

ADVANCED
DATABASE
SYSTEMS



Vectorized Query Execution

06

Andy Pavlo
CMU 15-721
Spring 2024

**Carnegie
Mellon
University**



LAST CLASS

We described how the DBMS divides up pipelines and runs the in parallel (**Task Parallelization**).

We also discussed how a DBMS will evaluate expressions and introduced the idea of query adaptivity.

TODAY'S AGENDA

Background

Implementation Approaches

Vectorization Fundamentals

Vectorized DBMS Algorithms

VECTORIZATION

The process of converting an algorithm's 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.

This technique is known as Data Parallelization.

WHY THIS MATTERS

Suppose the DBMS can parallelize some algorithm over 32 cores.

Assume each core has a 4-wide SIMD registers.

Potential Speed-up: $32x \times 4x = 128x$

SINGLE INSTRUCTION, MULTIPLE DATA

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.

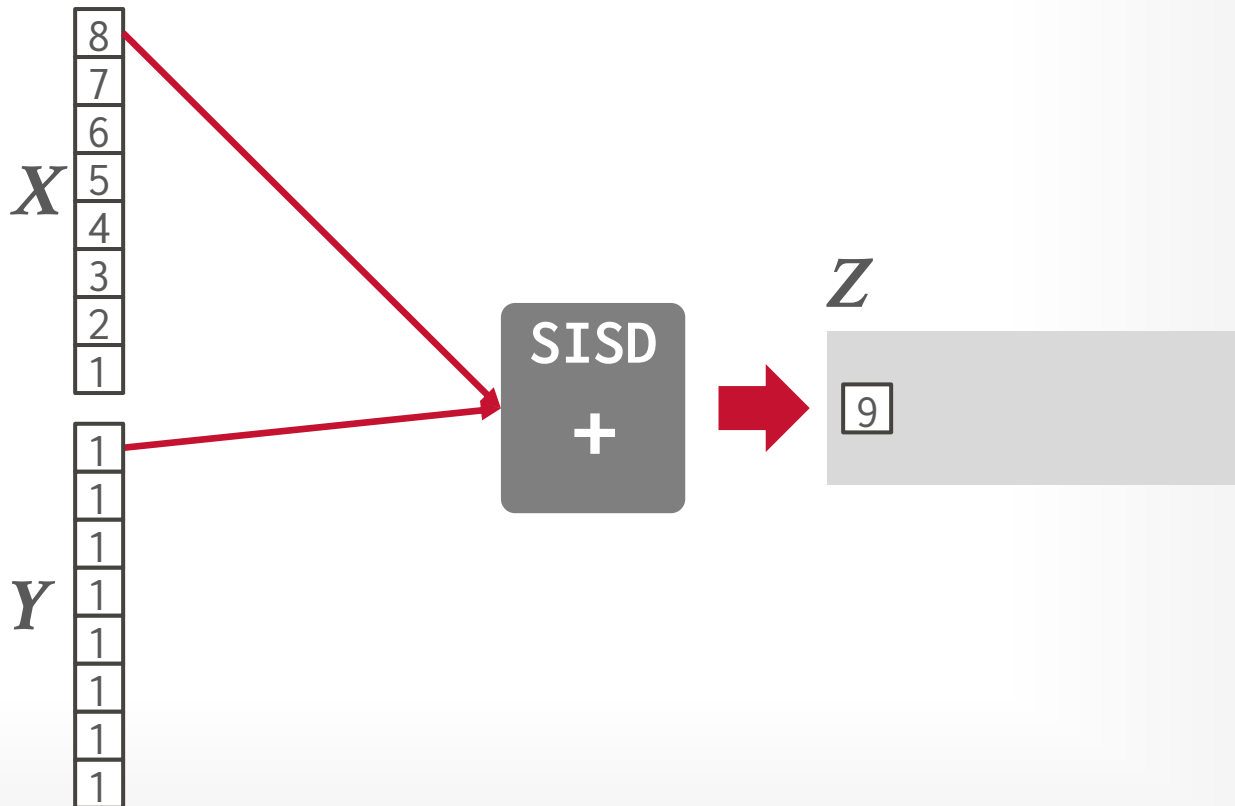
- **x86**: MMX, SSE, SSE2, SSE3, SSE4, AVX, AVX2, AVX512
- **PowerPC**: AltiVec
- **ARM**: NEON, SVE, SVE2
- **RISC-V**: RVV

SIMD EXAMPLE

$$X + Y = Z$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

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



SIMD EXAMPLE

$$X + Y = Z$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

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

X

8
7
6
5
4
3
2
1

Y

1
1
1
1
1
1
1
1

SISD
+



Z

9	8	7	6	5	4	3	2
---	---	---	---	---	---	---	---

SIMD EXAMPLE

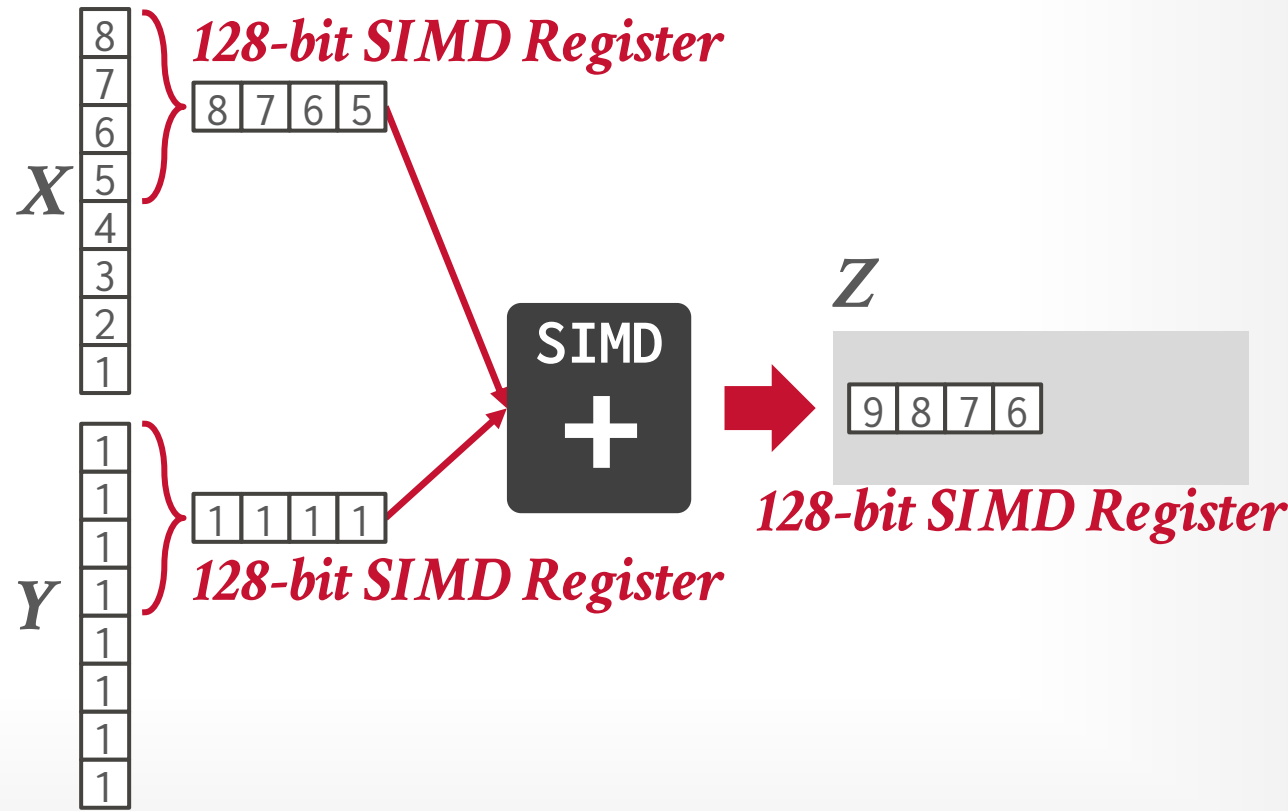
$$X + Y = Z$$

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} + \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} X_1 + Y_1 \\ X_2 + Y_2 \\ \vdots \\ X_n + Y_n \end{bmatrix}$$

```

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

```

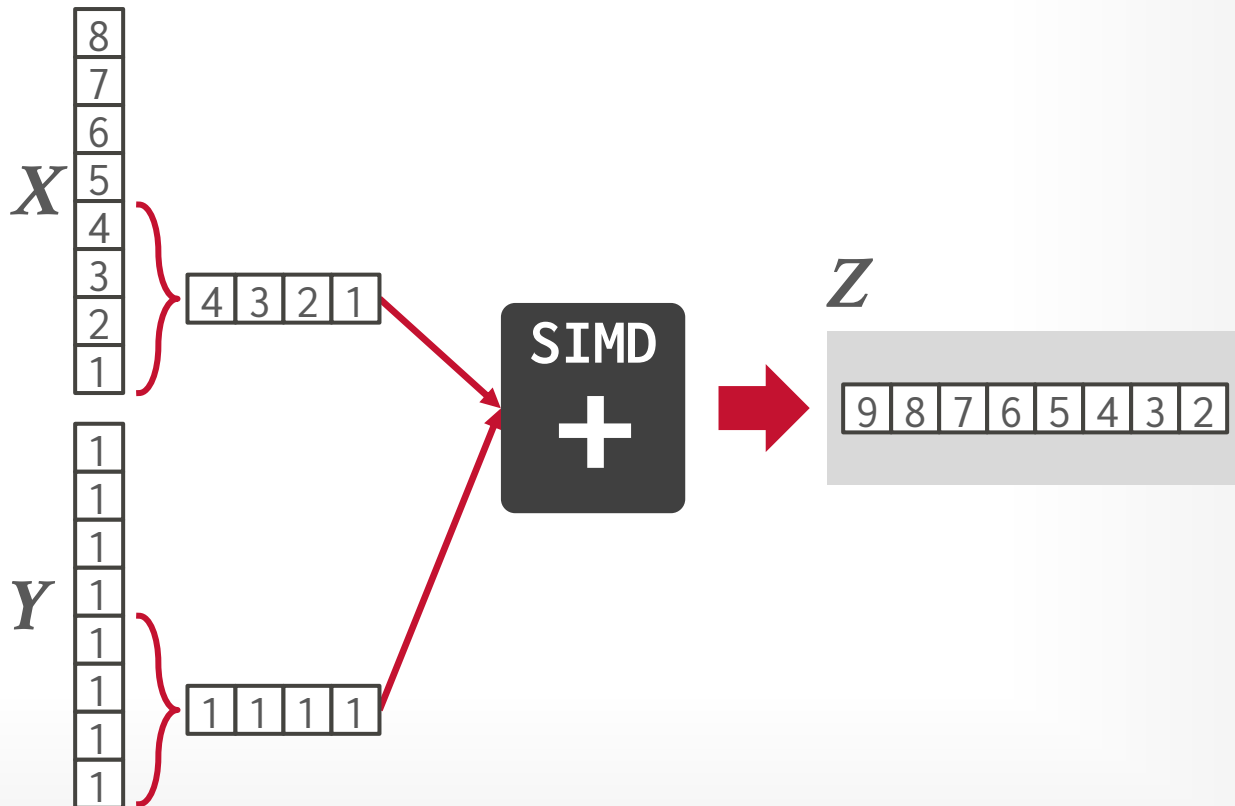


SIMD EXAMPLE

$$X + Y = Z$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

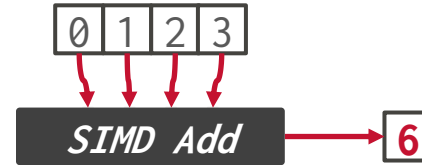
```
for (i=0; i<n; i++) {
    Z[i] = X[i] + Y[i];
}
```



VECTORIZATION DIRECTION

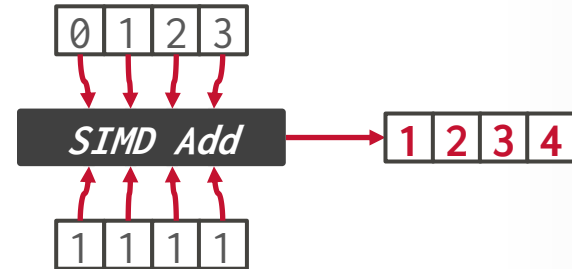
Approach #1: Horizontal

→ Perform operation on all elements together within a single vector.



Approach #2: Vertical

→ Perform operation in an elementwise manner on elements of each vector.



INTEL SIMD EXTENSIONS

		<i>Width</i>	<i>Integers</i>	<i>Single-P</i>	<i>Double-P</i>
1997	MMX	64 bits	✓		
1999	SSE	128 bits	✓	✓(×4)	
2001	SSE2	128 bits	✓	✓	✓(×2)
2004	SSE3	128 bits	✓	✓	✓
2006	SSSE 3	128 bits	✓	✓	✓
2006	SSE 4.1	128 bits	✓	✓	✓
2008	SSE 4.2	128 bits	✓	✓	✓
2011	AVX	256 bits	✓	✓(×8)	✓(×4)
2013	AVX2	256 bits	✓	✓	✓
2017	AVX-512	512 bits	✓	✓(×16)	✓(×8)

Source: [James Reinders](#)

AVX-512

Intel's 512-bit extensions to the AVX2 instructions.

→ Provides new operations to support data conversions, scatter, and permutations.

Unlike previous SIMD extensions, Intel split AVX-512 into groups that CPUs can selectively provide (except for "foundation" extension AVX-512F).

AVX-512

Intel's 512-bit extensions to the AVX2 instructions.
 → Provides new operations to support data conversions, scatter, and permutations.

Subset	F	CD	ER	PF	4FMAPS	4VNNIW	VPOPCNTDQ	VL	DQ	BW	IFMA	VBMI	VNNI	BF16	VBMI2	BITALG	VPCLMULQDQ	GFNI	VAES	VP2INTERSECT	FP16				
Knights Landing (Xeon Phi x200, 2016)	Yes		Yes	No	No																				
Knights Mill (Xeon Phi x205, 2017)					Yes				No																
Skylake-SP, Skylake-X (2017)					No				Yes	No	Yes	Yes	No	Yes	Yes	Yes	No	Yes	Yes	No	No	No	No	No	No
Cannon Lake (2018)					No																				
Cascade Lake (2019)					No				No																
Cooper Lake (2020)					No				No																
Ice Lake (2019)					No				No																
Tiger Lake (2020)					No				No																
Rocket Lake (2021)					No				No																
Alder Lake (2021)					No				No																
Zen 4 (2022)	No				No																				
Sapphire Rapids (2023)	No				No																				
Alder Lake (2021)	Partial ^{Note 1}				Partial ^{Note 1}																				
Zen 4 (2022)	Yes				Yes																				
Sapphire Rapids (2023)	Yes				Yes																				

AVX-512

instructions.
conversions,

	Nehalem (2009), Westmere (2010): Intel Xeon Processors (legacy)	Sandy Bridge (2012): Intel Xeon Processor E3/E5 family	Haswell (2014): Intel Xeon Processor E3 v3/E5 v3/E7 v3 Family	Knights Corner (2012): Intel Xeon Phi Coprocessor x100 Family	Knights Landing (2016): Intel Xeon Phi Processor x200 Family	Skylake (2017): Intel Xeon Scalable Processor Family
						AVX-512VL
						AVX-512DQ
						AVX-512BW
					512-bit	512-bit
					AVX-512ER	
					AVX-512PF	
					AVX-512CD	AVX-512CD
					AVX-512F	AVX-512F
				512-bit		
				IMCI		
			256-bit		AVX2	AVX2
		256-bit			AVX	AVX
					SSE*	SSE*
	128-bit					
	SSE*	SSE*	SSE*			

Legend: — primary instruction set — legacy instruction set

- Knights Landing
- Knights Mill (Xeon)
- Skylake-SP, SkyLake
- Cannon Lake (2021)
- Cascade Lake (2021)
- Cooper Lake (2022)
- Ice Lake (2019)
- Tiger Lake (2020)
- Rocket Lake (2021)
- Alder Lake (2021)
- Zen 4 (2022)
- Sapphire Rapids (2022)

	BMI2	BITALG	VPCLMULQDQ	GFNI	VAES	VP2INTERSECT	FP16
	No						
	No						
		No					
			No				
				No			
					No		
						No	
	Yes					Yes	No
							No
							No
							Yes

Note 1

IMPLEMENTATION

Choice #1: Automatic Vectorization

Choice #2: Compiler Hints

Choice #3: Explicit Vectorization

Ease of Use



*Programmer
Control*

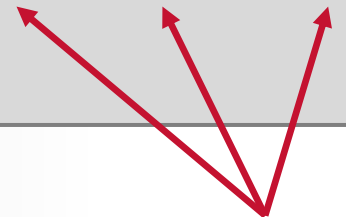
AUTOMATIC VECTORIZATION

The compiler can identify when instructions inside of a loop can be rewritten as a vectorized operation.

Works for simple loops only and is rare in database operators. Requires hardware support for SIMD instructions.

AUTOMATIC VECTORIZATION

```
void add(int *X,  
        int *Y,  
        int *Z) { ← *Z=*X+1  
    for (int i=0; i<MAX; i++) {  
        Z[i] = X[i] + Y[i];  
    }  
}
```



*These might point to
the same address!*

This loop is not legal to automatically vectorize because the code is written such that the addition is described sequentially.

COMPILER HINTS

Provide the compiler with additional information about the code to let it know that is safe to vectorize.

Two approaches:

- Give explicit information about memory locations.
- Tell the compiler to ignore vector dependencies.

COMPILER HINTS

```
void add(int *restrict X,  
        int *restrict Y,  
        int *restrict Z) {  
    for (int i=0; i<MAX; i++) {  
        Z[i] = X[i] + Y[i];  
    }  
}
```

The **restrict** keyword in C/C++ tells the compiler that the arrays are distinct memory locations for the lifetime of the pointers.

The compiler can then infer that it is safe to vectorize operations on those pointers.

COMPILER HINTS

```
void add(int *restrict X,
        int *restrict Y,
        int *restrict Z) {
    for (int i=0; i<N; i++)
        Z[i] = X[i] + Y[i];
}
```

The **restrict** keyword in C/C++ tells the compiler that the arrays

The screenshot shows a code editor window for a file named `vector_hash.cpp`. The code defines a template function `TightLoopHash` that takes several arguments, including `const T *__restrict ldata`, `hash_t *__restrict result_data`, and `const SelectionVector *__restrict sel_vector`. These arguments are highlighted with red boxes. The function body contains nested loops that iterate over data and selection vectors, performing hash operations. The editor interface includes tabs for 'Code' and 'Blame', and a status bar indicating '417 lines (377 loc) · 14.6 KB'.

```
duckdb / src / common / vector_operations / vector_hash.cpp
Code Blame 417 lines (377 loc) · 14.6 KB
28 template <bool HAS_RSEL, class T>
29 static inline void TightLoopHash(const T *__restrict ldata, hash_t *__restrict result_data, const SelectionVector *rsel,
30                                  idx_t count, const SelectionVector *__restrict sel_vector, ValidityMask &mask) {
31     if (!mask.AllValid()) {
32         for (idx_t i = 0; i < count; i++) {
33             auto ridx = HAS_RSEL ? rsel->get_index(i) : i;
34             auto idx = sel_vector->get_index(ridx);
35             result_data[ridx] = HashOp::Operation(ldata[idx], !mask.RowIsValid(idx));
36         }
37     } else {
38         for (idx_t i = 0; i < count; i++) {
39             auto ridx = HAS_RSEL ? rsel->get_index(i) : i;
40             auto idx = sel_vector->get_index(ridx);
41             result_data[ridx] = duckdb::Hash<T>(ldata[idx]);
42         }
43     }
44 }
```

HINTS

restrict keyword in C/C++
the compiler that the arrays

ClickHouse / src / AggregateFunctions / AggregateFunctionSum.h

```
Code Blame 645 lines (535 loc) · 21.2 KB
51 struct AggregateFunctionSumData
105 // Vectorized version
106 template <typename Value>
107 void NO_INLINE addMany(const Value * __restrict ptr, size_t start, size_t end)
108 {
109 #if USE_MULTITARGET_CODE
110     if (isArchSupported(TargetArch::AVX512BW))
111     {
112         addManyImplAVX512BW(ptr, start, end);
113         return;
114     }
115     if (isArchSupported(TargetArch::AVX512F))
116     {
117         addManyImplAVX512F(ptr, start, end);
118         return;
119     }
120     if (isArchSupported(TargetArch::AVX2))
121     {
122         addManyImplAVX2(ptr, start, end);
123         return;
124     }
125     if (isArchSupported(TargetArch::SSE42))
126     {
127         addManyImplSSE42(ptr, start, end);
128         return;
129     }
130 #endif
131 }
```

```
Raw Copy Download Edit
restrict ldata, hash_t * __restrict result_data, const SelectionVector * rsel,
const SelectionVector * __restrict sel_vector, ValidityMask &mask) {
```

```
lex(i) : i;
x);
n(ldata[idx], !mask.RowIsValid(idx));
```

```
x(i) : i;
);
data[idx];
```

COMPILER HINTS

```
void add(int *X,  
        int *Y,  
        int *Z) {  
    #pragma ivdep  
    for (int i=0; i<MAX; i++) {  
        Z[i] = X[i] + Y[i];  
    }  
}
```

This pragma tells the compiler to ignore loop dependencies for the vectors.

It is up to the DBMS developer to make sure that this is correct.

EXPLICIT VECTORIZATION

Use CPU intrinsics to manually marshal data between SIMD registers and execute vectorized instructions.

→ Not portable across CPUs (ISAs / versions).

There are libraries that hide the underlying calls to SIMD intrinsics.

→ [Google Highway](#)

→ [Simd](#)

→ [Expressive Vector Engine \(EVE\)](#)

→ [std::simd](#) (Rust Experimental)

EXPLICIT VECTORIZATION

```
void add(int *X,  
         int *Y,  
         int *Z) {  
    __mm128i *vecX = (__m128i*)X;  
    __mm128i *vecY = (__m128i*)Y;  
    __mm128i *vecZ = (__m128i*)Z;  
    for (int i=0; i<MAX/4; i++) {  
        _mm_store_si128(vecZ++,  
            ↪ _mm_add_epi32(*vecX++,  
                ↪ *vecY++));  
    }  
}
```

Store the vectors in 128-bit SIMD registers.

Then invoke the intrinsic to add together the vectors and write them to the output location.

AUTOMATIC VECTORIZATION

Evaluate how well the compiler can automatically vectorize the Vectorwise primitives.

→ Targets: GCC v7.2, Clang v5.0, ICC v18

ICC was able to vectorize the most primitives using AVX-512:

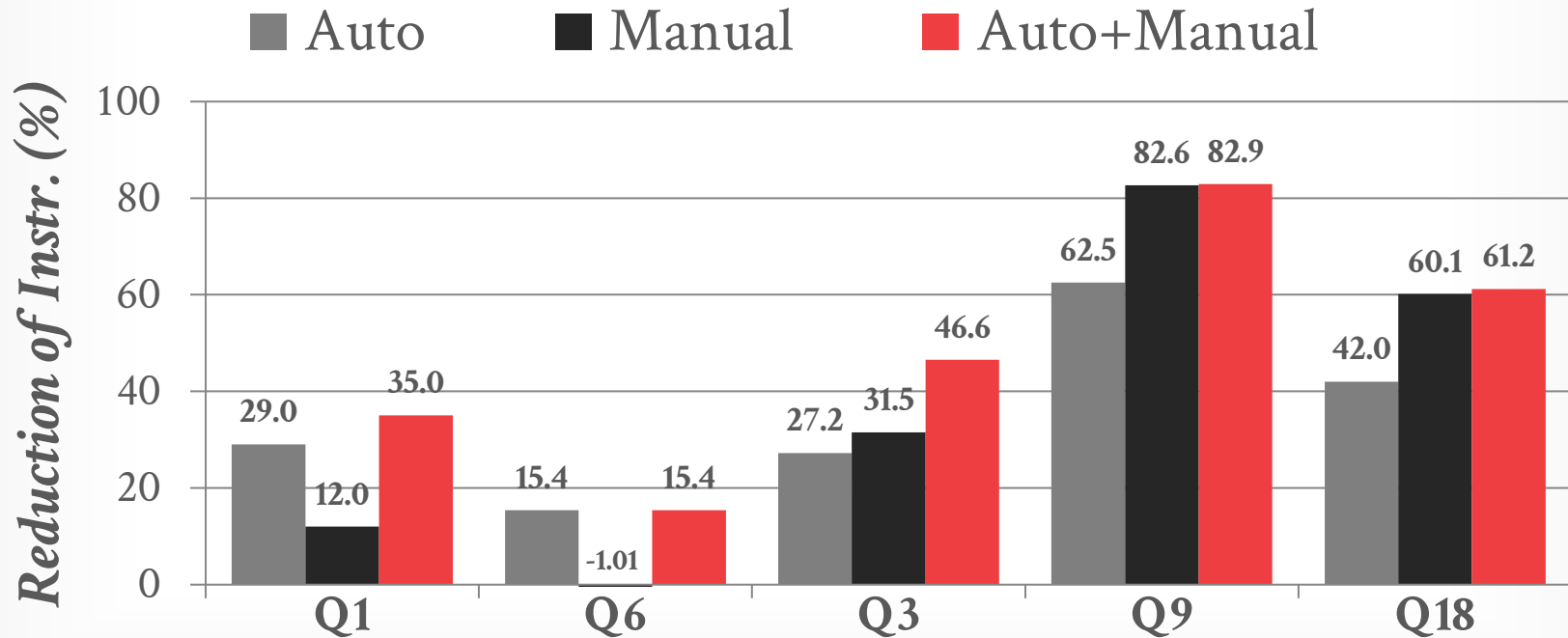
→ Vectorized: Hashing, Selection, Projection

→ Not Vectorized: Hash Table Probing, Aggregation

AUTOMATIC VECTORIZATION

Intel Core i9-7900X (10 cores × 2HT)

Compiler: ICC v18

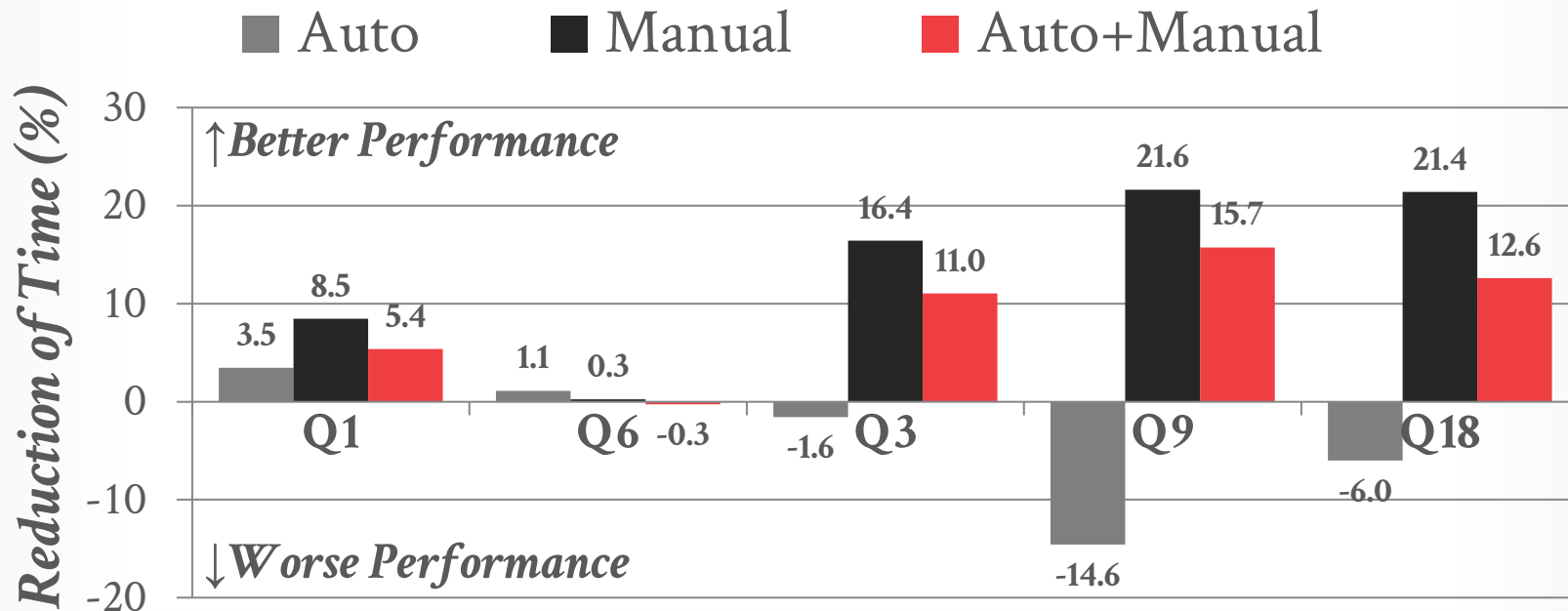


Source: [Timo Kersten](#)

AUTOMATIC VECTORIZATION

Intel Core i9-7900X (10 cores × 2HT)

Compiler: ICC v18



Source: [Timo Kersten](#)

VECTORIZATION FUNDAMENTALS

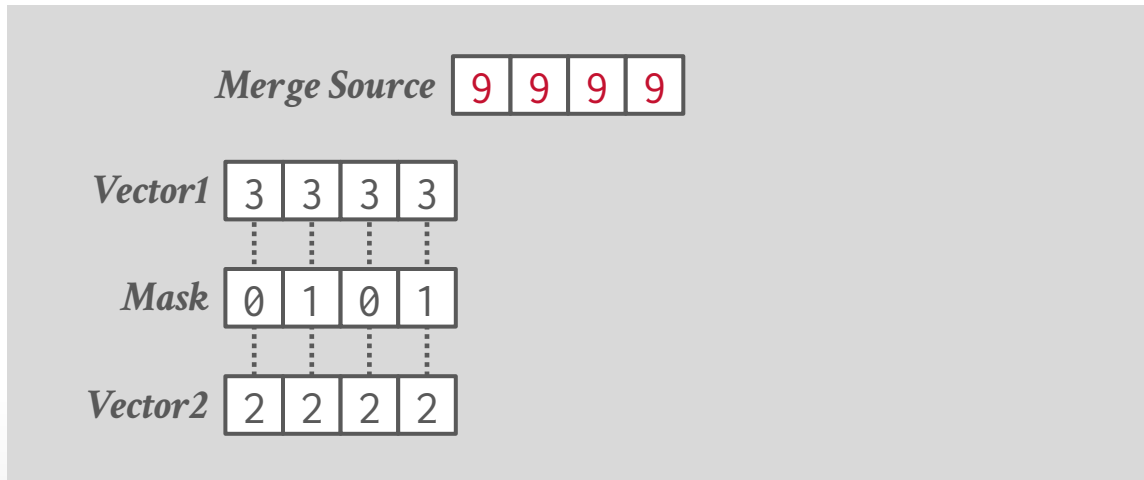
There are fundamental SIMD operations that the DBMS will use to build more complex functionality:

- Masking
- Permute
- Selective Load/Store
- Compress/Expand
- Selective Gather/Scatter



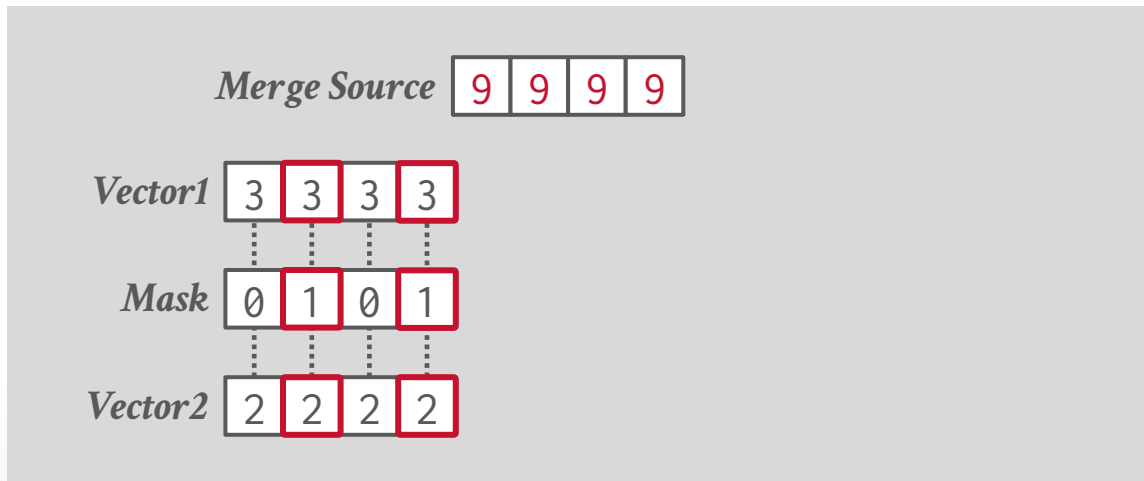
SIMD MASKING

Almost all AVX-512 operations support **predication** variants whereby the CPU only performs operations on lanes specified by an input bitmask.



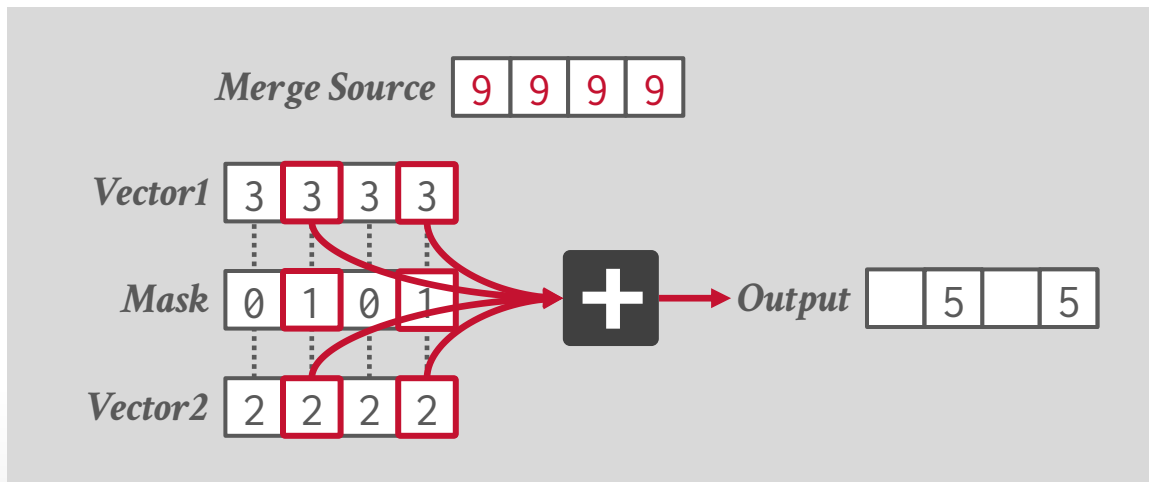
SIMD MASKING

Almost all AVX-512 operations support **predication** variants whereby the CPU only performs operations on lanes specified by an input bitmask.



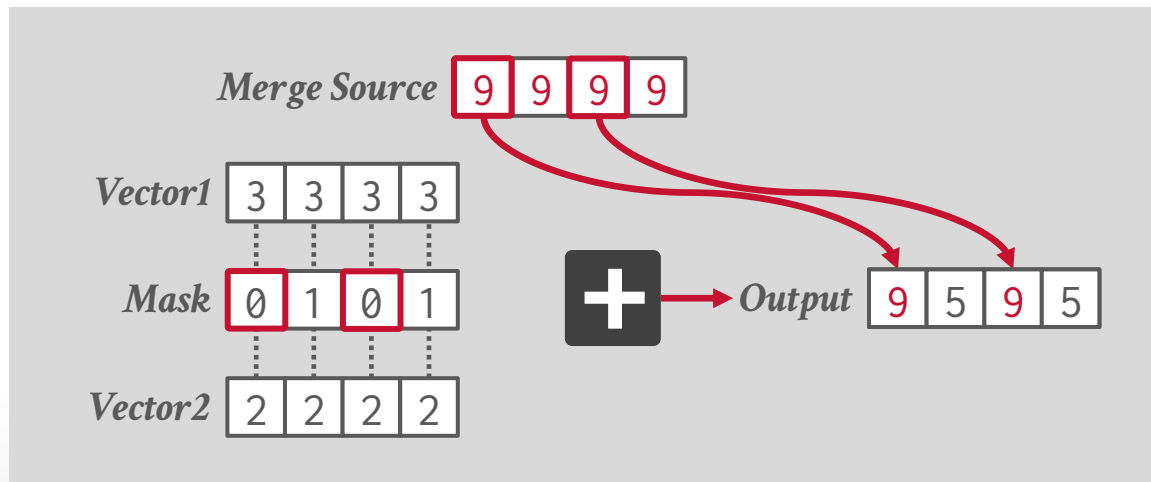
SIMD MASKING

Almost all AVX-512 operations support **predication** variants whereby the CPU only performs operations on lanes specified by an input bitmask.



SIMD MASKING

Almost all AVX-512 operations support **predication** variants whereby the CPU only performs operations on lanes specified by an input bitmask.

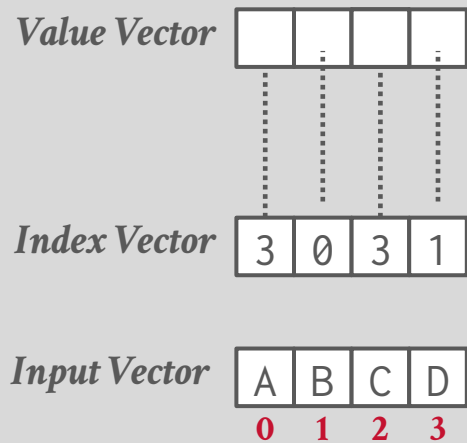


PERMUTE

For each lane, copy values in the **input vector** specified by the offset in the **index vector** into the **destination vector**.

Prior to AVX-512, the DBMS had to write data from the SIMD register to memory then back to the SIMD register.

Permute

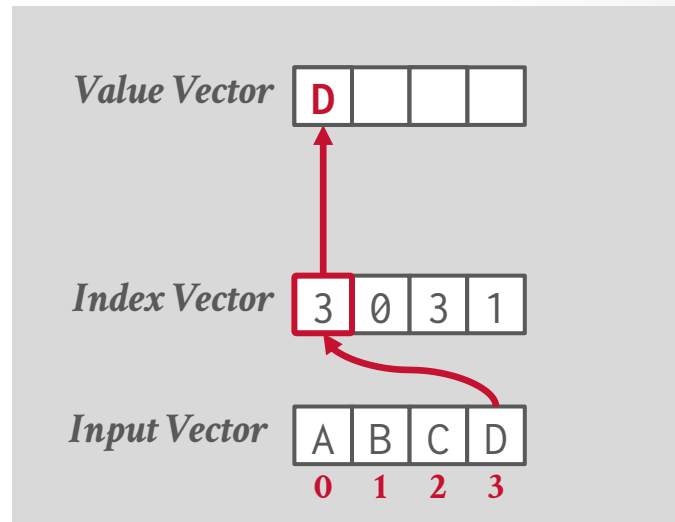


PERMUTE

For each lane, copy values in the **input vector** specified by the offset in the **index vector** into the **destination vector**.

Prior to AVX-512, the DBMS had to write data from the SIMD register to memory then back to the SIMD register.

Permute

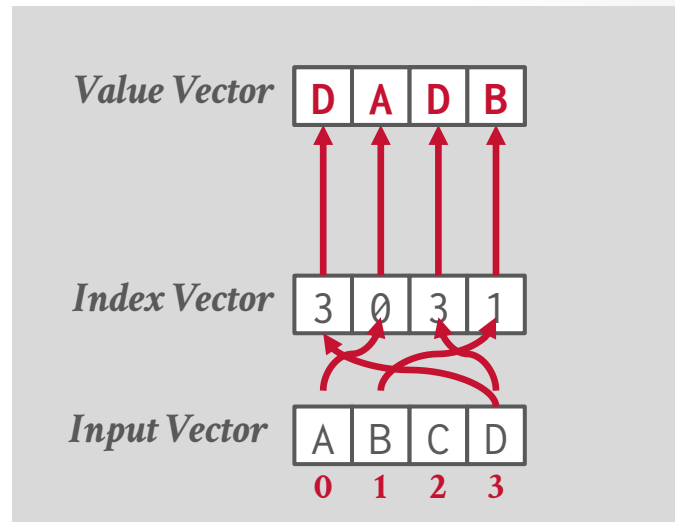


PERMUTE

For each lane, copy values in the **input vector** specified by the offset in the **index vector** into the **destination vector**.

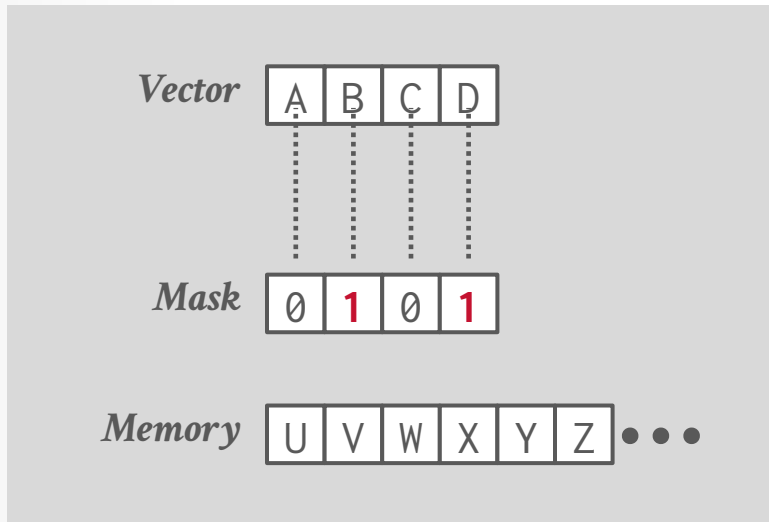
Prior to AVX-512, the DBMS had to write data from the SIMD register to memory then back to the SIMD register.

Permute



SELECTIVE LOAD/STORE

Selective Load



SELECTIVE LOAD/STORE

Selective Load

Vector

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

Mask

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

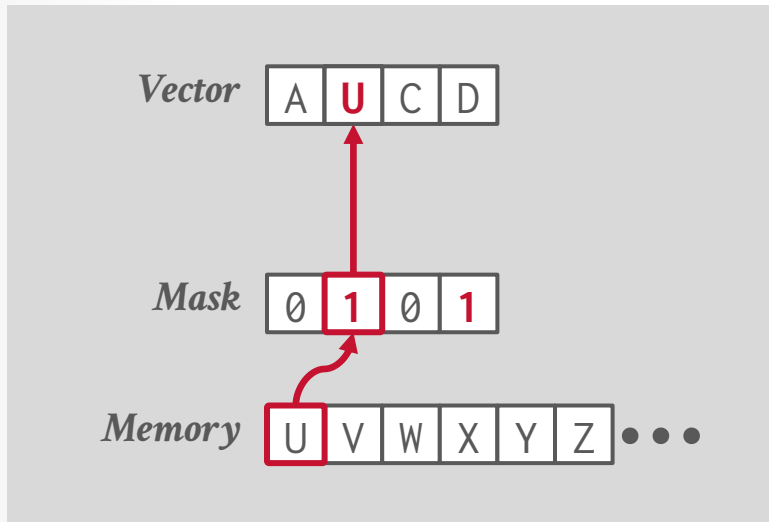
Memory

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

 ...

SELECTIVE LOAD/STORE

Selective Load



SELECTIVE LOAD/STORE

Selective Load

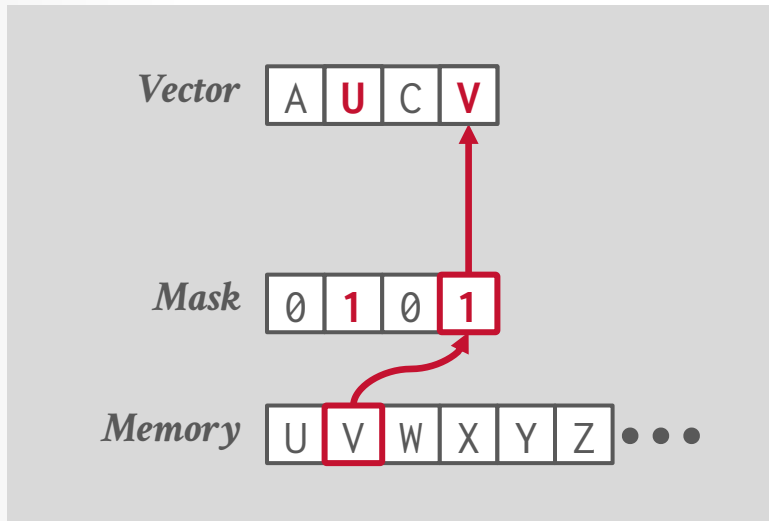
Vector A U C D

Mask 0 1 0 1

Memory U V W X Y Z ...

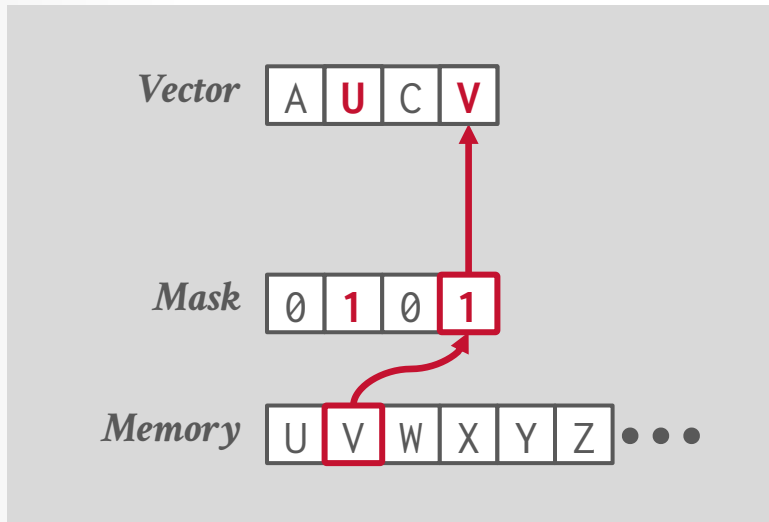
SELECTIVE LOAD/STORE

Selective Load

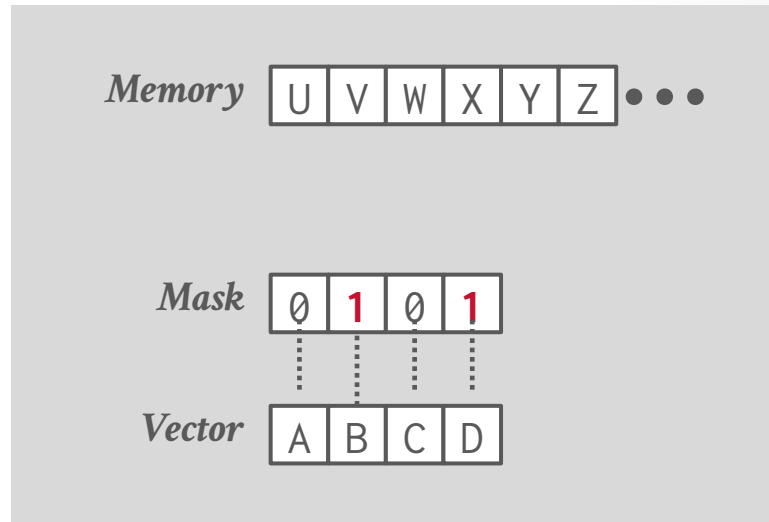


SELECTIVE LOAD/STORE

Selective Load

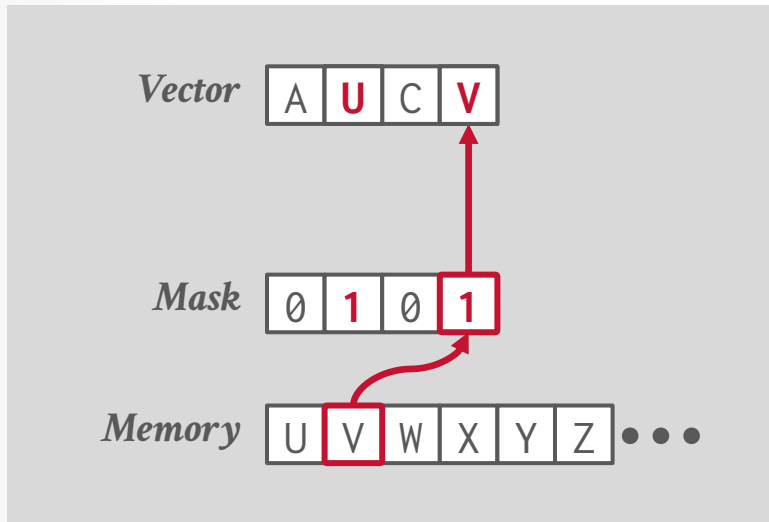


Selective Store

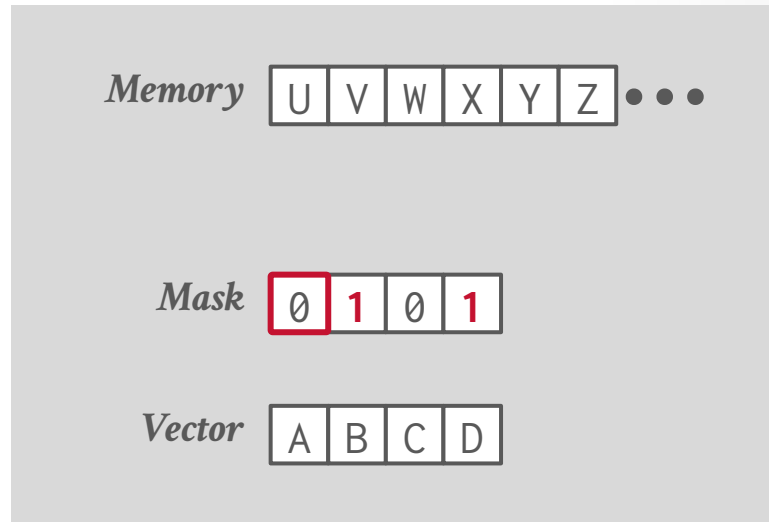


SELECTIVE LOAD/STORE

Selective Load

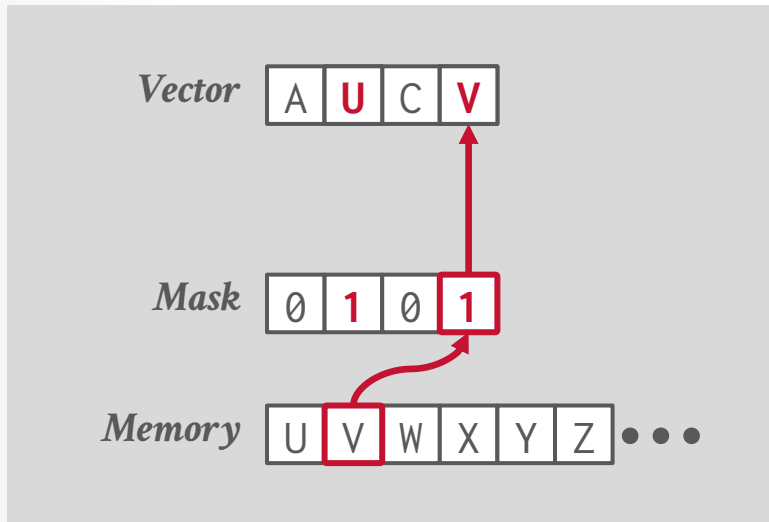


Selective Store

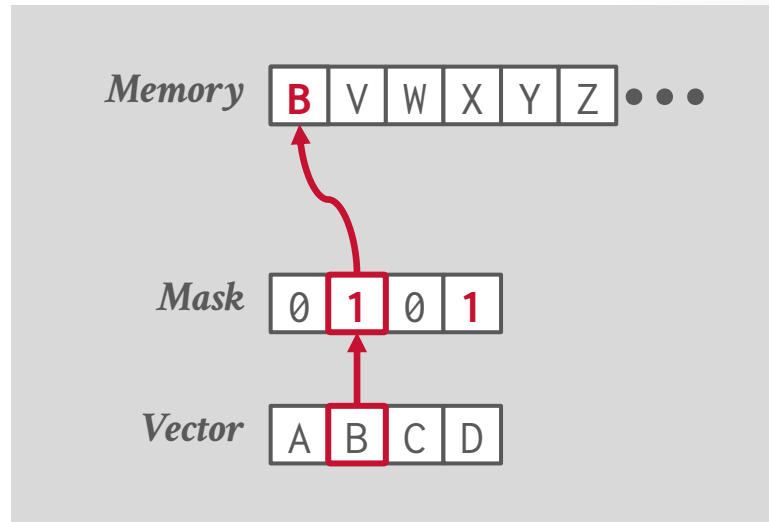


SELECTIVE LOAD/STORE

Selective Load

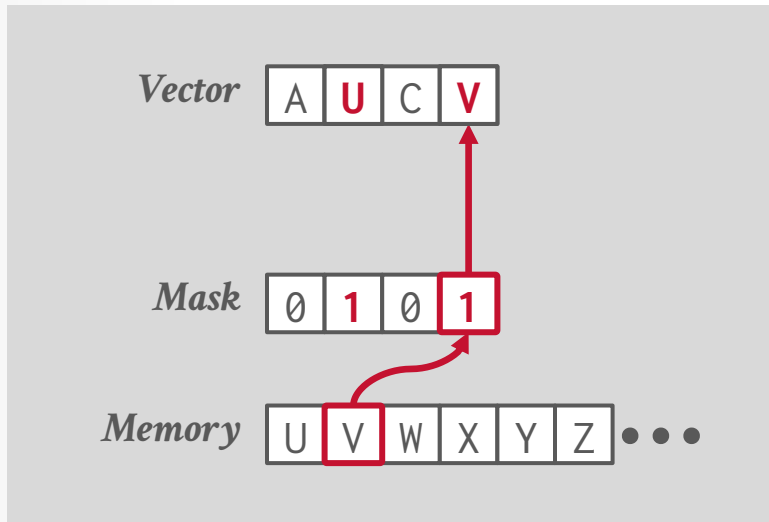


Selective Store

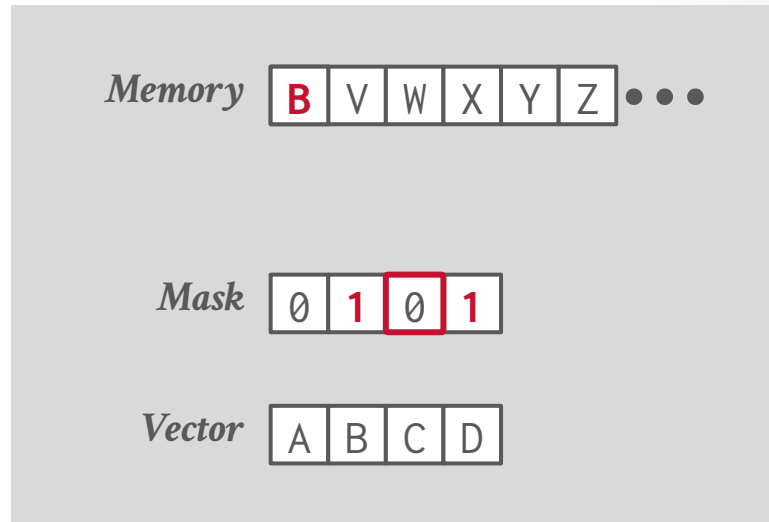


SELECTIVE LOAD/STORE

Selective Load

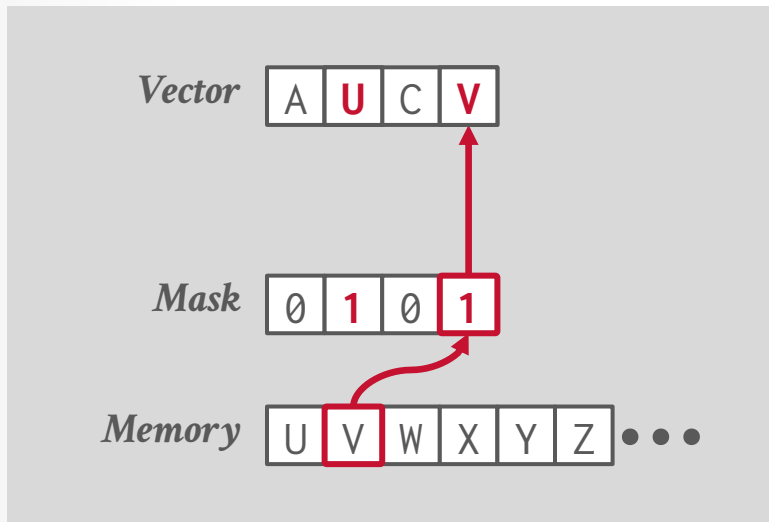


Selective Store

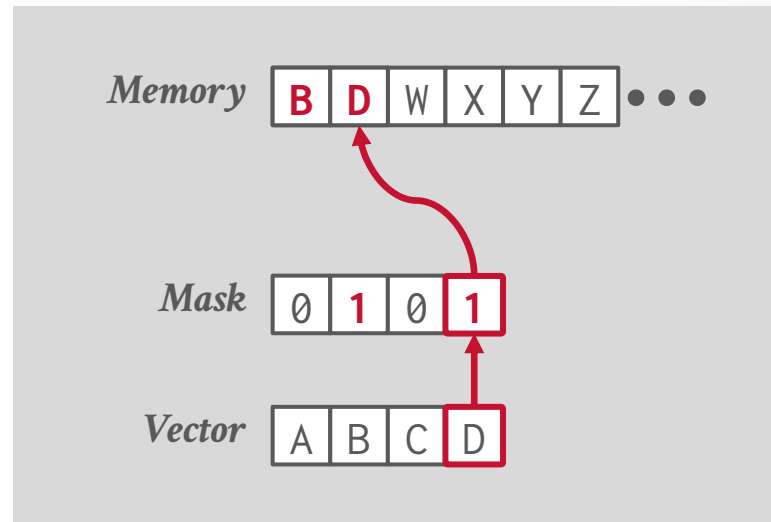


SELECTIVE LOAD/STORE

Selective Load



Selective Store



COMPRESS / EXPAND

Compress

Value Vector

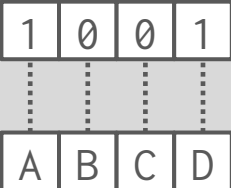
--	--	--	--

Index Vector

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

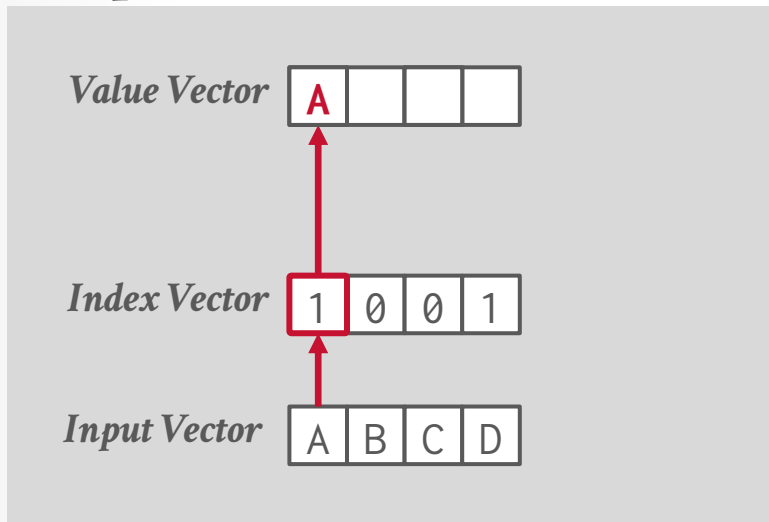
Input Vector

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



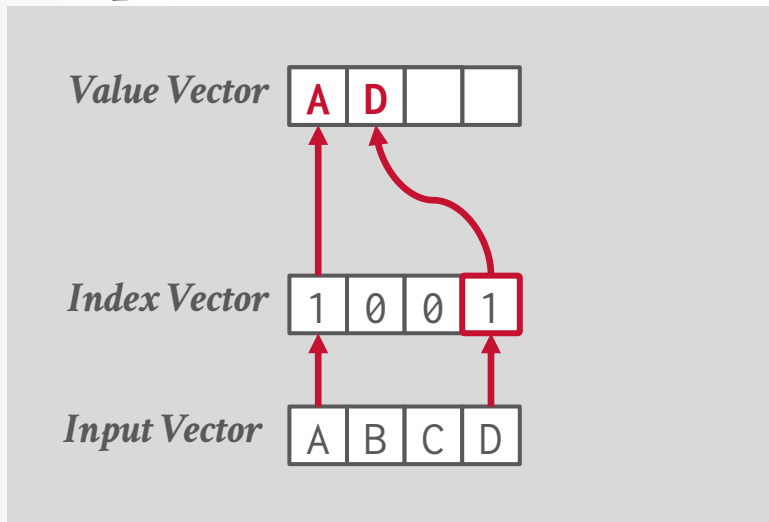
COMPRESS / EXPAND

Compress



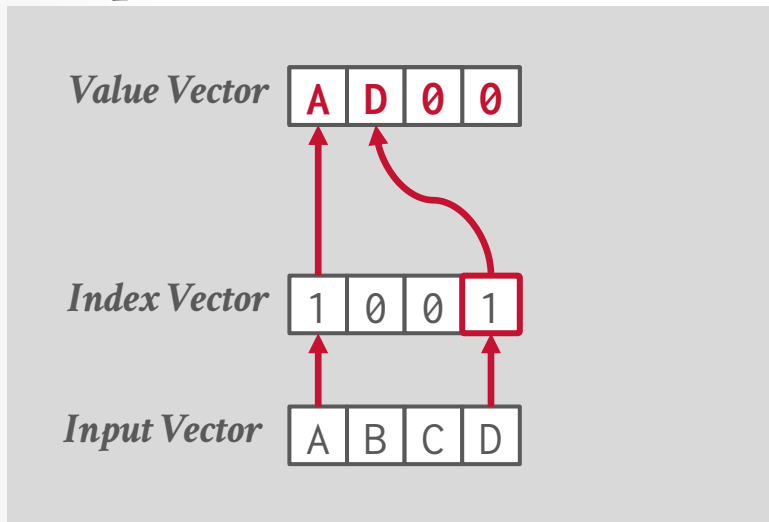
COMPRESS / EXPAND

Compress



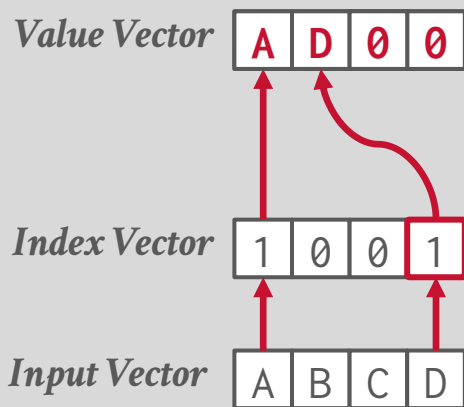
COMPRESS / EXPAND

Compress

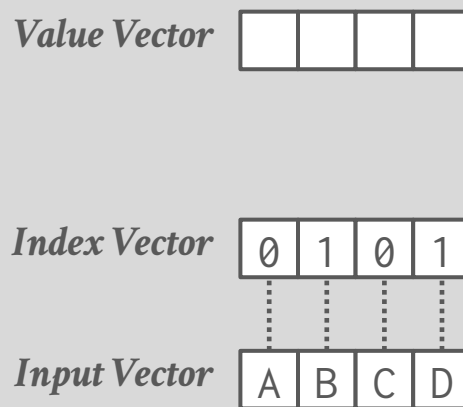


COMPRESS / EXPAND

Compress

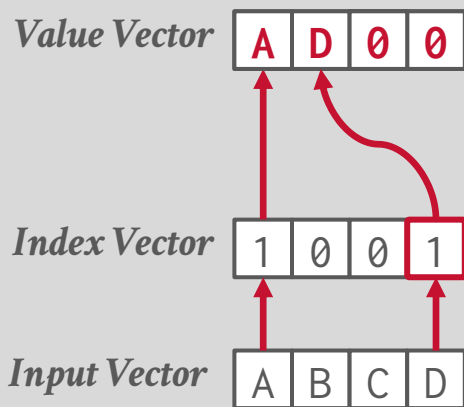


Expand

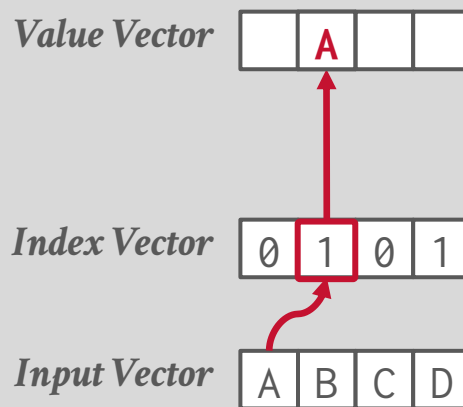


COMPRESS / EXPAND

Compress

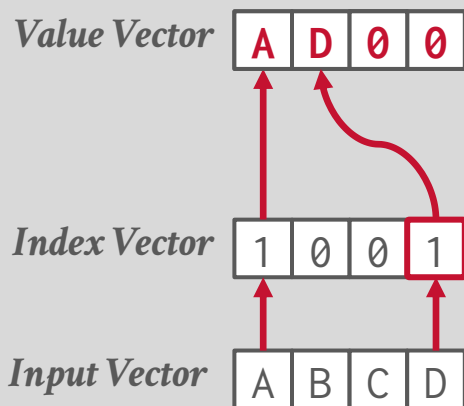


Expand

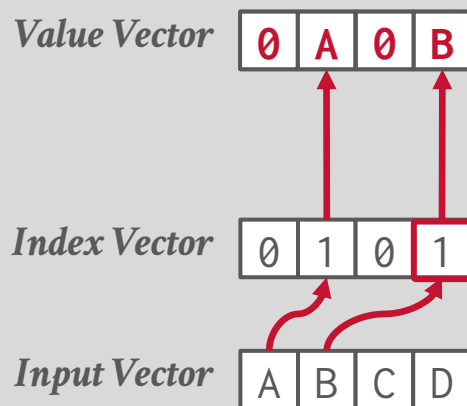


COMPRESS / EXPAND

Compress

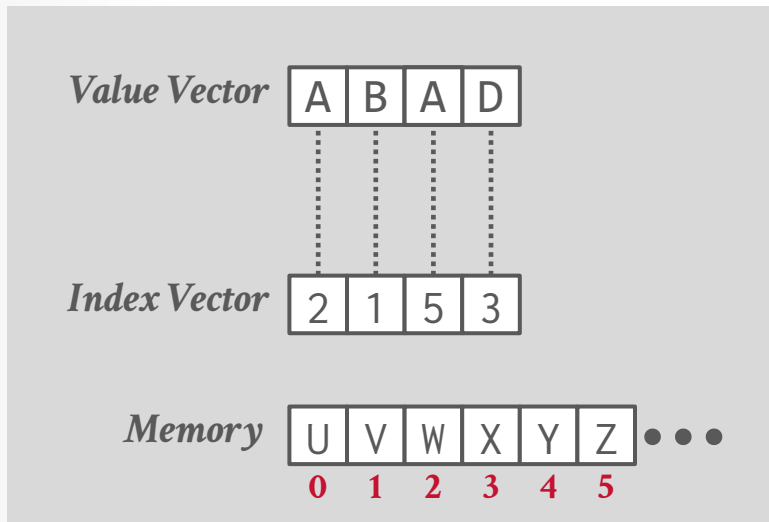


Expand



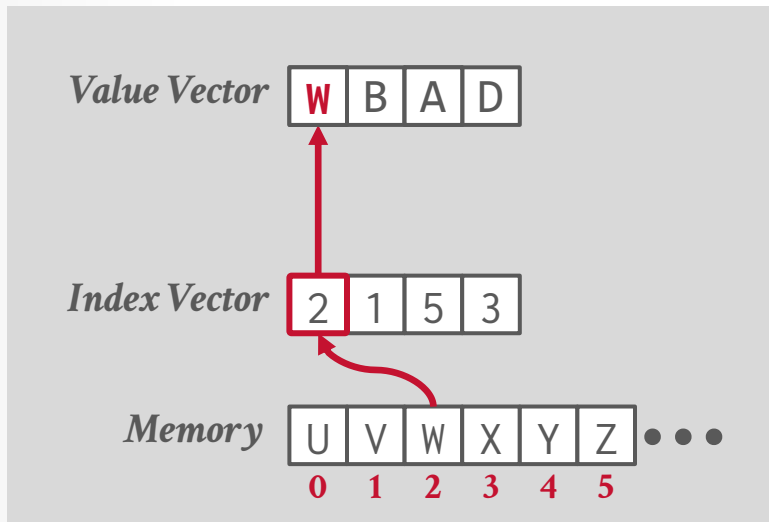
SELECTIVE SCATTER/GATHER

Selective Gather



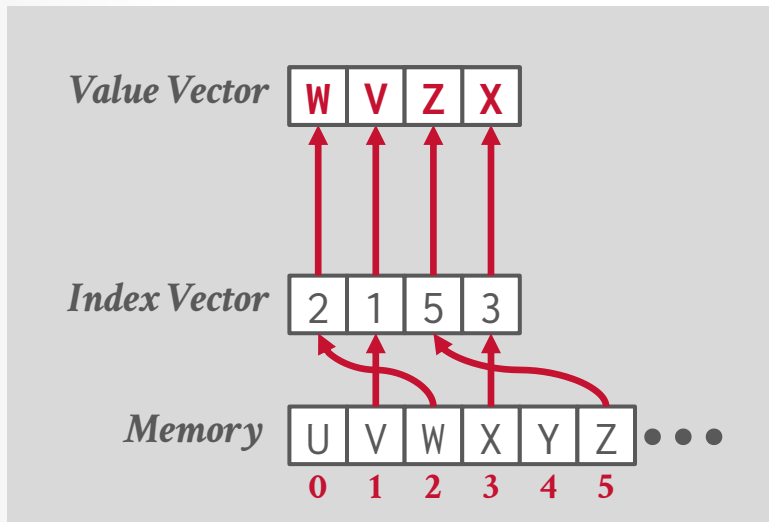
SELECTIVE SCATTER/GATHER

Selective Gather



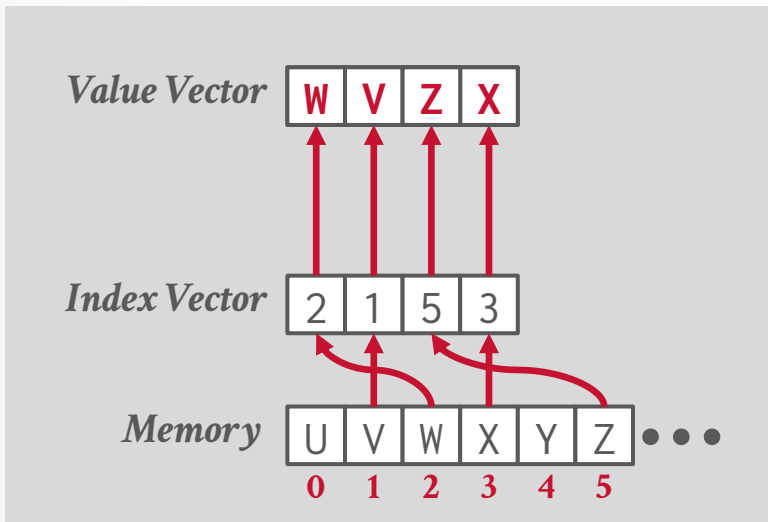
SELECTIVE SCATTER/GATHER

Selective Gather

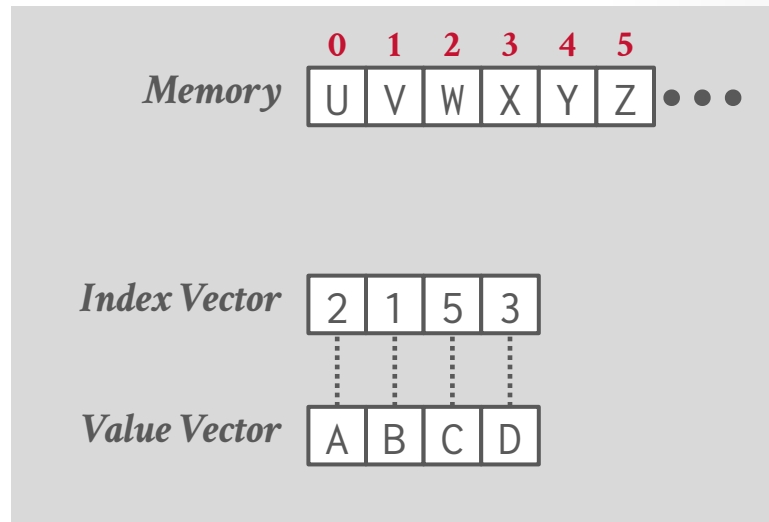


SELECTIVE SCATTER/GATHER

Selective Gather

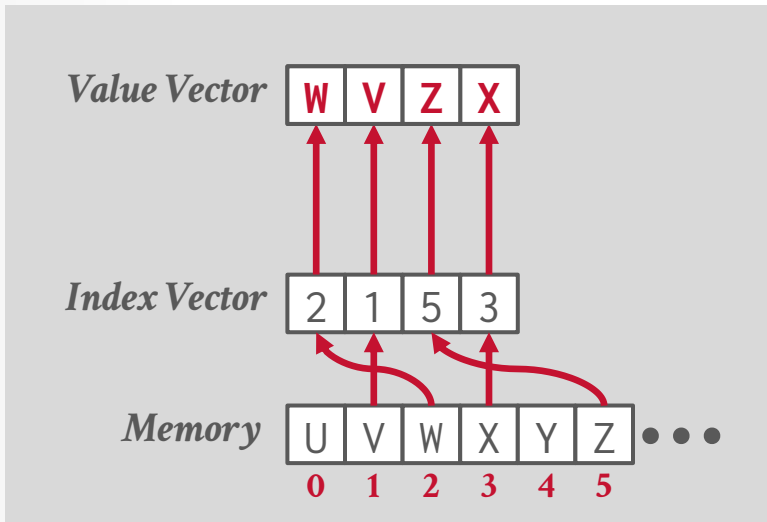


Selective Scatter

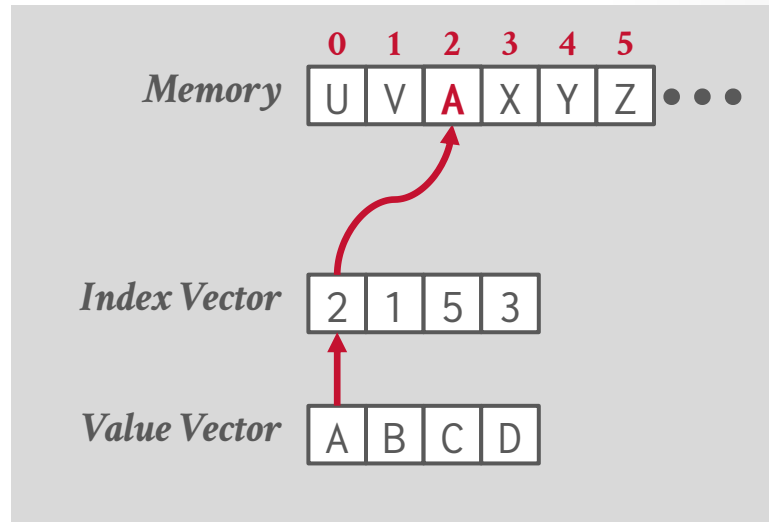


SELECTIVE SCATTER/GATHER

Selective Gather

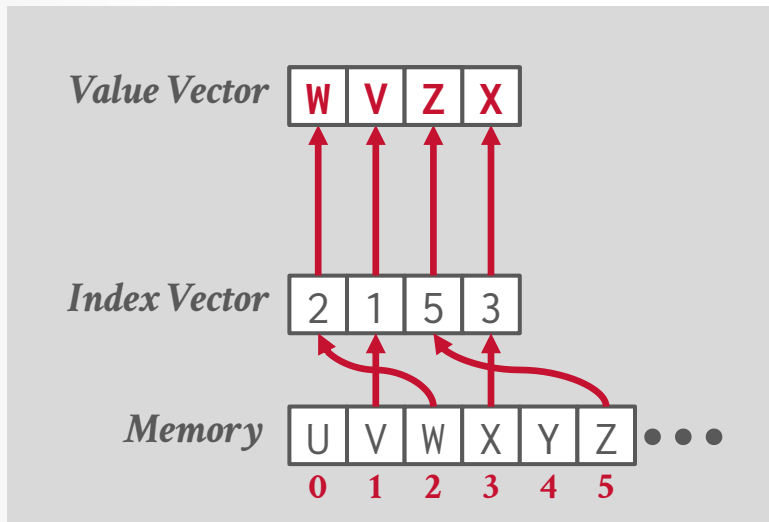


Selective Scatter

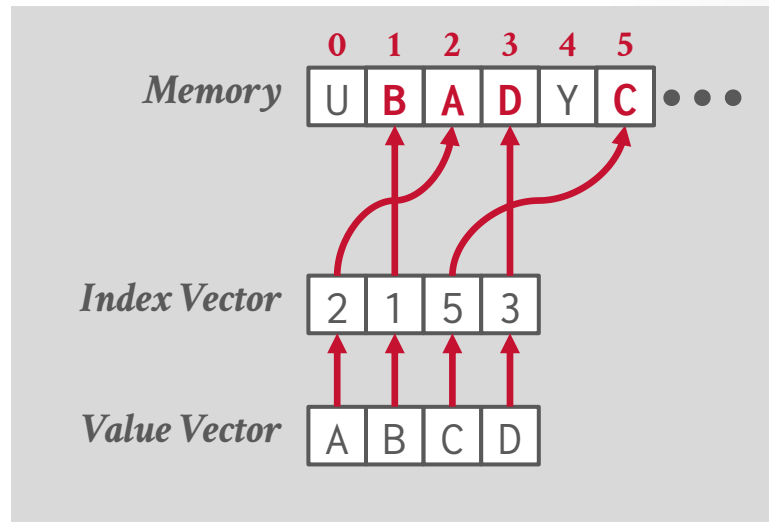


SELECTIVE SCATTER/GATHER

Selective Gather



Selective Scatter



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 unique data items per lane subset (i.e., no useless computations).

VECTORIZED OPERATORS

Selection Scans

Vector Refill

Hash Tables

Partitioning / Histograms



RETHINKING SIMD VECTORIZATION FOR
IN-MEMORY DATABASES
SIGMOD 2015

SELECTION SCANS

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
```

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

SELECTION SCANS

Vectorized

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

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

SELECTION SCANS

Vectorized

```

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

```

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

TID	KEY
100	A
101	N
102	D
103	Y
104	P
105	I
106	S
107	💩

Key Vector

A N D Y P I S 💩

SIMD Compare

Mask #1

0 1 0 1 1 0 1 0

SELECTION SCANS

Vectorized

```

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

```

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

TID	KEY
100	A
101	N
102	D
103	Y
104	P
105	I
106	S
107	💩

Key Vector

A N D Y P I S 💩

SIMD Compare

Mask #1

0 1 0 1 1 0 1 0

Mask #2

1 1 1 0 1 1 1 0

SELECTION SCANS

Vectorized

```

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

```

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

TID	KEY
100	A
101	N
102	D
103	Y
104	P
105	I
106	S
107	💩

Key Vector

A N D Y P I S 💩

SIMD Compare

Mask #1

0 1 0 1 1 0 1 0

Mask #2

1 1 1 0 1 1 1 0

SIMD AND

Mask #3

0 1 0 0 1 0 1 0

SELECTION SCANS

Vectorized

```

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

```

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

TID	KEY
100	A
101	N
102	D
103	Y
104	P
105	I
106	S
107	💩

Key Vector

A N D Y P I S 💩

SIMD Compare

Mask #1

0 1 0 1 1 0 1 0

Mask #2

1 1 1 0 1 1 1 0

SIMD AND

Mask #3

0 1 0 0 1 0 1 0

All Offsets

0 1 2 3 4 5 6 7

SELECTION SCANS

Vectorized

```

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

```

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

TID	KEY
100	A
101	N
102	D
103	Y
104	P
105	I
106	S
107	🤩

Key Vector

A N D Y P I S 🤩

SIMD Compare

Mask #1

0 1 0 1 1 0 1 0

Mask #2

1 1 1 0 1 1 1 0

SIMD AND

Mask #3

0 1 0 0 1 0 1 0

All Offsets

0 1 2 3 4 5 6 7

SIMD Compress

Matched Offsets

1 4 6

SELECTION SCANS

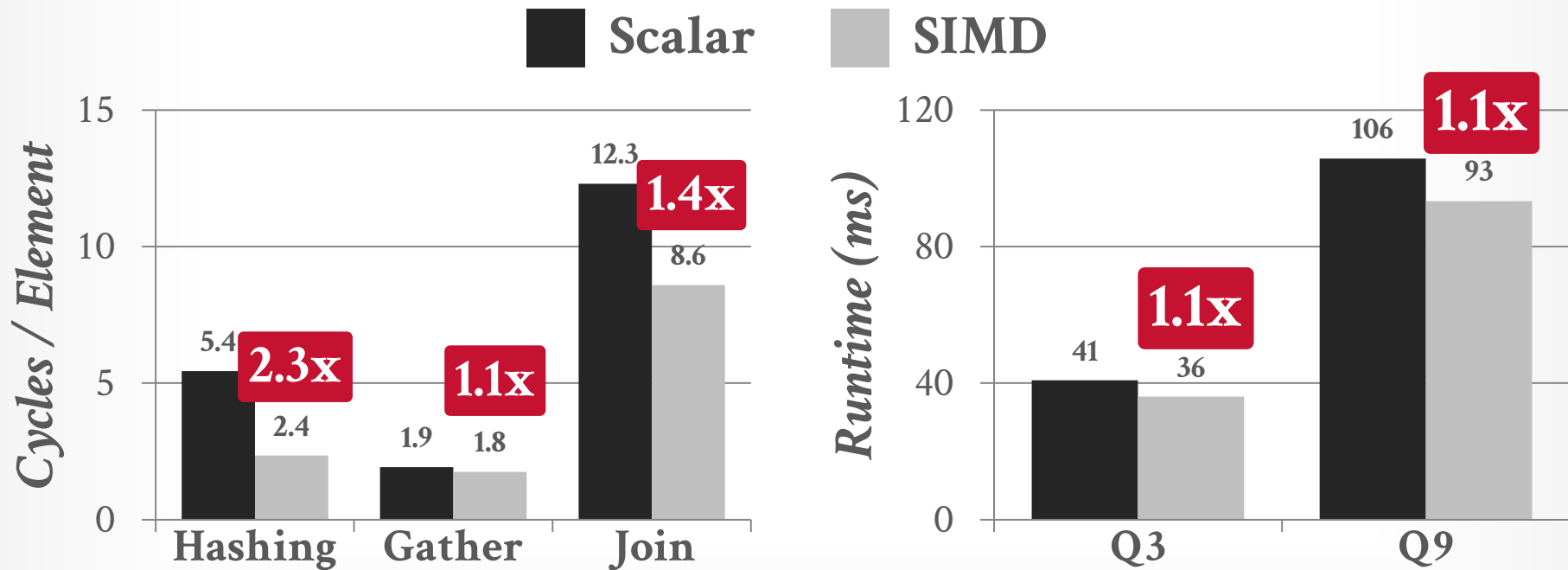
Evaluate branchless selection and hash probe in a [open-source implementation of Vectorwise](#).

Use AVX-512 because it includes instructions to make it easier to implement algorithms using vertical vectorization.

→ Selective operations using bitmask registers.

SIMD EVALUATION

Intel Core i9-7900X (10 cores × 2HT)
TPC-H Queries (Scalefactor=1)



Source: [Timo Kersten](#)

OBSERVATION

For each batch, the SIMD vectors may contain tuples that are no longer valid (they were disqualified by some previous check).

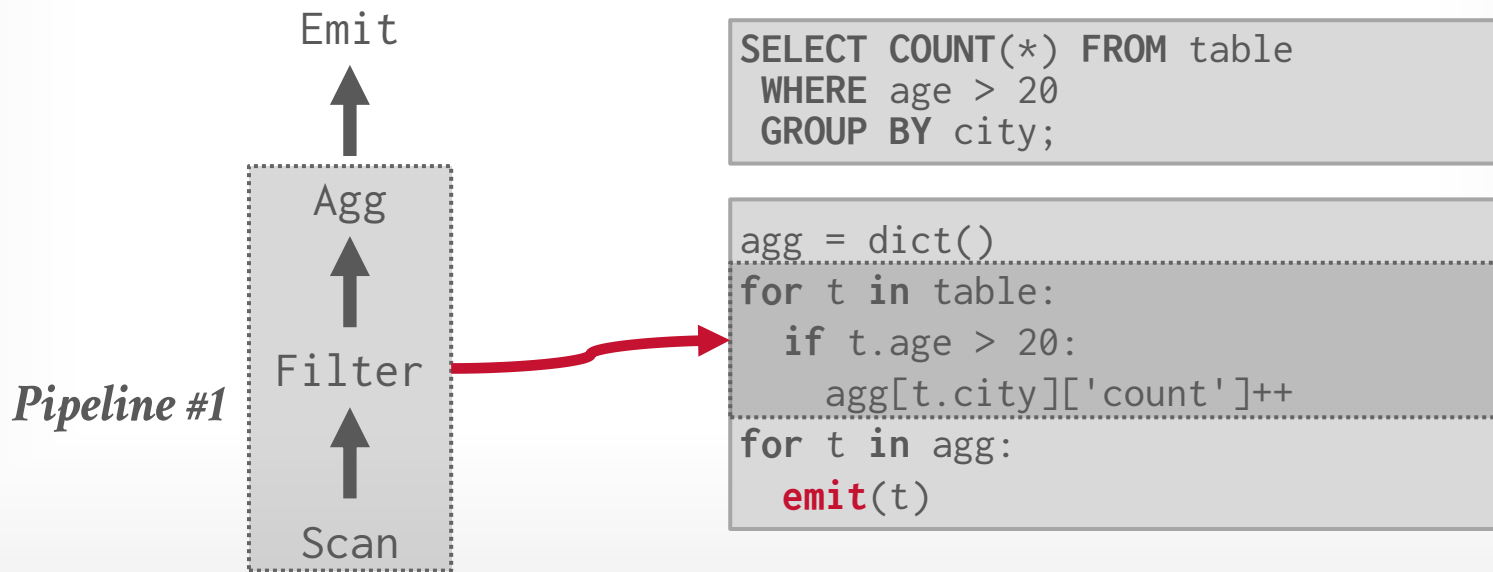


```
SELECT COUNT(*) FROM table
WHERE age > 20
GROUP BY city;
```

```
agg = dict()
for t in table:
    if t.age > 20:
        agg[t.city]['count']++
for t in agg:
    emit(t)
```

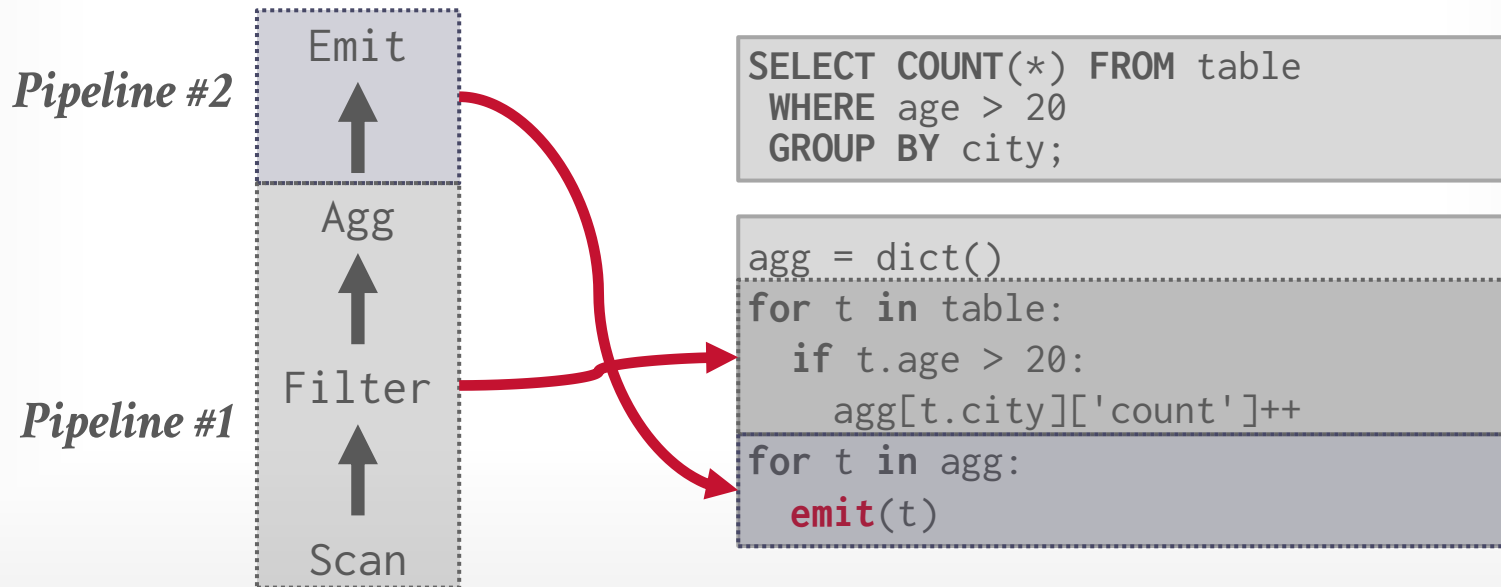
OBSERVATION

For each batch, the SIMD vectors may contain tuples that are no longer valid (they were disqualified by some previous check).



OBSERVATION

For each batch, the SIMD vectors may contain tuples that are no longer valid (they were disqualified by some previous check).



RELAXED OPERATOR FUSION

Vectorized processing model designed for query compilation execution engines.

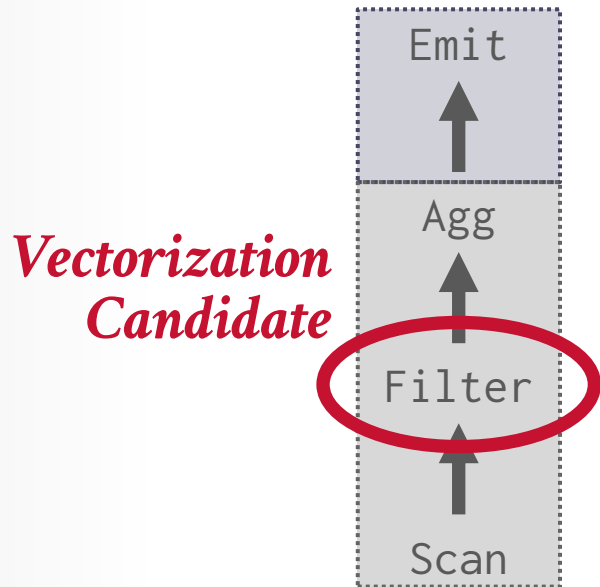
Decompose pipelines into stages that operate on vectors of tuples.

- Each stage may contain multiple operators.
- Communicate through cache-resident buffers.
- Stages are granularity of vectorization + fusion.



ROF EXAMPLE

```
SELECT COUNT(*) FROM table
WHERE age > 20 GROUP BY city;
```



ROF EXAMPLE

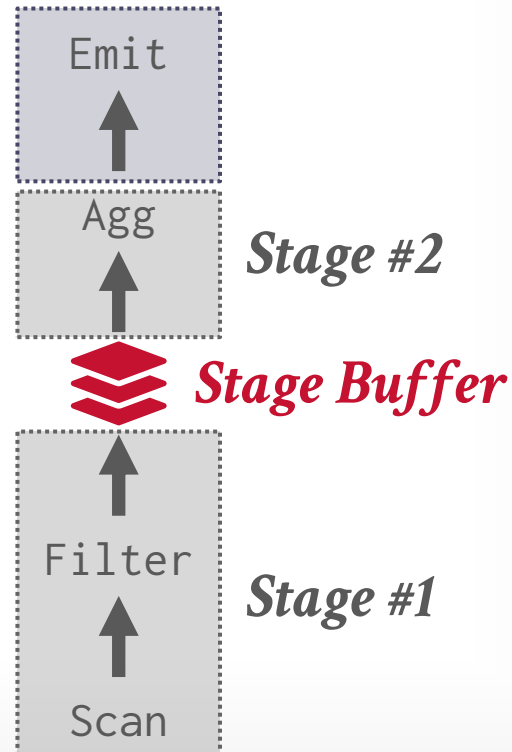
```
SELECT COUNT(*) FROM table
WHERE age > 20 GROUP BY city;
```



ROF EXAMPLE

```
SELECT COUNT(*) FROM table
WHERE age > 20 GROUP BY city;
```

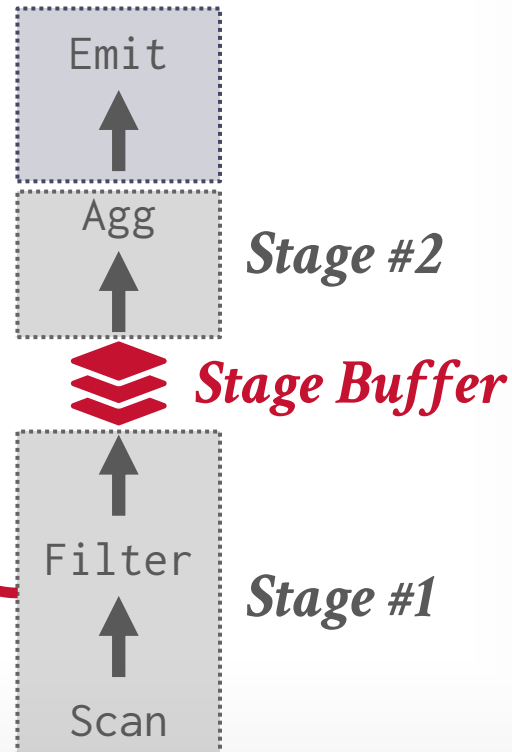
```
agg = dict()
for vt in table step 1024:
    buffer = simd_cmp_gt(vt, 20, 1024)
    if |buffer| >= MAX:
        for t in buffer:
            agg[t.city]['count']++
for t in agg:
    emit(t)
```



ROF EXAMPLE

```
SELECT COUNT(*) FROM table
WHERE age > 20 GROUP BY city;
```

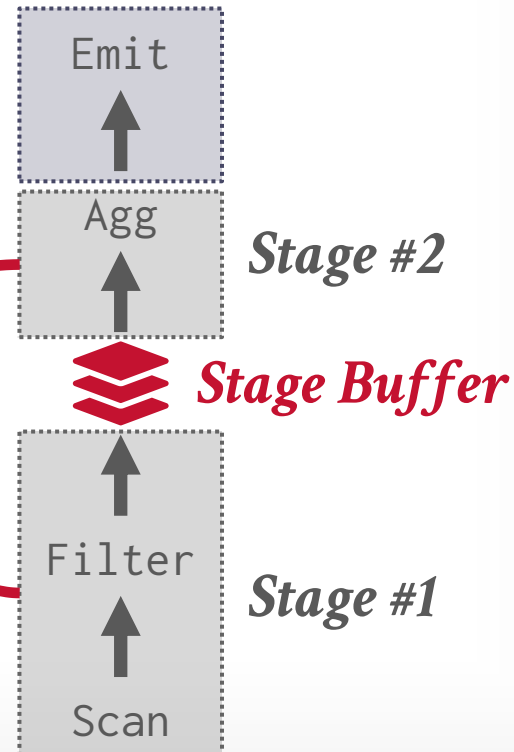
```
agg = dict()
for vt in table step 1024:
    buffer = simd_cmp_gt(vt, 20, 1024)
    if |buffer| >= MAX:
        for t in buffer:
            agg[t.city]['count']++
for t in agg:
    emit(t)
```



ROF EXAMPLE

```
SELECT COUNT(*) FROM table
WHERE age > 20 GROUP BY city;
```

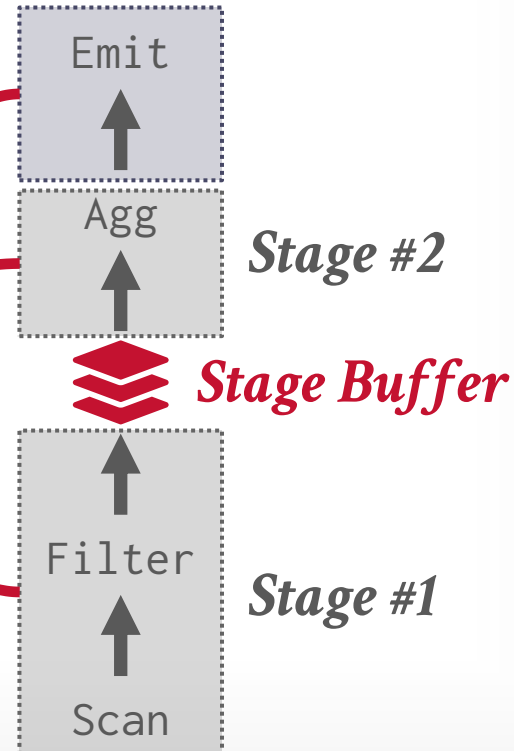
```
agg = dict()
for vt in table step 1024:
    buffer = simd_cmp_gt(vt, 20, 1024)
    if |buffer| >= MAX:
        for t in buffer:
            agg[t.city]['count']++
for t in agg:
    emit(t)
```



ROF EXAMPLE

```
SELECT COUNT(*) FROM table
WHERE age > 20 GROUP BY city;
```

```
agg = dict()
for vt in table step 1024:
    buffer = simd_cmp_gt(vt, 20, 1024)
    if |buffer| >= MAX:
        for t in buffer:
            agg[t.city]['count']++
for t in agg:
    emit(t)
```



ROF SOFTWARE PREFETCHING

The DBMS can tell the CPU to grab the next vector while it works on the current batch.

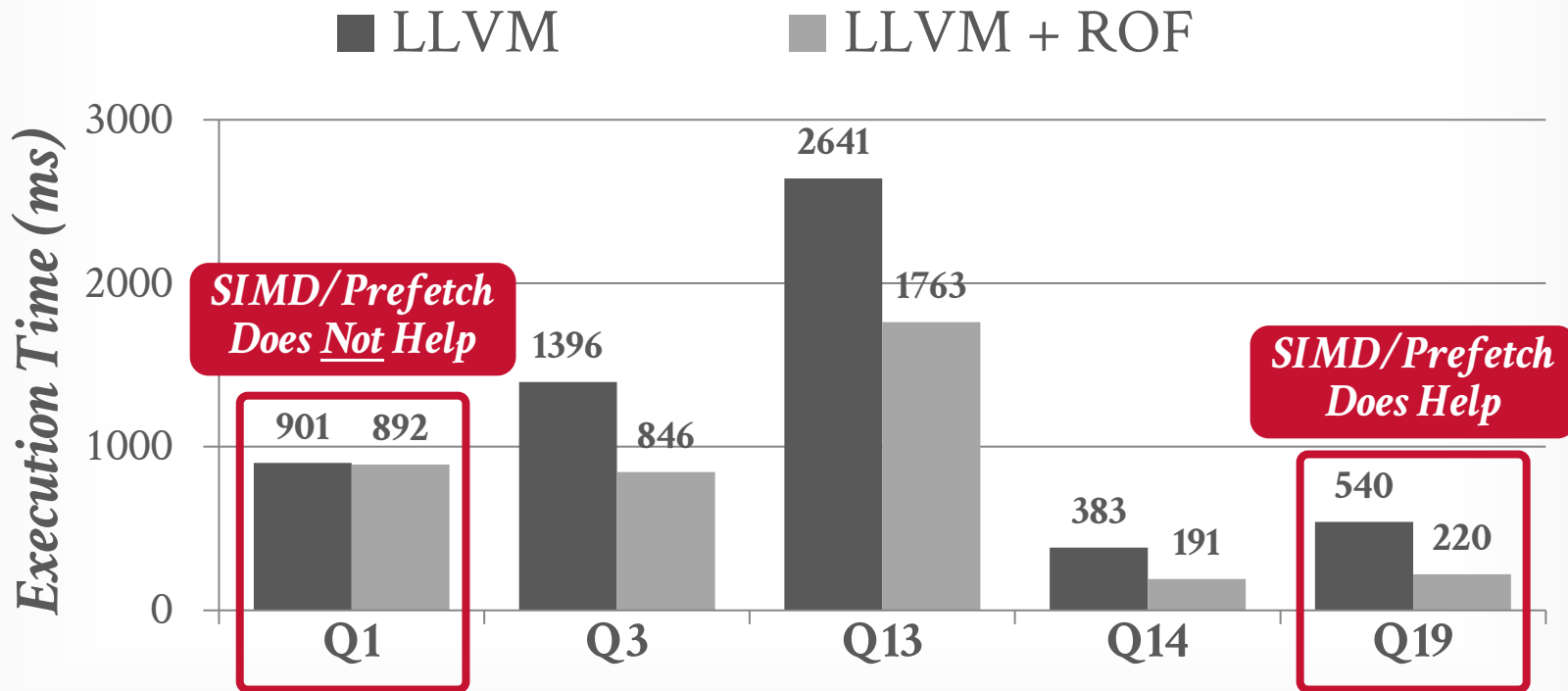
- Prefetch-enabled operators define start of new stage.
- Hides the cache miss latency.

Any prefetching technique is suitable

- Group prefetching, software pipelining, AMAC.
- Group prefetching works and is simple to implement.

ROF EVALUATION

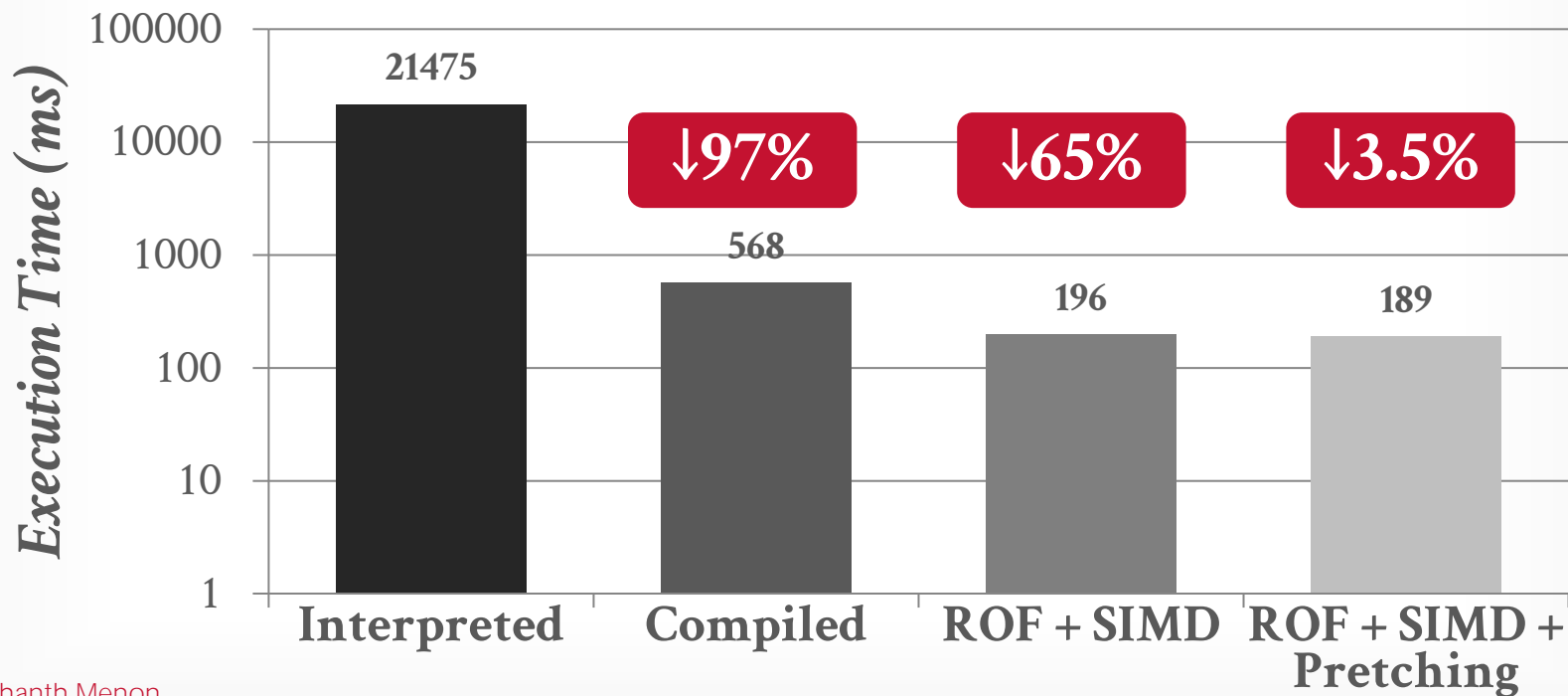
Dual Socket Intel Xeon E5-2630v4 @ 2.20GHz
TPC-H 10 GB Database



Source: [Prashanth Menon](#)

ROF EVALUATION - TPC-H Q19

Dual Socket Intel Xeon E5-2630v4 @ 2.20GHz
TPC-H 10 GB Database



Source: [Prashanth Menon](#)

VECTOR REFILL ALGORITHMS

Approach #1: Buffered

→ Use additional SIMD registers to stage results within an operator and proceed with next loop iteration to fill in underutilized lanes vectors.

Approach #2: Partial

- Use additional SIMD registers to buffer results from underutilized vectors and then return to previous operator to process the next vector.
- Requires fine-grained bookkeeping to make sure other operators do not clobber deferred vectors.

VECTORIZED OPERATORS

~~Selection Scans~~

~~Vector Refill~~

Hash Tables

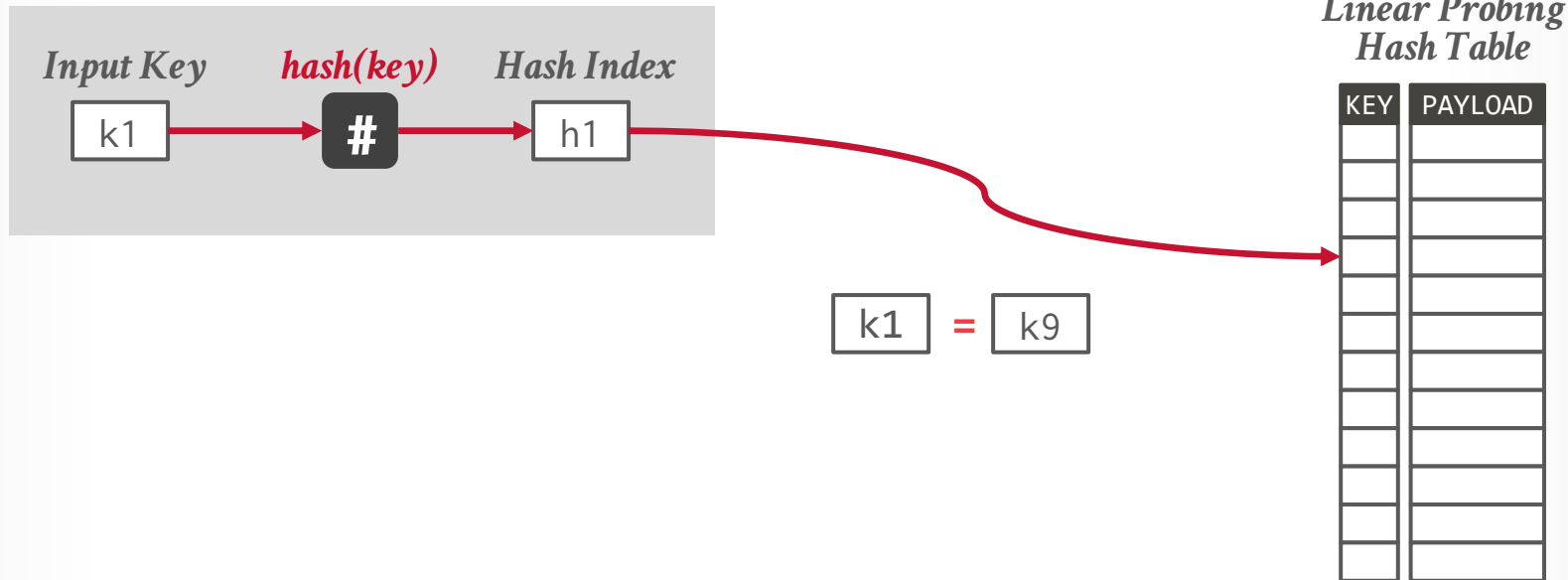
Partitioning / Histograms



RETHINKING SIMD VECTORIZATION FOR
IN-MEMORY DATABASES
SIGMOD 2015

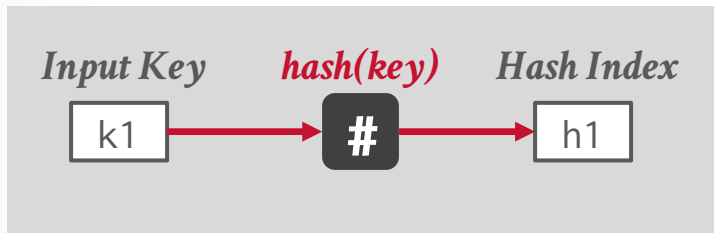
HASH TABLES - PROBING

Scalar

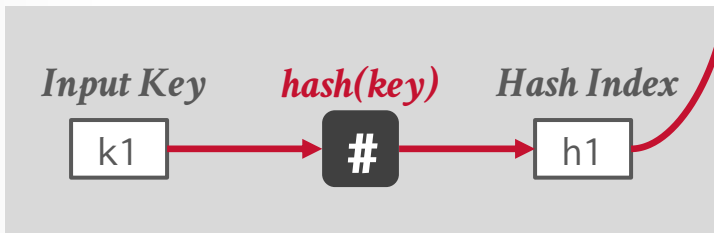


HASH TABLES - PROBING

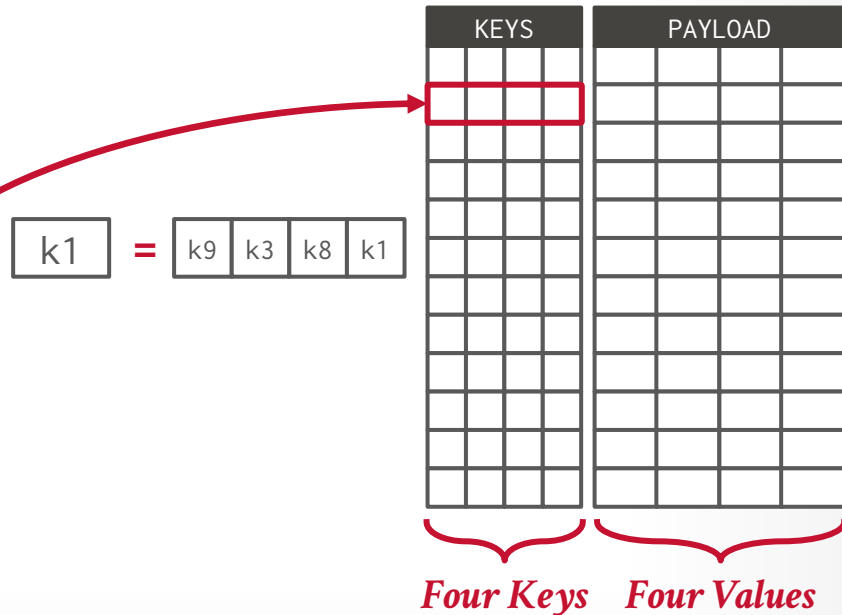
Scalar



Vectorized (Horizontal)

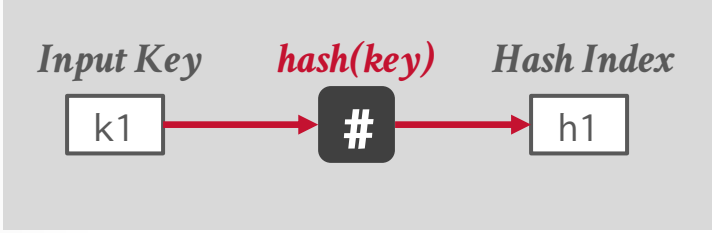


Linear Probing Bucketized Hash Table

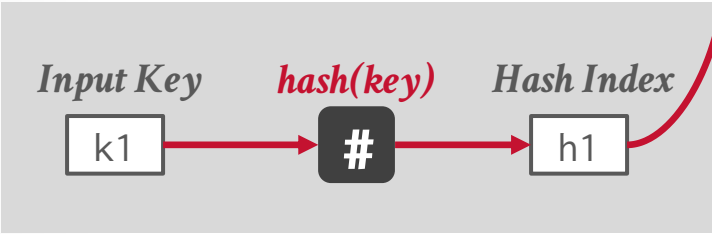


HASH TABLES - PROBING

Scalar

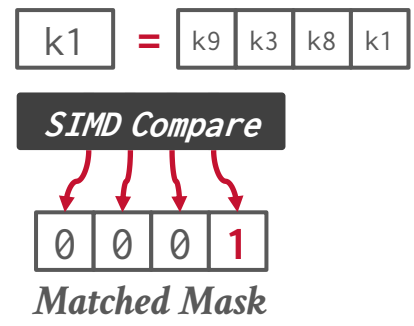


Vectorized (Horizontal)



Linear Probing Bucketized Hash Table

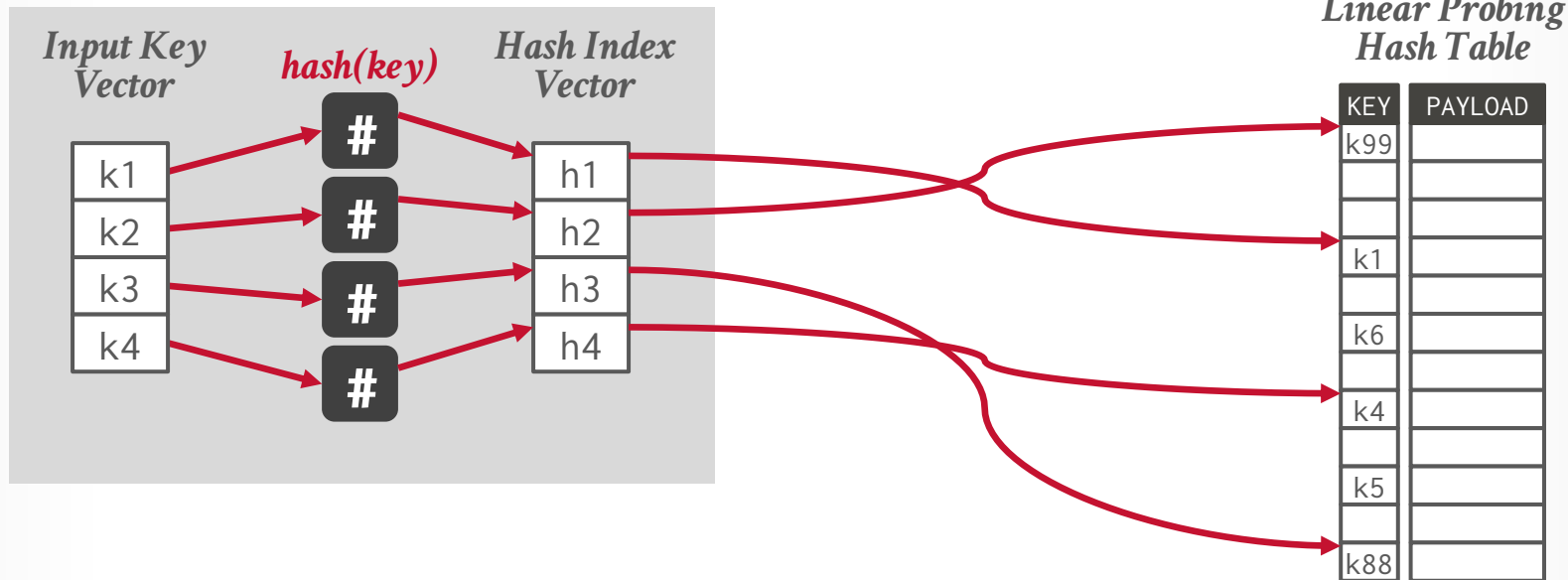
KEYS				PAYLOAD			



Four Keys Four Values

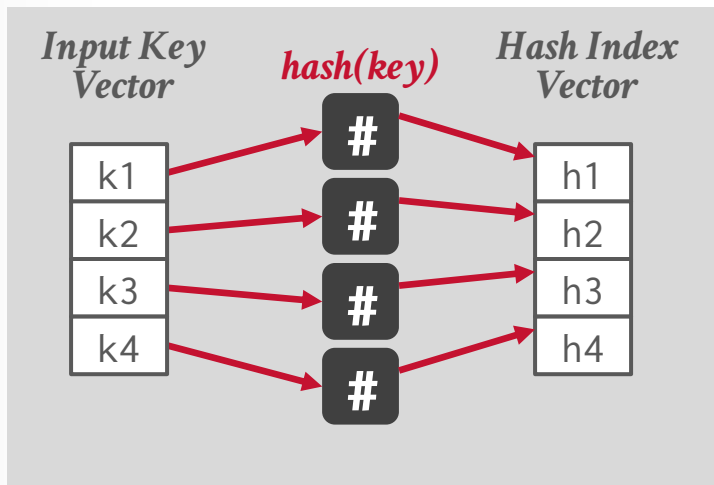
HASH TABLES - PROBING

Vectorized (Vertical)

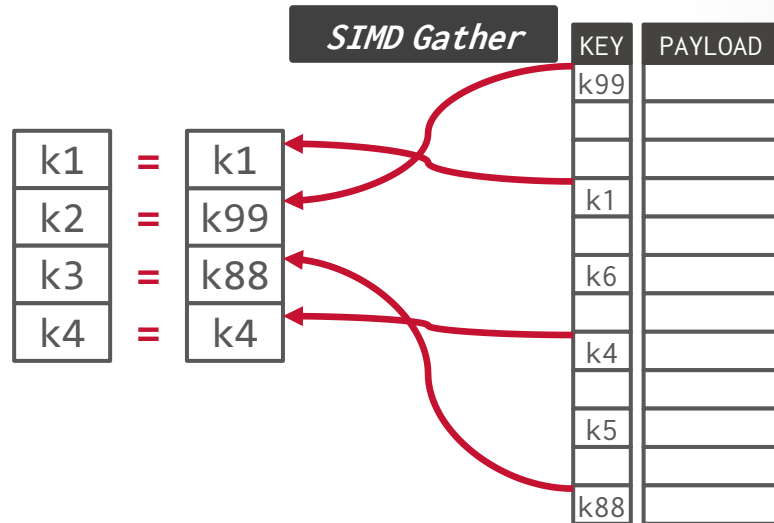


HASH TABLES – PROBING

Vectorized (Vertical)

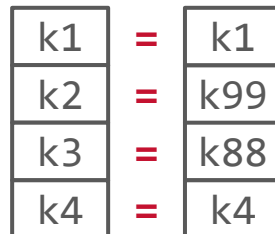
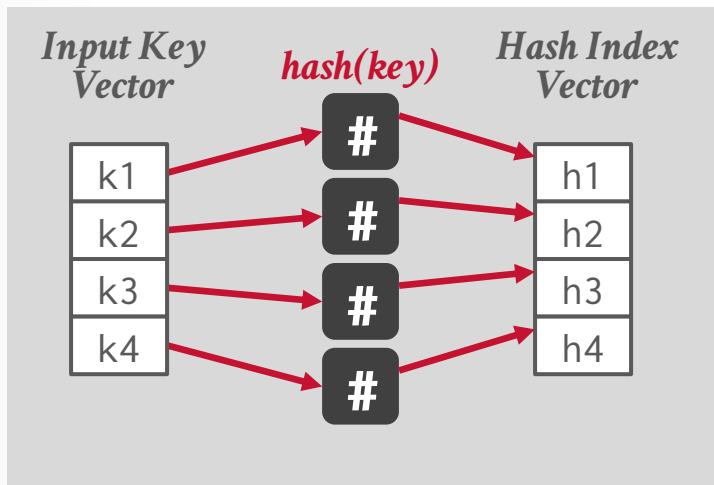


Linear Probing Hash Table



HASH TABLES – PROBING

Vectorized (Vertical)



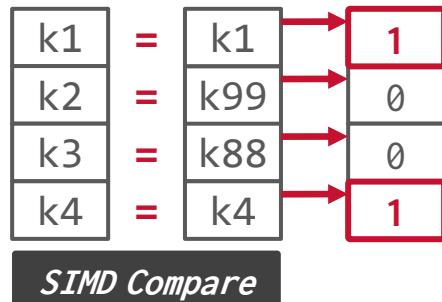
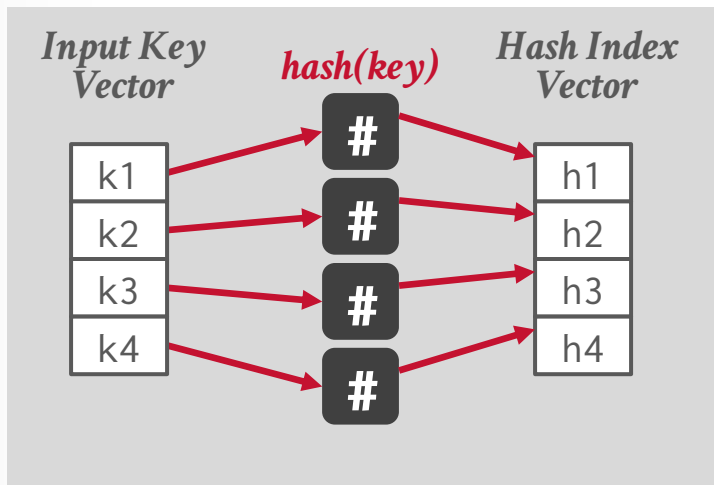
SIMD Compare

Linear Probing Hash Table

KEY	PAYLOAD
k99	
k1	
k6	
k4	
k5	
k88	

HASH TABLES – PROBING

Vectorized (Vertical)

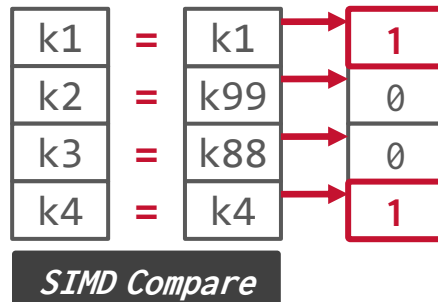
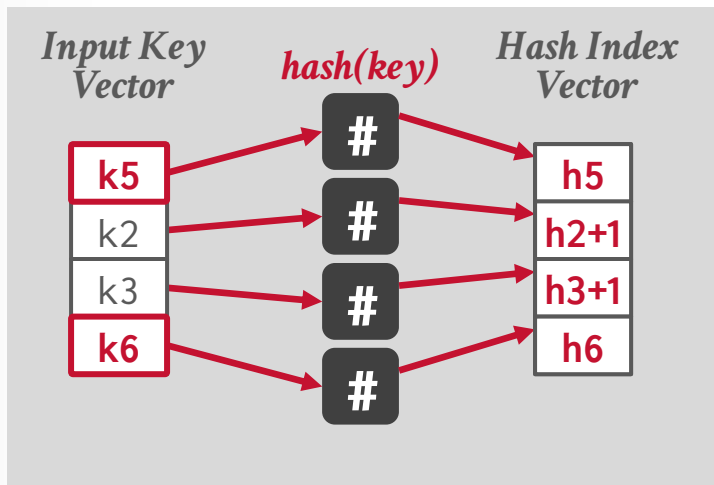


Linear Probing Hash Table

KEY	PAYLOAD
k99	
k1	
k6	
k4	
k5	
k88	

HASH TABLES – PROBING

Vectorized (Vertical)

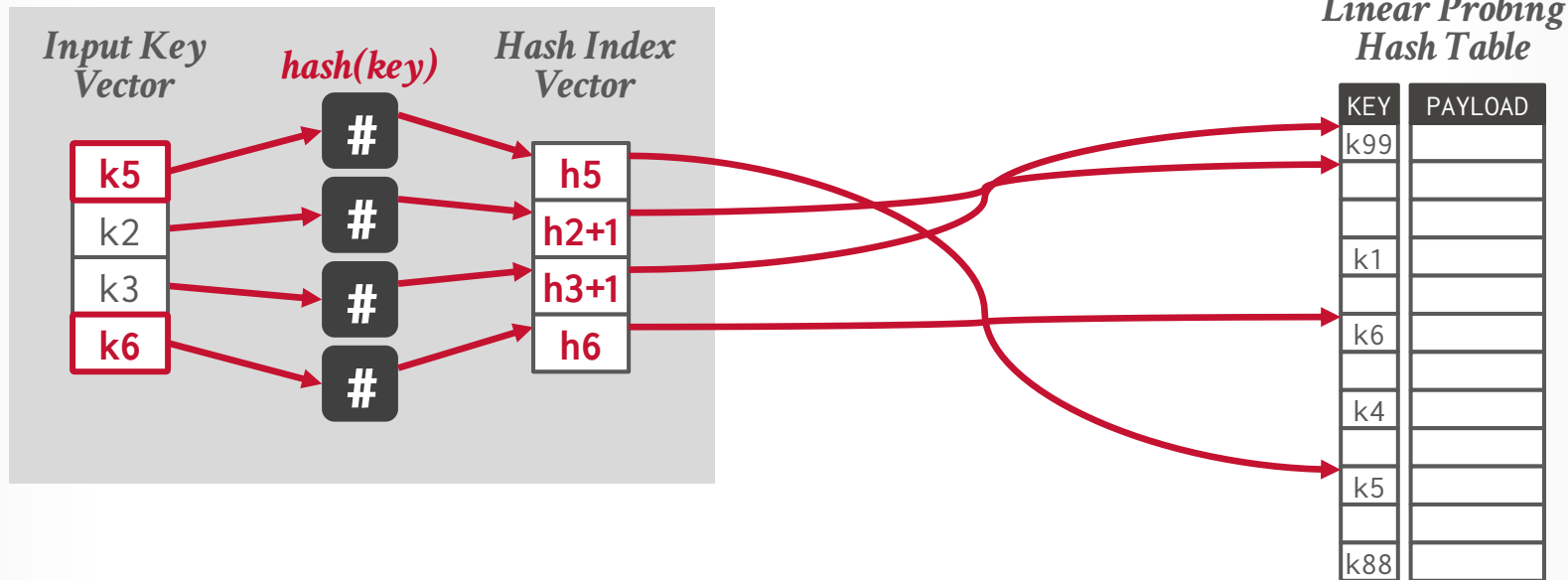


Linear Probing Hash Table

KEY	PAYLOAD
k99	
k1	
k6	
k4	
k5	
k88	

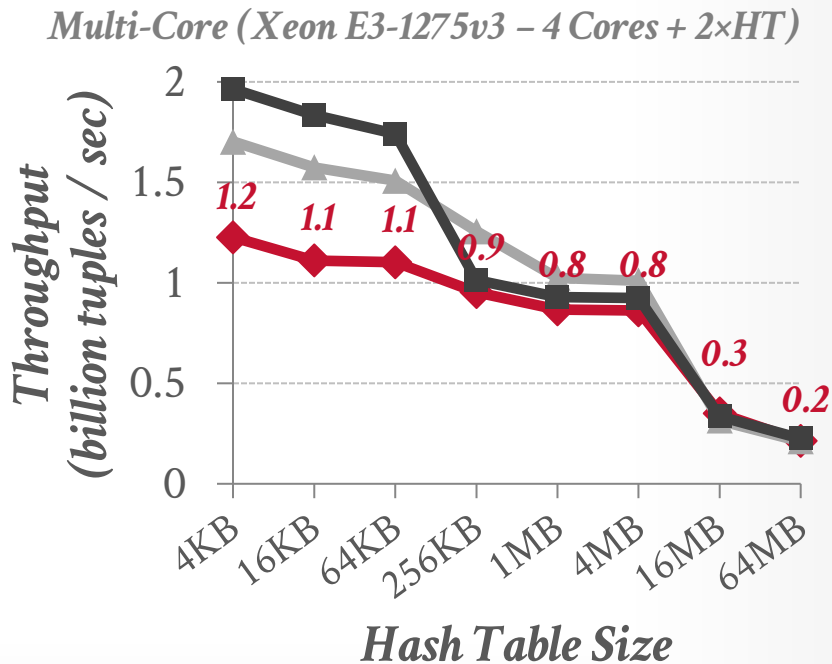
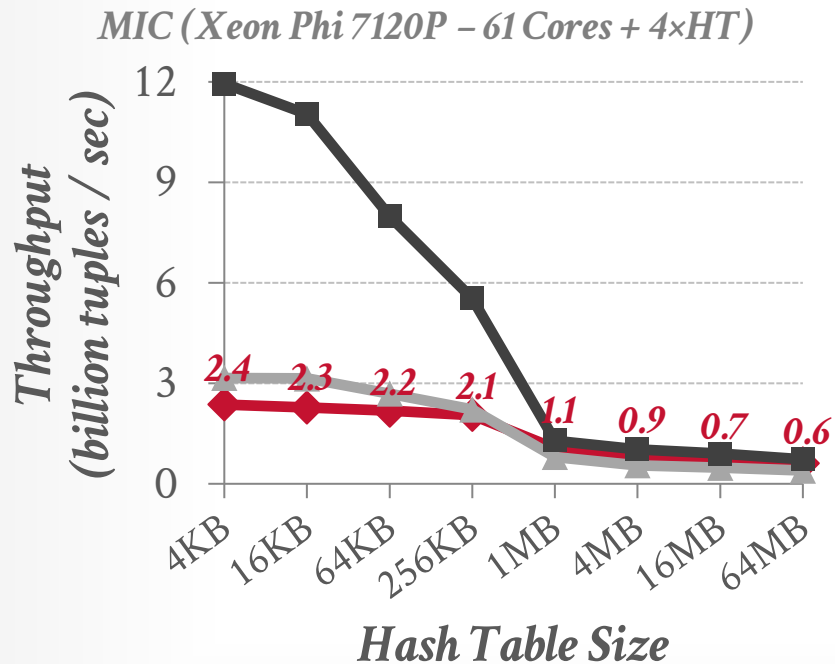
HASH TABLES - PROBING

Vectorized (Vertical)



HASH TABLES - PROBING

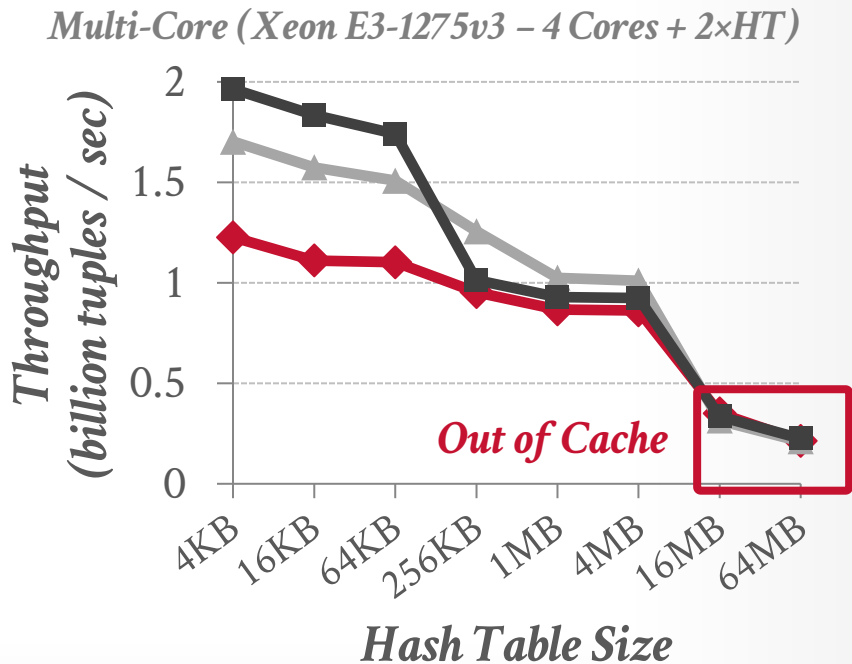
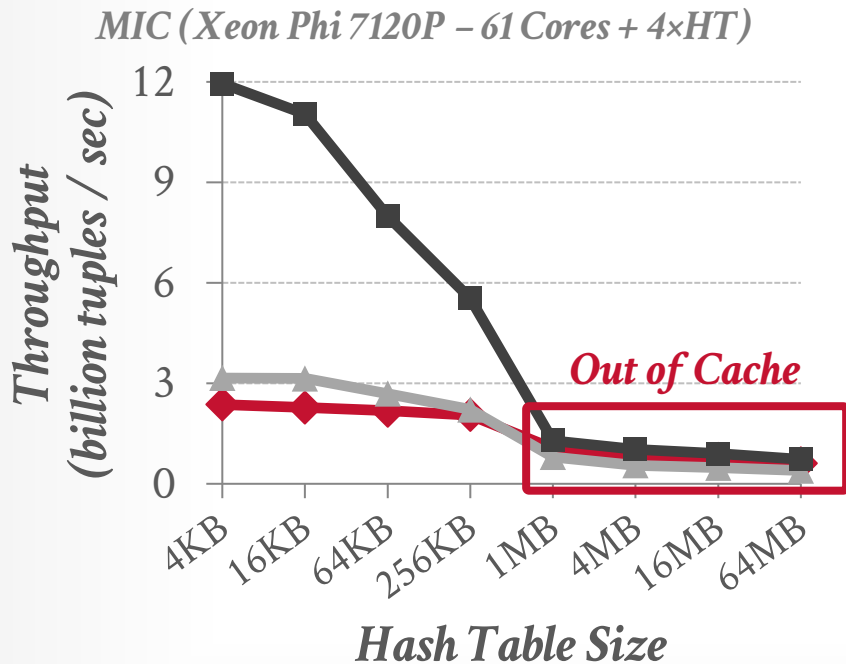
◆ Scalar ▲ Vectorized (Horizontal) ■ Vectorized (Vertical)



Source: [Orestis Polychroniou](#)

HASH TABLES - PROBING

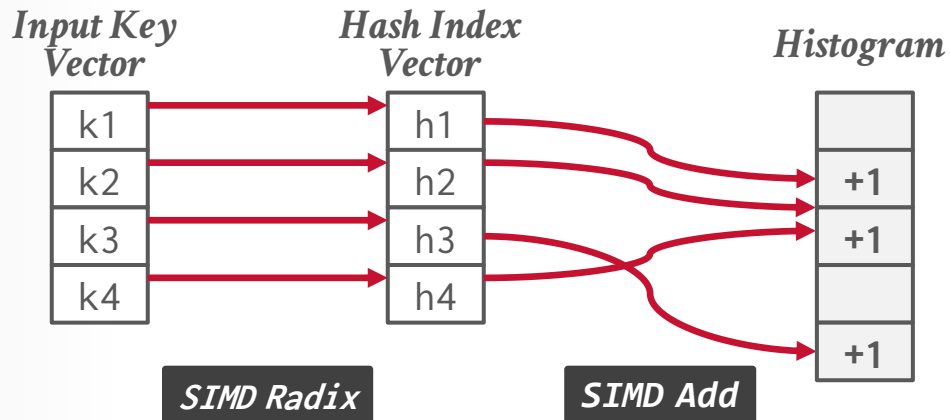
◆ Scalar ▲ Vectorized (Horizontal) ■ Vectorized (Vertical)



Source: [Orestis Polychroniou](#)

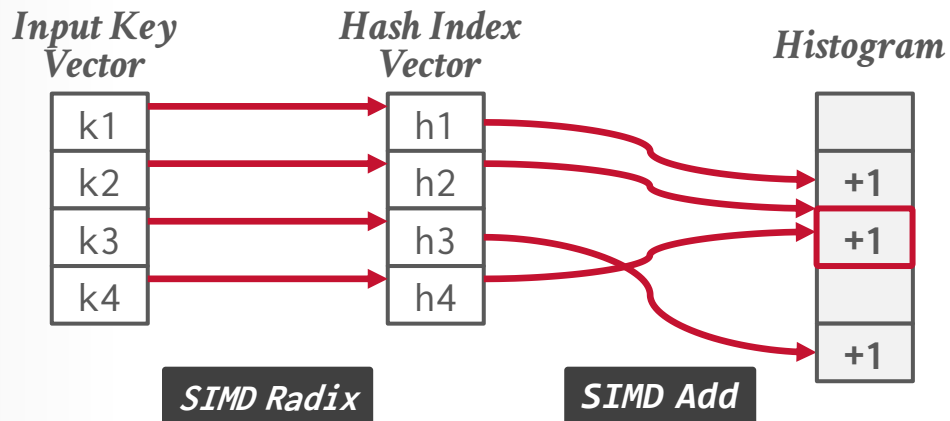
PARTITIONING - HISTOGRAM

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



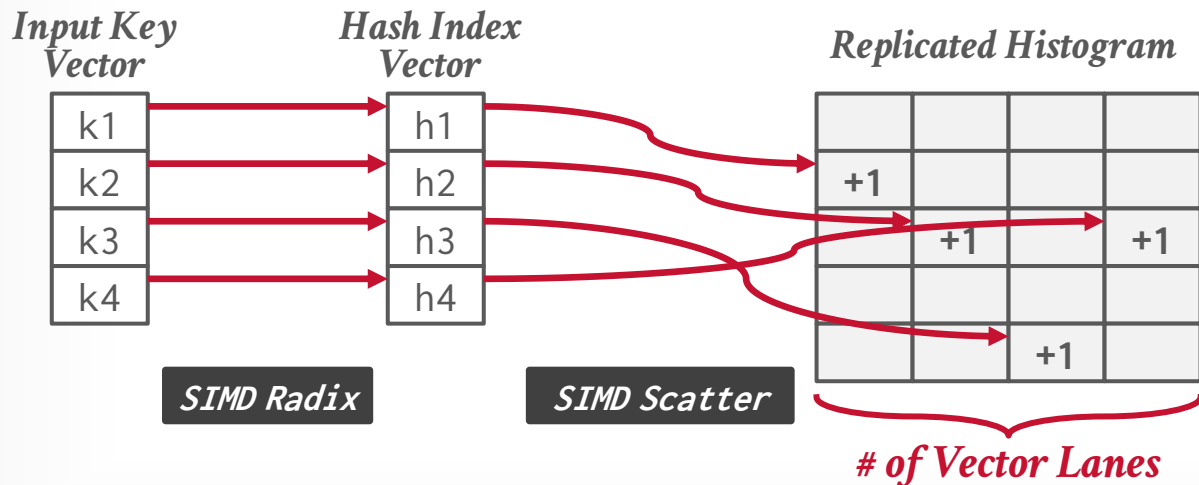
PARTITIONING - HISTOGRAM

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



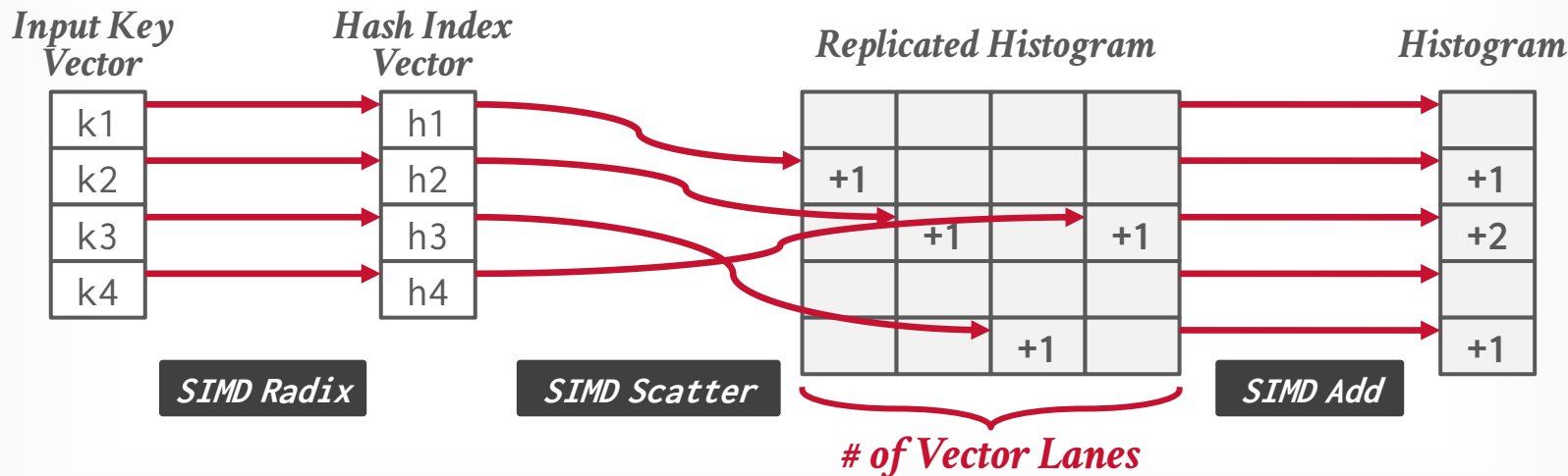
PARTITIONING - HISTOGRAM

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



PARTITIONING - HISTOGRAM

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



CAVEAT EMPTOR

AVX-512 is not always faster than AVX2.

AVX-512 is not always faster

cant underutilization in the preceding operators. We discuss this issue, among other things, in the following section.

5.3 Discussion and implications

The two strategies are not mutually exclusive. Within a single pipeline, both strategies can be applied to individual operators as long as buffering operators are aware of protected lanes (*mixed* strategy). Moreover, the query compiler might decide to *not* apply any refill strategy to certain operators. Especially, when a sequence of operators is quite cheap, divergence might be acceptable as long as the costs for refill operations are not amortized. Naturally, this is a physical query optimization problem that we will leave for future work. Nevertheless, we briefly discuss the advantages and disadvantages, as this is the first work in which we present the basic principles of vector-processing in compiled query pipelines.

As mentioned above, *consume everything* requires additional registers, which increases the register pressure and may lead to spilling. *partial consume* allocates additional registers as well, but these are restricted to (smaller) mask registers. Therefore, it is unlikely to be affected by (potential) performance degradation due to spilling.

The second major difference lies in the cost of refilling empty lanes. In a pipeline that follows the partial consume strategy, the very first operator, that is, the pipeline source, is responsible for refilling empty lanes. If other operators experience underutilization, they return the control flow to the previous operator while retaining ownership of the active lanes. This cascades downward until the source operator is reached, as shown in Fig. 6c. All operators between the pipeline source and the operator that returned the control flow may be subject to underutilization because all lanes in later stages are protected. The costs of refilling, therefore, depend on the length of the pipeline and the costs of the preceding operators. In general, the costs increase in the later stages. Nevertheless, *partial consume* can improve query performance if it is applied only to the very first operators. By contrast, the refilling costs of buffering operators do not depend on the pipeline length. Instead, the crucial factor governing these costs is the number of required buffer registers. The greater the number of buffers, the greater the number of *permute* instructions that need to be executed, whereas the number of required buffers depends on (i) the number of attributes passed along the pipeline and optionally on (ii) the number of registers required to save the internal state of the operator (e.g., a pointer to the current tree node).

6 Evaluation

We evaluate our approach with two major sources of control flow divergence, (i) predicate evaluation as part of a

Table 1 Hardware platforms

Model	Intel Knights landing (KNL)	Intel Skylake-X (SKX)
Model	Phi 7210	i9-7900X
Cores (SMT)	64 ($\times 4$)	10 ($\times 2$)
SIMD [bit]	2 \times 512	2 \times 512
Max. clock rate [GHz]	1.5	4.5
L1 cache	64 KiB	32 KiB
L2 cache	1 MiB	1 MiB
L3 cache	–	14 MiB

table scan and (ii) a hash join. Additionally, we experiment with a more complex operator, an approximate geospatial join. The experiments were conducted on an Intel Skylake-X (SKX) and an Intel Knights Landing (KNL) processor (cf., Table 1). The experiments were implemented in C++ and compiled with GCC 5.4.0 at optimization level three ($-O3$) and the target architecture set to `knl`. If not stated otherwise, we ran the experiments in parallel using two threads per core.² We dispatched the work in batches to the individual threads using batch sizes between 2^{16} and 2^{30} tuples. On the KNL platform, we placed the data in high-bandwidth memory (HBM); otherwise, the experiments would have been dominated by memory stalls. To measure the throughputs, we let each experiment run for at least three seconds, possibly consuming the input data multiple times.

6.1 Table scan

To evaluate the effects of divergence handling in table scans, we integrate our refill algorithms into the AVX-512 implementation of TPC-H Query 1 of Gubner et al. [6]. Additionally, we implemented and integrated the *materialization* approach as proposed by Menon et al. in [16].

From a high-level perspective, TPC-H Query 1 (or short Q1) is a structurally simple query that operates on a single fact table (`lineitem`) with a single scan predicate. It involves several fixed-point arithmetic operations in the aggregation based on the group by clause. In total, five additional attributes are accessed to compute eight aggregated values per group. Almost all tuples survive the selection (i.e., selectivity ≈ 0.98). Therefore, in its original form, Q1 does not suffer from control flow divergence. To simulate control flow divergence and the resulting underutilization of SIMD

² Please note that throughout our (multi-threaded) experiments, we did not observe any performance penalties through downclocking. Both processors KNL and SKX run stable at 1.4GHz and 4.0GHz, respectively.

AVX-512 is not always faster

cant underutilization in the preceding operators. We discuss this issue, among other things, in the following section.

5.3 Discussion and implications

The two strategies are not mutually exclusive. Within a single pipeline, both strategies can be applied to individual operators as long as buffering operators are aware of protected lanes (*mixed* strategy). Moreover, the query compiler might decide to *not* apply any refill strategy to certain operators. Especially, when a sequence of operators is quite cheap, divergence might be acceptable as long as the costs for refill operations are not amortized. Naturally, this is a physical query optimization problem that we will leave for future work. Nevertheless, we briefly discuss the advantages and disadvantages, as this is the first work in which we present the basic principles of vector-processing in compiled query pipelines.

As mentioned above, *consume everything* requires additional registers, which increases the register pressure and may lead to spilling. *partial consume* allocates additional registers as well, but these are restricted to (smaller) mask registers. Therefore, it is unlikely to be affected by (potential) performance degradation due to spilling.

The second major difference lies in the cost of refilling empty lanes. In a pipeline that follows the partial consume strategy, the very first operator, that is, the pipeline source, is responsible for refilling empty lanes. If other operators experience underutilization, they return the control flow to the previous operator while retaining ownership of the active lanes. This cascades downward until the source operator is reached, as shown in Fig. 6c. All operators between the pipeline source and the operator that returned the control flow may be subject to underutilization because all lanes in later stages are protected. The costs of refilling, therefore, depend on the length of the pipeline and the costs of the preceding operators. In general, the costs increase in the later stages. Nevertheless, partial consume can improve query performance if it is applied only to the very first operators. By contrast, the refilling costs of buffering operators do not depend on the pipeline length. Instead, the crucial factor governing these costs is the number of required buffer registers. The greater the number of buffers, the greater the number of `permute` instructions that need to be executed, whereas the number of required buffers depends on (i) the number of attributes passed along the pipeline and optionally on (ii) the number of registers required to save the internal state of the operator (e.g., a pointer to the current tree node).

6 Evaluation

We evaluate our approach with two major sources of control flow divergence, (i) predicate evaluation as part of a

Table 1 Hardware platforms

	Intel Knights landing (KNL)	Intel Skylake-X (SKX)
Model	Phi 7210	i9-7900X
Cores (SMT)	64 ($\times 4$)	10 ($\times 2$)
SIMD [bit]	2 \times 512	2 \times 512
Max. clock rate [GHz]	1.5	4.5
L1 cache	64 KiB	32 KiB
L2 cache	1 MiB	1 MiB
L3 cache	–	14 MiB

table scan and (ii) a hash join. Additionally, we experiment with a more complex operator, an approximate geospatial join. The experiments were conducted on an Intel Skylake-X (SKX) and an Intel Knights Landing (KNL) processor (cf., Table 1). The experiments were implemented in C++ and compiled with GCC 5.4.0 at optimization level three (`-O3`) and the target architecture set to `knl`. If not stated otherwise, we ran the experiments in parallel using two threads per core.² We dispatched the work in batches to the individual threads using batch sizes between 2^{16} and 2^{30} tuples. On the KNL platform, we placed the data in high-bandwidth memory (HBM); otherwise, the experiments would have been dominated by memory stalls. To measure the throughputs, we let each experiment run for at least three seconds, possibly consuming the input data multiple times.

6.1 Table scan

To evaluate the effects of divergence handling in table scans, we integrate our refill algorithms into the AVX-512 implementation of TPC-H Query 1 of Gubner et al. [6]. Additionally, we implemented and integrated the *materialization* approach as proposed by Menon et al. in [16].

From a high-level perspective, TPC-H Query 1 (or short Q1) is a structurally simple query that operates on a single fact table (`lineitem`) with a single scan predicate. It involves several fixed-point arithmetic operations in the aggregation based on the group by clause. In total, five additional attributes are accessed to compute eight aggregated values per group. Almost all tuples survive the selection (i.e., selectivity ≈ 0.98). Therefore, in its original form, Q1 does not suffer from control flow divergence. To simulate control flow divergence and the resulting underutilization of SIMD

² Please note that throughout our (multi-threaded) experiments, we did not observe any performance penalties through downclocking. Both processors KNL and SKX run stable at 1.4GHz and 4.0GHz, respectively.

AVX-512 is not always faster

² Please note that throughout our (multi-threaded) experiments, we did not observe any performance penalties through downclocking. Both processors KNL and SKX run stable at 1.4 GHz and 4.0 GHz, respectively.

cant underutilization in the preceding operators. We discuss this issue, among other things, in the following section.

5.3 Discussion and implications

The two strategies are not mutually exclusive. Within a single pipeline, both strategies can be applied to individual operators as long as buffering operators are aware of protected lanes (*mixed* strategy). Moreover, the query compiler might decide to *not* apply any refill strategy to certain operators. Especially, when a sequence of operators is quite cheap, divergence might be acceptable as long as the costs for refill operations are not amortized. Naturally, this is a physical query optimization problem that we will leave for future work. Nevertheless, we briefly discuss the advantages and disadvantages, as this is the first work in which we present the basic principles of vector-processing in compiled query pipelines.

As mentioned above, *consume everything* requires additional registers, which increases the register pressure.

Table 1 Hardware platforms

	Intel Knights landing (KNL)	Intel Skylake-X (SKX)
Model	Phi 7210	i9-7900X
Cores (SMT)	64 ($\times 4$)	10 ($\times 2$)
SIMD [bit]	2×512	2×512
Max. clock rate [GHz]	1.5	4.5
L1 cache	64 KiB	32 KiB
L2 cache	1 MiB	1 MiB
L3 cache	–	14 MiB

table scan and (ii) a hash join. Additionally, we experiment with a more complex operator, an approximate geospatial join. The experiments were conducted on an Intel Skylake-X (SKX) and an Intel Knights Landing (KNL) processor (cf., Table 1). The experiments were implemented in C++ and

at optimization level three (`-O3`) set to `kn1`. If not stated otherwise, in parallel using two threads. We work in batches to the individual between 2^{16} and 2^{30} tuples. On the the data in high-bandwidth memory experiments would have been. To measure the throughputs, for at least three seconds, possibly multiple times.

divergence handling in table algorithms into the AVX-512 query 1 of Gubner et al. [6]. Additionally, we integrated the *materialization* from Gubner et al. in [16].

ive, TPC-H Query 1 (or short query that operates on a single query with a single scan predicate.

The greater the number of buffers, the greater the number of *permute* instructions that need to be executed, whereas the number of required buffers depends on (i) the number of attributes passed along the pipeline and optionally on (ii) the number of registers required to save the internal state of the operator (e.g., a pointer to the current tree node).

6 Evaluation

We evaluate our approach with two major sources of control flow divergence, (i) predicate evaluation as part of a

it involves several fixed-point arithmetic operations in the aggregation based on the group by clause. In total, five additional attributes are accessed to compute eight aggregated values per group. Almost all tuples survive the selection (i.e., selectivity ≈ 0.98). Therefore, in its original form, Q1 does not suffer from control flow divergence. To simulate control flow divergence and the resulting underutilization of SIMD

² Please note that throughout our (multi-threaded) experiments, we did not observe any performance penalties through downclocking. Both processors KNL and SKX run stable at 1.4 GHz and 4.0 GHz, respectively.

CAVEAT EMPTOR

AVX-512 is not always faster than AVX2.

Some CPUs downgrade their clockspeed when switching to AVX-512 mode.

→ Compilers will prefer 256-bit SIMD operations.

If only a small portion of the process uses AVX-512, then it is not worth the downclock penalty.

EMPTOR

The frequency impact depends on the width of the operation and the specific instruction used.

71 There are three frequency levels, so-called licenses, from fastest to slowest: L0, L1 and L2. L0 is the "nominal" speed you'll see written on the box: when the chip says "3.5 GHz turbo", they are referring to the single-core L0 turbo. L1 is a lower speed sometimes called AVX turbo or AVX2 turbo⁵, originally associated with AVX and AVX2 instructions¹. L2 is a lower speed than L1, sometimes called "AVX-512 turbo".

The exact speeds for each license also depend on the number of active cores. For up to date tables, you can usually consult WikiChip. For example, the table for the Xeon Gold 5120 is here:

Mode	Base	Turbo Frequency/Active Cores													
		1	2	3	4	5	6	7	8	9	10	11	12	13	14
Normal	2,200 MHz	3,500 MHz	3,300 MHz	3,000 MHz	2,900 MHz	2,900 MHz	2,900 MHz	2,900 MHz	2,900 MHz	2,900 MHz	2,700 MHz	2,700 MHz	2,700 MHz	2,600 MHz	2,600 MHz
AVX2	1,800 MHz	3,100 MHz	3,100 MHz	2,900 MHz	2,900 MHz	2,700 MHz	2,700 MHz	2,700 MHz	2,700 MHz	2,300 MHz	2,300 MHz	2,300 MHz	2,200 MHz	2,200 MHz	
AVX512	1,200 MHz	2,900 MHz	2,900 MHz	2,500 MHz	2,500 MHz	1,900 MHz	1,900 MHz	1,900 MHz	1,900 MHz	1,600 MHz	1,600 MHz	1,600 MHz	1,600 MHz	1,600 MHz	

The Normal, AVX2 and AVX512 rows correspond to the L0, L1 and L2 licenses respectively. Note that the relative slowdown for L1 and L2 licenses generally gets worse as the number of cores increase: for 1 or 2 active cores the L1 and L2 speeds are 97% and 91% of L0, but for 13 or 14 cores they are 85% and 62% respectively. This varies by chip, but the general trend is usually the same.

Those preliminaries out of the way, let's get to what I think you are asking: which instructions cause which licenses to be activated?

Here's a table, showing the implied license for instructions based on their width and their categorization as light or heavy:

Width	Light	Heavy
Scalar	L0	N/A
128-bit	L0	L0
256-bit	L0	L1*
512-bit	L1	L2*

*soft transition (see below)

So we immediately see that all scalar (non-SIMD) instructions and all 128-bit wide instructions² always run at full speed in the L0 license.

... slower than AVX2.

... their clockspeed when de.

... SIMD operations.

... the process uses AVX-512, ... synclock penalty.

PARTING THOUGHTS

Vectorization is essential for OLAP queries.

But implementing an algorithm using SIMD is still mostly a manual process.

We can combine all the intra-query parallelism optimizations we've talked about in a DBMS.

- Multiple threads processing the same query.
- Each thread can execute a compiled plan.
- The compiled plan can invoke vectorized operations.

NEXT CLASS

Query Compilation

Project Status Discussion