

TagCurate:
Crowdsourcing the Verification of
Biomedical Annotations to Mobile Users

Bahar Sateli Sebastien Luong René Witte

Semantic Software Lab
Department of Computer Science and Software Engineering
Concordia University, Montréal, Canada

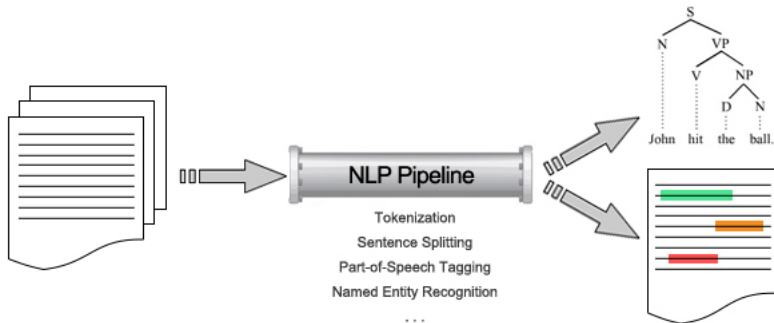
NETTAB 2013

- 1 Introduction
- 2 TagCurate System
- 3 Android-NLP Integration
- 4 Conclusion

Natural Language Processing (NLP)

Definition

A branch of Artificial Intelligence that uses various techniques to process content written in a natural language, e.g., English or German.



Bottleneck: Gold Standard Corpora

Manually annotated documents required for training & testing NLP pipelines (especially for machine learning components).

Can we 'crowdsource' some of this work to mobile users?

Challenge: Current Web-based annotation frameworks (e.g., GATE Teamware) not designed for mobile use

GATE Developer 6.1 build 3913

Messages | ANNIE | world-middle-east-206051

Annotation Sets | Annotations List | Annotations Stack | Co-reference Editor | Text

BBC News - Egypt crisis: Clashes in Cairo amid constitution row

Egypt crisis: Clashes in Cairo amid constitution row

Rival protesters have clashed outside the presidential palace in the Egyptian capital, Cairo, as unrest grows over a controversial draft constitution.

Stones were thrown and supporters of President Mohamed Morsi dismantled tents set up by anti-Morsi protesters.

Vice President Mahmoud Morsi has said a referendum on the draft will go ahead on 15 December despite the unrest.

But he indicated that changes could be made after the vote, saying the "door for dialogue" remained open.

He urged critics of the draft document to put their concerns in writing for future discussion.

The critics say the draft was rushed through parliament without proper consultation and that it does not do enough to protect political and religious freedoms and the rights of women.

The draft added to the anger generated by Mr Morsi passing a decree in late November which granted him wide-ranging new powers.

'Breakthrough'

Egyptian Vice-President 'Door open'

In a news conference broadcast live on state television, Mr Morsi said there was "real political will to pass the current period and respond to the demands of the public"

| Type | Set | Start | End | Id | Features |
|--------------|-----|-------|------|----|---|
| Organization | 0 | 8 | 5347 | | {matches=[5347, 5358, 5370, 5394, 5398, 5442, 5445, 5446, 5447, 5448, 5449, 5450, 5452], rule1=1} |
| Location | 11 | 16 | 5348 | | {locType=country, matches=[5348, 5356, 5360, 5375, 5388, 5395, 5396, 5407, 5416, 5437], rule1=1} |
| Location | 36 | 41 | 5349 | | {locType=city, matches=[5349, 5357, 5362, 5372, 5414], rule1=InLoc1, rule2=LocFinal} |
| Location | 65 | 70 | 5356 | | {locType=country, matches=[5348, 5356, 5360, 5375, 5388, 5395, 5396, 5407, 5416, 5437], rule1=1} |
| Location | 90 | 95 | 5357 | | {locType=city, matches=[5349, 5357, 5362, 5372, 5414], rule1=InLoc1, rule2=LocFinal} |
| Location | 306 | 311 | 5362 | | {locType=city, matches=[5349, 5357, 5362, 5372, 5414], rule1=InLoc1, rule2=LocFinal} |

146 Annotations (0 selected) Select:

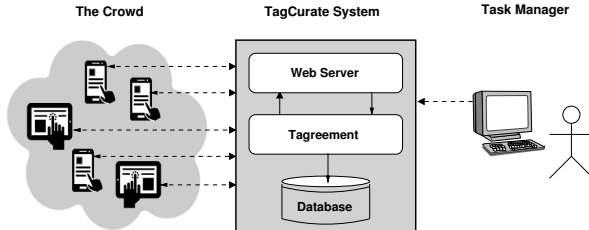
Document Editor | Initialisation Parameters

Address
☐ Address
☒ Date
☐ FirstPerson
☒ JobTitle
☒ Location
☐ Lookup
☒ Organization
☐ Percent
☒ Person
☐ Sentence
☐ SpaceToken
☐ Split
☐ Title
☐ Token
☐ Unknown
Original markups

New

- 1 Introduction
- 2 TagCurate System
 - System Architecture
 - Web-based Interface
 - Android App
- 3 Android-NLP Integration
- 4 Conclusion

System Architecture



Client-Server Model

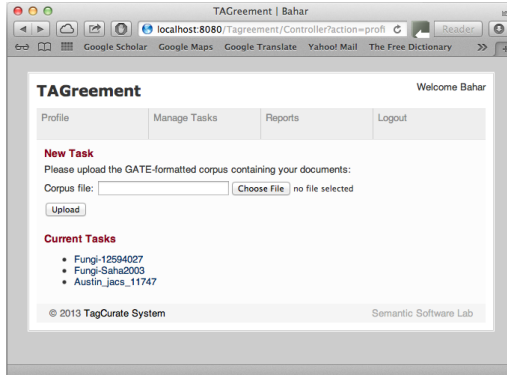
- RESTful communication over HTTP
- *Tagreement* component is responsible for managing the crowdsourcing as well as measuring (dis)agreements

User Groups

- **Task Managers**, define verification tasks using the web-based interface
 - e.g., *NLP pipeline developers, literature curators, ...*
- **The Crowd**, verify (biomedical) annotations using the Android app
 - i.e., *Virtually anyone with access to an Android-enabled device*

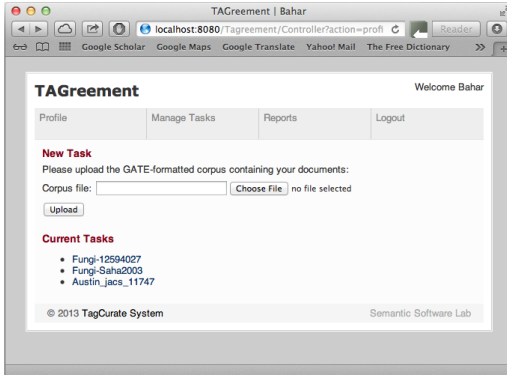
Tagreement Web-based Interface

- *Task Managers* can define and supervise crowdsourcing tasks
- Currently, only accepts GATE-formatted corpora



Tagreement Web-based Interface

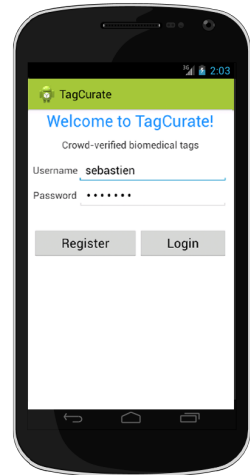
- *Task Managers* can define and supervise crowdsourcing tasks
- Currently, only accepts GATE-formatted corpora
- Stores an internal representation of each tag for distributed verification



| annot_id | task_id | user_id | representation | verified |
|----------|---------|---------|---|----------|
| 90 | 1 | 4 | <annotationInstance> <annotID>90</annotID> <taskID>1</taskID> <content>S. | 0 |
| 91 | 1 | 4 | <annotationInstance> <annotID>90</annotID> <taskID>1</taskID><content>S. mutans</content> | |
| 92 | 1 | 4 | <type>Organism</type> <startOffset>117</startOffset> <endOffset>126</endOffset><features> <feature key='abbrGenus' value='true' /><feature key='docName' value='S. mutans' /><feature key='scientificName' value='Streptococcus mutans' /> <feature key='class' value='Organism' /><feature key='ncbiid' value='1309' /><feature key='Rule' value='2' /><feature | |
| 93 | 1 | 4 | | |
| 94 | 1 | 4 | | |
| 95 | 1 | 4 | | |

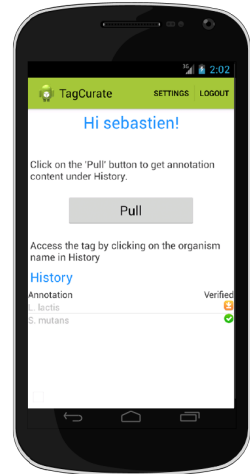
TagCurate Android App

- Developed based on the latest Android distribution (Jelly Bean version 4.3)
- Responsive design for phones and tablets



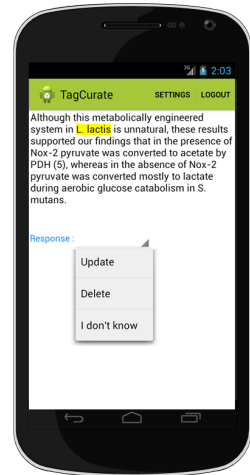
TagCurate Android App

- Developed based on the latest Android distribution (Jelly Bean version 4.3)
- Responsive design for phones and tablets
- Users authenticate themselves on the server
- Users *pull* tags from server
- Temporary storage of verification history



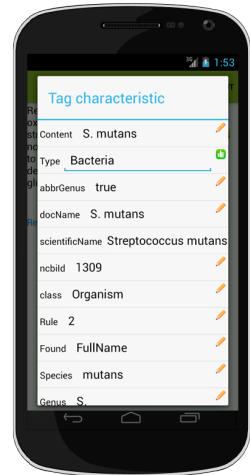
TagCurate Android App

- Developed based on the latest Android distribution (Jelly Bean version 4.3)
- Responsive design for phones and tablets
- Users authenticate themselves on the server
- Users *pull* tags from server
- Temporary storage of verification history
- View tags in context
- Verify whether a tag is a case of:
 - True Positive (correct)
 - False Positive (spurious)



TagCurate Android App

- Developed based on the latest Android distribution (Jelly Bean version 4.3)
- Responsive design for phones and tablets
- Users authenticate themselves on the server
- Users *pull* tags from server
- Temporary storage of verification history
- View tags in context
- Verify whether a tag is a case of:
 - True Positive (correct)
 - False Positive (spurious)
- Modify tags features
 - Pairs of < key, value >
 - Modifications reflect in the tag representation



What about the missing tags?

Manual Annotation

Users select a text span and assign type and features to the generated tag.

Pros

- Human-generated tags usually have a higher quality

Cons

- Difficult task on devices with small screen
- Difficult to achieve an adequate inter-annotator agreement
- Requires well-established annotation guidelines

Automatic Annotation

Users invoke domain-specific text mining pipelines that generate various tags from text.

Pros

- Reuse existing text mining pipelines

Cons

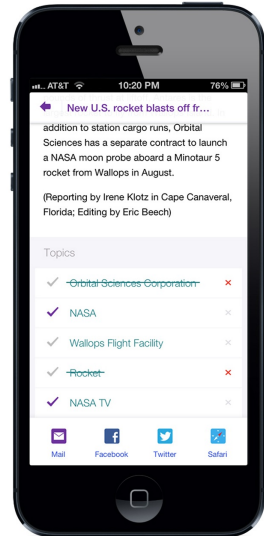
- Text mining techniques are resource-intensive

- 1 Introduction
- 2 TagCurate System
- 3 Android-NLP Integration
 - Mobile Applications of NLP
 - Semantic Assistants Framework
 - Developing NLP Android Apps
- 4 Conclusion

Mobile Applications of NLP

• Automatic Summarization

- Condensed version of document(s)
- Various types: Generic, Focused, Update
- e.g., Summly



(Image Courtesy of Yahoo!)

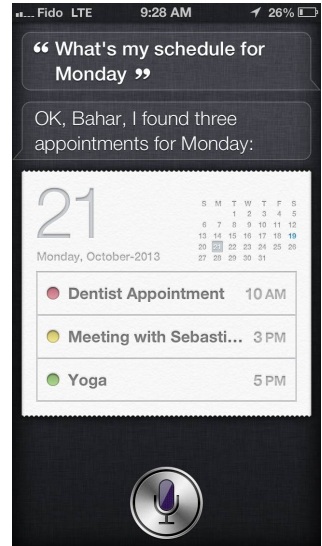
Mobile Applications of NLP

- **Automatic Summarization**

- Condensed version of document(s)
- Various types: Generic, Focused, Update
- e.g., Summly

- **Question Answering**

- Answering *factual* questions
- e.g., Apple's Siri App



Mobile Applications of NLP

- **Automatic Summarization**

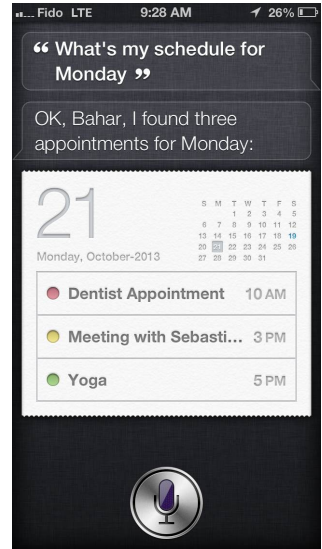
- Condensed version of document(s)
- Various types: Generic, Focused, Update
- e.g., Summly

- **Question Answering**

- Answering *factual* questions
- e.g., Apple's Siri App

- **Information Extraction (IE)**

- Identifying instances of specific classes
e.g., *Persons, Organization, Events, etc.*



Mobile Applications of NLP

• Automatic Summarization

- Condensed version of document(s)
- Various types: Generic, Focused, Update
- e.g., Summly

• Question Answering

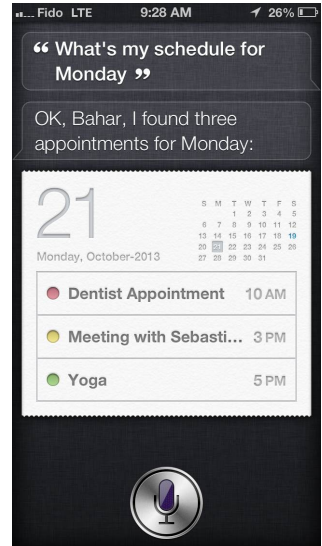
- Answering *factual* questions
- e.g., Apple's Siri App

• Information Extraction (IE)

- Identifying instances of specific classes
e.g., Persons, Organization, Events, etc.

• Content Development

- Combining other NLP services
- Generate new or complementary content



Mobile Applications of NLP

• Automatic Summarization

- Condensed version of document(s)
- Various types: Generic, Focused, Update
- e.g., Summly

• Question Answering

- Answering *factual* questions
- e.g., Apple's Siri App

• Information Extraction (IE)

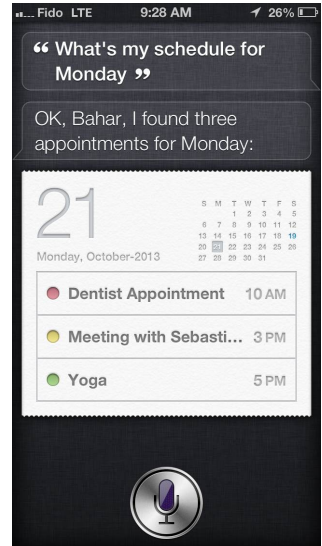
- Identifying instances of specific classes
e.g., Persons, Organization, Events, etc.

• Content Development

- Combining other NLP services
- Generate new or complementary content

• Other domain-specific services

- e-Health, e-Learning, etc.



Mobile Natural Language Processing

What we know

- Numerous mobile applications can benefit from NLP support
- Robust, open-source NLP frameworks are already available
- However, NLP analysis is a very **resource-intensive** task!

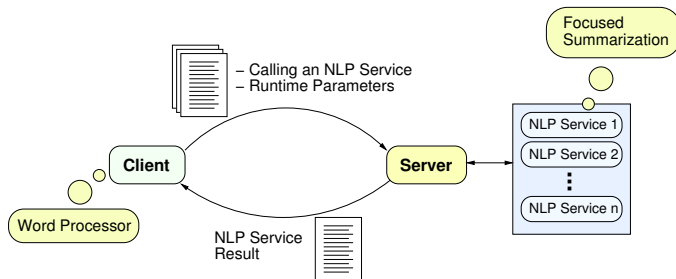
Semantic Assistants Android-NLP Integration

- Novel Android-NLP integration approach
- Provides Separation of Concerns
 - NLP developer does not need to know Android
 - Android app developer does not need to know NLP
- Android *library* for NLP service execution, rather than multiple apps
- Enable users to benefit from complex NLP services in their tasks

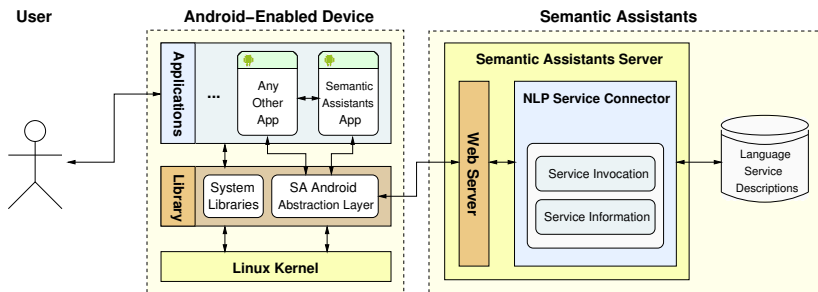
[B. Sateli, G. Cook, R. Witte, “*Smarter Mobile Apps through Integrated Natural Language Processing Services*”, MobiWIS 2013]

Semantic Assistants Framework

- Existing open-source (AGPL3) service-oriented architecture
- Brokers NLP pipelines as standard W3C Web services
- Avoids context-switching of user to external text mining applications
- Brings NLP analysis directly to various applications via plug-ins



Semantic Assistants NLP Intents



- Client-Server Model
 - Client is an Android app
 - Server-side component is the Semantic Assistants server
- RESTful communication over HTTP(S)
- Handles various NLP service result formats
 - *Annotation*, e.g., a person name in text
 - *Document*, e.g., summary of a long webpage
 - *Files*, e.g., an HTML document

Developing NLP Android Apps

Separation of Concerns

Android Developer

- Identify the NLP task
- Extend the SA intents by choosing a unique package name for this new service
- Embed the SA Android library in a new Android app
- Invoke the intent in app using the library

NLP Developer

- Develop the concrete NLP pipeline
- Deploy the pipeline on a SA server

Summary and Outlook

Summary

- Distribute annotation jobs to large user groups
- Expert annotators can focus on quality control and difficult cases
- Easily bring NLP pipelines to (Android) mobile apps

Ongoing work

- TagCurate app facelift
- Expanding the user profiles
- Finding incentives and introducing social aspects
- Add annotation capabilities (both manual and semi-automatic)

Find out more...

- Twitter: @SemSoft
- Web: <http://www.semanticsoftware.info/>