

# Compression for AGI

---

Yuxuan Lu

Nov. 3 2023

Northeastern Human-Centered AI Lab

- Generative models are *Lossless* compressors
  - “ChatGPT Is a Blurry JPEG of the Web”? No!
- LLMs are *State-of-the-Art* text compressors
  - comparing to deflate(gzip), Zstd, etc.
- Re-think about the training objective of foundation models

Revisit the Training Process of LLMs

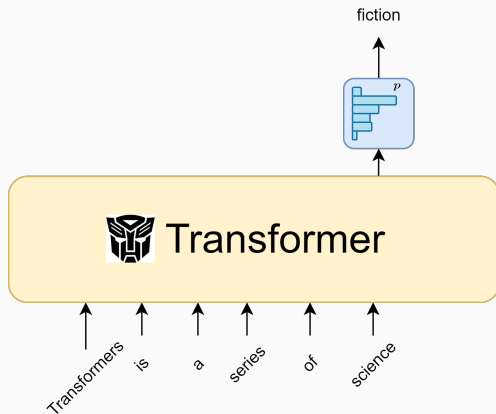
What are compressors

Why are LLMs *Lossless* compressors

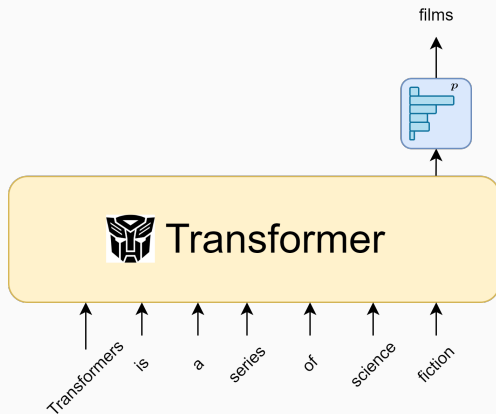
Rethink the goal for foundation models

# Revisit the Training Process of LLMs

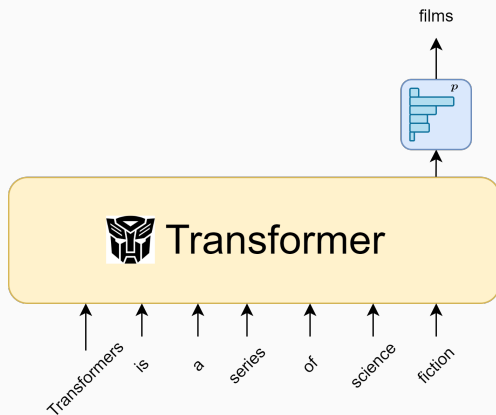
---



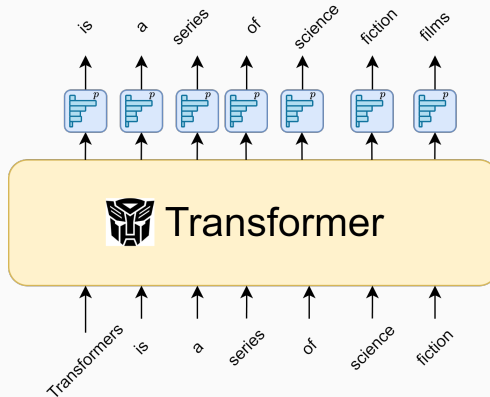
**Figure 1:** Transformer: A Black Box – step 1



**Figure 2:** Transformer: A Black Box – step 2



**Figure 3:** Transformer: A Black Box – step 2



**Figure 4:** Transformer: A Black Box – training



- We have a sequence:  $X = [x_1, x_2, \dots, x_n]$
- We put this sequence into LLM and get a sequence of probabilities:
  - $p_i = P(x_i | x_{<i})$
- We want  $p_i$  approaches 1, so we minimize the following loss during pre-training:
  - $L = \sum_{i=1}^n -\log P(x_i | x_{<i})$

**Questions?**

# What are compressors

---

- Compressors aims to find a better (smaller) way to express the same thing
-

**Why are LLMs *Lossless* compressors**

---

- Naturally, you would assume an LLM is a *lossy* compressor
  - That turns training corpus to model parameters
  - For LLaMa, the training dataset is 5.6TB
  - And the 65 billion parameters takes about 130GB of storage
    - So 43x compression rate?
    - Lossy!

- Naturally, you would assume an LLM is a *lossy* compressor
  - That turns training corpus to model parameters
  - For LLaMa, the training dataset is 5.6TB
  - And the 65 billion parameters takes about 130GB of storage
    - So 43x compression rate?
    - Lossy!
- Actually, LLaMa can compress the entire 5.6TB of training corpus to 397.3 GB losslessly
  - 14x compression
  - Best text compressor: 8.7x compression

- To compress the *entire training corpus*  $\mathcal{C}$  *losslessly*, we only need *the CODE to train the models* and  $\sum_{i \in \mathcal{C}} -\log P(x_i | x_{<i})$  bits of information.
  - The result size (after compression) *ISN'T* related the number of parameters!



- To compress the *entire training corpus*  $\mathcal{C}$  *losslessly*, we only need *the CODE to train the models* and  $\sum_{i \in \mathcal{C}} -\log P(x_i | x_{<i})$  bits of information.
  - The result size (after compression) *ISN'T* related the number of parameters!
- But, HOW?

- Imagine Alice is trying to send some text to Bob through a telephone wire
- Sending the raw text is very expensive
  - too large
- So Alice need to find a way to compress the data losslessly
- Alice needs to encode the data to “something”
- Bob needs to decode the “something” back to data
- So we need an “encoding” algorithm and a “decoding” algorithm

- Let's say we have a vocabulary of  $\mathcal{V}$  (which should be defined in the training code), and we want to encode the first token  $x_1$ .

- Let's say we have a vocabulary of  $\mathcal{V}$  (which should be defined in the training code), and we want to encode the first token  $x_1$ .
- First, we initialize the model and feed it the training data
  - and we can get the probability distribution of the first token  $p$  and an index  $x_1 \in \mathcal{V}$

- Let's say we have a vocabulary of  $\mathcal{V}$  (which should be defined in the training code), and we want to encode the first token  $x_1$ .
- First, we initialize the model and feed it the training data
  - and we can get the probability distribution of the first token  $p$  and an index  $x_1 \in \mathcal{V}$
- With arithmetic encoding, we can encode the index  $x_1$  with  $-\log_2 p_{x_1}$  bits
  - Rather than  $-\log_2 \frac{1}{|\mathcal{V}|}$  bits

- Let's say we have a vocabulary of  $\mathcal{V}$  (which should be defined in the training code), and we want to encode the first token  $x_1$ .
- First, we initialize the model and feed it the training data
  - and we can get the probability distribution of the first token  $p$  and an index  $x_1 \in \mathcal{V}$
- With arithmetic encoding, we can encode the index  $x_1$  with  $-\log_2 p_{x_1}$  bits
  - Rather than  $-\log_2 \frac{1}{|\mathcal{V}|}$  bits
- We send those  $-\log_2 p_{x_1}$  bits to Bob
  - Let's call it  $z_1$

- Let's say we have a vocabulary of  $\mathcal{V}$  (which should be defined in the training code), and we want to encode the first token  $x_1$ .
- First, we initialize the model and feed it the training data
  - and we can get the probability distribution of the first token  $p$  and an index  $x_1 \in \mathcal{V}$
- With arithmetic encoding, we can encode the index  $x_1$  with  $-\log_2 p_{x_1}$  bits
  - Rather than  $-\log_2 \frac{1}{|\mathcal{V}|}$  bits
- We send those  $-\log_2 p_{x_1}$  bits to Bob
  - Let's call it  $z_1$
- ... and update the model to minimize the loss, and repeat these steps for every data

- Bob retrieves the code for training the model



- Bob retrieves the code for training the model
- And he can initialize the exact same model (with random seed defined in the code)

- Bob retrieves the code for training the model
- And he can initialize the exact same model (with random seed defined in the code)
- So he will have the *exact same distribution*  $p$  for the first token
  - He then can decode the token with  $p$  and  $z_1$  with arithmetic decoding, get  $x_1$

- Bob retrieves the code for training the model
- And he can initialize the exact same model (with random seed defined in the code)
- So he will have the *exact same distribution  $p$*  for the first token
  - He then can decode the token with  $p$  and  $z_1$  with arithmetic decoding, get  $x_1$
- And he will run the training loop to update the model with  $x_1$
- Note that the model parameters isn't transmitted through the phone wire

- By sharing the same code, Bob can initialize the same model as Alice's
- With arithmetic encoding and decoding, bob can decode the token with the same probability distribution produced by the model and the encode transmitted from Alice
- And by training at the same time, the model is always synced between Alice and Bob

- With arithmetic encoding, we can use fewer bits to encode something have a higher probability in a distribution.
  - If the distribution is uniform,  $-\log_2 p_{x_i} = -\log_2 \frac{1}{|\mathcal{V}|}$ , which is same as naive storage
  - And as the model is continually training, it can predict the next token with higher and higher probability.
- Bob can re-construct the entire training corpus  $\mathcal{C}$  with  $\sum_{i \in \mathcal{C}} -\log P(x_i | x_{<i})$  bits of information
  - Which is exactly the same with the sum of training loss on all tokens!

**Questions?**

# Rethink the goal for foundation models

---

- If a computer program translates English and Chinese using an oracle comprising all possible combinations of Chinese and English combinations.
  - Does it have understanding of translation?

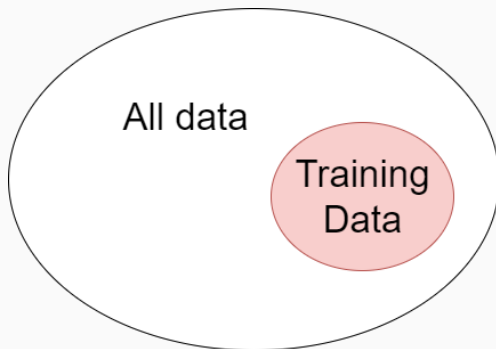


- If a computer program translates English and Chinese using an oracle comprising all possible combinations of Chinese and English combinations.
  - Does it have understanding of translation?
  - It has the *least* understanding of translation

- If a computer program translates English and Chinese using an oracle comprising all possible combinations of Chinese and English combinations.
  - Does it have understanding of translation?
  - It has the *least* understanding of translation
- What if it do it by a set of rules?
  - Tt have *some* understanding of translation.

- If a computer program translates English and Chinese using an oracle comprising all possible combinations of Chinese and English combinations.
  - Does it have understanding of translation?
  - It has the *least* understanding of translation
- What if it do it by a set of rules?
  - Tt have *some* understanding of translation.
- If we can make the rule set smaller, it will generalise better.

- “Lossy” compression means the model “remembers” everything in the training dataset
  - BAD generalization
- “Lossless” compression means the model can better predict unseen data samples
  - Better compression means GOOD generalization
    - Because *EVERY* example is unseen



**Figure 5:** All data V.S. Training Data

- For foundation models, what we want is good generalization ability
  - i.e. ability to generate or write *UNSEEN* samples.

## A recipe for perception

- Collect all useful perceptual information
- Learn to compress it as best as possible with a powerful foundation model
  - Better architecture
  - Scale
  - Tool use
  - ...

- How LLM works
- How LLM works as a compressor
  - And how should we use LLM
  - LLM isn't a search engine
- How compressor generates intelligence

**Questions?**