

GMP

1. []

1. [] AI []

1. []
2. []
3. []
4. []

2. []

2. []

1. []

1. [] 2 []
2. [] 4 []
 1. DRAM []
 2. [] LUT []
 3. []

3. []

1. [] **L0 Cache** []

2. []

1. MoE
2. []

3. []

1. []
 1. [] credit [] L1 []
 2. [] NoC []
 3. L2/L3 []
 4. []

2. outstanding / []

3. []
4. launch []

4. []

1. []
2. [] / []
 1. [] [] []
3. [] warp []
4. [] [] []

5. [] KV Cache ..) []

1. []
2. []

6. []

1. **Binary Lut** [] ?
2. **CIM** [] [] ?

7. 1D []

1. sequence
2.
3. LD MUL ST
 ---> DSA
4. L1 dma Mac dma Mac l1
5.
 1. 1D 2D
 1. VLD VST MLD MST VMUL VMUL_reduce MUL_join
 2. LD/ST mbarrier
 1. L1 bank bank count
 count

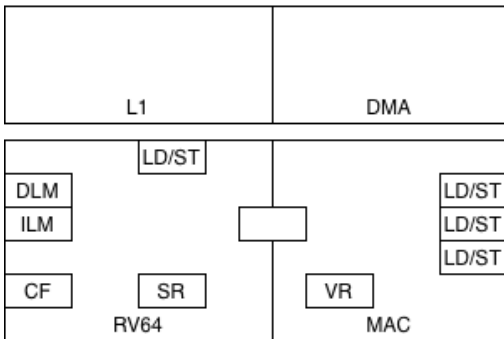
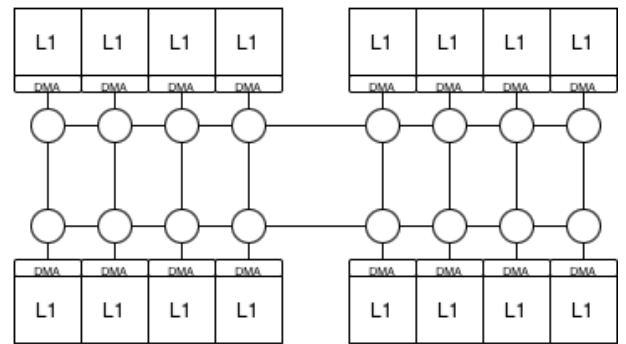
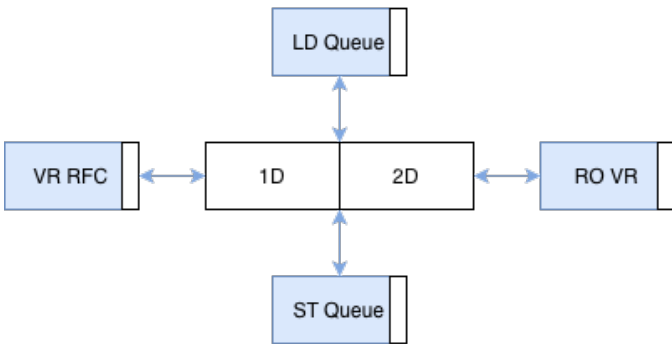
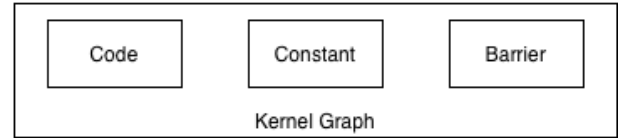
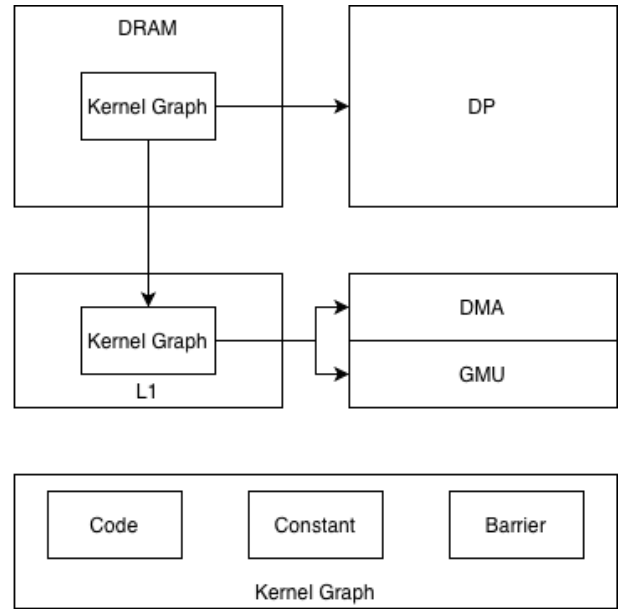
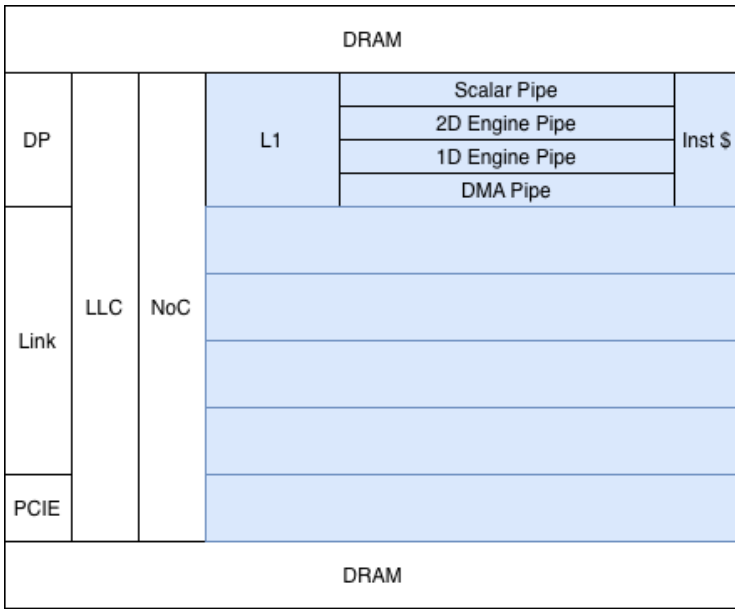
3.
 1.
 2.
 3. edga
 4.
 5.

4.
 1. Fork launch join sync
 2. sync
 3.

5.
6.
 1. load fetch
 2.

7. ISA
8. Launch
9. Sync

1. /
2.



Revision #46

Created 2025-03-16 06:37:22 UTC by Colin

Updated 2026-05-29 11:07:21 UTC by Colin