



# Few-shot learning in NLP

## SetFit Model

Moshe Wasserblat  
NLP/EAI/Intel Labs

April 2023

# BIO

- **NICE Systems**



- Led NLP research group
- First company to productize Speech2Text, ED, Text Analytics in Call-Center

- **INTEL**

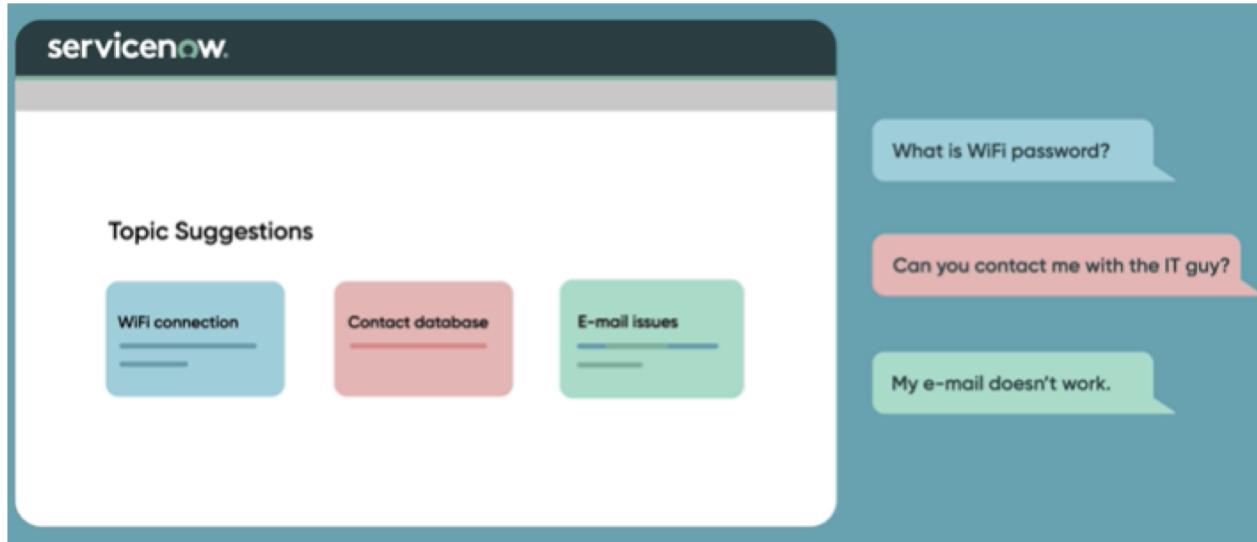


- Applied NLP/DL Research – Increase adoption of DL in the industry
- SoTA NLP research
- Explore compute features that disrupt our HW

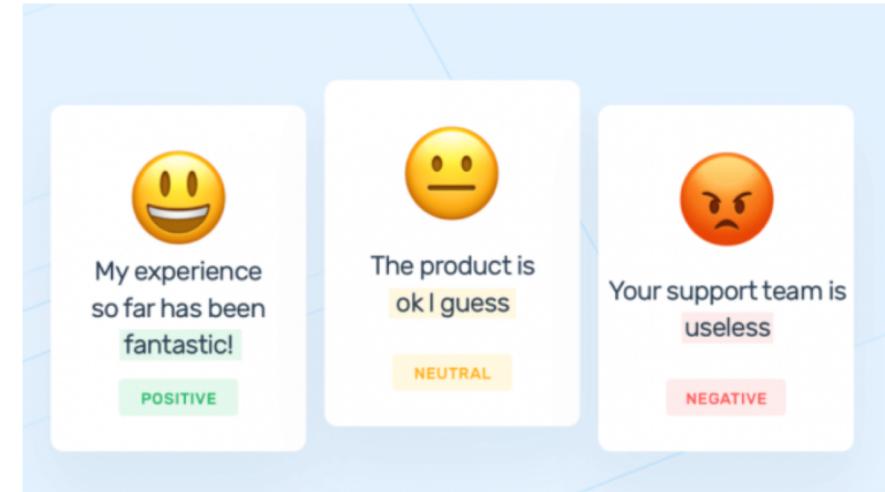
# Text Classification - Usages

48% of Large companies use NLP for document classification

## 1. Topic classification for routing



## 2. Sentiment Analysis



3. Spam detection
4. Hate speech detection
5. Etc.

# Text Classification in Healthcare

Label the following sentence based on whether it is related to an Adverse Drug Effect (\*ADE)



**{ Sentence:** ‘No regional side effects were noted’  
**Label:** ‘not ADE-related’ }

Label the following publication **abstracts** based on whether it is related to one or many of **HoC** 10 labels.



**{ Abstract:** ‘ZD1839 increased the expression of p27(KIP1) ... and maturation markers ( P < .001 ) ’  
**Label:** ‘evading growth suppressors, resisting cell death’ }

**HoC (Hallmarks of Cancer)**

# Agenda



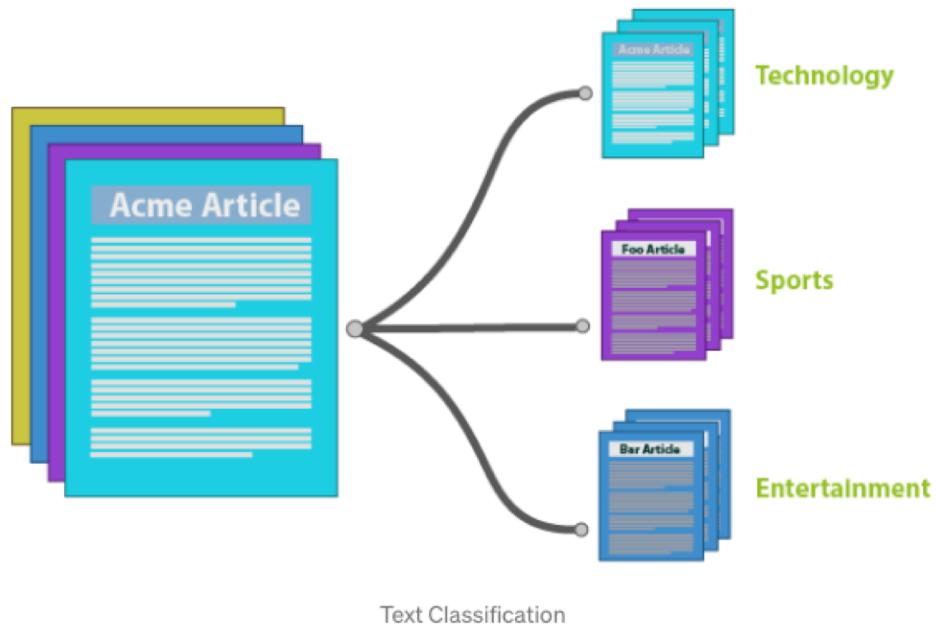
Few-shot learning overview



How does SetFit work?

# Few-shot Learning with Transformers

- Traditional fine-tuning -> **1000s** labeled training sample
- Fewshot learning -> **8-64** labeled training sample



# Method 1 - In-context Learning

## Language Models are Few-Shot Learners

### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

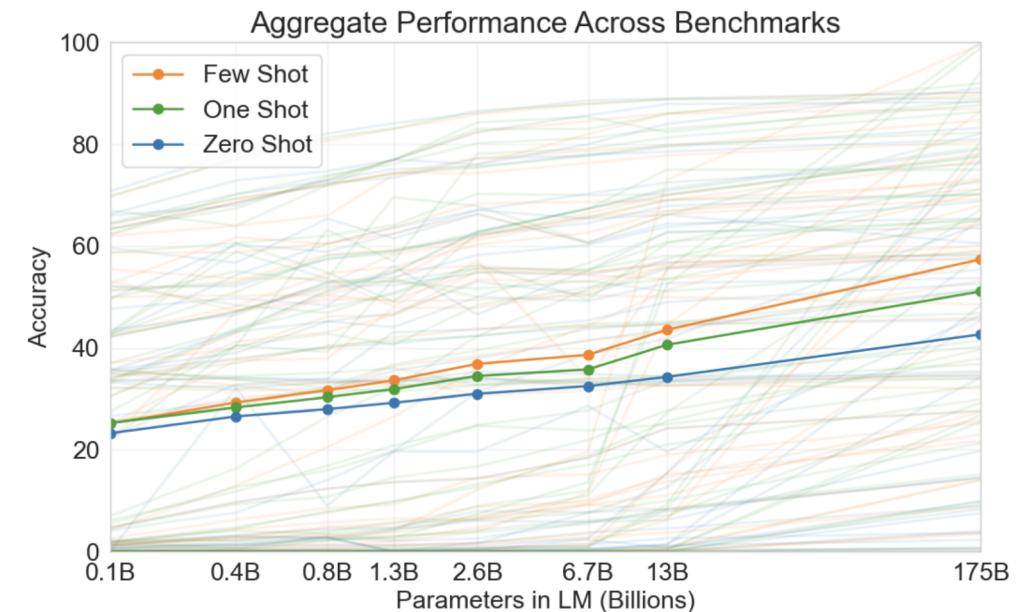
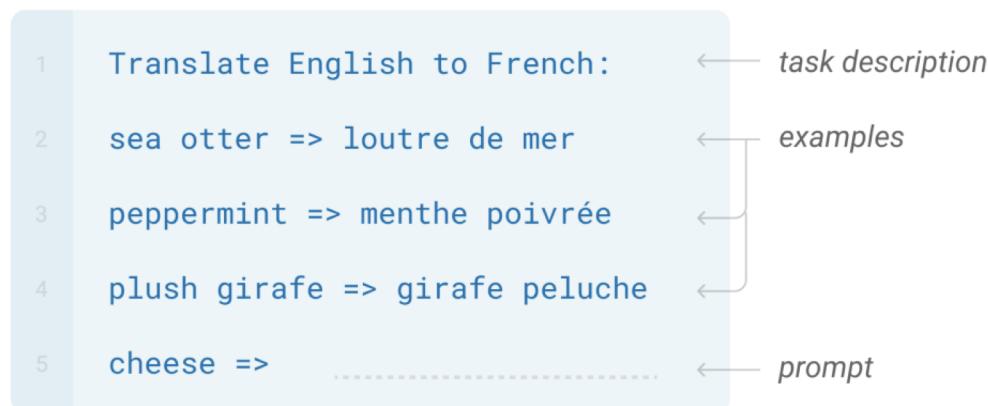


Image source: [Language Models are Few-Shot Learners](#)

- Popularised by GPT-3 with “in context learning”
- Performance increases with scale
- In practice:
  - performance sensitive to quality of “prompt engineering”
  - hard/expensive to deploy large models

# Method 2 - Parameter Efficient Fine-tuning

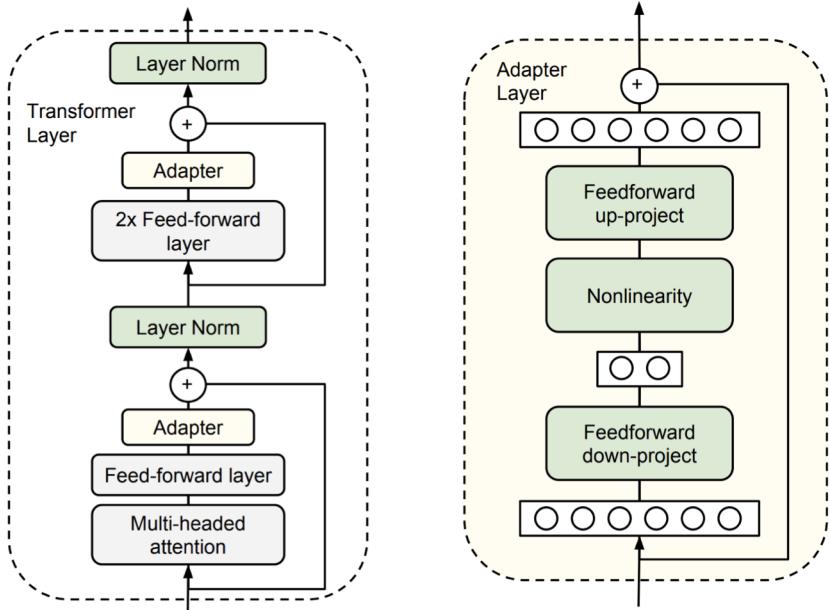


Image credit: [Parameter-Efficient Transfer Learning for NLP](#)

T-Few much more  
efficient / performant than GPT-3  
💪

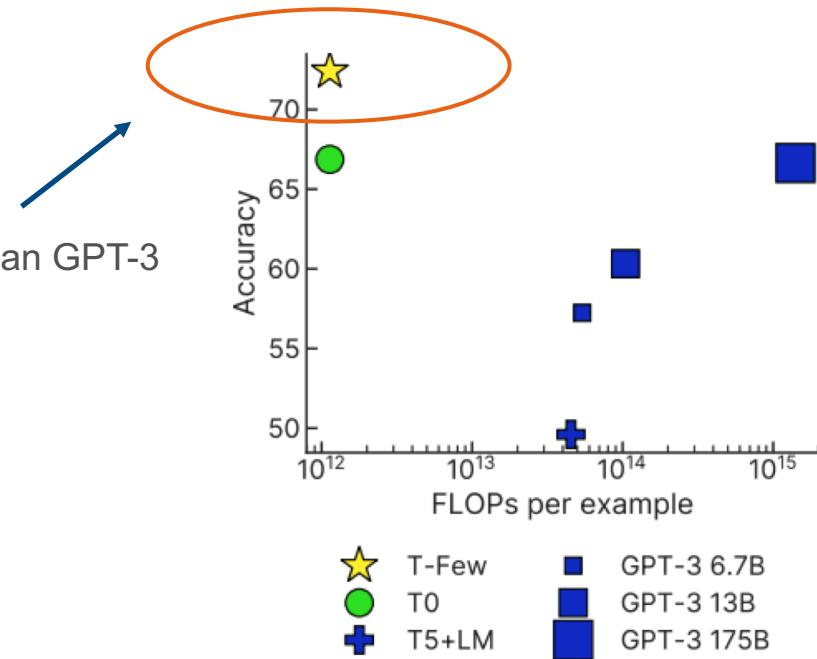


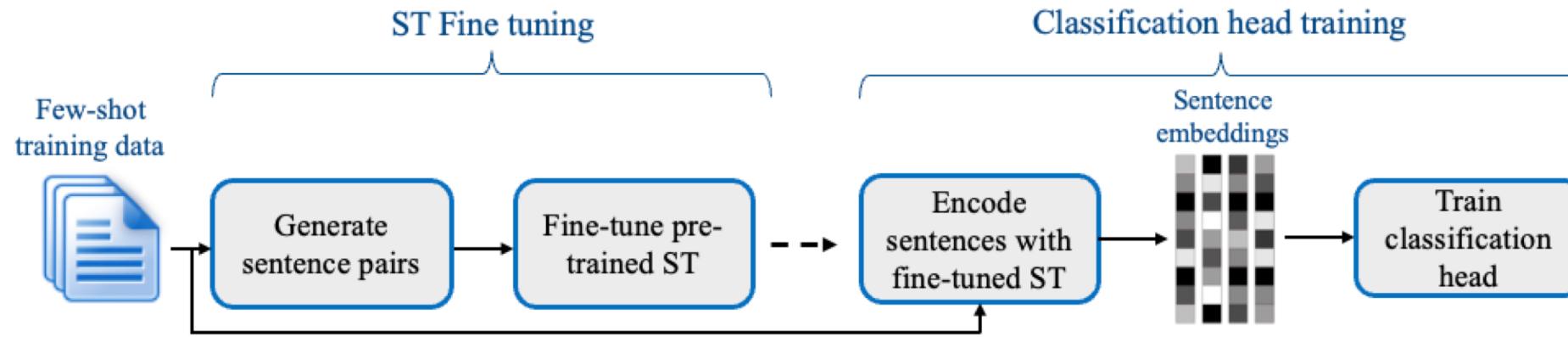
Image credit: [T-Few paper](#)

- Alternatives include “**parameter efficient fine-tuning**” (PEFT) which add/update a small number of parameters in a pretrained LM. **T-Few is current SOTA**
- In practice:
  - performance sensitive to quality of “**prompt engineering**” (like ICL)
  - SOTA performance relies on largish 11B (11.4GB) T0 model

Can we do better? 🤔

- No depending on large generative models
- Not depending on manually crafted prompts

# SetFit – Sentence Transformer Fine-Tuning



SetFit uses a simple, yet effective **2-stage training recipe**:

- Fine-tune a Sentence Transformer “body” with **contrastive learning** on small number of labelled examples (e.g. 8 examples per class)
- Train classification “head” using embeddings from fine-tuned Sentence Transformer

Works well for **text classification** (token classification WIP)

# Contrastive learning primer

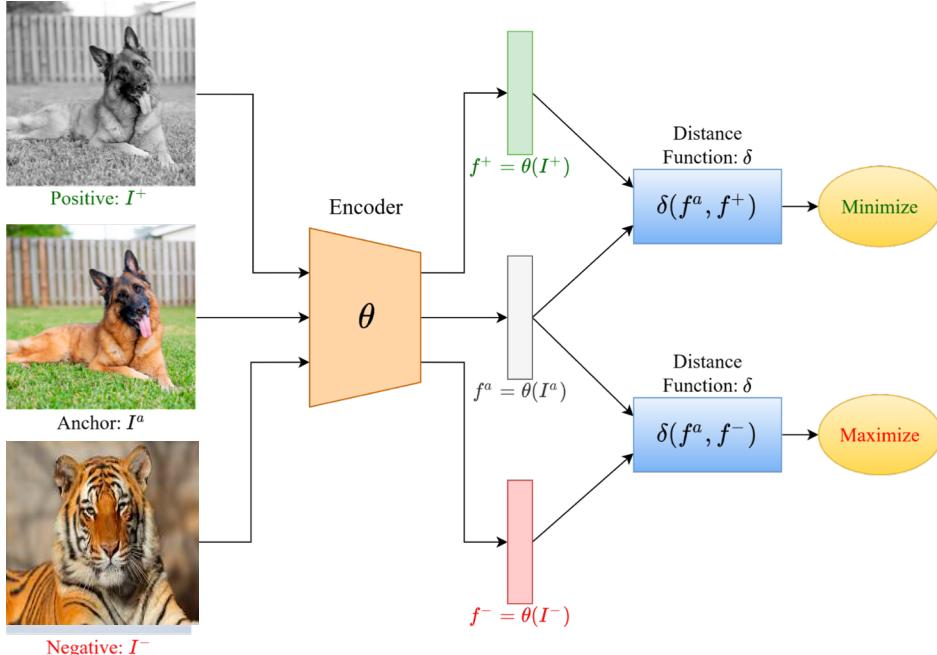
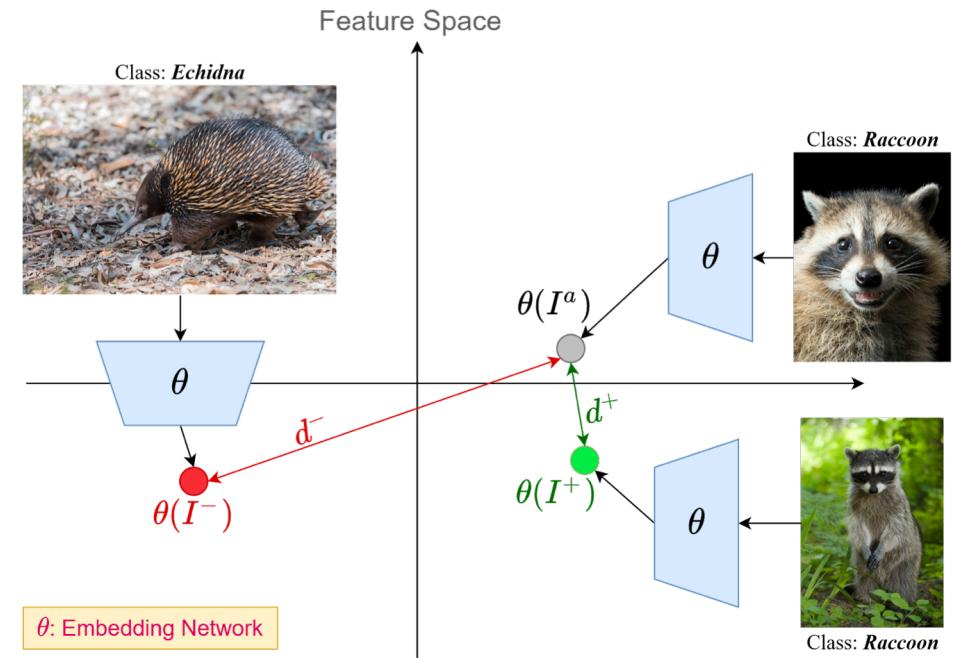


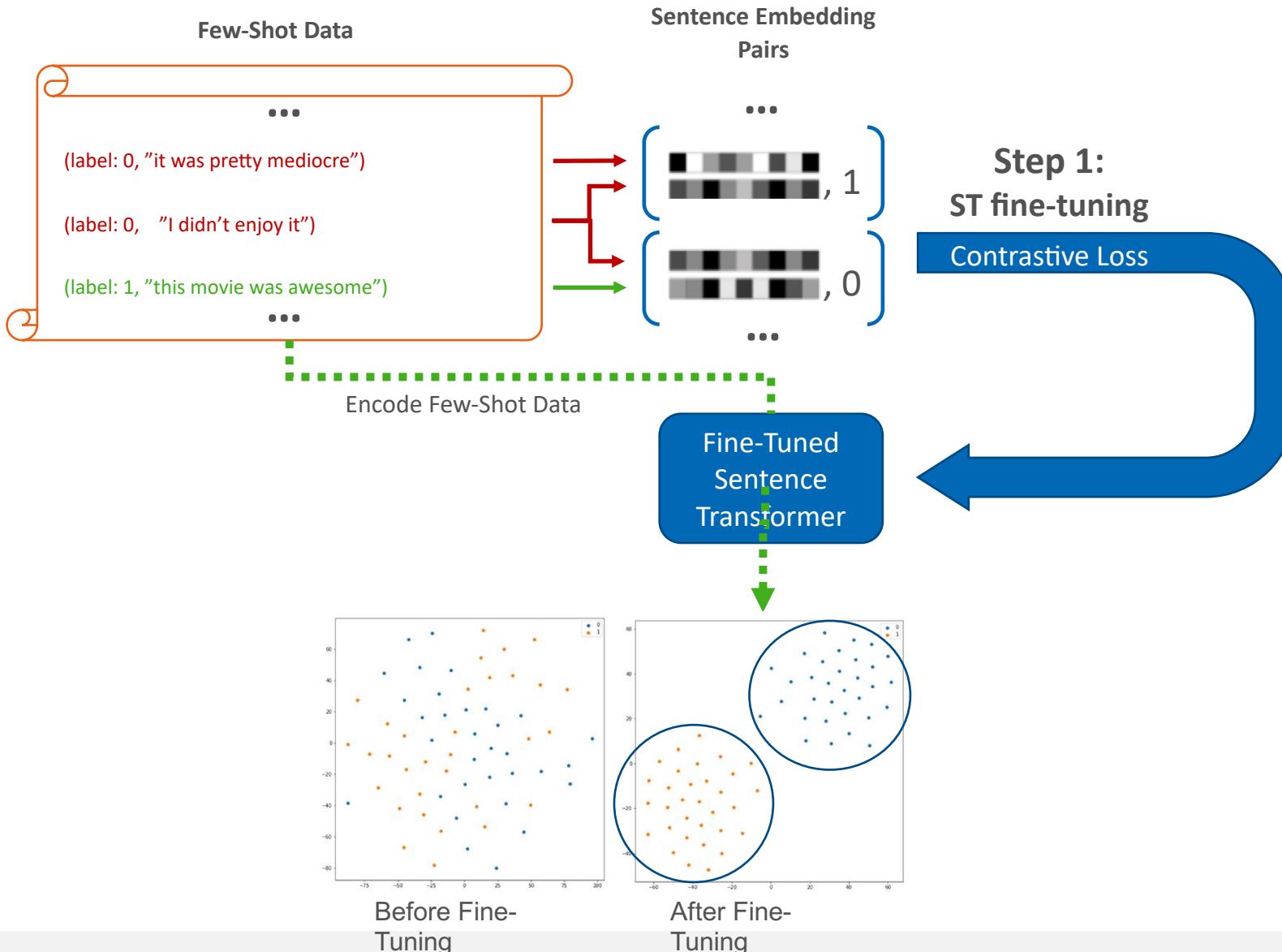
Image source: [V7 Blog](#)

Basic idea:

- Minimize (maximize) distance between “positive” (“negative”) examples
- Distance metric often Euclidean distance or cosine similarity
- Produces high quality representations for images / text

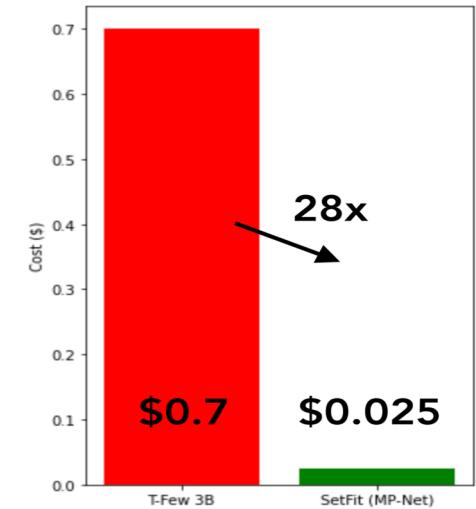


# SetFit - Training



# SetFit Results

Method	SST-5	AmazonCF	CR	Emotion	EnronSpam	AGNews	Average <sup>†</sup>
$ N  = 8^*$							
FINETUNE	33.5 <sub>2.1</sub>	9.2 <sub>4.9</sub>	58.8 <sub>6.3</sub>	28.7 <sub>6.8</sub>	85.0 <sub>6.0</sub>	81.7 <sub>3.8</sub>	43.0 <sub>5.2</sub>
PERFECT	34.9 <sub>3.1</sub>	18.1 <sub>5.3</sub>	81.5 <sub>8.6</sub>	29.8 <sub>5.7</sub>	79.3 <sub>7.4</sub>	80.8 <sub>5.0</sub>	48.7 <sub>6.0</sub>
ADAPET	50.0 <sub>1.9</sub>	19.4 <sub>7.3</sub>	91.0 <sub>1.3</sub>	46.2 <sub>3.7</sub>	85.1 <sub>3.7</sub>	85.1 <sub>2.7</sub>	58.3 <sub>3.6</sub>
T-FEW 3B	<b>55.0</b> <sub>1.4</sub>	19.0 <sub>3.9</sub>	<b>92.1</b> <sub>1.0</sub>	<b>57.4</b> <sub>1.8</sub>	<b>93.1</b> <sub>1.6</sub>	—	<b>63.4</b> <sub>1.9</sub>
SETFIT <sub>MPNET</sub>	43.6 <sub>3.0</sub>	<b>40.3</b> <sub>11.8</sub>	88.5 <sub>1.9</sub>	48.8 <sub>4.5</sub>	90.1 <sub>3.4</sub>	82.9 <sub>2.8</sub>	62.3 <sub>4.9</sub>
$ N  = 64^*$							
FINETUNE	45.9 <sub>6.9</sub>	52.8 <sub>12.1</sub>	88.9 <sub>1.9</sub>	65.0 <sub>17.2</sub>	95.9 <sub>0.8</sub>	88.4 <sub>0.9</sub>	69.7 <sub>7.8</sub>
PERFECT	49.1 <sub>0.7</sub>	<b>65.1</b> <sub>5.2</sub>	92.2 <sub>0.5</sub>	61.7 <sub>2.7</sub>	95.4 <sub>1.1</sub>	<b>89.0</b> <sub>0.3</sub>	72.7 <sub>1.9</sub>
ADAPET	54.1 <sub>0.8</sub>	54.1 <sub>6.4</sub>	92.6 <sub>0.7</sub>	72.0 <sub>2.2</sub>	96.0 <sub>0.9</sub>	88.0 <sub>0.6</sub>	73.8 <sub>2.2</sub>
T-FEW 3B	<b>56.0</b> <sub>0.6</sub>	34.7 <sub>4.5</sub>	<b>93.1</b> <sub>1.0</sub>	70.9 <sub>1.1</sub>	<b>97.0</b> <sub>0.3</sub>	—	70.3 <sub>1.5</sub>
SETFIT <sub>MPNET</sub>	51.9 <sub>0.6</sub>	61.9 <sub>2.9</sub>	90.4 <sub>0.6</sub>	<b>76.2</b> <sub>1.3</sub>	96.1 <sub>0.8</sub>	88.0 <sub>0.7</sub>	<b>75.3</b> <sub>1.3</sub>
$ N  = Full^{**}$							
FINETUNE	59.8	80.1	92.4	92.6	99.0	93.8	84.8



- SetFit performance comparable to other SOTA techniques like T-Few (3B) and ADAPET
- 1-2 orders of magnitude faster to train / run inference 🚗!

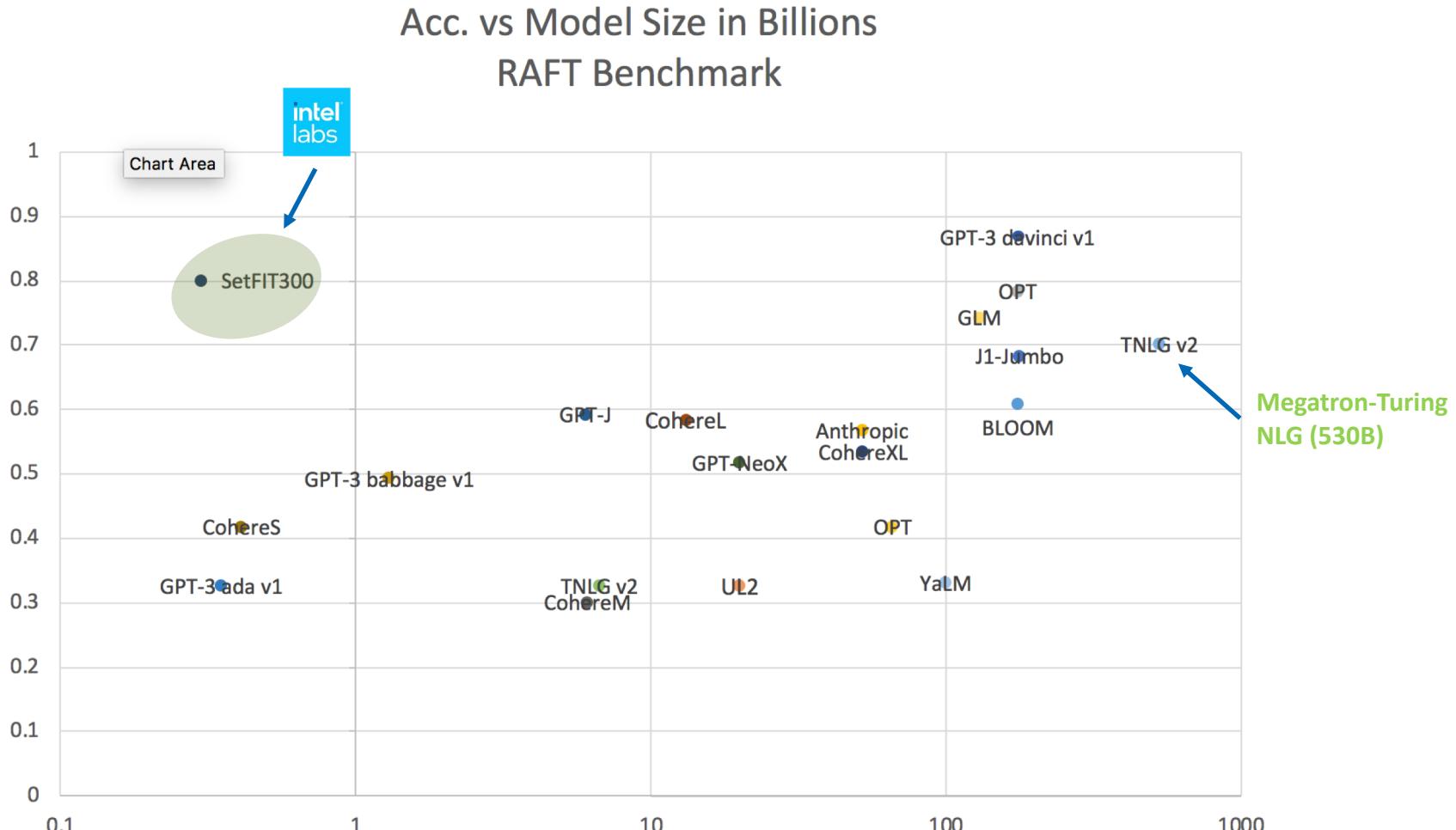
# RAFT Benchmark

Rank	Method	Score	Size*
1	YIWISE	76.8	-
2	T-FEW 11B	75.8	11B
4	Human baseline	73.5	-
350M	6 SETFIT <sub>ROBERTA</sub>	71.3	355M
	9 PET	69.6	235M
110M	11 SETFIT <sub>MPNET</sub>	66.9	110M
	12 GPT-3	62.7	175B

Table 3: SetFit compared to prominent methods on the RAFT leaderboard (as of Sept. 5, 2022). \*Number of parameters.

- RAFT is a real-world challenging few-shot benchmark
- SetFit outperforms GPT-3, while being 1600x smaller and not using prompts
- SetFit surpasses the Human baseline in 7 out of 11 tasks

# RAFT Benchmark – Healthcare use-case (EDA)



Stanford – Center for Research on Foundation Models  
[Large Models Benchmark](#)

# 🎉 SetFit SOTA on Bio Text Classification (full-shot) 🎉

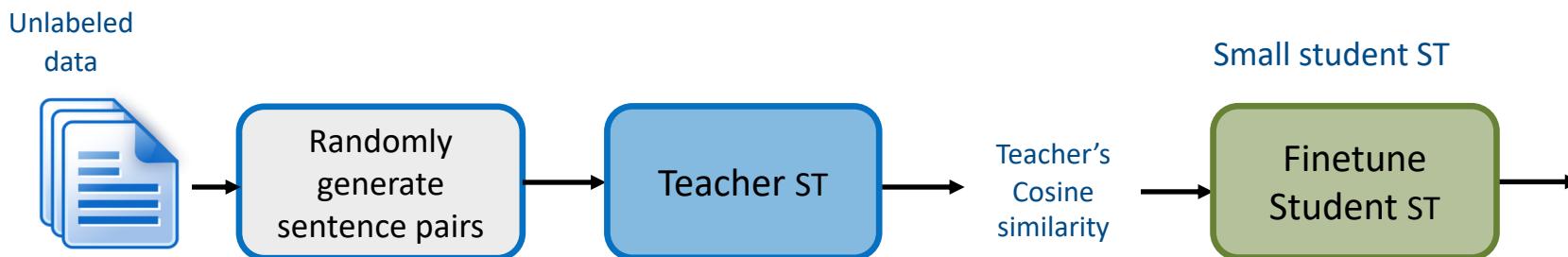
- 💪 Outperform models that were trained from scratch on Bio data
- 💾 3x times smaller
- 🏍️ Competitive with BioGPT

Model	#params[M]	F1	Pre-train Data
<b>BioBERT[1]</b>	110	81.5	Bio
<b>PubMedBERT[2]</b>	340	82.7	Bio
<b>BioLinkBERT[3]</b>	340	84.9	Bio
<b>GPT-2</b>	355	81.8	General
<b>BioGPT[4]</b>	347	85.1	Bio
<b>SetFit</b>	105	<b>85.1</b>	General



**HOC (Hallmarks of Cancer)**

# Experiments - Knowledge Distillation



Sentece1	Sentence2	Teacher's Cosine Sim
POS	POS	1
NEG	NEG	0.81
POS	POS	0.95
POS	NEG	0.05
POS	POS	0.91
NEG	POS	0.1
NEG	NEG	1
POS	NEG	0.02

Basic idea:

- Train teacher and student Sentence Transformer models with 16 examples / class
- Student learn to mimic teacher cosine-similarity
- \*Use student embeddings + teacher logits to train SetFit model

# Experiments - Knowledge Distillation

Method	Inf. FLOPs	Train FLOPs	Speed- up	Score
T-FEW 3B	1.6e11	3.9e15	1x	63.4 <sub>1.9</sub>
110M SETFIT <sub>MPNET</sub>	8.3e9	2.0e14	19x	62.3 <sub>4.9</sub>
15M SETFIT <sub>MINILM</sub> <sup>†</sup>	1.3e9	3.2e13	123x	60.3 <sub>1.6</sub>

Experiment details:

- 16 samples per class for teacher training + unlabeled data for distillation
- The score shown is an average across 5 datasets
- Setfit Student falls short in only 2 points from Setfit teacher

# SetFit: Few-shot Text Classification

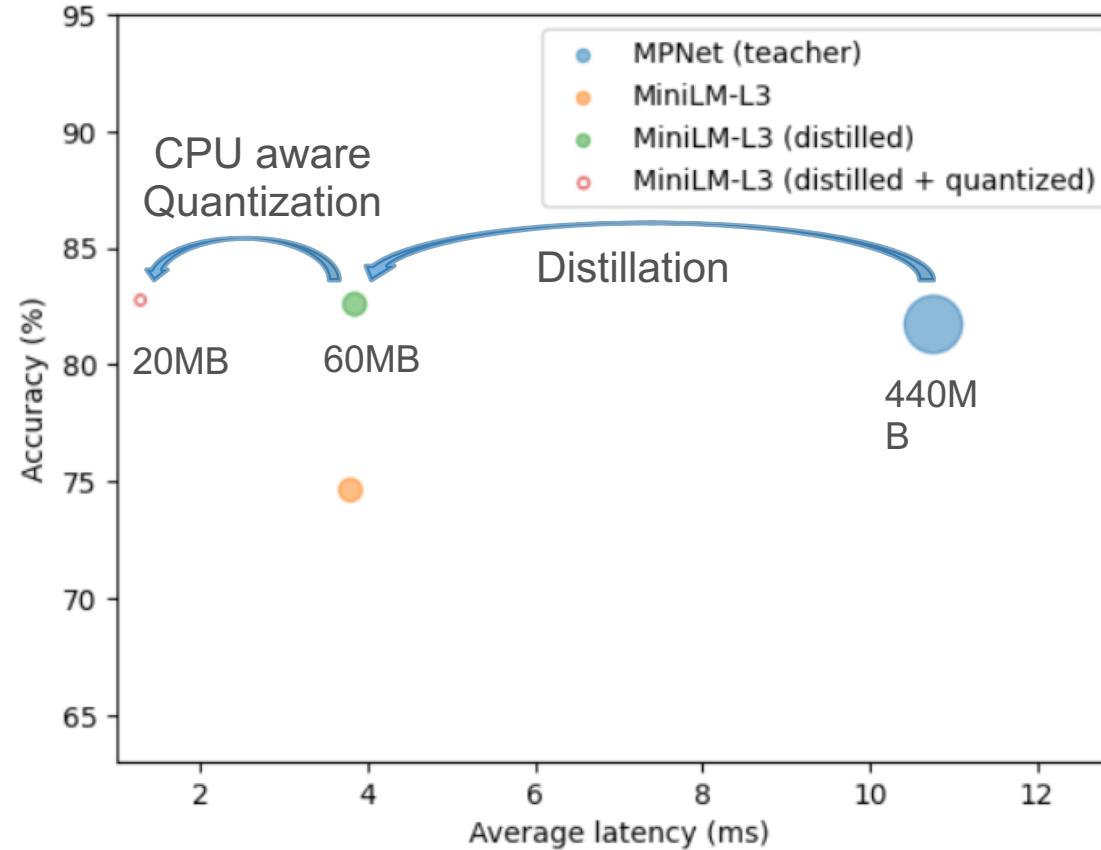
## Efficiency steps:

- Distillation
- CPU aware quantization

## Results:

- 22x model size reduction
- 10x latency reduction

CPU Efficiency



# Efficient Few-Shot Learning Without Prompts

Lewis Tunstall<sup>1</sup>, Nils Reimers<sup>2</sup>, Unso Eun Seo Jo<sup>1</sup>, Luke Bates<sup>3</sup>,  
Daniel Korat<sup>4</sup>, Moshe Wasserblat<sup>4</sup>, Oren Pereg<sup>4</sup>

<sup>1</sup>Hugging Face    <sup>2</sup>cohere.ai

<sup>3</sup>Ubiquitous Knowledge Processing Lab, Technical University of Darmstadt

<sup>4</sup>Emergent AI Lab, Intel Labs



+



+



UBIQUITOUS  
KNOWLEDGE  
PROCESSING

NLP/EAI Lab



Paper:

<https://arxiv.org/abs/2209.11055>



Blog:

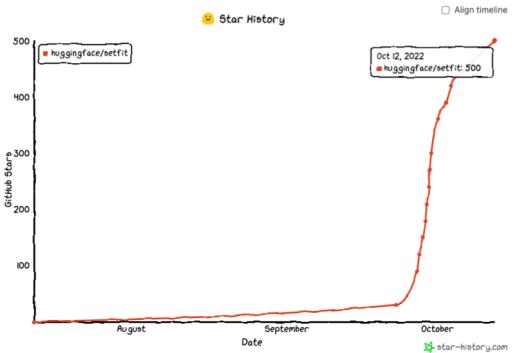
<https://huggingface.co/blog/setfit>



Code:

<https://github.com/huggingface/setfit>

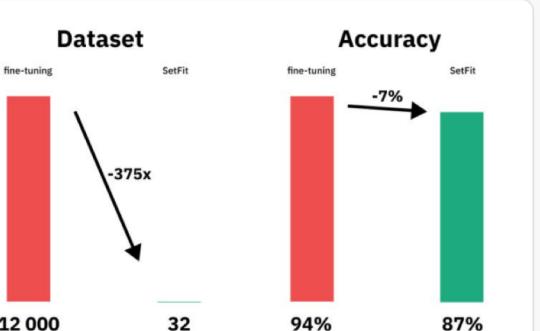
# Shout-out



150000+ downloads

[linkedin.com](#)  
Philipp Schmid on LinkedIn: Outperform OpenAI GPT-3 with SetFit for text-classification

OpenAI GPT-3 is one of the most known Large Language Models (LLM) today. Those LLMs are achieving good performance in classification, where you only have a few... (18 kB) ▾



*"This is huge! SetFit will help so many companies to get started with transformers, without the need to label a lot of data and compute power."*

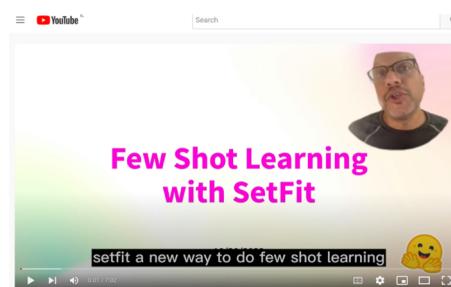
*"Now anyone can build a few shot classifier like a pro on only a handful of labelled data"*

R rubrix.readthedocs.io

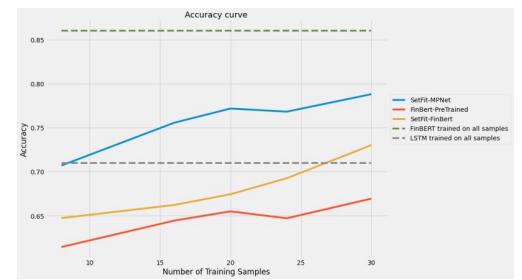
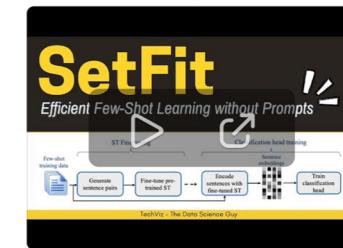
Few-shot classification with SetFit and a custom dataset

SetFit is an exciting open-source package for few-shot classification developed by teams at Hugging Face and Intel Labs. You can read all about it on the project repository. To showcase how powerfu... (485 kB) ▾

RUBRIX  
The open-source tool for data-centric NLP



YouTube TechViz - The Data Science Guy  
SetFit - Efficient Few-Shot Learning Without Prompts (Research Paper Walkthrough) ▾



YouTube Rohan-Paul-AI  
SetFit (Sentence Transformer Fine-tuning) - Paper Discussion with Code | NLP | Machine Learning ▾



# SetFit

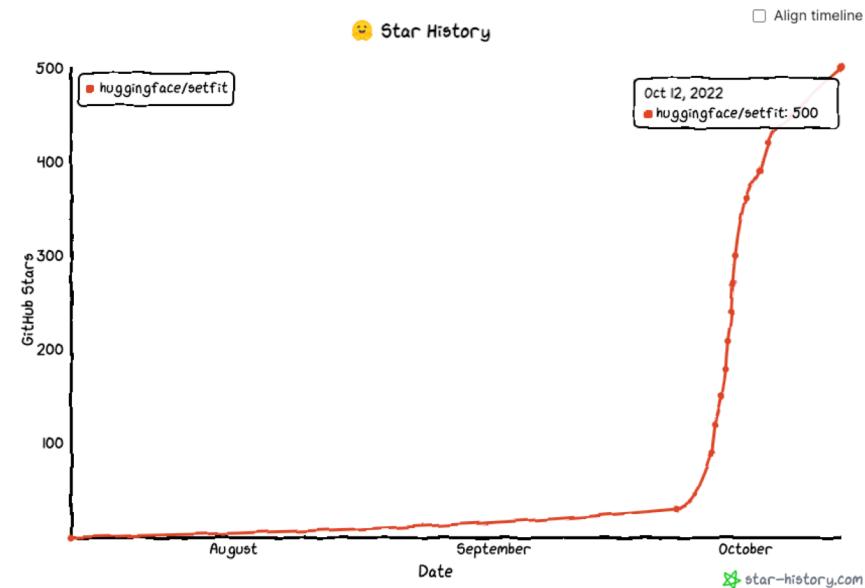


*150000+ installs*

*Outperform GPT-3, 1600x smaller*

*~20k faster inference*

*2min fine-tuning on SPR*



*Service providers that offer SetFit*



**Argilla**



**mantis**



**witty.works**



April 2023