

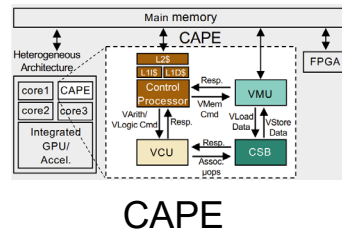
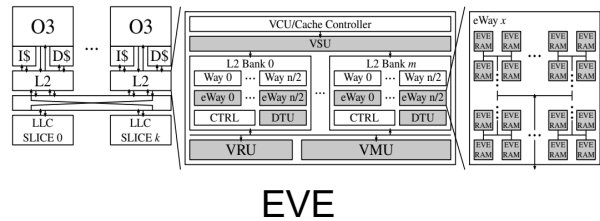
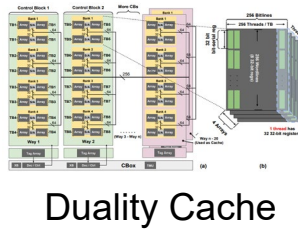
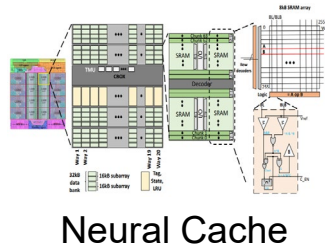
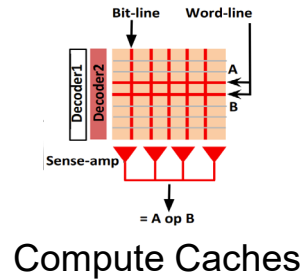


SUPPORTING A VIRTUAL VECTOR INSTRUCTION SET ON A COMMERCIAL COMPUTE-IN-SRAM ACCELERATOR

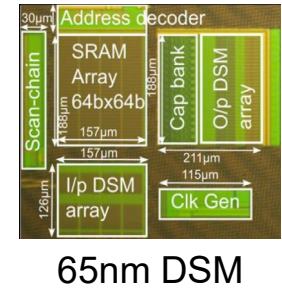
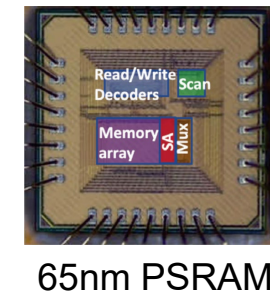
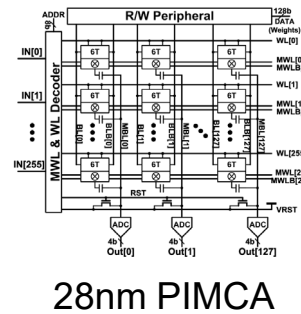
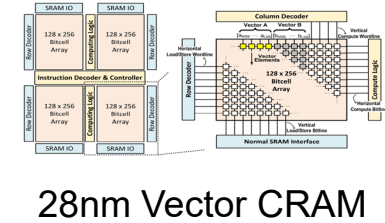
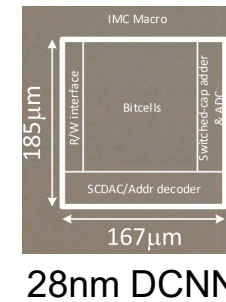
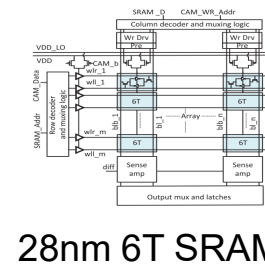
Courtney Golden^{1,3}, Dan Ilan², Caroline Huang¹, Niansong Zhang¹, Zhiru Zhang¹, and Christopher Batten¹

¹Cornell University, ²GSI Technology, ³Massachusetts Institute of Technology

Simulation



Test Chips

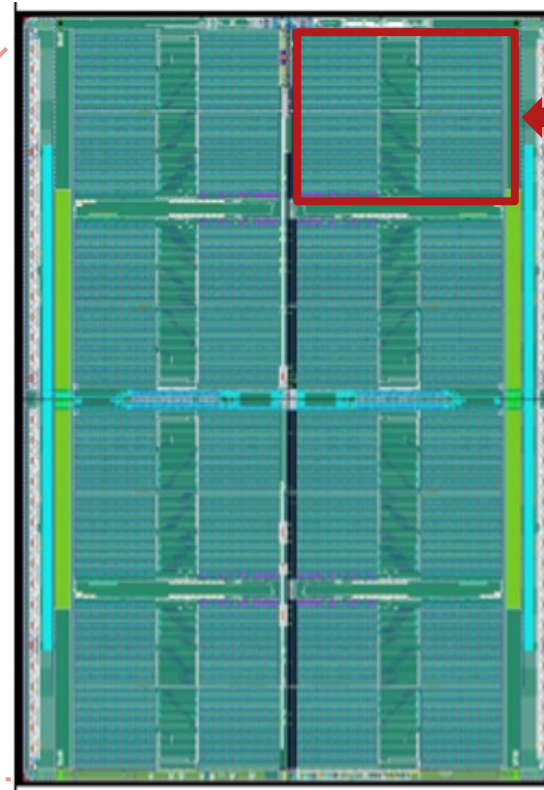
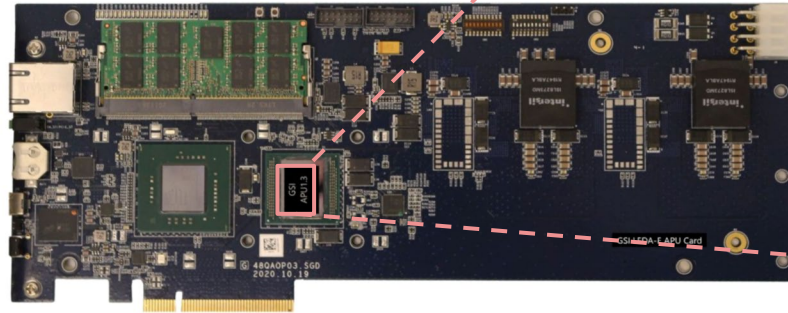


microcode

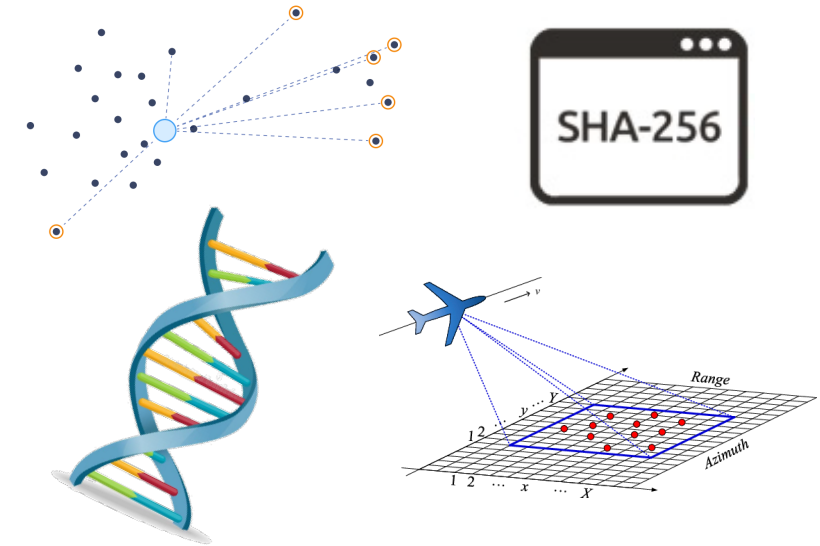
```
0xFFFF: RL = VRF[vsrc0];  
0xFFFF: RL |= VRF[vsrc1];  
0xFFFF: VRF[vdst] = RL;
```

higher-level abstraction

```
gvm1_add_u16(vr1, vr2, vr3);  
gvm1_conv_u32_gf16_8k(vr0, vm0, 4, 3, ..., vr7);  
gvm1_blk_sha1_folded(vm0, vm1, vm2, vm3, vr0);
```



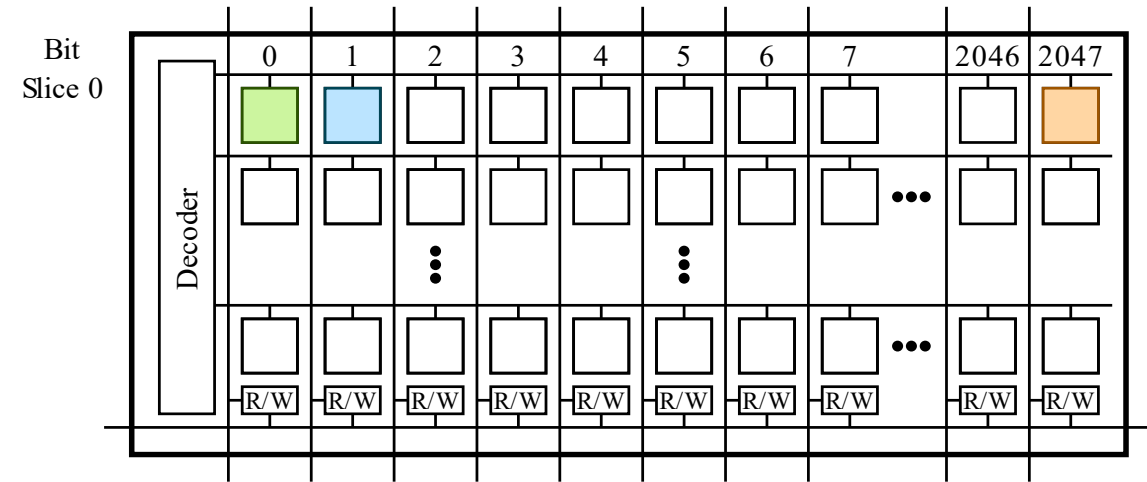
compute-enabled
SRAM array



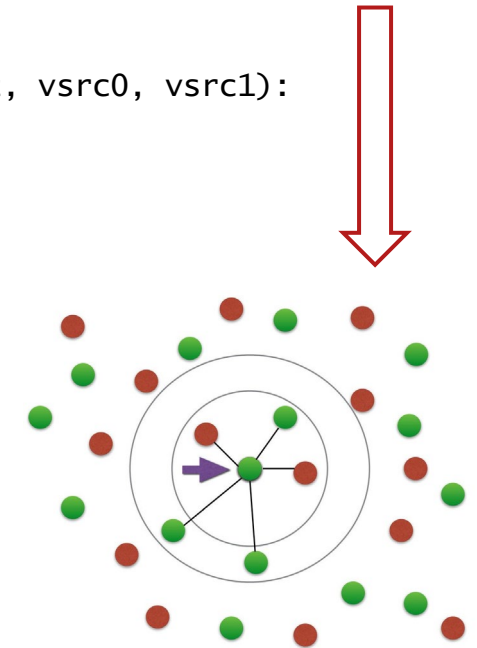
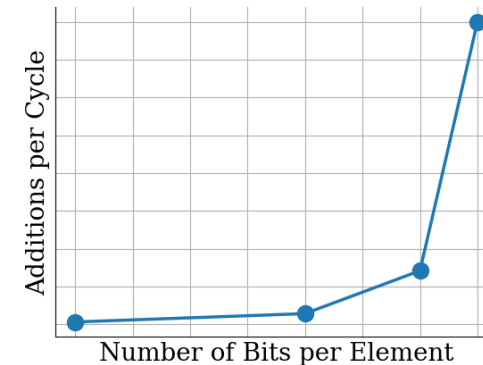
Our Goal: Explore the potential for supporting a virtual vector instruction set on a commercial compute-in-SRAM architecture to generalize beyond specific application domains

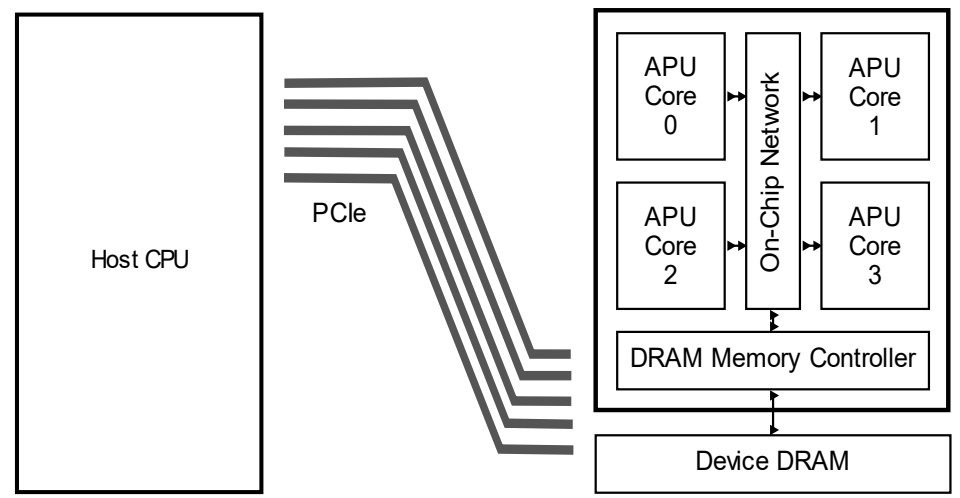
Supporting a Virtual Vector Instruction Set on a Commercial Compute-in-SRAM Accelerator

- Motivation
- **APU Microarchitecture**
- APU Microcoding
- Virtual Vector Instruction Set
- Microbenchmarking Results

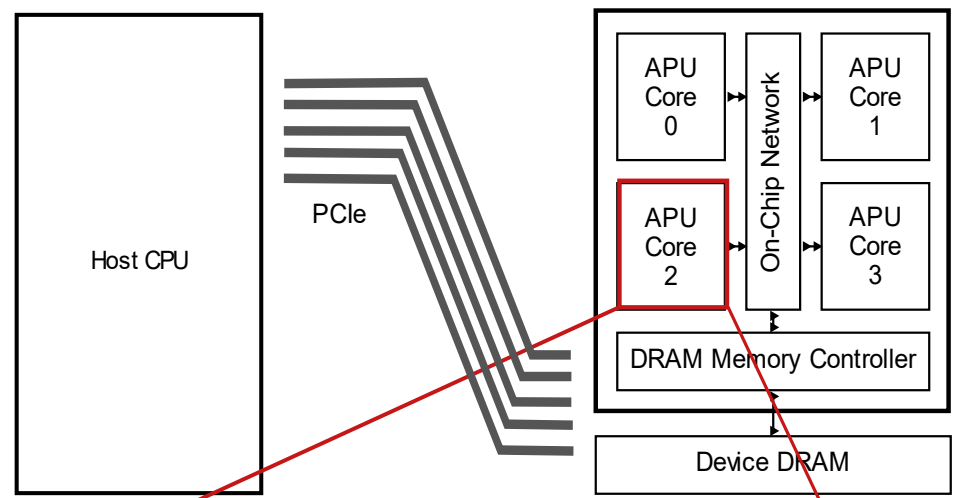


```
APL_FRAG _frag_bitwise_or(vdst, vsrc0, vsrc1):  
0xFFFF: RL = VRF[vsrc0];  
0xFFFF: RL |= VRF[vsrc1];  
0xFFFF: VRF[vdst] = RL;
```

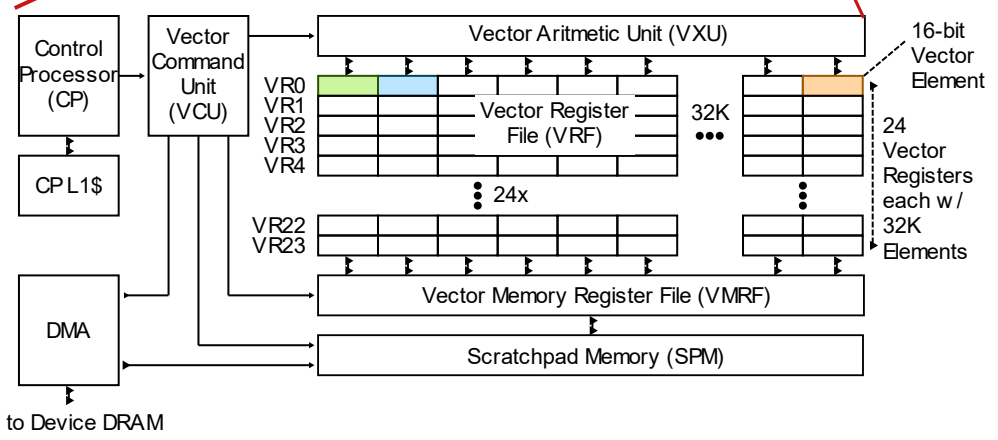




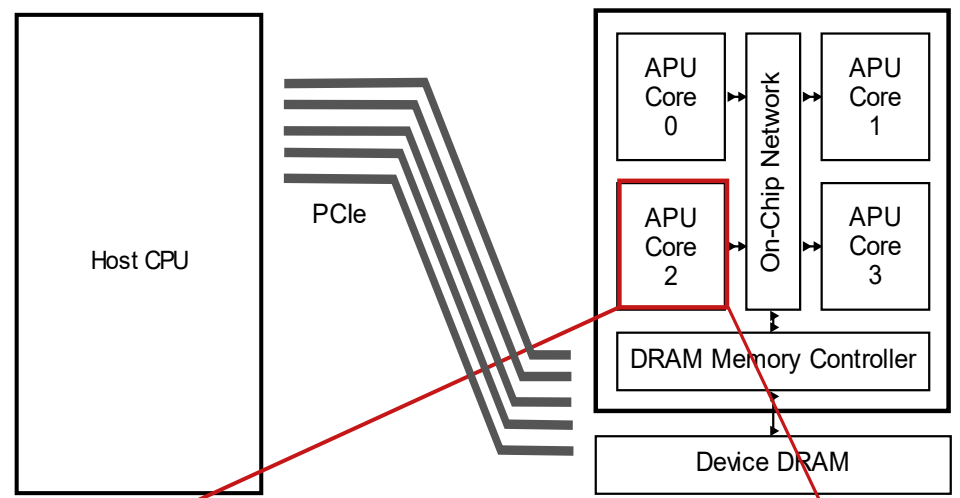
System Overview



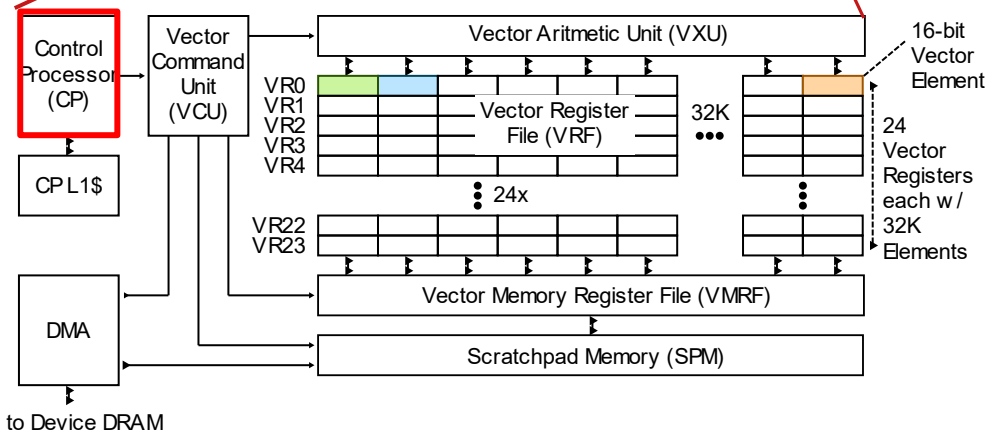
System Overview



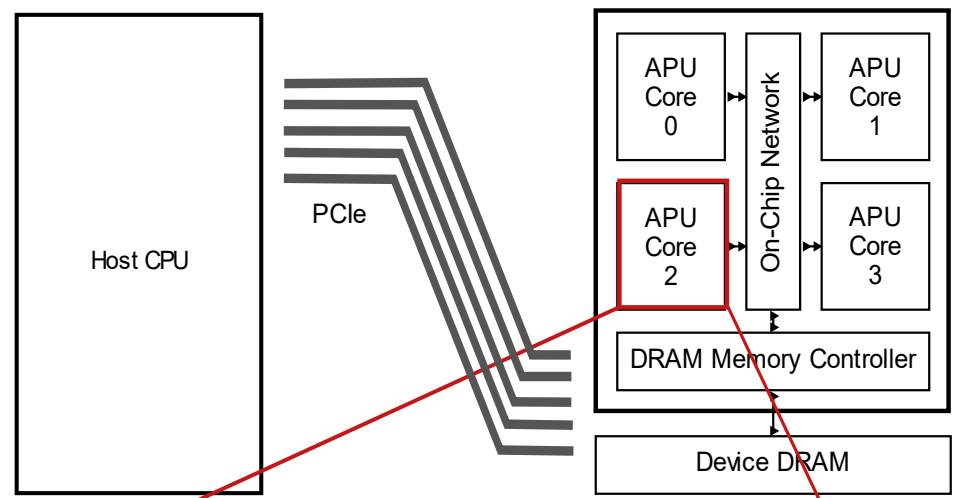
APU Core Logical View



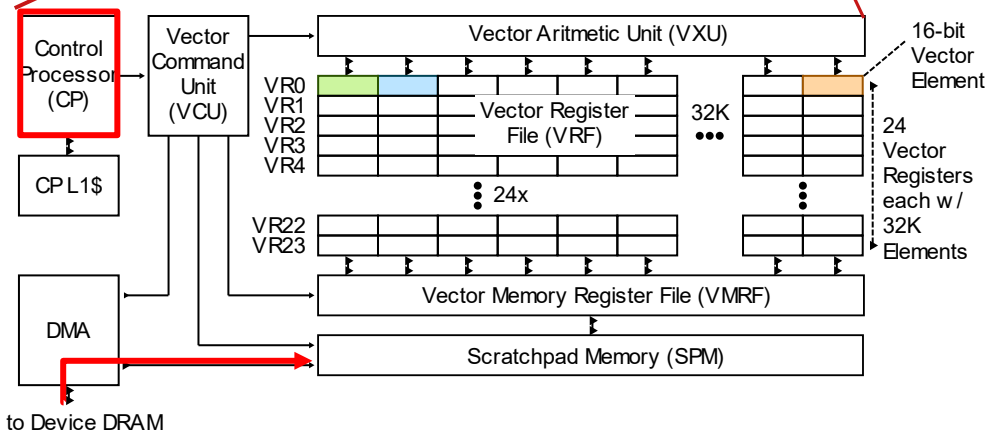
System Overview



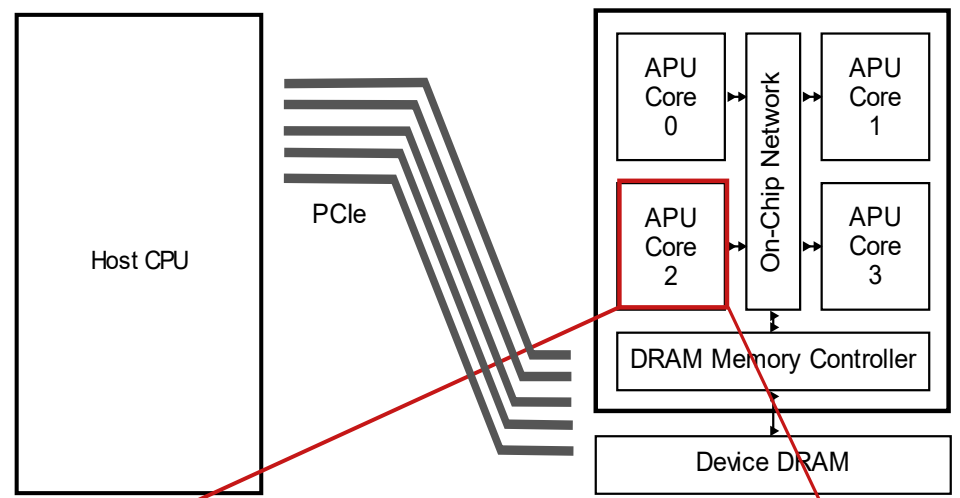
APU Core Logical View



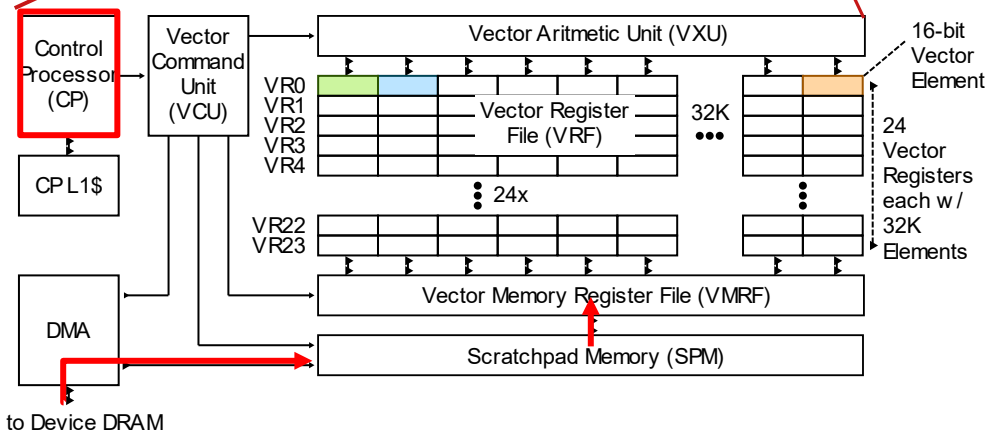
System Overview



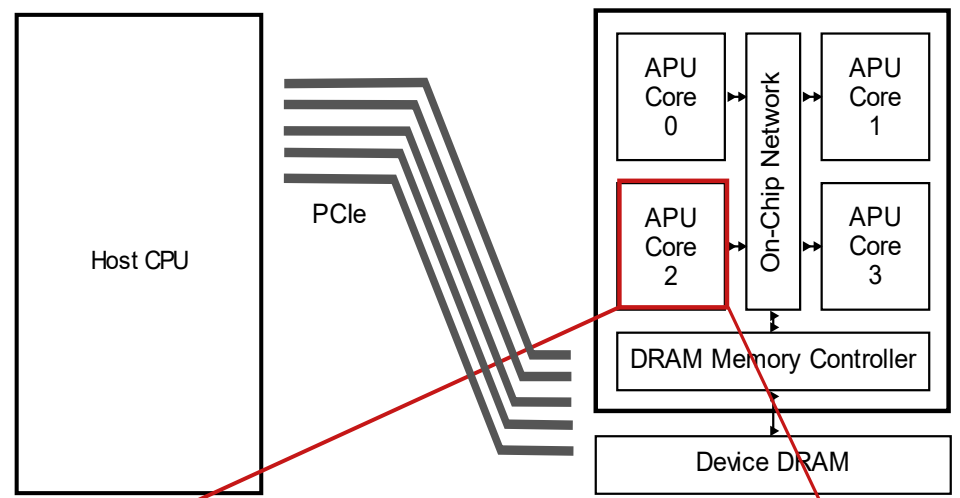
APU Core Logical View



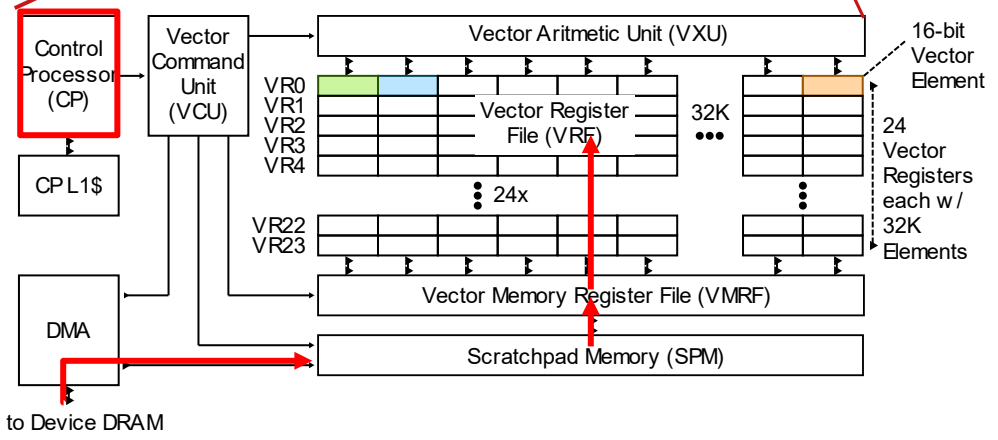
System Overview



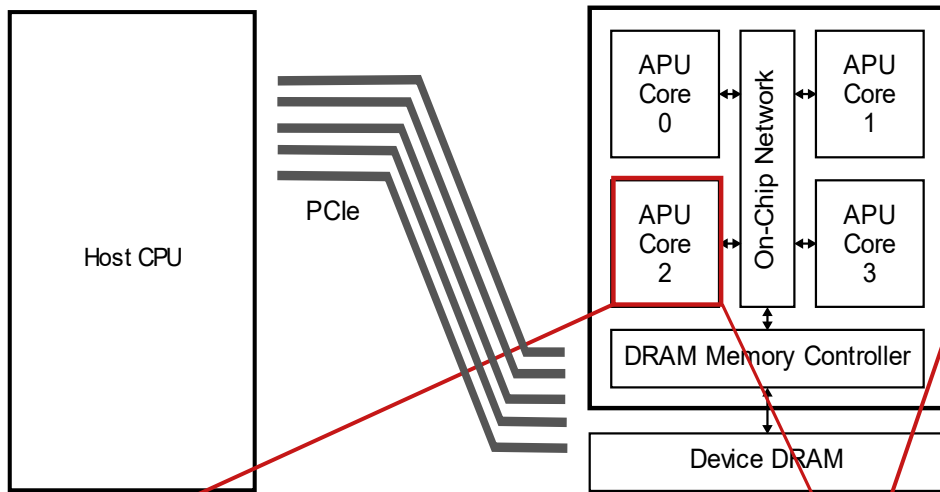
APU Core Logical View



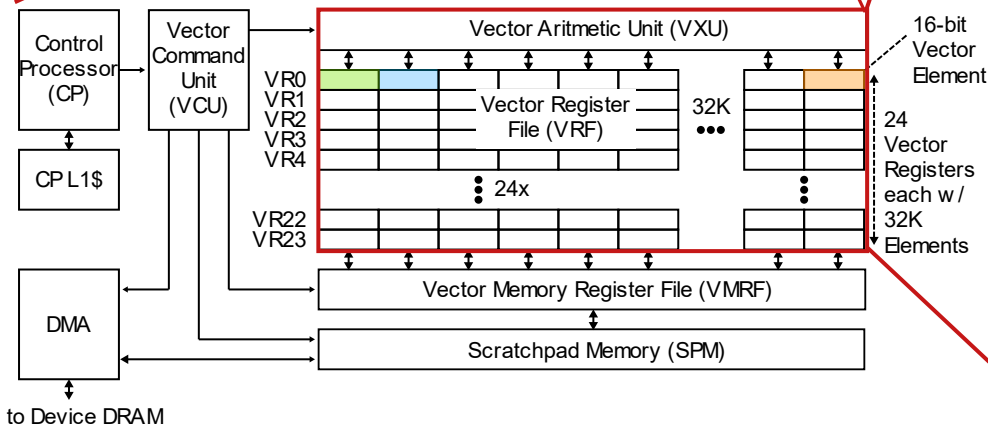
System Overview



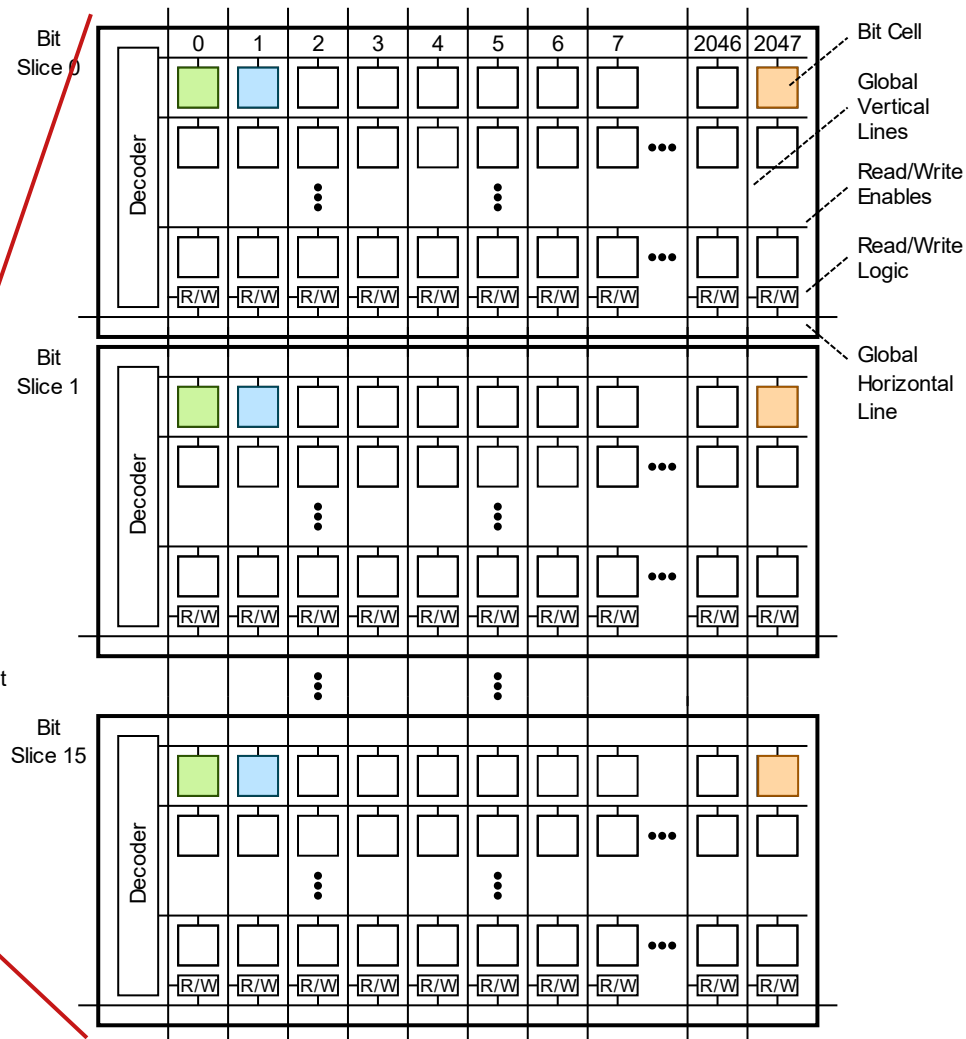
APU Core Logical View



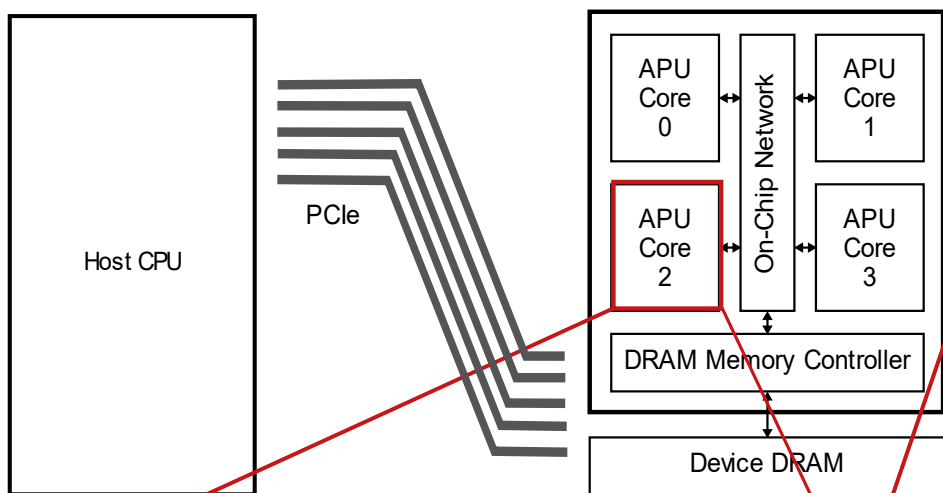
System Overview



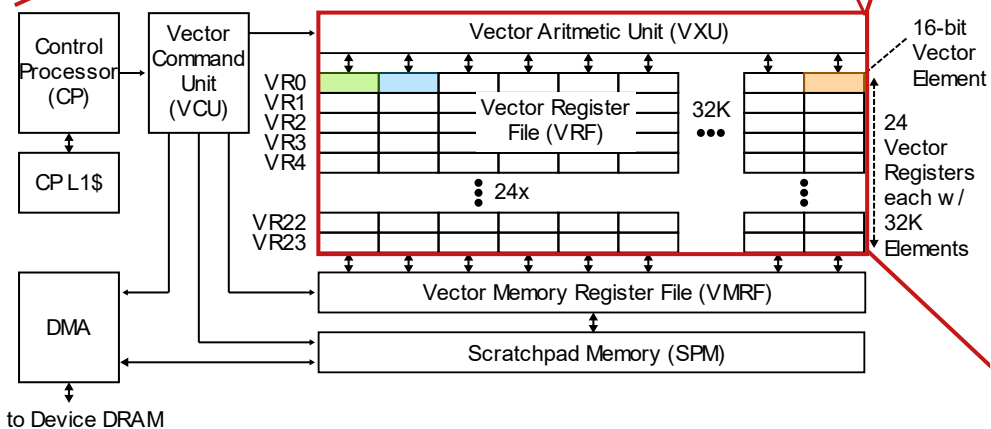
APU Core Logical View



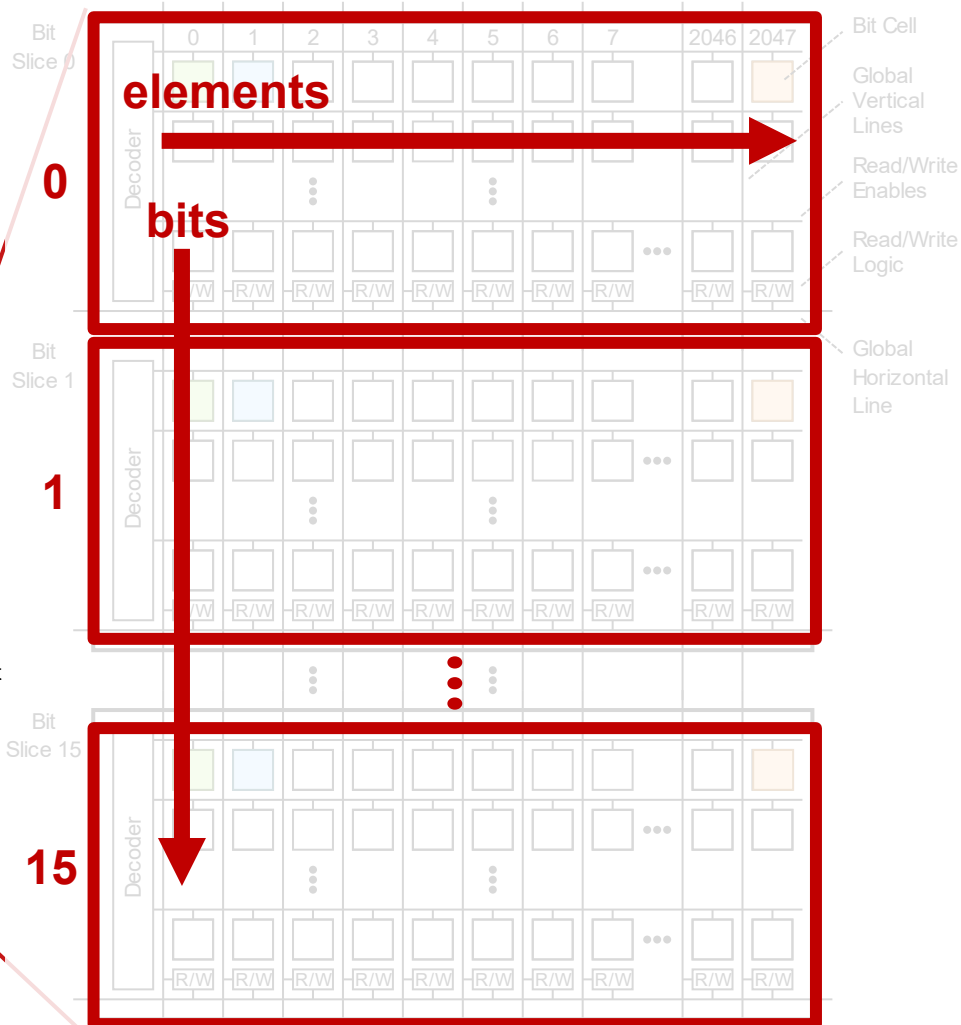
SRAM Bank Physical View



System Overview



APU Core Logical View



elements

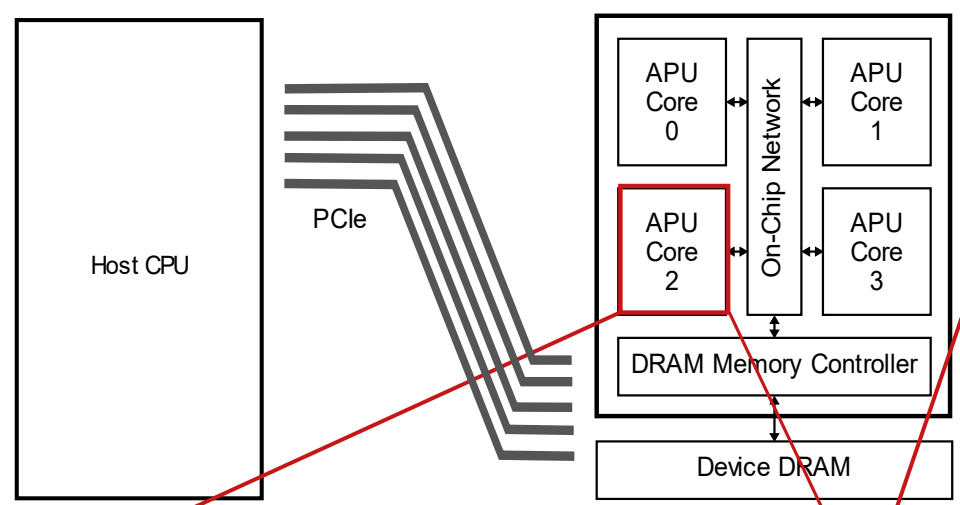
bits

0

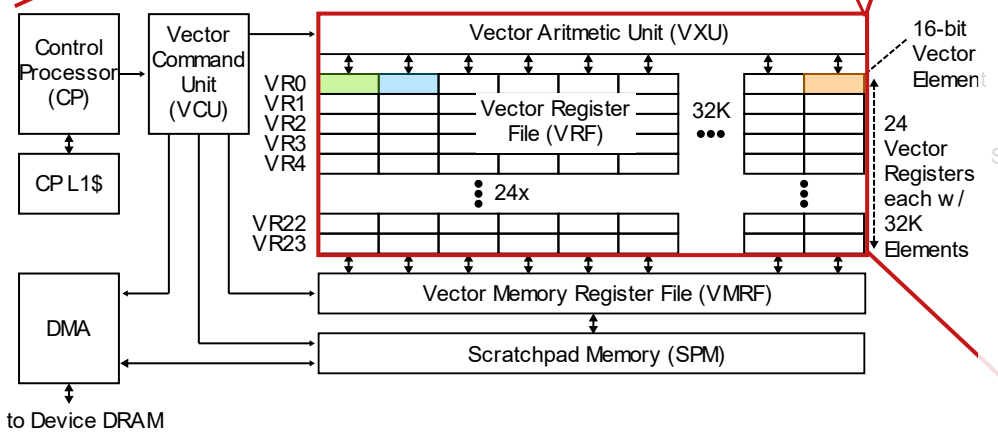
1

15

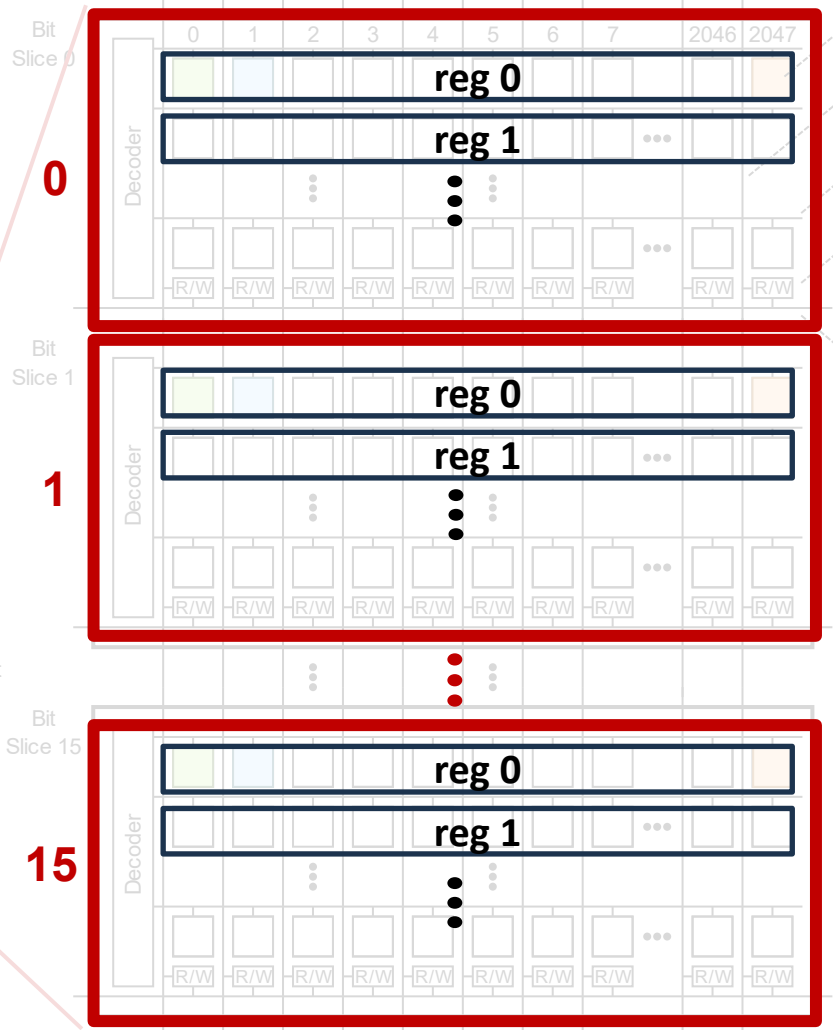
Bit Cell
Global Vertical Lines
Read/Write Enables
Read/Write Logic
Global Horizontal Line



System Overview



APU Core Logical View

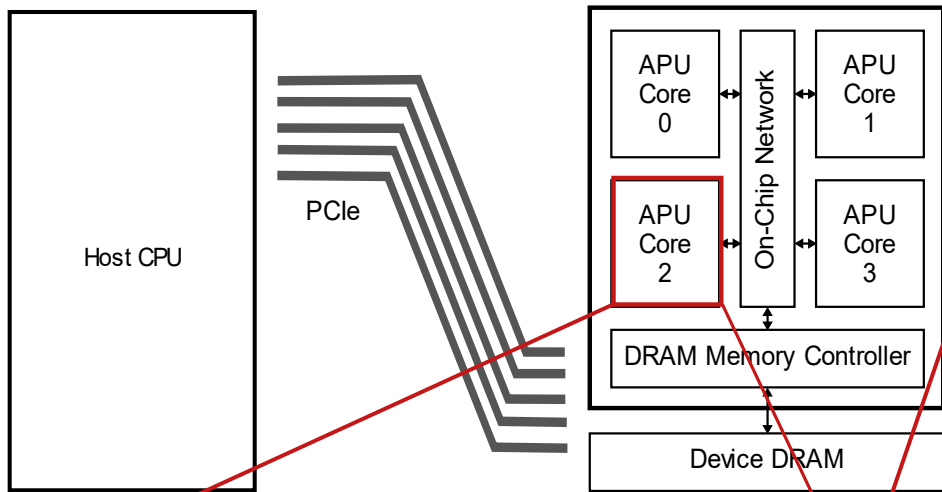


SRAM Bank Physical View

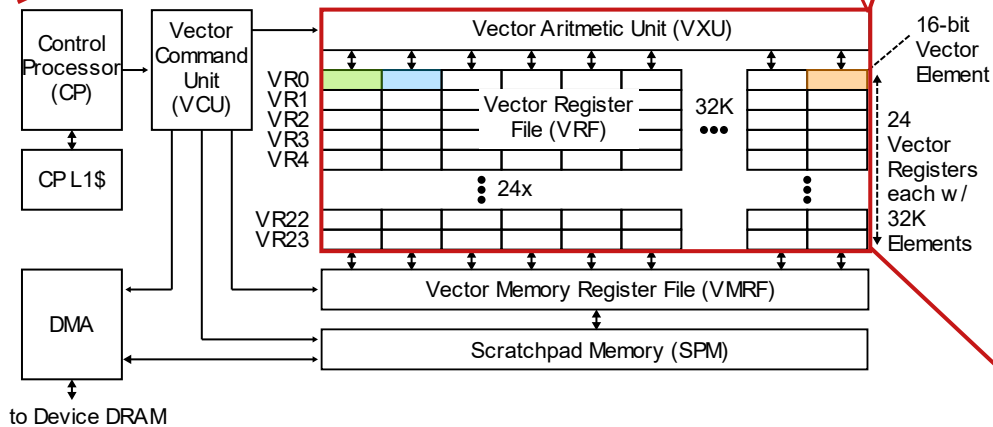
Bit Cell
Global Vertical Lines
Read/Write Enables
Read/Write Logic

registers interleaved

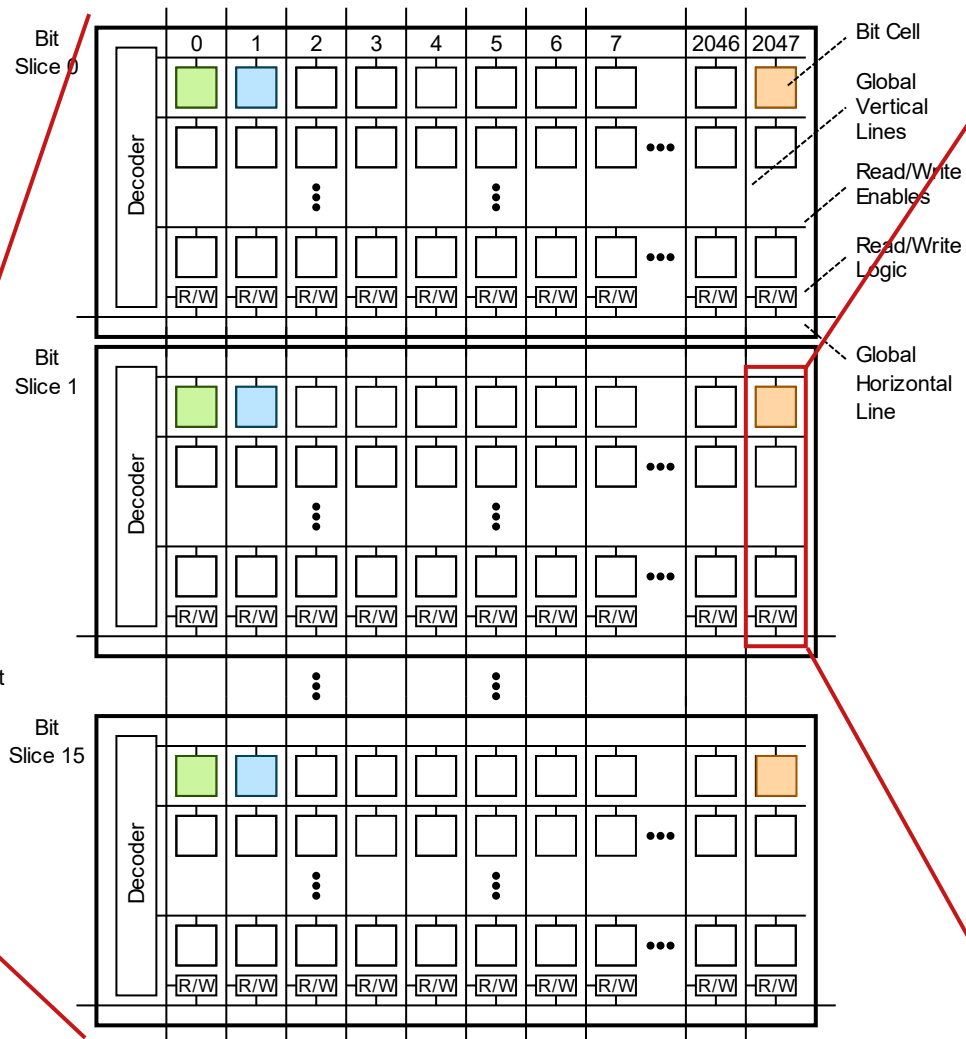
Global Horizontal Line



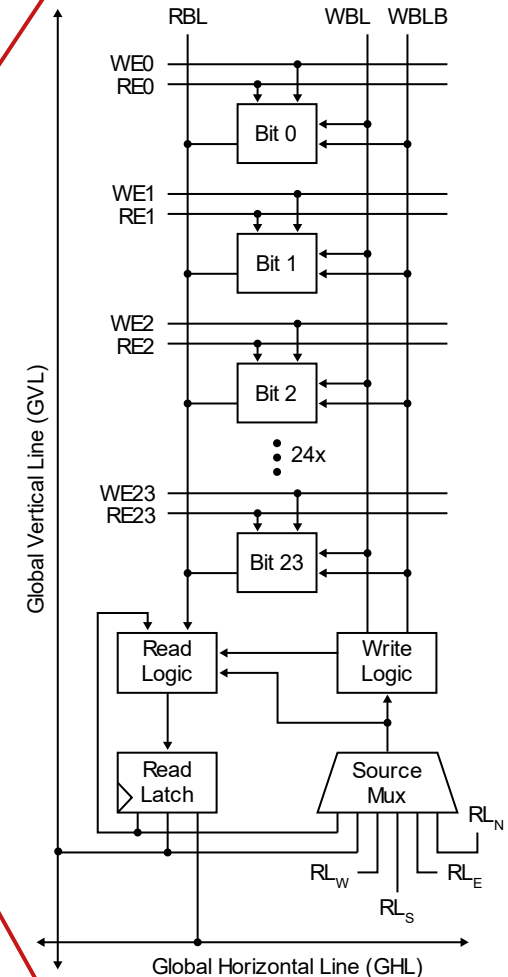
System Overview



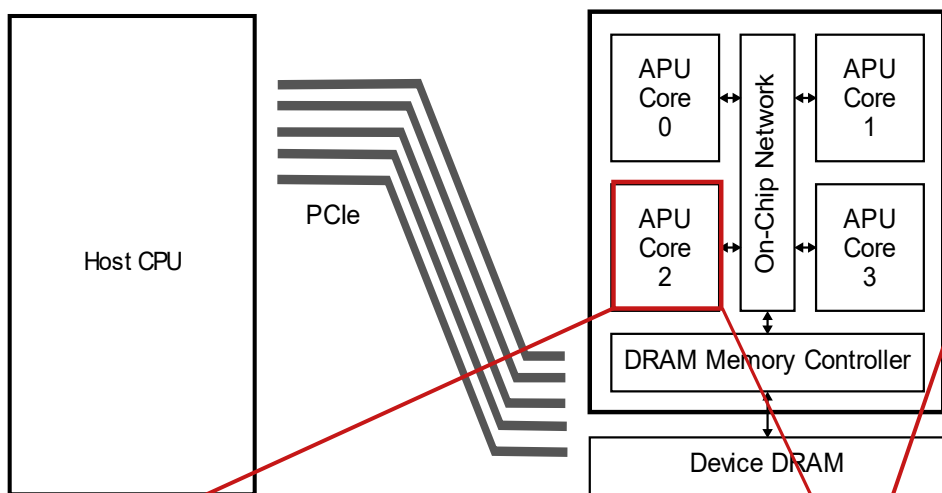
APU Core Logical View



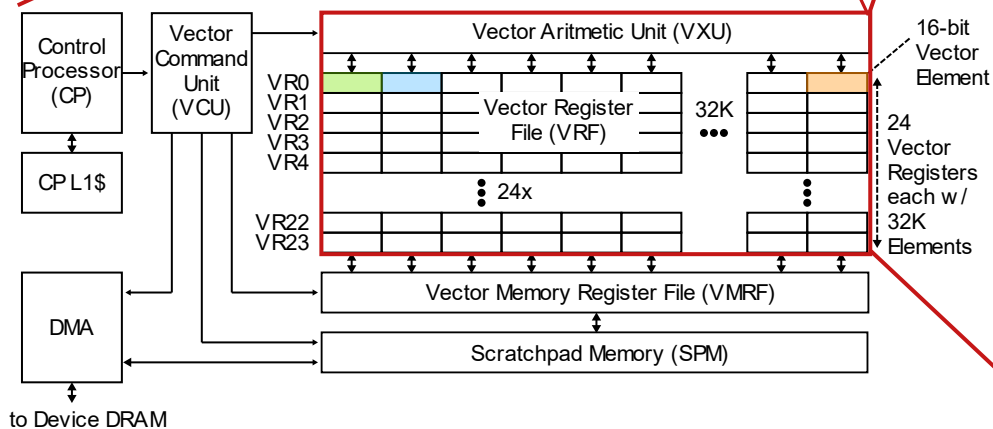
SRAM Bank Physical View



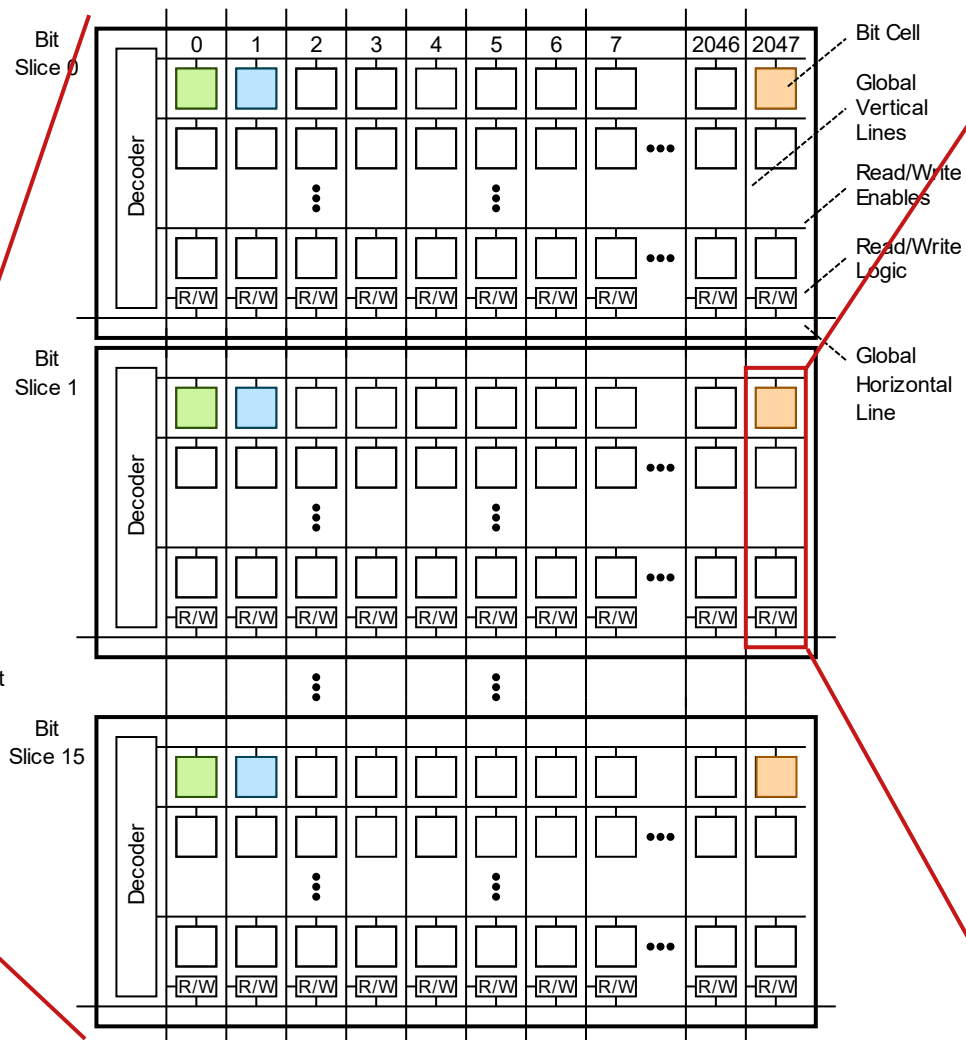
Bit Processor Circuitry



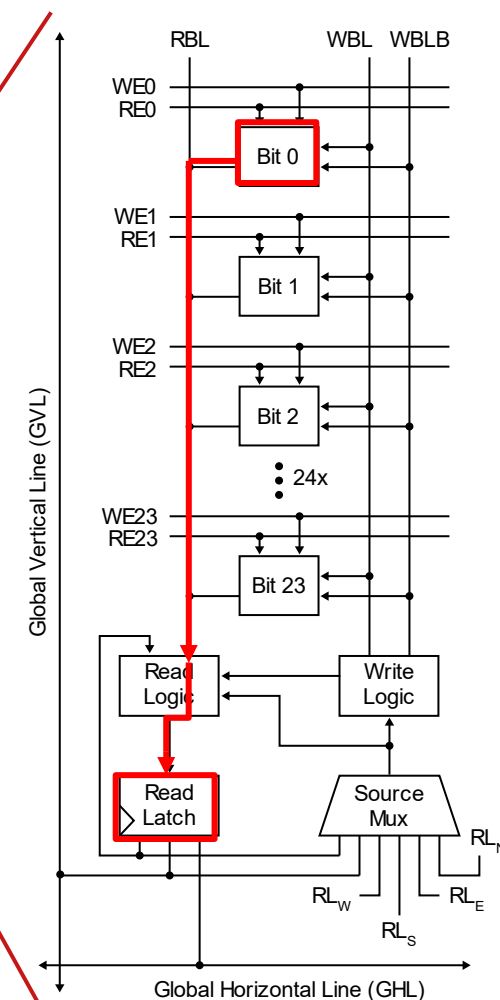
System Overview



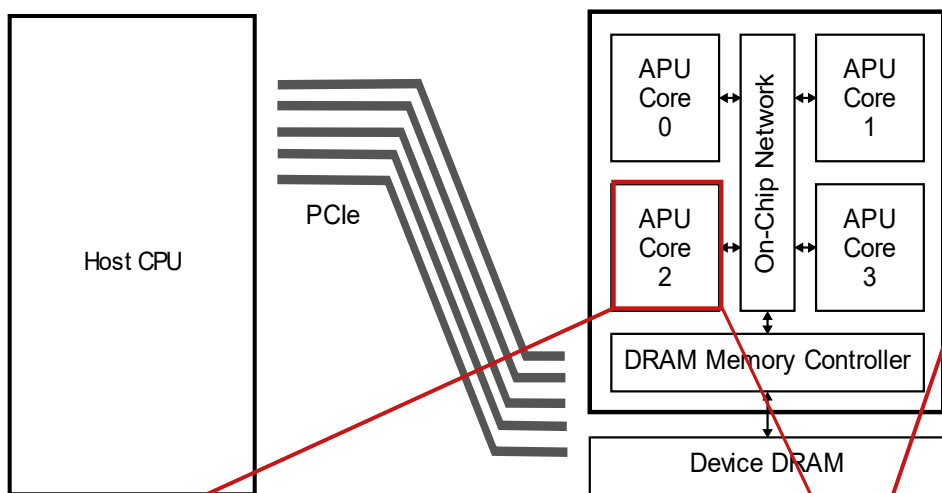
APU Core Logical View



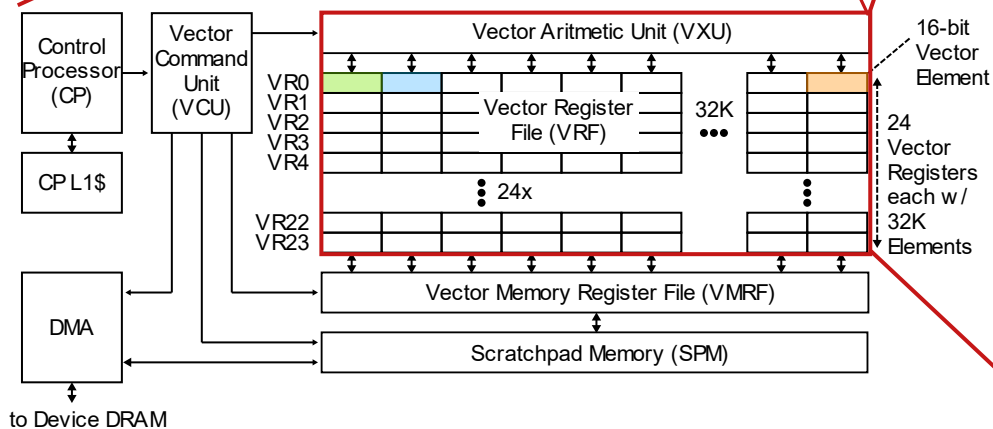
SRAM Bank Physical View



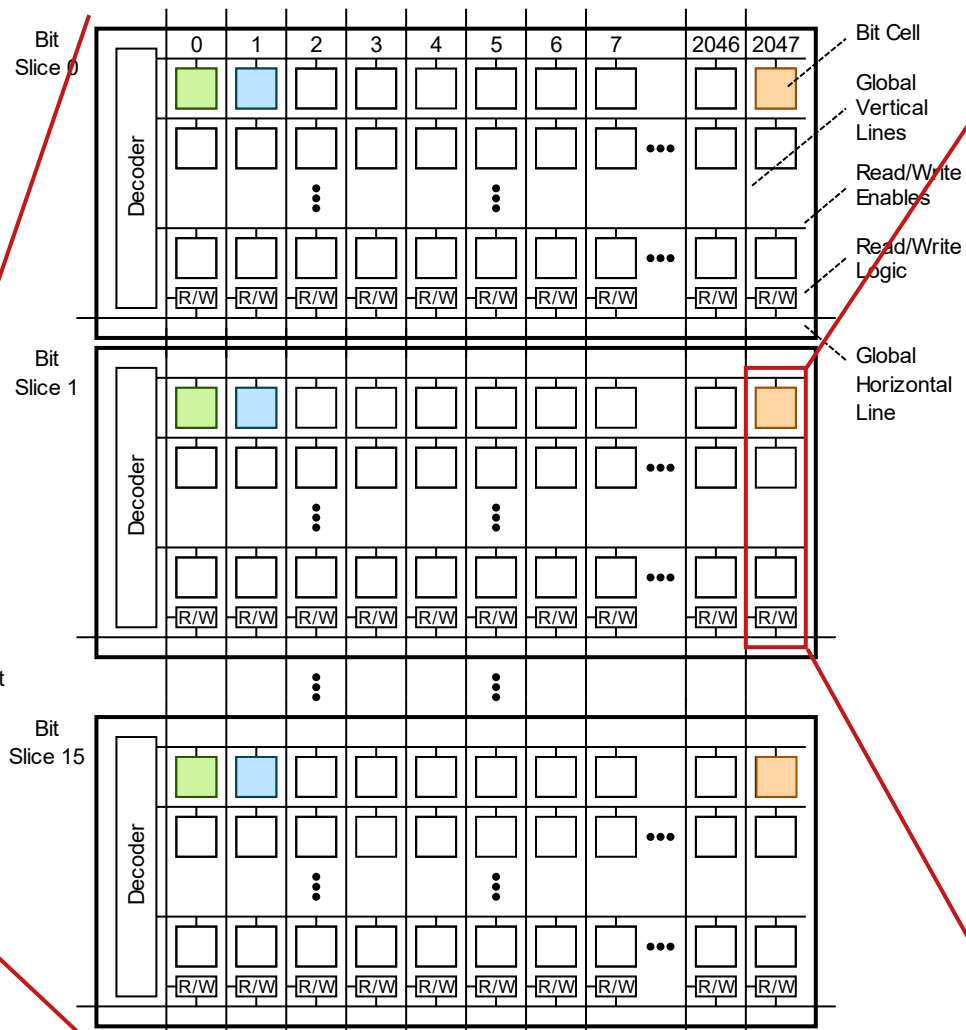
Bit Processor Circuitry



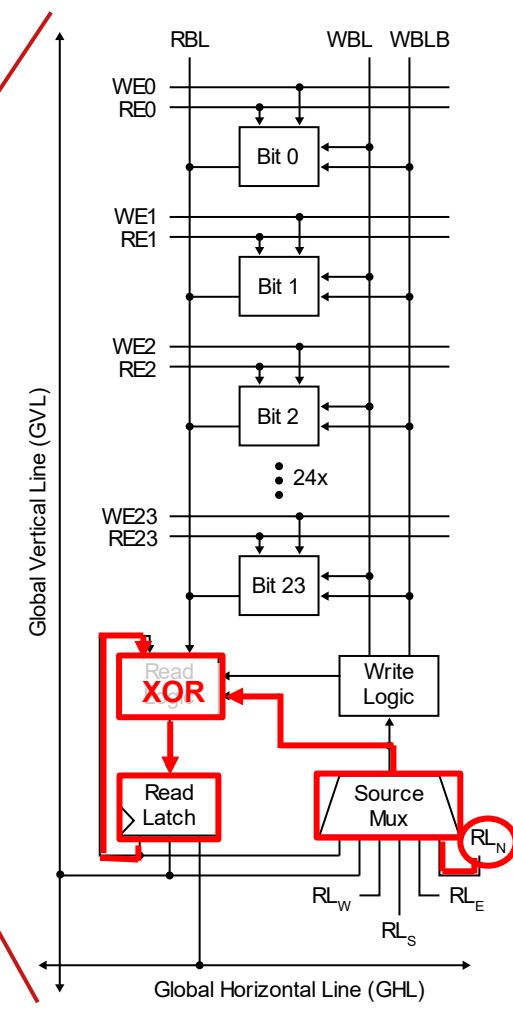
System Overview



APU Core Logical View



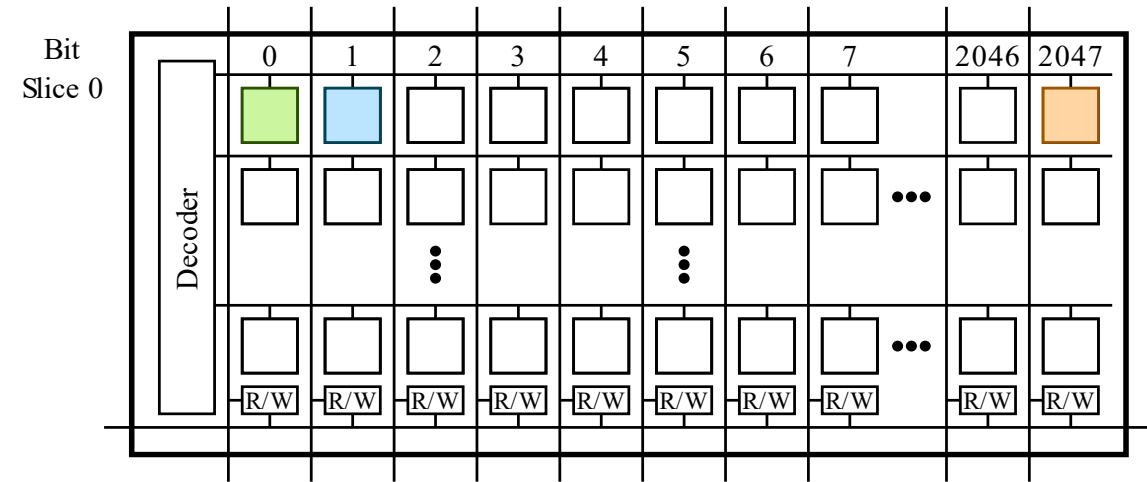
SRAM Bank Physical View



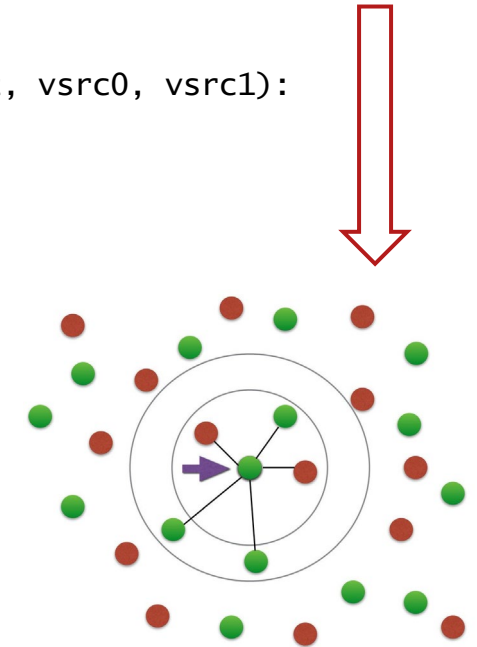
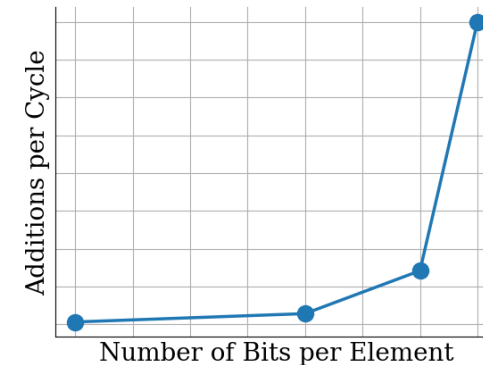
Bit Processor Circuitry

Supporting a Virtual Vector Instruction Set on a Commercial Compute-in-SRAM Accelerator

- Motivation
- APU Microarchitecture
- **APU Microcoding**
- Virtual Vector Instruction Set
- Microbenchmarking Results



```
APL_FRAG _frag_bitwise_or(vdst, vsrc0, vsrc1):  
0xFFFF: RL = VRF[vsrc0];  
0xFFFF: RL |= VRF[vsrc1];  
0xFFFF: VRF[vdst] = RL;
```



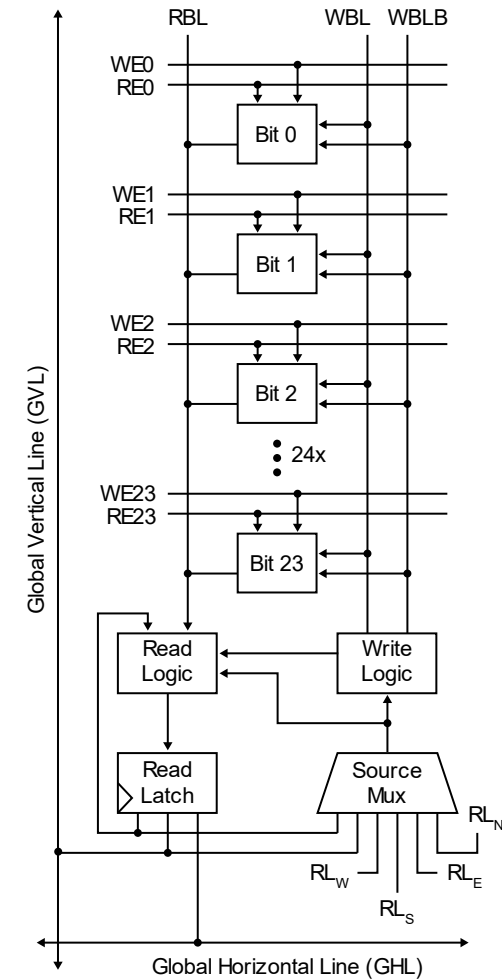
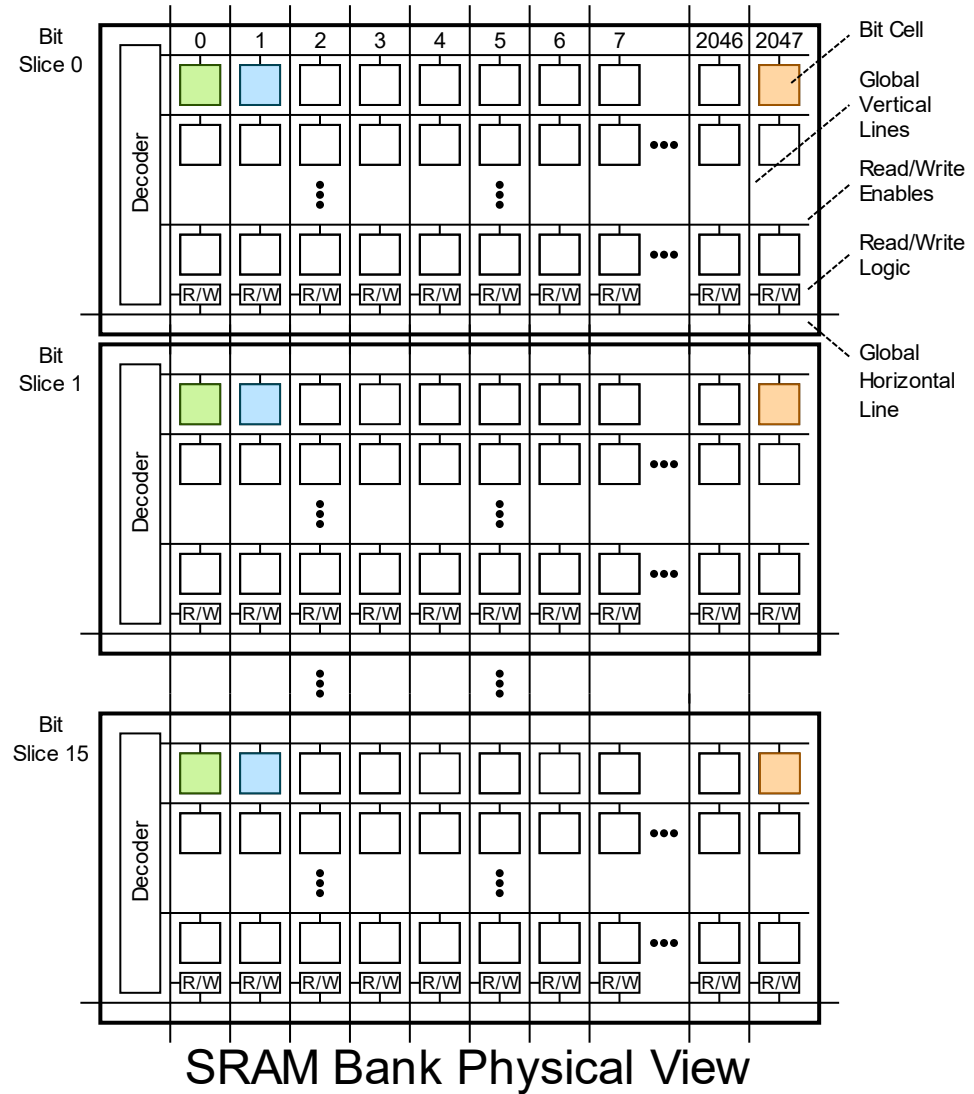
EXAMPLE #1: BITWISE OR

```

APL_FRAG _frag_bitwise_or(vdst, vsrc0, vsrc1):
    0xFFFF: RL = VRF[vsrc0];
    0xFFFF: RL |= VRF[vsrc1];
    0xFFFF: VRF[vdst] = RL;
    
```

```

GAL_TASK_ENTRY_POINT(apl_task_1, in, out)
{
    gvm1_init_once();
    struct commonif_struct *cmn_handle =
        (struct commonif_struct *)in;
    enum gvm1_vr16 a = GVM1_VR16_0;
    enum gvm1_vr16 b = GVM1_VR16_1;
    enum gvm1_vr16 output = GVM1_VR16_2;
    size_t vlen = 32768;
    vload(a, cmn_handle->ahndl, vlen);
    vload(b, cmn_handle->bhndl, vlen);
    bitwise_or(output, a, b);
    vstore(cmn_handle->output_hndl, output, vlen);
    return SUCCESS;
}
    
```



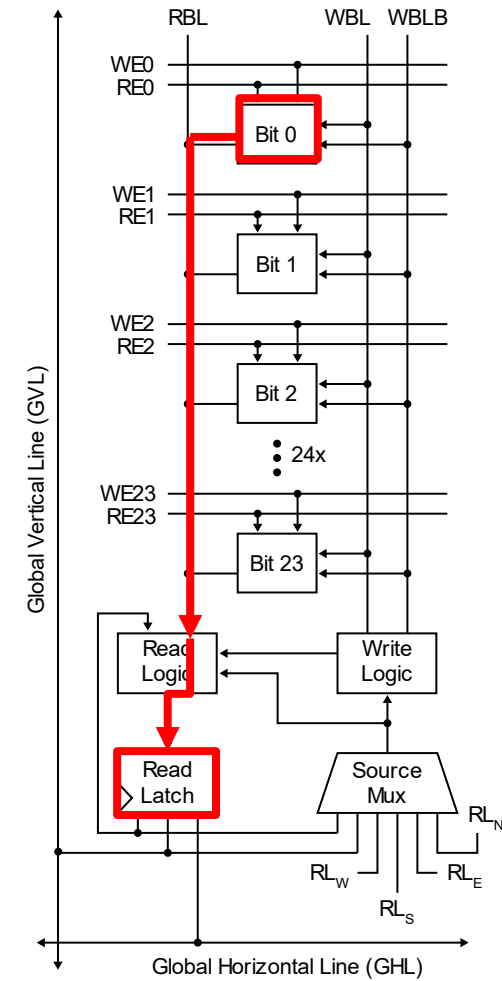
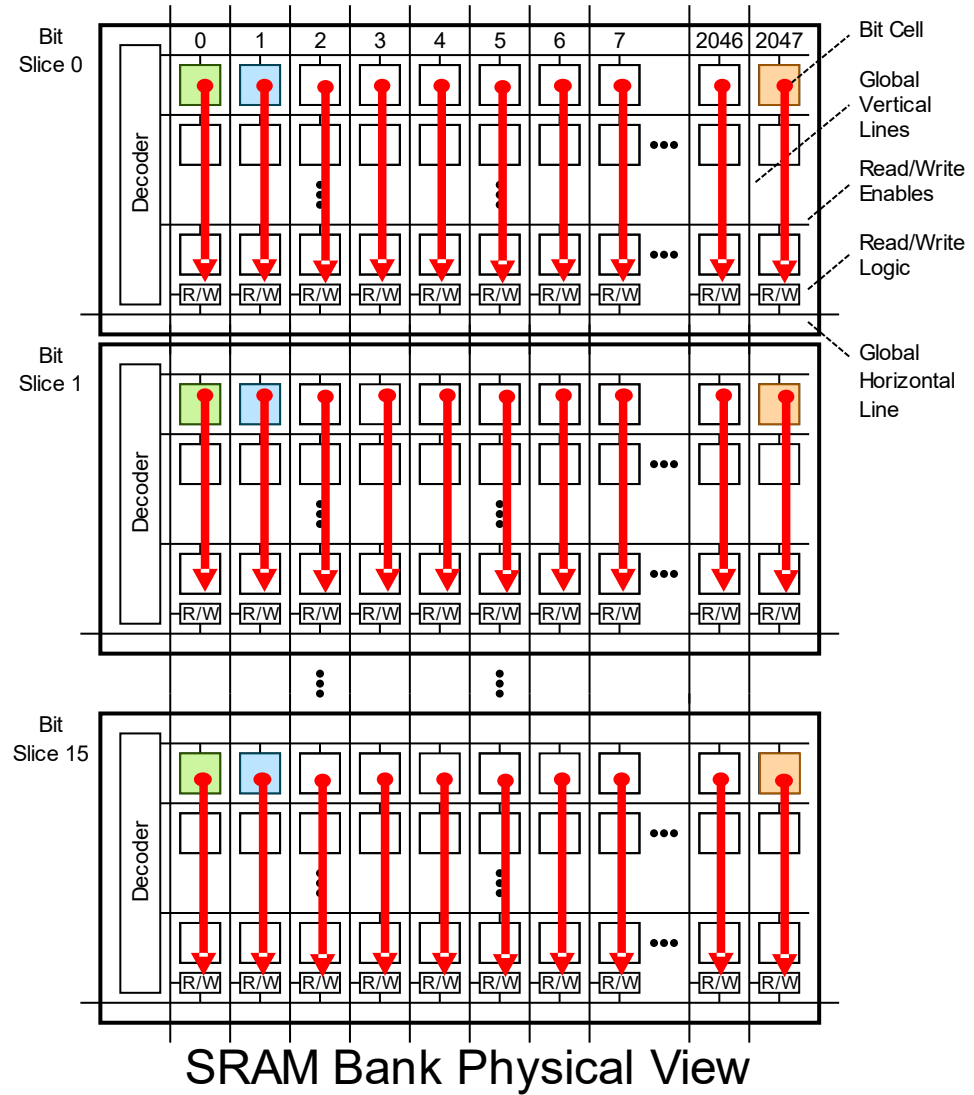
EXAMPLE #1: BITWISE OR

```
APL_FRAG frag_bitwise_or(vdst, vsrc0, vsrc1):
```

```
0xFFFF: RL = VRF[vsrc0];
```

```
0xFFFF: RL |= VRF[vsrc1];
```

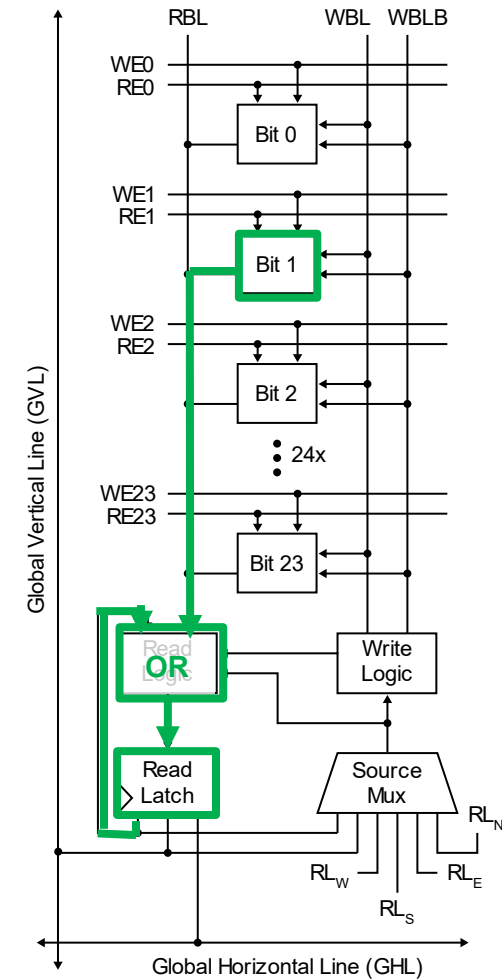
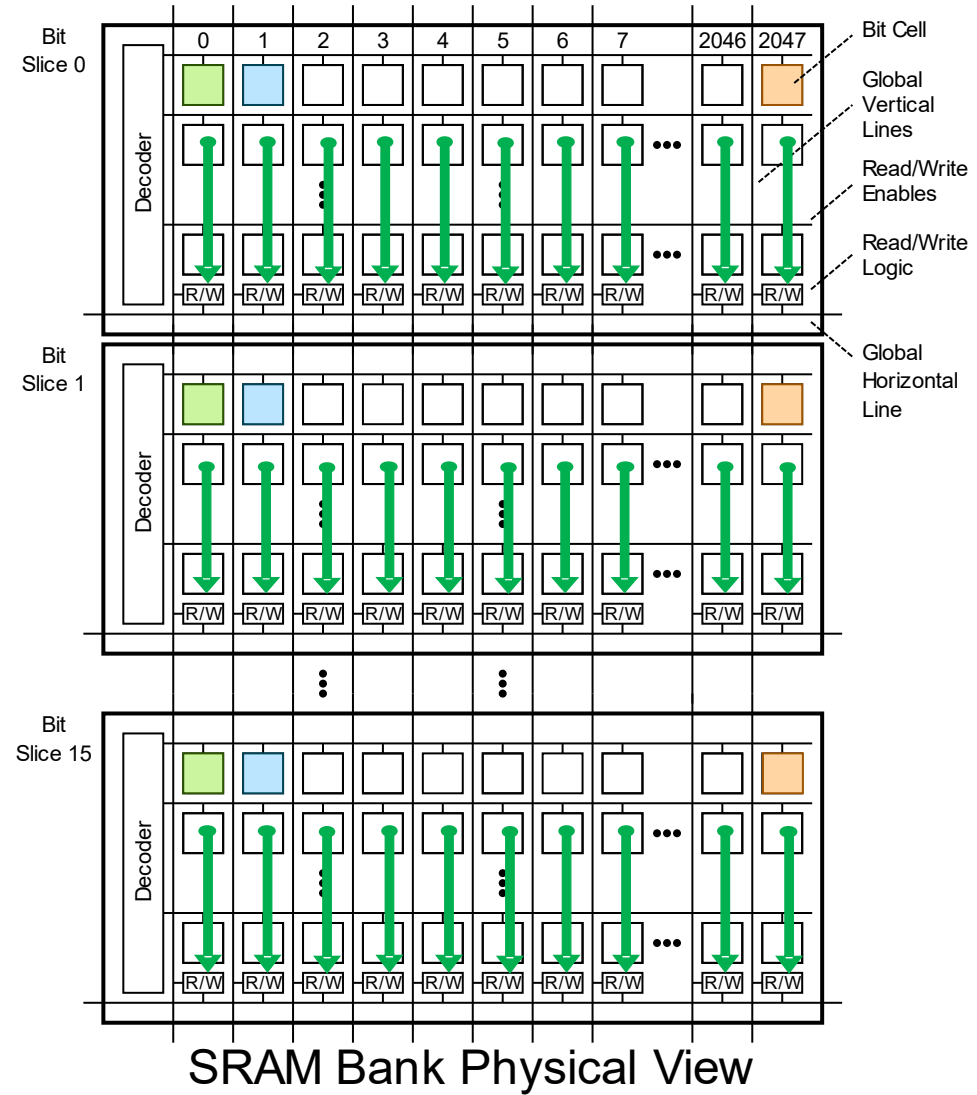
```
0xFFFF: VRF[vdst] = RL;
```



EXAMPLE #1: BITWISE OR

```

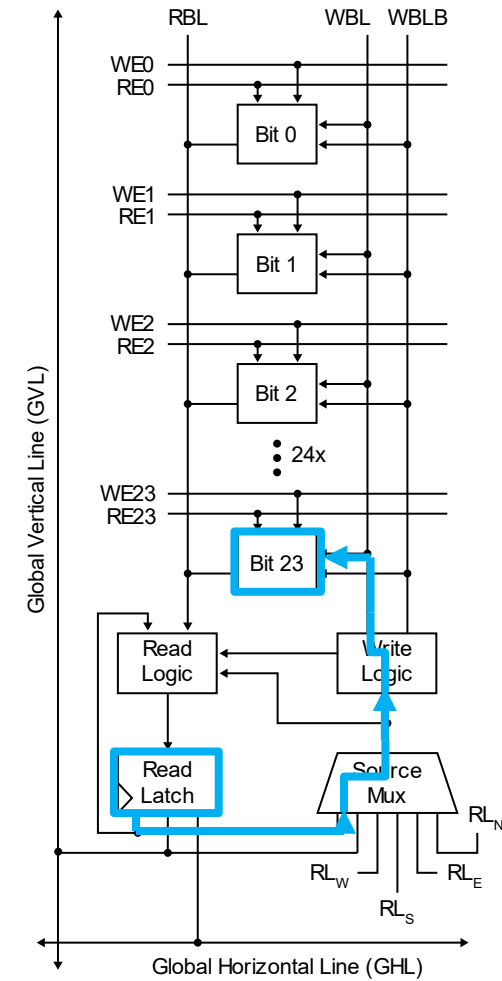
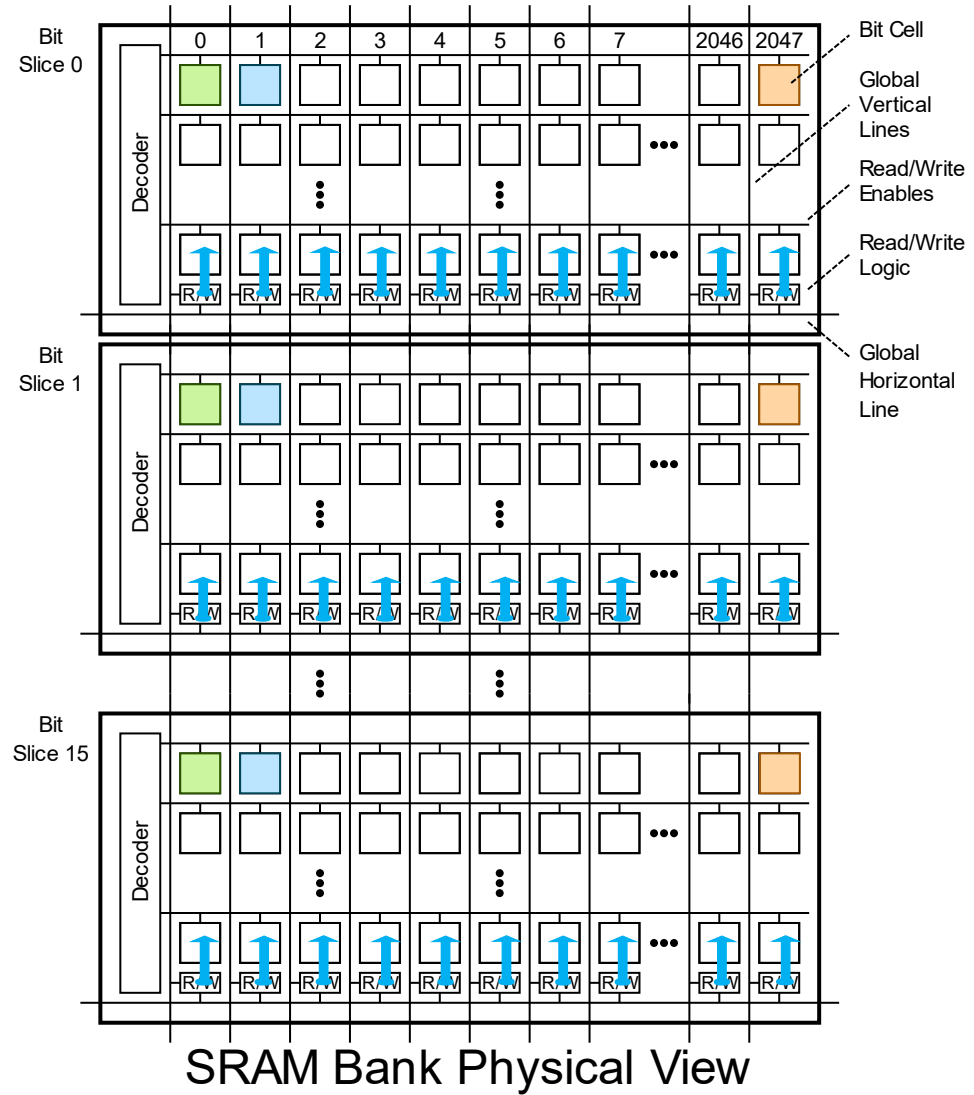
APL_FRAG _frag_bitwise_or(vdst, vsrc0, vsrc1):
  0xFFFF: RL = VRF[vsrc0];
  0xFFFF: RL |= VRF[vsrc1];
  0xFFFF: VRF[vdst] = RL;
  
```



EXAMPLE #1: BITWISE OR

```

APL_FRAG _frag_bitwise_or(vdst, vsrc0, vsrc1):
  0xFFFF: RL = VRF[vsrc0];
  0xFFFF: RL |= VRF[vsrc1];
  0xFFFF: VRF[vdst] = RL;
  
```



EXAMPLE #2: VECTOR-VECTOR ADD (BIT 0)

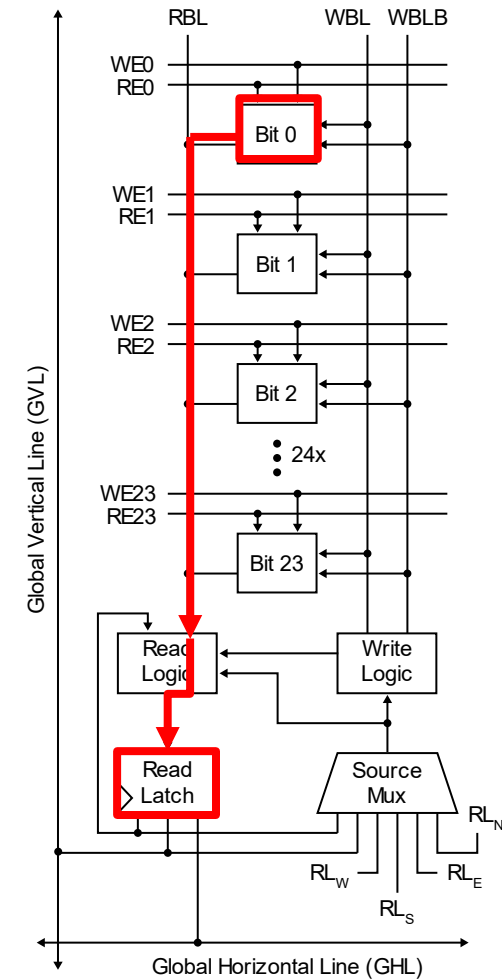
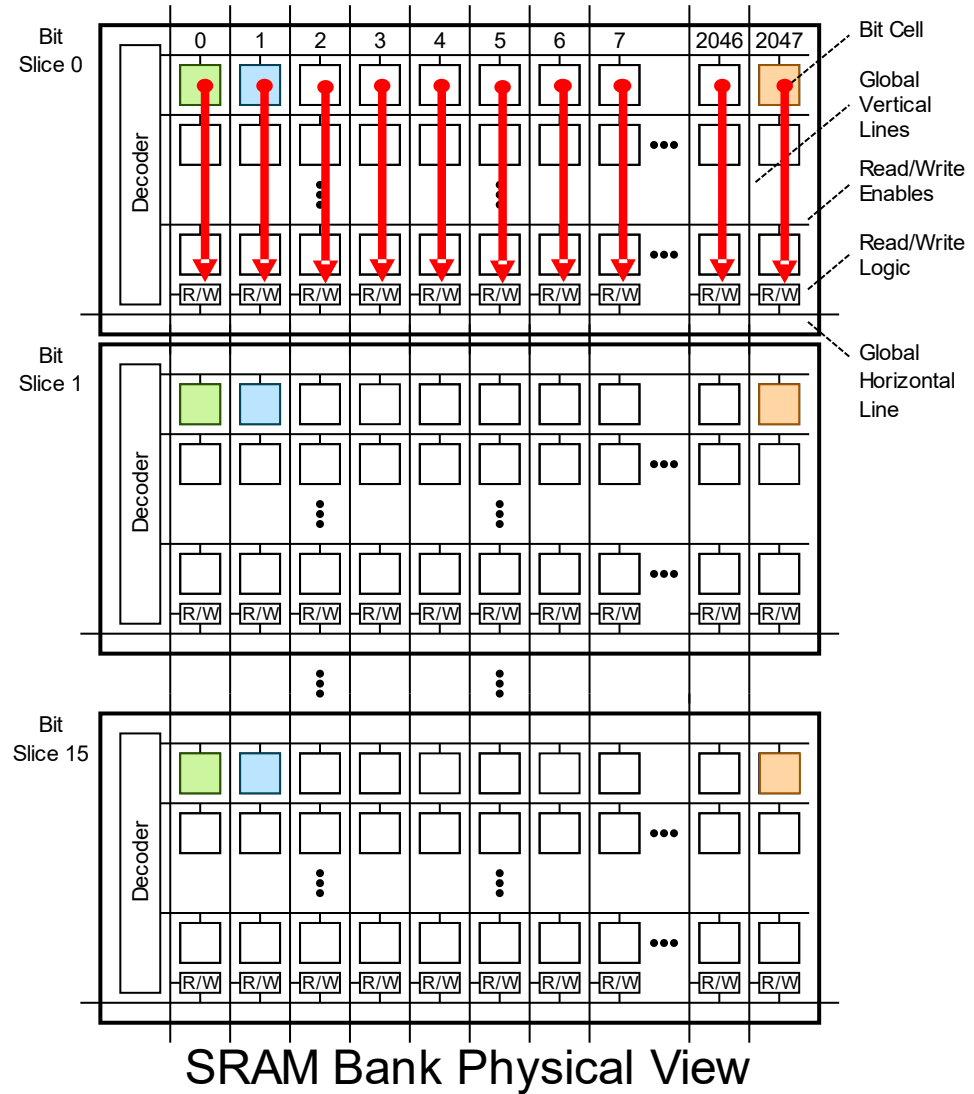
```
APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
```

```
// ---- bit 0 ----
// vdst = vsrc0 ^ vsrc1
0x0001: RL = VRF[vsrc0];
0x0001: RL ^= VRF[vsrc1];
0x0001: VRF[vdst] = RL;

// cout = vsrc0 * vsrc1
0x0001: RL = VRF[vsrc0, vsrc1];

...
```

vsrc0	vsrc1	vdst	c_out
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1



EXAMPLE #2: VECTOR-VECTOR ADD (BIT 0)

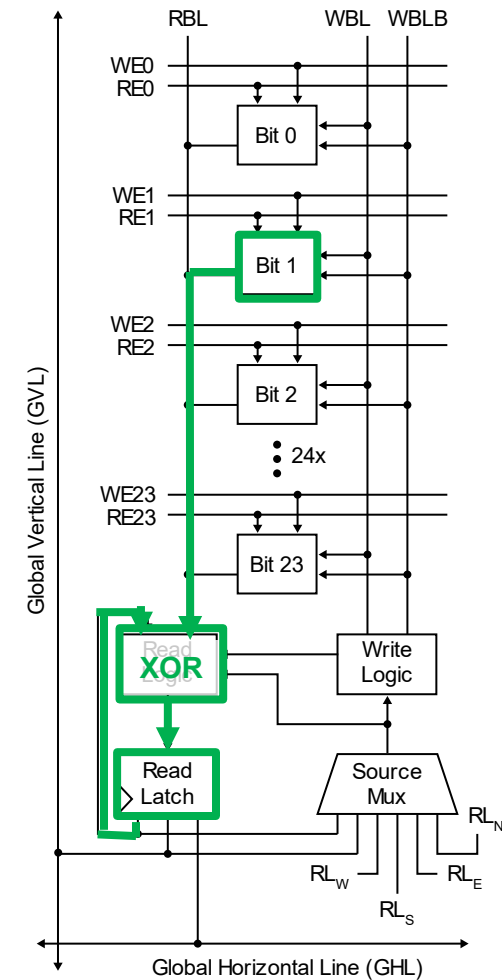
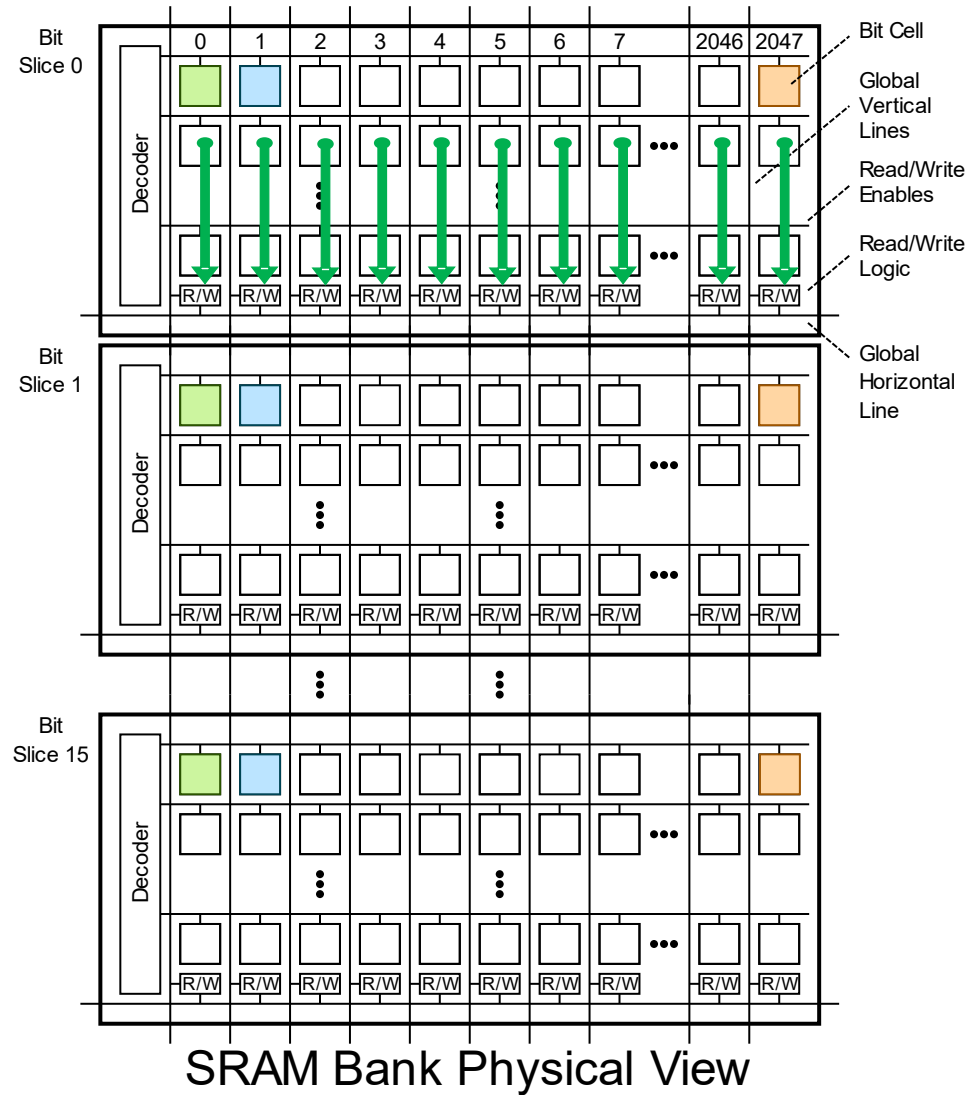
```
APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
```

```
// ---- bit 0 ----
// vdst = vsrc0 ^ vsrc1
0x0001: RL = VRF[vsrc0];
0x0001: RL ^= VRF[vsrc1];
0x0001: VRF[vdst] = RL;
```

```
// cout = vsrc0 * vsrc1
0x0001: RL = VRF[vsrc0, vsrc1];
```

...

vsrc0	vsrc1	vdst	c_out
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1



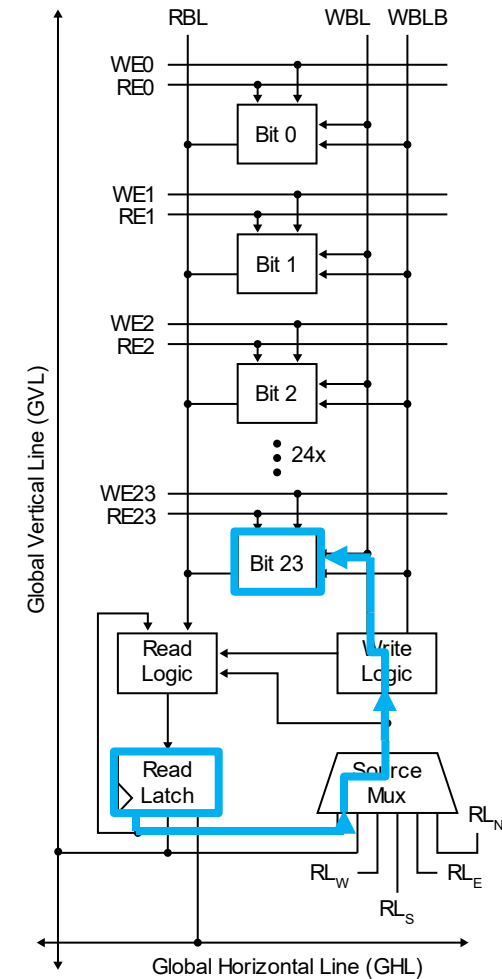
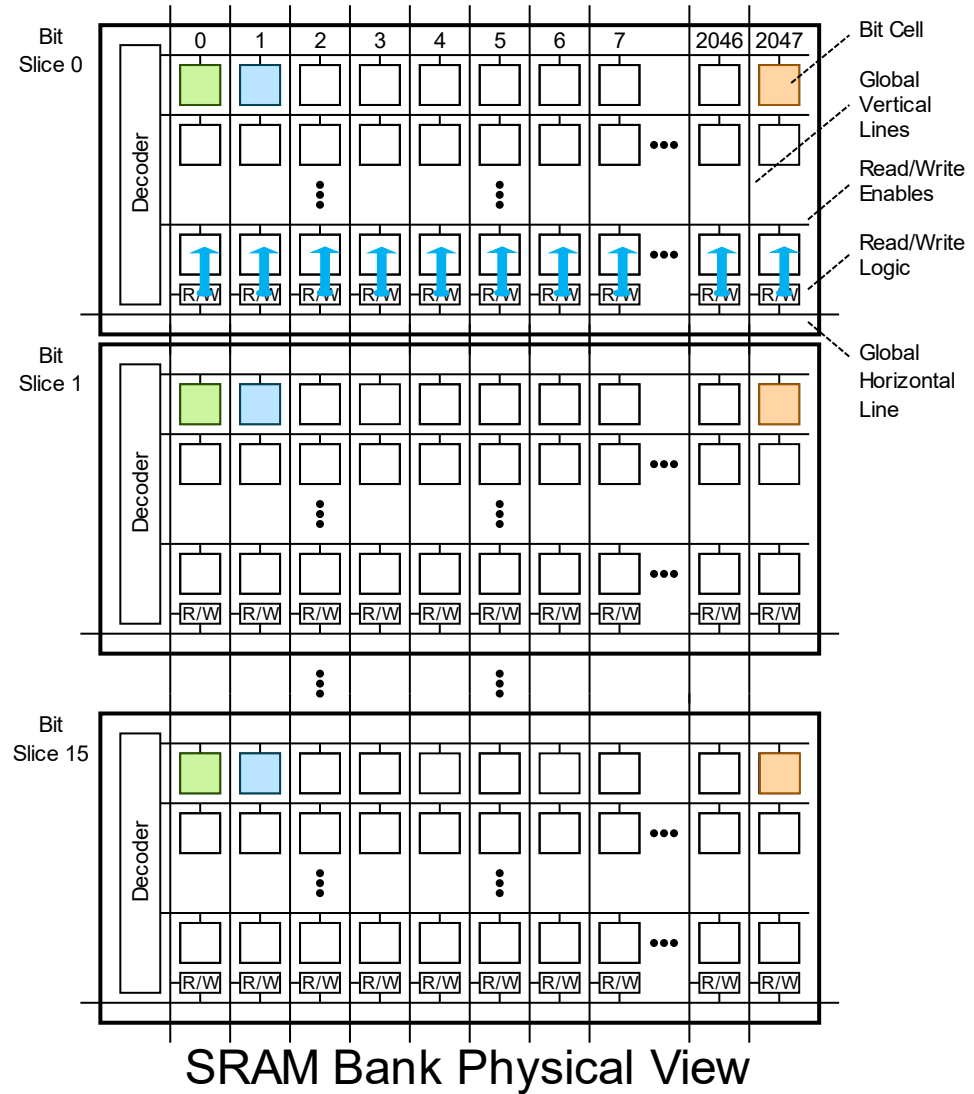
EXAMPLE #2: VECTOR-VECTOR ADD (BIT 0)

```
APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
```

```
// ---- bit 0 ----
// vdst = vsrc0 ^ vsrc1
0x0001: RL = VRF[vsrc0];
0x0001: RL ^= VRF[vsrc1];
0x0001: VRF[vdst] = RL;
```

```
// cout = vsrc0 * vsrc1
0x0001: RL = VRF[vsrc0, vsrc1];
...
```

vsrc0	vsrc1	vdst	c_out
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1



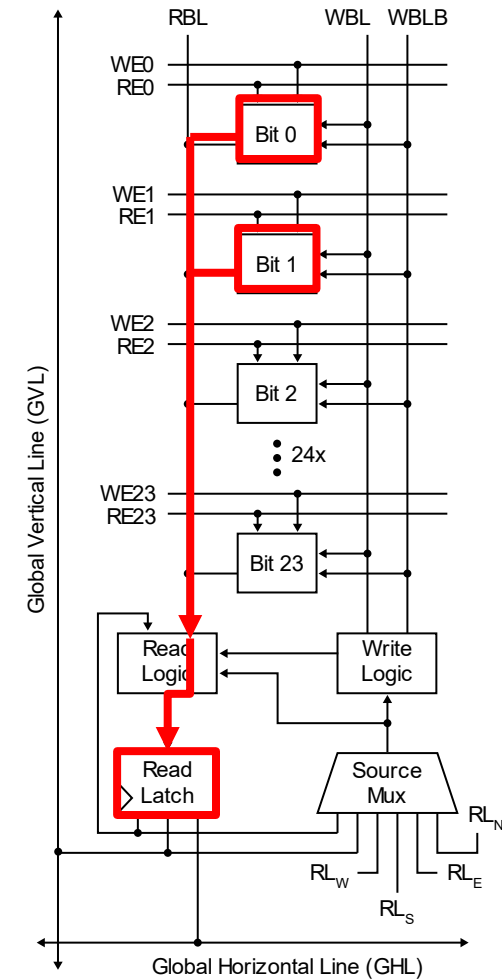
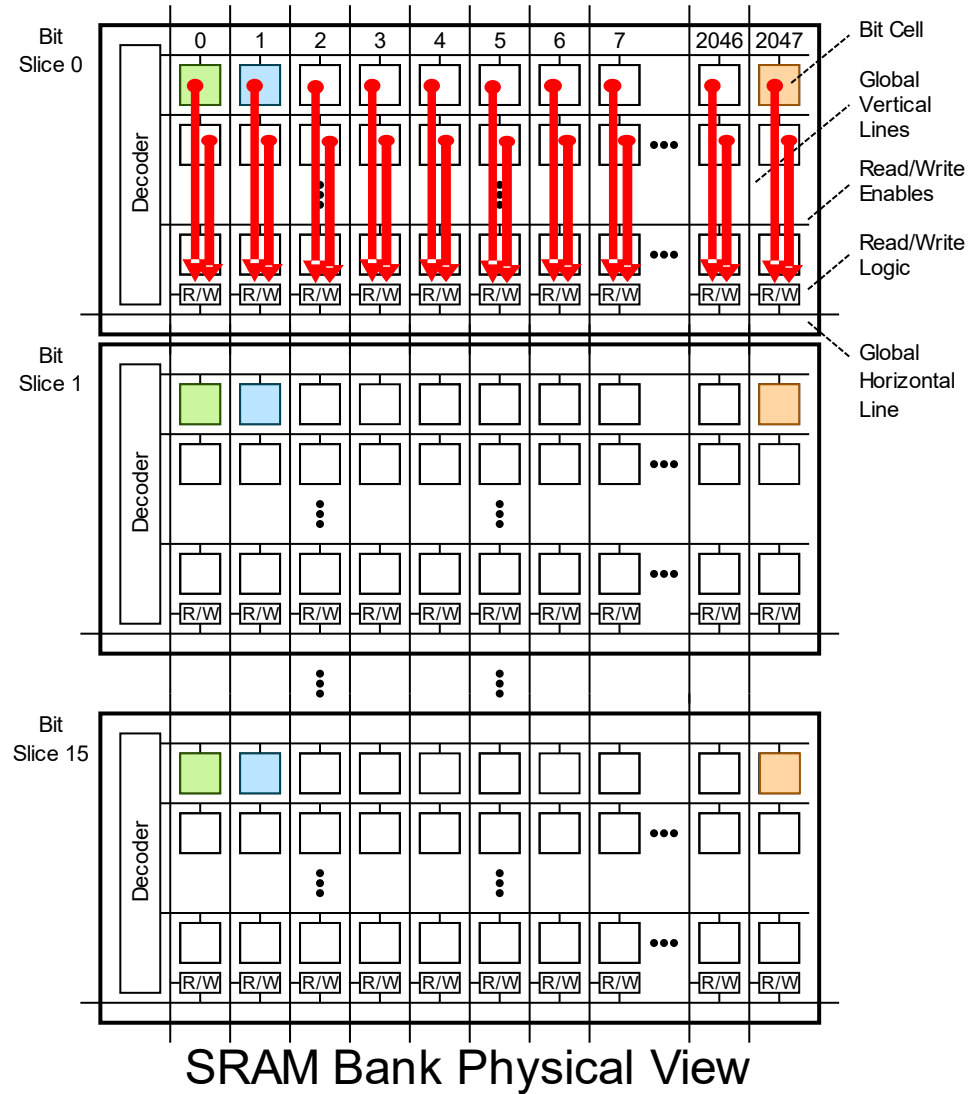
EXAMPLE #2: VECTOR-VECTOR ADD (BIT 0)

```
APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
```

```
// ---- bit 0 ----
// vdst = vsrc0 ^ vsrc1
0x0001: RL = VRF[vsrc0];
0x0001: RL ^= VRF[vsrc1];
0x0001: VRF[vdst] = RL;

// cout = vsrc0 * vsrc1
0x0001: RL = VRF[vsrc0, vsrc1];
...
```

vsrc0	vsrc1	vdst	c_out
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1



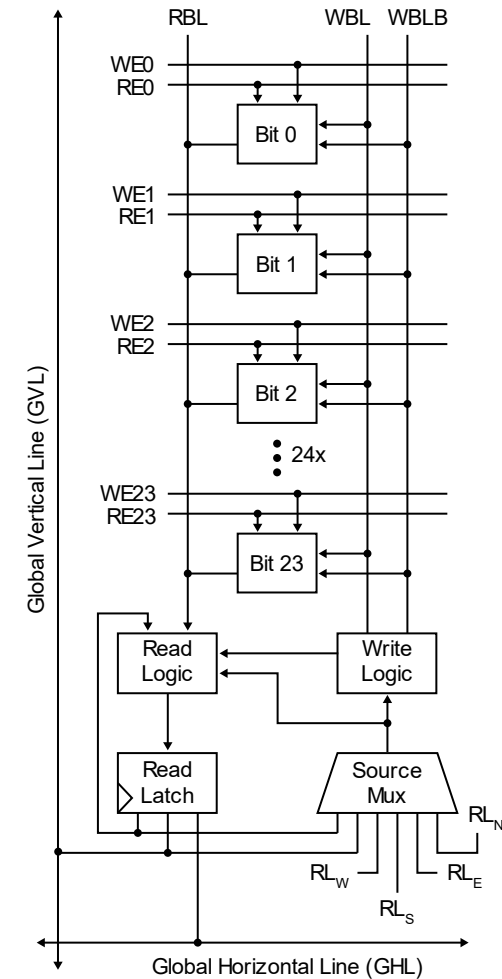
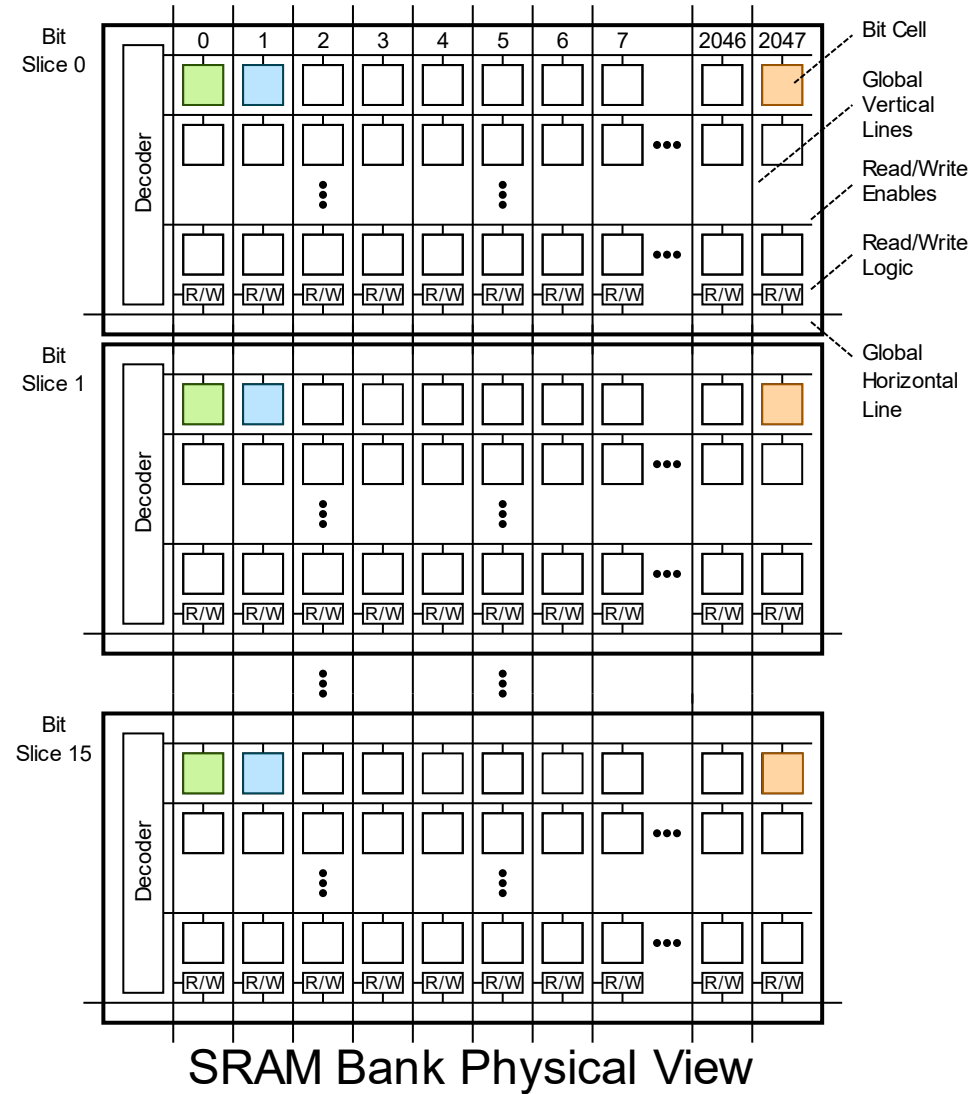
EXAMPLE #2: VECTOR-VECTOR ADD (BITS 1-15)

```

APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
...
// ---- bit 1 ----
// vdst = a ^ b ^ cin
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
(0x0001<<1): RL ^= RL_N;
(0x0001<<1): VRF[vdst] = RL;

// cout = a*b + b*cin + a*cin
(0x0001<<1): RL = VRF[vsrc0, vsrc1];
(0x0001<<1): VRF[temp_0] = RL;
(0x0001<<1): RL = VRF[vsrc1];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): VRF[temp_1] = RL;
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
...

```

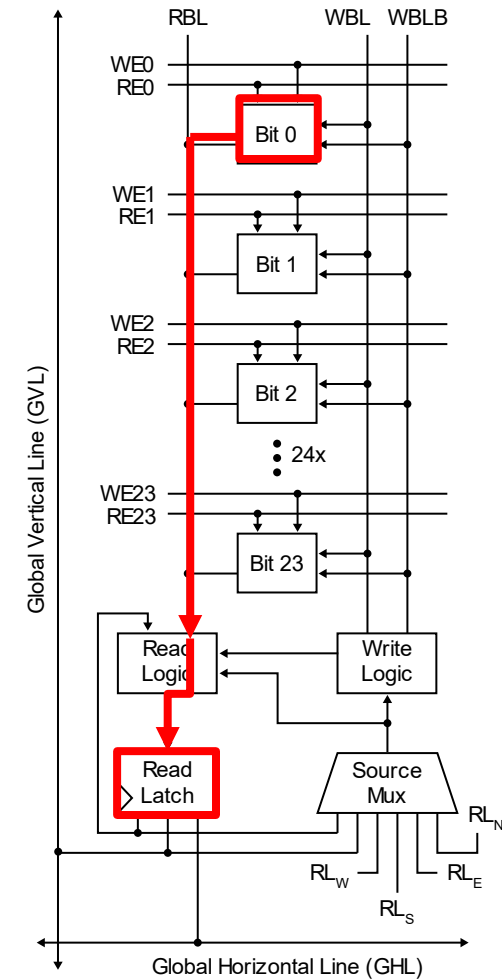
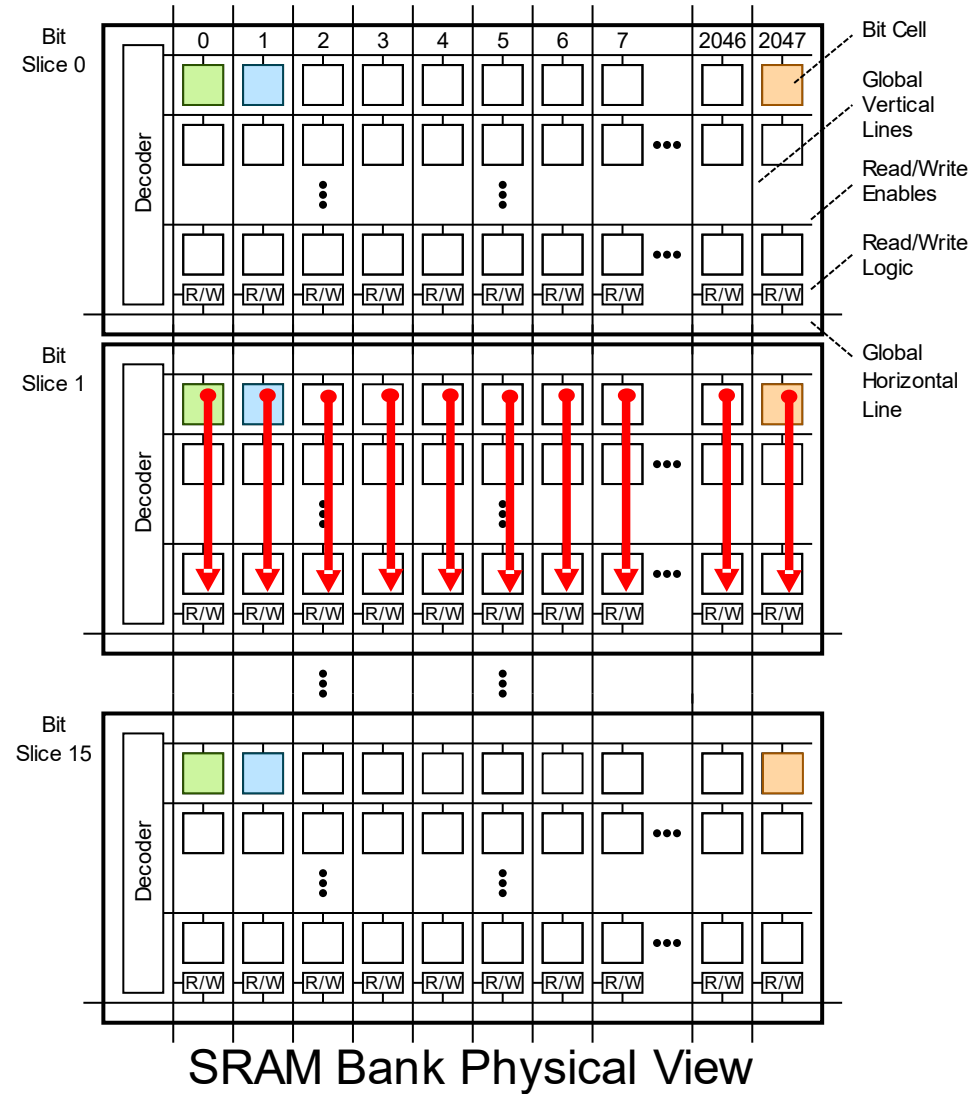


EXAMPLE #2: VECTOR-VECTOR ADD (BITS 1-15)

```

APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
...
// ---- bit 1 ----
// vdst = a ^ b ^ cin
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
(0x0001<<1): RL ^= RL_N;
(0x0001<<1): VRF[vdst] = RL;

// cout = a*b + b*cin + a*cin
(0x0001<<1): RL = VRF[vsrc0, vsrc1];
(0x0001<<1): VRF[temp_0] = RL;
(0x0001<<1): RL = VRF[vsrc1];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): VRF[temp_1] = RL;
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
...
    
```

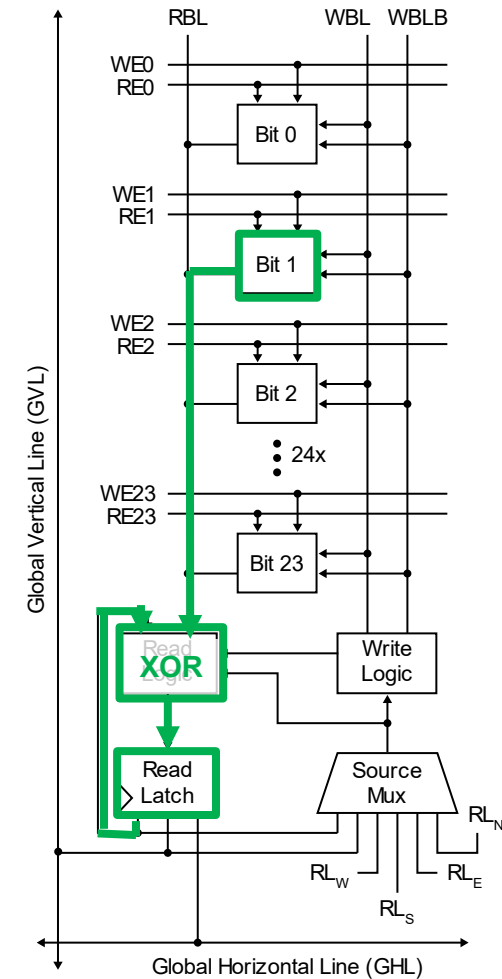
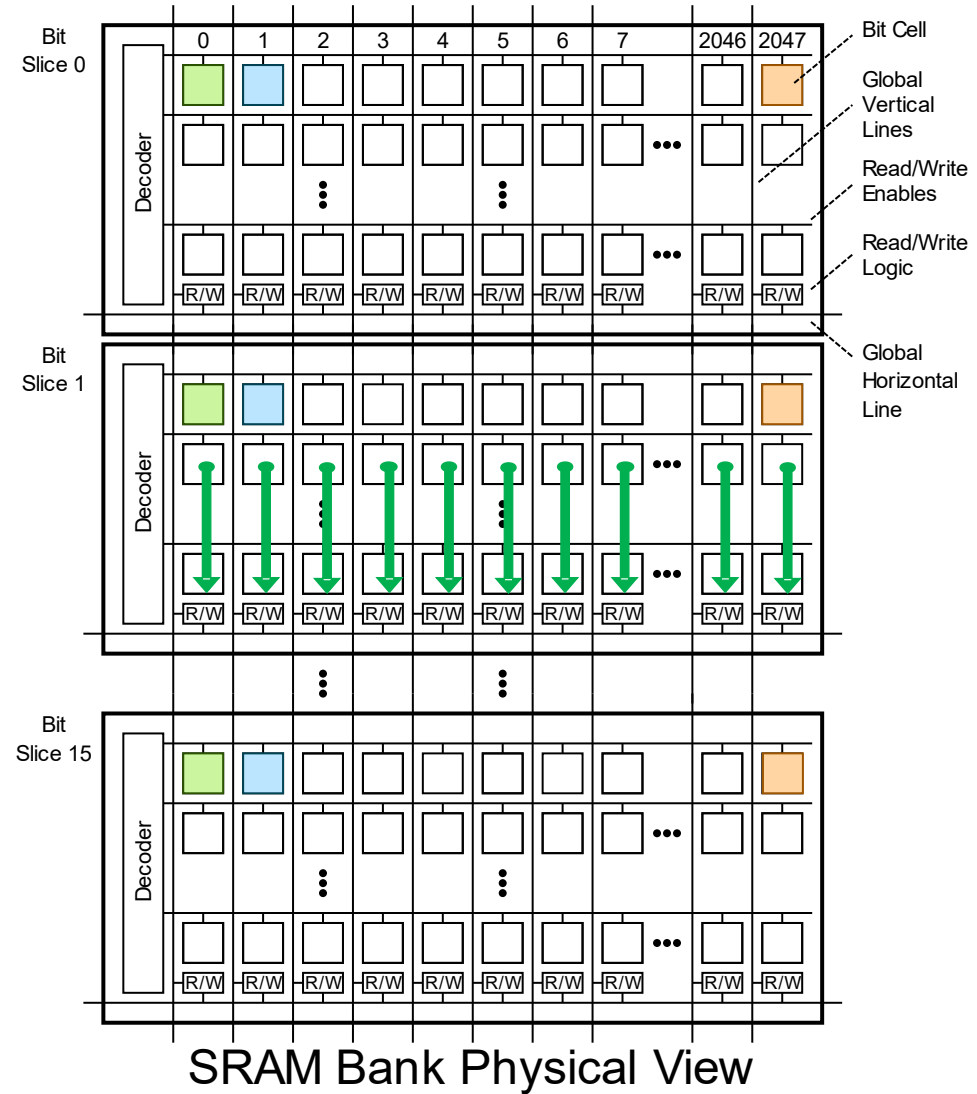


EXAMPLE #2: VECTOR-VECTOR ADD (BITS 1-15)

```

APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
...
// ---- bit 1 ----
// vdst = a ^ b ^ cin
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
(0x0001<<1): RL ^= RL_N;
(0x0001<<1): VRF[vdst] = RL;

// cout = a*b + b*cin + a*cin
(0x0001<<1): RL = VRF[vsrc0, vsrc1];
(0x0001<<1): VRF[temp_0] = RL;
(0x0001<<1): RL = VRF[vsrc1];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): VRF[temp_1] = RL;
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
...
    
```

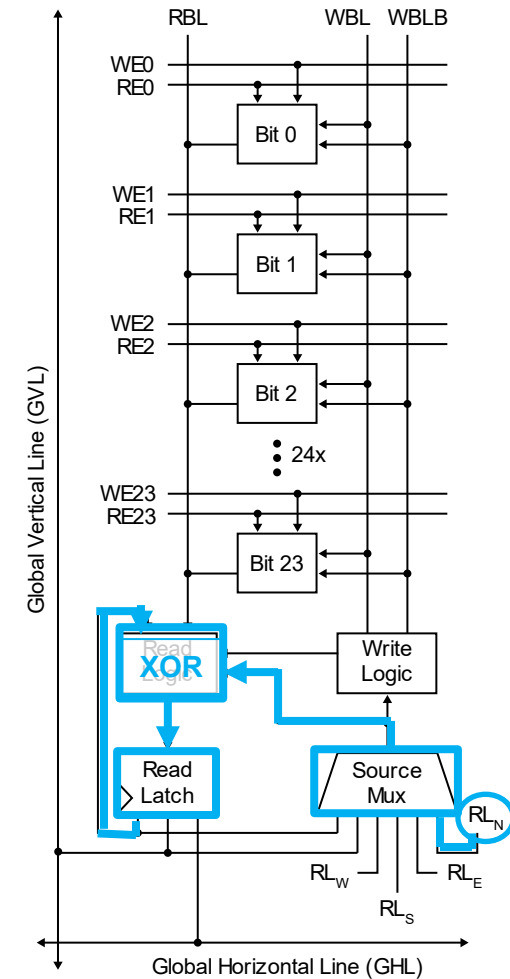
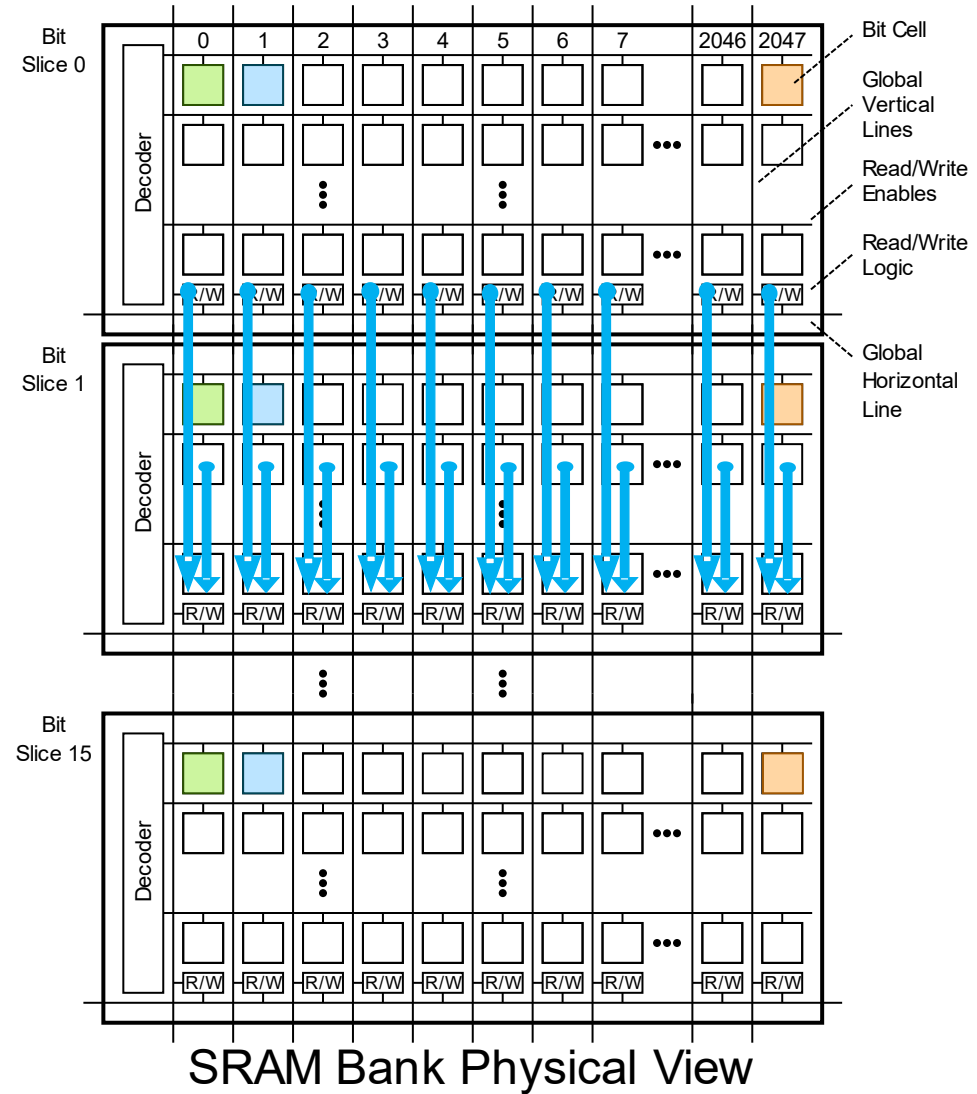


EXAMPLE #2: VECTOR-VECTOR ADD (BITS 1-15)

```

APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
...
// ---- bit 1 ----
// vdst = a ^ b ^ cin
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
(0x0001<<1): RL ^= RL_N;
(0x0001<<1): VRF[vdst] = RL;

// cout = a*b + b*cin + a*cin
(0x0001<<1): RL = VRF[vsrc0, vsrc1];
(0x0001<<1): VRF[temp_0] = RL;
(0x0001<<1): RL = VRF[vsrc1];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): VRF[temp_1] = RL;
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
...
    
```

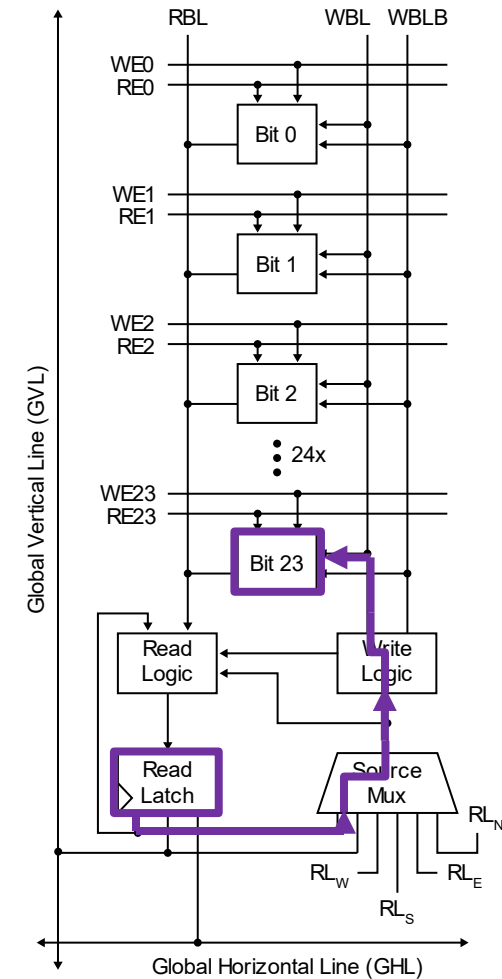
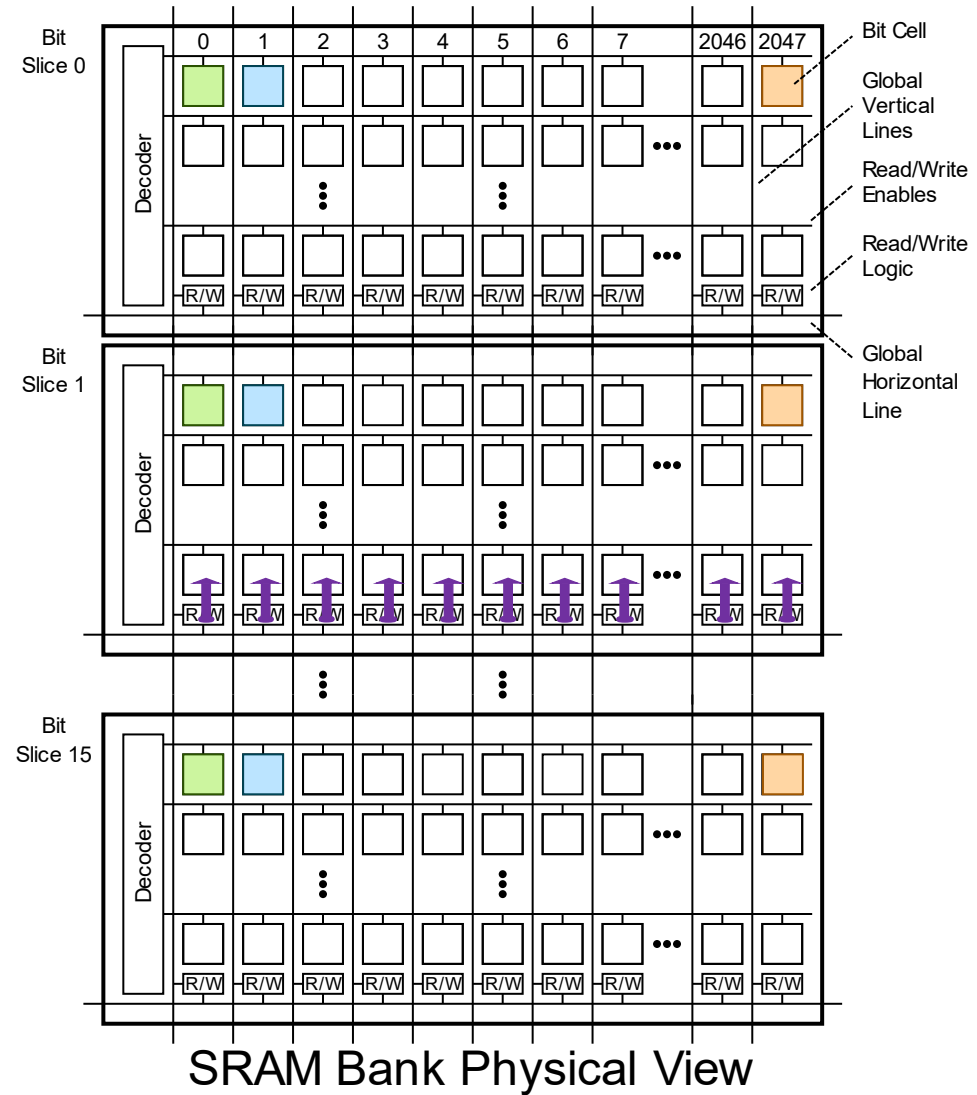


EXAMPLE #2: VECTOR-VECTOR ADD (BITS 1-15)

```

APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
...
// ---- bit 1 ----
// vdst = a ^ b ^ cin
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
(0x0001<<1): RL ^= RL_N;
(0x0001<<1): VRF[vdst] = RL;

// cout = a*b + b*cin + a*cin
(0x0001<<1): RL = VRF[vsrc0, vsrc1];
(0x0001<<1): VRF[temp_0] = RL;
(0x0001<<1): RL = VRF[vsrc1];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): VRF[temp_1] = RL;
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
...
    
```



EXAMPLE #2: VECTOR-VECTOR ADD (BITS 1-15)

```

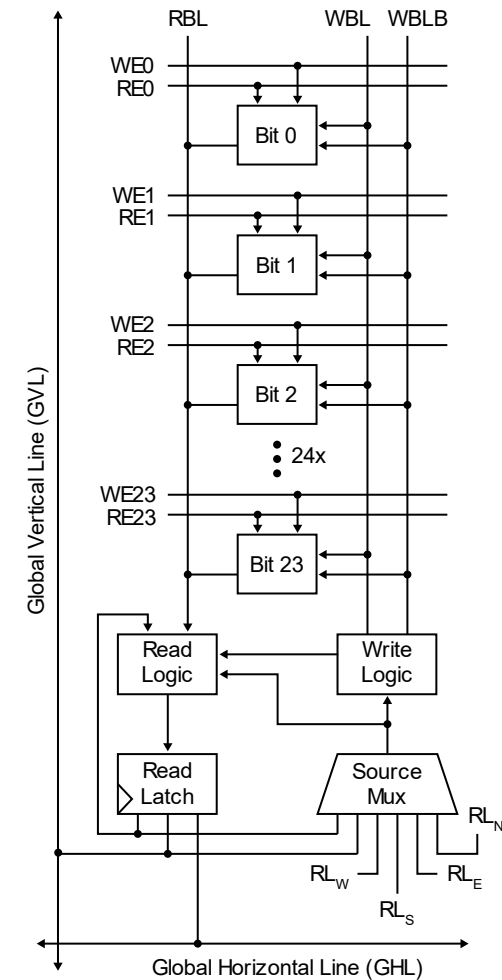
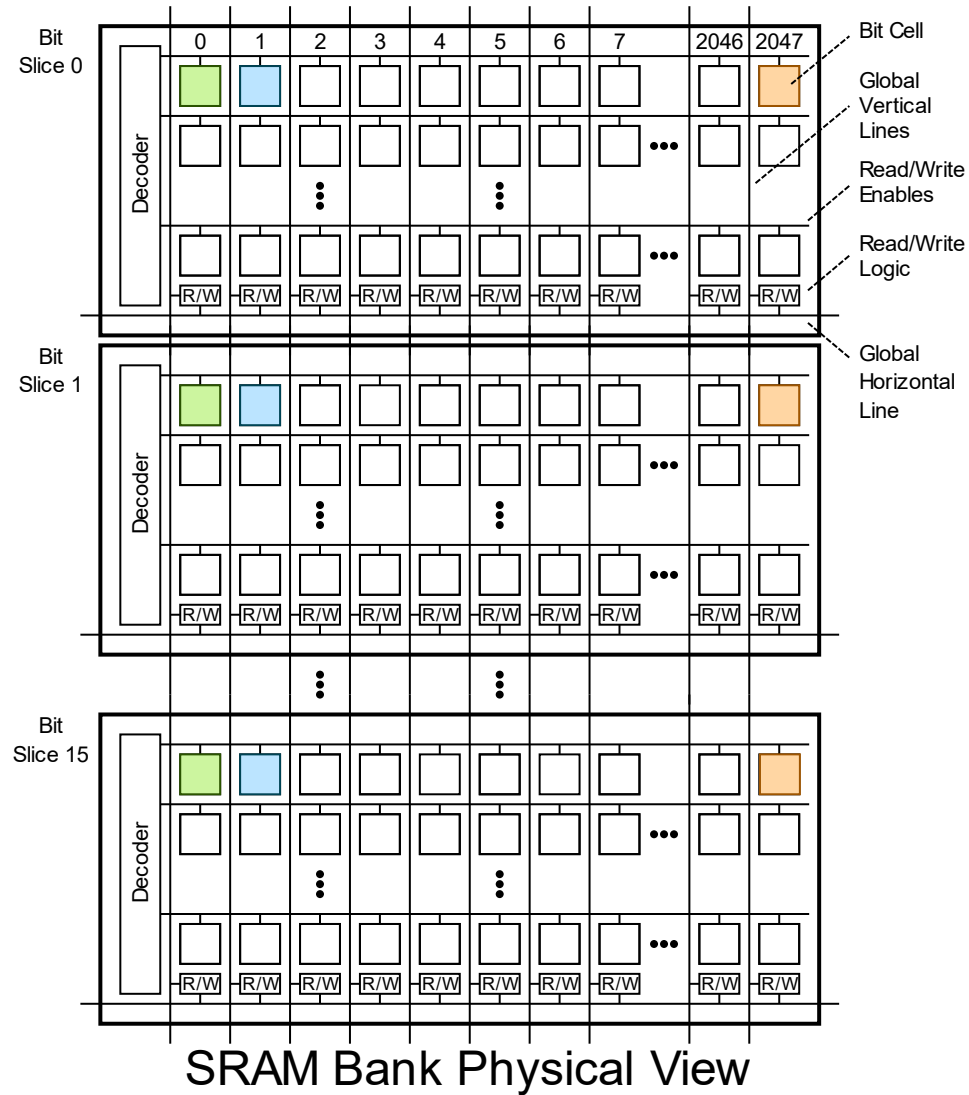
APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
...
// ---- bit 1 ----
// vdst = a ^ b ^ cin
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
(0x0001<<1): RL ^= RL_N;
(0x0001<<1): VRF[vdst] = RL;

// cout = a*b + b*cin + a*cin
(0x0001<<1): RL = VRF[vsrc0, vsrc1];
(0x0001<<1): VRF[temp_0] = RL;
(0x0001<<1): RL = VRF[vsrc1];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): VRF[temp_1] = RL;
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL &= RL_N;
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
...

```

Ripple-carry: 215 cycles

Carry-select: 12 cycles



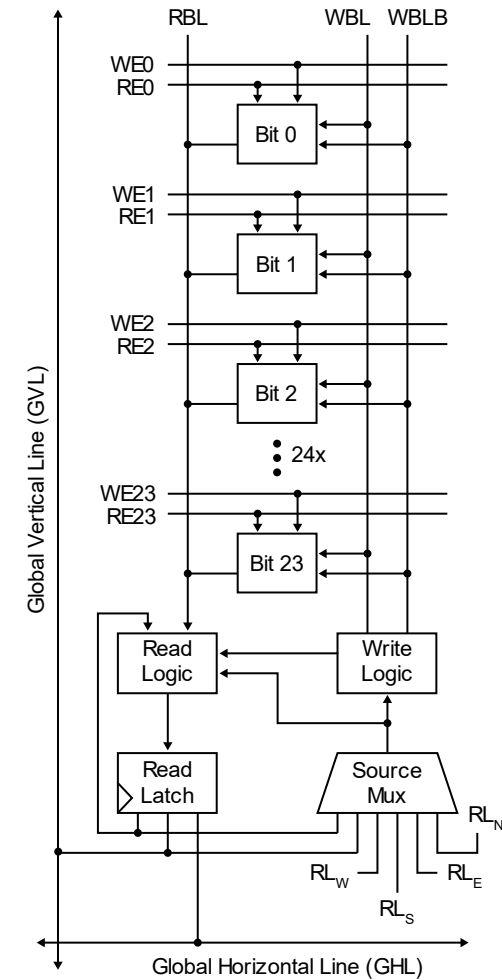
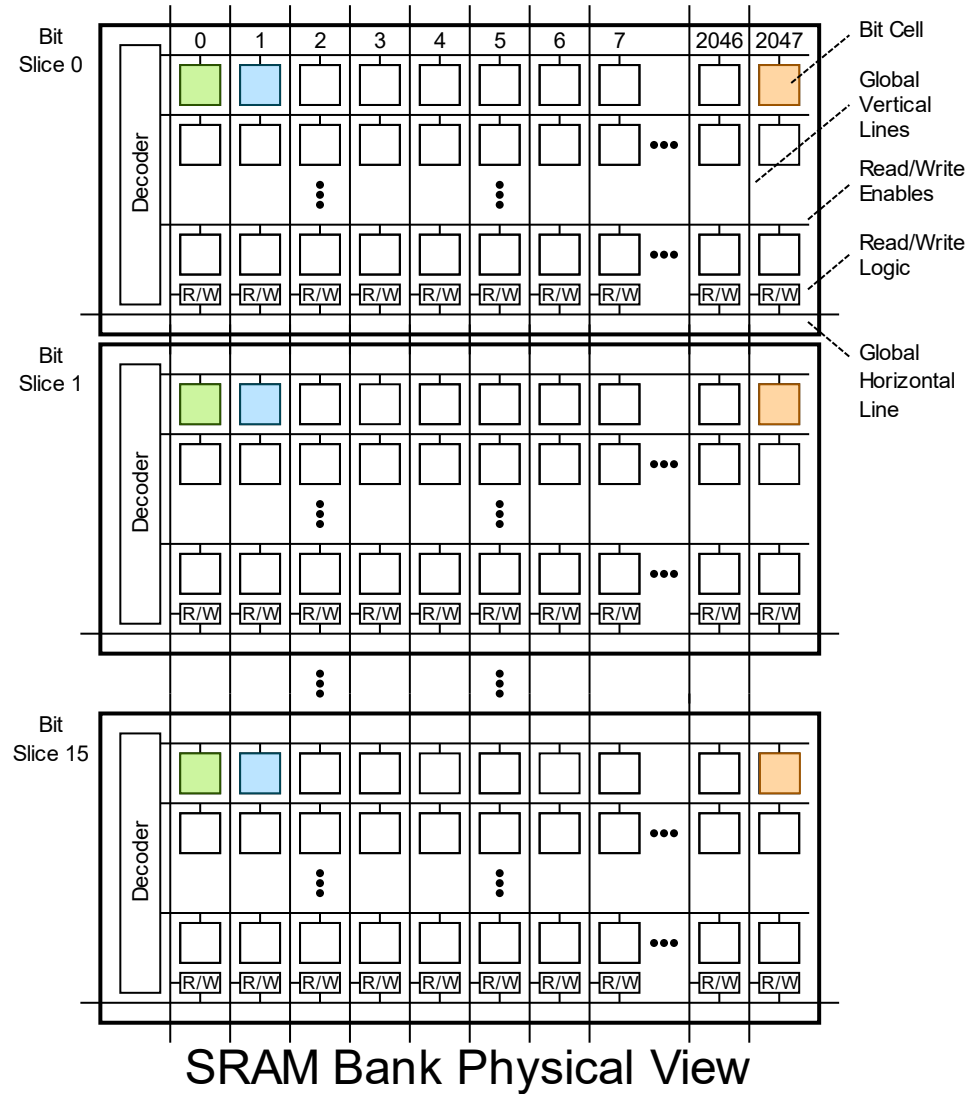
EXAMPLE #3: SEARCH

APL_FRAG search(vdst, vsrc, key):

```
key:      RL = VRF[vsrc];
~key:     RL = ~VRF[vsrc];
0xFFFF:  GVL = RL;
0xFFFF:  VRF[vdst] = GVL;
```

Key: 1 0 0 1 0 0 0 1 0 1 0 1 1 0 1 1

Element: 1 0 0 0 0 1 1 0 1 0 0 1 0 1 1 0



EXAMPLE #3: SEARCH

APL_FRAG search(vdst, vsrc, key):

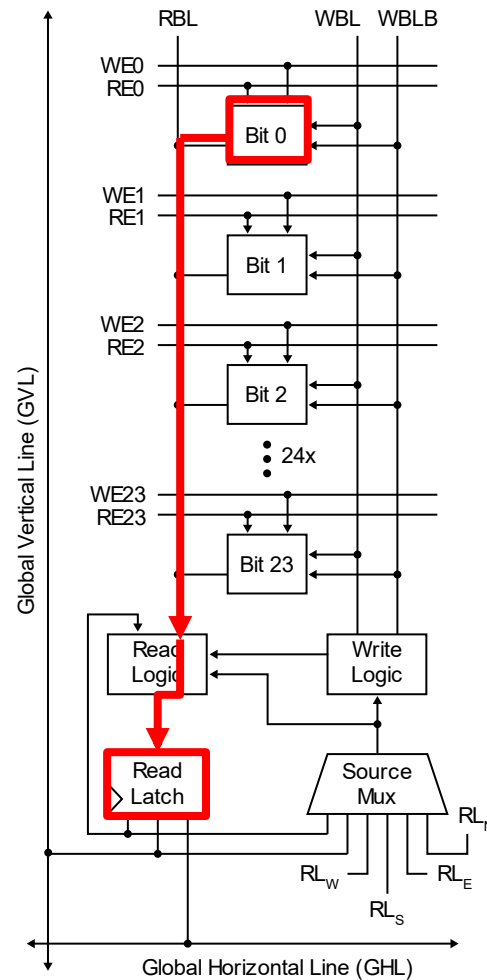
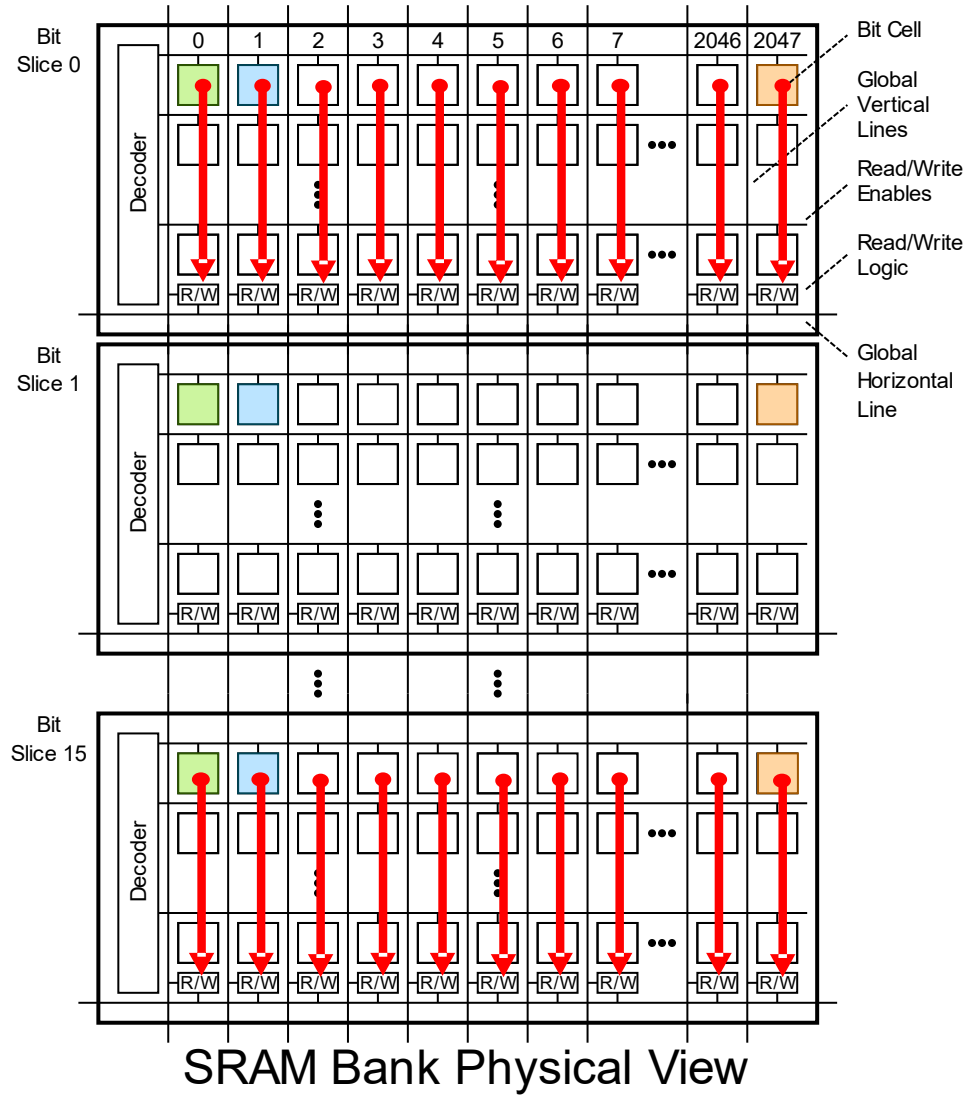
```

key:      RL = VRF[vsrc];
~key:    RL = ~VRF[vsrc];
0xFFFF:  GVL = RL;
0xFFFF:  VRF[vdst] = GVL;
    
```

Key: 1 0 0 1 0 0 0 1 0 1 0 1 1 0 1 1

Element: 1 0 0 0 0 1 1 0 1 0 0 1 0 1 1 0

Result: 1 0 0 0 1 0 1 0



EXAMPLE #3: SEARCH

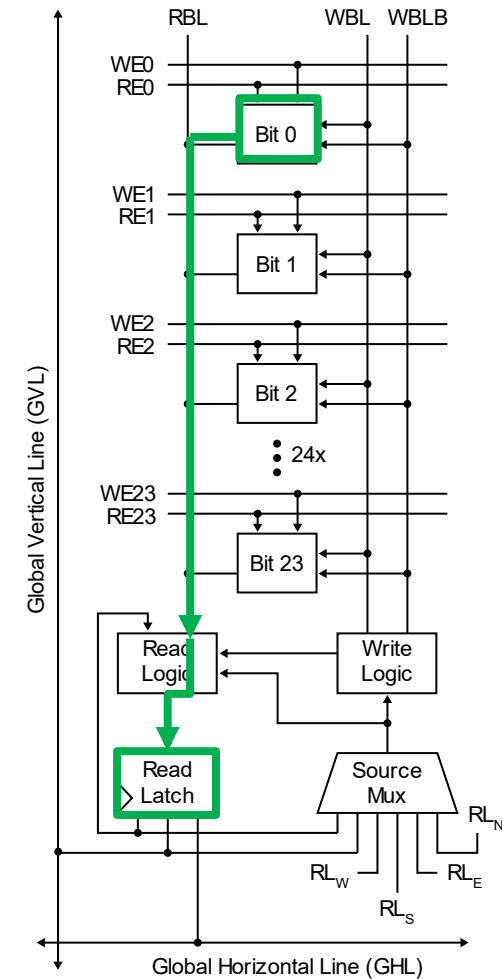
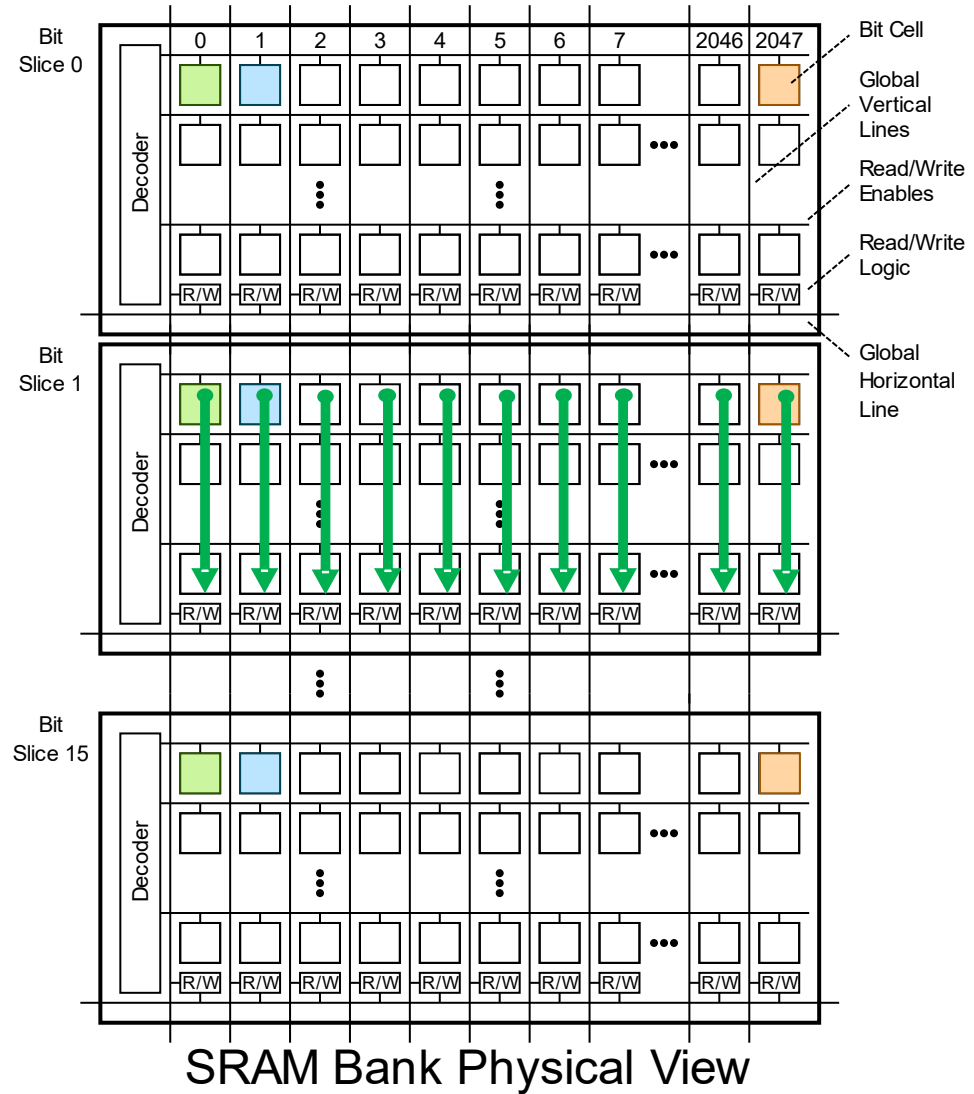
APL_FRAG search(vdst, vsrc, key):

```
key:      RL = VRF[vsrc];
~key:     RL = ~VRF[vsrc];
0xFFFF:  GVL = RL;
0xFFFF:  VRF[vdst] = GVL;
```

Key: 1 0 0 1 0 0 0 1 0 1 0 1 1 0 1 1

Element: 1 0 0 0 0 1 1 0 1 0 0 1 0 1 1 0

Result: 1 1 1 0 1 0 0 0 0 0 1 1 0 0 1 0



EXAMPLE #3: SEARCH

APL_FRAG search(vdst, vsrc, key):

```

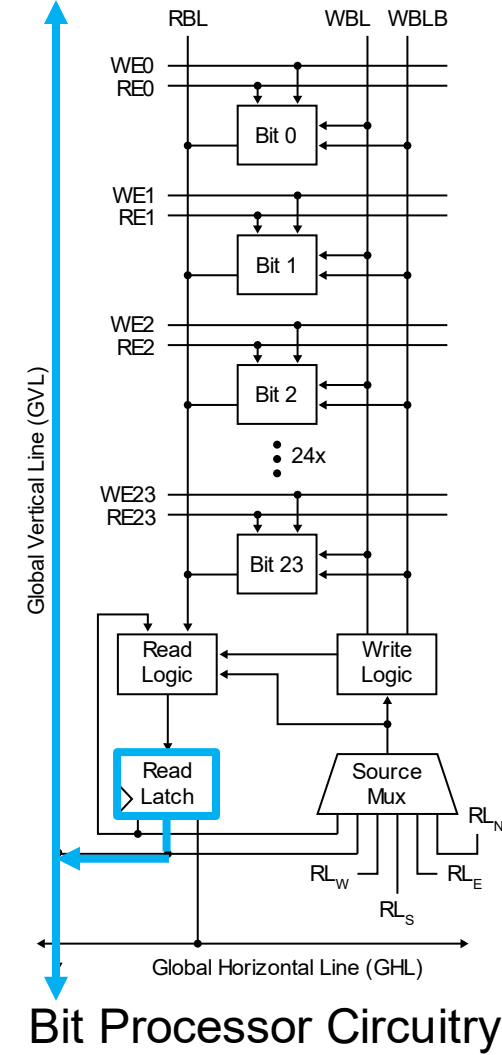
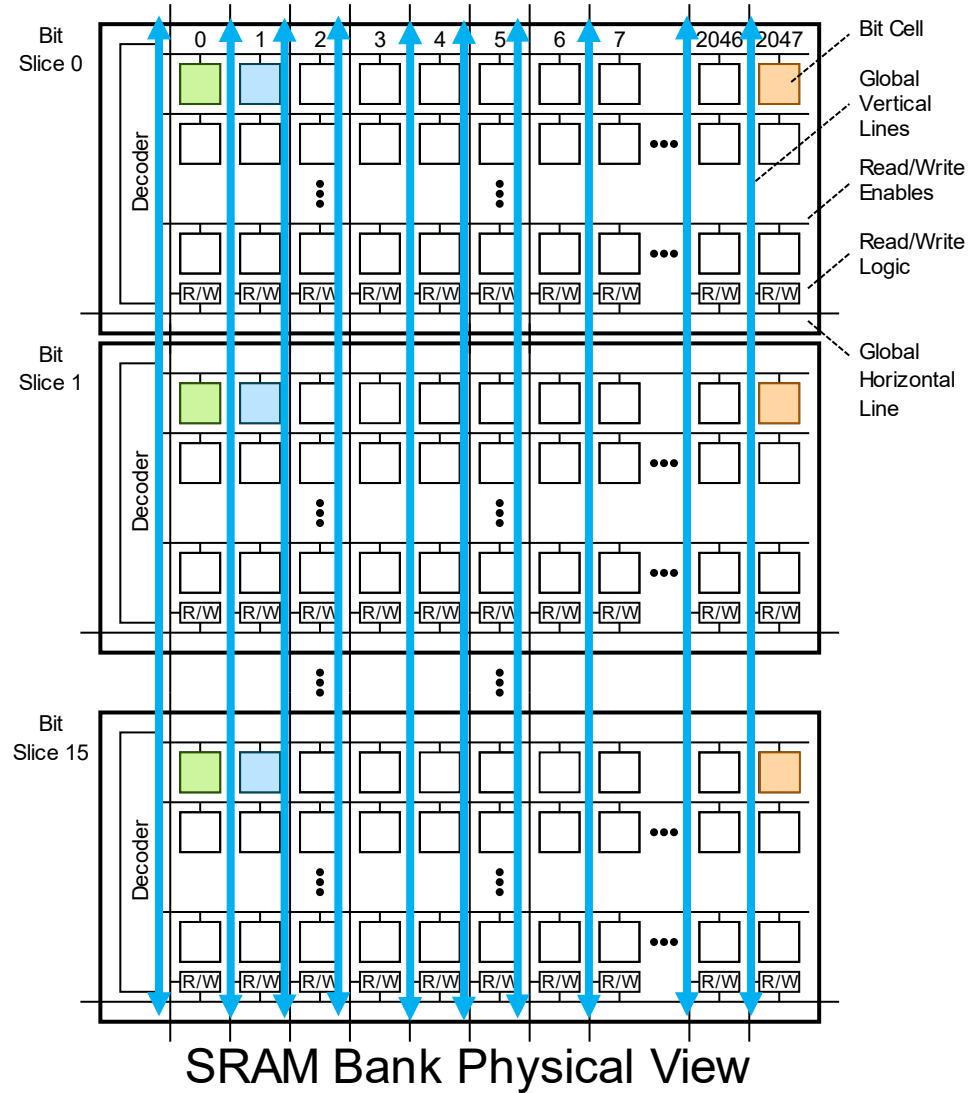
key:      RL = VRF[vsrc];
~key:     RL = ~VRF[vsrc];
0xFFFF:  GVL = RL;
0xFFFF:  VRF[vdst] = GVL;
    
```

Key: 1 0 0 1 0 0 0 1 0 1 0 1 1 0 1 1

Element: 1 0 0 0 0 1 1 0 1 0 0 1 0 1 1 0

Result: 1 1 1 0 1 0 0 0 0 0 1 1 0 0 1 0

AND



EXAMPLE #3: SEARCH

APL_FRAG search(vdst, vsrc, key):

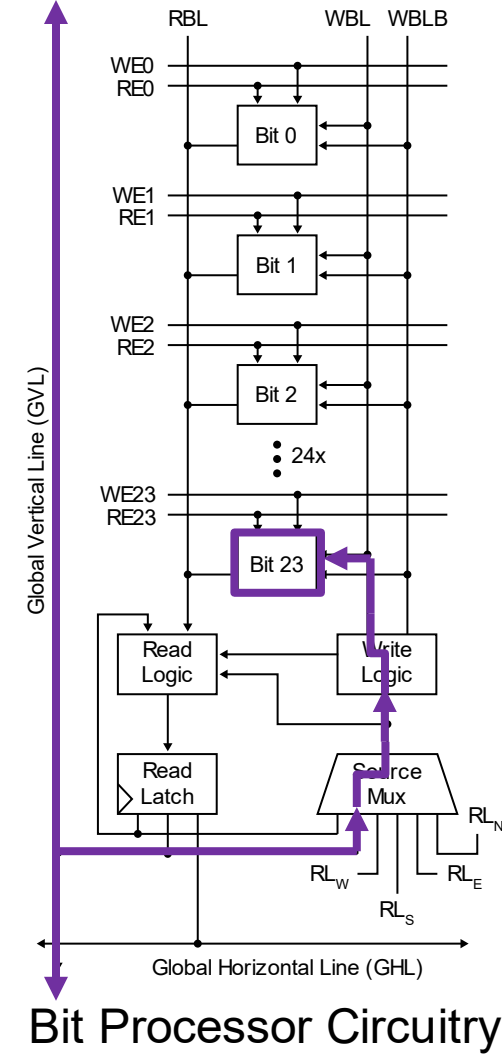
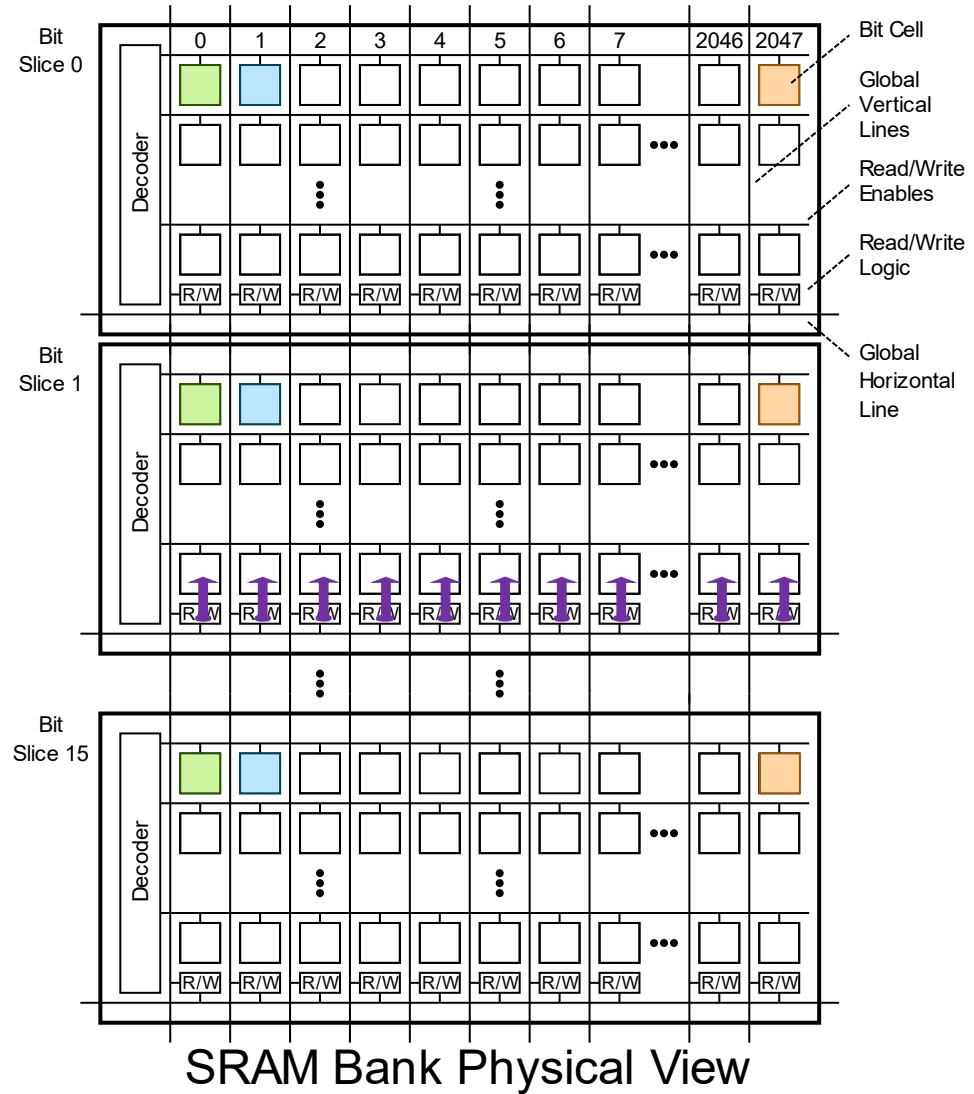
```
key:      RL = VRF[vsrc];
~key:     RL = ~VRF[vsrc];
0xFFFF:  GVL = RL;
0xFFFF:  VRF[vdst] = GVL;
```

Key: 1 0 0 1 0 0 0 1 0 1 0 1 1 0 1 1

Element: 1 0 0 0 0 1 1 0 1 0 0 1 0 1 1 0

Result: 1 1 1 0 1 0 0 0 0 0 1 1 0 0 1 0

AND = 0



EXAMPLE #3: SEARCH

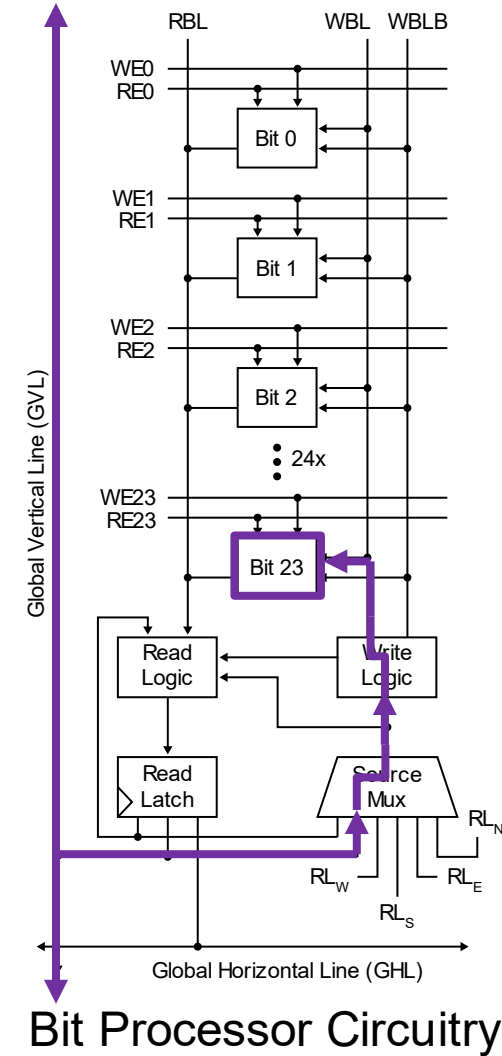
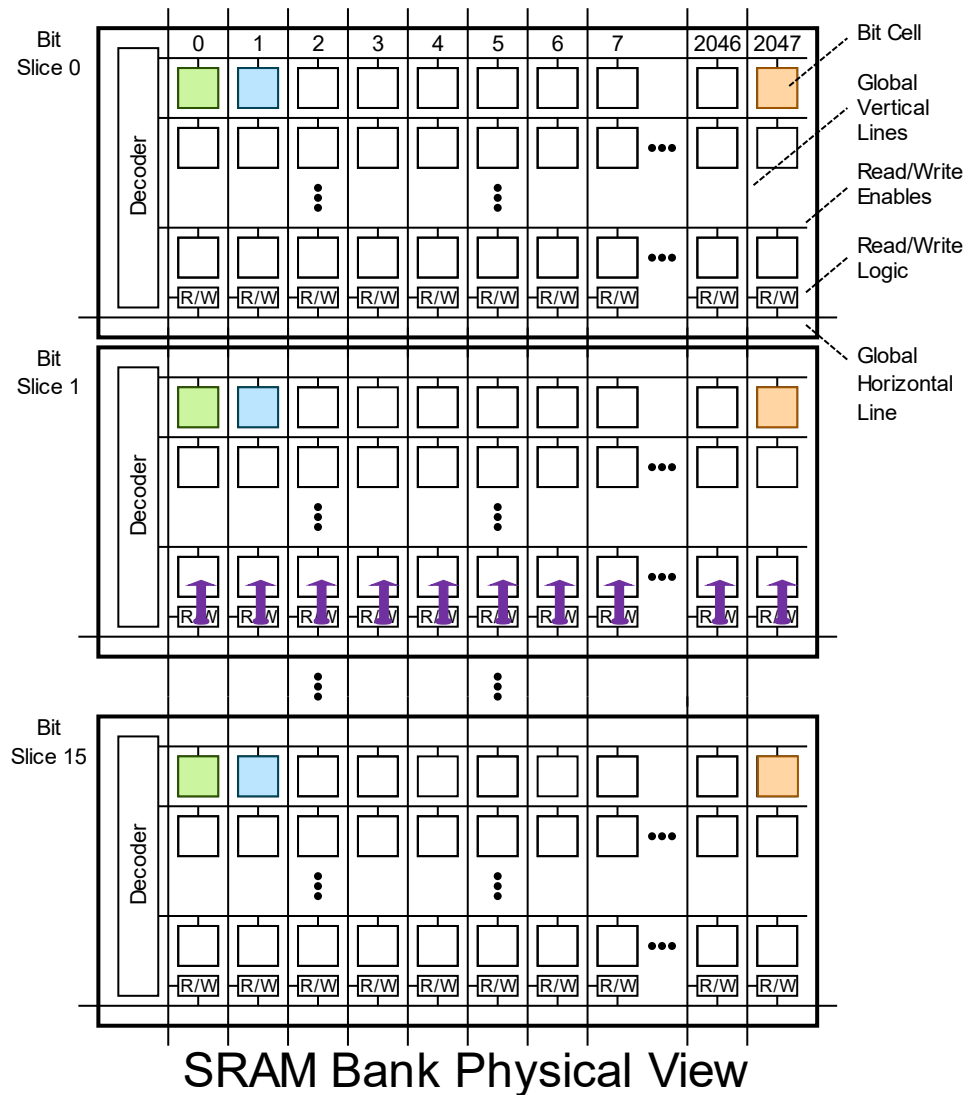
```
APL_FRAG search(vdst, vsrc, key):
```

```
key:      RL = VRF[vsrc];
~key:    RL = ~VRF[vsrc];
0xFFFF:  GVL = RL;
0xFFFF:  VRF[vdst] = GVL;
```

```
Key:      1 0 0 1 0 0 0 1 0 1 0 1 1 0 1 1
Element:  1 0 0 0 0 1 1 0 1 0 0 1 0 1 1 0
Result:   1 1 1 0 1 0 0 0 0 0 1 1 0 0 1 0
```

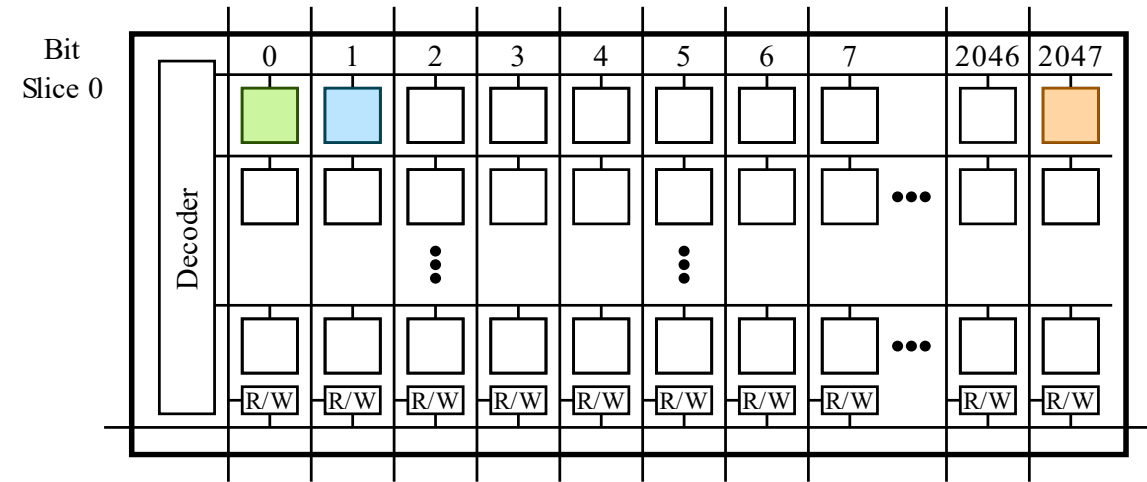
AND = 0

- Architecture doesn't directly support $\sim VRF[]$ operation
- More compact mask encodings are possible



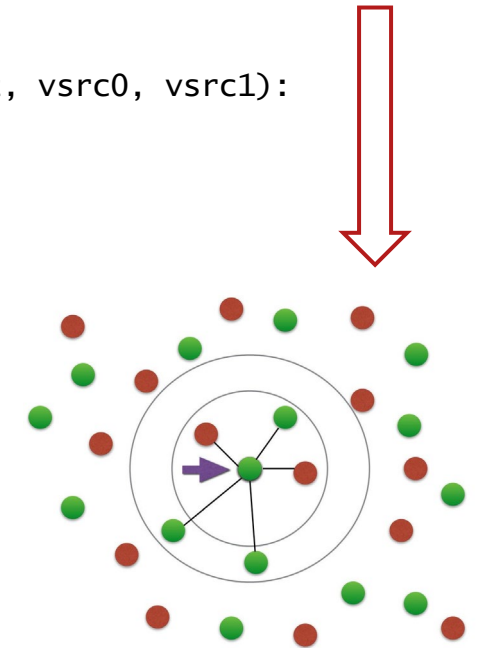
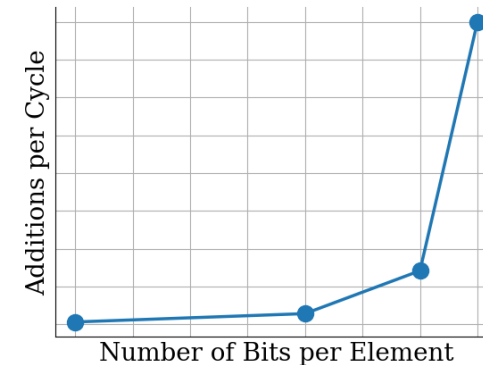
Supporting a Virtual Vector Instruction Set on a Commercial Compute-in-SRAM Accelerator

- Motivation
- APU Microarchitecture
- APU Microcoding
- **Virtual Vector Instruction Set**
- Microbenchmarking Results



```

APL_FRAG _frag_bitwise_or(vdst, vsrc0, vsrc1):
  0xFFFF: RL = VRF[vsrc0];
  0xFFFF: RL |= VRF[vsrc1];
  0xFFFF: VRF[vdst] = RL;
  
```



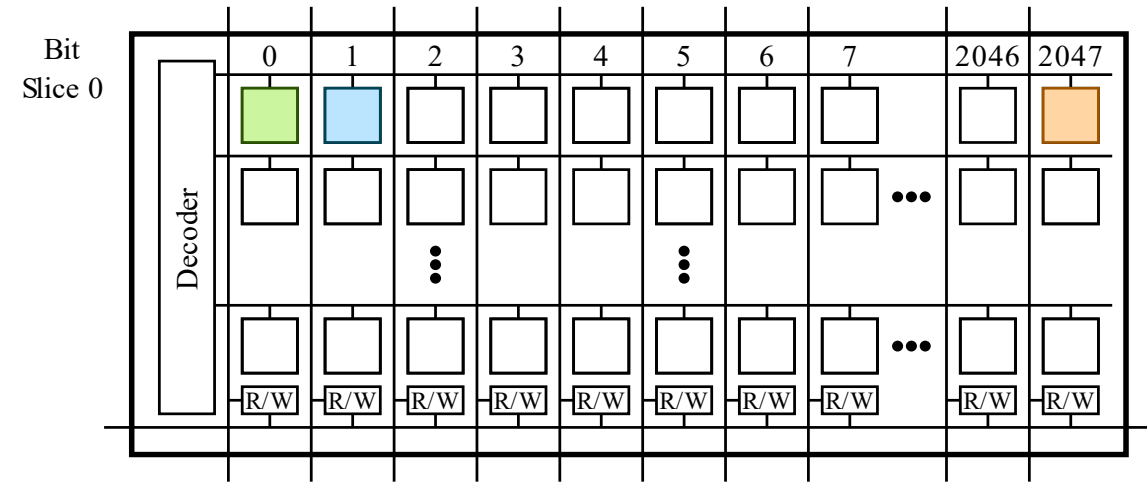
Instruction	Description	Operation Type			
		Element-Wise	Cross-Element	Mask	Memory Configuration
vor_vv	bit-wise or	✓			
vadd_vv	element-wise addition	✓			
vsub_vx	vector-scalar subtraction	✓			
vmul_vv	element-wise multiplication	✓			
vsll	bitwise left-shift logical	✓			
vsrl	bitwise right-shift logical	✓			
vmv_vv	copy vector	✓			
vmseq_vx	set-if-equal		✓		
vmseq_vx_m	set-if-equal (masked)		✓	✓	
vredmin_vs	find min element		✓		
vredmin_vs_m	find min element (masked)		✓	✓	
vmv_vx	broadcast scalar to vector		✓		
vmv_xs	extract element 0 of vector		✓		
vfirrst	index of first set mask bit		✓	✓	
vmsof	set only first set mask bit		✓	✓	
vmset	set all elements of mask		✓	✓	
vmnot	invert mask	✓		✓	
vmand	bitwise-and of two masks	✓		✓	
vmv_m	copy mask	✓		✓	
vload	loads vector from DRAM				✓
vstore	stores vector to DRAM				✓
vsetvl	set application vector length				✓

Limitations:

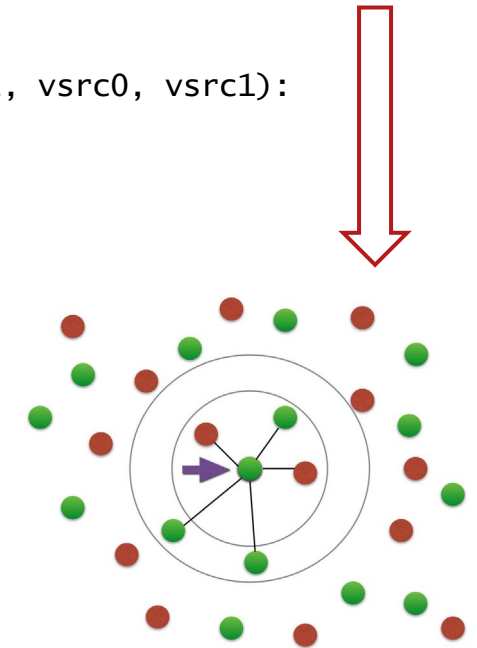
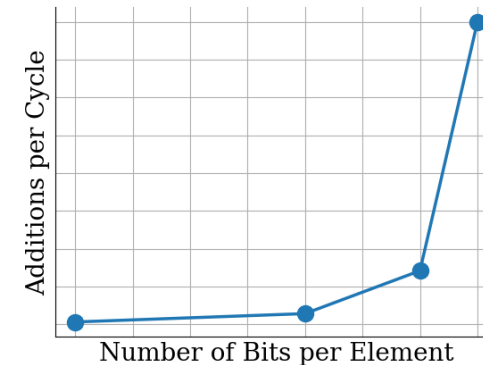
- Unsigned 16-bit integers
- No floating-point support
- Vector lengths up to 32,768 elements
- Assumes source operands are distinct
- Behavior on elements past the specified vector length is undefined

Supporting a Virtual Vector Instruction Set on a Commercial Compute-in-SRAM Accelerator

- Motivation
- APU Microarchitecture
- APU Microcoding
- Virtual Vector Instruction Set
- **Microbenchmarking Results**



```
APL_FRAG _frag_bitwise_or(vdst, vsrc0, vsrc1):  
0xFFFF: RL = VRF[vsrc0];  
0xFFFF: RL |= VRF[vsrc1];  
0xFFFF: VRF[vdst] = RL;
```



MICROBENCHMARKING NARROW-BITWIDTH VECTOR OPERATIONS

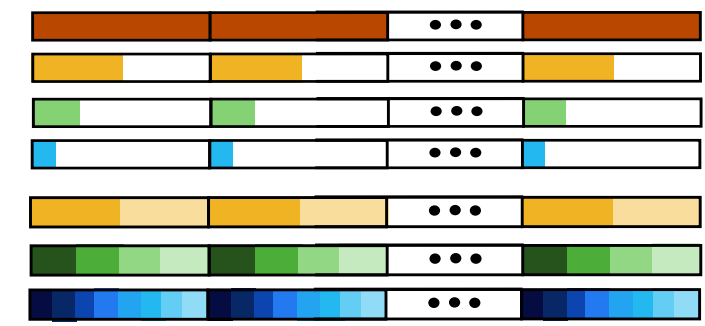


```

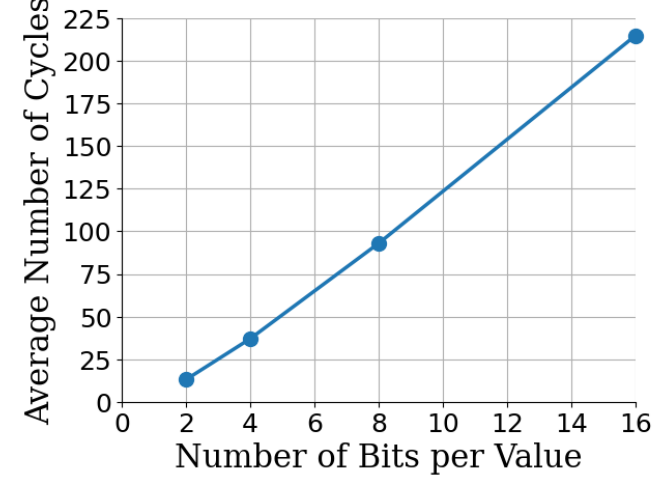
APL_FRAG _frag_vadd(vdst, vsrc0, vsrc1):
// ---- bit 0 ----
0x0001:    RL = VRF[vsrc0];
0x0001:    RL ^= VRF[vsrc1];
0x0001:    VRF[vdst] = RL;
0x0001:    RL = VRF[vsrc0, vsrc1];
// ---- bit 1 ----
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
...
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
// ---- bit 2 ----
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
...
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
// ---- bit 3 ----
(0x0001<<1): RL = VRF[vsrc0];
(0x0001<<1): RL ^= VRF[vsrc1];
...
(0x0001<<1): RL |= VRF[temp_0];
(0x0001<<1): RL |= VRF[temp_1];
...
    
```

vadd

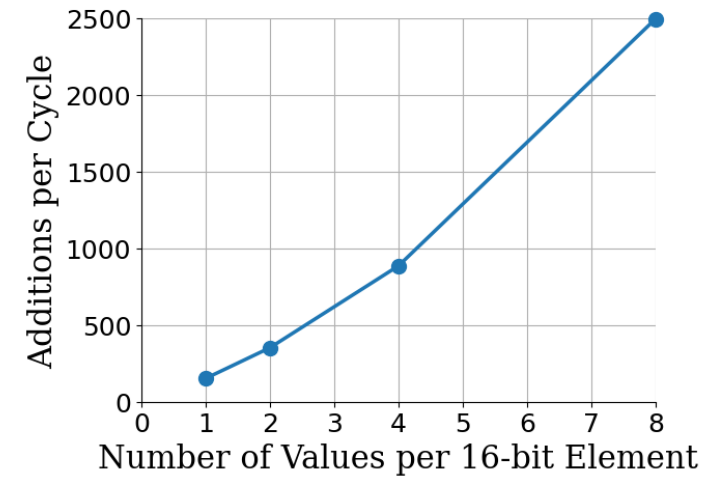
	packed operation	avg # cycles	ops/cycle
one	16-bit ADD	215	152
one	8-bit ADD	93	352
one	4-bit ADD	37	886
one	2-bit ADD	13	2497
two	8-bit ADDs	93	705
four	4-bit ADDs	37	3542
eight	2-bit ADDs	13	20010



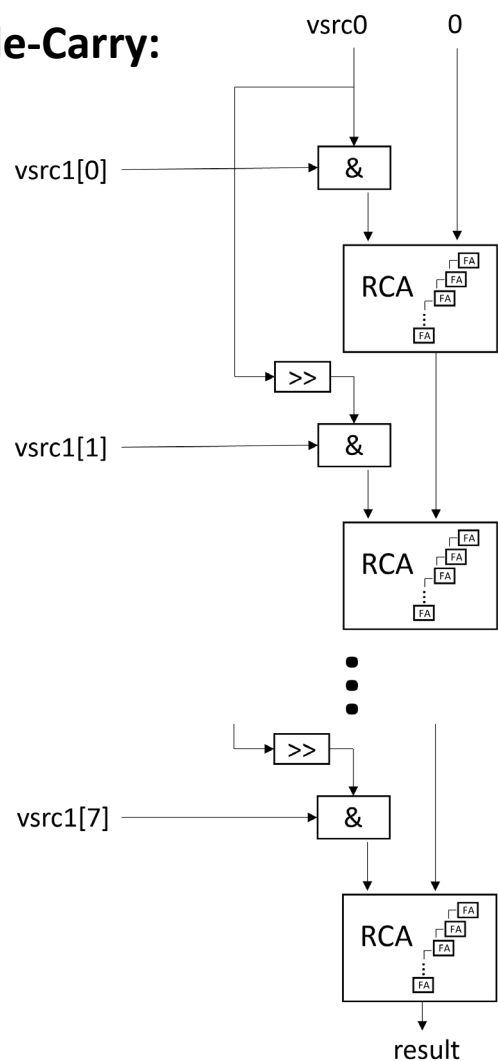
Cycle Counts for Single Vector ADD Operations



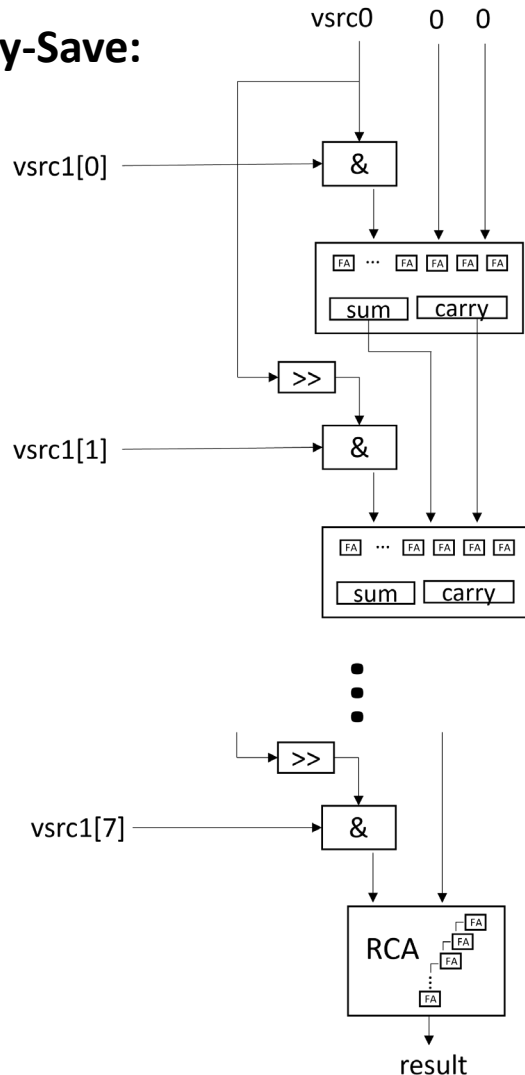
Throughput of Packed Vector ADD Operations



Ripple-Carry:



Carry-Save:



For each bit:

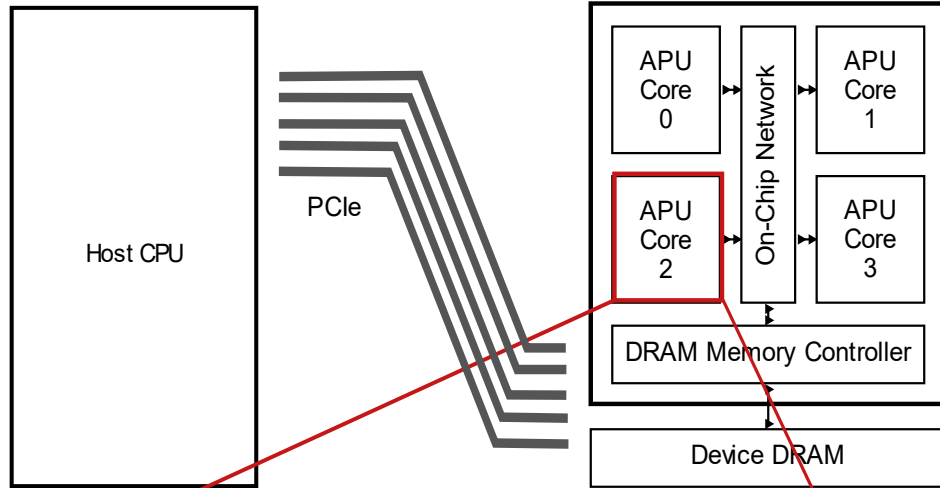
$$S_{out} = (X * Y_i) \wedge S_{in} \wedge C_{in}$$

$$C_{out} = (X * Y_i) * C_{in}$$

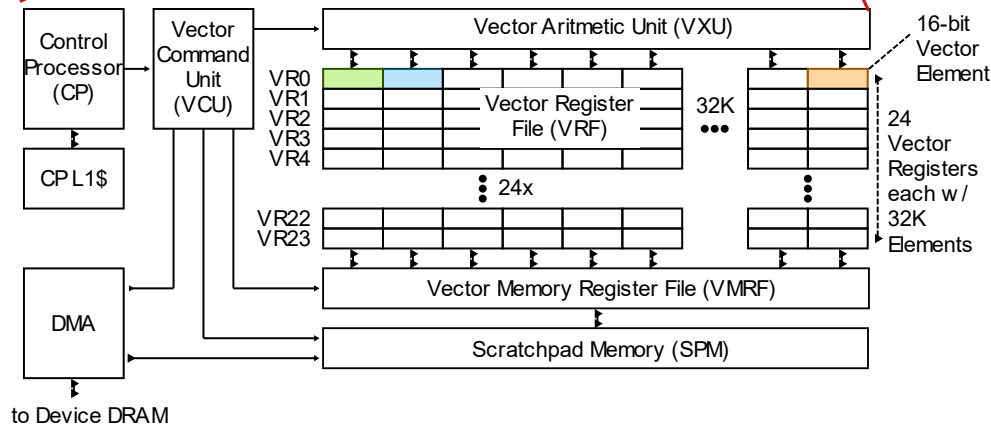
$$+ S_{in} * C_{in}$$

$$+ X * Y_i * S_{in}$$

Ripple-carry: 3540 cycles } 8.3x speedup
 Carry-save: 438 cycles }



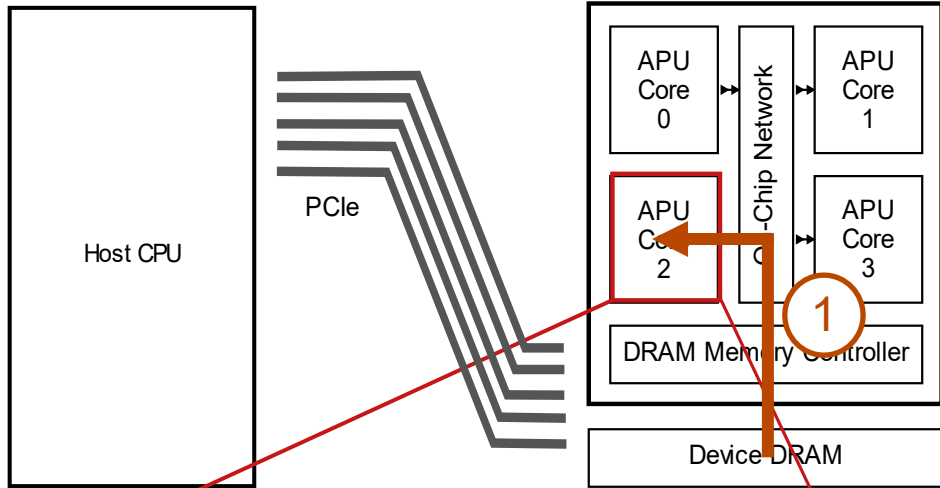
System Overview



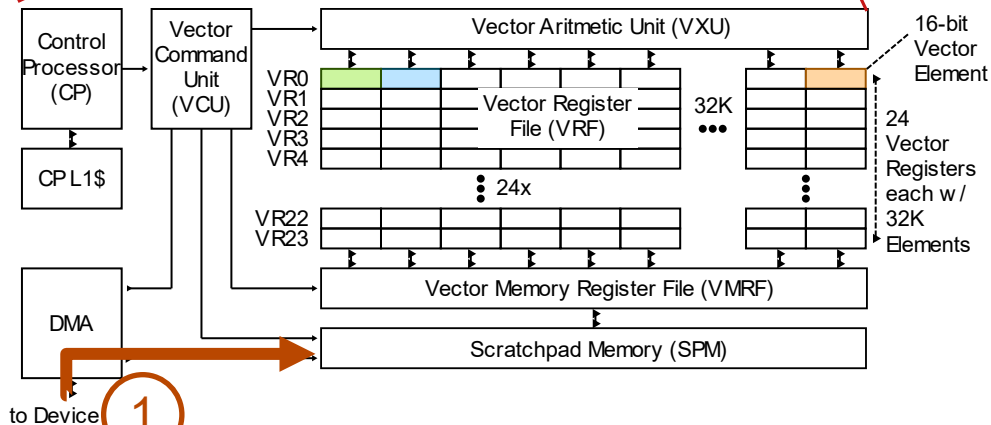
APU Core Logical View

Data Transfer	avg # cycles	elements/cycle
---------------	--------------	----------------

MICROBENCHMARKING DATA MOVEMENT AND FRAGMENT OVERHEAD



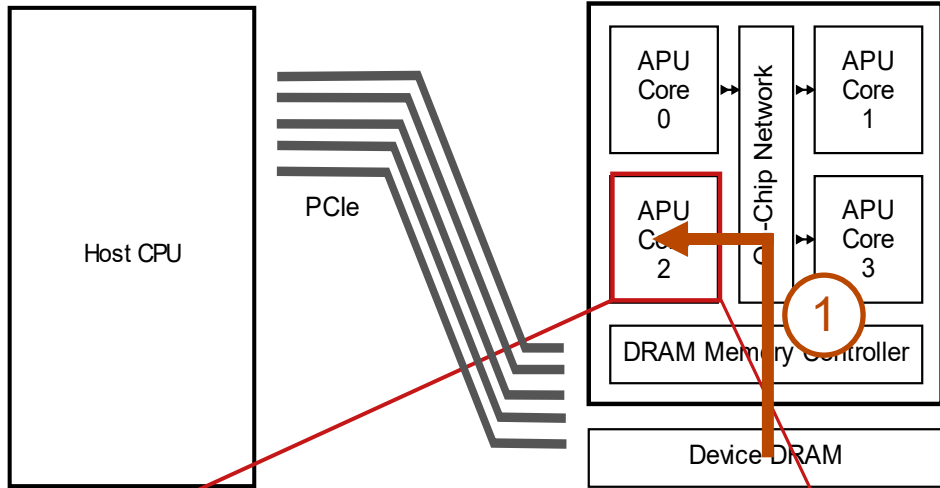
System Overview



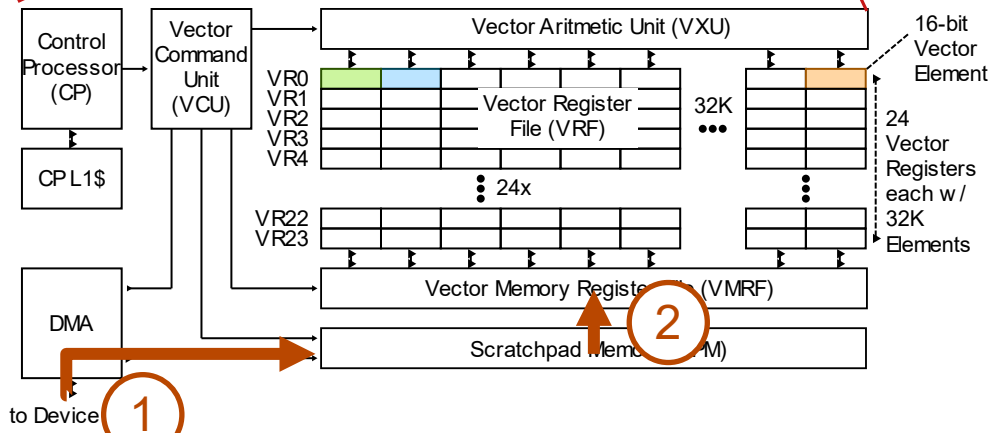
APU Core Logical View

Data Transfer	avg # cycles	elements/cycle
---------------	--------------	----------------

MICROBENCHMARKING DATA MOVEMENT AND FRAGMENT OVERHEAD



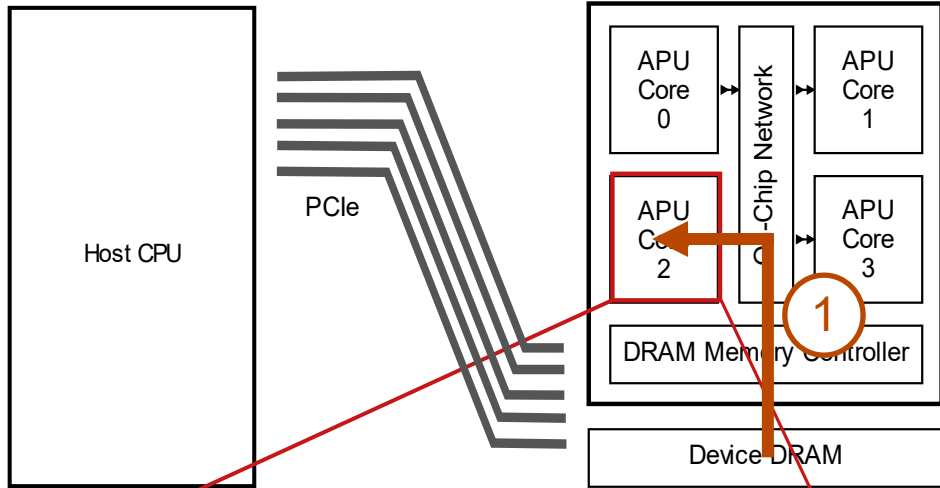
System Overview



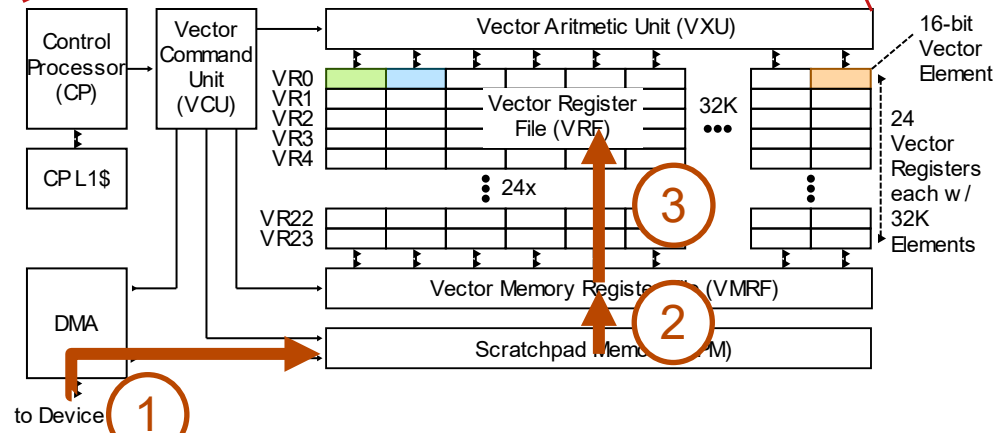
APU Core Logical View

Data Transfer	avg # cycles	elements/cycle
---------------	--------------	----------------

MICROBENCHMARKING DATA MOVEMENT AND FRAGMENT OVERHEAD



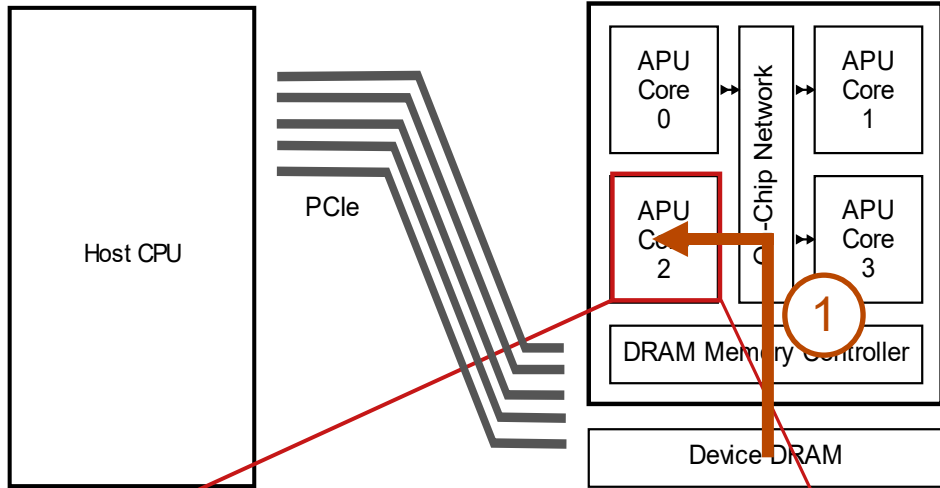
System Overview



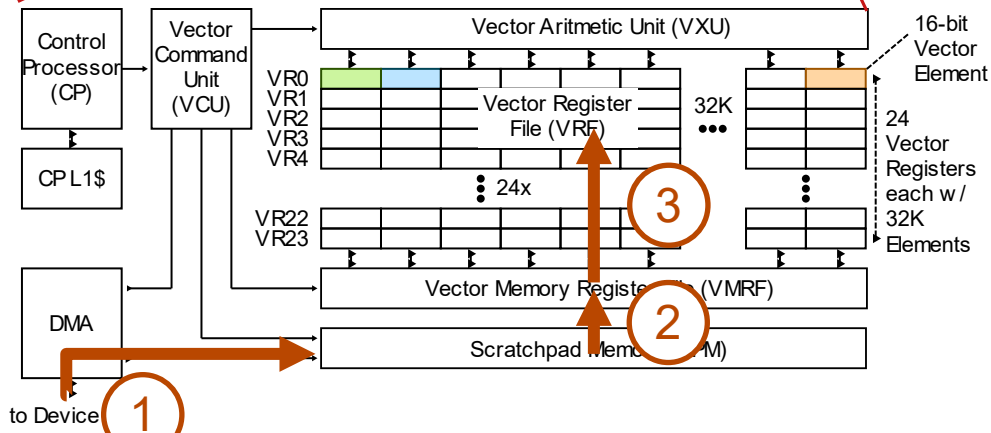
APU Core Logical View

Data Transfer	avg # cycles	elements/cycle
---------------	--------------	----------------

MICROBENCHMARKING DATA MOVEMENT AND FRAGMENT OVERHEAD



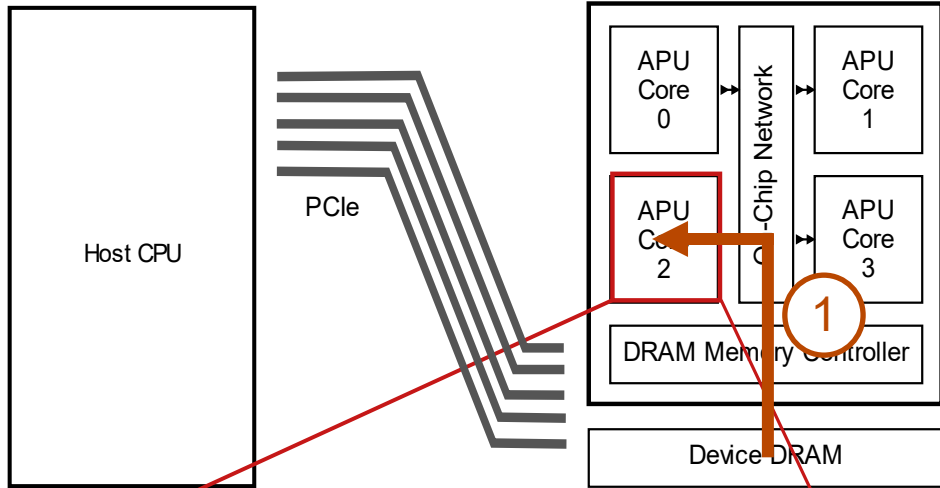
System Overview



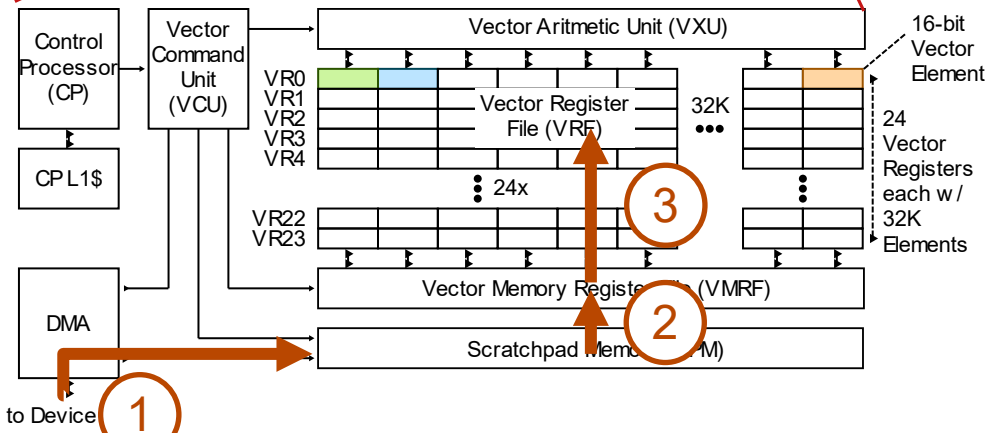
APU Core Logical View

Data Transfer	avg # cycles	elements/cycle
Total	22,066	1.4

MICROBENCHMARKING DATA MOVEMENT AND FRAGMENT OVERHEAD



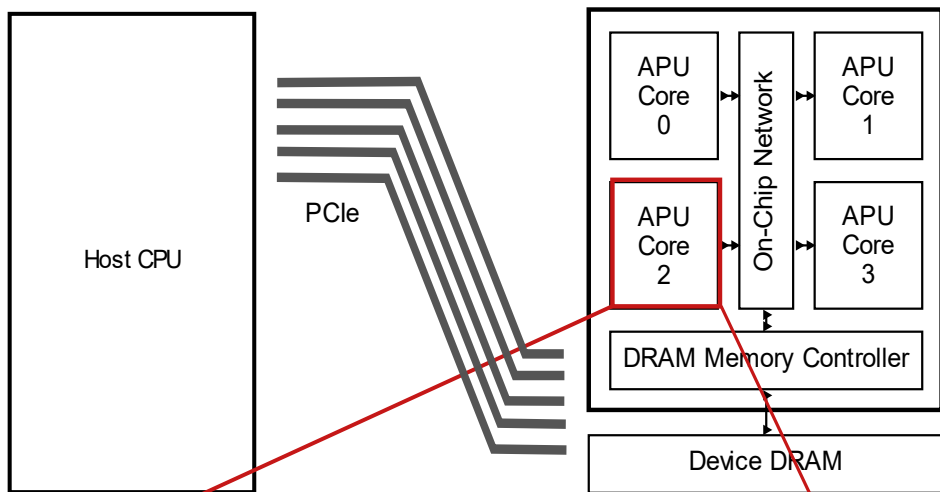
System Overview



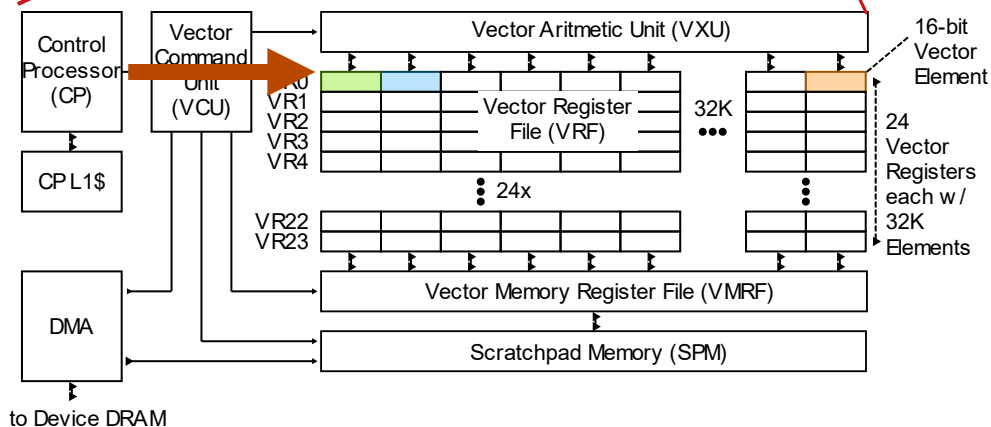
APU Core Logical View

Data Transfer	avg # cycles	elements/cycle
DRAM → SPM	21,575	1.5
SPM → VMRF	470	69.7
VMRF → VRF	21	1568.9
Total	22,066	1.4

MICROBENCHMARKING DATA MOVEMENT AND FRAGMENT OVERHEAD



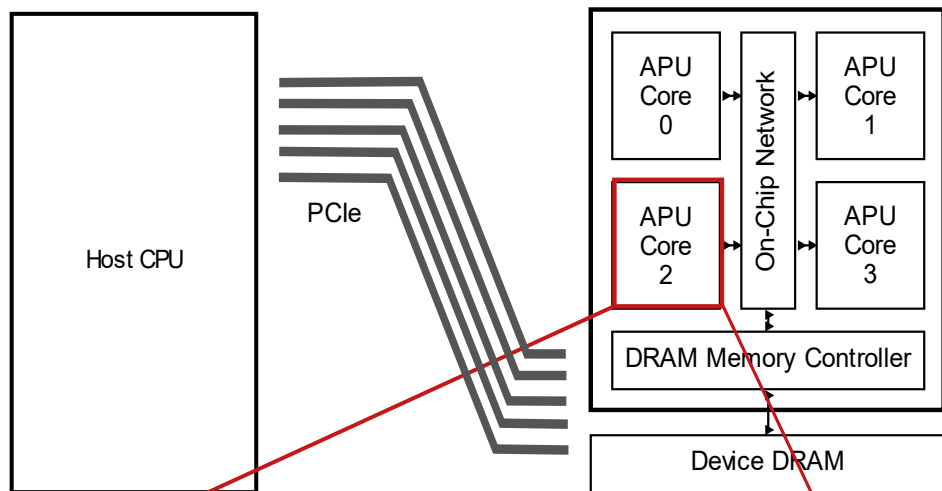
System Overview



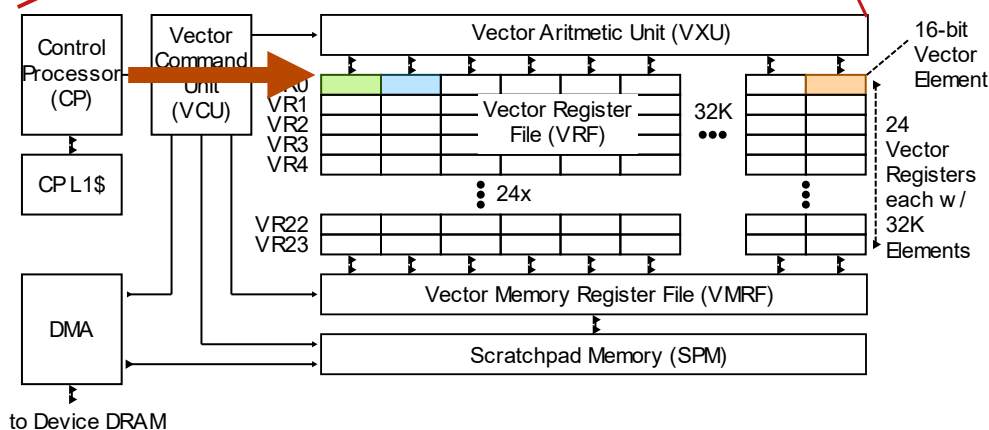
APU Core Logical View

Data Transfer	avg # cycles	elements/cycle
DRAM → SPM	21,575	1.5
SPM → VMRF	470	69.7
VMRF → VRF	21	1568.9
Total	22,066	1.4

MICROBENCHMARKING DATA MOVEMENT AND FRAGMENT OVERHEAD



System Overview

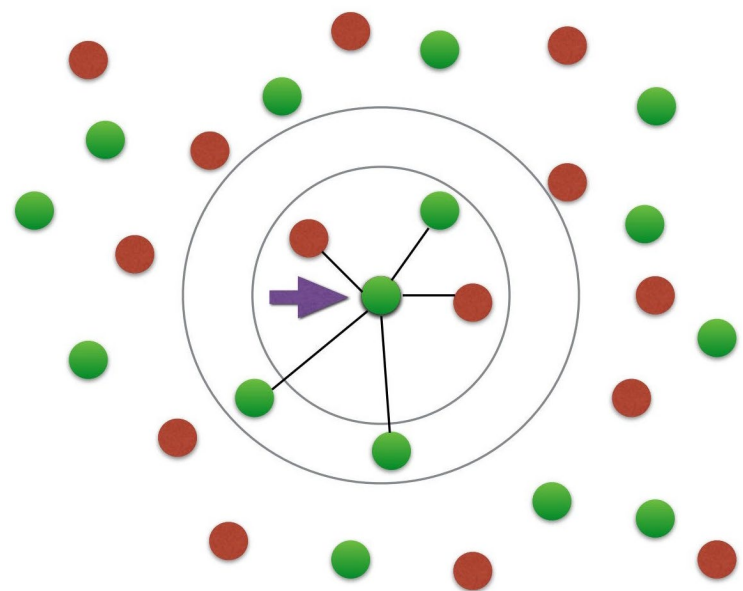


APU Core Logical View

Data Transfer	avg # cycles	elements/cycle
DRAM → SPM	21,575	1.5
SPM → VMRF	470	69.7
VMRF → VRF	21	1568.9
Total	22,066	1.4

Instruction	Description	Execution Time (cycles)	
		Theoretical (b=16)	Measured
<code>vor_vv</code>	bit-wise or	6	15
<code>vmseq_vx</code>	set-if-equal	7	13
<code>vmseq_vx_m</code>	set-if-equal (masked)	9	15
<code>vmset</code>	set all elements of mask	3	10
<code>vmnot</code>	invert mask	7	12
<code>vmand</code>	bitwise-and of two masks	9	13
<code>vmv_vx</code>	broadcast scalar to vector	5	12
<code>vmv_vv</code>	copy vector	4	11
<code>vmv_m</code>	copy mask	4	11
<code>vsll</code>	bitwise left-shift logical	3	13
<code>vsrl</code>	bitwise right-shift logical	3	13

4-10 cycles



ANN_SIFT10 dataset, 32K points

89 searches/second

Fusing operations to skip unnecessary computation

97 searches/second

Fusing operations to re-use intermediates

112 searches/second

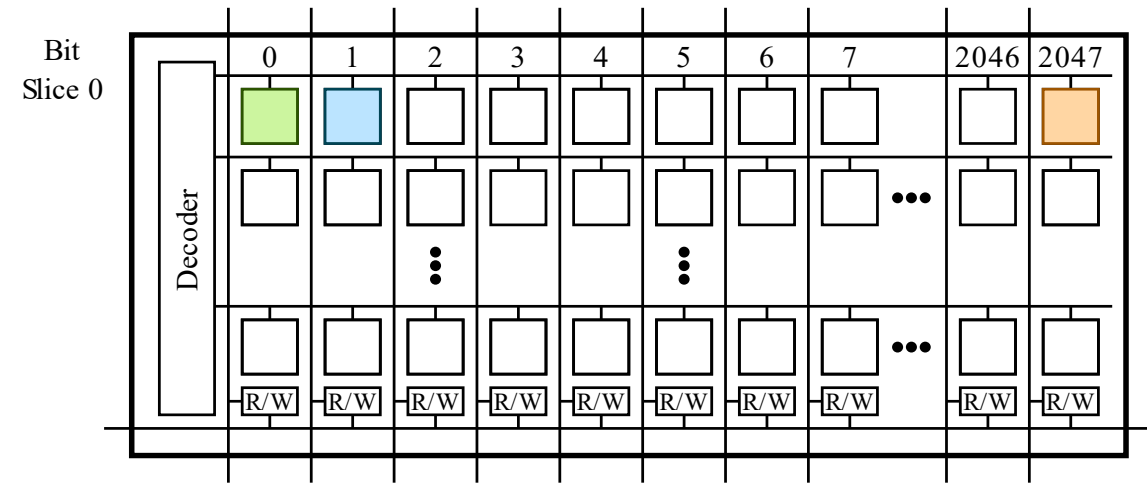
Asynchronous partial load from DRAM

129 searches/second

**1.4x speedup, but
68% more LOC**

Supporting a Virtual Vector Instruction Set on a Commercial Compute-in-SRAM Accelerator

1. APU is an interesting platform for applications with massive parallelism, bitwise operations, narrow bitwidths, and high data reuse
2. Proof-of-concept support for a RISC-V-like virtual vector instruction set
3. First step toward general-purpose computing on commercial domain-specific compute-in-SRAM accelerators



```

APL_FRAG _frag_bitwise_or(vdst, vsrc0, vsrc1):
  0xFFFF: RL = VRF[vsrc0];
  0xFFFF: RL |= VRF[vsrc1];
  0xFFFF: VRF[vdst] = RL;
  
```

