# REACTIVE STREAMS, linking REACTIVE APPLICATIONS to SPARK STREAMING

Luc Bourlier

Lightbend Inc.

# Agenda

- Spark Streaming

- Reactive Application

- Back pressure

- Reactive Streams
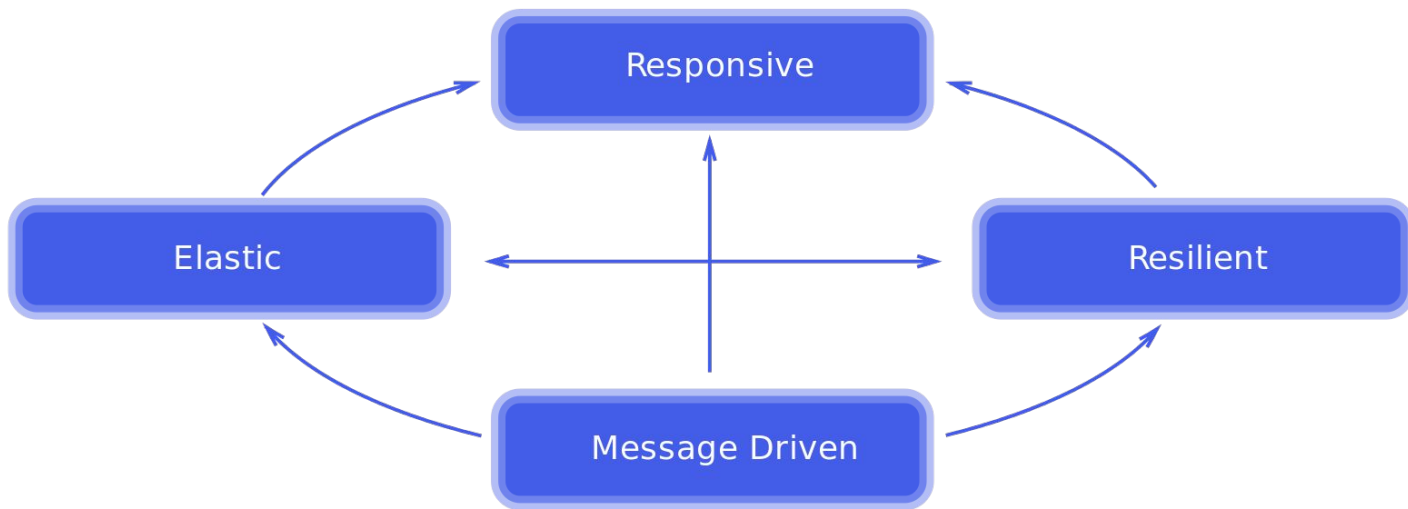
- Demo

# Spark Streaming

Lightbend

# Spark Streaming

# Reactive Application

# Reactive Application



http://www.reactivemanifesto.org

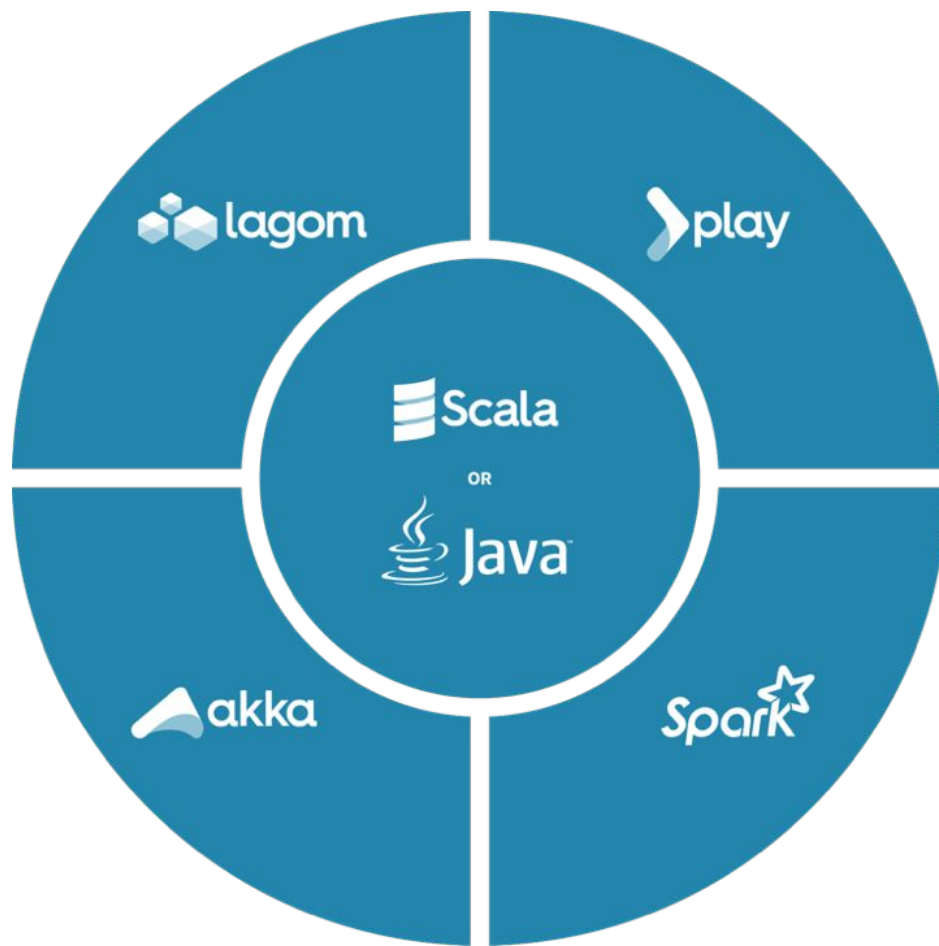# Reactive Application

Responsive                             responds in a timely manner

Resilient                    stays responsive in the face of failure

Elastic              stays responsive under varying workload

Message Driven                          relies on asynchronous
                                                    message-passing

# Resilience in Spark and Spark Streaming

- Support for all kinds of failures
  - Hardware
  - Software
  - Network

- Specific resilience for Spark Streaming
  - Recovery for continuous processing
  - Excess volume of data

# Resilience in Spark and Spark Streaming

- Support for all kinds of failures
  - Hardware
  - Software
  - Network
- Specific resilience for Spark Streaming
  - Recovery for continuous processing
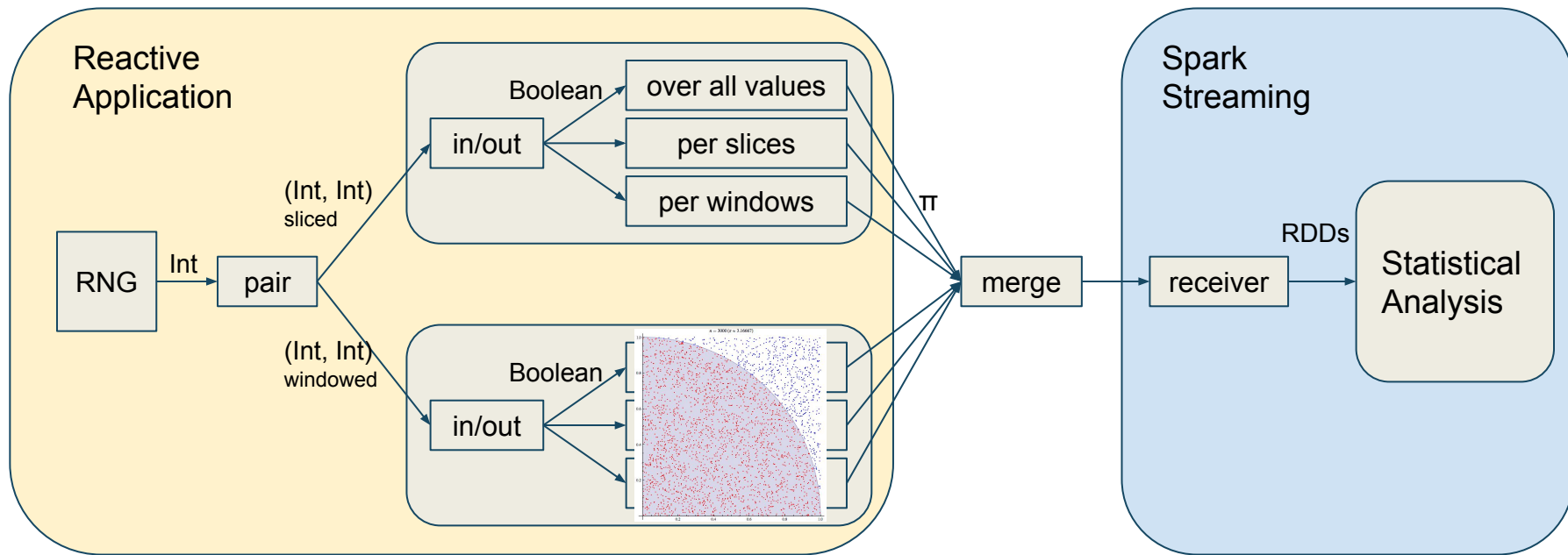  - Excess volume of data ← **the subject of this presentation**

Lightbend

# Demo

Lightbend

# Demo

# Back Pressure

Lightbend

# Back Pressure

- a slow consumer should slow down the producer
  - the produce applies pressure
  - the consumer applies back pressure
- the classic example: TCP

# Back Pressure in Spark Streaming

Lightbend

# Spark Streaming

Congestion support in Spark 1.4

Static rate limit

- `spark.streaming.receiver.maxRate`
- conservative
- difficult to find the right limit (depends on cluster size)
- one limit to all streams

Lightbend

# Spark Streaming

Back pressure in Spark 1.5

Dynamic rate limit

- rate estimator
  - estimates the number of element that can be safely processed by system during the batch interval
- rate sent to receivers
- rate limiter
  - relies on TCP to slow down producers

# Spark Streaming

Rate estimator

- each BatchCompleted event contains
  - processing delay, scheduling delay
  - number of element in mini-batch
- the rate is (roughly) elements / processingDelay
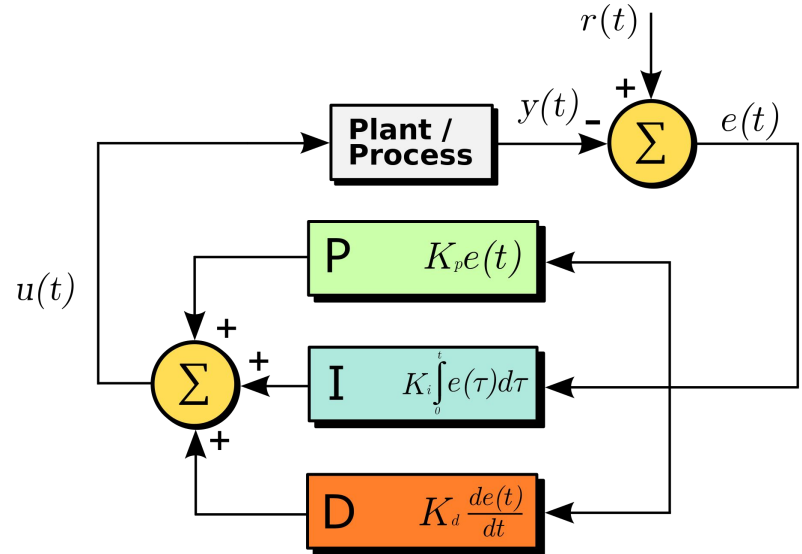- but what about accumulated delay?

Lightbend

# Spark Streaming

## Rate estimator

## Proportional-Integral-Derivative

- P, I, D constants change
  convergence, overshooting
  and oscillations



https://en.wikipedia.org/wiki/PID_controller
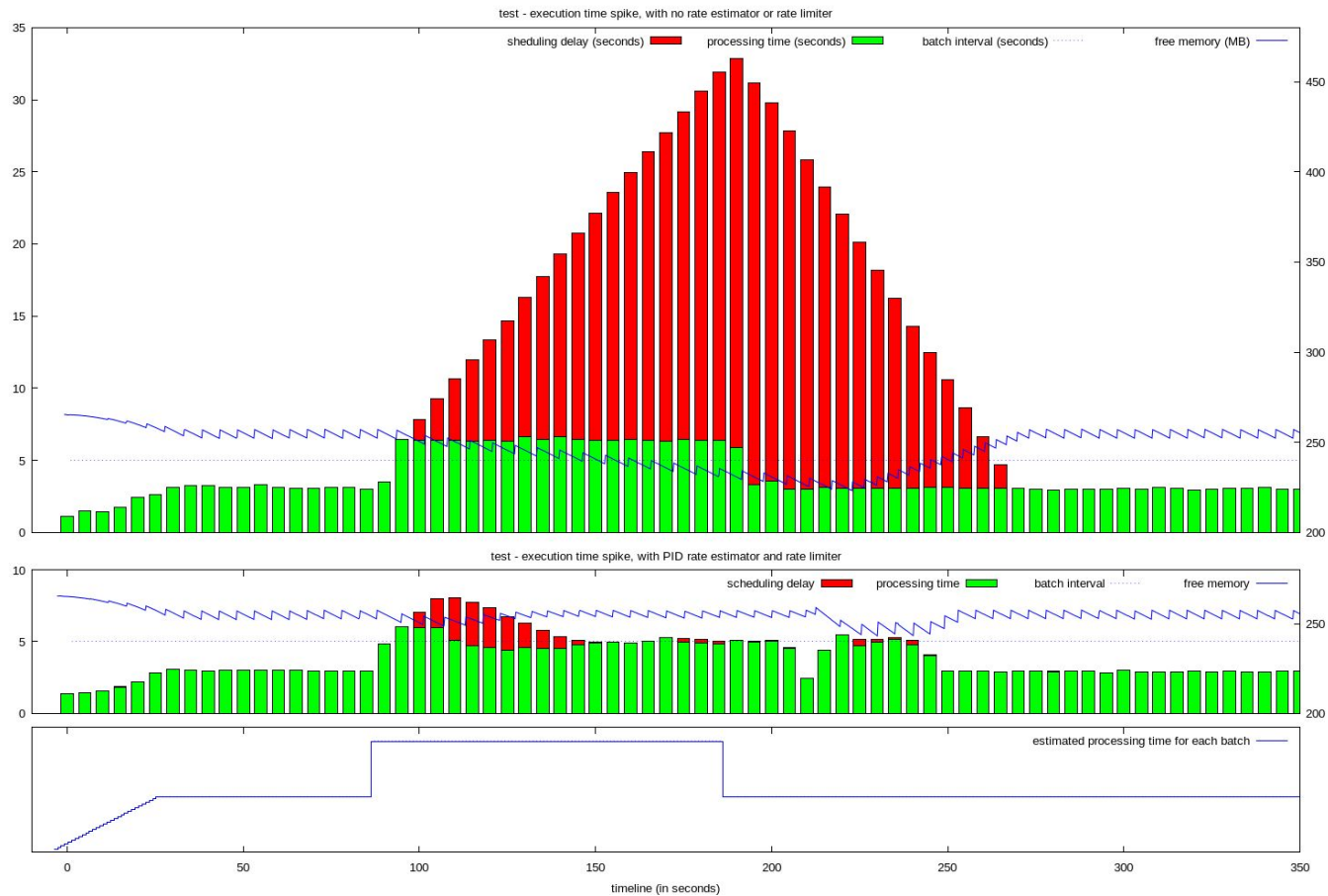
# Spark Streaming

Back pressure in Spark 1.5

- each input has its own estimator
- work with all stream receivers

  including KafkaDirectInputStream

- configuration
  - `spark.streaming.backpressure.enable    true`
  - `spark.streaming.backpressure.minRate    R`

test - execution time spike, with no rate estimator or rate limiter

test - execution time spike, with PID rate estimator and rate limiter

timeline (in seconds)

# Spark Streaming

Limitations

- linearity assumption

- records with similar execution times

- TCP back pressure accumulates in the TCP channel

Lightbend

# Reactive Streams

# Reactive Streams

- one tool to create reactive applications

- specification for back pressure interface to connect systems supporting back pressure in the JVM

  - small: 3 interfaces, 7 methods total

- subscriber controls rate by requesting elements from producers

http://www.reactive-streams.org

Lightbend

# End to end back pressure

# End to end back pressure

- Reactive application with reactive streams connector

  $\Rightarrow$ back pressure enabled

- Spark Streaming 1.5+

  $\Rightarrow$ back pressure enabled

- Reactive streams Spark Streaming receiver

  $\Rightarrow$ end to end back pressure

Lightbend

# Demo

Lightbend

# Spark 2.x ?

- Spark Streaming still available
  - same support

- Structured Streaming
  - experimental, no stable source API
  - different model
  - requires an updated solution

Lightbend

# THANK YOU.

luc.bourlier@lightbend.com