



# Recommendation at scale

*Simon Dollé*

*Berlin Buzzwords, June 2<sup>nd</sup>, 2015*



Ab  
**56 €**

**Axel Hotel Berlin**



Berlin

**Hotels anzeigen** >



Asus, X551mav-  
RcIn06, 15.6" ...

~~\$299.99~~

**Shop**



Hp 950xl/951  
High Yield...

**Shop**



Google Nexus 7  
Lte 7", 32gb ...

~~\$349.99~~

**Shop**



Rund vier Jahre von  
der Erde entfernt...

**» KAUFEN**



**Star Wars  
Trilogie...**



**Piqué-Polo aus  
100% Baumwolle**

**29,99 €**

**Shop**



# Where is the need for tech ?

## We buy

- Inventory ! (ad spaces)
- Billions of times a day
- All over the Internet
- For 95% of the population



## We sell

- Clicks !
  - (that convert)
  - (that convert a lot)



## We take the risk

*You pay only for what you get*





### Traffic

- » 800 k HTTP requests / sec (peak activity)
- » 29000 impressions /sec (peak activity)
- » Less than 10 ms to process an RTB request

### Physical infrastructure

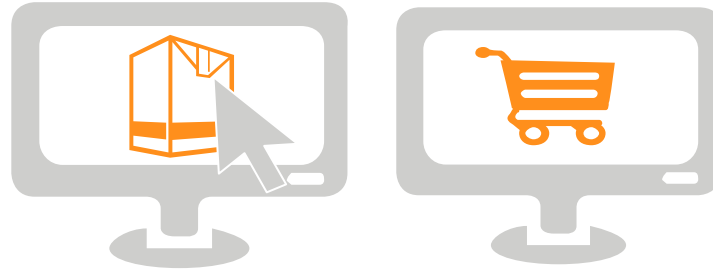
- » 6 Data centers on 3 continents operated and conceived in-house
- » ~ 12000 servers, largest Hadoop cluster in Europe
- » More than 35 PB of storage Big Data

# Data Sources

---



Catalog data



User behavior data



Ad display data

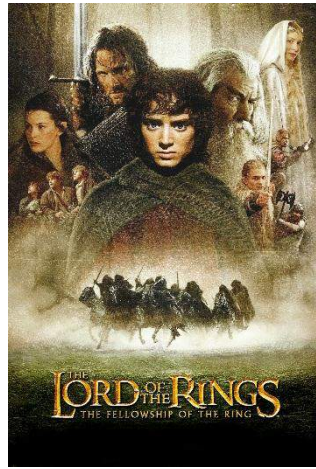
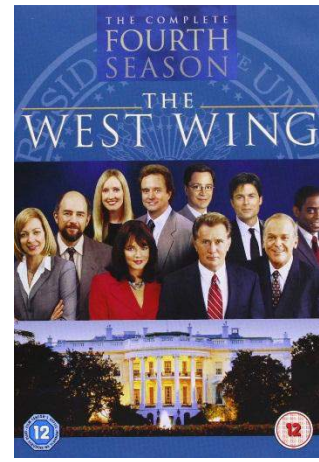
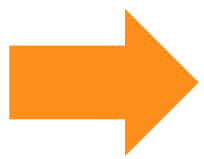


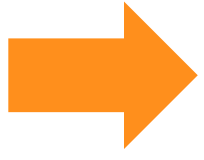
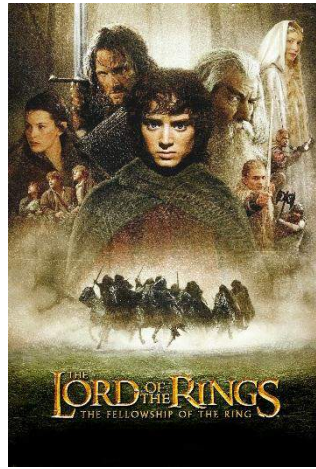
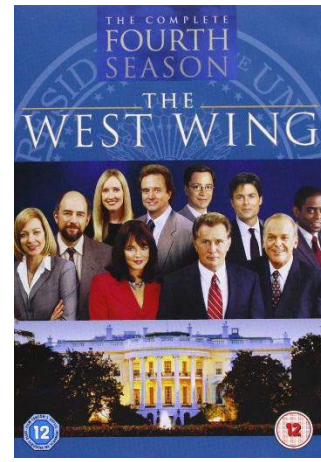
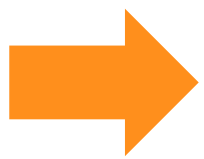




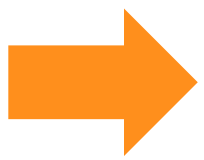
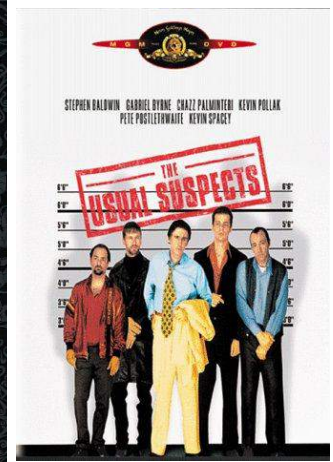
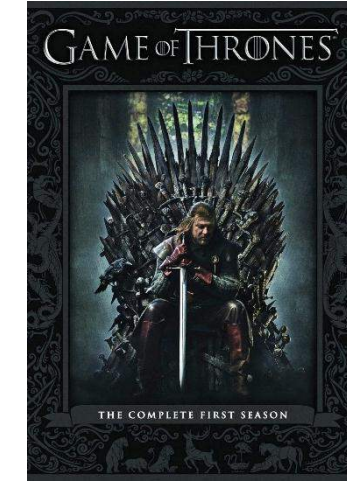
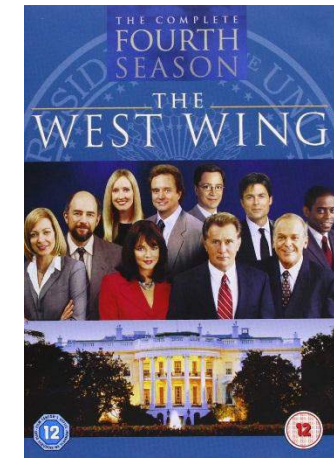
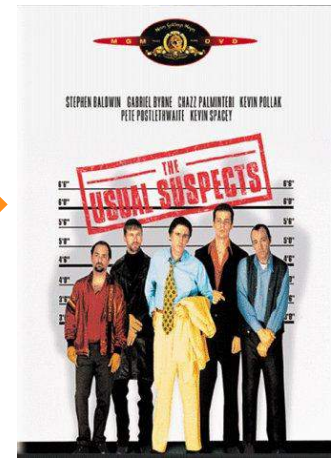
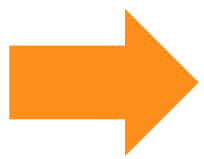


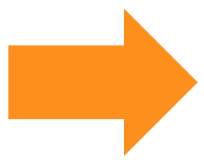
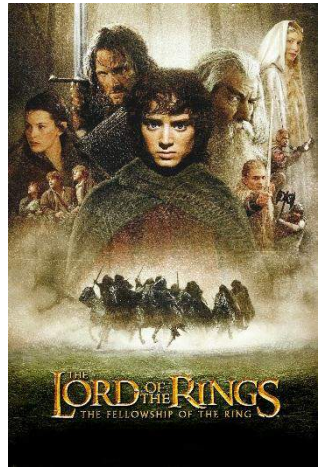
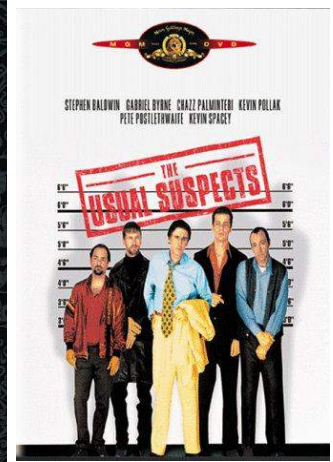
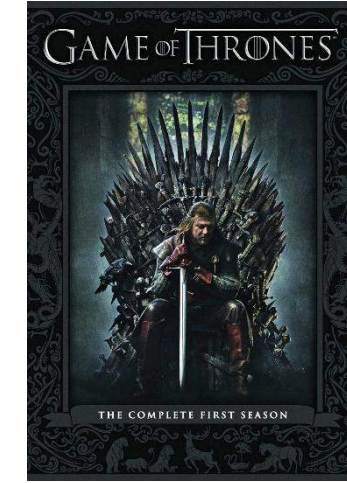
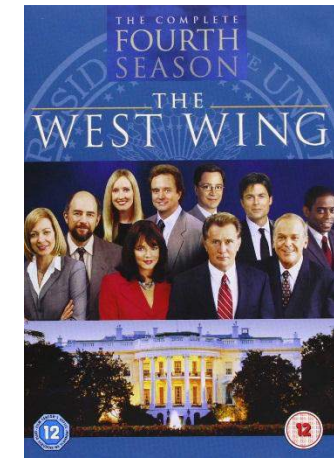
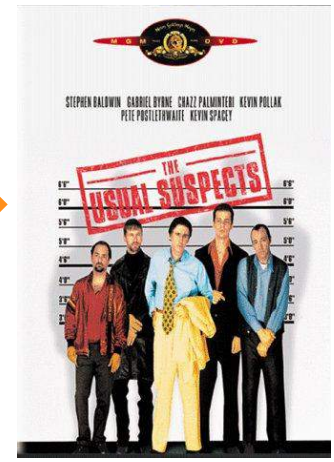
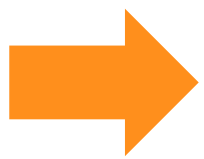










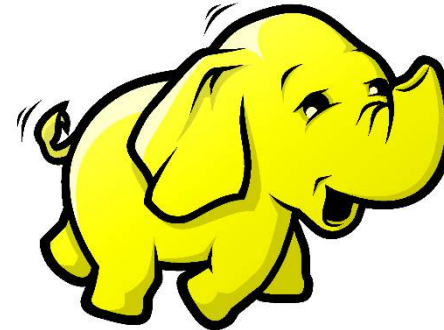


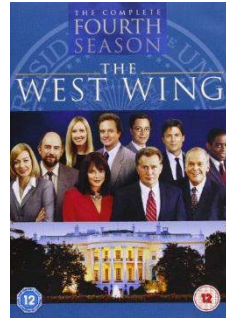


# Offline

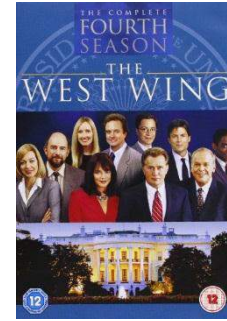
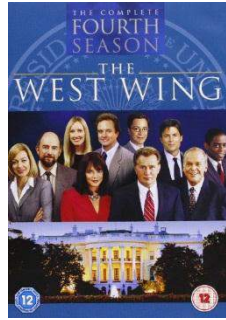
---

- Similarities computed on browsing data
- Based on coevents
- Computed on Hadoop cluster
- Map reduce jobs, pig
- Takes around 12 hours
- Tradeoff performance, speed







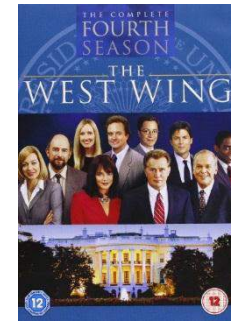
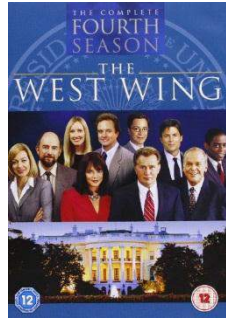


0.04

0.02

0.02

0.05

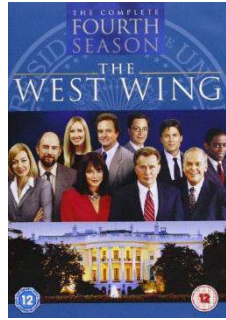


0.05

0.04

0.02

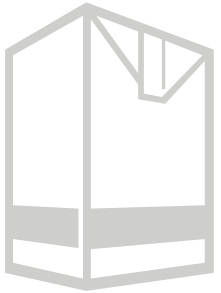
0.02



# ML model

---

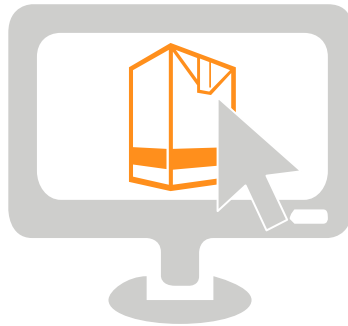
- Logistic regression models
- Features



Product-specific



User-specific



User-product interactions



Display-specific

# Online optimizations

---

- Algorithmic
  - Use simpler ML model
  - Quickly discard candidates
- Technical
  - Memcache + local cache
  - Async I/O

# Upcoming challenges

---

- Long(er)-term user profiles
- More and better product information (images, semantic, NLP)
- Instant-update of similarities
  - (because batch computation is soooo last year)



# Fancy a try ?

---

On your own:

- Our 1<sup>st</sup> public dataset is online: <http://bit.ly/1vgw2XC>
  - 4GB display and click data, Kaggle challenge in 2014
- NEW : 1TB dataset released a few weeks ago: <http://bit.ly/1PyH4Vq>
  - Hosted on Microsoft Azure, just waiting for you

With us !

<http://labs.criteo.com/jobs/>

# Questions?



Thank you !  
s.dolle@criteo.com •

