

Neural Network

Initialization: setting values of W & b for layers use paradigm depending on activation function for asymmetric & prevent vanishing gradient

→ log its: $\ln @W @b$

Normalization: apply batch / layer normalization to get desired properties \rightarrow mean & std dev

→ activation fxn.

→ loss calculation

→ zero the gradients

Backward pass: compute loss w.r.t. W & b

Optimizer: update the W & b

Why initialize with paradigm?

- Preventing Vanishing & Exploding gradient (too small or too large) gradients which make training difficult especially in deep NN
- Efficient convergence
- Avoiding symmetry breaking (if every w is same, every neuron will behave same i.e. every neuron will learn same features which means the entire network works as one neuron)

1. Input:

- We start with the **input data**. This could be structured data (like tabular features), image data (pixel intensities), text data (which needs to be embedded), etc.
- If the input is categorical, we embed it using an **embedding layer** (e.g., word embeddings for text), converting discrete tokens into dense, continuous-valued vectors.

2. Weights and Bias Initialization:

- Neural network layers have **weights** and **biases**, which are the trainable parameters of the model.
- These parameters are either initialized **randomly** or based on a specific **initialization scheme** (e.g., Xavier Initialization, He Initialization, etc.) depending on the activation function and network design.

3. Forward Pass:

- In the **forward pass**, we compute the output of the model layer by layer:
 1. Multiply the input (or embedding) with the **weights**, add the **biases**, and compute the **logits**:
Here, w are the weights, b are the biases, and x is the input (or the output from the previous layer).
 2. Pass the logits through an **activation function** (e.g., ReLU, Sigmoid, Tanh):
This introduces non-linearity into the network, allowing it to learn complex patterns.
logits
 3. (Optional) Pass the activations through a **normalization layer** (like BatchNorm or LayerNorm) to stabilize the training process.

4. Loss Calculation:

- The final output from the network is compared to the **ground truth labels** using a **loss function**.
- Common loss functions include:
 - Mean Squared Error (MSE) for regression,
 - Cross-Entropy Loss for classification tasks.
- The loss function computes a scalar value (the **loss**) that measures how far the model's predictions are from the ground truth.

5. Zero the Gradients:

- Before the **backward pass**, we zero out the gradients stored in the model's parameters to prevent accumulation of gradients from previous steps:

```
optimizer.zero_grad()
```

6. Backward Pass (Backpropagation):

- We perform backpropagation to compute the gradients of the loss with respect to each trainable parameter (weights and biases) in the network:

```
loss.backward()
```

During this process:

- The chain rule is applied to propagate the gradient of the loss from the output layer back to the earlier layers.
- The gradients of and are stored in their respective `.grad` attributes.

Step 7: Optimizer Step

1. The **optimizer** is responsible for updating the weights and biases using the computed gradients.

- Common optimizers include:
 - **Stochastic Gradient Descent (SGD)**: Updates weights based on the gradient direction and learning rate.
 - **Adam**: A more advanced optimizer that combines momentum and adaptive learning rates.
- Optimizer updates parameters using update rules like:
where α is the learning rate.

2. After the optimizer updates, the model's parameters are adjusted to (ideally) reduce the loss.

Step 8: Repeat for Multiple Iterations

1. Steps 3 to 7 (forward pass → loss calculation → backward pass → optimizer step) are repeated for many **iterations** or **epochs**.
 - An **iteration** refers to processing a single batch of data.
 - An **epoch** refers to processing the entire dataset once.
2. Over time, the loss decreases as the model's parameters are adjusted to minimize the difference between predictions and ground truth.

Normalization is done before activation

If we normalize the layer to mean & std (some value 0, 1) every epoch it will lead to:

- overwriting the learned representation
- loss of expressive power
- gradient destruction

causes:

- constraints learning
- convergence issues
- destroys specific patterns

We want the neurons to have varied parameters which allows network to learn patterns. ^(w & b)

Batch normalization leads to regularization as it forces each input to add entropy from batch

In training we use batch but when we want to deploy:
calculate the std dev & mean over the entire dataset
i) after training
ii) along training (Pytorch does this)

when using BatchNorm we don't need to use bias for its prior layers or we subtract mean which leads the effect of $b_1 = 0$

There are better options than BatchNorm

$$y = \frac{x - E(x)}{\sqrt{Var(x) + \epsilon}} \propto \gamma + \beta$$

Batch Norm(D), eps = to prevent variance from being 0

γ, β (bgain & bias) to shift the values
that allows learning

affine = True, allows γ & β to change with
epochs

track_running = calculates Unmean & Std
while training

Batch norm helps mitigate a few effects that arise from
improper layer initialization

Weight initialization for layer is dependent on the activation
fn and usually takes the form of
 $gain \times (\text{initialization paradigm formula})$

gain plays a major role as it ensures that the magnitude
of the weights of the layers remain balanced and
significant to maintain stable training & prevent vanishing
and exploding gradients

Activation Distribution:

Meaning: Activation functions like Tanh map inputs to a specific range, and neurons can saturate at the extreme ends of this range (near -1 or +1). Monitoring activation distributions helps understand whether neurons are operating in these saturated regions.

Problem Identification: If too many neurons are in the saturation region, their gradients become very small, leading to vanishing gradients. This means that the network learns very slowly or not at all, especially in deep networks.

Solution: Proper weight initialization techniques, like Kaiming initialization, prevent neurons from saturating too quickly. Additionally, Batch Normalization can help maintain a healthy distribution of activations throughout the training process by normalizing the outputs of each layer.

For every layer except the output print the mean, std and % of logits which lead to saturation, histogram of weights Activation distribution, gives an idea whether neurons are saturating or not, which could affect learning  
→ Proper weight initialization

Gradient Distribution:

Meaning: The gradient represents how much each parameter needs to be adjusted based on the loss. Monitoring the distribution of gradients allows us to assess whether backpropagation is working effectively.

Problem Identification: If gradients are very small (vanishing gradients), the weights will update minimally, causing slow learning. Conversely, if gradients are too large (exploding gradients), the updates can be too aggressive, causing unstable training or even NaN values in the network.

Solution: To mitigate this, you can use gradient clipping or ensure the use of appropriate initializations (e.g., Kaiming for ReLU or Xavier for Tanh). Additionally, using BatchNorm helps prevent gradients from either vanishing or exploding by normalizing the outputs across layers.

For every layer except the output print the mean & std of the gradient, histogram of gradients

Gradient distribution, ensures gradients are not too large (exploding) or small (vanishing)

→ Normalization + proper weight initialization

Weight Gradient Distribution:

Meaning: This focuses on the distribution of gradients specifically for the weights in the network. It helps you understand whether the optimizer is making reasonable updates to the parameters.

Problem Identification: If the gradients are too large (compared to the magnitude of the weights), the updates could be too aggressive, causing large jumps in the loss landscape and potentially missing the optimal solution (exploding gradients). If the gradients are too small, the learning rate may be too small or the gradients may be vanishing (too small to make effective updates).

Solution: Weight regularization (e.g., L2 regularization) and maintaining appropriate gradients (using BatchNorm and good weight initialization) can help mitigate these issues. You can also adjust the learning rate based on the distribution of the gradients to avoid large updates or stagnation.

For parameters print the mean, std & grad : data ratio

Weight gradient distribution, helps assess whether the gradients are appropriately scaled to parameters, crucial for stable weight updates

→ Adjusting lr, normalization

Update Magnitudes:

Meaning: This metric compares the standard deviation of weight updates (gradients scaled by the learning rate) with the standard deviation of the weights themselves. It helps assess if updates are too large or too small.

Problem Identification: A log10 ratio significantly greater than -3 indicates that the updates are too large, which could destabilize training. Conversely, a ratio much smaller than -3 suggests slow convergence, as the updates are too small.

Solution: Adjusting the learning rate is key here. If the update ratio is too large, reduce the learning rate. If it's too small, you may need to increase the learning rate or adjust the optimizer to make updates more efficient. A target ratio of about (0.1% of parameter value) is ideal, as it ensures updates are significant but not overwhelming.

$\log_{10} \left(\frac{\text{std of updates } (lr * p.grad)}{\text{std of parameter values (data)}} \right)$, should be 10^{-3}
if $>$, update overwhelms parameter
 $<$, slow convergence

10^{-3} , means updates are 0.1% of the parameter value

Update magnitude, helps monitor how effectively the optimizer is updating parameters, ensuring that updates are not too large or too small

→ Adjust lr

`sum()`

`0` (along the row)

↓ so sum is given for columns
[1, column]

`1` (along the column)

→ so sum is given for rows
[rows, 1]

Pytorch does gradient calculation for us making everything easy, but if one was to code everything keep in mind:

→ match the shapes

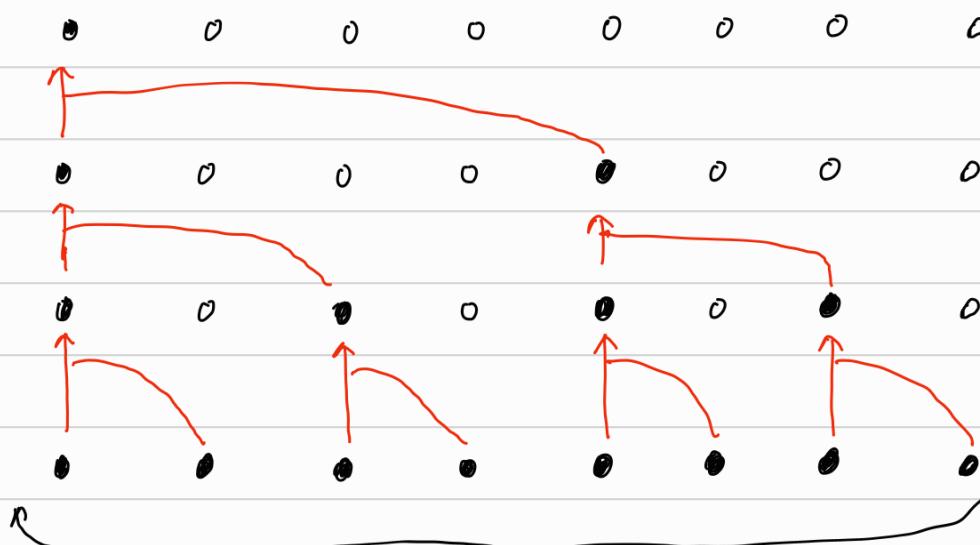
→ pay attention to broadcasting in operations

(if we broadcast [27, 1] to [27, 34] in forward pass we have to do opp. in backward pass)

→ one variable might be used multiple times so the gradient must be updated the same number of times

Go through calculus once

Wavenet architecture:



`context_size = 8`

decreases in each subsequent layer

Matrix multiplication (\otimes) works only on the last dimension of the first input all other dimensions are left unchanged

$$\text{Wavnet} \left\{ \begin{array}{l} 1 \\ 2 \\ 4 \\ 8 \end{array} \right\} \text{torch}(4, 4, 20) * \text{torch}(20, 200)$$

→ embeddings of bigram concatenated
Batch size = 4 { : bigrams (0,1) (2,3) (4,5) - -

`torch.view()` → a powerful tool that performs functions such as concatenation easily

Broadcasting, is a powerful feature but can cause subtle bugs which might affect the performance slightly

Convolution → “a for loop”, allows to forward linear layers efficiently over space

Building Neural Networks:

- Reading Documentation
- Seeing what layers do
- Pay attention to the shapes of input, output to the layers
- Build prototypes to see the changes in shapes through the layers

Language models, take block of n & predicts $n+1^{\text{th}}$ element till it reaches max. element size

Context size, for current element how many previous elements to pay attention to

Transformer Architecture:

- Encoder only: input level tasks like classification, NER
- Decoder only: autoregressive tasks such as text generation
- Encoder-decoder: seq2seq tasks such as MTL, summarization

Encoder:

- processes input into contextualized embeddings
- allows all tokens to attend to each other

Decoder:

- used for generating output tokens
- autoregressive using triangular masking

Embeddings: converting discrete tokens into dense continuous-valued vectors (word, contextual, subword, position, character)

Attention: allows tokens to "attend to" other tokens & gather relevant info dynamically, can be thought of as communication b/w tokens to gain info about each other and find which ones affect them and ones they affect

$$\text{Attention } (\mathbf{Q}, \mathbf{V}, \mathbf{K}) = \text{softmax} \left(\frac{\mathbf{Q} \mathbf{K}^T}{\sqrt{d_k}} \right) \cdot \mathbf{V}$$

\mathbf{Q} , query

\mathbf{V} , value

\mathbf{K} , key

d_k , dimension of key vector

} Each token is converted to three vectors using \mathbf{Q} : each token finds token with key, value valuable to them

It has no concept of space so positional encoding is important
 $\sqrt{d_k}$, scaled attention to diffuse the value to prevent high/low values

Self-attention: Q, K, V come from same sequence
Encoder, Decoder

Cross-attention: K, V can come from other source
Encoder-Decoder

Transformer model architecture:

- MHA: multiple single heads running parallelly each with its own self attention which allows the model to focus on different aspects of the sequence
- FFN: two linear layers & one non-linear activation fn.

$$FFN(x) = \text{relu}(xW_1 + b_1)W_2 + b_2$$

The next two components help with the problem of optimization of DNN:

- Residual connections: a highway of gradients from input to output add a few computation branches which are added to the highway.
 -  helps with flow of gradient during training, prevents vanishing gradients, allows DN to be trained more effectively.
 - It directly puts input in the output layers which is necessary to preserve original features while enabling transformation
- Normalization: normalize features to prevent exploding/vanishing gradient
 - Batch Norm: normalization of features across all samples of a batch together, focuses on general features of large dataset
 - LayerNorm: normalization of features of a sample at a time, preserves and highlights individuality of sample useful for sequence tasks
- Transformers use LayerNorm

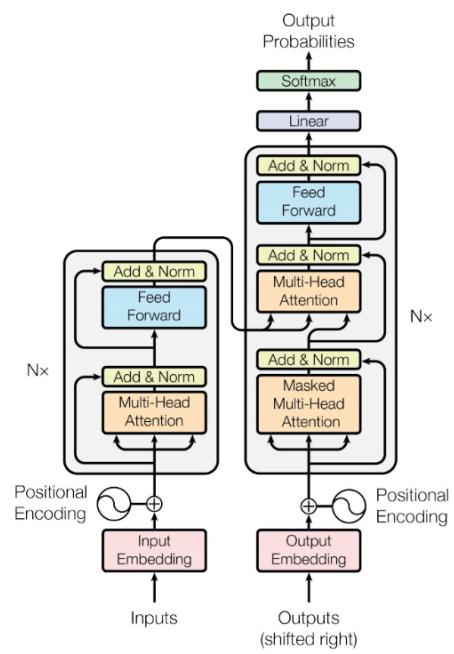


Figure 1: The Transformer - model architecture.

Dropout, regularization technique for NN that works by randomly switching off neurons forcing the network to learn without using neurons indirectly creating an ensemble of neural networks

GPT training process

1) Pre-training

- Goal: learn general language patterns using causal language modelling
- Task: predict next token in sequence based on previous tokens
- Data: billions of tokens
- Architecture: Decoder

2) Fine-tuning

→ Supervised: trains the pre-trained model on task-specific, labeled datasets to improve performance on specific tasks

→ RLHF:

- Human evaluators rank responses
- A reward model is trained to predict quality of responses
- Model is further fine-tuned using RL (PPO) to align with human preferences

Tokenization, process of breaking raw text into tokens (characters, words or subwords) depending on tokenizer.

String → Unicode code → utf-8 byte level binary representation
↓ Tokenizer

Why not use unicode code / utf-8 as vocabulary?

- large size
- lacks semantic indexing
- non-strategic

Byte Pair Encoding:

- starts with a vocab of 256 (all UTF-8 characters)
- iteratively merges the most frequent pair of characters / bytes into new tokens
 - ↑ Works across different categories (numeric, alphabets etc.)
 - ↑ reduces sequence length by creating multi-char / multi-byte token
 - ↓ Greedy frequency based merging can merge different categories, destroying semantic meaning

Tiktoken, converts text to utf-8 and performs BPE at byte level used for encoding efficiently in inference pipelines

Sentencpiece, works directly on unicode points, rare token handling

- character coverage: rare tokens can be replaced by UNR
 - byte fallback: rare tokens can fall back to byte-level encoding
- great for training new vocabularies from scratch for domain specific tasks

Training data for Tokenizer:

- representative data (what LLM will encounter in training & deployment)
- larger, domain-specific tokens help reduce sequence lengths while preserving semantic meaning
- Special tokens: provide metadata cues to the LLM e.g. separating documents, handling context or signalling task specific instruction

Optimal Vocab size:

Increasing vocab size reduces sequence length, but:

- increases the embedding table size, memory requirements & computational cost
 - beyond a point, tokens become so large (e.g. entire phrases) that semantic richness is lost, degrading the model's ability to learn meaningful representations
- Balance is key

Gisting, compresses large prompts by training specialized tokens that summarize or represent subsets of the text, reducing token count while maintaining semantic meaning.

Tokenizer splits:
• pre-process raw text using regex to create logical splits
• apply BPE on these splits, concatenate tokens

↑ mitigates semantic loss caused by arbitrary BPE merging

Fine-tuning special tokens:

- 1) Define special tokens
- 2) Modify the tokenizer (extend tokenizer's vocabulary)
- 3) Prepare the dataset (special tokens need to be used explicitly)
- 4) Resize model embedding layer
- 5) Fine tune model on dataset that includes special tokens

Optional: Freeze pre-trained parameters

Tokenizer-related LLM issues

1) Spelling errors

a large vocab packed with subwords means some valid words get split into parts, confusing the LLM when encountering unfamiliar spellings

2) Non-english language

3) Arithmetic

Tokenizers may merge numbers, which can make learning arithmetic patterns difficult unless explicitly handled

4) Unseen tokens in training but seen in tokenization (Solid Magikarp)

5) YAML > JSON, because it is token dense

6) Trailing whitespaces, models are trained with tokens with whitespace before

Training data is same for LLM & tokenizer:

↑ Optimized tokenizer for LLM

↑ Better efficiency in Model training (tokenizer produces fewer tokens)

↑ No "Dead" token in Vocab

↑ Improved Domain Adaptation (if training data is domain specific, the tokenizer makes domain specific terms or phrases part of vocab)

↓ Overfitting to training data

↓ Limited generalisation to out of domain tasks (rare or new words UNK or falling to byte-level code)

Balancing tradeoff (in case you plan to train on same data)

- Expand tokenizer training data (includes domain-specific & general purpose)
- Allow byte fallback for rare tokens
- Pre define special tokens / domain specific terms (even if rare in dataset)
- Evaluate generalization & regularization during LLM training

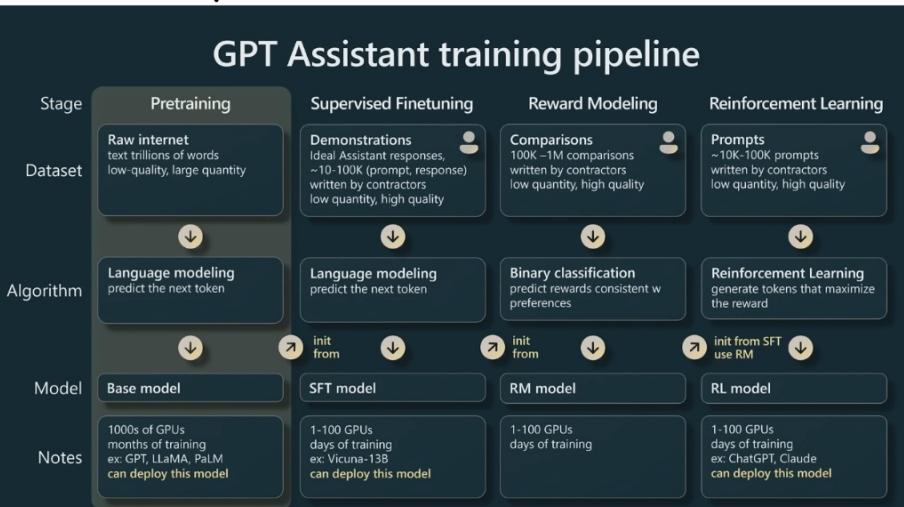
Training data should differ for:

- General purpose LLMs
- Multilingual & task adaptive Models

Training data can be same for:

- Highly domain specific model
- Efficient fine tuning

Stake of GPT



Pre-training trains the base model of LLM where it can only predict the next token no answering to prompts etc.

Tokenization

Transform all text into one very long list of integers.

Typical numbers:

~10-100K possible tokens
1 token ~ = 0.75 of word

Typical algorithm:

Byte Pair Encoding



Pretraining

The inputs to the Transformer are arrays of shape (B,T)

- B is the batch size (e.g. 4 here)
- T is the maximum context length (e.g. 10 here)

Training sequences are laid out as rows, delimited by special <|endoftext|> tokens

Row 1: Here is an example document 1 showing some tokens.

Row 2: Example document 2<|endoftext|>Example document 3<|endoftext|>Example document

Row 3: This is some random text just for example<|endoftext|>This

Row 4: 1,2,3,4,5

T = 10

One training batch, array of shape (B,T)

4342	318	281	1672	3188	352	4478	617	16326	13
16281	3188	362	50256	16281	3188	513	50256	16281	3188
1212	318	617	4738	2420	655	329	1672	50256	1212
16	11	17	11	18	11	19	11	20	11

B = 4 ↓

Pretraining

Each cell only "sees" cells in its row, and only cells before it (on the left of it), to predict the next cell (on the right of it)

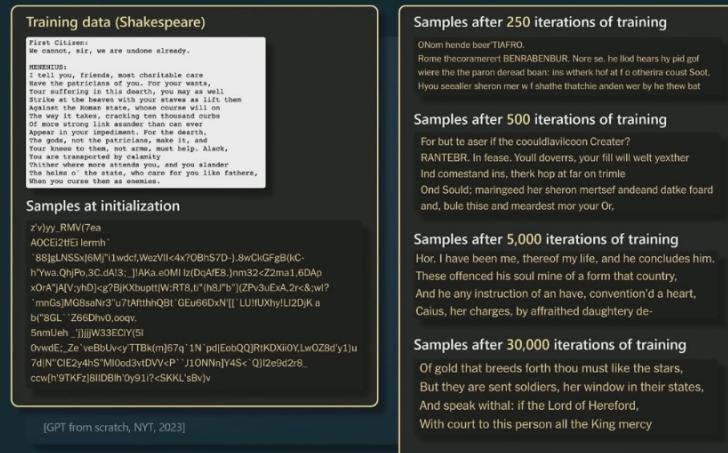
Green = a random highlighted token
Yellow = its context
Red = its target

One training batch, array of shape (B,T)

4342	318	281	1672	3188	352	4478	617	16326	13
16281	3188	362	50256	16281	3188	513	50256	16281	3188
1212	318	617	4738	2420	655	329	1672	50256	1212
16	11	17	11	18	11	19	11	20	11

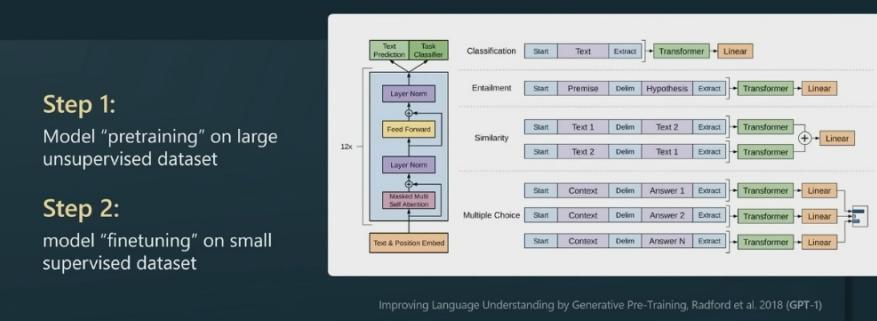
B = 4

Training process



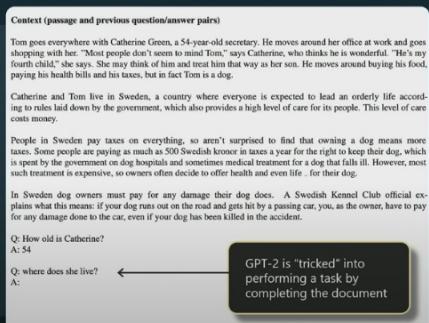
Base models learn powerful, general representations

- Step 1:
Model "pretraining" on large unsupervised dataset
- Step 2:
model "finetuning" on small supervised dataset

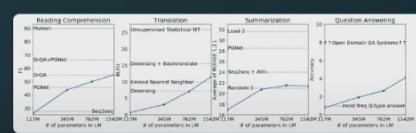


Base models can be prompted into completing tasks

Make your model look like a document!



GPT-2 kicked off the era of prompting over finetuning



Language Models are Unsupervised Multitask Learners, Radford et al. 2019 (GPT-2)

each token predicts the next token in sequence from the vocabulary and parameters are updated

Base models are good at general tasks, to fine-tune if we just need to train on a small supervised data set such as sentiment analysis

GPT Assistant training pipeline



RM Dataset

Three screenshots of the RM Dataset interface showing different examples of string manipulation tasks:

- Example 1:** Write a Python function that checks if a given string is a palindrome.
- Example 2:** Write a Python function that checks if a given string is a palindrome.
- Example 3:** Write a Python function that checks if a given string is a palindrome.

The interface includes code editor, output window, and explanatory text for each example.

RM Dataset

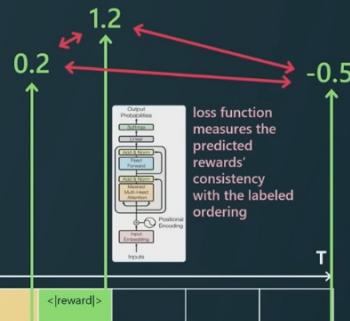
Three screenshots of the RM Dataset interface showing different examples of string manipulation tasks:

- Example 1:** Write a Python function that checks if a given string is a palindrome.
- Example 2:** Write a Python function that checks if a given string is a palindrome.
- Example 3:** Write a Python function that checks if a given string is a palindrome.

The interface includes code editor, output window, and explanatory text for each example.

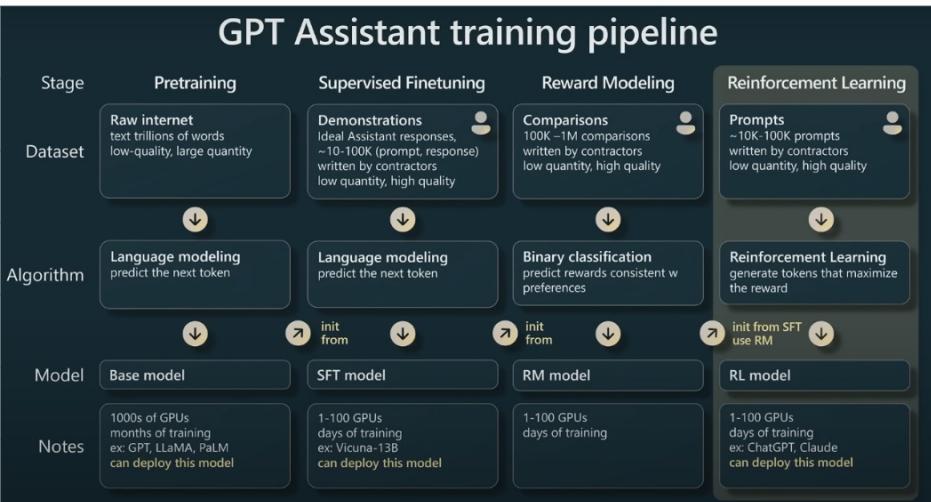
RM Training

Blue are the prompt tokens, identical across rows
 Yellow are completion tokens, different in each row
 Green is the special <reward> token "readout"
 Only the outputs at the green cells are used, the rest are ignored



prompt	completion 1	< reward >			
prompt	completion 2	< reward >
prompt	completion 3	...	< reward >				

$$\begin{array}{l} 1 \rightarrow 1.0 \\ 2 \rightarrow 0.5 \\ 3 \rightarrow 0 \\ 4 \rightarrow - \\ \end{array}$$

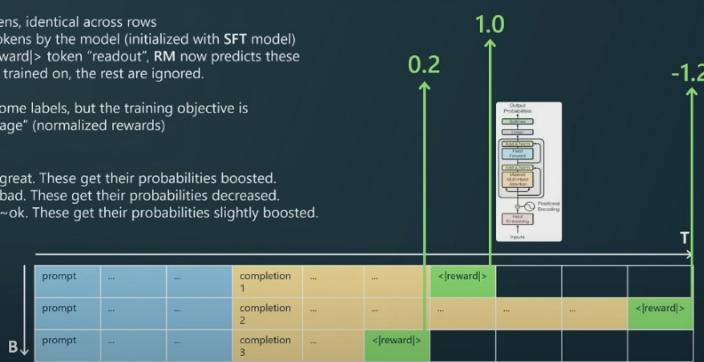


RL Training

Blue are the prompt tokens, identical across rows
Yellow are completion tokens by the model (initialized with SFT model)
Green is the special <reward> token "readout", RM now predicts these
Only the yellow cells are trained on, the rest are ignored.

The sampled tokens become labels, but the training objective is weighted by the "advantage" (normalized rewards)

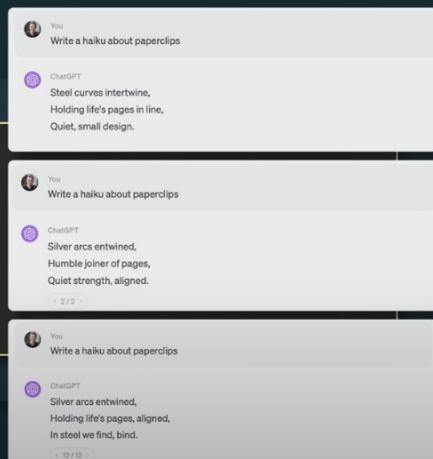
- Row #1 tokens were great. These get their probabilities boosted.
 - Row #2 tokens were bad. These get their probabilities decreased.
 - Row #3 tokens were -ok-. These get their probabilities slightly boosted.



Why RLHF?

It is easier to discriminate than to generate.

Simple example:
it's much easier to spot
a good haiku than
it is to generate one.



Humans find it easier
to discriminate b/w
responses than generate
one

Mode collapse

Finetuned models lose entropy

Base model entropy

Are bugs real? I mean, we've never seen any. But do butterflies turn into ~~crystalises~~ like caterpillars do? I bet it's gross to have grubs crawling around inside you. The bad bunny ~~lost~~ doesn't like to play fetch; I'd be able to take him hiking.

RLHF models might confidently output very few variations.
=> Base models can be better at tasks that require diverse outputs

RLHF model entropy

```
There is no one answer to this question as it depends on what you mean by "bug". If you are referring to insects, then yes, they are real. If you are referring to "are", -0.891 / 99.91% captain, the answer is yes, they're. "mean" -0.346 / 0.66% exist, "define" -11.669 / 0.66% simply -10.239 / 0.66% consider -10.545 / 0.66% "true" -10.559 / 0.66% "r" -12.986 / 0.66% "are" -13.193 / 0.66% "Mean" -13.527 / 0.66% "Are" -14.001 / 0.66% "definition" -14.355 / 0.66% "defined" -14.582 / 0.66% "ARE" -14.691 / 0.66%
```

Mode collapse

Finetuned models lose entropy

Base models can be better in tasks where you have N examples of things and want to generate more things.

Toy example:

Completion →

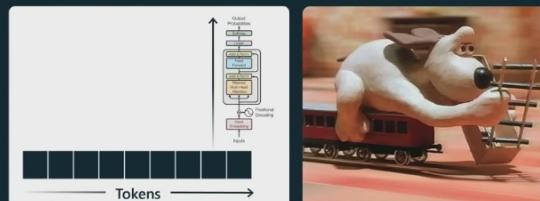
Here are 100 cool Pokemon names I made up:

Charizard
Bulbasaur
Pikachu
Venomoth
Slowpoke
Duranium
Ammnit
Mineboy
Poisonjaw
Nintin
Steelseer
Movemaker
Allsplode
Voltomb
BananaBreath
Shuriken

"California's population is 53 times that of Alaska."

- "For this next step of my blog let me compare the population of California and Alaska"
- "Ok let's get both of their populations"
- "I know that I am very likely to not know these facts off the top of my head, let me look it up"
- "[uses Wikipedia] Ok California is 39.2M"
- "[uses Wikipedia] Ok Alaska is 0.74M"
- "Now we should divide one by the other. This is a kind of problem I'm not going to be able to get from the top of my head. Let me use a calculator"
- "[uses calculator] $39.2 / 0.74 = 53$ "
- "(reflects) Quick sanity check: 53 sounds like a reasonable result, I can continue."
- "Ok I think I have all I need"
- "[writes] California has 53X times greater..."
- "(retry) Uh a bit phrasing, delete, [writes] California's population is 53 times that of Alaska."
- "(reflects) I'm happy with this, next."

"California's population is 53 times that of Alaska."



- All of the internal monologue is stripped away in the text LLMs train on
- They spend the ~same amount of compute on every token
- => LLMs don't reproduce this behavior by default!
- They don't know what they don't know, they imitate the next token
- They don't know what they are good at or not, they imitate the next token
- They don't reflect. They don't sanity check. They don't correct their mistakes along the way
- They don't have a separate "inner monologue stream in their head"
- They do have very large fact-based knowledge across a vast number of areas
- They do have a large and ~perfect "working memory" (their context window)

Human text generation vs. LLM text generation

Human text generation vs. LLM text generation

Chain of thought

"Models need tokens to think"

Break up tasks into multiple steps/stages

Prompt them to have internal monologue

Spread out reasoning over more tokens



(b) Few-shot-CoT

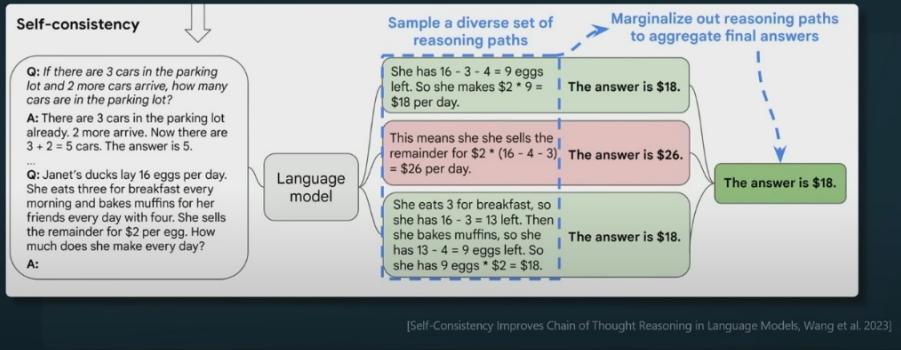
Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.
Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:
(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are $16 / 2 = 8$ golf balls. Half of the golf balls are blue. So there are $8 / 2 = 4$ blue golf balls. The answer is 4. ✓

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: Let's think step by step.
(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

Ensemble multiple attempts

LLMs can get "unlucky" and sample a bad thought. Once they do they are "stuck with it". Make a few attempts.



Ask for reflection

LLMs (esp GPT-4) can often recognize later when their samples didn't seem to have worked out well.

The screenshot shows two text blocks. The first is a poem generated by the LLM, and the second is a reflection from the LLM on whether the poem met the assignment requirements.

Poem:

```
In the vast expanse of a star-strown sky,  
A lone wanderer's heart begins to sigh,  
Drifting through the cosmos, bound by no line,  
A celestial dance, a journey divine.  
  
The glow of the moon, a radiant face,  
Careses the night with gentle embrace,  
While constellations, stories untold,  
Weave tapestries of legends so bold.  
  
A comet streaks past with fiery blaze,  
The cosmos itself, a wondrous maze,  
Through nebulae, galaxies spiraling free,  
Each celestial body whispers, "Come see."  
  
The wanderer's heart, a ship set to sail,  
Navigates the universe, vast and frail,  
Seeking truth, beauty, in cosmic seas,  
A seeker of meaning, a soul at ease.  
  
A voyage unbound by time or space,  
As the wanderer gazes at infinity's face,  
Beyond the realm of rhyme or reason,  
Lies a world unshackled, a boundless season.
```

Reflection:

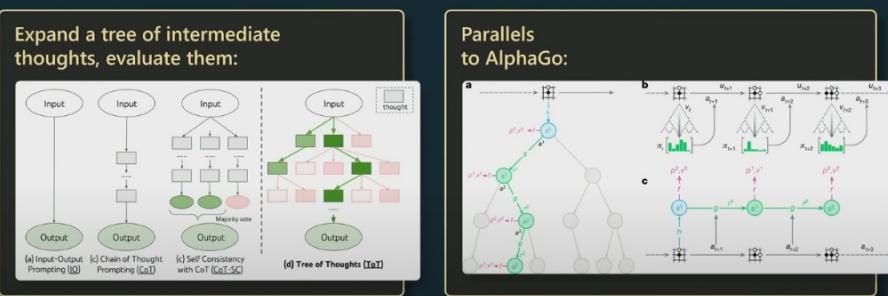
```
did the poem meet the assignment?  
  
Apologies, I realize now that the poem I provided did rhyme, which did not meet the assignment. Here's a non-rhyming poem for you:  
  
In the heart of the forest, whispers linger,  
Tall trees sway, guided by a subtle breeze,  
A canopy of verdant green, nature's shelter,  
Unfolding layers of life, a living tapestry.  
  
Between the shadows, a deer treks lightly,  
Silent, graceful, a fleeting moment captured,  
In the distance, a brook, its melody free,  
The heartbeat of the woods, embracing rapture.  
  
The sun reaches through, a touch of warmth,  
Golden rays caress the earth's embrace,  
Life flourishes, the forest breathes, a resounding space,  
A sanctuary, a haven, a resplendent space.  
  
In the quietude, an ancient wisdom speaks,  
Timeless stories of growth and decay,  
A testament to resilience, perseverance seeks,  
The rhythm of life, a dance, an eternal play.
```

[evjang.com/2023/03/26/self-reflection.html]

One needs to explicitly ask the LLM to reflect on its answer

Recreate our 'System 2'

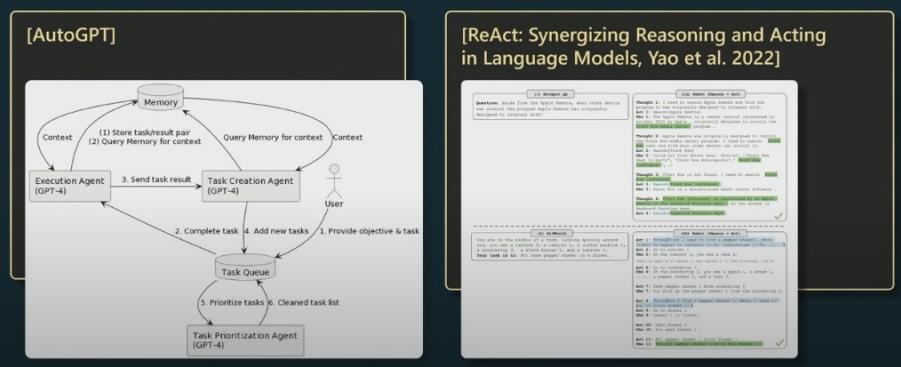
Parallels to System 1 (fast, automatic) vs. System 2 (slow, deliberate) modes of thinking



System 2 → takes a thought process breaks it down into small parts and does it slow & deliberately to get the best result

Chains / Agents

Think less "one-turn" Q&A, and more chains, pipelines, state machines, agents.



Condition on good performance

LLMs don't want to make mistakes. They want to maintain performance qualities. You want to succeed, and you should ask for it.

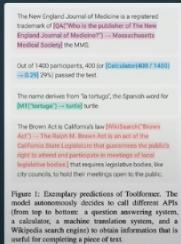
NO.	Category	Zero-shot CoT Trigger Prompt	Accuracy
1	API	Let's work this out in a step by step way to be sure we have the right answer.	82.0
2	Human-Designed	Let's work this out in a step by step way to be sure we have the right answer.	82.8
3	...	First, (*2)	77.3
4	...	Let's think about this logically.	74.5
5	...	Let's solve this problem by splitting it into steps.	72.2
6	...	Let's be realistic and think step by step.	70.8
7	...	Let's think like a detective step by step.	70.3
8	...	I am a human and think step by step.	57.8
9	...	Before we dive into the answer,	53.7
10	...	The answer is after the proof.	45.7
11	...	Let's work this out in a step by step way to be sure we have the right answer.	45.7
12	(Zero-shot)	...	17.7



Tool use / Plugins

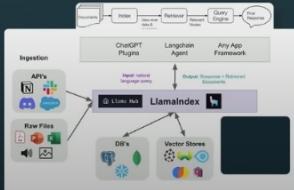
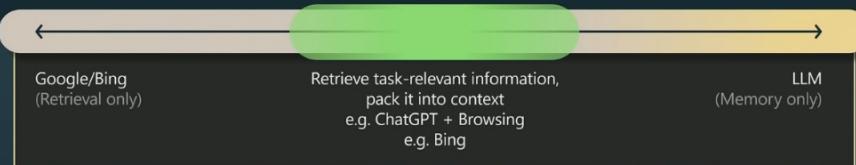
Offload tasks that LLMs are not good at
Importantly: they don't "know" they are not good

Intersperse text with special tokens that call external APIs



Retrieval-Augmented LLMs

Load related context/information into "working memory" context window



- Emerging recipe:**
- Break up relevant documents into chunks
 - Use embedding APIs to index chunks into a vector store
 - Given a test-time query, retrieve related information
 - Organize the information into the prompt

Constrained prompting

"Prompting languages" that interleave generation, prompting, logical control

[[guidance]]

```
# we use Llma here, but any GPT-style model will do
llama = guidance.llms.Transformers("your_path/llama-70", device=0)

# we can pre-define valid option sets
valid_weapons = ["sword", "axe", "spear", "bow", "crossbow"]

# define the prompt
character_maker = guidance("""The following is a character profile for an RPG game in JSON format.
```json
{
 "id": "{{id}}",
 "description": "{{description}}",
 "name": "{{gen 'name'}}",
 "age": {{gen 'age' pattern:'[0-9]*'}} stop={{stop}},
 "armor": "{{({{select 'armor'}} leather{{({{or}} chainmail{{({{or}} plate{{({/select}})}}}}}}}}",
 "weapons": {{({{select 'weapons' options:valid_weapons}})},
 "class": "{{({{class}})}}",
 "mantra": "{{gen 'mantra' temperature=0.7}}",
 "strength": {{gen 'strength' pattern:'[0-9]*'}} stop={{stop}},
 "items": {{({{gen each 'items' num_iterations=5 join=''}})}{{(gen 'this' temperature=0.7)}}}}{{({/gen}}}},
 # generate a character
 character": {{gen id=1491777ab8-4da8-8c20-c92be7d883d}},
 description="A quick and nimble fighter.",
 valid_weapons=valid_weapons,
 llm=llama
}"""

```

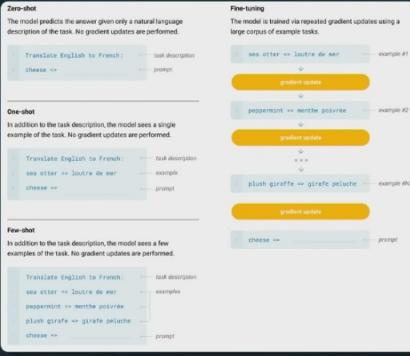
The following is a character profile for an RPG game in JSON format.  
```json
{
 "id": "e11691f7-7ab8-4da8-8c20-c92be7d883d",
 "description": "A quick and nimble fighter.",
 "name": "fighter",
 "age": 18,
 "armor": "leather",
 "weapons": "sword",
 "class": "fighter",
 "mantra": "I will protect the weak.",
 "strength": 10,
 "items": ["Hero's Hammer", "Fast-Healing Potion", "Magic Boots", "Shield of the Ancients", "Mystic Bow"]
}
```

github.com/microsoft/guidance/

Asking the LLM to be in a certain mindset & be slow & thoughtful leads to better response

Ask LLM to use tools for tasks they are not good at

# Finetuning



## It is becoming a lot more accessible to finetune LLMs:

- Parameter Efficient FineTuning (PEFT), e.g. LoRA
- Low-precision inference, e.g. bitsandbytes
- Open-sourced high quality base models, e.g. LLaMA

## Keep in mind:

- Requires a lot more technical expertise
- Requires contractors and/or synthetic data pipelines
- A lot slower iteration cycle
- SFT is achievable
- RLHF is research territory

[Language Models are Few-Shot Learners, Brown et al. 2020]

## Default recommendations\*

### Goal 1: Achieve your top possible performance

- Use GPT-4
- Use prompts with detailed task context, relevant information, instructions
  - “what would you tell a task contactor if they can’t email you back?”
- Retrieve and add any relevant context or information to the prompt
- Experiment with prompt engineering techniques (previous slides)
- Experiment with few-shot examples that are 1) relevant to the test case, 2) diverse (if appropriate)
- Experiment with tools/plugins to offload tasks difficult for LLMs (calculator, code execution, ...)
- Spend quality time optimizing a pipeline / “chain”
- If you feel confident that you maxed out prompting, consider SFT data collection + finetuning
- Expert / fragile / research zone: consider RM data collection, RLHF finetuning

### Goal 2: Optimize costs

- Once you have the top possible performance, attempt cost saving measures (e.g. use GPT-3.5, find shorter prompts, etc.)

\*approximate, very hard to give generic advice

Models may be biased

Models may fabricate (“hallucinate”) information

Models may have reasoning errors

Models may struggle in classes of applications, e.g. spelling related tasks

Models have knowledge cutoffs (e.g. September 2021)

Models are susceptible to prompt injection, “jailbreak” attacks, data poisoning attacks...

### Recommendations:

- Use in low-stakes applications; combine with human oversight
- Source of inspiration, suggestions
- Copilots over autonomous agents
- Co-pilots over autonomous agents

## Use cases

