

# Using physics-based machine learning to estimate unobserved quantities: A case study for landscape-scale soil and vegetation conductances to heat and water vapor

Andrew Bennett<sup>1</sup>, Maoya Bassiouni<sup>2</sup>, and Bart Nijssen<sup>3</sup>

<sup>1</sup>University of Arizona

<sup>2</sup>Swedish University of Agricultural Sciences

<sup>3</sup>University of Washington

November 22, 2023

## Abstract

While machine learning (ML) techniques have proven to have exceptional performance in prediction of variables that have long and varied observational records, it is not clear how to use such techniques to learn about intermediate processes which may not be readily observable. We build on previous work that found that encoding either known, or approximated, physical relationships into the machine learning framework can allow the learned model to implicitly represent processes that are not directly observed, but can be related to an observable quantity. Zhao et al. (2019) found that encoding a Penman-Monteith-like equation of latent heat in an artificial neural network could reliably predict the latent heat while providing an estimate of the resistance term, which is not readily observable at the landscape scale. Specifically, we advance this framework in two ways. First, we expand the physics-based layer to include the partitioning of both the latent and sensible heat fluxes among the vegetation and soil domains, each with their own resistance terms. Second, we couple a land-surface model (LSM), which provides information from simulated processes to the ML model. We thus effectively provide the ML model with both physics-informed inputs as well as maintain constraints such as mass and energy balance on outputs of the coupled ML-LSM simulations. Previously we found that coupling a LSM to the ML model could provide good predictions of bulk turbulent heat fluxes, and in this work we explore how incorporating the additional physics-based partitioning allows the model to learn more ecohydrologically-relevant dynamics in diverse biomes. Further, we explore what the model learned in predicting the unobserved resistance terms and what we can learn from the model itself. Zhao, W. L., Gentile, P., Reichstein, M., Zhang, Y., Zhou, S., Wen, Y., et al. (2019). Physics-Constrained Machine Learning of Evapotranspiration. *Geophysical Research Letters*, 46(24), 14496–14507. <https://doi.org/10.1029/2019GL085291>

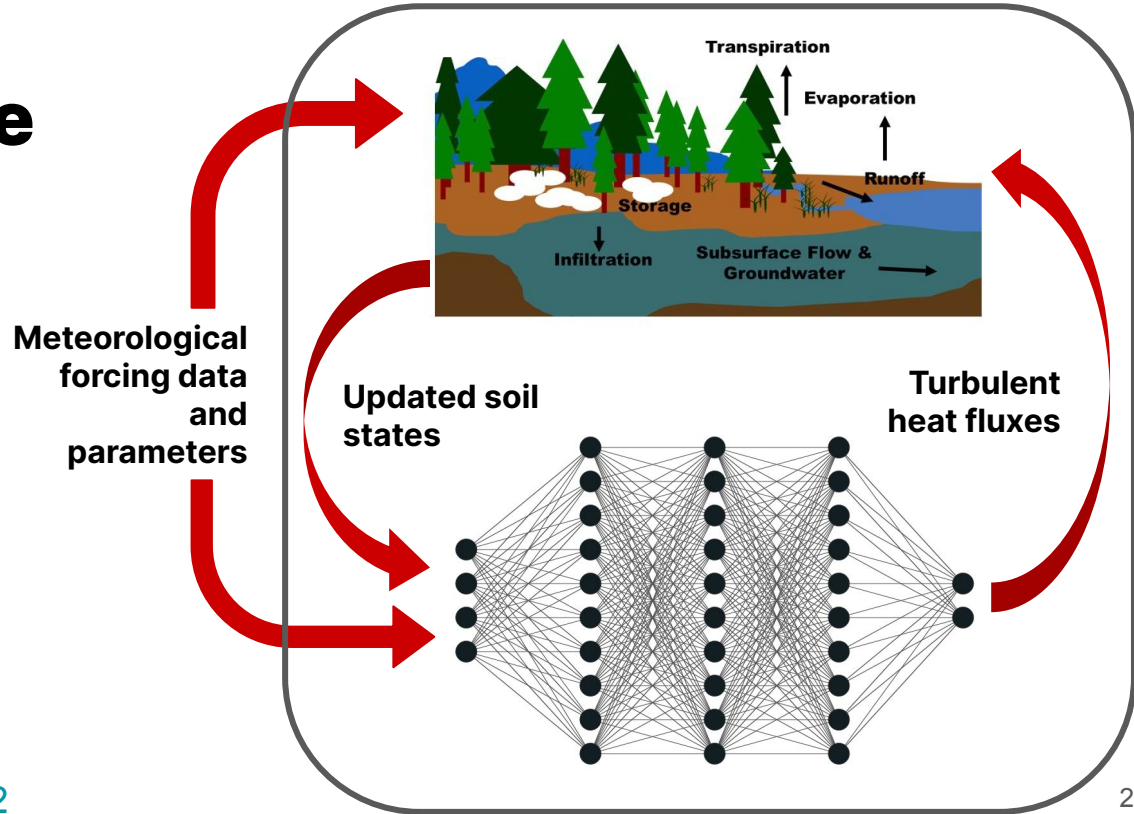
# Using physics-based machine learning to estimate unobserved quantities:

A case study for landscape-scale soil and vegetation conductances to heat and water vapor

**Andrew Bennett** (andrbenn@email.arizona.edu)  
**Maoya Bassiouni**  
**Bart Nijssen**

**We previously showed that a coupled DL-PBHM approach can make better predictions of turbulent heat fluxes than a PBHM alone**

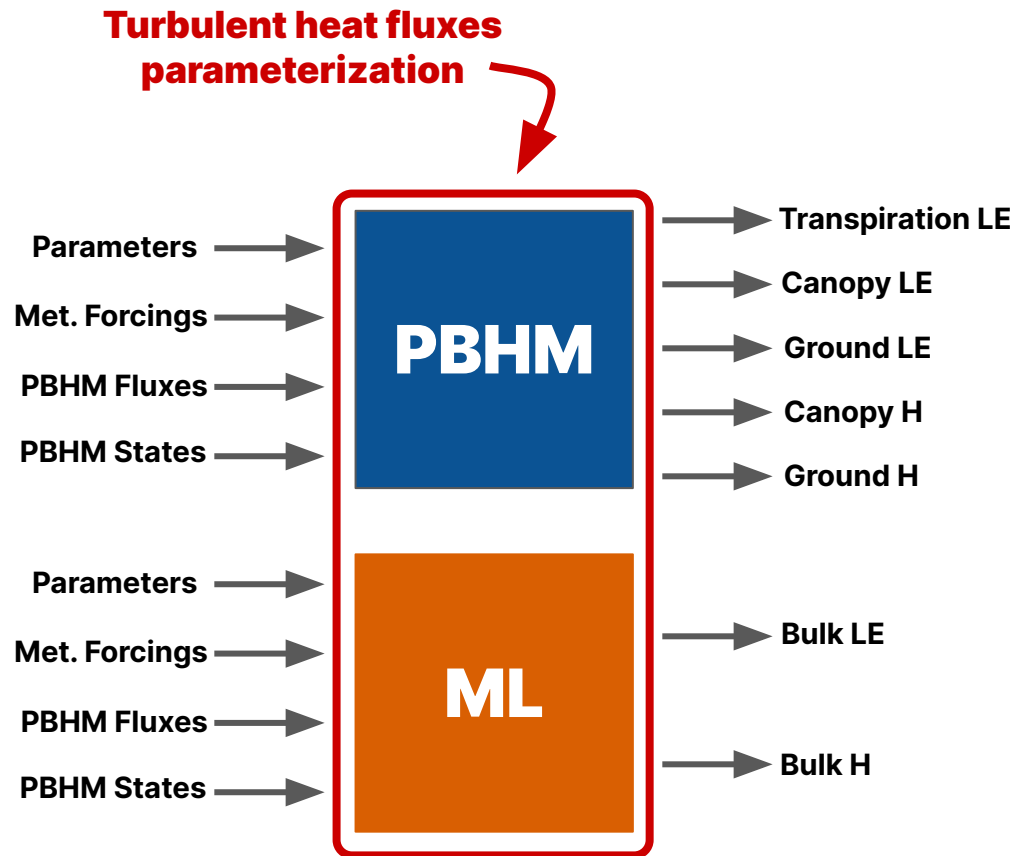
## **Process based hydrologic model (PBHM)**



Link to our previous work:

<https://doi.org/10.1002/essoar.10505081.2>

**One of the major shortcomings is a mismatch between process fidelity and the observed data for training**



Link to our previous work:

<https://doi.org/10.1002/essoar.10505081.2>

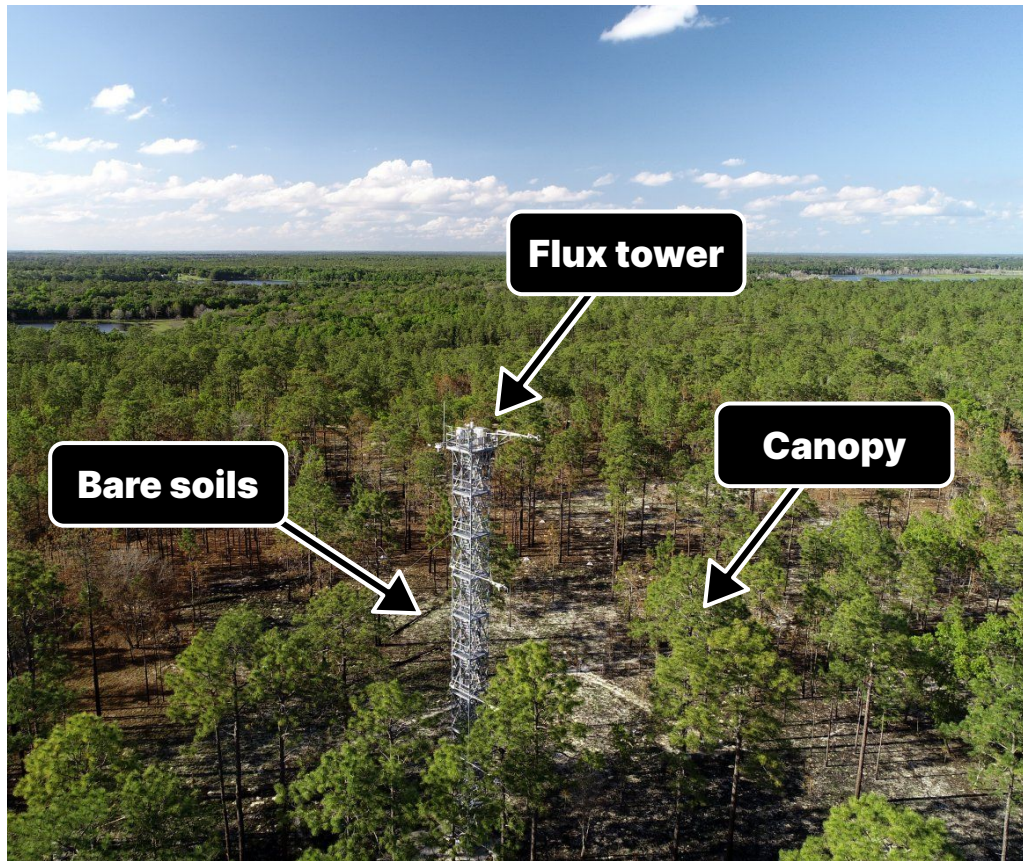
# A reminder: The land surface is heterogeneous!

Flux towers measure bulk fluxes





But we want to model the various  
components

Without fancy techniques supervised  
machine learning can only learn bulk  
fluxes from observations then

This presentation is about one of  
these “fancy techniques”







# So, we've got tradeoffs

	Process based model	Machine learned model
Superior performance		
Process fidelity		



Why don't we have both?

# So, we've got tradeoffs

	Process based model	Machine learned model
Superior performance		
Process fidelity		



# Why don't "physics" based models perform well?

These bulk transfer equations are very common in hydrologic and land surface modeling:

- Andreadis et al., 2009
- Bonan, 1991
- Inclan and Forkel, 1995
- Sellers et al., 1986
- Mahat et al., 2013
- Clark et al., 2015
- ...

$$\left. \begin{aligned} Q_h^{veg} &= -\rho_{air} c_p C_h^{veg} (T^{veg} - T^{cas}) \\ Q_h^{sfc} &= -\rho_{air} c_p C_h^{sfc} (T^{sfc} - T^{cas}) \end{aligned} \right\} \text{Sensible heat fluxes}$$

$$\left. \begin{aligned} Q_{evap}^{veg} &= -\frac{L_{vap} \rho_{air} \varepsilon}{P_{air}} C_{evap}^{veg} [e_{sat}(T^{veg}) - e^{cas}] \\ Q_{trans}^{veg} &= -\frac{L_{vap} \rho_{air} \varepsilon}{P_{air}} C_{trans}^{veg} [e_{sat}(T^{veg}) - e^{cas}] \\ Q_l^{sfc} &= -\frac{L_{vap} \rho_{air} \varepsilon}{P_{air}} C_w^{sfc} [\phi_{hum}^{sfc} e_{sat}(T^{sfc}) - e^{cas}] \end{aligned} \right\} \text{Latent heat fluxes}$$

# Why don't “physics” based models perform well?

These consist of three main parts

1. Constants & parameters
2. Temperature or moisture gradients
3. Conductance terms

$$\left. \begin{aligned} Q_h^{veg} &= -\underbrace{\rho_{air}}_{\text{orange}} \underbrace{c_p}_{\text{green}} \underbrace{C_h^{veg}}_{\text{purple}} \left( \underbrace{T^{veg}}_{\text{purple}} - \underbrace{T^{cas}}_{\text{purple}} \right) \\ Q_h^{sfc} &= -\underbrace{\rho_{air}}_{\text{orange}} \underbrace{c_p}_{\text{green}} \underbrace{C_h^{sfc}}_{\text{purple}} \left( \underbrace{T^{sfc}}_{\text{purple}} - \underbrace{T^{cas}}_{\text{purple}} \right) \end{aligned} \right\} \text{Sensible heat fluxes}$$

$$\left. \begin{aligned} Q_{evap}^{veg} &= -\underbrace{\frac{L_{vap} \rho_{air} \epsilon}{P_{air}}}_{\text{orange}} \underbrace{C_{evap}^{veg}}_{\text{green}} \left[ \underbrace{e_{sat}(T^{veg}) - e^{cas}}_{\text{purple}} \right] \\ Q_{trans}^{veg} &= -\underbrace{\frac{L_{vap} \rho_{air} \epsilon}{P_{air}}}_{\text{orange}} \underbrace{C_{trans}^{veg}}_{\text{green}} \left[ \underbrace{e_{sat}(T^{veg}) - e^{cas}}_{\text{purple}} \right] \\ Q_l^{sfc} &= -\underbrace{\frac{L_{vap} \rho_{air} \epsilon}{P_{air}}}_{\text{orange}} \underbrace{C_w^{sfc}}_{\text{green}} \left[ \underbrace{\phi_{hum}^{sfc} e_{sat}(T^{sfc}) - e^{cas}}_{\text{purple}} \right] \end{aligned} \right\} \text{Latent heat fluxes}$$

# Why don't “physics” based models perform well?

These consist of three main parts

1. Constants & parameters
2. Temperature or moisture gradients
3. Conductance terms



# Why don't "physics" based models perform well?

These consist of three main parts

1. **Constants & parameters**
2. **Temperature or moisture gradients**
3. **Conductance terms**

I'm going to argue these are either:


1. Pretty well known
2. Parts of other processes

# Why don't “physics” based models perform well?

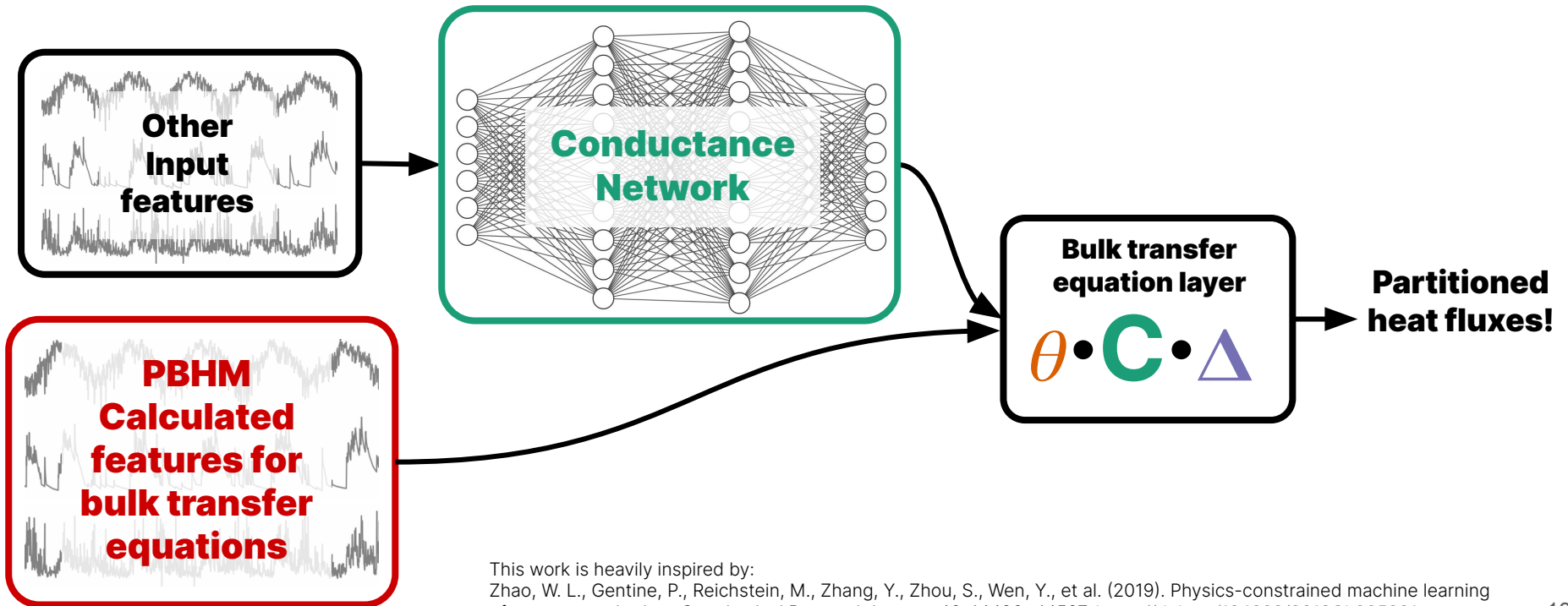
These consist of three main parts

1. Constants & parameters
2. Temperature or moisture gradients
3. Conductance terms

And likewise, this is  
where the model  
uncertainty really is...



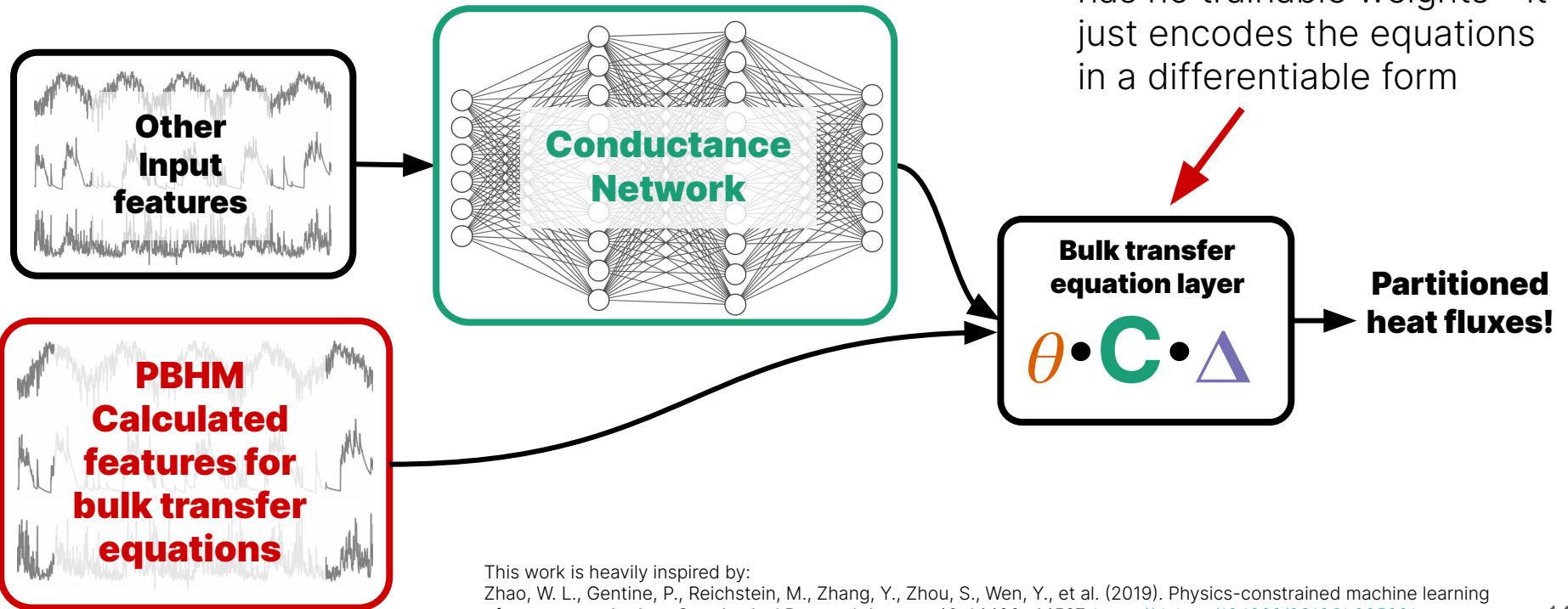
# The hybrid neural network architecture



This work is heavily inspired by:

Zhao, W. L., Gentine, P., Reichstein, M., Zhang, Y., Zhou, S., Wen, Y., et al. (2019). Physics-constrained machine learning of evapotranspiration. *Geophysical Research Letters*, 46, 14496–14507. <https://doi.org/10.1029/2019GL085291>

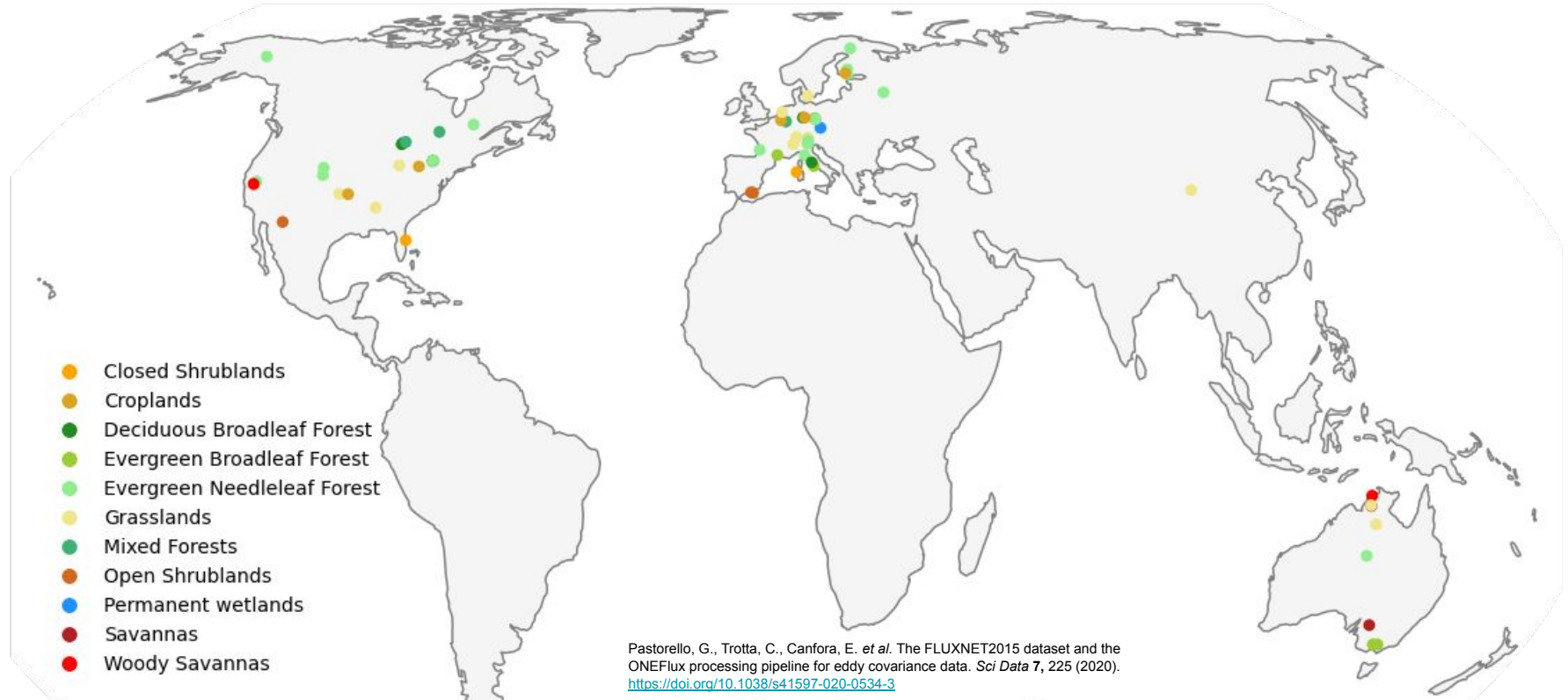
# The hybrid neural network architecture



Technical note: This "layer" has no trainable weights - it just encodes the equations in a differentiable form

This work is heavily inspired by:  
Zhao, W. L., Gentine, P., Reichstein, M., Zhang, Y., Zhou, S., Wen, Y., et al. (2019). Physics-constrained machine learning of evapotranspiration. *Geophysical Research Letters*, 46, 14496– 14507. <https://doi.org/10.1029/2019GL085291>

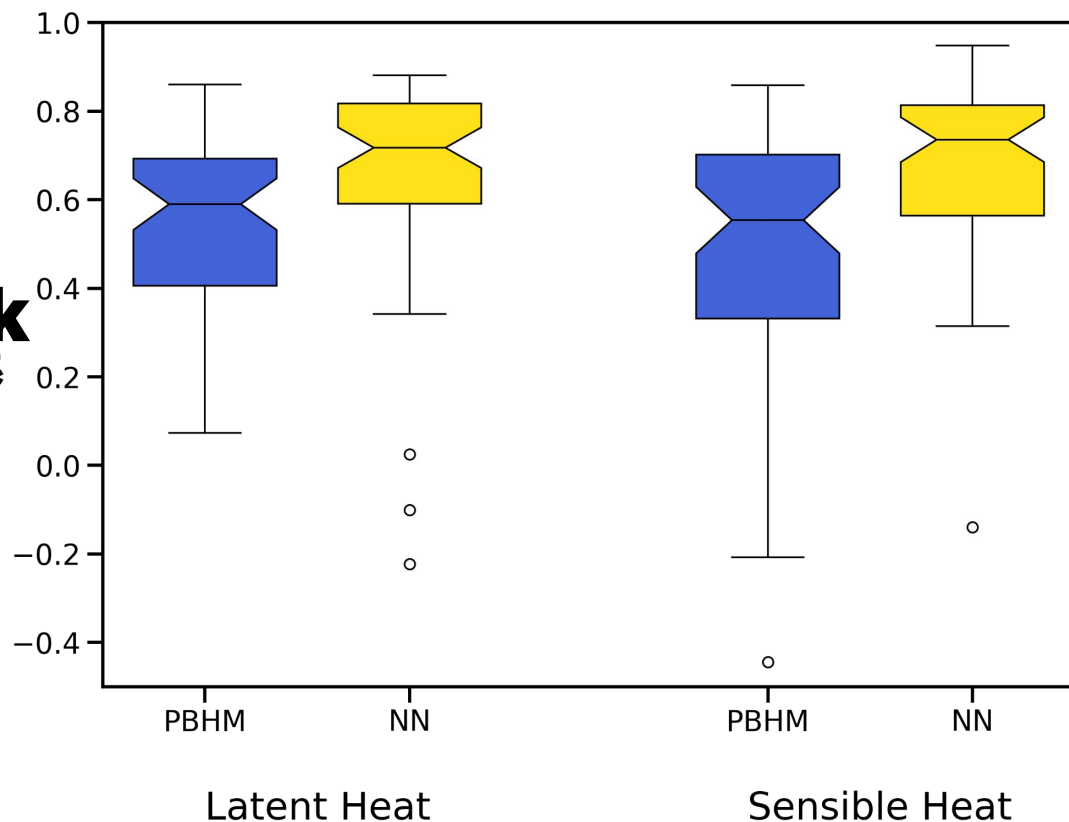
# We gathered data from 60 FluxNet sites, totalling over 500 site-years of half-hourly data





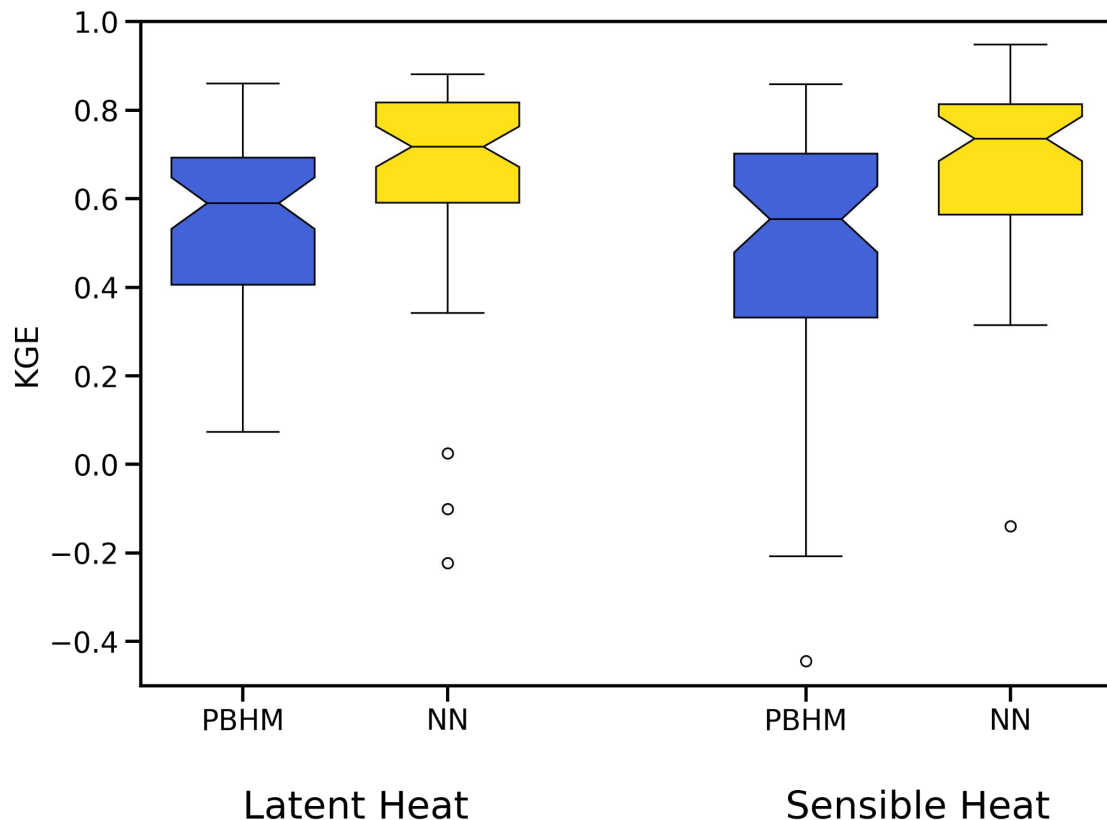
**We're still able to outperform a calibrated PBHM using the same bulk transfer equations**

*\*other pure ML based approaches outperform this but that's asking a different question*

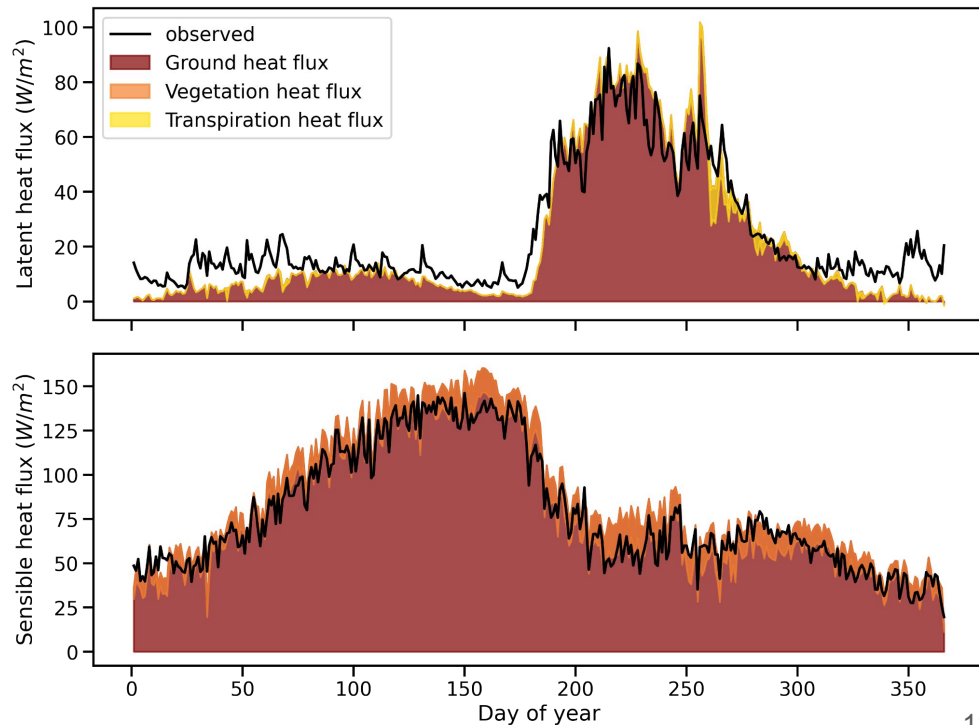


# Perhaps a rough upper bound on the performance of the “physics” based equations?

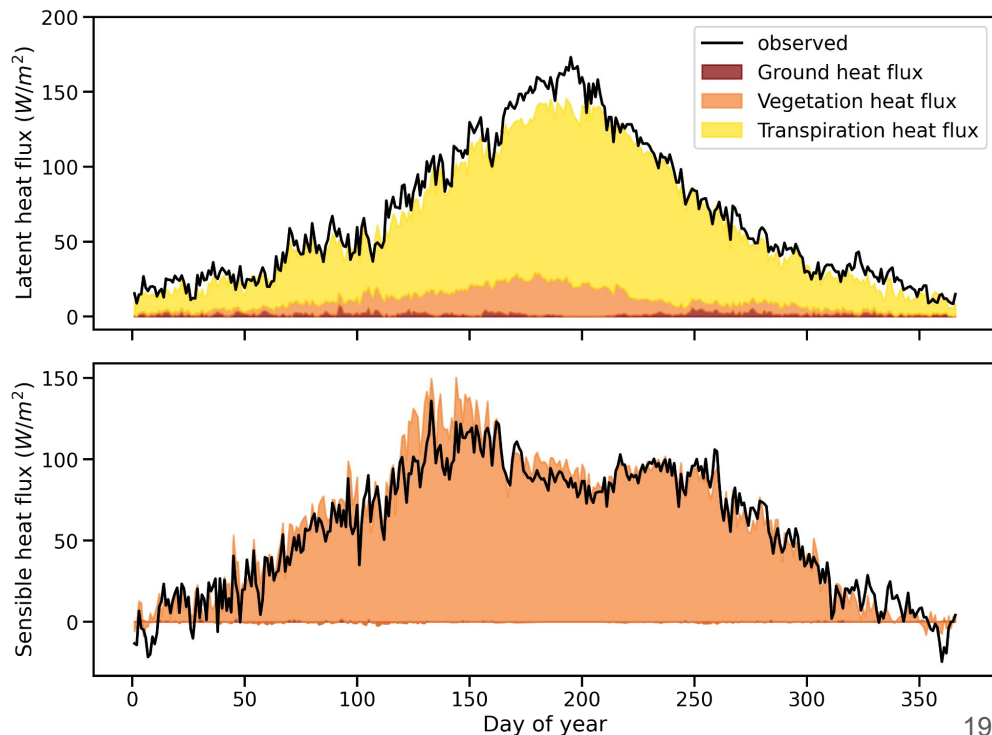
*\*more work remains to be done to ensure that these performance results are optimal (notably fixing outliers)*



# Walnut Gulch near Tombstone, AZ (US-Whs) shows ground component is largest

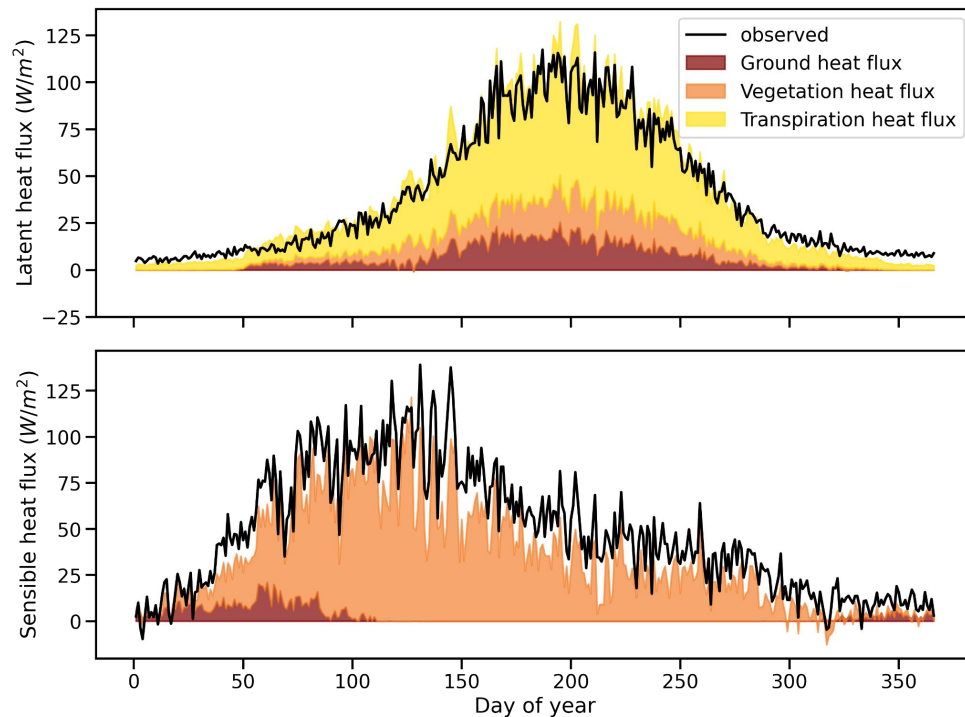


# Blodgett Forest near Sacramento, CA (US-Blo) shows vegetation components are largest



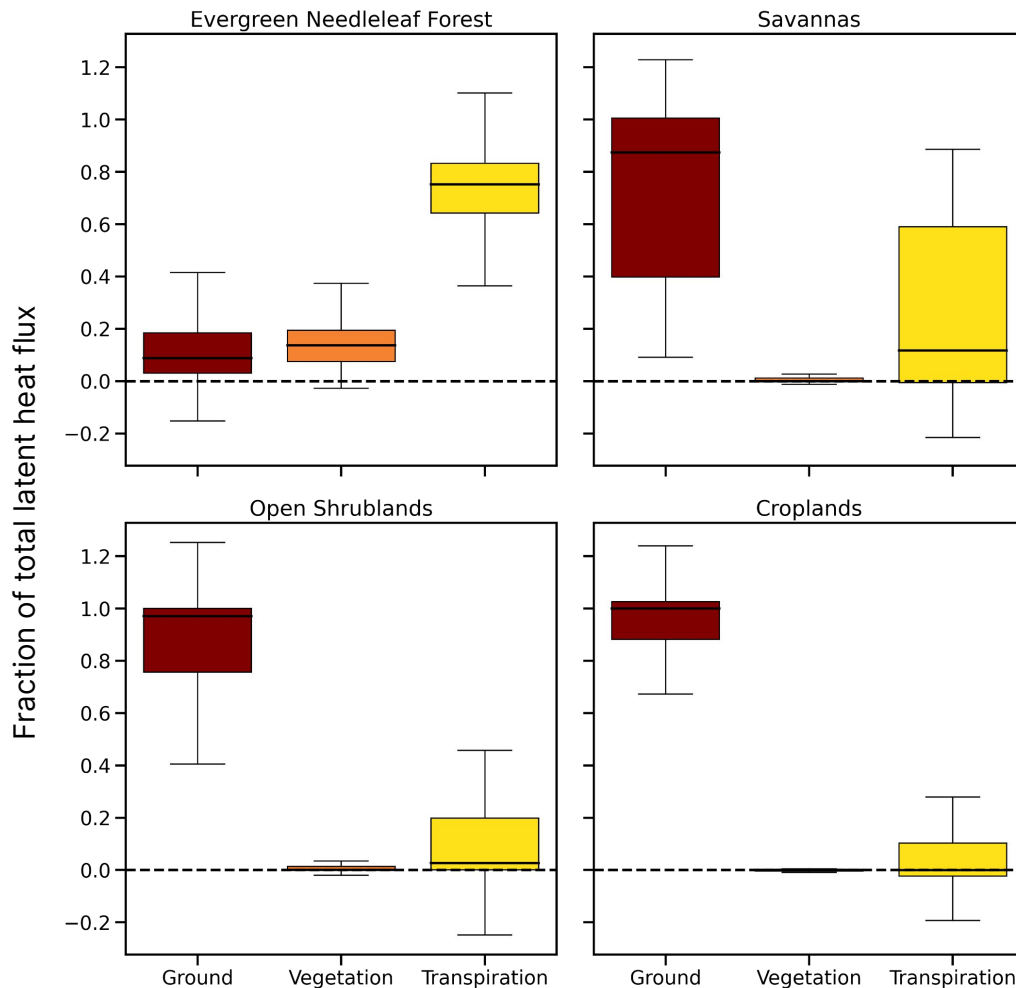


# Mixed forest near Vielsalm, Belgium shows a larger mixture between components



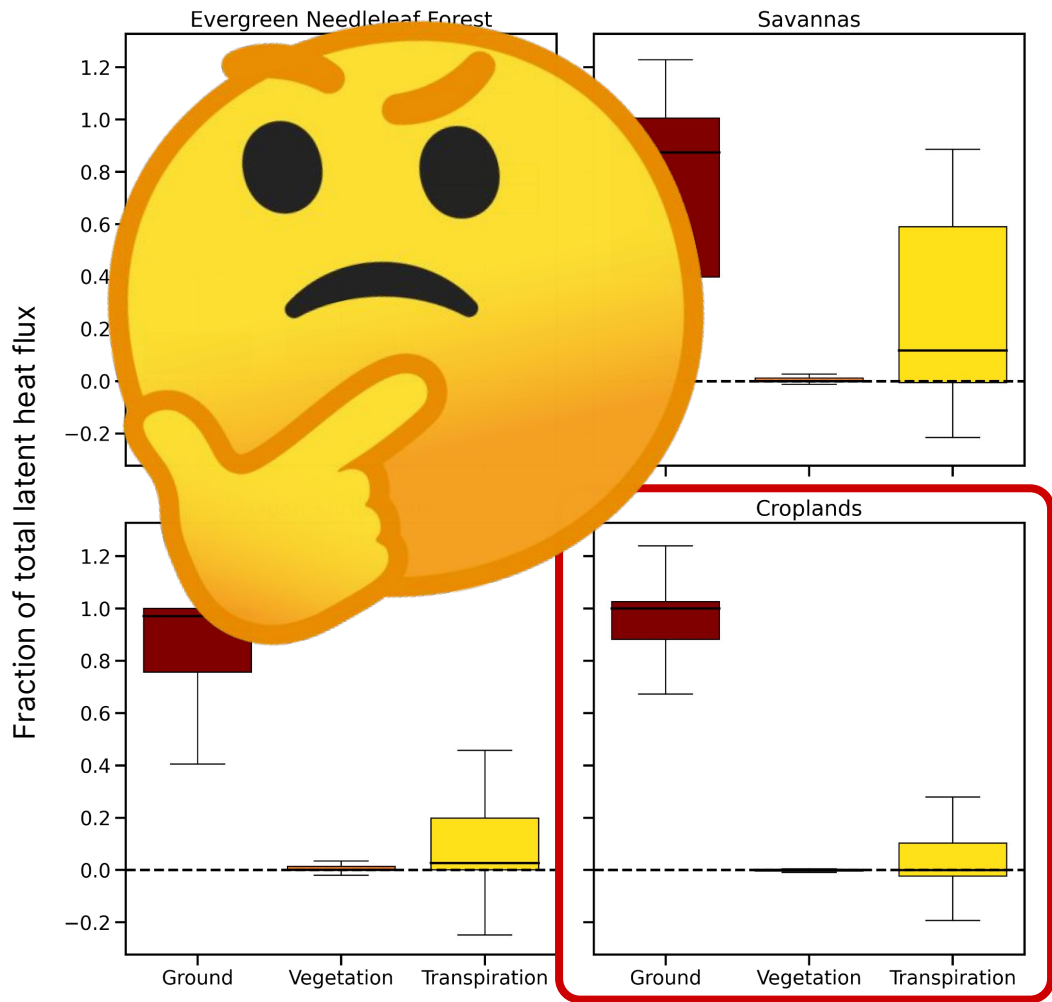
# The overall partitioning matches physical intuition to a first order

*Note: values can be  $<0$  and  $>1$  because condensation exists*

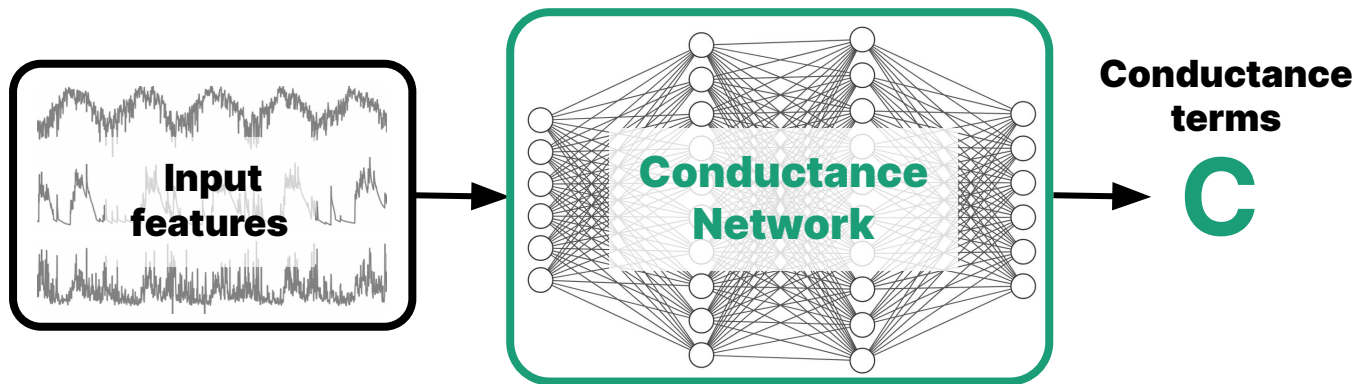


**The overall  
partitioning  
matches  
physical  
intuition to a  
first order,  
mostly**

*Note: values can be  $<0$  and  $>1$  because  
condensation exists*

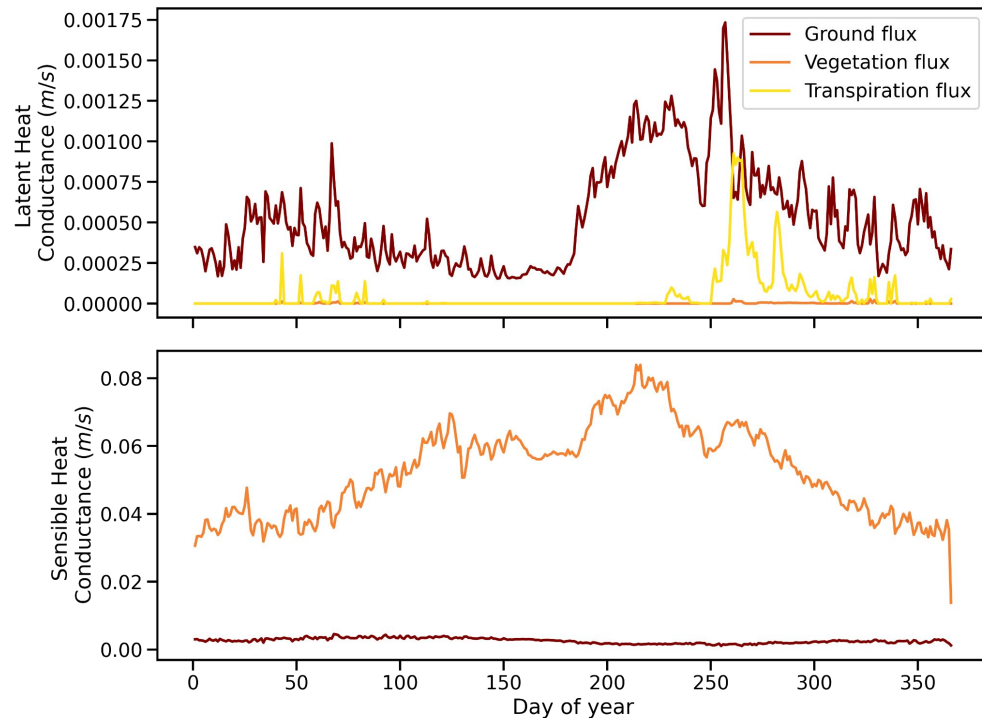
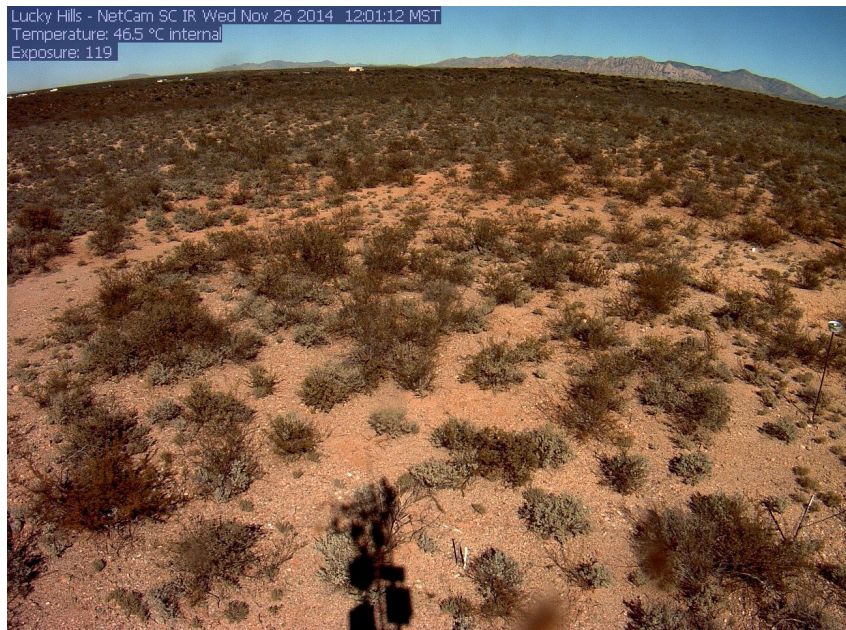


# We can also truncate the network to analyze the conductances!

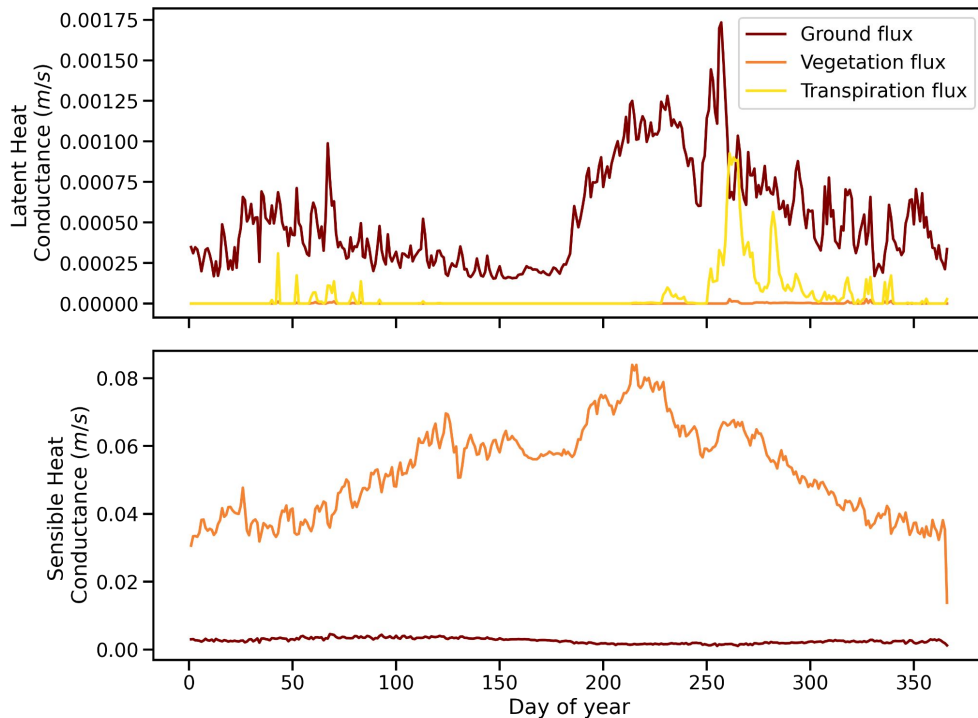
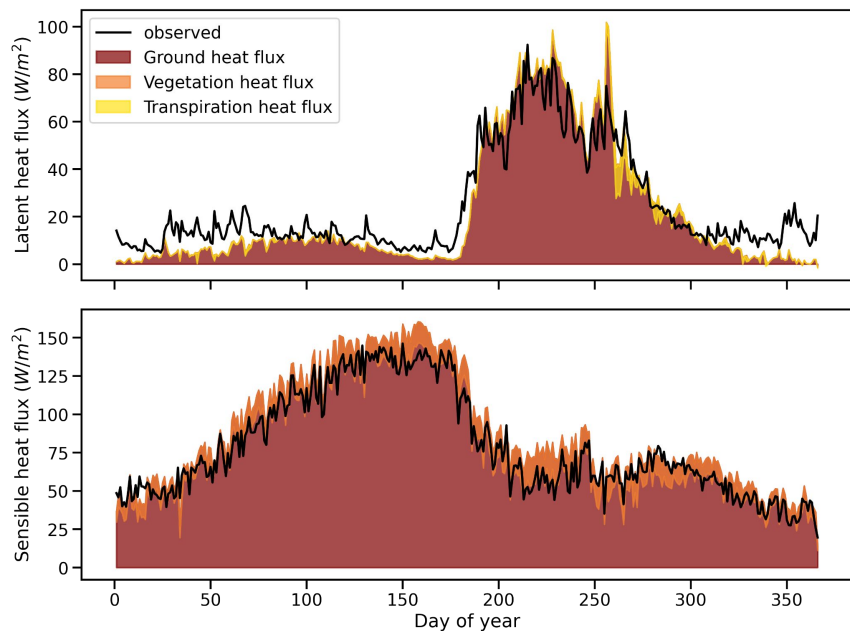




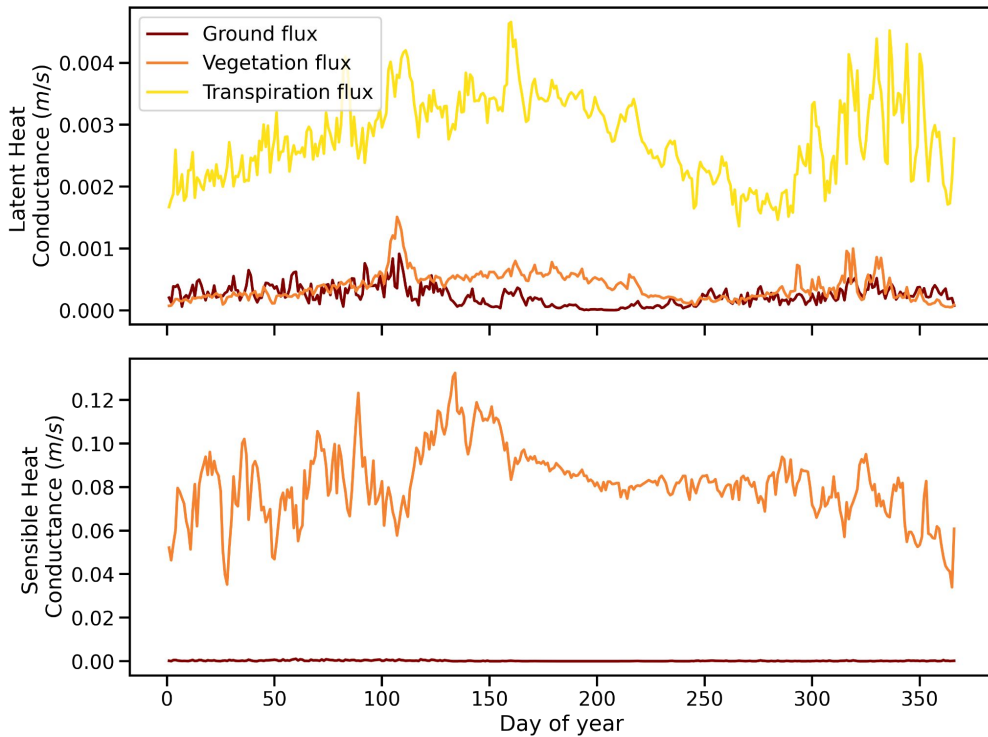
# Conductances at Walnut Gulch



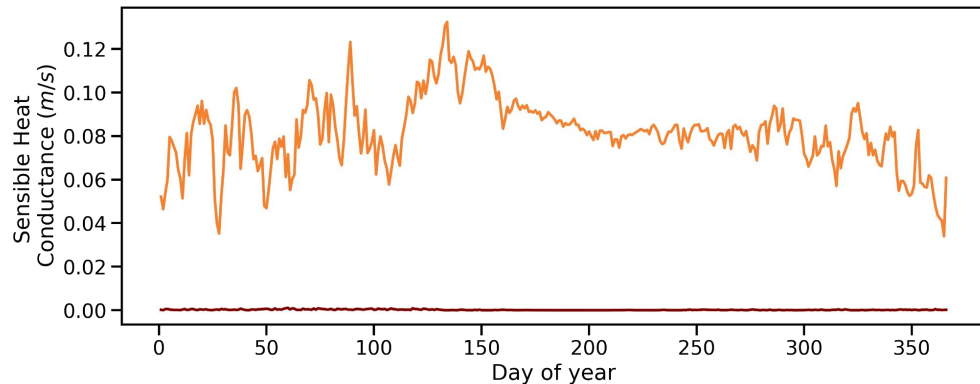
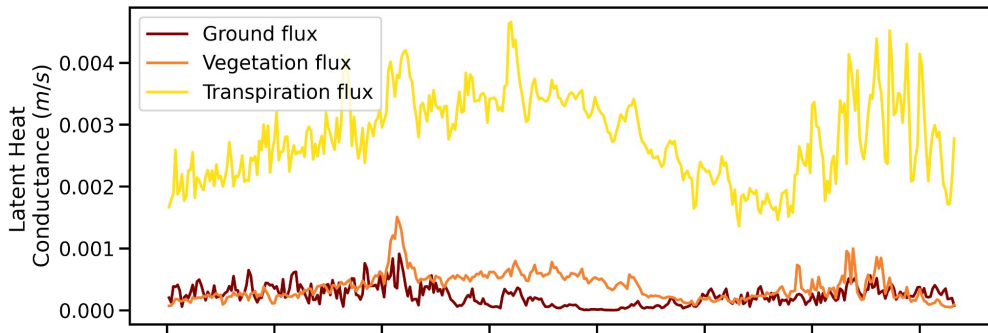
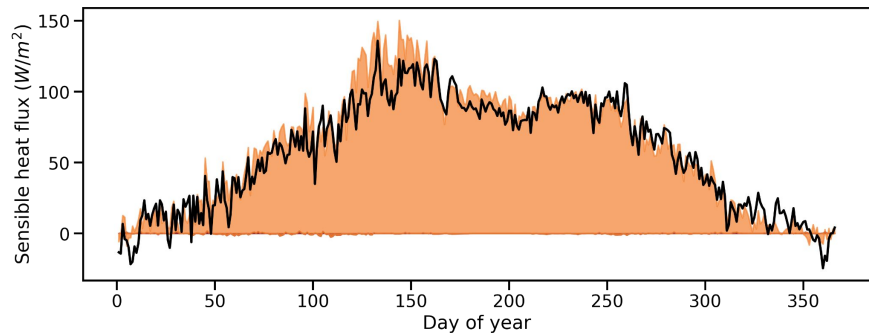
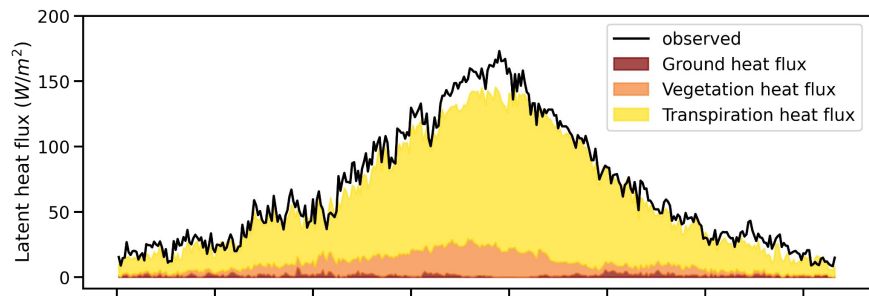
# Conductances are not 1-1 with heat fluxes



# Conductances at Blodgett Forest



# Conductances are not 1-1 with heat fluxes *still*



# Wrapup, future work, and some conceptual takeaways

We've quantified that a large amount of predictive performance is due to conductance terms

Need methods/data to better constrain the partitioning, particularly at sites with human interventions, like croplands

Coupling to the PBHM is still incomplete, but needed to analyze the effects on the full water cycle

**My big takeaway:** Process-based vs data-driven modeling should be a spectrum rather than a binary choice