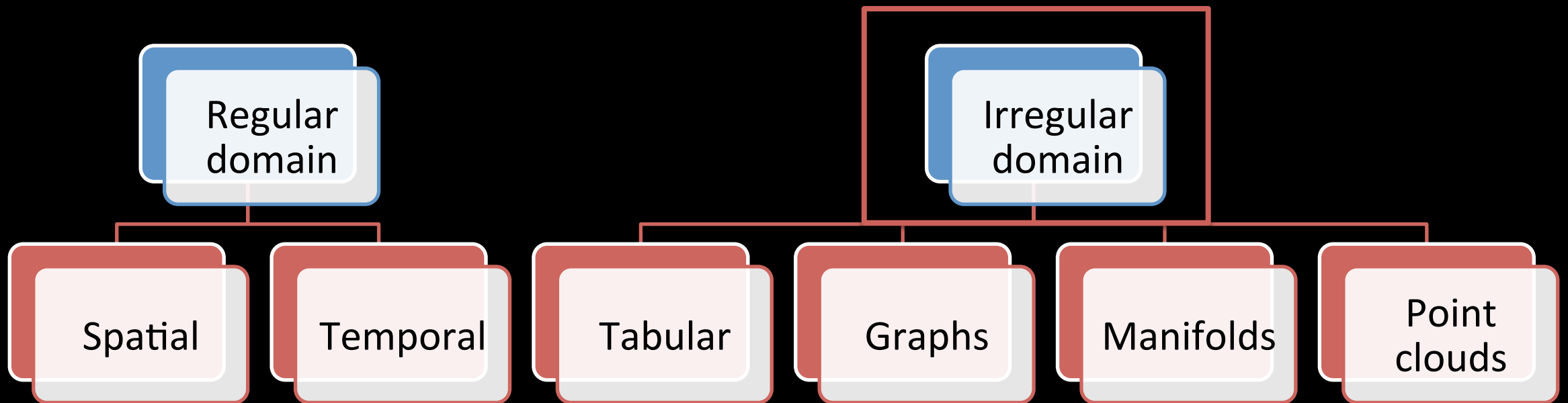# Scientific Domain-Informed Machine Learning

Prasanna Balaprakash
Computer Scientist

Mathematics and Computer Science Division and Leadership Computing Facility
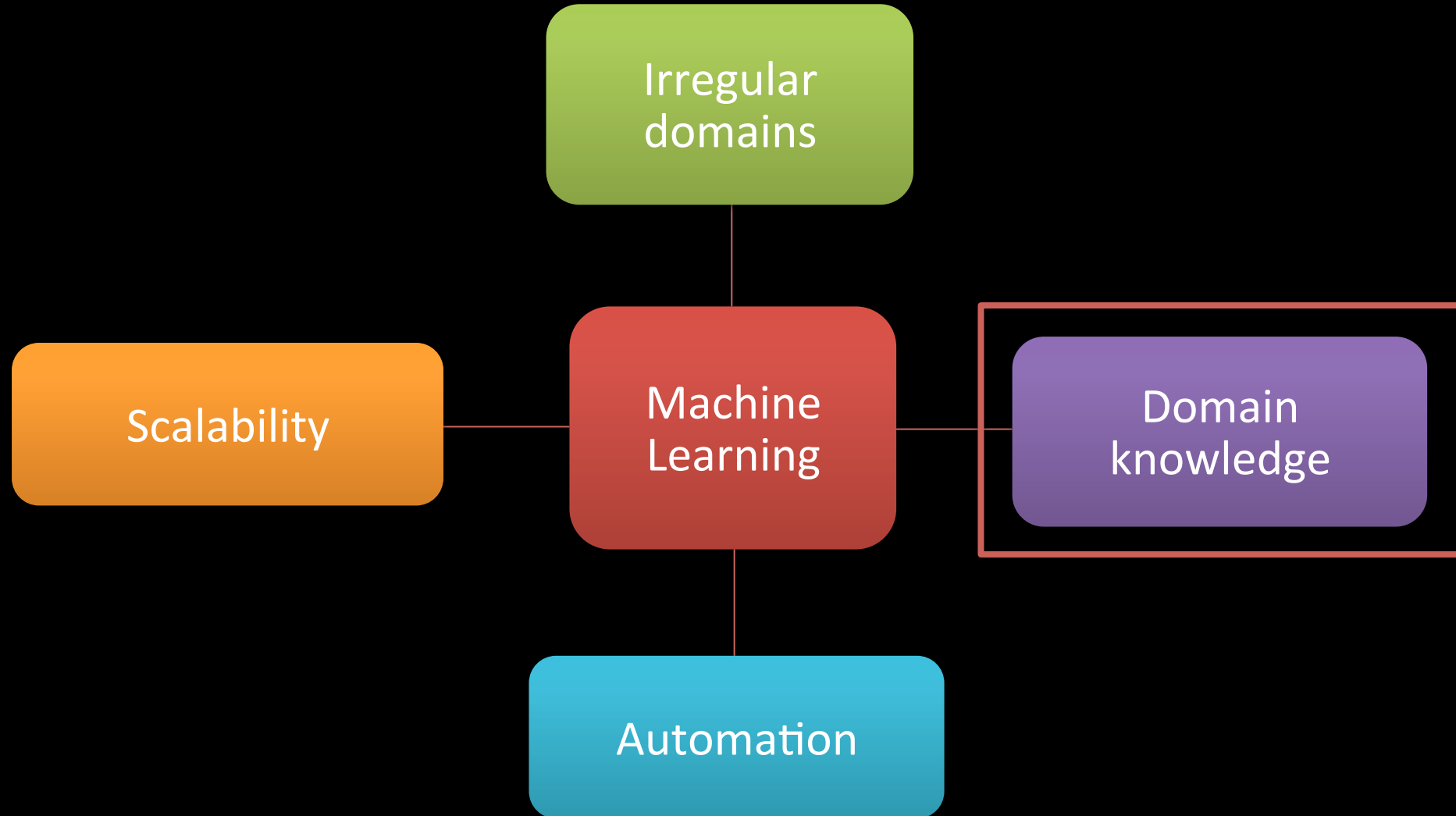Argonne National Laboratory

Joint work with P. Malakar, V. Vishwanath, K. Kumaran, V. Morozov, J. Wang, R. Kotamarthi

**ALCF Simulation, Data, and Learning Workshop**
**October 3, 2018**

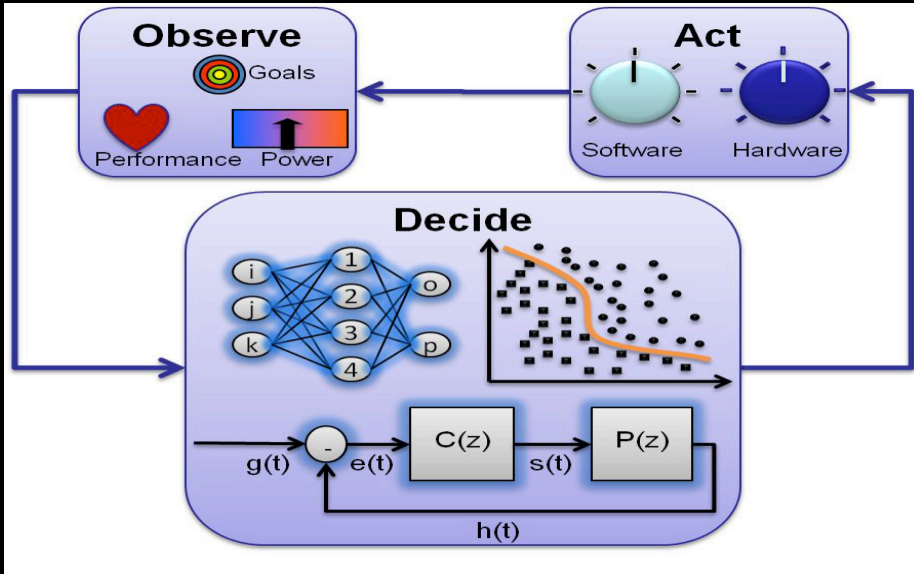# DOE scientific data for machine learning
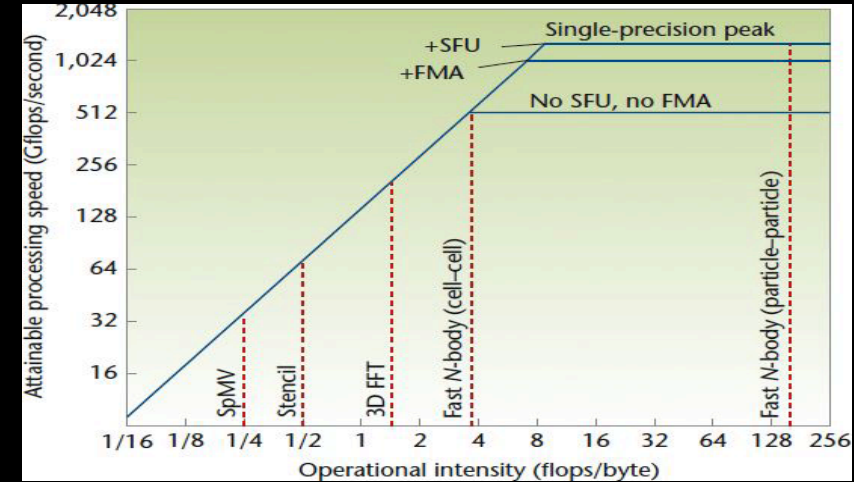
# Challenges for irregular domains

# Case studies

- Scientific application performance modeling
- Surrogate modeling in weather simulation
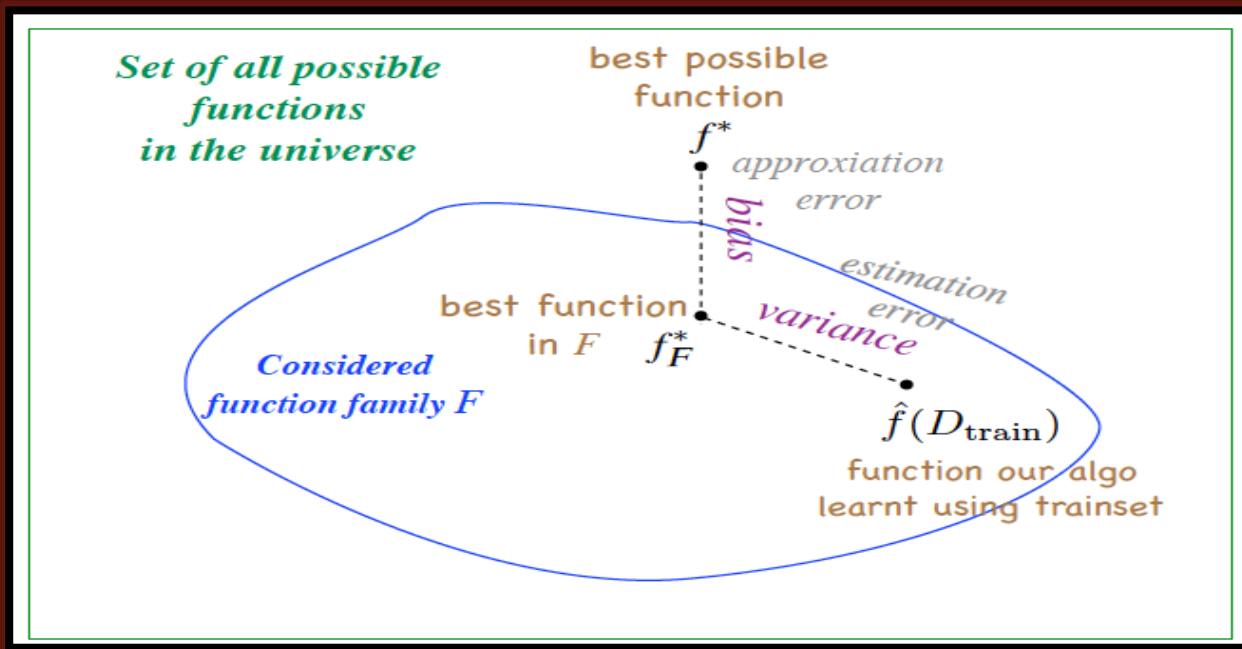
# Predictive models in HPC applications



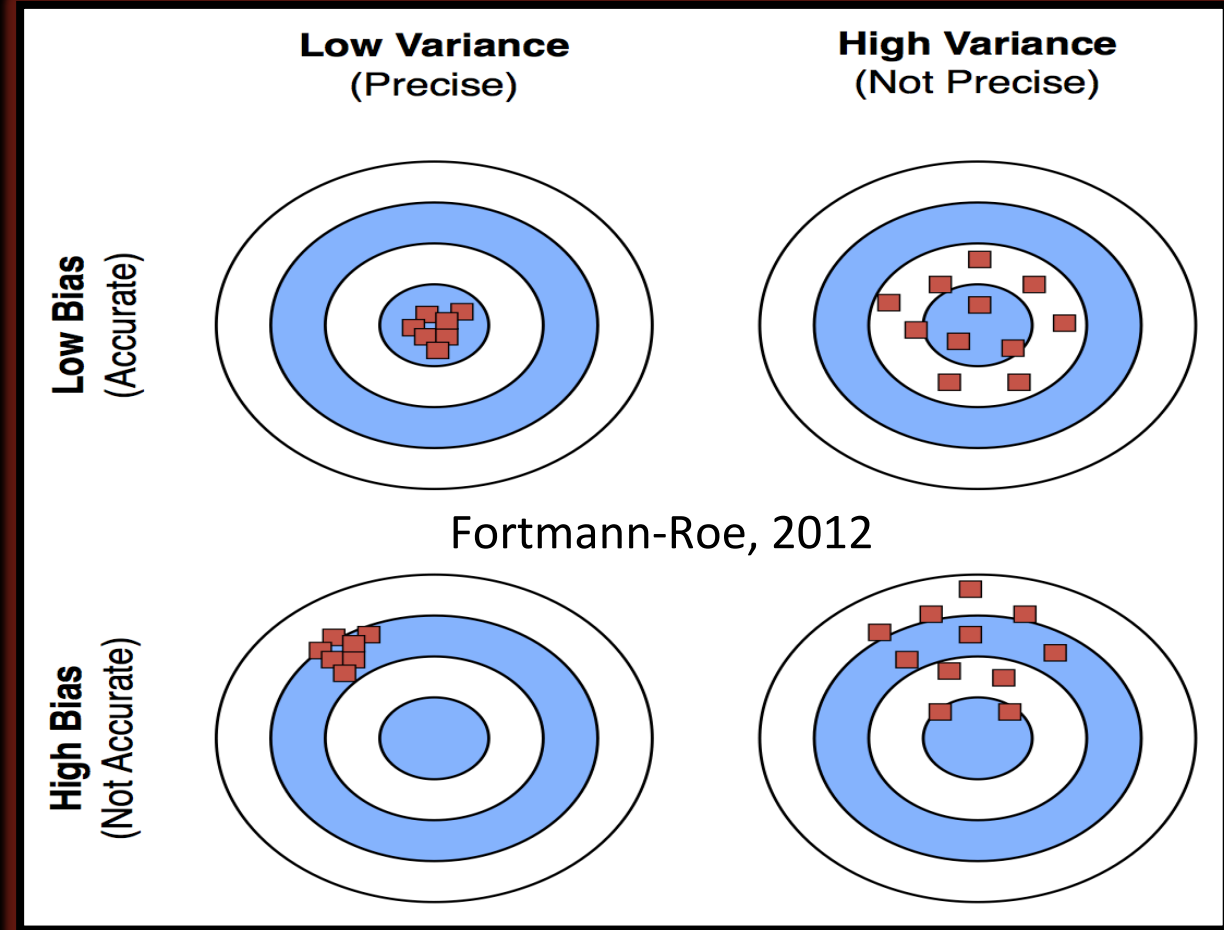[H. Hoffmann, World Changing Ideas, SA 2009]



[S. Williams et al., ACM 2009]

- Performance (run time) prediction still challenging
- ML-based performance modeling to bridge the gap
- Insights on important knobs that impacts performance
- Help prune large search spaces in performance tuning
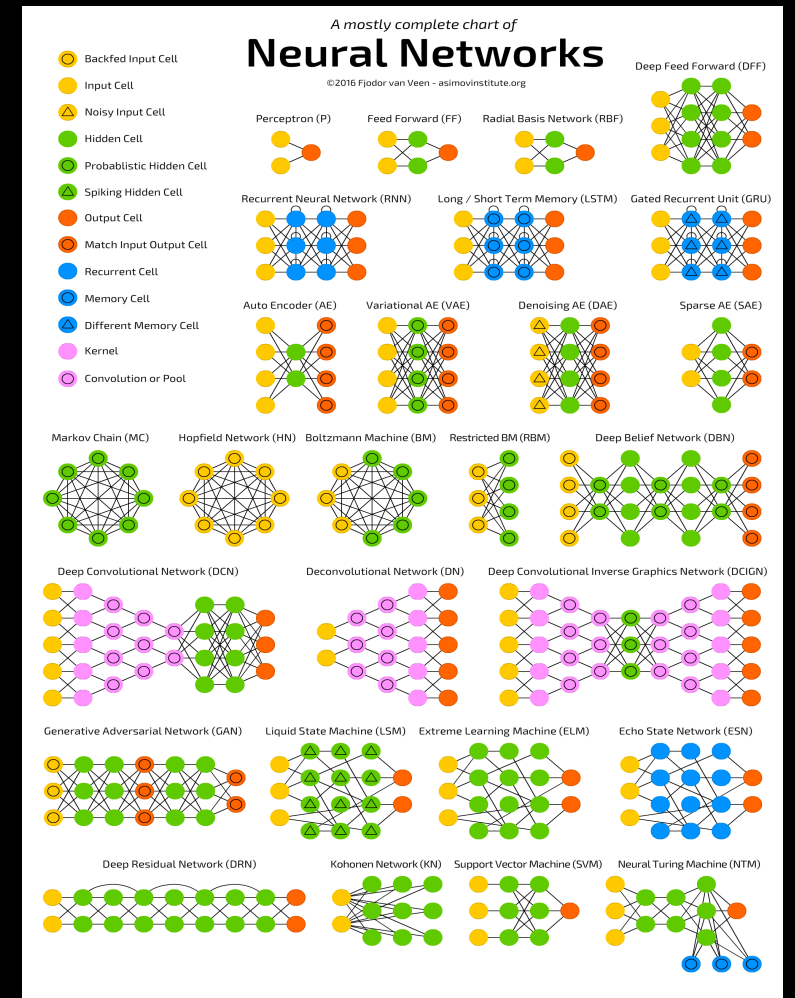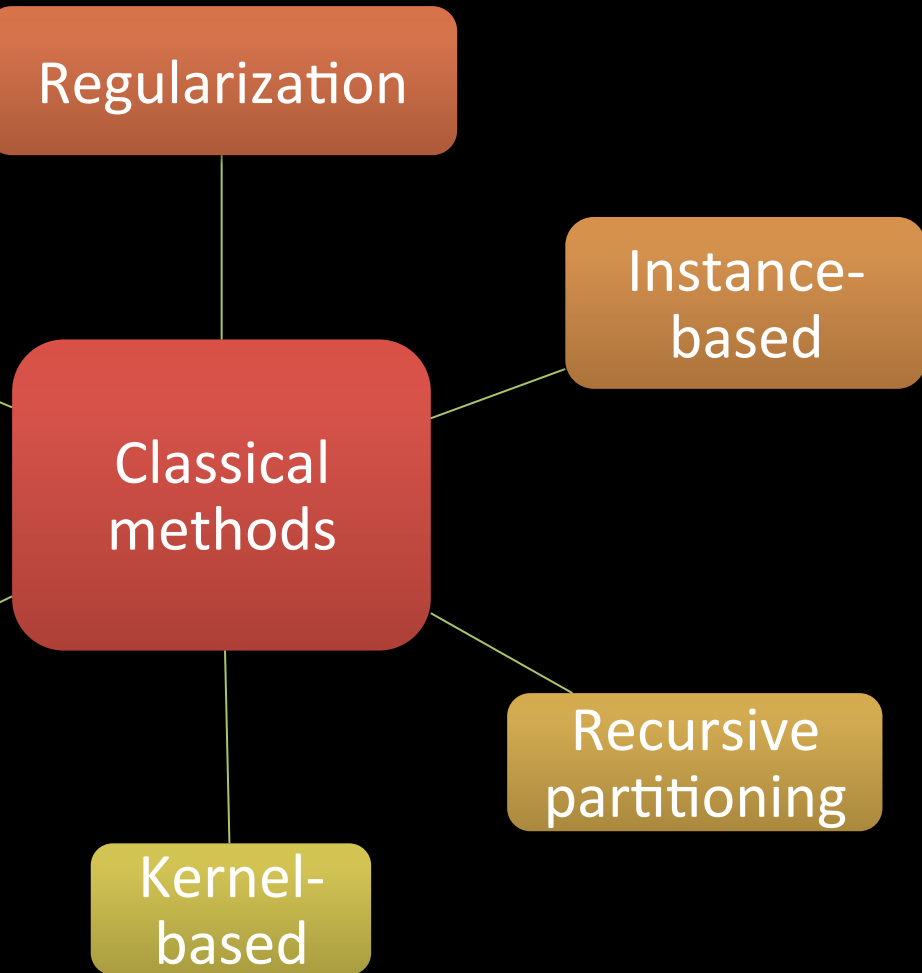
# Bias variance tradeoff



Set of all possible functions in the universe

best possible function $f^*$

approxiation error

bias

best function in $F$ $f_F^*$

estimation error

variance

Considered function family $F$

$\hat{f}(D_{train})$

function our algo learnt using trainset

Deep learning summer school lecture, CIFAR,  2016



**Low Variance** (Precise)

**High Variance** (Not Precise)

**Low Bias** (Accurate)

**High Bias** (Not Accurate)
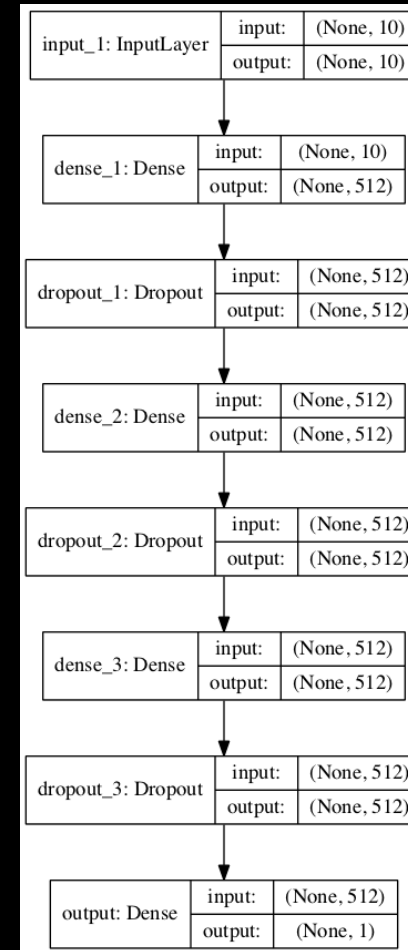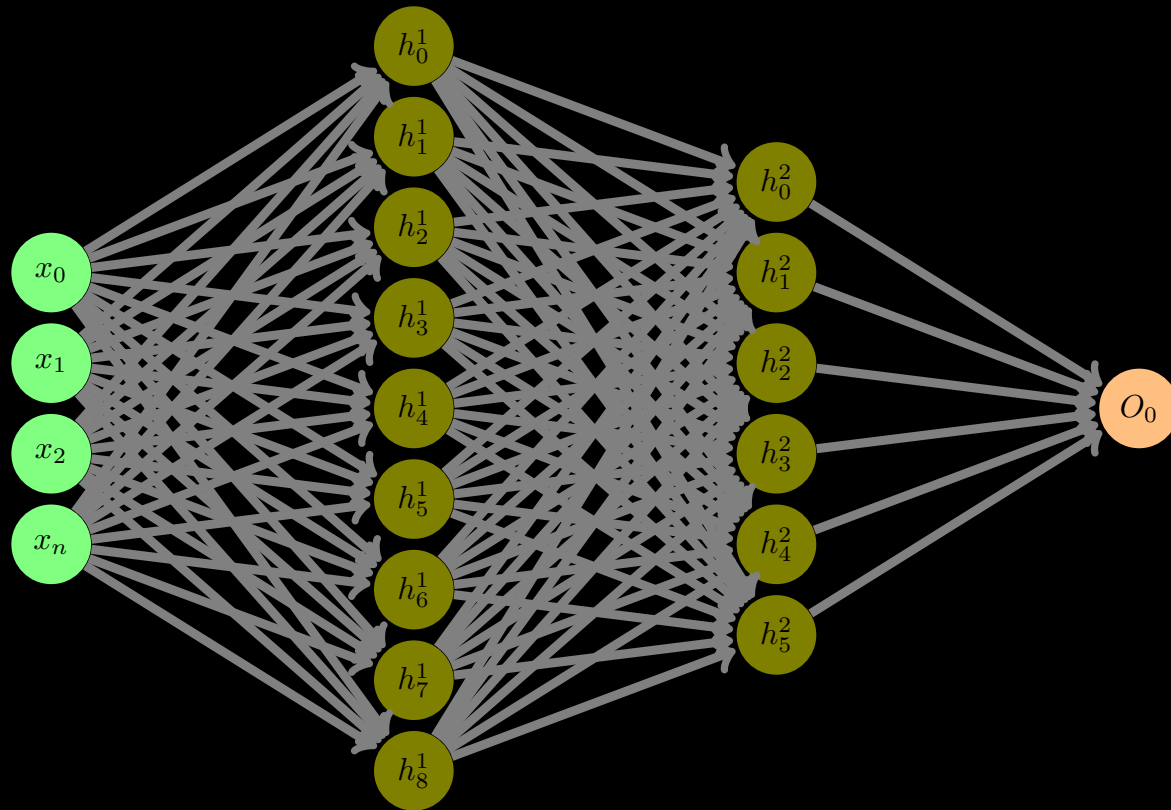
Fortmann-Roe, 2012

- ***No free lunch***: no single method will work well on all data set
- All supervised learning algorithms ***seek to reduce bias and variance*** in a different way

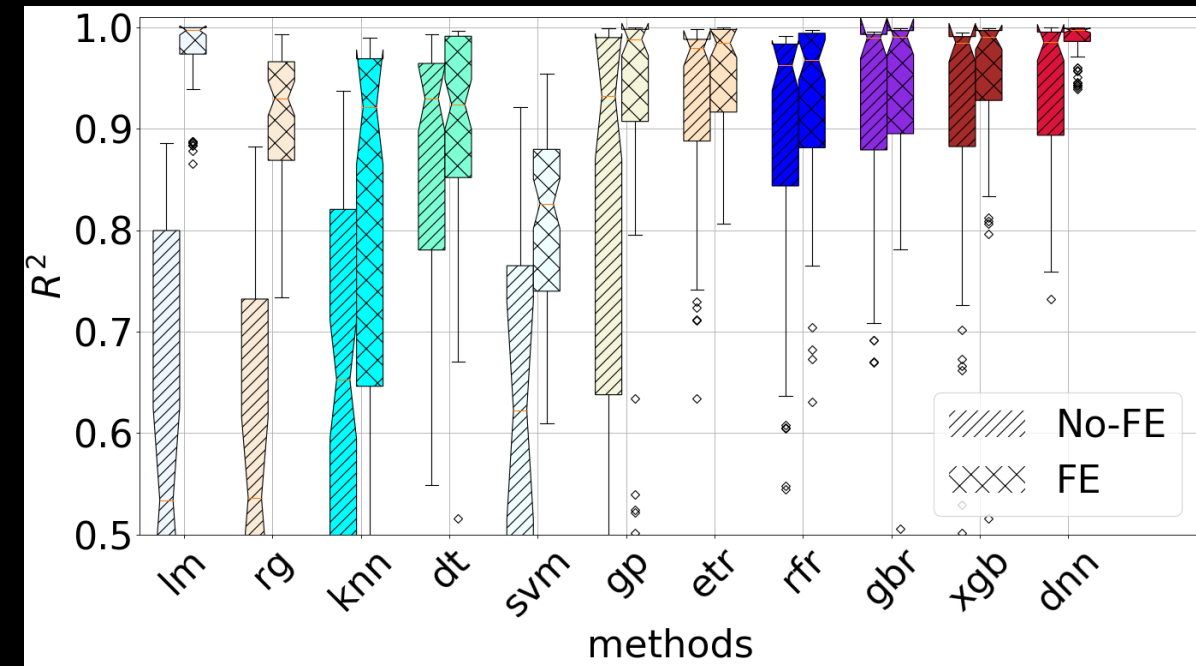# Supervised learning methods

# Deep neural networks

# Applications and platforms

| Name | Processor | Interconnect topology | Maximum # cores |
|------|-----------|----------------------|-----------------|
| Mira (Blue Gene/Q) | Power BQC 1.6 GHz | 5D torus | 131072 |
| Vesta (Blue Gene/Q) | Power BQC 1.6 GHz | 5D torus | 16384 |
| Edison (Cray XC30) | Intel Ivy Bridge 2.4 GHz | Aries with dragon-fly | 1728 |
| Hopper (Cray XE6) | AMD MagnyCours 2.1 GHz | Gemini with 3D torus | 12000 |

- Miniapps (# no of data points):
  - miniMD (< 2K); O(1024) nodes
  - miniAMR (< 1K); O(4096) nodes
  - miniFE (6K to 15K); O(8192) node
  - LAMMPS (< 1K ); O(1024) nodes

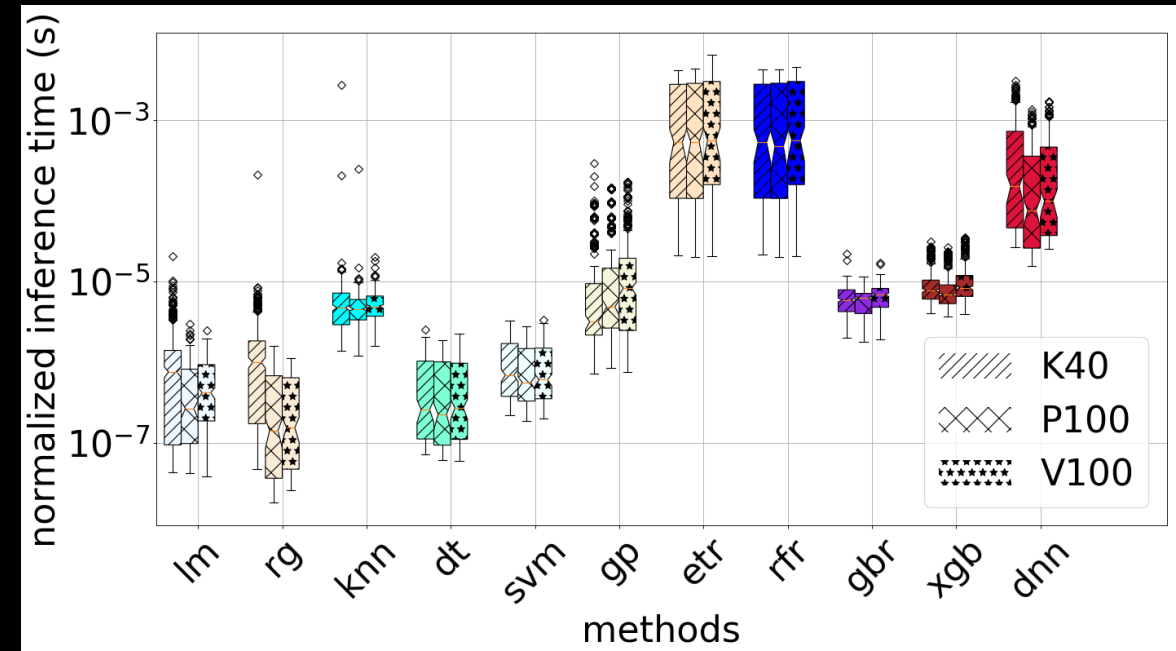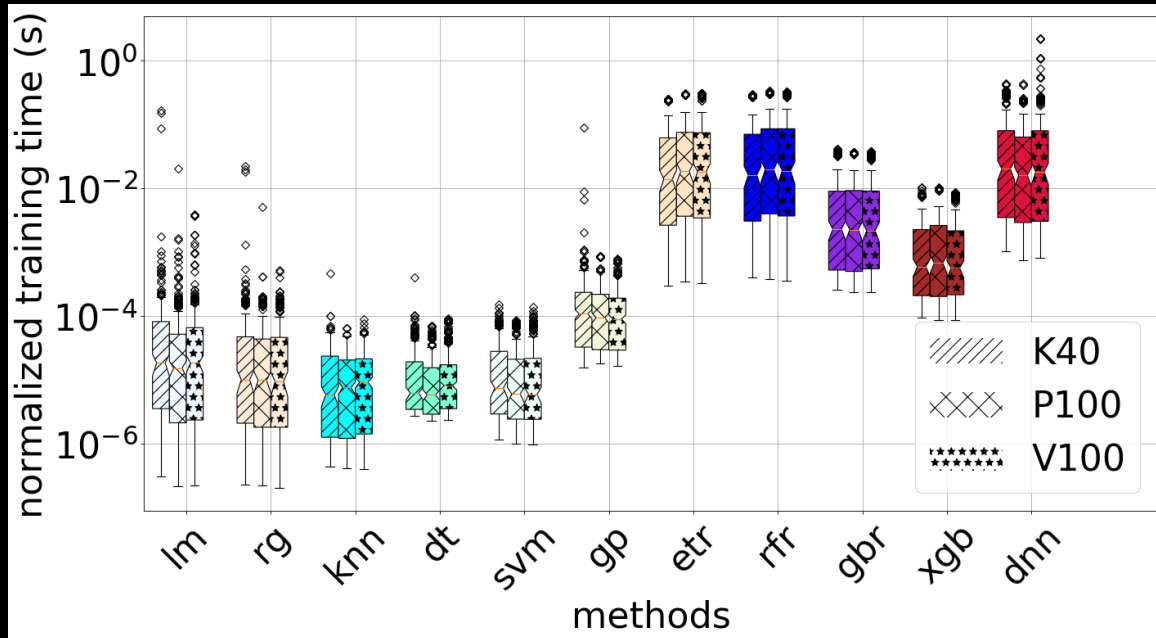# Impact of feature engineering

- No Feature Engineering (No-FE)
  - application input parameters
- Feature Engineering (FE)
  - application input parameters
  - ratio of the application problem size and the number of number of processes
  - inverse of the number of processes
  - binary logarithm of number of processes
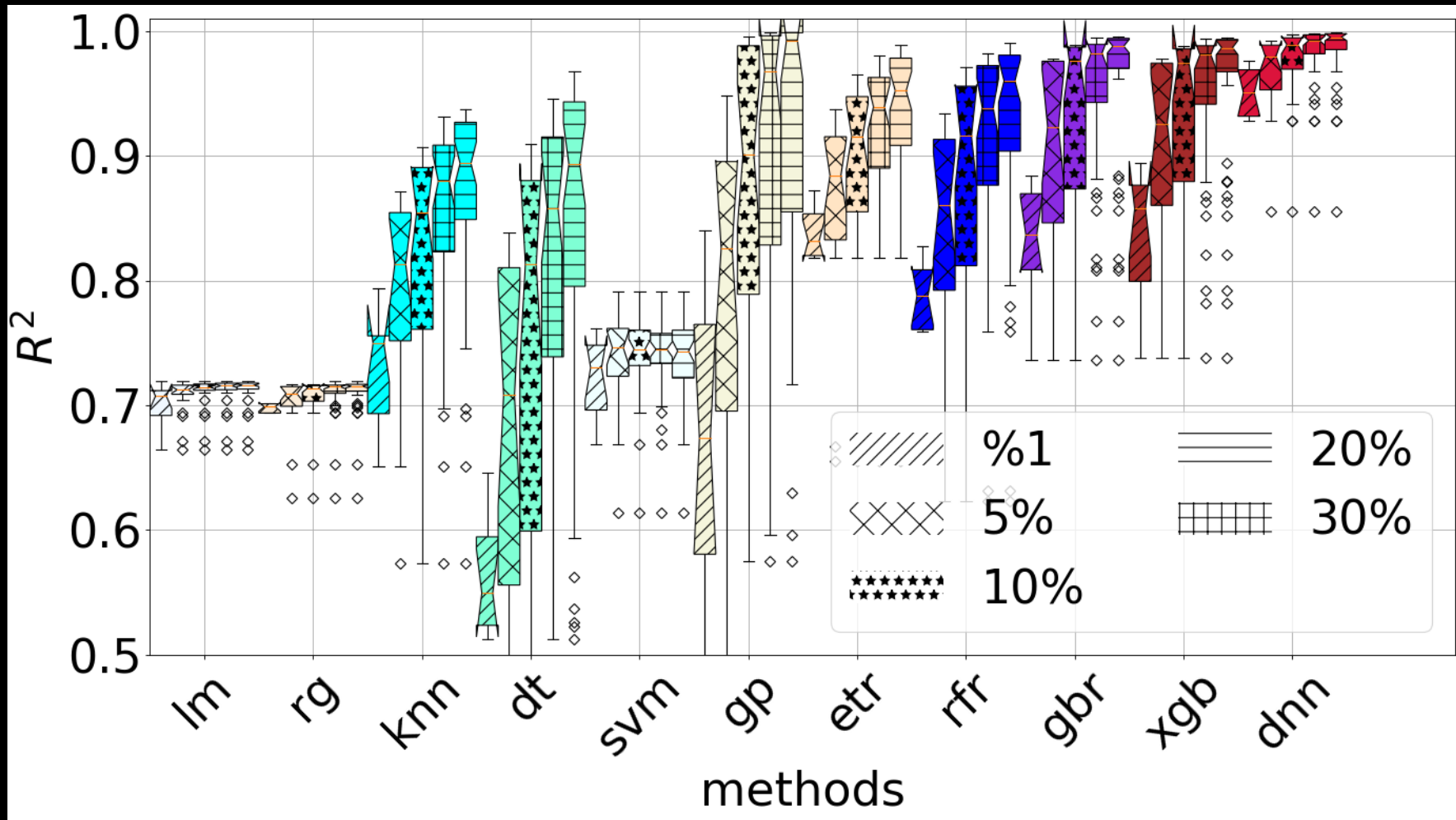


10 X 20:80 cross validation

Feature engineering has a significant impact on the accuracy
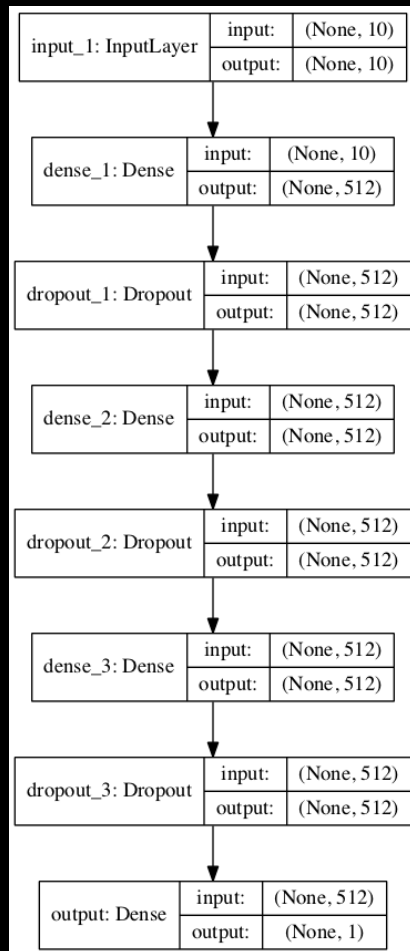
# Impact of hardware platforms



Algorithmic complexity has more impact than hardware platforms
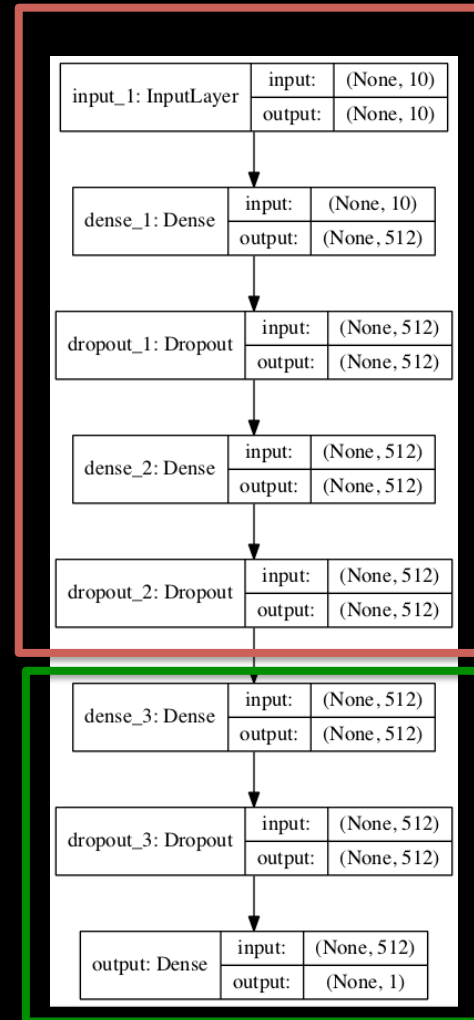
# Impact of training data size on accuracy



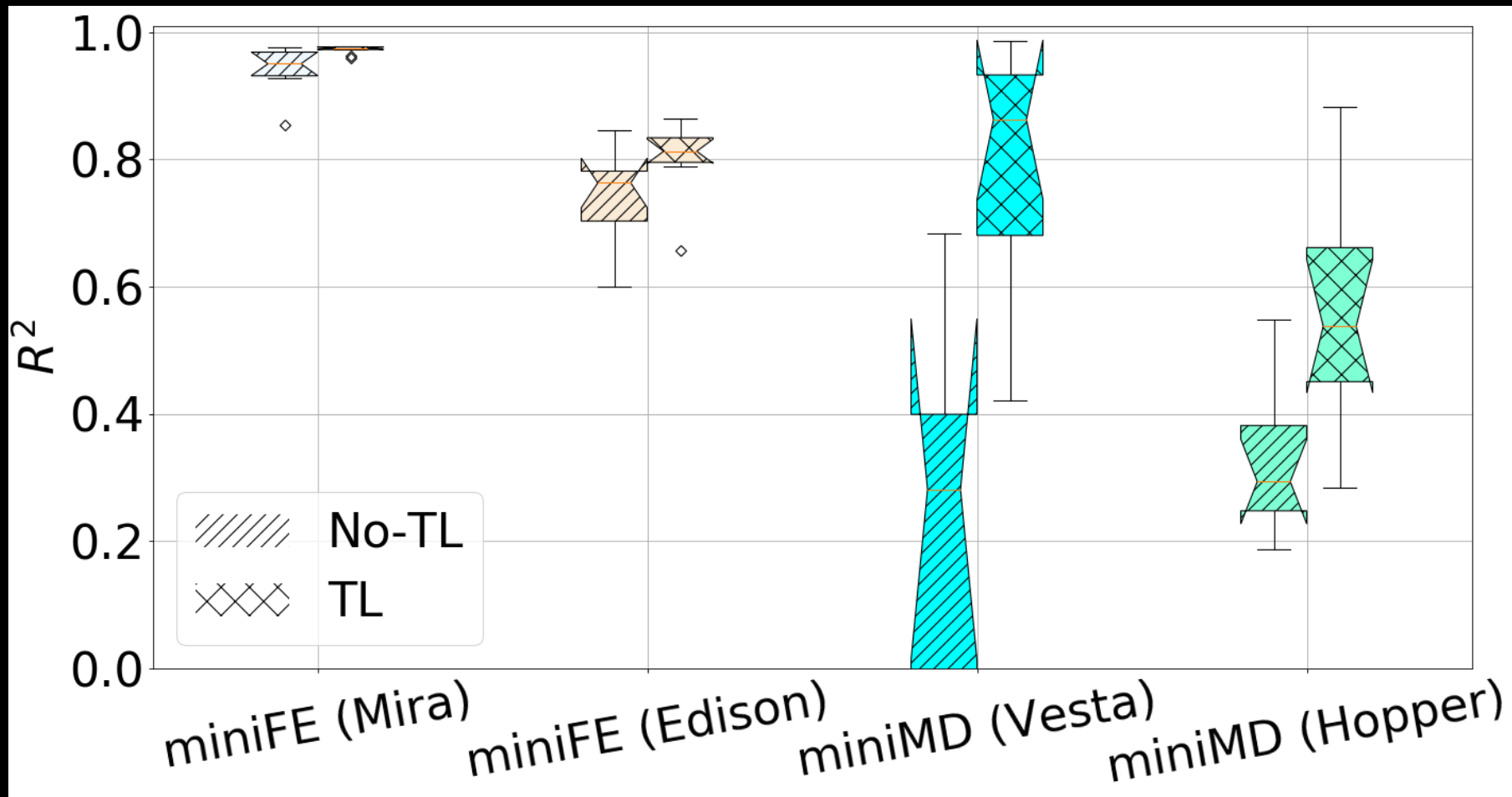Nonlinear methods leverage large training data size

# Transfer learning



Freeze weights
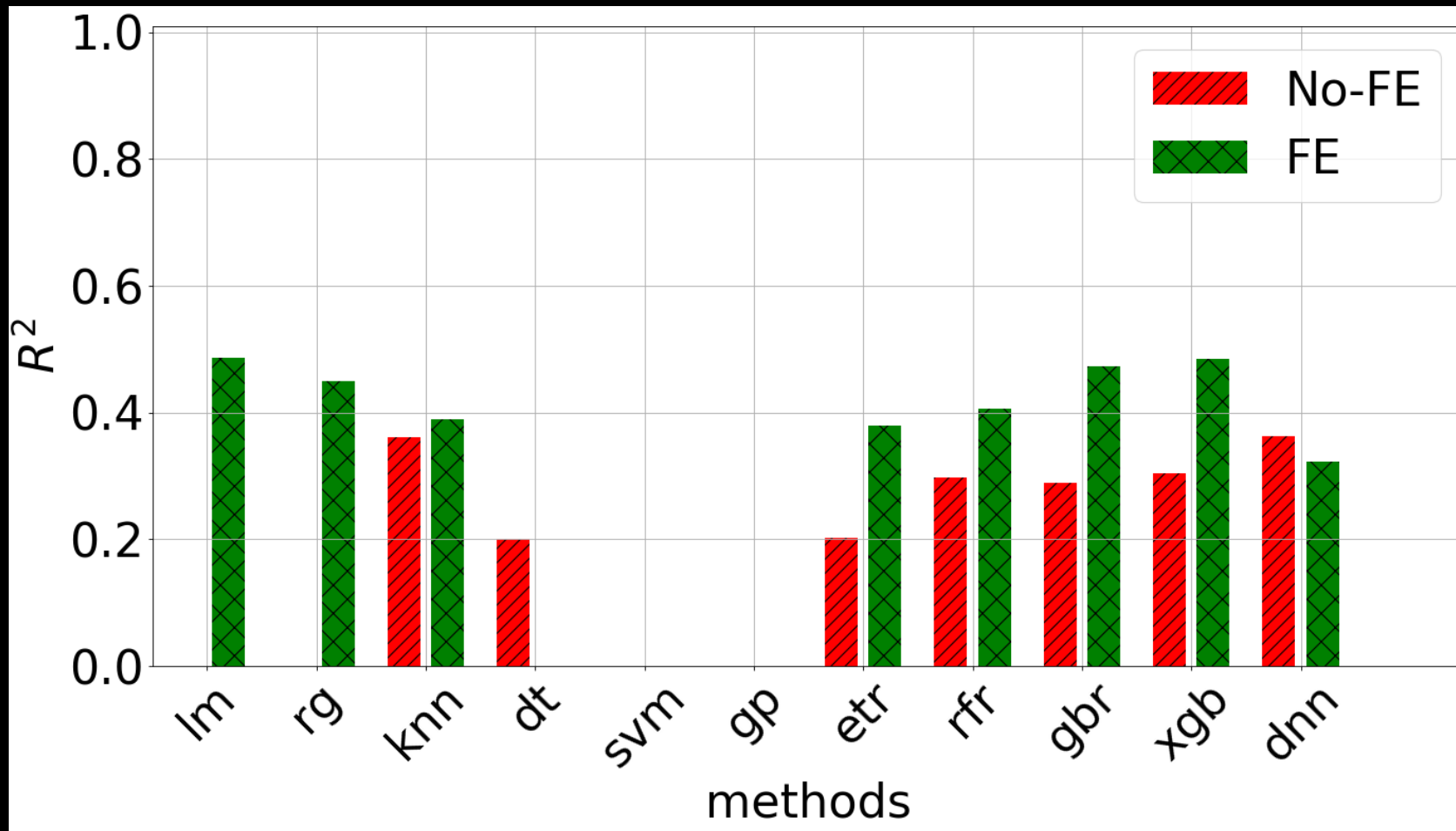
Retrain weights

Platform 1

Platform 2

# Transfer learning



Transfer learning significantly improves prediction accuracy
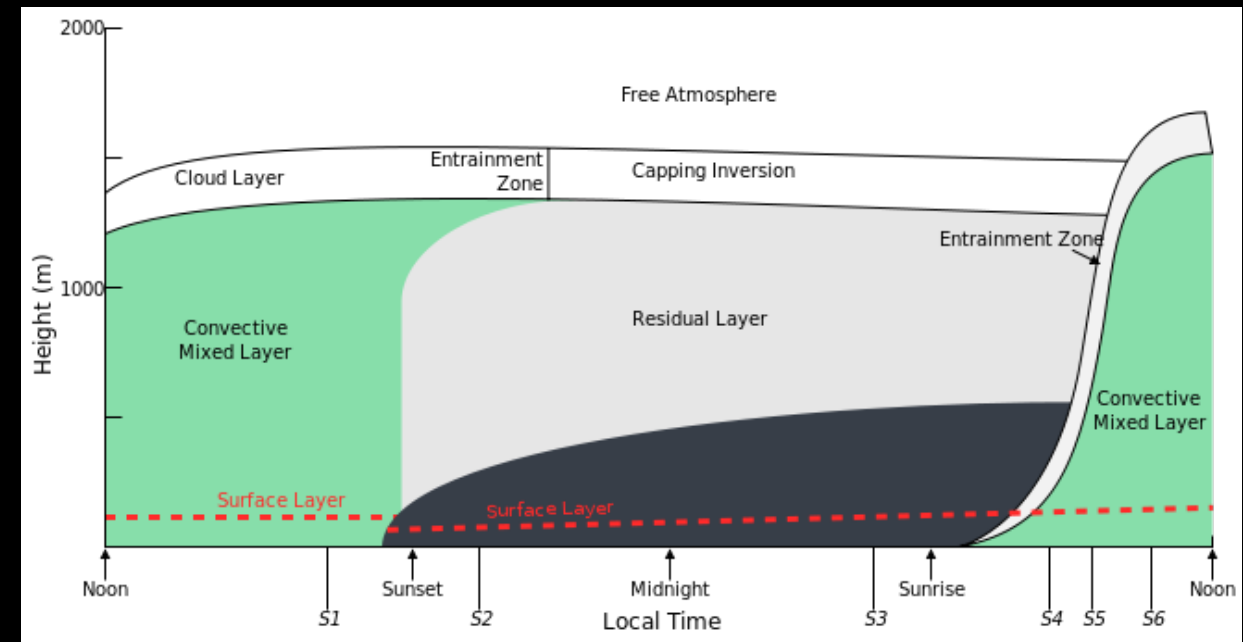
# Extrapolation
# from small to large problem sizes



Incorporating domain knowledge helps in exploration

# Case studies

- Scientific application performance modeling
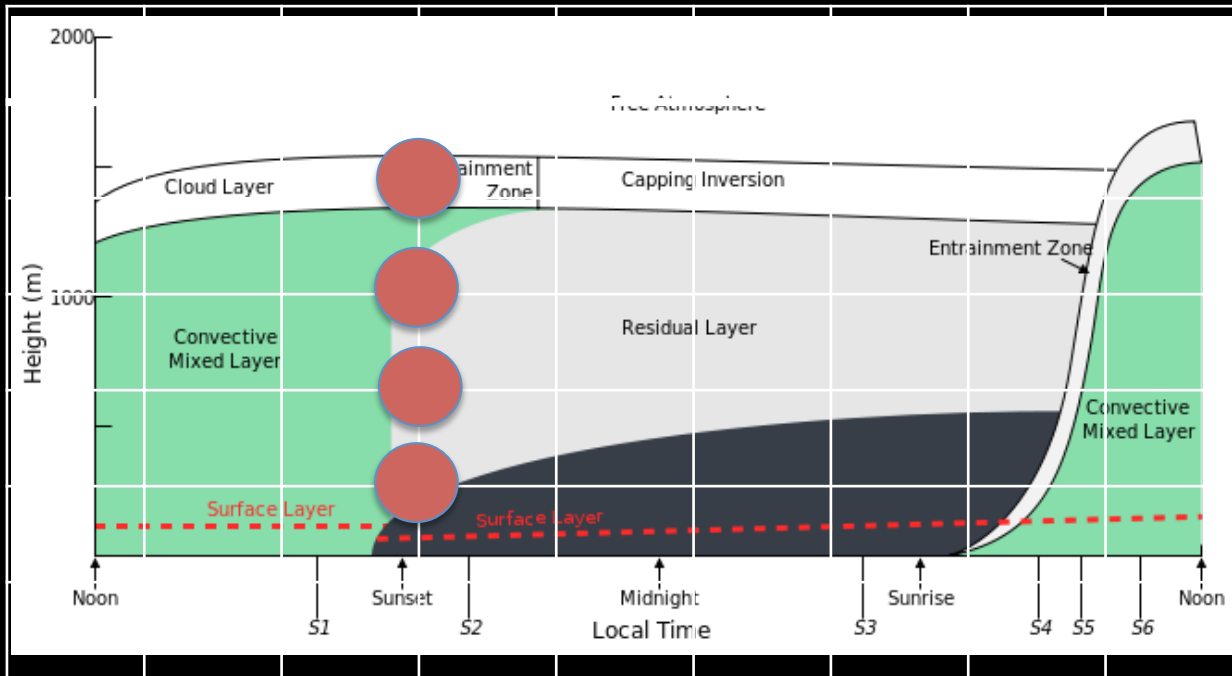- Surrogate modeling in weather simulation

# Planetary boundary layer

- Planetary boundary layer (PBL)
  – lowest part of the atmosphere
  – directly influenced by its contact with a planetary surface
  – responds to changes in surface radiative forcing
  – flow velocity, temperature, moisture, etc., display rapid fluctuations
  – computationally expensive in weather research and forecasting model



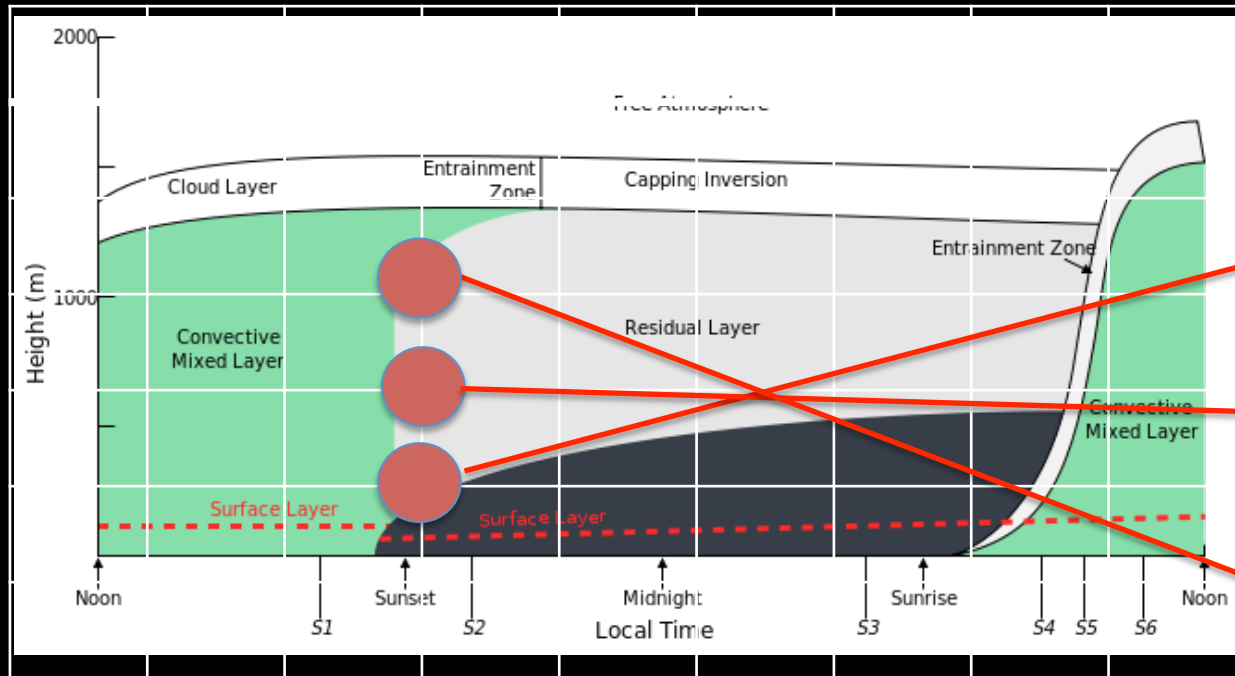https://en.wikipedia.org/wiki/Planetary_boundary_layer

# Planetary boundary layer
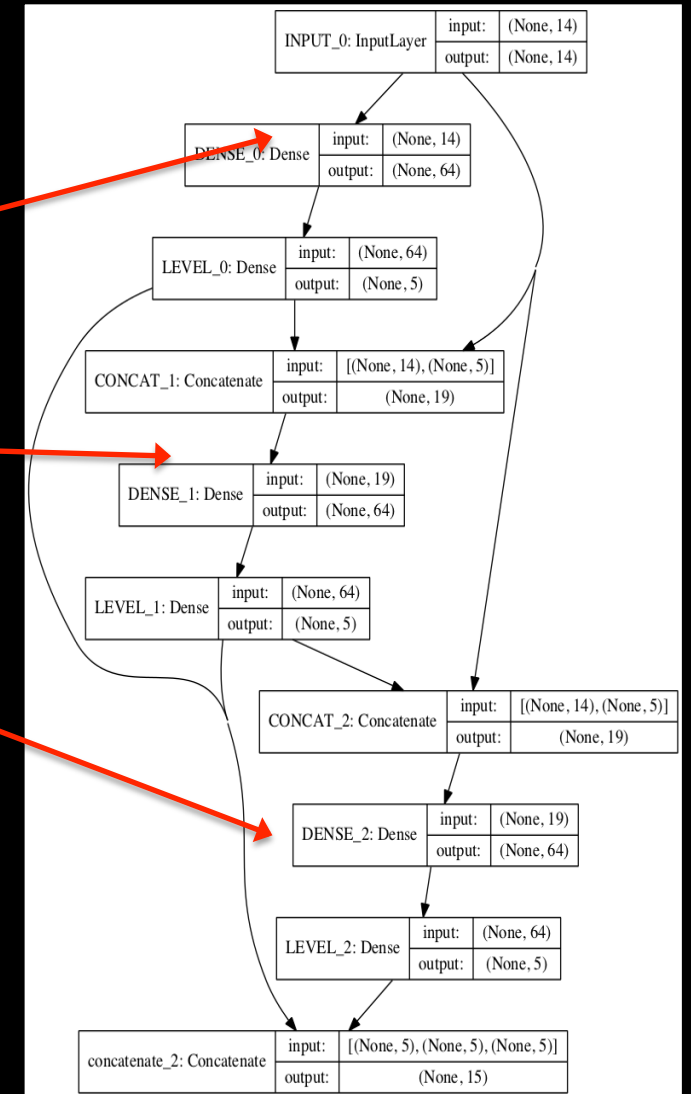


**Output:** profiles of wind, temperature, moisture

**Input:** Surface properties, fluxes, and ground temperature
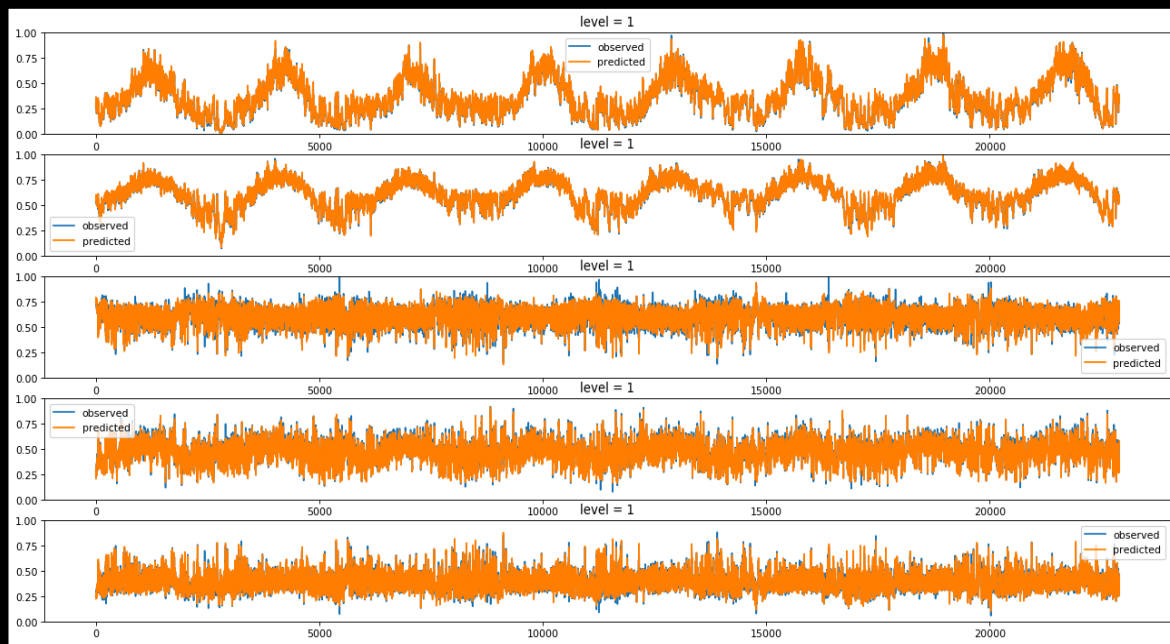
# Surrogate neural network



**Input:** Surface properties, fluxes, and ground temperature
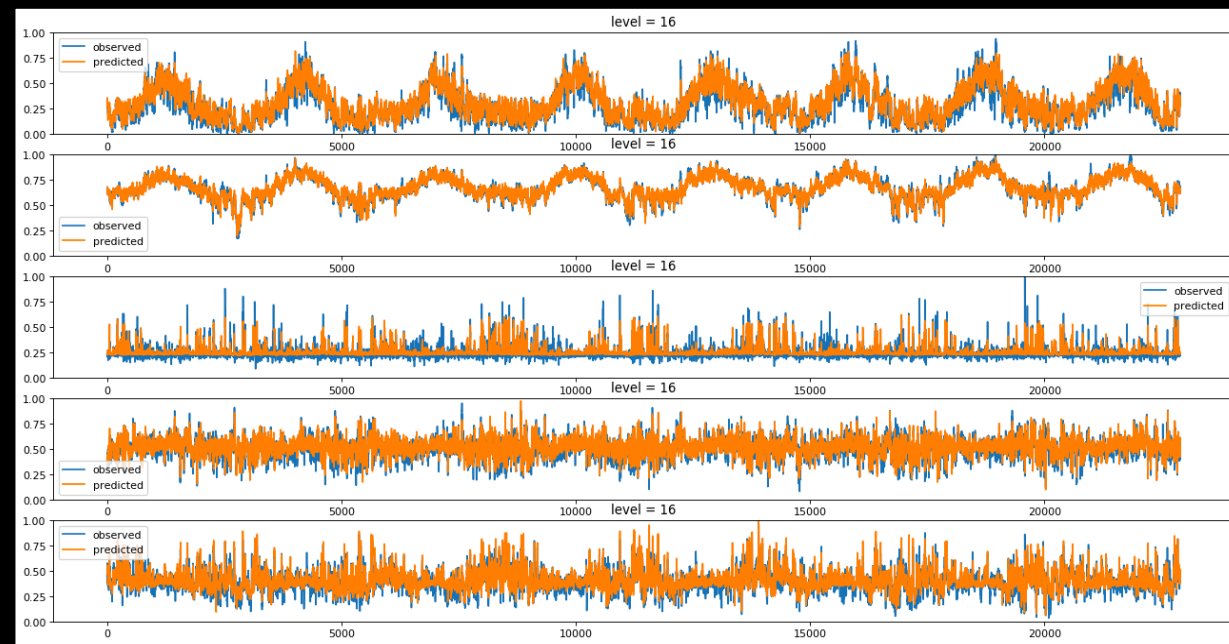
**Output:** profiles of wind, temperature, moisture

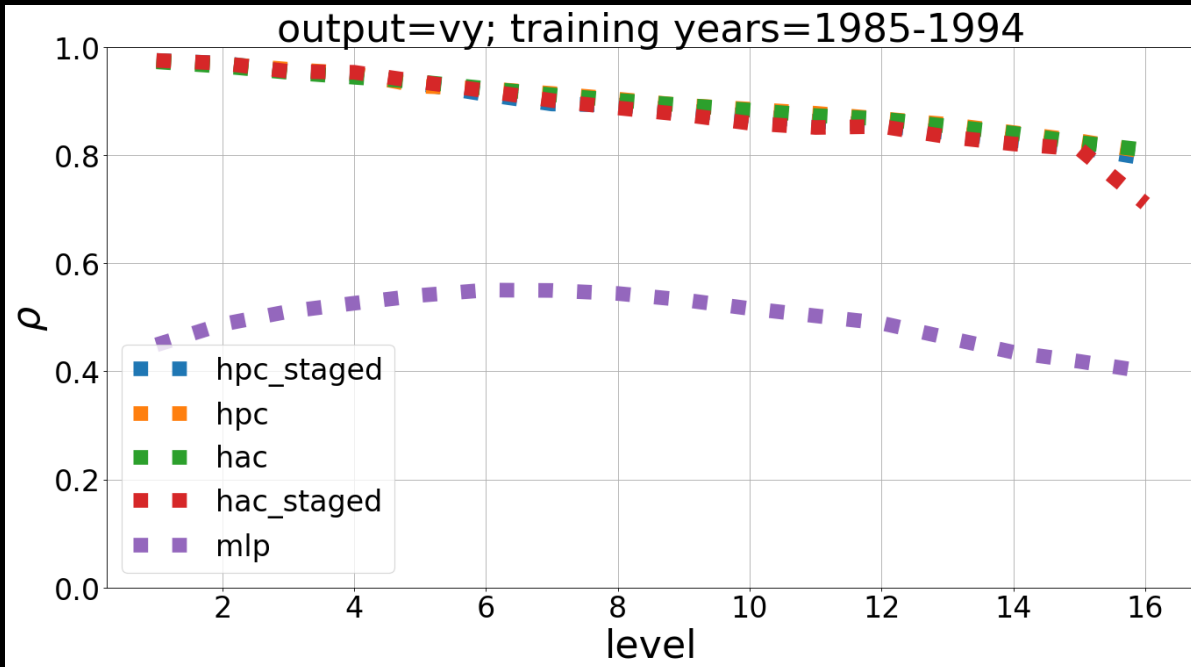# Prediction results
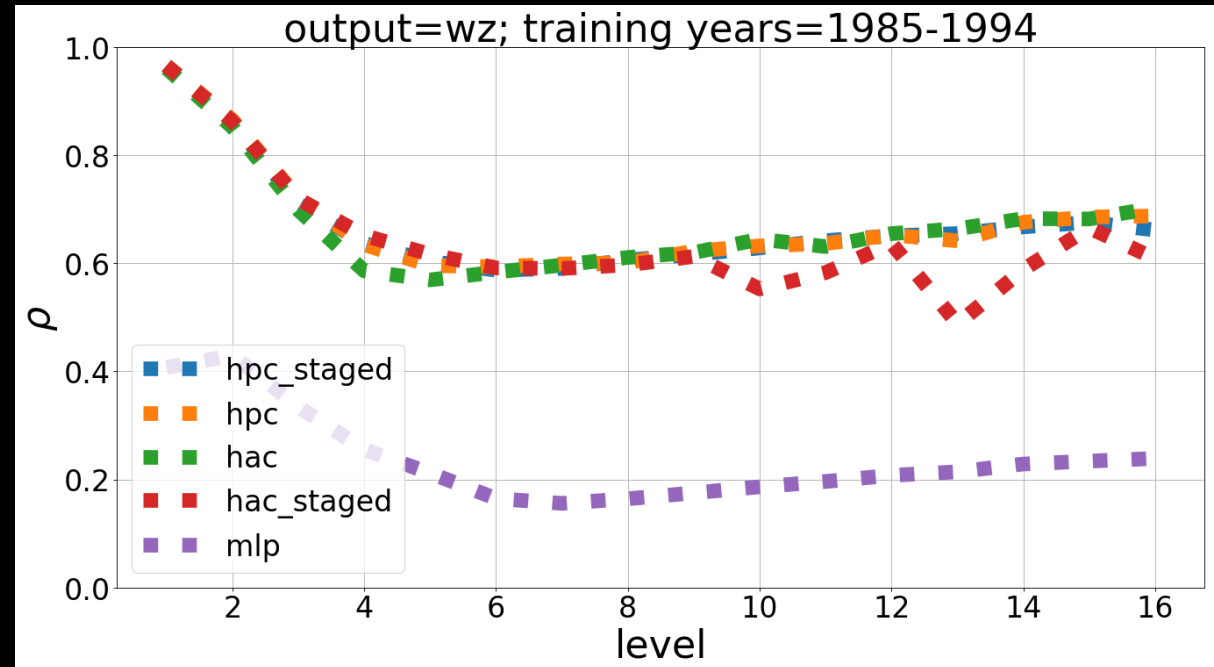
**(level 1: close to surface)**

**(level 16: ~2 km)**

# Prediction results



meridional wind (south-north wind)

vertical velocity (up and down motions)

# Prediction results

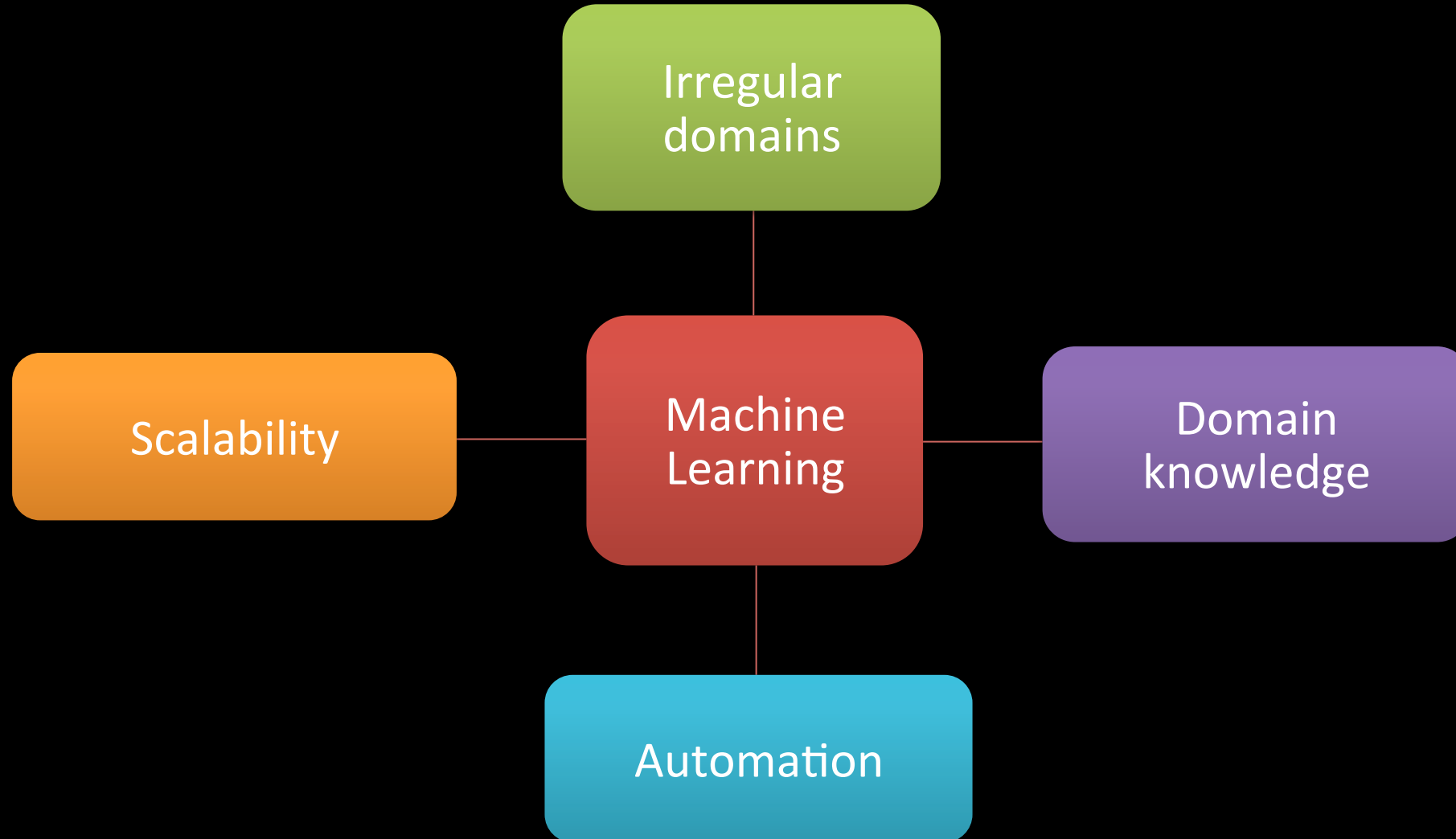# Challenges for irregular domains

# Acknowledgements

- U.S. DOE
  - ANL LDRD
  - ALCF
  - ASCR, Early Career Research Program