

Text Analysis for Social Media Cybersecurity: the AMiCA Project

Els Lefever

Language and Translation Technology Team (LT³)
Ghent University, Belgium



language and
translation
technology
team



LT³, LANGUAGE AND TRANSLATION TECHNOLOGY TEAM



- Dpt of Translation, Interpreting and Communication, Faculty of Arts and Philosophy, Ghent University
- fundamental and applied research in **language and translation technology**
- expertise in using **machine learning** for language technology problems (PoS-tagging and lemmatization, anaphora resolution, WSD, NER)
- Headed by Prof. Véronique Hoste





3 main research lines:

- Terminology & computational semantics
- Translation Technology
- Sentiment analysis and subjectivity detection



Terminology / computational semantics



- Lead: Prof. Els Lefever
- Automatic terminology extraction from monolingual, bilingual and comparable corpora (Ayla Rigouts Terryn)
- Automatic hypernym and synonym detection (Els Lefever)
- Term ambiguity in interdisciplinary research (Julie Mennes)
- Use of term extraction for translating documentaries (Sabien Hanouille)

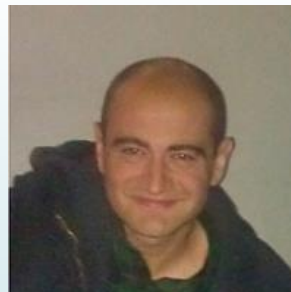




Translation Technology



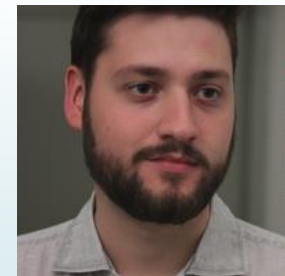
- Lead: Prof. Lieve Macken
- comparison of different methods of translation:
human vs. post-editing, human vs. CAT (Joke Daems)
- translation quality assessment and confidence
estimation for machine translation (Arda Tezcan)





Sentiment Analysis and Subjectivity detection

- Lead: Prof. Véronique Hoste
- automatic detection of cyberbullying (Cynthia Van Hee)
- suicide detection (Bart Desmet)
- Aspect-based sentiment Analysis (Orphée De Clercq)
- detection of subjectivity in annual reports (Nils Smeuninx)
- Irony detection (Cynthia Van Hee)
- Sentiment Analysis for economic events (Gilles Jacobs)





AMICA



Outline

- The context and goals of the AMiCA project
- Text normalization
- 3 Use cases:
 1. Detecting cyberbullying
 2. Suicide detection
 3. Age and gender profiling for detecting grooming



- IWT-SBO project, coordinated by CLiPS (UA)
- Partners:
 - CLiPS (text mining, UA)
 - MIOS (sociology, UA)
 - LT3 (text mining, UGent)
 - IBCN (software development, UGent)
 - VISICS (image processing, KUL)
- Combine text analytics, image and video analysis, and data mining





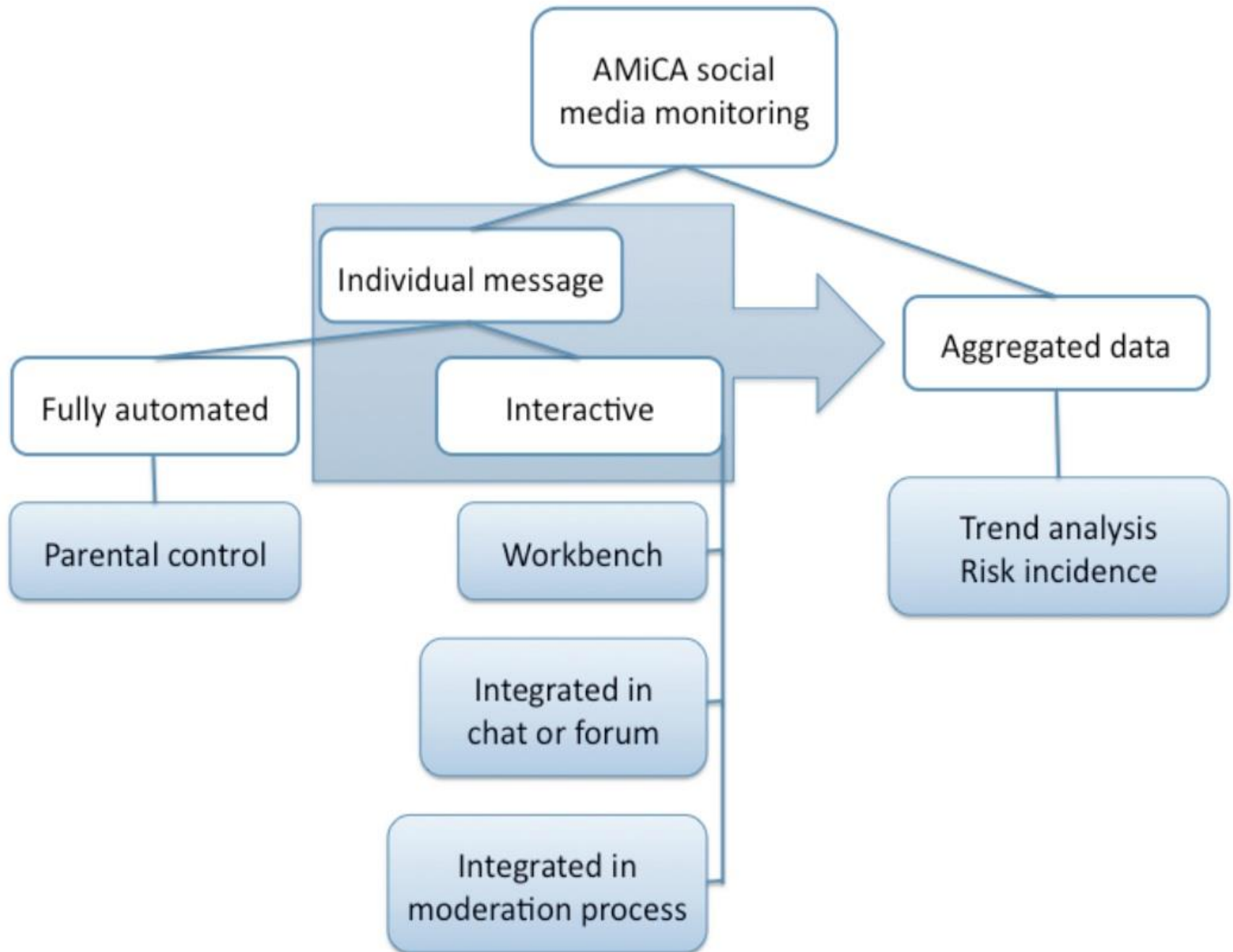
Goals

- Detect situations that are harmful or threatening to young people in social networks
 - Cyberbullying
 - Sexually transgressive behaviour (for example grooming by paedophiles)
 - Depression and suicide announcement
- Facilitate efficient action by moderators, police, parents, peer group, social services, ...
- Objective measurement, monitoring, trend analysis, ...



User Committee





How urgent is the problem?



- European “Kids online” study (EU, 2011)
 - Motivation for the project
 - Age 9-16 in 25 European countries
 - Results
 - Children are 90 minutes per day online
 - Half of them in their bedroom
 - 33% added strangers as friends
 - 15% shared personal information with strangers (Including photographs)
 - 12% felt they experienced harm
- www.eukidsonline.net



How urgent is the problem?

- European “Kids online” study: **update in 2014**
 - Age 9-16 in 25 European countries
 - Results since 2010 study, 9 to 16 year olds
 - Significant rise of use of social media
 - Rise of 23% to 43% of having contact with someone not met IRL before
 - Rise of 10% to 23% of having seen sexual images
 - Rise of 9% to 20% of having received sexual images
 - Rise of 13% to 17% are upset by something seen online
 - Rise of 13% to 20% of being exposed to hate messages
 - Rise of 7% to 11% of being exposed to self-harm sites
 - Rise of 7% to 12% of being exposed to cyberbullying

www.eukidsonline.net



Quick poll

- Who is in favor of software monitoring automatically your interactions in social media for risks and threats?

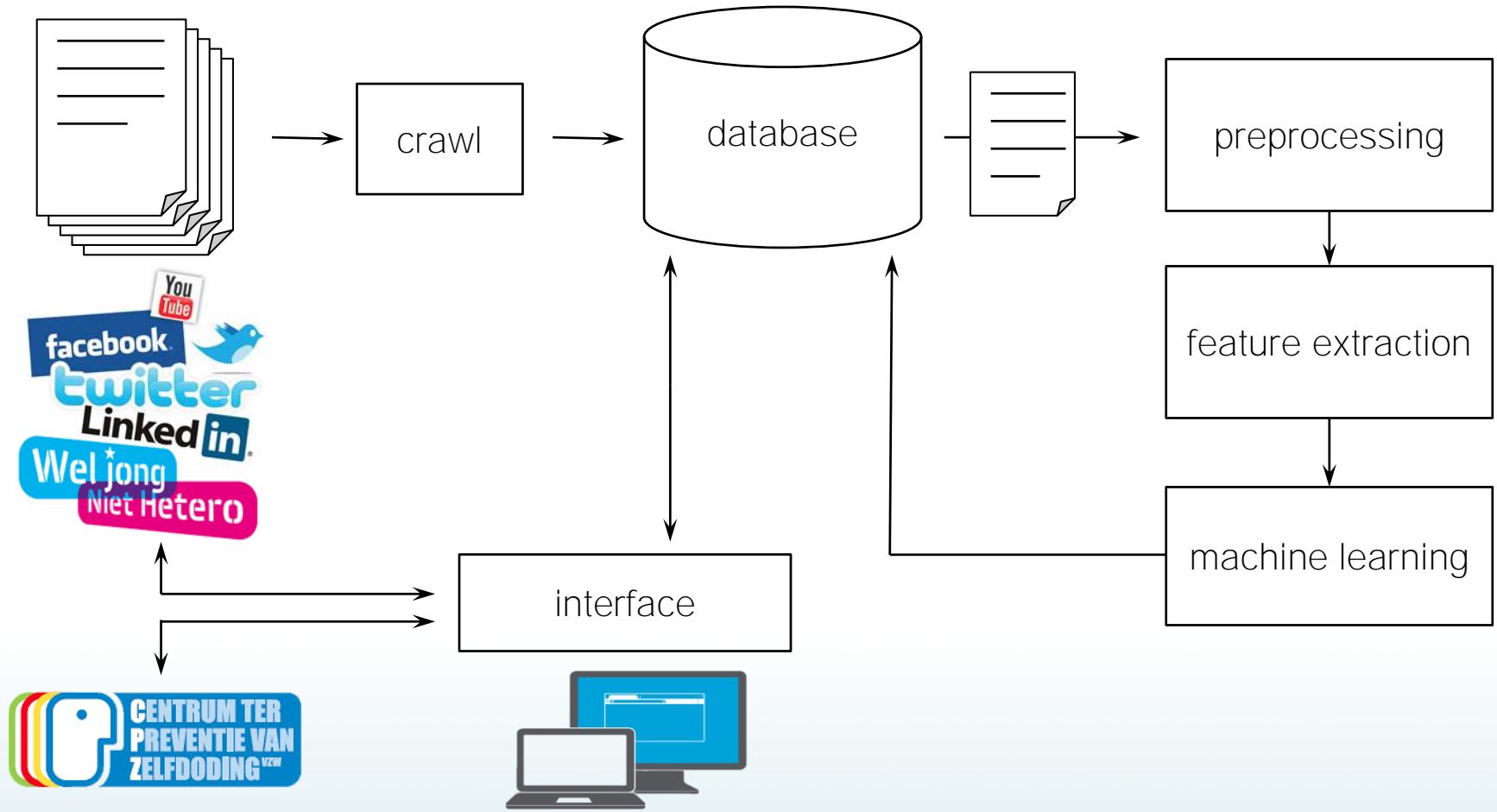
Should we do something about it?



- Majority of experts and adolescents is in favor of automatic monitoring
 - but only for situations they perceive as uncontrollable
 - with respect for privacy and with suitable follow-up, not involving too many parties, and giving control to the victim
- Mixed opinions with the parents depending on (negative) previous experience and level of trust in their children



Workflow





Crawl: example



Zwijg stomme trut! Gij
hebt geen leven tot op je
begravenis!!!

(English: Shut up stupid cow! You don't
have a life see you at your funeral!!!)



Crawl: example

Django administration

Welcome, **nlpapp** ▾

Recent Actions ▾

[Home](#) / [Nlp](#) / [Tweets](#) / 2015-09-29 22:52:47+00:00: Zwijg stomme trut! Gij hebt geen leven tot op je b...

Change tweet

History

Fields in **bold** are required.

Tweet url:

Timestamp:
Date: Today |
Time: Now |
Note: You are 2 hours ahead of server time.

Text:

User name:



PREPROCESSING / NORMALISATION OF USER-GENERATED TEXT



User Generated Content

Social media: blogs and microblogs (Twitter: 190 million tweets/day), wikis, podcasts, social networks (Facebook: 70 billion shares/month)

⇒ Enormous amount of UGC





UGC Normalization

Maxims of chat language:

- Write **as fast as you can** (fluent interaction)
 - » Abbreviations, letter omission, acronyms, flooding, concatenation, capitalization, punctuation, spelling and grammar errors, ...
- Write **as you speak** (informal character of the conversation)
 - » Dialectical, phonetic, emoticons, ...



Properties of chat language

- Omission of words / characters (spoke – spoken)
- Abbreviations, acronyms (LOL – laughing out loud)
- Deviations from standard spelling (luv – love, you iz – you are)
- Expression of emoticons:
 - Flooding (looooooooooove)
 - Emoticons (:p)
 - Capitalized letters (STUPID)
- Dutch-specific:
 - Concatenation of tokens (khou – ik hou)
 - Elimination of clitics and pronouns (edde – heb je)
 - Lot of dialects!



Example

	Example of Dutch SMS language
Original	Oguz ! Edde me Jana gesproke ? En ze flipt lyk omdak ghsmoord heb .. !
Normalized	Oh gods ! Heb je met Jana gesproken ? En ze flipt gelijk omdat ik gesmoord heb .. !
Translated	Oh god ! Did you speak to Jana ? And she's flipping because I smoked ... !



Problem for Text Analysis Tools

- Most NLP tools are developed for or trained on standard language
- They fail miserably on UGC
- Solutions
 - Develop new tools
 - E.g. Tweet NLP (CMU):
<http://www.cs.cmu.edu/~ark/TweetNLP/>
 - Normalize the ‘non-standard’ language
- On the positive side, non-standard language makes some analytics tasks easier!

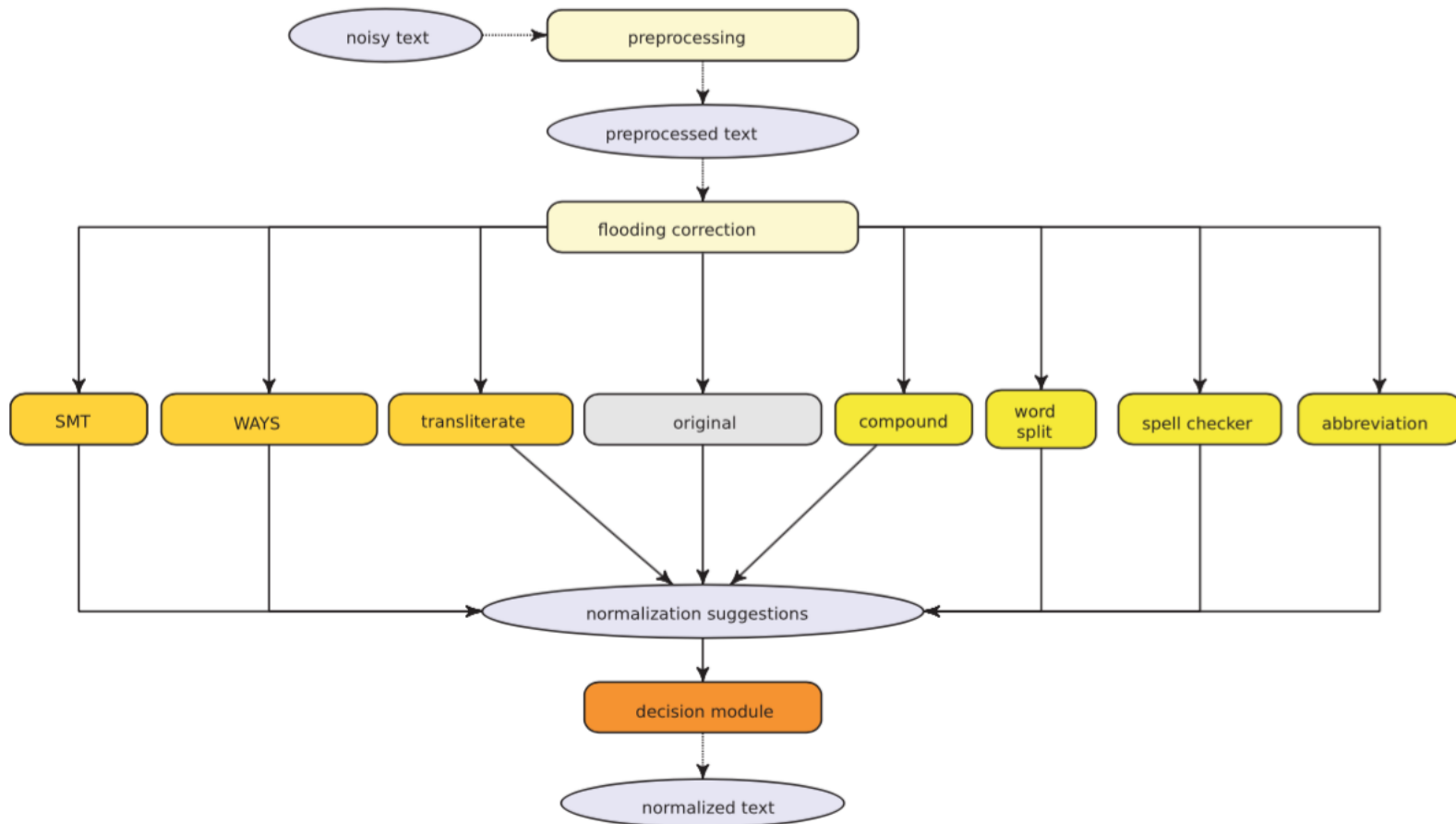


Normalization Approaches

- Three dominant approaches
 - Machine Translation: Source Language = non-standard and Target Language = standard
 - Spell Checking: Correct the incorrect words (statistical or dictionary-based)
 - Speech Recognition: Non-standard language = speech that has to be converted to text (HMMs)
- => We choose to follow an SMT approach and also go to the character-level



Ensemble Approach



Sarah Schulz, Guy De Pauw, Orphée De Clercq, Bart Desmet, Véronique Hoste, Walter Daelemans, and Lieve Macken. 2016. Multimodular text normalization of Dutch user-generated content. *ACM Trans. Intell. Syst. Technol.* 7, 4, (July 2016), 22 pages. DOI: <http://dx.doi.org/10.1145/2850422>



Modules

- Preprocessing
 - Tokenization and sentence splitting
 - includes emoticons, emojis etc.
 - Character flooooooooding
- Token-based modules
 - Abbreviations
 - Expansion dictionary (~ 350 abbrevs)
 - Spell checker
 - Levenshtein on dictionary (~ 2.3 million words)
 - Compound Module
 - Checks if a pair of words is actually one word
 - Word Splitter
 - 'misje' = 'mis je' (miss you)



Modules

- Context-based modules
 - SMT
 - Token-unigram, character unigram, character-bigram and combinations
 - Transliteration (supervised ML)
 - supervised ML, memory-based learning style
 - +da+_n i ++_ged -> iet
 - WAYS (Write As You Speak): G2P + P2G (memory-based learning)
 - ni (niet, *not*)
 - kem (ik heb, *I have*)
- “Original” Module
 - Many words are correct



Modules

- **Decision Module**
 - Moses decoder (SMT), dynamic search among the suggestions of the component modules
 - Uses (5-gram) language model and phrase table (dev. Set)



Evaluation

- Three types of UGC
 - Chat (Netlog)
 - SMS (Sonar corpus)
 - Microblog (Twitter)
- Train (60%) - Development (20%) - Test (20%)
- Total: 70,000 tokens, manually annotated
 - insertions, deletions, substitutions, transpositions
 - near-perfect annotator agreement
- Background corpora for language modeling

CGN (Spoken Dutch Corpus)	6,765,336
SoNaR (Balanced text corpus)	3,581,182
Open Subtitles Dutch (OSD)	90,147,315
Training set (TS)	56,523



Results

- Module level evaluation:
 - SMT and Transliterate modules perform best
 - Especially compounding and splitting problems remain
- Ensemble evaluation:
 - Best ensemble system: 92.9
- Extrinsic and Portability Evaluation
 - Tested on Ask FM for NLP tasks (with and without normalizing)
 - POS (+12%), LEM (+13%), NER (+8%)
- Problems remain especially in tokens with multiple normalization problems



USE CASE 1: CYBERBULLYING DETECTION



Research Motivation

- $\pm 20\text{-}40\%$ of all youth have been victimized online (Tokunaga, 2010)
- **Anonymity, lack of supervision and impact** make social media a convenient way for cyberbullies to target their victim (Hinduja & Patchin, 2006)
- Information overload on the Web has made **manual monitoring unfeasible**

more likely to be exposed to hate messages

13%	to
20%	

more likely to be exposed to pro-anorexia sites

9%	to
13%	

more likely to be exposed to self-harm sites

7%	to
11%	

more likely to be exposed to cyberbullying

7%	to
12%	

13%	to
17%	

European 9- to 16-year-olds say they are now: more likely to say they were **upset** by something seen online in 2014



Research Motivation

- Automatic detection systems allow for large-scale **social media monitoring**
- Goal => **reduce manual monitoring efforts** on social media



Related Research

- NLP applications for **automatic** cyberbullying prevention and detection
 - Cyberbullying detection (Yin et al., 2009; Reynolds et al., 2011; Nahar et al., 2013)
 - Sensitive topic identification (sexuality, race) (Dinakar et al., 2012)
 - Detection of bully profiles on social networks (Dadvar et al., 2013)

BUT:

- Focus on **posts from harassers**
- No distinction between different **types of cyberbullying**
- Datasets do not always follow a real-world **distribution**



Data set construction

- We need large data sets to train machine learning systems
- Data collection for Dutch and English

- Data from relevant social media
- BUT: few / private data



- Media campaign for donating examples of cyberbullying messages
- BUT: sensitive data!



- Cyberbullying simulations





Data set construction: media campaign



AMICA

Wat is AMICA?

Het AMICA-project (Automatic Monitoring for Cyberspace Applications) wil een mogelijk bedreigende situatie op sociale

netwerken automatisch herkennen door middel van taal- en beeldanalyse om zo de online veiligheid van kinderen te waarborgen.

In dit project wordt technologie ontwikkeld om deze situaties te herkennen en wordt nagegaan hoe de ontwikkelde technologie het best kan worden ingezet om slachtoffers van cyberpest te helpen.

Meer informatie over dit project vindt u op de website: <http://www.amicaproject.be>

Gezocht: cyberpestberichten

Beste,

In het kader van ons onderzoek naar cyberpest aan de universiteiten van Antwerpen, Leuven en Gent hebben we een oproep gelanceerd om zoveel mogelijk cyberpestberichten te verzamelen. De berichten willen we gebruiken als voorbeeldmateriaal om een systeem te ontwikkelen dat gevaarlijke situaties op sociale netwerken automatisch kan herkennen.

We willen onze oproep in het bijzonder richten naar jongeren die met cyberpest worden geconfronteerd worden. Graag vragen we aan u als hulpverleningsinstantie of school om onze vraag door te geven aan de jongeren waarmee u in contact komt. Concreet willen we u vragen dat u kort ons onderzoek schetst en vraagt of ze bereid zijn te helpen. Daarna kunt u de jongeren een brief meegeven waarin onze doelstellingen staan en waarin ze praktische richtlijnen vinden om de gezocht berichten op te sturen.

Onze dataverzameling loopt vanaf heden tot en met juni 2014, gedurende dit half jaar zijn alle initiatieven om berichten te verzamelen welkom. Elke bijdrage aan ons onderzoek is immers waardevol en kan ons helpen om een beter systeem te ontwikkelen waar we jongeren in de toekomst mogelijk mee kunnen helpen.

Wij bedanken u alvast voor uw hulp!
Het AMICA-onderzoeksteam

Voor meer informatie over deze oproep of ons project kan u steeds terecht bij:
Cynthia Van Hee Universiteit Gent cynthia.vanhee@ugent.be
Ben Verhoeven Universiteit Antwerpen ben.verhoeven@uantwerpen.be



altijd dicht bij jou

De Universiteit van Gent zoekt cyberpestberichten

Tags: cyberpesten, gent, universiteit

Een aantal onderzoeksgroepen van de universiteiten van Antwerpen, Gent en Leuven zijn een systeem aan het ontwikkelen dat automatisch cyberpestgedrag, en andere gevaarlijke situaties op het internet, kan herkennen op sociale netwerksites. Op die manier proberen zij een veilige internetomgeving te creëren voor jongeren.

De wetenschappers die hier mee bezig zijn, lanceren nu een oproep aan slachtoffers, en getuigen van cyberpesten, ouders en leerkrachten, om zoveel mogelijk berichten door te geven die zij als cyberpesten beschouwen. Dat kunnen e-mails zijn, sms'en, chatgesprekken, of berichten van Facebook, netlog, Twitter en Ask.fm. Alles wordt anoniem behandeld, en wordt alleen gebruikt voor intern onderzoek. Dus maak je geen zorgen, er wordt niks doorgegeven aan derden!

Dus als je dit wil, en kan doen, surf dan snel naar deze website: www.amicaproject.be of stuur meteen een e-mail: data@amicaproject.be. Uw e-mail wordt vertrouwelijk behandeld door het AMICA-onderzoeksteam.

26 nov 2013 - 16:18

Nieuw systeem spoort online

Beeldscherm van een computer met een bericht op Facebook.

Belgische onderzoekers willen cyberpesten tegenhouden door poging op Facebook en Twitter automatisch op te sporen. Het systeem moet op termijn ook ingezet worden bij zelfmoordgedachten.

Onderzoekers in samenwerking met het Vlaamse Instituut voor de Media (VIM) ontwikkelen een systeem dat automatisch berichten op sociale netwerken kan analyseren en zo potentiële slachtoffers kan helpen.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

pestgedrag automatisch op

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.



GEZOCHT: CYBERPESTBERICHTEN

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

Beeldscherm van een computer met een bericht op Facebook.

GEZOCHT: CYBERPESTBERICHTEN
Veilig internetten is belangrijk, maar niet altijd even gemakkelijk. Heb jij al een keer rare of onveilige dingen meegemaakt op het internet?

Verscheidene studies geven aan dat één Vlaamse jongere op tien recent het slachtoffer was van cyberpesten. De impact van dergelijke situaties is vaak erg groot doordat kwetsbare berichten soms langere tijd online blijven staan en slachtoffers alleen met dit probleem zitten omdat ze er niet over willen of durven te praten.

SAMEN STERK

In het wetenschappelijk project AMICA slaan een aantal onderzoeksgroepen van de universiteiten van Antwerpen, Leuven en Gent de handen in elkaar om een systeem te ontwikkelen dat automatisch kritische situaties zoals cyberpestgedrag en seksueel grensoverschrijdend gedrag herkent op sociale netwerksites om zo een veilige internetomgeving te kunnen garanderen voor jongeren.

BLIJF NIET BIJ DE PAKKEN ZITTEN

Om te kunnen bepalen wat cyberpesten is en hoe het kan worden herkend, zijn voldoende onderzoekers nodig van deze online berichten. Daarom doet AMICA een oproep aan u als ouder om zoveel mogelijk berichten aan ons beschikbaar te stellen waarin cyberpesten merkbaar is.

Als u kind getuige of slachtoffer is geweest van cyberpesten en als hiervan bewijsmateriaal beschikbaar is, dan zijn deze berichten meer dan welkom. De berichten kunnen bijvoorbeeld e-mails zijn, sms'jes, chatgesprekken of berichten van sociale media. De data worden anoniem behandeld en alleen gebruikt voor intern onderzoek, berichten worden in geen geval doorgegeven aan derden.

Wij in contact opnemen met AMICA (Automatic Monitoring for Cyberspace Applications)? Stuur dan een mailje naar data@amicaproject.be en blijf niet bij de pakken zitten. Meer informatie vind je op de website: www.amicaproject.be.



Data set construction: media campaign



Gezocht: cyberpestberichten

Beste,

In het kader van ons onderzoek naar cyberpesten aan de universiteiten van Antwerpen, Leuven en Gent hebben we een oproep gelanceerd om zoveel mogelijk



De Universiteit van Gent zoekt cyberpestberichten

Tags: cyberpesten, gent, universiteit

Een aantal onderzoeksgroepen van de universiteiten van Antwerpen, Gent en Leuven zijn een systeem aan het



RESULT: ± 30 reactions
± 368 messages (FB messages, hate pages, Netlog, mail, chat, etc.)



vindt u op de website:
<http://www.amicaproject.be>

Voor meer informatie over deze oproep of ons project kan u steeds terecht bij:
Cynthia Van Hee Universiteit Gent cynthia.vanhee@ugent.be
Ben Verhoeven Universiteit Antwerpen ben.verhoeven@uantwerpen.be



Nieuw systeem spoort online

door
Bart De Weert

**Belgische onderzoekers
vallen cyberpesten
op Facebook en Twitter
aan met een nieuw
systeem dat op zoek
gaat naar berichten
die op een manier
worden verspreid die
niet normaal is.**

Opmerkelijk is dat het systeem ook op zoek gaat naar berichten die op een manier worden verspreid die niet normaal is.

Van de Universiteit van Antwerpen en de Universiteit van Leuven zijn onderzoekers aan de slag met een nieuw systeem dat op zoek gaat naar berichten die op een manier worden verspreid die niet normaal is.

Van de Universiteit van Antwerpen en de Universiteit van Leuven zijn onderzoekers aan de slag met een nieuw systeem dat op zoek gaat naar berichten die op een manier worden verspreid die niet normaal is.

pestgedrag automatisch op

door
Bart De Weert

**'Elk bericht op
Facebook zal een
label meekrijgen'**

Opmerkelijk is dat het systeem ook op zoek gaat naar berichten die op een manier worden verspreid die niet normaal is.

Van de Universiteit van Antwerpen en de Universiteit van Leuven zijn onderzoekers aan de slag met een nieuw systeem dat op zoek gaat naar berichten die op een manier worden verspreid die niet normaal is.

Van de Universiteit van Antwerpen en de Universiteit van Leuven zijn onderzoekers aan de slag met een nieuw systeem dat op zoek gaat naar berichten die op een manier worden verspreid die niet normaal is.

GEZOCHT: CYBERPESTBERICHTEN

GEZOCHT: CYBERPESTBERICHTEN
Veilig internetten is belangrijk, maar niet altijd even gemakkelijk. Heb jij al een keer rare of onveilige dingen meegemaakt op het internet?

Verscheidene studies geven aan dat één Vlaamse jongere op tien recent het slachtoffer was van cyberpesten. De impact van dergelijke situaties is vaak erg groot doordat kwetsbare berichten soms langere tijd online blijven staan en slachtoffers alleen met dit probleem zitten omdat ze er niet over willen of durven te praten.

SAMEN STERK
In het wetenschappelijk project AMICA slaan een aantal onderzoeksgroepen van de universiteiten van Antwerpen, Leuven en Gent de handen in elkaar om een systeem te ontwikkelen dat automatisch kritische situaties zoals cyberpestgedrag en seksueel grensoverschrijdend gedrag herkent op sociale netwerken om zo een veilige internetomgeving te kunnen garanderen voor jongeren.

BLIJF NIET BIJ DE PAKKEN ZITTEN
Om te kunnen bepalen wat cyberpesten is en hoe het kan worden herkend, zijn voldoende onderzoekers nodig van deze online berichten. Daarom doet AMICA een oproep aan u als ouder om zoveel mogelijk berichten aan ons beschikbaar te stellen waarin cyberpesten merkbaar is.

Als uw kind getuige of slachtoffer is geweest van cyberpesten en als hiervan bewijsmateriaal beschikbaar is, dan zijn deze berichten meer dan welkom. De berichten kunnen bijvoorbeeld e-mails zijn, sms'en, chatgesprekken of berichten van sociale media. De data worden anoniem behandeld en alleen gebruikt voor intern onderzoek, berichten worden in geen geval doorgegeven aan derden.

Wij zijn contact opgenomen met AMICA (Automatic Monitoring for Cyberspace Applications)? Stuur dan een mailje naar data@amicaproject.be en blijf niet bij de pakken zitten. Meer informatie vind je op de website: www.amicaproject.be.

Dataset Construction: simulation experiments



- Role playing in secondary schools on social media platform: FB-like social network, scenarios, profile cards (roles), debriefing
- Additional goal: education (prevention)

The screenshot shows a user profile on a platform called AMiCA. The profile is for Dominique Verhaegen, who has a blonde avatar and a purple shirt. The page layout includes a top navigation bar with 'Nieuwsoverzicht' and '[Berichten]'. Below the profile picture, there are links for 'Zend Bericht', 'Blokkeer Gebruiker', 'Rapporteer', and 'Admin Settings'. The main content area shows a post by Dominique Verhaegen directed at Joni Claes, with a text about cancer and a link to 'Vind ik leuk'. Below this, there are several comments from other users like Emma Dewaele and Laura Van Boom, all dated '9 mei'. At the bottom right, there is a status 'Vrienden Online (0)'.

AMiCA

Nieuwsoverzicht [Berichten]

Dominique Verhaegen
net een zalig dagje gehad met sam! #blijftdemijnevoortaltijd

Updates Info Vrienden (5)

Dominique Verhaegen → **Joni Claes**:
Laat Sam nu eindelijk is met rust! Ik hoop echt dat ge een pijnlijke dood sterft loser. Vat vol miserie zijt gij en een ongelooflijk debiele kankermens zonder hart
Vind ik leuk · Reageer · 9 mei

Julie De Backer vindt dit leuk.

Emma Dewaele DOE IS RUSTIG GIJ VUILE BITCH
9 mei · Vind ik leuk · 1 vindt dit leuk

Joni Claes Vind je jezelf nu beter dan mij nu je dit allemaal zegt? Zoek een leven en scheld niet met kanker, dat is onrespectvol.
9 mei · Vind ik leuk · 1 vindt dit leuk

Laura Van Boom ja wa is u probleem?
9 mei · Vind ik leuk

Dominique Verhaegen wa moeide gij u nu weer! ga terug zulgen aan u tampons kankerhoer
9 mei · Vind ik leuk

Emma Dewaele GIJ ZIJT EEN DEBIELE KIND DOMINIQUE
9 mei · Vind ik leuk

Laura Van Boom Je echt een achterlijk kind
9 mei · Vind ik leuk

Zend Bericht
Blokkeer Gebruiker
Rapporteer
Admin Settings

GEMEENSCHAPPELIJKE VRIENDEN

Vrienden Online (0)



Data Annotation

- Brat rapid annotation tool (Stenetorp et al., 2012)
- Two annotation levels (Van Hee et al., 2015)
 - Post level
 - **Cyberbullying -vs- non-cyberbullying**
textual content that is published online by an individual and that is aggressive or hurtful against a victim.
 - **Harmfulness score**
 - 0 → the post does not contain indications of cyberbullying
 - 1 → the post contains indications of cyberbullying, although they are not severe
 - 2 → the post contains serious indications of cyberbullying
 - **Author's role**
 - Harasser
 - Bystander-defender
 - Victim
 - Bystander-assistant



Data Annotation

- (Sub)sentence level: identification of **fine-grained** text categories related to cyberbullying
 - Threat/blackmail
 - Insult
 - Curse/exclusion
 - Defamation
 - Sexual talk
 - Defense
 - Encouragements (to the harasser)

[Guidelines for the fine-grained analysis of cyberbullying, version 1.0](#) (2015)

Van Hee, C., Verhoeven, B., Lefever, E., De Pauw, G., Daelemans, W., & Hoste, V.



Data Annotation

Category	Brat annotation example	Translation
Threat/blackmail Expressions containing physical or psychological threats, or indications of blackmail.	<p>2_Har Threat or Blackmail als ik u tegen kom zieke rak op u gezicht x</p>	<i>I'll smash you in the face when I see you x</i>
Insult Expressions containing abusive, degrading or offensive language that are meant to insult the addressee.	<p>1_Har General Insult General Insult HAHAHAHA LOSER GIJ:(X AARDAPPELKOP</p>	<i>HAHAHAHA YOU LOSER :(X POTATO HEAD</i>
Curse/exclusion Expressions of a wish that some form of adversity or misfortune will befall the victim and expressions that exclude the victim from a conversation or a social group.	<p>2_Har Curse or Exclusion General Insult Pleeg zelfmoord niemand vindt u geestig ...</p>	<i>Just commit suicide, nobody thinks you're funny...</i>
Defamation Expressions that reveal confident, embarrassing or defamatory information about the victim to a large public.	<p>1_Har Defamation u mama versiert andere mannen hahahaha</p>	<i>Your mom is flirting with other men hahaha</i>
Sexual talk Expressions with a sexual meaning that are possibly harmful.	<p>1_Har Sexual harassment Stuur my u naaktfoto, nu!!</p>	<i>Send me a naked picture of yourself, now!!</i>
Defense Expressions in support of the victim, expressed by the victim himself or by a bystander.	<p>1_Bystander_defender General victim defense General victim defense Meid, koppie omhoog he! Laat je ni doen door die domme anoniempjes</p>	<i>Cheer up girl, don't let those stupid anon's make you feel bad</i>
Encouragements to the harasser Expressions in support of the harasser.	<p>2_Bystander_assistant General Insult Encouraging harasser inderdaad ze is geen leven waard !!</p>	<i>Indeed, she shouldn't be alive !!</i>

Ask.fm preliminary experiments



- Class
 - Binary (bullying or non-bullying)
 - Binary (for each fine-grained class)
- Features
 - Word unigrams and bigrams
 - Character trigrams
 - Sentiment features
- Classifier: SVM (Pattern) with linear kernel
- Data: ~85,000 posts
- Annotation agreement (kappa) 60-65%
- Very skewed data, scarce positive data (~10%)

Van Hee, C., Lefever, E., Verhoeven, B., Mennes, J., Desmet, B., De Pauw, G., Daelemans, W. & Hoste, V. (2015). Detection and fine-grained classification of cyberbullying events. Proceedings of RANLP, 672–680. Hissar, Bulgaria.



Results

	Precision	recall	F1-score
NL	76%	56%	65%
EN	74%	55%	63%

BUT:

- Ambiguity

“Hi bitches, anyone in for a movie tonight?”

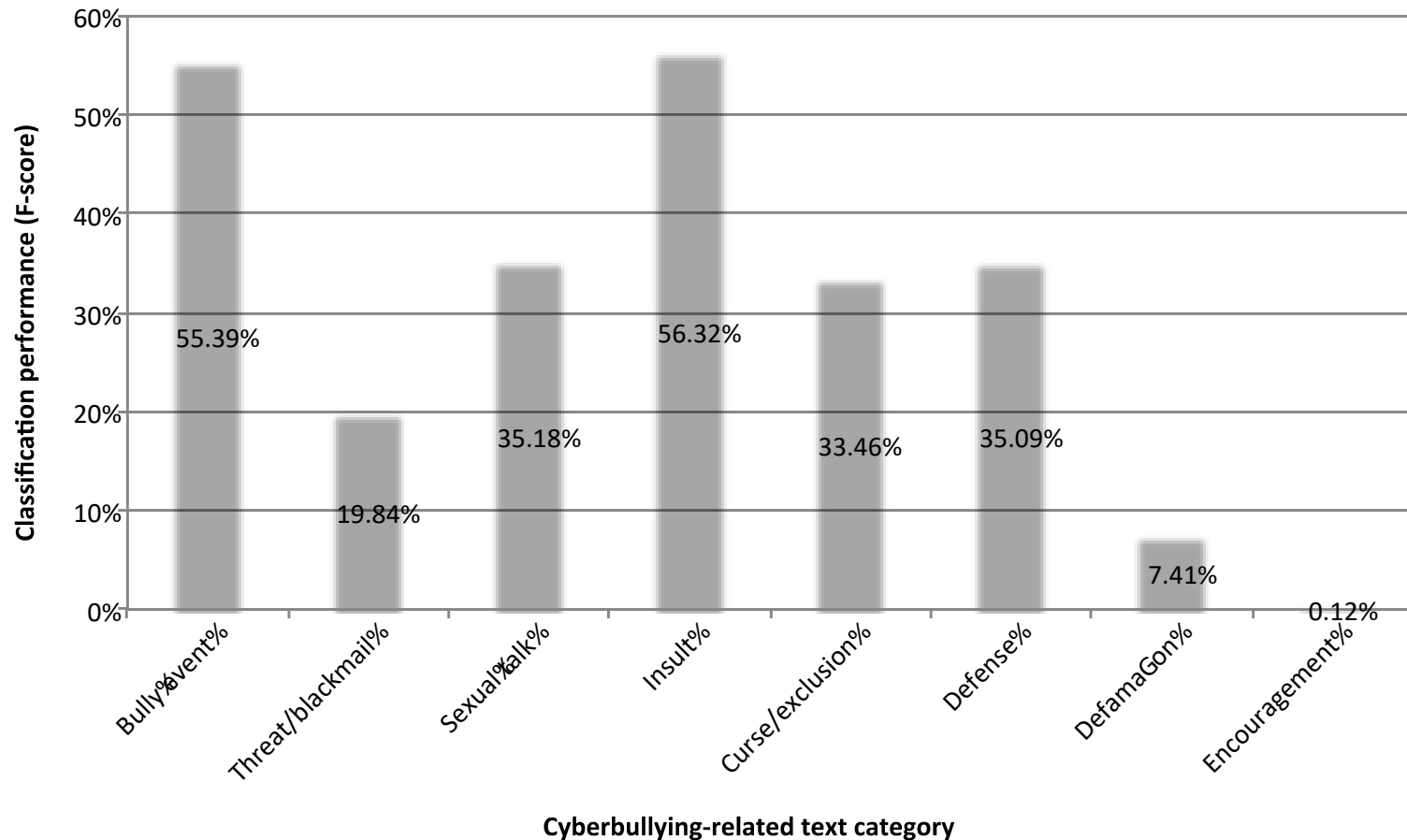
“Shut up, you bitch!”

- Implicit realizations of cyberbullying

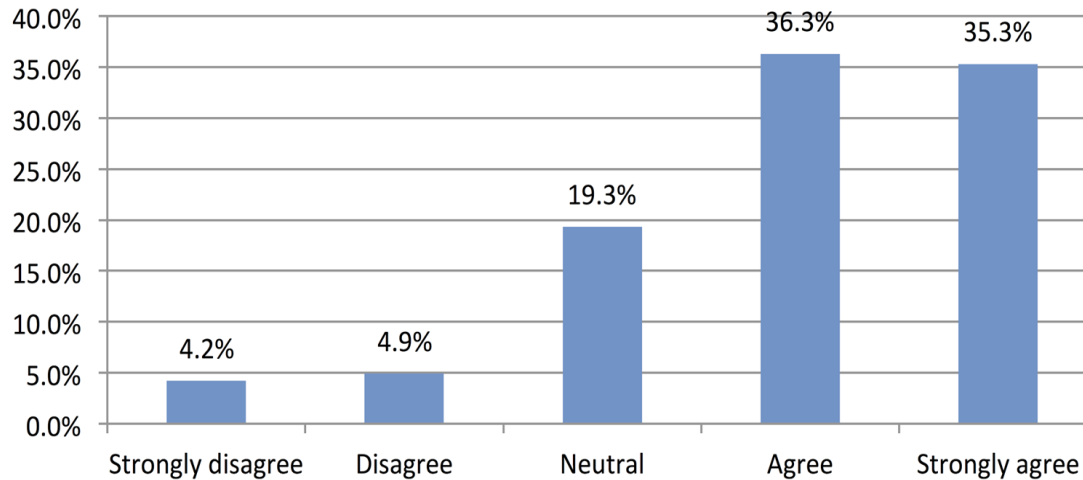
“You make my fists itch...”

- Data sparseness

Results (Van Hee et al. 2015)



Monitoring desirable?



- Follow-up is needed
- Privacy of youngsters should be respected
- Technical feasibility?

(Van Royen et al., 2014)

More info?



Cynthia Van Hee: cynthia.vanhee@ugent.be





USE CASE 2: SUICIDE DETECTION

Alarming figures Flemish adolescents



- **Self-mutilation:**
 - Every year by 7% at the age of 14-17
 - 2/3 through cutting & scratching
(Van Rijsselberghe et al., 2009)
- **Suicidal behaviour:**
 - 15-20% (age of 18) have thoughts of suicide (more than once) (Hublet et al., 2010)

Online self-harm behaviour



Kheb het al 3 keer geprobeerd,
ma kloop ier nog altijd rond... soms
zeg ik spijtig genoeg, soms ben ik
ook blij dat ik nog leef.





AMiCA technology: image analysis

- Automatic classification of images
- Object recognition in images
- Tekst recognition in images + OCR



If I jump now,
who will catch me?

If I jump now

who will catch me?

AMiCA technology: text analysis



Machine learning system **analyses** every message (word sequences, topic models, sentiment analysis, ...) and **answers** **two questions**:

- Is the message about suicide?



I never thought about cutting or suicide, because it leaves scars ...

- Is there a serious suicidal threat?



I already tried 3 times, but I'm still alive



Sometimes I feel bad, sometimes I'm glad I'm still alive



Text analysis: results

Experiments carried out on a data set of 10,000 messages, of which 851 are relevant and 257 are serious:

- Is the message about suicide? => recall: 9/10, 3% noise
- Is there a serious suicidal threat? => recall: 2/3, 25% noise



Does it work in practice?

What is the impact of the automatic detection system in a moderator setting?

Simulation of high work load of moderators:

- task: identify alarming messages that need a response (75)
- Lots of messages (1000)
- Limited moderation time (1 hour)
- Collaboration with CPZ (Flemish centre for suicide prevention) and moderators of the website “Wel Jong Niet Hetero” (LGBT web site)
- 1 group with / 1 group without system aid

Valorisation: interface



System

suicide-prevention.lt3.ugent.be/nlp/system/

Bart

Annotations for suicide prevention

log out

Marked as reaction

All

System

All

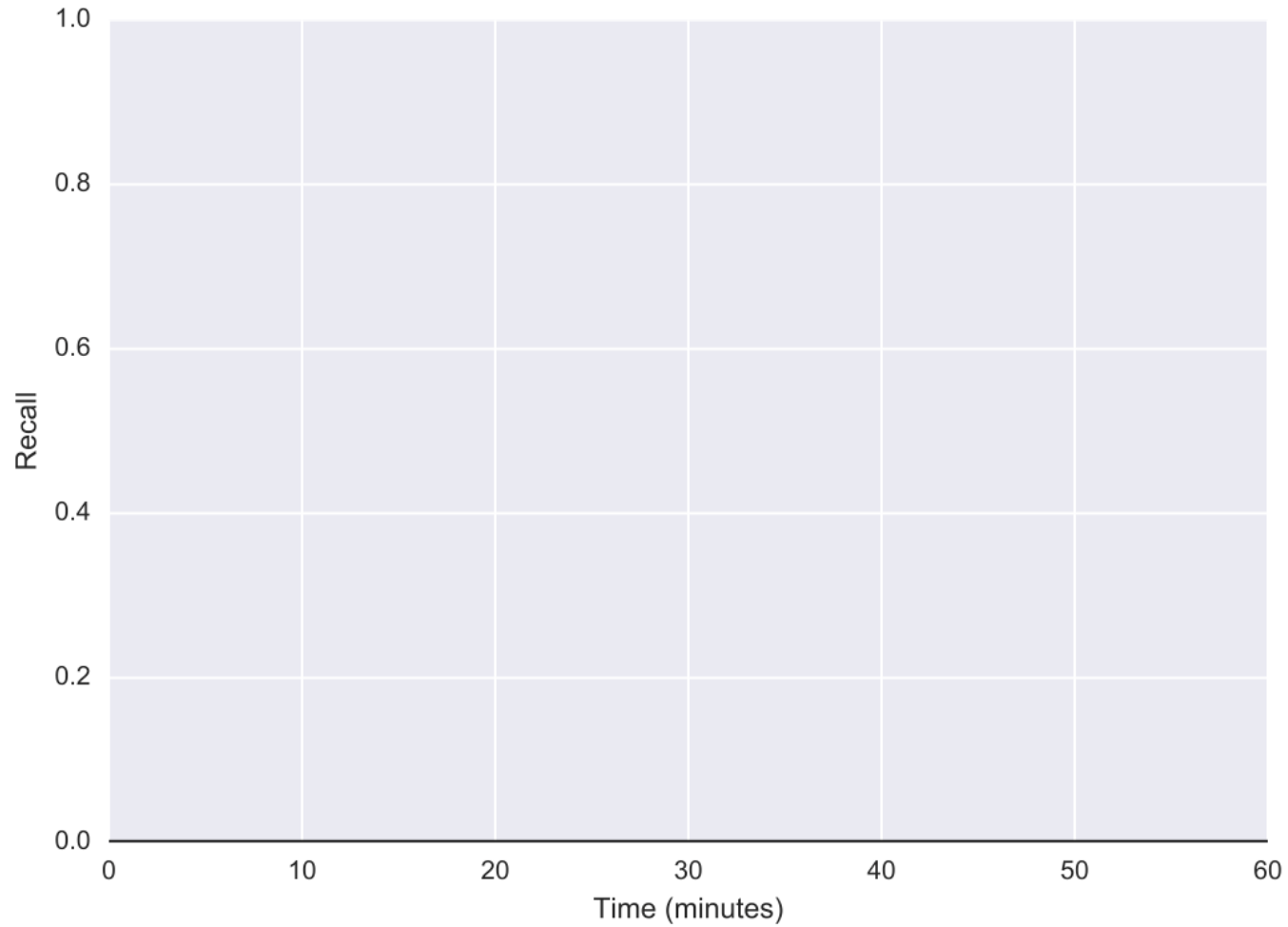
Messages 1000

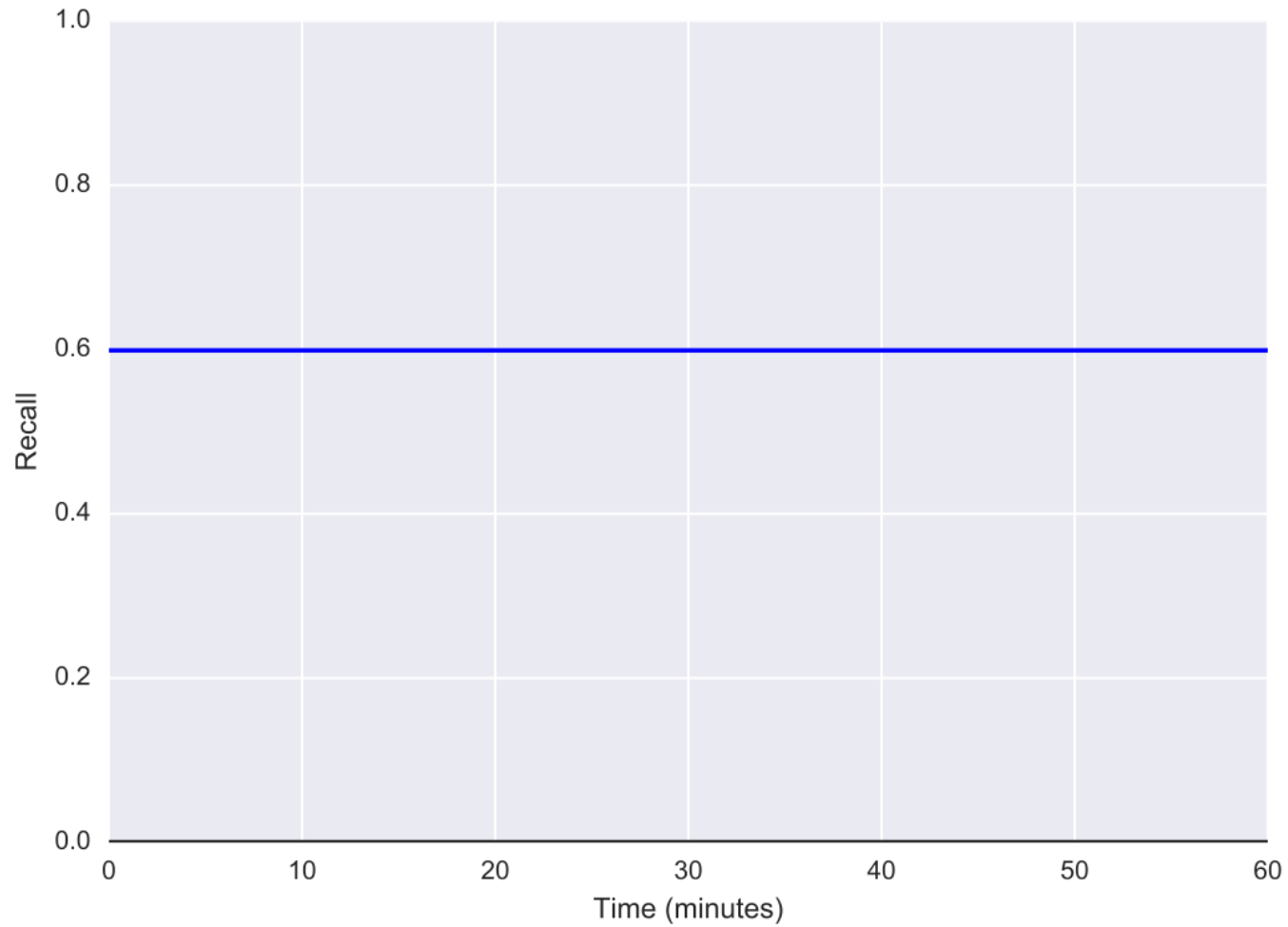
← Previous

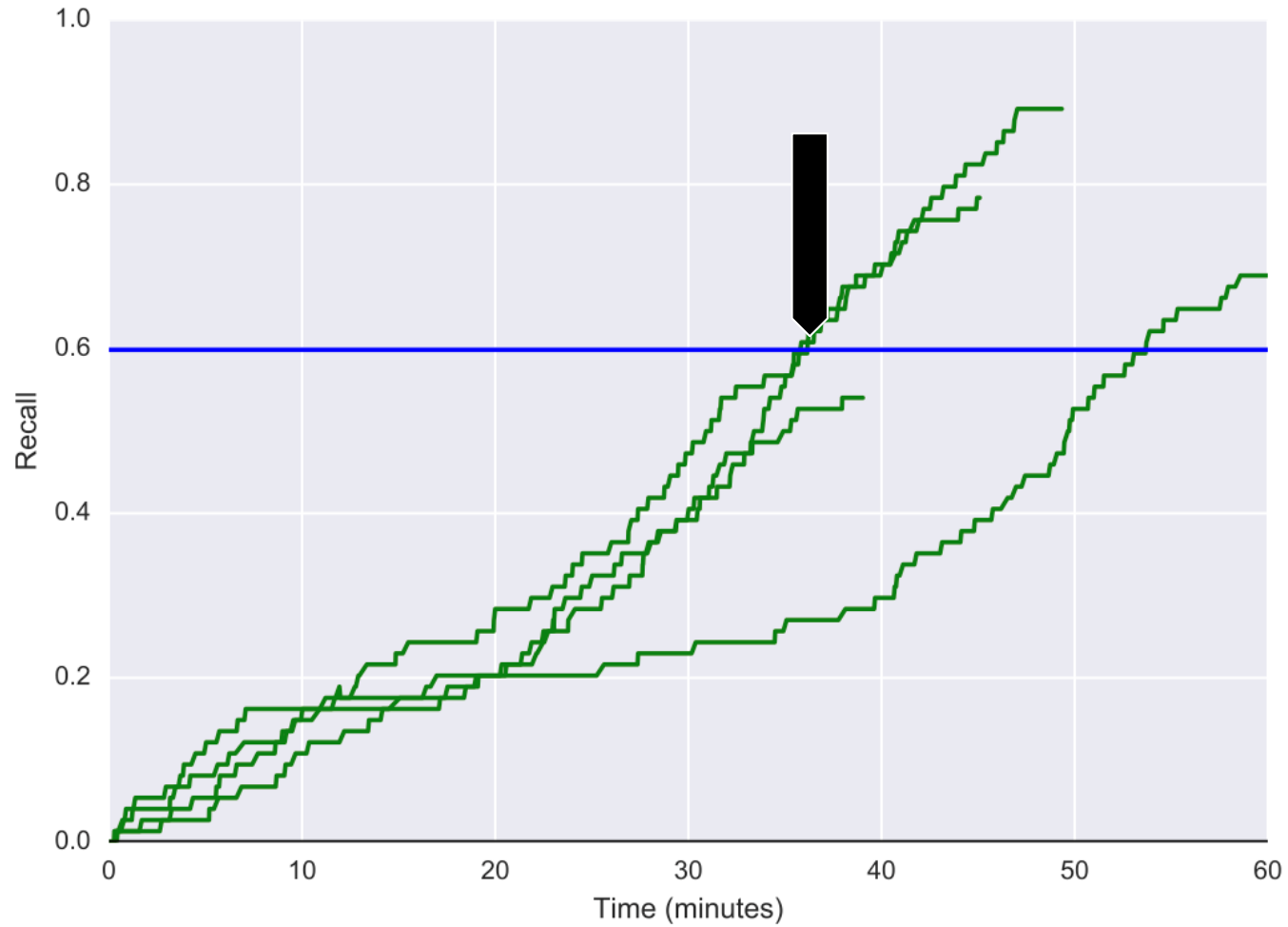
Page 1 of 20.

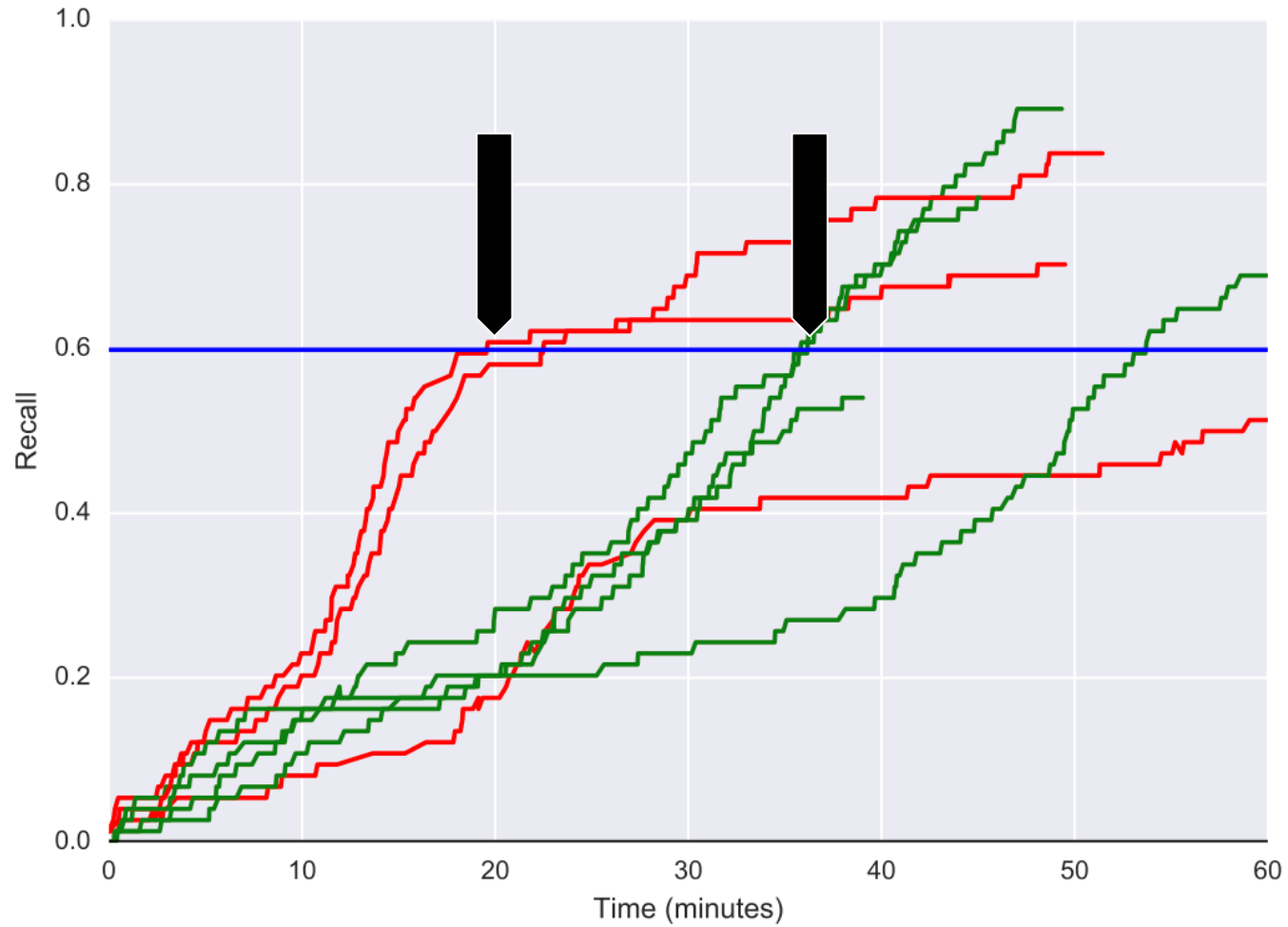
Next →

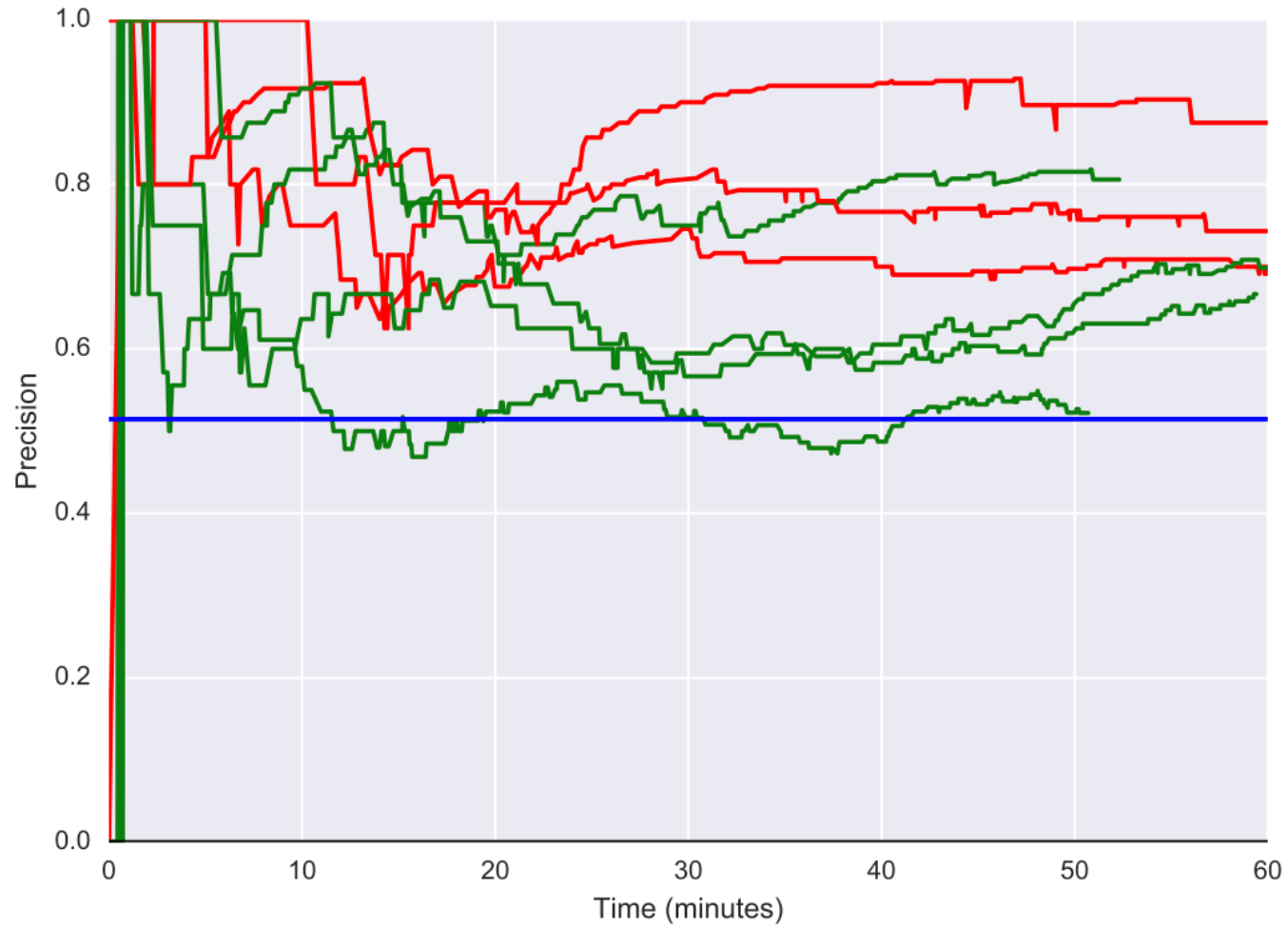
Date	Message	Reaction	System
2013-10-17 13:55	Heb je tips? Ikzelf droom om gitaar te kunnen spelen Er staat hier thuis een gitaar, maar behalve enkele domme dingskes lukt spelen niet echt Nog zo'n droom is drummen, dat lijkt me zooo zalig !	<div>Yes</div> <div>No</div>	No
2007-07-01 15:07	maar ik ken u ni, ze. (ik zat daar met mijn heteromaat in een hoekje. hij wou ni echt veel van het feestje meemaken, jammer genoeg)	<div>Yes</div> <div>No</div>	No
2015-06-26 15:46	Zal ik zeker doen! Ja, super nieuws tussen al die aanslagen!	<div>Yes</div> <div>No</div>	No
2015-06-26 23:24	Van alle homofobe reacties op de legalisering van het homohuwelijk in Amerika, vind ik dat deze toch wel de originaliteitsprijs verdient:	<div>Yes</div> <div>No</div>	No













More info?

Bart Desmet: bart.desmet@ugent.be

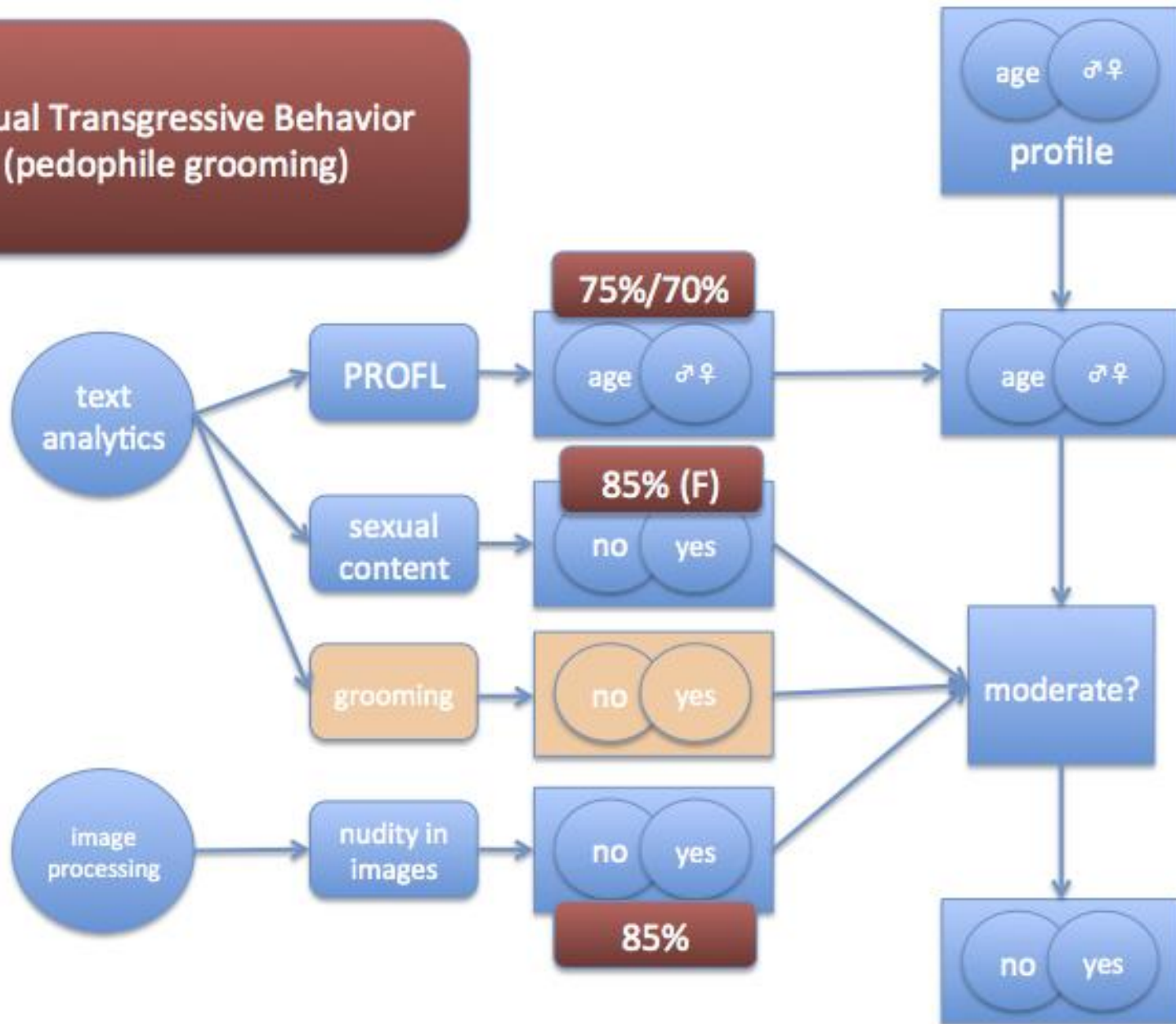




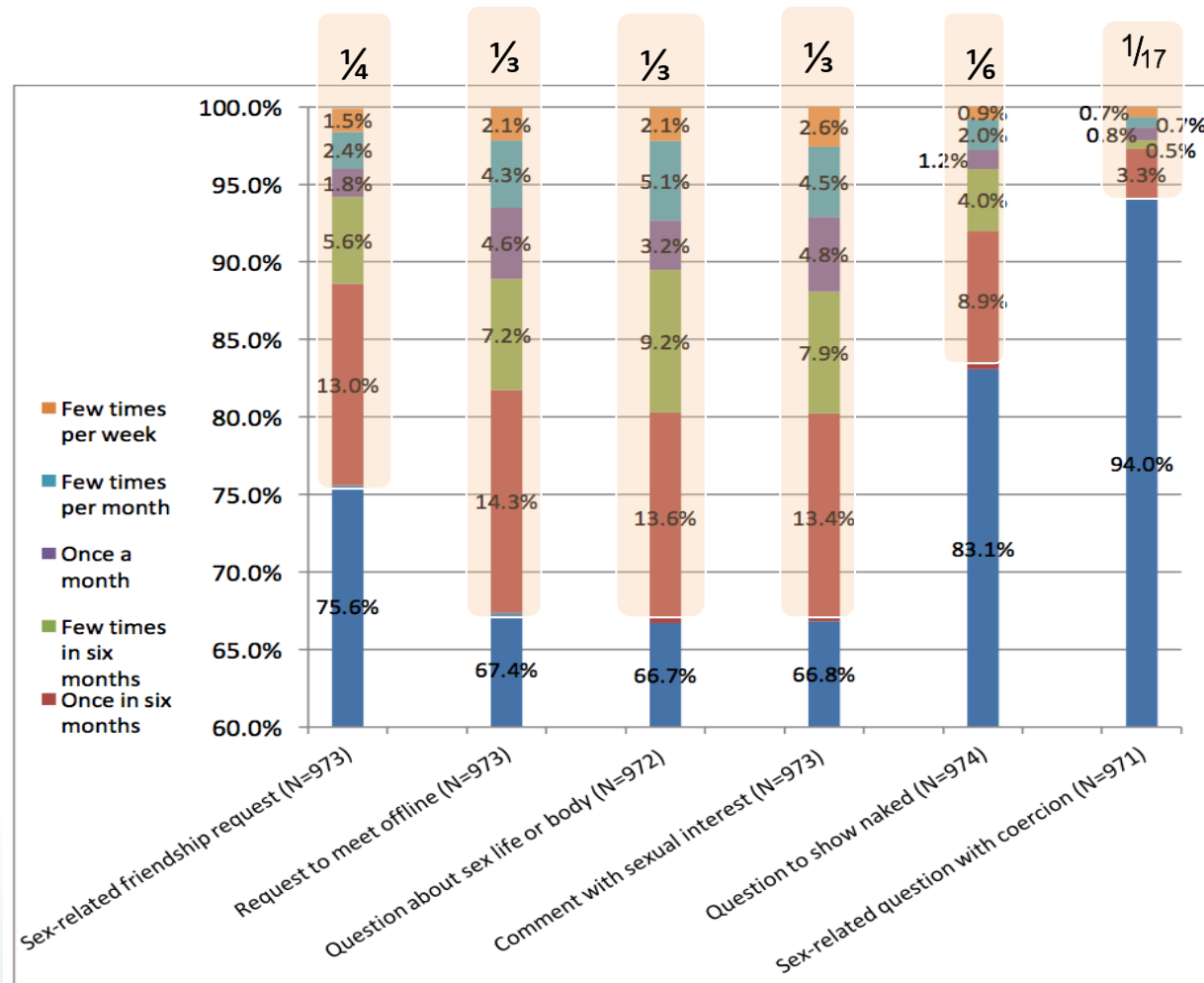
USE CASE 3: PROFILING FOR DETECTING PEDOPHILE GROOMING



Sexual Transgressive Behavior (pedophile grooming)



Motivation





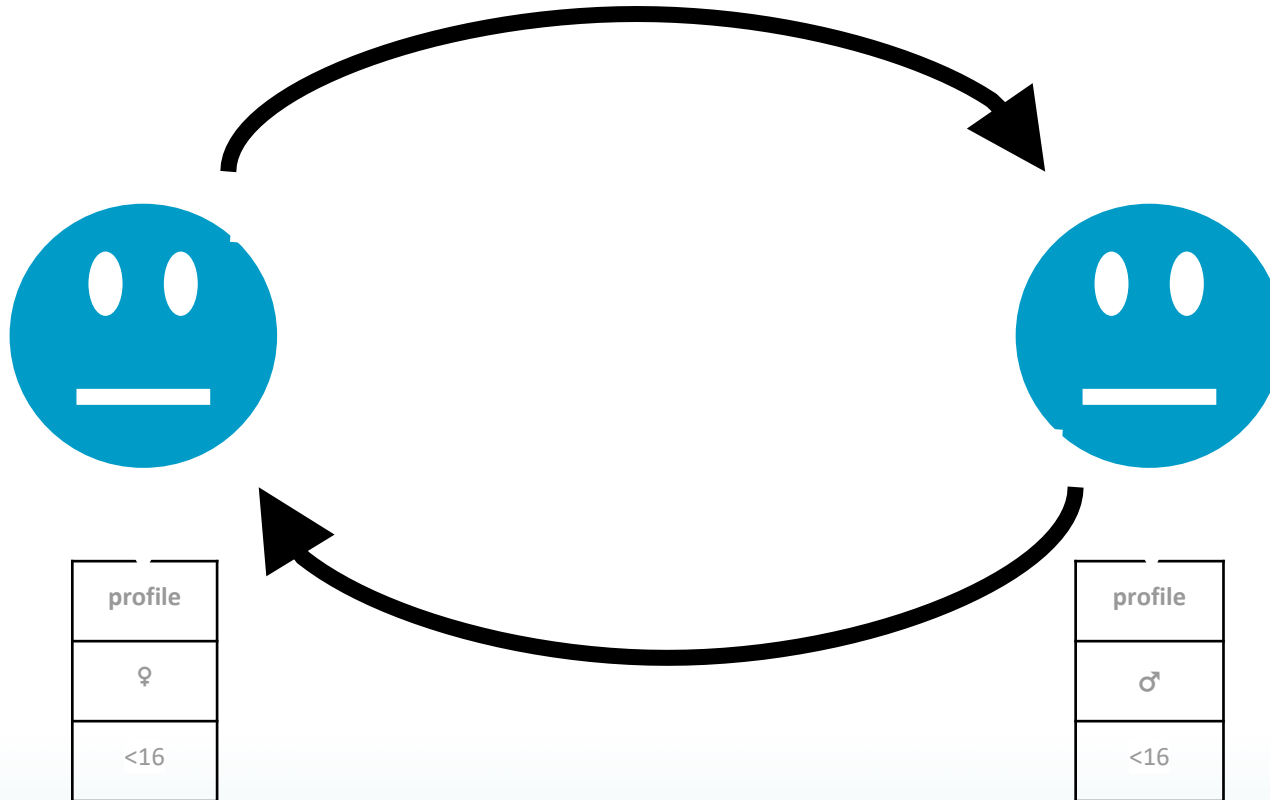
Motivation

- Survey: ± 1000 youngsters about the frequency, nature and appropriateness of sexual messages on social media
- Especially on Facebook
- Who?
 - 32% strangers
 - 29% friends IRL
 - 19% online friends

67% didn't like the message + 11% reported the incident

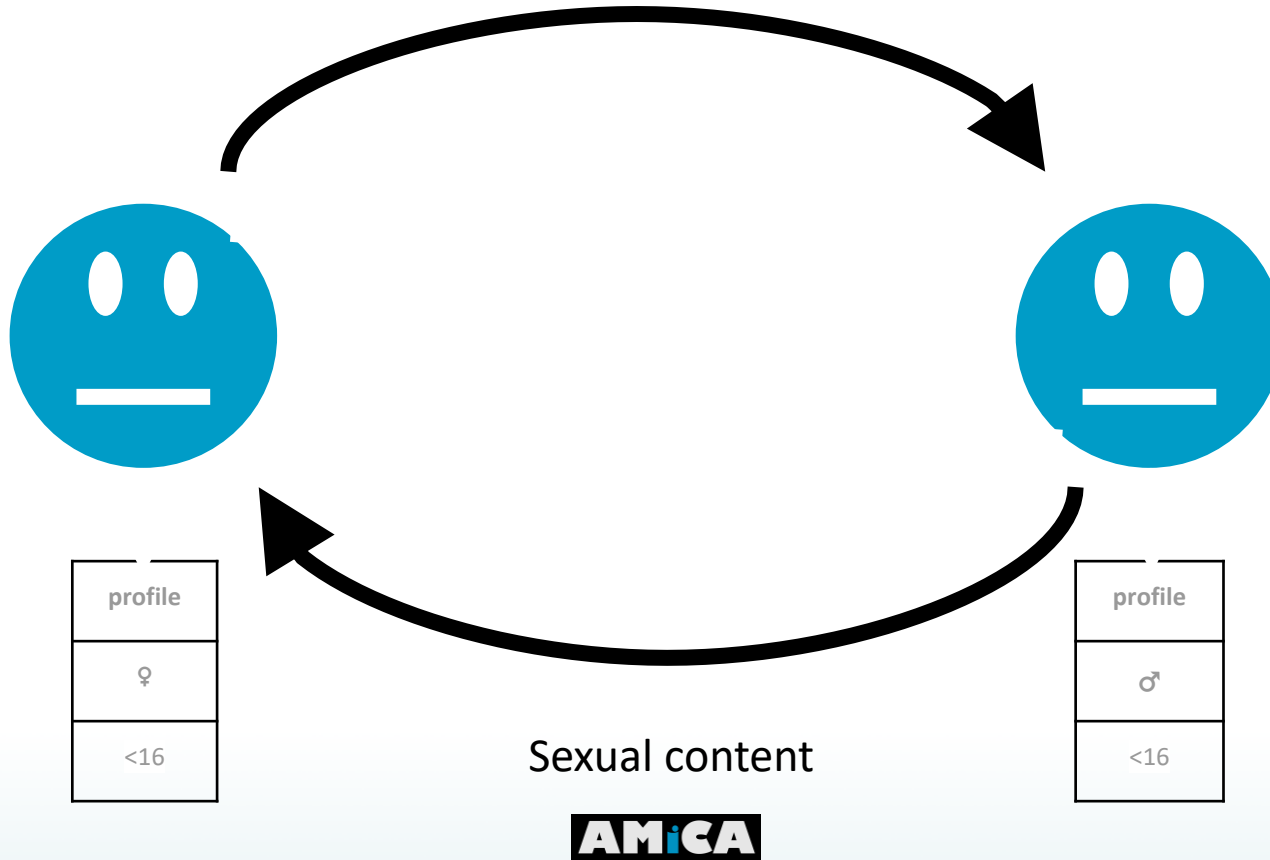


Approach



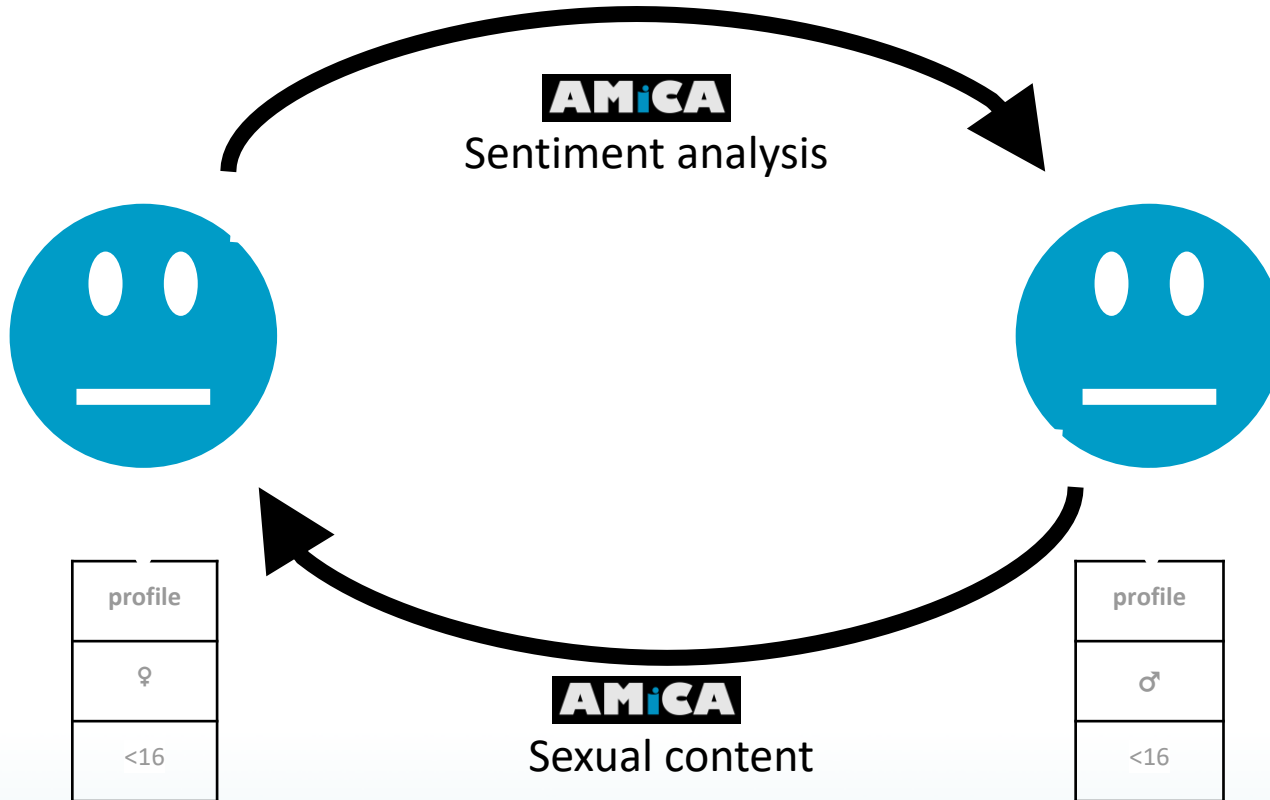


Approach

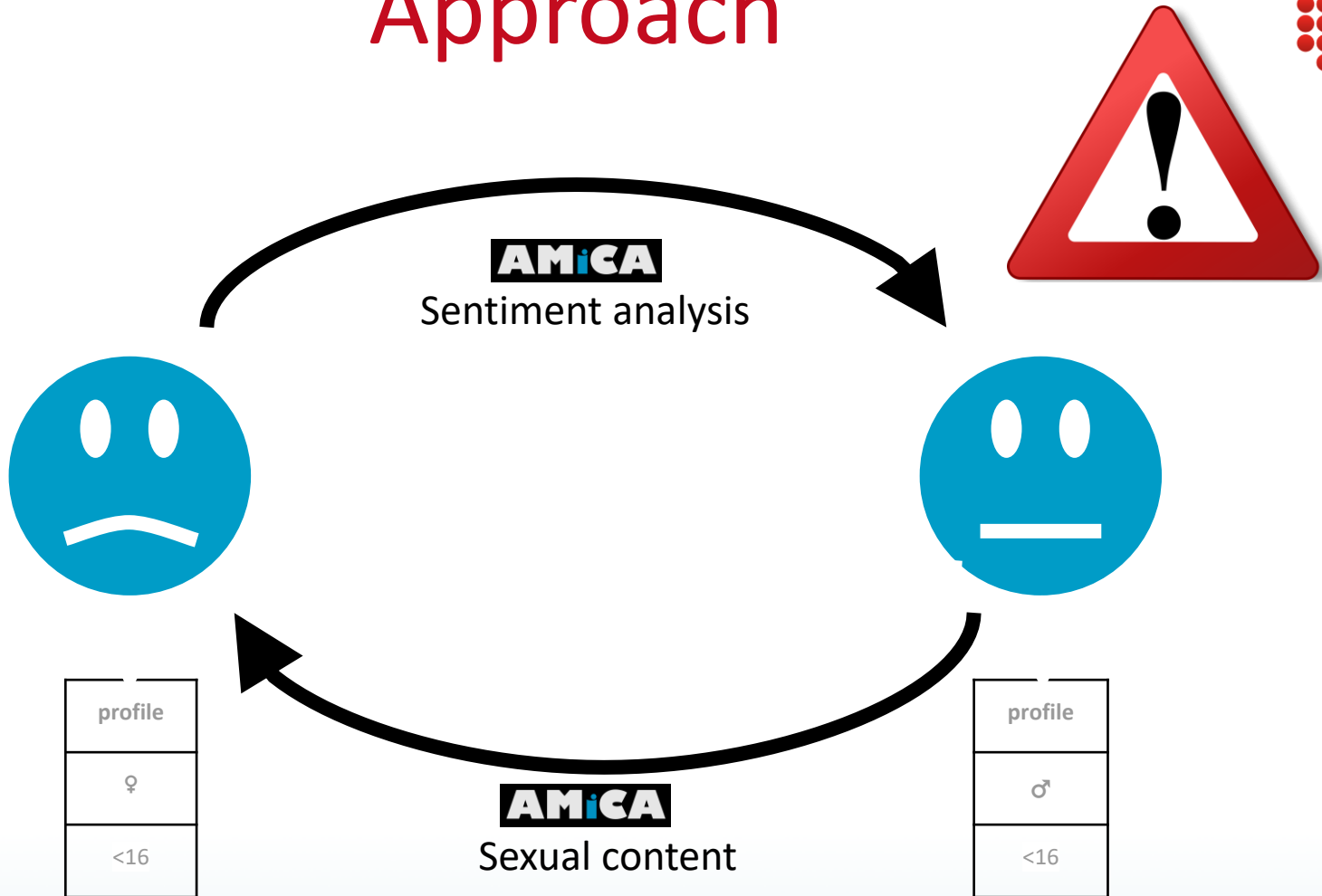




Approach

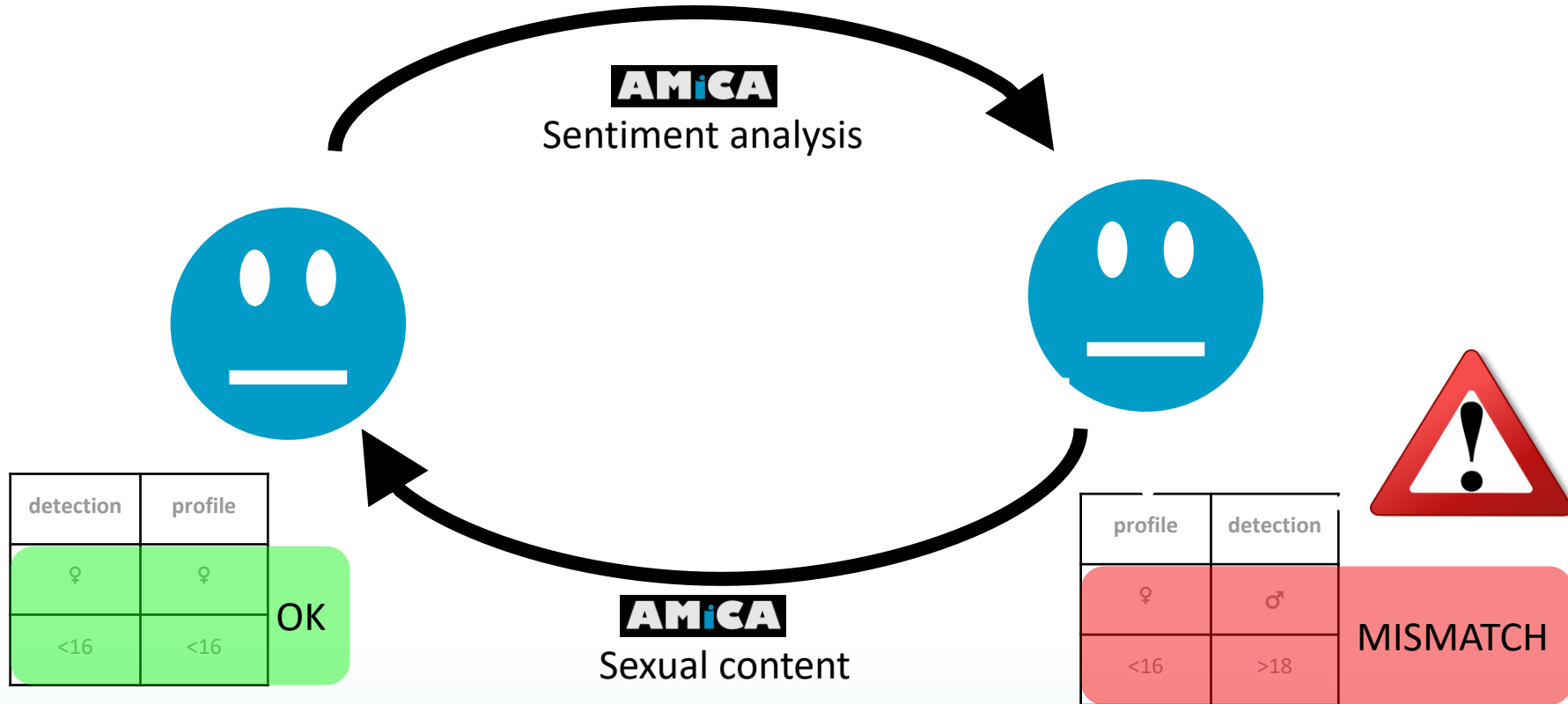


Approach





Approach





Profiling

- AMiCA profiler
 - Based on Chris Emmery's OMESA
 - <https://github.com/cmry/omesa>
- Age and Gender
 - Finding dubious SN profiles
 - Computed age and gender does not match given information
 - Optimizing recall (for moderator application)
 - Adapting to binary classification
 - Legally relevant age difference



Approach

- SN chat data (Netlog, 2010-2011)
 - 380k posts
 - 87k users
 - Data point = combined posts of a single user
 - Self-reported age, gender, and location
- Classes: age (binary), gender, age+gender
- 5-fold cross-validation
- SVM with linear kernel
- Features:
 - token n-grams
 - character n-grams



Results

- Gender
 - ~70%
 - Adding different types of features (LIWC, POS patterns, sentiment, etc) boosts f-scores slightly



Results

- Age:
 - Distinguish between users above and below age of consent (16 in Belgium), -16 versus +18 has priority
 - Optimize recall
 - Using cost and confidence parameters in SVMs
 - Up to 95% recall for -16; 92% recall for +18

Ref: Janneke van de Loo , Guy De Pauw, Walter Daelemans, Text-Based Age and Gender Prediction for Online Safety, International Journal of Cyber-Security and Digital Forensics (IJCSDF), 2016, 46-60.



Predator Detection

- Two classifiers
 - LiBSVM
 - Classify at the post level, aggregate at user level
 - Classify at the user level directly
 - Weighted voting of previous
 - Additional constraints
 - E.g. only one pedophile per conversation

Claudia Peersman, Frederik Vaassen, Vincent Van Asch, Walter Daelemans. Conversation Level Constraints on Pedophile Detection in Chat Rooms. CLEF 2012 (PAN), 2012.



Overall test results

- Grooming detection
 - Predator detection
 - 72 % f-score, 89% precision, 60% recall
 - Suspicious posts
 - 30% f-score, 36% precision, 26% recall



More info?

Walter Daelemans:
walter.daelemans@uantwerpen.be



Guy De Pauw:
guy.depauw@uantwerpen.be





DISCUSSION



discussion

- Is normalization and automatic detection accurate enough for applications in cybersecurity?
 - Precision - Recall trade-off
- Should we protect children and young people in social networks against their will?
 - Protection - privacy trade-off



Thank you!

Els.lefever@ugent.be

<http://www.amicaproject.be/>