# A Weather Foundation Model for the Power Grid

Cristian Bodnar[1], Raphaël Rousseau-Rizzi[2], Nikhil Shankar[1], James Merleau[2],
Stylianos Flampouris[1], Guillem Candille[2], Slavica Antic[2],
François Miralles[2], Jayesh K. Gupta[*,1]

[1] *Silurian AI*        [2] *Hydro-Québec*

[*] *corresponding author:* `jayesh@silurian.ai`

## Abstract

Weather foundation models (WFMs) have recently set new benchmarks in global forecast skill, yet their concrete value for the weather-sensitive infrastructure that powers modern society remains largely unexplored. In this study, we fine-tune Silurian AI's 1.5B-parameter WFM, Generative Forecasting Transformer (GFT), on a rich archive of Hydro-Québec asset observations—including transmission-line weather stations, wind-farm met-mast streams, and icing sensors—to deliver hyper-local, asset-level forecasts for five grid-critical variables: surface temperature, precipitation, hub-height wind speed, wind-turbine icing risk, and rime-ice accretion on overhead conductors. Across 6–72 h lead times, the tailored model surpasses state-of-the-art NWP benchmarks, trimming temperature mean absolute error (MAE) by 15%, total-precipitation MAE by 35%, and lowering wind speed MAE by 15%. Most importantly, it attains an average precision score of 0.72 for day-ahead rime-ice detection, a capability absent from existing operational systems, which affords several hours of actionable warning for potentially catastrophic outage events. These results show that WFMs, when post-trained with small amounts of high-fidelity utility data, can serve as a practical foundation for next-generation grid-resilience intelligence.

## 1   Introduction

Modern electricity networks are becoming markedly more weather sensitive. Rapid penetration of wind and solar introduces variability that is tightly coupled to mesoscale meteorology, while ageing transmission and distribution grids face mounting exposure to storms, icing, and heatwaves. Grid operators must now balance supply and demand in real time in ever more reactive way, safeguard assets against extreme events, and plan maintenance in a climate where "once-in-a-century" conditions recur with unsettling frequency. At the same time, the decarbonization of the economy increases social vulnerability to power outages, enhancing the need for reliability. Traditional numerical weather prediction (NWP) feeds remain indispensable, yet their spatial granularity (often $\geq$ 3 km) limits their usefulness for hyper-local, asset-level decisions such as dynamic line rating (IEEE Power and Energy Society, 2013), de-icing crew dispatch, or curtailment of turbine fleets in icing conditions.

Addressing these granular, asset-level challenges requires a new forecasting paradigm, one offered by recent advances in the world of large scale machine learning. The foundation-model paradigm of pretraining large transformers on petabyte-scale data, followed by fine-tuning for downstream tasks, has found great success across a variety of disciplines (Bommasani et al., 2021; Brown et al., 2020; OpenAI, 2023; Radford et al., 2021; Kirillov et al., 2023; Jumper et al., 2021; Rives et al., 2021; Wang et al., 2023). Weather foundation models (WFMs) apply this paradigm to the task of atmospheric prediction, leveraging vast reanalysis, satellite, and other climate data archives. By learning flow-consistent latent representations across lead times and scales, WFMs have surpassed operational NWP in global skill metrics while offering quick inference on commodity GPUs. Crucially, their parameter sharing and attention mechanisms allow regional adaptation with orders-of-magnitude fewer labelled samples than would be required to train a model from scratch.

While the AI weather models now match or surpass traditional operational NWP models for conventional forecasting tasks at broad, continent-spanning scales Bodnar et al. (2025); Chantry et al. (2025); Lam et al. (2023); Bi et al. (2023); Pathak et al. (2022); Chen et al. (2023), this work demonstrates a pivotal new capability brought about by WFM's regional adaptation capabilities: forecast products tailored to user needs: *both* flexibly and rapidly. This is a difficult endeavor using existing NWP models. This makes WFM technology invaluable for utility companies since the grid isn't much impacted by "averaged" and generic weather variables, but rather by the ways the weather interact with assets right here, right now: on this ridge, over these wind turbines, along this particular stretch of transmission line. More precisely, grid operation requires forecasts that can meet three challenging requirements at once:

- **Hyper-locality**: The forecast is sensitive enough to capture kilometer-scale pecularities like terrain, land-use, and other environmental factors in the vicinity.

- **Specificity**: Highly task-specific variables such as thermal ratings, renewable output, and icing risk can be forecast.

- **Actionability**: The forecast is timely, clear and expressed in terms that facilitate real-time decision making by grid operators (e.g., probability to exceed a critical threshold along with information on the false alarm rate).

WFMs are uniquely suited to meet these challenges, yet their potential remains largely untapped. This is not due to a failure of the technology, but rather a gap in awareness and a lack of established pathways for integrating their highly specific data into grid operations. This work aims to bridge this gap, showcasing how WFM technology can empower operators and build a more resilient energy infrastructure.

To illustrate WFM capabilities for the grid, in this paper we adapt Silurian's 1.5-billion-parameter WFM, Generative Forecasting Transformer (GFT) to Hydro-Québec's asset network, producing hyper-local hourly forecasts for five grid-critical variables: surface temperature, accumulated precipitation, hub-height wind speed, wind-turbine icing risk, and rime-ice accretion on overhead conductors. Against state-of-the-art NWP baselines over 1–72 h horizons, the fine-tuned model

1. Lowers temperature MAE by 15%,

2. Reduces total-precipitation MAE by 35%,

3. Reduces wind speed MAE by 15%, and

4. Achieves an unprecedented 0.72 average precision for day-ahead rime-ice detection, providing several additional hours of warning compared to existing process.

These gains translate into tangible operational benefits: earlier de-icing interventions, more reliable dynamic line ratings, and improved renewable dispatch (see Section A.3). The study thus demonstrates that WFMs, when enriched with utility-grade observations, can become a cornerstone of next-generation grid-resilience intelligence.

## 2 A Weather FM for the Grid

**Weather Foundation Models.** Recent foundation-model efforts such as ClimaX (Nguyen et al., 2023) and Aurora (Bodnar et al., 2025) have established that a single transformer-based backbone, pre-trained on heterogeneous climate data archives, can surpass global numerical weather prediction (NWP) skill. They also demonstrate how these models can be efficiently adapted for a wide variety of tasks in this space, ranging from wave modeling and air-pollution forecasting to regional downscaling. An overview of our WFM setup for grid applications is shown in Figure 1. In this work, we show how Silurian's foundation model, the Generative Forecasting Transformer (GFT), inheriting the same principles, can be adapted for a variety of energy infrastructure use cases that demand hyperlocal spatial detail and complex weather risk assessments.
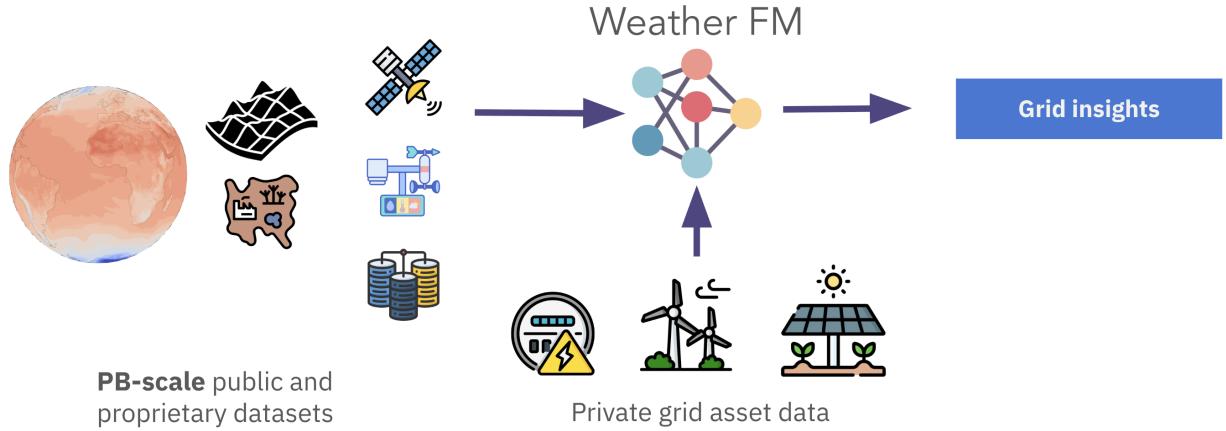
Figure 1: Weather foundation model for the power grid.

**Architecture Overview.** The model follows the general encoder–backbone–decoder paradigm (Figure 1). The encoder ingests operational weather inputs such as gridded and sparse observations and maps them into a latent representation. The backbone consists of a large transformer that acts as a neural simulator of the latent representation. Finally, the decoder translates the evolved latent representations into physical predictions. The GFT decoder contains two modules: a "dense" decoder and a "sparse" hyperlocal decoder. The dense decoder produces gridded outputs on a predefined coordinate grid, while the sparse decoder produces forecasts at a specific set of query latitude–longitude locations.

**Pretraining and Post-training.** The model is pretrained on a petabyte-scale dataset consisting of a mix of public and private data, with dense/gridded observations forming the bulk of this dataset. The general recipe is similar to the approach given in Bodnar et al. (2025). After pretraining on this large dataset, we post-train the model on a much smaller set of hyperlocal observations from sensors across the Hydro-Québec grid infrastructure.

**Post-training versus post-processing.** Utilities typically adapt fixed NWP guidance with Model Output Statistics, ensemble calibration, or site-specific regressors that operate one variable and one location at a time. Because the upstream flow solution never changes, these post-processing pipelines struggle to maintain multivariate consistency, cannot introduce new targets absent from the NWP feed, and scale poorly as the asset portfolio grows. Our post-training updates the core WFM weights so that a single backbone jointly forecasts temperature, wind, precipitation, and icing. The shared latent state enforces cross-variable coherence, enables the model to learn novel heads such as rime-ice probability directly from observations, and eliminates the need to maintain dozens of bespoke statistical correctors.

**Hydro-Québec post-training setup.** At inference time, the encoder ingests ECMWF-IFS analysis fields (ECMWF, 2021) and the sparse Hydro-Québec asset observations collected at the cycle issue time. The decoder is trained to produce five operational targets in one forward pass: 2 m temperature, hourly precipitation, hub-height wind speed, wind-farm icing probability, and rime-ice probability on transmission lines. We fine-tune on 2016–2023 observations and hold out January 2024–March 2025 for evaluation, balancing regression losses for the continuous variables with probabilistic classification losses for the icing heads. Unless noted otherwise, utility streams provide supervision only; they are not assimilated online.

## 3 Hyperlocal forecasting for grid operations

To demonstrate the flexibility of the foundation-model paradigm, we post-train GFT into GFT-HQ, on three important environmental tasks that have a significant impact on grid operations: (1) rime-ice forecasting

for transmission lines, (2) wind-farm wind and icing forecasting, and (3) temperature and precipitation forecasting.

Figure 2 illustrates Hydro-Québec's existing high-voltage network and corridors under study for future expansion, compiled from public planning materials and reports (Hydro-Québec, 2023a;b; 2024; 2025). These large-scale buildouts and prospective wind additions motivate asset-level forecasts that are fast, location-specific, and actionable.

There are two ways hyperlocal WFM forecast products can improve grid operations. The first is to improve the skill of forecast products already in use within the utility, while retaining a similar format and role in existing operational processes. This is the case for the temperature, precipitation, and wind forecasts in the present study. The second is to meet an operational need that is not currently addressed by existing structured process in the utility, as is the case for ahead-of-time transmission-line and wind-farm rime-ice forecasts. The effort required to integrate each product–and the benefits to the utility–differ considerably between these two types. Contributions of the first type are relatively easy to integrate into the value chain, and their benefit is determined by the improvement over available forecast products. Contributions of the second type are harder to integrate because, by definition, they require the creation of new activities within the company. The potential benefit of WFMs is much greater here, as it also comes from a qualitative change to the value chain—such as being able to pre-emptively deploy an operational team to de-ice lines ahead of the start of an event.

**Evaluation protocol.** Continuous targets (temperature, wind, precipitation) are scored with mean absolute error (MAE) and reported as fractional skill relative to ECMWF-IFS (ECMWF, 2021). Rare-event targets (icing) are assessed with precision, recall, F1, and average precision (PR-AUC), with lift computed against the relevant base rates: 3.7% for transmission-line rime ice and 13% for wind-farm icing. To reflect operational decision windows, we aggregate hourly probabilities into a 24 h "any icing" probability $q_t = 1 - \prod_h (1 - p_{t+h})$ alongside the hour-by-hour scores. ROC-AUC is provided for completeness but we emphasize PR metrics because they better reflect the cost of false dispatches in low-base-rate regimes.

### 3.1 Rime ice forecasting for transmission lines

Rime ice is formed when super-cooled water droplets freeze instantly upon contact with a sub-zero surface. It accumulates rapidly under fog, low-level cloud, and mountainous or coastal weather systems and poses a persistent reliability threat to overhead transmission lines. The ice accretion alters both the geometric and electrical characteristics of conductors: it increases effective diameter and surface roughness, thereby magnifying wind drag, modifying corona onset voltage, and adding eccentric mechanical loading. These translate into elevated risk of conductor clashing, flashovers, structural member fatigue, and ultimately line outages.

Despite the many operational headaches it causes, utilities lack reliable forecasting systems for rime ice. The gap traces chiefly to numerical weather prediction models: their kilometre-scale grids smooth out the shallow fog banks and ridge-top cloud filaments where rime forms, and their bulk microphysics schemes convert supercooled droplets to ice far too quickly, systematically erasing the liquid-water signal that drives accretion. With neither high-resolution physics nor a corridor-wide network of icing masts to assimilate or validate supercooled liquid water, model output remains both biased and unverifiable, leaving operators without a trustworthy baseline on which to build operational alerts.

**Data** For this task, we post-train GFT on hourly data collected from 14 high-elevation Sygivre rime-ice stations from Hydro-Québec's 40-station network; each measures ice-accretion cycles plus standard meteorological variables. See Section A.2 for details.

We convert the cumulative ice-accretion counts into a binary variable indicating the presence of rime ice within the last hour. In the resulting data, rime-ice events represent just 3.68% of total hours, underscoring the rarity of this phenomenon and the challenge it poses for any AI-based detection system. We train on 2016-2023 and validate on January 2024 to May 2025.
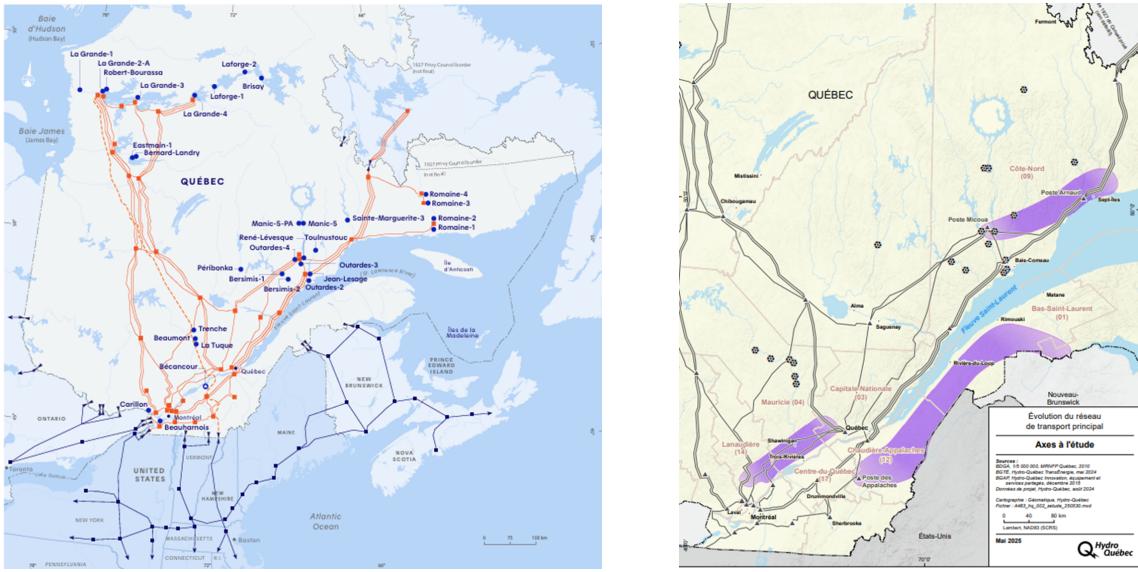
4

Figure 2: Hydro-Québec's major facilities and transmission line infrastructure: Current (Left), and planned (Right) expansion; compiled from public Hydro-Québec materials (Hydro-Québec, 2023a; 2025; 2023b; 2024).
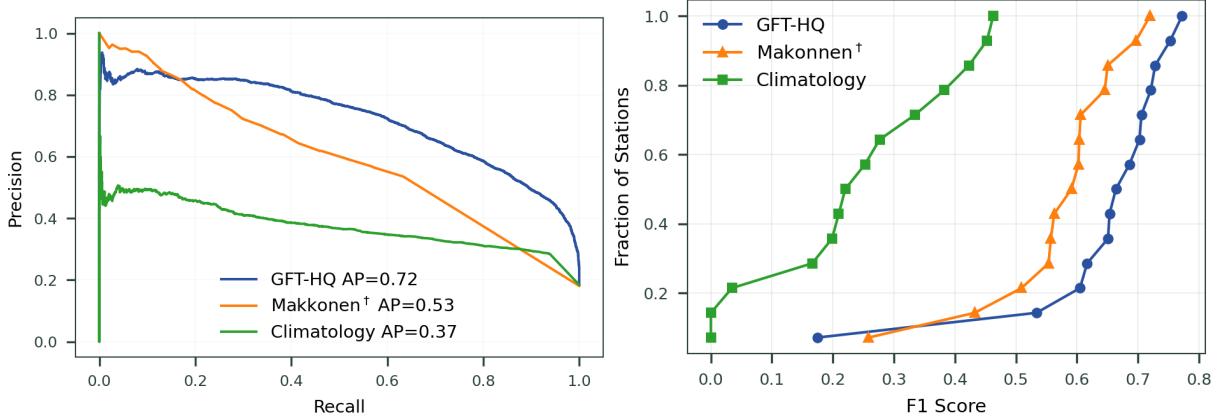


Figure 3: (a) Precision–recall curve for rime-ice forecast over the next 24 h forecasts. (b) Cumulative Distribution Function (CDF) of F1 scores of rime-ice over the next 24 h forecast over all stations. Note that Makkonen[†] is derived from ERA5 reanalysis and cannot be used operationally.

**Operational decision model** Dispatching de-icing crews carries a fixed cost $C_d$, while a missed icing event incurs a loss $L$ that can be partially avoided if action succeeds with effectiveness $\alpha$. Using a standard cost–loss analysis, we trigger action whenever the forecast probability exceeds

$$p^\star = \frac{C_d}{\alpha L}. \tag{1}$$

For windowed "any icing in the next 24 h" alerts we map hourly probabilities to $q_t = 1 - \prod_h (1 - p_{t+h})$ and apply the same rule, with hysteresis and watch/dispatch tiers to avoid thrashing in operations. This framing connects forecast quality directly to avoided outage risk and underpins the value analysis in Section A.3.

**Baselines** In addition to climatology, we compare against an ERA5-derived Makkonen index that interpolates reanalysis temperature, dewpoint, winds, and liquid-water content to each site and evaluates a physical

icing proxy (Hersbach et al., 2020). Because this proxy is generated from retrospective reanalysis rather than a true forecast, it serves as a strong physics-based reference rather than an operational baseline.

**Results**   Pinpointing the exact hour of icing is difficult, and hourly skill reflects this: GFT-HQ's F1 scores are typically in the 0.40–0.45 range during the first 24 h; see Appendix Figure 14a. Because operations often act on multi-hour windows, we also evaluate a 6-hour "any icing in window" target, which smooths timing errors and yields higher, more stable skill across lead times; see Appendix Figure 14b. In both settings, GFT-HQ provides the best performance and is comparable to or better than the ERA5-derived Makkonen reference out to 72 h. As can be seen in Section 3.1, GFT-HQ also detects the approximate icing period well. For instance, when forecasting whether any icing will occur over the next 24 h at the 14 locations of interest, GFT-HQ attains average precision AP = 0.72 ($\approx 8\times$ lift over base rate), outperforming an ERA5-derived Makkonen index (AP = 0.53; $\approx 6.1\times$ lift) and climatology (AP = 0.37; $\approx 3\times$ lift). Note that Makkonen[†] is computed from ERA5 reanalysis fields (see Section A.1.4; ERA5 described in Hersbach et al., 2020), so it is closer to observations than to a true forecast; we include it as a strong physics-based reference rather than an operational baseline. GFT-HQ exceeds this reference by $\approx 30\%$ relative AP (0.72 vs 0.53). Across 14 sites, station-level F1 scores are consistently higher for GFT-HQ, indicating improvements at the median and in the worst-case stations, not just at a few outliers.

**Case study**   Beyond aggregate skill, the finetuned model captured individual high-impact episodes that posed particular challenges for Hydro-Québec's operations. A prime example is the collapse of a transmission line from the Romaine hydro centre around November 19, 2024 due to rime ice (Gerbet, 2024). Figure 4 compares two nearby Sygivre stations (`ERIC_C` and `MONTAG_C`) separated by only 11 km using a Hovmöller-style view: rows are successive 6-hourly forecast initializations, columns are valid hours, colours show the forecast probability of 1-h rime ice, and the black strip below indicates the observed binary icing. At `ERIC_C`, a coherent high-probability swath is already present by 16 Nov and persists across cycles, peaking on 18–19 Nov when the longest observed episode occurs; subsequent bursts on 20–24 Nov are also captured. At `MONTAG_C`, icing was shorter and more intermittent, and the model response is correspondingly weaker— probabilities remain low most of the time with brief increases around 16–19 Nov. Despite their proximity, these stations exhibit markedly different behaviour, underscoring sharp micro-scale variability along the corridor; nevertheless the corridor-scale risk around 18–19 Nov was visible 1–3 days in advance, providing actionable lead time for grid operations.

**Initialization and inputs**   Unless otherwise stated, the case-study forecasts are produced by 6-hourly GFT-HQ cycles (issue times 00, 06, 12, 18 UTC). For each cycle, the encoder is conditioned on dense ECMWF-IFS Analysis fields (ECMWF, 2021) (0 h analyses) at the issue time—including multi-level temperature, humidity, and winds; surface pressure and near-surface fields; cloud liquid-water proxies; static orography and land–sea mask; and time encodings. These fields provide the initial conditions from which the backbone evolves the latent state forward in time. Hydro-Québec asset streams (Sygivre icing stations and wind-farm masts) are used for post-training supervision; unless explicitly noted, they are not assimilated at runtime during inference. The traces in Figure 4 therefore reflect differences solely from the changing synoptic initial state between successive 6-hourly initializations.

## 3.2   Wind-farm wind speed and ice forecasting

For utilities reliable forecasts of turbine-level icing and hub-height wind speed are mission-critical because they dictate how much variable generation will actually reach the grid, how much spinning reserve must be scheduled, and whether transmission lines risk overloads or voltage excursions when production suddenly collapses. An unexpected icing episode can slash a wind-plant's output by tens of megawatts within minutes, while excessive winds force cut-out shutdowns; both events create sharp ramps that the utility must offset with fast-responding thermal units, battery assets, or market purchases, often at premium prices. Day-ahead and intraday visibility into these weather-driven outages therefore underpins accurate load-generation balancing, congestion management, and reliable commitment of ancillary services, ultimately lowering imbalance penalties and safeguarding system stability.
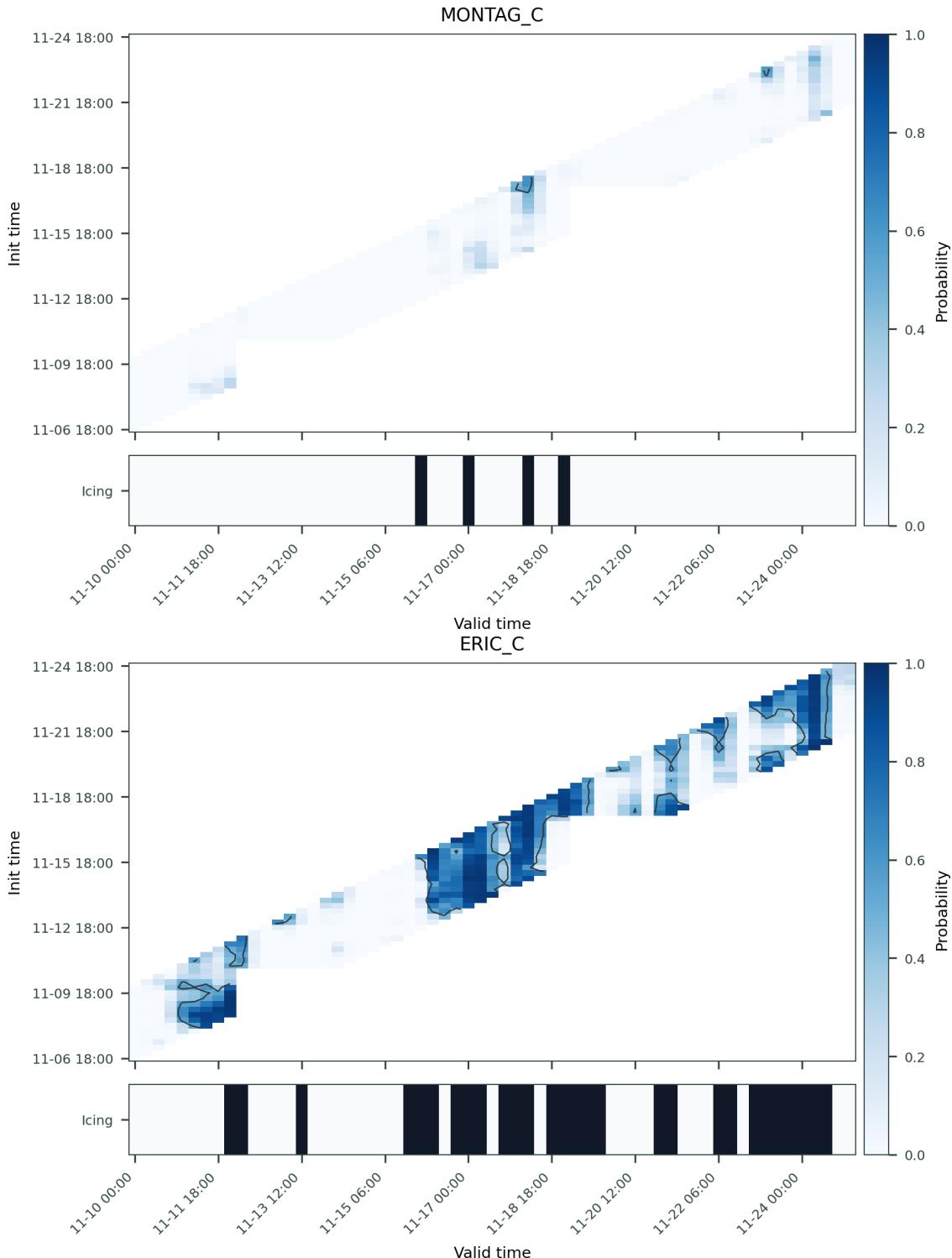
Figure 4: Romaine rime-ice event (Nov 2024) at two nearby Sygivre stations 11 km apart. Top: `MONTAG_C`; bottom: `ERIC_C`. Heatmap shows GFT-HQ 1-h rime-ice probabilities from successive 6-hourly initializations ($y$-axis) verifying at each valid time ($x$-axis; UTC) over 2024-11-06 to 2024-11-24. The black strip labelled "Icing" is the observed binary occurrence. At `ERIC_C`, a stable high-probability signal appears by 16 Nov and peaks on 18–19 Nov, aligning with the longest observed episode; several later bursts are also forecast. At `MONTAG_C`, the signal is weaker and more episodic, with only modest probabilities near the brief observed bursts, illustrating strong micro-scale differences despite close proximity.
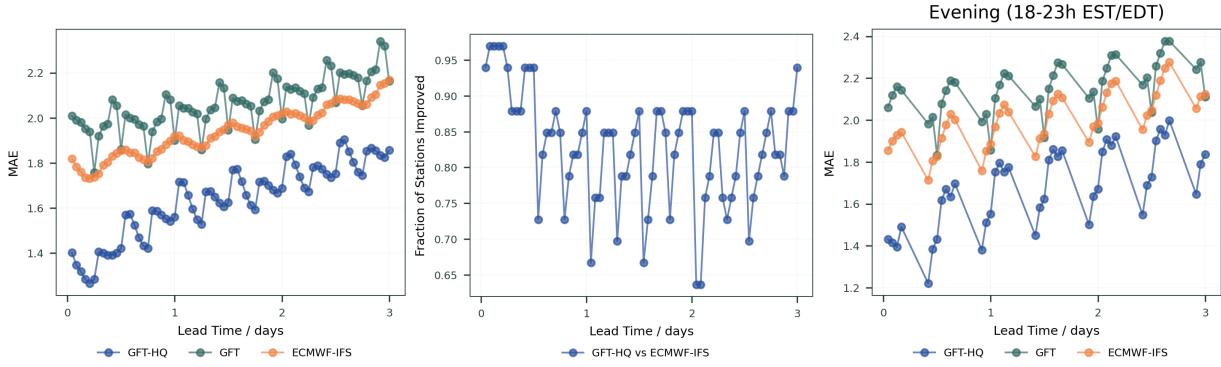
7

Figure 5: Wind-speed forecast improvements from finetuning. (a) Mean absolute errors of forecasts across all stations (b) Fraction of stations with lower forecasting errors than ECMWF-IFS (c) Mean absolute error across all the stations in the evening during peak load
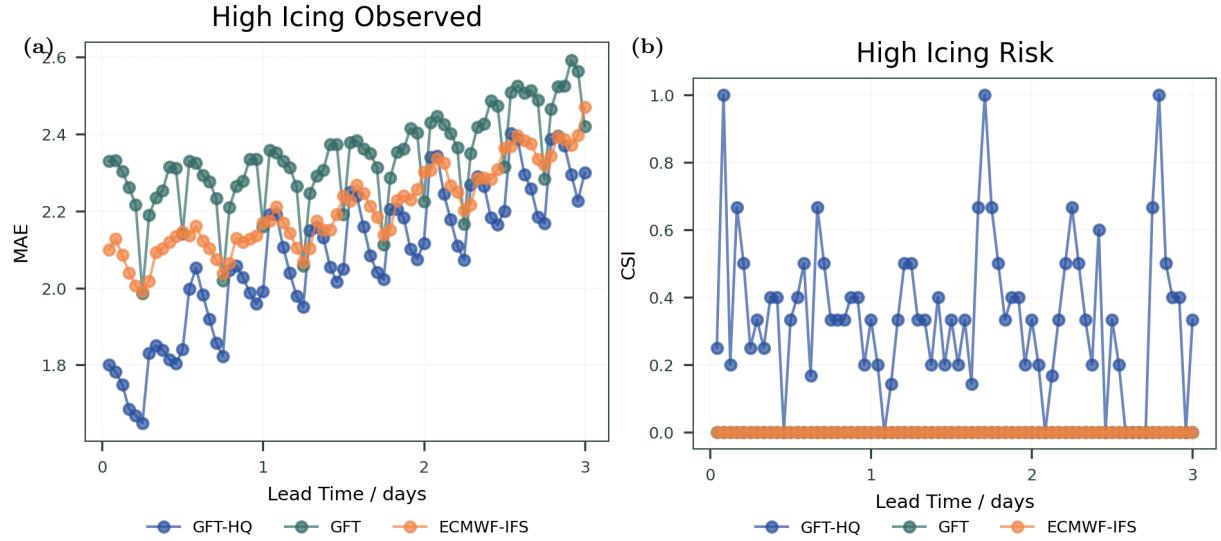


Figure 6: (a) Mean absolute error (m/s) across all stations during peak icing as measured by reduction in power production (b) Cut out windspeed ($> 25$m/s) detection during high icing

**Data** As the grid operator, Hydro-Québec does not have direct measurements of icing on wind-turbine blades. Instead, this label is inferred from substantial power-production declines observed at a wind farm given the measured wind speed and direction. The wind speed measurements come from met masts placed next to the wind farms such that potential wake effects are negligible. As before, we fine-tune the model on data from 2016–2023 and test on data from January 2024 to January 2025.

We use 65 meteorological masts near wind farms across Québec (10-minute data, multiple heights) and infer hourly wind-farm icing from production losses relative to potential output given the meteorological conditions. See Section A.2 for details.

**Results** Across the held-out period , the finetuned model (GFT-HQ) delivers strong and geographically consistent skill on both the *drivers* (hub-height wind) and the *impact* target (wind-farm icing).

**Icing risk** At a base rate of 13%, GFT-HQ attains AP = 0.73 ($\approx \times$lift over random), for detecting whether *any* icing will occur in the next 24 h, exceeding both Climatology (AP = 0.62) and an ERA5-derived Makkonen index (AP = 0.50). Evaluated strictly at 24 h lead time, GFT-HQ again leads (AP = 0.73) versus
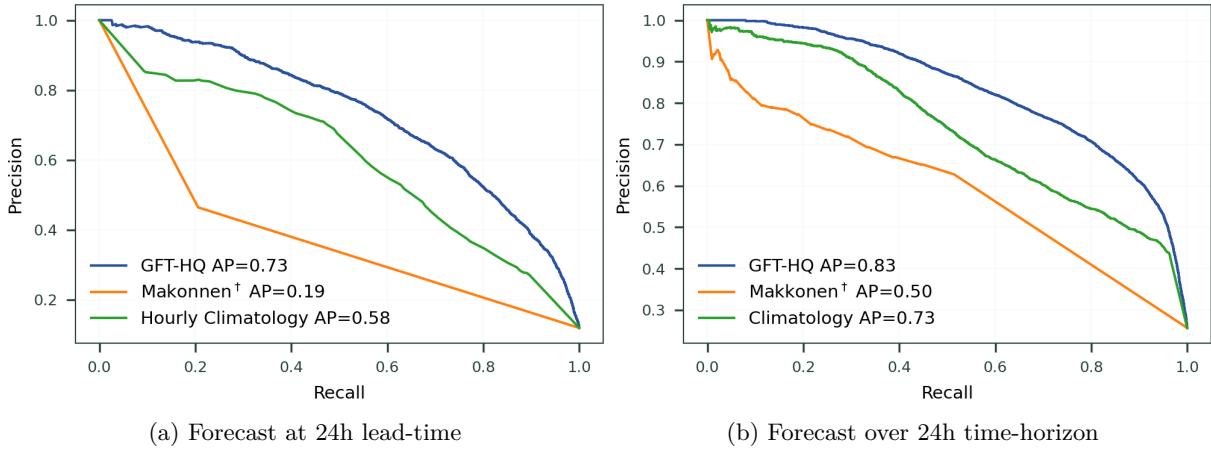
(a) Forecast at 24h lead-time      (b) Forecast over 24h time-horizon

Figure 7: Precision–recall curves for detecting wind-turbine icing across the wind-farm dataset. The fine-tuned model (GFT-HQ) is compared against climatology and an ERA5-derived Makkonen index. Panel (a) evaluates a fixed 24 h lead; panel (b) evaluates the windowed target of any icing within the next 24 h. Average precision (AP) summarizes performance; the evaluation base rate is 13%, so improvements correspond to substantial lift over random. Note that Makkonen[†] is included as a physics-based reference rather than an operational baseline.
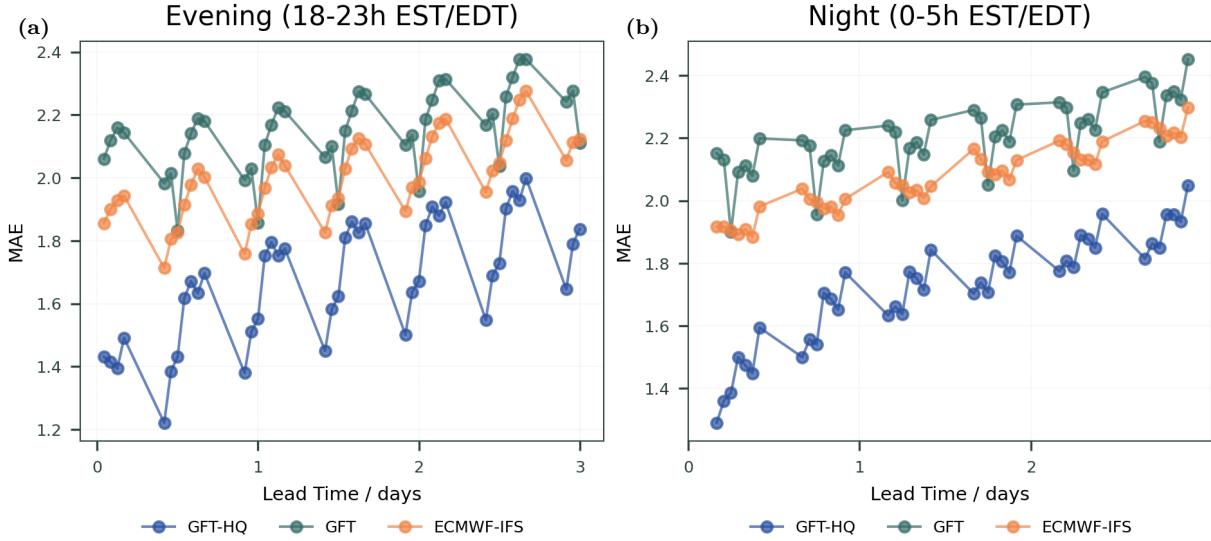


Figure 8: Diurnal hub-height wind MAE (m/s). GFT-HQ significantly narrows the mean absolute error during the evening/night periods when Hydro-Québec faces the highest load and icing exposure.

Climatology (AP=0.62), and Makkonen (AP = 0.19). Corresponding ROC performance is high as well (AUC = 0.93 vs 0.86–0.87 (Climatology) and 0.59–0.71 (Makkonen) in the two setups). *Note.* The Makkonen curve is computed from ERA5 reanalysis (see Section A.1.4) and is included as a strong physics-based *reference*, not an operational forecast baseline. Performance here is often worse than Climatology partially because icing is not directly measured at wind-farms, rather inferred from loss in power production and could include other factors as well.

**Wind-speed skill** Finetuning materially reduces hub-height wind MAE across 0–72 h relative to untuned GFT and ECMWF-IFS, with the majority of sites showing lower error than IFS at each lead time. Gains persist and are most pronounced during evening hours (18–23 EST/EDT), which are operationally critical for peak load.
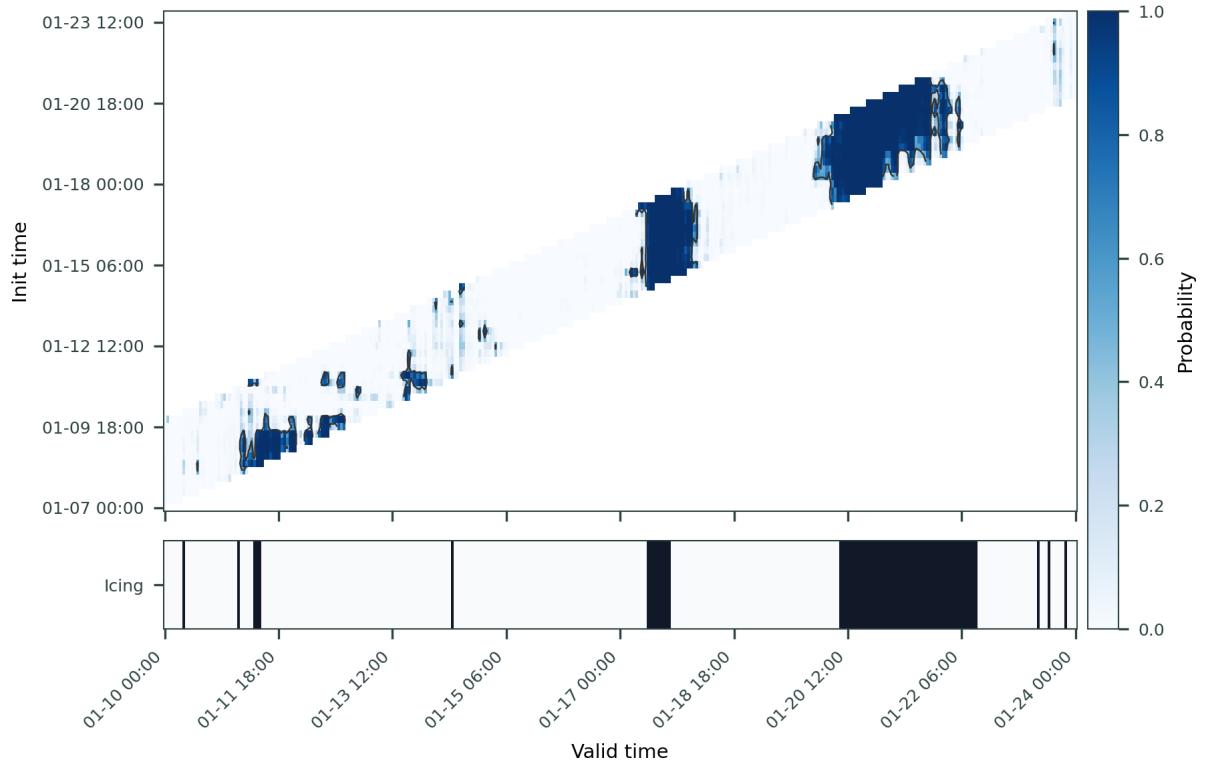
9

Figure 9: Wind-farm icing case study at an anonymized site ("AVA"). Rows are 6-hourly forecast initializations ($y$-axis); columns are valid hours (UTC, $x$-axis). The heatmap shows GFT-HQ probabilities of 1-h turbine icing; the black bar labelled "Icing" indicates observed occurrence inferred from production. A stable high-risk corridor appears several days in advance and peaks during the multi-day episode in late January, aligning with observations while also flagging shorter earlier bursts.

Lead-time averages remain high. Icing F1 stays above 0.6 through 60 h and ROC-AUC remains above 0.9 (plots in Section A.5.2)—while the diurnal slices in Figure 8 illustrate how GFT-HQ damps the errors compared to both baselines. Seasonal panels are provided in Section A.5.2; the evening/night margins are the largest, aligning with Hydro-Québec's highest-risk periods.

**Performance under high-icing conditions.** As can be seen in Figure 6, conditioned on hours flagged as high icing risk, GFT-HQ keeps the wind-speed accuracy edge and is the only model with consistently non-zero skill for turbine cut-out events ($> 25$m/s). Its CSI shows frequent high-skill episodes, whereas both GFT and ECMWF-IFS are near zero.

**Case study** To illustrate how the system presents actionable guidance to operators, Figure 9 shows a mid-winter icing episode at a representative Hydro-Québec wind farm (site code "AVA"). The Hovmöller-style plot stacks successive 6-hourly forecasts by initialization time ($y$-axis) against their valid hours ($x$-axis); colours denote the GFT-HQ probability of 1-h turbine icing, while the black strip beneath marks the observed binary icing derived from production losses. A coherent high-probability swath emerges and stabilizes 2–3 days before the prolonged mid-to-late January event, and shorter early-month bursts are also indicated. The cross-cycle persistence of this signal provides dependable early warning for curtailment planning and reserve scheduling.

## 3.3 Temperature and precipitation forecasting

Accurate temperature and precipitation forecasts sit at the heart of utility planning because they govern every major balance-sheet variable: demand, supply, and asset risk. Temperature is the dominant driver

(a) GFT-HQ Temperature skill vs. ECMWF-IFS



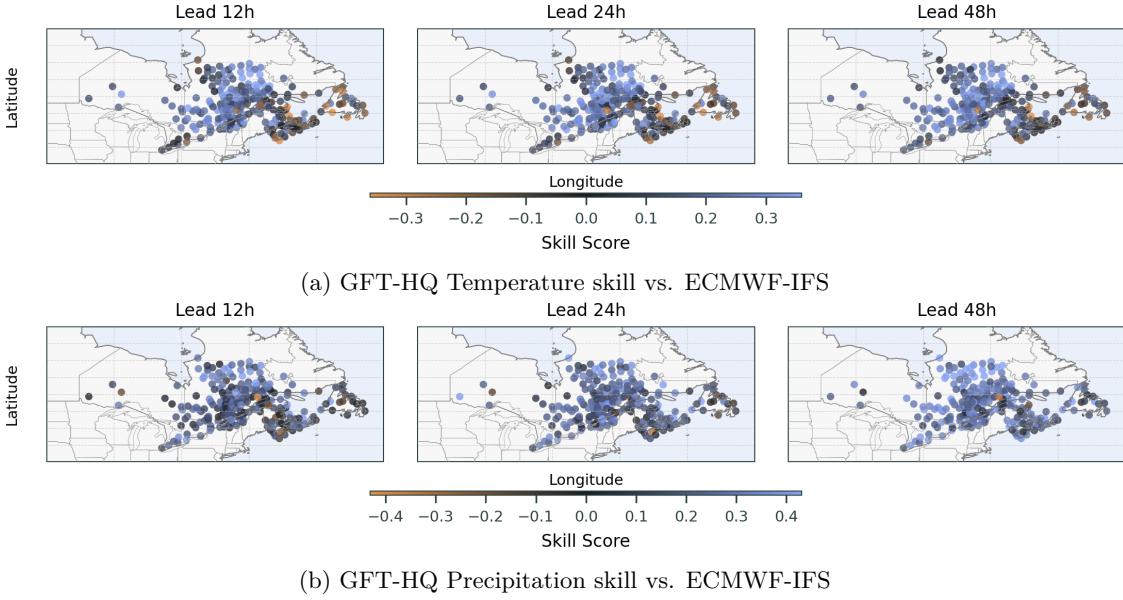(b) GFT-HQ Precipitation skill vs. ECMWF-IFS

Figure 10: Spatial coverage of positive fractional skill for GFT-HQ relative to ECMWF-IFS at 12 h, 24 h, and 48 h. Finetuning expands the footprint of improved stations for both temperature and hourly precipitation, supporting system-wide load and hydro inflow decisions.

of electric load—hot spells boost air-conditioning demand, cold snaps spike heating and resistive losses—so even a 1 °C error can translate into hundreds of megawatts of forecasting miss and costly imbalance charges. Precipitation matters just as much: rain, snow, or ice determine hydro-reservoir inflows, solar panel soiling, wind-turbine icing shutdowns, and outage probabilities from lightning or tree-fall. By anticipating these weather-dependent swings, a utility can commit generation and reserves economically, schedule maintenance when both load and storm risk are low, pre-stage crews and mobile transformers ahead of severe weather, and optimise fuel and emissions strategies. In short, precise temperature and precipitation outlooks convert meteorological uncertainty into actionable lead time, protecting reliability while trimming operating costs.

**Spatial footprint of gains.** Station-wise skill patterns (Section A.5.3) show that post-training amplifies rather than redistributes the model's strengths: GFT-HQ inherits the spatial structure of the untuned GFT but pushes most stations into the positive-skill regime by 24–48 h. Figure 10 highlights the expanded positive-skill footprint for both temperature and precipitation. The broad coverage means load-forecast errors decline system-wide instead of concentrating at a few sites, and precipitation improvements propagate across the river basins that drive Hydro-Québec's reservoir operations. Diurnal and seasonal breakouts in Figure 11 (additional panels in Section A.5.3) confirm that the gains hold during evening peaks and the extreme winter months that dominate risk planning.

**Data** We use hourly observations from federal/provincial networks and RMCQ (534 temperature stations and 241 automatically validated precipitation stations, each with $\geq 5$ years of data). We fine-tune the model on data from 2016-2023 and test on data from January 2024 to March 2025. See Section A.2 for details.

**Results** Figure 11a shows that GFT-HQ is consistently the lowest-error model for temperature forecasts across the entire 0–5-day horizon. Errors grow with lead time for all models, but the finetuned GFT-HQ model maintains a clear MAE gap to both GFT and ECMWF-IFS at every lead. The improvement is broadly distributed across sites: the fraction of stations with lower error than IFS rises from roughly a majority at short leads to $\approx 90-95\%$ by days 4-5. Gains are operationally aligned with HQ priorities: they are largest during evening peak hours (18-23 EST/EDT), and they persist through the anomalously hot 2024 summer, indicating robustness under regime shifts. In practical terms, trimming even a few-tenths of a degree in MAE
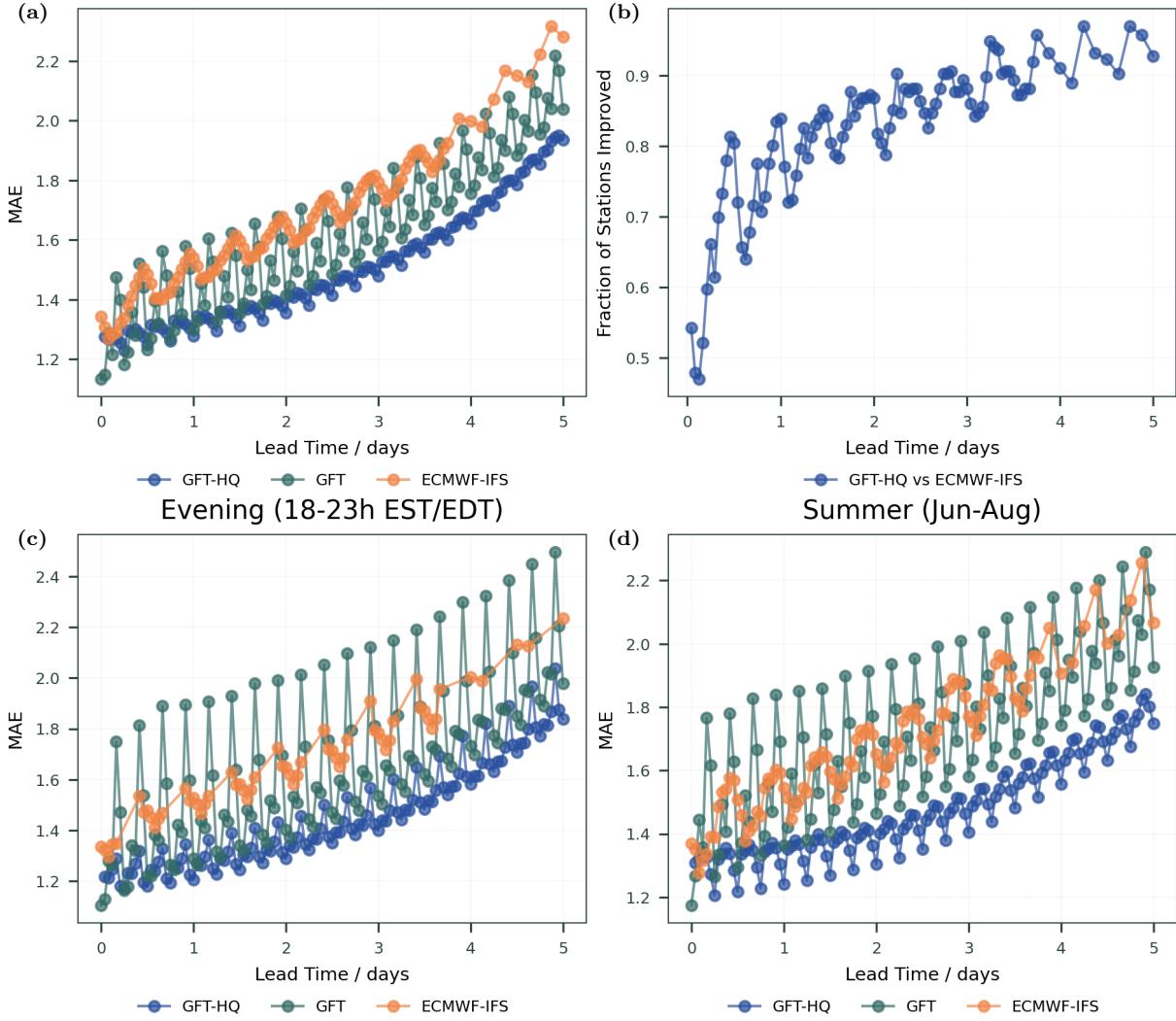
Figure 11: Temperature forecast improvements from fine-tuning. (a) Forecast errors across all stations (b) Fraction of stations with lower forecasting errors than ECMWF-IFS (c) Mean absolute error across all stations in the evening during peak load (d) Mean absolute error across all stations in the unusually hot 2024 summer

at these horizons materially reduces load-forecast error and the reserve over-commit or imbalance penalties that follow.

Figure 12a shows that GFT-HQ has the lowest MAE for precipitation forecasts at every lead from 0–5 days ($\approx 10\%$ better than GFT, $\approx 35\%$ better than ECMWF-IFS). The improvement is broadly distributed across sites: the fraction of stations beating IFS climbs to $\approx 90\%$ for day-ahead forecasts. Operationally, even a few-tenths of a tenth of a millimeter per hour compounds into meaningful 24-h accumulation error reductions, tightening day-ahead hydro-inflow and storm staffing plans.

## 4  Conclusion

Weather foundation models can materially improve weather intelligence for grid operations when post-trained on high-fidelity utility data. By adapting a pre-trained transformer based weather foundation model to Hydro-Québec assets, we produced hyper-local forecasts that outperformed strong NWP baselines on standard weather variables (temperature, precipitation, and wind). For rare but operationally critical hazards
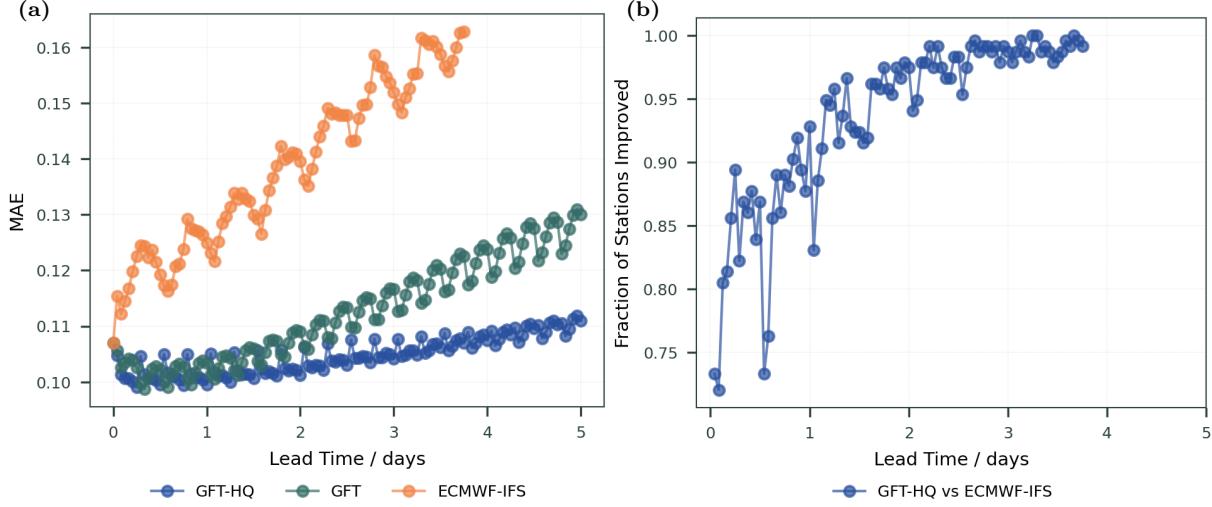
Figure 12: Hourly total precipitation (mm/h) forecast improvements from fine-tuning. (a) Forecast errors across all stations (b) Fraction of stations with lower forecasting errors than ECMWF-IFS

like rime ice, the model provides a new capability: reliable, hours-ahead alerts with usable precision-recall trade-offs. These gains translate into earlier and better-targeted interventions (e.g., de-icing; see Section A.3), more reliable renewable dispatch, and improved situational awareness at the asset level.

Looking ahead, we see three priorities for deployment: (i) decision-driven evaluation with utility-calibrated cost-loss parameters and crew constraints; (ii) tighter coupling with grid models (e.g., dynamic line ratings and outage risk) to convert forecast probabilities into asset risk; and (iii) continuous post-training and recalibration as new sensors come online. Beyond Québec, the same approach is directly transferable to other regions (Silurian AI, 2025) and infrastructure owners, requiring only modest amounts of local data to adapt the model.

## References

Zied Ben Bouallègue, Mariana C. A. Clare, Linus Magnusson, Estibaliz Gascón, Michael Maier-Gerber, Martin Janoušek, Mark Rodwell, Florian Pinault, Jesper S. Dramsch, Simon T. K. Lang, Baudouin Raoult, Florence Rabier, Matthieu Chevallier, Irina Sandu, Peter Dueben, Matthew Chantry, and Florian Pappenberger. The rise of data-driven weather forecasting: A first statistical assessment of machine learning–based weather forecasts in an operational-like context. *Bulletin of the American Meteorological Society*, 105(6):E864–E883, 2024. doi: 10.1175/BAMS-D-23-0162.1.

Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, 619(7970):533–538, 2023. doi: 10.1038/s41586-023-06185-3.

Cristian Bodnar, Wessel P Bruinsma, Ana Lucic, Megan Stanley, Anna Allen, Johannes Brandstetter, Patrick Garvan, Maik Riechert, Jonathan A Weyn, Haiyu Dong, et al. A foundation model for the Earth system. *Nature*, pp. 1–8, 2025.

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

Tom B. Brown, Benjamin Mann, Nick Ryder, et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, 2020.

Matthew Chantry, Simon Lang, Mihai Alexe, Jesper Dramsch, Baudouin Raoult, Mariana Clare, Mario Santa Cruz, Sara Hahner, Rilwan Adewoyin, Florian Pinault, et al. Aifs-ecmwf's data-driven forecasting system. In *105th Annual AMS Meeting 2025*, volume 105, pp. 449087, 2025.

Lei Chen, Xiaohui Zhong, Feng Zhang, Yuan Cheng, Yinghui Xu, Yuan Qi, and Hao Li. Fuxi: A cascade machine learning forecasting system for 15-day global weather forecast. *arXiv preprint arXiv:2306.12873*, 2023. doi: 10.48550/arXiv.2306.12873.

Martyn Clark, Subhrendu Gangopadhyay, Lauren Hay, Balaji Rajagopalan, and Robert Wilby. The Schaake shuffle: A method for reconstructing space–time variability in forecast fields. *Journal of Hydrometeorology*, 5(1):243–262, 2004.

Luca Delle Monache, Timothy D. Eckel, Daniel L. Rife, Bipin Nagarajan, and Kim Searight. Probabilistic weather prediction with an analog ensemble. *Monthly Weather Review*, 141(10):3498–3516, 2013.

ECMWF. *IFS Documentation CY47R3*, 2021. URL https://www.ecmwf.int/en/publications/ifs-documentation. Parts I–VII; see e.g. Part III: Dynamics and Numerical Procedures, doi:10.21957/b18qsx663.

Linyue Gao, Tao Tao, Yongqian Liu, and Hui Hu. A field study of ice accretion and its effects on the power production of utility-scale wind turbines. *Renewable energy*, 167:917–928, 2021.

Thomas Gerbet. Une ligne d'hydro-québec tombée, les centrales romaine-3 et 4 à l'arrêt. *Radio-Canada*, 2024. URL https://ici.radio-canada.ca/nouvelle/2125770/centrales-romaine-verglas-ligne-cote-nord.

Harry R. Glahn and D. A. Lowry. The use of model output statistics (mos) in objective weather forecasting. *Journal of Applied Meteorology*, 11(8):1203–1211, 1972.

Tilmann Gneiting, Adrian E. Raftery, Anton H. Westveld III, and Tom Goldman. Calibrated probabilistic forecasting using ensemble model output statistics. *Monthly Weather Review*, 133(5):1098–1118, 2005.

H. Hersbach, B. Bell, P. Berrisford, et al. The era5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020.

Hydro-Québec. Plan d'action 2035: Un québec vert et prospère. Technical report, Hydro-Québec, 2023a. Public action plan outlining generation and transmission expansion priorities.

Hydro-Québec. Plan d'approvisionnement 2023–2032. Technical report, Hydro-Québec, 2023b. Electricity supply plan submitted to the regulator.

Hydro-Québec. Plan stratégique 2024–2028. Technical report, Hydro-Québec, 2024. Corporate strategic plan with investment outlook.

Hydro-Québec. Évolution du réseau de transport principal: Axes à l'étude. Technical report, Hydro-Québec, May 2025. Map and study corridors for transmission expansion.

IEEE Power and Energy Society. IEEE standard for calculating the current–temperature relationship of bare overhead conductors, 2013.

International Organization for Standardization. Air quality meteorology siting classifications for surface observing stations on land, 2015. URL https://www.iso.org/standard/64287.html. Standard confirmed in 2020.

John Jumper, Richard Evans, Alexander Pritzel, et al. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.

Roger Koenker and Gilbert Bassett. Regression quantiles. *Econometrica*, 46(1):33–50, 1978.

J Kowles. Discussion of icing rate measurement and the rosemount icing rate system. *Rosemount Rep. A*, 67312:18, 1973.

JL Laforte, MA Allaire, and J Laflamme. Wind tunnel evaluation of a rime metering device using a magnetostrictive sensor. *Atmospheric research*, 36(3-4):287–301, 1995.

Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirnsberger, Meire Fortunato, Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, et al. Learning skillful medium-range global weather forecasting. *Science*, 382(6677):1416–1421, 2023.

Peng Li, Na Zhao, Donghua Zhou, Min Cao, Jingjie Li, and Xinling Shi. Multivariable time series prediction for the icing process on overhead power transmission line. *The Scientific World Journal*, 2014(1):256815, 2014.

Lasse Makkonen. Models for the growth of rime, glaze, icicles and wet snow on structures. *Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, 358(1776):2913–2939, 2000.

Petr Musilek, Dan Arnold, and Edward P Lozowski. An ice accretion forecasting system (iafs) for power transmission lines using numerical weather prediction. *Sola*, 5:25–28, 2009.

Tung Nguyen, Johannes Brandstetter, Ashish Kapoor, Jayesh K Gupta, and Aditya Grover. ClimaX: A foundation model for weather and climate. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 25904–25938. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/nguyen23a.html.

TR Oke. Guide to instruments and methods of observation. *WMO: Geneva, Switzerland*, 3:426, 2018.

OpenAI. GPT-4 technical report, 2023.

Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh Chattopadhyay, Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar Azizzadenesheli, Pedram Hassanzadeh, Karthik Kashinath, and Animashree Anandkumar. Fourcastnet: A global data-driven high-resolution weather model using adaptive fourier neural operators. *arXiv preprint arXiv:2202.11214*, 2022.

Alec Radford, Jong Wook Kim, Chris Hallacy, et al. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, 2021.

Adrian E. Raftery, Tilmann Gneiting, Fadoua Balabdaoui, and Michael Polakowski. Using Bayesian model averaging to calibrate forecast ensembles. *Monthly Weather Review*, 133(5):1155–1174, 2005.

Lennart Rasp and Sebastian Lerch. Neural networks for postprocessing ensemble weather forecasts. *Monthly Weather Review*, 146(11):3885–3900, 2018.

Suman Ravuri, Karel Lenc, Matthew Willson, Dmitry Kangin, Remi Lam, Piotr Mirowski, Megan Fitzsimons, Maria Athanassiadou, Sheleem Kashem, Sam Madge, et al. Skilful precipitation nowcasting using deep generative models of radar. *Nature*, 597(7878):672–677, 2021. doi: 10.1038/s41586-021-03854-z.

Alexander Rives, Joshua Meier, Tom Sercu, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021.

Roman Schefzik, Thorsteinn L. Thorarinsdottir, and Tilmann Gneiting. Uncertainty quantification in complex simulation models using ensemble copula coupling. *Statistical Science*, 28(4):616–640, 2013.

Silurian AI. Gft-us announcement. Silurian AI Blog, 2025. URL https://silurian.ai/blog/gft-us-announce.

W. C. Skamarock, J. B. Klemp, J. Dudhia, D. O. Gill, D. M. Barker, M. G. Duda, X.-Y. Huang, W. Wang, and J. G. Powers. A description of the advanced research wrf version 4. Technical Report NCAR/TN-556+STR, NCAR, 2019.

Casper Kaae Sønderby, Lasse Espeholt, Jonathan Heek, Mostafa Dehghani, Avital Oliver, Tim Salimans, Shreya Agrawal, Jason Hickey, and Nal Kalchbrenner. Metnet: A neural weather model for precipitation forecasting. *arXiv preprint arXiv:2003.12140*, 2020. doi: 10.48550/arXiv.2003.12140.

Daniela Szwarcman, Sujit Roy, Paolo Fraccaro, Þorsteinn Elí Gíslason, Benedikt Blumenstiel, Rinki Ghosal, Pedro Henrique de Oliveira, Joao Lucas de Sousa Almeida, Rocco Sedona, Yanghui Kang, et al. Prithvi-eo-2.0: A versatile multi-temporal foundation model for earth observation applications. *arXiv preprint arXiv:2412.02732*, 2024.

Maxime Taillardat, Olivier Mestre, Michaël Zamo, and Philippe Naveau. Calibrated ensemble forecasts using quantile regression forests and ensemble model output statistics. *Monthly Weather Review*, 144(6): 2375–2393, 2016.

Stéphane Vannitsem, Daniel S. Wilks, and Jakob W. Messner (eds.). *Statistical Postprocessing of Ensemble Forecasting*. Elsevier, 2018.

Guoyu Wang, Jie Shen, Minghong Jin, Shuai Huang, Zhong Li, and Xinchun Guo. Prediction model for transmission line icing based on data assimilation and model integration. *Frontiers in Environmental Science*, 12:1403426, 2024.

Tien Wang, Yogesh Balaji, Ting Chen, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023.

Daniel S. Wilks. *Statistical Methods in the Atmospheric Sciences*. Academic Press, 3 edition, 2011.

Yang Yang, Hao Lou, Jinran Wu, Shaotong Zhang, and Shangce Gao. A survey on wind power forecasting with machine learning approaches. *Neural Computing and Applications*, 36(21):12753–12773, 2024.

## A  Appendix

### A.1  Methods

#### A.1.1  WFM post-training vs. NWP post-processing

This subsection expands the discussion in Section 2 by detailing how we adapt a pretrained weather foundation model (WFM), namely GFT, to Hydro-Québec assets via *post-training* on utility observations. Conceptually, let $f_\theta$ denote the pretrained forecaster producing multi-variable fields and site forecasts at hourly lead times. Post-training updates $\theta$ using supervised losses defined on asset-level targets, so that the model's latent dynamics and decoders jointly produce hyper-local predictions. This approach fundamentally differs from traditional Numerical Weather Prediction (NWP) *post-processing*.

The primary distinction lies in where the learning occurs. Conventional post-processing methods—such as Model Output Statistics (MOS), Ensemble Model Output Statistics (EMOS), or various machine learning models (Glahn & Lowry, 1972; Gneiting et al., 2005; Raftery et al., 2005; Koenker & Bassett, 1978; Delle Monache et al., 2013)—learn a statistical mapping from a *fixed* upstream NWP model's output to local observations. The NWP model itself remains unchanged. In contrast, our post-training approach directly updates the core forecasting model, $f_\theta$, enabling its internal spatiotemporal representations to better capture local physical phenomena like terrain and human infrastructure induced effects like icing.

This distinction leads to several key advantages:

- **Multivariate Coherence:** Post-processing is typically performed on a per-variable, per-site basis, making it challenging to maintain physical and statistical consistency across different weather variables and lead times. While methods like the Schaake Shuffle (Clark et al., 2004) or Ensemble Copula Coupling (ECC) (Schefzik et al., 2013) can reconstruct some multivariate dependencies, our post-training framework is inherently a multi-task learning problem. By jointly predicting variables such as temperature, wind, precipitation, and icing, the model learns their interdependencies, ensuring greater physical coherence.

- **Generation of Novel Variables:** Post-processing is constrained by the output variables of the upstream NWP model. It cannot generate forecasts for phenomena not explicitly predicted by the NWP system. Our WFM, however, can be trained to predict new variables (here e.g., probability of rime ice on transmission lines) directly from local observations, leveraging the rich, shared latent state of the foundation model.

- **Scalability and Operational Value:** A traditional post-processing approach often requires developing and maintaining a large portfolio of individual models for each site and variable. The post-trained WFM provides a more scalable solution, with a single set of model weights serving hundreds of assets and multiple forecasting tasks. Because the fine-tuning adjusts the model's core dynamics rather than merely correcting surface-level biases, it yields more significant improvements in forecasting rare, high-impact events that are critical for operational decision-making.

### A.1.2   Experimental design

**Model Inputs and Targets**. At inference time, the encoder is initialized with dense ECMWF-IFS Analysis fields at the cycle issue time. The model is fine-tuned to decode to five targets used in Hydro-Québec operations: 2 m temperature, hourly precipitation, hub-height wind-speed, wind-farm icing indicator, and transmission-line rime-ice indicator. Unless explicitly stated, Hydro-Québec asset streams are used for supervision only and are not assimilated at runtime.

**Supervision and Data**. We employ standard regression losses for the continuous targets (temperature, wind, precipitation) and probabilistic classification losses for the binary icing indicators. The model is trained on observational data spanning from 2016 to 2023 and evaluated on a hold-out test set covering 2024 and 2025. Forecasts are evaluated at hourly cadence and also aggregated into decision windows (e.g., any icing in next 24 h) for operational metrics.

**Inference**. The post-trained model produces both dense gridded outputs and sparse hyper-local site forecasts in a single forward pass. This unified workflow replaces the complex, multi-stage process of running an NWP model followed by a large suite of separate post-processing models.

### A.1.3   Evaluation metrics and interpretation

We summarize the scores used in this work and their operational meaning for AI-based weather forecasting, building on the protocol introduced in the main text.

**Base rate (class imbalance).**   For a binary event indicator $y \in \{0, 1\}$, the *base rate* is $\pi = \frac{1}{N} \sum y$, i.e., the fraction of positive hours (or windows). Rime-ice at Sygivre sites is rare (example: $\pi \approx 3.68\%$); wind-farm icing windows can be more frequent (example: $\pi \approx 13\%$). Low $\pi$ makes precision–recall metrics more informative than accuracy.

**MAE (mean absolute error).**   For continuous variables (temperature, wind, precipitation), $MAE = \frac{1}{N} \sum |\hat{y} - y|$ measures average magnitude of errors in native units (K, m/s, mm/h). Lower is better; it is robust and directly maps to operational tolerances (e.g., thermal ratings (IEEE Power and Energy Society, 2013), curtailment thresholds).

17

**Precision, Recall, F1.** For thresholded probabilistic alerts $\hat{y}_\tau = \mathbb{K}[p \geq \tau]$:

$$\text{Precision} = \tfrac{\text{TP}}{\text{TP+FP}}, \qquad\qquad\qquad \text{(false-alarm burden)}$$
$$\text{Recall (POD/TPR)} = \tfrac{\text{TP}}{\text{TP+FN}}, \qquad\qquad\qquad \text{(missed-event rate)}$$
$$\text{F1} = \tfrac{2\,\text{Precision}\cdot\text{Recall}}{\text{Precision+Recall}}. \qquad\qquad \text{(balance under imbalance)}$$

In operations, precision captures costly false dispatches; recall captures avoided misses; F1 summarizes the trade-off for rare hazards like icing.

**PR curves, AP (PR-AUC), and lift.** Sweeping the alert threshold yields a precision–recall (PR) curve. *Average Precision* (AP) is the area under the PR curve (also referred to as PR-AUC). AP increases when high-probability hours correspond to observed events. *Lift* contextualizes AP under imbalance: lift $= \frac{\text{AP}}{\pi}$; a lift of 7× means a random hour drawn from the model's top-ranked alerts is seven times likelier to be an event than an arbitrary hour.

**ROC and AUC.** The ROC curve plots TPR vs FPR as $\tau$ varies; its area (AUC) lies in $[0, 1]$. ROC-AUC is threshold-free and widely used, but with very low base rates it can overstate usefulness; we therefore emphasize PR/AP for icing and report ROC-AUC for completeness (e.g., wind-farm icing).

**CSI (Critical Success Index) and IoU.** CSI (a.k.a. *threat score*) measures categorical event performance at a fixed threshold: CSI $= \frac{\text{TP}}{\text{TP+FP+FN}}$. In computer vision terms this is the *Intersection over Union* (IoU) between the predicted and observed event sets; the two are equivalent for binary events in time or space.

**Relative skill vs baseline.** When comparing to a baseline $b$ on an error metric $E$ (e.g., MAE) we report fractional skill

$$\text{Skill} = 1 - \frac{E_{\text{model}}}{E_b} \ \in (-\infty, 1],$$

so that positive values indicate improvement ("positive = lower error than IFS" in our maps).

**Windowed events.** For "any icing in the next $\Delta$ h" decisions we aggregate hourly probabilities $p_{t+h}$ into a window probability $q_t = 1 - \prod_h (1 - p_{t+h})$ and then evaluate the same metrics on the windowed binary events; see *Windowed event probability* in Section A.3.

### A.1.4 ERA5-derived Makkonen index

We construct a simple, physically motivated reference for rime-icing risk from ERA5 reanalysis fields (hourly, 0.25° grid; Hersbach et al., 2020). This "Makkonen" index (Makkonen, 2000) is used only as a retrospective, physics-based benchmark in our comparisons and not as an operational forecast.

**Sites and interpolation.** We take the locations of Hydro-Québec wind farms and Sygivre transmission-line sites and bilinearly interpolate ERA5 to each site at hourly cadence (Hersbach et al., 2020). Where site height above ground is unknown, we assign a fallback height of 80 m for wind farms and 50 m for transmission lines. To limit computation and reflect relevant icing layers, we use pressure levels from 800 to 1000 hPa and the time range 2024-01-01 to 2025-06-01.

**Variables.** The following fields are used: 2 m temperature, 2 m dewpoint, 10 m and 100 m wind components (to form wind speed), surface pressure, surface geopotential, pressure-level geopotential, and specific cloud liquid water content on pressure levels from ERA5 at a site (Hersbach et al., 2020). Given the hour, we derive:

- **Wind speed at site height** $v(z)$ via a power-law profile using 10 m and 100 m winds:

$$v(z) = v_{10} \left(\tfrac{z}{10}\right)^\alpha, \quad \alpha = \frac{\ln(v_{100}/v_{10})}{\ln(100/10)}. \tag{2}$$

- **Air temperature at site height** using a standard environmental lapse rate of 6.5 K km$^{-1}$ from 2 m: $T(z) = T_{2\,\mathrm{m}} - 6.5\,\mathrm{K\,km^{-1}}\,\frac{z}{1000}$.

- **Liquid water content (LWC) at site height** from the pressure-level specific cloud liquid water content, sampled at the site height using the geopotential field and converted to volumetric units (kg m$^{-3}$).

**Icing and proxy rate.** Following operational heuristics, rime icing is considered *feasible* when all of the following hold:

$$v(z) > 0 \text{ m s}^{-1}, \quad T(z) \in [260, 275] \text{ K}, \quad \text{LWC}(z) > 0.001 \text{ g m}^{-3}. \tag{3}$$

We then define a simple *rate-of-icing proxy*

$$r(t) = \begin{cases} v(z)\,\text{LWC}(z), & \text{if icing feasible}, \\ 0, & \text{otherwise}, \end{cases} \tag{4}$$

which intentionally omits geometry- and material-dependent collection efficiencies and thus serves only as a relative index. For 24 h events, we aggregate $r(t)$ over the decision window and threshold to produce a binary "icing in next 24 h" signal used in our comparisons.

## A.2 Datasets

### A.2.1 Sygivre icing sensors

The Sygivre network comprises 40 atmospheric icing measurement stations, designed to measure either freezing rain or rime-ice accumulation in the vicinity of electricity transmission infrastructure. The main instrument at the stations is an ice-accretion measurement system. The central component of the system is a Rosemount 871 magnetorestrictive oscillation probe excited to vibrate at a natural frequency of 40 kHz (Kowles, 1973). As ice accretes on the probe, the inertia increases, and the corresponding decrease in vibration frequency is measured. While the Rosemount 871 is used widely to assess the occurrence of icing conditions, it is not meant for quantifying large accumulations of ice, which can encase the sensor and isolate it from the environment. For that reason, the second component of the system is a cylindrical mount for the probe, developed at the Hydro-Québec Research Center (Laforte et al., 1995). The cylindrical mount contains a heater and an electromechanical system to shake the probe and remove ice completely, thus enabling continued measurements in large events. The final component is a controller that triggers the shaking and heating functions to rid the probe of ice at predetermined intervals. After measuring a decrease of 200 Hz, the probe heats up and shakes to remove the accumulated ice, and a "de-icing cycle" is logged. The power-line ice mass accretion corresponding to a 200 Hz decrease in probe frequency is calibrated separately in controlled experiments for dense freezing-rain ice and porous rime ice (Laforte et al., 1995). An icing event is characterized by the number of "de-icing cycles" the probe undergoes, which provides an estimate of the ice mass accumulation. Of the 40 stations of the Sygivre network, 26 are in valley or coastal areas that are more exposed to freezing-rain events. The remaining 14 stations, used in the present study, are located at altitude near steep slopes and mostly record rime-ice events. Sygivre stations also include temperature, humidity, wind speed, and wind-direction sensors.

### A.2.2 Wind-farm meteorological masts

The wind data used here come from 65 out of a network of 70 permanent meteorological masts located near wind farms in Québec. The wind masts have been operational since wind-farm commissioning and have collectively recorded 800 mast-years of 10-minute-resolution data since 2006. Each mast includes at least three measurement levels between 10 and 130 meters. The masts are equipped with both heated and unheated sensors from manufacturers such as NRG Systems, Thies Clima, R.M. Young, Risø, WindSonic (Gill Instruments), and Vaisala. These sensors measure wind speed and direction, temperature, humidity, and atmospheric pressure. Data undergo automated quality control to remove invalid or erroneous values, though some uncertainty remains due to the absence of manual inspection. When multiple sensors are present at the same height, their validated measurements are combined into a single representative value per level.

Ice-loss production for a given wind farm is not directly measured, but inferred from its potential production considering the wind speed and the temperature observed over all its turbines. The potential production is evaluated using the wind-farm power curve, which accounts for non-atmospheric losses (electrical losses, wake effects, etc.). The ice-loss production fraction is then taken as the ratio of the observed production of the wind farm to its potential production.

### A.2.3 Temperature and precipitation stations

The weather data used for post-training and evaluation of precipitation and temperature forecasts is gathered by the RMCQ (Quebec Collaborative Weather Network) as well as by Environment and Climate Change Canada and the Ministry of the Environment, the Fight Against Climate Change, Wildlife and Parks of Quebec. RMCQ centralizes weather station data shared by members, including utilities, mining companies and a wildfire prevention organization, that exploit weather stations spread over Quebec. The hourly temperature data used here are not validated and come from 534 stations, while the hourly precipitation data come from a subset of 241 stations and are automatically validated. Each station used in the study has at least the equivalent of 5 years of available data (43800 hours). Precipitation measuring stations are set up following WMO standards (International Organization for Standardization, 2015; Oke, 2018). Precipitation data is validated automatically by the nonprofit Solutions Mesonet, using Oklahoma Climatological Survey guidelines. Validation includes steps like comparing nearest neighbours, flagging time discontinuities, thresholding for values that are unphysical or out of climatological range, and comparing similar instruments at a given station (e.g., Weighing Precipitation Gauge vs Tipping Bucket Rain Gauge).

## A.3 De-icing decision making

Grid operators care less about generic skill and more about whether a forecast triggers the right crew action at the right time. We therefore evaluate the rime-ice head using a standard *cost–loss* framework, adapted to partial loss avoidance by pre-emptive de-icing.

**Setup.** For each asset (or line segment) $i$ and time $t$, let $y_{t,i} \in \{0,1\}$ indicate whether at least one rime-ice accretion cycle occurs in the window $[t, t + \Delta]$ (e.g., $\Delta = 24$ h). The model outputs a probability $p_{t,i} = \Pr(y_{t,i} = 1 \mid \mathcal{F}_t)$. The operator chooses $a_{t,i} \in \{\text{DISPATCH}, \text{HOLD}\}$. Dispatching incurs a cost $C_{d,i}(t)$ (crew hours, travel, overtime), independent of $y_{t,i}$. If icing occurs and no dispatch is made, the operator faces a loss $L_i(t)$ (outage risk, repair, penalties). If icing occurs and dispatch happens, only a fraction $(1 - \alpha_i(t))$ of this loss remains; $\alpha_i(t) \in [0,1]$ is the *mitigation effectiveness* (fraction of loss avoided by pre-emptive action). Travel-time constraints $\tau_i(t)$ impose a lead-time requirement: dispatch is only effective if initiated before $t + \Delta - \tau_i(t)$.

**Single-asset optimal threshold.** Under these assumptions, the expected cost of DISPATCH is $\text{EC}_{\text{DISPATCH}} = C_{d,i}(t) + p_{t,i}(1 - \alpha_i(t)) L_i(t)$ and the expected cost of HOLD is $\text{EC}_{\text{HOLD}} = p_{t,i} L_i(t)$. DISPATCH is optimal when $\text{EC}_{\text{DISPATCH}} \leq \text{EC}_{\text{HOLD}}$, yielding the probability threshold $p_i^\star(t)$ reported in Equation (1). Intuitively, crews mobilize when the icing probability exceeds the ratio "cost of acting" over "avoidable loss."

**Windowed event probability.** Decisions target "any icing in the next $\Delta$ hours." If $p_{t+h,i}$ is the per-hour probability for $h = 0, \ldots, H - 1$ with $H = \Delta/1\,\text{h}$, the event probability over the window is the union

$$q_{t,i} \;=\; 1 - \prod_{h=0}^{H-1} \left(1 - p_{t+h,i}\right).$$

When using the window mean $\bar{p}_{t,i} = H^{-1} \sum_h p_{t+h,i}$ as a score, a Poisson/independence approximation yields $q_{t,i} \approx 1 - \exp(-H\,\bar{p}_{t,i})$. The optimal threshold on $q$ maps to a threshold on $\bar{p}$ via

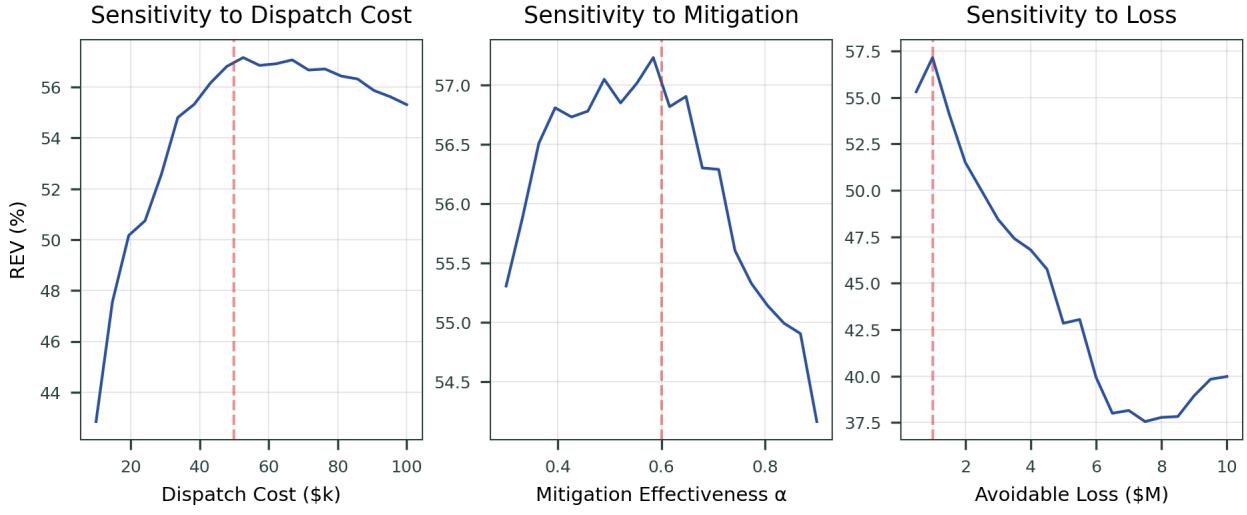$$\bar{p}_i^\star(t) \;=\; -\frac{\ln\left(1 - p_i^\star(t)\right)}{H}.$$

Figure 13: REV sensitivity: one-at-a-time variation of dispatch cost $C_d$, mitigation effectiveness $\alpha$, and avoidable loss $L$; vertical dashed lines mark the base setting.

Table 1: Default de-icing decision parameters and example thresholds.

| Parameter | Default |
|---|---|
| Decision window $\Delta$ | 24 h |
| Issuance cadence | 6 h |
| Lead-time (mobilization) $\tau$ | 0 h |
| Mitigation effectiveness $\alpha$ | 0.6 |
| Crew capacity $H(t)$ | 3 crews |
| Per-asset crew-hours $h_i(t)$ | 12 h |
| Dispatch cost $C_d$ (ground) | $50k |
| Helicopter adder | $75k |
| Avoidable loss $L$ | $400k |
| Hysteresis/persistence | $p_{\text{on}} = p^\star,\ p_{\text{off}} = 0.6\,p^\star,\ k = 2$ |

**Lead-time aware policy.** Let $\tilde{p}_{t,i}$ denote the forecast issued at time $t$ for the next $\Delta$ hours. Define the effective window as those issue times that satisfy the travel constraint $t \leq t_{\text{event}} - \tau_i(t)$ for any $t_{\text{event}} \in [t, t+\Delta]$. We evaluate at a 6-hour issuance cadence (00, 06, 12, 18), trigger a WATCH when $\tilde{p}_{t,i} \geq \eta\,p_i^\star(t)$ and a DISPATCH when $\tilde{p}_{t,i} \geq p_i^\star(t)$, with $0 < \eta < 1$ (e.g., $\eta = 0.7$) providing earlier mobilization. To avoid "thrashing," we apply hysteresis with two thresholds $p_{\text{on}}^\star > p_{\text{off}}^\star$ and a persistence requirement across issuance steps (e.g., exceedance for $k = 2$ consecutive issuances).

**Scope and limitations.** The dollar values and parameter ranges used here are *illustrative* and chosen to make the framework concrete rather than to prescribe operational thresholds. Real systems exhibit wide variability by terrain and access (helicopter vs. ground), voltage class, outage externalities, work rules, and weather. We also abstract from several complexities: (i) heavy–tailed and network–coupled losses (e.g., cascading redispatch and penalties), (ii) routing, staging, backlog and queueing across multi–day events, (iii) shared crews and switching clearances.

**Measuring decision value.** We report the *relative economic value* (REV) of the forecast policy against both a climatology policy and a perfect-forecast oracle:

$$\text{REV} = \frac{\mathbb{E}[C_{\text{clim}}] - \mathbb{E}[C_{\text{fcst}}]}{\mathbb{E}[C_{\text{clim}}] - \mathbb{E}[C_{\text{perf}}]} \in [0, 1], \tag{5}$$

21

where expected costs are averaged over the held-out period. Here $C_{\text{fcst}}$ is the realized cost under the threshold/hysteresis policy above, $C_{\text{clim}}$ is an instance-wise climatology baseline obtained by applying the same decision rule to climatology probabilities $C_{\text{perf}}$ assumes perfect knowledge of $y_{t,i}$.

We include a one-way sensitivity analysis to $(C_d, \alpha, L)$ to show robustness of REV and the optimal threshold.

**Base scenario results.** With $C_d = \$50,000$, $\alpha = 0.6$, and $L = \$400,000$ (yielding $p^\star \approx 20.8\%$ on the 24 h event probability), we obtain:

- REV = 57.2%;

- $C_{\text{fcst}} = \$51,913$; $C_{\text{clim}} = \$69,064$; $C_{\text{perf}} = \$39,055$;

- base rate $\mathbb{E}[y] = 8.7\%$.

Thus the forecast reduces expected cost by \$17,151 relative to the climatology baseline and closes about 57% of the gap to perfect information: $(69{,}064 - 51{,}913)/(69{,}064 - 39{,}055) \approx 0.572$.

## A.4 Related Work

**Traditional NWP vs. AI-based Weather Forecasting**: Weather forecasting has long been dominated by physics-based numerical weather prediction (NWP) models, which solve discretized fluid dynamics equations but at enormous computational cost. Operational global models like ECMWF's IFS require supercomputers and still make simplifying assumptions (e.g., subgrid parametrizations) that limit accuracy. In recent years, data-driven systems have advanced rapidly, offering orders-of-magnitude faster inference with increasingly competitive skill. More recently, systems including DeepMind's GraphCast Lam et al. (2023) and ECMWF's AIFS Chantry et al. (2025) demonstrated that learned surrogates can match or surpass operational guidance on many metrics at a fraction of the compute, catalyzing a shift toward hybrid and ML-first forecasting; see also (Ben Bouallègue et al., 2024) for an operational-like statistical assessment. Short-range precipitation nowcasting has also seen strong ML gains (Ravuri et al., 2021; Sønderby et al., 2020).

**Foundation Models for Weather and Climate**: Inspired by the paradigm of large-scale pretraining in NLP and vision, researchers began developing foundation models for Earth science that can be adapted to many downstream tasks Bodnar et al. (2025); Nguyen et al. (2023); Bi et al. (2023); Pathak et al. (2022); Chen et al. (2023); Szwarcman et al. (2024). Earliest example was ClimaX, a deep learning model trained on heterogeneous weather/climate datasets in 2023. ClimaX demonstrated that a single pre-trained model could be fine-tuned to diverse tasks, from short-range forecasts of standard variables to climate projections and downscaling. More importantly, it often matched conventional specialized models on benchmark tests. More recently, Aurora pushed this concept further: a 1.3-billion-parameter transformer that unified multiple Earth-system domains within one model. After pretraining on over a million hours of geophysical data, Aurora was post-trained (fine-tuned) for specific high-value applications, achieving state-of-the-art accuracy in each. Notably, Aurora outperformed operational NWP systems in 10-day global weather forecasting (improving error on 92% of targets vs. ECMWF's IFS), 5-day tropical cyclone tracking (20% lower track error than the multi-agency consensus), 5-day air quality (beat ECMWF's chemistry model on 74% of metrics), and 10-day ocean wave prediction (beat ECMWF's wave model on 86% of metrics), all at a fraction of the computing cost. Crucially, these feats were achieved by fine-tuning the same pretrained model for only a few hours per task, without retraining the core architecture. This demonstrates the power of foundation models: once a model has "learned the grammar of the Earth system" in its latent representations, it can be rapidly adapted to new domains with limited data. The Generative Forecasting Transformer (GFT) family extends this idea even further. GFT surpasses Aurora's global weather skill on all metrics, add finer temporal resolution (hourly forecasts), and expand the set of predicted variables. A regional variant, GFT-US, now delivers kilometer-scale 0–30 h forecasts tailored to the United States. In summary, the field has rapidly evolved from single-purpose neural nets to general-purpose weather foundation models (WFMs) that can democratize high-quality forecasts by post-training on specific tasks at modest cost. Our work builds directly on these advances—particularly Aurora and its GFT successors—by exploring how a WFM can be leveraged in the power grid domain.

**NWP post-processing vs. WFM post-training**: Utilities have long relied on statistical post-processing to adapt NWP outputs to sites and assets. Common approaches include bias correction and calibration such as Model Output Statistics (MOS), ensemble MOS (EMOS), Bayesian model averaging (BMA), quantile-based methods, analog ensembles, and neural/ML regressors (Glahn & Lowry, 1972; Gneiting et al., 2005; Raftery et al., 2005; Koenker & Bassett, 1978; Delle Monache et al., 2013; Rasp & Lerch, 2018; Taillardat et al., 2016; Wilks, 2011; Vannitsem et al., 2018). These methods learn a mapping from fixed NWP fields to local targets, typically one variable at a time and one site at a time, and they do not alter the upstream atmospheric dynamics. By construction, they cannot introduce truly new target variables that are absent from the NWP (e.g., rime-ice risk) and often struggle to maintain cross-variable, cross-lead physical coherence; dependency-preserving techniques like the Schaake Shuffle and Ensemble Copula Coupling (ECC) address some multivariate aspects but do not modify the underlying flow predictions (Clark et al., 2004; Schefzik et al., 2013). In contrast, our approach adapts the forecasting *model itself*. We fine-tune a pretrained WFM end-to-end on utility observations so that its latent dynamics and decoders jointly produce hyper-local, multi-variable predictions. This confers several advantages for operations: (i) multi-task training enforces consistency across variables (e.g., wind, temperature, precipitation, icing) and lead times; (ii) the model's sparse hyper-local decoder targets specific assets without training a separate model per site; (iii) new targets like rime-ice probability can be learned directly from observations; and (iv) updates propagate into the model's spatiotemporal representation, improving fine-scale phenomena rather than merely correcting biases at the output. Empirically, we observe that a single post-trained WFM can replace a portfolio of site-specific post-processing models while delivering higher skill, especially for rare, high-impact events.

**Weather Forecasting for Power Systems**: The electric power grid is highly weather-sensitive, with assets and operations affected by temperature, wind, precipitation, icing, etc. Traditionally, utilities have relied on NWP model outputs combined with statistical corrections for site-specific forecasting. For instance, wind power forecasting methods have long used physical NWP feeds to predict turbine wind speeds, often augmented by local regression models Yang et al. (2024). Beyond generation forecasts, weather hazards pose a major threat to grid infrastructure. Icing of transmission lines is a prime example: accreted ice from freezing rain or rime can add massive weight and cause line sag or breakage, leading to catastrophic outages (Makkonen, 2000). Icing also degrades wind-turbine production and availability in the field (Gao et al., 2021). This has made accurate icing forecasts increasingly crucial for grid reliability and disaster preparedness.

Forecasting line icing is notoriously challenging. Many icing forecast methods build on numerical weather prediction (NWP) models to provide the meteorological inputs for ice accretion models. High-resolution regional models (e.g., WRF (Skamarock et al., 2019)) can simulate temperature, wind, humidity, and precipitation fields that drive icing formation. Physical ice accretion models use these weather parameters to estimate ice growth on structures based on thermodynamic and empirical relationships (Makkonen, 2000). Musilek et al. (2009) developed an Ice Accretion Forecasting System (IFAS) that ingests NWP model data (including a precipitation-type algorithm by Ramer) to predict freezing rain occurrence and ice loads on power lines. Data-driven approaches have also been explored for transmission-line icing prediction (Li et al., 2014).

However, purely physical models face limitations. NWP models have constrained resolution and tend to smooth out terrain effects, so they may miss localized microclimate conditions (e.g., rime icing on a ridge or heavier glaze in a valley). This mismatch means a weather model might predict the general area of an icing event but still underestimate or misplace the peak ice loads. Wang et al. (2024) used a WRF model with 3D variational data assimilation to simulate a power line icing event. Local weather observations improved the accuracy of temperature and humidity fields (especially in the first 24 hours), which in turn made the icing forecasts more realistic near instrumented sites. However, WRF models are significantly more expensive to tune and operate compared to WFMs.
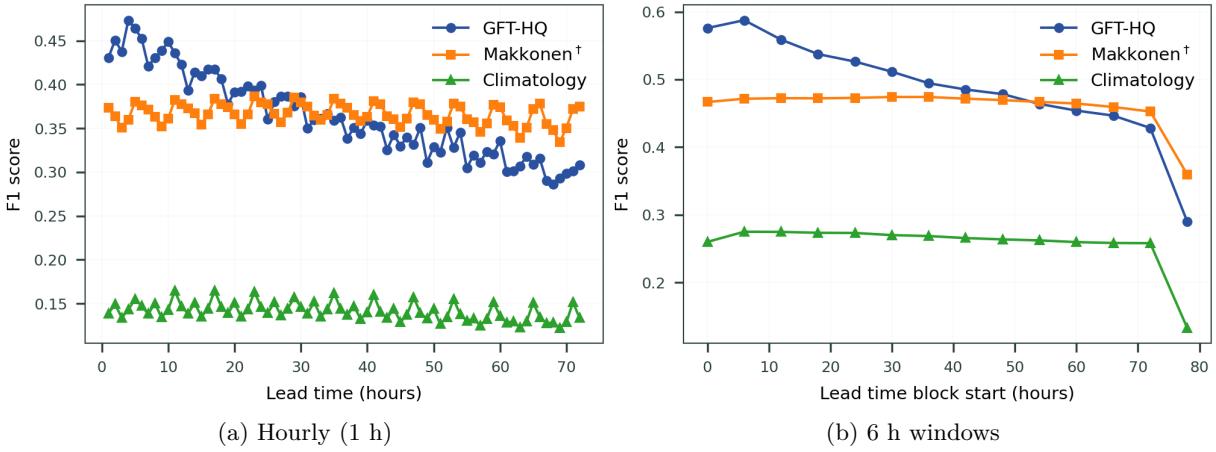
(a) Hourly (1 h)

(b) 6 h windows

Figure 14: Lead-time dependence of transmission-line rime-ice detection skill. Left (**a**): hourly F1 vs. lead time (1–72 h). Right (**b**): F1 for 6-hour windowed targets ("any icing in window") vs. window start time. Aggregating to 6 h smooths timing errors and better reflects current HQ operations, producing higher and more stable scores. In both settings, GFT-HQ achieves the highest skill and is comparable to or better than the ERA5-derived Makkonen reference through 72 h; climatology trails substantially. All scores are pooled across the 14 Sygivre stations.



Figure 15: Lead-time dependence of wind-farm icing risk forecasts. F1 score vs. lead time (1-72h), pooled across all windfarms. GFT-HQ maintains the highest skill at all leads ($\approx 0.65$ at short leads, tapering to $\approx 0.58$ by 72 h).

## A.5   Additional Results

### A.5.1   Rime ice forecasting for transmission lines

Figure 14a shows hour-by-hour F1. Figure 14b evaluates 6-hour windows, which is closer to current HQ operations; this aggregation reduces penalties from small timing offsets and leads to higher, more stable scores across lead times while preserving the ranking between methods.

### A.5.2   Wind-farm wind speed and ice forecasting

**Consistent icing risk forecasts**   Figure 15 compares the F1 score over different lead-times of the forecast with GFT-HQ maintaining high skill for more than 2-day lead times.

As can be observed in Figure 16, corresponding ROC performance is high as well (AUC = 0.93 vs 0.86–0.87 (Climatology) and 0.59–0.71 (Makkonen) in the hourly and 24-hour windowed setups.

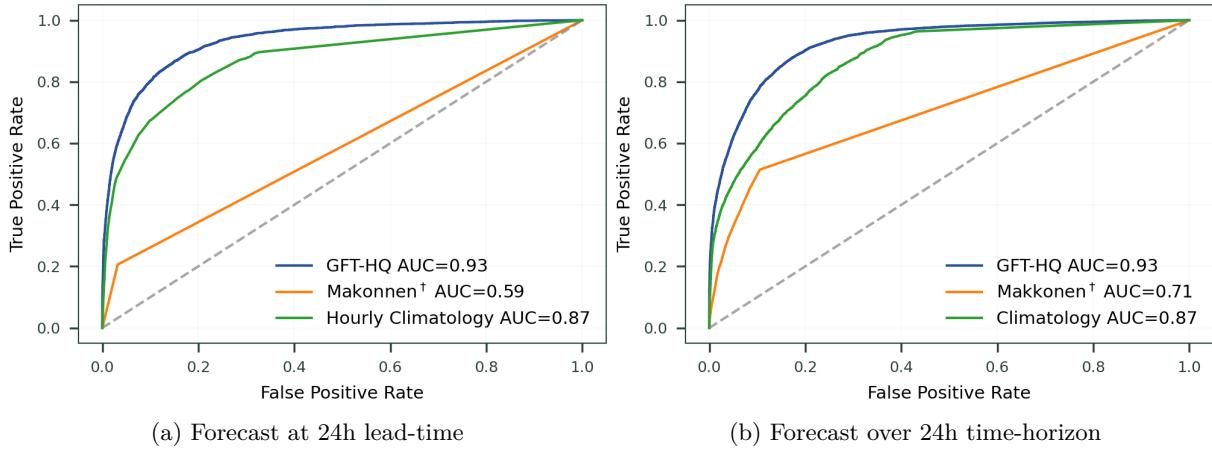(a) Forecast at 24h lead-time  (b) Forecast over 24h time-horizon

Figure 16: Wind-farm icing ROC-AUC curves.

**Consistent Performance Across Diurnal and Seasonal Cycles.** We analyzed the hub-height wind forecast skill of GFT-HQ against baseline models across different times of day (Figure 17) and seasons (Figure 18). The results demonstrate a robust and consistent performance advantage for the finetuned model.

A key finding is that GFT-HQ significantly dampens the error oscillations present in the baseline models. While the original GFT exhibits a pronounced diurnal error ripple, GFT-HQ reduces both the amplitude and the mean of this cycle. This leads to a more stable error profile that persists across all forecast leads.

This benefit is not confined to a particular time of day or season. GFT-HQ achieves the lowest Mean Absolute Error (MAE) in all diurnal periods (Morning, Afternoon, Evening, Night) and across all four seasons. Notably, its performance gains are largest during evening and nighttime hours, which are operationally critical periods of high energy demand and icing risk. While absolute forecast errors are naturally highest in winter and lowest in summer, GFT-HQ's performance advantage over ECMWF-IFS and GFT remains constant year-round.

Operationally, this multi-faceted consistency is highly valuable. The reduced diurnal and seasonal error volatility translates directly into more reliable and trustworthy day-ahead wind power forecasts, supporting more efficient grid management and resource planning throughout the year.

### A.5.3  Temperature and precipitation forecasting

**Spatial structure of temperature gains**  Figure 19 maps station-wise skill vs ECMWF-IFS (positive = lower error than IFS) at successive leads. Two patterns stand out:

- Finetuning amplifies, it doesn't relocate, the skill. The spatial pattern for GFT-HQ mirrors that of the untuned GFT, but with larger positive skill at virtually every lead. By 24-48 h the maps are dominated by positive (blue) stations for GFT-HQ, while the untuned GFT shows a similar but noticeably weaker footprint. This is expected given the finetuning setup here.

- Gains grow with lead time and remain geographically broad. Isolated neutral/negative pockets appear at short leads, but the fraction of improved stations increases monotonically with lead, consistent with Figure 11b.

*Operationally*, this means the temperature MAE reduction is system-wide rather than site-specific: load-forecast errors should fall across the control area, not just at a handful of stations, and the advantage persists into the day-ahead window where commitment and hedging decisions are made.
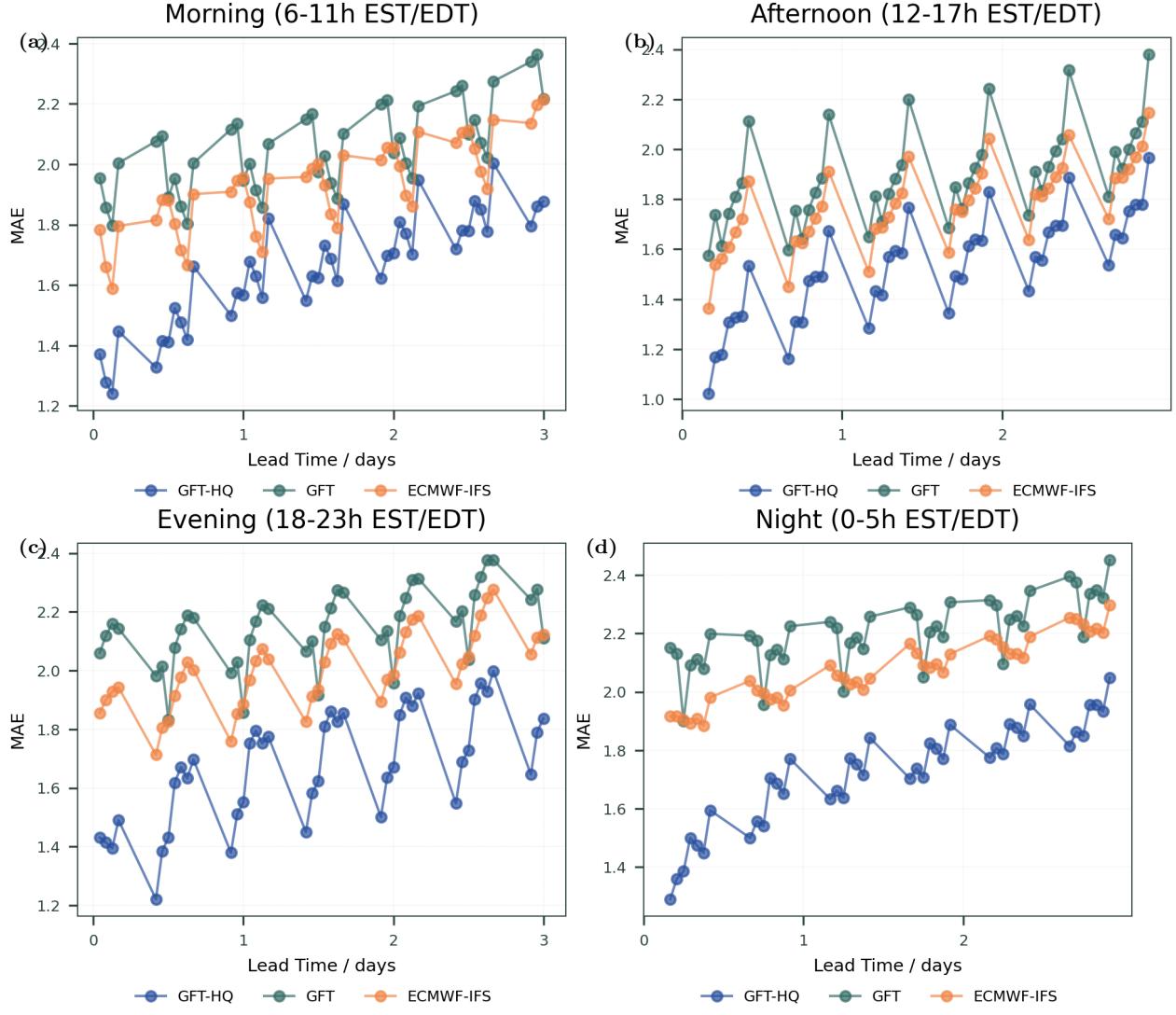
Figure 17: Hourly hub-height wind speed forecast improvements from fine-tuning: Station-averaged MAE for GFT-HQ, GFT, and ECMWF-IFS, grouped by local time of the day. GFT-HQ is lowest-error across all panels and leads, with the largest improvements in *Evening/Night* and a sustained gap through the day-ahead window. (MAE units are m s$^{-1}$; vertical oscillations reflect the diurnal verification hour across lead time.

**Spatial structure of precipitation gains** Figure 20 maps station-level skill vs. ECMWF-IFS for hourly total precipitation at successive leads (positive = lower error than IFS). The pattern is similar to improvements in temperature forecasts. Operationally, this means tighter basin-wide accumulation signals for hydro-inflow at the horizons HQ actually commits resources.
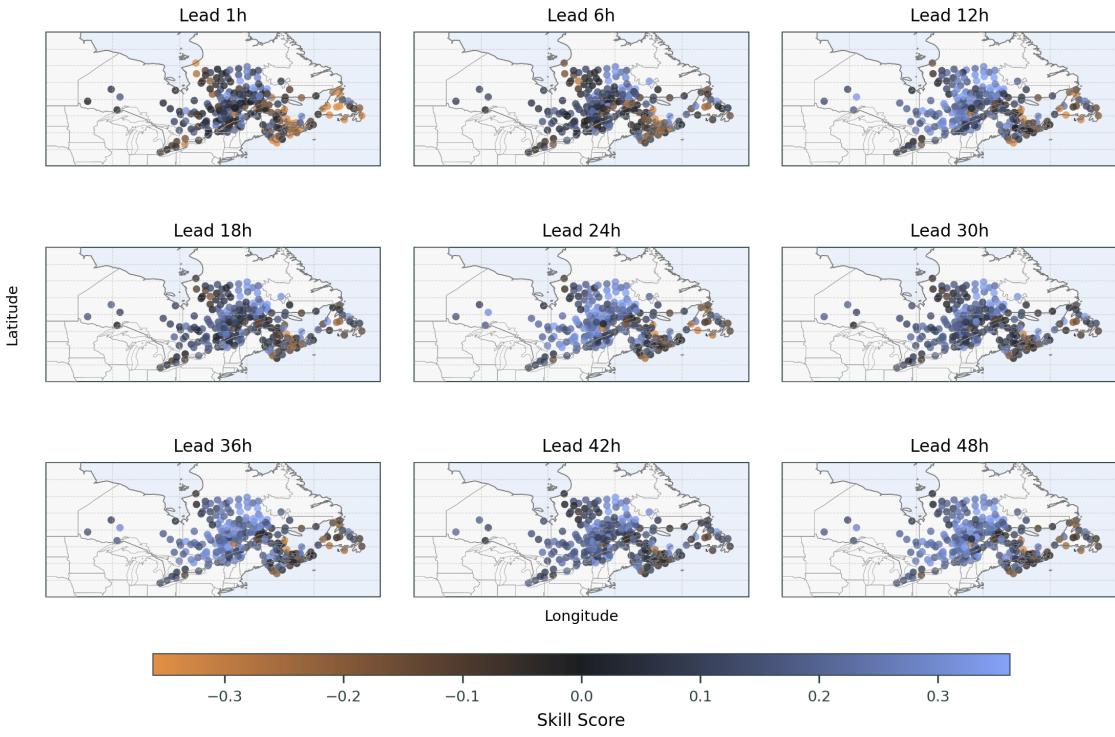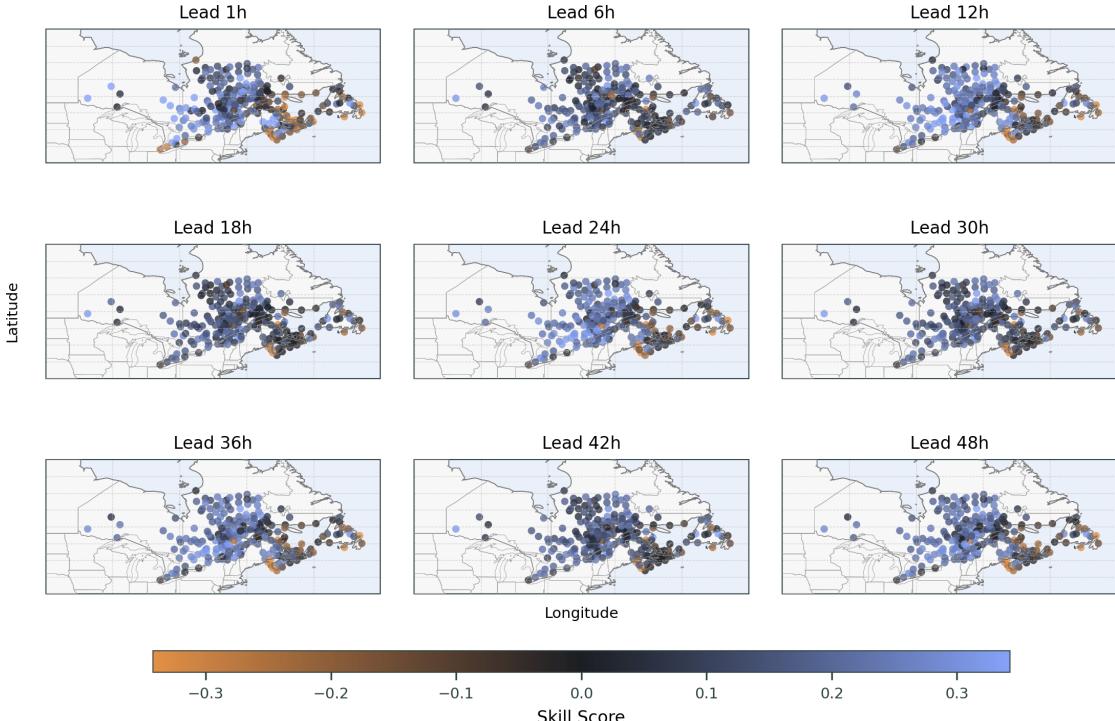
Figure 18: Hub-height wind speed MAE by season: Station-averaged MAE for GFT-HQ, GFT, and ECMWF-IFS, grouped by season. GFT-HQ improves on skill in all seasons. However, the error patterns suggest that it could be helpful to oversample winter season to further improve performance in that time-period.

(a) GFT-HQ temperature forecast skill improvement compared to ECMWF-IFS at different lead-times.



(b) GFT temperature forecast skill improvement compared to ECMWF-IFS at different lead-times.

Figure 19: Spatial distribution of GFT-HQ relative skill improvements compared to original GFT's skill against ECMWF-IFS. Higher score is better. GFT-HQ improved skill everywhere at longer lead-times, but it had a similar spatial distribution as the original GFT.

(a) GFT-HQ hourly total precipitation forecast skill improvement compared to ECMWF-IFS at different lead-times.



(b) GFT hourly total precipitation forecast skill improvement compared to ECMWF-IFS at different lead-times.
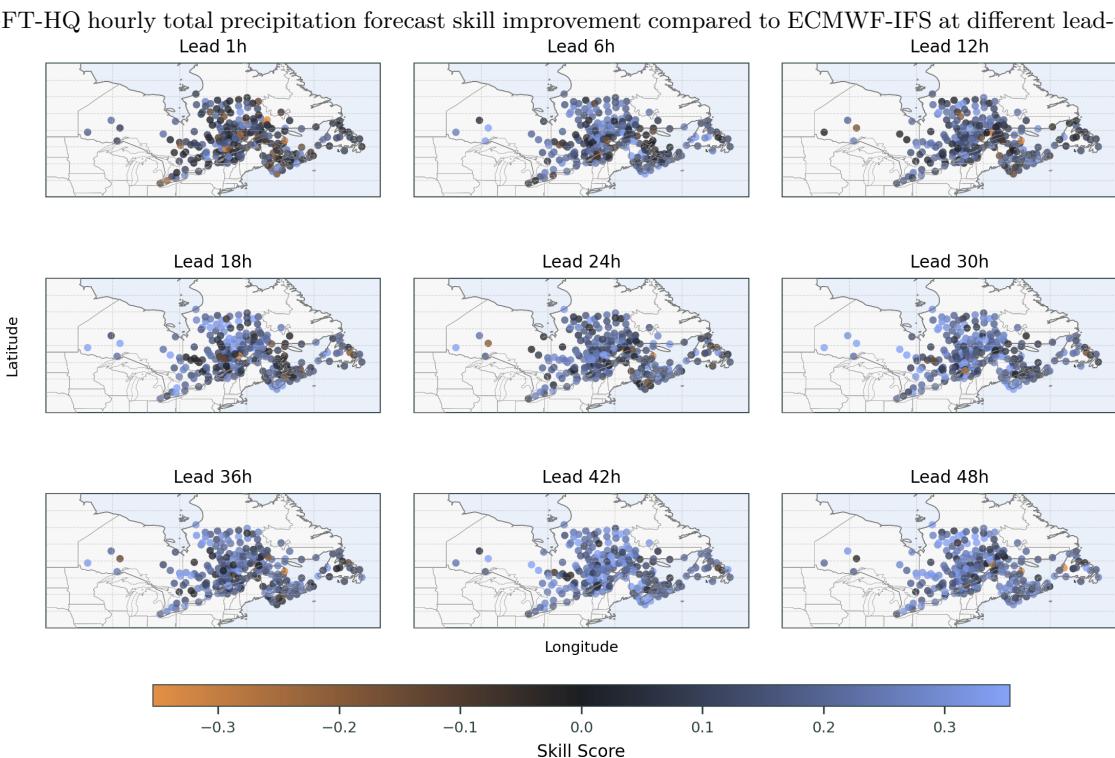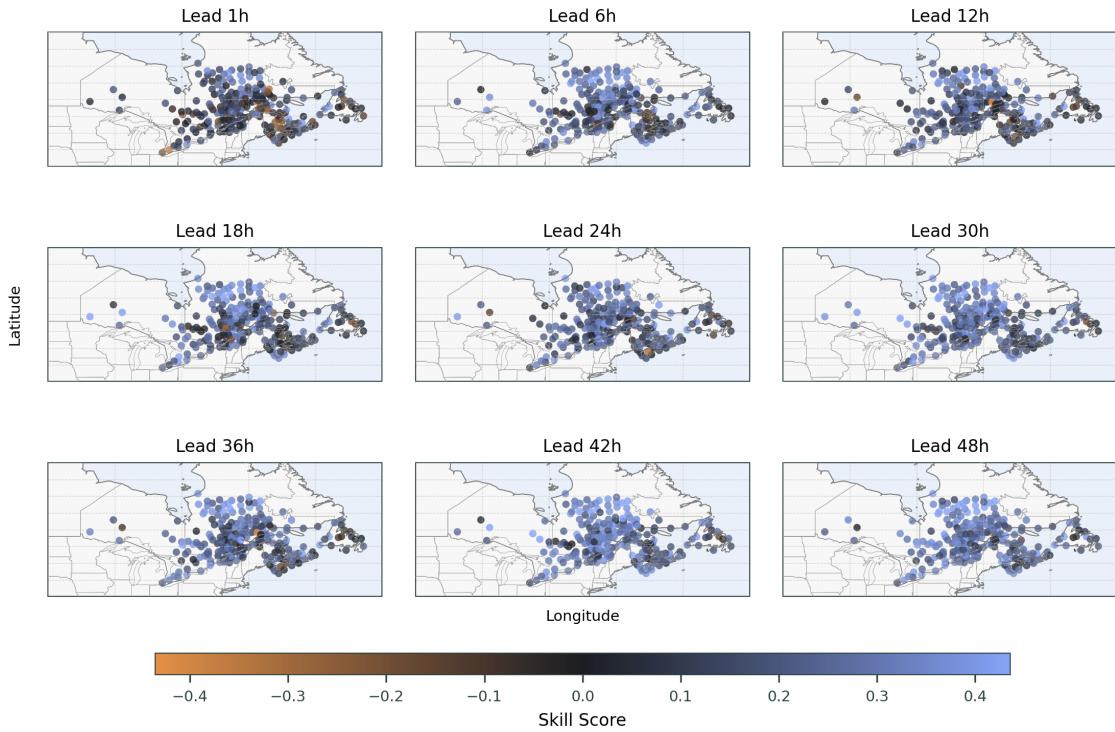
Figure 20: Spatial distribution of GFT-HQ's hourly precipitation skill improvements compared to original GFT's skill against ECMWF-IFS. Higher score is better. Finetuning increases the magnitude and coverage of positive skill at all leads.
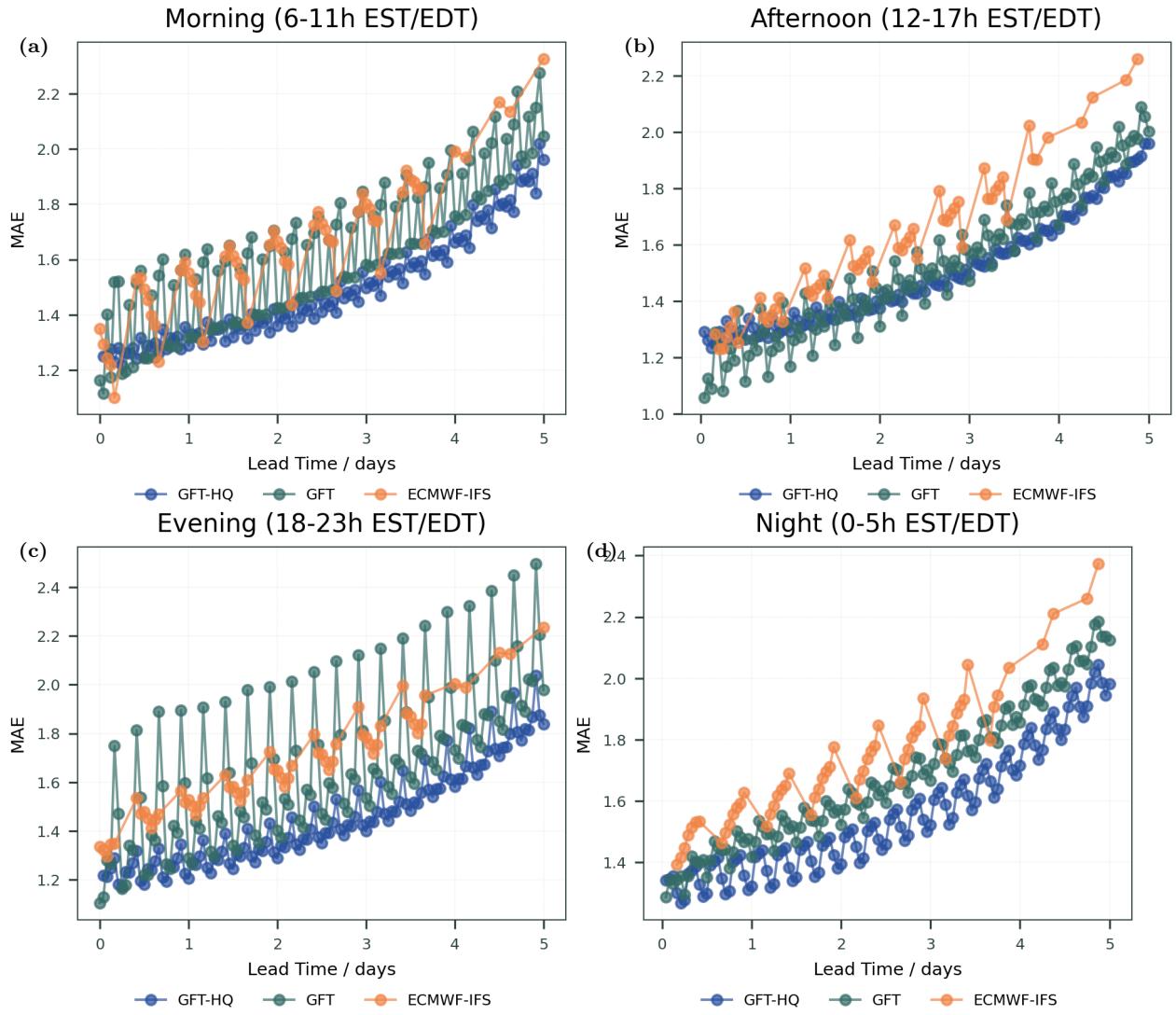
Figure 21: Surface temperature MAE by time of day. GFT-HQ outperforms other baselines under all conditions, except in the afternoons against raw GFT.
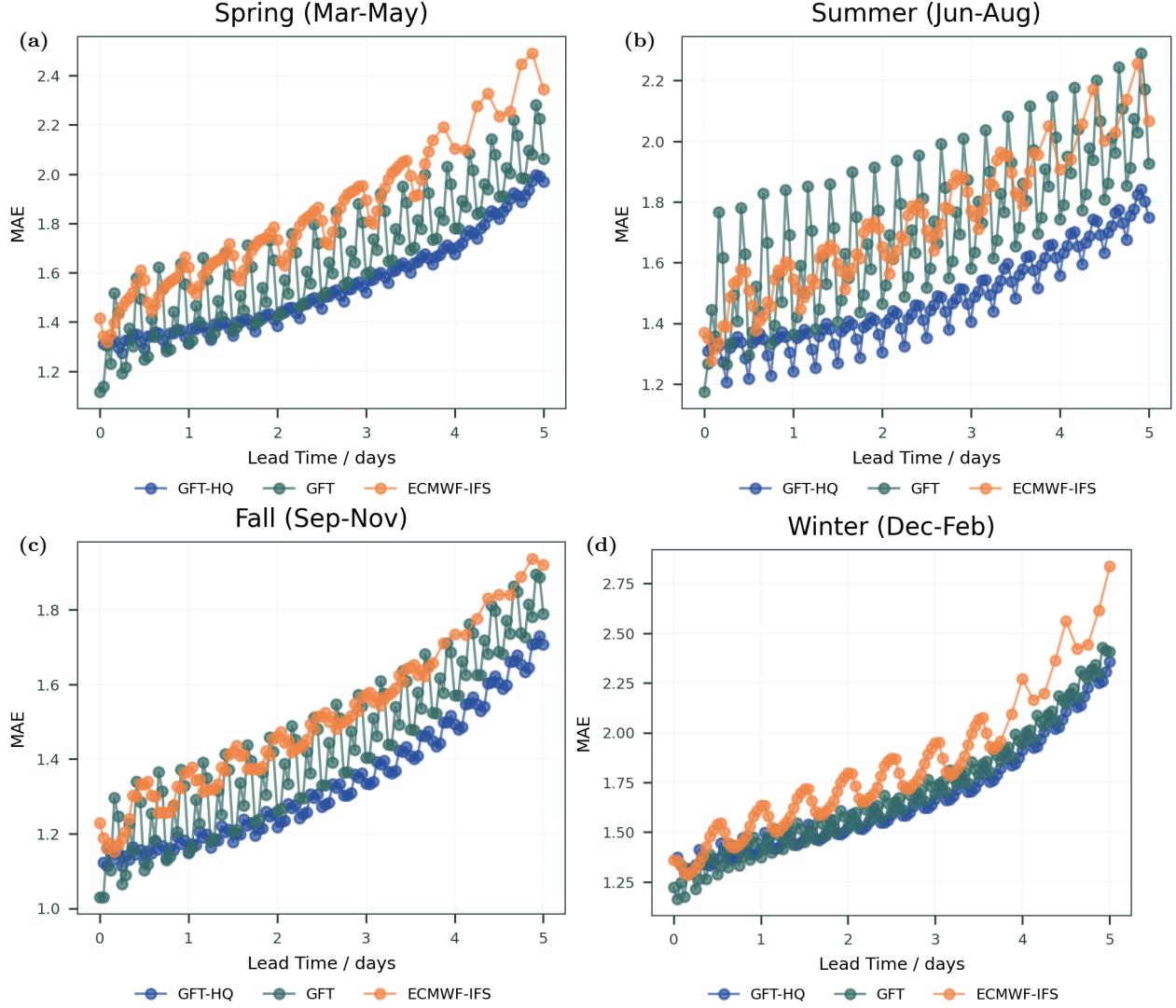
Figure 22: Surface temperature MAE by season. GFT-HQ outperforms other baselines under all conditions for 6-120h lead-times. However, the error patterns suggest that it could be helpful to oversample winter season during future fine-tuning.