

Measuring 2.0

ir. drs. Niels Basjes

Lead IT Architect Scalable Solutions



bol.com[®]
de winkel van ons allemaal

Agenda

- **Introduction**
- **Bol.com**
 - Personalization
- **Measuring 1.0**
 - Limitations
- **Measuring 2.0**
 - What do we really want
 - What is possible
- **Stream Processing**
 - Design
 - Implementation
 - Operations



Niels Basjes

nbasjes@bol.com

@nielsbasjes

<https://github.com/nielsbasjes>

**TU-Delft Computer Science
Nyenrode Business School**

**Software developer
Research Scientist (NLR)
Infra Architect (NLR)
WebAnalytics Architect
Lead IT Architect (Bol.com)**

**Contributor Apache Hadoop,
Pig, HBase, Storm, Flink, ...**

Committer Apache Avro



Bol.com

A collage of several screenshots from the Bol.com website, overlapping each other. The screenshots show various parts of the site, including the navigation menu, promotional banners, product listings, and a list of best-selling PC games. The screenshots are tilted and layered, creating a sense of depth and showcasing different aspects of the user interface.

bol.com
Kies een categorie

- Boeken
- Muziek, Film & Games
- Computer & Elektronica
- Speelgoed & Hobby
- Baby & Kind
- Mooi & Gezond
- Sieraden & Accessoires
- Sport & Vrije tijd
- Wonen & Koken
- Tuin & Klussen
- Dier
- Cadeaus & Inspiratie
- Aanbiedingen**

Verkopend Zakelijk Fotoservice Welkom Niels Bestelstatus Alle artikelen Zoeken

Kies een categorie
Gratis verzending vanaf 20
Game- Klantenservice
Lijstjes

30% korting op Pampers Tot 50% op bedtextiel
Gratis verzending vanaf 20 euro, gratis retourneren en 30 dagen bedenktijd*

Carnaval
Dag- & weekaanbiedingen

Ik wil jou
een mooi cadeau geven

Sieraden met een hartje al vanaf 9,-

Tot 60% korting op parfums

Cry Primal
Special edition met 3 extra missies
Vanaf leeftijd 7+

Bestverkochte pc games
Iedere 24 uur up-to-date

- 1 De Sims 4 PC € 49,99
- 2 Rise of the Tomb Raider - PC € 44,99
- 3 Grand Theft Auto V - PC € 42,99
- 4 Football Manager 2017 PC € 49,99
- 5 Fallout 4 - PC € 42,99

Meer bestverkochte pc games

Accessoires
Bestverkochte pc accessoires

- 1 Logitech G430 - 7.1 Virtueel Su... € 74,99
- 2 Microsoft Xbox 360 Controller - ... € 34,99
- 3 Razer Kraken USB Essential 7.1 ... € 59,99
- 4 Turtle Beach Ear Force PX22 Wir... € 69,-

Dagaanbieding
Je hebt nog 14 uur, 15 minuten, 2 seconden

Diffenyz Plato Fonteinset - 36 x 18 cm - Keramiek
De Diffenyz Plato Fonteinset is door ~~149,99~~ **59,99**

ANNO
Anno 2205 - Special Edition PC € 47,99

BREIN DIDACTIEK
Breindidactiek € 24,95
Just Cause 3 - Day One Edition - PC € 87,49

MINDSET
Mindset in het dagelijks leven € 17,95

Mindset, de weg naar



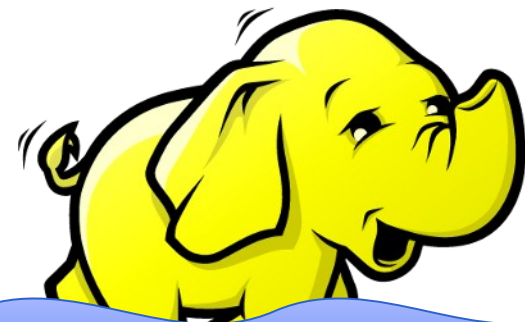
> 11 million products
for sale
> 40 million in catalog



> 6 million active
customers
> 27 million visits per
month



> 2500 million
pageviews/ye
ar



Hadoop in production
since 2010





Happy customers

Change the website to fit the needs of the visitor



Search Suggestions

Verkopen Zakelijk Fotoservice

Welkom Niels ▾ Bestelstatus Lijstjes Klantenservice 

bol.com

Kies een categorie

Boeken

Muziek, Film & Games

Computer & Elektronica

Speelgoed

Baby & Kind

Mooi & Gezond

Sieraden & Accessoires

Sport & Vrije tijd

Wonen & Koken

Tuin & Klussen

Dier

Cadeaubonnen

Aanbiedingen

fifth e|

Alle artikelen ▾

Zoeken



Zoeksuggesties

fifth element in Alle artikelen

fifth element in Dvd

fifth element dvd

fifth element blu ray

fifth elephant

fifth estate

fifth era

fifth essence

Populaire zoekresultaten

 The Fifth Element - Luc Besson
Dvd | 1 disk

 Fifth Element - Charlie Creed-Miles
Blu-ray | 1 disk

Tot 50% op dekbedovertrekken Dag- & weekaanbiedingen

20 euro, gratis retourneren en 30 dagen bedenktijd*

...j ons hebt besteld? Daarom geven we je graag een cadeau. Profiteer direct van

...middelen, babyvoeding, 2de hands en artikelen van externe verkopers. [Bekijk de volledige](#)

sale tot 60% korting
solden



Computerdeals

kortingen op veel
PC artikelen



50% korting

op was- en
schoonmaakmiddelen



Personalization

Long term history

Verkopen Zakelijk Fotoservice

Welkom Niels Bestelstatus Lijstjes Klantenservice

bol.com

Browsing history

Kies een categorie

Boeken

Muziek, Film & Games

Computer & Elektronica

Speelgoed

Baby & Kind

Mooi & Gezond

Sieraden & Accessoires

Sport & Vrije tijd

Wonen & Koken

Tuin & Klussen

Dier

Cadeaubonnen

Aanbiedingen

Frisse Start

Beste van 2015

Tot 50% op dekbedovertrekken

Dag- & weekaanbiedingen

Gratis verzending vanaf 20 euro, gratis retourneren en 30 dagen bedenktijd*

Alle artikelen

Zoeken



Beste Niels,

Wist je dat je al meer dan 100 maal bij ons hebt besteld? Daarom geven we je graag een cadeau. Profiteer direct van € 7,50 korting bij bol.com.

Bij besteding boven €7.50. M.u.v. Nederlandse boeken, geneesmiddelen, babyvoeding, 2de hands en artikelen van externe verkopers. Bekijk de volledige actievoorwaarden.

Wishlist

sale

tot 60% korting

solden



Computerdeals

kortingen op veel PC artikelen

Browsing history



50% korting

op was- en schoonmaakmiddelen



Based on Measuring 1.0

We have some issues
here...



Measurements are...

- **too old.**
 - Available once every 24 hours.
 - So personalization is a 'day behind'.

Je bekeek



Thunderbirds Are Go
- V1

Ook interessant:



Thunderbirds - De Ultieme
Collectie
€ 13,99



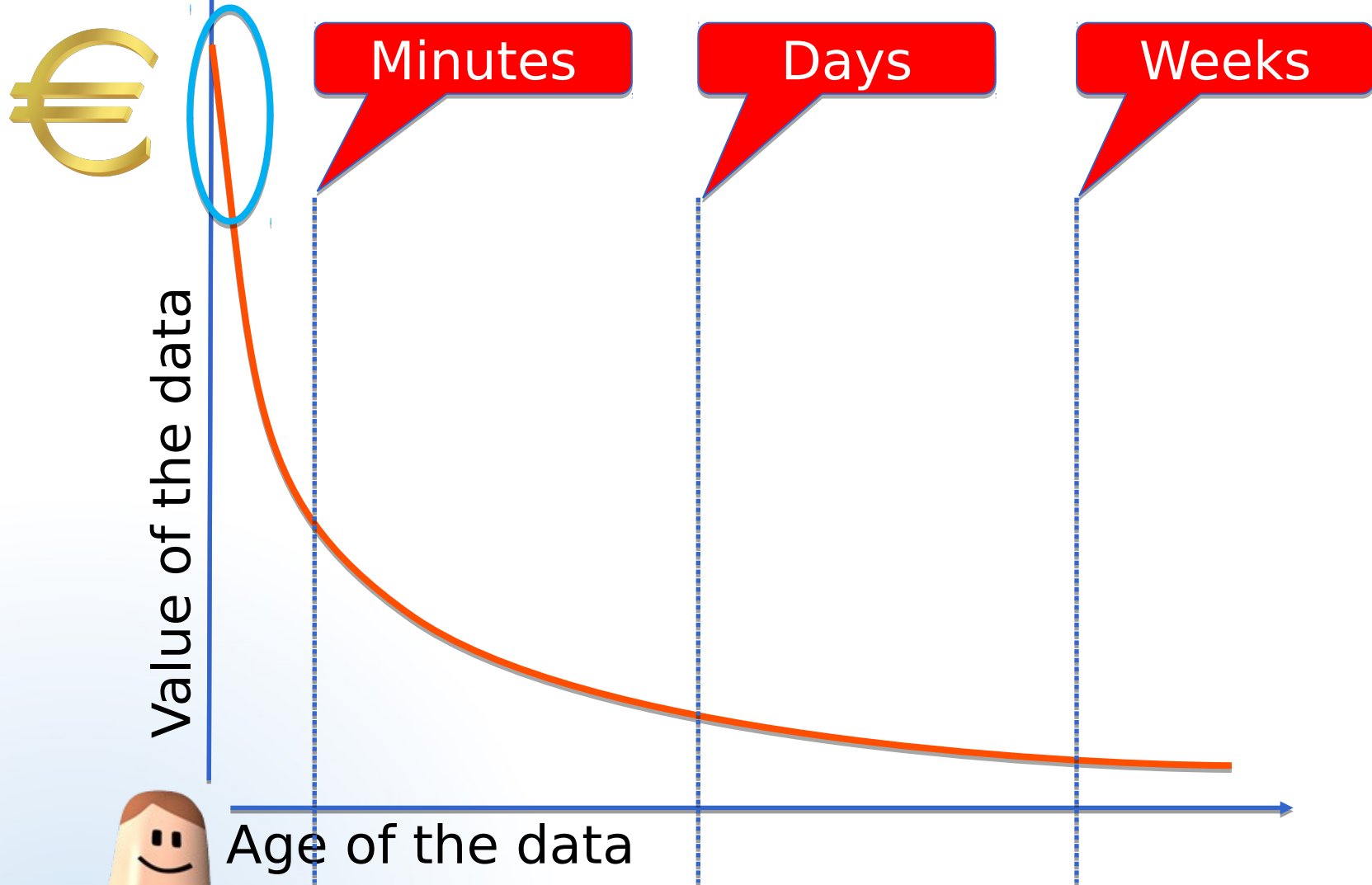
Thunderbirds Are Go
€ 18,99



Thunderbirds 2004
€ 14,99

Useless inspiration:
I was interested in this
YESTERDAY

Data relevance decay



Measurements are...

done using JavaScript ...

```
<script type="text/javascript">
  if(typeof(s)=='undefined'){s = {};s.t = function(){};}

var link = document.querySelector && document.querySelector('link[rel:
link && link.href; if (canonical) { s.prop3 = canonical; s.eVar26 = 'I
s.eVar4=unescape('ps43');
s.pageName=unescape('\n\n/index.html');
s.eVar3=unescape('main');
s.prop2=unescape('Home');
s.prop18=unescape('0');
s.prop1=unescape('main:algemeen');
s.prop29=unescape('VqeYkQpid0IAAELxnY8AAAb4');
s.eVar12=unescape('not logged in');
s.prop25=unescape('niet ssl');
s.prop50=unescape('\n\n/');
s.eVar51=unescape('www.bo1.com,DESKTOP');
s.measureCookies=unescape('false');
```

Measurements are...

- **Heavy on clients**
 - Everything is javascript
- **We cannot measure “Everything”.**
 - Sometimes we run out of eVars.
 - No personalized banners.
- **Unclear/unspecified/complicated.**
 - Errors because eVars get ‘reused’ over time.



Measurements are...

- **Incomplete/Inaccurate.**
 - JavaScript is sometimes slow
 - Everything goes on the URL

```
s.products=unescape(';1001004007516  
421;;;;evar3=books|  
evar7=1001004007516421_Pluk van de  
Petteflet|evar35= 171|evar41=KI1200|  
evar47=B|evar9=New');
```



All special deals page

<http://www.bol.com/nl/aanbiedingen/index.html>

De allerbeste aanbiedingen!

Vind ik leuk

553

Tweeten

18

Volgen



Dagaanbieding!

Ontvang elke dag de
dagaanbieding in je mail

Je hebt nog 13 uur, 36 minuten, 43 seconden

Aquaplay draagbare Waterbaan AquaLock 616 - Waterbaan

Prachtige draagbare waterbaan van Aquaplay. Dirigeer jouw bootjes door de kanalen en sluisen. Leer spelenderwijs hoe water beweegt.



€49,99

€ 29,99

Je bespaart: 40%

Vandaag voor 23:00 uur
besteld, morgen in huis.

+ In winkelwagentje

★★★★★ (7 reviews) > Lees reviews

Waterbaan | 88 cm | 3 - 6 jaar | 3 - 6 jaar | Waterbaan

De draagbare AquaLock 616 is een van de bestsellers van AquaPlay. Neem de draagbare AquaLock 616 waterwereld

Weekaanbiedingen!

Je hebt nog 4 dagen, 13 uur, 36 minuten, 43 seconden

34 weekaanbiedingen



Dead Or Alive (3DS)

★★★★★ (3 reviews)

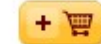
Nintendo 3DS

Vandaag voor 23:00 uur besteld,
morgen in huis.

€29,99

€ 9,99

Je bespaart: 67%!



Electrolux UltraCaptic Stofzuiger

Zakloos | Hepa | 1400 W |
Stofzuiger | 1400 W

Vandaag voor 23:00 uur besteld,
morgen in huis.

€379,00

€ 179,00

Je bespaart: 53%!



Sam Smith - In The Lonely

€17,99



So...

things need to change



Business goal

**We want to be able to know
and process everything
about our webshop
NOW**

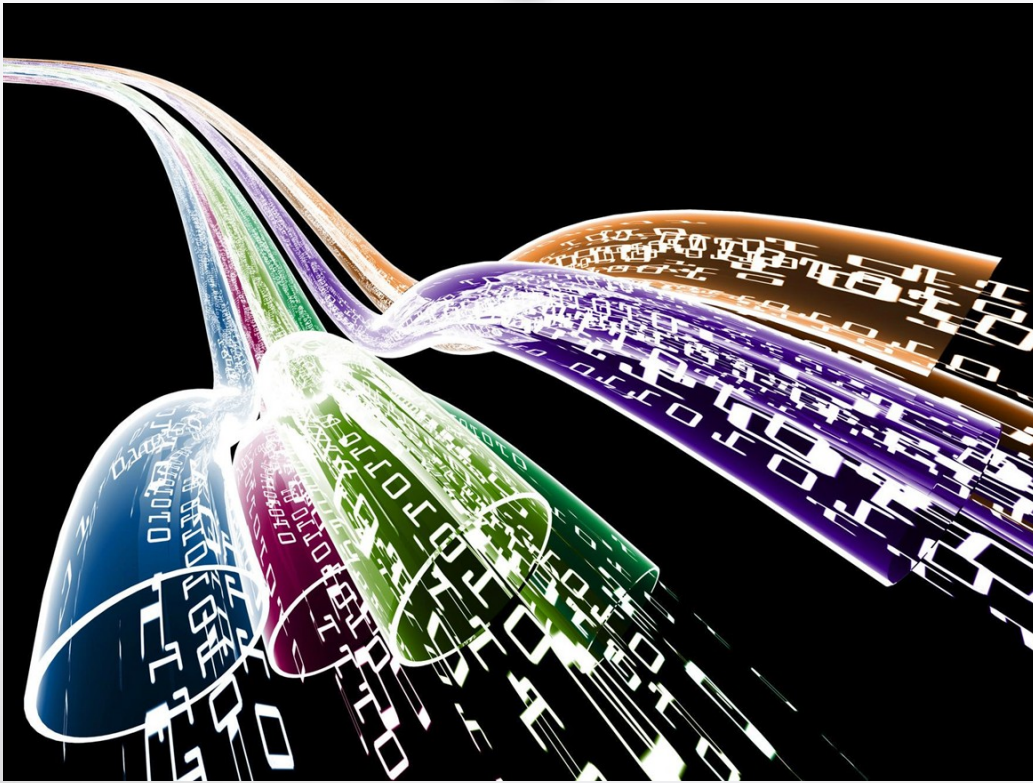
**Because we want to be able to
assist our visitors
NOW**



Business goal

- **While the visitor is still on the site**
 - Respond to what they have done so far.
 - Last months/year
 - Within seconds after the 'click'.
- **Batches multiple times per day**
 - Search suggestions
 - Search rank
 - Recommendations





Measuring 2.0



It's really all about...

- **Measuring**
 - Better
- **Processing**
 - Faster
- **Applications**
 - More relevant



Goals of “Measuring 2.0”

- **Measure everything of our website**
 - All interactions (also AJAX)
 - All channels (also mobile, email, ..)
 - All countries
 - All details
 - All visitors (also Googlebot)
- **More reliable data**
- **Lowest possible load on the client**
- **Lowest possible latency (< 1 second)**

AMBITIE

NOU GEWOON

ALLES

Loesje



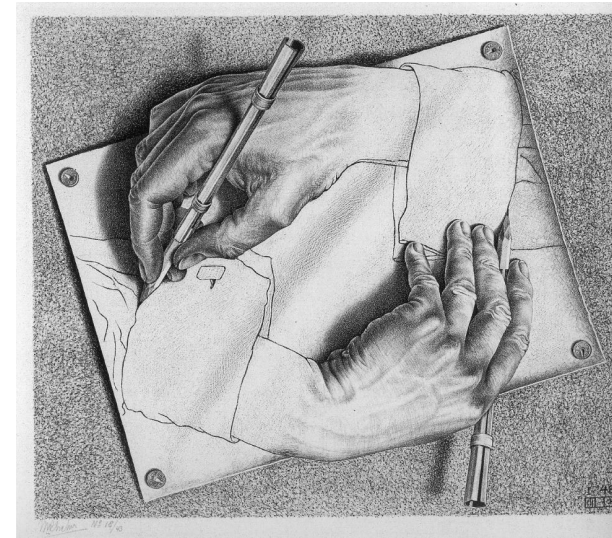
Goals of “Measuring 2.0”

- **Developer**
 - Easy to build
 - Easy to validate
 - Test automation
- **Business**
 - Always measure everything
 - Data is “independent”
 - New questions are allowed



Goals of “Measuring 2.0”

- **Privacy**
 - Personal Data Protection Act
 - *Wet bescherming persoonsgegevens*
 - No profiling beyond 2 years
- **Security**
 - Avoid storing “login” info.
- **Business**
 - Do “profiling” on everything
 - For many years (>2)



Measure everything

Boeken > Nederlandse boeken > Kind & Jeugd > Voorlees- & Sprookjesboeken > Pluk van de Petteflet

Pluk van de Petteflet

Herziene Editie.Nu Met Nieuwe Illustraties

Auteur: Annie M.G. Schmidt | ★★★★★ 39 reviews | Stel een vraag

Favoriet 8



Inkijkexemplaar

Bindwijze: Hardcover

Samenvatting

Pluk had een klein rood kraanwagentje. Hij reed ermee door de hele stad en zocht naar een huis om in te wonen. Af en toe stopte hij. En dan vroeg hij aan de mensen: 'Weet u niet een huis voor me?'

Uiteindelijk vindt hij wat: het torentje van de Petteflet. Daar maakt hij kennis met de Stampertjes, mevrouw Helderder, Aagje en Zaza, maar ook met andere dieren als Dollie, Langhorns en de Krullenaar kan hij het goed

 **Annie M.G. Schmidt**
Auteur

Blijf op de hoogte

Co-auteur: Fiep Westendorp
Uitgever: Querido Kinderboek

- Nederlands
- 200 pagina's
- Querido Kinderboek
- april 2010

> Alle productspecificaties

€49,95 **16^{.96}**

Je bespaart 15%

Vandaag voor 23:00 uur besteld, morgen in huis

Prijs en levertijd

Verkoop door bol.com

- ✓ Gratis verzending vanaf 20 euro
- ✓ 30 dagen bedenktijd en gratis retourneren
- ✓ Ophalen bij een bol.com afhaalpunt mogelijk
- ✓ Dag en nacht klantenservice
- ✓ Cadeautje? Laat het voor je inpakken en bezorgen

+ In winkelwagentje Zet op verlanglijstje

Andere verkopers (11)

Tweedehands

vanaf € 11,95

> Bekijk en vergelijk alle verkopers

vanaf € 16,96

Anderen bekeken ook:



Jip en Janneke
Annie M.G. Schmidt
€ 25,95



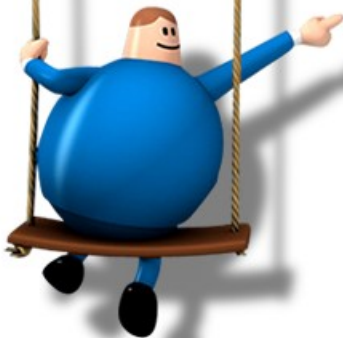
Floddertje
Annie M.G. Schmidt & Fiep Westendorp
€ 15,95



Pluk redt de dieren
Annie M.G. Schmidt & Fiep Westendorp
€ 15,95



Otje
Annie M.G. Schmidt
€ 17,50





Cause and Effect

is what really interests us



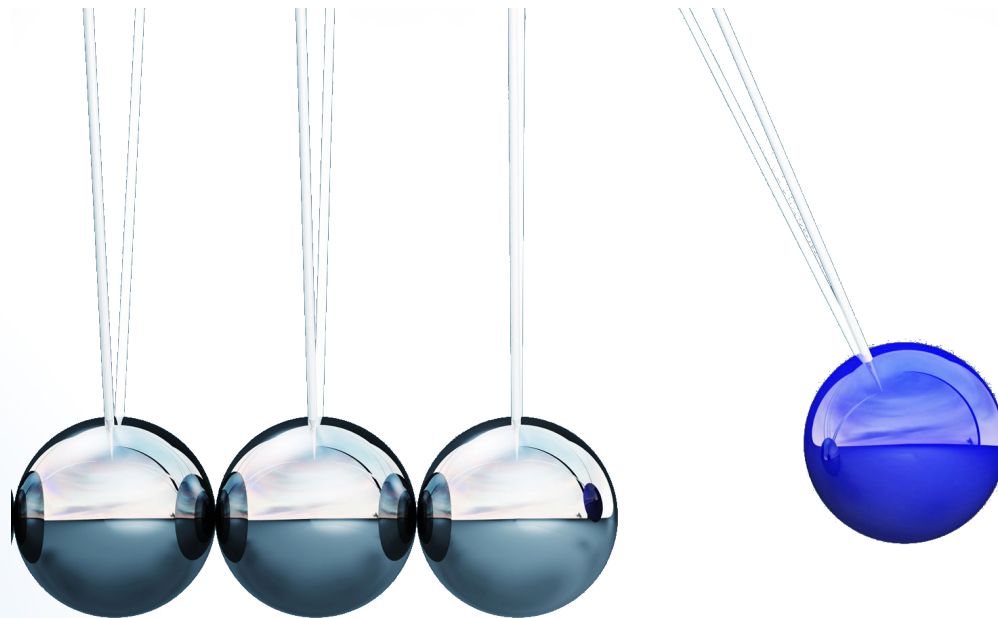
Our usecases

- **Banner optimization**
 - Look / Search → Next page better banner
- **A/B testing**
 - Show feature → Use → Buy product
- **Search Suggestions**
 - Search → Find → Choose → Buy product
- **Attribution modeling**
 - Show ad → Click → Buy product
- ...

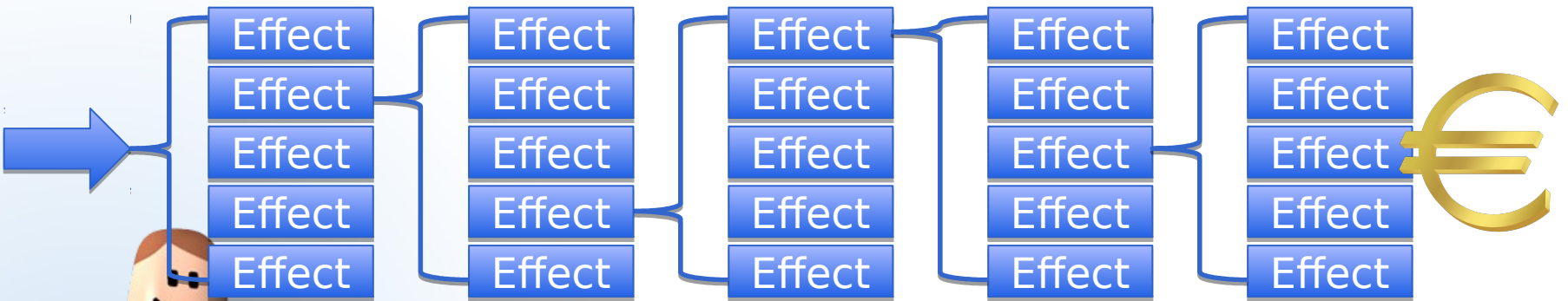
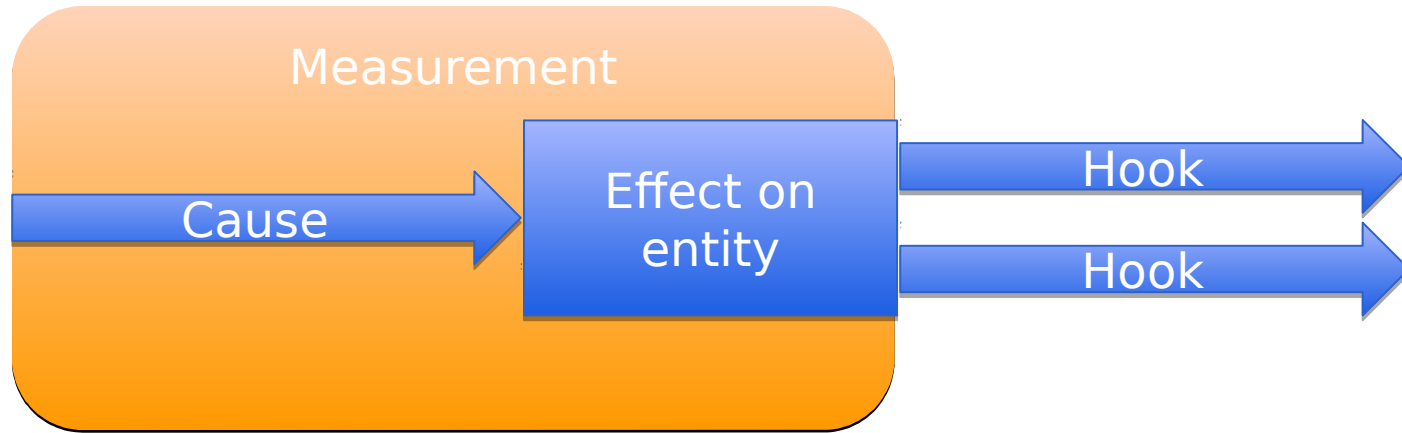


Behavioral analytics

- **Cause and effect**
 - **Action:** We show something
 - **Reaction:** To click or not to click



The “measurement”

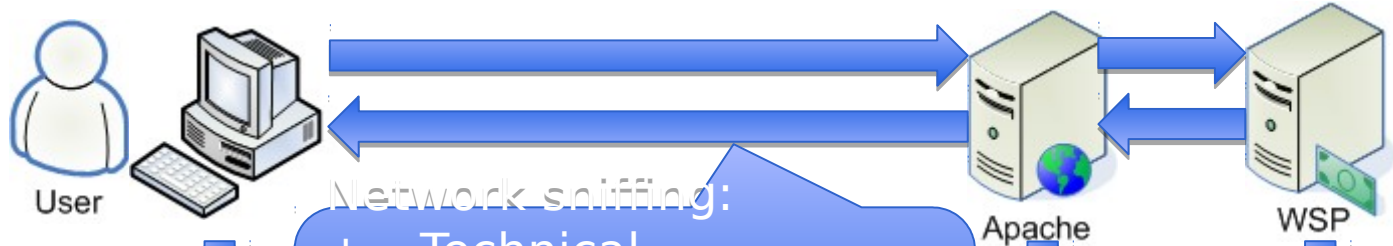




How do we measure?



How can we measure?



Network sniffing:
+ Technical performance.
+ All visitors
- Not "what has been seen"

Access logs:
+ Technical performance.
+ All visitors
- Only 'URLS', no details.

JavaScript:
+ What has been seen.
+ Client side performance
- Slow & Unreliable
- Limited details

Application logs:
+ All details: what and why
+ All visitors
- Not "what has been seen"



Hybrid approach



**We measure everything
serverside**
(unless it really cannot be measured that way)

**If we need to measure
clientside then it must be as
lightweight as possible.**

JavaScript:

- + What has been seen.
- + Client side performance
- Slow & Unreliable
- Limited details

Application logs:

- + All details: what and why
- + All visitors
- Not "what has been seen"

- Not all visitors





Processing the data



Requirements

- **Online:**
 - Near real-time (< 1 second)
 - Have long history available
 - Seamless integration history + real-time
- **Offline:**
 - Incremental batch jobs during the day
 - Have long history available



Lambda Architecture

- **Basic idea:**

- Batch tools are good at 'a lot of data'
- Realtime tools are good at 'low latency'

- **Lambda Architecture:**

- Combine the two at query time.

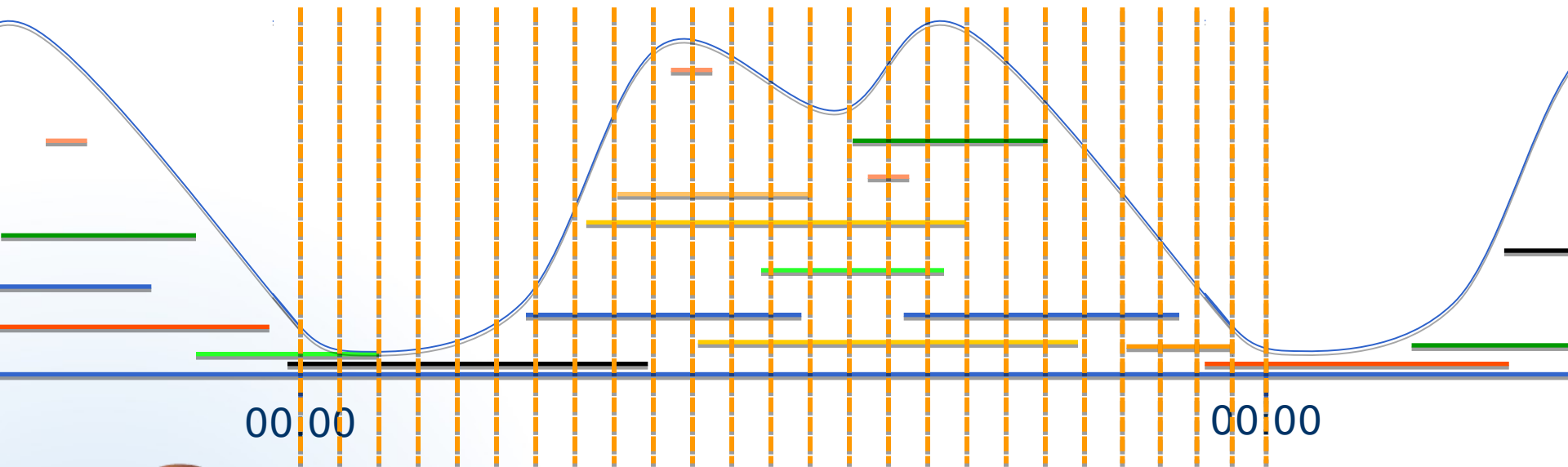


Cut boundary at a millisecond.



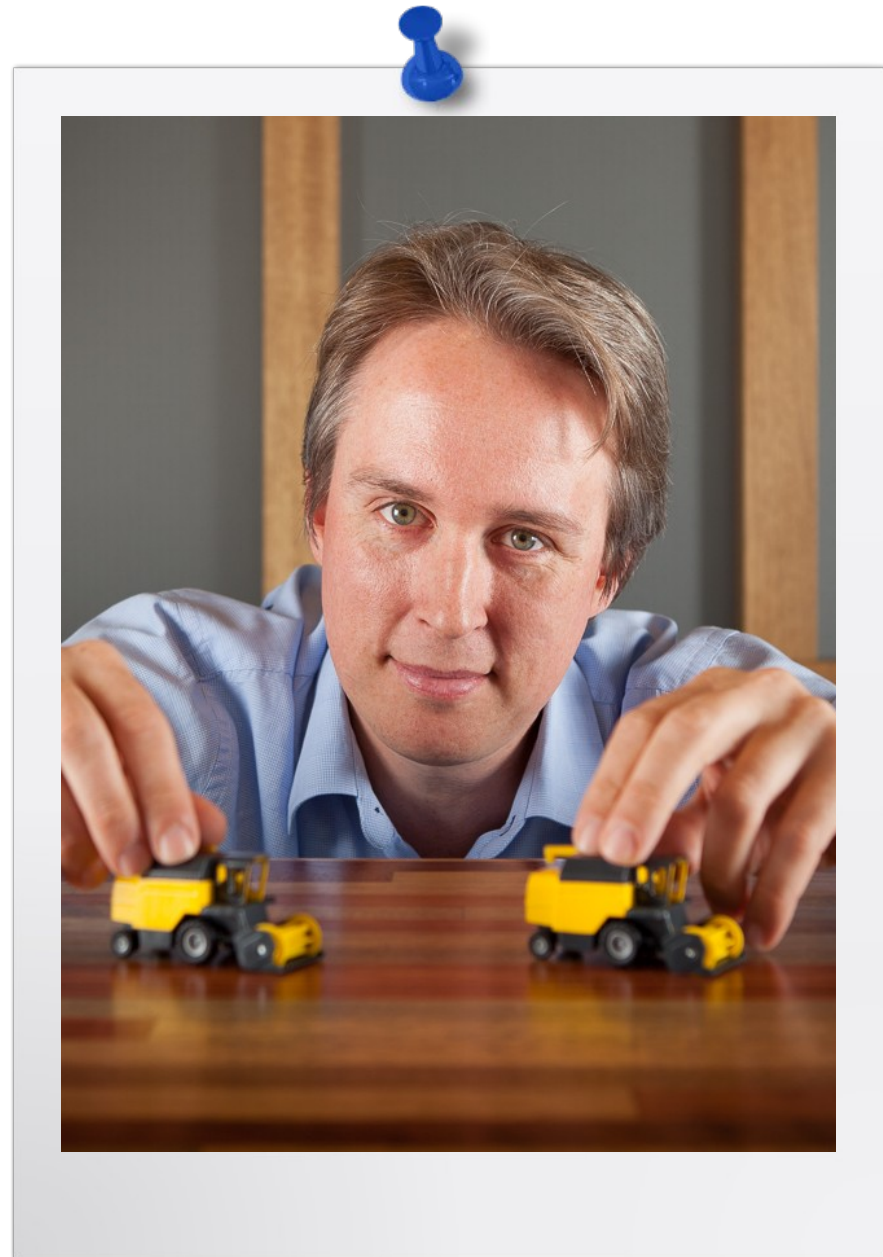
Many batches per day

Joining the pieces becomes very hard



Sessionized Lambda Architecture

by Niels Basjes (Bol.com)



Sessionized Lambda Architecture

Extension of the Lambda Architecture

- **Bounded event streams**
 - that stay together
- **Queries take time**
 - and multiple can overlap
- **Service orientation**
 - because it is part of a bigger thing



Focus on Visits

- **Browsers**
 - The software installed on a computer
- **Sessions**
 - Start: Visit website
 - End: Close browser
 - Can last for days ... weeks ... !
 - Device is 'suspended'
- **Visits**
 - Start: Visit website
 - End:
 - 30 minutes idle
 - max 12 hours active

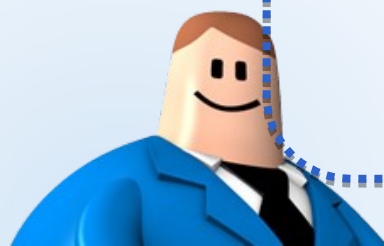
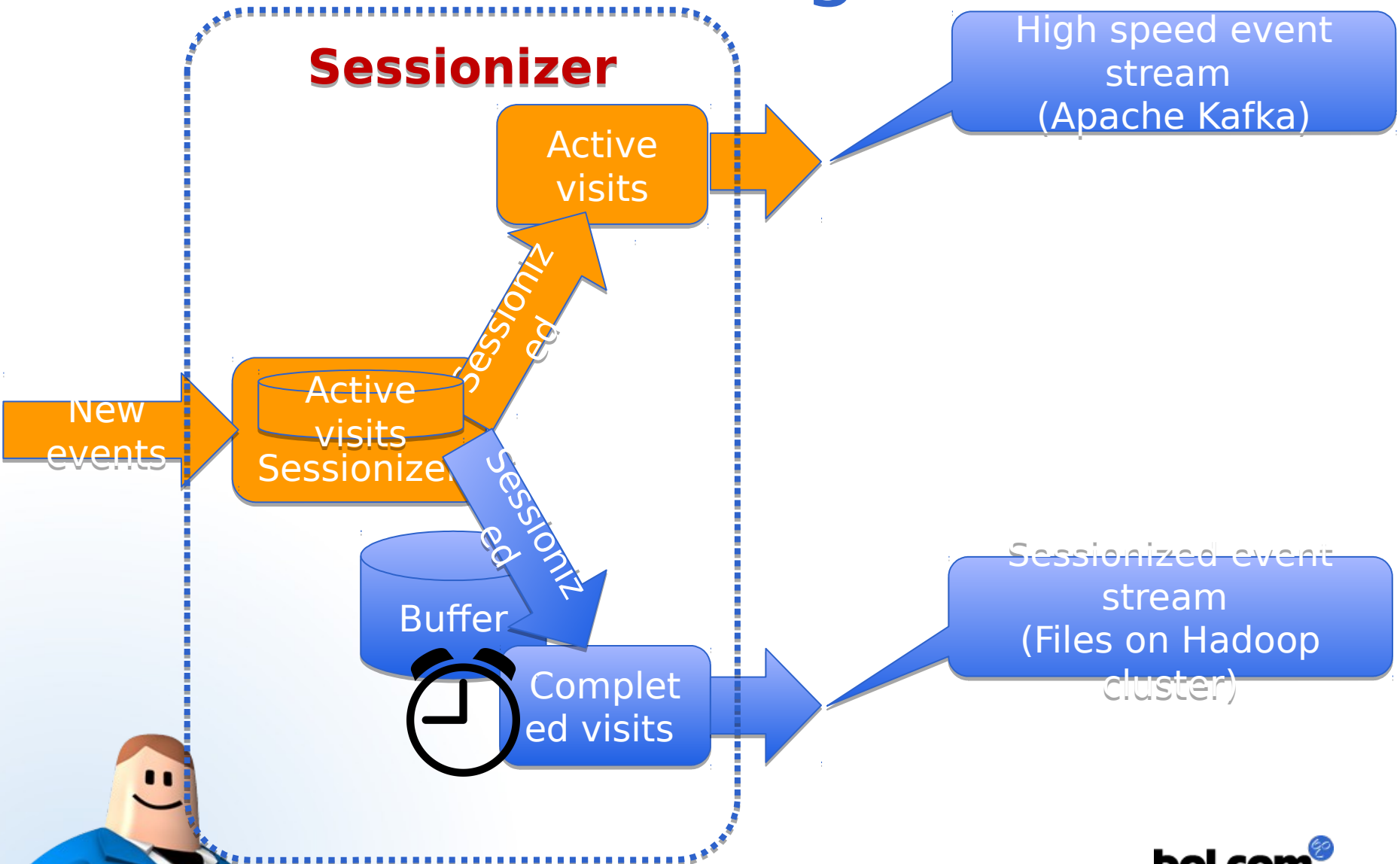




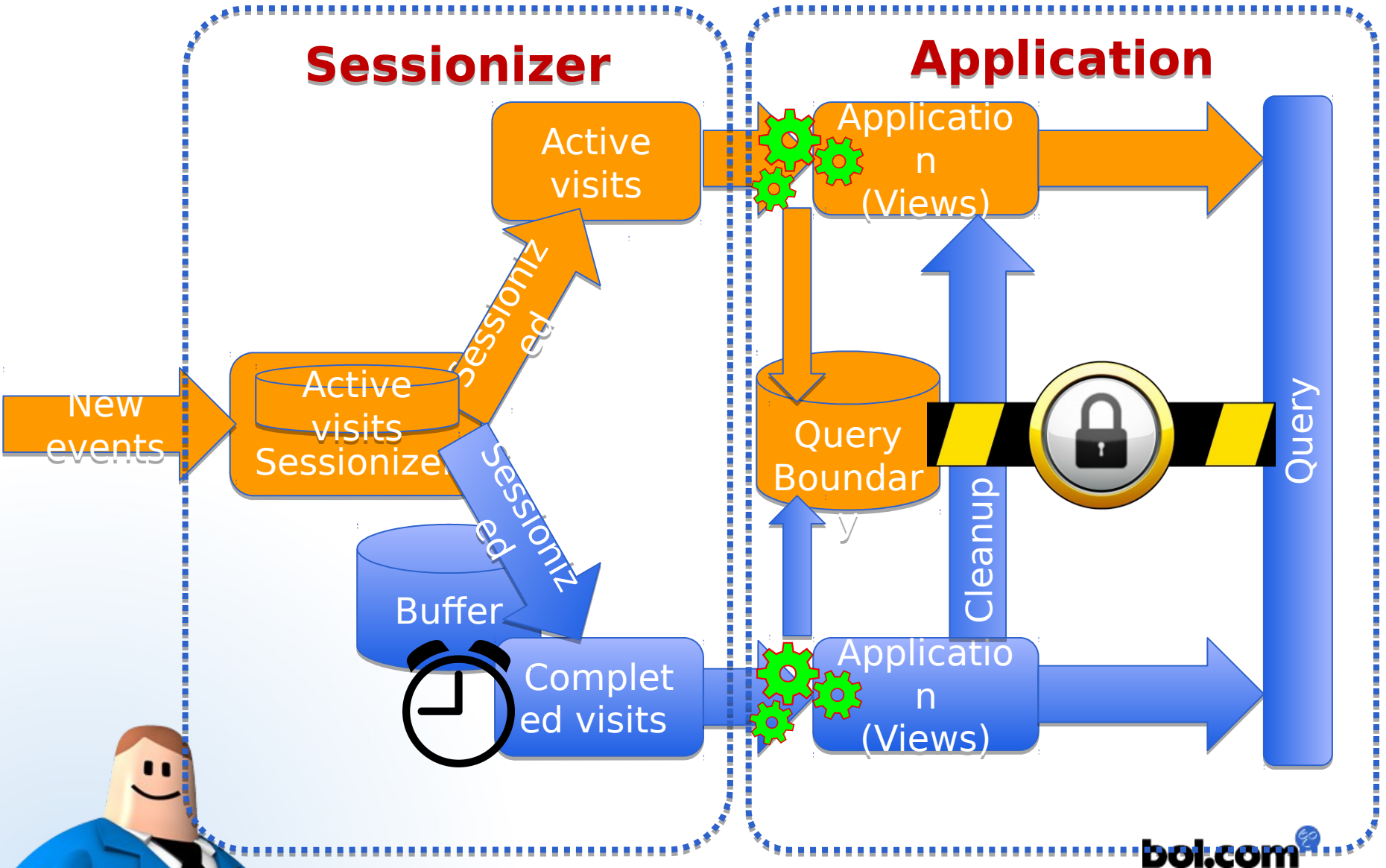
Processing Architecture



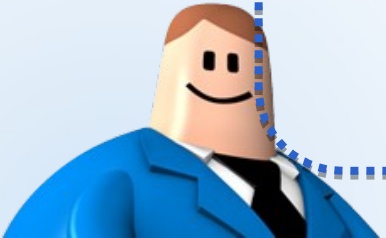
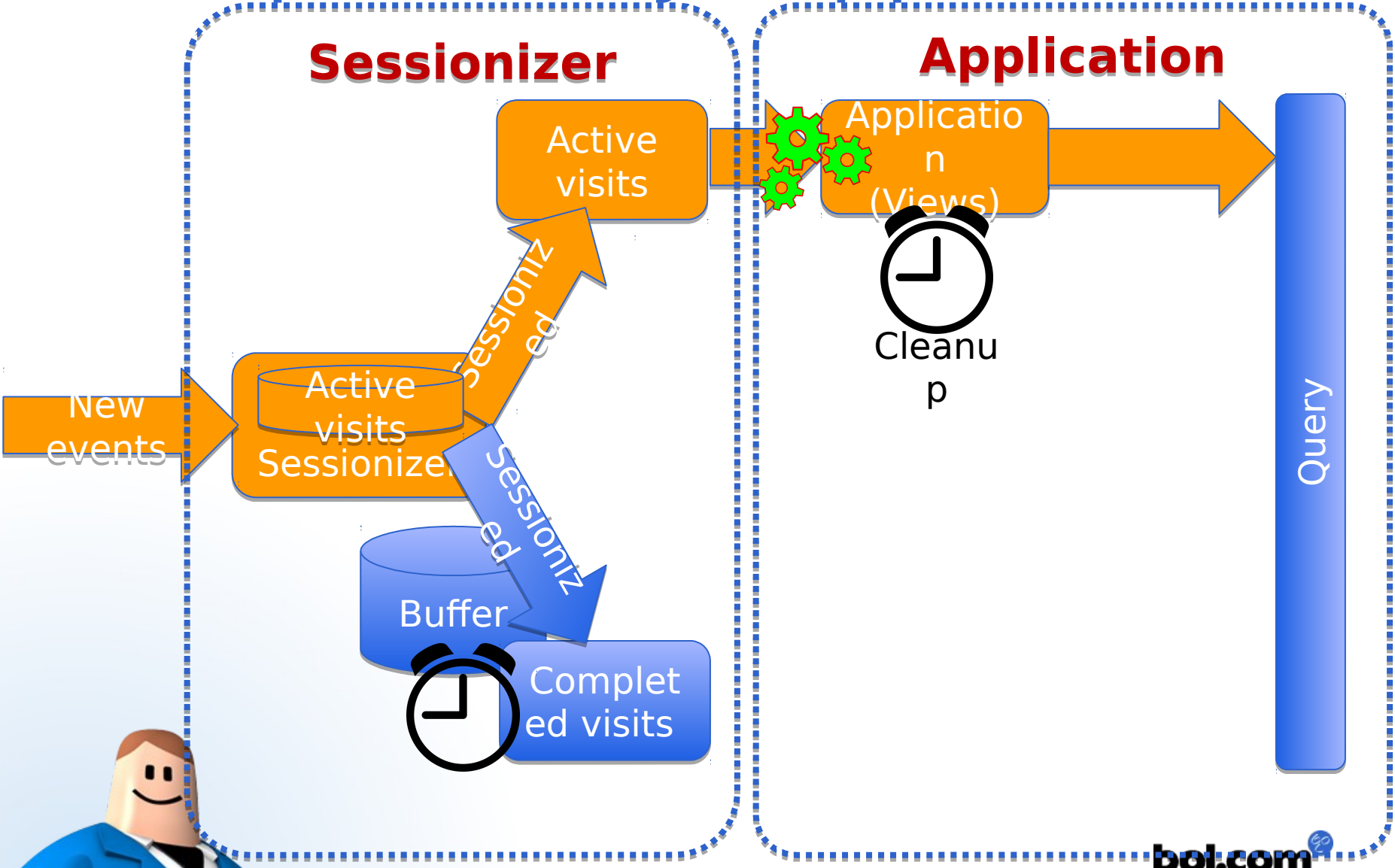
Measuring 2.0



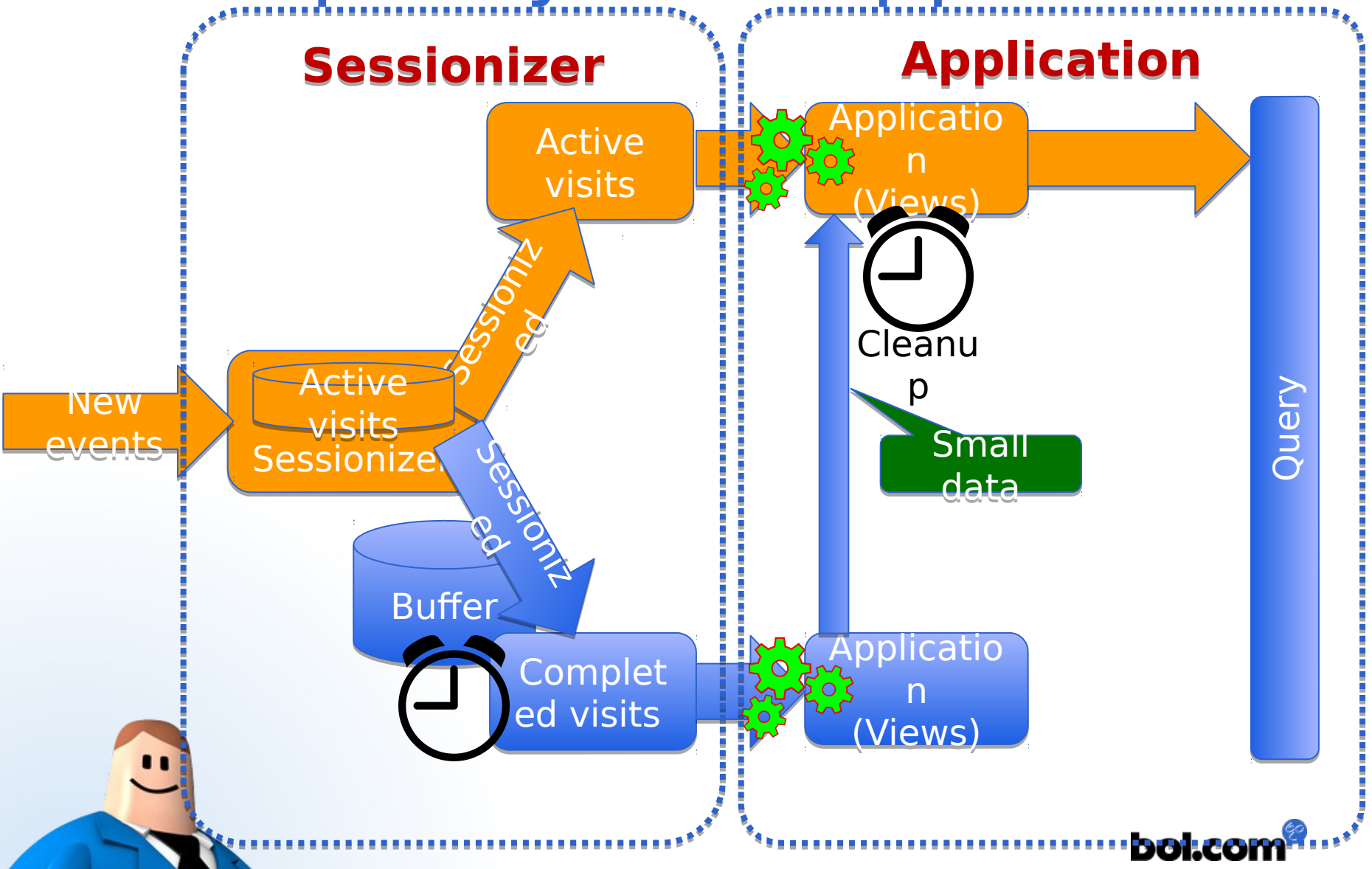
“Full” Sessionized Lambda



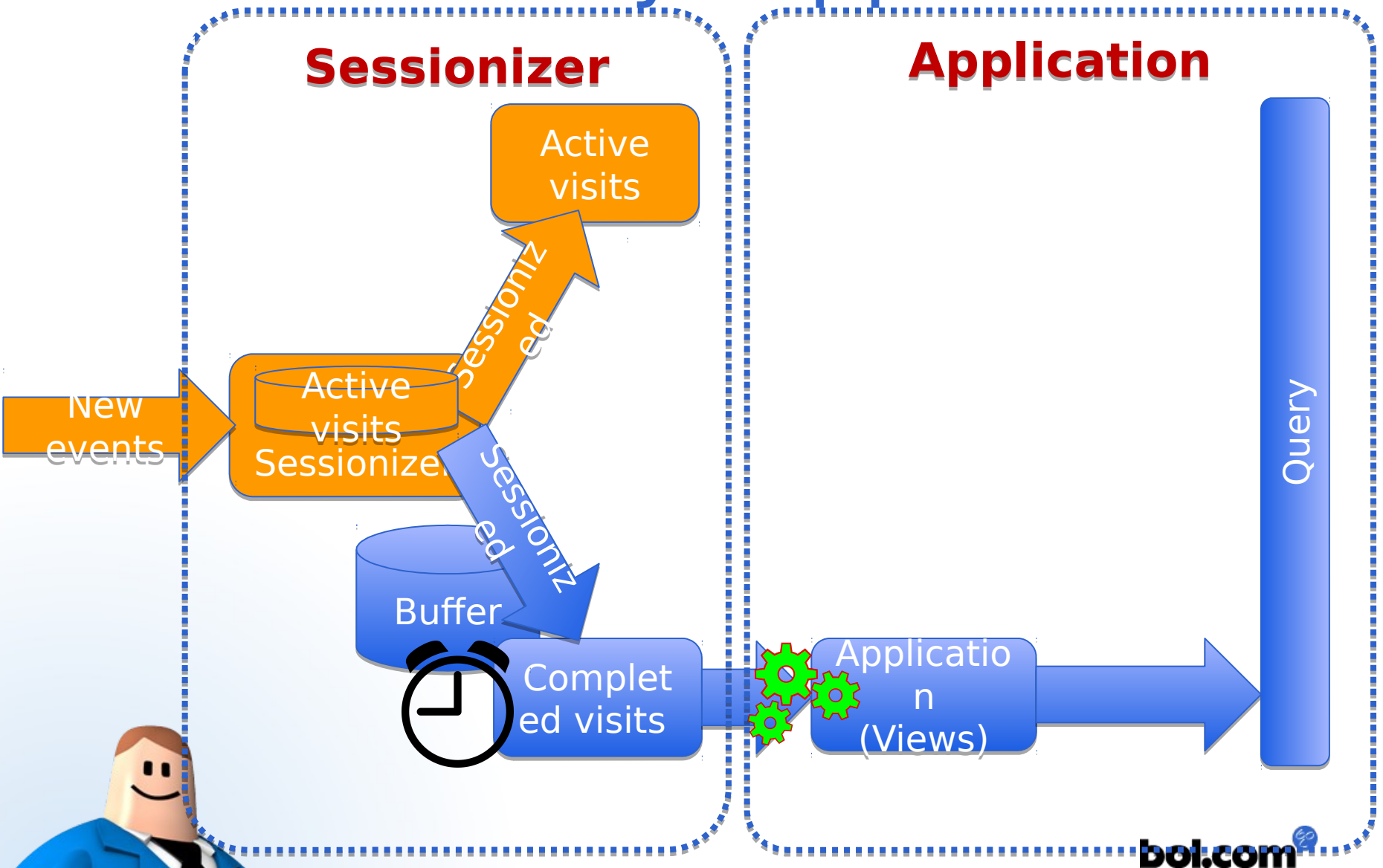
“Speed only” application



“Simple hybrid” application



“Batch only” application





Privacy vs Profiling



Privacy protection

Identifiable data may not be used for profiling for more than 2 years *)



*) My wording & understanding, I'm no lawyer.

Identifiable data

(Groups of) data elements that can identify a single person

In the clickstream:

- **IP Address**
- **Customer number**
 - *After logging in*
- **Unique browser ID**
 - Random value in cookie
 - *Visitor can clean this*



Long term “profiling”?

“Looking for behavior patterns”

- **Query pattern:**
 - GROUP BY `${CustomerId}`
- **Same result**
 - GROUP BY **ENCRYPT**(`${CustomerId}`)
 - GROUP BY **HASH**(`${CustomerId}`)
 - GROUP BY **HASH**(`${CustomerId}`-**SALTED**)



Hash collisions
are not important



Our solution

- **Data in motion (streaming)**
 - Identifiable & Anonymous
 - Kafka expires after 4 weeks
- **Data in rest (files on disk)**
 - Identifiable
 - Delete files after 2 years
 - Anonymous
 - Hash with yearly secret salt.
 - *So only 'group by' within a specific year*



What can we analyze?

We can analyze the exact behavior of

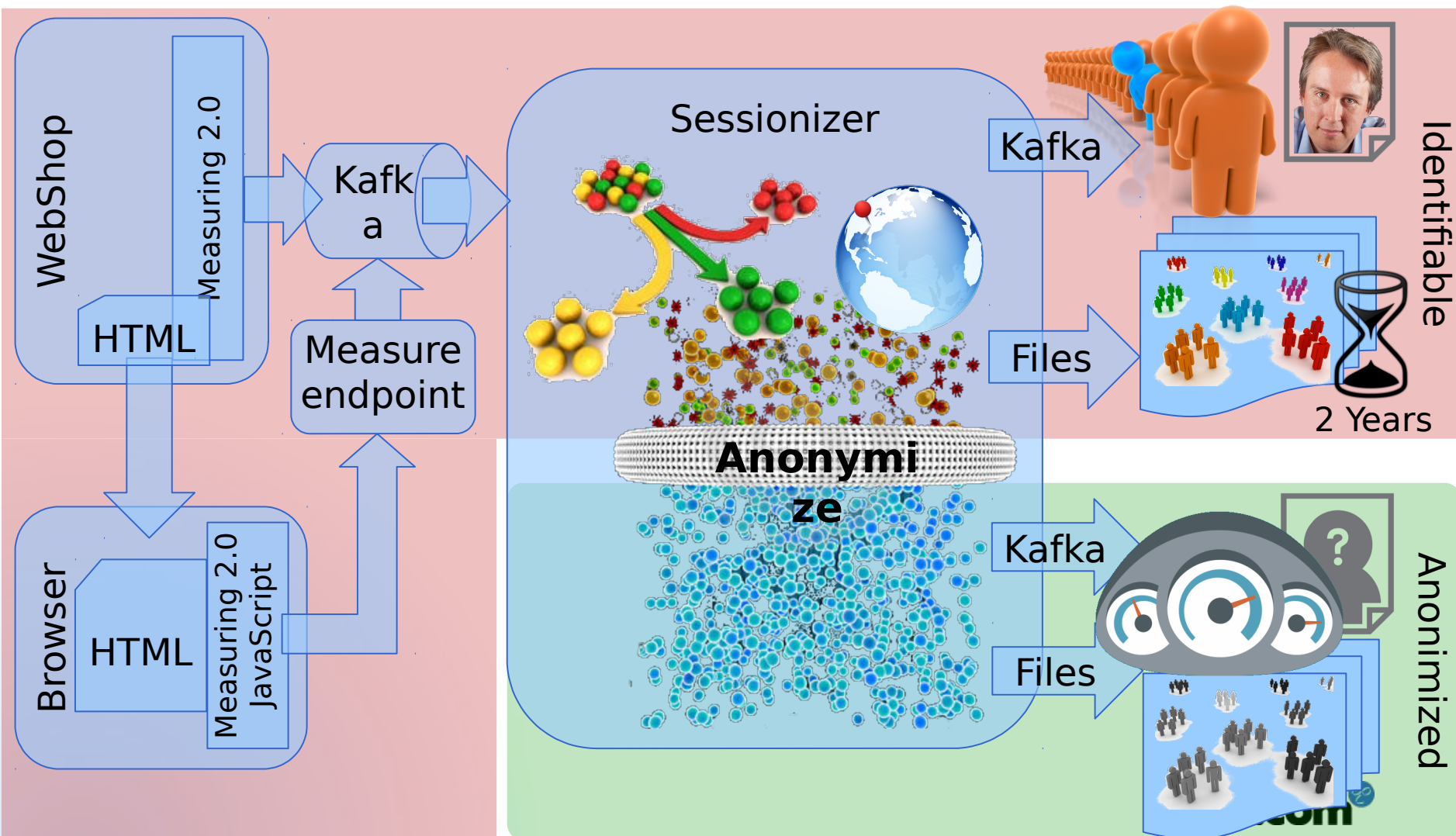
- **Named individuals**
 - in the last 24 months.
- **Anonymous individuals**
 - per year for many years.



Implementation



Implementation design



2014: First prototype



Apache Storm

- **Stream processing**
- **Low latency**
- **Many users**



Challenges

```
TopologyBuilder builder = new TopologyBuilder();
builder.setSpout("words", new TestWordSpout(), 10);
builder.setBolt("exclaim1", new ExclamationBolt(), 3)
    .shuffleGrouping("words");
builder.setBolt("exclaim2", new ExclamationBolt(), 2)
    .shuffleGrouping("exclaim1");
```

- **Hard to program**
 - Unfriendly API
- **“At least once”**
- **Need “Exactly once”?**
 - Trident □ Micro batches
 - No more low latency
- **Storm is stateless**
 - State is an application problem... □
 - No recovery of state after failure
- **Run on an existing cluster?**
 - I never got it to run on Yarn (in 2014!).

Make sure
these names
match!



Spark Streaming ?

- **What about Spark Streaming?**
 - Micro batches
 - Too much latency





Apache Flink



Apache Flink

- **Low latency**
- **Exactly once**
 - With recovery after failure
- **Runs on Yarn**



Flink supports our goals

- **Manages state in the framework**
 - and saves it in case of failure
 - CheckPoints & SavePoints
- **Windows**
 - The basic component for ‘visits’
- **“Event time”**
 - A ‘time out’ is based on the events.



Challenges

- **Running on Kerberos secured cluster**
 - I fixed that for HBase (FLINK-2977)
 - It dies after 173.5 hours
 - common problem on secured Yarn
 - Delegation Tokens problem (HDFS-9276 ?)
- **My “Windows” are too big to fit in memory.**
 - Only keep the visit ‘state’ in memory
 - Persist the events in HBase
 - TODO: Evaluate RocksDB
- **Exactly once**
 - Not on Kafka output!

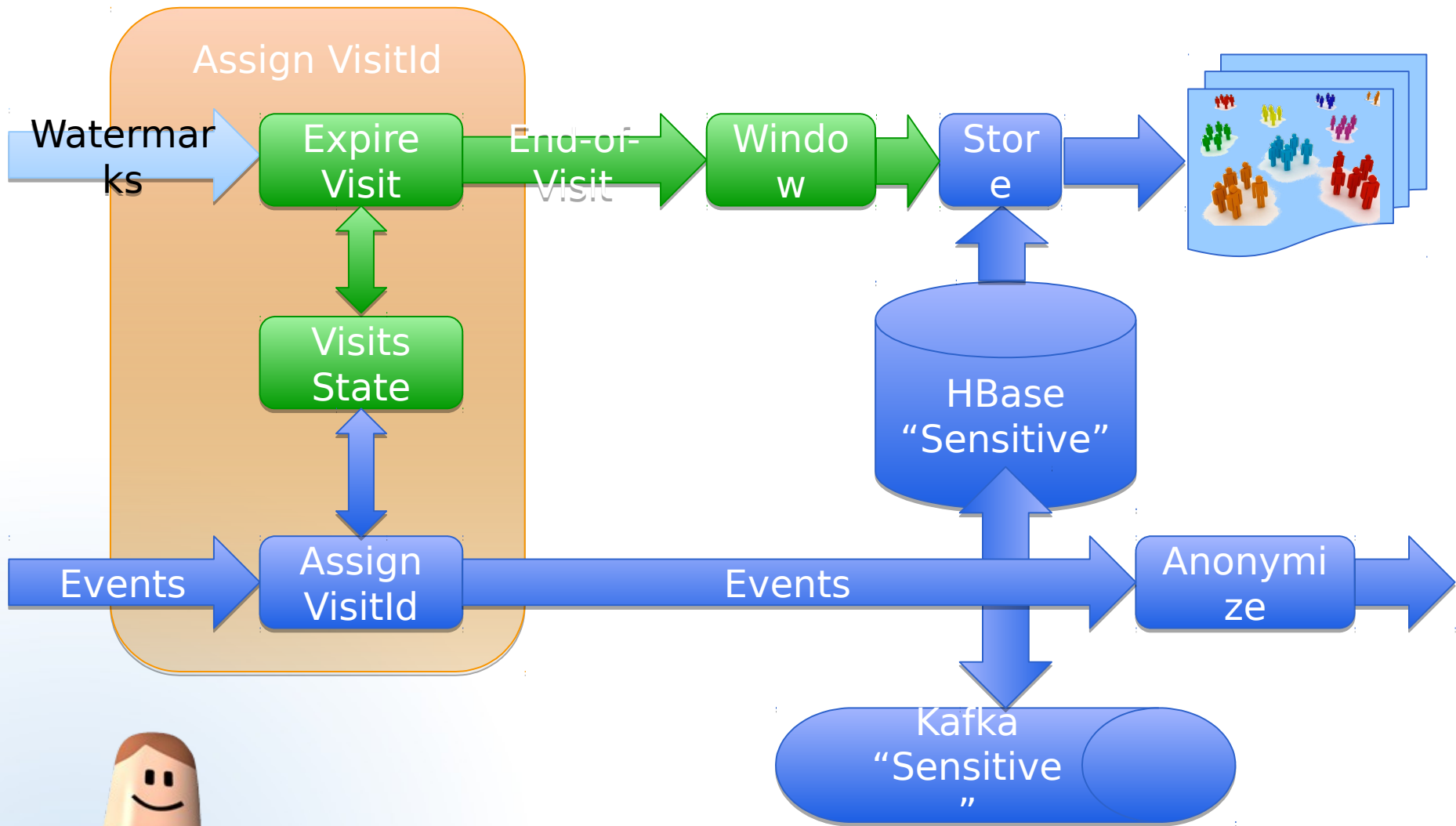


Assign VisitId?

- **Window**
 - Only releases the records after “PURGE”
 - Realtime stream?
- **Custom operator**
 - Must handle everything manually



Assign VisitId Operator



Implementation

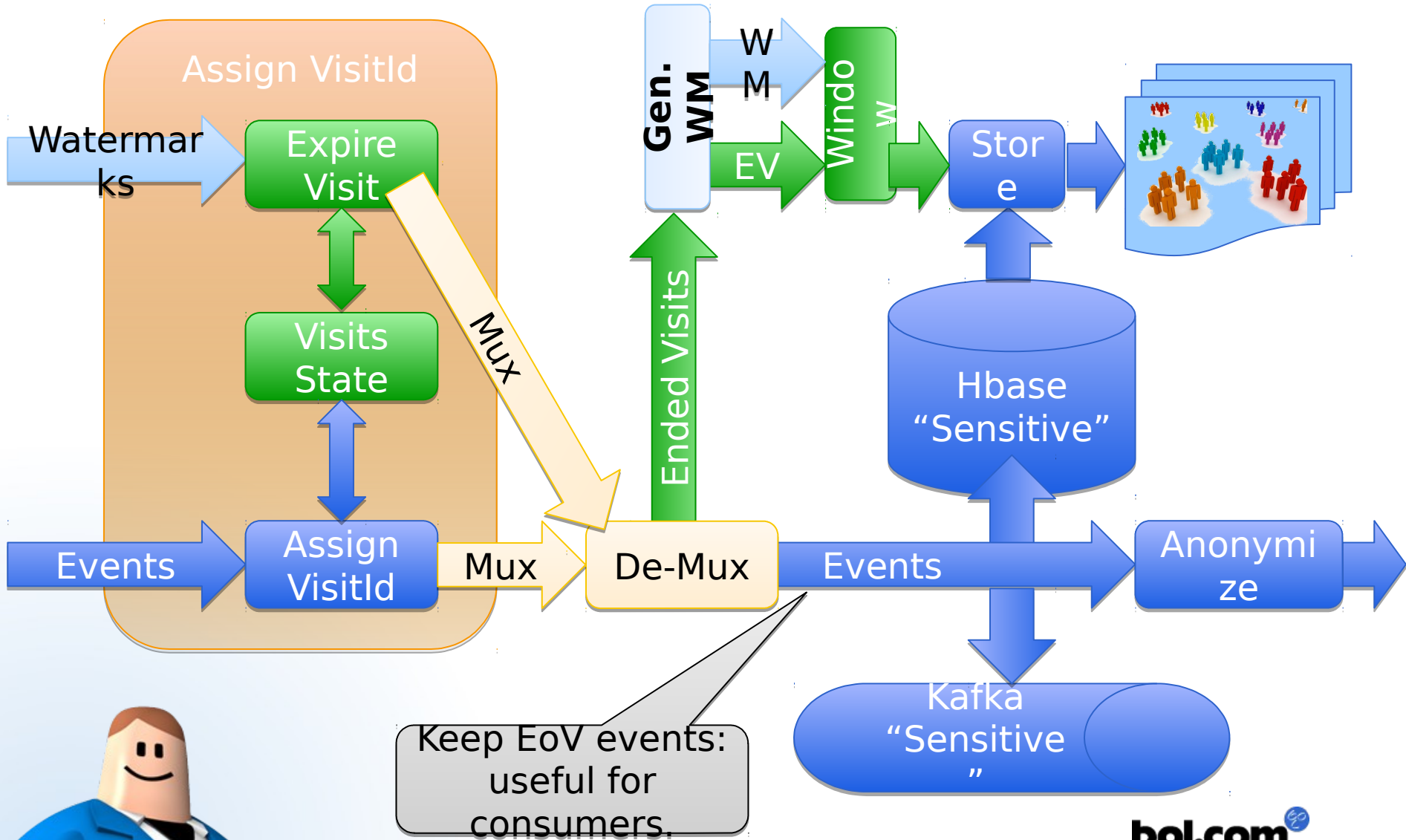
- **Caveats**
 - An operator has only 1 type of output
 - Watermarks are essential for Windows



Watermarks trigger window evaluation.



Final plumbing



Keep EoV events:
useful for
consumers.

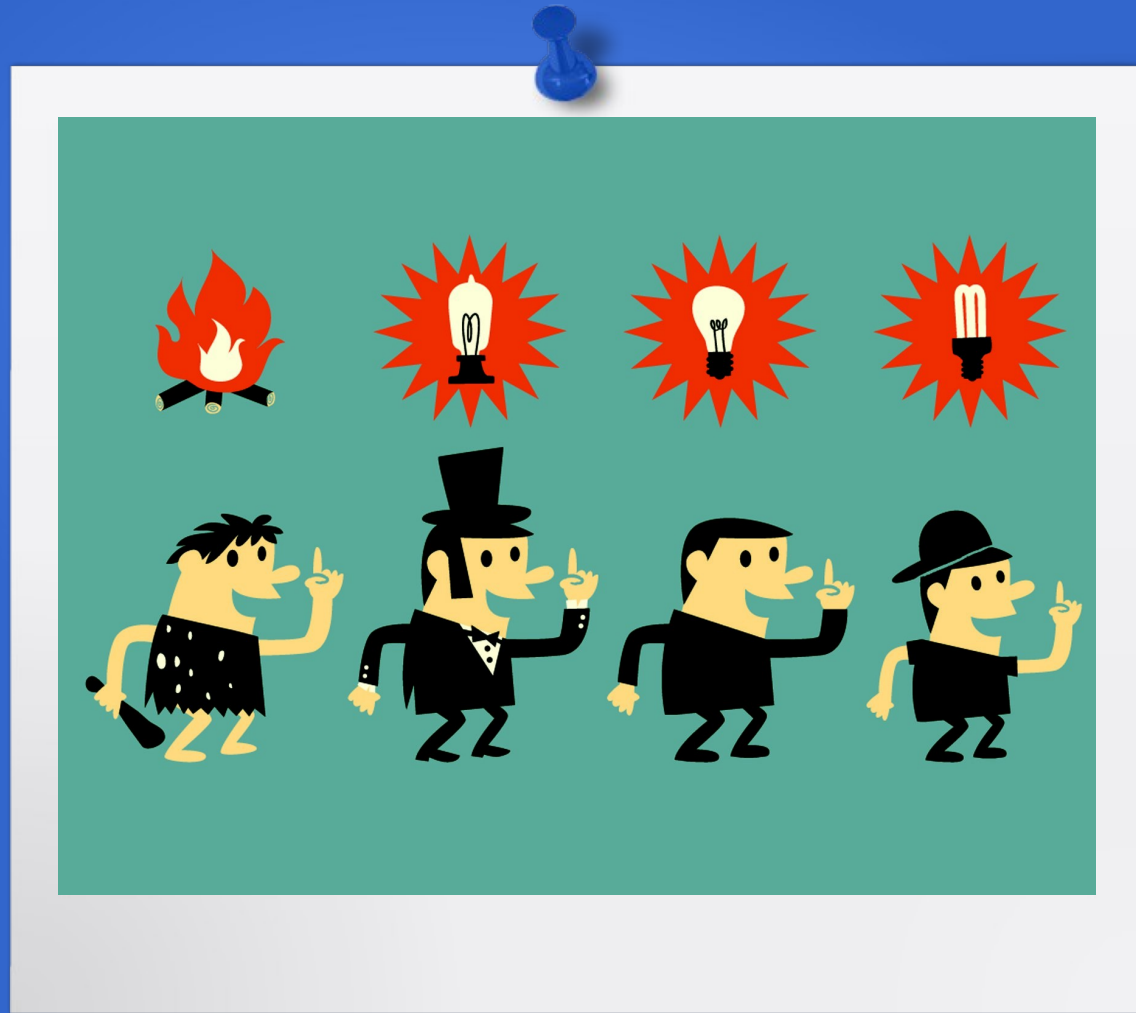
Debugging watermarks

- **Watermarks are**
 - The core of building Windows
 - Tricky to get right
 - Input data is 'messy'
 - You can get a second 'same window'
- **So**
 - Need debugging for the combination of watermarks & events.



Evolution in streaming data

Because requirements change over time.



Persisting data

- **Define Schema in: Apache AVRO**

- Nice schema language (IDL)
- Generates Java classes
- Supports schema evolution



- **Persist in files: Apache Parquet**

- Compresses well
- Great on read use cases
- Write straight from Avro classes.
 - See `o.a.parquet.avro.AvroParquetWriter`



Persisting data

- **Stream: Apache Kafka**
 - High throughput low latency event pipe
 - Persists records for a few weeks
 - Requires the record to be a byte[]
- **Question:**
 - How to serialize record into byte[]?



Persisting data

- **Apache AVRO**
 - Standard byte[] for the record.
 - Needs the original Schema on read !!
- **What if the Schema changes ???**
 - Kafka keeps the events for weeks
- **AVRO-1704 (Work in progress)**
 - Serialize record to a “Message” byte[]
 - schema fingerprint (hash)
 - pluggable schema ‘database’



This is no hype

Thursday March 27th 2014



Thuiswinkel awards

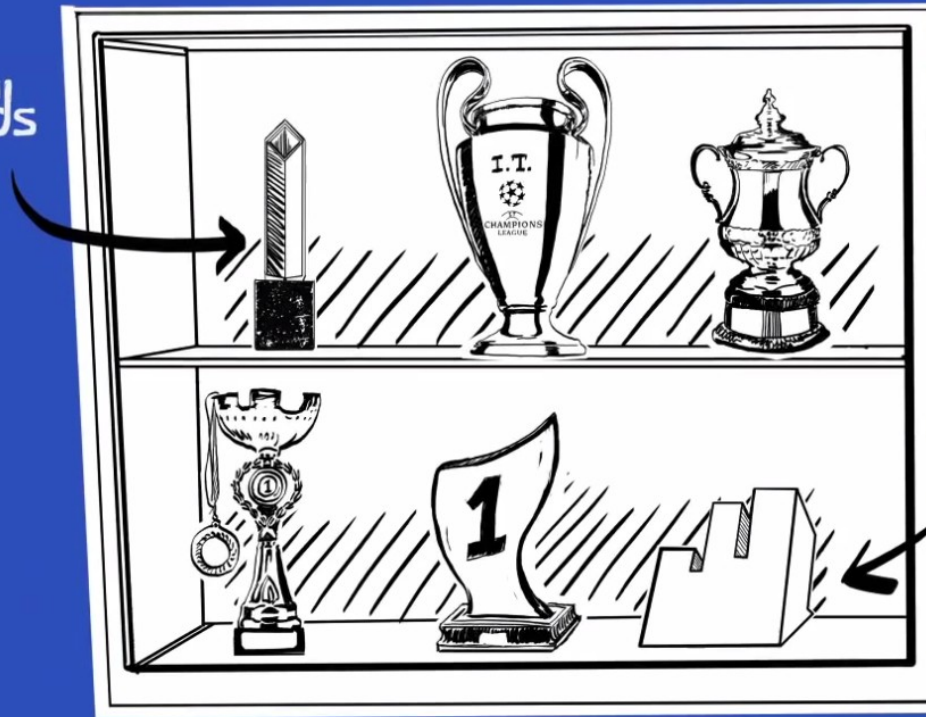
**Bol.com was elected as the
'Best webshop of the Netherlands'**

During the last year bol.com made great steps forwards by successfully applying customer profiling and BigData.



Join us

Thuiswinkel
awards



UITDAGEND
IT-TRAININGSPROGRAMMA



ON THE
JOB
COACHING

effie
awards

PERSONAL
TRAININGS
PROGRAMMA



40 SCRUMTEAMS



www.bol.com/banen

Q&A

You have

Questions

We have

Answers