# Parallelism Continued 6.S079 Lecture 19

Sam Madden

4/13/2022

Lab 5 Due

Topics:
Dask distributed
Spark
Pushdown & preaggregation
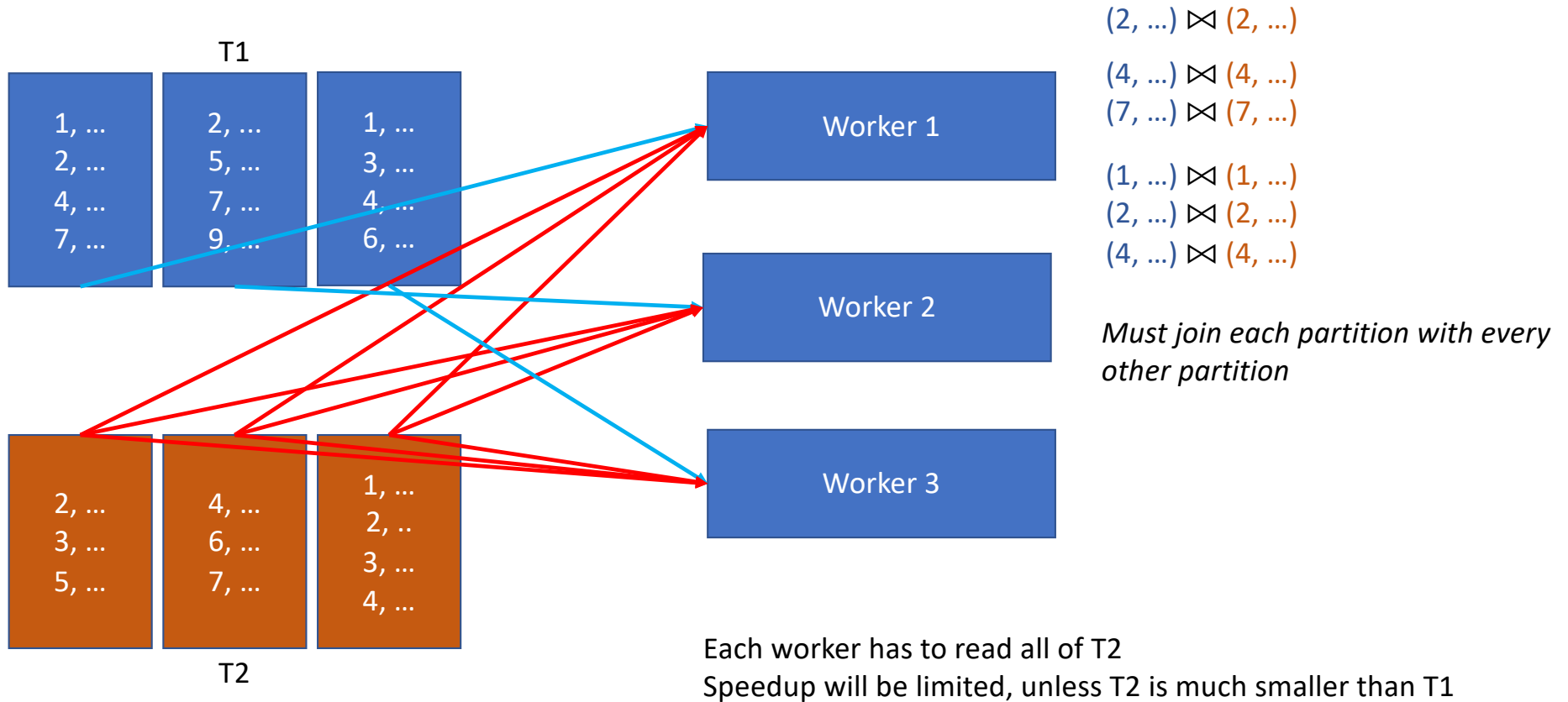Hands-on Dask

# Last Time

- Introduced Parallel Processing
- Look at Parallel Dataflow as a common set of operations that can be readily parallelized
- Studied parallel join and parallel aggregation
- Introduced Dask, a parallel implementation of Python pandas (and numpy and scikit learn)
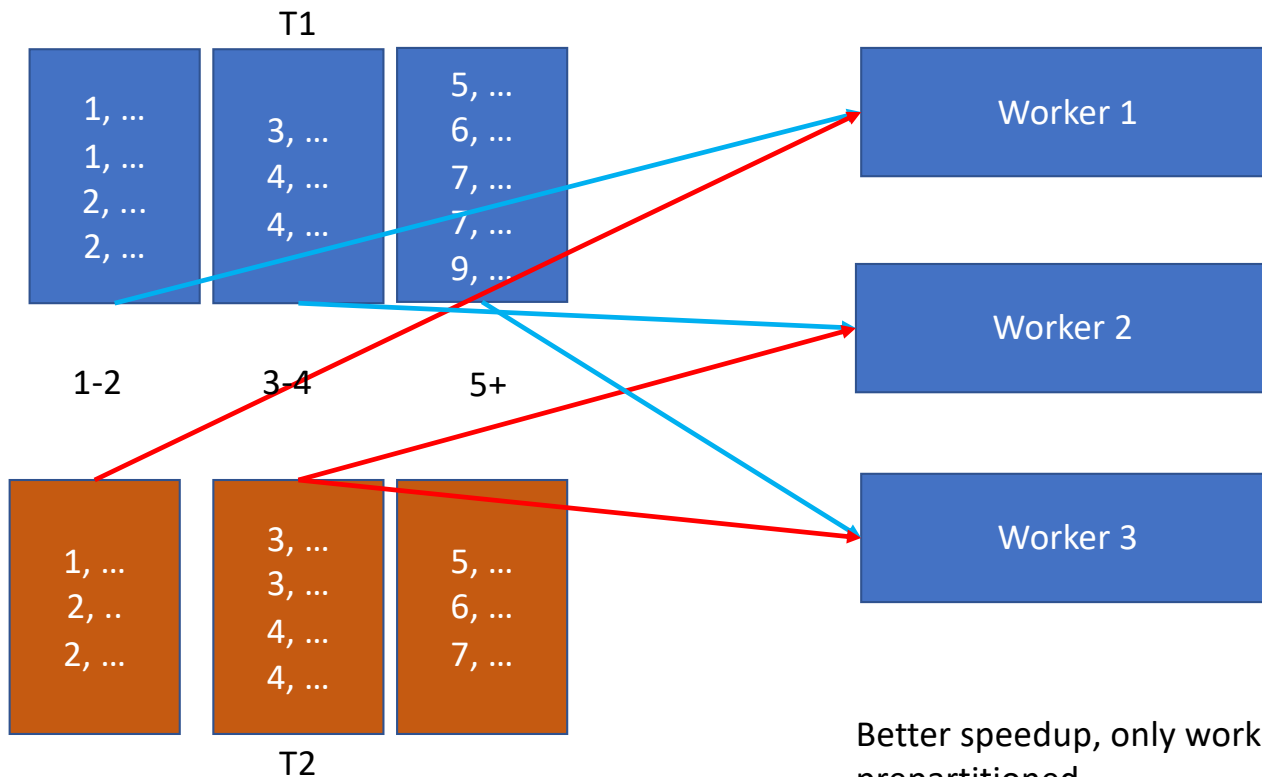
# Parallel Join – Random Partitioning Naïve Algo
## (1, …) indicates value of join attribute

T1

| 1, … | 2, … | 1, … |
| 2, … | 5, … | 3, … |
| 4, … | 7, … | 4, … |
| 7, … | 9, … | 6, … |

| 2, … | 4, … | 1, … |
| 3, … | 6, … | 2, .. |
| 5, … | 7, … | 3, … |
|      |      | 4, … |

T2

Worker 1

Worker 2

Worker 3

$(2, …) \bowtie (2, …)$
$(4, …) \bowtie (4, …)$
$(7, …) \bowtie (7, …)$

$(1, …) \bowtie (1, …)$
$(2, …) \bowtie (2, …)$
$(4, …) \bowtie (4, …)$

*Must join each partition with every other partition*

Each worker has to read all of T2
Speedup will be limited, unless T2 is much smaller than T1

# Parallel Join – Prepartitioned
## (1, …) indicates value of join attribute

*Only need to join partitions that match*

T1

| 1, … | 3, … | 5, … |
| 1, … | 4, … | 6, … |
| 2, … | 4, … | 7, … |
| 2, … |  | 7, … |
|  |  | 9, … |

1-2          3-4          5+

| 1, … | 3, … | 5, … |
| 2, .. | 3, … | 6, … |
| 2, … | 4, … | 7, … |
|  | 4, … |  |

T2

Worker 1

Worker 2

Worker 3

(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

T1

| 1, ... | 2, ... | 1, ... |
| 2, ... | 5, ... | 3, ... |
| 4, ... | 7, ... | 4, ... |
| 7, ... | 9, ... | 6, ... |

Worker 1

Worker 2

Worker 3

| 1, ... | 4, ... | 7, ... |
| 2, ... | 3, ... | 5, ... |
| 2, ... | 4, ... | 7, ... |
| 1, ... |   | 9, ... |
|   |   | 6, ... |

Resulting partitions are divided by range

Worker 1

Worker 2

Worker 3

| 2, ... | 3, ... | 5, ... |
| 1,... | 4,... | 6,... |
| 2,... | 3,... | 7,.... |
|   | 4,... |   |

T2

| 2, ... | 4, ... | 1, ... |
| 3, ... | 6, ... | 2, .. |
| 5, ... | 7, ... | 3, ... |
|   |   | 4, ... |

Following repartitioning, can run prepartitioned join
Here, partitioning can be done in parallel, so better than naïve
No worker has to operate on all of T2

# Recap: Large Join In Dask

```python
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

```python
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 resuse previously peristed rsult
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 pands-on-spark projects, e.g., koalas provide pandas API compatibility

# Example

Demo!

```python
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)
```
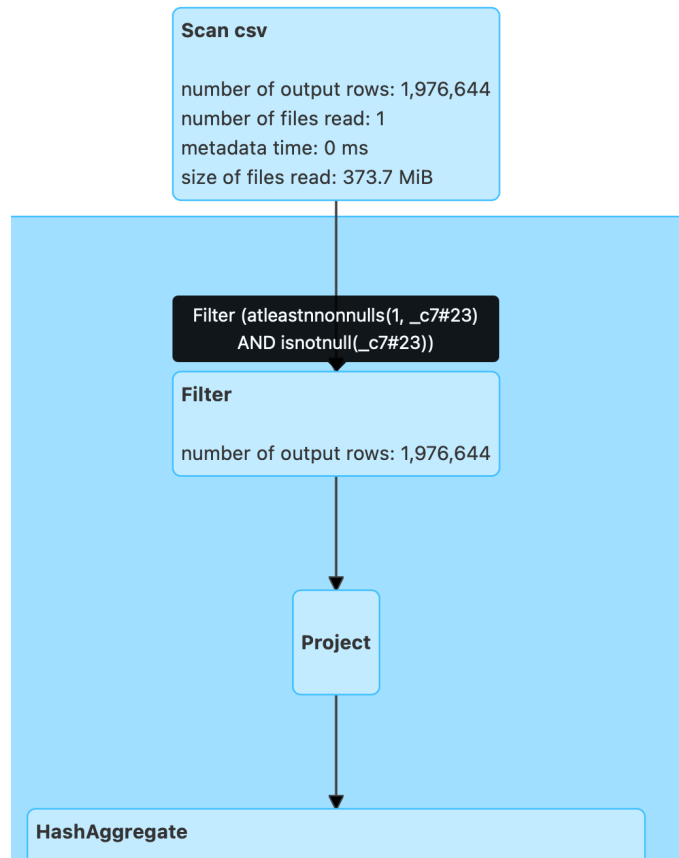
*This is a way to run spark locally; most people run a cluster of machines and submit jobs, like the dask distributed demo before*
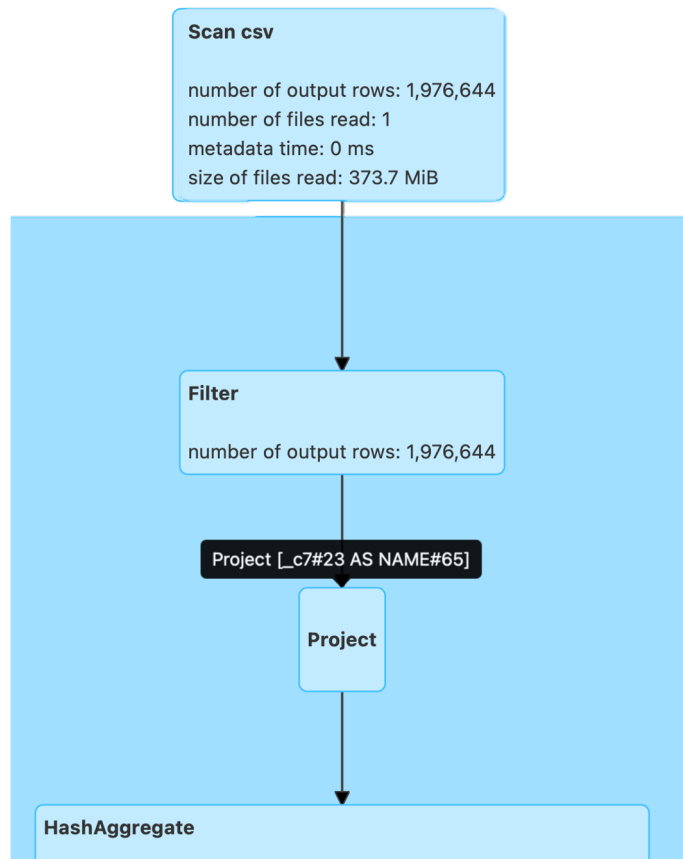
# Spark Under the Hood

- 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



**Scan csv**

number of output rows: 1,976,644
number of files read: 1
metadata time: 0 ms
size of files read: 373.7 MiB

Filter (atleastnnonnulls(1, _c7#23)
AND isnotnull(_c7#23))

**Filter**

number of output rows: 1,976,644

**Project**

**HashAggregate**

# Projection Push Down

**Scan csv**

number of output rows: 1,976,644
number of files read: 1
metadata time: 0 ms
size of files read: 373.7 MiB

**Filter**

number of output rows: 1,976,644
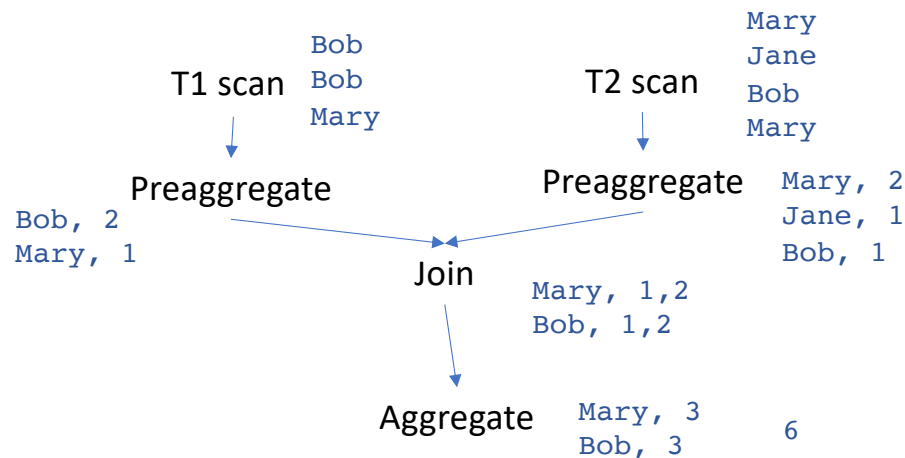
Project [_c7#23 AS NAME#65]

**Project**

**HashAggregate**

# 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

T1 scan
```
Bob
Bob
Mary
```

T2 scan
```
Mary
Jane
Bob
Mary
```

Preaggregate
```
Bob, 2
Mary, 1
```

Preaggregate
```
Mary, 2
Jane, 1
Bob, 1
```

Join
```
Mary, 1,2
Bob, 1,2
```

Aggregate
```
Mary, 3
Bob, 3
```
6

*In spark, preaggregate, join and aggregate can all be done massively in parallel*

# 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