

# TweetRank

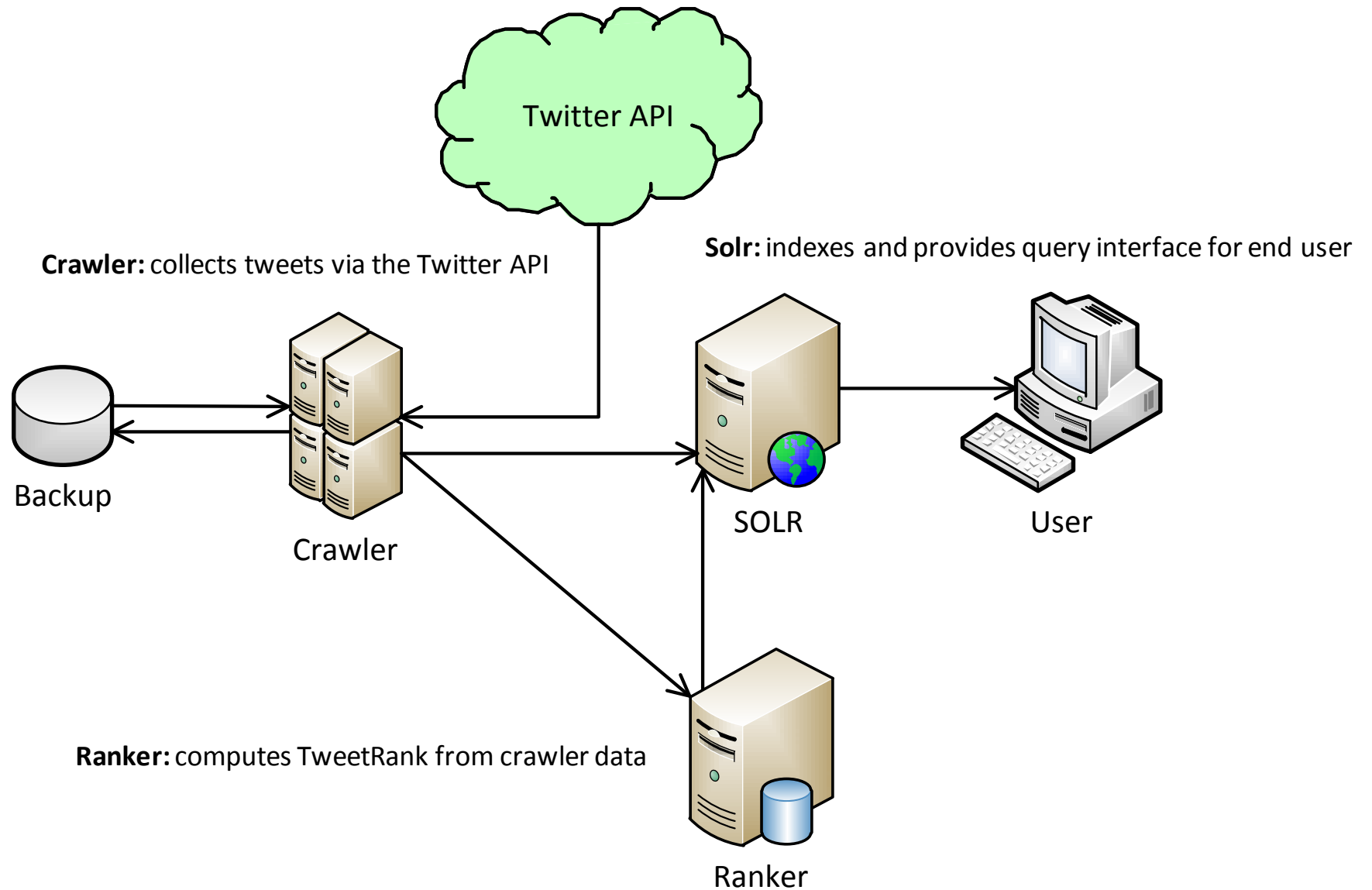
TweetRank is an attempt to apply the PageRank algorithm on Twitter statuses (tweets).

It uses a different rank calculation which considers attributes such as number of replies/retweets and hashtags.

Victor Hallberg	<a href="mailto:victorha@kth.se">victorha@kth.se</a>
Johan Stjernberg	<a href="mailto:stjer@kth.se">stjer@kth.se</a>
Joan Puigcerver I Perez	<a href="mailto:joanpip@kth.se">joanpip@kth.se</a>
Alexander Hjalmarsson	<a href="mailto:alehja@kth.se">alehja@kth.se</a>
Christoffer Rydberg	<a href="mailto:chrryd@kth.se">chrryd@kth.se</a>

# Components

---



# Crawler

- Uses the Twitter *HTTP REST API*

- Twitter limits the number of queries to 150 per hour
- Crawler gathers as much data as possible from each query
- Use multiple proxies to bypass the query limit
- Runs on multiple threads in multiple machines

- How does it work?

1. Start with a queue of some users
2. Pop the first user from the queue and query the Twitter API for tweets and friends for it
3. Add friends and user mentions in each tweet to the user queue
4. Send tweet data to Solr and the ranker
5. Go to 2

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <statuses type="array">
3   <status>
4     <created_at>Tue May 15 18:24:39 +0000 2012</created_at>
5     <id>202464528154365953</id>
6     <text>Get updated data about smartphone usage around the
world http://t.co/IDreU8Fd - more info: http://t.co/
mQNm02yc</text>
7     <source>web</source>
8     <truncated>>false</truncated>
9     <possibly_sensitive>>false</possibly_sensitive>
10    <user>
11      <id>20536157</id>
12    </user>
13    <retweet_count>83</retweet_count>
14    <favorited>>false</favorited>
15    <retweeted>>false</retweeted>
16  </status>
17 </statuses>
18 <status>
19   <created_at>Tue May 15 16:54:53 +0000 2012</created_at>
20   <id>202441934985761025</id>
21   <text>Congrats to the team behind the amazing 3D models
of Getaria, Spain - winners of the 2012 Model Your Town
competition http://t.co/VFpn697t</text>
22   <source>web</source>
23   <truncated>>false</truncated>
24   <possibly_sensitive>>false</possibly_sensitive>
25   <user>
26     <id>20536157</id>
27   </user>
28   <retweet_count>125</retweet_count>
```

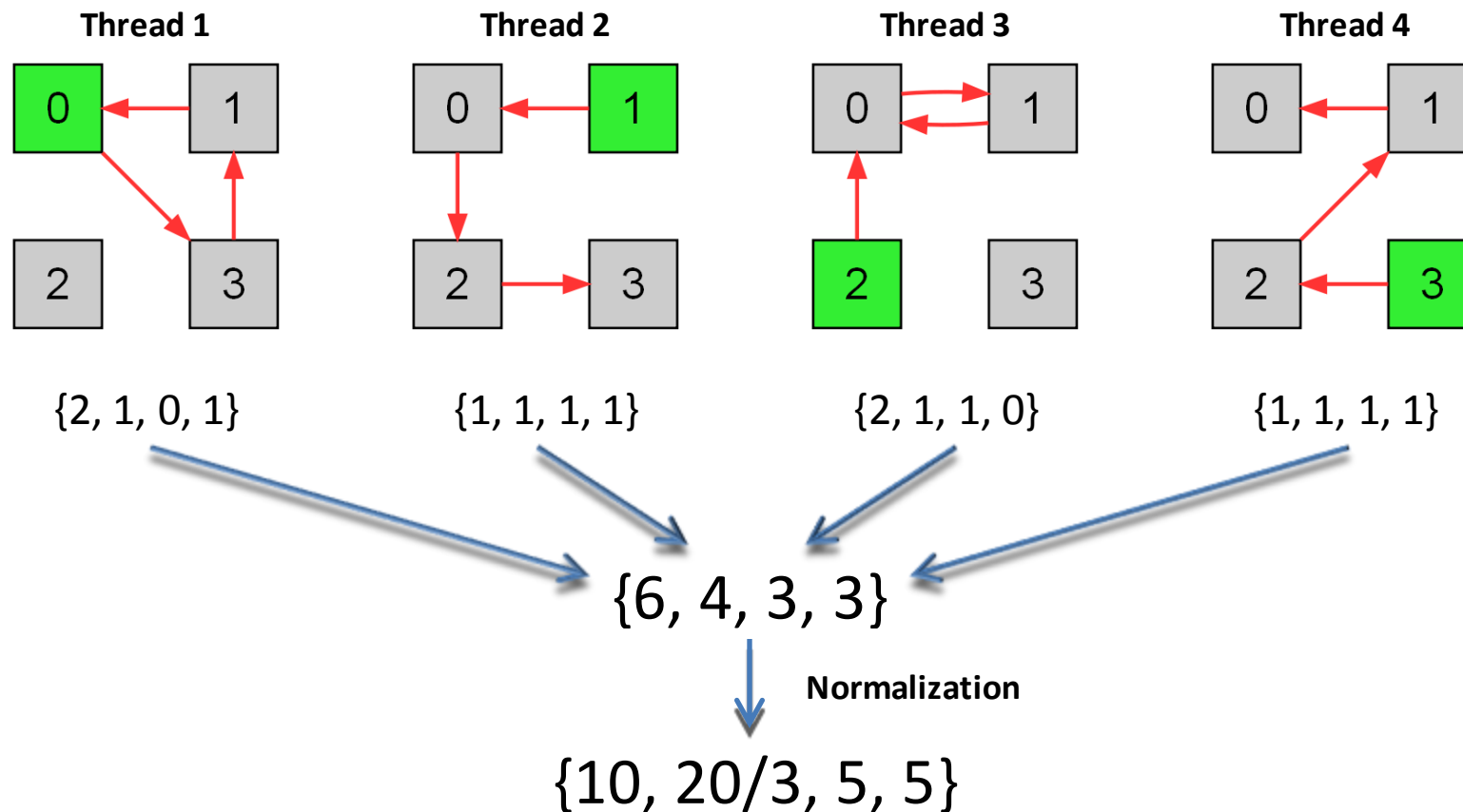
# Rank algorithm

---

- Uses the complete path *Monte Carlo* algorithm, stopping at dangling nodes
- Starts a **randomized walk from each node**
  - at least 100 times
  - at most total number of tweets / 100
- Random path length
  - 20% chance at each node that the surfer stops
- **Probability of visiting tweet  $x$  from  $y$**  estimated from:
  1. Random access
  2. Retweeted or replied
  3. Author of  $x$  mentioned by tweet  $y$
  4. Author of  $x$  followed by author of  $y$
  5. Hashtag shared by tweets  $y$  and  $x$
- A stochastic matrix is built with these probabilities
- TweetRank is the eigenvector of this matrix

# Ranker

- Ranker runs on multiple threads, where each thread computes one walk at a time
- Rank for each tweet is calculated as the sum of visits for each node from every walk
- Normalization - divide by the maximum number of visits for a node and then multiply by 10



# Solr / Lucene

- Handles indexing and searching
- Crawler sends tweets to be indexed through *HTTP POST* requests
- **Scores** are calculated as a product of:
  - TweetRank
  - *tf-idf* (hashtag matches are boosted)
- **TweetRank data** in text file on the server
  - Enables rank updates without having to replace (re-index) existing documents
  - Uses *ExternalFileField* feature in Solr



# Conclusions

---

- **PageRank can successfully be adapted to Twitter statuses**
- **Tweaking & optimization required for good results**
  - More variables than the original PageRank
    - Friend count
    - Mentions
    - Retweets
    - Hashtags
  - Can be simplified through statistical analysis of tweets
- **Hashtags are worth considering**
  - Easy way to find similar tweets on Twitter
  - Only 8.4% of the crawled tweets contain hashtag(s)
  - Bridge between tweets in TweetRank = increases rank
  - Boost query terms matching hashtags = more relevant
  - Hard to determine relevance between individual tweets - we didn't