# Collaborative Accelerators for In-Memory MapReduce on Scale-up Machines

**Abraham Addisie**, Valeria Bertacco
University of Michigan

ASP-DAC 2019
January 24, 2019 Tokyo, Japan

UNIVERSITY OF
MICHIGAN

ada
Applications Driving Architectures

1

# Data-intensive applications challenges



**Challenge #1 :**



- Complex parallel programming models (e.g. Pthreads, OpenCL, OpenMP)

Solution: ⬇

- Simplified parallel programming models (e.g. MapReduce)

# Data-intensive applications challenges
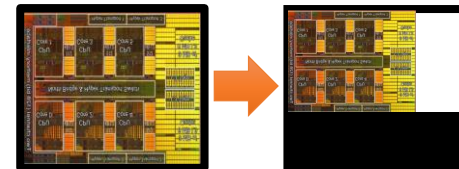


**Challenge #1 :**



- Complex parallel programming models (e.g. Pthreads, OpenCL, OpenMP)

Solution:

- Simplified parallel programming models (e.g. MapReduce)

**Challenge #2 :**



New features

Dark silicon

- Inefficiency in new hardware systems due to poor device scaling

Solution:

Accelerator-rich architectures

**Input**

car bus
car car
bus

car bus
bus car
bus

car car
car car

# MapReduce – WordCount example

**Input**       **Split**

car bus
car car
bus

→

car   bus
car car
bus

car bus
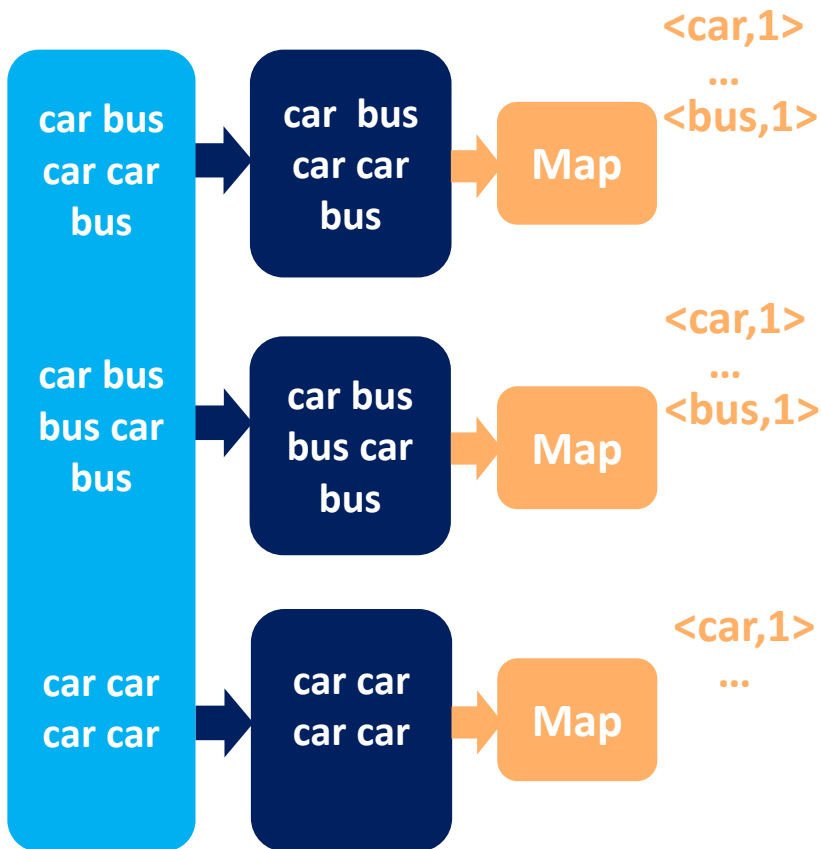bus car
bus

→

car bus
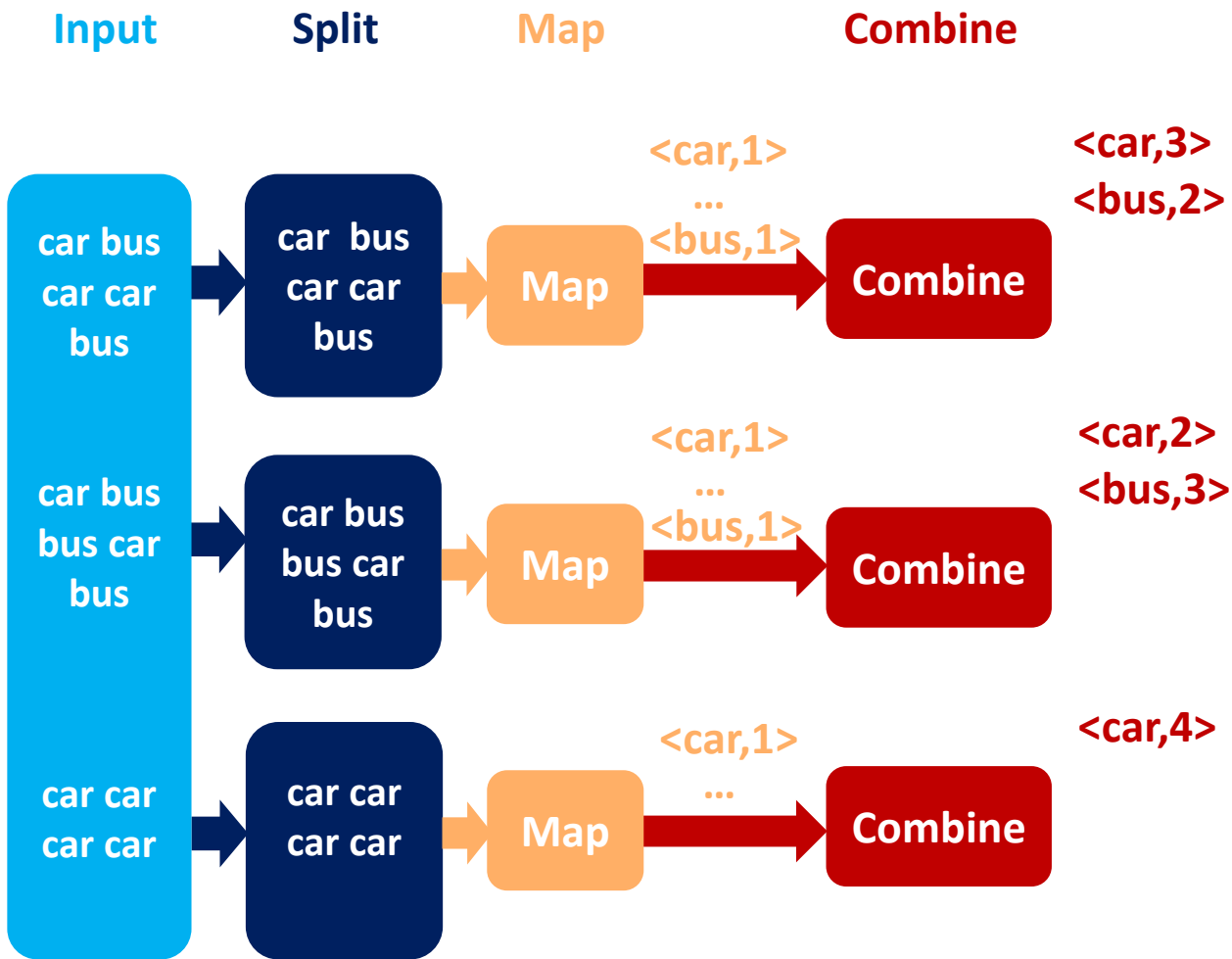bus car
bus
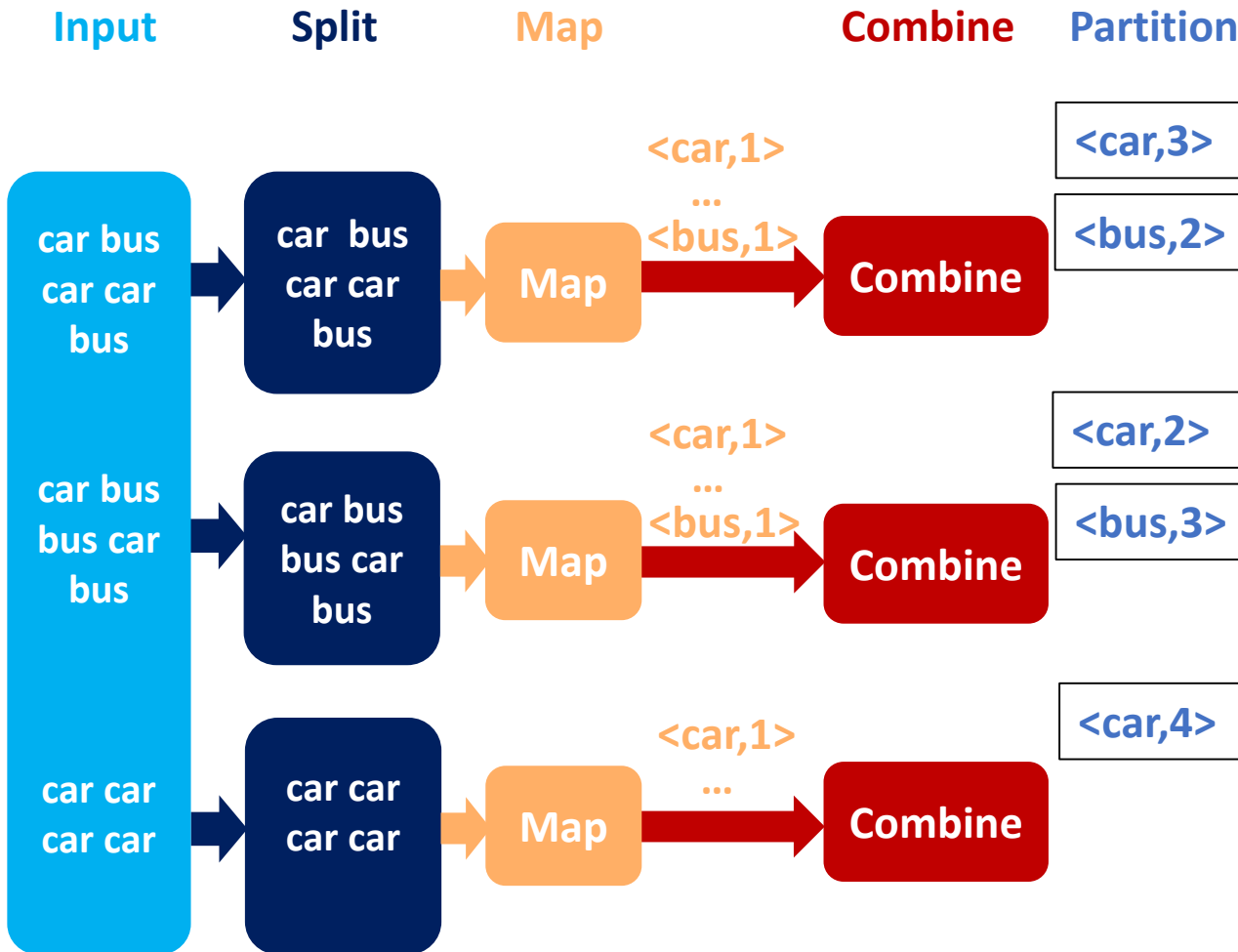
car car
car car

→

car car
car car

**Input**   **Split**   **Map**

| Input | Split | Map | |
|---|---|---|---|
| car bus car car bus | car   bus car car bus | Map | <car,1> ... <bus,1> |
| car bus bus car bus | car bus bus car bus | Map | <car,1> ... <bus,1> |
| car car car car | car car car car | Map | <car,1> ... |

# MapReduce – WordCount example

**Input**  **Split**  **Map**  **Combine**

car bus car car bus

car  bus car car bus → Map → <car,1> … <bus,1> → Combine → **<car,3>** **<bus,2>**

car bus bus car bus

car bus bus car bus → Map → <car,1> … <bus,1> → Combine → **<car,2>** **<bus,3>**

car car car car

car car car car → Map → <car,1> … → Combine → **<car,4>**

# MapReduce – WordCount example

| Input | Split | Map | Combine | Partition |
|-------|-------|-----|---------|-----------|



**Input:** car bus car car bus

**Split:** car bus car car bus

**Map:** Map

`<car,1>` ... `<bus,1>`

**Combine:** Combine

**Partition:** `<car,3>` `<bus,2>`

**Input:** car bus bus car bus

**Split:** car bus bus car bus

**Map:** Map

`<car,1>` ... `<bus,1>`

**Combine:** Combine

**Partition:** `<car,2>` `<bus,3>`

**Input:** car car car car

**Split:** car car car car

**Map:** Map

`<car,1>` ...

**Combine:** Combine

**Partition:** `<car,4>`

# MapReduce – WordCount example



Input | Split | Map | Combine | Partition | Shuffle/Reduce

car bus car car bus

car bus bus car bus

car car car car
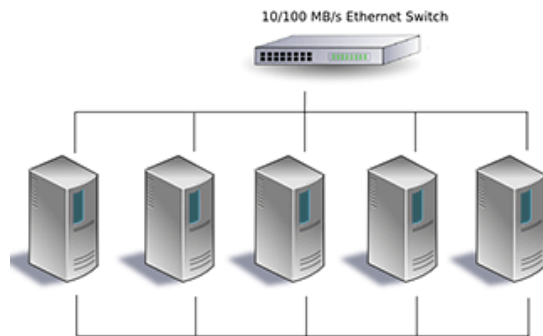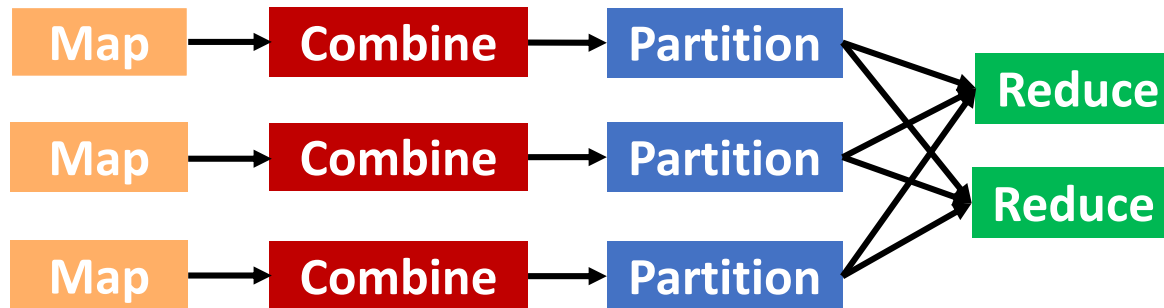
<car,1> ... <bus,1>

<car,3>
<bus,2>

<car, 3>
<car, 2>
<car, 4>

<car,2>
<bus,3>

<bus, 2>

<car,4>  <bus, 3>

# MapReduce – WordCount example

**Input**  **Split**  **Map**  **Combine**  **Partition**  **Shuffle/Reduce**  **Output**

car bus car car bus

car bus car car bus → **Map** → <car,1> ... <bus,1> → **Combine**

<car,3>

<bus,2>

<car, 3>

**Reduce** → car, 9

car bus bus car bus

car bus bus car bus → **Map** → <car,1> ... <bus,1> → **Combine**

<car,2>

<bus,3>

<car, 2>

<car, 4>

car car car car

car car car car → **Map** → <car,1> ... → **Combine**

<car,4>  <bus, 3>

<bus, 2>

**Reduce** → bus, 5
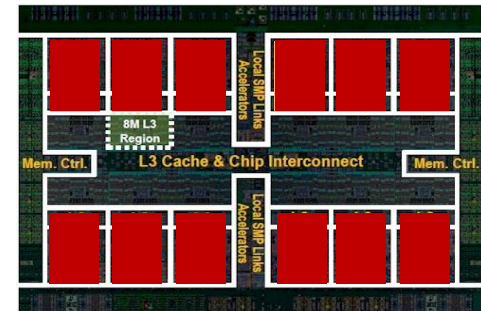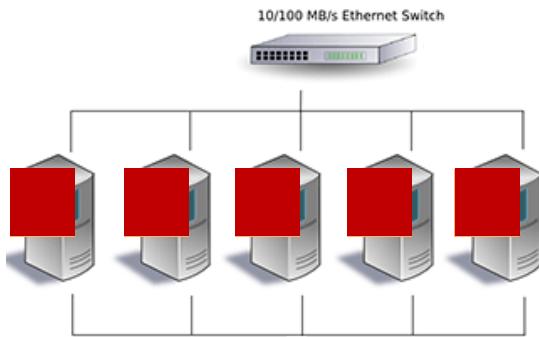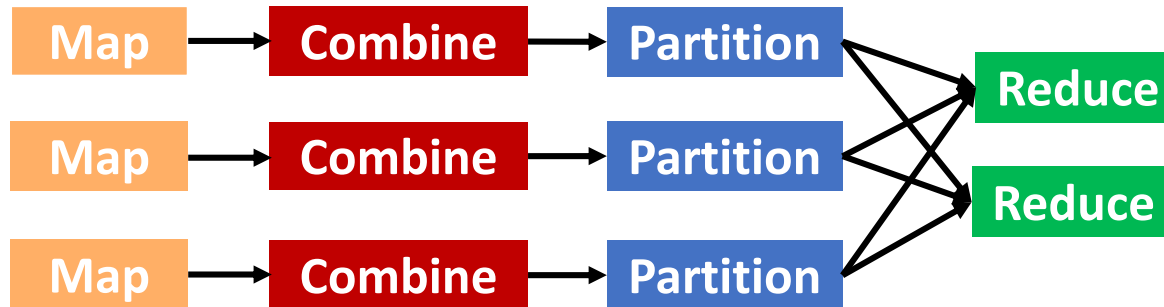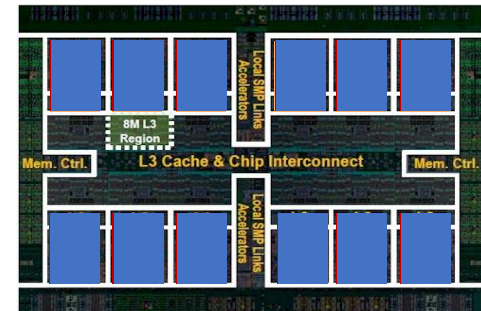
# MapReduce implementations



Scale-out:
Hadoop @ Yahoo

GPU:
Mars @ HKUST

CMP/Scale-up:
Phoenix++ @ Stanford

# MapReduce implementations



Scale-out:
Hadoop @ Yahoo

GPU:
Mars @ HKUST

CMP/Scale-up:
Phoenix++ @ Stanford

# MapReduce implementations

Map → Combine → Partition

Map → Combine → Partition → Reduce
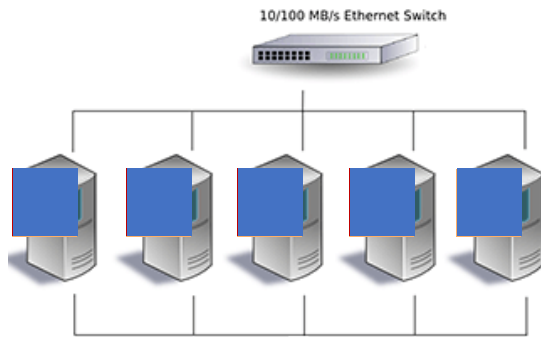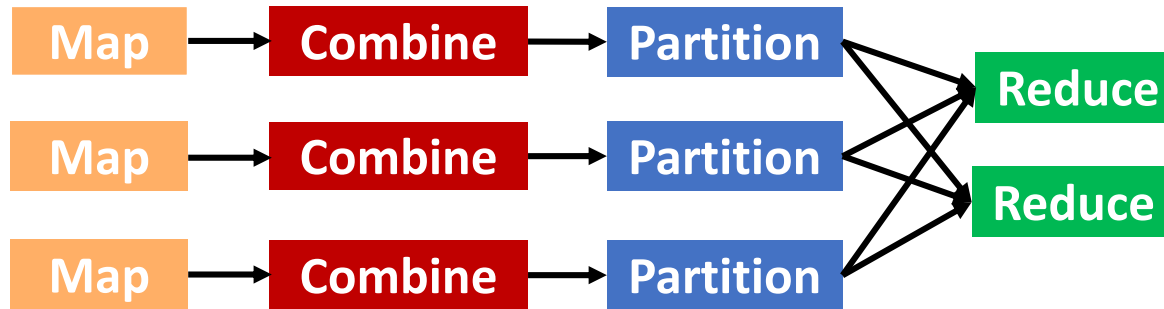
Map → Combine → Partition → Reduce

Scale-out:
Hadoop @ Yahoo

GPU:
Mars @ HKUST

CMP/Scale-up:
Phoenix++ @ Stanford

# MapReduce implementations

Map → Combine → Partition

Map → Combine → Partition → Reduce

Map → Combine → Partition → Reduce

**Scale-out:**
Hadoop @ Yahoo

**GPU:**
Mars @ HKUST

**CMP/Scale-up:**
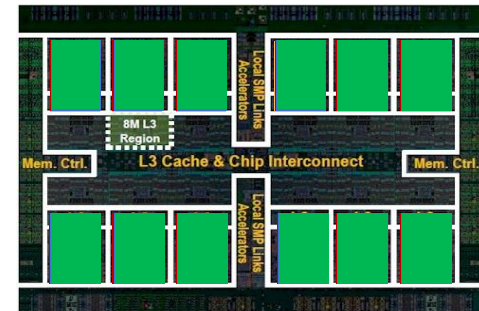Phoenix++ @ Stanford
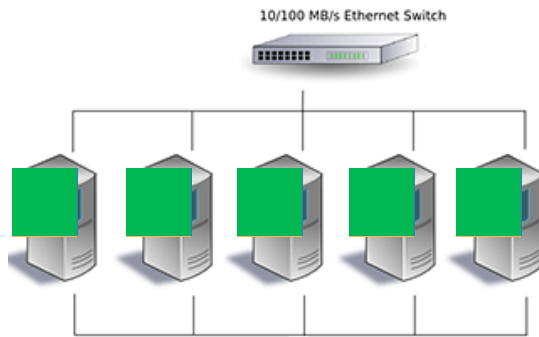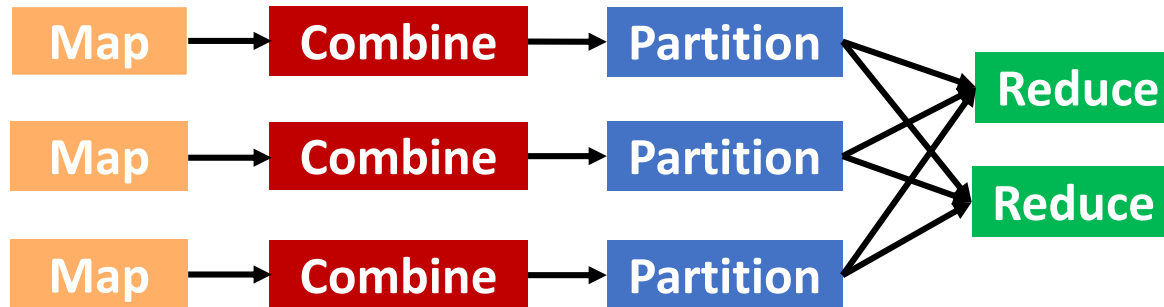
# MapReduce implementations
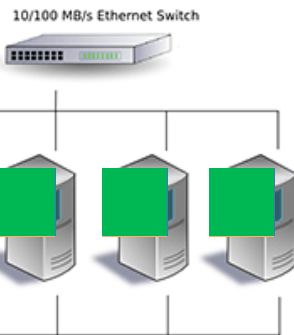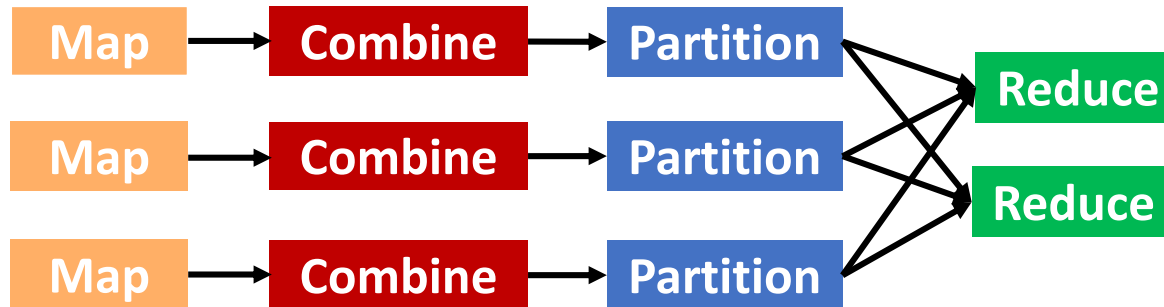


Scale-out:
Hadoop @ Yahoo

GPU:
Mars @ HKUST

CMP/Scale-up:
Phoenix++ @ Stanford

# MapReduce implementations



Scale-out:
Hadoop @ Yahoo

GPU:
Mars @ HKUST

focus of this work

CMP/Scale-up:
Phoenix++ @ Stanford

- Execution breakdown

## 16-core, Intel Xeon E5-2630 Server



Legend:
- map
- combine
- partition
- shuffle/reduce

Y-axis: execution time (ms), values: 0, 500, 1,000, 1,500, 2,000

X-axis: wc, mm, avg, pvc, h-img, lr, h-rt, h-mv, sc

Workload description
**wc**: word count

**mm**: min-max

**avg**: average

**pvc**: page view count

**h-img**: histogram image

**h-rt**: histogram user

**h-mv**: histogram movie

**sc**: sequence count

- Execution breakdown



16-core, Intel Xeon E5-2630 Server

map & combine dominate

Workload description

**wc**: word count

**mm**: min-max

**avg**: average

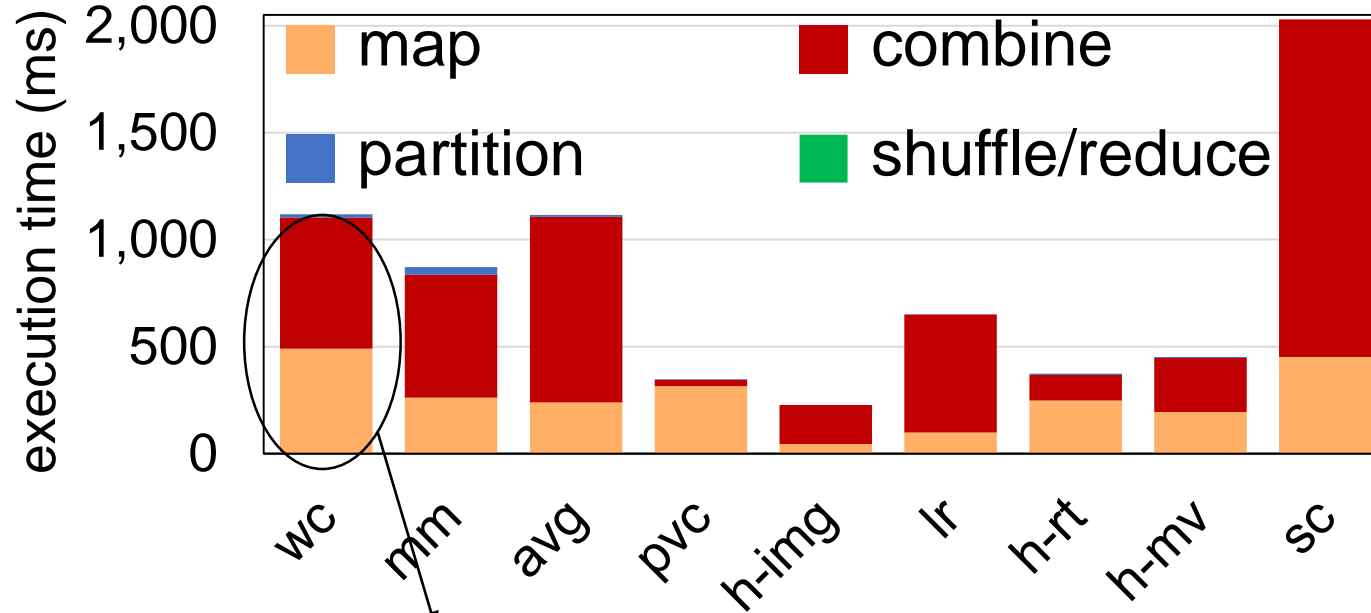**pvc**: page view count

**h-img**: histogram image

**h-rt**: histogram user

**h-mv**: histogram movie

**sc**: sequence count

# Phoenix++ studies

- Execution breakdown

## 16-core, Intel Xeon E5-2630 Server



Legend:
- map
- combine
- partition
- shuffle/reduce

y-axis: execution time (ms) — 0, 500, 1,000, 1,500, 2,000

x-axis: wc, mm, avg, pvc, h-img, lr, h-rt, h-mv, sc

map & combine dominate

Workload description

**wc**: word count

**mm**: min-max

**avg**: average

**pvc**: page view count

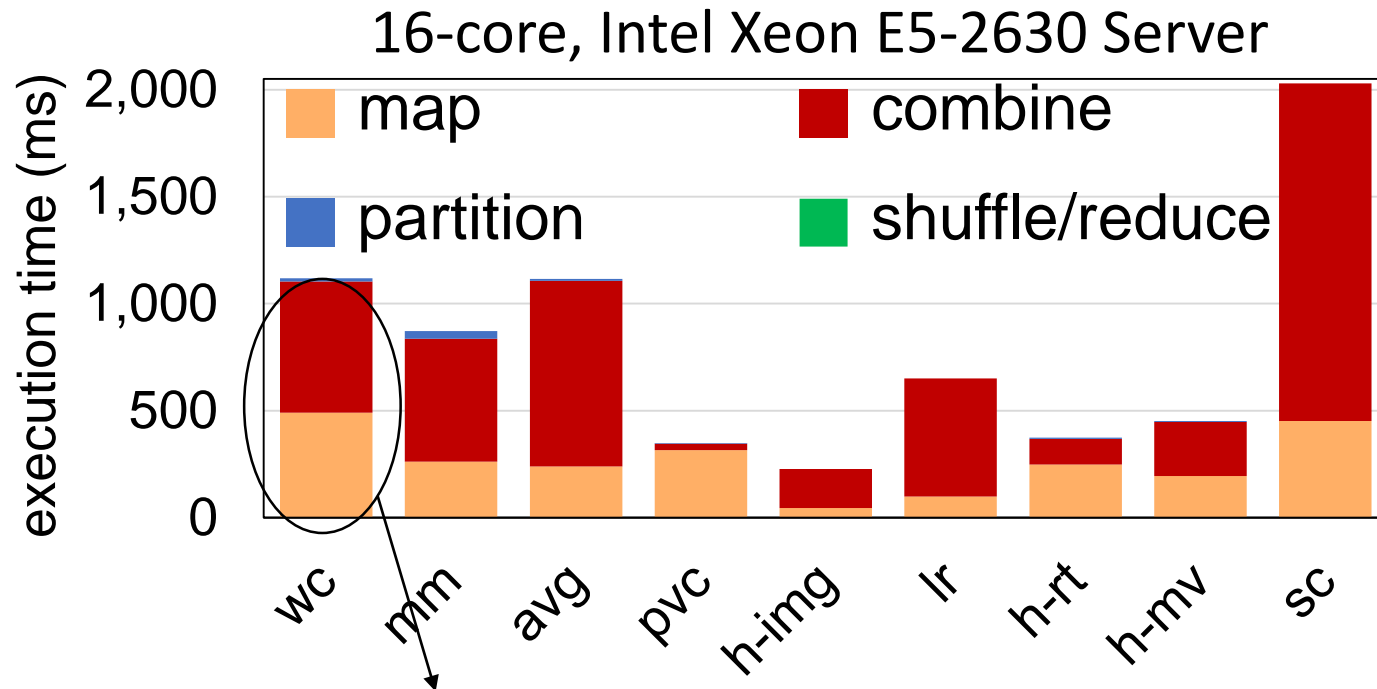**h-img**: histogram image

**h-rt**: histogram user

**h-mv**: histogram movie

**sc**: sequence count

- Inefficiencies:
  - Serial execution of map and combine phases
  - Inefficient key-value lookup during combine phase

# CASM overview

- Execution flow

# CASM overview

- Execution flow



Node #1      Node #2

Phoenix++ in CMP    Phoenix++ in CMP + CASM

Time

# CASM overview

- Execution flow



Node #1                                                          Node #2

Phoenix++ in CMP    Phoenix++ in CMP + CASM

Time

# CASM overview

- Execution flow



Node #1

Node #2

Phoenix++ in CMP          Phoenix++ in CMP + CASM

**Time**

Executed by cores in parallel

Executed by the collaborative accelerators

# Key contributions

- Scalable and collaborative accelerators

- Parallel execution of map and combine phases

- Faster execution of combine phase by the accelerators

- In-hardware hash function

- Each accelerator has two main storage structures

• Each accelerator has two main storage structures

**home scratchpad**

| key | value |
|-----|-------|
| plane | 1 |
| truck | 1 |

**local scratchpad**

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

high frequency keys aggregated locally

low frequency keys aggregated at their home

# Storage design space exploration

- **Option one**: local-only
  key-value pairs replicated across multiple accelerators

- **Option two**: home-only
  a key-value pair exists only in one location

accel0

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

accel1

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

accel2

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

accel0

| key | value |
|-----|-------|
| car | 3 |
| truck | 1 |
| plane | 1 |

accel1

| key | value |
|-----|-------|
| train | 2 |
| ship | 1 |
| bicycle | 2 |

accel2

| key | value |
|-----|-------|
| bus | 2 |
| motor | 1 |
| rocket | 1 |

- **Option one**: local-only
  key-value pairs replicated
  across multiple accelerators

- **Option two**: home-only
  a key-value pair exists only in
  one location

accel0

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

accel1

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

accel2

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

unified large memory
with no replication

| key | value |
|---------|-------|
| car | 3 |
| truck | 1 |
| plane | 1 |
| train | 2 |
| ship | 1 |
| bicycle | 2 |
| bus | 2 |
| motor | 1 |
| rocket | 1 |

# Storage design space exploration

- **Option one**: local-only
  key-value pairs replicated across multiple accelerators

- **Option two**: home-only
  a key-value pair exists only in one location

accel0

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

accel1

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

accel2

| key | value |
|-----|-------|
| car | 3 |
| bus | 2 |
| train | 2 |

accel0

| key | value |
|-----|-------|
| car | 3 |
| truck | 1 |
| plane | 1 |

accel1

| key | value |
|-----|-------|
| train | 2 |
| ship | 1 |
| bicycle | 2 |

accel2

| key | value |
|-----|-------|
| bus | 2 |
| motor | 1 |
| rocket | 1 |

| storage options | network traffic | memory traffic |
|-----------------|-----------------|----------------|
| local-only | low | high |
| home-only | high | low |
| local + home | low | low |

✓

# Key-value pair eviction policy

- "frequency" and "collision" bits are stored in scratchpads

  Scratchpad structure:

  | key | value | **frequency** | **collision** |
  |---|---|---|---|

  ➤ Frequency and collision update units:

  ➤ Simple heuristic function to identify frequently occurring keys:

  is match found?
  — No → increment collision
  — Yes ↓ increment frequency

  is freq. >= coll.?
  — No → stored kv is evicted
  — Yes ↓ new kv is sent to destination

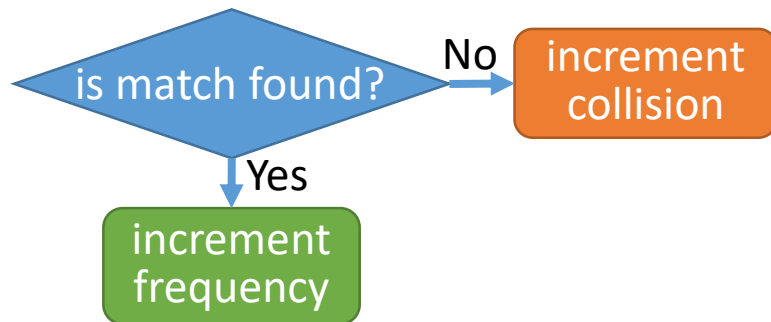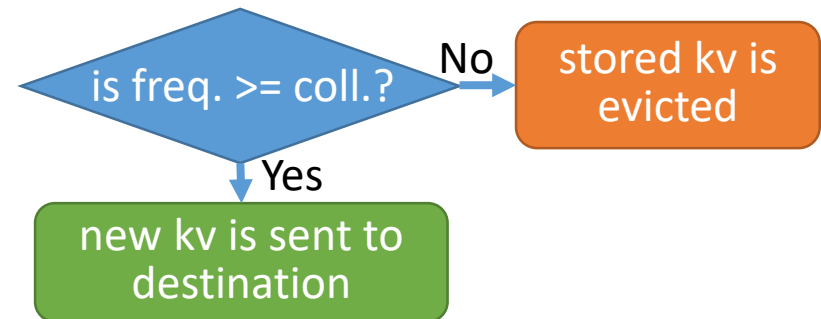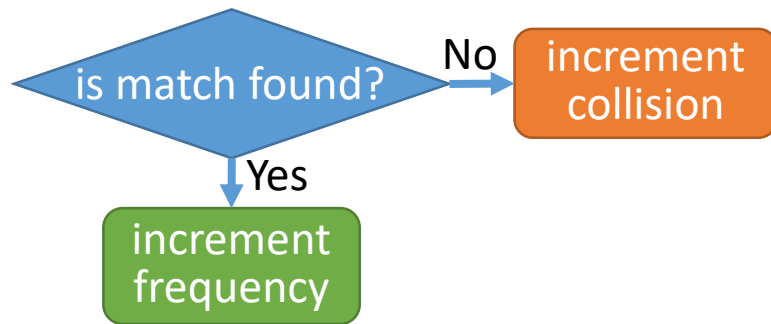# Key-value pair eviction policy

- "frequency" and "collision" bits are stored in scratchpads

Scratchpad structure:

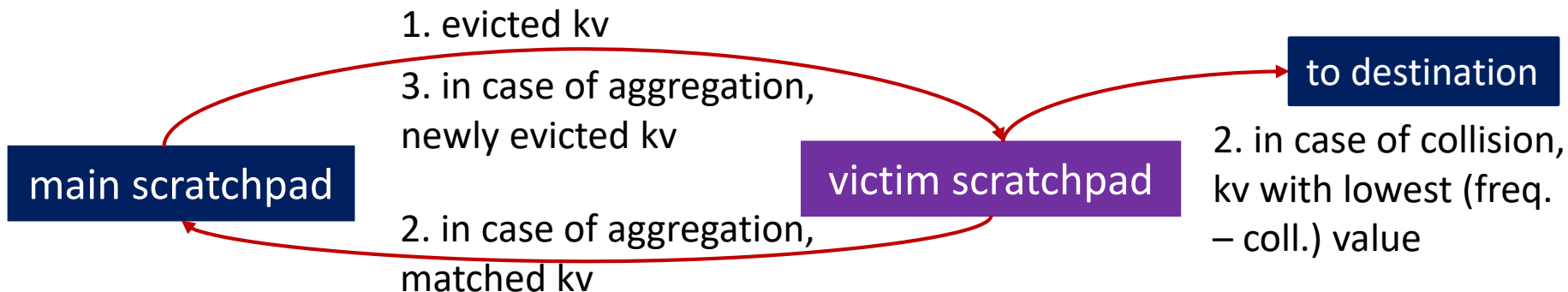| key | value | **frequency** | **collision** |
|-----|-------|---------------|---------------|

➢ Frequency and collision update units:

is match found? — No → increment collision

↓ Yes

increment frequency

➢ Simple heuristic function to identify frequently occurring keys:

is freq. >= coll.? — No → stored kv is evicted

↓ Yes

new kv is sent to destination

- Each scratchpad is augmented with victim scratchpad

1. evicted kv

3. in case of aggregation, newly evicted kv

to destination

main scratchpad

victim scratchpad

2. in case of collision, kv with lowest (freq. – coll.) value

2. in case of aggregation, matched kv

# Experimental framework

- **Scale-up CMP** configuration (Gem5/Garnet)

| Parameter | Value |
|---|---|
| Core | 64 cores, OoO, 8-wide |
| L1 D&I caches | 16KB |
| L2 cache (shared) | 128KB per core/slice |
| Coherence | MOESI directory-based |
| Memory | 4xDDR3-1600, 12GB/s |

- **CASM** configuration (Gem5/Garnet)

| Parameter | Value |
|---|---|
| Scratchpad size | 16KB |
| #entries per victim scratchpad | 8 |
| Max. key size | 64 bits |
| Max. value size | 64 bits |
| Freq. & coll. size | 8 bits |

# Experimental framework

- Scale-up CMP configuration (Gem5/Garnet)

| Parameter | Value |
|---|---|
| Core | 64 cores, OoO, 8-wide |
| L1 D&I caches | 16KB |
| L2 cache (shared) | 128KB per core/slice |
| Coherence | MOESI directory-based |
| Memory | 4xDDR3-1600, 12GB/s |

- CASM configuration (Gem5/Garnet)

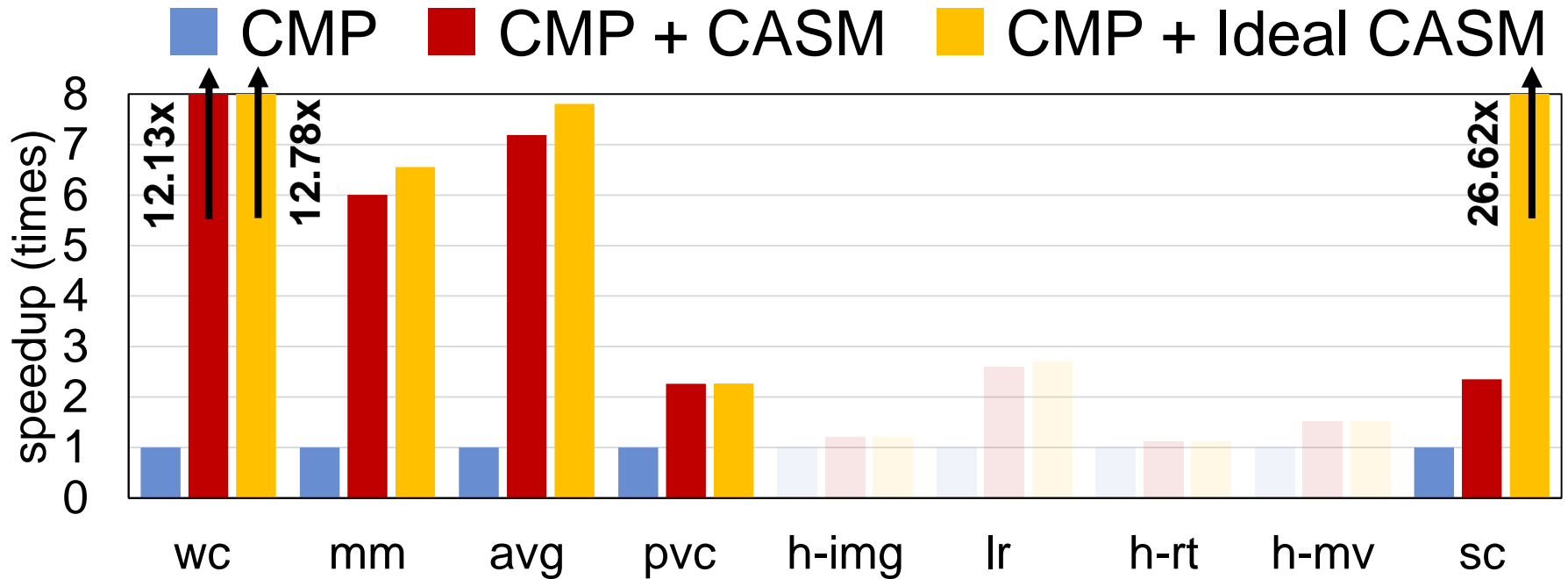| Parameter | Value |
|---|---|
| Scratchpad size | 16KB |
| #entries per victim scratchpad | 8 |
| Max. key size | 64 bits |
| Max. value size | 64 bits |
| Freq. & coll. size | 8 bits |

- Workload characteristics

| workload | wc | mm | avg | pvc | h-img | lr | h-rt | h-mv | sc |
|---|---|---|---|---|---|---|---|---|---|
| #unique keys | 257K | 28K | 28K | 10K | 768 | 5 | 5 | 20K | 3.5M |
| cache locality | low | low | low | low | high | high | high | high | low |

- > 4x speedup on average
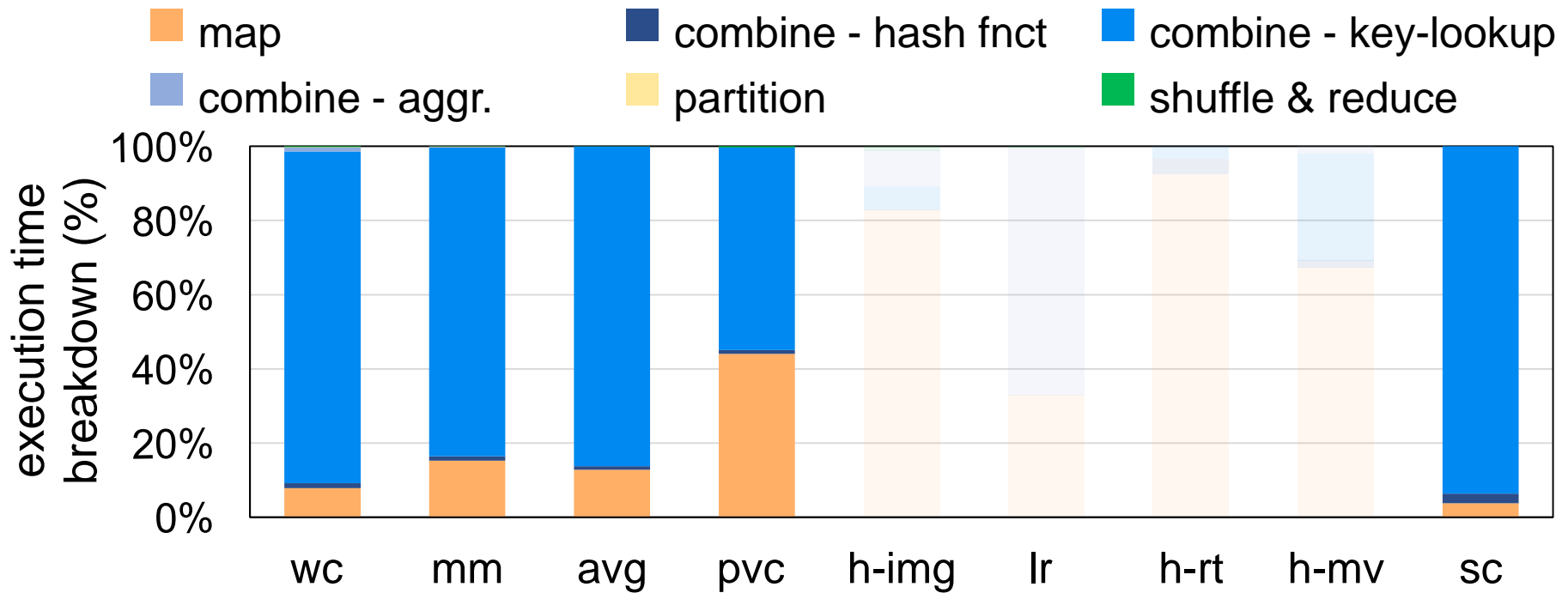
- > 3.5x energy saving on average

- > **4x** speedup on average

- > **3.5x** energy saving on average

**Large #unique keys** & **no cache locality** ➡ **high speedup**

# Sources of performance benefits

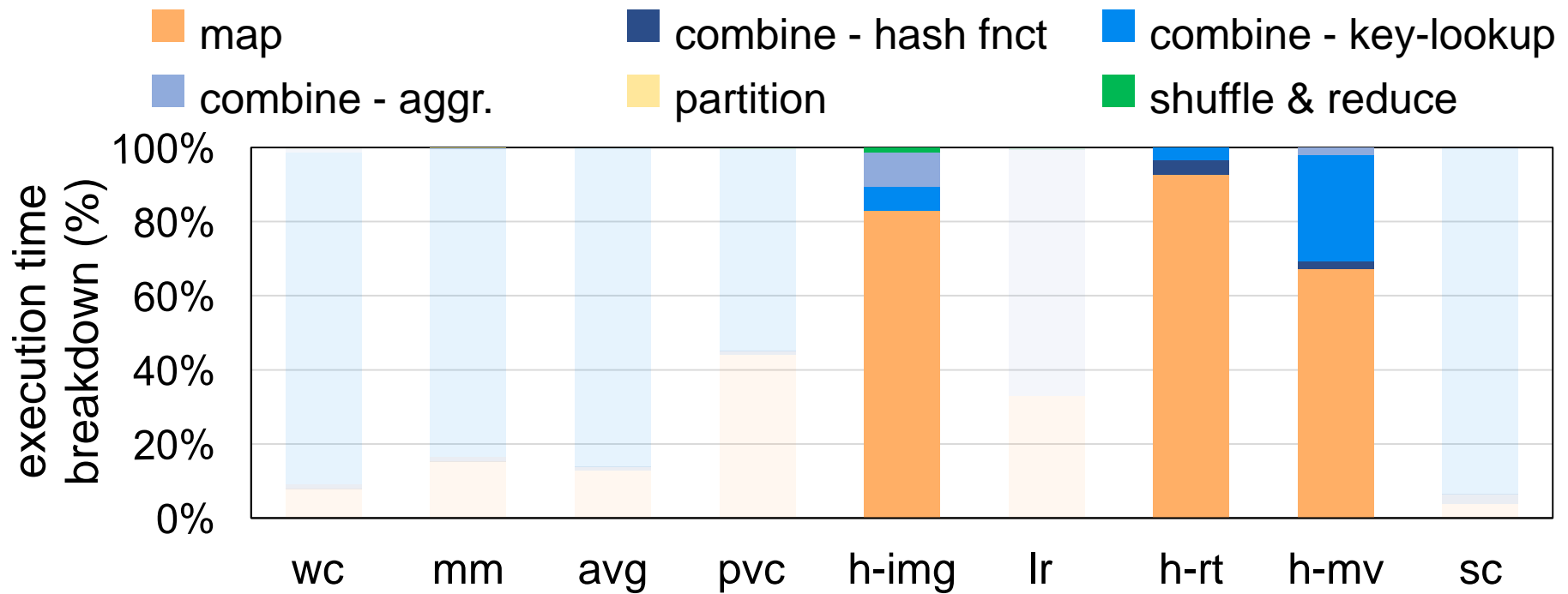- Mainly due to offloading the combine phase to CASM



**Dominant phase:**

combine – key-lookup

- Mainly due to offloading the combine phase to CASM

map  ■ combine - hash fnct  ■ combine - key-lookup
■ combine - aggr.  partition  ■ shuffle & reduce

execution time breakdown (%)

100%
80%
60%
40%
20%
0%

wc   mm   avg   pvc   h-img   lr   h-rt   h-mv   sc

**Dominant phase:**

map

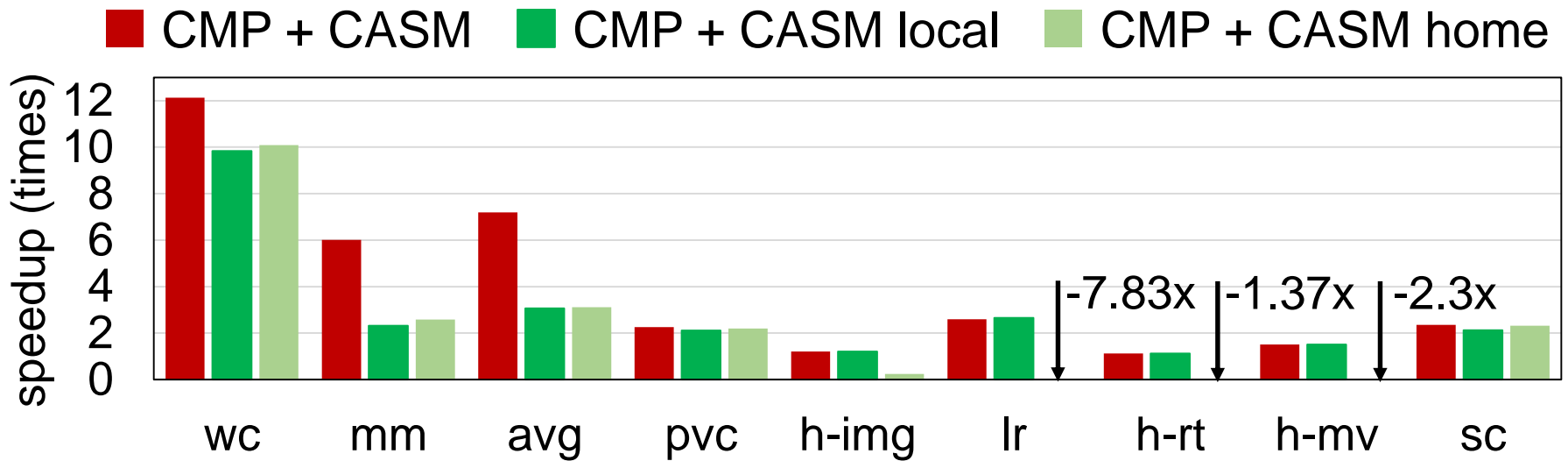- Mainly due to offloading the combine phase to CASM



**Dominant phase:**

combine – aggregation

# Speedup contribution: local vs home



Legend: CMP + CASM (dark red), CMP + CASM local (green), CMP + CASM home (light green)

Y-axis: speedup (times), 0 to 12

Categories: wc, mm, avg, pvc, h-img, lr, h-rt, h-mv, sc

Annotations: -7.83x, -1.37x, -2.3x

- A hybrid of local and home accelerators provides significant benefits across applications

# Conclusion

- MapReduce on scale-up machines suffers from:
  - serial execution of map and combine phases
  - inefficient key-value lookup

- Solution:
  - Parallel execution of map and combine phases
  - Local/home partitioned on-chip storage
  - Aggregation near on-chip storage

- CASM provides:
  - >4x in performance on average
  - >3.5x in energy saving on average
  - < 6% of area overhead