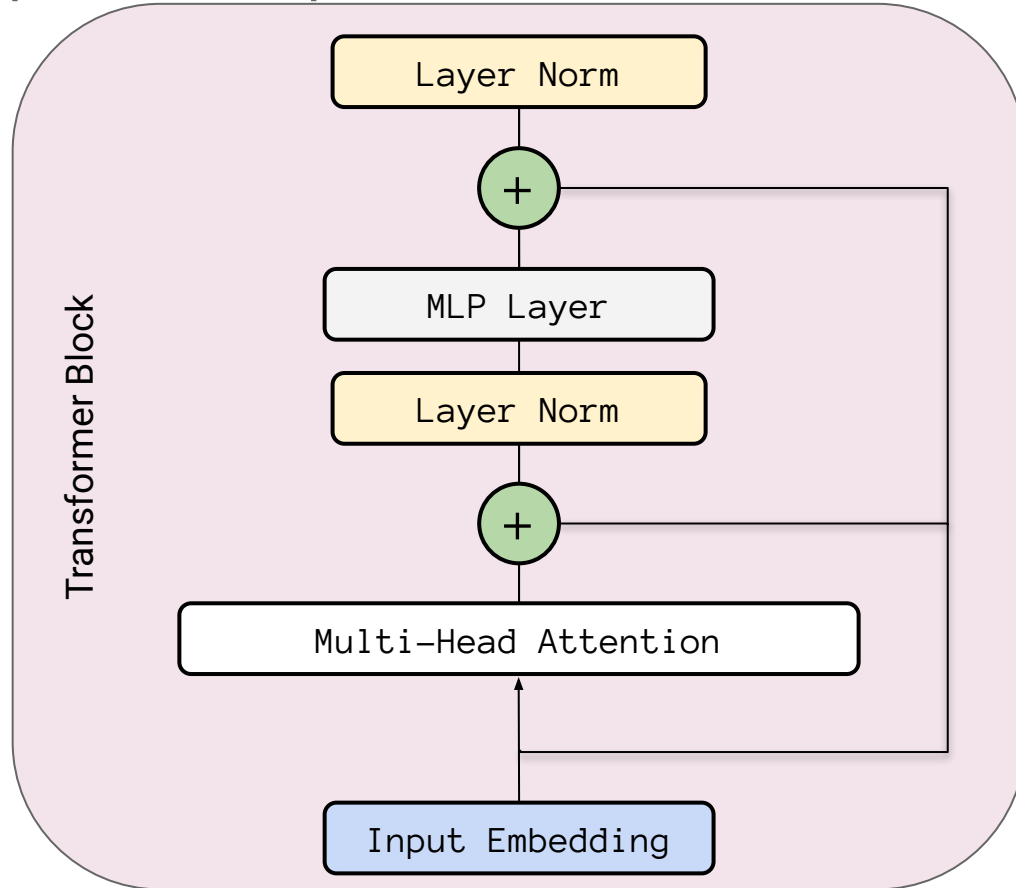# Training Large Language Models
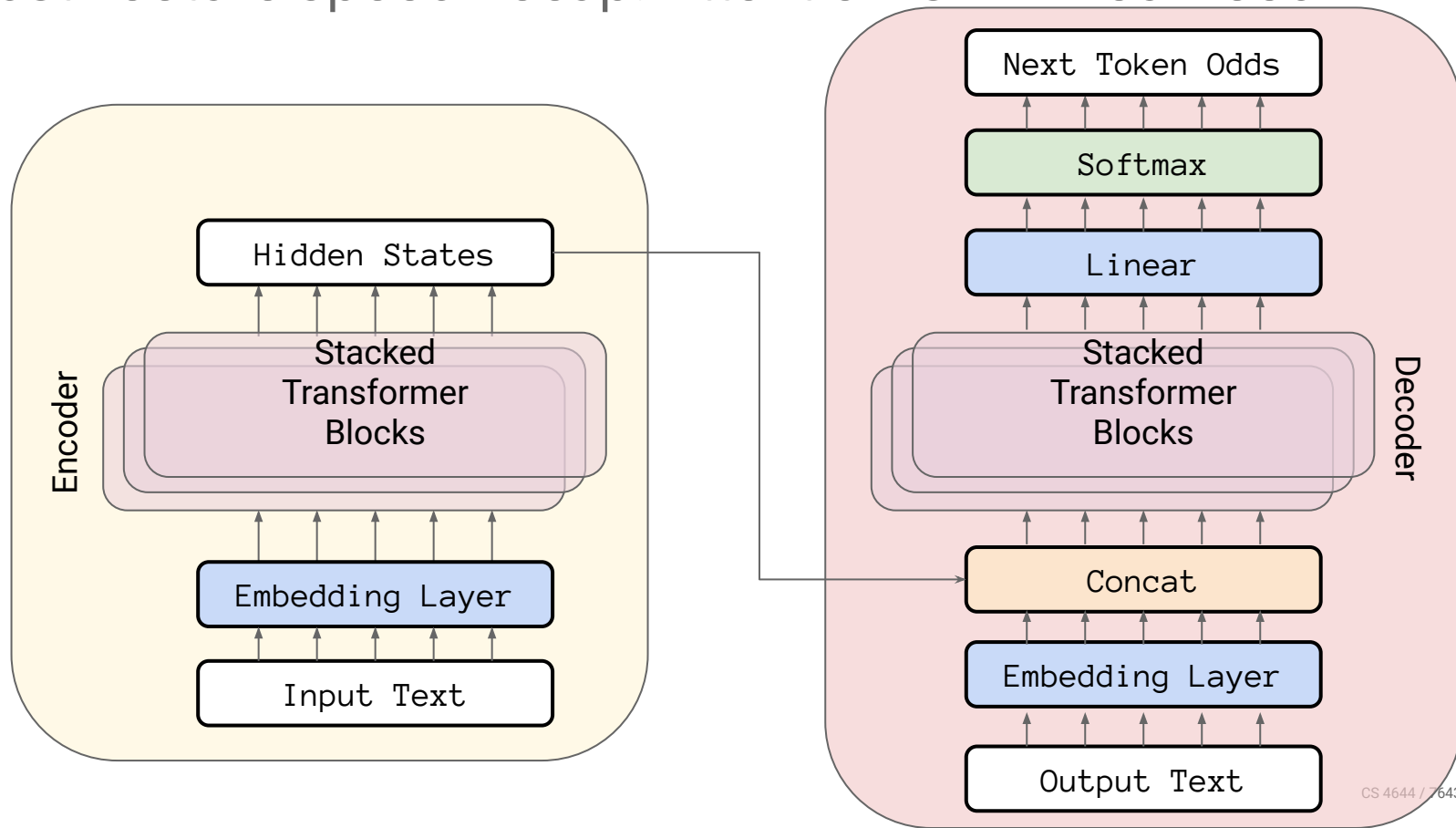
CS 4644 / 7643: Deep Learning

William Held
School of Interactive Computing
Georgia Institute of Technology
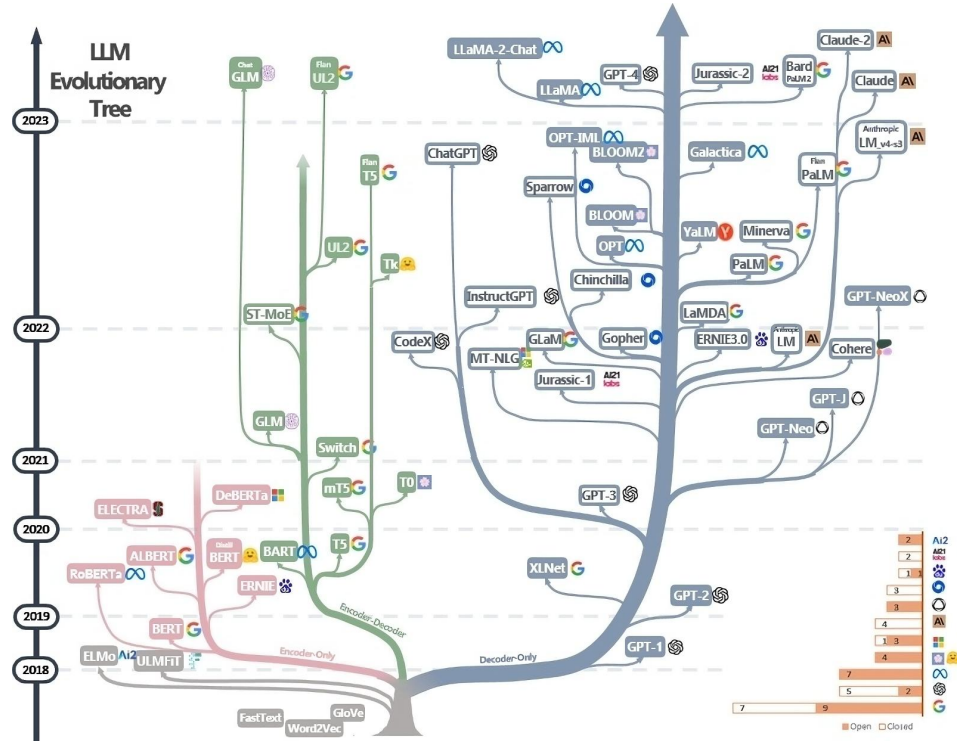
# Last Lecture Speed Recap: The Transformer Block

# Last Lecture Speed Recap: Attention is "All" You Need
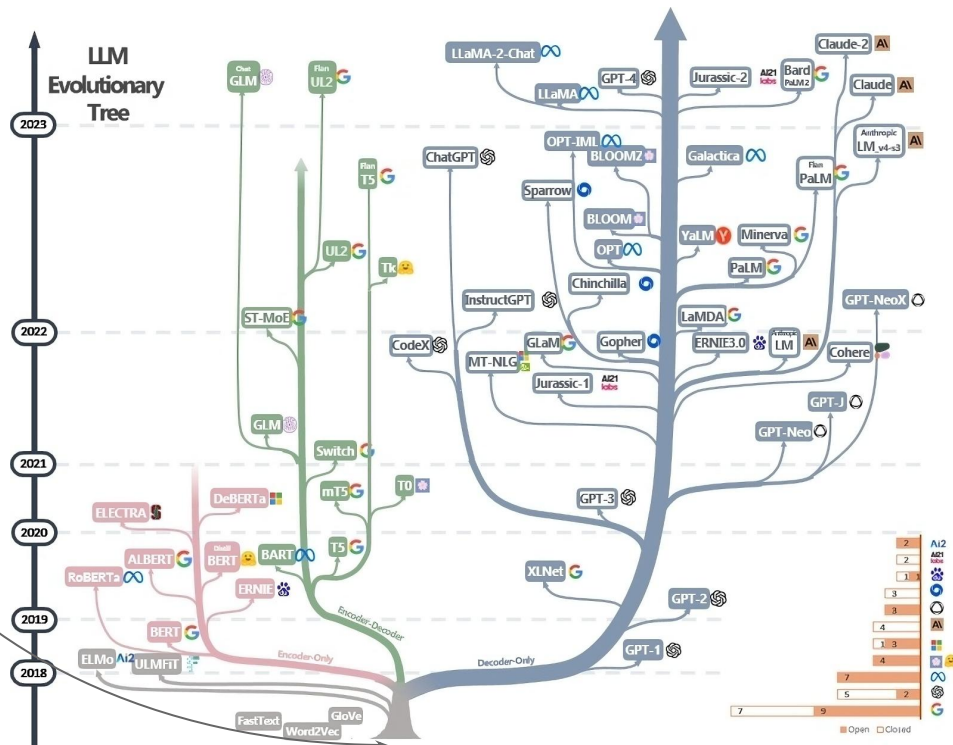
# How do we go from purpose driven models to LLMs?



LLM Evolutionary Tree

# How do we go from purpose driven models to LLMs?

**Self-Supervised Learning**
How do we most effectively turn
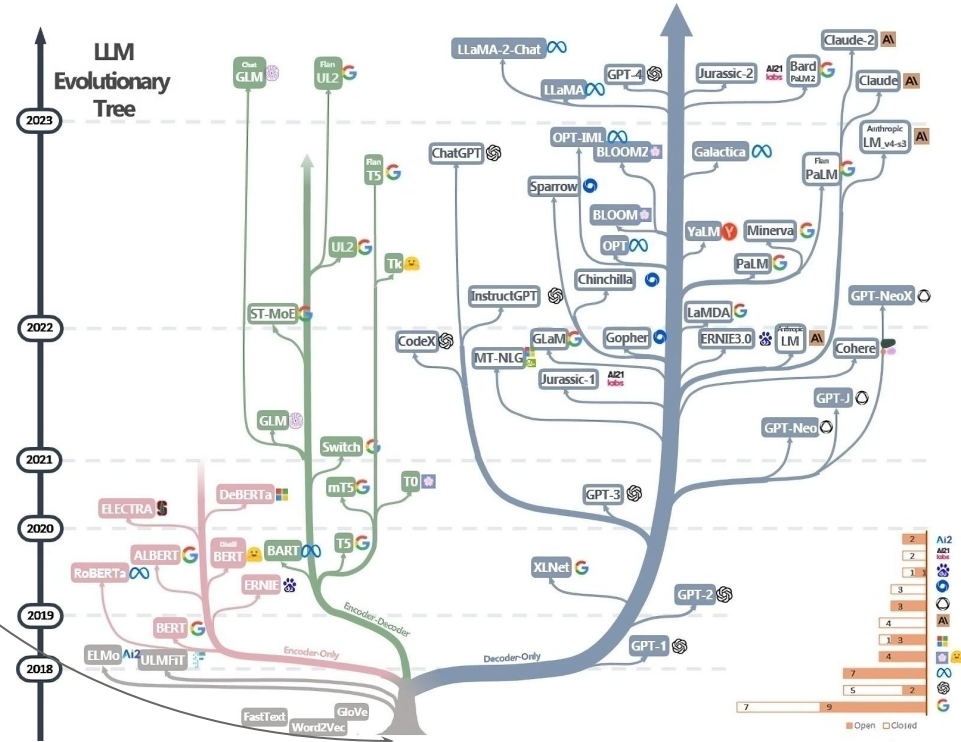raw text into meaningful loss?

# How do we go from purpose driven models to LLMs?
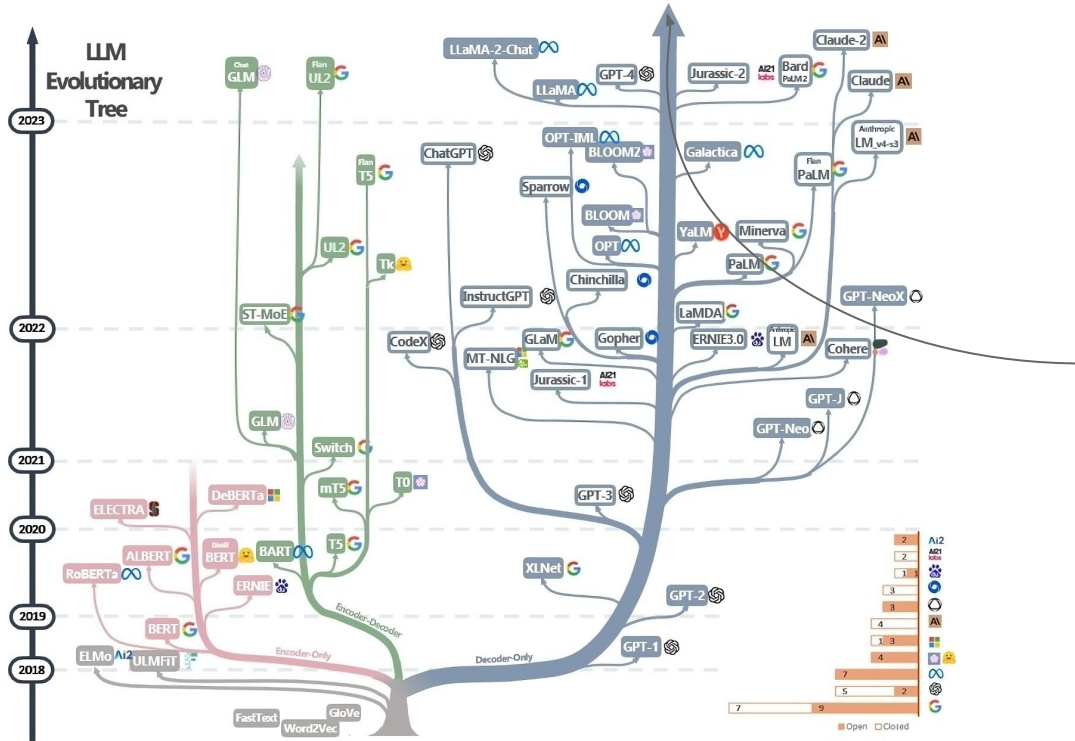


**Self-Supervised Learning**
How do we most effectively turn raw text into meaningful loss?
<u>Covered Today</u>
- Encoder Only
- Decoder Only
- Encoder-Decoder

CS 4644 / 7643 Deep Learning - William Held

# How do we go from purpose driven models to LLMs?



**Data Scaling**
How do we source and train on high-quality data at scale?
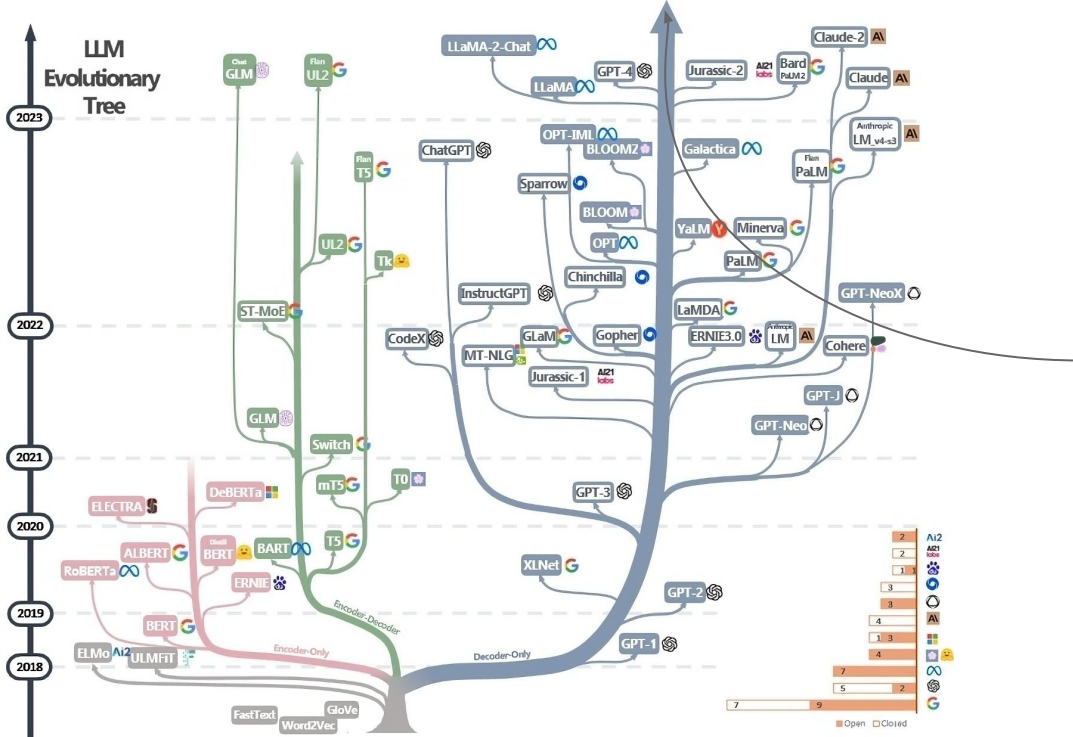
CS 4644 / 7643 Deep Learning - William Held

# How do we go from purpose driven models to LLMs?



**Data Scaling**
How do we source and train on high-quality data at scale?
<u>Covered Today</u>
- Data Curation Over Time
- Distributed Training

# LLM Advancements have been driven primarily by these two

**Self-Supervised Learning**
How do we most effectively turn
raw text into meaningful loss?

**Data Scaling**
How do we source and train on
high-quality data at scale?

# SSL | From raw text to loss!

Input         Masking

Hello         1

World         0

!             1

[PAD]         1

Masked Language Model

Devlin et al. 2018 (BERT)

# SSL | What is the "Mask" in a Masked Language Model?

Input          Masking

Hello          1

World          0

!              1

[PAD]          1

Masked Language Model

Devlin et al. 2018 (BERT)

# SSL | What is the "Mask" in a Masked Language Model?

Input          Masking

Hello            1

World            0

!                1

[PAD]            1

## Recall

```
Similarities: E = QXT / sqrt(DQ)
Attention Matrix: A = softmax(E,dim=1)
Output vectors: Y = AX
Y_i = ∑_j A_i,j X
```

$Y_i = \sum_j A_{i,j} X$

Masked Language Model

Devlin et al. 2018 (BERT)

# SSL | What is the "Mask" in a Masked Language Model?

Input     Masking

Hello     1

World     0

!         1

[PAD]     1

## Masked Attention

Similarities: E = (QXT / sqrt(DQ)) * MASK
Attention Matrix: A = softmax(E,dim=1)
Output vectors: Y = AX
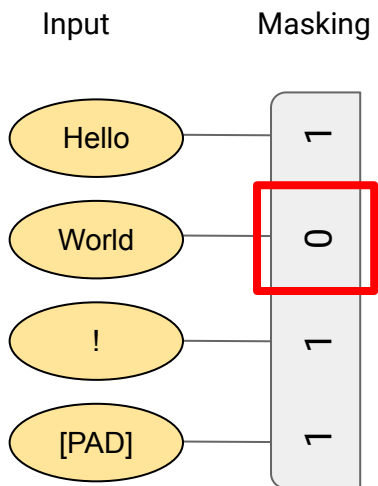$Y_i = \sum_j A_{i,j} X$

Masked Language Model

Devlin et al. 2018 (BERT)

# SSL | What is the "Mask" in a Masked Language Model?

Input        Masking

Hello        1

[MASK]       0

!            1

[PAD]        1

## Intuition

If $\text{MASK}_i = 0$, then $Y_i = \sum_{j, j \mathrel{!}= i} A_{i,j} X$

a.k.a the representation of the masked token is created purely from context

Masked Language Model

Devlin et al. 2018 (BERT)

# SSL | Masked Token Prediction

Input          Masking          Transformer

Hello ──── 1 ──── Encoder

[MASK] ──── 0 ────

! ──── 1 ────

[PAD] ──── 1 ────

Masked Language Model

Devlin et al. 2018 (BERT)

# SSL | Masked Token Prediction



Masked Language Model

Devlin et al. 2018 (BERT)

# SSL | Masked Token Prediction

Optimize Negative Log Likelihood
loss = -log(P("World" | Context))



P("World" | Context)

Softmax

Hidden States

Stacked
Transformer
Blocks

Encoder

Embedding Layer

Input Text

# **SSL** | Masked Token Prediction

Optimize Negative Log Likelihood
```
loss = -log(P("World" | Context)
```

Equivalent to the Cross-Entropy
Loss term from Lecture 3!

# Side Note | Tokens v.s. Words

Languages have a lot of words!

```
If V = Number of Words:
        O(V) Memory Scaling
```



```
P("World" | Context)
```

Encoder

Softmax

Hidden States

Stacked
Transformer
Blocks

Embedding Layer

Input Text

# Side Note | Tokens v.s. Words

Languages have a lot of words!

```
If V = Number of Words:
        O(V) Memory Scaling
        O(V) Runtime Scaling
```



P("World" | Context)

Softmax

Hidden States

Stacked Transformer Blocks

Encoder

Embedding Layer

Input Text

# Side Note | Tokens v.s. Words

Languages have a lot of words!

```
If V = Number of Words:
        O(V) Memory Scaling
        O(V) Runtime Scaling
```

This limits our vocabulary size a lot.

```
Tokenizers:
    Pre-processing to split words into
smaller chunks called "Tokens" so that we
    can cover all words with smaller V
```

P("World" | Context)

Encoder

Softmax

Hidden States

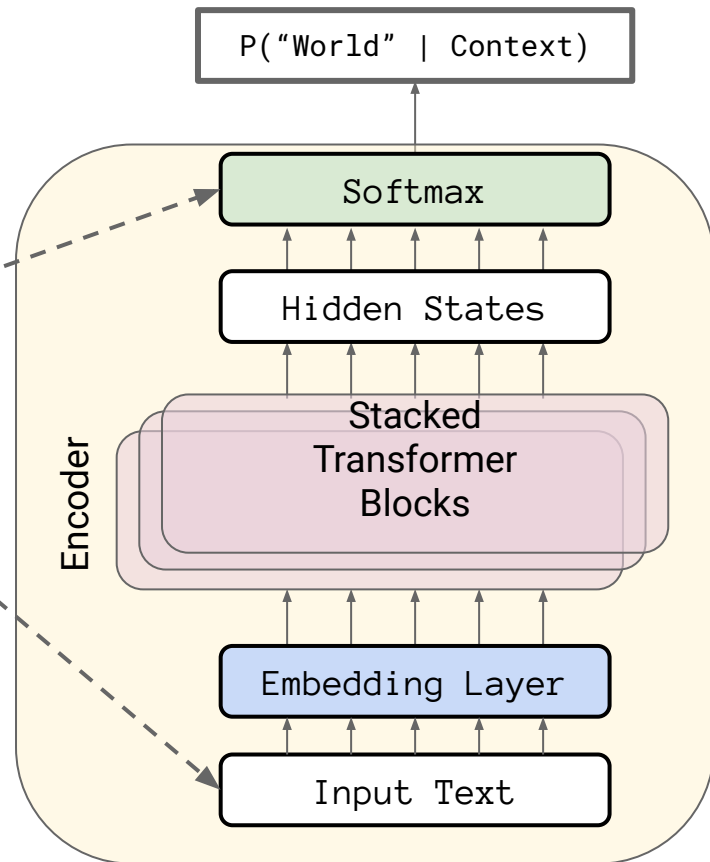Stacked Transformer Blocks

Embedding Layer

Input Text

# Side Note | Tokens v.s. Words

Languages have a lot of words!

```
If V = Number of Words:
        O(V) Memory Scaling
        O(V) Runtime Scaling
```
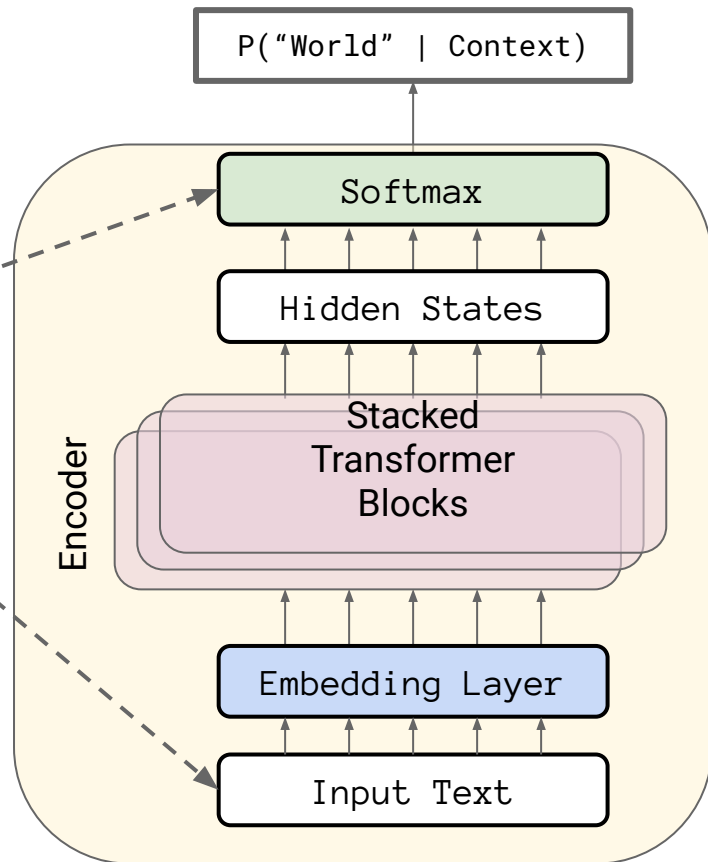
This limits our vocabulary size a lot.

```
Tokenizers:
    Pre-processing to split words into
smaller chunks called "Tokens" so that we
    can cover all words with smaller V

    Important but outside of Course Scope
        HuggingFace Tokenizer Summary
```



P("World" | Context)

Encoder

Softmax

Hidden States

Stacked Transformer Blocks

Embedding Layer

Input Text

# **Data** | BERT used existing curation!

**BERT Corpus**
English Wikipedia + BooksCorpus

**Size**
~3 Billion Tokens

**Quality**
High quality text,
Broad "Academic" Knowledge,
Limited Diversity

Devlin et al. 2018 (BERT)

# **Applications** | Encoders as "Foundation" Language Models

| Input | Masking | Transformer | Sequence Classification |
|---|---|---|---|

[CLS]

I

hate

you

Encoder

[CLS] Hidden State

MLP Layer

Softmax

Sentiment Label

Devlin et al. 2018 (BERT)

# **Applications** | Encoders as "Foundation" Language Models



Devlin et al. 2018 (BERT)

# **Applications** | Encoders as "Foundation" Language Models



Inputs

Transformer

Candidate Retrieval

Pretrained Encoder

[CLS] Hidden States

Sents.

Attentional State Neighborhood

True Coreference Label

Intel
or
2009 Earthquake

Nvidia
or
2013 Earthquake

AMD
or
2010 Earthquake

Held et al. 2021

# Questions?

Input      Masking      Transformer      Mask Prediction

Hello    1

[MASK]    0

!    1

[PAD]    1

Encoder

World

Masked Language Model

# SSL | "How does GPT work?"



Input      Masking      Transformer      Next Token Prediction

| Input | Causal Mask | Decoder | Next Token Prediction |
|-------|-------------|---------|-----------------------|
| Hello | | | World |
| World | | | ! |
| ! | | | [EOS] |
| [PAD] | | | |

Radford et al. 2019 (GPT-2)

# SSL | Autoregressive Language Modeling

Masking

Causal Mask

| Hello | 1 | 0 | 0 | 0 |
| World | | | | |
| ! | | | | |
| [PAD] | | | | |

Radford et al. 2019 (GPT-2)

# SSL | Autoregressive Language Modeling

Masking

Causal Mask

| Hello | | | |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| World | | | |
| 1 | 1 | 0 | 0 |
| ! | | | |
| [PAD] | | | |

Radford et al. 2019 (GPT-2)

# SSL | Autoregressive Language Modeling

Masking

Hello

World

!

[PAD]

### Causal Mask

| | | | |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

Radford et al. 2019 (GPT-2)

# SSL | Autoregressive Language Modeling

Masking

Causal Mask



| | | | |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 |

Hello
World
!
[PAD]

## Masked Attention Again!

```
Similarities: E = (QXT / sqrt(DQ)) * MASK

Attention Matrix: A = softmax(E,dim=1)

Output vectors: Y = AX
```
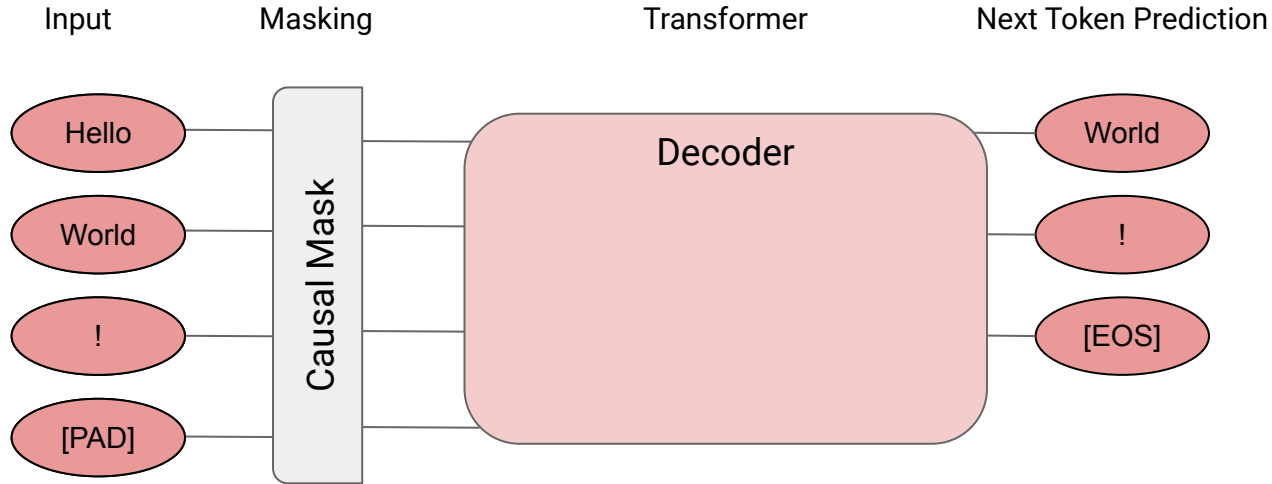
$Y_i = \sum_j A_{i,j} X$

Tokens only affected by preceding tokens

Radford et al. 2019 (GPT-2)

# SSL | First successful GPT Model, Purely Autoregressive



Input     Masking     Transformer     Next Token Prediction

Hello

World

!

[PAD]

Causal Mask

Decoder

World

!

[EOS]

Radford et al. 2019 (GPT-2)

# SSL | First successful GPT Model, Purely Autoregressive

| Input | Masking | Transformer | Next Token Prediction |
|-------|---------|-------------|----------------------|



## Optimize Negative Log Likelihood of Whole Sequence

```
loss = -(log(P("World" | "Hello") + log(P("!" | "Hello World") +
       log(P("[EOS]" | "Hello World!")))
```

Radford et al. 2019 (GPT-2)

# **Data** | Increasing Token Count via Human Curation Heuristics

## GPT-2 Corpus
All Reddit Outbound links with at least 3 karma

## Size
~10 Billion Tokens

## Quality
High quality text,
Broad Knowledge,
Improved Diversity

| URL Domain | # Docs | % of Total Docs |
|---|---|---|
| bbc.co.uk | 116K | 1.50% |
| theguardian.com | 115K | 1.50% |
| washingtonpost.com | 89K | 1.20% |
| nytimes.com | 88K | 1.10% |
| reuters.com | 79K | 1.10% |
| huffingtonpost.com | 72K | 0.96% |
| cnn.com | 70K | 0.93% |
| cbc.ca | 67K | 0.89% |
| dailymail.co.uk | 58K | 0.77% |
| go.com | 48K | 0.63% |

Radford et al. 2019 (GPT-2)

Ok, but what should I use?

# SSL | Classification Comparison

| Model | MNLI | CoLA | SST-2 | MRPC | STS-B | QQP | QNLI | RTE | Avg |
|-------|------|------|-------|------|-------|-----|------|-----|-----|
| GPT-2-original | 85.9/85.6 | 54.8 | 94.5 | 86.9/82.2 | 86.3/85.2 | 72.5/89.3 | 91.2 | 69.8 | 80.9 |
| GPT-2-finetuned | 85.8/85.5 | 40.9 | 94.5 | 87.0/81.0 | 85.6/84.3 | 71.4/88.5 | 91.5 | 69.0 | 78.8 |
| RoBERTa-large | 90.1/89.7 | 63.8 | 96.1 | 91.2/88.3 | 90.9/90.7 | 72.5/89.6 | 94.5 | 85.9 | 86.5 |

He et al. 2021

# SSL | Pretrained Retrieval Comparison



GPT-2 separates into two clusters

Left | PCA | Right                    UMAP

manifold with other positional IDs →

contains only position ID 0 →

Bert consists of multiple small clusters

PCA                                     UMAP

absolute position ID

260
240
220
200
180
160
140
120
100
80
60
40
20
0

each cloud consists of unsorted posIDs →

Encoders can't generate!

# SSL | Encoder-Only vs. Decoder-Only

## Encoder
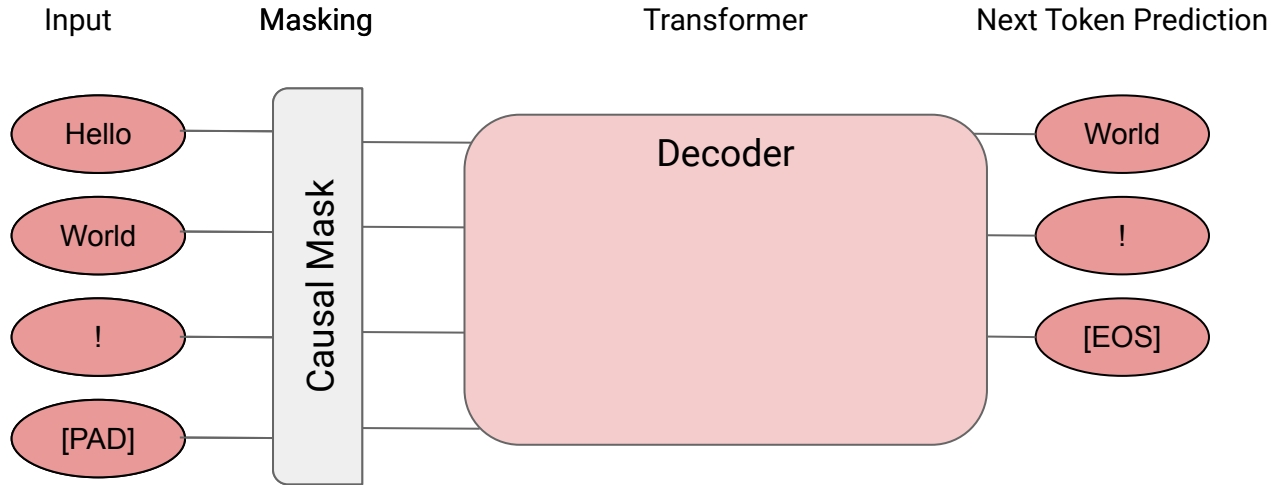
+ Retrieval
+ Classification
- No Generative Abilities

## Decoder

+ Generative Abilities
- Retrieval
- Classification

Wang et al. 2022

# SSL | Encoder-Only vs. Decoder-Only

**Encoder**

+   Retrieval
+   Classification
-   No Generative Abilities

**Decoder**

+   Generative Abilities ———— This is pretty essential!
-   Retrieval
-   Classification

Wang et al. 2022

# Questions?

Input        Masking        Transformer        Next Token Prediction



Autoregressive Language Model

# SSL | Encoder-Only vs. Decoder-Only

**Encoder**

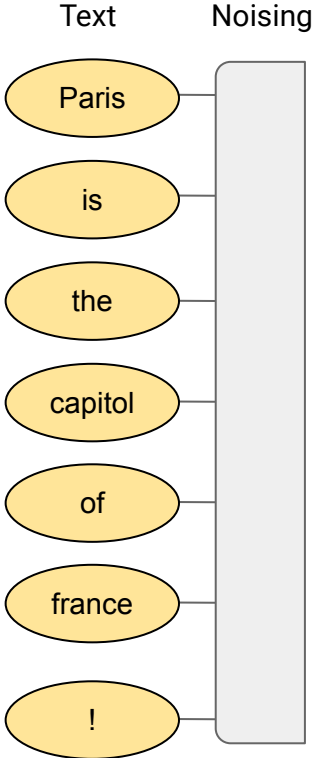How to keep this? — + Retrieval
+ Classification
- No Generative Abilities

**Decoder**

+ Generative Abilities — This is pretty essential!
- Retrieval
- Classification

# SSL | Encoder-Decoder Returns

# SSL | Universal Text-to-Text

Text          Noising

Paris

is

the

capitol

of

france

!

Raffel et al. 2019

# SSL | Universal Text-to-Text

Text    Noising

Paris

is

the

capitol

of

france

!

Original text

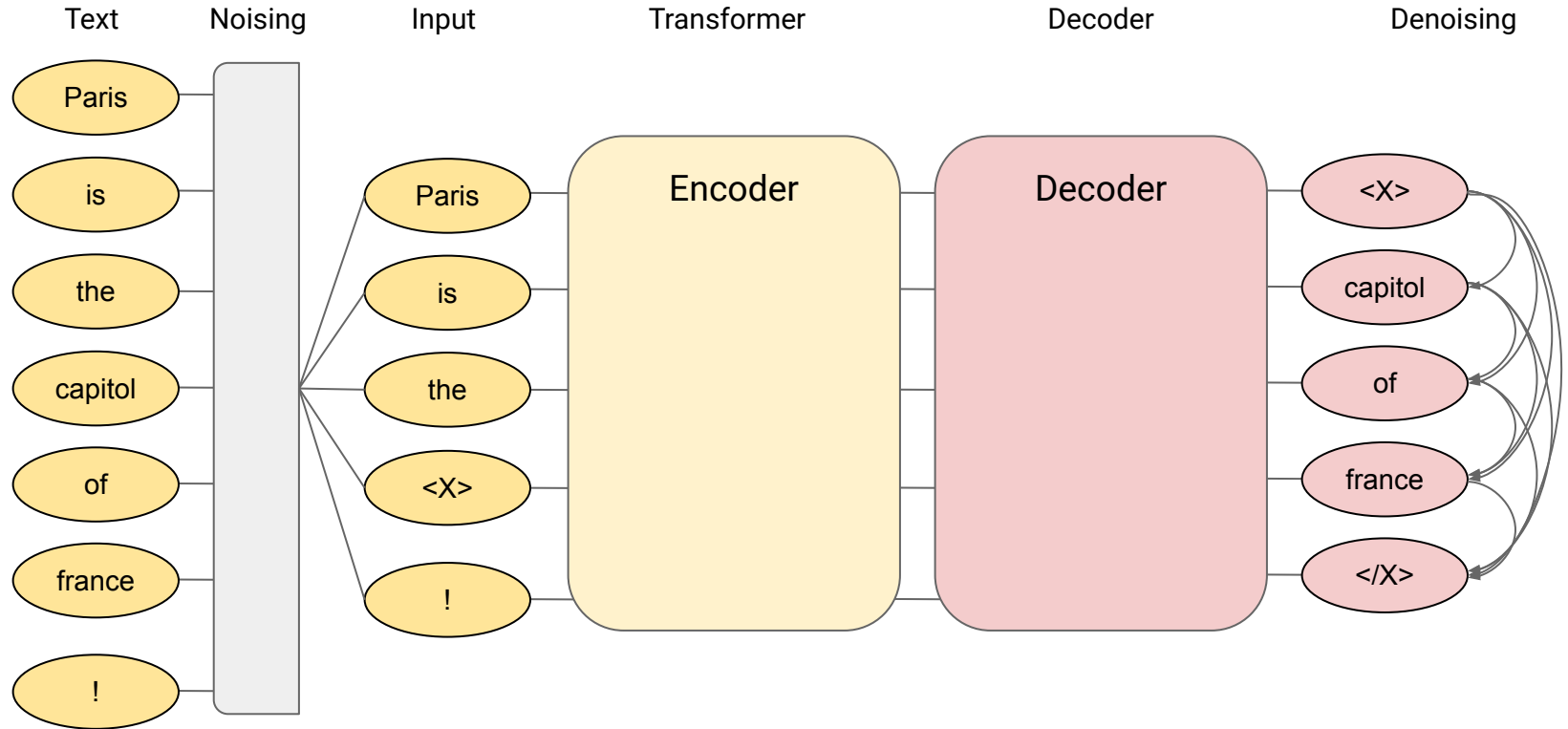Thank you ~~for inviting~~ me to your party ~~last~~ week.

Inputs

Thank you <X> me to your party <Y> week.

Targets

<X> for inviting <Y> last <Z>

Raffel et al. 2019

# SSL | Universal Text-to-Text



Text   Noising   Input   Transformer   Decoder   Denoising

Paris — is — the — capitol — of — france — !

Paris — is — the — <X> — !

Encoder

Decoder

<X> — capitol — of — france — </X>

Raffel et al. 2019

CS 4644 / 7643 Deep Learning - William Held

# SSL | Universal Text-to-Text

Text       Noising

Paris
is
the
capitol
of
france
!

A_C._E.
Token Masking

D E . A B C .
Sentence Permutation

C . D E . A B
Document Rotation

A . C . E .
Token Deletion

A B C . D E .

A _ . D _ E .
Text Infilling

Lewis et al. 2020

# SSL | UL2 - Text-to-Text Pushed to Limits

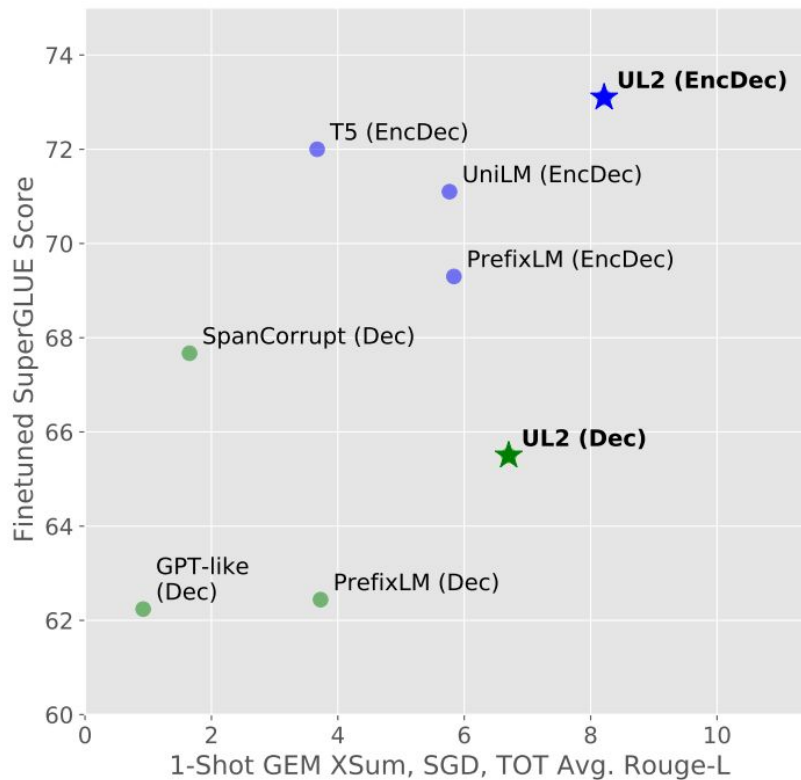# SSL | Universal Text-to-Text

Regardless of noise, Loss Function remains the same still!

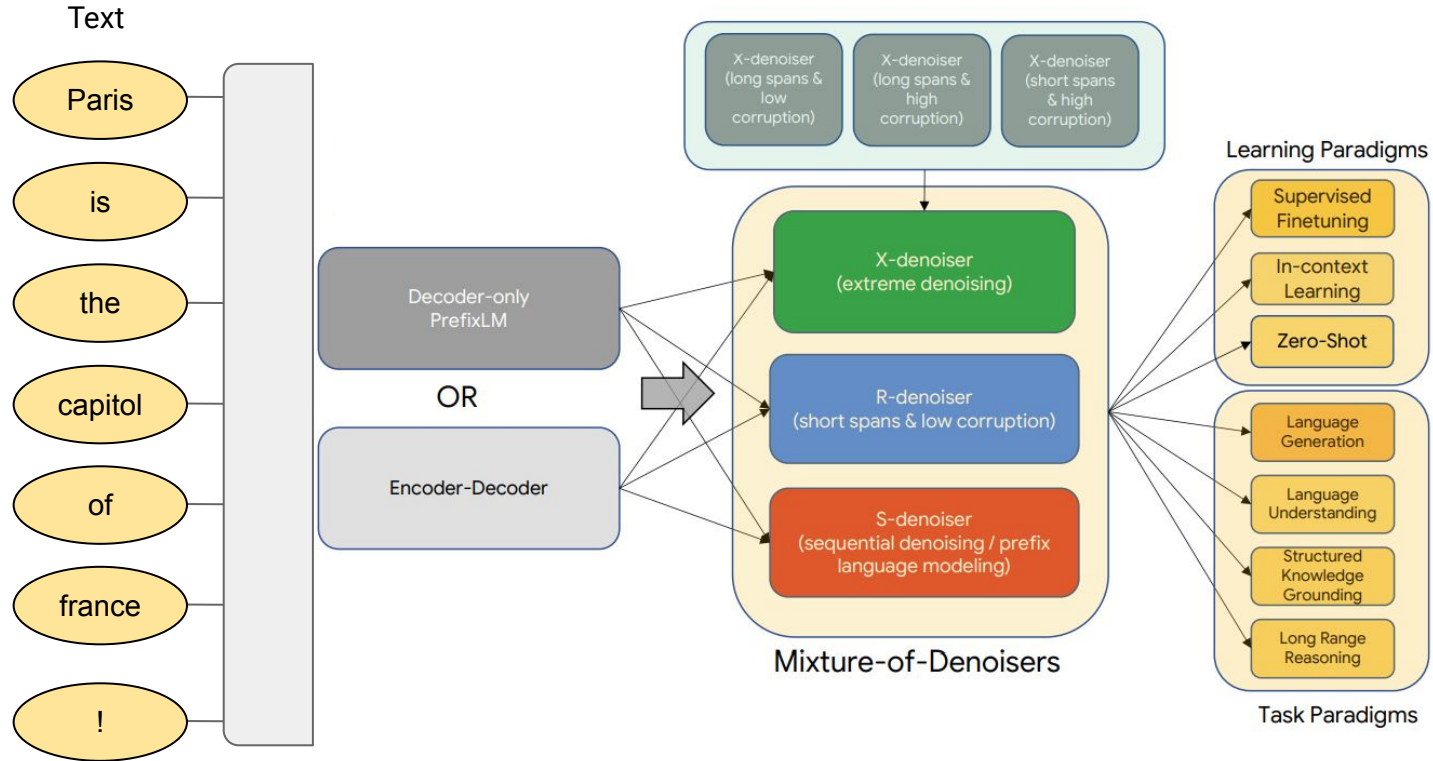Continue using Negative Log Likelihood

```
loss = -(log(P(Denoised Sequence | Noised Sequence))
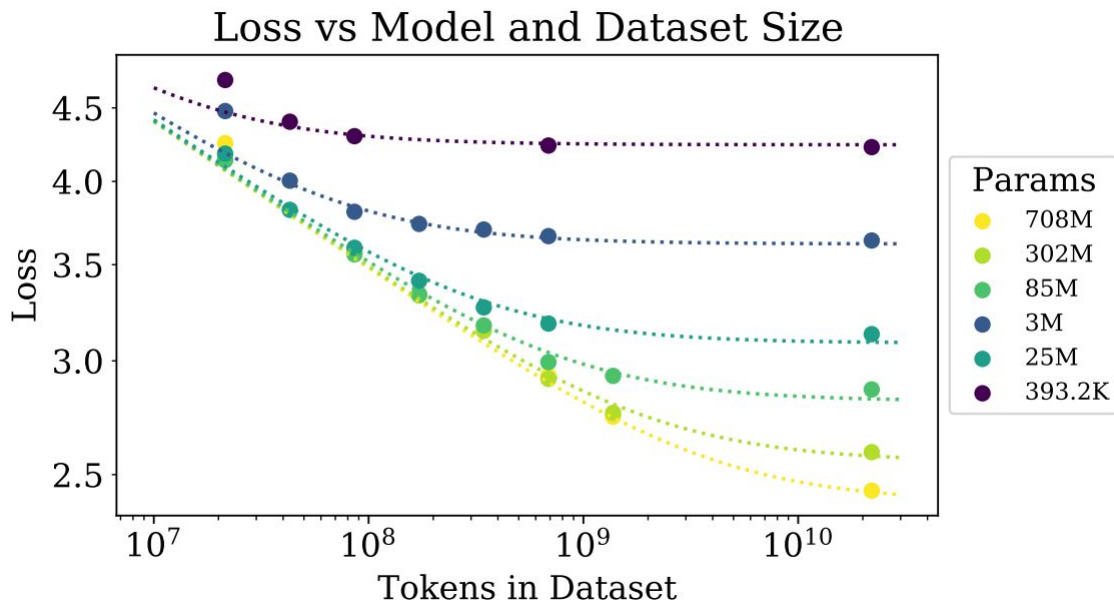```

# SSL | Universal Text-to-Text Is Architecture Agnostic



Tay et al. 2023

# Questions?



Text

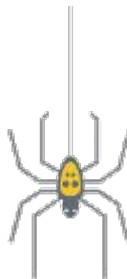Paris is the capitol of france !

Decoder-only PrefixLM

OR

Encoder-Decoder

X-denoiser (long spans & low corruption)   X-denoiser (long spans & high corruption)   X-denoiser (short spans & high corruption)

X-denoiser (extreme denoising)

R-denoiser (short spans & low corruption)

S-denoiser (sequential denoising / prefix language modeling)

Mixture-of-Denoisers

Learning Paradigms
- Supervised Finetuning
- In-context Learning
- Zero-Shot

- Language Generation
- Language Understanding
- Structured Knowledge Grounding
- Long Range Reasoning

Task Paradigms

# **Data** | Moving to truly Large Language Models

Today's LLMs are driven data and model scaling
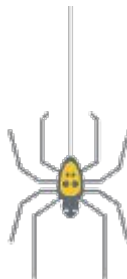
Loss vs Model and Dataset Size



Kaplan et al. 2020
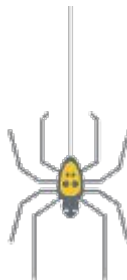
# Data | Moving to truly Large Language Models



We could get a lot more data from CommonCrawl!

# **Data** | Moving to truly Large Language Models



We could get a lot more data from CommonCrawl!
A lot of it is spam though…

# **Data** | Moving to truly Large Language Models

We could get a lot more data from CommonCrawl!
A lot of it is spam though…
How do we get "useful" data?

# Data | C4 - First Scaling of Data Via Common Crawl

**T5 Corpus (AKA C4)**
All Common Crawl Text Which Meets Heuristics
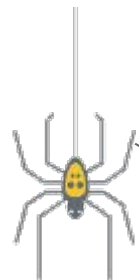
**Size**
~350 Billion Tokens

**Quality**
Varying quality text,
Broad Knowledge,
Improved Diversity

- We only retained lines that ended in a terminal punctuation mark (i.e. a period, exclamation mark, question mark, or end quotation mark).

- We discarded any page with fewer than 3 sentences and only retained lines that contained at least 5 words.

- We removed any page that contained any word on the "List of Dirty, Naughty, Obscene or Otherwise Bad Words".[6]

- Many of the scraped pages contained warnings stating that Javascript should be enabled so we removed any line with the word Javascript.

- Some pages had placeholder "lorem ipsum" text; we removed any page where the phrase "lorem ipsum" appeared.

- Some pages inadvertently contained code. Since the curly bracket "{" appears in many programming languages (such as Javascript, widely used on the web) but not in natural text, we removed any pages that contained a curly bracket.

- Since some of the scraped pages were sourced from Wikipedia and had citation markers (e.g. [1], [citation needed], etc.), we removed any such markers.

- Many pages had boilerplate policy notices, so we removed any lines containing the strings "terms of use", "privacy policy", "cookie policy", "uses cookies", "use of cookies", or "use cookies".

- To deduplicate the data set, we discarded all but one of any three-sentence span occurring more than once in the data set.

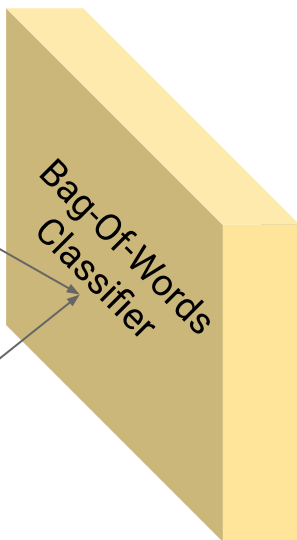Raffel et al. 2019

# Data | GPT-3 - Increased Scaling Via Curation

Training

Low-Quality, High Volume

Bag-Of-Words Classifier

Distinguish High and Low Quality

| URL Domain | # Docs | % of Total Docs |
|---|---|---|
| bbc.co.uk | 116K | 1.50% |
| theguardian.com | 115K | 1.50% |
| washingtonpost.com | 89K | 1.20% |
| nytimes.com | 88K | 1.10% |
| reuters.com | 79K | 1.10% |
| huffingtonpost.com | 72K | 0.96% |
| cnn.com | 70K | 0.93% |
| cbc.ca | 67K | 0.89% |
| dailymail.co.uk | 58K | 0.77% |
| go.com | 48K | 0.63% |

High Quality, Medium Volume

Brown et al. 2020

# **Data** | GPT-3 - Increased Scaling Via Curation

Filtering

Bag-Of-Words
Classifier

Keep "False" Positives

"False" positive ~= High Quality

Brown et al. 2020

# Data | GPT-2 to Original GPT-3 was mostly data scaling

**GPT-3 Corpus**
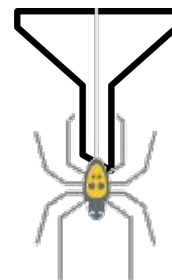Common-Crawl Filtered using
GPT-2 Training Data

**Size**
~400 Billion Tokens

**Quality**
High-ish quality text,
Broad Knowledge,
Web-scale Diversity

Brown et al. 2020

# **Data** | Recent Open Source models focus heavily on data scaling

**Llama 1 Corpus**

**Size**
~1.4 Trillion Tokens

**Quality**
Varying quality text,
Broad Knowledge,
Web-scale Diversity,
Includes Code!

| Dataset | Sampling prop. | Epochs | Disk size |
|---|---|---|---|
| CommonCrawl | 67.0% | 1.10 | 3.3 TB |
| C4 | 15.0% | 1.06 | 783 GB |
| Github | 4.5% | 0.64 | 328 GB |
| Wikipedia | 4.5% | 2.45 | 83 GB |
| Books | 4.5% | 2.23 | 85 GB |
| ArXiv | 2.5% | 1.06 | 92 GB |
| StackExchange | 2.0% | 1.03 | 78 GB |

Touvron et al. 2023

# Data | Recent Open Source models focus heavily on data scaling

**Falcon Refined Web Corpus**

**Size**
5 Trillion Tokens

**Quality**
Varying quality text,
Broad Knowledge,
Web-scale Diversity,
Includes Code

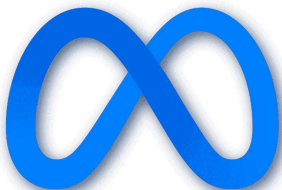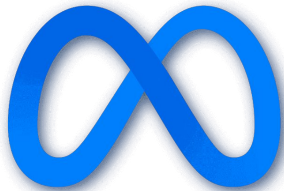# **Data** | Data Mixture has become the biggest "secret"

### Llama 2 Corpus

**Size**
> 2 Trillion Tokens

**Quality**
Minimal details known

Touvron et al. 2023 (b)

### PALM-2 Corpus

**Size**
> 3.6 Trillion Tokens

**Quality**
No details known

Anil et al. 202320

### GPT-4 Corpus

**Size**
Unknown (Est. 11T Tokens)
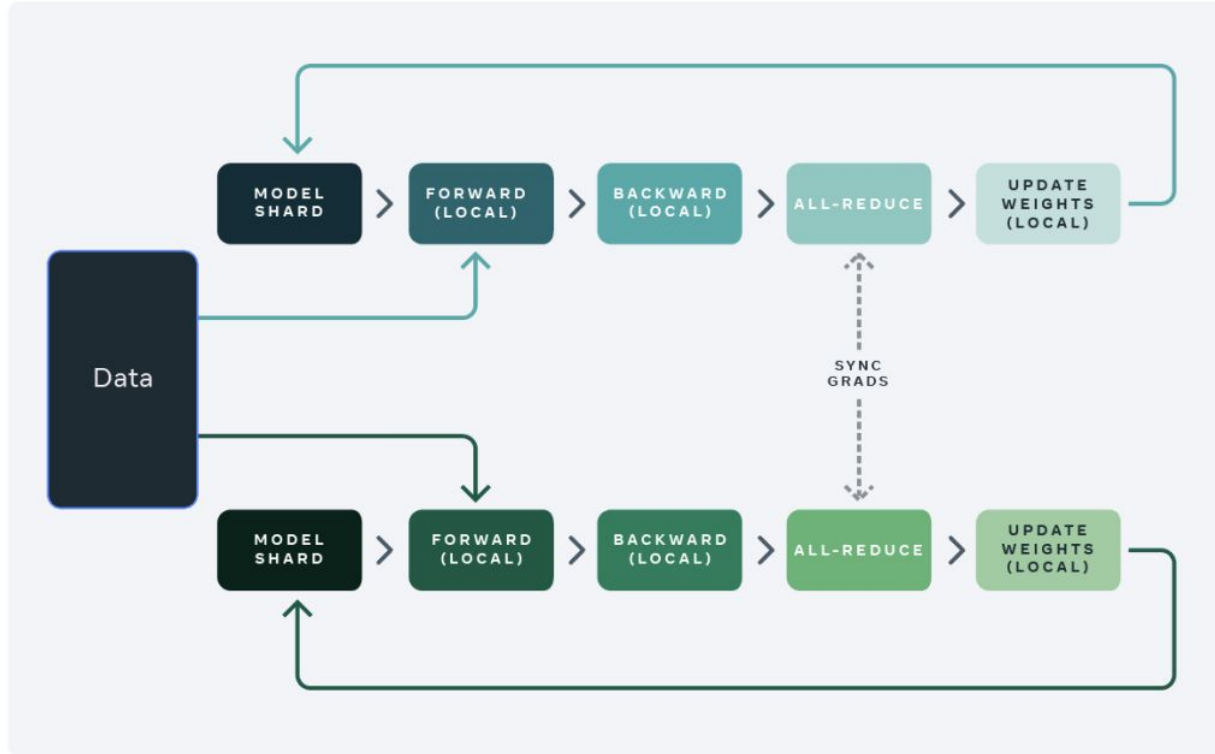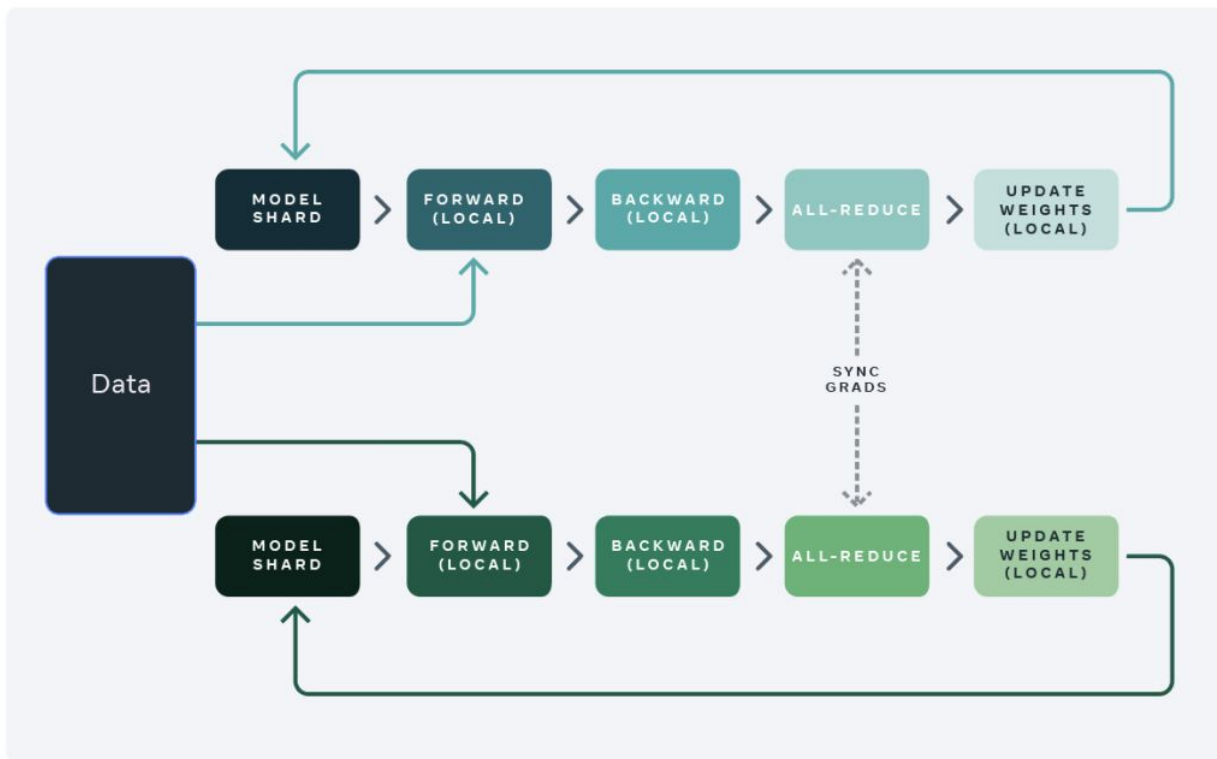
**Quality**
No details known

OpenAI 2023

# Questions?

| **Llama 2 Corpus** | **PALM-2 Corpus** | **GPT-4 Corpus** |
|:---:|:---:|:---:|
| **Size** | **Size** | **Size** |
| > 2 Trillion Tokens | > 3.6 Trillion Tokens | Unknown (Est. 11T Tokens) |
| **Quality** | **Quality** | **Quality** |
| Minimal details known | No details known | No details known |

Touvron et al. 2023 (b)  Anil et al. 202320  OpenAI 2023

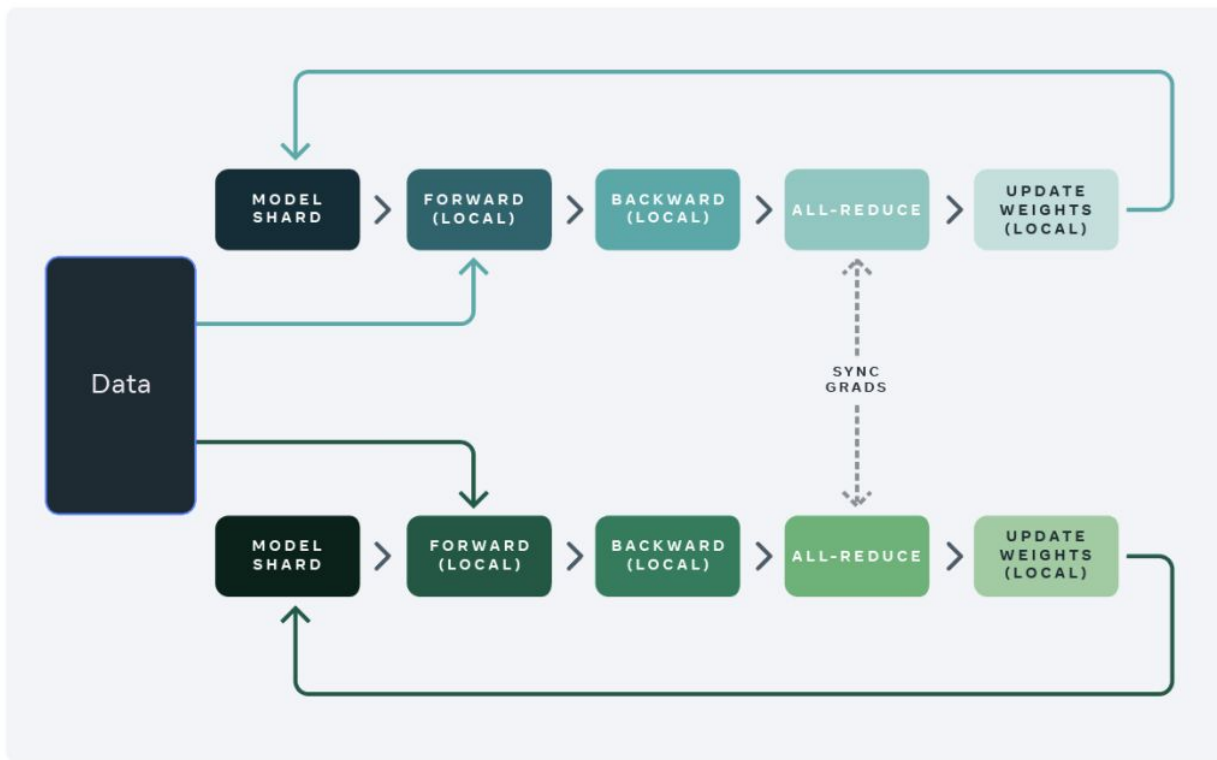# Scaling Parameters | Data Parallel Training

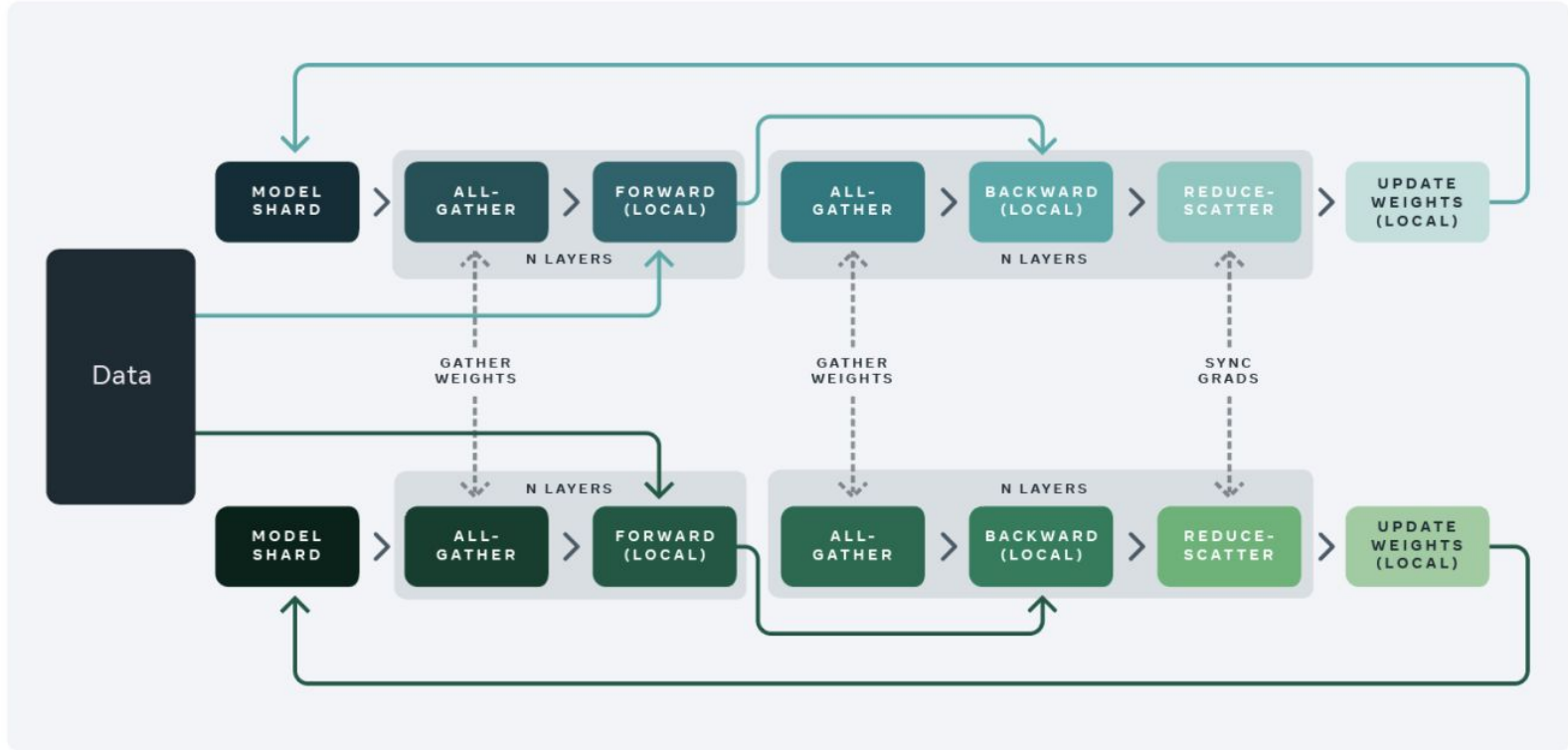# Scaling Parameters | Data Parallel Training



Total memory increases linearly with shards

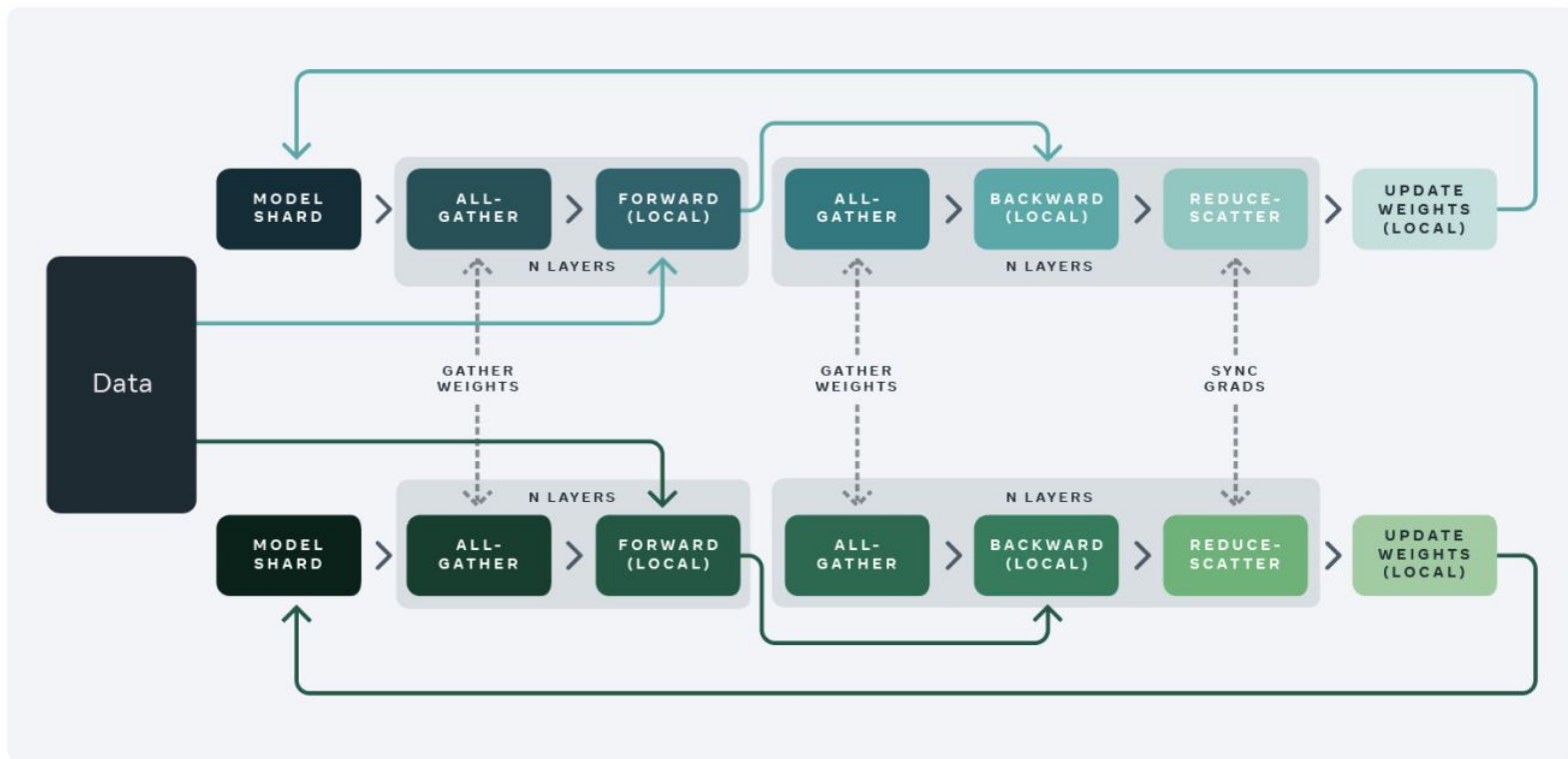# Scaling Parameters | Data Parallel Training



Max memory constrains model size

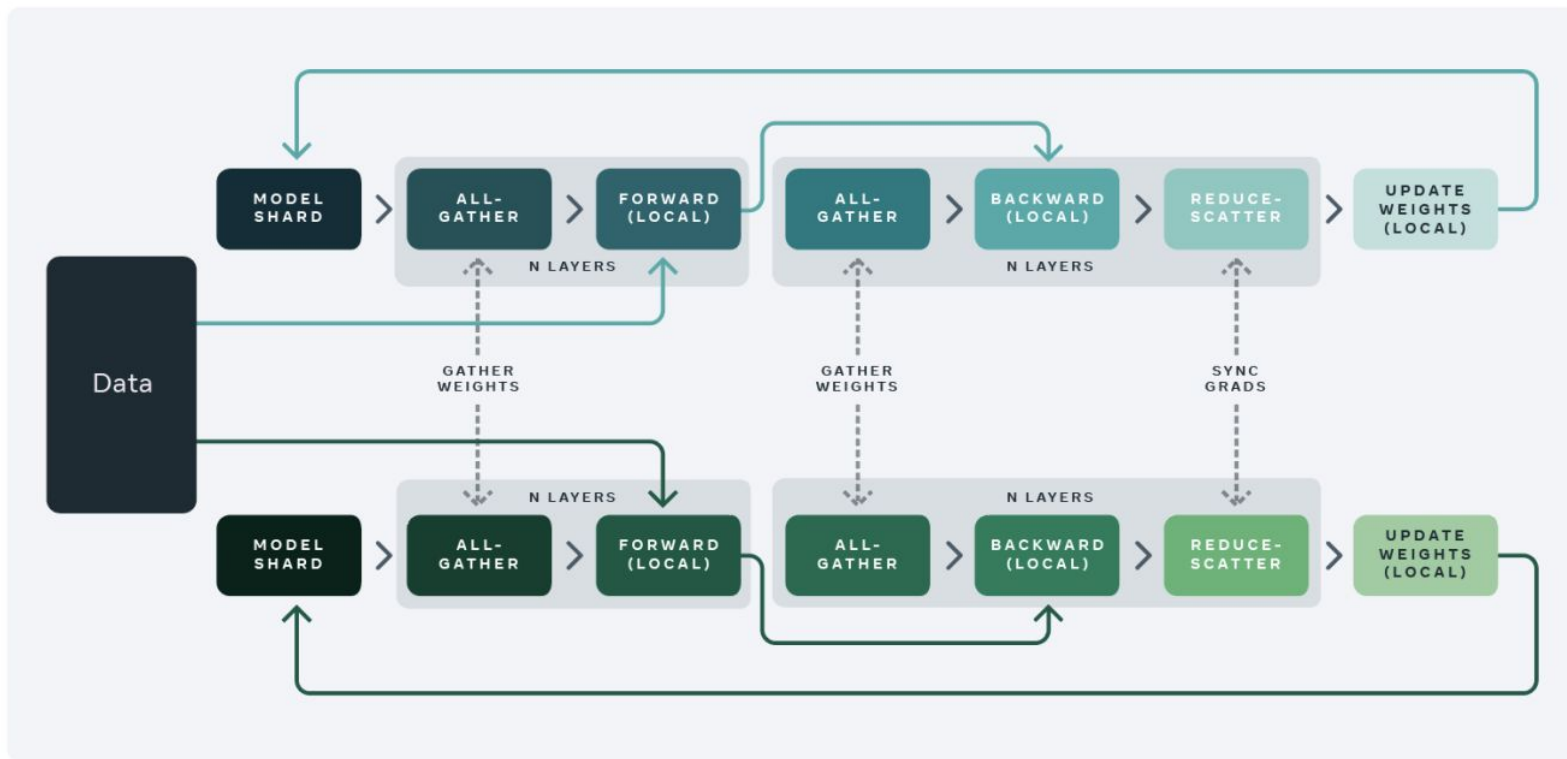# Scaling Parameters | *Fully* Sharded Data Parallel Training

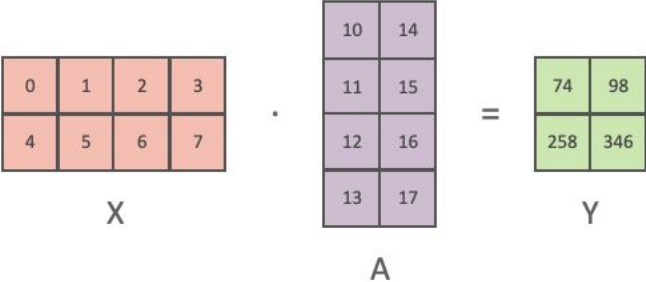# Scaling Parameters | *Fully* Sharded Data Parallel Training



Total memory is constant

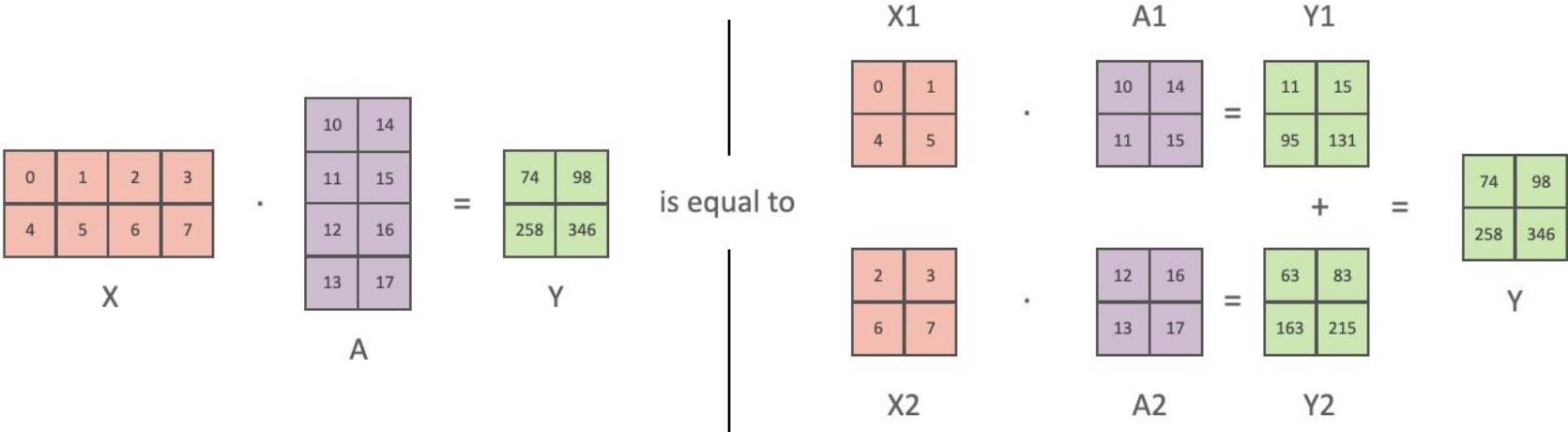# Scaling Parameters | *Fully* Sharded Data Parallel Training



Max single GPU memory constrains layer size

# Scaling Parameters | Tensor Parallel Training

https://huggingface.co/docs/transformers/v4.15.0/parallelism#tensor-parallelism

# Scaling Parameters | Tensor Parallel Training

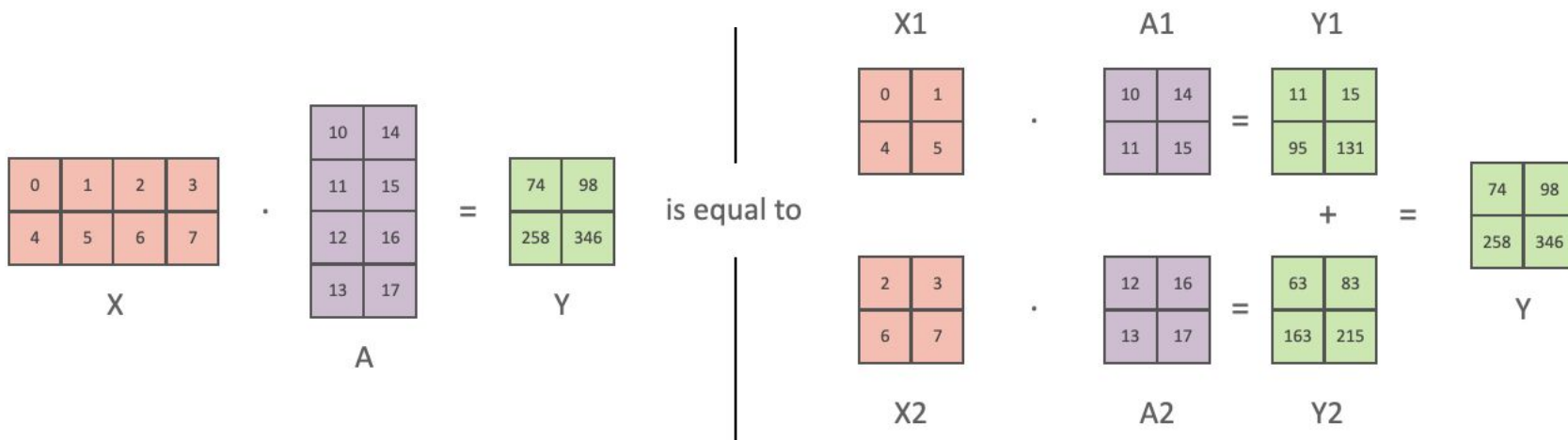CS 4644 / 7643 Deep Learning - William Held
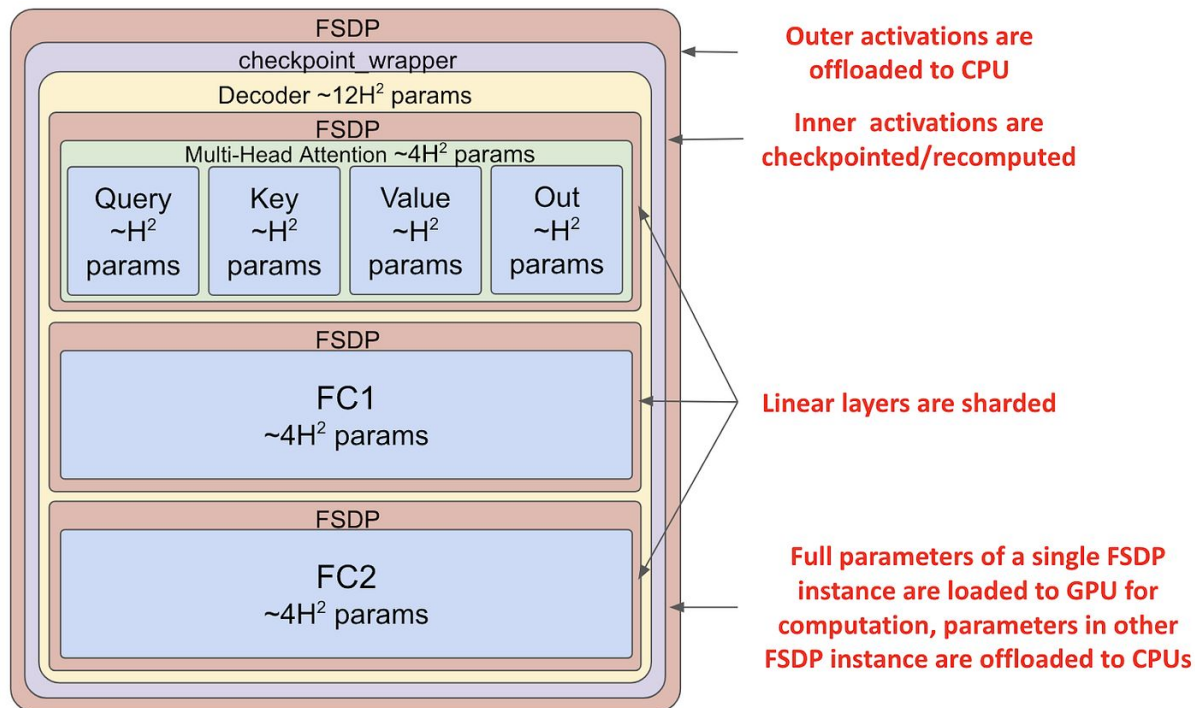
# Scaling Parameters | Tensor Parallel Training



Don't need to sync gradients!

# Scaling Parameters | Tensor Parallel Training



Don't need to sync gradients!
Max GPU memory constrains a layer shard

# Scaling Parameters | FSDP + TP = ~Limitless Scaling



FSDP

checkpoint_wrapper

Decoder ~12H² params

FSDP

Multi-Head Attention ~4H² params

| Query ~H² params | Key ~H² params | Value ~H² params | Out ~H² params |

FSDP

FC1 ~4H² params

FSDP

FC2 ~4H² params

**Outer activations are offloaded to CPU**

**Inner activations are checkpointed/recomputed**

**Linear layers are sharded**

**Full parameters of a single FSDP instance are loaded to GPU for computation, parameters in other FSDP instance are offloaded to CPUs**

[1 Trillion Parameter Model with Tensor Parallelism and FSDP](#)

# Final Questions?

## Fill out my anonymous feedback form