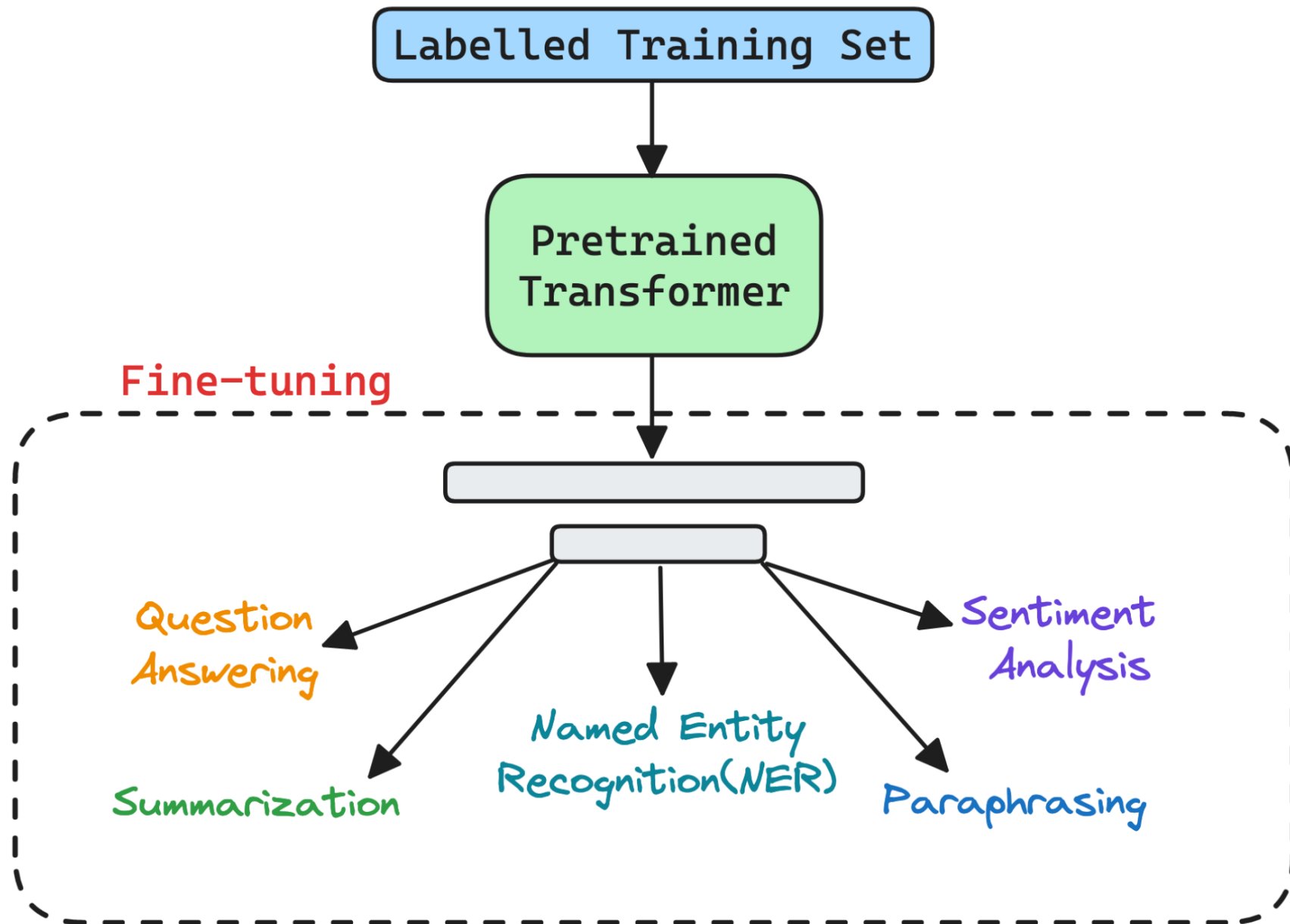
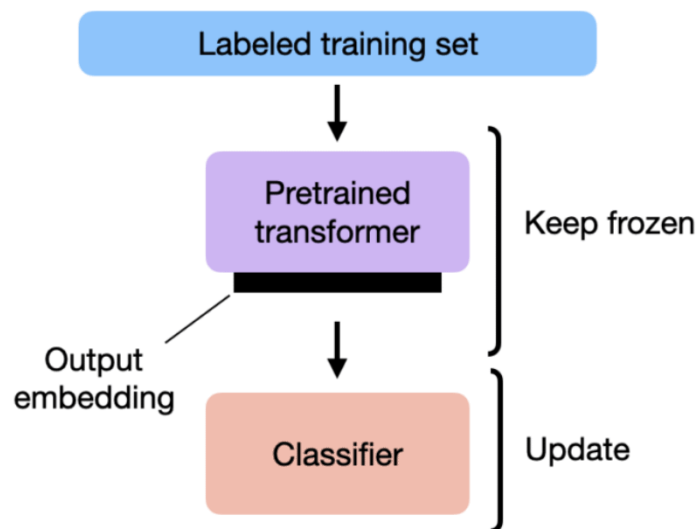


Fine-tuning an LLM! 🚀

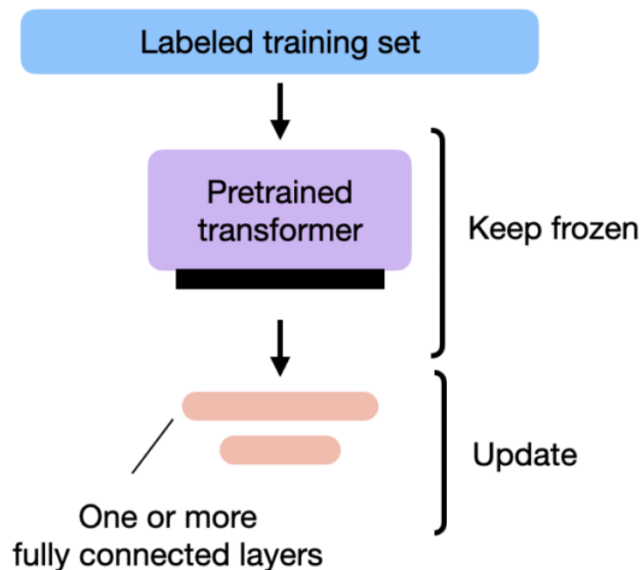


Fine-tuning Strategies

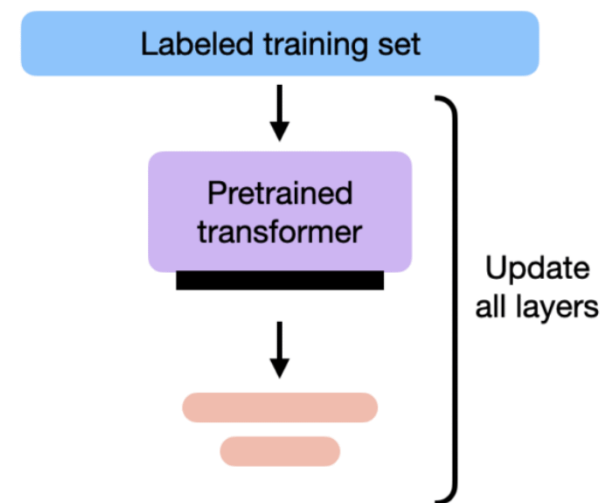
1) FEATURE-BASED APPROACH



2) FINETUNING I



3) FINETUNING II

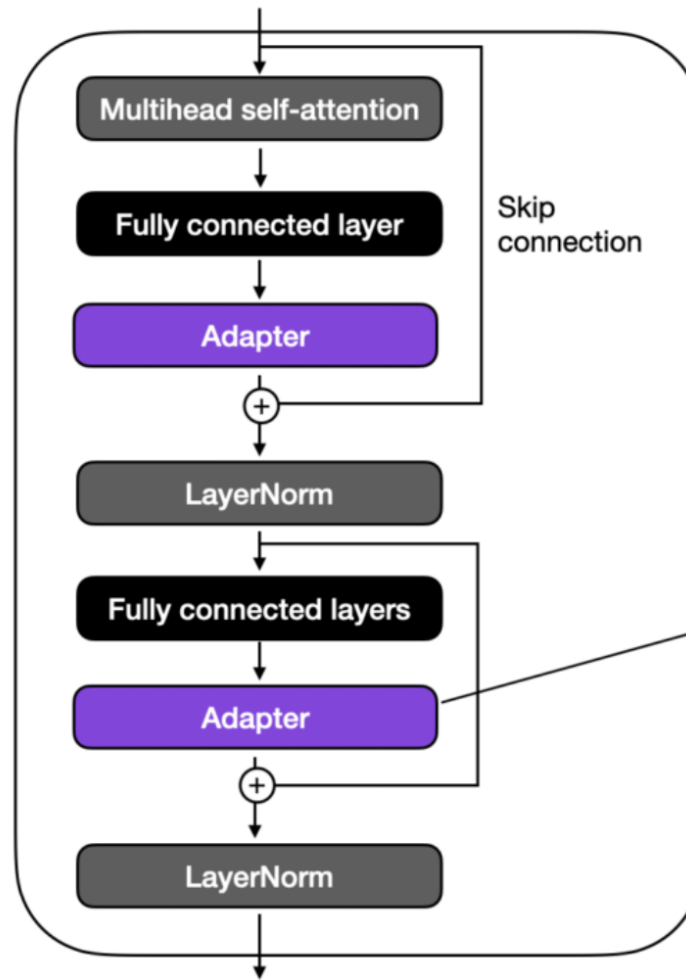
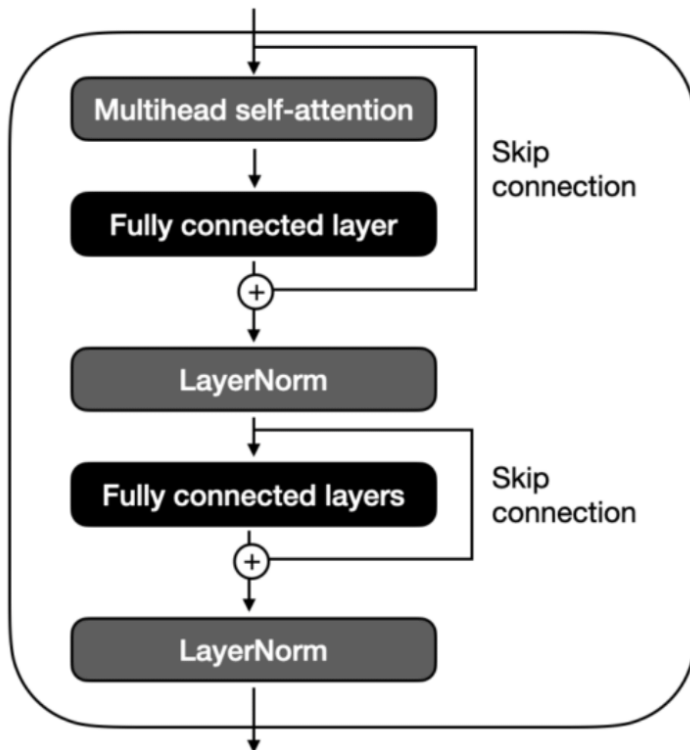


Transformer with Adapters

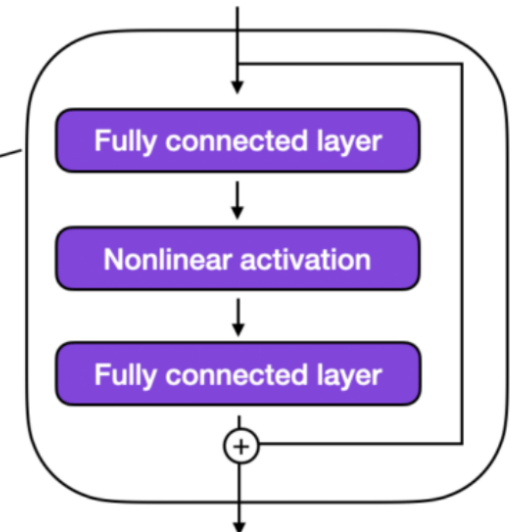
Only the adapter layers in each block are trainable

TRANSFORMER BLOCK WITH ADAPTERS

REGULAR TRANSFORMER BLOCK



ADAPTER LAYERS



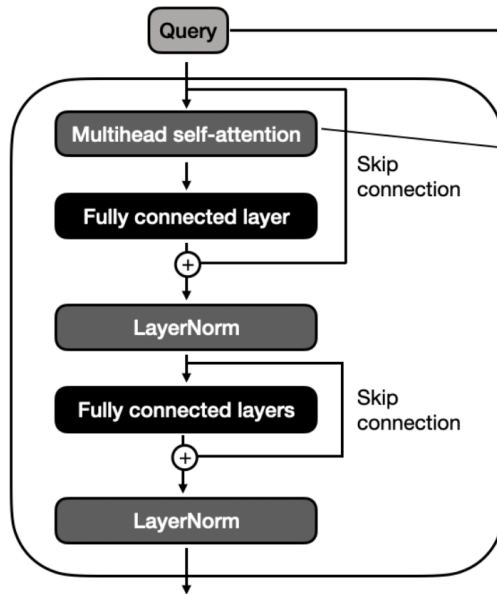


adapter.py

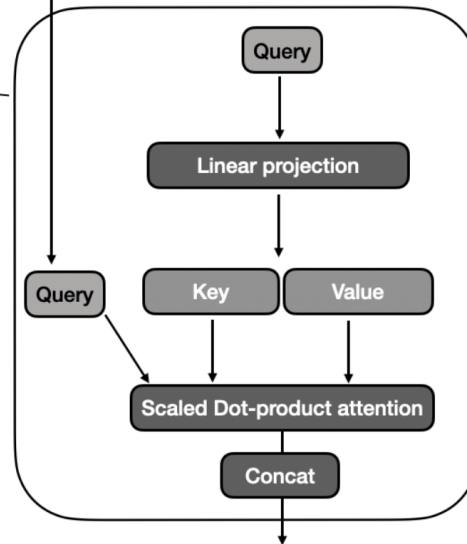
 @akshay_pachaar

```
def transformer_block_with_adapter(x):  
    residual = x  
    x = self_attention(x)  
    x = AdapterLayers(x) # adapter  
    x = LayerNorm(x + residual)  
    residual = x  
    x = FullyConnectedLayers(x)  
    x = AdapterLayers(x) # adapter  
    x = LayerNorm(x + residual)  
    return x
```

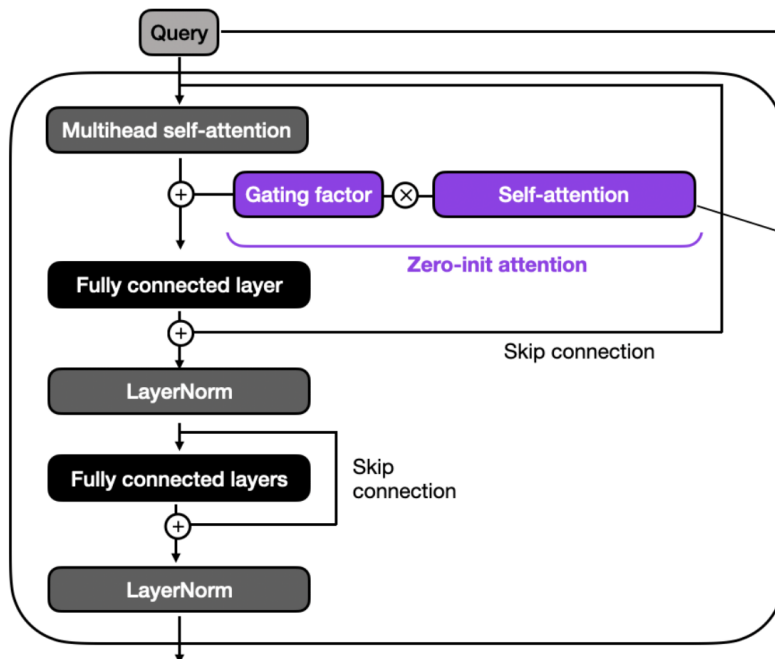
REGULAR TRANSFORMER BLOCK



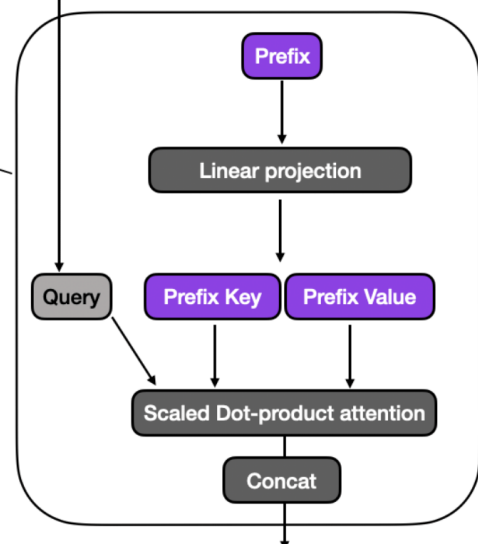
REGULAR SELF-ATTENTION




TRANSFORMER BLOCK WITH LLAMA ADAPTER




SELF-ATTENTION WITH PREFIX-MODIFIED KEY AND VALUE SEQUENCES












 llama_adapter.py

```
def transformer_block_with_llama_adapter(x, gating_factor, soft_prompt):
    residual = x
    y = zero_init_attention(soft_prompt, x) #llama-adapter: prepend prefix
    x = self.attention(x)
    x = x + gating_factor * y #llama-adapter: apply zero_init_attention
    x = LayerNorm(x + residual)
    residual = x
    x = FullyConnectedLayers(x)
    x = LayerNorm(x + residual)
    return x
```

 @akshay_pachaar

That's a wrap!

If you interested in:

- Python 
- Data Science 
- Machine Learning 
- MLOps 
- NLP 
- Computer Vision 
- LLMs 

Follow me on LinkedIn 

Everyday, I share tutorials on above topics!

Cheers!! 



@akshay_pachaar