

Quantization: less bits per weight

Pruning: less number of weights

Huffman Encoding

original network
→
original size

Train Connectivity

Prune Connections

Train Weights

same accuracy
→
9x-13x reduction

Cluster the Weights

Generate Code Book

Quantize the Weights
with Code Book

Retrain Code Book

same accuracy
→
27x-31x reduction

Encode Weights

Encode Index

same accuracy
→
35x-49x reduction

