

Arbitrary Stateful Processing



Objective

Aggregate data in an arbitrary way

Manage stream data "manually"



Scenario

You're running a social network and people share updates

Your social network needs to store everything

You want to compute the average storage used by every type of post

```
text,3,4812
video,1,5328585
text,4,5326
audio,5,523662

// type of posts, count of posts, total storage used
```

Average storage = total storage / total count, grouped by post type

Takeaways

Define a function to update state based on new records

```
socialPosts
    .groupByKey(_ postType) ← group first
    .mapGroupsWithState(GroupNameTimeout.NoTimeout())(updateAvgStorage)
    .writeStream
    .format("console")
    .outputMode("update") ← append not supported
    .start()
    .awaitTermination()
```

```
def updateAvgStorage(
    postType: String, ← key by which the records are grouped
    group: Iterator[SocialPostRecord], ← new records to aggregate
    state: GroupState[SocialPostBulk] ← current state
): AveragePostStorage = { ← final value to be returned for the whole group
    // your code here
}
```

Also flatMapWithState when a group produces multiple outputs

Spark rocks

