



Optimizing Apache Spark

The Five Most Common  
Performance Problems

Skew

# The 5 Most Common Performance Problems (The 5 Ss)

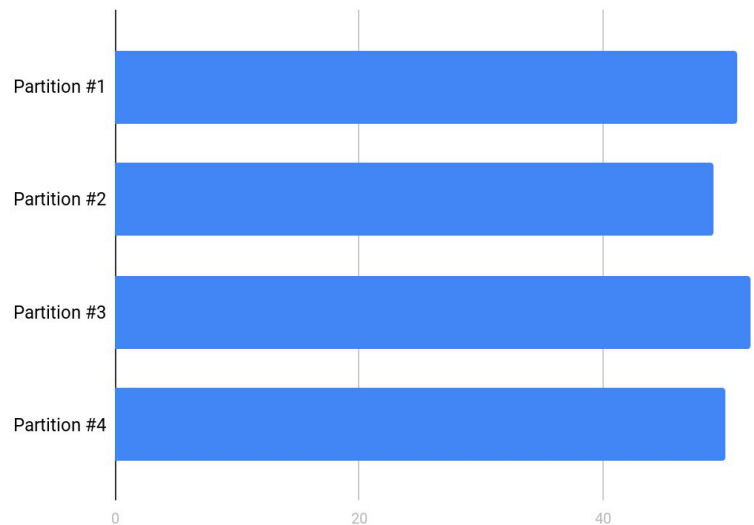
## Skew

- Data is typically read in as 128 MB partitions and evenly distributed
- As the data is transformed (e.g. aggregated), it's possible to have significantly more records in one partition than another
- A small amount of skew is ignorable
- But large skews can result in spill or worse, hard to diagnose OOM Errors

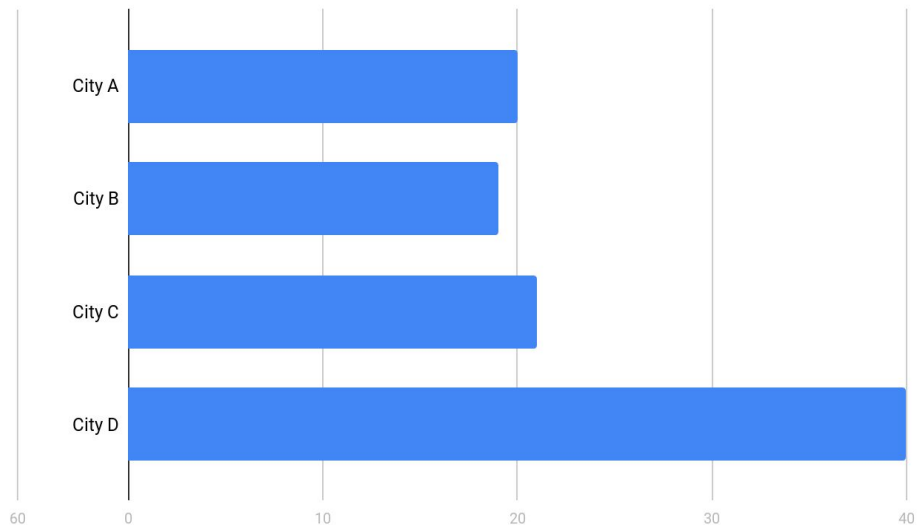
# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Before & After

Before aggregation



After aggregation by city



# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Ramifications

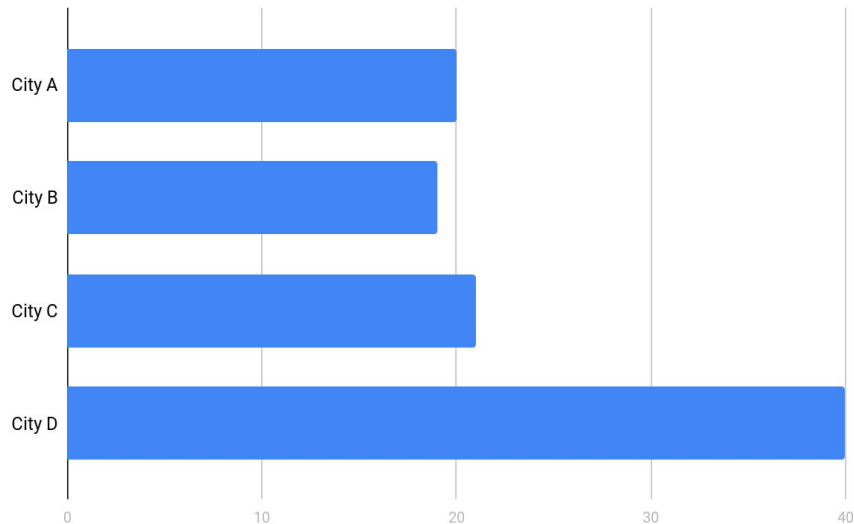
If **City D** is 2x larger than A, B or C...

- It takes 2x as long to process
- It requires 2x as much RAM

The ramifications of that is...

- The entire stage will take as long as the longest running task
- We may not have enough RAM for these skewed partitions

After aggregation by city



What can we do to mitigate skew?

# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Time vs RAM

We need to first ask which problem are we solving for?

- Solving for the RAM problem is only treating the symptoms and not the root cause.
- The RAM problem manifests itself as **Spill** and/or OOM Errors and should **not** be the first thing we solve for...*more on spill later*
- The first problem to solve for is the uneven distribution of records across all partitions which manifests itself as proportionally slower tasks

# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Mitigation

There are several strategies for fixing skew:

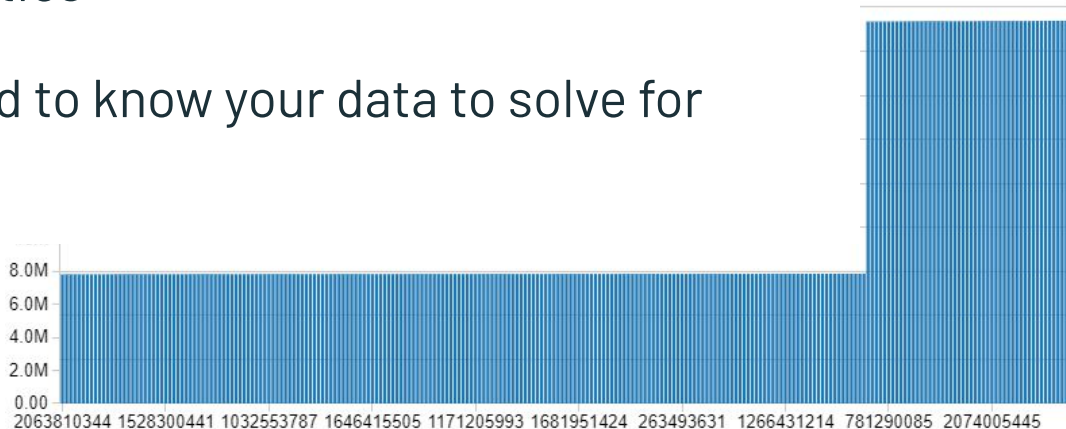
- Employ a Databricks-specific skew hint (see [Skew Join optimization](#))
- Enable Adaptive Query Execution in Spark 3
- Salt the skewed column with a random number creating better distribution across each partition at the cost of extra processing

# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - How Skewed?

See [Experiment #1596](#), **Step B**

- Our perfectly engineered data has a skew in US cities that is ~3x larger than all other countries
- Counts come in at 23 million for skewed cities vs 8 million for other cities
- As we will see, you really need to know your data to solve for this...*maybe not with AQE*





# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Baseline vs Hint

See [Experiment #1596](#), **Step C** and **Step D**

- Contrast the **last stage** of the **last job** for the two commands
  - Note the key code differences
  - Note the total execution time of the corresponding commands
  - Note the total number of tasks
  - In the **Spark UI, Stage Details**
    - Note the “health” of the stage as seen in the **Event Timeline**
    - Note the min/median/max **Shuffle Read Size** under **Summary Metrics**
    - Note the total amount of spill under **Aggregated Metrics by Executor**

# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Baseline vs Hint, Review

Step	Code	Duration	Tasks	Health	Shuffle	Spill
C	Standard	~30 min	832	Bad	0 / 0 / ~100 KB / ~400 MB / ~3 GB	~50 GB
D	Skew Hint	~35 min	832	Mostly OK	134 MB / 174 MB / 184 MB / 195 MB / 1.1 GB	~4 GB

- This scenario introduces the Databricks-specific skew hint (see [Skew Join optimization](#))
- Note the call `.hint("skew", "city_id")`

# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Baseline vs Hint vs w/AQE

See [Experiment #1596](#), **Step E** with **Step C** and **Step D**

- Contrast the **last stage** of the **last job** for the two commands
  - Note the key code differences
  - Note the total execution time of the corresponding commands
  - Note the total number of tasks
  - In the **Spark UI, Stage Details**
    - Note the “health” of the stage as seen in the **Event Timeline**
    - Note the min/median/max **Shuffle Read Size** under **Summary Metrics**
    - Note the total amount of spill under **Aggregated Metrics by Executor**

# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Baseline vs Hint, Review

Step	Code	Duration	Tasks	Health	Shuffle	Spill
C	Standard	~30 min	832	Bad	0 / 0 / ~100 KB / ~400 MB / ~3 GB	~50 GB
D	Skew Hint	~35 min	832	Mostly OK	~135 MB / ~175 MB / ~180 MB / ~200 MB / ~1 GB	~4 GB
E	w/AQE	~25 min	1489	Excellent	0 / ~115 MB / ~115 MB / ~125 MB / ~130 MB	0

- **Step E** uses **Spark 3's** new feature **Adaptive Skewed Join**
  - See **`spark.sql.adaptive.skewJoin.enabled`**
  - See **`spark.sql.adaptive.advisoryPartitionSizeInBytes`**
- The first two **jobs** are read in parallel

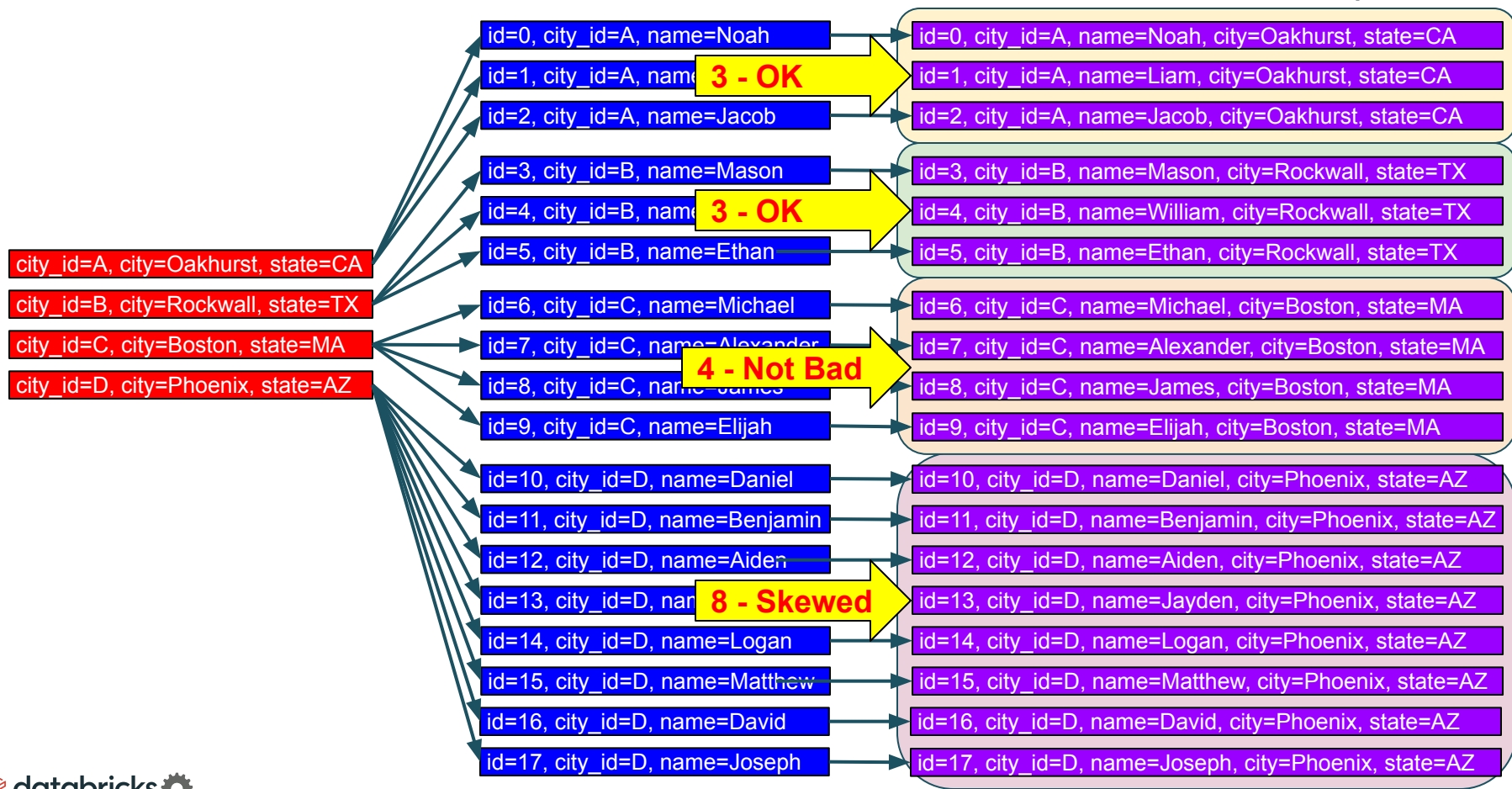
# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Salted Join

- This approach is by far the most complicated to implement
- It can actually take longer to execute in some cases
- Remains a viable option when other solutions are not available
- The idea is to split large partitions into smaller ones using a “salt”
- Has the side effect of splitting small partitions into even smaller ones
- It's more about guaranteeing execution of all tasks  
And **not** a uniform duration for each task

Let's review how a "standard" join works...

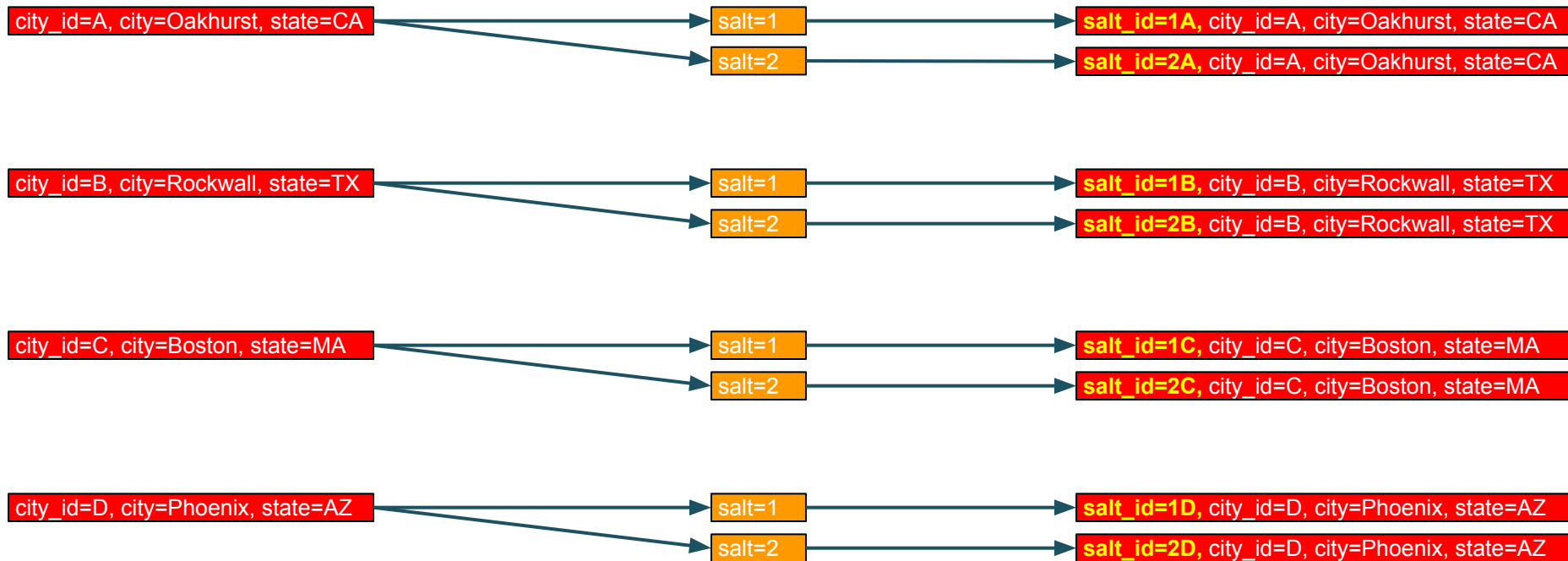
...4 distinct partitions



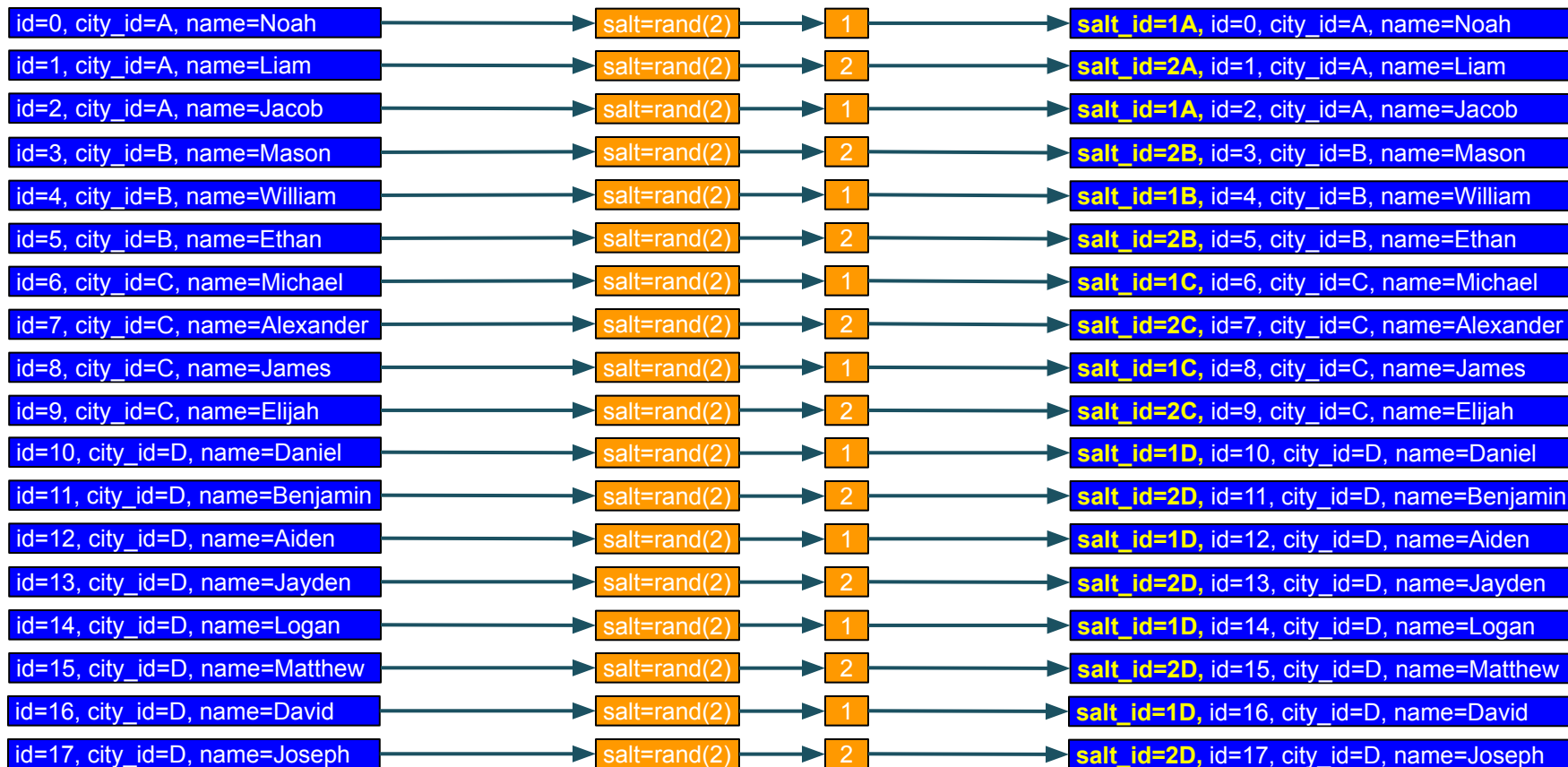
Let's review how a "salted" join works...



# Step #1: Cross join the dimensions to the salt value

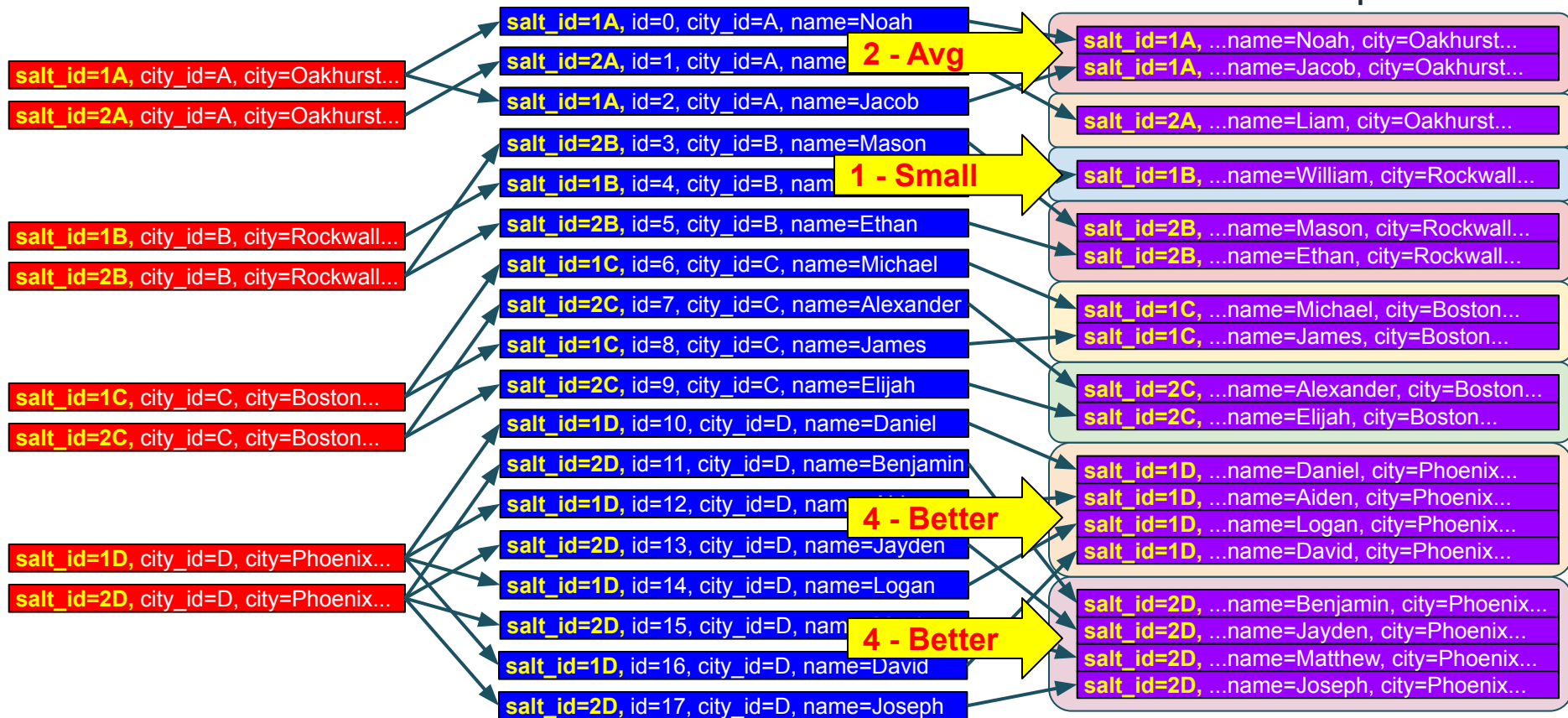


## Step #2: Randomly salt the fact table



# Step #3: Join the salted tables

... 8 distinct partitions



# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Skew Join, in Action

- Step F-1: Create a DataFrame based on the range of our “skew factor”
  - In the visual example, we used “2”
  - In this code example, we are using “7”
  - You can estimate this based on how many times larger the maximum partition is compared to the median partition size
- Step F-2: For the dimension table, cross join the salts with the city table (repartitioning can help mitigate spills and evenly redistributes the new dimension table across all partitions)
- Step F-3: For the fact table, randomly assign a salt to each record
- Step F-4: Join the two tables based on the **salted\_city\_id**

# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Baseline vs Hint vs w/AQE vs Salted

See [Experiment #1596](#), **Step F-4** with **Step C** through **Step E**

- Contrast the **last stage** of the **last job** for the two commands
  - Note the key code differences
  - Note the total execution time of the corresponding commands
  - Note the total number of tasks
  - In the **Spark UI, Stage Details**
    - Note the “health” of the stage as seen in the **Event Timeline**
    - Note the min/median/max **Shuffle Read Size** under **Summary Metrics**
    - Note the total amount of spill under **Aggregated Metrics by Executor**

# The 5 Most Common Performance Problems (The 5 Ss)

## Skew - Baseline vs Hint, Review

Step	Code	Duration	Tasks	Health	Shuffle	Spill
C	Standard	~30 min	832	Bad	0 / 0 / ~100 KB / ~400 MB / ~3 GB	~50 GB
D	Skew Hint	~35 min	832	Mostly OK	~135 MB / ~175 MB / ~180 MB / ~200 MB / ~1 GB	~4 GB
E	w/AQE	~25 min	1489	Excellent	0 / ~115 MB / ~115 MB / ~125 MB / ~130 MB	0
F	Salted	>35 min	832	Better/Bad	~400 KB / ~75 MB / ~150 MB / ~300 MB / ~800 MB	0

- Salting a skewed dataset has a number of problems
- You don't want to salt on the fly - it should be a persisted view of the data
- Consider instead, invest the energy to salt only the skewed keys
- In our example, that would mean salting US cities only

