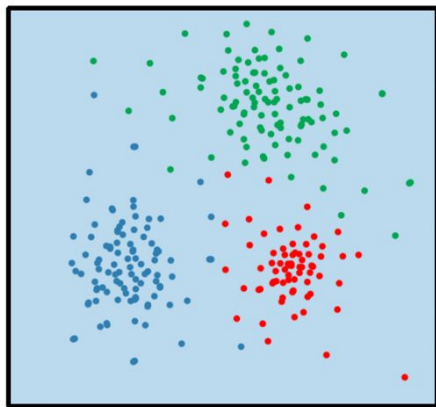


LLMs: Large Language Models

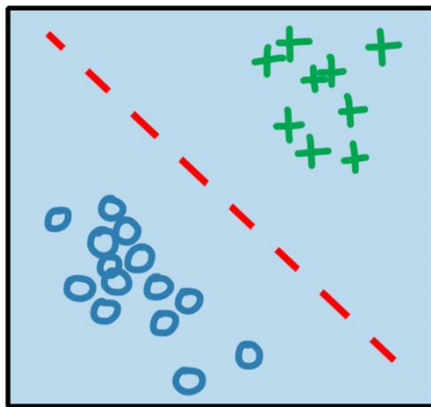
Supervised vs Unsupervised learning

machine learning

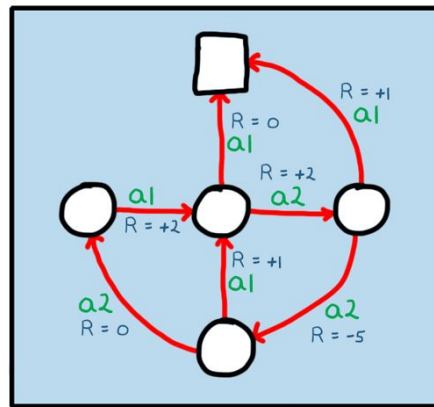
unsupervised
learning



supervised
learning



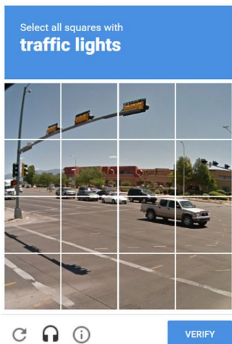
reinforcement
learning



Supervised vs Unsupervised learning

Supervised learning:

- Known labels
- Examples:
 - credit, fraud,
 - image labels



- Labels are expensive

Unsupervised learning:

- No labels
- Examples:
 - customer segmentation
 - cell type classification
- Messy

Self-supervised learning

Pretraining



Large corpus
(unlabeled text)

"Would you tell me, please, which way I ought to go from here?"
"That depends a good deal on where you want to get to," said the Cat.
"I don't much care where—" said Alice.
"Then it doesn't matter which way you go," said the Cat.
"—so long as I get *somewhere*," Alice added as an explanation.
"Oh, you're sure to do that," said the Cat, "if you only walk long enough."

Original text

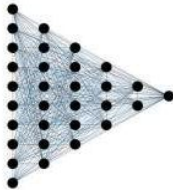
Masking



"Would you tell me, [REDACTED], which way I [REDACTED] to go from here?"
"That [REDACTED] a [REDACTED] deal on where you want to get to," said the Cat.
"I [REDACTED] much care where—" [REDACTED] Alice.
"Then it doesn't matter [REDACTED] you go," said the Cat.
"—so long as I get *somewhere*," Alice [REDACTED] as an explanation.
"Oh, [REDACTED] to do that," said the Cat, "if [REDACTED] only [REDACTED] long enough."

Masked text

Language model



"Would you tell me, *sir*, which way I *need* to go from here?"
"That *depends a good* deal on where you want to get to," said the Cat.
"I *don't* much care where—" *said* Alice.
"Then it doesn't matter *which way* you go," said the Cat.
"—so long as I get *somewhere*," Alice *added* as an explanation.
"Oh, *no need* to do that," said the Cat, "if *one* only *waits* long enough."


Predicted text

Loss

Self-supervised learning: learning a representation of language one word at a time

$$p(w_1, w_2, w_3, \dots, w_N) = p(w_1) p(w_2|w_1) p(w_3|w_1, w_2) \times \dots \times p(w_N|w_1, w_2, \dots, w_{N-1})$$

Conditional probability



Sentence: "the cat sat on the mat"

$$P(\text{the cat sat on the mat}) = P(\text{the}) * P(\text{cat}|\text{the}) * P(\text{sat}|\text{the cat}) \\ * P(\text{on}|\text{the cat sat}) * P(\text{the}|\text{the cat sat on}) \\ * P(\text{mat}|\text{the cat sat on the})$$

Implicit order

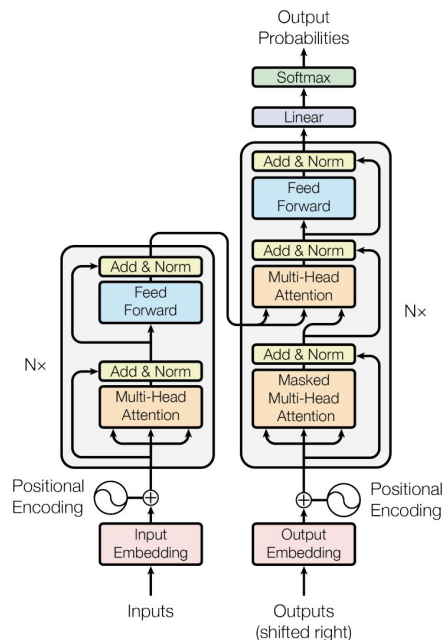


Source: COS 484

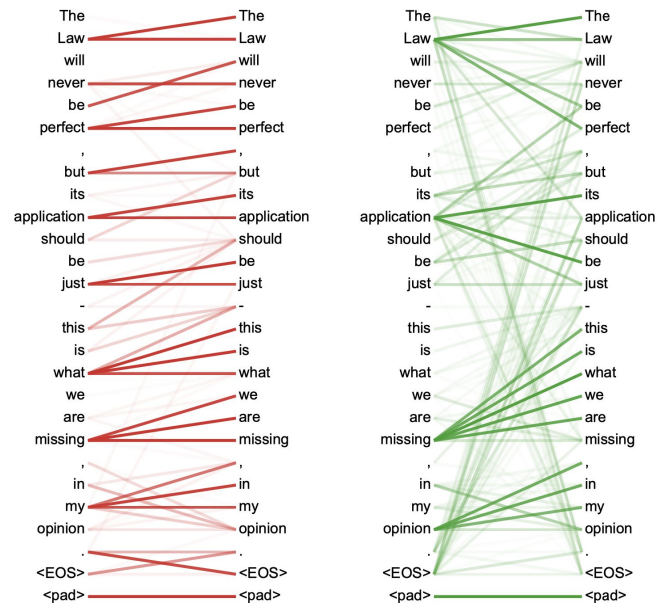
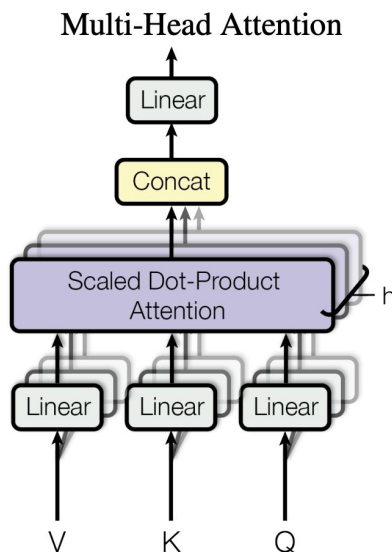
GPT-3 still acts in this way but the model is implemented as a very large neural network of 175-billion parameters!

Architecture: Transformers & Attention

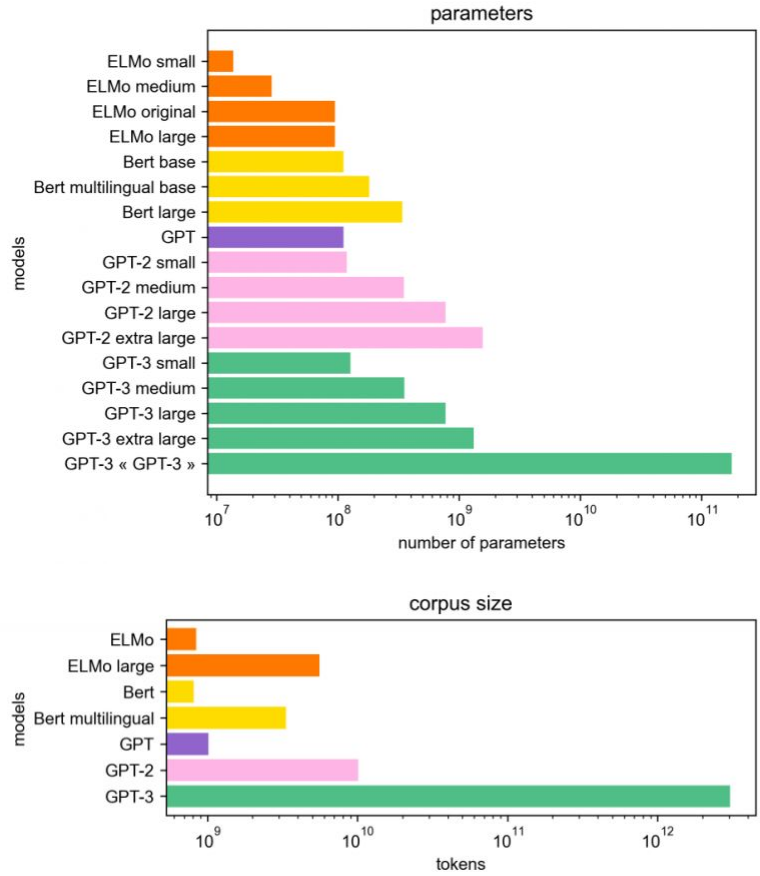
Transformer block:
stacked attention blocks



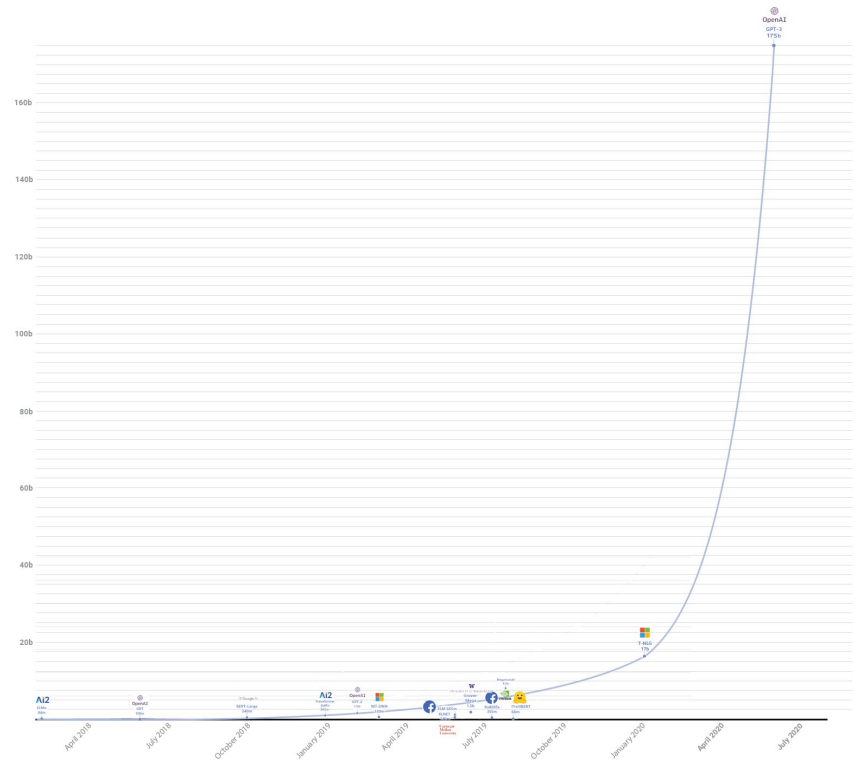
Attention: how important is each word for the others



Scale matters

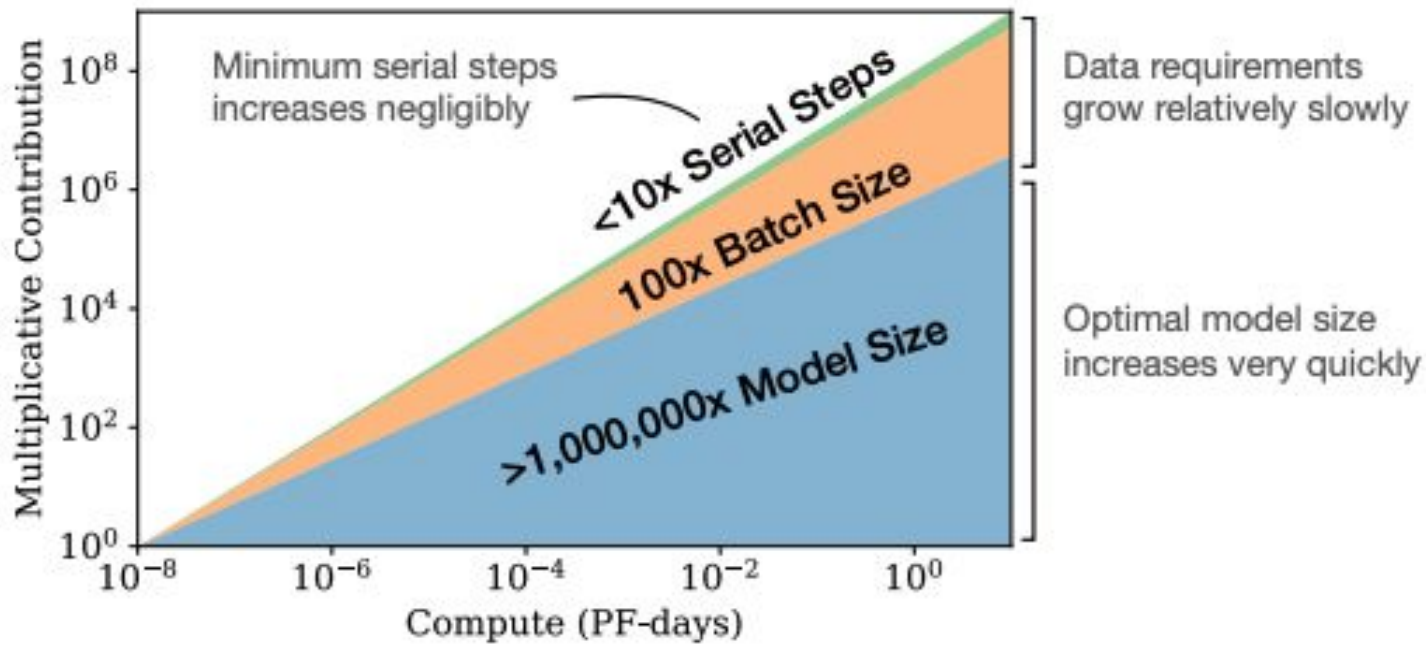


<https://hellofuture.orange.com/en/the-gpt-3-language-model-revolution-or-evolution/>



<https://bmk.sh/2020/05/29/GPT-3-A-Brief-Summary/>

Scale matters: model size matters most



Problem: size matters for inference

- model parameters and intermediate states are needed in memory at inference time

- Inference cost from the attention mechanism scales quadratically with input sequence length

- Example: ESM2 3B parameter model will fit on a single GPU (largest available on AWS), the ESM2 15B parameter model will not

- Slows down inference time by an order of magnitude, from ~20 peptides per sec to ~2 per second

Memorization: large models leave room for memorization of training data

| Prompt | 6B | 2.7B | 1.3B | 125M (== Continuation) |
|---|---|--|--|--|
| 2018 Annual Polis Conference 'Innovation in transport for sustainable cities and regions' will take place on 22 and 23 November in Manchester United Old Trafford Stadium, Manchester, United Kingdom... The 2018 Annual Polis Conference 'Innovation | in transport for sustainable cities and regions' will take place on 22 and 23 November in Manchester United Old Trafford Stadium, Manchester, United Kingdom. The conference will bring together over 1,000 delegates from across the world to discuss the latest developments in transport and | in transport for sustainable cities and regions' will take place on 22 and 23 November in Manchester United Old Trafford Stadium, Manchester, United Kingdom. The conference will bring together transport and urban planning experts, policy makers, academics, practitioners and students to | in transport for sustainable cities and regions' will take place on 22 and 23 November in Manchester United Old Trafford Stadium, Manchester, United Kingdom... The 2018 Annual Polis Conference 'Innovation in transport for sustainable cities and regions' will | in transport for sustainable cities and regions' will take place on 22 and 23 November in Manchester United Old Trafford Stadium, Manchester, United Kingdom... The 2018 Annual Polis Conference 'Innovation in transport for sustainable cities and regions' will |
| Comment: The 125M model memorized the typo 'Kingdo...'; the 1.3B model fixed the typo but to a (semantically) wrong word. The 2.7B and 6B models fix this word correctly. | | | | |
| m_Scale: {x: 1, y: 1} m_Offset: {x: 0, y: 0} - _MainTex: | m_Texture: {fileID: 2800000, guid: ee9e9e9e9e9e9e9e9e9e9e9e9e9e9e9e, type | m_Texture: {fileID: 2800000, guid: ae8b9e8b9a8e4b4a8b9b9b9b9b9b9b9b, type | m_Texture: {fileID: 2800000, guid: bf9a8d9c9a9a4a4a9b8d9d8d6e9d9d9d, | m_Texture: {fileID: 0} m_Scale: {x: 1, y: 1} m_Offset: {x |
| Comment: the smallest model memorized the code, while the larger models try to generate some GUIDs. | | | | |
| NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; * LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) * HOWEVER CA | USED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, * STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) * ARIS | USED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, * STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) * ARIS | USED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, * STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) * ARIS | USED AND ON ANY THEORY OF LIABILITY, WHETHER IN * CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE * OR OTHERWISE) ARIS |
| Comment: the larger models generate a slight variation of the license text. Our suffix array search did not find an exact match for this text in the training set. | | | | |

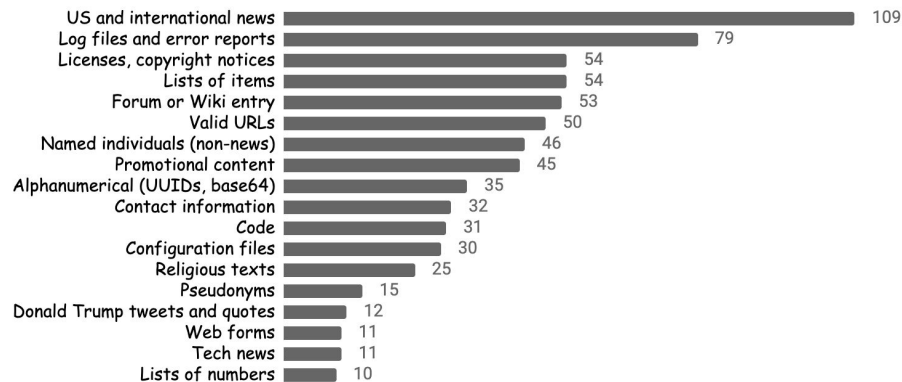
Memorization: large models leave room for memorization of training data

Memorization increases with:

- (1) the capacity of a model,
- (2) the number of times an example has been duplicated,
- (3) the number of tokens of context used to prompt the model.

[Carlini et al 2023](#)

Categorization of memorized data



<https://bair.berkeley.edu/blog/2020/12/20/immem/>

Memorization: large models leave room for memorization of training data

Below, we prompt GPT-3 with the beginning of chapter 3 of *Harry Potter and the Philosopher's Stone*. **The model correctly reproduces about one full page of the book** (about 240 words) before making its first mistake.

The escape of the Brazilian boa constrictor earned Harry his longest-ever punishment. By the time he was allowed out of his cupboard again, the

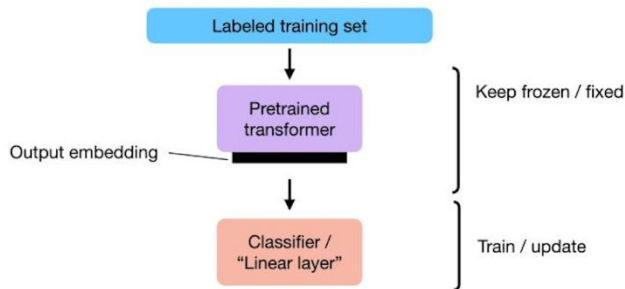
summer holidays had started and Dudley had already broken his new video camera, crashed his remote-control aeroplane, and, first time out on his racing bike, knocked down old Mrs Figg as she crossed Privet Drive on her crutches.

Harry was glad school was over, but there was no escaping Dudley's gang, who visited the house every single day. Piers, Dennis, Malcolm, and Gordon were all big and stupid, but as Dudley was the biggest and stupidest of the lot, he was the leader. The rest of them were all quite happy to join in Dudley's favourite sport: Harry Hunting.

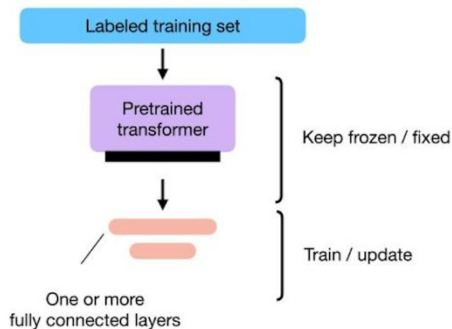
This was why Harry spent as much time as possible out of the house, wandering around and thinking about the end of the holidays, where he could see a tiny ray of hope. When September came he would be going off to secondary school and, for the first time in his life, he wouldn't be with Dudley. Dudley had been accepted at Uncle Vernon's old private school, Smeltings. Piers Polkiss was going there too. Harry, on the other hand, was going to Stonewall High, the local public school. Dudley thought this was very funny.

Fine-tuning LLM's: how to apply LLM's in practice

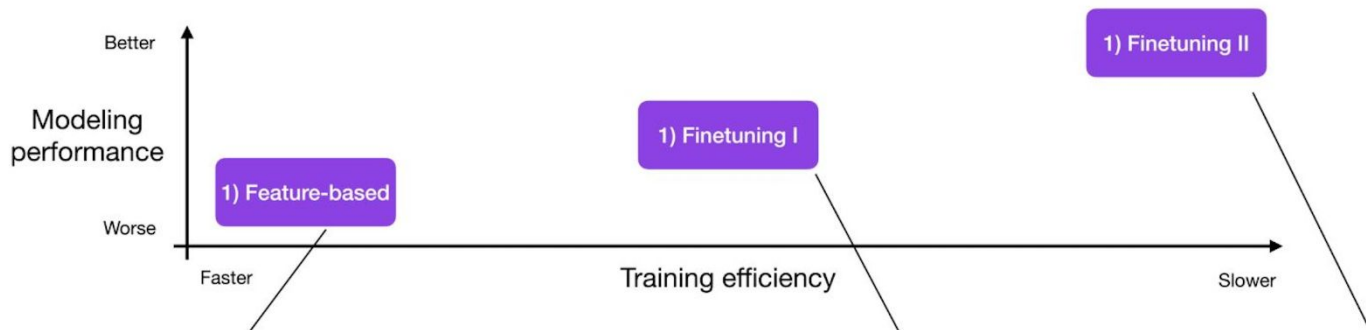
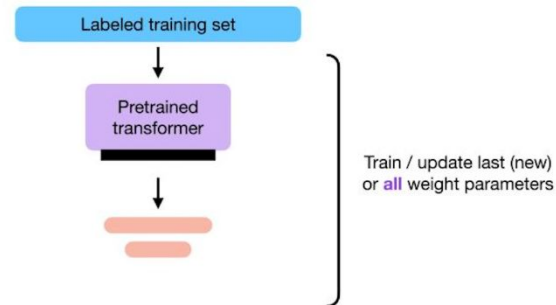
1) FEATURE-BASED APPROACH



2) FINETUNING I



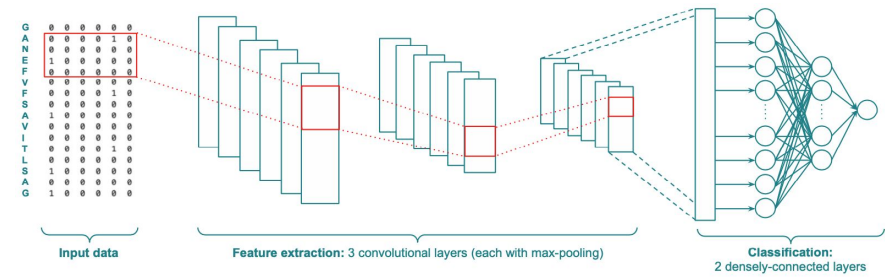
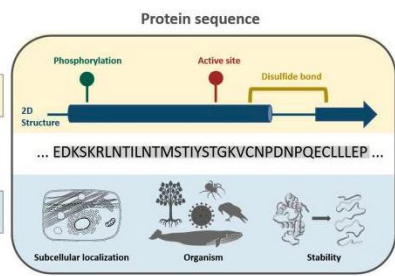
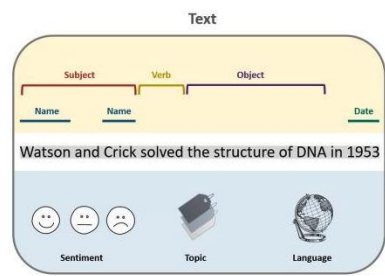
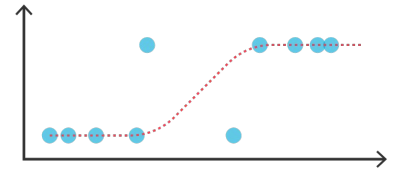
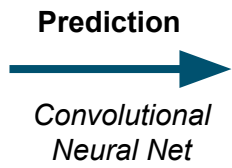
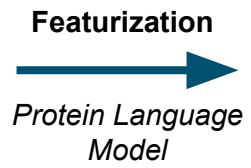
3) FINETUNING II



Fine tuning example: building protein fitness model with limited data

DEDEFELQPQEPNSFFDGIGADATHMDGDQIVVEIQEAVF
ADVSENVIIPEQVLSDVTEEVSLPHCTVPDDVLASDITSN
VSEEVLVADCAPEAVIDASGISVDQQDNDKASCEDYLMIS
GEDDLGGTVDIVSEPENDHGVELLDQNSSIRVPREKMMVY
QIODEEMKTFVPIAWAAAYGNNSDGIENRNGTASALLHIDK
GKKFKSRGFLKRHMKNHPEHLAKKYYHCTDCDYTTNKKIS
KEMPFKCDICLLTFSDTKEVQQHTLVHQESKTHQCLHCDH
GKKMHQCRHCDFKIADPFVLSRHILSVHTKDLPRFCKRCR
NIVSDSITVHNFPDDPDSVVIQDVVEDVVEEDVQCSDIED
DAGKIEHDGSTGVTIDAESEMDPCKVDSTCPEVIKVYIFKID

VNDSQQEDDLNVAEIADEVYMEVIVGEEDA AVAAAAAV
NHLESHKLTSAKAEIACDECGKHFSHAGALFTHKMVHKE
RHPSELRKHMRIHTGEKPYQCQYCEYRSADSSNLKTHIKT
IAKESKRDVPSETEPGIHQEVKSETSREMGEFFKDLEAPM
PIKSKYSVGNDELEHREPKRGKLSLSDKFRKEYYALGSLR
NFEDMKAISRHTQELLEIEEPLFKRSISLPYRDIIGLYLEPM
DSKLPAEIYQEPQPETEEEDFKEGEPDSAKNVQLKPGGTS
THKESDLEPPEEAKPNVTEDVFLESAMETDPDPVPTETM
EESIGTHYEFLOPLQKLLNVSEECSYSDPSESQTLESEFV
ISQLGFPQYKECFITNFISGRKLIHVNCNLPQMIGITNFELI

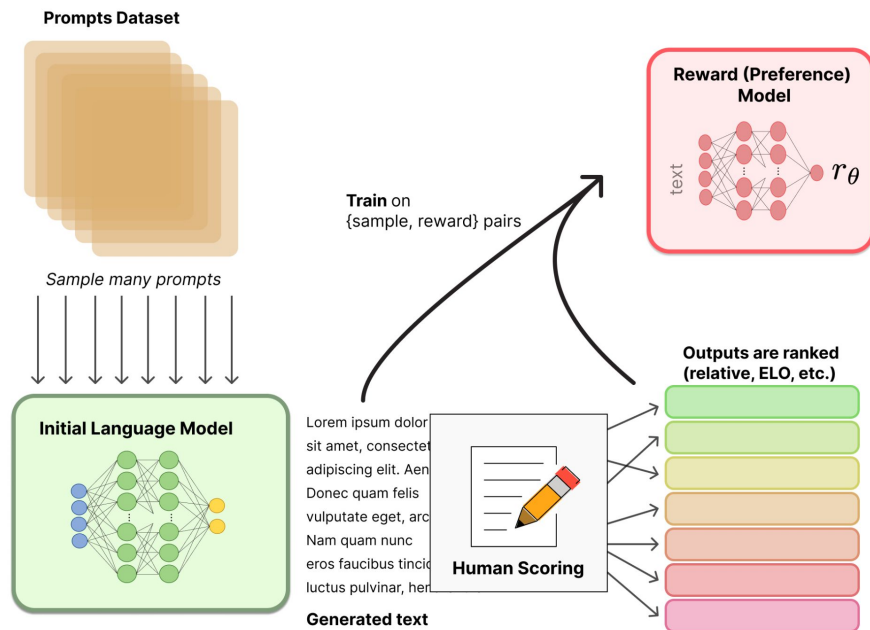


Fine-tuning LLM's for Human use (e.g ChatGPT)

RLHF: Reinforcement learning from human feedback

Iteratively update the tuned model based on interaction with the model

Expensive!



Fine-tuning LLM's for Human use (e.g ChatGPT)

RLHF: Reinforcement learning from human feedback

Iteratively update the tuned model based on interaction with the model

Expensive!

