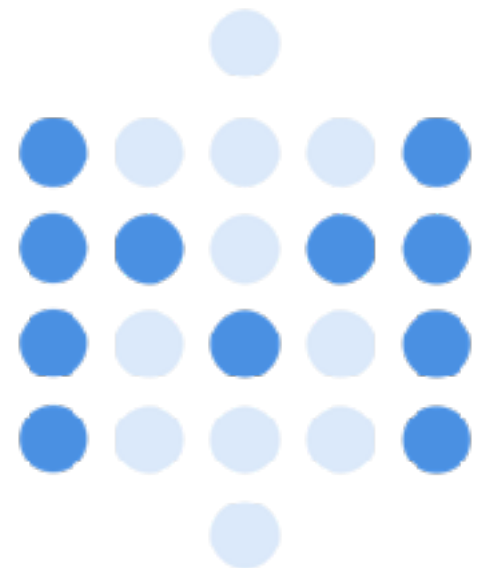


Transducing for fun and profit

simon@metabase.com
@sbelak



Clojure at a glance

- `(lisp (running-on :JVM))`
- Functional, dynamic, immutable
- Excellent concurrency and state-management primitives
- Unparalleled data manipulation

Anatomy of a transducer

- Transducers decomplex recursion mechanism, transformation, building the output, and access mechanism

```
(transduce (map inc) conj (range 10))
```

- 3 user-facing “protocols”: transducer, reducing fn, CollReduce

transducer and reducing function

```
1 (defn my-map [f]
2   (fn [rf]
3     (fn
4       ([_] (rf))
5       ([result] (rf result))
6       ([result input]
7         (rf result (f input))))))
```

```
1 (defn my-min
2   ([_ Double/POSITIVE_INFINITY]
3   ([acc] acc)
4   ([acc e]
5     (min acc e)))
```

transducer and reducing function

```
1 (defn my-map [f]
2   (fn [rf]
3     (fn
4       ([acc (rf))
5        ([result] (rf result))
6         ([result input]
7          (rf result (f input))))))
```

```
1 (defn my-min
2   ([acc Double/POSITIVE_INFINITY]
3    ([acc] acc)
4    ([acc e]
5     (min acc e))))
```

Using a transducer to
wrap/keep state



```
1 (defn my-drop [n]
2   (fn [rf]
3     (let [nv (volatile! n)]
4       (fn
5         ([acc (rf))
6          ([result] (rf result))
7           ([result input]
8            (let [n @nv]
9              (vswap! nv dec)
10             (if (pos? n)
11                result
12                (rf result input))))))))))
```

Wrap Java

```
1 (import 'com.clearspring.analytics.stream.cardinality.HyperLogLogPlus)
2
3 (defn cardinality
4   (□ (HyperLogLogPlus. 14 25))
5   ([^HyperLogLogPlus acc] (.cardinality acc))
6   ([^HyperLogLogPlus acc x]
7    (.offer acc x)
8    acc))
```

CollReduce protocol

- Get the next element
- Makes transducing **data structure-agnostic** allowing us to (re)use transducers for things such as clojure.async channels

Transducing an async channel

```
1 (let [ch (async/chan 1)]
2   (async/onto-chan ch (range 10))
3   (async/<!! (async/transduce identity + 0 ch)))
```

```
1 (let [ch (async/chan 1 my-count)]
2   (async/onto-chan ch (range 10))
3   (async/<!! (async/into [] ch)))
```


Composing transducers

1. comp transducers

```
1 (transduce (comp (filter odd?)
2             (map inc))
3           conj
4           (range 10))
```

2. Reducing function and transducer

```
1 ((remove nil?) +)
```

3. github.com/henrygarner/redux

post-complete

Data structure that can be manipulated like any other

fuse

```
1 (transduce identity
2   (redux/post-complete
3     +
4     (fn [sum]
5       (Math/sqrt sum)))
6   (range 10))
```

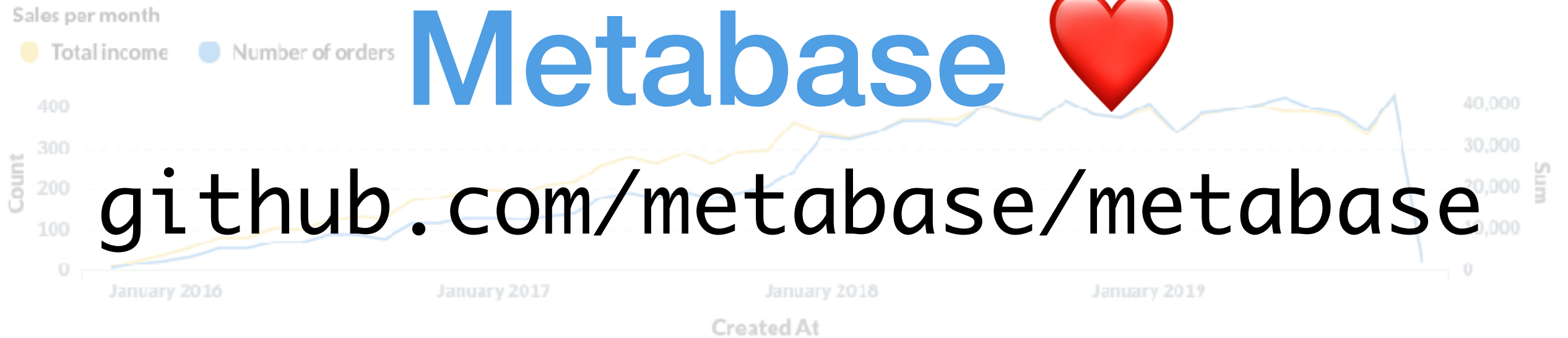
```
1 (transduce identity
2   (redux/fuse {:sum +
3               :conj conj})
4   (range 10))
```

On-line/streaming analysis

Metabase

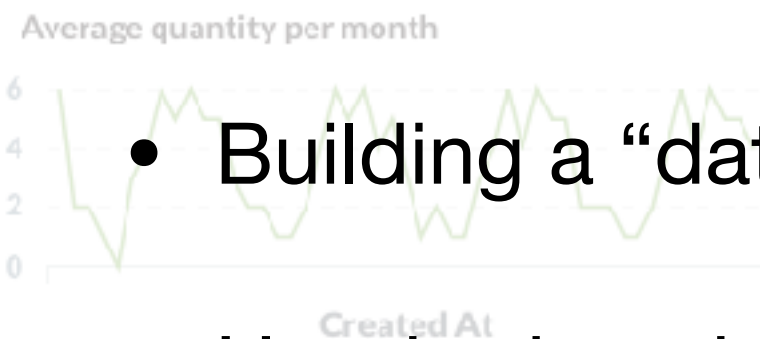


github.com/metabase/metabase



How these transactions are distributed

- Open source analytics tool (runs on-premises)

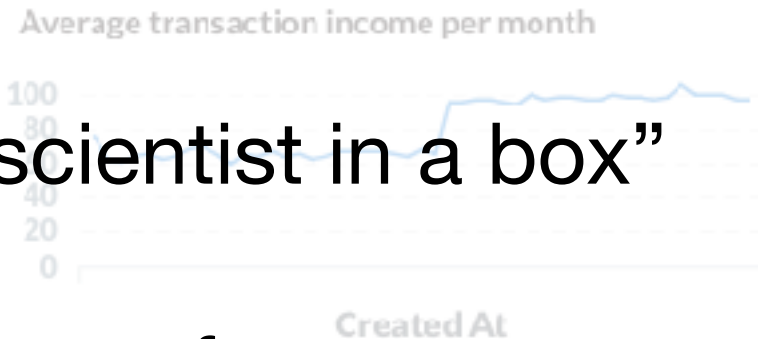


- Building a “data scientist in a box”

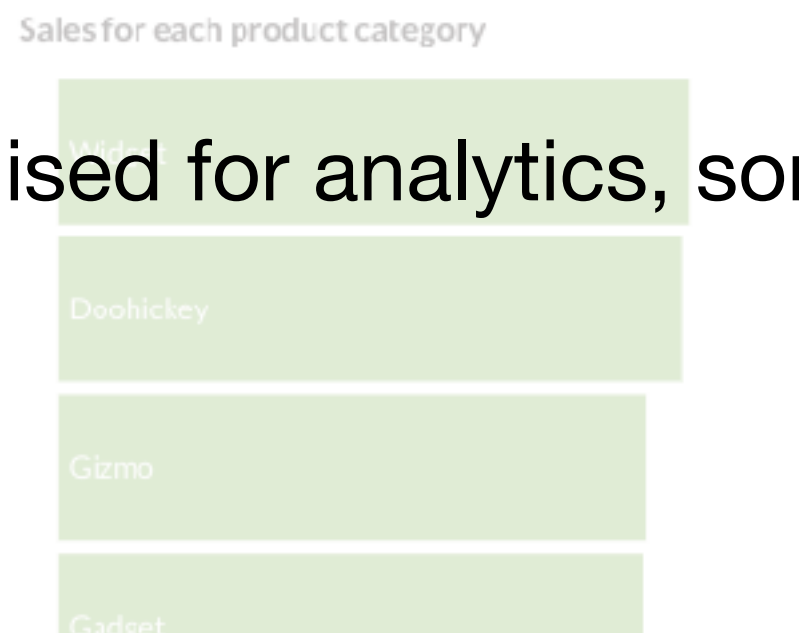
- Hundreds to billions of rows

Sales per product

Title	Count
Rustic Leather Plate	126
Incredible Copper Lamp	121
Rustic Rubber Bottle	88
Incredible Bronze Gloves	87
Durable Concrete Shoes	84



- Some DBs optimised for analytics, some not



Many batch algorithms can be turned into online ones

$$\bar{x} = \frac{1}{n} \left(\sum_{i=1}^n x_i \right) = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\sigma^2 = \frac{\sum_{i=1}^N (X - \mu)^2}{N}$$

```
1 (defn mean
2   ([[] [0 0]])
3   ([[s c :as acc] e]
4     [(+ s e) (inc c)]])
5   ([[s c]]
6     (when-not (zero? c)
7       (/ s c))))
```

```
1 (defn variance
2   ([[] [0 0 0]])
3   ([[c m ss :as acc] e]
4     (let [c' (inc c)
5           m' (+ m (/ (- e m) c'))]
6       [c' m' (+ ss (* (- e m') (- e m)))]))
7   ([[c m ss]
8     (when-not (zero? c)
9       (let [c' (dec c)]
10          (if (pos? c')
11              (/ ss c') 0))))))
```

Parallelize independent computations

Find a recursive relation

github.com/MastodonC/kixi.stats

- Count
- (Arithmetic) mean
- Geometric mean
- Harmonic mean
- Median
- Variance
- Interquartile range
- Standard deviation
- Standard error
- Skewness
- Kurtosis
- Covariance
- Covariance matrix
- Correlation
- Correlation matrix
- Simple linear regression
- Standard error of the mean
- Standard error of the estimate
- Standard error of the prediction
- ...

Data = code

**Using transducers is
worth it for the
composition alone**

Annoyances

- Can only transduce one coll at a time
- Always have to pass in an xf (especially annoying when using redux)
- Having functions that return a transducer or not is error prone
- Inconsistent support for transducers in core library

Questions

simon@metabase.com
@sbelak

