

Parallelism

6.S079 Lecture 18

Sam Madden

4/11/2022

Lab 5 Due Weds

Parallelism Goal

- Make a job faster by running on multiple processors
- What do we mean by faster?

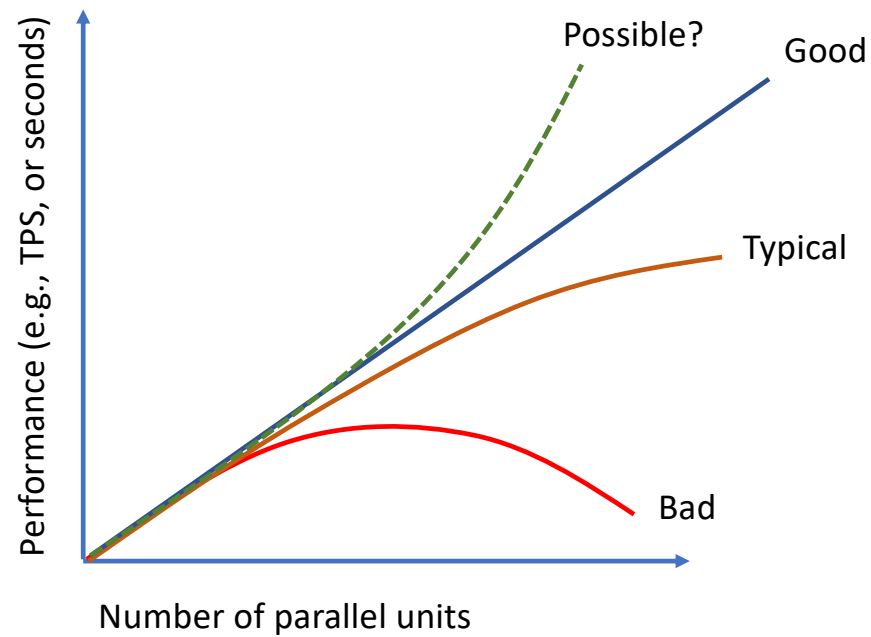
$speed\ up = \frac{old\ time}{new\ time}$ on same problem, with N times more hardware

$scale\ up = \frac{1x\ larger\ problem\ on\ 1x\ hardware}{Nx\ larger\ problem\ on\ Nx\ hardware}$

- Not necessarily the same: smaller problem may be harder to parallelize

Speedup Goal

- Linear?

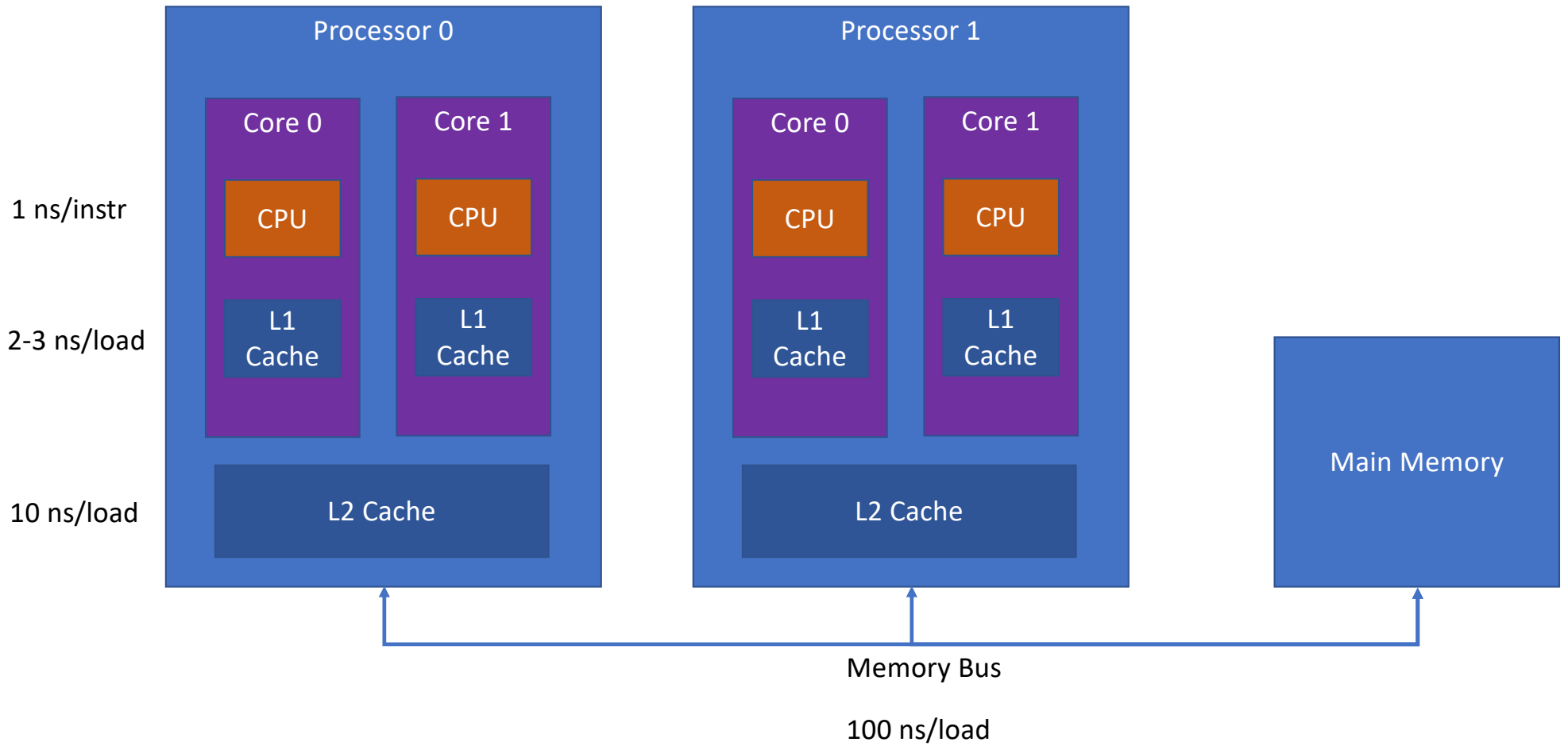


Barriers to Linear Scaling

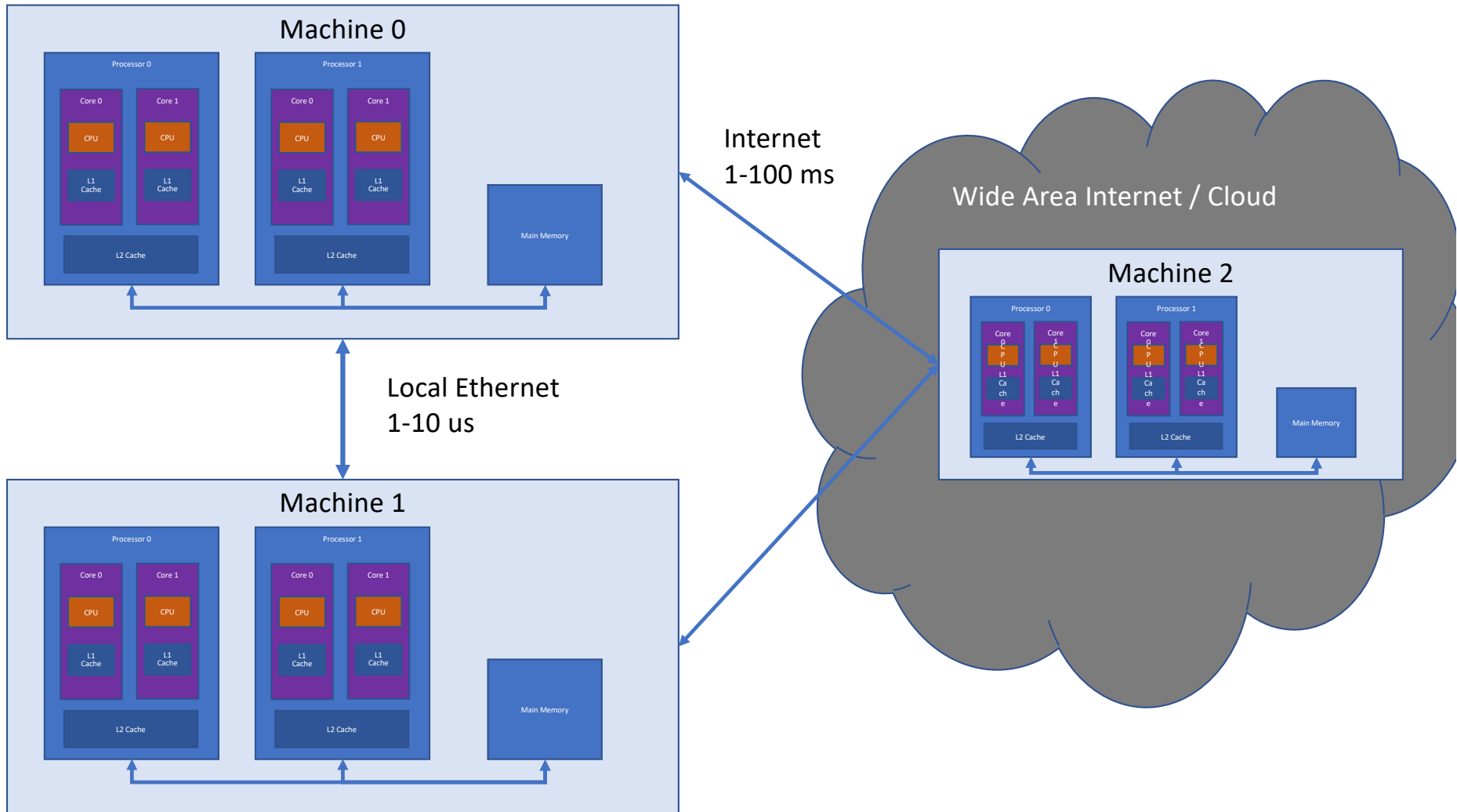
- Startup times
 - e.g., may take time to launch each parallel executor
- Interference
 - processors depend on some shared resource
 - E.g., input or output queue, or other data item
- Skew
 - workload not of equal size on each processor
- *Almost all workloads will stop scaling at some point!*
- What are some barriers in data science workloads?

Properties of Parallelizable Workloads

- Provide linear speedup
- Usually can be decomposed into small units that can be executed independently
 - "embarrassingly parallel"
- As we will see, SQL-style operations generally provide this
- Some ML algorithms support it, but often tricky

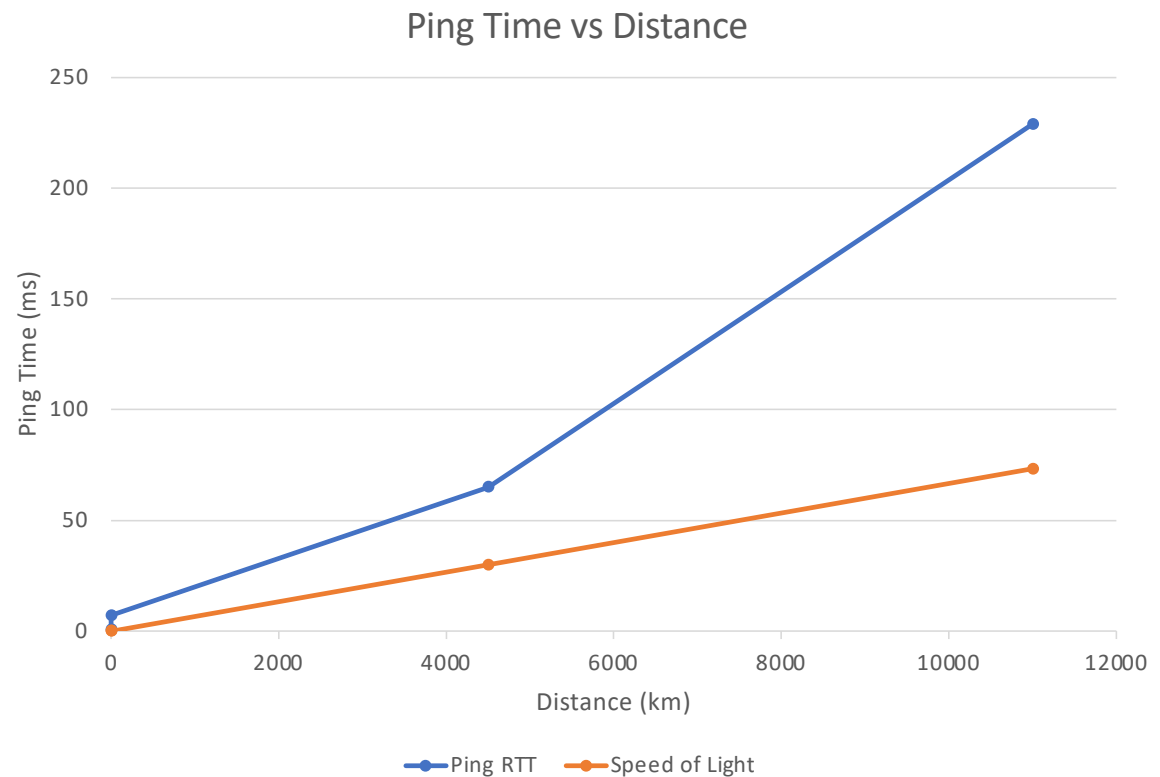


Some machines may have 2 levels of cache per core

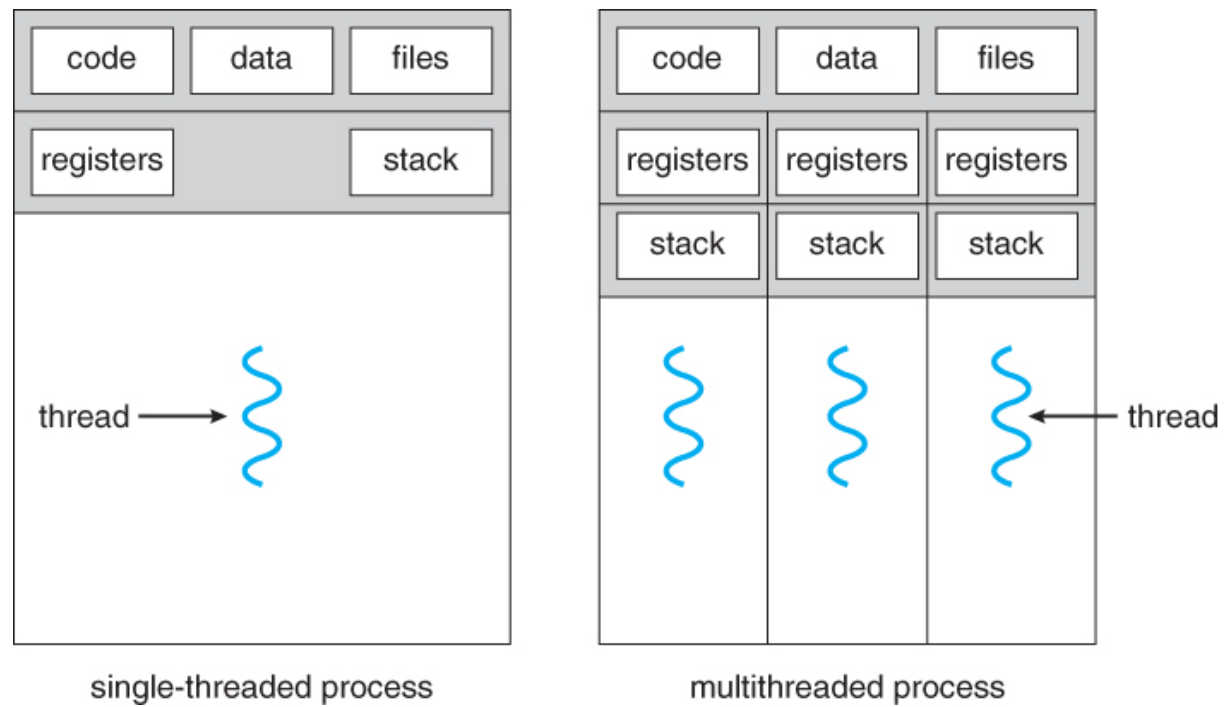


Ping Test (Ethernet inside CSAIL)

- csail.mit.edu
 - 0.7 ms
- mit.edu
 - 14.0 ms
- harvard.edu
 - 7.0 ms
- berkeley.edu
 - 65.1 ms
- tsinghua.edu
 - 229.5 ms



Threads vs Processes



https://www.cs.uic.edu/~jbell/CourseNotes/OperatingSystems/4_Threads.html

Python Threads API

```
import threading
```

```
t = threading.Thread(target=func_name, args=(a1,a2,...))  
t.start()    #start thread running – main thread continues  
t.join()    #wait for thread to finish
```

```
lock = threading.Lock()    #create a lock object  
lock.acquire() #acquire the lock; block if another thread has it  
lock.release() #release the lock
```

Problem: Python Global Interpreter Lock (GIL)
Only one thread can be executing python code at once

Python Multiprocessing API

```
import multiprocessing
```

```
p = multiprocessing.Process(target=func_name, args=(a1,a2,...))
```

```
p.start()    #start thread running – main thread continues
```

```
p.join()    #wait for thread to finish
```

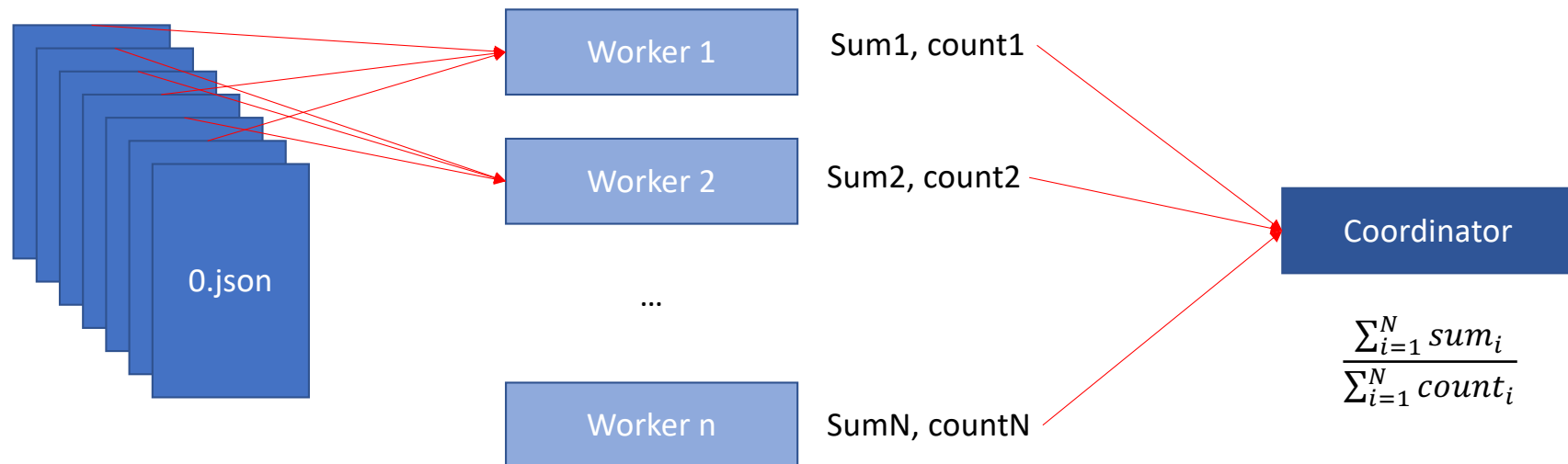
```
lock = multiprocessing.Lock()    #create a lock object
```

```
lock.acquire() #acquire the lock; block if another thread has it
```

```
lock.release() #release the lock
```

Parallel Aggregation

Task: compute average age across all people



```
{"age": 30, "name": ["Michal", "Sharpe"],  
"occupation": "Archivist", "telephone":  
"285.290.9033", "address": {"address":  
"458 Girard Plantation", "city":  
"Wentzville"}, "credit-card": {"number":  
"5384 0033 6904 0042", "expiration-date":  
"06/23"}}
```

Parallel Aggregation Implementation

- Use multiprocessing, not threading
- Main thread creates a work queue

```
q = multiprocessing.Queue()
```

- Puts work on it, as pointers to files

```
q.put(file1); q.put(file2)
```

- Starts threads, passing them the work queue, as well as a result queue
- Threads pull from queue in a loop:

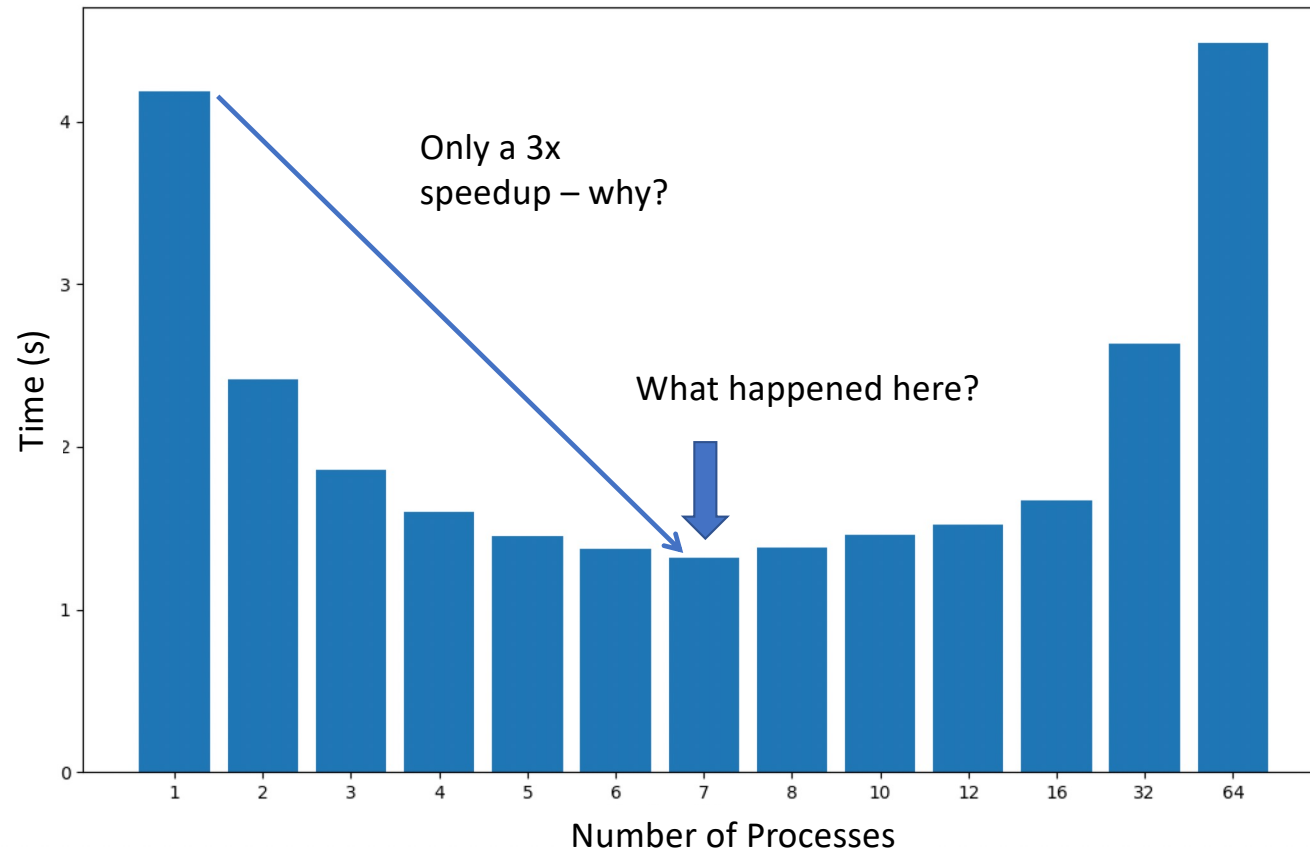
```
while True:  
    f = q.get(block=False)  
    process(f)
```

- Threads compute running sum and average
- Once complete, threads put their running sum and average on the result queue:

```
out_q.put((age_sum, age_cnt))
```

- Main thread blocks on result queue to read a result from each worker:

```
for p in procs:  
    (p_sum, p_count) = out_q.get()
```



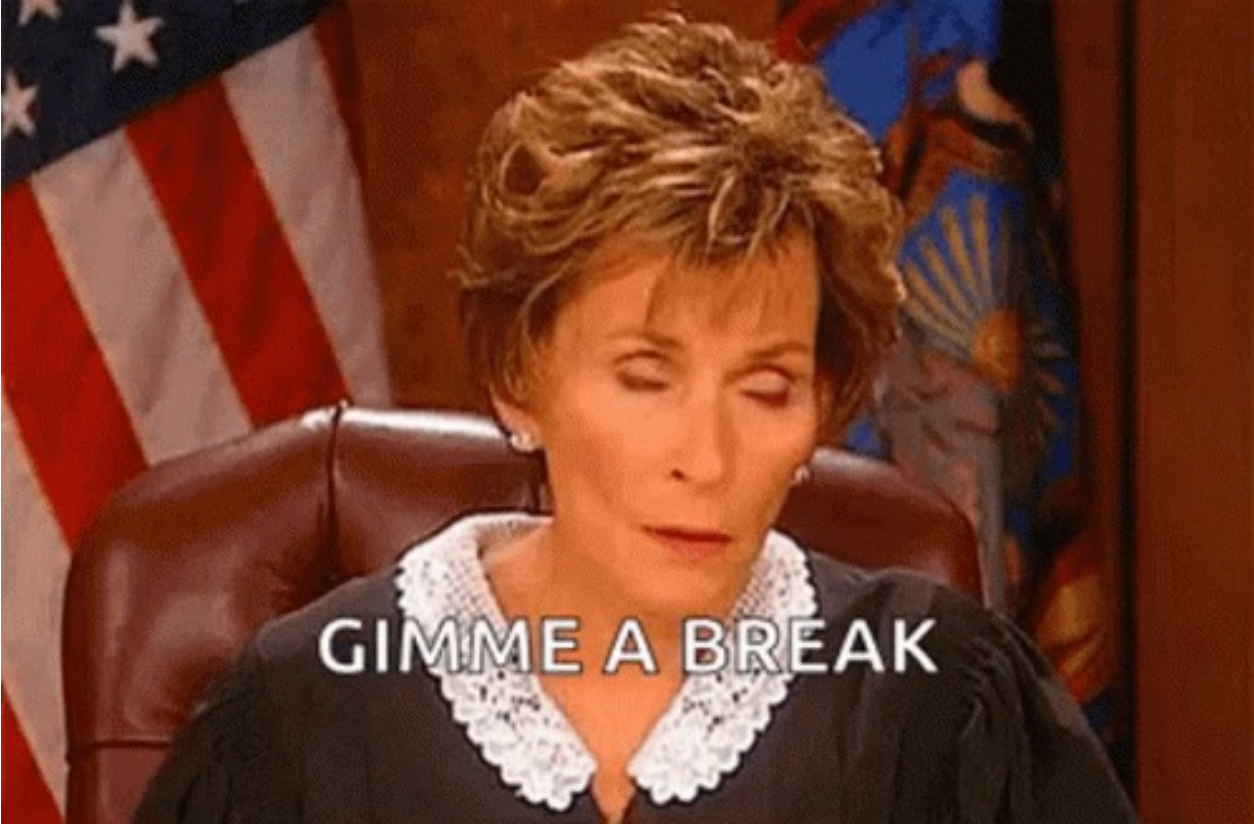
Clicker

Why didn't this program speed up beyond 8 processes? Choose all that apply

- a) Not enough memory
- b) Not enough processors
- c) Startup overheads of launching processes
- d) Too much coordination between processes

<https://clicker.mit.edu/6.S079>

Break



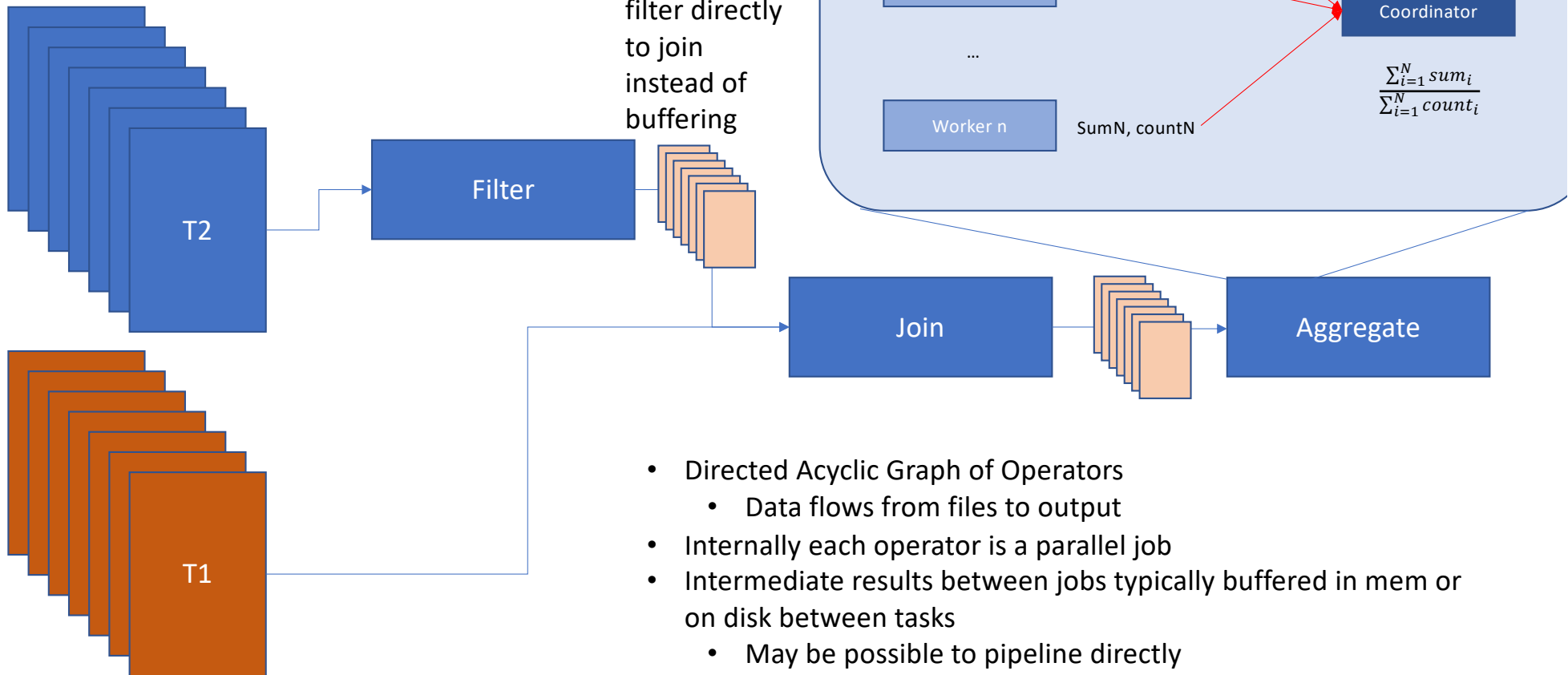
Parallelism Approach

Split given data set split into N partitions

Use M processors to process this data in parallel

We will need to come up with parallel implementations of common operators

Parallel Dataflow Example



Parallel Dataflow Operations

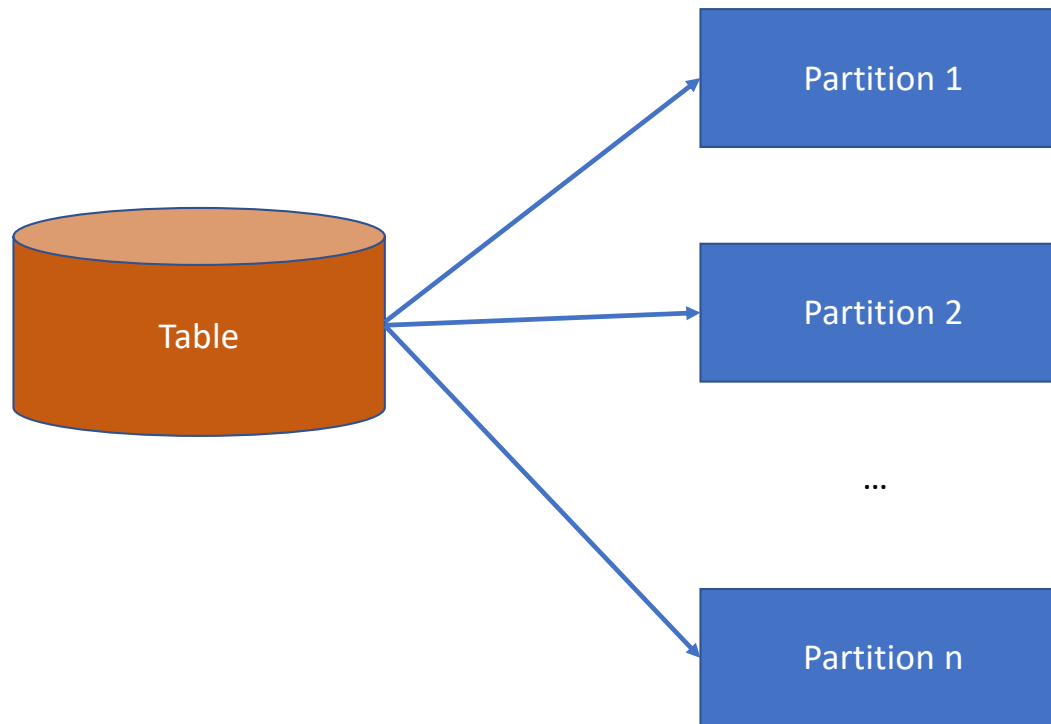
- Filter
- Project
- Element-wise or row-wise transform
- Join
 - Repartition vs broadcast
- Aggregate
- Sort
- Train an ML model
- Arbitrary python "UDF"s

Which of these are easy to parallelize?

Partitioning Strategies

- Random / Round Robin
 - Evenly distributes data (no skew)
 - Requires us to repartition for joins
- Range partitioning
 - Allows us to perform joins/merges without repartitioning, when tables are partitioned on join attributes
 - Subject to skew
- Hash partitioning
 - Allows us to perform joins/merges without repartitioning, when tables are partitioned on join attributes
 - Only subject to skew when there are many duplicate values

Round Robin Partitioning



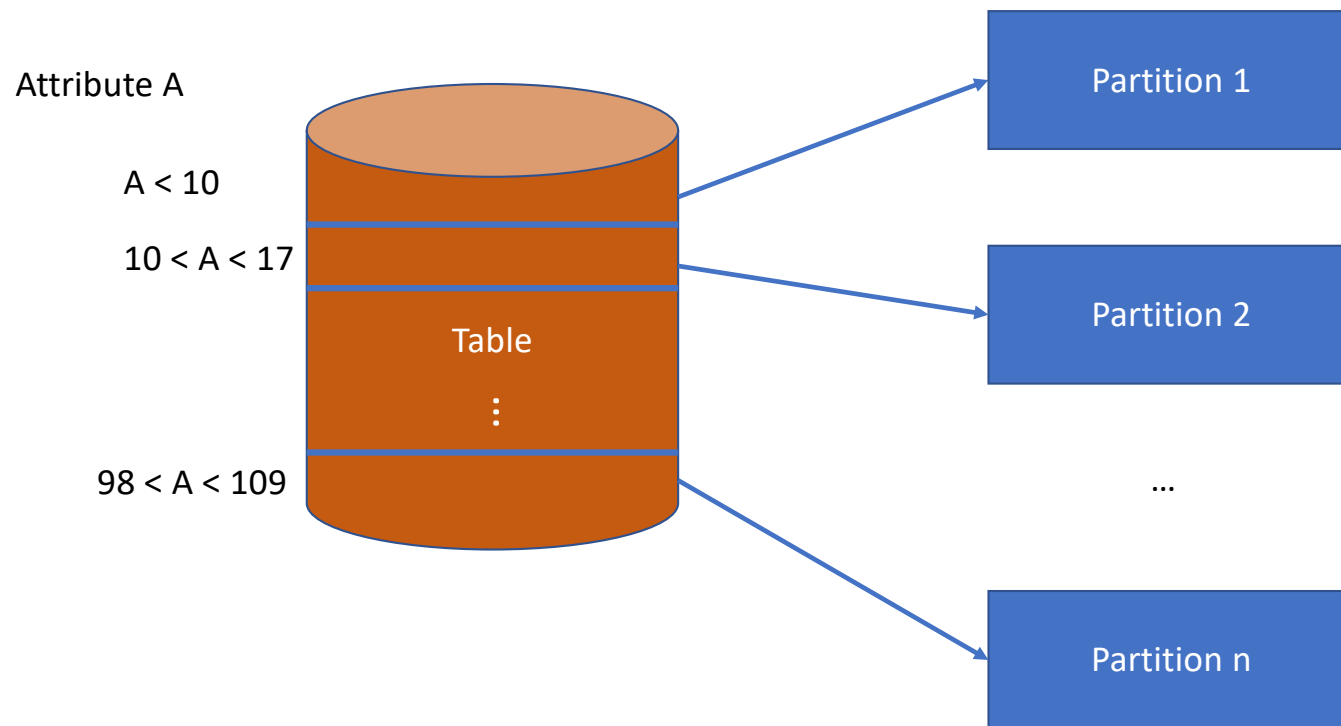
Advantages:

Each partition has the same number of records

Disadvantage:

No ability to push down predicates to filter out some partitions

Range Partitioning



Advantages:

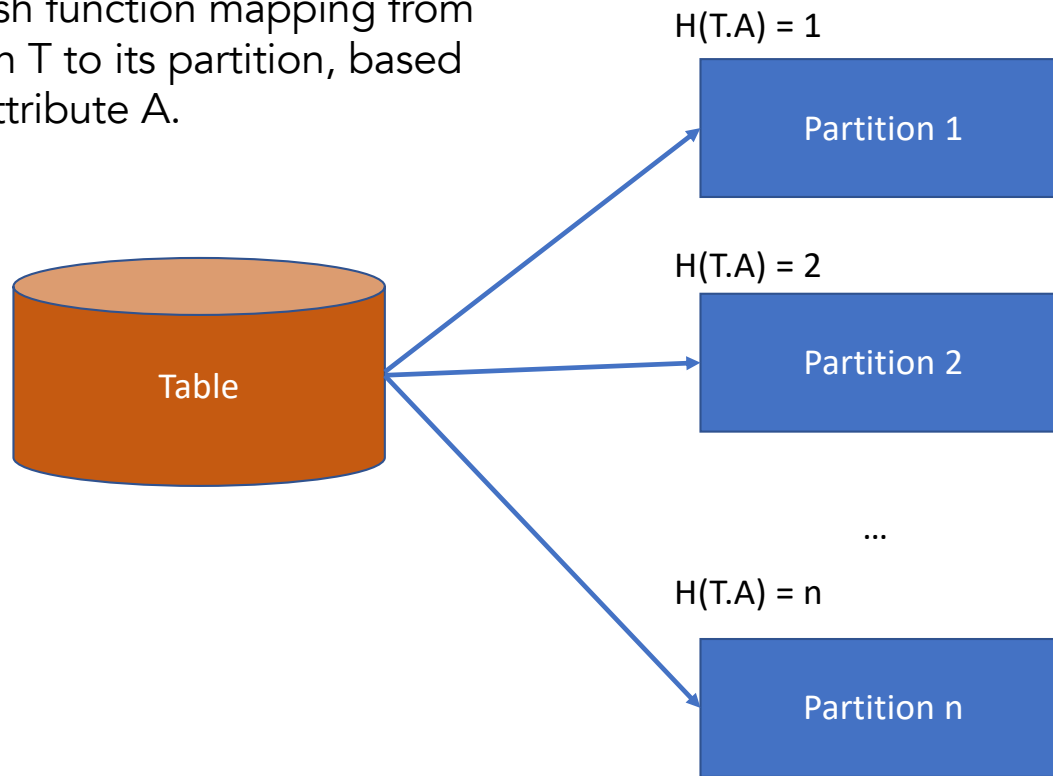
Easy to push down predicates (on partitioning attribute)

Disadvantage:

Difficult to ensure equal sized partitions, particularly in the face of inserts and skewed data

Hash Partitioning

$H(T.A)$ is a hash function mapping from each record in T to its partition, based on value of attribute A .



Advantages:

Each partition has about the same number of records, unless one value is very frequent

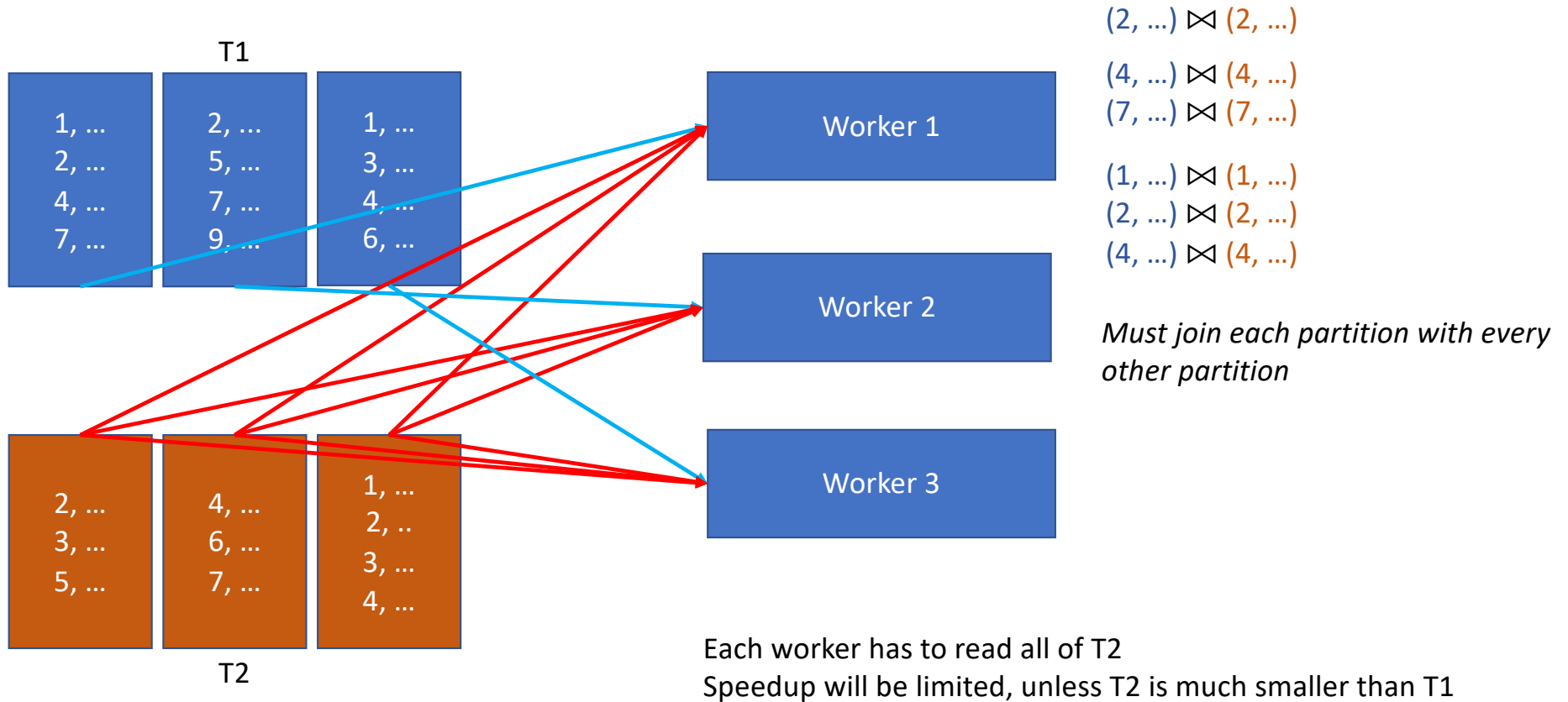
Possible to push down equality predicates on partitioning attribute

Disadvantages:

Can't push down range predicates

Parallel Join – Random Partitioning Naïve Algo

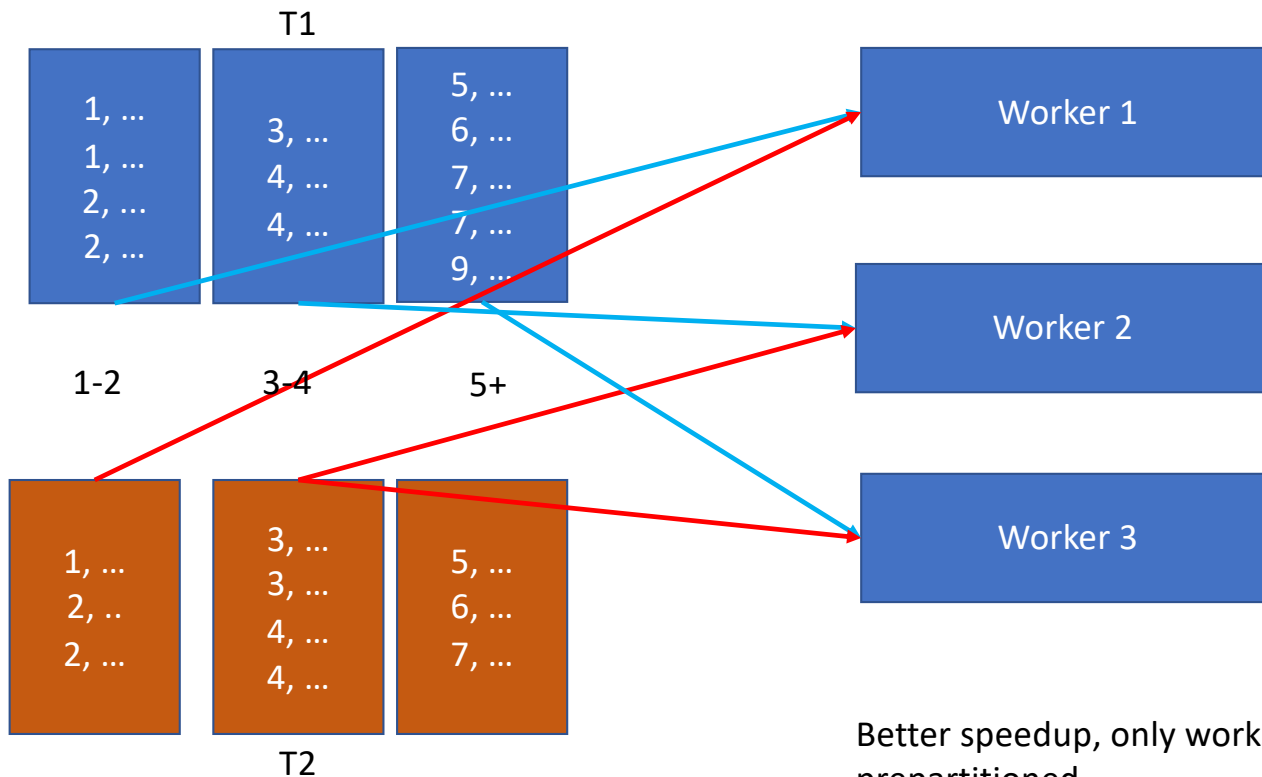
(1, ...) indicates value of join attribute



Parallel Join – Prepartitioned

(1, ...) indicates value of join attribute

Only need to join partitions that match



(1, ...) ⋈ (1, ...)
(1, ...) ⋈ (1, ...)
(2, ...) ⋈ (2, ...)
(2, ...) ⋈ (2, ...)
(2, ...) ⋈ (2, ...)
(2, ...) ⋈ (2, ...)
(2, ...) ⋈ (2, ...)

This is what our Postgres example showed

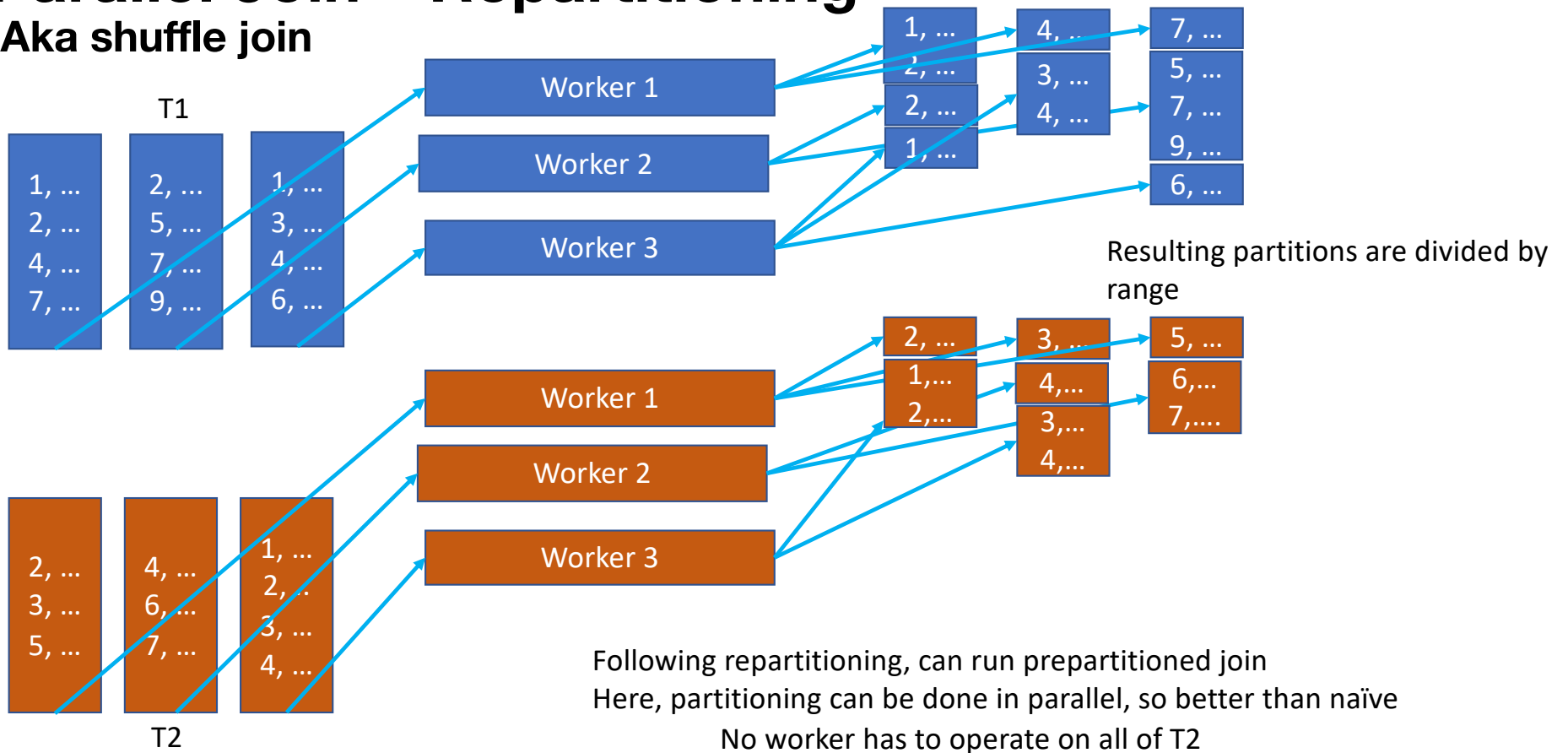
Better speedup, only works if data is properly prepartitioned

Should be 3x faster than single node join

Skew problem (hashing may help)

Parallel Join – Repartitioning

Aka shuffle join



Dask

<https://dask.org>

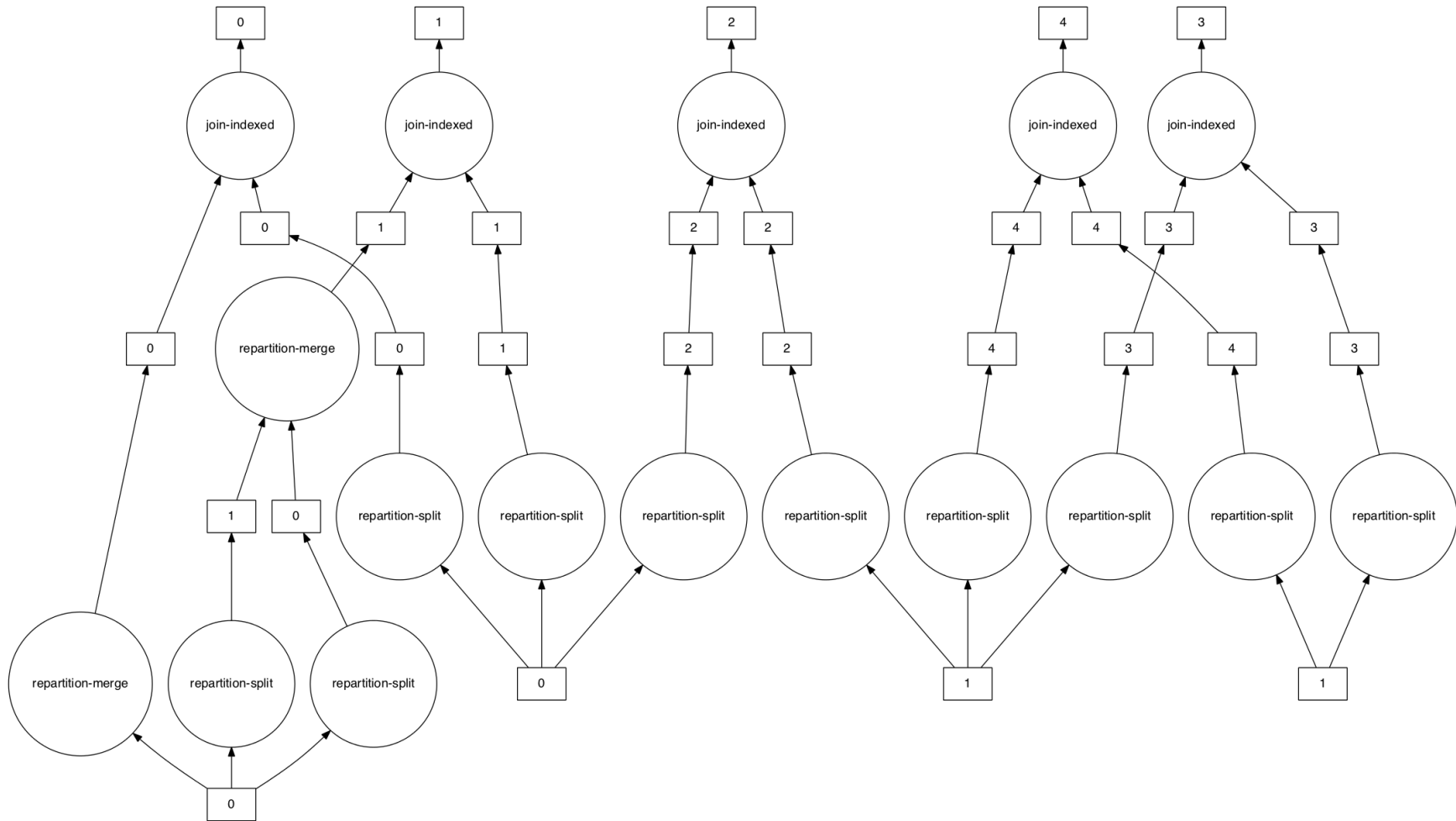


- General purpose python parallel / distributed computation framework
- Includes parallel implementation of Pandas dataframes
- Usually straightforward to translate a pandas program into a parallel implementation
 - Just use `dask.dataframe` instead of `pandas.dataframe`
 - Have to specify a parallel configuration to run on, via `Client()` object
 - Can be a local machine or distributed cluster
- Also has support for other types of parallelism, e.g., `dask.bag` class that allows parallel operation on collections of python objects

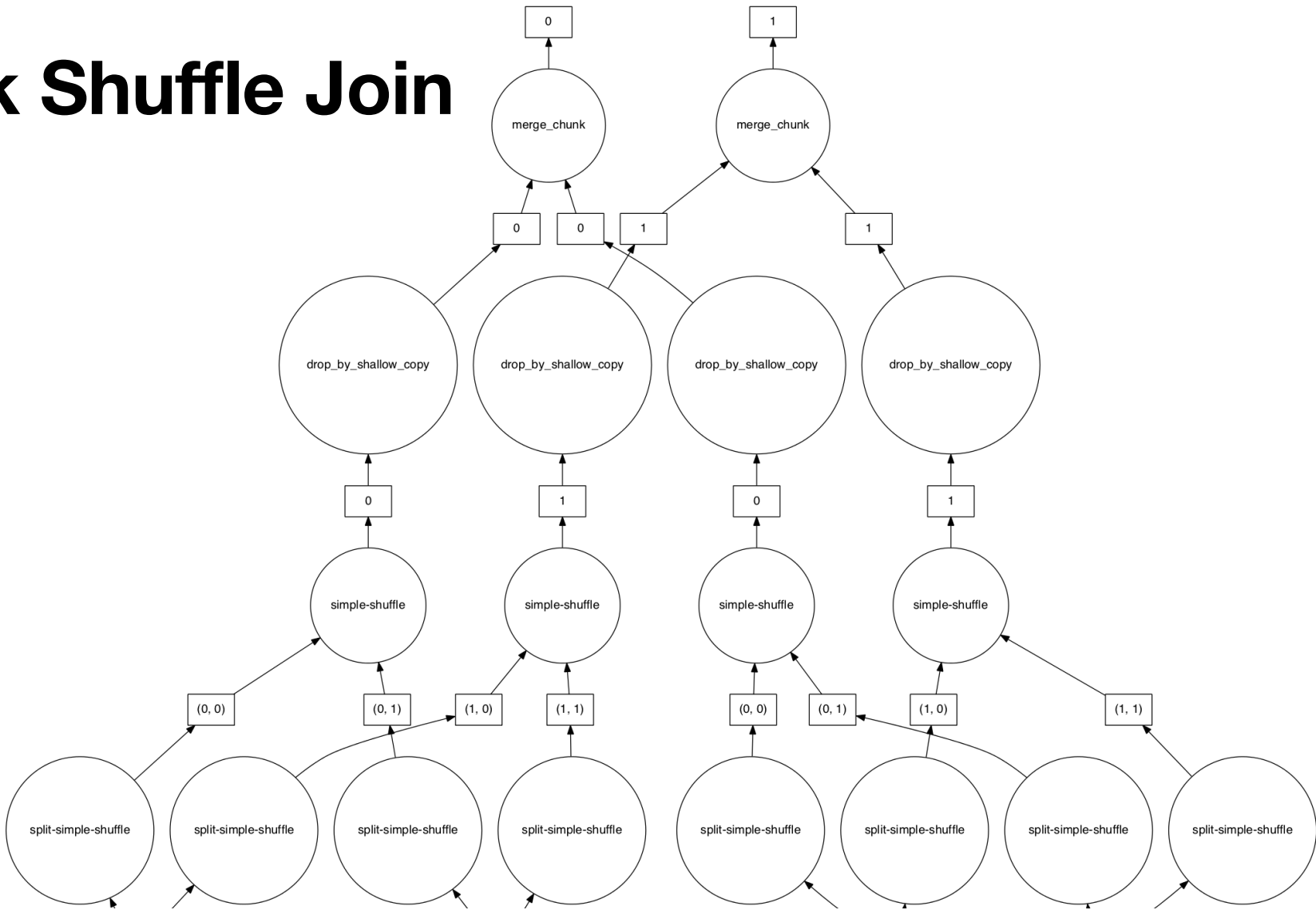
Large Join Demo

- Changing number of nodes
- Changing join algorithm

Dask Partitioned Join



Dask Shuffle Join



Many alternatives

- MapReduce / Hadoop
 - Rewrite you program as collection of parallel map() and reduce() jobs
 - Hard to do, slow()
- Spark
 - Popular library -- similar to dask, more focused on large scale distributed
 - Includes parallel implementations of ML and other operations
 - Difficult to use

Summary

- Parallelism is a good way to improve performance
- Ideal: linear speedup
 - Difficult to achieve in practice
- Some operations can be trivially parallelized with partitioned parallelism, e.g., filters and maps
- Other operations – like joins – are more difficult
- Dask is a popular open-source parallel programming library for Python
 - Next time – you'll get to try it out as a part of Lab 6