

n-TangentProp: A Quasilinear Algorithm for
Taking Derivatives of Neural Networks

Kyle R. Chickering
UC Davis

1. Motivation
2. PINNs + Neural Networks
3. Autodiff
4. Tangent Prop + n-Tangent Prop
5. Results
6. Conclusions

1. Motivation
2. PINNs + Neural Networks
3. AutoDiff
4. Tangent Prop + n-Tangent Prop
5. Results
6. Conclusions



Paper online, work under
Dr. P. Simard
in Feb. of last year

Motivation I

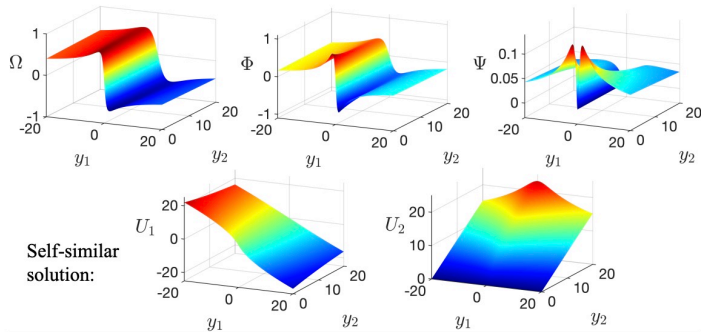
- The Navier-Stokes Existence & Smoothness
(T. Tao in Notices)

Motivation I

- The Navier-Stokes Existence & Smoothness
(T. Tao in Notices)
- Self-Similar Solutions tell us about Finite-time
Breakdown

Motivation I

- The Navier-Stokes Existence & Smoothness (T. Tao in Notices)
- Self-Similar Solutions tell us about finite-time Breakdown



Boussinesq

$$\begin{aligned} \Omega + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Omega &= \Phi, \\ (2 + \partial_{y_1} U_1)\Phi + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Phi &= -\partial_{y_1} U_2 \Psi, \\ (2 + \partial_{y_2} U_2)\Psi + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Psi &= -\partial_{y_2} U_1 \Phi, \\ \operatorname{div} \mathbf{U} &= 0. \end{aligned} \quad (3)$$

Hard to solve numerically!

From Wang, Lai, Gómez-Serrano, Buckmaster '23

Motivation II

- The "boundary conditions" for self-similar solutions look like smoothness constraints, i.e. $U \in C^5(B_\varepsilon(x_0))$ for a point x_0

Motivation II

- The "boundary conditions" for self-similar solutions look like smoothness constraints, i.e. $U \in C^5(B_\varepsilon(x_0))$ for a point x_0
- Exact equation is unknown! Must find $\lambda, U, \Omega, \Psi, \Phi$ all at once!

$$\Omega + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Omega = \Phi,$$

$$(2 + \partial_{y_1} U_1)\Phi + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Phi = -\partial_{y_1} U_2 \Psi, \quad (3)$$

$$(2 + \partial_{y_2} U_2)\Psi + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Psi = -\partial_{y_2} U_1 \Phi,$$

$$\operatorname{div} \mathbf{U} = 0.$$



We don't know what λ is!

Motivation II

- The "boundary conditions" for self-similar solutions look like smoothness constraints, i.e. $U \in C^5(B_\varepsilon(x_0))$ for a point x_0
- Exact equation is unknown! Must find $\lambda, U, \Omega, \Psi, \Phi$ all at once!

$$\Omega + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Omega = \Phi,$$

$$(2 + \partial_{y_1} U_1)\Phi + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Phi = -\partial_{y_1} U_2 \Psi, \quad (3)$$

$$(2 + \partial_{y_2} U_2)\Psi + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla \Psi = -\partial_{y_2} U_1 \Phi,$$

$$\operatorname{div} \mathbf{U} = 0.$$



We don't know what λ is!

• PINNs!

PINN & Neural Networks

Neural Networks I

- Repeated functional composition

$$u = \sigma(W^n \cdot \sigma(W^{n-1} \cdot \sigma(\dots \sigma(W^1 \cdot x))))$$

- Nonlinear activation σ

Neural Networks I

- Repeated functional composition

$$u = \sigma(W^n \cdot \sigma(W^{n-1} \cdot \sigma(\dots \sigma(W^1 \cdot x))))$$

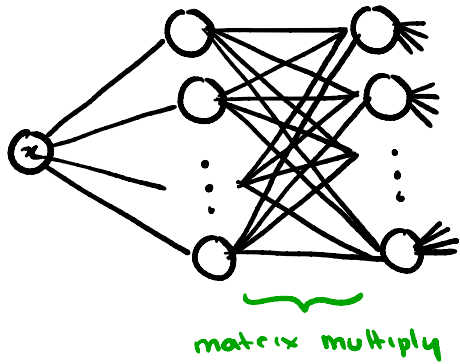
- Nonlinear activation σ
- Nonlinear is important!

Neural Networks I

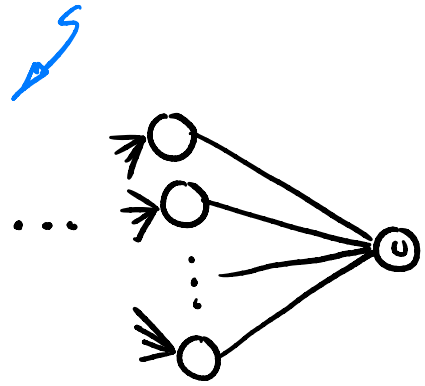
- Repeated functional composition

$$U = \sigma(W^n \cdot \sigma(W^{n-1} \cdot \sigma(\dots \sigma(W^1 \cdot x))))$$

- Nonlinear activation σ
- Nonlinear is important!



Inspired by brains



Neural Networks II

- Universal Functional approximator

Neural Networks II

- Universal Functional approximator
- As width $\rightarrow \infty$ we can approximate any continuous function on a compact domain

Neural Networks II

- Universal Functional approximator
- As width $\rightarrow \infty$ we can approximate any continuous function on a compact domain
- Does not tell us how to get approximation

Neural Networks III

- Suppose I have a labeled dataset $X = \{(x_i, y_i)\}_{i \in \mathcal{X}}$,
I want a function to approximate or generalize X ,
i.e. $f(x_i) \sim y_i \quad \forall i \in \mathcal{X}$.

Neural Networks III

- Suppose I have a labeled dataset $X = \{(x_i, y_i)\}_{i \in \mathcal{X}}$,
I want a function to approximate or generalize X ,
i.e. $f(x_i) \sim y_i \quad \forall i \in \mathcal{X}$.
- Neural net: update weights W .

Neural Networks III

- Suppose I have a labeled dataset $X = \{(x_i, y_i)\}_{i \in \mathcal{X}}$,
I want a function to approximate or generalize X ,
i.e. $f(x_i) \sim y_i \quad \forall i \in \mathcal{X}$.

- Neural net: update weights W .

- Loss:
$$L(u_w) = \frac{1}{|\mathcal{X}|} \sum_{i \in \mathcal{X}} |u_w(x_i) - y_i|^2$$

Neural Networks III

- Suppose I have a labeled dataset $X = \{(x_i, y_i)\}_{i \in \mathcal{X}}$,
I want a function to approximate or generalize X ,
i.e. $f(x_i) \sim y_i \quad \forall i \in \mathcal{X}$.
- Neural net: update weights W .
- Loss:
$$L(u_w) = \frac{1}{|\mathcal{X}|} \sum_{i \in \mathcal{X}} |u_w(x_i) - y_i|^2$$
- Update params:
$$W \leftarrow W - \alpha \nabla_W L(u_w)$$
- Gradient descent: $\nabla_W L(u_w) \rightarrow 0$ means critical pt.

PINNs

- Neural Networks are C^∞ Functions

PINNs

- Neural Networks are C^∞ Functions
- If u_θ is a neural network, we can compute $\frac{d^n}{dz^n} u_\theta$ exactly.

PINNs

- Neural Networks are C^∞ Functions
- If u_θ is a neural network, we can compute $\frac{d^n}{dz^n} u_\theta$ exactly.
- If $F(\partial^2 u; x)$ a differential equation, run gradient descent

$$L(u_\theta) = \frac{1}{N} \sum_{k=1}^N |F(\partial^2 u_\theta; x_k)|^2 + BC$$

PINNs

- Neural Networks are C^∞ Functions
- If u_θ is a neural network, we can compute $\frac{d^n}{dz^n} u_\theta$ exactly.
- If $F(\partial^2 u; x)$ a differential equation, run gradient descent

$$L(u_\theta) = \frac{1}{N} \sum_{k=1}^N |F(\partial^2 u_\theta; x_k)|^2 + BC$$

- NNs as functional approximator, $u_\theta \rightarrow u$ during training!

Accuracy Gains

- Sobolev Loss

$$L_{\text{Sob.}}^{(m)}(u_\theta) = \frac{1}{N} \sum_{k=1}^N \sum_{j=0}^m \omega_j |\nabla_x^j F(\partial^\alpha u_\theta; x_k)|^2 + \text{BC}$$

Accuracy Gains

- Sobolev Loss

$$L_{\text{Sob.}}^{(m)}(u_\theta) = \frac{1}{N} \sum_{k=1}^N \sum_{j=0}^m \omega_j |\nabla_x^j F(\partial^\alpha u_\theta; x_k)|^2 + \text{BC}$$

key takeaway: More enforced derivatives gives
better accuracy & training efficiency
(Maddu et al. '22)

Accuracy Gains

- Sobolev Loss

$$L_{\text{Sob.}}^{(m)}(u_\theta) = \frac{1}{N} \sum_{k=1}^N \sum_{j=0}^m \mathcal{Q}_j |\nabla_x^j F(\partial^\alpha u_\theta; x_k)|^2 + \text{BC}$$

Key takeaway: More enforced derivatives gives better accuracy & training efficiency

(Maddu et al. '22)

→ More Derivatives!

Do PINNs Just Suck?

- Hard to train (Hessian)

Do PINNs Just Suck?

- Hard to train (Hessian)
- Inefficient to train, expensive to train

Do PINNs Just Suck?

- Hard to train (Hessian)
- Inefficient to train, expensive to train
- No rigorous error analysis

Do PINNs Just Suck?

- Hard to train (Hessian)
- Inefficient to train, expensive to train
- No rigorous error analysis
- "Cold Water" Papers

PINN Use Cases

- Combined Forward-Inverse Problems
 - measurements from a fluid experiment
 - Solve N.S. & viscosity simultaneously

PINN Use Cases

- Combined Forward-Inverse Problems
 - measurements from a fluid experiment
 - Solve N.S. & viscosity simultaneously
- Auxillary Conditions
 - For example: Enforcing smoothness

PINN Use Cases

- Combined Forward-Inverse Problems
 - measurements from a fluid experiment
 - Solve N.S. & viscosity simultaneously
- Auxillary Conditions
 - For example: Enforcing smoothness
- If you can do it easily w/o a PINN, you probably don't need a PINN!

Autodiff

Autodiff

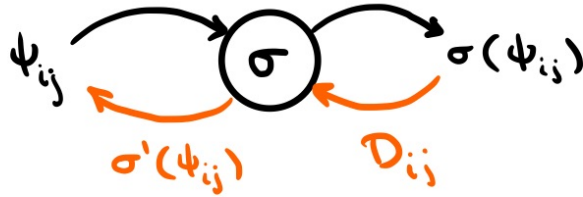
- How do we actually compute the derivatives?

Autodiff

- How do we actually compute the derivatives?
- Computational Graphs (DAG)

Autodiff

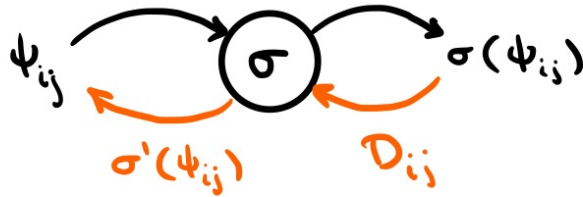
- How do we actually compute the derivatives?
- Computational Graphs (DAG)



- Forward : $\sigma(\psi)$

Autodiff

- How do we actually compute the derivatives?
- Computational Graphs (DAG)



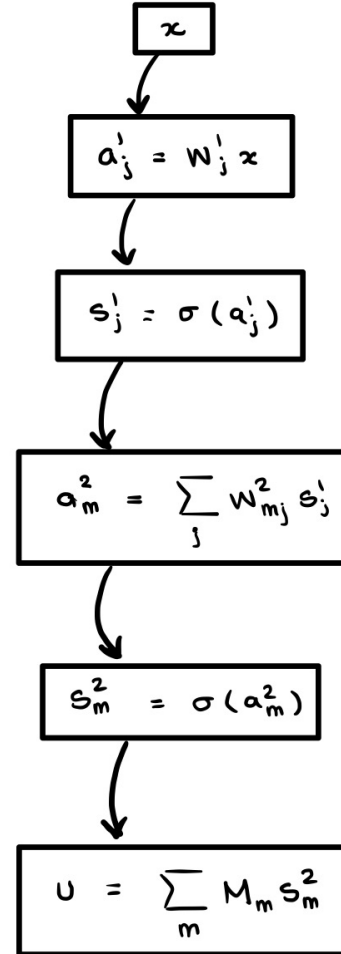
- Forward: $\sigma(\psi)$
- Backwards: $D \odot \sigma'(\psi)$

Computation Graph for NN

$$u = M \cdot \sigma(W^2 \cdot \sigma(W^1 x))$$

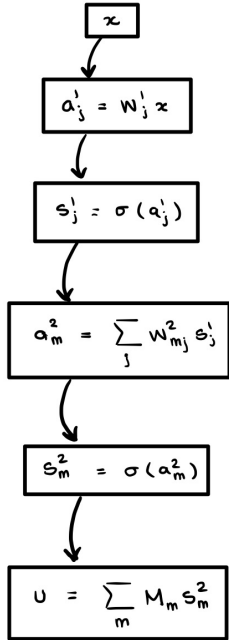
Computation Graph for NN

$$u = M \cdot \sigma(W^2 \cdot \sigma(W^1 x))$$



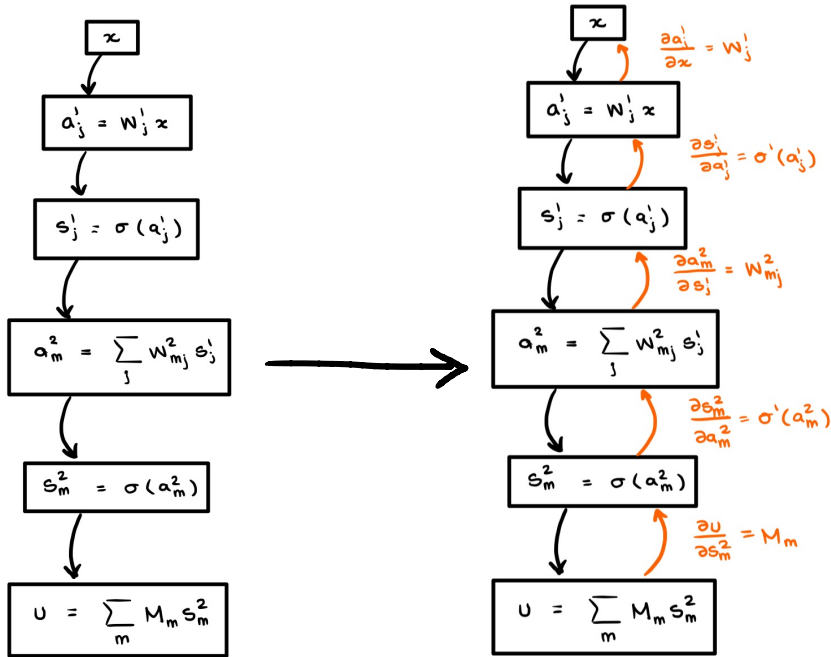
Autodiff Algorithm

- Forwards \rightarrow Backwards \rightarrow New Graph



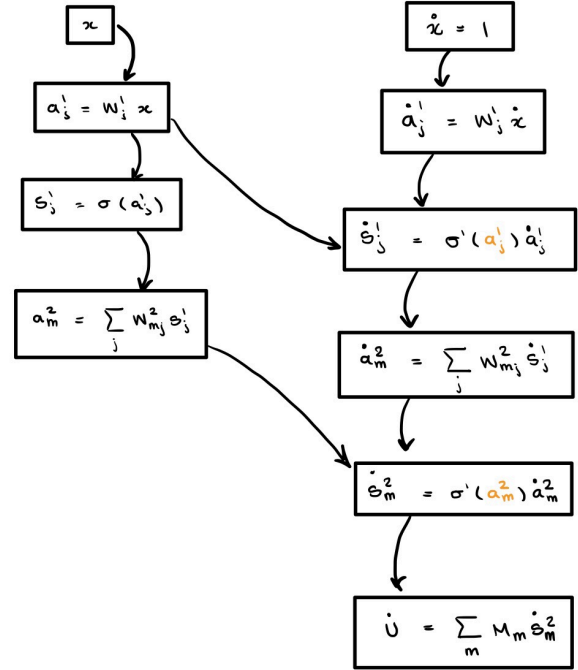
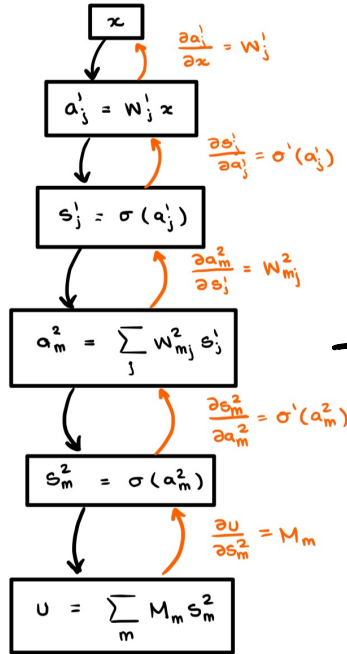
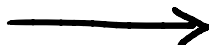
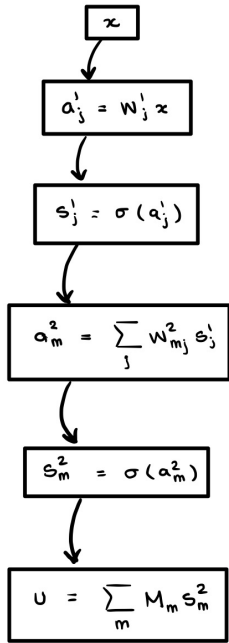
Autodiff Algorithm

- Forwards \rightarrow Backwards \rightarrow New Graph



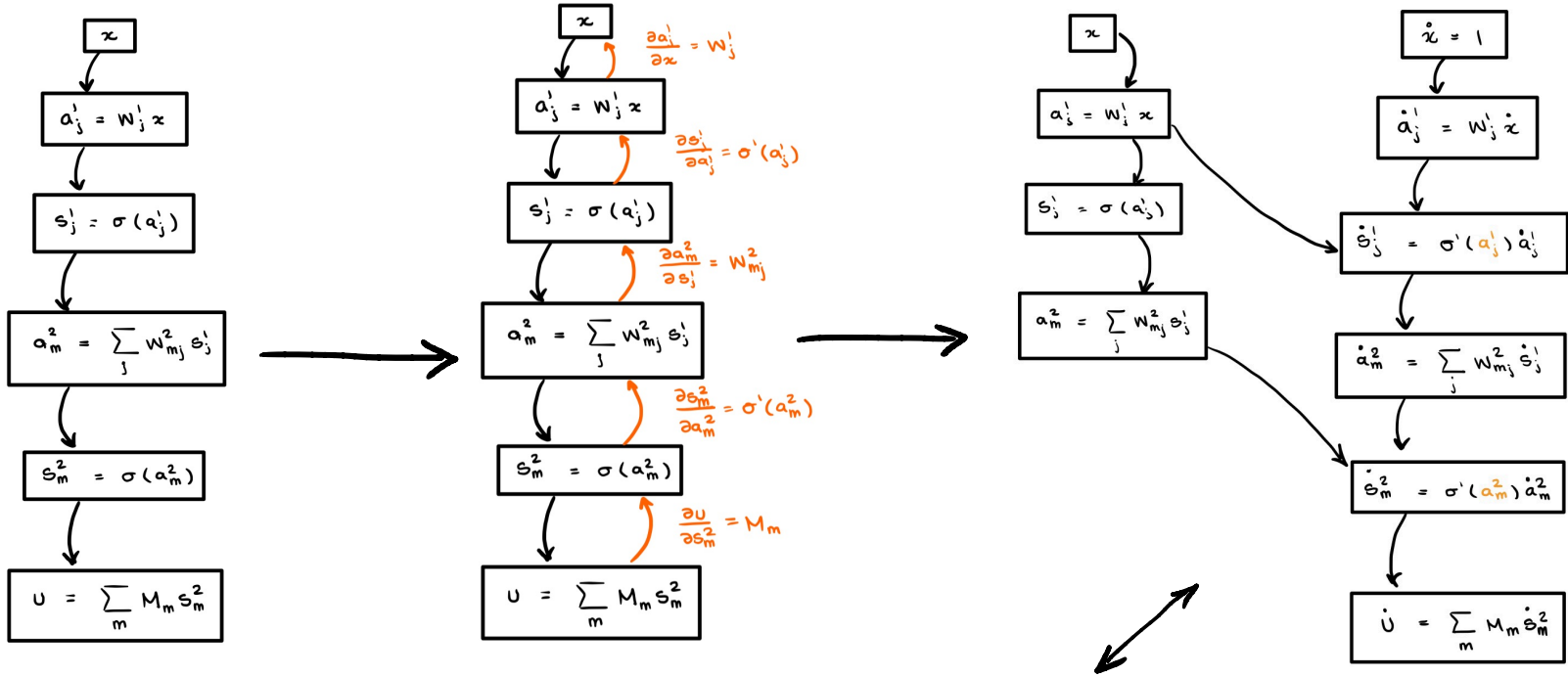
Autodiff Algorithm

- Forwards \rightarrow Backwards \rightarrow New Graph



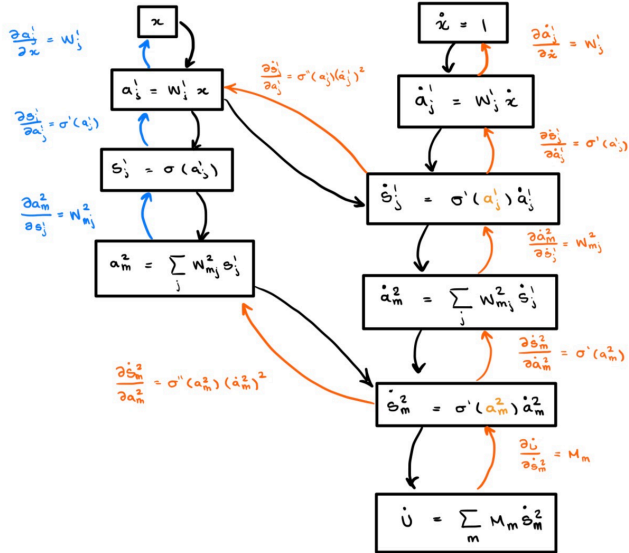
Autodiff Algorithm

- Forwards \rightarrow Backwards \rightarrow New Graph



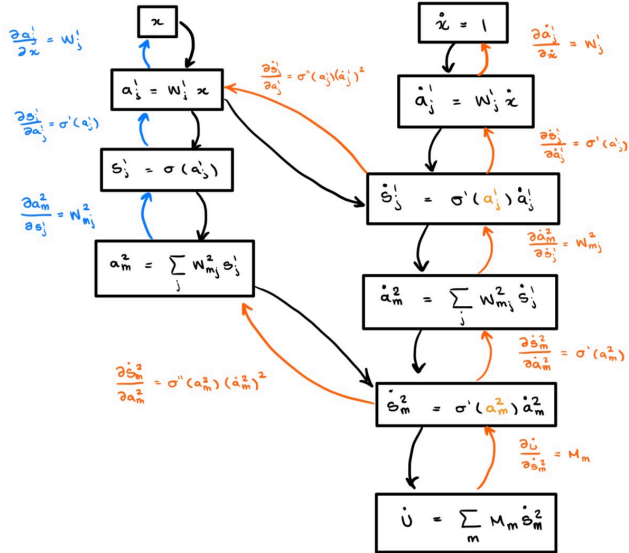
$$U' = M \cdot \sigma'(W^2 \cdot \sigma(W^1 x)) \cdot (W^2 \cdot \sigma'(W^1 x)) \cdot W^1$$

Computation is Repeated Through Chain Rule



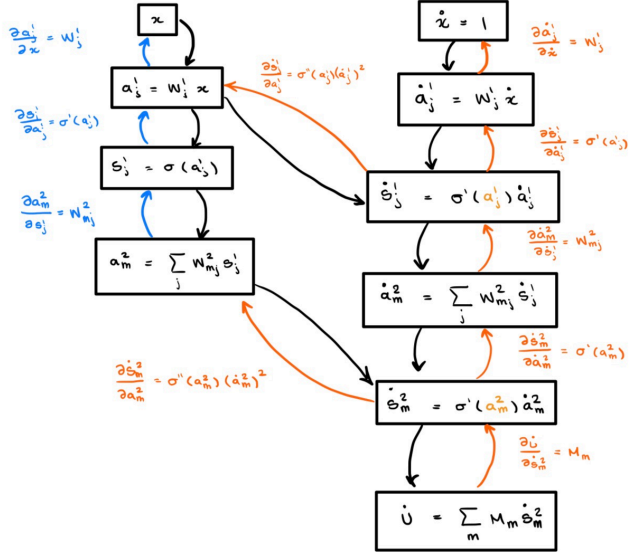
Computation is Repeated Through Chain Rule

We already computed these!

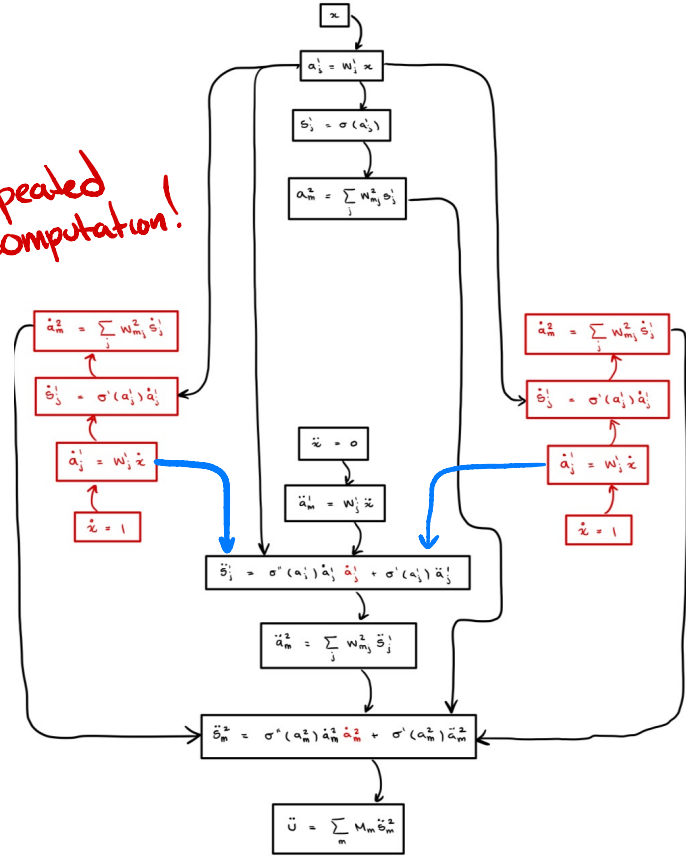


Computation is Repeated Through Chain Rule

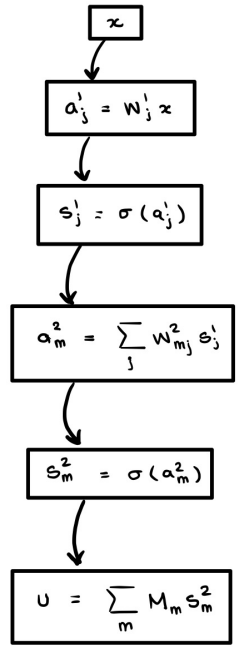
We already computed these!



Repeated Computation!

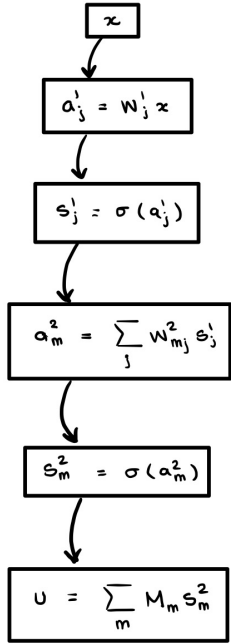


AutoDiff Graph Grows Exponentially!

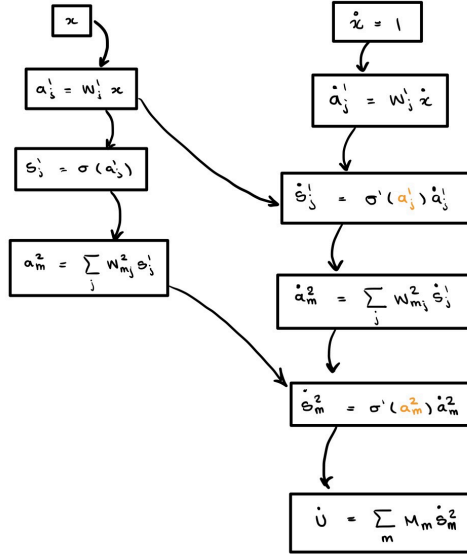


Zeroth

AutoDiff Graph Grows Exponentially!

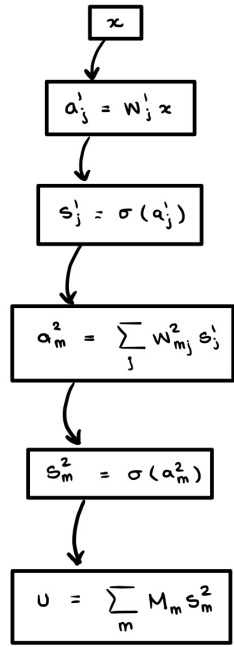


Zeroth

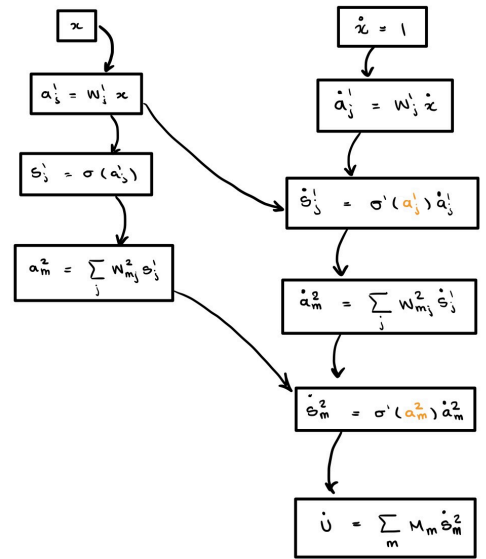


First

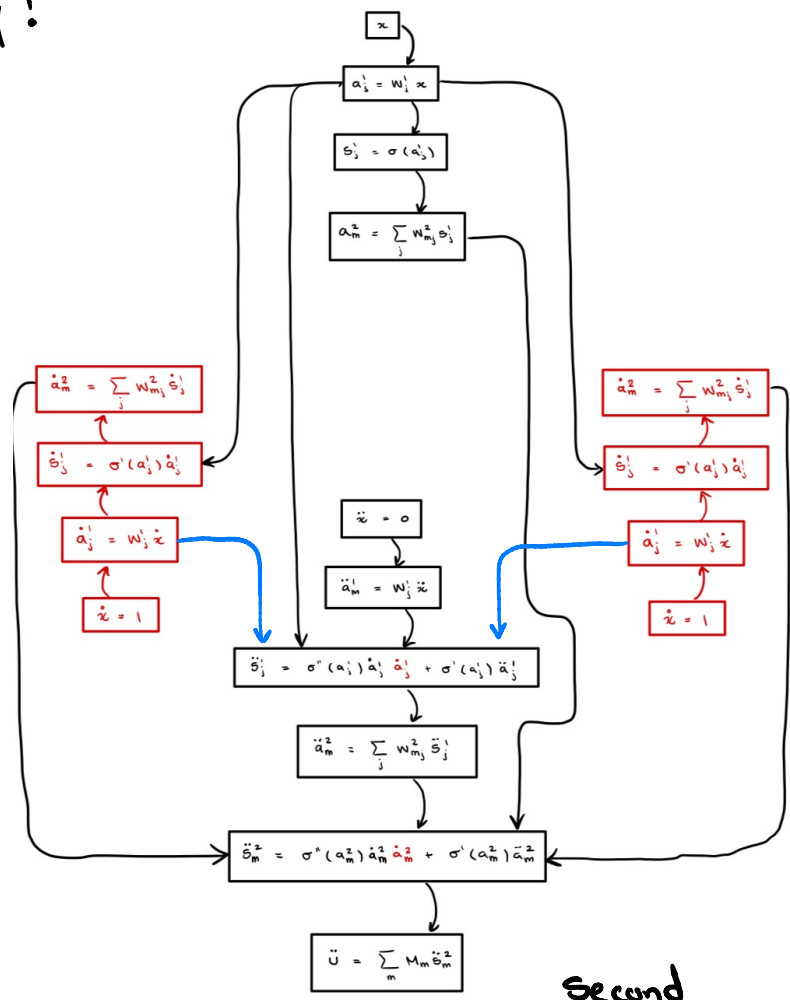
Auto diff Graph Grows Exponentially!



Zeroth

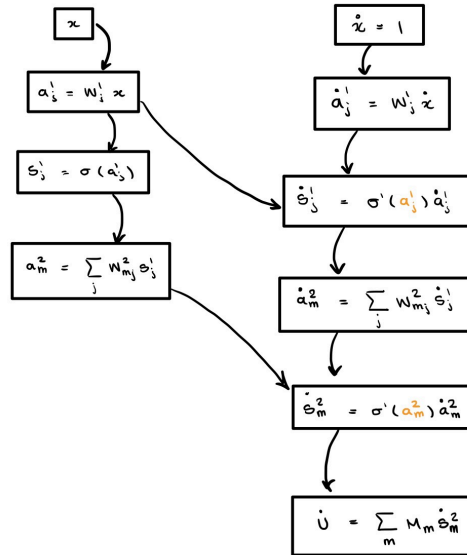
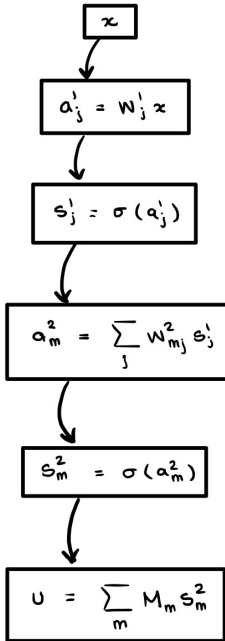


First

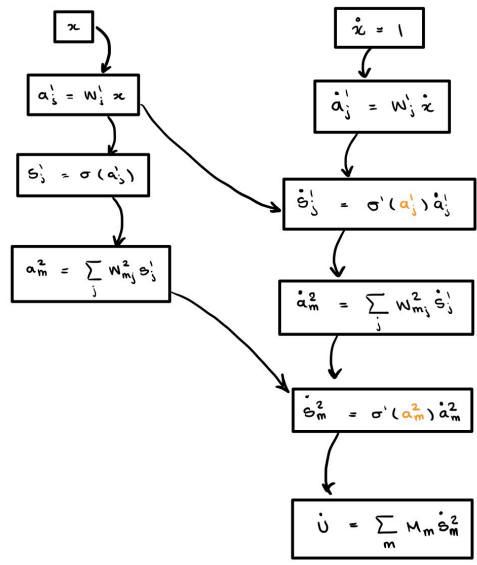
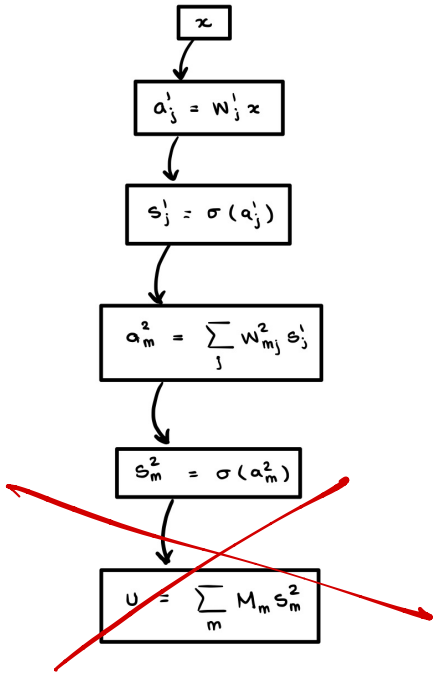


Second

Counting Nodes I

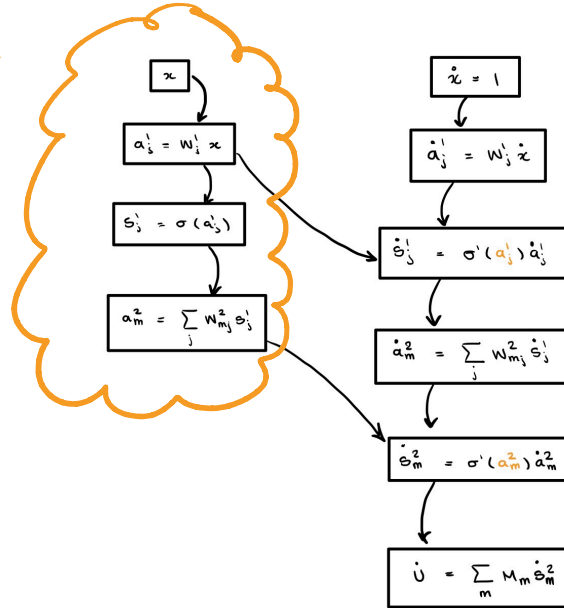
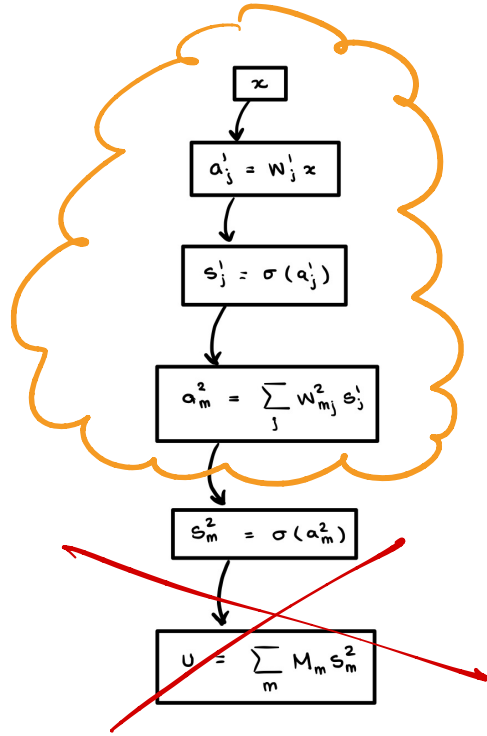


Counting Nodes I



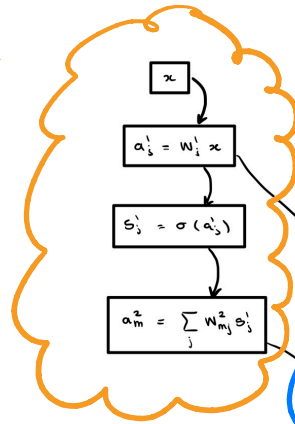
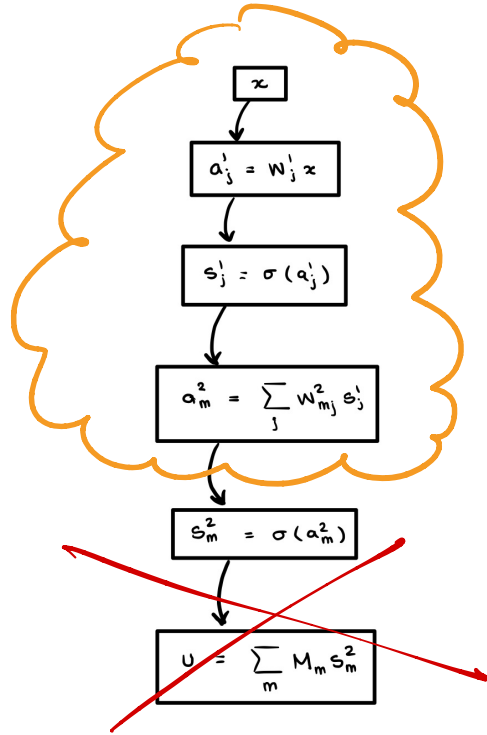
Counting Nodes I

old: 4 nodes

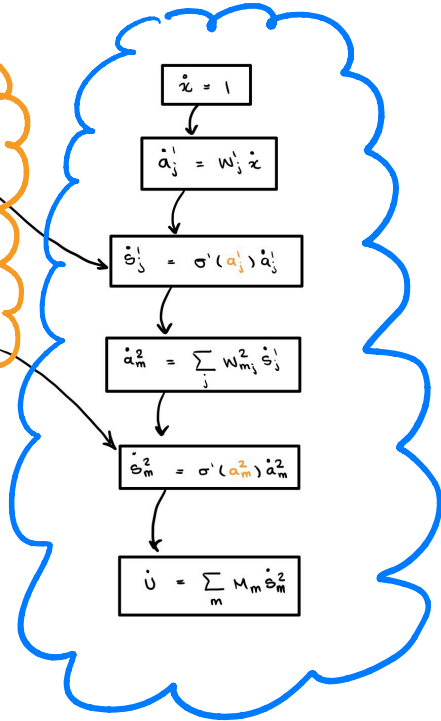


Counting Nodes I

old: 4 nodes

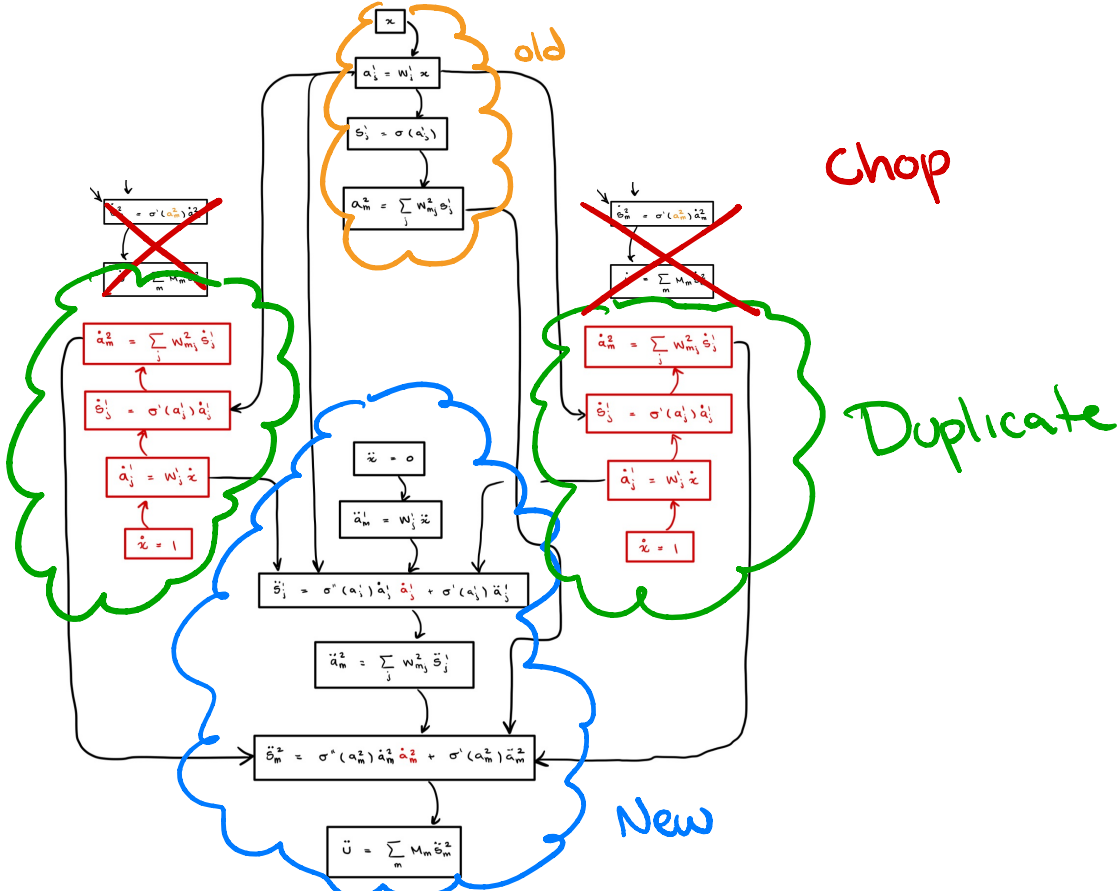


New: 6 nodes

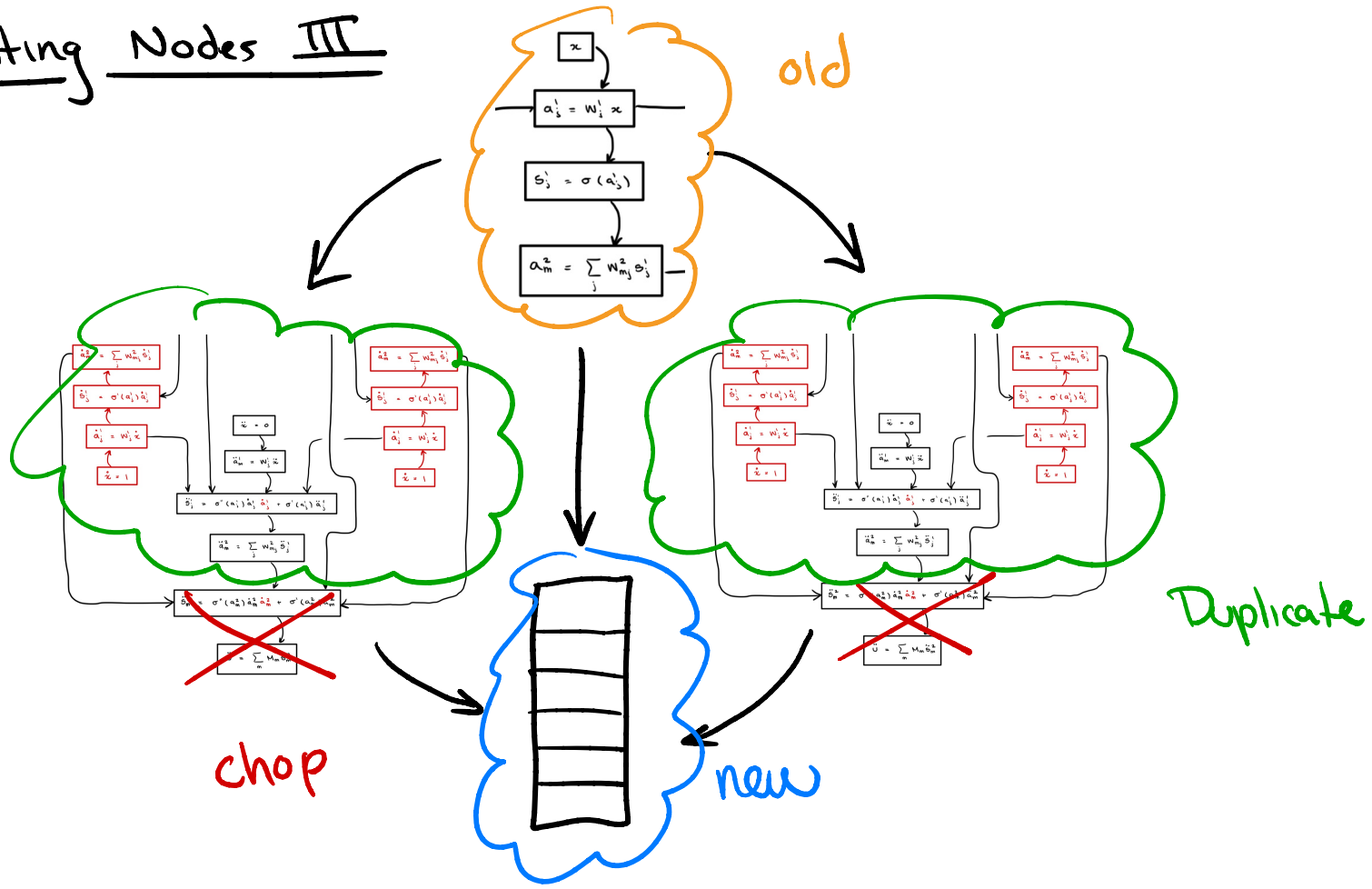


Counting Nodes II

- 6 new, 4 old, chop & duplicate the rest



Counting Nodes III



Counting Nodes III:

- $\# \text{Nodes}(n) = 10 + 2 \times (\# \text{Nodes}(n-1) - 6)$
- $\# \text{Nodes}(0) = 6$

Counting Nodes III:

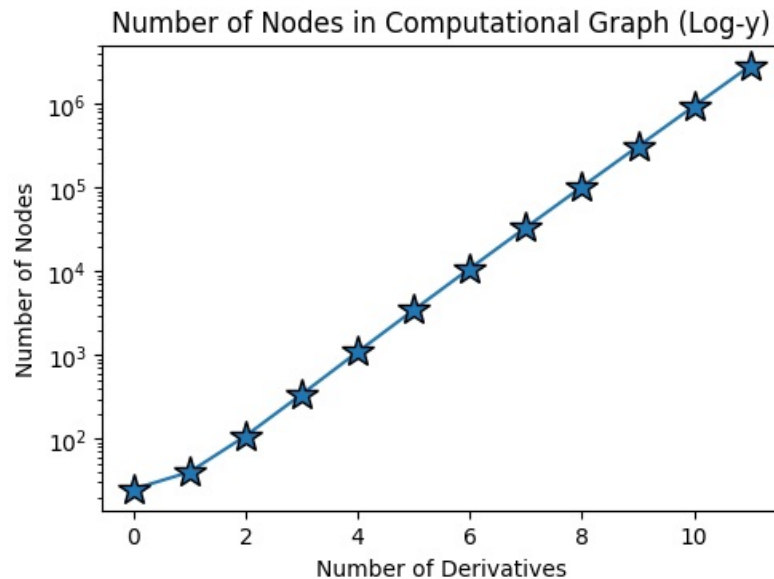
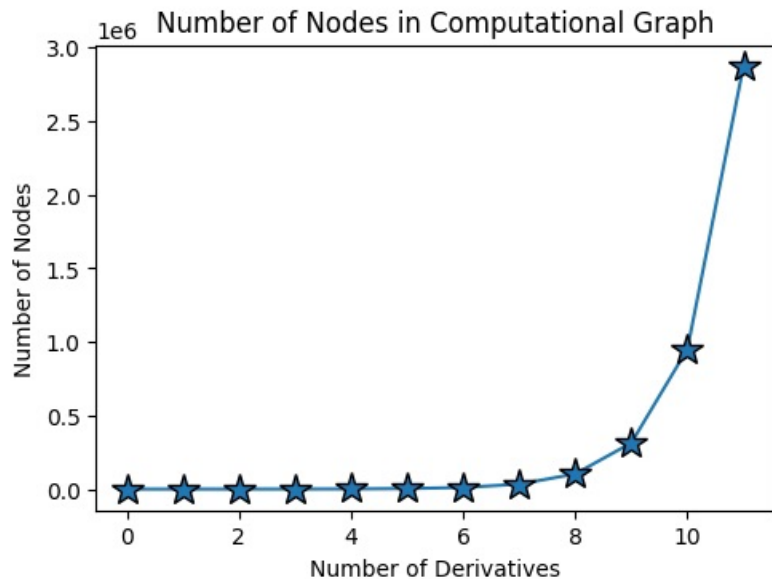
- $\# \text{ Nodes } (n) = 10 + 2 \times (\# \text{ Nodes } (n-1) - 6)$
- $\# \text{ Nodes } (0) = 6$
- $\Rightarrow \# \text{ Nodes } (n) = 2^{n+1} - 2$

Counting Nodes III:

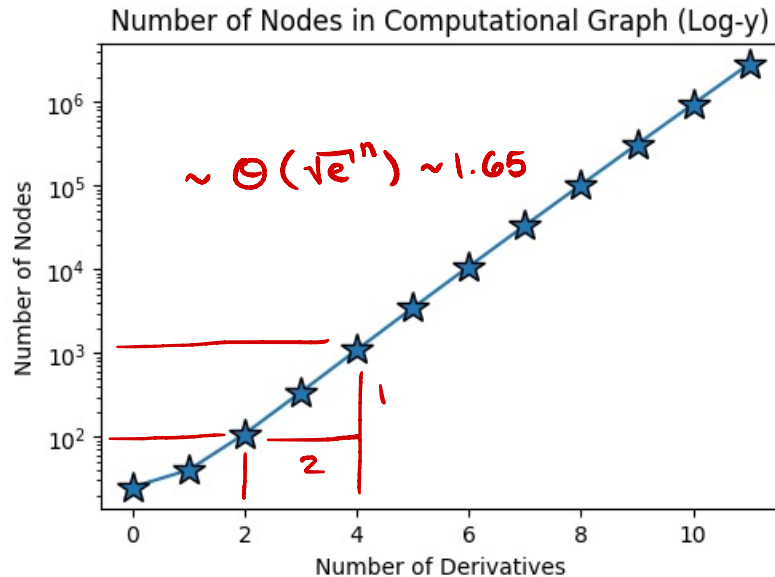
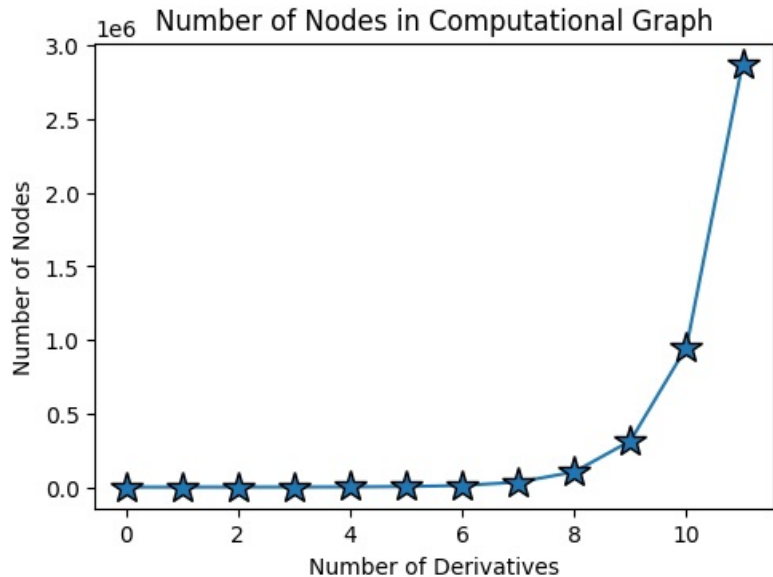
- $\# \text{Nodes}(n) = 10 + 2 \times (\# \text{Nodes}(n-1) - 6)$
- $\# \text{Nodes}(0) = 6$
- $\Rightarrow \# \text{Nodes}(n) = 2^{n+1} - 2$
- Worse if we have to evaluate all derivatives up to n

$$\sum_{k=0}^n \# \text{Nodes}(k) = 2^{n+2} - 2$$

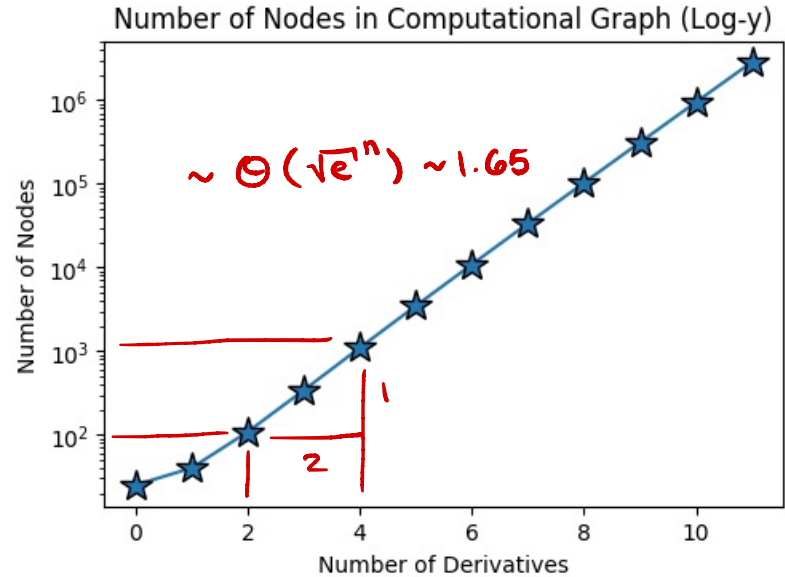
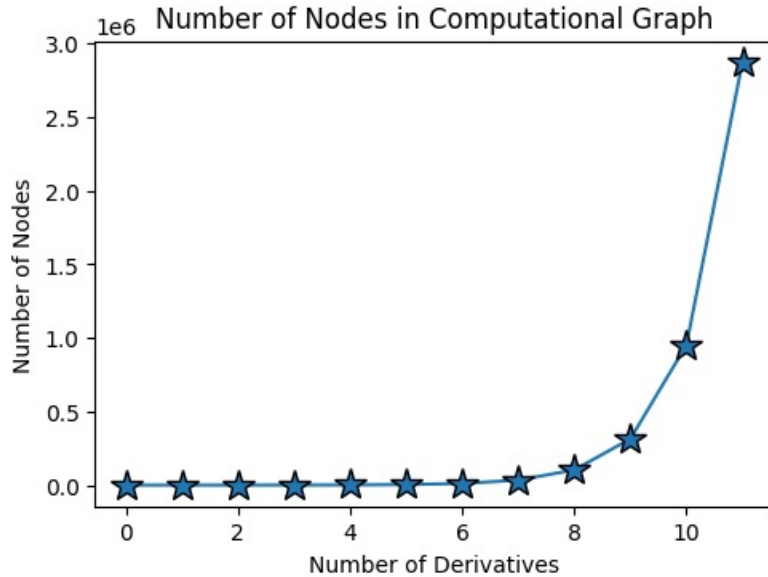
Counting Nodes V : Empirical



Counting Nodes V : Empirical

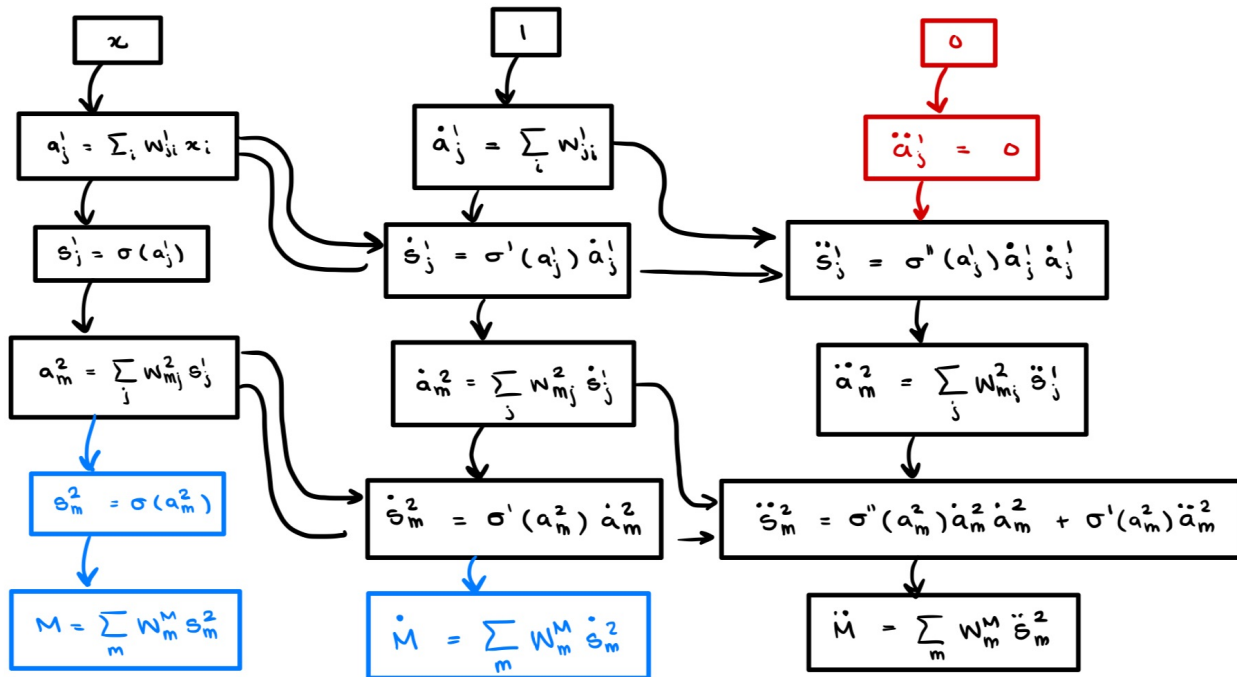


Counting Nodes V : Empirical

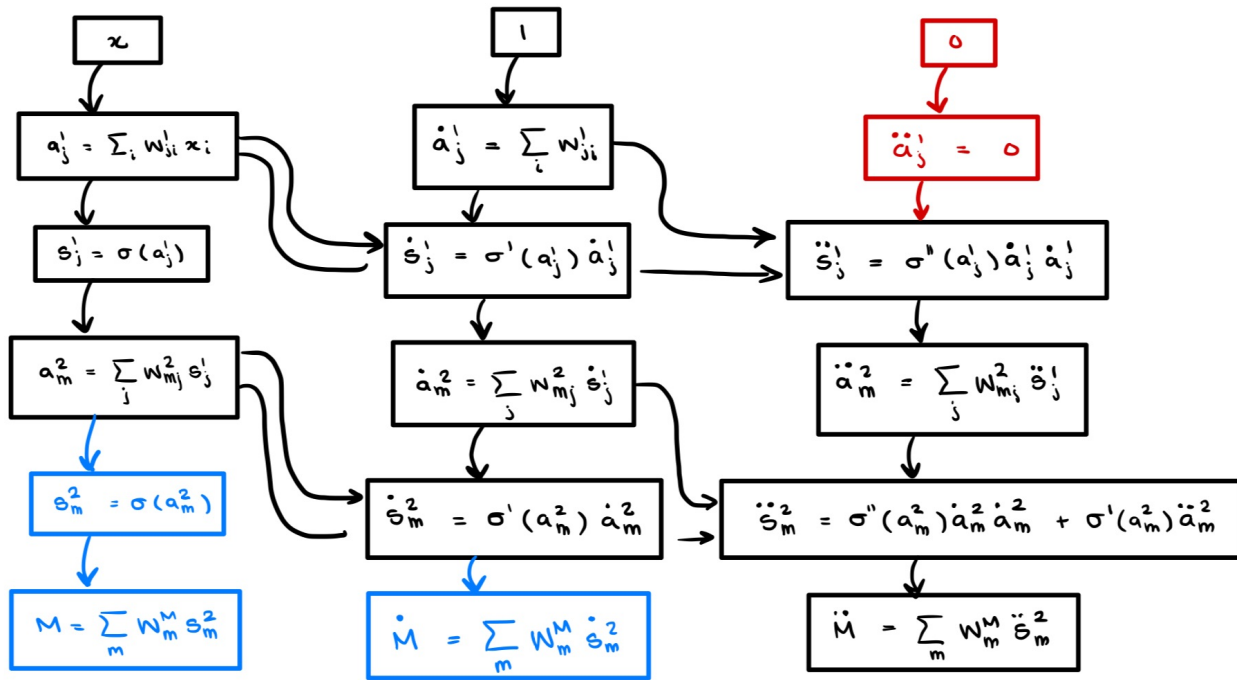


- Graph optimizations beat 2^n , but still exponential!

What Does the Minimal Graph for u'' look like?



What Does the Minimal Graph for u'' look like?



→ n-TangentProp will construct this instead!

Tangent Prop $\hat{=}$ n-Tangent Prop

Tangent Prop (Simard, Victorri, LeCun, Denker, 1991)

- Developed before Autodiff was in widespread use

Tangent Prop (Simard, Victorri, LeCun, Denker, 1991)

- Developed before Autodiff was in widespread use
- First order constraint improves training convergence (c.f. Sobolev Loss!)

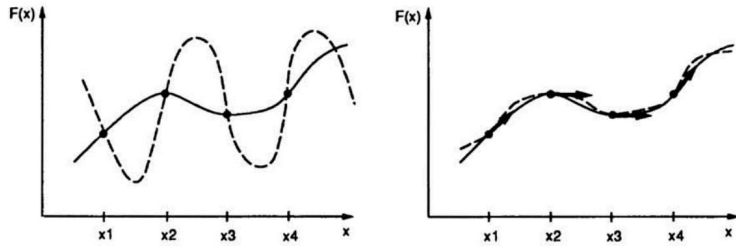


Figure 2: Learning a given function (solid line) from a limited set of example (x_1 to x_4). The fitted curves are shown in dotted line. Top: The only constraint is that the fitted curve goes through the examples. Bottom: The fitted curves not only goes through each examples but also its derivatives evaluated at the examples agree with the derivatives of the given function.

From SVLCO 1991

Tangent Prop (Simard, Victorri, LeCun, Denker, 1991)

- Developed before Autodiff was in widespread use
- First order constraint improves training convergence (c.f. Sobolev Loss!)

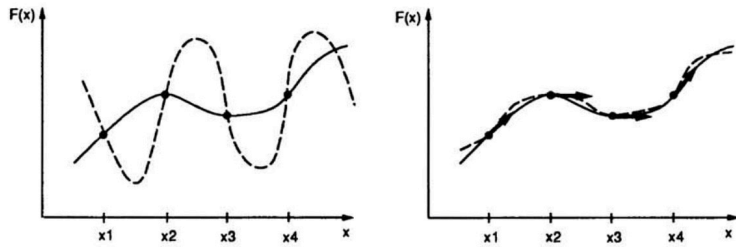


Figure 2: Learning a given function (solid line) from a limited set of example (x_1 to x_4). The fitted curves are shown in dotted line. Top: The only constraint is that the fitted curve goes through the examples. Bottom: The fitted curves not only goes through each examples but also its derivatives evaluated at the examples agree with the derivatives of the given function.

Forward

$$\bullet a_i^l = \sum_j W_{ij}^l z_j^{l-1}, \quad z_j^l = \sigma(a_j^l)$$

Derivative

$$\bullet \dot{a}_i^l = \sum_j W_{ij}^l \dot{z}_j^{l-1}, \quad \dot{z}_j^l = \sigma'(a_j^l) \dot{a}_j^l$$

From SVLCO 1991

Tangent Prop (Simard, Victorri, LeCun, Denker, 1991)

- Developed before Autodiff was in widespread use
- First order constraint improves training convergence (c.f. Sobolev Loss!)

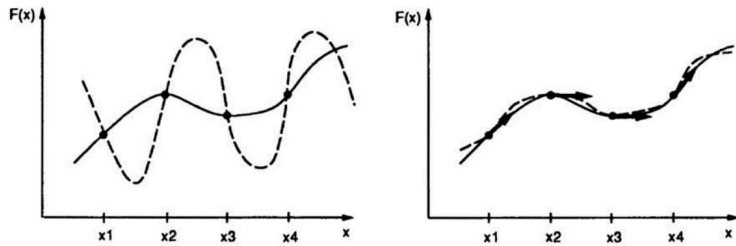


Figure 2: Learning a given function (solid line) from a limited set of example (x_1 to x_4). The fitted curves are shown in dotted line. Top: The only constraint is that the fitted curve goes through the examples. Bottom: The fitted curves not only goes through each examples but also its derivatives evaluated at the examples agree with the derivatives of the given function.

Forward

$$\bullet a_i^l = \sum_j w_{ij}^l z_j^{l-1}, \quad z_j^l = \sigma(a_j^l)$$

Derivative

$$\bullet \dot{a}_i^l = \sum_j w_{ij}^l \dot{z}_j^{l-1}, \quad \dot{z}_j^l = \sigma'(a_j^l) \dot{a}_j^l$$

- Both computed during forward

From SVLCO 1991

n-Tangent Prop I

- Leibniz formula: Product Rule

$$\frac{d^n}{dx^n} (fg) = \sum_{k=0}^n C^L(k, n) f^{(k)} g^{(n-k)}$$

n-Tangent Prop I

- Leibniz formula: Product Rule

$$\frac{d^n}{dx^n} (fg) = \sum_{k=0}^n C^L(k, n) f^{(k)} g^{(n-k)}$$

- Faà di Bruno formula: Chain Rule

$$\frac{d^n}{dx^n} (f \circ g) = \sum_{\bar{p} \in \mathcal{P}(n)} C_{\bar{p}, n}^F f^{(|\bar{p}|)} \circ g \prod_{j=1}^n (g^{(j)})^{p_j}$$

n - Tangent Prop II

- Use F_{in} to compute all derivs in a single forward pass

N-Tangent Prop II

- Use Faà to compute all derivs in a single forward pass

- Replace $\dot{a}_i^l = \sum_j W_{ij}^l \dot{z}_j^{l-1}$, $\dot{z}_j^l = \sigma'(a_j^l) \dot{a}_j^l$ with

$$\partial^m a_i^l = \sum_j W_{ij}^l \partial^m z_j^{l-1},$$

$$\partial^m z_j^l = \sum_{\bar{p} \in \mathcal{P}(m)} C_{\bar{p}, m}^F \sigma^{(|\bar{p}|)}(a_j^l) \prod_{\nu=1}^m (\partial^{\nu} a_j^l)^{\bar{p}_\nu}$$

n - Tangent Prop II

- Use $F_{\bar{a}i}$ to compute all derivs in a single forward pass

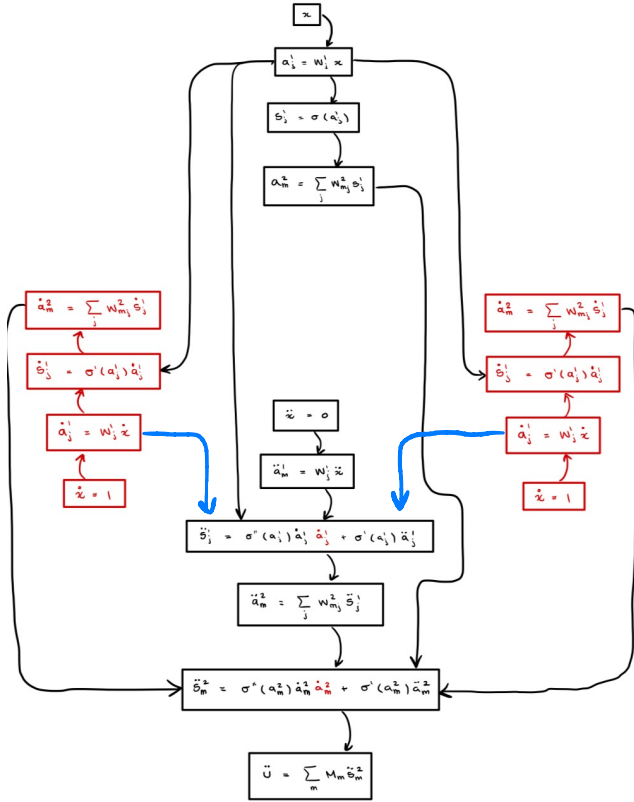
- Replace $\dot{a}_i^l = \sum_j W_{ij}^l z_j^{l-1}$, $z_j^l = \sigma'(a_j^l) \dot{a}_j^l$ with

$$\partial^m a_i^l = \sum_j W_{ij}^l \partial^m z_j^{l-1},$$

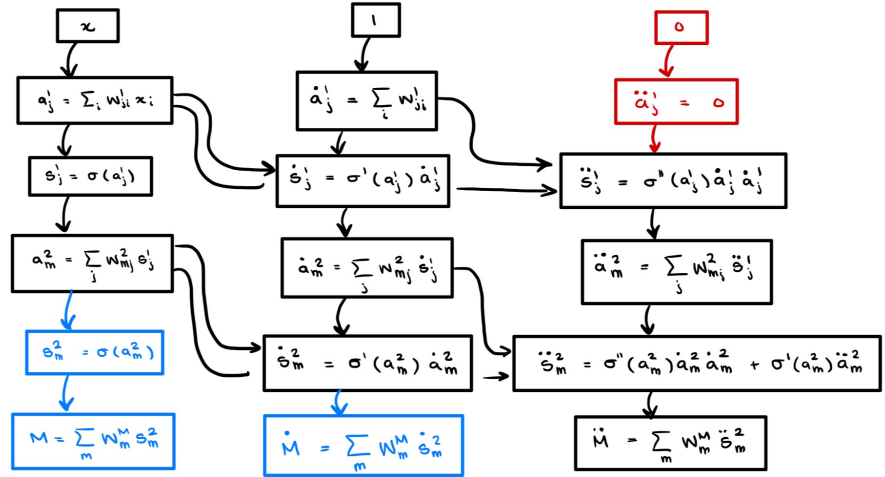
$$\partial^m z_j^l = \sum_{\bar{p} \in \mathcal{P}(m)} C_{\bar{p}, m}^F \sigma^{(|\bar{p}|)}(a_j^l) \prod_{v=1}^m (\partial^v a_j^l)^{\bar{p}_v}$$

* Forces us to compute derivatives in order!

Computational Graphs

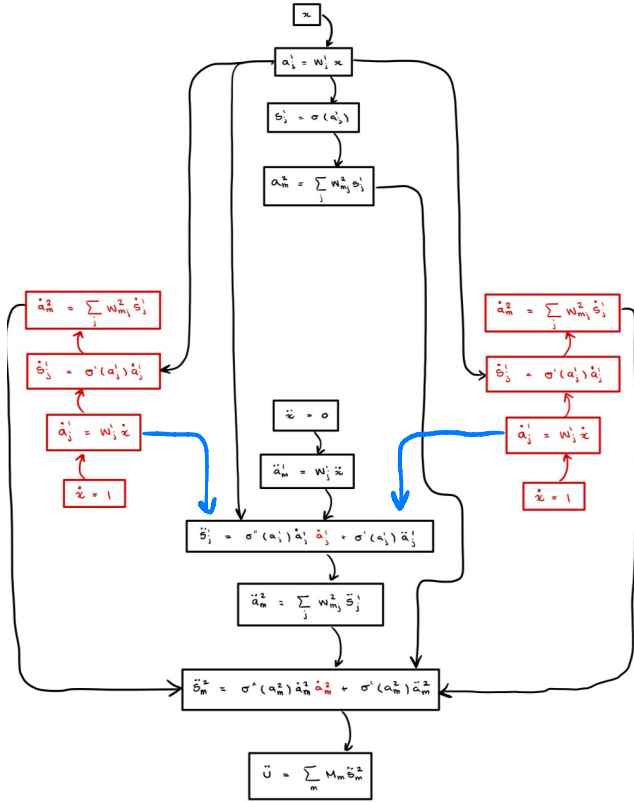


Autodiff

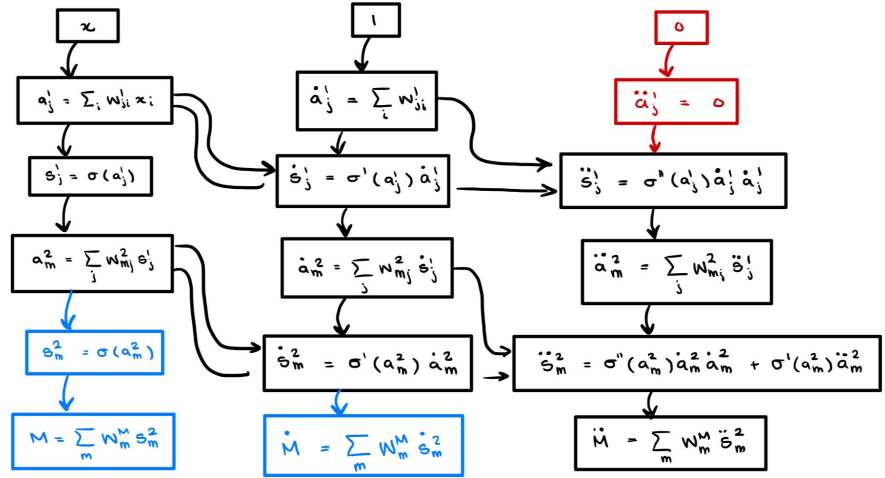


n-Tangent Prop

Computational Graphs



Autodiff



n-Tangent Prop

To compute 0,1,2 derivs:

AD: $18 + 10 + 6 = 34$ nodes

nTP: 18 nodes

Runtime Analysis

- Linear # of nodes in computation graph

Runtime Analysis

- Linear # of nodes in computation graph
- Assume original computation time is $\Theta(c)$, then

Ad: $\Theta(2^n c)$,
exponential

n-TP: $\Theta(nc)$
linear!

Runtime Analysis

- Linear # of nodes in computation graph
- Assume original computation time is $\Theta(c)$, then

Ad: $\Theta(2^n c)$,
exponential

n-TP: $\Theta(nc)$
linear!

- Ignores a lot: We will do empirical analysis below

Pedantic Runtime Analysis I

- Faia di Bruno introduces a non-trivial dependence on $|P(n)|$.

Pedantic Runtime Analysis I

- Faà di Bruno introduces a non-trivial dependence on $|P(n)|$.
- Hardy - Ramanujan 1917: $|P(n)| = \Theta\left(\frac{1}{n} e^{\sqrt{n}}\right)$

Pedantic Runtime Analysis I

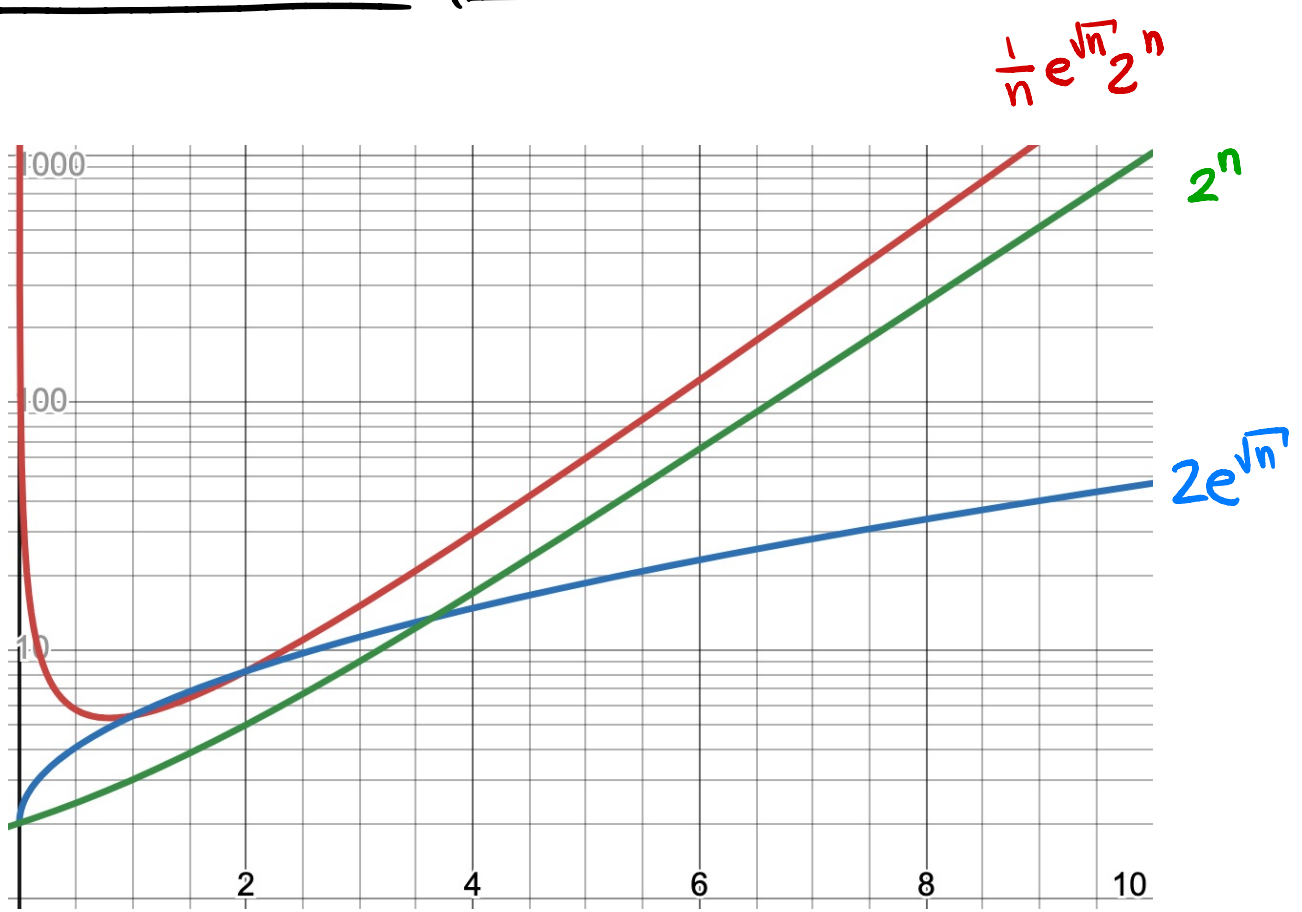
- Faà di Bruno introduces a non-trivial dependence on $|P(n)|$.
- Hardy - Ramanujan 1917: $|P(n)| = \Theta\left(\frac{1}{n} e^{\sqrt{n}}\right)$
- But this will appear in both n-Tangent Prop's Auto diff!

Pedantic Runtime Analysis I

- Faà di Bruno introduces a non-trivial dependence on $|P(n)|$.
- Hardy - Ramanujan 1917: $|P(n)| = \Theta\left(\frac{1}{n} e^{\sqrt{n}}\right)$
- But this will appear in both n-Tangent Prop's Auto diff!

$$Ad = \Theta\left(\frac{e^{\sqrt{n}}}{n} Q^n\right), \quad nTP = \Theta\left(e^{\sqrt{n}}\right)$$

Pedantic Runtime Analysis II



Results

Results: Neural Network I

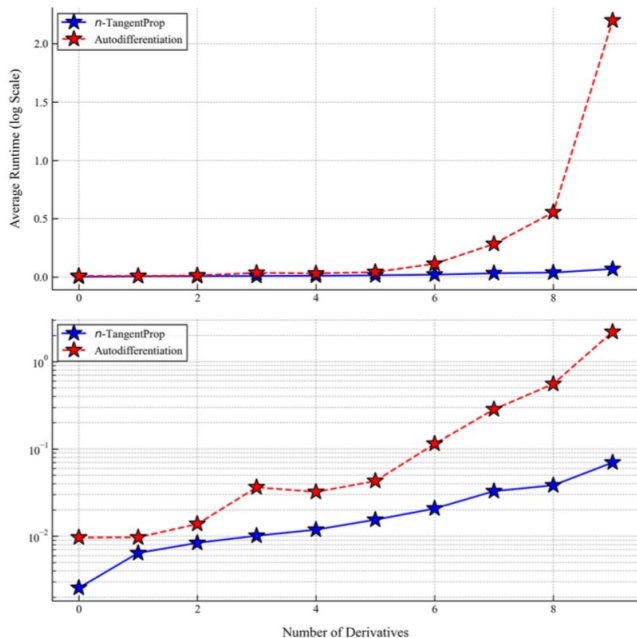


Fig. 1. Average runtime for a combined forward and backwards pass using autodifferentiation (red) and n -TangentProp (blue). The top and bottom frames show the same data, however the bottom frame is plotted with a logarithmic y -axis. Each model is run 100 times and the average for each trial is plotted. The network has 3 hidden layers of 24 neurons each, a common PINN architecture. The batch size is $2^8 = 256$ samples. The forward and backwards pass times are shown separately in Figures [2](#) and [3](#) respectively.

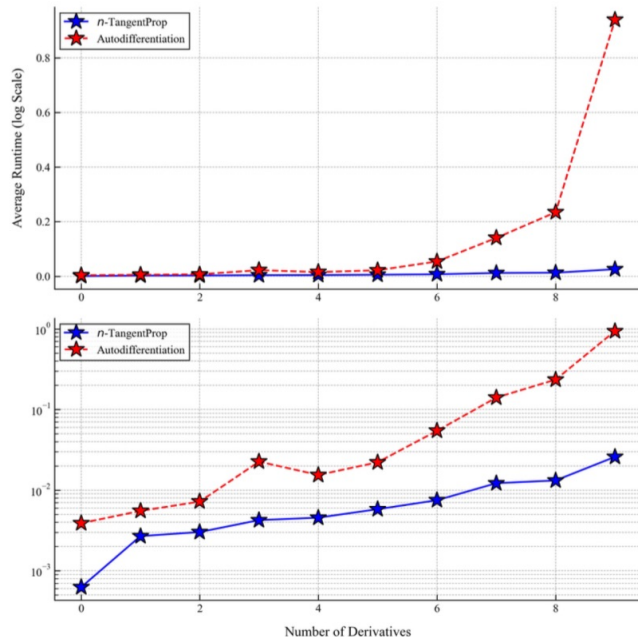
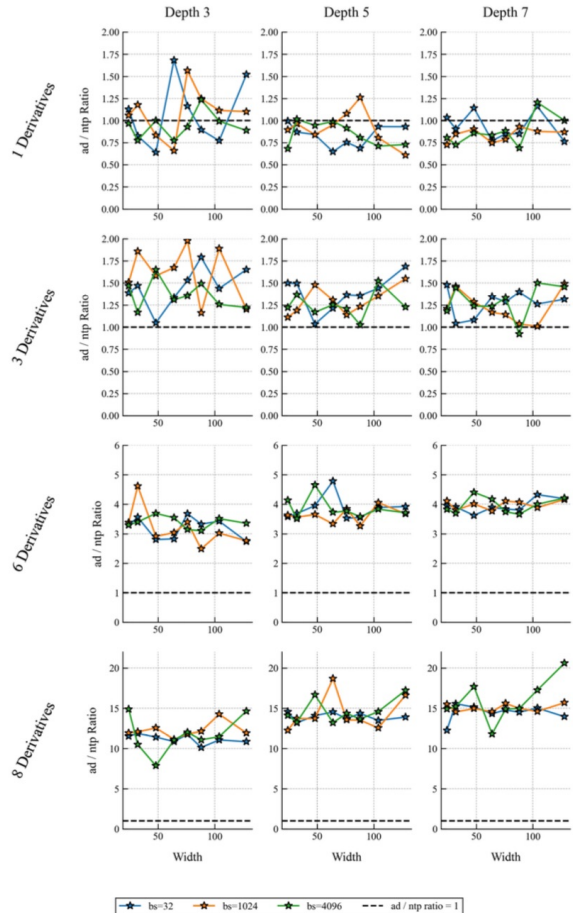


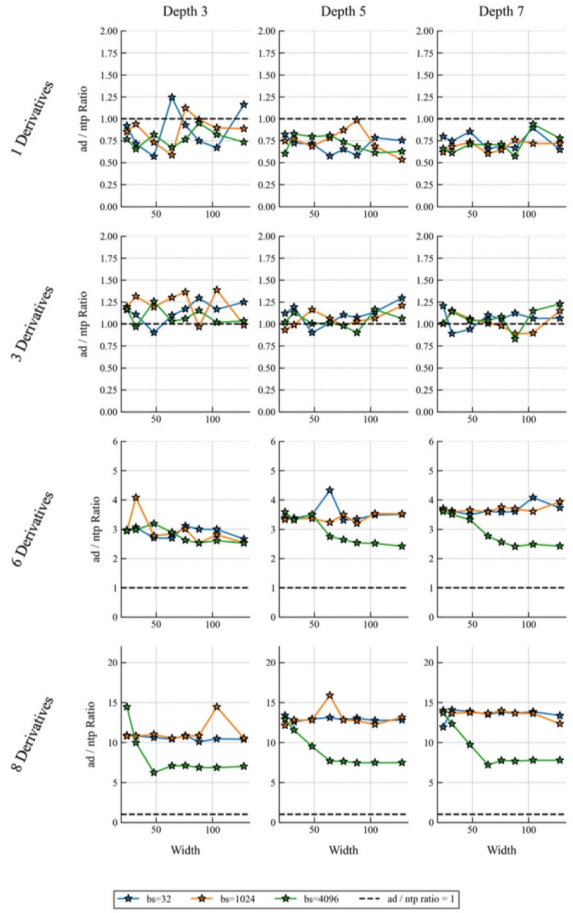
Fig. 2. Forward pass times for the model shown in Figure [1](#). The top and bottom frames show the same data, however the bottom frame is plotted with a logarithmic y -axis. Each model is run 100 times and the average for each trial is plotted. The network has 3 hidden layers of 24 neurons each, a common PINN architecture. The batch size is $2^8 = 256$ samples.

Results: Neural Networks II

Forward



Combined



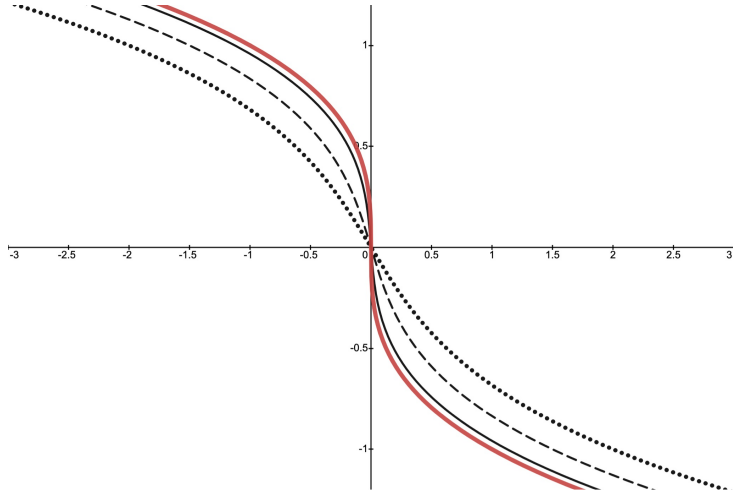
★ bs=32
 ★ bs=1024
 ★ bs=4096
 - - - ad / ntp ratio = 1

★ bs=32
 ★ bs=1024
 ★ bs=4096
 - - - ad / ntp ratio = 1

Self-Similar Burgers I

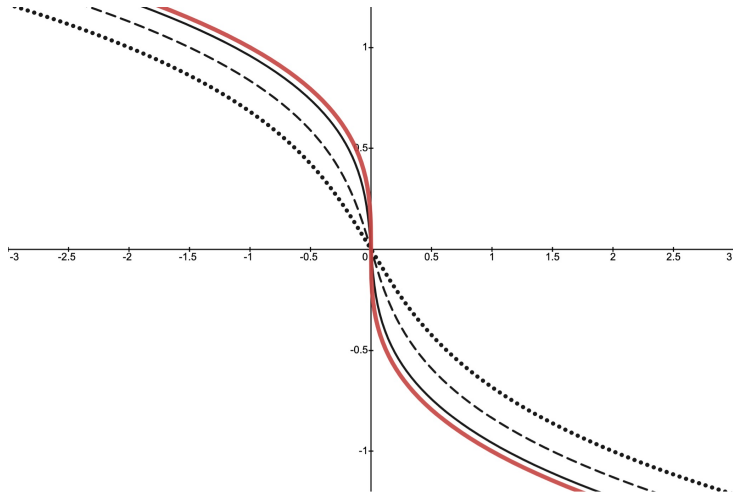
- Burgers : $\partial_t U + U \partial_x U = 0$

Self-Similar Burgers I



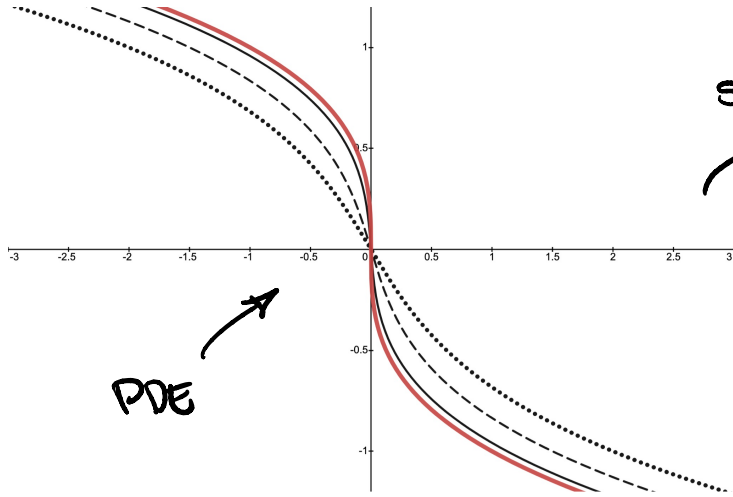
- Burgers : $\partial_t U + U \partial_x U = 0$
- shocks : $|\partial_x U| \rightarrow \infty$ in finite time

Self-Similar Burgers I

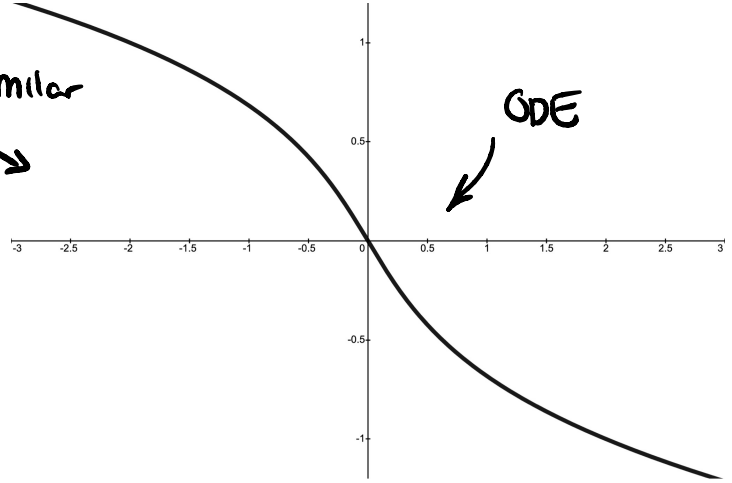


- Burgers : $\partial_t U + U \partial_x U = 0$
- Shocks : $|\partial_x U| \rightarrow \infty$ in finite time
- This makes analysis hard

Self-Similar Burgers I



self-similar



- Burgers : $\partial_t U + U \partial_x U = 0$
- shocks : $|\partial_x U| \rightarrow \infty$ in finite time
- This makes analysis hard

Self-Similar Burgers II

- Self-Similar transform takes advantage of Power law relation between spatial & temporal scales

Self-Similar Burgers II

- Self-Similar transform takes advantage of Power law relation between spatial & temporal scales

- $X = \frac{x}{(1-t)^{1+\lambda}}$, $U(x,t) = (1-t)^\lambda W(X)$

$$\Rightarrow -\lambda W + ((1+\lambda)X + W)W' = 0$$

Self-Similar Burgers II

- Self-Similar transform takes advantage of Power law relation between spatial & temporal scales

$$\bullet \quad X = \frac{x}{(1-t)^{1+\lambda}}, \quad u(x,t) = (1-t)^\lambda W(X)$$

$$\Rightarrow -\lambda W + ((1+\lambda)X + W)W' = 0$$

- ODE in X , Allows us to study solution AT the shock, since X^{-1} captures the shock formation dynamics

Self-Similar Burgers II

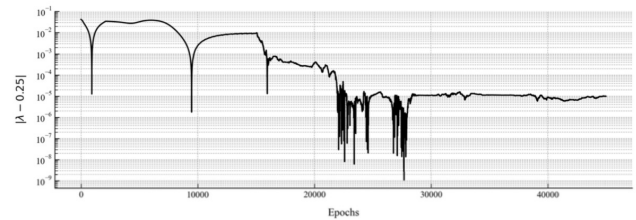
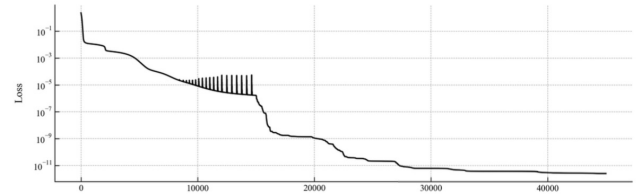
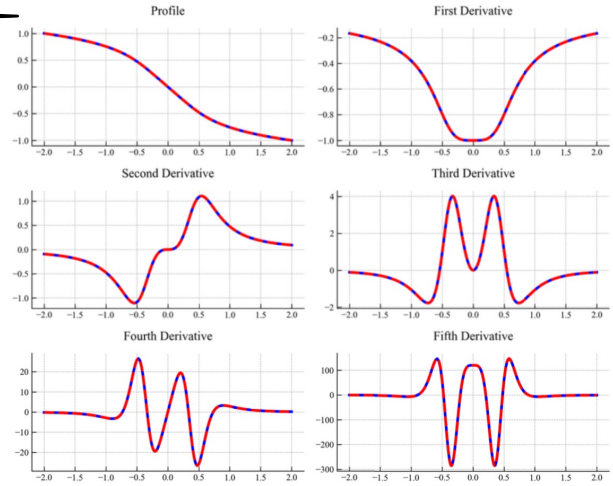
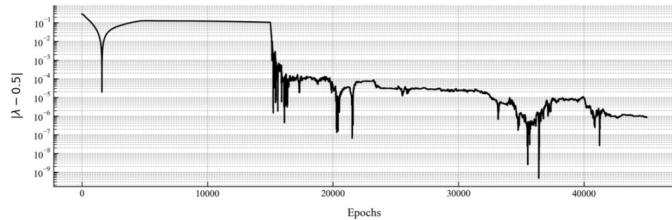
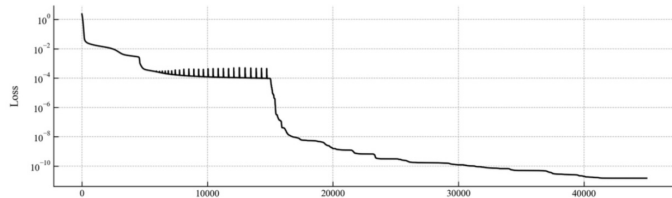
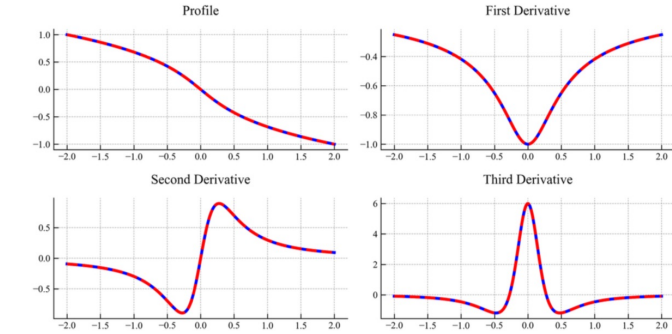
- Self-Similar transform takes advantage of Power law relation between spatial & temporal scales

- $$X = \frac{x}{(1-t)^{1+\lambda}}, \quad u(x,t) = (1-t)^\lambda W(X)$$

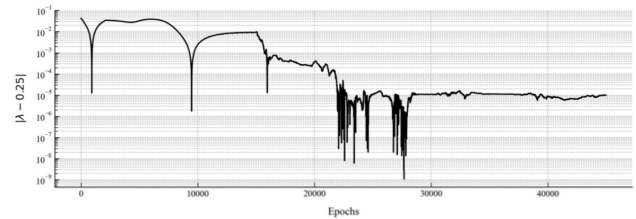
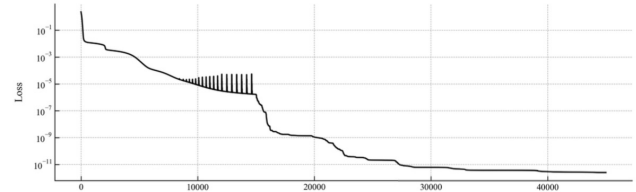
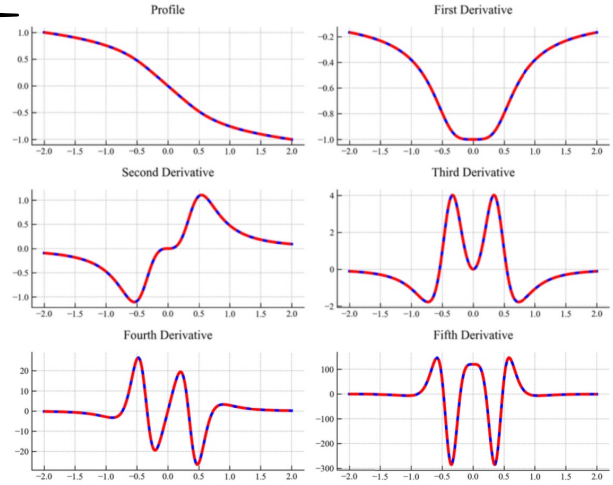
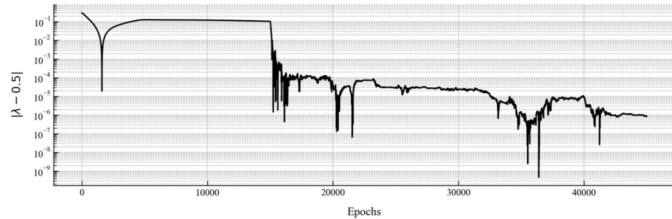
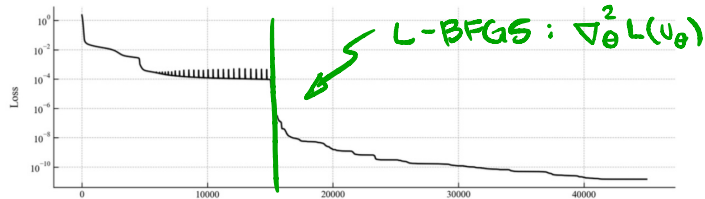
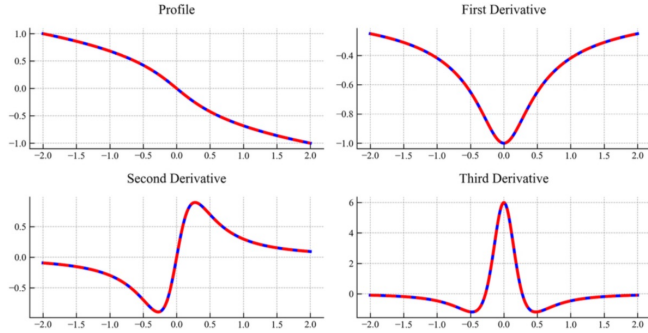
$$\Rightarrow -\lambda W + ((1+\lambda)X + W)W' = 0$$

- ODE in X , Allows us to study solution AT the shock, since X^{-1} captures the shock formation dynamics
- PINN: Solve for W, λ simultaneously

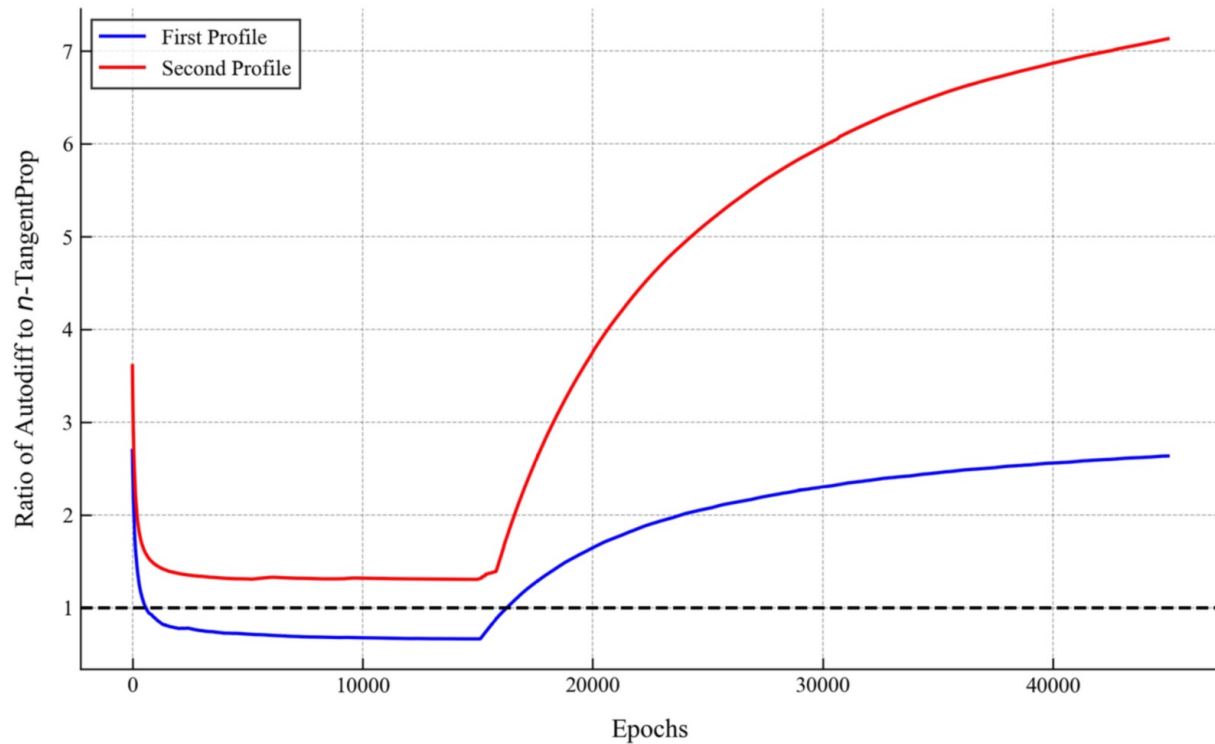
Results: Self-Similar Burgers I



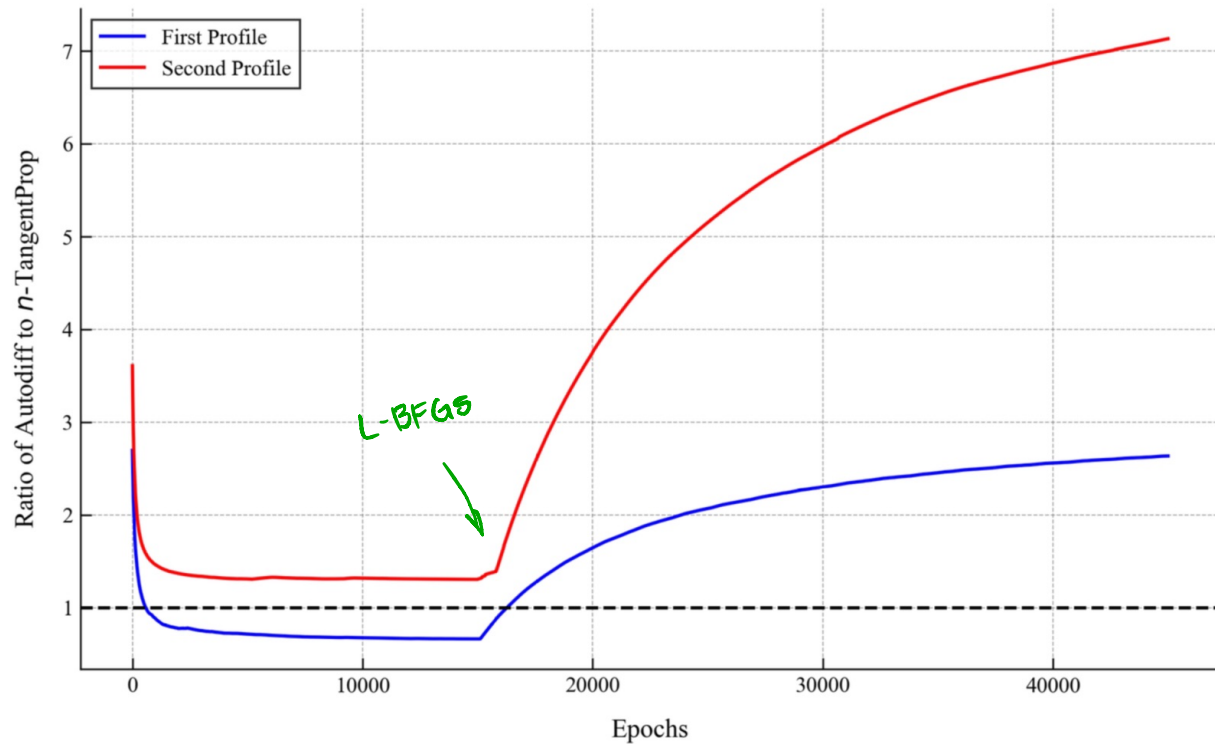
Results: Self-Similar Burgers I



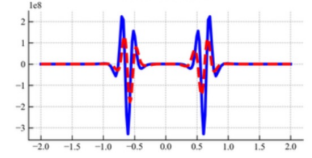
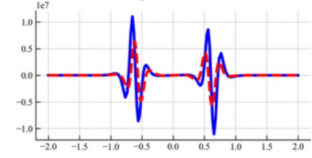
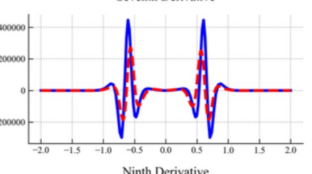
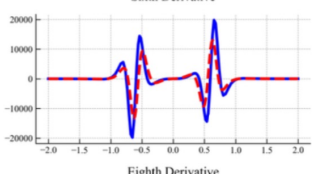
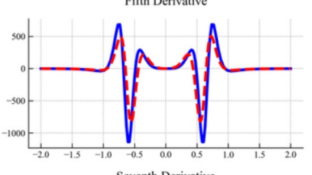
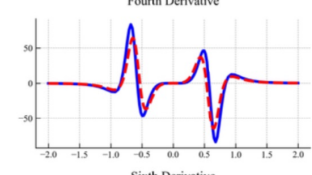
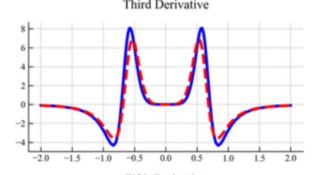
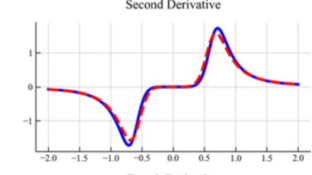
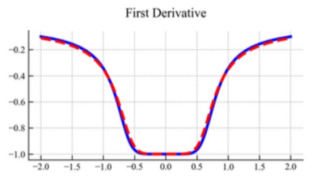
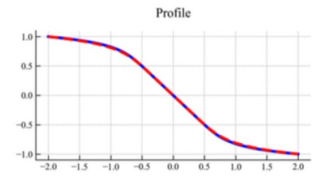
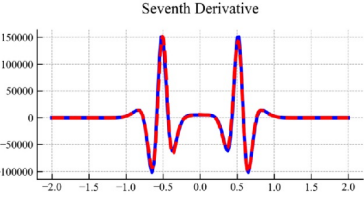
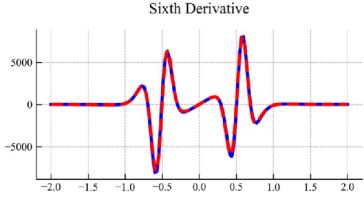
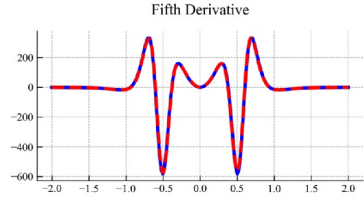
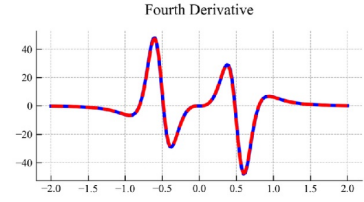
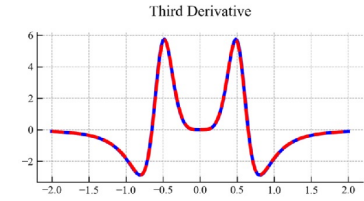
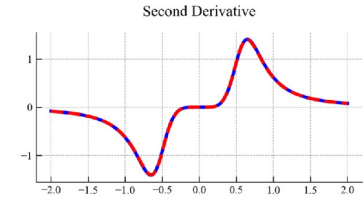
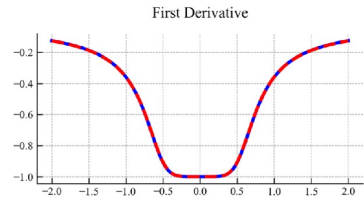
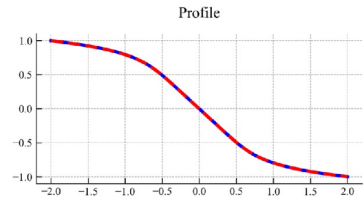
Results: Self-Similar Burgers II



Results: Self-Similar Burgers II

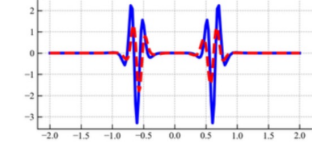
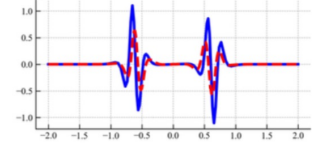
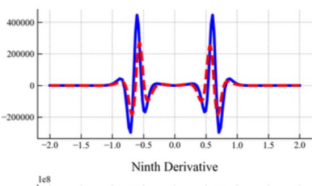
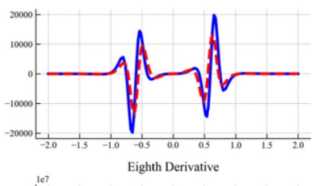
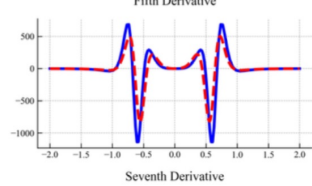
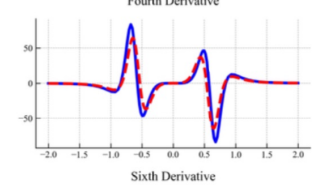
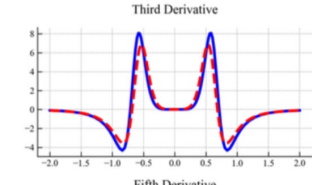
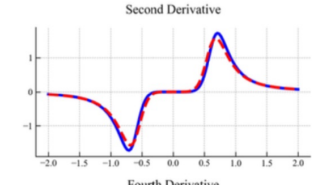
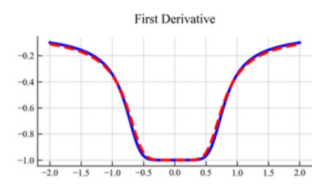
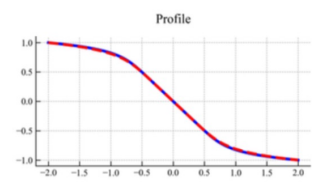
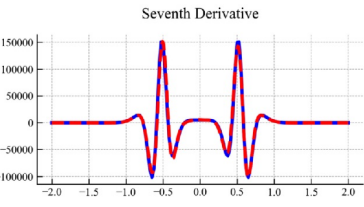
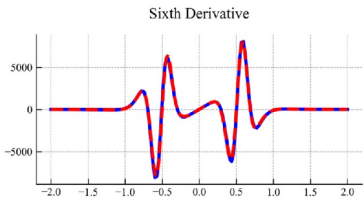
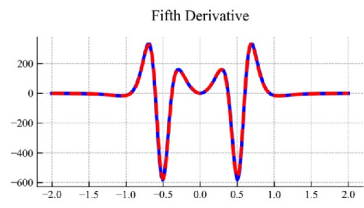
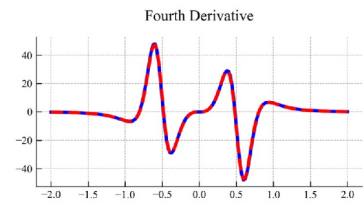
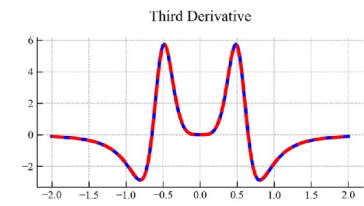
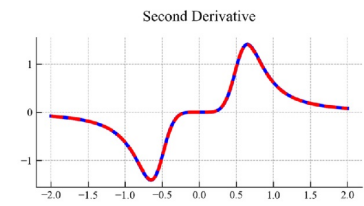
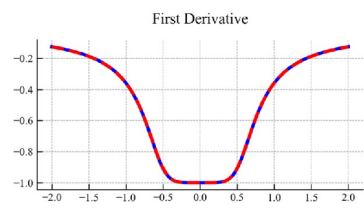
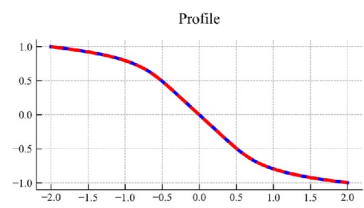


Results : Self-Similar Burgers III



Results : Self-Similar Burgers III

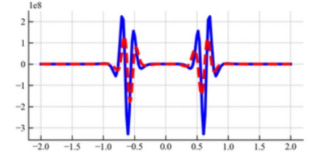
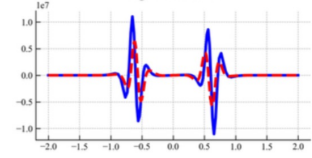
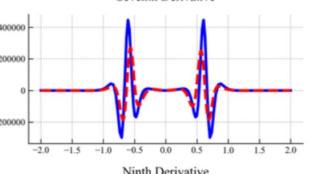
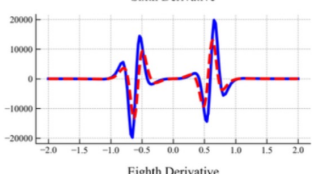
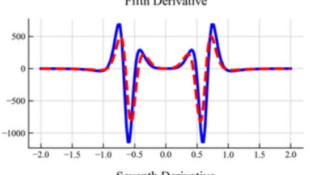
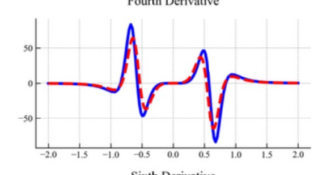
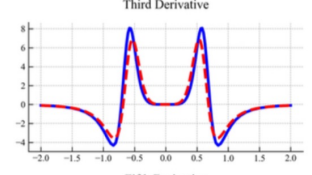
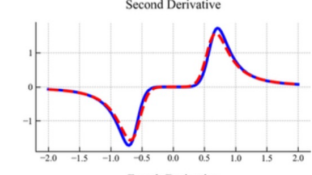
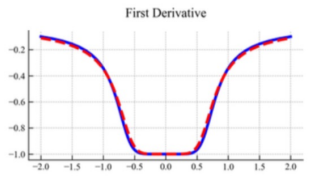
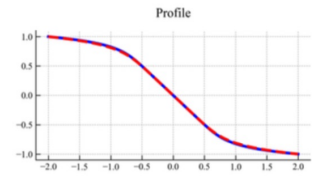
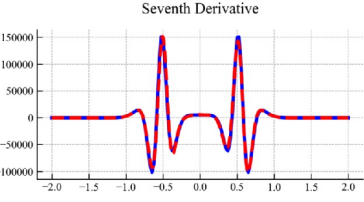
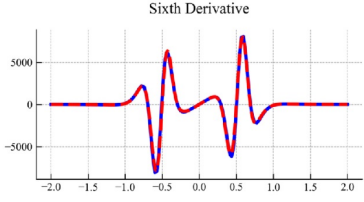
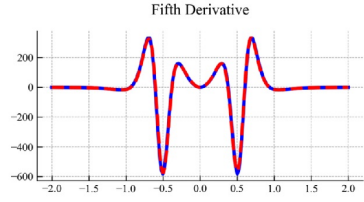
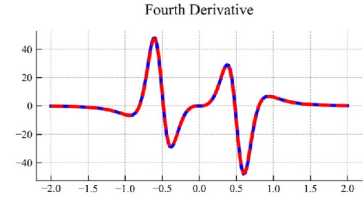
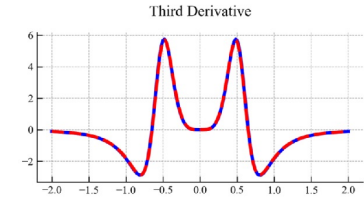
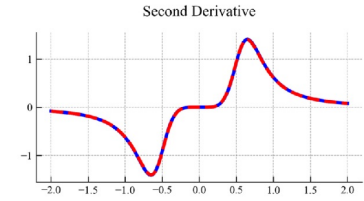
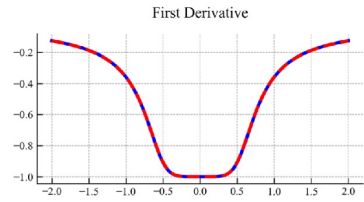
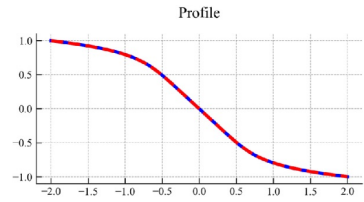
26x Faster



Results : Self-Similar Burgers III

26x Faster

65x Faster



Conclusions

Conclusions

- Use Sobolev Loss + n-TP for more efficient training

Conclusions

- Use Sobolev Loss + n-TP for more efficient training
- Faster training, less compute

Conclusions

- Use Sobolev Loss + n-TP for more efficient training
- Faster training, less compute
- Democratization of PINN training

Conclusions

- Use Sobolev Loss + n-TP for more efficient training
- Faster training, less compute
- Democratization of PINN training
- Multiprop, other applications?

Review

- Neural Networks → ANN

Review

- Neural Networks \rightarrow PINN
- Autodiff is exponential in # of derivatives

Review

- Neural Networks \rightarrow PINN
- Autodiff is exponential in # of derivatives
- n-Tangent Prop computes minimal comp. graph

Review

- Neural Networks \rightarrow PINN
- Autodiff is exponential in # of derivatives
- n-Tangent Prop computes minimal comp. graph
- More accurate solutions, $\sim 3x$ faster for simple problem, $65+x$ faster on harder problem.



Thank You!

