

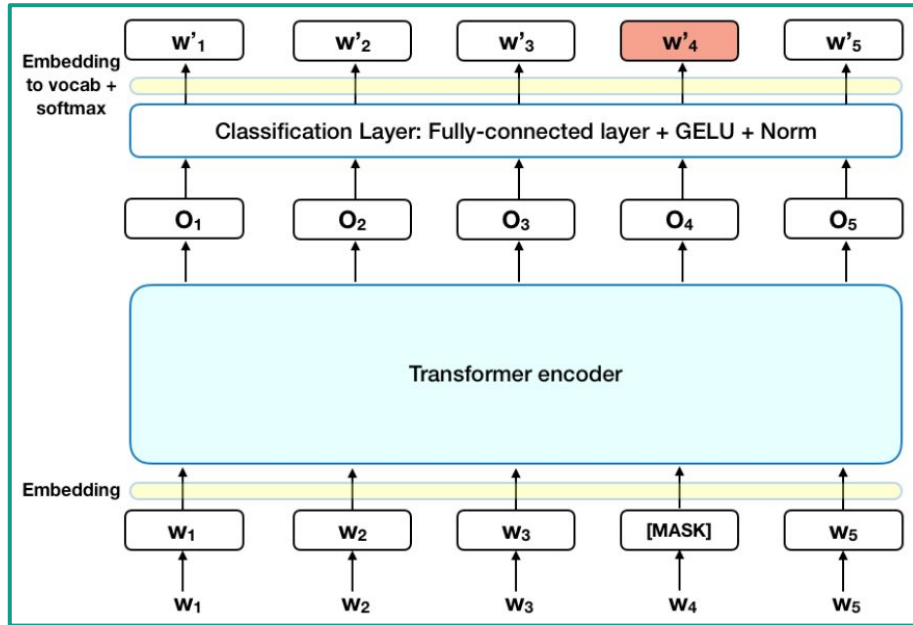
# **Tiny CamemBERT and CamemBERT - A Comparative Study of Two French BERT-Based Models**

Bouchet Elliot, Brun Samuel, Gauci Corentin, Lithaud Alexandre

---

# Tiny BERT and CamemBERT

# The BERT Model



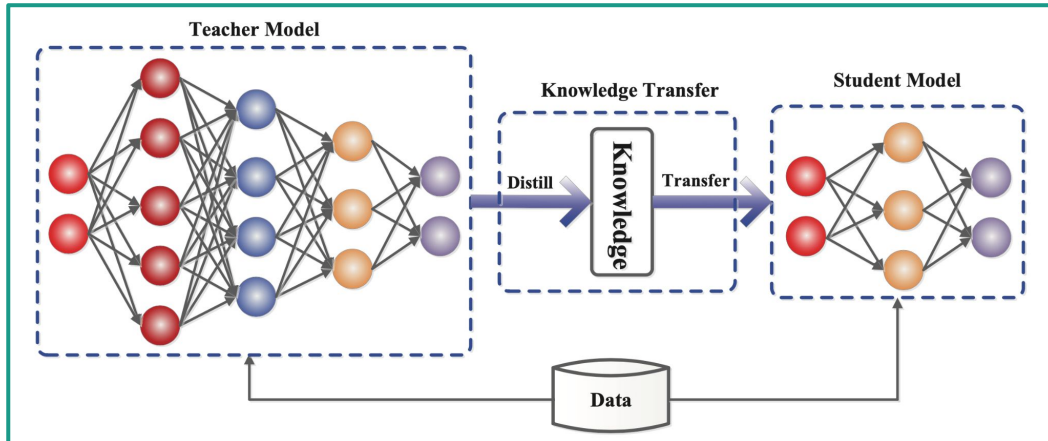
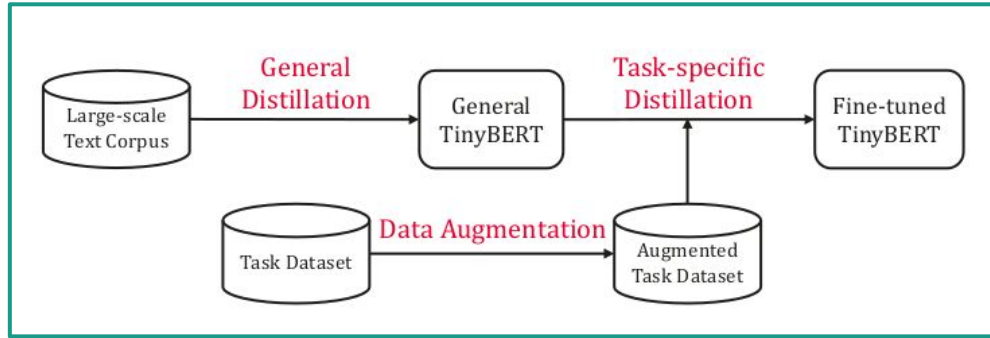
- BERT is pre-trained on large corpora of text data using unsupervised learning tasks like masked language modeling and next sentence prediction
- BERT is built upon the Transformer architecture, which enables it to capture long-range dependencies in text efficiently through self-attention mechanisms.

# The CamemBERT Model

- CamemBERT is a state-of-the-art language model for French based on the RoBERTa architecture
- A model that aims to investigate the feasibility of training monolingual Transformer-based language models for other languages
- Shows that a relatively small web crawled dataset (4GB) leads to results that are as good as those obtained using larger datasets (130+GB).



# The arrangements of TinyBERT

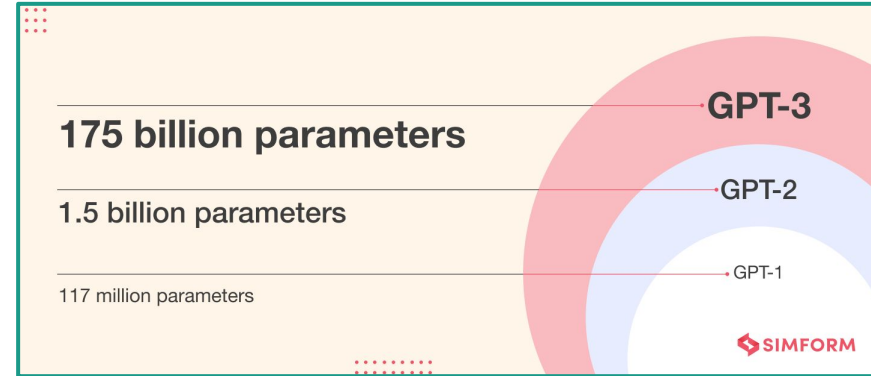
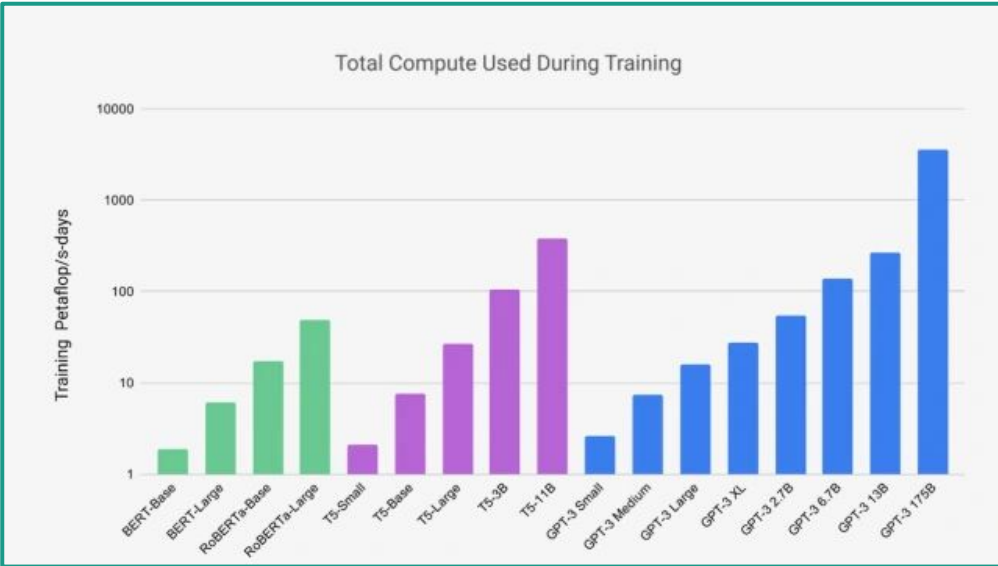


- Use a technique of model compression known as Knowledge Distillation.
- Use a teacher in order to transfer the linguistic knowledge with fewer parameters.
- Faster inference speed and overall steady accuracy.

---

# Context of our project

# The Problem of Newer and Bigger models



**Sam Altman Seeks Trillions of Dollars to Reshape Business of Chips and AI**

---



















# Organisation



# GIT



## Active branches

<b>main</b>  <span>default</span> <span>protected</span>			 
<a href="#">104e064d</a> · Merge remote-tracking branch 'origin/camenBERT' into mergeFinal · 2 minutes ago			
<b>tinyCamemBERT</b> 	3   0	 New	  
<a href="#">78f2cd5f</a> · TinyCamemBERT · 4 minutes ago			
<b>camenBERT</b> 	13   0	 New	  
<a href="#">52a40e6d</a> · Results of accuracy with camemBERT and a fine tuned versions on NLI tasks · 16 hours ago			
<b>tinyBERT</b> 	10   0	 New	  
<a href="#">3196e6eb</a> · TinyBert model, lock issue · 3 days ago			

---

# Experiments with CamemBERT

# Exploration of CamemBERT

For CamemBERT we aimed to reproduce the Natural Language Inference performances described in the paper. We used XNLI which is a subset of a few thousand examples from MNLI.

<b>P<sup>a</sup></b>	A senior is waiting at the window of a restaurant that serves sandwiches.	Relationship
<b>H<sup>b</sup></b>	A person waits to be served his food.	Entailment
	A man is looking to order a grilled cheese sandwich.	Neutral
	A man is waiting in line for the bus.	Contradiction
<sup>a</sup> P, Premise. <sup>b</sup> H, Hypothesis.		

Example of Natural language Inference

Model	Acc.	#Params
mBERT (Devlin et al., 2019)	76.9	175M
XLM <sub>MLM-TLM</sub> (Lample and Conneau, 2019)	<u>80.2</u>	250M
XLM-R <sub>BASE</sub> (Conneau et al., 2019)	80.1	270M
CamemBERT (fine-tuned)	<b>82.5</b>	110M
<i>Supplement: LARGE models</i>		
XLM-R <sub>LARGE</sub> (Conneau et al., 2019)	<u>85.2</u>	550M
CamemBERT <sub>LARGE</sub> (fine-tuned)	<b>85.7</b>	335M

Table showing NLI accuracy on the French XNLI test set

CamemBERT NLI	Accuracy
Before fine-tuning	33.25%
After fine-tuning	81.48%

# CamemBERT performances with sentiment analysis

## French Twitter Sentiment Analysis

1.5 million tweets in French and their sentiment.  
label: Polarity of the tweet (0 = negative, 1 = positive)

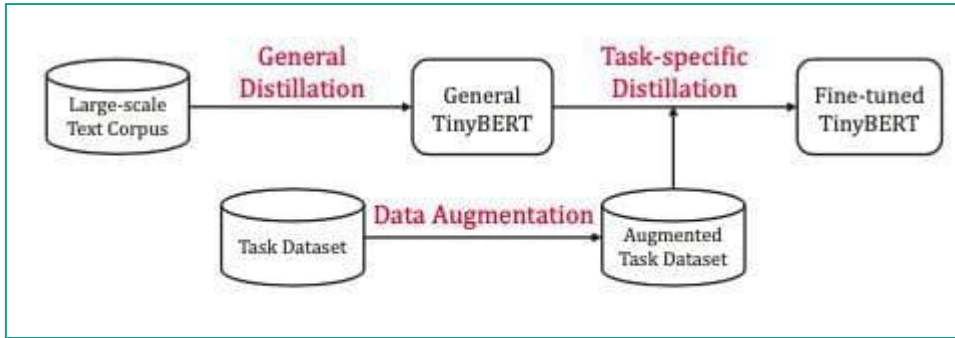
CamemBERT (sentiment analysis)	Loss	Accuracy
Before fine-tuning	0.6937	49.43%
After fine-tuning	0.3702	83.62%
CamemBERT accuracy in the paper concerning NLI	Loss	Accuracy
After fine-tuning	...	82.50%

---

# Experiments with TinyBERT

# Exploration of TinyBERT : Source Code

How to use :



Steps :

- 1) `pregenerate_training_data.py`
- 2) `general_distill.py` -> pre-train
- 3) `data_augmentation.py`
- 4) `task_distill.py` -> train & evaluate

Source :

<https://github.com/huawei-noah/Pretrained-Language-Model/tree/master/TinyBERT>

Advantages :

- Precise control over the tools
- Choice of datasets/parts of datasets
- Architecture strictly following those described in the paper

Disadvantages :

- Heavy and complex code
- Obsolete code (last update 3 years ago)
- Very little documentation
- Software and hardware requirements

# Exploration of TinyBERT : Hugging Face Version

## Advantages :

- Classic hugging face interface
- Easily compatible with notebooks
- Already pre-train

## Disadvantages :

- Reduced architecture control

## Model Name :

“huawei-noah/TinyBERT\_General\_4L\_312D”



# Hugging Face

## Hugging Face Provide :

- Tokenizer (english version)
- Model (pre-train on some data)
- Easy access to databases

---

# Experiments with TinyCamemBERT



# TinyCamemBERT : Our buildings

## Already Pre-trained Version :

- pre-train on english data

## Results :

Parameters : 14 M

Time : 20 minutes

TinyBERT (sentiment analysis)	Accuracy
Before fine-tuning	49.23%
After fine-tuning	79.62%

## Our Pre-trained Version :

- pre-train on few french data  
> same data as CamemBERT

## Results :

pre-train:

1k data, 15 min

Parameters : 14 M

Time : 25 minutes

TinyCamemBERT (sentiment analysis)	Accuracy
Before fine-tuning	50.20%
After fine-tuning	75.62%

---

# CamemBERT Vs TinyCamemBERT

# Comparison CamemBERT & TinyBERT



TinyCamemBERT (sentiment analysis)	Accuracy	CamemBERT (sentiment analysis)	Accuracy
Before fine-tuning	50.20%	Before fine-tuning	49.43%
After fine-tuning	75.62%	After fine-tuning	83.62%

TinyBERT (sentiment analysis)	Accuracy
Before fine-tuning	49.23%
After fine-tuning	79.62%

**Differences :**

- Training time
- Datasets
- Tokenizer

---

# Conclusion

# Future of TinyBERT

Knowledge distillation remains a vibrant area of research in machine learning, with ongoing efforts to develop more effective distillation techniques and apply them to a wide range of tasks and domains

In today's experience of NLP, it is likely that knowledge distillation and model compression techniques will be very important

About 6,160 results (0.38 seconds)

Symbolic **Knowledge Distillation**: from General Language Models to ...

[ACL Anthology > 2022.naacl-main.341](#)



Empirical results demonstrate that, for the first time, a human-authored commonsense **knowledge** graph is surpassed by our automatically d

Lifelong Language **Knowledge Distillation** - ACL Anthology

[ACL Anthology > 2020.emnlp-main.233](#)



To address this issue, we present Lifelong Language **Knowledge Distillation** (L2KD), a simple but efficient method that can be easily applic

Making Monolingual Sentence Embeddings Multilingual using ...

[ACL Anthology > 2020.emnlp-main.365](#)



In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4512–4525, Online. Associa

Sequence-Level **Knowledge Distillation** - ACL Anthology

[ACL Anthology > ...](#)



Yoon Kim and Alexander M. Rush. 2016. Sequence-Level **Knowledge Distillation**. In Proceedings of the 2016 Conference on Empirical Met

# The Other Model Compression Techniques



- Quantization
- Weights Pruning
- Low-rank approximation
- Sparse matrices

Y. Gong, L. Liu, M. Yang, and L. Bourdev. 2014. Compressing deep convolutional networks using vector quantization. arXiv preprint arXiv:1412.6115.

S Han, J. Pool, J. Tran, and W. Dally. 2015. Learning both weights and connections for efficient neural network. In NIPS.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, & Weizhu Chen. (2021). LoRA: Low-Rank Adaptation of Large Language Models.

....

# Bibliography



- Martin, L., Muller, B., Ortiz Suárez, P., Dupont, Y., Romary, L., Clergerie, É., Seddah, D., & Sagot, B. (2020). CamemBERT: a Tasty French Language Model. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, & Qun Liu. (2020). TinyBERT: Distilling BERT for Natural Language Understanding.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, & Dario Amodei. (2020). Language Models are Few-Shot Learners.

	Claude 3 Opus	Claude 3 Sonnet	Claude 3 Haiku	GPT-4	GPT-3.5	Gemini 1.0 Ultra	Gemini 1.0 Pro
Undergraduate level knowledge <i>MLLU</i>	<b>86.8%</b> 5-shot	<b>79.0%</b> 5-shot	<b>75.2%</b> 5-shot	<b>86.4%</b> 5-shot	<b>70.0%</b> 5-shot	<b>83.7%</b> 5-shot	<b>71.8%</b> 5-shot
Graduate level reasoning <i>GPQA, Diamond</i>	<b>50.4%</b> 0-shot CoT	<b>40.4%</b> 0-shot CoT	<b>33.3%</b> 0-shot CoT	<b>35.7%</b> 0-shot CoT	<b>28.1%</b> 0-shot CoT	—	—
Grade school math <i>GSM8K</i>	<b>95.0%</b> 0-shot CoT	<b>92.3%</b> 0-shot CoT	<b>88.9%</b> 0-shot CoT	<b>92.0%</b> 5-shot CoT	<b>57.1%</b> 5-shot	<b>94.4%</b> Maj1@32	<b>86.5%</b> Maj1@32
Math problem-solving <i>MATH</i>	<b>60.1%</b> 0-shot CoT	<b>43.1%</b> 0-shot CoT	<b>38.9%</b> 0-shot CoT	<b>52.9%</b> 4-shot	<b>34.1%</b> 4-shot	<b>53.2%</b> 4-shot	<b>32.6%</b> 4-shot
Multilingual math <i>MGSM</i>	<b>90.7%</b> 0-shot	<b>83.5%</b> 0-shot	<b>75.1%</b> 0-shot	<b>74.5%</b> 8-shot	—	<b>79.0%</b> 8-shot	<b>63.5%</b> 8-shot
Code <i>HumanEval</i>	<b>84.9%</b> 0-shot	<b>73.0%</b> 0-shot	<b>75.9%</b> 0-shot	<b>67.0%</b> 0-shot	<b>48.1%</b> 0-shot	<b>74.4%</b> 0-shot	<b>67.7%</b> 0-shot
Reasoning over text <i>DROP, F1 score</i>	<b>83.1</b> 3-shot	<b>78.9</b> 3-shot	<b>78.4</b> 3-shot	<b>80.9</b> 3-shot	<b>64.1</b> 3-shot	<b>82.4</b> Variable shots	<b>74.1</b> Variable shots
Mixed evaluations <i>BIG-Bench-Hard</i>	<b>86.8%</b> 3-shot CoT	<b>82.9%</b> 3-shot CoT	<b>73.7%</b> 3-shot CoT	<b>83.1%</b> 3-shot CoT	<b>66.6%</b> 3-shot CoT	<b>83.6%</b> 3-shot CoT	<b>75.0%</b> 3-shot CoT
Knowledge Q&A <i>ARC-Challenge</i>	<b>96.4%</b> 25-shot	<b>93.2%</b> 25-shot	<b>89.2%</b> 25-shot	<b>96.3%</b> 25-shot	<b>85.2%</b> 25-shot	—	—
Common Knowledge <i>HellaSwag</i>	<b>95.4%</b> 10-shot	<b>89.0%</b> 10-shot	<b>85.9%</b> 10-shot	<b>95.3%</b> 10-shot	<b>85.5%</b> 10-shot	<b>87.8%</b> 10-shot	<b>84.7%</b> 10-shot

Credit: Anthropic