

# INTRO TO INFERENCE

How to Run AI Models on a GPU



A Google Cloud x NVIDIA Learning Path





Dr. Nova, I've been using ChatGPT and LLMs all day... but I have no idea how they work behind the scenes.

Great question, Kai!  
Every AI model has two phases: Training and Inference.

Training is when the model learns from data — adjusting its weights through millions of examples.

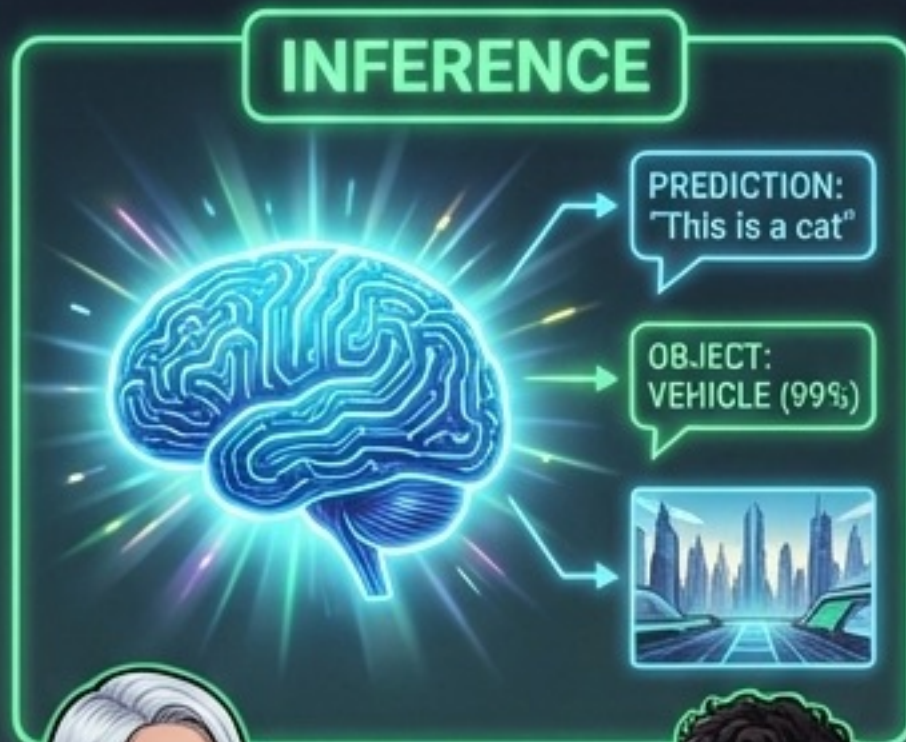
## TRAINING



Think of training as studying for an exam...

...and inference is taking the exam. The model applies what it learned.

## INFERENCE



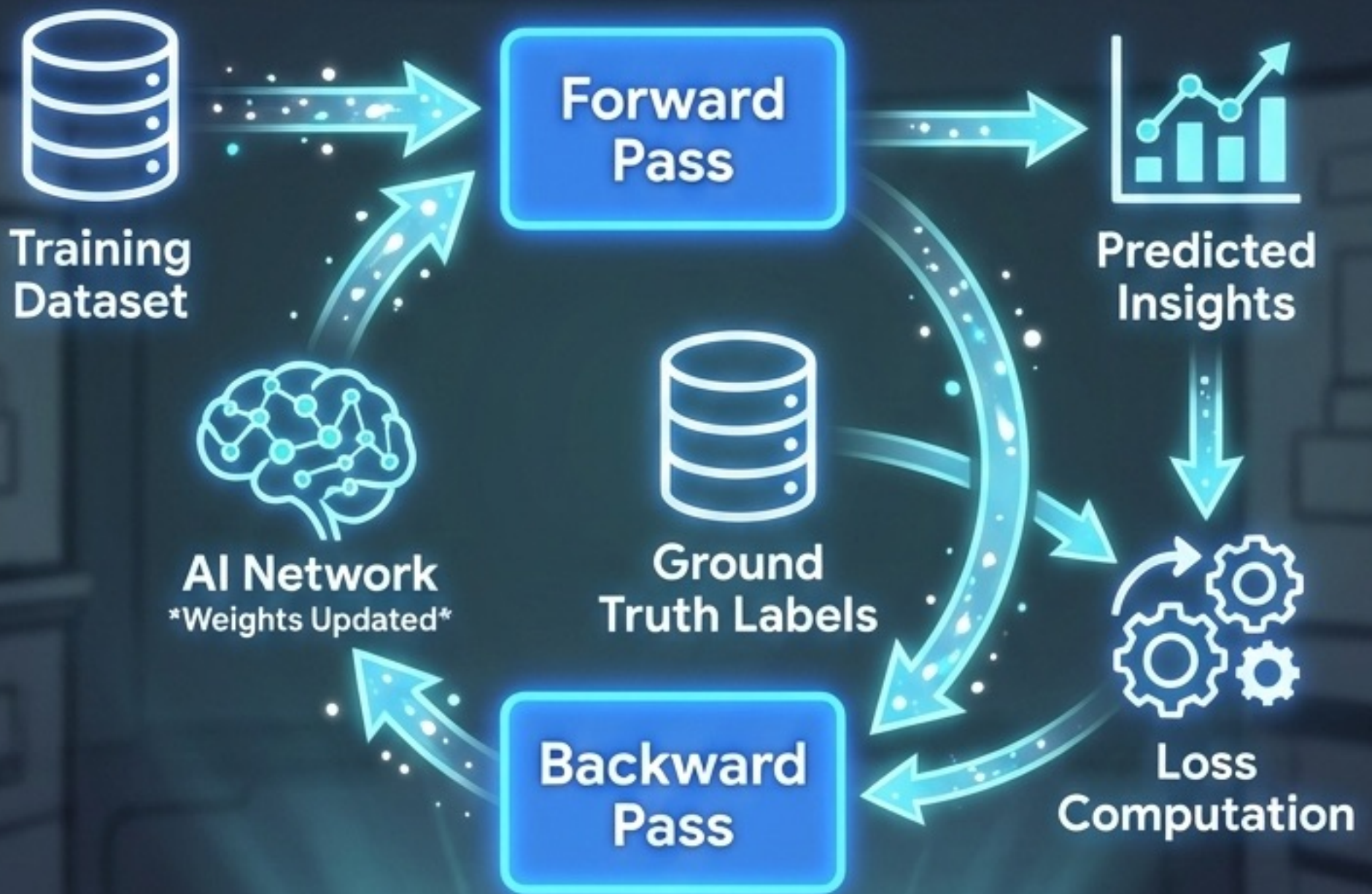
So inference is the part users actually interact with!



Training has two key steps: a Forward Pass and a Backward Pass.



The forward pass makes predictions. The backward pass compares them to the truth, computes the error, and adjusts the weights.



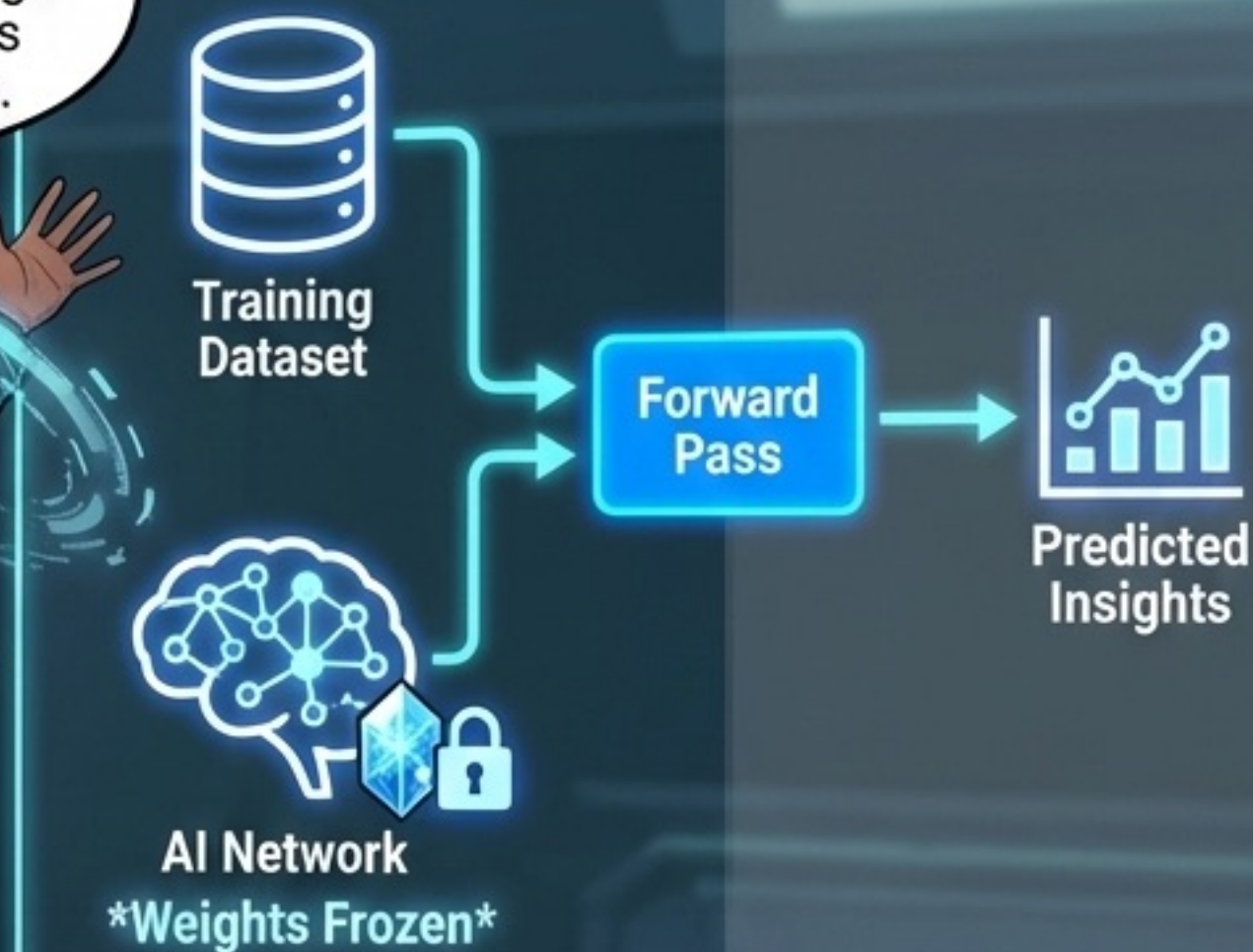
So the model keeps going around this loop, getting better each time?

Exactly! But once training is done, inference is much simpler...





Inference is just the forward pass. No backward pass, no weight updates. The model's knowledge is frozen.



But if millions of people are using an LLM at once, how does it handle all that?



Those are the three big questions of serving AI in production. And that's where GPUs come in!

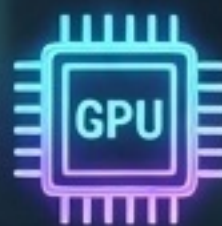
How **FAST** can we respond?



How **MANY** users at once?

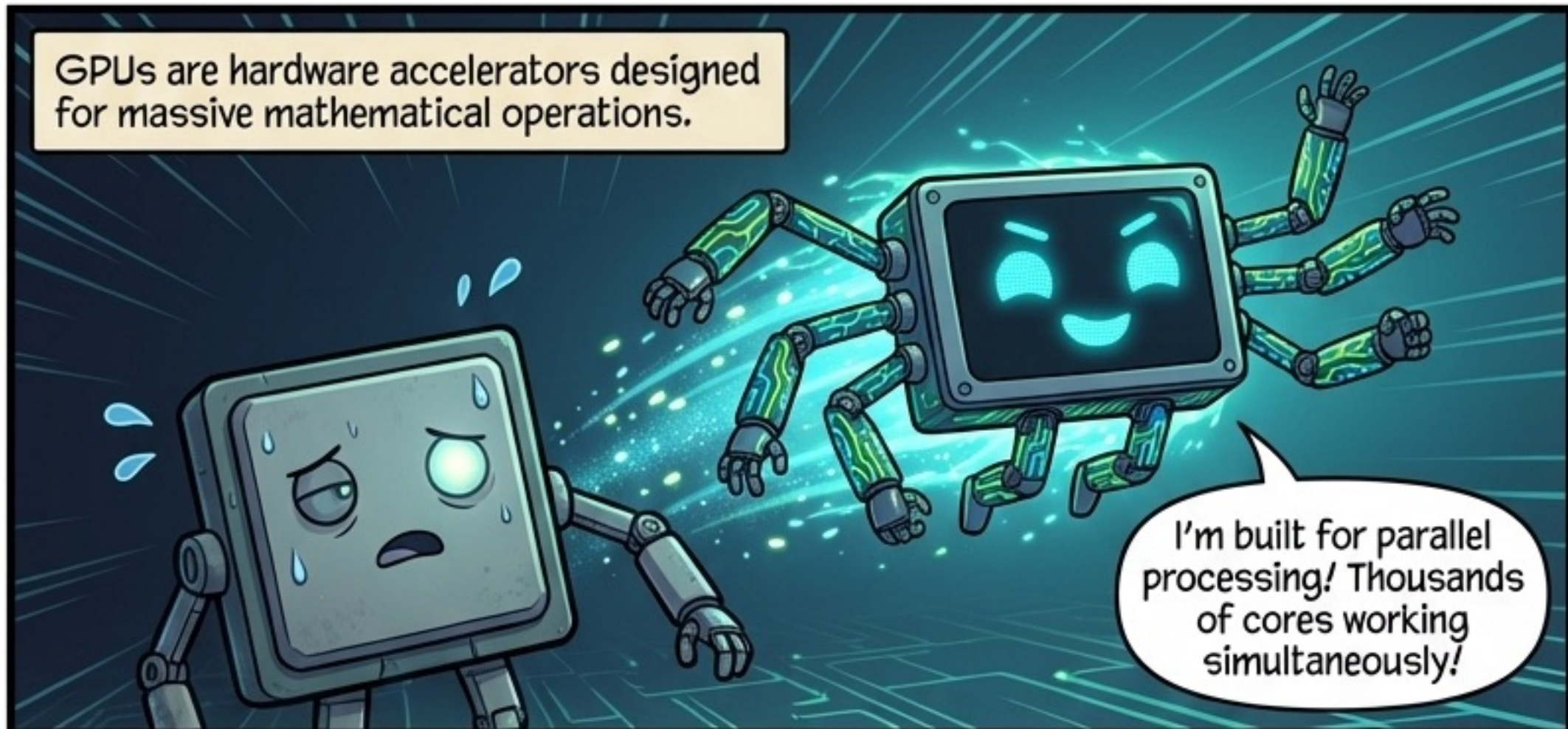


How **MUCH** compute?

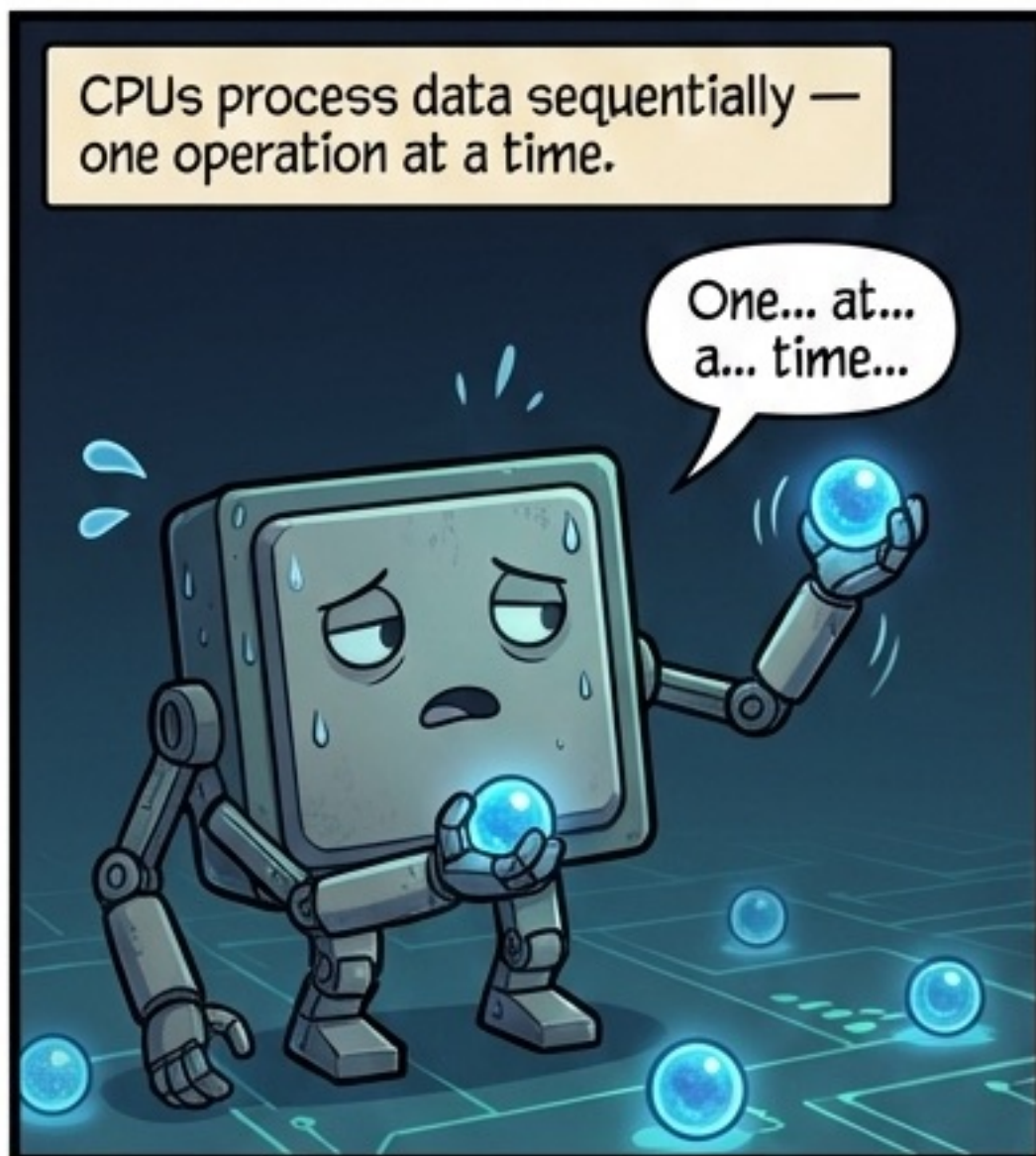




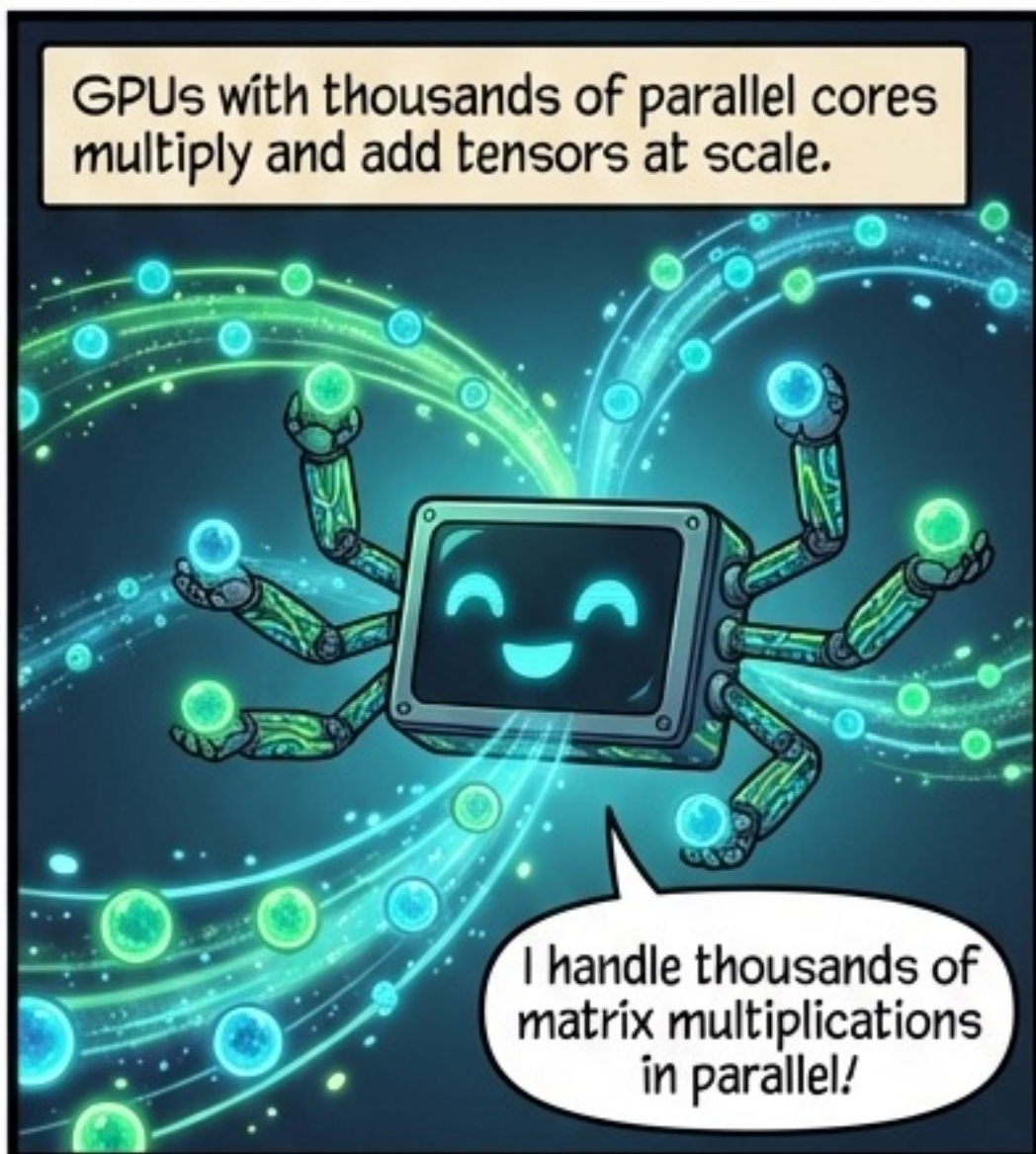
GPUs are hardware accelerators designed for massive mathematical operations.



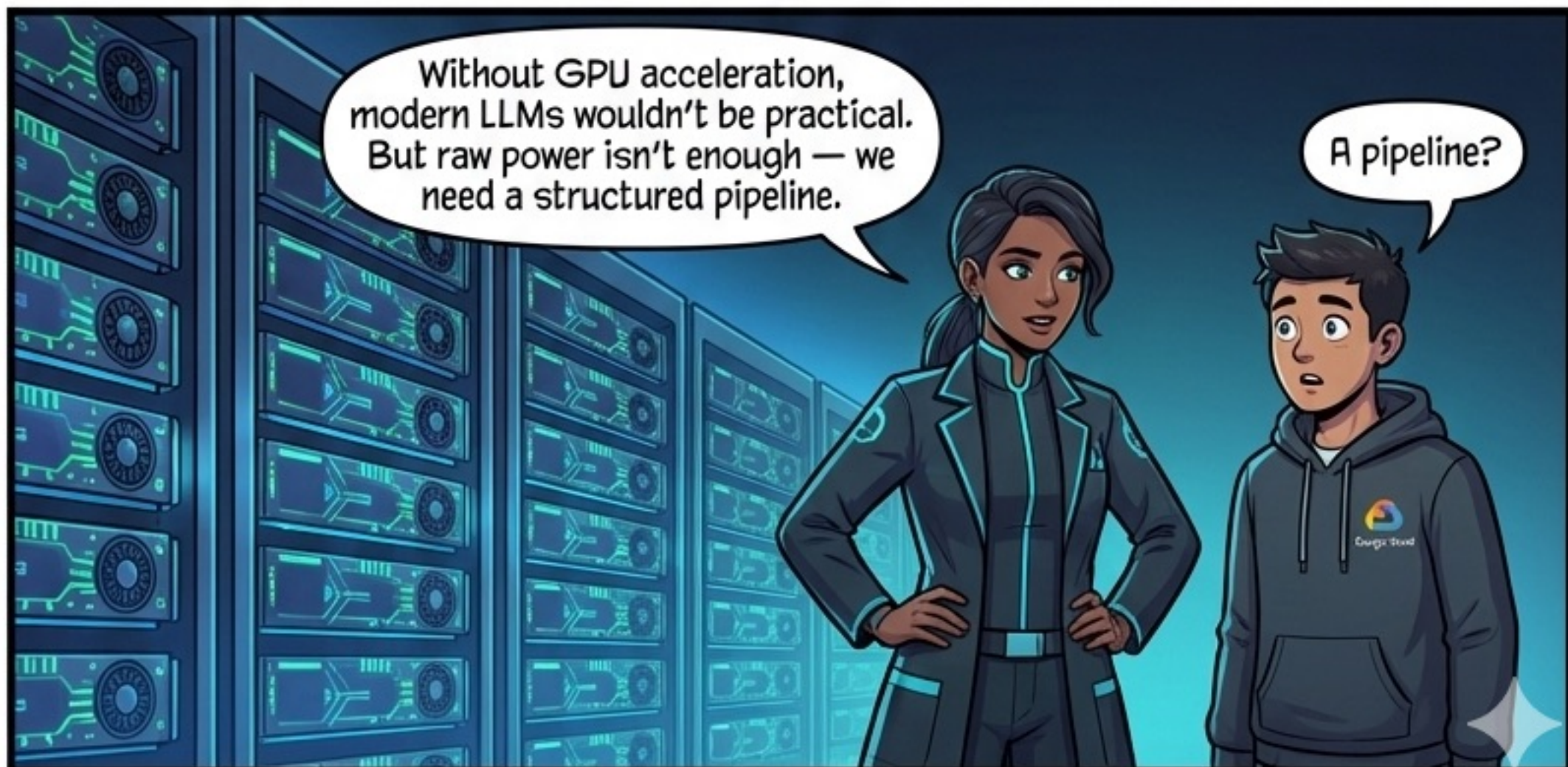
CPUs process data sequentially — one operation at a time.



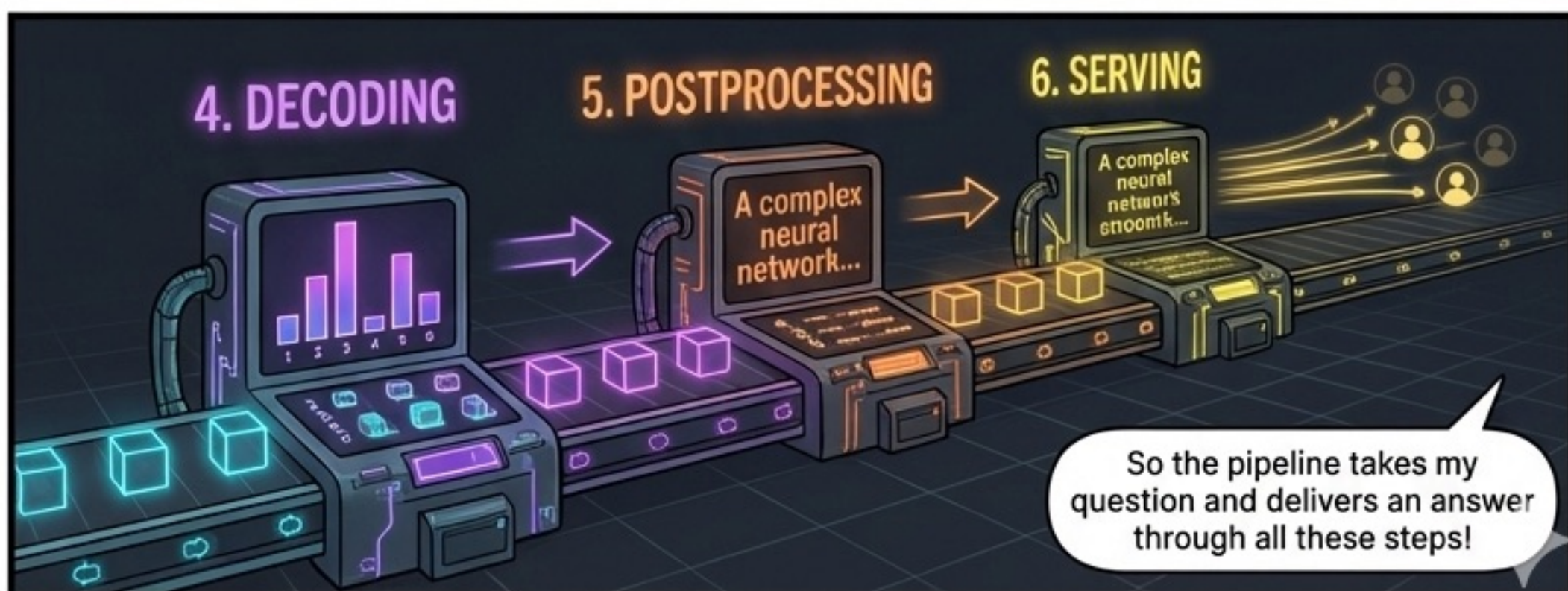
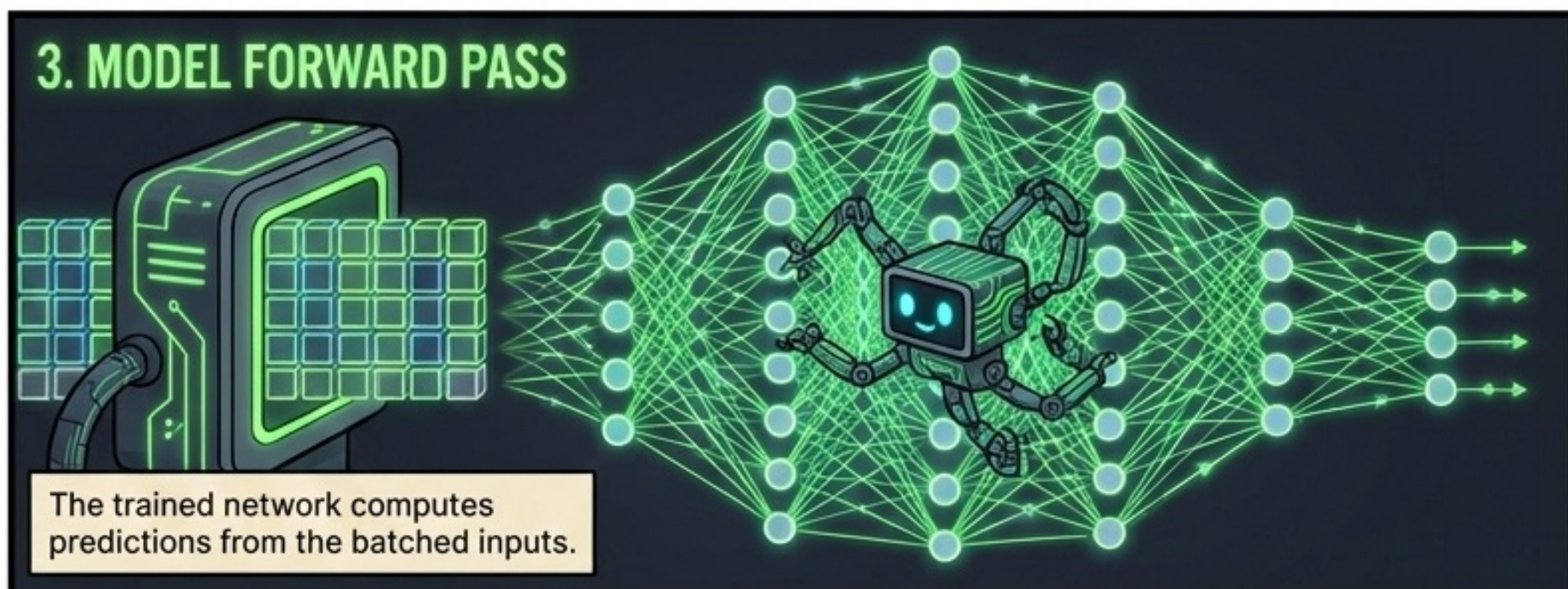
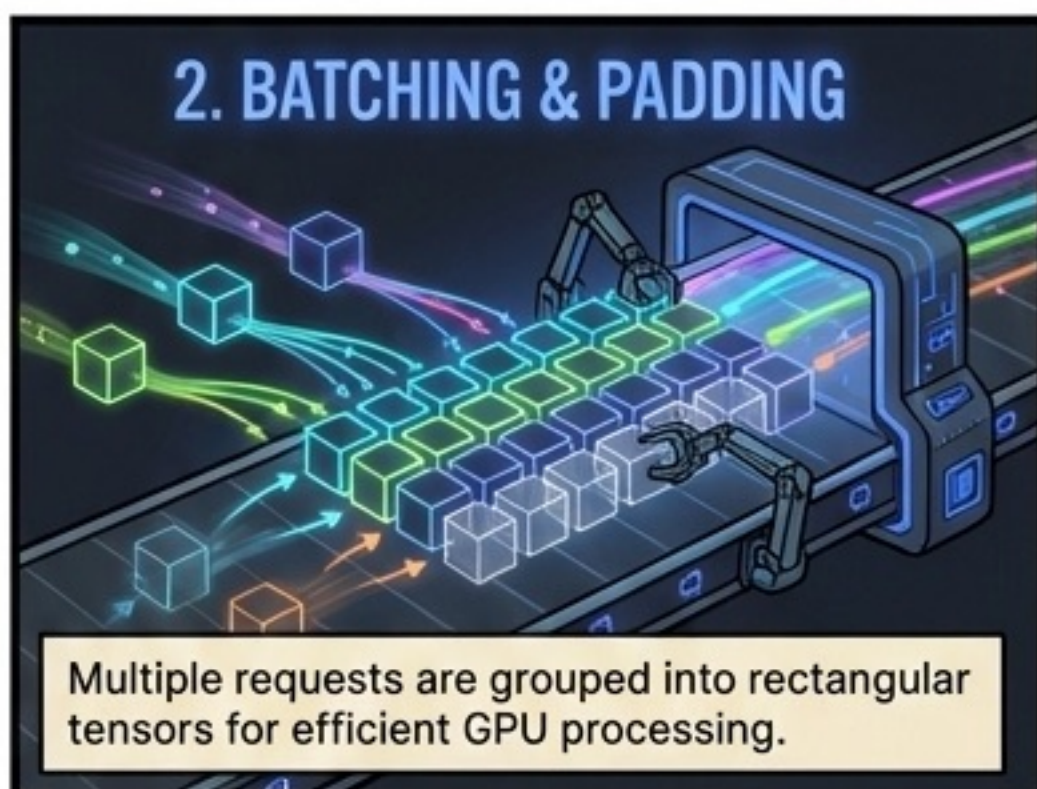
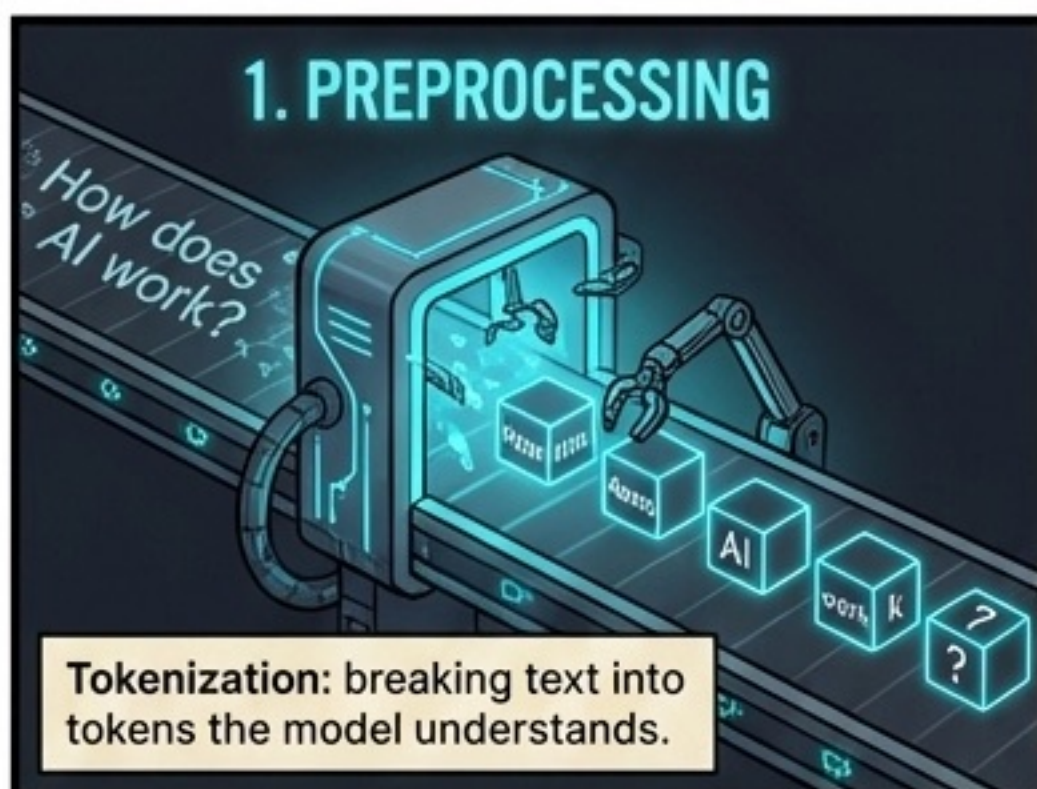
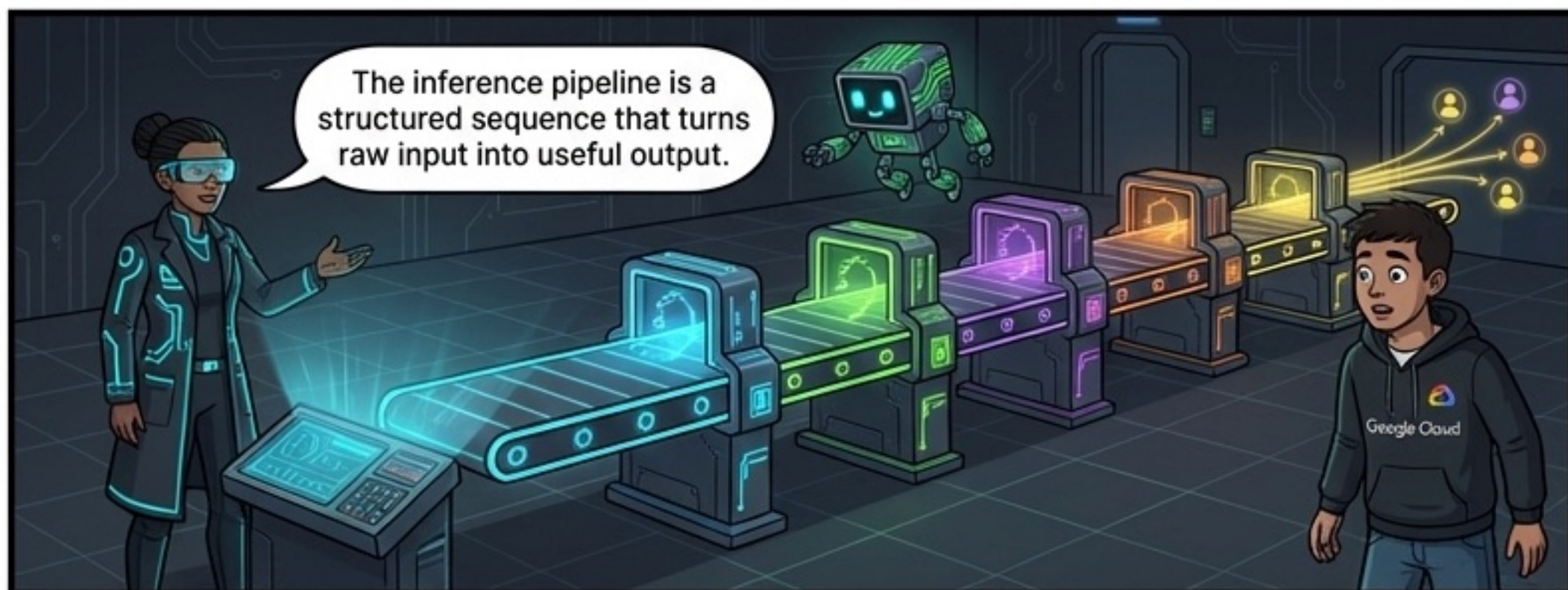
GPUs with thousands of parallel cores multiply and add tensors at scale.



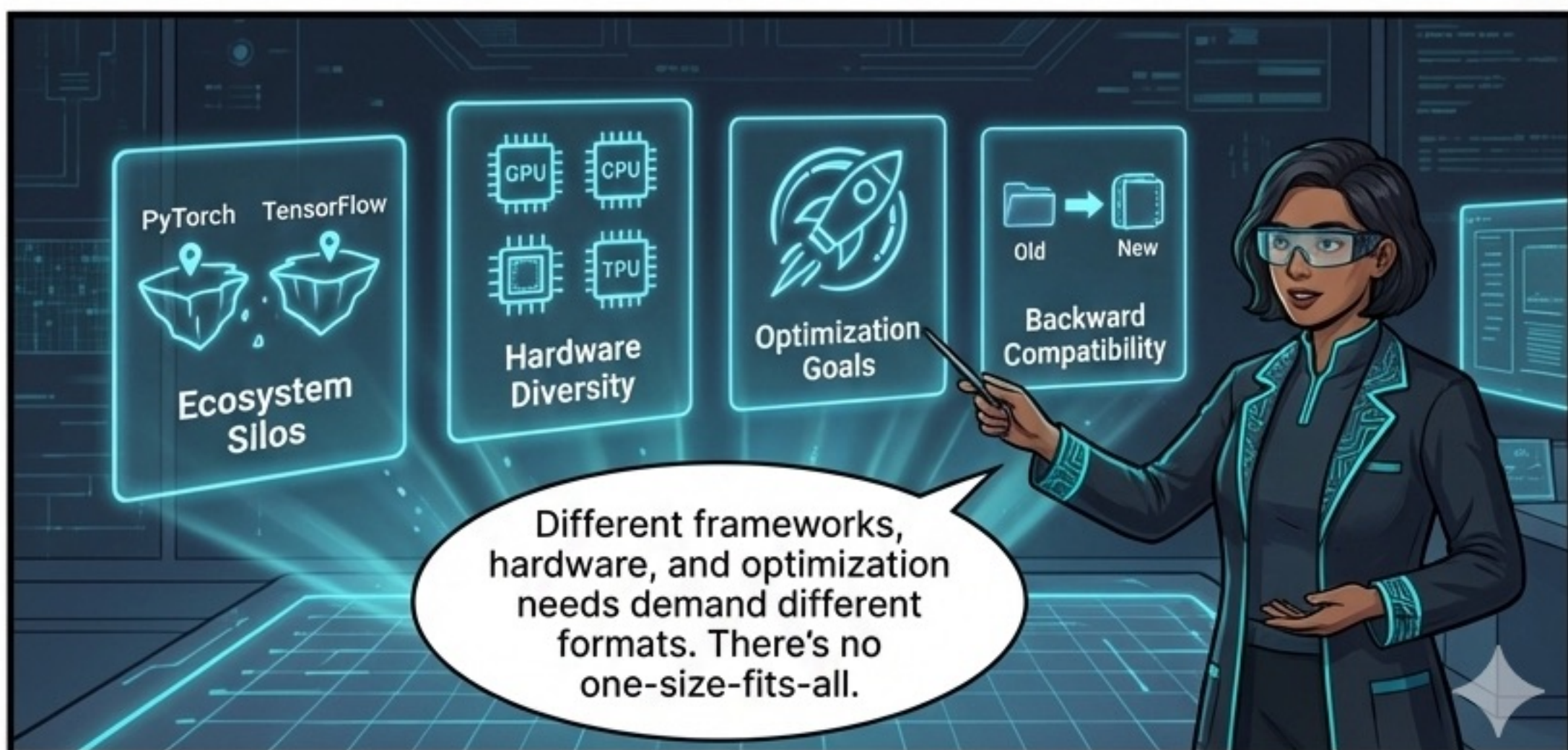
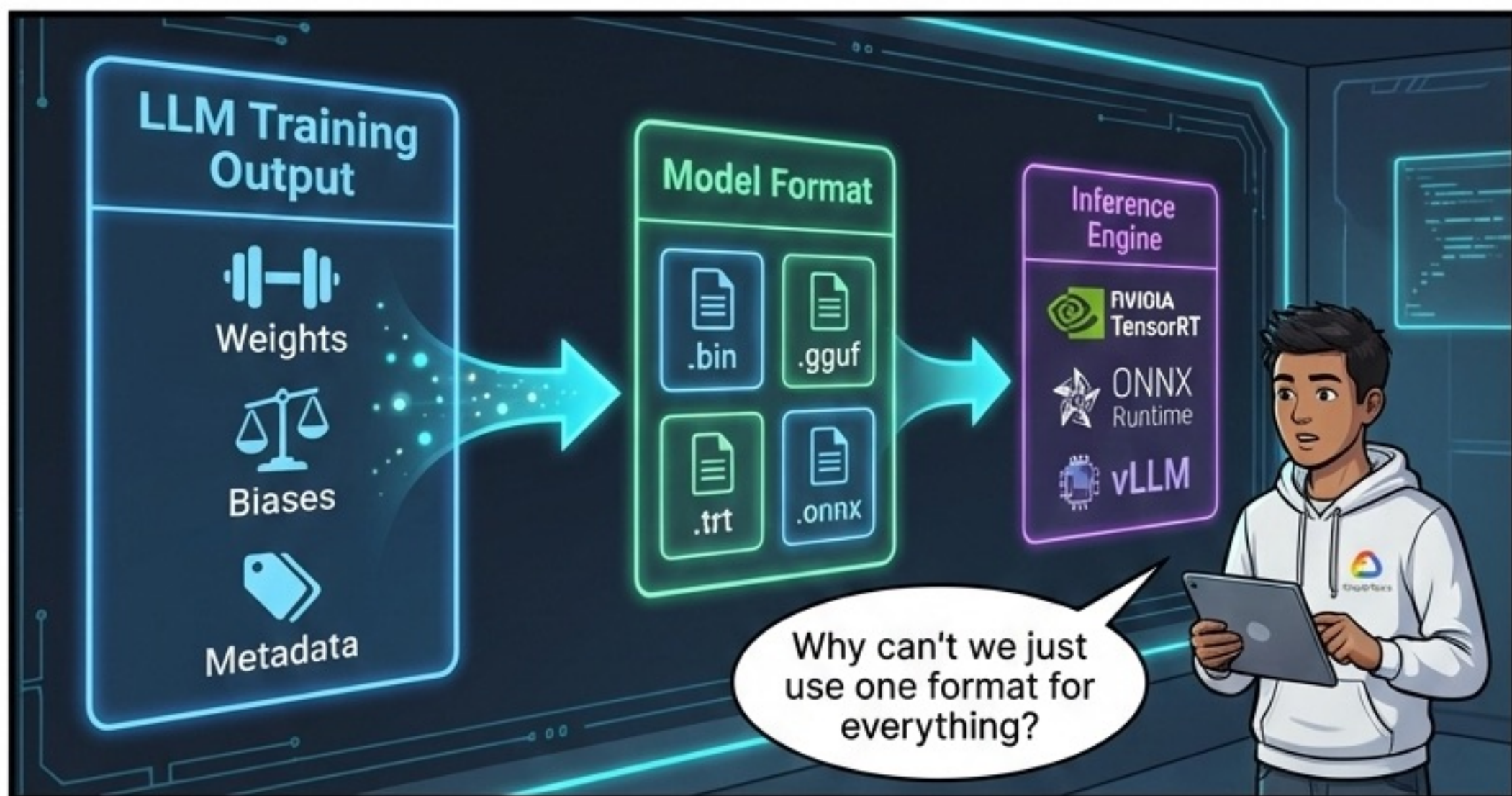
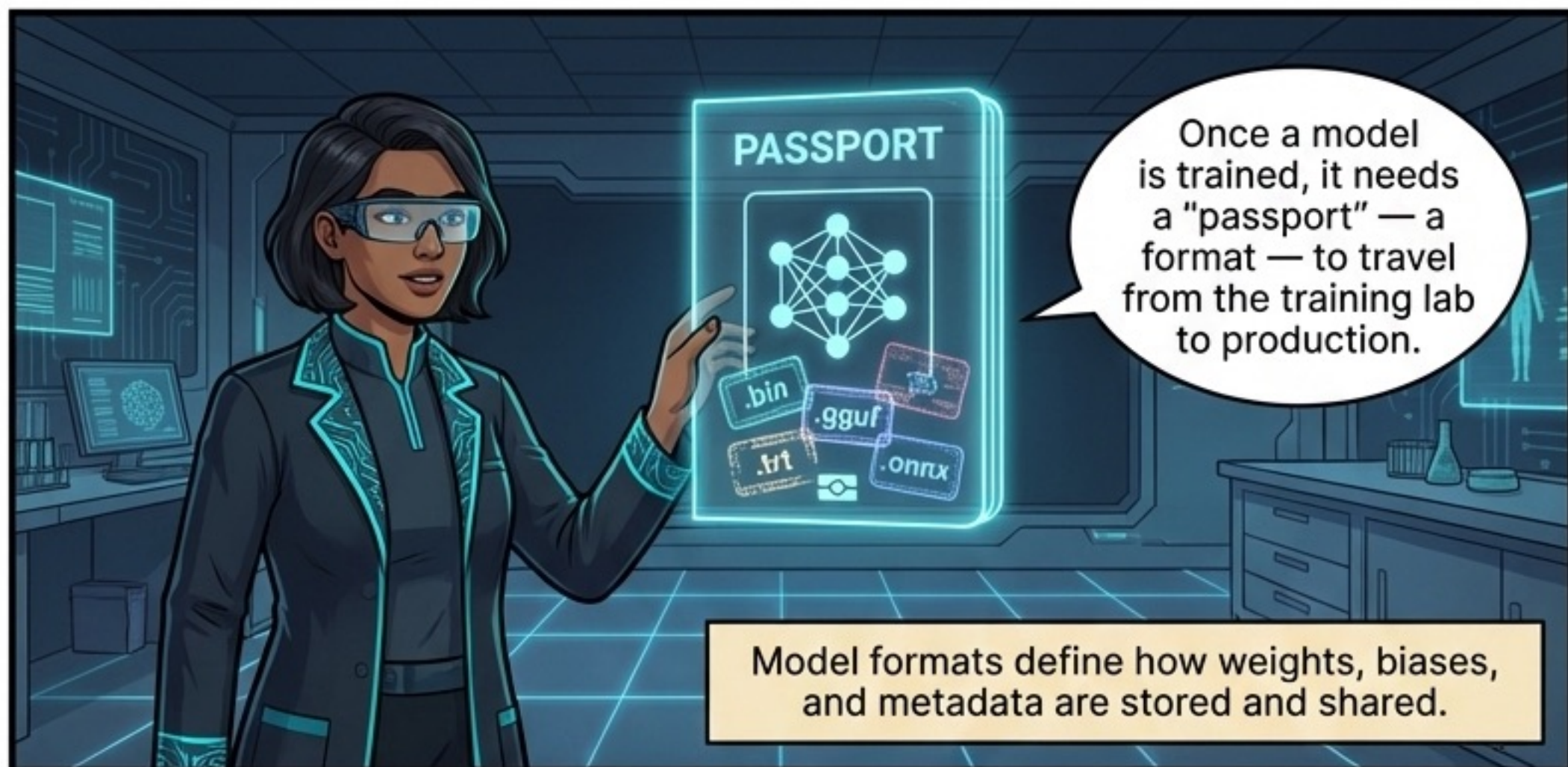
Without GPU acceleration, modern LLMs wouldn't be practical. But raw power isn't enough — we need a structured pipeline.













## SAFETENSORS



Fast, secure. The default for sharing on Hugging Face. Prevents code execution on load.

## GGUF



Compact, quantized. Powers 'LLMs on your laptop' with llama.cpp.

## TENSORRT ENGINES



Compiled for NVIDIA GPUs. Maximum performance, lowest latency.

## ONNX



The universal bridge. Best for mixed-hardware environments.

Research → Framework-native

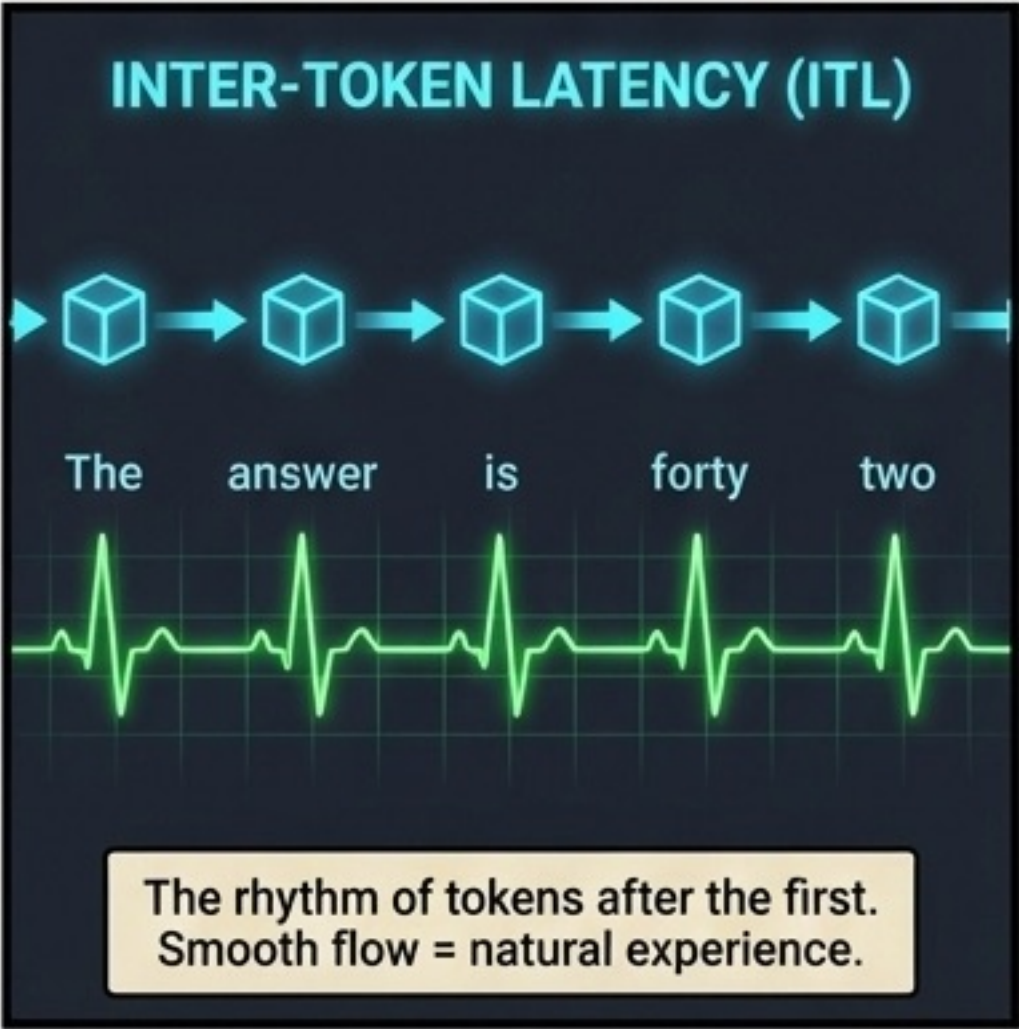
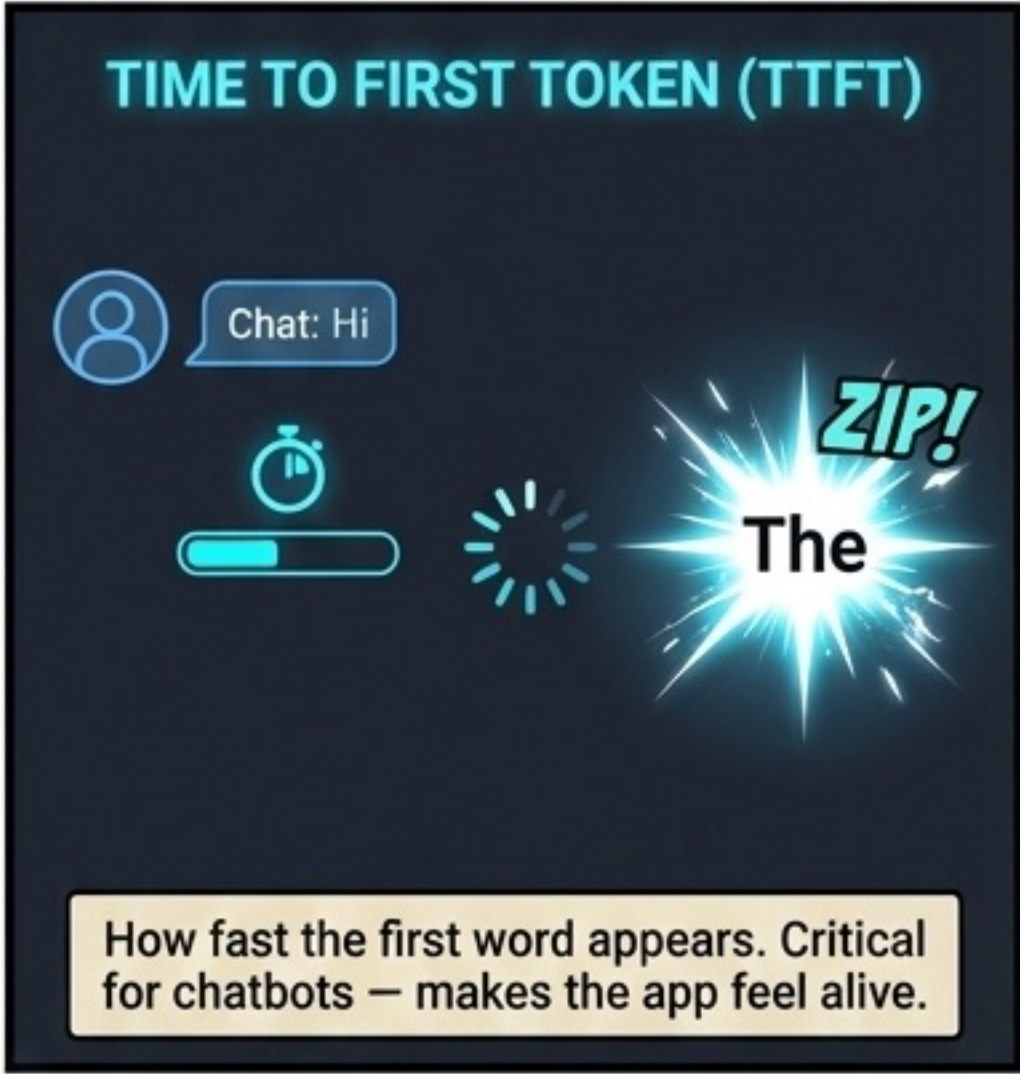
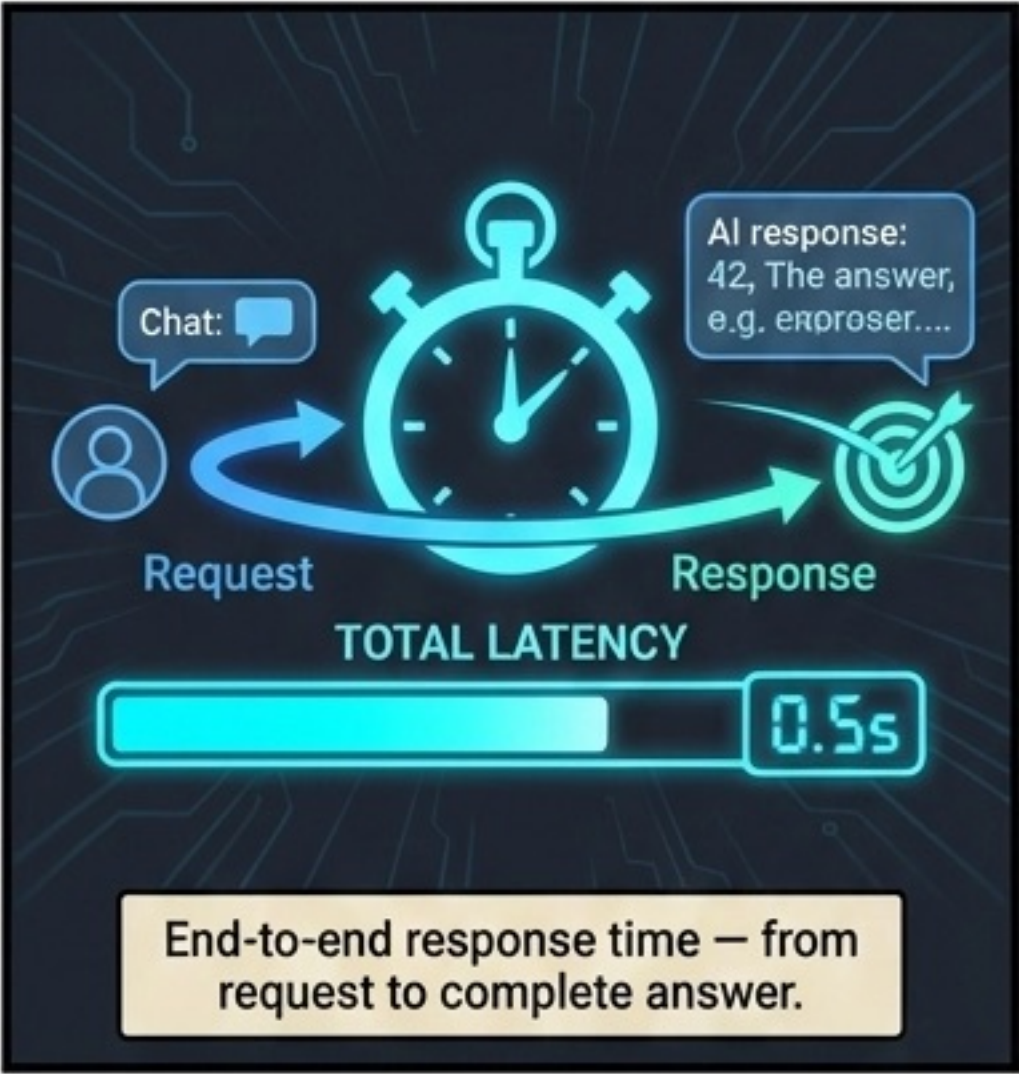
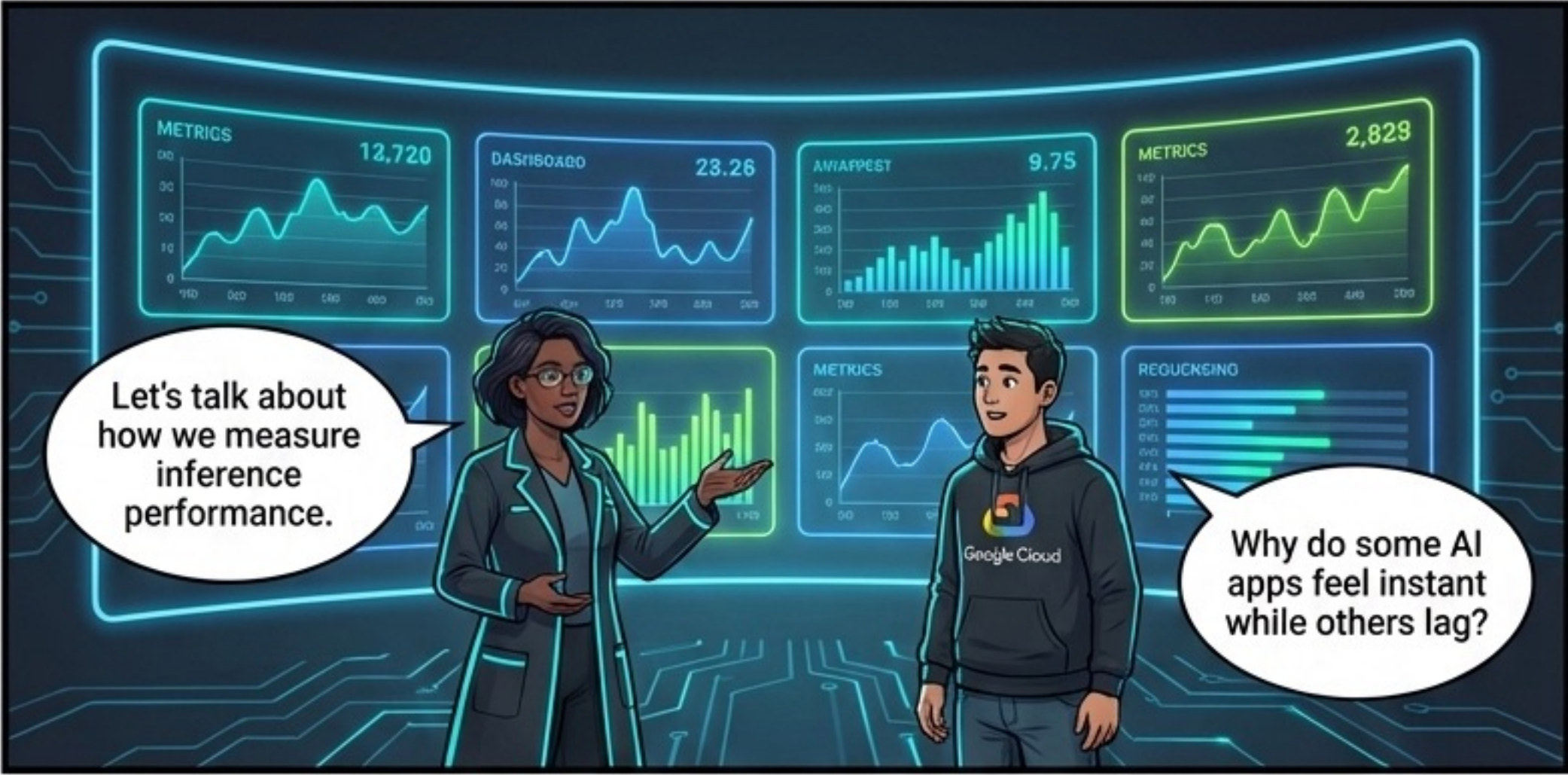
Sharing → Safetensors

Local → GGUF

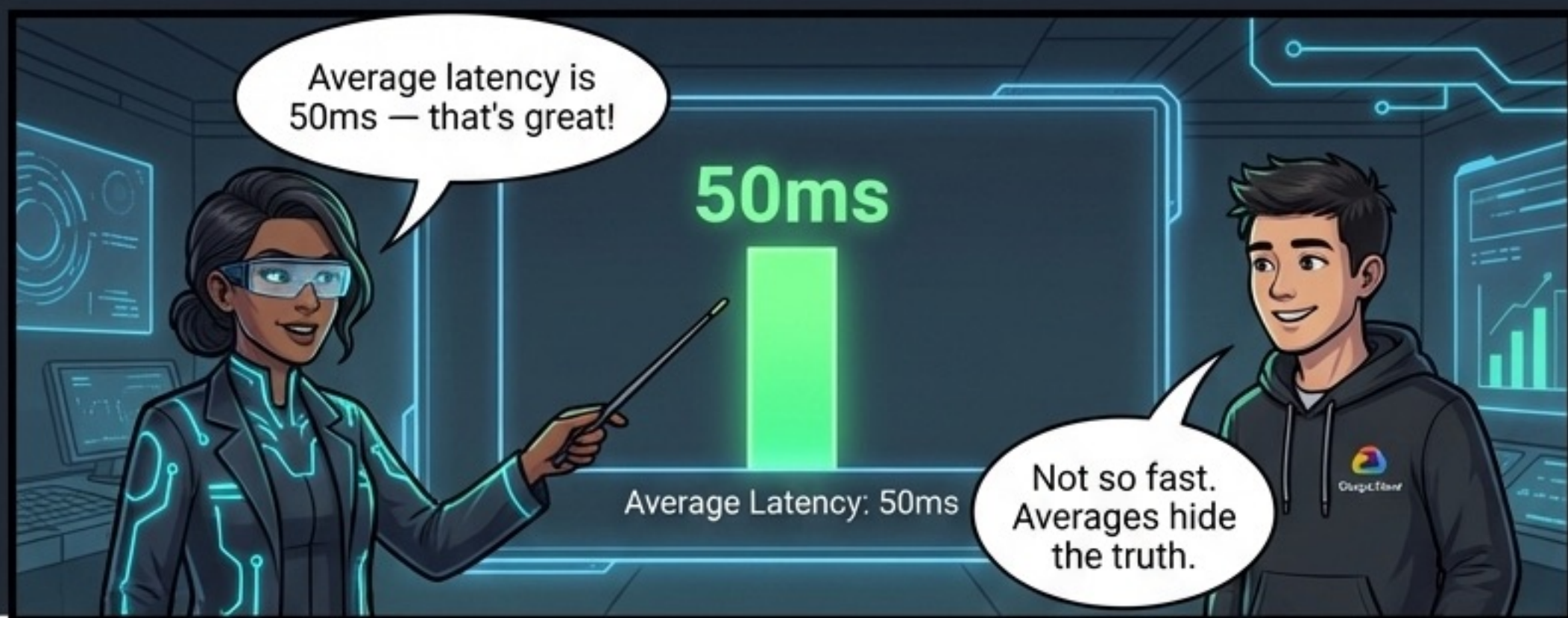
Production → TensorRT/ONNX

The right format depends on where the model is in its lifecycle!

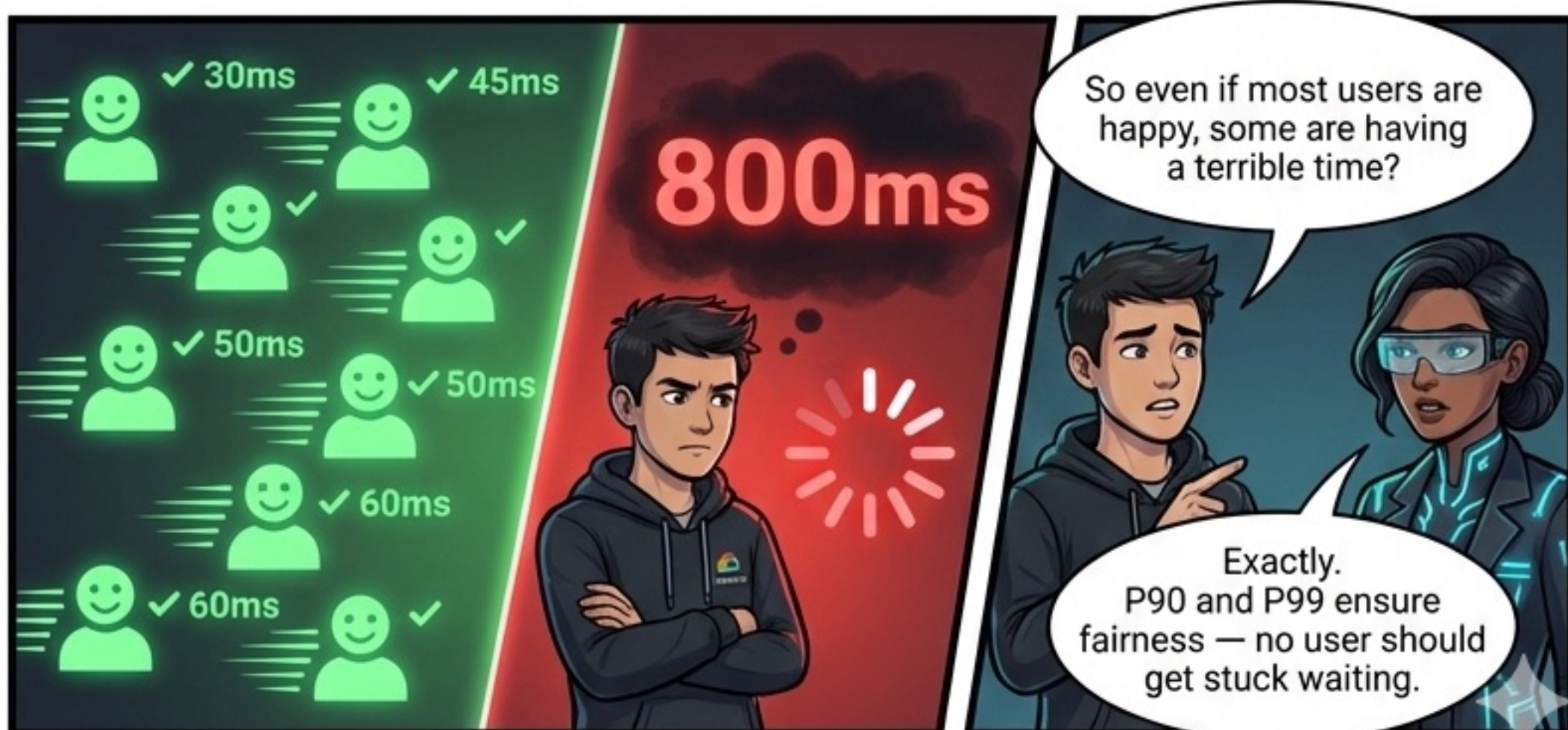
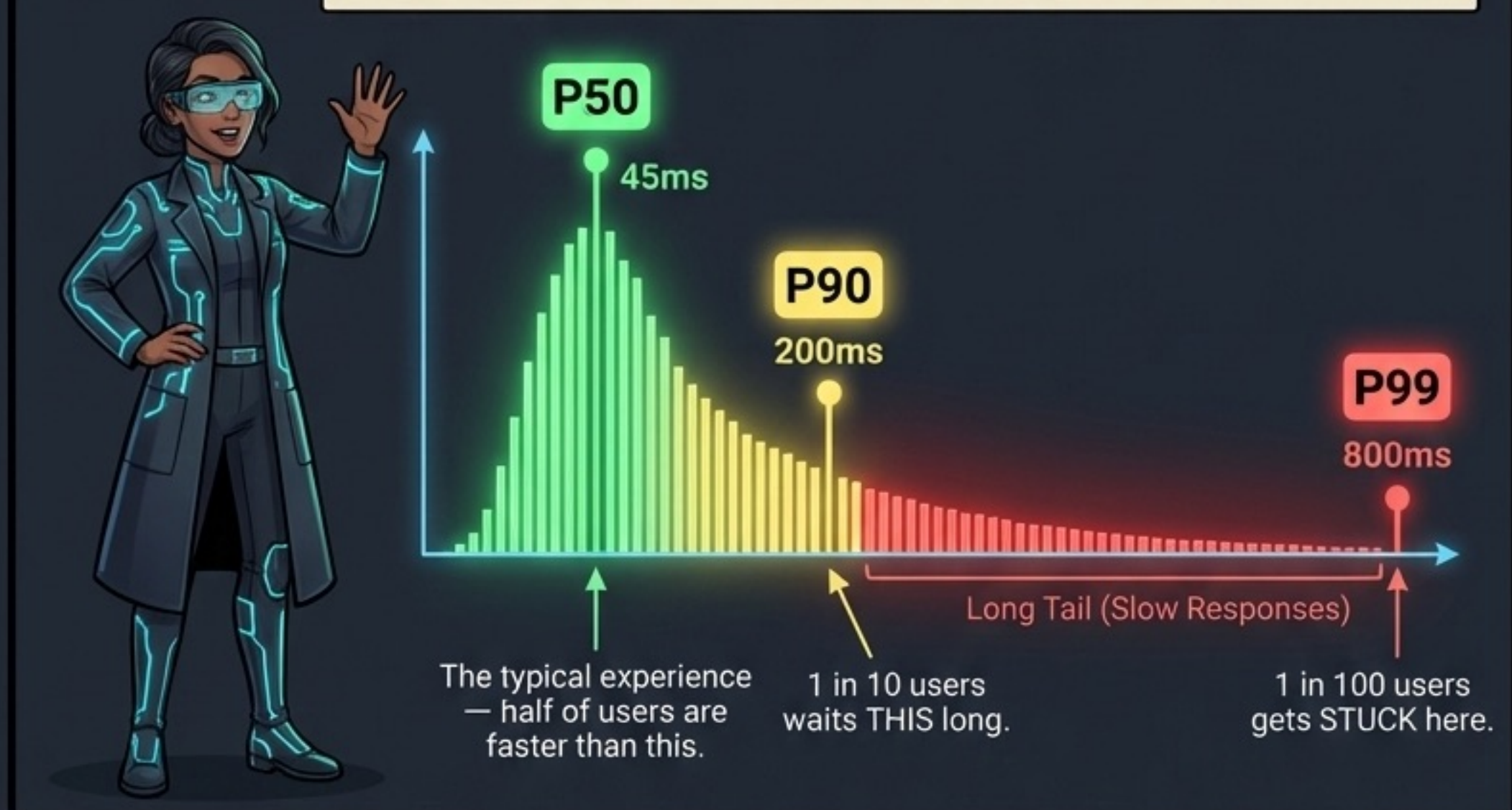






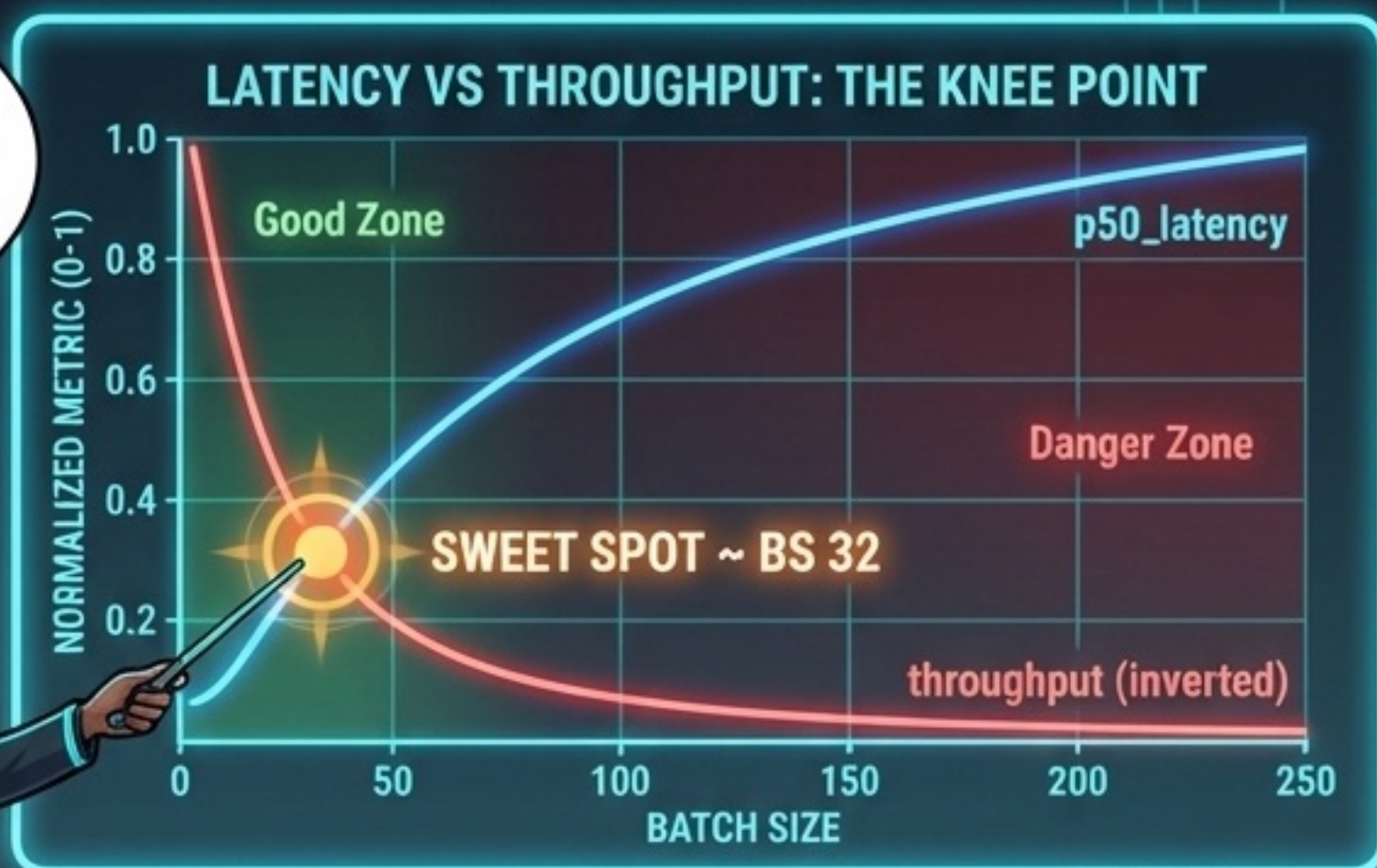


Percentiles reveal the real distribution. P50 is the median. P90 and P99 show the worst-case experiences that averages hide.





Throughput and latency pull in opposite directions. The Knee Point is where you get the best balance.



Past the knee point, throughput gains flatten while latency skyrockets. Users suffer.



**Chatbots:**  
Low TTFT:  
instant response

**Translation:**  
High Throughput:  
handle massive  
workloads

**Streaming:**  
Steady ITL:  
natural flow

**On-device AI:**  
Small Footprint:  
speed + low power

Different applications optimize for different metrics.

Percentiles reveal the real distribution. The knee point balances throughput and latency. That's how we make AI fast AND reliable.

I'm ready to build my own inference pipeline!

Let's go parallel!

Training teaches the model. Inference applies it. Now go build something amazing.