

# Towards Real-Time Summarization of Scheduled Events from Twitter Streams



City University of  
New York  
New York, NY, USA

Arkaitz Zubiaga, Damiano Spina, Enrique Amigó, Julio Gonzalo

arkaitz.zubiaga@qc.cuny.edu, {damiano,enrique,julio}@lsi.uned.es

23rd ACM Conference on Hypertext and Social Media (HT2012)  
Milwaukee, WI, USA. June 25-28, 2012.



UNED NLP & IR Group  
Madrid, Spain

## Introduction

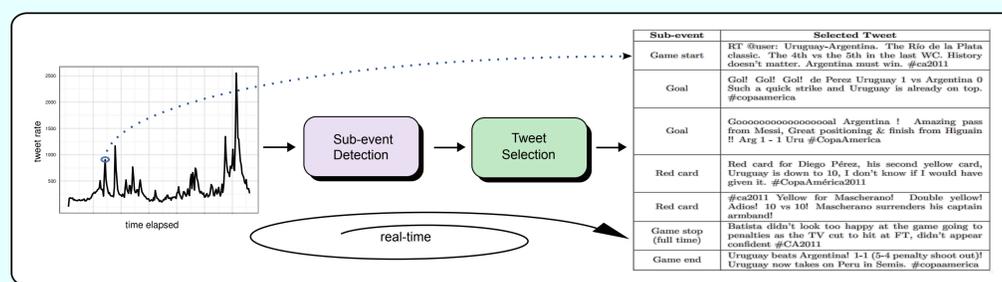
- Users **live-tweet** about events they are following in real-time.
- The **overwhelming amount of tweets** makes it difficult for the user:
  - to **follow the full stream** while finding out about new sub-events
  - to **retrieve** from Twitter the **main**, summarized **information** about which are the **key things happening** at the event.
- **Real-time summarization of events** enables to **follow streams** and to **find out about main sub-events** happening during the event.

## Dataset

- Soccer competition: **#copaamerica 2011**.
- **26 soccer games - 1.4M+ tweets**.
- **Multilingual** (majority of es, pt & en).
- **Reference summaries from Y! Sports**:
  - Key sub-events (goals, red cards, penalties,...)
  - Timestamps.



## System



### Sub-event detection

By looking at the histogram in real-time, are we observing an outstanding activity right now?

1. Increase (based on recent activity).
2. **Outliers (based on full histogram).**

### Tweet selection

If a sub-event happened right now, what tweet describes it?

1. TF (representative tweet right in the moment).
2. **KLD (representative tweet comparing to whole game).**

## Results

### Sub-event detection

	Precision	Recall	F1	Redundancy
<b>Increase</b>	0.29	0.81	0.41	45.4
<b>Outliers</b>	<b>0.51</b>	<b>0.84</b>	<b>0.63</b>	<b>25.6</b>

### Tweet selection

	Recall		
	es	en	pt
<b>TF</b>	0.79	0.74	0.78
<b>KLD</b>	<b>0.84</b>	<b>0.77</b>	<b>0.82</b>

$$\text{Recall} = \frac{\text{sub-events in reference}}{\text{correct tweets in summary}}$$

	Precision		
	es	en	pt
<b>TF</b>	0.79	0.74	0.79
<b>KLD</b>	<b>0.84</b>	<b>0.79</b>	<b>0.83</b>

$$\text{Precision} = \frac{\text{correct + novel tweets in summary}}{\text{total tweets in summary}}$$

## Conclusions

- **Real-time summaries** with **precision and recall** values above **80%**.
- The fact that **users tweet at the same time with overlapping vocabulary** helps not only **detect that a sub-event** occurs, but also **select a tweet** to describe it.
- Considering previous **information seen during the event is helpful** to this end, rather than only consider the most recent activity.
- **Regardless of the audience** tweeting about an event, **our method effectively reports** the key sub-events occurred during a game.
- The **most important sub-events**, such as goals and game ends, are **reported almost perfectly**.