

# ASCELLA: Accelerating Sparse Computation by Enabling Stream Accesses to Memory

Bahar Asgari, Ramyad Hadidi, Hyesoon Kim



comparch

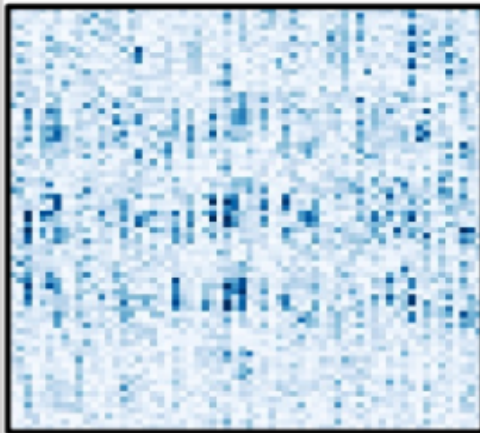


# Sparse matrices are everywhere!

2

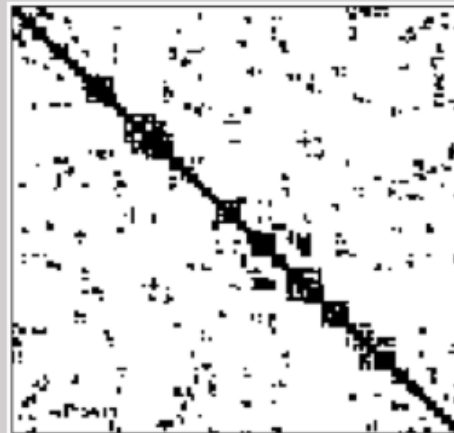
Sparse matrix vector multiplication (SpMV) is a main operation in:

## Neural Networks



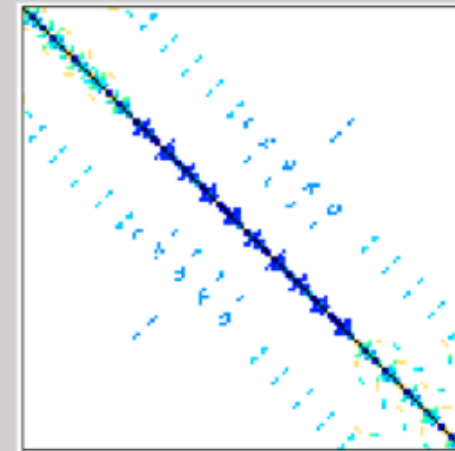
The weight matrix  
is sparse

## Graph Analytics



The adjacency matrix  
is sparse

## Differential Equations

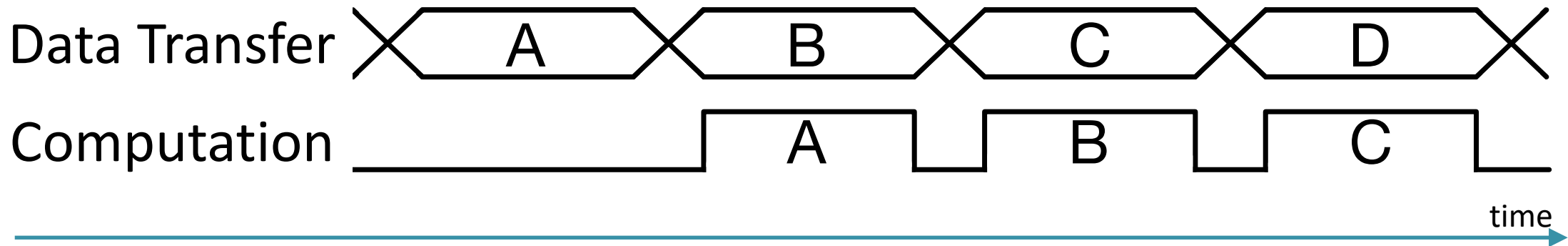


The coefficient matrix  
is sparse



# An ideal hardware accelerator for SpMV

- ▶ SpMV can be accelerated:
  - ▶ By stream data from memory
  - ▶ By using a parallel dot product engine
- ▶ Ideally, we want compute time and data-transfer time for blocks of data (e.g., A, B, C, and D) to be equal:

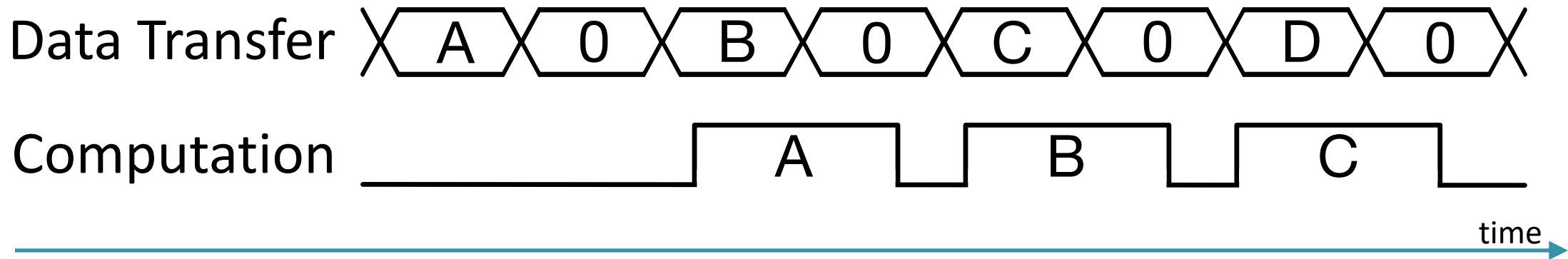


\*a block is the unit of streaming data from memory



# Decompression can cause a bottleneck

- ▶ Sparse matrices are often stored and transferred in compressed forms
  - ▶ Example: compressed sparse row (CSR) and blocked CSR (BCSR)
- ▶ The compressed data must be first decompressed
- ▶ Decompression is slow and causes bottleneck





# Why decompression is slow?

---

5

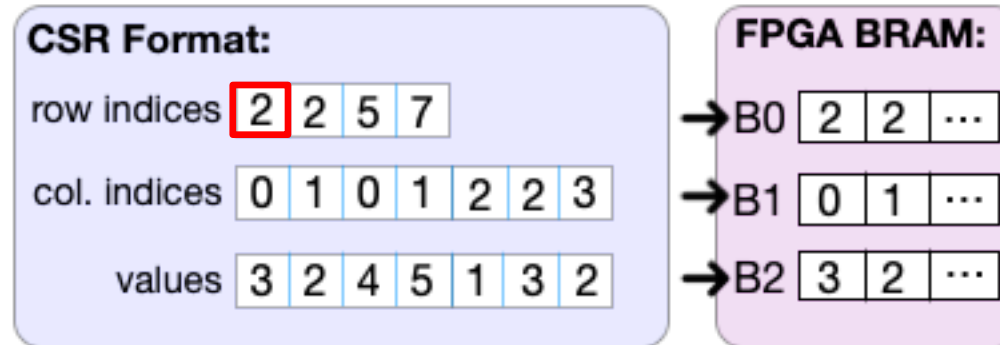
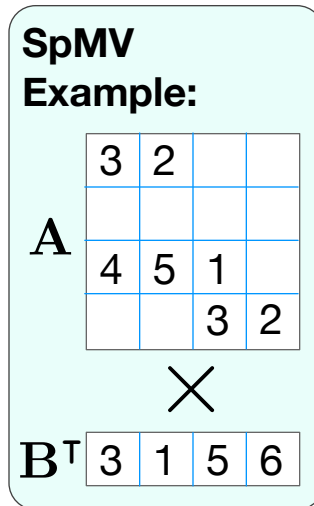
CSR and BCSR use three vectors to represent a sparse matrix

- ▶ Row indices (offsets)
- ▶ Column indices
- ▶ Values

To decompress a non-zero row, we need to

- ▶ First, read one element of row indices
- ▶ Then, read column indices and values as required

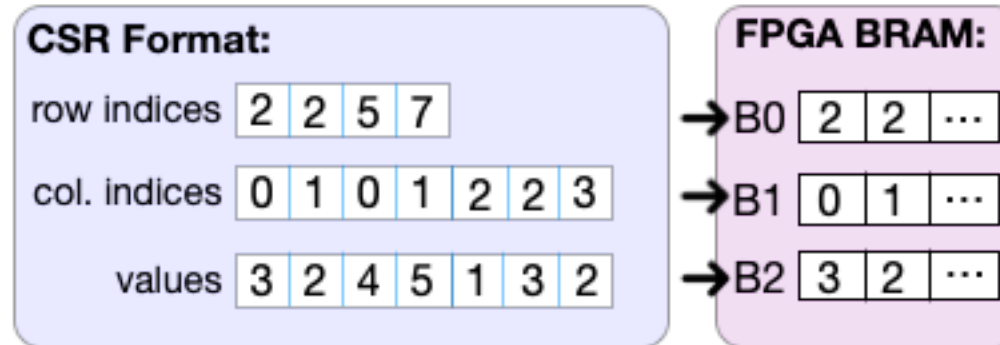
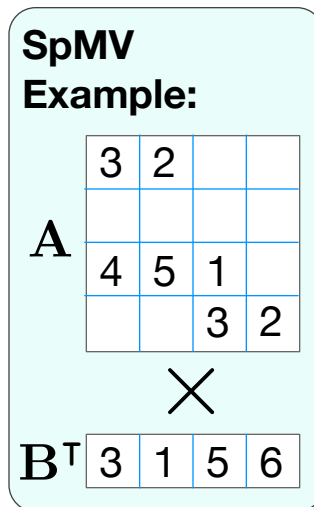
# Decompression from CSR format



**BRAM Accesses Timeline:**

	R Read Operation			C Compute Operation			
<i>cycle: 0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>
<b>R</b> B0 → 2							
<i>cycle: 7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>11</i>	<i>12</i>	<i>13</i>	

# Decompression from CSR format



**BRAM Accesses Timeline:**

	R Read Operation		C Compute Operation			
cycle: 0	1	2	3	4	5	6
R B0 → 2	C 2-0=2					
cycle: 7	8	9	10	11	12	13

# Decompression from CSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

×

**CSR Format:**

row indices 2 2 5 7

col. indices 0 1 0 1 2 2 3

values 3 2 4 5 1 3 2

**FPGA BRAM:**

→ B0 2 2 ...

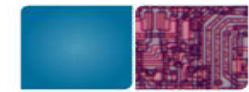
→ B1 0 1 ...

→ B2 3 2 ...

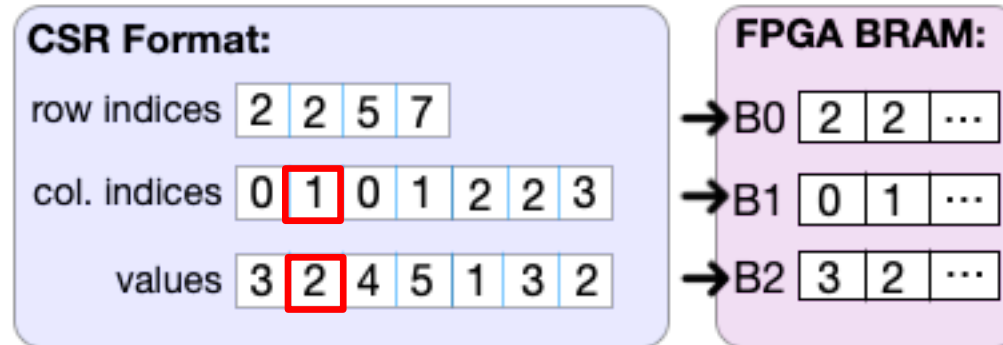
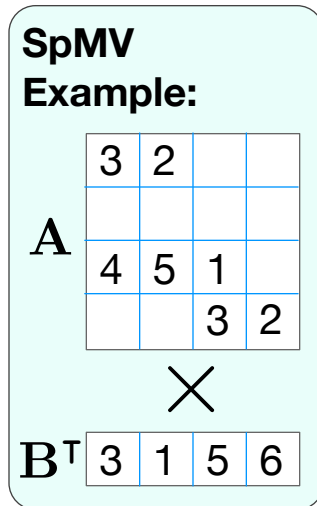
**BRAM Accesses Timeline:**

			<b>R</b> Read Operation		<b>C</b> Compute Operation	
<i>cycle: 0</i>	1	2	3	4	5	6
<b>R</b> B0 → 2	<b>C</b> 2-0=2	<b>R</b> B1 → 0 <b>R</b> B2 → 3				
<i>cycle: 7</i>	8	9	10	11	12	13





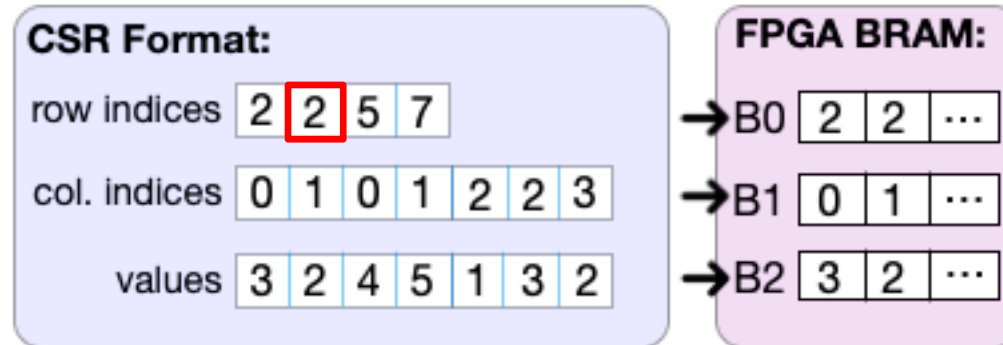
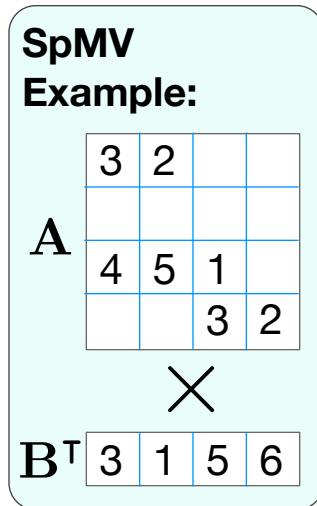
# Decompression from CSR format



**BRAM Accesses Timeline:**

			<b>R</b> Read Operation		<b>C</b> Compute Operation	
<i>cycle: 0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>
<b>R</b> B0 → 2	<b>C</b> 2-0=2	<b>R</b> B1 → 0 <b>R</b> B2 → 3	<b>R</b> B1 → 1 <b>R</b> B2 → 2			
<i>cycle: 7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>11</i>	<i>12</i>	<i>13</i>

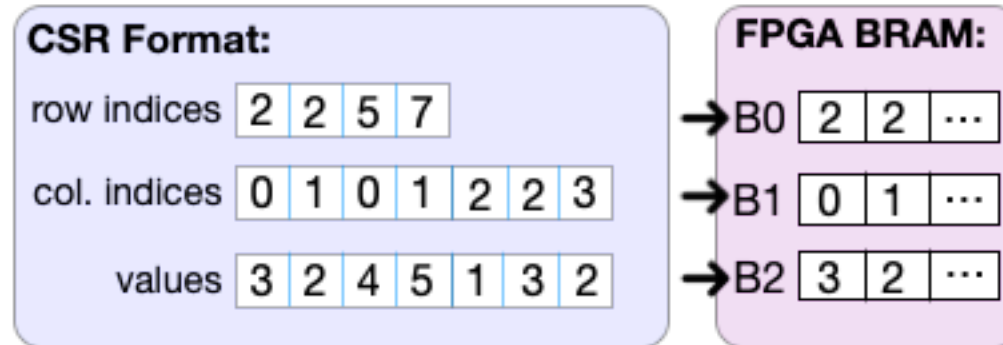
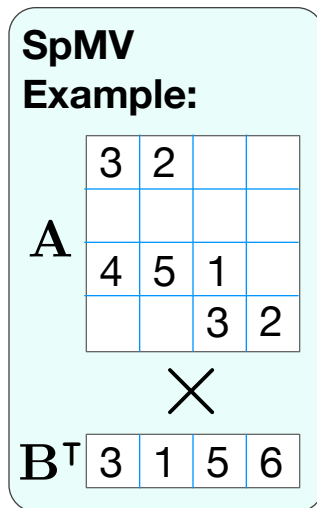
# Decompression from CSR format



**BRAM Accesses Timeline:**

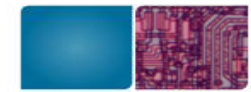
<i>cycle: 0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>
<b>R</b> B0 → 2	<b>C</b> 2-0=2	<b>R</b> B1 → 0 <b>R</b> B2 → 3	<b>R</b> B1 → 1 <b>R</b> B2 → 2	<b>R</b> B0 → 2		
<i>cycle: 7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>11</i>	<i>12</i>	<i>13</i>

# Decompression from CSR format



**BRAM Accesses Timeline:**

		R Read Operation		C Compute Operation		
cycle: 0	1	2	3	4	5	6
R B0 → 2	C 2-0=2	R B1 → 0 R B2 → 3	R B1 → 1 R B2 → 2	R B0 → 2	C 2-2=0	
cycle: 7	8	9	10	11	12	13



# Decompression from CSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

×

**CSR Format:**

row indices 

2	2	5	7
---	---	---	---

col. indices 

0	1	0	1	2	2	3
---	---	---	---	---	---	---

values 

3	2	4	5	1	3	2
---	---	---	---	---	---	---

**FPGA BRAM:**

→ B0 

2	2	...
---	---	-----

→ B1 

0	1	...
---	---	-----

→ B2 

3	2	...
---	---	-----

**BRAM Accesses Timeline:**

				<b>R</b> Read Operation	<b>C</b> Compute Operation	
<i>cycle: 0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>
<b>R</b> B0 → 2	<b>C</b> 2-0=2	<b>R</b> B1 → 0 <b>R</b> B2 → 3	<b>R</b> B1 → 1 <b>R</b> B2 → 2	<b>R</b> B0 → 2	<b>C</b> 2-2=0	<b>R</b> B0 → 5
<i>cycle: 7</i>	<i>8</i>	<i>9</i>	<i>10</i>	<i>11</i>	<i>12</i>	<i>13</i>
<b>C</b> 5-2=3	<b>R</b> B1 → 0 <b>R</b> B2 → 4	<b>R</b> B1 → 1 <b>R</b> B2 → 5	<b>R</b> B1 → 2 <b>R</b> B2 → 1	<b>C</b> 7-5=2	<b>R</b> B1 → 2 <b>R</b> B2 → 3	<b>R</b> B1 → 3 <b>R</b> B2 → 2

# Decompression from BCSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

$\times$

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

**BCSR Format (2x2 sub-blocks) :**

row indices 

1	3
---	---

col. indices 

0	0	2
---	---	---

values 

3	2	0	0
4	5	0	0
1	0	3	2

**FPGA BRAM:**

→ B0 

1	3	
---	---	--

→ B1 

0	0	2
---	---	---

→ B2 

3	4	1
---	---	---

→ B3 

2	5	0
---	---	---

→ B4 

0	0	3
---	---	---

→ B5 

0	0	2
---	---	---

**BRAM Accesses Timeline:**

**R** Read Operation      **C** Compute Operation

cycle: 0	1	2	3	4	5	6

# Decompression from BCSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

×

**BCSR Format (2x2 sub-blocks) :**

row indices **1 3**

col. indices **0 0 2**

values

3	2	0	0
4	5	0	0
1	0	3	2

**FPGA BRAM:**

B0	1	3	
B1	0	0	2
B2	3	4	1
B3	2	5	0
B4	0	0	3
B5	0	0	2

**BRAM Accesses Timeline:**

**R** Read Operation      **C** Compute Operation

cycle: 0	1	2	3	4	5	6
<b>R</b> B0 → 1						

# Decompression from BCSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

×

**BCSR Format (2x2 sub-blocks) :**

row indices 

1	3
---	---

col. indices 

0	0	2
---	---	---

values 

3	2	0	0
4	5	0	0
1	0	3	2

**FPGA BRAM:**

→ B0	1	3	
→ B1	0	0	2
→ B2	3	4	1
→ B3	2	5	0
→ B4	0	0	3
→ B5	0	0	2

**BRAM Accesses Timeline:**

cycle: 0	1	2	3	4	5	6
<b>R</b> B0 → 1	<b>C</b> 1-0=1					

**R** Read Operation      **C** Compute Operation

# Decompression from BCSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

×

**BCSR Format (2x2 sub-blocks) :**

row indices 1 3

col. indices 0 0 2

values 3 2 0 0

4 5 0 0

1 0 3 2

**FPGA BRAM:**

→ B0 1 3

→ B1 0 0 2

→ B2 3 4 1

→ B3 2 5 0

→ B4 0 0 3

→ B5 0 0 2

**BRAM Accesses Timeline:**

**R** Read Operation      **C** Compute Operation

cycle: 0	1	2	3	4	5	6
<b>R</b> B0 → 1	<b>C</b> 1-0=1	<b>R</b> B1 → 0				
		<b>R</b> B2 → 3				
		<b>R</b> B3 → 2				
		<b>R</b> B4 → 0				
		<b>R</b> B5 → 0				



# Decompression from BCSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

×

**BCSR Format (2x2 sub-blocks) :**

row indices **1 3**

col. indices **0 0 2**

values

3	2	0	0
4	5	0	0
1	0	3	2

**FPGA BRAM:**

B0	1	3	
B1	0	0	2
B2	3	4	1
B3	2	5	0
B4	0	0	3
B5	0	0	2

**BRAM Accesses Timeline:**

**R** Read Operation      **C** Compute Operation

cycle: 0	1	2	3	4	5	6
<b>R</b> B0 → 1	<b>C</b> 1-0=1	<b>R</b> B1 → 0	<b>R</b> B0 → 3			
		<b>R</b> B2 → 3				
		<b>R</b> B3 → 2				
		<b>R</b> B4 → 0				
		<b>R</b> B5 → 0				

# Decompression from BCSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

×

**BCSR Format (2x2 sub-blocks) :**

row indices 

1	3
---	---

col. indices 

0	0	2
---	---	---

values 

3	2	0	0
4	5	0	0
1	0	3	2

**FPGA BRAM:**

→ B0	1	3	
→ B1	0	0	2
→ B2	3	4	1
→ B3	2	5	0
→ B4	0	0	3
→ B5	0	0	2

**BRAM Accesses Timeline:**

cycle: 0	1	2	3	4	5	6
<b>R</b> B0 → 1	<b>C</b> 1-0=1	<b>R</b> B1 → 0	<b>R</b> B0 → 3	<b>C</b> 3-1=2		
		<b>R</b> B2 → 3				
		<b>R</b> B3 → 2				
		<b>R</b> B4 → 0				
		<b>R</b> B5 → 0				

**R** Read Operation      **C** Compute Operation

# Decompression from BCSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

×

**BCSR Format (2x2 sub-blocks) :**

row indices 

1	3
---	---

col. indices 

0	0	2
---	---	---

values 

3	2	0	0
4	5	0	0
1	0	3	2

**FPGA BRAM:**

→ B0	1	3	
→ B1	0	0	2
→ B2	3	4	1
→ B3	2	5	0
→ B4	0	0	3
→ B5	0	0	2

**BRAM Accesses Timeline:**

cycle: 0	1	2	3	4	5	6
<b>R</b> B0 → 1	<b>C</b> 1-0=1	<b>R</b> B1 → 0	<b>R</b> B0 → 3	<b>C</b> 3-1=2	<b>R</b> B1 → 0	
		<b>R</b> B2 → 3			<b>R</b> B2 → 4	
		<b>R</b> B3 → 2			<b>R</b> B3 → 5	
		<b>R</b> B4 → 0			<b>R</b> B4 → 0	
		<b>R</b> B5 → 0			<b>R</b> B5 → 0	

# Decompression from BCSR format

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

×

**B<sup>T</sup>**

3	1	5	6
---	---	---	---

**BCSR Format (2x2 sub-blocks) :**

row indices 

1	3
---	---

col. indices 

0	0	2
---	---	---

values

3	2	0	0
4	5	0	0
1	0	3	2

**FPGA BRAM:**

→ B0	1	3	
→ B1	0	0	2
→ B2	3	4	1
→ B3	2	5	0
→ B4	0	0	3
→ B5	0	0	2

**BRAM Accesses Timeline:**

			<b>R</b> Read Operation		<b>C</b> Compute Operation	
<i>cycle: 0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>
<b>R</b> B0 → 1	<b>C</b> 1-0=1	<b>R</b> B1 → 0	<b>R</b> B0 → 3	<b>C</b> 3-1=2	<b>R</b> B1 → 0	<b>R</b> B1 → 2
		<b>R</b> B2 → 3			<b>R</b> B2 → 4	<b>R</b> B2 → 1
		<b>R</b> B3 → 2			<b>R</b> B3 → 5	<b>R</b> B3 → 0
		<b>R</b> B4 → 0			<b>R</b> B4 → 0	<b>R</b> B4 → 3
		<b>R</b> B5 → 0			<b>R</b> B5 → 0	<b>R</b> B5 → 2



# Key Challenge of Decompressing CSR and BCSR

---

21

Creating each row of data has following overheads:

- ▶ One access to the meta data
- ▶ One computation

Reading the column indices and values is sequential because

- ▶ We do not know in advance which elements of column indices and values are going to be accessed.
- ▶ We cannot partition and allocate those two vectors across the blocks of BRAM to guarantee parallel reads.



# Key Insights and Solutions

---

22

To address the challenge we propose Ascella

Ascella achieves the ideal streaming for sparse problems by

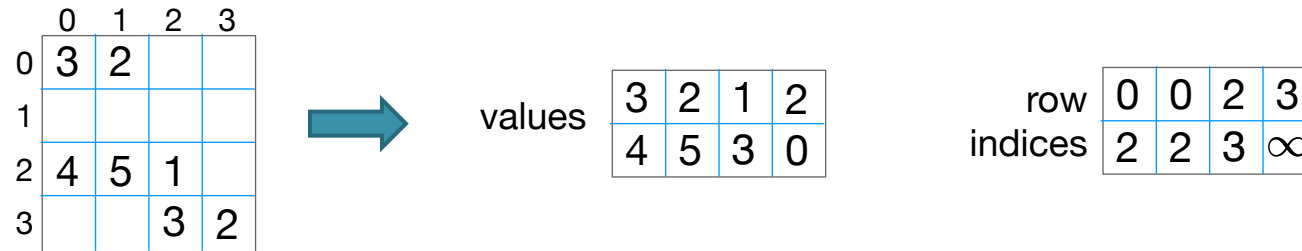
- ▶ Avoiding extra accesses to meta data
- ▶ Providing deterministic parallel accesses to data

Ascella is an accelerator for SpMV that sustains a balance between computation and data-transfer time



# Contributions of Ascella

Ascella uses a compressed format similar to list of lists (LIL<sup>1</sup>):

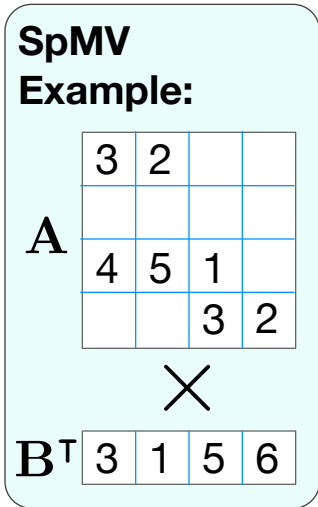


Ascella implements a lightweight microarchitecture that

- ▶ Connects the streamlines of memory to the parallel dot-product engine
- ▶ Enables ideal data streaming

<sup>1</sup>[https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.lil\\_matrix.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.lil_matrix.html)

# Decompression mechanism of Ascella



**LIL Format**

row indices	0	0	2	3
values	3	2	1	2
	4	5	3	0

**Mapping to BRAM Blocks**

B0	B1	B2	B3	B4	B5	B6	B7
0	0	2	3	3	2	1	2
2	2	3	∞	4	5	3	0
row indices				values			

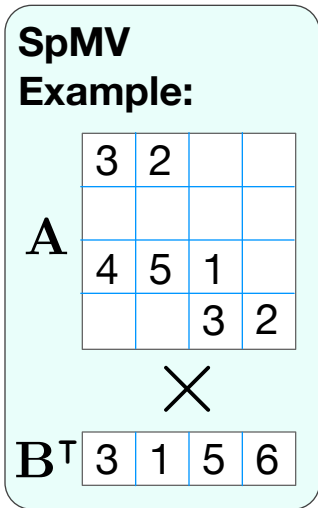
**BRAM Accesses for Decompression:**

		cycle: 0		1		2	
R	B0	→	0	R	B4	→	3
R	B1	→	0	R	B5	→	2
R	B2	→	2	R	B6	→	1
R	B3	→	3	R	B7	→	2

**R** Read Operation



# Decompression mechanism of Ascella



**LIL Format**

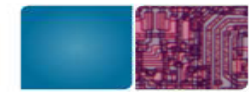
row indices	0	0	2	3
values	3	2	1	2
	4	5	3	0

**Mapping to BRAM Blocks**

B0	B1	B2	B3	B4	B5	B6	B7
0	0	2	3	3	2	1	2
2	2	3	∞	4	5	3	0
row indices				values			

**BRAM Accesses for Decompression:**

cycle: 0				1				2			
R B0 → 0	R B1 → 0	R B2 → 2	R B3 → 3	R B0 → 2	R B1 → 2	R B2 → 2	R B3 → 3	R B4 → 4	R B5 → 5	R B6 → 1	R B7 → 2
R B4 → 3	R B5 → 2	R B6 → 1	R B7 → 2	R B4 → 4	R B5 → 5	R B6 → 1	R B7 → 2				



# Decompression mechanism of Ascella

**SpMV Example:**

**A**

3	2		
4	5	1	
		3	2

×

**B<sup>T</sup>** 3 1 5 6

**LIL Format**

row indices	0	0	2	3
values	3	2	1	2
	4	5	3	0

**Mapping to BRAM Blocks**

B0	B1	B2	B3	B4	B5	B6	B7
0	0	2	3	3	2	1	2
2	2	3	∞	4	5	3	0
row indices				values			

**BRAM Accesses for Decompression:**

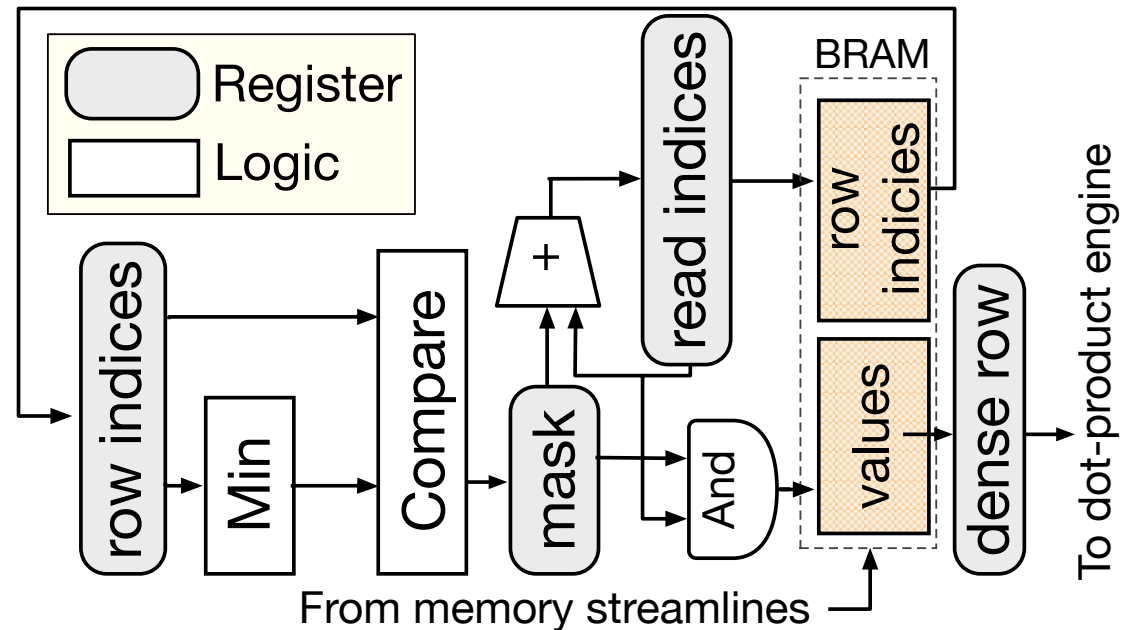
cycle: 0				1				2							
R B0 → 0	R B1 → 0	R B2 → 2	R B3 → 3	R B4 → 3	R B5 → 2	R B6 → 1	R B7 → 2	R B0 → 2	R B1 → 2	R B2 → 2	R B3 → 3	R B4 → 4	R B5 → 5	R B6 → 1	R B7 → 2
												R B2 → 3	R B3 → 3	R B6 → 2	R B7 → 3

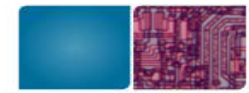


# Microarchitecture of Ascella

At each step of decompression:

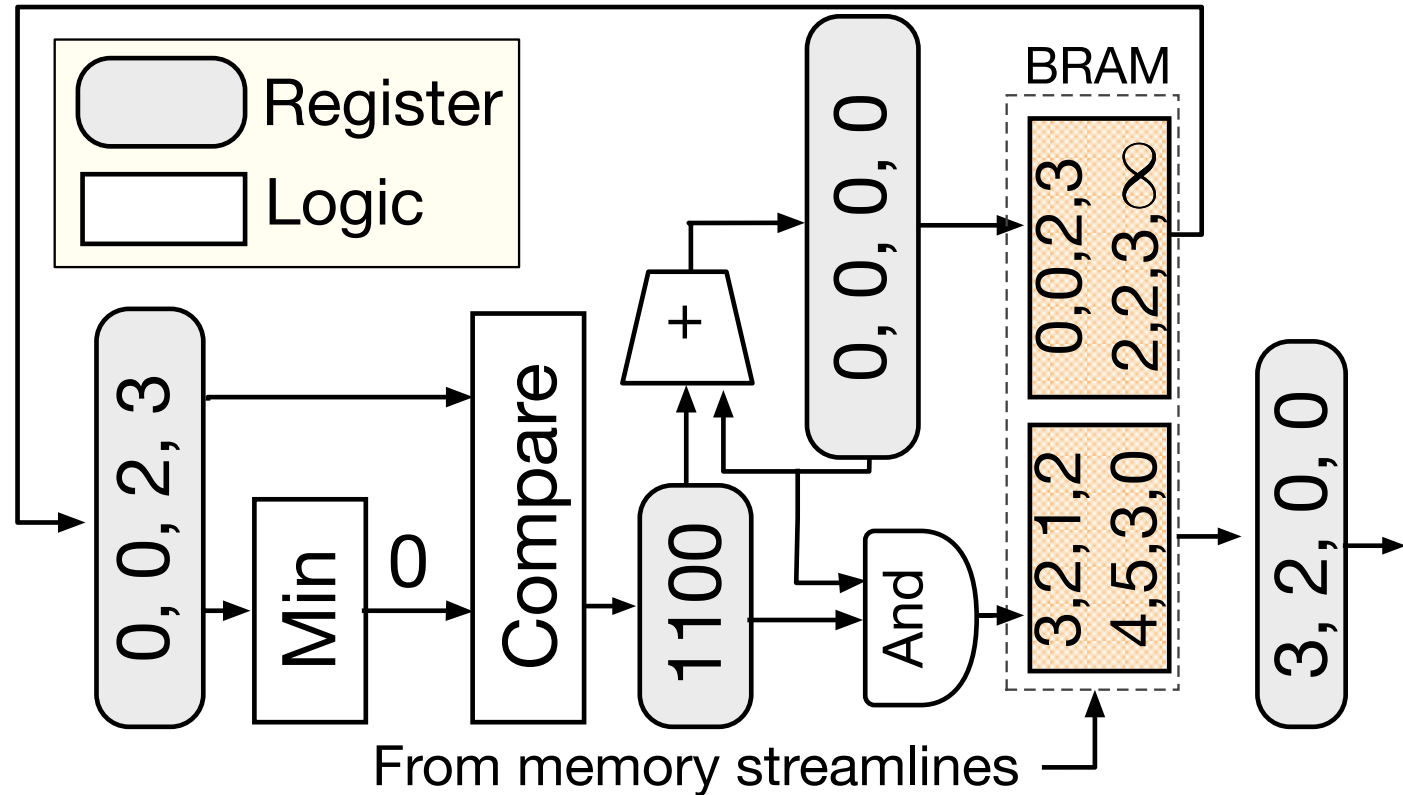
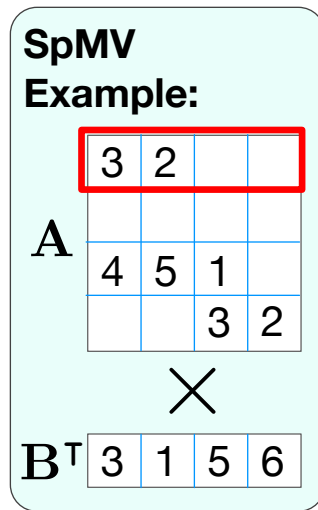
- ▶ ***read indices*** are used to to read the ***row indices***
- ▶ the minimum of ***row indices*** is used to create a ***mask***
- ▶ the mask
  - ▶ Selects values of ***dense row***
  - ▶ Updates the ***read indices***





# Microarchitecture of Ascella

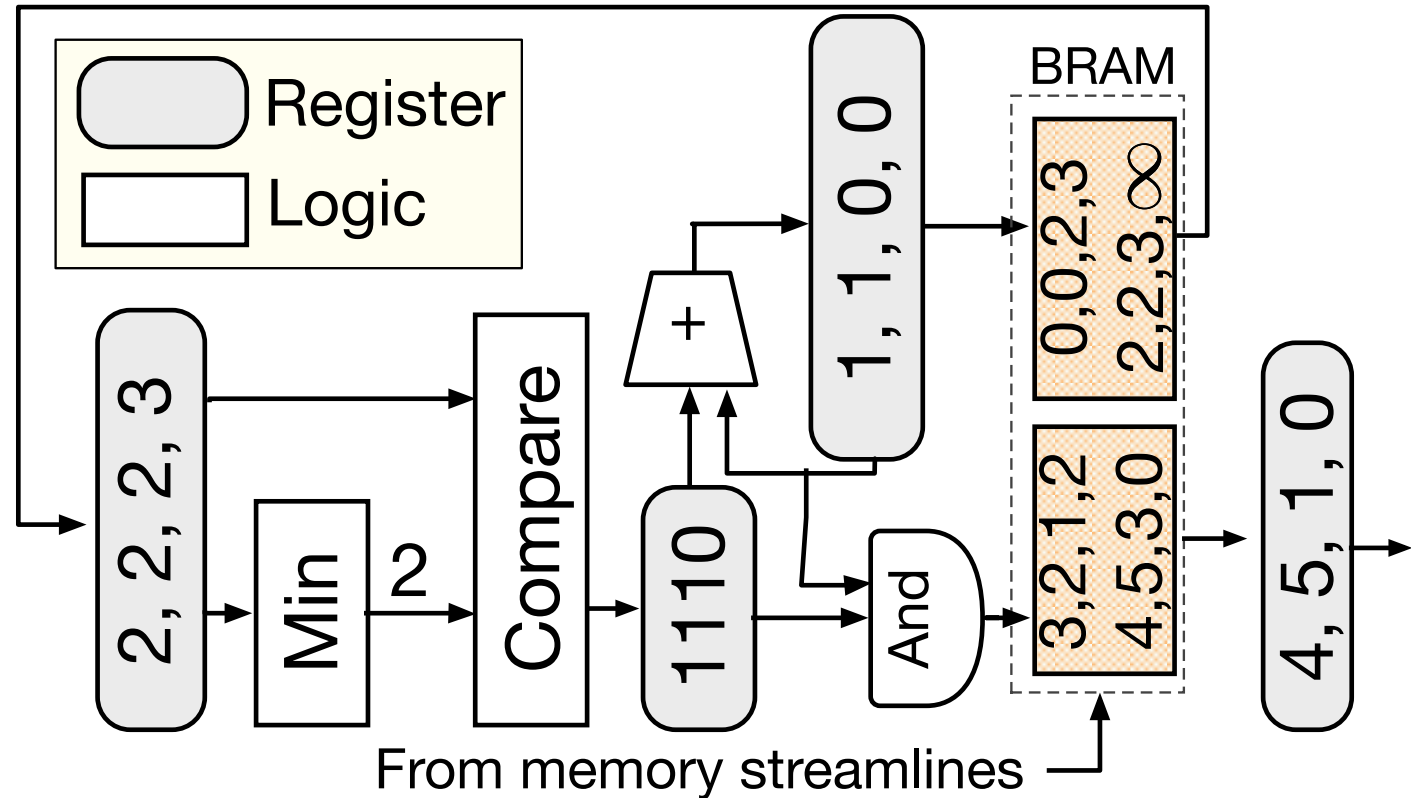
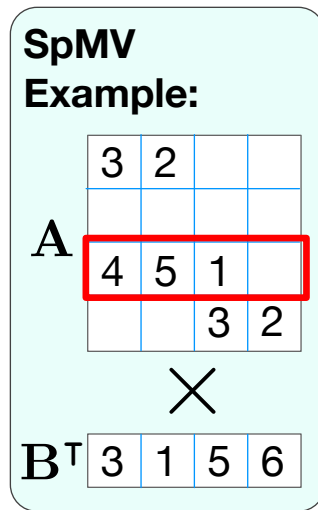
Creating the first row:

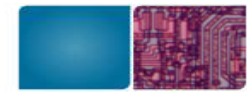




# Microarchitecture of Ascella

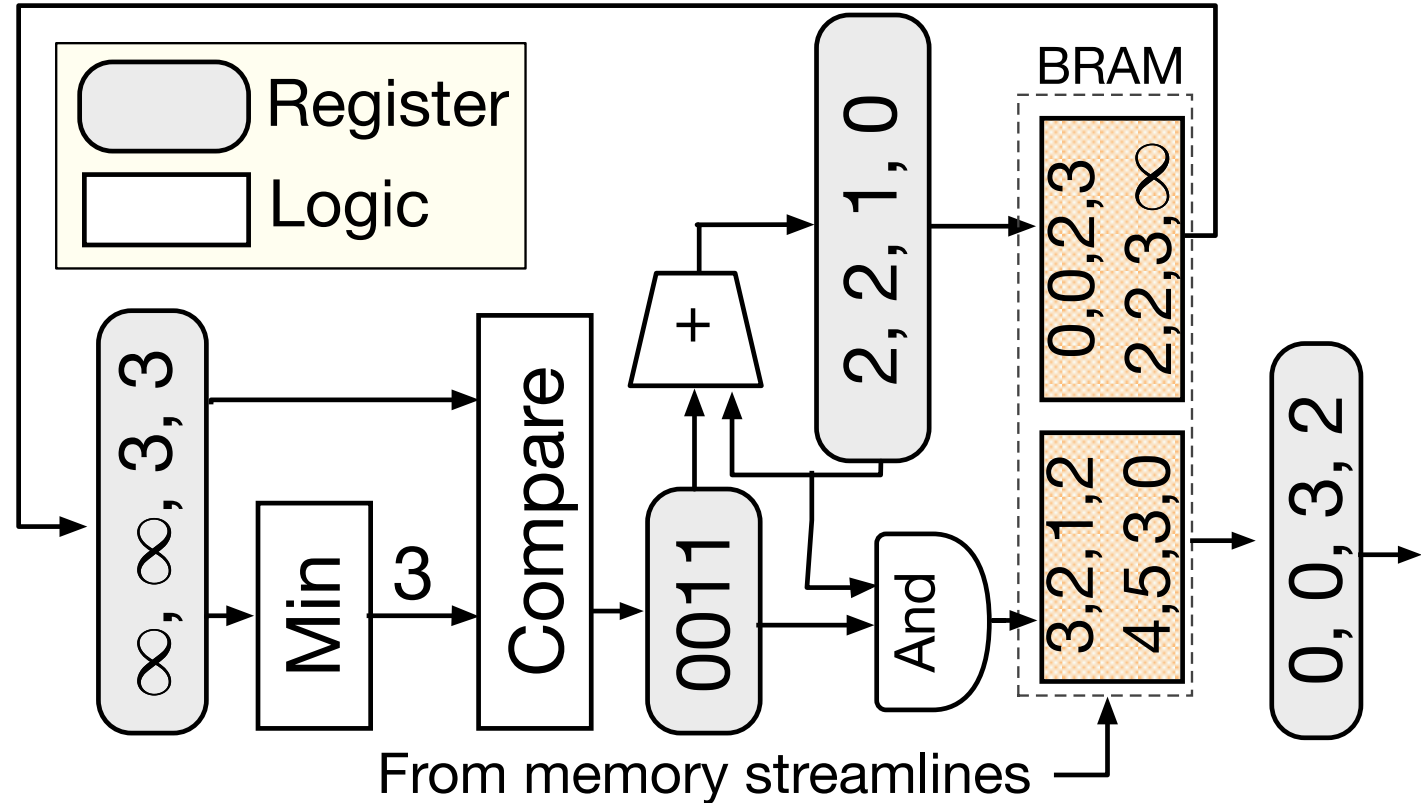
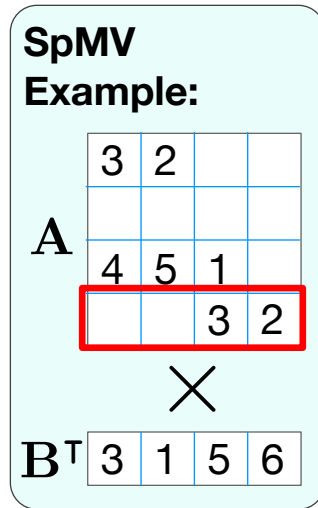
Creating the third row:





# Microarchitecture of Ascella

Creating the last row:





# Experimental Setup

31

We apply SpMV on a group of matrices from SuiteSparse collection

We implement Ascella and the baselines

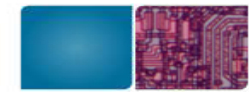
- ▶ Using Xilinx Vivado HLS
- ▶ On ZYNQ XC7Z020FPGA.

Our comparison metrics are

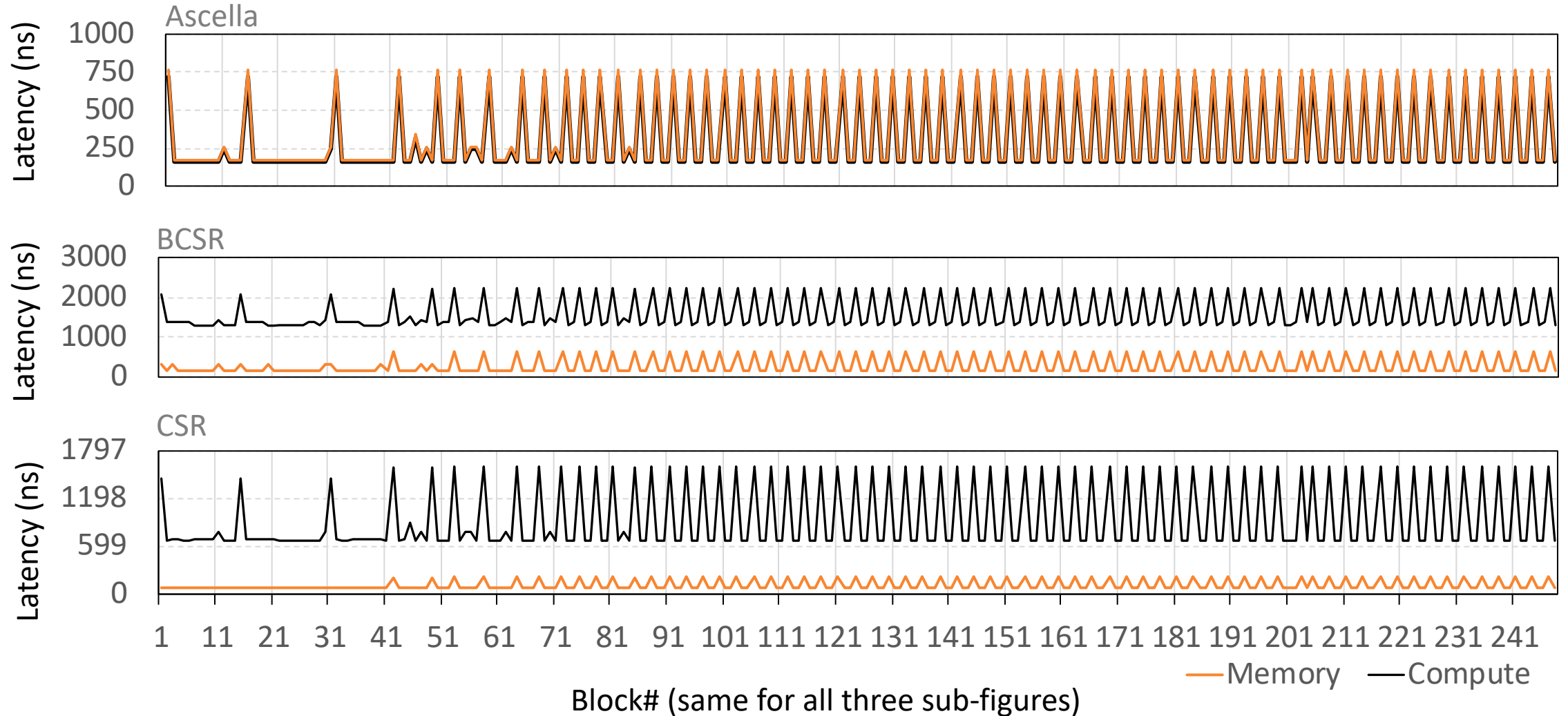
- ▶ Latency
- ▶ Recourse utilization

The configurations of Ascella and baselines include

- ▶ AXI stream interfaces
- ▶ DDR3 memory
- ▶ 100 MHz frequency
- ▶ 32-bit integers



# Performance Evaluation



Dataset: thermomech-TC

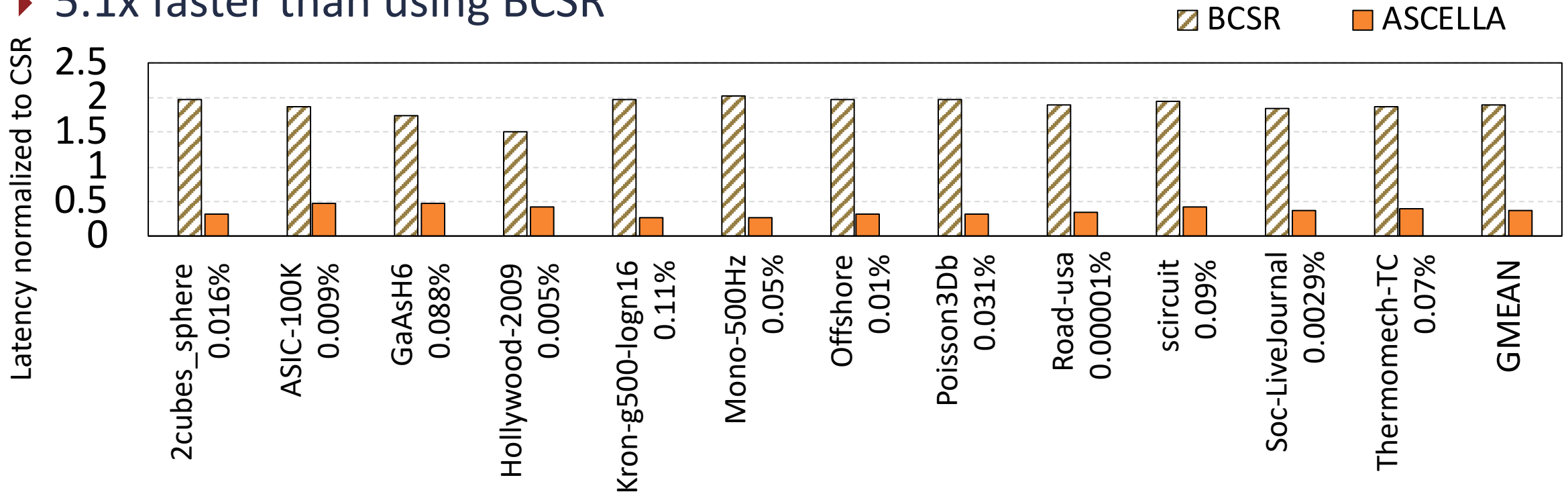




# Performance Evaluation

On average, Ascella executes SpMV

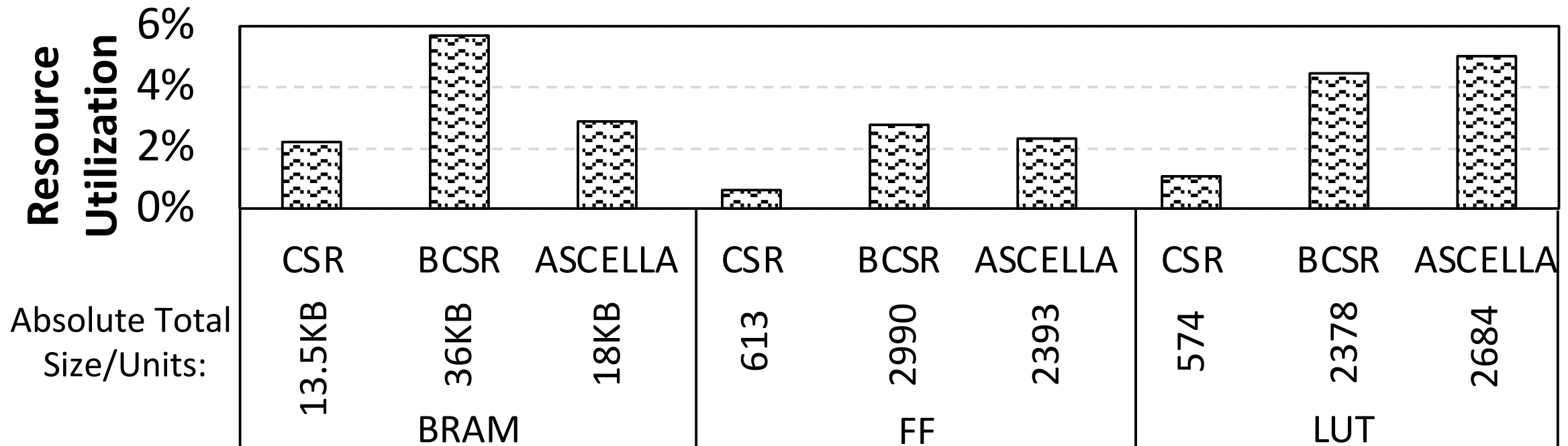
- ▶ 2.7x faster than using CSR
- ▶ 5.1x faster than using BCSR





# Resource Utilization

- ▶ BCSR and Ascella use more BRAM because of partitioning.
- ▶ CSR has the lowest flip-flop and look-up table (LUT) utilization because it does not implement any parallelism.





# Conclusions

---

35

Ascella is a streaming accelerator for sparse problems that

- ▶ Streams the non-zero values of sparse matrices and processes them as they come at the same pace
- ▶ Is a significant step towards accelerating larger sparse problems because its storage format
  - ▶ facilitates partitioning large matrices
  - ▶ is supported in Python libraries, which makes the implementation straightforward