

# Harnessing Social Media Streams for Local Information Needs

Dyaa Albakour

University of Glasgow

UCREL CRS, Lancaster University, 12 March 2015



@dyaaa



# Social media

As of December 2012 <sup>1</sup>:

- #users on Facebook **1.2 billion**
- #tweets per day **190 million**
- #pictures to Flickr **3,000/min**



- #people accessing the Web from mobiles **818.4 million**
- **26%** of mobile app usage is social networking <sup>2</sup>

The “2013 Q1 report” of the Global Web Index:<sup>3</sup>

- **A rise in active engagement** across all social platforms with Twitter the fastest growing (access from mobile phones)

1 <http://www.statisticbrain.com/social-networking-statistics/>

2 <http://www.pswebdesign.com/social-media-and-mobile-phones/>

3 <http://www.globalwebindex.net/Stream-Social>

# Local Information Needs and Social Media



University  
of Glasgow

- **Local Search is attracting more demand**
  - Local Search constitutes 43% of Google Queries <sup>1</sup>
  - *What is happening near me now?*  
*“near me”, “in Lancaster”, “on campus”*
  - Activities I can do now or later today
- **People are using social media to reflect on real-world events in real-time [1]**
  - Communicating to their social circle (what is happening? what are they doing? where are they? ..)
  - *Sporting events, earthquakes, protests, riots..*

<sup>1</sup> <http://chitika.com/insights/2012/local-search-study/>

[1] Yardi, S., and boyd, D. Tweeting from the town square: Measuring geographic local networks. In ICWSM'10.

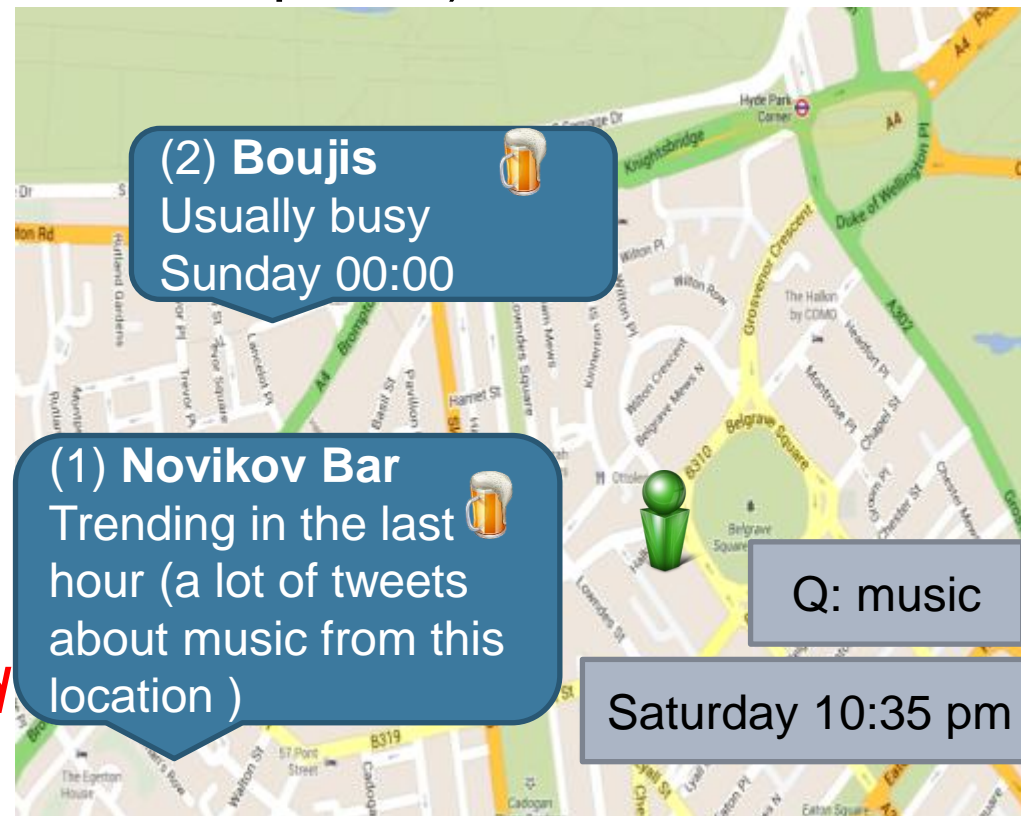
# Interaction Scenarios

## Input

- Keyword queries or zero-queries;
- Context (time, location and/or user profile)

## Output

- **Retrieve** and rank events that has currently started from social media posts
- **Filter** social media content about the event
- **Anticipate and recommend** locations that may have interesting activities for the user





***In this talk***



University  
of Glasgow

**Local Event Retrieval using Twitter as a Social Sensor**

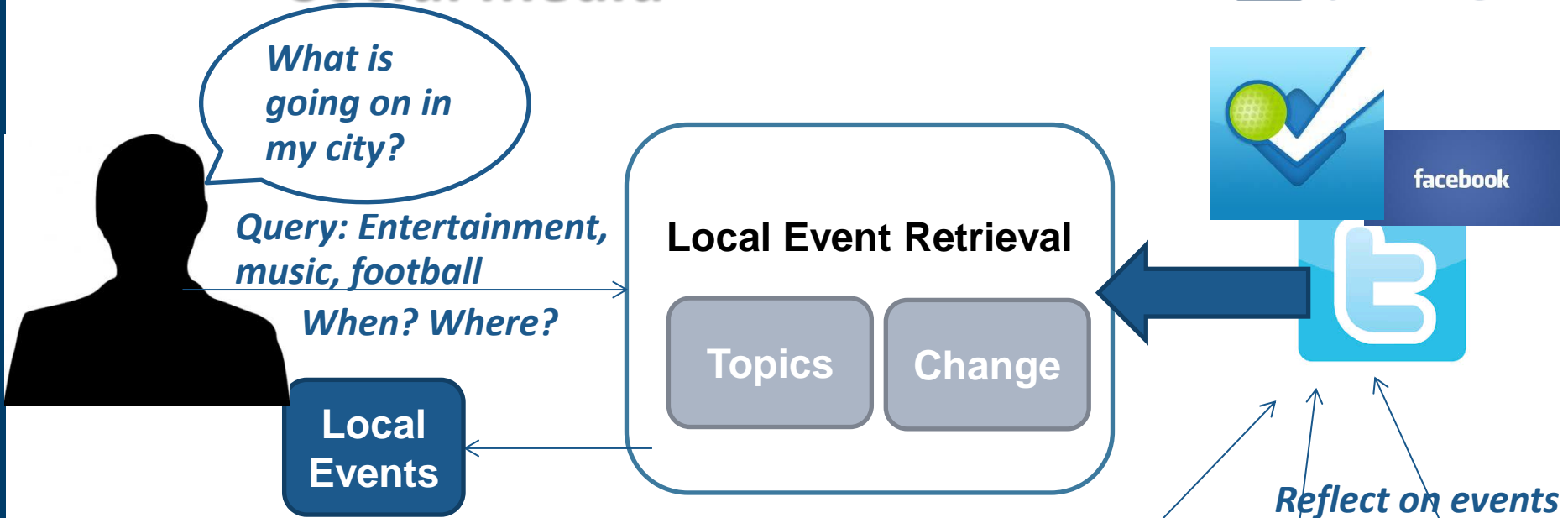
**Twitter Real-time Filtering**

**Anticipation and Personalised Venue Recommendation  
using Location-based Social Networks (LBSNs)**

Using Twitter as a Social Sensor

# LOCAL EVENT RETRIEVAL FROM SOCIAL MEDIA

# Local Event Retrieval from Social media



# Contributions



University  
of Glasgow

- **The new task of Local Event Retrieval from Twitter (Twitter as a social sensor)**
- **A framework for Local Event Retrieval**
- **Evaluation with a newly created dataset using crowdsourcing and local news feeds**

*M-D. Albakour, C. Macdonald and I. Ounis. Identifying Local Events by Using Microblogs as Social Sensors. In proceedings of OAIR 2013.*

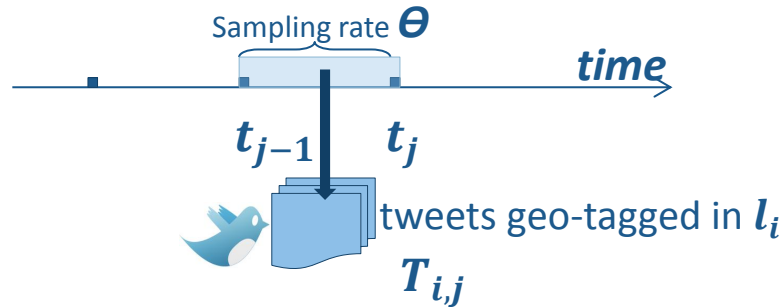


# Local Event Retrieval using Twitter



University of Glasgow

- Given a user query ( $q$ ):
  - Retrieve a ranked list of **local events** that are **relevant** to the user query ( $q$ )
- We model a location as a time series

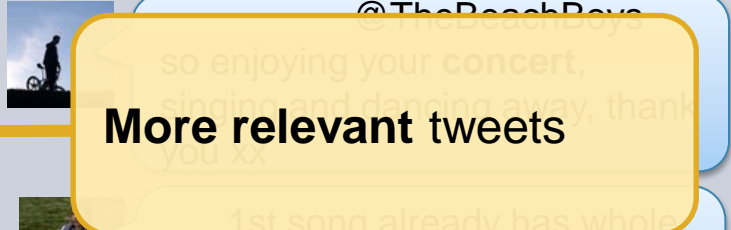
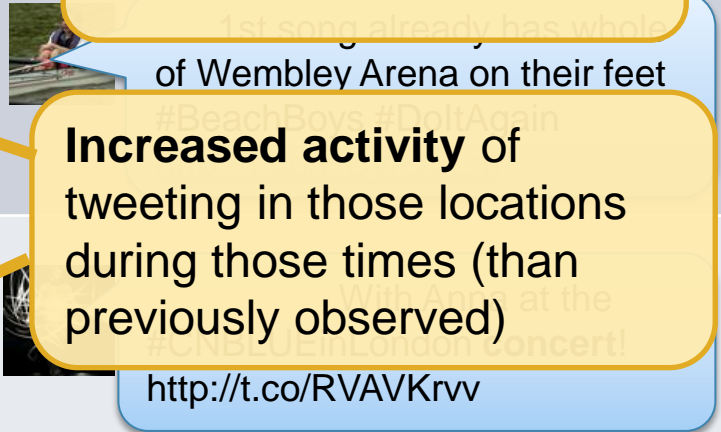
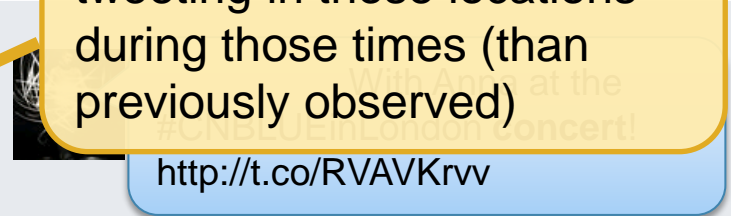


- What people tweet reflect what is happening in a location at a certain time
- A local event has (1) a starting time and ; (2) a location

Ranking function  $R(q, \langle l_i, t_j \rangle)$ :

Rank tuples  $\langle l_i, t_j \rangle$  according to how likely  $t_j$  represents a starting time of a matching event that occurred in  $l_i$  using the tweets

## Example of responses for query (concert)

Rank	Start Time	Location	Description (Tweets)
1	Today 19:15	Wembley	 
2	Yesterday 20:00	London O2	
3		..	..

# A Framework for Local Event Retrieval

## Two Components:

- Topically related tweets to  $q$  in location  $l_i$  at around  $t_i$
- Increasing tweeting activity

$$R(q, \langle l_i, t_j \rangle) \sim (1 - \lambda) \cdot S(q, T_{i,j}) + \lambda \cdot E(q, \langle l_i, t_j \rangle)$$

Q  
T  
qu  
(1

The voting model to aggregate ranking of individual tweets

Quantifies the **change** in the **tweeting activity** at time  $t_j$  in location  $l_i$

**(2) The change component**

## How do we estimate the change in the tweeting activity?

### Change point Analysis

- Quantify how likely is the tweeting activity (about a topic) is an outlier with respect to previous observations.
- Apply the Grubb Test [2]  
*Normalised score (0..1)*



**The tweeting activity:** is measured by the topical component score  $S(q, T_{ij})$

[2] F. Grubb. Procedures for detecting outlying observations. *Technometrics*, 11, 1969

## Research Question:

- What is the ***impact*** of the different components, in our framework, on the **ranking effectiveness**?



# Datasets



Code	Tweets	Events	Queries
Entire London	1.28m geotagged within London  12 days (2012-03-10 to 2012-03-22) → 3/10/12	<ul style="list-style-type: none"><li>• <i>Young believers choir concert</i></li></ul>	
4 boroughs	864k geotagged within 4 central boroughs of London  12 days (2012-03-10 to 2012-03-22) → 3/10/12	<ul style="list-style-type: none"><li>• From local news</li><li>• <i>Hospital volunteers honoured</i></li></ul>	The title of the

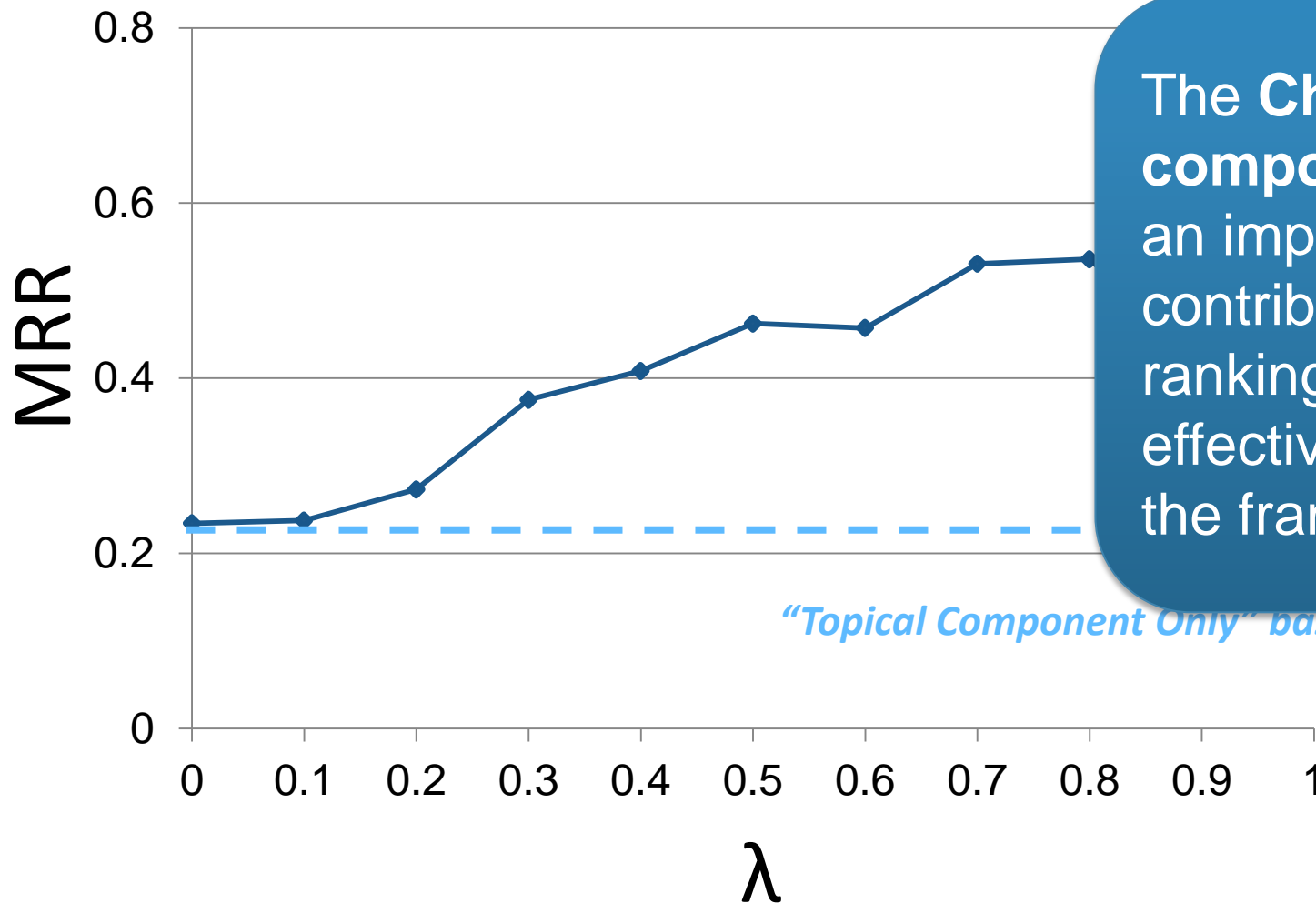
- Coarse-grained
- Single event for each query
- Popular events

- Finer-grained
- Single event for each query
- More localised events

# *Experimental Setup*

- A sampling rate of 15 minutes
- *DFReeKLIM* for ranking tweets in the voting model
- **Baseline:** using the topical component only ( $\lambda=0$ )
- Evaluation methodology inspired by the video segmentation evaluation for assessing the accuracy of *correctly identifying the starting time* of an event

# Results (Entire London)

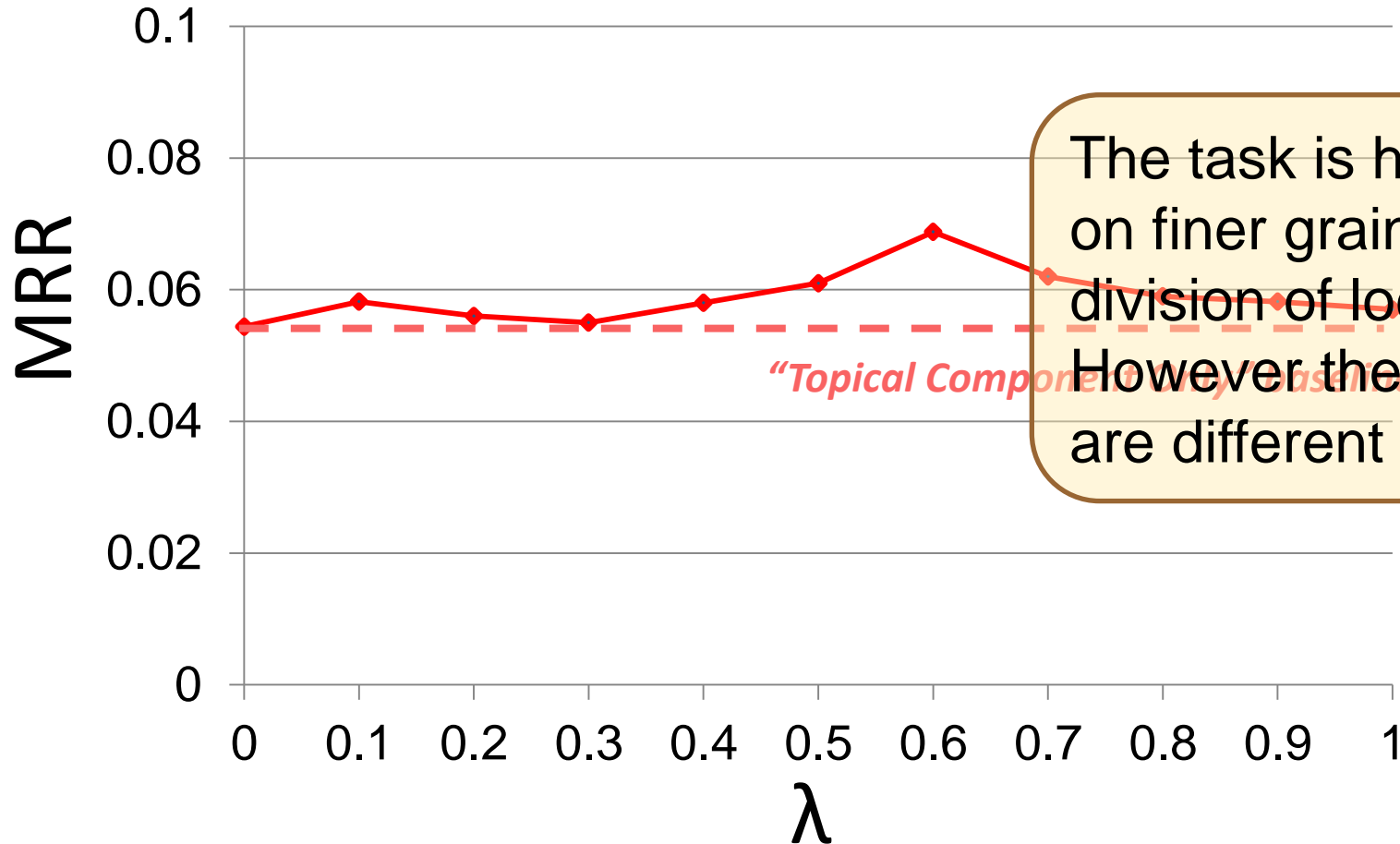


The **Change component** has an important contribution on the ranking effectiveness of the framework

*"Topical Component Only" baseline*

*Scores obtained for the MRR measure*

# Results (4 boroughs)



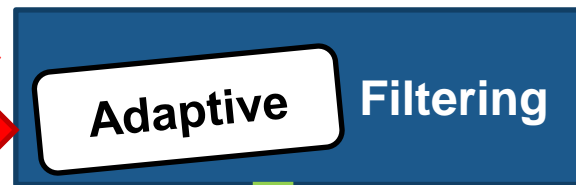
Scores obtained for the MRR measure

# REAL-TIME TWITTER FILTERING



# Real-Time Tweet Filtering

- **Producers:** Huge activity around the globe (on average around 5700 published tweets per second)<sup>1</sup>
- **Consumers:** want to stay up-to-date with **relevant** content (not everything!)



Filtered tweets



Thousands of tweets per second

Interested in a Topic(s):  
a **query**, and a starting  
time

<sup>1</sup> <https://blog.twitter.com/2013/new-tweets-per-second-record-and-how>

# Challenges

Tweets are very short documents

→ vocabulary mismatch



US Unemployment



The tweet does NOT contain any of the terms in the query or the user profile

1- Sparsity  
(Brevity)

Google News #RonPaul Chairman Ron Paul to Tackle the Fed and Jobs - The New American <http://goo.gl/fb/CLjEp>

Liberty  
Click  
Click



Thu Feb 03 2011 16:26:30

# Challenges

Twitter is a **highly dynamic** social medium (it reflects a highly dynamic world!)

*Long-term interest*



Classic



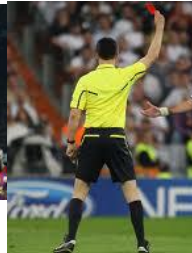
Who is going to play/miss?

*short-term interest*

Goals / cards / chances ?

*short-term interest*

2- Drift of interests



*Before the match*

*During the match*

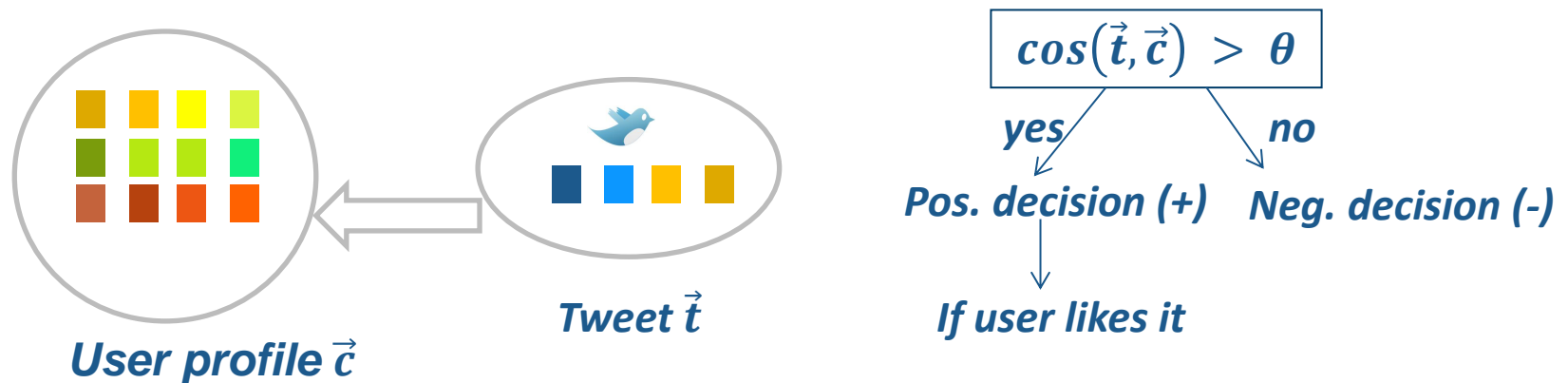
The interests swing between different aspects (subtopics) of the more general topic (due to **events** in the real world)

- **Build on news filtering approaches to tackle the problem of adaptive tweet filtering**
- **Address the unique challenges in filtering tweets:**
  - Address **Sparsity** by deriving a richer representation of the user profile
  - Address **Drift** by balancing between short-term and long-term interests

*M.-Dyaa Albakour, Craig Macdonald, Iadh Ounis: On sparsity and drift for effective real-time filtering in microblogs. CIKM 2013: 419-428*

# Tweet Filtering with Incremental Rocchio

- We build on a common technique for News Filtering: the popular Incremental Rocchio's classifier (RC) [3]
  - Build a profile online (vector of terms)



- We considered another state-of-the-art news filtering approach of Regularised Logistic Regression (LR) [4]
  - Evaluation suggests that Incremental Rocchio (RC) significantly outperforms LR (full details in the paper).

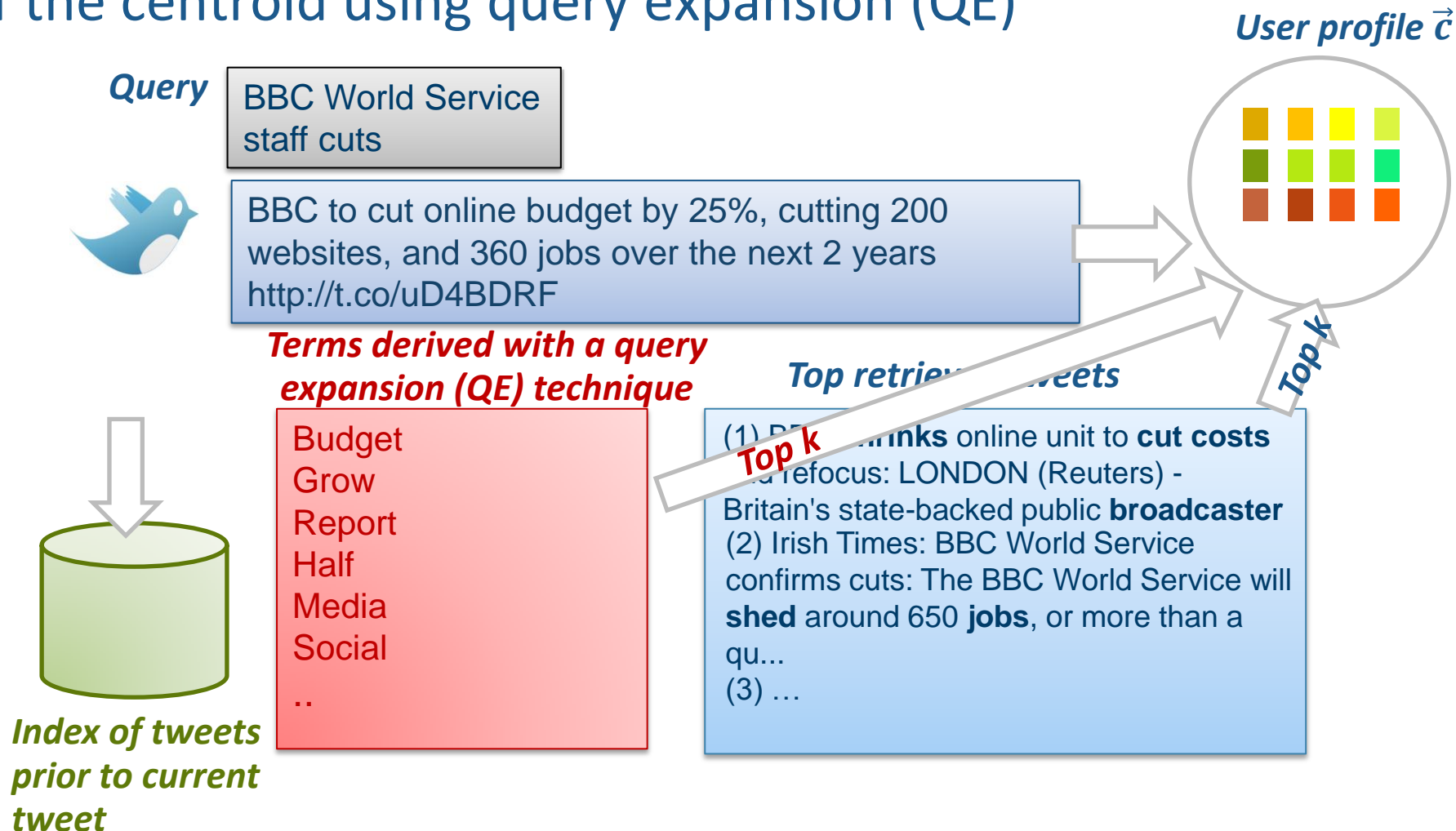
[3] J. Allan. Incremental relevance feedback for information filtering. In Proc. of SIGIR, 1996.

[4] Y. Zhang. Using bayesian priors to combine classifier for adaptive filtering. In Proc. of SIGIR, 2004



# Handling Sparsity

Derive **relevant** and **timely** terms for a richer representation of the centroid using query expansion (QE)



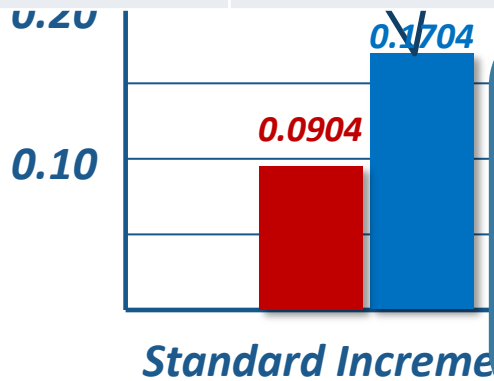
- TREC 2012 Microblog Track – Real-time Filtering task
  - Tweets2011 (around 10m tweets over 16 days)
- We have built a real-time filtering infrastructure
  - using Storm and Terrier<sup>1</sup>
- **Experimental Setup**
  - Standard stopword removal and Porter stemming
  - Dirichlet language modelling to weight terms in the vectors
  - Threshold *tuned* on the 10 TREC training topics (38 testing topics)
  - Bo1 DFR for query expansion (as provided by Terrier)

**Research Question:** Are our adaptations for tackling sparsity, using QE, successful in **improving filtering effectiveness?**

# Results: Sparsity

Significantly improves F and utility ✓

	Set_Prec	Set_Recl	F_0.5	T11SU
RC + Qe + Te	0.4206	<b>0.3370</b>	<b>0.3435</b>	0.3615
TREC 2012 Best approach	<b>0.6219</b>	0.1740	0.3338	0.4117

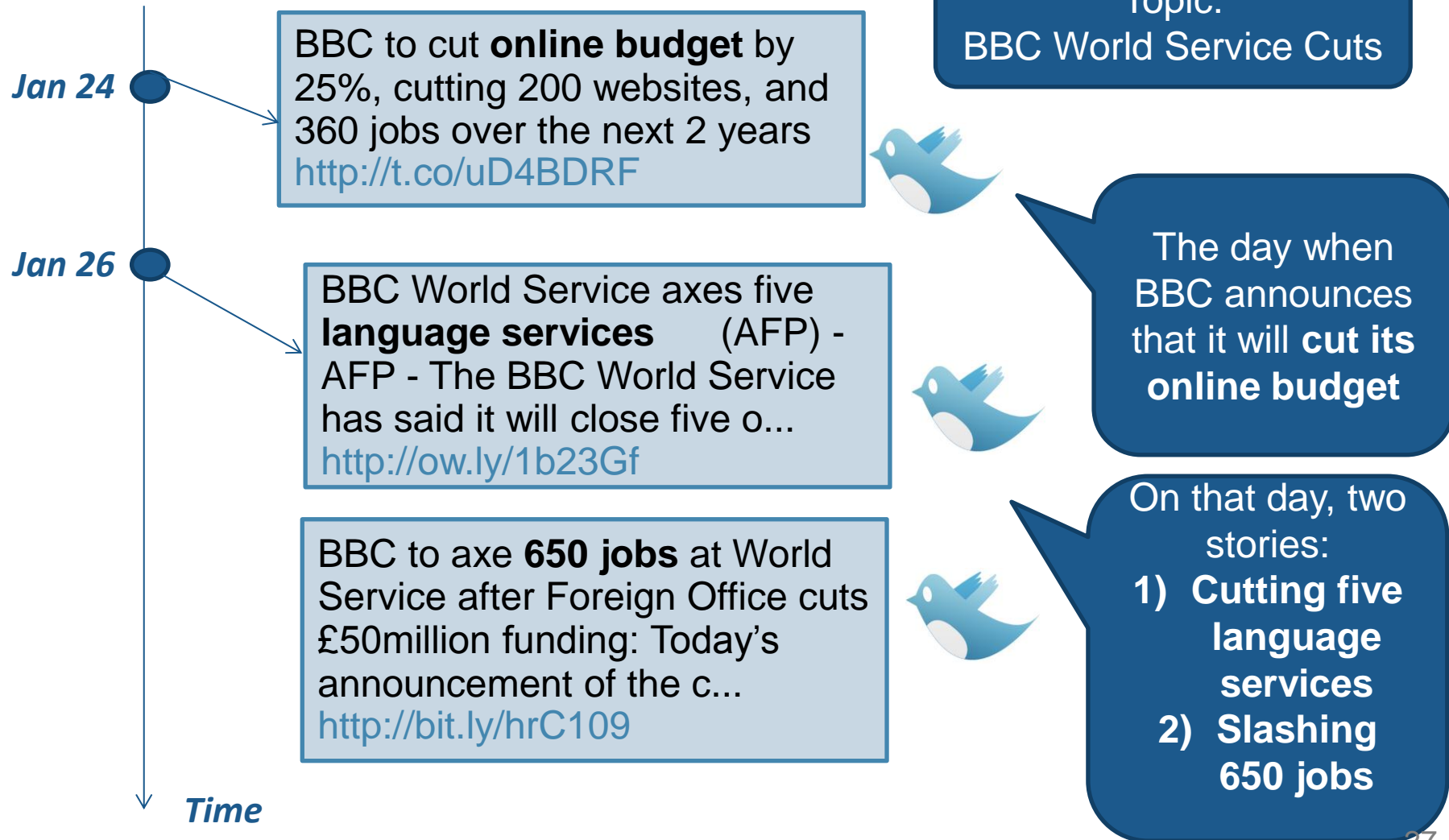


Our approach is *more balanced* as opposed to the *conservative* best TREC 2012 approach!

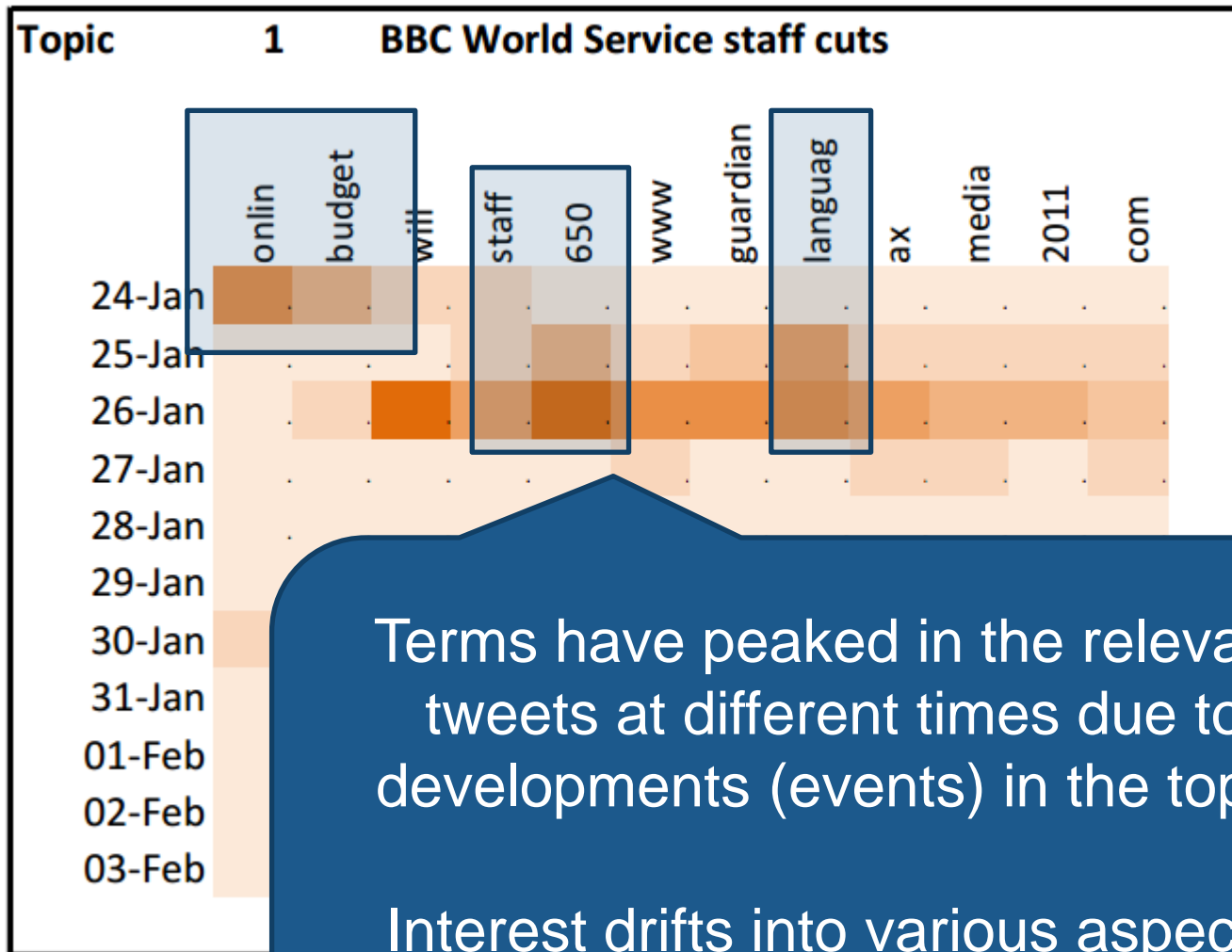
TREC 2012 Best Approach  
 - - - F\_0.5  
 - - - T11SU

# What is Drift?

## Illustrative Example



# Empirical Viewpoint of Drift

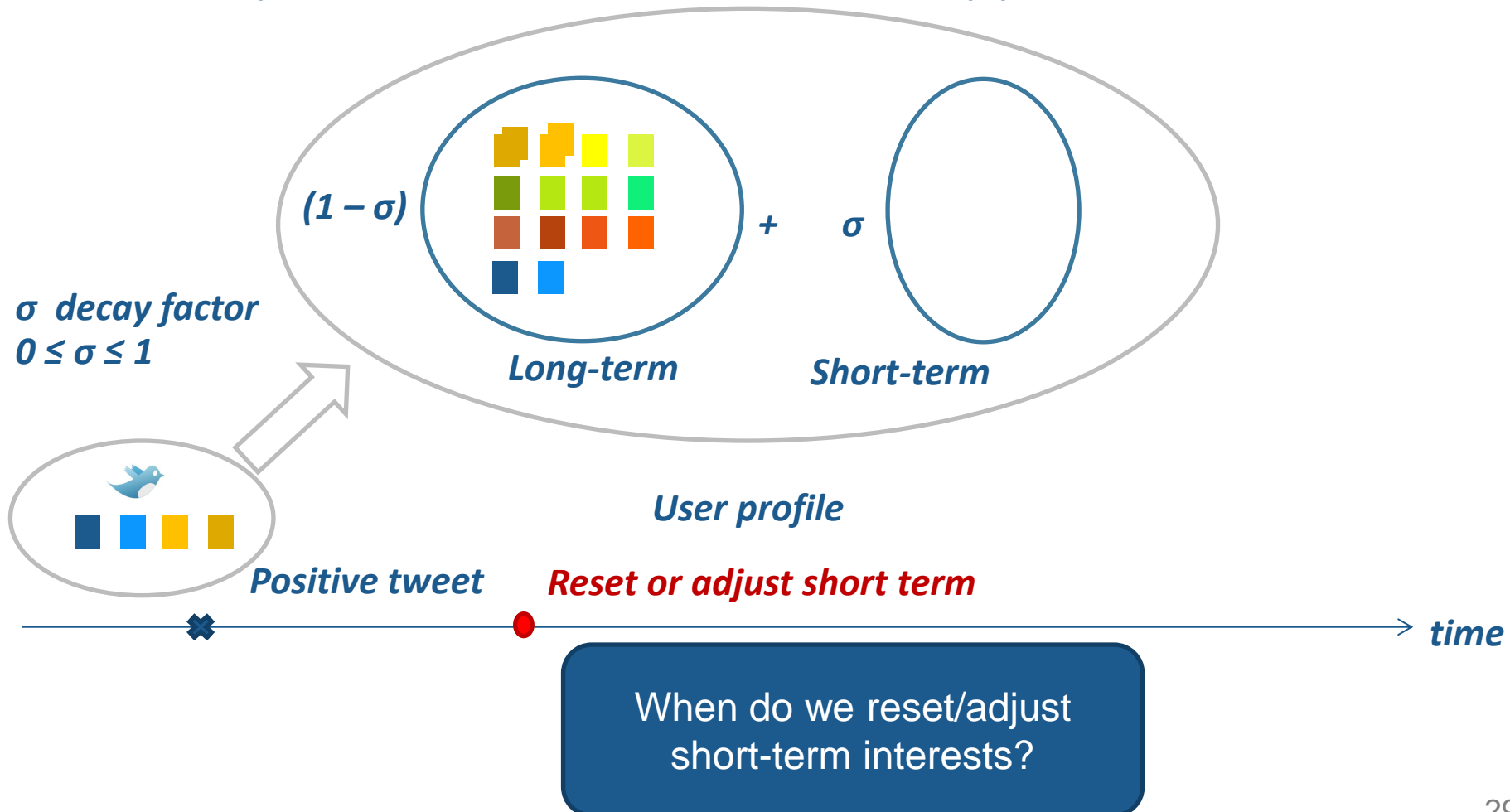


Terms have peaked in the relevant tweets at different times due to developments (events) in the topic

Interest drifts into various aspects (sub-topics) over time

# Handling Drift

- Dynamically changing the centroid over time to represent both **short-term** interests and **long-term** interests in the overall topic (combined with the QE approach)

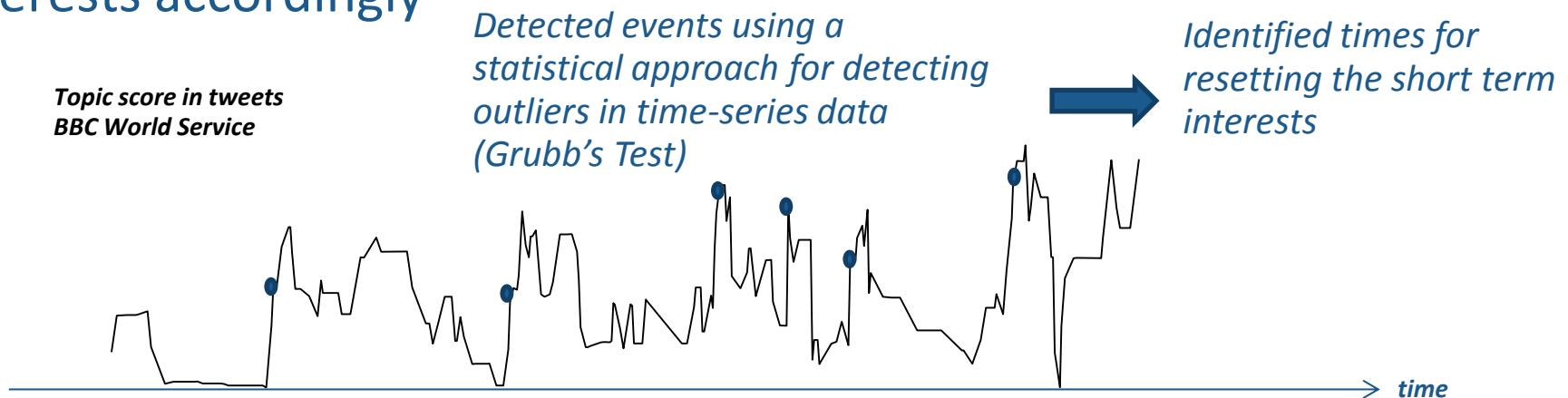


# When does drift occur?

When do we reset/adjust short-term interests?

1. **Arbitrary adjustments:** The most recent  $n$  positive tweets
2. **Daily adjustments:** The tweets in the current calendar day
3. **Event detection** [5] to automatically identify times when events related to the topic occurred and reset the short-term interests accordingly

Adhoc



Event detection can be applied on the Twitter stream itself or external news streams

# Experiments: Drift

Identical setup to the one used before

The QE approach for handling sparsity as a baseline

The newswire stream

- BBC, CNN, Google News, New York Times, Guardian, Reuters, The Register and Wired

## Research Questions:

- (1) **Adhoc methods** vs. **event detection** for handling drift?
- (3) **sensitivity** of the filtering performance to the **decay factor**  $\sigma$ ?



# Results: Spars

- ✗ Adhoc methods **failed**
- ✓ The recall is slightly improved.

- ✓ The **increase in recall is significant.**

- ✓ **Event detection** is helping!

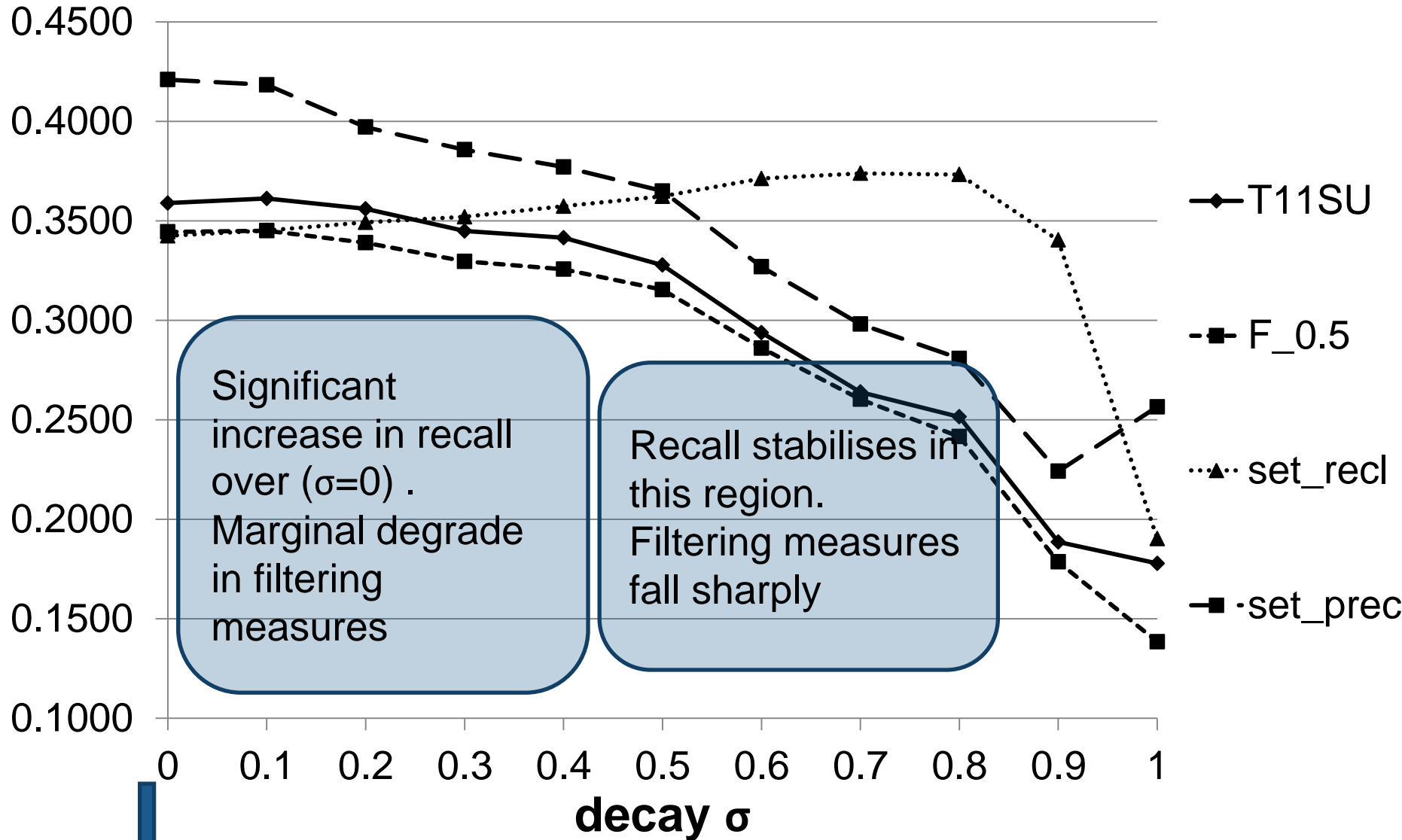
- Differences are marginal when using a different stream for events! (Events overlap in both streams)

Bas  
RC+  
Arbi  
(n=1  
Dail  
(σ =  
Eve  
Twee  
(σ = 0.4)  
Event detection using  
News Streams  
(σ = 0.4)

	0.3724 ▼	<b>0.3598 ▲ ▲</b>	0.3198 ▼	0.3351 ▼

A single triangle means the differences are not statistically significant using a paired t-test at  $p < 0.05$ . Double triangles mean the differences are statistically significant

# Sensitivity to decay



Baseline: no difference between short-term and long-term

More emphasis on short-term interests

# Conclusions

- Tackled sparsity and drift for real-time twitter filtering
- **State-of-the-art** for real-time twitter filtering
- With an **event detection** approach to tackle drift, we can **significantly improve** the filtering **recall** while only marginally harming the filtering utility

# ANTICIPATION AND PERSONALISED VENUE RECOMMENDATION

# Venue Recommendation

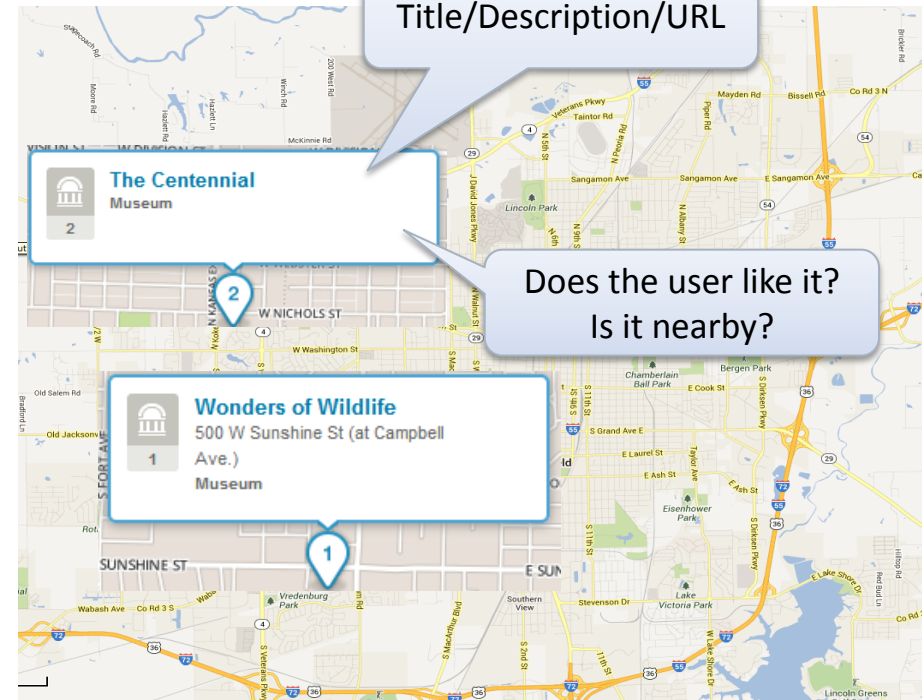


Entertain me!

Zero-query



Location ( **Springfield** ) + time ( 2 pm ) + ..



Venue recommendation has different potential use cases:

- **Tourists** use case: *"I have one day in this city, what should I see?"*
- **Residents** use case: discover/explore new venues, avoid noisy or polluted places, ...

# Existing Services

## What do people currently use?

–A tourist guide, The List, Yelp, FourSquare?

## No anticipation of venue popularity...

Suggestions for **Best Nearby** near **Kelvingrove Park, Glasgow**

Or try: **Food, Nightlife, Coffee, Shops, Arts, Outdoors**

Show me: **Specials** Haven't Been Friends Price Open Now Saved

### 1. Kelvingrove Park

8.2  
14 Parkgrove Tce.  
Park

Popular with out-of-towners

"... A CARGO! ITZ MARE FUN!! WOOP WOOPZI!" (1 tip)

Save



### 2. Marks & Spencer

7.5  
165-169 Great George St. (Byres Rd)  
Grocery Store

You haven't been here yet

Stuart · August 7, 2011  
Don't arrive hungry!

Save

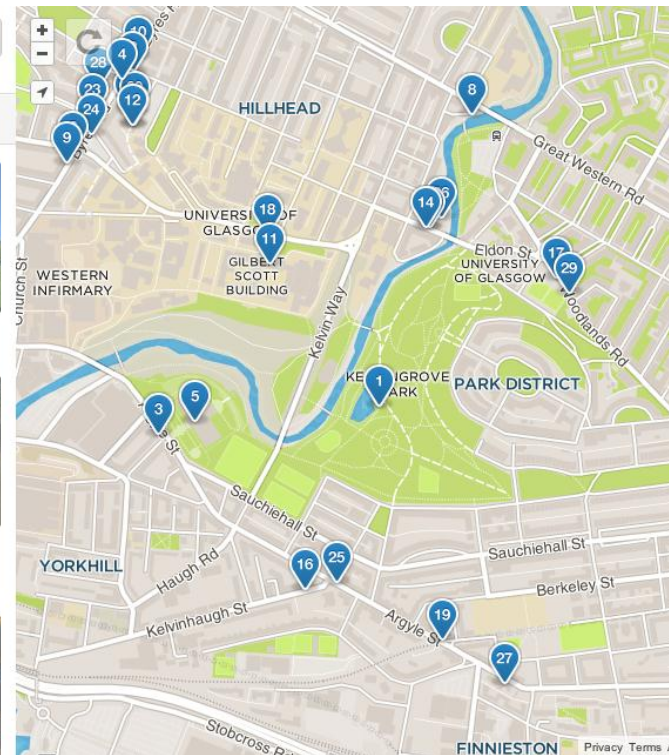


### 3. BrewDog Glasgow

8.6  
1397 Argyle St  
Bar · 1 · ££££

You haven't been here yet

"Try the Killa Morchilla burger. Yummy yummy...." (3 tips)



Recommendations from Foursquare at 10pm, in March

## Venue recommendation: help users decide where to go

*“I’m new to the city. What should I visit?”*

**We argue that effective venue suggestions should encompass:**

- Cold-start: we don’t know where you have been before
- Personalised: recommend venues that I would **like**
- Time-aware: Quality venues will be **popular**

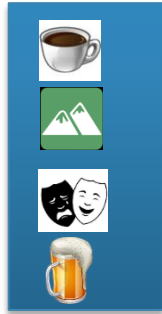
**We developed and evaluated a probabilistic model for time-aware and personalised venue recommendation**

# Ranking Venues

Input



Output



Ranked list of venues

$$P(\text{building icon} | \text{person icon}, \text{stopwatch icon})?$$

Not available!

 Location ( *Springfield* )





# Venue Popularity

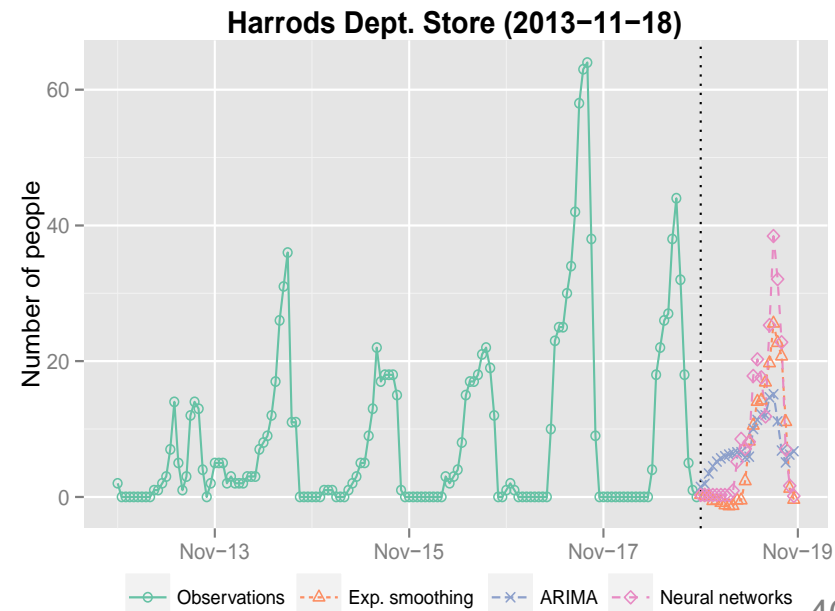


How busy a venue will be later in the near future (in the next few hours)

- we anticipate how **popular** the venue will be

Popularity – we **forecast** the attendance of venues based on past Foursquare checkins

- *Anticipating the future attendance*
- Foursquare API as a *social sensor* of the level of venue attendance (“check-ins”)
- time series forecasting models



# ***Personalisation***



University  
of Glasgow

Not available!

# *Personalisation*



University  
of Glasgow

Not available!

# *Evaluation – venue popularity*



University  
of Glasgow

Not available!

# *User Study*



University  
of Glasgow

Not available!

# *User Study*



University  
of Glasgow

Not available!

# *User Study*



University  
of Glasgow

Not available!

# *Results of the User Study*



University  
of Glasgow

Not available!



***Thanks!***



University  
of Glasgow

## Acknowledgments

This work has been carried out in the scope of the EC co-funded project SMART (FP7-287583).

Co-authors:

Romain Deveaud, Craig Macdonald, Iadh Ounis



*[dyaa.albakour@glasgow.ac.uk](mailto:dyaa.albakour@glasgow.ac.uk)*