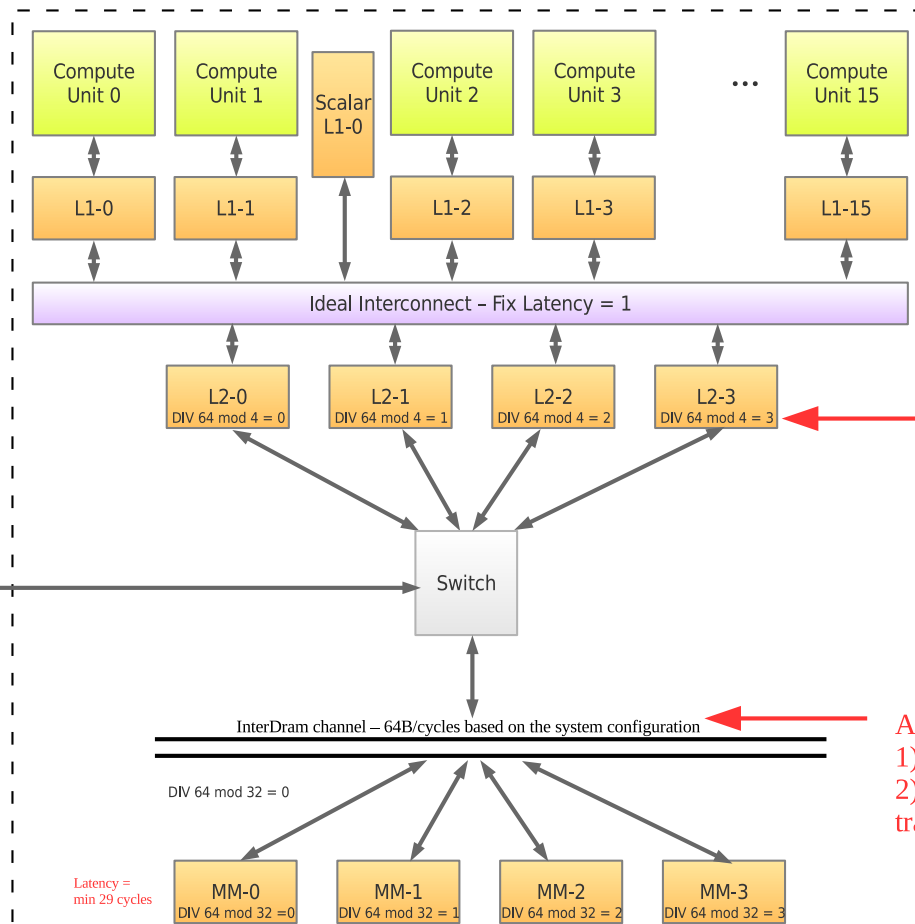


GPU0



The GPUs are connected to each other through either PCIe bus, NVLink, or another network suggested by Kim. Kim network is not considering the **Wire Latency** so first course of action is to calculate this latency, and the NVLink.

* Important. The first data copy is also uses one of these networks. So whatever the latency is we should measure the latency of a mem-copy of each Partition over the connections, and add them up to the whole execution time (in a serialize fashion)

GPU1, GPU2,
GPU3 ,

Yifan, would you review this and see if I am missing something? Thank you limitless, K: