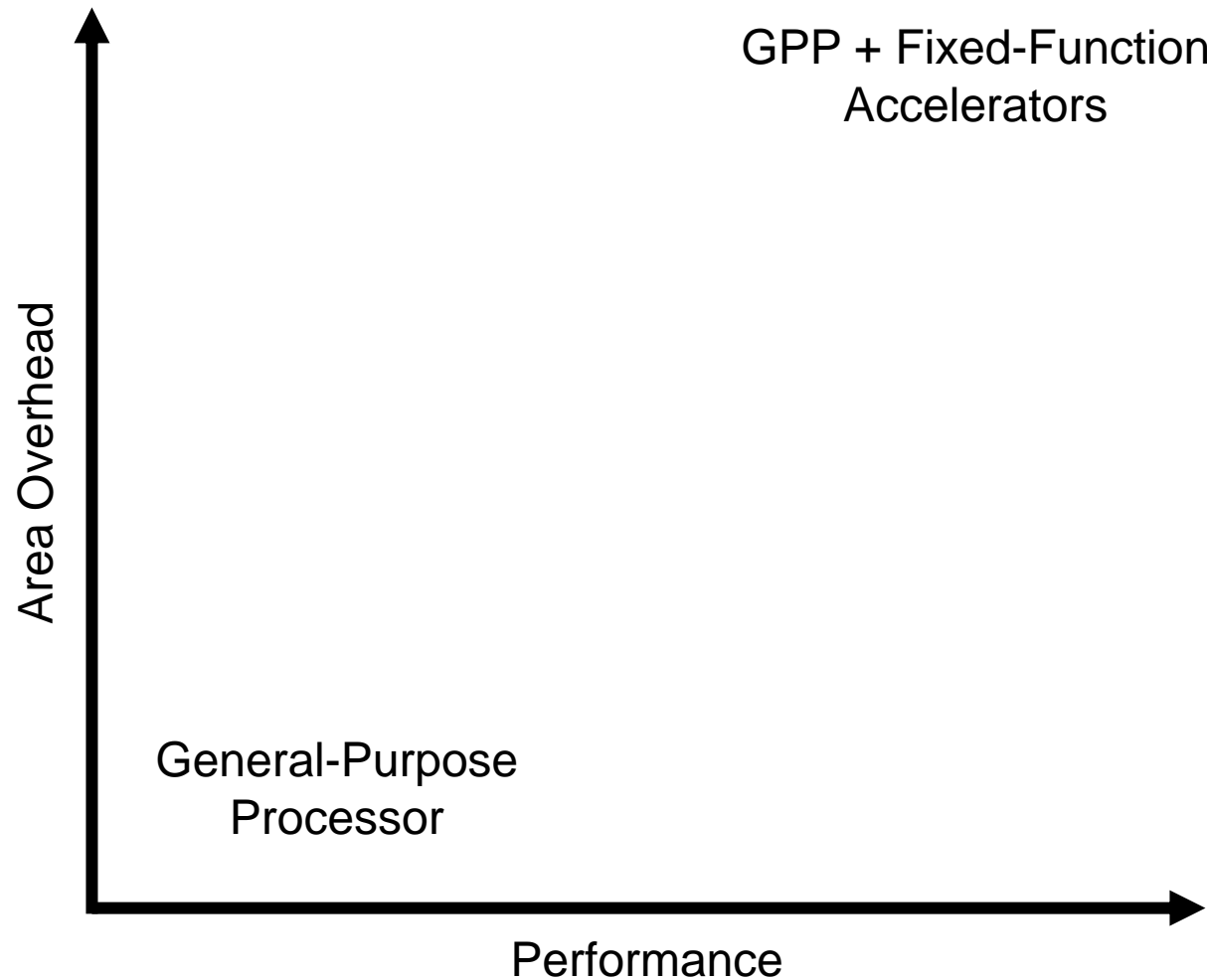**CSL**

# EVE: Ephemeral Vector Engines

*Khalid Al-Hawaj*, Tuan Ta, Nick Cebry, Shady Agwa, Olalekan Afuye, Eric Hall, Courtney Golden, Alyssa B. Apsel, Christopher Batten
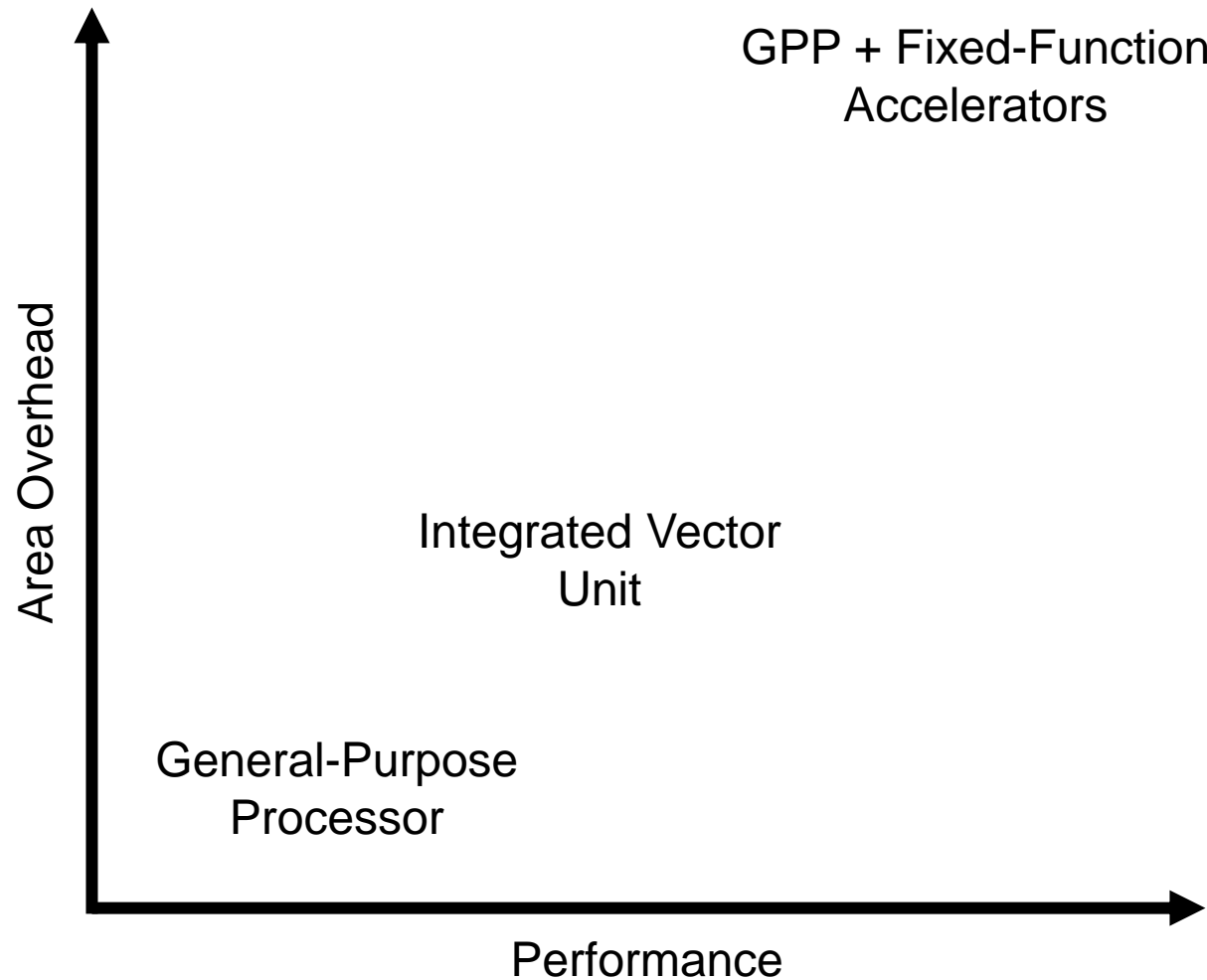
**Electrical and Computer Engineering**

**Cornell University**

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 1 of 15

Area Overhead

Performance

GPP + Fixed-Function Accelerators
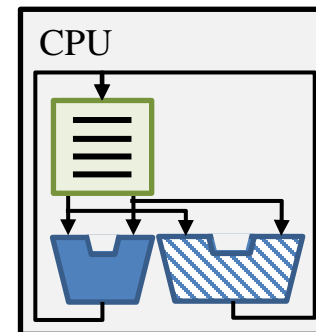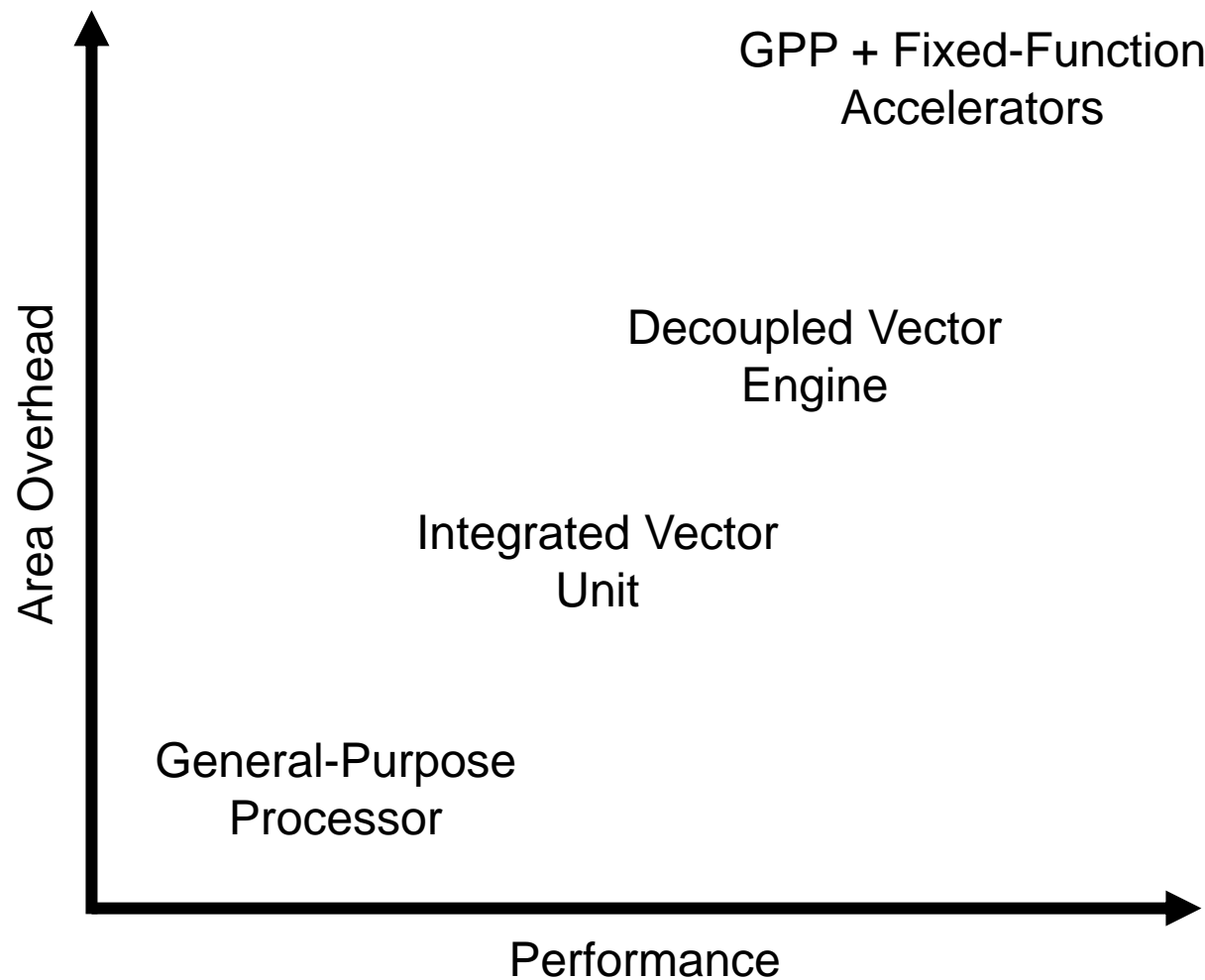
General-Purpose Processor

- Computer architects rely on specialization to increase performance and efficiency

- Vector micro-architectures to tackle regular data-parallel applications

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 1 of 15

■ Computer architects rely on specialization to increase performance and efficiency

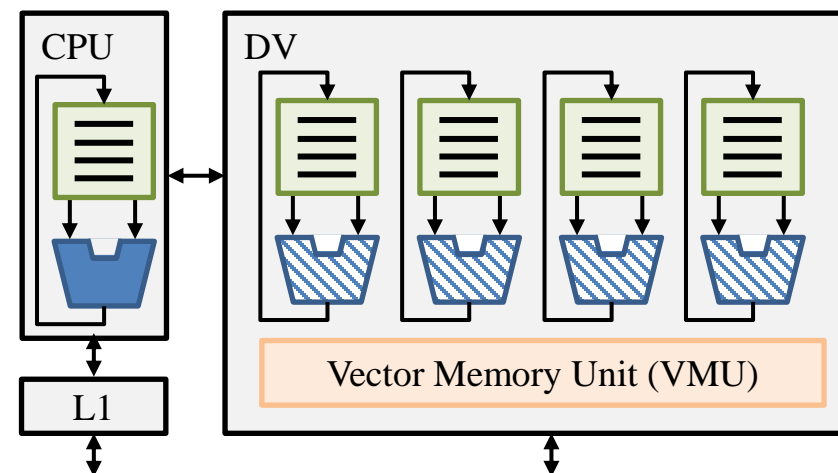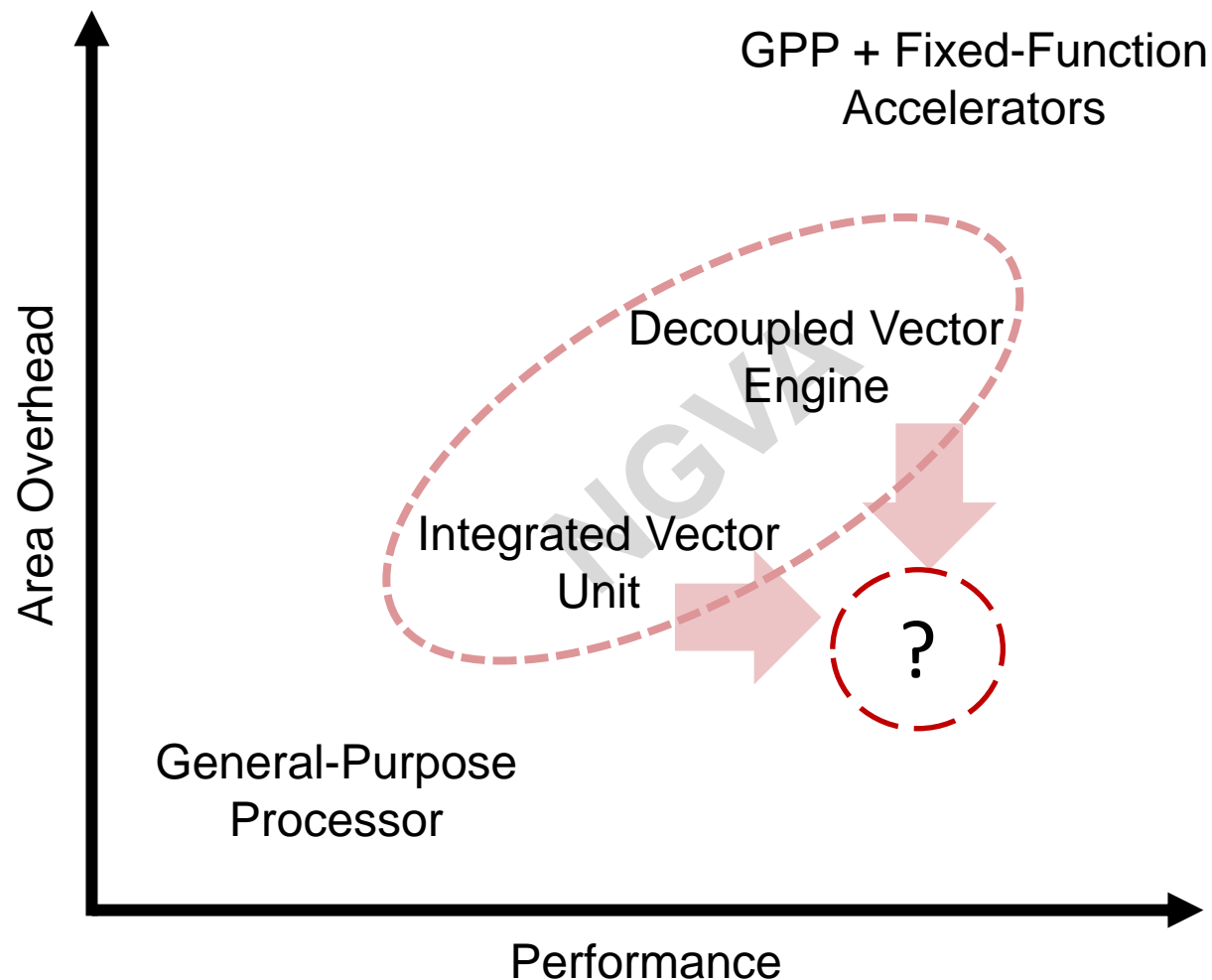■ Vector micro-architectures to tackle regular data-parallel applications

• Integrated vector unit (IV)

Area Overhead

GPP + Fixed-Function Accelerators

Integrated Vector Unit

General-Purpose Processor

Performance

CPU

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 1 of 15

- Computer architects rely on specialization to increase performance and efficiency

- Vector micro-architectures to tackle regular data-parallel applications
  - Integrated vector unit (IV)
  - Decoupled vector engine (DV)

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 1 of 15

- Computer architects rely on specialization to increase performance and efficiency

- Vector micro-architectures to tackle regular data-parallel applications
  - Integrated vector unit (IV)
  - Decoupled vector engine (DV)

- Emerging trend of next-generation vector architectures (NGVA)

**Is it possible to achieve performance comparable to a _DV_ while incurring an area overhead equivalent to an _IV_?**

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 1 of 15

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
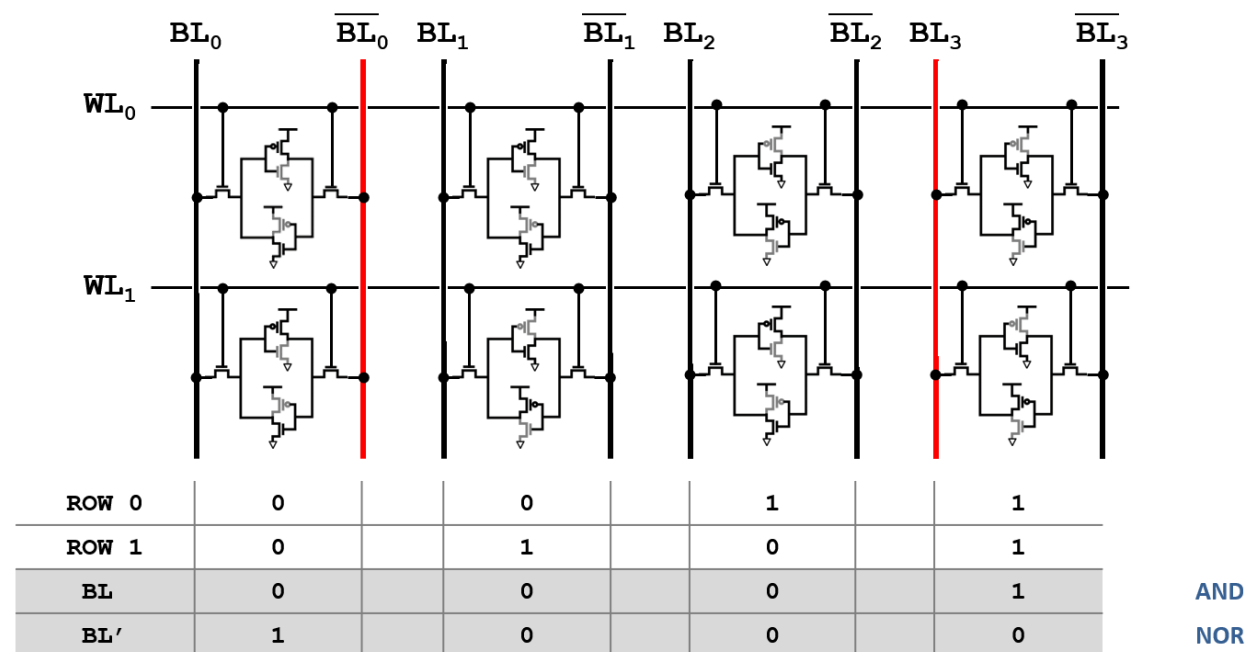
Page 2 of 15

- Recent work on SRAM-based compute-in-memory have shown promise in alleviating the area-overhead often associated with vector execution

- Subsequent work has explored implementing more complex operations on-top of bit-line compute

**What <u>abstraction</u> would be more suitable to enable high-programmability of an SRAM-based compute-in-memory micro-architecture?**

**A Configurable TCAM / BCAM / SRAM using 28nm push-rule 6T bit cell**
Supreet Jeloka[1], Naveen Akesh[2], Dennis Sylvester[1], and David Blaauw[1]
[1]University of Michigan, Ann Arbor, MI, [2]Oracle, Santa Clara, CA



| | $BL_0$ | $BL_1$ | $BL_2$ | $BL_3$ | |
|---|---|---|---|---|---|
| ROW 0 | 0 | 0 | 1 | 1 | |
| ROW 1 | 0 | 1 | 0 | 1 | |
| BL | 0 | 0 | 0 | 1 | AND |
| BL' | 1 | 0 | 0 | 0 | NOR |

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 2 of 15

**Fixed-Function Accelerators**

**Decoupled Vector Engine**

**Integrated Vector Unit**

**?**

**General-Purpose Processor**

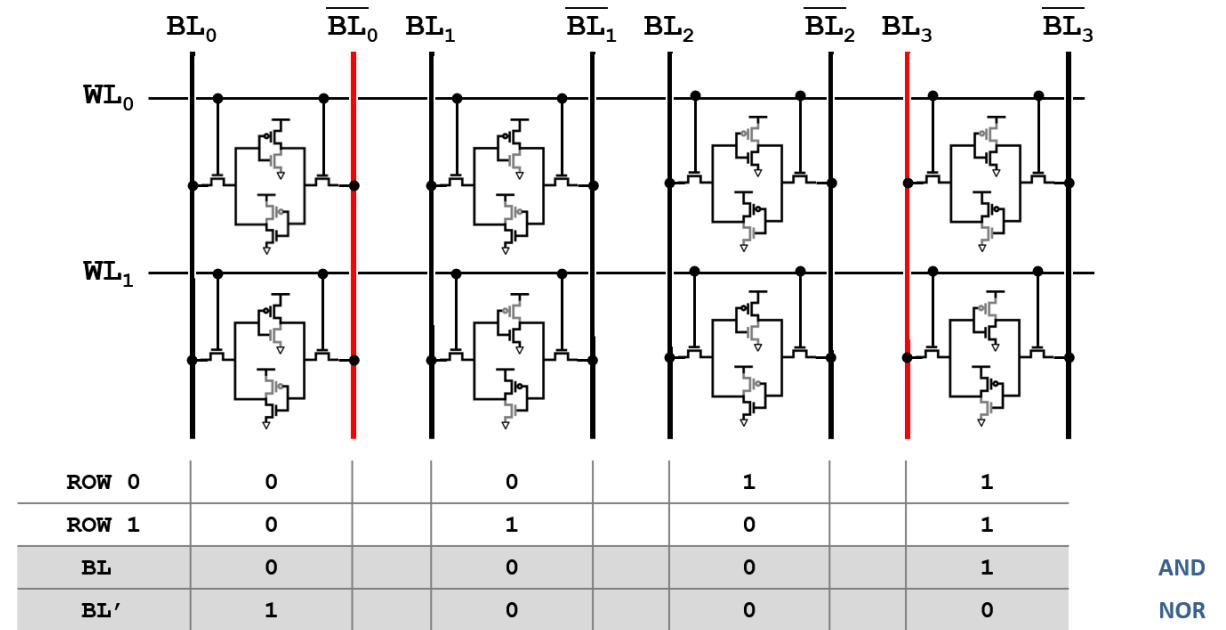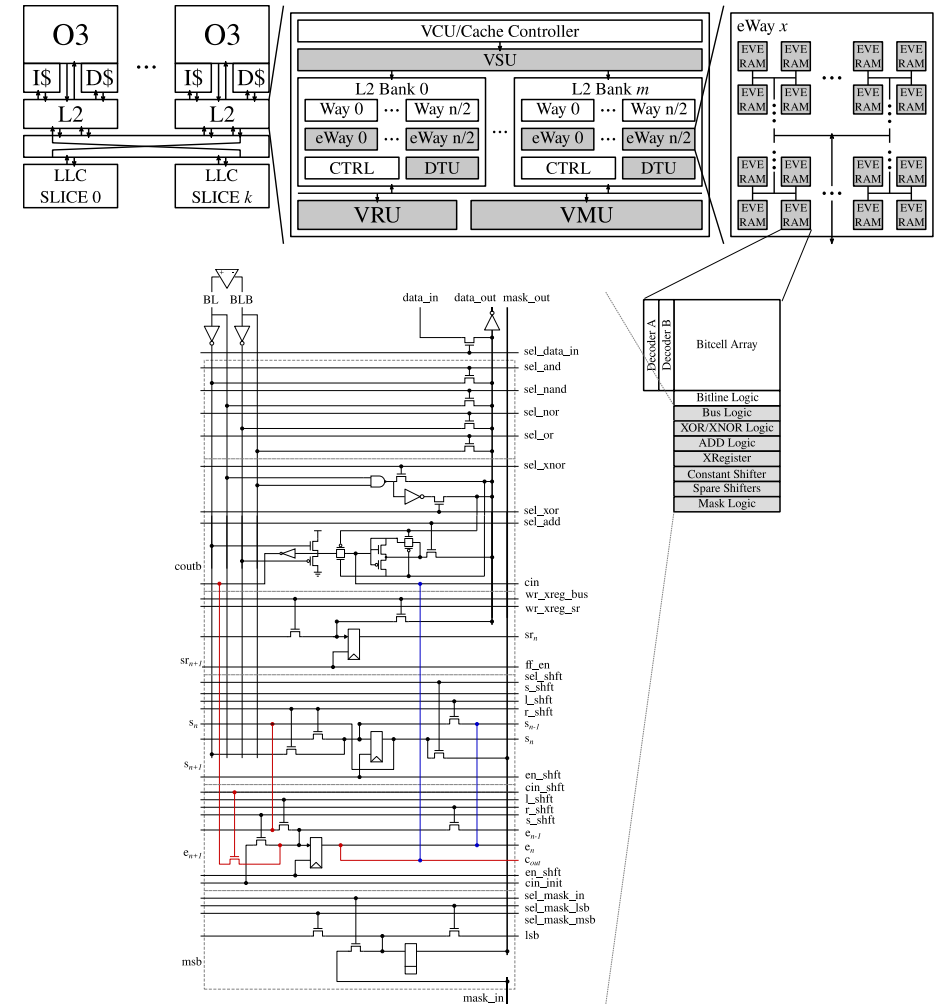Area Overhead (y-axis) vs Performance (x-axis)

**Is it possible to achieve performance <u>comparable</u> to a _DV_ while incurring an area overhead <u>equivalent</u> to an _IV_?**

**A Configurable TCAM / BCAM / SRAM using 28nm push-rule 6T bit cell**

Supreet Jeloka[1], Naveen Akesh[2], Dennis Sylvester[1], and David Blaauw[1]
[1]University of Michigan, Ann Arbor, MI, [2]Oracle, Santa Clara, CA

| | | | | | | |
|---|---|---|---|---|---|---|
| ROW 0 | 0 | | 0 | | 1 | 1 |
| ROW 1 | 0 | | 1 | | 0 | 1 |
| BL | 0 | | 0 | | 0 | 1 | **AND** |
| BL' | 1 | | 0 | | 0 | 0 | **NOR** |

**What <u>abstraction</u> would be more suitable to enable high-programmability of an SRAM-based compute-in-memory micro-architecture?**

Cornell University
Computer Systems Laboratory
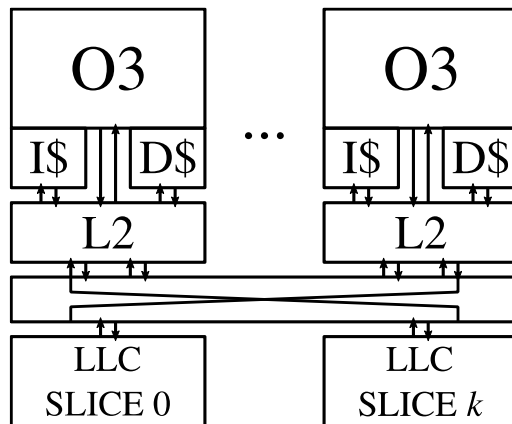
Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 3 of 15

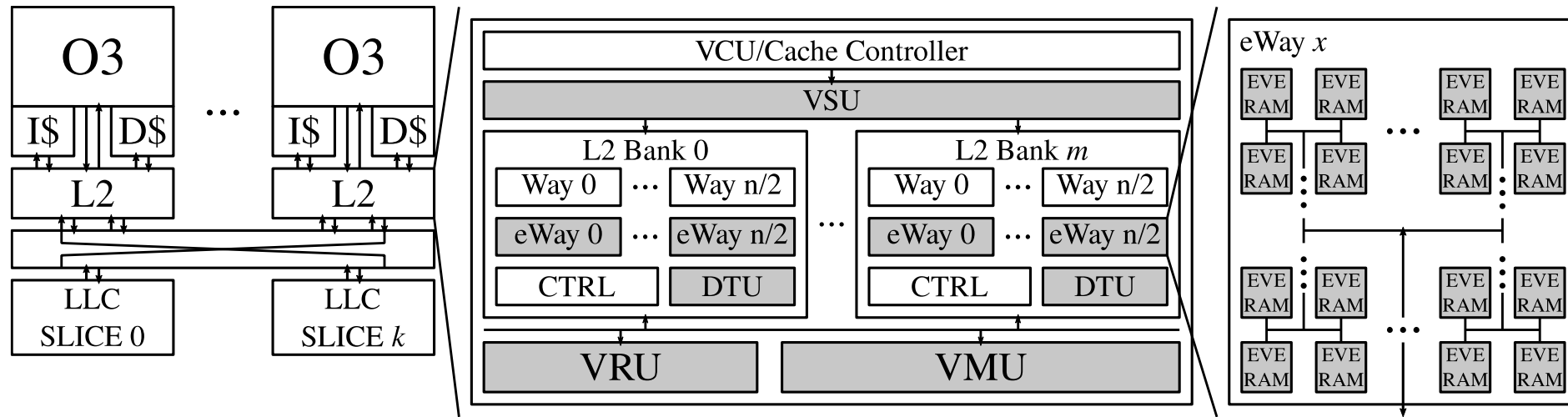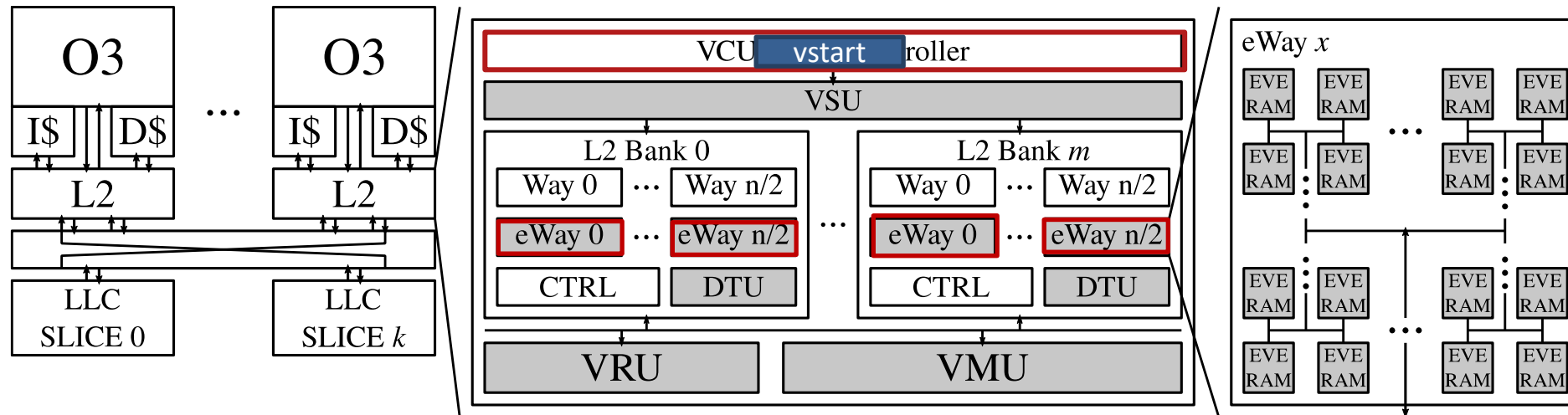| | Duality Cache [ISCA2019] | EVE |
|---|---|---|
| **Abstraction** | SIMT/Threading | Next-Generation Vector |
| **Cache Hierarchy** | Last Level Cache | Level-2 Cache |
| **Private or Shared** | Shared | Private |
| **Execution Paradigm** | Bit-Serial Execution | Bit-Hybrid Execution |

- Architectural template for a novel SRAM-based compute-in-memory next-gen vector engine that supports the full RISC-V RVV specifications

- Bit-hybrid execution to balance throughout and latency by alleviating row and column under-utilization

- Detailed evaluation of EVE show-casing the impacts and benefits of bit-hybrid execution on an SRAM-based compute-in-memory micro-architecture

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 4 of 15

**Motivation**

- **EVE Micro-Architecture**

- **EVE Bit-Hybrid Execution Paradigm**

- **EVE Micro-Programming & Circuits**

- **EVE Evaluation**

**Conclusion**

**Motivation**

- **EVE Micro-Architecture**

- EVE Bit-Hybrid Execution Paradigm

- EVE Micro-Programming & Circuits

- EVE Evaluation

**Conclusion**

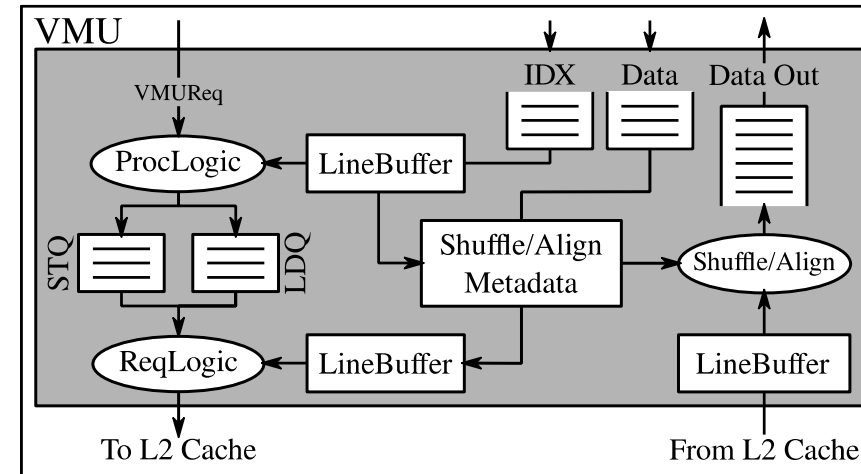EVE reconfigure parts of the private **L2 cache** to act as the vector execution hardware **on-demand**

Motivation • **Micro-Architecture** • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 5 of 15

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 5 of 15

Motivation • **Micro-Architecture** • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
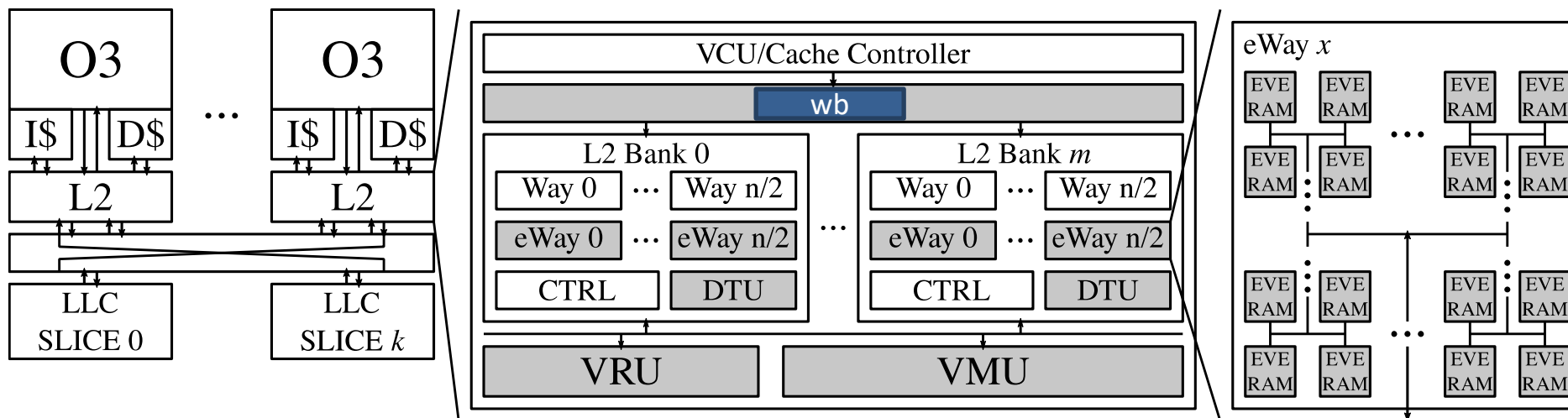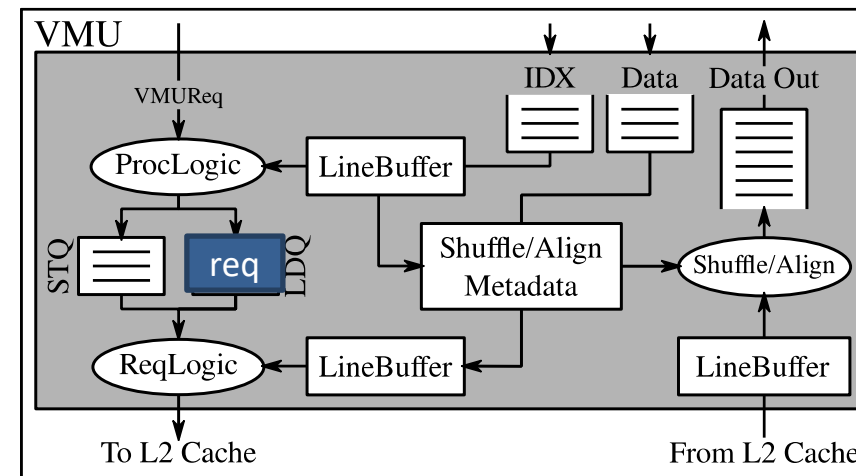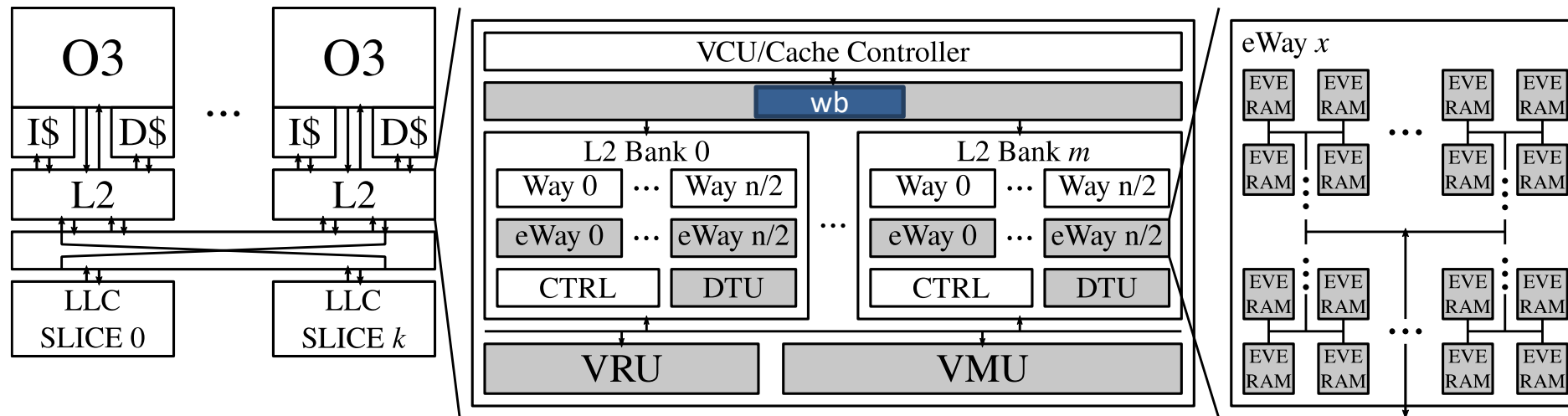
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
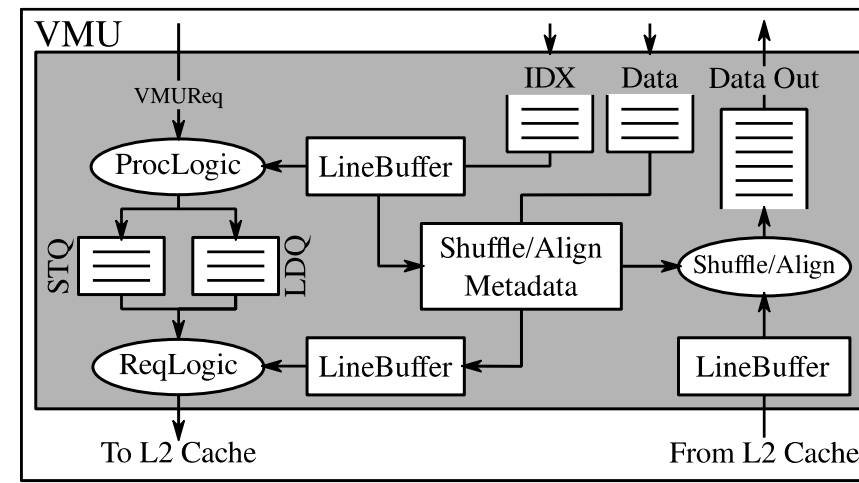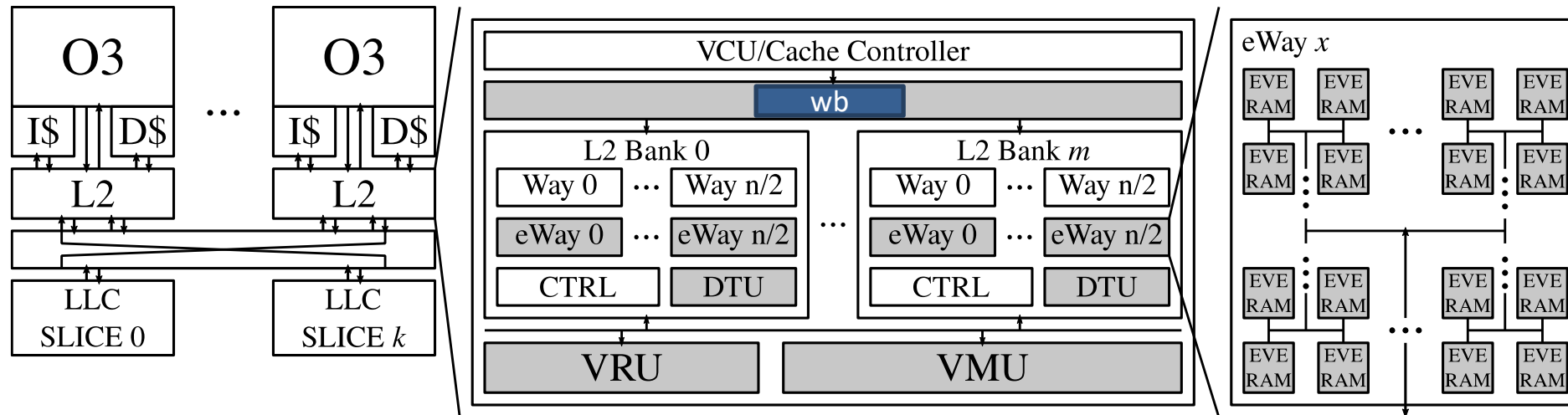
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • **Micro-Architecture** • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
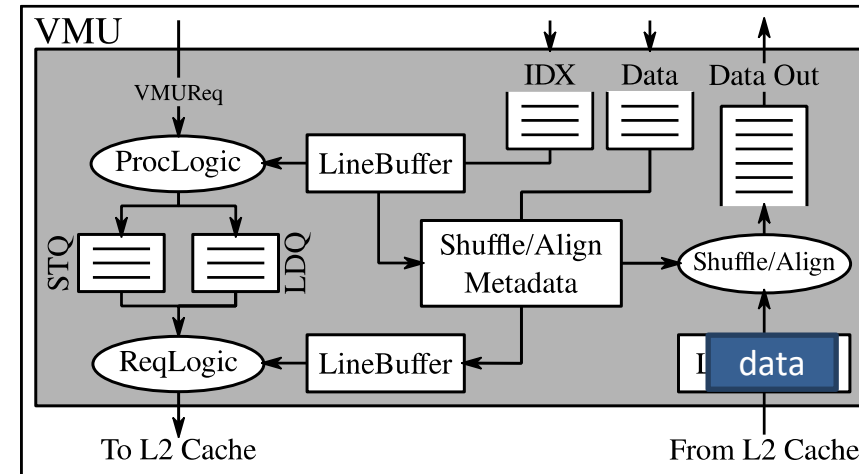
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
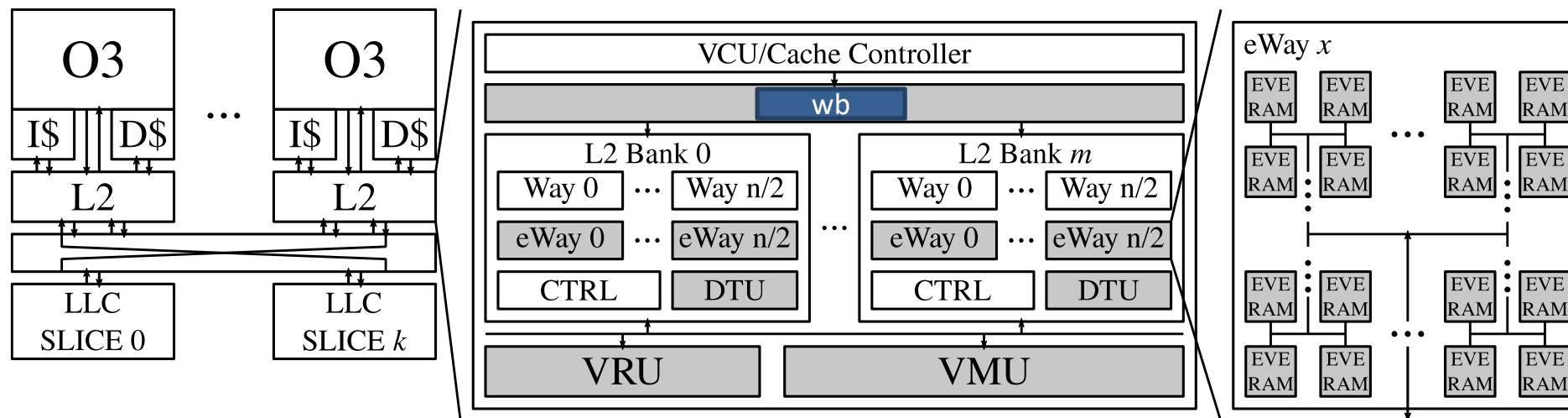
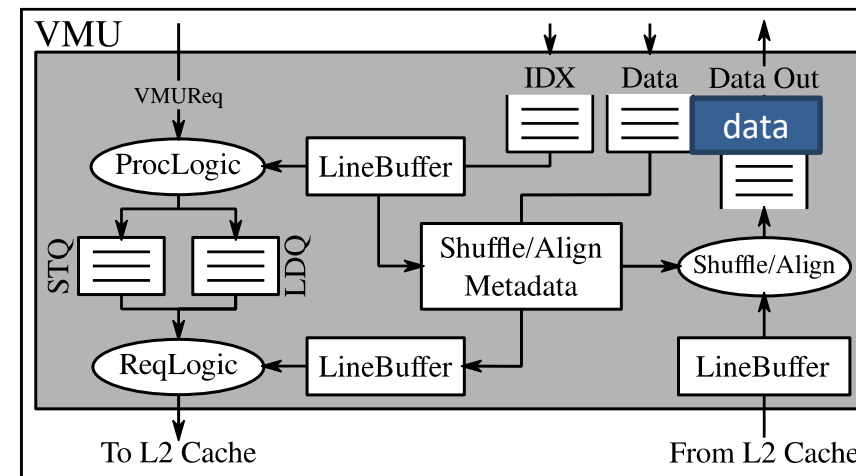Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Cornell University
Computer Systems Laboratory
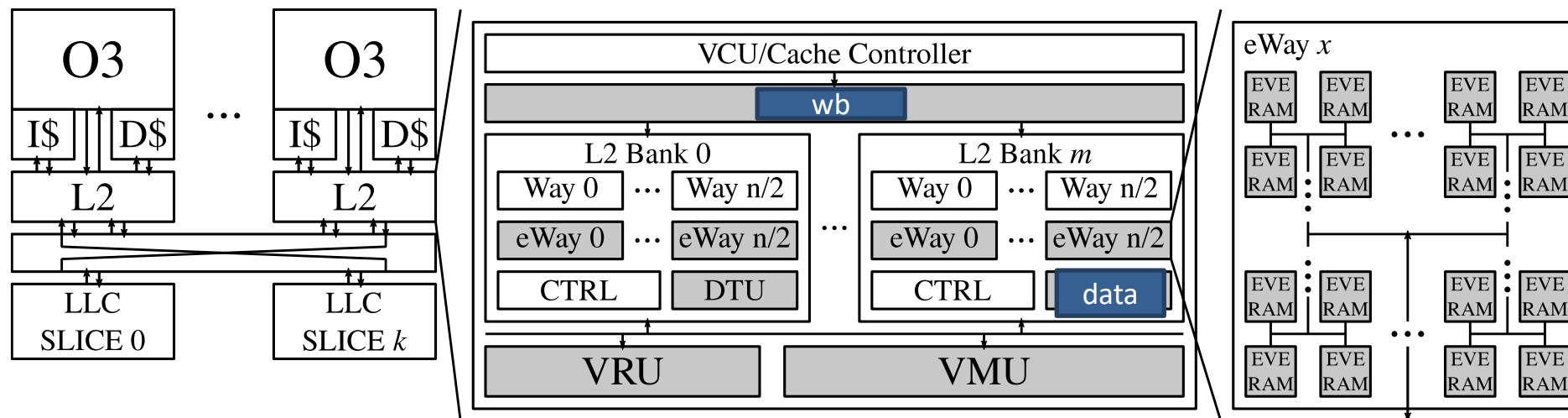
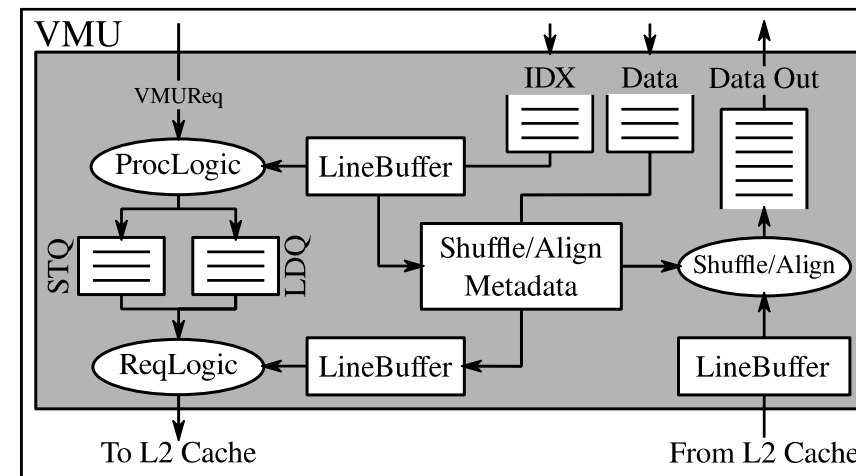Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 6 of 15

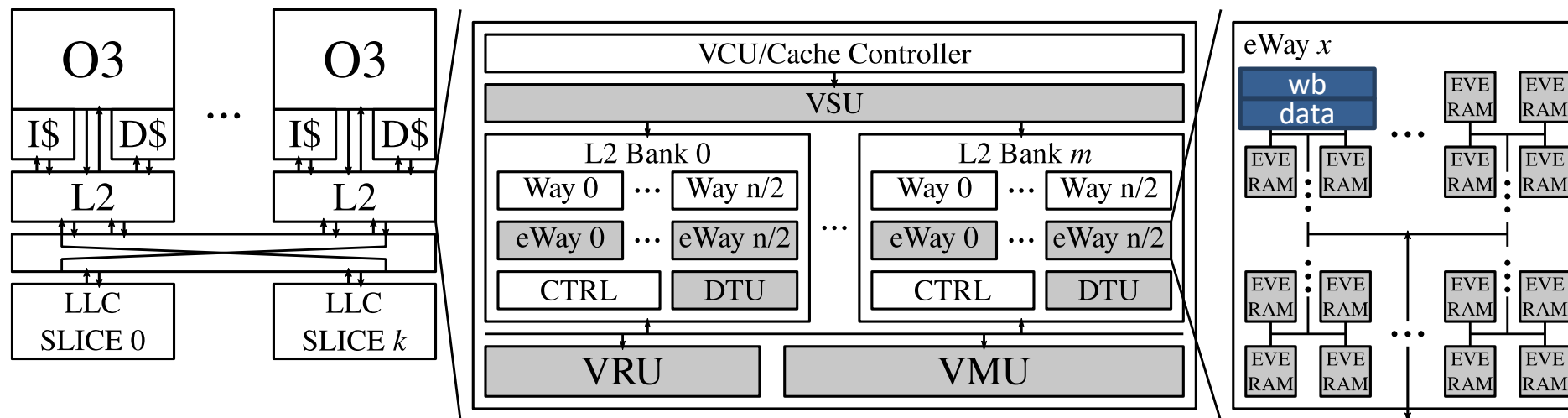Cornell University
Computer Systems Laboratory

```
void vvadd( int* c, int* a, int * b) {
  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
→   vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 6 of 15
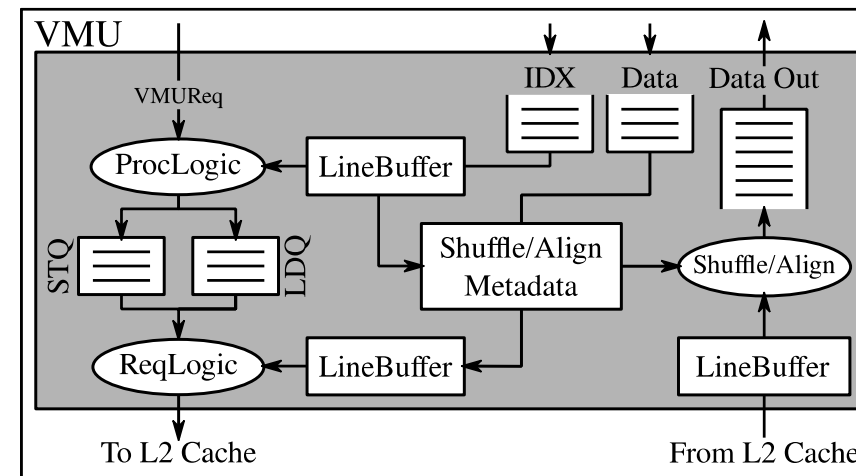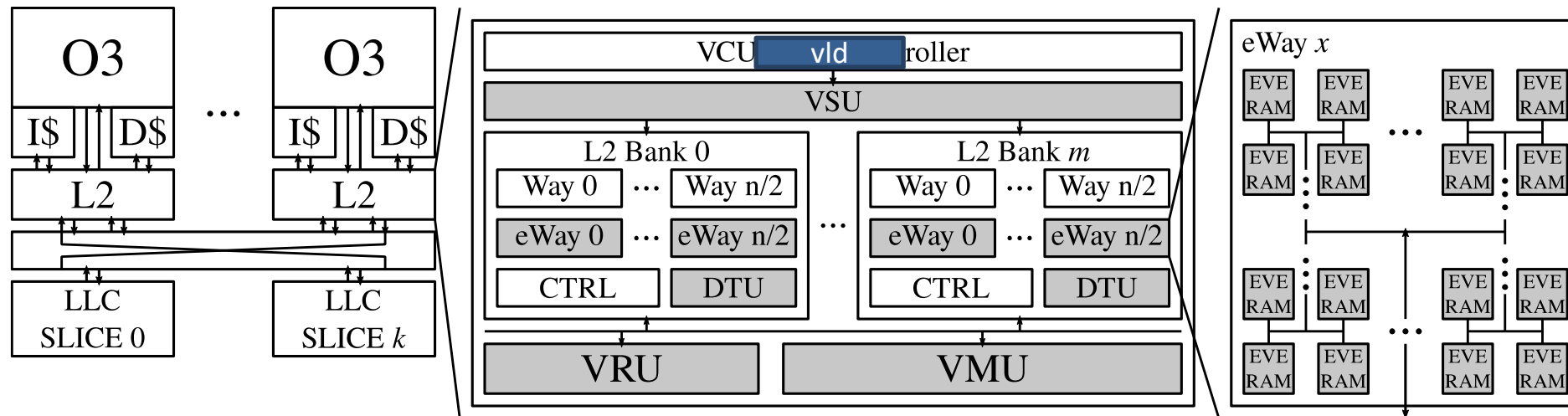
```
void vvadd( int* c, int* a, int * b) {

    vstart();

    for (int i = 0; i < len; i += vsetvl(len)) {
        vld( v0, a + i );
        vld( v1, b + i );
        vadd.vv( v2, v0, v1 );
        vst( v2, c + i );
    }

    vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
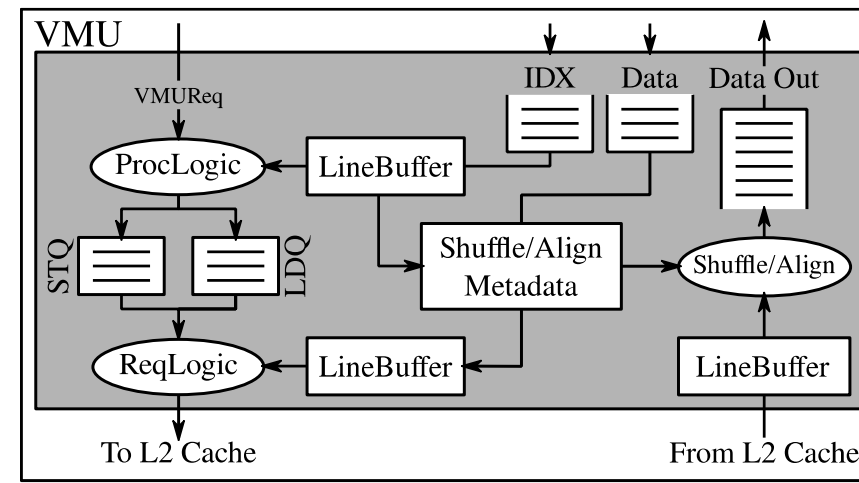
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Cornell University
Computer Systems Laboratory
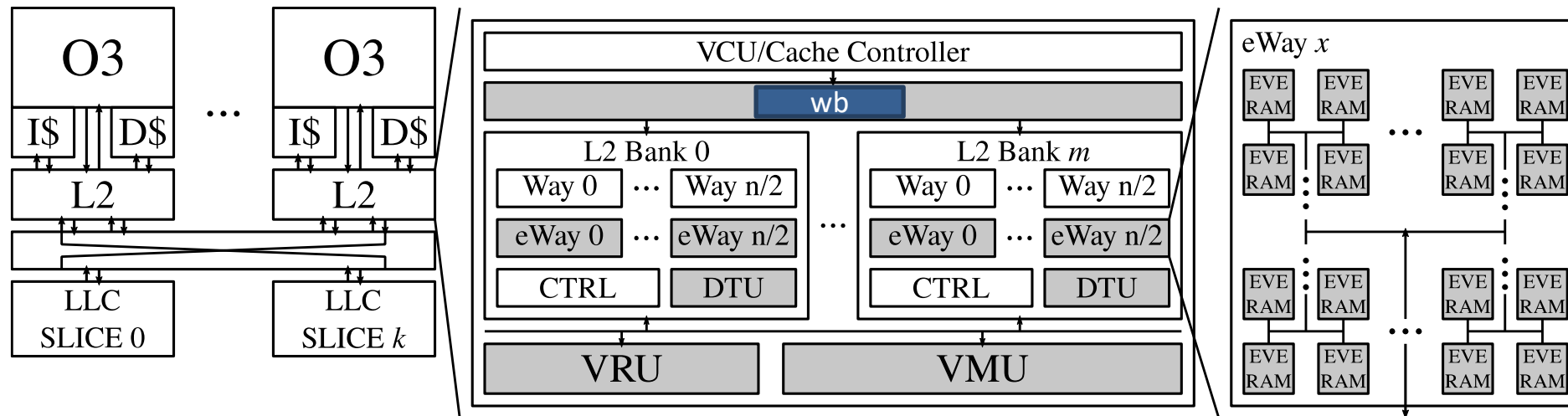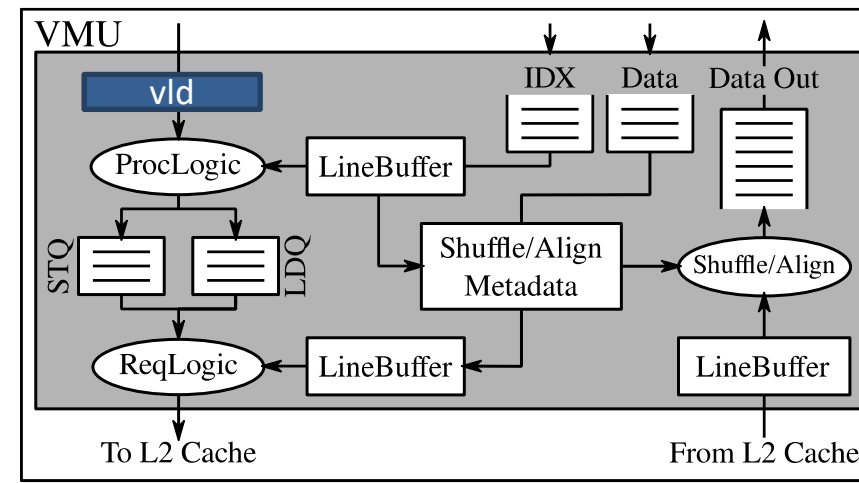
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();

}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
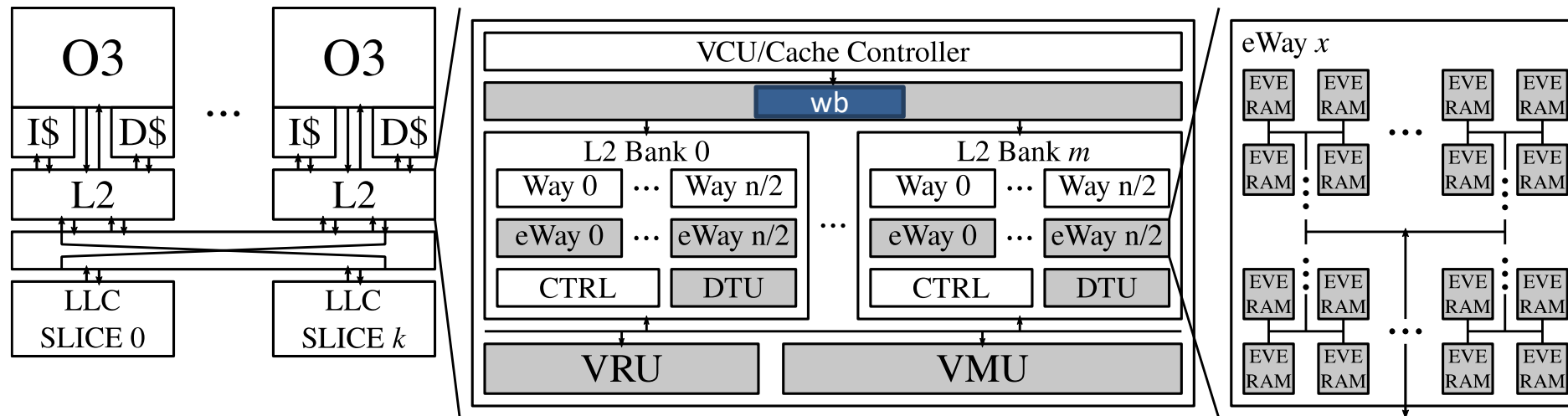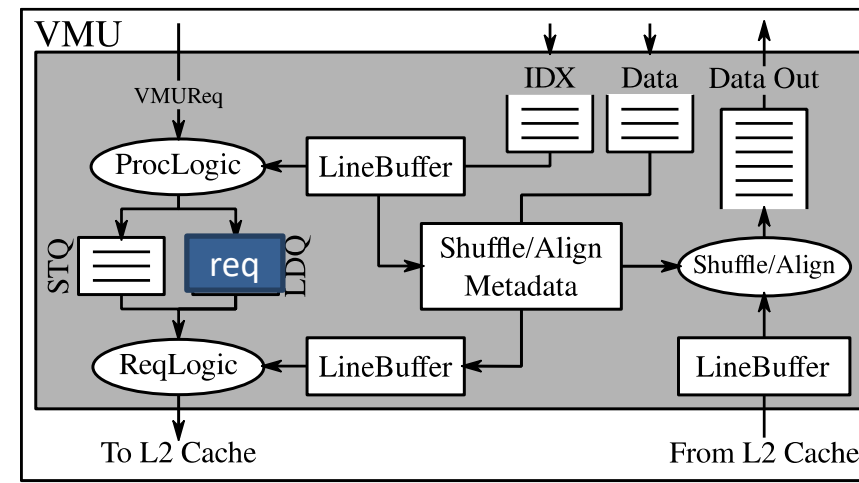
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
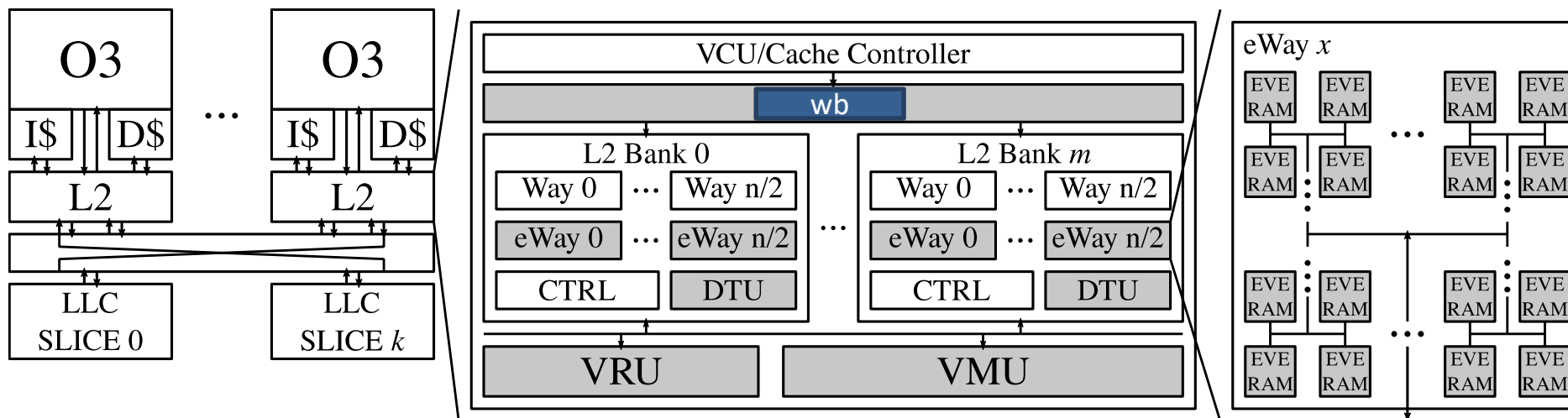
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 6 of 15
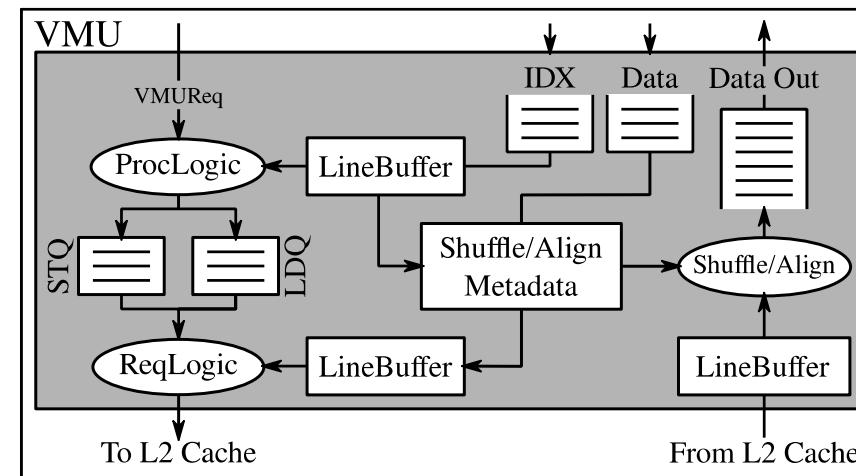
# EVE MICRO-ARCHITECTURE



```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```
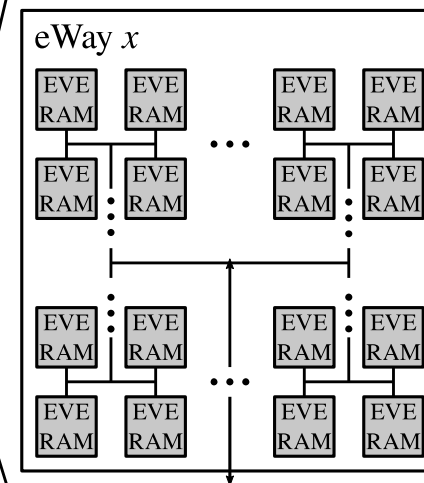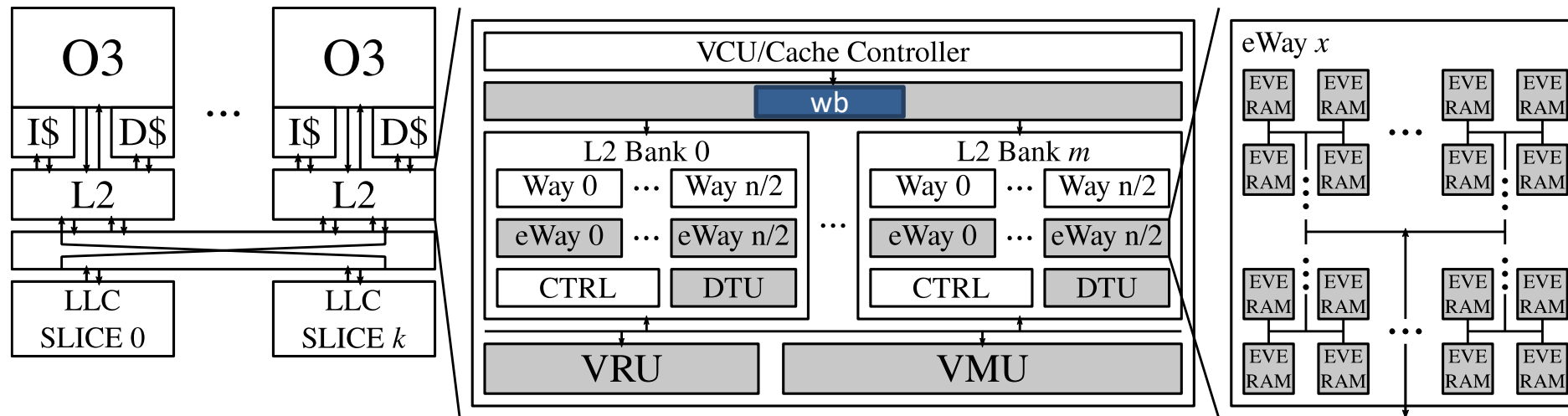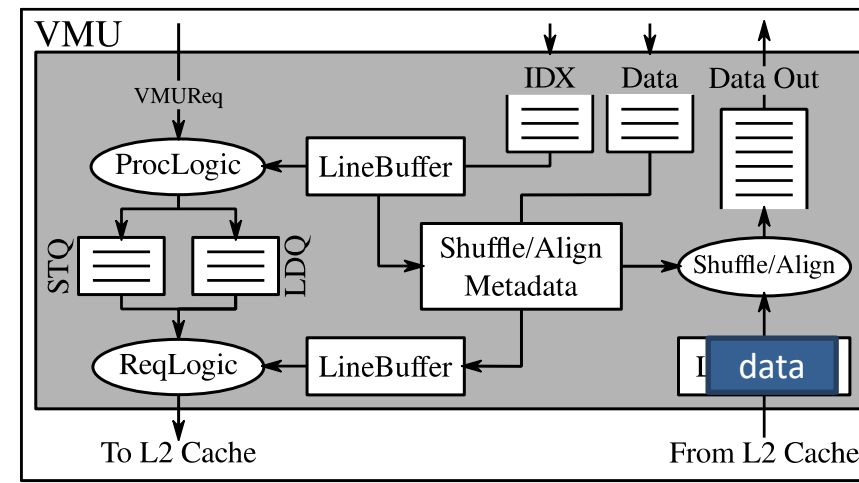
```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 6 of 15

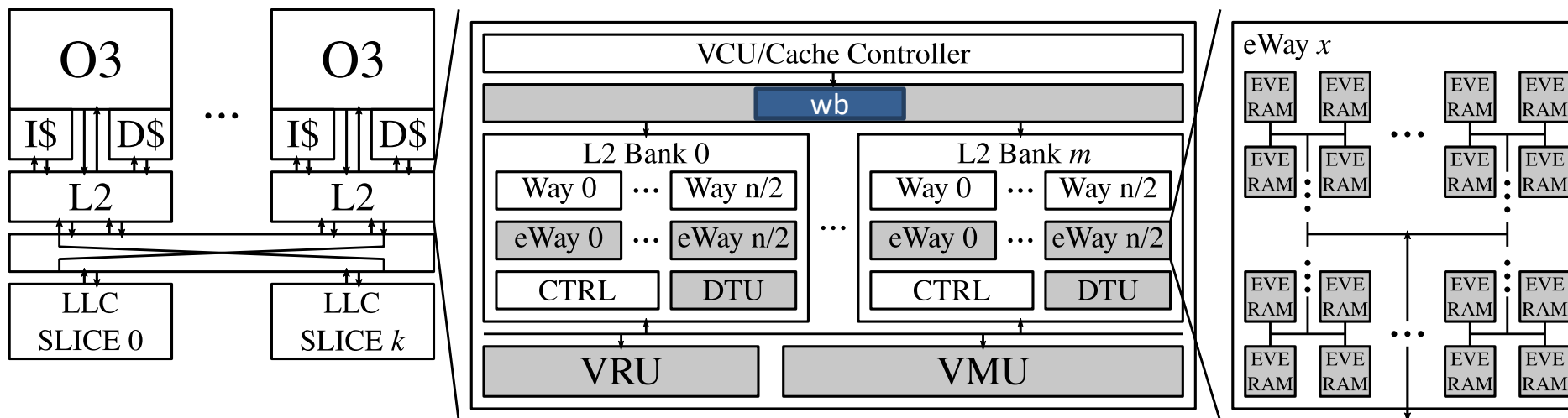Cornell University
Computer Systems Laboratory

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
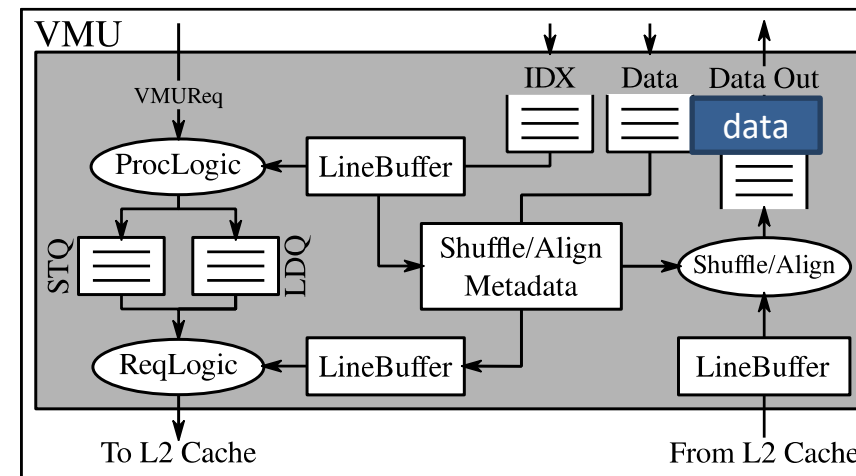
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
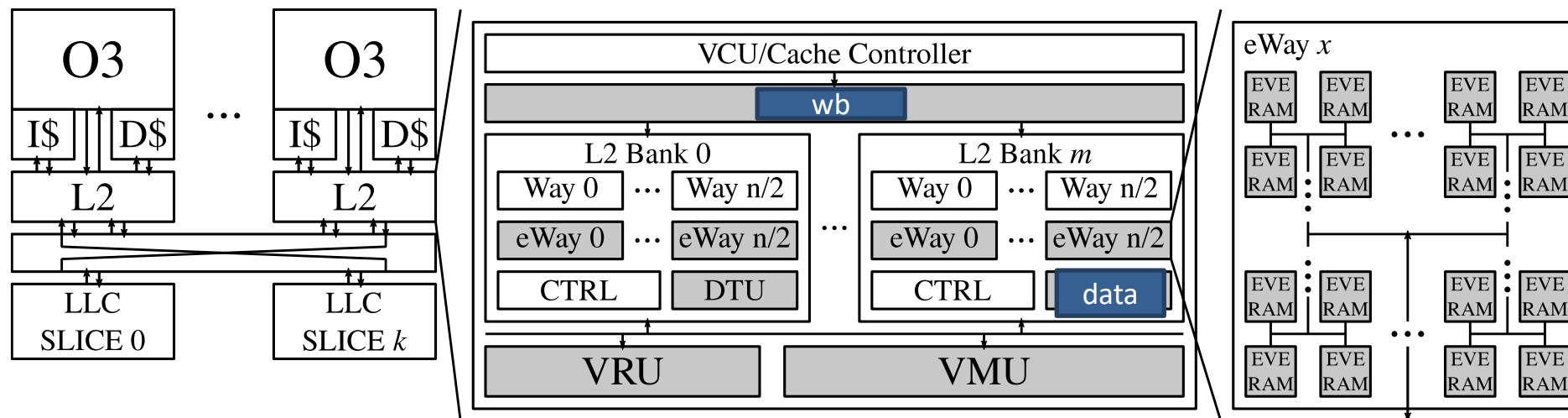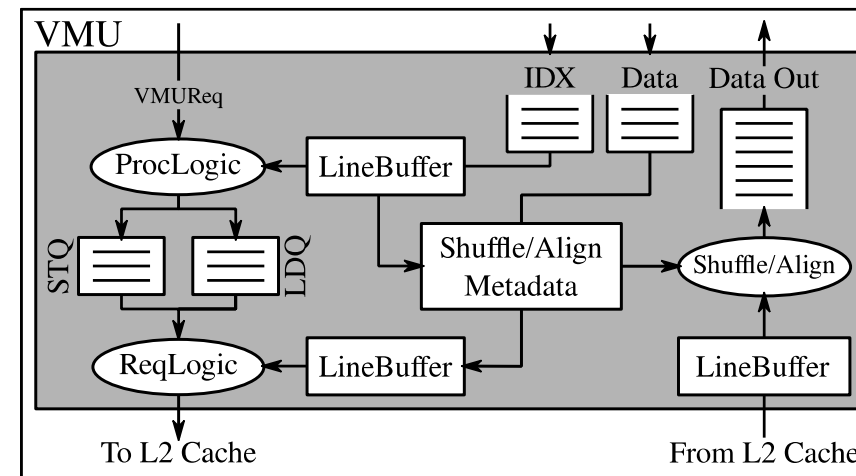
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
→   vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
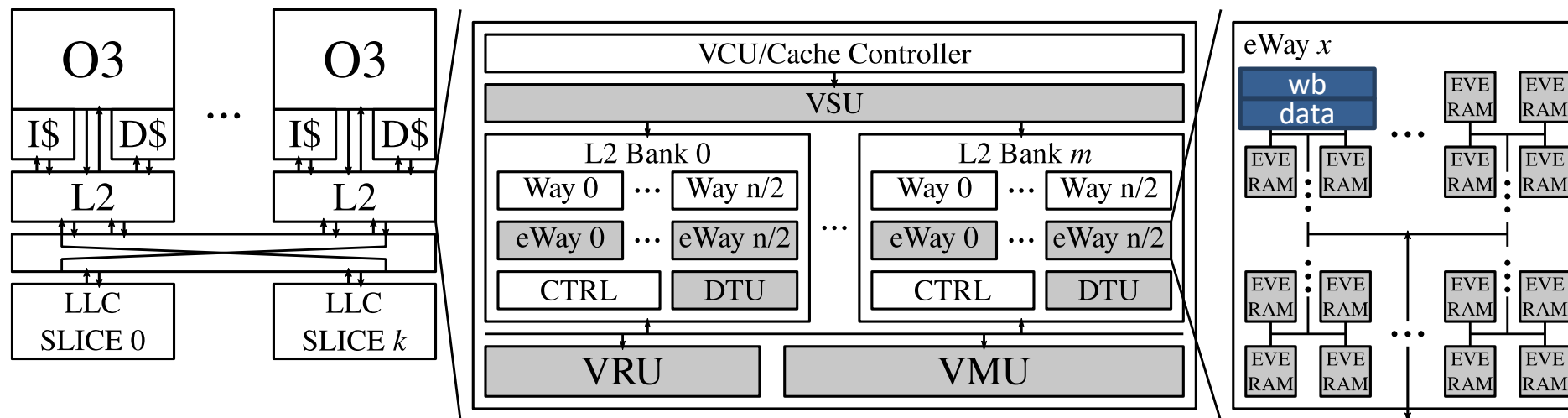
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();

}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 6 of 15

Cornell University
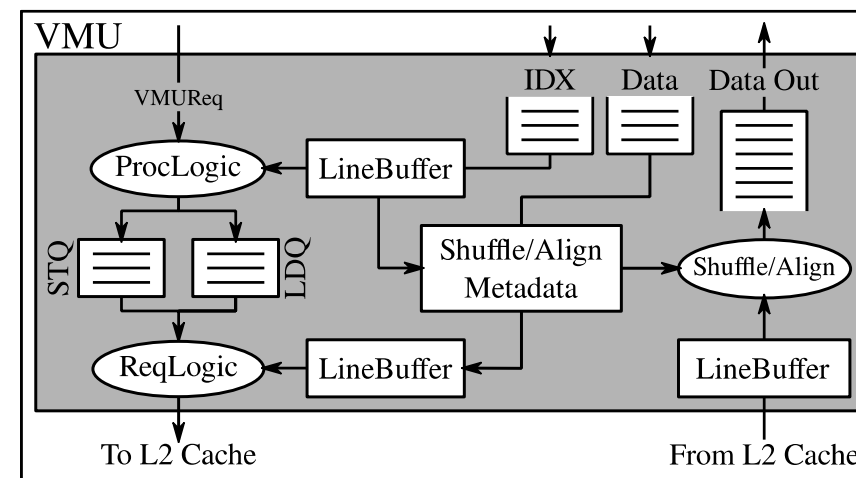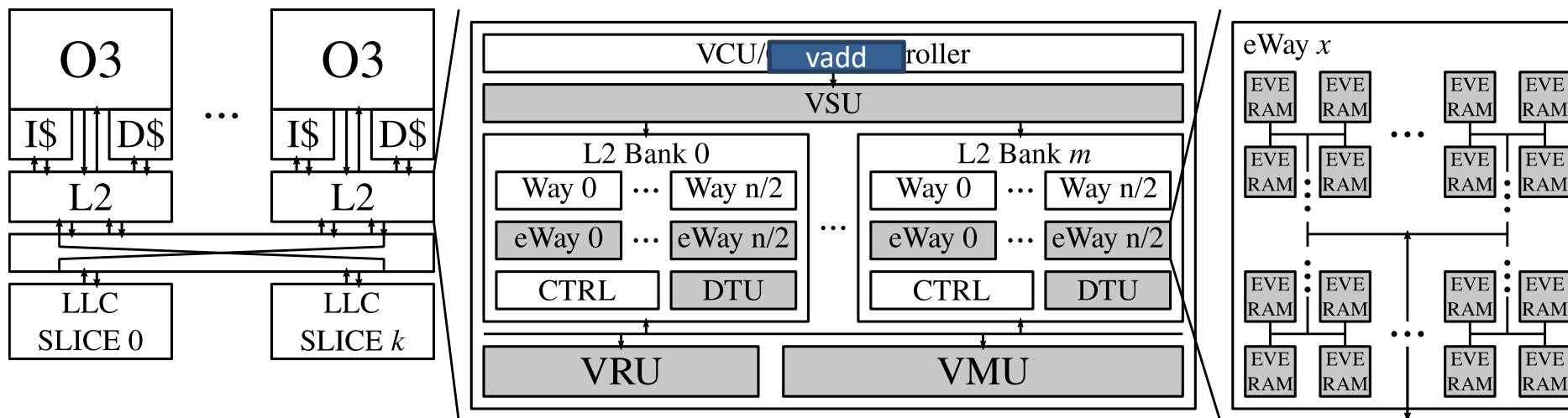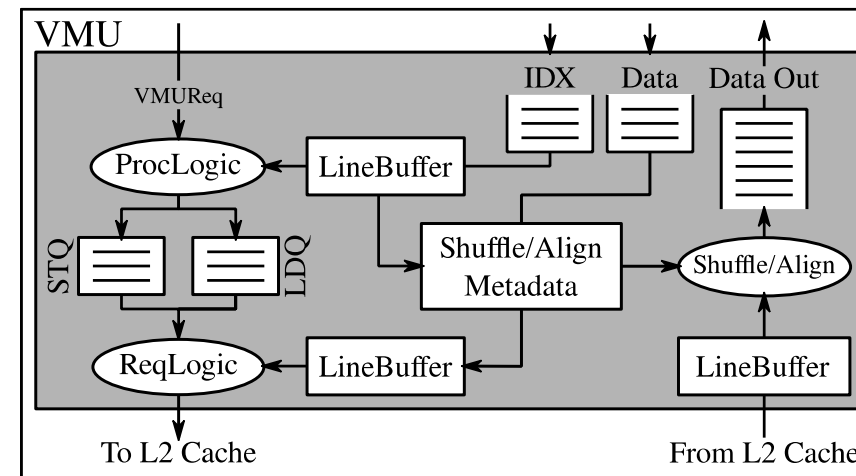Computer Systems Laboratory

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
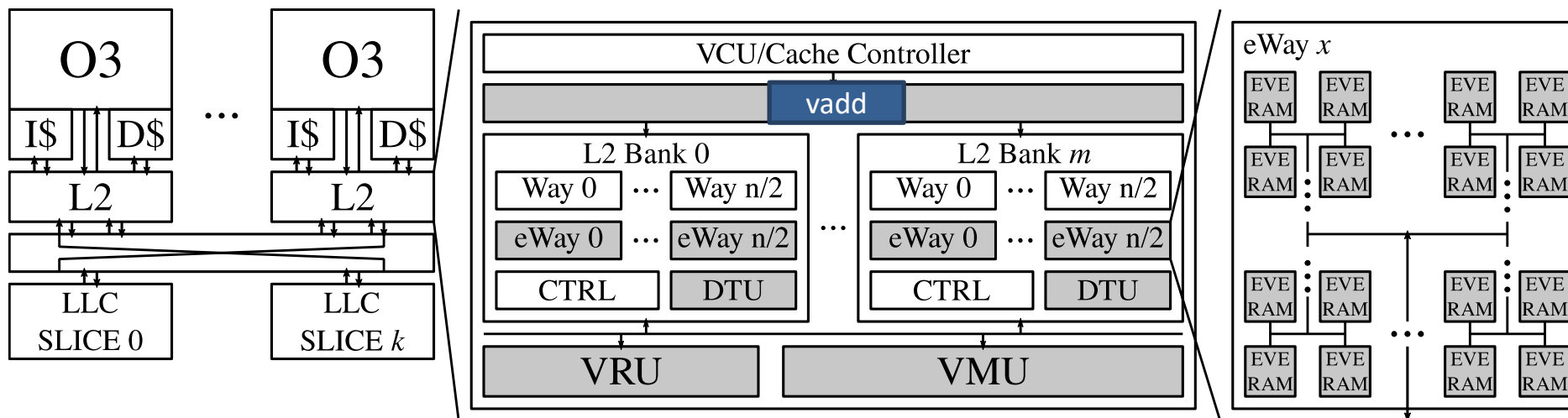
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
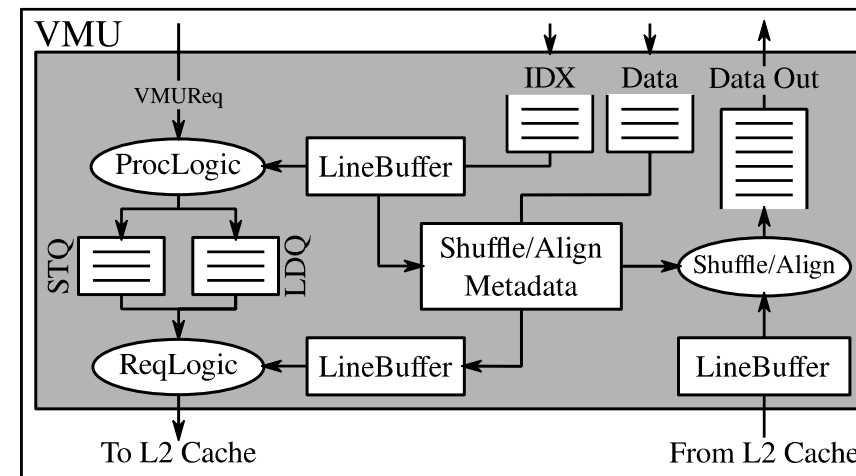
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
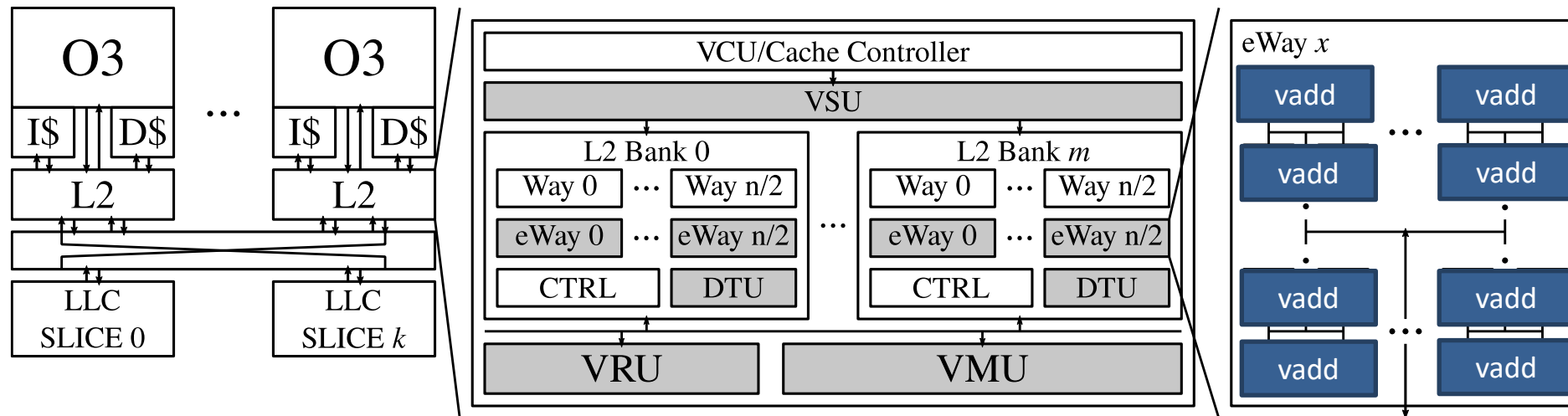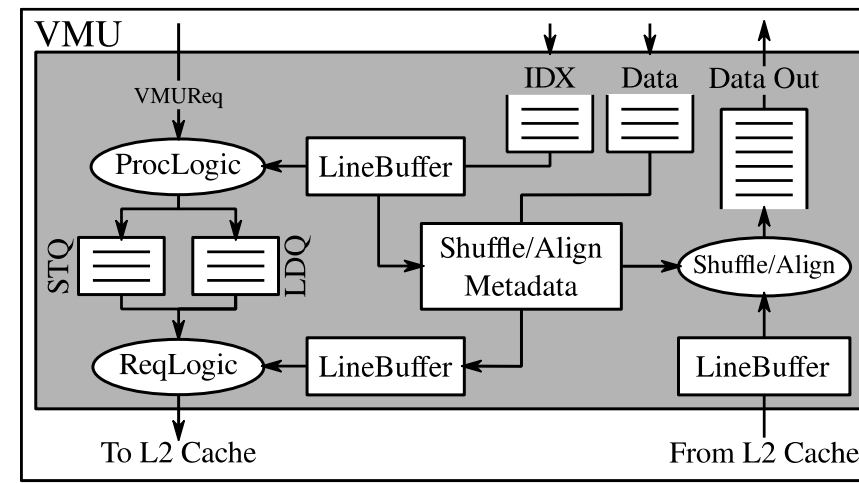
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();

}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Cornell University
Computer Systems Laboratory
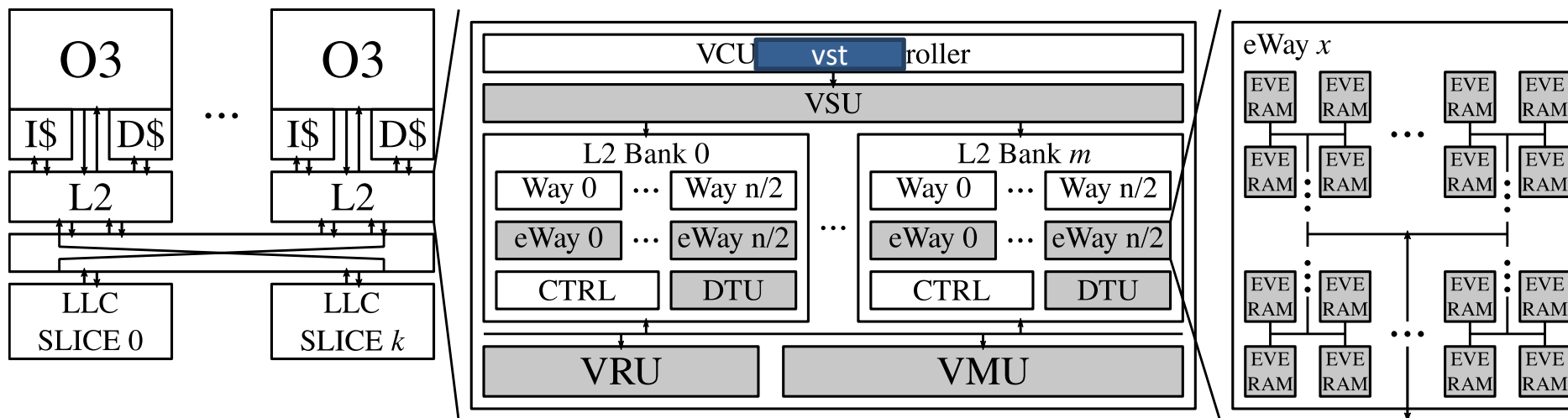
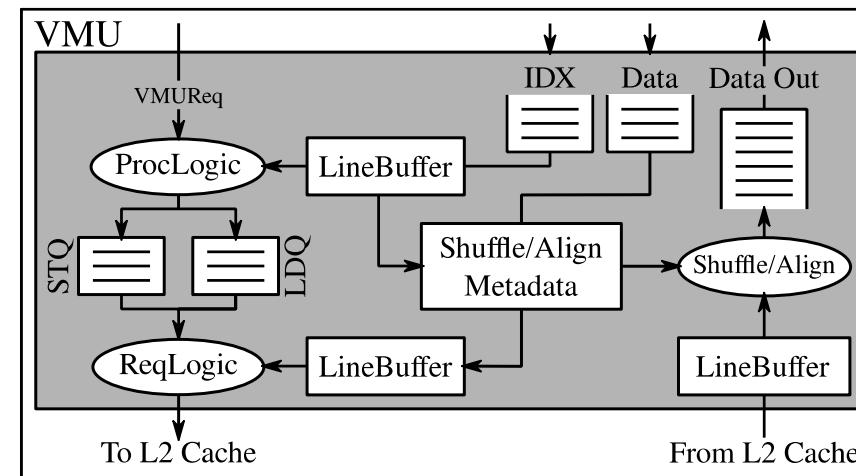Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();

}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
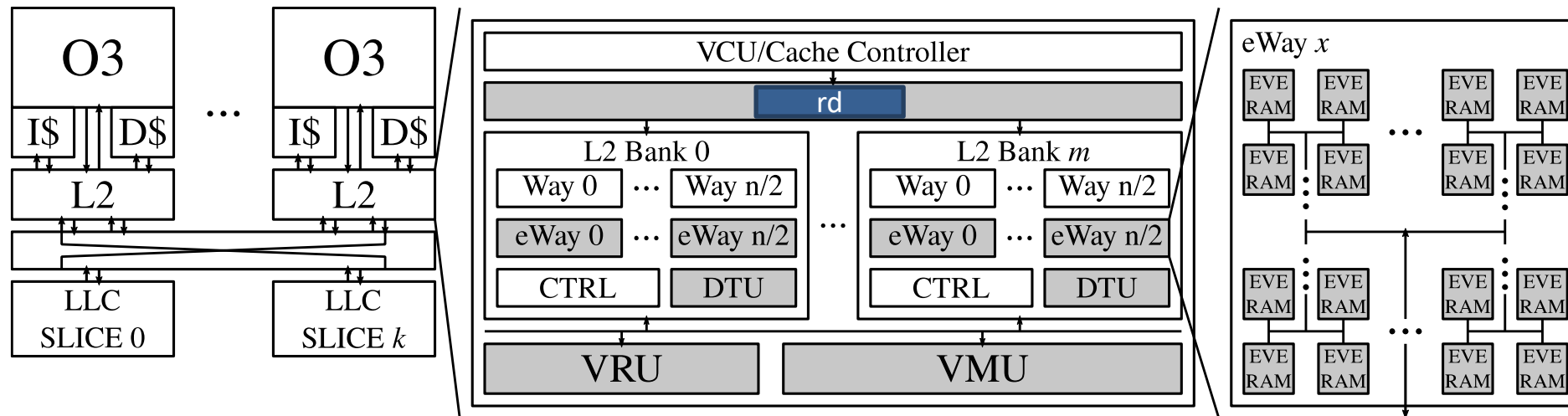
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();
}
```

Motivation • **Micro-Architecture** • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
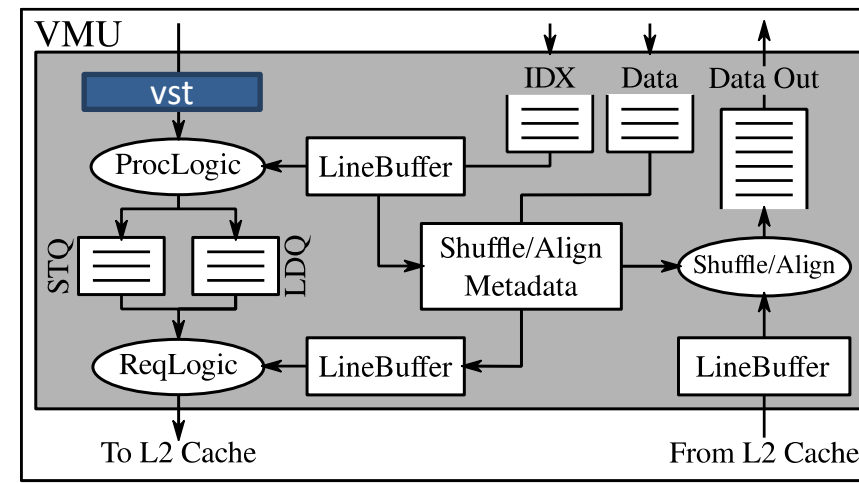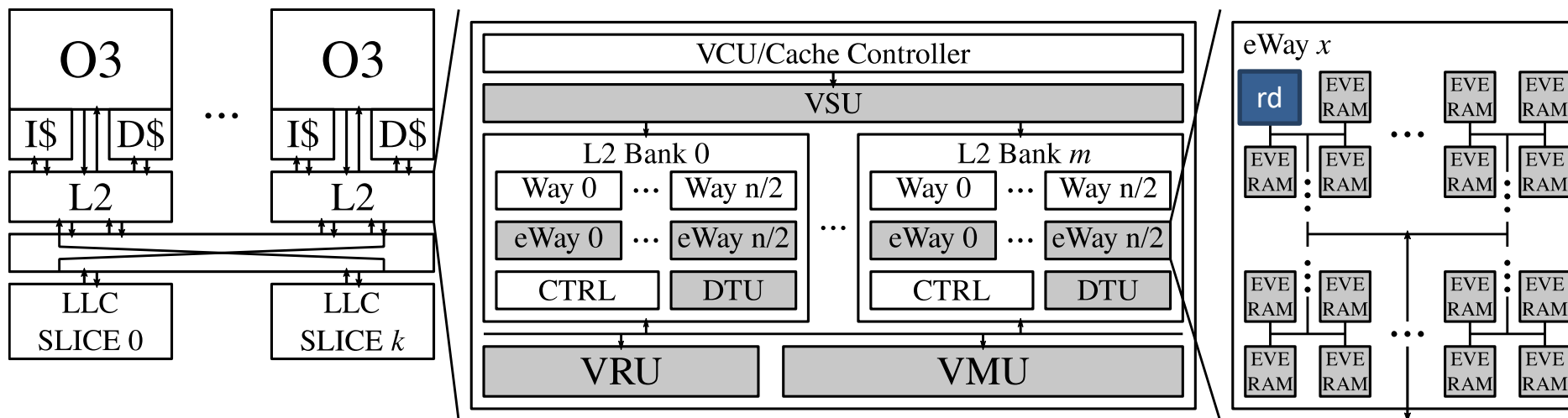
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • **Micro-Architecture** • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
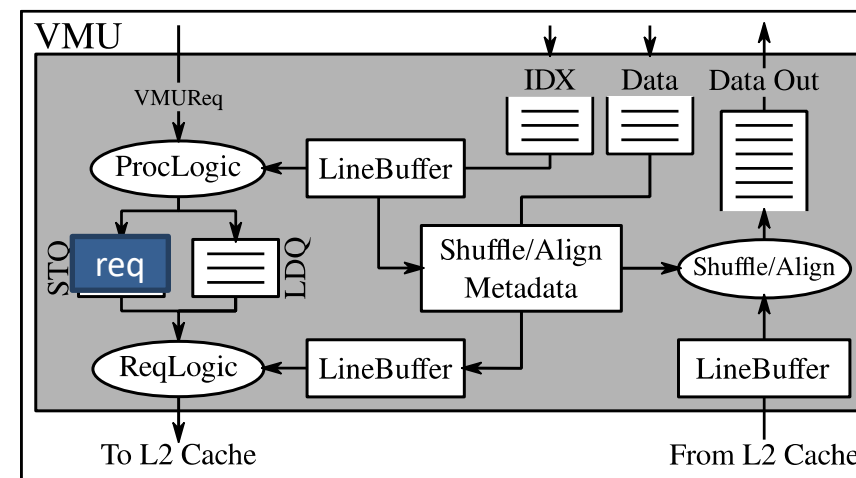
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
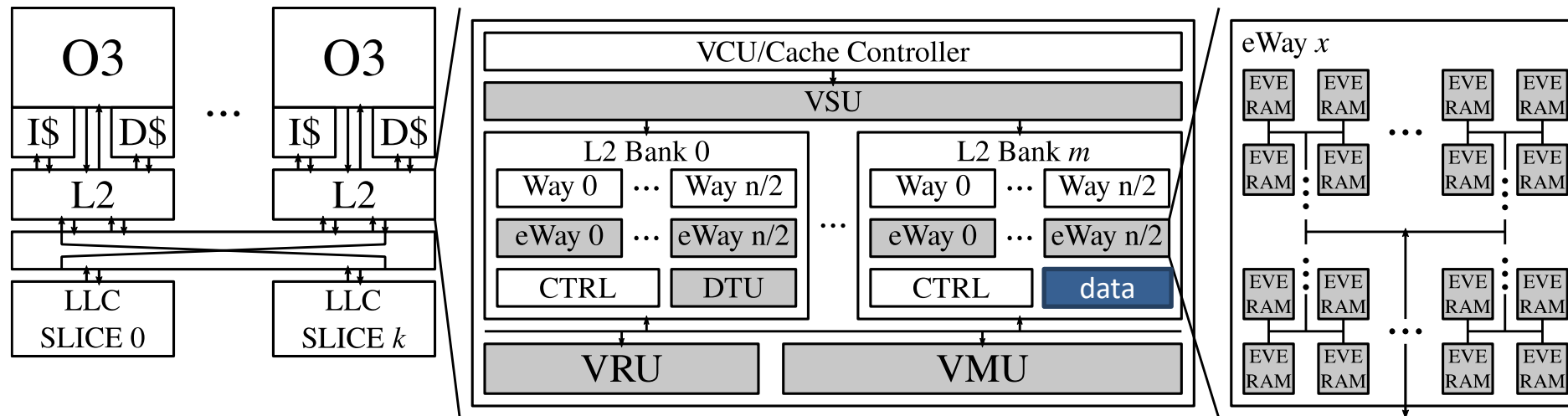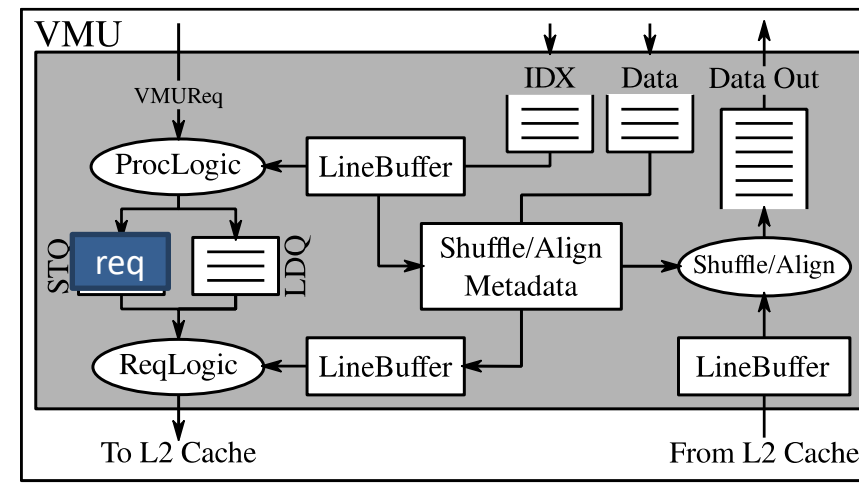
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
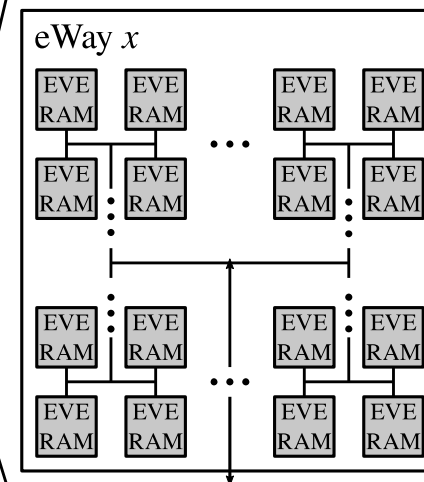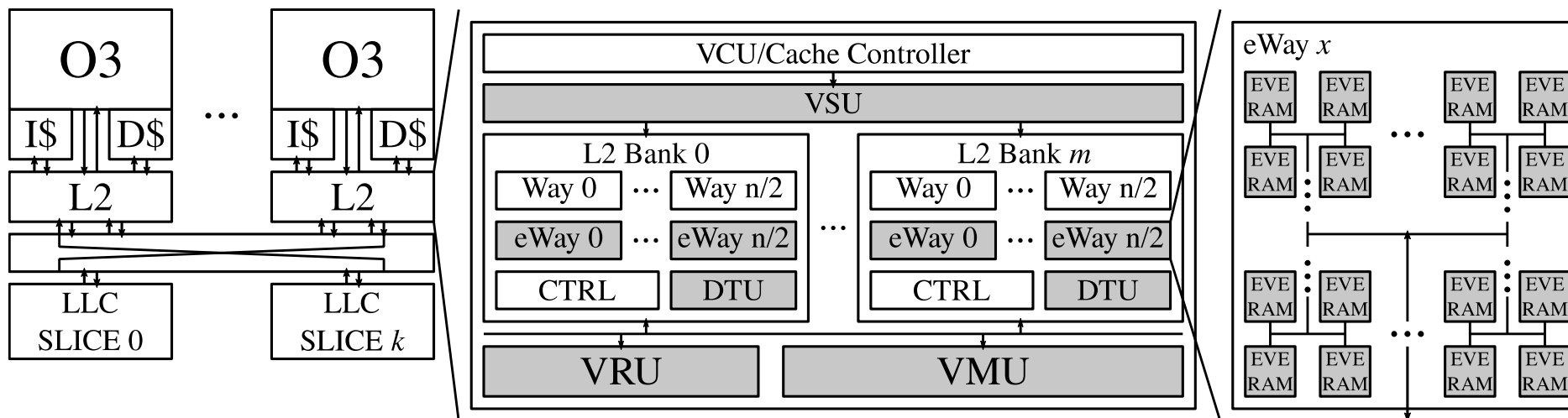
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 6 of 15

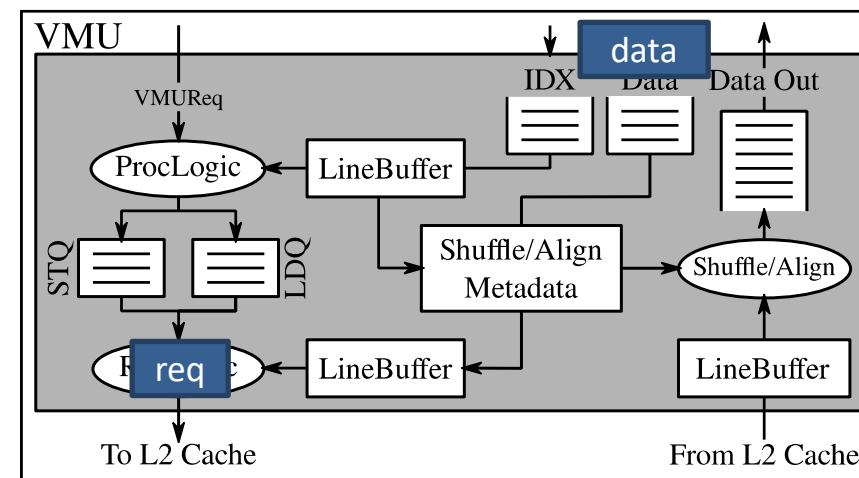Cornell University
Computer Systems Laboratory

```
void vvadd( int* c, int* a, int * b) {

  vstart();

  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }

  vend();
}
```

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
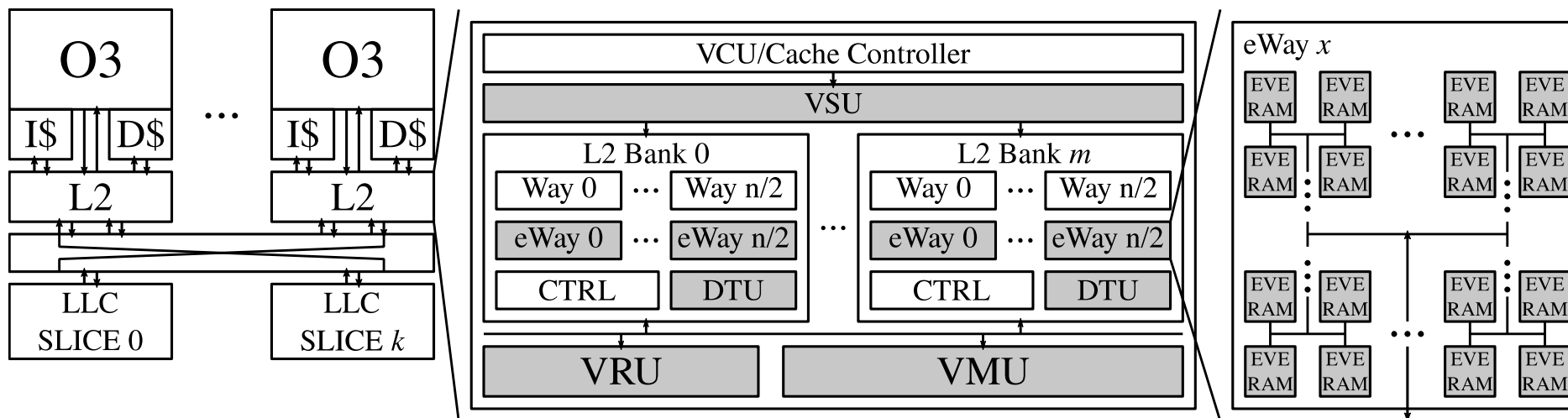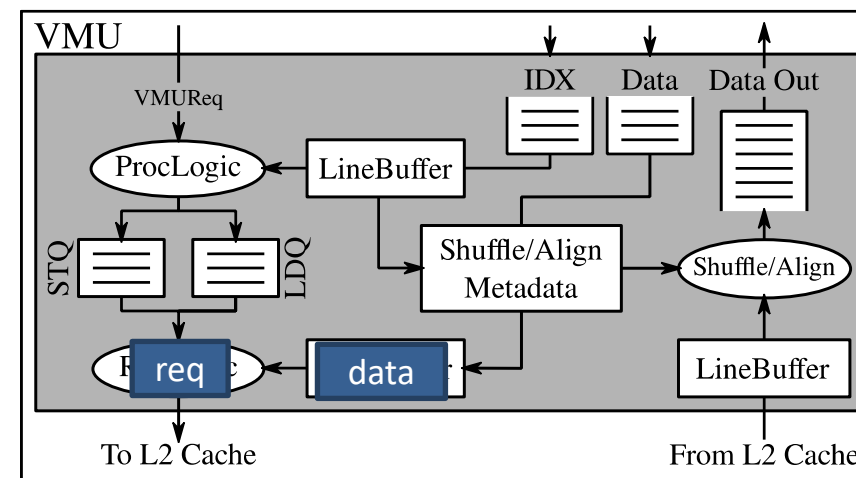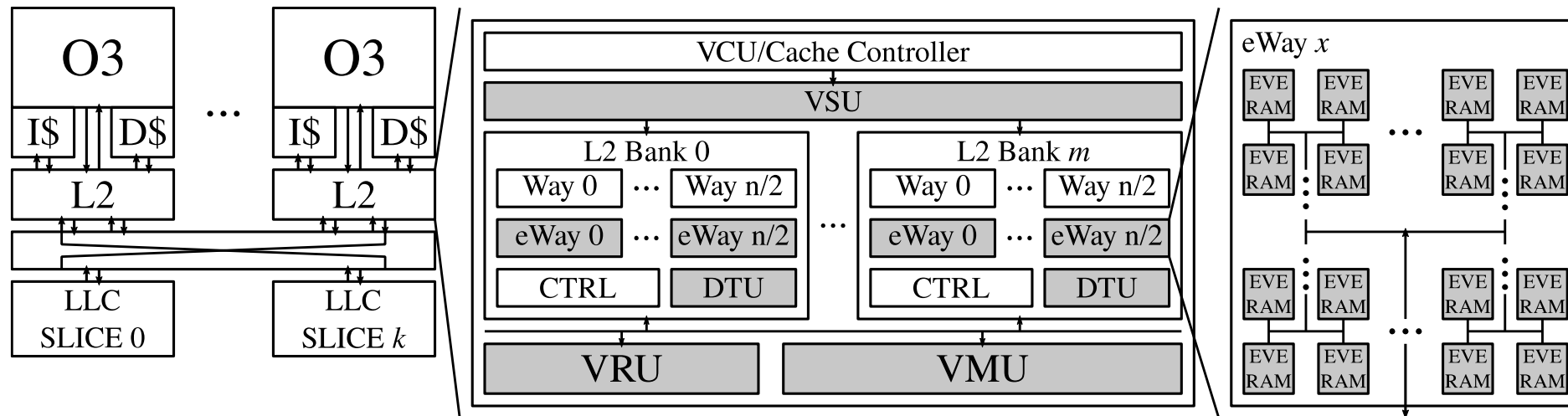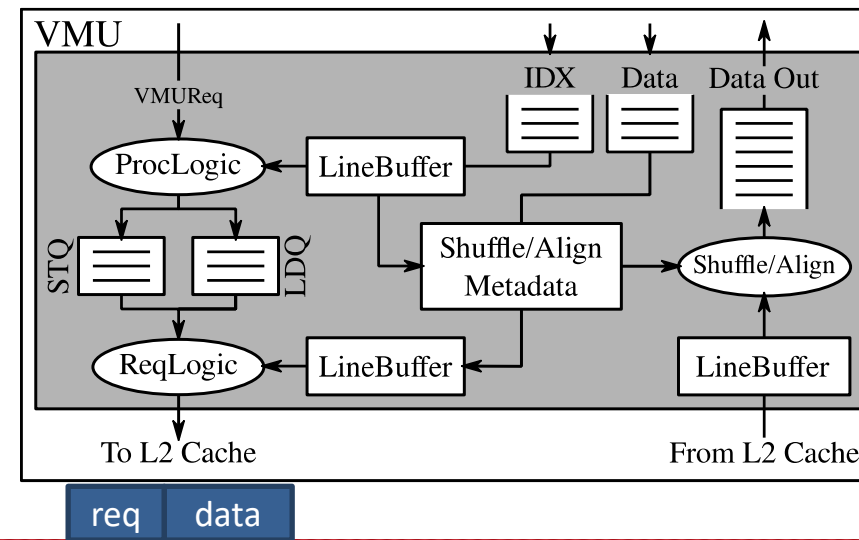
Page 6 of 15

```
void vvadd( int* c, int* a, int * b) {
  vstart();
  for (int i = 0; i < len; i += vsetvl(len)) {
    vld( v0, a + i );
    vld( v1, b + i );
    vadd.vv( v2, v0, v1 );
    vst( v2, c + i );
  }
  vend();
}
```

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 6 of 15

Motivation

- EVE Micro-Architecture

- **EVE Bit-Hybrid Execution Paradigm**

- EVE Micro-Programming & Circuits

- EVE Evaluation

Conclusion

**Assumptions:**

- 8-bit elements
- 16x16 SRAM array

**Definitions:**

- *Segment:* Elements are broken down into segments. A segment size can vary from 1 bit to 8 bit.
- *Parallelization factor:* Size of the segment (to be processed in parallel) in bits.
- Bit-Serial has parallelization factor of 1
- Bit-Parallel has parallelization factor of 8

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • **Bit-Hybrid Execution Paradigm** • Micro-Programming & Circuits • Evaluation • Conclusion

Page 7 of 15

# EVE Bit-Hybrid Execution Paradigm

## Assumptions:

- 8-bit elements
- 16x16 SRAM array

## Definitions:

- *Segment:* Elements are broken down into segments. A segment size can vary from 1 bit to 8 bit.

- *Parallelization factor:* Size of the segment (to be processed in parallel) in bits.

- Bit-Serial has parallelization factor of 1

- Bit-Parallel has parallelization factor of 8



## Balanced Utilization:
Perfect utilization of all bit-cells and in-situ ALUs in the SRAM array

## Column Under-Utilization:
Reducing the number of in-situ ALU in-favor of storage in the SRAM array

Cornell University
Computer Systems Laboratory
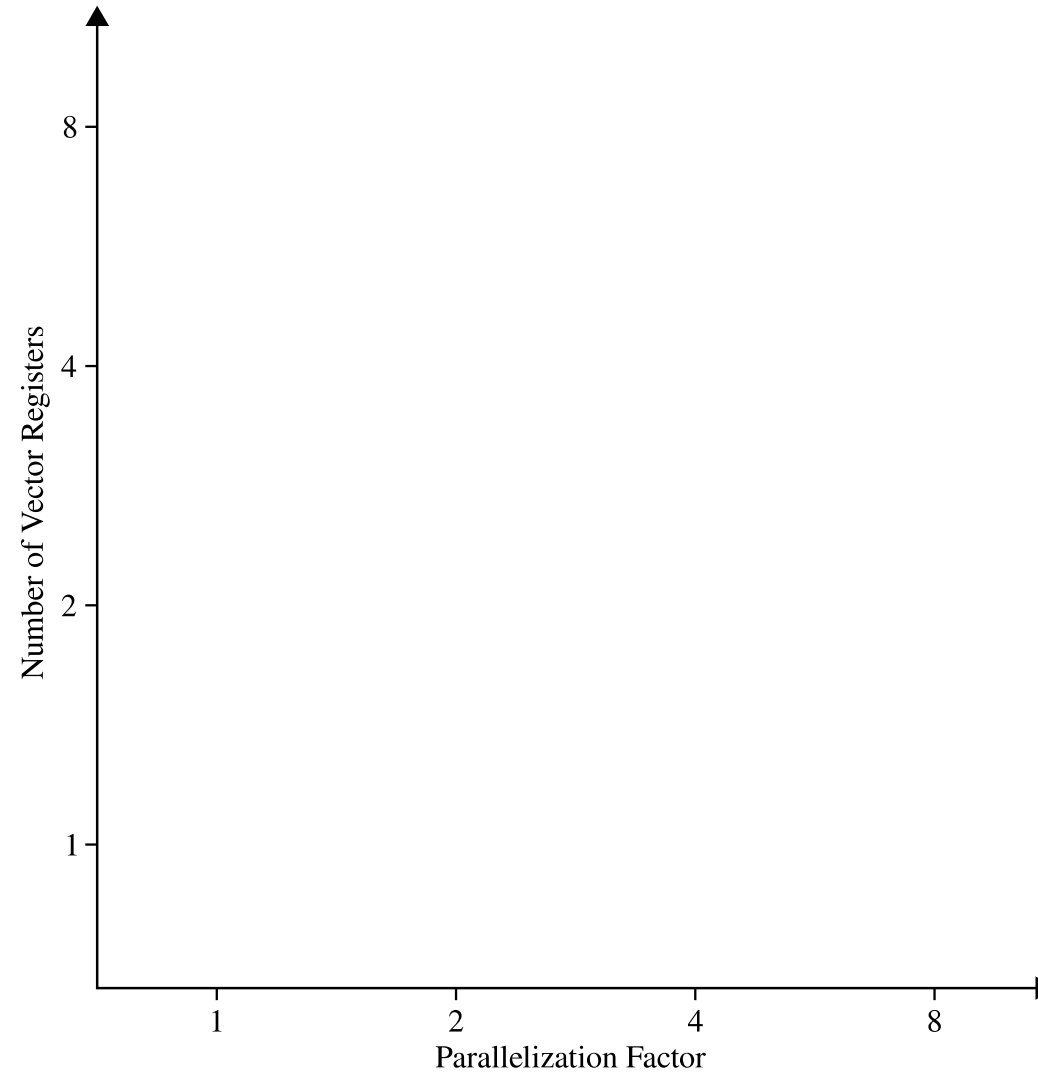
Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 8 of 15

# EVE Bit-Hybrid Execution Paradigm

**Assumptions:**

- 8-bit elements
- 16x16 SRAM array

**Definitions:**

- *Segment:* Elements are broken down into segments. A segment size can vary from 1 bit to 8 bit.
- *Parallelization factor:* Size of the segment (to be processed in parallel) in bits.
- Bit-Serial has parallelization factor of 1
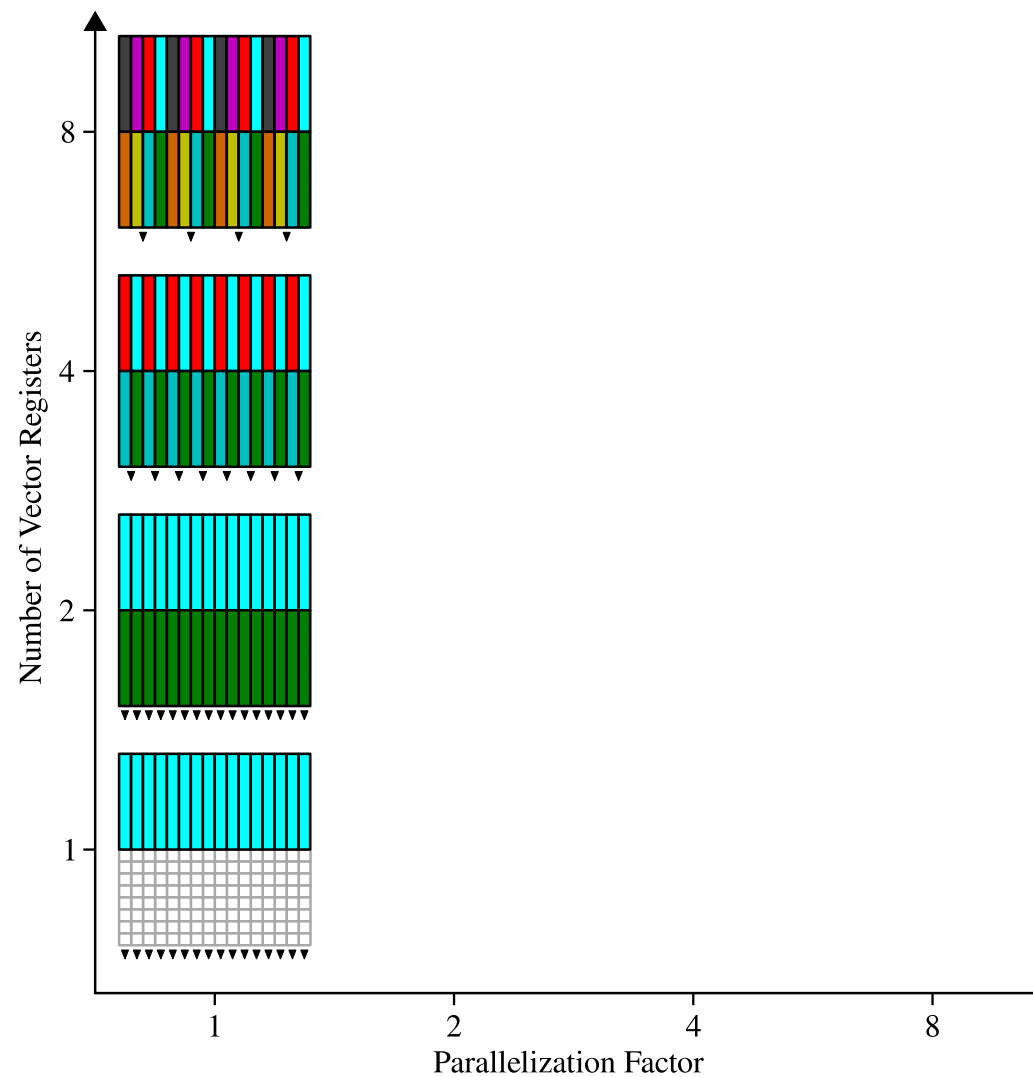- Bit-Parallel has parallelization factor of 8



Number of Vector Registers (y-axis)

Parallelization Factor (x-axis)

**Balanced Utilization:**
Perfect utilization of all bit-cells and in-situ ALUs in the SRAM array

**Column Under-Utilization:**
Reducing the number of in-situ ALU in-favor of storage in the SRAM array

**Row Under-Utilization:**
Under-utilized bit-cells in the SRAM array

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 8 of 15

## Assumptions:

- 8-bit elements
- 16x16 SRAM array

## Definitions:

- *Segment:* Elements are broken down into segments. A segment size can vary from 1 bit to 8 bit.
- *Parallelization factor:* Size of the segment (to be processed in parallel) in bits.
- Bit-Serial has parallelization factor of 1
- Bit-Parallel has parallelization factor of 8



## Balanced Utilization:
Perfect utilization of all bit-cells and in-situ ALUs in the SRAM array

## Column Under-Utilization:
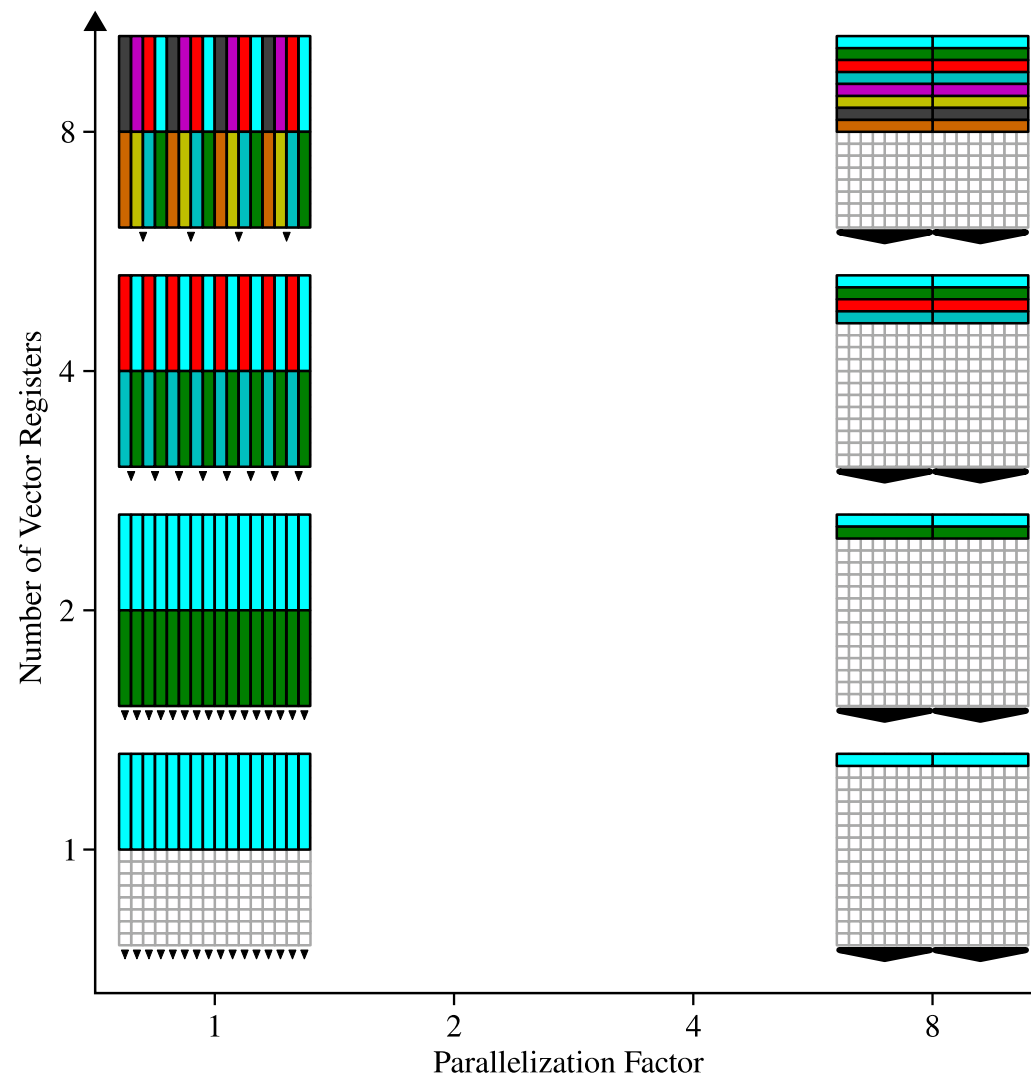Reducing the number of in-situ ALU in-favor of storage in the SRAM array

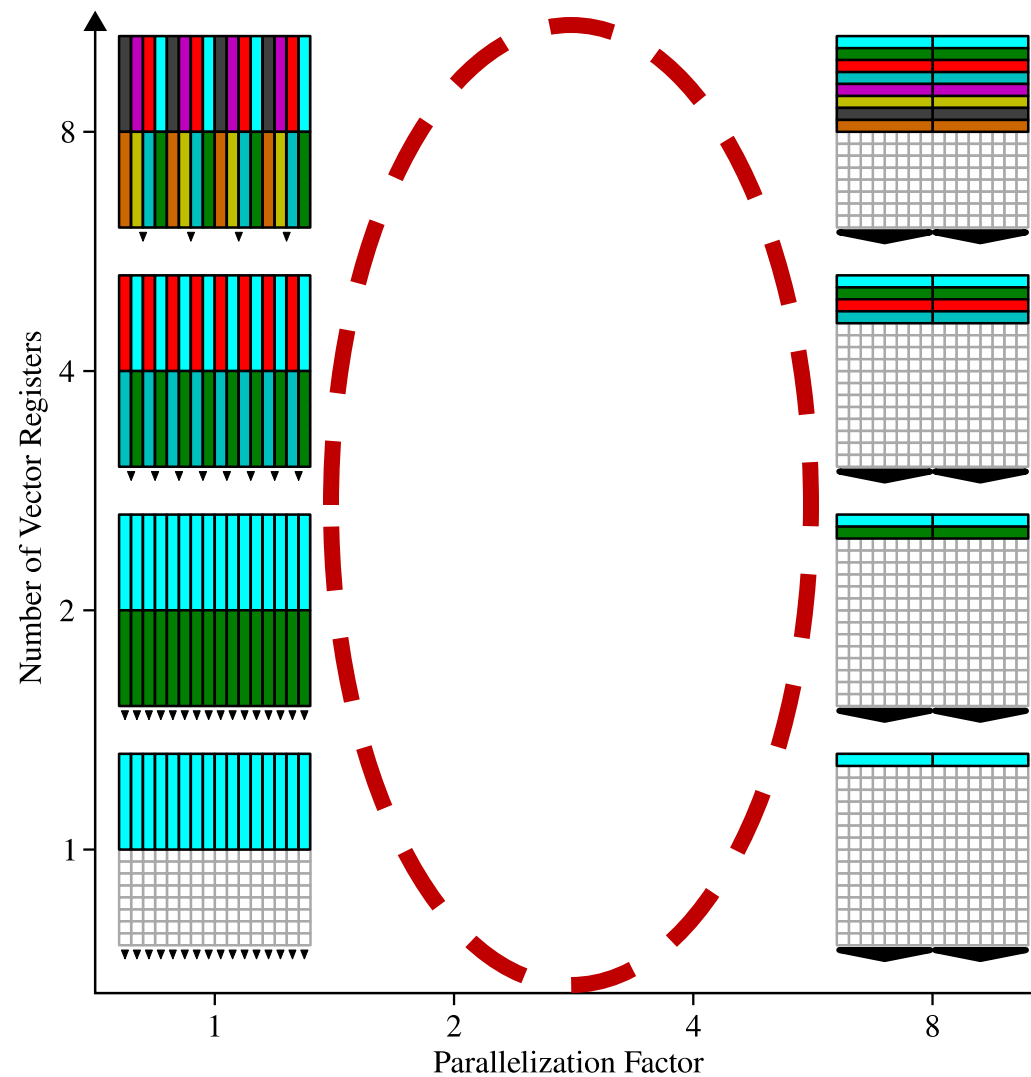## Row Under-Utilization:
Under-utilized bit-cells in the SRAM array

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • **Bit-Hybrid Execution Paradigm** • Micro-Programming & Circuits • Evaluation • Conclusion

Page 9 of 15

# EVE BIT-HYBRID EXECUTION PARADIGM

**Assumptions:**

- 8-bit elements
- 16x16 SRAM array

**Definitions:**

- *Segment:* Elements are broken down into segments. A segment size can vary from 1 bit to 8 bit.
- *Parallelization factor:* Size of the segment (to be processed in parallel) in bits.
- Bit-Serial has parallelization factor of 1
- Bit-Parallel has parallelization factor of 8



**Balanced Utilization:**
Perfect utilization of all bit-cells and in-situ ALUs in the SRAM array

**Column Under-Utilization:**
Reducing the number of in-situ ALU in-favor of storage in the SRAM array

**Row Under-Utilization:**
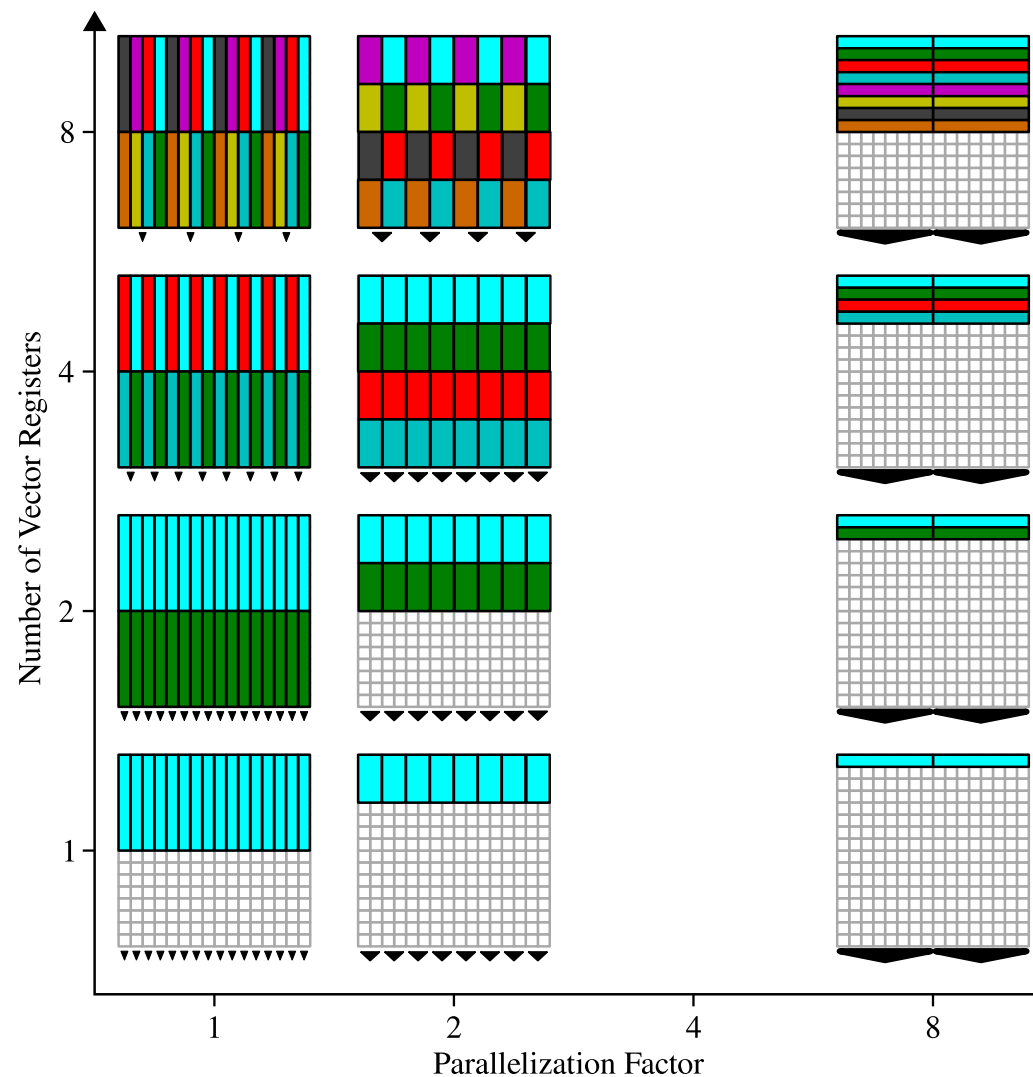Under-utilized bit-cells in the SRAM array

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • **Bit-Hybrid Execution Paradigm** • Micro-Programming & Circuits • Evaluation • Conclusion

Page 9 of 15

**Assumptions:**

- 8-bit elements
- 16x16 SRAM array

**Definitions:**

- *Segment:* Elements are broken down into segments. A segment size can vary from 1 bit to 8 bit.
- *Parallelization factor:* Size of the segment (to be processed in parallel) in bits.
- Bit-Serial has parallelization factor of 1
- Bit-Parallel has parallelization factor of 8



**Balanced Utilization:**
Perfect utilization of all bit-cells and in-situ ALUs in the SRAM array

**Column Under-Utilization:**
Reducing the number of in-situ ALU in-favor of storage in the SRAM array

**Row Under-Utilization:**
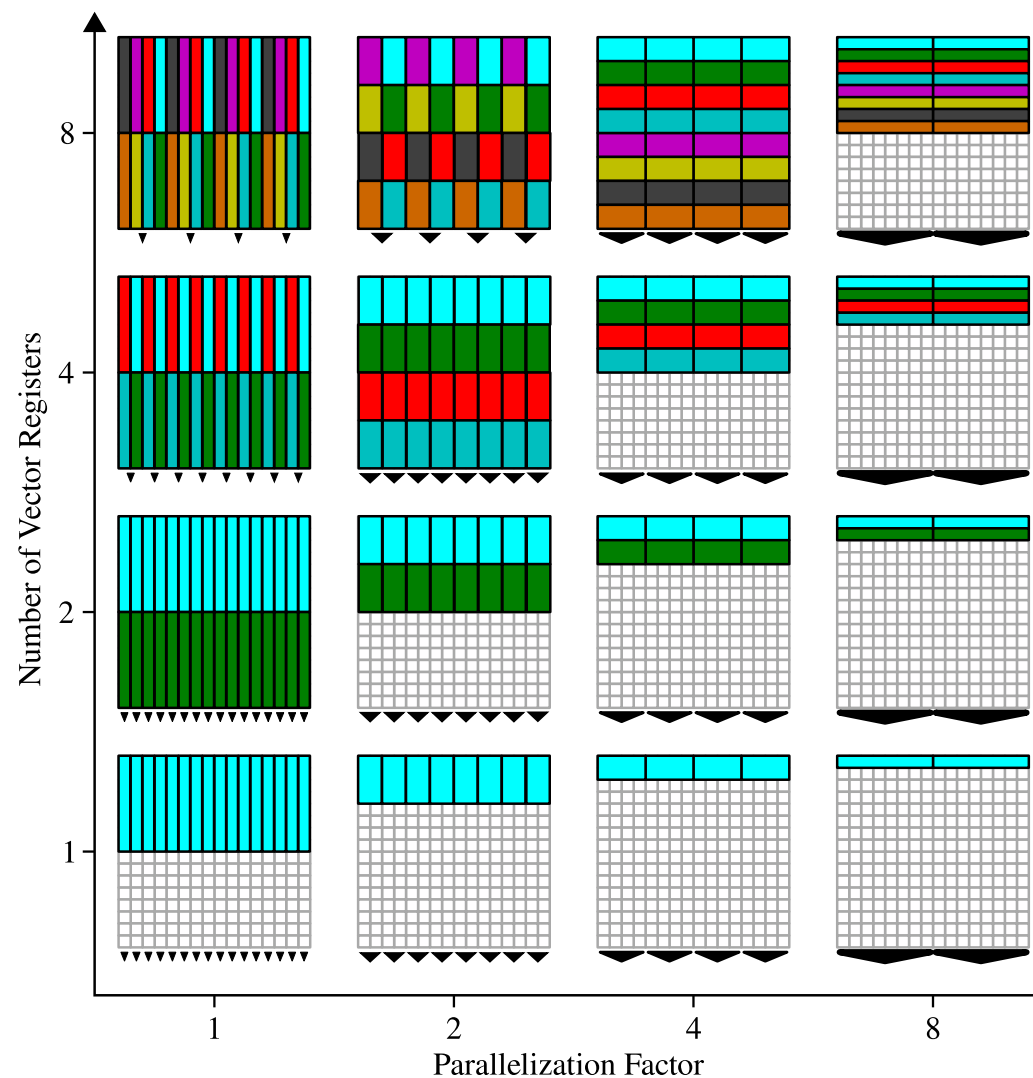Under-utilized bit-cells in the SRAM array

Motivation • Micro-Architecture • **Bit-Hybrid Execution Paradigm** • Micro-Programming & Circuits • Evaluation • Conclusion

Page 9 of 15

Cornell University
Computer Systems Laboratory

**Assumptions:**

- 8-bit elements
- 16x16 SRAM array

**Definitions:**

- *Segment:* Elements are broken down into segments. A segment size can vary from 1 bit to 8 bit.

- *Parallelization factor:* Size of the segment (to be processed in parallel) in bits.

- Bit-Serial has parallelization factor of 1

- Bit-Parallel has parallelization factor of 8



**Balanced Utilization:**
Perfect utilization of all bit-cells and in-situ ALUs in the SRAM array

**Column Under-Utilization:**
Reducing the number of in-situ ALU in-favor of storage in the SRAM array

**Row Under-Utilization:**
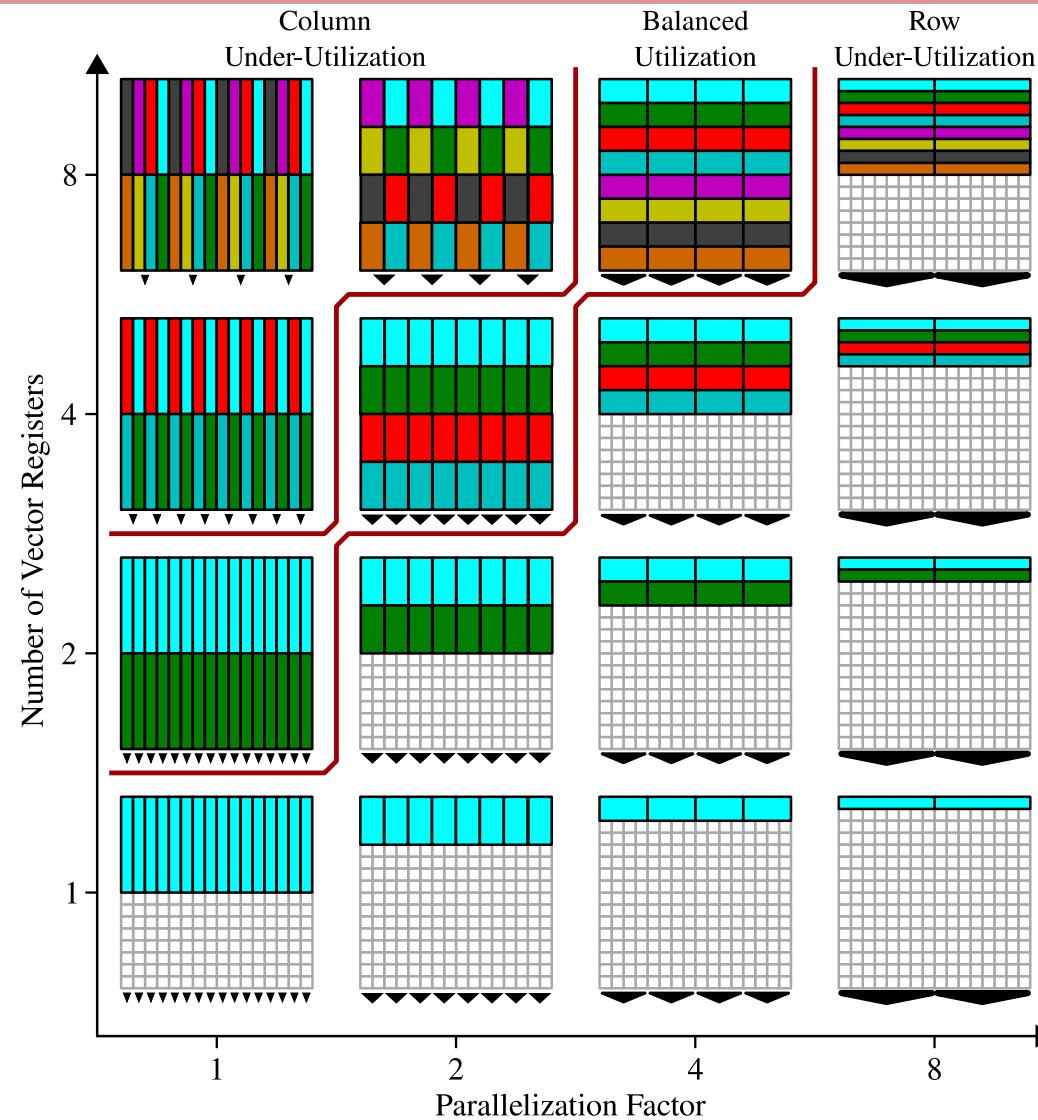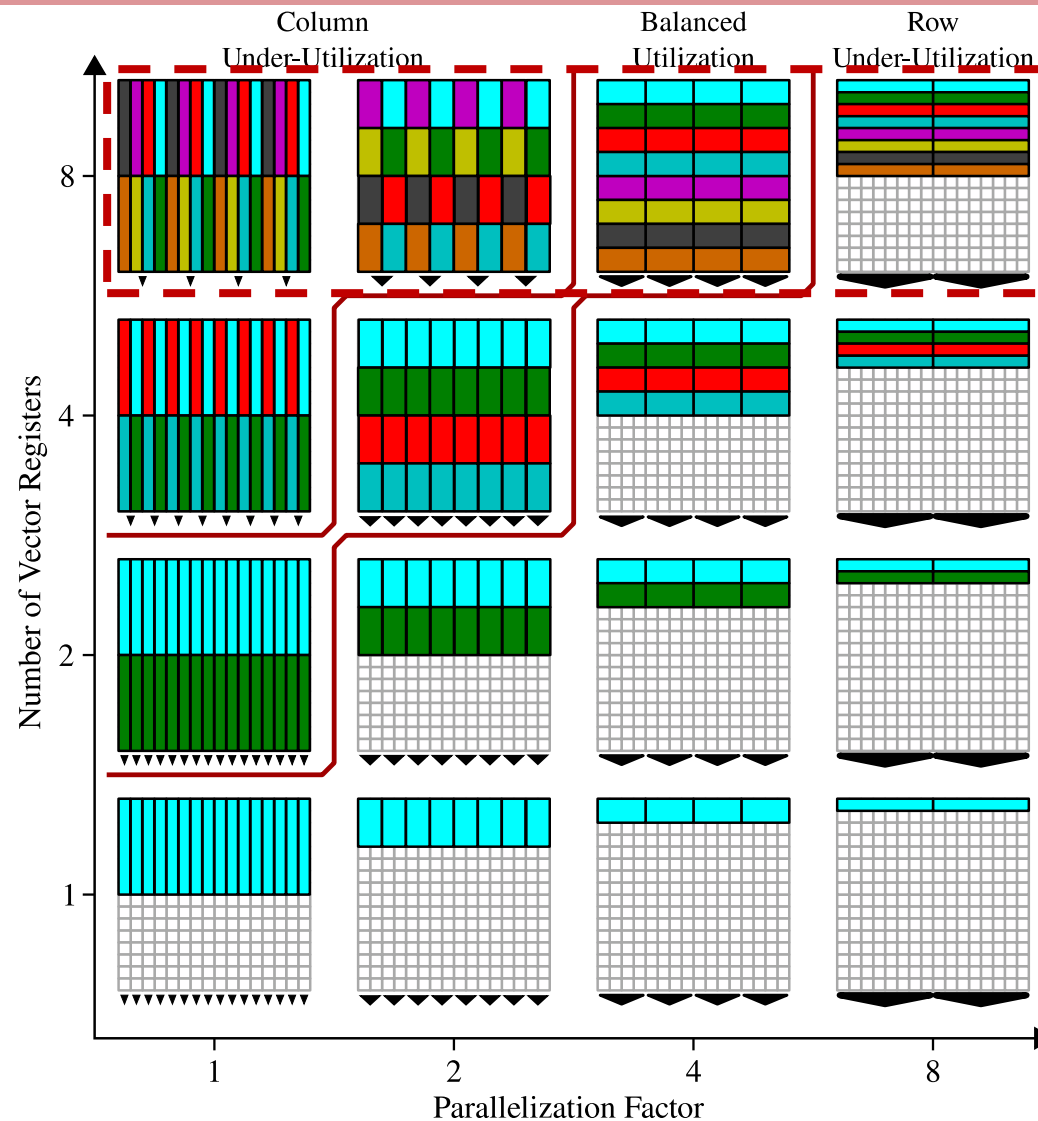Under-utilized bit-cells in the SRAM array

Motivation • Micro-Architecture • **Bit-Hybrid Execution Paradigm** • Micro-Programming & Circuits • Evaluation • Conclusion

Cornell University
Computer Systems Laboratory

Page 10 of 15

**Assumptions:**

- 8-bit elements
- 16x16 SRAM array

**Definitions:**

- *Segment:* Elements are broken down into segments. A segment size can vary from 1 bit to 8 bit.
- *Parallelization factor:* Size of the segment (to be processed in parallel) in bits.
- Bit-Serial has parallelization factor of 1
- Bit-Parallel has parallelization factor of 8

**Balanced Utilization:**
Perfect utilization of all bit-cells and in-situ ALUs in the SRAM array

**Column Under-Utilization:**
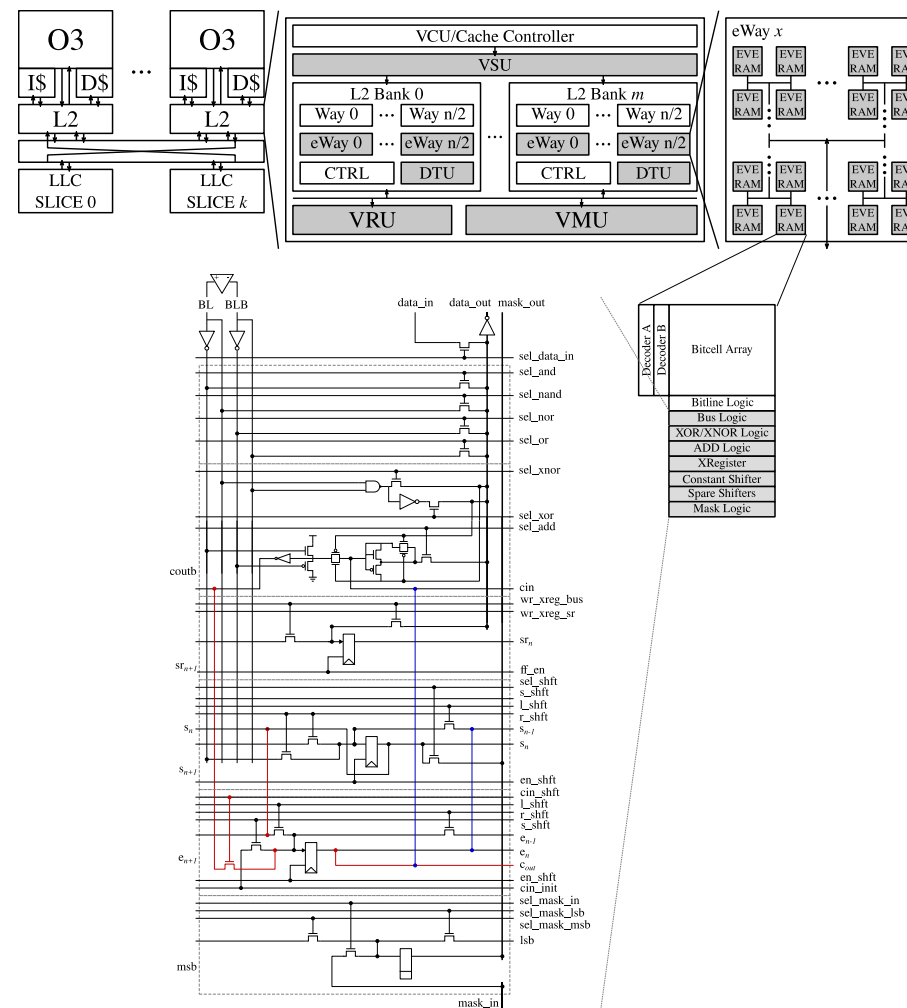Reducing the number of in-situ ALU in-favor of storage in the SRAM array

**Row Under-Utilization:**
Under-utilized bit-cells in the SRAM array

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • **Bit-Hybrid Execution Paradigm** • Micro-Programming & Circuits • Evaluation • Conclusion
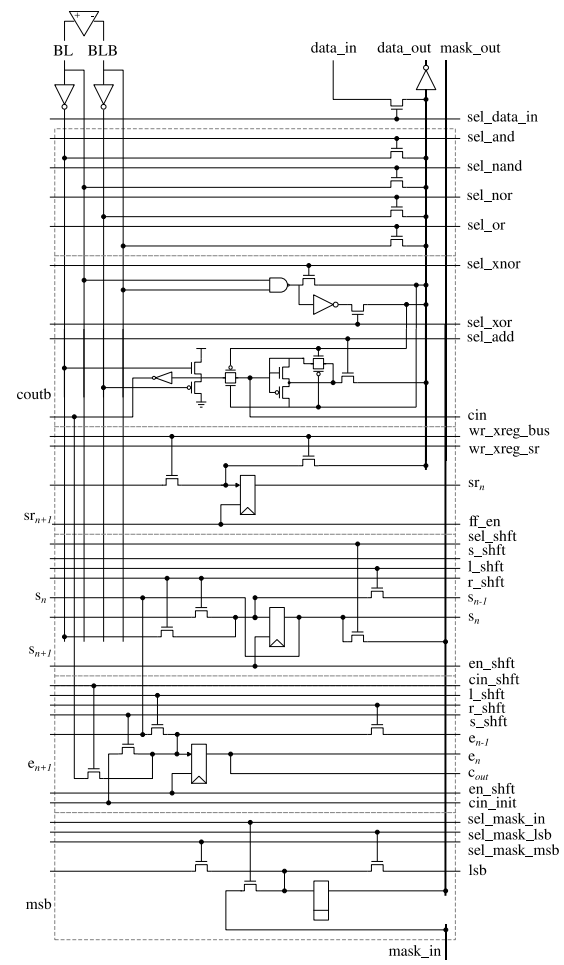
Page 10 of 15

Motivation

- EVE Micro-Architecture

- EVE Bit-Hybrid Execution Paradigm

- **EVE Micro-Programming & Circuits**

- EVE Evaluation

Conclusion
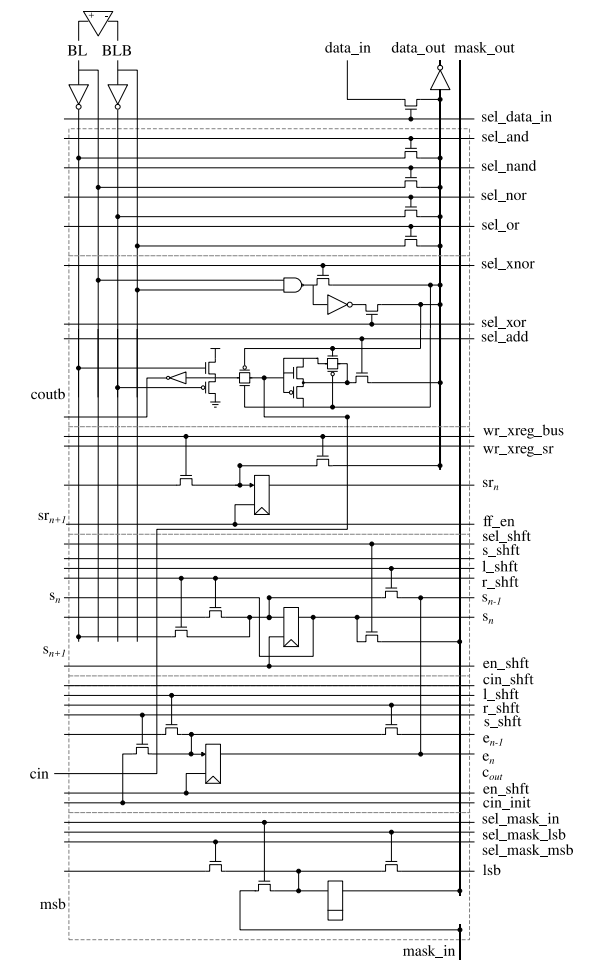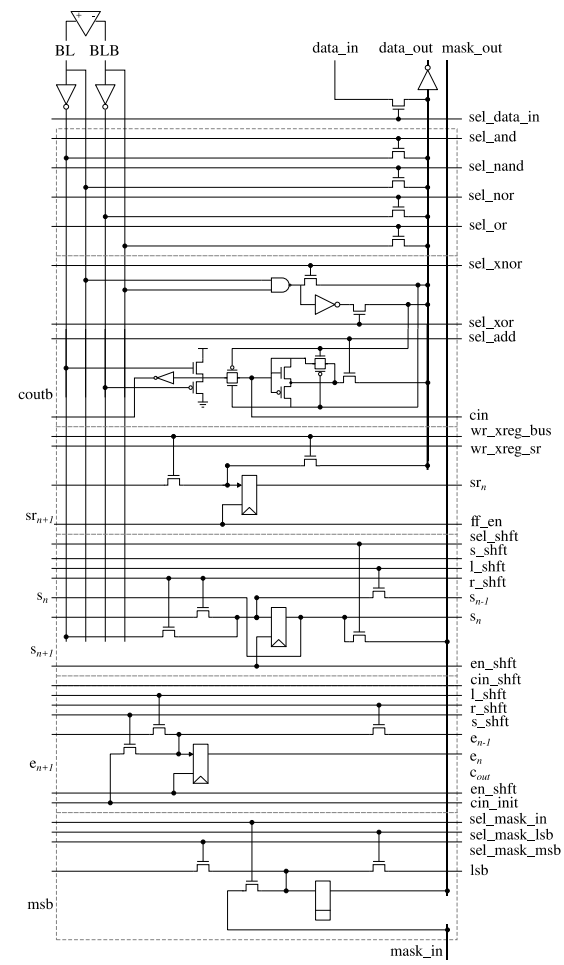
```
1   init    cnt_0, N
  loop:
1   blc     addr_a, addr_b
    ; decr  cnt_0
2   wb      addr_c, add
    ; bnz   cnt_0, loop
```
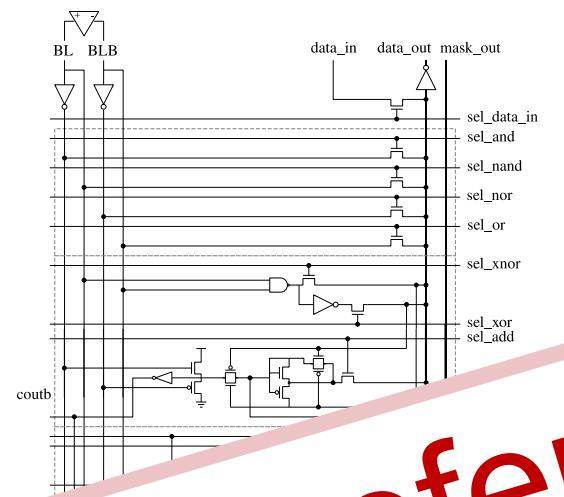
MSB

. . . . .

LSB

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 11 of 33

MSB                    . . . . .                    LSB

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Cornell University
Computer Systems Laboratory

Page 11 of 33

**Motivation**

■ **EVE Micro-Architecture**

■ **EVE Bit-Hybrid Execution Paradigm**

■ **EVE Micro-Programming & Circuits**

■ **EVE Evaluation**

**Conclusion**

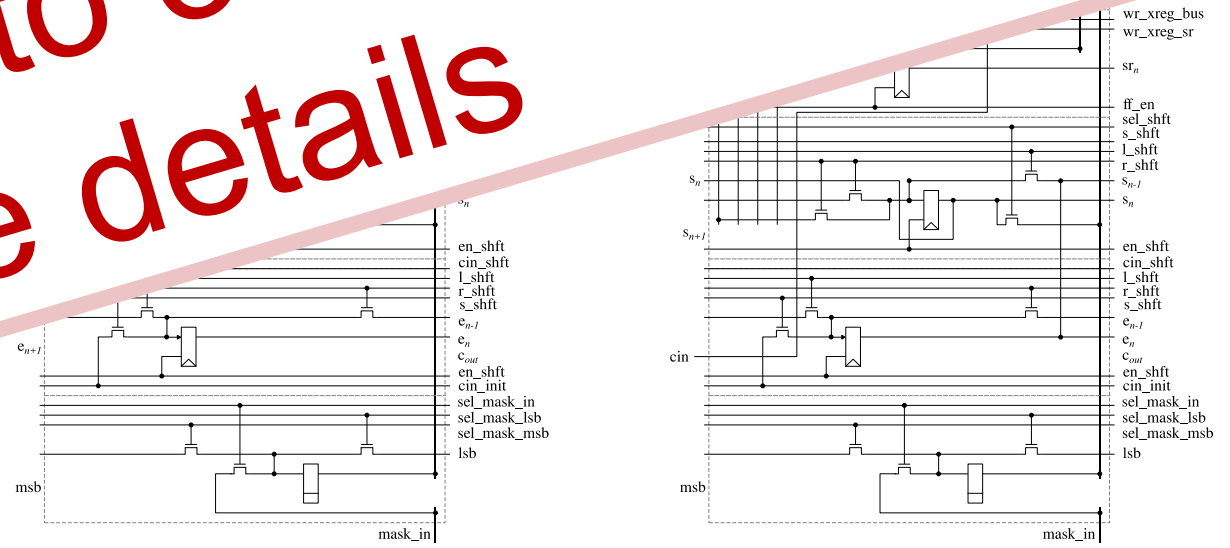- Utilized a simplified version of EVE circuits to estimate area and cycle-time of EVE SRAM:
  - Area: OpenRAM-generated layout
  - Cycle-time: Extracted SPICE-netlist
- Area for other components (VCU, VSU, DTU, VRU, and VMU) was estimated through SRAM-array equivalence
- EVE incurs around 12% area overhead over O3 core
- EVE incurs no cycle-time overhead for parallelization factors of eight or less



| EVE-1 | EVE-2 | EVE-4 | EVE-8 | EVE-16 | EVE-32 |
|---|---|---|---|---|---|
| 10.1% | 11.7% | 11.7% | 11.7% | 11.7% | 11.0% |

Area Overhead Over O3

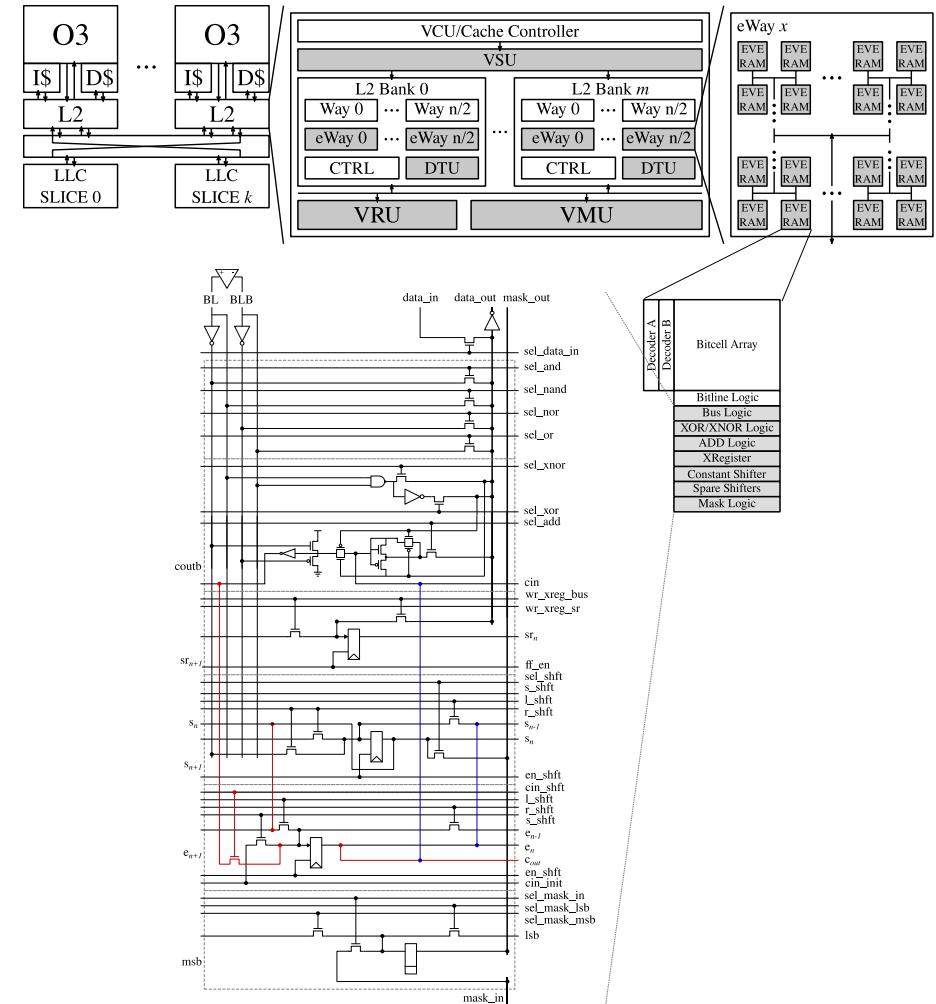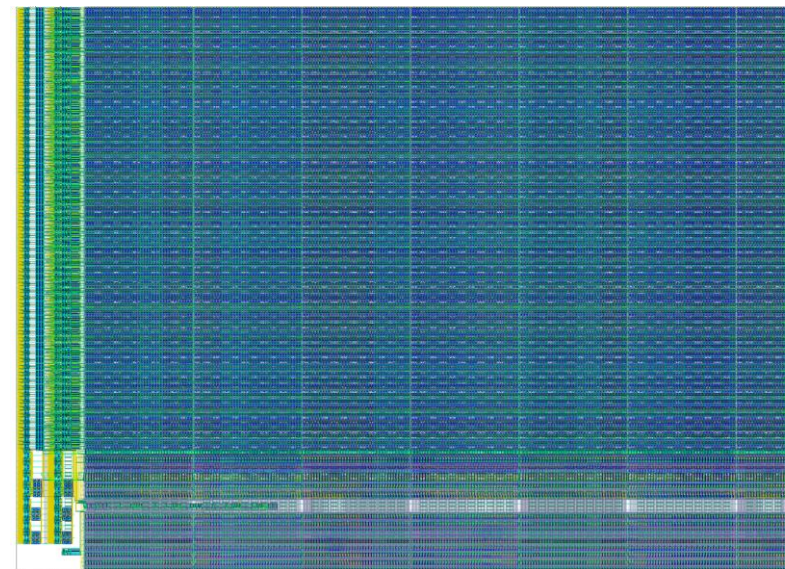| EVE-1 | EVE-2 | EVE-4 | EVE-8 | EVE-16 | EVE-32 |
|---|---|---|---|---|---|
| 1.0x | 1.0x | 1.0x | 1.0x | 1.15x | 1.51x |

Cycle-Time Overhead Over O3

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 12 of 15

**O3+IV**

- 4 elements
- OoO execution
- 3 exec. pipes

**O3+DV**

- 64 elements
- In-order execution
- 4 exec. pipes

**O3+EVE**

- 2048-256 elements
- In-order execution
- 1 exec. pipe

Cornell University Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 13 of 15

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

**IV unit offers modest speedups at a very low area-overhead cost (~10%).**

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

**DV engine achieves much higher speedups at the cost of larger area-overhead (~100%).**

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

**We reach balanced utilization with EVE-4 (i.e., parallelization factor of 4).**

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

**Due to sub-optimized cache-subsystem for vector memory operations, EVE-8 achieves better performance as the lower latency of EVE-8 is preferred to the longer vector lengths of EVE-4.**
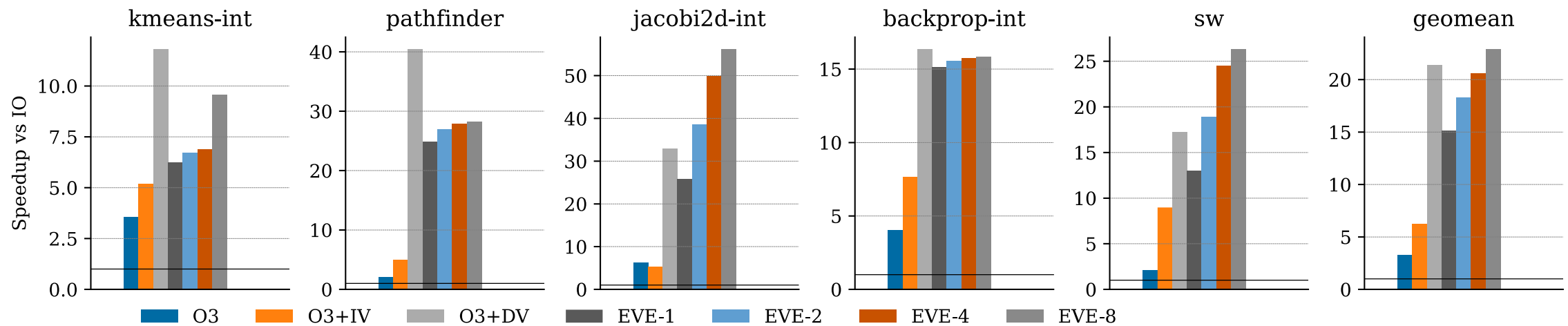
Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

**EVE-16 takes a cycle-time penalty of 15%; thus, it performs slower than EVE-8 despite the decreased latency**

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

Cornell University
Computer Systems Laboratory

**EVE-32 takes a cycle-time penalty of over 50%; despite having no transpose overhead, it performs very poorly.**

Cornell University
Computer Systems Laboratory

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion
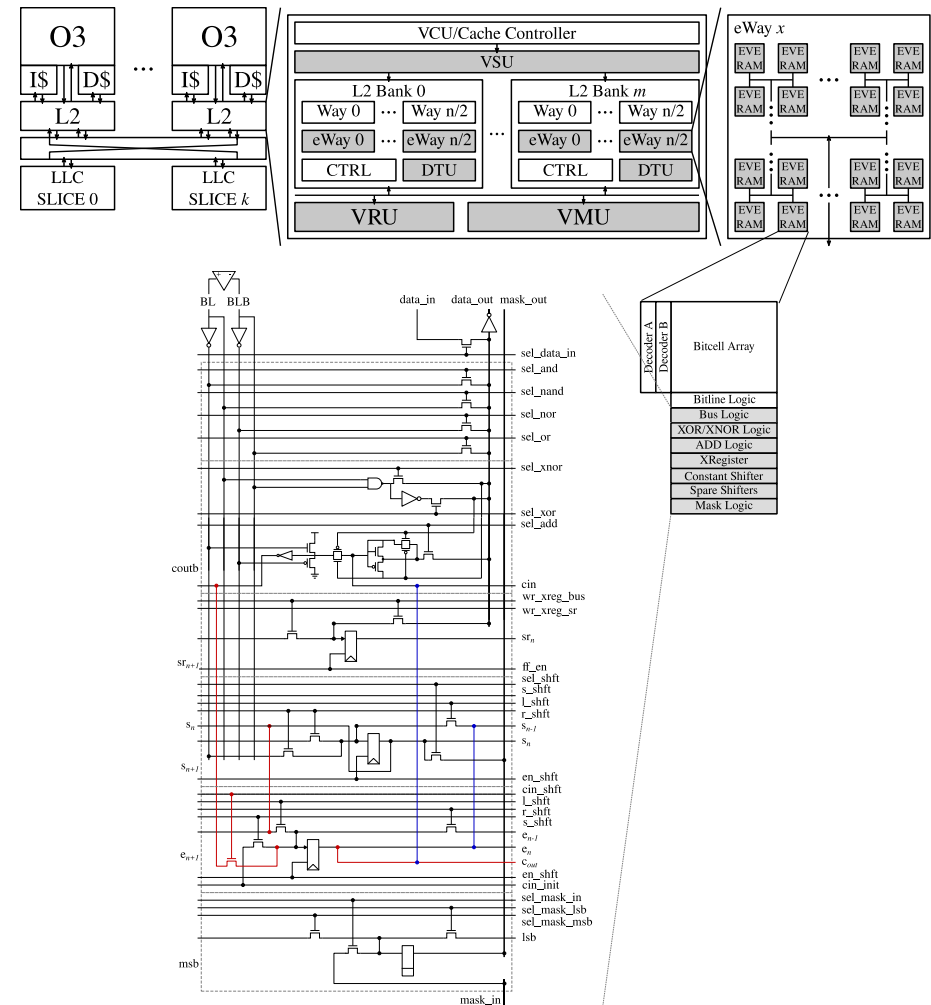
Page 14 of 15

**EVE-8 achieves comparable performance to an aggressive DV (~26x), while incurring an area-overhead equivalent to an IV (~12%).**

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 14 of 15

**Motivation**

■ **EVE Micro-Architecture**

■ **EVE Bit-Hybrid Execution Paradigm**

■ **EVE Circuits and Micro-Programming**

■ **EVE Evaluation**

**Conclusion**

- Architectural template for a novel SRAM-based compute-in-memory next-gen vector engine that supports the full RISC-V RVV specifications

- Bit-hybrid execution to balance throughout and latency by alleviating row and column under-utilization

- Detailed evaluation of EVE show-casing the impacts and benefits of bit-hybrid execution on an SRAM-based compute-in-memory micro-architecture

Motivation • Micro-Architecture • Bit-Hybrid Execution Paradigm • Micro-Programming & Circuits • Evaluation • Conclusion

Page 15 of 15