

[Day 1] - Intro to Unit 1

"foundational LLM & text generation" Reading podcast

• Foundation = transformer arch.

→ Encoder + decoder

takes input intenc. gen. output piece by piece.

• Transformer layer

- input txt needs to be preproc into tokens.

- tokens → dense vector (embedding) which captures meaning → pos. encoding (makes structure)

• Self-attention: lets model figure out relationships between words & phrases.

1) creating query, keys & values

2) calc scores

3) normalization

4) weighted value.

general parallelizing

- prefix caching

- spec. decoding

Inference techs:

• Output approx: change output slightly (refining) - quantization - distillation

• Output pres: keep output same, optimize computation.

• Multi-head attention: mult sets of which run in parallel

• Layer normalization + residual connection

• feed-forward layer (learns funct. of input)

• Decoder Only → masked self attention + gen txt +
• Mixin of experts: MoE; experts, gating network, combination mechanism

• keeps activity level @ each level, training faster
• short cuts, way for network to remember what it learned earlier.

• directly gen output token by token (simpler & gen next)

• sequence to sequence model capable of translating seq. from one domain into seq. in another domain.

- org only contained encoder & decoder

• Pretraining → fine tuning (smaller, bias on task) → reinforcement learning (human feedback) →

• Parameter efficient fine tuning: train small part of model, makes fine tuning less expensive. → adapters based fine tuning.

• Designing input given to model to get desired output
- techs: zero-shot, few shot, chain of thought

• Sampling tech: greedy search, random samp., temp.

stop, top p, best of end.

speed up inference (inference)

quality vs. speed.

Evaluating model performance

• Pure multi. designed for pure
- consider whole system
- what does good mean?

Methods:
- human evaluation (\$\$\$)
- AI (determinator)
- need to be calibrated.

Transformer

Training

Parameter Eff.

Prompt Engineering

Sampling

[Day 1] - Unit 1 (comp.)

"Prompt Engineering"
Podcast (Reddit)

Sampling

- Config LLM output tail output length (affects costs / pr. time)
 - low token limit.
 - engineer prompt to be targeted to get concise response (ex: react technique)
- How LLM chooses next word based on prob.
- **Temp:** "dial" it randomness
 - lower: pushes model towards most likely word
 - * predictable, deterministic output. **Precision**
 - higher: introduces randomness, diverse results.
 - + exploring new ideas
- **Top k / Top P (nucleus sampling):** Limit the next word to most prob. candidate
 - ↑ k = variety in output - ↓ k = focus - k = budget
 - **Top P (prob.):** selects smallest set of words whose value adds up to P. (lower: n/more; higher: all pos)
- Filter w/ Top k & P, temp makes final choice
 - * Ideals: temp = 2, top P = 0.95, k = 30
- Model gets stuck repeating phrase, happens @ low/high temp
 - * too predictable or too random

Rec. Loop Buy

Prompt Eng. Techs

Automatic prompt eng: ai automation
creation of effective prompts.
- effective for writing code.

Multi-modal prompting: Using inputs like
image/audio

Best Practices:
1) provide examples - design
2) be specific, use instructions
3) max token length!

- **General (zero shot):** giving model task descpt input w/o exs
- **Doc. prompts:** keeping track of what worked / didn't
- **One/hw shot:** provide examples within prompt to guide model
 - quality of examples is crucial (edge cases)
- **System contextual + role prompting:** giving model all details.
 - ↳ context + purpose
 - ↳ giving LLM personal (clear perspective)
- **Step back prompting:** starts w/ asking broader q, before specific
 - + mitigate biases
- **Chain of thought:** boosts reasoning capabilities by prompting model to generate reasoning steps. (uses more tokens)
- **Self-consistency:** get model to generate more reasoning paths for same prompt, choose most consistent answer. (expansion)
- **The of Thought:** LLM explores mult. reasoning paths
- **React:** reason & act, combining LLM reasoning cap. w/ ability to use external tools (search engine)

Codelabs

[prompting fund.]

Output length: affects cost & performance.

- more tokens = ↑ computation (↑ energy, latency, cost)

Temperature: controls degree of randomness in token selection. ↑ temp = ↑ # of candidate tokens from which next output is selected. = diverse results.

- temp @ 0 = greedy decoding (no diversity)

TOP P: defines prob. threshold that once cum. exceeded tokens stop being selected as candidates.

0 → greedy decoding ; 1 → every token selected in top N words.

Prompting: zero shot (direct), enum (label assist), chain of thought (intermediate reasoning steps), reason & Act, thinking mode (generative thinking process)

Define an evaluator: "instruction following, groundness, fluency, verbosity or quality"

• UFG: quickly train through a prompt w/ small set of test documents, compare diff models, unify system does not regress

• Pointwise eval: eval a single input/output pair against criteria - useful for eval. singular outputs in absolute sense (good/bad)

• Pairwise eval: comparing 2 outputs against each other

• LLM limits: LLMs have problems on certain tasks, an evaluator will not be able to accurately eval. this task.

- solutions: conducting focus, understand limitations, include human evals to calibrate & determine baseline, not relying on "internal knowledge" (provide w/ input / tool)

• Improve confidence w/ mult. evals (disc), sum prompts.

• GPU mem needed for inference on 3B param. model: 12GB

• Reinforcement learning from human FB: reward model to incentivize human - support context w/ up to 2m tokens

• human config that controls randomness: temp.

• Models w/ no try to pull in more requests from user

• UX/UI

• AI system design: prompt caching, on device vs. ping-pong a server.

• prompt eng: input, temp.,

• evaluation - crucial

Evaluation & Structured output

Challenges

Conclusion