

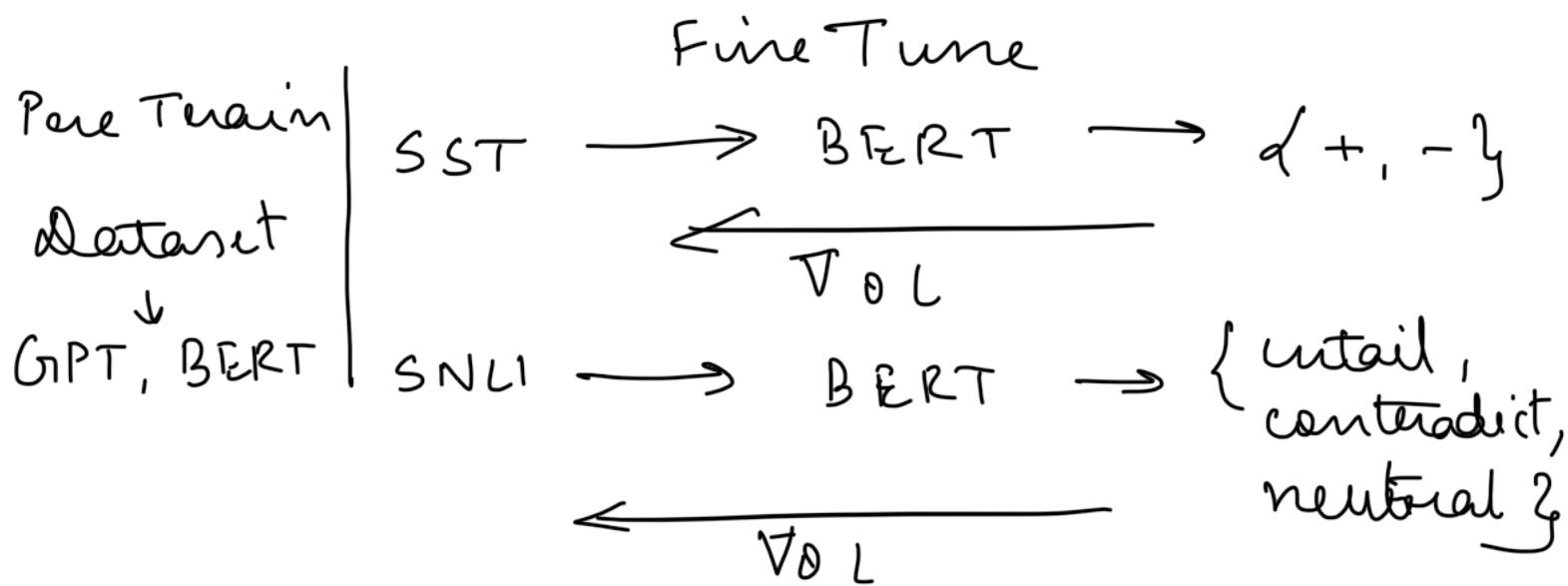
## Week-4 DLP

### L1 Task Specific Fine Tuning & It's Limitations

- Trained the GPT-2 model with a CLM training objective.

$$\text{Min } L = - \sum_{D} \sum_{i=1}^T y_i (\log(\hat{y}_i))$$

- But, we have to refine it for the downstream tasks like "classification", "text genera" & conversa".
- earlier, we had to fine tune for each downstream tasks.



- finetuning → makes the model to adapt for downstream tasks.
- Initialize the params with the params learned by solving pre-trained objectives
- At the input, add additional tasks based on type of downstream Task. At output replace the pre-training LN head with classification head of a linear layer.

- Our obj. to predict the label  
 $\hat{y} = P(y|x_1, \dots, x_m) = \text{softmax}(W_y h_e^m)$   
 $\text{min. } L = - \sum_{(x_i, y_i)} \log(\hat{y}_i)$
- Now, the pre-trained model acts as a feature extractor & the classification head acts as a simple classifier.  
 Do you have to train all the params of model known as full fine-tuning.
- GPT-2 small  $\rightarrow$  3 to 4 GB of GPU  
 of batch size 1 to train the model.
- GPT-2 large  $\rightarrow$  > 22 or 24 GB GPU to train the model.
- So think of memory engg. to fine-tune huge models?

## L2 Emerging Abilities

- There are many approaches used so far to get memory engg. to fine-tune LM.
  - Additives, Selectives, Adapters, Soft Prompts, Reparameterizing-based
- PEFT  $\rightarrow$  Parameter efficient fine tuning

- In PEFT we have LORA, GLoRe, AdaLoRa are most used (+ quantization)
- even for downstream fine tune you need lots of data points. Sometimes you don't have enough labelled samples. But what about humans, we don't see samples, right?
- ∴ we don't need supervised fine-tuning at all for most of the tasks.

Currently,  $p(\text{output} \mid \text{input})$

Now,

$p(\text{output} \mid \text{input}, \text{task})$

task is just an instance in plain needs to the input seq.

- In this setting, to get a good performance, we have to scale up both model size & data size.
- push the limits: 1.5B to 175B. The ability to learn to learn from few examples ↑ as model size increases.
- Since the model learns a new task from samples within the context window, this approach → in context learning.
- the adapt occurs during inference, there is no need to share model wts. for

fine tuning. This method led to paucity of prompt engineering.

↳ 0-shot, few-shot, chain of thought, prompt chaining.

### L3 Instruction & Preference Tuning

- 0 shot learning performance is often poor → whatever be the model size. Fine-tune the model on the instruction (1 approach is to reformat the available datasets into instruction sets)  
↳ refer to FLAN dataset.
- We have to fine tune the model to align with user's intent. For this we need human labeller for a collection of prompts. Use it to fine-tune supervised fine tuning, the model using RLHF.

#### Fine-Tuning

SFT

task-specific full  
instruction tuning  
preference tuning <sup>RLHF</sup>  
memory efficient  
(PEFT, Quantization)

Prompt

0 shot, few  
shot, chain of  
thought, prompt  
chaining, meta  
prompting.

- till now, we have trained the model  
now we need to fine tune it (HF  
will help you → peft, teel, SFT, bits  
and bytes, unslotter)

t4-l6 → demos