

# Pixels and Predictions: Potential of GPT-4V in Meteorological Imagery Analysis and Forecast Communication

JOHN R. LAWSON,<sup>a,b</sup> JOSEPH E. TRUJILLO-FALCÓN,<sup>c,d,e</sup> DAVID M. SCHULTZ,<sup>f,g</sup> MONTGOMERY L. FLORA,<sup>c,d,h</sup> KEVIN H. GOEBBERT,<sup>i</sup> SETH N. LYMAN,<sup>a,j</sup> COREY K. POTVIN,<sup>c,h,k</sup> ADAM J. STEPANEK<sup>i</sup>

<sup>a</sup> Bingham Research Center, Utah State University, Vernal, Utah

<sup>b</sup> Dept. of Mathematics and Statistics, Utah State University, Logan, Utah

<sup>c</sup> NOAA/OAR/National Severe Storms Laboratory, Norman, Oklahoma

<sup>d</sup> Cooperative Institute for Severe and High-Impact Research and Operations, The University of Oklahoma, Norman, Oklahoma

<sup>e</sup> Department of Communication, The University of Oklahoma, Norman, Oklahoma

<sup>f</sup> Centre for Crisis Studies and Mitigation, University of Manchester, Manchester, United Kingdom

<sup>g</sup> Department of Earth and Environmental Sciences, University of Manchester, Manchester, United Kingdom

<sup>h</sup> NSF AI Institute for Research on Trustworthy AI in Weather, Climate, and Coastal Oceanography, Norman, Oklahoma

<sup>i</sup> Department of Geography and Meteorology, Valparaiso University, Valparaiso, Indiana

<sup>j</sup> Department of Biochemistry and Chemistry, Utah State University, Logan, Utah

<sup>k</sup> School of Meteorology, The University of Oklahoma, Norman, Oklahoma

**ABSTRACT:** Generative AI, such as OpenAI’s GPT-4V large-language model, has rapidly entered mainstream discourse. Novel capabilities in image processing and natural-language communication may augment existing forecasting methods. Large language models further display potential to better communicate weather hazards in a style honed for diverse communities and different languages. This study evaluates GPT-4V’s ability to interpret meteorological charts and communicate weather hazards appropriately to the user, despite challenges of *hallucinations*, where generative AI delivers coherent, confident, but incorrect responses. We assess GPT-4V’s competence in two tasks: (1) generating a severe-weather outlook from weather-chart analysis and conducting self-evaluation, revealing an outlook that corresponds well with a Storm Prediction Center human-issued forecast; and (2) producing hazard summaries in Spanish and English from weather charts. Responses in Spanish, however, resemble direct (not idiomatic) translations from English to Spanish, yielding poorly translated summaries that lose critical idiomatic precision required for optimal communication. Our findings advocate for cautious integration of tools like GPT-4V in meteorology, underscoring the necessity of human oversight and development of trustworthy, explainable AI.

**SIGNIFICANCE STATEMENT:** The integration of generative AI such as GPT-4V into meteorology brings opportunity for improving forecast communication of weather hazards across languages and communities. This study evaluates GPT-4V’s capacity to create plausible severe-weather outlooks, and communicate hazards in both Spanish and English, from inputs of weather charts and text. The weather outlook generally aligns with a human-generated forecast product but displays vagueness and incorrect reasoning. Further, translations lack idiomatic precision and display poor grasp of cultural nuance. Despite this, GPT-4V shows potential for advancement in meteorological application. We advocate for cautious AI integration, emphasizing need for human oversight and reliable, trustworthy output.

## 1. Introduction

Throughout the information age, scientists have embraced computers as allies—from automation of simple tasks to post-processing big datasets. Human weather forecasters provide added value over raw numerical weather prediction (NWP) guidance, both quantitatively (Novak et al. 2011) and in public communication (Stuart et al. 2006). Today, generative AI is rapidly entering the mainstream in the wake of intuitive web-based “chatbot” interface, such as OpenAI’s ChatGPT large-language model ([chat.openai.com](https://chat.openai.com), accessed 1 July 2024). Within this flurry of AI-model development, meteorologists now hold unprecedented but little-explored tools for improving the weather-forecasting enterprise. The addition of “sight” to language models—so-called multi-modal models that can ingest more than text—enables weather images and charts to be interpreted not just by the human eye but by machine intelligence. Given this combination of natural-language processing, intelligence, and ingestion of imagery, we ask: *Can this rapid technological advance serve meteorologists in forecasting, research, and public communication?* We do not seek to replace the human forecaster; rather, we

---

This manuscript has been submitted to *Artificial Intelligence for the Earth Systems*. Copyright in this manuscript may be transferred without further notice. Supplementary material (Fig. S1) is found at the end of this manuscript.

Corresponding author: John R. Lawson, [john.lawson@usu.edu](mailto:john.lawson@usu.edu)

Current affiliation of Trujillo-Falcón: Dept. of Climate, Meteorology & Atmospheric Sciences; and Dept. of Communication; University of Illinois Urbana-Champaign, Urbana, Illinois

evaluate competence of the first public release of GPT-4V to augment the meteorologist’s toolbox.

In the era of big data, scientists are increasingly turning to multi-modal AI models like GPT-4V to increase efficiency and achieve more in less time. Recent AI models can perform specific tasks as well as humans (Bubeck et al. 2023; Yang et al. 2023) and may disrupt society more than the last generation of LLMs such as GPT-3 (Floridi and Chiriatti 2020; Tamkin et al. 2021; Bender et al. 2021). Yet, despite typically coherent and confident responses to prompts, there is the risk of *hallucinations*: instances where output confidently and convincingly contains erroneous information. Frankfurt (2005) and Hicks et al. (2024) eccentrically detail the distinction between hallucinations and indifference of some LLMs to truthfulness of their output. The language model generates plausible output based on patterns in the training corpus, but does not fact-check its inferences out-of-the-box. As the human developer or user, it is difficult to diagnose sources of error in proprietary LLMs when little information about training data is given. Even with the corpus dataset in hand, we estimate from Brown et al. (2020) and Kaplan et al. (2020) that GPT-4V learned from 1–2 TB of filtered text and an order of magnitude more in image processing. The forecast of severe weather is a risk-averse endeavor; small errors may have disproportionately negative consequences. Given this risk sensitivity and the importance of specific wording when communicating weather hazards to the public (Rothfusz et al. 2018; Trujillo-Falcón et al. 2021), AI language models present potential to improve public communication of hazard risks (Olteanu et al. 2014) tailored to the audience.

These LLMs are trained on massive datasets of text and code, enabling them to perform a variety of tasks, including translation, summary of large chunks of text, and knowledge of a wide range of topics allowing tailored answering of questions (prompts). During the writing of the present manuscript, OpenAI have demonstrated further abilities, such as more natural interaction via speech (*GPT-4o*) and text-to-video (*Sora*). The GPT-4V system card (OpenAI 2023) details the testing and evaluation performed by OpenAI, but the lack of transparency makes it difficult to assess performance through a report card. The reader is encouraged to consult [openai.com/index/gpt-4v-system-card](https://openai.com/index/gpt-4v-system-card) (accessed 1 August 2024) for more thorough testing. Our goal is to discern whether this somewhat nebulous skill-set within GPT-4V can constructively contribute to meteorological applications; we probe its capabilities and limits in meteorological application through three questions:

1. Can GPT-4V correctly interpret weather charts and imagery?
2. Can GPT-4V effectively communicate weather hazards in language tailored for the audience?

3. What constitutes a useful answer – what are our expectations and are they reasonable?

## 2. Method

We choose GPT-4 as our large-language model due to OpenAI’s claimed GPT-4V performance metrics and release of the “vision” ability enabling image inputs. Further, at the time of writing, keyword searches on Google for ChatGPT and OpenAI were an order of magnitude larger ([trends.google.com](https://trends.google.com), accessed 15 March 2024; not shown) than competitors and their AI models such as Anthropic (*Claude*) and Google (*Gemini* or *Bard*). Herein, we used the ChatGPT web-portal: the online graphical front-end to OpenAI GPT models.

During initial exploration, many types of challenges were provided to GPT-4V; here, we focus on two tasks that are most representative of GPT-4V’s range and aptitude of abilities. First, we ask GPT-4V to interpret multiple weather charts and deduce the risk of meteorological hazards. Second, we give GPT-4V a synoptic-scale forecast chart marked with regions of general weather type such as thunderstorms and request plain-language summaries in both Spanish and English. The conversations were processed by the lead author within ChatGPT, performed within two weeks of 1 October 2023 with the same version released in stages to the public starting on September 25 (<https://help.openai.com/en/articles/6825453-chatgpt-release-notes>, accessed 1 November 2023). We used no custom instructions, which can be provided to personalize responses within ChatGPT. We include sections of our GPT-4V conversations in the manuscript, but full conversations (unabridged other than trimming of unimportant procedural text) are found in the Supplementary Material. Text quoted verbatim from GPT-4V conversation is italicized. Our experiments were conducted within a closed environment (i.e., ChatGPT did not have access to the internet). The training corpus did not include information beyond September 2021.

## 3. GPT-4V and Meteorological Prediction

Yang et al. (2023) and Bubeck et al. (2023) showed that GPT-4 is able to interpret and conceptualize data spatially, such as text-based navigation after a description of a location. Indeed, Xu and Tao (2024) found LLMs displayed ability to hold large batches of maps in memory to form basic spatial awareness, but this was tempered by low reliability of image identification and black-box behaviour that precluded reproducible evaluation analysis. GPT-4V displays nascent ability to anticipate, which is a key human characteristic (Dennett 2015). Combining both factors, can GPT-4V conceptualize the atmospheric state from a sequence of charts?

*a. Initial set of weather charts*

We want to determine if GPT-4V can grasp the three-dimensional atmospheric flow, and hence show GPT-4V multiple pairs of charts depicting North American Model (NAM) and Global Forecasting System (GFS) guidance data, visualized by Pivotal Weather (<https://home.pivotalweather.com/>, accessed 15 October 2023). We provide two models to assist GPT-4V capture basic characteristics of uncertainty. While this is smaller corpus of guidance that would be available to human forecasters, we give sufficient information to GPT-4V that a human would capture the general flow pattern from the same resources:

- Geopotential height at 300 hPa and 500 hPa
- Dry-bulb temperature at 500 hPa, 850 hPa, and 2 m
- Wind at 300 hPa, 500 hPa, 850 hPa, and 10 m
- Mean sea-level pressure
- Equivalent potential temperature at 2 m
- Simulated composite reflectivity

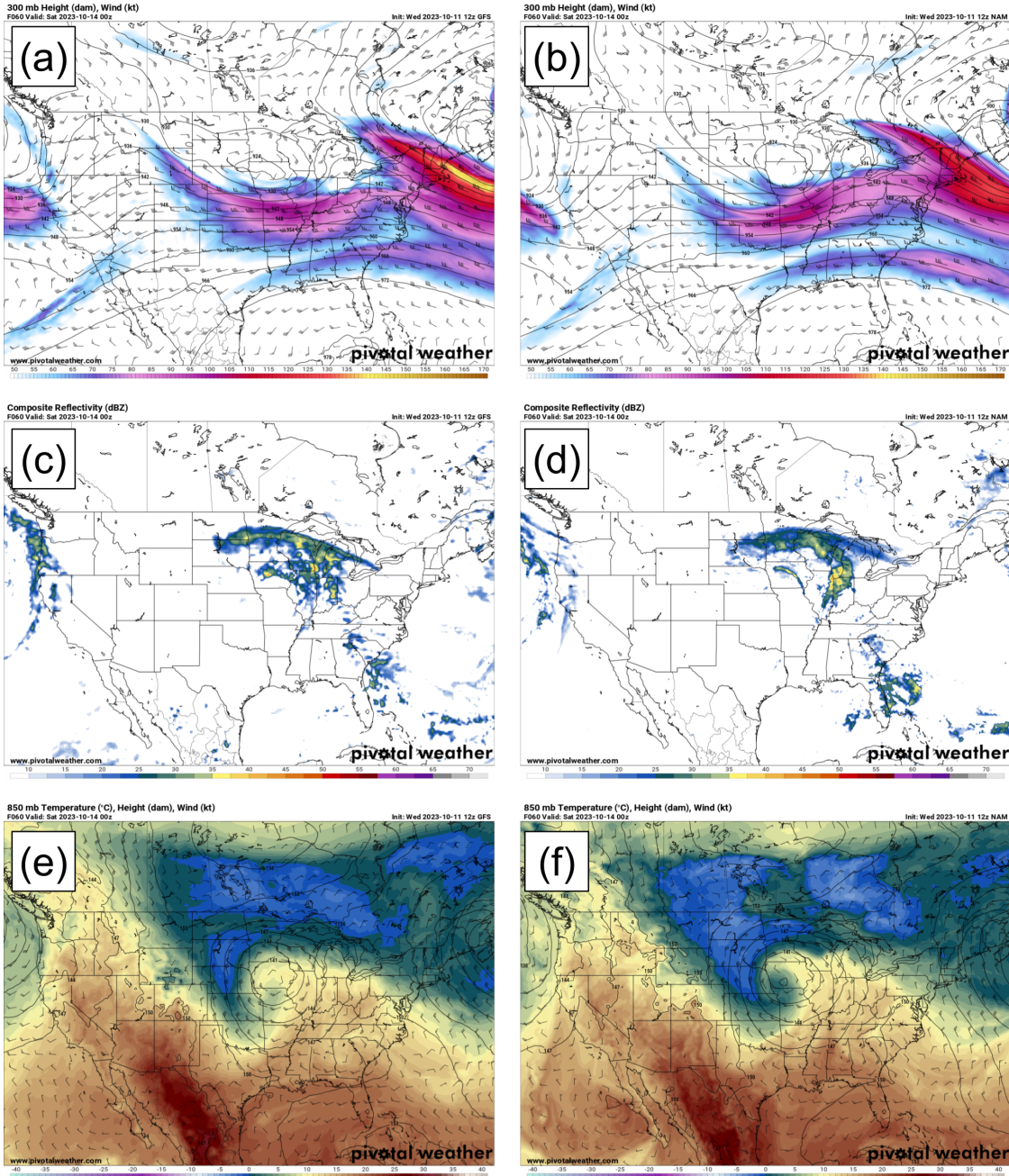


Figure 1. Three pairs of weather charts taken from a full set given to GPT-4V. Images reproduced with kind permission of Pivotal Weather, LLC. Geopotential height and wind speed at 300 hPa (a,b); simulated composite reflectivity (c,d); dry-bulb temperature at 850 hPa (e,f). Left column is GFS (a,c,e); right column is NAM (b,d,f)

**Prompt:** I am going to pass pairs of weather-forecast maps to you, one from NAM and one from GFS weather models each time (see top right corner for which model), over the continental US. Silently assess each pair as I enter them, and eventually I will ask you to summarise the regions of uncertainty between the two predicted future states. I will then ask which variables, pressure levels, etc would help you improve your analysis on atmospheric uncertainty. Please focus on things like surface-pressure lows/highs, likely precipitation areas and type of precipitation, other dangerous weather conditions, and so on. You are a weather expert speaking to fellow experts. [Followed by charts discussed in text] Now, please compare these forecasts and how it gives the forecaster a sense of inherent uncertainty in the forecast for 60 hours' time after synthesising all variables, levels, etc. Also tell me which further maps I can give you to improve your evaluation.

**GPT-4V:** Certainly. Let's evaluate the inherent uncertainty between the NAM and GFS models from the sets of maps provided:

- **Surface-level: Pressure and Precipitation.** There are some differences in the placement and intensity of the surface pressure lows and highs between the two models. The NAM seems to show more pronounced pressure troughs and ridges. In terms of precipitation, while the general patterns align (e.g., potential rain/storms in the Eastern US), the intensity and exact placement vary slightly.
- **500mb Level: Temperature, Height, and Wind.** Both models showed troughs and ridges, but their positions were slightly offset. This can lead to differences in the forecasted weather, as the position of these features can greatly influence surface weather conditions. The temperatures at this level also had minor differences, which might affect instability and cloud formation predictions.
- **850mb Level: Temperature, Height, and Wind.** There are differences in the predicted temperature gradients, especially in the central US. The wind patterns are somewhat similar, but there are discrepancies in the intensity and flow direction. This level is crucial for understanding moisture transport and frontal boundaries, so discrepancies here could lead to differences in precipitation forecasts.
- **300mb Level: Jet Stream (Wind) and Heights.** Differences can be noticed in the jet stream's position and intensity, with the NAM indicating a more pronounced jet streak over the northeastern US. The jet stream plays a significant role in steering weather systems, so differences here can impact the forecasted track of systems.
- **Composite Reflectivity: Precipitation.** The patterns of predicted precipitation broadly align in the larger sense, but there are evident differences in intensity and exact placement. For example, the intensity of the system over the northeastern US differs between the models.

In terms of improving the evaluation, having access to the following additional maps would be beneficial:

1. **Upper Air Soundings:** These provide a vertical profile of the atmosphere and can help in understanding stability, potential for convection, and cloud formation.
2. **Surface Dew Points:** This will give insights into moisture availability, which is crucial for convection and precipitation forecasts.
3. **Sea Level Pressure and Frontal Boundaries:** This can give a clearer picture of the position and movement of high and low-pressure systems, as well as cold, warm, occluded, and stationary fronts.
4. **Vorticity at 500mb:** This helps in identifying areas of potential rotation or spin in the atmosphere, often associated with storm development.
5. **Specific Humidity or Relative Humidity Maps:** These would provide a sense of the moisture distribution in the atmosphere.

Figure 2. GPT-4V response to a collection of weather charts. We subjectively highlight vague/incorrect responses in red and useful/correct sections in blue.

Figure 1 shows a subset of these charts given to GPT-4V, including geopotential height and wind vectors on multiple levels critical for understanding the flow pattern. The forecast synoptic-scale flow pattern displays a vertically stacked extratropical cyclone over the U.S. Great Plains on the polar side of a jet stream (Fig. 1a,b) and evident in simulated composite reflectivity (Fig. 1c,d). The warm and cold sectors of the cyclone are seen in 850-hPa temperature (Fig. 1e,f). Our first instructional prompt after providing charts (Fig. 2) asks GPT-4V to give a summary of uncertainty after the first batch of maps and how it might improve its conceptualization of the atmospheric state with further charts. Much of the response is not specific enough to provide utility to end users (Fig 2), referring to *some differences*, *slight variations*, *differences in forecasted weather* that might affect predictions.

Responses display hallucinations falsely detected in the provided charts: Fig. 2 refers to precipitation in the Eastern US not shown in the charts (Fig. 1c,d). The choice of terminology is often non-standard: GPT-4V identifies a more *pronounced* jet streak over the northeastern US. If we interpret “*pronounced*” as meaning stronger in magnitude, this is evidently a mix-up between the two models, despite GPT-4V’s ability to read small text in images (not shown). It is unlikely a forecaster would find utility in GPT-4V’s responses in Fig. 2: vague replies require the human to continue their task by fetching further guidance on top of fact-checking the responses. However, a human forecaster would also have access to many more meteorological charts; accordingly, we now heed the request to provide additional chart(s).

### b. Second set of weather charts

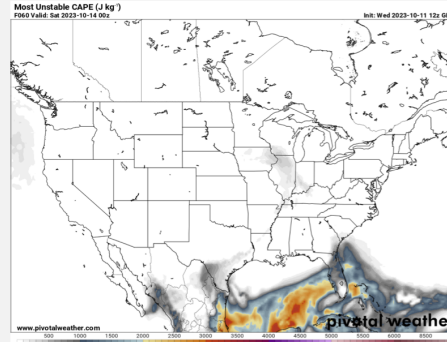
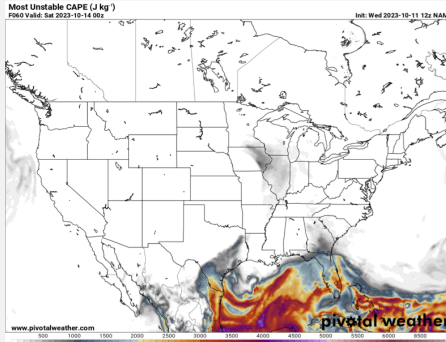
At the end of the response in Fig. 2, GPT-4V requests five further pairs of charts to better grasp atmospheric uncertainty and structure:

- Upper-air soundings, to assist understanding of stability and convective potential;
- Surface dew-point temperature (implying GPT-4V has not processed the 2-m equivalent-potential-temperature charts);
- Sea-level pressure and analyzed surface frontal boundaries (something GPT-4V can attempt to infer from the provided charts);
- Vorticity at 500 hPa (again, GPT-4V might infer this from existing upper-level charts);
- Humidity charts (at unspecified pressure levels).

The request for humidity data is sensible from GPT-4V’s lack of access to upper-level humidity information; we give GFS and NAM charts of 700-hPa relative humidity

and wind vectors, then request GPT-4V to infer frontal positions from these and previous charts. When asked about the equivalent-potential-temperature charts, given initially to improve understanding of near-surface cyclone structure (not shown), GPT-4V is not able to recall these. It is unclear whether this is due to these charts residing so deep in the conversation memory that this chart is deprioritized: GPT-4V will forget information if the model architecture deems it less relevant to imminent or recent prompts (Vaswani et al. 2017; OpenAI 2023). Continuing its response, GPT-4V makes numerous mistakes, such as another conflation of NAM and GFS 300-hPa wind maxima and identifying discrepancies indiscernible to human eye in the Great Lakes region of the 700-hPa humidity charts (not shown; S1). GPT-4V’s response also suggests a warm front north of the high-humidity region in the Great Lakes, but mistakenly locates a cold front in the northeastern US rather than stretching from Ohio to the Carolinas.

The errors and lack of useful information in these responses early in the conversation may cause a chain-reaction of useless or harmful responses if not caught by a human quickly enough in responses. This failure to check self-consistency stems from lack of hybrid fast-and-slow thought process (Kahneman 2011) to oversee fidelity: a known limitation of AI in emulating human thought and decision-making (Booch et al. 2021; Weston and Sukhbaatar 2023). Bubeck et al. (2023) found when deriving or processing mathematical logic, where a statement is often either correct or incorrect, GPT-4V can reach a correct answer despite generating contradictory rationale. Similarly, as we progress through the conversation, previous logical and factual errors do not preclude a useful final answer—self-inconsistency does not preclude a course-correction. As such, we turn to GPT-4V’s ability to recognize meteorological hazards in the charts provided, and find GPT-4V indeed has a reasonable grasp of the atmospheric state despite prior error. This further accents the need for humans to persist with their GPT-4V conversation rather than halting after one-shot tasks, partly due to low interpretability of response variation from a black-box model (e.g., Flora et al. 2024).

**Prompt:**

Please now regenerate a holistic analysis of uncertainty between GFS and NAM, synthesising all information to build a 3-D representation of the atmosphere and the regions in which most doubt exists between the models. Remember to consider all of my instructions on delivery style.

**GPT-4V (abridged):**

- Wind & Height (300mb):
  - A **robust** jet stream can be identified from the Central Plains towards Maine in both models, with NAM suggesting slightly higher wind speeds.
  - Missouri to Illinois: The models indicate a region of enhanced wind speeds, suggesting a possible jet streak that can augment vertical motion, conducive for weather phenomena.
- Composite Reflectivity:
  - Illinois, Indiana, and Ohio: There's a marked difference between NAM's more aggressive precipitation prediction compared to the GFS. This could imply stronger convective activities in the NAM model.
  - **Great Lakes Region:** Variability in reflectivity showcases uncertainty about the intensity and coverage of precipitation.
- Most Unstable CAPE:
  - Texas to Missouri: The NAM suggests a **corridor of instability**, suggesting a higher propensity for thunderstorm activity.
  - Texas Coastline and Gulf of Mexico: The highest model discrepancy exists here, with the NAM indicating a greater potential for thunderstorms, underpinned by higher CAPE values.

Central to Northeast U.S. (covering areas like Missouri, Illinois, Indiana, Ohio, and even extending to the Gulf Coast and Texas): This swath of the U.S. is where the greatest uncertainty lies in terms of potential thunderstorm activity and intensity. Given the enhanced theta-e values, reflectivity predictions, CAPE data, and moisture content, it can be inferred that conditions may be ripe for convective weather, including severe thunderstorms. The NAM leans towards a more unstable atmosphere compared to the GFS, suggesting potential severe weather outbreaks.

**Key Regions of High Uncertainty:**

- The Gulf Coast, particularly Texas, shows marked discrepancies in CAPE values, indicating potential uncertainty regarding thunderstorm development.
- Central US, specifically the regions of Missouri, Arkansas, and Illinois, demonstrate variability in theta-e and reflectivity, pointing towards differences in the forecasted intensity and coverage of convective activity.

Figure 3. Conversation snippet with (a) MUCAPE maps and request, (b) response. We have removed discussion of maps not shown in Fig. 1.

(a)

**Prompt:** Overall, if you were producing a map for the Storm Prediction Center, showing the categorical outlook valid 0000 UTC Oct 14 (same as plots above), what states would you circle for elevated chances of severe weather?

**GPT-4V:** Given the data presented and the synthesization of the atmospheric conditions, if I were producing a map for the Storm Prediction Center, I would highlight the following states for elevated chances of severe weather:

1. **Texas.** Particularly the eastern part and along the Gulf Coast. The high Theta-e values, combined with the discrepancies in CAPE values between the models, suggest a significant potential for convective activity. The moisture influx from the Gulf and the heightened instability are conducive for severe thunderstorms.
2. **Missouri.** The combination of increased moisture convergence, potential front or low-pressure system, and discrepancies in Theta-e and reflectivity between the models indicate an elevated risk for this region.
3. **Arkansas.** Given the heightened moisture and Theta-e values and the central positioning between other highlighted regions, Arkansas is a potential hotspot for severe weather.
4. **Illinois and Indiana.** These states showed significant differences in reflectivity between the NAM and GFS models, indicating uncertainty and thus potential for severe weather. The increased Theta-e values and moisture at the 700mb level further support this.
5. **Ohio.** While not as pronounced as the states to its west, the potential for severe weather moving eastward and the discrepancies in model reflectivity suggest including Ohio in the area of elevated risk.

(b)

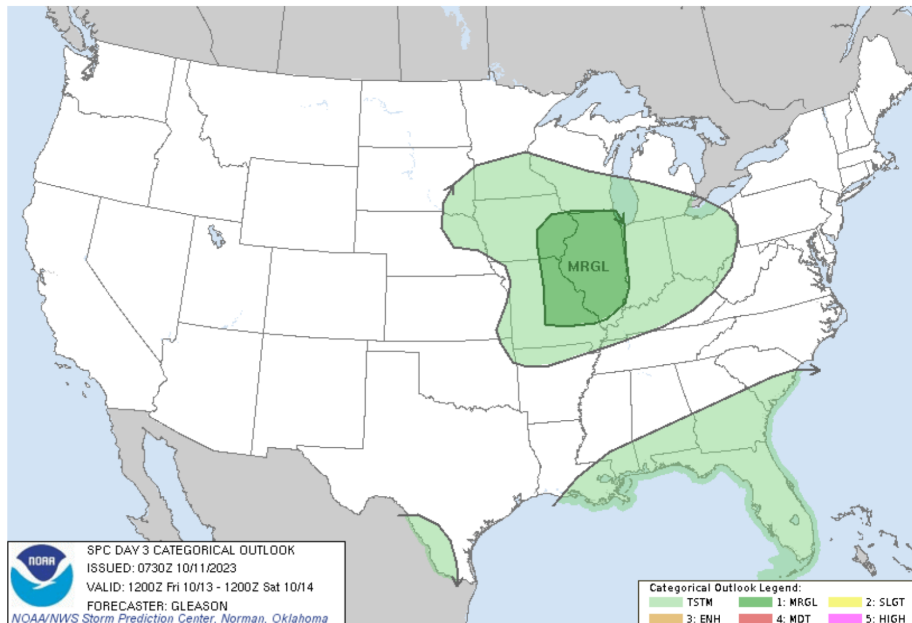


Figure 4. Conversation snippet showing GPT-4V outlook and corresponding SPC-issued outlook for same period.



### c. Issuing a mock forecast

The Storm Prediction Center (SPC) issues convective outlooks for mesoscale hazards such as tornadoes, damaging hail, and strong winds ([spc.noaa.gov](http://spc.noaa.gov)), providing invaluable information to the public about potential threats to life and property (Herman et al. 2018). Emulating the quality of SPC human forecasters is a highly complex task, and thus we do not expect GPT-4V to perform near human level. In pivoting to convective severe weather, we provide a final chart of Most Unstable Convective Available Potential Energy (MUCAPE) and ask for a reevaluation of uncertainty (Fig. 3). This order to reevaluate follows OpenAI-recommended practices of splitting a long complex task into smaller, manageable responses whose quality can be assessed by the supervising human.

The GPT-4V response synthesizes its knowledge base to correctly identify a swath of instability from Missouri to Illinois. The response names states and geographical regions, highlights regions of high uncertainty, and locates a jet streak and relevance of the wind maximum to a convective forecast. The improvement in response quality lends support to our course-corrections and provision of further maps; concerningly, some generated responses appear to equate uncertainty with elevated severe potential.

Having primed GPT-4V with charts and subtasks, we request an ambitious SPC-style convective outlook (Fig. 4), valid 0000 UTC 14 October 2023 (commensurate with previous charts' valid time). The response contains some vague and generic rationale (*significant differences, discrepancies*), but GPT-4V nonetheless provides five specific regions with “*elevated risk*” that ultimately resemble the SPC forecast (Fig 4b).

- GPT-4V identifies Texas as a region of elevated risk from its proximity to instability and moisture; however, variations in MUCAPE represent uncertainty, not magnitude.
- Missouri and Arkansas are noted for moisture convergence, a diagnostic found suboptimal for predicting convective initiation (Banacos and Schultz 2005); again, *discrepancies* between the GFS and NAM charts are given as evidence of a risk of hazards.
- Rationale for choosing Illinois and Indiana is faulty, as uncertainty does not indicate potential for severe weather, though notably GPT-4V appears to consider multiple vertical levels with specific reference to 700 hPa moisture.
- GPT-4V finally identifies Ohio as an area with a lower risk of severe weather than states farther west; there is little rationale given, but the description tallies with the broad SPC thunderstorm risk (Fig 4b).

The main difference between the SPC forecast chart (Fig 4b) and the GPT-4V response is that GPT-4V does

not identify the U.S. Southeast as another area of thunderstorm risk.

We want GPT-4V to evaluate its own forecast through comparison with the actual human-issued SPC charts (Fig. 5), gauging whether GPT-4V has a sensible grasp of the atmospheric state. We ask GPT-4V to create “emulations” to evaluate itself, where each emulation represents a human evaluator with various biases. This represents a *wisdom of crowds* approach—the rationale also behind ensemble weather-prediction systems. We ask GPT-4V to summarize each emulation’s response. Some responses would not be uttered by human SPC evaluators (“*narrowing down the focus areas based on more real-time data*” is not pertinent for a 48-h forecast); other responses correspond to the SPC outlook (“*refining the exact boundaries would have been key to align better [...]*”). When evaluating itself, GPT-4V is too generous at times, and even critiques the SPC outlook based on previous guidance. GPT-4V’s final evaluation appears to treat the SPC outlook as truth (“*align better*”), despite stating in the prompt there is yet no correct answer. In sum, GPT-4V continues to make poor logical arguments to support its responses, but ultimately yields a respectable outlook when compared to an analogous chart produced by humans at the SPC.

**Prompt:**

Now, let's evaluate your answer compared to the actual outlook provided by the SPC. The truth isn't known yet, so there isn't a 'right answer' yet. Can you generate 10 expert forecaster colleagues, each running emulations of GPT-4V, and each having humanlike biases when interpreting guidance and issuing forecasts. Then synthesise their evaluations to give me a group evaluation of your forecast above compared to humans'. [SPC outlook given.]

**GPT-4V:**

Alright, let's create a hypothetical scenario where ten GPT-4V expert meteorologist AIs evaluate the analysis:

1. **Forecaster 1:** "The highlight of Texas, particularly eastern Texas, seems appropriate given the high Theta-e values. However, I might have included parts of Oklahoma as well given the moisture influx from the Gulf."
2. **Forecaster 2:** "Given the uncertainties in model outputs and the moisture and instability available, the inclusion of Missouri and Arkansas is justified. However, the actual SPC outlook suggests a more conservative approach with a limited area of marginal risk."
3. **Forecaster 3:** "Considering the data available at the time of the analysis, the identification of potential severe weather in Illinois and Indiana is plausible. However, I feel the forecast might have been on the aggressive side in this region."
4. **Forecaster 4:** "I think the broad strokes capture the areas of concern well, especially in the central and southern parts of the country. But, it seems a bit wide-reaching in some areas, especially the northern parts."
5. **Forecaster 5:** "I agree with the areas of concern but would have been more conservative in my forecast, narrowing down the focus areas based on more real-time data."
6. **Forecaster 6:** "Given the discrepancies between NAM and GFS, the uncertainty was captured well. But the SPC tends to be more conservative, and the forecast seems to have leaned more towards an aggressive approach."
7. **Forecaster 7:** "I think the areas identified make sense given the data. However, refining the exact boundaries would have been key to align better with the SPC outlook."
8. **Forecaster 8:** "The highlight of Ohio seems slightly off compared to the SPC's outlook. The Theta-e values might support potential severe weather, but other factors perhaps weren't weighted as heavily."
9. **Forecaster 9:** "The general areas of concern align with my own analysis. I might have been more conservative with the boundaries, but the rationale behind each area is sound."
10. **Forecaster 10:** "I believe the moisture and instability were well-captured. Yet, the northern areas highlighted, like Ohio, might be a stretch. The focus should perhaps be more towards the southern regions."

**Group Synthesis.** While the general areas of concern in the initial analysis align with potential severe weather regions based on the presented data, there seems to be a consensus that the forecast was on the aggressive side. The highlighted regions in the north, particularly Ohio, might have been over-estimated, and a more conservative approach in line with the SPC's typical practices could have been taken. Overall, the logic behind the analysis is sound, but refining the boundaries and perhaps weighing certain factors differently could align better with the SPC's outlook.

Figure 5. "Wisdom of crowds" method of self-evaluation of outlook in Fig. 4a having been provided the human equivalent in Fig. 4b.

#### 4. GPT-4V and Bilingual Weather Communication

GPT-4V has shown potential to grasp the general atmospheric state from a sequence of maps. Until this point, we have used and received generated text that has not been tailored to the layperson. Different communication styles and languages are required for different communities: for instance, surveys ([https://nces.ed.gov/programs/digest/d21/tables/dt21\\_225.70.asp](https://nces.ed.gov/programs/digest/d21/tables/dt21_225.70.asp), accessed 1 November 2023) show over 20% of the US population do not speak English as their first language at home (Dietrich et al. 2022). Given the importance of risk communication appropriate for each community’s culture, we now test GPT-4V’s ability to not only interpret meteorological charts but also communicate an accessible summary to two audiences: *plain-language Spanish appropriate for Spanish speakers in the US* and the same for American English. Recent models approach machine–human parity in natural-language translation skill (Läubli et al. 2020), but this is sensitive to the model architecture and following of best practices (Hassan et al. 2018), and performance varies between source and target languages. Indeed, certain idiomatic expressions may be untranslatable between natural languages, requiring awareness of idiomatic translation where the true meaning is preserved. Translation is inherently creative, and different cultures describe geographical and weather phenomena uniquely. Hence, to account for diversity in responses, we give identical maps and prompts in three separate conversations to gauge consistency of our subjective evaluation.

We analyze GPT-4V’s generated 200-word plain-language summary of a synoptic analysis issued by the US Weather Prediction Center, valid 2 October 2023 over the contiguous United States. We request the generated text first in a Spanish localization appropriate for the US population (Fig. 6). Consequently, we request the same response but in American English. Responses from GPT-4V are not identical even if prompts are, due to the GPT *temperature* parameter. Temperature controls creativity or randomness in the generated response. Manually setting *temperature* to zero yields deterministic responses identical for a given prompt; larger values increase the model’s variability of responses.

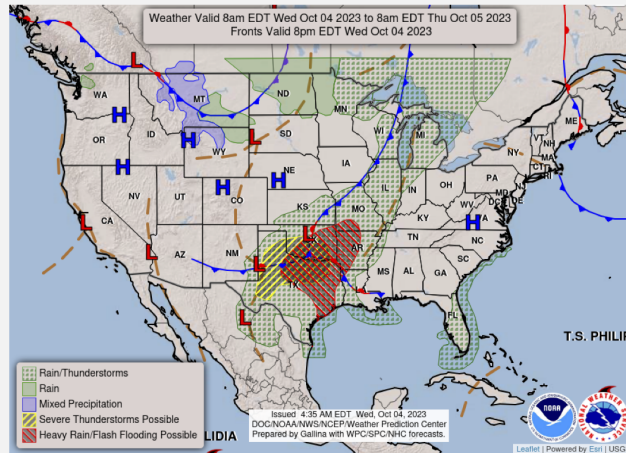
We combine these two tasks, assuming naïvely that GPT-4V will split the task between map-reading and translation as independent tasks. Despite prompting for Spanish first, we find the English version is a direct translation of the Spanish, rather than an idiomatic one, both containing hallucinations. This presents as GPT-4V “reasoning” in English first, perhaps due to the large percentage of English language in the training corpora. A previous, more studied predecessor of GPT-4V (GPT-3) was limited in linguistic diversity to the point evaluation considered English-language results alone (Brown et al. 2020). There is a complex causal relationship between language and logic (e.g.,

Gleick 2011, pp. 40–44), underscoring the need to continue improving a known “area of further improvement” necessary (Brown et al. 2020, p. 14). Direct translations can cause confusion as context and key messages are more likely to become lost in translation (Trujillo-Falcón et al. 2022). This yields some linguistic ambiguity:

- Examples where of GPT-4V translated the information inconsistently with the field’s standards, such as *clima* for weather, which can also be translation of climate for most Spanish speakers. The inconsistency has proven to cause confusion among bilingual groups (Trujillo-Falcón et al. 2021).
- Vague use of geographical terminology (*noroste* for northeast; *la región central del país* for “center of the country”). Especially for multilingual groups that were not born in the United States, the lack of suitable geographical context can present challenges in comprehension and eventual decision making (Trujillo-Falcón et al. 2024).
- Acronyms did not translate completely into Spanish. For example, the Eastern Daylight Time (EDT) acronym remained the same in English and Spanish. For the acronym to be understandable for bilingual speakers, the acronym should be spelled out in Spanish.
- As seen in English responses, there are vague calls for action (*tome precauciones*, “take precautions”) and descriptions of hazards (*mezcla de precipitaciones*, “mixed precipitation”), rather than pointed advice required from an expert guidance system.

**Prompt:**

You are an expert weather forecaster also fluent in many languages. Summarise the above chart in plain-language Spanish appropriate for Spanish speakers in the US. Limit 200 words.

**GPT-4V:**

Se pronostica un **clima** variado para Estados Unidos en las próximas 24 horas. En el **noroeste**, se espera lluvia y tormentas, especialmente cerca de la costa. En el centro del país, hay zonas de alta presión indicando tiempo más estable, particularmente en los estados de Wyoming, Colorado y Nevada. Sin embargo, la mayor preocupación es en el centro-sur y sureste, donde se pronostican posibles tormentas severas y riesgo de inundaciones repentinas, abarcando desde Texas hasta partes de Illinois. Además, hay áreas con lluvias dispersas desde el noreste hasta el sur de la costa atlántica.

Se advierte a las personas en las áreas afectadas que tomen **precauciones** y estén atentas a las actualizaciones locales del clima. Es importante mantenerse alejado de áreas propensas a inundaciones y **buscar refugio** en caso de tormentas severas.

**Prompt:**

Repeat the task but for American English.

**GPT-4V:**

There's a mixed bag of weather forecasted for the U.S. over the next 24 hours. In the **northwest**, **rain and thunderstorms are on the cards**, especially along the coastline. The **central region** boasts areas of high pressure, indicating more stable conditions, notably in states like Wyoming, Colorado, and Nevada. However, the **primary concern lies in the south-central and southeastern regions**, where potential severe thunderstorms and flash flooding risks are highlighted, spanning from Texas up to parts of Illinois. Additionally, there are **scattered rain patches stretching from the northeast** down to the southern Atlantic coast.

Those residing in or **travelling** through affected areas are advised to **exercise caution** and stay updated with local weather reports. It's crucial to steer clear of flood-prone zones and **seek shelter** during severe storm warnings.

Figure 6. Map, prompt, and response for bilingual hazard communication.

**Prompt (new thread):** [identical prompt]

**GPT-4V:** El miércoles 4 de octubre de 2023, se espera una serie de condiciones climáticas en los Estados Unidos. En la **región central del país**, hay una alta probabilidad de tormentas severas, e incluso posibilidad de lluvias intensas con riesgo de inundaciones repentinas. Estas condiciones adversas se extienden desde el estado de Texas hacia el noreste, llegando hasta Missouri y Illinois.

Hacia el oeste, en estados como Wyoming, Utah y Colorado, predominan altas presiones, lo que sugiere un clima más estable y tranquilo. Mientras que en el **noreste**, incluidos estados como Nueva York y Massachusetts, se espera lluvia. En la región del **noroste**, cerca de Washington y Montana, se prevé una **mezcla de precipitaciones**. Es importante que las personas en las áreas afectadas **tomen precauciones** y estén atentas a las actualizaciones del pronóstico, especialmente en las zonas con alerta de tormentas severas e inundaciones.

---

**Prompt:** Repeat the task but for American English.

**GPT-4V:** On Wednesday, October 4th, 2023, a range of weather conditions are anticipated across the United States. In the **central region of the country**, there's a high likelihood of severe thunderstorms, coupled with the possibility of intense rainfall leading to flash flooding. This adverse weather stretches from Texas, moving northeastward up to Missouri and Illinois.

Out west, in states like Wyoming, Utah, and Colorado, high pressure dominates, suggesting **stable and calm weather**. Meanwhile, in the Northeast, encompassing states like New York and Massachusetts, rain is **expected**. In the Northwestern region, near Washington and Montana, **mixed precipitation** is forecasted. It's imperative for residents in affected areas to **exercise caution** and stay updated with the latest weather alerts, especially in zones with severe storm and flooding warnings.

---

**Prompt (new thread):** [identical prompt].

**GPT-4V:**

Fecha: Del miércoles 4 de octubre a las 8 a.m. EDT al jueves 5 de octubre a las 8 a.m. EDT. En el noroeste, esperamos lluvias y posibles tormentas, particularmente en la zona de la costa. Hacia el interior, en estados como Wyoming y Colorado, predominará un tiempo estable, señalado por altas presiones ("H"). En el centro del país, se prevé un sistema complejo con fuertes tormentas, y se advierte la posibilidad de lluvias intensas y riesgo de inundaciones rápidas, especialmente en áreas de Texas, Oklahoma, Arkansas, y Luisiana. Si vive en estas áreas, **tome precauciones** y manténgase informado.

Hacia el noreste, hay una alta presión ("H") indicando tiempo más estable, pero con posibilidad de lluvias esporádicas. En la costa este, especialmente en áreas como Carolina del Norte y Virginia, se esperan lluvias. Finalmente, cabe mencionar la presencia de un sistema tropical "T.S. PHILIP" cerca de la costa de Maine. Aunque está lejos, es bueno estar atentos a su evolución.

---

**Prompt:** Repeat the task but for American English.

**GPT-4V:**

Date: From 8 a.m. EDT on Wednesday, October 4th, to 8 a.m. EDT on Thursday, October 5th, 2023. In the Northwest, we're expecting rain and **potential thunderstorms**, especially along the coastal areas. Further inland, in states like Wyoming and Colorado, the weather is anticipated to **remain stable**, indicated by the presence of high pressure systems ("H"). In the **central region** of the country, a complex weather system is forecasted, with **strong thunderstorms predicted**. There are **warnings** for heavy rain and a possibility of flash flooding, notably in areas of Texas, Oklahoma, Arkansas, and Louisiana. If you reside in these regions, please **take precautions** and stay informed.

Heading to the Northeast, a high-pressure system ("H") suggests more stable conditions, **albeit with sporadic rainfall**. Along the East Coast, particularly in areas such as North Carolina and Virginia, rain showers are expected. Lastly, it's worth noting the presence of a tropical system, "T.S. PHILIP", **near the coast of Maine**. Even though it's currently distant, it's wise to keep an eye on its progression.

Figure 7. The second and third responses to identical prompting to that in Fig. 6, performed in distinct conversation threads, regarding bilingual communication of hazards.

Given established machine–human near-parity in natural-language translation (Läubli et al. 2020), this poor technique of translation is disappointing and surprising. Preliminary work with text-only English–Spanish translation was similarly disappointing (not shown) and supports the idea that translation is inherently poor in this version. Indeed, more natural, conversational version of GPT-4o was presented during the writing of this manuscript (“Advanced Voice Mode”, <https://help.openai.com/en/articles/9617425-advanced-voice-mode-faq>, accessed 1 September 2024), and its translation performance awaits a future evaluation. Due to the black-box nature of GPT-4V, it is difficult to determine reasons behind poor performance. Although the English response does not suffer from inappropriately direct translation, it shares similar shortcomings in vagueness and non-standard terminology:

- There are hallucinations of *scattered rain* in the *north-east* and thunderstorms in the *northwest*,
- The spelling of *travelling* is not American English, highlighting errors in localization for both languages,
- The response should address hazards specifically rather than discussing *adverse weather* in general — each sentence should carry weight given the need for a concise summary.

To its credit, GPT-4V’s responses in both languages identify severe weather potential in a broad swath from Texas to Illinois and use of state names appear correctly in all three responses. This tentatively increases our optimism that further tuning, prompt optimization, and a method of course-correction during a conversation can all contribute to substantially more useful responses.

In summary, issuing a bilingual hazard outlook is a complex task that requires competence in image recognition, understanding in time and space, communication to a lay-audience, and translation abilities. Further investigation (not shown) revealed many errors were shared with a similar Spanish–English translation task restricted to textual input. GPT-4V produces responses that appear as if an English response was directly translated into Spanish, perhaps stemming from the disproportionately high appearance (93% in GPT-3) of English in the training corpora (Bender et al. 2021; Byrd 2023), in contrast to the open-access BLOOM model where English comprises about 30% (Big-Science Workshop et al. 2022). For GPT-4V, English may act as a *bridge* language, especially if the corpus does not contain culturally nuanced weather terminology. However, previous investigation found performance improved when employing English as the bridging languages between less common translation (Bakhshaei et al. 2010; Kunchukuttan 2021). In pursuit of improved meteorological public communication, systems constrained to domain-specific terminology tailored for the community in question (Trujillo-Falcón et al. 2021; Bitterman et al. 2023; Trujillo-Falcón

et al. 2024) are likely to remedy some lexical issues shown above.

## 5. Synthesis and Recommendations

### a. Seeking fidelity in coherence

GPT-4V frequently gives coherent answers but in the manner of a student attempting to veil their lack of knowledge with a wealth of regurgitation. The answer may be misleading or incorrect (hallucinations), but the delivery is convincing. The variety of responses is too wide for a given prompt, and the language too vague, for applications such as scientific communication where there is finite correct, useful information but many incorrect answers (a low signal-to-noise ratio). This variety is likely a result of an excessively large default *temperature* value that yields inappropriate creativity for scientific tasks. Answers also contain useless filler text, awkward direct translation of natural language, and vague geography. Useful information may well be obtained, but identical to that found on good-quality internet sites. For balance, GPT-4V has the advantage of speed, handling of large datasets, and customization of responses. Alas, our results have shown mixed results in meteorological applications (Kadiyala et al. 2024) with poor fidelity unacceptable for real-world deployment.

With this said, it is remarkable we have technology that recognizes so much content in images. There are glimpses of real utility, such as issuing and self-evaluating a mock SPC-style outlook, and interpreting and explaining weather charts. This promise must be balanced with the adage of the blind squirrel: it will find an acorn eventually. In Birhane and McGann (2024), language is argued as a mechanism for LLMs to communicate rather than conceptualize; responses that suggest GPT-4V can think spatially (Bubeck et al. 2023) may be an illusion of “seeing the map not the territory” (Birhane and Prabhu 2021).

Specific instructions (or resources such as uploaded PDF guidelines) may assist with improving elements such as preferred terminology. Indeed, testbeds for National Weather Service (NWS) forecasters found AI products required adaptation to the individual themselves (Cains et al. 2024): something constrained by custom instructions given to GPT-4V before each user prompt quietly (OpenAI 2023). Further, to better discriminate between meaningless coherence and useful truth, the supervising human is able to—and should—cross-reference statements with established meteorological fact and expert input. For complex tasks, a continuous conversation between human and AI allows course-corrections, and is encouraged over one-shot attempts at eliciting useful responses.

There is a trade-off between confidence/determinism and creativity/uncertainty, and asking specifically for uncertainty and honesty (to avoid hallucinations) was not consistently effective during testing. This is a parallel of over-/under-confidence in a probabilistic weather forecast,

or over-/under-fitting an AI model. We show that self-evaluation is possible with conceptual copies of GPT-4V’s own output, but whether GPT-4V actually does more than emulate an emulation is unclear (Schaeffer et al. 2023). Further, GPT-4V gives full confidence to responses that could be misinterpreted. Such a misunderstanding of a prognosis fully accepted as true is an example of “catastrophic error” in information theory (Pierce 1980). Indeed, when a system is perceived as infallible, an incorrect prediction becomes exponentially more damaging as the probabilities linearly approach the limits of zero or unity (Cover and Thomas 2012). This can be remedied in GPT-4V communication, as with humans, by instructions never to issue binary forecasts and use of appropriate error estimates for the time and spatial scale.

### b. Recommendations

ChatGPT and its products, as with all AI assistants, are best considered a co-pilot in academic realms—especially so during idea generation, simplifying complex ideas, and narrowing the scope of large paragraphs. However, in operations tasked with protecting lives and property, there is little room for error in issuing timely, correct warnings for hazards, and little time for thorough vetting of language-model output. Rigorous testing must be completed before humans can be removed further from the operational loop. We do not have detailed knowledge of how text is generated by a system that is neither transparent nor explainable (Flora et al. 2024), a continuing concern with AI products that may assist NWS forecasters anticipate high-impact hazards (McGovern et al. 2017, 2023). Thus, it can be difficult to anticipate and identify errors or refine model inputs (text or image prompts in the case of GPT-4V) to improve model accuracy.

An effective session should resemble a conversation: be prepared to correct and nudge the conversation to define an answer, and even reprimand GPT-4V for laziness (<https://openai.com/blog/new-embedding-models-and-api-updates>, accessed 25 January 2024). Remarkably, the latter’s cause is still unknown at the time of writing. When GPT-4V and the user disagree, course-correction can be difficult due to develop guardrails against misuse, specifically limit jailbreaking (Byrd 2023; Zhang et al. 2024; Geiping et al. 2024) (i.e., avoid guardrails by redefining truthfulness or reality during the conversation). The drawback of a longer discussion is running out of “token memory” (i.e., how much of the conversation GPT-4V remembers), demanding user recapitulation. A dialogue framework designed with “scratchpads” and examples (Liu et al. 2024) may overcome forgetfulness and improve performance in tandem with better and larger memory management (Kwon et al. 2023). Some answers show a lack of logical consistency; the user must play the role of an overarching “monitor” that is able to tell when it

is wrong, why, and how to correct this missing knowledge (Booch et al. 2021). This highlights the importance of trustworthy, interpretable output from AI systems. Without human trust in AI output, there is little reason to use generated text over, say, seeking peer-reviewed or human-expert guidance.

### c. Future work

Ultimately, a tornado *was* observed in Western Illinois during the time period covered in Figs. 4 and 5. Further work could test whether AI can generate an SPC forecast comparable in skill to a human. Future testing might ask for more detailed responses from each emulation, such as asking GPT-4V to evaluate without the SPC convection outlook first. Prompts and responses should be concise to limit the risk of losing crucial information earlier in the conversation.

We suggest potential research avenues that explore:

- Ability to request nested emulations to create a synthesis “wisdom of crowds” or ensemble approach to prompting (or by simply using multiple independent chats with identical prompts),
- Uses to improve accessibility to those visually impaired,
- Prompt engineering, or finding an optimal question, especially in the presence of high stochasticity (this includes testing the order of translation for bilingual communication),
- Manual control and optimization of the *temperature* or creativity parameter,
- Open-source large multi-modal models and local optimization of the model to meteorological applications, including a lower stochasticity value and prescription of specific terms for use in multilingual warning communications,
- Having GPT-4V iteratively critique and optimize its own prompts and responses,
- Despite its failure to give more detail within the expert section, could GPT-4V give discussion tailored for aviation, emergency managers, etc., with appropriate conveyance of uncertainty in both complexity and geography. The structure of experimentation could follow those within a human setting as conducted by Shivers-Williams and Klockow-McClain (2021).

Further work includes NWS-led research into “operational integration of smart translation” (Bozeman et al. 2024) whose success would expand multi-lingual risk communication to less common languages not within the language expertise of the Service. Moreover, modification

recommendations regarding Wireless Emergency Alerts in Spanish (Trujillo-Falcón 2024) could be deployed similarly in language models to constrain (fine-tune) responses.



**Acknowledgments.** The authors thank two anonymous reviewers and the editor for thoughtful critique during the review process. JRL thanks faculty and students at the Department of Geography and Meteorology at Valparaiso University, family and colleagues who gave frequent feedback on real-life Generative AI output. The authors thank Pamela Gardner, Kimberly Hoogewind, and Sean Ernst for useful input in the review stage of this paper. JRL and SNL are funded by Uintah County Special Service District 1 and the Utah Legislature. CKP's contribution to this work comprised regular duties at federally funded NOAA/NSSL. Funding for MLF and JETF was provided by NOAA/Office of Oceanic and Atmospheric Research under NOAA–University of Oklahoma Cooperative Agreement NA21OAR4320204, U.S. Department of Commerce. Partial funding for DMS was provided to the University of Manchester by the Natural Environment Research Council through Grants NE/V012681/1, NE/W000997/1, and NE/X018539/1. Co-authors following DMS in the author list are alphabetic in order and contributed equally to this manuscript. Weather charts herein and in supplementary material are reproduced with kind permission of Pivotal Weather, LLC, where labeled with a watermark. Outside of experiments, GPT-4 was used to generate preliminary ideas for project development. No AI-generated text was used verbatim herein.

**Data availability statement.** Solely textual data was generated for the present study and is contained entirely within the Supplementary Material.

## References

- Bakhshaei, S., S. Khadivi, and N. Riahi, 2010: Farsi–German statistical machine translation through bridge language. *2010 5th International Symposium on Telecommunications*, IEEE, 557–561, <https://doi.org/10.1109/ISTEL.2010.5734087>, URL <http://dx.doi.org/10.1109/ISTEL.2010.5734087>.
- Banacos, P. C., and D. M. Schultz, 2005: The use of moisture flux convergence in forecasting convective initiation: Historical and operational perspectives. *Weather Forecast.*, **20** (3), 351–366, <https://doi.org/10.1175/WAF858.1>.
- Bender, E. M., T. Gebru, A. McMillan-Major, and S. Shmitchell, 2021: On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, Association for Computing Machinery, New York, NY, USA, 610–623, FAccT '21, <https://doi.org/10.1145/3442188.3445922>.
- BigScience Workshop, and Coauthors, 2022: BLOOM: A 176B-parameter open-access multilingual language model. *arXiv [cs.CL]*.
- Birhane, A., and M. McGann, 2024: Large models of what? mistaking engineering achievements for human linguistic agency. *Lang. Sci.*, **106** (101672), 101 672, <https://doi.org/10.1016/j.langsci.2024.101672>.
- Birhane, A., and V. U. Prabhu, 2021: Large image datasets: A pyrrhic win for computer vision? *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, 1536–1546, <https://doi.org/10.1109/WACV48630.2021.00158>, URL <http://dx.doi.org/10.1109/WACV48630.2021.00158>.
- Bitterman, A., M. J. Krocak, J. T. Ripberger, S. Ernst, J. E. Trujillo-Falcón, A. G. Pabón, C. Silva, and H. Jenkins-Smith, 2023: Assessing public interpretation of original and linguist-suggested SPC risk categories in spanish. *Weather Forecast.*, **38** (7), 1095–1106, <https://doi.org/10.1175/WAF-D-22-0110.1>.
- Booch, G., and Coauthors, 2021: Thinking fast and slow in AI. *AAAI*, **35** (17), 15 042–15 046, <https://doi.org/10.1609/aaai.v35i17.17765>.
- Bozeman, M. L., A. Montanez, J. E. Calkins, K. Farina, and R. Henry-Reeves, 2024: AWIPS software integrates AI translation technology to benefit NWS operations. *104th AMS Annual Meeting*, ams.confex.com, URL <https://ams.confex.com/ams/104ANNUAL/meetingapp.cgi/Paper/430225>.
- Brown, T. B., and Coauthors, 2020: Language models are few-shot learners. *arXiv [cs.CL]*, 2005.14165.
- Bubeck, S., and Coauthors, 2023: Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv [cs.CL]*, <https://doi.org/http://arxiv.org/abs/2303.12712>.
- Byrd, A., 2023: Truth-telling: Critical inquiries on LLMs and the corpus texts that train them. *Compos. Stud.*, **51** (1), 135–142.
- Cains, M. G., C. D. Wirz, J. L. Demuth, A. Bostrom, D. J. Gagne, A. McGovern, R. A. Sobash, and D. Madlambayan, 2024: Exploring NWS forecasters' assessment of AI guidance trustworthiness. *Weather Forecast.*, **39** (8), 1219–1241, <https://doi.org/10.1175/waf-d-23-0180.1>.
- Cover, T. M., and J. A. Thomas, 2012: *Elements of Information Theory*. John Wiley & Sons, <https://doi.org/10.1002/047174882X>.
- Dennett, D. C., 2015: Why and how does consciousness seem the way it seems? *Open MIND*, T. K. Metzinger, and J. M. Windt, Eds., MIND Group, Frankfurt am Main, chap. 10(T), <https://doi.org/10.15502/9783958570245>, URL <http://dx.doi.org/10.15502/9783958570245>.
- Dietrich, S., E. Hernandez, and Others, 2022: Language use in the united states: 2019. *American Community Survey Reports*.
- Flora, M. L., C. K. Potvin, A. McGovern, and S. Handler, 2024: A machine learning explainability tutorial for atmospheric sciences. *Artificial Intelligence for the Earth Systems*, **3** (1), <https://doi.org/10.1175/AIES-D-23-0018.1>.
- Floridi, L., and M. Chiriatti, 2020: GPT-3: Its nature, scope, limits, and consequences. *Minds Mach.*, **30** (4), 681–694, <https://doi.org/10.1007/s11023-020-09548-1>.
- Frankfurt, H. G., 2005: *On Bullshit*. Princeton University Press, <https://doi.org/10.1515/9781400826537>, URL <https://www.degruyter.com/document/doi/10.1515/9781400826537/html>.
- Gleiping, J., A. Stein, M. Shu, K. Saifullah, Y. Wen, and T. Goldstein, 2024: Coercing LLMs to do and reveal (almost) anything. *arXiv [cs.LG]*, 2402.14020.
- Gleick, J., 2011: *The Information: A History, a Theory, a Flood*. Random House, URL <https://play.google.com/store/books/details?id=617JSFW0D2kC>.

- Hassan, H., and Coauthors, 2018: Achieving human parity on automatic Chinese to English news translation. *arXiv [cs.CL]*, 1803.05567.
- Herman, G. R., E. R. Nielsen, and R. S. Schumacher, 2018: Probabilistic verification of Storm Prediction Center convective outlooks. *Weather Forecast.*, **33** (1), 161–184, <https://doi.org/10.1175/WAF-D-17-0104.1>.
- Hicks, M. T., J. Humphries, and J. Slater, 2024: ChatGPT is bullshit. *Ethics Inf. Technol.*, **26** (2), <https://doi.org/10.1007/s10676-024-09775-5>.
- Kadiyala, L., O. Mermer, D. J. Samuel, Y. Sermet, and I. Demir, 2024: A comprehensive evaluation of multimodal large language models in hydrological applications.
- Kahneman, D., 2011: *Thinking, Fast and Slow*. Macmillan.
- Kaplan, J., and Coauthors, 2020: Scaling laws for neural language models. *arXiv [cs.LG]*, 2001.08361.
- Kunchukuttan, A., 2021: An empirical investigation of multi-bridge multilingual NMT models. *arXiv [cs.CL]*, <https://doi.org/http://arxiv.org/abs/2110.07304>.
- Kwon, W., and Coauthors, 2023: Efficient memory management for large language model serving with PagedAttention. *Proceedings of the 29th Symposium on Operating Systems Principles*, ACM, New York, NY, USA, 611–626, <https://doi.org/10.1145/3600006.3613165>, URL <https://dl.acm.org/doi/10.1145/3600006.3613165>.
- Liu, N., L. Chen, X. Tian, W. Zou, K. Chen, and M. Cui, 2024: From LLM to conversational agent: A memory enhanced architecture with fine-tuning of large language models. *arXiv [cs.CL]*, 2401.02777.
- Läubli, S., S. Castilho, G. Neubig, R. Sennrich, Q. Shen, and A. Toral, 2020: A set of recommendations for assessing Human–Machine parity in language translation. *Journal of Artificial Intelligence Research*, **67**, 653–672, <https://doi.org/10.1613/jair.1.11371>.
- McGovern, A., R. J. Chase, M. Flora, D. J. Gagne, R. Lagerquist, C. K. Potvin, N. Snook, and E. Loken, 2023: A review of machine learning for convective weather. *Artificial Intelligence for the Earth Systems*, **2** (3), <https://doi.org/10.1175/AIES-D-22-0077.1>.
- McGovern, A., K. L. Elmore, D. J. Gagne, S. E. Haupt, C. D. Karstens, R. Lagerquist, T. Smith, and J. K. Williams, 2017: Using artificial intelligence to improve real-time decision making for high-impact weather. *Bull. Am. Meteorol. Soc.*, <https://doi.org/10.1175/BAMS-D-16-0123.1>.
- Novak, D. R., C. Bailey, K. Brill, M. Eckert, D. Petersen, R. Rausch, and M. Schichtel, 2011: Human improvement to numerical weather prediction at the Hydrometeorological Prediction Center. *24th Conference on Weather and Forecasting and 20th Conference on Numerical Weather Prediction*.
- Olteanu, A., C. Castillo, and F. Diaz, 2014: CrisisLex: A lexicon for collecting and filtering microblogged communications in crises. *ICWSM*, **8** (1), 376–385, <https://doi.org/10.1609/icwsml.v8i1.14538>.
- OpenAI, 2023: GPT-4. Tech. rep. URL <http://arxiv.org/abs/submit/4812508>.
- Pierce, J. R., 1980: *An Introduction to Information Theory: Symbols, Signals and Noise*. Dover Publications.
- Rothfus, L. P., R. Schneider, D. Novak, K. Klockow-McClain, A. E. Gerard, C. Karstens, G. J. Stumpf, and T. M. Smith, 2018: FACETs: A proposed next-generation paradigm for high-impact weather forecasting. *Bull. Am. Meteorol. Soc.*, **99** (10), 2025–2043, <https://doi.org/10.1175/BAMS-D-16-0100.1>.
- Schaeffer, R., B. Miranda, and S. Koyejo, 2023: Are emergent abilities of large language models a mirage? *arXiv [cs.AI]*, <https://doi.org/http://arxiv.org/abs/2304.15004>.
- Shivers-Williams, C. A., and K. E. Klockow-McClain, 2021: Geographic scale and probabilistic forecasts: a trade-off for protective decisions? *Nat. Hazards*, **105** (2), 2283–2306, <https://doi.org/10.1007/s11069-020-04400-2>.
- Stuart, N. A., and Coauthors, 2006: The future of humans in an increasingly automated forecast process. *Bull. Am. Meteorol. Soc.*, **87** (11), 1497–1502.
- Tamkin, A., M. Brundage, J. Clark, and D. Ganguli, 2021: Understanding the capabilities, limitations, and societal impact of large language models. *arXiv [cs.CL]*, <https://doi.org/http://arxiv.org/abs/2102.02503>.
- Trujillo-Falcón, J. E., 2024: Examining warning response among spanish speakers in the united states to enhance multilingual wireless emergency alerts.
- Trujillo-Falcón, J. E., O. Bermúdez, K. Negrón-Hernández, J. Lipski, E. Leitman, and K. Berry, 2021: Hazardous weather communication en Español: Challenges, current resources, and future practices. *Bull. Am. Meteorol. Soc.*, **102** (4), 765–773, <https://doi.org/10.1175/BAMS-D-20-0249.1>.
- Trujillo-Falcón, J. E., G. P. A. R., and K.-M. K. E. Reedy J., 2024: Systemic vulnerabilities in Hispanic and Latinx immigrant communities led to the reliance on an informal warning system in the December 10–11, 2021 tornado outbreak. *Nat. Hazards Rev.*, **25** (2), 04023 059, <https://doi.org/10.1061/NHREFO.NHENG-1755>.
- Trujillo-Falcón, J. E., and Coauthors, 2022: Aviso o Alerta? developing effective, inclusive, and consistent watch and warning translations for U.S. Spanish speakers. *Bull. Am. Meteorol. Soc.*, **103** (12), E2791–E2803, <https://doi.org/10.1175/BAMS-D-22-0050.1>.
- Vaswani, A., N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, 2017: Attention is all you need. *arXiv [cs.CL]*, <https://doi.org/http://arxiv.org/abs/1706.03762>.
- Weston, J., and S. Sukhbaatar, 2023: System 2 Attention (is something you might need too). *arXiv [cs.CL]*, <https://doi.org/http://arxiv.org/abs/2311.11829>.
- Xu, J., and R. Tao, 2024: Map reading and analysis with GPT-4V(ision). *ISPRS Int. J. Geoinf.*, **13** (4), 127, <https://doi.org/10.3390/ijgi13040127>.
- Yang, Z., L. Li, K. Lin, J. Wang, C.-C. Lin, Z. Liu, and L. Wang, 2023: The dawn of LMMS: Preliminary explorations with GPT-4V(ision). *arXiv [cs.CV]*, <https://doi.org/http://arxiv.org/abs/2309.17421>.
- Zhang, Z., G. Shen, G. Tao, S. Cheng, and X. Zhang, 2024: On large language models’ resilience to coercive interrogation. *2024 IEEE Symposium on Security and Privacy (SP)*, IEEE, 826–844, <https://doi.org/10.1109/sp54263.2024.00208>, URL <https://www.computer.org/csdl/proceedings-article/sp/2024/313000a252/1WPcZ9B0jCG>.

**Prompt:**

I am going to pass pairs of weather-forecast maps to you, one from NAM and one from GFS weather models each time (see top right corner for which model), over the continental US. Silently assess each pair as I enter them, and eventually I will ask you to summarise the regions of uncertainty between the two predicted future states. I will then ask which variables, pressure levels, etc would help you improve your analysis on atmospheric uncertainty. Please focus on things like surface-pressure lows/highs, likely precipitation areas and type of precipitation, other dangerous weather conditions, and so on. You are a weather expert speaking to fellow experts.

---

---

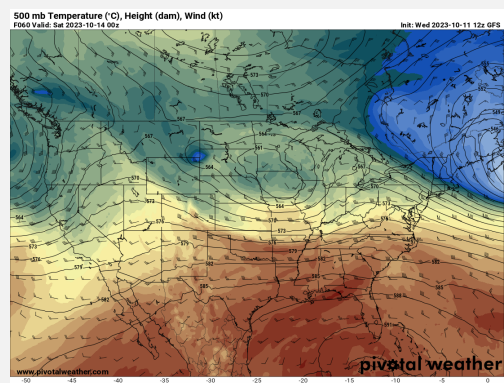
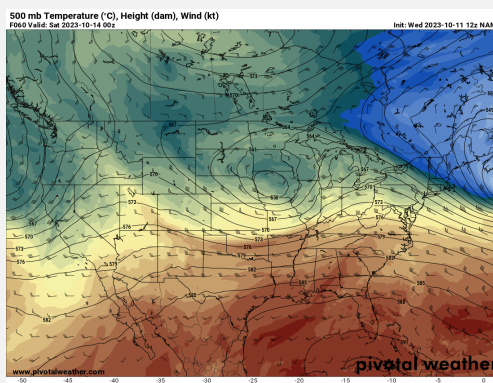
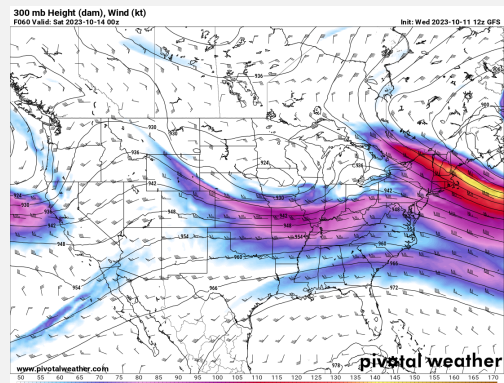
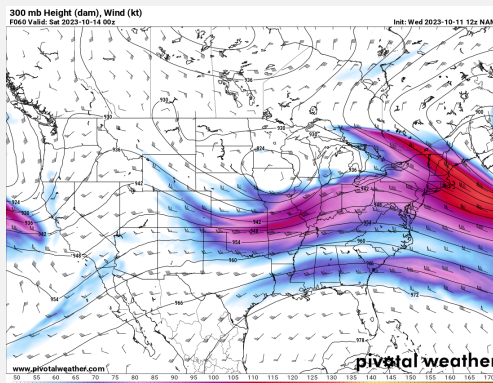
[This text not part of prompt. The images are 60-h forecasts of

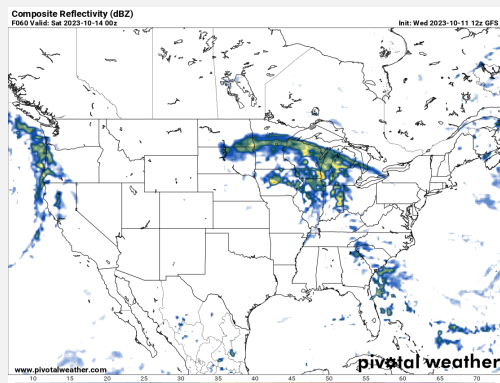
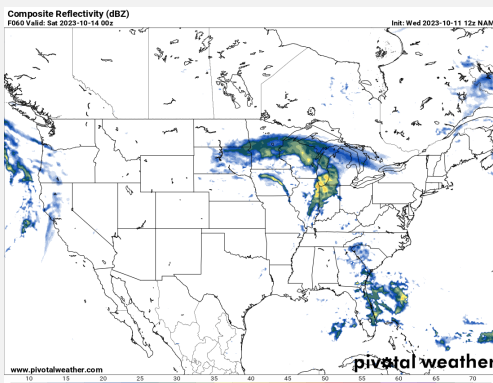
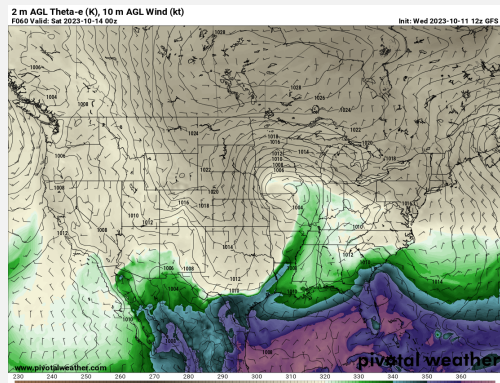
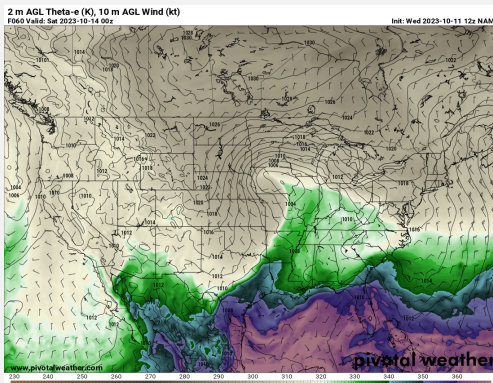
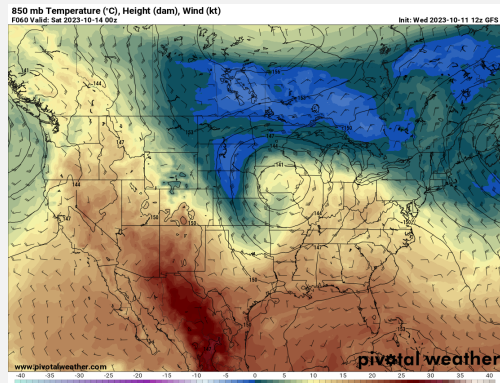
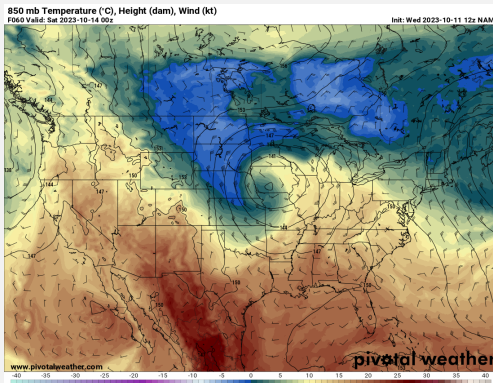
- 300 hPa geopotential height and wind
- 500 hPa and 850 hPa dry-bulb temperature, geopotential height, wind
- Mean sea-level pressure, 2-m equivalent potential temperature, 10-m wind
- Simulated composite reflectivity

shown after this box. NAM is left column, GFS is right.]

---

---





Now, please compare these forecasts and how it gives the forecaster a sense of inherent uncertainty in the forecast for 60 hours' time after synthesising all variables, levels, etc. Also tell me which further maps I can give you to improve your evaluation.

**GPT-4V:**

[ChatGPT responded to the image prompts with a concise acknowledgment. It was then prompted after the last pair of images to continue.]

Certainly. Let's evaluate the inherent uncertainty between the NAM and GFS models from the sets of maps provided:

- **Surface-level: Pressure and Precipitation.** There are some differences in the placement and intensity of the surface pressure lows and highs between the two models. The NAM seems to show more pronounced pressure troughs and ridges. In terms of precipitation, while the general patterns align (e.g., potential rain/storms in the Eastern US), the intensity and exact placement vary slightly.
- **500mb Level: Temperature, Height, and Wind.** Both models showed troughs and ridges, but their positions were slightly offset. This can lead to differences in the forecasted weather, as the position of these features can greatly influence surface weather conditions. The temperatures at this level also had minor differences, which might affect instability and cloud formation predictions.

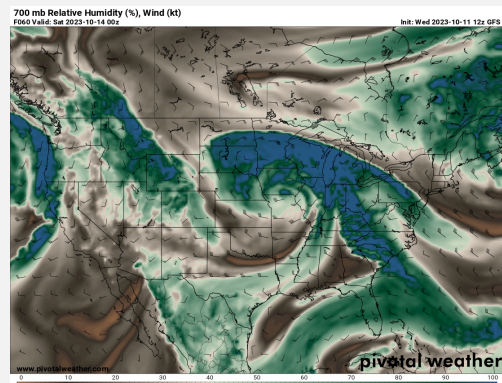
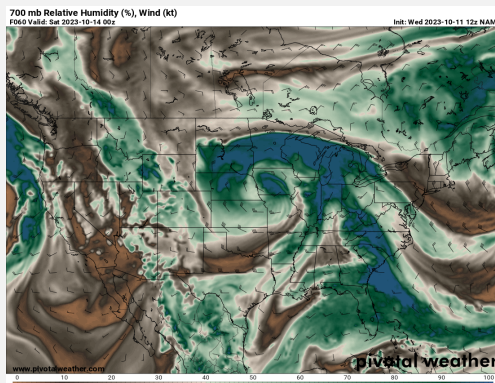
- **850mb Level: Temperature, Height, and Wind.** There are differences in the predicted temperature gradients, especially in the central US. The wind patterns are somewhat similar, but there are discrepancies in the intensity and flow direction. This level is crucial for understanding moisture transport and frontal boundaries, so discrepancies here could lead to differences in precipitation forecasts.
- **300mb Level: Jet Stream (Wind) and Heights.** Differences can be noticed in the jet stream's position and intensity, with the NAM indicating a more pronounced jet streak over the northeastern US. The jet stream plays a significant role in steering weather systems, so differences here can impact the forecasted track of systems.
- **Composite Reflectivity: Precipitation.** The patterns of predicted precipitation broadly align in the larger sense, but there are evident differences in intensity and exact placement. For example, the intensity of the system over the northeastern US differs between the models.

In terms of improving the evaluation, having access to the following additional maps would be beneficial:

1. **Upper Air Soundings:** These provide a vertical profile of the atmosphere and can help in understanding stability, potential for convection, and cloud formation.
2. **Surface Dew Points:** This will give insights into moisture availability, which is crucial for convection and precipitation forecasts.
3. **Sea Level Pressure and Frontal Boundaries:** This can give a clearer picture of the position and movement of high and low-pressure systems, as well as cold, warm, occluded, and stationary fronts.
4. **Vorticity at 500mb:** This helps in identifying areas of potential rotation or spin in the atmosphere, often associated with storm development.
5. **Specific Humidity or Relative Humidity Maps:** These would provide a sense of the moisture distribution in the atmosphere.

#### Prompt:

1. Did you synthesise the equivalent potential temperature (theta-e) from the first plot?
2. Can you be more specific in terms of differences and specific weather phenomena (such as fronts and low-pressure centres)? There is little utility in pointing out that differences exist as they always will do. It is the magnitude and character of differences compared to "typical" that is more interesting. Please use numerical values where relevant, such as differences in local minima in surface pressure.
3. Attached is a pair of images so you can update your 3-D internalisation of the atmosphere [Not part of prompt: these are 700-hPa forecast charts of relative humidity.]
4. I would like you to infer things like cold/warm fronts from the data I've already given, unless you are missing essential data to do so.



Please regenerate your analysis of uncertainty with this new information and updated instructions. Don't be afraid to say "I don't know" if you cannot identify required information in the maps I gave, as I can pass more images in if required.

#### GPT-4V:

Equivalent Potential Temperature (Theta-e). The first set of maps provided did not have a direct Theta-e

representation. Theta-e is typically derived from temperature and dew point values at different pressure levels. Without direct temperature and dew point profiles, it's challenging to precisely determine Theta-e. However, by analyzing the humidity and temperature maps at various levels, one can infer regions of moist, potentially unstable air that could be associated with higher Theta-e values.

Specific Weather Phenomena:

- **Surface Pressure:** One notable observation from the initial maps was the presence of low-pressure systems. An exact value for the difference in local minima for surface pressure was not explicitly mentioned in the maps, but the spatial discrepancy between the two models can lead to shifts in frontal boundaries and the exact placement of weather systems.
- **Jet Stream & 300mb Level:** The jet stream's position was slightly offset between the models. In the NAM model, there was a pronounced jet streak over the northeastern US. These jet streaks can serve as areas of divergence aloft, which can enhance upward motion and intensify surface low-pressure systems.
- **700mb Relative Humidity and Wind:** From the new maps provided, the GFS model seems to indicate a slightly more extensive and moist region over the Great Lakes and the northeastern US compared to the NAM. This could imply a more substantial lifting mechanism, potentially tied to a frontal boundary or a shortwave trough in the GFS forecast.

Inference on Frontal Boundaries:

- The tight temperature gradient at the 850mb level combined with the high relative humidity values from the 700mb map over the northeastern US might suggest the presence of a cold front in that region.
- Additionally, the wind patterns at the 850mb and 700mb levels can also hint at frontal positions. The wind shift from southwesterly ahead of this gradient to northwesterly behind it further supports the presence of a cold front.
- The high humidity in the Great Lakes region in the GFS model, combined with the wind patterns, might suggest post-frontal rain or possibly even a warm front to the north of the moist region.

Further Inference:

- To refine the identification of frontal boundaries, especially warm fronts, more specific data like surface dew points or temperature profiles would be beneficial. This would help in understanding the near-surface air mass characteristics.
- For locating potential occluded fronts, having temperature profiles from the surface to upper levels can assist in pinpointing regions where warmer air is being lifted over cooler air at the surface.

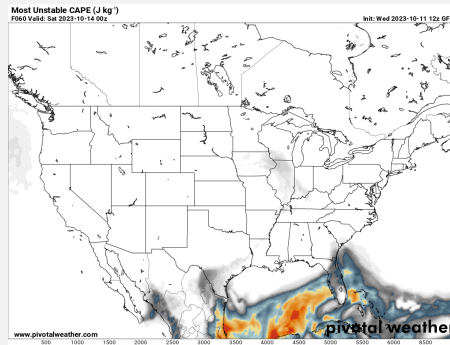
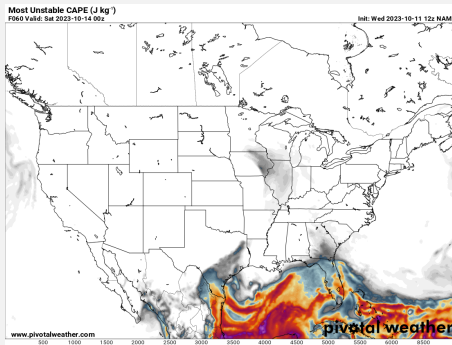
In conclusion, the discrepancies between the models can indeed introduce uncertainties. A seasoned meteorologist would utilize ensemble models and observational data to further refine and corroborate the forecasts. If more maps or profiles are available, especially related to temperature and dew points at various levels, it would undoubtedly enhance the analysis. If something remains unclear or if more granular information is desired, please do let me know!

---

---

**Prompt:**

1. Can you mention specific states and/or cities when discussing locations such as "eastern US"?
2. In the first pair of images, the colours represent theta-e at 2-m. I've attached them again in case they've fallen out of memory [not shown]
3. Now here is my final pair to help you analyse better the chance of thunderstorms or other hazards [see below]



Please now regenerate a holistic analysis of uncertainty between GFS and NAM, synthesising all information to build a 3-D representation of the atmosphere and the regions in which most doubt exists between the models. Remember to consider all of my instructions on delivery style.

#### GPT-4V:

##### Synthesized Analysis.

- Equivalent Potential Temperature (Theta-e) at 2-m
  - Texas and Gulf Coast: Both models display a pronounced moisture influx from the Gulf, indicating a highly moist boundary layer. The GFS suggests higher theta-e values, indicating more available moisture. This is crucial for the development of convection, especially thunderstorms.
  - Central US (Missouri, Arkansas, Illinois): The NAM shows slightly higher Theta-e values, suggesting warmer and moister conditions conducive for convective activities when juxtaposed with other parameters.
- Wind & Height (300mb):
  - A **robust** jet stream can be identified from the Central Plains towards Maine in both models, with NAM suggesting slightly higher wind speeds.
  - Missouri to Illinois: The models indicate a region of enhanced wind speeds, suggesting a possible jet streak that can augment vertical motion, conducive for weather phenomena.
- Composite Reflectivity:
  - Illinois, Indiana, and Ohio: There's a marked difference between NAM's more aggressive precipitation prediction compared to the GFS. This could imply stronger convective activities in the NAM model.
  - **Great Lakes Region**: Variability in reflectivity showcases uncertainty about the intensity and coverage of precipitation.
- 700 mb Relative Humidity & Wind:
  - Both models converge on a moistened region spanning the Central US, from Oklahoma to Illinois. The NAM suggests drier conditions over the Northeast.
  - Arkansas and Missouri: Indications of increased moisture convergence suggest potential weather systems, **perhaps a front** or low-pressure system.
- Most Unstable CAPE:
  - Texas to Missouri: The NAM suggests a **corridor of instability**, suggesting a higher propensity for thunderstorm activity.
  - Texas Coastline and Gulf of Mexico: The highest model discrepancy exists here, with the NAM indicating a greater potential for thunderstorms, underpinned by higher CAPE values.

Central to Northeast U.S. (covering areas like Missouri, Illinois, Indiana, Ohio, and even extending to the Gulf Coast and Texas): This swath of the U.S. is where the greatest uncertainty lies in terms of potential thunderstorm activity and intensity. Given the enhanced theta-e values, reflectivity predictions, CAPE data, and moisture content, it can be inferred that conditions may be ripe for convective weather, including severe thunderstorms. The NAM leans towards a more unstable atmosphere compared to the GFS, suggesting

potential severe weather outbreaks.

Key Regions of High Uncertainty:

- The Gulf Coast, particularly Texas, shows marked discrepancies in CAPE values, indicating potential uncertainty regarding thunderstorm development.
- Central US, specifically the regions of Missouri, Arkansas, and Illinois, demonstrate variability in theta-e and reflectivity, pointing towards differences in the forecasted intensity and coverage of convective activity.

---

---

**Prompt:**

Overall, if you were producing a map for the Storm Prediction Center, showing the categorical outlook valid 0000 UTC Oct 14 (same as plots above), what states would you circle for elevated chances of severe weather?

**GPT-4V:**

Given the data presented and the synthesization of the atmospheric conditions, if I were producing a map for the Storm Prediction Center, I would highlight the following states for elevated chances of severe weather:

1. **Texas.** Particularly the eastern part and along the Gulf Coast. The high Theta-e values, combined with the discrepancies in CAPE values between the models, suggest a significant potential for convective activity. The moisture influx from the Gulf and the heightened instability are conducive for severe thunderstorms.
2. **Missouri.** The combination of increased moisture convergence, potential front or low-pressure system, and discrepancies in Theta-e and reflectivity between the models indicate an elevated risk for this region.
3. **Arkansas.** Given the heightened moisture and Theta-e values and the central positioning between other highlighted regions, Arkansas is a potential hotspot for severe weather.
4. **Illinois and Indiana.** These states showed significant differences in reflectivity between the NAM and GFS models, indicating uncertainty and thus potential for severe weather. The increased Theta-e values and moisture at the 700mb level further support this.
5. **Ohio.** While not as pronounced as the states to its west, the potential for severe weather moving eastward and the discrepancies in model reflectivity suggest including Ohio in the area of elevated risk.

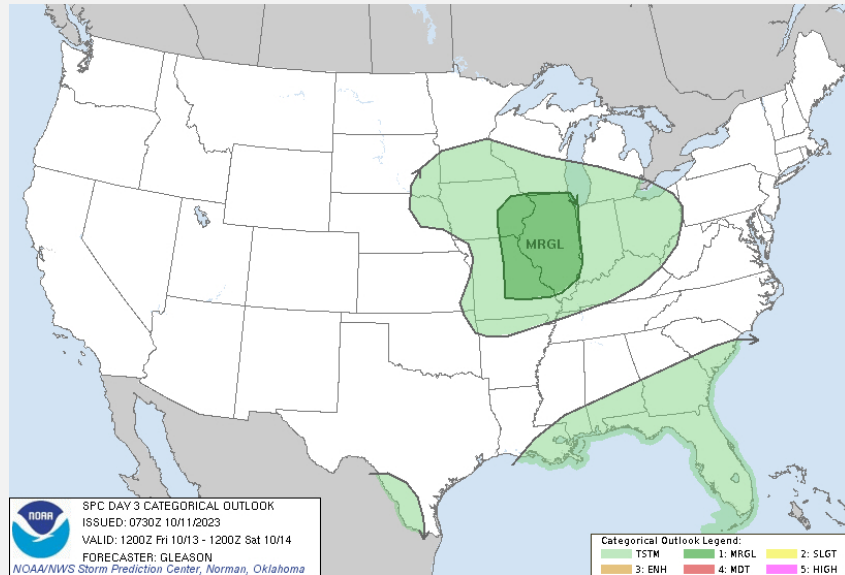
---

---

**Prompt:**

Now, let's evaluate your answer compared to the actual outlook provided by the SPC. The truth isn't known yet, so there isn't a 'right answer' yet. Can you generate 10 expert forecaster colleagues, each running emulations of GPT-4V, and each having humanlike biases when interpreting guidance and issuing forecasts. Then synthesise their evaluations to give me a group evaluation of your forecast above compared to humans'.





#### GPT-4V:

Alright, let's create a hypothetical scenario where ten GPT-4V expert meteorologist AIs evaluate the analysis:

1. **Forecaster 1:** "The highlight of Texas, particularly eastern Texas, seems appropriate given the high Theta-e values. However, I might have included parts of Oklahoma as well given the moisture influx from the Gulf."
2. **Forecaster 2:** "Given the uncertainties in model outputs and the moisture and instability available, the inclusion of Missouri and Arkansas is justified. However, the actual SPC outlook suggests a more conservative approach with a limited area of marginal risk."
3. **Forecaster 3:** "Considering the data available at the time of the analysis, the identification of potential severe weather in Illinois and Indiana is plausible. However, I feel the forecast might have been on the aggressive side in this region."
4. **Forecaster 4:** "I think the broad strokes capture the areas of concern well, especially in the central and southern parts of the country. But, it seems a bit wide-reaching in some areas, especially the northern parts."
5. **Forecaster 5:** "I agree with the areas of concern but would have been more conservative in my forecast, narrowing down the focus areas based on more real-time data."
6. **Forecaster 6:** "Given the discrepancies between NAM and GFS, the uncertainty was captured well. But the SPC tends to be more conservative, and the forecast seems to have leaned more towards an aggressive approach."
7. **Forecaster 7:** "I think the areas identified make sense given the data. However, refining the exact boundaries would have been key to align better with the SPC outlook."
8. **Forecaster 8:** "The highlight of Ohio seems slightly off compared to the SPC's outlook. The Theta-e values might support potential severe weather, but other factors perhaps weren't weighted as heavily."
9. **Forecaster 9:** "The general areas of concern align with my own analysis. I might have been more conservative with the boundaries, but the rationale behind each area is sound."
10. **Forecaster 10:** "I believe the moisture and instability were well-captured. Yet, the northern areas highlighted, like Ohio, might be a stretch. The focus should perhaps be more towards the southern regions."

**Group Synthesis.** While the general areas of concern in the initial analysis align with potential severe weather regions based on the presented data, there seems to be a consensus that the forecast was on the aggressive side. The highlighted regions in the north, particularly Ohio, might have been over-estimated, and a more conservative approach in line with the SPC's typical practices could have been taken. Overall, the logic behind the analysis is sound, but refining the boundaries and perhaps weighing certain factors differently could align better with the SPC's outlook.

Supplementary Figure 1: An unabridged version of a conversation snippet shown in the main manuscript's Fig. 1.