# Introduction

Vision-capable large language models (LLMs) have recently emerged, enabling AI to interpret images in addition to text. Notable examples include **OpenAI's GPT-4 Vision (GPT-4V)**, Google's **Gemini** multimodal models, and Anthropic's **Claude 3** family with vision (e.g. the "Opus" model). These systems are designed to analyze visual data like charts and graphs and answer questions or provide explanations about them. This report reviews studies and evaluations of such models on **interpreting data visualizations** – from business charts (financial performance, KPIs, revenue trends) to scientific and medical graphs – and assesses their performance on key tasks: identifying trends and patterns, spotting anomalies, generating narrative summaries, and even making forecasts based on chart data. We highlight documented successes and limitations, compare their performance to humans and older tools, and note techniques (prompting or fine-tuning) that improve their chart-reading abilities.

## Trend and Pattern Recognition in Charts

One fundamental skill is reading the overall **trends and patterns** in a chart. Studies show GPT-4V can reliably recognize high-level trends, comparisons, and extremes in visualized data. For example, an academic evaluation of GPT-4V found it capably **identifies overall trends and extreme values** in charts, and makes correct comparisons between data points [1] . In practice, GPT-4V often succeeds at describing whether a time-series line is rising or falling, which category dominates a bar chart, or which point is highest/lowest in a plot. It can even combine visual cues with its world knowledge – for instance, when shown a complex AI research scatter plot, GPT-4V correctly described an increasing trend in compute over time and noted notable points like "GPT-4" and "AlphaGo" at the high end [2] . The model interpreted axes and log-scales properly and even explained a policy threshold line on the graph (a dashed line indicating a government disclosure mandate) in context [3] . This demonstrates the model's strength in summarizing pattern information and integrating relevant context.

However, performance can vary with task complexity and input detail. A comprehensive benchmark called **ChartInsights** tested 18 multimodal models (including GPT-4V and other open and closed models) on low-level chart analysis tasks (e.g. reading distributions, correlations, etc.). Under a basic prompt, the average accuracy was only 36%, but GPT-4V achieved the **highest accuracy ~56%** on these tasks [4] . This suggests GPT-4V currently leads its peers in pattern recognition from charts, but still only gets about half of such detailed questions correct without special prompting. Google's **Gemini** models are also designed for multimodal reasoning and "can help make sense of complex written and visual information," according to Google [5] . In internal testing, **Gemini Ultra** has shown state-of-the-art performance on many benchmarks, presumably including visual tasks, though detailed public results on chart interpretation are sparse. Anecdotally, early users found **Gemini's chart-reading** ability to be similar to GPT-4V's – in one report, both GPT-4V and an early Gemini vision model struggled with a simple bar chart, misidententing which bar was highest and drawing inconsistent conclusions [6] . This indicates that out-of-the-box, even cutting-edge models might confuse basic comparisons in charts at times.

Anthropic's vision-enabled **Claude 3** has rapidly improved on visual reasoning tasks as well. The latest *Claude 3.5 Sonnet* model (2024) specifically touts "step-change improvements" in tasks requiring visual

reasoning, **"like interpreting charts and graphs."** In fact, Claude 3.5's vision was noted to surpass the earlier Claude 3 Opus on chart interpretation tasks [7] . The Claude v3 model lineup (Haiku, Sonnet, Opus) indicates that the largest model *Opus* "demonstrates exceptional performance in… **chart and graph comprehension**" [8] . This suggests Claude's top-tier model can reliably discern trends and structure in visual data. While direct head-to-head data is limited, one academic test reported a smaller Claude 3 model (Haiku) reached about 49.5% accuracy on low-level chart questions, a bit behind GPT-4V's 56% [9] . We can expect Claude's newest Opus model to narrow this gap, given Anthropic's improvements.

In summary, vision-enabled LLMs **can indeed pick up major patterns** in charts: rising or falling trends, comparative magnitudes, and outliers. GPT-4V in particular excels at high-level insights ("rough value" judgments), sometimes even handling multi-step reasoning about a chart better than humans. Researchers observed GPT-4V struggling with a straightforward value lookup (reading an exact bar height for "Japan") that most humans answered easily, yet the same model correctly answered a more complex question requiring counting multiple bars ("how many bars are shorter than Thailand's") [10] [11] . This counter-intuitive result – easier handling of a multi-comparison task than a single precise reading – highlights that GPT-4V grasps general patterns well (what's larger/smaller overall, etc.) while faltering at exact numeric perception. Overall, for broad **trend analysis**, these models show strong performance approaching human-level understanding in many cases, especially when obvious patterns are present.

## Identifying Anomalies and Key Data Points

Identifying **anomalies, outliers, or key data points** in visual data is another important capability. Vision-LLMs can often point out maxima, minima, or any category that stands out. GPT-4V, for instance, reliably finds extreme values when they are visually evident (tallest bar, highest peak on a line). In the ChartInsights benchmark, one of the tasks was to **find extrema**, which GPT-4V could do with moderate success even from the image alone – about 52% accuracy without any data provided [12] . (With exact data given, it reached 87% on extrema-finding [13] .) This indicates the model can catch the highest or lowest points a bit more than half the time by eyeballing the chart, which, while far from perfect, outperforms chance and earlier models. In fact, GPT-4V's accuracy on finding extrema or making comparisons was higher than on precise value lookups, confirming it is better at relative judgments than at reading exact figures [14] [15] .

For more subtle anomalies (e.g. a data point that doesn't fit a trend or an unexpected spike), the models can flag them if prompted. In finance or KPI dashboards, a **sudden dip or surge** is an anomaly of interest. There is evidence that GPT-4V can notice such points – one industry analysis notes GPT-4V can help identify "emerging patterns [and] market anomalies" in financial charts [16] . The model's broad training could allow it to recognize that, say, one quarter's revenue is unusually low compared to a steady historical trend, and flag that as noteworthy. However, users have also observed that these models sometimes **over-interpret or hallucinate reasons for anomalies**. A Reddit user testing simple charts found GPT-4 and Gemini not only misread some values but also *"hallucinated backstories about the data"* when asked about the chart [17] . This means the model might detect an outlier and invent an explanation (e.g. claiming a reason for a sales dip that isn't actually given by the chart). Such behavior is a known limitation – the model is biased to provide coherent narratives, so it may attribute causality or context that wasn't actually present, leading to misinformation.

Another challenge is distinguishing true anomalies from normal variance. In scientific and medical charts, identifying an outlier data point or an abnormal pattern (like an anomalous spike in a vital sign chart) requires careful reading. While we did not find a specific medical-chart evaluation, it's likely similar issues

arise: the models can describe obvious outliers but might miss subtle deviations if not explicitly asked. They also struggle with **color-coded anomalies**. For example, GPT-4V cannot consistently interpret multiple colors in complex charts [18] . In a stacked bar chart, an unusual segment color or a very small segment might be "anomalous," but GPT-4V often misidentifies colors or misses small segments entirely. Researchers noted GPT-4V had a 0% success rate on questions about one stacked bar chart, largely due to confusion over the color legend and segment sizes [19] [20] . Claude 3 and others likely share this limitation to some extent, though Claude's documentation emphasizes "accurate image recognition," which might help with reading legends or highlighting different-colored outliers.

In summary, **key data points and anomalies** are partly within reach: these models do highlight maxima/ minima well and can call out obvious odd-one-out points in a visualization. But their reliability is not yet at human level. They might miss subtle outliers or conversely over-report something as significant. Additionally, without careful prompting, an LLM may offer an *explanation* for an anomaly that sounds plausible but is not backed by the chart (a form of hallucination). Caution is needed – a human analyst would verify any AI-identified anomaly and its cause. On the positive side, as multimodal models improve, they could become useful assistants to quickly surface "notable points" in a dashboard or report, leaving final judgment to humans. Research like ChartInsights shows that new prompting techniques can greatly improve low-level analysis accuracy (more on this later), which would include better anomaly detection performance.

## Narrative Summaries and Explanations of Charts

One of the most valuable capabilities of these vision-LLMs is producing **natural-language summaries** of charts – essentially doing "chart captioning" or explanation. GPT-4V and similar models have demonstrated that they can generate **detailed, coherent narratives** about visualized data. For instance, GPT-4V can look at a business performance graph and produce a paragraph describing the trends and what they mean. In a demonstration, GPT-4V was given a complex scatter plot (computing power vs. date for notable AI systems) and it returned a rich explanation: it described the axes, noted the logarithmic scale, identified clusters of points (earlier systems vs. modern systems), named specific examples (GPT-4, AlphaGo at the high end; Perceptron at the low end), and even explained a policy annotation on the chart [2] [3] . This kind of contextual, insightful summary goes well beyond simple image captioning – the model is effectively acting like a data analyst, **translating visual data into a written report**.

Such successes are reported across domains. Google's Gemini, integrated into tools like Google Sheets, can not only create charts but also *"provide insights and explanations"* about them [21] . For example, a user could ask Gemini (via a chatbot in Google Cloud or Workspace) "What are the key trends in this sales chart?" and receive an answer highlighting that sales grew steadily, with a seasonal peak, etc. Open-source efforts have also tackled this: researchers have built specialized models like **ChartLlama** to improve chart understanding and generation. ChartLlama was fine-tuned on a large, diverse set of chart images and their descriptions; as a result, it **outperforms prior methods on chart-to-text tasks**, generating more accurate and detailed explanations of charts than both older dedicated models and general models like LLaVA [22] . In qualitative comparisons, ChartLlama produced coherent multi-sentence summaries of charts while baselines like LLaVA-1.5 sometimes lapsed into repetition or nonsense when describing a complex figure [23] . This underscores how fine-tuning on visualization data can significantly enhance the narrative abilities of an LLM for charts.

Despite these advances, **limitations in chart explanations** remain. A major issue is **hallucination** – the model might include details not actually present. The Georgia Tech evaluation of GPT-4V noted occasional hallucinated or inconsistent statements when the model explained charts [24] [25] . For example, GPT-4V might assert a chart has a certain label or trend that isn't there upon close inspection. Another study in the medical domain (neuro-ophthalmology images) found GPT-4V gave mostly accurate interpretations but still made some factual errors [26] (domain-specific chart interpretation can suffer if the model's training data or knowledge is lacking). Furthermore, if a chart is poorly designed or has deceptive elements, the model's summary can be led astray. Researchers tested GPT-4V on **misleading charts** (e.g. truncated axes, inverted axes) and found the model *can* be fooled by these just like humans [27] . For instance, an inverted y-axis (low values at top) might trick GPT-4V into saying a trend is decreasing when in reality (with normal axes) it's increasing. This is a limitation in the model's visual literacy – it does not have an innate "sense" for normal conventions if the chart deviates from them, unless it catches that clue in the image.

When it comes to **identifying key insights and explaining them clearly**, the models sometimes also show a lack of focus. They may enumerate everything visible rather than zeroing in on the most important insight. A human analyst looking at a KPI dashboard will pick out the 1–2 critical takeaways ("Revenue hit an all-time high, but customer churn spiked unusually"). An LLM might instead produce a verbose description covering all bars or lines unless prompted to summarize concisely. That said, there are promising reports: one evaluation praised GPT-4V's ability to produce **"detailed reports, identifying patterns and providing insights that might be missed"**, which suggests it could surface subtle points an analyst overlooks [28] . In education, GPT-4V has been used to generate explanations of charts to students, adjusting the language to be understandable [29] , showing the potential for AI to make data stories more accessible.

Overall, **narrative chart summarization is a strong suit** of these models, especially GPT-4V. They combine visual reading with language generation to not only state what is in the chart but also to contextualize it (sometimes pulling in outside knowledge). Successes include generating multi-paragraph explanations of business graphics and scientific figures that are largely accurate. The key caveat is that a human should verify these summaries for any subtle mistakes or invented details. As fine-tuned models (like ChartLlama or others) reduce hallucinations and as prompting techniques improve, AI-generated chart narratives could become reliable enough to use in business intelligence reports or news articles as first drafts. In fact, media and industry reports have begun exploring using GPT-4V or Gemini to automatically write commentary on charts, though always with a human editor in the loop for now.

## Forecasting from Visual Data

**Predicting future values or trends** based on a chart is an especially challenging task – it requires not just reading the data but extrapolating it (and sometimes applying domain knowledge). There is limited formal research on LLMs making forecasts from charts, since this veers into speculative or analytical territory beyond the given data. However, some case studies give insight into what these models can and cannot do.

A notable example comes from meteorology: *King et al.* (2024) evaluated GPT-4V's ability to analyze **weather maps and produce a forecast** outlook. They tasked GPT-4V with generating a severe weather outlook (in English and Spanish) after showing it multiple meteorological charts [30] . The result was a **plausible forecast that aligned well with the official human-issued forecast** for that period [31] . In other words, GPT-4V successfully interpreted features on pressure maps and radar charts (fronts, instability indices, etc.) to predict where severe storms were likely, and its outlook was broadly on target. This is impressive, suggesting the model synthesized visual weather patterns correctly. However, the study also noted the AI's

forecast **"displayed vagueness and incorrect reasoning"** in parts [32] . For instance, it might have identified the right region of risk but been unable to articulate the meteorological reasoning clearly, or it hedged with generic language. Moreover, when producing bilingual forecasts, GPT-4V's Spanish output was a word-for-word translation of English that sounded unnatural and lost some meaning [33] . The authors conclude that while GPT-4V shows potential to aid forecasters, it should be used with caution and always under human supervision [34] [35] . This cautious stance likely applies to **any forecasting scenario**: the model can extrapolate trends (e.g. continue a line) and even use background knowledge (seasonality, typical patterns), but it is not guaranteed to apply these correctly or to know when a straight-line extrapolation is inappropriate.

In the business domain, one might ask an LLM with vision: "Here is our sales trend for the last 5 years (with a chart). What do you forecast for next year?" The model will certainly attempt an answer, often by extending the current trend. If the chart shows steady linear growth, GPT-4V may predict "It looks like sales could reach around X next year if the trend continues" (perhaps computing an approximate growth rate). If the data is volatile, it might express uncertainty or pick an average. **No published study** yet quantifies how accurate such AI forecasts are for business metrics, which are influenced by many external factors the model cannot see in the chart. It's reasonable to assume the model's forecast is at best a simplistic projection. In finance, even human analysts hesitate to forecast purely from past-chart patterns ("past performance is not indicative of future results"). An LLM doesn't have a secret statistical engine for time-series forecasting; it mostly guesses based on the shape of the curve and any hints from its training data. Indeed, the **risk of confident-but-wrong predictions ("hallucinated" forecasts)** is real. The meteorology study highlighted how small reasoning errors in a forecast can have disproportionate consequences [36] . In a business context, an AI might misread a seasonal peak as a growth trend and forecast erroneously high numbers. Therefore, while GPT-4V or Gemini can be asked to make forecasts, their outputs should be taken as broad suggestions rather than rigorous predictions.

On the flip side, these models could be useful in generating **scenario-based narratives**. For example: "Given this revenue chart, what factors might drive next year's change?" The model could say, "If the current growth rate of ~10% quarterly continues, revenue would reach $X; however, if there's a downturn similar to 2020, it could flatten or fall." This isn't a single forecast but a qualitative discussion of possibilities, which the model's knowledge enables it to do (mixing the visual trend with general economic insight). Such use-cases blur the line between pure chart reading and applying world knowledge – something these large models are uniquely positioned to attempt, though again with the **need for human oversight** to vet the soundness of those arguments.

In summary, **forecasting based on charts is an emerging and experimental capability** of vision-LLMs. Early results (like GPT-4V's weather outlook matching a human forecast) are encouraging in demonstrating the model's comprehension of temporal progressions in data. Yet, consistent accuracy in predictions is not demonstrated, and there are clear warnings about overconfidence and errors. At present, it is safest to view any AI-generated forecast as a rough conjecture – potentially helpful as one input among others, but not a substitute for domain-specific forecasting models or expert analysis. Future fine-tuning or tool integration (for example, linking an LLM to a statistical forecasting library) could enhance this ability. For now, the **models excel more at retrospective analysis** ("tell me what the chart shows and why it matters") than at prospective analysis ("tell me what comes next").

# Techniques to Improve Visual Chart Interpretation

Researchers and practitioners have been actively exploring ways to **boost the performance** of LLMs on chart interpretation tasks. Two major approaches have emerged: clever prompting strategies and fine-tuning the models on chart-specific data. Both have shown substantial improvements in accuracy and reliability.

**Prompting strategies**: One innovative method introduced is the *"Chain-of-Charts"* prompting, analogous to chain-of-thought in text. In the ChartInsights study, Wu et al. (2024) designed this strategy specifically for low-level data analysis questions. Instead of asking a single question about a chart, they prompt GPT-4V to break down the problem into a step-by-step visual reasoning process (almost like guiding it through reading the chart piecewise). This technique yielded a remarkable improvement – GPT-4V's accuracy jumped from ~56% to **80.5%** on their benchmark tasks when using the Chain-of-Charts prompts [37] . That is a **24% absolute improvement** simply by structuring the prompt better. Moreover, they combined this with a *visual prompt* technique: essentially directing the model's attention to relevant parts of the image (for example, by describing or marking the region of interest in the prompt). With both textual and visual prompting enhancements, GPT-4V reached **83.8%** accuracy [38] on the low-level tasks – closing the gap to near expert-level performance on their test set. This demonstrates that *how* you ask the model to analyze the chart can significantly affect outcomes. Similarly, Anthropic notes that Claude's vision system employs *chain-of-thought reasoning* internally for math in visuals [39] , suggesting that prompting it to "show its work" on charts might help. As a practical tip, users have found it useful to ask multi-step questions: e.g. "First list the values you see for each category, then answer the question." Such prompts force the model to extract raw data (via its OCR ability) before reasoning, reducing mistakes. If the chart text (labels or numbers) can be read, GPT-4V will often transcribe it correctly when asked, providing a more solid basis for subsequent analysis. This two-step Q&A approach is essentially mimicking how a person might first write down chart data then analyze it.

Another straightforward workaround is to **provide the underlying data** in addition to the image. If a user has access to the data table or can copy values from the chart, giving those to the LLM will dramatically increase accuracy. The Georgia Tech evaluation showed that GPT-4V with the dataset provided hit **87% accuracy** on a chart question set, versus 31% with image alone [40] . It went from underperforming a classic chart QA system to far *outperforming* it, once the exact numbers were available [41] [42] . Essentially, GPT-4V could then do precise calculations and lookups without guessing from pixels. This suggests a practical enhancement: if you can extract chart text via OCR or a data export, feed it to the model (either in the prompt or by a tool) – you combine the model's reasoning with the data's precision. Some users advocate this approach: *"Better to convert the chart into textual data for LLMs than rely on pure image analysis,"* since current models aren't fully reliable on visual estimates [43] . OpenAI's system card also notes GPT-4V has trouble with fine-grained visual details like tiny text or very close bar heights [24] [44] , so eliminating that uncertainty helps.

**Fine-tuning models on chart data**: While GPT-4V and others are generalists, researchers have started creating specialist models tuned for charts. **ChartLlama** (2023) is one example – by generating a large synthetic training set of charts with associated questions and explanations (using GPT-4 to help generate this data), the authors trained ChartLlama to develop a deeper "understanding" of visualizations [45] [22] . The result was state-of-the-art performance across several chart interpretation benchmarks: it beat previous models on **ChartQA (question answering)**, on **chart-to-text** description tasks, and on **data extraction from charts** [22] . In qualitative tests, the fine-tuned model made noticeably fewer errors in long

explanations or complex queries than base multimodal models [46] . Another effort called **LLaVA-Chart** adapted an open-source vision-LLM (LLaVA) specifically for charts via instruction tuning [47] . These fine-tuned models incorporate "knowledge" of how to read axis tick marks, how to interpret legends, and even some best practices of visualization (things that a general model might not have seen enough of). They also tend to be less prone to hallucination in this domain, because the training pairs reinforce grounding answers in the visual content. For instance, ChartLlama's authors noted that a baseline model would often repeat words or add irrelevant text when describing a chart, whereas ChartLlama's outputs were concise and relevant [48] . Fine-tuning effectively teaches the model to be a **better behaved chart reader**.

It's worth mentioning that **model size and architecture** also play a role. Larger models (with more parameters or training data) usually perform better at these tasks. Google's Gemini Ultra, being a very large model, presumably has an edge in complex reasoning about visuals. Anthropic's Claude 3 Opus is similarly a large model specialized for tough cases (and as noted, it excels at technical diagrams and charts [8] ). Meanwhile, smaller models (like Claude Instant or open-source 7B-parameter ones) struggle more with charts out-of-the-box. But with focused training (even a smaller model can improve via fine-tuning on relevant data), their gap can be narrowed.

In conclusion, improvements in vision-LLM chart interpretation are rapidly being made through *prompt engineering and fine-tuning*. By structuring prompts to guide the model's visual reasoning – e.g. breaking a question into steps, highlighting relevant regions, or simply ensuring the model "sees" all necessary details – users can get significantly more accurate answers. And by training models on curated chart datasets (or using synthetic data generation pipelines), researchers have achieved specialist systems that outperform general models on chart-heavy tasks. These techniques address many of the current limitations, reducing errors like missed values, color confusion, or invented facts. As these enhancements trickle into production (either via updated base models or as optional "analysis modes"), we can expect vision-capable LLMs to become much more **reliable assistants for data visualization interpretation**.

## Comparison of Model Performance and Key Findings

The table below summarizes key findings about several prominent vision-capable LLMs and their performance on chart interpretation tasks:

| Model | Strengths in Chart Interpretation | Documented Limitations |
|---|---|---|
| **GPT-4 Vision (GPT-4V)** | - Recognizes high-level trends, extremes, and comparisons well [1] . <br>- Can generate rich narrative explanations of charts, referencing axes, legends, and context [2] [3] . <br>- With underlying data provided, achieves very high accuracy (87% on a chart QA benchmark) [40] and can do correct computations/derivations [49] . <br>- State-of-the-art performance (~56% accuracy) on low-level chart questions among 18 tested models [4] . | - Struggles with precise value reading from images alone (only 31% accuracy without data on the same QA benchmark) [41] . Often cannot **exactly** read chart scales or tiny differences [44] . <br>- Cannot reliably distinguish multiple colors or segments in complex charts (e.g. stacked bars) [44] [19] . <br>- Prone to **hallucinations** or inconsistent answers about chart details [25] – may invent explanations or misidentify labels. <br>- Falls short of human-level on overall visual literacy (scored ~19.7 vs humans' 28.8 on a visualization literacy test) [50] , mainly due to errors on "simple" lookup questions. |
| **Google Gemini (Vision capabilities)** | - Designed as a multimodal model from the ground up, showing strong performance on vision-language tasks. Touted as **skilled at making sense of complex visual info** (e.g. combining text and charts) [5] . <br>- Integrates with Google tools: can generate charts from data and then explain or answer questions about them in natural language [51] [52] . <br>- Handles a variety of chart types and can apply customizations (trend lines, facets) and then interpret those visualizations for the user [53] [54] . <br>- Early usage indicates ability to extract structured information from charts (e.g., reading values, plotting trends) and produce insights for business analytics [55] . | - **Anecdotal tests** show current Gemini (e.g. via Bard) can make similar mistakes as GPT-4V on basic charts [6] – misidentifying the largest category, logical inconsistencies in comparisons, etc. <br>- Still under evaluation: no detailed public benchmarks yet on its chart QA accuracy. Likely faces similar issues with precise quantitative reading and hallucination if not explicitly constrained. <br>- Being a proprietary model, less is published about failure modes; user reports suggest it sometimes gives overly general answers or needs very clear prompts to extract specific chart details. <br>- As with GPT-4V, may not account for external factors in forecasts and could give a straight-line extrapolation without caveats (requires user to prompt for uncertainty). |

| Model | Strengths in Chart Interpretation | Documented Limitations |
|---|---|---|
| **Anthropic Claude 3 (Vision, e.g. Opus & Sonnet)** | - **Improved visual reasoning in newest version**: Claude 3.5 Vision (Sonnet) outperforms Claude 3 Opus on image tasks, especially those needing reasoning like chart/graph interpretation [56] . Noticeably better at these tasks than previous Claude. <br>- **Claude 3 Opus** (the largest model) excels at technical visuals: strong at reading diagrams, charts, and mathematical plots [8] . It can handle embedded math or complex layouts that require step-by-step logic (uses internal chain-of-thought). <br>- Good at extracting text from images (OCR) even if imperfect quality [57] , which helps in reading chart titles, labels, and data annotations (useful for accuracy). <br>- Tends to maintain a polite, hedged tone in explanations, which can be useful for not overstating chart conclusions (helpful in business settings to discuss trends with uncertainty). | - Earlier Claude vision models were slightly behind GPT-4V in chart QA performance; e.g., **Claude 3 (Haiku)** scored ~49.5% vs GPT-4V's 56% on low-level tasks [9] . The gap may persist in some areas without further fine-tuning. <br>- Color and fine detail issues likely similar to GPT-4V (e.g., could mix up legend entries if many colors). No specific study published on this, but the underlying challenges of vision apply. <br>- As a large model, it can produce very verbose answers; sometimes Claude's chart explanations include extraneous details or reiteration. Users might need to prompt it to focus. <br>- Limited public data on **forecasting** with Claude vision – presumably it has no special training for predictions, so similar caution as GPT-4V (it might guess trends but with no guarantee of accuracy). |
| **Specialized Chart Models (e.g. ChartLlama, LLaVA-Chart)** | - Fine-tuned specifically on chart data, achieving **state-of-the-art** results on multiple evaluation sets [22] . ChartLlama, for example, topped prior systems in ChartQA (Q&A about charts), chart-to-text description, and data extraction tasks. <br>- Better **visual grounding**: These models stick closer to the actual chart content. ChartLlama showed fewer hallucinations (it didn't add unwarranted text as some general models did) [48]  and followed long multi-part chart instructions accurately [58] . <br>- Can handle a **wide variety of chart types and tasks** because their training data was diversified (bar, line, pie, scatter; tasks from captioning to editing charts) [45] [22] . This makes them versatile within the chart domain. <br>- Often smaller and open-source, allowing deployment in specialized applications (e.g., an analyst's assistant that is fine-tuned on the company's own dashboard styles). | - Being specialized, they may not perform as well outside the chart domain (they trade generality for skill in visualization tasks). One would use them in tandem with a general LLM. <br>- They still have limits – if a chart is extremely complex or novel in design, a fine-tuned model might struggle if that specific scenario wasn't in its training. <br>- Fine-tuned models might inherit any biases or errors present in the synthetic training data. For instance, if the data generation had systematic quirks, the model might mirror them. <br>- Many of these are research prototypes (ChartLlama is from 2023); not as battle-tested in real business environments as GPT-4 or Claude. Support and continuous improvement depend on the community or developers. |

**Table:** Summary of capabilities and limitations of leading vision-enabled LLMs on chart interpretation tasks. (GPT-4V = GPT-4 with Vision; Gemini = Google's multimodal model; Claude 3 Opus = Anthropic's vision model; specialized models = fine-tuned academic models for charts.) Each shows strong abilities to describe and analyze visual data, with different weaknesses to consider.

## Conclusion

Vision-capable LLMs like GPT-4V, Google's Gemini, and Claude 3 are **opening new possibilities** in how we extract insights from charts and graphs. They have demonstrated an ability to interpret visual data in ways that approach how a human analyst would – identifying trends, comparing values, pointing out outliers, and narrating the "story" behind the data. In business contexts, this means such AI can quickly digest a dashboard of KPIs and summarize performance or flag key changes; in scientific and medical contexts, they can read figures or plots and explain the findings in plain language. These successes, documented in both academic studies and early industry use cases, show that the technology is maturing fast.

At the same time, **important limitations persist**. Models often struggle with precise reading of chart details (e.g. exact numeric values, small text, color encodings) and can be *led into error by tricky visuals* or simply by the ambiguity of images. They also have a tendency to **hallucinate** – producing plausible-sounding but incorrect explanations or causal stories that aren't actually supported by the chart. When benchmarked against expert humans, current models still lag on overall "visual literacy" scores, especially for tasks like careful value lookup or interpretation of deceptive graphics [50] [59] . In comparisons with traditional AI systems, GPT-4V and peers outperform older chart QA pipelines in many areas (especially when given data), but can underperform them on pure image parsing in some cases [60] [61] – indicating there is room to blend approaches (for example, using computer vision techniques to supplement the LLM).

A bright spot is that **rapid progress is being made**. Techniques like better prompting (which guide the model's reasoning) and fine-tuning (which teaches the model using domain-specific examples) have dramatically improved accuracy on chart interpretation tasks [37] [22] . We can expect the next generation of models (e.g. future Gemini versions or GPT-5) to naturally incorporate some of these improvements, making them more reliable out-of-the-box. Additionally, the competitive landscape – OpenAI, Google, Anthropic, and academic groups – means there is active investment in overcoming these limitations, since the ability for AI to understand business data visuals is highly valued.

In practical terms, **human analysts and decision-makers should see these AI tools as aides, not replacements**. They excel at quickly summarizing what's in a chart and even providing a first pass at "why" (drawing on their vast information). This can save time and offer new perspectives ("did you notice this correlation?"). But the human must verify the findings and provide judgment, especially for high-stakes decisions. For now, an ideal workflow might involve the AI doing the initial visual analysis and narrative, and the human reviewing and correcting it – much like an AI junior analyst.

In conclusion, vision-capable LLMs have made impressive strides in interpreting data visualizations: they can **spot patterns, anomalies, explain charts in natural language, and even venture forecasts**, all of which were very hard for AI just a few years ago. Academic evaluations highlight both the **successes** (e.g. GPT-4V doubling performance on multi-step chart questions vs. prior systems [62] ) and the **pitfalls** (e.g. confusion with simple tasks, hallucinated details [59] ). Industry case studies echo these findings and show growing adoption in BI tools and analytics workflows. As models improve and learn to be more "chart-literate," we may soon have AI systems that can parse a company's quarterly report visuals or a scientific

paper's graphs as proficiently as a human expert – but until then, a collaborative approach that **combines AI speed and human expertise** will yield the best results in interpreting data from images.

---

1  10  11  12  13  14  15  18  19  20  24  25  27  40  41  42  44  49  50  59  60  61  62  An Empirical Evaluation of the GPT-4 Multimodal Language Model on Visualization Literacy Tasks

https://faculty.cc.gatech.edu/~john.stasko/papers/vis24-llm.pdf

2  3  Getting Started with GPT-4 Vision for Data Analysis

https://blog.mlq.ai/gpt-4-vision-data-analysis/

4  9  37  38  Evaluating Task-based Effectiveness of MLLMs on Charts

https://arxiv.org/html/2405.07001v1

5  Introducing Gemini: Google's most capable AI model yet

https://blog.google/technology/ai/google-gemini-ai/

6  17  43  Vision models that can read charts correctly? : r/LocalLLaMA

https://www.reddit.com/r/LocalLLaMA/comments/1bm7wsz/vision_models_that_can_read_charts_correctly/

7  56  57  Introducing Claude 3.5 Sonnet \ Anthropic

https://www.anthropic.com/news/claude-3-5-sonnet

8  39  Claude v3 Vision - Relevance AI

https://relevanceai.com/llm-models/utilize-claude-v3-vision-for-effective-image-analysis

16  28  29  55  GPT-4 Vision: Overview, capabilities, use cases and benefits

https://www.leewayhertz.com/gpt-4-vision/

21  51  52  53  54  Data Analysis and Charts with Google Gemini | by Leon Nicholls | Medium

https://leonnicholls.medium.com/data-analysis-and-charts-with-google-gemini-0833ba0c7e4f

22  23  45  46  48  58  ChartLlama: A Multimodal LLM for Chart Understanding and Generation

https://tingxueronghua.github.io/ChartLlama/

26  Accuracy of the Image Interpretation Capability of ChatGPT-4 Vision …

https://pubmed.ncbi.nlm.nih.gov/39508800/

30  31  32  33  34  35  36  Pixels and Predictions: Potential of GPT-4V in Meteorological Imagery Analysis and Forecast Communication

https://arxiv.org/html/2404.15166v1

47  GitHub - zengxingchen/ChartQA-MLLM: [IEEE VIS 2024] LLaVA-Chart

https://github.com/zengxingchen/ChartQA-MLLM