# MoEShard



$W_i$ $W_o$ Partial outputs
Replicate
Input tokens
GPU 0
Output tokens
GPU 1
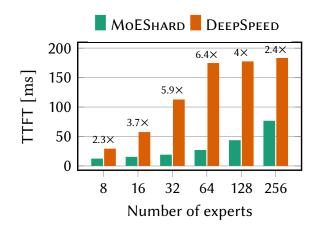column-wise split   row-wise split

- Tensor Sharding
- Block Sparse Matrix Multiplication[1]

- Perfect load balancing
- No token dropping

- Up to 6.4× faster than DeepSpeed



[1] Gale, Trevor, et al., MLSys 2023