



SIMD Vectorization

Rethinking SIMD Vectorization for In-Memory Databases

Sung-Soo Kim

sungsoo@etri.re.kr

Data Platform Research Section

ETRI

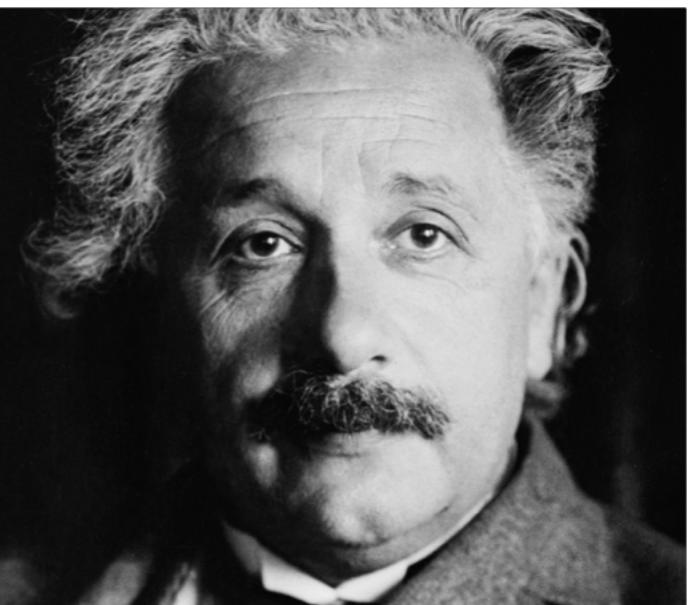
Outline

2

- **Background**
- **SIMD Vectorization Algorithms**
 - **Selective Load/Store, Selective Gather/Scatter**
 - **Selective Scans**
 - **Hash Tables - Probing**
 - **Partitioning - Histogram**
- **Summary**

*“Reading, after a certain age,
diverts the mind too much from its creative pursuits.
Any man who reads too much and uses his own brain
too little fails into lazy habits of thinking.”*

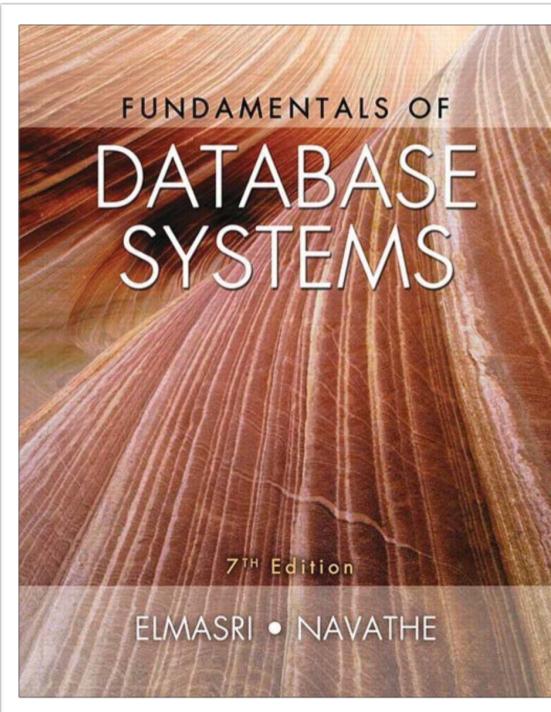
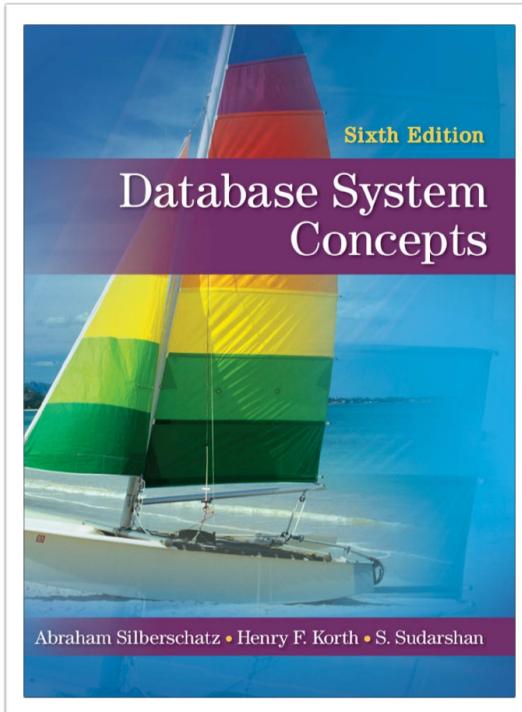
- ALBERT EINSTEIN



References & Slide Credits

3

- *Database System Concepts, 6th Edition*
- *Fundamentals of Database Systems, 7th Edition*
- *Rethinking SIMD Vectorization for In-Memory Databases, SIGMOD 2015.*



Rethinking SIMD Vectorization for In-Memory Databases

Orestis Polychroniou^{*} Columbia University Arun Raghavan Oracle Labs Kenneth A. Ross^{*} Columbia University
orestis@cs.columbia.edu arun.raghavan@oracle.com kar@cs.columbia.edu

ABSTRACT
Analytical databases are continuously adapting to the underlying hardware in order to saturate all sources of parallelism. At the same time, hardware evolves in multiple directions to explore different trade-offs. The MIC architecture, one such example, strays from the mainstream CPU design paradigm by utilizing SIMD instructions to fill the performance gap. Databases have been attempting to utilize the SIMD capability of CPUs for years. However, CPUs have only recently added wider SIMD registers and more SIMD instructions, since they do not rely primarily on SIMD for efficiency. In this paper, we present novel designs and implementations for analytical databases that utilize advanced SIMD operations, such as gathers and scatters. We study selections, hash tables, and partitioning, and combine them with SIMD and joins. Our evaluation on the MIC-based Xeon Phi coprocessor shows that the latest stream CPUs show that our vectorization designs are up to an order of magnitude faster than state-of-the-art SIMD and vectorized designs. We also highlight the benefits of efficient vectorization on the algorithmic design of in-memory database operators, as well as the architectural design and power usage of hardware, by showing SIMD cores compute quickly fast to memory access. This work is applicable to CPUs and co-processors with advanced SIMD capabilities, using either many simple cores or fewer complex cores.

1. INTRODUCTION
Real time analytics are the steering wheels of big data driven business intelligence. Database customer needs have extended beyond OLTP to high ACID transaction throughput, to support OLAP queries over a massive data warehouse. As a consequence, vendors offer fast OLAP solutions, either by rebuilding a new DBMS for OLAP [37], or by improving within the existing OLTP-focused DBMS.

^{*}Work partly done when first author was at the Oracle Labs.
Supported by NSF grants IIS-1422888 and an Oracle gift.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are made available for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from <http://www.acm.org>.
SIGMOD '15, May 31–June 4, 2015, Melbourne, Victoria, Australia.
Copyright 2015 ACM 978-1-4503-2784-9/15/05 ...\$15.00.
<http://doi.acm.org/10.1145/2723572.2747945>.

The advent of large main-memory capacity is one of the reasons that kind-of-a-way analytical query execution has become possible. Query optimization used to measure blocks fetched from disk as the primary unit of query cost. Today, the entire database can often remain main-memory resident while being processed by the query optimizer and execution environment. The prevalent shift is database design for the new era are column stores [19, 28, 37]. They allow for higher data compression and easier parallelization to minimize the number of column accesses per query, and use column oriented execution coupled with late materialization [9] to eliminate unnecessary access to RAM resident columns. Furthermore, column stores have a wide range of parallelism: thread parallelism, instruction level parallelism, and data parallelism. Analytical databases have evolved to take advantage of all forms of parallelism. Thread parallelism is adopted for index creation operations, by splitting the input equally among threads [3, 4, 5, 8, 14, 31, 40], and in the case of queries that combine multiple operators, by splitting the input into multiple parallel streams to split the materialized data in chunks that are distributed to threads dynamically [18, 28]. Instruction level parallelism is achieved by applying the same operation to a block of tuples [9] and then splitting the block into smaller blocks [18, 22]. Data parallelism is achieved by implementing each operator to use SIMD instructions effectively [7, 15, 26, 30, 39, 41]. The latest Intel Xeon Phi coprocessor [1] is a means to deliver more performance within the same power budget available per chip. Mainstream CPUs have evolved on all sources of parallelism, featuring massively superscalar pipelines, multiple levels of instruction and data cache, and advanced SIMD capabilities, all replicated on multiple cores per CPU chip. For example, the latest Intel Haswell architecture [2] has 16 cores, each with a 64-bit register file with 192 register buffer entries for out-of-order execution, 256-bit SIMD instructions, two levels of private caches per core and a large shared cache, and scales up to 18 cores per chip.

Concurrently, a new approach to processor design has surfaced. The design, named the many-integrated-core (MIC) architecture, uses cores that are designed to be highly specialized for a certain task by removing the massively superscalar pipeline, out-of-order execution, and the large L3 cache. Each core is based on a Pentium 4 processor with a single in-order pipeline, but is augmented with larger SIMD registers, advanced SIMD instructions, and simultaneous multithreading to hide load and instruction latency. Since each core has a smaller area and power footprint, more cores can be packed in the chip.

Multi-Core CPUs

4

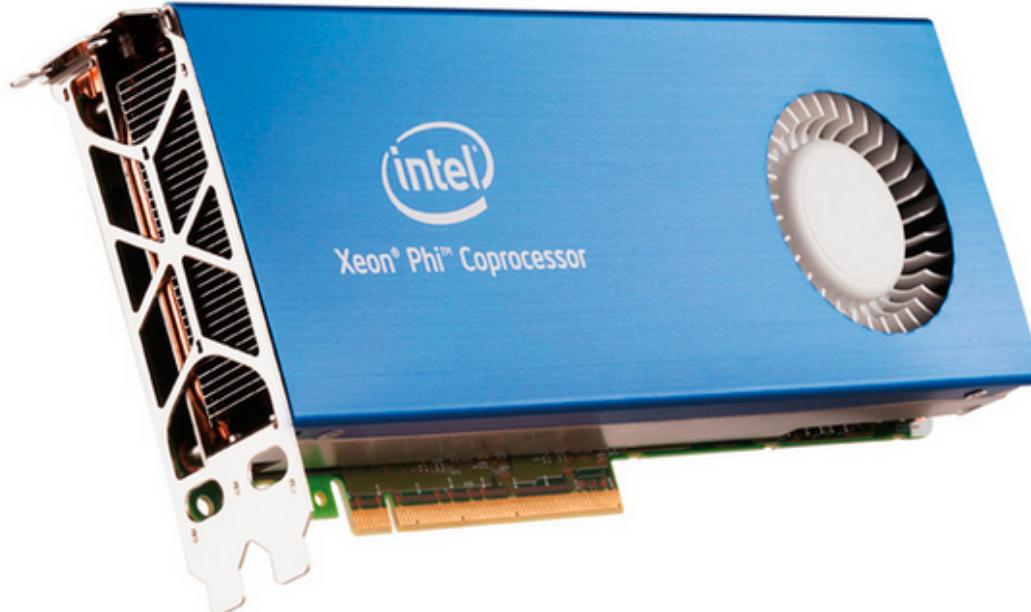
- Use a small number of *high-powered* cores.
 - Intel Haswell / Skylake
 - High power consumption and area per core.

- Massively superscalar and aggressive out-of-order execution
 - Instructions are issued from a *sequential* stream.
 - Check for *dependencies* between instructions.
 - Process *multiple instructions* per clock cycle.

Many Integrated Cores (MIC)

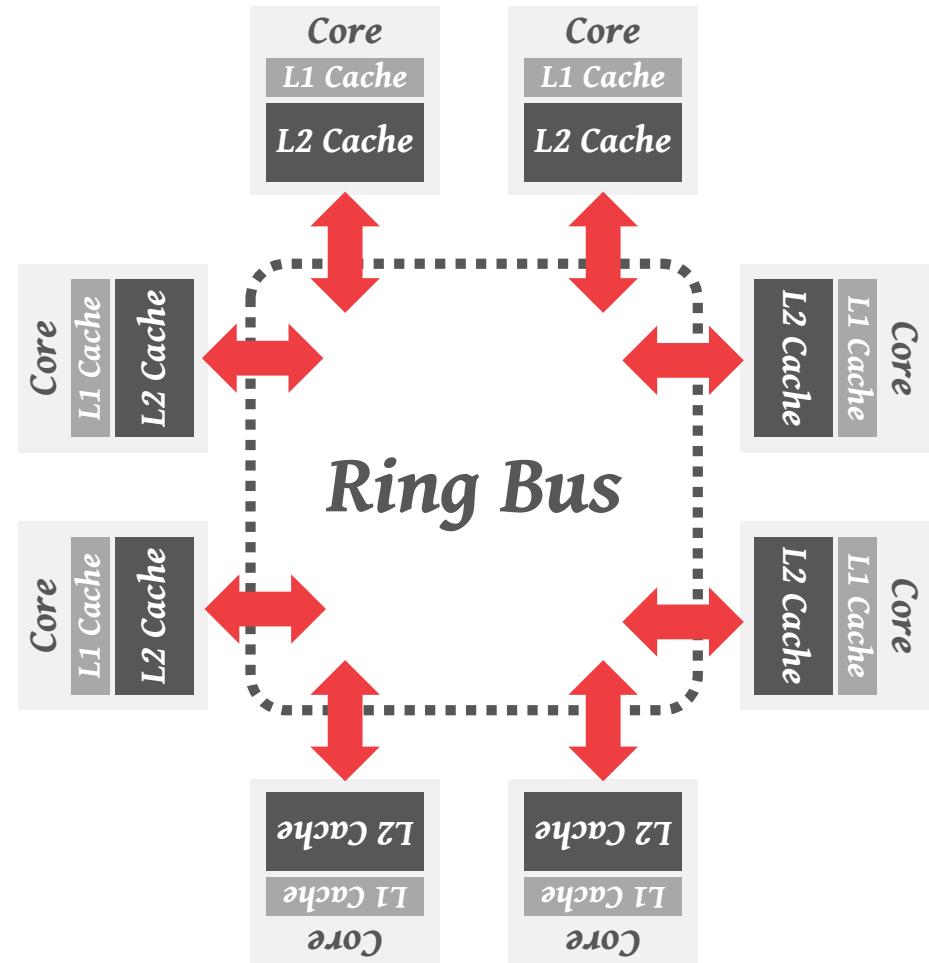
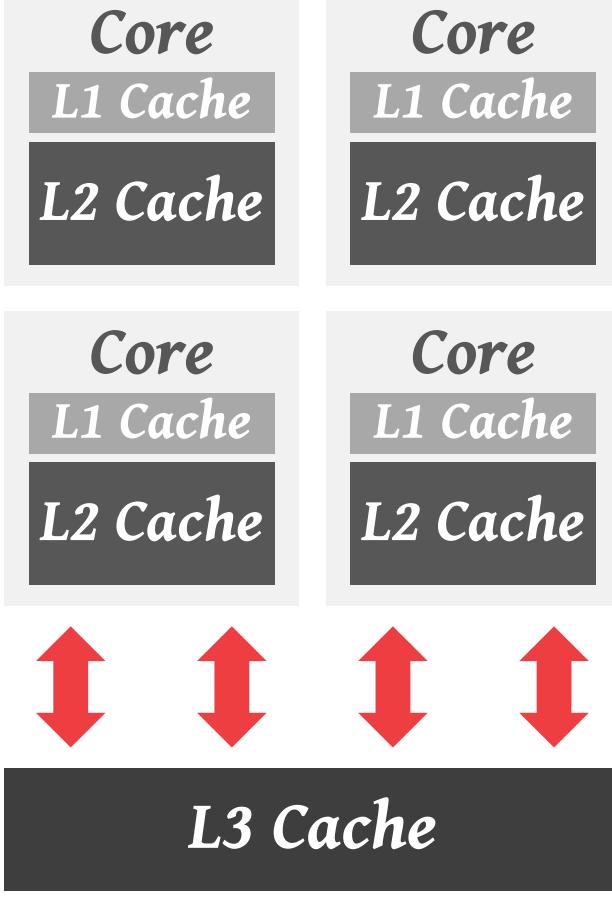
5

- Use a larger number of *low-powered* cores.
 - Intel Xeon Phi
 - Cores = Intel P54C (aka Pentium from the 1990s).
 - *Low* power consumption and area per core.
- Non-superscalar and in-order execution but with expanded SIMD capabilities.
 - More *advanced instructions* with larger register sizes.



Multi-Core vs. MIC

6



Vectorization

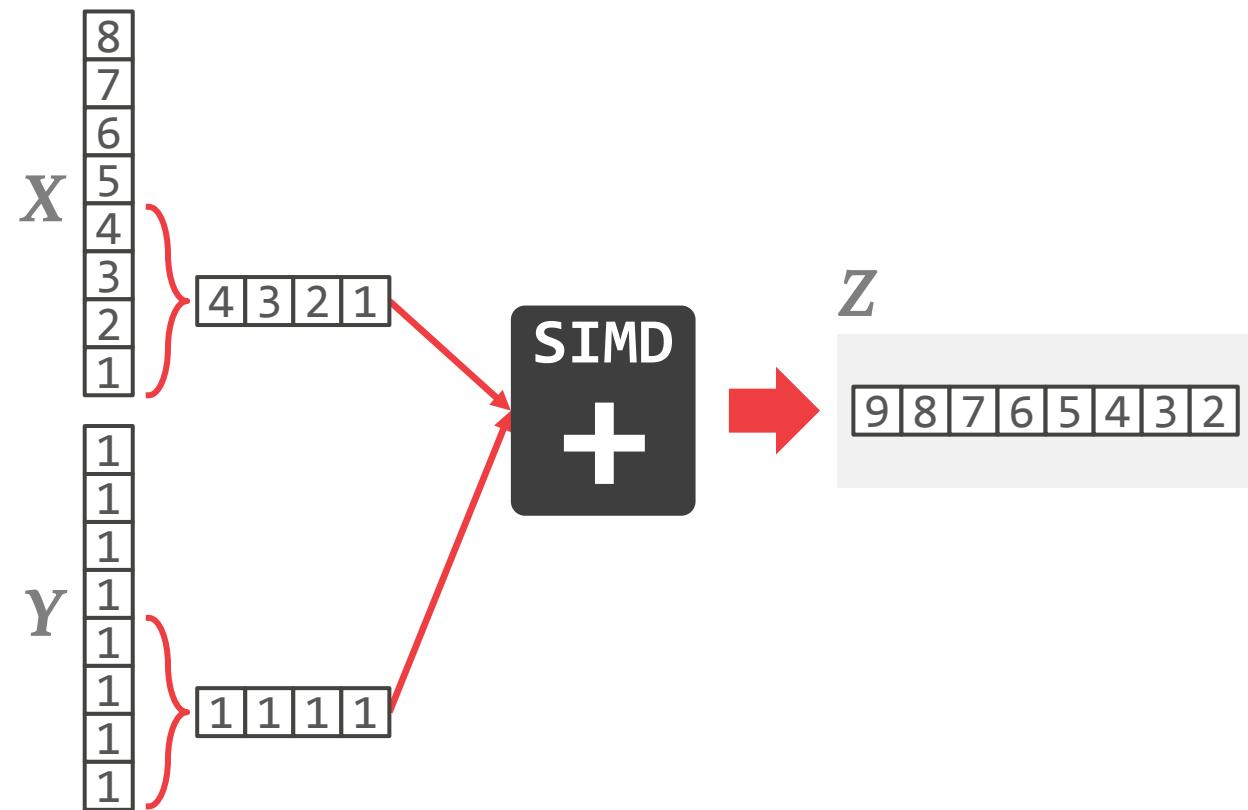
7

- A program is converted from a **scalar** implementation that processes a **single pair** of operands **at a time**, to a **vector** implementation that processes **one operation on multiple pairs of operands at once**.

$$X + Y = Z$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \dots \\ x_n + y_n \end{pmatrix}$$

```
for (i=0; i<n; i++) {  
    Z[i] = X[i] + Y[i];  
}
```



Automatic Vectorization

8

- The compiler can identify when instructions inside of a loop can be *written* as a *vectorized* operation.
 - e.g., Intel ICC compiler
- Works for **simple loops only** and is rare in database operators.
- Requires hardware support for SIMD instructions.

Manual Vectorization

9

■ Linear Access Operators

- Predicative evaluation
- Compression

■ Ad-hoc Vectorization

- Sorting
- Merging
 - e.g., Sort-Merge Join

■ Composable Operations

- Multi-way trees
- Bucketized hash tables

Single Instruction, Multiple Data (SIMD)

10

- A class of CPU instructions that allow the processor to perform **the same operation on multiple data points simultaneously**.
- All major ISAs have microarchitecture support SIMD operations.
 - **x86**: MMX, SSE, SSE2, SSE3, SSE4, AVX
 - **PowerPC**: Altivec
 - **ARM**: NEON

SIMD Example

11

$$X + Y = Z$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \dots \\ x_n + y_n \end{pmatrix}$$

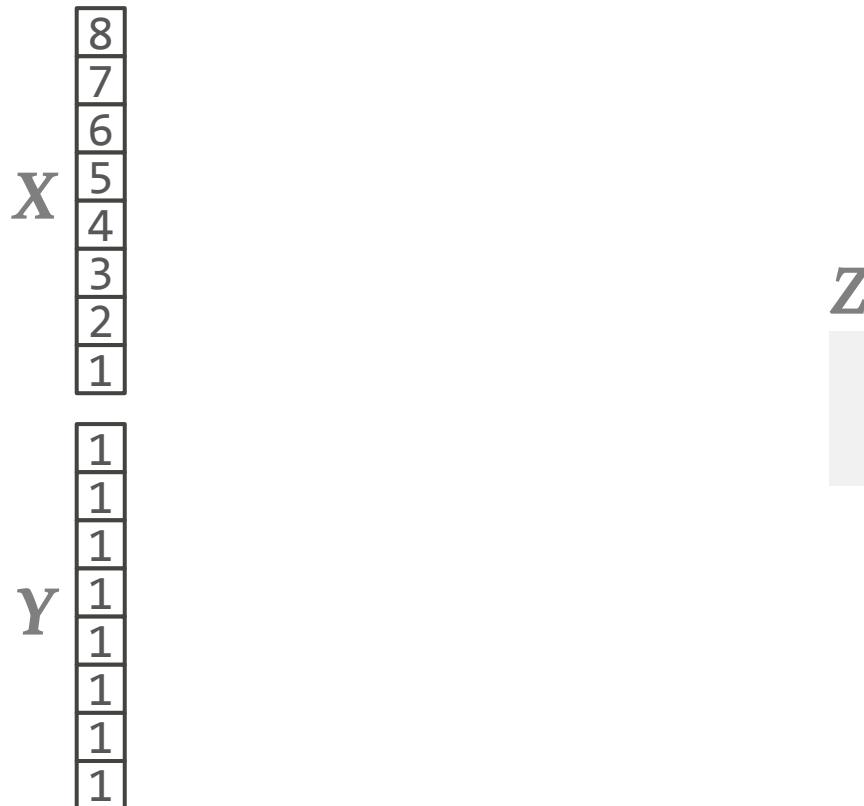
SIMD Example

12

$$X + Y = Z$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \dots \\ x_n + y_n \end{pmatrix}$$

```
for (i=0; i<n; i++) {  
    Z[i] = X[i] + Y[i];  
}
```



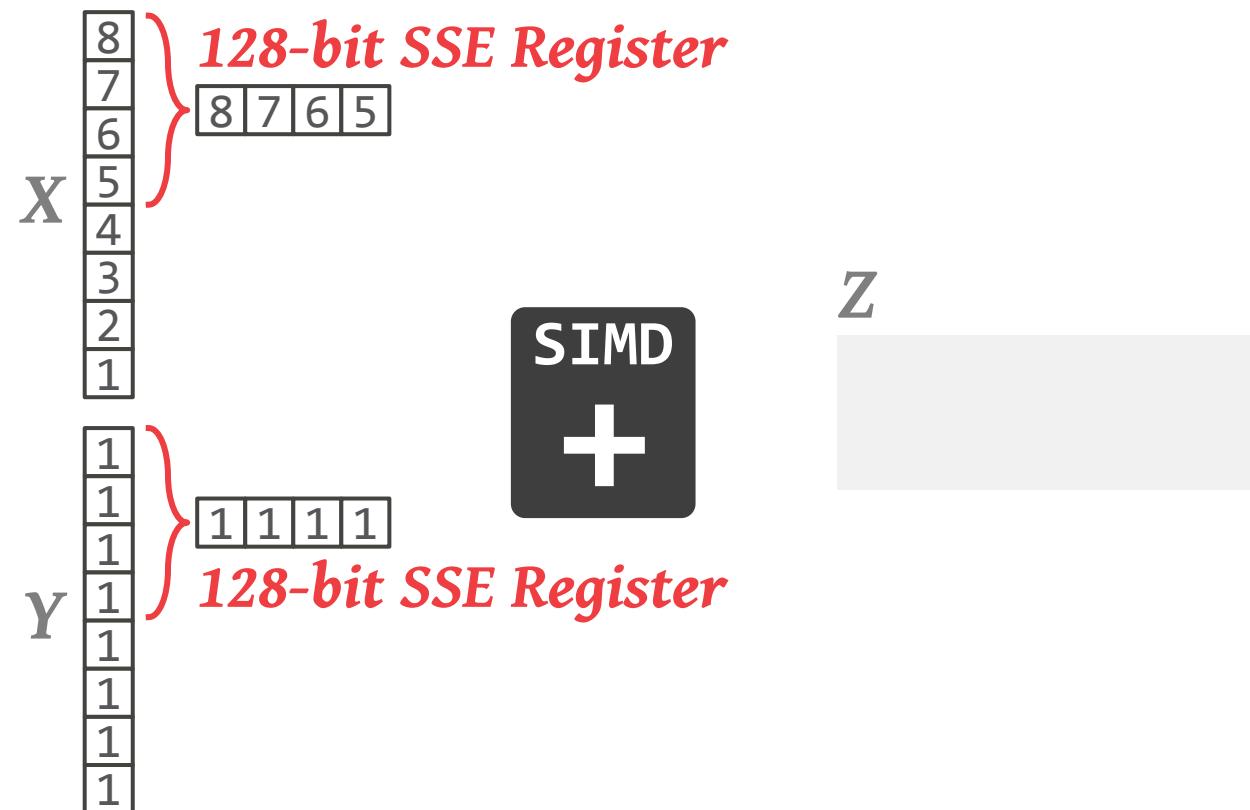
SIMD Example

13

$$X + Y = Z$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \dots \\ x_n + y_n \end{pmatrix}$$

```
for (i=0; i<n; i++) {  
    Z[i] = X[i] + Y[i];  
}
```



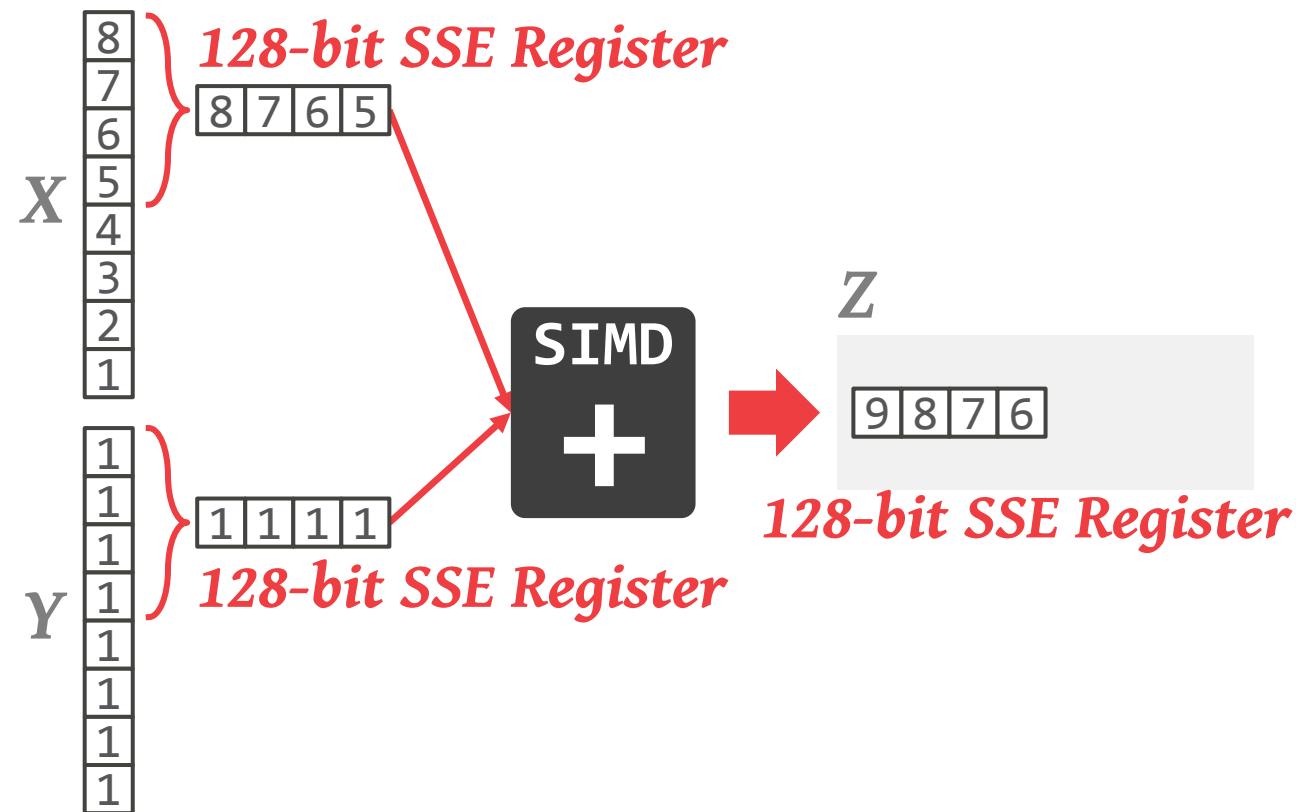
SIMD Example

14

$$X + Y = Z$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \dots \\ x_n + y_n \end{pmatrix}$$

```
for (i=0; i<n; i++) {  
    Z[i] = X[i] + Y[i];  
}
```



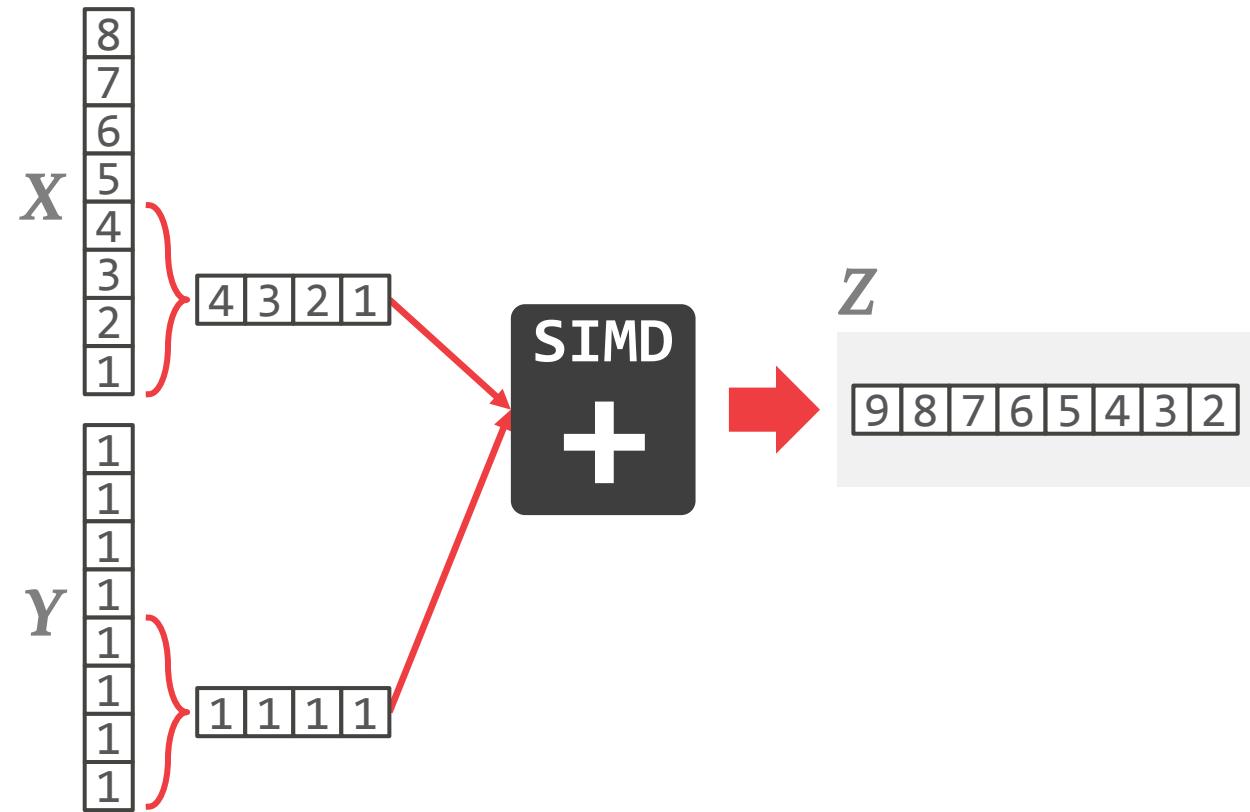
SIMD Example

15

$$X + Y = Z$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} = \begin{pmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \dots \\ x_n + y_n \end{pmatrix}$$

```
for (i=0; i<n; i++) {  
    Z[i] = X[i] + Y[i];  
}
```



Streaming SIMD Extensions (SSE)

16

- SSE is a collection SIMD instructions that target special 128-bit SIMD registers.
- These registers can be packed with four 32-bit scalars after which an operation can be performed on each of four elements simultaneously.
- First introduced by Intel in 1999.

SSE Instructions

17

■ Data Movement

- Moving data in and out of vector registers

■ Arithmetic Operations

- Apply operation on multiple data items (e.g., 2 doubles, 4 floats, 16 bytes)
- Example: **ADD, SUB, MUL, DIV, SQRT, MAX, MIN**

■ Logical Instructions

- Logical operations on multiple data items
- Example: **AND, OR, XOR, ANDN, ANDPS, ANDNPS**

■ Comparison Instructions

- Comparing multiple data items (**==, <, <=, >, >=, !=**)

■ Shuffle Instructions

- Move data in between SIMD registers

■ Miscellaneous

- Conversion: Transform data between x86 and SIMD
- Cache Control: Move data directly from SIMD registers to memory (bypassing CPU cache.)

Vectorized DBMS Algorithms

- Principles for **efficient vectorization** by using **fundamental vector operations** to construct more advanced functionality.

- Favor **vertical vectorization** by processing different input data **per lane**.
 - Maximize **lane utilization** by executing different things **per lane subset**.



Contributions

19

■ Full vectorization

- From $O(f(n))$ scalar to $O(f(n)/W)$ vector operations
 - Random accesses excluded
- Principles for good (*efficient*) vectorization
 - Reuse fundamental *operations* across multiple vectorizations
 - Favor *vertical vectorization* by processing different input data per lane
 - Maximize *lane utilization* by executing different things per lane subset

■ Vectorize *basic* operators & build *advanced* operators (in memory)

- Selection scans
- Hash Tables
- Partitioning
- Sorting
- Joins

■ Show impact of *good* vectorization

- On *software* (database system) & *hardware* design

Fundamental Operations

20

- Selective *Load*
- Selective *Store*
- Selective *Gather*
- Selective *Scatter*

Fundamental Vector Operations

21

Selective Load

Vector

A	B	C	D
---	---	---	---

Mask

0	1	0	1
---	---	---	---

Memory

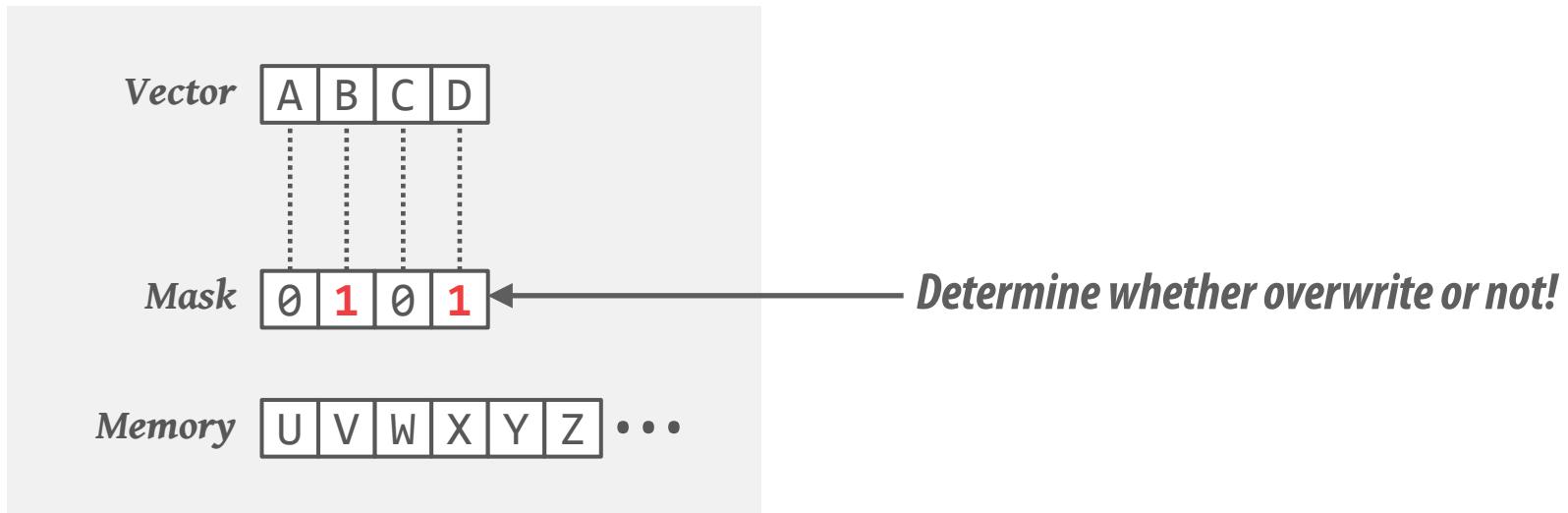
U	V	W	X	Y	Z
---	---	---	---	---	---

 ...

Fundamental Vector Operations

22

Selective Load



Fundamental Vector Operations

23

Selective Load

Vector

A	B	C	D
---	---	---	---

Mask

0	1	0	1
---	---	---	---

Memory

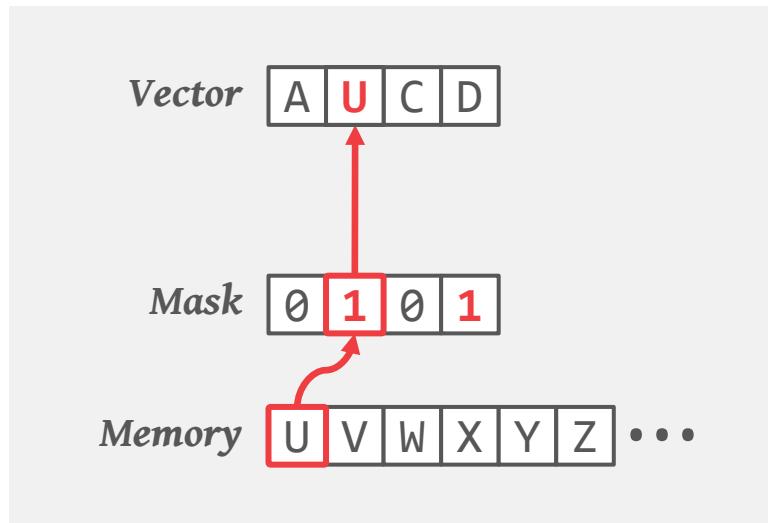
U	V	W	X	Y	Z
---	---	---	---	---	---

 ...

Fundamental Vector Operations

24

Selective Load



Fundamental Vector Operations

25

Selective Load

Vector

A	U	C	D
---	----------	---	---

Mask

0	1	0	1
---	----------	----------	----------

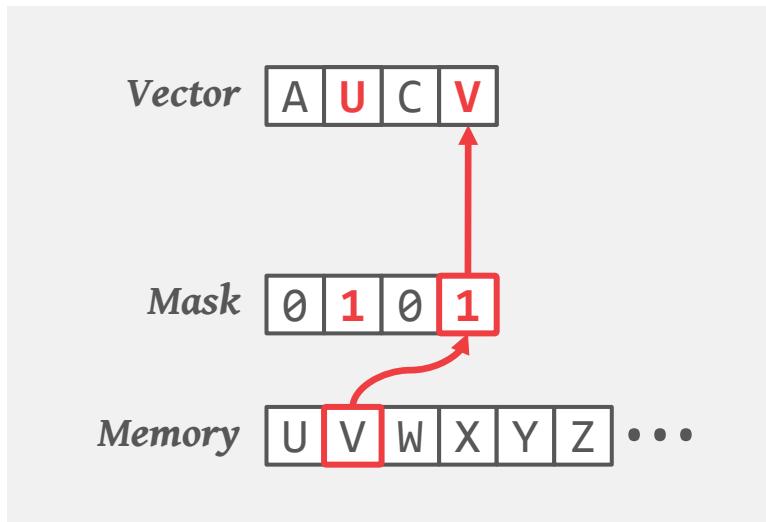
Memory

U	V	W	X	Y	Z	...
---	---	---	---	---	---	-----

Fundamental Vector Operations

26

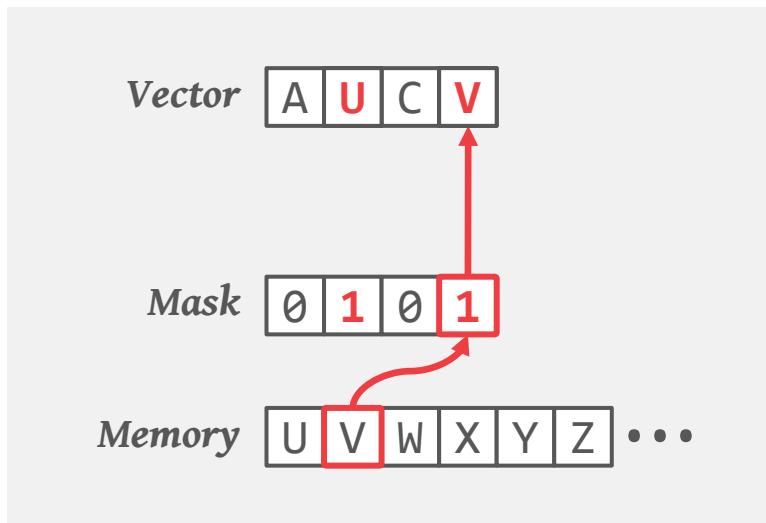
Selective Load



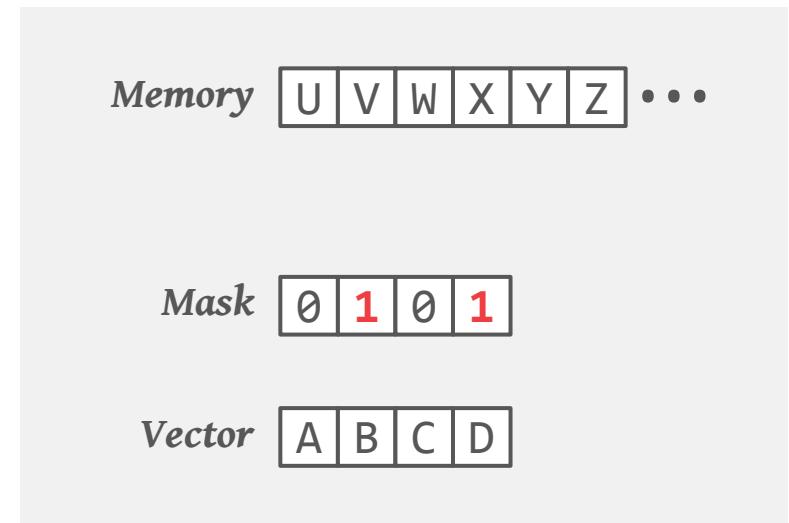
Fundamental Vector Operations

27

Selective Load



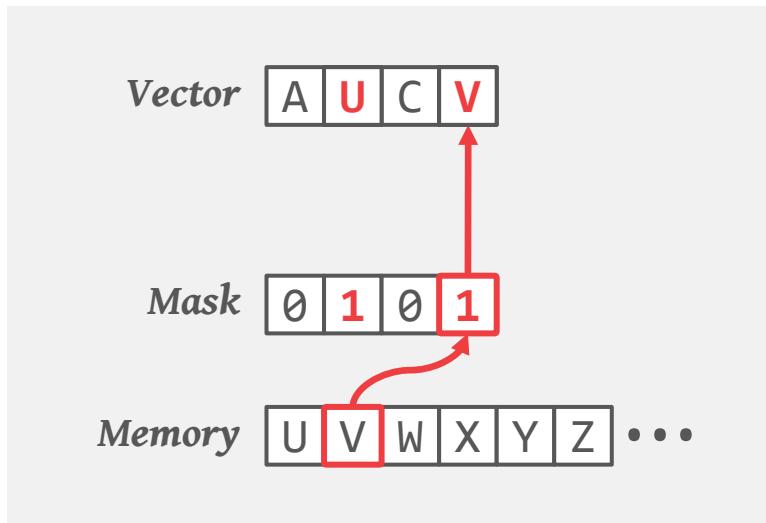
Selective Store



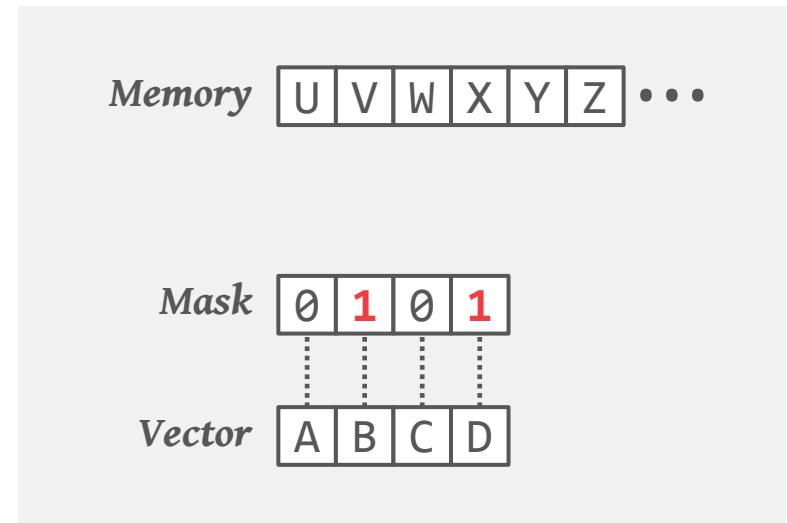
Fundamental Vector Operations

28

Selective Load



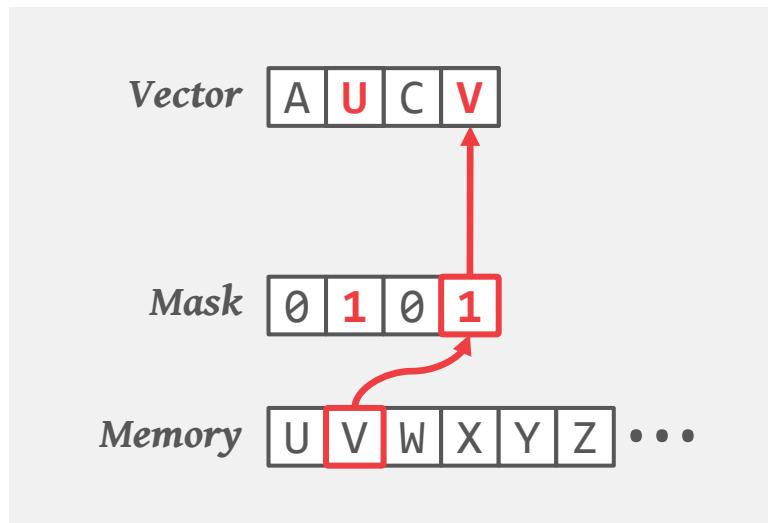
Selective Store



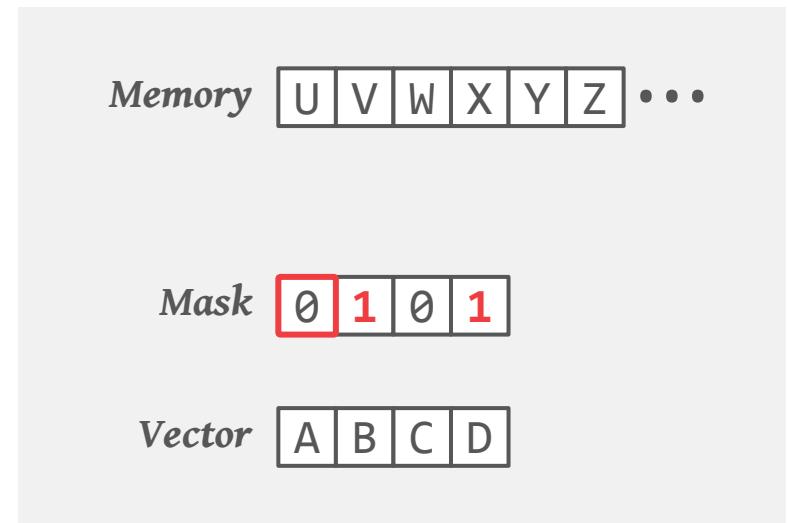
Fundamental Vector Operations

29

Selective Load



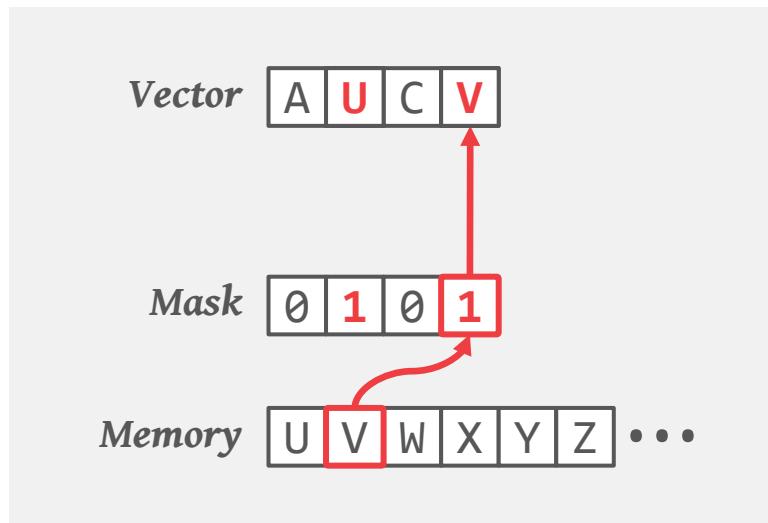
Selective Store



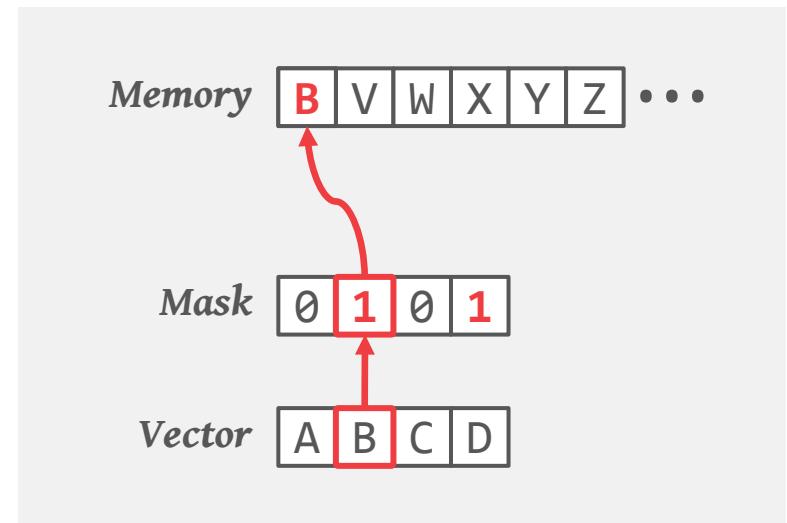
Fundamental Vector Operations

30

Selective Load



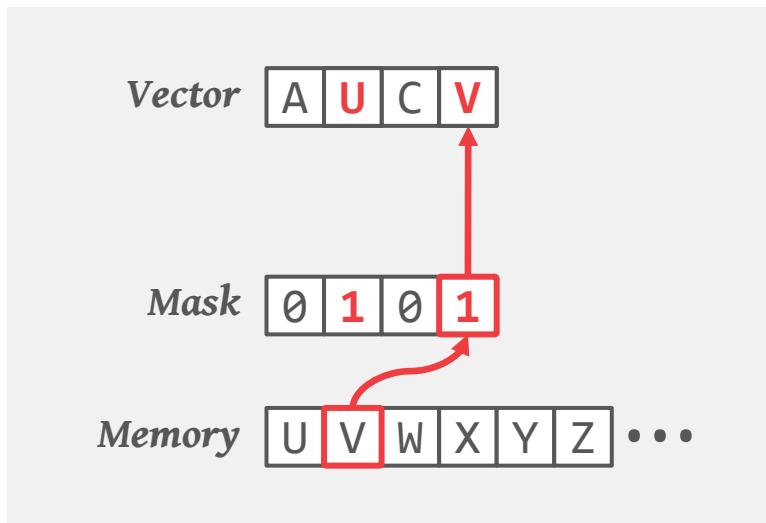
Selective Store



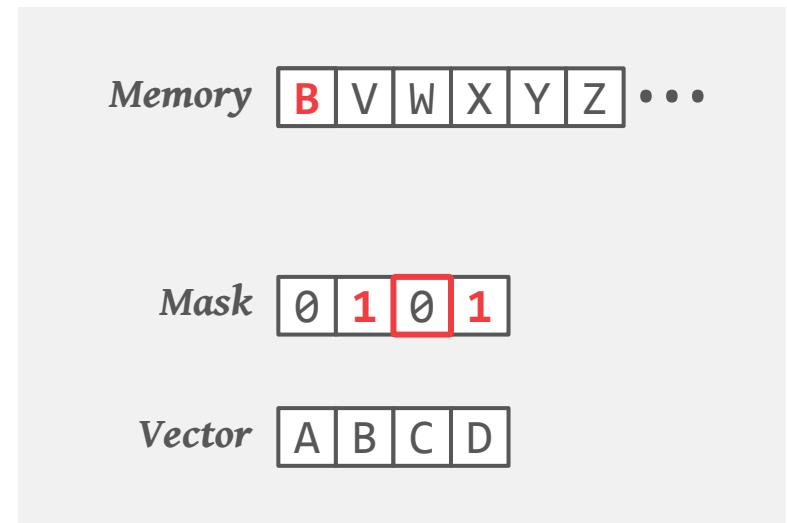
Fundamental Vector Operations

31

Selective Load



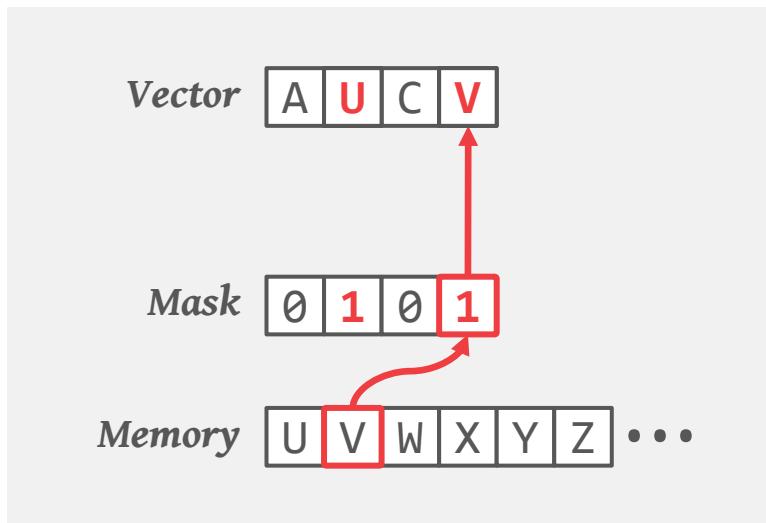
Selective Store



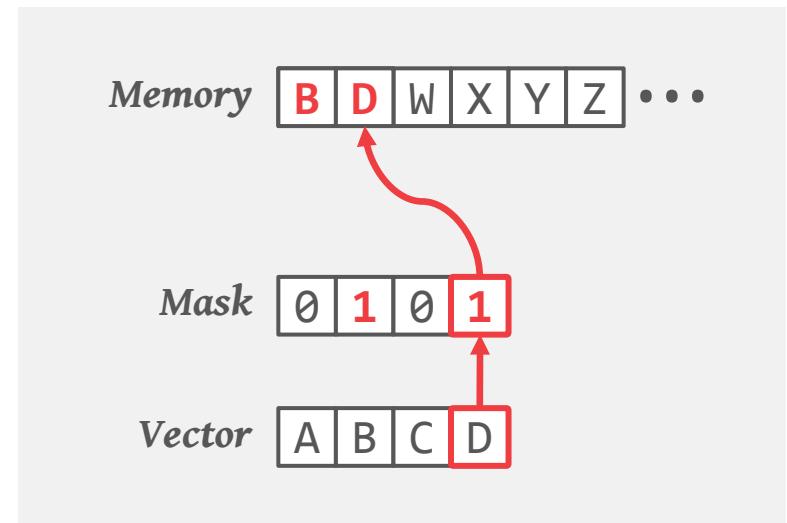
Fundamental Vector Operations

32

Selective Load



Selective Store



Fundamental Vector Operations

33

Selective Gather

Value Vector

A	B	A	D
---	---	---	---

Index Vector

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

Memory

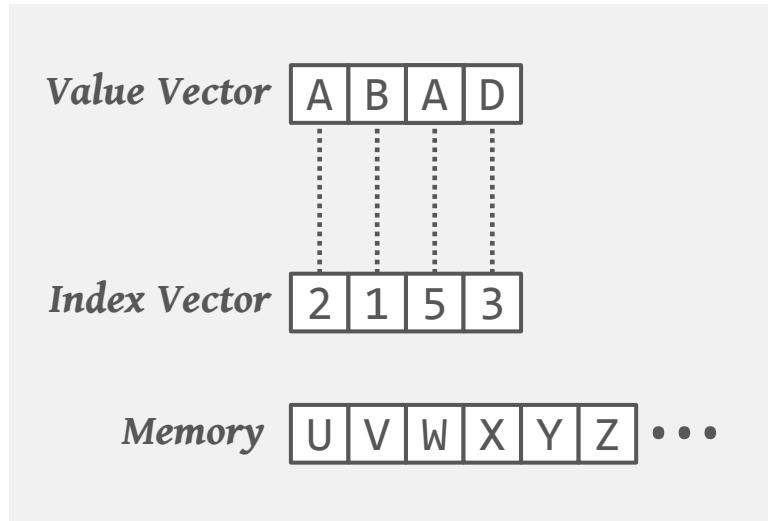
U	V	W	X	Y	Z
---	---	---	---	---	---

 ...

Fundamental Vector Operations

34

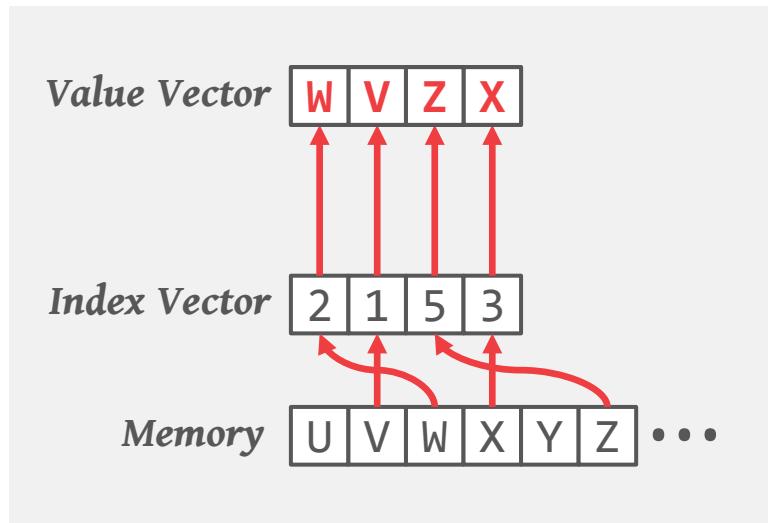
Selective Gather



Fundamental Vector Operations

35

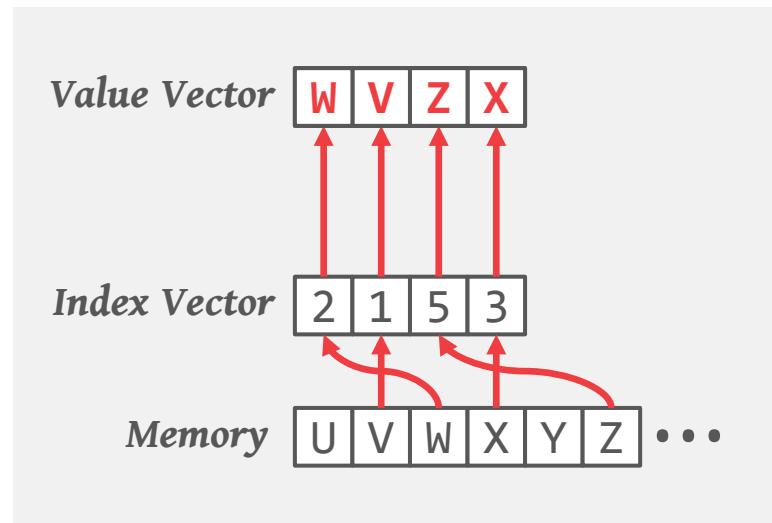
Selective Gather



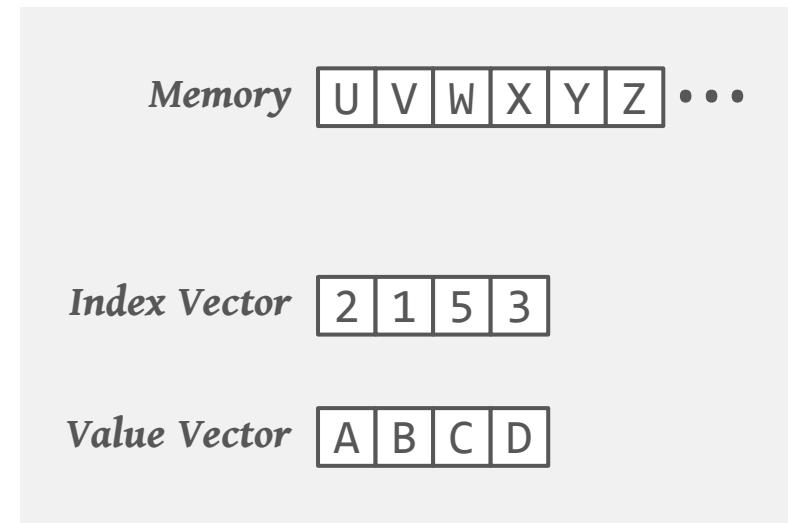
Fundamental Vector Operations

36

Selective Gather



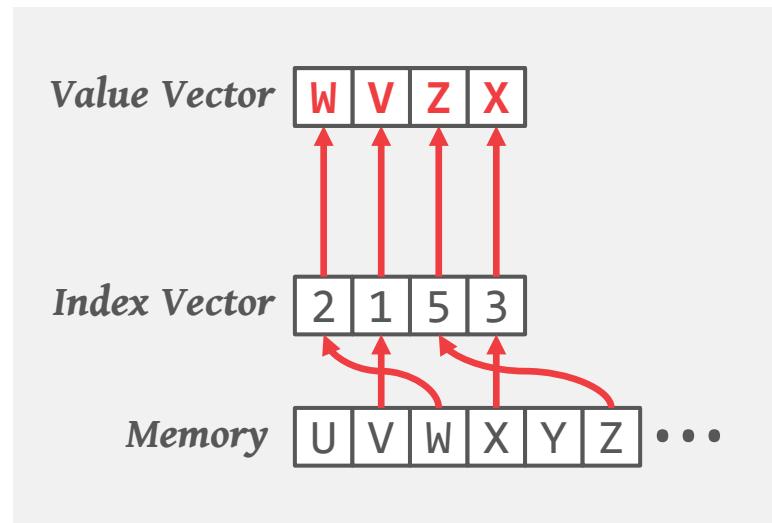
Selective Scatter



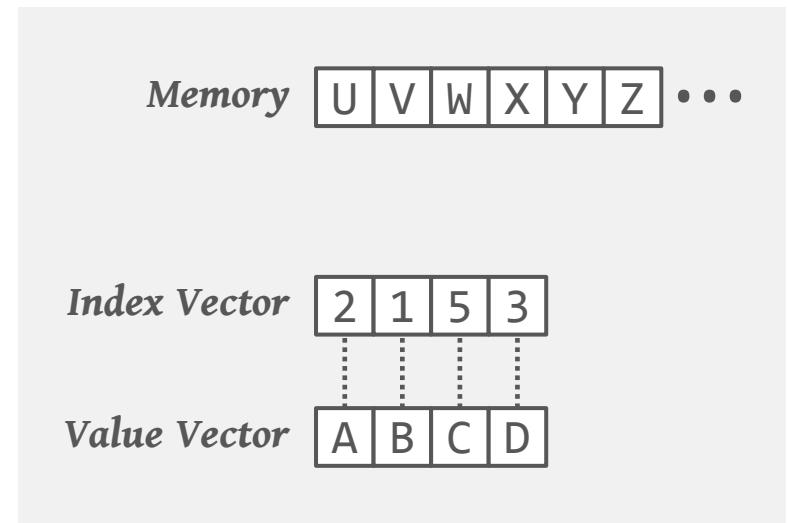
Fundamental Vector Operations

37

Selective Gather



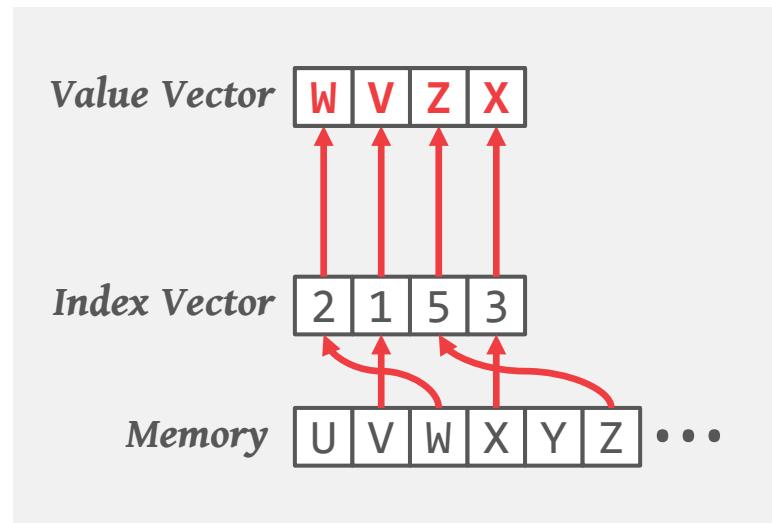
Selective Scatter



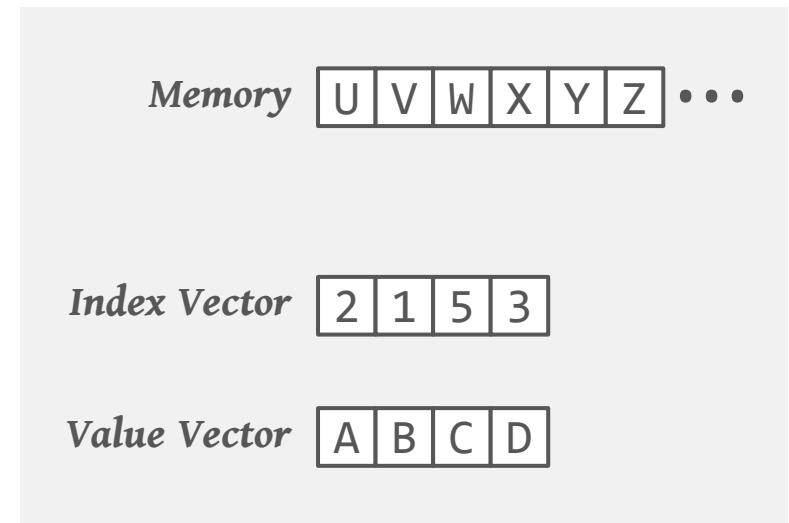
Fundamental Vector Operations

38

Selective Gather



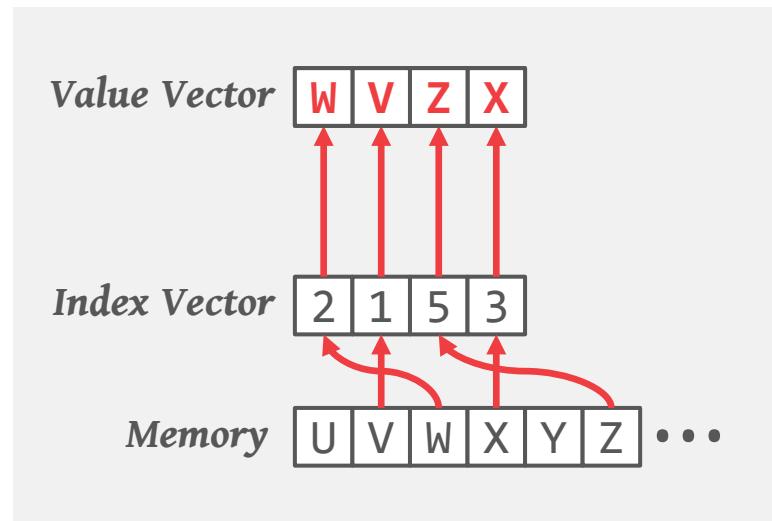
Selective Scatter



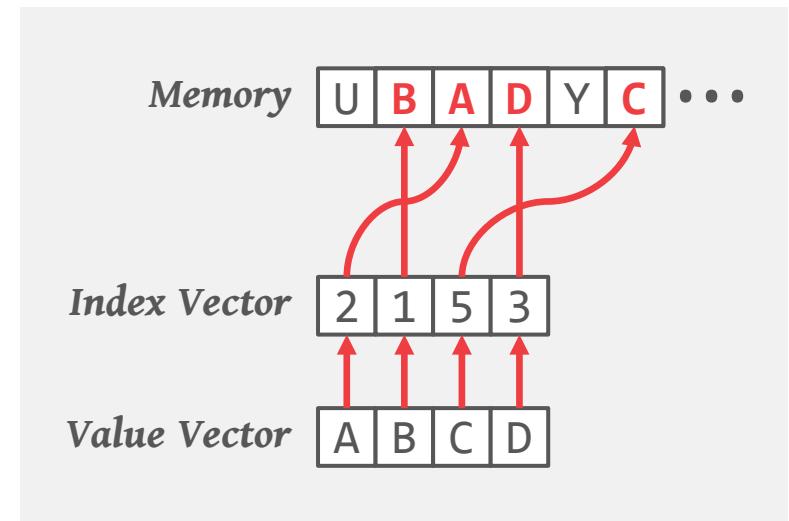
Fundamental Vector Operations

39

Selective Gather



Selective Scatter



Issues

40

- Gathers and scatters are **not** really executed **in parallel** because the *L1 cache only allows one or two distinct accesses per cycle.*
- Gathers are **only supported** in **modern CPUs**.
- Selective loads and stores are also **emulated** in **Xeon CPUs** using **vector permutations**.

Vectorized Operators

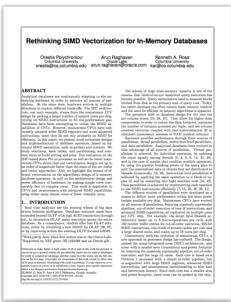
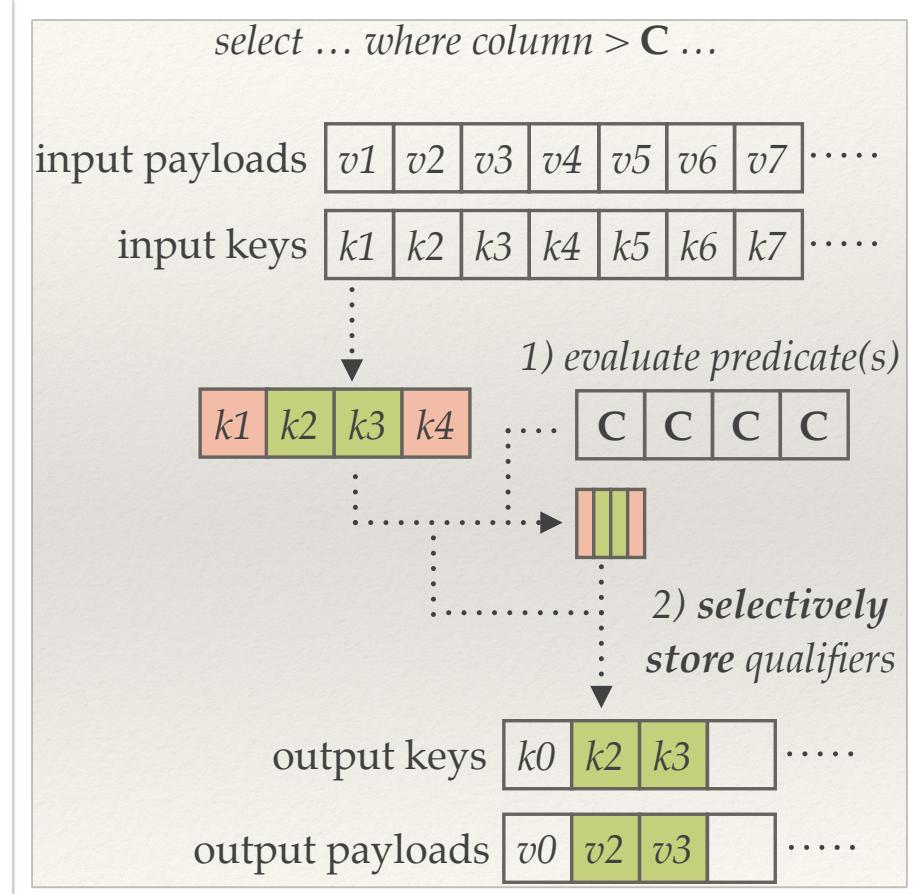
41

■ Selection Scans

■ Hash Tables

■ Partitioning

■ Paper provides additional Info: — *Joins, Sorting, Bloom filters.*



RETHINKING SIMD VECTORIZATION
FOR IN-MEMORY DATABASES
SIGMOD 2015

Selection Scans

42

```
SELECT * FROM table  
WHERE key >= $(low)  
AND key <= $(high)
```

Selection Scans

43

Scalar (Branching)

```
i = 0
for t in table:
    key = t.key
    if (key≥low) && (key≤high):
        copy(t, output[i])
        i = i + 1
```

Selection Scans

44

Scalar (Branching)

```
i = 0
for t in table:
    key = t.key
    if (key≥low) && (key≤high):
        copy(t, output[i])
        i = i + 1
```

branch misprediction penalty!

Selection Scans

45

Scalar (Branching)

```
i = 0
for t in table:
    key = t.key
    if (key≥low) && (key≤high):
        copy(t, output[i])
        i = i + 1
```

Scalar (Branchless)

```
i = 0
for t in table:
    copy(t, output[i])
    key = t.key
    m = (key≥low ? 1 : 0) &&
        ↳(key≤high ? 1 : 0)
    i = i + m
```

Selection Scans

46

Scalar (Branching)

```
i = 0
for t in table:
    key = t.key
    if (key≥low) && (key≤high):
        copy(t, output[i])
        i = i + 1
```

Scalar (Branchless)

```
i = 0
for t in table:
    copy(t, output[i])
    key = t.key
    m = (key≥low ? 1 : 0) &&
        ↳(key≤high ? 1 : 0)
    i = i + m
```

avoid branch misprediction penalty!

Selection Scans

47

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk ≥ low ? 1 : 0) &&
          (vk ≤ high ? 1 : 0)
    if vm ≠ false:
        simdStore(vt, vm, output[i])
        i = i + |vm≠false|
```

Selection Scans

48

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk ≥ low ? 1 : 0) &&
          (vk ≤ high ? 1 : 0)
    if vm ≠ false:
        simdStore(vt, vm, output[i])
        i = i + |vm≠false|
```

v_t is a vector of tuples in table

Selection Scans

49

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk) ←
        vm = (vk ≥ low ? 1 : 0) &&
            ↳(vk ≤ high ? 1 : 0)
    if vm ≠ false:
        simdStore(vt, vm, output[i])
        i = i + |vm≠false|
```

v_k is a key vector for predicate comparison

Selection Scans

50

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk ≥ low ? 1 : 0) &&
          (vk ≤ high ? 1 : 0)
    if vm ≠ false:
        simdStore(vt, vm, output[i])
        i = i + |vm≠false|
```

produce vector mask

Selection Scans

51

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk ≥ low ? 1 : 0) &&
          (vk ≤ high ? 1 : 0)
    if vm ≠ false:
        simdStore(vt, vm, output[i])
        i = i + |vm≠false|
```

SIMD store

Selection Scans

52

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk ≥ low ? 1 : 0) &&
          (vk ≤ high ? 1 : 0)
    if vm ≠ false:
        simdStore(vt, vm, output[i])
        i = i + |vm≠false|
```

Selection Scans (Example)

53

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk ≥ low ? 1 : 0) &&
          ↳(vk ≤ high ? 1 : 0)
    if vm ≠ false:
        simdStore(vt, vm, output[i])
        i = i + |vm≠false|
```

```
SELECT * FROM table
WHERE key >= "O" AND key <= "U"
```

Selection Scans (Example)

54

Vectorized

```
i = 0
for vt in table:
    SIMDLoad(vt.key, vk)
    vm = (vk ≥ low ? 1 : 0) &&
        ↳(vk ≤ high ? 1 : 0)
    if vm ≠ false:
        SIMDStore(vt, vm, output[i])
        i = i + |vm≠false|
```

ID	KEY
1	J
2	O
3	Y
4	S
5	U
6	X

← *Column-oriented layout!*

```
SELECT * FROM table
WHERE key >= "O" AND key <= "U"
```

Selection Scans (Example)

55

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk ≥ low ? 1 : 0) &&
        ↳(vk ≤ high ? 1 : 0)
    if vm ≠ false:
        simdStore(vt, vm, output[i])
        i = i + |vm≠false|
```

ID	KEY
1	J
2	O
3	Y
4	S
5	U
6	X

Key Vector

J O Y S U X

```
SELECT * FROM table
WHERE key >= "O" AND key <= "U"
```

Selection Scans (Example)

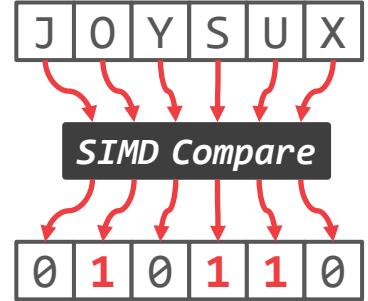
56

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk >= low ? 1 : 0) &&
        (vk <= high ? 1 : 0)
    if vm != false:
        simdStore(vt, vm, output[i])
        i = i + |vm!=false|
```

ID	KEY
1	J
2	O
3	Y
4	S
5	U
6	X

Key Vector



Mask

```
SELECT * FROM table
WHERE key >= "O" AND key <= "U"
```

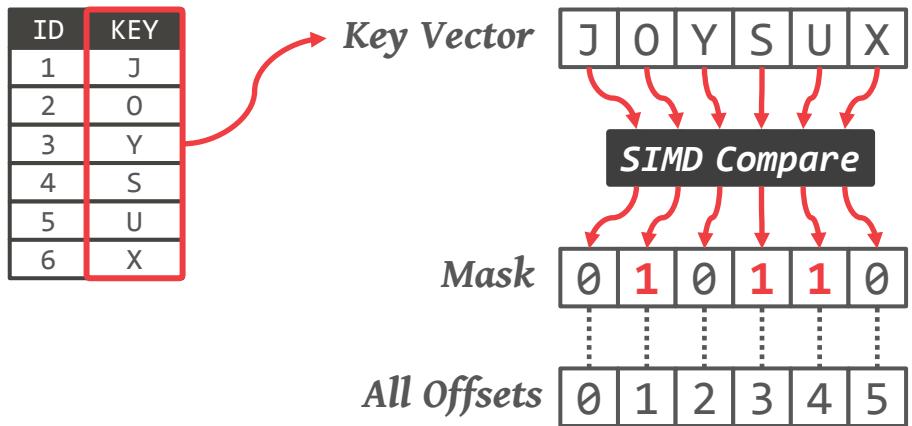
Selection Scans (Example)

57

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk >= low ? 1 : 0) &&
         (vk <= high ? 1 : 0)
    if vm != false:
        simdStore(vt, vm, output[i])
        i = i + |vm!=false|
```

```
SELECT * FROM table
WHERE key >= "O" AND key <= "U"
```



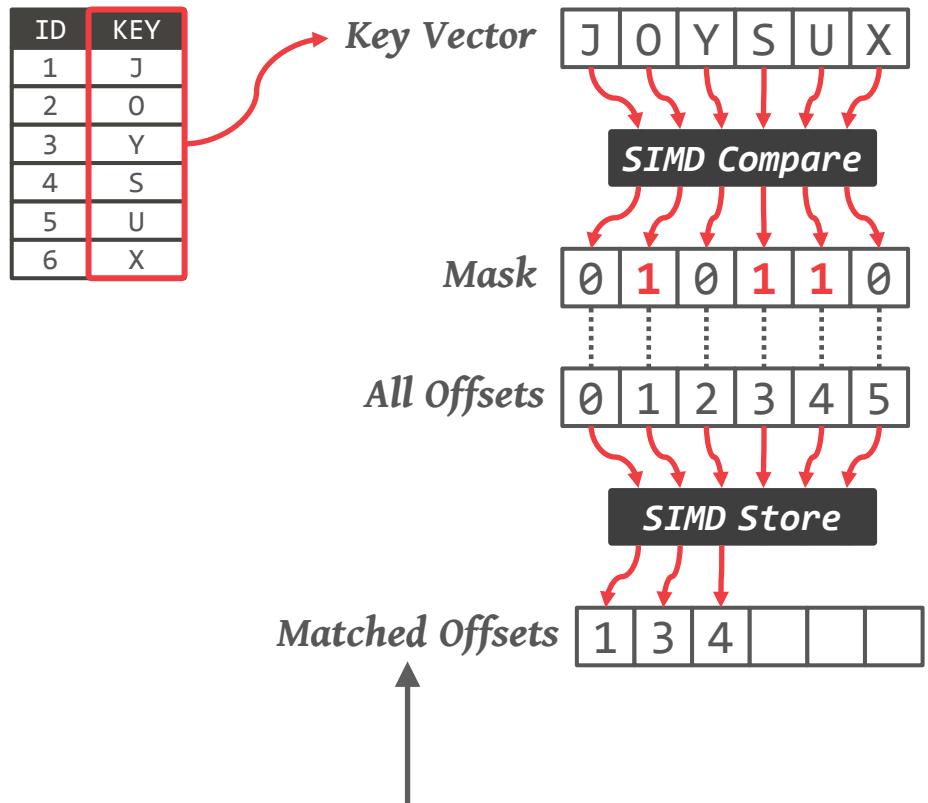
Selection Scans (Example)

58

Vectorized

```
i = 0
for vt in table:
    simdLoad(vt.key, vk)
    vm = (vk >= low ? 1 : 0) &&
         (vk <= high ? 1 : 0)
    if vm != false:
        simdStore(vt, vm, output[i])
        i = i + |vm!=false|
```

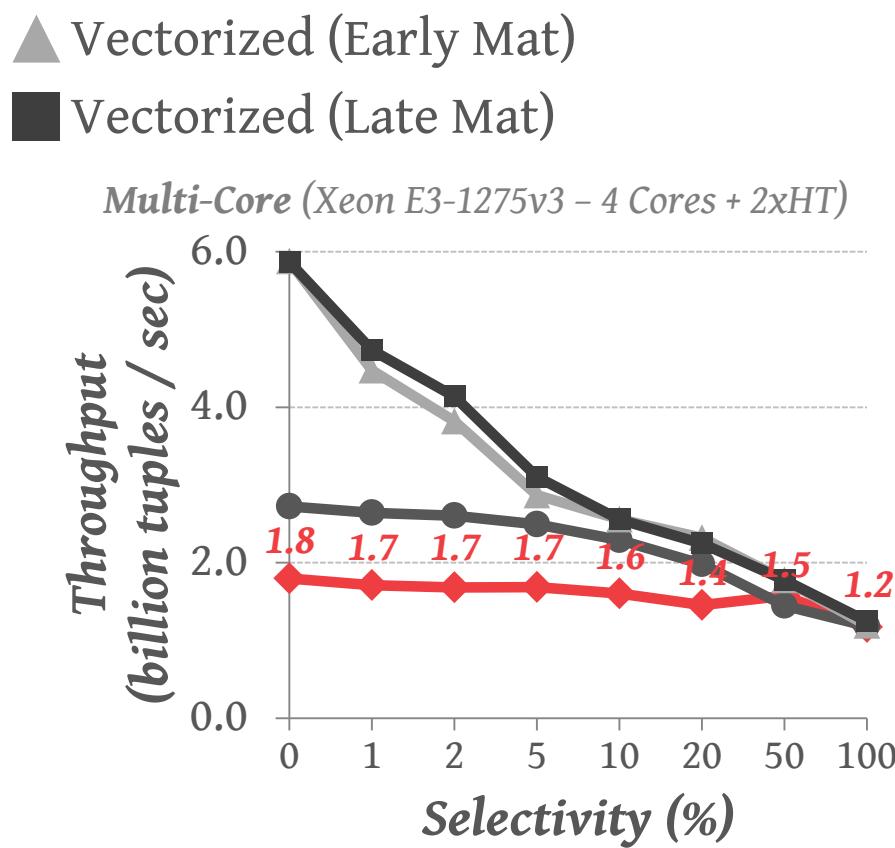
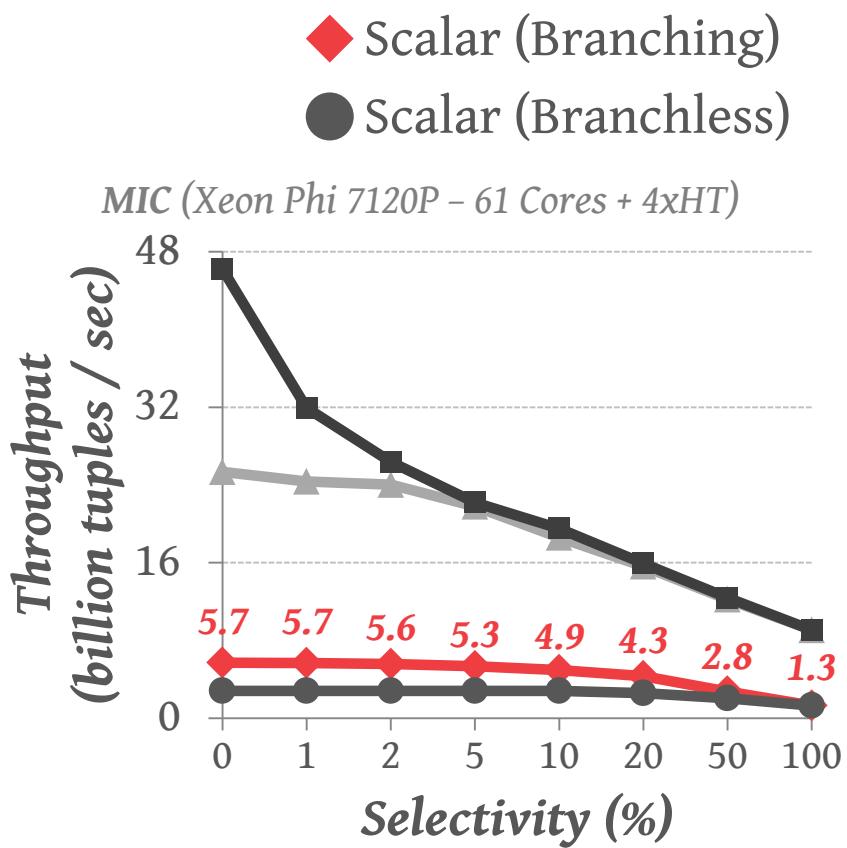
```
SELECT * FROM table
WHERE key >= "O" AND key <= "U"
```



*Further SIMD processing is required or not.
It depends on materialization strategy...*

Performance Comparison: Selection Scans

59



Hash Tables - Probing

60

*selectively load
input keys
(2nd iteration)*

The diagram illustrates a mapping between input keys and hash indexes. On the left, four input keys are listed vertically: *k5*, *k2*, *k3*, and *k6*. An arrow points from this list to a column of hash indexes on the right. The hash indexes are also listed vertically: *h5*, *h2*, *h3*, and *h6*. The mapping is as follows: *k5* maps to *h5*, *k2* maps to *h2*, *k3* maps to *h3*, and *k6* maps to *h6*.

conflicting lanes kept

non-conflicting lanes replaced

linear probing hash table

$k1$	$v1$
$k0$	$v0$
$k4$	$v4$

Assumption:

Use open addressing scheme

Linear Probing Hash Table

Hash Tables - Probing (Scalar)

61

Scalar

Input Key

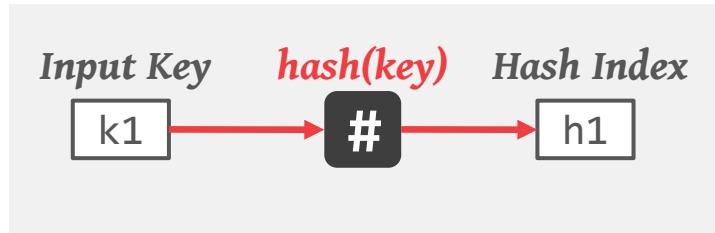
k1

Linear Probing Hash Table

Hash Tables - Probing (Scalar)

62

Scalar

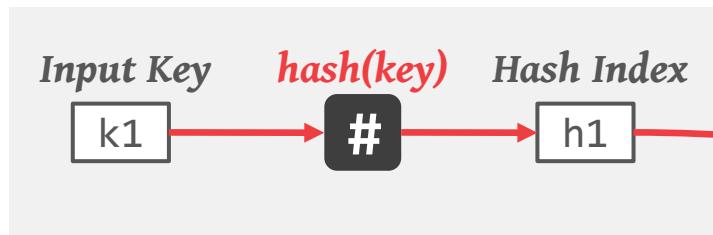


Linear Probing Hash Table

Hash Tables - Probing (Scalar)

63

Scalar



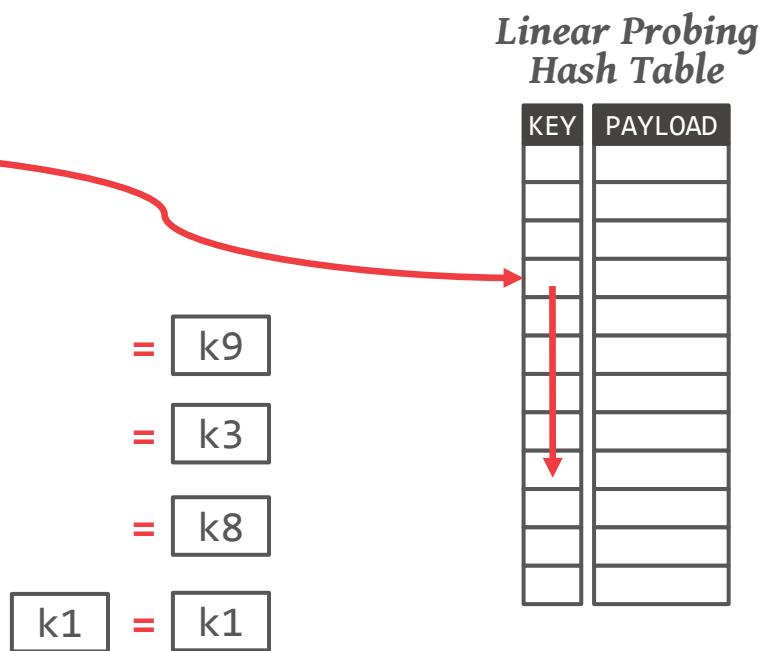
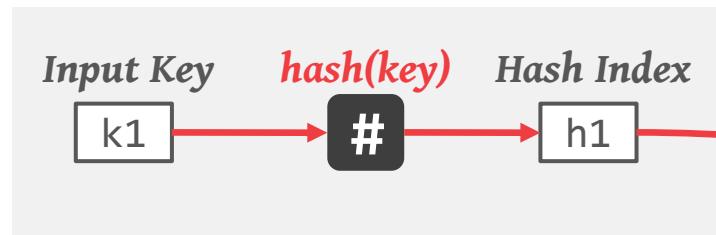
$$\boxed{k1} = \boxed{k9}$$

Linear Probing Hash Table

Hash Tables - Probing (Scalar)

64

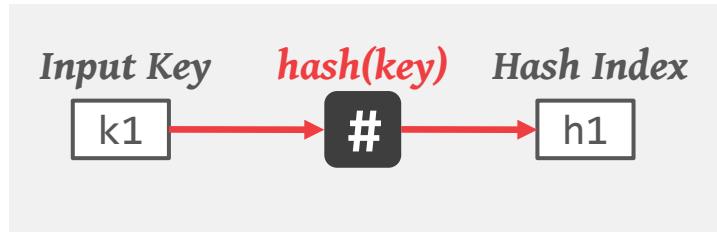
Scalar



Hash Tables - Probing (Vector)

65

Scalar



Vectorized (Horizontal)

Linear Probing Bucketized Hash Table

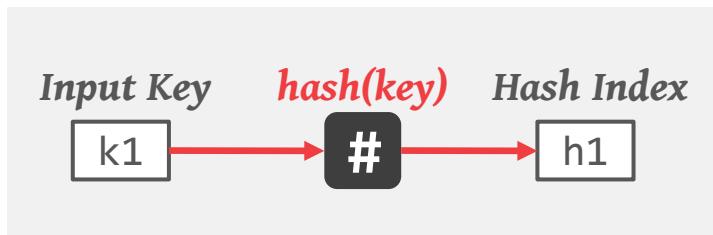
Hash Tables - Probing (Vector)

66

Scalar



Vectorized (Horizontal)

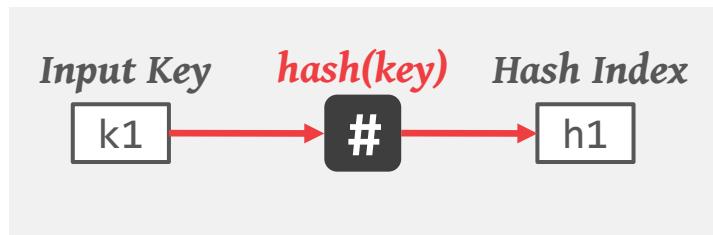


Linear Probing Bucketized Hash Table

Hash Tables - Probing (Vector)

67

Scalar



Vectorized (Horizontal)



The diagram illustrates a **Linear Probing Bucketized Hash Table**. On the left, a red arrow points from a box labeled **k1** to a row of four boxes containing **k9**, **k3**, **k8**, and **k1**. Below this row is a box labeled **SIMD Compare**. To the right is a grid representing the hash table, divided into two columns: **KEYS** and **PAYOUT**. The **KEYS** column has 10 rows, and the **PAYOUT** column has 10 rows. The first row of the **KEYS** column contains four boxes, the second row contains four boxes, and the third row contains four boxes, all highlighted with a red border. The remaining seven rows of the **KEYS** column are empty.

Hash Tables - Probing: Vertical Vectorization

68

Vectorized (Vertical)

*Input Key
Vector*

k1
k2
k3
k4

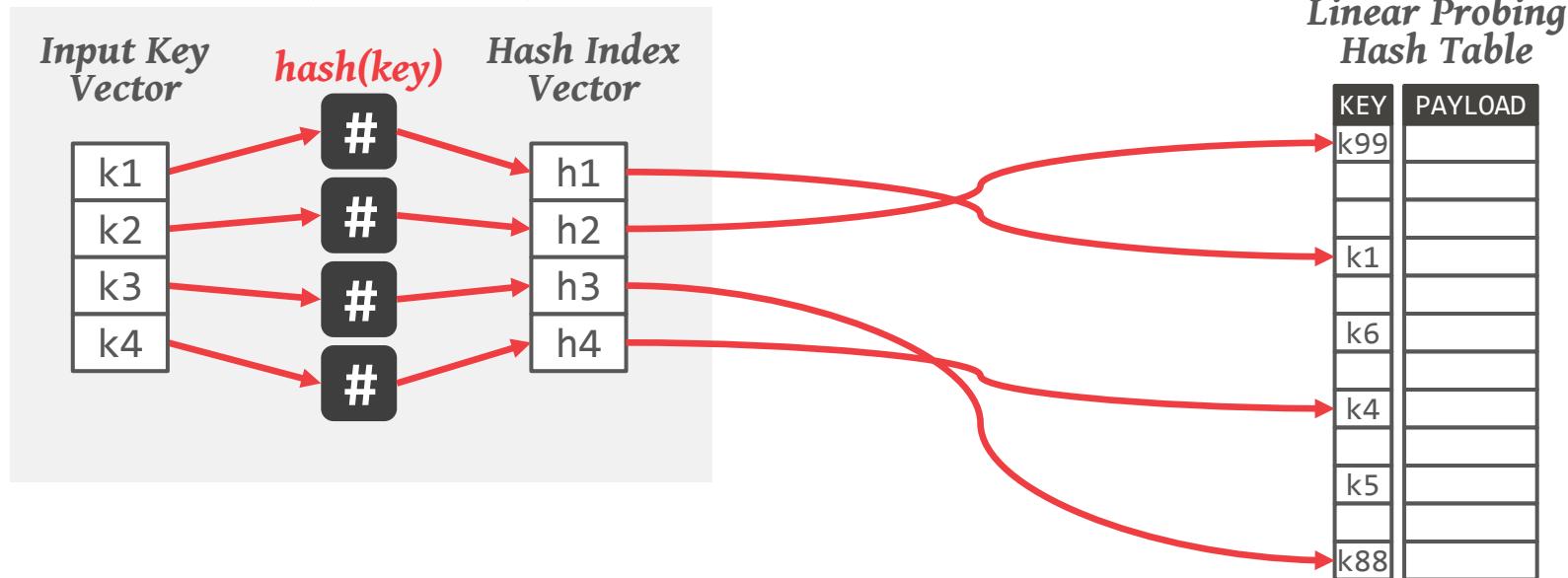
*Linear Probing
Hash Table*

KEY	PAYOUT
k99	
k1	
k6	
k4	
k5	
k88	

Hash Tables - Probing: Vertical Vectorization

69

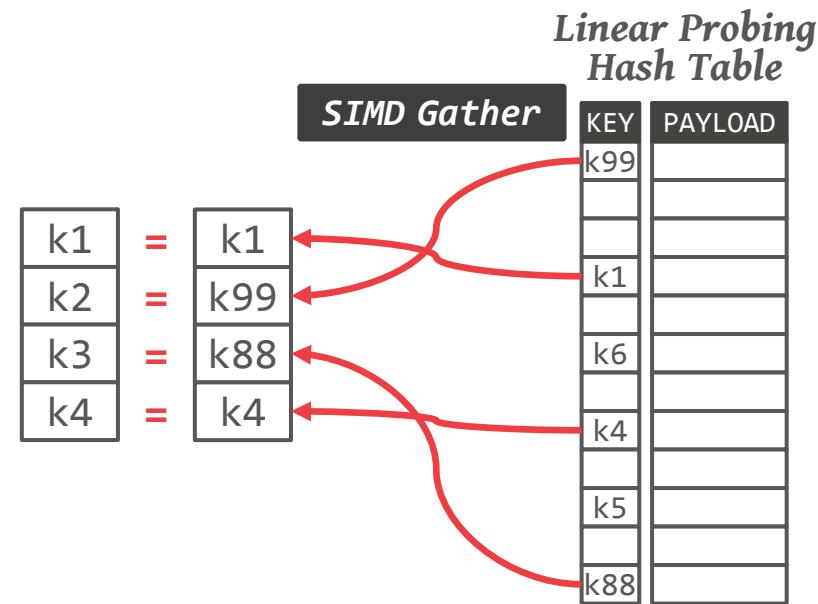
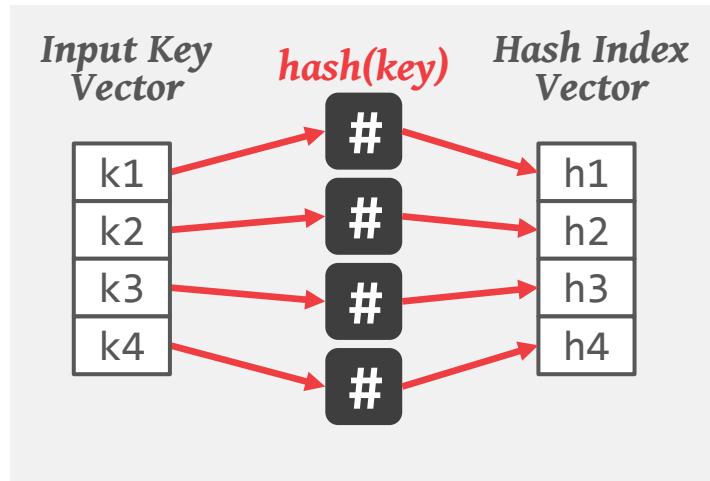
Vectorized (Vertical)



Hash Tables - Probing: Vertical Vectorization

70

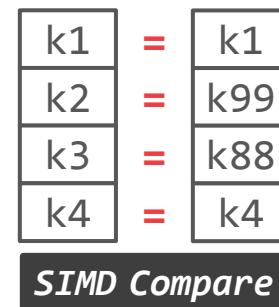
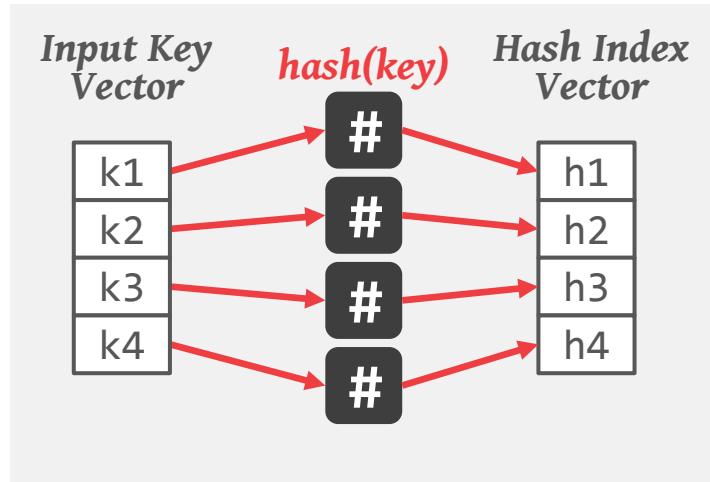
Vectorized (Vertical)



Hash Tables - Probing: Vertical Vectorization

71

Vectorized (Vertical)



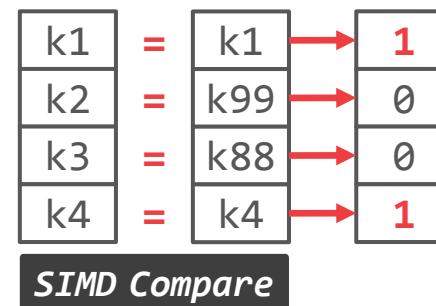
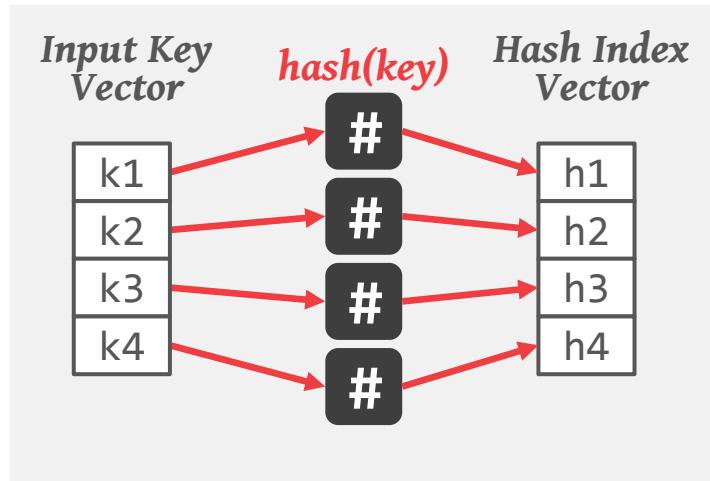
*Linear Probing
Hash Table*

KEY	PAYOUT
k99	
k1	
k6	
k4	
k5	
k88	

Hash Tables - Probing: Vertical Vectorization

72

Vectorized (Vertical)



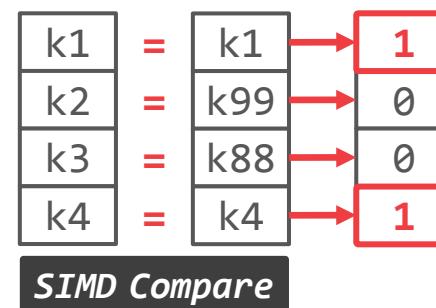
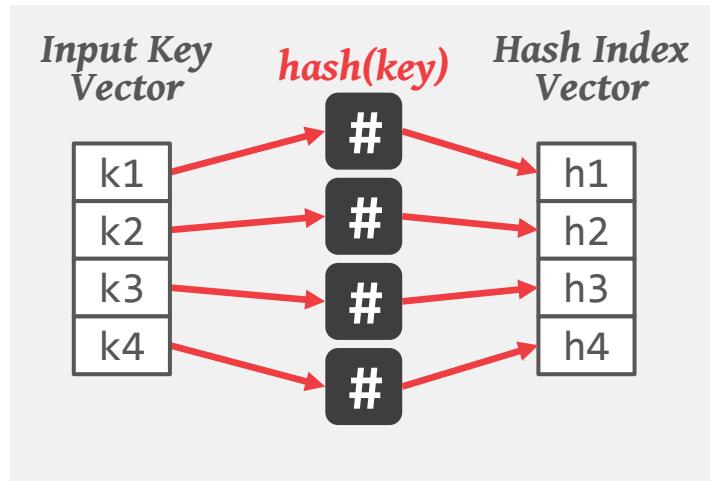
Linear Probing
Hash Table

KEY	PAYOUT
k99	
k1	
k6	
k4	
k5	
k88	

Hash Tables - Probing: Vertical Vectorization

73

Vectorized (Vertical)



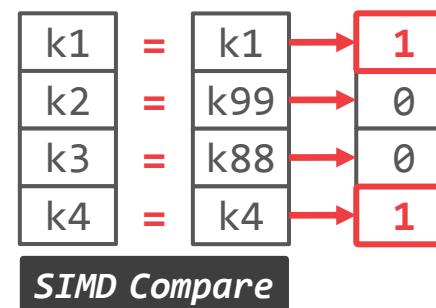
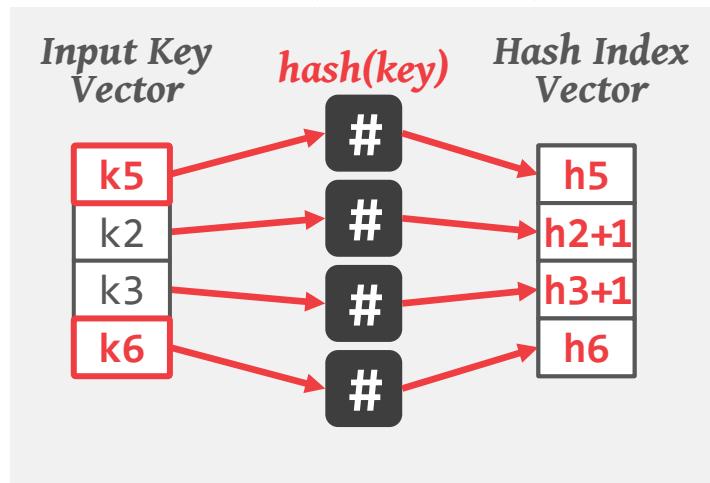
Linear Probing Hash Table

KEY	PAYOUT
k99	
k1	
k6	
k4	
k5	
k88	

Hash Tables - Probing: Vertical Vectorization

74

Vectorized (Vertical)



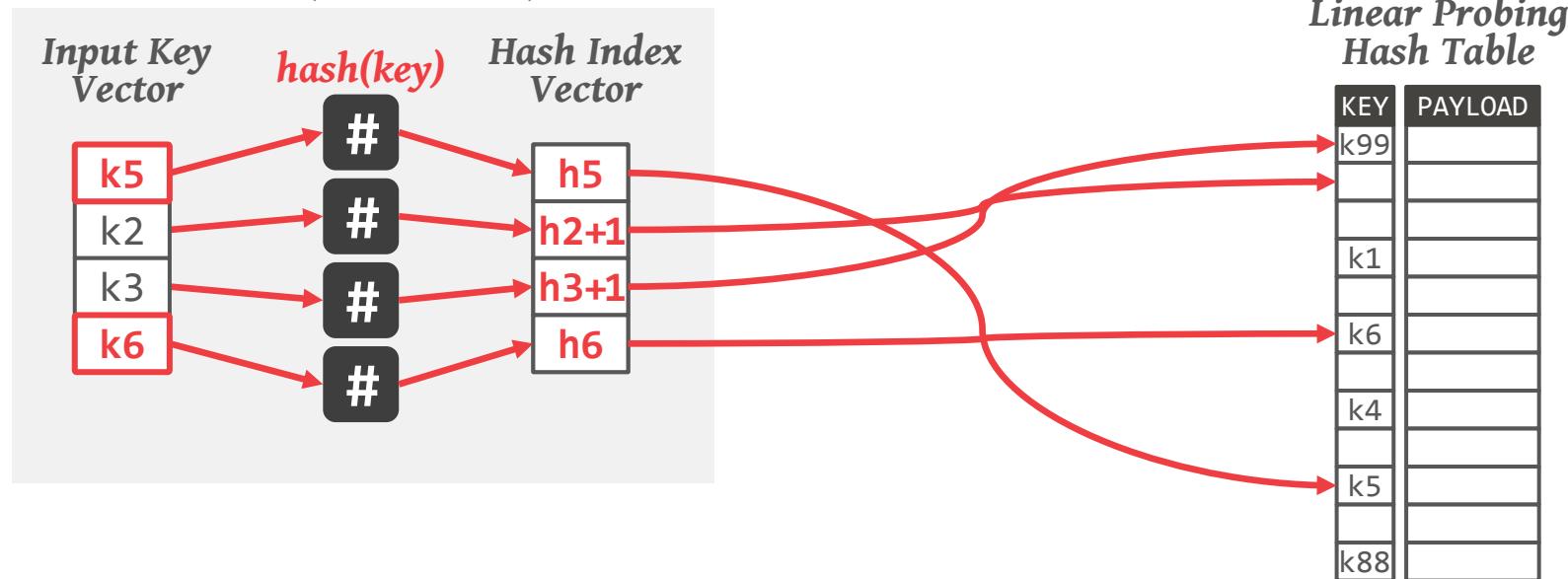
Linear Probing Hash Table

KEY	PAYOUT
k99	
k1	
k6	
k4	
k5	
k88	

Hash Tables - Probing: Vertical Vectorization

75

Vectorized (Vertical)



Hash Tables - Probing

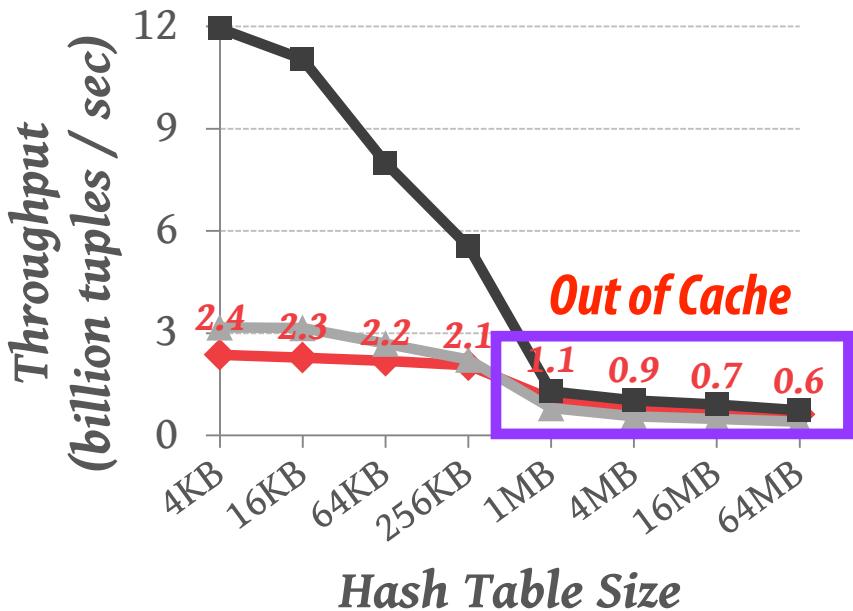
76

◆ Scalar

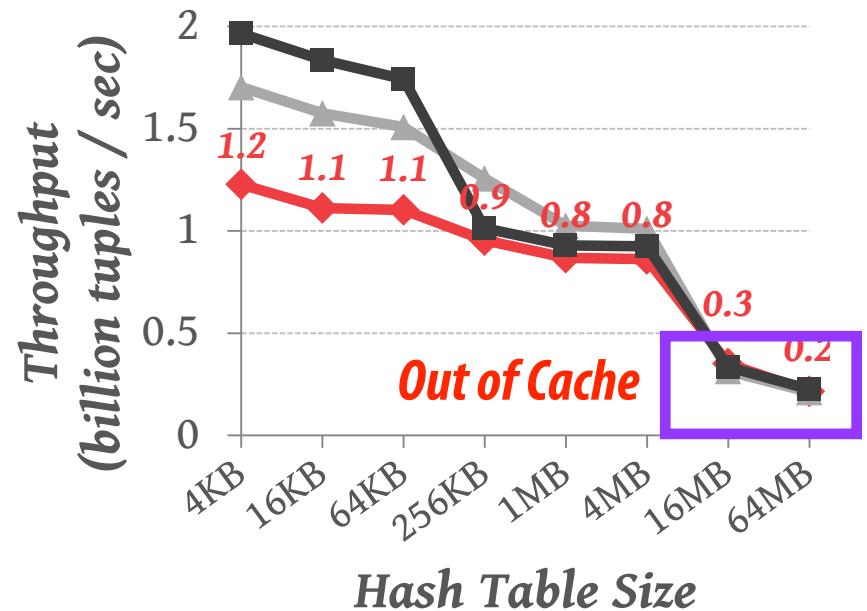
▲ Vectorized (Horizontal)

■ Vectorized (Vertical)

MIC (Xeon Phi 7120P – 61 Cores + 4xHT)



Multi-Core (Xeon E3-1275v3 – 4 Cores + 2xHT)



Partitioning - Histogram

77

- Use *scatter* and *gathers* to **increment counts**.
- Replicate the *histogram* to handle **collisions**.

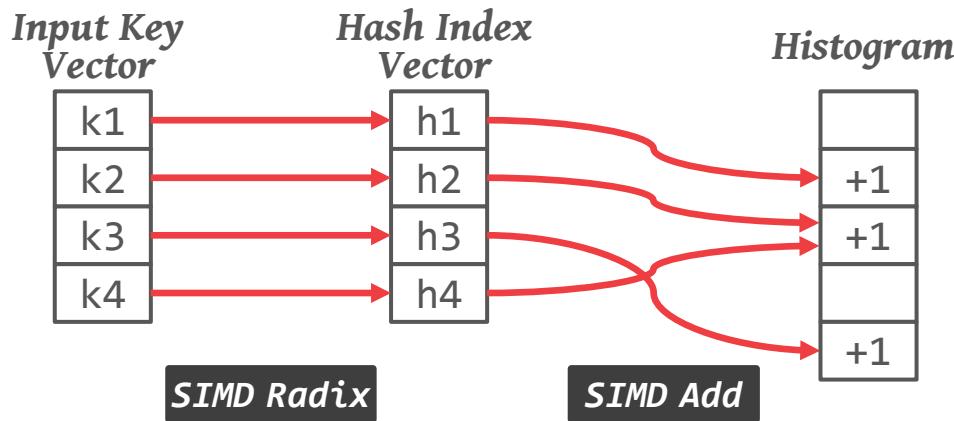
*Input Key
Vector*

k1
k2
k3
k4

Partitioning - Histogram

78

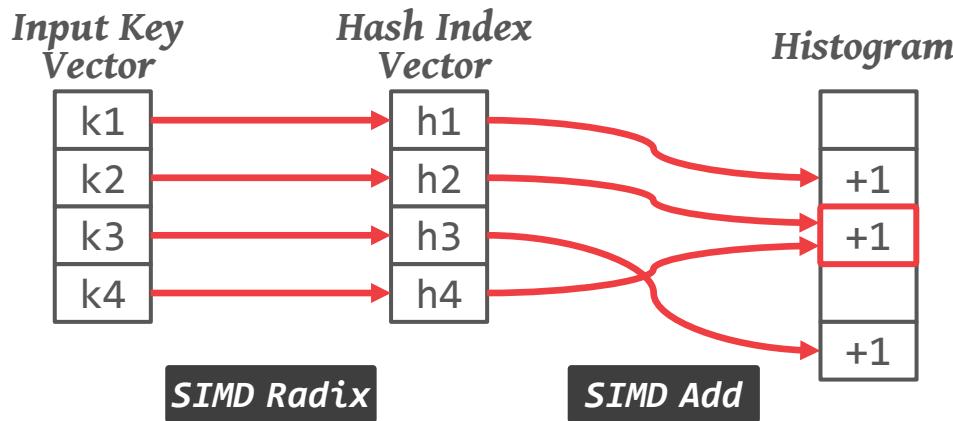
- Use *scatter* and *gathers* to increment counts.
- Replicate the *histogram* to handle collisions.



Partitioning - Histogram

79

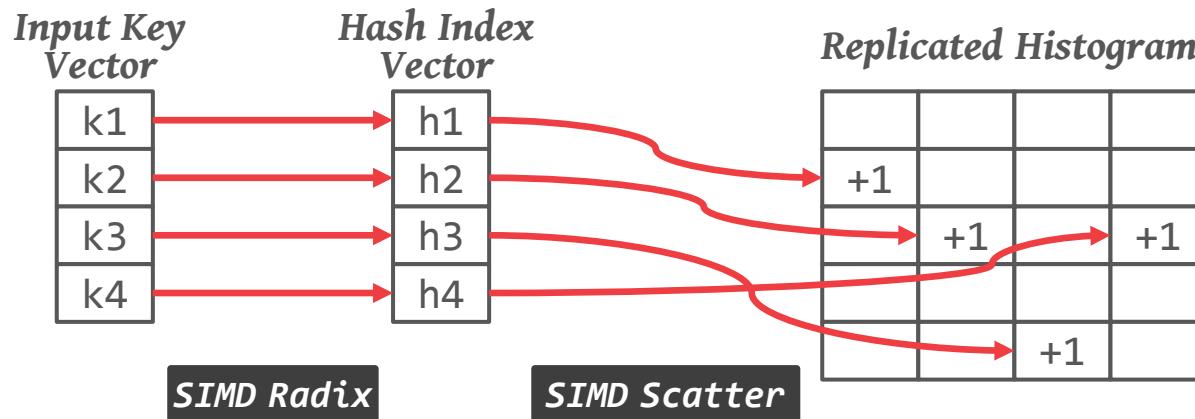
- Use *scatter* and *gathers* to increment counts.
- Replicate the *histogram* to handle collisions.



Partitioning - Histogram

80

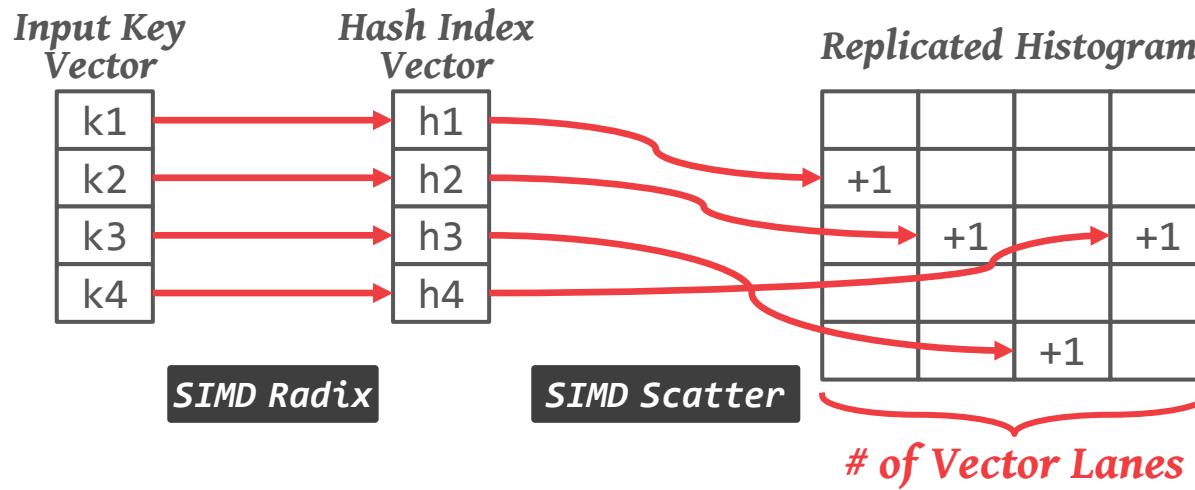
- Use *scatter* and *gathers* to increment counts.
- Replicate the *histogram* to handle collisions.



Partitioning - Histogram

81

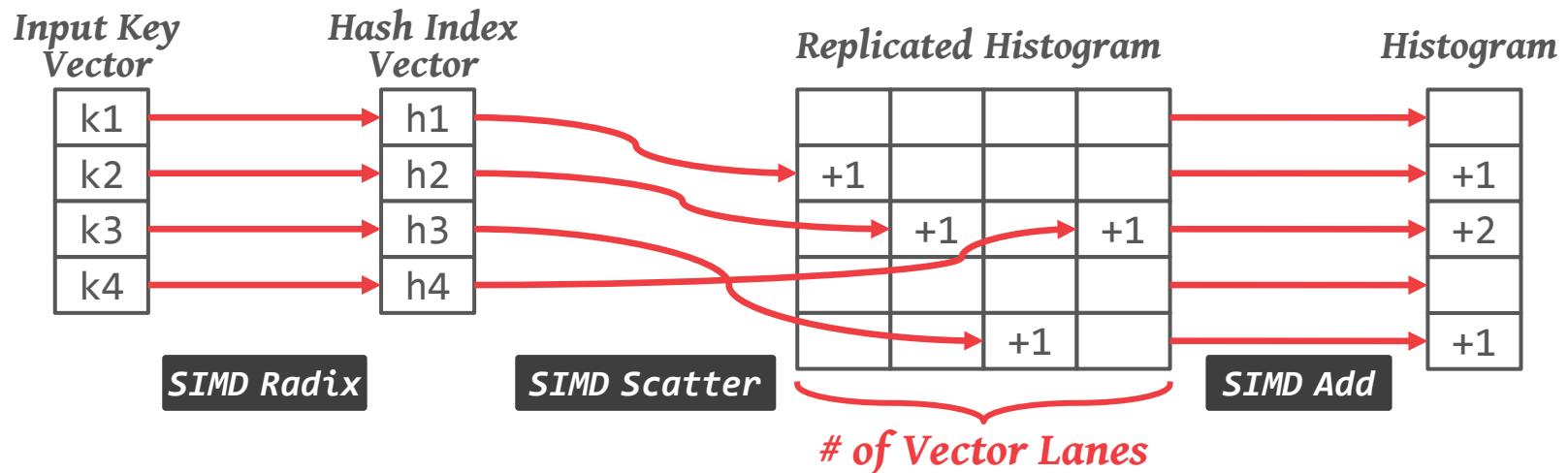
- Use *scatter* and *gathers* to increment counts.
- Replicate the *histogram* to handle collisions.



Partitioning - Histogram

82

- Use *scatter* and *gathers* to increment counts.
- Replicate the *histogram* to handle collisions.



■ No Partitioning

- Build one shared hash table using *atomics*
- *Partially* vectorized

■ Min Partitioning

- *Partition* building table
- Build hash table *per thread*
- *Fully* vectorized

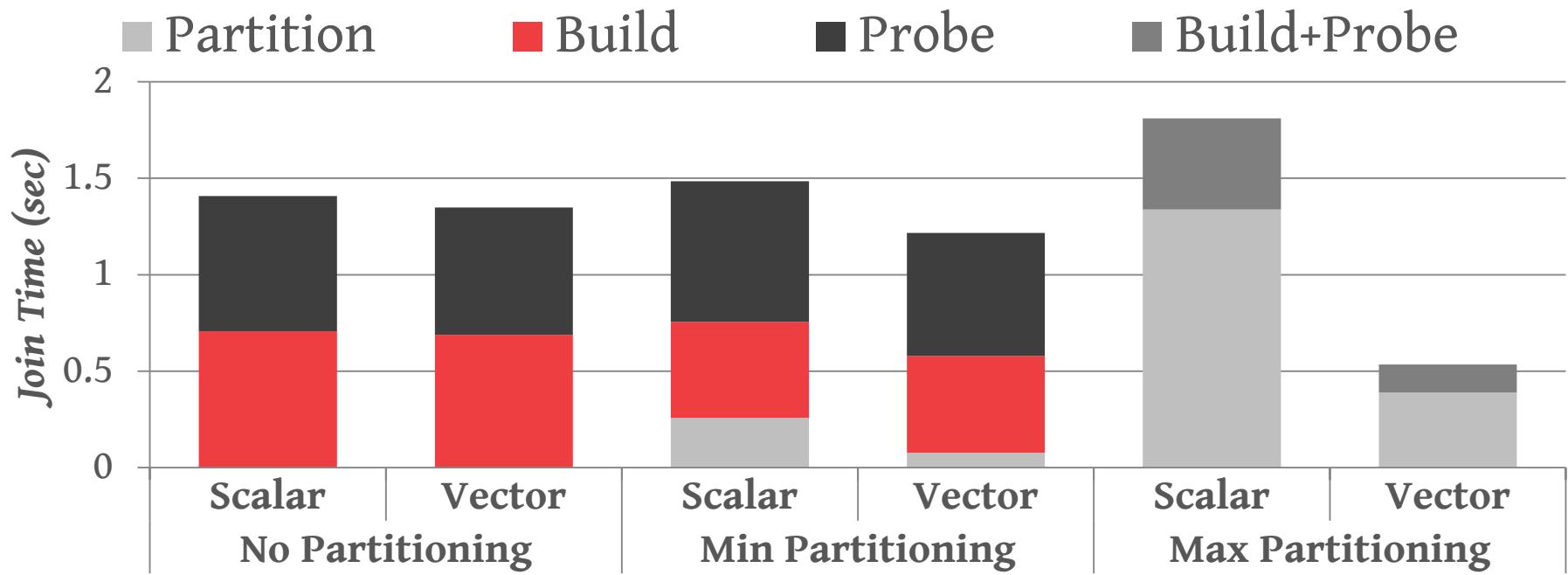
■ Max Partitioning

- Partition both tables repeatedly
- Build and probe *cache-resident hash tables*
- *Fully* vectorized

Joins

84

$200M \bowtie 200M$ tuples (32-bit keys & payloads)
Xeon Phi 7120P - 61 Cores + 4xHT



Main Takeaways

85

- **Vectorization** is essential for **OLAP** queries.
- **Impact on hardware design**
 - *Simple* cores almost as fast as *complex* cores for OLAP
 - 61 simple P54C cores ~ 32 complex Sandy Bridge cores
 - Improved *power efficiency* for analytical databases
- We can combine all the *intra-query parallelism* optimizations.
 - *Multiple threads* processing the same query.
 - Each thread can execute a *compiled plan*.
 - The compiled plan can invoke *vectorized* operations.

Summary

86

- **Background**
- **SIMD Vectorization Algorithms**
 - **Selective Load/Store, Selective Gather/Scatter**
 - **Selective Scans**
 - **Hash Tables - Probing**
 - **Partitioning - Histogram**

