

## 承影 GPGPU 架构文档手册 v1.95

编写人：杨轲翔 (yangkx20@mails.tsinghua.edu.cn)

清华大学集成电路学院 dsp-lab

承影是基于 RISC-V 向量扩展实现的开源 GPGPU，开源项目主页面见 [THU-DSP-LAB/ventus-gpgpu: GPGPU processor supporting RISCV-V extension, developed with Chisel HDL \(github.com\)](https://github.com/THU-DSP-LAB/ventus-gpgpu)

本文档主要描述了承影 GPGPU 在软硬件视角下的功能描述，部分功能直接以承影下一版本标定，本文档中会备注说明目前未支持到的情况。

如果在 opengl 描述、驱动、编译器功能上有问题，或者硬件设计上不足之处，欢迎在 github 上提 issue，或邮件联系作者。

### 简介

本文档描述了承影 GPGPU 的设计内容，包括 OpenCL 编程视角和微架构视角。承影 GPGPU 指令集以 RISC-V 向量扩展（后文简称为 RVV）为核心设计 GPGPU，相比 RISC-V 标量指令，具有更丰富的表达含义，可以实现访存特性表征、区分 workgroup 和 thread 操作等功能。核心思想是在编译器层面以 v 指令作为 thread 的行为描述，并将 thread->warp/workgroup 的公共数据合并为标量指令。硬件上一个 warp 就是一个 RVV 程序，通常向量元素长度为 num\_thread，同时又将 workgroup 中统一执行的公共地址计算、跳转等作为标量指令执行，即 Vector-Thread 架构。硬件将 warp 分时映射到 RVV 处理器的 lane 上去执行。相比其它 SIMT 架构，在硬件上的折中是无法实现完全的 per thread per pc，仍然需要以 workgroup（或分支状态下的 warp\_split）执行。RVV 指令集在变长上有三个方面的体现：硬件 vlen 改变；SEW 元素宽度改变；LMUL 分组改变。本架构特点在于这三个参数在编译期都已固定，元素数目大部分情况也固定为 num\_thread，**本架构本质上是 SIMT。**

## 术语表

- SM: streaming multiprocessor, 流多处理器单元
- sGPR: scalar general purpose register, 标量寄存器
- vGPR: vector general purpose register, 向量寄存器

memory 划分和编程模型定义: (在本文档中局部内存、共享内存可能同时使用 localmem 和 sharedmem 来指代; 线程束和线程采用 cuda 的说法 warp 和 thread; 其它词均采用 opengl 中的说法)

cuda	opengl	解释
globalmem	globalmem	全局内存, 用 __global 描述, 可以被 kernel 的所有线程访问到
constantmem	constantmem	常量内存, 用 __constant 描述, 是全局地址空间的一部分
localmem	privatemem	私有内存, 各 thread 自己的变量, 和内核参数, 是全局内存的一部分
sharedmem	localmem	局部内存, 用 __local 描述, 供同一 work-group 间的线程进行数据交换
grid	NDRange	一个 kernel 由多个 NDRange 组成, 一个 NDRange 由多个 workgroup 组成
block/CTA	workgroup	工作组, 在 SM 上执行的基本单位
warp	wavefront(AMD)	32 个 thread 组成一个 warp, 仅对硬件可见
thread	work-item	线程/工作项, 是 OpenCL C 编程时描述的最小单位。

## 参数表

变量名	解释
num_thread	一个 warp 里的 thread 数, 默认值 32
num_warp	硬件上一个 SM 里允许的最大 warp 数 (可以来自不同

workgroup)  
num\_block 硬件上一个 SM 里允许的最大 workgroup 数  
num\_lane 硬件上一个 SM 的运算单元里一次能同时处理的 thread 数  
localmem\_max 硬件上一个 SM 里提供的 localmem 的最大空间

因此,  $\text{num\_thread} \times \text{num\_warp}$  就代表了一个 SM 里的最大 thread 数目。

## 编程模型和驱动程序功能

从 OpenCL 视角来看这个 device。

### 硬件上的对应关系

整个 GPU 作为一个 compute device, SM 对应 Compute Unit(CU), SM 内部多个执行单元对应多个 PE。

### 任务执行模型

与 OpenCL 一致, 将 workgroup 映射到 CU 上执行, 各个 thread 映射到 PE 上执行, 硬件上会将 thread 以 warp(相邻 32 个 thread 一组)为单位打包, 呈现出 SIMD 的执行效果。目前 NDRange 拆分为 workgroup 在驱动上进行, workgroup 拆分为 warp 在硬件上进行。

### 驱动提供的功能

由 opencl 驱动(poctl)来管理 command queue, 创建和分配 buffer, poctl 以 kernel 为单位分配任务, 并为每个任务创建 metadata buffer。共享内存空间、任务间的顺序和事件同步机制也由 poctl 管理。poctl 传递给硬件驱动后, 硬件驱动将以 workgroup 为单位, 把任务发送到 CTA-scheduler 处理。

在 poctl 后端添加基于 verilator 的 ventus device 和 ISS spike ventus device, 以完成物理地址分配和任务启动。

目前 poctl 创建的 buffer 包括:

- NDRange 的 metadata buffer 和 kernel 程序



- kernel 的 argument buffer
- kernel argument 中显式引用的 buffer
- 为 private mem、print buffer 分配的空间

任务启动时，由硬件驱动直接传递的信号为：

- PTBR // page table base addr
- CSR\_KNL // metadata buffer base addr
- CSR\_WGID // 当前 workgroup 在 SM 中的 id，仅供硬件辨识
- CSR\_WID // warp id，当前 warp 属于 workgroup 中的位置
- LDS\_SIZE // localmem\_size，编译器提供 workgroup 需要占用的 localmem 空间。privatemem\_size 默认按照每个线程 1kB 来分配。
- VGPR\_SIZE // vGPR\_usage，编译器提供 workgroup 实际使用的 vGPR 数目 (对齐 4)
- SGPR\_SIZE // sGPR\_usage，编译器提供 workgroup 实际使用的 sGPR 数目 (对齐 4)
- CSR\_GIDX/Y/Z // workgroup idx in NDRange
- host\_wf\_size // 一个 warp 中 thread 数目
- host\_num\_wf // 一个 workgroup 中 warp 数目

## runtime 行为

约定 kernel 启动时，NDRange 的参数通过 metadata 的 buffer 传递，该 buffer 的内容为：

```
cl_int clEnqueueNDRangeKernel(cl_command_queue command_queue,  
                               cl_kernel kernel,           //kernel_entry_ptr & kernel_arg_ptr  
                               cl_uint work_dim,           //work_dim  
                               const size_t *global_work_offset, //global_work_offset_x/y/z  
                               const size_t *global_work_size,  //global_work_size_x/y/z  
                               const size_t *local_work_size,   //local_work_size_x/y/z  
                               cl_uint num_events_in_wait_list,  
                               const cl_event *event_wait_list,
```

```
cl_event *event)
```

```
/*
```

```
#define KNL_ENTRY 0
```

```
#define KNL_ARG_BASE 4
```

```
#define KNL_WORK_DIM 8
```

```
#define KNL_GL_SIZE_X 12
```

```
#define KNL_GL_SIZE_Y 16
```

```
#define KNL_GL_SIZE_Z 20
```

```
#define KNL_LC_SIZE_X 24
```

```
#define KNL_LC_SIZE_Y 28
```

```
#define KNL_LC_SIZE_Z 32
```

```
#define KNL_GL_OFFSET_X 36
```

```
#define KNL_GL_OFFSET_Y 40
```

```
#define KNL_GL_OFFSET_Z 44
```

```
#define KNL_PRINT_ADDR 48
```

```
#define KNL_PRINT_SIZE 52
```

```
*/
```

kernel 的参数由另一块 kernel\_arg\_buffer 传递，该 buffer 中会按顺序准备好 kernel 的 argument，包括具体参数值或其它 buffer 的地址。在 NDRange 的 metadata 中仅提供 kernel\_arg\_buffer 的地址 knl\_arg\_base。

kernel 函数执行前会先执行 start.S:

```
# start.S
```

```
start:
```

```
csrr sp, CSR_LDS # set localmemory pointer
```

```
addi tp, x0, 0 # set privatememory pointer
```

```
# clear BSS segment
```

```
#
```

```
# clear BSS complete
```

```
csrr t0, CSR_KNL
```

```
lw t1, KNL_ENTRY(t0)
```

```
lw a0, KNL_ARG_BASE(t0)
jalr t1
# end.S
end:
endprg
```

约定 kernel 的打印信息通过 print buffer 向 host 传递。print buffer 的地址和大小在 metadata\_buffer 中提供，运行中的 thread 完成打印后，将所属 warp 的 CSR\_PRINT 置位。host 轮询到有未处理信息时，将 print buffer 从设备侧取出处理，并将 CSR\_PRINT 复位。

## 架构说明

硬件上的 ABI、指令集、寄存器接口的部分，以及对内存系统的说明，以应对 OpenCL kernel 的编程需求。

## 指令集范围

### RV32V

实际选择的指令集范围为：RV32 I M A zfinx zve32f

V 里面支持的主要是独立数据通路的指令，当前 RVV 原有的 shuffle widen narrow gather reduction 都不支持。

下表列举出了目前支持的标准指令的范围，有变化的指令已经声明。

	承影支持情况	指令变化
RV32I	不支持 ecall ebreak, 支持 SV39 虚拟地址	
RV32M F	支持 RV32M zfinx zve32f	
RV32A	支持	
RV32V-Register State	仅支持 LMUL=1 和 2	



RV32V-ConfigureSetting	支持计算 vl, 可通过该选项配置支持不同宽度元素	
RV32V-LoadsAndStores	支持 vle32.v vlse32.v vluxe32.v 访存模式	vle8 等指令语义改为 “各 thread 向向量寄存器元素位置写入”, 而非连续写入
RV32V-IntegerArithmetic	支持绝大多数 int32 计算指令	vmv.x.s 语义改为 “各 thread 均向标量寄存器写入”, 而非总由向量寄存器 idx_0 写入, 多线程同时写入是未定义行为, 正确性由程序员保证; vmv.s.x 语义改为与 vmv.v.x 一致
RV32V-FixedPointArithmetic	添加 int8 支持, 视应用需求再添加其它类型	
RV32V-FloatingPointArithmetic	支持绝大多数 fp32 指令, 添加 fp64 fp16 支持	
RV32V-ReductionOperations	视应用和编译器需求再考虑添加, 例如需要支持 OpenCL2.0 中的 work_group_reduce 时	
RV32V-Mask	支持各 lane 独立计算和设置 mask 的指令	vmsle 等指令语义改为 “各 thread 向向量寄存器元素位置写入”, 而非连续写入
RV32V-Permutation	不支持, 视应用和编译器需求再考虑添加	

RV32V- 不支持，视应用和编译器需  
ExceptionHandling 求再考虑添加

## 自定义指令

### *barrier* 线程同步指令

`barrier x0,x0,imm # meet barrier`

`barriersub x0,x0,imm # barrier for subgroup`

`barrier` 对应 `openc1` 的 `barrier(cl_mem_fence_flags flags)` 和 `work_group_barrier(cl_mem_fence_flags flags, [memory_scope scope])` 函数，实现同一 `workgroup` 内的 `thread` 间数据同步。`memory_scope` 缺省值为 `memory_scope_work_group`。

`imm` 为 5bit，具体编码如下：

<code>imm[4:3]</code>	00	01	10	11
<code>memory_scope</code>	<code>work_group</code> (default)	<code>work_item</code>	<code>device</code>	<code>all_svm_devices</code>
<code>imm[2:0]</code>	<code>imm[2]=1</code>	<code>imm[1]=1</code>	<code>imm[0]=1</code>	000
<code>CLK_X_MEM_FENCE</code>	<code>IMAGE</code>	<code>GLOBAL</code>	<code>LOCAL</code>	<code>USE_CNT</code>

开启 `_openc1_c_subgroups` feature 后，则改为 `barriersub` 指令，对应 `memory_scope=subgroup` 的情况，此时 `imm[4:3]` 固定为 0，`cl_mem_fence_flags` 为 `imm[2:0]`，与 `barrier` 指令一致。

### *endprg* 任务结束指令

`endprg x0,x0,x0 # meet the end of the kernel`

需要显式插入到 `kernel` 末尾，表明当前 `warp` 执行结束。只能在无分支的情况下使用。

### *vbeq/join* 线程分支控制指令

`vbeq vs2, vs1, offset # set predicate vs2==vs1, and set branch address  
pc+4+offset`



join v0, v0, offset # set join address pc+4+offset 隐式 SIMT-stack 实现分支控制。每个分支块的结尾需要用 join 指示汇合点地址。join 的源操作数默认为 0。vbeq 参照 beq 提供了 vbne vblt vbge vbgtu vbgeu 版本，指令编码修改了 func3 段。

### *regext{i}* 寄存器扩展指令

regext x0,x0,imm12 # (x3,x2,x1,x0)=imm12, register index extend  
regexti x0,x0,imm12 # (imm[10:5],x2,x0)=imm12, imm and register index extend  
用宏指令扩展可用寄存器数目，该指令表明下一条指令的寄存器编号（和立即数）会扩展。

regexti 只对 V 扩展的 VI 类 OP-imm5 指令生效，将其视为 imm11。对其它立即数指令不支持立即数扩展，可用 regext 进行寄存器扩展。

当前版本编译器对寄存器扩展指令，支持按需要进行扩展。对于除自定义指令外的立即数指令，默认使用 11 位立即数，即认为由 regexti + vi 指令总是组成 64bit 长指令。在编译器视角下，立即数指令将跳过 regext 这一阶段。

### *vlw.v/vsw.v privatemem* 访存指令

仅用于访问 privatemem，以标量访存指令形式提供的指令，但访问向量地址、写入向量寄存器。为便于编译器和编程等使用，在 per-thread 视角下该地址为 0-1k 的连续地址，由硬件根据 thread\_id 和 CSR\_PDS 完成偏移。vlw.v vd, imm11(rs1) # vd <- mem[(rs1+imm11)\*num\_thread\*num\_warp+thread\_idx]  
vsw.v vs2, imm11(rs1) # mem[(rs1+imm11)\*num\_thread\*num\_warp+thread\_idx] <- vs2

### *vlw12.v/vsw12.v* 带 12 位立即数地址的向量访存指令

vlw12.v vd, imm12(vs1) # vd <- mem[vs1+imm12]  
vsw12.v vs2, imm12(vs1) # mem[vs1+imm12] <- vs2 vlw12 参照 lw 提供了 vlh12 vlb12 vlhu12 vlbu12 版本，vsw12 也有 vsh12 vsb12 版本。

### *vadd12.vi* 带 12 位立即数向量整数加减指令

vadd12.vi vd, vs1, imm12 # vd <- vs1 + imm12

## vftta.vv 浮点卷积指令

vftta.vv vd, vs2, vs1, v0.mask # vd <- vs2 conv vs1 + vd

## vfexp.v 浮点指数指令

vfexp.v vd, vs2, v0.mask # vd <- exp(vs2)

完整指令编码及描述如下:

31	25	24	20	19	15	14	12	11	7	6	0	assemble	description	type
off[12:10:5]	vs2	vs1	0 0 0	0 0 1	off[4:1:11]	1 0 1 1 0 1 1	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	vbeq vs2, vs1, offset vbneg vs2, vs1, offset vblt vs2, vs1, offset vbge vs2, vs1, offset vbltu vs2, vs1, offset vbgeu vs2, vs1, offset	if(vs2 == vs1) PC(in_stack) += sext(offset) if(vs2 != vs1) PC += sext(offset) if(vs2 < vs1) PC += sext(offset) //signed if(vs2 >= vs1) PC += sext(offset) //signed if(vs2 < u vs1) PC += sext(offset) //unsigned if(vs2 >= u vs1) PC += sext(offset) //unsigned	分支控制指令
off[12:10:5]	vs2	vs1	0 1 1	0 1 1	off[4:1:11]	1 0 1 1 0 1 1	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	join vs2, vs1, offset	分支汇合, vs1和vs2默认为零	
imm[11:0](x3,x2,x1,x0)		rs1	0 1 0	0 1 1	rd		0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	regext x0,x0,imm regexti x0,x0,imm endprg x0,x0,x0	扩展寄存器编号。rs1和rd默认为零 扩展vop.vi指令的寄存器和立即数。rs1和rd默认为零 kernel运行结束	寄存器扩展指令
0 0 0 0 0 0 0 0	rs2	imm[4:0]	1 0 0	rd		0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	barriersub x0,x0,imm	对应barrier(),imm提供memory_scope和cl_mem_fence_flags	同步和任务控制指令
0 0 0 0 0 0 1 1		imm[4:0]	0 0 0	vd		0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	vaddi2.vi vd, vs1, imm	vd = vs1 + imm	自定义计算指令
0 0 0 0 1 0 m	vs2	0	1 1 0	vd		0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	vfexp vd, vs2, v0.mask	vd = exp(vs2)	
0 0 0 0 1 1 m	vs2	vs1	1 0 0	vd		0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	0 0 0 0 1 0 1 1	vftta.vv vd, vs2, v0.mask	vd = v2 conv v1 + vd	
imm[11:0]		vs1	0 1 0	vd		1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	vld12.v vd, offset(vs1) vld12.v vd, offset(vs1) vld12.v vd, offset(vs1) vld12.v vd, offset(vs1)	vdc-mem[addr], addr=vs1+offset	自定义访存指令, 对目标量访存的长立即数版本
imm[11:5]	vs2	vs1	0 1 1	imm[4:0]	1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	1 1 1 1 0 1 1	vsh12.v vs2, offset(vs1) vsh12.v vs2, offset(vs1) vsh12.v vs2, offset(vs1)	mem[addr]<-vs2, addr=vs1+offset	
0	imm[10:0]	rs1	0 0 1	vd		0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	vld1.v vd, offset(rs1) vld1.v vd, offset(rs1) vld1.v vd, offset(rs1) vld1.v vd, offset(rs1)	vdc-mem[addr] addr=(rs1+imm)*num_thread_in_wg+thread_idx funct3字段编码与1w 1h 1b等访存指令一致	自定义访存指令, 专用于访问
1	imm[10:5]	rs1	0 1 0	imm[4:0]	0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	0 1 0 1 0 1 1	vsw12.v vs2, offset(rs1) vsh.v vs2, offset(rs1) vsh.v vs2, offset(rs1) vsh.v vs2, offset(rs1)	mem[addr]<-vs2 addr=(rs1+imm)*num_thread_in_wg+thread_idx funct3字段编码与sw sh sb等访存指令一致	private memory

## 寄存器和 ABI

### 寄存器设置

架构寄存器数目: sGPR 64 个, vGPR 256 个, 元素宽度均为 32bit。64bit 数据使用 register pair。

物理寄存器数目: sGPR 256 个, vGPR 1024 个, 由硬件实现架构寄存器到物理寄存器的映射。

编译器提供 GPR 的使用量 (vGPR 和 sGPR 的实际使用数目, 是 4 的倍数), 硬件根据实际使用情况分配更多的 workgroup 同时调度。

从 RVV 视角下看, 向量寄存器宽度 vlen 固定为 num\_thread\*32bit, 硬件上相当于 vsetvli 指令的 SEW=32bit, ma, ta, LMUL=1。从 SIMT 编程视角看, 每个 thread 拥有至多 256 个宽度为 32bit 的 vGPR, 而 workgroup 拥有 64 个 sGPR。整个 workgroup 只需要做一次的操作,

如 kernel 和非 kernel 函数中的地址计算，就使用 sGPR；如果有分支的情况，则使用 vGPR，例如非 kernel 函数的参数传递。

一个warp(32thread)拥有的寄存器资源：vgpr0-255, sgpr0-127，每个格子代表32bit

	每个thread私有的一个float16 x变量						thread间有一个公共的float4 y变量	
	thread31	thread30	...	thread2	thread1	thread0		所有thread共有
v0	x.0	x.0		x.0	x.0	x.0	s0	
v1	x.1	x.1		x.1	x.1	x.1	s1	
v2	x.2	x.2		x.2	x.2	x.2	s2	
v3	x.3	x.3		x.3	x.3	x.3	s3	
v4	x.4	x.4		x.4	x.4	x.4	s4	
v5	x.5	x.5		x.5	x.5	x.5	s5	
v6	x.6	x.6		x.6	x.6	x.6	s6	
v7	x.7	x.7		x.7	x.7	x.7	s7	
v8	x.8	x.8		x.8	x.8	x.8	s8	
v9	x.9	x.9		x.9	x.9	x.9	s9	
v10	x.A	x.A		x.A	x.A	x.A	s10	
v11	x.B	x.B		x.B	x.B	x.B	s11	
v12	x.C	x.C		x.C	x.C	x.C	s12	
v13	x.D	x.D		x.D	x.D	x.D	s13	
v14	x.E	x.E		x.E	x.E	x.E	s14	
v15	x.F	x.F		x.F	x.F	x.F	s15	
v16							s16	

v0的[31:0]是thread0私有的，[63:32]是thread1私有的，...

所以OpenCL的向量类型只能使用分组寄存器表达

分组寄存器由编译器展开，以提供对 OpenCL 向量类型的支持。

分组寄存器在硬件上需要多周期发射，且寄存器依赖不便于检测，在编译器上完成这一步要容易很多。后续可以考虑针对标量 load/store 提供分组寄存器操作，标量访存指令有其特殊性：对连续地址访问的提升很大，如 GCN3 中的 S\_LOAD\_DWORDX8。

### 特殊寄存器

- x0: 0 寄存器
- x1: ra 返回 pc 寄存器
- x2: sp stack pointer - localmem baseaddr
- x4: tp privatemem baseaddr

### 栈空间说明

由于 OpenCL 不允许在 Kernel 中使用 malloc 等动态内存函数，也不存在堆，因此可以让栈空间向上增长。tp 用于各 thread 私有寄存器不足时压栈（即 vGPR spill stack slots），sp 用于公共数据压栈，以及在编程中显式声明了 \_\_local 标签



的数据（即 sGPR spill stack slots 和 localmem 的访问，实际上 sGPR spill stack slots 和 local data 都将作为 localmem 的一部分）。  
编译器提供 localmem 的数据整体使用量（按照 sGPR spill 1kB，结合 local 数据的大小，共同作为 localmem\_size），供硬件完成 workgroup 的分配。

### 参数传递 ABI

对于 kernel 函数，a0 是参数列表的基址指针，第一个 clSetKernelArg 设置的显存起始地址存入 a0 register，kernel 默认从该位置开始加载参数。  
对于非 kernel 函数，使用 v0-v31 和 stack pointer 传递参数，v0-v15 作为返回值。

### 自定义 CSR

注：在汇编器中可以使用小写后缀来表示对应的 CSR，例如用 tid 代替 CSR\_TID。

description	name	addr
该 warp 中 id 最小的 thread id，其值为 CSR_WID*CSR_NUMT，配合 vid.v 可计算其它 thread id。	CSR_TID	0x800
该 workgroup 中的 warp 总数	CSR_NUMW	0x801
一个 warp 中的 thread 总数	CSR_NUMT	0x802
该 workgroup 的 metadata buffer 的 baseaddr	CSR_KNL	0x803
该 SM 中本 warp 对应的 workgroup id	CSR_WGID	0x804
该 workgroup 中本 warp 对应的 warp id	CSR_WID	0x805
该 workgroup 分配的 local memory 的 baseaddr，同时也是该 warp 的 xgpr spill stack 基址	CSR_LDS	0x806
该 workgroup 分配的 private memory 的 baseaddr，同时是该 thread 的 vgpr spill stack 基址	CSR_PDS	0x807
该 workgroup 在 NDRange 中的 x id	CSR_GIDX	0x808
该 workgroup 在 NDRange 中的 y id	CSR_GIDY	0x809

该 workgroup 在 NDRange 中的 z id	CSR_GIDZ	0x80a
向 print buffer 打印时用于与 host 交互的 CSR	CSR_PRINT	0x80b

## 内存系统

承影 GPGPU 每个 SM 有单独的 L1 内存子系统，包括指令缓存、共享内存 (local memory/scratchpad memory)、数据缓存、常量缓存。其中指令 Cache 由前端取指级直接访问，其他存储器由指令通过 LSU 访问。所有 SM 访问同一块 L2 内存子系统。

### 地址空间

目前按照 32 位地址空间设计。privatemem 的访问需要使用专门的 vlw.v 指令，该指令会为每个 thread 自动计算其地址偏移量。在编译器视角，每个线程可用的 privatemem 空间是从 CSR\_PDS 开始连续的 0-1kB 空间，硬件会自动转换以便于 warp 的连续访存。

localmem 和 globalmem 的访问使用 vle32.v vloexi32.v vlw32.v 指令访问。二者使用地址字段进行区分，小于 local\_mem\_max 的地址均认为是访问 localmem 的。

localmem 使用实地址，访问片上 SRAM；globalmem、privatemem、print buffer 使用的地址由驱动实现分配和管理（物理上这三者都将映射到 ddr）。

### 一致性和连贯性特征

SM 间数据缓存的一致性 (memory coherence) 由程序员显式通过同步指令维护 (Consistency-Directed Coherence)，硬件不维护一致性。同步指令对应到 RISC-V A 扩展中的特定指令，硬件上 cache 提供 flush 和 invalidate 功能。

连贯性 (memory consistency) 在承影中可以理解为当前 SM 访存指令执行结果对其他 SM 可见的顺序是否和当前 SM 执行这些访存指令的顺序相同。在 RVWMO (RISC-V Weak Memory Ordering) 中，普通访存操作的连贯性要求由 PPO (Preserved Program Order) 规则 1 和规则 2 定义。对于常规读写操作，承影 L1 数据缓存的微架构表现如下表：

	【相 同 地 址】				【不 同 地 址】			
操作顺序类型	R-R	W-W	W-R	R-W	R-R	W-W	W-R	R-W
(对其他 SM)是否保序	是	是	是	是	是	否	否	是
RVWMO PPO 的保序要求	要求	要求	不要 求	要求	不要 求	不要 求	不要 求	不要 求

\*本表格中，A-B 表示访存操作 A 的程序顺序早于访存操作 B。

## L1 指令缓存

主要特征如下：

- （尚未固定）2 路组相联，缓存行大小为 blocksize（单位 Byte），总容量 64\*blocksize；
- 每次 LSU 访问的最大返回宽度为 4B；
- 每次 L1-L2 之间数据读写的最大宽度均为 128B（一个缓存行）；
- 替换策略支持 LRU、FIFO 替换策略；
- 支持对正在进行的请求进行无效化；
- （尚未固定）最多同时存在 128 个未完成 L2 访问请求；支持相同缓存行的 L2 访问请求合并，最大合并数 8。

## L1 数据缓存

主要特征如下：

- （尚未固定）2 路组相联，缓存行大小为 blocksize（单位 Byte），总容量 64\*blocksize；
- 一般的，blocksize = num\_thread；
- 虚拟地址索引、虚拟地址标记（VIVT）；
- 每次 LSU 访问的最大返回宽度为 blocksize；
- 每次 L1-L2 之间数据读写的最大宽度均为 blocksize（一个缓存行）；



- 写策略为写回-写不分配;
- 替换策略支持 LRU、FIFO 替换策略;
- 支持对整个数据高速缓存的无效和清除操作, 支持对单条缓存行的无效和清除操作;
- (尚未固定) 最多同时存在 128 个未完成 L2 访问请求; 支持相同缓存行的 L2 访问请求合并, 最大合并数 8。

数据缓存支持的指令类型包括:

- 基础指令集 I:
  - 访存指令: LOAD 和 STORE
  - 访存排序指令: FENCE
- 原子指令扩展 A:
  - 原子操作指令: AMO
  - 预留性读/条件写指令: LR/SC
- 向量指令扩展 V:
  - 向量读指令: VL, VLS, VLX
  - 向量写指令: VS, VSS, VSX
  - 自定义向量访存指令

### share memory

主要特征如下:

- (尚未固定) 2 路组相联, 缓存行大小为 blocksize (单位 Byte), 总容量  $64 * \text{blocksize}$ ;
- 一般的,  $\text{blocksize} = \text{num\_thread}$ ;
- 每次 LSU 访问的最大返回宽度为 blocksize;

share memory 支持的指令类型包括:

- 基础指令集 I:
  - 访存指令: LOAD 和 STORE

- 原子指令扩展 A:
  - 原子操作指令：AMO
- 向量指令扩展 V:
  - 向量读指令：VL, VLS, VLX
  - 向量写指令：VS, VSS, VSX
  - 自定义向量访存指令

## L2 缓存

主要特征如下：

- 多路组相连，组数和路数可通过 CacheParameters 配置，缓存行大小以及最小可写单位可通过 InclusiveCacheMicroParameters 配置，目前分别为 128Bytes 和 4Bytes
- 物理地址索引，物理地址标记 (PIPT)
- 写策略为写回-写分配
- 替换策略支持随机、伪 LRU
- 支持对整个数据高速缓存的无效和清除操作，支持对单条缓存行的无效和清除操作
- 具备 MSHR，支持非阻塞访存
- 除普通 load、store 操作外，支持 AMO 原子操作，通过 Tilelink 协议支持 LR、SC 操作
- 支持 stride 预取
- L2cache 通过一个 AXI Adapter 发送符合 AXI 4 协议的数据给主存

## 处理器模式

仅支持机器模式，程序启动前即会配置好 CSR。内存管理由驱动完成。  
暂不支持异常和中断处理。

## 总线接口

承影包含一个 AXI4 主设备接口和一个 AXI4-Lite 从设备接口。L2Cache 通过 AXI4 主设备接口访问片外存储器。host 通过 AXI4-Lite 接口，对 workgroup 进行配置。

## 互连网络

若干个 SM 组成一组 SM\_cluster，SM\_cluster 与 L2Cache bank 间通过 crossbar 连接。

## 微架构（硬件视角）

### 任务分配和汇编相关

#### CTA 任务分配

在硬件层面，将按照 32 个 thread 组成一个 warp 的形式，作为整体在 SM 硬件上进行调度。同一个 block 的 warp 只能在同一个 SM 上运行，但是同一 SM 可以容纳来自不同 block 甚至不同 grid 的若干个 warp。

CPU 发送给 GPU 的任务以 workgroup 为基本单位，由 CTA scheduler 接收，CTA scheduler 会按 block 中包含的总 warp 数信息，以及需要占用的 local memory、sharedmemory 大小，将 block 对应分配到空闲（即剩余资源足够）的 SM 上。

CTA scheduler 以 warp 为单位逐个发送给 SM，同一 workgroup 的 warp 会分配到同一 SM 中，warp\_slot\_id 的低位即表明了该 warp 在当前 workgroup 中的 id，高位表明了 workgroup 本身所属的 id。相应的，SM 会通过此 id，计算出当前 warp 在所属 workgroup 中的位置，并将该值置于 CSR 寄存器中，供软件使用。分配的 localmem baseaddr 需要通过 CSR 读取，而 privatemem baseaddr 和 register base 则由硬件隐式映射。由于一个 warp 只有一套 CSR，thread\_id 需要用 vid.v+CSR\_TID 计算出。



## 汇编编程说明

1. `get_global_id()`通过 `vid.v + csrr tid` 三条指令实现。
2. 输入参数和访存地址需要按照预设参数传递方式，从 CSR 读取使用
3. 自定义指令的使用：
  1. `predicate`: 我们在支持 `rvv` 定义的软件控制 `mask` 的同时，也支持用自定义指令来启动隐式的硬件 `predicate`，详见自定义指令一节。
  2. `warp` (即 `kernel`) 运行结束时需要显式使用 `endprg` 指令。
  3. 同一 `block` 内 `thread` 同步，使用 `barrier` 指令。

其余行为与 `rvv` 编程一致。

对于超过单组硬件处理能力长度的向量数据，支持使用 `rvv` 中定义的 `stripmining` 方式执行，默认单次处理 `num_thread` 个数据。与向量处理器不同的是，可以用软件 `mask` 实现，也可以用 `SIMT-stack` 实现，也可以在 `block` 大小允许时拆分为更多 `warp` 去调度。

`SIMT-stack` 补充：目前从软件视角看单 `warp` 执行，功能与 `rvv mask` 完全一致。优势在于 1)所有 `thread` 方向一致时，可跳过 `if` 分支或 `else` 分支；2)减少对寄存器堆(`v0`)访问次数，减少获取操作数时的 `bankconflict`；3)实现快速嵌套分支，目前硬件支持最坏情况下的 `num_thread-1` 层嵌套；4)为后续硬件实现 `multi-path IPDOM`、独立线程调度、`warp` 合并等提供便利。

## 指令集架构

### RVV 与 GPGPU 的结合

《量化研究方法》中提到了向量处理单元与多线程 GPU 在 `SIMD` 层面上的工作形式十分相似，向量处理器的车道与多线程 `SIMD` 的线程是相似的。区别在于通常 GPU 的硬件单元更多，`chime` (钟鸣) 更短，向量处理器通过深度流水线化的访问来隐藏延迟，GPU 则是通过同时多 `warp` 切换来隐藏延迟。因此在向量层面的操作上，`RVV` 足以覆盖住 `GPGPU` 中的操作。此外，形如 `AMD` 和 `turing` 后的 `NV`，提供了标量 `ALU`，也是借鉴了向量处理器的方式。

因此，在 RVV 的基础上添加自定义的分支控制指令（实际上沿用 RVV 本身的 mask 也能实现）、线程同步指令、线程控制指令，就能实现 GPGPU 的功能。为了最大限度的保留对 RVV 开源工具链的兼容性，我们对 RVV 中的大部分指令都进行了支持。少数不支持的指令包括：1. 涉及线程间数据交换的 shuffle 等指令，在 GPGPU 中线程间通常是独立操作，数据依赖需要用 atomic 或 barrier 显式操作 2. 向量寄存器长度和宽度变化的指令，GPGPU 中几乎不会触及（少有的几条向量或者量化相关的功能会需要类似的功能） 3. 64bit 相关的指令，在后续版本将支持。

在 RVV 的 stripmining 基础上添加 warp 级别并行，或许能在更优尺度上裁剪向量/SIMT 指令，探索划分和调度空间。

### 寄存器设计

单个 SM 上能同时承载的最大 warp 数为 num\_warp，每个 warp 由 num\_thread 个线程组成。每个 warp 都有一套自己的寄存器，每个标量寄存器的宽度为 32bit，每个向量寄存器的宽度为 32\*num\_thread，并归属于各个 thread 私有。

物理寄存器堆采用统一方式，根据各 warp 的实际使用情况动态分配。

虽然所用指令形式和意义相同，但区别于 rvv，我们目前实现的 GPGPU 中并不支持向量寄存器的长度和数量变化（长度固定为 32bit，数量固定为 num\_thread），因此对于 vsetl 系列指令，只有返回的剩余元素数量是有效的。

### 地址映射

用地址范围来区分 localmem 和 globalmem，其中 localmem 由 CTA-scheduler 管理，globalmem、privatemem、constantmem 由驱动分配和管理。

GPU 中的物理地址空间包括 localmem 和 globalmem，早期版本 cuda 和 OpenCL 编程中需要显式声明地址属性，在 PTX 中每个地址都带有属性声明了其类型。目前承影尚未实现 MMU，所有 SM 共用一个 4GB 的全局地址空间，其中地址为 0-128kB 的字段将被映射到 SM 内部各自的 sharedmem 上，从 global\_baseaddr 到 4G 的空间则映射到同一块

ddr 的相同地址上（即该部分使用物理地址）。coherence 在 GPU 中是由软件显式管理的，通过 flush 和 invalidate 实现，对应 fence 指令。

## 指令集范围

目前支持 RVV 中的指令包括以下类型：

1. 计算类，包括整数运算（加减、比较、移位、位运算、乘、乘加、除）、单精度浮点运算（加减、乘、乘加、除、整型转换、比较），支持带 mask 执行
2. 访存类，包括三种访存模式，以及 byte 级读写，支持带 mask 执行
3. mask 类，包括比较及逻辑运算，但不包括 vmsbf 等涉及 thread 间通信的指令，也不支持 gather 等操作

对于改变向量位宽的指令暂不支持，但 vsetl 系列指令可以返回 vl 供 stripmining 使用。RV32I 中支持除 ecall ebreak 外的指令，M F 中支持 32bit 相关指令。

目前支持的自定义指令包括：

1. predicate：在支持 rvv 定义的软件控制 mask 的同时，也支持用自定义指令来启动隐式的硬件 predicate，由 vbeq 等启动的隐式硬件 predicate 将默认对后续指令生效，直到触发新的 vbeq 或 join 时才会改写。启动的方式是使用自定义的 vbeq 系列线程分支指令，该指令会启动一个 split，计算出当前分支的 mask 情形，并将 else 对应的 mask 压入 SIMT-stack，然后带有 mask 执行 if 段，待 if 段末尾遇到 join 后再将 else 段及其对应 mask 出栈，待 if 段执行完成后合并恢复 mask，分支结束。该过程可以嵌套，压栈的最坏情况深度等同于单 warp 中的线程数 32（如果每次总选择最多的方向去压栈而非默认压 if，那么栈深度只需要 5 即可）。此外，如果分支计算出全走其中一条，将不会压栈并跳过另一条分支，带有 branch divergence 的 for 循环也可以用此机制实现。
2. barrier：用该指令可以实现 warp 间的同步，每个遇到该 barrier 指令的 warp 都会等待，直到该 workgroup 中所有未结束的 warp 都遇到此指令才会再次继续。
3. endprg：warp 运行结束时需要显式使用 endprg 指令。
4. regext{i}：为后一条指令扩展使用的寄存器数目和立即数宽度。
5. 其它为提高代码效率的自定义指令，可参见“kernel 内汇编指令说明 - 指令集范围 - 自定义指令”一节



## 驱动接口

遵循 CTA 调度器提供的接口信息支持，配置好使用的寄存器数目、localmem baseaddr、warp id 等。

## 微架构设计

总体分为前端和后端，前端包括了取指、译码、指令缓冲、寄存器堆、发射、记分板、warp 调度，后端包括了 ALU、vALU、vFPU、LSU、SFU、CSR、SIMT-stack、warp 控制等。涉及寄存器连接的模块间均采用握手机制传递信号。

### warp 调度

warp 调度器主要的功能包括：

1. 接收 CTA 调度器提供的 warp 的信息，分配 warp 所属的硬件单元并预设 CSR 寄存器值，激活该 warp 并标记所属的 block 信息。在 warp 执行完成后，将该信息返回给 CTA 调度器并释放对应硬件。
2. 接收流水线发送的 barrier 指令信息，将指定的 warp 锁住直到其所属的 block 的所有活跃 warp 到达此 barrier。
3. 选择发给 icache 的 warp，通常采用贪婪的策略，但当 icache 发生 miss，或 ibuffer 已满时会将该 warp 的 pc 回退并切换到下一个 warp 发射。
4. 选择发给执行单元的 warp，通常采用轮询的策略，选择出当前指令缓冲有效、记分板未显示冲突、执行单元空闲的指令。切换 warp 仅需要一个周期。

### 取指

每个 warp 存储各自的 pc，被选中送入 icache 的 pc 会+4，其余保留原值，遇到跳转时则替换为目标地址。

### 译码

icache 命中的指令进行译码，译码器根据指令内容转换出对应的控制信号，并送入对应的指令缓冲中。

## 指令缓冲

指令缓冲是一系列的 FIFO，每个 warp 有各自的 ibuffer，接收译码后的输入并等待选中发射。

## 操作数收集器

操作数收集器会接收来自指令缓冲中的请求，依据所需数据类型，经由 crossbar 向寄存器堆访问获取数据，获取后即进入待发射状态。

寄存器堆采用 unified 方案设计：硬件上单个 SM 里共有 1024 个向量寄存器和 512 个标量寄存器，单个 warp 可以使用多达 256 个向量寄存器和 64 个标量寄存器，由硬件完成物理寄存器的分配。编译器指明使用的寄存器数目（4 的倍数），单个 warp 使用寄存器数目越少，就能在硬件上分配更多的 warp 同时运行。

寄存器堆硬件上分了 4-bank，每个 bank 一读一写 port。

目前版本标量和向量寄存器不支持同时访问，后续会修改。

## 发射

发射仲裁由 warp 调度器进行，被选中的控制信号与源操作数一起，依据识别其所需运算单元的类型，发送到对应运算单元执行。

## 记分板

每个 warp 有各自的记分板，当一条指令发射成功后，所写入的寄存器将被记分板标记，下一条指令若会读写已被标记的寄存器，则不允许发射，直到指令执行完成记分板释放对应寄存器。

分支、线程分歧、跳转、barrier 指令也会被锁住，只有这些指令完成判断且能继续顺序执行时才会释放。执行 barrier 期间，会同时清空除 ibuffer 外的流水线

（属于该 warp 的部分）；若发生跳转，会同时清空所有流水线（属于该 warp 的部分），此后记分板解除锁定。

## 写回

各个运算单元在输出级均有 FIFO，等待写回寄存器的结果暂存在其中，由 Arbiter 选中后写回寄存器。写回标量寄存器与向量寄存器是完全独立的数据通路。

## ALU

ALU 中进行标量运算，包括 warp 间共用的数据，以及跳转控制等。

## vALU

是单个 ALU 的复制，是核心的整型运算单元，供 warp 的多线程车道进行运算。典型运算消耗 1cycle。支持折叠为 num\_lane 个。

## vFPU

是核心的浮点运算单元，供 warp 的多线程车道进行运算，标量浮点运算也在此进行。全流水设计，乘法和乘加有复用数据通路。典型乘加消耗 5cycle，乘法消耗 3cycle，并支持折叠为 num\_lane 个。

## MUL

乘法运算单元，供 warp 多线程车道进行运算。使用 wallace tree 结构，全流水设计，典型乘法消耗 2cycle，并支持折叠为 num\_lane 个。

## LSU

是核心的访存单元，会根据地址范围判断将读写请求发送给 sharedmem 或 dcache（再由 L1 dcache 访问 L2cache，再访问 ddr 即 globalmemory）。LSU 中有 MSHR 形式的结构，可以一次存储和记录多个 LSU 请求，也会收集 dcache 和 sharedmem 返回的数据，集齐后再返回给流水线。LSU 中会完成 strided 及任意模式下向量地址的计算，并根据地址范围以 cacheline 为单位进行合并访问。最理想的情况（指地址连续且对齐 cacheline）一次访存即可取出单个 warp 所需结果。bank conflict 由 sharedmem 和 dcache 自行处理。在 LSU 中还会记录现有的访存请求信息，以实现 fence 指令。当遇到该指令后，会让所属 warp 的所有访存请求处理完成（读数据返回数据，写数据返回写响应）后再发送新的访存请求。



## CSR

在 warp 启动时，对应的 CSR 会设置好应用程序所需的一些值，包括 thread id 等。vsetl 也在此计算。其余与 riscv cpu CSR 的功能一致。

## SFU

区别于 PTX 中提供的 sin cos 等函数，目前的 SFU 只支持了 rv 定义的整数除法取余、浮点除法、浮点平方根功能。

本身运算就需要多周期完成，加上 SFU 中的运算单元数量少于 lane 数，因此如果未 mask 的线程较多时，chime 会更长。

## TC

tensorcore，全流水，支持特定格式下的张量计算。

## warp 控制

barrier、endprg 会发送给 warp 调度器处理。  
后续会考虑支持动态并行，添加 warpgen 指令。

## SIMT-stack

SIMT stack 的主要功能如下：维护分支嵌套控制流，保障程序运行的正确性；在实际没有分支分歧发生时，跳过不必要程序段的执行。

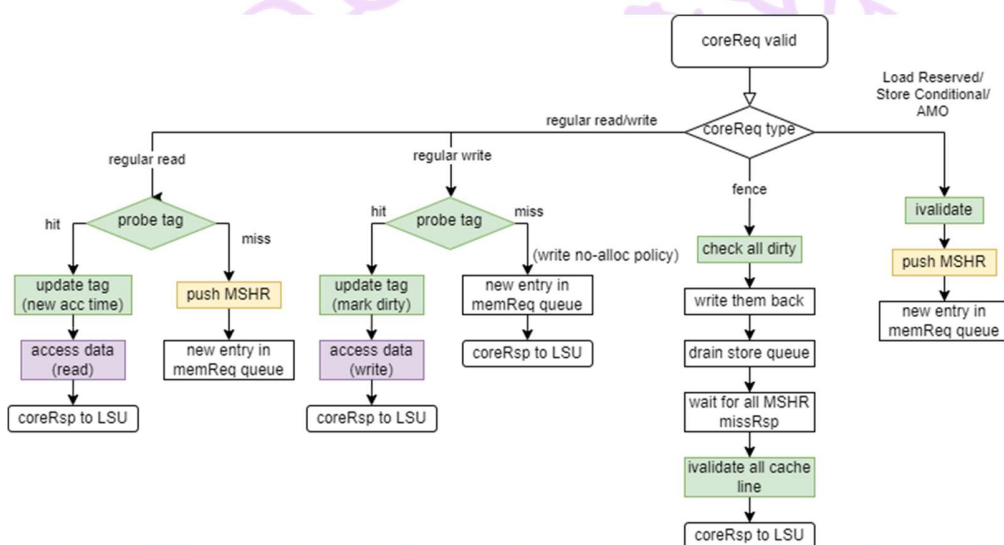
由 SIMT-stack 设置的隐式 mask 会在该 warp 执行过程中一直生效，直到有其它分支管理支持对其进行修改。该 mask 与 rvv 软件形式的 mask 可以叠加生效。

与分支管理相关的自定义扩展指令集有上表所示 7 条，其中 1-6 条为分支指令。

以 vbeq 指令为例，需要完成的功能为：取源操作数 vs2 与 vs1，valu 模块对这两个向量寄存器中的元素一一进行比较，对于第 i 个元素，若  $vs2[i]=vs1[i]$ ，则计算结果 out[i] 为 1，最终 valu 的输出结果 out 为分支指令对应的 else 路径掩码，同时译码模块将向分支管理模块发送分支发生标记以及 else 路径 PC 起始值 PC branch。

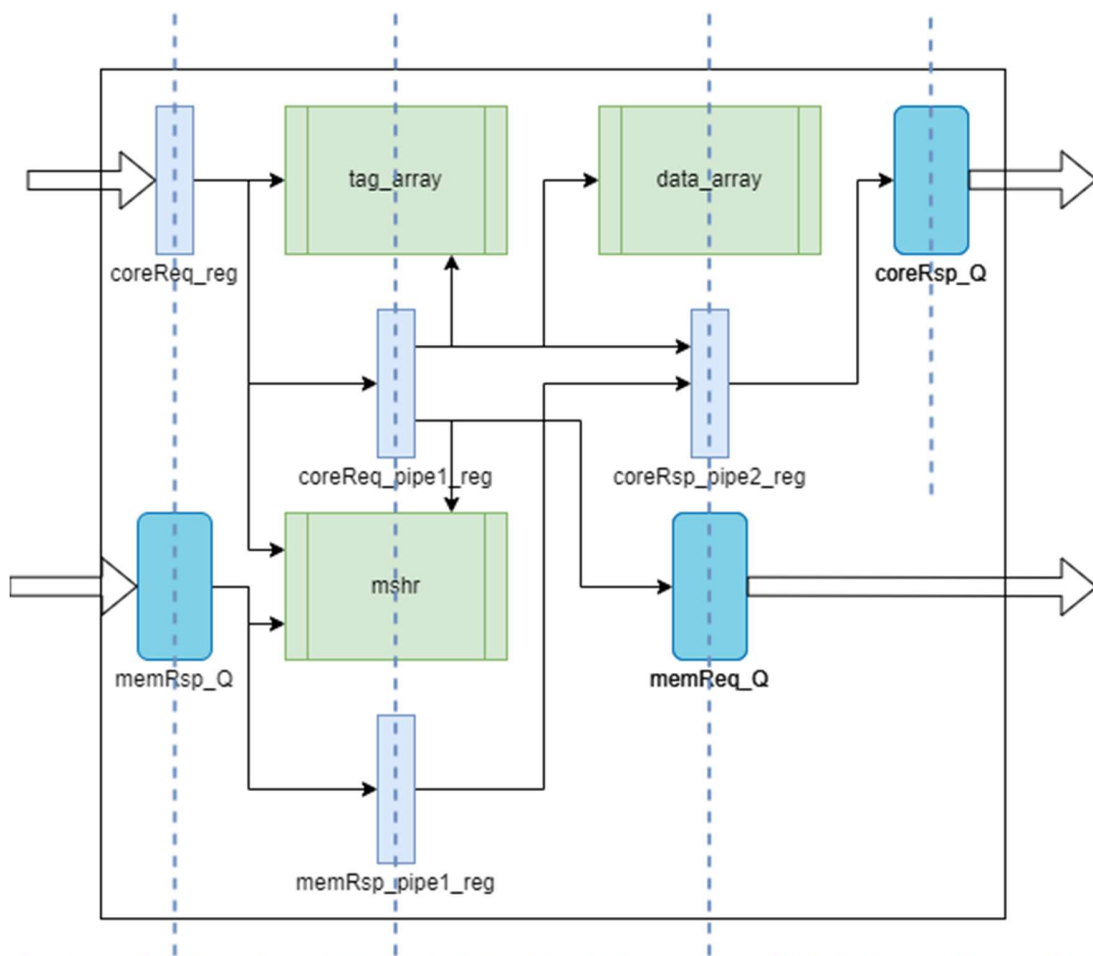
## DCache

整个 cache 分为对 LSU 接入的两组接口 coreReq 和 coreRsp，以及对更高层 cache (L2) 接入的两组接口 memReq 和 memRsp。其中 coreRsp 直接接入 SM 核心流水线写回级。coreReq 支持的请求类型包括：普通读请求（标量/向量）、普通写请求（标量/向量）、fence 请求、预留性读请求、条件写请求、原子操作请求。这些请求在 cache 中经历的处理操作如下图：



这些请求操作大多数与指令——对应。AMO 指令是一个例外。AMO 指令允许带有.aq 和/或.rl 标识符 (acquire 和 release)，未来 RV 标准指令集/扩展指令集也可能会添加普通访存指令携带该标识符的支持。对于携带此类标识符的访存指令，LSU 会将其分割为一个不携带标识符的访存请求和一条对应的 fence 请求。

整个 cache 由三个流水级构成，示意图如下：



对 coreReq 路径来说:

S0: 发出 tag 和 mshr SRAM 访问请求。

S1: 根据请求类型和 tag 返回值决定后续操作方式: 压栈 memReq\_Q 或/与 发出 data 访问请求。

S2: 压栈 coreRsp\_Q。

对 memRsp 路径来说:

S0: 发出 mshr SRAM 访问请求。

S1: 根据请求类型和 mshr 返回值, 更新 tag、组织 coreRsp。

S2: 为了避免流水线冲突引入的冗余流水级。



整体微架构方案如图所示。

