer. The integration of GPT-4 into Microsoft’s Copilot suite means a manager could type: “*Explain the key drivers of sales growth this quarter*” and get a narrative referencing the data (e.g. “Sales grew 5% mainly due to a 12% increase in Region A, offset by a slight decline in Region B...”). These real-world tools show LLMs adding a **conversational interface to business analytics**, enabling non-experts to glean insights from numerical data without poring over spreadsheets.
- **Finance and Trading Use Cases:** There are reports of hedge funds and banks using LLMs to parse market data and even assist in trading strategies. One case study (cited in **Financial Times**, 2023) described how an asset management firm used GPT-3 to summarize quarterly earnings call transcripts and cross-reference those summaries with financial metrics to decide trading positions. While this is primarily text analysis, it feeds into numeric decisions (forecasting stock moves). Additionally, startups are offering AI-driven financial research services where an LLM reads SEC filings or economic reports and highlights implications (e.g., “The Fed’s tone in the statement was hawkish, which could mean slower growth – consider adjusting GDP forecast downwards.”). These applications blend textual and numerical reasoning and show the broad interest in deploying LLMs in the finance domain.

Conclusion

Large Language Models have rapidly moved from intriguing AI demos to practical tools in business data analysis and forecasting. Since 2020, a growing body of evaluations indicates that LLMs like GPT-3.5/4 and their peers can analyze numerical business data in a human-like manner, offering coherent explanations and even competitive predictions in certain cases. They have successfully identified trends in time-series data, generated forward-looking statements from financial reports, and automated routine analytical tasks. In comparative benchmarks, LLMs have matched or exceeded traditional statistical models on some forecasting tasks and even outperformed human analysts at specific predictive challenges ² ³⁰ . These successes underscore the *transformative potential* of LLMs in fields like finance, operations, and market research – they serve as tireless, ultra-fast assistants that can draft insights from raw numbers.

However, current LLMs are **far from infallible or sufficient on their own**. Clear failure modes have been documented, especially when precise numerical accuracy or rigorous statistical reasoning is required ²⁵ ²⁶ . LLMs tend to be most reliable as *analytical summarizers* (telling the story in data) rather than *exact forecasters* (pinpointing future values). As such, the consensus emerging from academia and industry is that LLMs will not replace traditional analytics or expert forecasters, but rather **augment** them. A likely best practice is to use LLMs to handle the heavy lifting of data digestion and initial hypothesis generation, and then have human analysts or solid quantitative models refine and validate the outputs. In real-world

deployments like ForecastGPT and BloombergGPT, we already see this principle – the AI provides draft forecasts and narratives, while humans provide guidance, ground truth data, and critical oversight.

Looking ahead, ongoing research is closing some of the gaps. Techniques like fine-tuning LLMs on domain-specific data (e.g. financial time series), integrating calculation engines or databases (to verify facts and do math), and improving prompt strategies (e.g. better chain-of-thought for finance) are all showing promise ²⁰. New models continue to push the envelope of both size and training data freshness, which may enhance their quantitative reasoning. There is also increasing emphasis on evaluation protocols that truly test an LLM's forecasting capability in realistic conditions (to ensure we aren't misled by it memorizing the past) ²⁹ ³⁷.

In summary, **LLMs have proven to be powerful tools for business data analysis, delivering human-like insights at scale, but they must be applied with awareness of their limitations.** Successful use cases so far pair the pattern-recognition strength of LLMs with the rigor of traditional analytics and human expertise. As research and practice evolve, we can expect LLMs to become an integral part of business intelligence workflows – not as oracles that perfectly predict the future, but as versatile aides that enhance our ability to interpret data and make informed forecasts.

¹ ⁷ ⁸ ⁹ ¹⁰ ¹¹ ¹² ¹³ ¹⁴ ³⁴ ³⁵ ³⁶ ³⁸ [aclanthology.org](https://aclanthology.org/2024.emnlp-main.1204.pdf)

<https://aclanthology.org/2024.emnlp-main.1204.pdf>

² ³ ⁴ ⁵ ⁶ ³⁹ [arxiv.org](https://arxiv.org/pdf/2310.07820)

<https://arxiv.org/pdf/2310.07820>

¹⁵ ¹⁷ ¹⁸ [\[2407.17866\] Financial Statement Analysis with Large Language Models](https://arxiv.org/abs/2407.17866)

<https://arxiv.org/abs/2407.17866>

¹⁶ ¹⁹ ³⁰ ⁴⁴ ⁴⁵ [GPT-4 Turbo Outperforms Humans in Financial Predictions | by Glenn Hopper | Medium](https://medium.com/@glenn_53777/gpt-4-turbo-outperforms-humans-in-financial-predictions-cceeea24c4cf)

https://medium.com/@glenn_53777/gpt-4-turbo-outperforms-humans-in-financial-predictions-cceeea24c4cf

²⁰ ²¹ ²² ²³ ²⁴ ²⁵ ²⁶ [Can Large language model analyze financial statements well?](https://aclanthology.org/2025.finnlp-1.19.pdf)

<https://aclanthology.org/2025.finnlp-1.19.pdf>

²⁷ ²⁸ ²⁹ ³⁷ ⁴⁰ ⁴¹ [LLM Forecasting Evaluations Need Fixing | SPY Lab](https://spylab.ai/blog/forecasting-pitfalls/)

<https://spylab.ai/blog/forecasting-pitfalls/>

³¹ ³² ³³ [Total Recall? Evaluating the Macroeconomic Knowledge of Large Language Models](https://www.federalreserve.gov/econres/feds/files/2025044pap.pdf)

<https://www.federalreserve.gov/econres/feds/files/2025044pap.pdf>

⁴² ⁴³ [HARMAN Introduces ForecastGPT, a Generative AI powered Forecasting Platform | HARMAN](https://news.harman.com/releases/harman-introduces-forecastgpt-a-generative-ai-powered-forecasting-platform)

<https://news.harman.com/releases/harman-introduces-forecastgpt-a-generative-ai-powered-forecasting-platform>