

Text Analytics in the Wild

Trends and Technologies

Mark Cieliebak

20.11.2013

ZHAW

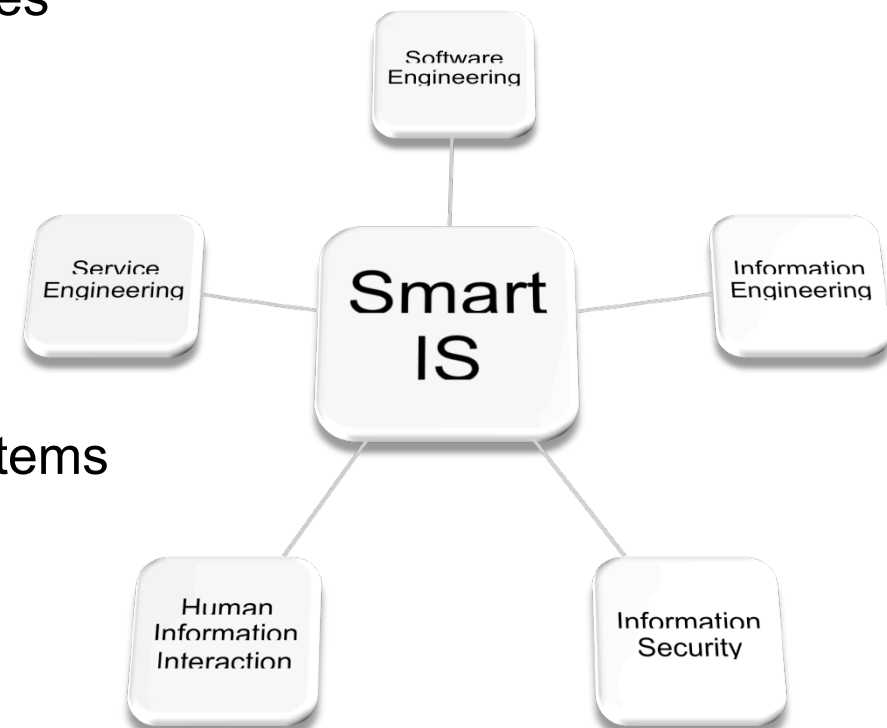
Zurich University of Applied Sciences

- more than 10'000 students
- ~230 Professors (FTE)

InIT

Institute of Applied Information Systems

- since 2005
- ~35 professors
- Labs on Cloud Computing, Data Science, Visual Computing etc.

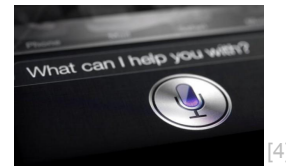




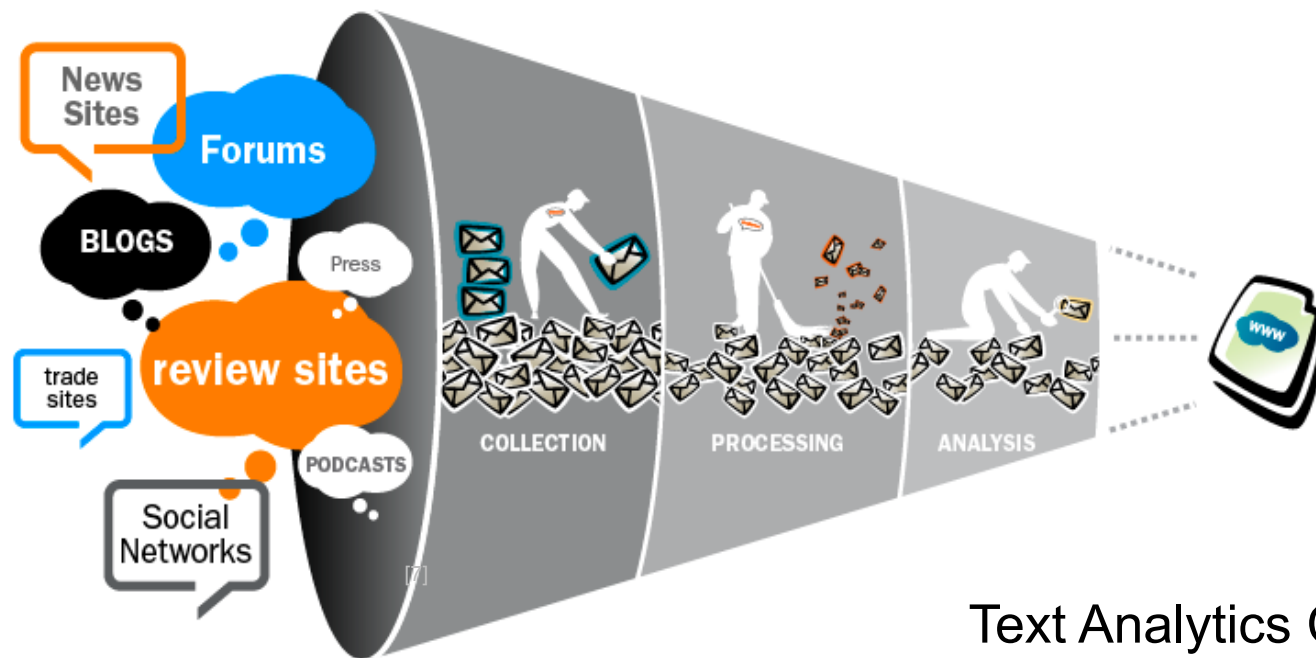
Executive Summary



[5]



Sample Application: Social Media Monitoring



Text Analytics Components:

- Find relevant documents
- Hot topic detection
- Sentiment detection

"Netgear's WiFi Analytics is a free Android app that I find very handy when it comes to troubleshooting and monitoring a home network.

It's simple, but that simplicity is its strength and gives home users insight into a home wireless network, without need to know a lot about networking." [8]

Simple Sentiment Detection

Idea: Count number of positive and negative words

"This analysis is good [+1]."	1 (pos)
"I find it beautiful [+1] and good [+1]."	2 (pos)
"It looks terrible [-1]."	-1 (neg)
"This car has a blue color."	0 (neu)

Use Sentiment-Dictionary:

POSITIVE:	NEUTRAL:	NEGATIVE:
good	hello	bad
love	see	hate
nice	I	ugly
...

Improvements

Sample Text

- "This analysis is **good**."
"This analysis is **excellent**."
- "The car **really very expensive**."
- "This car has an **appealing** design and **comfortable** seats, but it is **expensive**."

Solution

- *Fine-Grained Dictionaries:*
"This analysis is **good** [+2]. "
"This analysis is **excellent** [+3]."
- *Detect Booster Words:*
"This car **really very expensive**[-1 -1 -2] ."
- *New Category "Mixed":*
"This car has an **appealing**[+1] design and **comfortable**[+1] seats, but it is **expensive**[-1]. "
→ MIX(+2/-1)

Improvements 2

Sample Text

- This analysis is **not** **good**.
- The car is **appealing** and I do **not** find it **expensive**.
- I do **not** find the car **expensive** and it is **appealing**.

Solution

- *Invert all Scores:*
"This analysis is **not**^[*-1] **good**^[+3]."
 - *Invert only score of words occurring after the negation:*
"The car is **appealing**^[+3] and I do **not**^[*-1] find it **expensive**^[-2]"
- Need to "understand" the sentence
→ Need linguistic analysis!

Sentence										
Sentence							Conj.	Sentence		
Noun Phrase	Verb Phrase							Noun Phrase	Verb Phrase	
	Verb	Adverb	Verb	Noun Phrase		Adj.			Det.	Verb
Det.				Det	Noun					
I	do	not	find	the	car	expensive	and	it	is	appealing

-> Invert scores of words being in the same phrases as negation:

“I do not find the car **expensive**[+2]
 and it is **appealing**[+3].” → +5 (pos)

Rule-based Sentiment Detection

"This movie was **like** a **horror** event."

Noun Phrase		Verb Phrase				
Determiner	Noun	Verb	Prepositional Phrase			
			Preposition	Noun Phrase		
				Determiner	Noun	Noun
This	movie	was	like	a	horror	event

Sentiment Dictionary: *like* (Verb): +2, *like* (PP): 0

⇒ "This movie was like[0] a **horror**[-2] event." → -2 (**neg**)

Technology Stack

for Sentiment Detection on Web Documents using Linguistic Analysis

Linguistic Framework

- Gate: text analysis framework with experimentation GUI; includes sentence splitter, tokenizer, chunker, pos-tagger, etc; plugin-mechanism for new libraries; very widely used. <http://gate.ac.uk/>
- OpenNLP: Java-based, easy to use; includes tokenization, sentence segmentation, part-of-speech tagging, named entity extraction, etc; by Apache Foundation. <http://opennlp.apache.org/>

Web Boilerplate Extraction

- Boilerpipe: removes ads, menus, HTML-tags etc. <http://boilerpipe-web.appspot.com/>

Language Detection

- Nutch language identifier plugin: 14 languages preimplemented, „easy“ to add more. <http://wiki.apache.org/nutch/LanguagelIdentifierPlugin>

Sentence Splitting (various programming languages)

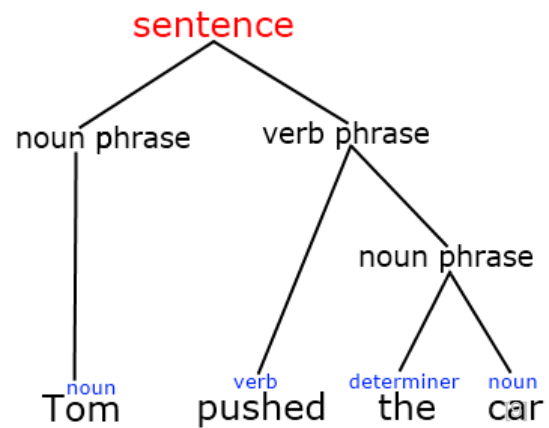
- OpenNLP S 2.2 <http://opennlp.apache.org/>
- CoreNLP R 18.7 <http://nlp.stanford.edu/software/corenlp.shtml>
- GATE R 60.6 <http://gate.ac.uk>
- LingPipe R 1.7 <http://alias-i.com/lingpipe/>
- MxTerminator S 2.8 <ftp://ftp.cis.upenn.edu/pub/adwait/jmx/>
- Punkt U 7.1 <http://nltk.org/api/nltk.tokenize.html>
- RASP R 0.3 <http://ilexir.co.uk/applications/rasp/>
- Splitta S 31.5 <http://code.google.com/p/splitta>
- tokenizer R 0.3 <http://www.cis.uni-muenchen.de/~wastl/misc/>
- opennlp ist java und easy to use.

Part-of-Speech Tagger

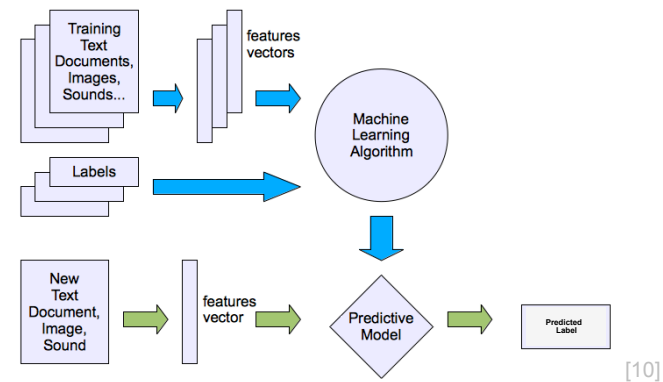
- TreeTagger: language-independent, several languages already implemented; also for lemmatization. <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>
- TweetNLP: POS-tagger just for Twitter. <http://www.ark.cs.cmu.edu/TweetNLP/>

Approaches to Sentiment Detection

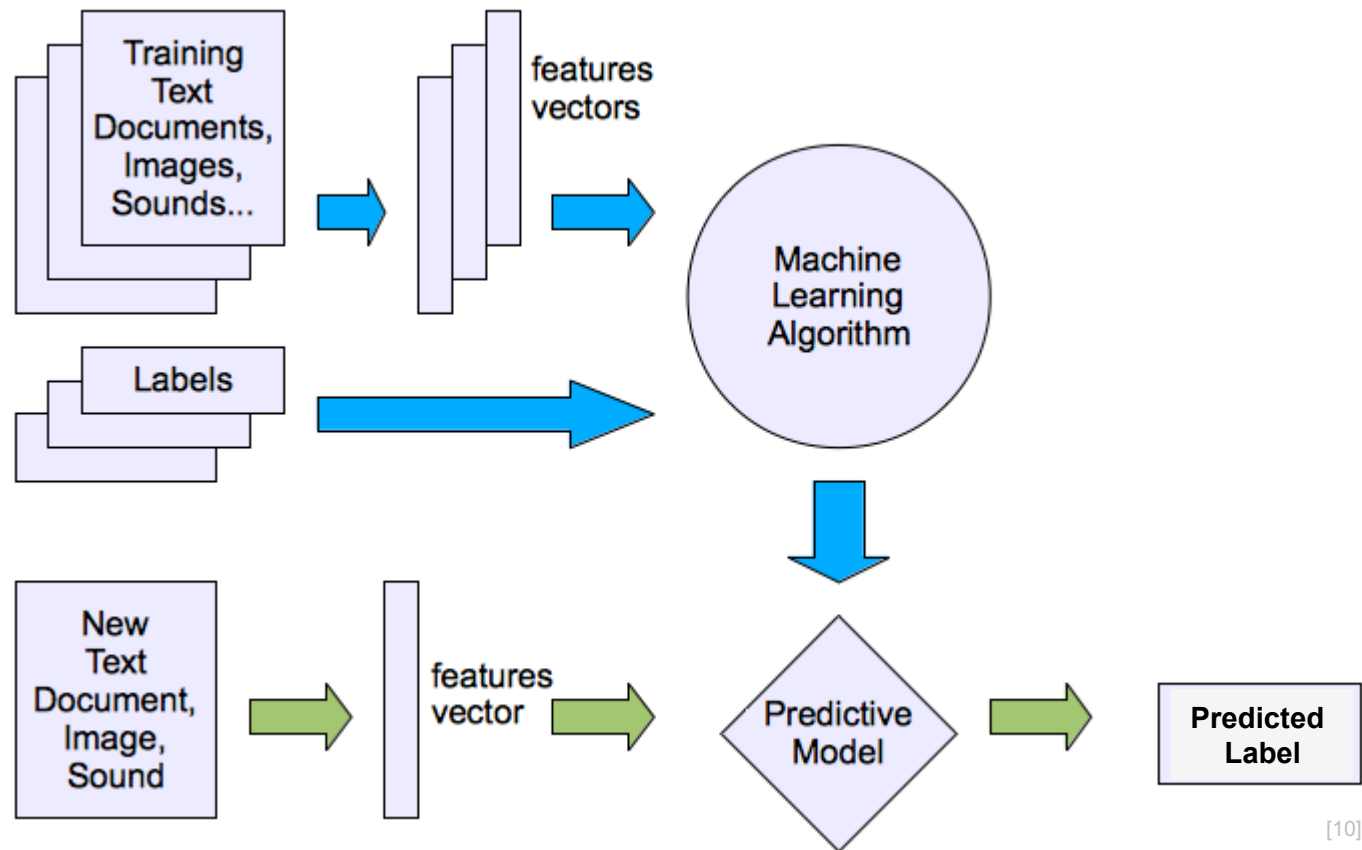
Rule-Based



Corpus-Based



Corpus-Based Sentiment Detection



Corpus-Based Sentiment Detection

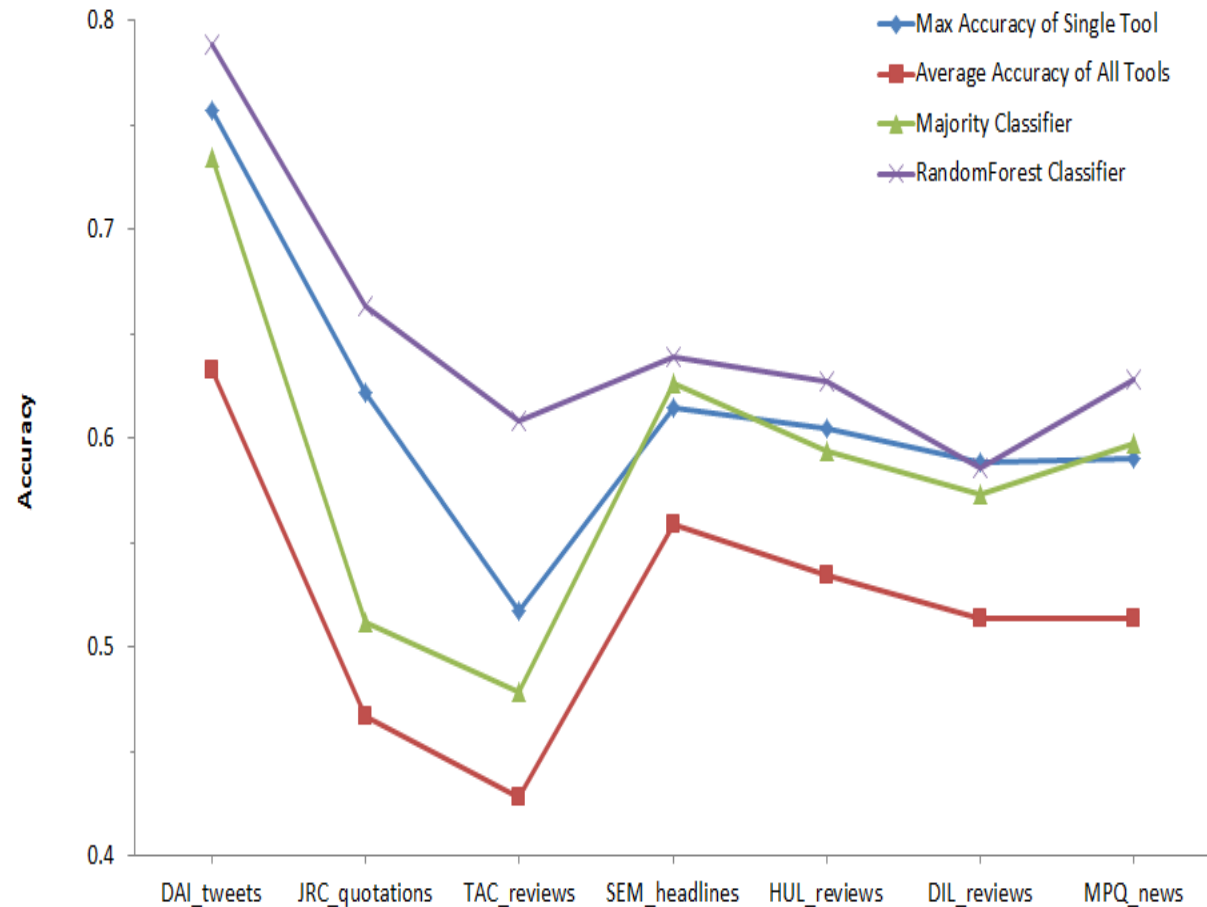
Annotated Corpus

Sentence	Polarity
This analysis is good.	Pos
It looks awful.	Neg
This car has a blue color.	Neu
This car has an appealing design, comfortable seats, but it is expensive.	Mix
This car has a very appealing design, comfortable seats, but it is really expensive.	Mix
This analysis is not good.	Neg
This car has an appealing design, comfortable seats and it is not expensive.	Mix
This movie was like a horror event.	Neg
This car is appealing and is not expensive.	Mix
...	...

Challenges in Sentiment Detection

- **Sarcastic Statements**
“What a great car, it stopped working in the second day.”
- **Conditional Statements**
“If I can find a good camera, I will buy it.”
“If you are looking for a good phone, buy Nokia”
- **Context Dependencies**
"Salaries for software engineers are extremely high."
- **User Expectations vs. Reality**
one "sentScore" vs. ambivalent docs
- **Opinion Spam Detection**
- **Huge Effort for New Languages**

Sentiment Detection: State of the Art



Evaluation of 7 commercial tools on 30'000 documents [11]

Main Results:

- Best tool has 60% accuracy on all documents
- "Easiest" document type is Tweets: up to 78% accuracy
- Meta-classifier beats each single tool

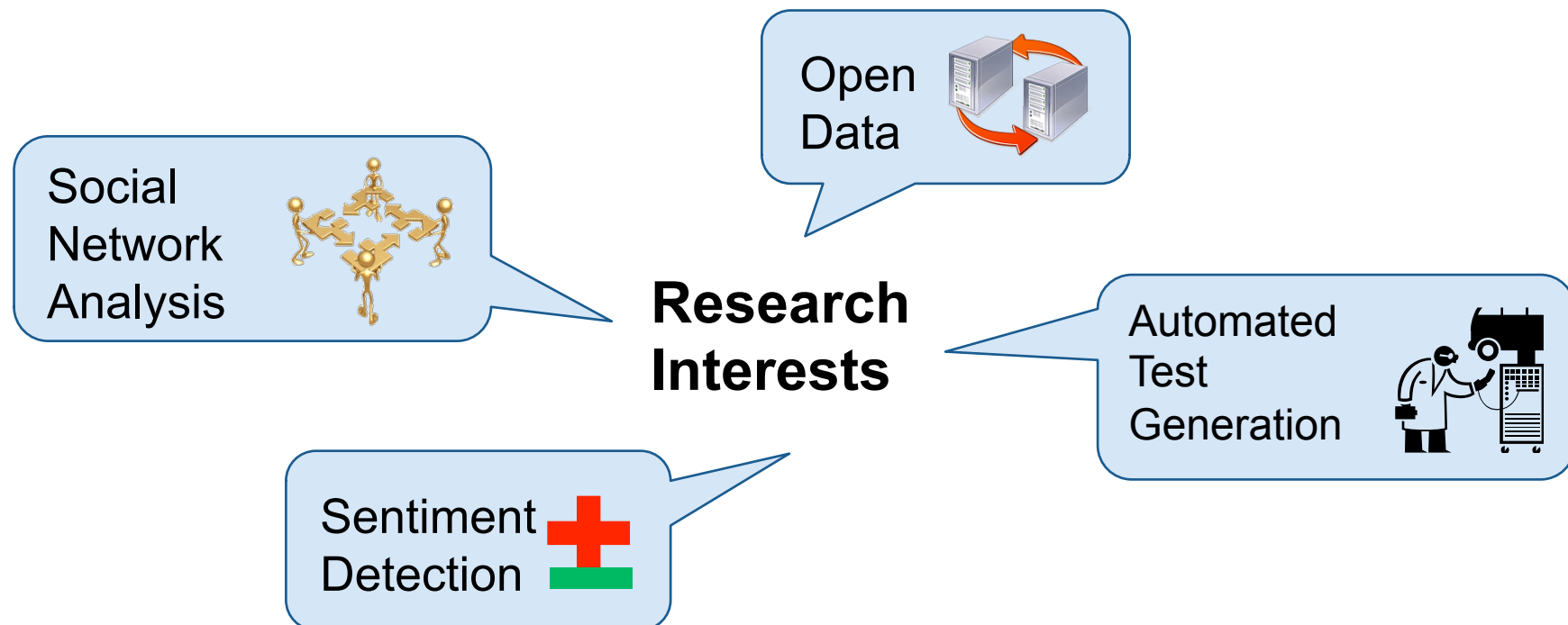


Mark Cieliebak

Institute of Applied Information Technology (InIT)

ZHAW, Winterthur

Email: ciel@zhaw.ch, Website: www.zhaw.ch/~ciel



References

- [1] http://www.teleportmyjob.com/blog/wp-content/uploads/2013/05/find_job.jpg
- [2] <http://screenshots.de.sftcdn.net/de/scrn/3340000/3340669/swiftkey-23-630x535.png>
- [3] https://lh4.ggpht.com/Fs4HLaMCgmY6O7BWGKZB8LdOM8GmiRAbKUcdJysNp0k74xmPMTSGVHigwvaLy_Tw=w300
- [4] <http://revaxmedia.co.uk/wp-content/uploads/2012/09/Siri-What-Can-I-Help-You-With.jpg>
- [5] <http://www.betadaily.com/wp-content/uploads/2011/02/BusinessPlanExecutiveSummary.jpg>
- [6] <http://everycook.org/cms/en/hardware-en/pictures-en#>
- [7] <http://cdn.fulltraffic.net/images/blog/social-media-monitoring-process.png>
- [8] <http://www.pcmag.com/article2/0,2817,2426416,00.asp>
- [9] <http://i1059.photobucket.com/albums/t424/shmi11y/UoS%20LingSite/tompushedthecar.png>
- [10] <http://bigsnarf.files.wordpress.com/2013/04/supervised.png>
- [11] Cieliebak et al., ESSEM 2013.

"McAfee Antivirus & Security for Android ... is free for All Access users.

Even if you don't have All Access, McAfee is certainly attractive with its well-known name, ... and enormous slew of security features.

But I was disappointed that the app didn't live up to its incredible potential, and would recommend looking at lower-priced alternatives if you're not an All Access user."