# Understanding Deep Architectures with Reasoning Layer

Xinshi Chen[1], Yufei Zhang[2], Christoph Reisinger[2], Le Song[1]

[1]Georgia Institute of Technology,   [2]University of Oxford

NeurIPS 2020

# Unrolled Algorithm As A Layer

$$y^* = \mathbf{argmin}_y E^*(y)$$

---
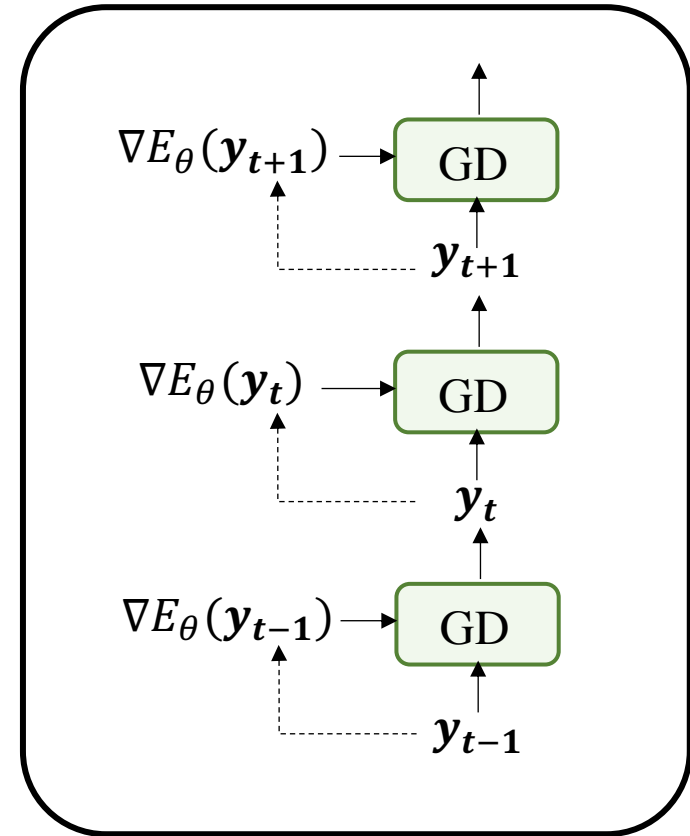Iterative Algorithm (Gradient Descent)

For $k = 1, 2, \ldots$ do
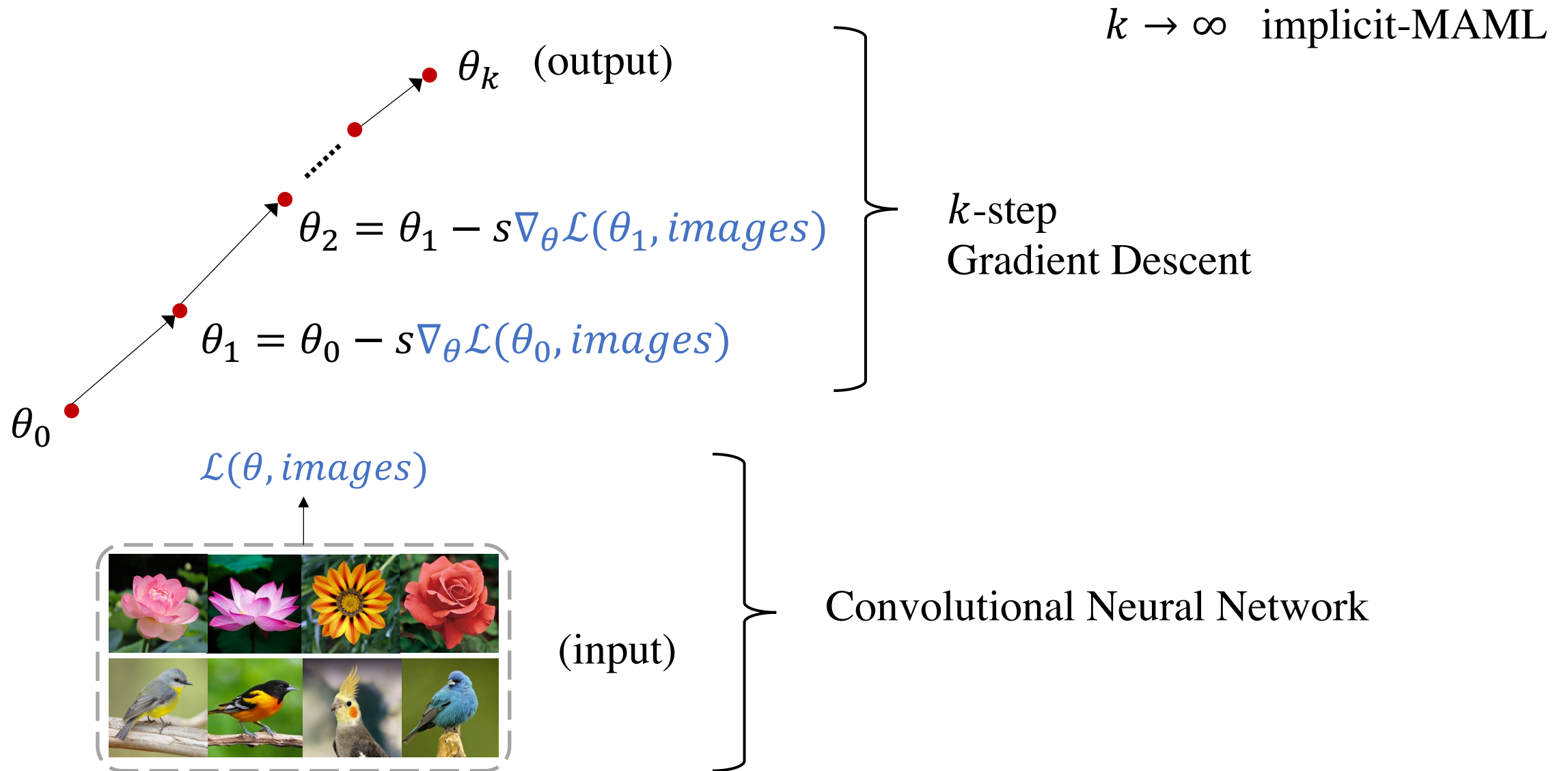
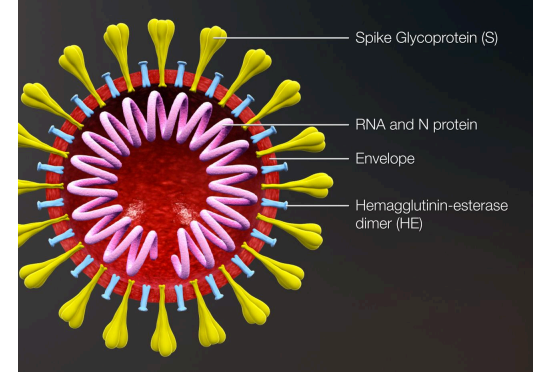$$y_{t+1} \leftarrow y_k - s\nabla E(y_k)$$

Done

---

Unrolling



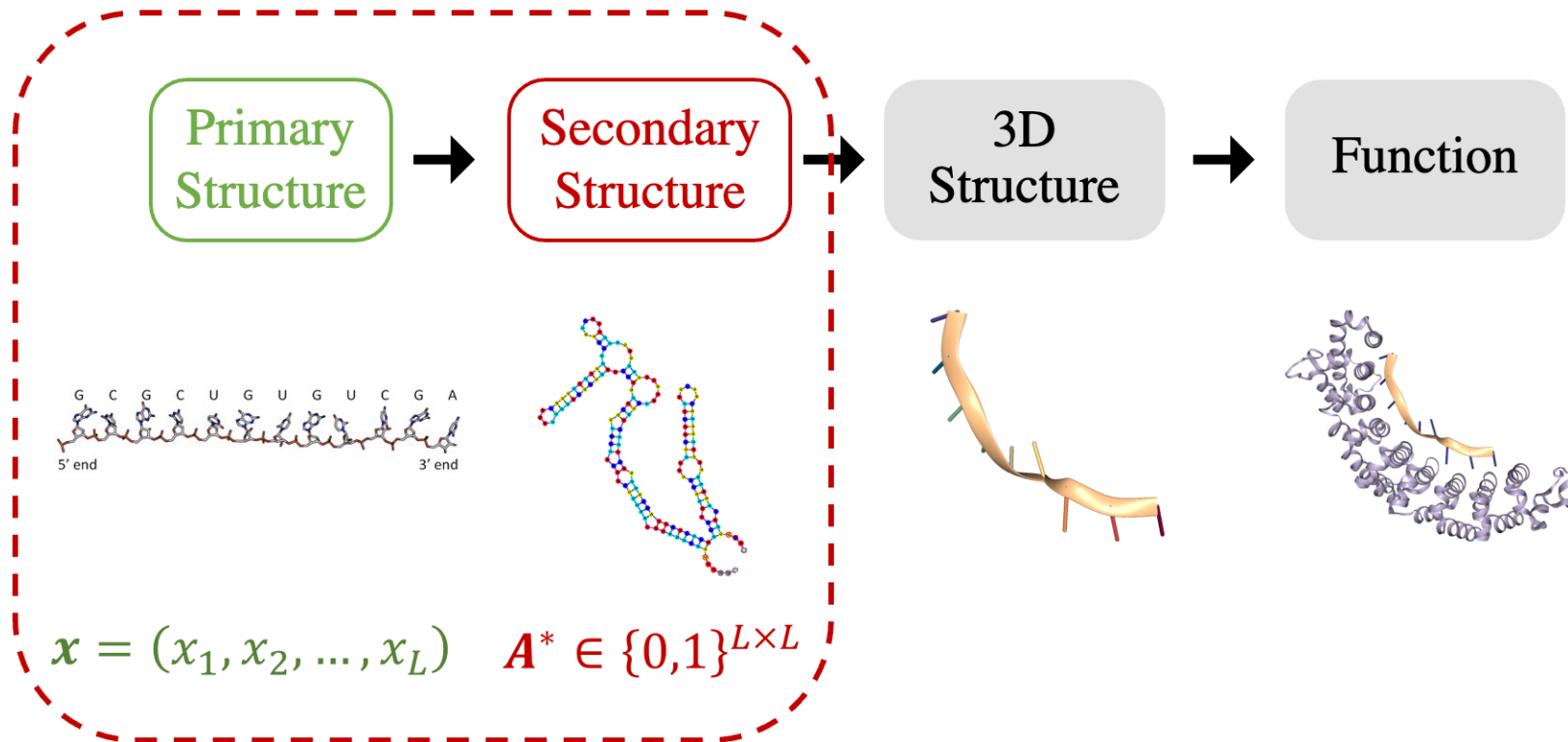An algorithm can be unrolled and truncated and then used as a specialized layer in the deep learning model.

# Ex 1: MAML (Model-Agnostic Meta-Learning)



$k \to \infty$   implicit-MAML

$\theta_k$   (output)

$\theta_2 = \theta_1 - s\nabla_\theta \mathcal{L}(\theta_1, images)$

$\theta_1 = \theta_0 - s\nabla_\theta \mathcal{L}(\theta_0, images)$

$k$-step
Gradient Descent

$\theta_0$

$\mathcal{L}(\theta, images)$

(input)

Convolutional Neural Network

# Ex 2: RNA Secondary Structure Prediction



RNA secondary structure prediction



Primary Structure → Secondary Structure → 3D Structure → Function

$$\boldsymbol{x} = (x_1, x_2, \ldots, x_L)$$

$$\boldsymbol{A}^* \in \{0,1\}^{L \times L}$$

# Ex 2: E2Efold -- Constrained Optimization Solver as a Layer

$$x = (x_1, x_2, \ldots, x_L)$$

Transformer & Convolution $\longrightarrow$ $E_\theta(x, A)$

Unrolled Algorithm

$$A^* \approx \operatorname{argmin}_{A \in \mathcal{A}} E_\theta(x, A)$$

$A^*$

G C G C U G U G U C G A
5' end          3' end

Highly **expressive** module

Encode complex sequence information and dependency

Highly **structured** module

- Enforces the constraints
- Restrict the output space

$$\min_{A \in [0,1]^{L \times L}} E_\theta(x, A) + \rho \|A\|_1$$

$$\text{s.t.} \quad M(x) \circ A = A$$
$$A^\top = A$$
$$A\mathbf{1} \leq \mathbf{1}$$
$$A \geq 0$$

*Chen, Xinshi, et al. "RNA Secondary Structure Prediction By Learning Unrolled Algorithms." ICLR 2019*

# Hybrid Architecture

$$Alg_{\phi}^{k}(E_{\theta}(x, \cdot))$$

**reasoning module**

**(algorithm layer)**

- Unrolled iterative algorithms
- Executes prescribed operations
- Interpretable

$$E_{\theta}(x, \cdot)$$

**neural module**

- Model complex information of the inputs

$x$

End-to-end differentiable architecture trained with $(x, y^*)$ pairs

$\nabla E_{\theta}(y_{t+1}) \rightarrow$ GD $\rightarrow y_{t+1}$

$\nabla E_{\theta}(y_{t}) \rightarrow$ GD $\rightarrow y_{t}$

$\nabla E_{\theta}(y_{t-1}) \rightarrow$ GD $\rightarrow y_{t-1}$

G C G C U G U G U C G A

5' end          3' end

# Hybrid Architecture

$$Alg_{\phi}^{k}(E_{\theta}(x, \cdot ))$$



**reasoning module**

$$E_{\theta}(x, \cdot )$$

**neural module**

$x$

Approximate

$$y^* = \mathbf{argmin}_y \, E^*(x, y)$$

optimization over $y$

$E^*(x, y)$  some unknown energy function

$x$

# Questions

- Different algorithms can solve the SAME reasoning task
- How are they differ~~~ each other when use~~~ reasoning module?

$$Alg_{\phi}^{k}(E_{\theta}(x, \cdot))$$

**reasoning module**

$\nabla E_{\theta}(y_{t+1}) \rightarrow$ GD
$y_{t+1}$

$\nabla E_{\theta}(y_{t}) \rightarrow$ GD
$y_{t}$

$\nabla E_{\theta}(y_{t-1}) \rightarrow$ GD
$y_{t-1}$

$$E_{\theta}(x, \cdot)$$

**neural module**

$x$

G C G C U G U G U C G A

5' end          3' end

Algorithm properties
- convergence
- stability
- sensitivity

Effects? Relation?

Representation & generalization ability of the overall hybrid model

# Problem Setting: Optimization Module + Neural Energy Module

Unrolled **Optimization** Algorithm

$$Alg_{\phi}^{k}\left(\underbrace{E_{\theta}(x, y)}_{\text{Neural Energy}}\right)$$

- $\theta$ : parameters in the neural module
- $\phi$ : step size in the unrolled algorithm
- $k$ : number of unrolled iterations

- We restrict to the case when $E_{\theta}(x, y)$ is a quadratic function in $y$.

  - $E_{\theta}(x, y) = \frac{1}{2} y^{\top} Q_{\theta}(x) y + y^{\top} b$, where $Q_{\theta}(x)$ is a neural network.

- Ground-truth model is $y^* = \mathbf{argmin}_y E^*(x, y)$ for some unknown energy function $E^*$

- Training dataset contains $n$ many input-output pairs $(x, y^*)$, without intermediate supervision on $E^*$

# How To Design The Reasoning Module (Algorithm Layer)?

- Different optimization algorithms, which one is better?

$Alg = $ Gradient Descent

$$GD_\phi^k\left(E_\theta(x, y)\right)$$

$Alg = $ Nesterov's Accelerated Gradient

$$NAG_\phi^k\left(E_\theta(x, y)\right)$$

?

- More iterations $k$, the better?

$$GD_\phi^k\left(E_\theta(x, y)\right)$$

?

Equilibrium model

$$GD_\phi^\infty\left(E_\theta(x, y)\right) = \mathbf{argmin}_y\left(E_\theta(x, y)\right)$$

# Algorithm Property

1) Convergence
   - portrays how fast the optimization error decreases ask the number of iterations $k$ grows.

$$\| Alg_\phi^k(E(x,y)) - \text{argmin}_y(E(x,y)) \| \leq \boldsymbol{Cvg}(\boldsymbol{k},\boldsymbol{\phi}) \| Alg_\phi^0(E(x,y)) - \text{argmin}_y(E(x,y)) \|$$

2) Stability
   - characterizes its robustness to small perturbations in the optimization objective $\boldsymbol{E}_\theta(\boldsymbol{x},\boldsymbol{y})$.

$$\| Alg_\phi^k(E(x,y)) - Alg_\phi^k(\hat{E}(x,y)) \| \leq \boldsymbol{Stab}(\boldsymbol{k},\boldsymbol{\phi}) \| E - \hat{E} \|_\infty$$

3) Sensitivity
   - characterizes its robustness to small perturbations in the step size $\phi$ in the algorithm

$$\| Alg_\phi^k(E(x,y)) - Alg_\varphi^k(E(x,y)) \| \leq \boldsymbol{Sens}(\boldsymbol{k}) \, |\phi - \varphi|$$
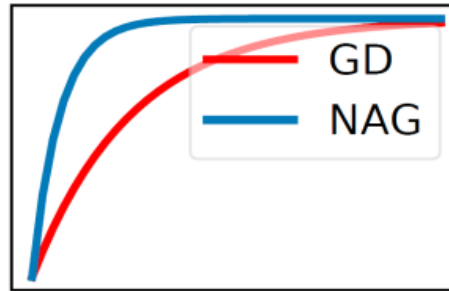
   - Robustness to perturbations in parameters is referred in the deep learning community to "parameter perturbation error" or "sharpness".
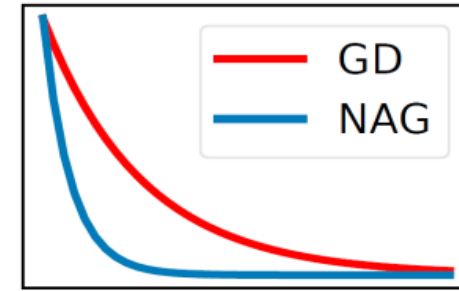
# GD and NAG: Algorithm Property Comparison

| Alg | $Cvg(k, \phi)$ | $Stab(k, \phi)$ | $Sens(k)$ |
|-----|-----|-----|-----|
| $\mathtt{GD}_\phi^k$ | $\mathcal{O}\left((1 - \phi\mu)^k\right)$ | $\mathcal{O}\left(1 - (1 - \phi\mu)^k\right)$ | $\mathcal{O}\left(k(1 - c_0\mu)^{k-1}\right)$ |
| $\mathtt{NAG}_\phi^k$ | $\mathcal{O}\left(k(1 - \sqrt{\phi\mu})^k\right)$ | $\mathcal{O}\left(1 - (1 - \sqrt{\phi\mu})^k\right)$ | $\mathcal{O}\left(k^3(1 - \sqrt{c_0\mu})^k\right)$ |



NAG converges faster

GD more stable

NAG less sensitive

Faster algorithm less stable

# Main Theorem: Local Rademacher Complexity

Local Rademacher complexity of $\boldsymbol{Alg_\phi^k}\left(\boldsymbol{E_\theta(x,y)}\right)$

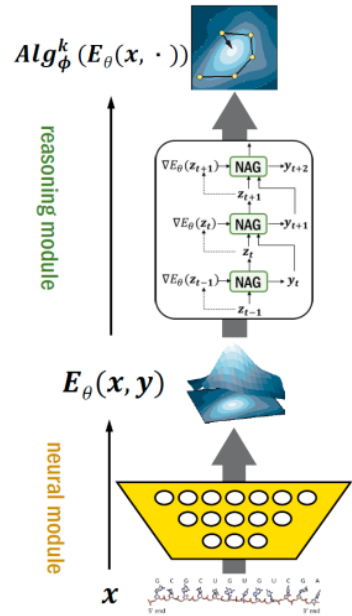**Theorem 3.1.** *Assume the problem setting in Sec 2. Then we have for any $t > 0$, it holds true that*

$$\mathbb{E}R_n\ell_{\mathcal{F}}^{loc}(r) \leq \sqrt{2}dn^{-\frac{1}{2}}Stab(k)\left(\sqrt{(Cvg(k)M + \sqrt{r})^2C_1(n) + C_2(n,t)} + C_3(n,t) + 4\right) + Sens(k)B_\Phi$$
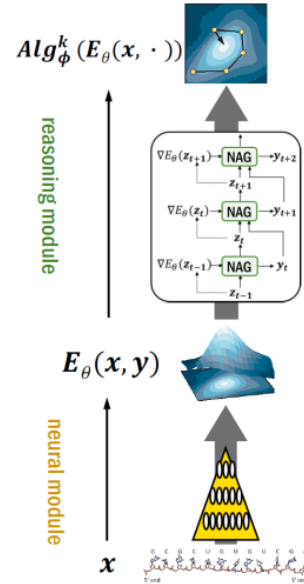
**stability**    **convergence**

**sensitivity**

complexity of
neural module

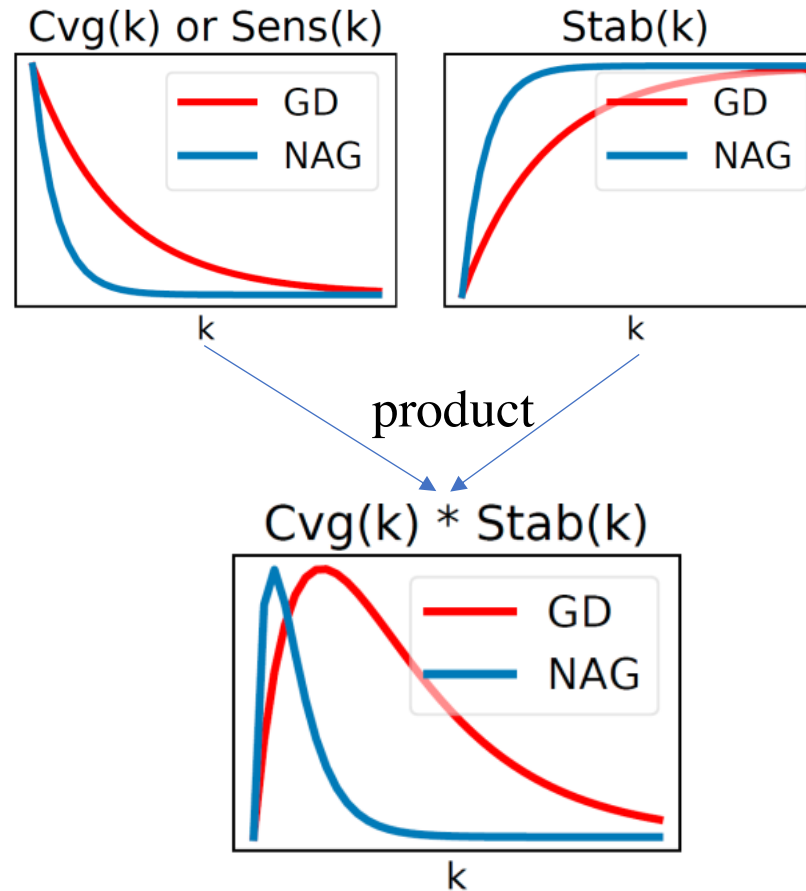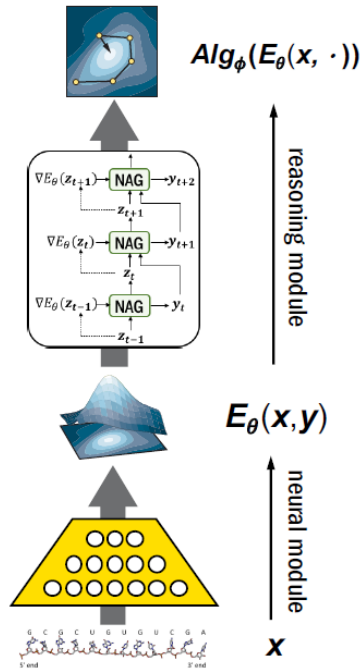# Implication I



**Over-parametrization**
$C_1, C_2, C_3$ large



**Under-parametrization**
$C_1, C_2, C_3$ small

- Bound is dominated by $Stab(k)$

- More iterations ($k \to \infty$), worse generalization

- Fix $k$, GD generalize better than NAG

# Implication II



**About-right parameterization**
$C_1, C_2, C_3$ not large or small

- Bound is dominated by the product $Stab(k) * Cvg(k)$

- More iterations ($k$ large) better generalization

# Good Fit between Experiments and Theory
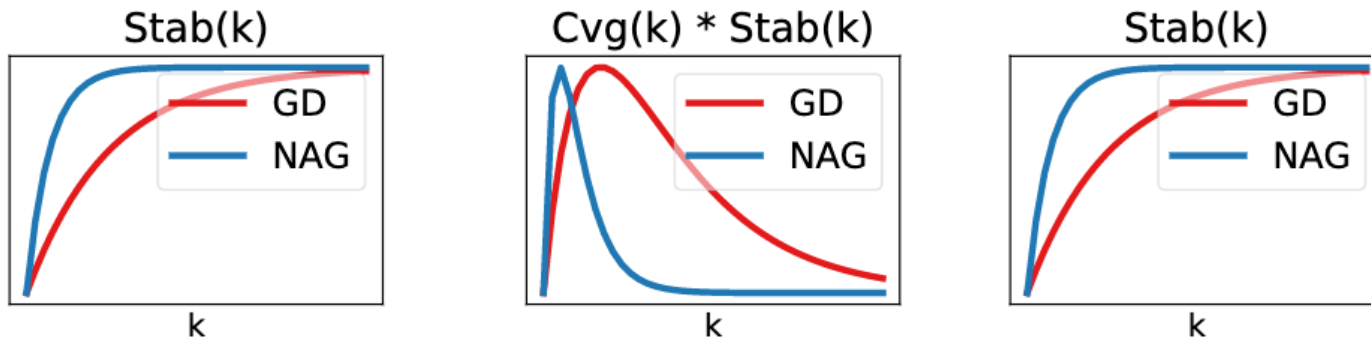
- Generalization gaps, when varying the *hidden dimension* of the neural module.



- Corresponds to the theoretically analyzed algorithm properties:



Align well with the implication of our theorem!

See more details in our paper:

**Understanding Deep Architecture With Reasoning Layer**
**https://papers.nips.cc/paper/2020/file/0d82627e10660af39ea7eb69c3568955-Paper.pdf**

Q&A!