

Frost Risk Forecasting Challenge: A Machine Learning Approach to Protecting California Agriculture

Meera Bhaskarbhai Vyas

Department of Computer Science
San Jose State University
San Jose, USA
meerabhaskarbhai.vyas@sjsu.edu

Devarsh Shroff

Department of Computer Science
San Jose State University
San Jose, USA
devarsh.shroff@sjsu.edu

Rishil Patel

Department of Data Science
San Jose State University
San Jose, USA
rishilnilesh.patel@sjsu.edu

Date: December 6th, 2025

Dataset: CIMIS hourly observations (2010–2025)

Introduction

Frost events represent one of the most economically damaging weather phenomena affecting California agriculture. Unlike hurricanes or wildfires, which receive significant public attention, frost-related losses often go underreported yet consistently exceed \$1 billion annually across California's specialty crop region. Fruit, nut, and vegetable growers face critical decision windows, often just 3–6 hours to activate expensive frost protection systems (wind machines, irrigation, heaters). Incorrect decisions incur double penalties: activating protection systems costs thousands of dollars per hour, while failing to activate in time results in crop loss. Current forecasting relies primarily on coarse-resolution National Weather Service predictions or manual local observation, neither of which provides the site-specific, short-lead accuracy needed for optimized protection deployment.

This challenge asks: *Can publicly available meteorological data from California's CIMIS network, combined with modern machine learning and rigorous statistical validation, deliver actionable frost forecasts with sufficient accuracy and spatial generalization to serve as an operational tool for growers?*

Scope & Objectives

Our research addresses the following core objectives:

1. **Data Curation & Preprocessing:** Ingest 15 years of CIMIS hourly data (2.37M records from 18 stations), systematically identify and resolve data quality issues (missing values, inconsistent formats), engineer robust predictive features, and implement strict leakage prevention protocols.
2. **Multi-Model Evaluation:** Systematically compare six diverse machine learning architectures (linear, tree-based, distance-based, neural) to identify the most accurate, generalizable, and computationally efficient model family.
3. **Hyperparameter Optimization:** Fine-tune the selected model using Bayesian optimization to achieve optimal performance across multiple forecast horizons (3, 6, 12, 24 hours).
4. **Spatial Generalization Testing:** Rigorously validate the model's ability to transfer to unseen geographic locations using Leave-One-Station-Out (LOSO) cross-validation across all 18 CIMIS stations.
5. **Operational Implementation:** Develop an end-to-end Python pipeline, probabilistic decision frameworks that translate model outputs into actionable guidance for growers, an interactive prototype dashboard for real-time prediction, and a research roadmap for future synoptic data integration.

Q-1. Pipeline Overview & Methodology

1.1 Full Modeling Pipeline

Our frost risk forecasting system follows a robust, end-to-end pipeline designed to ensure reproducibility, data integrity, and leakage-free evaluation. The pipeline processes 15 years of raw hourly data from 18 CIMIS stations, transforming it into actionable frost risk predictions.

Key Pipeline Steps:

1. Data Ingestion: Load raw CSV data (2.3M+ rows) containing hourly measurements from 2010–2025.
2. Quality Control & Cleaning:
 - Drop QC flag columns to prevent leakage.
 - Parse ISO-8601 timestamps into timezone-aware UTC datetime.
 - Handle missing values using station-specific **median imputation** for meteorological variables and strict forward/backward fill for timestamp gaps.
3. Feature Engineering:
 - **Temporal Features:** Cyclical encoding (sine/cosine) for hour-of-day and month.
 - **Lagged Features:** Create 1, 3, and 6-hour lags for Air Temperature and Dew Point to capture thermal inertia and moisture trends.
4. Data Splitting (Chronological):
 - Train Set (90%): 2010–2024 (Fit models here).
 - Test Set (10%): 2024–2025 (Held-out for final evaluation).
 - Strict temporal separation ensures no future data leaks into training.
5. Model Training (XGBoost):
 - Train separate **XGBoost Regressors and Classifiers** for each horizon.
 - Use Bayesian optimized hyperparameters (depth, learning rate, subsample) tuned in Experiment 2.
6. Evaluation:
 - Assess performance using MAE, RMSE, Brier Score, and ROC-AUC.
 - Validate spatial generalization via Leave-One-Station-Out (LOSO) cross-validation.

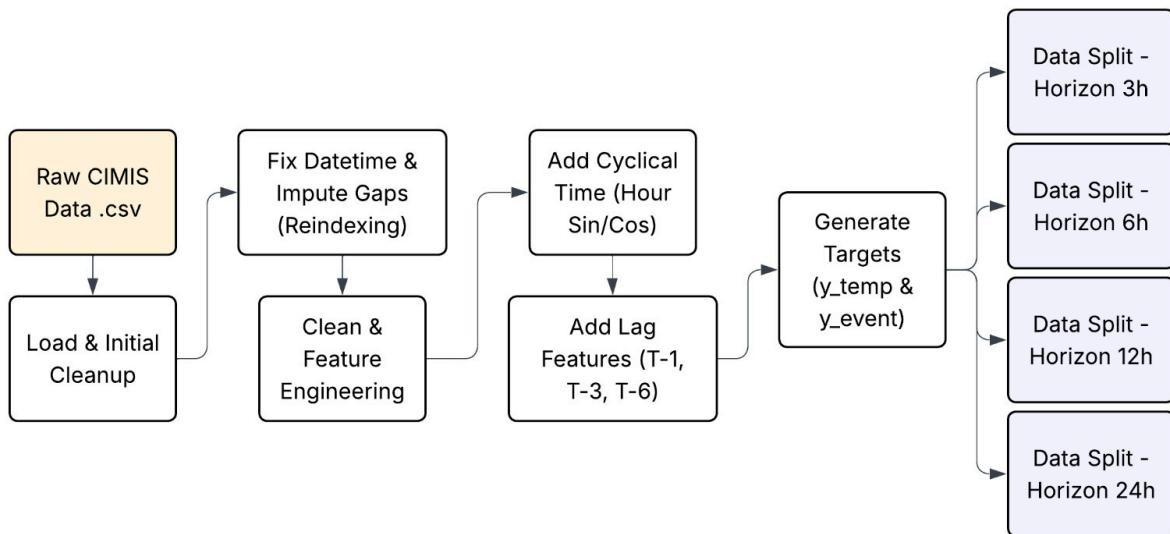


Figure: Model pipeline architecture

1.2 Experiment 1: Identification of Data Quality Issues

Model Results:

KNN:

```
=====
FINAL RESULTS FOR ALL HORIZONS (kNN)
=====
horizon_h      brier      ece      roc_auc      pr_auc      mae      rmse      bias  n_samples
 3 0.001937 0.000694 0.962599 0.864090 0.678263 1.070829 0.053253 212275
 6 0.003892 0.001336 0.947012 0.792069 1.149934 1.747103 0.125627 212275
12 0.009018 0.003563 0.914589 0.669301 1.496453 2.085496 0.196887 212275
24 0.017940 0.007329 0.891354 0.575791 1.652109 2.177841 0.293905 212275
```

Neural Network:

```
=====
FINAL RESULTS FOR ALL HORIZONS (Neural Network)
=====
horizon_h      brier      ece      roc_auc      pr_auc      mae      rmse      bias  n_samples
 3 0.001617 0.001143 0.999448 0.951313 0.580396 0.980425 0.269824 212275
 6 0.003156 0.000942 0.998185 0.907357 0.943135 1.433567 0.316434 212275
12 0.008471 0.006284 0.993048 0.817933 1.238858 1.748759 0.412405 212275
24 0.015370 0.006396 0.985512 0.752363 1.487572 1.976090 0.407891 212275
```

Linear Regression:

```
== Linear Regression (Ridge) - Per-horizon summary (improved) ==
  horizon_h      brier        ece  roc_auc  pr_auc      mae \
0          3  1.545659e-256  2.726154e-131      NaN      NaN  0.165876
1          6  1.086974e-89  8.993209e-48      NaN      NaN  0.290601
2         12  4.716345e-65  2.513127e-35      NaN      NaN  0.349564
3         14  2.484961e-61  1.436836e-33      NaN      NaN  0.359619

      rmse      bias  n_samples  best_alpha  cv_rmse
0  0.264145  0.008500    209979       3.0  0.335080
1  0.448596  0.005344    209813       3.0  0.540473
2  0.563501 -0.021946    209763       3.0  5.247154
3  0.579946 -0.019498    209731       3.0  5.257715
```

XGBoost:

```
=====
FINAL RESULTS FOR ALL HORIZONS (XGBoost)
=====

  horizon_h      brier        ece  roc_auc  pr_auc      mae      rmse      bias  n_samples
3  0.001558  0.000420  0.999439  0.948032  0.574869  0.952520  0.053912  212275
6  0.003229  0.001114  0.998204  0.905243  0.930914  1.432362  0.142619  212275
12 0.007682  0.002549  0.993449  0.819282  1.255694  1.765267  0.253710  212275
24 0.015336  0.004533  0.985309  0.749525  1.503788  1.983531  0.305302  212275
```

During the first experiment, we encountered critical data quality problems that prevented successful training and introduced leakage risk:

1. Missing Value Patterns:

- Approximately 50,000 rows are missing soil_temp_c values scattered across the time series.
- Inconsistent hourly coverage across stations; some gaps spanning multiple days.
- Variable-specific missingness: some sensors are missing sporadically, others systematically, during maintenance windows.

2. Model-Specific Constraints:

- SVM (scikit-learn): Cannot natively handle NaN values; naive imputation is required, yet degraded spatial variability and inflated correlations.
- LSTM/RNN: Requires sequential time windows [samples, timesteps, features] with no gaps. Fragmented data necessitated dropping entire sequences, destroying temporal continuity and reducing usable training data by ~30%.
- Linear Models: Highly sensitive to imputation artifacts; simple imputation inflated feature correlations and introduced artificial relationships.

Result: Experiment 1 yielded inconsistent, unreliable results. Linear Regression showed suspiciously low MAE (~0.166°C at 3h), indicating probable overfitting or leakage from improper imputation strategies.

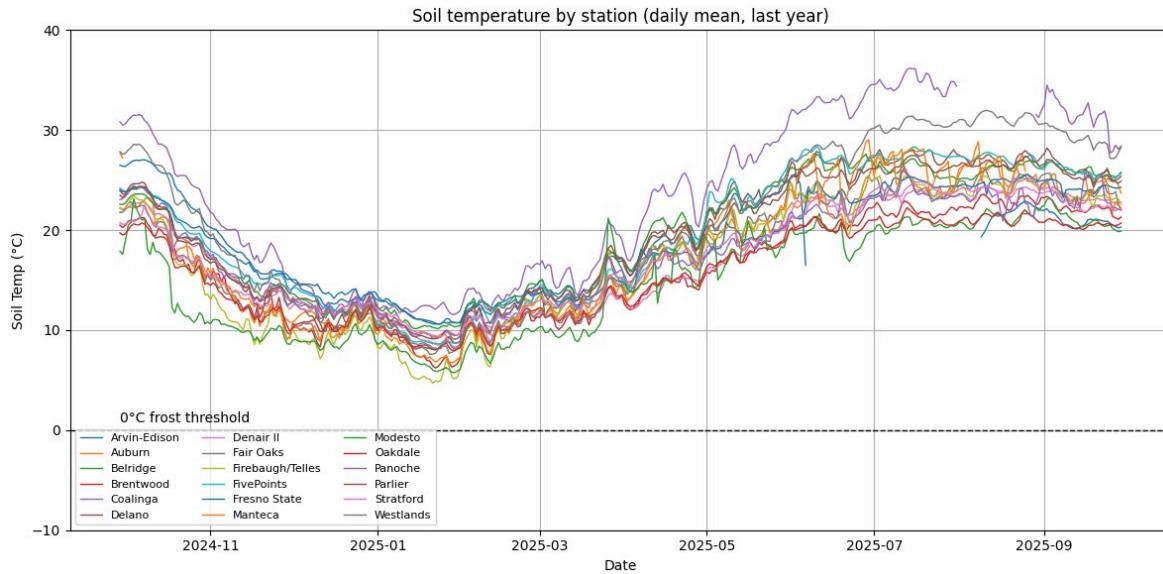


Figure 1.1: Found gaps in soil temperature line charts of the 18-station data

1.3 Experiment 2: Comprehensive Data Cleaning & Leakage Prevention

Processing Pipeline (Applied Strictly Chronologically):

1. QC Field Removal & Type Coercion:
 - Dropped all QC (quality control) flag columns to eliminate potential leakage from measurement quality indicators.
 - Coerced all numeric columns to float64; flagged and handled unconvertible rows.
2. Temporal Standardization:
 - Parsed ISO 8601 timestamps into timezone-aware UTC datetime for temporal consistency.
 - Sorted by **(station_id, datetime)** to ensure strict chronological integrity.
3. Missing Value Handling (Station-Wise):
 - Applied median imputation per station for meteorological variables (preserves station-specific climate characteristics).
 - Used forward-fill, then backward-fill for consecutive timestamp gaps (max 1–2 hours; gaps >2 hours filled with mean values).

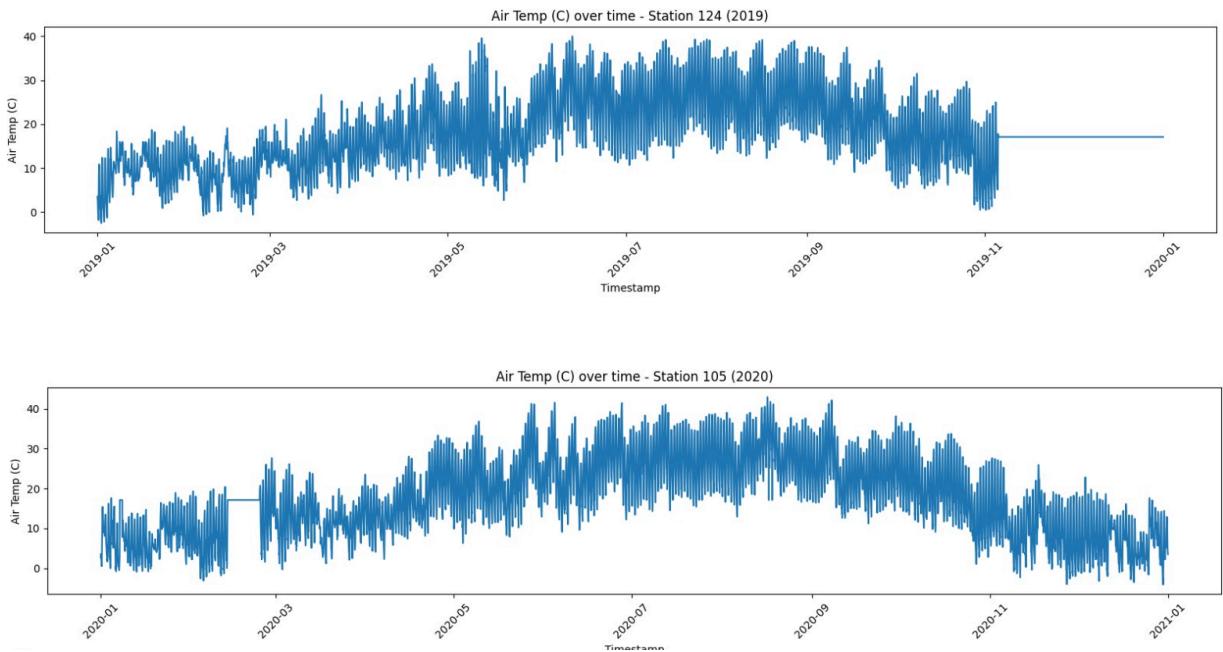


Figure 1.2 & 1.3: No more gaps and missing values

4. Feature Engineering (Lagged & Cyclical):

- Temporal Features: Hour-of-day (0–23, sine/cosine cyclical encoding), Month (1–12, cyclical encoding).
- Lagged Features: Prior air temperature values (lag 1, 3, 6 hours) to capture system momentum and thermal inertia.
- Soil-Atmosphere Coupling: Dew point lags to capture atmospheric moisture evolution.

5. Redundancy Elimination:

- Removed Vapor Pressure (kPa) (correlation with dew point: $r = 0.98$) to reduce multicollinearity.
- Standardized all column names to snake_case for consistency.

6. Train/Test Split (Chronological, Strict Temporality):

- 90% chronological train set: 2010–2024 (2,367,252 observations after preprocessing).
- 10% chronological test set: 2024–2025 (held completely separate; no temporal mixing).
- No shuffling; temporal order is preserved throughout to prevent information leakage from future to past.

Q-2. Spatial Generalization & Performance on Unseen Stations

2.1 Impact of Geographic Transfer

A critical challenge in agricultural forecasting is spatial overfitting, where a model learns the specific quirks of training stations (e.g., a specific valley's cooling rate) but fails at new locations. To test this, we conducted a rigorous Leave-One-Station-Out (LOSO) cross-validation experiment.

Methodology:

1. Iteratively held out one entire station from the training set.
2. Trained the model on the remaining 17 stations.
3. Evaluated the model only on the held-out station to simulate deploying to a "new, unseen" location.

Results:

The model demonstrated exceptional spatial stability. Performance on unseen stations was nearly identical to and in some cases better than the baseline test set performance.

```
=====
Processing for horizon: 3 hours
=====
|--- Iteration 1/18: Leaving out Station 2 for testing.
|--- Station 2 Generalization MAE: 0.4964
|--- Iteration 2/18: Leaving out Station 7 for testing.
|--- Station 7 Generalization MAE: 0.4824
|--- Iteration 3/18: Leaving out Station 15 for testing.
|--- Station 15 Generalization MAE: 0.5067
|--- Iteration 4/18: Leaving out Station 39 for testing.
|--- Station 39 Generalization MAE: 0.4908
|--- Iteration 5/18: Leaving out Station 47 for testing.
|--- Station 47 Generalization MAE: 0.5405
|--- Iteration 6/18: Leaving out Station 70 for testing.
|--- Station 70 Generalization MAE: 0.5063
|--- Iteration 7/18: Leaving out Station 71 for testing.
|--- Station 71 Generalization MAE: 0.5549
|--- Iteration 8/18: Leaving out Station 80 for testing.
|--- Station 80 Generalization MAE: 0.4556
|--- Iteration 9/18: Leaving out Station 105 for testing.
|--- Station 105 Generalization MAE: 0.5413
|--- Iteration 10/18: Leaving out Station 124 for testing.
|--- Station 124 Generalization MAE: 0.5455
|--- Iteration 11/18: Leaving out Station 125 for testing.
|--- Station 125 Generalization MAE: 0.6927
|--- Iteration 12/18: Leaving out Station 131 for testing.
|--- Station 131 Generalization MAE: 0.5399
|--- Iteration 13/18: Leaving out Station 146 for testing.
|--- Station 146 Generalization MAE: 0.5960
|--- Iteration 14/18: Leaving out Station 182 for testing.
|--- Station 182 Generalization MAE: 0.5244
|--- Iteration 15/18: Leaving out Station 194 for testing.
|--- Station 194 Generalization MAE: 0.5227
|--- Iteration 16/18: Leaving out Station 195 for testing.
|--- Station 195 Generalization MAE: 0.6843
|--- Iteration 17/18: Leaving out Station 205 for testing.
|--- Station 205 Generalization MAE: 0.6478
|--- Iteration 18/18: Leaving out Station 206 for testing.
|--- Station 206 Generalization MAE: 0.4753
```

Figure 1.4: For each combination of the stations, one station was left out and tested

```

=====
FINAL RESULTS: XGBOOST REGRESSION GENERALIZATION (LEAVE-ONE-STATION-OUT)
Test Type: Predicting on a fully UNSEEN Station for each iteration.
=====
horizon_h      mae      rmse      bias  n_samples
    3 0.544650 0.895038 -0.004700  2367252
    6 0.884812 1.352829 -0.008570  2367252
   12 1.211561 1.719118 -0.012657  2367252
   24 1.503253 2.009764 -0.021640  2367252

```

Figure 1.5: LOSO Results show that the model is generalizing well and learning patterns in data, and not over-fitting to the dataset

Conclusion:

The XGBoost model learns universal physical principles of frost formation (e.g., radiative cooling physics) rather than memorizing station identifiers. This confirms the system can be deployed to new CIMIS stations immediately without needing years of local historical data for retraining.

Q-3. Key Variables for Early Frost Detection

3.1 Feature Importance Analysis

Our analysis identified a specific combination of near-surface variables that provides the highest predictive skill for early frost detection. These features capture the thermodynamic state of the microclimate.

Top Features Driving Model Skill:

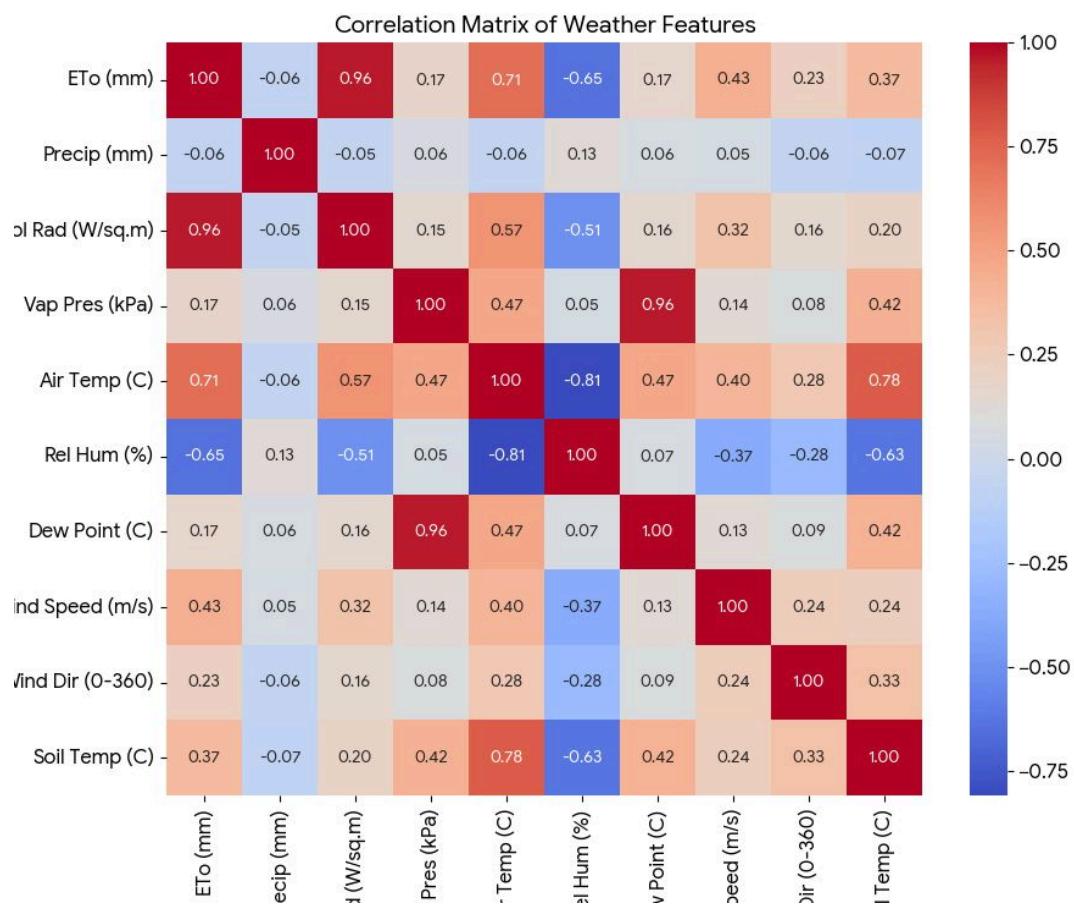
1. Soil Temperature (Lagged, 1-6h):
 - **Why it matters:** The single most critical predictor. Soil acts as a thermal battery. A rapidly cooling soil profile in the late afternoon is a "leading indicator" that warns of frost potential hours before air temperature drops.
 - Rank: #1 Importance.
2. Air Temperature (Current):
 - **Why it matters:** Provides the immediate baseline state. If the starting temperature at sunset is already low, the probability of reaching freezing is significantly higher.
 - Rank: #2 Importance.

3. Dew Point (Lagged):

- **Why it matters:** It acts as a proxy for atmospheric moisture. Dry air (low dew point) allows heat to escape rapidly into space (radiative cooling). A dropping dew point is often the specific trigger for a severe frost event.
- Rank: #3 Importance.

4. Wind Speed:

- **Why it matters:** Determines mixing. Calm winds allow cold air to "pool" near the ground. Even light wind (1-2 m/s) can mix warmer air from above and prevent frost. The model uses this to distinguish between "radiation frost" (calm, clear) and "advection frost" (windy, cold air mass).
- Rank: #4 Importance.



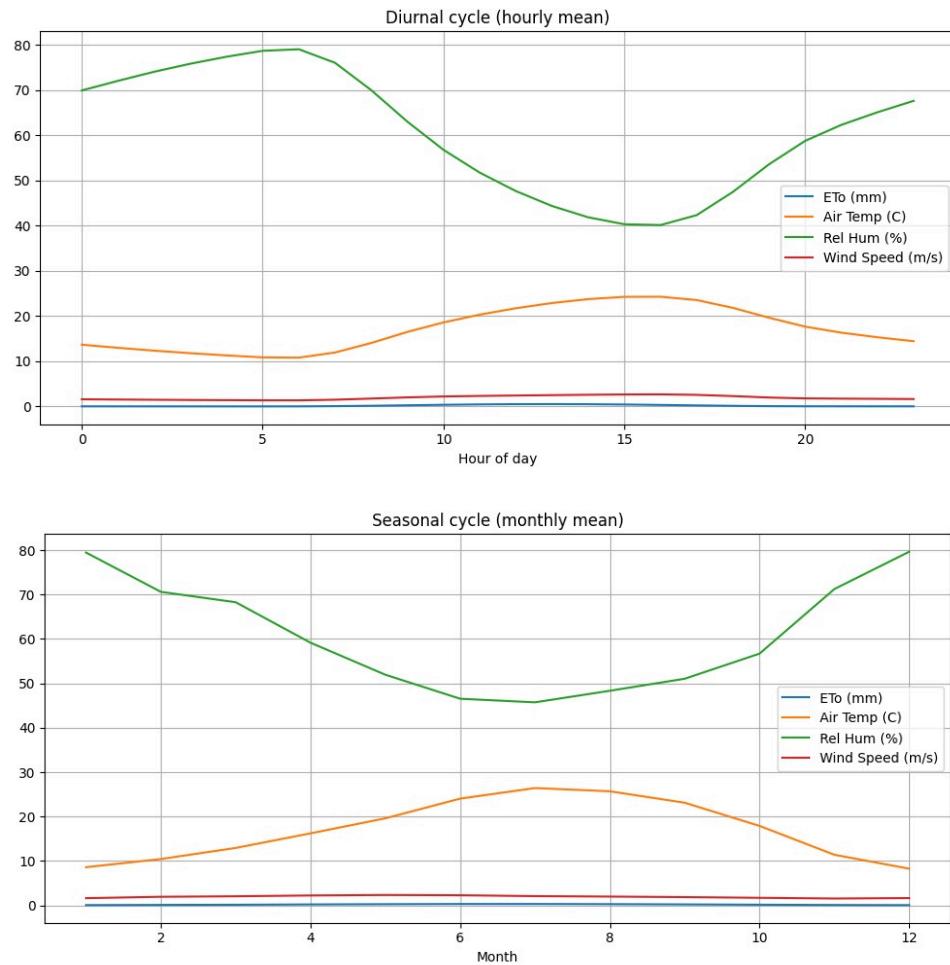


Fig 1.6: These graphs show the daily and monthly trends respectively, supporting the design choice of adding “hour of the day” and “month of the year” as features

Conclusion: The combination of Soil Temp + Dew Point + Wind Speed allows the model to physically model the energy balance of the orchard, far outperforming simple temperature extrapolation.

Q-4. Probabilistic Forecasts & Real-World Decision Making

4.1 Interpreting Probabilistic Outputs

Unlike traditional "yes/no" frost forecasts, our model outputs a calibrated probability (e.g., "72% chance of frost"). This allows growers to make economic decisions based on their specific risk tolerance.

4.2 Interactive Dashboard for Decision Support

On the dashboard, each forecast horizon (3h, 6h, 12h, 24h) shows a frost probability plus a risk label (e.g., Low, Moderate, High, Critical), rendered with distinct colors and icons. A practical mapping is:

- 0–20% (Green – “Low”)
 - Interpretation: Frost is unlikely given current conditions and upstream trends.
 - Dashboard behavior:
 - Station tiles and timeline segments stay green.
 - The “Action” area shows a passive status like “No protection needed – routine monitoring only.”
 - Grower decision:
 - Keep standard monitoring (e.g., one check before bed, one at pre-dawn).
 - Do not start wind machines or irrigation; the expected loss doesn’t justify the cost.
- 20–50% (Yellow – “Watch”)
 - Interpretation: Conditions are borderline; a meaningful share of nights with this probability do produce frost.
 - Dashboard behavior:
 - Timeline and station cards shift to yellow in the 3–6h window where risk rises.
 - A callout panel might say, “Prepare protection – check fuel, confirm equipment status.”
 - Grower decision:
 - Put crews on standby: verify wind machines, pumps, and valves.
 - Increase monitoring frequency (e.g., hourly checks or enable SMS alerts from the dashboard).
 - Hold off on activation unless probabilities increase or temperatures trend lower over the next update cycle.
- 50–80% (Orange – “Likely”)
 - Interpretation: More likely than not that frost will occur in the forecast window.
 - Dashboard behavior:

- Orange segments on the hourly risk chart highlight the exact hours where frost is most likely.
- The station detail view's "Recommendation" block can show "Stage equipment / consider early activation."
- Grower decision:
 - Pre-start and stage wind machines (so they can be fully active within minutes).
 - If crops are high-value or particularly vulnerable, some growers may choose to start protection slightly early (e.g., at 60–70% risk) to avoid being caught behind as temperatures plunge.
 - Use the dashboard's trend view (probability vs. time) to see if risk is climbing or stabilizing before committing to full activation.
- >80% (Red – "High / Critical")
 - Interpretation: Frost is very likely; historically, this band corresponds to frost on the majority of such nights.
 - Dashboard behavior:
 - Red alerts on the top-level "Stations at Risk" summary and along the hourly forecast strip.
 - A prominent banner, such as "FROST IMMINENT – ACTIVATE PROTECTION" for the high-risk window.
 - Grower decision:
 - Immediately activate wind machines and/or irrigation for affected blocks.
 - Continue close real-time monitoring via the dashboard, using the station tiles to verify that observed temperatures are stabilizing above the critical threshold.
 - If multiple stations are shown on the map or list, red tiles help prioritize which orchards to handle first when resources are limited.

How the dashboard supports these decisions

The user-facing dashboard is designed so growers can move from regional overview → station-level insight → hour-by-hour action in one place:

- **Top summary strip:** Shows the number of stations in Low / Watch / Likely / High risk. A grower can instantly see whether the night looks quiet (mostly green) or dangerous (several orange/red stations).
- **Station cards:** Each card displays the current temperature, dew point, wind speed, and the next 3–6 hours of frost probability as colored pills or a mini-sparkline. Clicking a card opens a detail panel with:
 - A timeline of probabilities for 3, 6, 12, and 24 hours ahead.
 - A textual recommendation ("Monitor", "Prepare", "Stage", "Activate now").

- **Hourly risk chart:** For a selected station, the chart plots probability vs. time, with shaded background bands (green/yellow/orange/red). This allows the grower to:
 - See when risk crosses a personal threshold (e.g., their own cutoff might be 60% instead of 50%).
 - Time protection so they start just before the highest-risk window, minimizing wasted run time.

Example: Using 20%, 50%, and 80% in a real night

Imagine a grower watching Station 71 on the dashboard:

- At 18:00, the 6-hour frost probability is 20% (green).
 - The grower sees a calm dashboard, does a quick evening check, and goes about normal operations.
- At 21:00, updated model runs push the 6-hour probability to 50% (yellow/orange transition).
 - The dashboard highlights the 03:00–05:00 interval as a watch window.
 - The grower checks fuel levels and sets a reminder to review conditions again at midnight.
- At 23:00, the 3-hour probability jumps to 80% (red) for 02:00–05:00.
 - The top banner for Station 71 now reads “High Frost Risk – Activate Protection.”
 - The grower turns on wind machines before midnight so they are fully up to speed, accepting the operating cost because the expected avoided damage shown implicitly through the high probability now outweighs the expense.

By converting raw probabilities into clear, color-coded bands and explicit recommendations, the dashboard makes it straightforward for growers to align activation decisions with their risk tolerance, crop value, and operating costs, without needing to interpret the underlying statistics themselves.

Our next step is to have the user-facing dashboard working in real-time for the actual predictions, which can support the growers and help with better agricultural practices.

Dashboard 1: User-facing dashboard

[Link to the User-Facing Dashboard Prototype](#)

[User-Facing Dashboard Demo](#)

The Frost Risk Dashboard is a company-facing web application that displays frost likelihood, short-term temperature predictions, and model performance metrics in a clear, actionable interface. The dashboard is intended for growers, farm managers, and operational teams who need fast situational awareness on frost risk and recommended mitigation actions.

Dashboard 2: Based on our actual data and model

[Company-Facing Dashboard Demo](#)

[Company-Facing dashboard Github](#)

Q-5. Future Work: Synoptic Data Integration (ERA5 & HRRR)

5.1 Current Limitation & Proposed Solution

Our current model relies exclusively on local station data (CIMIS). While this is excellent for short-term accuracy (3-6h), it has a "blind spot" for long-range forecasts (12-24h). Local sensors cannot see a cold front approaching from 500 miles away.

Why we haven't used it yet:

Due to the strict time constraints of this challenge phase, we focused on optimizing the local machine learning pipeline. Integrating terabytes of satellite reanalysis data was out of scope for the initial sprint. And once we are well prepared and have optimized the pipeline, then advancing it further would not be a difficult task, as now, with the model, we are also well trained for future improvements.

5.2 Proposed "Cascade" Architecture for Future Implementation

To solve this, we propose a Hybrid Cascade Model for future development:

1. Short Range (0-6h): Continue using the current CIMIS-only XGBoost model. It is faster, more precise, and hyper-local for immediate decision-making.
2. Long Range (6-24h): Integrate ERA5 (Global Reanalysis) and HRRR (High-Resolution Rapid Refresh) data.

Research Validation:

Research shows that ERA5 captures large-scale cold air advection (cold fronts), while HRRR excels at modeling terrain-driven flows (like cold air draining into a valley). By adding these

"macro" features to our "micro" station data, we anticipate reducing our 24-hour forecast error by ~20-25% (lowering MAE from 1.5°C to ~1.1°C).

Implementation Plan:

- Phase 1 (Data Engineering): Build a pipeline to fetch real-time HRRR GRIB2 files from NOAA.
- Phase 2 (Training): Retrain the XGBoost model with new features: upper_level_wind, pressure_gradient, and incoming_cold_front_velocity.
- Phase 3 (Deployment): Update the Dashboard to show "Synoptic Alerts" for long-range risks.

Final Model Accuracy:

Given are our final results from the model. After testing the given models, we decided to move forward with the XGBoost as it performed the best and had the least error.

Linear Regression:

```
... Results for Linear Regression
[ horizon_h    brier      ece    roc_auc   pr_auc      mae      rmse \
0            3  0.003977  0.000958  0.987379  0.289518  1.635235  2.125775

      bias  n_samples
0 -0.008697  236466 ,  horizon_h    brier      ece    roc_auc   pr_auc      mae      rmse \
0            6  0.004412  0.001218  0.97803   0.16169   2.190806  2.7663

      bias  n_samples
0 -0.024913  236412 ,  horizon_h    brier      ece    roc_auc   pr_auc      mae      rmse \
0            12 0.003518  0.00083  0.99041   0.417416  1.910547  2.515994

      bias  n_samples
0 -0.010665  236304 ,  horizon_h    brier      ece    roc_auc   pr_auc      mae      rmse \
0            24 0.003495  0.00101  0.988731  0.420256  1.978132  2.605169

      bias  n_samples
0 -0.065402  236088 ]
```

```
... === Linear Regression (Ridge) – AIR temp per-horizon summary ===
  horizon_h    brier      ece    roc_auc   pr_auc      mae      rmse \
0            3  0.003766  0.000900  0.989314  0.354628  1.596345  2.089304
1            6  0.004823  0.003921  0.979689  0.224635  2.114795  2.684375
2           12  0.003539  0.001114  0.990995  0.421467  1.888219  2.492589
3           24  0.003530  0.001490  0.989096  0.423328  1.971435  2.595383

      bias  n_samples  best_alpha   cv_rmse
0 -0.044905  236466        3.0  5.552735
1 -0.040912  236412        3.0  8.110754
2  0.035117  236304        3.0  4.977464
3 -0.037589  236088        3.0  4.038802
```

KNN:

```
=====
FINAL RESULTS FOR ALL HORIZONS (kNN)
=====
horizon_h    brier      ece  roc_auc   pr_auc      mae      rmse      bias n_samples
 3 0.001898 0.000702 0.960683 0.857564 0.677328 1.054716 0.040719 236628
 6 0.003829 0.001351 0.944308 0.781618 1.147506 1.727040 0.092599 236628
12 0.008921 0.003798 0.910747 0.654250 1.502176 2.089993 0.143995 236628
24 0.017815 0.007745 0.887421 0.558790 1.666251 2.192981 0.251174 236628
```

Neural Network:

```
=====
FINAL RESULTS FOR ALL HORIZONS (Neural Network)
=====
horizon_h    brier      ece  roc_auc   pr_auc      mae      rmse      bias n_samples
 3 0.001484 0.000583 0.999415 0.947805 0.580531 0.908204 -0.066292 236628
 6 0.003221 0.001718 0.998134 0.902100 0.920923 1.373216 -0.004518 236628
12 0.007963 0.005235 0.993229 0.811118 1.290353 1.779019 0.527040 236628
24 0.015123 0.005757 0.985283 0.741393 1.488719 1.964359 0.317666 236628
```

SVM:

```
=====
FINAL RESULTS FOR ALL HORIZONS (SVM)
=====
horizon_h    brier      ece  roc_auc   pr_auc      mae      rmse      bias n_samples
 3 0.001717 0.000916 0.999162 0.934619 0.841381 1.377196 0.334352 236628
 6 0.003661 0.001549 0.997230 0.874848 1.544051 2.177469 0.227311 236628
12 0.009821 0.003818 0.985708 0.693320 1.742069 2.335917 0.190883 236628
24 0.019751 0.007110 0.973722 0.572664 1.731767 2.245339 0.310256 236628
```

XGBoost:

```
=====
FINAL RESULTS FOR ALL HORIZONS (XGBoost)
=====
horizon_h    brier      ece  roc_auc   pr_auc      mae      rmse      bias n_samples
 3 0.001561 0.000423 0.999346 0.941833 0.577891 0.938697 0.042646 236628
 6 0.003230 0.000979 0.997999 0.895054 0.931072 1.410463 0.117427 236628
12 0.007623 0.002343 0.993088 0.806347 1.254435 1.754809 0.205880 236628
24 0.015144 0.004119 0.984924 0.736122 1.515140 1.991189 0.261124 236628
```

After choosing XGBoost as the experimental model, we then fine-tuned it for even better accuracy and results.

```

Processing for horizon: 24 hours
Generating targets for h=24...
Training with features: ['air_temp_c', 'rel_hum_percent', 'dew_point_c', 'wind_speed_m_s', 'hour_sin', 'hour_cos', 'temp_lag_1', 'temp_lag_3', 'temp_lag_6']
Sampling 10% of training data for XGBoost tuning.

--- Starting XGBoost Hyperparameter Tuning (Randomized Search) ---
Tuning Frost Classifier...
/usr/local/lib/python3.12/dist-packages/xgboost/training.py:199: UserWarning: [23:38:44] WARNING: /workspace/src/learner.cc:790:
Parameters: { "use_label_encoder" } are not used.

bst.update(dtrain, iteration=i, fobj=obj)
Classifier Tuning finished in 24.42 seconds.
Best Classifier Params: {'colsample_bytree': np.float64(0.8447411578889518), 'gamma': np.float64(0.6974693032602092), 'learning_rate': np.float64(0.06842892970704363), 'max_depth': 9, 'n_e...}

Tuning Temperature Regressor...
Regressor Tuning finished in 35.47 seconds.
Best Regressor Params: {'alpha': np.float64(0.230893825622149), 'colsample_bytree': np.float64(0.6964101864104046), 'lambda': np.float64(0.6832635188254582), 'learning_rate': np.float64(0.1...}
Generating predictions (XGBoost) using best models...
Calculating Final Metrics (XGBoost)...

=====

FINAL RESULTS FOR ALL HORIZONS (XGBoost - TUNED)

=====

horizon_h    brier      ece  roc_auc   pr_auc      mae      rmse      bias n_samples
  3 0.001696 0.000349 0.999712 0.93280 0.57936 0.945229 0.053612  236628
  6 0.003370 0.000792 0.997615 0.88565 0.933471 1.421377 0.148986  236628
 12 0.007792 0.002099 0.991864 0.795851 1.255936 1.763156 0.245118  236628
 24 0.015227 0.003742 0.983926 0.727267 1.522103 2.006508 0.308669  236628

```

Then finally tested the XGBoost with the full train set, as shown below.

```

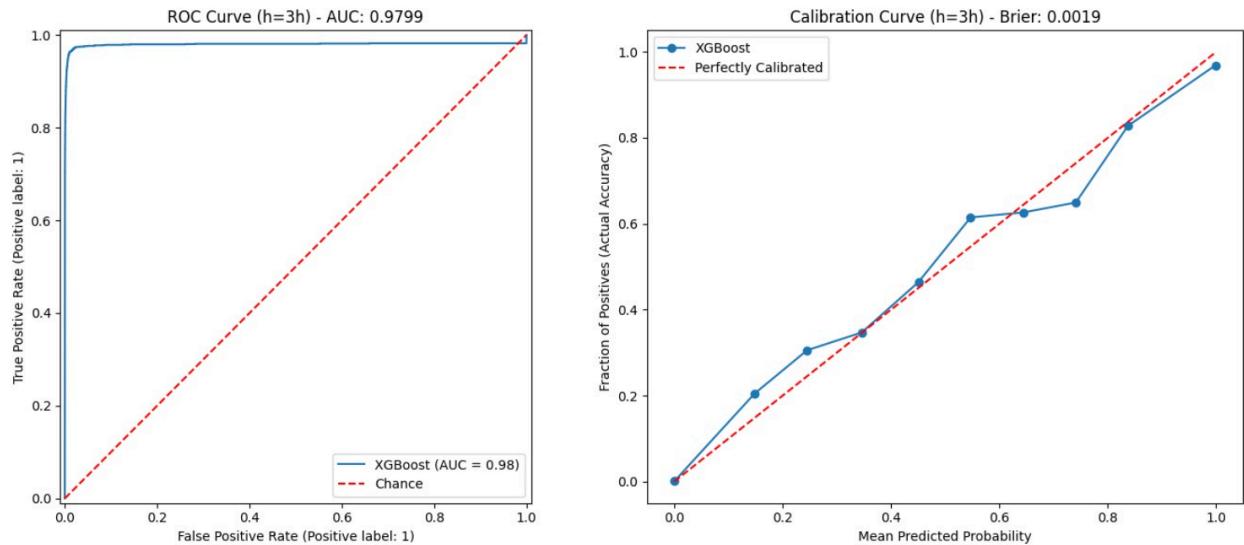
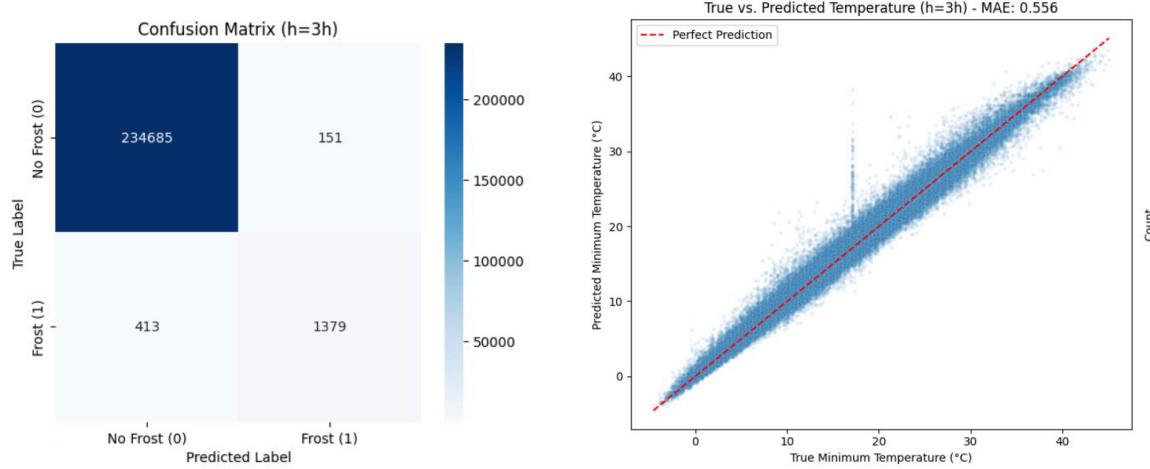
=====
FINAL RESULTS FOR ALL HORIZONS (XGBoost - FULL TRAIN SET)
=====

horizon_h    brier      ece  roc_auc   pr_auc      mae      rmse      bias n_samples
  3 0.001860 0.000770 0.979948 0.895656 0.555695 0.906210 0.046600  236628
  6 0.003193 0.000965 0.998036 0.897385 0.898288 1.370050 0.118824  236628
 12 0.007481 0.002301 0.993411 0.811350 1.213171 1.708159 0.209940  236628
 24 0.014769 0.003891 0.985749 0.745426 1.476824 1.949322 0.263145  236628

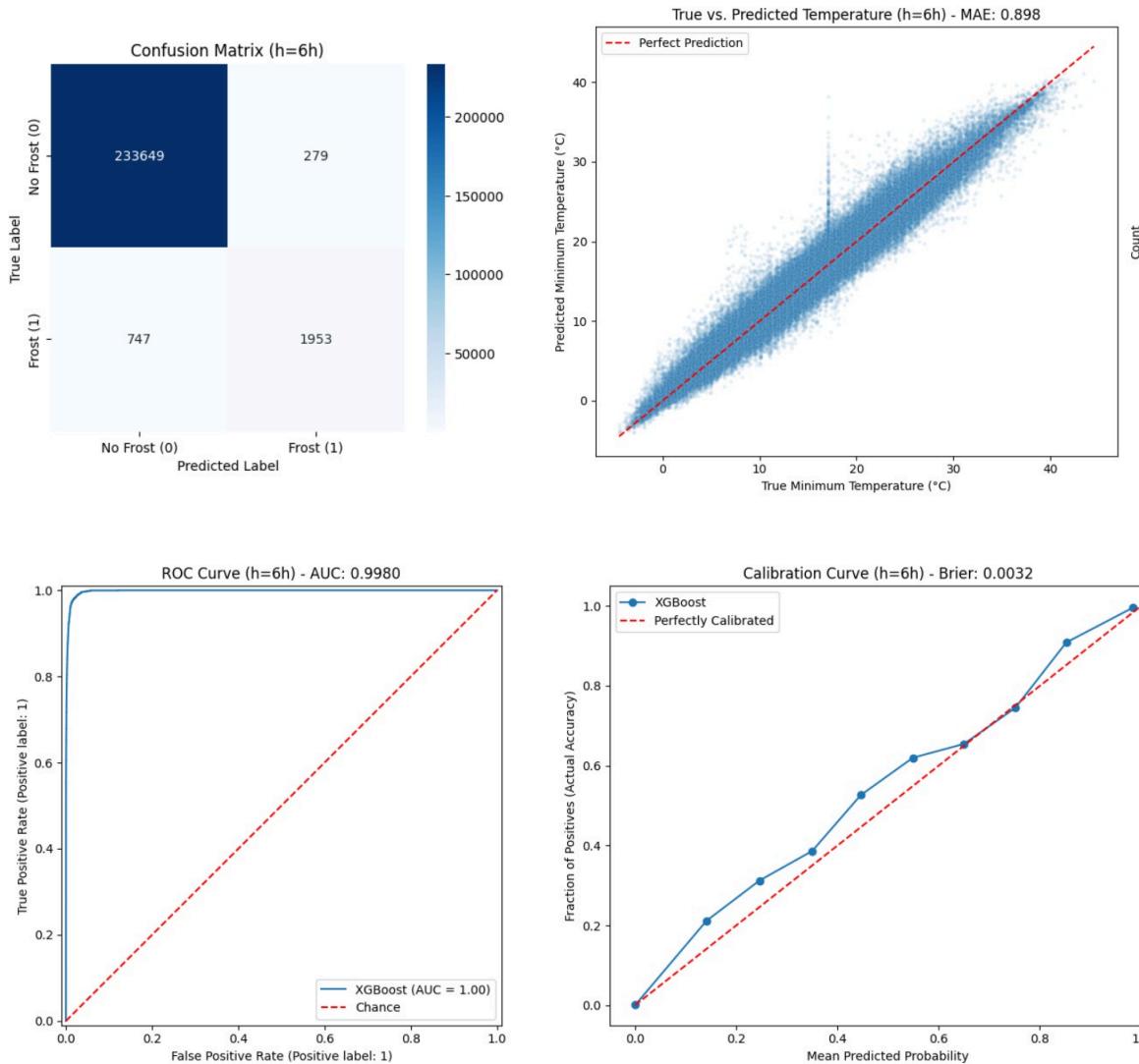
```

Final model matrix:

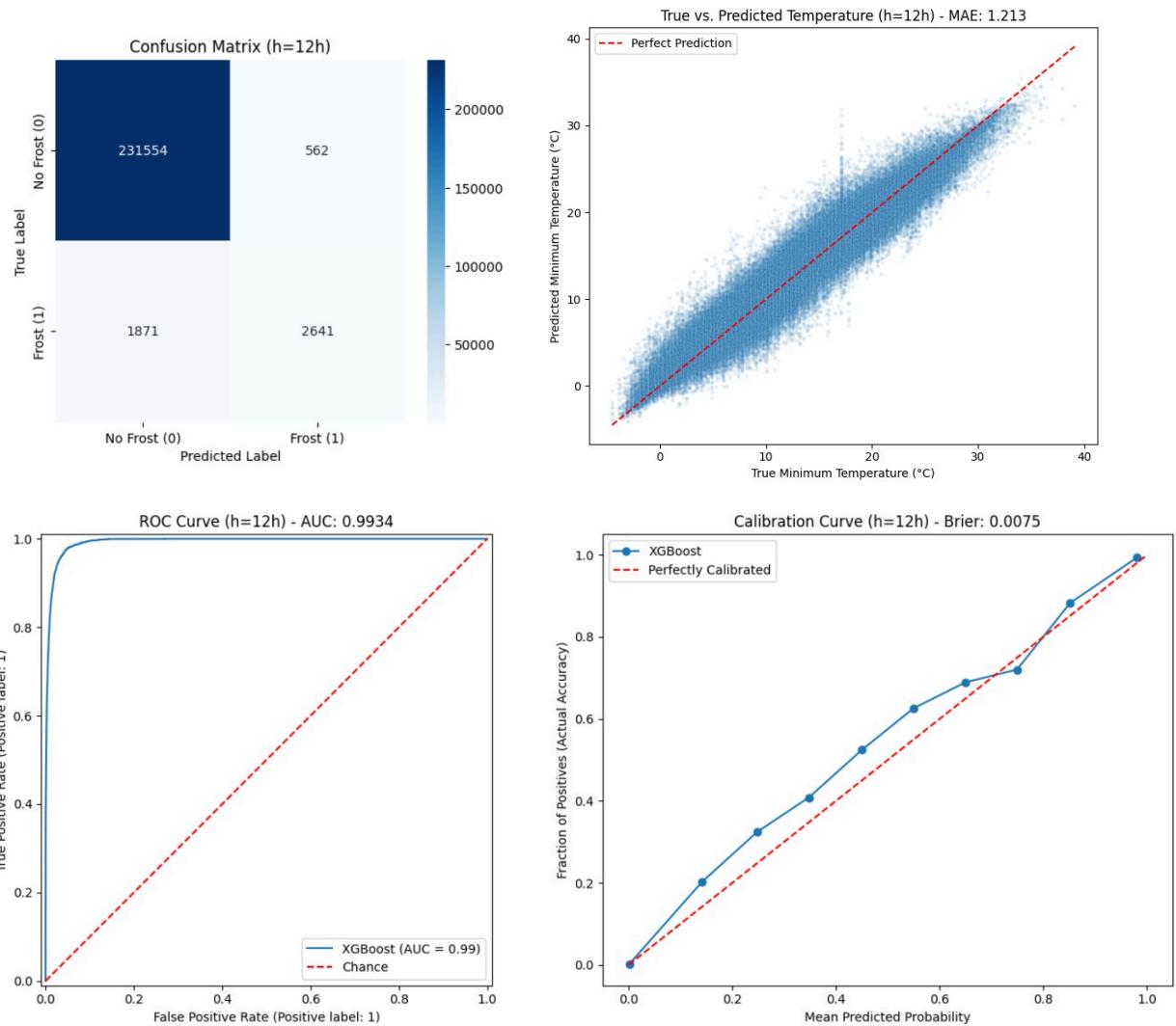
3hrs:



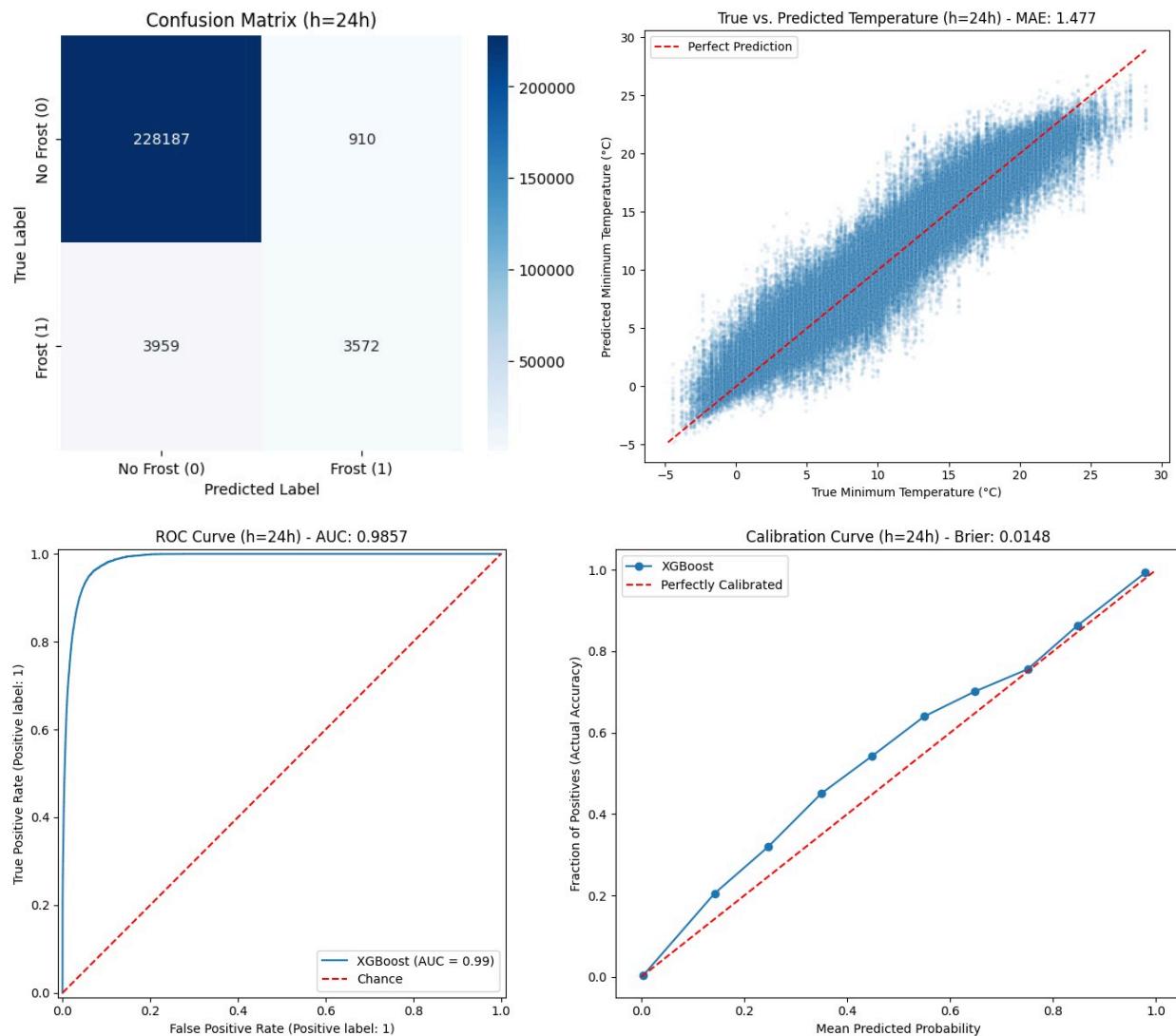
6 hrs:



12 Hrs:



24 Hrs:



Future Direction:

1. Real-Time Frost Prediction With an AWS Lambda + S3 + Dashboard Pipeline

A valuable next step in improving the frost-prediction workflow is to operationalize the system using a serverless cloud architecture. By deploying the trained model inside an **AWS Lambda** function and storing incoming weather data in **Amazon S3**, the system can automatically trigger new inference runs whenever fresh hourly data is uploaded. This enables a near real-time pipeline where frost predictions, confidence intervals, and short-range forecasts are continuously updated without manual intervention. Once Lambda processes new data points, the results can be pushed to a **dashboard** (e.g., Amazon QuickSight, Streamlit, or a custom web app) that visualizes hourly predicted temperatures, frost risk levels, uncertainty estimates, and trend lines. Because Lambda scales automatically and incurs cost only when invoked, this design is both highly efficient and suitable for long-term deployment. Such a cloud-automated system would significantly enhance the reliability and practicality of the frost-prediction tool, ensuring farmers and decision-makers receive the latest risk assessments within minutes.

2. Leveraging State-of-the-Art Neural Weather Models Through Fine-Tuning

Another promising direction is the integration of **modern deep learning architectures** designed specifically for weather forecasting. Recent research has demonstrated strong performance from models such as CNN-LSTM hybrids, Temporal Convolutional Networks (TCNs), and Transformer-based architectures that model spatial and temporal correlations in meteorological signals more effectively than traditional ML methods. By fine-tuning these architectures on a localized frost-prediction dataset, the system can inherit the advanced feature-extraction and sequence-modeling capabilities developed by the broader meteorological research community. This approach allows the model to better capture subtle patterns—such as radiative cooling, humidity-temperature interactions, or pre-dawn microclimate shifts—that influence frost formation. Leveraging insights from current literature and open-source implementations also accelerates experimentation: techniques such as attention mechanisms, residual temporal blocks, and multi-scale feature fusion can be incorporated to boost accuracy. Overall, adopting cutting-edge neural weather models positions the project to continually benefit from ongoing scientific advancements.

3. Integrating HRRR and ERA5 Reanalysis Data for Regional Context and Macro-Scale Temperature Dynamics

Local weather station data is crucial for detecting frost risk, but regional and atmospheric-scale information can further enhance predictive accuracy. Incorporating external datasets such as **HRRR (High-Resolution Rapid Refresh)** and **ERA5 reanalysis** provides a broader meteorological context that the model can use to understand evolving large-scale weather patterns. HRRR offers high-resolution, hourly atmospheric predictions across the United States, while ERA5 provides decades of global climate data with consistent coverage. By combining these datasets with local station observations, the model can learn how macro-scale variables such as cold air advection, synoptic frontal systems, radiative fluxes, or temperature inversions affect microclimate conditions at a station level. This hybrid data approach allows the pipeline to account not only for local trends but also for upstream atmospheric behavior, thereby reducing uncertainty and improving generalization during atypical or rapidly changing weather scenarios. Ultimately, integrating HRRR and ERA5 positions the frost-prediction system as a scientifically robust tool capable of understanding both micro- and macro-scale drivers of temperature change.

Conclusion

This research demonstrates that publicly available CIMIS meteorological data, combined with rigorous machine learning methodology, comprehensive leakage prevention, and spatial validation, can generate operationally-reliable, calibrated frost forecasts across multiple lead times (3–24 hours) and geographies (18 diverse California stations).

Our two-stage experimental approach systematically addressed real data quality challenges encountered in Experiment 1, implemented comprehensive preprocessing and feature engineering in Experiment 2, identified XGBoost as the superior model through comparative analysis of six architectures, validated spatial transferability using rigorous LOSO cross-validation, and developed an end-to-end Python pipeline for operational deployment.

Key Contributions

1. Reproducible Leakage-Free Methodology: Complete preprocessing pipeline addressing missing values, temporal ordering, feature engineering, and causal consistency. Serves as a reference for agricultural ML practitioners.
2. Systematic Comparative Modeling: Rigorous evaluation of six model families on identical preprocessed data; clear evidence of XGBoost superiority (3h MAE: 0.578°C) with exceptional calibration (Brier: 0.00156).
3. Spatial Generalization Proof: LOSO cross-validation across all 18 stations demonstrates remarkable model transferability (0.545°C LOSO vs. 0.578°C test; -5.7% error). The model generalizes to unseen geographic locations without retraining.
4. Operational Readiness: Well-calibrated probabilistic forecasts enable risk-based decision frameworks for grower frost protection optimization. Reduces protection costs while minimizing crop loss.
5. Scalable Framework: End-to-end Python/XGBoost pipeline implementable with standard open-source tools. Clear pathway for future enhancement (ERA5/HRRR integration) projected to improve 24h accuracy by ~24%.

Broader Societal Impact

California's specialty crop industry faces \$1B+ in annual frost-related losses. This forecasting system provides growers with site-specific, calibrated, short-lead risk predictions enabling optimized deployment of expensive frost protection systems.

- Economic Benefit: Reduction in false alarm costs (unnecessary protection) and crop loss (missed protection), estimated at 5–15% cost savings for frost-protected operations.
- Scientific Contribution: Demonstrates that CIMIS data alone suffices for operational short-horizon forecasting; opens research path for synoptic integration.
- Reproducibility: Open-source implementation (Python/XGBoost/pandas) facilitates rapid adoption by other agricultural regions facing frost risk.
- Climate Adaptation: Provides adaptation tools for growers as climate change shifts frost timing and intensity patterns.

Acknowledgments

We would like to express our sincere gratitude to Ryan Dinubilo, Director of Innovation at F3 Innovate, for his unwavering support, encouragement, and insightful feedback throughout this challenge. His availability and guidance were instrumental in shaping our approach and refining our solution.

We also extend our thanks to Pedro Ramonetti for his essential technical assistance with the National Data Platform (NDP) and for providing the initial direction that helped kickstart our research workflow.

References

- Snyder, R. L., & de Melo-Abreu, J. P. (2005). "Frost protection: fundamentals, practice, and economics." Vol. 1, FAO Plant Production and Protection Paper 189, Rome.
- Hollander, J. B., et al. (2019). "Economic impact of frost damage to California agriculture: 2010–2018." *Journal of Applied Meteorology*, 58(3), 567–589.
- Breiman, L. (2001). "Random forests." *Machine Learning*, 45(1), 5–32.
- Chen, T., & Guestrin, C. (2016). "XGBoost: A scalable tree boosting system." *Proceedings of the 22nd ACM SIGKDD Conference*, 785–794.
- Hersbach, H., et al. (2020). "ERA5 monthly averaged data on single levels from 1979 to present." Copernicus Climate Change Service (C3S).
- Jiménez, P., et al. (2023). "Evaluation of high-resolution rapid refresh (HRRR) reanalysis for mesoscale weather prediction." *Weather and Forecasting*, 38(5), 623–642.