



*Peeking Inside the Black Box:*  
**Explainable AI Methods for a Precipitation-type Model**

*Belen Saavedra,  
Dhamma Kimpara, David John Gagne II, John Schreck, Charlie Becker, Gabrielle Gantos*

*NCAR Machine Integration and Learning of Earth Systems*



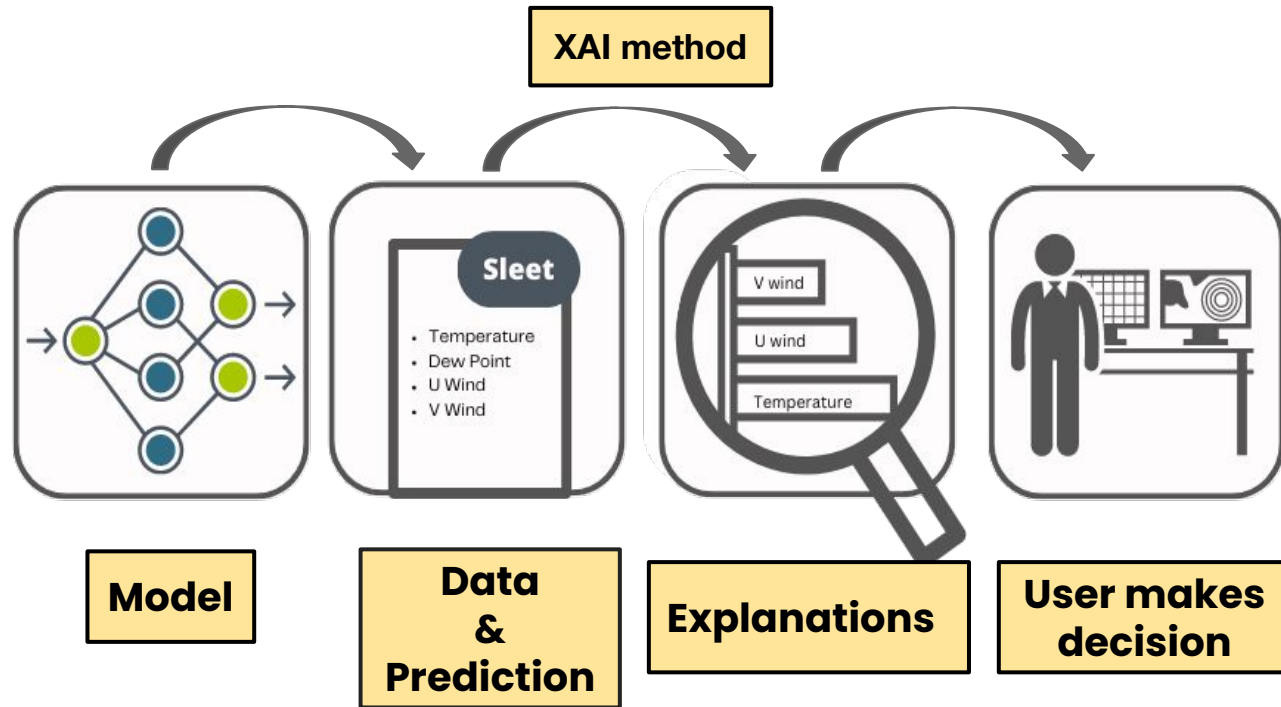
August 2, 2023



# Why do we need to understand ML models?



# XAI Pipeline

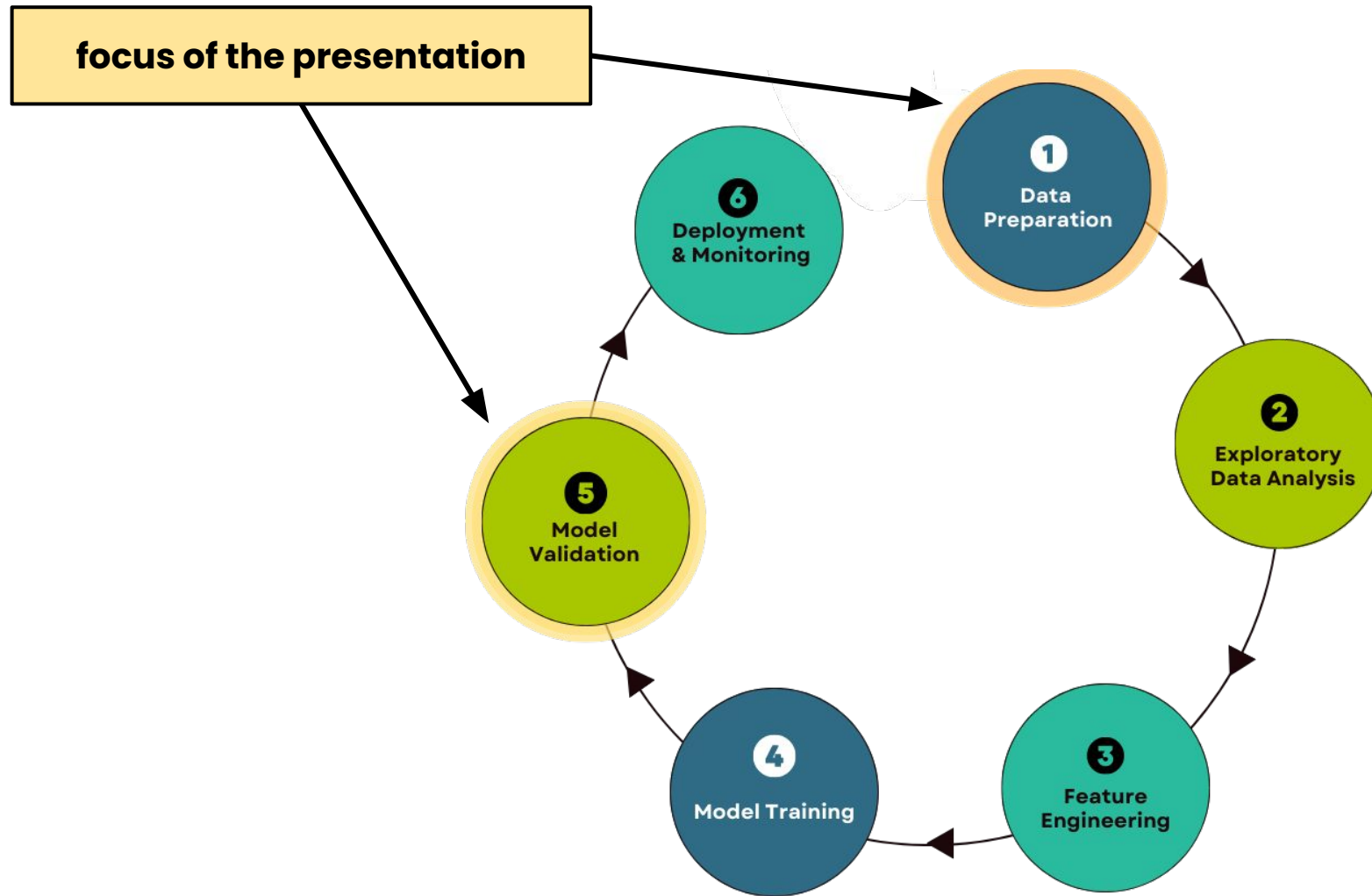


XAI methods aim to:

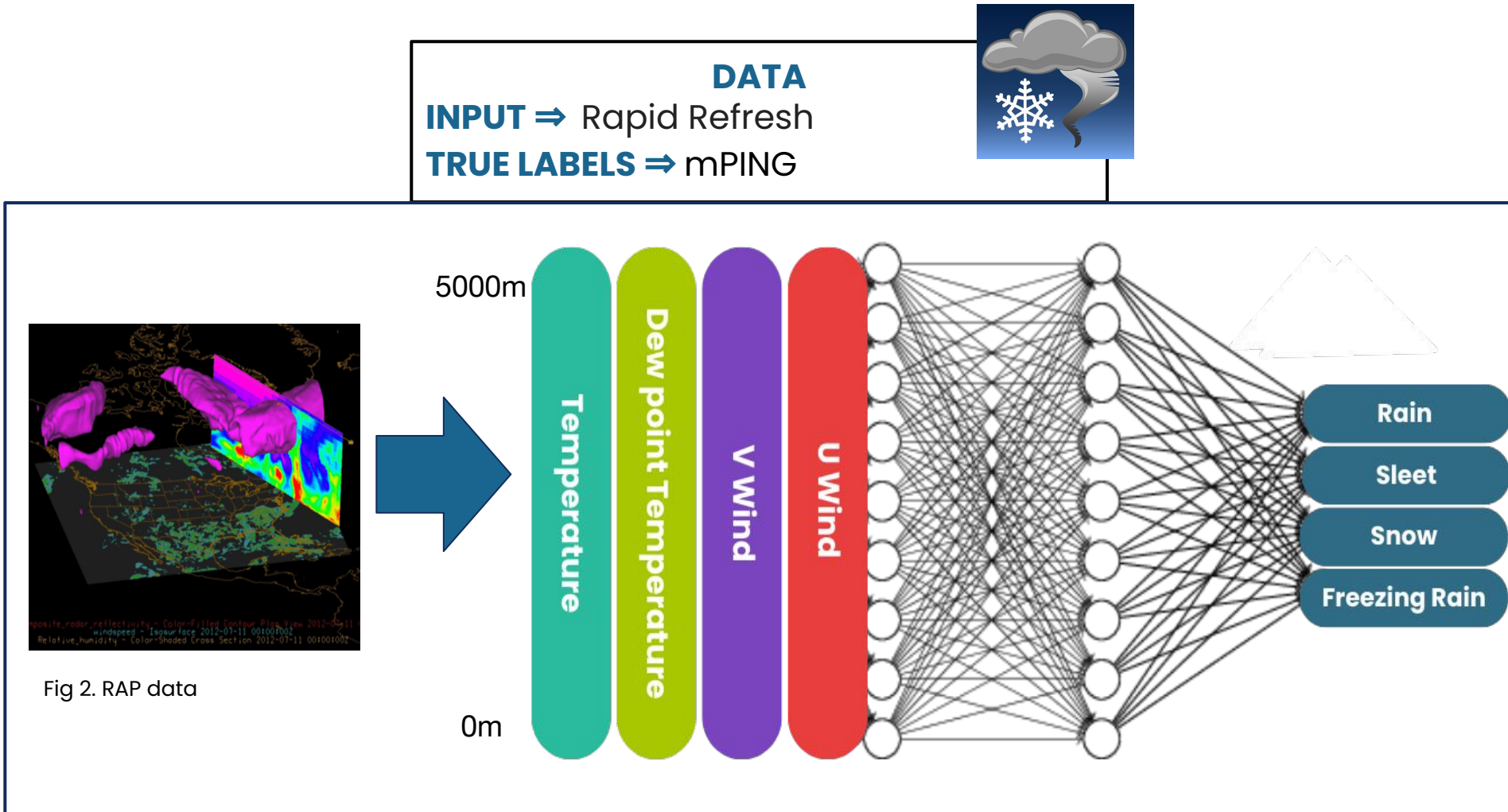
- Verify if the predictions of the ML models are consistent with the real-world
- Increase the credibility of machine learning models for both technical and non-technical users

Figure 1. XAI pipeline

# Machine Learning Pipeline



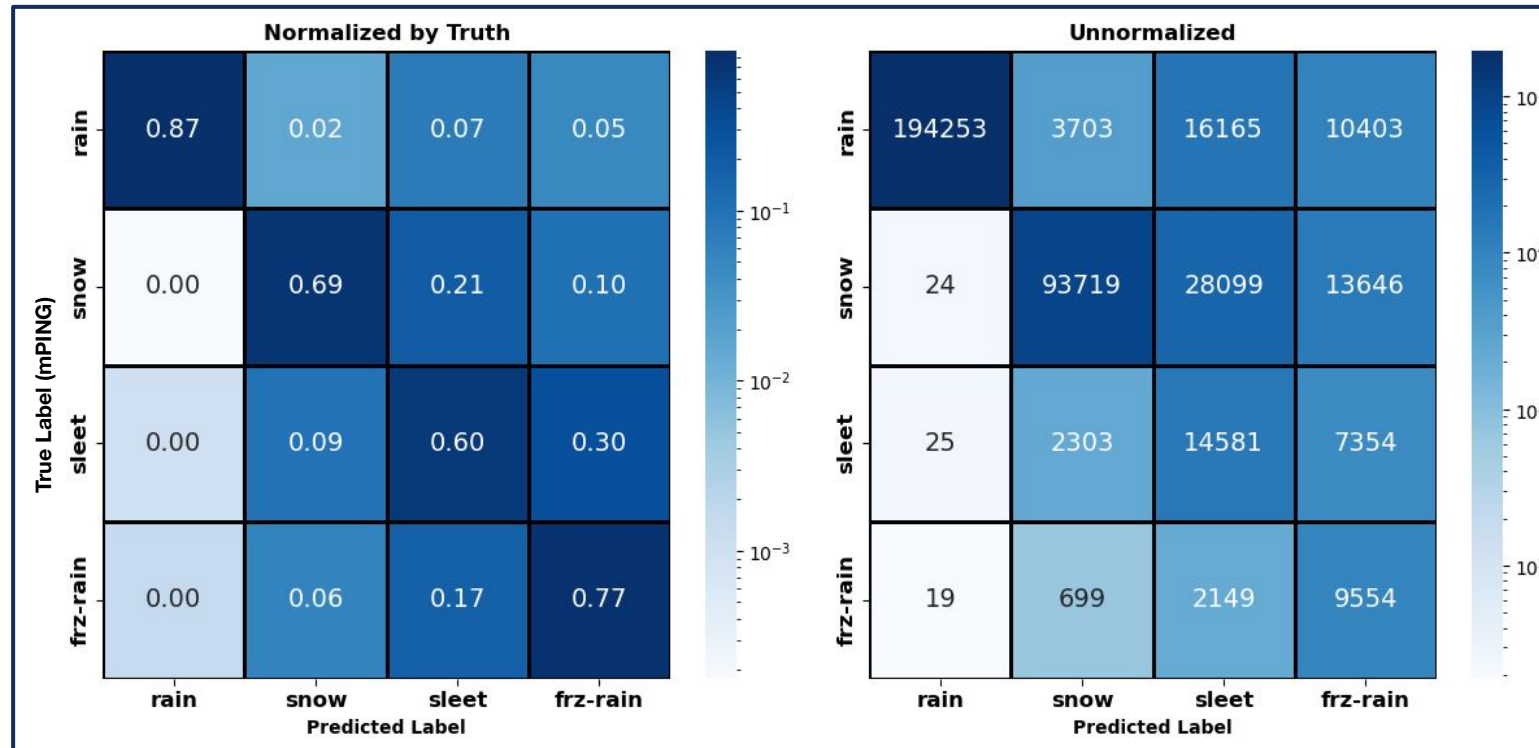
# Precipitation-type model



# Precipitation-type model

## PERFORMANCE:

- mPING vs ML
- Overprediction of rain
- Under prediction of sleet and freezing rain



# Post hoc XAI methods

Gradient * Input	Which <b>features are most influential</b> in predicting the model's output?
Shapley Additive Explanations (SHAP)	How much does <b>each feature contribute to the model's predictions</b> ?
Permutation Feature Importance	How does the <b>performance of the model change</b> when the information content of a feature is destroyed?

$$\mathbf{A}_{\text{Gradient} \odot \text{Input}}^c = \frac{\partial S_c(\mathbf{x})}{\partial \mathbf{x}} \odot \mathbf{x}.$$

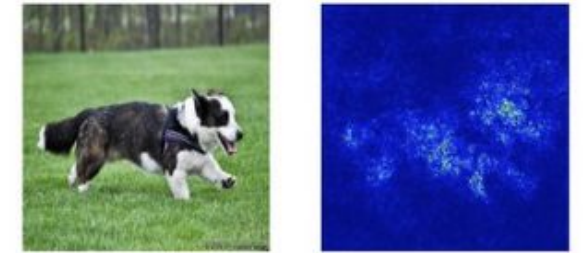
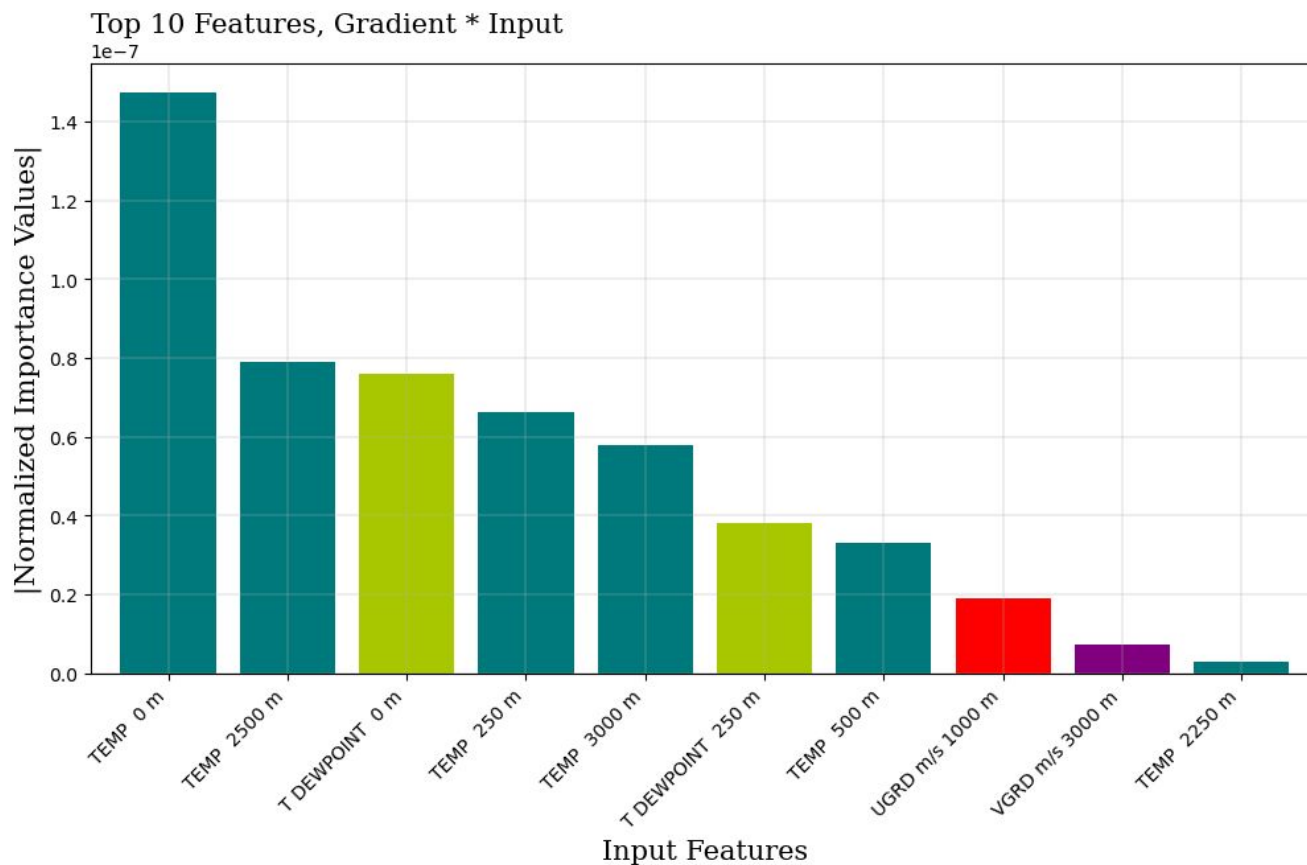


Fig. 3 Input \* Gradient attribution method

# Gradient \* Input

Which **features are most influential** in predicting the model's output?



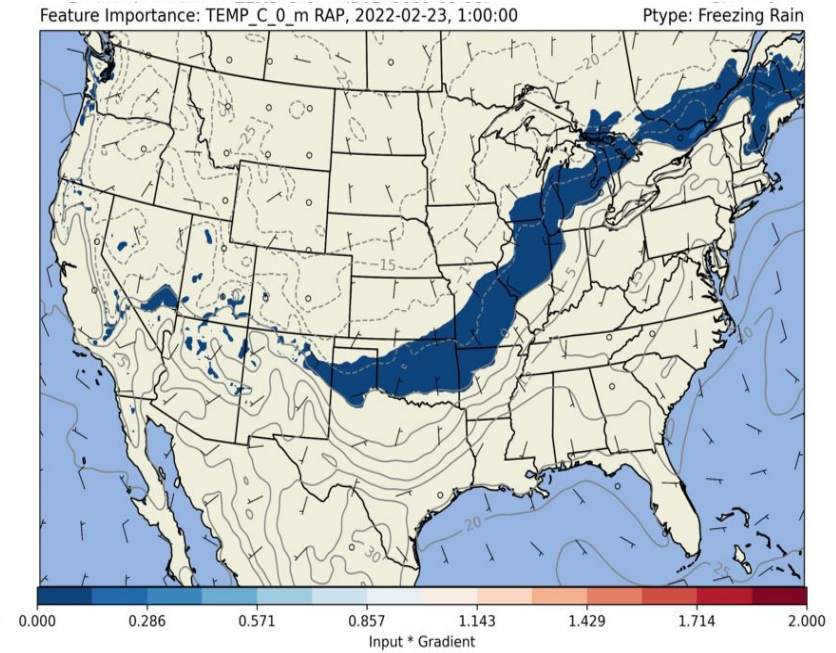
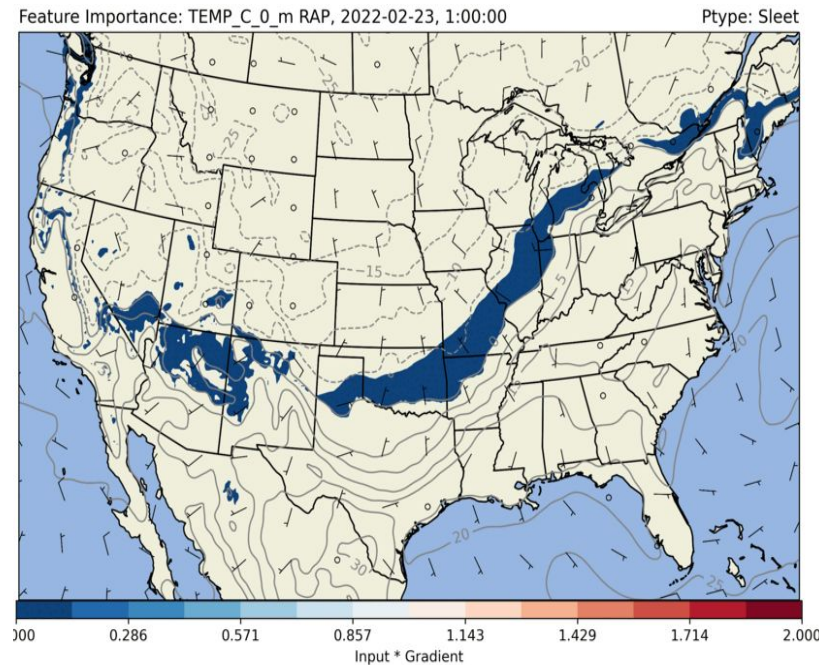
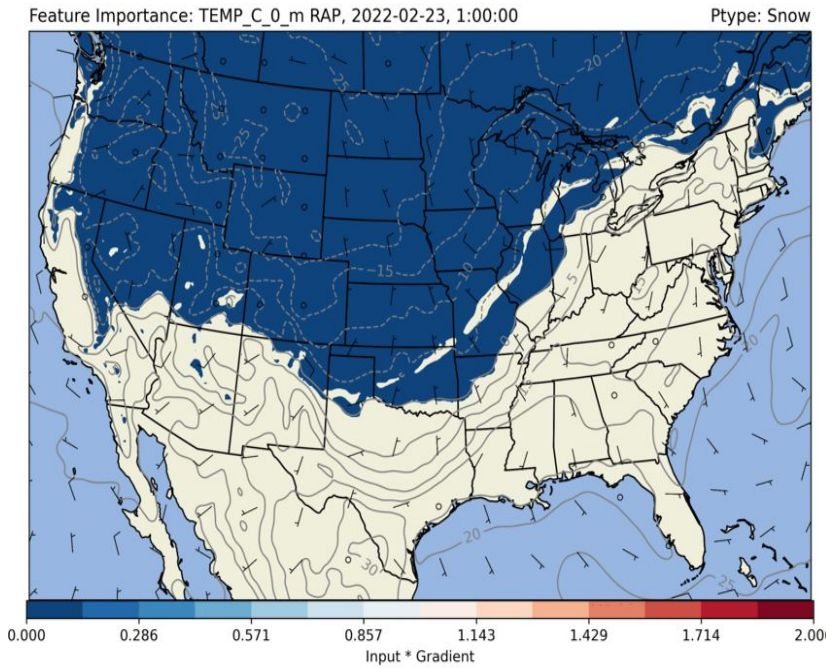
Gradient \* Input works by multiplying the gradient of the model's output with the input features.

$$\mathbf{A}_{\text{Gradient} \odot \text{Input}}^c = \frac{\partial S_c(\mathbf{x})}{\partial \mathbf{x}} \odot \mathbf{x}.$$



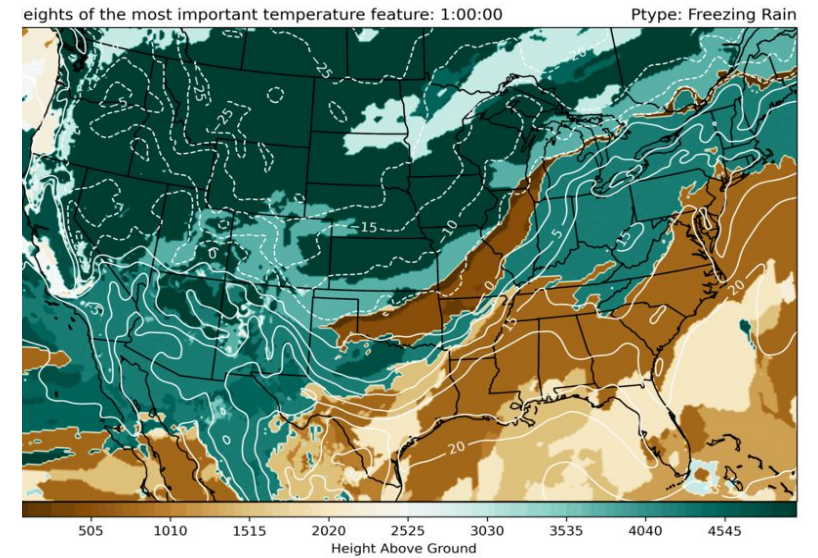
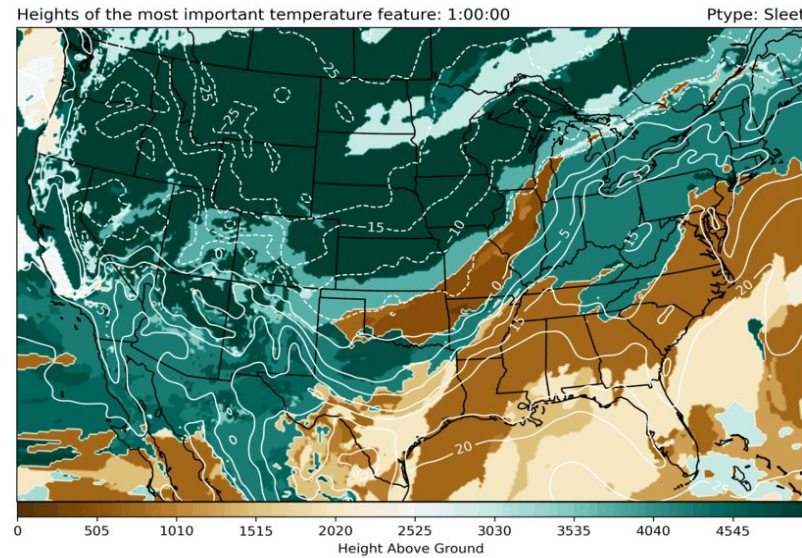
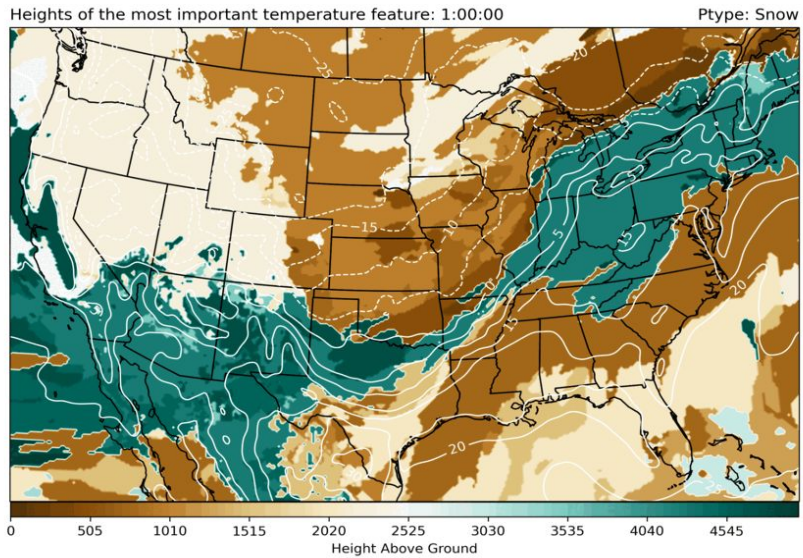
# Gradient \* Input: CONUS plots

Which **features** are most influential in predicting the model's output?



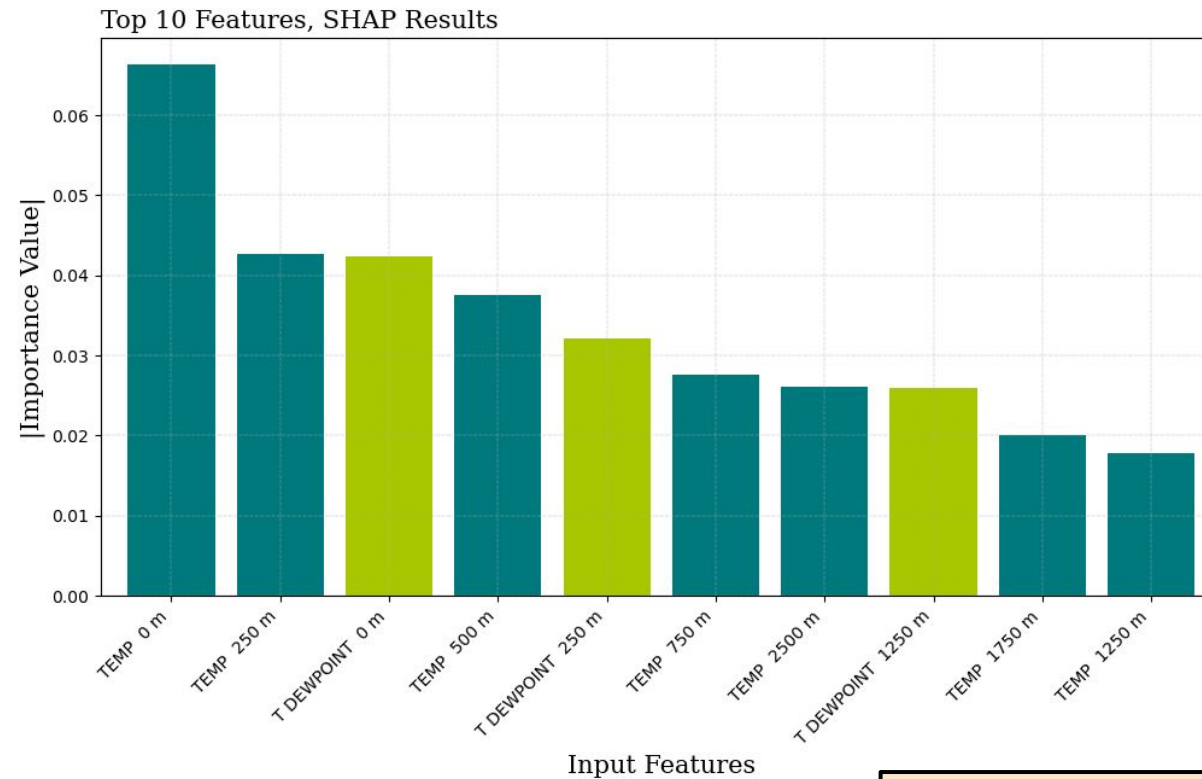
# Gradient \* Input: CONUS plots

Which **features are most influential** in predicting the model's output with respect to their height?



# Shapley Additive Explanations (SHAP)

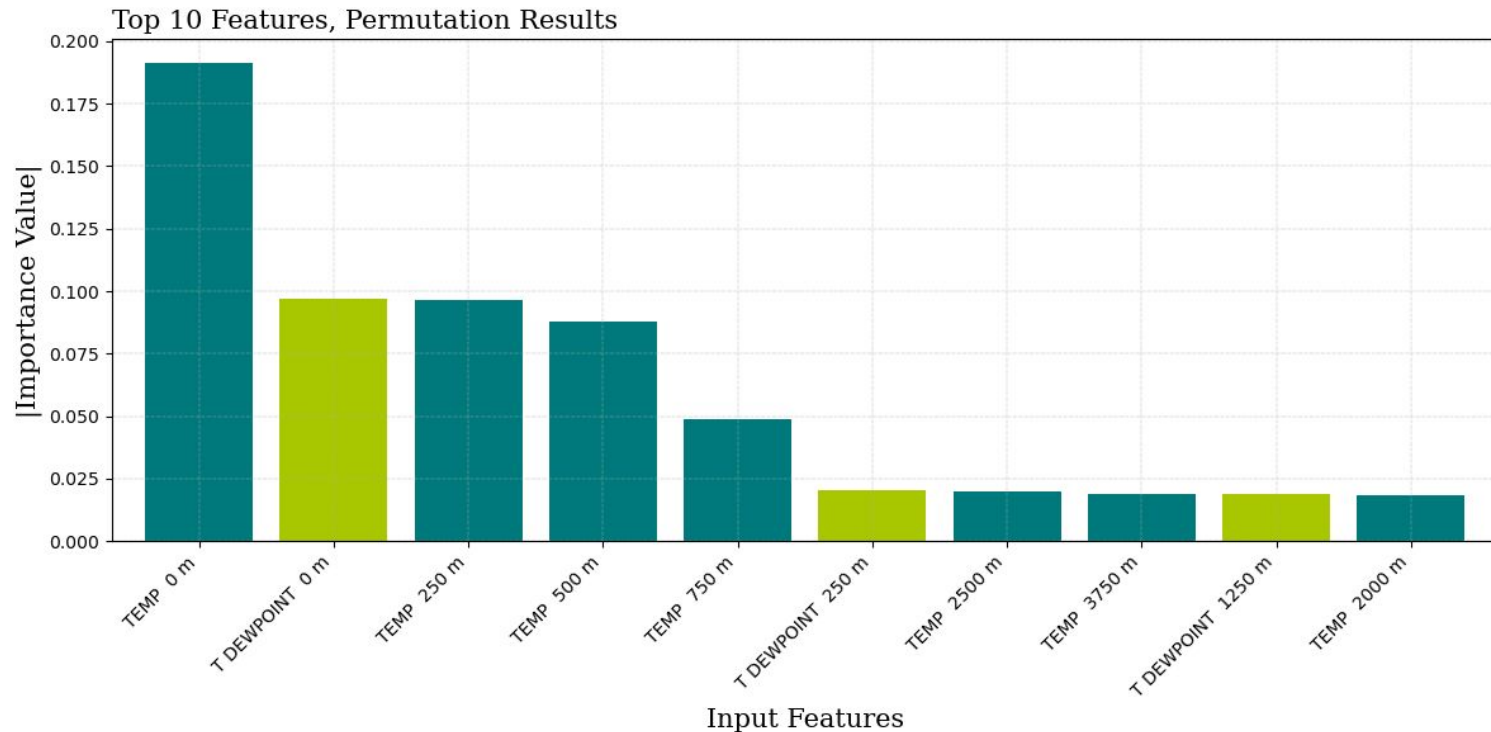
How much does **each feature contribute to the model's predictions?**



SHAP calculates the average contribution of each feature, representing how much each feature influences the model's prediction

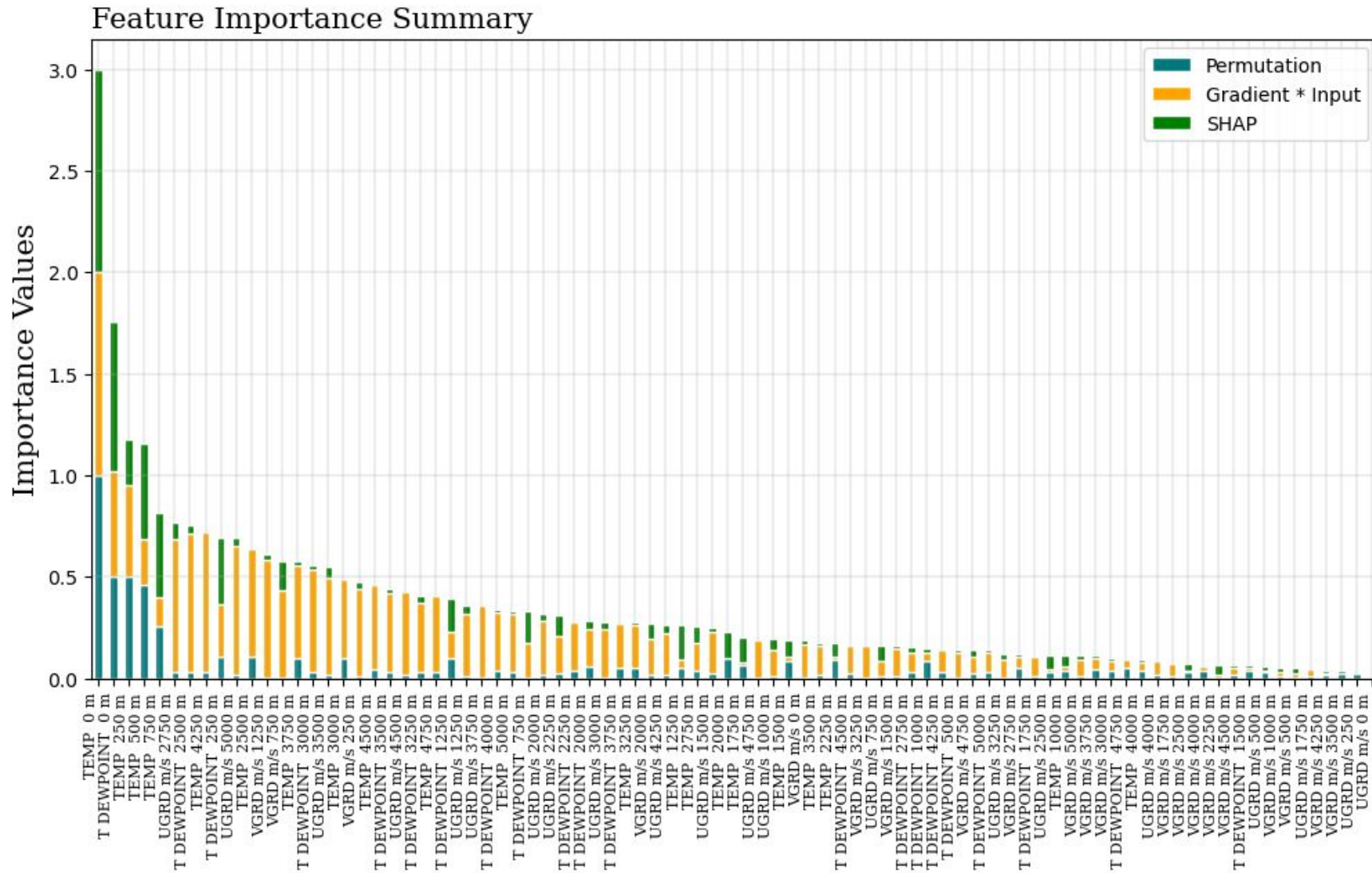
# Permutation Feature Importance

What is the importance of each feature in predicting the model's output when the feature values are randomly shuffled?



Permutation feature importance works by randomly shuffling the values of a single feature and measuring the resulting change in the model's performance. The feature with the largest change in performance is considered to be the most important feature.

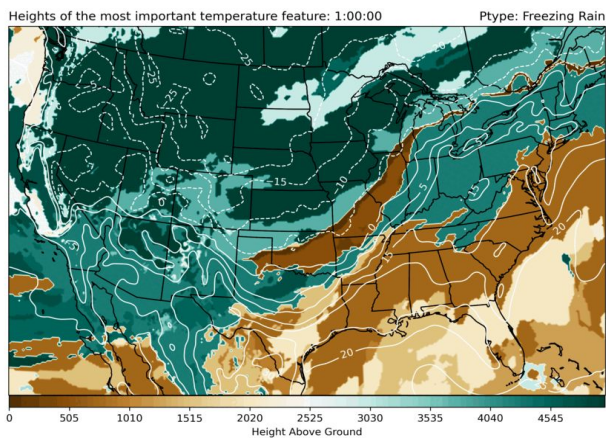
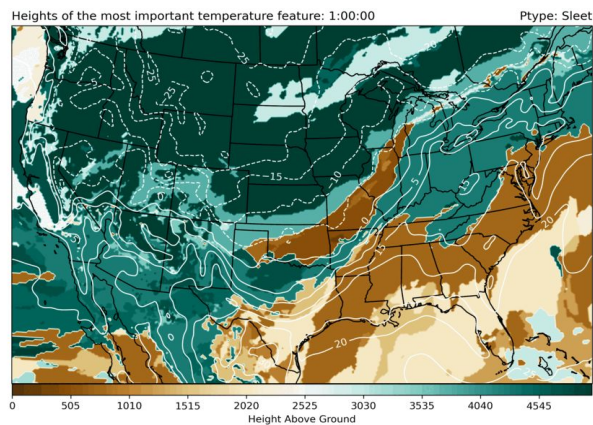
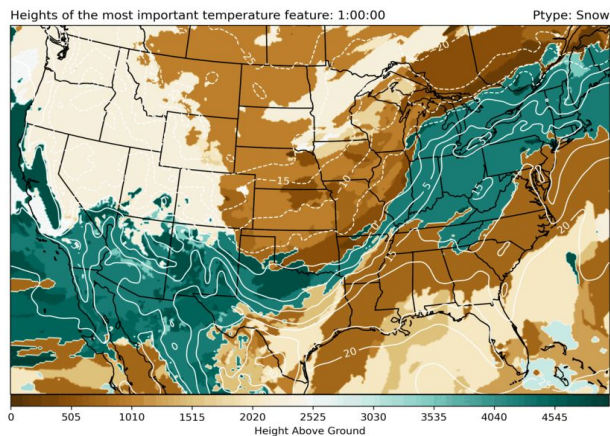
# XAI Results Summary



# Limitations of XAI methods

- There is not a XAI method that works for every explainability task  
*Some factors to consider:*
  - ◆ model type
  - ◆ scope of the explanation
  - ◆ audience - who needs to understand the model?
- They can be computationally expensive
- XAI methods often rely on simplification techniques may not capture the nuances of the decision-making process
- The results of XAI might be hard to interpret

# Main Takeaways and Future Steps



## XAI Methods:

- Temperature at 0m is the top feature for across the three XAI methods
- The Input features that are near the surface tend to be the most important
- Each XAI method provides slightly different results.

## Broader implications:

- Support the communication of the predictions of this model to a wide audience (decision makers, forecasters and general users)

# References

[1] McGovern, A., Lagerquist, R., Gagne, D. J., Jergensen, G. E., Elmore, K. L., Homeyer, C. R., & Smith, T. (2019). Making the black box more transparent: Understanding the physical implications of machine learning.

## Figures:

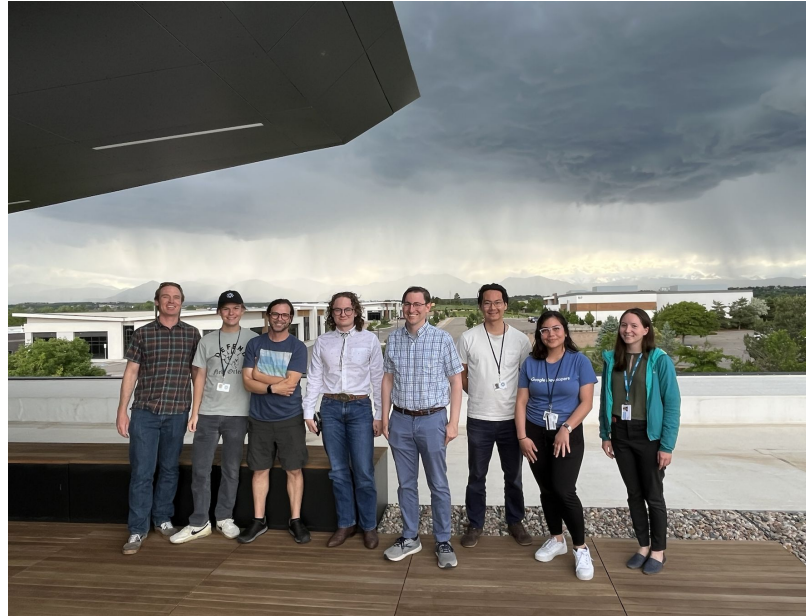
Fig 1. Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. arXiv preprint arXiv:1602.04938, 2016.

Fig 2. RAP, NOAA, Rapid Refresh/Rapid Update Cycle,  
<https://www.ncei.noaa.gov/products/weather-climate-models/rapid-refresh-update#:~:text=The%20National%20Centers%20for%20Environmental,for%20smaller%20regions%20of%20interest>.

Fig 3. <https://i.stack.imgur.com/Nxhrr.jpg>



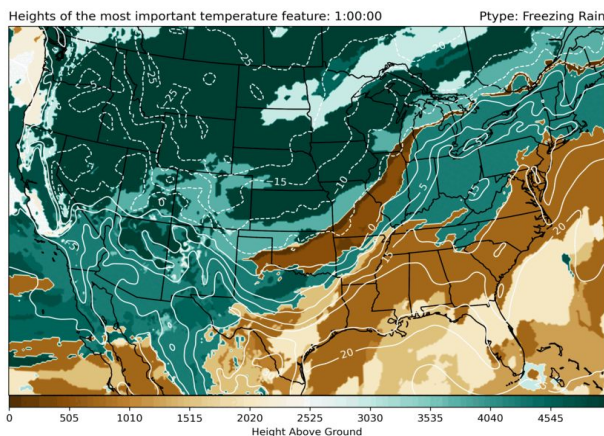
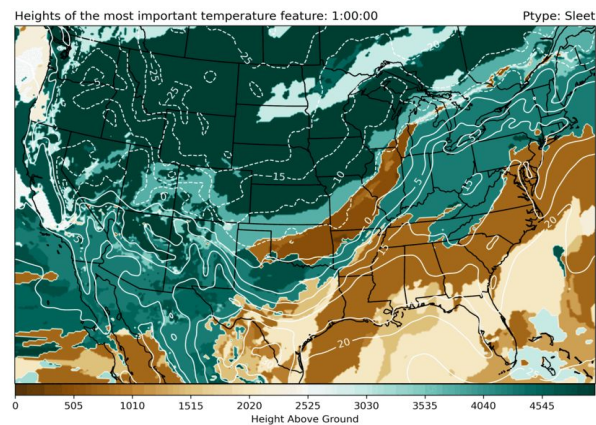
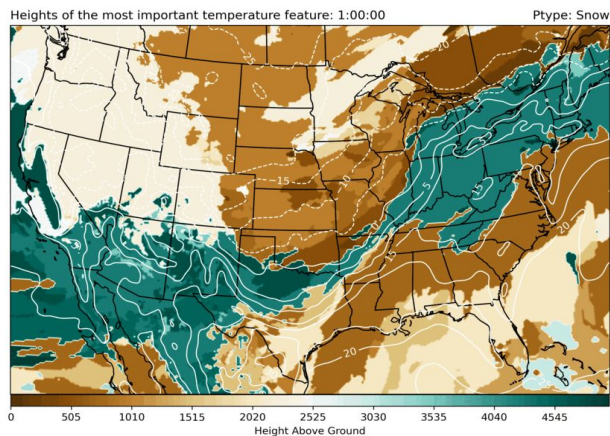
# Acknowledgements



*Part of the MILES group and Vaisala interns*

- **MILES group:**  
David John Gagne II, John Schreck, Charlie Becker, Gabrielle Gantos, Dhamma Kimpara, Hayden Outlaw
- **SIParCS:**  
Virginia Do, Julius Owusu Afriyie, and the intern cohort
- **NESSI:**  
Jerry Cyccone, Ben Fellman, and the NESSI 2023 cohort
- **Funding:**  
National Science Foundation under Grant No. ICER-2019758.

# Main Takeaways and Future Steps



## XAI Methods:

- Temperature at 0m is the top feature for each of the methods
- The Input features that are below 1000 meters above ground tend to be the most important
- Each XAI method provides slightly different results.

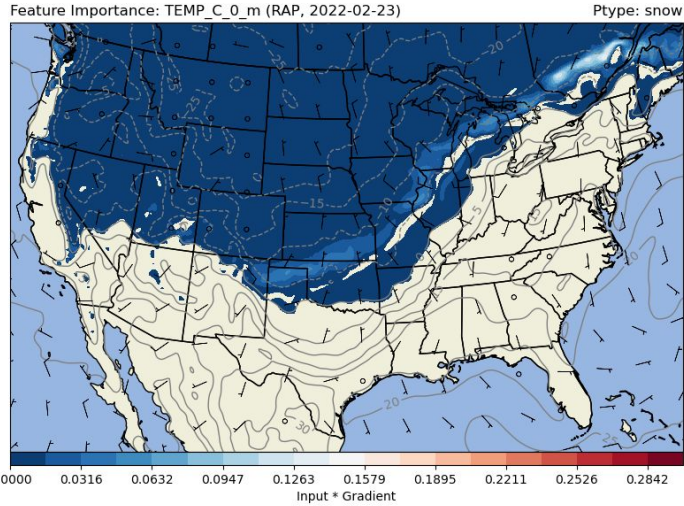
## Broader implications:

- Support the communication of the predictions of this model to a wide audience (decision makers, forecasters and general users)

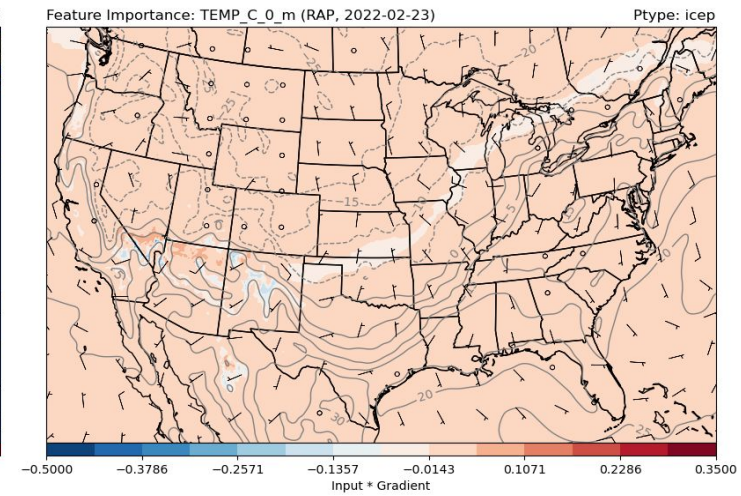
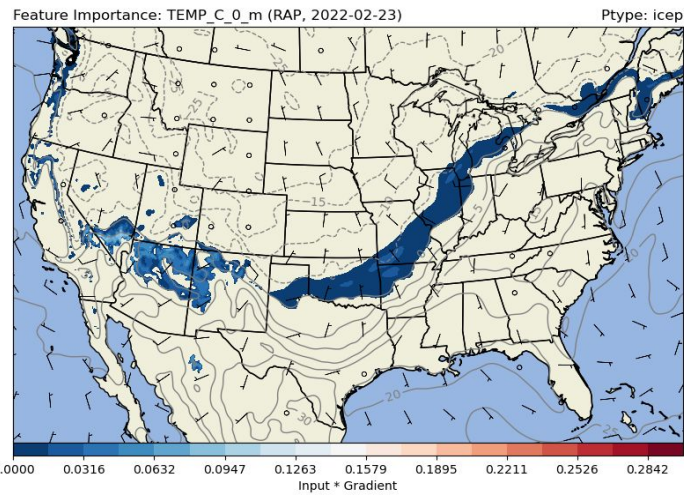
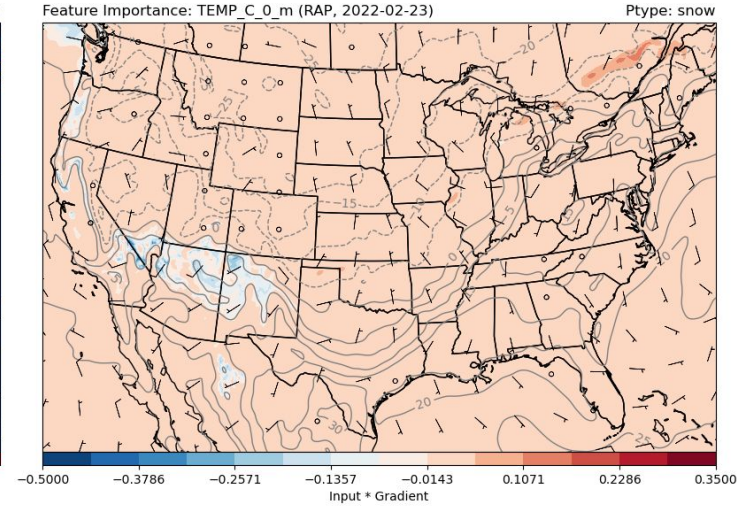
Questions/Feedback?

# Appendix: Input \* Gradient CONUS plots

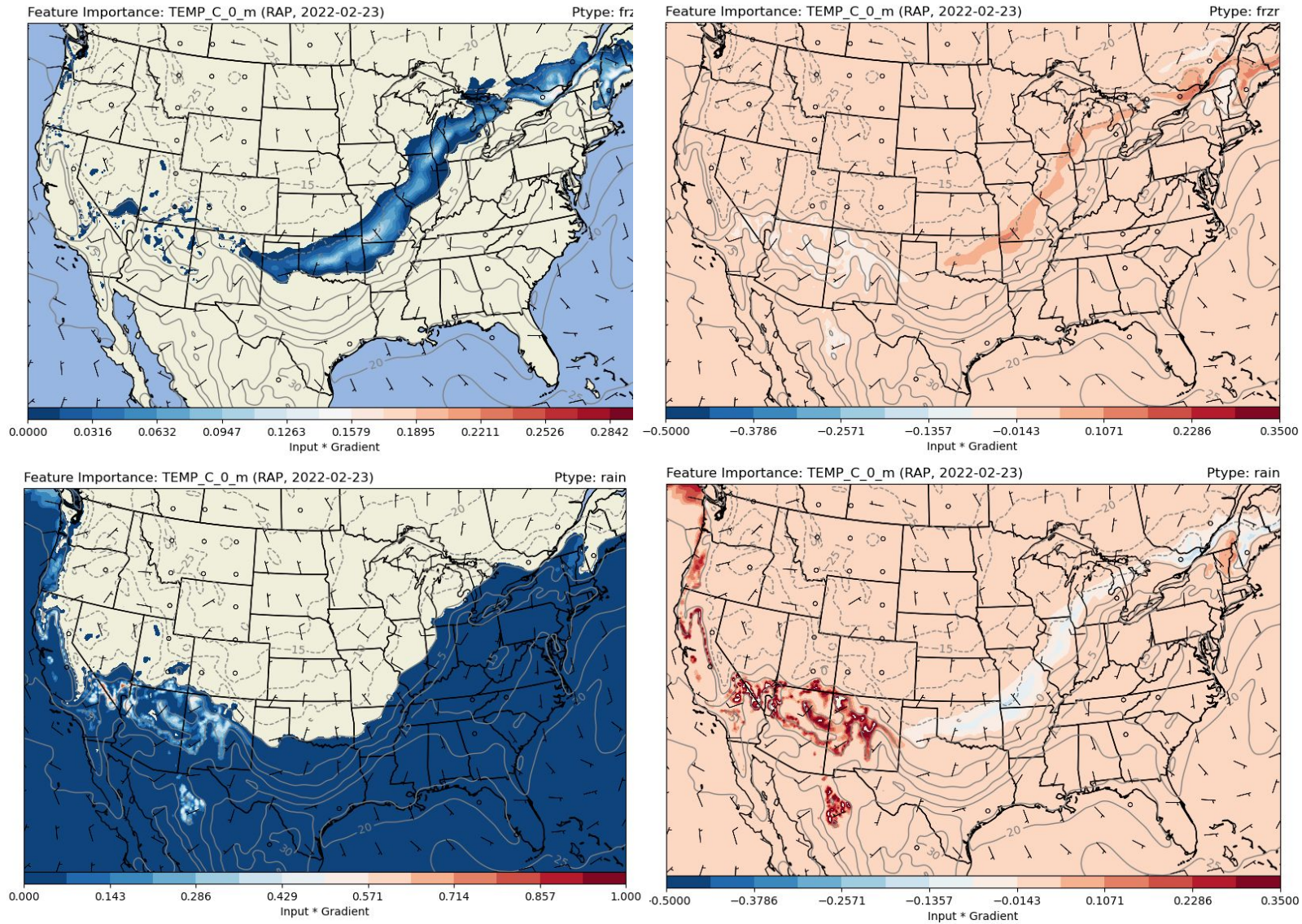
## Absolute Values



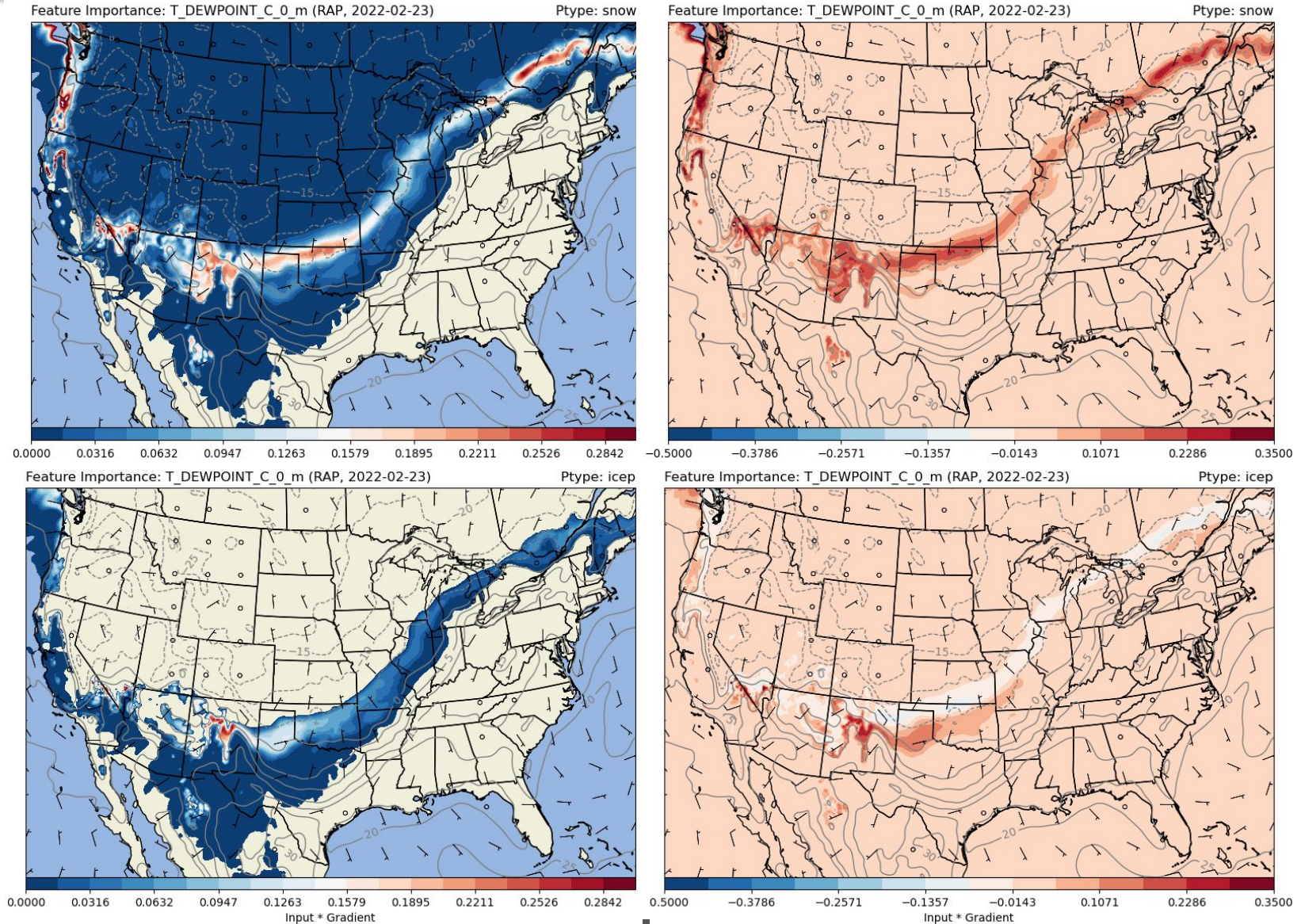
## Raw Values



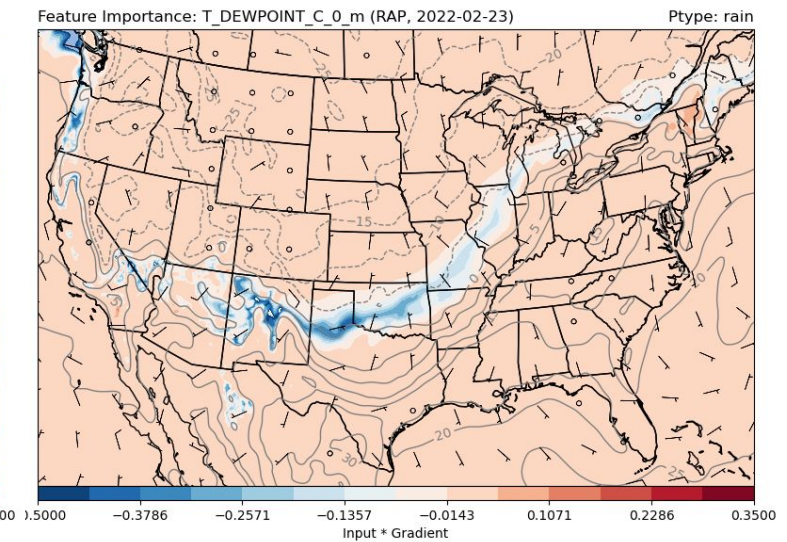
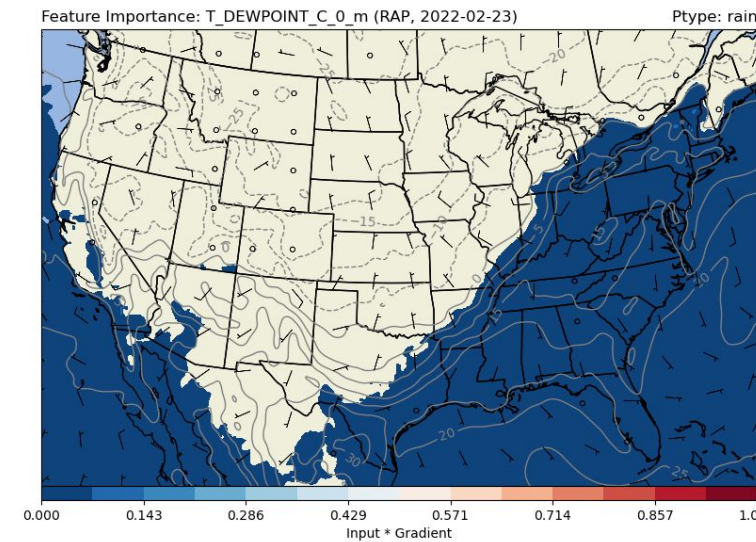
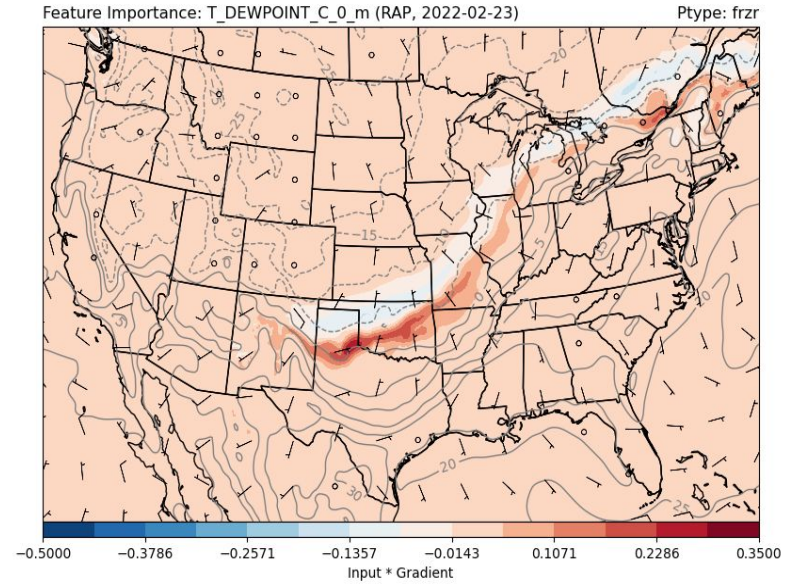
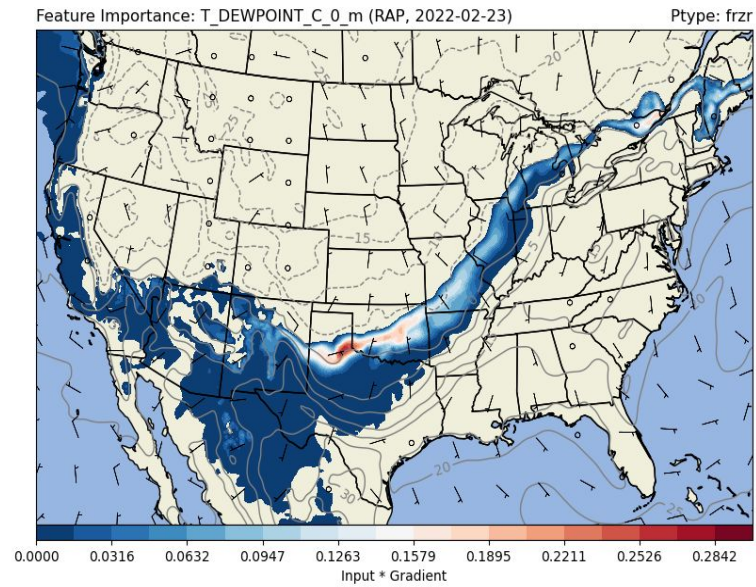
# Feature Importance by precipitation type: Temperature



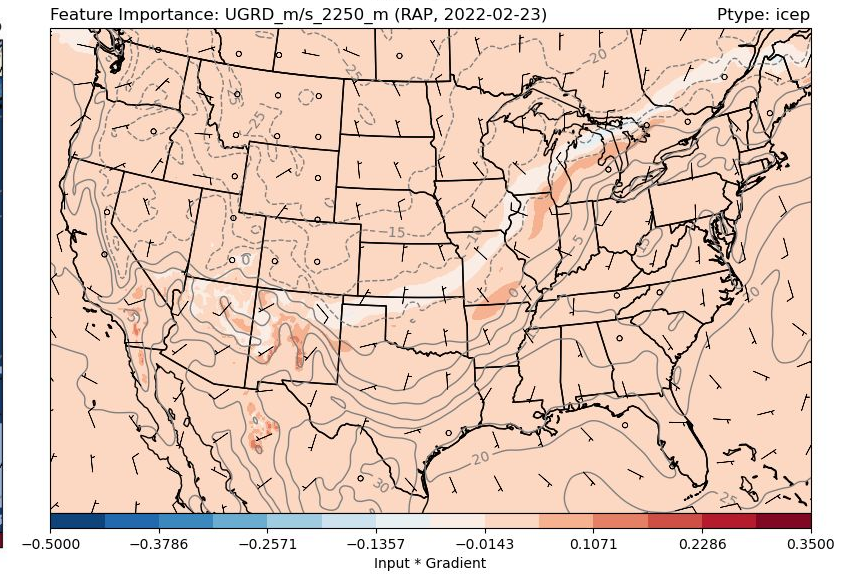
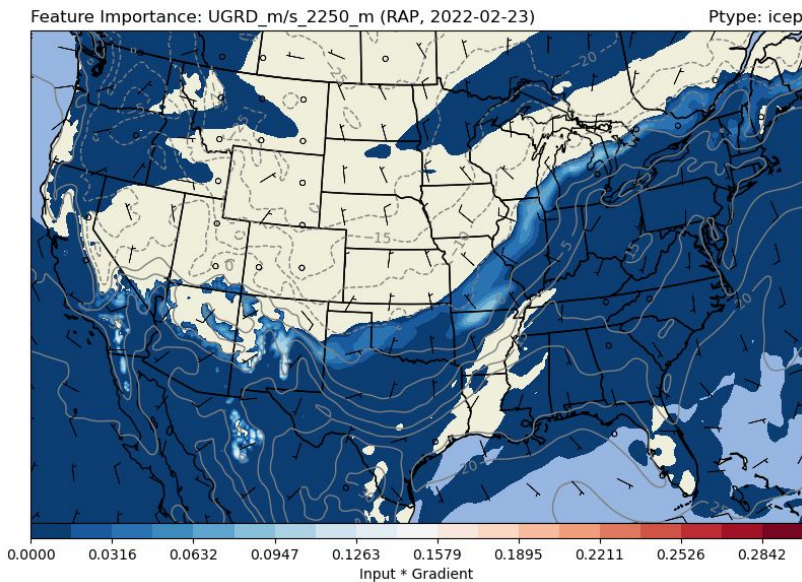
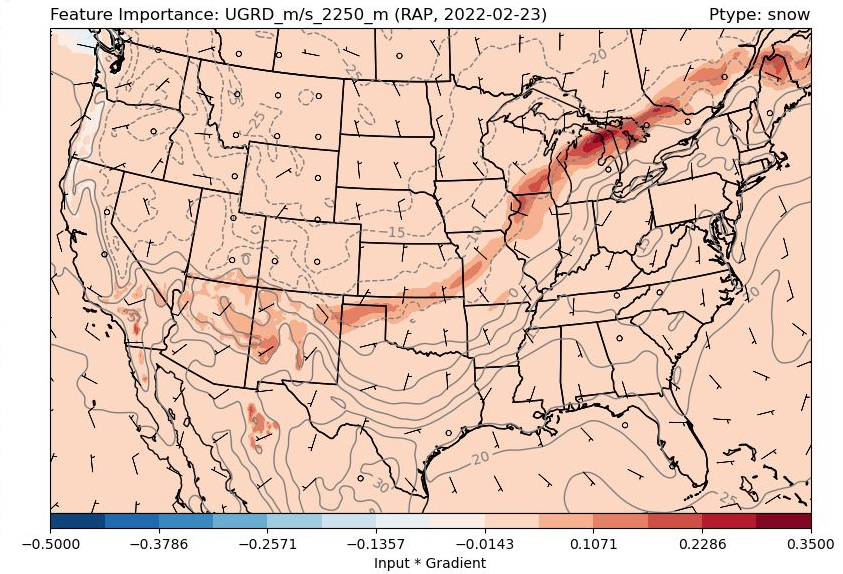
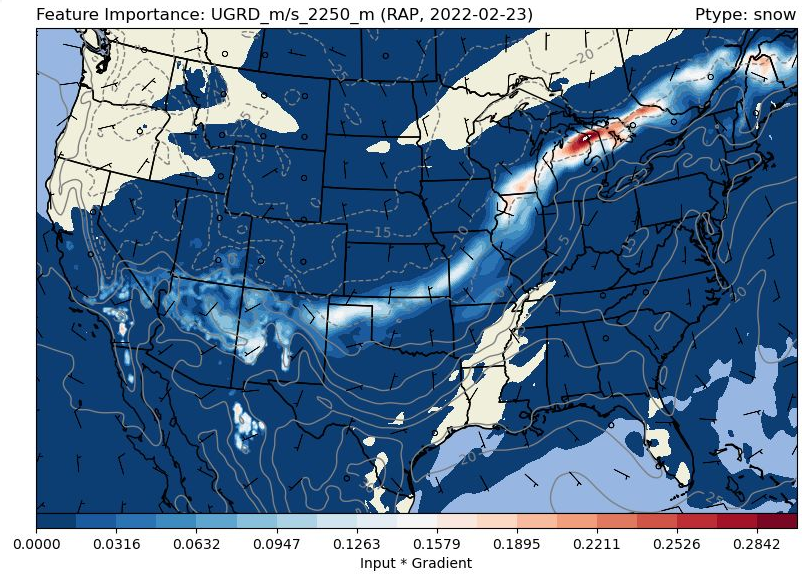
# Feature Importance by precipitation type: Dew Point



# Feature Importance by precipitation type: Dew Point

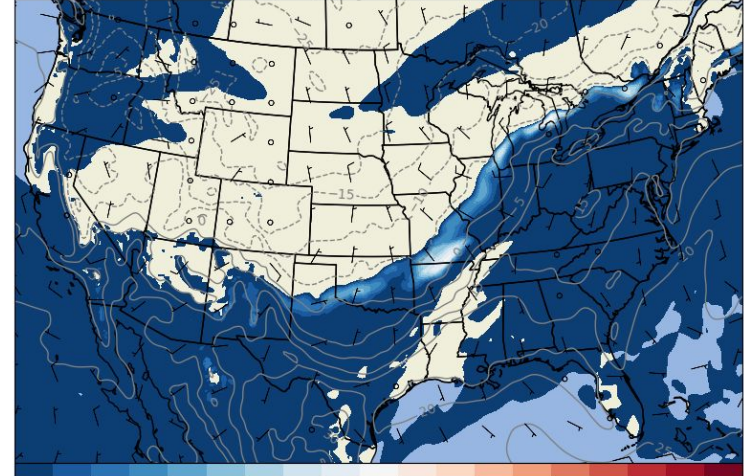


# Feature Importance by precipitation type: UGRD



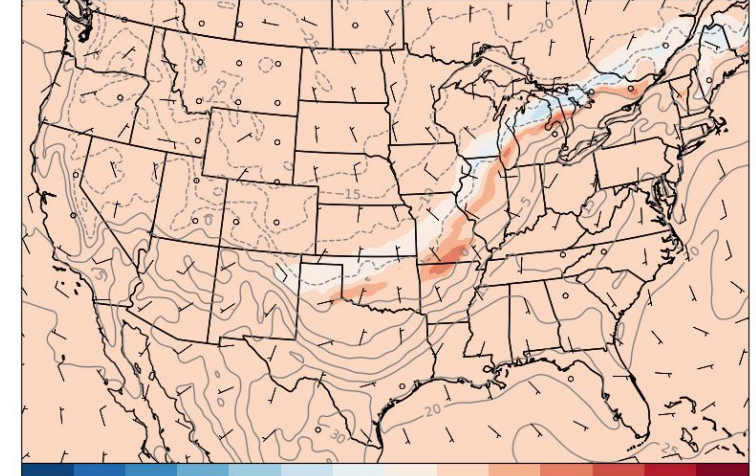
# Feature Importance by precipitation type: UGRD

Feature Importance: UGRD\_m/s\_2250\_m (RAP, 2022-02-23) Ptype: frzr



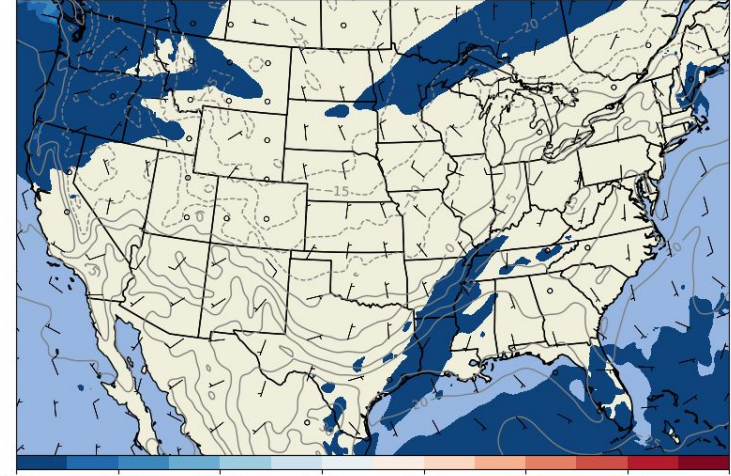
0.0000 0.0316 0.0632 0.0947 0.1263 0.1579 0.1895 0.2211 0.2526 0.2842  
Input \* Gradient

Feature Importance: UGRD\_m/s\_2250\_m (RAP, 2022-02-23) Ptype: frzr



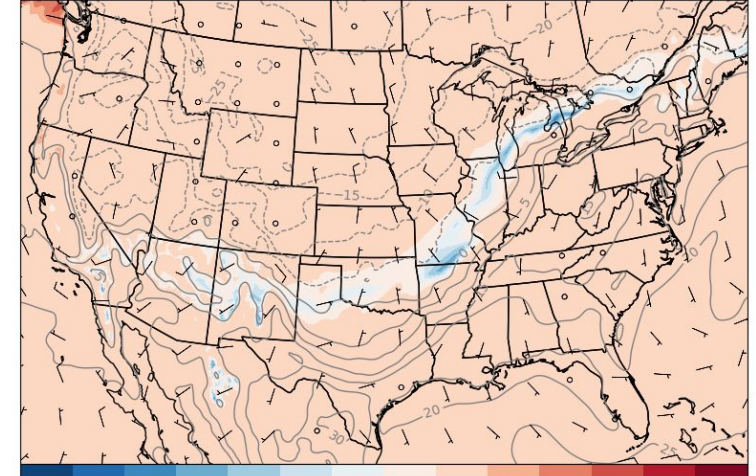
-0.5000 -0.3786 -0.2571 -0.1357 -0.0143 0.1071 0.2286 0.3500  
Input \* Gradient

Feature Importance: UGRD\_m/s\_2250\_m (RAP, 2022-02-23) Ptype: rain



0.000 0.143 0.286 0.429 0.571 0.714 0.857 1.000  
Input \* Gradient

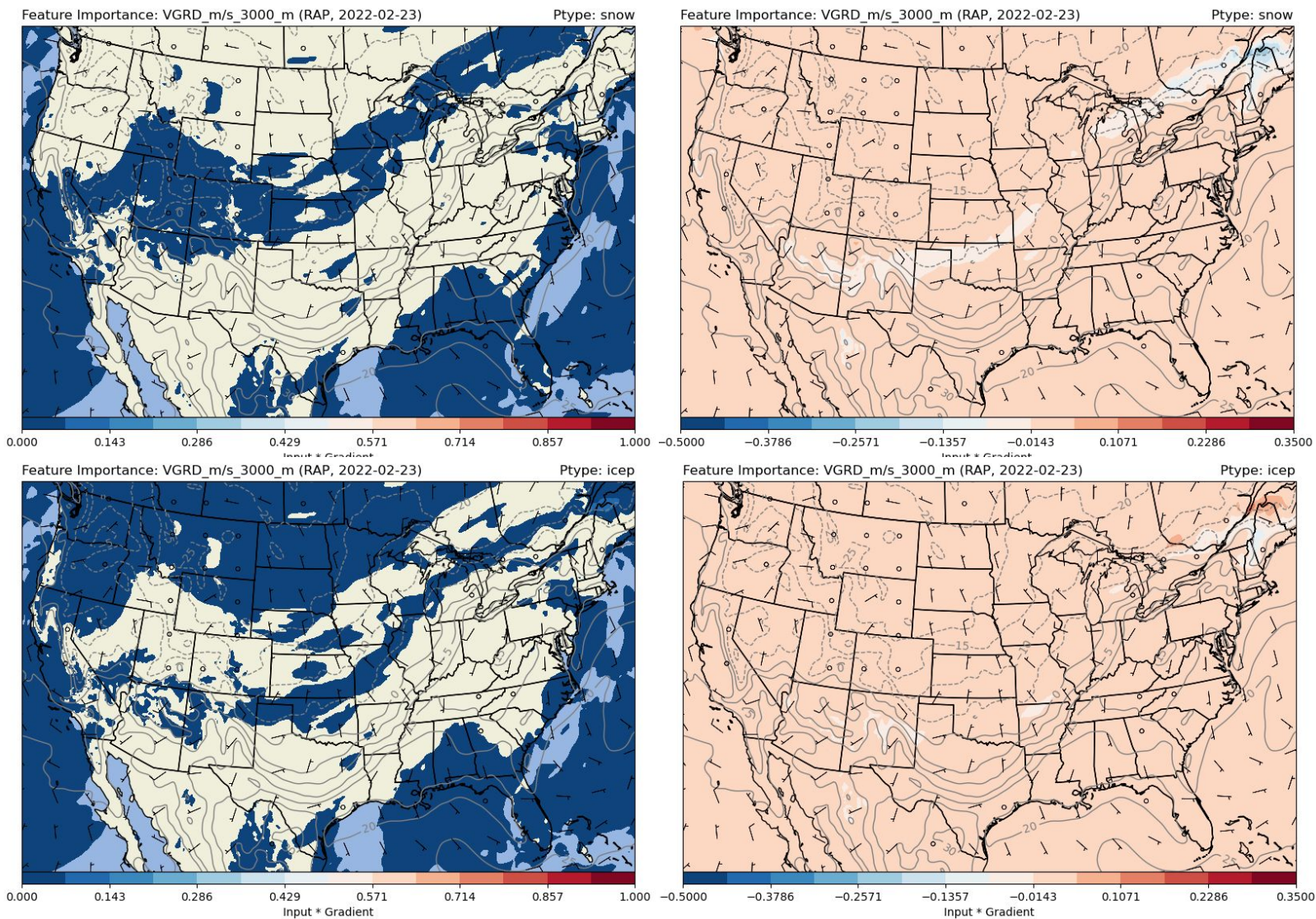
Feature Importance: UGRD\_m/s\_2250\_m (RAP, 2022-02-23) Ptype: rain



-0.5000 -0.3786 -0.2571 -0.1357 -0.0143 0.1071 0.2286 0.3500  
Input \* Gradient



# Feature Importance by precipitation type: VGRD



# Feature Importance by precipitation type: VGRD

