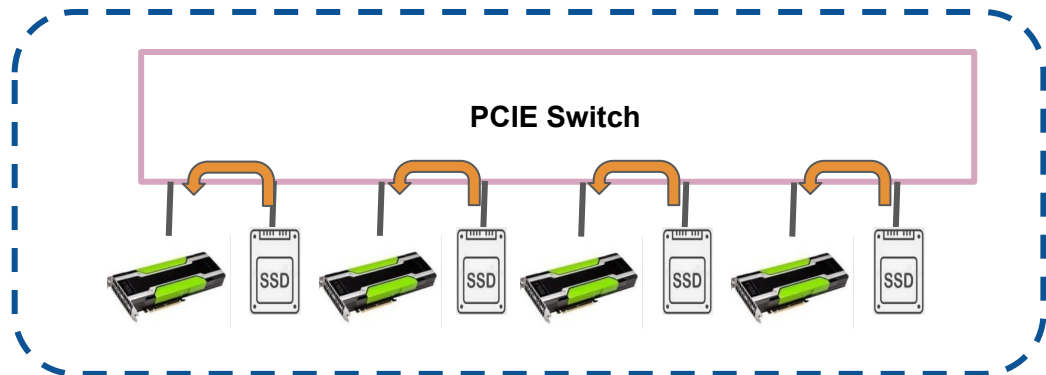
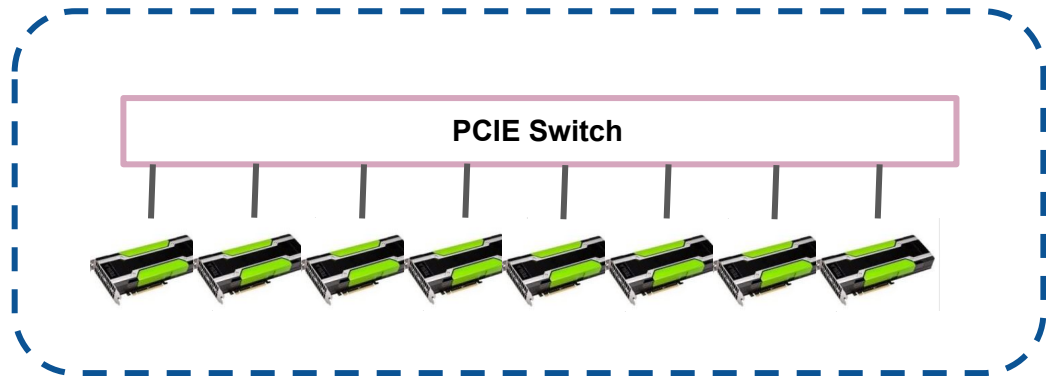


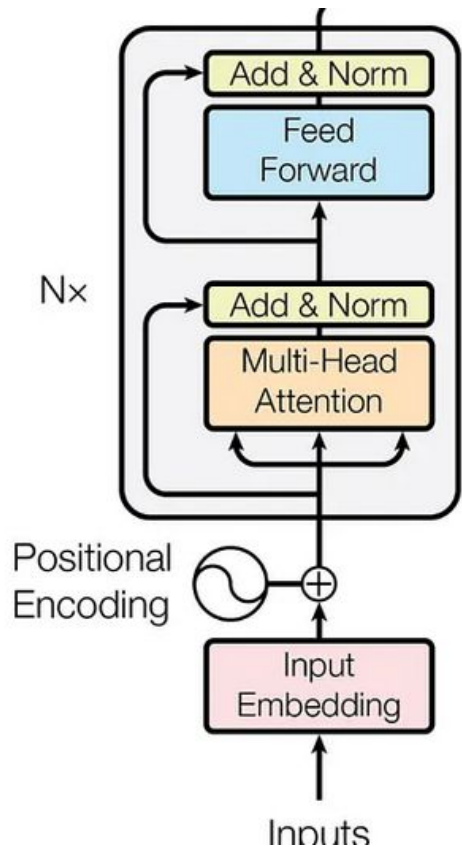
Beyond KV-Cache

GPU Server Restructured



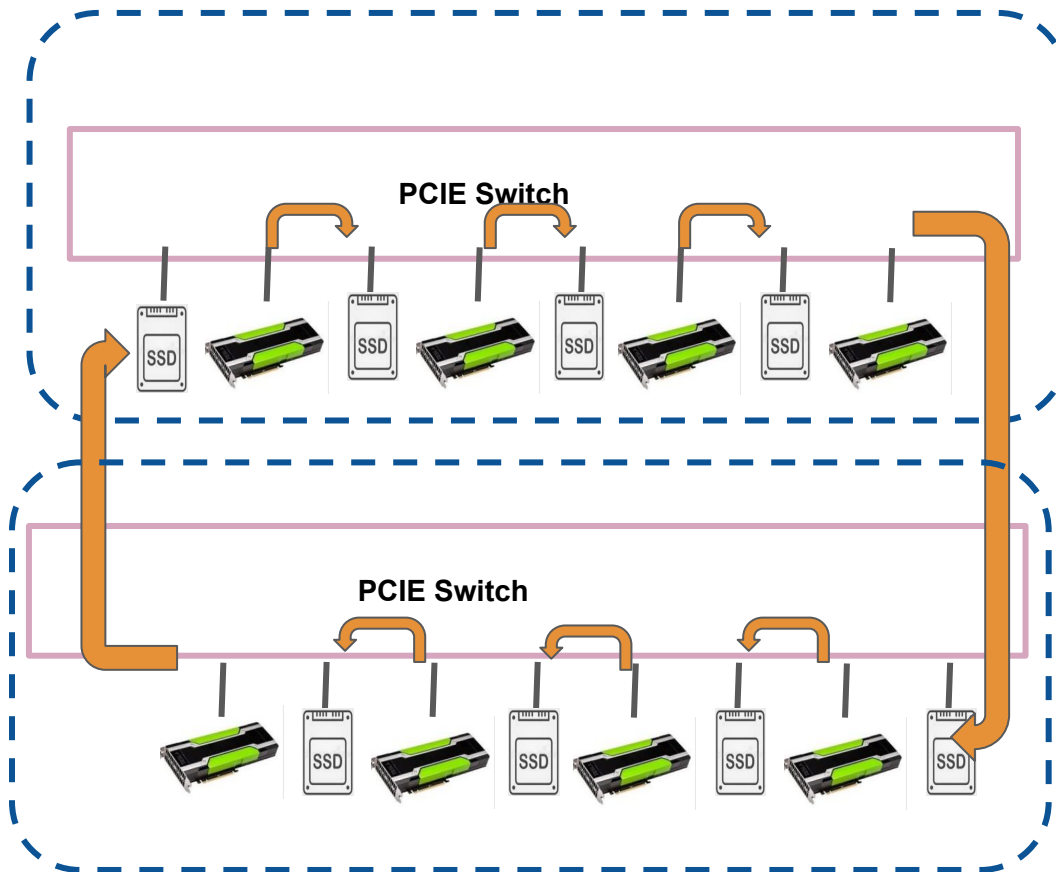
- GPU paired with SSD
- GPU now have dedicated computing scratch pad
- SSD will hold all model parameters
- GPU exchange data through neighbour SSD

SSD LLM inference Streaming Computing



- GPT-3 175B, 96 layer
- Each layer $\sim 1.8\text{B}$, 3.6G Bytes
- Gen5x16 SSD bandwidth $\sim 50\text{GB/s}$
- Model loading time for one layer: $3.6/50 \sim 70\text{ms}$
- Computing: $100\text{K} \times 1.8\text{B} / 100\text{T} \sim 1.8\text{s}$

GPU Server Ring Extension



- Commodity GPU, less HBM dependency
- HBM could be fully used for KV-Cache and computing
- Better parallelism could be achieved