

BigFrame

A Benchmark Tool for Big Data Analytics

Andy

Hong Kong PolyU



THE HONG KONG
POLYTECHNIC UNIVERSITY
香港理工大學



What is BigFrame

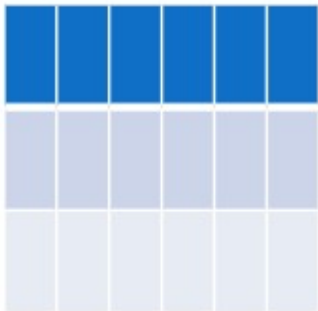
- A benchmark tool that can:
 - generate a large amount of data (PB of data)
 - generate a mixture of social network (graph), tweet (text) and e-commerce (relational) data.
 - provide benchmark workflows (hadoop implementation)
- Open source, free to use

How BigFrame helps hadoop

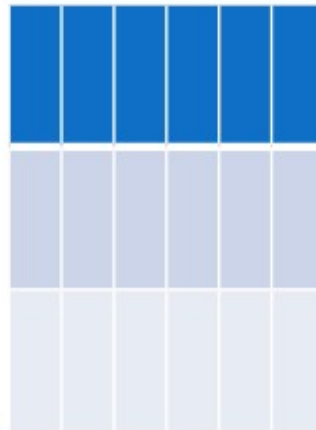
- Testing hadoop:
 - Can it handle multiple data at the same time?
 - What is the bottleneck?
 - How should we tune hadoop's system parameters?

Data Generated

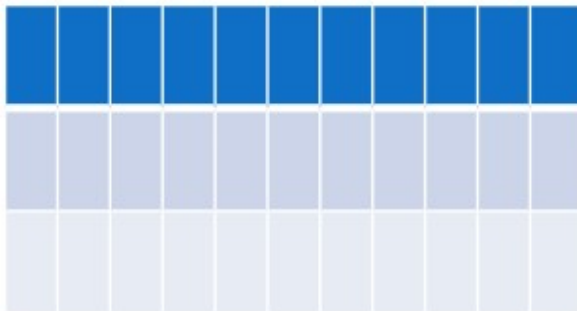
Item



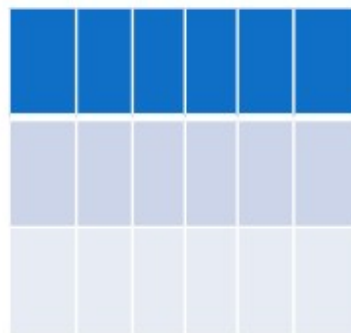
Web_sales



Customer



Promotion



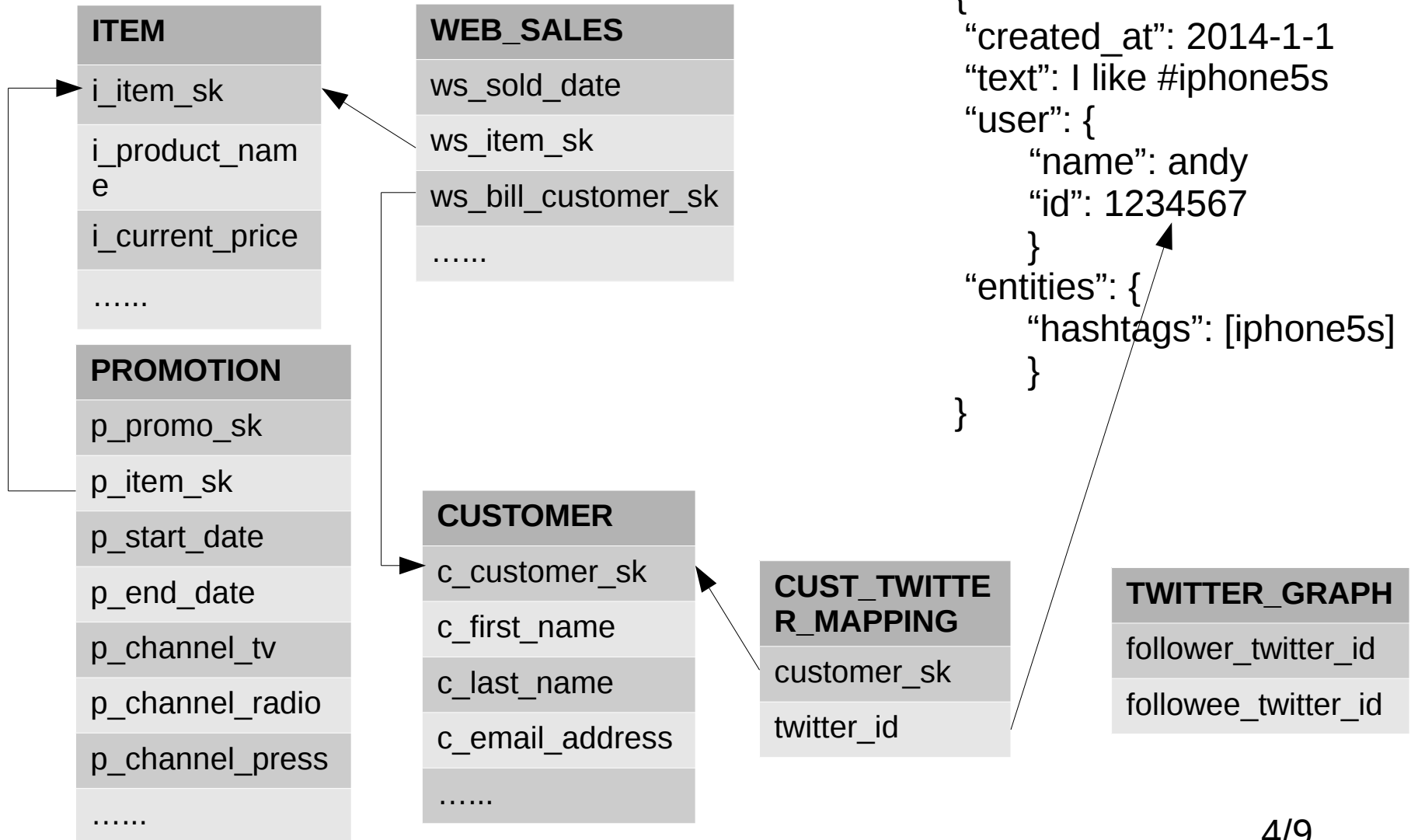
Tweets



Relationships



Data Schema

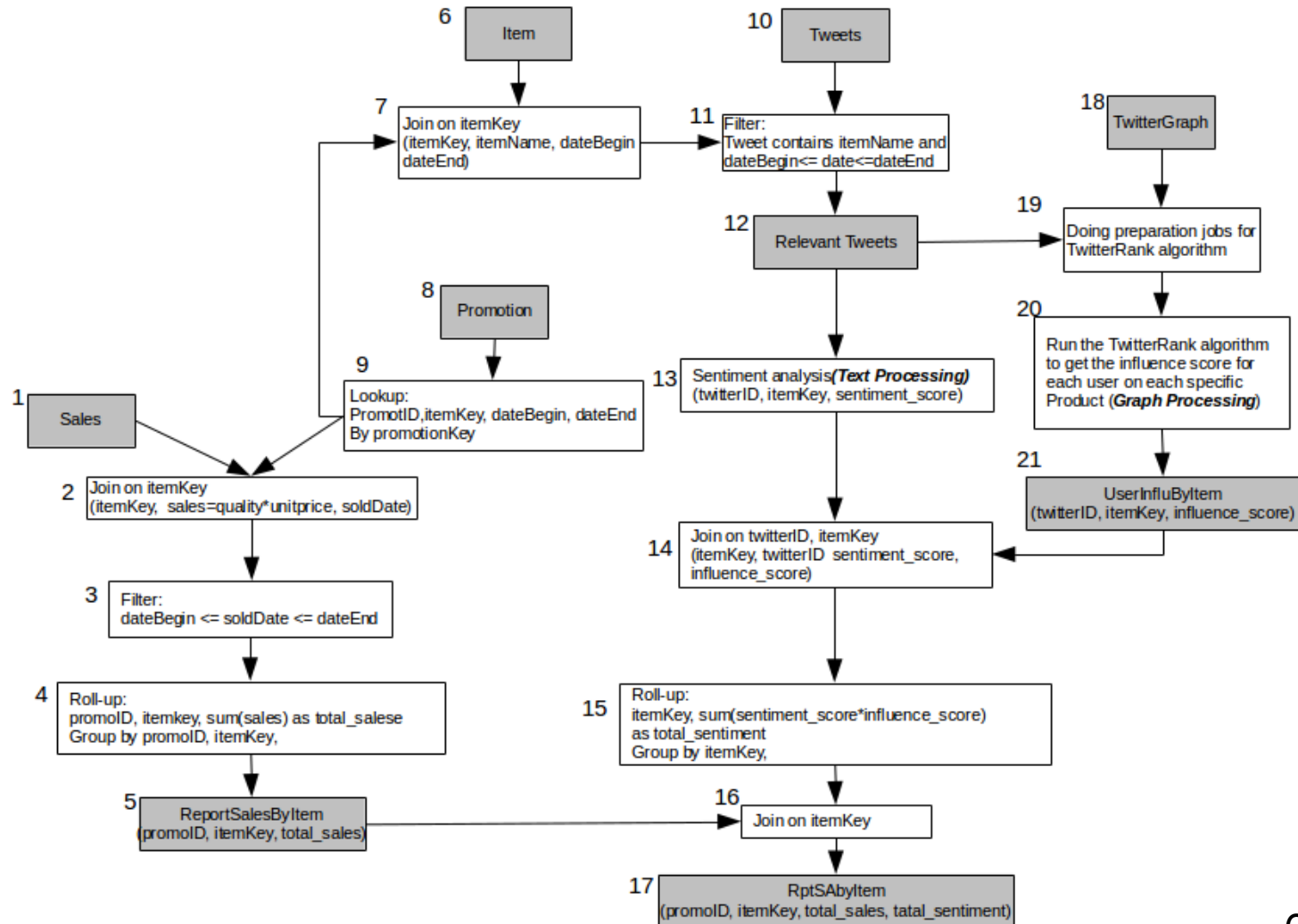


Benchmark Workflow

- The benchmark workflow relates the sales report to the sentiment analysis of all promoted products after a set of promotions is ended
 - The example result is:

promo_sk	item_sk	total_sales	sentiment_score
1	10	10,000	45
2	30	100	-67

Benchmark Workflow



Install BigFrame

- Download it from GitHub:
 - git clone <https://github.com/bigframeteam/BigFrame>
- Run the following command in bigframe's root directory to compile it:
 - *sbt/sbt assembly*

Generate the benchmark data

- Specify the data size by setting the value of property “***bigframe.datavolume***” in the “***conf/bigframe-core.xml***” file
- Run the command in the root directory of bigframe:
 - *./bin/datagen -mode datagen*

Run the benchmark workflow

- The implementation is provided: [workflow](#)
- Run the command in the root directory of bigframe to start the benchmark:
 - *./bin/qgen -mode runqueries*