



WeatherBench ≥ 2

What's next for AI-weather models and evaluation?

Stephan Rasp, Google Research

With help from many others:

Stephan Hoyer, Peter Battaglia, Alex Merose, Ian Langmore, Tyler Russell, Alvaro Sanchez, Rob Carver, Vivian Yang, Shreya Agrawal, Carla Bromberg, Jared Sisk, Luke Barrington, Aaron Bell, and Fei Sha (Google)

Matthew Chantry, Zied Ben Bouallegue and Peter Dueben (ECMWF)

Google Research



WeatherBench 2 - Status Quo

WeatherBench 2 is a benchmark for **global, medium -range weather prediction** .

It consists of:

1. **Data** : Relevant data freely available as Zarr on GCS (ERA5, IFS HRES and ENS, ML forecasts).
2. **Code** : Parallelizable and reproducible evaluation code on GitHub.
3. **Website** : Up-to-date platform showing state-of-the-art of AI-weather modeling.

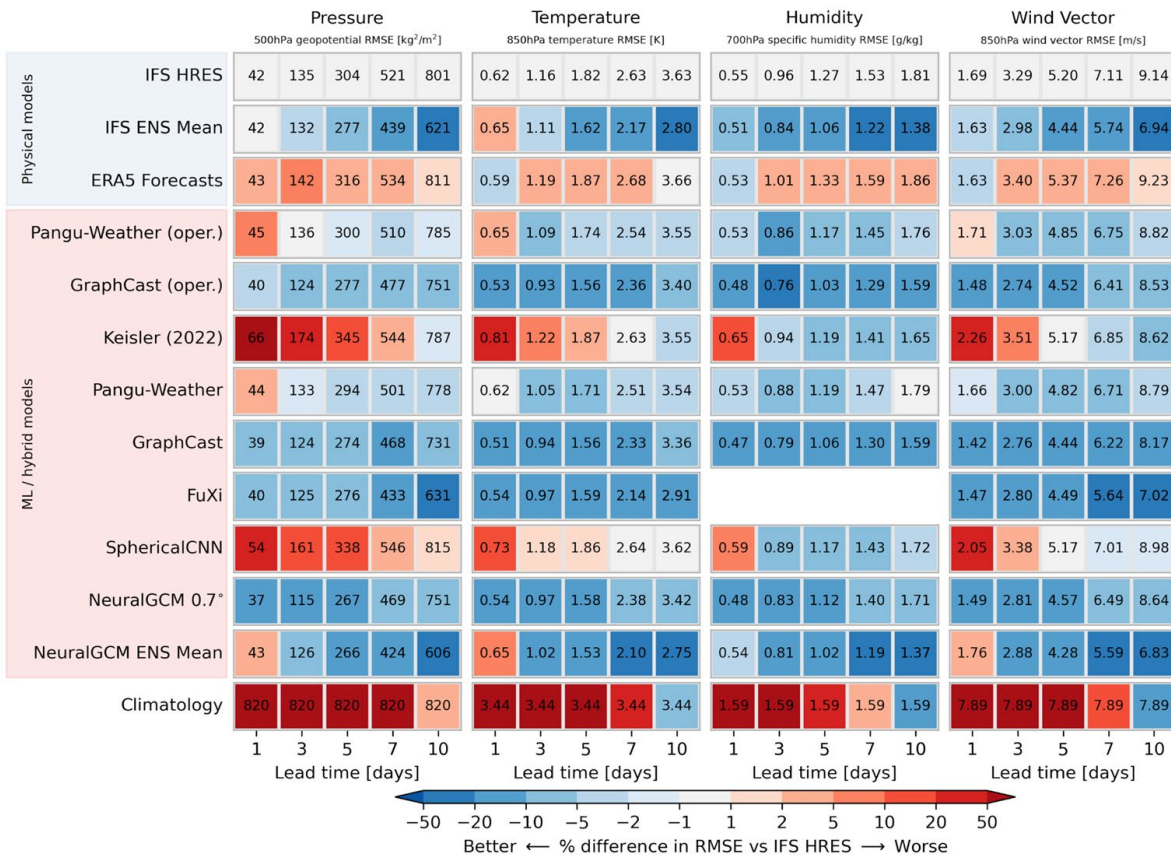
For background information, check out our **paper** (Arxiv, soon to be published in JAMES).

For technical information, visit the **GitHub page** and the **documentation** .



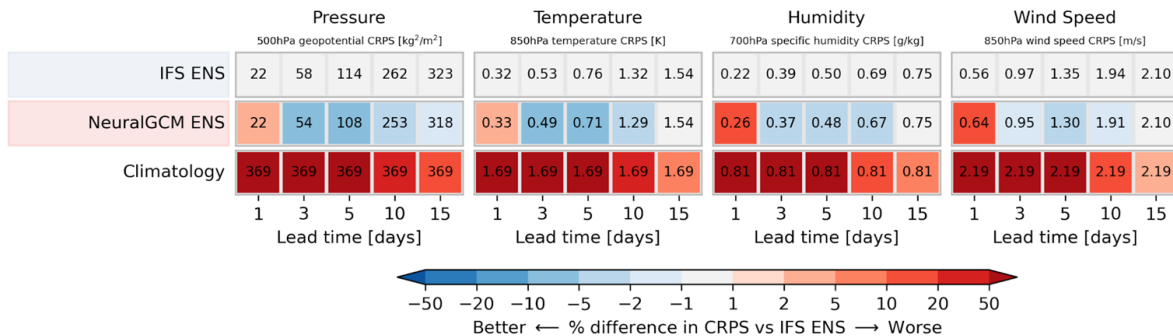
Lessons from the leaderboard

- ML models are roughly on par with physics-based models.
- Many deterministic ML models blur. Spectra are somewhere in-between HRES and ENS.
- Overall, most ML models have similar performance. ERA5 seems to be the main limiting factor.





Lessons from the leaderboard



- Less progress on ensemble methods.
- Existing models roughly on par with IFS ENS.

+ GenCast (and Pangu ensemble)



From research to operations

ERA5 is not available for initialization in real time.

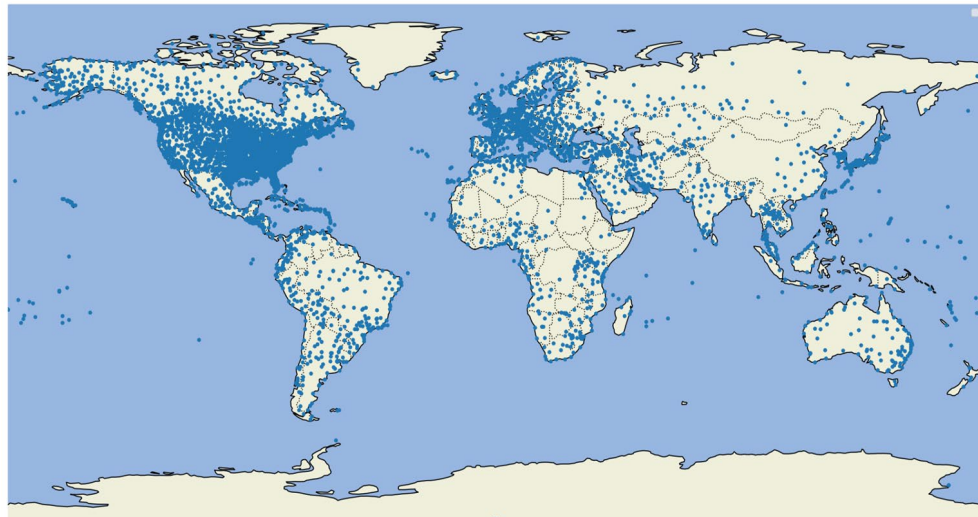
ERA5 is often not the best ground truth for impactful weather (see ECMWF’s “lower half of the scorecard”).

Many real-world applications require post-processing to higher-quality datasets.

Forecast latency matters.

→ The next step in evaluation: operational conditions and “best” ground truth.

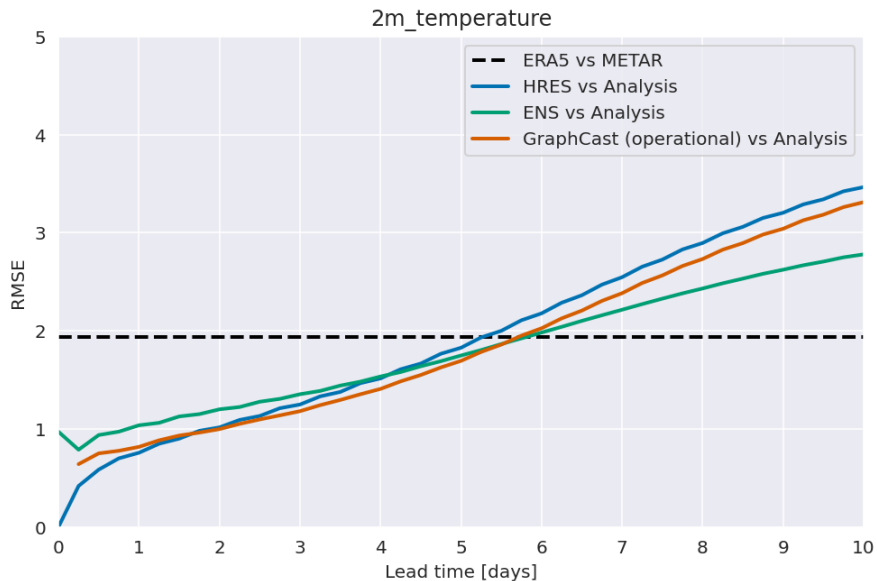
Station evaluation



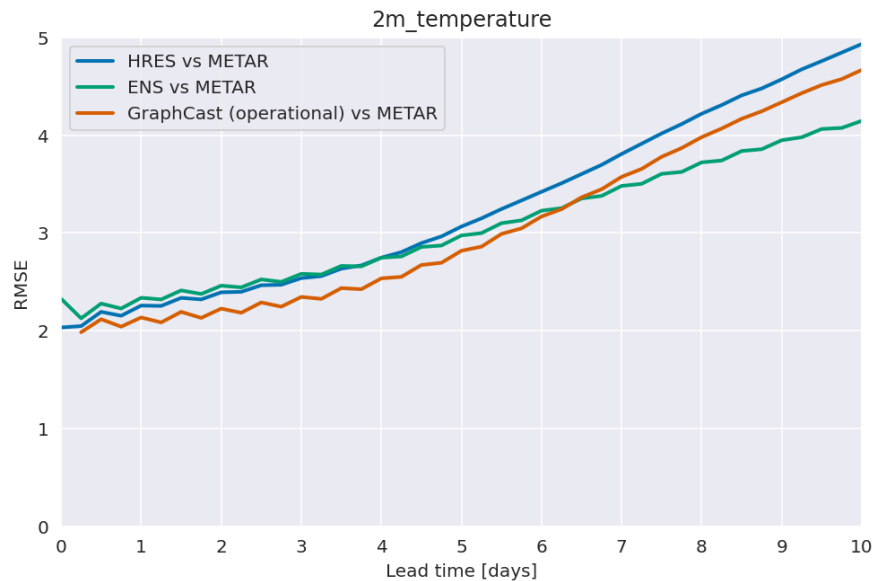
- Evaluation against ~5000 METAR stations.
- All 00/12 initializations for 2020.
- Gridded fields are bilinearly interpolated to station locations.



Station evaluation



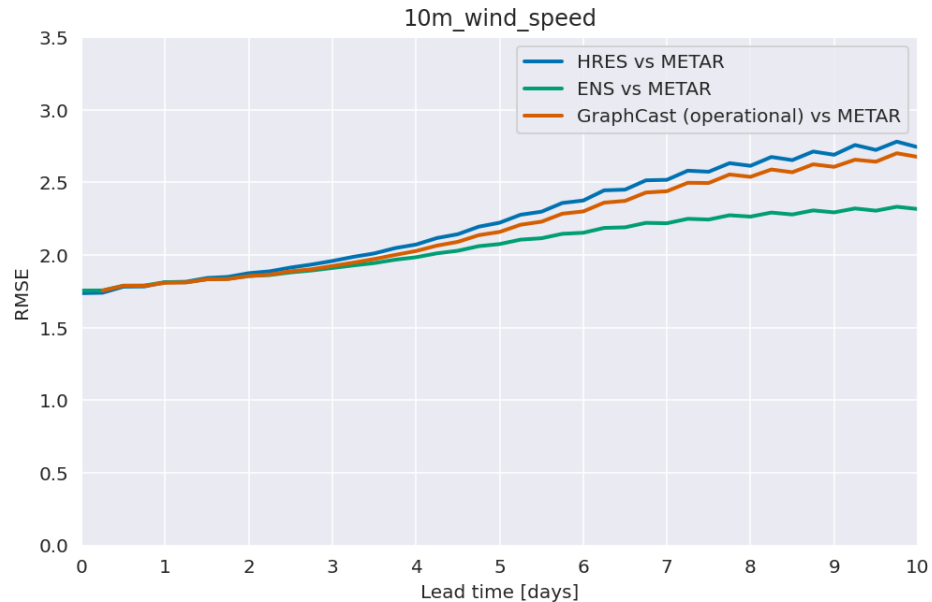
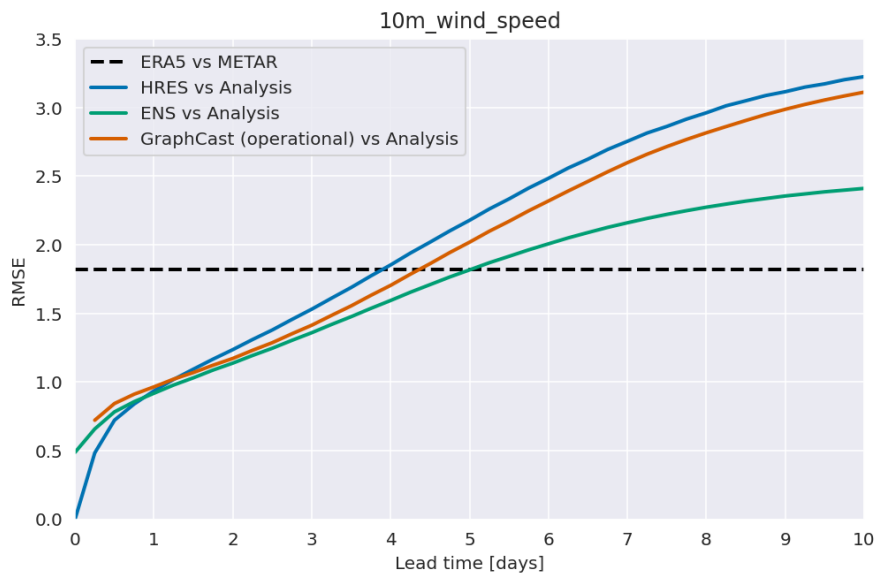
→ ERA5 error \approx 5 day forecast error.



→ Relative score of models largely unchanged.



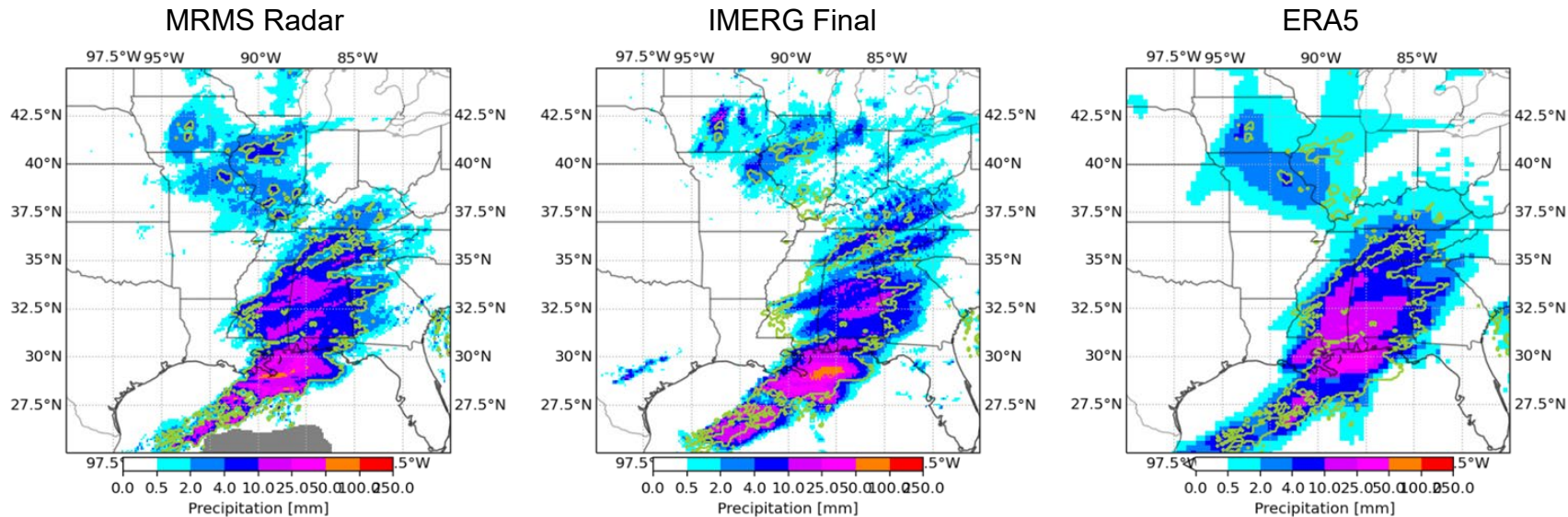
Station evaluation



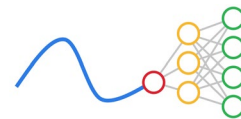
→ Same applies to wind speed.



Precipitation

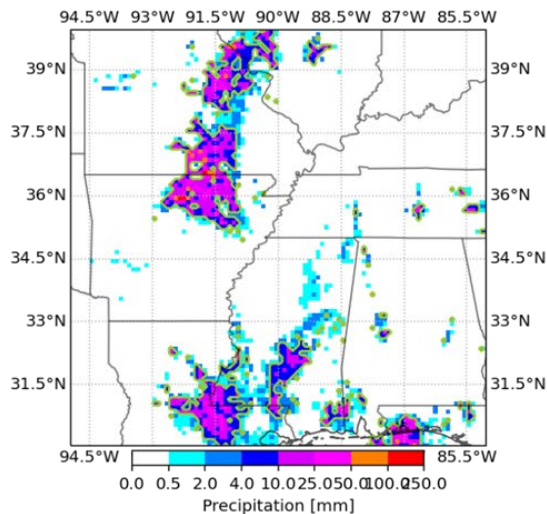


- No single “best” precipitation ground truth.
- Rain gauges are sparse and noisy.
- Radar derived products (e.g. MRMS in the US) are only regional.

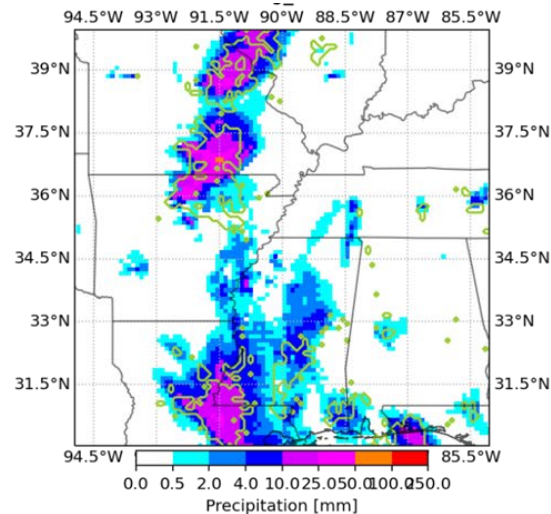


Precipitation

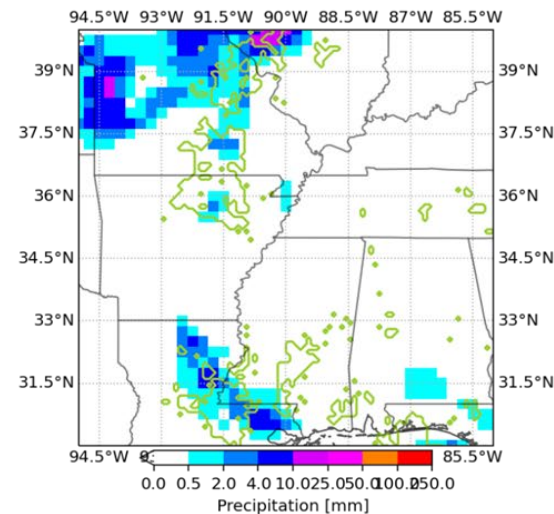
MRMS Radar



IMERG Final



ERA5



- IMERG (and other satellite derived products) are global but not perfectly accurate (CSI_{4mm/6hr} \approx 0.4 for IMERG vs 0.35 for ERA5).
- No shortcut to evaluating against a range of “ground truths”.



An observation benchmark?

Next-gen ML models will likely be trained directly against observations.

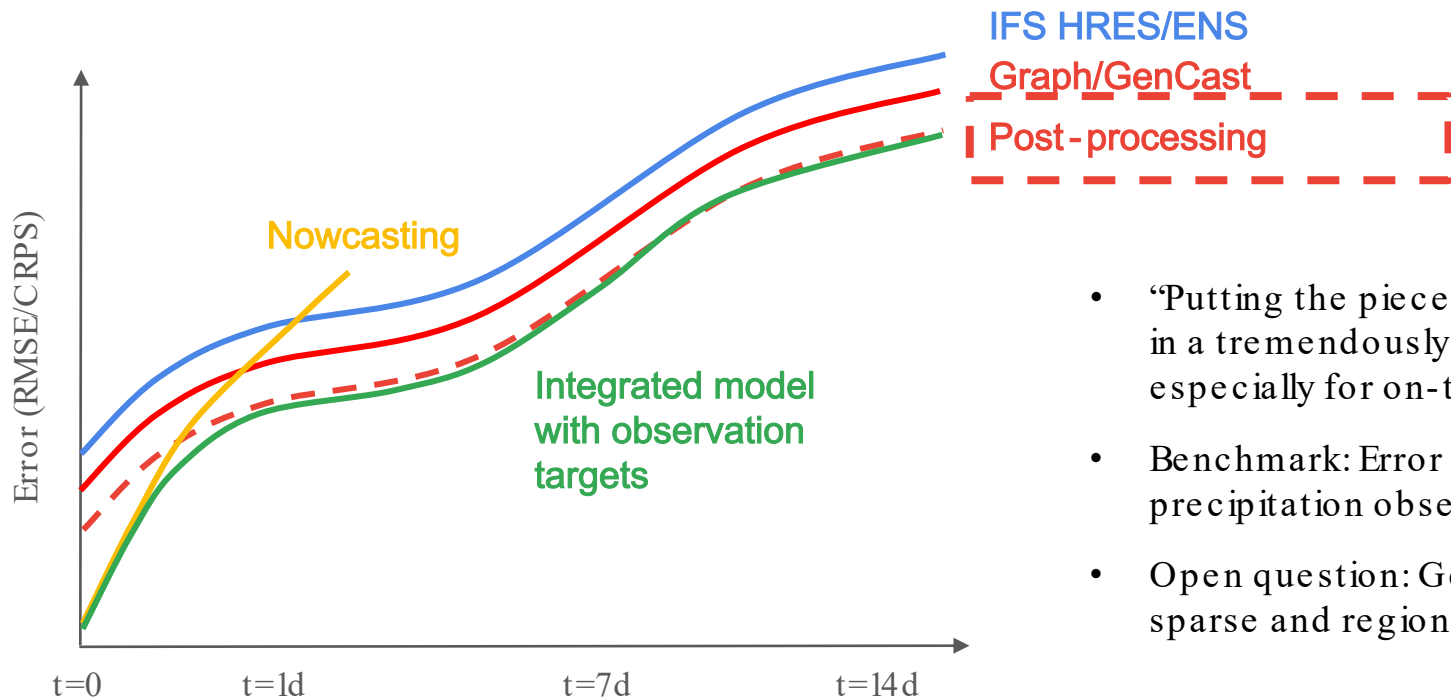
Therefore, WeatherBench 3(?) should be an observation benchmark but ...

- Is there agreement on the “best” ground truth? Especially, for precipitation?
- Sparse observations (e.g. weather stations) require generalization. Therefore, we need a hold-out set of stations (and agree what that would be).
- What should the test period be? Many high-quality observations only available for recent years.
- How can we compare against the current “state-of-the-art”, i.e. commercial forecast providers?

Please let me know: What should an observation benchmark look like?



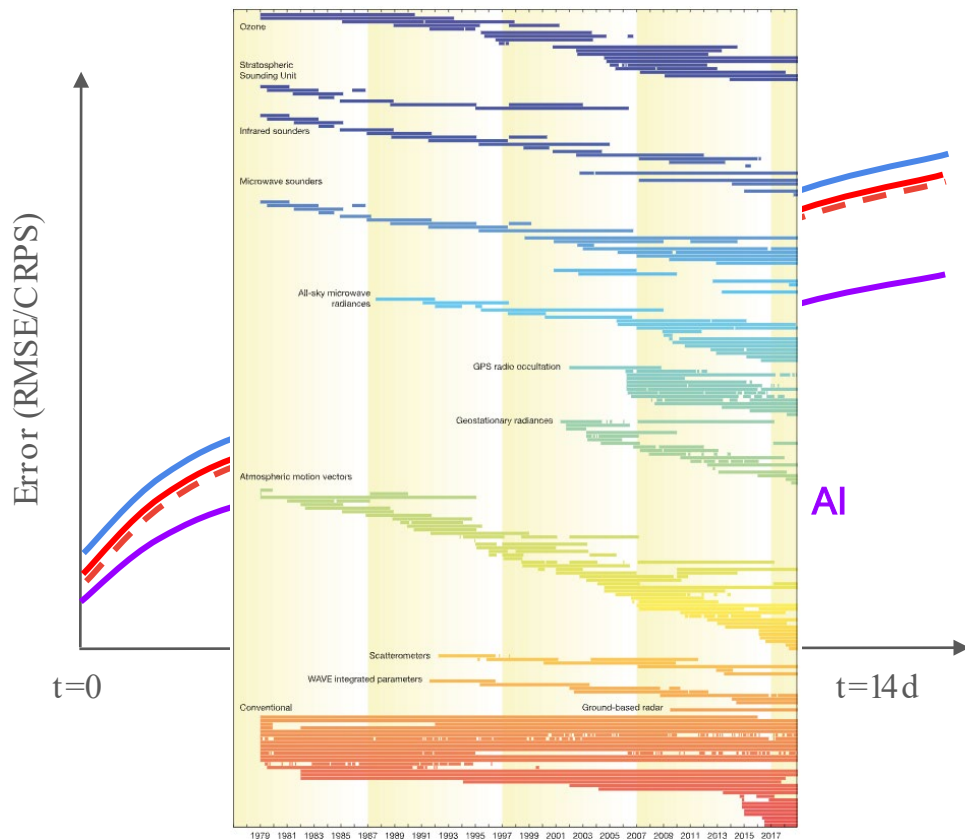
The grand challenge for AI models



- “Putting the pieces together” can result in a tremendously useful model, especially for on-the-ground weather.
- Benchmark: Error vs surface and precipitation observations
- Open question: Generalization of sparse and regional observations



The grand challenge for AI models



- Post-processing/nowcasting has little impact on large scale.
- >50% of potential improvements in initial conditions.
- Challenge: Exploit existing observations to improve ICs and large-scale forecasts.
- Benchmark: Z500, TC track and intensity, etc.
- Requires significant investment in data infrastructure.