

## Sparse Activation: Computational Efficiency through Selective Expert Activation

### Dense Model (All Active)



**All 8 Experts Active (100%)**

High computation & memory usage

### Sparse MoE (Top-2 Active)



**Only 2 Experts Active (25%)**

75% computation reduction!

Active Experts	<b>8 / 8 (100%)</b>
Computation	<b>Full (100%)</b>
Memory Usage	<b>High</b>
Latency	<b>High</b>

Active Experts	<b>2 / 8 (25%)</b>
Computation	<b>75% Reduction</b>
Memory Usage	<b>Low</b>
Latency	<b>Low</b>

### ⚡ Efficiency Gains

Sparse activation enables models with billions of parameters to run with the computational cost of much smaller models. A 64-expert MoE with Top-2 routing achieves **32x parameter scaling with only 2x computation increase!**