# Feature Learning in Infinite-Width Neural Networks

#### Greg Yang

Microsoft Research AI

Presenting the 4<sup>th</sup> Paper of the *Tensor Programs* Series

Joint work with ex-Microsoft AI Resident Edward Hu

#### Feature Learning is Crucial in Deep Learning

Imagenet and Resnet



#### Pretraining and Transfer via Finetuning



- Discard readout layer from pretraining
- *Linear Finetuning*: Train new readout layer only
- *Total Finetuning*: Train the entire network
  - Our conclusions will apply here as well

#### Current Theory: Neural Tangent Kernel

- Linear evolution easy to analyze
- Yields optimization and generalization results



#### But NTK Limit Does Not Learn Features!

#### Word2Vec example



Figure 1: PCA of Word2Vec embeddings of common US cities and states, for NTK, width-64, and width- $\infty$  feature learning neural networks. NTK embeddings are essentially random, while cities and states get naturally separated in embedding space as width increases in the feature learning regime.

#### Feature Learning $\infty$ -Width NN on Real Tasks

Word2Vec

- Pretraining: learn word embeddings so that similar words have close embeddings, in an unsupervised manner.
- Transfer: word analogy task, e.g. "what to a queen is as a man to a woman?"



Similar Results on Metalearning (MAML)

#### Overview

Our experiments verify this is the right limit

#### **Our Theoretical Contributions**

- Classify all "viable" ∞-width limits into feature learning and kernel limits
- Identify "the maximal" feature learning limit
- Propose the *Tensor Programs* technique for deriving its equations, and more generally, the limit of any neural computation

#### Significance

- Framework for studying feature learning in overparametrized NN
- Formulas for training featurelearning ∞-width NN in variety of settings (e.g. pretraining, metalearning, reinforcement learning, GANs, etc)
- Mostly solves the problem of taking ∞-width limits

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The *Tensor Programs* technique for deriving the limit

#### • Parametrizations of Neural Networks

- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The *Tensor Programs* technique for deriving the limit

#### abc-Parametrizations

- Consider *L*-hidden-layer perceptron
  - Input dim d, output dim 1, width n, nonlinearity  $\phi$ , input  $\xi \in \mathbb{R}^d$
  - Weights  $W^1 \in \mathbb{R}^{n \times d}$  and  $W^2, \dots, W^L \in \mathbb{R}^{n \times n}$  and  $W^{L+1} \in \mathbb{R}^{1 \times n}$

• 
$$h^1(\xi) = W^1 \xi \in \mathbb{R}^n$$

- $x^{l}(\xi) = \phi(h^{l}(\xi)) \in \mathbb{R}^{n}, h^{l+1}(\xi) = W^{l+1}x^{l}(\xi) \in \mathbb{R}^{n}$  for l = 1, ..., L-1
- Network output (i.e. logits)  $f(\xi) = W^{L+1}x^{L}(\xi)$

An *abc-parametrization* is given by a set of numbers  $\{a_l, b_l\}_l \cup \{c\}$  s.t.

- a) Parametrize each  $W^{l} = n^{-a_{l}}w^{l}$  where  $w^{l}$  is trained instead of  $W^{l}$
- b) Initialize each  $w_{\alpha\beta}^{l} \sim \mathcal{N}(0, n^{-2b_{l}})$
- c) SGD learning rate is  $\eta n^{-c}$  for some width-independent  $\eta$ .

#### Examples

 $h^{1}(\xi) = W^{1}\xi \in \mathbb{R}^{n}, \ x^{l}(\xi) = \phi\left(h^{l}(\xi)\right) \in \mathbb{R}^{n}, \ h^{l+1}(\xi) = W^{l+1}x^{l}(\xi) \in \mathbb{R}^{n}, \ f(\xi) = W^{L+1}x^{L}(\xi)$ 

	Definition	ΝΤΚ	Standard (Pytorch default)	Mean Field ( $L=1$ )		
a <sub>l</sub>	$W^l = n^{-a_l} w^l$	$\begin{cases} 0 & \text{if } l = 1 \\ \frac{1}{2} & \text{if } l > 1 \end{cases}$	0	$\begin{cases} 0 & \text{if } l = 1 \\ 1 & \text{if } l = 2 \end{cases}$		
b <sub>l</sub>	$w_{\alpha\beta}^l \sim \mathcal{N}(0, n^{-2b_l})$	0	$\begin{cases} 0 & \text{if } l = 1 \\ \frac{1}{2} & \text{if } l > 1 \end{cases}$	0		
С	$LR = \eta n^{-c}$	0	0	-1		
	Ĩ					
We	are ignoring dependences should be thought of	on input dim $d$ (which as a constant)	By changing $a_l, b_l$ while give each layer	By changing $a_l, b_l$ while fixing $a_l + b_l$ , we effectively give each layer $l$ its own learning rate.		

#### Why Do We Care About Parametrizations?

- Each parametrization admits an ∞width limit
  - abc-parametrizations correspond to natural class of limits, including NTK, GP, Mean Field, etc
- Parametrizations are important practically
  - E.g. Glorot, He, Fixup



#### • Parametrizations of Neural Networks

- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The *Tensor Programs* technique for deriving the limit

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The *Tensor Programs* technique for deriving the limit

#### Instability and Triviality

$$LR = \eta n^{-c}$$

- If LR too large w.r.t width (i.e. c too small), then logits or preactivations blow up with width during training
  - We say this abc-parametrization is *unstable*
- If LR too small w.r.t width (i.e. c too big), then the NN function doesn't evolve in the ∞-width limit
  - We say this abc-parametrization is *trivial*
- Only *nontrivial stable* abc-parametrizations are meaningful

#### Dynamical Dichotomy Theorem

- Any *nontrivial stable* abc-parametrization yields a (discrete-time) ∞-width limit
- When trained for any finite time, this limit either
  - 1. (Feature learning limit) allows the embedding  $x^{L}(\xi)$  to evolve nontrivially or
  - 2. (Kernel limit) is described by kernel gradient descent in function space, but not both.

Example: Mean Field when L = 1

Example: NTK

for some kernel K,

 $f_{t+1} - f_t = -K\mathcal{L}'(f_t, y)$ 

What is ruled out? e.g. higher order generalizations of NTK dynamics like  $f_{t+1} - f_t = -\langle K^{(2)}, \mathcal{L}'(f_t, y)^{\otimes 2} \rangle$ 

*Interesting consequence:* the NN function must be identically 0 at init in any feature learning limit

#### A Caricature of Space of abc-Parametrizations



 The nontrivial stable params form a high dimensional polytope

Will explain soon

- Feature learning params form 2 facets
- Everything else is in kernel regime

#### NTK, Standard Param Don't Learn Features

- Intuition why
  - The last layer weights get too much gradient, relative to weights in the body
  - We want to use larger learning rate to enable feature learning, but then the logits would blow up.



- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The *Tensor Programs* technique for deriving the limit

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The *Tensor Programs* technique for deriving the limit

#### Maximal Update Parametrization ( $\mu$ P)

- Modify Standard Param to get Maximal Update Param
  - Last layer: divide logits by  $\sqrt{n}$  and use  $\Theta(1)$  learning rate

• i.e. 
$$a_{L+1} = \frac{1}{2}$$
,  $c = 0$ 

- i.e.  $f(\xi) = \frac{1}{\sqrt{n}} w^{L+1} x^L(\xi)$  where  $w_{\alpha\beta}^{L+1} \sim \mathcal{N}\left(0, \frac{1}{n}\right)$
- This alone suffices to enable feature learning
- First layer: increase the gradient by n by setting  $a_1 = -\frac{1}{2}$ ,  $b_1 = 1/2$ 
  - i.e.  $h^1(\xi) = \sqrt{n}w^1\xi$  where  $w^1_{\alpha\beta} \sim \mathcal{N}\left(0, \frac{1}{n}\right)$
  - Needed to enable feature learning in *every* layer

#### Maximality of Maximal Update Param

• Maximal Update Param is *Maximally Feature Learning*: *Every* layer learns features  $(h^l(\xi), x^l(\xi))$  for every l will evolve during training).



#### 1-Dimension Degeneracy in abc

• For any  $\theta \in \mathbb{R}$  and  $t \ge 0$ ,  $f(\xi)$  at time t stays fixed for all input  $\xi$  if  $a_l \leftarrow a_l + \theta$ ,  $b_l \leftarrow b_l - \theta$ ,  $c \leftarrow c - 2\theta$ 



Space of abc-Parametrizations

#### Summary

Equivalent when L = 1 because, for  $\theta = 1/2$ ,  $a_l \leftarrow a_l + \theta, b_l \leftarrow b_l - \theta, c \leftarrow c - 2\theta$ 

 $h^{1}(\xi) = W^{1}\xi \in \mathbb{R}^{n}, \ x^{l}(\xi) = \phi\left(h^{l}(\xi)\right) \in \mathbb{R}^{n}, \ h^{l+1}(\xi) = W^{l+1}x^{l}(\xi) \in \mathbb{R}^{n}, \ f(\xi) = W^{L+1}x^{L}(\xi)$ 

	Definition		ΝΤΚ	Standard	b	Standard ( $1/n$ LR)	Mean Field	(L=1)	Maximal Update
a <sub>l</sub>	$W^{l} = n^{-a_{l}} w^{l}$	$\begin{cases} 0\\ \frac{1}{2} \end{cases}$	if $l = 1$ if $l > 1$	0		0	$\begin{cases} 0 & \text{if } l = \\ 1 & \text{if } l = \end{cases}$	= 1 = 2	$\begin{cases} -\frac{1}{2} & \text{if } l = 1\\ 0 & \text{if } 2 \le l \le L\\ \frac{1}{2} & \text{if } l = L + 1 \end{cases}$
b <sub>l</sub>	$w_{\alpha\beta}^{l} \sim \mathcal{N}(0, n^{-2b_{l}})$		0	$\begin{cases} 0 & \text{if } l = \\ \frac{1}{2} & \text{if } l > \end{cases}$	1 1	$\begin{cases} 0 & \text{if } l = 1 \\ \frac{1}{2} & \text{if } l > 1 \end{cases}$	0		1/2
С	$LR = \eta n^{-c}$		0	0		1	-1		0
Nontrivial									
	Stable								
Feature Learning				-					
Any depth									

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The *Tensor Programs* technique for deriving the limit

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

#### Key Theoretical Idea: Tensor Programs

When width is large, every activation vector has roughly iid coordinates, at **any** time during training. Using Tensor Programs, we can recursively calculate such coordinate distributions, and consequently understand how the neural network function evolves.

#### Key Theoretical Idea: Tensor Programs



- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All













Assume input and output dim = 1



For any *t*, there are coefficients  $A_t, B_t, C_t, D_t \in \mathbb{R}$ , which converge deterministically, such that

Assume input and output dim = 1



For any *t*, there are coefficients  $A_t, B_t, C_t, D_t \in \mathbb{R}$ , which converge deterministically, such that









#### Maximal Update Limit of Linear 1-Hidden-Layer NN with random init

 $Width-(d_{in} + d_{out})$ Linear 1-Hidden-Layer NN with "diagonal" init

> This is what we compute for largescale experiments like Word2Vec

 $\begin{array}{l} d_{in} = d_{out} = \text{vocab size} \\ d_{in} + d_{out} = 140k \text{ on text8} \\ d_{in} + d_{out} = 280k \text{ on fil9} \end{array}$ 

#### 1-Hidden-Layer: Summary

- The weight matrices have iid coordinates at initialization
- The function output converges due to Law of Large Numbers
- Gradients have approx. iid coordinates
- So after gradient update, weight coordinates are still approx. iid
- Repeat
- In the linear case, we express weights at any time as linear combinations of weights from initialization
  - This allows us to have efficient calculation of limit

#### *L*-Hidden-Layer: An Appetizer

- $n \times n$  Gaussian random matrix W in the middle of network
  - Central limit behavior
    - *Wx* is a Gaussian vector if *x* independent of *W*
  - Correlation with  $W^{\top}$ 
    - Appears after 1 step of SGD
      - no effect in 1<sup>st</sup> step due to Gradient Independence Phenomenon
- See paper for details of the *L*-hidden-layer limit

The Tensor Program Framework Automates All of These Derivations

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

#### What is a Tensor Program?

 A set of inductively generated vectors and scalars, given an initial set of matrices, vectors, and scalars

Initial



## What is a Tensor Program?

- A set of inductively generated vectors and scalars, given an initial set of matrices, vectors, and scalars
- Generation rules



MatMul









#### SGD as a Tensor Program

- 1-hidden-layer
  - Only need to use Nonlin and Moment
- *L*-hidden-layer
  - Need to use MatMul because of  $n \times n$  Gaussian matrix in the middle of network

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

#### Master Theorem

# Tells you how a Tensor Program behaves in the $n \rightarrow \infty$ limit.

#### Master Theorem

If 1) initial scalars converge deterministically and 2) initial matrices & vectors are sampled as iid Gaussians, then

- all vectors generated in the program have iid coordinates in the n → ∞ limit, and there are rules to calculate such limit distributions.
- all scalars generated in the program converge to deterministic values, and there are rules to calculate such limit scalars.

Embedding  $x^{l}(\xi)$  of input  $\xi$  of trained network

Output (i.e. logits) of trained network

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All



#### Church-Turing Thesis for Deep Learning

Any "useful" deep learning computation can be expressed as a Tensor Program.

#### Consequence: Infinite-Width Limits for All

- SOTA architectures: ResNet, Transformers, etc
- SGD with momentum, weight decay
- Adam, Adadelta, Adafactor, etc
- Natural Gradient Descent
- Pretraining and Finetuning
- Metalearning (e.g. MAML)
- Deep Reinforcement Learning (DQN, DDPG, etc)
- Image Generation (GANs, VAEs, etc)

# *Tensor Programs* "Compile" Finite-Width Computation to Infinite-Width Computation

To derive the infinite-width limit of **any** neural computation (e.g. SGD training),

- 1) express it as a Tensor Program, and
- 2) mechanically apply the Master Theorem.



- Parametrizations of Neural Networks
- Dichotomy of Parametrizations
- The "Maximal" Feature Learning Limit
- The Tensor Programs technique for deriving the limit
  - Key Idea
  - Motivating Example: Linear 1-Hidden-Layer NN
  - What is a Tensor Program?
  - Master Theorem
  - Infinite-Width Limit for All

## Summary

#### **Our Contributions**

- Classify all abc-parametrizations and their ∞-width limits
- Identify Maximal Update Parametrization for maximizing feature learning
- Propose the *Tensor Programs* technique for deriving its equations, and more generally, the limit of any neural computation

#### Significance

- Framework for studying feature learning in overparametrized NN
- Formulas for training featurelearning ∞-width NN in variety of settings (e.g. pretraining, metalearning, reinforcement learning, GANs, etc)
- Mostly solves the problem of taking ∞-width limits

#### Looking Ahead

- What kinds of representations are learned?
- How does it inform us about finite neural networks?
- How does this feature learning affect training and generalization?
- How does this jibe with the scaling law of language models?
- Can we train an infinite-width GPT...so GPT∞?
- ... and so many more questions are now ripe for the picking!

Thank you! Question?



Scan to link to paper

#### Max Learning Rates



Verifying Max Learning Rate for  $\mu P$  and SP

#### Same for Deep MLP

Let  $x(\xi)$  denote the last layer features (before output)

Time

index

#### **Neural Tangent Limit**

- At init, features  $x_0(\xi)$  has coordinates of size  $\Theta(1)$
- After t steps,  $x_t(\xi) x_0(\xi)$  has coordinates of size  $\Theta(1/\sqrt{n})$ .
- Feature kernel  $x_t(\xi)^{\top} x_t(\zeta)/n$  is fixed wrt t in  $n \to \infty$  limit
- Pretraining does not improve transfer learning

#### <mark>Maximal Update</mark> Limit

- At init, features  $x_0(\xi)$  has coordinates of size  $\Theta(1)$
- After t steps,  $x_t(\xi) x_0(\xi)$  has coordinates of size  $\Theta(1)$ .
- Feature kernel  $x_t(\xi)^T x_t(\zeta)/n$ changes with t in  $n \to \infty$  limit
- Pretraining improves transfer learning