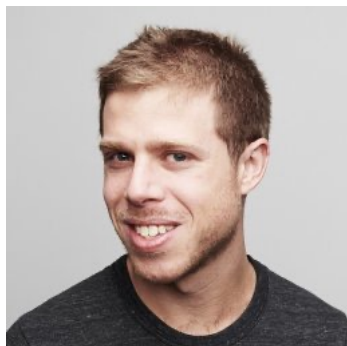


# Effective Theory of Transformers



Emily Dinan, Sho Yaida, Susan Zhang [[arXiv:2304.02034](https://arxiv.org/abs/2304.02034)]

# (Spinoff from) Effective Theory of Deep Neural Networks



Dan Roberts, Sho Yaida, Boris Hanin [[arXiv:2106.10165](https://arxiv.org/abs/2106.10165); [Cambridge University Press](https://www.cambridge.org/9781009048111)]

# (Branch of) Physics of Machine Learning

Order-one fraction of people in this room

# Historical Motivation

Steam Engine



Wonderful Things

# Historical Motivation

Steam Engine

```
graph LR; A[Steam Engine] --> B[Thermodynamics  
Statistical Mechanics]; B --> C[Wonderful Things]
```

Thermodynamics  
Statistical Mechanics

Wonderful Things

# Historical Motivation

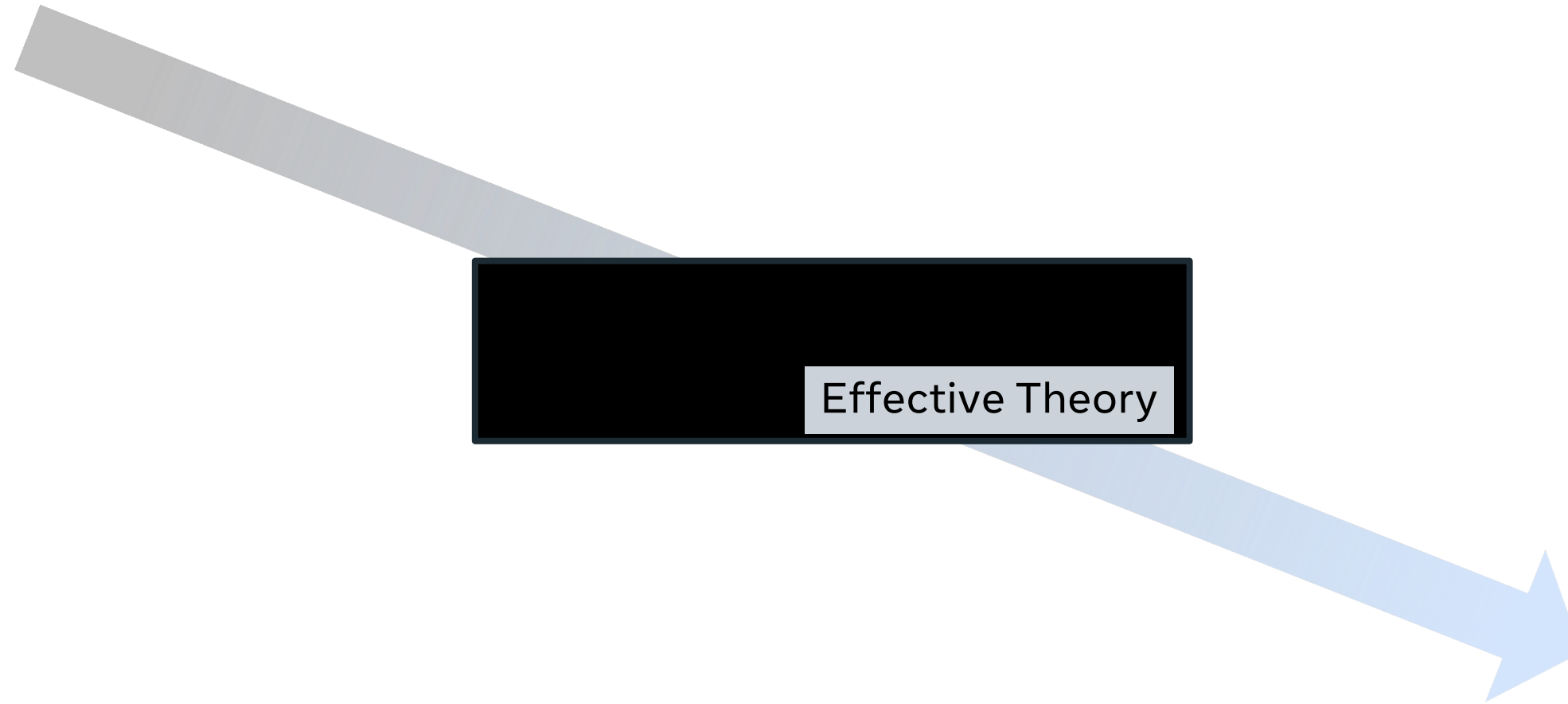
Deep Learning



Wonderful Things

# Historical Motivation

Deep Learning



Wonderful Things

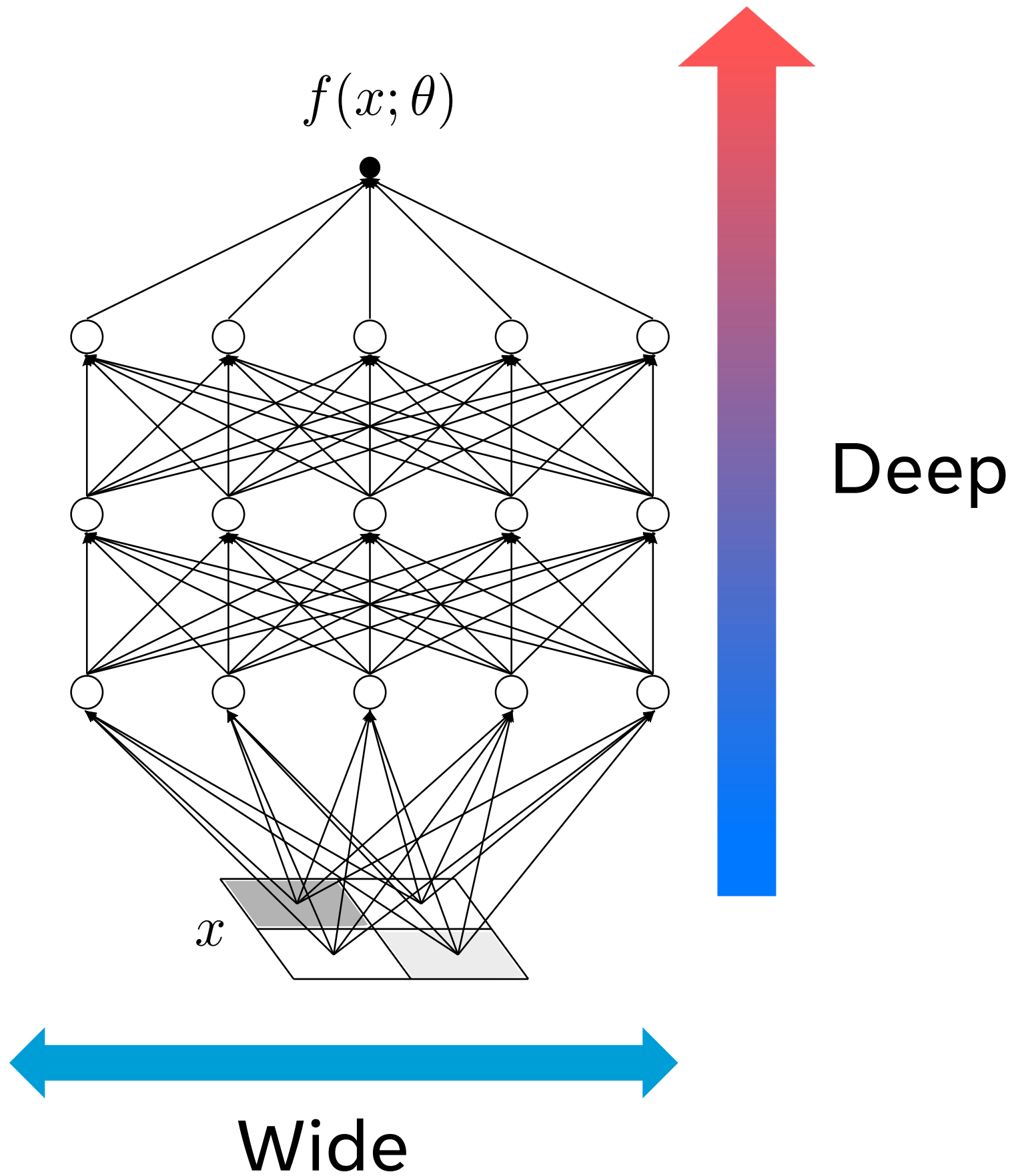
# Outline

1. Simplification @Scale
2. Hyperparameter Scaling for Transformers

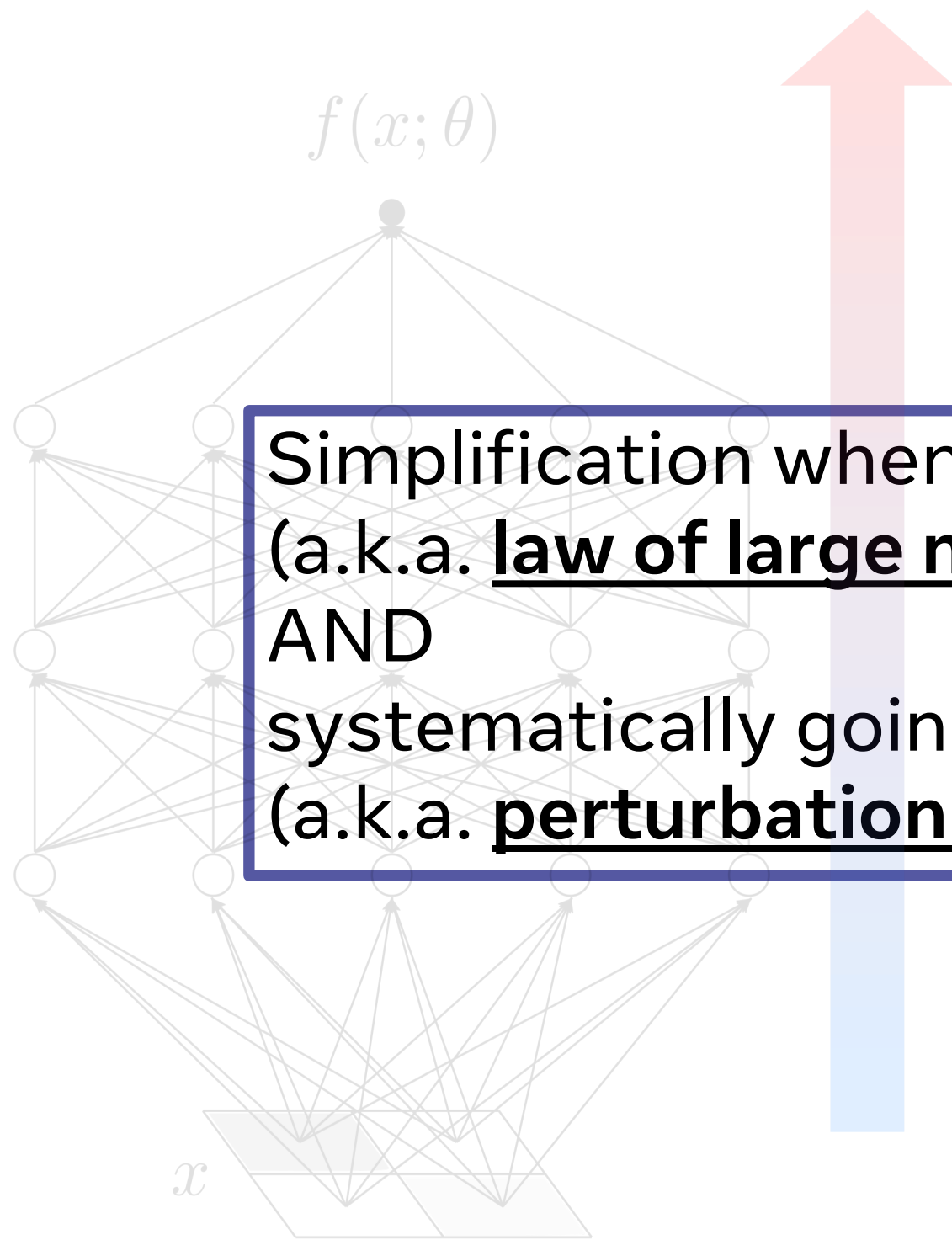


# 1. Simplification @Scale

# Physics of Wide & Deep Neural Networks

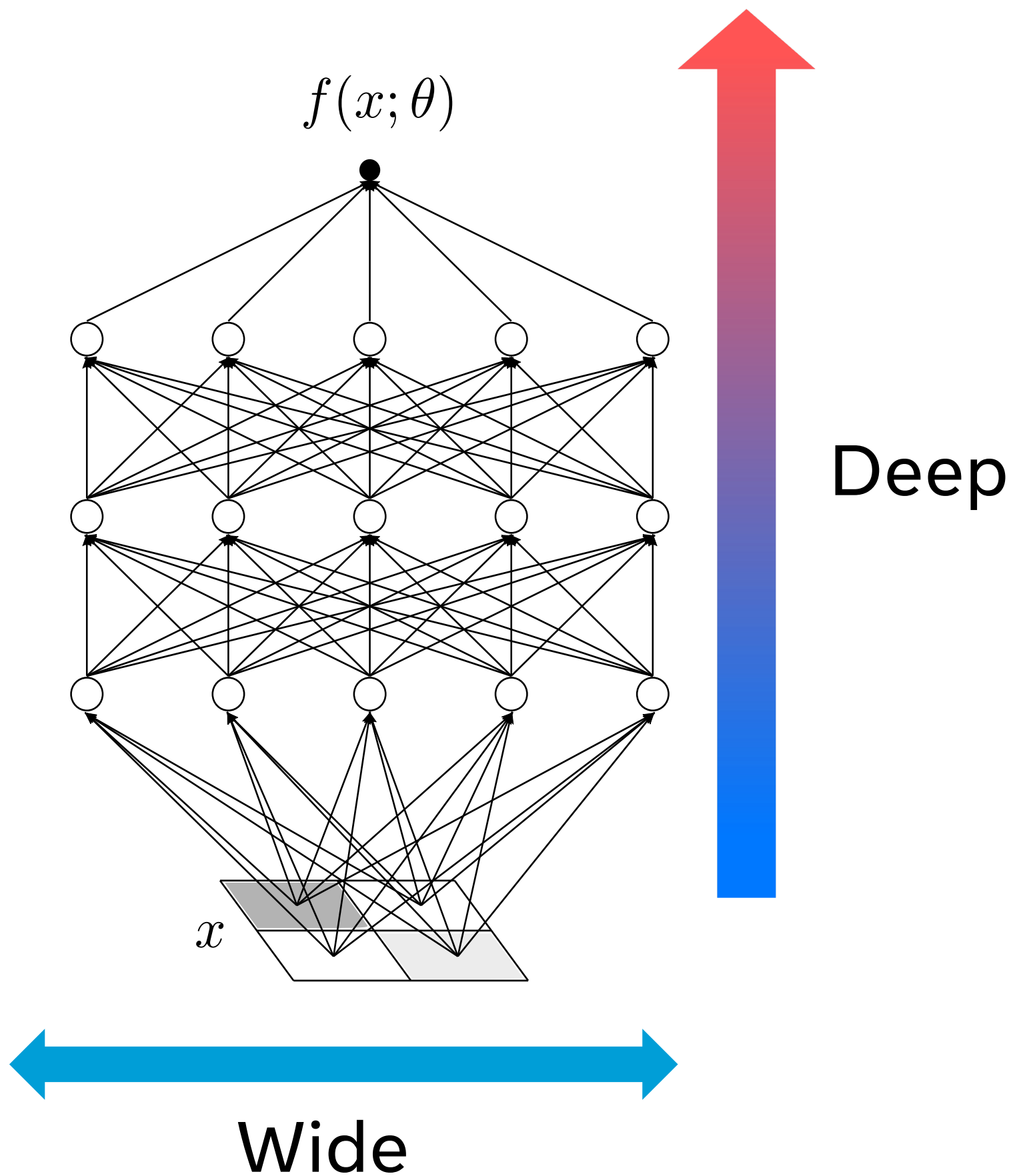


# Physics of Wide & Deep Neural Networks



Simplification when there are infinitely-many neurons in hidden layers  
(a.k.a. **law of large numbers**; a *free theory* at width =  $\infty$  )  
AND  
systematically going beyond that idealized limit  
(a.k.a. **perturbation theory**; a *weakly-interacting theory* when width  $\gg$  depth )

# Physics of Wide & Deep Neural Networks

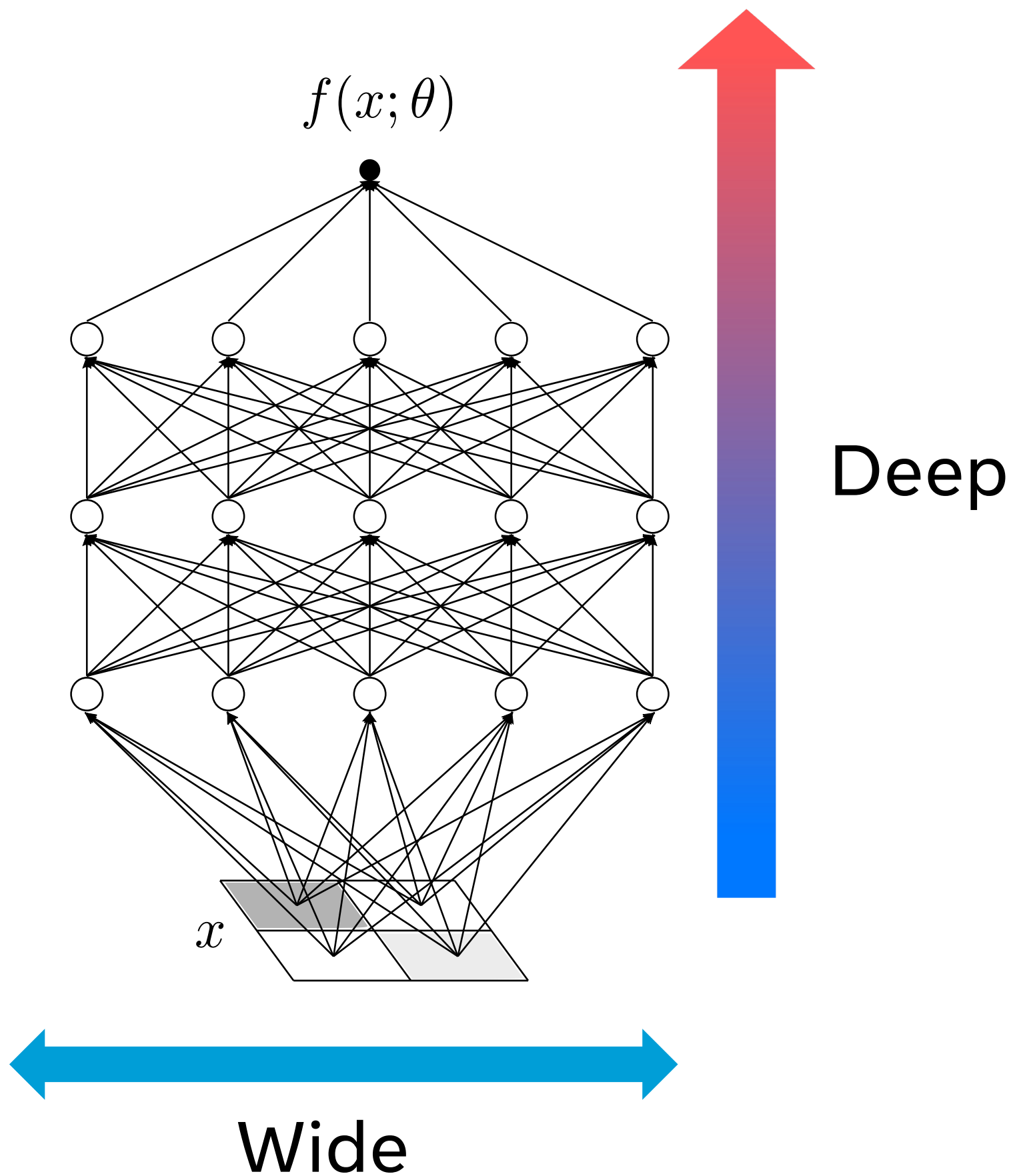


Effective theory works well when

$$\text{width, depth} \gg 1 \quad \& \quad \left( \frac{\text{depth}}{\text{width}} \ll 1 \right)$$

i.e., typical large-scale setups!

# Physics of Wide & Deep Neural Networks



Effective theory works well when

$$\text{width, depth} \gg 1 \quad \& \quad \left( \frac{\text{depth}}{\text{width}} \ll 1 \right)$$

i.e., typical large-scale setups!

[D. Roberts, S. Yaida, B. Hanin: [arXiv:2106.10165](https://arxiv.org/abs/2106.10165); [Cambridge University Press](https://www.cambridge.org/9781009051000) +many many many others]

# Applied Physics of Wide & Deep Neural Networks

Effective theory → how to scale initialization/training hyperparameters

→ Bringing theory & practice closer together @ large scales

# Applied Physics of Wide & Deep Neural Networks

Effective theory → how to scale initialization/training hyperparameters

→ Bringing theory & practice closer together @ large scales

- For Transformers: w/ Emily Dinan and Susan Zhang [[arXiv:2304.02034](https://arxiv.org/abs/2304.02034)]



- ...

# #TL;DR

- **(Meta-)Principle of Criticality:**

“every preactivation/activation should be of order one, i.e., not too big, not too small”

—————→ `stddev ~ np.sqrt(1.0/fan_in)` for weights @ initialization

`fan_in` = number of signals coming into a layer

`fan_out` = number of signals coming out of a layer



# #TL;DR

- **(Meta-)Principle of Criticality:**

“every preactivation/activation should be of order one, i.e., not too big, not too small”

—————→ `stddev ~ np.sqrt(1.0/fan_in)` for weights @ initialization

- **(Meta-)Principle of Equivalence:**

“every layer should contribute equally to learning (that is, to the neural tangent kernel)”

—————→ `lr ~ 1.0/fan_in` for SGD

—————→ `lr ~ 1.0/(fan_in*np.sqrt(fan_out))` for AdamW

`fan_in` = number of signals coming into a layer

`fan_out` = number of signals coming out of a layer

# #TL;DR

- **(Meta-)Principle of Criticality:**

“every preactivation/activation should be of order one, i.e., not too big, not too small”

—————→ `stddev ~ np.sqrt(1.0/fan_in)` for weights @ initialization

- **(Meta-)Principle of Equivalence:**

“every layer should contribute equally to learning (that is, to the neural tangent kernel)”

—————→ `lr ~ 1.0/fan_in` for SGD

—————→ `lr ~ 1.0/(fan_in*np.sqrt(fan_out))` for AdamW

\* Caveat 1: we need to be careful with “fan\_in” for sparse inputs.

\* Caveat 2: stem—head weight tying requires us to rescale output logits.

\* Caveat 3: we can further fine-tune “1.0” but that’s largely moot due to LayerNorms.

\* Meta-Caveat:

preceding and following discussions assume the *neural-tangent* scaling of various hyperparameters; for more general cases including the *maximal-update* scaling (proposed by Greg Yang & Edward Hu), see

[arXiv:2210.04909](https://arxiv.org/abs/2210.04909) (S. Yaida, “Meta-Principled Family of Hyperparameter Scaling Strategies”)

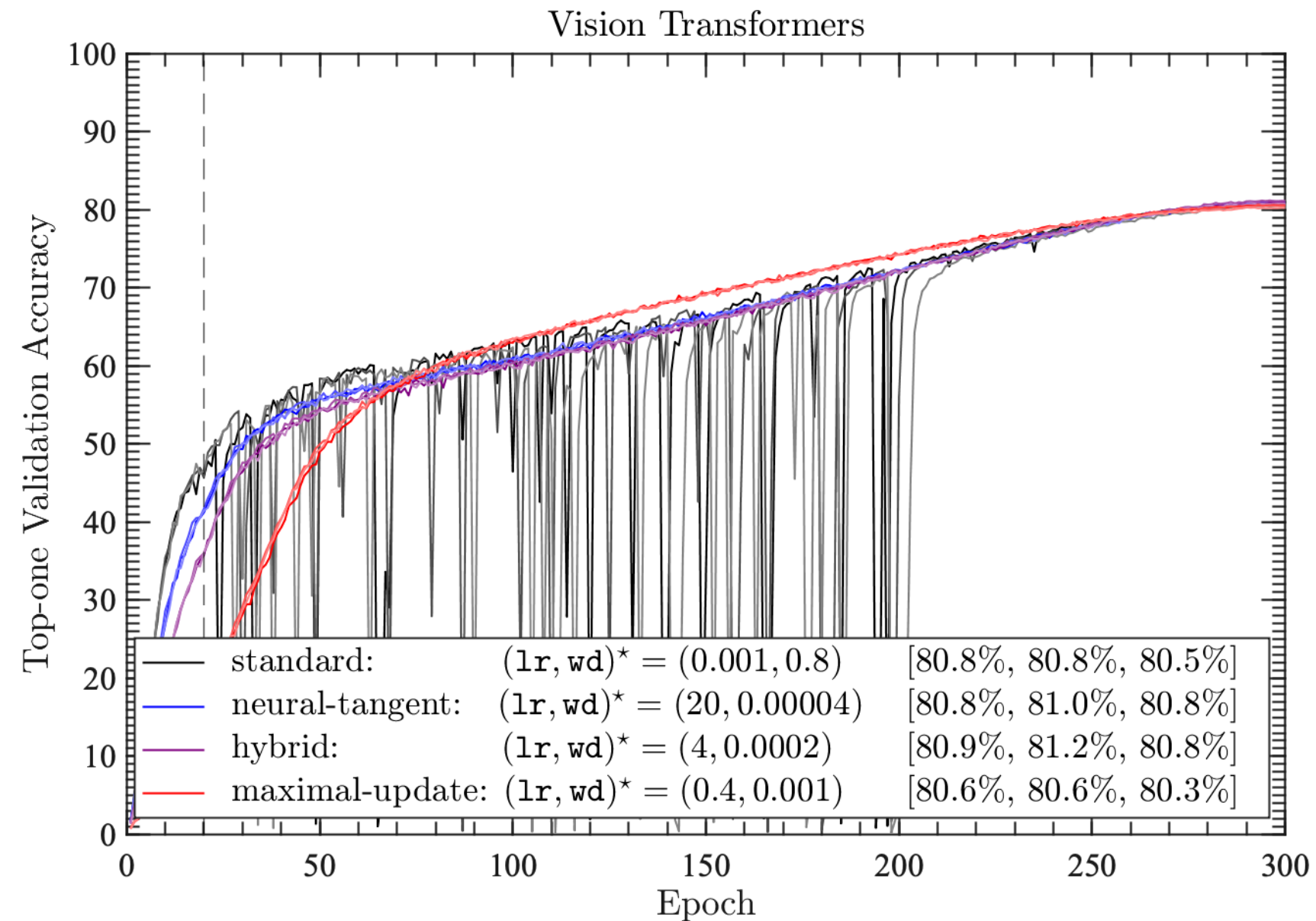
or

various footnotes in our paper [[arXiv:2304.02034](https://arxiv.org/abs/2304.02034)] under presentation

# 2. Hyperparameter Scaling for Transformers

`bias=elementwise_affine=False`

# Does it work?

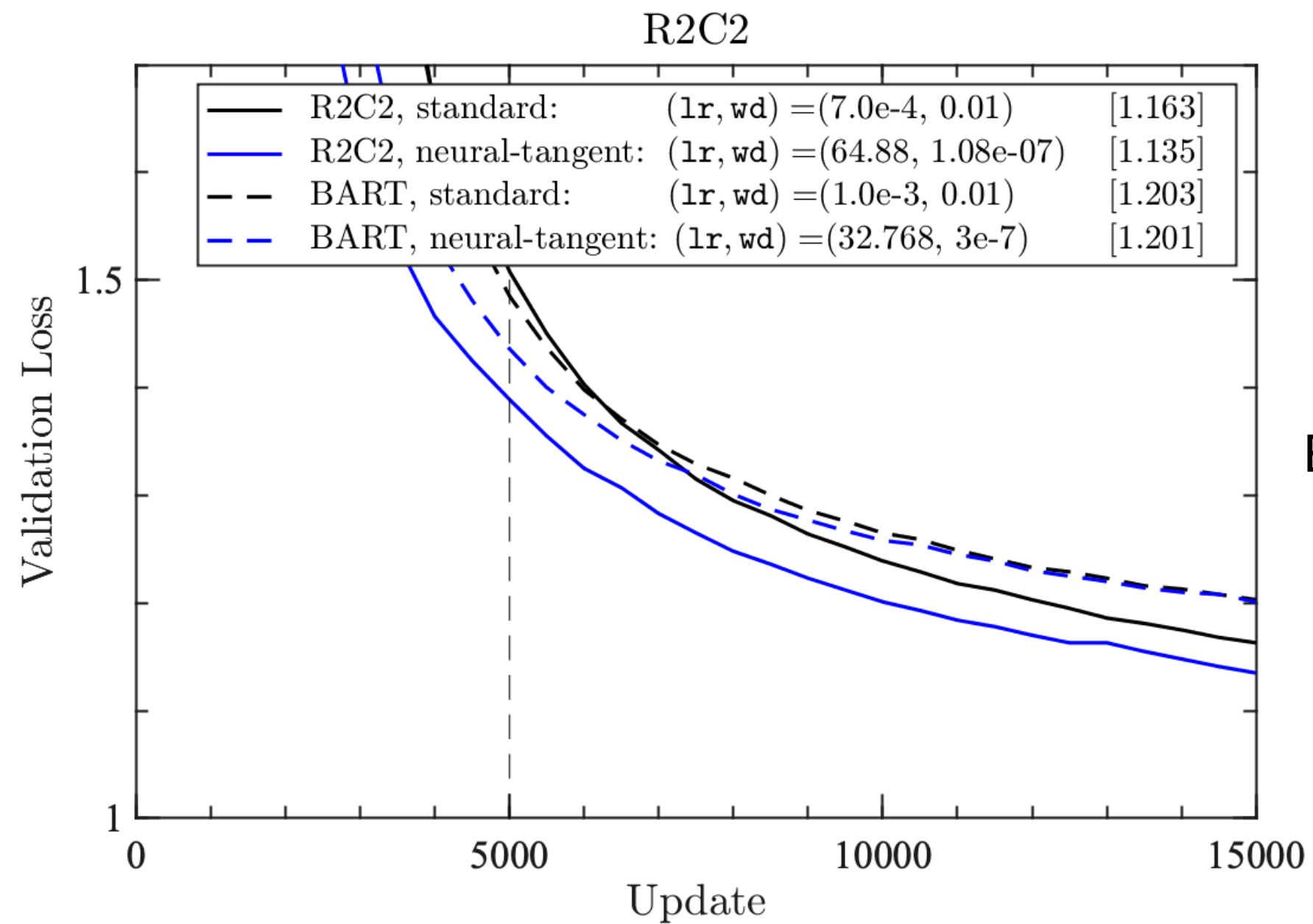


It seems to work well here. 😊

One curious observation:

neural-tangent scaling *increases* the learning rate for positional-embedding parameters by  $\sim 1000$  relative to the standard uniform scaling, and it makes the training *stabler*.

# Does it work?



Encouraging but not conclusive...

# Technical-looking Summary 1: Initialization Hyperparameters

`stddev ~ np.sqrt(1.0/fan_in)`

that their covariances should be scaled as

$$\mathbb{E} [W_{i_1 j_1}^{\text{patch}} W_{i_2 j_2}^{\text{patch}}] = \left( \frac{C_{\text{patch}}}{n_{\text{patch}}} \right) \delta_{i_1 i_2} \delta_{j_1 j_2} \quad \parallel \quad \mathbb{E} [W_{i_1 j_1}^{\text{WE}} W_{i_2 j_2}^{\text{WE}}] = (C_{\text{WE}}) \delta_{i_1 i_2} \delta_{j_1 j_2}, \quad (1.25)$$

$$\mathbb{E} [b_{t_1; i_1}^{\text{PE}} b_{t_2; i_2}^{\text{PE}}] = (C_{\text{PE}}) \delta_{t_1 t_2} \delta_{i_1 i_2}, \quad (1.26)$$

$$\mathbb{E} [Q_{c_1 i_1}^{h_1} Q_{c_2 i_2}^{h_2}] = \left( \frac{C_Q}{n} \right) \delta_{c_1 c_2} \delta_{i_1 i_2} \delta^{h_1 h_2}, \quad (1.27)$$

$$\mathbb{E} [K_{c_1 i_1}^{h_1} K_{c_2 i_2}^{h_2}] = \left( \frac{C_K}{n} \right) \delta_{c_1 c_2} \delta_{i_1 i_2} \delta^{h_1 h_2}, \quad (1.28)$$

$$\mathbb{E} [V_{c_1 i_1}^{h_1} V_{c_2 i_2}^{h_2}] = \left( \frac{C_V}{n} \right) \delta_{c_1 c_2} \delta_{i_1 i_2} \delta^{h_1 h_2}, \quad (1.29)$$

$$\mathbb{E} [U_{i_1 c_1}^{h_1} U_{i_2 c_2}^{h_2}] = \left( \frac{C_U}{n} \right) \delta_{i_1 i_2} \delta_{c_1 c_2} \delta^{h_1 h_2}, \quad (1.30)$$

$$\mathbb{E} [W_{i_1 j_1} W_{i_2 j_2}] = \left( \frac{C_W}{n} \right) \delta_{i_1 i_2} \delta_{j_1 j_2}, \quad (1.31)$$

$$\mathbb{E} [X_{i_1 j_1} X_{i_2 j_2}] = \left( \frac{C_X}{Mn} \right) \delta_{i_1 i_2} \delta_{j_1 j_2}, \quad (1.32)$$

$$\mathbb{E} [W_{i_1 j_1}^{\text{head}} W_{i_2 j_2}^{\text{head}}] = \left( \frac{C_{\text{head}}}{n} \right) \delta_{i_1 i_2} \delta_{j_1 j_2}, \quad b_i^{\text{head}} = 0 \quad \parallel \quad \mathcal{N}_{\text{rescale}} = \sqrt{\frac{1}{n}}, \quad (1.33)$$

where  $\mathbb{E}[\cdot]$  denotes an expectation value with respect to the initialization distribution,  $\delta_{ij}$  is the Kronecker delta (i.e.,  $\delta_{ij} = 1$  when  $i = j$  and  $\delta_{ij} = 0$  when  $i \neq j$ ), and all the initialization hyperparameters  $C_G$  are order-one numbers for each group  $G$  of model parameters. Here, by

# Technical-looking Summary 2: Training Hyperparameters (LR)

of  $\sim |g^\mu|$  modifies the above relative learning-rate factors to  $\boxed{\text{lr} \sim 1.0 / (\text{fan\_in} * \text{np.sqrt}(\text{fan\_out}))}$

$$\tilde{\lambda}_{W^{\text{patch}}} = \left( \frac{1}{n_{\text{patch}} \sqrt{n}} \right) \tilde{\Lambda}_{W^{\text{patch}}} \quad || \quad \tilde{\lambda}_{W^{\text{WE}}} = \left( \frac{1}{\sqrt{n}} \right) \tilde{\Lambda}_{W^{\text{WE}}}, \quad (1.108)$$

$$\tilde{\lambda}_{b^{\text{PE}}} = \left( \frac{1}{\sqrt{n}} \right) \tilde{\Lambda}_{b^{\text{PE}}}, \quad (1.109)$$

$$\tilde{\lambda}_Q = \left( \frac{1}{n \sqrt{n}} \right) \tilde{\Lambda}_Q, \quad (1.110)$$

$$\tilde{\lambda}_K = \left( \frac{1}{n \sqrt{n}} \right) \tilde{\Lambda}_K, \quad (1.111)$$

$$\tilde{\lambda}_V = \left( \frac{1}{n \sqrt{n}} \right) \tilde{\Lambda}_V, \quad (1.112)$$

$$\tilde{\lambda}_U = \left( \frac{1}{n \sqrt{n}} \right) \tilde{\Lambda}_U, \quad (1.113)$$

$$\tilde{\lambda}_W = \left( \frac{1}{n \sqrt{Mn}} \right) \tilde{\Lambda}_W, \quad (1.114)$$

$$\tilde{\lambda}_X = \left( \frac{1}{Mn \sqrt{n}} \right) \tilde{\Lambda}_X, \quad (1.115)$$

$$\tilde{\lambda}_{W^{\text{head}}} = \left( \frac{1}{n \sqrt{n_{\text{out}}}} \right) \tilde{\Lambda}_{W^{\text{head}}}, \quad \tilde{\lambda}_{b^{\text{head}}} = \left( \frac{1}{\sqrt{n_{\text{out}}}} \right) \tilde{\Lambda}_{b^{\text{head}}}, \quad (1.116)$$

where – like the order-one initialization hyperparameters  $C_G$ 's –  $\tilde{\Lambda}_G$ 's are order-one training hyperparameters which we could in principle tune but won't.<sup>25</sup>



```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

# Attention Blocks

$$\text{stddev} \sim \text{np.sqrt}(1.0/\text{fan\_in})$$

$$\text{lr} \sim 1.0/(\text{fan\_in} * \text{np.sqrt}(\text{fan\_out}))$$

The query, key, value, and out\_proj weights all have (\*)

`fan_in=embedding_dim`

`fan_out=embedding_dim`

(\*Ok, with the traditional usage of terminology, strictly speaking, it depends on how you code QKV, but it doesn't really matter much here.)

# Attention Blocks


```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

The query, key, value, and out\_proj weights all have (\*)

```
fan_in=embedding_dim
```

```
fan_out=embedding_dim
```



```
stddev=np.sqrt(1.0/embedding_dim)  
lr=np.sqrt(1.0/embedding_dim**3)
```

(\*Ok, with the traditional usage of terminology, strictly speaking, it depends on how you code QKV, but it doesn't really matter much here.)

# Attention Blocks

```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

One comment about the famous query—key dot product normalization:

$$\frac{1}{\sqrt{C}} \mathbf{q} \cdot \mathbf{k} \quad \text{where } C = \# \text{ of channels per head} = \frac{\text{embedding\_dim}}{\text{n\_heads}}$$

# Attention Blocks

```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

One comment about the famous query—key dot product normalization:

$$\frac{1}{\sqrt{C}} \mathbf{q} \cdot \mathbf{k} \quad \text{where } C = \# \text{ of channels per head} = \frac{\text{embedding\_dim}}{\text{n\_heads}}$$

This is theoretically correct (while Tensor Program V paper does something else).

# FC Blocks

$$\text{stddev} \sim \text{np.sqrt}(1.0/\text{fan\_in})$$

$$\text{lr} \sim 1.0/(\text{fan\_in} * \text{np.sqrt}(\text{fan\_out}))$$

The 1st layer

`fan_in=embedding_dim`

`fan_out=M*embedding_dim`

The 2nd layer

`fan_in=M*embedding_dim`

`fan_out=embedding_dim`


# FC Blocks

$$\text{stddev} \sim \text{np.sqrt}(1.0/\text{fan\_in})$$

$$\text{lr} \sim 1.0/(\text{fan\_in} * \text{np.sqrt}(\text{fan\_out}))$$

The 1st layer


$\text{fan\_in} = \text{embedding\_dim}$   
 $\text{fan\_out} = M * \text{embedding\_dim}$



$\text{stddev} = \text{np.sqrt}(1.0/\text{embedding\_dim})$   
 $\text{lr} = \text{np.sqrt}(1.0/(M * (\text{embedding\_dim} ** 3)))$

The 2nd layer

$\text{fan\_in} = M * \text{embedding\_dim}$   
 $\text{fan\_out} = \text{embedding\_dim}$



$\text{stddev} = \text{np.sqrt}(1.0/(M * \text{embedding\_dim}))$   
 $\text{lr} = \text{np.sqrt}(1.0/((M ** 2) * (\text{embedding\_dim} ** 3)))$

# Word-Embedding Parameters

$$\text{stddev} \sim \text{np.sqrt}(1.0/\text{fan\_in})$$

$$\text{lr} \sim 1.0/(\text{fan\_in} * \text{np.sqrt}(\text{fan\_out}))$$

Naively...,

`fan_in=n_vocab`

`fan_out=embedding_dim`



# Word-Embedding Parameters

Because of one-hot-ness,

~~fan\_in=n\_vocab~~ fan\_in=1

fan\_out=embedding\_dim

```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

# Word-Embedding Parameters

```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

Because of one-hot-ness,

~~fan\_in=n\_vocab~~ fan\_in=1

fan\_out=embedding\_dim



```
stddev=1.0
```

```
lr=np.sqrt(1.0/embedding_dim)
```

# Word-Embedding Parameters

```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

Because of one-hot-ness,

~~fan\_in=n\_vocab~~ fan\_in=1

fan\_out=embedding\_dim



```
stddev=1.0
```

```
lr=np.sqrt(1.0/embedding_dim)
```

\* Lesson 1: be mindful with sparse features.

\* Lesson 2: the learning rate is relatively much higher than those for other layers.

# (Word-Embedding Parameters)<sup>T</sup>

```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

Because of stem—head weight tying,

$$\text{logit}_i = \sum_{j=1}^{\text{embedding\_dim}} (\text{WE.weight})_{ji} (\text{last\_features})_j$$

# (Word-Embedding Parameters)<sup>T</sup>

$$\text{stddev} \sim \text{np.sqrt}(1.0/\text{fan\_in})$$

$$\text{lr} \sim 1.0/(\text{fan\_in} * \text{np.sqrt}(\text{fan\_out}))$$

Because of stem—head weight tying,

$$\text{logit}_i = \sum_{j=1}^{\text{embedding\_dim}} (\text{WE.weight})_{ji} (\text{last\_features})_j \sim O(\sqrt{\text{embedding\_dim}})$$

# (Word-Embedding Parameters)<sup>T</sup>

$$\text{stddev} \sim \text{np.sqrt}(1.0/\text{fan\_in})$$

$$\text{lr} \sim 1.0/(\text{fan\_in} * \text{np.sqrt}(\text{fan\_out}))$$

Because of stem—head weight tying,

$$\text{logit}_i = \sum_{j=1}^{\text{embedding\_dim}} (\text{WE.weight})_{ji} (\text{last\_features})_j \sim O(\sqrt{\text{embedding\_dim}})$$



$$\text{logit} = \text{logit} / \text{np.sqrt}(\text{embedding\_dim})$$

# Positional-Embedding Parameters

```
stddev ~ np.sqrt(1.0/fan_in)
```

```
lr ~ 1.0/(fan_in*np.sqrt(fan_out))
```

Finally, quickly, for positional-embedding (additive, like bias), we have

```
fan_in=1
```

```
fan_out=embedding_dim
```



```
stddev=1.0 or 0.0 or something  
lr=np.sqrt(1.0/embedding_dim)
```

# Technical-looking Summary 1: Initialization Hyperparameters

`stddev ~ np.sqrt(1.0/fan_in)`

that their covariances should be scaled as

$$\mathbb{E} [W_{i_1 j_1}^{\text{patch}} W_{i_2 j_2}^{\text{patch}}] = \left( \frac{C_{\text{patch}}}{n_{\text{patch}}} \right) \delta_{i_1 i_2} \delta_{j_1 j_2} \quad \parallel \quad \mathbb{E} [W_{i_1 j_1}^{\text{WE}} W_{i_2 j_2}^{\text{WE}}] = (C_{\text{WE}}) \delta_{i_1 i_2} \delta_{j_1 j_2}, \quad (1.25)$$

$$\mathbb{E} [b_{t_1; i_1}^{\text{PE}} b_{t_2; i_2}^{\text{PE}}] = (C_{\text{PE}}) \delta_{t_1 t_2} \delta_{i_1 i_2}, \quad (1.26)$$

$$\mathbb{E} [Q_{c_1 i_1}^{h_1} Q_{c_2 i_2}^{h_2}] = \left( \frac{C_Q}{n} \right) \delta_{c_1 c_2} \delta_{i_1 i_2} \delta^{h_1 h_2}, \quad (1.27)$$

$$\mathbb{E} [K_{c_1 i_1}^{h_1} K_{c_2 i_2}^{h_2}] = \left( \frac{C_K}{n} \right) \delta_{c_1 c_2} \delta_{i_1 i_2} \delta^{h_1 h_2}, \quad (1.28)$$

$$\mathbb{E} [V_{c_1 i_1}^{h_1} V_{c_2 i_2}^{h_2}] = \left( \frac{C_V}{n} \right) \delta_{c_1 c_2} \delta_{i_1 i_2} \delta^{h_1 h_2}, \quad (1.29)$$

$$\mathbb{E} [U_{i_1 c_1}^{h_1} U_{i_2 c_2}^{h_2}] = \left( \frac{C_U}{n} \right) \delta_{i_1 i_2} \delta_{c_1 c_2} \delta^{h_1 h_2}, \quad (1.30)$$

$$\mathbb{E} [W_{i_1 j_1} W_{i_2 j_2}] = \left( \frac{C_W}{n} \right) \delta_{i_1 i_2} \delta_{j_1 j_2}, \quad (1.31)$$

$$\mathbb{E} [X_{i_1 j_1} X_{i_2 j_2}] = \left( \frac{C_X}{Mn} \right) \delta_{i_1 i_2} \delta_{j_1 j_2}, \quad (1.32)$$

$$\mathbb{E} [W_{i_1 j_1}^{\text{head}} W_{i_2 j_2}^{\text{head}}] = \left( \frac{C_{\text{head}}}{n} \right) \delta_{i_1 i_2} \delta_{j_1 j_2}, \quad b_i^{\text{head}} = 0 \quad \parallel \quad \mathcal{N}_{\text{rescale}} = \sqrt{\frac{1}{n}}, \quad (1.33)$$

where  $\mathbb{E}[\cdot]$  denotes an expectation value with respect to the initialization distribution,  $\delta_{ij}$  is the Kronecker delta (i.e.,  $\delta_{ij} = 1$  when  $i = j$  and  $\delta_{ij} = 0$  when  $i \neq j$ ), and all the initialization hyperparameters  $C_G$  are order-one numbers for each group  $G$  of model parameters. Here, by



# Technical-looking Summary 2: Training Hyperparameters (LR)

of  $\sim |g^\mu|$  modifies the above relative learning-rate factors to  $\boxed{\text{lr} \sim 1.0 / (\text{fan\_in} * \text{np.sqrt}(\text{fan\_out}))}$

$$\tilde{\lambda}_{W^{\text{patch}}} = \left( \frac{1}{n_{\text{patch}} \sqrt{n}} \right) \tilde{\Lambda}_{W^{\text{patch}}} \quad || \quad \tilde{\lambda}_{W^{\text{WE}}} = \left( \frac{1}{\sqrt{n}} \right) \tilde{\Lambda}_{W^{\text{WE}}}, \quad (1.108)$$

$$\tilde{\lambda}_{b^{\text{PE}}} = \left( \frac{1}{\sqrt{n}} \right) \tilde{\Lambda}_{b^{\text{PE}}}, \quad (1.109)$$

$$\tilde{\lambda}_Q = \left( \frac{1}{n \sqrt{n}} \right) \tilde{\Lambda}_Q, \quad (1.110)$$

$$\tilde{\lambda}_K = \left( \frac{1}{n \sqrt{n}} \right) \tilde{\Lambda}_K, \quad (1.111)$$

$$\tilde{\lambda}_V = \left( \frac{1}{n \sqrt{n}} \right) \tilde{\Lambda}_V, \quad (1.112)$$

$$\tilde{\lambda}_U = \left( \frac{1}{n \sqrt{n}} \right) \tilde{\Lambda}_U, \quad (1.113)$$

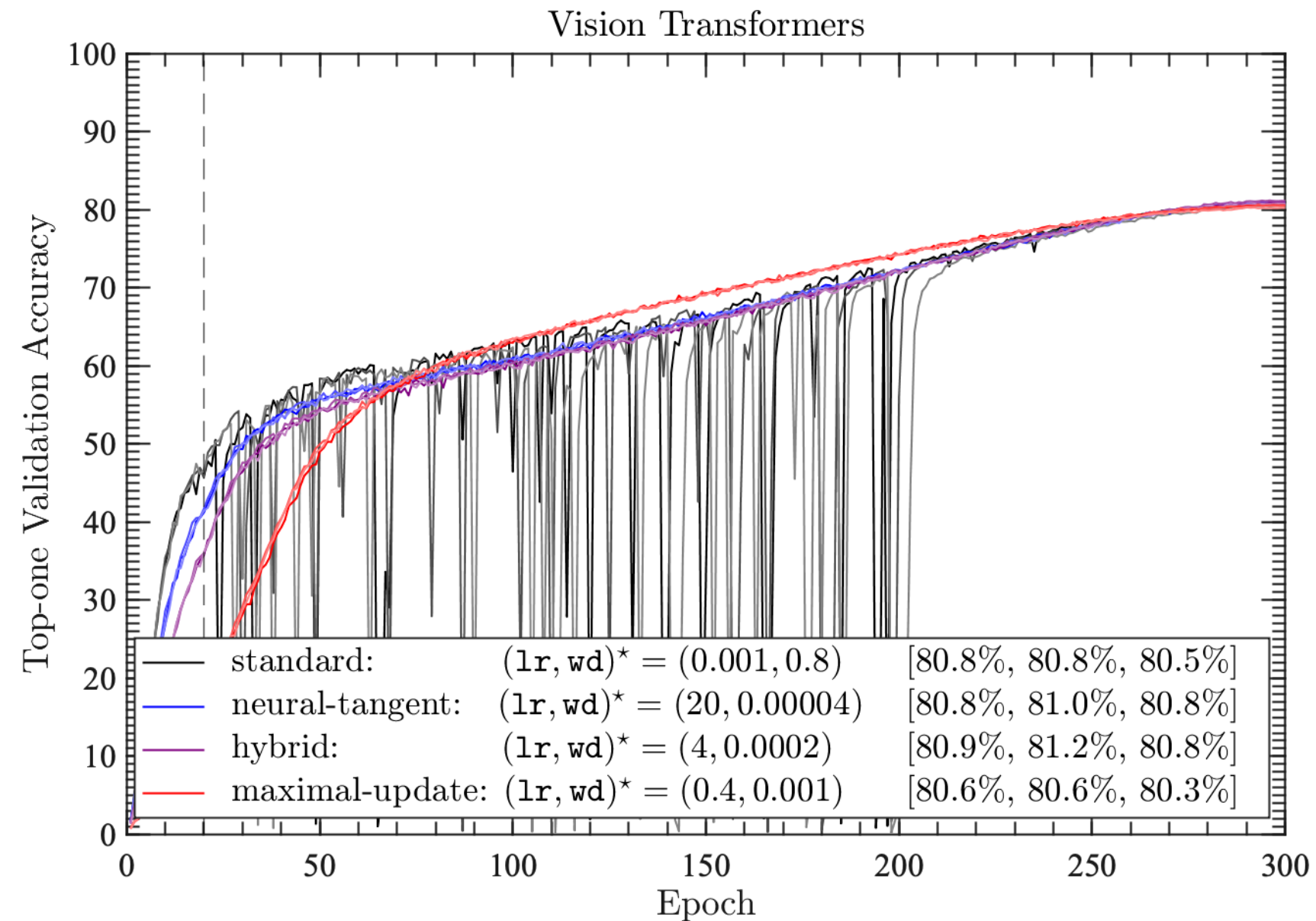
$$\tilde{\lambda}_W = \left( \frac{1}{n \sqrt{Mn}} \right) \tilde{\Lambda}_W, \quad (1.114)$$

$$\tilde{\lambda}_X = \left( \frac{1}{Mn \sqrt{n}} \right) \tilde{\Lambda}_X, \quad (1.115)$$

$$\tilde{\lambda}_{W^{\text{head}}} = \left( \frac{1}{n \sqrt{n_{\text{out}}}} \right) \tilde{\Lambda}_{W^{\text{head}}}, \quad \tilde{\lambda}_{b^{\text{head}}} = \left( \frac{1}{\sqrt{n_{\text{out}}}} \right) \tilde{\Lambda}_{b^{\text{head}}}, \quad (1.116)$$

where – like the order-one initialization hyperparameters  $C_G$ 's –  $\tilde{\Lambda}_G$ 's are order-one training hyperparameters which we could in principle tune but won't.<sup>25</sup>

# Does it work?

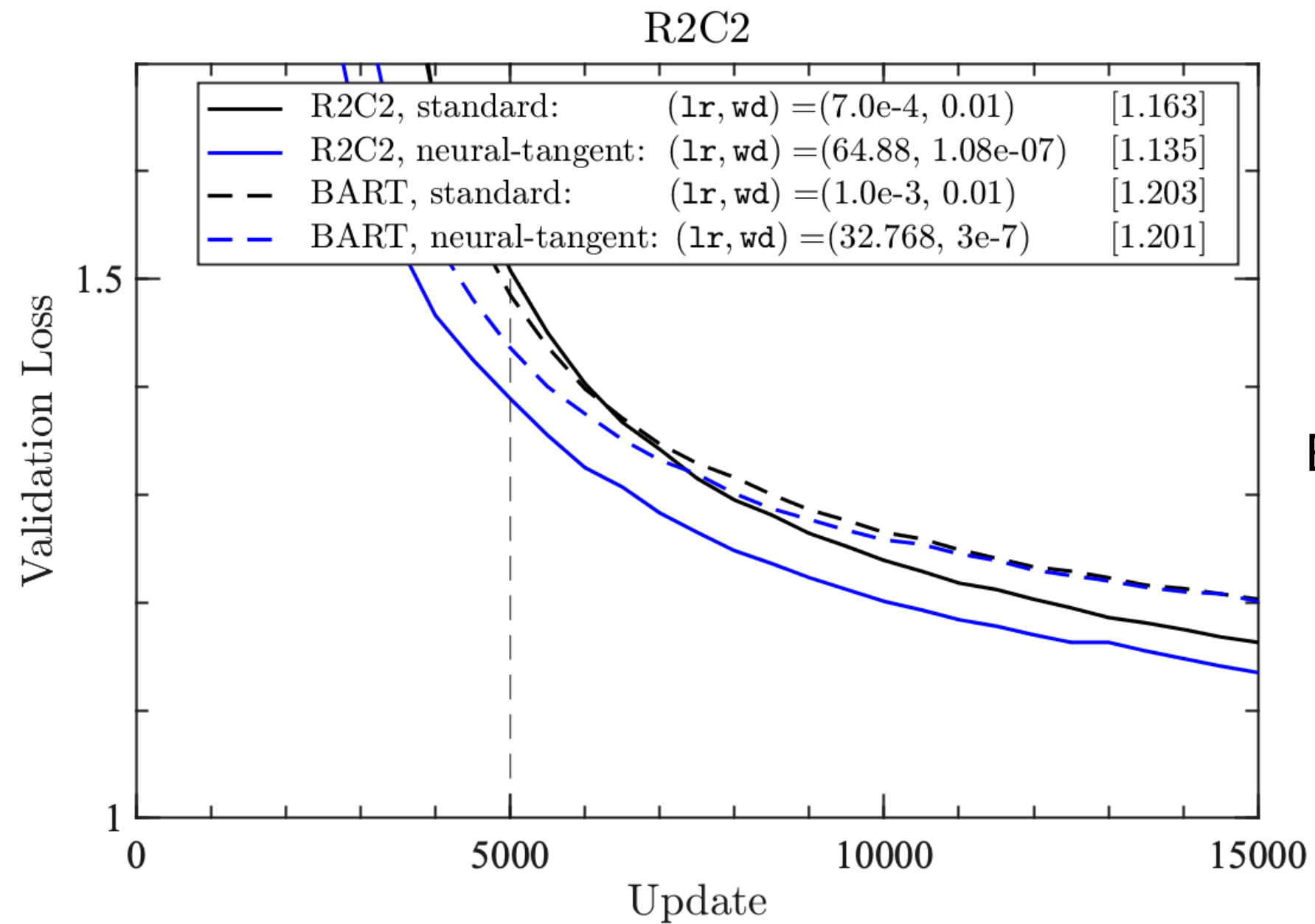


It seems to work well here. 😊

One curious observation:

neural-tangent scaling *increases* the learning rate for positional-embedding parameters by  $\sim 1000$  relative to the standard uniform scaling, and it makes the training *stabler*.

# Does it work?



Encouraging but not conclusive...

Does it work?

**lots of  
encouraging/discouraging anecdotal  
to come — from you**