# City Data Fusion

http://citydatafusion.org

DeRiVE 2013 Workshop
21.10.2013, ISWC 2013, Sydney, Australia

# Listening to the pulse of our cities during City Scale Events

Emanuele Della Valle

emanuele.dellavalle@polimi.it

http://emanueledellavalle.org

# Share, Remix, Reuse — Legally

- This work is licensed under the Creative Commons Attribution 3.0 Unported License.

- **Your are free:**

  **to Share** — to copy, distribute and transmit the work

  **to Remix** — to adapt the work

- **Under the following conditions**

  **Attribution** — You must attribute the work by inserting
  - "[source http://citydatafusion.org]" at the end of each reused slide
  - a credits slide stating
    - These slides are partially based on "Listening the pulse of our cities during City Scale Events" by E. Della Valle

- To view a copy of this license, visit http://creativecommons.org/licenses/by/3.0/
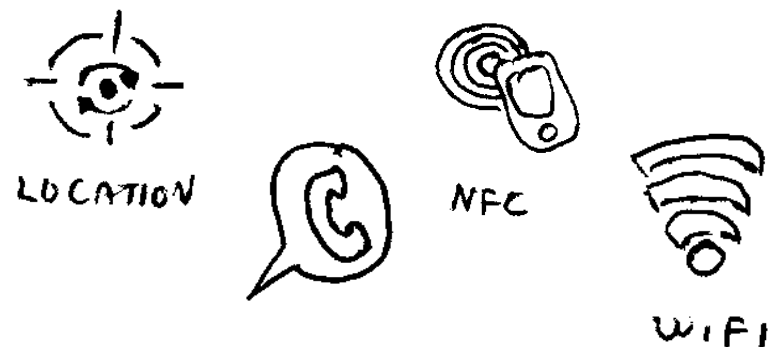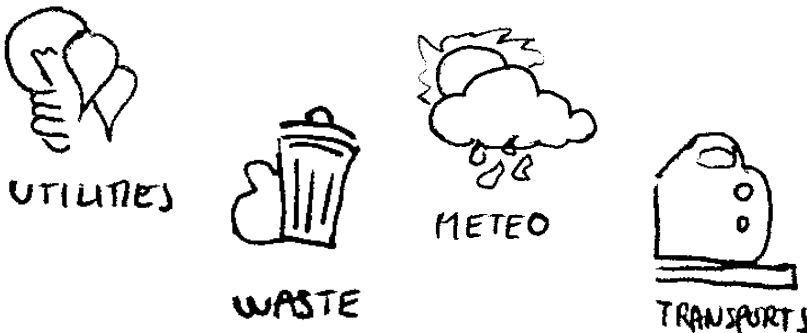
# Agenda

- Context

- Ingredients and challenges

- Research hypothesis

- Streaming Linked Data Framework
  - design principles
  - Architecture
  - Components

- Testing the research hypothesis
  - London Olympic Games 2013
  - Milano Design Week 2013

- Conclusions

Streams of information flows through our cities thanks to:

the pervasive deployment of sensors in our cities

UTILITIES

WASTE

METEO

TRANSPORTS

the wide adoption of smart phones (equipped with sensors)

LOCATION

NFC

WIFI

the usage of (location-based) social networks

the availability of datasets about urban environment

open Street Map

DBpedia

Open data

Free base

WIKI DATA

- **We can feel the pulse of our cities by**

**fusing** all those data sources

**making sense** of the fused information

# City Scale Events as test beds

- Characteristics
  - Lasting days
  - Hundreds of venues
  - Thousands of events
  - Hundreds of thousands of visitors

- Questions
  - Which are the most attractive events?
  - What do visitors think about the events they join?
  - What is their mood before, during and after the event they join?

- Ground truth
  - The program of the event
  - News about the event

# Example of City Scale Event
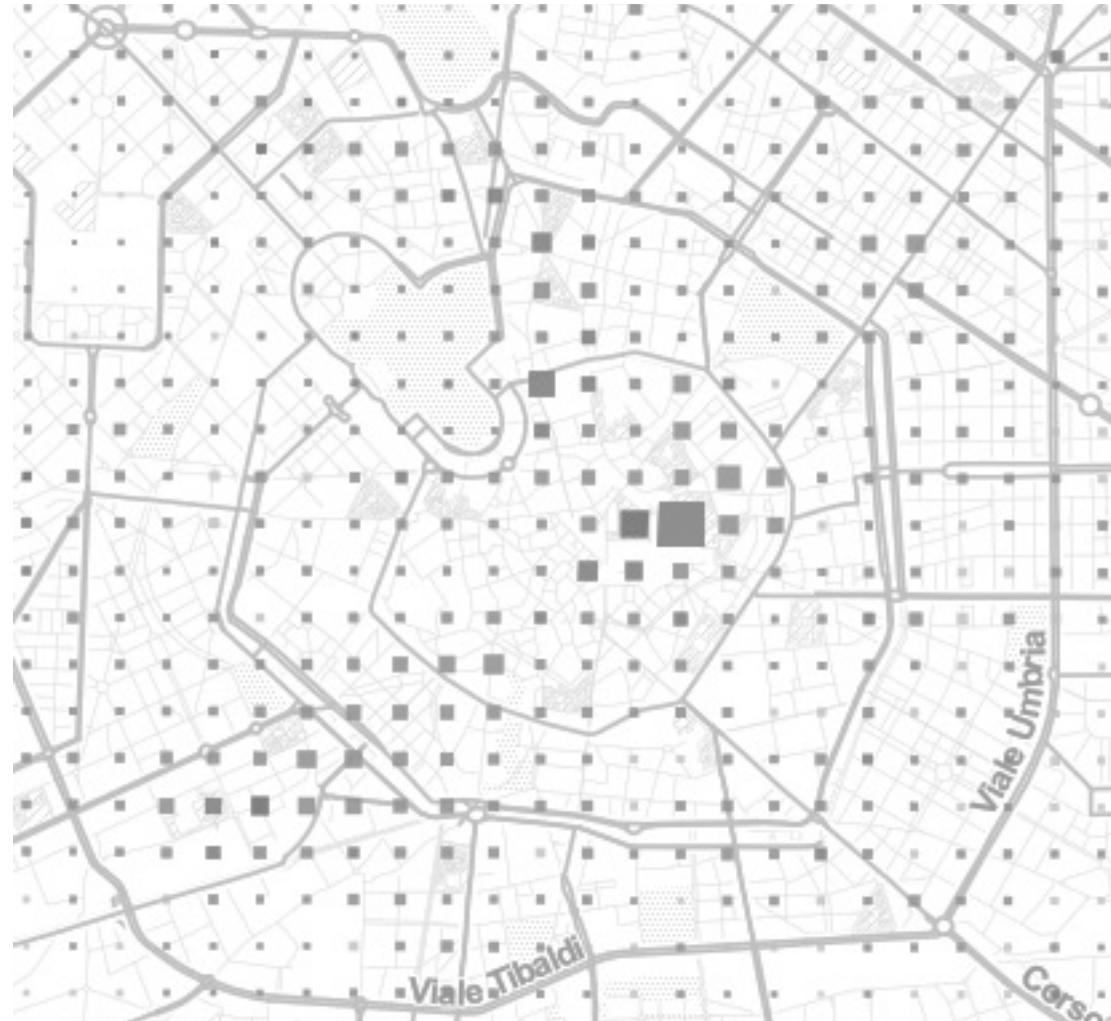## The Milano Design Week 2013 (MDW'13)

- Problem
  - Sponsor and organizer of a city scale event needs to quantify the return on investment

- Existing solutions
  - Spread people around the hundreds of event to asses the success of the various events is effective, but expensive

- Challenge
  - Obtaining comparable results by analysing public social streams

# E.g., is Milano Design Week perceivable?

Step 1: associate mobile traffic to urban areas



Real data recorded on 13 April 2013 between 13:00 and 00:00

Step 2: subtract what is systematic



Real data recorded on 13 April 2013 between 13:00 and 00:00

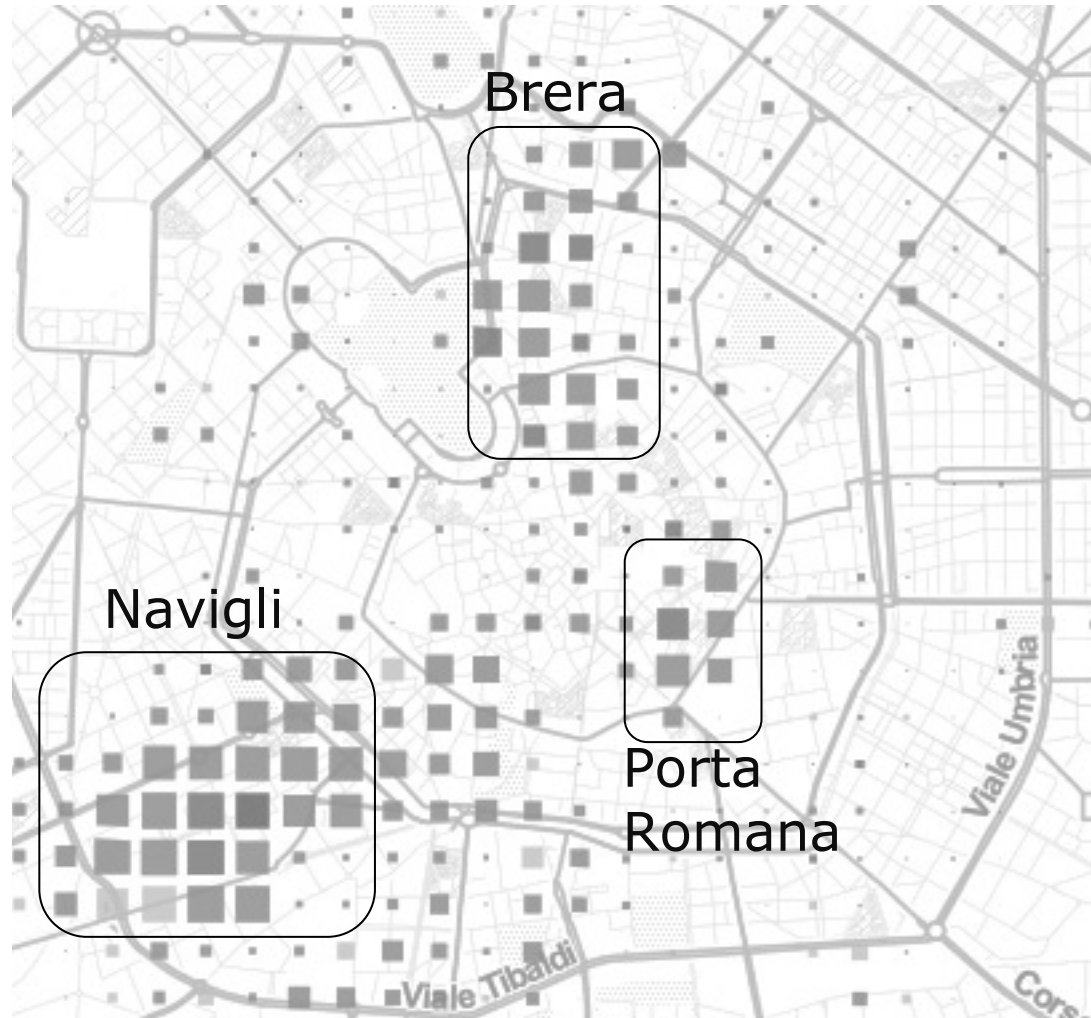## Step 3: Identify interesting areas



Brera

Navigli

Porta Romana

Real data recorded on 13 April 2013 between 13:00 and 00:00

# E.g., is Milano Design Week perceivable?

Step 4: retrieve the top hashtags



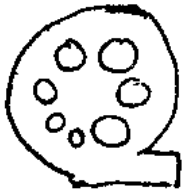Real data recorded on 13 April 2013 between 13:00 and 00:00

Step 5: exclude what is systematic



Brera

Navigli

Porta Romana

fuorisalone2013
fuorisalone
breradesigndistrict
casamadrefoodbar
zerofuorisalone

Real data recorded on 13 April 2013 between 13:00 and 00:00

# Ingredients to be combined

- **semantic technologies**
    - Address "*variety*" using Ontology Based Data Access
    - Named Entity recognition and linkage
    - Knowledge discovery (e.g., detecting systematicy)

- **streaming algorithms**
    - Address "*velocity*" of data stream
    - Address "*volume*" by being able to process data that do not fit in main memory

- **crowd-sourcing techniques**
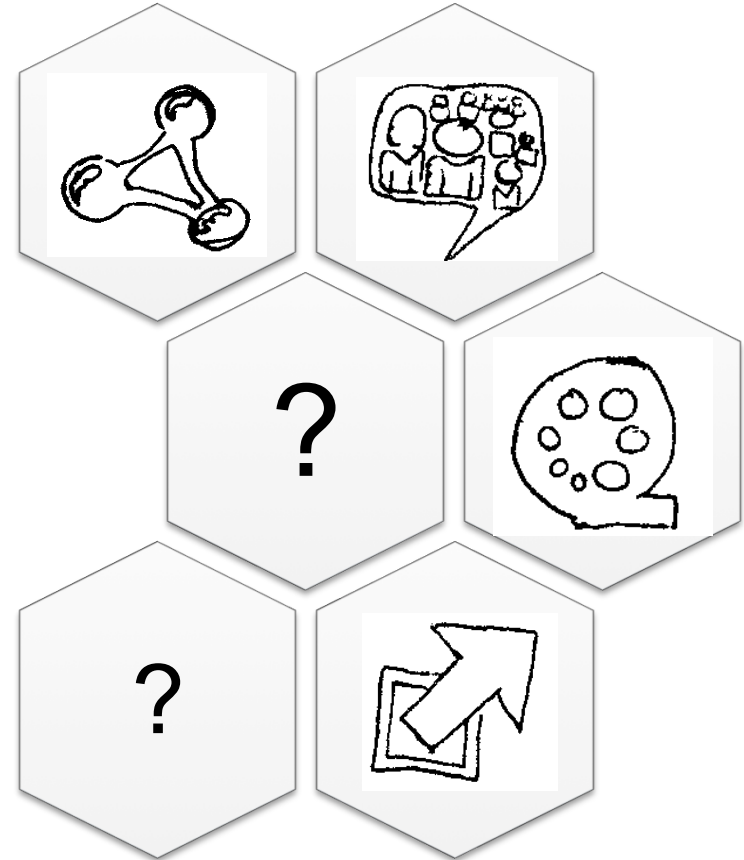    - Address "*veracity*" by cleansing and enriching data

- **Visual analytics**
    - Allow no-expert access to data
    - Tell stories out of data

# Limitation of current systems

- Insufficient methods for making sense in **real-time** of **heterogeneous** data and social streams w.r.t. the vast collections of (open) data

- Lack of crowd-sourcing techniques whose **incentives** leverage needs of people in the urban environment

- Lack of visualization techniques tailored to **non-experts**

?

?

# Research hypothesis

1. To scale order matters

2. Crowdsourcing needs the urban-centric incentives

3. Visualization must tell stories

# Research hypothesis: order matters!

- Observation: order reflects recency, relevance, trustability …

<table>
<tr><td rowspan="5"><strong>Types of orders</strong></td><td>Combinations</td><td>Continuous top-k Q/A</td><td>Order-aware reasoning</td></tr>
<tr><td>Relevance, Trustability, etc.</td><td>Top-k Q/A</td><td>Top-k Reasoning</td></tr>
<tr><td>Recency</td><td>DSMS/CEP</td><td>Stream reasoning</td></tr>
<tr><td>Indexes</td><td>Traditional solutions</td><td>Scalable reasoning</td></tr>
<tr><td></td><td>No</td><td>Yes</td></tr>
</table>

**Types of reasoning**

- harnessing orders is key to make sense in real-time of heterogeneous, massive and volatile data

# Research hypothesis: urban-centric incentives!

[source: http://www.behance.net/gallery/Maslows-Hierarchy-of-Needs-Infographic/4376921 ]

Maslow's hierarchy of needs

- **incentives** designed **for urban environment and life styles** are key

[Source: http://www.densitydesign.org/2013/04/whatever-the-weather/ ]

# Testing the research hypothesis

1. To scale order matters
   - Stream Reasoning
     – RDF Stream
     – Continuous SPARQL
     – Incremental Materialization for RDF Streams (IMaRS)
     – C-SPARQL Engine
     – RESTFul Services for RDF Stream Processors
     – Streaming Linked Data Framework
   - SPARQL Rank
     – Rank aware SPARQL algebra
     – ARQ-Rank

2. Crowdsourcing needs the urban-centric incentives
   - Urban Games With A Purpose
     – UrbanMatch
     – Urbanopoly

3. Visualization must tell stories
   - On going work

# Focus of this key note

1. To scale order matters
   - Stream Reasoning
     - RDF Stream
     - Continuous SPARQL
     - Incremental Materialization for RDF Streams (IMaRS)
     - C-SPARQL Engine
     - RESTFul Services for RDF Stream Processors
     - Streaming Linked Data Framework
   - SPARQL Rank
     - Rank aware SPARQL algebra
     - ARQ-Rank

2. Crowdsourcing needs the urban-centric incentives
   - Urban Games With A Purpose
     - UrbanMatch
     - Urbanopoly

3. Visualization must tell stories
   - On going work

# Streaming Linked Data Framework

- Input Data Formats
  - Streaming information: RDF streams
  - Background Information: RDF graphs

- Query Language
  - Continuous SPARQL

- Features
  - Adapters to access the social streams, e.g., twitter
  - Ability to record and replay portions of the social stream
  - Possibility to decorate the social stream with sentiment information
  - Possibility to continuously analyzing the social stream
  - Possibility to built complex application composing complex networks of decorators and analyzers
  - Possibility to publishing and visualizing results of continuous analysis
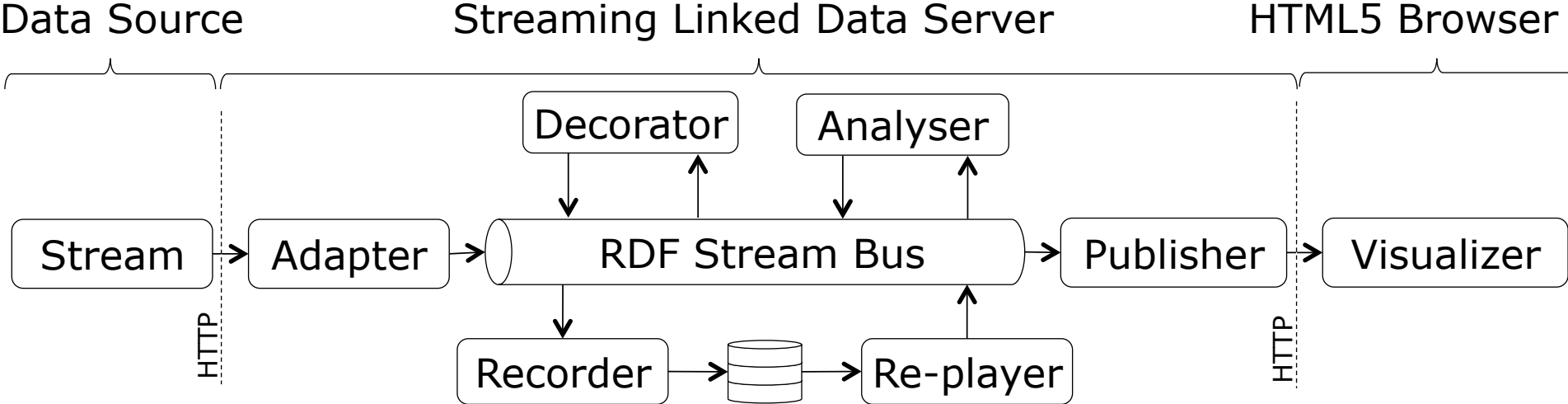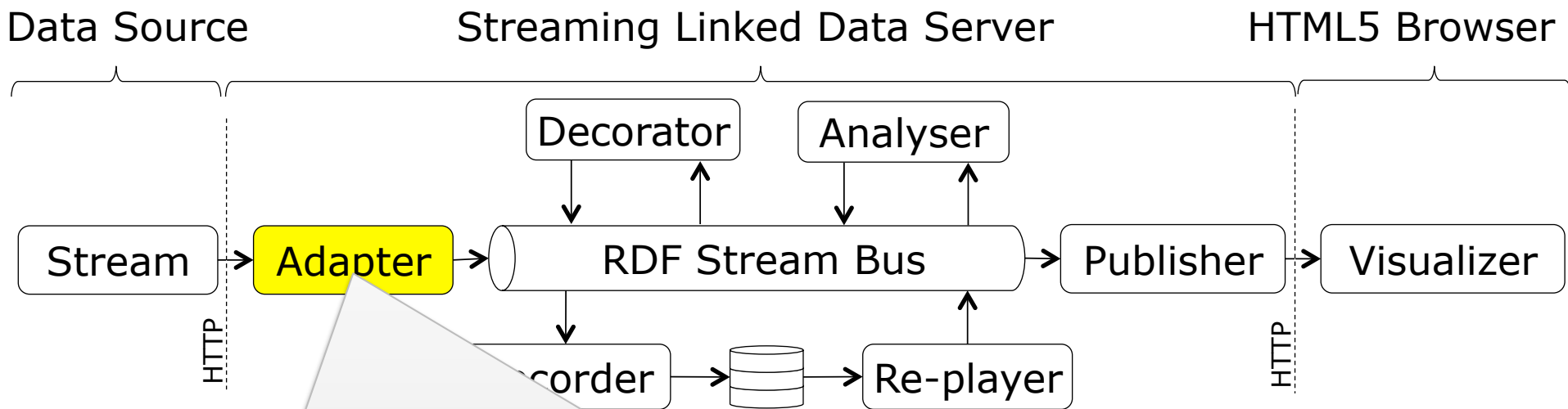
# Design Principles

- Follow publish/subscribe patter
  - Publisher and subscribers do not have to know each other
  - Subscribers can join and leave in any moment

- Adopt a reliable message-passing
  - Guarantees delivery order

- Minimise latency by using main memory
  - Avoiding disk I/O bottleneck

# Architecture

Data Source      Streaming Linked Data Server      HTML5 Browser



```
                    ┌───────────┐   ┌───────────┐
                    │ Decorator │   │ Analyser  │
                    └───────────┘   └───────────┘
                         │  ↑           │  ↑
┌────────┐   ┌─────────┐ ↓  │           ↓  │ ┌───────────┐   ┌────────────┐
│ Stream │→→ │ Adapter │→→( RDF Stream Bus )→→│ Publisher │→→ │ Visualizer │
└────────┘   └─────────┘     │          ↑    └───────────┘   └────────────┘
   HTTP                      ↓          │          HTTP
                    ┌──────────┐    ┌───────────┐
                    │ Recorder │→→[DB]→→│ Re-player │
                    └──────────┘    └───────────┘
```

# Architecture - Adapters

Data Source          Streaming Linked Data Server          HTML5 Browser

Decorator          Analyser

Stream → **Adapter** → ◯ RDF Stream Bus → Publisher → Visualizer

HTTP

...corder → ▭ → Re-player

HTTP
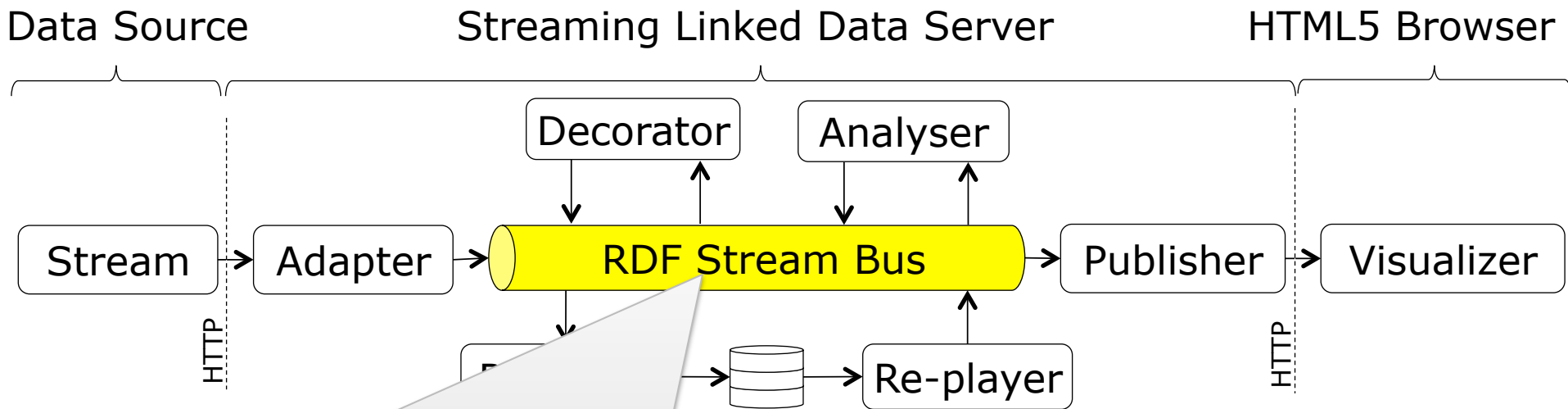
An adapter streams time-stamped RDF graphs

[] sioc:content "This is for everyone #london2012 #oneweb #openingceremony";
   sioc:has_creator :timberners_lee;
   sioc:topic :london2012, :oneweb, :openingceremony .

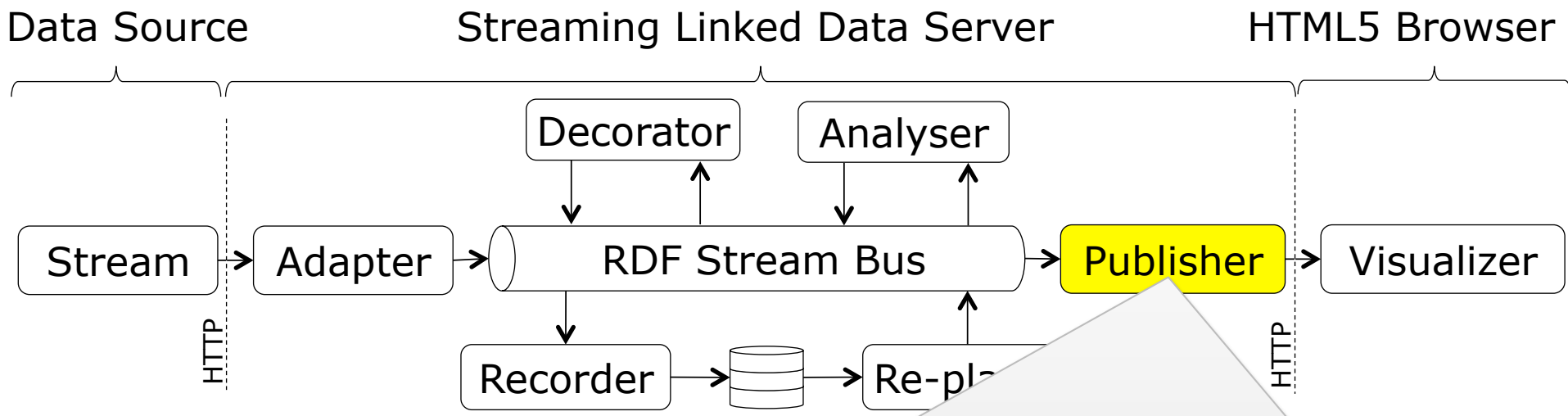Available adapters include: twitter, instagram, pachube, and linked sensor data

Data Source     Streaming Linked Data Server     HTML5 Browser

Decorator     Analyser

Stream → Adapter → RDF Stream Bus → Publisher → Visualizer

HTTP

Re-player

HTTP

The RDF Stream Bus supports the publish/subscribe communication

Data Source   Streaming Linked Data Server   HTML5 Browser

Decorator | Analyser

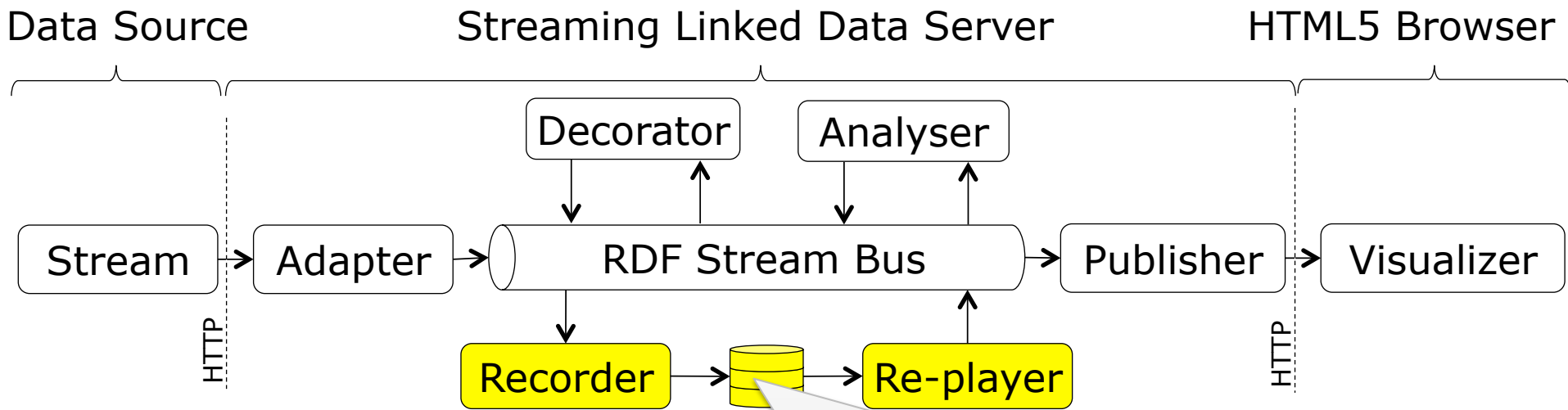Stream → Adapter → ( RDF Stream Bus → **Publisher** → Visualizer

HTTP

Recorder → ⬚ → Re-pl...

HTTP

A publisher make available on the Web the last RDF graphs of an
RDF stream following a variation of Streaming Linked Data format*:

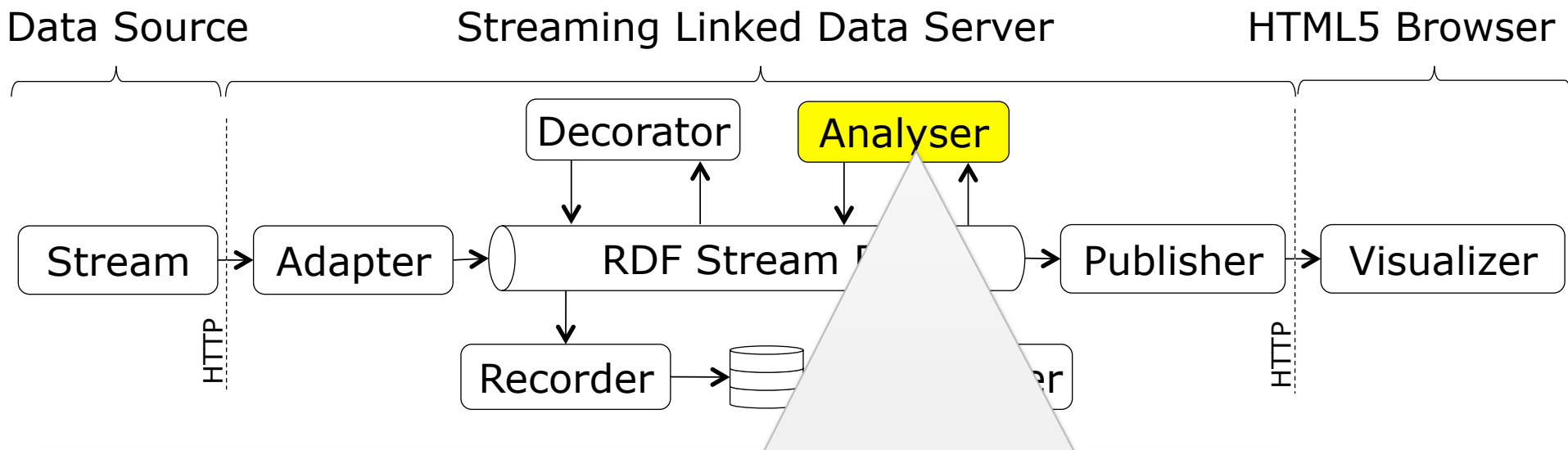<#sGraph> —sld:contains→ ○ —rdf:rest→ ○ —rdf:rest→ rdf:nil

rdf:first ↓        rdf:first ↓

<#iGraph1>   <#iGraph2>

* Barbieri, D.F., Della Valle, E.: A proposal for publishing data streams as linked data - a position
paper. In: LDOW. (2010)

# Architecture – Recorders and Re-players

Data Source          Streaming Linked Data Server          HTML5 Browser



A recorder records the content of an RDF stream following a variation of Streaming Linked Data format.
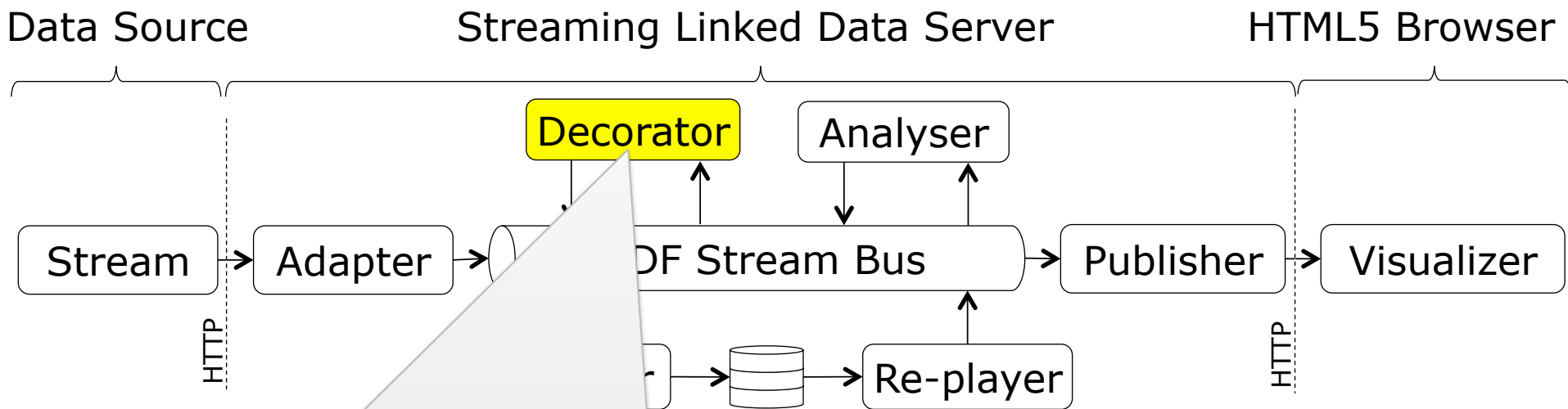A re-player re-plays a recorded stream. It can vary the speed.

Data Source        Streaming Linked Data Server        HTML5 Browser

Decorator    Analyser

Stream → Adapter → ( RDF Stream ) → Publisher → Visualizer

HTTP

Recorder →

HTTP

An analysers continuously execute C-SPARQL queries.
E.g., *count the number of times each hashtag is used in the last 15 minutes updating the counting every minute*.

```
REGISTER STREAM HashtagAnalysis AS
CONSTRUCT { [] sld:about ?tag ; sld:count ?n . }
FROM STREAM <http://.../London2012> [RANGE 15m STEP 1m]
WHERE { { SELECT ?tag (COUNT(?tweet) AS ?n)
          WHERE { ?tweet sioc:topic ?tag . } GROUP BY ?tag } }
```

Data Source          Streaming Linked Data Server          HTML5 Browser

**Decorator**   **Analyser**

**Stream** → **Adapter** →   DF Stream Bus   → **Publisher** → **Visualizer**
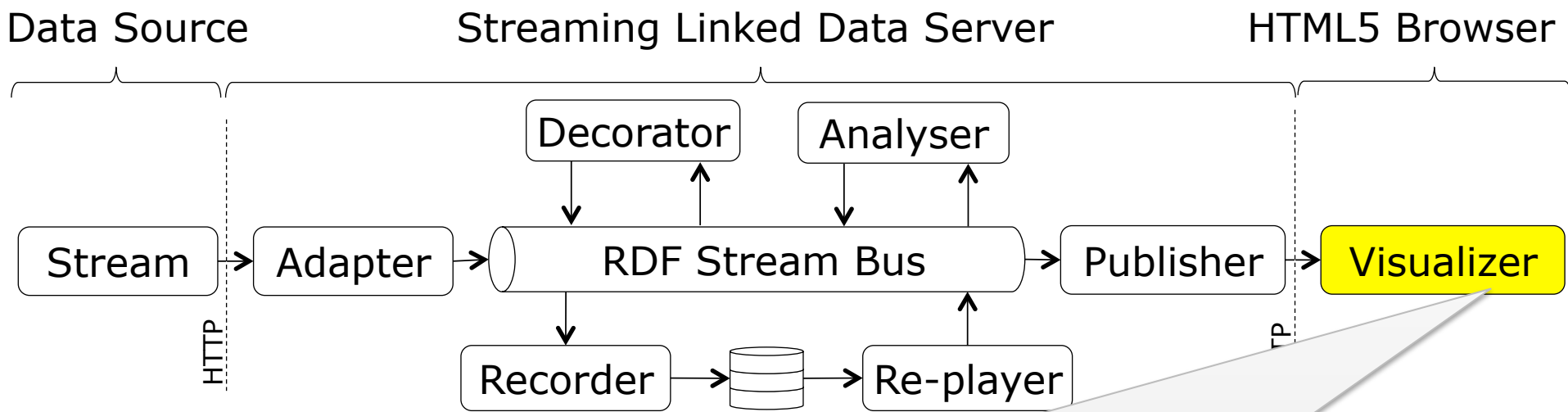
HTTP

→ **Re-player**

HTTP

A decorator adds information to streamed RDF graphs that match
a given patter.
E.g., Dictionary-based sentiment classifiers* (known to be efficient for
short texts concentrating on a single topic, such as tweets) was used in
this work to decorate each tweet.

* Tsytsarau, M., Palpanas, T., Denecke, K.: Scalable Detection of Sentiment-Based
  Contradictions. In: DiversiWeb workshop, WWW, Hyberabad, India (2011)

Data Source          Streaming Linked Data Server          HTML5 Browser

Stream → Adapter → ( RDF Stream Bus ) → Publisher → **Visualizer**

Decorator    Analyser

Recorder → Re-player

HTTP          HTTP

A visualizer displays the published linked data.
Available visualizers include:

| Heatmaps | Bar charts | Area charts | Dot charts |
| --- | --- | --- | --- |

# The London Olympic Games 2012 (LOG'12)

- Problem
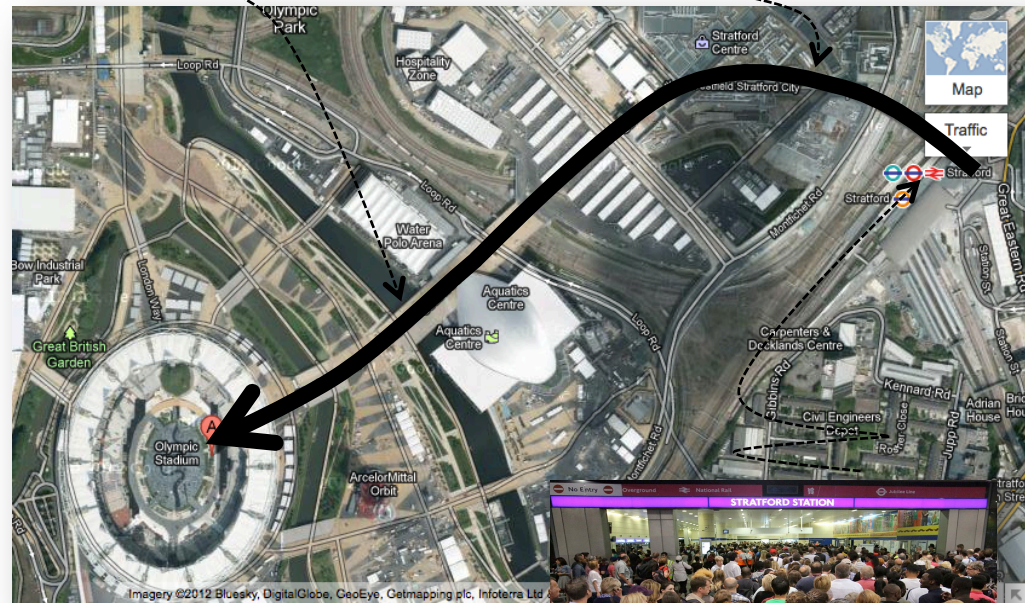  - To manage a big event requires tracking in real time the movement of crowds.

- Existing solutions
  - CCTV, and mobile network data analysis are effective, but expensive

- Challenge
  - Obtaining comparable results by analysing public social streams

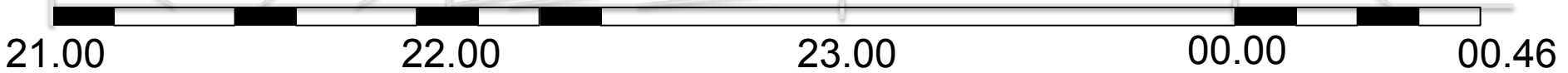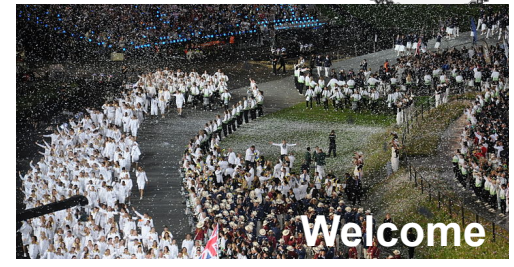# Case study #1 tracking attention of the crowds

- ## The problem
  - Managers of Big event want to track in real time if the event is capturing the attention of the audience.

- ## Input
  - 20 thousand of tweets streamed by Twitter between 9pm July 27th and 1am July 28th

- ## Ground truth
  - The Opening Ceremony Broadcast
    - http://www.youtube.com/watch?v=4As0e4de-rI
  - The wikipedia page describing the Opening Ceremony
    - http://en.wikipedia.org/wiki/ 2012_Summer_Olympics_opening_ceremony

Happy and Glorious

Interlude

Welcome

Pandemonium

Straight on till morning

... thanks Tim

There Is a Light

21.00          22.00          23.00          00.00          00.46

- Interesting phenomena are visible at different scales



World          Continent          City
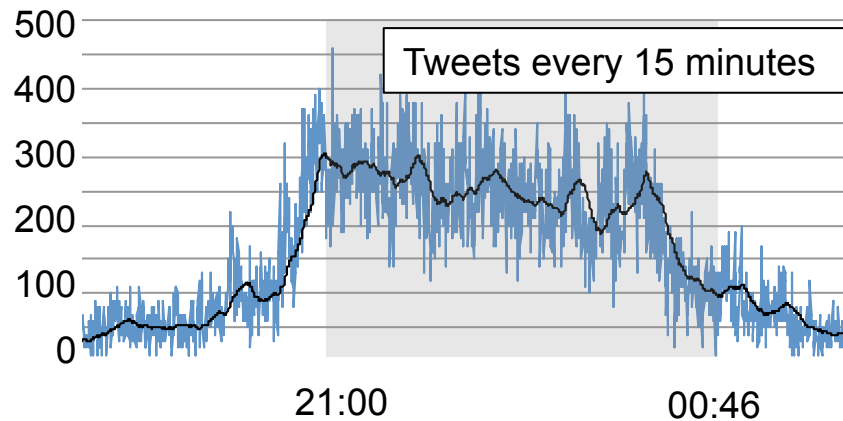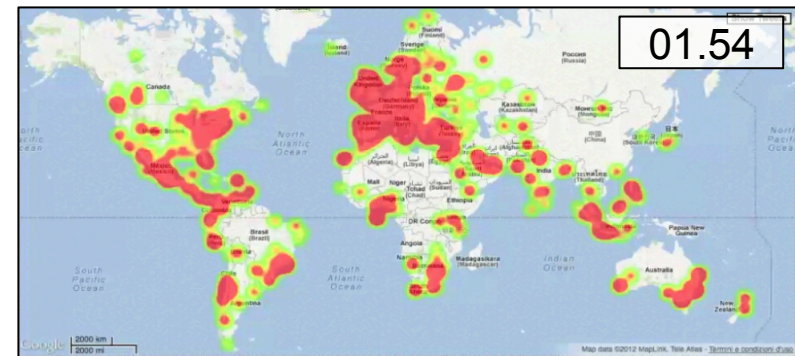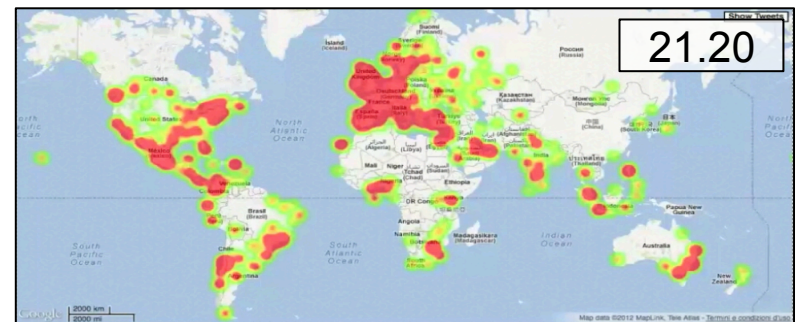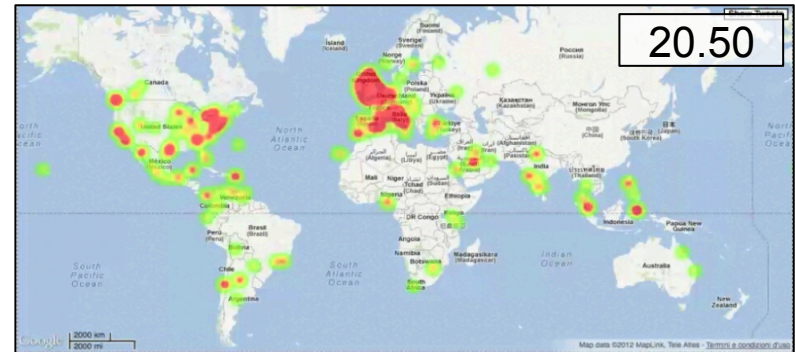
- The Opening Ceremony is clearly visible in the volume of tweets containing LOG'12 related hashtags
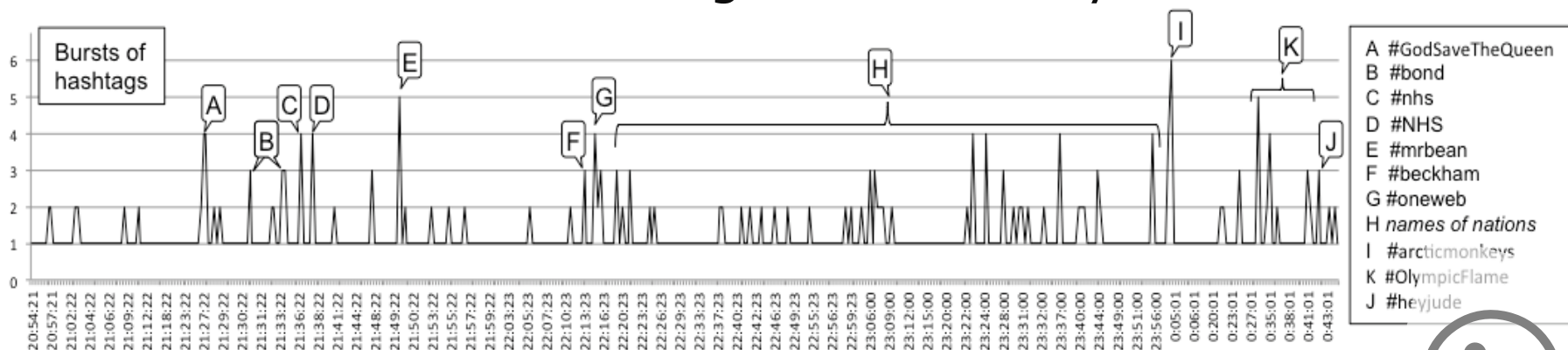
- Global scope



Tweets every 15 minutes

21:00    00:46

# Findings

- Bursts of hashtags usage capture what attracts the attention of those watching the ceremony world wide



- Detailed Analysis

| Moments of Ceremony | # of hashtags | Fraction |
|---|---|---|
| Total | 189 | 100.00% |
| Hashtagged with an emotional state | 34 | 17.99% |
| Correctly hashtagged | 72 | 38.10% |
| - Right on time (1 min tolerance) | 50 | 26.46% |
| - After the event (15 min tolerance) | 13 | 6.88% |
| - Before the event (15 min tolerance) | 9 | 4.76% |

# Findings

- Audience loosed attention while the ceremony was progressing

- Audience emotions where more evident in the first and the last part of the ceremony

| | | |
|---|---|---|
| 1. Countdown | (21:00-21:04) |
| 2. Green and Pleasant Land | (21:04–21:09) |
| 3. Pandemonium | (21:09–21:25) |
| 4. Happy and Glorious | (21:25–21:35) |
| 5. Second to the right | (21:35–21:47) |
| 6. Interlude | (21:47–21:52) |
| 7. Frankie and June say… | (21:52–22:09) |
| 8. Abide with Me | (22:09–22:20) |
| 9. Welcome | (22:20–00:00) |
| 10. Bike a.m. | (00:00–00:07) |
| 11. Let the Games Begin | (00:07–00:24) |
| 12. There Is a Light | (00:24–00:38) |
| 13. And in the end | (00:38–00:46) |

Tweets per minute

unreleated
related
emotion

Fraction of hashtags
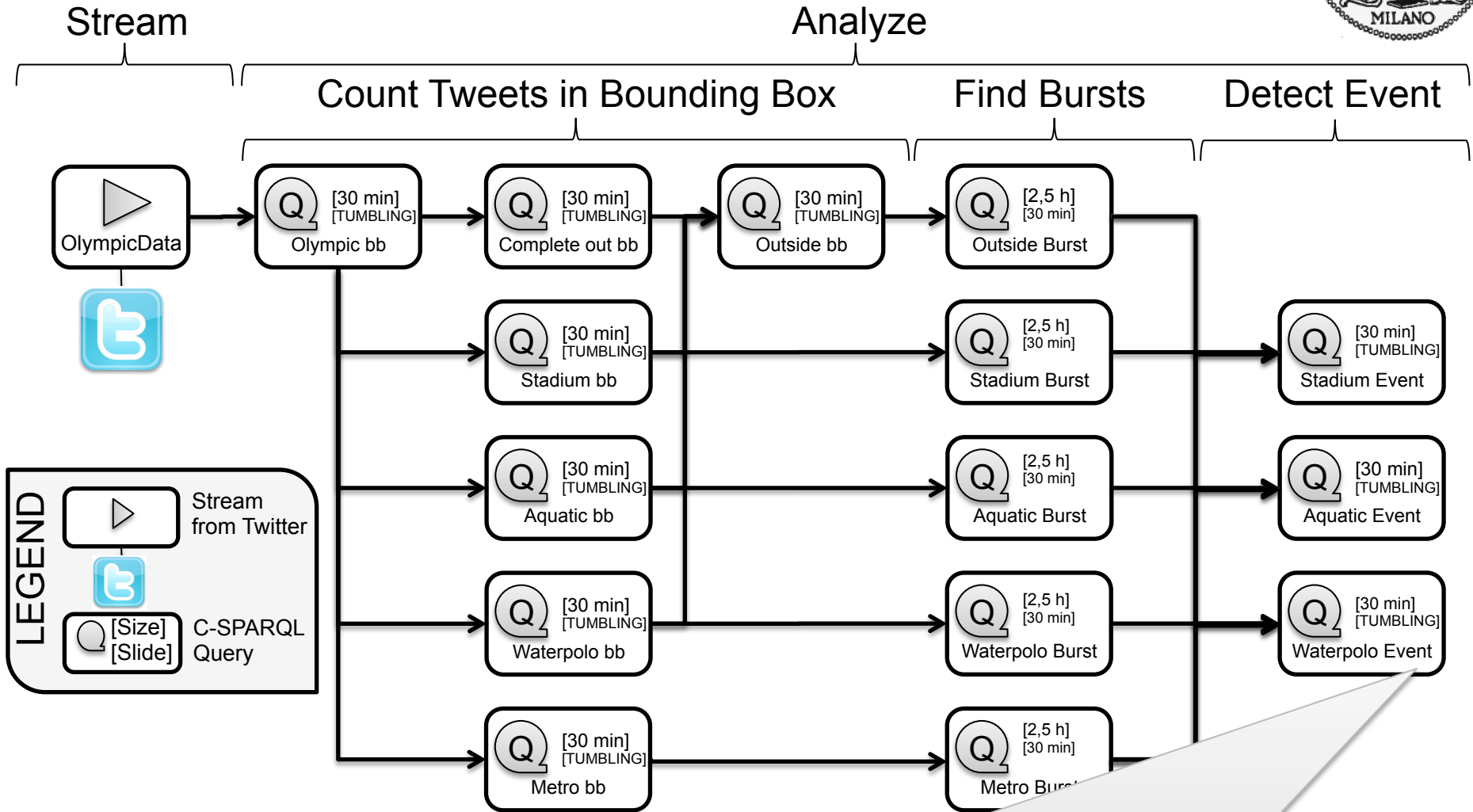
unreleated
related
emotion

- ## The problem
  - detect the events given the position of a set of venues and socially listening their surroundings

- ## Input
  - three million tweets streamed by Twitter between July 25[th] and August 13[th] 2012

- ## Conditions

| Type | Venue | Events | Capacity (seats) |
|---|---|---|---|
| Large | Olympic stadium | athletic games | 80,000 |
| Medium | Aquatic centre | swimming, diving and synchronized swimming | 17,500 |
| Small | Water polo arena | Water polo | 5,000 |

- ## Ground truth
  - calendar of Olympic Games

# SLD application

## Stream

## Analyze

### Count Tweets in Bounding Box

### Find Bursts

### Detect Event

OlympicData

Q Olympic bb — [30 min] [TUMBLING]

Q Complete out bb — [30 min] [TUMBLING]

Q Outside bb — [30 min] [TUMBLING]

Q Outside Burst — [2,5 h] [30 min]

Q Stadium bb — [30 min] [TUMBLING]

Q Stadium Burst — [2,5 h] [30 min]

Q Stadium Event — [30 min] [TUMBLING]

Q Aquatic bb — [30 min] [TUMBLING]

Q Aquatic Burst — [2,5 h] [30 min]

Q Aquatic Event — [30 min] [TUMBLING]

Q Waterpolo bb — [30 min] [TUMBLING]

Q Waterpolo Burst — [2,5 h] [30 min]

Q Waterpolo Event — [30 min] [TUMBLING]

Q Metro bb — [30 min] [TUMBLING]

Q Metro Burst — [2,5 h] [30 min]

### LEGEND

▷ Stream from Twitter
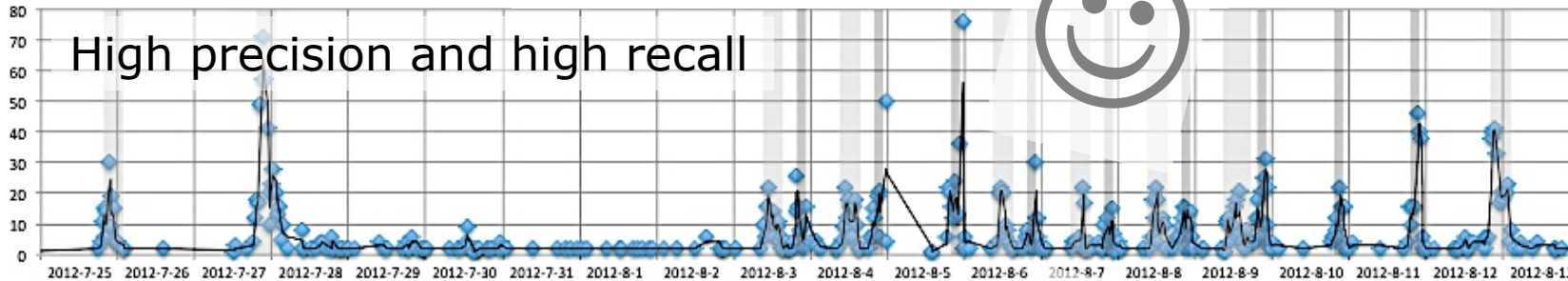
Q [Size] [Slide] C-SPARQL Query

(Metro Burst -> Outside Burst -> In ?venue Burst) within 30 min => event in ?venue

# Evaluation



Stadium
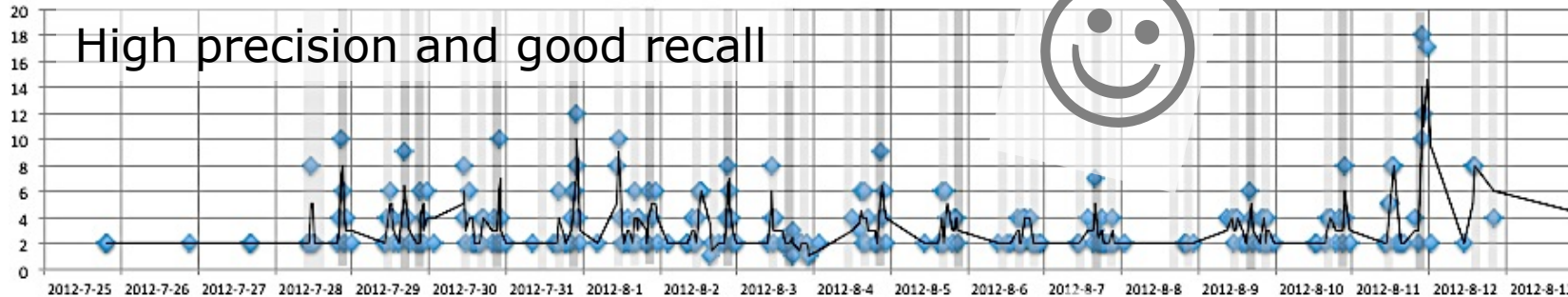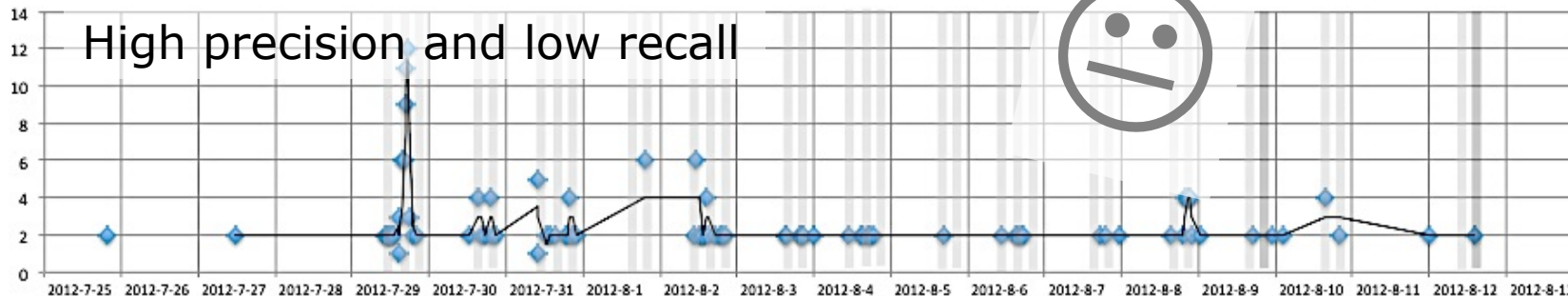
High precision and high recall

Water Centre

High precision and good recall

Water-polo Arena

High precision and low recall

Capacity (seats)

# Case study #3 visualizing crowds movements

- ## The problem
  - Visualize the movement of the crowds

- ## Input
  - three million tweets streamed by Twitter between July 25th and August 13th 2012

- ## Conditions

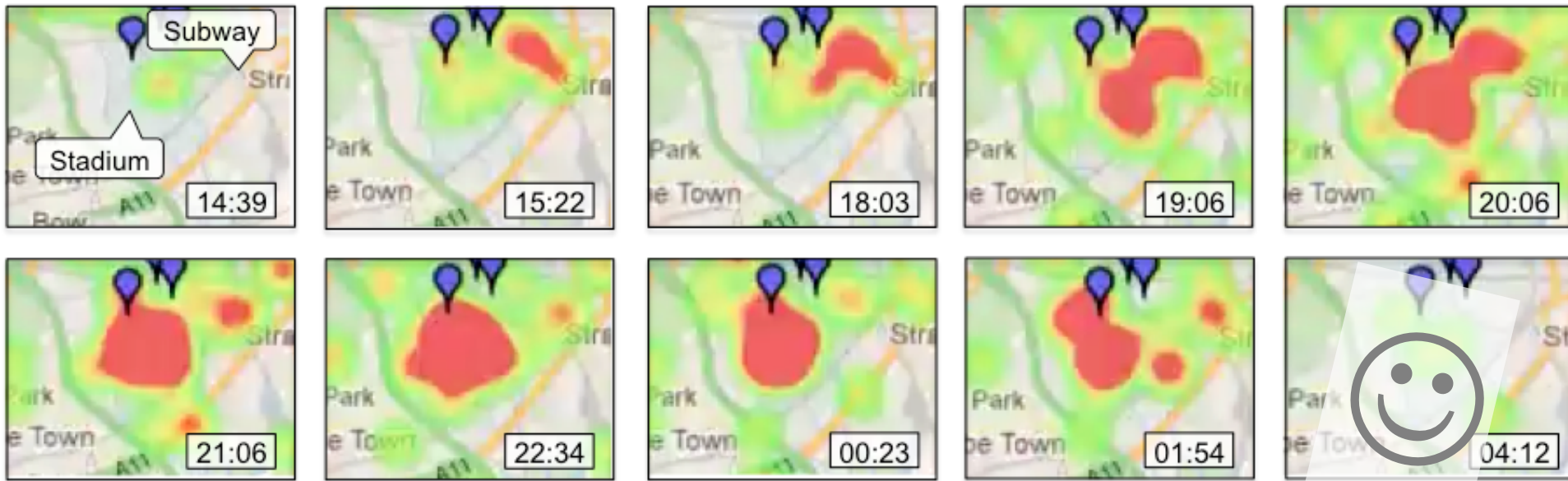| Type | Venue | Events | Capacity (seats) |
|---|---|---|---|
| Large | Olympic stadium | athletic games | 80,000 |
| Medium | Aquatic centre | swimming, diving and synchronized swimming | 17,500 |
| Small | Water polo arena | Water polo | 5,000 |

- ## Ground truth
  - Expert judgement ;-)

## Opening Ceremony at the Olympic Stadium



## A crowded event at the Aquatic Centre (July 31st, 2012)

- Problem
  - Sponsor and organizer of a city scale event needs to quantify the return on investment

- Existing solutions
  - Spread people around the hundreds of event to asses the success of the various events is effective, but expensive

- Challenge
  - Obtaining comparable results by analysing public social streams
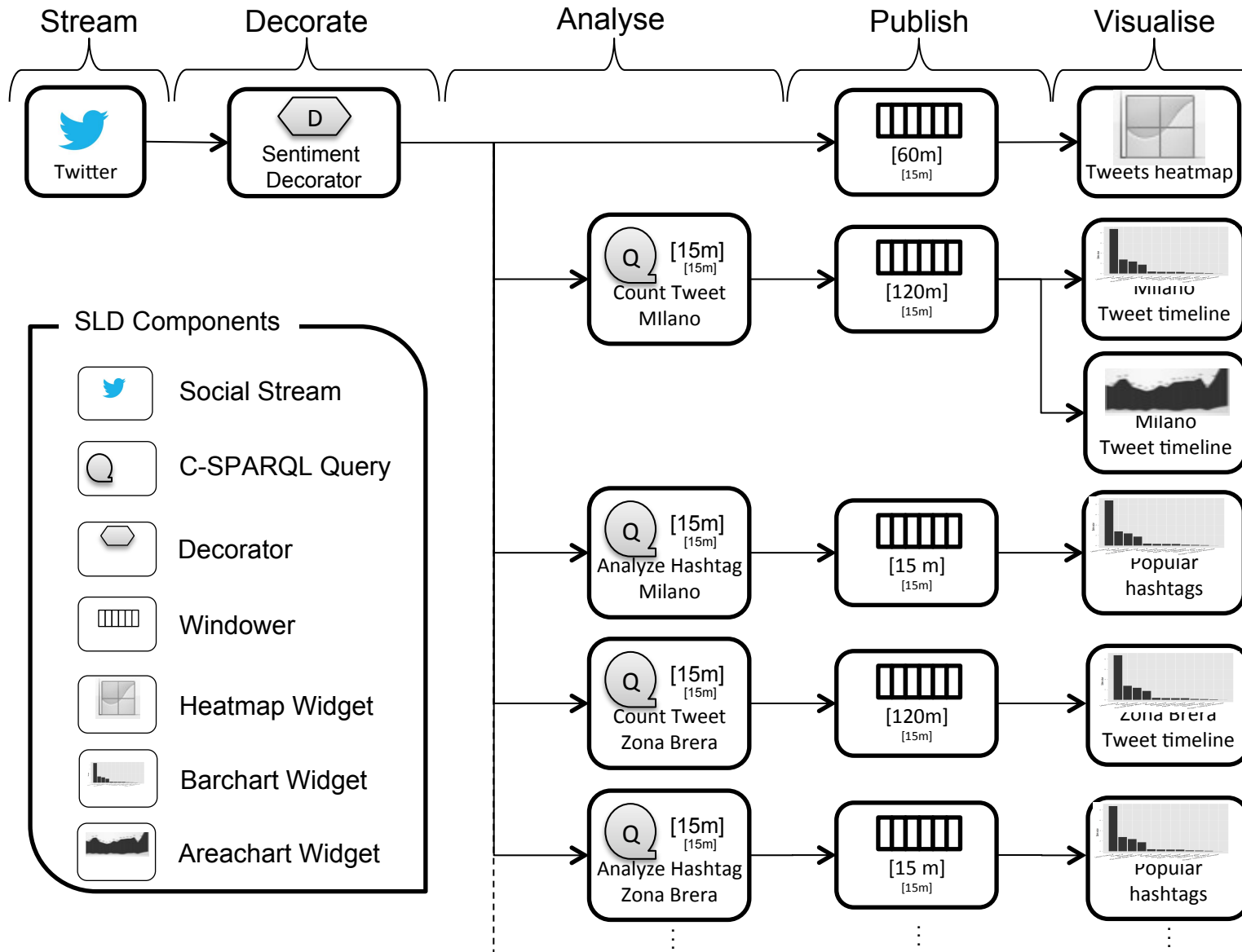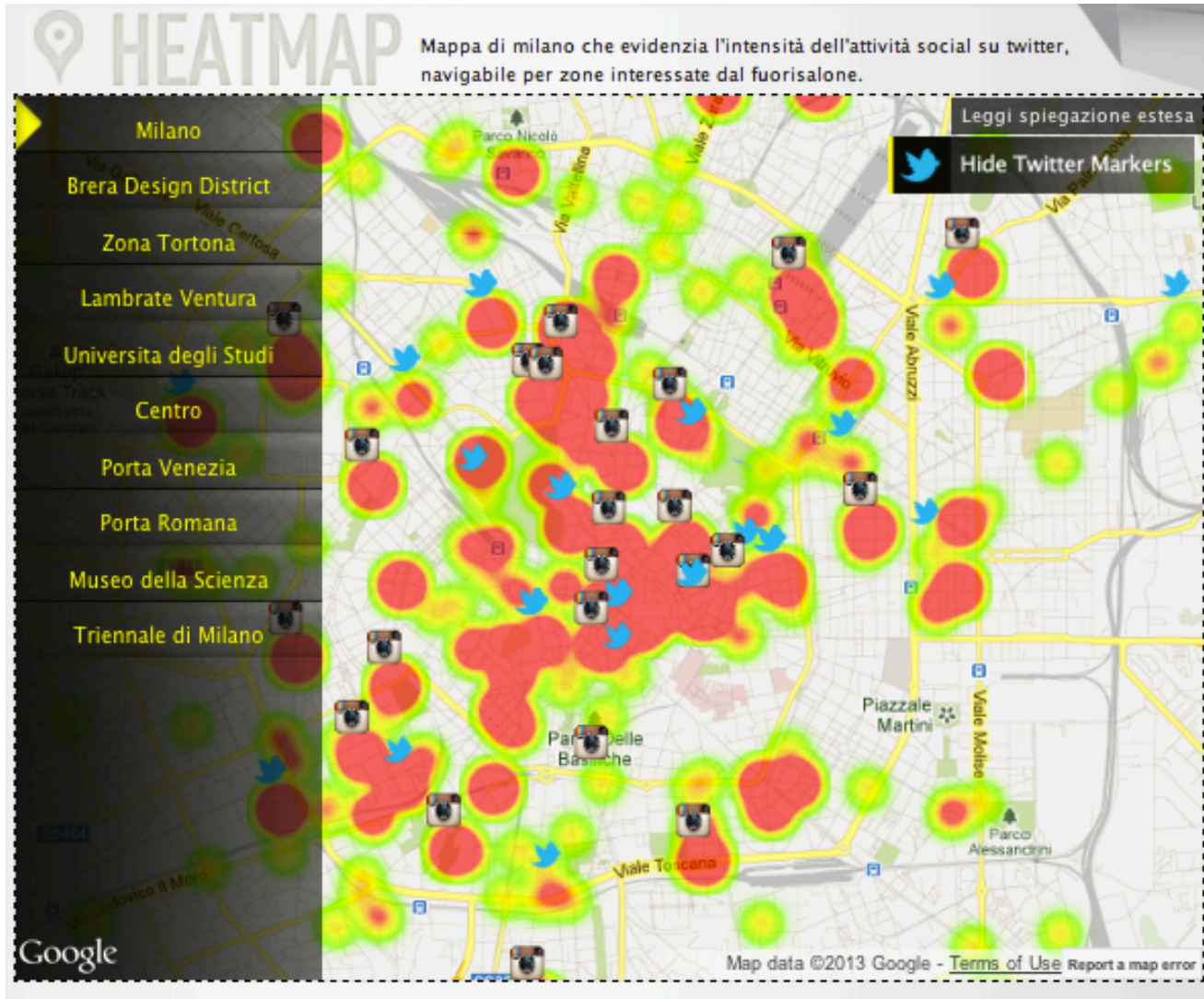
- The problem
  - Is MDW visible in the social streams posted by people in Milano area? If yes in real-time,
    1. What are the districts from which MDW visitors post the most?
    2. What are the most frequently used hashtags?
    3. How do people feel before, during and after the event they join?
  - Can these question answered at a cost a SME can afford?

- Input
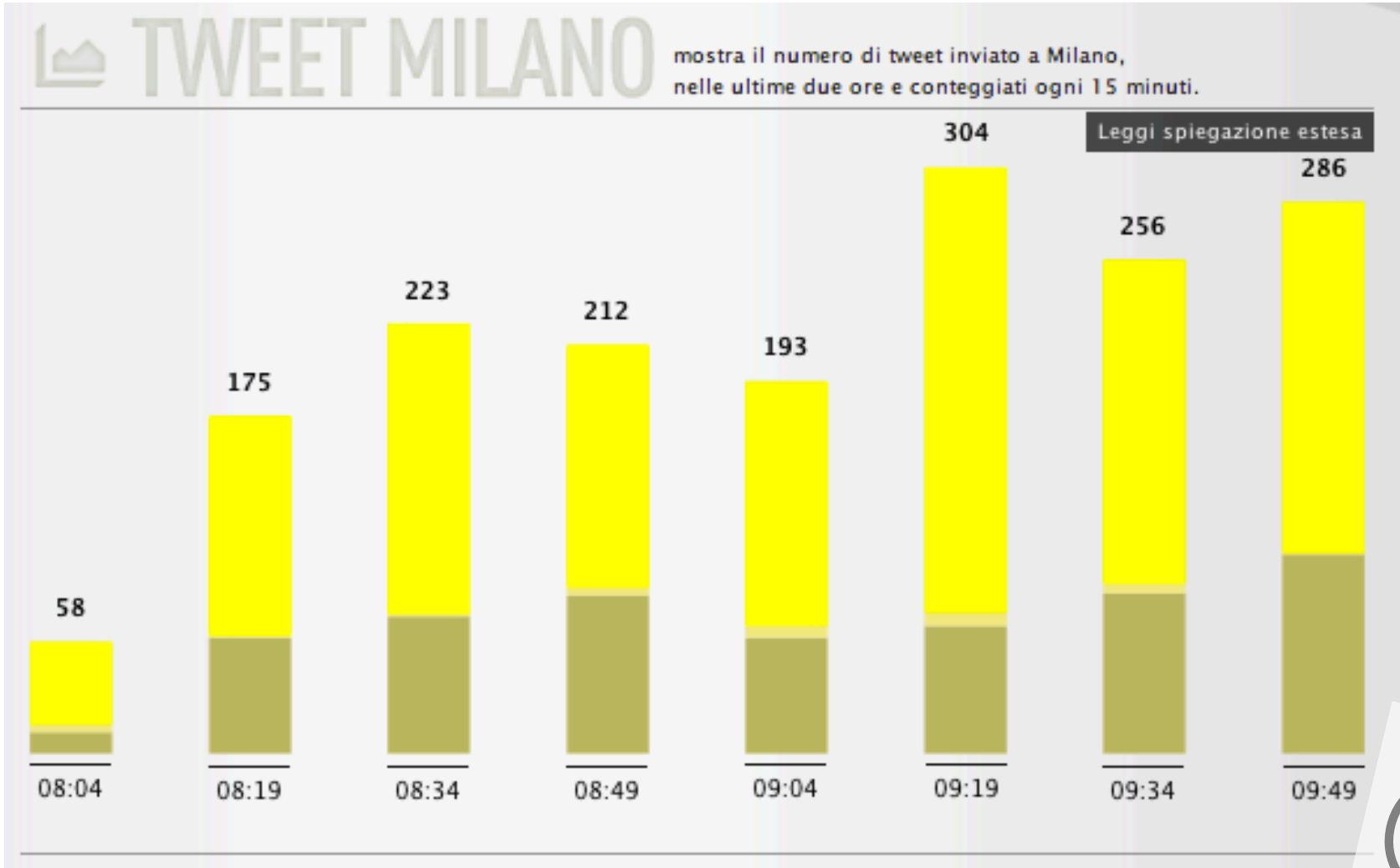  - 106,770 tweets streamed by Twitter between April 9th and April 14th 2013

# SLD application

April 9-14, 2013

- Distinct users
  - 12.031
- Invocation of Linked Data Publisher
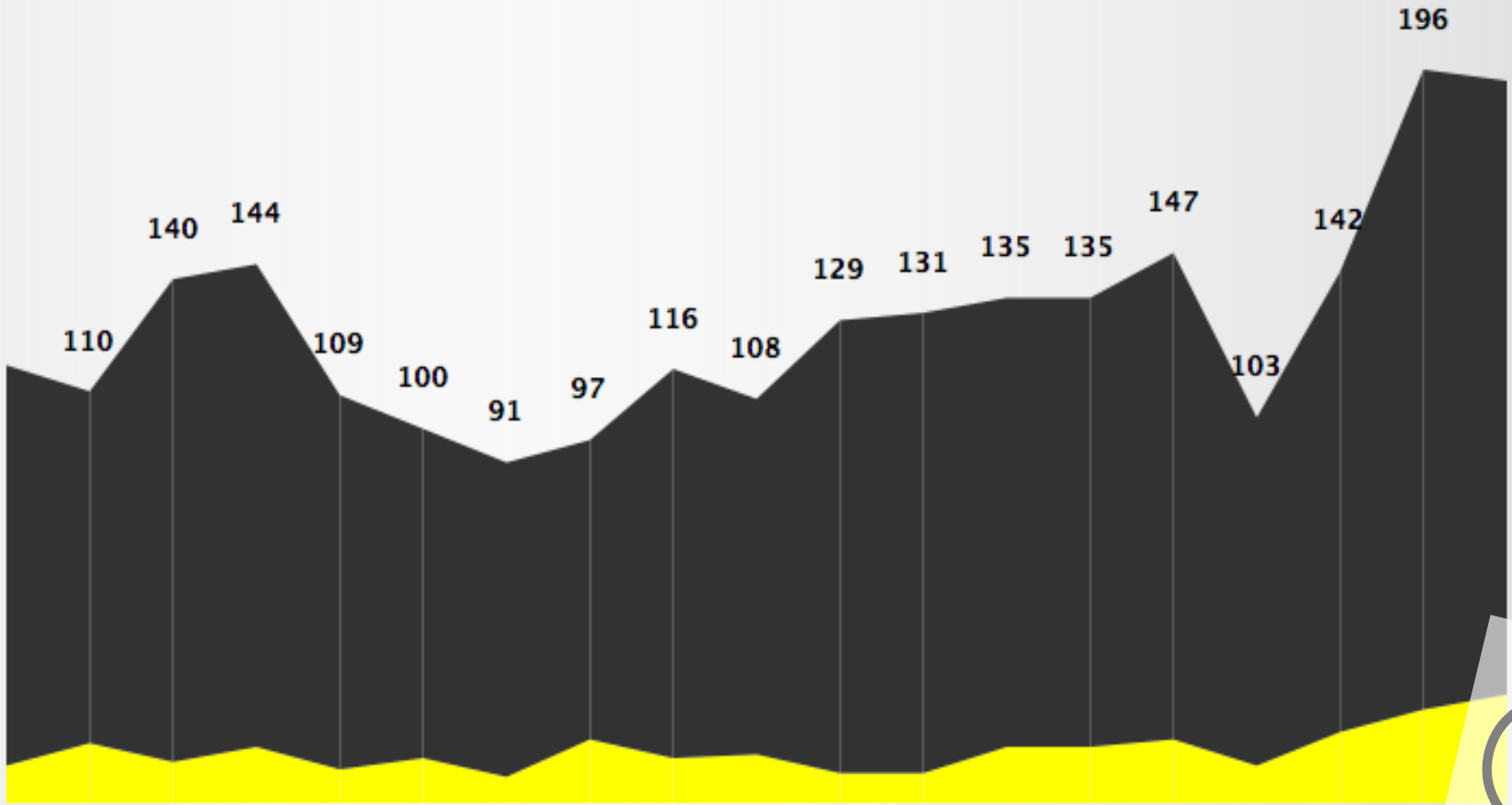  - 1,136,052
- Cost
  - 25 €/month
  - 2 cores, 2 GB

TWEET COMPARE

Confronta il numero di tweet di Milano (grigio) con i tweet che contengono # collegati all'evento fuorisalone (giallo)

# TOP HASHTAG

i 10 hashtag maggiormente usati nei tweet a Milano, conteggiati ogni 15 minuti.

| Hashtag | Count |
|---|---|
| triennale | 15 |
| fuorisalone | 11 |
| inaugurazione | 11 |
| milano | 11 |
| design | 10 |
| designweek | 9 |
| milan | 9 |
| TDM6 | 5 |
| sindromeinfluenza | 3 |
| museum | 2 |

DOT CHART

Griglia che confronta le zone di milano piu attive, conteggiando i tweet ogni 15 minuti.

MDW 2013 is visible in the volume of micro-posts

| April 9th, 2013 at 18.00 | posts |
|---|---|
| fuorisalone | 30 |
| designweek | 28 |
| nabasalone | 20 |
| milano | 9 |
| design | 6 |

| April 11th, 2013 at 18.00 | posts |
|---|---|
| milano | 25 |
| fuorisalone | 22 |
| design | 10 |
| designweek | 6 |
| 32giornata | 6 |

| April 13th, 2013 at 18.00 | posts |
|---|---|
| fuorisalone | 28 |
| designweek | 21 |
| nabasalone | 17 |
| milano | 10 |
| inter | 8 |

| April 15th, 2013 at 18.00 | posts |
|---|---|
| inter | 20 |
| diretta | 11 |
| cagliarii | 6 |
| milan | 4 |
| seriea | 3 |

MDW 2013 is visible in the top-5 hashtags used in the micro-posts

☺

| Venue | posts |
|---|---:|
| cesati antiques & works of art | 16653 |
| Porta nuova 46/b | 13416 |
| Circolo Filologico | 9891 |
| Adele Svettini Antichità | 7366 |
| ALTAI | 5592 |
| Bigli19 | 5175 |
| Dudalina | 4875 |
| Galleria DadaEast | 3550 |
| borronichemicals | 1078 |
| Antonio Lupi Showroom Milano | 995 |
| Instituto Cervantes Milano | 752 |
| GALLERIA D'ARTE CONTEMPORANEA CINESE | 560 |

The most attractive venues are found

☺

- The number of tweets posted by the same user where not enough to answer question 3 (How do people feel before, during and after the event they join?)
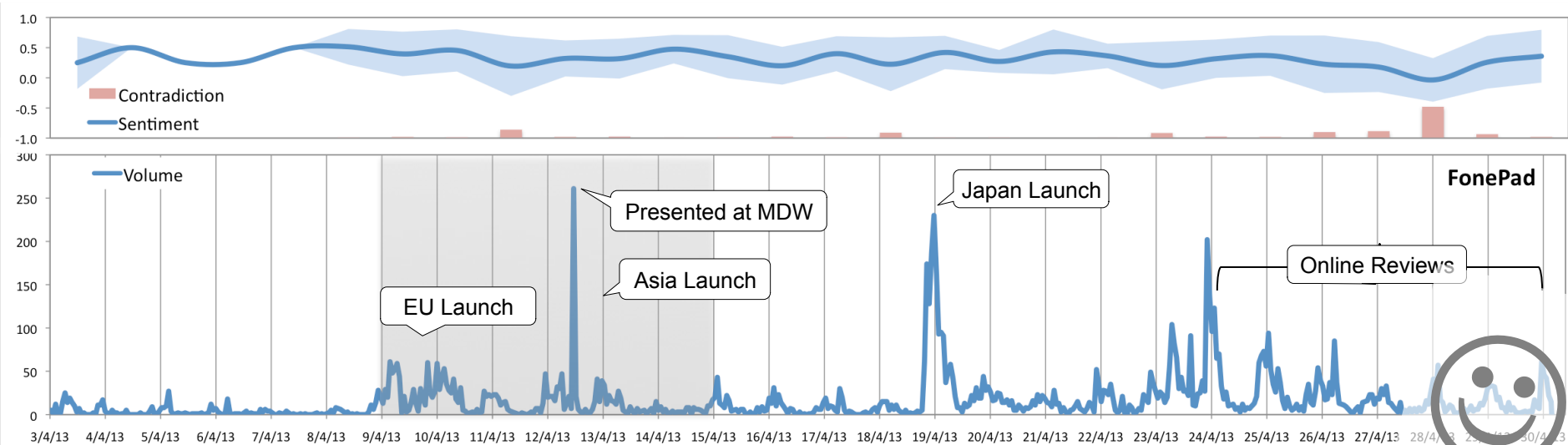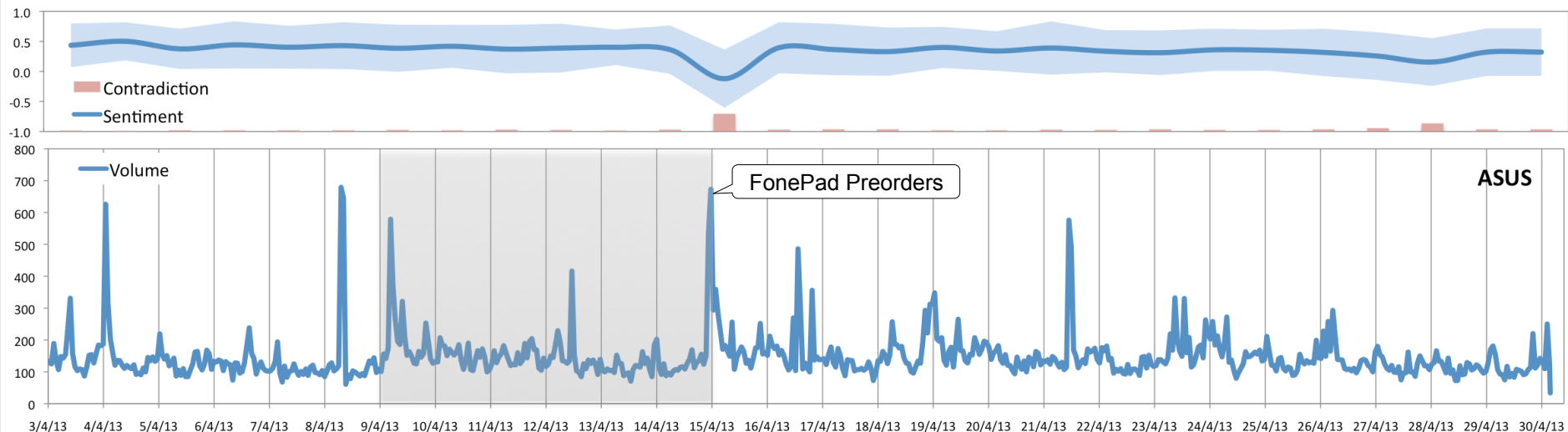
# Test bed MDW'13 – case study 2

- The problem
  - Is the launch of ASUS products during MDW visible in the social streams posted by people around the world?
  - If yes, not necessarily in real-time,
    1. What are the products that attract more attention?
    2. What is the global sentiment before, during and after the launch?

- Input
  - 107,044,487 tweets registered with SLD between April 3rd and April 30th, 2013 asking Twitter to send to SLD tweets containing 300 words related to MDW, ASUS and its products

- Ground Truth
  - News, movies, and other media published on the Web in the same time period
  - Tool: google advance search

- As expected, the method did not handle sarcasm in a satisfactory manner

- some tweets about FonePad contained sarcastic sentences
  - e.g., *"wanna buy it so bad!"*
  - It was classified as negative
  - it was expressing positive sentiment.

# Wrapping up

- Listening to social streams and proper visual analytics methods can **unveil interesting city scale phenomena**
  - where people are gathering
  - what people are interested in
  - where people interested in a given topic are
  - if an event is running
  - what people feels about some topic
  - if the people are perceiving the message an event organizer want to transmit
  - how the feeling of the people evolves over time
  - …

# Problematic aspects

- **lack of volume** in social streams prevents to perform meaningful analysis

- **Sarcasm and idioms** hinder the possibility to capture the opinion of people from highly volatile social streams

- **limited a priori knowledge** about the event hampers the ability to link social content to background data

- Basic research
  - continuous top-k query answering
  - crowdsource data cleansing and linking
  - determining what is systematic is difficult

- Applied research
  - profile a crowd
  - find opinion makers
  - predict social trends

# On-going collaborations

| Who | Semantic techs | Streaming algorithms | Crowd-sourcing | Visual analytics |
|---|:---:|:---:|:---:|:---:|
| CEFRIEL | ■ | | ■ | |
| Density Design Lab – PoliMi | | | | ■ |
| DISI – University of Trento | ■ | ■ | | |
| KDD Lab – ISTI, CNR, Pisa | ■ | | ■ | |
| ML Group – SIEMENS | ■ | | | |
| Ontology Eng. Group – UPM | ■ | ■ | | |
| Saltlux – Korea | ■ | ■ | | ■ |
| SKIL Lab - Telecom Italia | ■ | ■ | | |
| Studio Labo | | | | ■ |
| Web IS - TU Delft | ■ | | ■ | |

# City Data Fusion

http://citydatafusion.org

DeRiVE 2013 Workshop
21.10.2013, ISWC 2013, Sydney, Australia

# Thank you! Any question?

Emanuele Della Valle

emanuele.dellavalle@polimi.it

http://emanueledellavalle.org