- Large Language Models (LLMs) are powerful but resource-hungry. Running them efficiently, especially on consumer GPUs, requires a game-changer: Quantization.

# What is Quantization?

- Quantization is a technique that reduces a model's memory footprint and computational requirements by representing weights and activations with lower precision (e.g., 4-bit or 8-bit instead of 16-bit or 32-bit).



LLM.int8()

8-bit Vector-wise Quantization

(1) Find vector-wise constants: $C_W$ & $C_X$

(2) Quantize
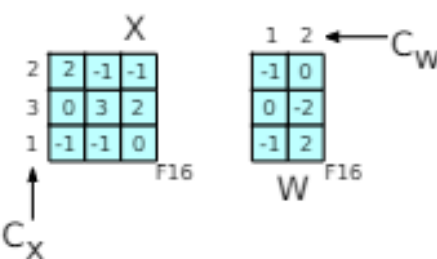$$X_{F16} * (127/C_X) = X_{I8}$$
$$W_{F16} * (127/C_W) = W_{I8}$$

(3) Int8 Matmul
$$X_{I8} \ W_{I8} = Out_{I32}$$

(4) Dequantize
$$\frac{Out_{I32} * (C_X \otimes C_W)}{127*127} = Out_{F16}$$

16-bit Decomposition

(1) Decompose outliers
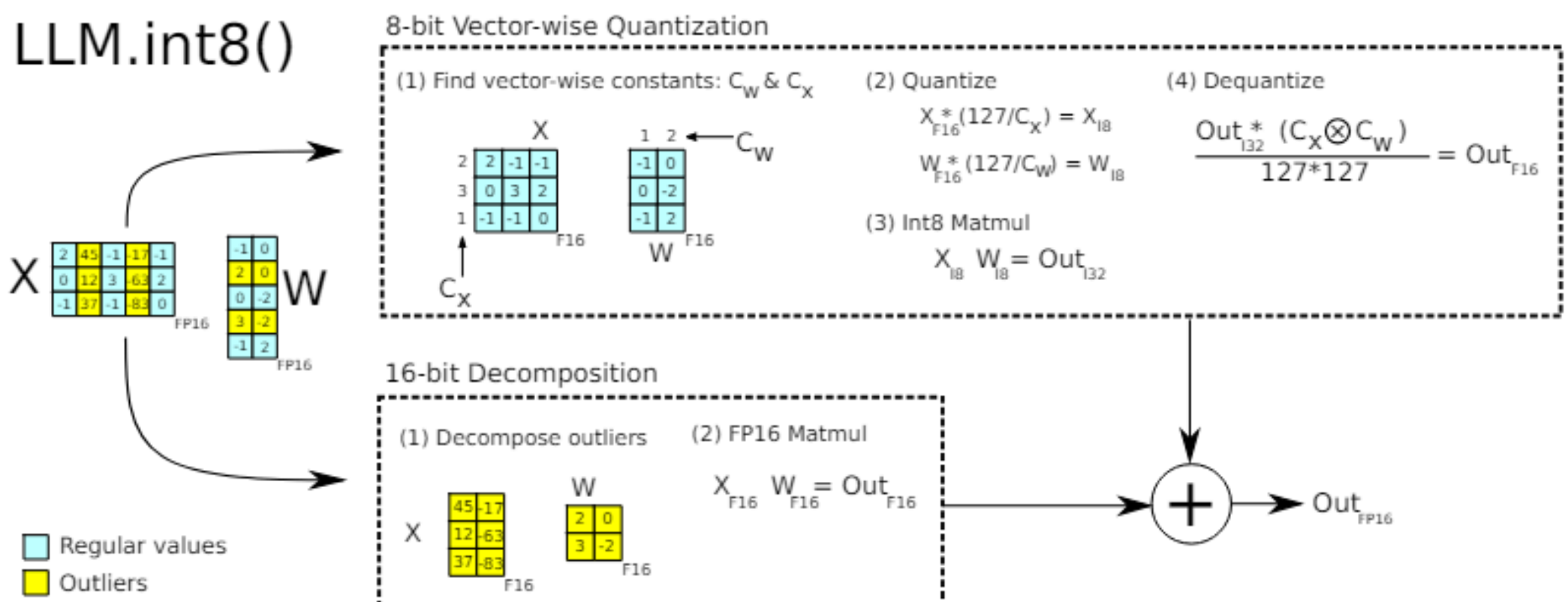
(2) FP16 Matmul
$$X_{F16} \ W_{F16} = Out_{F16}$$

Regular values
Outliers

$Out_{FP16}$

# Why Quantization?

- Reduced Memory Usage – Load larger models on smaller hardware

- Faster Inference – Less computation = quicker responses

- Lower Power Consumption – Crucial for edge AI & mobile deployment

- Cheaper Deployment – Enables cost-efficient scaling

# Types of Quantization for LLMs

## Post-Training Quantization (PTQ)

- Convert a trained FP16/FP32 model to a lower precision format
- Fast & easy, but may lead to slight accuracy loss

## Quantization-Aware Training (QAT)

- Train the model while simulating quantization effects
- Better accuracy but requires retraining

# Popular Quantization Methods

- **GPTQ** (Generalized Post-Training Quantization) – Optimized for minimal performance loss

- **AWQ** (Activation-aware Weight Quantization) – Preserves accuracy in low-bit models

- **Bitsandbytes** (8-bit / 4-bit Quantization) – Hugging Face integration for easy model loading

# How to Use Quantization in LLMs?

- Run models with 4-bit quantization using bitsandbytes

```python
from transformers import AutoModelForCausalLM, AutoTokenizer
import torch
import bitsandbytes as bnb

model_name = "mistralai/Mistral-7B-v0.1"

model = AutoModelForCausalLM.from_pretrained(
    model_name,
    device_map="auto",
    load_in_4bit=True,   # Enables 4-bit quantization
    bnb_4bit_compute_dtype=torch.float16
)

tokenizer = AutoTokenizer.from_pretrained(model_name)
```

# Speed up inference with GPTQ quantized models

```python
from auto_gptq import AutoGPTQForCausalLM

quantized_model =
AutoGPTQForCausalLM.from_quantized("TheBloke/Llama-2-
13B-GPTQ")
```

## Does Quantization Affect Accuracy?

- Yes, but modern quantization methods (GPTQ, AWQ) minimize loss while offering huge efficiency gains.

## Who Benefits from Quantization?

✅ Researchers needing faster experiments
✅ Startups with limited cloud budgets
✅ Developers running LLMs locally on consumer GPUs
✅ AI on mobile & edge devices

Stay Tuned for **Day 32** of

**Mastering LLMs**