



Parallel, Functional & Streaming Programming with Scala

Introduction to STREAM COMPUTING

Map reduce



The overall MapReduce word count process

Map reduce

Good

- Hadoop distributed file system (HDFS)
- Distributed & massively parallel
- Move the compute to the data
- Accomplish what was impossible before

Not so good

- Finite data set
- Batch process
- Begin to chain jobs together
- Failure?
- Recovery?
- Idempotent?

Workflows

0ozie

Luigi

Azkaban

Airflow

Pinball

Cascading

Taskflow





Streams are all around us

Stock market trading

Social media

Sensors

Weather

Transportation traffic

Video



The opposite of stream is BATCH processing

Has a beginning and an end

The data is time bounded and finite

Example:

- Sales from a single store in a month
- One day trading on a single stock exchange

32			33
Covelle 191	. Batter 1911	almen 1811 Besc	
that II for for	11 the let & for fet "	In 24 for farm 5 and file 13 year far	7.00
ale I for March	11 10 may 5 for may 12	fet a for gove 12 a for the to goo	
May 1 for afere	14 as for 1 for for 12	feb 20 for feb 10 or Mor 11 for for	1 5 50
may 28 for may	1400 July 3 fe July 10	Dan 9 for fet 2 2 - af 2 for for	2 175
Kaly 29 So Cales	Gon sque la sep- 10	marge de tra por may 1 for mos	en 5 pr
thing I do July	(you ear 5 fee ear 12)	Amiso for the lite fight for the	- 150
ang 21 for ang	(Tot men 5 for man 12	ap 22 for a for for going to ap may to	ne 4000
all I by with	The out of the dec 1 his	and it do the barrie land to be	a 250
ort- bit be toll	(9 m	This of far may 550 and ? for for	iy 10.00
vr21 fr at	6 m 1912	for 17 for fine 1 to set 9 for an	\$ 1000
nor 3 for set	(Toc Lowello E	for 29 ste cam stor della	1112 60
der II de ner	(Seven 9 la nea 9	29 la luce (500 ments do as	r 11. Ce
der 3' for der	S 1 for 15 for dec 7.	any 1 to goly 150 de 1 for m	ac 12 50
3m 7 for dec	Profo 25 for dec 7	agel for eng Box for 4 for de	2 12 50
1 1 1 0	14	14 1 (2.50 H / 6	12 12
me to let	14	suto do sup how me for the	1 1050
ap forma	1.*	1430 to set Sor up 3 to m	a 1210
my 4 da apr	7 e-	or I for white may a for an	1250
may 11 for all	6.00	and a not Boot Supply the Pro-	~ 12 50
a may	. 5.11	now I be set 500 and I for for	4 1200
San So for for	7 -	must for not 2th any 1 to an	12 12 50
July 12 for fore	6.44	and to may Set and 2 to a	21-12.50
. A G Sula	8-04	les : for mon 5 on . She ? for 1.	101250
by the ang	6	dec 9 for now 2 to All a for ch	2 1210
All fe an	.S.a.	all to the Sec god of the g.	
sorts to set	5 00 C 00	14 23 ander Son	
man 2 la sett	400		
Moo 16 for oct	5.00		
12. 22 Lorald	2.64		

Stocks are ticking somewhere



The real world never shuts up, never pauses, never stops.

So how do we program it?

Source, sink, operator



Directed acyclic graph (DAG)



Cyclic Bad!



Streams are different

- 1. Data is unbounded or infinite
- 2. Continuous processing
- 3. There is no <u>now</u>
- 4. Eventual consistency vs false sense of consistency
- 5. Closer to reality

#1 Unbounded



#2 Batch vs stream





Credit: Tyler Akidau.

Calculate closing (end of day) balance

Batch

Start with yesterday's closing balance

Add up all the events for the day

Publish a new closing balance

Stream

Recalculate with each new event

Emit new balance with each event Tag a balance event as closing balance

#3 Time

Event time:

the time at which events actually occurred

Ingest time:

the time the event is written to the system

Process time:

the time at which events are observed in the system

Time skew



Credit: Tyler Akidau.

When is now?





#4 Consistency

Trade data arrives at end of day (EOD)

Processing runs to create EOD status of trades

Corrections exist for previous days

Previous EOD is also changed

Batched processes give us a false sense of consistency

Eventual consistency



Eventual consistency

Applies to organizations too

Eventually all the downstream processes will have consistent values

Batch:

Around 6 pm all the systems have the same value, but it is probably wrong

Stream:

At 6:00:00.0 pm all systems have a value, but perhaps not the exact same value

Considerations

- 1. Message delivery
 - a. At most once
 - b. At least once
 - c. Exactly once
- 2. Windowing
 - a. Fixed
 - b. Sliding
 - c. Session
- 3. Joins
 - a. Inner
 - b. Left
 - c. Outer
- 4. Scaling

At most once delivery

"Best effort" approach

Messages may be lost



At least once delivery

"Guaranteed delivery"

Duplicates will occur



Exactly once delivery

Expensive!

- Distributed snapshot/state checkpointing
- At-least-once event delivery plus message deduplication



Hurst's Law

Complexity can neither be created nor destroyed; it can only be displaced.

Pay attention to where it went!

Fixed window

Time is partitioned into same-length, non-overlapping chunks.

Each event belongs to exactly one window



Window by process time



Credit: Tyler Akidau.

Window by event time



Credit: Tyler Akidau.

Session window

Sequences of events terminated by a gap of inactivity greater than some timeout



Window into session by event time



Credit: Tyler Akidau.

Sliding window

Fixed length and fixed period

Eg, the last 10 seconds reviewed every 2 seconds



Process vs event time


03:06 process time

















"All input data with event times less than X have been observed"

Heuristic

Best guess, never perfect

Crossing Joining the streams





Customer data with web traffic data

Securities data with trading data

Weather data with transportation traffic data

Two or more sensors

Join data

. . .

- 1000001: 12501 : 45.0
- 1000009: 12508 : 45.1
- 1000012: 12406 : 45.3
- 1000013: 12518 : 45.2
- 1000021: 12508 : 45.4
- 1000007: 12501 : 45.2
- 1000026: 12409 : 45.5
- 1000029: 12402 : 46.0
- 1000035: 12502 : 46.4

- 1000004: 12501 : S.L.RW
- 1000005: 12502 : ..FW
- 1000005: 12503 : S.L.RW
- 1000018: 12504 : S.L.RW
- 1000021: 12505 : S.L.RW
- 1000021: 12401 : ..FW
- 1000022: 12402 : M..G
- 1000025: 12403 : M..G
- 1000027: 12404 : S.L.G

• . . .

Join data

- 1000001: 12501 : 45.0
- 1000009: 12508 : 45.1
- 1000012: 12406 : 45.3
- 1000013: 12518 : 45.2
- 1000021: 12508 : 45.4
- 1000007: 12501 : 45.2 •
- 1000026: 12409 : 45.5
- 1000029: 12402 : 46.0
- 1000035: 12502 : 46.4 •

- 1000004: 12501 : S.L.RW
- 1000005: 12502 : ..FW

- 1000005: 12503 : S.L.RW
- 1000018: 12504 : S.L.RW
- 1000021: 12505 : S.L.RW
- 1000021: 12401 : ..FW
- 1000022: 12402 : M..G
- 1000025: 12403 : M..G
- 1000027: 12404 : S.L.G

. . .

• . . .

Input stream example



Inner join

Emits an output when both input sources have records with the same key.



INNER STREAM-STREAM JOIN







Left join

Emits an output for each record in the left or primary input source. If the other source does not have a value for a given key, it is set to null.









Outer join

Emits an output for each record in either input source. If only one source contains a key, the other is null.







Scaling

۴.,

Parallelism



Fundamental challenge of parallel processing?

Two approaches to parallel processing?

Partition by source



Partition round robin



Partition by key



Keys

Types of keys

- 1. Natural key such as a place name
- 2. Generated key (aka surrogate key)
- 3. Composite key

One man's generated key is another man's natural key

Never create keys with metadata!



Universally unique identifier

Sequence identifiers will NOT work in a distributed environment

128-bit

4fedefdb-4c7c-42b5-ae00-b6d286034b2c

Collision: generate 1 billion UUIDs per second for about 85 years

https://www.uuidgenerator.net/

Partition by hash-mod

hash(key) % partitions

hash("Amsterdam") % 10

8763125 % 10

5

























Partitions of partitions

Initial

- 1. 0..9
- 2. 10..19
- 3. 20..29
- 4. 30..39

Expanded

- 1. 0..6
- 2. 10..16
- 3. 20..26
- 4. 30..36
- 5. 7..9,17..19
- 6. 27..29,37..39
Real world example

System example



Sample rates



Legacy processing



Legacy batch processing



Stream processing

Time	Sensor reading
2456	Latitude, longitude, altitude
2635	
2762	Pitch, roll, heading
2789	
2812	Latitude, longitude, altitude
2893	
2910	Pitch, roll, heading
2998	
3015	Latitude, longitude, altitude
3075	
3147	Pitch, roll, heading
3222	

Sample rates



Stream processing

 $G(2789) = G_{2762} + ((2789-2762)/(2910-2762))G_{2910}$

	Time	Sensor reading
	2456	Latitude, longitude, altitude
	2635	
\langle	2762	Pitch, roll, heading
	2789	
	2812	Latitude, longitude, altitude
	2893	
\langle	2910	Pitch, roll, heading
	2998	
	3015	Latitude, longitude, altitude
	3075	
	3147	Pitch, roll, heading
	3222	

THANK YOU!



Calculate the average

• . . .

- 100001: 45.0
- 100009: 45.1
- 100012: 45.3
- 100013: 45.2
- 100021: 45.4
- 100007: 45.2
- 100026: 45.5
- 100029: 46.0
- 100035: 46.4

. . .

Process vs event time

