

GMP

- [Sync And Async](#)
- [AI !\[\]\(1207edb9a08751d3d55970560645ed23_img.jpg\) !\[\]\(d7a34a706cfa4ef37c62a369101e1b36_img.jpg\) 2D !\[\]\(7325769475e8f4bf67f57a0cbebc8ab9_img.jpg\)](#)
- [GMP](#)
- [GPU / GPU](#)
- [GPU GPU](#)
- [GPU](#)
- [DynamicGraphMultiProcessor !\[\]\(1a468f12cdfc63dc07896d0781cf55ec_img.jpg\)](#)
- [GPU !\[\]\(a9a0baec8ceb7d7c04180806eca8d32a_img.jpg\)](#)

Sync And Async

?????

1. []
2. []
3. N to N []
4. []
5. [] pulling [] latency
6. []
 1. transformer [] Flash-attention [] L1 [] fusion [] DMA []
7. []

??Global????????????????

1. []
2. []
3. []
 1. DMA [] Kernel []
4. []

????????????

1. [] Sync [] mailbox [] N to N []
2. [] producer consumer [] Pipeline
3. []
 1. []
 2. []
4. []
5. []
 1. [] write mbx [] wait mbx
[] ack [] latency []

???????????? mbarrier

1. [] M barrier
[]
 1. []
 2. []
2. []
 1. []
[]
 2. [] fence
 3. [] Core0 Store() Core0 Fence() memory Ack() Core0 Signal(Core1 mbx)

1. SRAM
 1. bank
 2. cycle by cycle
 3. latency

NV Blackwell

PTX SASS AI

<https://docs.nvidia.com/cuda/parallel-thread-execution/index.html#tensorcore-5th-generation-family-instructions> <https://docs.nvidia.com/cuda/cuda-binary-utilities/index.html#blackwell-instruction-set>

1. L0 memory
 1. tensor memory 2D/1D NPU/DSA
 2. Tensor memory to memory Tensor
 3. 1D load/store memory
 4. L0 2D/1D L0 fusion code Fusion
2. 128/256 2D
3. TPC 2 SM 2D 2 core
4. OCP-MX micro-scaling
5. thread issue Tensor SIMT style thread
6. weight-stationary GEMM mask bit padding 0

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

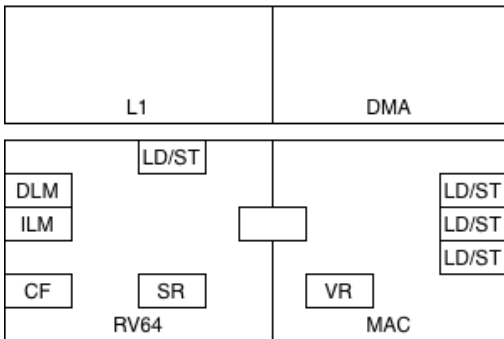
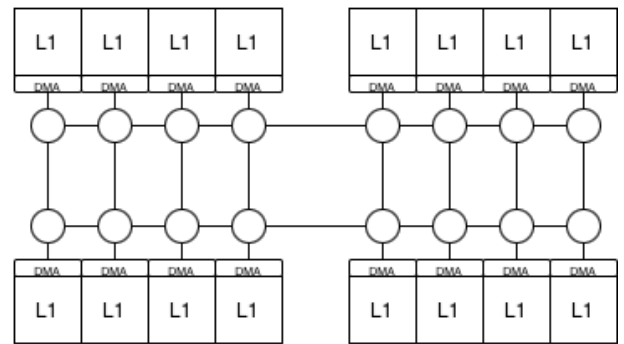
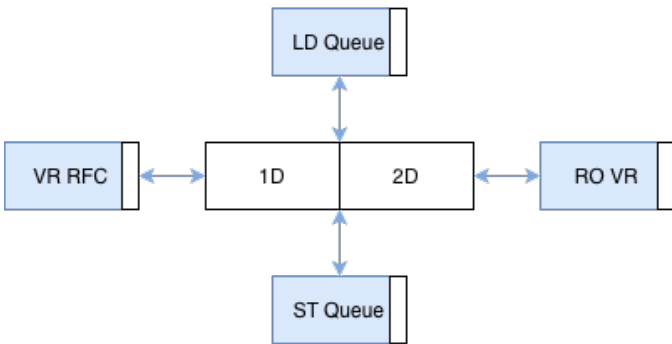
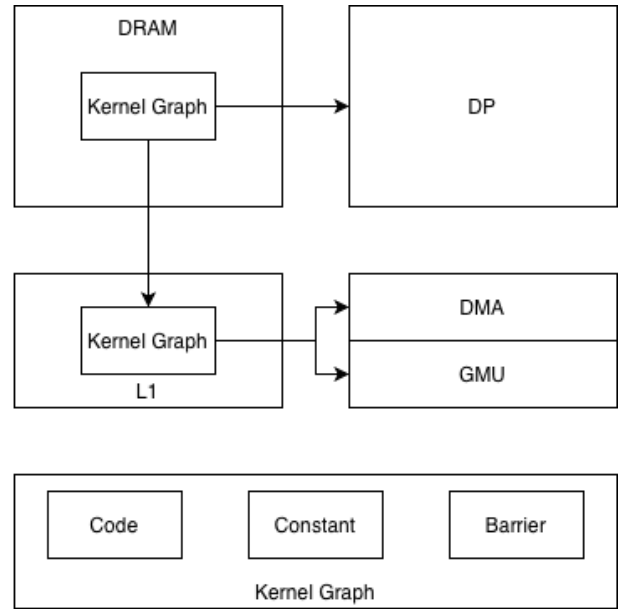
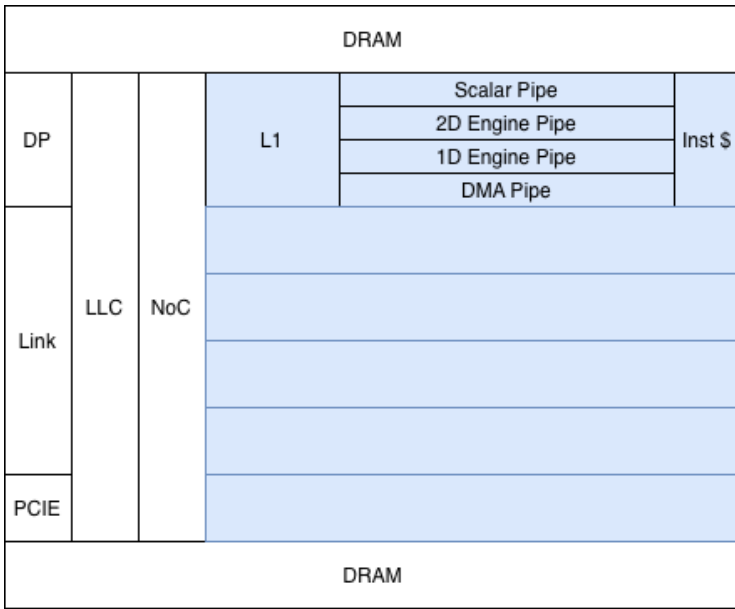
3.
1.
2.
3. edga
4.
5.

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

5.
6.
1. load fetch
2.

7. ISA
8. Launch
9. Sync

1. /
2.



??/???

??

1. []
2. [] SP-PU-L1-DMA []
 1. []
 2. []
 3. []
 4. LD/ST [] fence []
3. Launch [] fork
 1. []
4. launch/signal/wait:join
 1. launch pu instruction: write_back_id local_id
 1. write_back atomic add/sub
 2. wait instruction: local_id
 1. local_id GE LE counter
 3. wait remote instruction: remote_id
 1. local_id GE LE counter
5. [] & []
 1. []
 2. []
 1. [] scale
 2. 1D/2D []
 3. [] Reduce []
 1. L3 Atomic
 3. []
 4. []
 1. leading tailing latency
 1. LD/ST inflight delay [] DMA Delay
 1. ping-pong
 2. pipeline
 2. kernel [] latency
 1. kernel []
 3. kernel [] gap
 1. kernel []
 2. [] fusion
 2. [] ready [] fence []
 1. []
 1. m-barrier
 3. []

MAC

1. Vector []
2. [] VR [] RO WO []
3. [] Vector []
4. [] VR []
5. Vector []

Software

1. []

ISA

1. []
 1. RV64i
2. []
 1. VLD VST VMUL VADD REDUCE_ADD REDUCE_MAX REDUCE_MIN VMUL_REDUCE_ADD
 2. MLD MST
3. []
 1. GEMM
 1. [] / [] 128 [] = 7 [] + 3 [] = 10 []
 2. [] 7 []
 3. Opcode
4. fence
 1. L1 cache line []
5. VR
 1. [] VR data hazard [] VR []
 1. []
 2. [] VR [] [] rename ing [] bank []
 3. [] VR count []
 1. LD ST [] []
 1. []
 2. [] GEMM [] LD []
 3. [] ST []
6. L1
 1. [] cache line
 2. **cache line** []
 3. [] **L1 CacheLine** []
 1. [] **Load** [] **L1 cache line** []
 2. [] **cache line** [] **mask** []
 3. [] **L1-L1** []
 4. [] **MNK** []
7. DMA

1. ROM Launch Kernel kernel
2. cmd
3. DMA kernel linear copy kernel launch
1. launch constructor
8. 128 RV64
- 9.
10. / RiscV-V

??_???

??/???

??

1. []
2. [] SP-PU-L1-DMA []
 1. []
 2. []
 3. []
 4. LD/ST [] fence []
3. Launch [] fork
 1. []
4. launch/signal/wait:join
 1. launch pu instruction: write_back_id local_id
 1. write_back atomic add/sub
 2. wait instruction: local_id
 1. local_id GE LE counter
 3. wait remote instruction: remote_id
 1. local_id GE LE counter
5. [] & []
 1. []
 2. []
 1. [] scale
 2. 1D/2D []
 3. [] Reduce []
 1. L3 Atomic
 3. []
 4. []
 1. leading tailing latency
 1. LD/ST inflight delay [] DMA Delay
 1. ping-pong
 2. pipeline
 2. kernel [] latency
 1. kernel []
 3. kernel [] gap
 1. kernel []
 2. [] fusion
 2. [] ready [] fence []

1.
1. m-barrier
3.

MAC

1. Vector
2. VR RO WO
3. Vector
4. VR
5. Vector

Software

1.

ISA

1.
 1. RV64i
2.
 1. VLD VST VMUL VADD REDUCE_ADD REDUCE_MAX REDUCE_MIN VMUL_REDUCE_ADD
 2. MLD MST
3.
 1. GEMM
 1. $\frac{128}{7} = 18 + 3 = 21$
 2. 7
 3. Opcode
4. fence
 1. L1 cache line
5. VR
 1. VR data hazard VR
 1.
 2. VR rename ing bank
 3. VR count
 1. LD ST
 1.
 2. GEMM LD
 3. ST
6. L1
 1. cache line
 2. **cache line**
 3. **L1 CacheLine**
 1. **Load** **L1 cache line**

????

????????

[[Pipe]] 6 \$208 RTL = engine pipe

0. ? DMA ?????

[[DMA]] "engine pipe + sync pool + fabric master"

	DMA	
	→ transform →	→ →
	L1 DRAM	L1 DRAM
	transform buffer	L0
cmd_desc	xform_flags, xform_params	compute_op, loop, dtype, accum

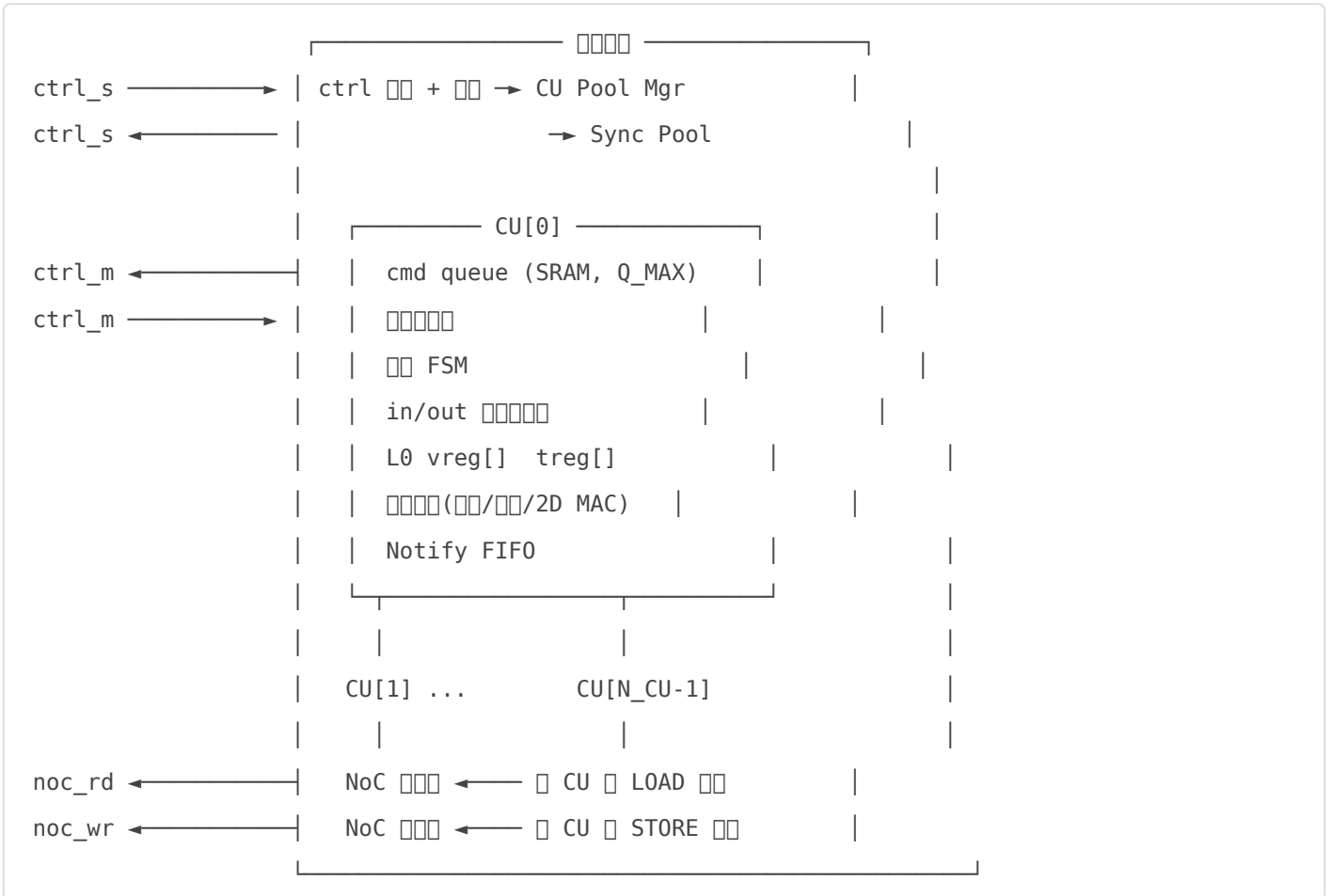
DMA [[DMA]]

1. ????

N_CU	engine pipe =	4	\$13
N_SYNC	sync counter	32	
Q_MAX		16	2
K_IN / K_OUT / K_SIG	cmd_desc	4/4/4	
N_L0_VEC	L0	8	
N_L0_TEN	L0	4	
VEC_LANE	lane	16	INT8 lane
MAC_M / MAC_K	2D MAC	8 / 8	systolic
MAC_ACC_W	MAC	32	INT32
OUT_TYPE_SET		{8, 16, 32}	round/sat
LOOP_W	inner-loop	16	

□□	□□	□□	□□
STRIDE_DIM / STRIDE_W / SHAPE_W	□ [[DMA]]	3 / 32 / 24	
□□ COUNT_W / HANDLE_W / ADDR_W / NOC_W	□ [[DMA]]		

2. ????



3. ????

□□□□ [[DMA]] \$3 □□□□□□□□

- noc_rd_* / noc_wr_* □□□□□□□□ L1 □ DRAM □□□□□□ cmd_desc □□ in/out handle.unit_addr □□
- ctrl_s_* / ctrl_m_* □□□□□□

4. cmd_desc ????

□□ ≤ 1024 bit □□ □□□□□□ □□□□□□ □□□□□□ op + □□□□□□ □

□□	□□	□□	□□
----	----	----	----

0	4	op	0=SCALAR_ALU, 1=VEC_ALU, 2=VEC_REDUCE, 3=MAC_2D, 4=ELEMENTWISE_ACT(ReLU/GELU), 5=SOFTMAX_PARTIAL, ...
4	4	dtype_in	0=I8, 1=I16, 2=I32, 3=FP8(), 4=FP16()
8	4	dtype_out	
12	2	accum_mode	0=, 1= C += A*B, 2= +
14	2	act_kind	ELEMENTWISE
16	16	loop_count	inner-loop step
32	48	src_a_base	byte_addr pipe
80	96	src_a_stride[3]	
176	72	src_a_shape[3]	
248	48	src_b_base	MAC_2D / VEC_ALU
296	96	src_b_stride[3]	
368	48	dst_base	
416	96	dst_stride[3]	
512	72	dst_shape[3]	
584	8	l0_a_reg	A L0
592	8	l0_b_reg	
600	8	l0_c_reg	
608	4+4+4+4	n_in / n_out / n_sig / _	
624	96	in_list[4]	{handle:16, delta:8}
720	96	out_list[4]	
816	96	sig_list[4]	
912	112	reserved/padding	

src_b_* MAC / VEC_ALU reduce 0

5. ???????

[[DMA]] \$5 beat PSEND 8 × 128 bit opcode

6. CU Pool ???

[[DMA]] \$6

- meta[N_CU] cu_class / / MAC
- PALLOC params.cu_class class

7. CU ?????xN_CU?

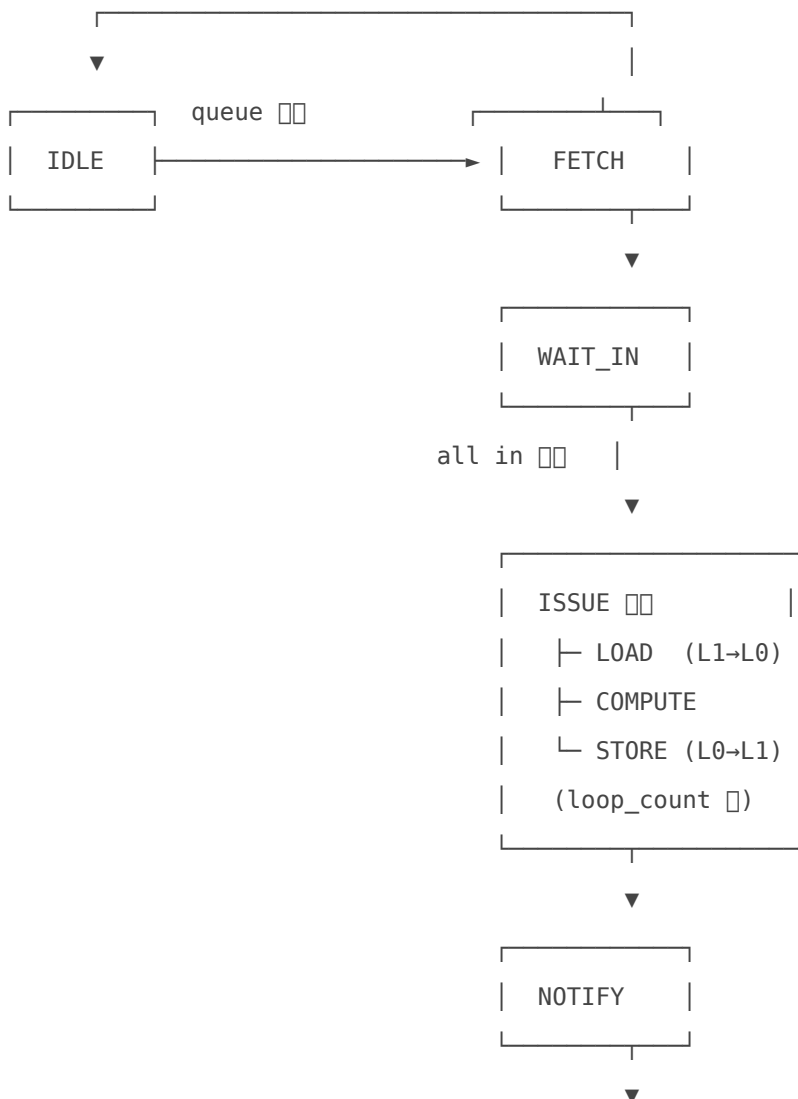
7.1 ????

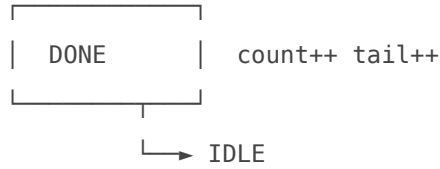
[[DMA]] \$7.1

7.2 ?????

[[DMA]] \$7.2

7.3 ?? FSM





State	Duration	Dependencies
IDLE	1 cycle	
FETCH	1 + (STRIDE_DIM-1) cycle	cmd_desc + stride wrap_delta DMA §7.4
WAIT_IN		COUNT_QUERY in_list
ISSUE	§7.6	
NOTIFY	(n_in + n_out + n_sig) cycle	Notify FIFO
DONE	1 cycle	

7.4 LO ????

Register	Capacity
vreg[N_LO_VEC]	VEC_LANE × dtype_in_bits vector /
treg[N_LO_TEN]	MAC_M × MAC_K × dtype_in_bits tile
	vreg: 2R/1W treg: 2R/1W FMA
	undefined LOAD
	CU

vreg register file small treg SRAM macro treg ≈ 8×8×8 = 512 bit N_LO_TEN = 4 × 2 KiB treg SRAM 2R/1W MAC 2 1

vreg treg

7.5 ????

cu_class

7.5.1 ??

- 1 lane max(dtype_in, dtype_out) ≤ 32 bit
- op ADD/SUB/MUL/SHIFT/CMP/SELECT/BIT
- latency 1 cycle throughput 1/cycle

7.5.2 ??

- VEC_LANE INT8 ALU + 16-bit accumulator option
- op lane-wise ADD/SUB/MUL/MAC/CMP lane REDUCE_SUM/MAX/MIN
- latency 1-2 cycle throughput 1/cycle
- elementwise activation LUT 4 KiB ROM/SRAM ReLU / GELU / sigmoid

7.5.3 2D MAC

- `MAC_M × MAC_K` systolic `[[[`
- `[[[` `INT8 × INT8 → INT16 partial → [[[` `MAC_ACC_W [[[`
- `[[[` A treg (`M × K`), B treg (`K × N`) `[[[` `N = MAC_M`) `[[[` `cmd_desc`
`[[[`
- `[[[` C treg (`M × N`) `[[[` treg `accum_mode`
- pipeline depth = `MAC_K + 2` cycle
- `[[` cycle `[[[` A `[[[` B `MAC_M = MAC_K = 8` `[[` 8 cycle `[[[` 9 cycle `[[` cycle
`[[[`

`[[` / `[[` / round `cmd_desc.dtype_out` `[[` `accum_mode` `[[[` saturate + truncate /
round_to_nearest `[[[` stage `[[[`

7.6 LOAD/COMPUTE/STORE inner-loop ??

`[[[` `cmd` `[[` `loop_count` `[[` micro-op `[[[`

cycle:	N	N+1	N+2	N+3	N+4	N+5
LOAD	op0	op1	op2	op3	op4	op5
COMPUTE		op0	op1	op2	op3	op4
STORE			op0	op1	op2	op3

- L0 `[[[` buffer `LOAD` `[[` `COMPUTE` `[[` + `[[` `STORE` `[[`
- `[[[` `l0_a_reg / l0_b_reg / l0_c_reg` `[[[` op `[[` reg `[[[` `op_id % N_L0_REG` `[[`
cyclic `[[[`
- leading 2 cycle `[[` `LOAD` `COMPUTE/STORE` `[[[` tailing 2 cycle `[[` `STORE`
- inner-loop `[[[` / `[[[` stall `[[` NoC backpressure `[[`

7.7 in / out ?????

`[[[` 3 `[[[` A `[[` B `[[` C `[[[` `[[DMA]] $7.4` `[[[`

- `[[[` L1 / DRAM `[[[`
- inner-loop `[[` `loop_count` `[[[` `shape[3]` `[[[`
- `[[[` NoC `[[` burst `[[` = `[[` L0 `[[[` / `NOC_W` `[[[` 1 beat `[[`

7.8 Notify FIFO

`[[[` `[[DMA]] $7.6`

8. Sync Pipe Pool

`[[[` `[[DMA]] $8`

9. ???????

9.1 NoC ? / ?

[[DMA]] §9.1 NoC [] [] L1 / DRAM [] L1 [] outstanding []
DRAM [] DRAM []

9.2 ctrl_m ??

[[DMA]] §9.2 []

9.3 ctrl_s ????

[[DMA]] §9.3 [] PALLOC [] `params.cu_class` [] `cu_class` []

10. ??????????

10.1 CU ? pending wait CAM

[[DMA]] §10.1 []

10.2 CU ? ctrl_m

[[DMA]] §10.2 []

10.3 CU ? NoC?????????

- LOAD CU [] `noc_rd_req` [] L1 / DRAM [] master [] §7.7 [] `resp` []
`req_id` [] CU [] L0 []
- STORE CU [] `noc_wr_req + wr_data` [] `ack` [] fire-and-forget [] NOTIFY []
`COUNT_DELTA` [] CU [] FSM [] STORE [] NoC []
`id` []

10.4 ? cmd_desc ? in ? out ??? pipe

[[Pipe]] §289 []

1. WAIT_IN [] pipe [] `count` \geq `in` []
2. LOAD []
3. COMPUTE
4. STORE []
5. NOTIFY [] pipe [] `COUNT_DELTA(-in_delta)` [] `COUNT_DELTA(+out_delta)`

[] In-place [] NoC [] `id` []

11. ?????

[] / []	[]
[[DMA]] §11 []	
L0 vreg / treg	[] [] undefined []

□□ / □□	□□□
□□□□ pipeline □□□	□□
MAC □□□	0

L0 □□□□□□□□ LOAD □ COMPUTE □□□ bug□□□□□□□□□□□□

12. ???????

□ [[DMA]] §12□□□

□□	□□	□□
dbg_cu_state[N_CU]	$N_CU \times 3$	□ CU FSM □□
dbg_mac_active[N_CU]	N_CU	MAC □□□□
dbg_inner_loop_idx[N_CU]	$N_CU \times LOOP_W$	□□ micro-op □□
dbg_l0_read_hazard[N_CU]	N_CU	LOAD/COMPUTE □ reg □□□□□□ 0□
dbg_pipeline_stall_cycles[N_CU]	$N_CU \times 16$	NoC backpressure □□ stall

13. ???????

1. N_CU □□□□□□□□ 1 □ 2D MAC □□ + 2 □□□□□□ + 1 □□□□□□ vs □□□
2. $VEC_LANE / MAC_M / MAC_K$ □□□□ — □□□□□□
3. dtype □□□□□□ INT4 / FP8 / FP16 / BF16□
4. MAC □□□□□□ systolic 1D vs 2D□ weight-stationary vs output-stationary□
5. L0 □□□□□□□□□□ FMA □ 3R+1W □□□□□□ 2R+1W□
6. inner-loop □□□□□□ 3 □ vs 4 □□□□□□□□□□ stall □□□
7. □□ round / saturate □□□□□□
8. □□□□□□□□ LUT □□□□□□ per-CU □□

14. ??? Pipe ??

□□□□□	□□□□□
engine pipe □□	§7.1 □□ + §7.2 count + §7.3 FSM
sync pipe □□	§8
□□ A	§6 + §8 □□□□□□ bitmap
□□ B	§3 ctrl vs noc □□□□
□□ D	§8 sync □□
□□ 3□□ pipe □□□□	§7.3 FSM □ channel □□□□
□□ 6□□□□□□□	§5 / §6 / §7.4 L0 □□□□□□
□□ 7□□□□□□□□	§6 alloc □ cu_class / queue_depth □□

□□□□	□□□□
□□ 8□□□ pipe□	engine □ \$7□ sync □ \$8
\$289□ in/out □ pipe□	\$10.4 NOTIFY □□

DynamicGraphMultiProcessor??

Dynamic Graph Multi Processor ??

??

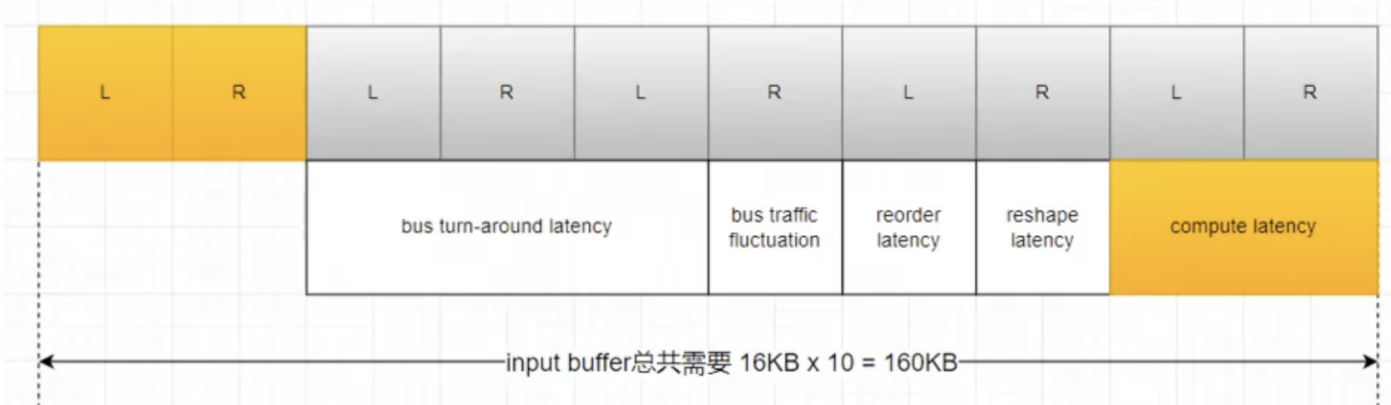
1. Etched GPU TFLOPS
2. + +
3. Etched H100 800 3.3%
Transformer Sohu FLOPS
90% GPU 30%
4. Scaling Law Scale Up
Scaling Law Scale Down

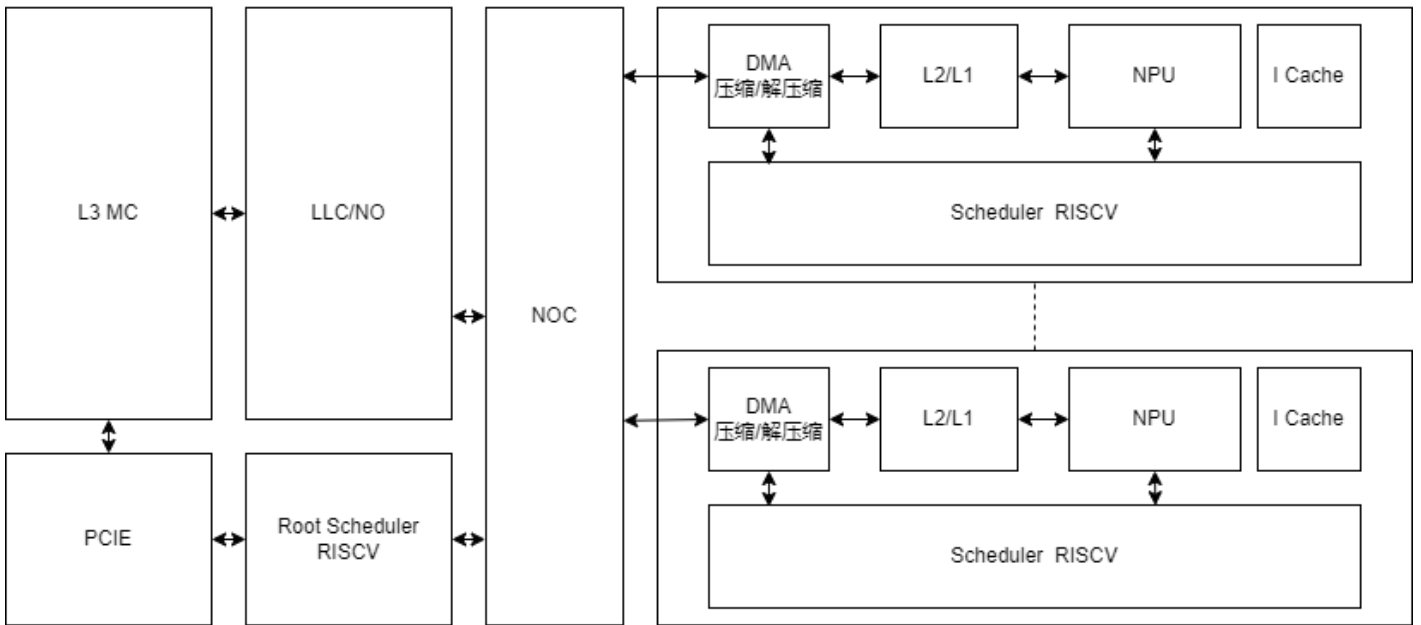
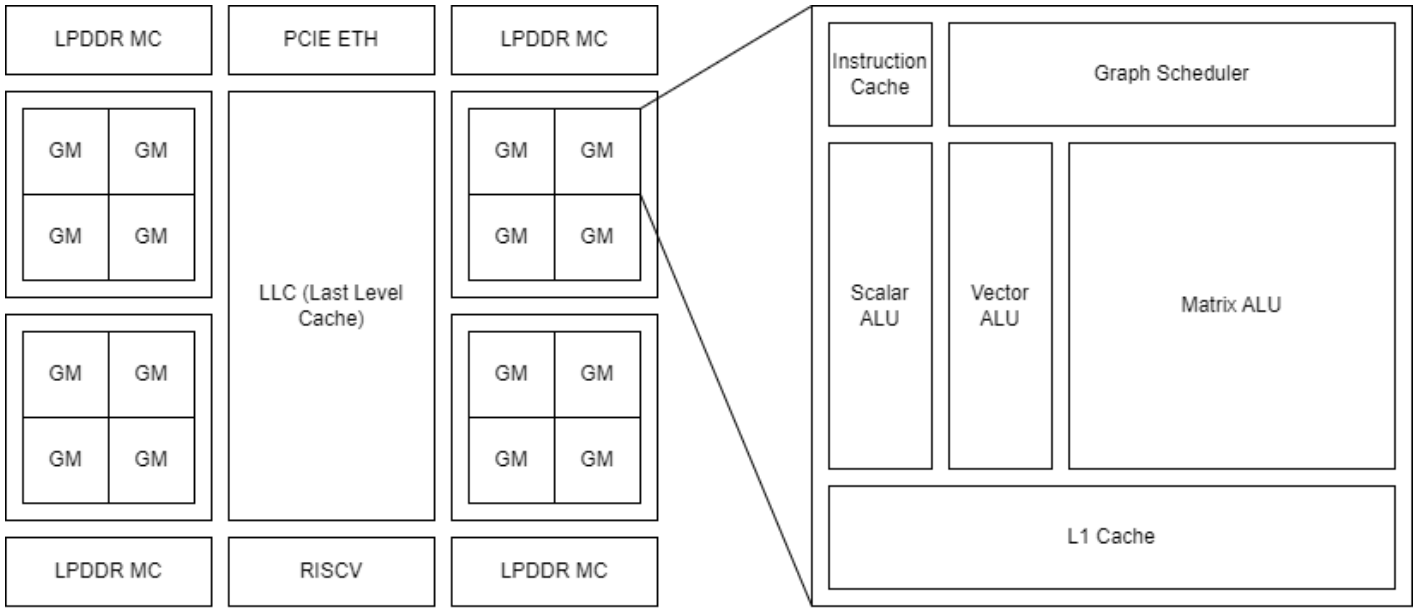
????

- AI
- Int8 pattern
- 2D =>
- GPU DRAM
- NoC Cache Fork/Join
NoC flow control
- NoC/Cache
- consistency coherency fence+sync
- 1. GMP fork/join
- 2. fence
- 3. fence
- IP cpu cache
IP IP
cache_hint
- DSA
- SOC
- NOC
- x

3. Launch/Sync/DataFlow

1. Launch Sync Fork/Join
 1. **Fork** **Join**
 2. fork launch join wait
 3. fork join
 4. fork
2. Launch Sync
3. NoC Launch
4. Graph Kernel
5. **latency**
4. launch
 1. -> -> in flight pipeline
 2.
 3.
 4.
 5. SRAM
 1. SRAM L1 L1
 2. L1 cache GEMM
 3.
 4. L1 bank thread
 5. Engine ALU RegRead MALU Id Cal
 6. St
6. graph launch
7. DSA
 1. PPA gpu gs
 2. simt simi
 3. data structure hazard
8. DSA
 1. fence TensorALU LD ST
 2. latency standby





1. `fork(kernel, fork(gemm, for launch))`
2. `kernel, fork(gemm, for launch)`
2. `32bit, + +2D`

2. launch
3. launch launch
4. launch sub graph graph group
"signal"
9. delay
1. data hzd
10. stall stall

??

1. Fork Join

1. fork
 - 1.
 2. join
 3. index
 4. launch ID ID IP memory mapping
2. fork load
3. join
 1. launch ID
2. VR L0.5 L0.5
3. VR L0 4
4. TR L0.5, L0.5
5. DTE transpose pad slice deslice
6. PU
- 7.

1.

Width (bits)	4	6	3	3	1	
Meaning	Reuse flags	Wait barrier mask	Read barrier index	Write barrier index	Yield flag	St cycl

2. Reuse flags 4 register cache
3. 6 barrier thread cuda
6
8. async_group
 1. async copy bulk
9. mbarrier fence
10. credit valid/ready
 - 1.
 - 2.
 - 3.

????????

design/ .md [[LLC]] + LLC.pipeline.html

0. ??

- RTL / [logix] .md RTL
- .md L1.md DMA.md LLC.md
- HTML <unit>.pipeline.html git LFS

1. ?????

.md ##### N.

\$1		/ /
\$2		ASCII mermaid
\$3		/ /
\$4		
\$5		opcode
\$6 - \$M		/
\$M+1		\$2
\$M+2 -	Pin / Lock / Bypass DGMP	

H4 ##### N. H5 ##### N.M H6 ##### N.M.K

2. ?????

" " N LLC.md \$11

2.1 ??

- §a-§b
- SRAM 1RW vs 1R1W
- `<unit>.pipeline.html`
- Mermaid / WaveDrom / YAML

2.2 §N.1

<code><name></code>	<code><arb / sram_port / ff_rw / cam / table / comb></code>	<code>§x.y</code>

SRAM FF

2.3 §N.2

mermaid `flowchart LR` stage

2.4 §N.3 ~ §N.X

H6 ##### §N.X.Y

1. mermaid `flowchart LR` stage +
2. WaveDrom JSON cycle stage / latch / edge
3. **latch** stage→stage latch +
4. ×

2.5 §N.{?? 2} ????????????

- × / /
- **stall** stall
- **bank**

2.6 §N.{?? 1} ??? schema

YAML logix / RTL schema

```

unit: <name>
clock: <main_clk>

resources:

```

```
<name>: { kind: ..., ports/slots/width/banks: ... }
```

paths:

```
<path_name>:
```

```
stages:
```

```
- { id: <S0/S1/...>, in: [...], uses: [...], out: [...] }
```

arbitration:

```
- { resource: <res>, contenders: [...], rule: "<text>" }
```

stall_conditions:

```
- { src: <res.cond>, effect: "<text>" }
```

XXXXXXXXXXXXXXXXXXXX

<unit>.pipeline.yaml XXX

3. ???????

□	□	□
□□	S0 / S1 / S2 / S3 / S4	□□
□□□□	S3' / S4'□□□□	□□□ S2 □□ mshr
□□□□	F0 / F1 / F2	□□ fill□
□□ FSM	W0 / W1 / W2 / W3	□□□ walker
□□□□	B0 / B1	□□ bypass
□□	S2'□□□□□□□□ + □□	□□ master □□ S2

XXXXXXXXXXXXXXXXXXXX

4. ?????????

- □□□ □ \$N □ \$N.M □ \$N.M.K
- □□ □ [[□□□]] \$N □ Obsidian wiki link □□□□□□□□ .md □□
- □□□□□□□□□□ "□□□□ □[[L1]] \$3.2 □ NoC slave □□ □□□□□□□□

5. ?????????

□	□	□
Mermaid <code>flowchart</code>	□□□□□□ + □□	□□ GitHub / Obsidian / VS Code □□□□
WaveDrom JSON	□□□□□□ + □□□□ + □□□□	□□ wavedrom.com/editor.html □□ Obsidian + WaveDrom □□

