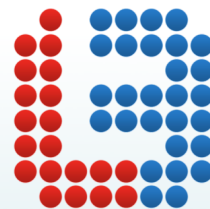# Text Analysis for Social Media Cybersecurity: the AMiCA Project

Els Lefever

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

UNIVERSITEIT GENT

language and translation technology team

# LT³, LANGUAGE AND TRANSLATION TECHNOLOGY TEAM

# LT³

- 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)
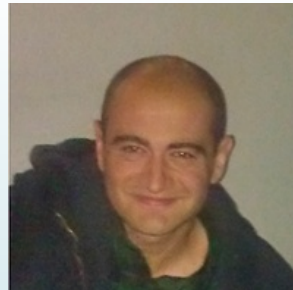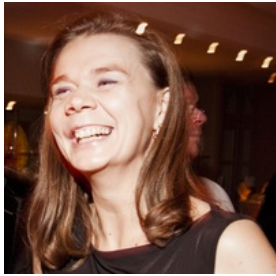- Use of term extraction for translating documentaries (Sabien Hanoulle)
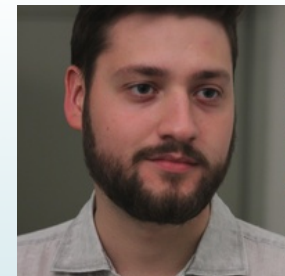
# Translation Technology

- 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
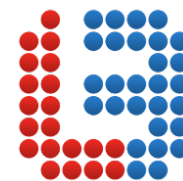
# www.amicaproject.be

- 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)

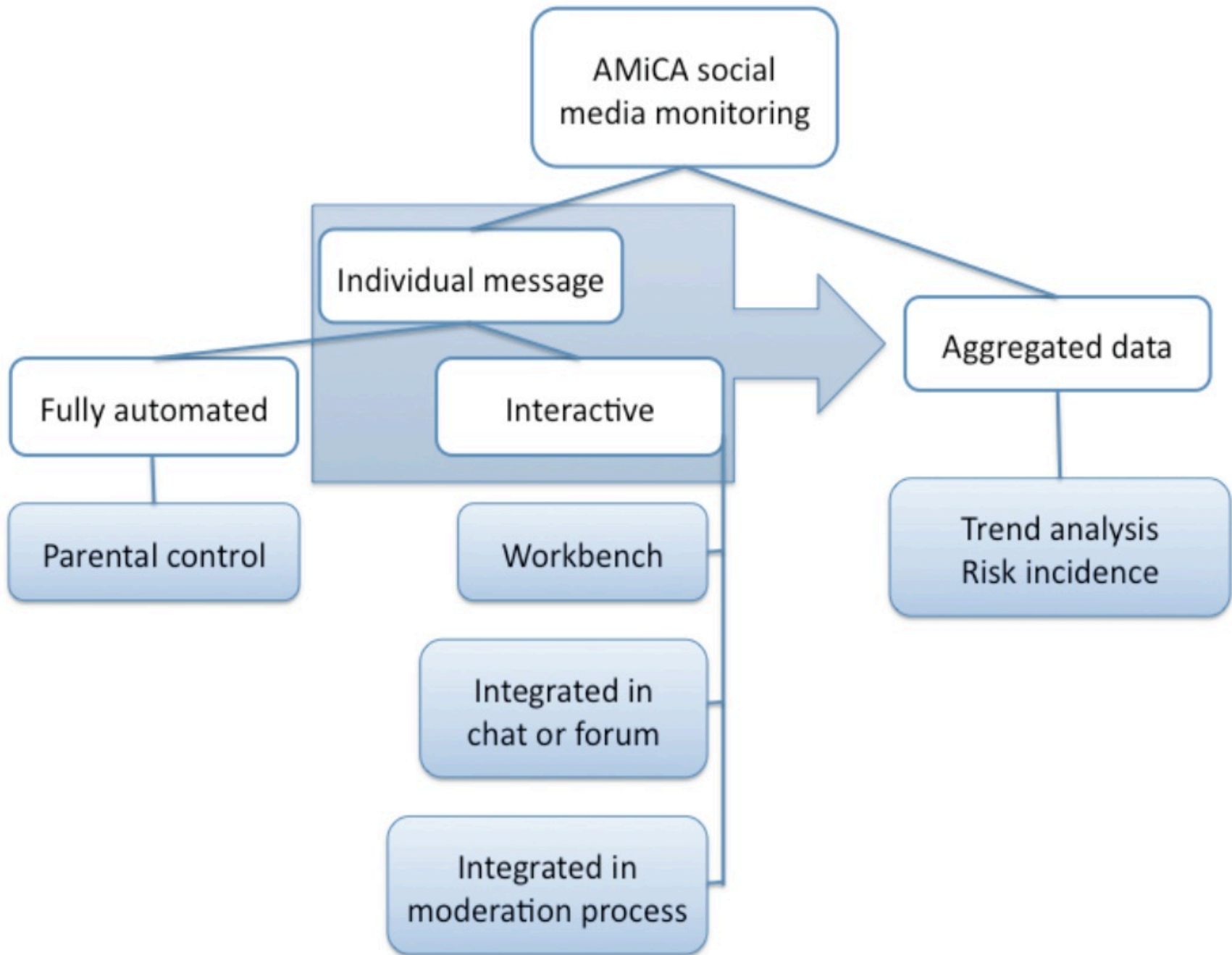- Combine text analytics, image and video analysis, and data mining

# 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?
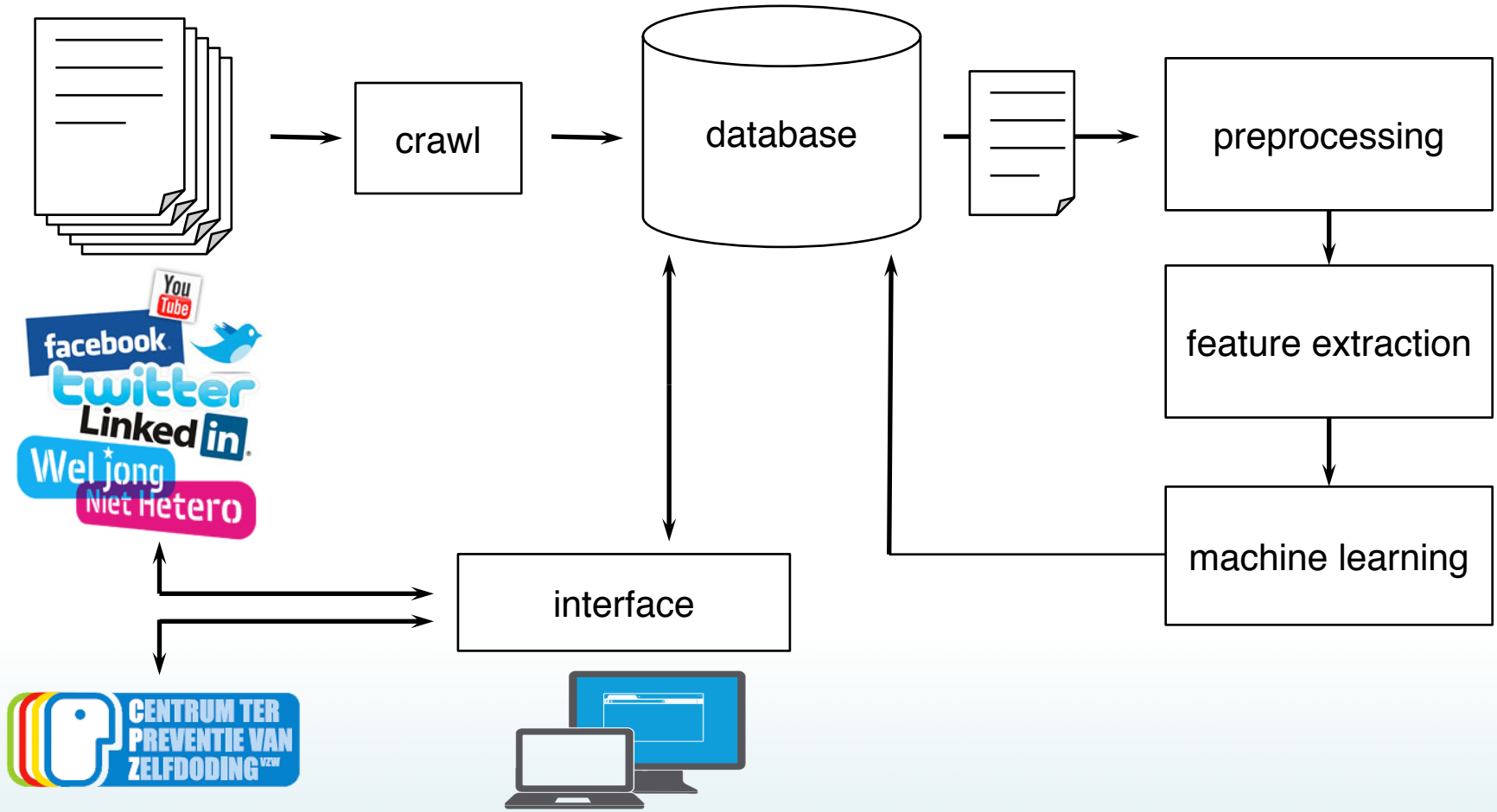
# Should we do something about it?

- 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:**

https://twitter.com/SonicStef/statuses/62

**Timestamp:**

Date: 2015-09-29    Today | 📅

Time: 22:52:47    Now | 🕐

Note: You are 2 hours ahead of server time.

**Text:**

Zwijg stomme trut! Gij hebt geen leven tot o

**User name:**

Pestertje1998

20

# 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

# 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 emotion:
  • Flooding (loooooooove)
  • 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
  - <u>Machine Translation</u>: Source Language = non-standard and Target Language = standard
  - <u>Spell Checking</u>: Correct the incorrect words (statistical or dictionary-based)
  - <u>Speech Recognition</u>: 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

# Modules

- <u>Preprocessing</u>
  - Tokenization and sentence splitting
    - includes emoticons, emojis etc.
  - Character floooooooding
- <u>Token-based modules</u>
  - 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

- ± 20-40% 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



**RESULT:** ± 30 reactions

± 368 messages (FB messages, hate pages, Netlog, mail, chat, etc.)

# 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)

# 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 <u>not</u> contain indications of cyberbullying
      - 1 → the post contains <u>indications</u> of cyberbullying, although they are <u>not severe</u>
      - 2 → the post contains <u>serious indications</u> of cyberbullying
    - Author's role
      - Harasser
      - Victim
      - Bystander-defender
      - 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. | **2_Har** [Threat or Blackmail] ¶ als ik u tegen kom zieke rak op u gezicht x | *I'll smash you in the face when I see you x* |
| **Insult** Expressions containing abusive, degrading or offensive language that are meant to insult the addressee. | **1_Har** [General insult] [General insult] ¶ HAHAHAHA LOSER GIJ:( X AARDAPPELKOP | *HAHAHAHA YOU LOSER :( X POTATO HEAD* |
| **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. | **2_Har** [Curse or Exclusion] [General insult] ¶ Pleeg zelfmoord niemand vindt u geestig ... | *Just commit suicide, nobody thinks you're funny...* |
| **Defamation** Expressions that reveal confident, embarrassing or defamatory information about the victim to a large public. | **1_Har** [Defamation] ¶ u mama versiert andere mannen hahahaha | *Your mom is flirting with other men hahahaha* |
| **Sexual talk** Expressions with a sexual meaning that are possibly harmful. | **1_Har** [Sexual harassment] ¶ Stuur my u naaktfoto, nu!! | *Send me a naked picture of yourself, now!!* |
| **Defense** Expressions in support of the victim, expressed by the victim himself or by a bystander. | **1_Bystander_defender** [General victim defense] [General victim defense] ¶ Meid, koppie omhoog he! Laat je ni doen door die domme anoniempjes | *Cheer up girl, don't let those stupid anons make you feel bad* |
| **Encouragements to the harasser** Expressions in support of the harasser. | **2_Bystander_assistant** [General insult] [Encouraging harasser] ¶ inderdaad ze is geen leven waard !! | *Indeed, she shouldn't be alive !!* |

# 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

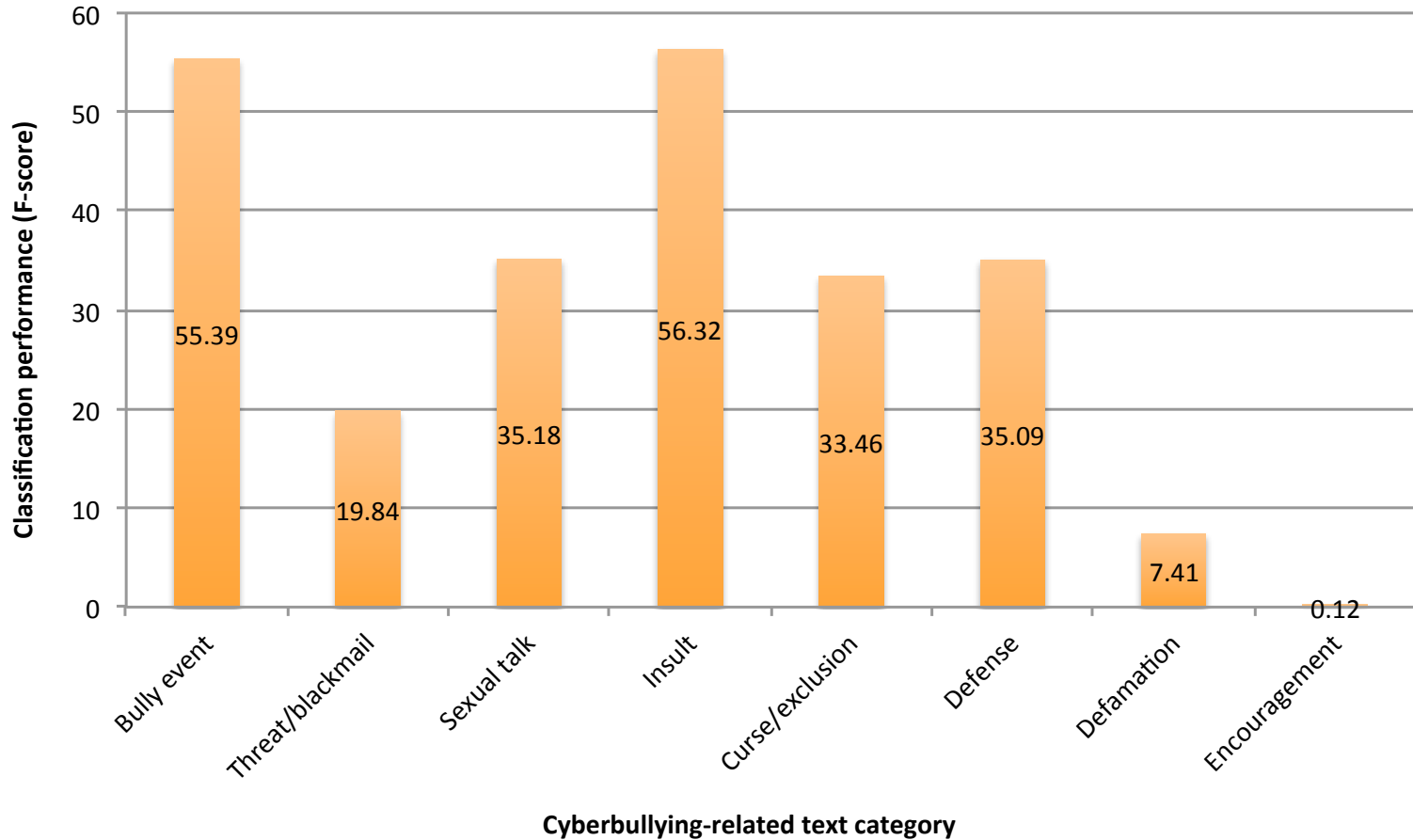  *"Hi bitches, anyone in for a movie tonight?"*
  *"Shut up, you bitch!"*

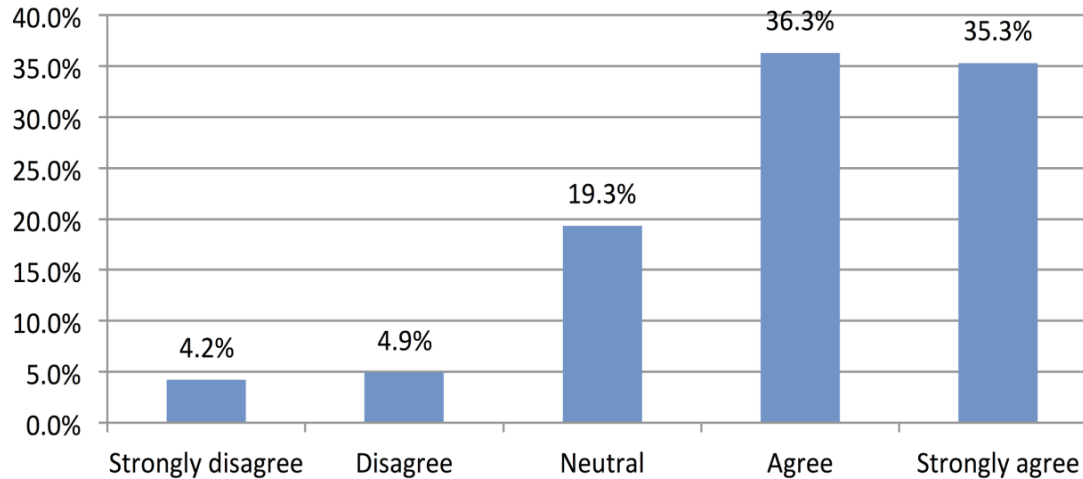- Implicit realizations of cyberbullying

  *"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?

# 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

# More info?
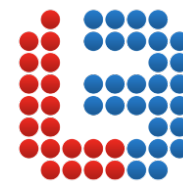
Bart Desmet: bart.desmet@ugent.be

# USE CASE 3: PROFILING FOR DETECTING PEDOPHILE GROOMING

Sexual Transgressive Behavior (pedophile grooming)

profile: age ♂♀

text analytics → PROFL → 75%/70% (age ♂♀) → age ♂♀

text analytics → sexual content → 85% (F) (no / yes)

text analytics → grooming → no / yes

image processing → nudity in images → no / yes → 85%

moderate? → no / yes

# 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

AMiCA
Sentiment analysis

AMiCA
Sexual content

| detection | profile | |
|-----------|---------|----|
| ♀ | ♀ | OK |
| <16 | <16 | |

| profile | detection | |
|---------|-----------|----------|
| ♀ | ♂ | MISMATCH |
| <16 | >18 | |

# 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/