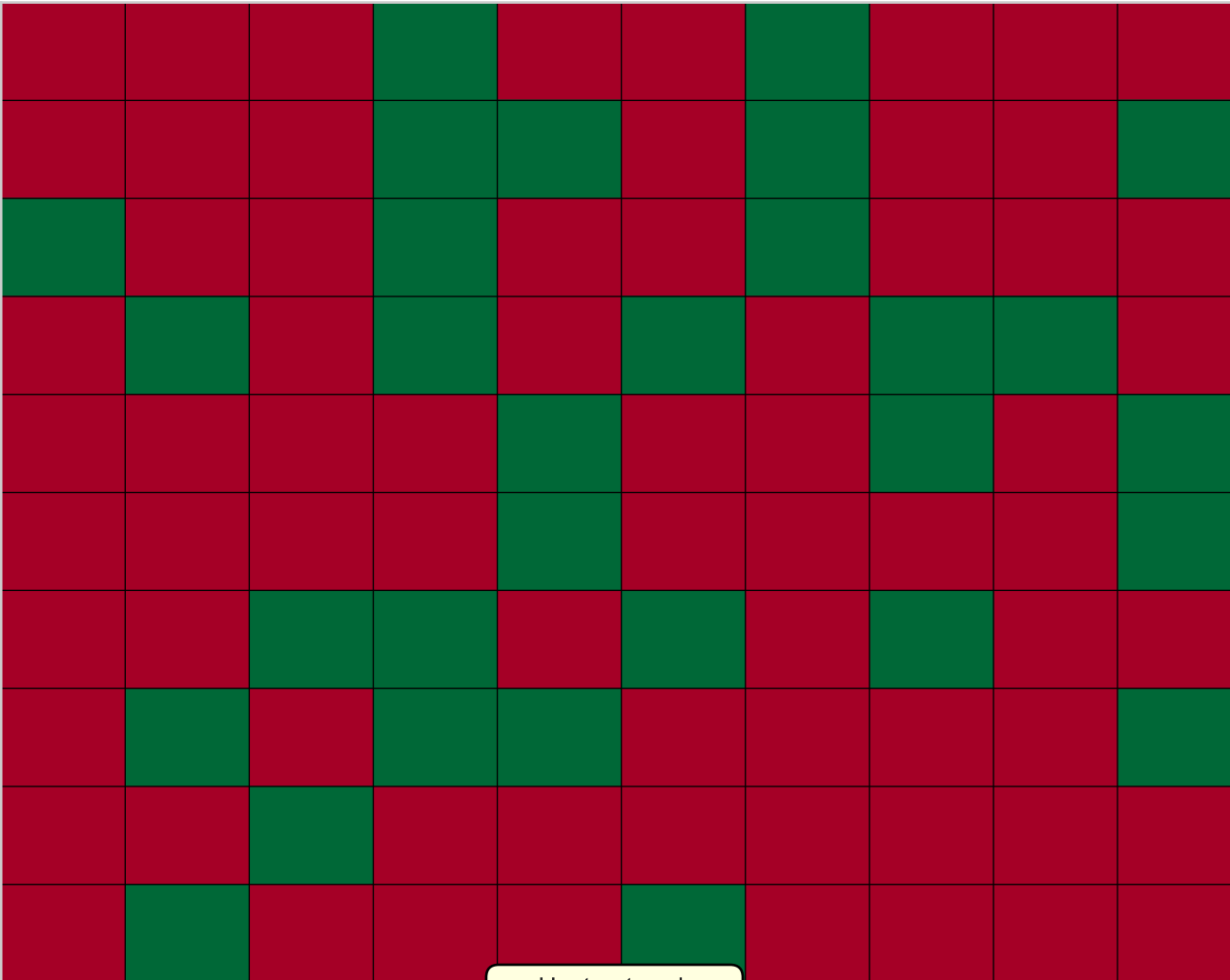# Pruning Strategies: Random vs Structured
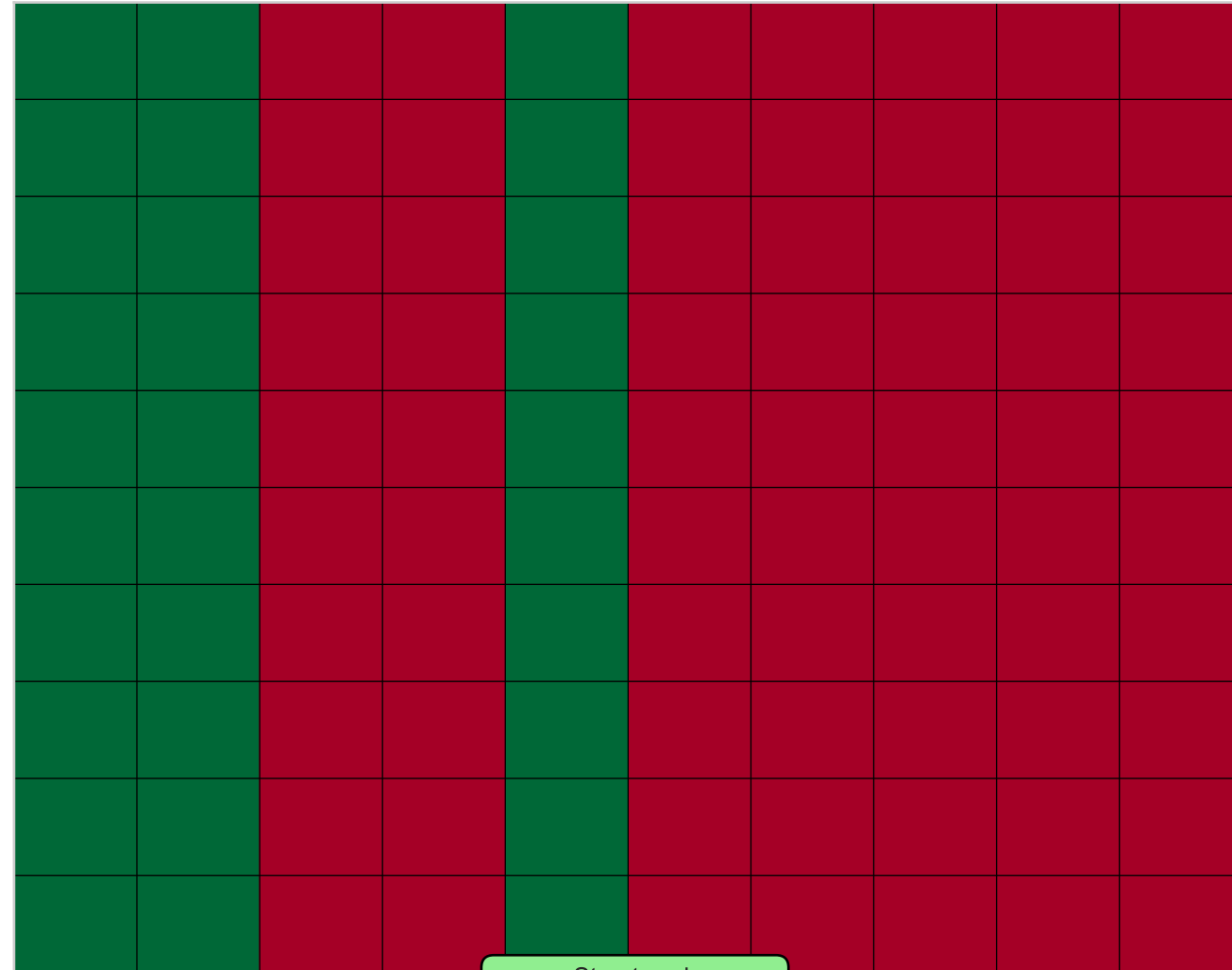## Both reduce parameters, but structured is hardware-friendly

**Unstructured Pruning**
**(Remove Individual Weights)**

**Structured Pruning**
**(Remove Entire Neurons)**



Unstructured:
☐ Better accuracy
☐ Needs sparse ops

Structured:
☐ Real speedup
☐ Slightly lower accuracy