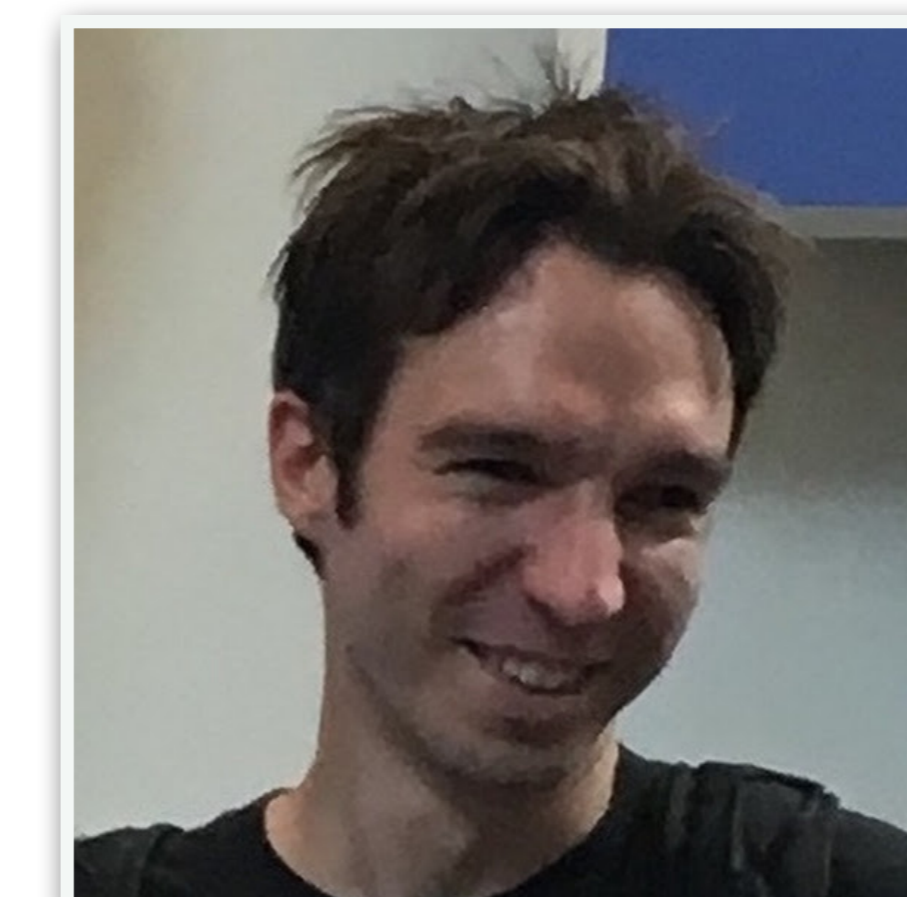
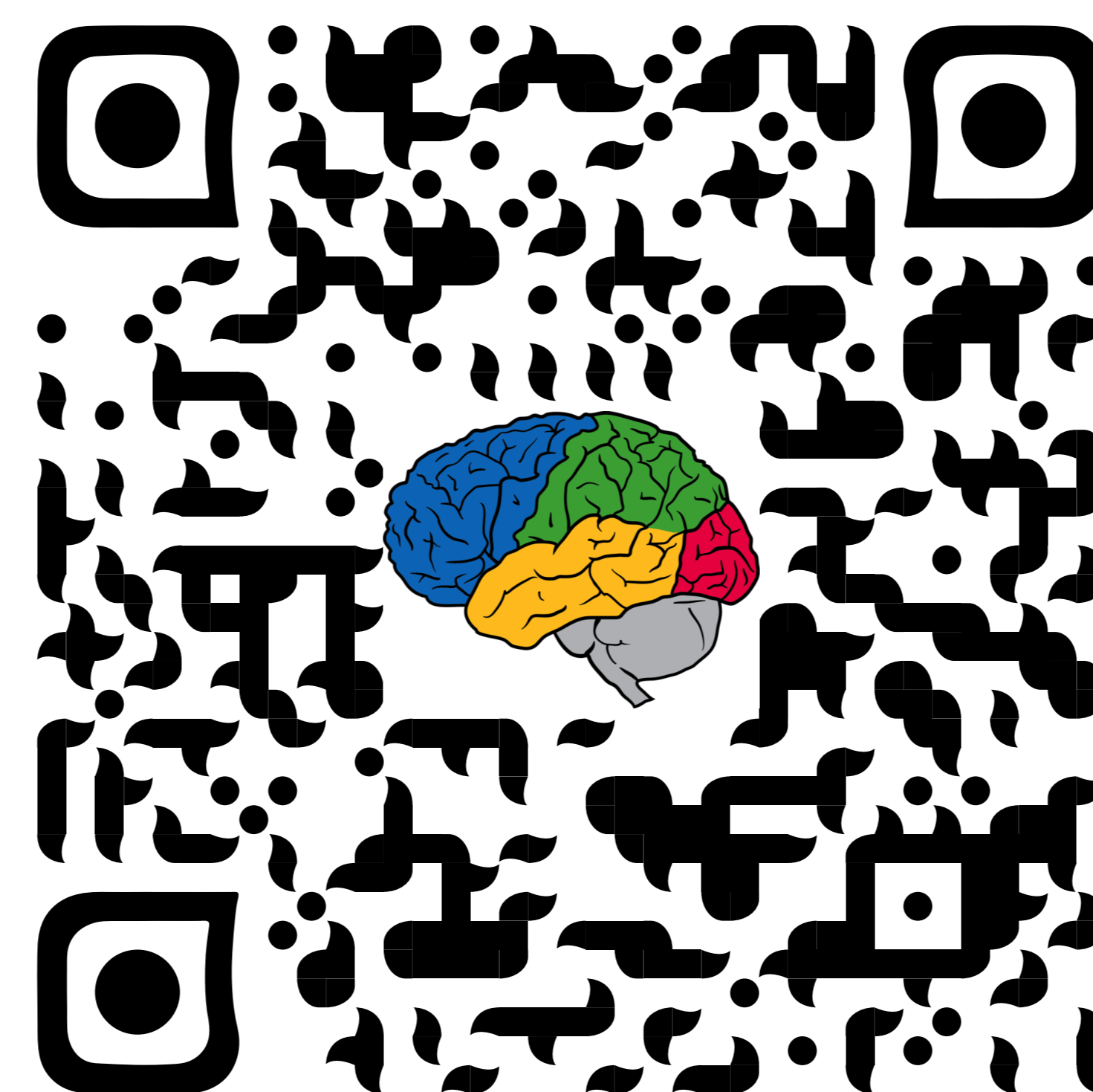


Stop the war
Stand with Ukraine

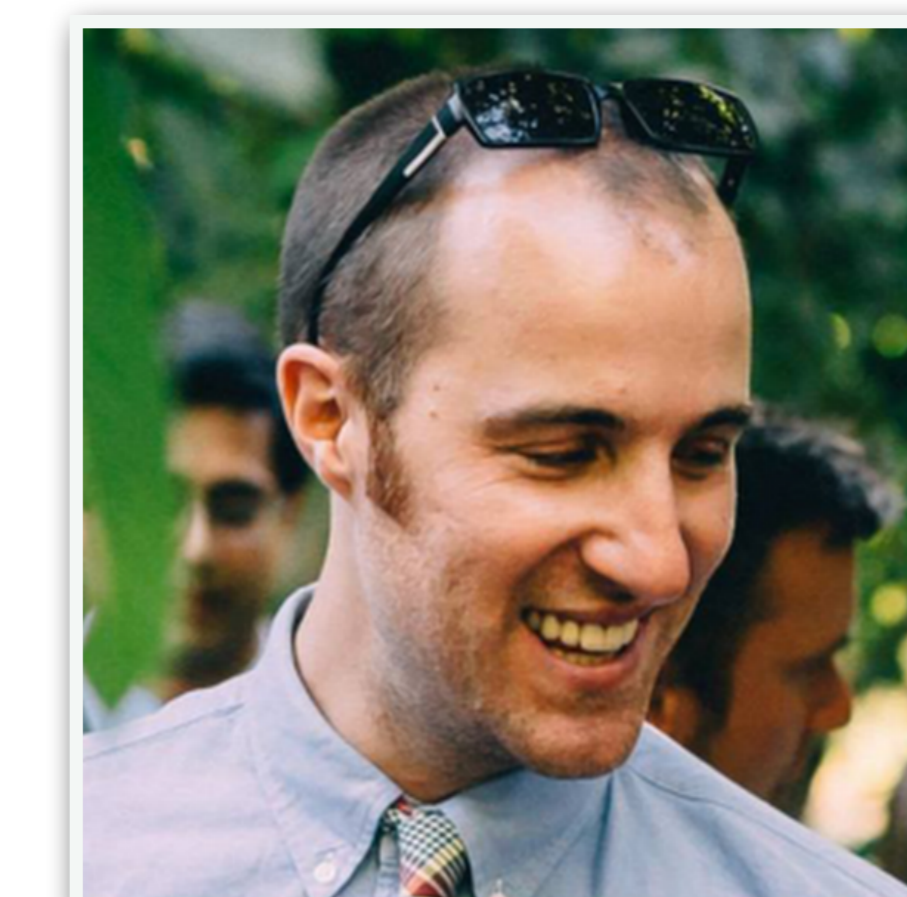
Fast Finite Width (Neural) Tangent Kernel (NTK)

ICML 2022

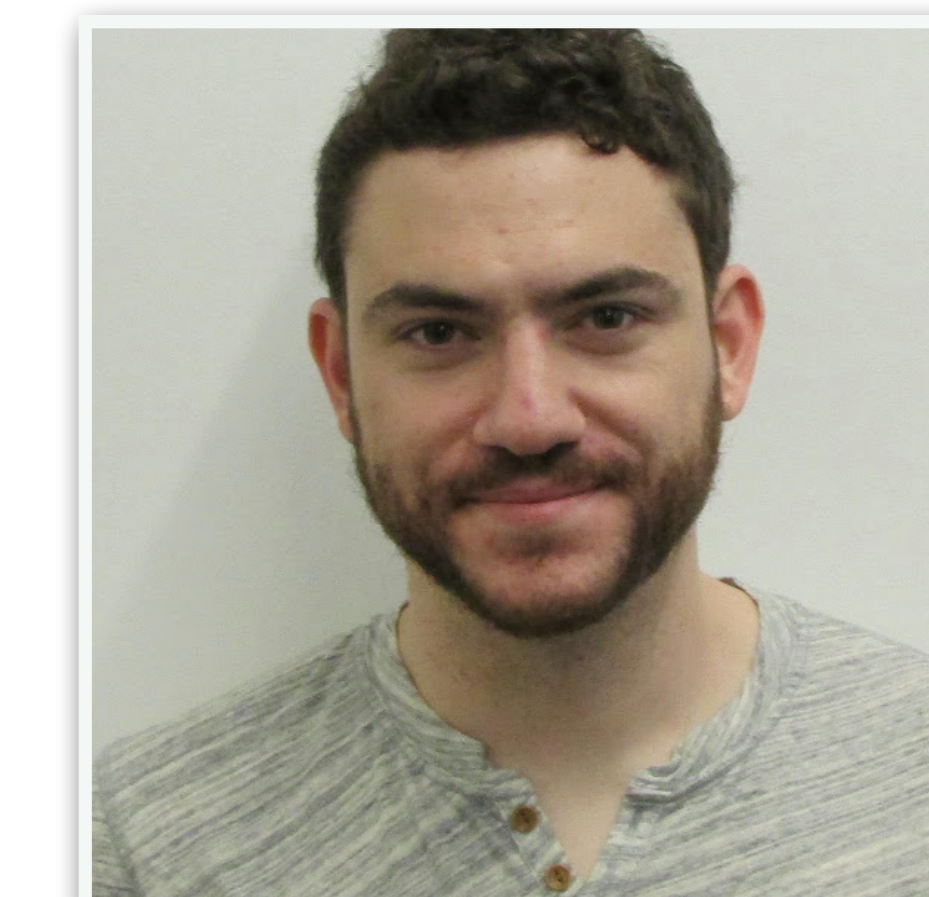
$$\frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T$$



[Roman Novak](#)



[Jascha Sohl-Dickstein](#)



[Samuel S. Shoenholz](#)

June 2022

Presentation by Roman Novak

github.com/google/neural-tangents

Plan

1. **NTK definition and notation**

2. Recap of Automatic Differentiation (AD)

3. Baseline NTK computation complexity

4. Two new algorithms for computing the NTK

5. Benchmarks

Definition

Neural network $f(\theta, x) \in \mathbb{R}^O$ Parameters $\theta \in \mathbb{R}^P$

Definition

Neural network $f(\theta, x) \in \mathbb{R}^O$ Parameters $\theta \in \mathbb{R}^P$

The NTK is the outer product of Jacobians:

$$\Theta_{\theta}(x_1, x_2) := \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T$$

Definition

Neural network $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$ Parameters $\theta \in \mathbb{R}^{\mathbf{P}}$

The NTK is the outer product of Jacobians:

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{\mathbf{O} \times \mathbf{O}} := \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{\mathbf{O} \times \mathbf{P}} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta}^T}_{\mathbf{P} \times \mathbf{O}}$$

Definition

Neural network $f(\theta, x) \in \mathbb{R}^O$ Parameters $\theta \in \mathbb{R}^P$

Batch size N

The NTK is the outer product of Jacobians:

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{N^O \times N^O} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{N^O \times P} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta}^T}_{P \times N^O}$$

Definition

Neural network $f(\theta, x) \in \mathbb{R}^O$ Parameters $\theta \in \mathbb{R}^P$

Batch size N

The NTK is the outer product of Jacobians:

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{\mathbf{NO} \times \mathbf{NO}} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{\mathbf{NO} \times \mathbf{P}} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta}^T}_{\mathbf{P} \times \mathbf{NO}}$$

Applications in approximate inference, deep learning theory, NAS, MAML, etc.

Definition

Neural network $f(\theta, x) \in \mathbb{R}^O$ Parameters $\theta \in \mathbb{R}^P$ Batch size N

The NTK is the outer product of Jacobians:

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{N \times O \times N \times O} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{N \times O \times P} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta}^T}_{P \times N \times O}$$

Very costly contraction: $N^2 O^2 P$ time!

Definition

Neural network $f(\theta, x) \in \mathbb{R}^O$ Parameters $\theta \in \mathbb{R}^P$ Batch size N

The NTK is the outer product of Jacobians:

1000² for ImageNet

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{N \times O \times N \times O} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{N \times O \times P} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta}}_{P \times N \times O}^T$$

Tens of millions

Very costly contraction: $N^2 O^2 P$ time!

Definition

Neural network $f(\theta, x) \in \mathbb{R}^O$ Parameters $\theta \in \mathbb{R}^P$ Batch size N

The NTK is the outer product of Jacobians:

This work: reducing the contraction cost.

$$\Theta_{\theta}(x_1, x_2) = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{N \times O} \underbrace{\frac{\partial f(\theta, x_2)}{\partial \theta}^T}_{P \times N}$$

The diagram shows the contraction of two Jacobian matrices. The left Jacobian is of size $N \times O$ and the right Jacobian is of size $P \times N$. Their outer product results in a $N \times O \times P \times N$ contraction.

Very costly contraction: ~~$N^2 O^2 P$~~ time!

Plan

1. NTK definition and notation

2. **Recap of Automatic Differentiation (AD)**

3. Baseline NTK computation complexity

4. Two new algorithms for computing the NTK

5. Benchmarks

Recap of Automatic Differentiation (AD)

Two fundamental AD building blocks:

- **Jacobian-vector product** $JVP_{(f,\theta,x)} : \theta_t \in \mathbb{R}^{\mathbf{P}} \mapsto \frac{\partial f(\theta, x)}{\partial \theta} \theta_t \in \mathbb{R}^{\mathbf{O}}$

- **vector-Jacobian product** $VJP_{(f,\theta,x)} : f_c \in \mathbb{R}^{\mathbf{O}} \mapsto \frac{\partial f(\theta, x)}{\partial \theta}^T f_c \in \mathbb{R}^{\mathbf{P}}$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Recap of Automatic Differentiation (AD)

Every other AD operation is implemented with them:

- Gradient $\text{grad}(f)(\theta, x) = \text{VJP}_{f, \theta, x}(1)$
- Reverse-mode Jacobian $\text{jacrev}(f)(\theta, x) = \text{vmap}(\text{VJP}_{f, \theta, x})(I_{\mathbf{O}})$
- Forward-mode Jacobian $\text{jacfwd}(f)(\theta, x) = \text{vmap}(\text{JVP}_{f, \theta, x})(I_{\mathbf{P}})$
- Hessian $\text{hessian}(f)(\theta, x) = \text{jacfwd}(\text{jacrev}(f)(\theta, x))(\theta, x)$
- ...

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Recap of Automatic Differentiation (AD)

- Time complexities of the forward pass (FP), JVP and VJP:

Op	Forward Pass (FP)	Jacobian-vector Product (JVP)	Vector-Jacobian Product (VJP)
Time	N [FP]		

Neural Network

$$f(\theta, x) \in \mathbb{R}^O$$

Output size

O

Parameters

$$\theta \in \mathbb{R}^P$$

Batch size

N

Forward pass

[FP]

Recap of Automatic Differentiation (AD)

- Time complexities of the reverse-mode Jacobian:

$$\underbrace{\frac{\partial f(\theta, x)^T}{\partial \theta}}_{\mathbf{P} \times \mathbf{O}} := \frac{\partial f(\theta, x)^T}{\partial \theta} I_{\mathbf{O}} = \underbrace{[\text{VJP}_{(f, \theta, x)}(e_1), \dots, \text{VJP}_{(f, \theta, x)}(e_{\mathbf{O}})]}_{\mathbf{O} \text{ VJPs}}$$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Recap of Automatic Differentiation (AD)

Given Jacobian definitions, we can compute its cost via JVP/VJP costs.

Op	Forward Pass (FP)	Jacobian-vector Product (JVP)	Vector-Jacobian Product (VJP)	Jacobian (reverse-mode)
Time	N [FP]			NO [FP]

Neural Network

$$f(\theta, x) \in \mathbb{R}^O$$

Output size

O

Parameters

$$\theta \in \mathbb{R}^P$$

Batch size

N

Forward pass

[**FP**]

Plan

1. NTK definition and notation
2. Recap of Automatic Differentiation (AD)
3. **Baseline NTK computation complexity**
4. Two new algorithms for computing the NTK
5. Benchmarks

Baseline: Jacobian contraction

Recall that we want to compute

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{\mathbf{NO} \times \mathbf{NO}} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{\mathbf{NO} \times \mathbf{P}} \underbrace{\frac{\partial f(\theta, x_2)^T}{\partial \theta}}_{\mathbf{P} \times \mathbf{NO}}$$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

$[\mathbf{FP}]$

Baseline: Jacobian contraction

Baseline algorithm: compute two Jacobians and contract them.

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{\mathbf{NO} \times \mathbf{NO}} = \left[\underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{\mathbf{NO} \times \mathbf{P}} \quad \underbrace{\frac{\partial f(\theta, x_2)^T}{\partial \theta}}_{\mathbf{P} \times \mathbf{NO}} \right]$$

Reverse-mode Jacobian + contraction

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Baseline: Jacobian contraction

Baseline algorithm: compute two Jacobians and contract them.

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{\mathbf{NO} \times \mathbf{NO}} = \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{\mathbf{NO} \times \mathbf{P}} \underbrace{\frac{\partial f(\theta, x_2)^T}{\partial \theta}}_{\mathbf{P} \times \mathbf{NO}}$$

\uparrow Reverse-mode Jacobian + $\mathbf{N}^2 \mathbf{O}^2 \mathbf{P}$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

$[\mathbf{FP}]$

Baseline: Jacobian contraction

Baseline algorithm: compute two Jacobians and contract them.

$$\underbrace{\Theta_{\theta}(x_1, x_2)}_{\mathbf{NO} \times \mathbf{NO}} = \begin{matrix} \underbrace{\frac{\partial f(\theta, x_1)}{\partial \theta}}_{\mathbf{NO} \times \mathbf{P}} & \underbrace{\frac{\partial f(\theta, x_2)^T}{\partial \theta}}_{\mathbf{P} \times \mathbf{NO}} \end{matrix}$$

$$\mathbf{NO} [\mathbf{FP}] + \mathbf{NO} [\mathbf{FP}] + \mathbf{N}^2 \mathbf{O}^2 \mathbf{P}$$

Neural Network	Output size	Parameters	Batch size	Forward pass
$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$	\mathbf{O}	$\theta \in \mathbb{R}^{\mathbf{P}}$	\mathbf{N}	$[\mathbf{FP}]$

Baseline: Jacobian contraction

Method	Jacobian Contraction
Time	\mathbf{NO} [FP] + $\mathbf{N^2O^2P}$

Neural Network
 $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$

Output size
 \mathbf{O}

Parameters
 $\theta \in \mathbb{R}^{\mathbf{P}}$

Batch size
 \mathbf{N}

Forward pass
[FP]

Plan

1. NTK definition and notation
2. Recap of Automatic Differentiation (AD)
3. Baseline NTK computation complexity
4. **Two new algorithms for computing the NTK**
5. Benchmarks

Idea 1: NTK-vector products

Similarly to Jacobians, compute the NTK via NTK-vector products applied to O columns of the identity matrix I_O .

Neural Network

$$f(\theta, x) \in \mathbb{R}^O$$

Output size

O

Parameters

$$\theta \in \mathbb{R}^P$$

Batch size

N

Forward pass

[FP]

Idea 1: NTK-vector products

Similarly to Jacobians, compute the NTK via NTK-vector products applied to \mathbf{O} columns of the identity matrix $I_{\mathbf{O}}$. For $\mathbf{N} = 1$:

Evaluate $\Theta_{\theta}(x_1, x_2) = \Theta_{\theta}(x_1, x_2)I_{\mathbf{O}}$ row-by-row.

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Idea 1: NTK-vector products

Similarly to Jacobians, compute the NTK via NTK-vector products applied to \mathbf{O} columns of the identity matrix $I_{\mathbf{O}}$. For $\mathbf{N} = 1$:

Evaluate $\Theta_{\theta}(x_1, x_2) = \Theta_{\theta}(x_1, x_2)I_{\mathbf{O}}$ row-by-row.

For a single input column v :

$$\Theta_{\theta}(x_1, x_2)v$$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Idea 1: NTK-vector products

Similarly to Jacobians, compute the NTK via NTK-vector products applied to \mathbf{O} columns of the identity matrix $I_{\mathbf{O}}$. For $\mathbf{N} = 1$:

Evaluate $\Theta_{\theta}(x_1, x_2) = \Theta_{\theta}(x_1, x_2)I_{\mathbf{O}}$ row-by-row.

For a single input column v :

$$\Theta_{\theta}(x_1, x_2)v = \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v$$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Idea 1: NTK-vector products

Similarly to Jacobians, compute the NTK via NTK-vector products applied to \mathbf{O} columns of the identity matrix $I_{\mathbf{O}}$. For $\mathbf{N} = 1$:

Evaluate $\Theta_{\theta}(x_1, x_2) = \Theta_{\theta}(x_1, x_2)I_{\mathbf{O}}$ row-by-row.

For a single input column v :

$$\Theta_{\theta}(x_1, x_2)v = \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v = \frac{\partial f(\theta, x_1)}{\partial \theta} \text{VJP}_{(f, \theta, x_2)}(v)$$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Idea 1: NTK-vector products

Similarly to Jacobians, compute the NTK via NTK-vector products applied to \mathbf{O} columns of the identity matrix $I_{\mathbf{O}}$. For $\mathbf{N} = 1$:

Evaluate $\Theta_{\theta}(x_1, x_2) = \Theta_{\theta}(x_1, x_2)I_{\mathbf{O}}$ row-by-row.

For a single input column v :

$$\Theta_{\theta}(x_1, x_2)v = \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v = \frac{\partial f(\theta, x_1)}{\partial \theta} \text{VJP}_{(f, \theta, x_2)}(v) = \text{JVP}_{(f, \theta, x_1)}[\text{VJP}_{(f, \theta, x_2)}(v)]$$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Idea 1: NTK-vector products

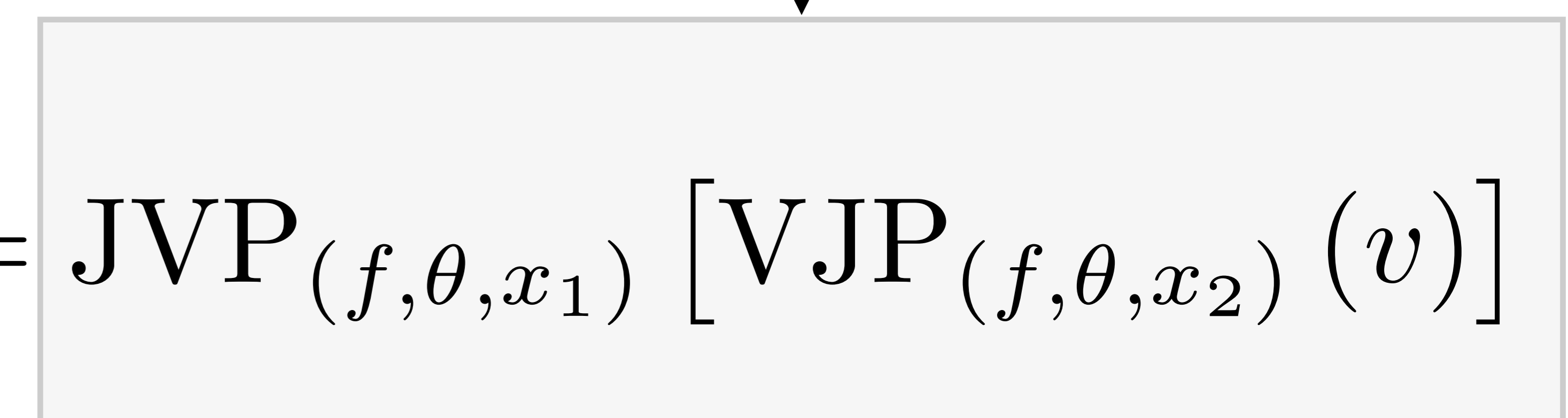
Similarly to Jacobians, compute the NTK via NTK-vector products applied to \mathbf{O} columns of the identity matrix $I_{\mathbf{O}}$. For $\mathbf{N} = 1$:

Evaluate $\Theta_{\theta}(x_1, x_2) = \Theta_{\theta}(x_1, x_2)I_{\mathbf{O}}$ row-by-row.

For a single input column v :

$$\Theta_{\theta}(x_1, x_2)v = \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v = \frac{\partial f(\theta, x_1)}{\partial \theta} \text{VJP}_{(f, \theta, x_2)}(v) = \text{JVP}_{(f, \theta, x_1)}[\text{VJP}_{(f, \theta, x_2)}(v)]$$

[FP]



Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Idea 1: NTK-vector products

Similarly to Jacobians, compute the NTK via NTK-vector products applied to \mathbf{O} columns of the identity matrix $I_{\mathbf{O}}$. For $\mathbf{N} = 1$:

Evaluate $\Theta_{\theta}(x_1, x_2) = \Theta_{\theta}(x_1, x_2)I_{\mathbf{O}}$ row-by-row. (repeat \mathbf{O} times)

For a single input column v : \mathbf{O} [FP]

$$\Theta_{\theta}(x_1, x_2)v = \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v = \frac{\partial f(\theta, x_1)}{\partial \theta} \text{VJP}_{(f, \theta, x_2)}(v) = \text{JVP}_{(f, \theta, x_1)}[\text{VJP}_{(f, \theta, x_2)}(v)]$$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Idea 1: NTK-vector products

Batched setting: needs to be done for every pair of x_1 and x_2 !

N²**O** [**FP**]

$$\Theta_{\theta}(x_1, x_2)v = \frac{\partial f(\theta, x_1)}{\partial \theta} \frac{\partial f(\theta, x_2)}{\partial \theta}^T v = \frac{\partial f(\theta, x_1)}{\partial \theta} \text{VJP}_{(f, \theta, x_2)}(v) = \text{JVP}_{(f, \theta, x_1)}[\text{VJP}_{(f, \theta, x_2)}(v)]$$

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

O

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

N

Forward pass

[**FP**]

Idea 1: NTK-vector products

Method	Jacobian Contraction	NTK-vector products
Time	\mathbf{NO} [FP] + $\mathbf{N^2O^2P}$	$\mathbf{N^2O}$ [FP]

Neural Network

$$f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$$

Output size

\mathbf{O}

Parameters

$$\theta \in \mathbb{R}^{\mathbf{P}}$$

Batch size

\mathbf{N}

Forward pass

[FP]

Idea 2: Structured derivatives

Exploit the structure in primitive derivatives for computing the NTK.

Neural Network

$$f(\theta, x) \in \mathbb{R}^O$$

Output size

O

Parameters

$$\theta \in \mathbb{R}^P$$

Batch size

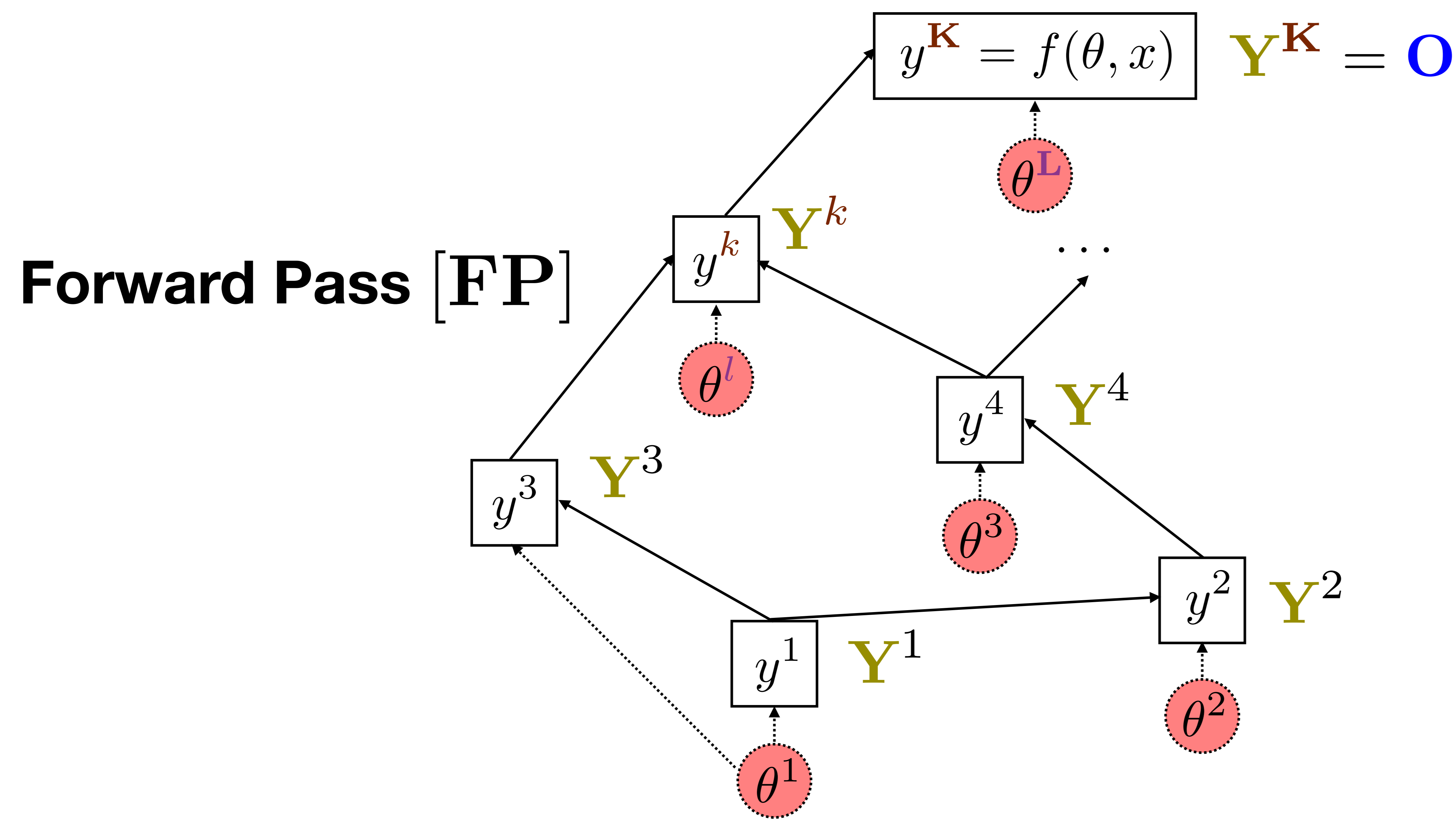
N

Forward pass

[FP]

Idea 2: Structured derivatives

Consider the computation graph of the function f :

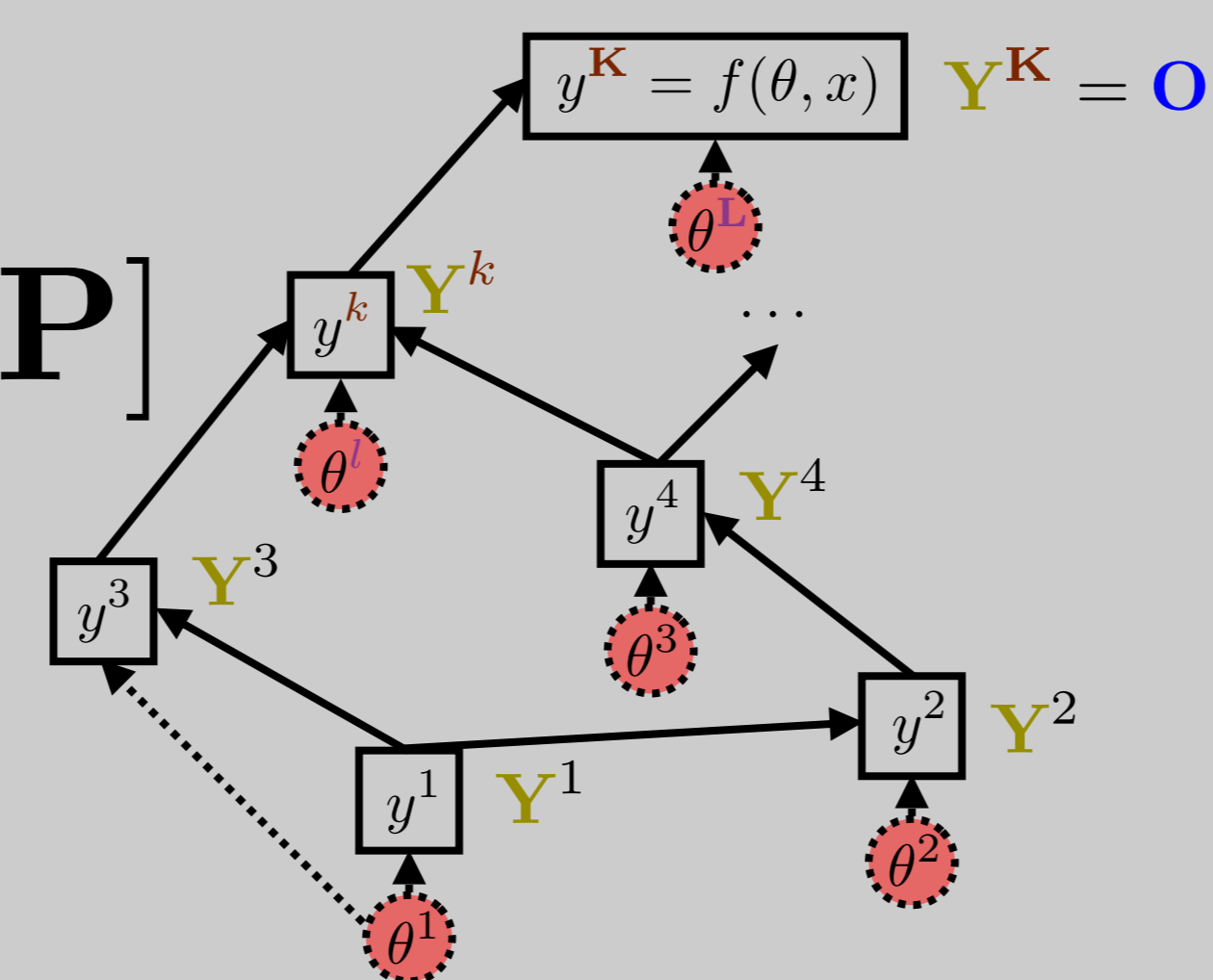


$$\Theta_{\theta}(x_1, x_2) = \frac{\partial f_1}{\partial \theta} \frac{\partial f_2}{\partial \theta}^T$$

Neural Network $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$

Batch size \mathbf{N} Output size \mathbf{O}

Forward Pass [FP]



Parameters $\theta = \text{vec}[\theta^0, \dots, \theta^L]$,

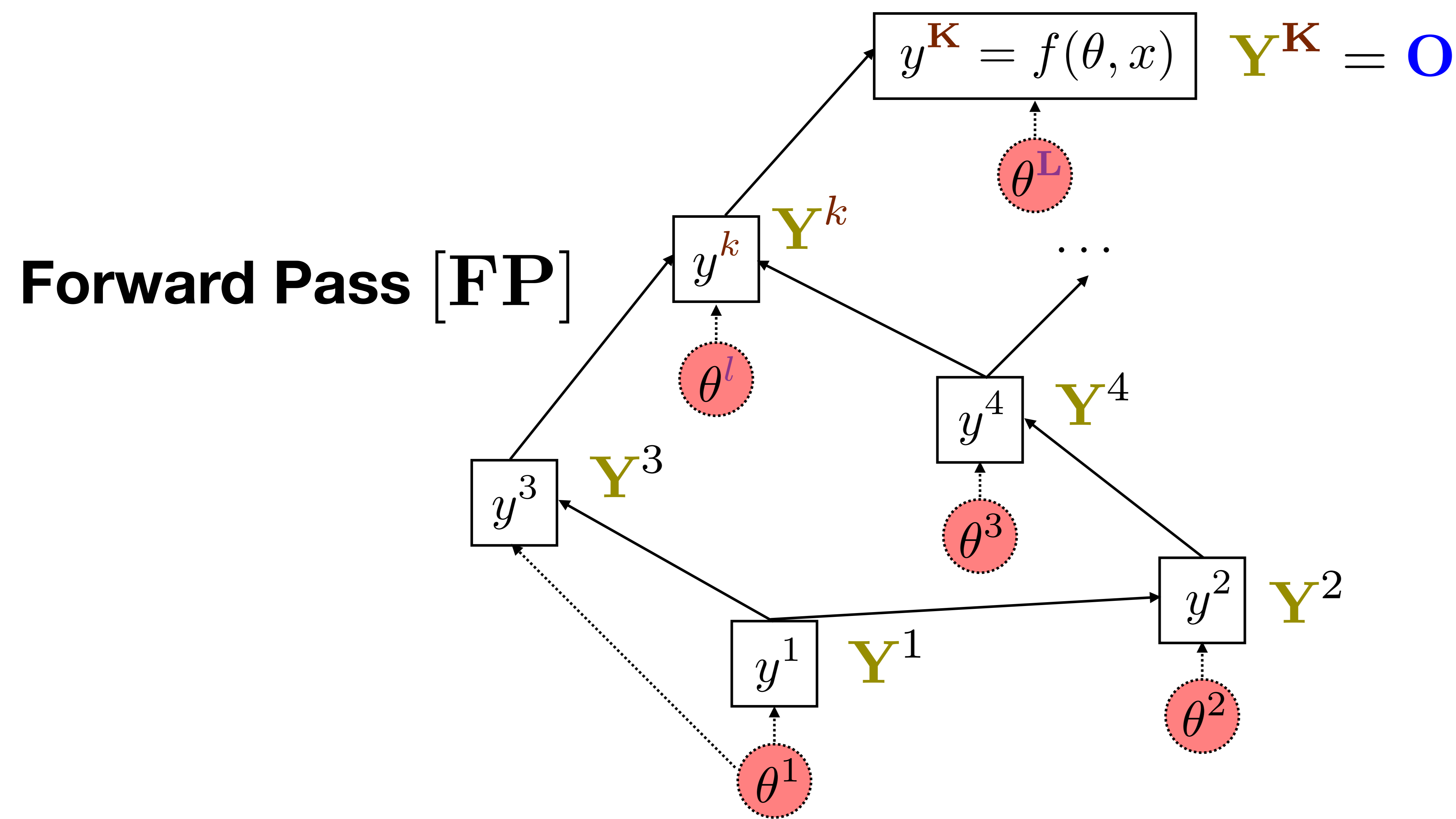
$$\mathbf{P} = \sum_{l=0}^L \mathbf{P}^l$$

Intermediate primitive outputs y^1, \dots, y^K ,

$$\mathbf{Y} = \sum_{k=1}^K \mathbf{Y}^k$$

Idea 2: Structured derivatives

Apply the chain rule to the NTK expression:

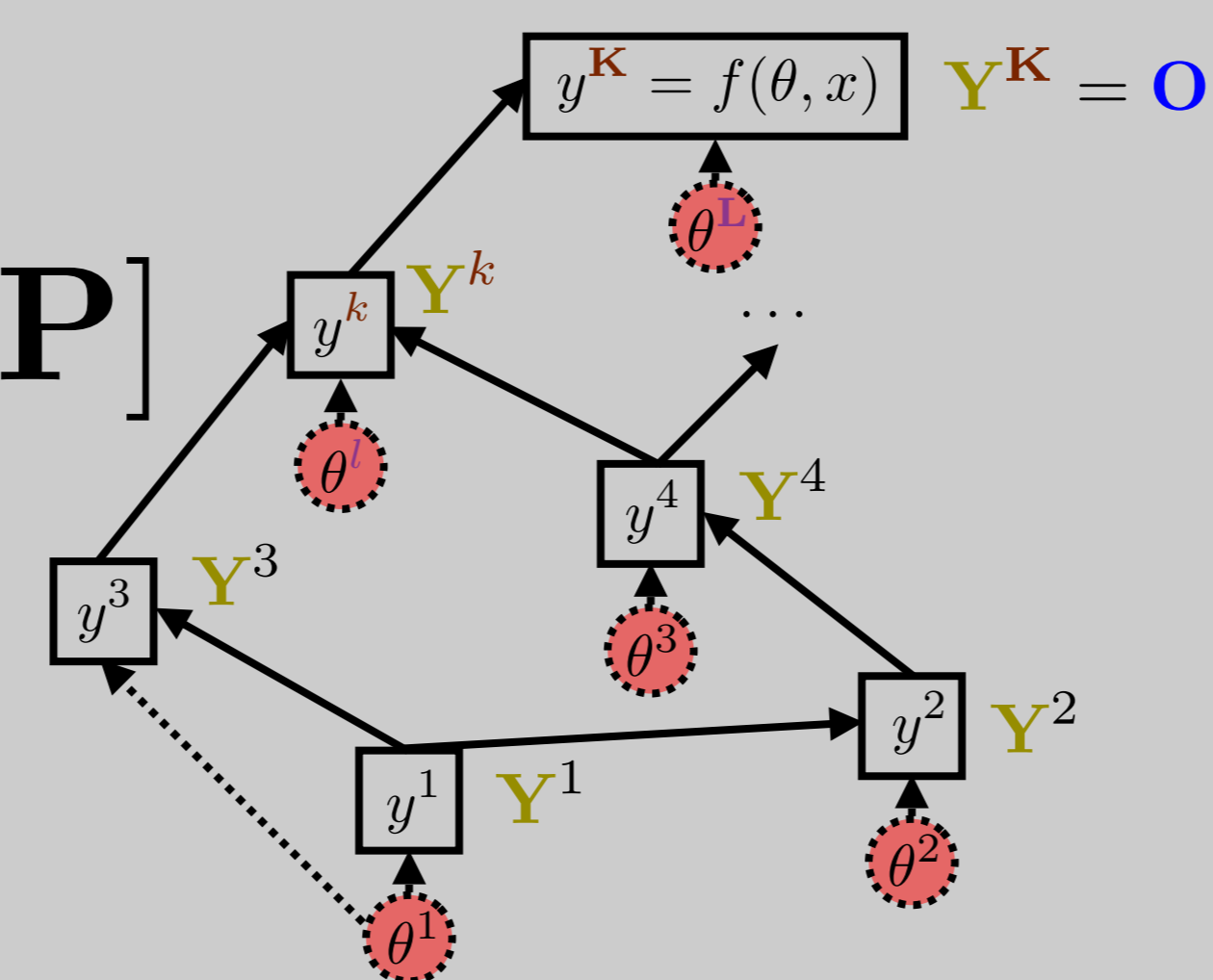


$$\Theta_{\theta}(x_1, x_2) = \frac{\partial f_1}{\partial \theta} \frac{\partial f_2}{\partial \theta}^T = \sum_{l=0}^L \frac{\partial f_1}{\partial \theta^l} \frac{\partial f_2}{\partial \theta^l}^T =$$

Neural Network $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$

Batch size \mathbf{N} Output size \mathbf{O}

Forward Pass [FP]



Parameters

$$\theta = \text{vec} [\theta^0, \dots, \theta^L]$$

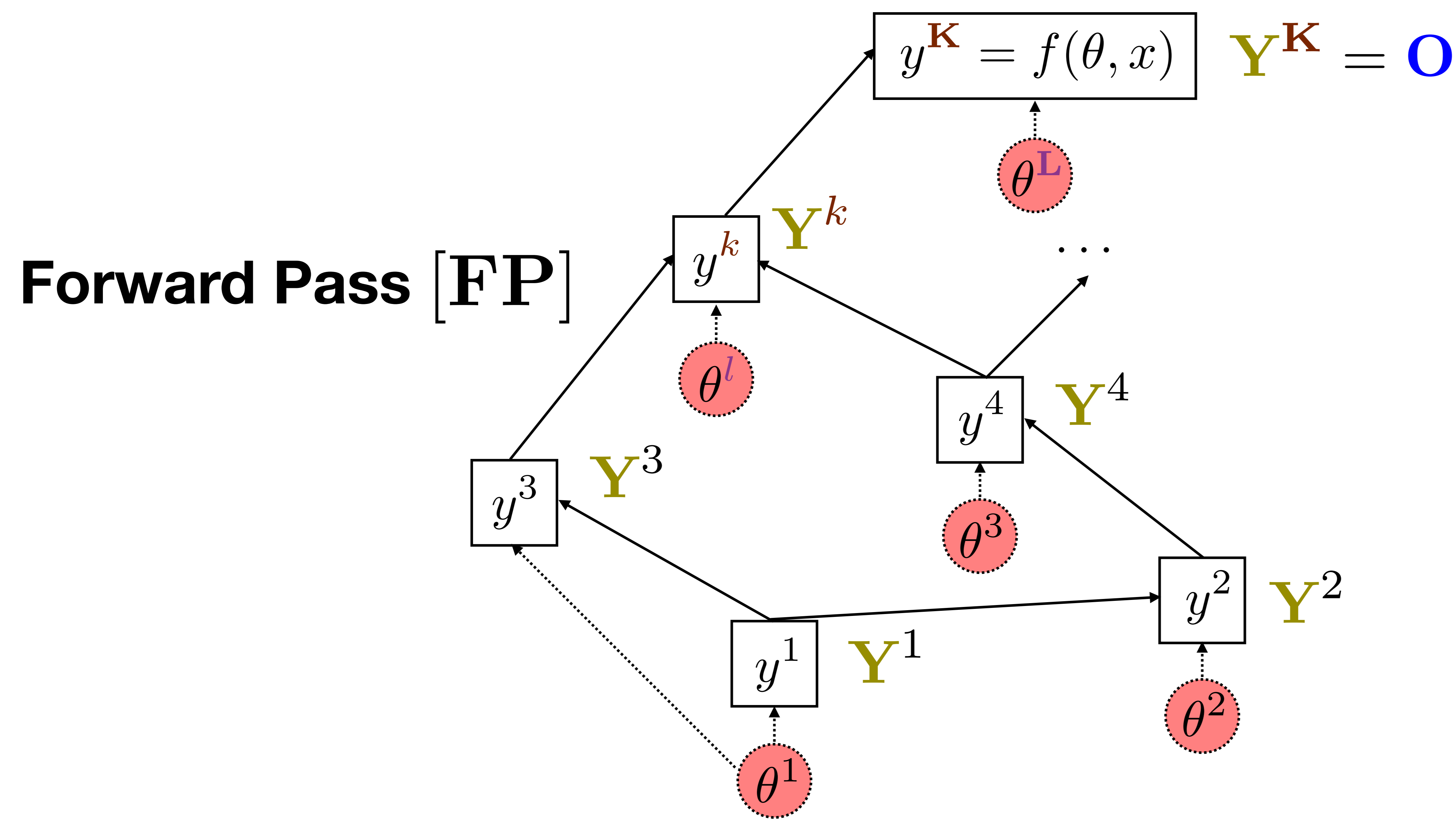
$$\mathbf{P} = \sum_{l=0}^L \mathbf{P}^l$$

Intermediate primitive outputs y^1, \dots, y^K ,

$$\mathbf{Y} = \sum_{k=1}^K \mathbf{Y}^k$$

Idea 2: Structured derivatives

Apply the chain rule to the NTK expression:



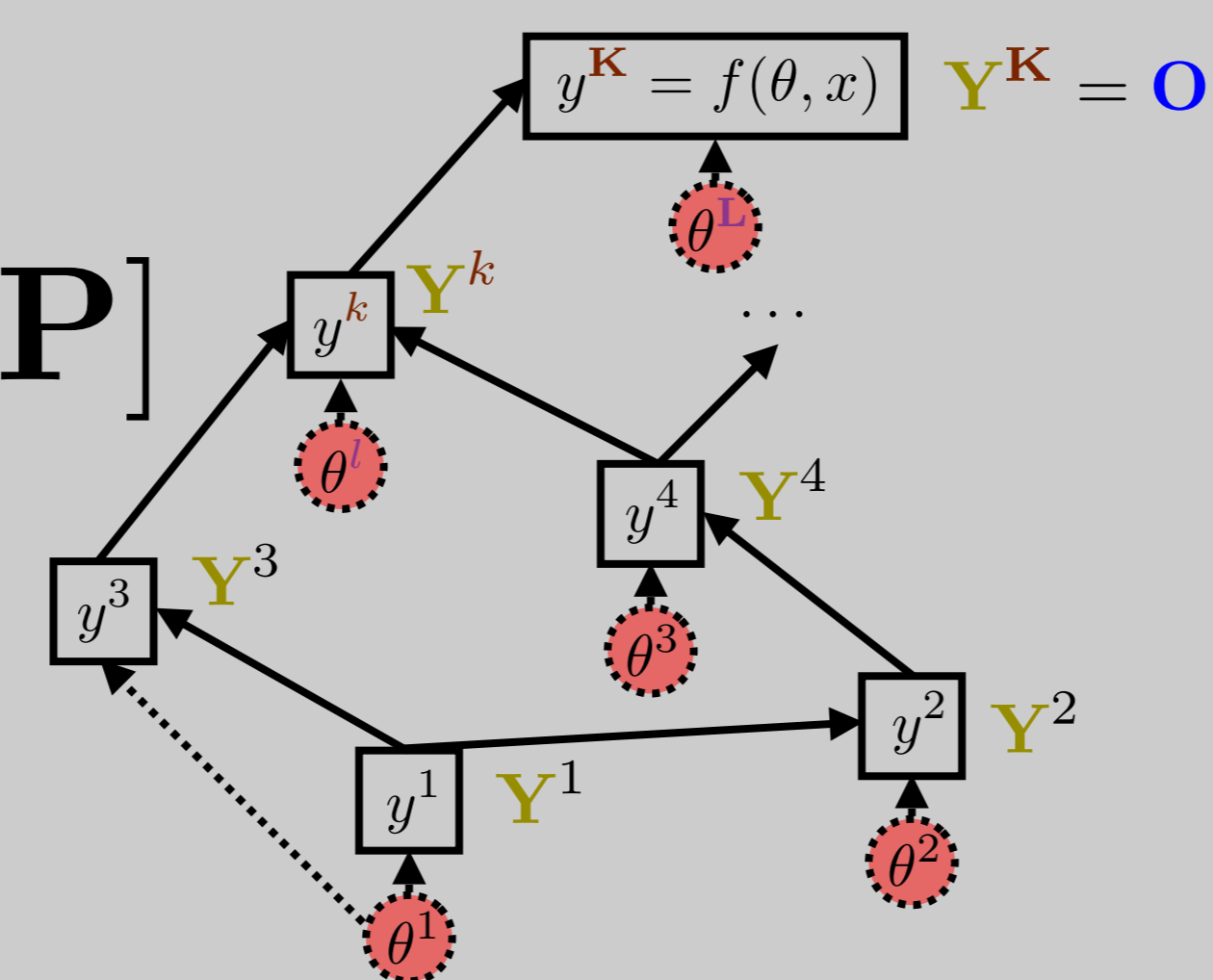
$$\Theta_{\theta}(x_1, x_2) = \frac{\partial f_1}{\partial \theta} \frac{\partial f_2}{\partial \theta}^T = \sum_{l=0}^{\mathbf{L}} \frac{\partial f_1}{\partial \theta^l} \frac{\partial f_2}{\partial \theta^l}^T =$$

$$= \sum_{l, k_1, k_2} \left(\frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \right) \left(\frac{\partial f_2}{\partial y_2^{k_2}} \frac{\partial y_2^{k_2}}{\partial \theta^l} \right)^T =$$

Neural Network $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$

Batch size \mathbf{N} Output size \mathbf{O}

Forward Pass [FP]



Parameters

$$\theta = \text{vec} [\theta^0, \dots, \theta^{\mathbf{L}}],$$

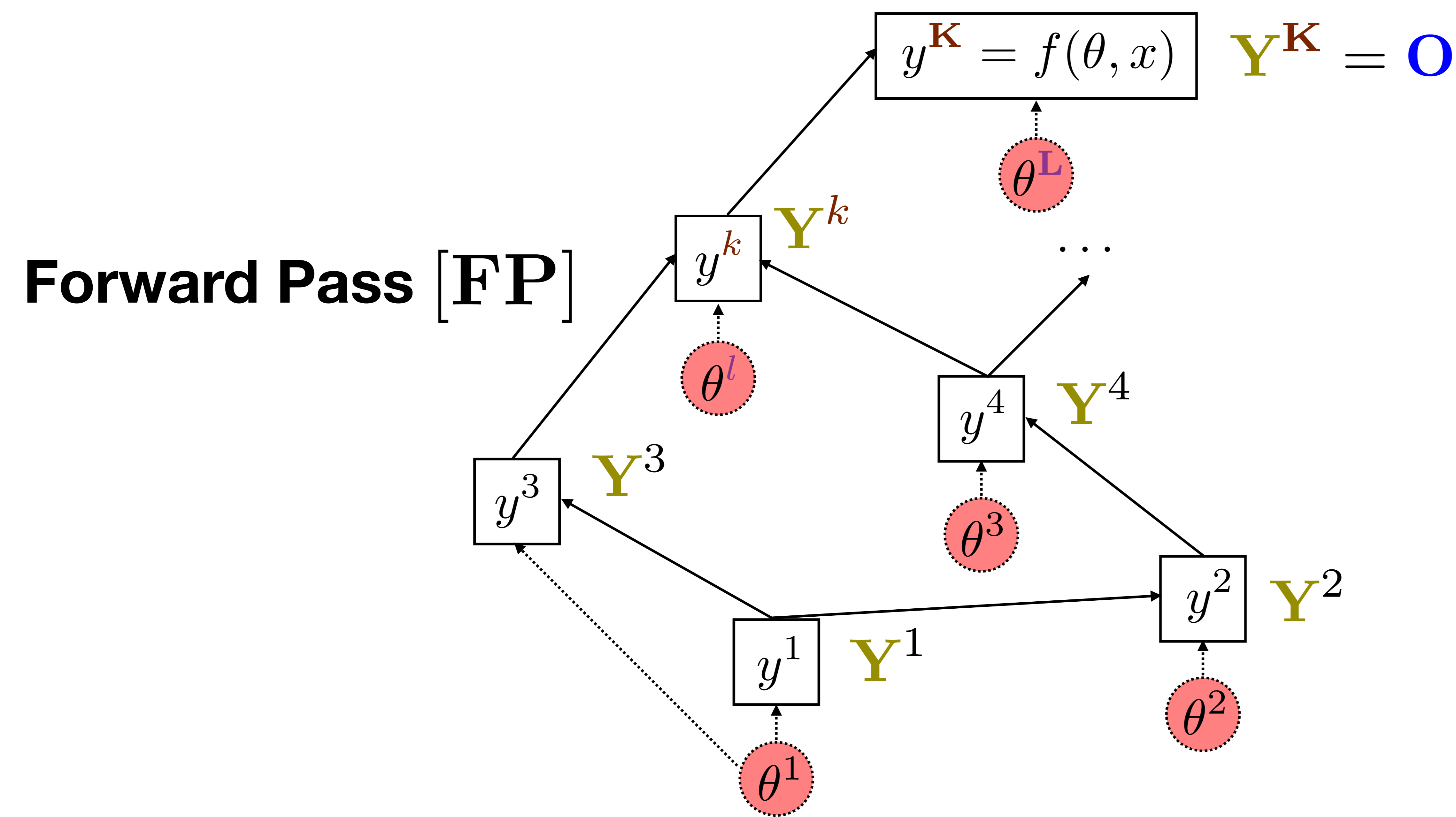
$$\mathbf{P} = \sum_{l=0}^{\mathbf{L}} \mathbf{P}^l$$

Intermediate primitive outputs $y^1, \dots, y^K,$

$$\mathbf{Y} = \sum_{k=1}^{\mathbf{K}} \mathbf{Y}^k$$

Idea 2: Structured derivatives

Apply the chain rule to the NTK expression:

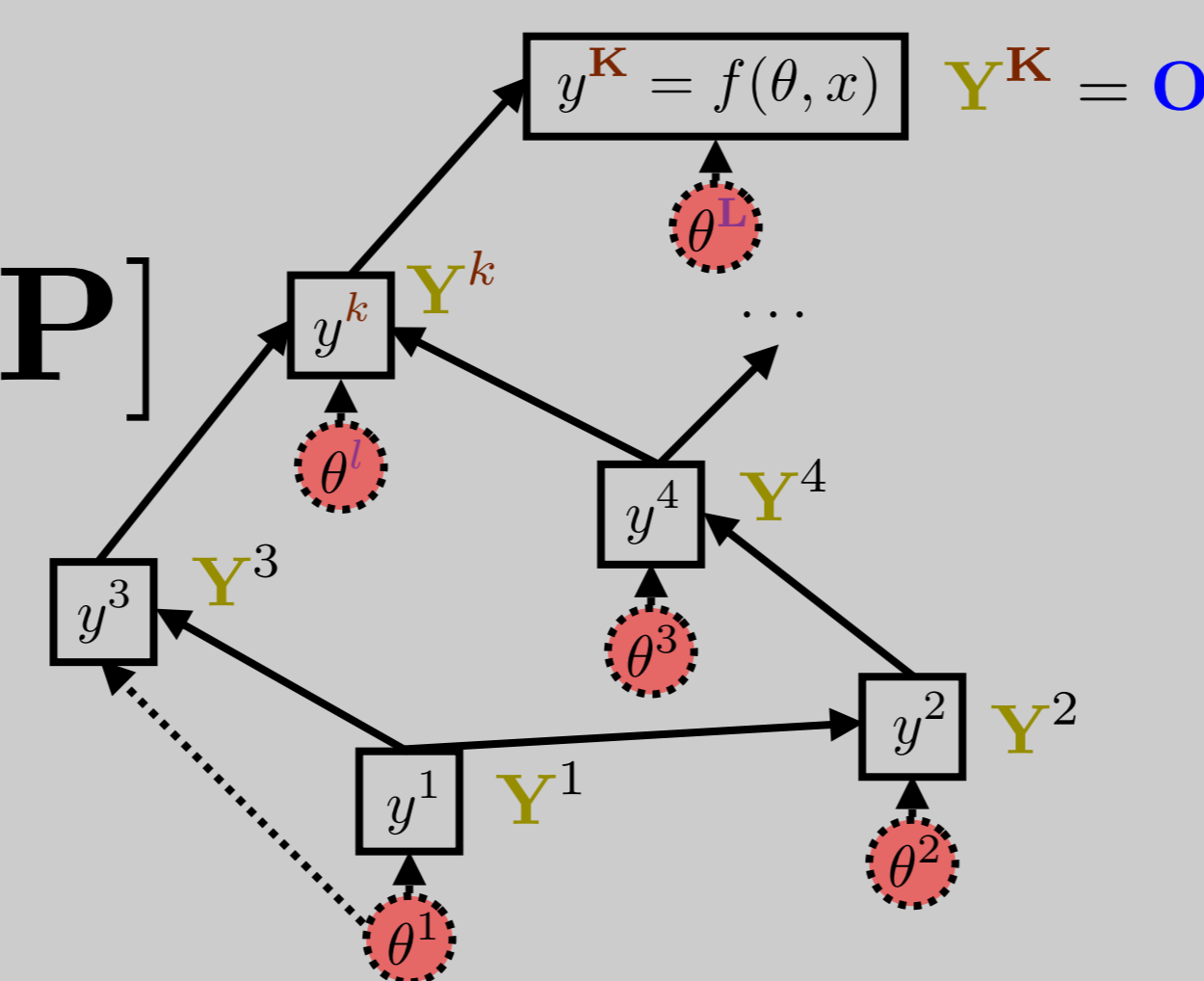


$$\begin{aligned} \Theta_{\theta}(x_1, x_2) &= \frac{\partial f_1}{\partial \theta} \frac{\partial f_2}{\partial \theta}^T = \sum_{l=0}^{\mathbf{L}} \frac{\partial f_1}{\partial \theta^l} \frac{\partial f_2}{\partial \theta^l}^T = \\ &= \sum_{l, k_1, k_2} \left(\frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \right) \left(\frac{\partial f_2}{\partial y_2^{k_2}} \frac{\partial y_2^{k_2}}{\partial \theta^l} \right)^T = \\ &= \sum_{l, k_1, k_2} \frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \frac{\partial y_2^{k_2}}{\partial \theta^l}^T \frac{\partial f_2}{\partial y_2^{k_2}}^T \end{aligned}$$

Neural Network $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$

Batch size \mathbf{N} Output size \mathbf{O}

Forward Pass [FP]



Parameters

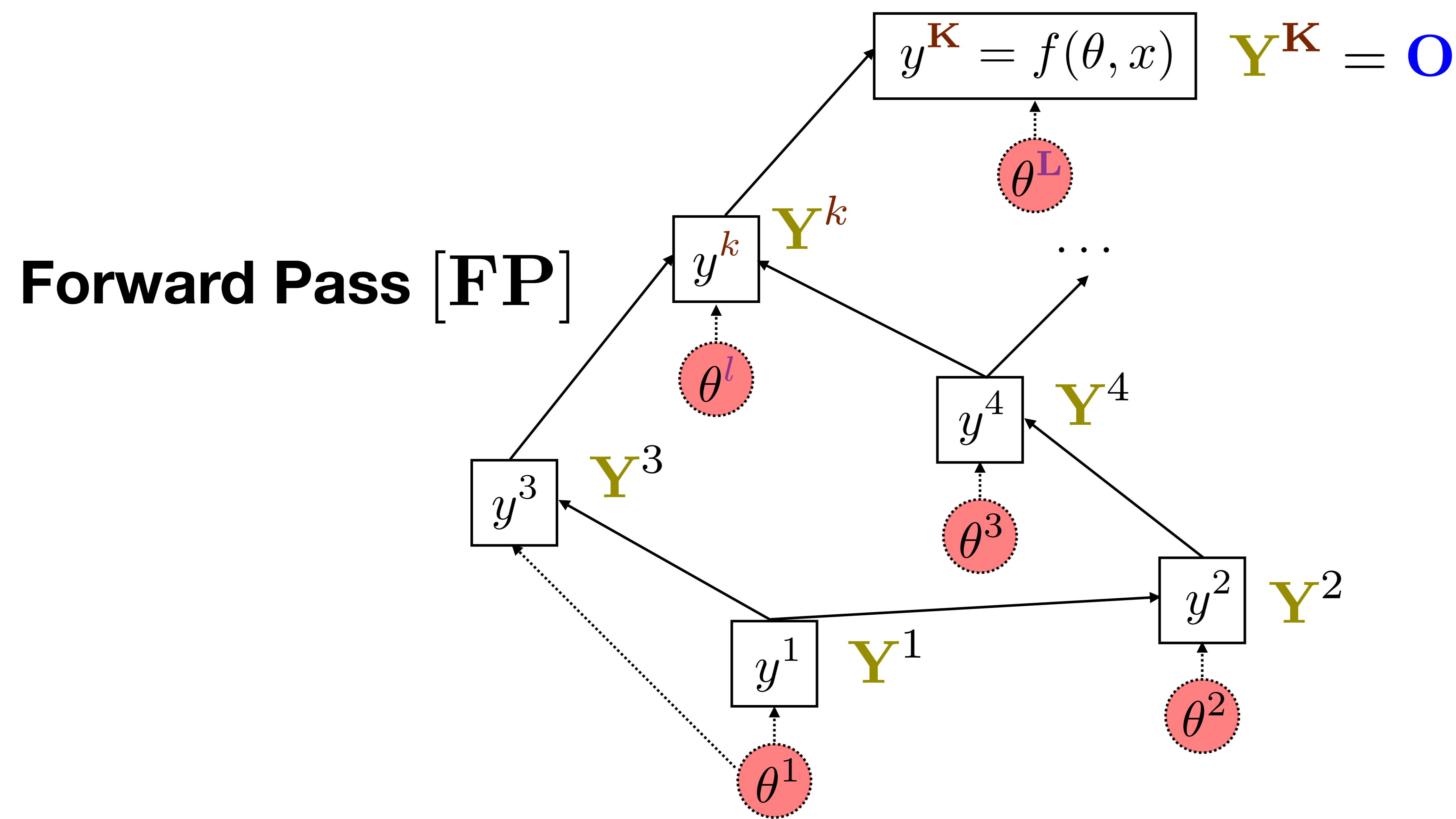
$$\theta = \text{vec} [\theta^0, \dots, \theta^{\mathbf{L}}], \quad \mathbf{P} = \sum_{l=0}^{\mathbf{L}} \mathbf{P}^l$$

Intermediate primitive outputs $y^1, \dots, y^{\mathbf{K}}$,

$$\mathbf{Y} = \sum_{k=1}^{\mathbf{K}} \mathbf{Y}^k$$

Idea 2: Structured derivatives

Apply the chain rule to the NTK expression:



$$\Theta_{\theta}(x_1, x_2) = \frac{\partial f_1}{\partial \theta} \frac{\partial f_2}{\partial \theta}^T = \sum_{l=0}^L \frac{\partial f_1}{\partial \theta^l} \frac{\partial f_2}{\partial \theta^l}^T =$$

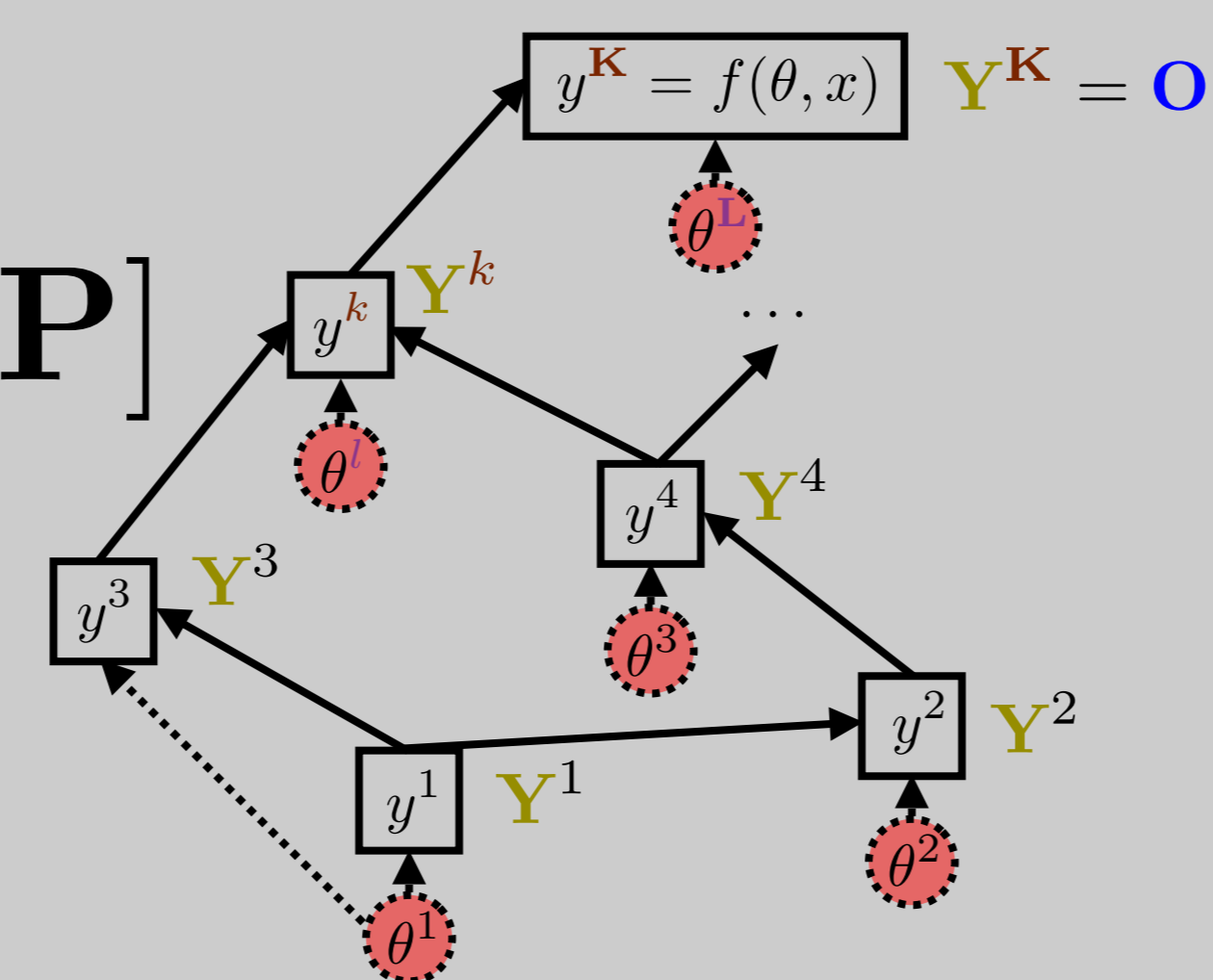
$$= \sum_{l, k_1, k_2} \left(\frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \right) \left(\frac{\partial f_2}{\partial y_2^{k_2}} \frac{\partial y_2^{k_2}}{\partial \theta^l} \right)^T =$$

$$= \sum_{l, k_1, k_2} \boxed{\frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \frac{\partial y_2^{k_2}}{\partial \theta^l}^T \frac{\partial f_2}{\partial y_2^{k_2}}^T}$$

Neural Network $f(\theta, x) \in \mathbb{R}^O$

Batch size N Output size O

Forward Pass [FP]

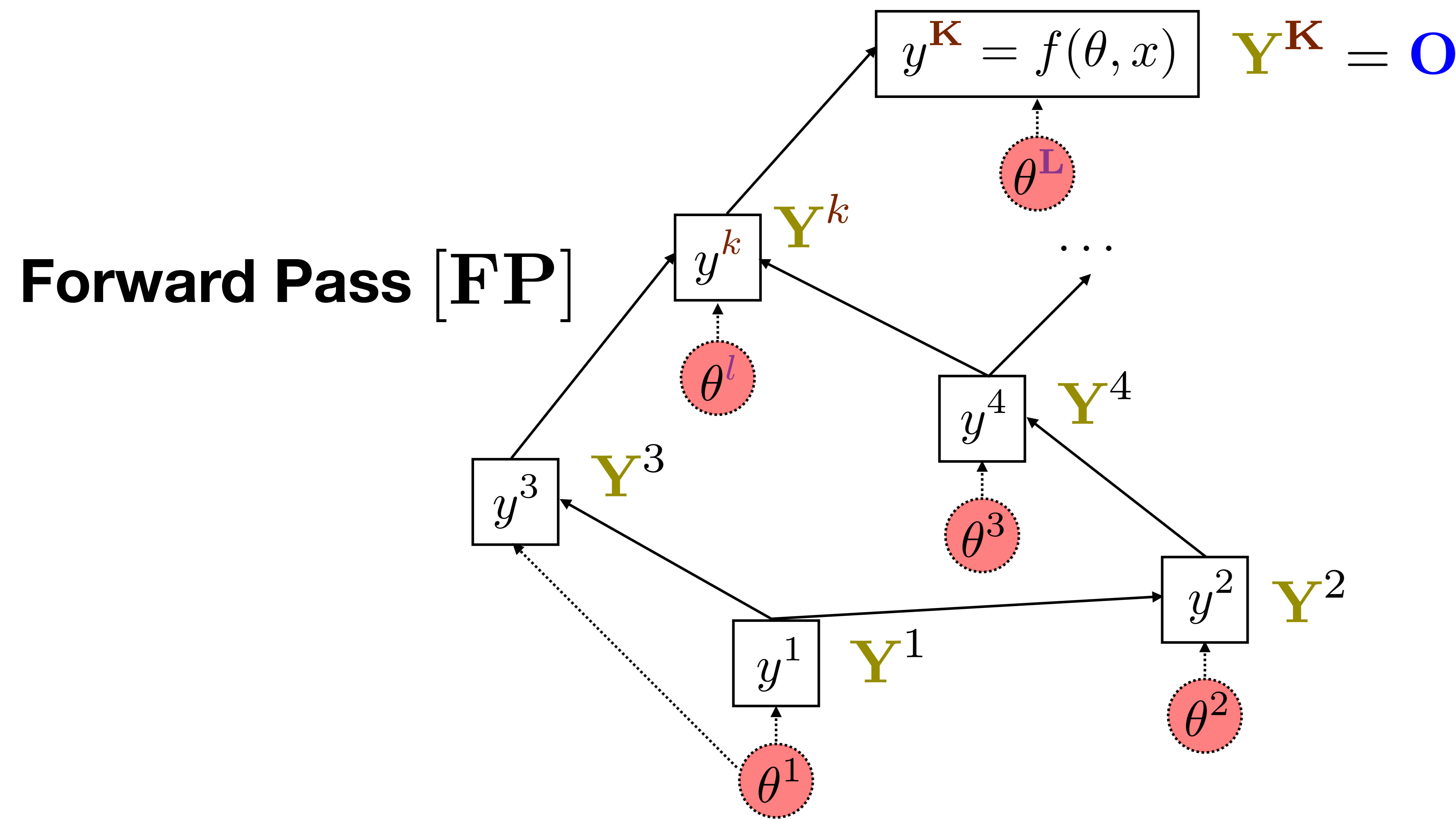


Parameters $\theta = \text{vec} [\theta^0, \dots, \theta^L]$, $\mathbf{P} = \sum_{l=0}^L \mathbf{P}^l$

Intermediate primitive outputs y^1, \dots, y^K , $\mathbf{Y} = \sum_{k=1}^K \mathbf{Y}^k$

Idea 2: Structured derivatives

NTK is a sum of matrix-Jacobian-Jacobian-matrix products.

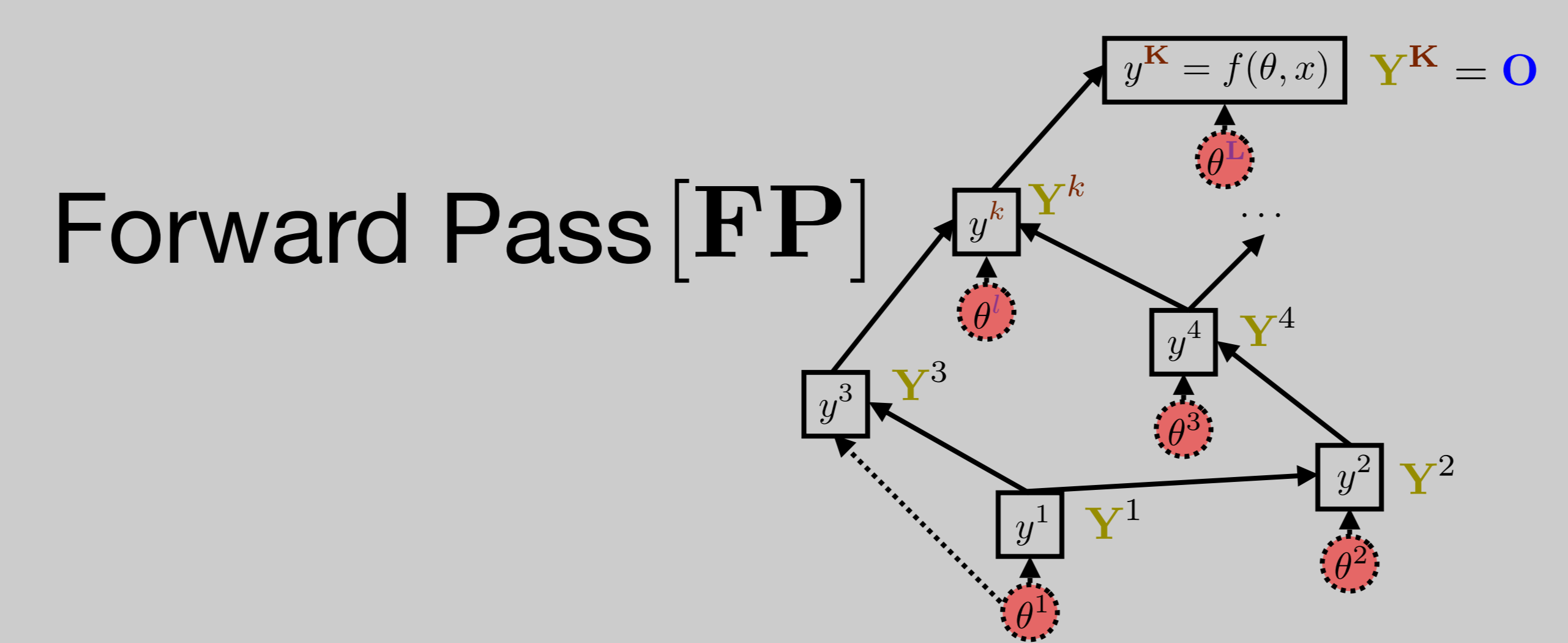


$$\Theta_{\theta}(x_1, x_2) = \frac{\partial f_1}{\partial \theta} \frac{\partial f_2}{\partial \theta}^T = \sum_{l=0}^L \frac{\partial f_1}{\partial \theta^l} \frac{\partial f_2}{\partial \theta^l}^T =$$

"matrix-Jacobian-Jacobian-matrix product" (MJJMP)

$$= \sum_{l, k_1, k_2} \left(\frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \right) \left(\frac{\partial f_2}{\partial y_2^{k_2}} \frac{\partial y_2^{k_2}}{\partial \theta^l} \right)^T$$

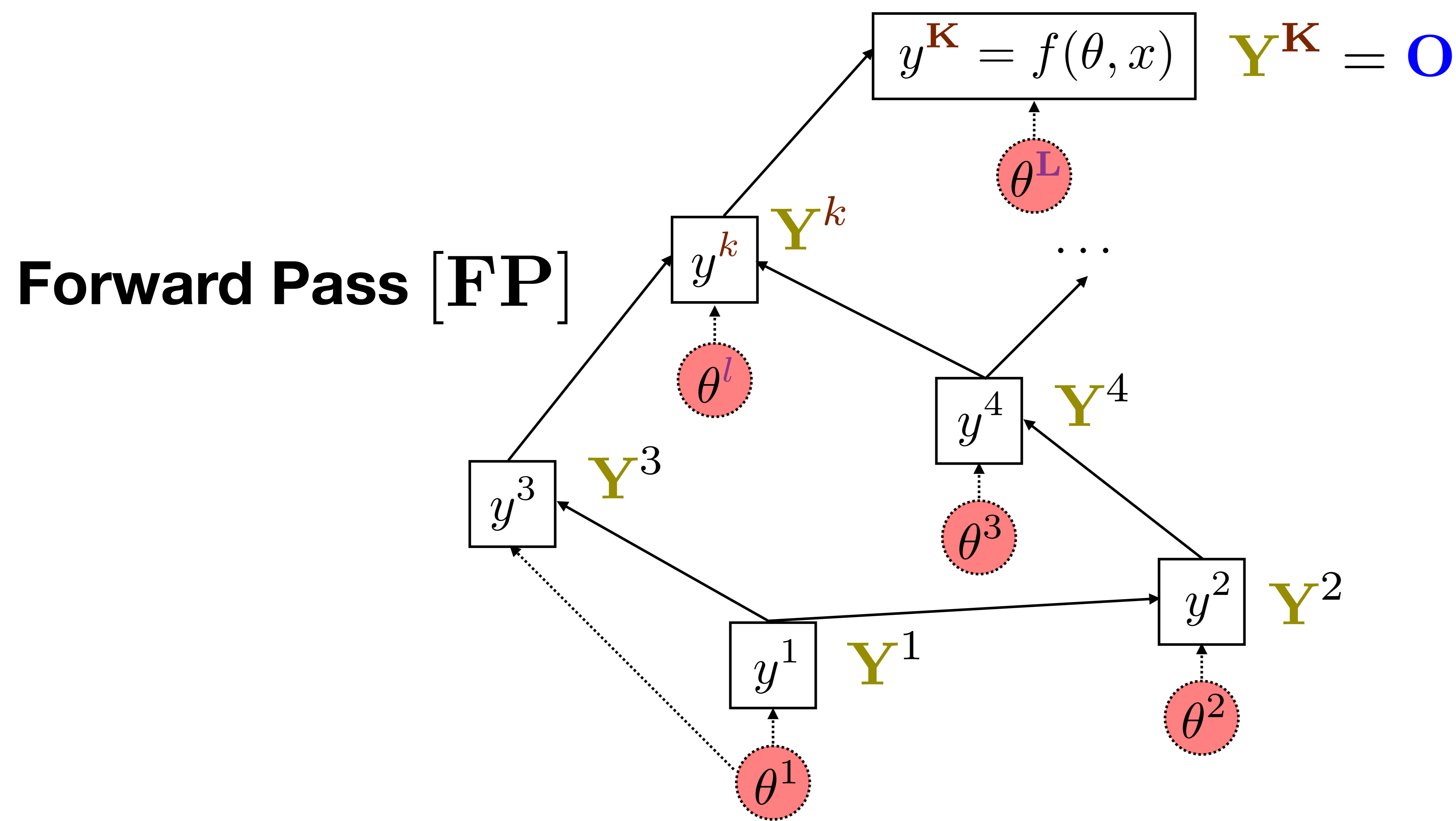
Neural Network $f(\theta, x) \in \mathbb{R}^O$
 Batch size N Output size O



Parameters $\theta = \text{vec} [\theta^0, \dots, \theta^L]$, $\mathbf{P} = \sum_{l=0}^L \mathbf{P}^l$
 Intermediate primitive outputs y^1, \dots, y^K , $\mathbf{Y} = \sum_{k=1}^K \mathbf{Y}^k$

Idea 2: Structured derivatives

Idea: for every pair of primitives y_1, y_2 , implement efficient MJJMPs:



$$\left(\begin{array}{c} c_1, c_2 \\ \underbrace{\mathbf{0} \times \mathbf{Y}_1 \quad \mathbf{0} \times \mathbf{Y}_2}_{\mathbf{0} \times \mathbf{O}} \end{array} \right) \mapsto \underbrace{c_1 \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T}_{\mathbf{0} \times \mathbf{O}} c_2^T$$

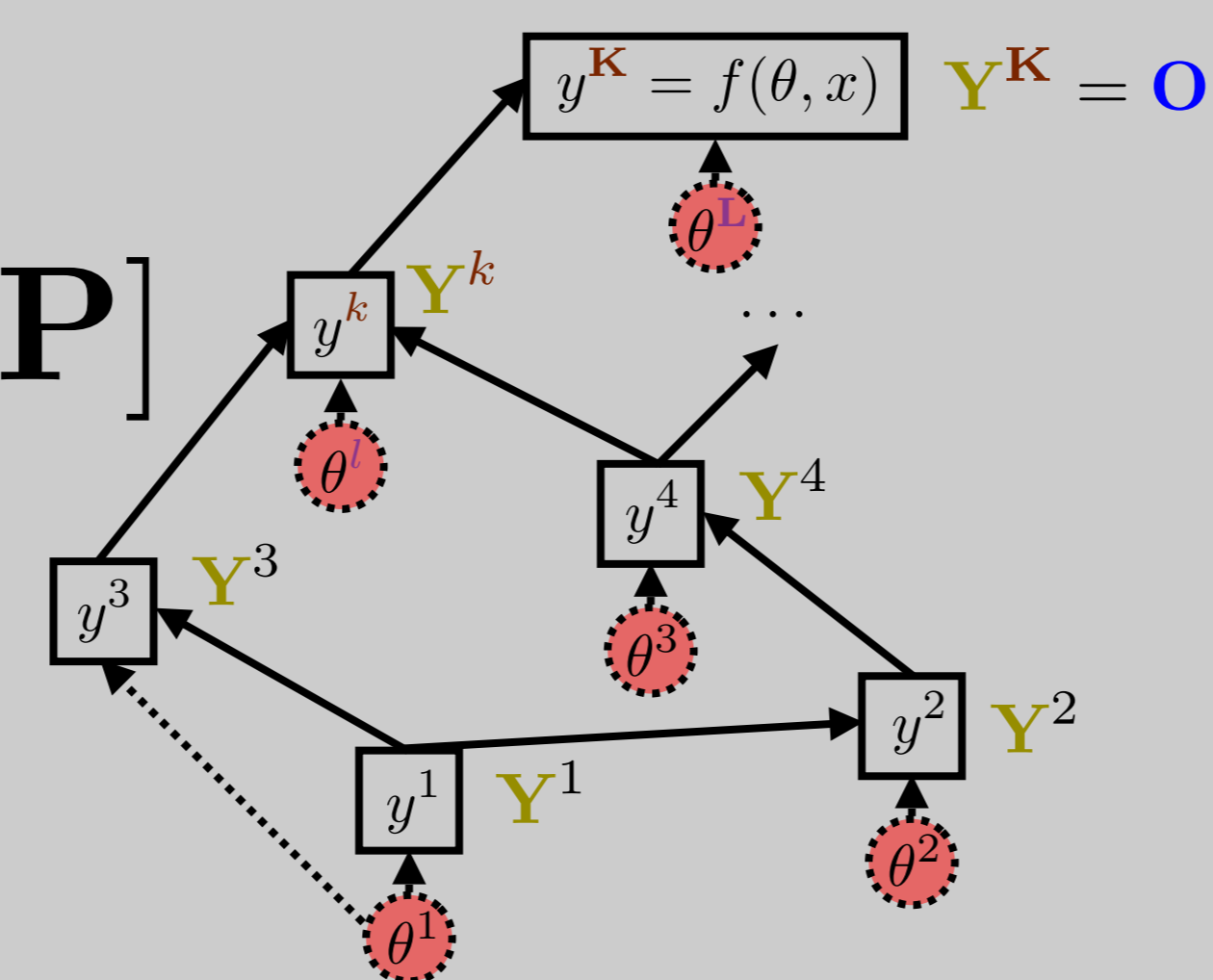
"matrix-Jacobian-Jacobian-matrix product" (MJJMP)

$$= \sum_{l, k_1, k_2} \frac{\partial f_1}{\partial y_1^{k_1}} \frac{\partial y_1^{k_1}}{\partial \theta^l} \frac{\partial y_2^{k_2}}{\partial \theta^l} \frac{\partial f_2}{\partial y_2^{k_2}}^T$$

Neural Network $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$

Batch size \mathbf{N} Output size \mathbf{O}

Forward Pass [FP]



Parameters

$$\theta = \text{vec} [\theta^0, \dots, \theta^{\mathbf{L}}], \quad \mathbf{P} = \sum_{l=0}^{\mathbf{L}} \mathbf{P}^l$$

Intermediate primitive outputs

$$y^1, \dots, y^{\mathbf{K}}, \quad \mathbf{Y} = \sum_{k=1}^{\mathbf{K}} \mathbf{Y}^k$$

Idea 2: Structured derivatives

Idea: for every pair of primitives y_1, y_2 , implement efficient **MJJMPs**:

$$\left(\underbrace{c_1}_{\mathbf{O} \times \mathbf{Y}_1}, \underbrace{c_2}_{\mathbf{O} \times \mathbf{Y}_2} \right) \mapsto \underbrace{c_1 \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^T c_2^T}_{\mathbf{O} \times \mathbf{O}}$$

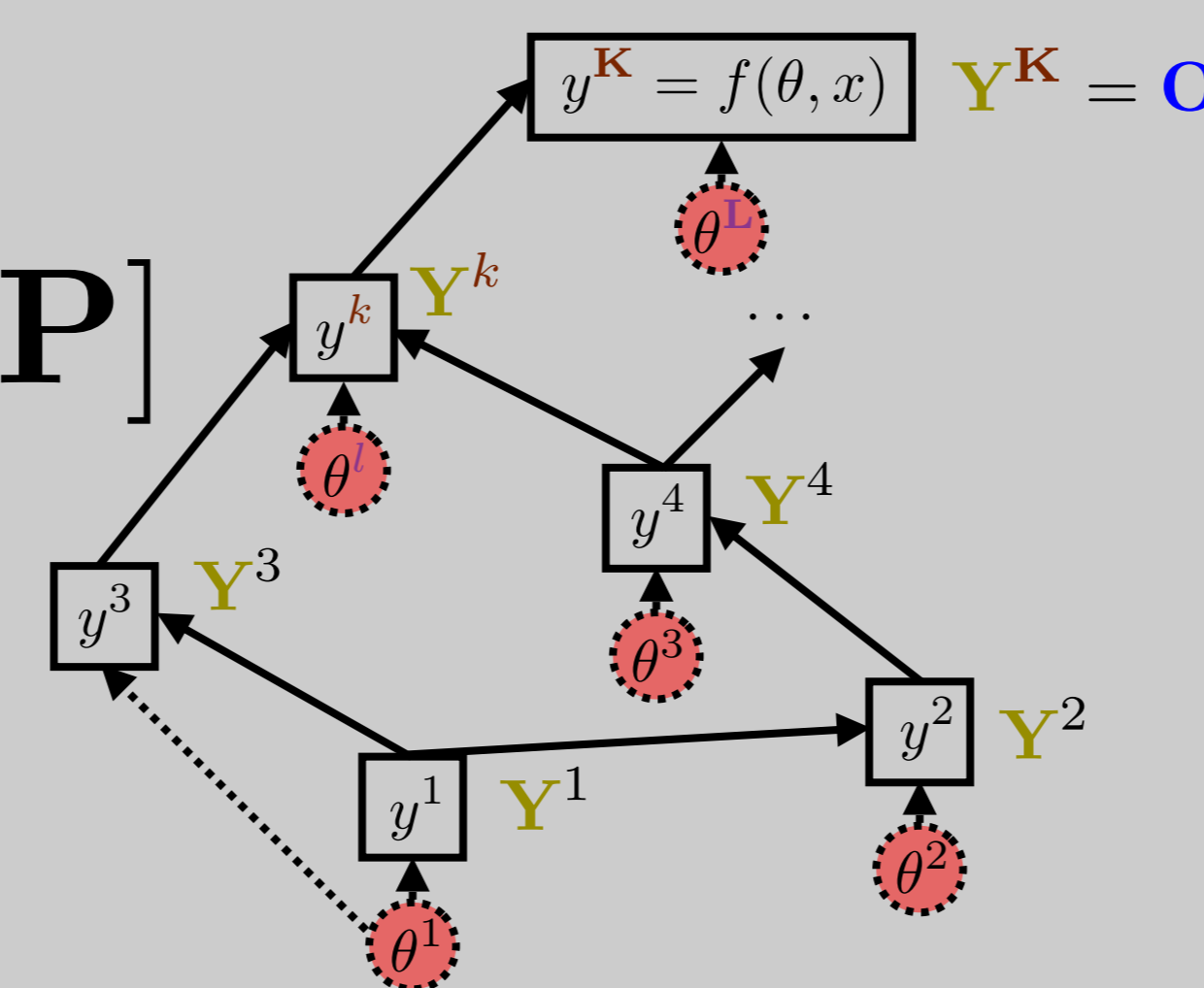
Similar to existing AD tools: for every primitive y , implement efficient

$$\text{JVP}_{(y,\theta)} : \theta_t \in \mathbb{R}^{\mathbf{P}} \mapsto \frac{\partial y}{\partial \theta} \theta_t \in \mathbb{R}^{\mathbf{Y}} \quad \text{VJP}_{(y,\theta)} : y_c \in \mathbb{R}^{\mathbf{Y}} \mapsto \frac{\partial y}{\partial \theta}^T y_c \in \mathbb{R}^{\mathbf{P}}$$

Neural Network $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$

Batch size \mathbf{N} Output size \mathbf{O}

Forward Pass [FP]



Parameters

$$\theta = \text{vec} [\theta^0, \dots, \theta^{\mathbf{L}}], \quad \mathbf{P} = \sum_{l=0}^{\mathbf{L}} \mathbf{P}^l$$

Intermediate primitive outputs $y^1, \dots, y^{\mathbf{K}}$,

$$\mathbf{Y} = \sum_{k=1}^{\mathbf{K}} \mathbf{Y}^k$$

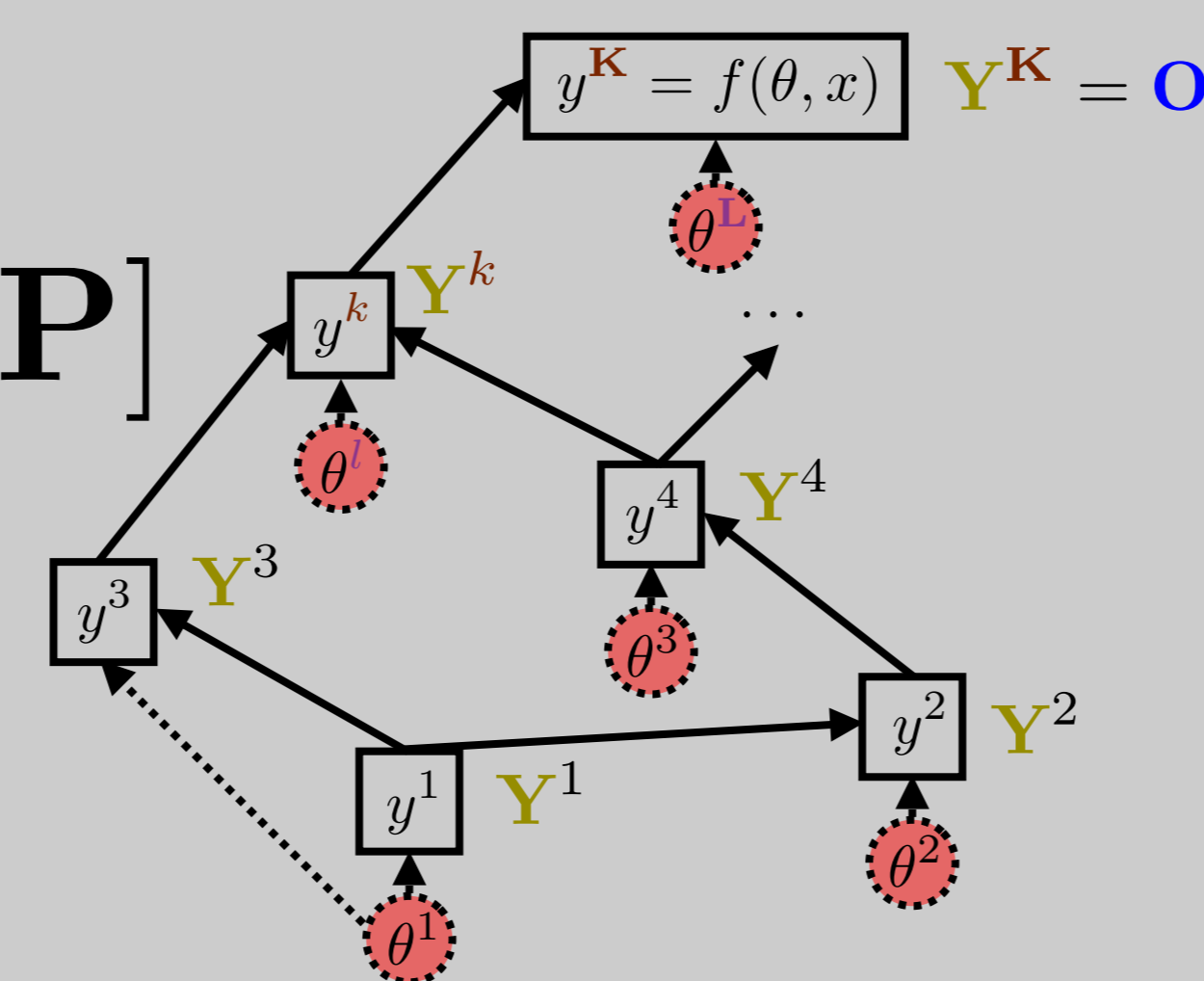
Methods summary

Method	Jacobian Contraction	NTK-vector products	Structured derivatives
Time	\mathbf{NO} [FP] + $\mathbf{N}^2\mathbf{O}^2\mathbf{P}$	$\mathbf{N}^2\mathbf{O}$ [FP]	\mathbf{NO} [FP] + MJJMP

Neural Network $f(\theta, x) \in \mathbb{R}^{\mathbf{O}}$

Batch size \mathbf{N} Output size \mathbf{O}

Forward Pass [FP]



Parameters $\theta = \text{vec} [\theta^0, \dots, \theta^{\mathbf{L}}]$, $\mathbf{P} = \sum_{l=0}^{\mathbf{L}} \mathbf{P}^l$

Intermediate primitive outputs $y^1, \dots, y^{\mathbf{K}}$, $\mathbf{Y} = \sum_{k=1}^{\mathbf{K}} \mathbf{Y}^k$

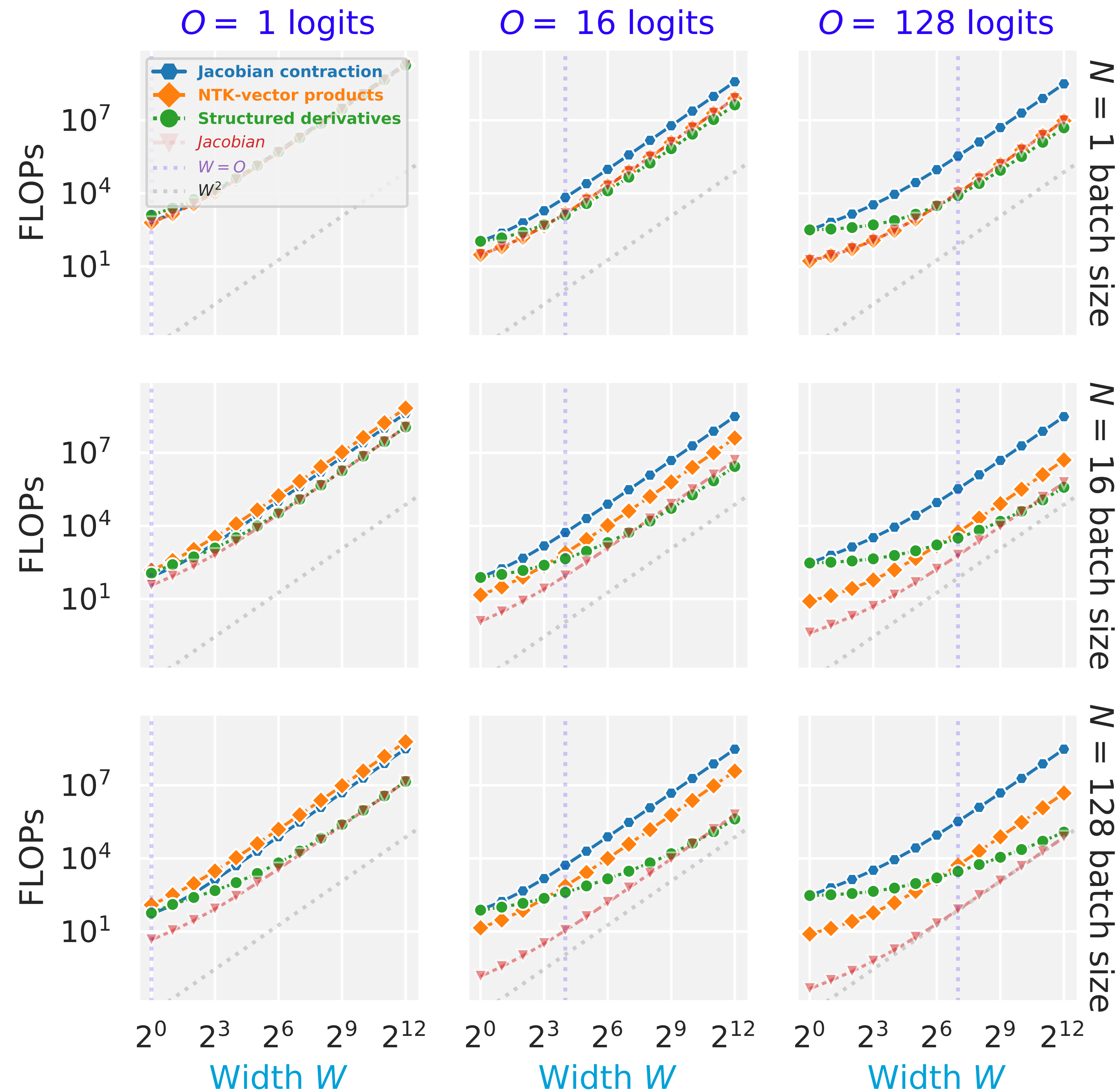
Plan

1. NTK definition and notation
2. Recap of Automatic Differentiation (AD)
3. Baseline NTK computation complexity
4. Two new algorithms for computing the NTK
5. **Benchmarks**

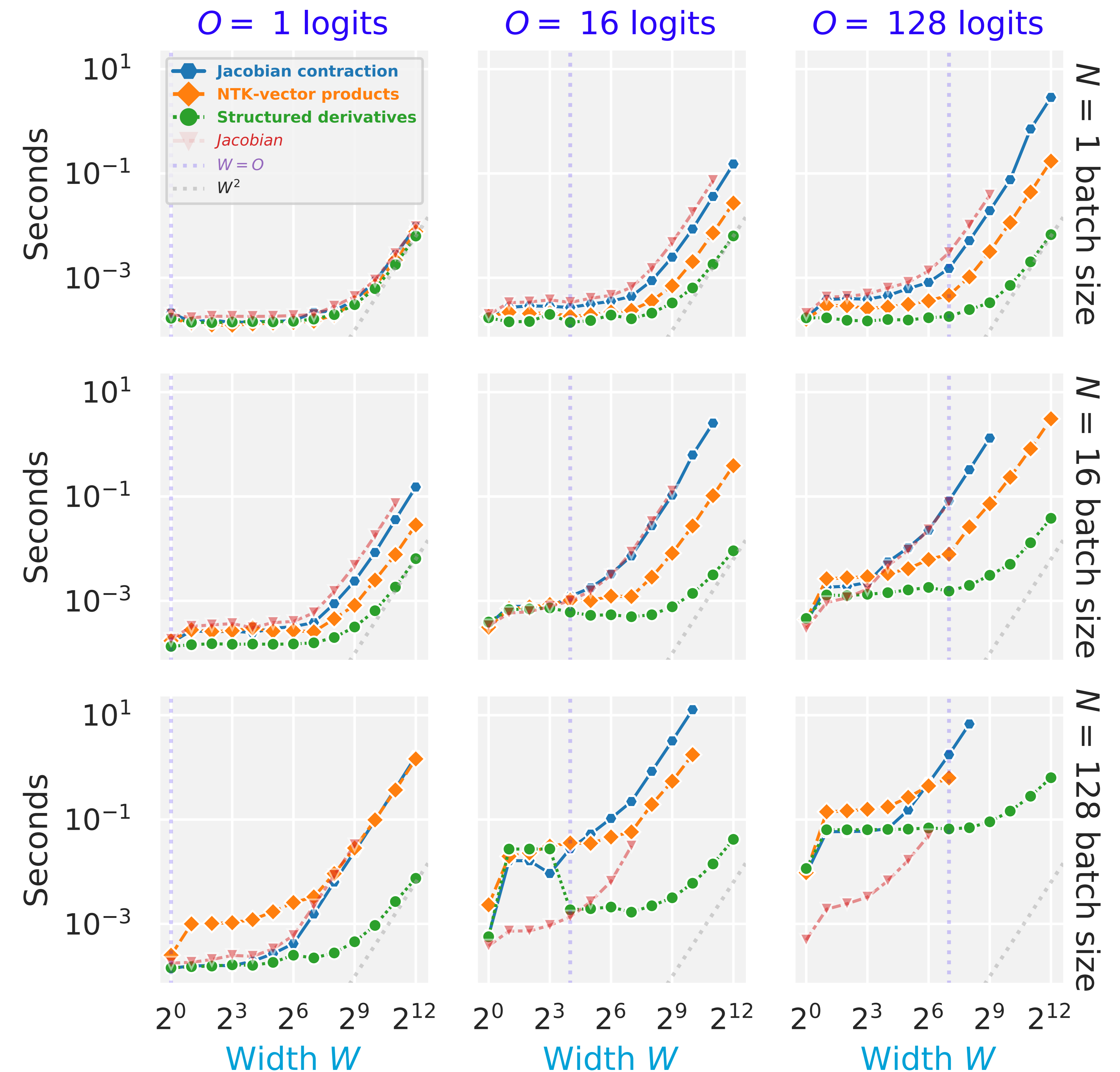
Benchmarks on FCNs

Stop the war
Stand with Ukraine

FLOPs (per NTK entry)

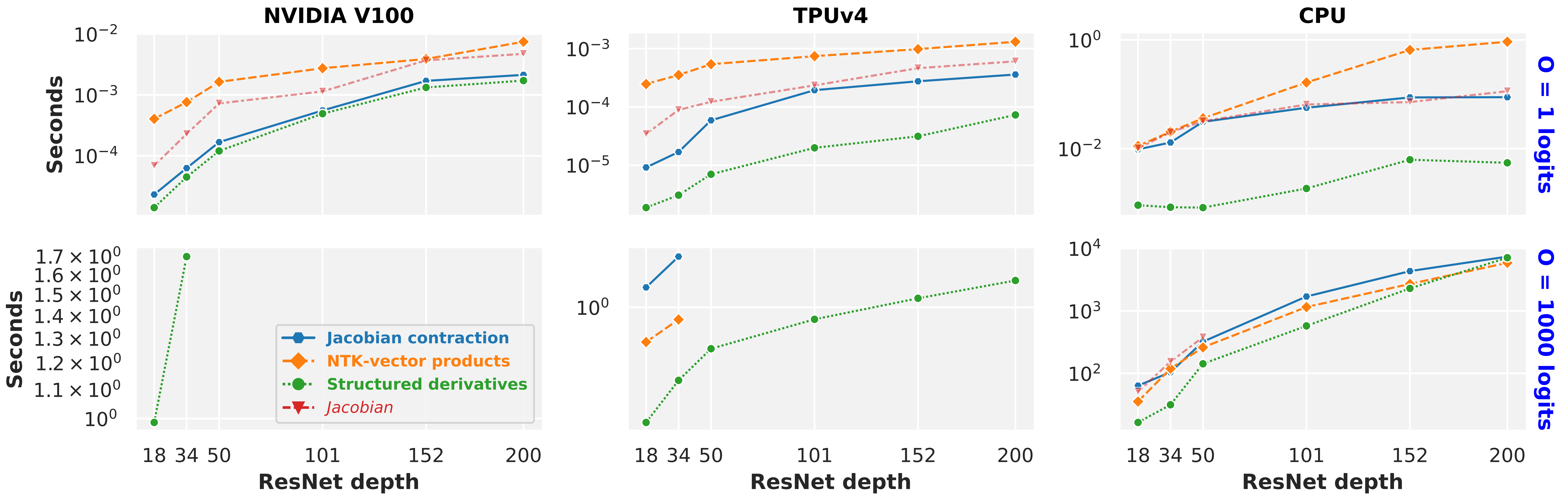


Wall-clock time (TPUv3)

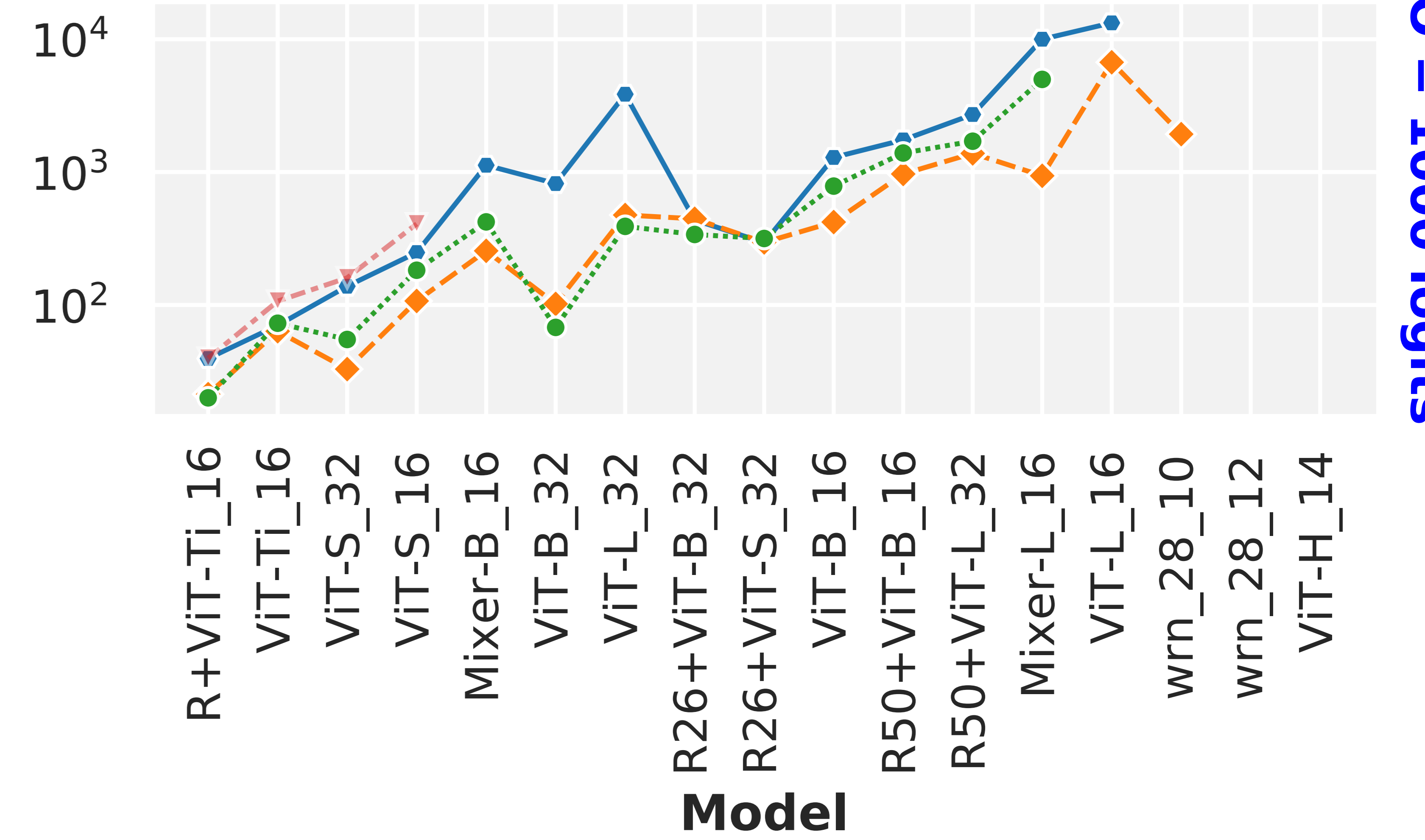
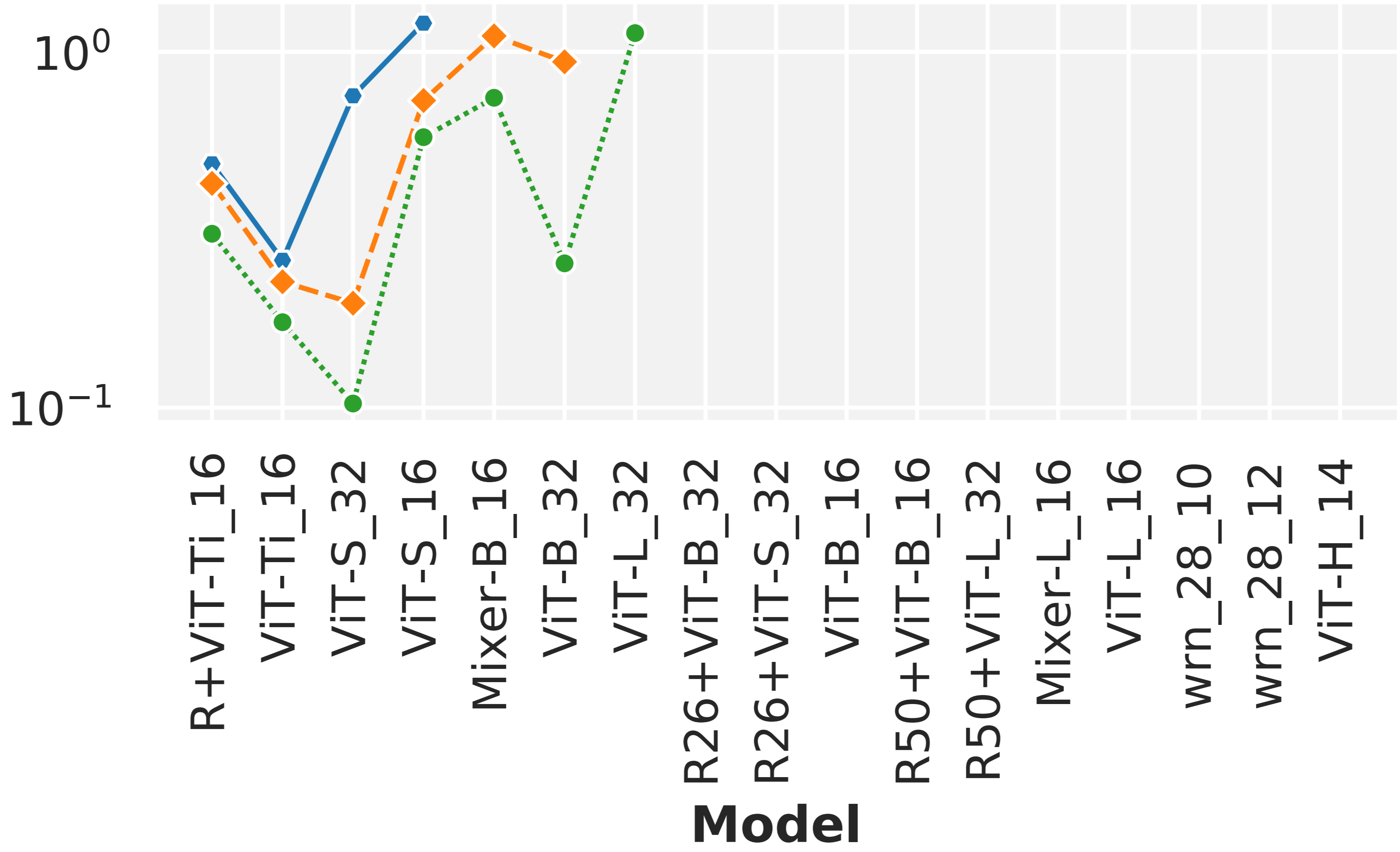
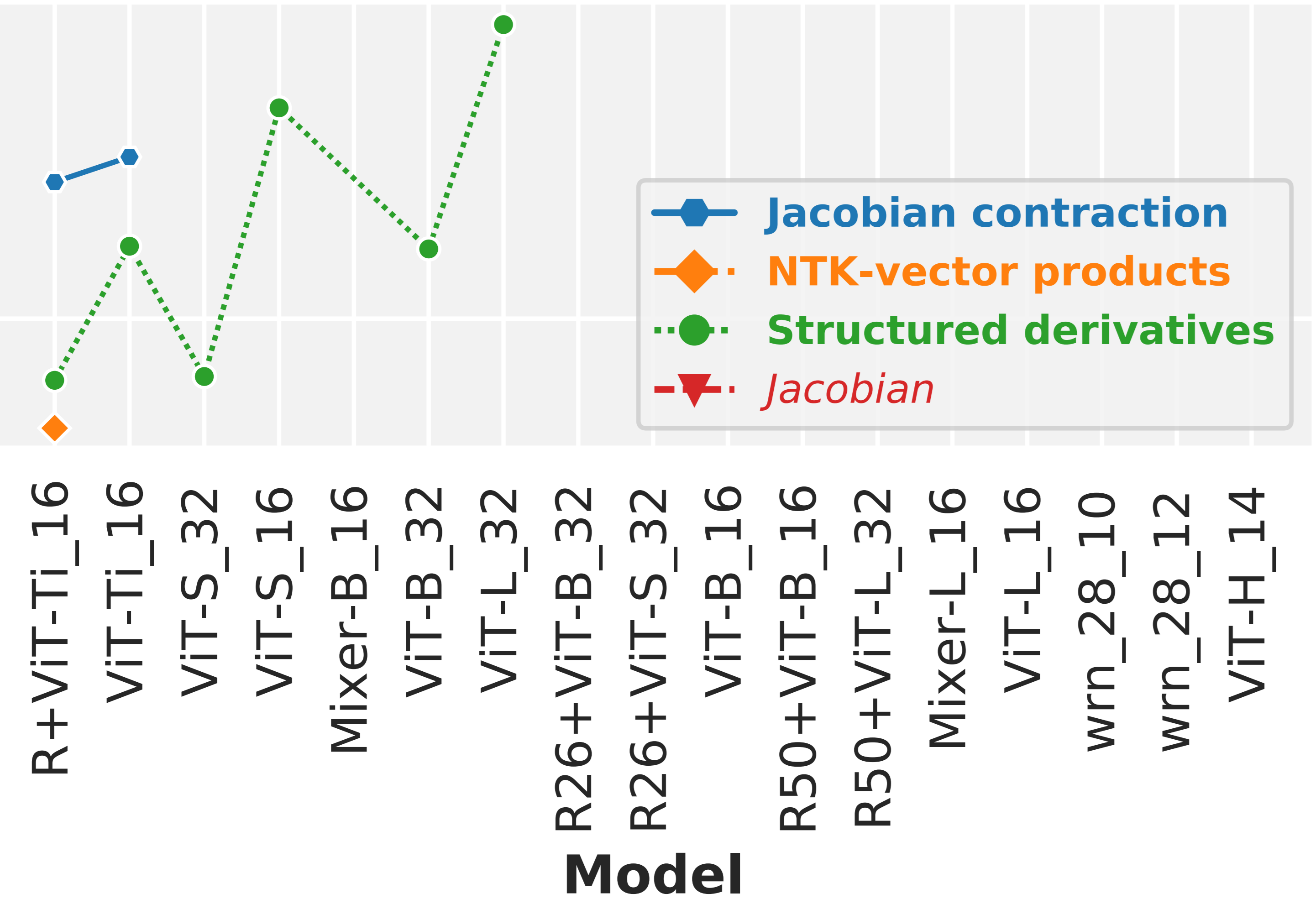
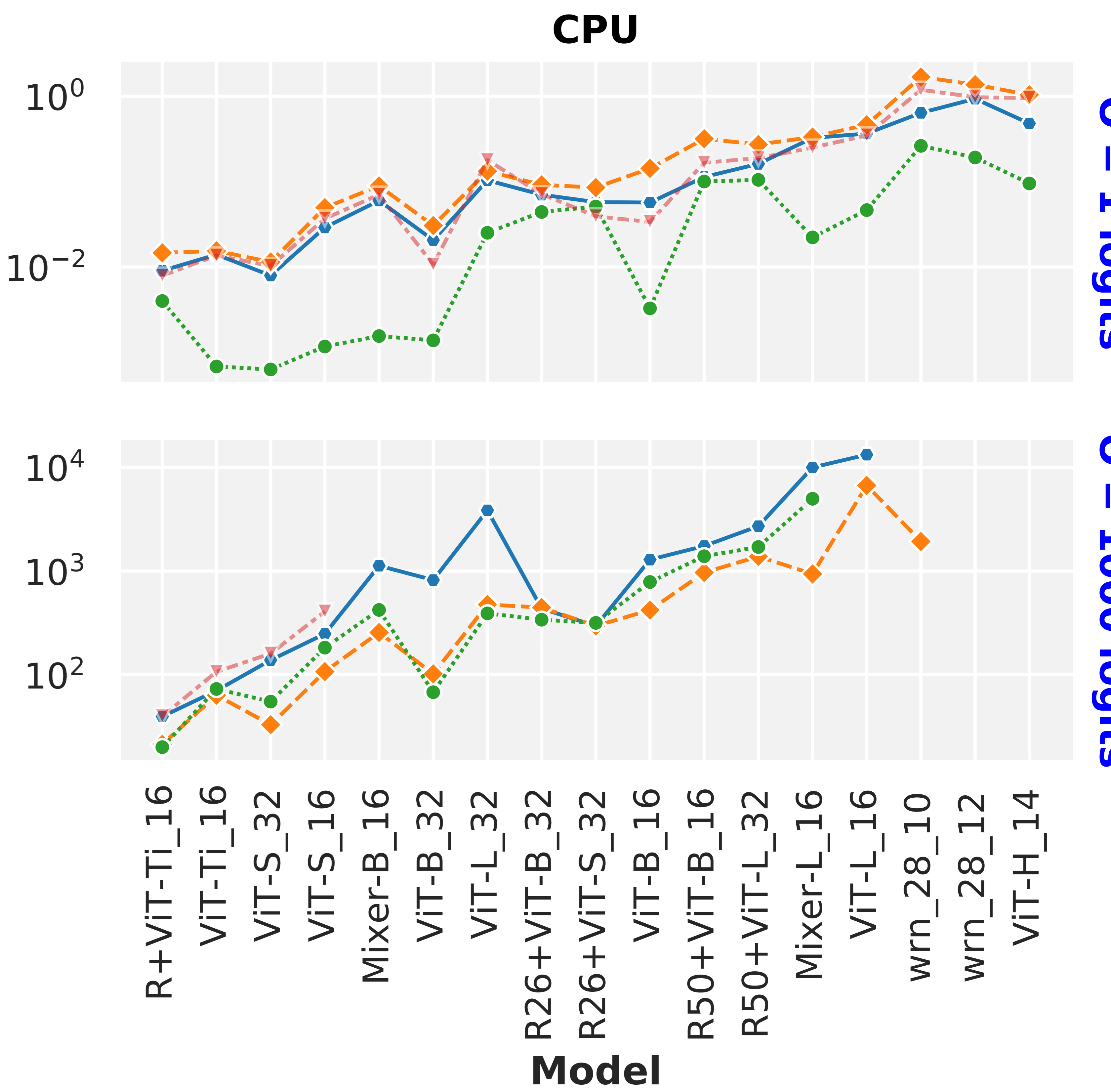
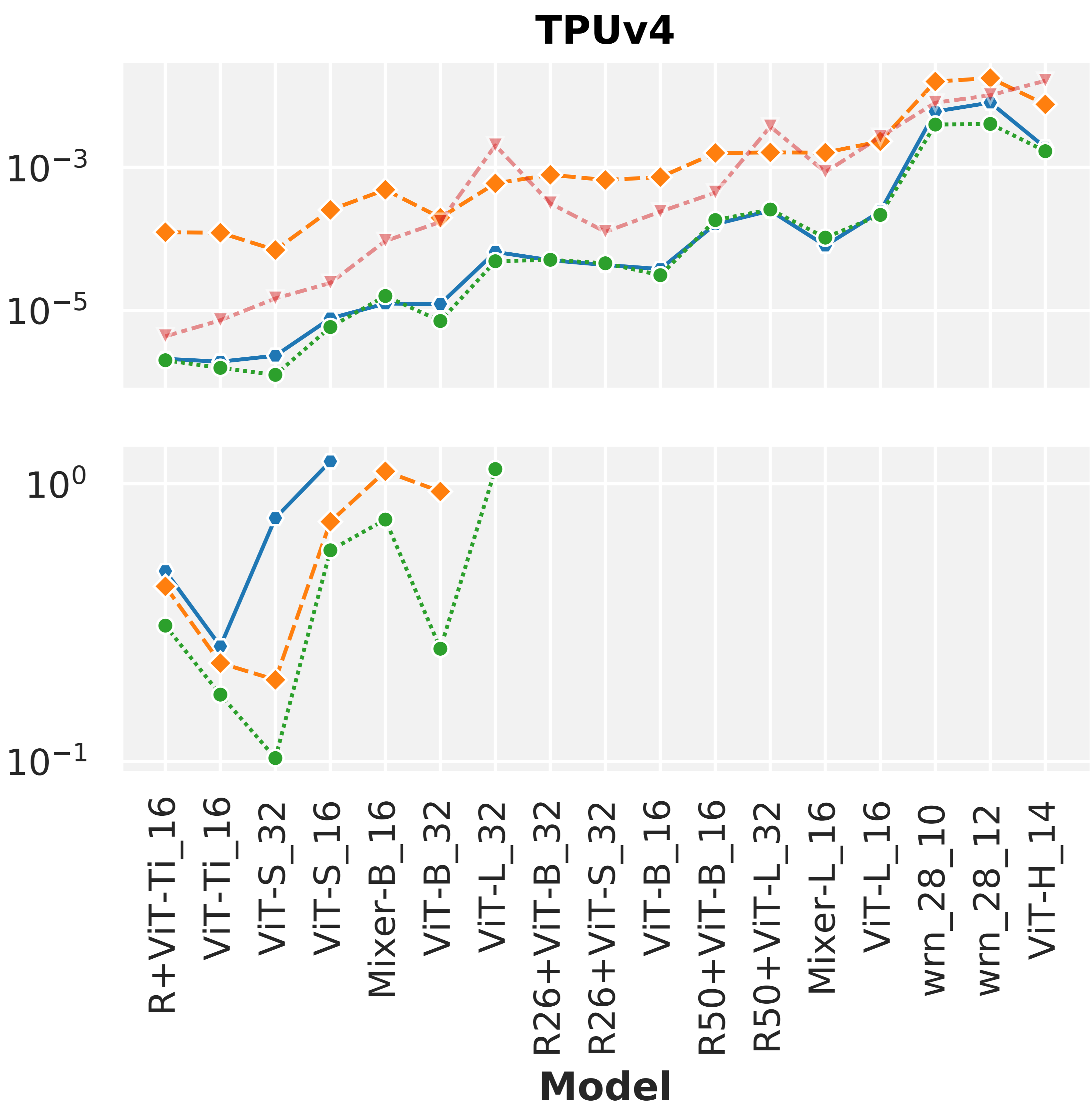
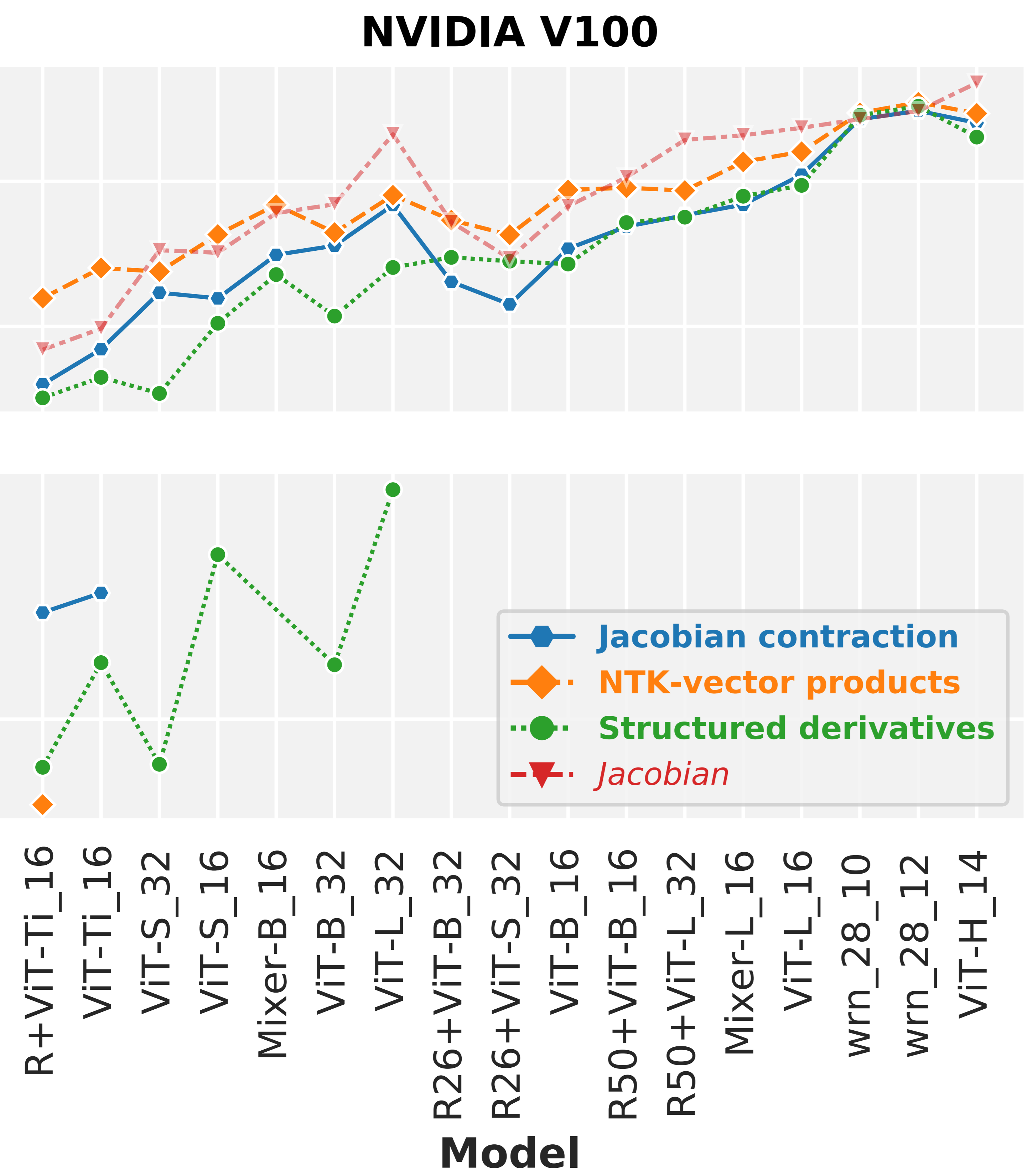
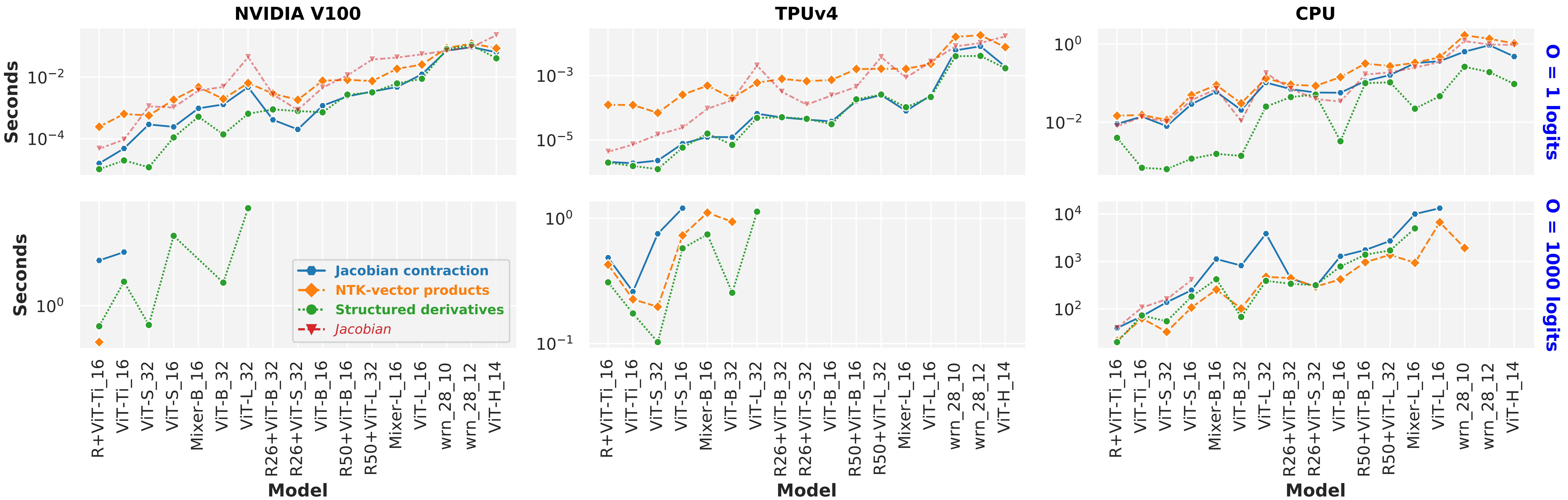


Benchmarks on ImageNet ResNets

Stop the war
Stand with Ukraine



ImageNet Transformers/Mixers/Hybrids



R+ViT-Ti_16
ViT-Ti_16
ViT-S_32
ViT-S_16
Mixer-B_16
ViT-B_32
ViT-L_32
R26+ViT-B_32
R26+ViT-S_32
ViT-B_16
R50+ViT-B_16
R50+ViT-L_32
Mixer-L_16
ViT-L_16
wrn_28_10
wrn_28_12
ViT-H_14

R+ViT-Ti_16
ViT-Ti_16
ViT-S_32
ViT-S_16
Mixer-B_16
ViT-B_32
ViT-L_32
R26+ViT-B_32
R26+ViT-S_32
ViT-B_16
R50+ViT-B_16
R50+ViT-L_32
Mixer-L_16
ViT-L_16
wrn_28_10
wrn_28_12
ViT-H_14

R+ViT-Ti_16
ViT-Ti_16
ViT-S_32
ViT-S_16
Mixer-B_16
ViT-B_32
ViT-L_32
R26+ViT-B_32
R26+ViT-S_32
ViT-B_16
R50+ViT-B_16
R50+ViT-L_32
Mixer-L_16
ViT-L_16
wrn_28_10
wrn_28_12
ViT-H_14

2-line code sample

```
pip install neural_tangents
```

```
import neural_tangents as nt

# Given any f: params, x -> f(params, x)
ntk_fn = nt.empirical_ntk_fn(f, implementation=1) # 1, 2, or 3
ntk = ntk_fn(x1, x2, params)
```

github.com/google/neural-tangents



Stop the war
Stand with Ukraine

SupportUkraineNow.org

