



Working with Skewed Data: The Iterative Broadcast

Rob Keevil - KnowIT
Fokko Driesprong - GoDataDriven

#EUde11



About Rob Keevil

- Philosophy, Politics & Economics
- Big Data Architect
- Financial and Anti-Fraud Domains
- Scala Programmer



KNOWIT
services

About Fokko Driesprong

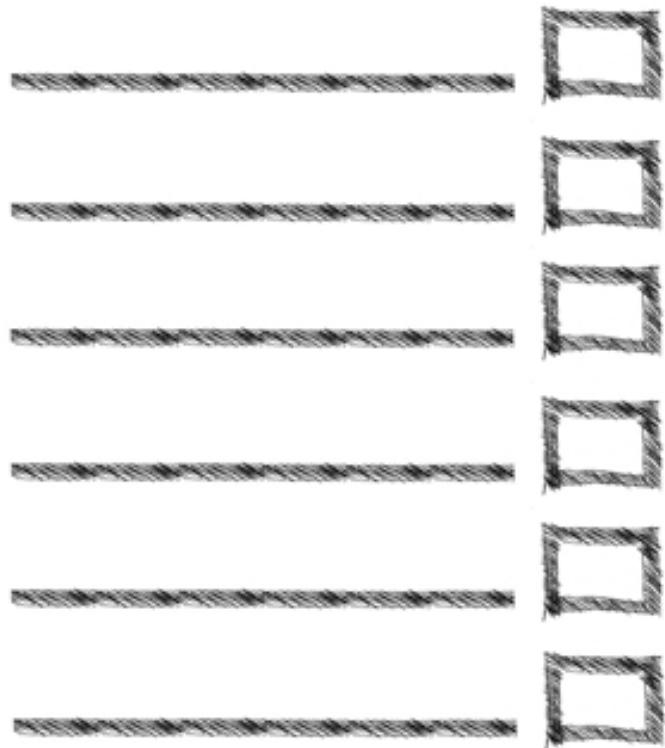
- Software Engineering and Distributed Systems
- Data Engineer at GoDataDriven
- Committer & PPMC Member, Apache Airflow
- Open source enthusiast



GoDataDriven
proudly part of Xebia Group

Overview

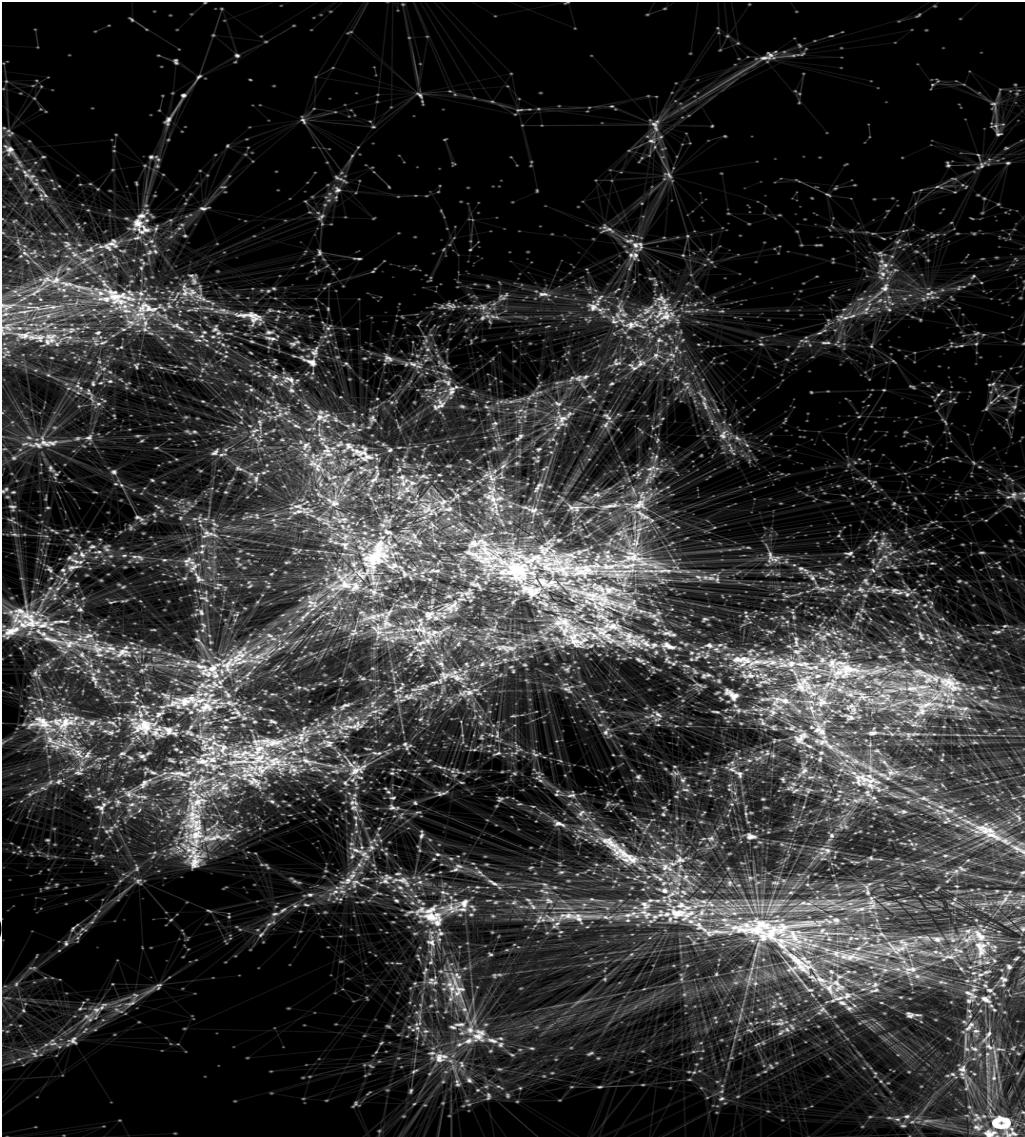
- The Skew Problem at ING
- Join Strategies in Spark
- Our Solution
- Performance Analysis
- Considerations
- Future Work



The Skew Problem at



- Billions of transactions...
- Billions of accounts...
- Millions of companies...
- ...with a very uneven distribution



Familiar?

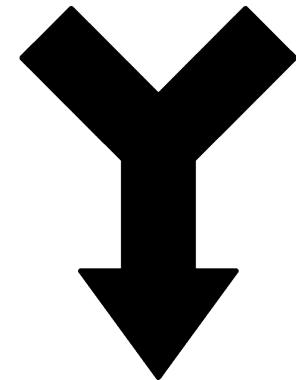
ed	Duration	Tasks: Succeeded/Total	Input	Output	Shuff
18 19:41:01	5.1 min	199/200			6.1 G

Summary Metrics for 199 Completed Tasks

Metric	Min	25th percentile	Median	75th percentile	Max
Duration	4 s	7 s	9 s	13 s	2.2 min
GC Time	0.2 s	0.4 s	0.6 s	2 s	8 s
Shuffle Read Size / Records	9.7 MB / 3542436	14.8 MB / 6104493	16.0 MB / 7704481	18.4 MB / 10158526	40.0 MB / 105456519

Join Types Native to Spark

- Sort Merge Join
- Broadcast Hash Join
- Shuffled Hash Join
- Broadcast Nested Loop Join
- Cartesian Product

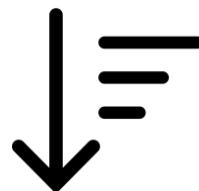


The Sort Merge Join

Key	Column A
1	HELLO
2	SPARK
4	SPAM
3	GOOD

Stage 1:

Determine partitions by hashing the key(s)



Key	Column B
1	WORLD!
4	SPAM
3	MORNING
2	SUMMIT
1	DUBLIN
4	SPAM
3	EVENING
2	STREAMING
1	EVERYBODY
4	SPAM
3	NIGHT
2	SQL

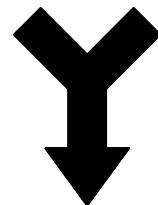
Executor 1

Key	Column A
1	HELLO
3	GOOD

Key	Column A	Column B
1		WORLD!
1		DUBLIN
1		EVERYBODY
3		MORNING
3		EVENING
3		NIGHT

Stage 2:

Merge



Executor 2

Key	Column A
2	SPARK
4	SPAM

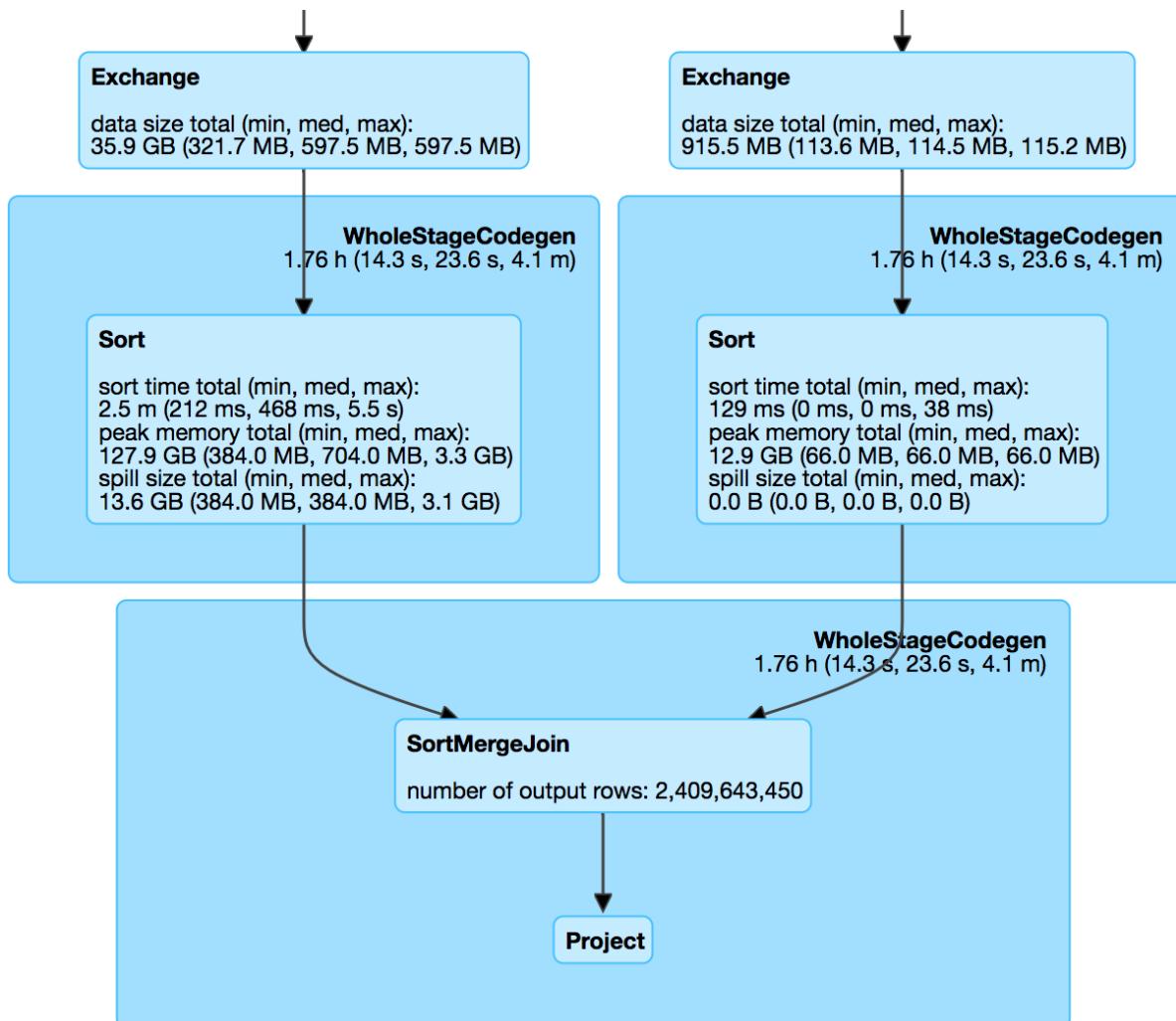
Key	Column A	Column B
2		SUMMIT
2		STREAMING
2		SQL
4		SPAM
4		SPAM
4		SPAM

Executor 1

Key	Column A	Column B
1	HELLO	WORLD!
1	HELLO	DUBLIN
1	HELLO	EVERYBODY
3	GOOD	MORNING
3	GOOD	EVENING
3	GOOD	NIGHT

Executor 2

Key	Column A	Column B
2	SPARK	SUMMIT
2	SPARK	STREAMING
2	SPARK	SQL
4	SPAM	SPAM
4	SPAM	SPAM
4	SPAM	SPAM



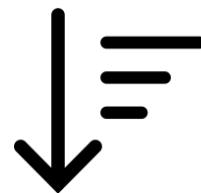
Regular Table

Key	Column A
1	HELLO
2	SPARK
4	SPAM
3	GOOD

Skewed Table

Stage 1:

Determine partitions by hashing the key(s)



Key	Column B
4	SPAM
4	SPAM
4	SPAM
2	SUMMIT
3	MORNING
4	SPAM
4	SPAM
4	SPAM
1	WORLD!
4	SPAM
4	SPAM
4	SPAM

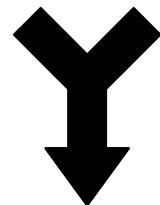
Executor 1

Key	Column A
1	HELLO
3	GOOD

Key	Column A	Column B
1		WORLD!
3		MORNING

Stage 2:

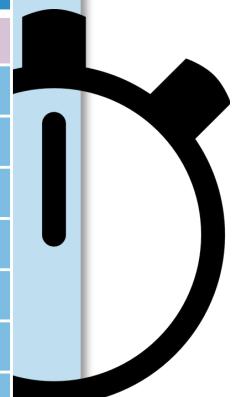
Merge



Executor 2

Key	Column A
2	SPARK
4	SPAM

Key	Column A	Column B
2		SUMMIT
4		SPAM



Executor 1

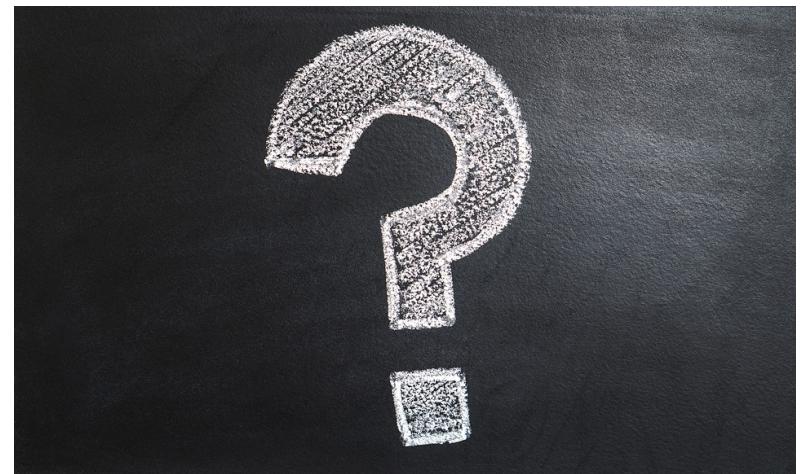
Key	Column A	Column B
1	HELLO	WORLD!
3	GOOD	MORNING

Executor 2

Key	Column A	Column B
2	SPARK	SUMMIT
4	SPAM	SPAM

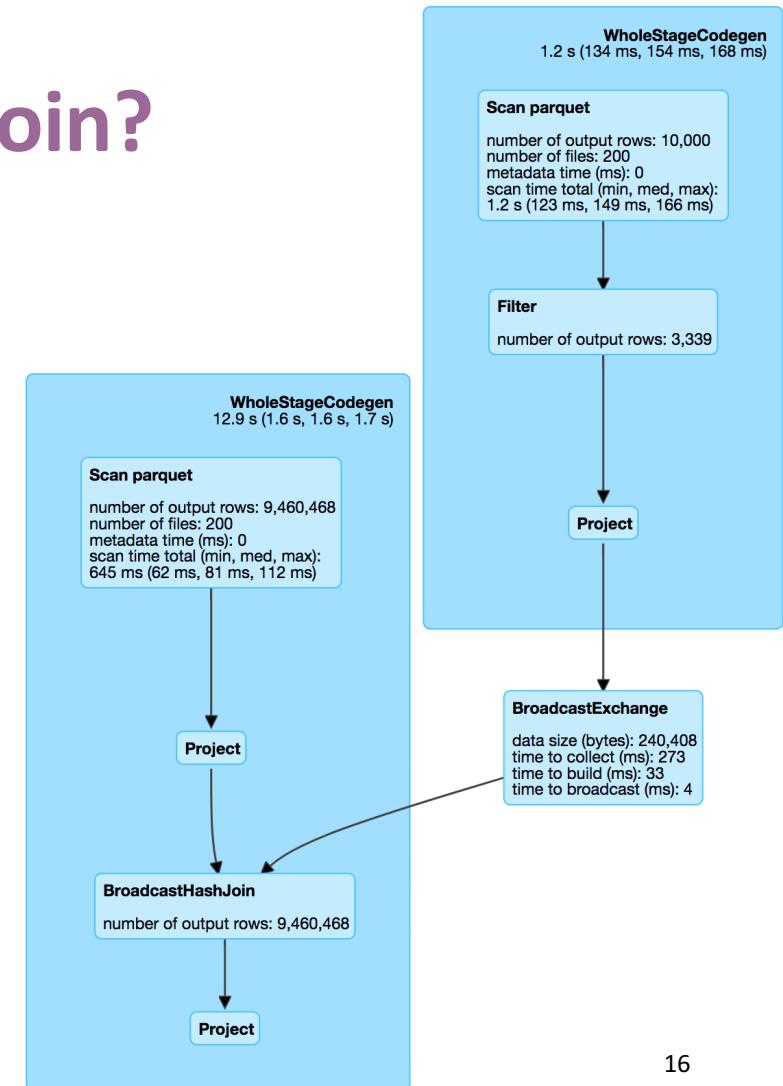
Potential Solutions

- Give it more resources?
- Repartition the data?



What about Broadcast Hash Join?

- Leaves the partitions larger table untouched
- Copies the entire smaller table to every executor
- Iterate over the large table, and join using a HashMap

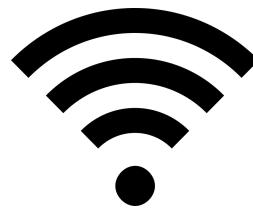


Smaller Table

Key	Column A
1	HELLO
4	SPAM
2	SPARK
3	GOOD

Stage 1:

Broadcast
smaller
dataset to
all
executors



Larger Table

Key	Column B
4	SPAM
4	SPAM
4	SPAM
2	SUMMIT
3	MORNING
4	SPAM
4	SPAM
4	SPAM
1	WORLD!
4	SPAM
4	SPAM
4	SPAM

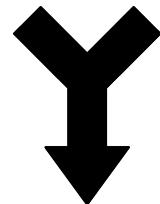
Executor 1

Key	Column A
1	HELLO
4	SPAM
2	SPARK
3	GOOD

Key	Column A	Column B
4		SPAM
4		SPAM
4		SPAM
2		SUMMIT
3		MORNING
4		SPAM

Stage 2:

Left join



Executor 2

Key	Column A
1	HELLO
4	SPAM
2	SPARK
3	GOOD

Key	Column A	Column B
4		SPAM
4		SPAM
1		WORLD!
4		SPAM
4		SPAM
4		SPAM

Executor 1

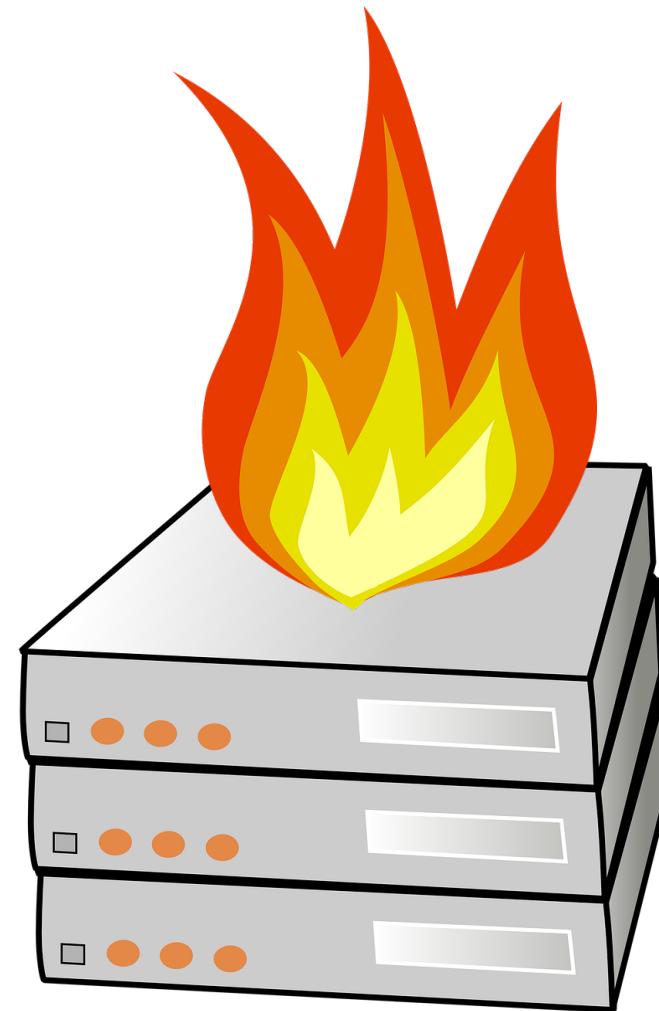
Key	Column A	Column B
4	SPAM	SPAM
4	SPAM	SPAM
4	SPAM	SPAM
2	SPARK	SUMMIT
3	GOOD	MORNING
4	SPAM	SPAM

Executor 2

Key	Column A	Column B
4	SPAM	SPAM
4	SPAM	SPAM
1	HELLO	WORLD!
4	SPAM	SPAM
4	SPAM	SPAM
4	SPAM	SPAM

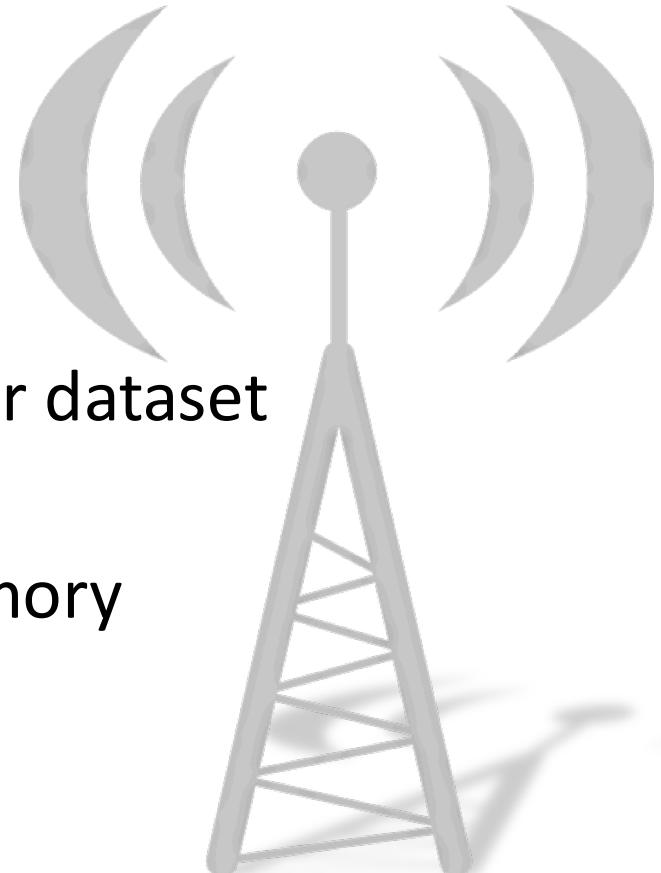
But...

- ...What if your “smaller” table isn’t so small?
- The whole small table needs to fit in driver memory and every executor



Introducing: The Iterative Broadcast!

- Divide the smaller table into “passes”
- Broadcast a pass & left join to the larger dataset
- Clear the broadcast partition from memory
- Repeat until all passes are processed



Stage 1:

Assign Pass Number
to Smaller Table

Key	Column A	Pass
1	HELLO	1
4	SPAM	2
2	SPARK	2
3	GOOD	1

Executor 1

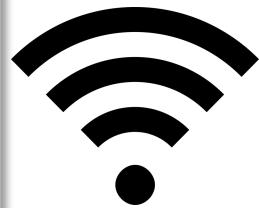
Key	Column A	Pass
1	HELLO	1
3	GOOD	1

Executor 1

Key	Column B
4	SPAM
4	SPAM
4	SPAM
2	SUMMIT
3	MORNING
4	SPAM

Stage 2:

Broadcast
first pass



Executor 2

Key	Column A	Pass
1	HELLO	1
3	GOOD	1

Executor 2

Key	Column B
4	SPAM
4	SPAM
1	WORLD!
4	SPAM
4	SPAM
4	SPAM

Executor 1

Key	Column A	Pass
1	HELLO	1
3	GOOD	1

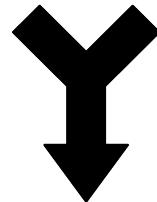
Key	Column A	Column B
4		SPAM
4		SPAM
4		SPAM
2		SUMMIT
3		MORNING
4		SPAM

Stage 3:

Left join
first pass

Executor 2

Key	Column A	Pass
1	HELLO	1
3	GOOD	1



Key	Column A	Column B
4		SPAM
4		SPAM
1		WORLD!
4		SPAM
4		SPAM
4		SPAM

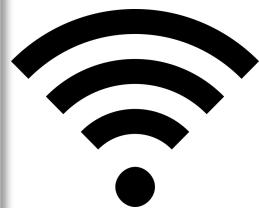
Executor 1

Key	Column A	Pass
2	SPARK	2
4	SPAM	2

Key	Column A	Column B
4		SPAM
4		SPAM
4		SPAM
2		SUMMIT
3	GOOD	MORNING
4		SPAM

Stage 4:

Broadcast second
pass



Executor 2

Key	Column A	Pass
2	SPARK	2
4	SPAM	2

Key	Column A	Column B
4		SPAM
4		SPAM
1	HELLO	WORLD!
4		SPAM
4		SPAM
4		SPAM

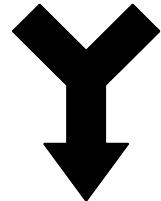
Executor 1

Key	Column A	Pass
2	SPARK	2
4	SPAM	2

Key	Column A	Column B
4		SPAM
4		SPAM
4		SPAM
2		SUMMIT
3	GOOD	MORNING
4		SPAM

Stage 5:

Left join
second
pass



Executor 2

Key	Column A	Pass
2	SPARK	2
4	SPAM	2

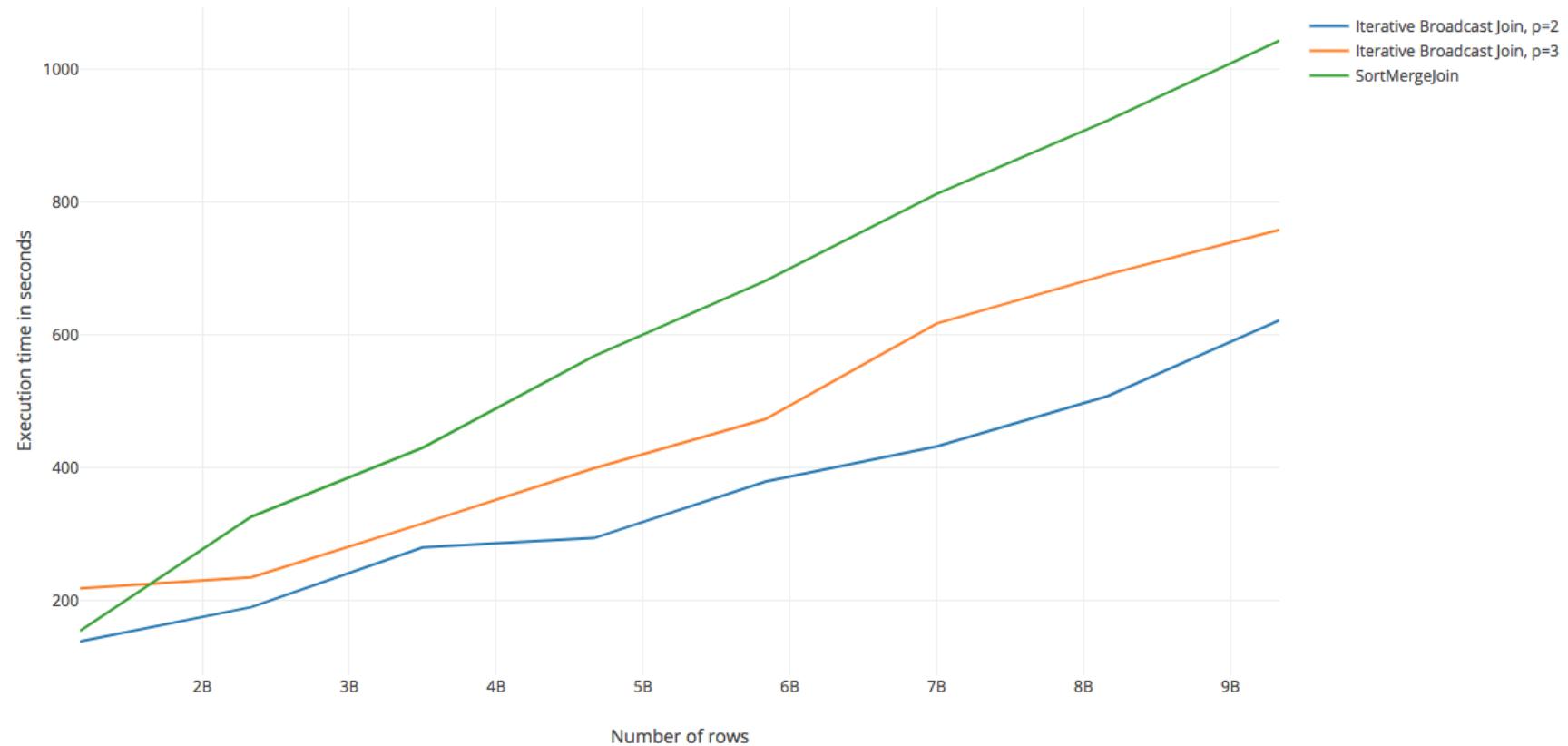
Key	Column A	Column B
4		SPAM
4		SPAM
4		SPAM
1	HELLO	WORLD!
4		SPAM
4		SPAM
4		SPAM

Performance Analysis

- Setup notes
 - Tests on Amazon EMR using 4x m4.4xlarge workers
 - 5 cores and 18gb memory per executor
 - Spark 2.1
 - Benchmark on GitHub
 - Skewed dataset, 10000 keys
 - Largest key 10000 rows, 2nd key 5000 rows, 3rd key 2500 rows ... 10000th key 1 row
 - Uniform dataset, evenly spread

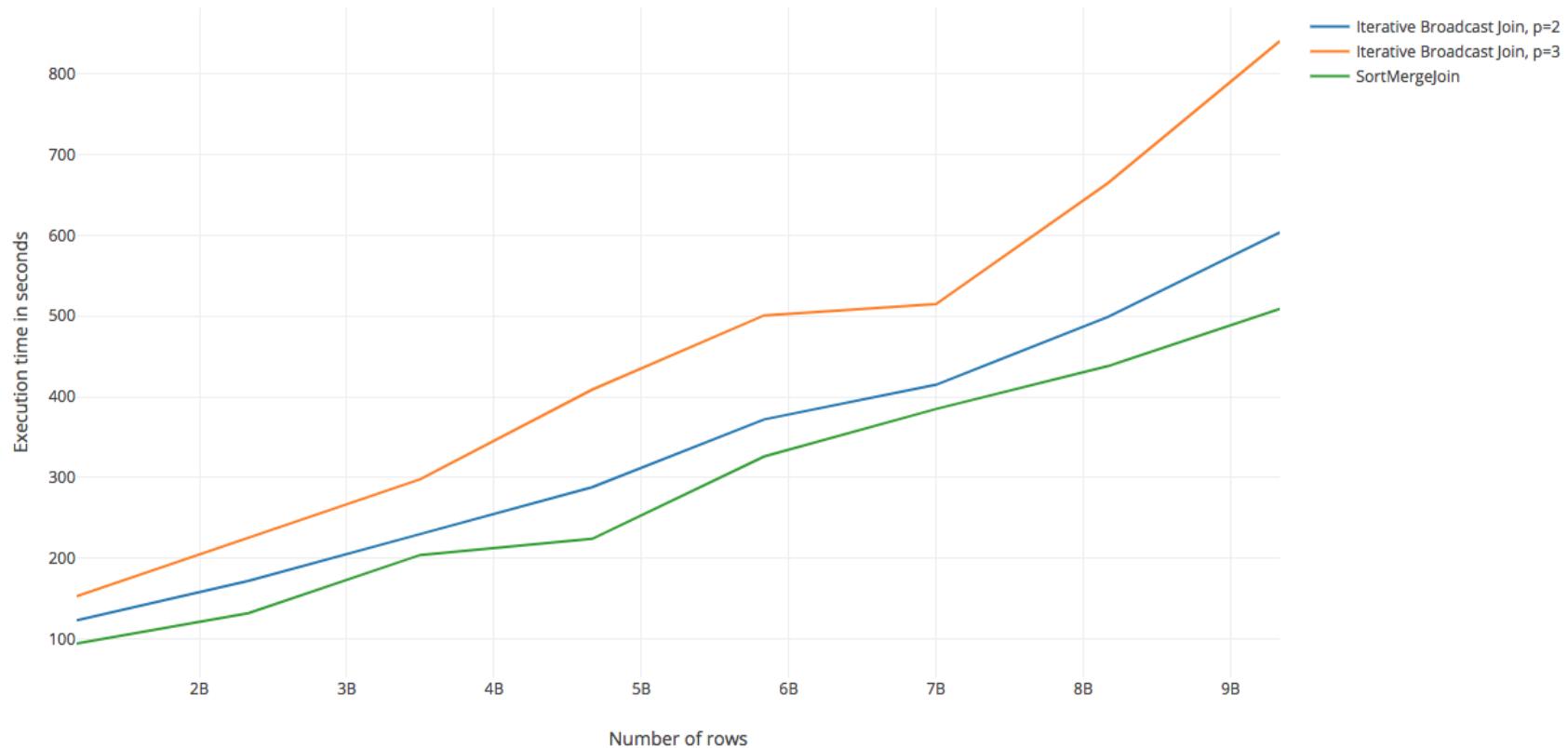
Performance Analysis – Skewed Data

Skewed partitions, 100.000 keys



Performance Analysis – Uniform Data

Uniform partitions, 100.000 keys

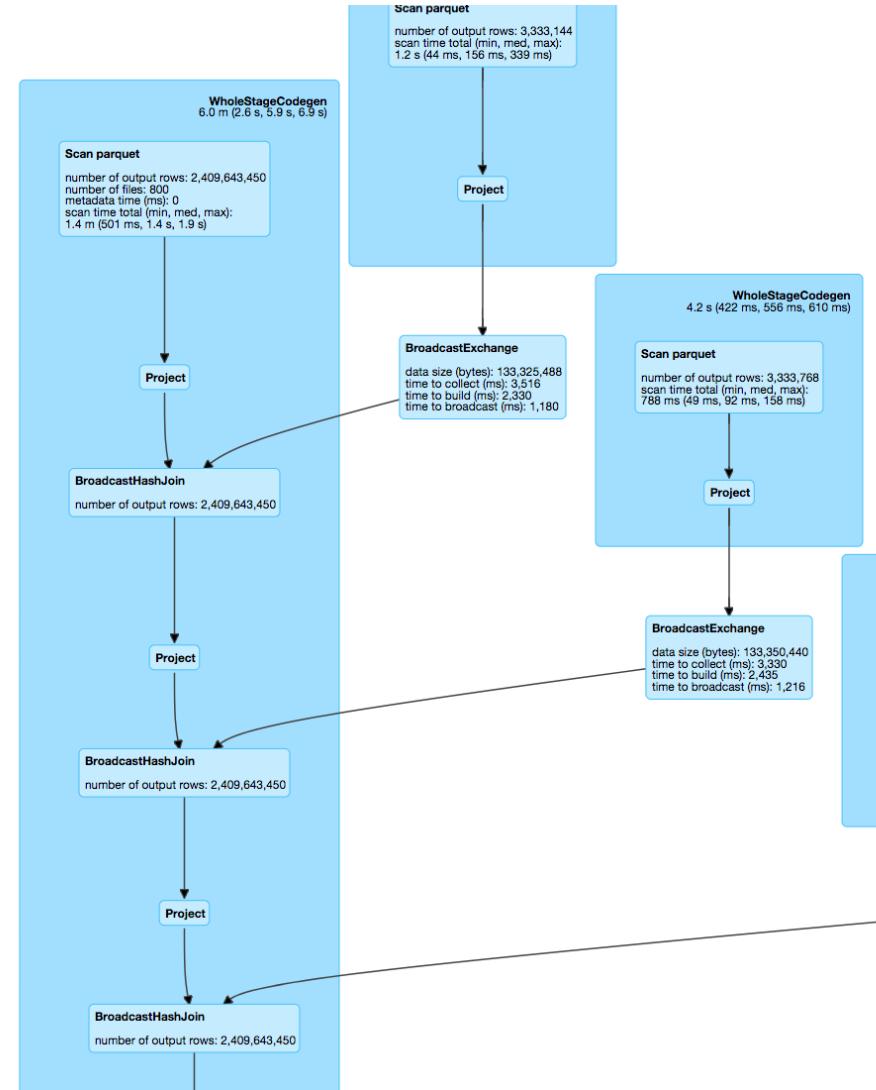


Considerations

- Hyper parameters
 - Size of the broadcast variable
 - The number of passes
- Additional complexity

Future Work

- In-memory iterations
- Native to Spark



Recap

- Skew hurts Spark parallelism and stability
- Default join types have issues with skewed data
- Iterative Broadcast Join can be used to process skewed data while maintaining parallelism

That's All Folks!

The screenshot shows a GitHub repository page for `godatadriven / iterative-broadcast-join`. The repository has 10 commits, 1 branch, 0 releases, and 1 contributor. The latest commit was made 10 hours ago. The repository is licensed under Apache-2.0. The commit history includes:

- Fokko First version of the iterative broadcast join (10 hours ago)
- project Melt all cpu cores instead of just one (2 days ago)
- src First version of the iterative broadcast join (10 hours ago)
- build.sbt (5 days ago)
- .travis.yml Melt all cpu cores instead of just one (2 days ago)
- LICENSE Initial commit (2 months ago)
- README.md Extended the generator a bit (9 days ago)
- build.sbt Got some decent testset now :-) (4 days ago)

The repository description is "The iterative broadcast join example code." and the title of the README file is "The iterative broadcast join".

https://github.com/godatadriven/iterative-broadcast-join