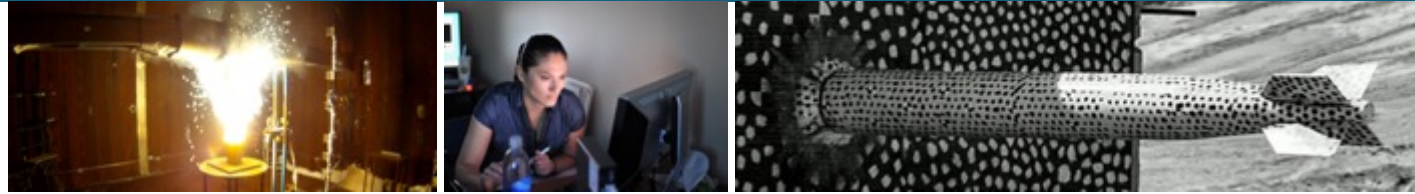


Training and Generalization of Residual Neural Networks as Discrete Analogues of Neural ODEs



Khachik Sargsyan (SNL-CA)

SNL-CA : Joshua Hudson, Oscar Diaz-Ibarra, Marta D'Elia, Habib Najm
Emory U : Lars Ruthotto, Haley Rosso

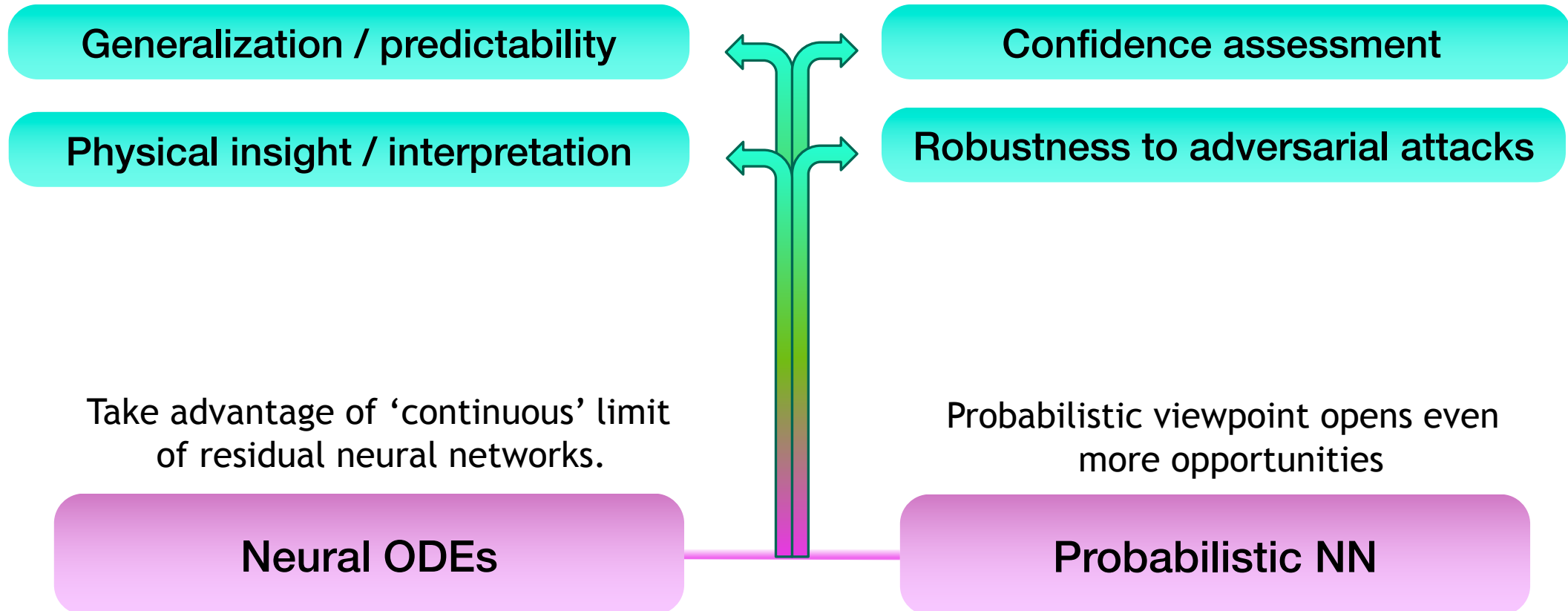
July 26, 2022

MLDL Workshop

SNL LDRD Project: Analysis of Neural Networks as Random Dynamical Systems



Despite all the success, there are many recognized challenges and unknowns in neural network behavior





Arguably, the two most important hurdles along the way

Goals:

Generalization / predictability

Confidence assessment

Tools:

Neural ODEs / ResNets

Probabilistic NN

Take advantage of legacy knowledge in ODEs and UQ to achieve

- Improved architectures
- Generalizable models
- Confidence assessment
- Robustness to noise

Main building block: ResNets



Neural Networks (NNs) layer-to-layer function

$$h_{t+1} = F(h_t, \theta)$$

state weights

Main building block: ResNets



Neural Networks (NNs) layer-to-layer function

$$h_{t+1} = F(h_t, \theta)$$

state weights

Residual NN: learn the residual, not the state

$$h_{t+1} = h_t + F(h_t, \theta)$$

Main building block: ResNets



Neural Networks (NNs) layer-to-layer function

$$h_{t+1} = F(h_t, \theta)$$

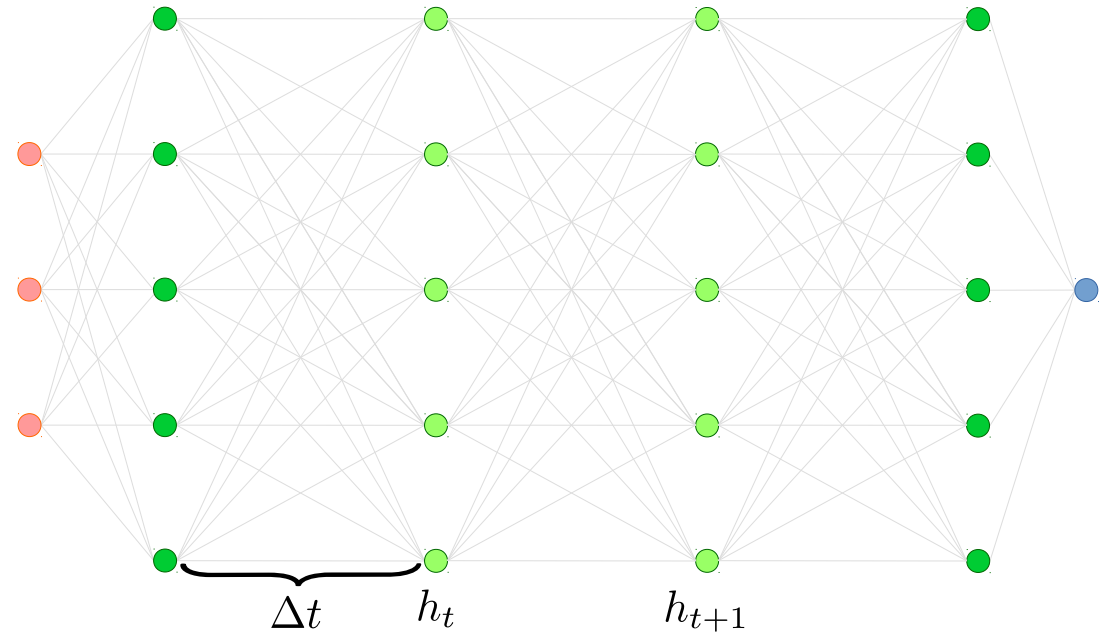
state weights

Residual NN: learn the residual, not the state

$$h_{t+1} = h_t + F(h_t, \theta)$$

Now, take the limit of infinite layers

$$\frac{dh(t)}{dt} = F(h(t), \theta)$$



Main building block: ResNets



Neural Networks (NNs) layer-to-layer function

$$h_{t+1} = F(h_t, \theta)$$

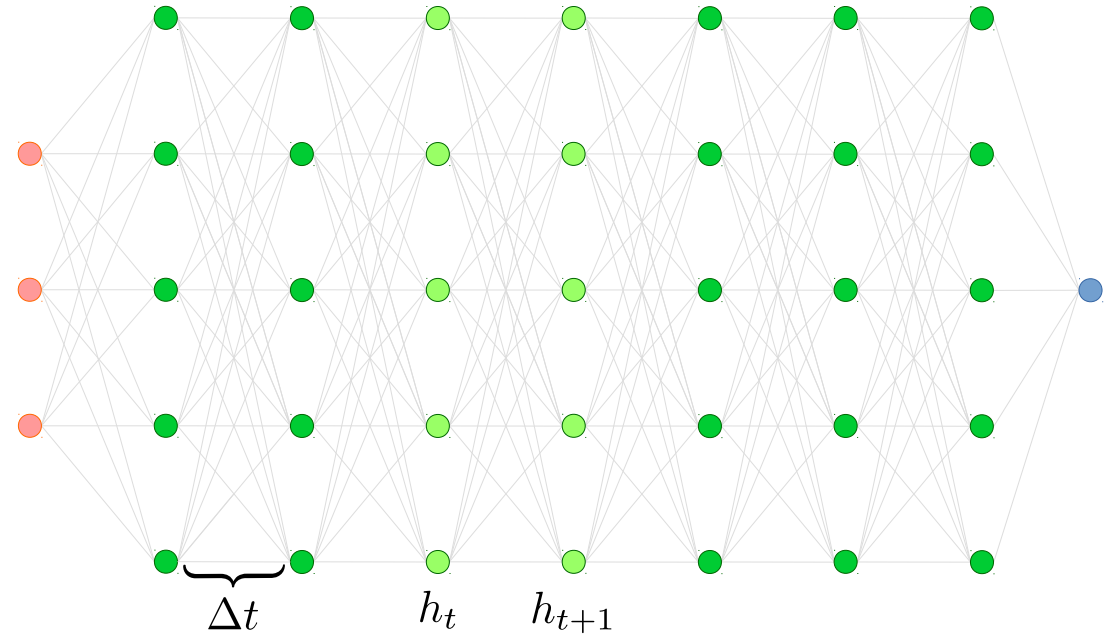
state weights

Residual NN: learn the residual, not the state

$$h_{t+1} = h_t + F(h_t, \theta)$$

Now, take the limit of infinite layers

$$\frac{dh(t)}{dt} = F(h(t), \theta)$$



Main building block: ResNets



Neural Networks (NNs) layer-to-layer function

$$h_{t+1} = F(h_t, \theta)$$

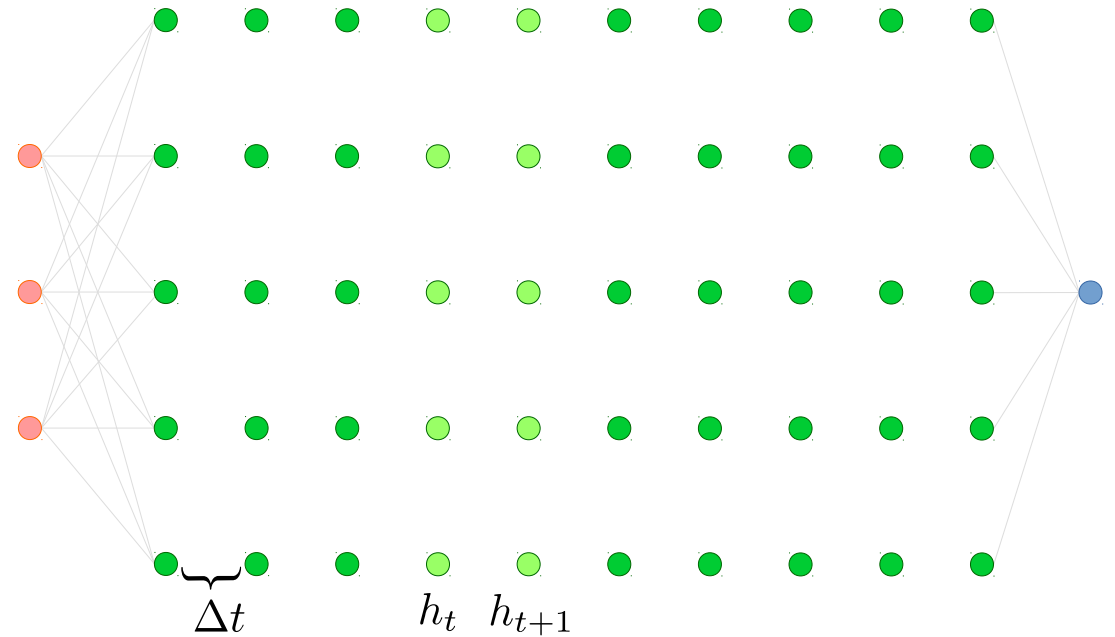
state weights

Residual NN: learn the residual, not the state

$$h_{t+1} = h_t + F(h_t, \theta)$$

Now, take the limit of infinite layers

$$\frac{dh(t)}{dt} = F(h(t), \theta)$$



Main building block: ResNets



Neural Networks (NNs) layer-to-layer function

$$h_{t+1} = F(h_t, \theta)$$

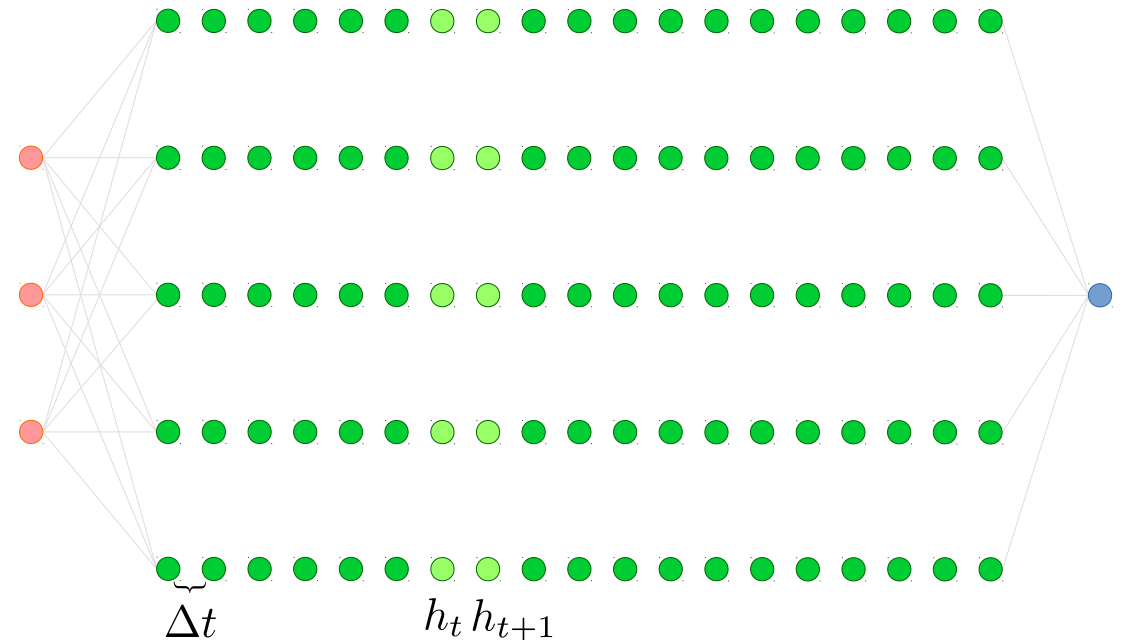
state weights

Residual NN: learn the residual, not the state

$$h_{t+1} = h_t + F(h_t, \theta)$$

Now, take the limit of infinite layers

$$\frac{dh(t)}{dt} = F(h(t), \theta)$$



Main building block: ResNets



Neural Networks (NNs) layer-to-layer function

$$h_{t+1} = F(h_t, \theta)$$

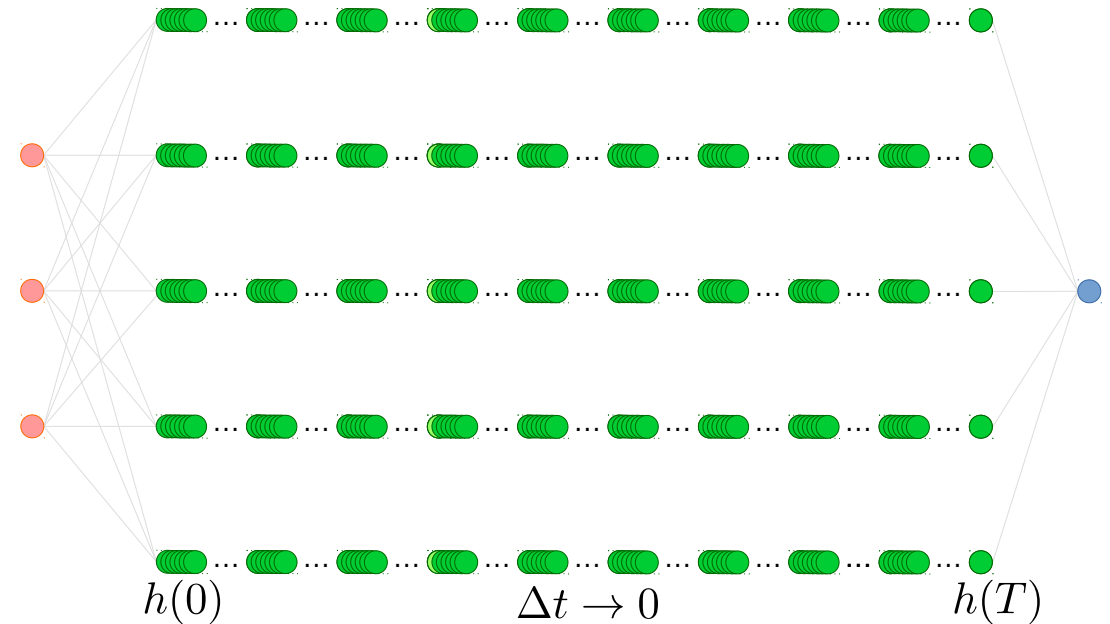
state weights

Residual NN: learn the residual, not the state

$$h_{t+1} = h_t + F(h_t, \theta)$$

Now, take the limit of infinite layers

$$\frac{dh(t)}{dt} = F(h(t), \theta)$$



Focus on: ResNet and NODE in a regression setting (supervised ML)



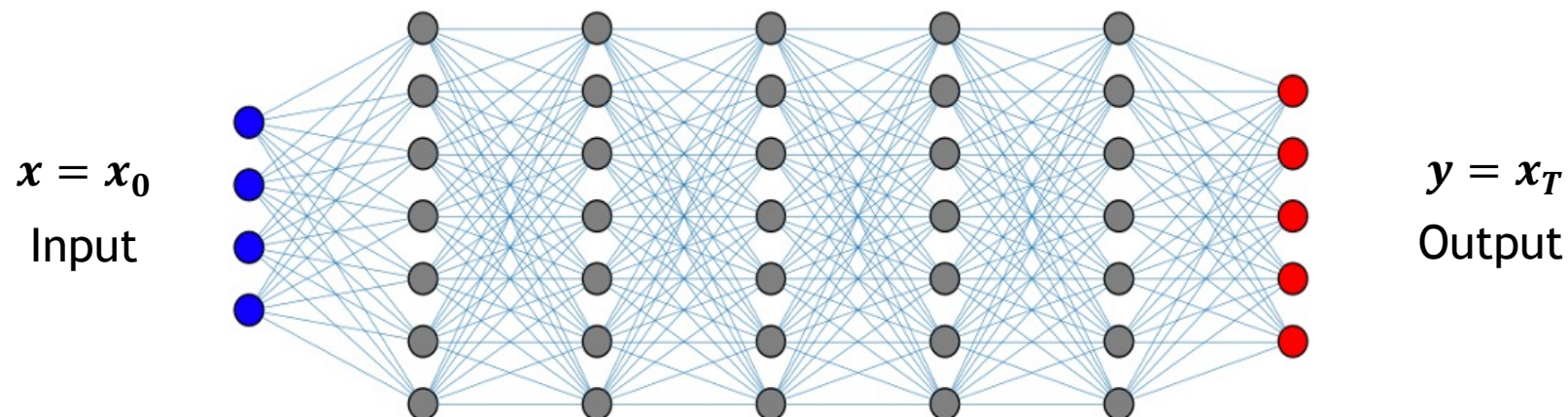
ResNet (discrete)

$$\left\{ \begin{array}{l} x_1 = x + \alpha_0 \sigma(W_0 x_0 + b_0) \\ \vdots \\ x_{n+1} = x_n + \alpha_n \sigma(W_n x_n + b_n) \\ \vdots \\ y = x_{L-1} + \alpha_{L-1} \sigma(W_{L-1} x_{L-1} + b_{L-1}) \end{array} \right.$$

Neural ODE (continuous)

$$\frac{dx}{dt} = \sigma(W(t)x + b(t))$$

$$x(0) = x \quad x(T) = y$$

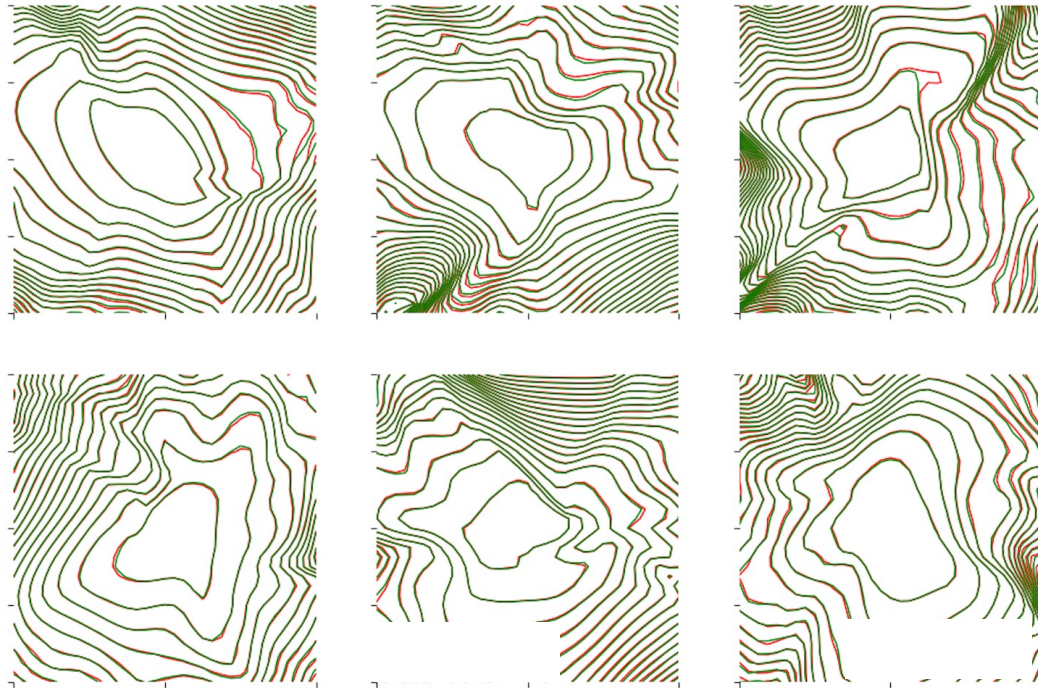


ResNets regularize loss landscape compared to MLPs



MLP NN: $x_{n+1} = \sigma(W_n x_n + b_n)$

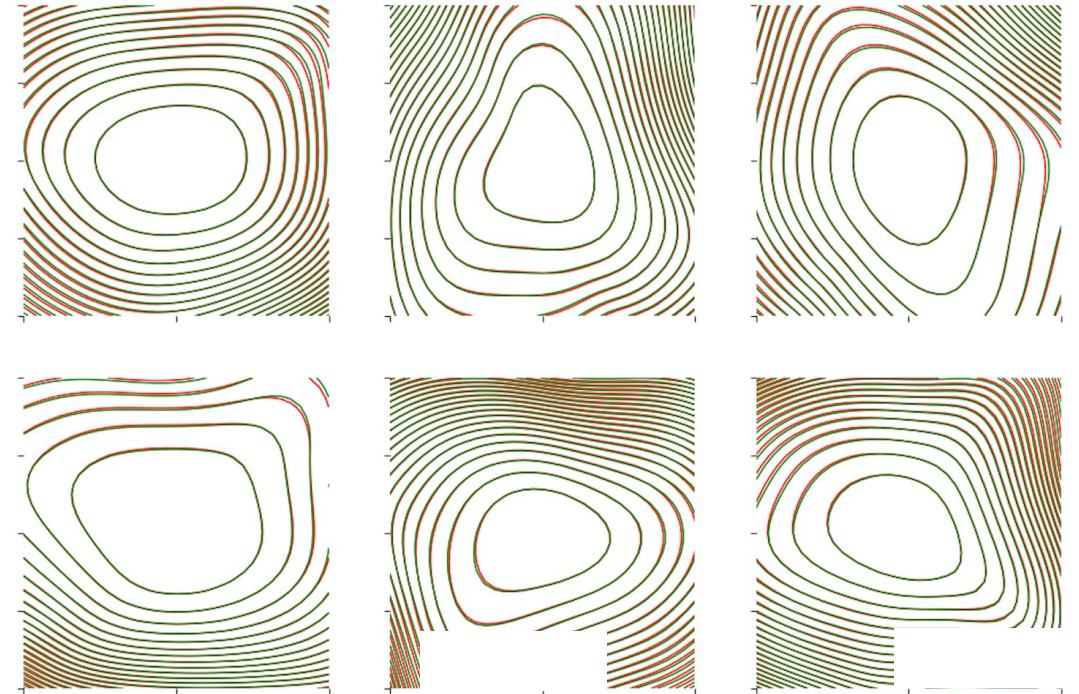
Multilayer Perceptron (learning the layer)



— Training (N=1988) — Testing (N=498)

ResNet: $x_{n+1} = x_n + \alpha_n \sigma(W_n x_n + b_n)$

ResNets (learning the layer diff.)



— Training (N=1988) — Testing (N=498)

Weight parameterization as a regularization tool, inspired by ODEs



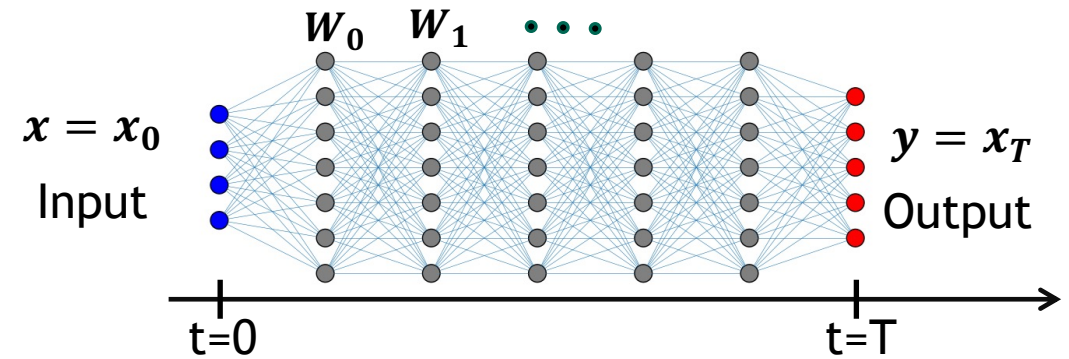
ResNet: $x_{n+1} = x_n + \alpha_n \sigma(W_n x_n + b_n)$

Training for weight matrices W_0, W_1, \dots

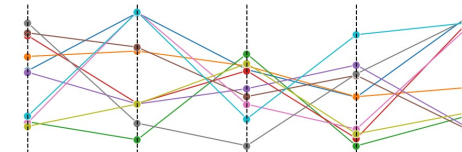
Heavily overparameterized,
does not generalize well

Parameterize $W(t; \alpha)$ and train for α 's.

Parameterization of weight functions
reduces capacity and
improves generalization

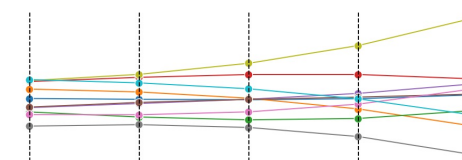


Business
as usual

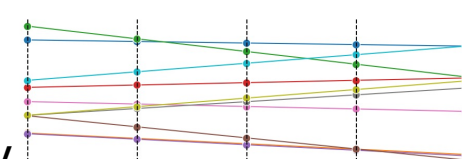


NonPar $W(t; \alpha)$
 $= W_{tL/T}$

Dial down
complexity



Cubic $W(t; \alpha)$
 $= \alpha t^3 + \beta t^2 + \dots$

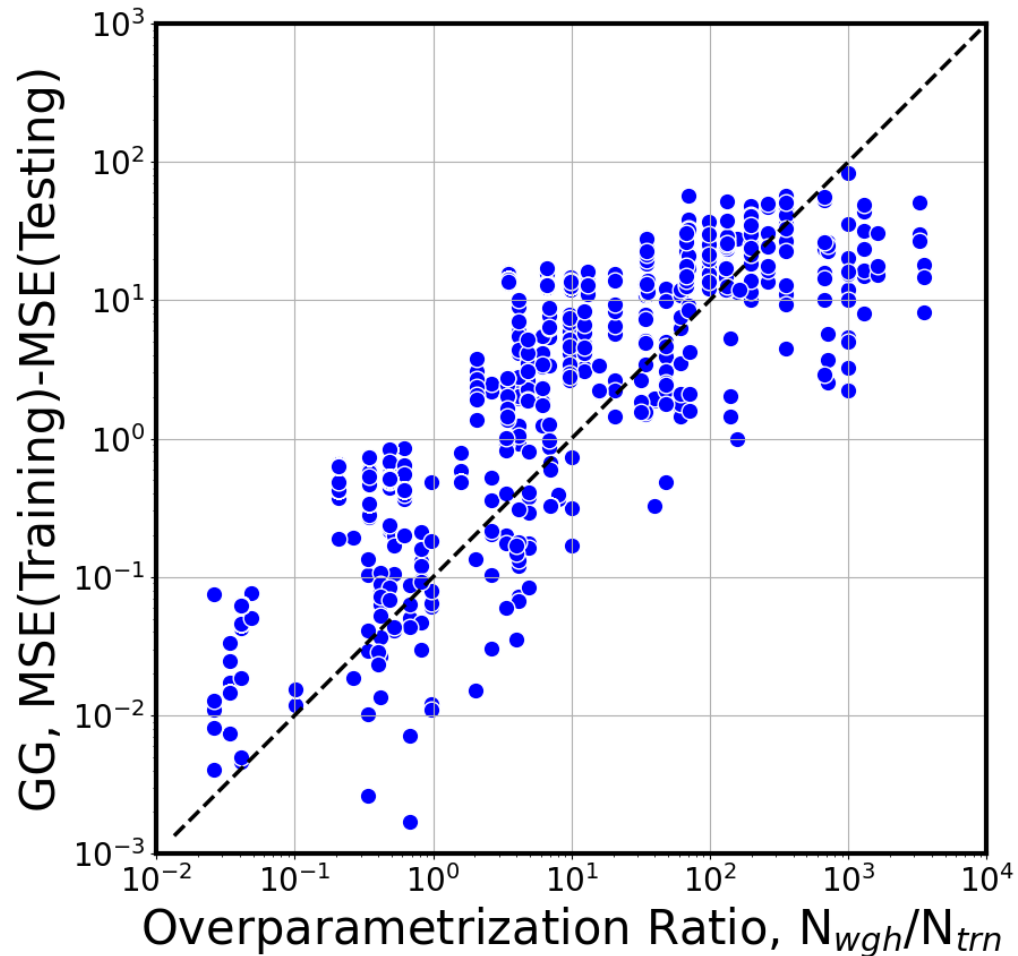


Linear $W(t; \alpha)$
 $= \alpha t + \beta$

Weight parameterization (WP) improves generalization



Better Generalization

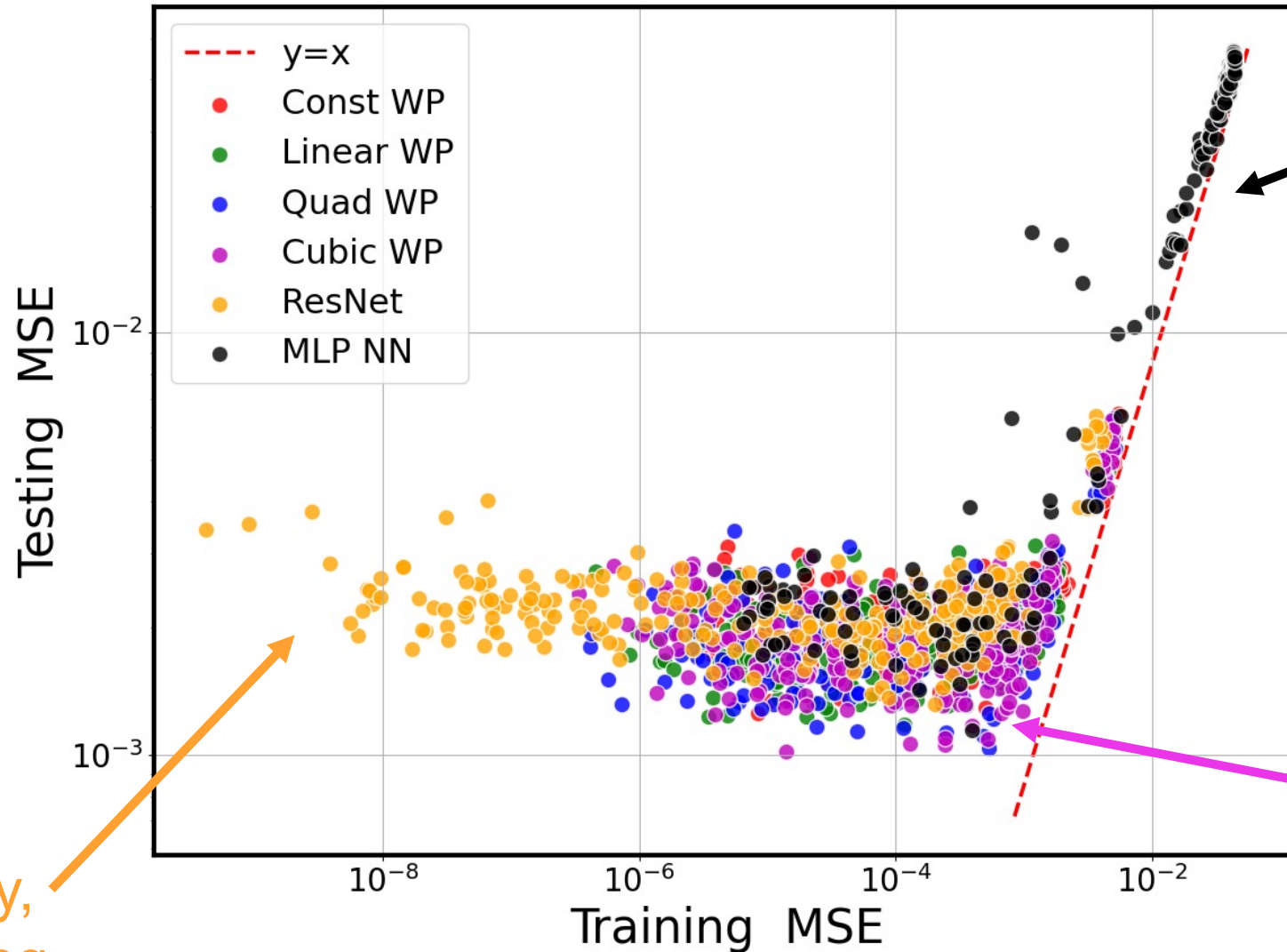


- Generalization Gap correlates with overparameterization
- Weight-parameterized ResNets reduce Generalization Gap

Each dot is a training run with varying weight parameterization functions

Weight Parameterization

ResNet + WP improves accuracy

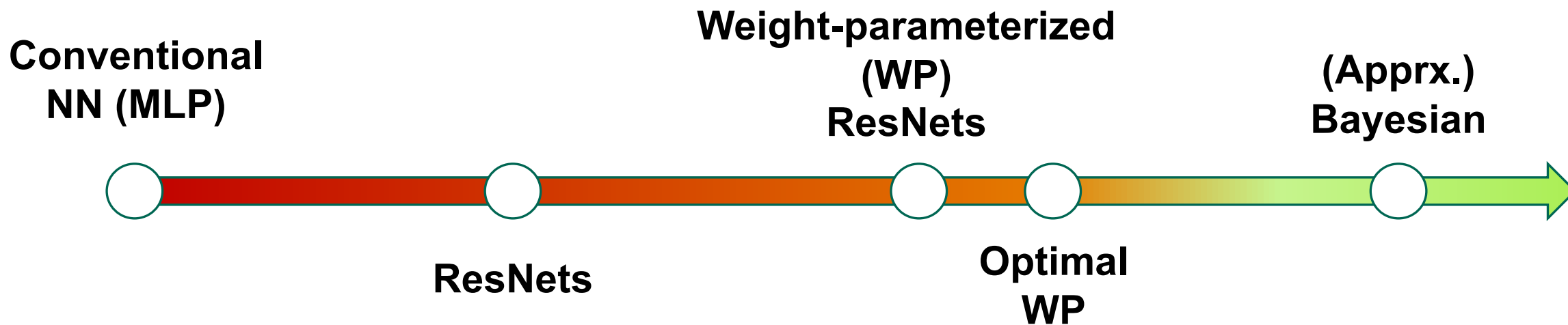


ResNet:
ok accuracy,
but overfitting

MLP:
bad accuracy

WP:
perfect tradeoff

Progressive improvements: path toward generalization and confidence assessment



Probabilistic Learning: Bayesian NN



- Conventional NN: training for deterministic weight matrices W_0, W_1, \dots
- Probabilistic approach: training for probability distributions $p(W_0), p(W_1), \dots$
- Three classes of options:

Full Bayesian



Approximate Bayesian



Ensemble methods

- Markov chain Monte Carlo (MCMC)
 - Typically, infeasible for overparameterized NNs
 - With weight parameterization loss functions are better behaved (lower-dimensional, fewer symmetries), hence MCMC path more feasible
- Variational methods
 - Practically feasible, but many hyperparameters to tune
 - Typically underestimates extrapolative predictions
- Heuristic, but
 - ... works best for complex models
 - Deep ensembles, committee of models
 - Many recent papers connecting as a Bayesian approximation

Q*U*iNN: Quantifying Uncertainties in NN

software soon-to-be-released on github



Deterministic

torch.nn.module

Probabilistic

wrapper(torch.nn.module)

Usage:

uqnet = MCMC_NN(nnet)

```
class MCMC_NN(QUiNNBase):
    def __init__(self, nnmodule, verbose=True):
        super(MCMC_NN, self).__init__(nnmodule)
        self.verbose = verbose
```

Option 1: MCMC

uqnet = VI_NN(nnet)

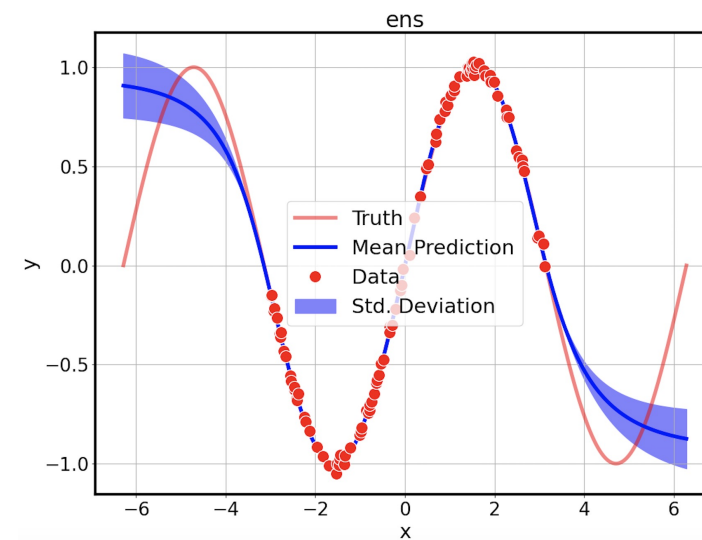
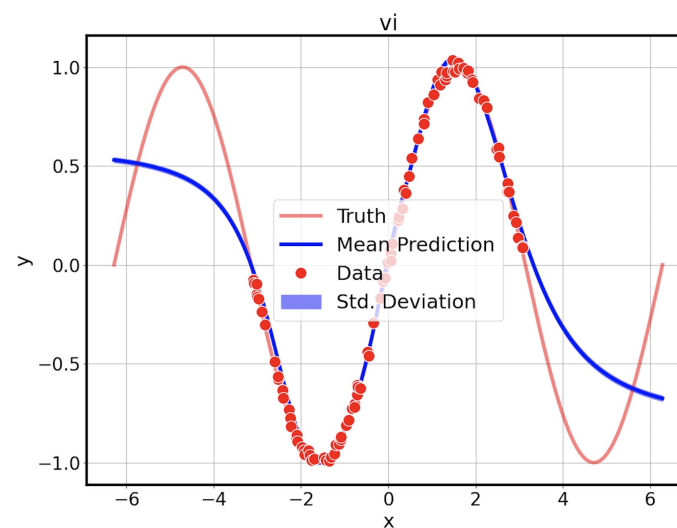
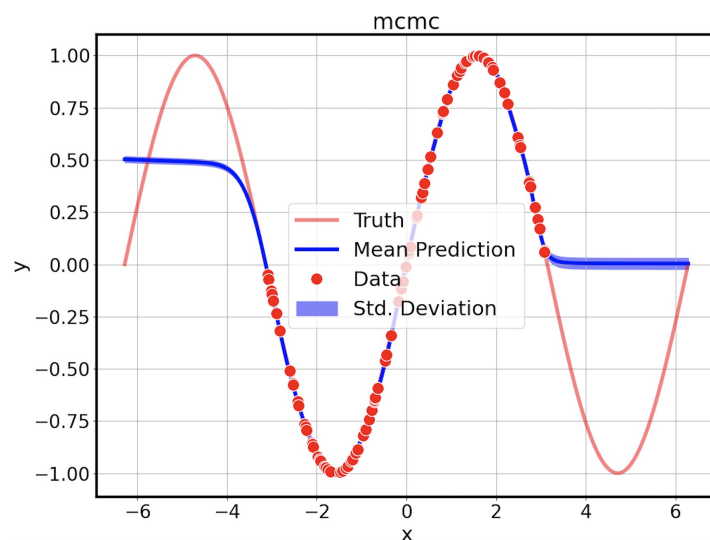
```
class VI_NN(QUiNNBase):
    def __init__(self, nnmodule, verbose=False):
        super(VI_NN, self).__init__(nnmodule)
        self.bmodel = BNet(nnmodule)
        self.verbose = verbose
```

Option 2: Variational Inference

uqnet = Ens_NN(nnet, nens=nmc)

```
class Ens_NN(QUiNNBase):
    def __init__(self, nnmodule, nens=1, verbose=False):
        super(Ens_NN, self).__init__(nnmodule)
        self.verbose = verbose
        self.nens = nens
```

Option 3: Ensembling



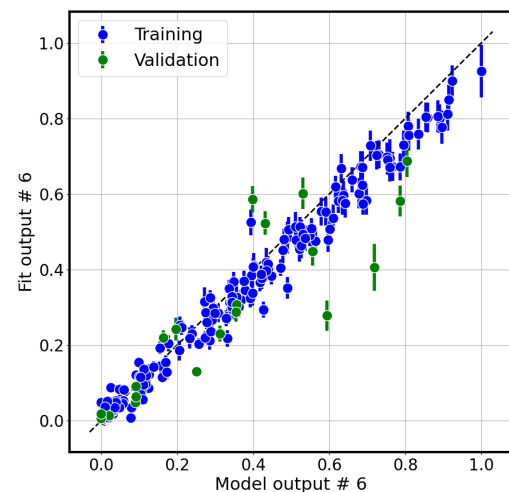
Apps:

- Multiple applications are informing the development of foundational research
- None of these applications have been previously exposed to NN prediction uncertainties, particularly in the context of ResNets and weight parameterization



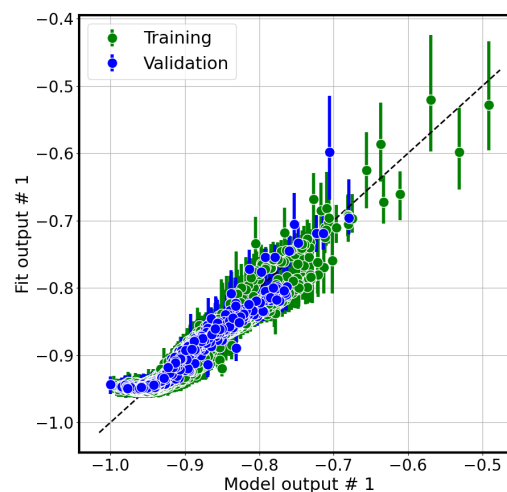
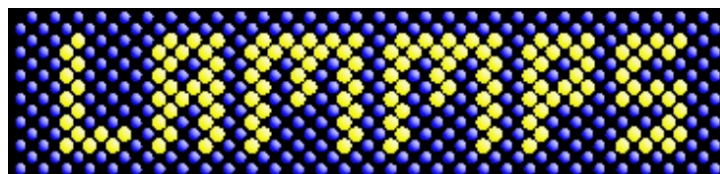
E3SM Vegetation Dynamics

- 15 input parameters
- 10 static output QoIs
- 2K training simulations



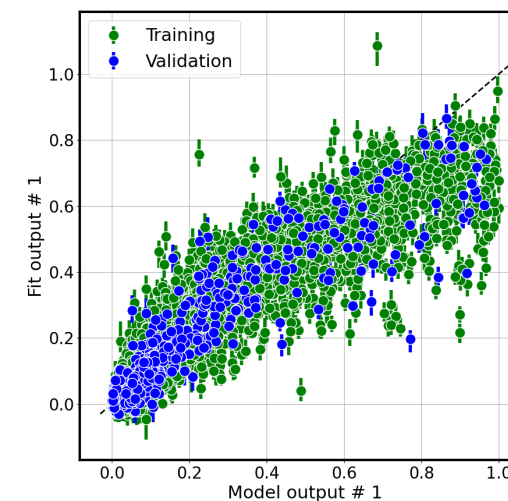
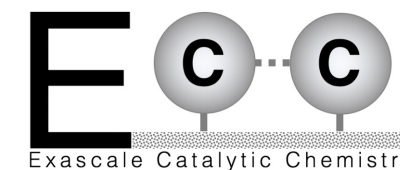
FitSNAP Entropy Dataset

- 30 input bases
- 1 output (Energy/Force/Stress)
- 20K training DFT simulations

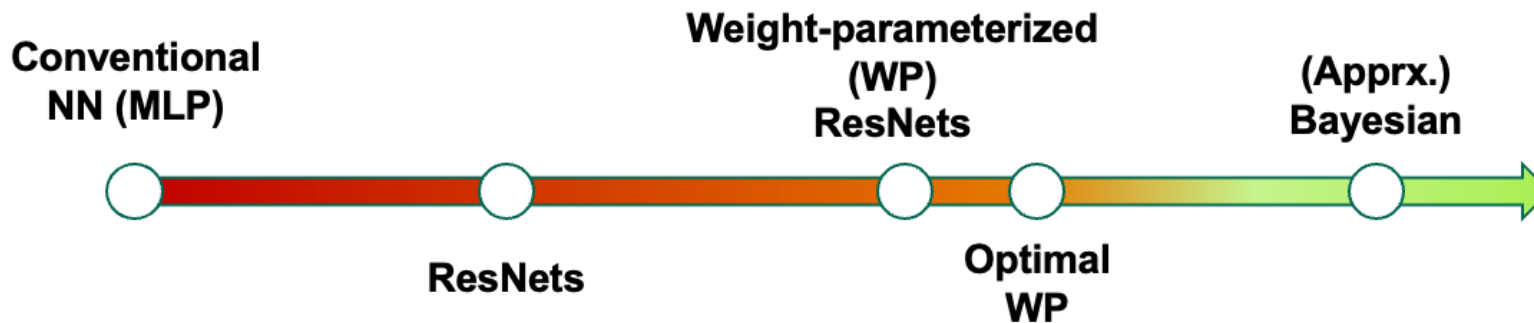


CO-on-Pt(111) Adsorbate

- 6 input d.o.f.
- 1 output (Energy)
- 10K training DFT simulations



Summary

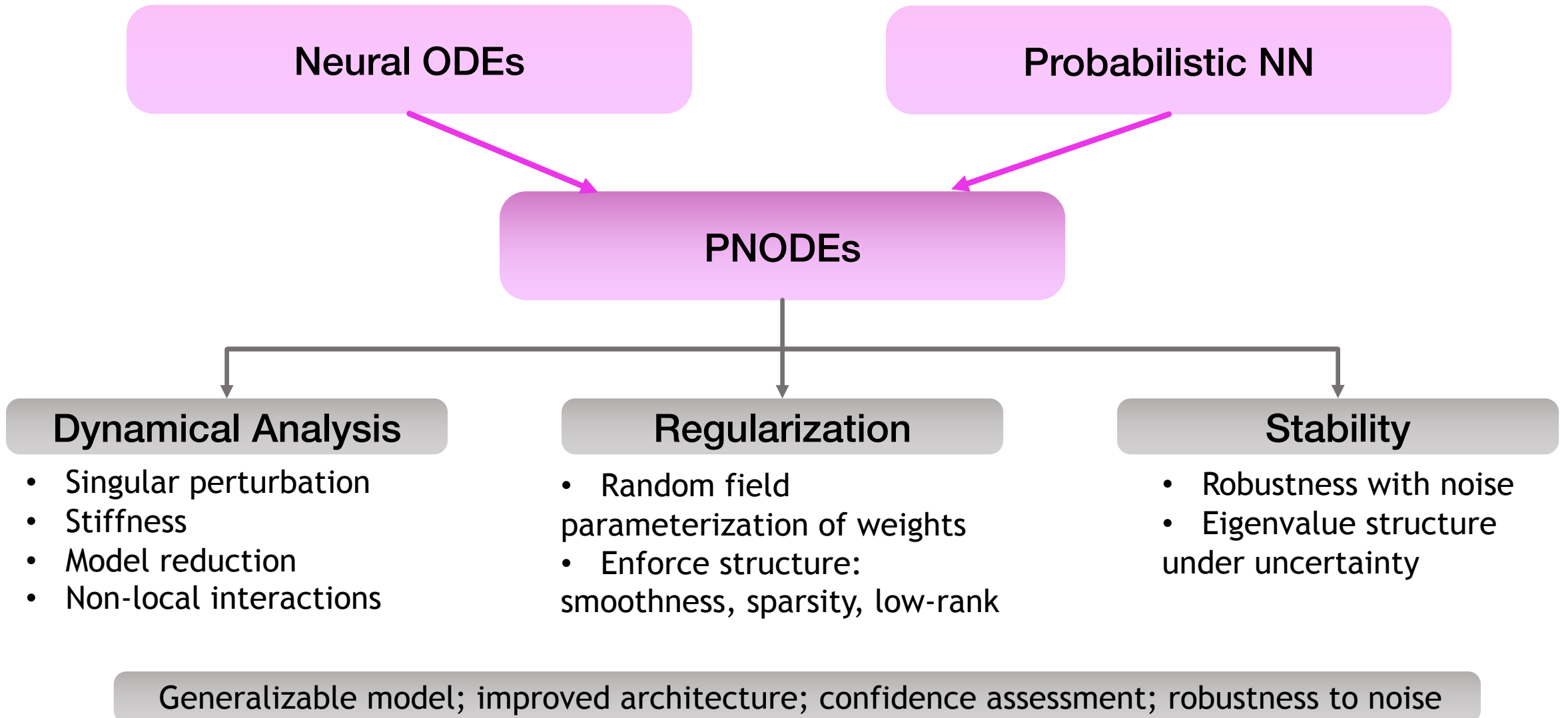


- Focus on ResNets and draw inspiration from ODEs
- ResNets regularize the learning problem, smoother loss surface
- **Weight parameterization (WP)** allows further regularization
- Optimal (e.g., ortho basis) WP for better training and more accuracy
- **Probabilistic approaches** more feasible with weight-parameterized ResNets
- Need to find sweet spot between empirical to fully Bayesian



Extra Materials

Analysis of Neural Networks as Random Dynamical Systems



Foundational capabilities impacting multiple applications



Predictive capability of Neural Networks (NNs) hinges on generalization (ability to predict well outside training data).

Regularization of NNs as a way to achieve generalization.

- ✓ Stiffness Penalization
- ✓ Weight Parameterization
- ✓ Probabilistic Weights

Methods

- ✓ Climate Land Modeling
- ✓ Catalytic Chemistry
- ✓ Materials Science

Applications

DTO vs OTD



ResNet

$$\mathbf{x}_{n+1} = \mathbf{x}_n + \alpha_n \sigma(\mathbf{W}_n \mathbf{x}_n + \mathbf{b}_n)$$

NODE

$$\frac{d\mathbf{x}}{dt} = \sigma(\mathbf{W}(t)\mathbf{x} + \mathbf{b}(t))$$

Forward equivalence:

Neural ODE discretized using explicit Euler and ResNet produce identical outputs choosing time step: $\Delta t = \frac{T}{L}$, $\alpha_n := \Delta t$, $\mathbf{W}_n := \mathbf{W}(n\Delta t)$ and $\mathbf{b}_n := \mathbf{b}(n\Delta t)$ for all n .

Backward not so much:

Gradient computations differ!

Consider $W(t) \equiv W$ and $b \equiv 0$:

Discretized Neural ODE with adjoint method:

$$\nabla loss = 2((1 + \delta t W)^L x - y)(1 + \delta t W)^L x$$

ResNet with backpropagation:

$$\nabla loss = 2((1 + \delta t W)^L x - y)(1 + \delta t W)^{L-1} x$$

- Gradients converge as $L \rightarrow \infty$ but differences can be large for small L ,
- Optimize then discretize (adjoint method) \neq discretize then optimize (backpropagation).

Prior Work on Probabilistic NN



- Probabilistic NN have been around since 90s [MacKay, 1992; Neal, 1997]
 - Full probabilistic treatment was infeasible back then (and still is, generally)
 - Recent work showed avenues via variational methods with clever tricks:
 - Bayes by Backprop [Blundell, 2015]; Probabilistic backprop [Hernandez-Lobato 2015]
 - Ghahramani, “Probabilistic Machine Learning and Artificial Intelligence”. *Nature*, 2015
 - “Nearly all approaches to probabilistic programming are Bayesian since
it is hard to create other coherent frameworks for automated reasoning about uncertainty”
 - Industry *is* catching up: Bayesflow at Google, infer.NET at Microsoft, Uber has shown interest
 - Still not industry-standard: expensive, not well understood.