



# Scaling Prior-Data Fitted Networks for Physical System Learning

Parikshit Pareek, Ph.D.

Assistant Professor

Department of Electrical Engineering

Indian Institute of Technology Roorkee (IIT Roorkee)

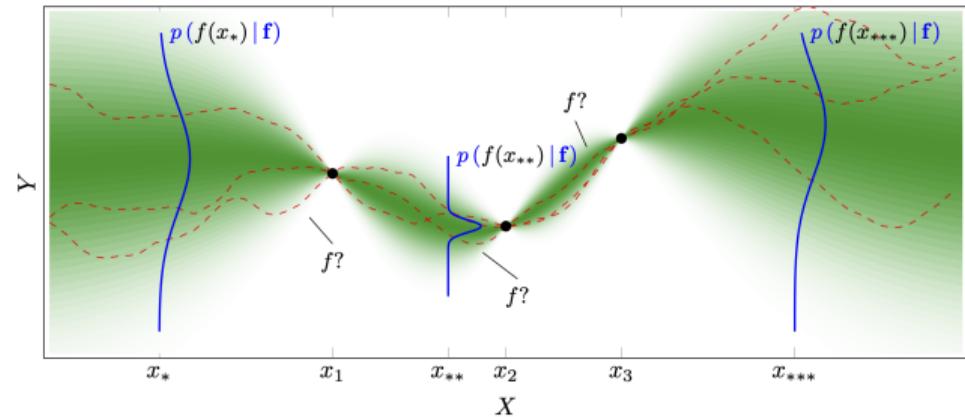
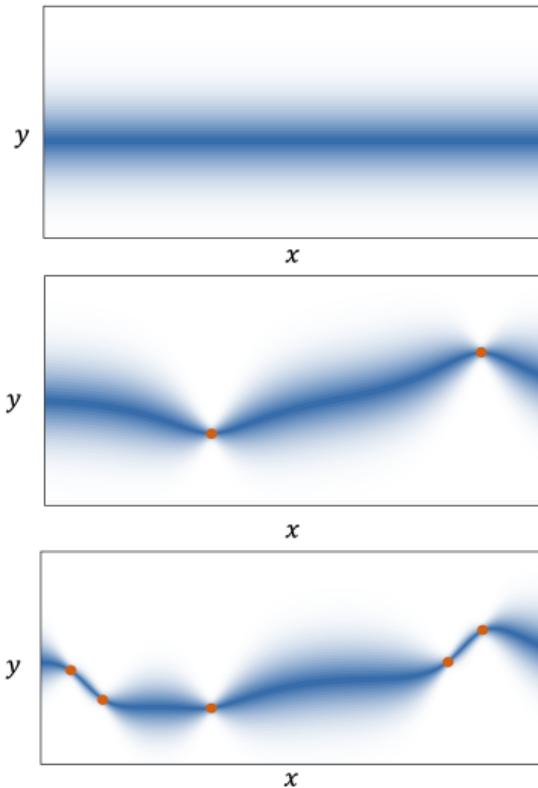
## Time Equation for ML Models

$$T_{data} + T_{train} + T_{predict}$$

$$\mathcal{M}(\theta_i | \mathcal{D}_i) : \text{Need Retraining for Each } \mathcal{D}.$$

Training Data  
↓  
 $\mathcal{M}(\theta_i | \mathcal{D}_i)$   
↑  
Model Weights

# Gaussian Process Working Idea: In Brief



The distribution of  $f(x^*)$ ,  $f(x^{**})$  and  $f(x^{***})$  for the three inputs  $x^*$ ,  $x^{**}$  and  $x^{***}$ , conditional on the observed values  $f(x_1)$ ,  $f(x_2)$  and  $f(x_3)$ .

# Prior-Data Fitted Networks (PFNs)

**Goal:** Perform BAYESIAN INFERENCE in a single forward pass.

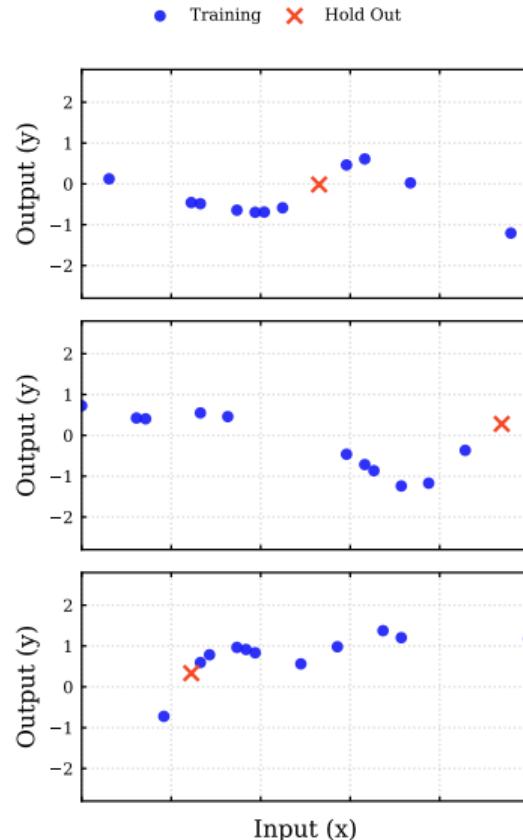
**Why?** Bayesian inference is slow, and repeated retraining makes the computational cost prohibitive

**How?** Learn to approximate the *Posterior Predictive Distribution (PPD)*  $p(y|x, \mathcal{D})$

## Amortized Bayesian Inference

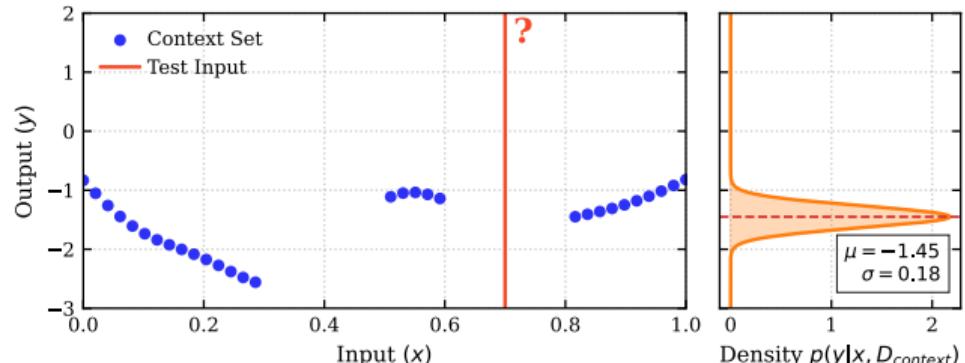
Given a new dataset  $\mathcal{D}_{\text{context}}$  and query point  $x_{\text{test}}$ , it outputs  
 $q_{\theta^*}(y_{\text{test}} | x_{\text{test}}, \mathcal{D}_{\text{context}}) \approx p(y_{\text{test}} | x_{\text{test}}, \mathcal{D}_{\text{context}})$   
in a single forward pass, where  $\theta^*$  are the optimal PFN parameters.

# PFN Idea



Train PFN using Negative Log-Likelihood Loss over hold-out samples

$$\ell_{\theta} = \sum_{k=1}^K [ -\log q_{\theta}(\mathbf{y}^k | \mathbf{x}^k, \mathcal{D}^k) ]$$



# PFN: Training & Inference

Sample prior datasets  $D^{(i)} \sim p(\mathcal{D})$

$$D^{(1)} = D_{train}^{(1)} \cup \{(x_{test}^{(1)}, y_{test}^{(1)})\}$$

⋮

$$D^{(K)} = D_{train}^{(K)} \cup \{(x_{test}^{(K)}, y_{test}^{(K)})\}$$

Train the PFN by minimizing

$$-\sum_{i=1}^K \log q_\theta(y_{test}^{(i)} | x_{test}^{(i)}, D_{train}^{(i)})$$

Actual context dataset & test input

$$(D_{context}, x_{test})$$

$$q_{\theta^*}(y_{test} | x_{test}, D_{context}) \approx p(y_{test} | x_{test}, D_{context})$$

PFN with parameters  $\theta^*$

# One + Two Questions

## Scalability

Scalability: Can PFNs perform GP-style inference for 50D regression?  
No - standard PFNs don't scale beyond  $\sim 10$ D

## Attention Encoding: Localization

Why is everyone encoding  $x + y$  in PFNs? How does this affect localization?  
**Localization:** Nearby points provide more info about the value at an unknown point than distant ones.

## Backbone vs Attention

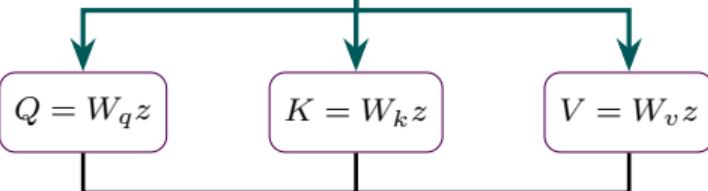
Are transformers all we need, or can PFNs also be built using CNNs etc.?

# Main Idea: Decouple Input and Output

## Vanilla Attention

joint embedding

$$z_i = \phi_x(x_i) + \phi_y(y_i)$$



$$H = \text{softmax}\left(\frac{Q(z)K(z)^\top}{\sqrt{d_k}}\right)V(z)$$

Mixing Input & Output

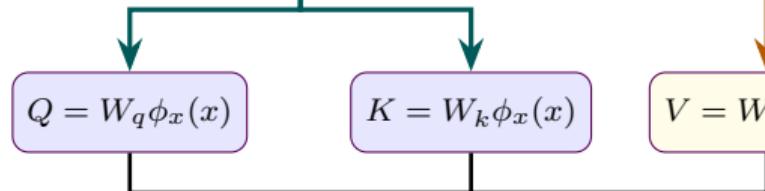
## Decoupled-Value Attention (DVA)

input encoder

$$\phi_x(x_i)$$

value encoder

$$\phi_y(y_i)$$



$$H = \text{softmax}\left(\frac{Q(x)K(x)^\top}{\sqrt{d_k}}\right)V(y)$$

Keeping Input & Output Separate

# Theorem: Enforcing Localization

## Theoretical Result

For DVA with linear embeddings, the attention weight  $\alpha_i(x_*)$  (for a test input  $x_*$ ) assigned to a context point  $x_i$  is proportional to a **Mahalanobis RBF Kernel**:

$$\alpha_i(x_*) \propto \exp\left(-\frac{1}{2\tau}\|(x_* - x_i)\|_A^2\right)$$

## Implication:

- As the distance  $\|x_* - x_i\|$  increases, the attention weight decays **exponentially**.
- This mathematically forces the Model to behave like a GP Kernel: localization.
- This holds true for any backbone architecture.

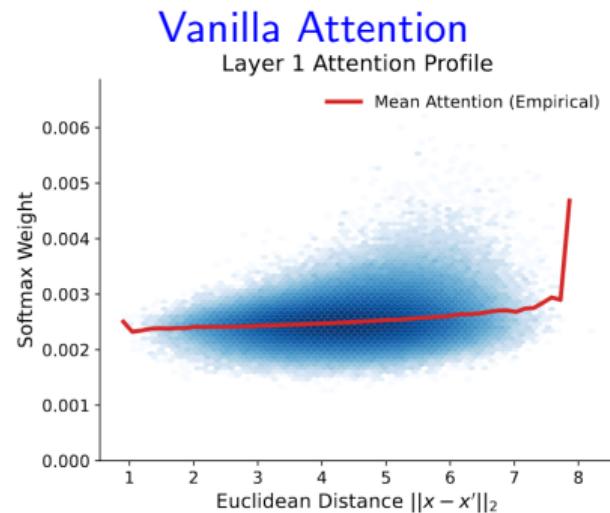
---

$$\|u\|_A^2 = u^T A u$$

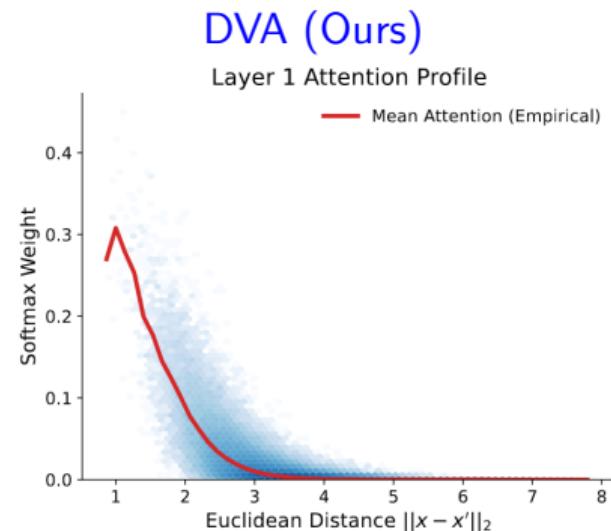
( $\mathcal{P}^2$ -LAB, EE, IITR)

# Empirical Evidence: 10D Input Space

Does the theory hold in practice? We look at the attention weights in a 10D task.

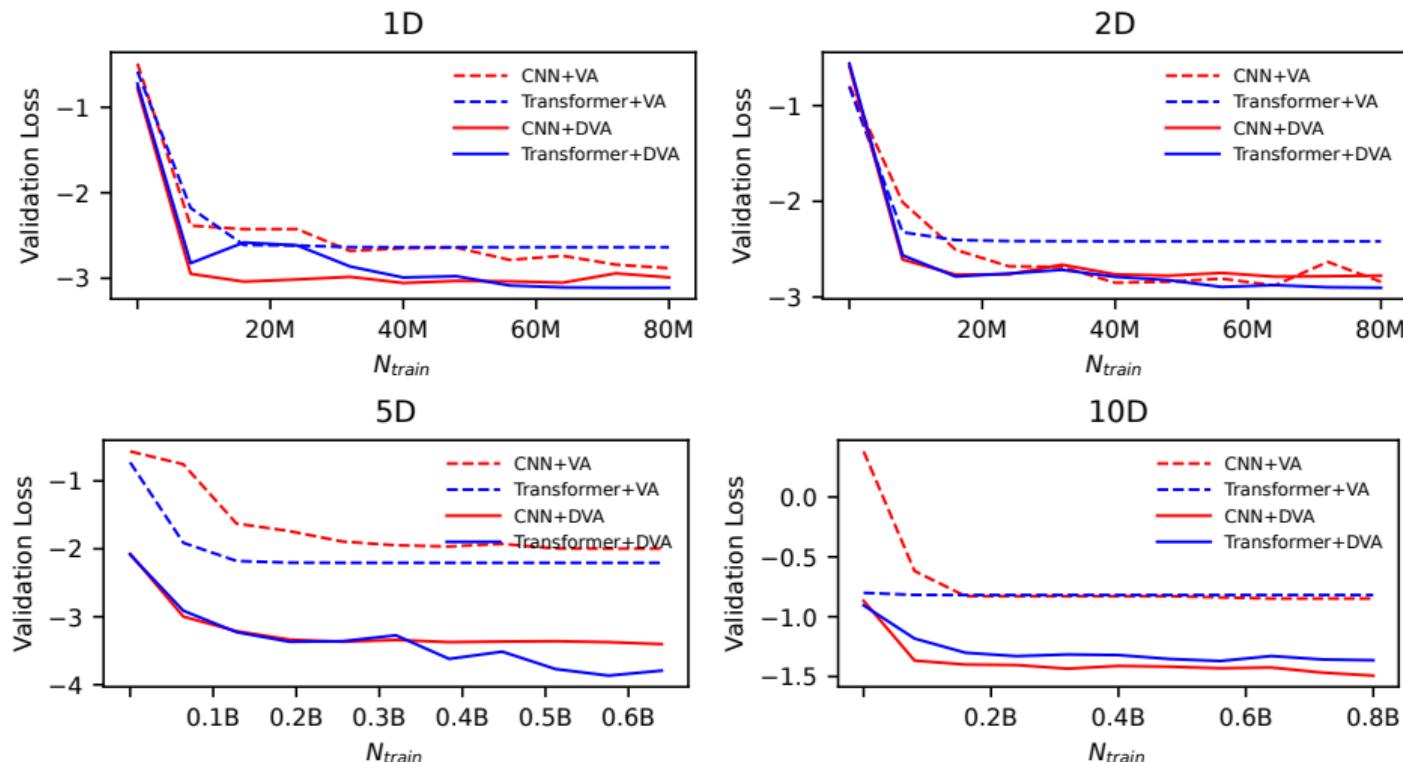


**No Localization**  
Weights are flat/uniform.  
The model fails to localize inputs.

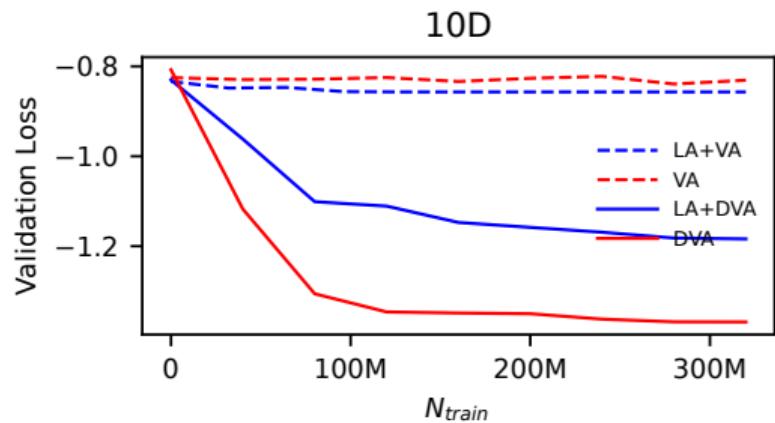
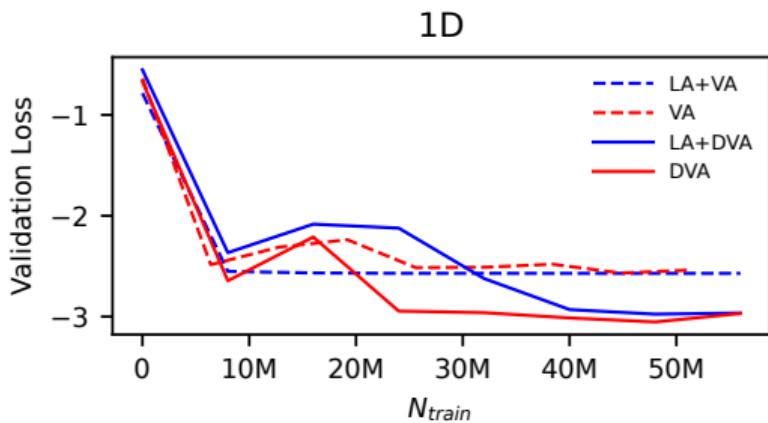


**Exponential Decay**  
Weights drop as distance increases.  
(Matches Theorem 1)

# DVA Outperforms VA: Scale and Error

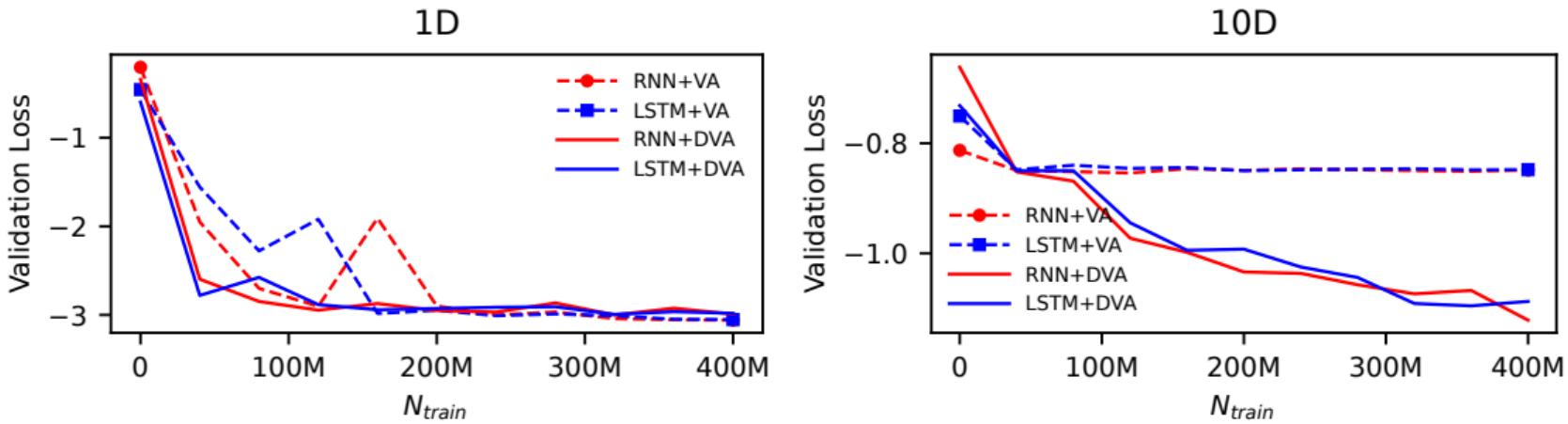


# Decoupling Helps with Linear Attention Too



Softmax DVA  $\gg$  Linear DVA  $\gg$  Softmax VA  $\gg$  Linear VA

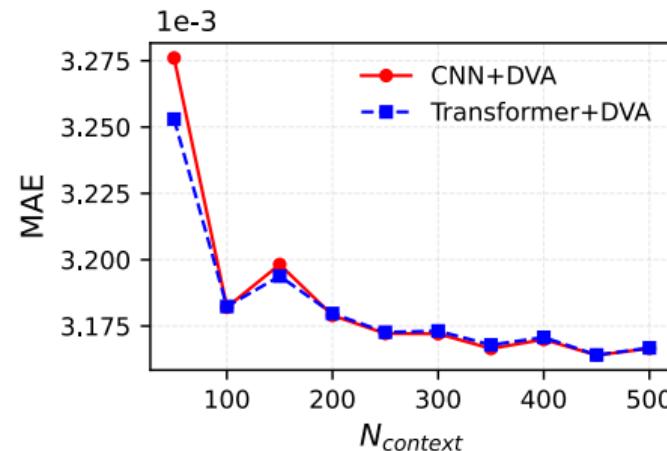
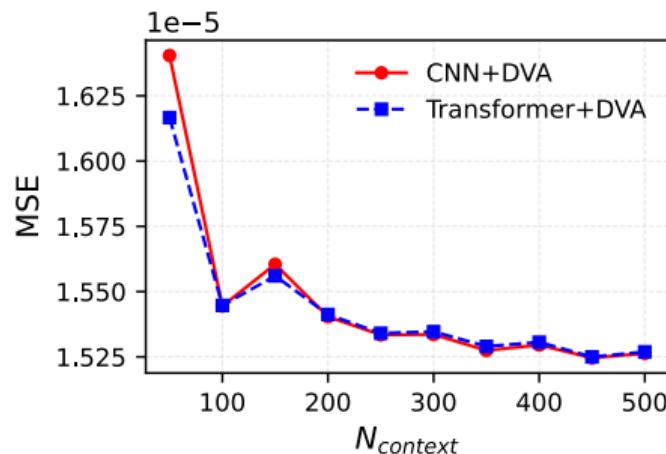
# Attention Agnostic: RNN/LSTM Results



With DVA, PFN training performance is nearly identical across backbones  
: CNN, RNN, LSTM, Transformer.

# 64D Power-Flow Learning

- Model: IEEE 33-bus AC power flow (64D inputs: 32 real + 32 reactive loads) → 32 bus voltages.
- Exact GP (per-bus) yields lowest MSE/MAE but is *slow* for real-time (needs 32 GPs).
- CNN+DVA and Transformer+DVA PFNs trade small accuracy loss ( $\text{MAE} \sim 10^{-3}$ ) for  $\sim 80\times$  faster inference.



# The PPD of Power Flow

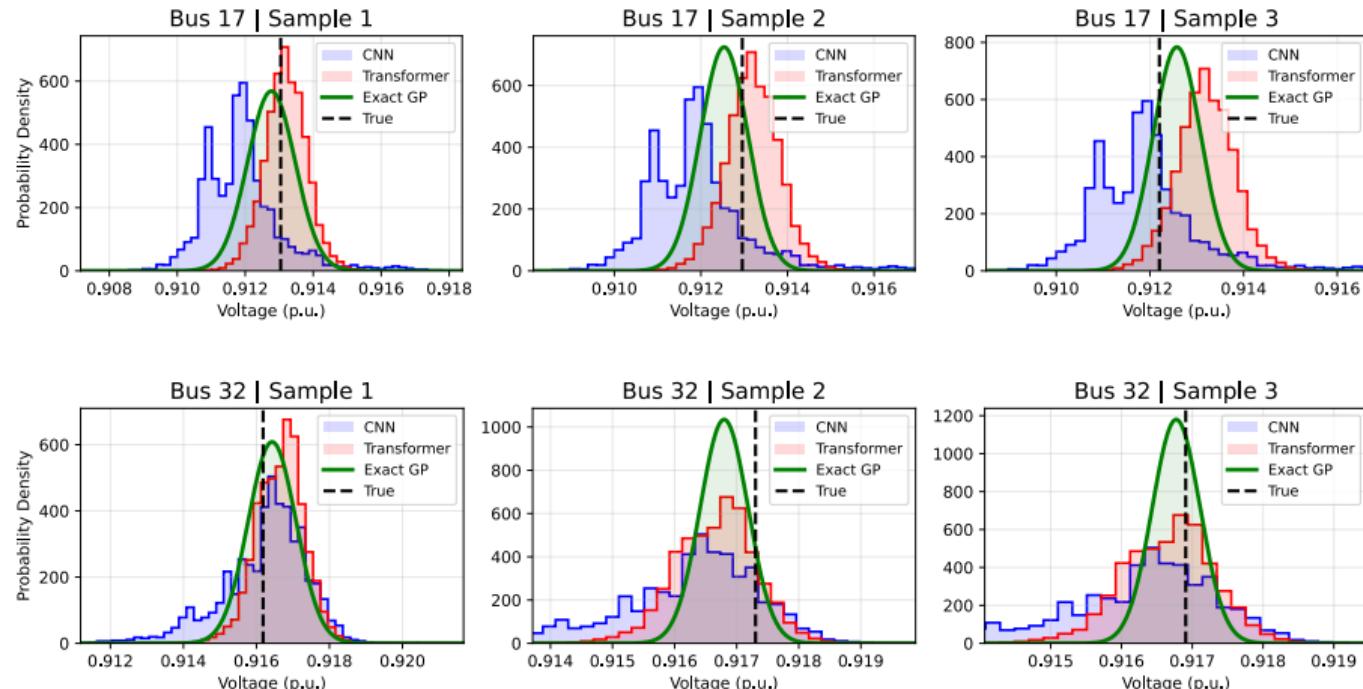
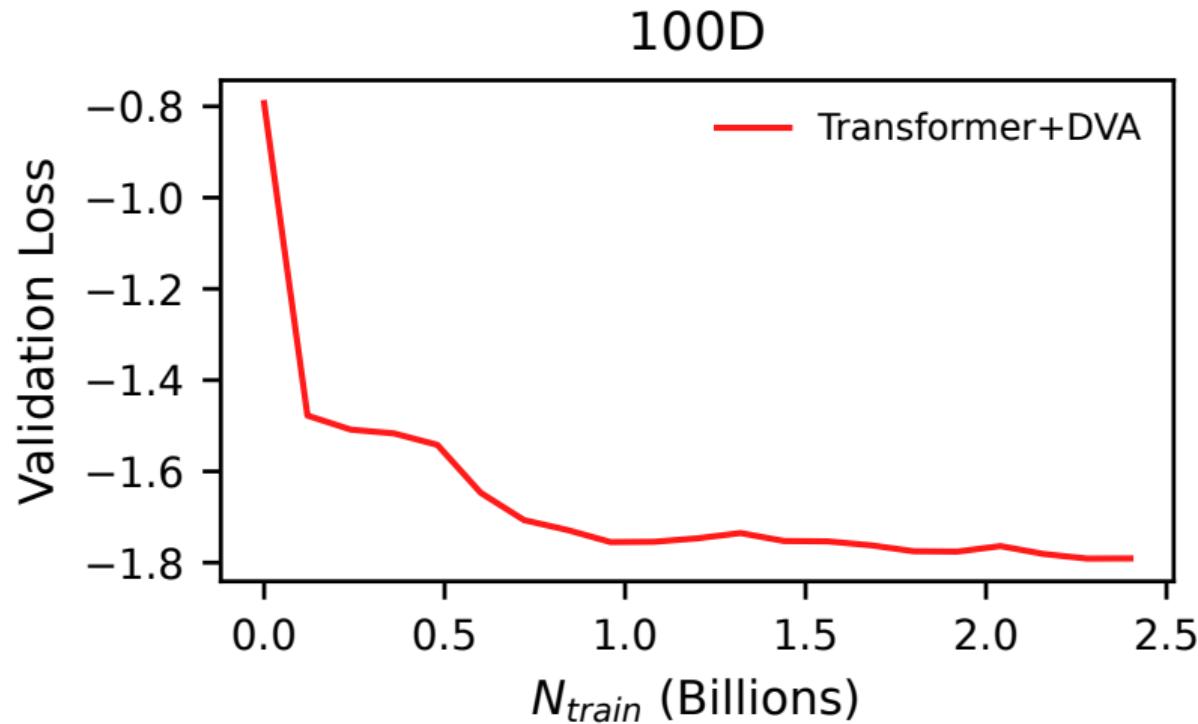


Figure: Comparing the PPD for three distinct samples for (Top) Bus 18. (Bottom) Bus 33

## 100D is Possible Too



## Summary & Takeaways

- **Decoupled-Value Attention:** A specialized attention rule that preserves input locality and mirrors GP inference.
- **Bias Reduction:** DVA cuts PFN bias by  $> 50\%$  in high-D tasks (5D, 10D) compared to vanilla attention.
- **Attention  $\gg$  Backbone:** CNN/RNN/LSTM PFNs with DVA perform on par with Transformers, despite far fewer parameters.
- **Scaling to 64D:** DVA-enabled PFNs successfully learn 64D physical equations at  $\sim 80\times$  GP speed.

### Big Question

What if we could make a PFN which works for all ND GP Regressions?  
A FOUNDATION MODEL FOR REGRESSION

# Best Part of The Work!

My Coauthors are Third Year UG Students @IITR!



**Kaustubh Sharma**

B.Tech Electrical Engg.  
kaustubh202.github.io

Lab Website



**Simardeep Singh**

B.Tech Metallurgical Engg.  
linkedin.com/in/simar7220

Full Pre-Print



Funding Support



Faculty Initiation Grant



PM Early Career Research  
Grant-2025