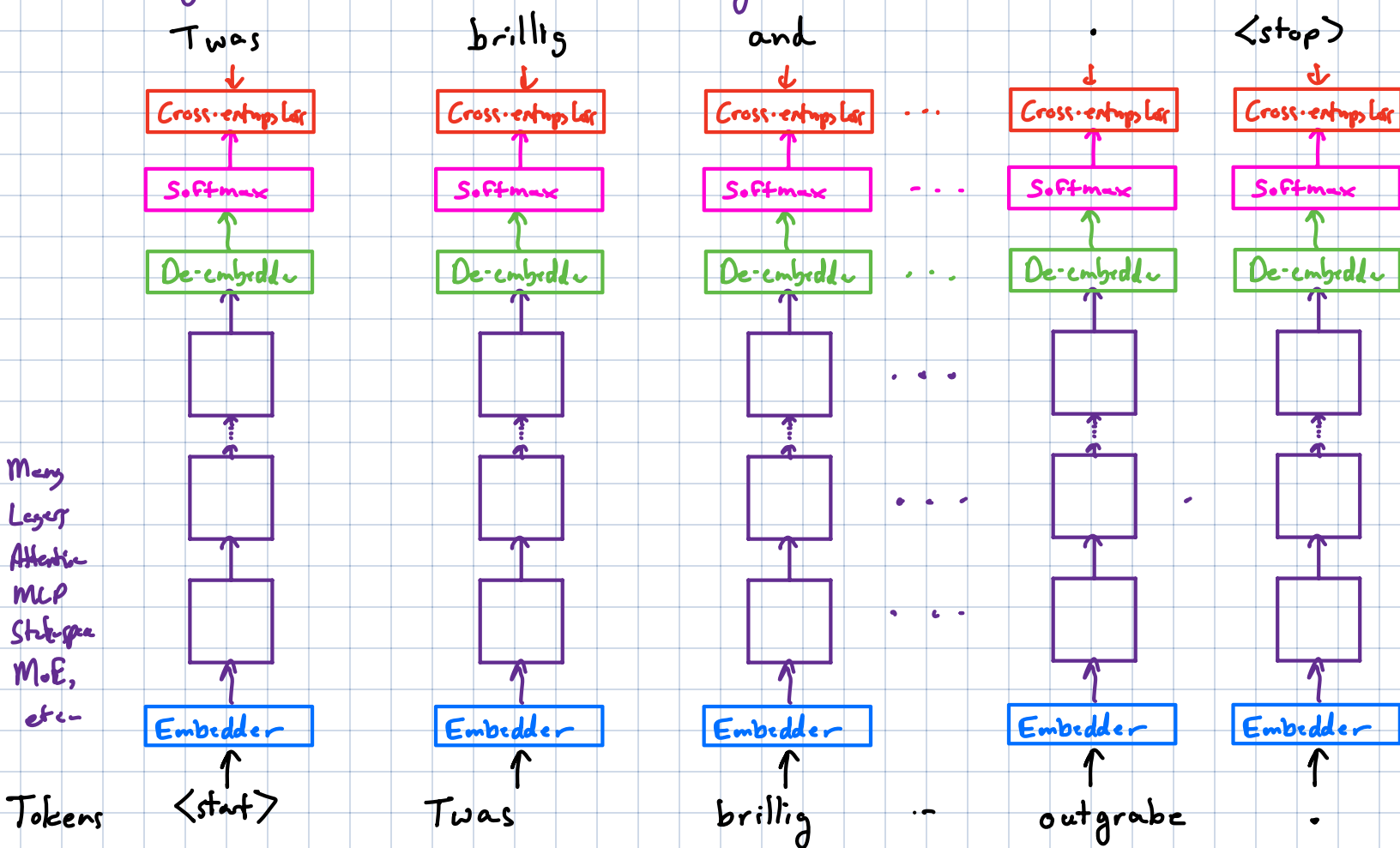Today: ICL/Prompting
Fine-tuning
Parameter-efficient Fine-tuning
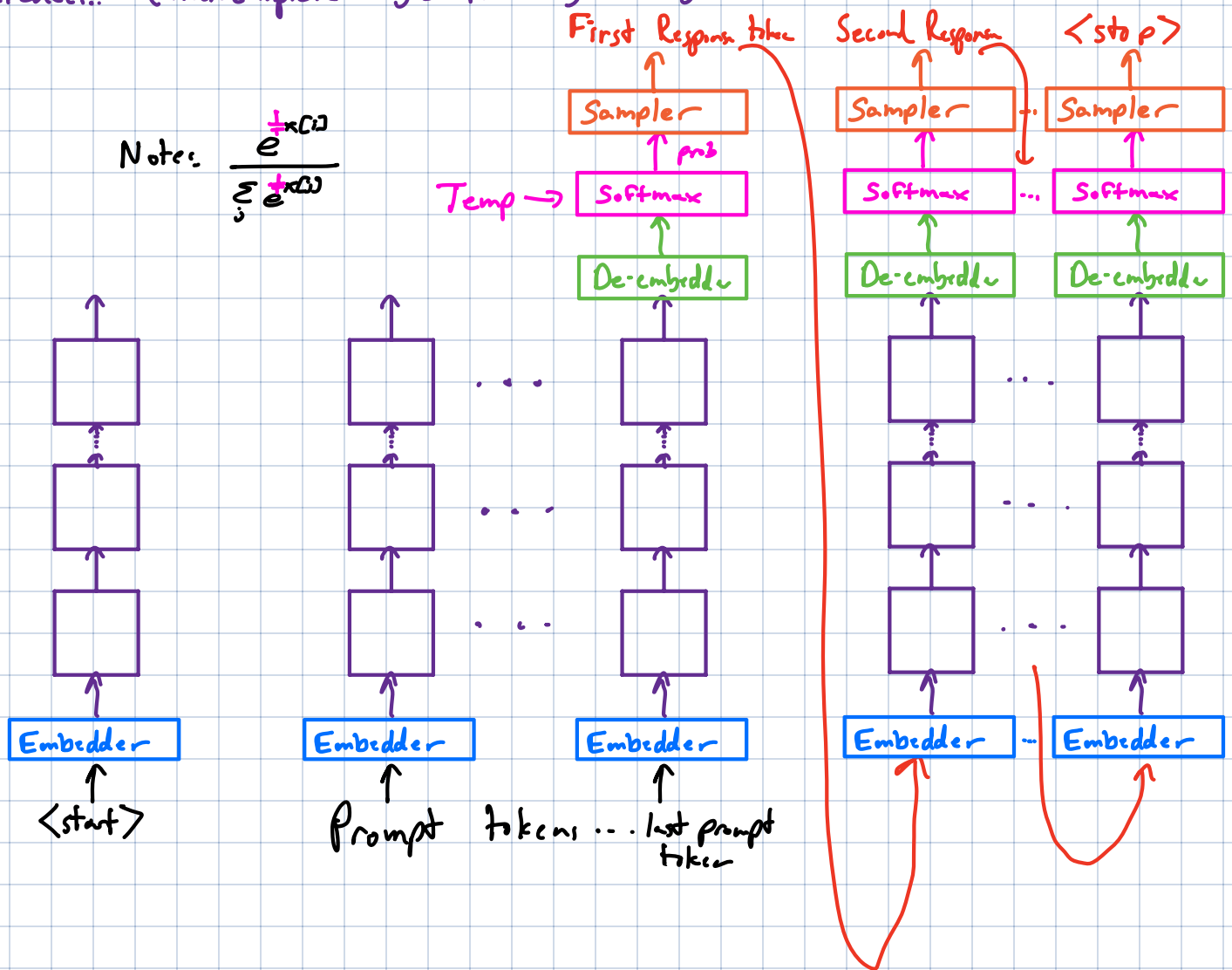
Announce: Start on Projects
Compute Offered
(NAIRR & Tinker)

GPT-style models... Recall Pretraining: (Next token Prediction)



| Twas | brillig | and | | · | \<stop\> |

Tokens: \<start\> | Twas | brillig | .. | outgrabe | .

Many
Layers
Attention
MLP
Stitcherper
MoE,
etc~

Data mix matters.

vs Inference... (Autocomplete Style Autoregressive generation)

First Response token    Second Response    <stop>



Notes: $\dfrac{e^{\frac{1}{T} \times C_i}}{\sum_j e^{\frac{1}{T} \times C_j}}$

Temp →

Some history:

GPT-1 : Next-token prediction gives good embeddings
(2018)

GPT-2 : With bigger models, AR completions result in natural prose/poetry
(2019)  and the model knows facts about the world that can be
elicited by the right prompt. ("Zero-shot learning")

Example: The capital of France is Paris.

GPT-3 : With even bigger models, also get in-context learning. ("Few Shot" learning)
(2020)

Example: Red: Ruby ; Blue: Saphire
        vs
        Red: Raspberry; Blue: Blueberry.

Revolutionary... Can now solve a learning problem with a few examples without any fine-tuning or gradient descent. Pure Prompting.

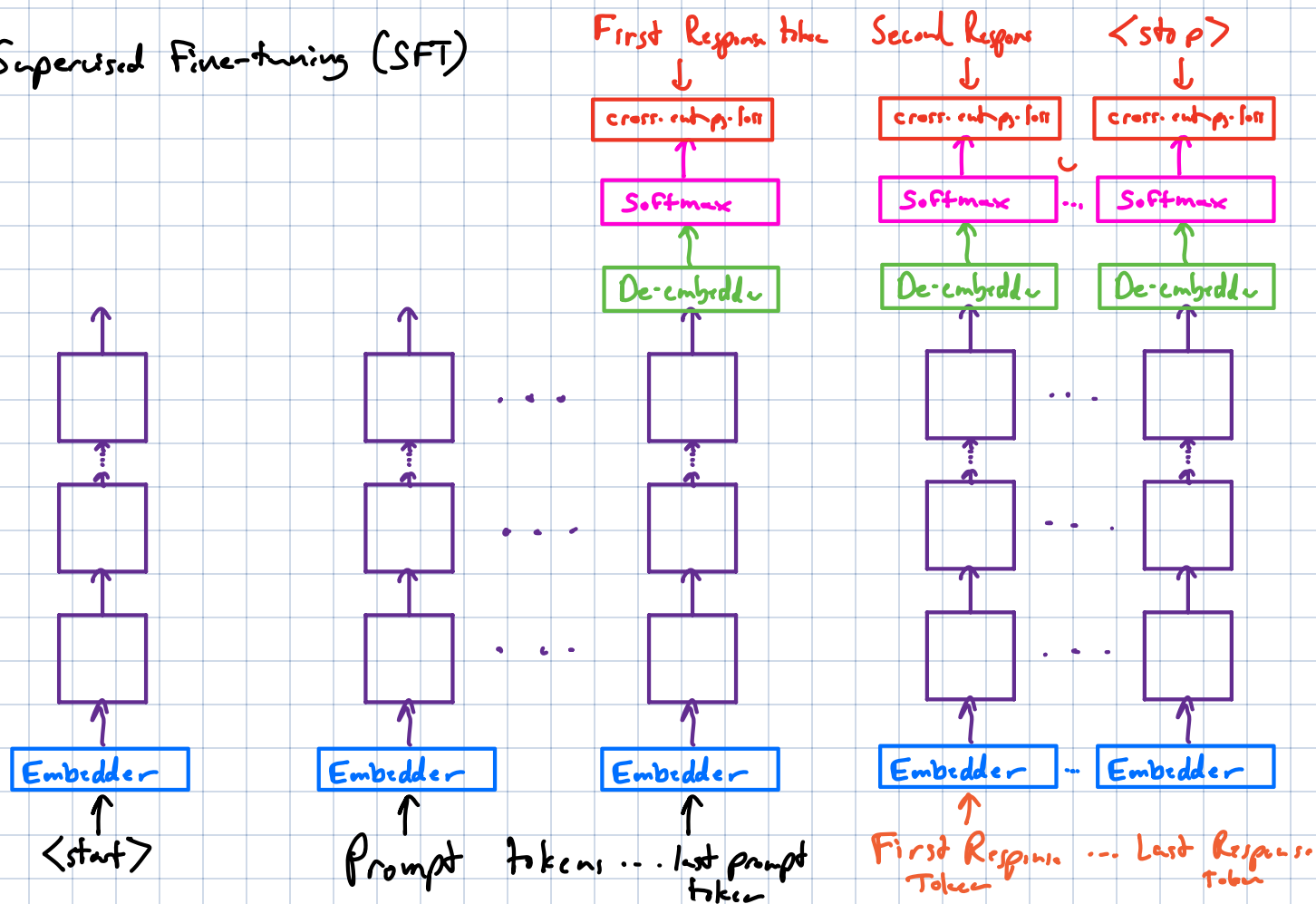Challenge: all training data must fit in context.

Other natural questions: (A) How can we find better prompts?

(B) How can we make more usefully promptable models?

Start with (B): Critical Paper: Instruct GPT (2021)

# Fine-tuning for Instruction Following or Question Answering...

Supervised Fine-tuning (SFT)

First Response token    Second Response    <stop>

cross. entropy. loss    cross. entropy. loss    cross. entropy. loss

Softmax    Softmax  ...   Softmax

De-embedder    De-embedder    De-embedder

Embedder    Embedder  ...  Embedder    Embedder  ..  Embedder

<start>     Prompt tokens ... last prompt token    First Response Token  ...  Last Response token

In 2023, Llama model released. ⟶ Alpaca (showed Instruction following examples) are key

Dolly ⟶ showed it was enough examples that moderate

Step Back: Fine-tuning or using a pretrained model for a new task:

0) Pure prompting — No gradient-descent steps

1) Treat as embedding — a feature extractor. Do classical ML to train.
   Also called "linear probing"                          a separate model

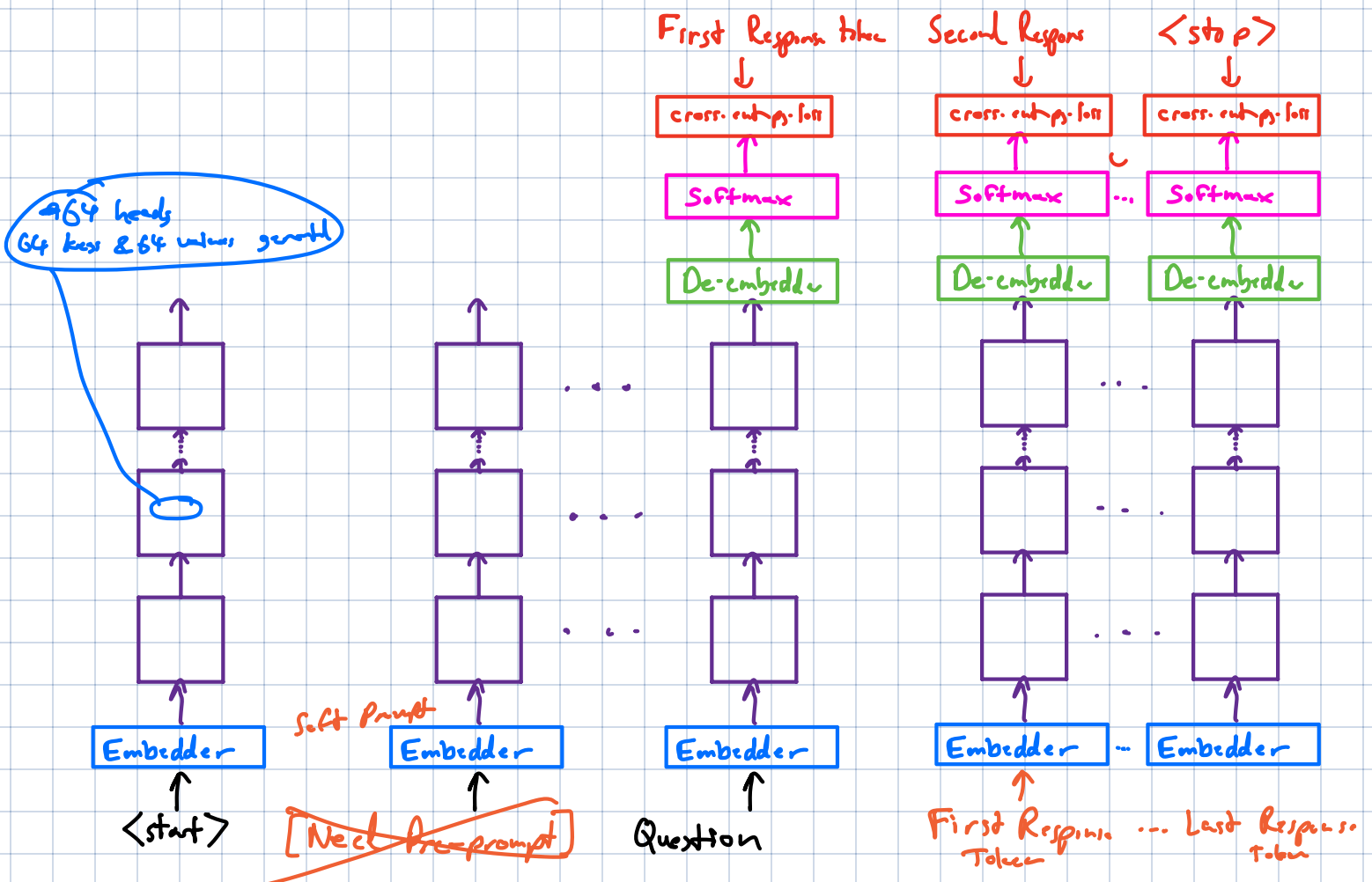2) Full fine-tune    — adjust the weights of the pretrained model


Return to (A): How to get better prompts?

Aside: If you want to do this, and only have API access to the model, tools like
       DSPy are what you want to use.
       New algorithms inside like GEPA (this summer) that use LLMs to
       reflect on the execution of prompts (related to Text Grad in spirit)
                                         "textual gradients"

       Genetic Algorithm Style...


Assume we have model...



First Response token    Second Response    <stop>
        ↓                      ↓               ↓
   cross. entropy. loss    cross. entropy. loss    cross. entropy. loss

   Softmax              Softmax     ...,   Softmax

   De-embedder          De-embedder        De-embedder

64 heads
64 keys & 64 values, generate

Embedder      Embedder      Embedder      Embedder ~ Embedder

<start>    Soft Prompt        Question      First Response ... Last Response
           [Next Free prompt]                Token              Token

To be precise: If Pre-prompt is length 100

Call it $\langle special_1 \rangle$ ... — — — $\langle special_{100} \rangle$

when the Embedder has 100 new columns added for

These 100 new columns are learnable parameters

# Soft Prompts & Prefixes work pretty well.

And far far fewer parameter than the entire model.

example for a soft-prompt, Embedding 4096 (Llama 3, Olmo 2, etc.)

Prompt-length 100

Total Params: 40K vs 7 or 8 Billion

Soft Prefix, (Just the k-v cache) Assume 32 layers

$$100 \times \underbrace{(4096}_{keys} + \underbrace{4096)}_{values} \times \underbrace{32}_{layers} = 262K \times 100 = 26M$$