# BERT-Base: Parameter and Memory Breakdown
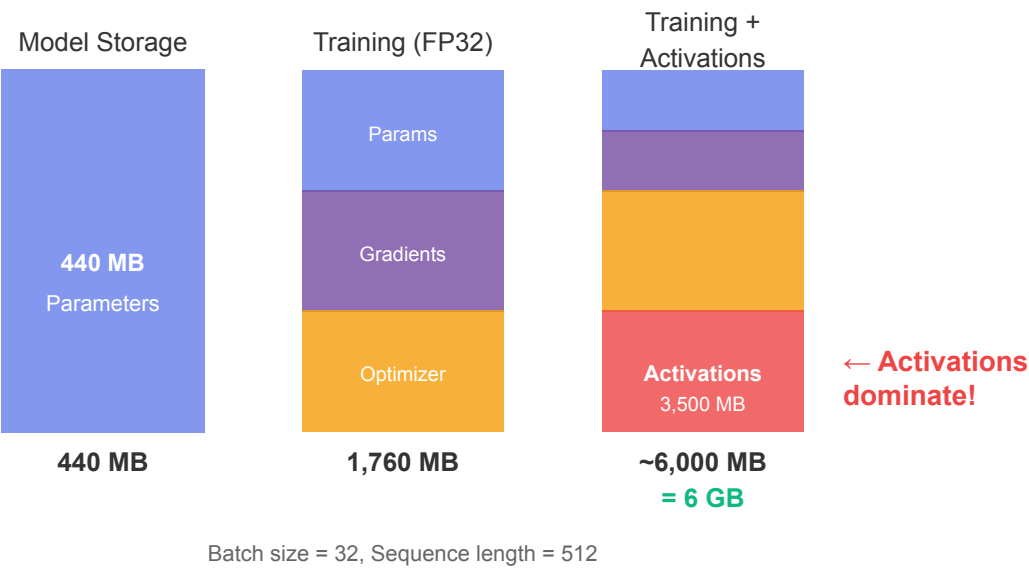
## Parameter Distribution (110M Total = 440 MB)

Feed-Forward Networks: 57M parameters (52%)

Embeddings: 23M (21%)

Attention: 21M (19%)

Other: 9M

| 0M | 55M | 110M |

## Per-Layer Breakdown (12 Layers)

Layer 12
Layer 11
Layer 10
Layer 9
Layer 8
Layer 7
Layer 6
Layer 5
Layer 4
Layer 3
Layer 2
Layer 1

Parameters (millions)

■ Attention: 1.77M
■ Feed-forward: 4.72M

**Each layer: ~6.5M parameters**

## Training Memory Breakdown

Model Storage

**440 MB**
Parameters

**440 MB**

Training (FP32)

Params

Gradients

Optimizer

**1,760 MB**

Training + Activations

Activations
3,500 MB

**~6,000 MB**
**= 6 GB**

← **Activations dominate!**

Batch size = 32, Sequence length = 512

## Key Insights for Leaders

• Feed-forward layers: 52% of parameters

• Activations: 60% of training memory

• Doubling batch size ≈ doubles memory

• Doubling sequence length ≈ doubles memory

• Optimizer states: 2× parameter memory

• Total training memory >> model size

• Memory is often the limiting factor

• Gradient checkpointing trades speed for memory