

LLM Engineering

MASTER AI & LARGE LANGUAGE MODELS



Training Results

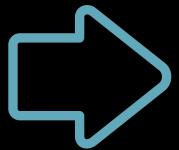
What you can now do

- Generate text and code with Frontier Models and Open Source models using APIs and HuggingFace, including tools, assistants and RAG
- Follow a 5 step strategy to solve problems, including dataset curation, making a baseline model, and fine-tuning a Frontier model
- Run QLoRA for fine-tuning open-source models including defining and choosing hyper-parameters and running and monitoring training

By end of this session you'll be able to

- Explain how Training works
- Run inference on a QLoRA fine-tuned model
- Confidently carry out the end-to-end process for selecting and training open source models to solve a business problem

The Four Steps in Training



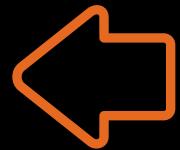
Forward pass

Predict the next token in training data



Loss calculation

How different was it to the true next token



Backward pass

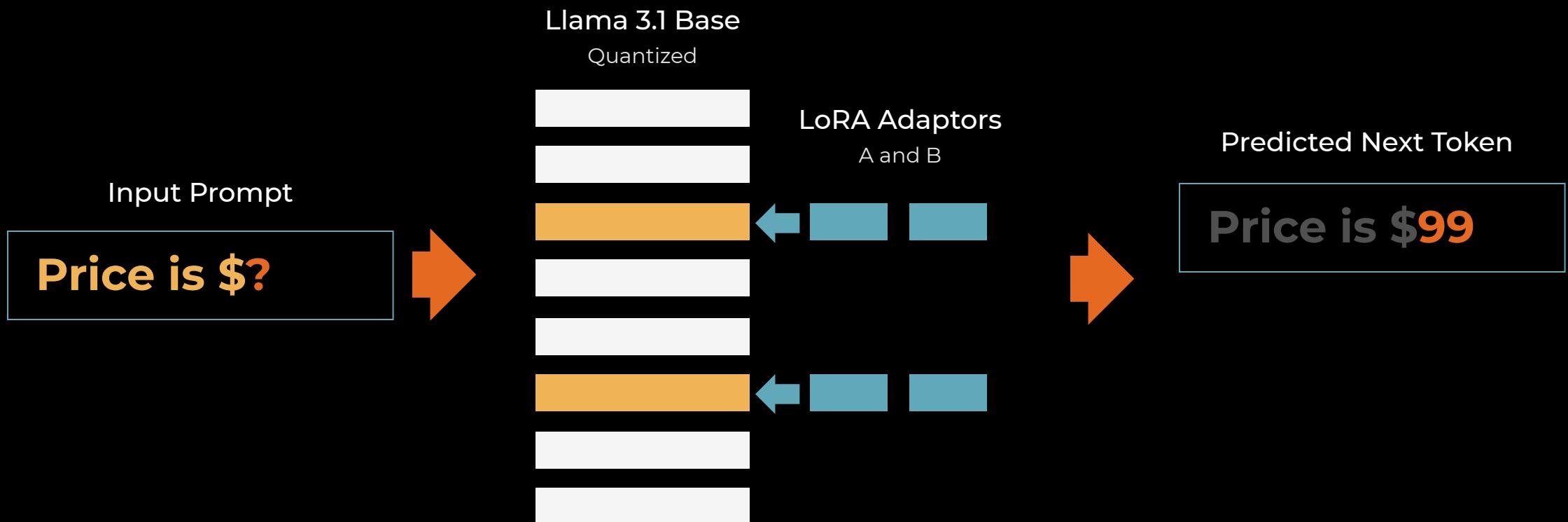
How much should we tweak parameters to do better next time (the "gradients")



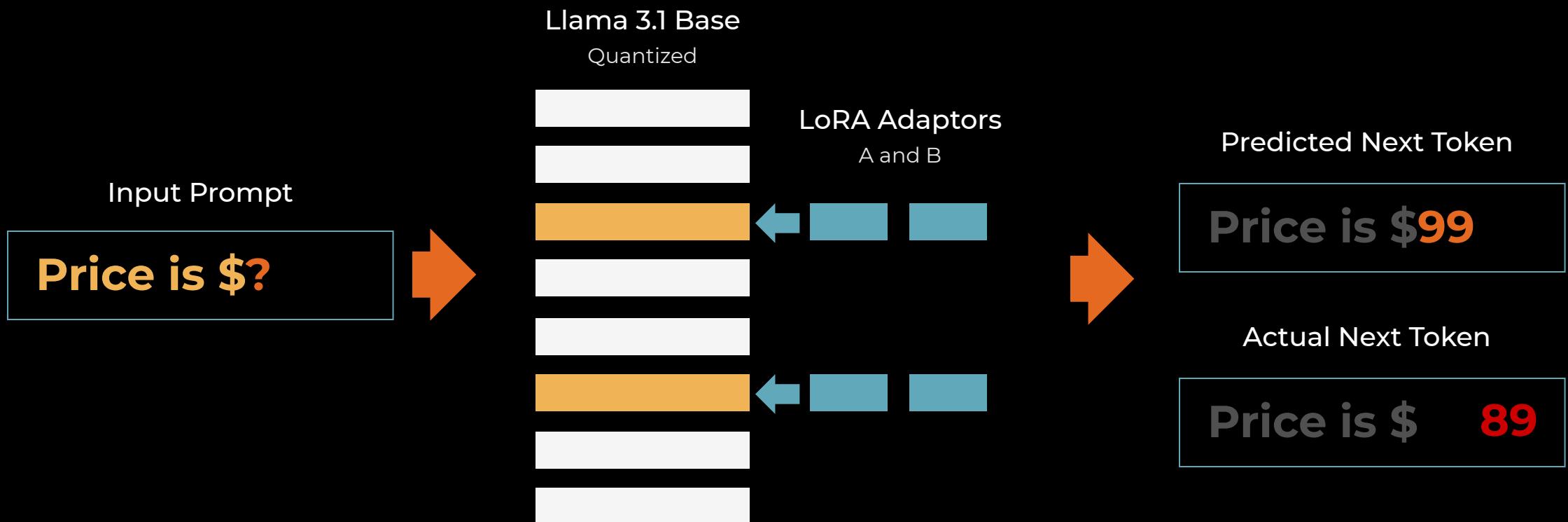
Optimization

Update parameters a tiny step to do better next time

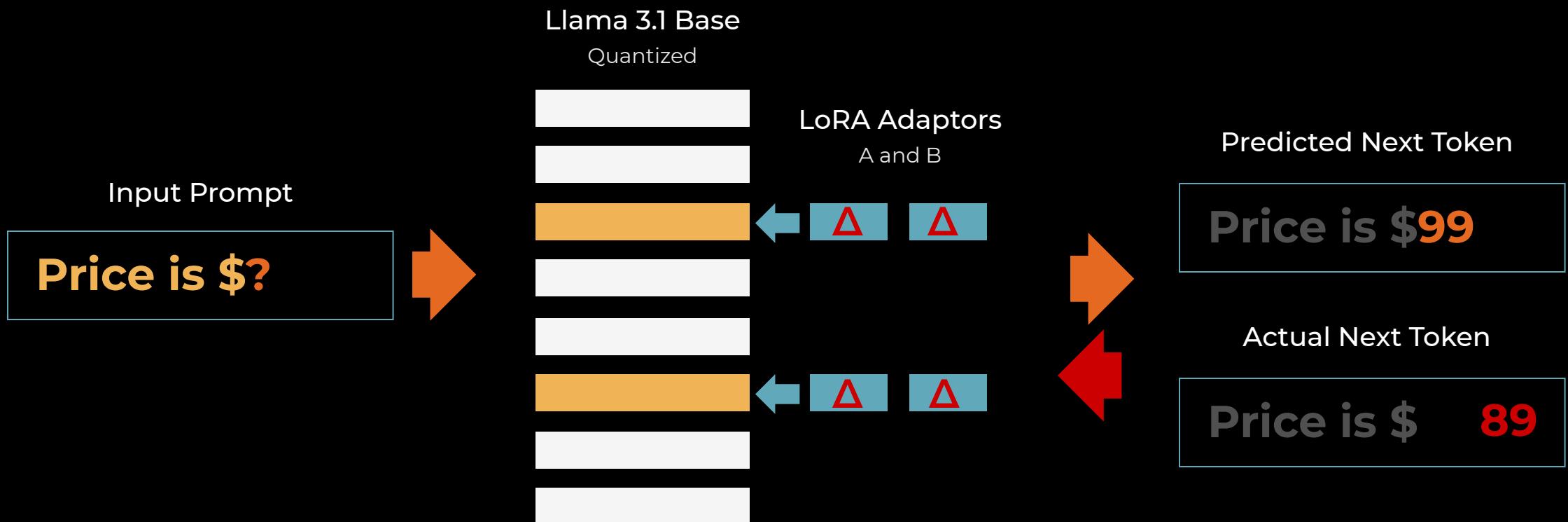
1. The Forward Pass



2. The Loss Calculation

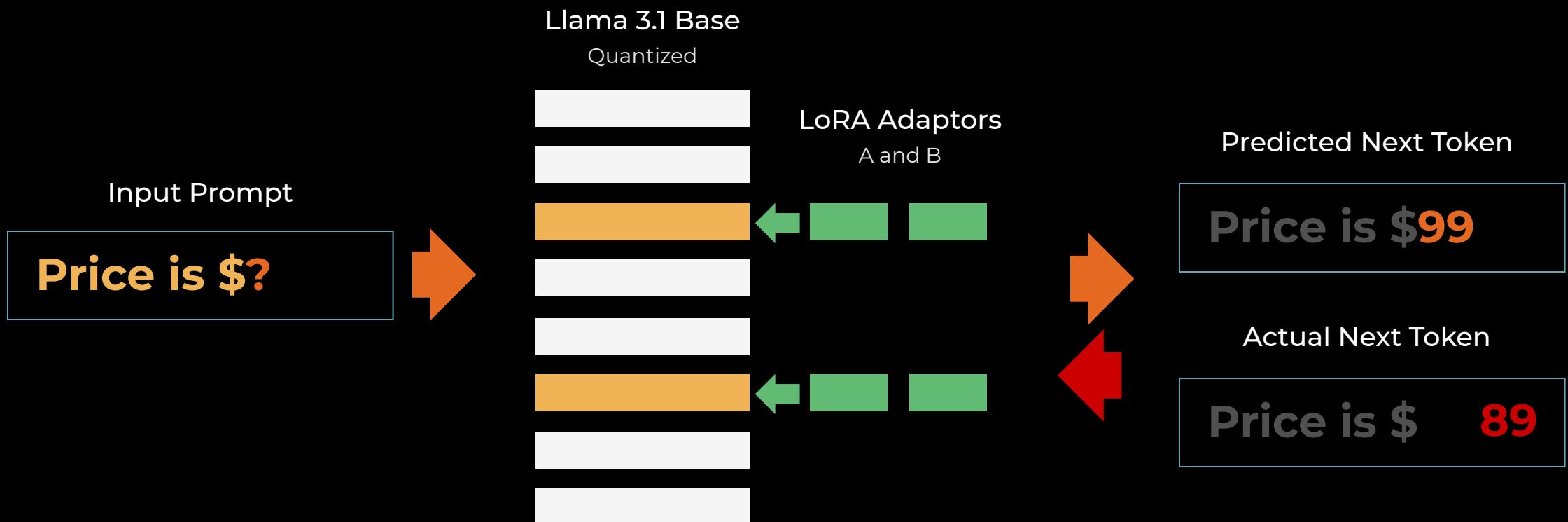


3. The Backward Pass ("Backprop")



4. Optimization

Shift weights a tiny amount (the Learning Rate) for a slightly higher chance of predicting the right token next time



An important technical detail

Next Token Prediction and Cross Entropy Loss

The Model Output

- The model doesn't simply "Predict the next token"
- Rather, it outputs the probabilities of all possible next tokens
This is the result of using the 'softmax' function over the output from the last layer
- During inference, you can pick the token with highest probability, or sample from possible next tokens

The Loss Function

- The approach for calculating loss is quite simple:
 - Just ask: what probability did the model assign to the token that actually was the correct next token?
 - In practice we then take the log of this probability and times by -1
 - So 0 means we were 100% confident of the right result; higher numbers mean lower confidence
- This is called **cross-entropy loss**



A reminder of where we are



RESULTS!





PROGRESS

I've kept the best to last

What you can now do

- Generate text and code with Frontier Models and Open Source models using APIs and HuggingFace, including tools, assistants and RAG
- Follow a 5 step strategy to solve problems, including dataset curation, making a baseline model, and fine-tuning a Frontier model
- Confidently carry out the end-to-end process for selecting and training specialized open source models that can outperform the Frontier

Next week is the finale - you will be able to

- Deploy customized models behind an API
- Create production products that use custom models
- Create an end-to-end solution to a commercial problem with groundbreaking LLMs