

# Programming with Data Bootcamp: Lecture 8

<http://dsg.csail.mit.edu/6.S079/>

Slides courtesy of Sam Madden  
/ Tim Kraska (6.S079)

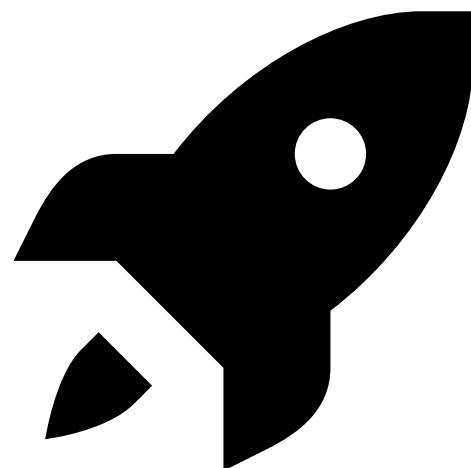


## Key ideas:

Single-node Parallelism

Multi-node Parallelism

- Dask
- Spark
- Ray



# Parallelism Goal

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

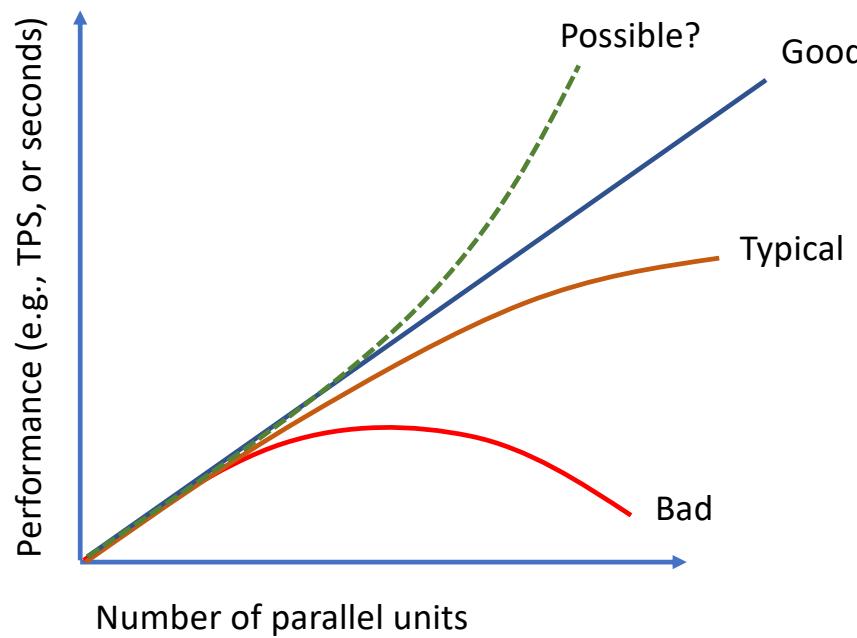
$speed up = \frac{old\ time}{new\ time}$  on same problem, with N times more hardware

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

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

# Speedup Goal

- Linear?

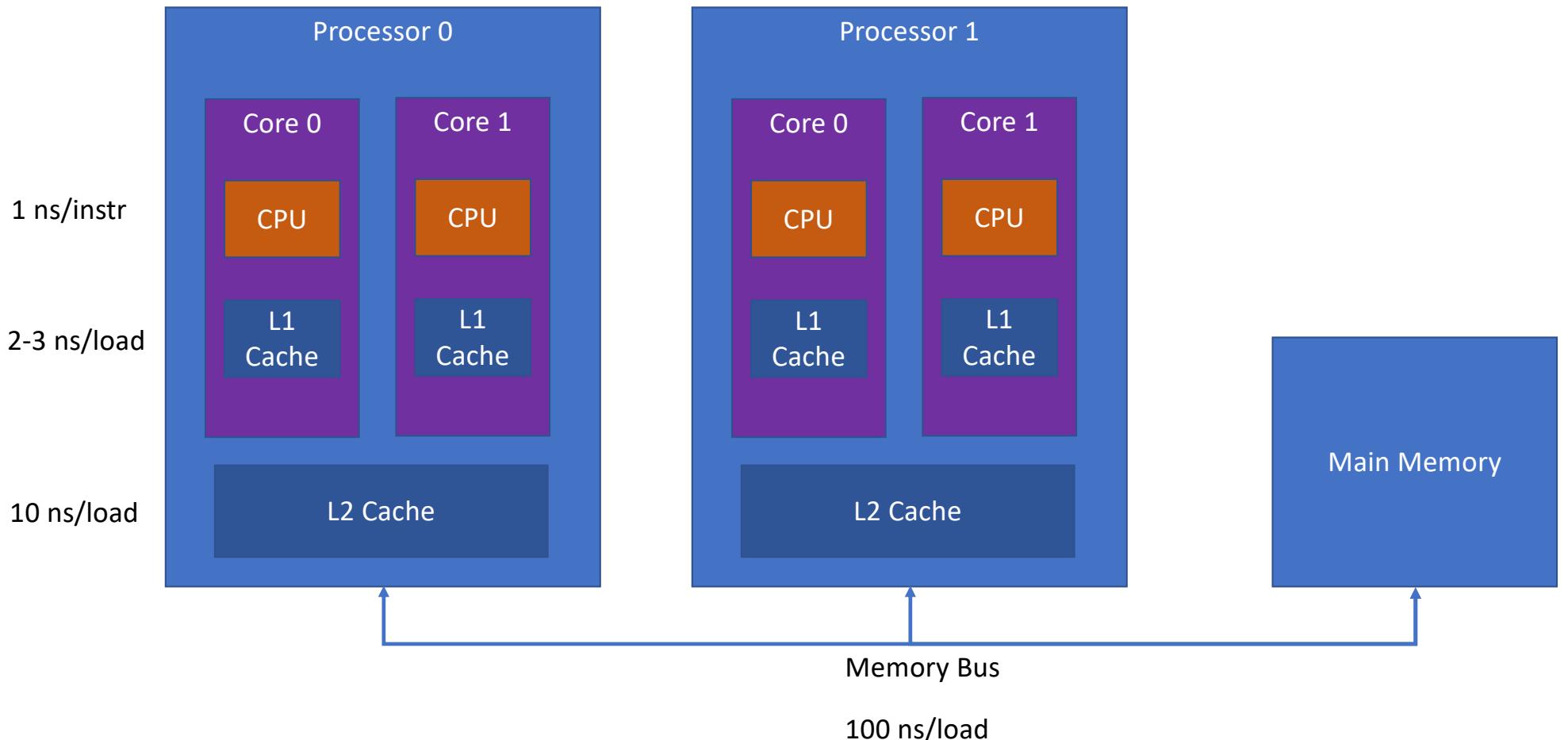


# Barriers to Linear Scaling

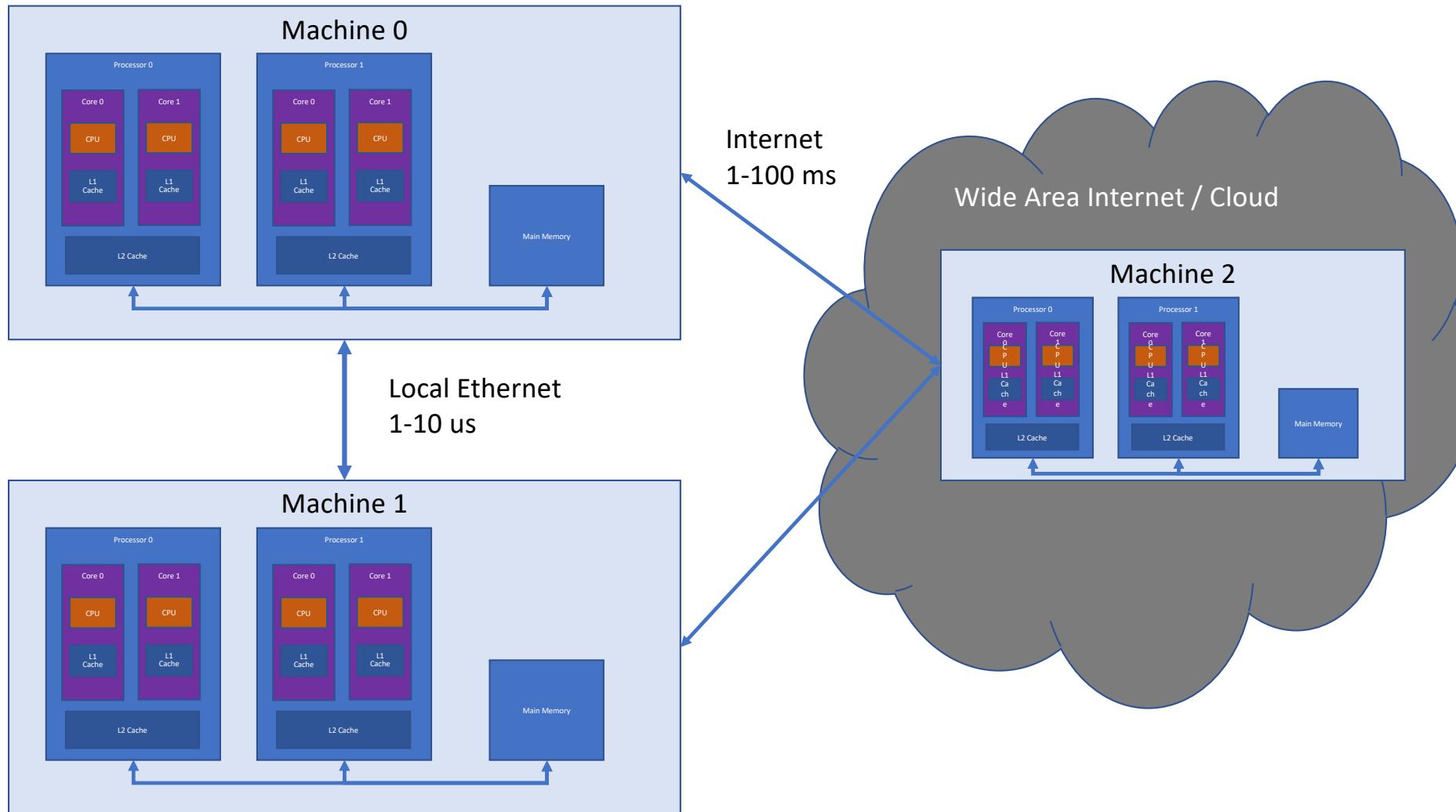
- Startup times
  - e.g., may take time to launch each parallel executor
- Interference
  - processors depend on some shared resource
  - E.g., input or output queue, or other data item
- Skew
  - workload not of equal size on each processor
- *Almost all workloads will stop scaling at some point!*
-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

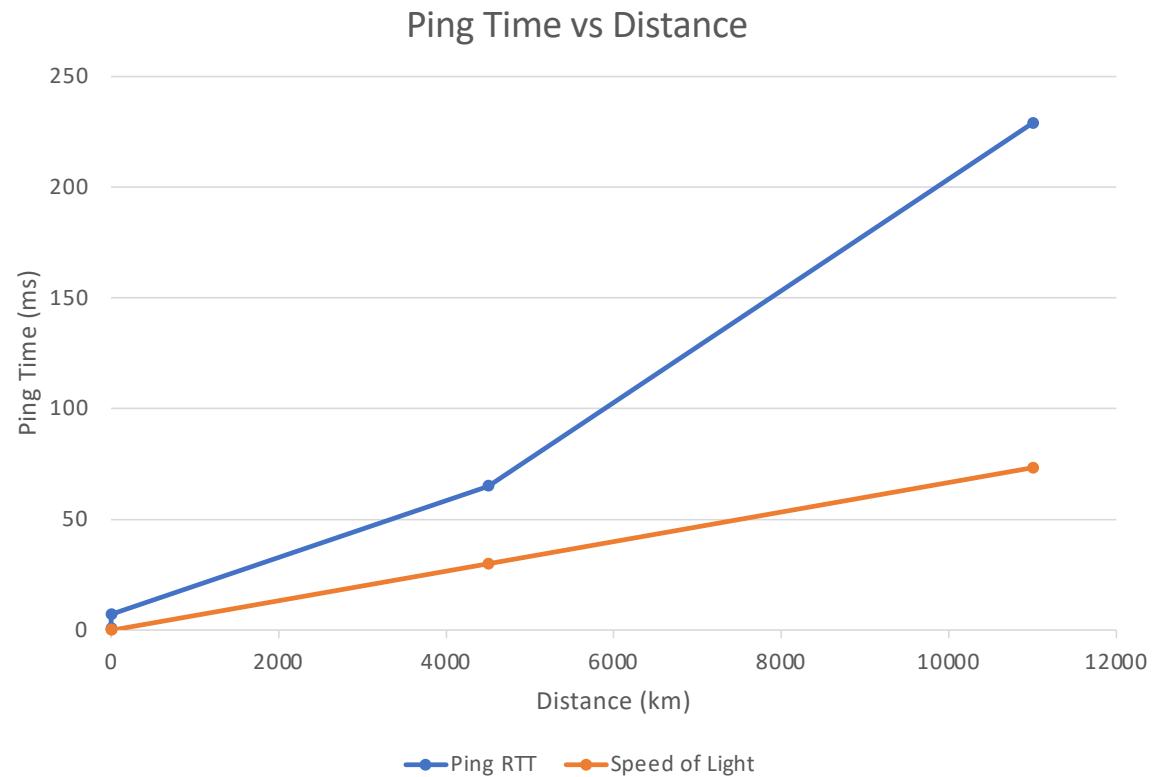


Some machines may have 2 levels of cache per core

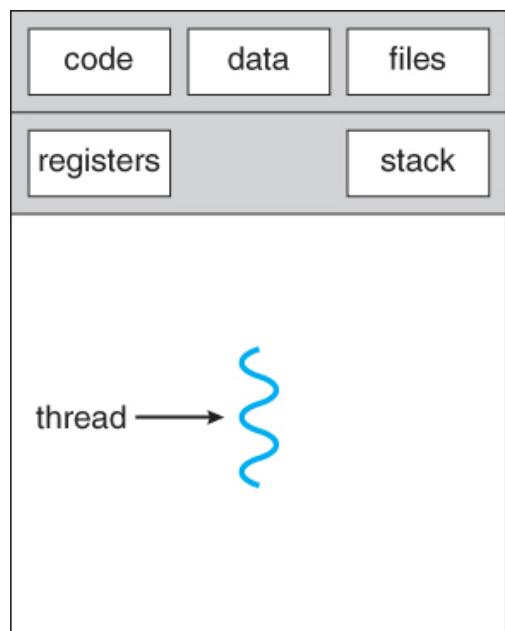


# Ping Test (Ethernet inside CSAIL)

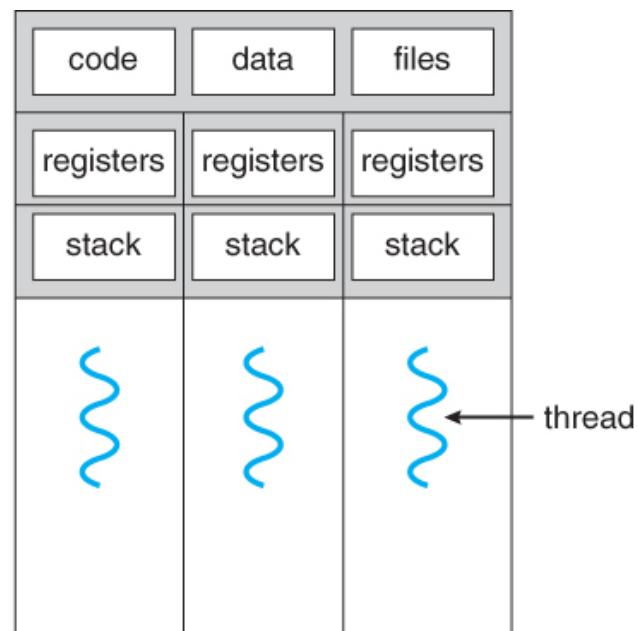
- csail.mit.edu
  - 0.7 ms
- mit.edu
  - 14.0 ms
- harvard.edu
  - 7.0 ms
- berkeley.edu
  - 65.1 ms
- tsinghua.edu
  - 229.5 ms



# Threads vs Processes



single-threaded process



multithreaded process

[https://www.cs.uic.edu/~jbell/CourseNotes/OperatingSystems/4\\_Threads.html](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

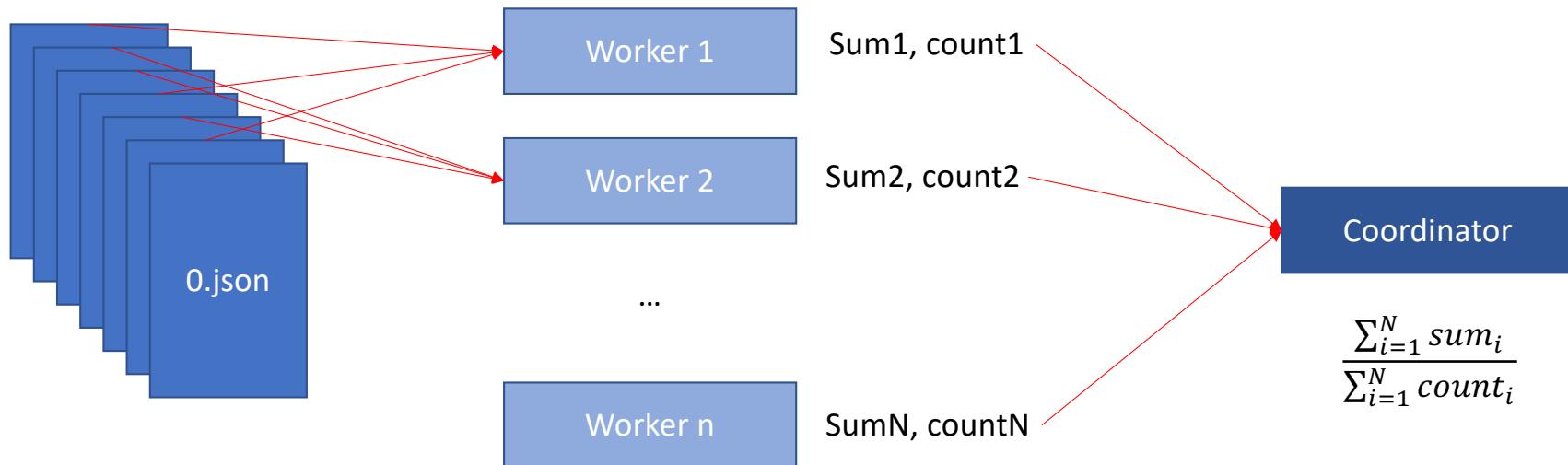
```
import multiprocessing

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:

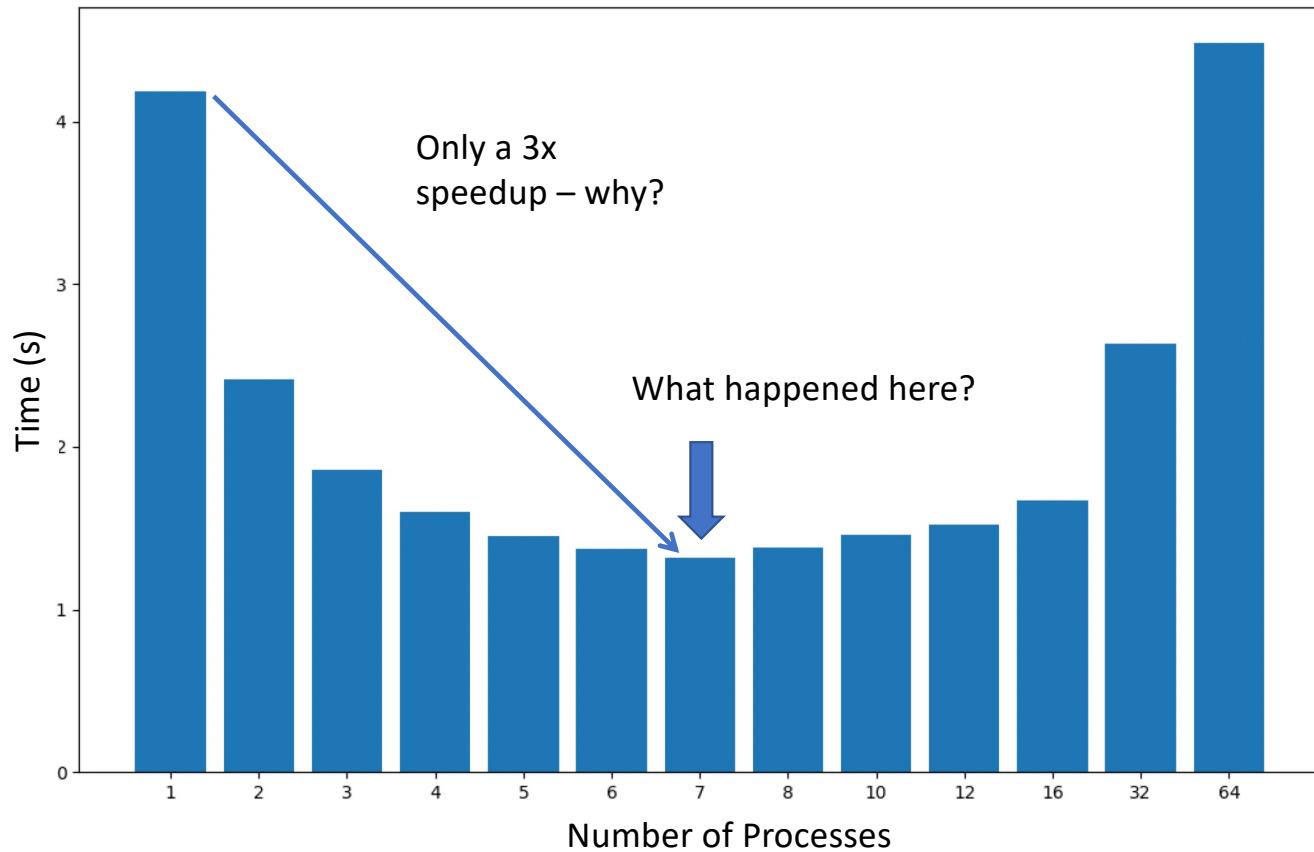
```
while True:  
    f = q.get(block=False)  
    process(f)
```

- 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()
```



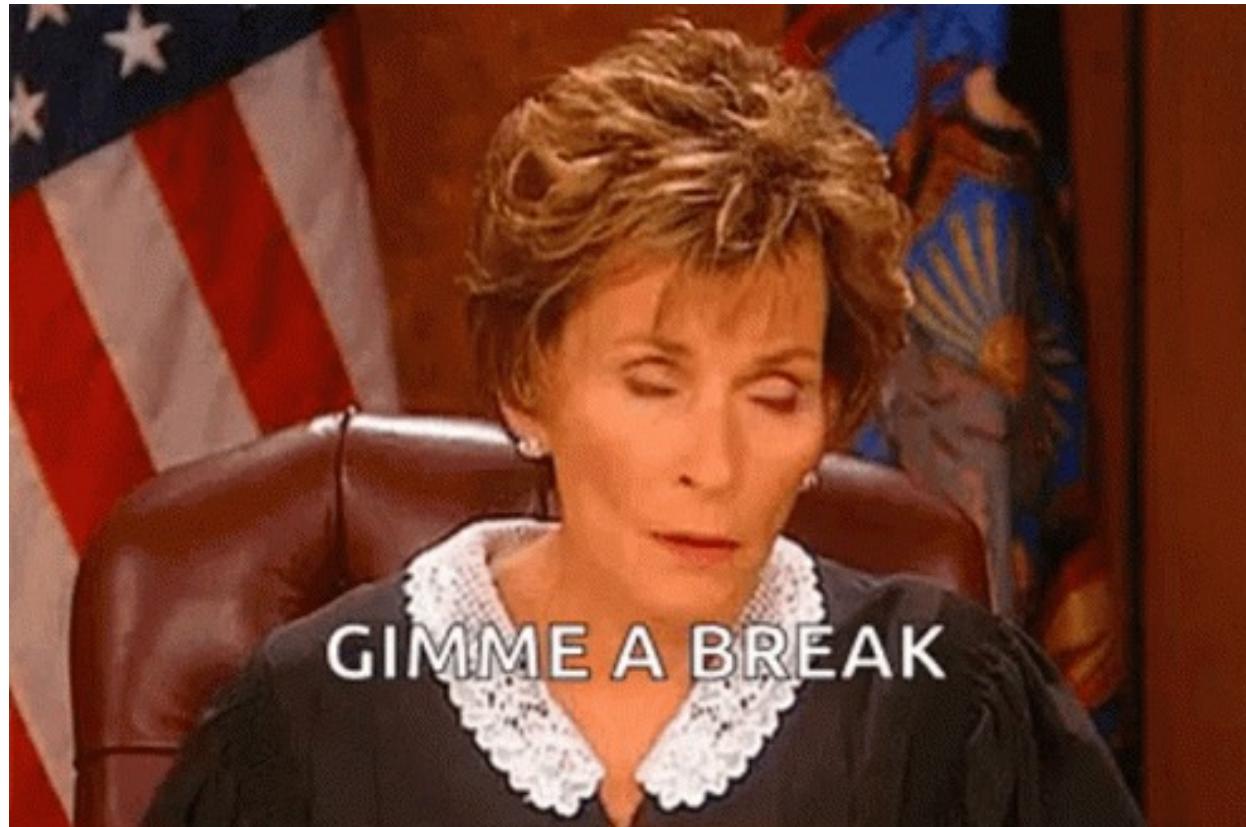
# Clicker

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



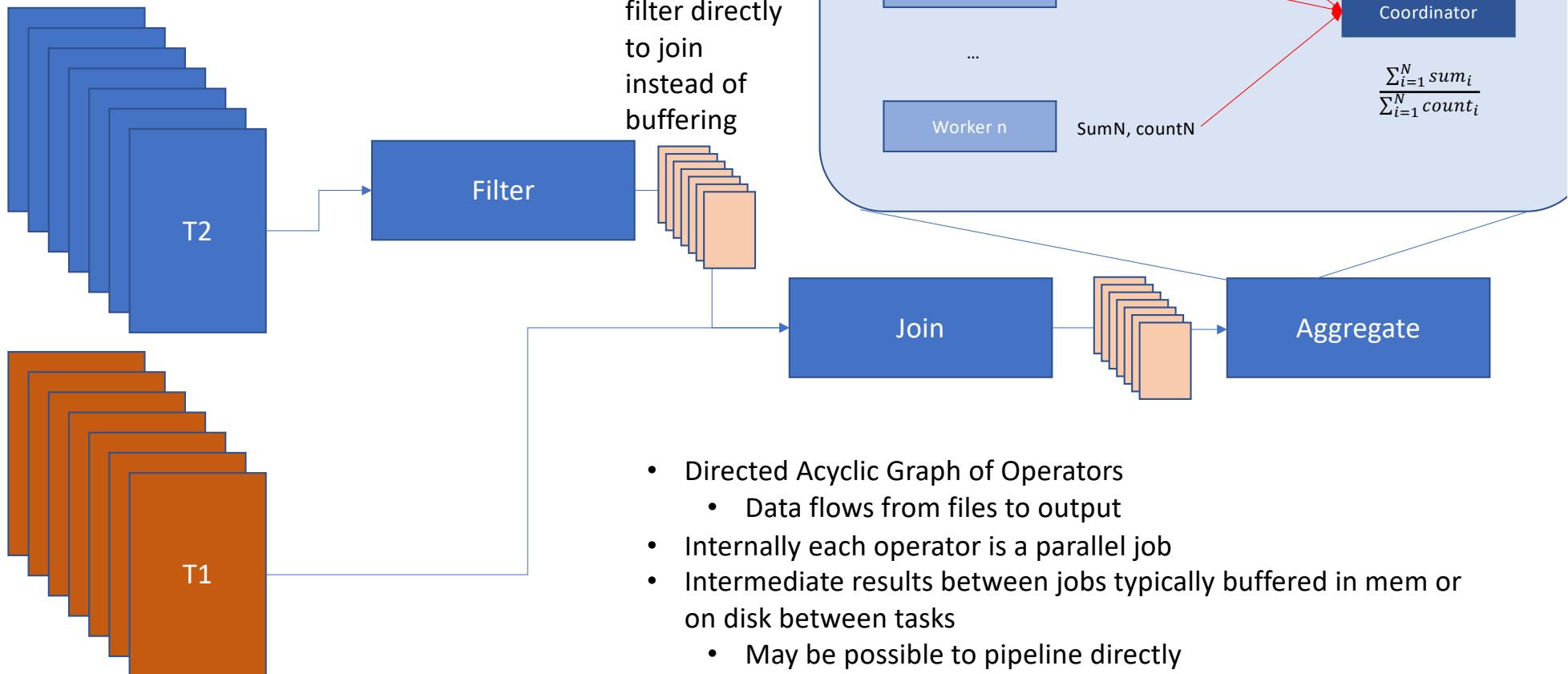
# **Parallelism Approach**

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

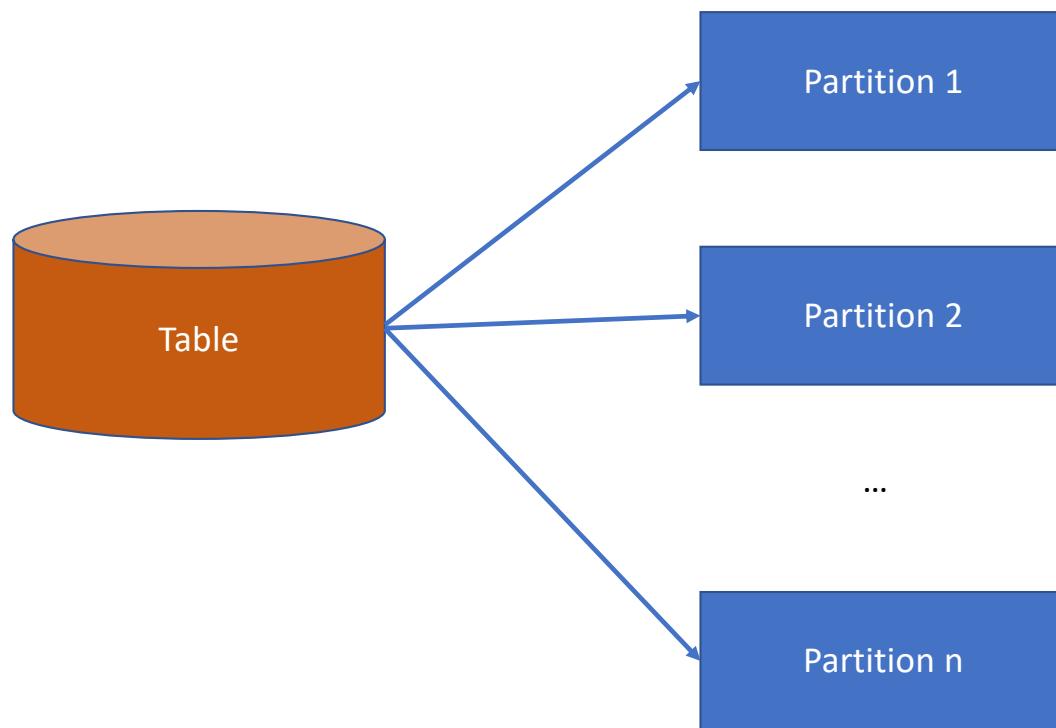
- Filter
- Project
- Element-wise or row-wise transform
- Join
  - Repartition vs broadcast
- Aggregate
- Sort
- Train an ML model
- Arbitrary python "UDF"s

*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



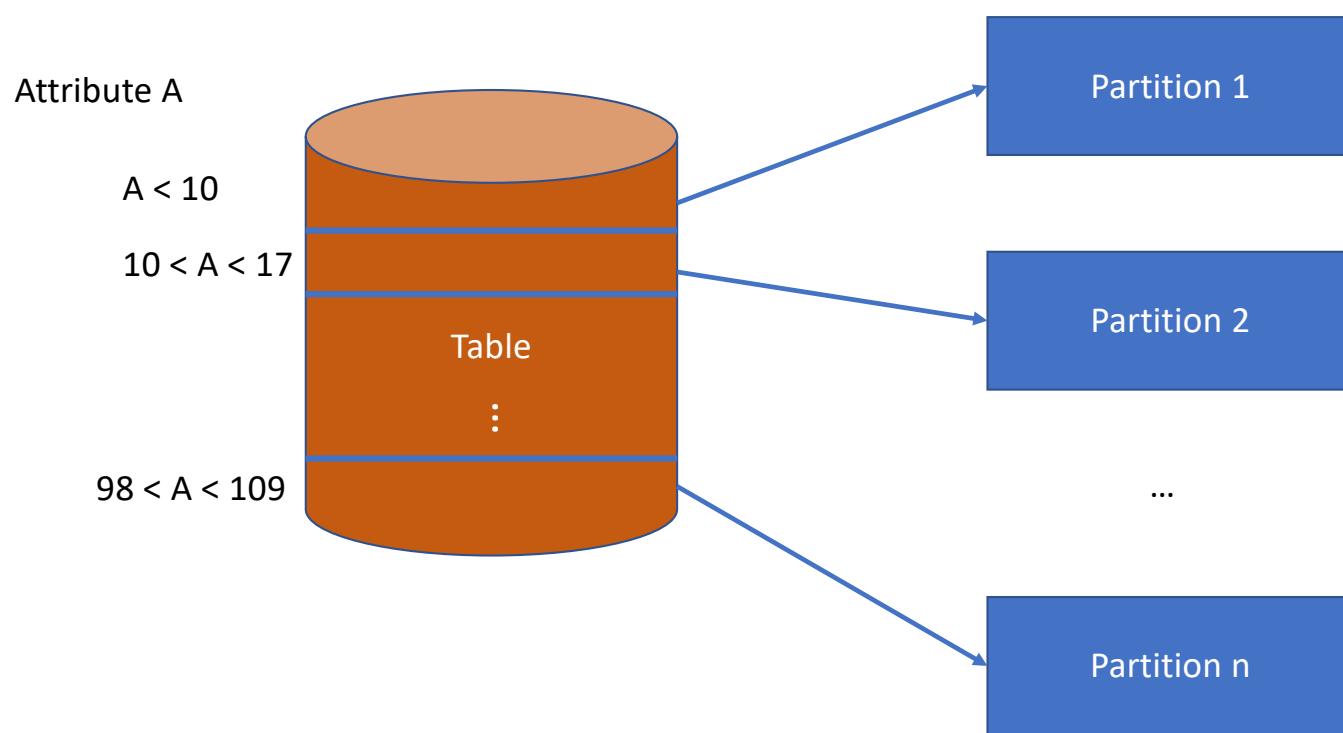
## Advantages:

Each partition has the same number of records

## Disadvantage:

No ability to push down predicates to filter out some partitions

# Range Partitioning



Advantages:

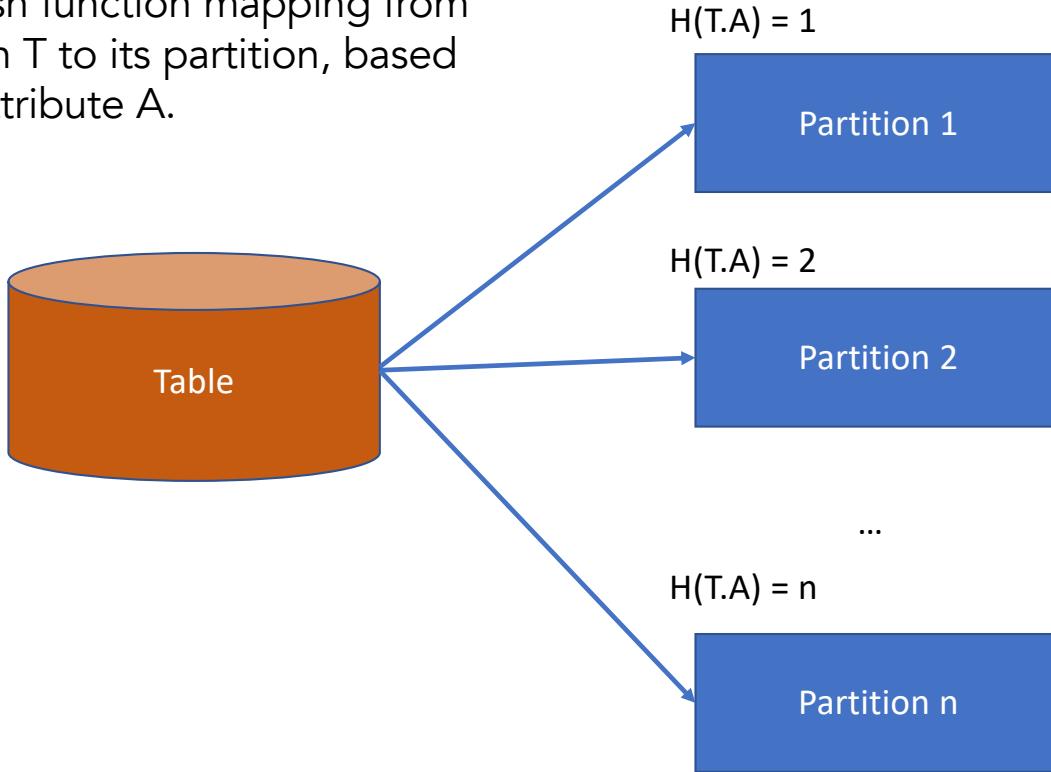
Easy to push down predicates (on partitioning attribute)

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

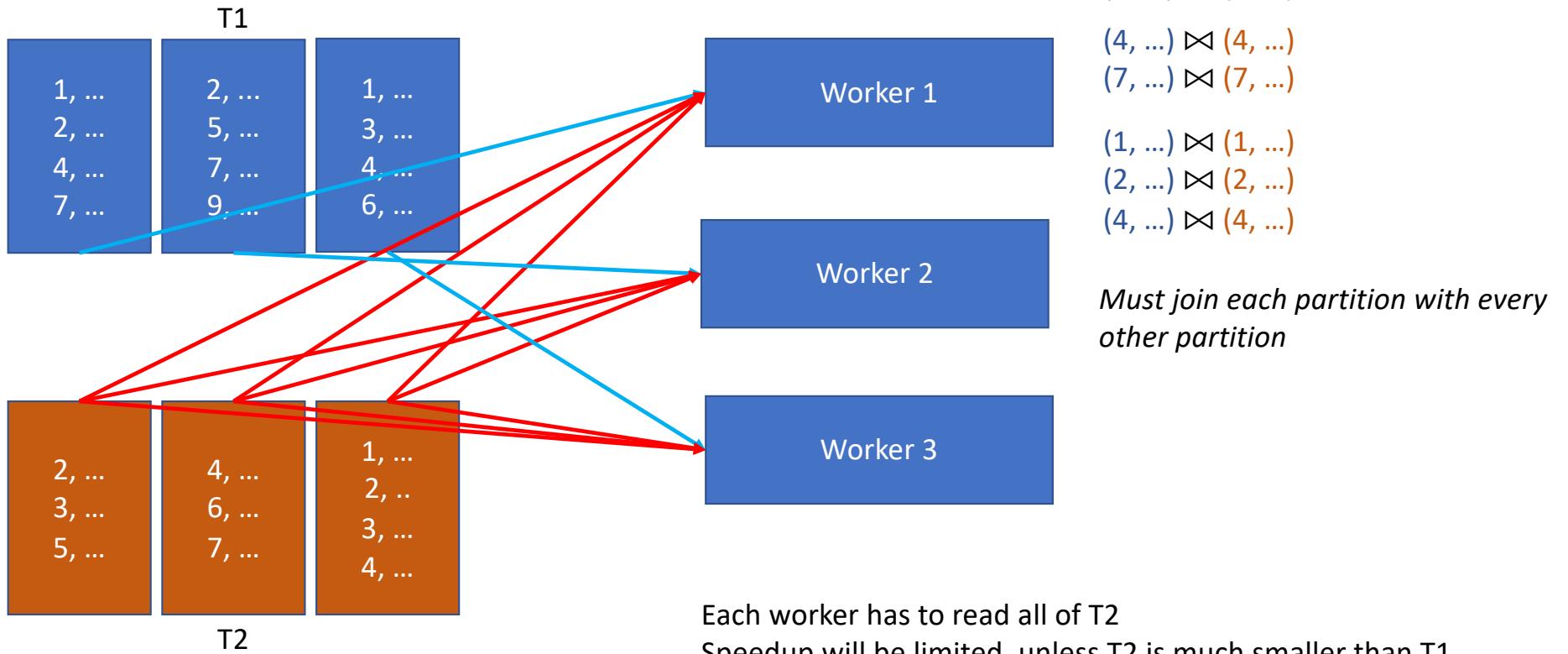
Possible to push down equality predicates on partitioning attribute

## Disadvantages:

Can't push down range predicates

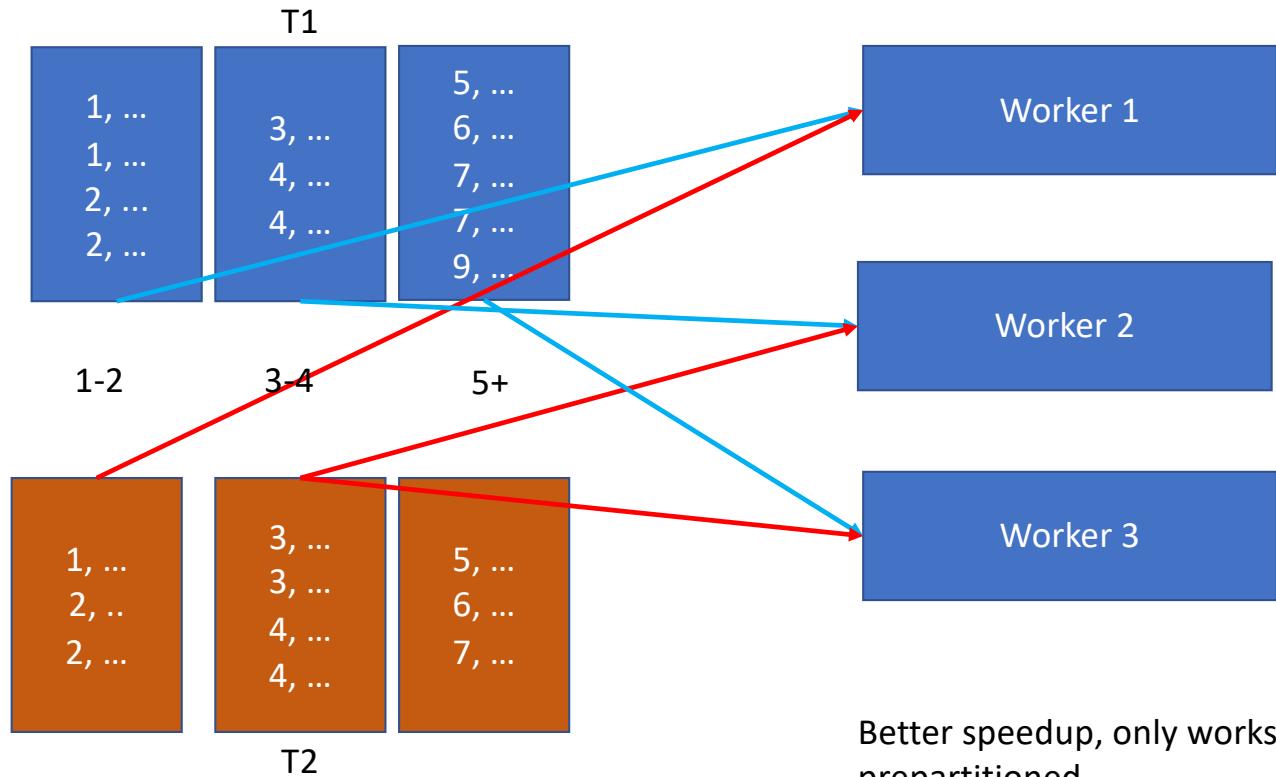
# Parallel Join – Random Partitioning Naïve Algo

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



# Parallel Join – Prepartitioned

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



*Only need to join partitions that match*

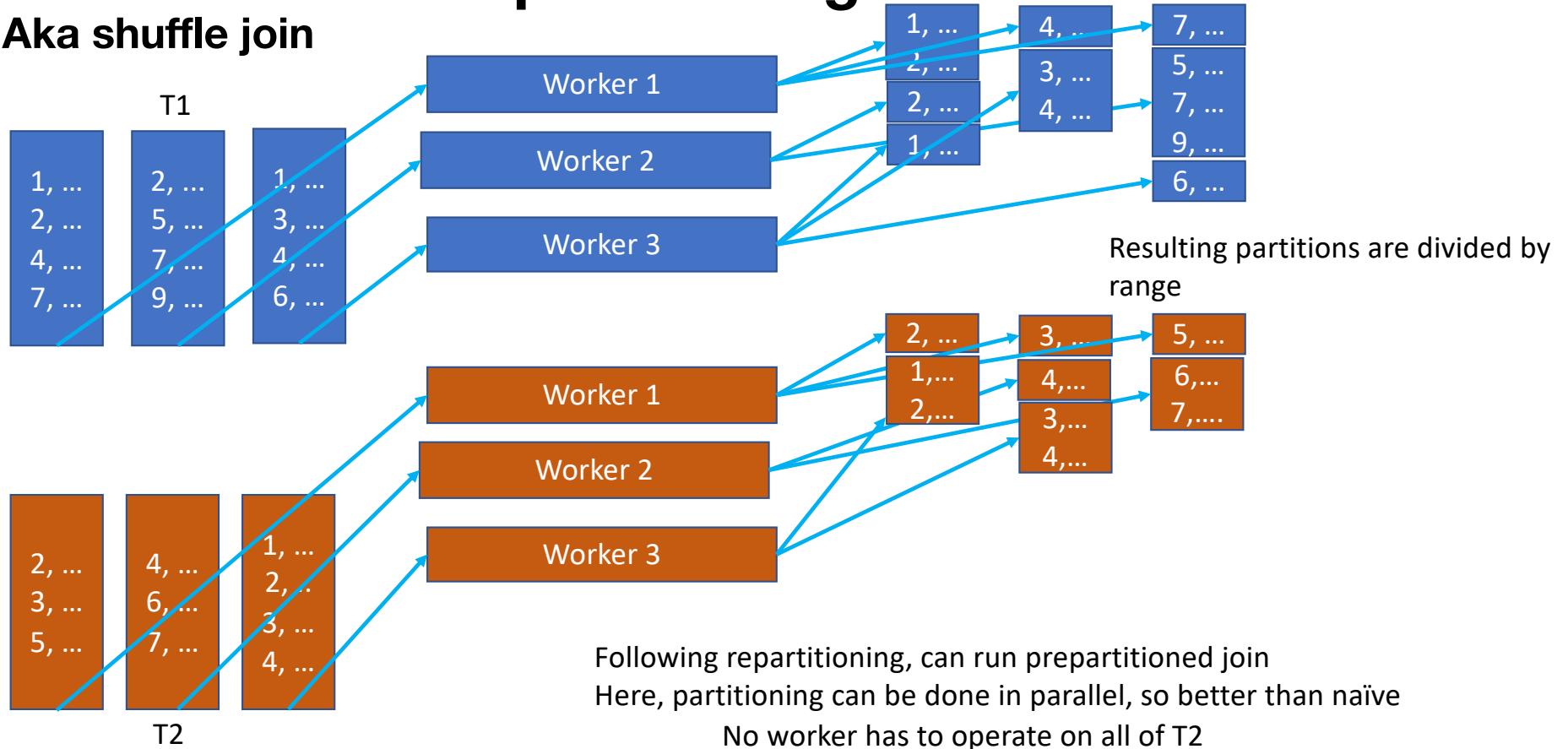
$(1, \dots) \bowtie (1, \dots)$   
 $(1, \dots) \bowtie (1, \dots)$   
 $(2, \dots) \bowtie (2, \dots)$

*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>

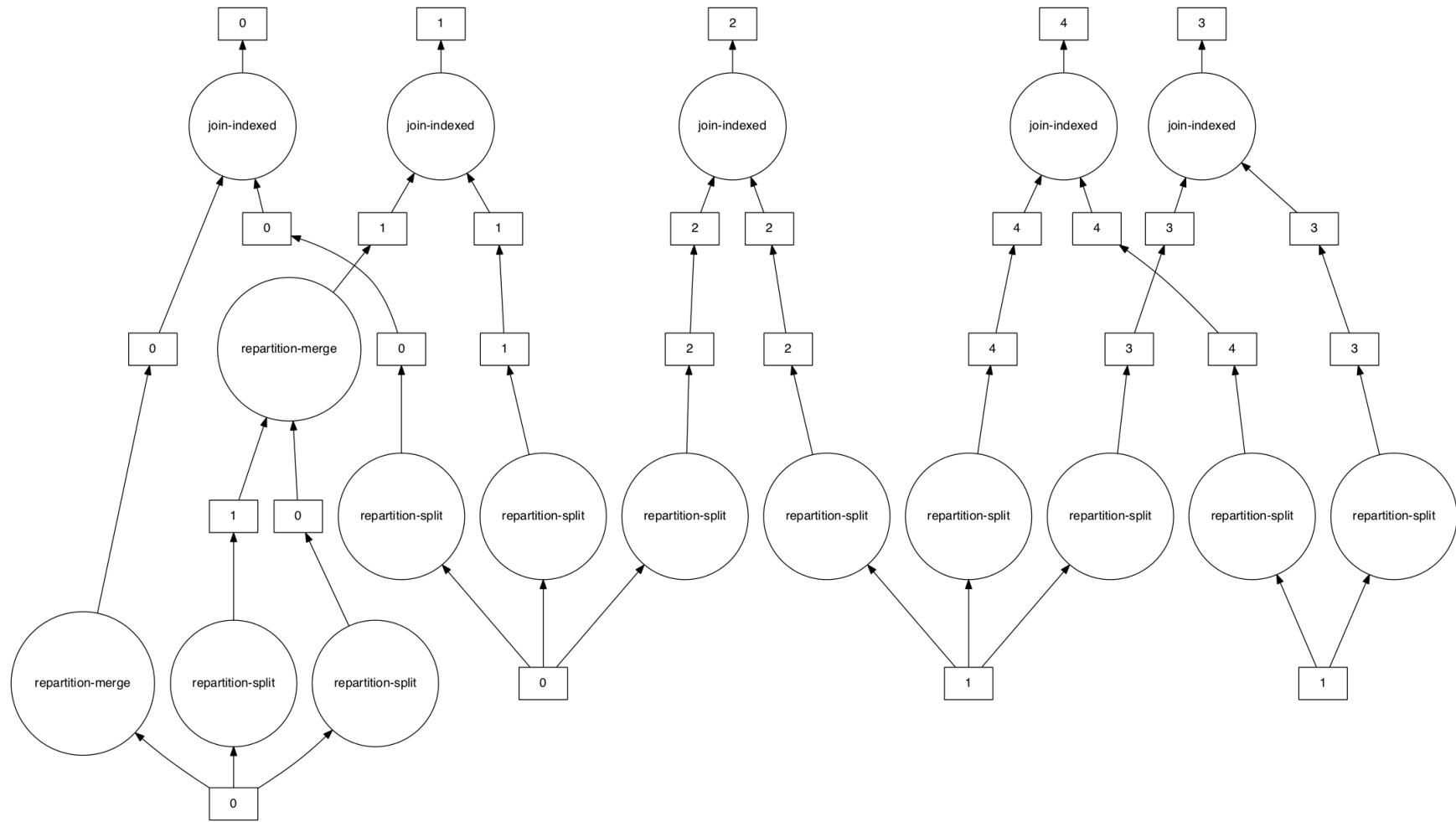


- General purpose python parallel / distributed computation framework
- Includes parallel implementation of Pandas dataframes
- Usually straightforward to translate a pandas program into a parallel implementation
  - Just use `dask.dataframe` instead of `pandas.dataframe`
  - Have to specify a parallel 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

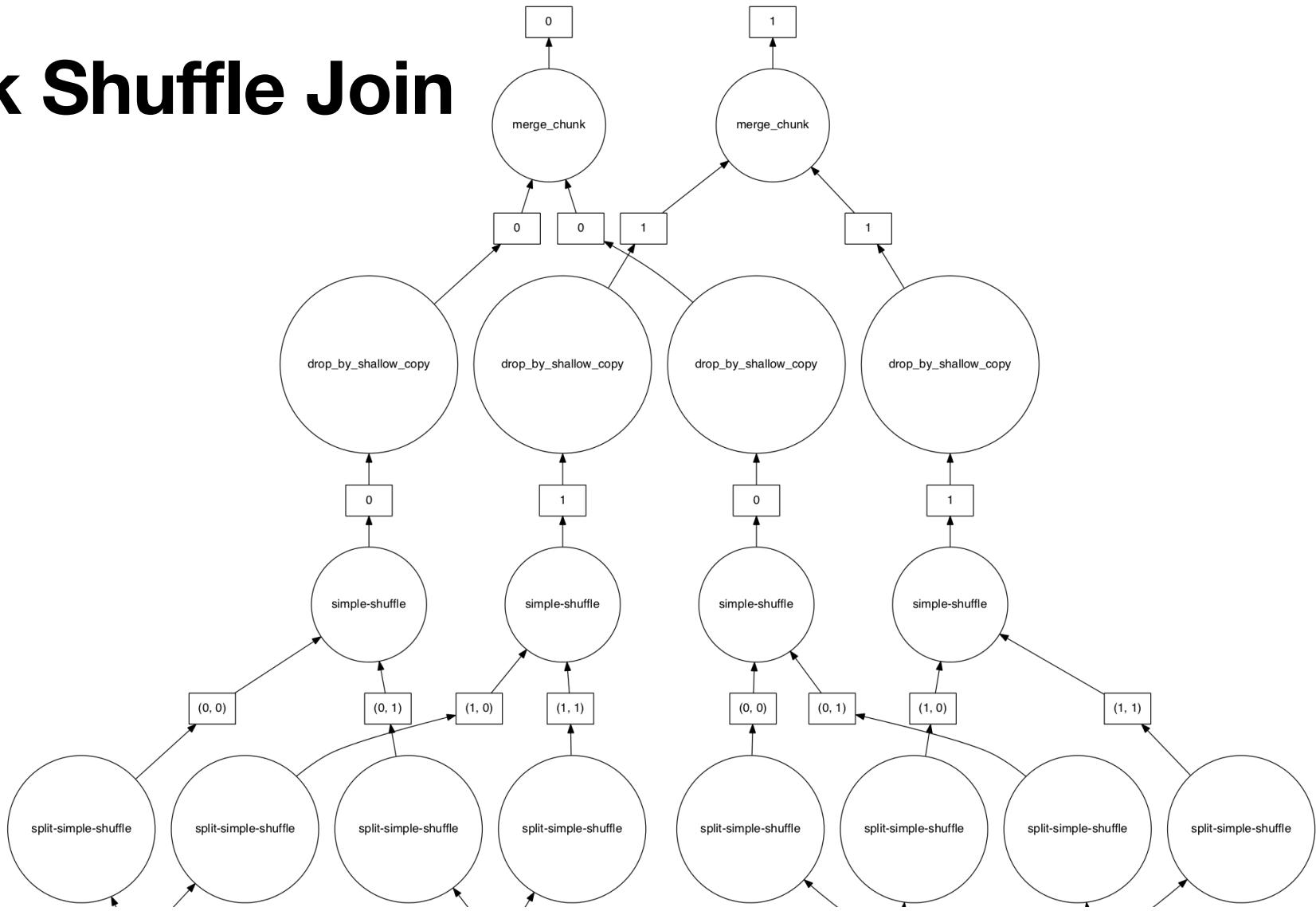
# **Large Join Demo**

- Changing number of nodes
- Changing join algorithm

# Dask Partitioned Join



# Dask Shuffle Join

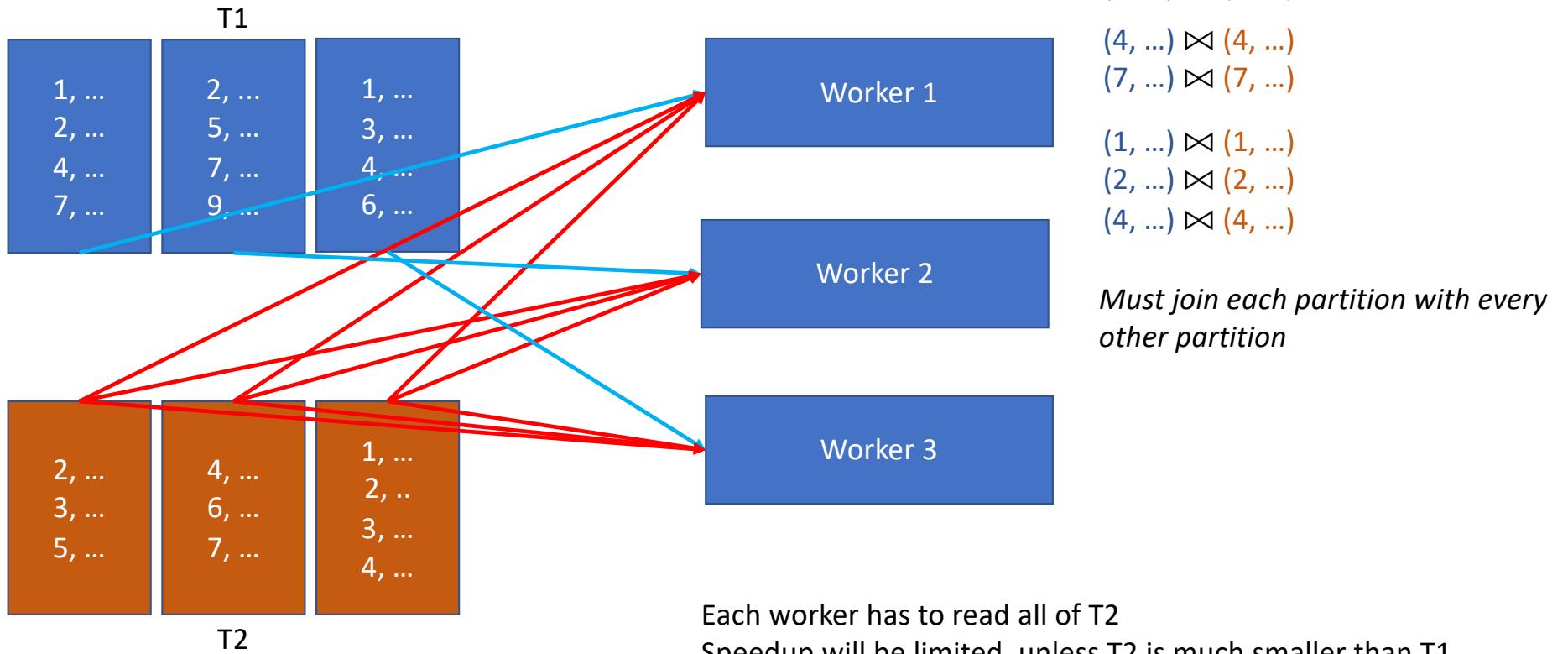


# Many alternatives

- MapReduce / Hadoop
  - Rewrite your 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

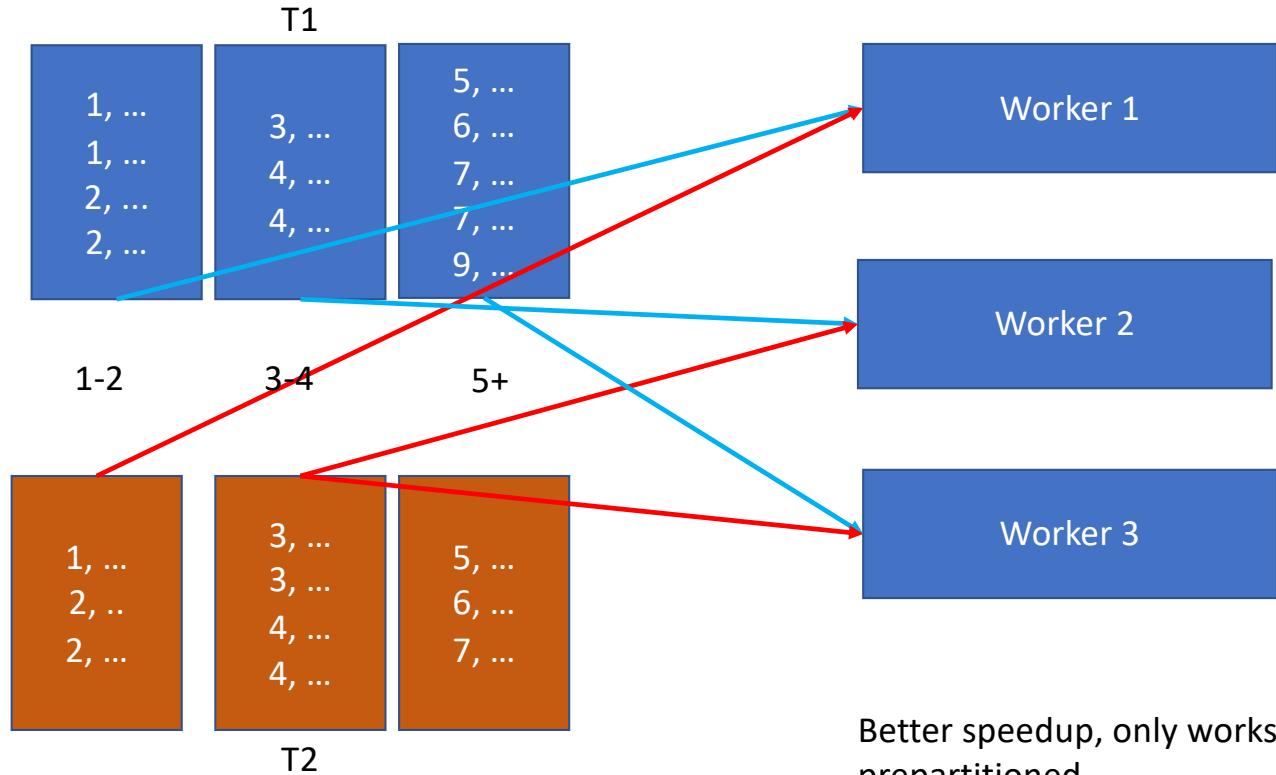
# Parallel Join – Random Partitioning Naïve Algo

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



# Parallel Join – Prepartitioned

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



*Only need to join partitions that match*

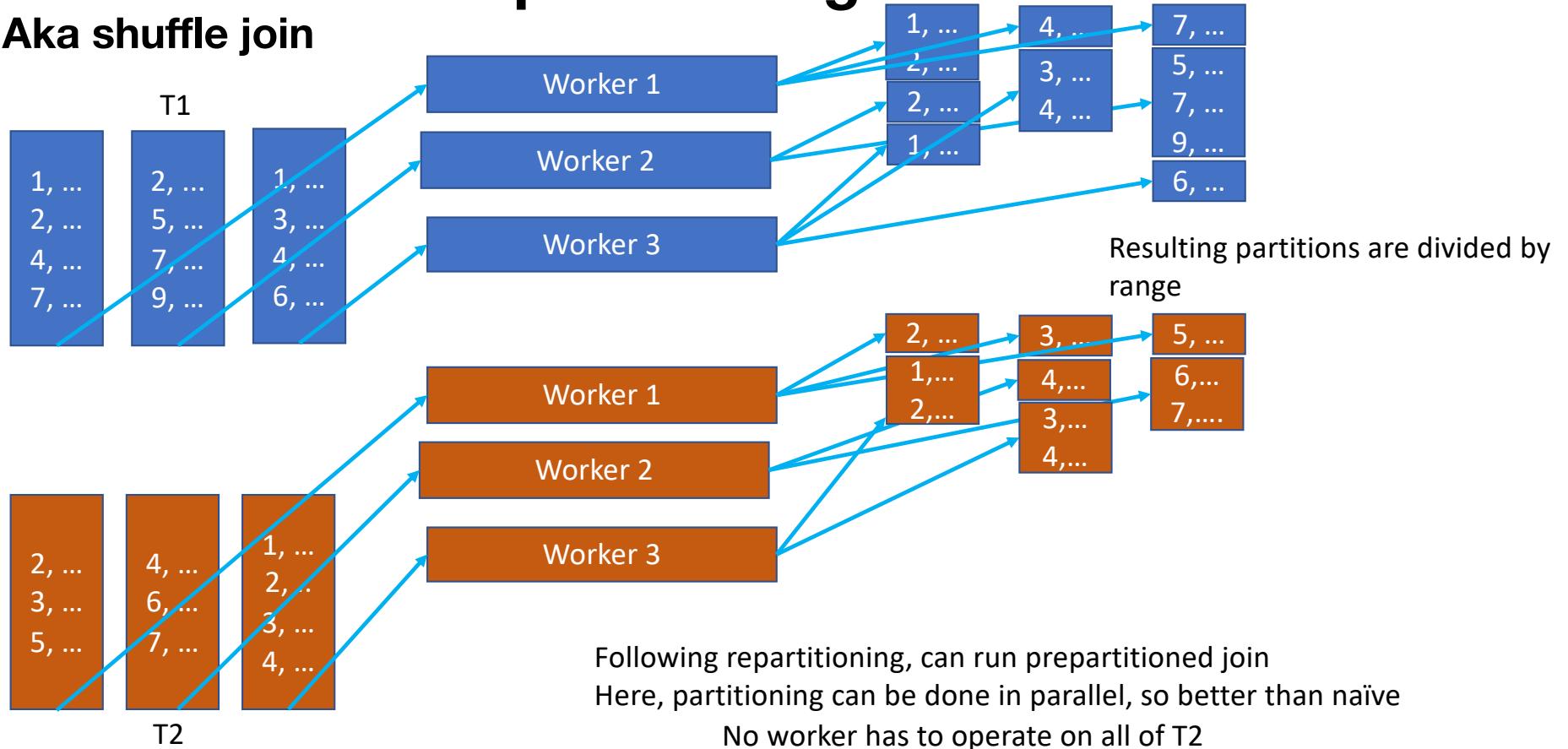
$(1, \dots) \bowtie (1, \dots)$   
 $(1, \dots) \bowtie (1, \dots)$   
 $(2, \dots) \bowtie (2, \dots)$

*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



# Recap: Large Join In Dask

```
client = Client(n_workers=8, threads_per_worker=1, memory_limit='16GB')

header = "CMTE_ID,AMNDT_IND,RPT_TP,TRANSACTION_PGI,IMAGE_NUM,TRANSACTION_TP ..."
PATH = "indiv20/by_date/itcont_2020_20010425_20190425.txt"
PATH2 = "indiv20/by_date/itcont_2020_20190426_20190628.txt"

df = dask.dataframe.read_csv(PATH, low_memory=False, delimiter='|', header=None ...)
df2 = dask.dataframe.read_csv(PATH2, low_memory=False, delimiter='|', header=None ...)
df = df.dropna(subset=['NAME']).drop_duplicates(subset=['NAME'])
df2 = df2.dropna(subset=['NAME']).drop_duplicates(subset=['NAME'])

# make 3 copies
df = df.append(df)
df = df.append(df)
df = df.append(df)

df2 = df2.append(df2)
df2 = df2.append(df2)
df2 = df2.append(df2)

ans = df.merge(df2, on='NAME').count()

ans = ans.compute()      Execution is deferred until compute is called

print(f"found {ans} matches")
```

# Dask Distributed

*“Distributed” = multiple machine*

*“Parallel” = multiple processors on same machine*

- Demo on Amazon
  - Much slower than laptop, t3.large machines (8GB RAM, 2x vCPU ~30% performance / CPU)
- Single local executor: 174.3 s
- Single distributed worker: 200.5
- Three distributed workers: 78.5 s (2.2x/2.6 speedup)

# Subgraph Caching via “Persist”

- Can “persist” a subresult to cause it to be stored in memory
- Avoids recomputing

```
n1 = df.loc[:, ["NAME"]].persist()
n2 = df2.loc[:, ["NAME"]].persist()

#will compute the count and persist n1 and n2
ans = n1.merge(n2, on='NAME').count()
print(ans.compute())

#will reuse previously persisted result
ans2 = n1.merge(n2, on='NAME').max()
print(ans2.compute())
```

# Fault Tolerance Model

- Retries tasks that fail
- Resilient to the failure of any one worker
- Demo

# Spark

- Distributed / parallel data processing system
- pyspark.sql engine very similar to dask in functionality
  - Slightly different API
  - Other pandas-on-spark projects, e.g., koalas provide pandas API compatibility

# Example

Demo!

```
spark = SparkSession.builder.appName("SimpleApp").getOrCreate()

path = "indiv20/by_date/itcont_2020_20010425_20190425.txt"
path2 = "indiv20/by_date/itcont_2020_20190426_20190628.txt"
header = "CMTE_ID,AMNDT_IND,RPT_TP,TRANSACTION_PGI,IMAGE_NUM,TRANSACTION_TP, ...

df_spark = spark.read.csv(path, sep='|', header = False)
df_spark = df_spark.toDF(*header)
df_spark = df_spark.dropna(subset=["NAME"]).dropDuplicates(subset=["NAME"])
df_spark = df_spark.union(df_spark)
df_spark = df_spark.union(df_spark)
df_spark = df_spark.union(df_spark)

df_spark2 = spark.read.csv(path2, sep='|', header = False)
df_spark2 = df_spark2.toDF(*header)
df_spark2 = df_spark2.dropna(subset=["NAME"]).dropDuplicates(subset=["NAME"])
df_spark2 = df_spark2.union(df_spark2)
df_spark2 = df_spark2.union(df_spark2)
df_spark2 = df_spark2.union(df_spark2)

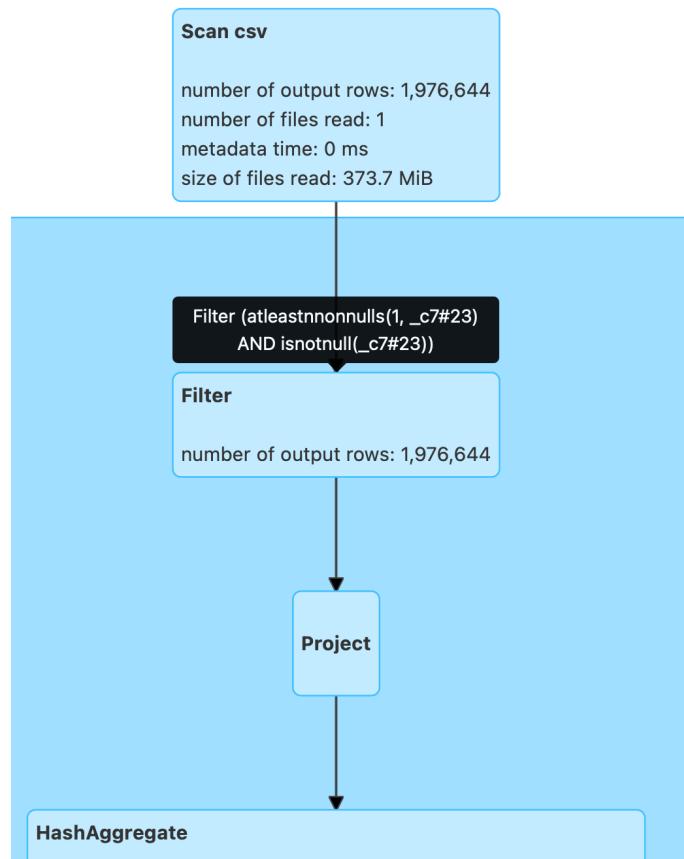
ans = df_spark.join(df_spark2, on='NAME').count()
print(ans)
```

*This is a way to run spark locally;  
most people run a cluster of machines  
and submit jobs, like the dask  
distributed demo before*

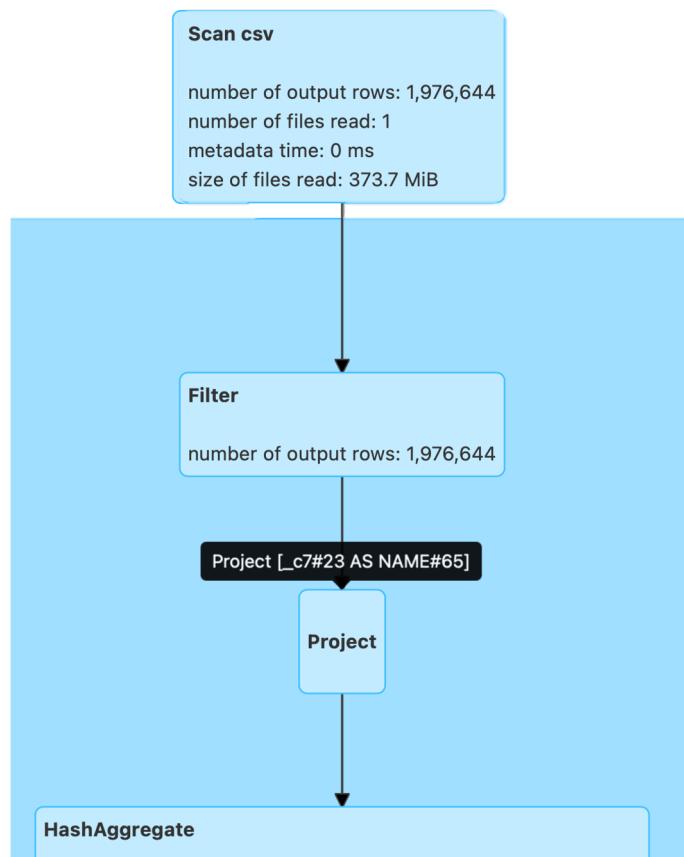
# Spark Under the Hood

- Compiles to Java/Scala
  - Makes understand what tasks are doing and debugging messages somewhat confusing
- Query optimizer much smarter than Dask
  - Projection push down
  - Pre-aggregation

# Projection Push Down

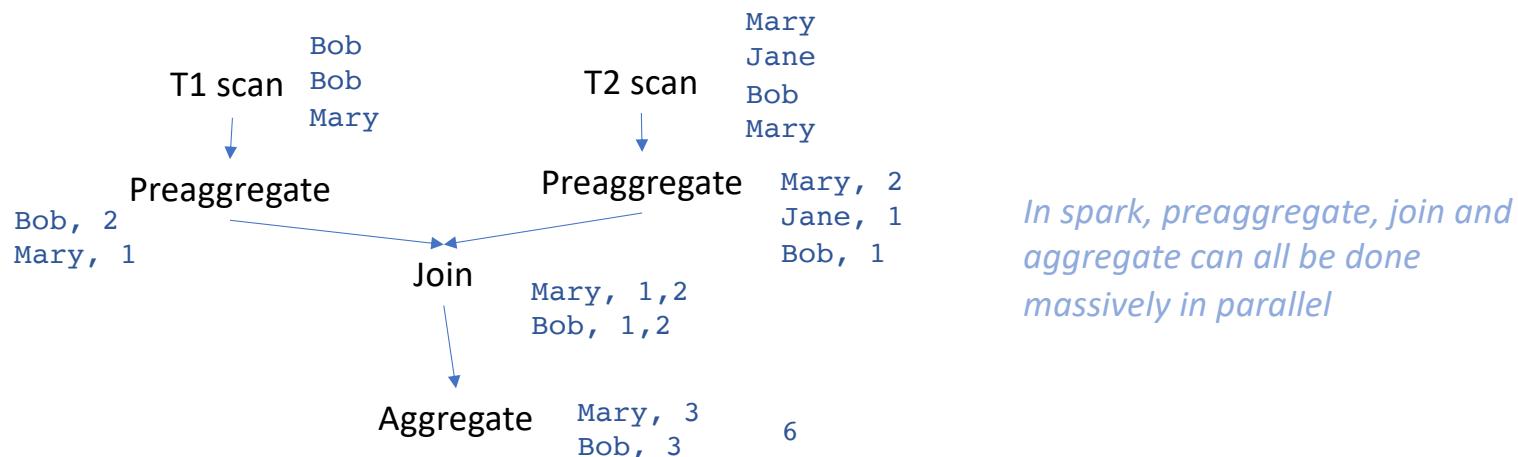


# Projection Push Down



# Preaggregation

- Goal: count the number satisfying records in the join
- Idea: count records in each table before the join
- Join {record, count} pairs from tables to compute final join
- Eliminates the number of records that need to join



# **Spark vs Dask**

- Dask is much smaller, more pythonic
- Spark generally performs better
  - More optimized for very large datasets on S3 / cloud storage
  - Dask lacks query optimization
- Spark is harder to use and debug
  - Compilation down to Java makes it hard to understand what is happening, sometimes
- Many other packages in spark, including
  - SparkML
  - Spark Streaming
  - A variety of data lake / storage tools

# Summary

- Dask and Spark both support parallel and distributed computation over data
  - Both scale from a few processors to hundreds of machines
- Dask is good for parallelizing pandas/numpy code
- Spark more sophisticated, less tied to python ecosystem

# Ray: Distributed Framework for Emerging AI Applications

- Designed for training and inference of reinforcement learning (RL) models
- Needed to support three distinct workloads:
  - Model serving (inference)
  - Simulation
  - Model training
- Python native

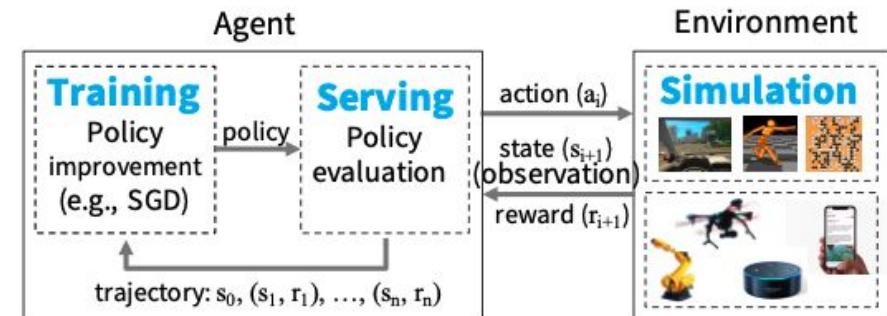


Figure 1: Example of an RL system.

# Support for RL Workloads $\leftrightarrow$ Support for Most ML Workloads

- Training and inference are two key workloads across all of ML
  - RL uses both by definition
- Supporting (distributed) low-latency simulation effectively means supporting (distributed) parallel computation
  - Very useful for feature generation
- Ray's Task/Actor programming model gives tremendous flexibility and control over computation

# Example RL Program

```
// evaluate policy by interacting with env. (e.g., simulator)
rollout(policy, environment):
    trajectory = []
    state = environment.initial_state()
    while (not environment.has_terminated()):
        action = policy.compute(state) // Serving
        state, reward = environment.step(action) // Simulation
        trajectory.append(state, reward)
    return trajectory

// improve policy iteratively until it converges
train_policy(environment):
    policy = initial_policy()
    while (policy has not converged):
        trajectories = []
        for i from 1 to k:
            // evaluate policy by generating k rollouts
            trajectories.append(rollout(policy, environment))
            // improve policy
            policy = policy.update(trajectories) // Training
    return policy
```

Figure 2: Typical RL pseudocode for learning a policy.

# Ray Programming Model

- Task/Actor model
  - **Task:** a **stateless** computation
    - Ex: compute action in env. and get observation
    - Ex: pre-process training input
  - **Actor:** a **stateful** computation
    - Ex: maintain state of simulation environment
    - Ex: update model weights
- Both execute on a remote worker (i.e. node) and return a **future**
- *Roughly speaking:* task  $\leftrightarrow$  function; actor  $\leftrightarrow$  class

Tasks (stateless)	Actors (stateful)
Fine-grained load balancing	Coarse-grained load balancing
Support for object locality	Poor locality support
High overhead for small updates	Low overhead for small updates
Efficient failure handling	Overhead from checkpointing

Table 2: Tasks vs. actors tradeoffs.

# Ray Programming Model

Name	Description
<code>futures = f.remote(args)</code>	Execute function $f$ remotely. <code>f.remote()</code> can take objects or futures as inputs and returns one or more futures. This is non-blocking.
<code>objects = ray.get(futures)</code>	Return the values associated with one or more futures. This is blocking.
<code>ready_futures = ray.wait(futures, k, timeout)</code>	Return the futures whose corresponding tasks have completed as soon as either $k$ have completed or the timeout expires.
<code>actor = Class.remote(args)</code> <code>futures = actor.method.remote(args)</code>	Instantiate class $Class$ as a remote actor, and return a handle to it. Call a method on the remote actor and return one or more futures. Both are non-blocking.

Table 1: Ray API

# Ray Programming Model

```
@ray.remote
def create_policy():
    # Initialize the policy randomly.
    return policy
```

Task

```
@ray.remote(num_gpus=1)
class Simulator(object):
    def __init__(self):
        # Initialize the environment.
        self.env = Environment()
    def rollout(self, policy, num_steps):
        observations = []
        observation = self.env.current_state()
        for _ in range(num_steps):
            action = policy(observation)
            observation = self.env.step(action)
            observations.append(observation)
        return observations
```

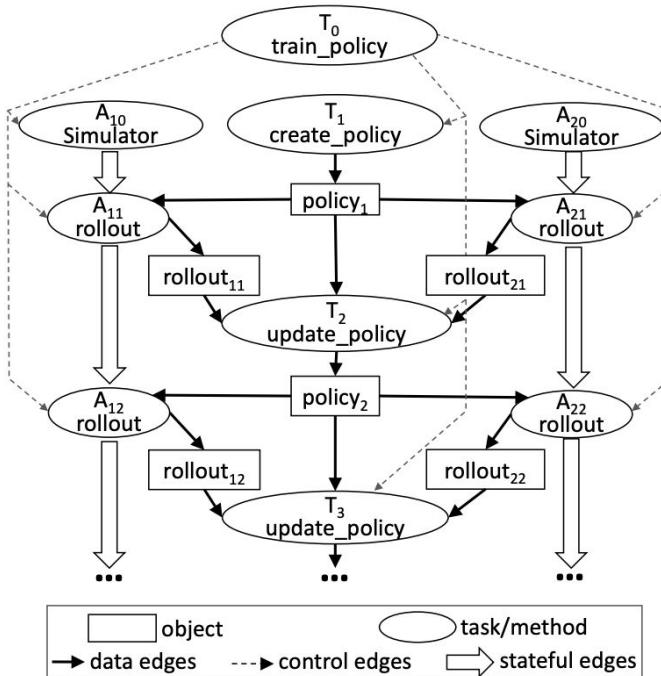
Actor

```
@ray.remote(num_gpus=2)
def update_policy(policy, *rollouts):
    # Update the policy.
    return policy
```

Task

```
@ray.remote
def train_policy():
    # Create a policy.
    policy_id = create_policy.remote()
    # Create 10 actors.
    simulators = [Simulator.remote() for _ in range(10)]
    # Do 100 steps of training.
    for _ in range(100):
        # Perform one rollout on each actor.
        rollout_ids = [s.rollout.remote(policy_id)
                      for s in simulators]
        # Update the policy with the rollouts.
        policy_id =
            update_policy.remote(policy_id, *rollout_ids)
    return ray.get(policy_id)
```

# Ray Computation Model



```
@ray.remote
def train_policy():
    # Create a policy.
    policy_id = create_policy.remote()
    # Create 10 actors.
    simulators = [Simulator.remote() for _ in range(10)]
    # Do 100 steps of training.
    for _ in range(100):
        # Perform one rollout on each actor.
        rollout_ids = [s.rollout.remote(policy_id)
                      for s in simulators]
        # Update the policy with the rollouts.
        policy_id =
            update_policy.remote(policy_id, *rollout_ids)
    return ray.get(policy_id)
```

# Ray: Underneath the Hood

- Application Layer vs. System Layer
- Bottom-up distributed scheduler
- GCS
- In-Memory Distributed Object Store

# Ray: Underneath the Hood

- **Driver** executes program (one driver per program)
- **Workers** execute tasks; **Actors** execute actor methods & maintain state
- In-Memory Object Store is distributed across nodes and manages inputs + intermediate results
  - Implemented w/shared memory using Apache Arrow's Plasma Store

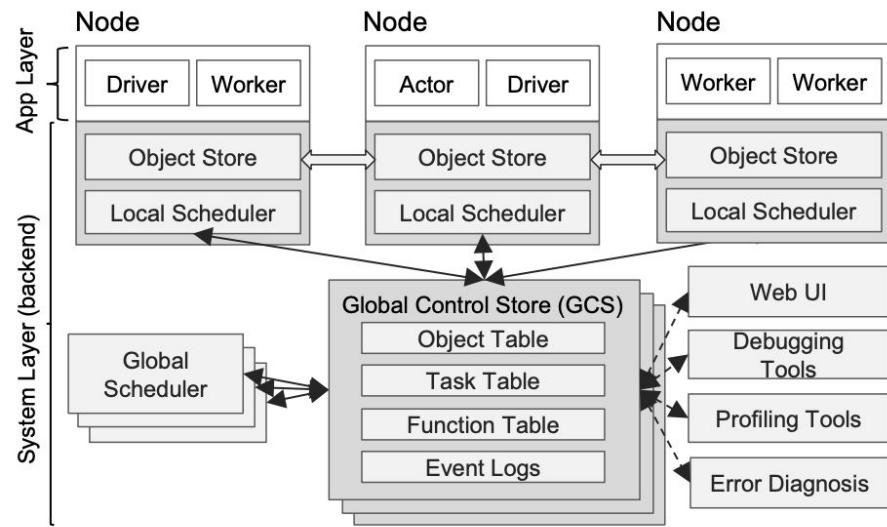


Figure 5: Ray's architecture consists of two parts: an *application* layer and a *system* layer. The application layer implements the API and the computation model described in Section 3, the system layer implements task scheduling and data management to satisfy the performance and fault-tolerance requirements.

# Ray: Underneath the Hood

- **GCS** is designed to support scheduling millions of tasks / sec\*
  - \* = mileage may vary
- “Bottom-up scheduling” → first try to schedule tasks locally
  - Global scheduling only happens if/when node overloaded
- GCS is a (sharded and replicated) key-value store
  - Key: Object / Task IDs
  - Value: node location

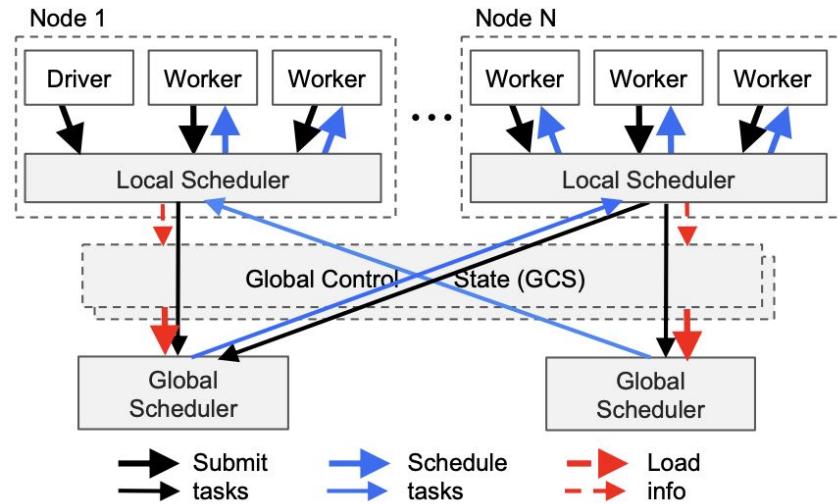
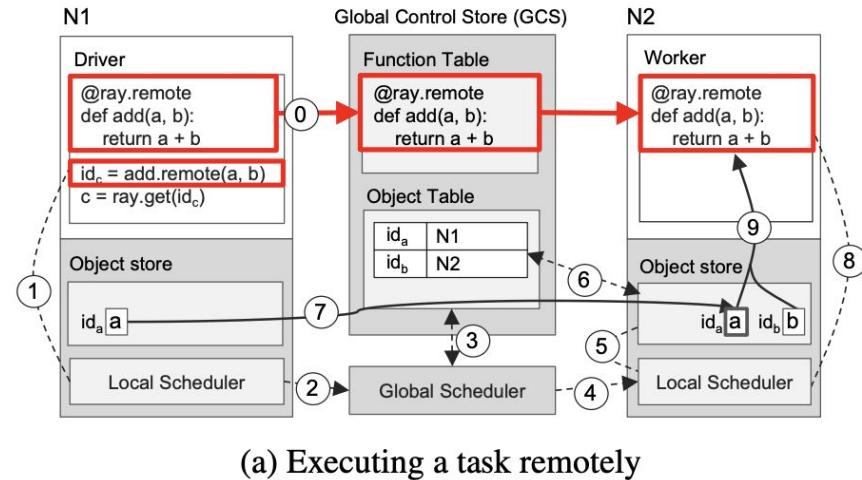


Figure 6: Bottom-up distributed scheduler. Tasks are submitted bottom-up, from drivers and workers to a local scheduler and forwarded to the global scheduler only if needed (Section 4.2.2). The thickness of each arrow is proportional to its request rate.

# Example Task Execution

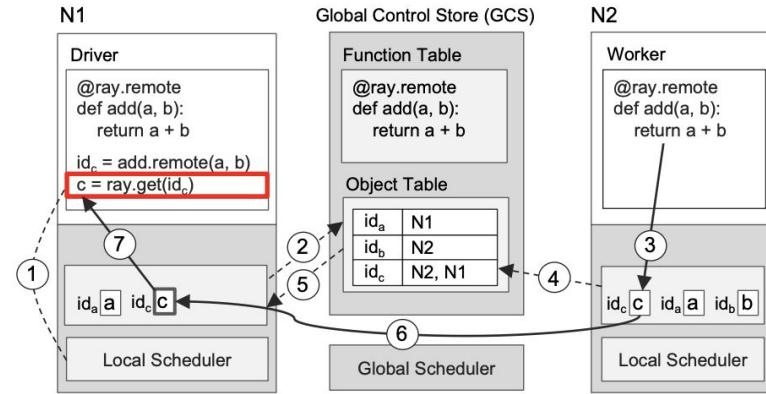
- Obj. A and B start on nodes N1 and N2 respectively
- Remote fcn. **add** is registered w/GCS upon init. and distributed to every worker (step 0)



- Steps (1-4): task submitted at N1 gets scheduled on N2 (for sake of ex.)
- Steps (5-7): input A is copied to N2 to bring all inputs to N2
- Steps (8-9): local sched. invokes task once inputs are ready

# Example Returning Result of Execution

- Same setup as in previous slide
- Steps (1-2): lookup for value C results in N1's obj. store registering callback w/Object Table



(b) Returning the result of a remote task

- Steps (3-5): N2 completes execution of **add** and adds entry for C in Object Table
- Steps (5-7): callback is triggered; C is copied to N1 and returned

# Feature Extraction Example

```
# create futures
futures = [
    batch_ids.options(resources={"c5.xlarge": 0.5})
        .remote(s3_bucket, corpus_name, list(trips))
    for trips in trip_chunks
]

# block on results from futures
task_fraction_factor = 100
num_tasks = len(futures)
results = []
while futures:
    # get futures that have finished
    num_returns = min(int(num_tasks/task_fraction_factor), len(futures))
    done_futures, not_done_futures = ray.wait(futures, num_returns=num_returns)

    # fetch results for finished futures
    done_results = ray.get(done_futures)
    results.extend(done_results)
    logging.info(f"Num tasks finished: {len(results)}/{num_tasks}")

    # update futures that still need to finish
    futures = not_done_futures
```

# Summary

- Framework designed for RL
- Turned out to be suitable for most ML workloads
- Python native → easy to parallelize existing Python code across cluster
  - (Sometimes just need to add `@ray.remote` and handle futures)
- Simple but powerful:
  - Allows for fine-grained control of computation + data placement
  - Extensive libraries developed for distributed:
    - Training
    - Inference
    - Dataset transformation