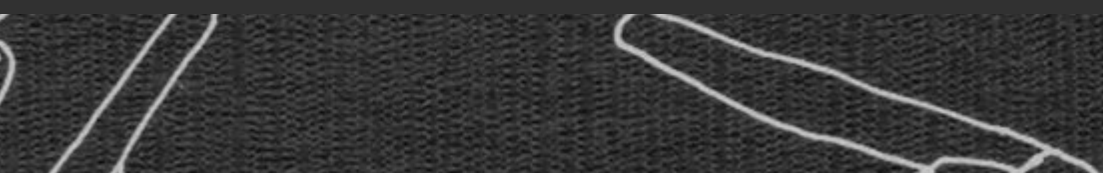
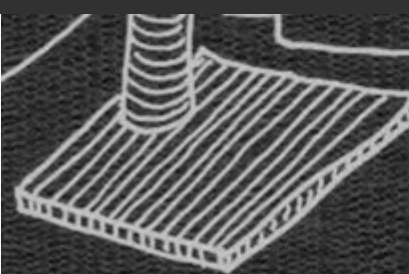




Hatbudskap och våldsbejakande extremism i digitala miljöer



Lisa Kaati

lisa.kaati@dsv.su.se

Katie Cohen

katie.cohen@foi.se

Toxiskt språk i digitala miljöer

- kommunikationshandlingar som förgiftar samtalsklimatet:
- hets mot folkgrupp, förtal, förgripelse mot tjänsteman
- integritetskränkningar, respektlöshet
- näthat, hat och hot
- hate speech, dangerous speech

Vad är toxiskt?

- Kontextberoende
- Koder
- Vi har alla olika uppfattning om vad som är toxiskt

Toxiskt språk

most government and media offices shut down or run on a skeleton crew on weekends. Financial markets don't move on weekends and thus don't drive news. Weekends are basically the downtown in which they plot the plan and agenda and how the media will promote it for the coming week. News and info always moves slower on weekends.

Toxiskt språk

most government and media offices shut down or run on a skeleton crew on weekends. Financial markets don't move on weekends and thus don't drive news. Weekends are basically the downtown in which ((they)) plot the plan and agenda and how the media will promote it for the coming week. News and info always moves slower on weekends.

Automatiska metoder för att känna igen toxiskt språk

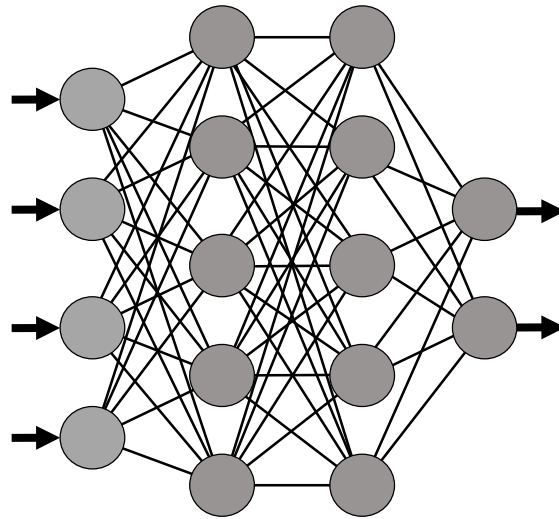
- Flertal olika metoder: från enkla ordlistebaserade metoder till avancerade AI-system
- Språkberoende– vissa språk är lågresurs språk
- Maskininlärningsbaserade metoder behöver massor av träningsdata

Transfer learning med språkmodell

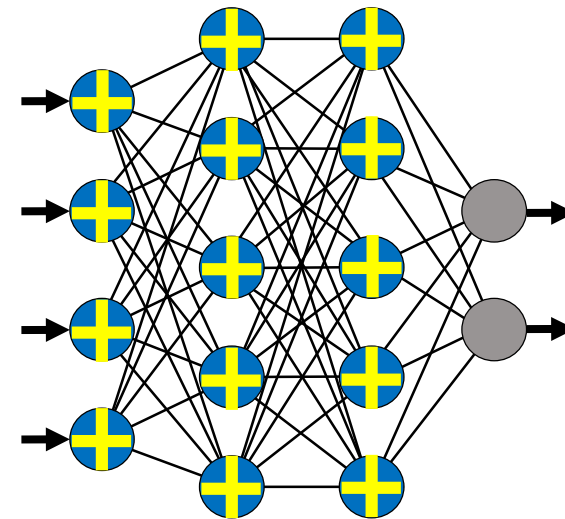


WIKIPEDIA
Den fria encyklopedin

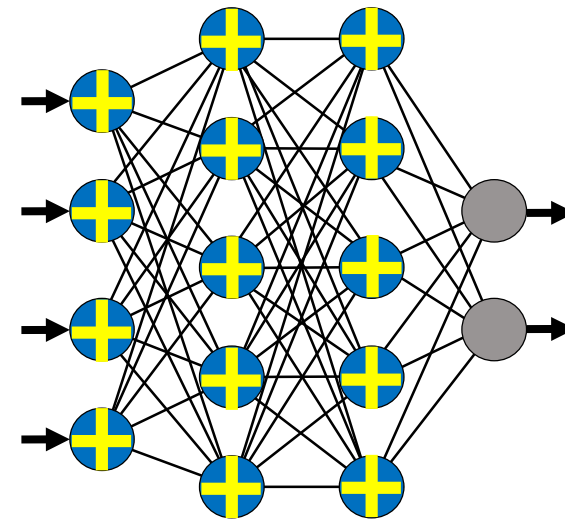
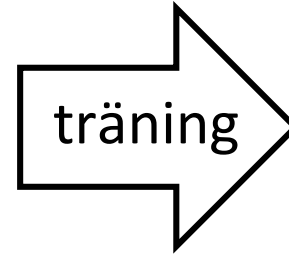
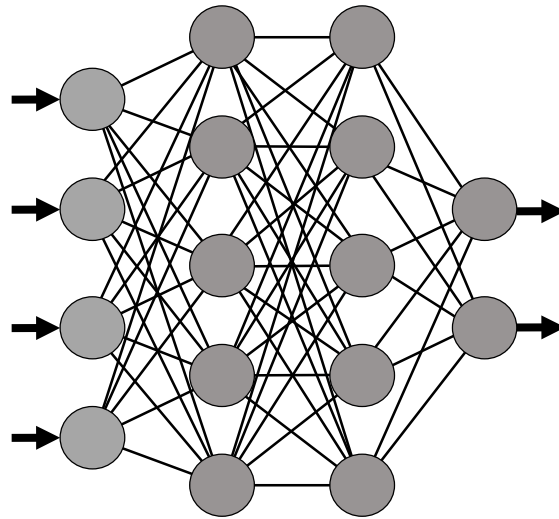
FamiljeLiv.se



träning

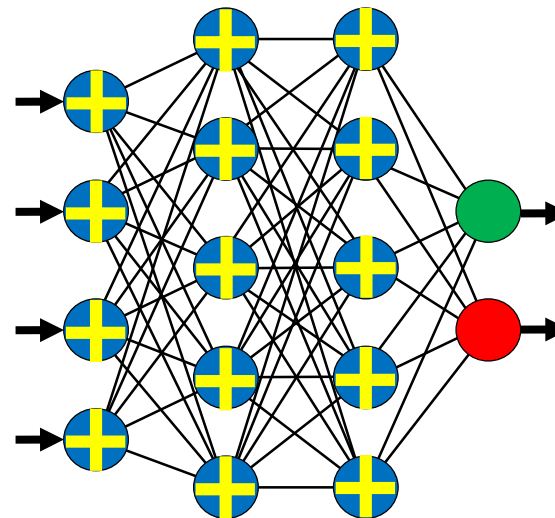


Transfer learning med språkmodell



Fine-tuning: träna vidare för speciell uppgift med annoterade exempel:

Nån borde knuffa C framför bussen!

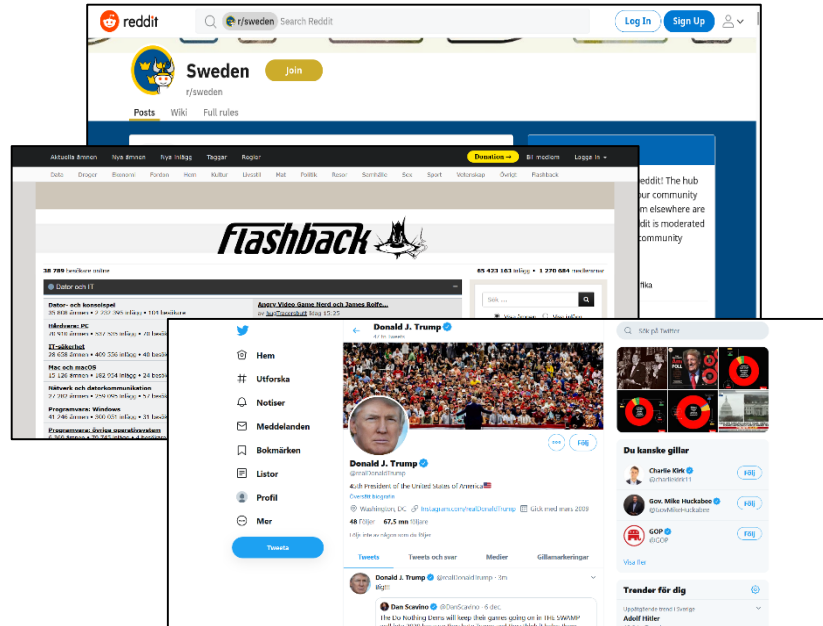


ja, toxiskt

Automatisk detektion av toxiskt språk

- Alla automatiska metoder för textanalys går fel ibland
- Inga metoder automatiska metoder kan förstå text lika bra som en människa...
- ...men dom är mycket snabbare och tröttnar inte.
- Våra automatiska metoder är lämpliga för:
 - kvantitativa analyser av stora textmängder
 - gallring: t.ex. reducera 1 miljon tweets till 3 000 som vi kan kolla för hand

Toxiskt språk i svenska digitala miljöer

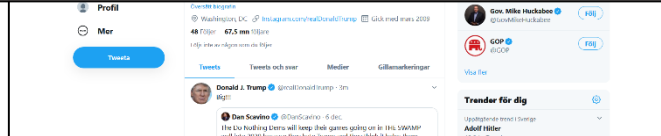


Toxiskt språk i svenska digitala miljöer

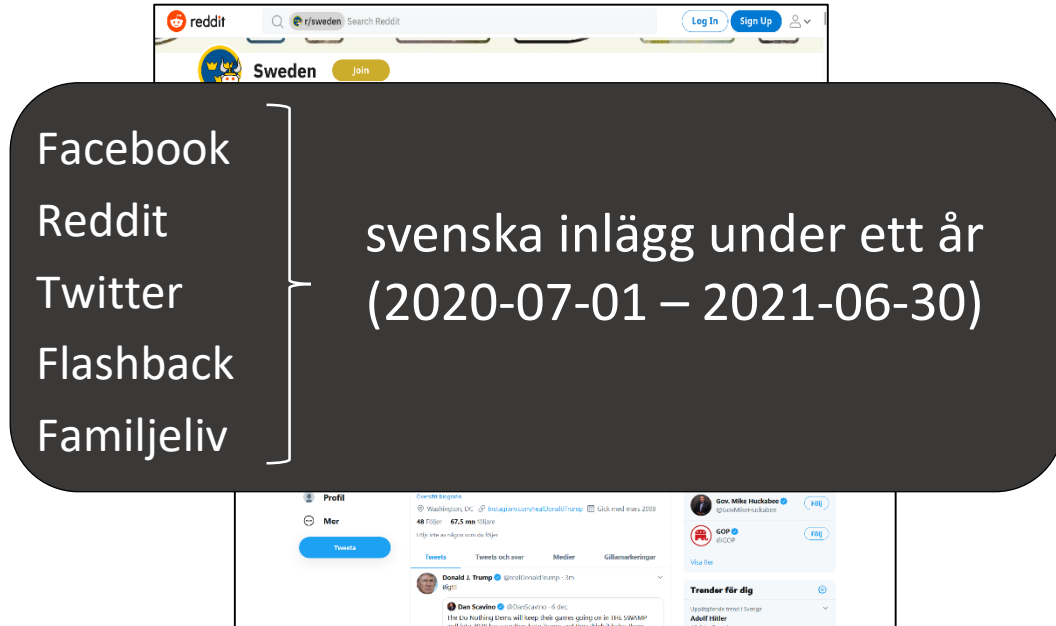


Facebook
Reddit
Twitter
Flashback
Familjeliv

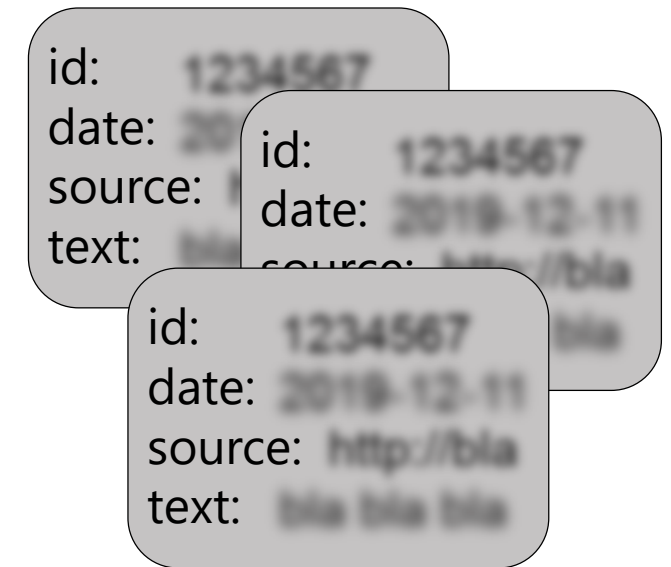
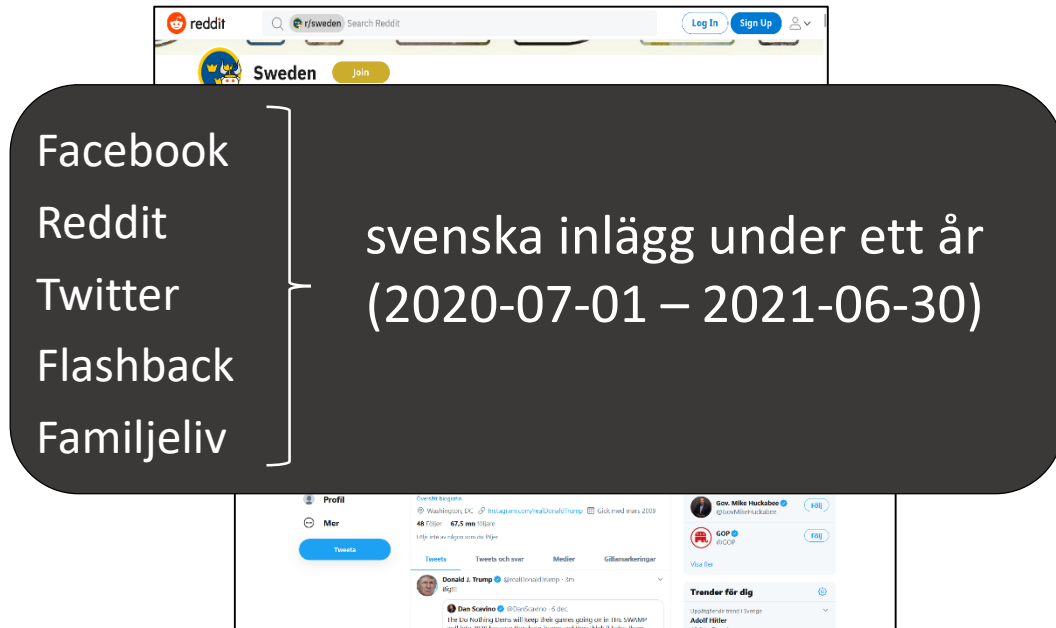
svenska inlägg under ett år
(2020-07-01 – 2021-06-30)



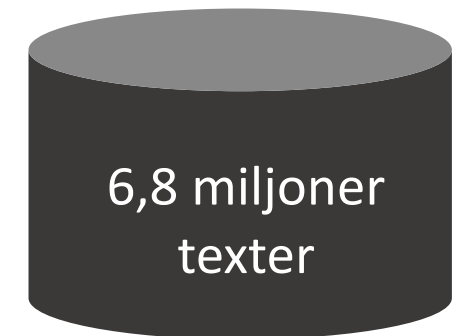
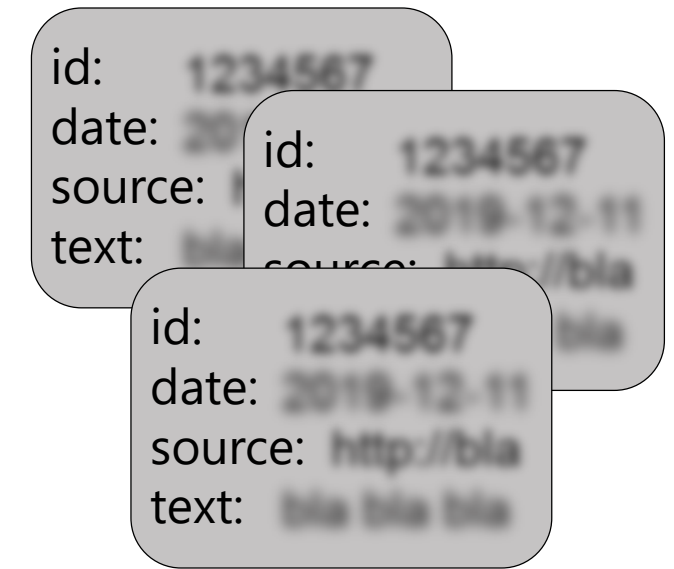
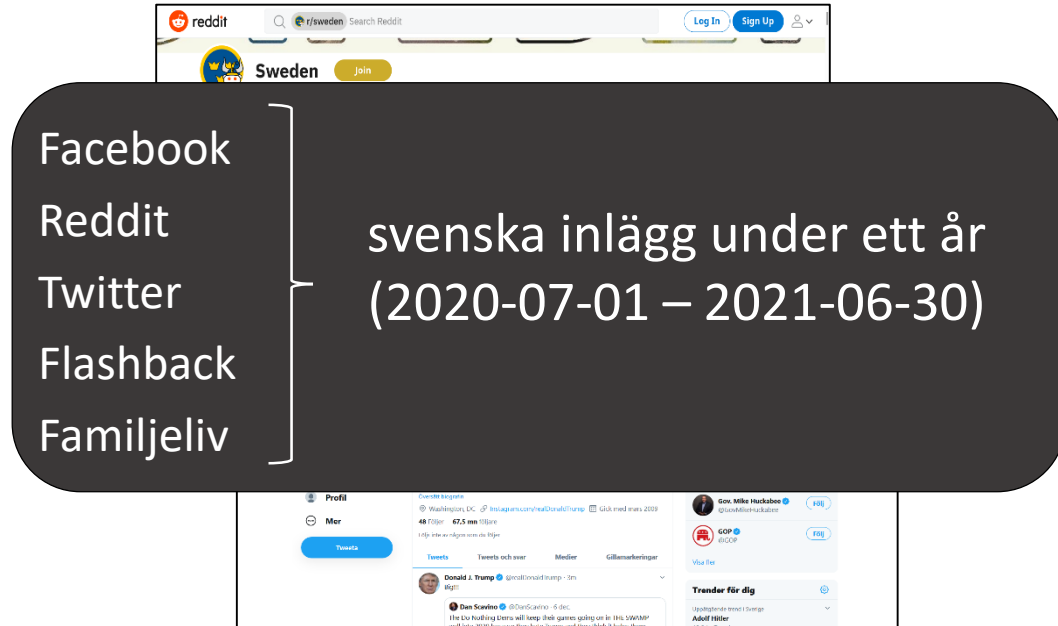
Toxiskt språk i svenska digitala miljöer



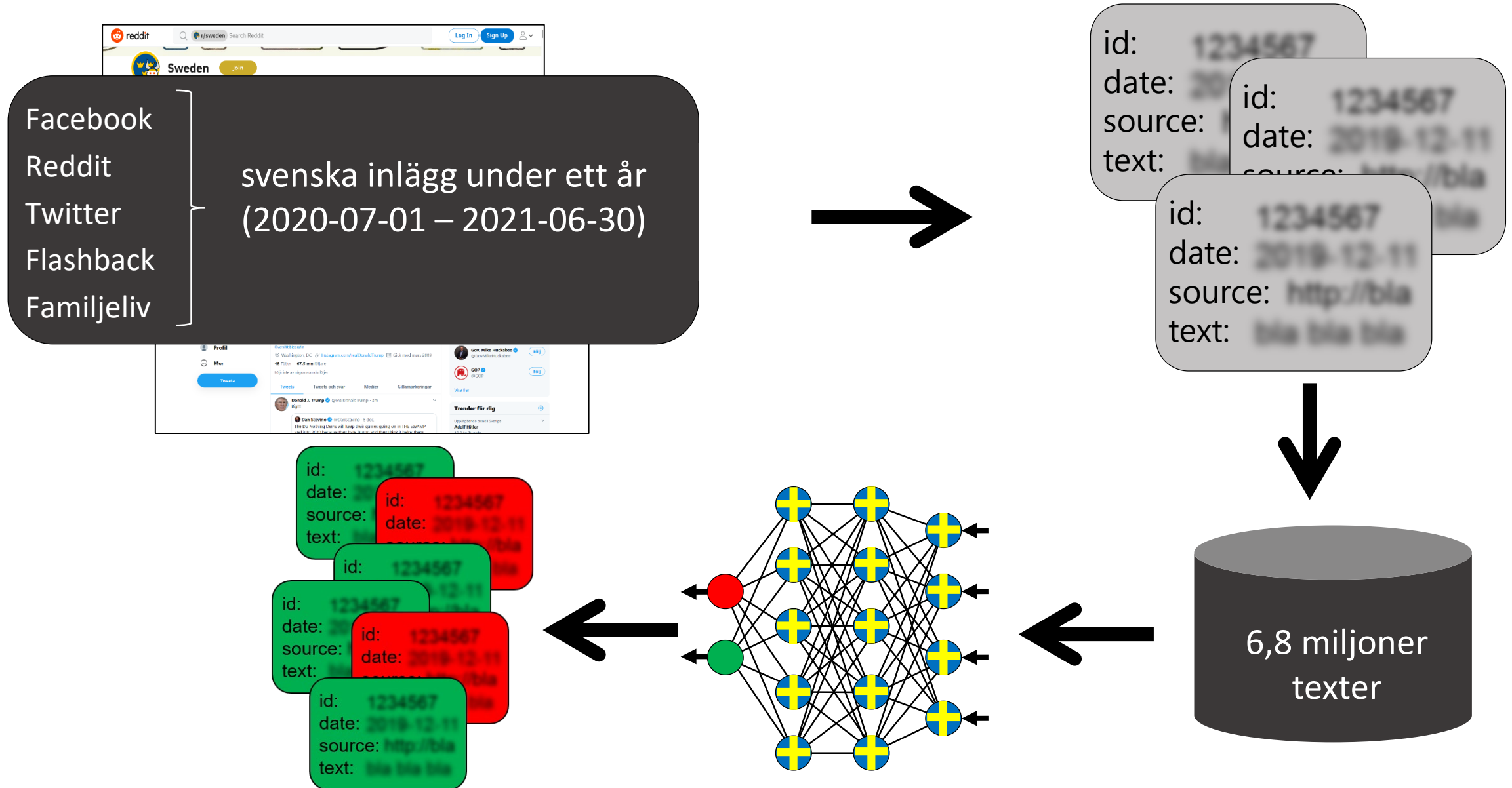
Toxiskt språk i svenska digitala miljöer



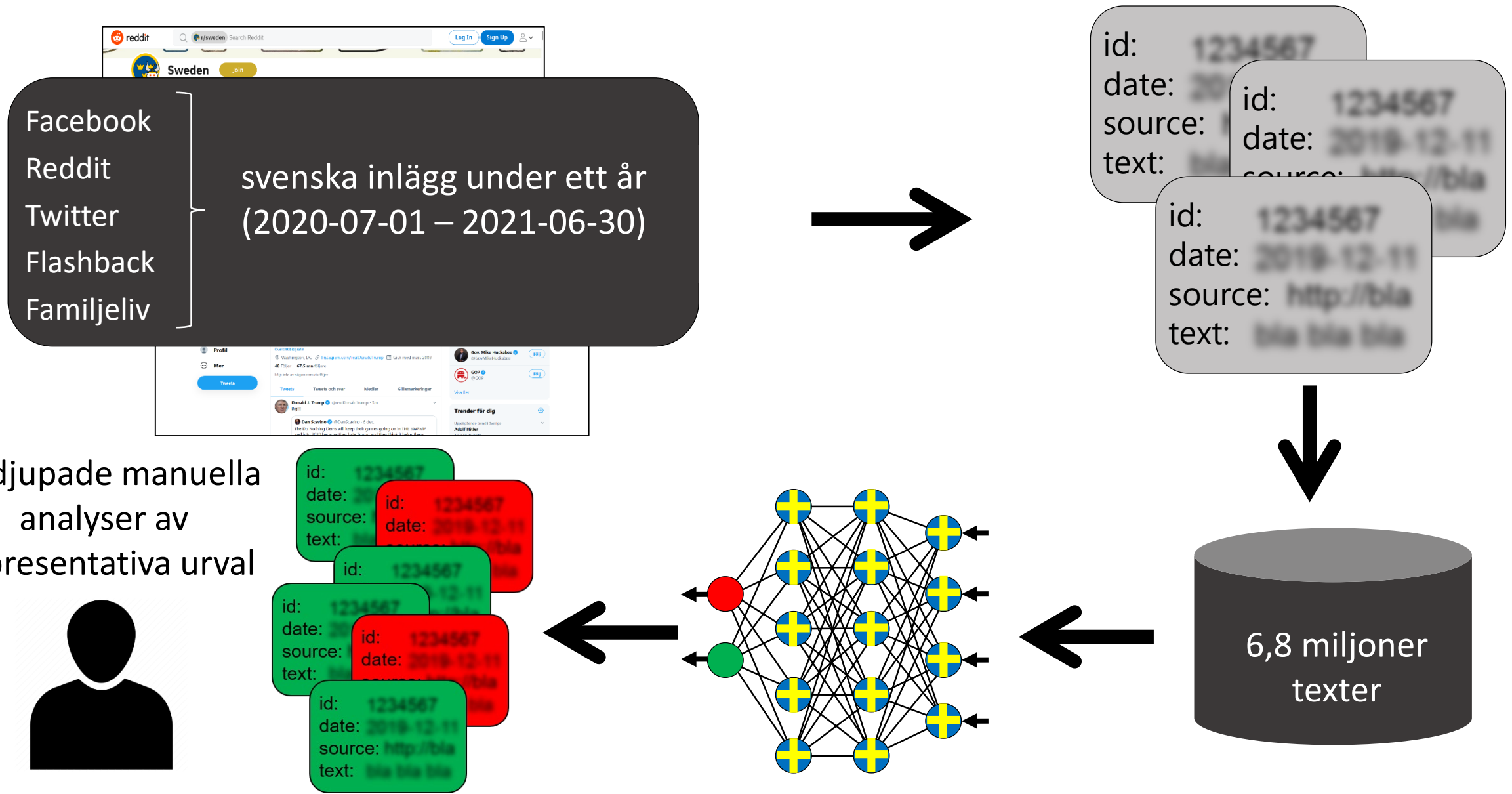
Toxiskt språk i svenska digitala miljöer



Toxiskt språk i svenska digitala miljöer

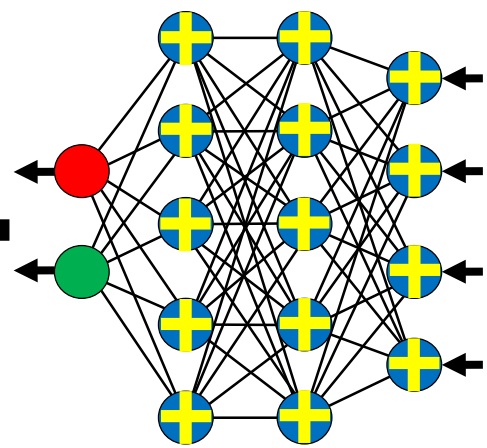
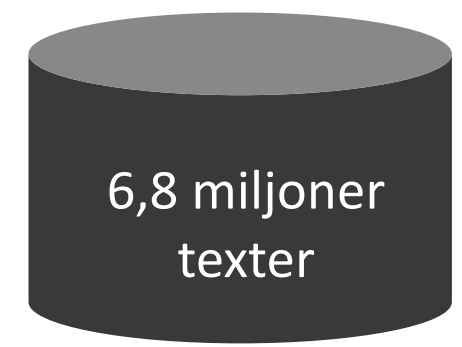
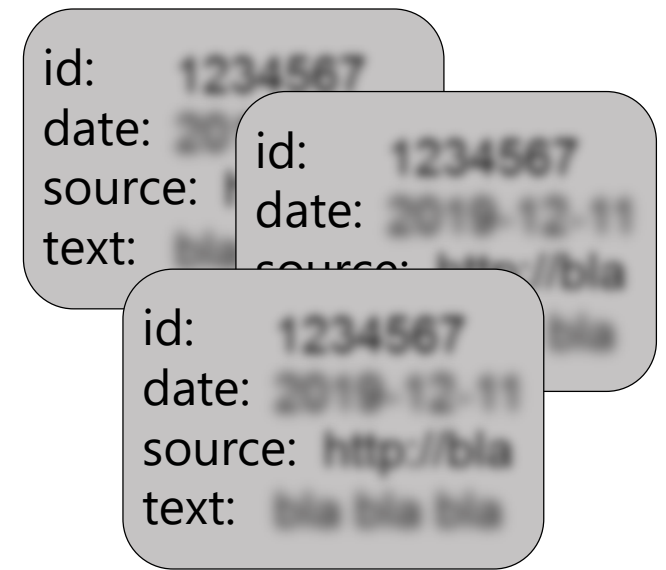
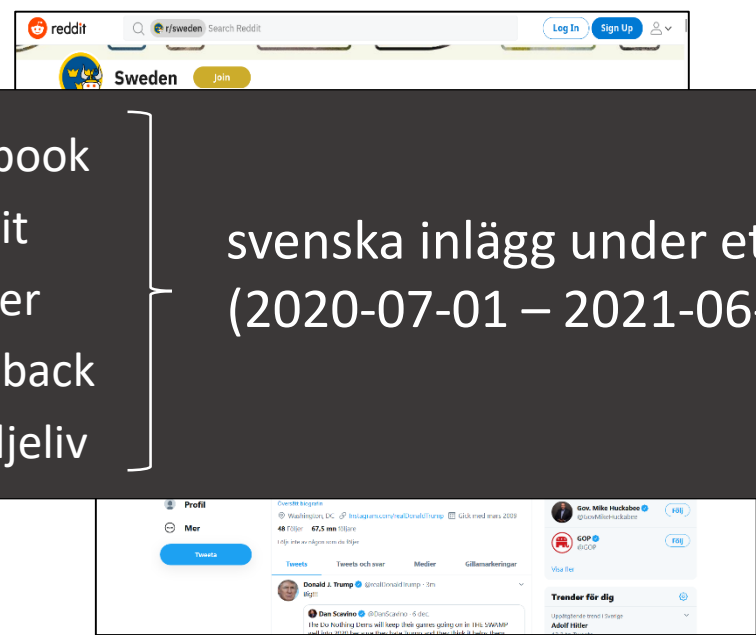


Toxiskt språk i svenska digitala miljöer

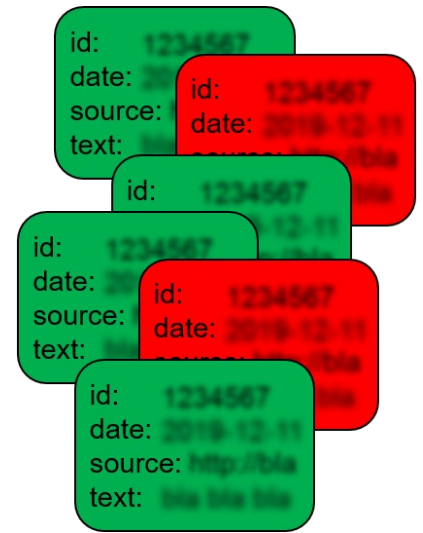
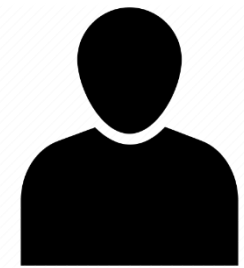


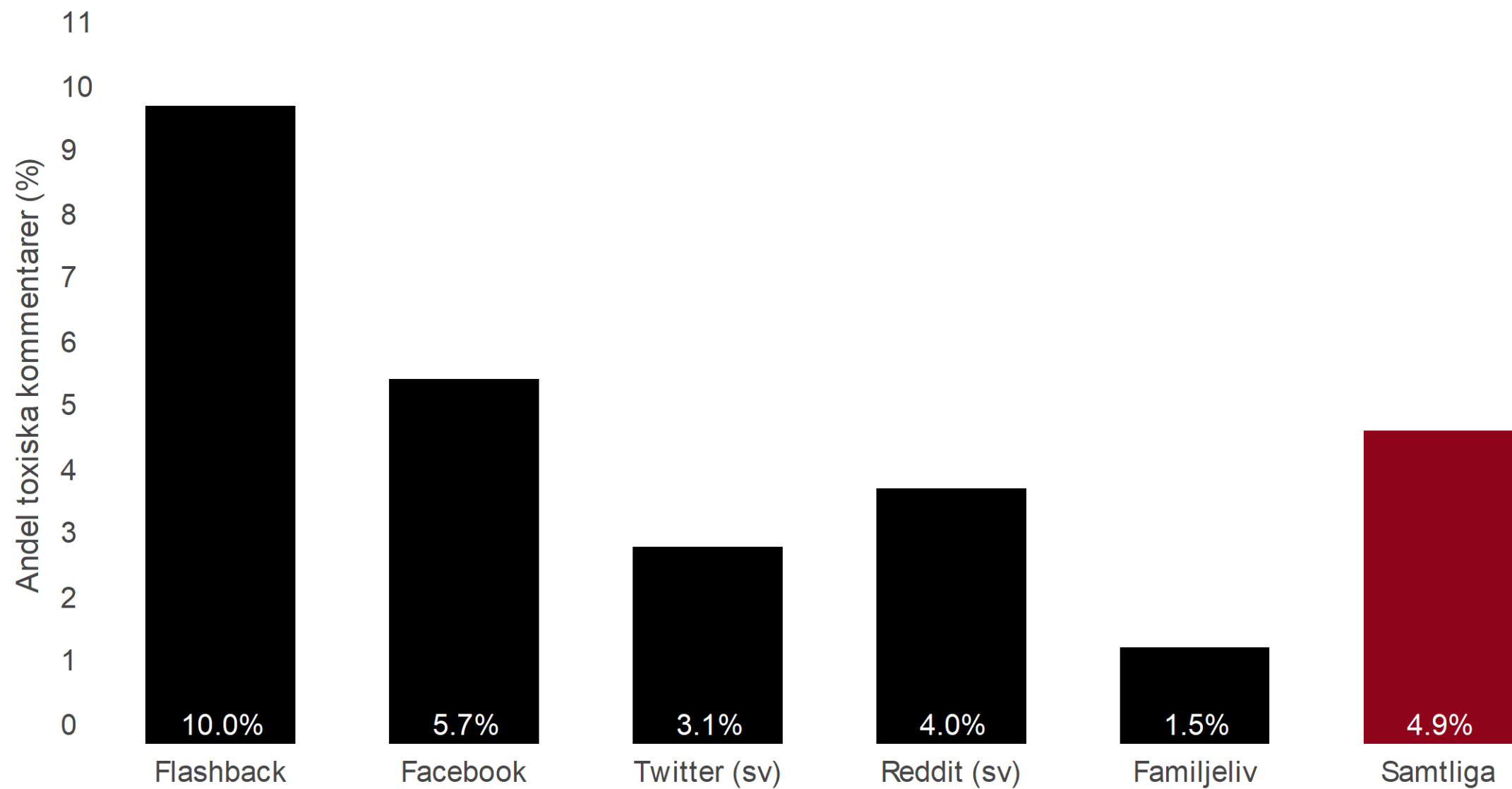
Facebook
Reddit
Twitter
Flashback
Familjeliv

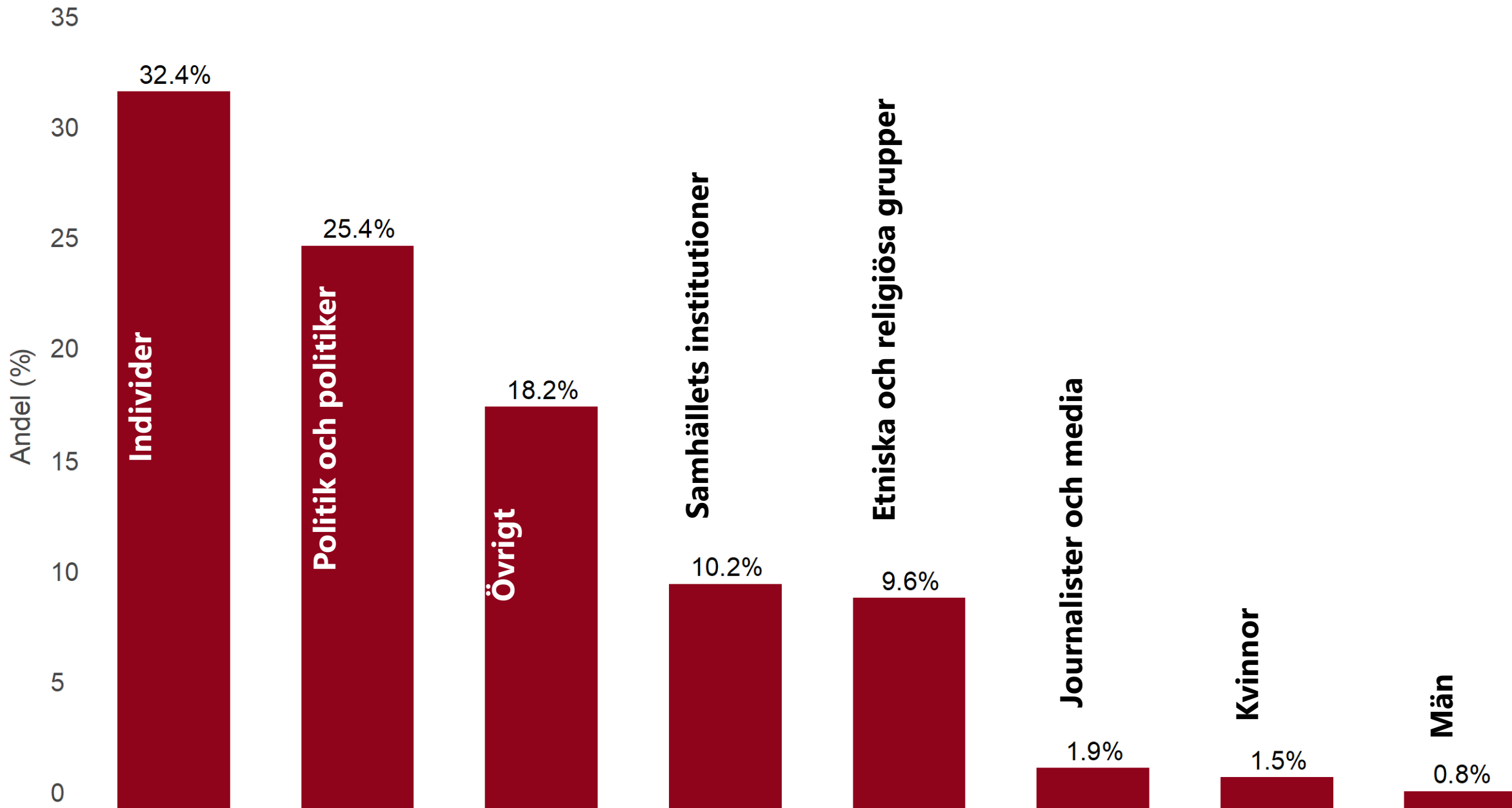
svenska inlägg under ett år
(2020-07-01 – 2021-06-30)

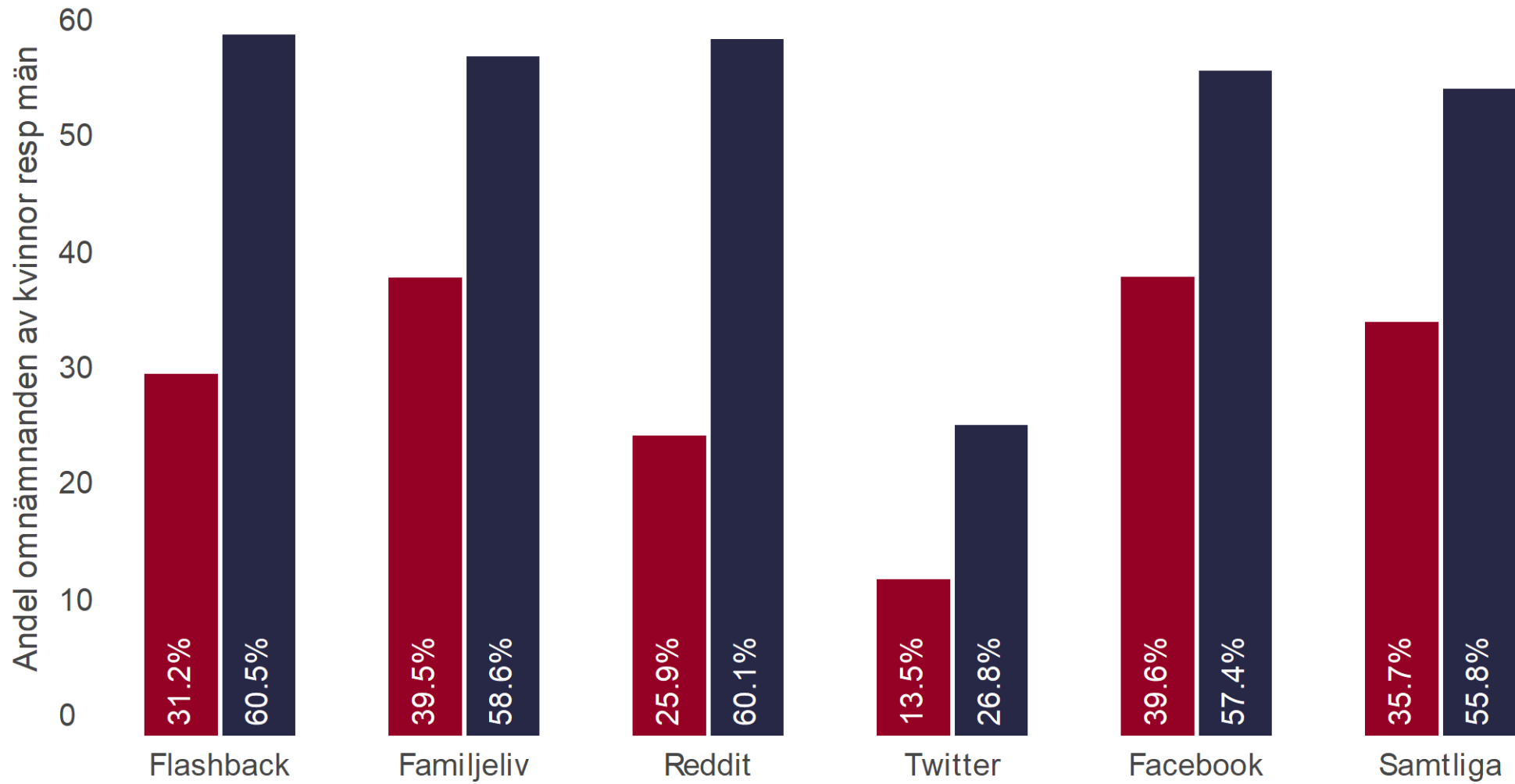


fördjupade manuella analyser av representativa urval

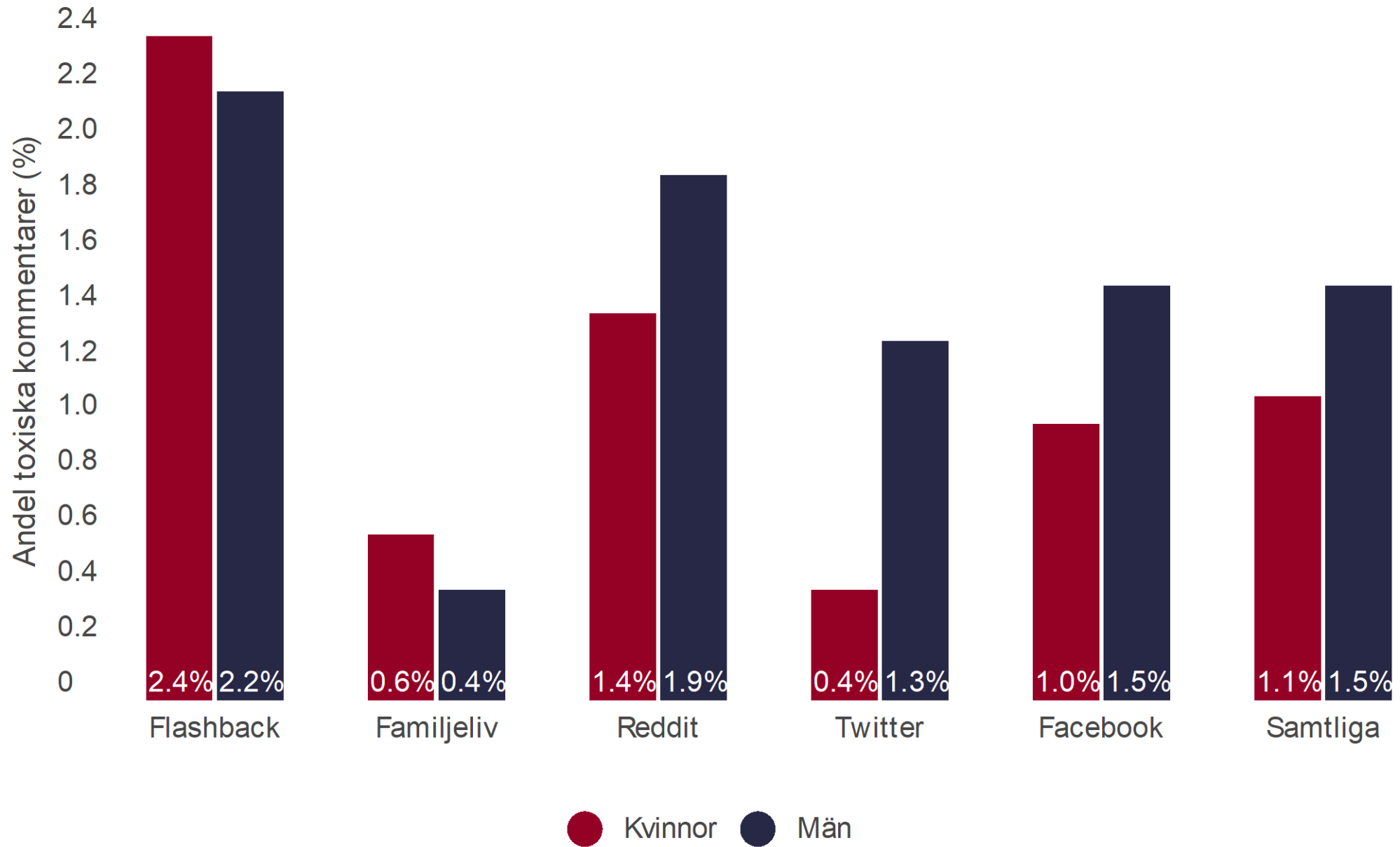


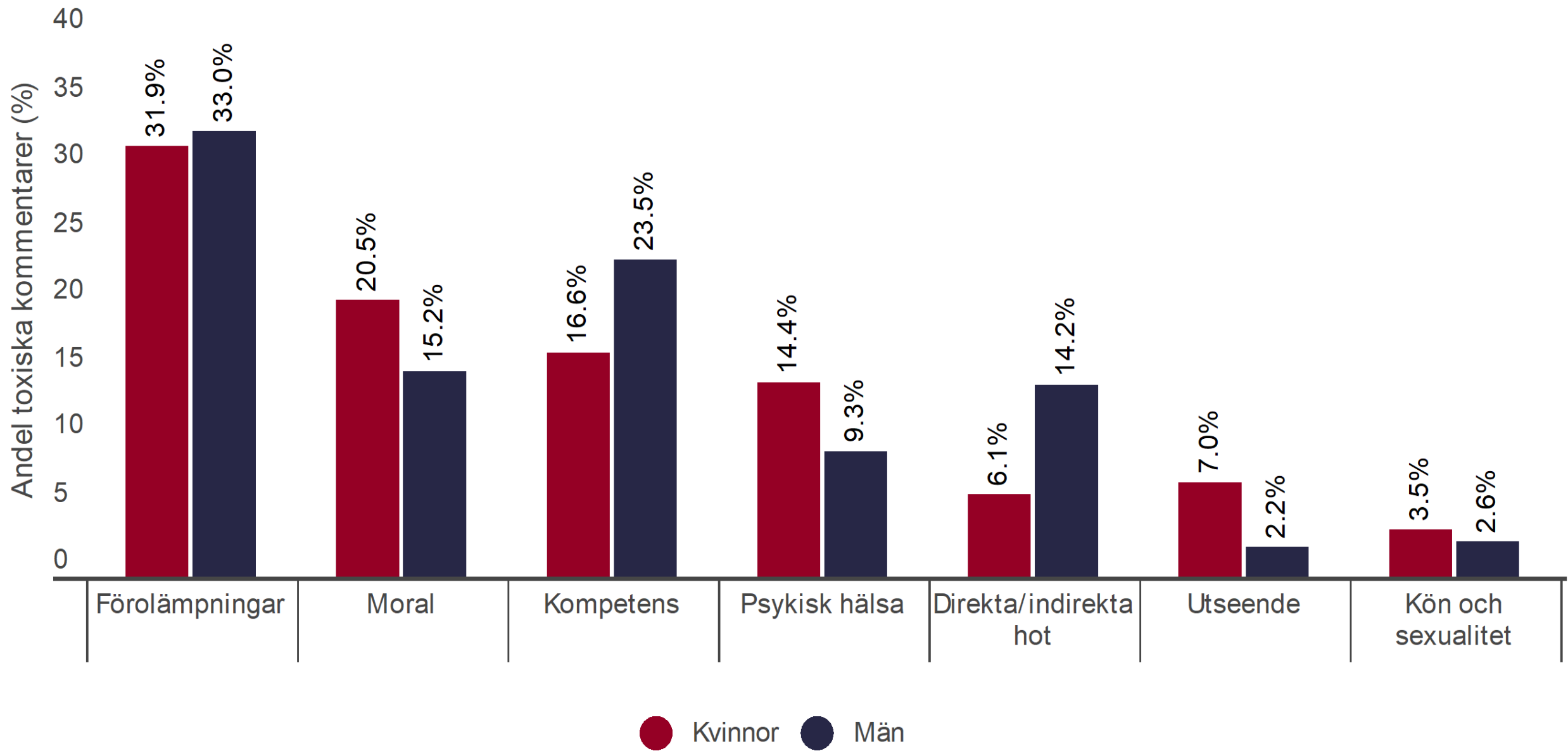






● Kvinnor ● Män





Toxiskt språk i svenska digitala miljöer

Lisa Kaati, Björn Pelzer, Katie Cohen, Daniel Wallgren, Jenny Yourstone, Nazar Akrami.

När samtalsklimatet i digitala miljöer blir infekterat av kränkande kommunikation äventyras även det demokratiska samtalet. I den här studien undersöker vi några svenskspråkiga digitala miljöer för att få en uppfattning om förekomsten av toxiskt språk, det vill säga kommunikationshandlingar som medför kränkningar mot mottagaren eller en tredje part. Vi undersöker även vilka grupper eller företeelser som de toxiska kommentarerna riktas mot samt förändring över tid i två av Sveriges största diskussionsforum.



6 800 000

Vi har analyserat kommentarer från olika sociala medieplattformar som producerats under ett års tid.



4,9%

Av alla kommentarerna innehåller toxiskt språk



32,4%

Av de toxiska kommentarerna riktas mot individer



25,4%

Av de toxiska kommentarerna handlar om politik eller om enskilda politiker.



10,2%

Av de toxiska kommentarerna riktas mot samhällets institutioner



9,6%

Av de toxiska kommentarerna riktas mot etniska och religiösa grupper



1,9%

Av de toxiska kommentarerna riktas mot media och journalister



Kvinnor som grupp utsätts nästan dubbelt så ofta för toxiska kommentarer än män

Könsskillnader i utsatthet för toxiskt språk online

Lisa Kaati, Katie Cohen, Björn Pelzer, Daniel Wallgren, Jenny Yourstone, Nazar Akrami.

Ett flertal studier har visat att kvinnor och män utsätts i ungefär lika stor omfattning för kränkande kommentarer på nätet, men oftast på olika sätt. I den här studien har vi undersökt skillnader mellan kvinnors och mäns utsatthet för toxiskt språk, liksom eventuella karaktärsskillnader på toxiskt språk riktat mot kvinnor respektive män. Undersökningen baseras på ett års data från några av de största svenskspråkiga sociala medierna.



6 800 000 kommentarer

Vi har analyserat kommentarer från olika sociala medieplattformar producerade under ett år.



Mansnamn och manliga pronomen nämns nästan dubbelt så ofta som kvinnonamn och kvinnliga pronomen i de undersökta källorna.



Psykisk hälsa

Kvinnor blir i större utsträckning än män utsatta för nedvärderande kommentarer om psykisk hälsa och förmåga.



Kompetens

Män utsätts för fler nedvärderade kommentarer om bristande kompetens eller prestation inom sin yrkesgärning eller allmänt.



Utseende

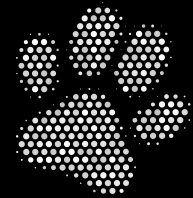
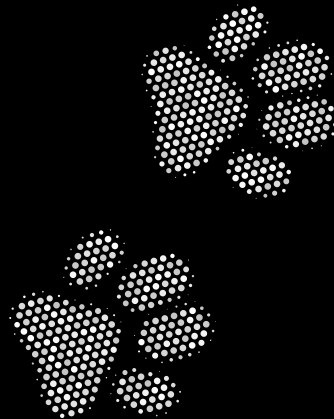
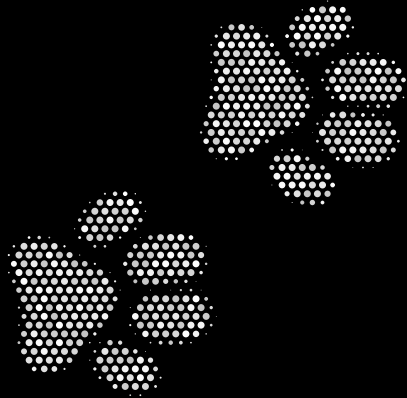
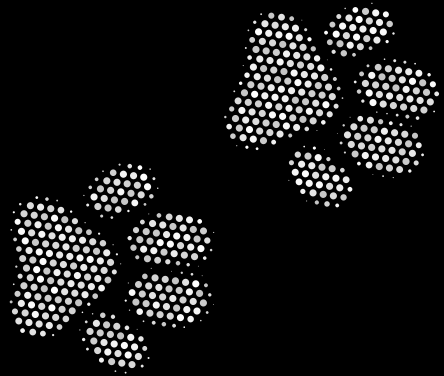
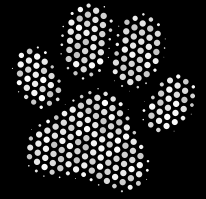
Kvinnor blir i större omfattning utsatta för nedvärderande kommentarer om utseende som innefattar påståenden om att en person är ful eller oattraktiv, eller objektifierande och närgångna kommentarer om en persons kropp oavsett värdeledning.



Hot om våld och bestraffning

Män blir i större omfattning utsatta för kommentarer som innehåller hot/vålds- och bestraffningsidéer, direkt eller indirekt formulerade avsikter eller önskingar om att en person eller grupp ska, dö, försvinna eller utsättas för våld.

Hotbedömningar i digitala miljöer



Attentatspersoner använder internet

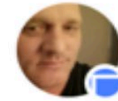


Textanalys i hotbedömningar 2012 – 2022



- Kan man stoppa en ensamagerande terrorist i tid?
- Finns det varningssignaler?
- Finns varningssignalerna på internet?
- Är textanalys användbart?

Textanalys i hotbedömningar 2012 – 2022



Robert Bowers @onedingo

2 hours ago

HIAS likes to bring invaders in that kill our people.
I can't sit by and watch my people get slaughtered.
Screw your optics, I'm going in.

(/784.30 KB 3956X1284 /3b/8c/64fcea3dc388/e3cd027ed4cf3cb...)



bump race war thread Philip Manshaus 08/10/2019 (Sat) 13:56:
[33e786 \(1\)](#) [Preview] No. 73841 [Hide User] [X] >>73844 >>73853

well cobbers it's my time, i was elected by saint tarrant after all



Alek Minassian

2 hrs · 🌐

Private (Recruit) Minassian Infantry 00010, wishing to speak to Sgt 4chan please. C23249161. The Incel Rebellion has already begun! We will overthrow all the Chads and Stacys! All hail the Supreme Gentleman Elliot Rodger!

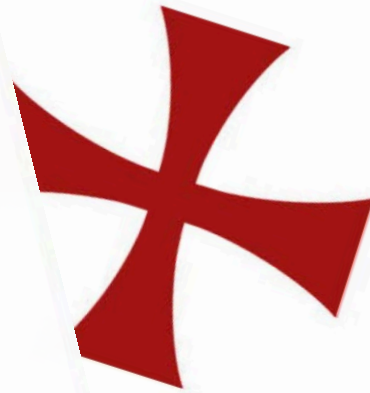
The Great Replacement

TOWARDS A NEW SOCIETY



WE MARCH EVER FORWARDS

2083



of Independence

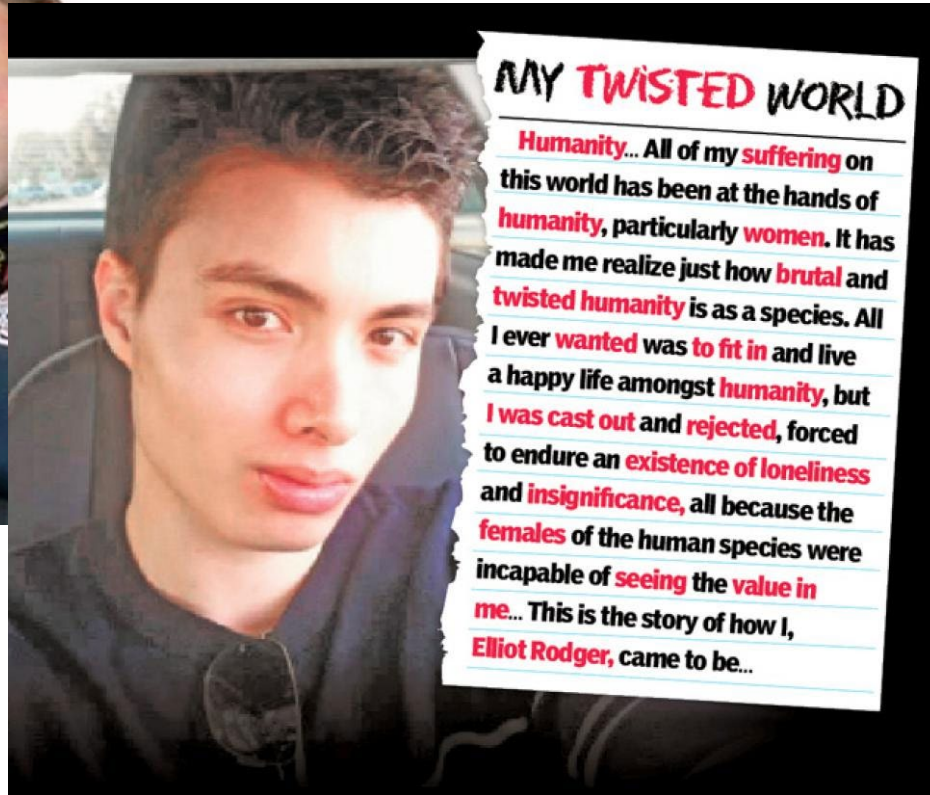
via
ique Solomonic

rick, London - 2011

Utvidgad forskningsfråga



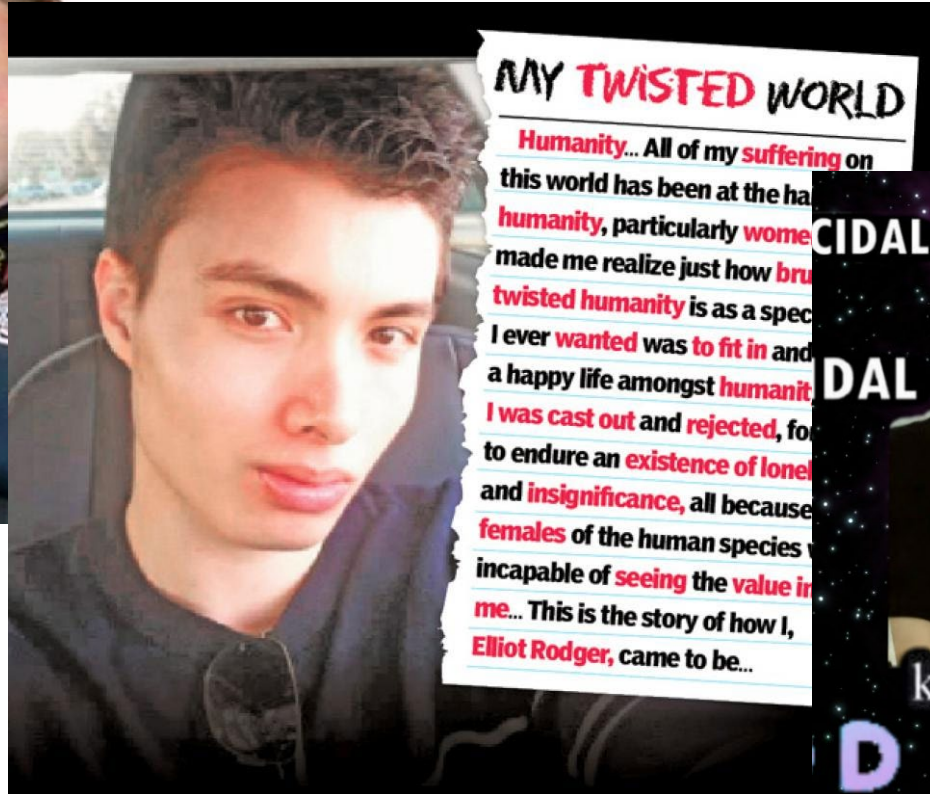
Utvidgad forskningsfråga



MY TWISTED WORLD

Humanity... All of my suffering on this world has been at the hands of humanity, particularly women. It has made me realize just how brutal and twisted humanity is as a species. All I ever wanted was to fit in and live a happy life amongst humanity, but I was cast out and rejected, forced to endure an existence of loneliness and insignificance, all because the females of the human species were incapable of seeing the value in me... This is the story of how I, Elliot Rodger, came to be...

Utvidgad forskningsfråga



MY TWISTED WORLD

Humanity... All of my suffering on this world has been at the hands of humanity, particularly women. This made me realize just how brutal and twisted humanity is as a species. I ever wanted was to fit in and have a happy life amongst humanity. I was cast out and rejected, forced to endure an existence of loneliness and insignificance, all because of the females of the human species who are incapable of seeing the value in me... This is the story of how I, Elliot Rodger, came to be...

ICIDAL

DAL

MONSTERS
HUMAN...

kill me now

D BYE

We will be free

WEIRD
IS
RAD

NIN



I LOST MYSELF

Varför datoriserade analyser?

- Sparar mänskliga resurser
- Skyddar integritet
- Möjlighet att se mönster som ett mänskligt öga inte kan överblicka

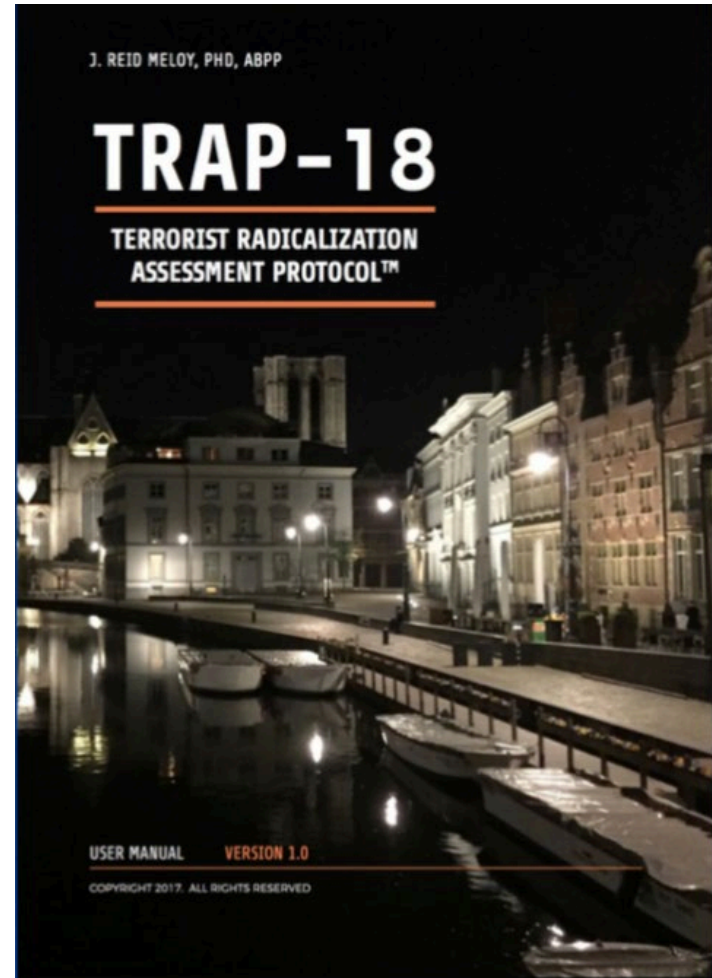
Datoriserade analyser kräver människor för:

- 1) rätt frågeställning
- 2) experttolkning



TRAP -18

- Strukturerad professionell bedömning
- Bedömer förstagångsförbrytare
- Ekologisk validitet, kan även bedöma skolskjutare och stalkers
- Baserat på 10 riskfaktorer och 8 beteendeindikatorer



Risikfaktorer och indikatorer



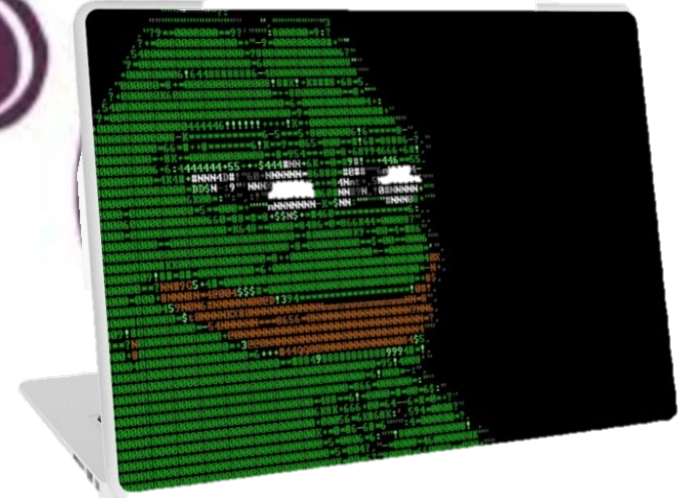
Beteendeindikatorer i digitala miljöer

Indikatorerna i TRAP-18 är tänkta att observeras hos individer som fysiskt är under bevakning, men...



Beteendeindikatorer i digitala miljöer

...flera av beteendeindikatorerna kan observeras där potentiella gärningspersoner ofta finns: i digitala miljöer.



Indikatorer i TRAP-18

- Förberedelse
 - Testkörning
 - Ökad aktivitet
 - Identifikation
 - Fixering
 - Tunnelseende
 - Läckage
 - Explicit avsikt
- } konkreta handlingar
- } psykologiska tillstånd
- } (oftast) verbala handlingar

Att observera psykologiska tillstånd...

- Identifikation
 - Fixering
 - Tunnelseende
- } psykologiska tillstånd

...i skriven text

- Sentimentanalys – maskininlärningsmodeller
- Komplexitet i resonemang
 - Läsbarhetsmått (t ex LiX eller Flesch-Kincaid)
 - Ordlistebaserad approach (t ex LIWC)

Exempel: Identifikation

- Ser sig själv som en krigare/militär
- Ser sig själv som nyckelperson i ideologisk kamp
- Ser sig själv som efterträdare till känd våldsverkare



Exempel: Identifikation

- Ser sig själv som en krigare/militär
Vapen, militär terminologi
- Ser sig själv som nyckelperson i ideologisk kamp
Uttryck relaterade till handling, ex "göra något"
- Ser sig själv som efterträdare till kända våldsverkare
**Tar efter någon annans språkbruk
"saint"/"chad" etc**



Vadå saint och chad?



Saint

Chad Brenton Tarrant



Chad

LIWC – Linguistic inquiry and word count

- Ordräkningsbaserat verktyg för stora textmängder
- Ursprungligen utvecklad av psykologer för att mäta effekt av terapeutiskt skrivande
- Bygger på forskning om hur ordbruk korrelerar med kön, ålder, personlighet, utbildning, känslomässiga tillstånd...
- Ordräkningsbaserade verktyg ger mycket information men har begränsningar:

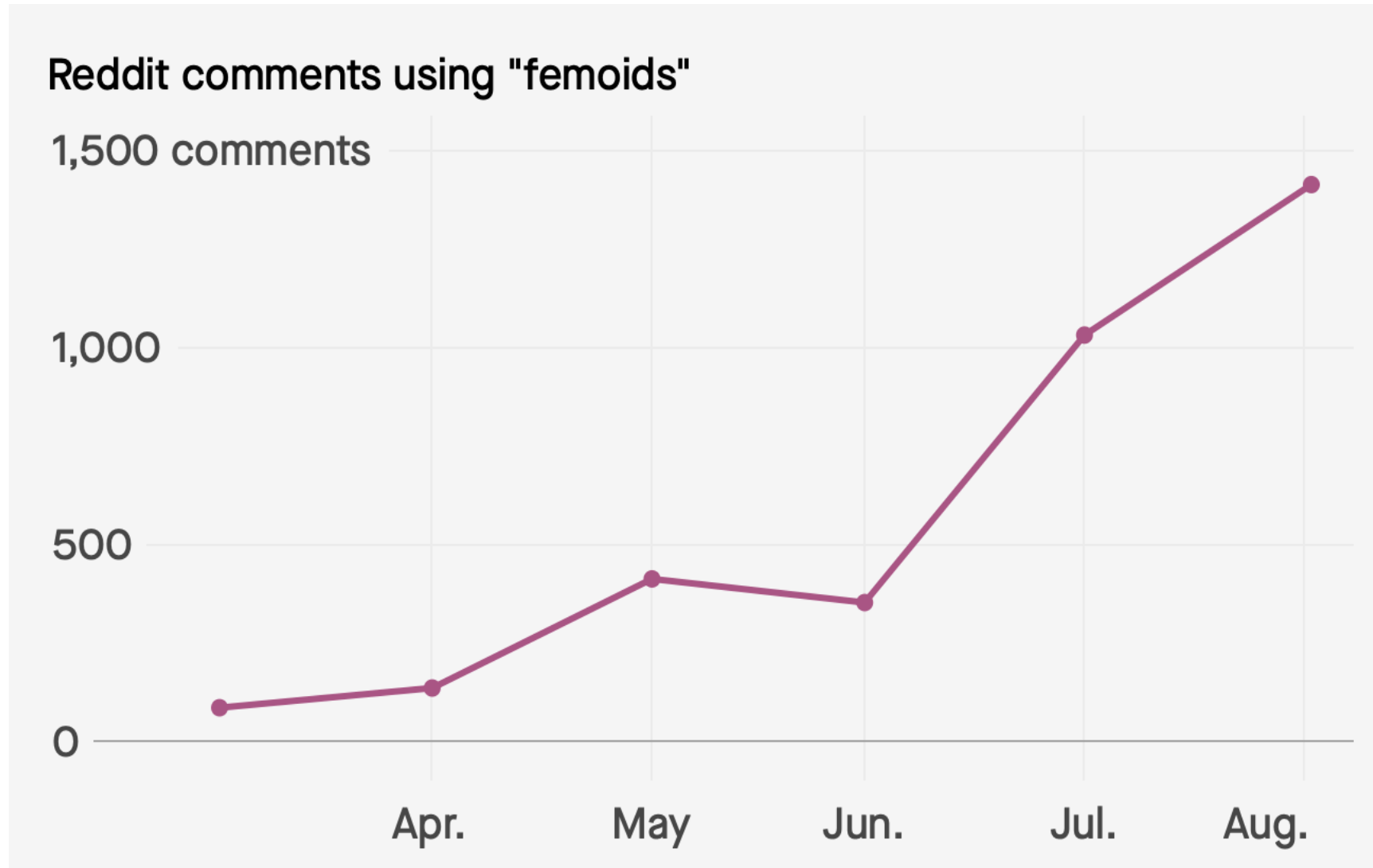
Kontext och jargong

C. What the fuck did you just fucking say about me, you little bitch? I'll have you know I graduated top of my class in the Navy Seals, and I've been involved in numerous secret raids on Al-Quaeda, and I have over 300 confirmed kills. I am trained in gorilla warfare and I'm the top sniper in the entire US armed forces. You are nothing to me but just another target. I will wipe you the fuck out with precision the likes of which has never been seen before on this Earth, mark my fucking words. You think

Navy seal cypypasta meme



Lingvistiska normer förändras



Anpassade semi-automatiska analysverktyg

2083



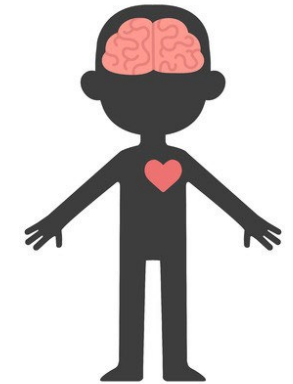
A European Declaration of Independence

De Laude Novae Militiae
Pauperes communitates Christi Templique Solomonici

By Andrew Berwick, London - 2011

The Great Replacement

TOWARDS A NEW SOCIETY



Tack!

Lisa.kaati@dsv.se
katie.cohen@foi.se