

# CSCSE 638 Natural Language Processing Foundation and Techniques

## Lecture 10: Pre-Training and Model Distillation

Kuan-Hao Huang

Spring 2025



(Some slides adapted from Vivian Chen)

# SIGIR 2025 LiveRAG Challenge

- <https://sigir2025.dei.unipd.it/live-rag-challenge.html>
- RAG: Retrieval Augmented Generation
- Advance RAG research and compare the performance of their solutions with other teams on a fixed corpus



# SIGIR 2025 LiveRAG Challenge

Date (2025)	Details
<b>Feb 24</b>	<b>Application submission deadline</b> - <a href="#">SIGIR2025 easychair site</a> (Select: SIGIR2025 LiveRAG Challenge track)
Mar 12	<ul style="list-style-type: none"><li>• Application submission notification to selected teams</li><li>• Opening of easychair site for short paper submission</li><li>• AWS and Pinecone resources and credits made available to selected teams together with detailed operational instructions</li></ul>
Mar 15	Training and testing tool ( <a href="#">DataMorgana</a> ) made available to teams
May 8	"Dry" test for participants of live service on a small question set
<b>May 12</b>	<b>Live Challenge Day</b> – test questions shared and live service for answers submission opens
<b>May 19</b>	<b>Short paper submission deadline</b> - <a href="#">SIGIR2025 easychair site</a> (Select: SIGIR2025 LiveRAG Challenge track)
May 29	Short paper notification and announcement of finalists
July 17	<ul style="list-style-type: none"><li>• LiveRAG Workshop at SIGIR'2025 in Padua, Italy</li><li>• Presentation of research by selected teams</li><li>• Announcement of winner and runner(s)-up</li></ul>

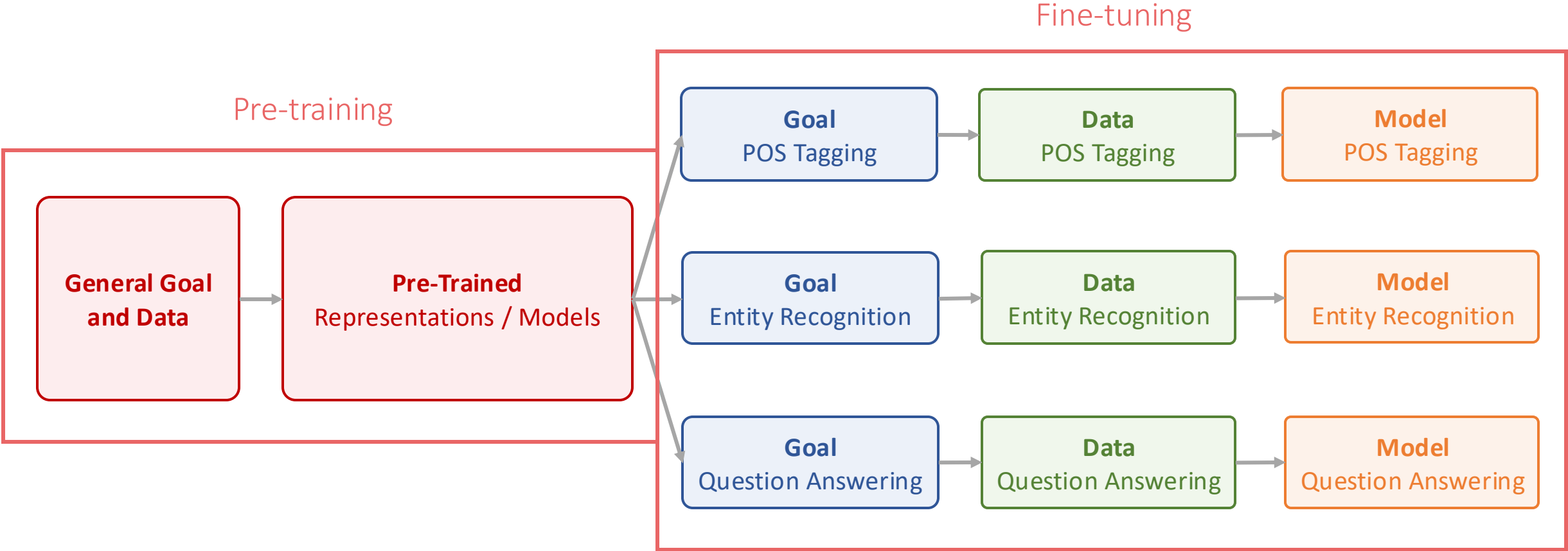
## PRIZES

- First Prize: \$5000
- Second Prize: \$3000
- Third Prize: \$2000

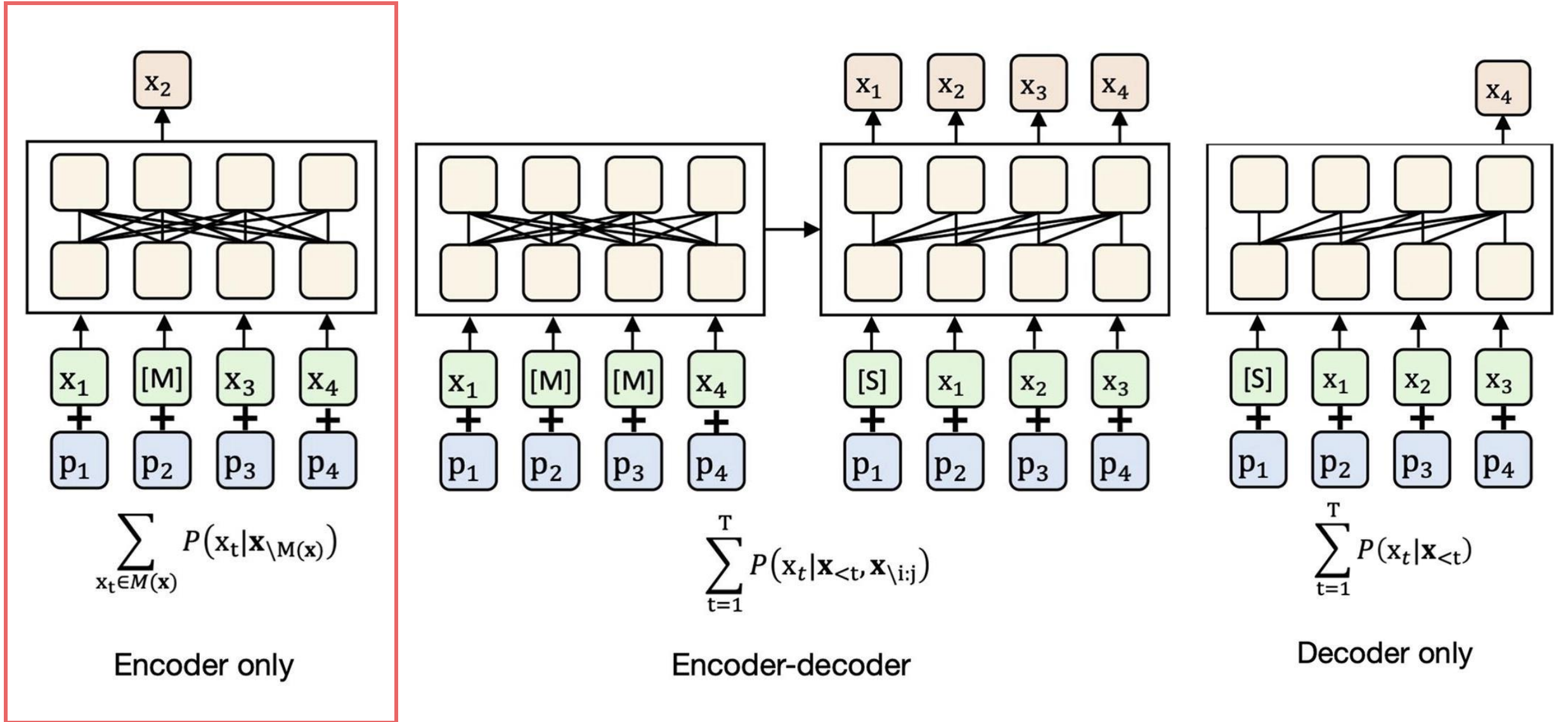
# Lecture Plan

- Pre-Training
  - Encoder-Only Pre-Training
  - Encoder-Decoder Pre-Training
  - Decoder-Only Pre-Training
- Model Distillation

# Recap: Fine-Tuning with Pre-Training

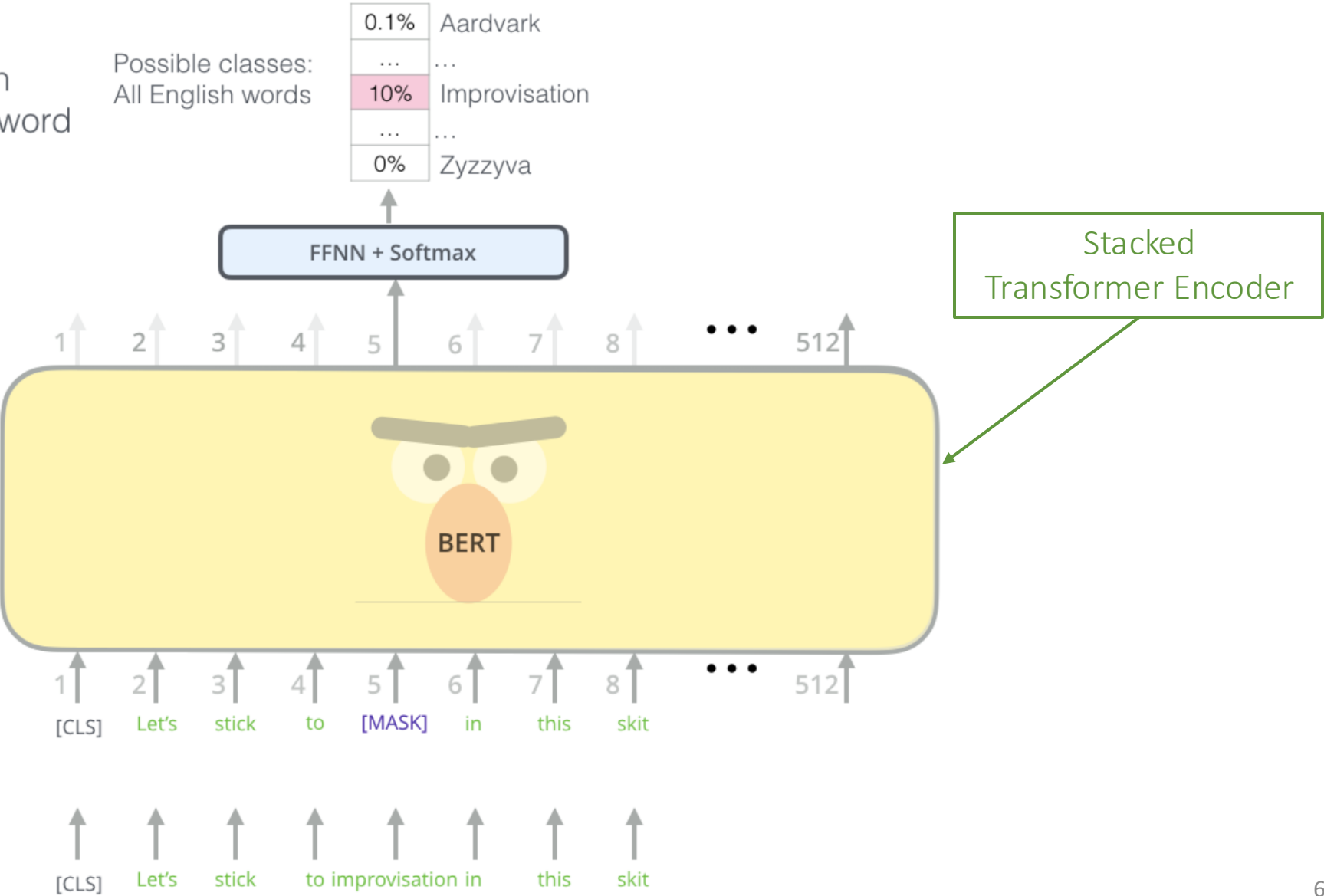


# Recap: Types of Pre-Training



# Recap: BERT – Masked Language Modeling

Use the output of the masked word's position to predict the masked word

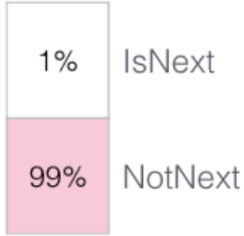


Randomly mask 15% of tokens

Input

# Recap: BERT – Next Sentence Prediction

Predict likelihood that sentence B belongs after sentence A

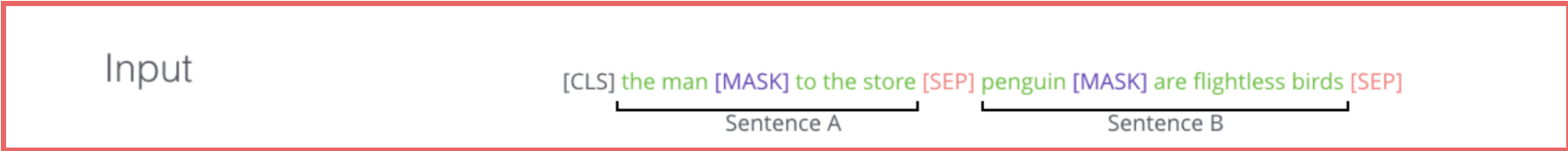


FFNN + Softmax



Positive example: real next sentence  
Negative example: random sentence

Tokenized Input

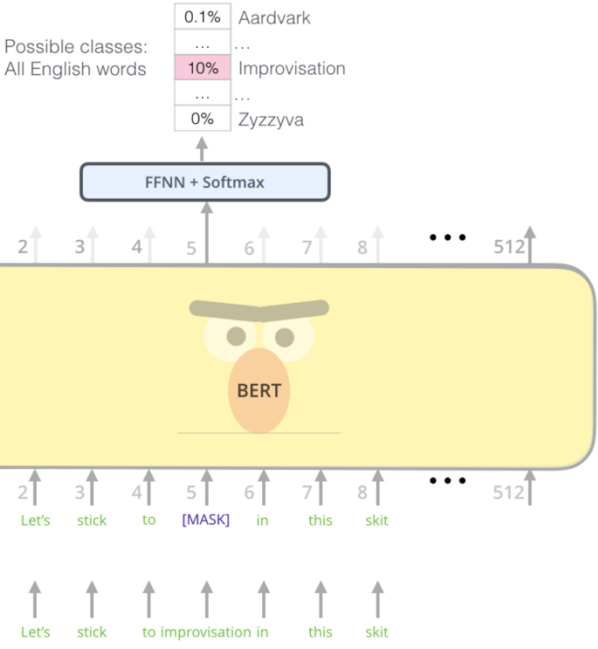




# Recap: Other Encoder-Only Pre-Trained Models

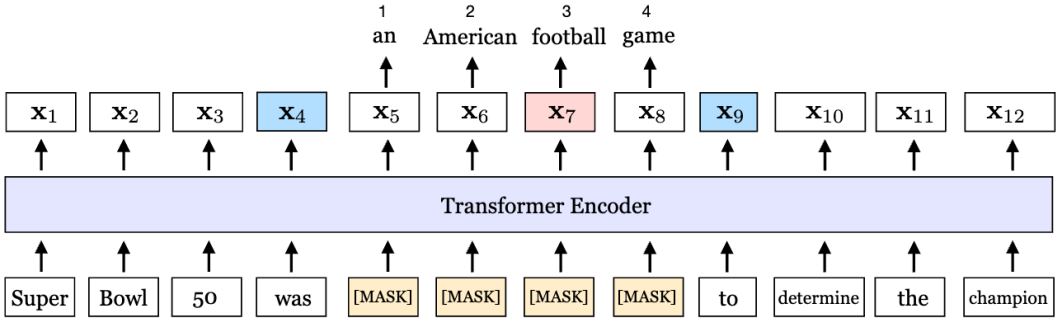
- RoBERTa
- SpanBERT

Use the output of the masked word's position to predict the masked word

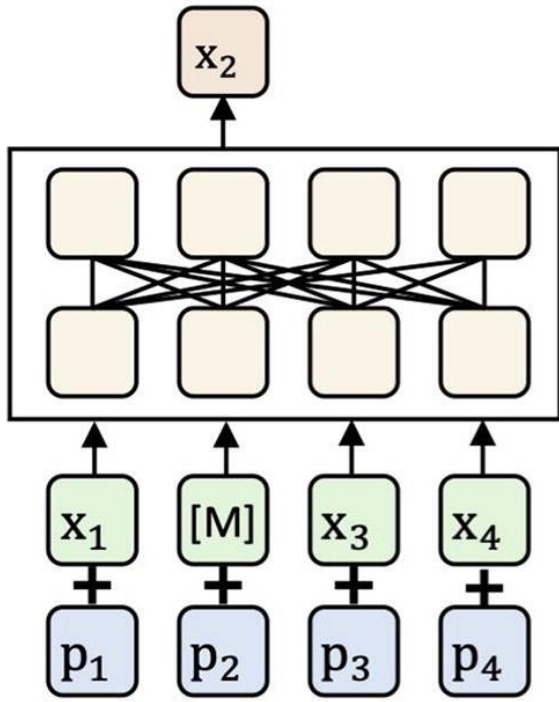


$$\mathcal{L}(\text{football}) = \mathcal{L}_{\text{MLM}}(\text{football}) + \mathcal{L}_{\text{SBO}}(\text{football})$$

$$= -\log P(\text{football} | \mathbf{x}_7) - \log P(\text{football} | \mathbf{x}_4, \mathbf{x}_9, \mathbf{p}_3)$$

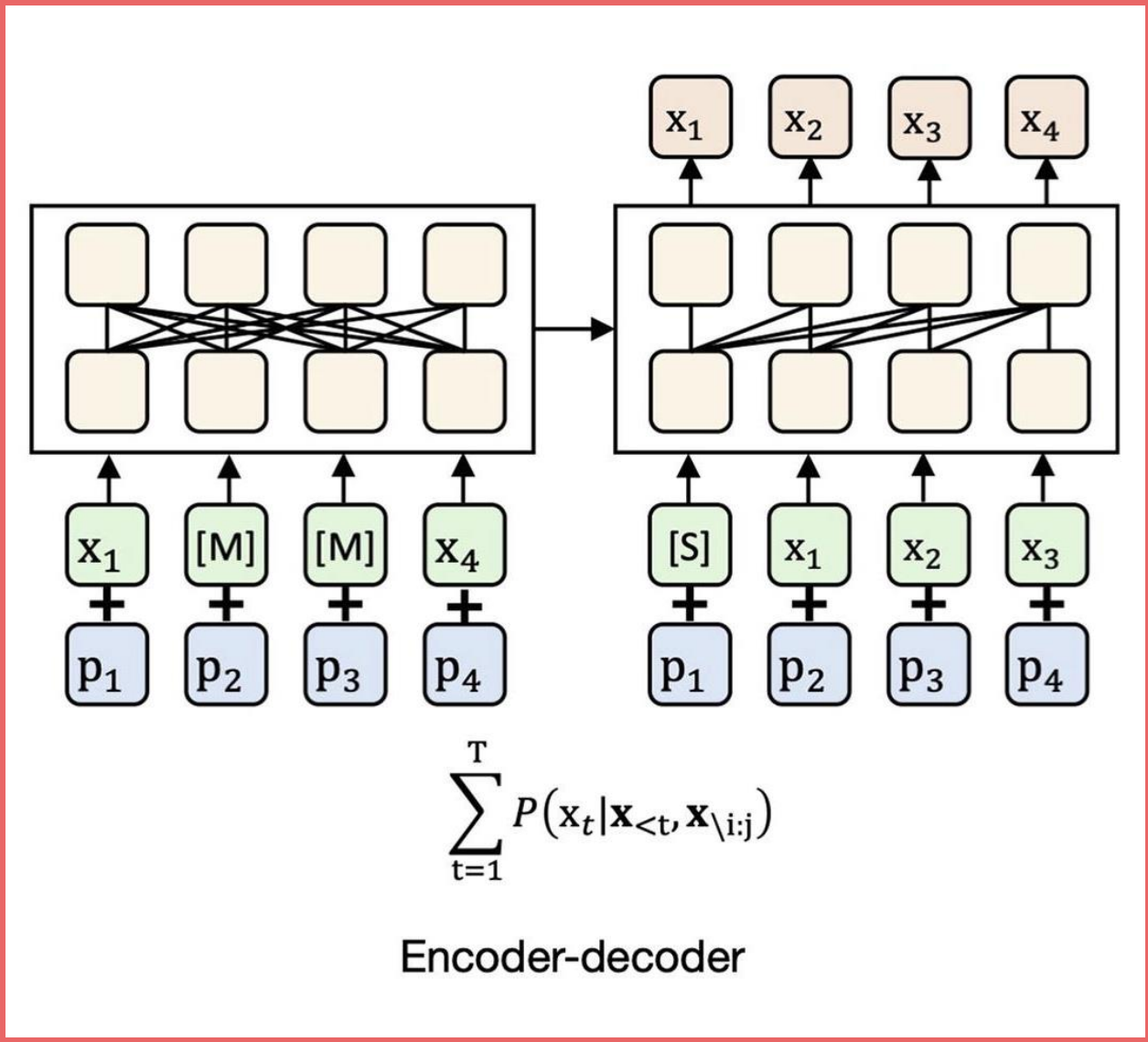


# Recap: Types of Pre-Training



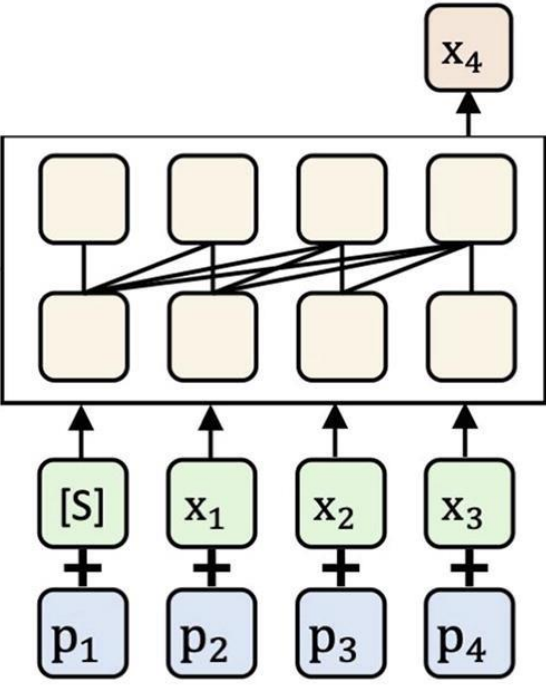
$$\sum_{x_t \in M(x)} P(x_t | \mathbf{x}_{\setminus M(x)})$$

Encoder only



$$\sum_{t=1}^T P(x_t | \mathbf{x}_{<t}, \mathbf{x}_{\setminus i;j})$$

Encoder-decoder

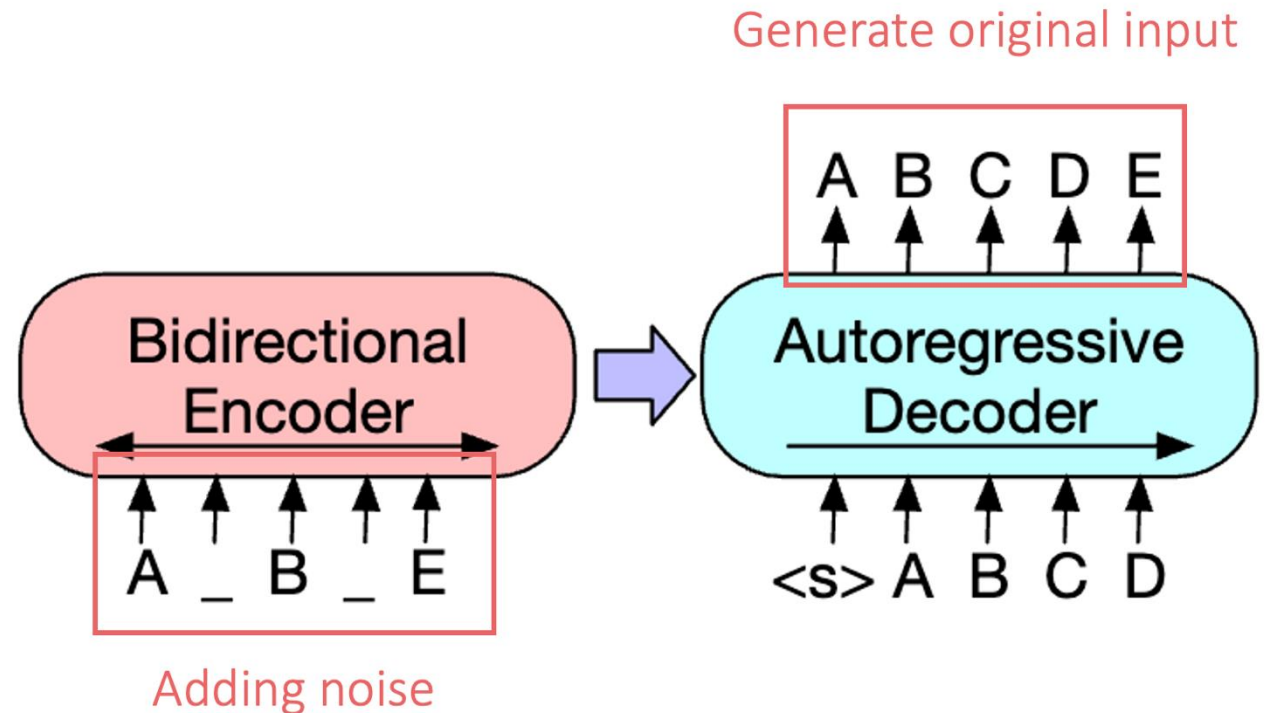


$$\sum_{t=1}^T P(x_t | \mathbf{x}_{<t})$$

Decoder only

# Recap: BART – Denoising Objective

- Token Masking
  - A<mask>CD<mask>F. → ABCDEF
- Token Deletion
  - ACDF. → ABCDEF.
- Text Infilling
  - A<mask>D<mask>F. → ABCDEF.
- Sentence Permutation
  - FG. ABC. DE. → ABC. DE. FG.
- Document Rotation
  - E. FG. ABC. D → ABC. DE. FG.



# Encoder-Decoder: T5

- Text-to-Text Transfer Transformer (T5)

## Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer

**Colin Raffel\***

CRAFFEL@GMAIL.COM

**Noam Shazeer\***

NOAM@GOOGLE.COM

**Adam Roberts\***

ADAROB@GOOGLE.COM

**Katherine Lee\***

KATHERINELEE@GOOGLE.COM

**Sharan Narang**

SHARANNARANG@GOOGLE.COM

**Michael Matena**

MMATENA@GOOGLE.COM

**Yanqi Zhou**

YANQIZ@GOOGLE.COM

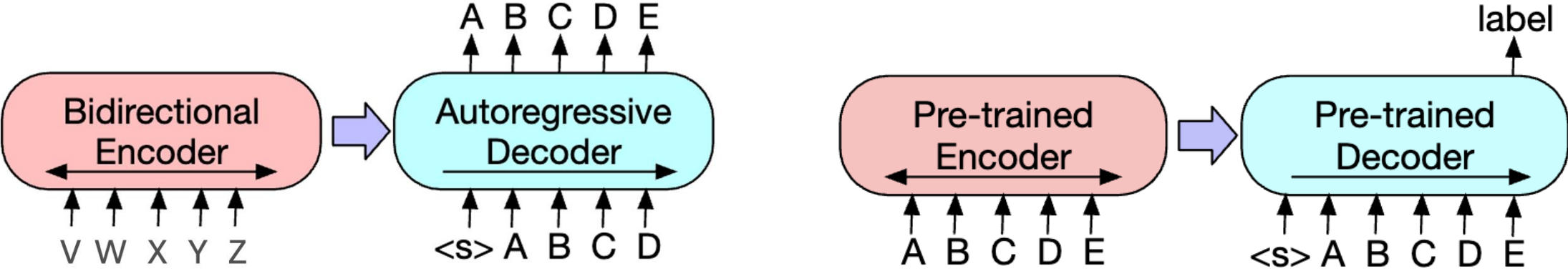
**Wei Li**

MWEILI@GOOGLE.COM

**Peter J. Liu**

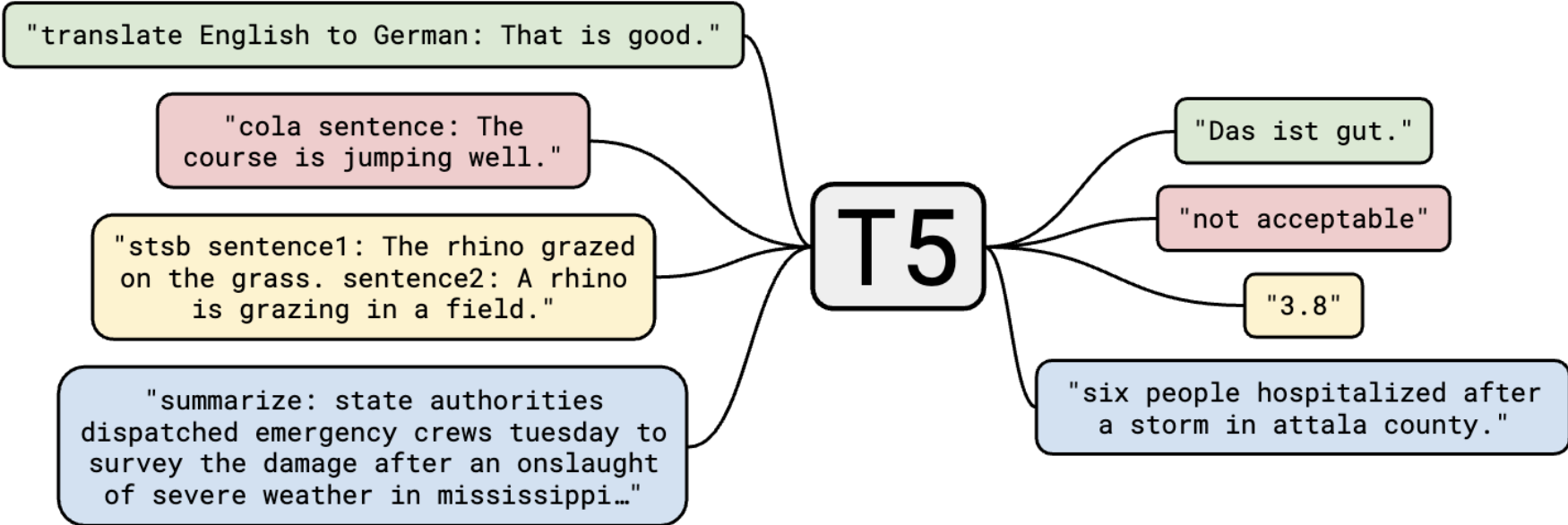
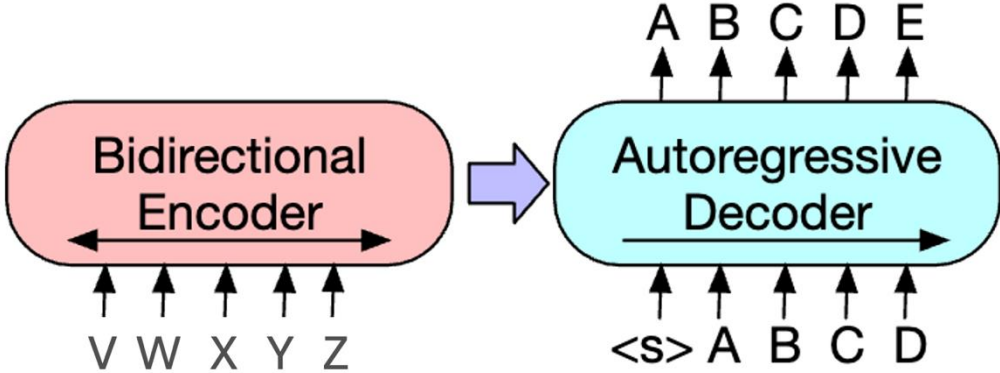
PETERJLIU@GOOGLE.COM

# Motivation: BART



Different ways when considering classification and seq2seq generation

# Convert Everything to Text-to-Text Tasks



# Masked Span Reconstruction (Seq2Seq Version)

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>

# Multi-Task Learning

- Convert everything to text-to-text tasks
- Jointly fine-tune them together



# Multi-Task Learning

## D.7 SST2

**Original input:**

**Sentence:** it confirms fincher 's status as a film maker who artfully bends technical know-how to the service of psychological insight .

**Processed input:** sst2 sentence: it confirms fincher 's status as a film maker who artfully bends technical know-how to the service of psychological insight .

**Original target:** 1

**Processed target:** positive

# Multi-Task Learning

## D.4 MRPC

### Original input:

**Sentence 1:** We acted because we saw the existing evidence in a new light ,  
through the prism of our experience on 11 September , " Rumsfeld said .

**Sentence 2:** Rather , the US acted because the administration saw " existing  
evidence in a new light , through the prism of our experience on September  
11 " .

**Processed input:** mrpc sentence1: We acted because we saw the existing evidence  
in a new light , through the prism of our experience on 11 September , " Rumsfeld  
said . sentence2: Rather , the US acted because the administration saw "  
existing evidence in a new light , through the prism of our experience on  
September 11 " .

**Original target:** 1

**Processed target:** equivalent

# Multi-Task Learning

## D.16 WMT English to German

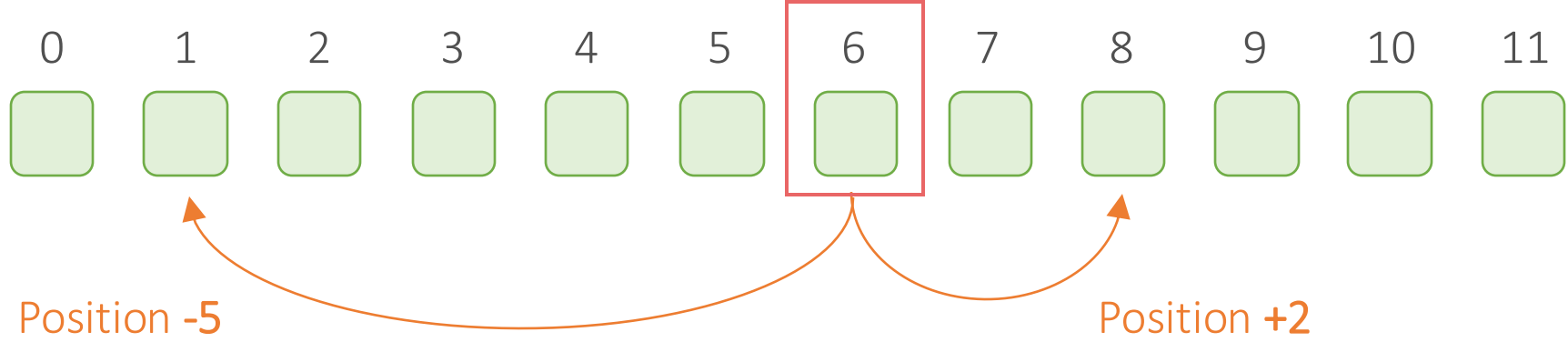
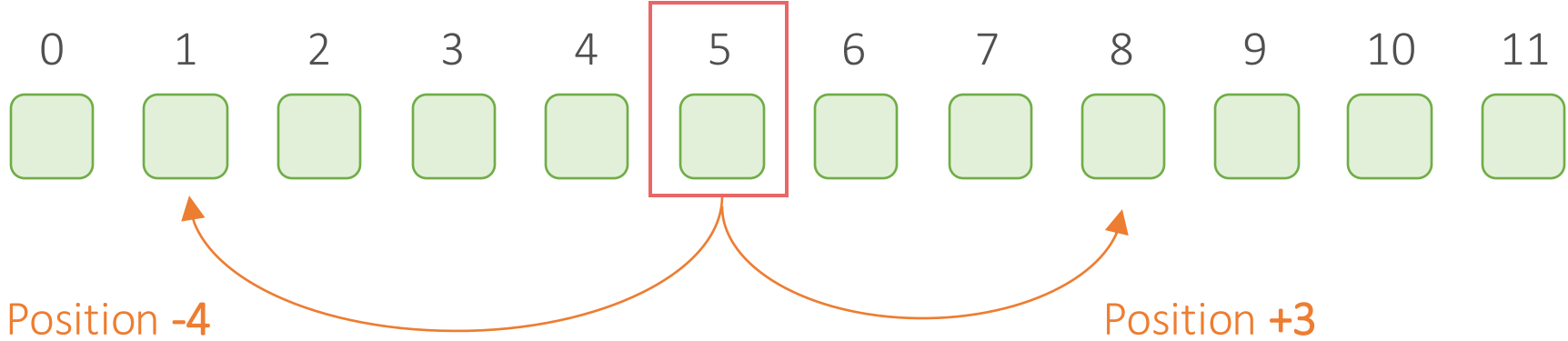
**Original input:** "Luigi often said to me that he never wanted the brothers to end up in court," she wrote.

**Processed input:** translate English to German: "Luigi often said to me that he never wanted the brothers to end up in court," she wrote.

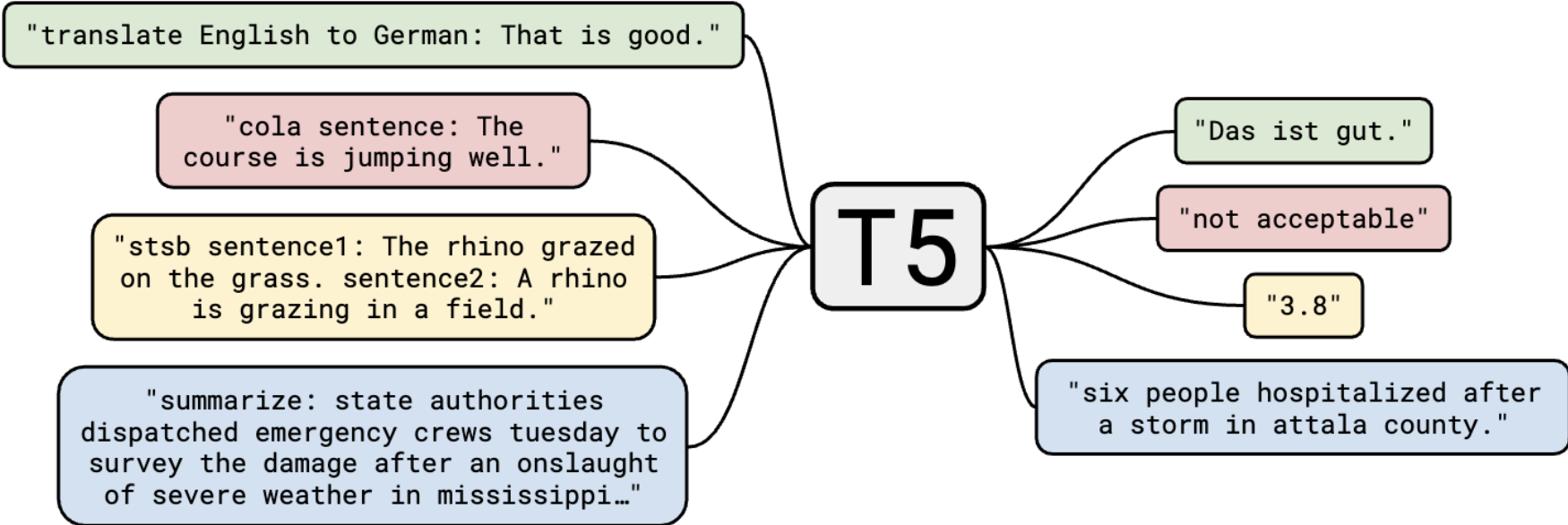
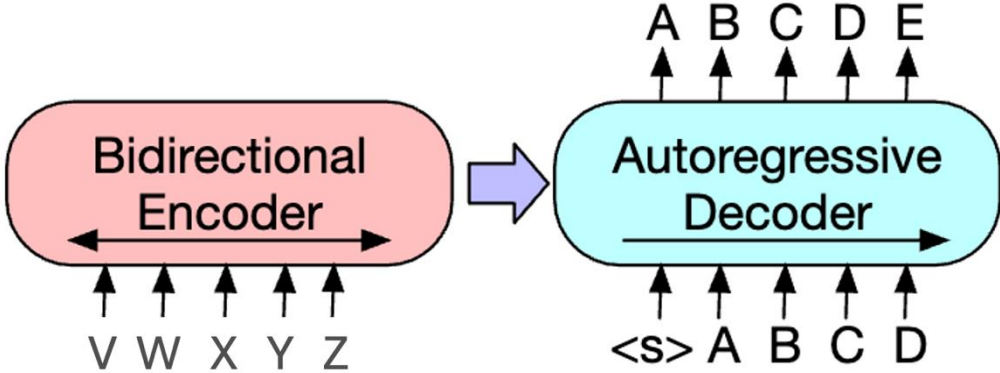
**Original target:** "Luigi sagte oft zu mir, dass er nie wollte, dass die Brüder vor Gericht landen", schrieb sie.

**Processed target:** "Luigi sagte oft zu mir, dass er nie wollte, dass die Brüder vor Gericht landen", schrieb sie.

# Relative Position



# Fine-Tuning: Text-to-Text For Everything



# Promising Results

Model	QQP F1	QQP Accuracy	MNLI-m Accuracy	MNLI-mm Accuracy	QNLI Accuracy	RTE Accuracy	WNLI Accuracy
Previous best	74.8 <sup>c</sup>	<b>90.7<sup>b</sup></b>	91.3 <sup>a</sup>	91.0 <sup>a</sup>	<b>99.2<sup>a</sup></b>	89.2 <sup>a</sup>	91.8 <sup>a</sup>
T5-Small	70.0	88.0	82.4	82.3	90.3	69.9	69.2
T5-Base	72.6	89.4	87.1	86.2	93.7	80.1	78.8
T5-Large	73.9	89.9	89.9	89.6	94.8	87.2	85.6
T5-3B	74.4	89.7	91.4	91.2	96.3	91.1	89.7
T5-11B	<b>75.1</b>	90.6	<b>92.2</b>	<b>91.9</b>	96.9	<b>92.8</b>	<b>94.5</b>

Model	SQuAD EM	SQuAD F1	SuperGLUE Average	BoolQ Accuracy	CB F1	CB Accuracy	COPA Accuracy
Previous best	90.1 <sup>a</sup>	95.5 <sup>a</sup>	84.6 <sup>d</sup>	87.1 <sup>d</sup>	90.5 <sup>d</sup>	95.2 <sup>d</sup>	90.6 <sup>d</sup>
T5-Small	79.10	87.24	63.3	76.4	56.9	81.6	46.0
T5-Base	85.44	92.08	76.2	81.4	86.2	94.0	71.2
T5-Large	86.66	93.79	82.3	85.4	91.6	94.8	83.4
T5-3B	88.53	94.95	86.4	89.9	90.3	94.4	92.0
T5-11B	<b>91.26</b>	<b>96.22</b>	<b>88.9</b>	<b>91.2</b>	<b>93.9</b>	<b>96.8</b>	<b>94.8</b>

Model	MultiRC F1a	MultiRC EM	ReCoRD F1	ReCoRD Accuracy	RTE Accuracy	WiC Accuracy	WSC Accuracy
Previous best	84.4 <sup>d</sup>	52.5 <sup>d</sup>	90.6 <sup>d</sup>	90.0 <sup>d</sup>	88.2 <sup>d</sup>	69.9 <sup>d</sup>	89.0 <sup>d</sup>
T5-Small	69.3	26.3	56.3	55.4	73.3	66.9	70.5
T5-Base	79.7	43.1	75.0	74.2	81.5	68.3	80.8
T5-Large	83.3	50.7	86.8	85.9	87.8	69.3	86.3
T5-3B	86.8	58.3	91.2	90.4	90.7	72.1	90.4
T5-11B	<b>88.1</b>	<b>63.3</b>	<b>94.1</b>	<b>93.4</b>	<b>92.5</b>	<b>76.9</b>	<b>93.8</b>

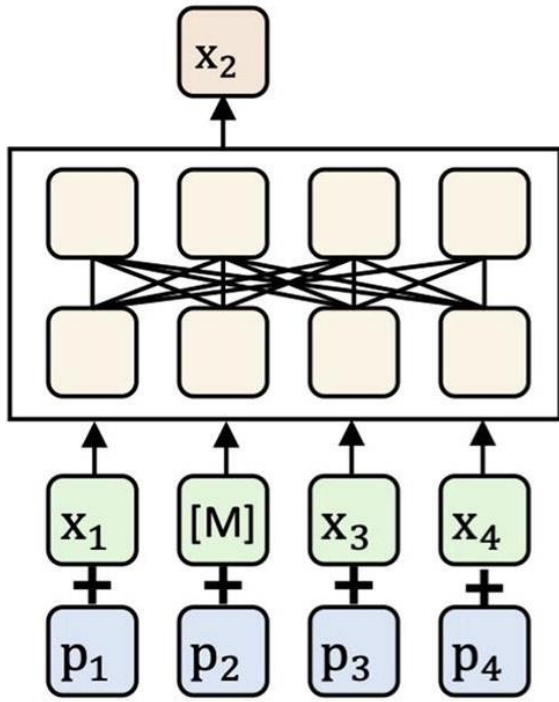
# Use T5



## Hugging Face

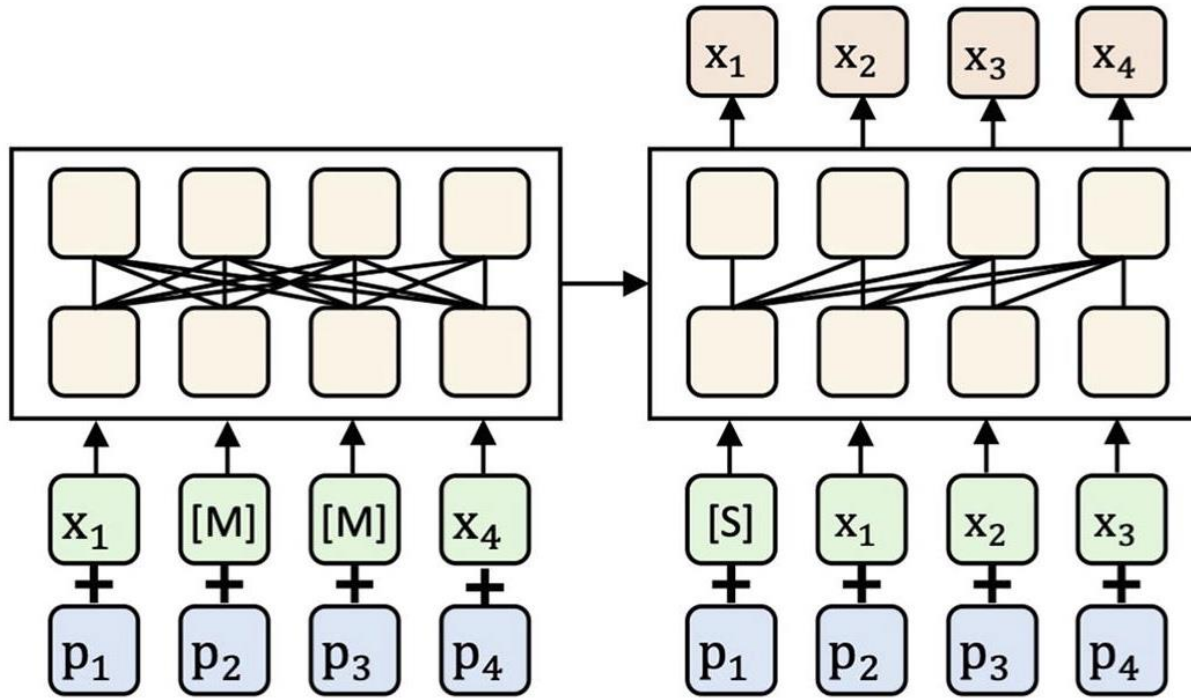
- T5-small:
  - # parameters  $\approx$  60M
- T5-base:
  - # parameters  $\approx$  220M
- T5-large:
  - # parameters  $\approx$  770M
- T5-3B: #
  - parameters  $\approx$  3B
- T5-11B:
  - # parameters  $\approx$  11B

# Types of Pre-Training



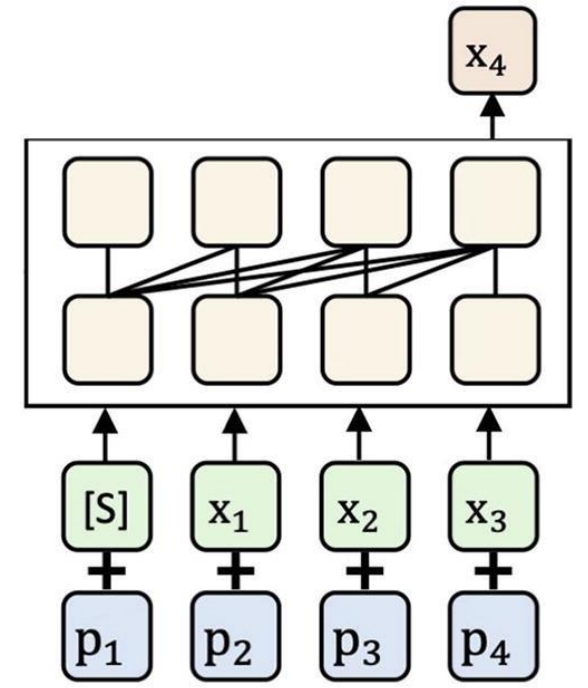
$$\sum_{x_t \in M(x)} P(x_t | \mathbf{x}_{\setminus M(x)})$$

Encoder only



$$\sum_{t=1}^T P(x_t | \mathbf{x}_{<t}, \mathbf{x}_{\setminus i;j})$$

Encoder-decoder



$$\sum_{t=1}^T P(x_t | \mathbf{x}_{<t})$$

Decoder only

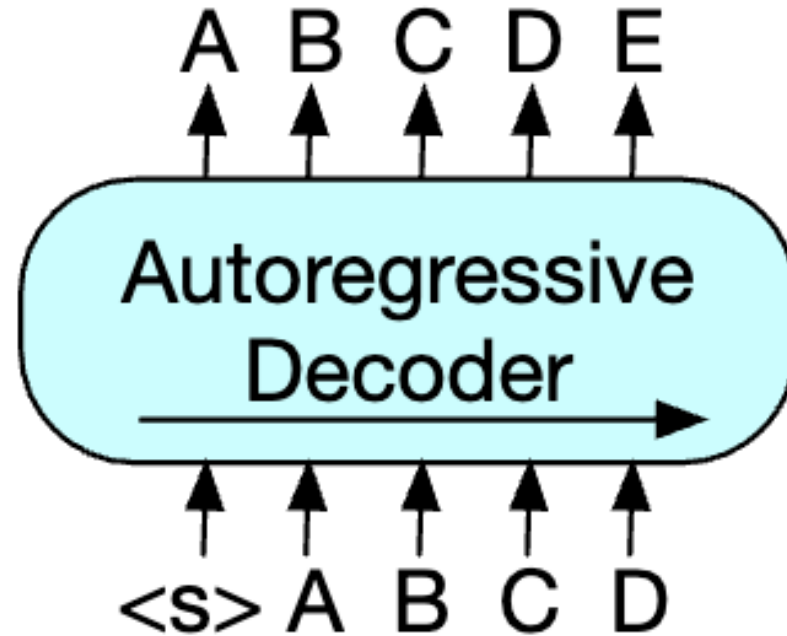


# Decoder-Only: GPT

- Improving Language Understanding by Generative Pre-Training, OpenAI 2018
  - **Generative Pre-trained Transformer (GPT)**
- Language Models are Unsupervised Multitask Learners, OpenAI 2019
  - GPT-2
- Language Models are Few-Shot Learners, OpenAI 2020
  - GPT-3

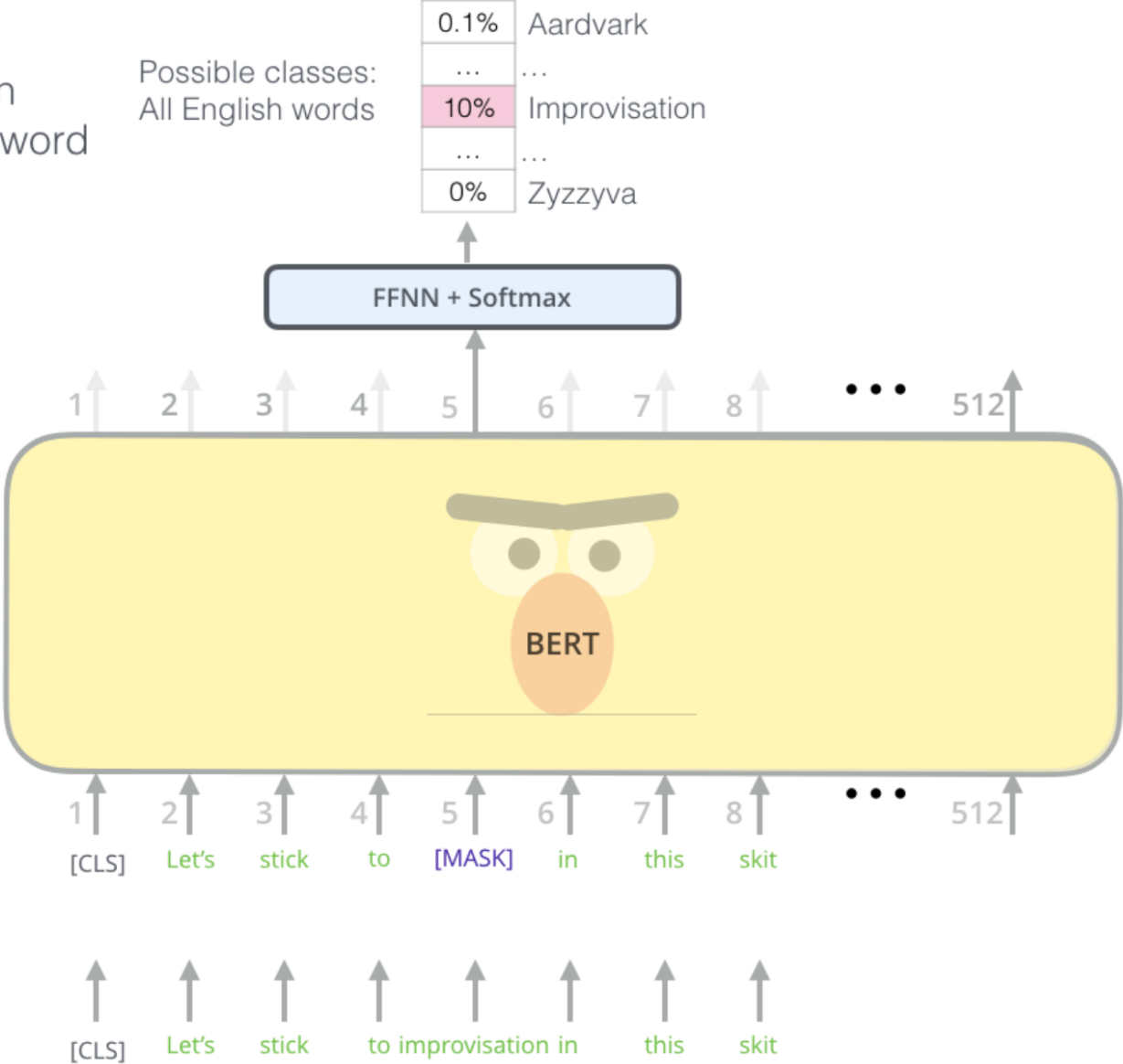
# Language Modeling

- Next word prediction
- Trained with large corpus

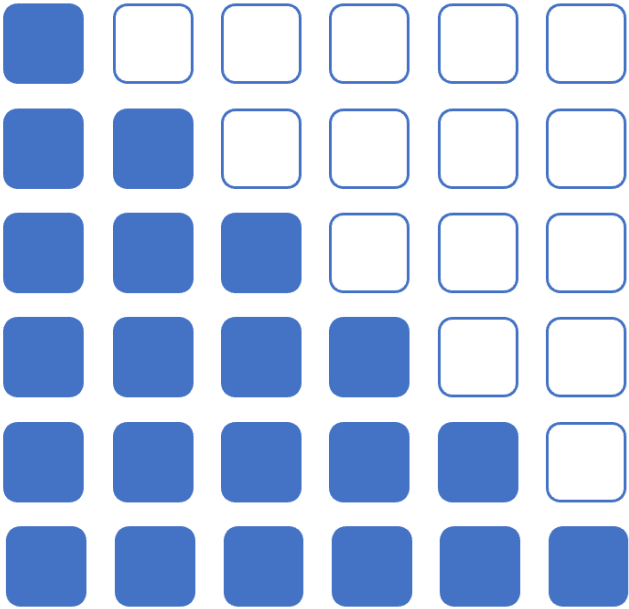
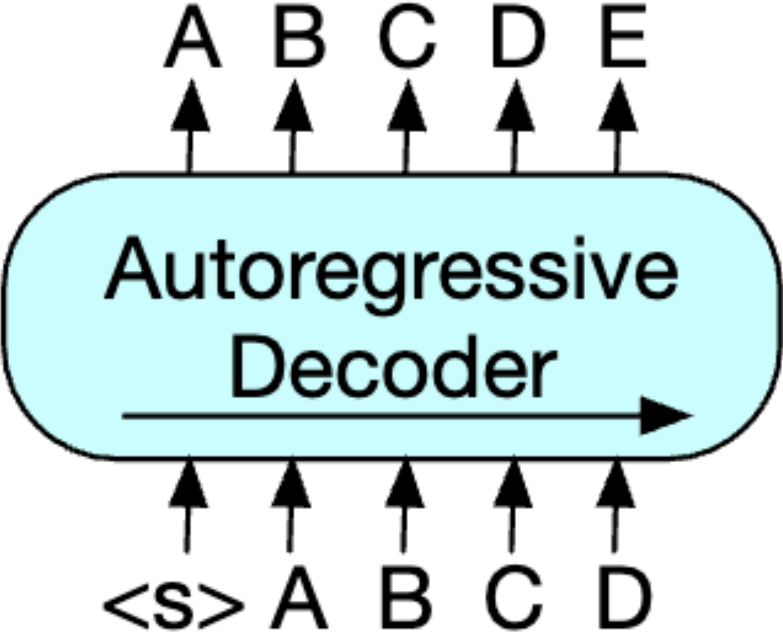


# Comparison: Masked Language Models

Use the output of the masked word's position to predict the masked word



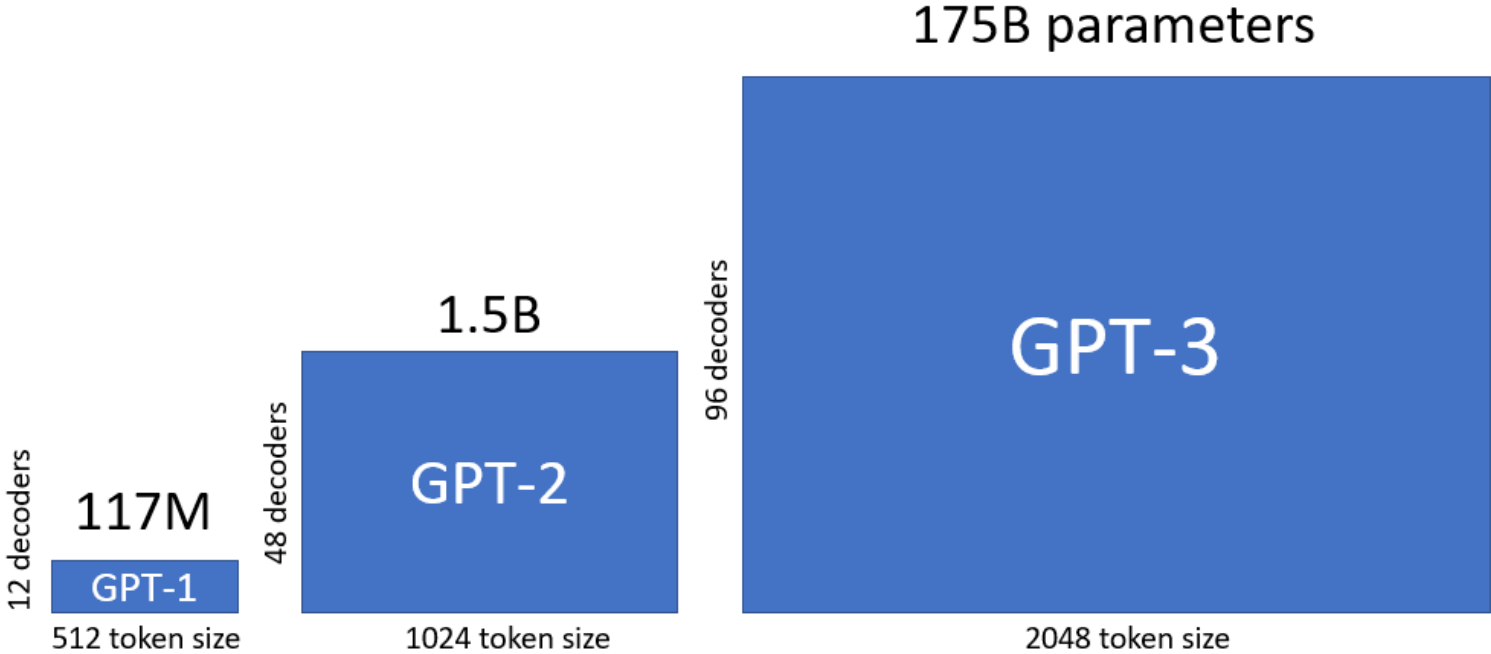
# Comparison: Causal Language Models



Causal Masking

# GPT-3: From Fine-Tuning to Few-Shot Learning

- Even larger training data, even larger model size



# GPT-3: From Fine-Tuning to Few-Shot Learning

- Solve entirely new tasks by **few-shot learning (in-context learning)**

Circulation revenue has increased by 5% in Finland. // Positive

Panostaja did not disclose the purchase price. // Neutral

Paying off the national debt will be extremely painful. // Negative

The company anticipated its operating profit to improve. // \_\_\_\_\_



Circulation revenue has increased by 5% in Finland. // Finance

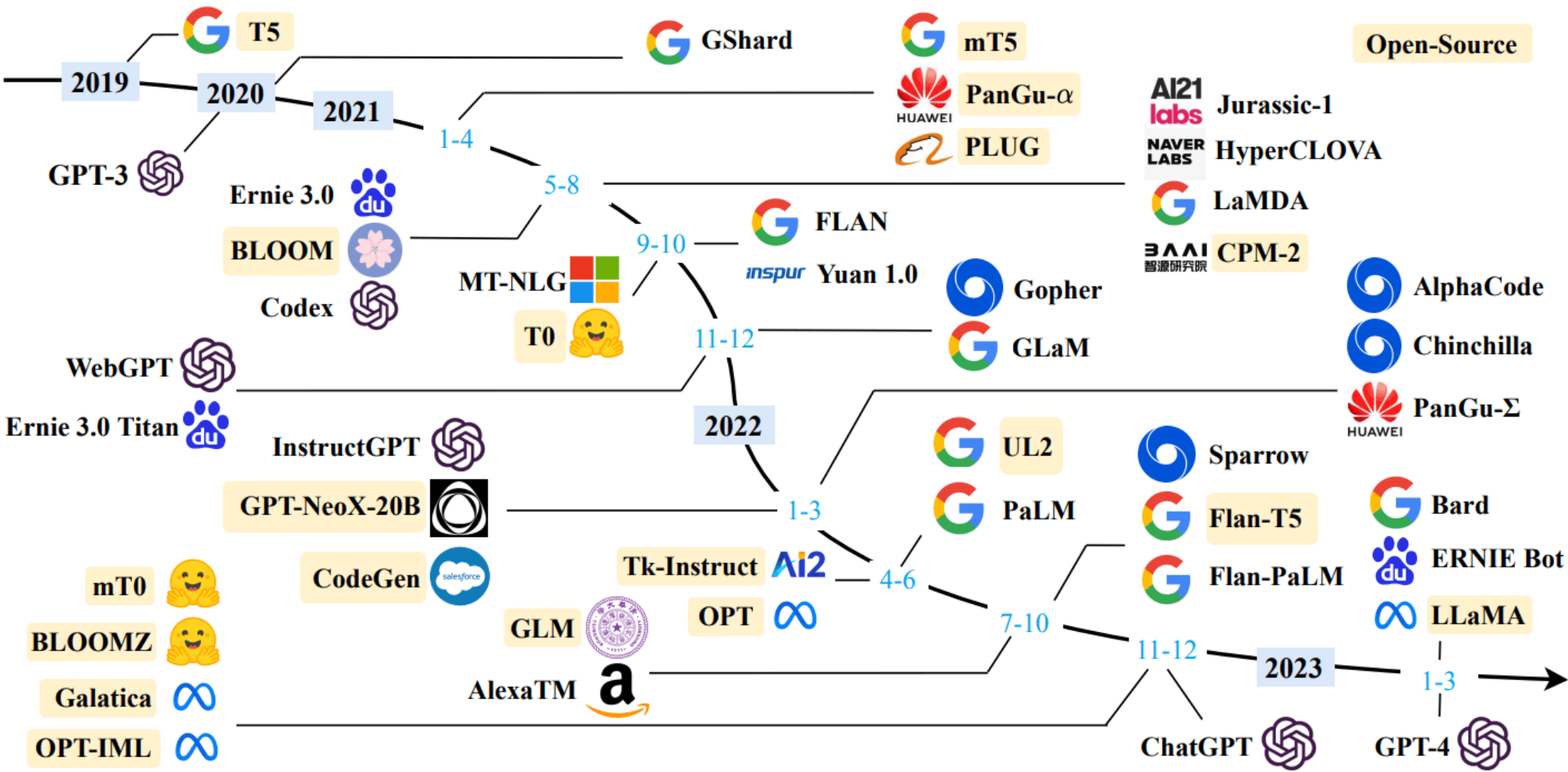
They defeated ... in the NFC Championship Game. // Sports

Apple ... development of in-house chips. // Tech

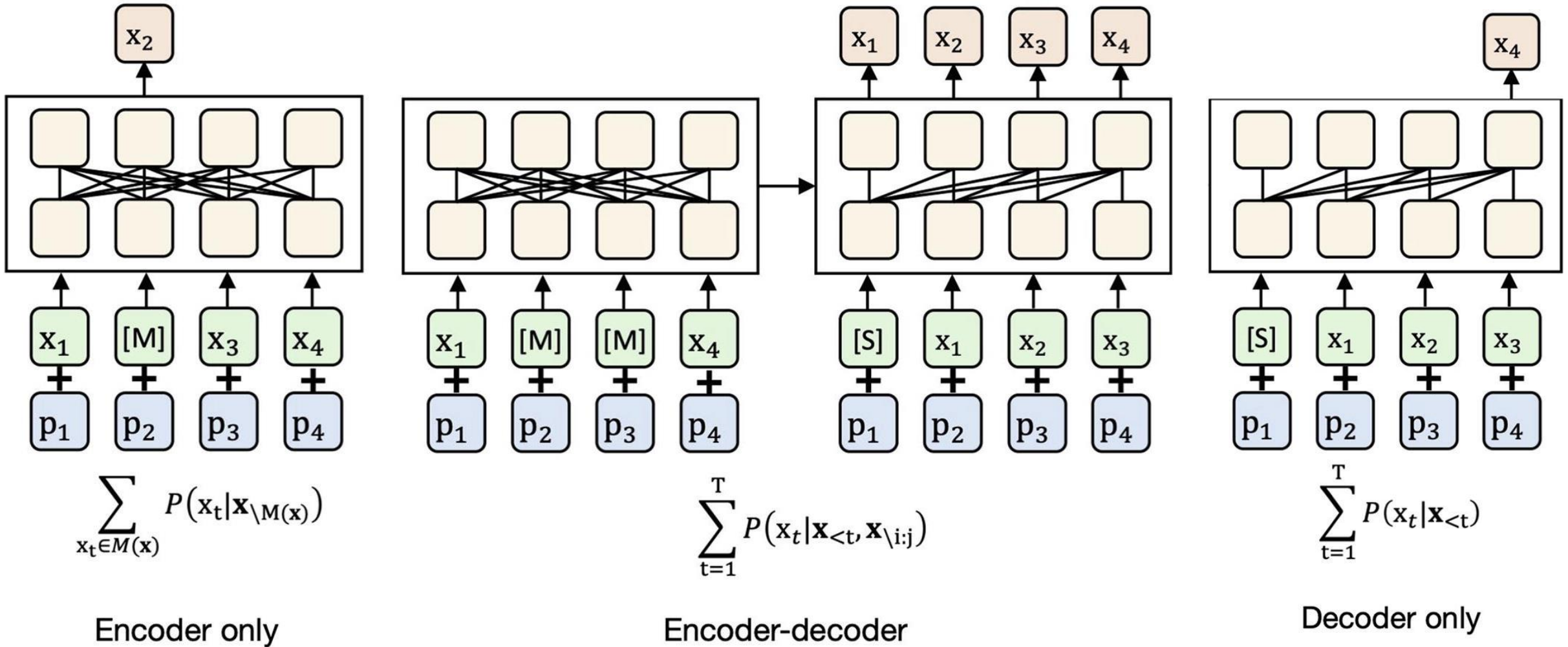
The company anticipated its operating profit to improve. // \_\_\_\_\_



# Large Language Models



# Types of Pre-Training





# Use GPT



## Hugging Face

- GPT-2-small
  - # parameters  $\approx$  117M
- GPT-2-medium
  - # parameters  $\approx$  345M
- GPT-2-large
  - # parameters  $\approx$  762M
- GPT-2-xl
  - # parameters  $\approx$  1.5B

# Lecture Plan

- Pre-Training
  - Encoder-Only Pre-Training
  - Encoder-Decoder Pre-Training
  - Decoder-Only Pre-Training
- Model Distillation

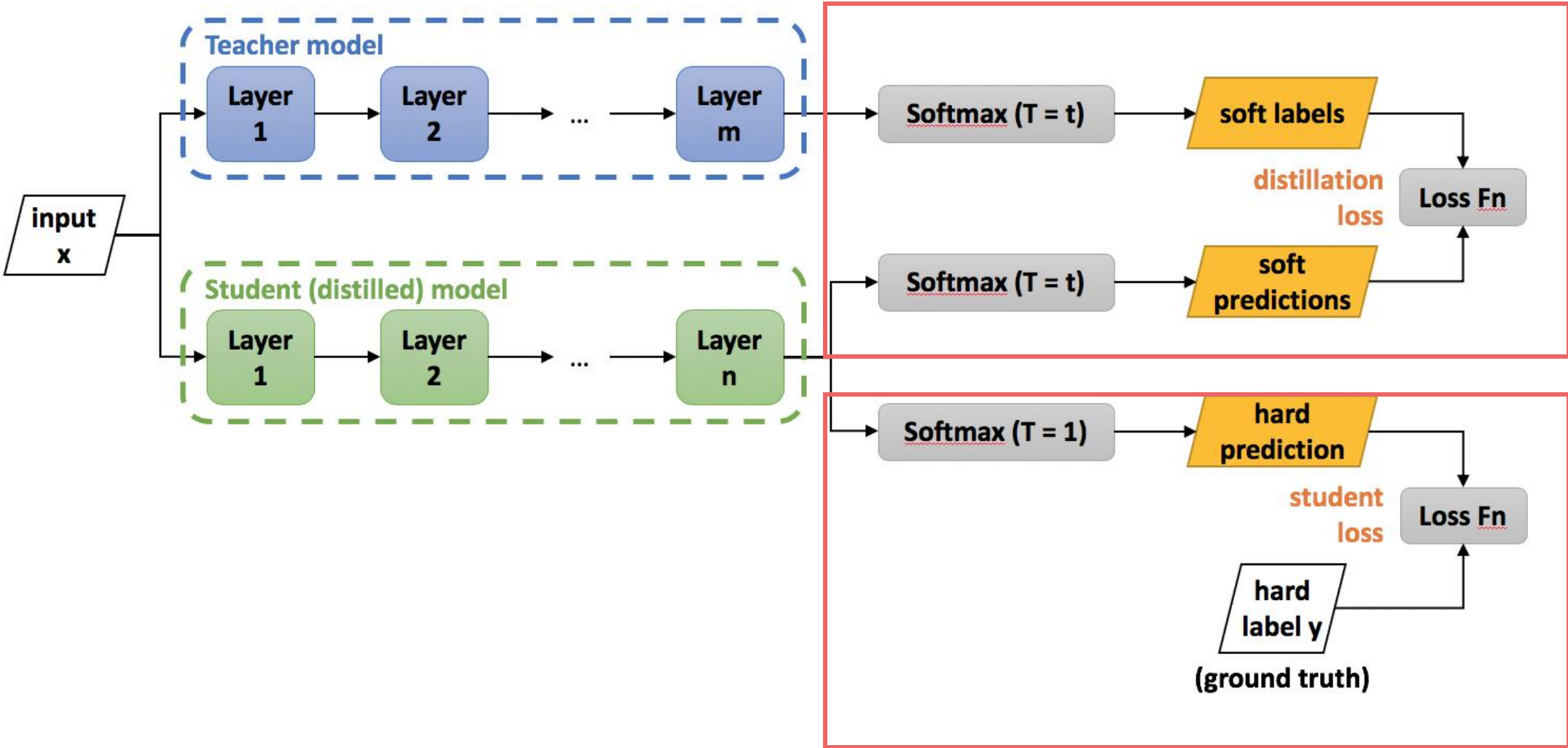
# Model Distillation

- Distill knowledge from a large model to a small model while maintaining similar capability
  - **Large model:** teacher model
  - **Small model:** student model
- Train a student model to mimic the behavior of the teacher model
- Reduce the number of parameters

Why don't we train a student model directly from data?

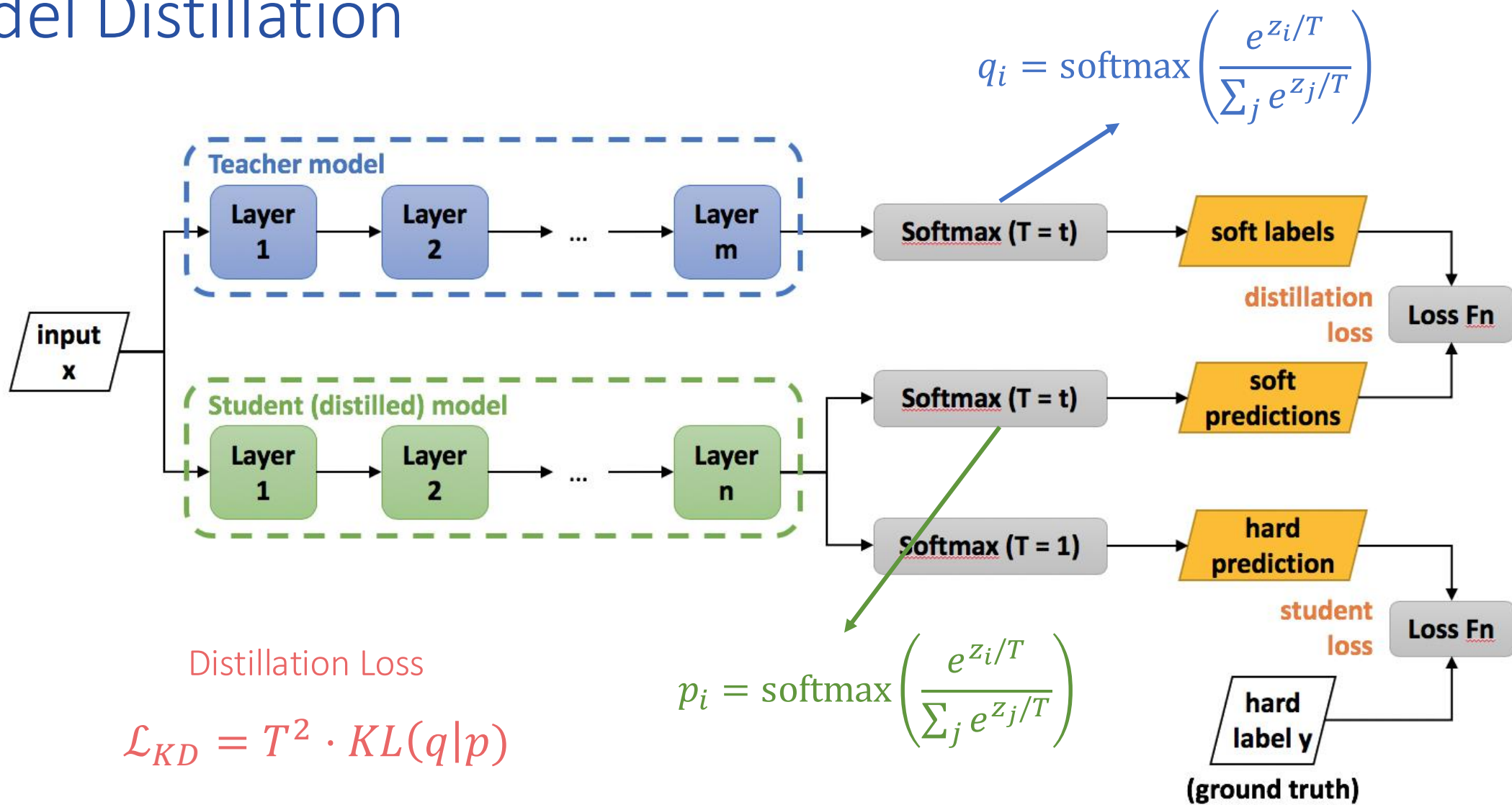
# Model Distillation

Mimic teacher's behavior

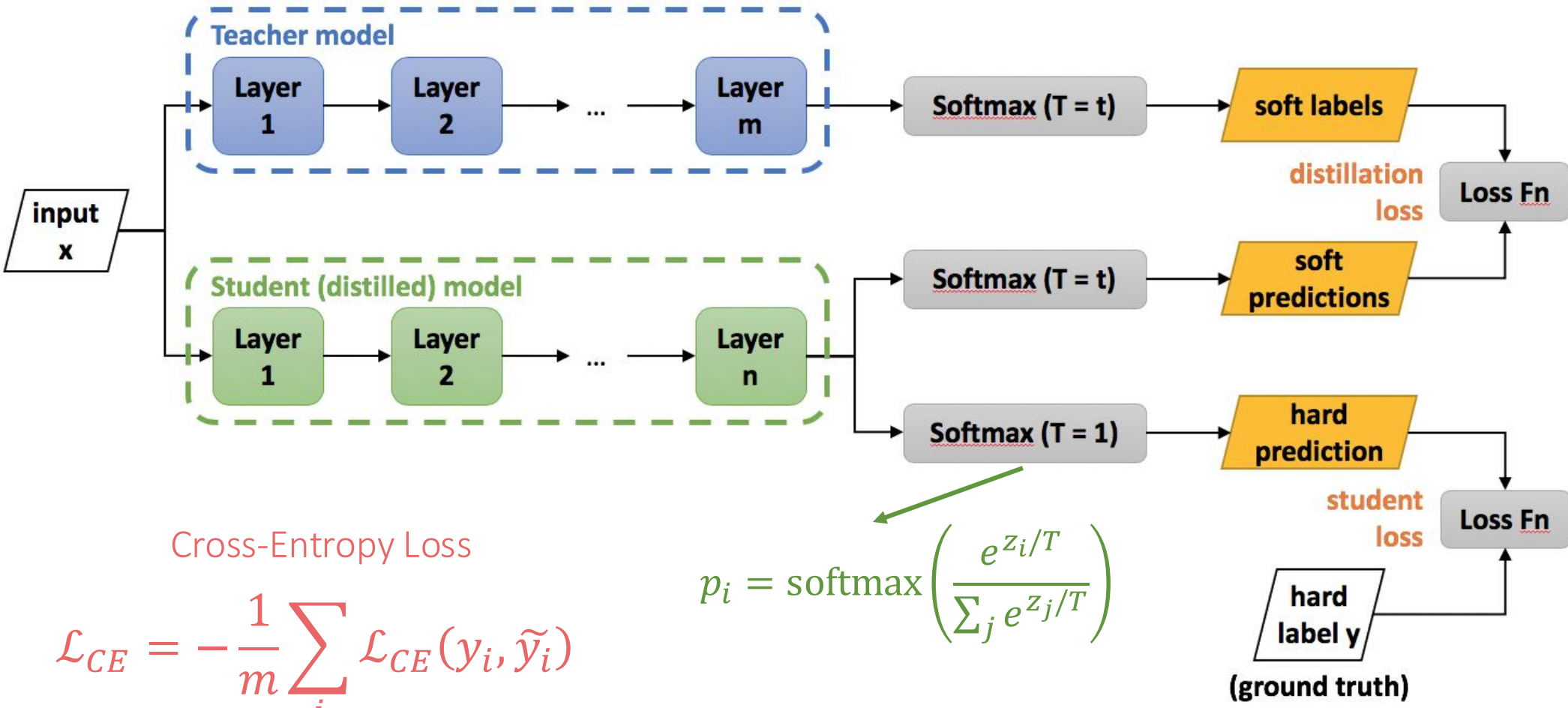


Learn from data

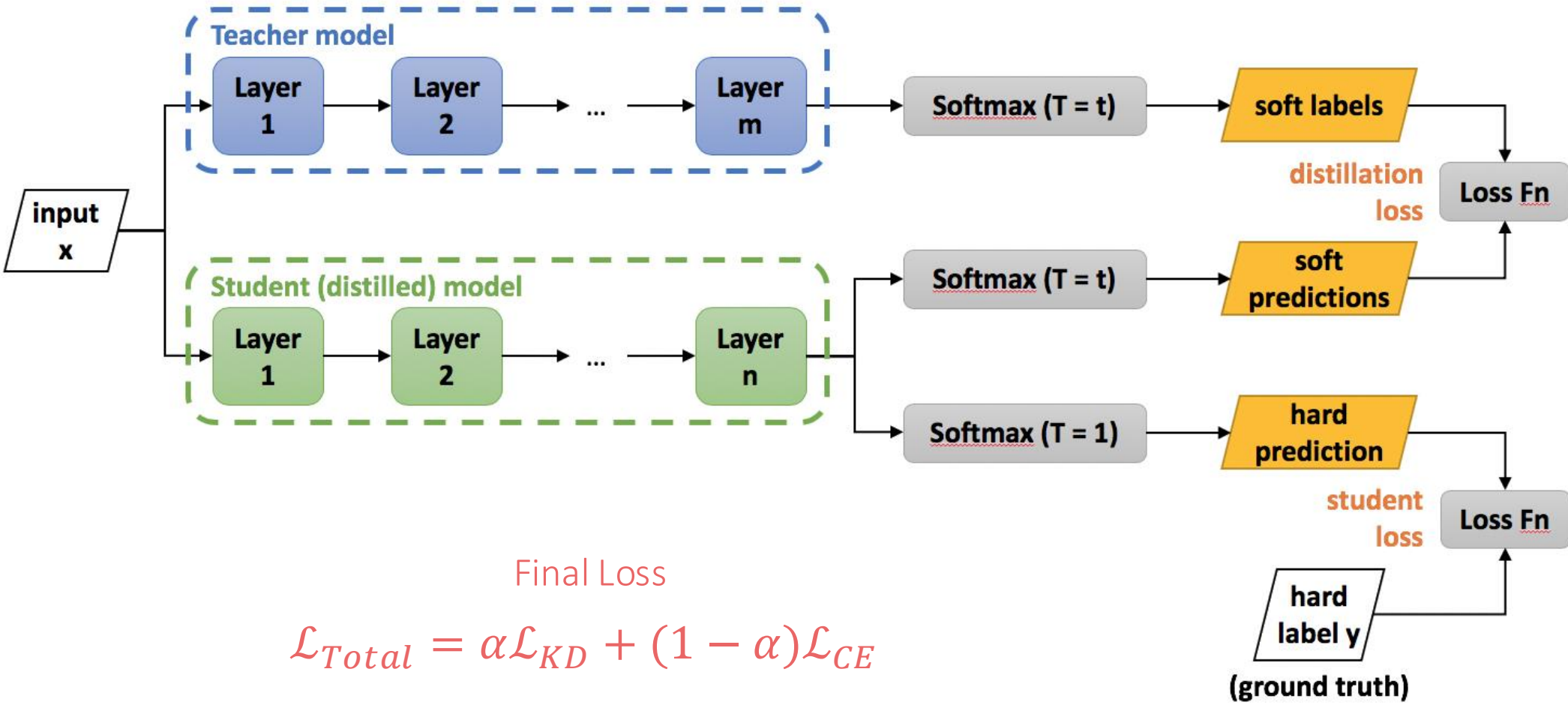
# Model Distillation



# Model Distillation



# Model Distillation



# DistilBERT

---

**DistilBERT, a distilled version of BERT: smaller,  
faster, cheaper and lighter**

---

**Victor SANH, Lysandre DEBUT, Julien CHAUMOND, Thomas WOLF**  
Hugging Face  
`{victor,lysandre,julien,thomas}@huggingface.co`



# DistilBERT

Smaller Size

Model	# param. (Millions)	Inf. time (seconds)
ELMo	180	895
BERT-base	110	668
DistilBERT	66	410

- BERT-base
  - 12 layers, hidden size = 768, 12 attention heads
- DistilBERT
  - 6 layers, hidden size = 768, 12 attention heads

Almost similar performance

Model	Score	CoLA	MNLI	MRPC	QNLI	QQP	RTE	SST-2	STS-B	WNLI
ELMo	68.7	44.1	68.6	76.6	71.1	86.2	53.4	91.5	70.4	56.3
BERT-base	79.5	56.3	86.7	88.6	91.8	89.6	69.3	92.7	89.0	53.5
DistilBERT	77.0	51.3	82.2	87.5	89.2	88.5	59.9	91.3	86.9	56.3

# MobileBERT

## **MobileBERT: a Compact Task-Agnostic BERT for Resource-Limited Devices**

**Zhiqing Sun<sup>1\*</sup>, Hongkun Yu<sup>2</sup>, Xiaodan Song<sup>2</sup>, Renjie Liu<sup>2</sup>, Yiming Yang<sup>1</sup>, Denny Zhou<sup>2</sup>**

<sup>1</sup>Carnegie Mellon University {zhiqings, yiming}@cs.cmu.edu

<sup>2</sup>Google Brain {hongkunyu, xiaodansong, renjieliu, dennyzhou}@google.com

# MobileBERT

- Instead of less layers, reduce the hidden size

			BERT <sub>LARGE</sub>	BERT <sub>BASE</sub>	IB-BERT <sub>LARGE</sub>	MobileBERT	MobileBERT <sub>TINY</sub>
embedding	h <sub>embedding</sub>		1024	768	128		
			no-op	no-op	3-convolution		
	h <sub>inter</sub>	1024	768	512			
body	Linear	h <sub>input</sub> h <sub>output</sub>			$\begin{bmatrix} \begin{pmatrix} 512 \\ 1024 \end{pmatrix} \\ \begin{pmatrix} 512 \\ 4 \end{pmatrix} \\ \begin{pmatrix} 1024 \\ 4096 \\ 1024 \end{pmatrix} \\ \begin{pmatrix} 1024 \\ 512 \end{pmatrix} \end{bmatrix} \times 24$	$\begin{bmatrix} \begin{pmatrix} 512 \\ 128 \end{pmatrix} \\ \begin{pmatrix} 512 \\ 4 \end{pmatrix} \\ \begin{pmatrix} 128 \\ 128 \end{pmatrix} \times 4 \\ \begin{pmatrix} 128 \\ 512 \end{pmatrix} \end{bmatrix} \times 24$	$\begin{bmatrix} \begin{pmatrix} 512 \\ 128 \end{pmatrix} \\ \begin{pmatrix} 128 \\ 4 \end{pmatrix} \\ \begin{pmatrix} 128 \\ 128 \end{pmatrix} \times 2 \\ \begin{pmatrix} 128 \\ 512 \end{pmatrix} \end{bmatrix} \times 24$
	MHA	h <sub>input</sub> #Head h <sub>output</sub>	$\begin{bmatrix} \begin{pmatrix} 1024 \\ 16 \\ 1024 \end{pmatrix} \\ \begin{pmatrix} 1024 \\ 4096 \\ 1024 \end{pmatrix} \end{bmatrix} \times 24$	$\begin{bmatrix} \begin{pmatrix} 768 \\ 12 \\ 768 \end{pmatrix} \\ \begin{pmatrix} 768 \\ 3072 \\ 768 \end{pmatrix} \end{bmatrix} \times 12$			
	FFN	h <sub>input</sub> h <sub>FFN</sub> h <sub>output</sub>					
	Linear	h <sub>input</sub> h <sub>output</sub>					
#Params			334M	109M	293M	25.3M	15.1M

# MobileBERT

	#Params	#FLOPS	Latency	CoLA	SST-2	MRPC	STS-B	QQP	MNLI-m/mm	QNLI	RTE	GLUE
				8.5k	67k	3.7k	5.7k	364k	393k	108k	2.5k	
ELMo-BiLSTM-Attn	-	-	-	33.6	90.4	84.4	72.3	63.1	74.1/74.5	79.8	58.9	70.0
OpenAI GPT	109M	-	-	47.2	93.1	87.7	84.8	70.1	80.7/80.6	87.2	69.1	76.9
BERT <sub>BASE</sub>	109M	22.5B	342 ms	<b>52.1</b>	<b>93.5</b>	<b>88.9</b>	<b>85.8</b>	71.2	<b>84.6/83.4</b>	90.5	66.4	78.3
BERT <sub>BASE</sub> -6L-PKD*	66.5M	11.3B	-	-	92.0	85.0	-	70.7	81.5/81.0	89.0	65.5	-
BERT <sub>BASE</sub> -4L-PKD <sup>†</sup> *	52.2M	7.6B	-	24.8	89.4	82.6	79.8	70.2	79.9/79.3	85.1	62.3	-
BERT <sub>BASE</sub> -3L-PKD*	45.3M	5.7B	-	-	87.5	80.7	-	68.1	76.7/76.3	84.7	58.2	-
DistilBERT <sub>BASE</sub> -6L <sup>†</sup>	62.2M	11.3B	-	-	92.0	85.0		70.7	81.5/81.0	89.0	65.5	-
DistilBERT <sub>BASE</sub> -4L <sup>†</sup>	52.2M	7.6B	-	32.8	91.4	82.4	76.1	68.5	78.9/78.0	85.2	54.1	-
TinyBERT*	14.5M	1.2B	-	43.3	92.6	86.4	79.9	<b>71.3</b>	82.5/81.8	87.7	62.9	75.4
MobileBERT <sub>TINY</sub>	15.1M	3.1B	40 ms	46.7	91.7	87.9	80.1	68.9	81.5/81.6	89.5	65.1	75.8
MobileBERT	25.3M	5.7B	62 ms	50.5	92.8	88.8	84.4	70.2	83.3/82.6	90.6	66.2	77.7
MobileBERT w/o OPT	25.3M	5.7B	192 ms	51.1	92.6	88.8	84.8	70.5	84.3/ <b>83.4</b>	<b>91.6</b>	<b>70.4</b>	<b>78.5</b>