# Lesson 3: Pretraining and finetuning for building LLM

1 token ≈ 1 word Lets say.        **Pretraining (First Stage)**

Trained on 410 billion tokens from CommonCrawl ⟶ Different weights.

webtext2 → 19 billion

Books → 67 billion

wiki → 3 billion

LLM perform well b/c trained on Huge amounts of data
is overcomplction

The icon is in the _ferret_; trained on this type of task (Translating)
is can do a whole range of tasks as well.

Can begin to perform tasks w/out ever training on them. (MC Summarize, etc.)

Pretraining → Train LLM on huge data so can do wide range of tasks

## Finetuning (Second Stage)

- CBO online; you want chatbot so users can ask questions maybe?    Fraud detection
  - you want specific responses, not generic.
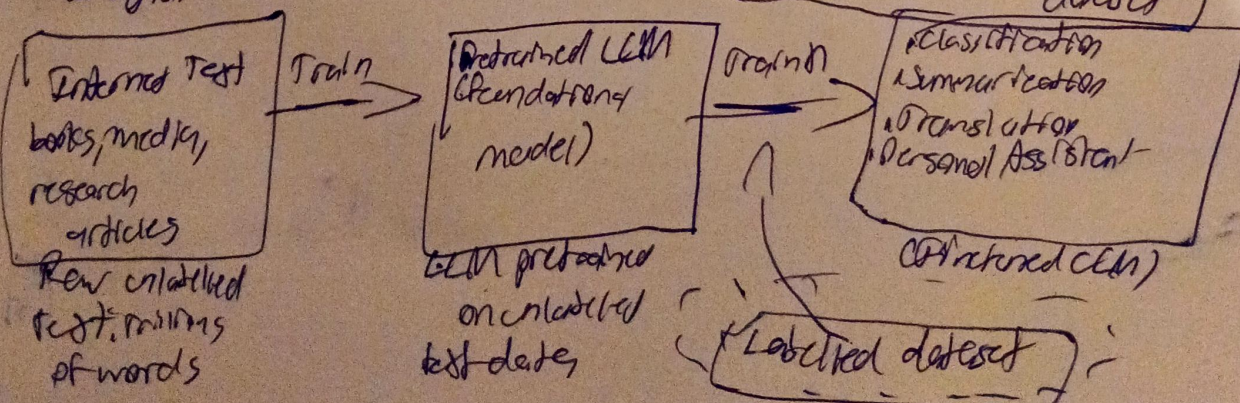  - want be specific to company or application
- High-quality is why ppl finetune; refine on est on smaller dataset; train model on own data.
- OG training data may lack extensive knowledge of certain topics, like legal information, etc. $5million cost for GPT 3 pretraining
- Otherw/ai is for legal teams

Pretraining LLM on labelled dataset

| Internet Text<br>books, media,<br>research<br>articles | Train → | Pretrained LLM<br>(foundational<br>model) | Train → | • classification<br>• Summarization<br>• Translation<br>• Personal Assistant |
|---|---|---|---|---|
| Raw unlabelled text millions of words | | LLM pretrained on unlabelled text data | | (Finetuned LLM) |

Labelled dataset

# Steps for building LLM

1. Train on large corpus of text data (Raw text)
   - Regular text w/ out any labelling info

2. First stage of LLM (pretraining)
   - GPT-3 model is pretrained model which is capable of text completion

3. After obtaining pretrained LLM, we can further train LLM on labelled data - finetuning

4. Two popular finetuning categories

   Instruction finetuning          Classification tasks finetuning

   Labelled dataset consists of     Labelled dataset
   Instruction-answer pairs         consists of text
                                    and associated
   - text translation,              labels,
   online customer                   - emails → spam or no spam
   support

---

## Basic Intro to Transformers

---

- Deep neural network arch. most modern LLMs use
  - 2017 paper introduced
  - Attention is all you need
- originally proposed for english → german, english → french,
- Decoder