

a stripe talk

Aggregator: MapReduce in the type system

by Dan Frank | @danielhfrank

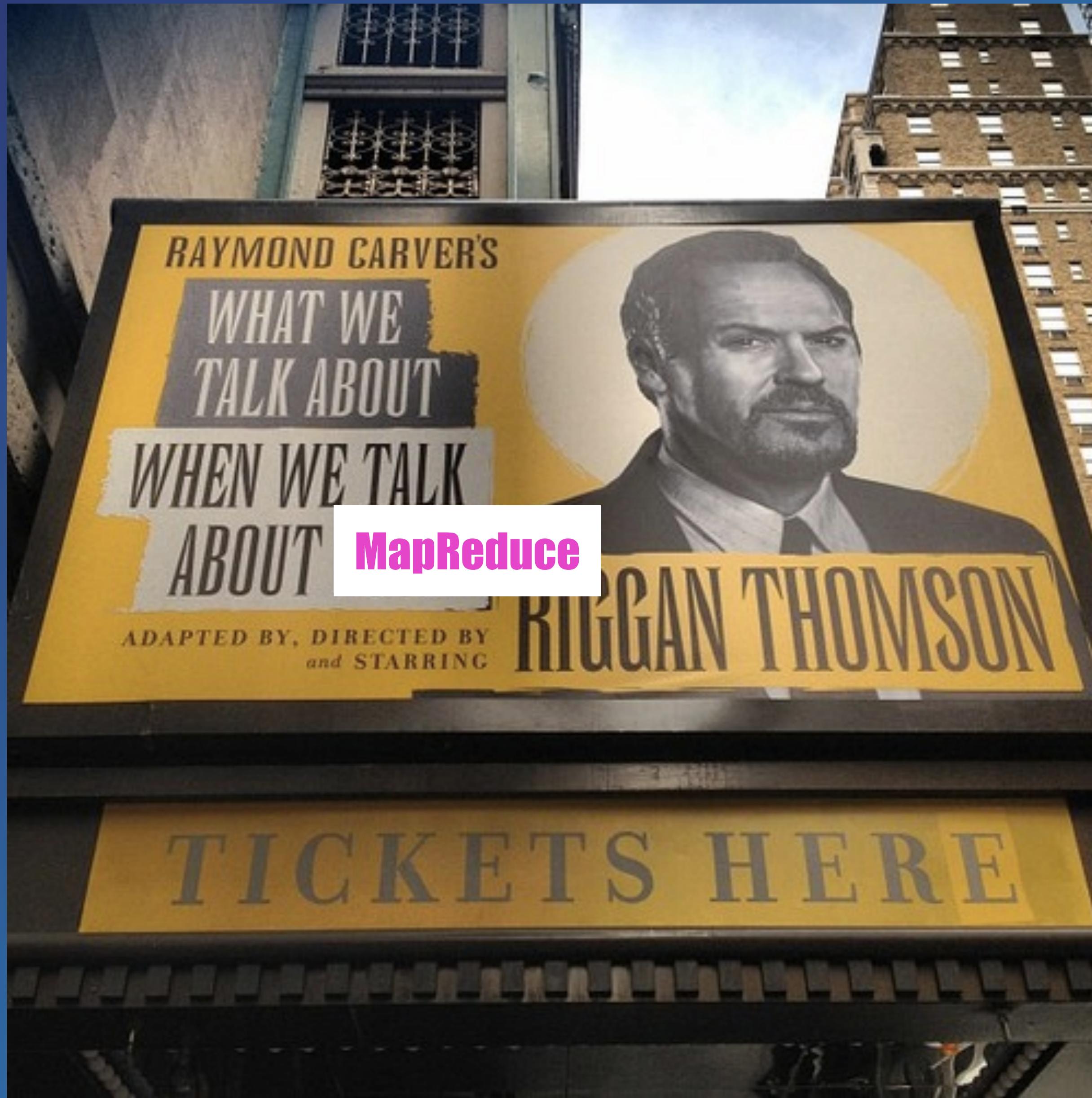
stripe

stripe

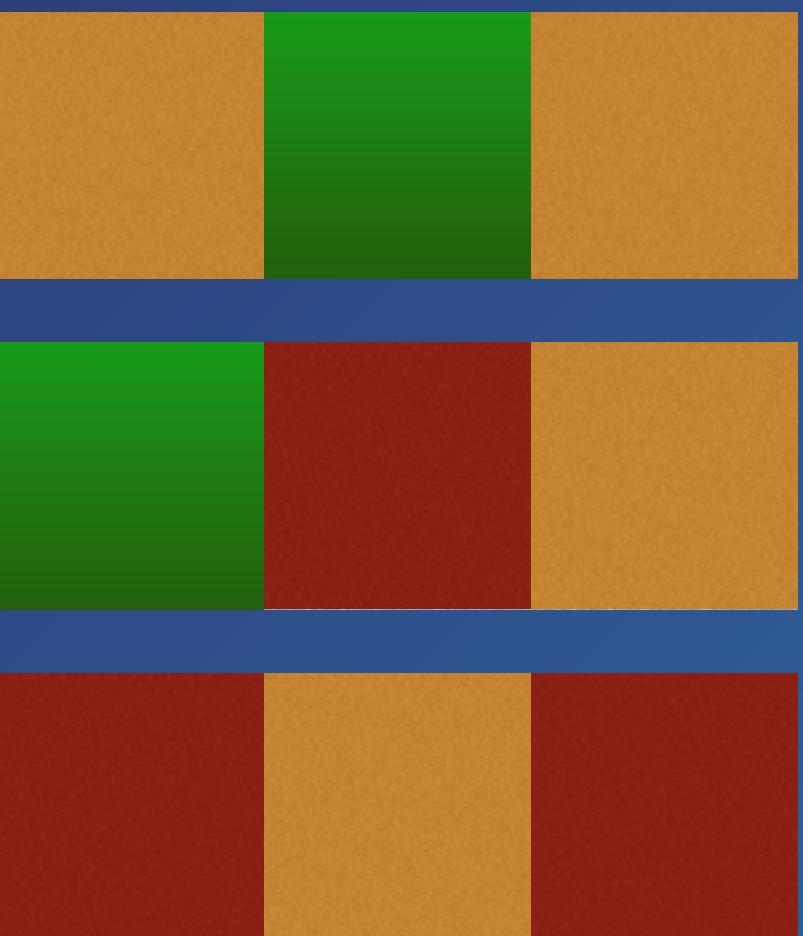
Data at Stripe

- Units: cards, users, customers, transactions
- Product usage analytics for business / product development, machine learning for fraud prevention, and much more
- Today: Use Aggregator on this data to predict fraud

stripe



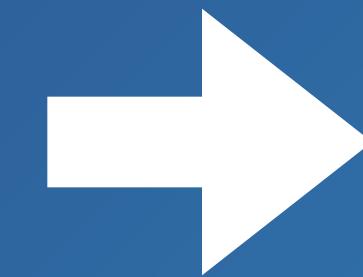
Raw Data



Properties



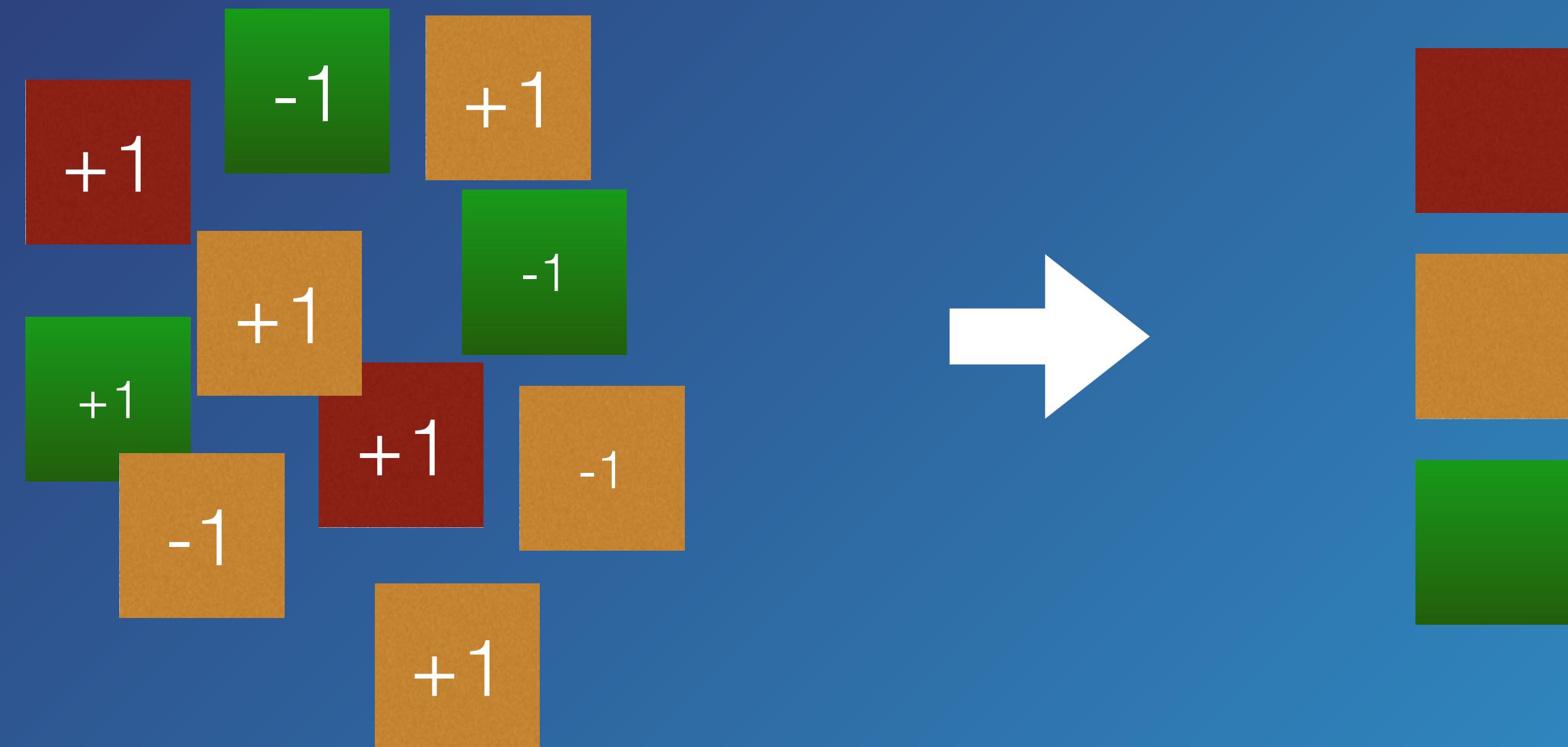
Further Reduction



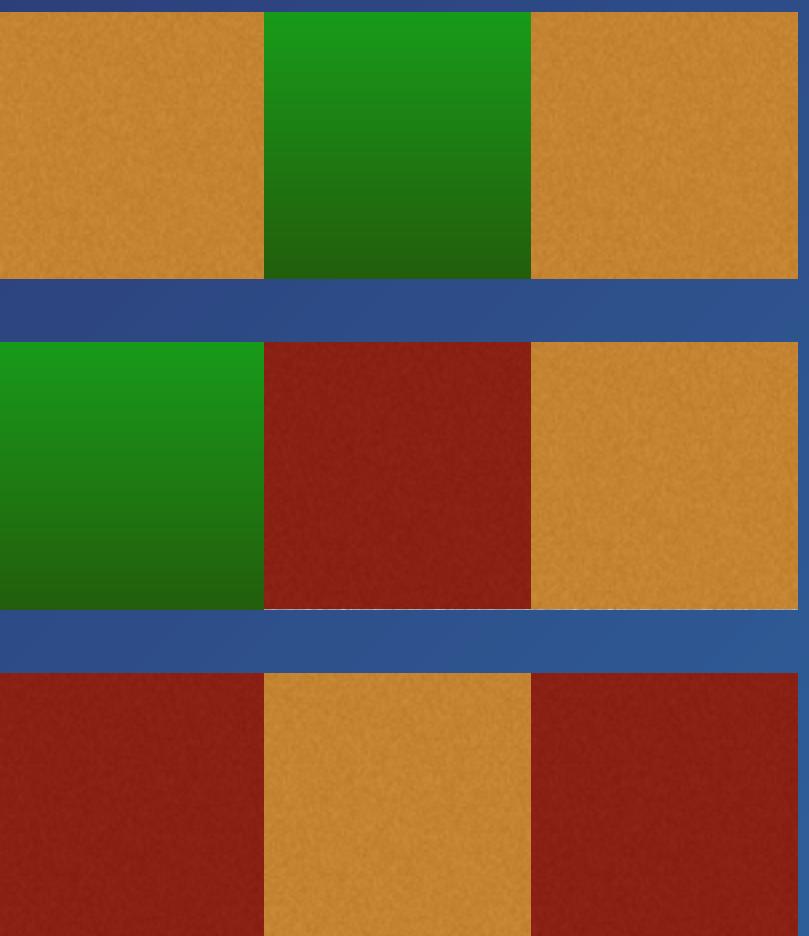
Max = 4
Mean = 3
Std = 0.82



Example: Comment Ranking



Raw Data

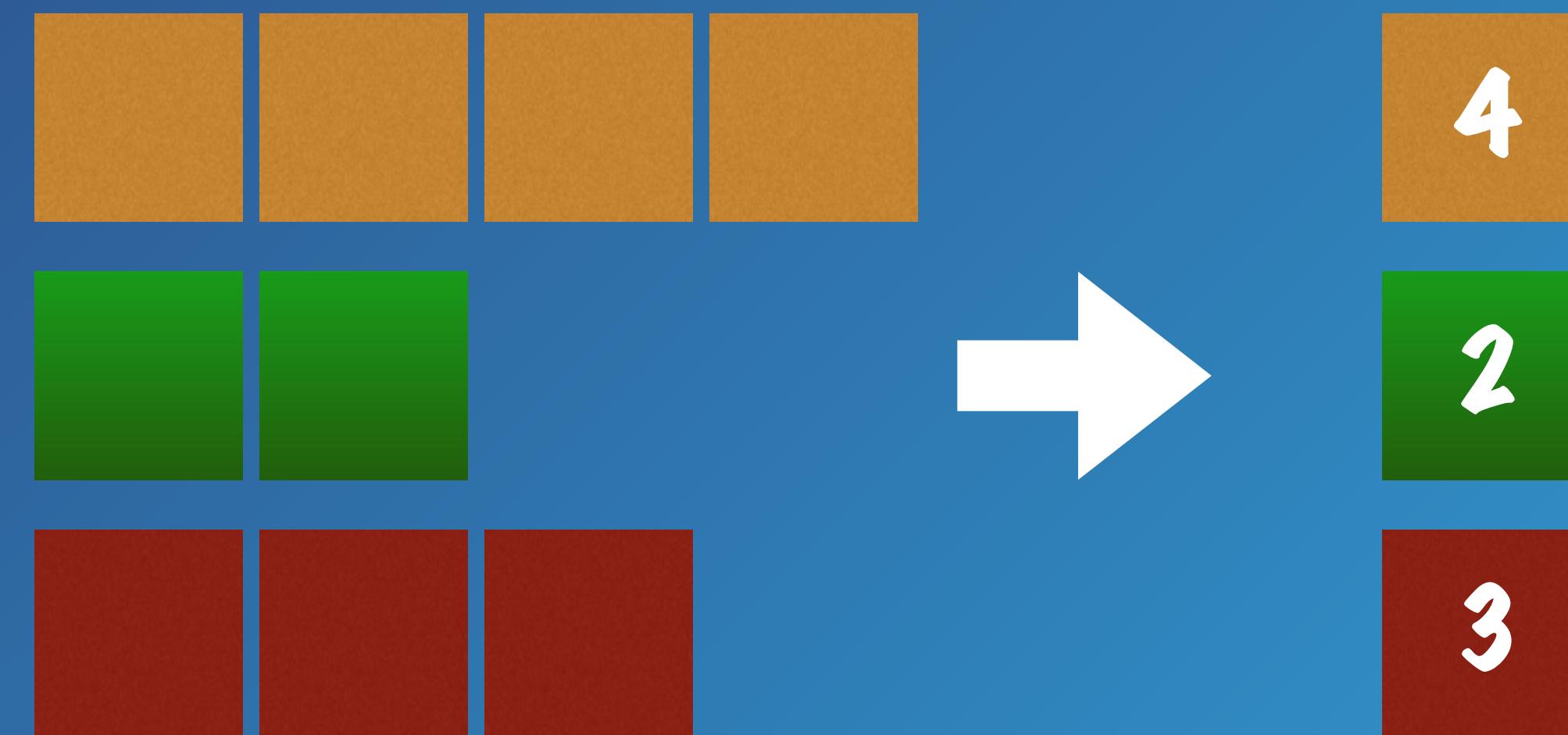
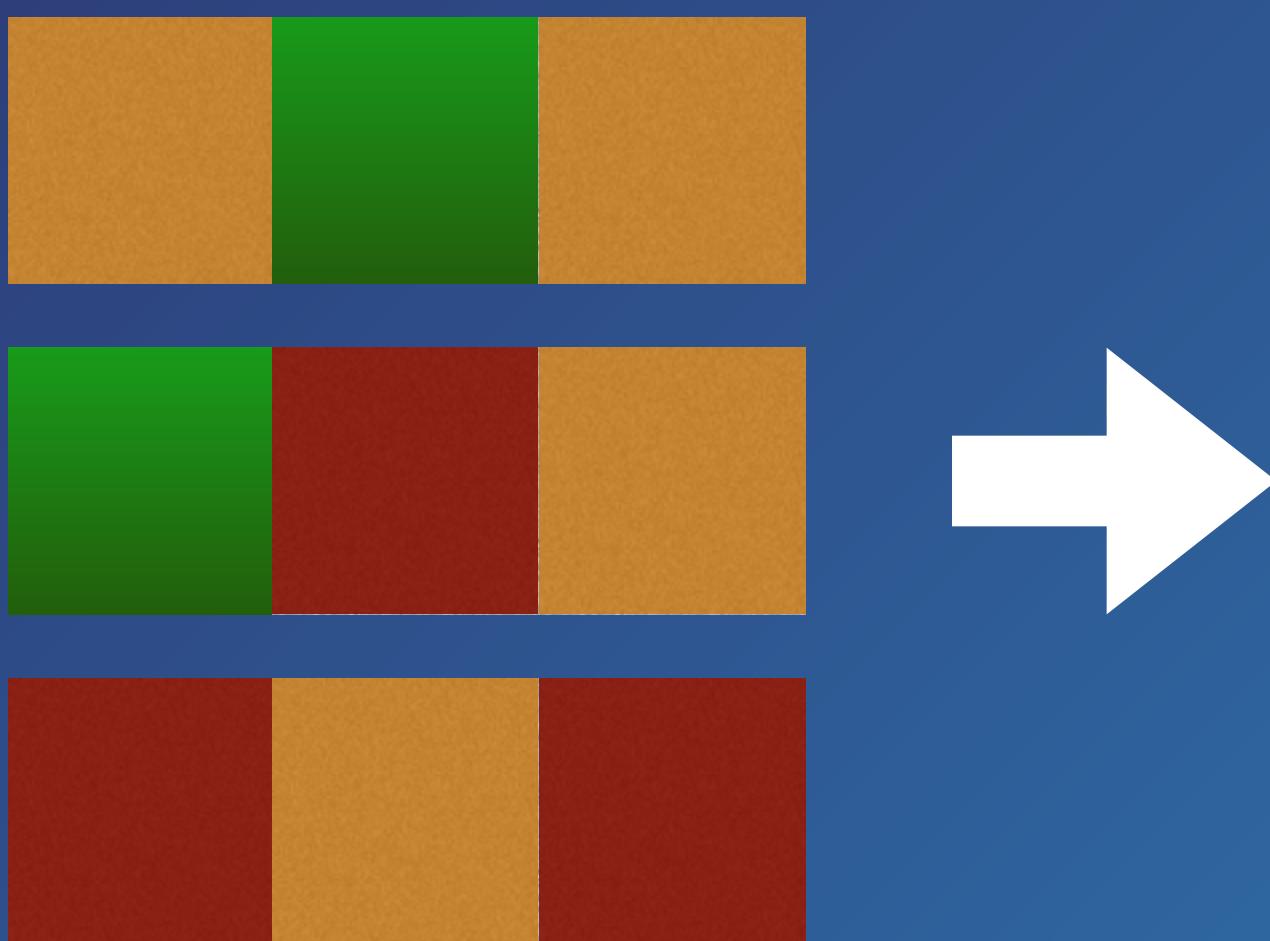


Properties



Example: How many charges per merchant?

```
1 SELECT merchant, COUNT(*)  
2 FROM charges  
3 GROUP BY 1|
```





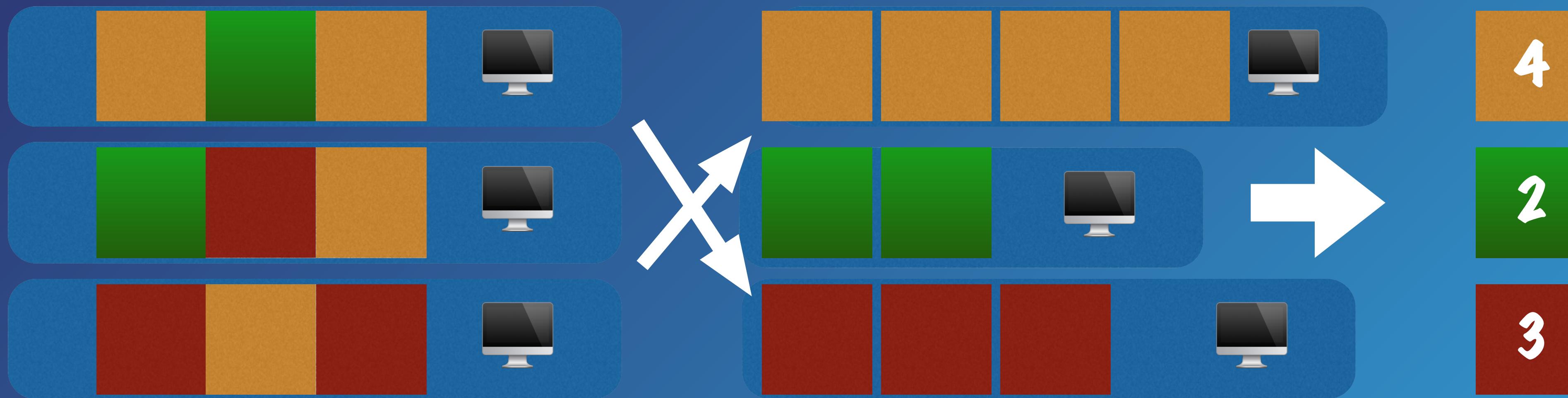
map($K_1, V_1 \Rightarrow (K_2, V_2)$
reduce($K_2, Seq[V_2] \Rightarrow (K_3, V_3)$)



Example: How many charges per merchant?

```
1 def map(charge):  
2     yield (charge['merchant'], 1)  
3  
4 def reduce(key, values):  
5     yield (key, sum(values))
```

Why MapReduce?





Breaking the decomposition boundary



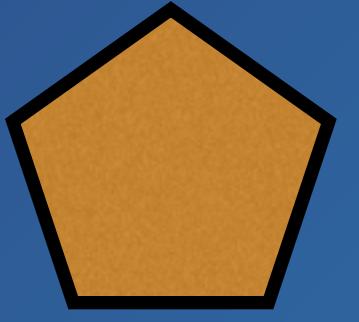
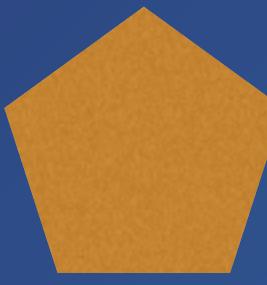
(Orange, 1)



The "real" components

- Group data points by entities of interest (keys)
- Reduce a sequence of data points to a quantity of interest

stripe



3



Example: How many cards per merchant?

```
1 SELECT merchant, COUNT(DISTINCT(card))  
2 FROM charges  
3 GROUP BY 1
```



Example: How many cards per merchant?

```
1 def map(_, charge):  
2     yield (charge['merchant'], charge['card'])  
3  
4 def reduce(key, values):  
5     yield (key, len(set(values)))
```

Aggregator

- Scala trait in Algebird: "Abstract Algebra for Scala"
- Use types to guide composition and reuse
- Agnostic to execution platform



Example: How many cards per merchant?

$\text{Seq}[\text{Charge}] \Rightarrow \text{Int}$



Computation type

$\text{Seq}[\text{Charge}] \Rightarrow \text{Set}[\text{Card}] \Rightarrow \text{Int}$

```
val chars = List('S', 'C', 'A', 'L', 'A')

chars.map(_ => 1).reduce(_ + _) // => 5

chars.map(Set(_)).reduce(_ ++ _).size // => 4
```

Aggregator Components

Aggregator[A, B, C]

prepare: A => B

reduce: (B, B) => B

present: B => C

stripe

A

B

C

Example: How many cards per merchant?

```
1 class MerchantCardCount extends Aggregator[Charge, Set[Card], Long] {  
2     def prepare(a: Charge) = Set(a.card)  
3  
4     def reduce(b1: Set[Card], b2: Set[Card]) = b1 + b2  
5  
6     def present(b: Set[Card]) = b.size  
7 }
```

The Next Level

- Composition, reuse, and the "standard library"
- Aggregator execution patterns
- Aggregator and fraud at Stripe

Leveraging the standard library

Aggregator.uniqueCount[Card]

[Card, Set[Card], Int]

composePrepare

```
Aggregator.uniqueCount[Card]  
  .composePrepare(a: Charge => a.card)
```

[Charge, Set[Card], Int]

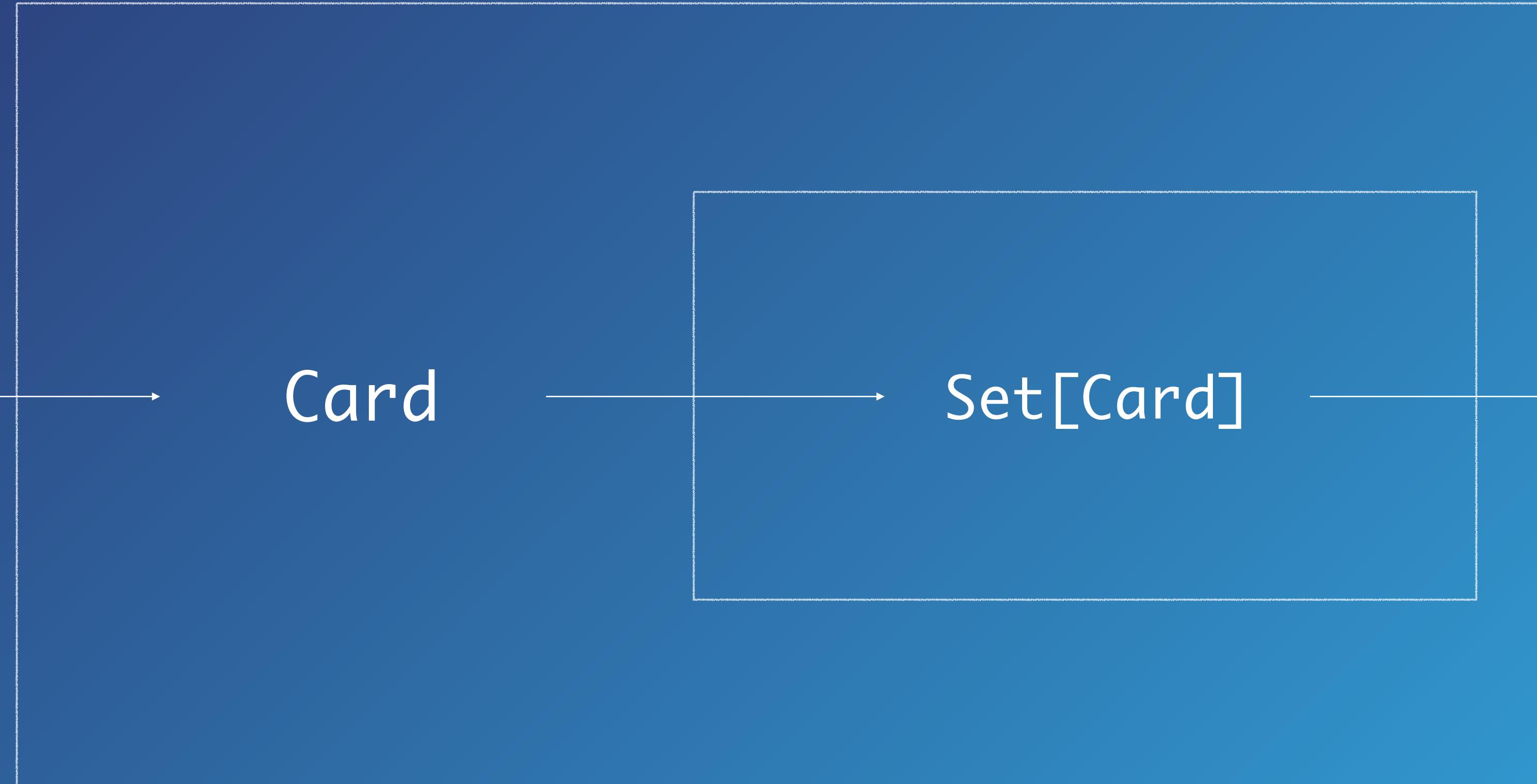
stripe

Charge

Card

Set[Card]

Int





Really leveraging the standard library

```
SetSizeAggregator[Card](10) // bits for hll  
  .composePrepare(a: Charge => a.card)
```

[Charge, HLL[Card] , Int]

Approx Set Size -> HyperLogLog

Approx Item Count -> CountMinSketch

Approx Quantiles -> QTree

Heavy Hitters, Decayed Value, etc..

Aggregator Execution

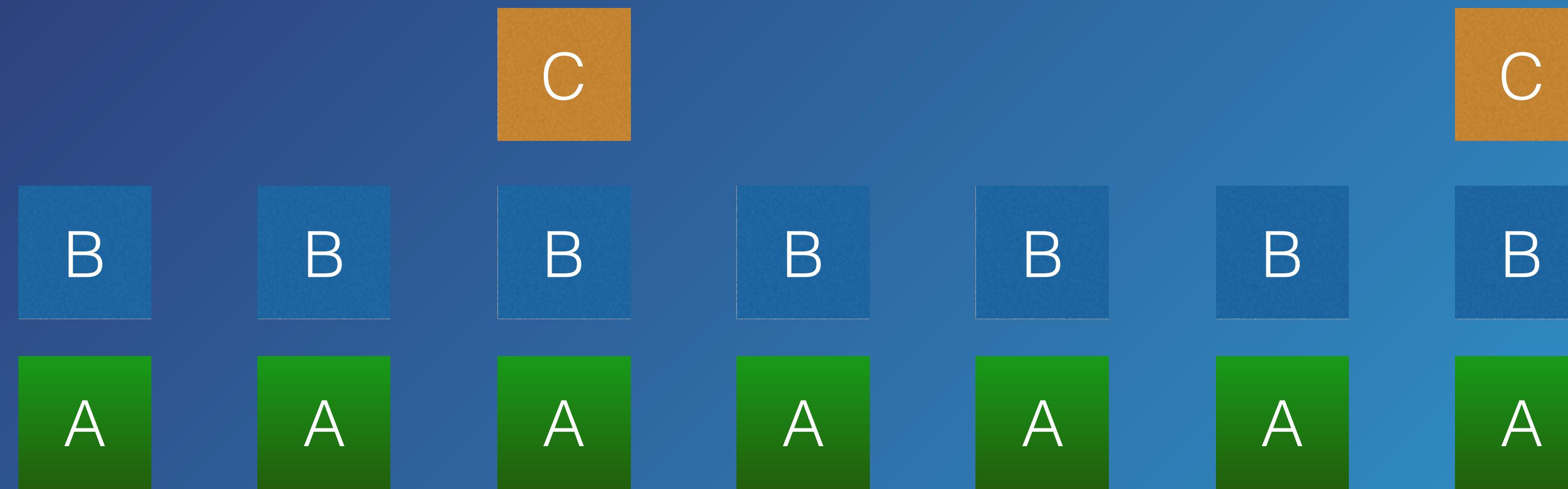
```
apply(inputs: TraversableOnce[A]): C
```

MapReduce with Aggregator

```
def map(v: V) = (key(v), v)
```

```
def reduce(k: K, vs: Seq[V]) = agg(vs)
```

More execution



Fraud Prevention at Stripe

- Merchants as bad actors
- Detect via activity on Stripe (charges)

Predictive Model

When evaluating a merchant, compute score like:

$$\text{sum}(\text{weight1} * \text{feature1} + \text{weight2} * \text{feature2} + \dots)$$



Static Features: Location, Time of Day...

Dynamic Features: Mean Charge Amount, Number of Cards



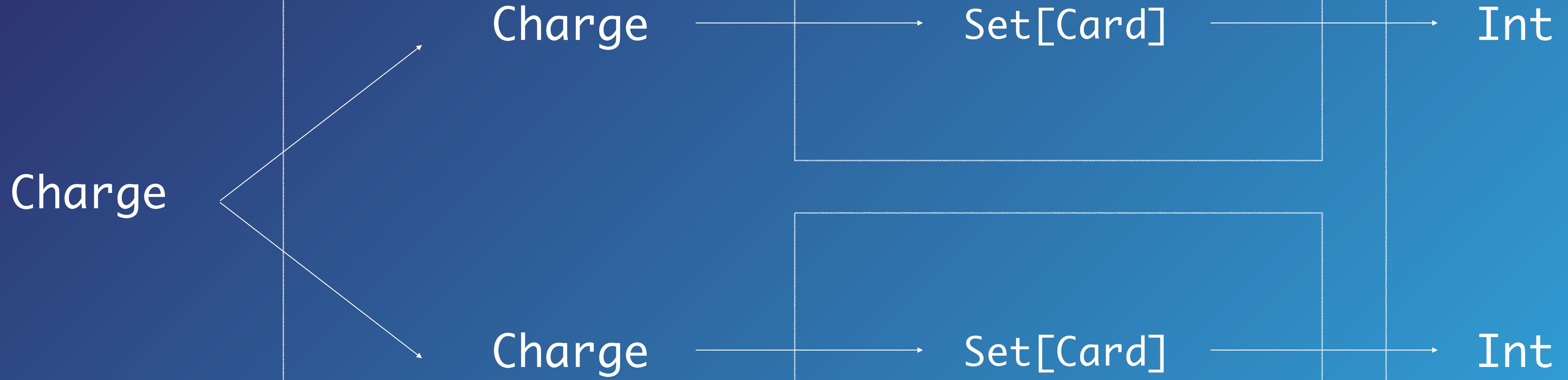
Feature: Number of international cards

```
def internationalCard(c: Charge): Option[Card] = ???  
  
SetSizeAggregator[Card](10) // our normal SetSize  
  .sumBefore() // make it work with Option  
  .composePrepare(a: Charge => a.internationalCard)
```

Composing Aggregators

```
val intlRatio: Aggregator[A, (B1, B2), (C1, C2)] =  
  from2(intlCards, totalCards)
```

stripe



andThenPresent

```
val intlRatio: Aggregator[A, (B1, B2), Double] =  
  from2(intlCards, totalCards)  
    .andThenPresent(case (intl, total) => intl / total)
```

stripe



More features, no problem

```
val features = from2(intlRatio, prepaidRatio)  
  .andThenPresent(toJson(_))
```

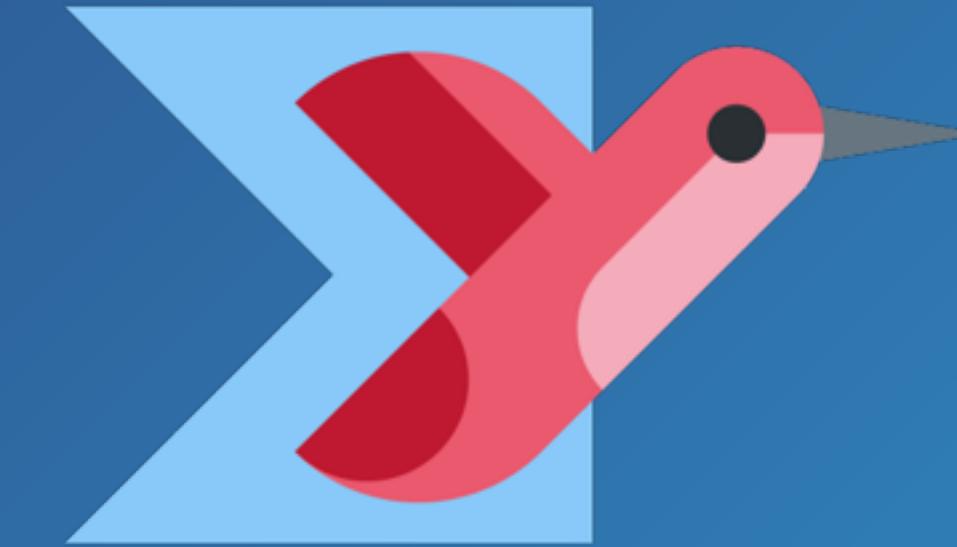
stripe



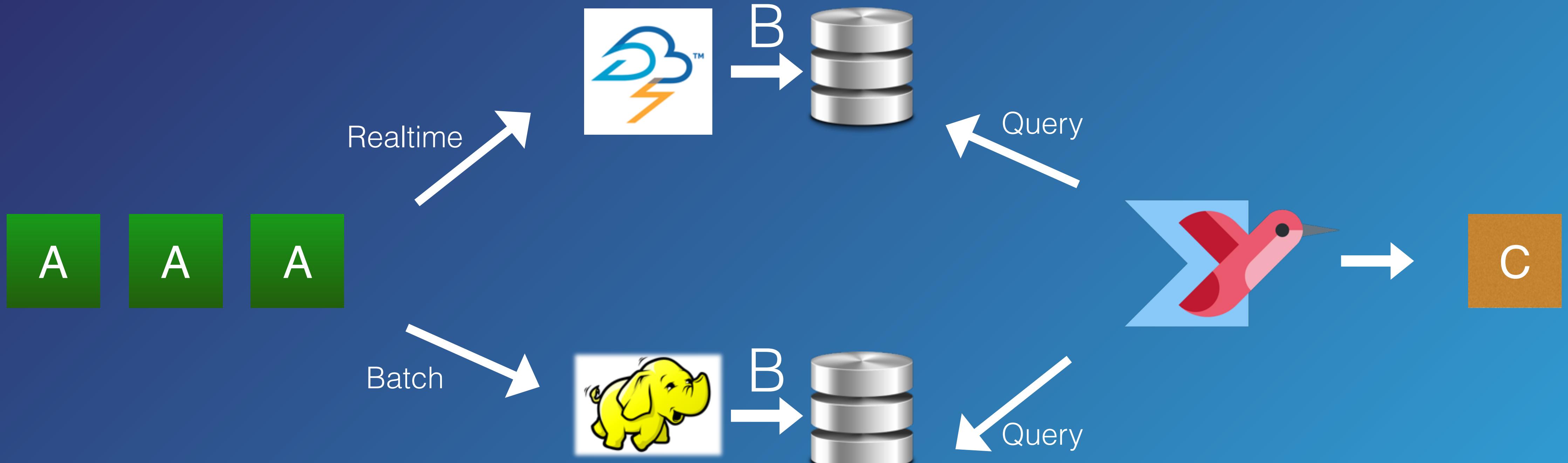
Fraudsters Don't Run on Cron



Beyond Hadoop



- Summingbird: run MapReduce-like code on Storm or Hadoop
- Combine results from two platforms at read-time
- Run complex calculations with Algebird + Aggregator



Recap

- Reduction (and MapReduce) at heart of data analysis
- Aggregator defines reduction in the type system
- (De)composability encourages reuse
- Execution-agnosticism encourages flexible deployments

stripe



<https://stripe.com/jobs/>

stripe

Thanks.

Dan Frank | @danielhfrank | df@stripe.com