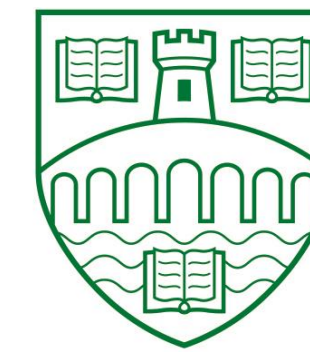# Continuous Applications using Apache Spark

## Online Machine Learning with Streaming Data

Robert Hamlet

MSc Big Data

Supervised by Andreas Lang and Jonathan Forbes (Aquila Insight),
Dr. Kevin Swingler (University of Stirling)

## Background

### Streaming Data

Applying models to streaming data reduces 'time-to-insight', critical for time-critical applications of machine learning such as anomaly detection and e-commerce.
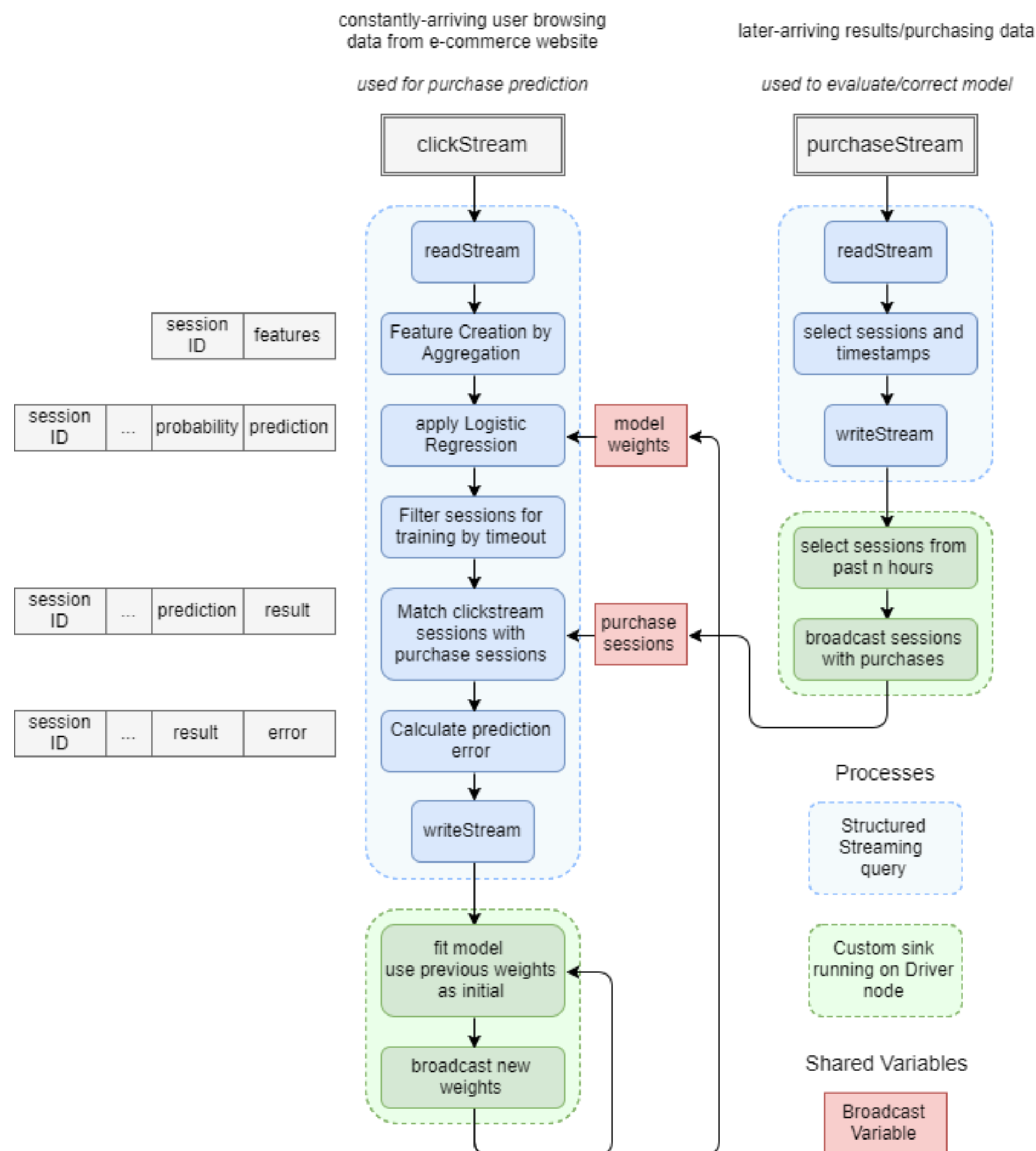
### Online Machine Learning

Continuously-adapting models can keep abreast of drifting environmental conditions, representing an improvement in mitigating model degradation over regular batch rebuilds.

### Apache Spark

Apache Spark, along with its new Structured Streaming functionality, is a rapidly-growing distributed processing platform that could potentially allow continuous architectures within a single environment.

### Scenario

This architecture has been developed using e-commerce user browsing and purchasing data from the 'RecSys Challenge 2015'. The imagined case study explored predicting the likelihood of a user making a purchase, perhaps to trigger offering threshold users a voucher to sway them.
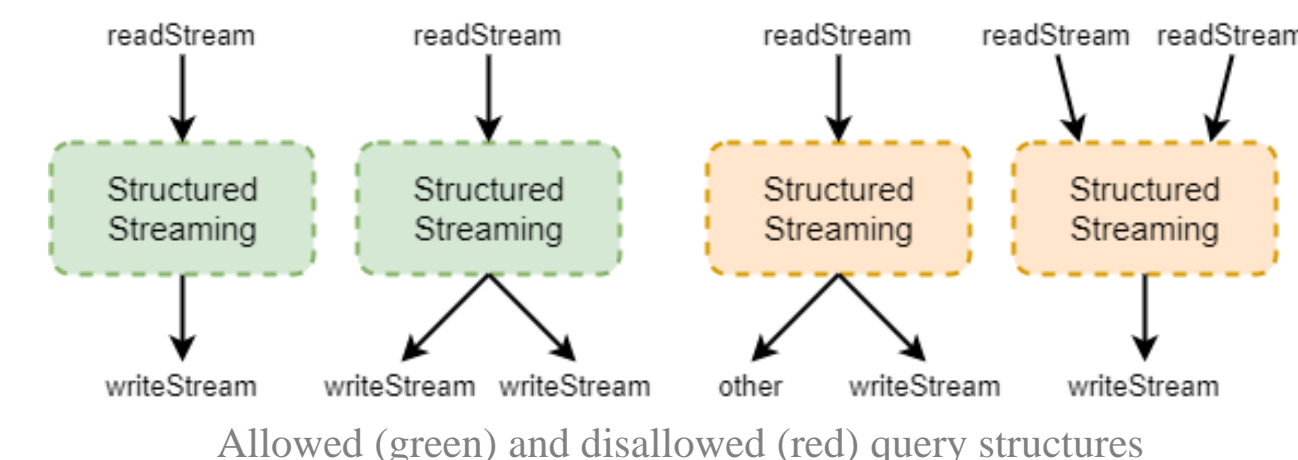


## Challenges

### Crossing the Streams

Structured Streaming queries have a rigidly defined and enclosed structure and *do not yet support stream-stream joins*. Moving data between streams can be achieved using *Broadcast Variables*, but triggering the broadcast of these *breaks query linearity* (see fig below, queries must end with a dataframe, and a writeStream).

### Model Training

The new spark.ml machine learning library *does not yet have a streaming logistic regression model*, and a model transformer *cannot be re-trained* at each trigger of a Structured Streaming query.



Allowed (green) and disallowed (red) query structures

## Solution

Solving both challenges could be possible by performing the model training and variable broadcast outside of the query and its restrictions by using a *custom sink*.

This could train a new model each time it is triggered, without the constraint of having to output a dataframe by writeStream.

Since this can be executed on the driver node it will have access to the SparkContext (necessary to broadcast variables), although this introduces a performance bottleneck.