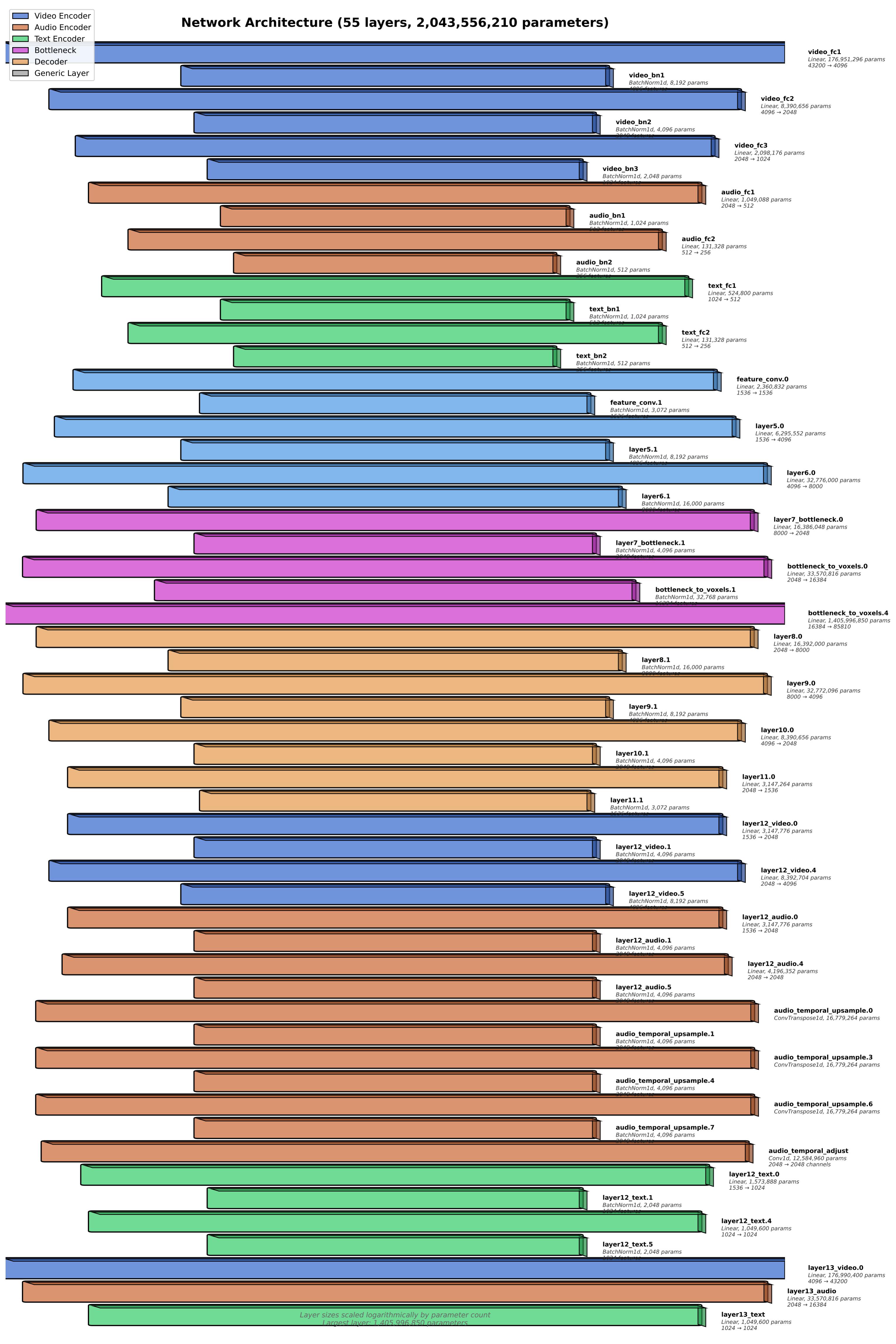


# Network Architecture (55 layers, 2,043,556,210 parameters)



Layer sizes scaled logarithmically by parameter count  
Largest layer: 1,405,996,850 parameters

**video.fc1**  
Linear, 176,951,296 params  
43200 → 4096

**video.fc2**  
Linear, 8,390,656 params  
4096 → 2048

**video.fc3**  
Linear, 2,098,176 params  
2048 → 1024

**audio.fc1**  
Linear, 1,049,088 params  
2048 → 512

**audio.fc2**  
Linear, 131,328 params  
512 → 256

**text.fc1**  
Linear, 524,800 params  
1024 → 512

**text.fc2**  
Linear, 131,328 params  
512 → 256

**feature.conv.0**  
Linear, 2,360,832 params  
1536 → 1536

**layer5.0**  
Linear, 6,295,552 params  
1536 → 4096

**layer6.0**  
Linear, 32,776,000 params  
4096 → 8000

**layer7\_bottleneck.0**  
Linear, 16,386,048 params  
8000 → 2048

**bottleneck\_to\_voxels.0**  
Linear, 33,570,816 params  
2048 → 16384

**bottleneck\_to\_voxels.1**  
Linear, 1,405,996,850 params  
16384 → 85810

**layer8.0**  
Linear, 16,392,000 params  
2048 → 8000

**layer9.0**  
Linear, 32,772,096 params  
8000 → 4096

**layer10.0**  
Linear, 8,390,656 params  
4096 → 2048

**layer11.0**  
Linear, 3,147,264 params  
2048 → 1536

**layer12\_video.0**  
Linear, 3,147,776 params  
1536 → 2048

**layer12\_video.4**  
Linear, 8,392,704 params  
2048 → 4096

**layer12\_audio.0**  
Linear, 3,147,776 params  
1536 → 2048

**layer12\_audio.4**  
Linear, 4,196,352 params  
2048 → 2048

**audio\_temporal\_upsample.0**  
ConvTranspose1d, 16,779,264 params, 2048 → 4096 channels

**audio\_temporal\_upsample.3**  
ConvTranspose1d, 16,779,264 params, 4096 → 2048 channels

**audio\_temporal\_upsample.6**  
ConvTranspose1d, 16,779,264 params, 2048 → 2048 channels

**audio\_temporal\_adjust**  
Conv1d, 12,584,960 params, 2048 → 2048 channels

**layer12\_text.0**  
Linear, 1,573,888 params, 1536 → 1024

**layer12\_text.1**

BatchNorm1d, 2,048 params, 1024 features → 1024 channels

**layer12\_text.4**  
Linear, 1,049,600 params, 1024 → 1024

**layer12\_text.5**

BatchNorm1d, 2,048 params, 1024 features → 1024 channels

**layer13\_video.0**  
Linear, 176,990,400 params, 4096 → 43200

**layer13\_audio**  
Linear, 33,570,816 params, 2048 → 16384

**layer13\_text**  
Linear, 1,049,600 params, 1024 → 1024