

# THE VALUE OF DATA FOR THE DEVELOPMENT OF TOP QUALITY LT

[ADAM FELDMANN UNIVERSITY OF PÉCS]



## WHY WE NEED MORE DATA FOR LANGUAGE MODELS?

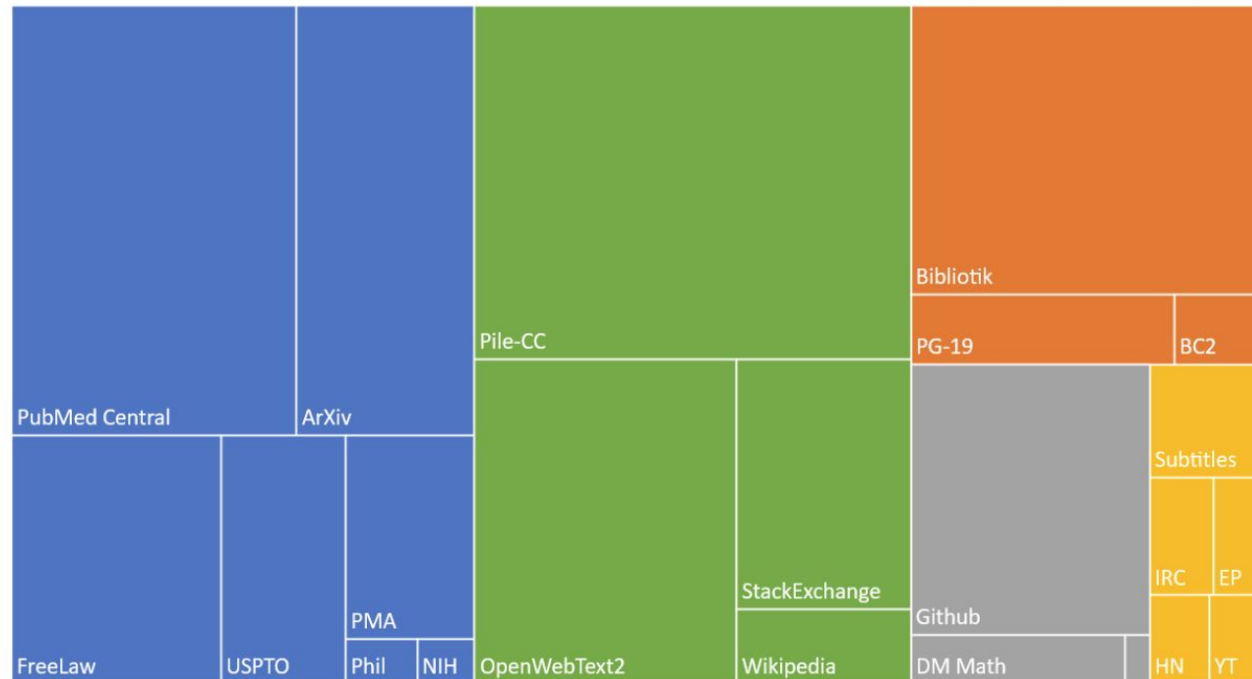
- Language models are the engine of language technology
- Era of large language models
- S.O.T.A solutions need more data  
(aka. Pains and promises of BERT,GPT-3, etc.)

# HOW MUCH DATA IS ENOUGH DATA?

GPT-3, Wu Dao 2.0, Megatron-LM 530B

Composition of the Pile by Category

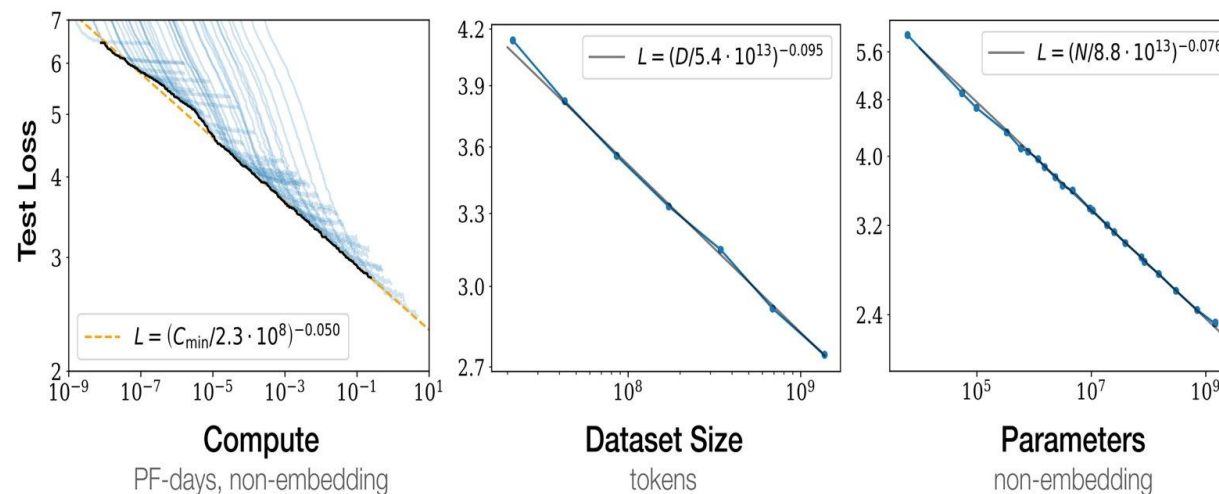
■ Academic ■ Internet ■ Prose ■ Dialogue ■ Misc



(Source: [\[2101.00027\] The Pile: An 800GB Dataset of Diverse Text for Language Modeling \(arxiv.org\)](#))

# HOW MUCH DATA IS ENOUGH DATA?

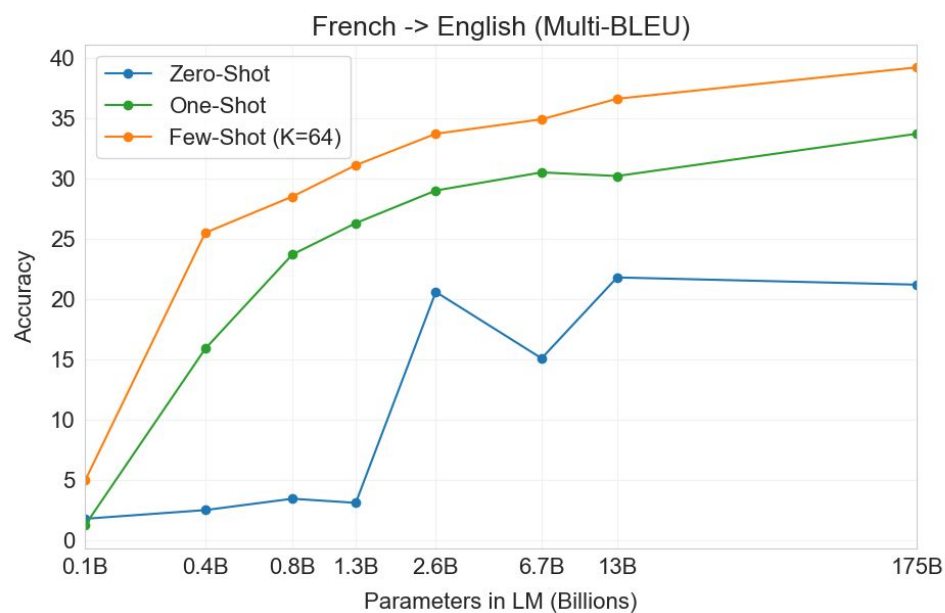
- Hard to tell....but it is calculable now based on empirical researches.



**Figure 1** Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute<sup>2</sup> used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

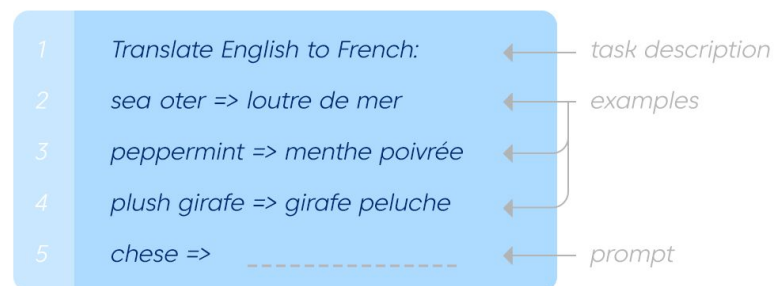
(Source: [2001.08361.pdf \(arxiv.org\)](#) Jared Kaplan et al; 2020; Scaling Laws For Neural Language Models)

# MACHINE TRANSLATION (MT) USING PROMPT PROGRAMMING PARADIGM



## Few-shot

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



# AN UNEXPECTED BENEFICIAL PHENOMENON: PROMPT PROGRAMMING

- New paradigm in AI.
- A new method for using multitask metalearner models like GPT.
- Originated from zero-shot or few-shots learning.
- Prompting allows the user to generate a very specific behavior from the system.

# THANK YOU FOR YOUR ATTENTION!

Website: [www.lr-coordination.eu](http://www.lr-coordination.eu)

Twitter: @LR\_Coordination

Facebook: [www.facebook.com/EuropeanLanguageResourceCoordination](https://www.facebook.com/EuropeanLanguageResourceCoordination)

LinkedIn: [www.linkedin.com/in/lrcoordination](https://www.linkedin.com/in/lrcoordination)

Email: [info@lr-coordination.eu](mailto:info@lr-coordination.eu)

