

#### **Carnegie Mellon University** ADVANCED DATABASE SYSTEMS

# Parallel Hash Join Algorithms

Andy Pavlo // 15-721 // Spring 2023

# ADMINISTRIVIA

#### Project #2:

- $\rightarrow$  Feedback Submission: Saturday April 1<sup>st</sup>
- $\rightarrow$  Final Submission: Monday May 1<sup>st</sup>
- $\rightarrow$  Sign up for a system if you haven't yet!

#### Project #3

- $\rightarrow$  Proposal Presentation: Wednesday March 1<sup>st</sup>
- → Status Update Presentation: Monday April 3<sup>rd</sup>
- $\rightarrow$  Final Presentations: TBA

## **TODAY'S AGENDA**

Background Parallel Hash Join Hash Functions Hashing Schemes Evaluation

# PARALLEL JOIN ALGORITHMS

Perform a join between <u>two</u> relations on multiple threads simultaneously to speed up operation.  $\rightarrow$  We will discuss multi-way joins in <u>Lecture #13</u>.

Two main approaches:

- $\rightarrow$  Hash Join
- $\rightarrow$  Sort-Merge Join

We won't discuss nested-loop joins because an OLAP DBMS almost never wants to use this...



#### 5

# OBSERVATION

Many OLTP DBMSs do not implement hash join.

# But an **index nested-loop join** is conceptually equivalent to a hash join.

- $\rightarrow$  Index NL joins typically means using an existing B+Tree.
- $\rightarrow$  Hash join will build a hash table (index) on the fly and then discard immediately after the operation is complete.

#### HASHING VS. SORTING JOINS

- 1970s Sorting
- 1980s Hashing
- 1990s Equivalent
- 2000s Hashing

2010s – Hashing (Partitioned vs. Non-Partitioned)2020s – Non-Partitioned Hashing



# PARALLEL JOIN ALGORITHMS



## JOIN ALGORITHM DESIGN GOALS

These goals matter whether the DBMS is using a hardware-conscious vs. hardware-oblivious algorithm for joins.

#### **Goal #1: Minimize Synchronization**

 $\rightarrow$  Avoid taking latches during execution.

#### Goal #2: Minimize Memory Access Cost

- $\rightarrow$  Ensure that data is always local to worker thread.  $\rightarrow$  Reuse data while it exists in CPU cache.



## IMPROVING CACHE BEHAVIOR

Factors that affect cache misses in a DBMS:

- $\rightarrow$  Cache + TLB capacity.
- $\rightarrow$  Locality (temporal and spatial).

#### Non-Random Access (Scan):

- $\rightarrow$  Clustering data to a cache line.
- $\rightarrow$  Execute more operations per cache line.

#### Random Access (Lookups):

 $\rightarrow$  Partition data to fit in cache + TLB.

#### PARALLEL HASH JOINS

Hash join is one of the most important operators in a DBMS for OLAP workloads.  $\rightarrow$  But it is still not the dominant cost.

It is important that we speed up our DBMS's join algorithm by taking advantage of multiple cores.
→ We want to keep all cores busy, without becoming memory bound.

# HASH JOIN (R⊳S)

#### Phase #1: Partition (optional)

 $\rightarrow$  Divide the tuples of **R** and **S** into disjoint subsets using a hash function on the join key.

#### Phase #2: Build

 $\rightarrow$  Scan relation **R** and create a hash table on join key.

#### Phase #3: Probe

 $\rightarrow$  For each tuple in **S**, look up its join key in hash table for **R**. If a match is found, output combined tuple.





# PARTITIONING PHASE

#### Approach #1: Implicit Partitioning

- $\rightarrow$  The data was partitioned on the join key when it was loaded into the database.
- $\rightarrow$  No extra pass over the data is needed.

#### Approach #2: Explicit Partitioning

- → Divide only the outer relation and redistribute among the different CPU cores.
- $\rightarrow$  Can use the same radix partitioning approach we talked about last time.

# PARTITION PHASE

Split the input relations into partitioned buffers by hashing the tuples' join key(s).

- $\rightarrow$  Ideally the cost of partitioning is less than the cost of cache misses during build phase.
- → Sometimes called *Grace Hash Join / Radix Hash Join*.

Contents of buffers depends on storage model:

- $\rightarrow$  **NSM**: Usually the entire tuple.
- $\rightarrow$  **DSM**: Only the columns needed for the join + offset.

# PARTITION PHASE

#### Approach #1: Non-Blocking Partitioning

- $\rightarrow$  Only scan the input relation once.
- → Produce output incrementally and let other threads build hash table at the same time.

#### Approach #2: Blocking Partitioning (Radix)

- $\rightarrow$  Scan the input relation multiple times.
- $\rightarrow$  Only materialize results all at once.
- $\rightarrow$  Sometimes called *radix hash join*.

## NON-BLOCKING PARTITIONING

Scan the input relation only once and generate the output on-the-fly.

#### **Approach #1: Shared Partitions**

- $\rightarrow$  Single global set of partitions that all threads update.
- $\rightarrow$  Must use a latch to synchronize threads.

#### **Approach #2: Private Partitions**

- $\rightarrow$  Each thread has its own set of partitions.
- $\rightarrow$  Must consolidate them after all threads finish.



#### Data Table







**ECMU-DB** 15-721 (Spring 2023)



**CMU·DB** 15-721 (Spring 2023)



**CMU-DB** 15-721 (Spring 2023)



**ECMU-DB** 15-721 (Spring 2023)







#### **ECMU-DB** 15-721 (Spring 2023)



**CMU·DB** 15-721 (Spring 2023)



**CMU·DB** 15-721 (Spring 2023)

## RADIX PARTITIONING

Scan the input relation multiple times to generate the partitions.

Two-pass algorithm:

- $\rightarrow$  Step #1: Scan R and compute a histogram of the # of tuples per hash key for the radix at some offset.
- $\rightarrow$  Step #2: Use this histogram to determine per-thread output offsets by computing the <u>prefix sum</u>.
- $\rightarrow$  **Step #3:** Scan **R** again and partition them according to the hash key.

# The radix of a key is the value of an integer at a position (using its base).

 $\rightarrow$  Efficient to compute with bitshifting + multiplication.

#### *Keys* 19 12 23 08 11 04



# The radix of a key is the value of an integer at a position (using its base).

 $\rightarrow$  Efficient to compute with bitshifting + multiplication.





# The radix of a key is the value of an integer at a position (using its base).

 $\rightarrow$  Efficient to compute with bitshifting + multiplication.





The radix of a key is the value of an integer at a position (using its base).

 $\rightarrow$  Efficient to compute with bitshifting + multiplication.

Compute radix for each key and populate histogram of counts per radix.



**ECMU-DB** 15-721 (Spring 2023)

















# PREFIX SI

Scan Primitives for Vector Computers\* Siddhartha Chatterjee

Guy E. Blelloch

Marco Zagha

School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213-3890

#### Abstract

This paper describes an optimized implementation of a set of scan (also called all-prefix-sums) primitives on a single processor of a CRAY Y-MP, and demonstrates that their use leads to greatly improved performance for several applications that cannot be vectorized with existing compiler technology. The algorithm used to implement the scans is based on an algorithm for parallel computers and is applicable with minor modifications to any register-based vector computer. On the CRAY Y-MP, the asymptotic running time of the plus-scan is about 2.25 times that of a vector add, and is within 20% of optimal. An important aspect of our implementation is that a set of segmented versions of these scans are only marginally more expensive than the unsegmented versions. These segmented versions can be used to execute a scan on multiple data sets without having to pay the vector startup cost (n1/2) for each set.

The paper describes a radix sorting routine based on the scans that is 13 times faster than a Fortran version and within 20% of a highly optimized library sort routine, three operations on trees that are between 10 and 20 times faster than the corresponding C versions, and a connectionist learning algorithm that is 10 times faster than the corresponding C version for sparse and irregular

#### 1 Introduction

Vector supercomputers have been used to supply the high computing power needed for many applications. However, the performance obtained from these machines critically depends on the ability to produce code that vectorizes well. Two distinct approaches have been taken to meet this goal-vectorization of "dusty decks" [18], and language support for vector intrinsics, as seen in the proposed Fortran 8x standard [1]. In both cases, the focus of the work has been in speeding up "scientific" computations, characterized by regular and static data structures and predominantly regular access patterns within these data structures. These alternatives are not very effective for problems that

author and should not be interpreted as representing the official policies, either expressed or implied, of DARPA or the U.S. government.

#### CH2916-5/90/0000/0666/\$01.00 © IEEE

Routine	/, (clks/elt)	
	Scan version	Scalar version
segmented plus scan (int) parallel radix sort (64 bits) branch_sums	2.5 3.7 896.8	7.5 (*) 74.8 (*) 11730.0 (*)
root_sums delete_vertices	9.8 19.2	206.5 (†) 208.6 (†) 276.2 (†)

Table 1: Incremental processing times per element for primitives and applications discussed in this paper, for both scan and scalar versions. All numbers are for a single processor of a CRAY Y-MP. 1 clock tick = 6 ns. Items marked with (+) were written in Fortran and those marked with (†) were written in C.

create and manipulate more irregular and dynamically varying data structures such as trees and graphs.

Elsewhere, scan (prefix sum) operations have been shown to be extremely powerful primitives in designing parallel algorithms for manipulating such irregular and dynamically changing data structures [3]. This paper shows how the scan operations can also have great benefit for such algorithms on pipelined vector machines. It describes an optimized implementation of the scan primitives on the CRAY Y-MP1, and gives performance numbers for several applications based on the primitives (see Table 1). The approach in the design of these algorithms is similar to that of the Basic Linear Algebra Subprograms (BLAS) developed in the context of linear algebra computations [14] in that the algorithms are based on a set of primitives whose implementations are optimized rather than having a compiler try to vectorize existing

The remainder of this paper is organized as follows. Section 2 introduces the scan primitives and reviews previous work on vectorizing scans. Section 3 discusses our implementation in detail and presents performance numbers. Section 4 presents other primitives used in this paper, and three applications using these primitives. Finally, future work and conclusions are given

CRAY Y-MP and CFT77 are trademarks of Cray Research, Inc.

2All Fortran code discussed in this paper was compiled with CFT77 Version 3.1. All C code was compiled with Cray PCC Version XMP/YMP 4.1.8.

The prefix sum of a sequence of  $(x_0, x_1, ..., x_n)$ is a second sequence of numbers  $(y_0, y_1, ..., y_n)$ that is a running total of the inp



#### SECMU-DB 15-721 (Spring 2023)

<sup>&</sup>quot;This research was supported by the Defense Advanced Research Projects Agency (DOD) and monitored by the Avionics Laboratory, Air Force Wright Aeronautical Laboratories, Aeronautical Systems Division (AFSC), Wright-Patterson AFB, Ohio 45433-6543 under Contract F33615-87-C-1499, ARPA The views and conclusions contained in this document are those of the

#### **RADIX PARTITIONS**

#### Step #1: Inspect input, create histograms





**ECMU·DB** 15-721 (Spring 2023)














35

Step #3: Read input and partition 35



Recursively repeat until target number of partitions have been created







# **OPTIMIZATIONS**

#### Software Write Combine Buffers:

- $\rightarrow$  Each worker maintains local output buffer to stage writes.
- $\rightarrow$  When buffer full, write changes to global partition.
- $\rightarrow$  Similar to private partitions but without a separate write phase at the end.

#### Non-temporal Streaming Writes

 $\rightarrow$  Workers write data to global partition memory using streaming instructions to bypass CPU caches.



#### 46

# BUILD PHASE

The threads are then to scan either the tuples (or partitions) of R.

For each tuple, hash the join key attribute for that tuple and add it to the appropriate bucket in the hash table.

 $\rightarrow$  The buckets should only be a few cache lines in size.

## HASH TABLES

#### **Design Decision #1: Hash Function**

- $\rightarrow$  How to map a large key space into a smaller domain.
- $\rightarrow$  Trade-off between being fast vs. collision rate.

#### **Design Decision #2: Hashing Scheme**

- $\rightarrow$  How to handle key collisions after hashing.
- $\rightarrow$  Trade-off between allocating a large hash table vs. additional instructions to find/insert keys.

## HASH FUNCTIONS

We do not want to use a cryptographic hash function for our join algorithm.

We want something that is fast and will have a low collision rate.

- $\rightarrow$  **Best Speed:** Always return '1'
- → **Best Collision Rate:** Perfect hashing

See <u>SMHasher</u> for a comprehensive hash function benchmark suite.

We do not wan function for ou

We want some collision rate.

 $\rightarrow \textbf{Best Speed: A} \\ \rightarrow \textbf{Best Collision}$ 

See <u>SMHasher</u> benchmark sui

	Public			🏹 Sponsor	A Notifications
Code 🕢 Issues 26	1 Pull requests 1	Actions	Projects	🕮 Wik	di <sup>1</sup> Security <sup>1</sup> Insights
ਿ master - smhasher /	README.md				3.12
rurban results of fwojcik's	fixes ×				Go to file
🞗 7 contributors 🏼 🗿 🦱 🖪					Latest commit 9a5e665 on Aug 19 🕚 History
	) 🕲 🗲 🔟 🎝				
444 lines (395 sloc)	37.4 KB				
					<> 🗅 Raw Blame 🖉 🕶 🖞
SMhasher					
build passion O to un					
build passing 📀 build failing	build failing				
build passing ③ build failing Hash function	build failing MiB/sec	cycl./hash	cycl./map	size	<b>0</b> . III
build passing @ build failing Hash function donothing32	build   falling     MiB/sec   15316474.36	cycl./hash 6.00	cycl./map	size	Quality problems
build passing Duild failing Hash function donothing32 donothing64	built   talling     MiB/sec   15316474.36     15330019.19   15330019.19	cycl./hash 6.00 6.00	cycl./map -	<b>size</b> 13	Quality problems bad seed 0, test NOP
build passing Duild tailing Hash function donothing32 donothing64 donothing128	built   MiB/sec     15316474.36   15330019.19     15278983.09   15278983.09	<b>cycl./hash</b> 6.00 6.00 6.00	cycl./map - -	size 13 13	Quality problems bad seed 0, test NOP bad seed 0, test NOP
build passing tuild tailing Hash function donothing32 donothing64 donothing128 NOP_OAAT_read64	built   Yalling     MiB/sec   15316474.36     15330019.19   15278983.09     15278983.09   28467.50	<b>cycl./hash</b> 6.00 6.00 6.00 18.48	cycl./map - - -	size 13 13 13	Quality problems bad seed 0, test NOP bad seed 0, test NOP bad seed 0, test NOP
build passing Duild failing Hash function donothing32 donothing64 donothing128 NOP_OAAT_read64 BadHash	built   talling     MiB/sec   15316474.36     15330019.19   15278983.09     15278983.09   28467.50     28467.51   524.81	cycl./hash 6.00 6.00 18.48 96.20	cycl./map - - - -	size 13 13 13 13 47	Quality problems bad seed 0, test NOP bad seed 0, test NOP bad seed 0, test NOP
build passing @ build failing Hash function donothing32 donothing64 donothing128 NOP_OAAT_read64 BadHash sumhash	built   MiB/sec     15316474.36   15330019.19     15278983.09   15278983.09     28467.50   524.81     524.81   7169.08	<b>cycl./hash</b> 6.00 6.00 18.48 96.20 27.12	cycl./map - - - - - -	size 13 13 13 13 47 47	Quality problems bad seed 0, test NOP bad seed 0, test NOP bad seed 0, test FAIL
build passing tailing Hash function donothing32 donothing64 donothing128 NOP_OAAT_read64 BadHash sumhash Sumhash32	built   Alling     MiB/sec   15316474.36     15330019.19   15278983.09     15278983.09   28467.50     28467.50   524.81     7169.08   22556 18	cycl./hash 6.00 6.00 6.00 18.48 96.20 27.12	cycl./map - - - - - - - - -	size 13 13 13 13 47 47 363	Quality problems   bad seed 0, test NOP   bad seed 0, test NOP   bad seed 0, test NOP   test NOP   bad seed 0, test FAIL   bad seed 0, test FAIL
build passing Duild failing Hash function donothing32 donothing64 donothing128 NOP_OAAT_read64 BadHash sumhash sumhash32 multiply_shift	built   talling     MiB/sec   15316474.36     15316474.36   15330019.19     15278983.09   15278983.09     28467.50   28467.50     524.81   524.81     22556.18   22556.18	cycl./hash 6.00 6.00 18.48 96.20 27.12 22.98	cycl./map - - - - - - - - - - - -	size 13 13 13 13 47 47 363 863	Quality problems   bad seed 0, test NOP   bad seed 0, test NOP   bad seed 0, test NOP   bad seed 0, test FAIL   bad seed 0, test FAIL   UB, test FAIL
build passing Duild failing Hash function donothing32 donothing64 donothing128 NOP_OAAT_read64 BadHash sumhash sumhash32 multiply_shift pair_multiply_shift	built   MiB/sec     I5316474.36   15330019.19     I53330019.19   15278983.09     I5278983.09   28467.50     I5278983.09   28467.50     I5278983.09   28467.50     I5278983.09   28467.50     I5278983.09   28467.50     I5278983.09   28467.50     I5278983.09   15278983.09     I5278983.09   15278983.09     I5278983.09   15278983.09     I5278983.09   15278983.09     IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII	cycl./hash     6.00     6.01     6.02     6.03     6.04     6.05     6.05     6.06     6.07     6.08     6.09     6.00     6.00     6.00     6.00     6.00     6.00     6.00     6.00     6.00     6.00     7.12     22.98     28.69	cycl./map 	size 13 13 13 13 47 47 363 863 345	Quality problems   bad seed 0, test NOP   bad seed 0, test NOP   bad seed 0, test NOP   bad seed 0, test FAIL   bad seed 0, test FAIL

# HASH FUNCTIONS

#### <u>CRC-64</u> (1975)

 $\rightarrow$  Used in networking for error detection.

#### MurmurHash (2008)

 $\rightarrow$  Designed to a fast, general purpose hash function.

#### **Google CityHash** (2011)

 $\rightarrow$  Designed to be faster for short keys (<64 bytes).

#### Facebook XXHash (2012)

 $\rightarrow$  From the creator of zstd compression.

#### **Google FarmHash** (2014)

 $\rightarrow$  Newer version of CityHash with better collision rates.

## HASH FUNCTION BENCHMARK

#### Intel Core i7-8700K @ 3.70GHz



15-721 (Spring 2023)

#### HASHING SCHEMES

**Approach #1: Chained Hashing** 

**Approach #2: Linear Probe Hashing** 

Approach #3: Robin Hood Hashing

Approach #4: Hopscotch Hashing

Approach #5: Cuckoo Hashing



Maintain a linked list of <u>buckets</u> for each slot in the hash table.

Resolve collisions by placing all elements with the same hash key into the same bucket.

- $\rightarrow$  To determine whether an element is present, hash to its bucket and scan for it.
- $\rightarrow$  Insertions and deletions are generalizations of lookups.



















Single giant table of slots.

Resolve collisions by linearly searching for the next free slot in the table.

- $\rightarrow$  To determine whether an element is present, hash to a location in the table and scan for it.
- $\rightarrow$  Must store the key in the table to know when to stop scanning.
- $\rightarrow$  Insertions and deletions are generalizations of lookups.




























#### 73

# OBSERVATION

To reduce the number of wasteful comparisons during the build/probe phases, it is important to avoid collisions of hashed keys.

This requires a hash table with  $\sim 2 \times$  the number of slots as the number of elements in **R**.

Variant of linear probe hashing that steals slots from "rich" keys and give them to "poor" keys.

- $\rightarrow$  Each key tracks the number of positions they are from where its optimal position in the table.
- $\rightarrow$  On insert, a key takes the slot of another key if the first key is farther away from its optimal position than the second key.























Variant of linear probe hashing where keys can move between positions in a **<u>neighborhood</u>**.

- $\rightarrow$  A neighborhood is contiguous range of slots in the table.
- $\rightarrow$  The size of a neighborhood is a configurable constant (ideally a single cache-line).
- $\rightarrow$  A key is guaranteed to be in its neighborhood or not exist in the table.

The goal is to have the cost of accessing a neighborhood to be the same as finding a key.



#### Sec MU·DB

15-721 (Spring 2023)



**ECMU·DB** 15-721 (Spring 2023)

84



**SECMU-DB** 15-721 (Spring 2023)



Neighborhood Size = 3 Neighborhood #1 Neighborhood #6 Neighborhood #2 Neighborhood #3 Neighborhood #4 Neighborhood #6















Neighborhood Size = 3

#### Neighborhood #4



**ECMU·DB** 15-721 (Spring 2023)

#### Neighborhood #4



**ECMU·DB** 15-721 (Spring 2023)

Neighborhood Size = 3

#### Neighborhood #4











Neighborhood #4



Neighborhood #4





Neighborhood Size = 3

**EXAMPLE DB** 15-721 (Spring 2023)







**ECMU·DB** 15-721 (Spring 2023)

Neighborhood #6



Neighborhood #6



Neighborhood #6



Neighborhood #6


### HOPSCOTCH HASHING



Neighborhood #6

Use multiple tables with different hash functions.

- $\rightarrow$  On insert, check every table and pick anyone that has a free slot.
- $\rightarrow$  If no table has a free slot, evict the element from one of them and then re-hash it find a new location.

Look-ups are always O(1) because only one location per hash table is checked.







#### **ECMU-DB** 15-721 (Spring 2023)

#### Hash Table #1



**ECMU-DB** 15-721 (Spring 2023) Hash Table #2

#### Hash Table #1



 $\begin{array}{c} Insert \ X \\ hash_1(X) & hash_2(X) \end{array}$ 

Insert Y hash<sub>1</sub>(Y) hash<sub>2</sub>(Y) -

Hash	T	able	#2
------	---	------	----





#### Hash Table #1



Insert X  $hash_1(X) hash_2(X)$ 

Insert Y hash<sub>1</sub>(Y) hash<sub>2</sub>(Y)

Insert Z - hash<sub>1</sub>(Z) hash<sub>2</sub>(Z) -

#### Hash Table #2



**ECMU-DB** 15-721 (Spring 2023)

#### Hash Table #1



 $\begin{array}{c} Insert \ X \\ hash_1(X) & hash_2(X) \end{array}$ 

Insert Y hash<sub>1</sub>(Y) hash<sub>2</sub>(Y)

 $\frac{\text{Insert } Z}{\text{hash}_1(Z) \quad \text{hash}_2(Z)} \sim$ 

#### Hash Table #2





#### Hash Table #1



Insert X  $hash_1(X) hash_2(X)$ 

Insert Y hash<sub>1</sub>(Y) hash<sub>2</sub>(Y)

Insert Z  $hash_1(Z)$   $hash_2(Z)$  $hash_1(Y)$ 

### Hash Table #2



#### Hash Table #1



Insert X  $hash_1(X) hash_2(X)$ 

Insert Y hash<sub>1</sub>(Y) hash<sub>2</sub>(Y)

 $\begin{array}{c} \text{Insert } Z \\ \text{hash}_1(Z) & \text{hash}_2(Z) \\ \hline \text{hash}_1(Y) \end{array}$ 

### Hash Table #2



**ECMU-DB** 15-721 (Spring 2023)

#### Hash Table #1



 $\begin{array}{c} \text{Insert } X\\ hash_1(X) & hash_2(X) \end{array}$ 

Insert Y hash<sub>1</sub>(Y) hash<sub>2</sub>(Y)

Insert Z  $hash_1(Z) \quad hash_2(Z)$   $hash_1(Y)$   $hash_2(X)$ 

#### Hash Table #2



# PROBE PHASE

For each tuple in S, hash its join key and check to see whether there is a match for each tuple in corresponding bucket in the hash table constructed for R.

- $\rightarrow$  If inputs were partitioned, then assign each thread a unique partition.
- $\rightarrow$  Otherwise, synchronize their access to the cursor on **S**.

Create a Bloom Filter during the build phase when the key is likely to not exist in the hash table.

→ Threads check the filter before probing the hash table. This will be faster since the filter will fit in CPU caches.

 $\rightarrow$  Sometimes called *sideways information passing.* 



MICRO ADAPTIVITY IN VECTORWISE

122

Create a Bloom Filter during the build phase when the key is likely to not exist in the hash table.

→ Threads check the filter before probing the hash table. This will be faster since the filter will fit in CPU caches.

→ Sometimes called *sideways information passing*.

RO ADAPTIVITY IN VECTORWISE

IGMOD 2013



Create a Bloom Filter during the build phase when the key is likely to not exist in the hash table.

→ Threads check the filter before probing the hash table. This will be faster since the filter will fit in CPU caches.

→ Sometimes called *sideways information passing*.



122



Create a Bloom Filter during the build phase when the key is likely to not exist in the hash table.

→ Threads check the filter before probing the hash table. This will be faster since the filter will fit in CPU caches.

→ Sometimes called *sideways information passing*.



122



## HASH JOIN VARIANTS

	No-P	Shared-P	Private-P	Radix
Partitioning	No	Yes	Yes	Yes
Input scans	0	1	1	2
Sync during partitioning	_	Spinlock per tuple	Barrier, once at end	Barrier, 4 · #passes
Hash table	Shared	Private	Private	Private
Sync during build phase	Yes	No	No	No
Sync during probe phase	No	No	No	No

# BENCHMARKS

Implemented multiple variants of hash join algorithms based on previous literature and compare unoptimized vs. optimized versions.

Core approaches:

- $\rightarrow$  No Partitioning Hash Join
- $\rightarrow$  Concise Hash Table Join
- $\rightarrow$  2-pass Radix Hash Join (Chained vs. Linear)

Special Case: Arrays for monotonic primary keys.







15-721 (Spring 2023)

128



129

15-721 (Spring 2023)

# PARTING THOUGHTS

Partitioned-based joins outperform no-partitioning algorithms in most settings, but it is non-trivial to tune it correctly.

AFAIK, every DBMS vendor picks one hash join implementation and does not try to be adaptive.

#### 131

### **NEXT CLASS**

#### Parallel Sort-Merge Joins

