# BUILDING A RECOMMENDATION ENGINE USING DIVERSE FEATURES

Divyanshu Vats
*Sailthru*

# Joint work with

Max Sperlich
(Sailthru)

David Glueck
(Sailthru)

Alex Gaudio
(Alluvium)

Jeremy Stanley
(Instacart)

SAILTHRU

# Powering More Than 400 Ecommerce & Media Brands

Sailthru allows marketers to manage consumer relationships at the individual user level across all channels. It's the realization of a promise that has been made to marketers for more than a decade. Sailthru's ability to deliver personal communications and experiences to every unique individual is driving lift and creating revenue where marketers, like you, want it most.

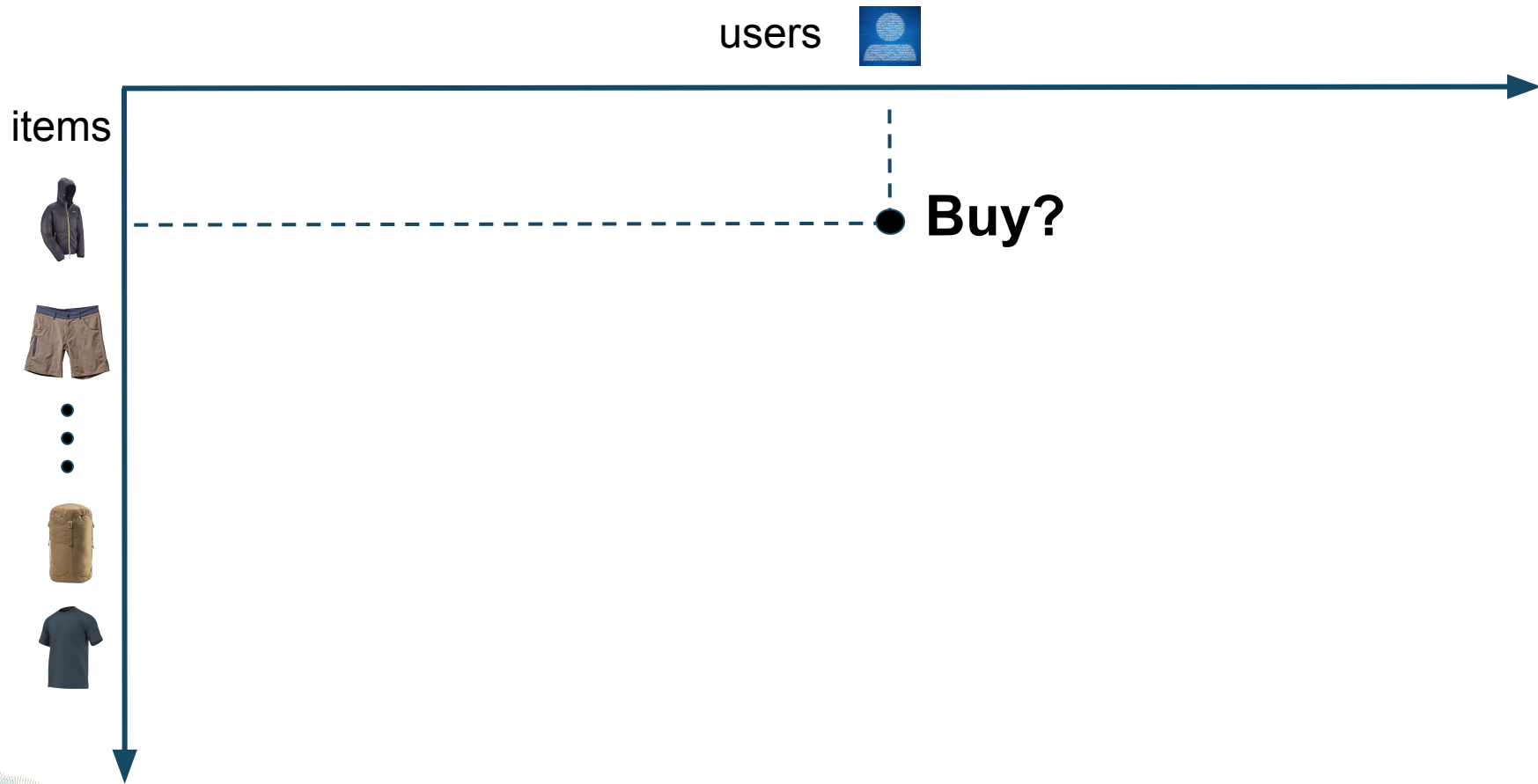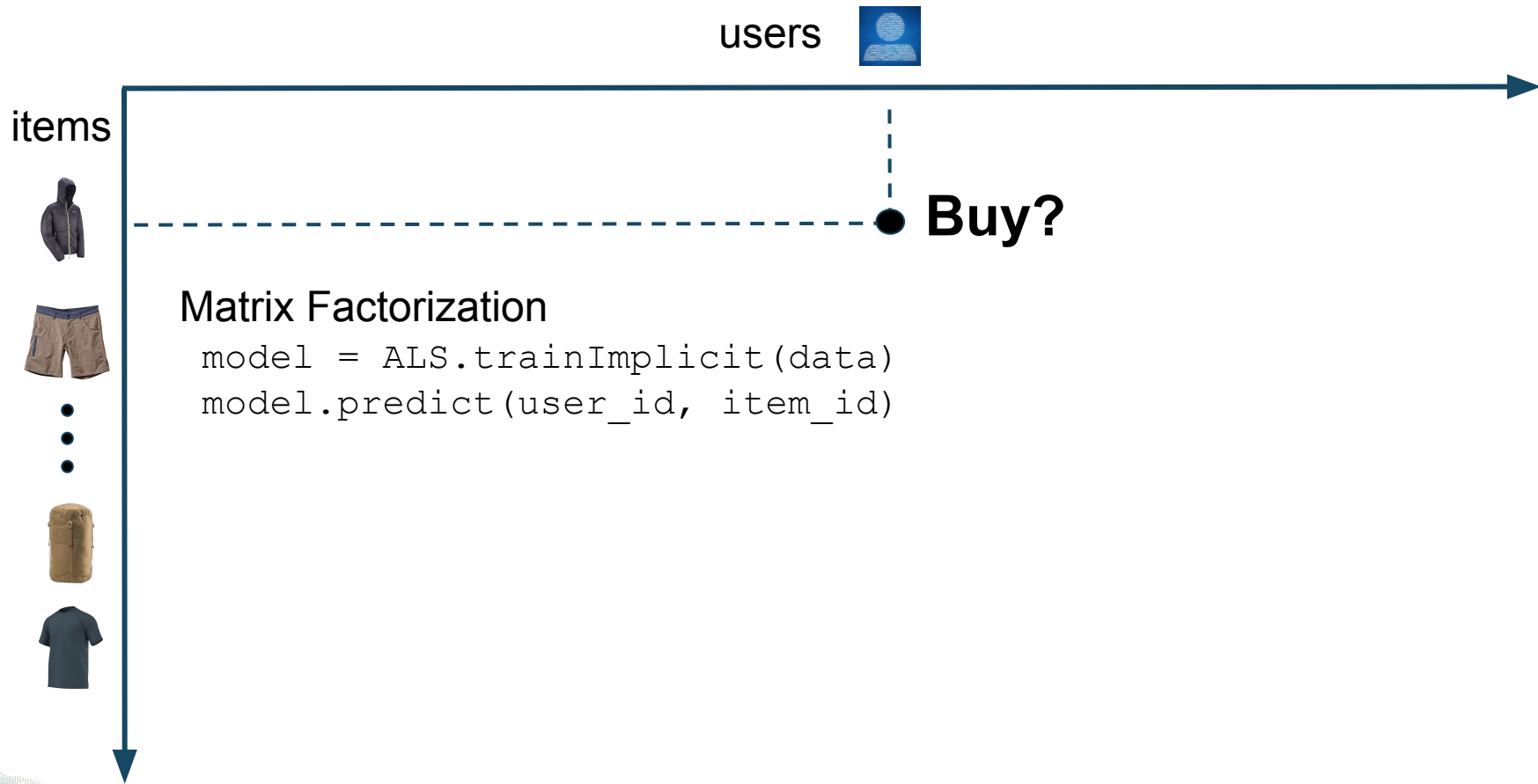**Mashable**   The Economist   **BIRCHBOX**   ALEX AND ANI   FRANK & OAK
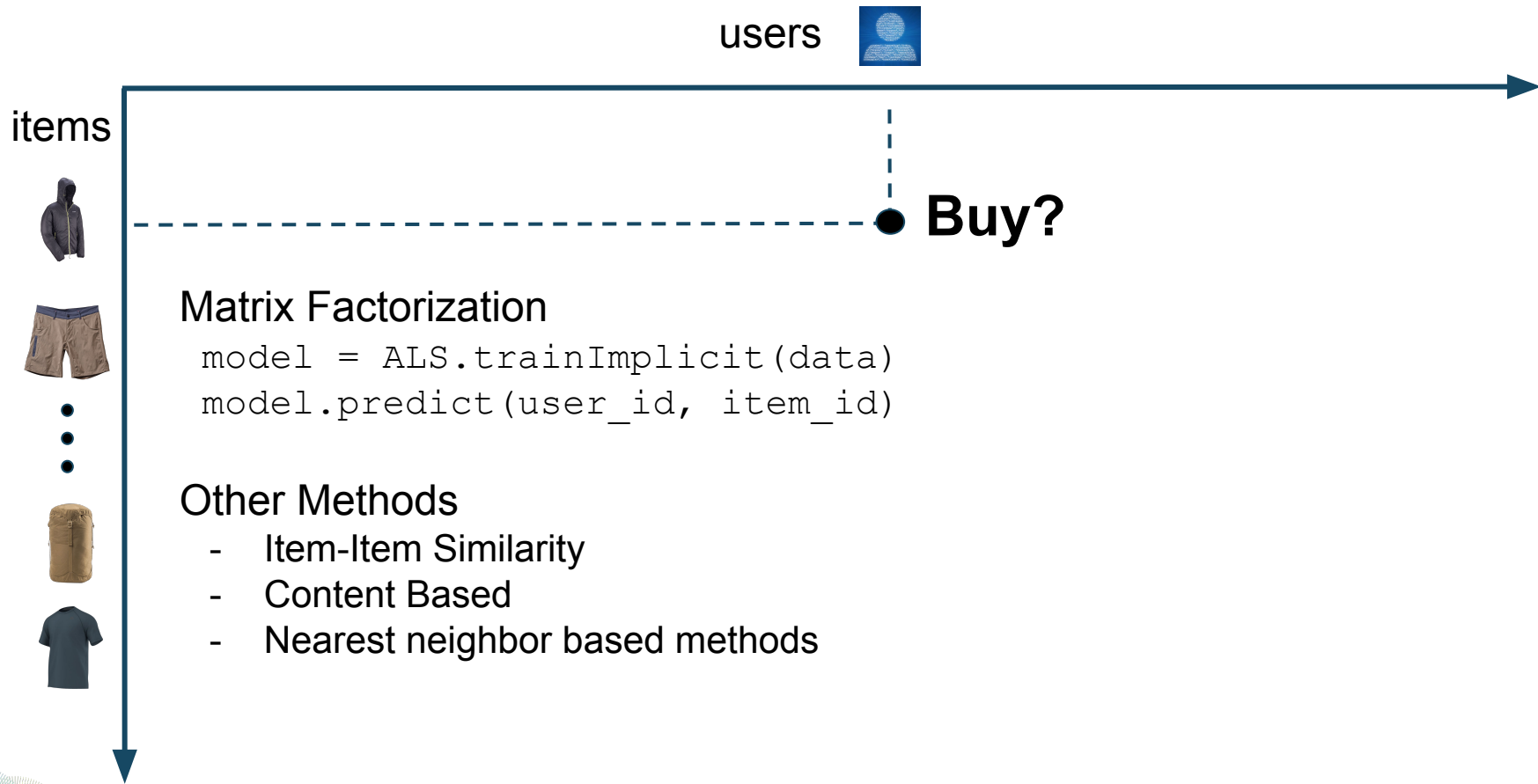
# Customer Retention Cloud

SAILTHRU

users

items

**Buy?**

Matrix Factorization

```
model = ALS.trainImplicit(data)
model.predict(user_id, item_id)
```

SAIL7HRU

users

items

**Buy?**

Matrix Factorization
```
model = ALS.trainImplicit(data)
model.predict(user_id, item_id)
```

Other Methods
- Item-Item Similarity
- Content Based
- Nearest neighbor based methods

SAILTHRU

users

items

**Buy?**

Matrix Factorization
```
model = ALS.trainImplicit(data)
model.predict(user_id, item_id)
```

Other Methods
- Item-Item Similarity
- Content Based
- Nearest neighbor based methods

👍 Open source scalable implementations
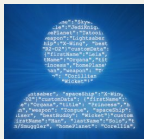
👎 Easily Incorporate **additional features**?

SPARK SUMMIT EAST
2016

SAIL**THRU**

# Additional Features?

# Additional Features?

# Additional Features?



user

location   # clicks

browser

item

image   title

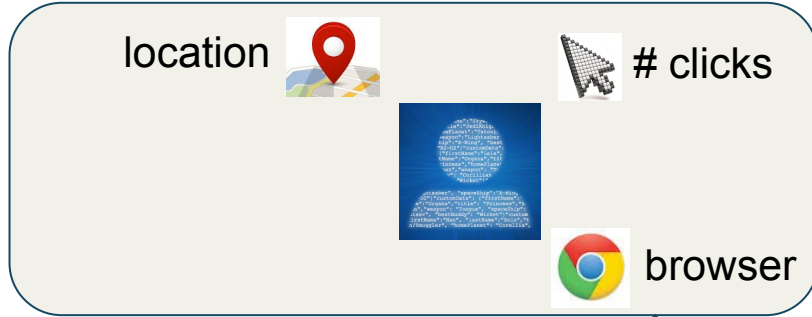Marled Yarn Wool-Blend Sweater in Gray
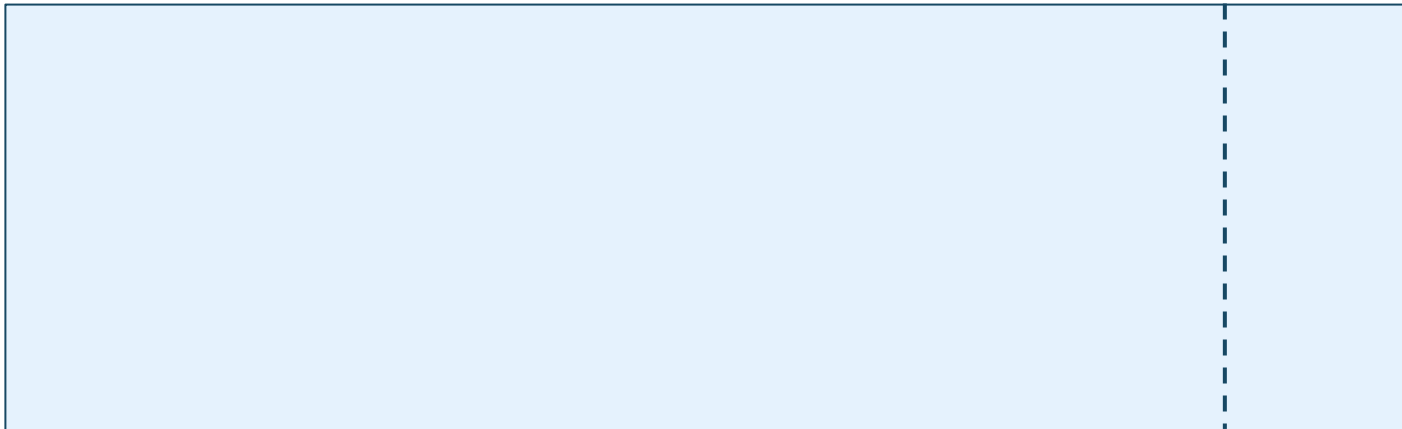
68 USD

description

user-item

- purchases
- item views

Goal: Framework to **easily** integrate **user**, **item**, and **user-item** interactions for meaningful recommendations
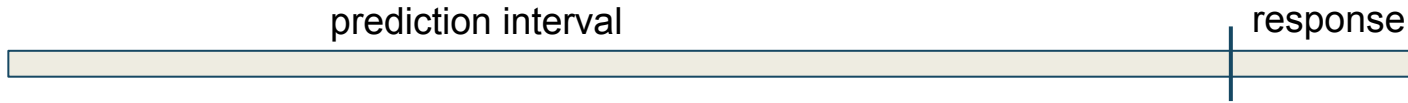
SAIL*T*HRU

prediction interval          response

prediction interval                                    response

features                                          | 1

                                                  | 0

SAIL7HRU

prediction interval | response

| user | item | user-item interactions | user-item features | response |
| --- | --- | --- | --- | --- |
| | | | | 1 |
| | | | | 0 |

Can use a combination of collaborative filtering algorithms!
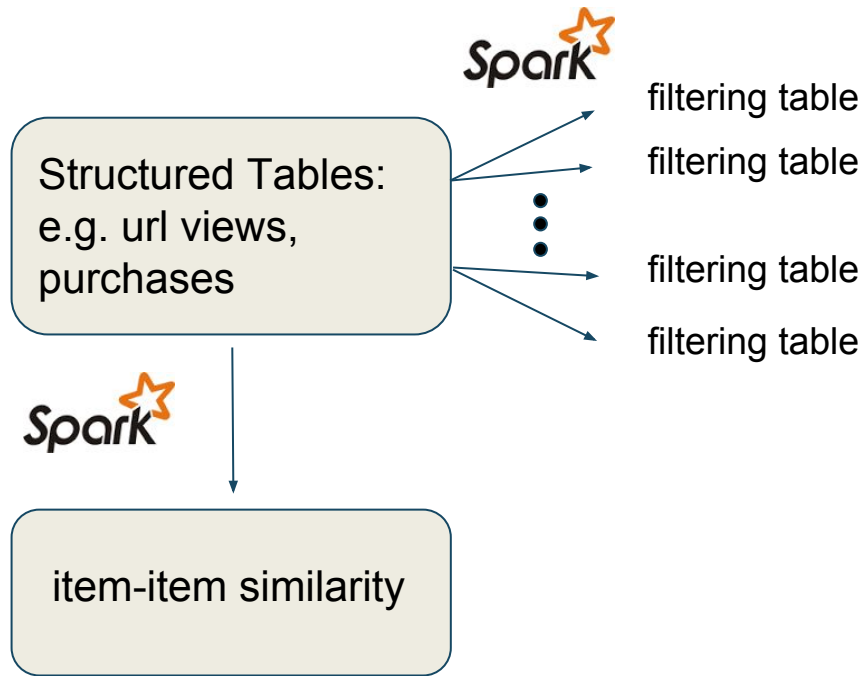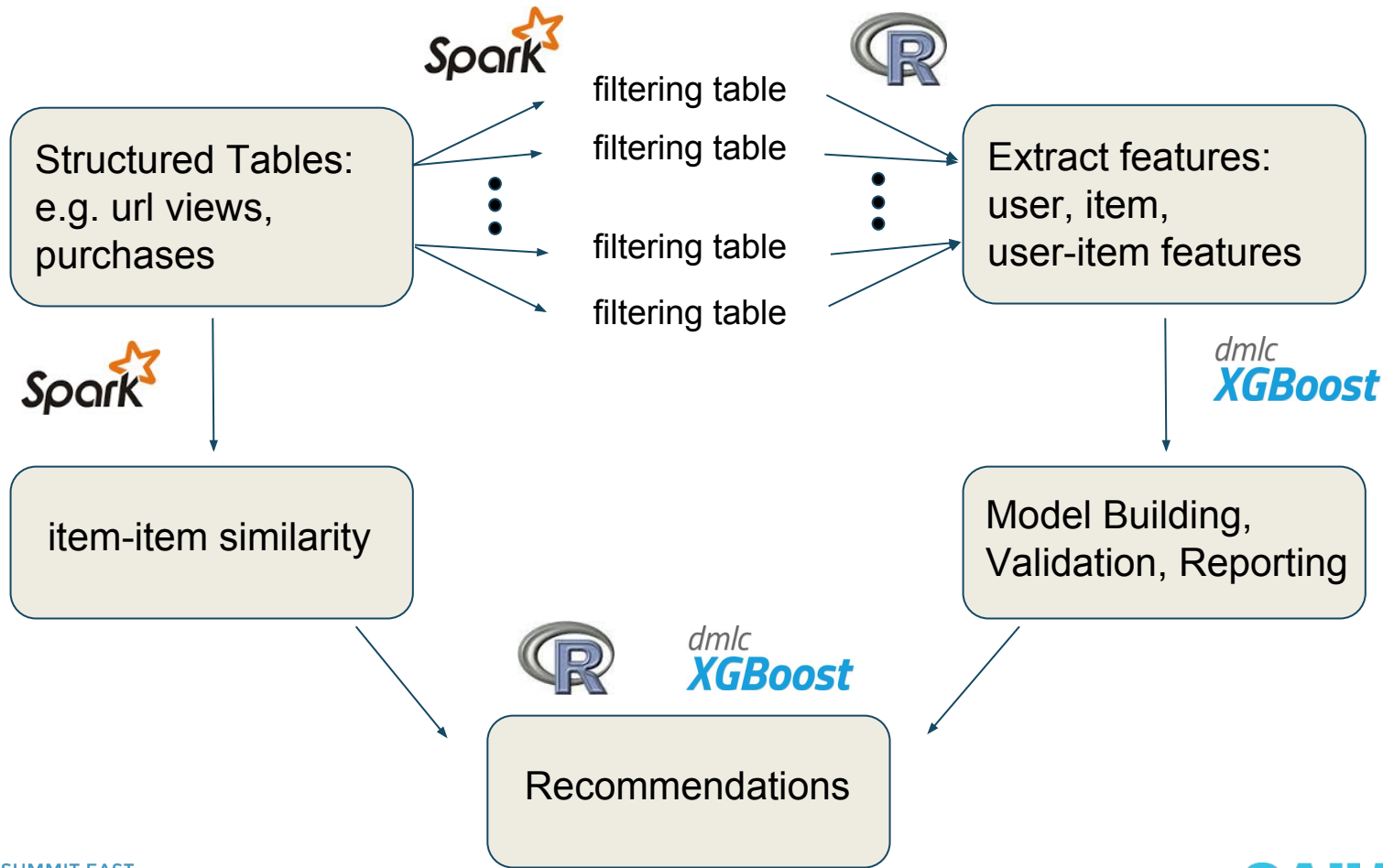
SAILTHRU

# So far….



users

items

**Buy?**

- Easily integrate user, item, and user-item features

- Recommendation → Binary classification

- Next: Implementation and use of Spark!

SAIL7HRU

Structured Tables:
e.g. url views,
purchases

SAIL**7**HRU

Spark

filtering table
filtering table
⋮
filtering table
filtering table

R

Structured Tables:
e.g. url views,
purchases

Extract features:
user, item,
user-item features

dmlc
XGBoost

Spark

item-item similarity

Model Building,
Validation, Reporting

dmlc
XGBoost
R

Recommendations

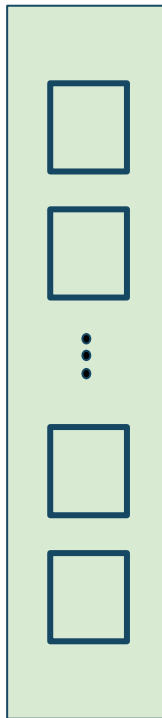SPARK SUMMIT EAST
2016

SAILTHRU
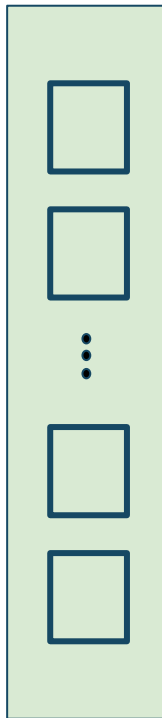
# Know Thy Partitions!

```
data.select('user_id', 'url')
    .groupBy('user_id')
    .count()
    .filter(count > 10)
```

SAILTHRU

# Know Thy Partitions!

```
data.select('user_id', 'url')
     .groupBy('user_id')
     .count()
     .filter(count > 10)
```
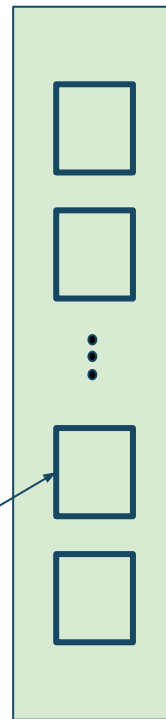
Inefficient!

SAIL7HRU

# Know Thy Partitions!

```
data.select('user_id', 'url')
    .groupBy('user_id')
    .count()
    .filter(count > 10)
```
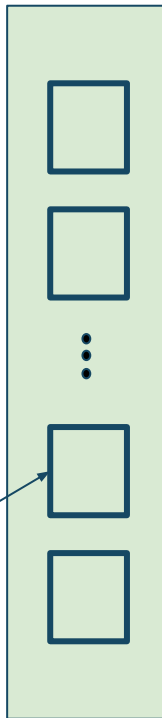
partition by user_id

# Know Thy Partitions!

```
data.select('user_id', 'url')
    .mapPartitions(filter_url)
```

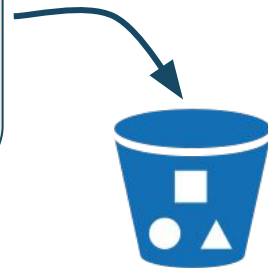partition by user_id

# Break up application into small chunks!

```
users = data.select('user_id', 'url')
            .mapPartitions(filter_url)
            .collect()



BigData.filter(lambda x: x.user_id in users)
```

SAIL7HRU

# Break up application into small chunks!

```
users = data.select('user_id', 'url')
            .mapPartitions(filter_url)
            .collect()
```
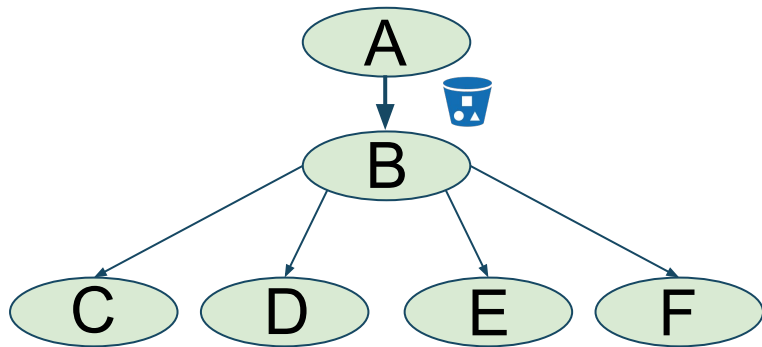
# Break up application into small chunks!

```
users = data.select('user_id', 'url')
              .mapPartitions(filter_url)
              .collect()
```

```
BigData.filter(lambda x: x.user_id in users)
```
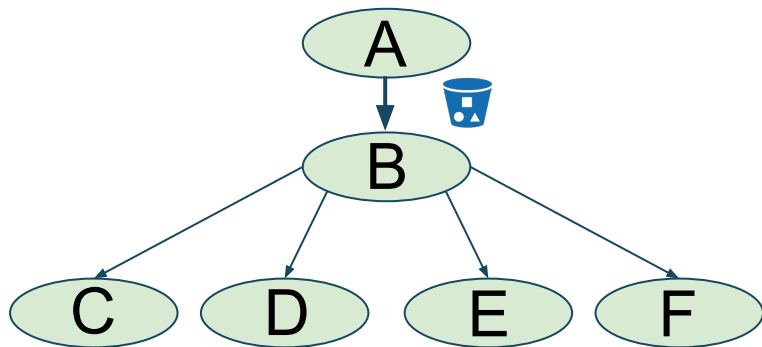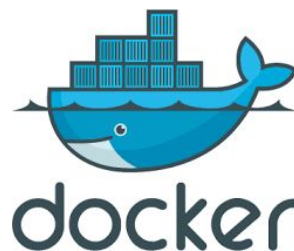
SAIL7HRU

# Mesos + Scheduler + Docker + Spark

- Carefully define applications and state a dependency graph

- Manage graph using: [github.com/sailthru/stolos](github.com/sailthru/stolos)
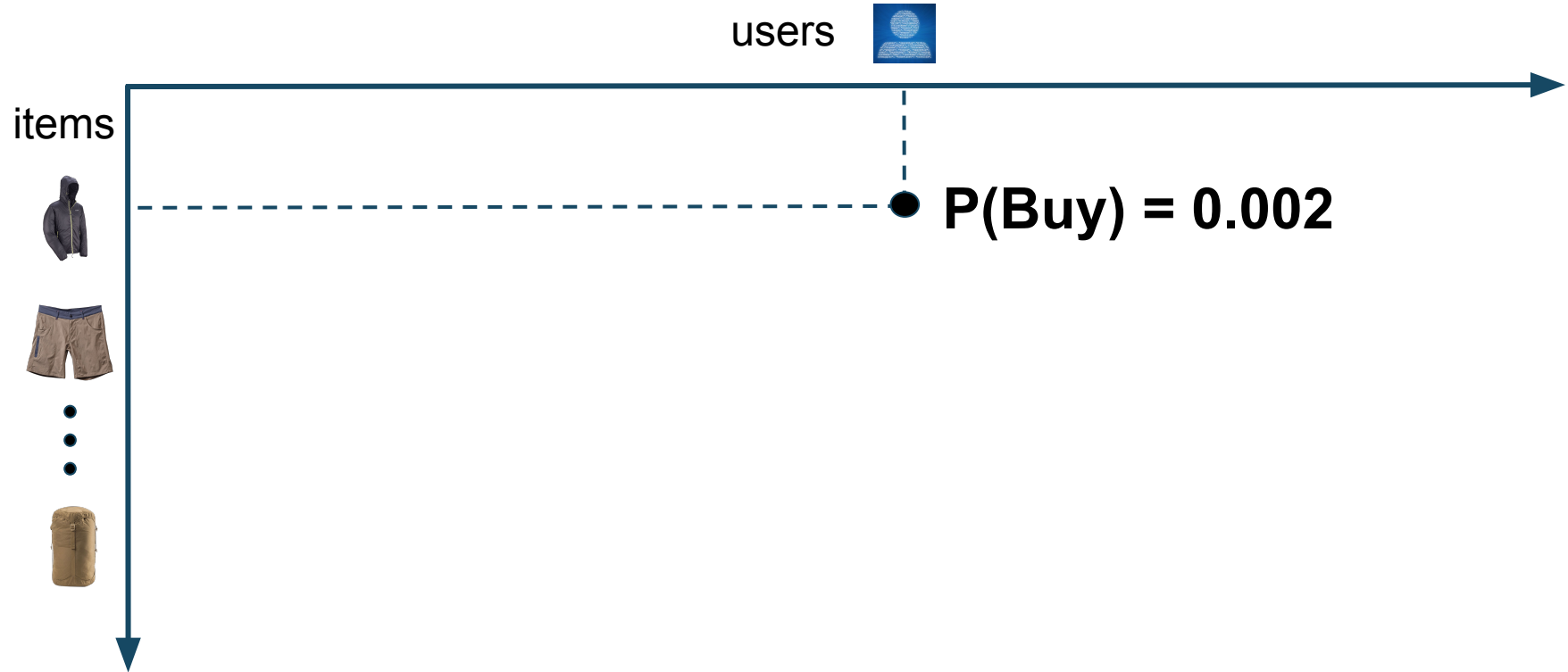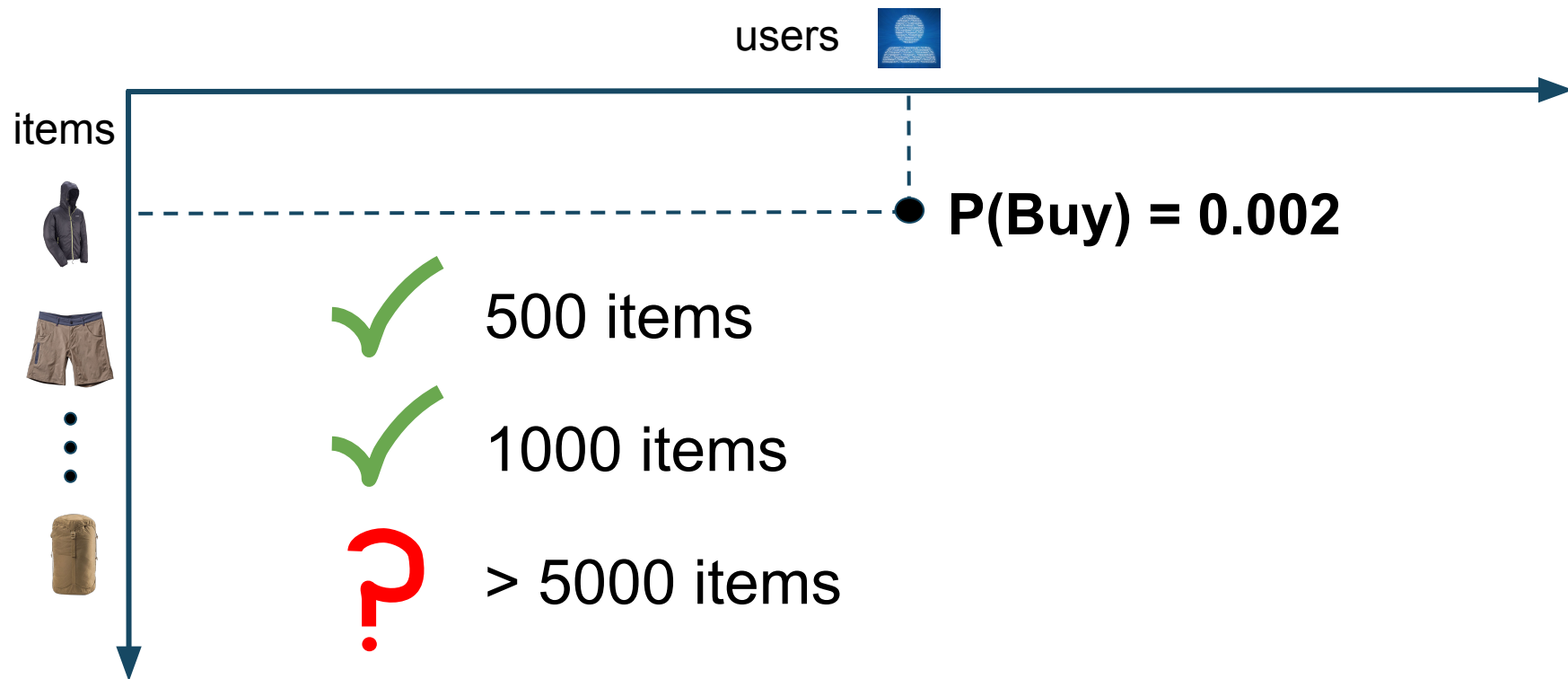
# Mesos + Scheduler + Docker + Spark



- Carefully define applications and state a dependency graph

- Manage graph using:
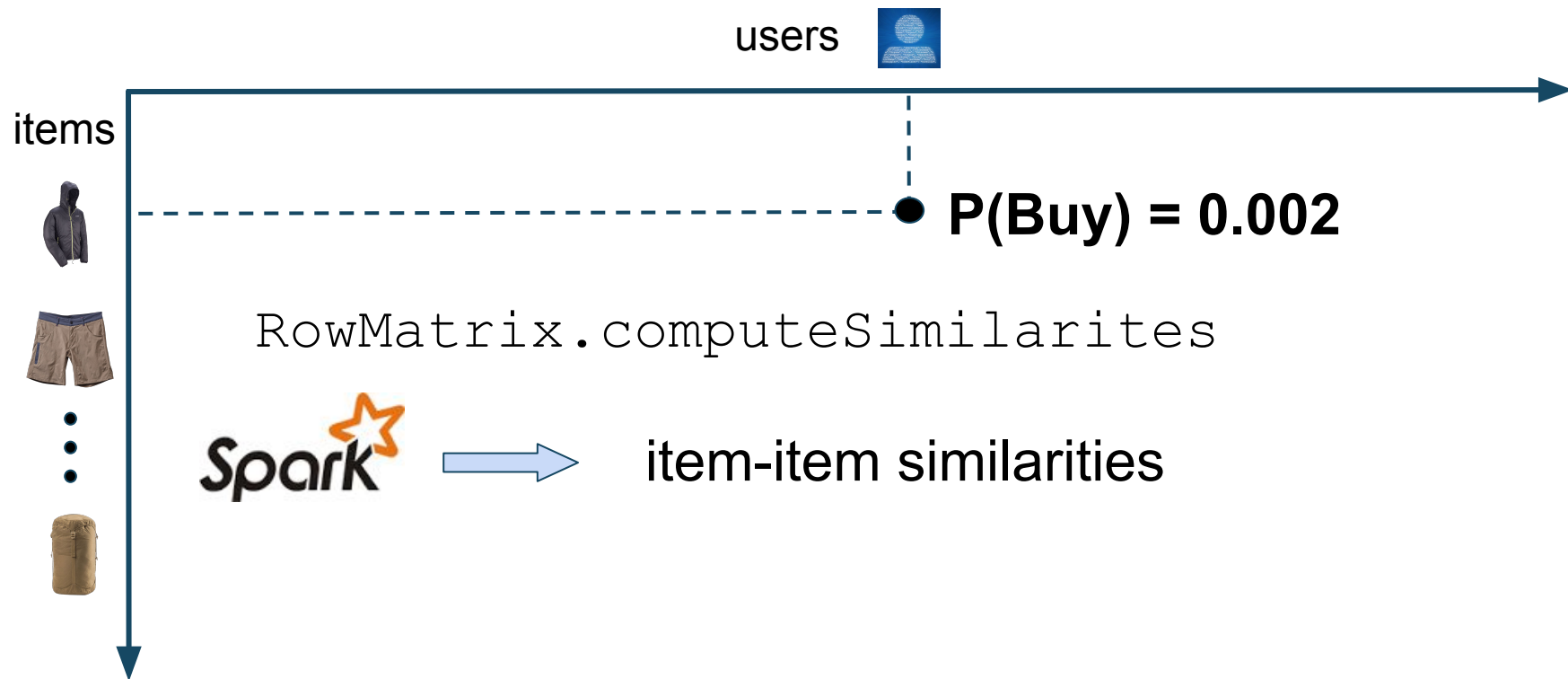  [github.com/sailthru/stolos](github.com/sailthru/stolos)

# Item-Item Similarities

users

items

**P(Buy) = 0.002**
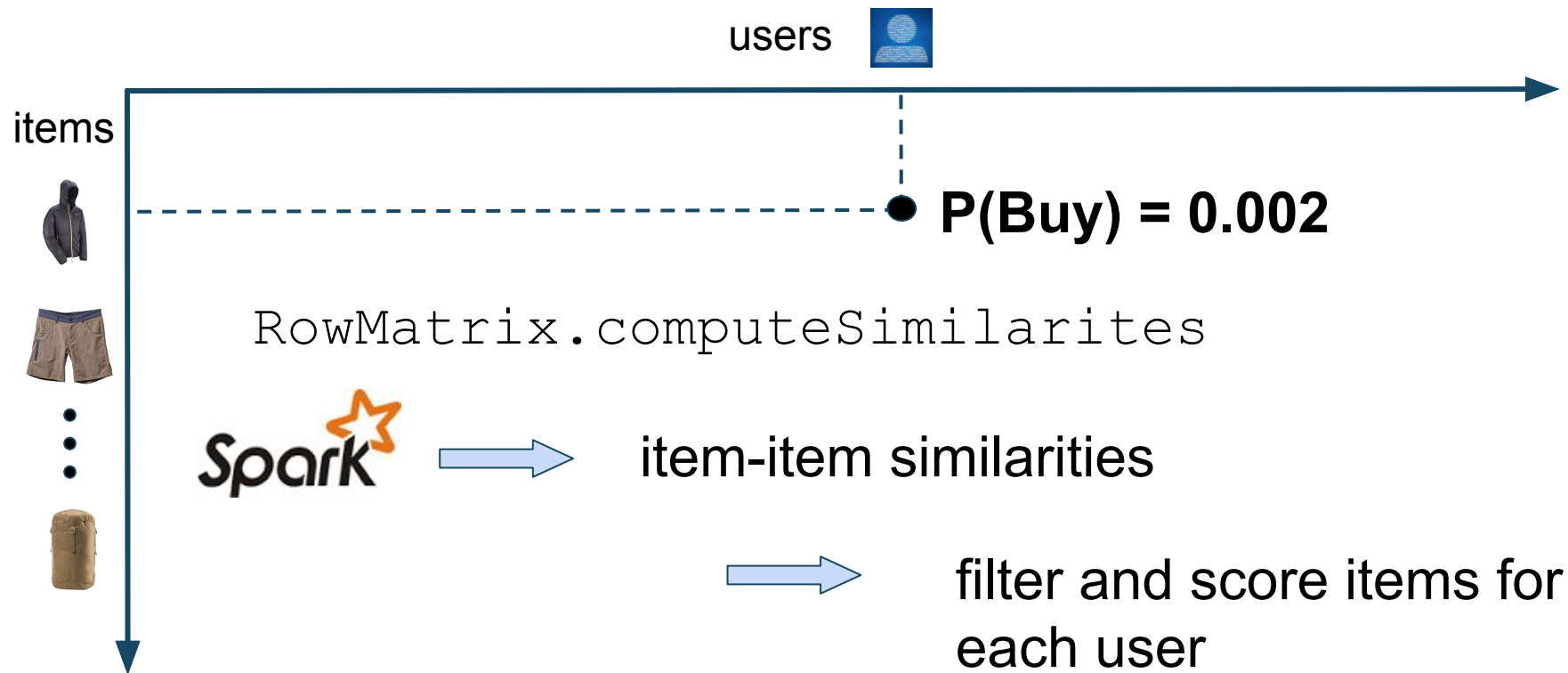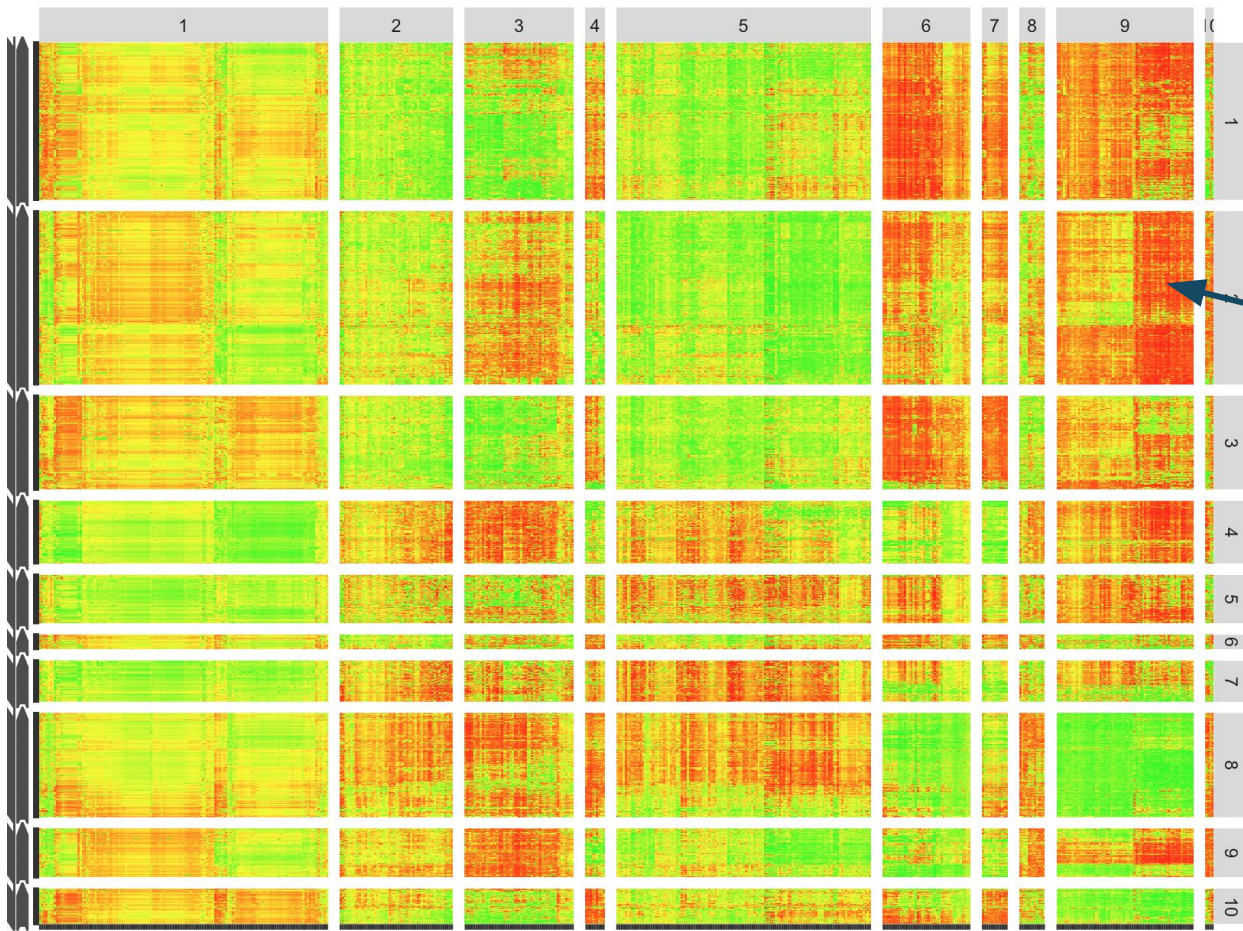
SAIL*7*HRU

# Item-Item Similarities



users

items

P(Buy) = 0.002

✓ 500 items

✓ 1000 items

? > 5000 items

SAIL7HRU

# Item-Item Similarities

users

items

**P(Buy) = 0.002**

`RowMatrix.computeSimilarites`
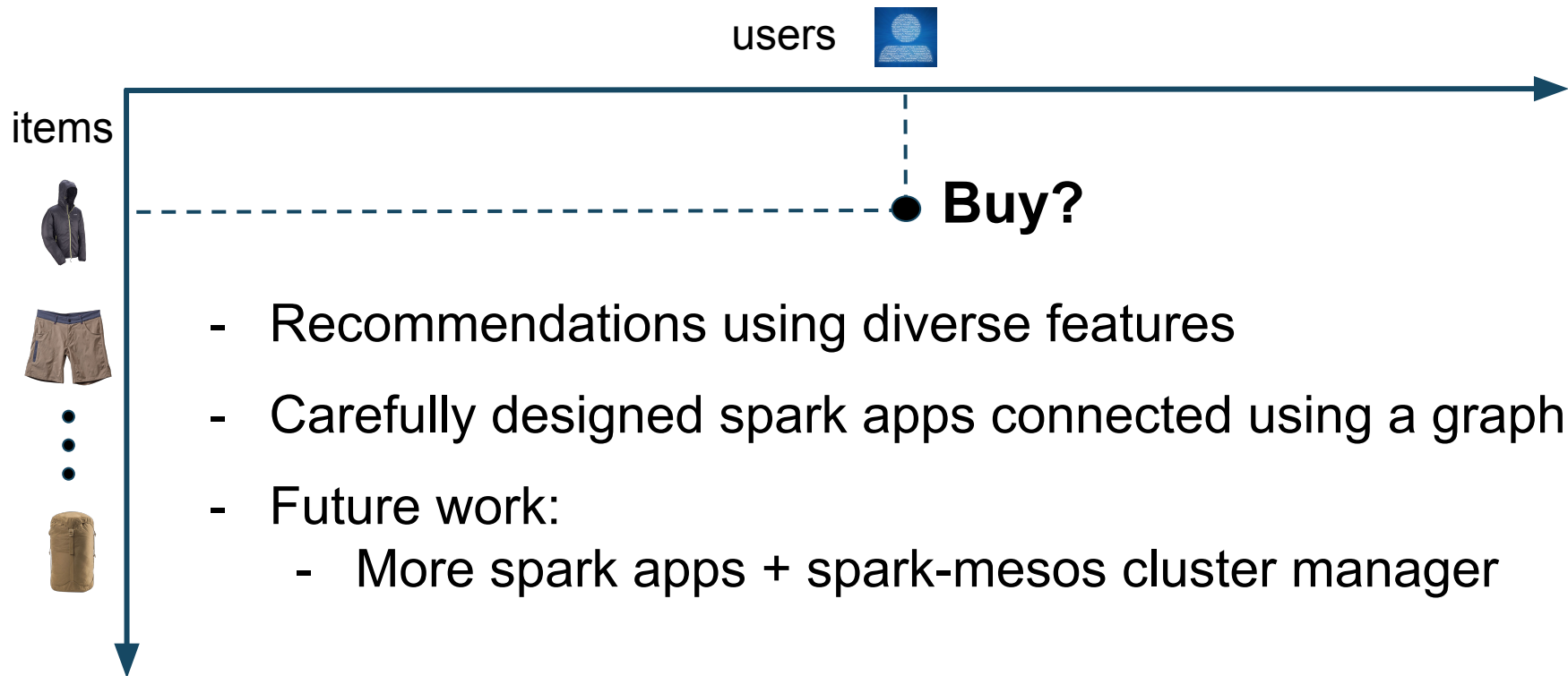
item-item similarities

# Item-Item Similarities

users

items

P(Buy) = 0.002

RowMatrix.computeSimilarites

Spark ➡ item-item similarities

➡ filter and score items for each user

SAILTHRU

cluster of users and items

# To Summarize..

users

items

**Buy?**

- Recommendations using diverse features

- Carefully designed spark apps connected using a graph

- Future work:
    - More spark apps + spark-mesos cluster manager

**SAIL7HRU**

# THANK YOU.

email: dvats@sailthru.com