

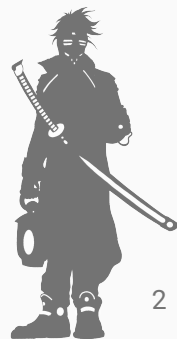
# Neuro-symboliczna sztuczna inteligencja w walce z cyberprzemocą w Internecie

dr inż. Michał Marcińczuk (Samurai Labs)  
[michal.marcinczuk@samurailabs.ai](mailto:michal.marcinczuk@samurailabs.ai)

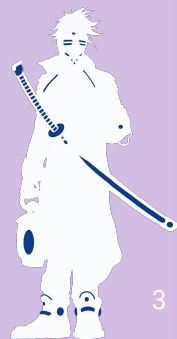
MLGdańsk #126  
6 czerwca 2022 r.



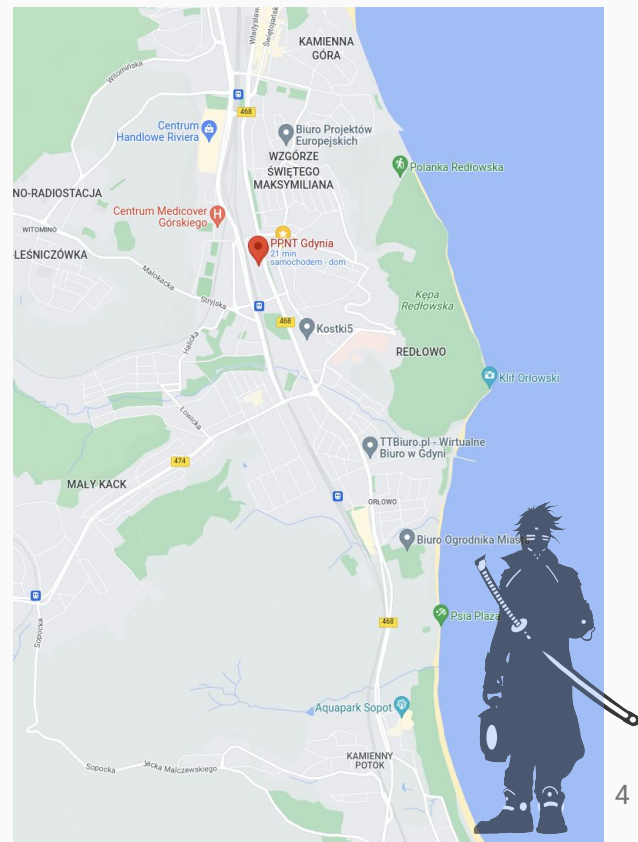
1. O nas
2. Główne usługi
3. Badania
4. Metody symboliczne + maszynowe uczenie



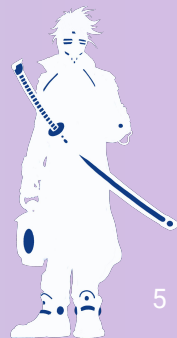
# Wprowadzenie

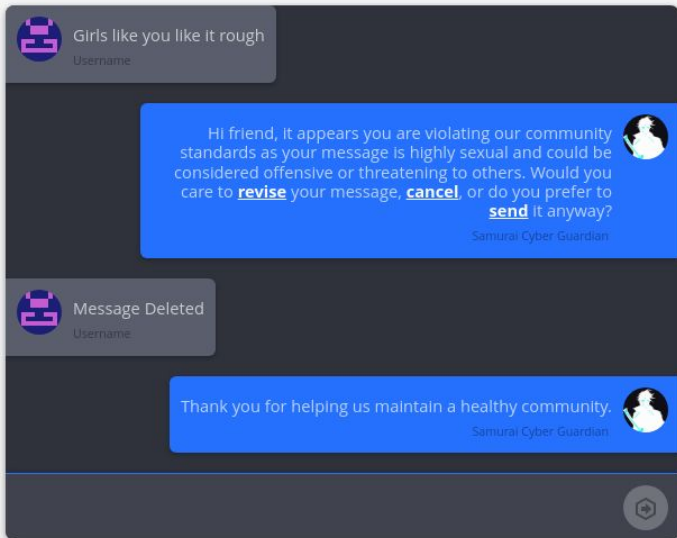


- **Lokalizacja** — Gdynia (PPNT)
- **Misja** — “Wierzimy, że Internet powinien być bezpieczny dla każdego. Naszą misją jest ochrona dzieci przed przemocą. Jesteśmy pionierami technologii, która sprawia, że internet będzie bezpiecznym miejscem przez walkę z cybernękaniem, prześladowaniem i przemocą internetową.”
- **Interdyscyplinarny zespół** — sztuczna inteligencja, kognitywistyka, lingwistyka, pedagogika, psychologia, matematyka, informatyka, biznes.
- **Obszary zainteresowania** — analiza treści pod kątem wykrywania: [cyberprzemocy](#), [mowy nienawiści](#), [treści suicydalnych](#), [pedofilii](#).



# Główne usługi





76% of cyberbullies voluntarily change their tone after being educated and redirected by Samurai's Cyber Guardian

## Blackmail Detection Example:

*Send me more nudes or I'll publish the ones I already have*



Request

Alternative

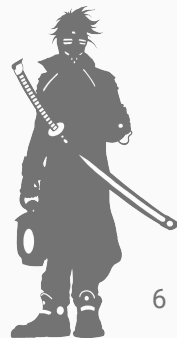
Negative Consequences

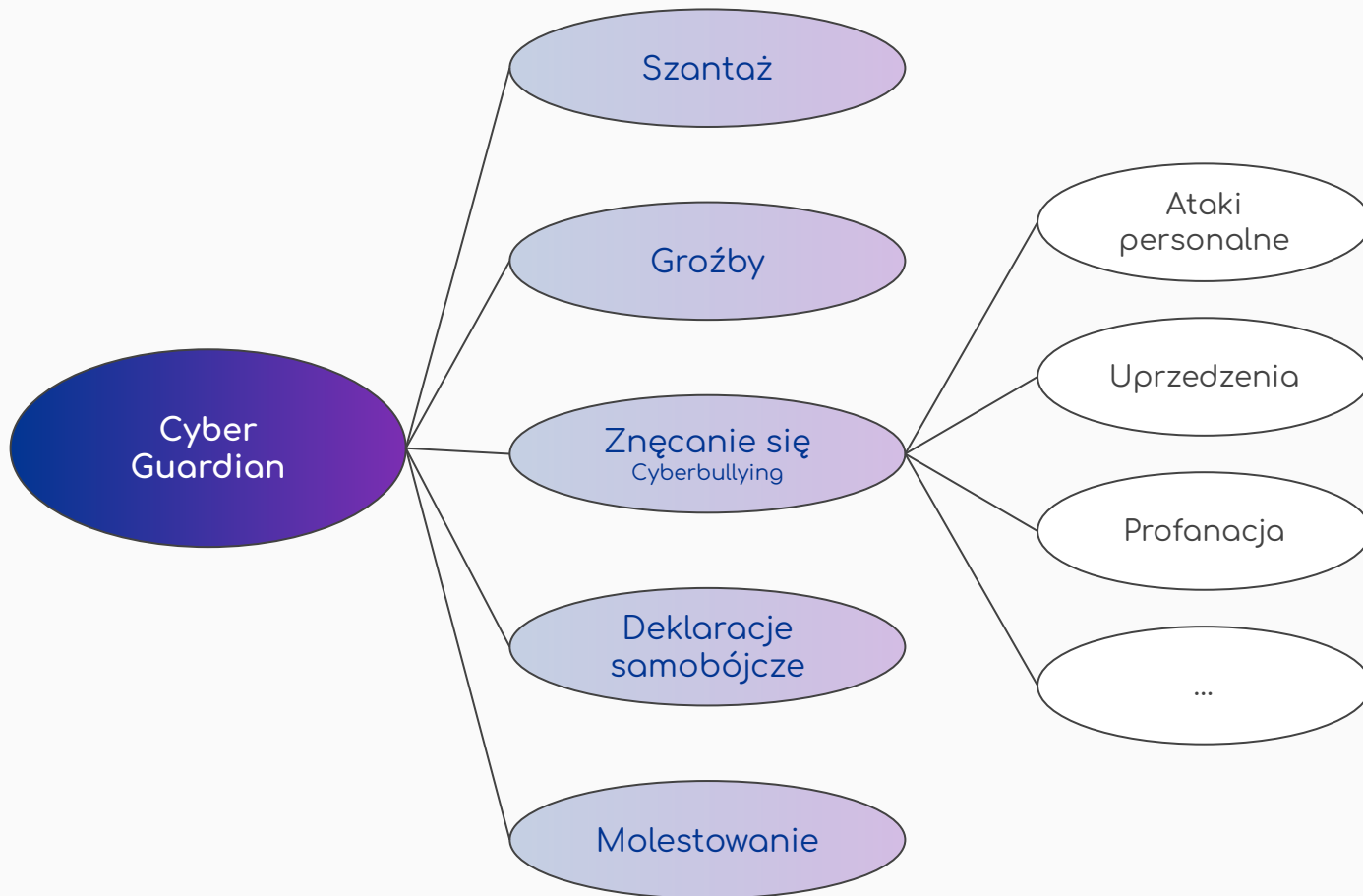
Samurai Cyber Guardian Detection

**Blackmail** threat to reveal information

**Sexual harassment** attempt to solicit intimate photographs

<https://username.samurailabs.ai/#cyberguardian>

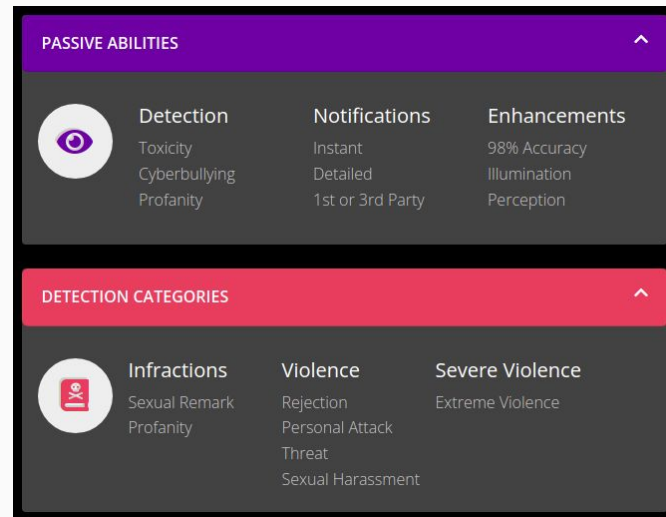
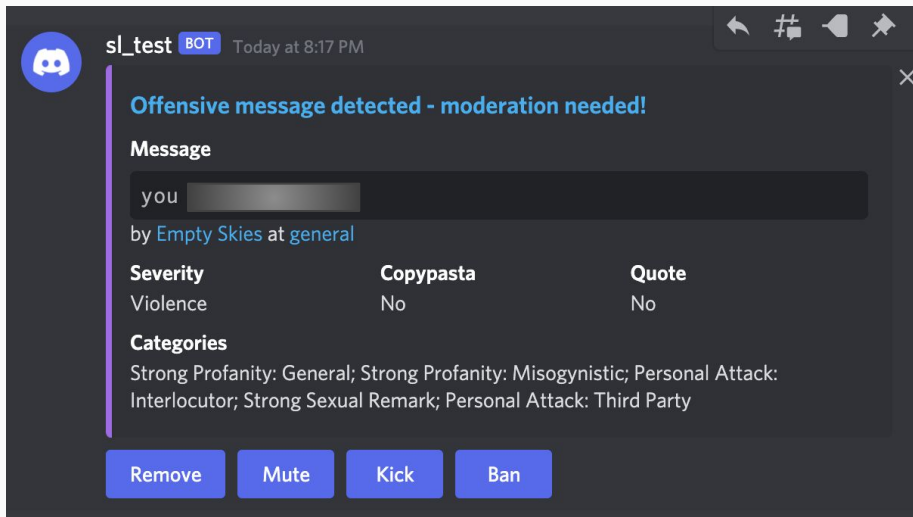




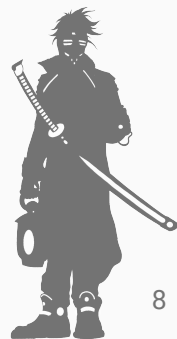
**Podział ze względu na odbiorcę:** druga osoba, trzecia osoby, grupy mniejszościowe, politycy, itd.

**Podział ze względu na siłę:** poważne, silne, łagodne





<https://s4m.ai>





## Setting the Tone for Protecting Your Communities

### Trust & Safety Starts with Username Moderation

The most effective solution for guarding your communities against users with disruptive identities. Built for automation. Scalable for any application.

TRY IT NOW

LEARN MORE

Please input your username.

**daehDcixot**

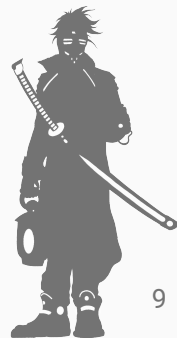
This username violates our community standards.

Decomposition: **TOXIC** **DHEAD**

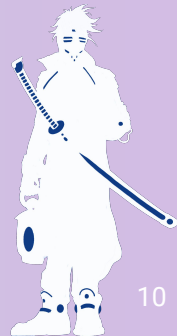
Categories: **INAPPROPRIATE** **PROFANITY**

Toxic Elements: **TOXIC** **DHEAD**

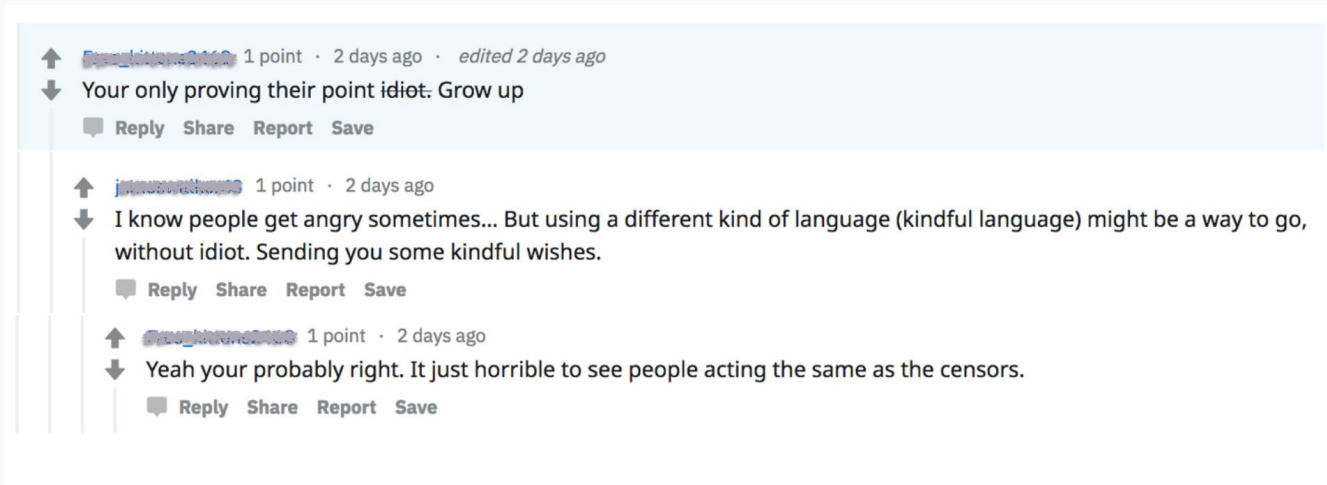
<https://username.samurailabs.ai>



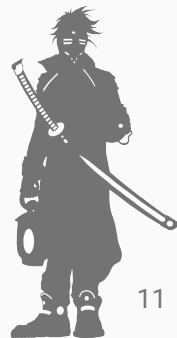
# Badania



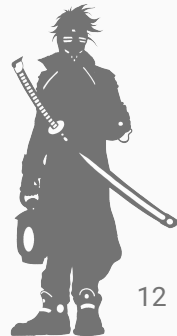
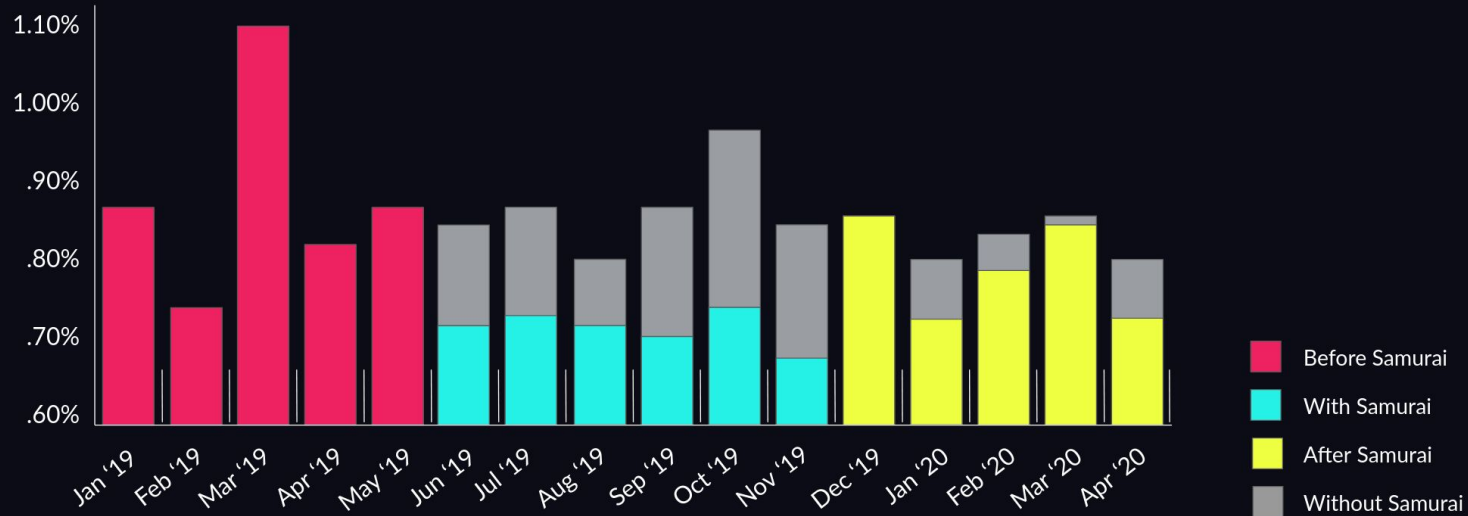
- Bot działający na redditcie **r/MensRights**
- Wykrywał i reagował na przemoc słowną
- Reakcją była odpowiedź na post, której celem było “ostudzenie” emocji



Michał Bilewicz, Patrycja Tempska, Gniewosz Leliwa, Maria Dowgiałło, Michalina Tańska, Rafał Urbaniak, Michał Wroczyński  
[Artificial intelligence against hate: Intervention reducing verbal aggression in the social network environment](https://onlinelibrary.wiley.com/doi/abs/10.1002/ab.21948), Aggressive Behavior.  
2021;47:260–266



% of comments  
w/ violent language



- Treści suicydalne na Reddicie

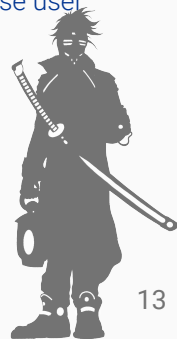
Michał Ptaszynski, Monika Zasko-Zielinska, Michał Marcinczuk, Gniewosz Leliwa, Marcin Fortuna, Kamil Soliwoda, Ida Dziublewska, Olimpia Hubert, Paweł Skrzek, Jan Piesiewicz, Paula Karbowska, Maria Dowgiałło, Juuso Eronen, Patrycja Tempa, Maciej Brochocki, Marek Godny and Michał Wroczynski. [Looking for Razors and Needles in a Haystack: Multifaceted Analysis of Suicidal Declarations on Social Media—A Pragmalinguistic Approach](https://www.mdpi.com/1660-4601/18/22/11759), Int. J. Environ. Res. Public Health 2021, 18(22), 11759

<https://www.mdpi.com/1660-4601/18/22/11759>

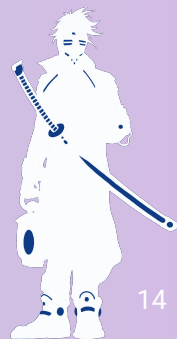
- Wpływ ataków personalnych na aktywność użytkowników Reddita

Rafał Urbaniak, Michał Ptaszyński, Patrycja Tempa, Gniewosz Leliwa, Maciej Brochocki, Michał Wroczynski. [Personal attacks decrease user activity in social networking platforms](https://www.sciencedirect.com/science/article/abs/pii/S0747563221002958). Computers in Human Behavior Volume 126, January 2022, 106972

<https://www.sciencedirect.com/science/article/abs/pii/S0747563221002958>



# Modele symboliczne i maszynowe uczenie



Metody  
symboliczne



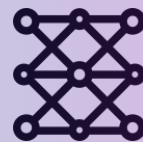
## Zalety

- Pełna kontrola nad działaniem modelu,
- Przewidywalne zachowanie,
- Wsparcie dla wyjaśniania decyzji,
- Możliwość dostrajania modeli

## Wady

- Czasochłonny proces budowy systemu,
- Ryzyko utraty przejrzystości reguł ze względu na ich wzrastającą złożoność,

Maszynowe  
uczenie

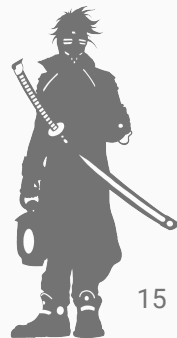


## Zalety

- Relatywna "łatwość" budowy modelu,
- Możliwości automatyzacji procesu,
- W przypadku pretrenowanych modeli dostęp do wiedzy ogólnej wykraczającej poza zbiór uczący,

## Wady

- Brak lub ograniczone wsparcie dla wyjaśniania decyzji, szczególnie w przypadku głębokich sieci neuronowych,



Please input your username.

Ash Hull

Please input your username.

Ash Hull

This username violates our community standards.

Dlaczego?

Please input your username.

Ash Hull

This username violates our community standards.

English:

OFFENSIVE

PROFANITY

Due to similarity to "asshole"

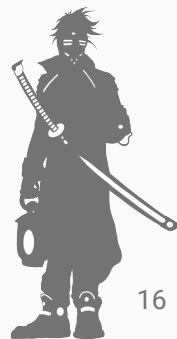


## Detailed output

Decomposition of the username - DigBildo - Big D\*Ido

Categories - Offensive, Sexual

Toxic Element - D\*Ido





Przykład dla modelu do wykrywania gróźb

“ok then we will kill those black nationalists, terrorists and pedophile, but you will need to die with them”



groźba

Dlaczego?



poważna groźba względem 2 osoby

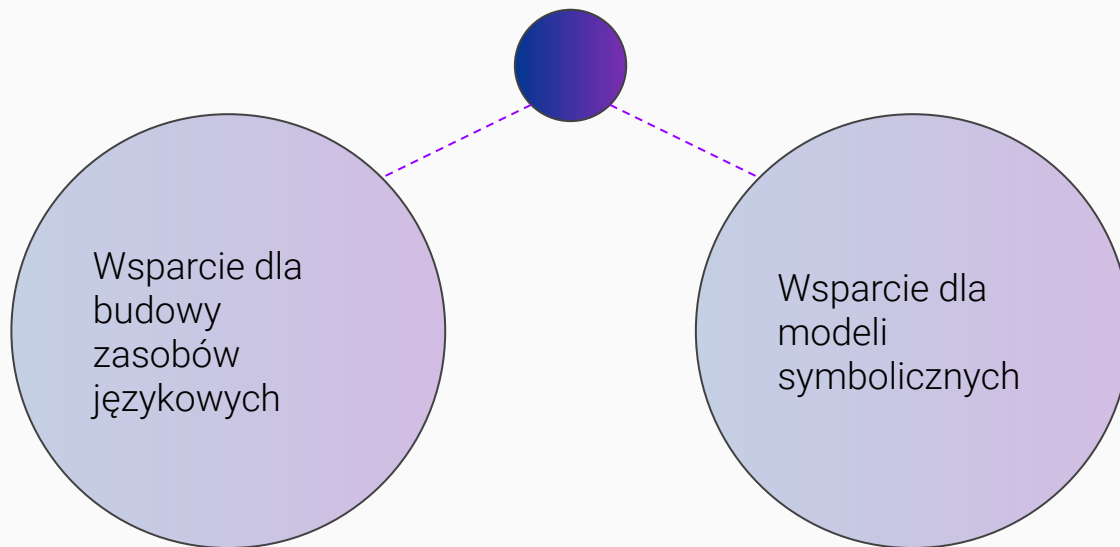
**you will need to die with them**

@you @future @must @threat\_verbs\_die\_suffer  
@preposition @third\_object\_pronoun

[14:15]	you
[15:16]	future
[16:18]	must
[18:19]	threat_verbs_die_suffer
[19:20]	preposition
[20:21]	third_object_pronoun

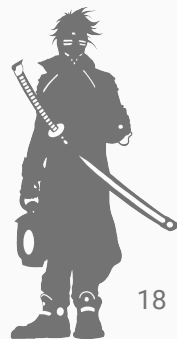
you
will
need to
die
with
them



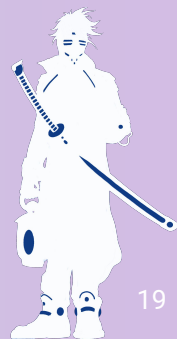


- Odkrywanie danych do anotacji;
- Rozszerzanie zasobów językowych;
- Weryfikacja ręcznej anotacji.

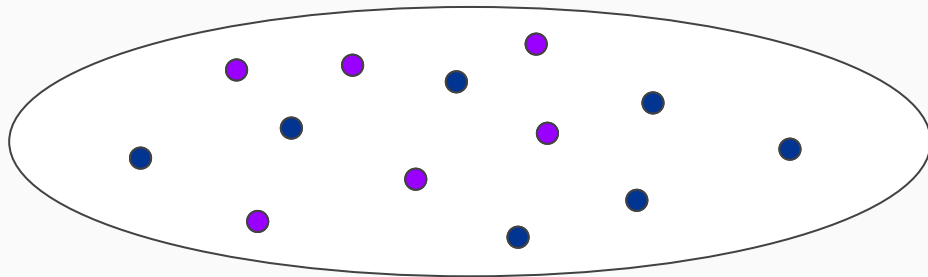
- Realizacja określonych zadań na potrzeby modelu symbolicznego;
- Filtrowanie wyników modelu symbolicznego.



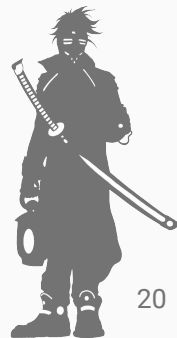
# Odkrywanie danych



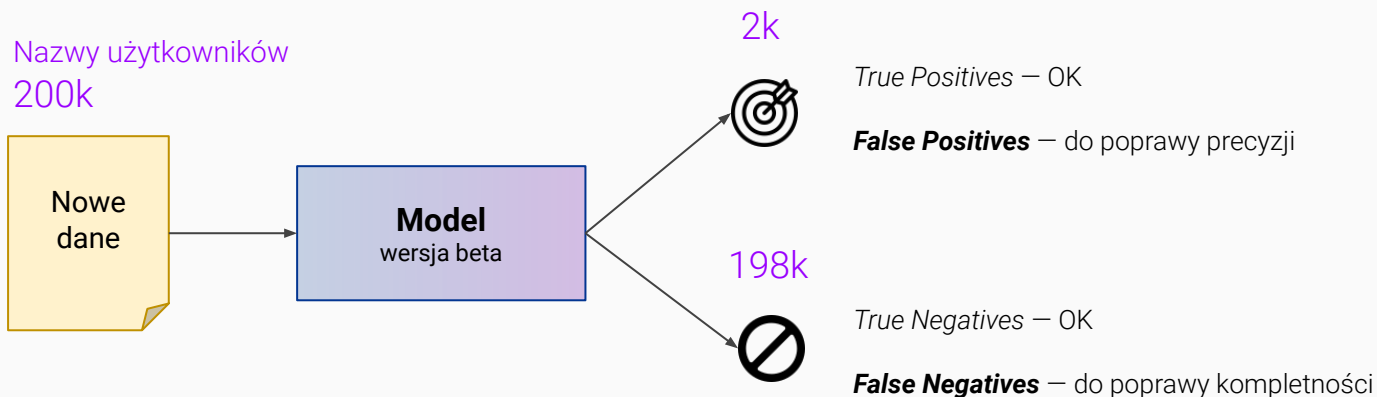
## Budowa zbioru referencyjnego – preferowane (idealne) podejście



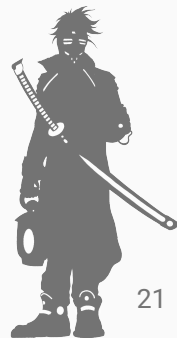
- **Zaleta** – rzeczywiste odzwierciedlenie populacji
- **Wada** – problem z danymi o nierównomiernej dystrybucji klas
  - zjawiska, które chcemy wykrywać są rzadkie, np. 1 na 100/1000,
  - wydłuża to proces anotacji,
  - z czasem coraz trudniej znaleźć nowe, unikalne przykłady.



## Procedura pracy z nowymi danymi

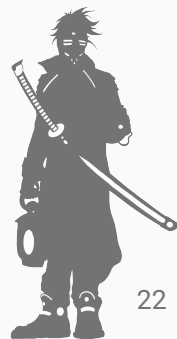


Tylko **1 na 100** nazw jest toksyczna



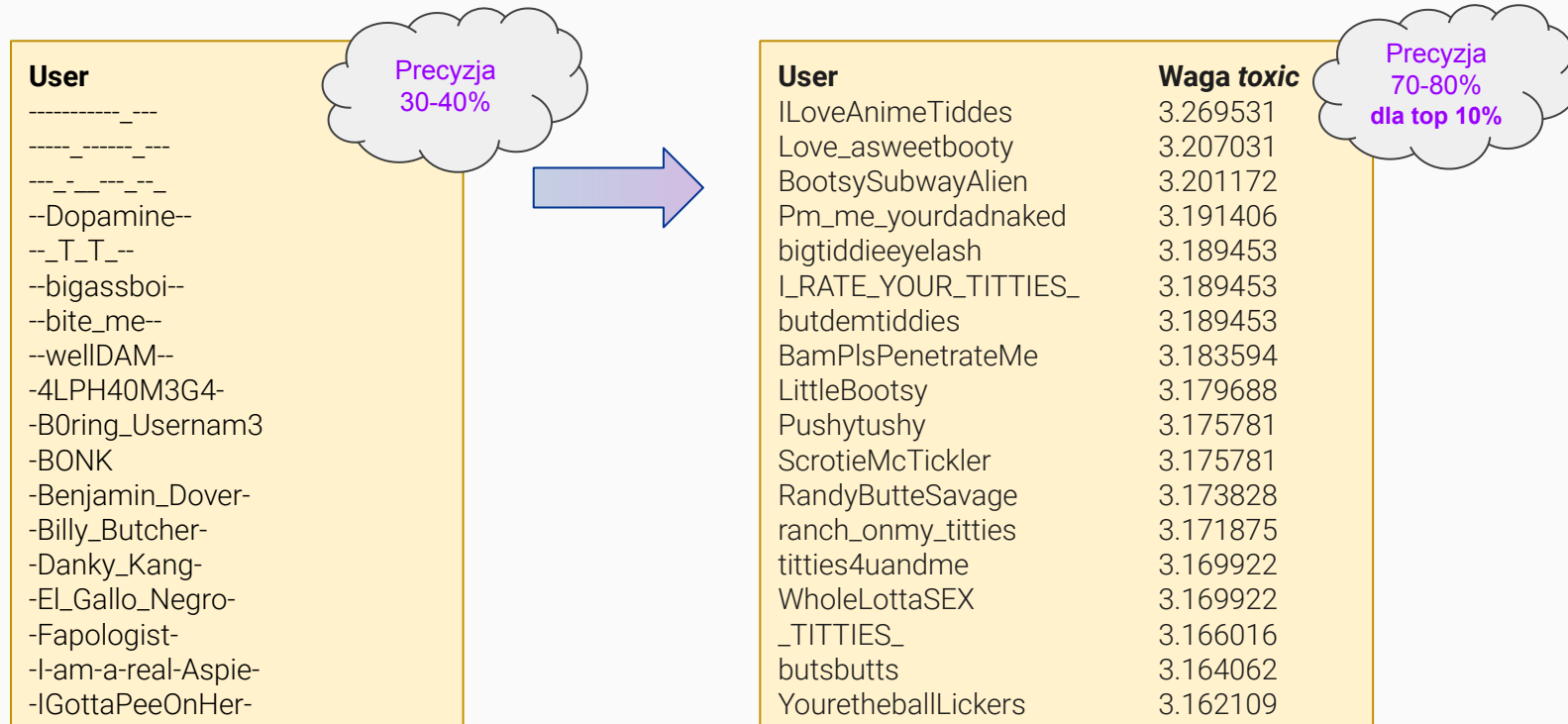
## Wykorzystanie pretrenowanych modeli ML

- **Modele MLM** (masked-language model) uczone w sposób nienadzorowany na dużych kolekcjach danych — *BERT*, *RoBERTa*, *DistilBERT*, *BigBird*, *Muppet*, itp. **potrafią uchwycić wiedzę ogólną, znaczenie sekwencji znaków (słów).**
- **Docelowy model ML** jest dostrajany na danych referencyjnych dla danego zadania.
- **Przykład:**
  - **Zadanie:** wykrywanie toksycznych nazw użytkowników
  - **Zbiór referencyjny:** ~25k toxic, ~60k non-toxic,
  - **Wynik na eval:** **precyzja i kompletność powyżej 80%,**
  - **Wynik na nowym zbiorze:** **precyzja 38%, kompletność 80%,**
  - **Częstość toxic w próbce:** **17%** vs **1%** — **17 x** więcej istotnych nowych przykładów.

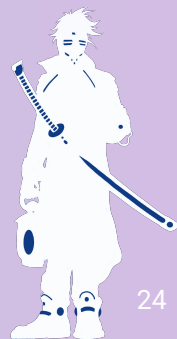


## Wykorzystanie pewności klasyfikacji do sortowania wyników

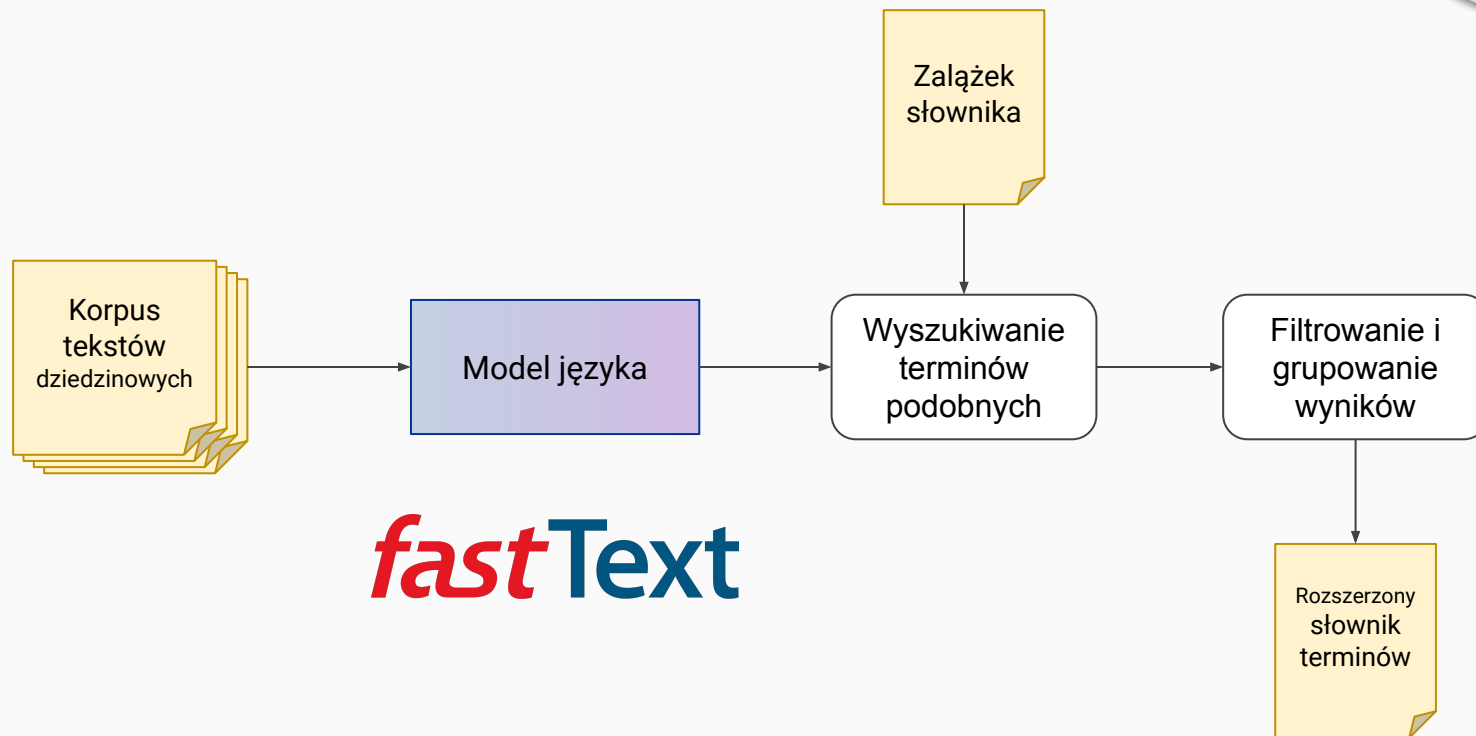
— im wyższa pewność modelu, tym wyższa precyzja klasyfikacji



# Rozszerzanie zasobów językowych







*fast*Text

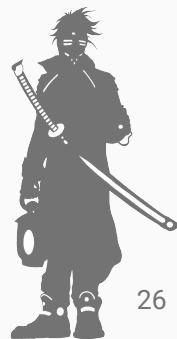


## Terminy ofensywne z uwzględnieniem morfologii

Liczba pojedyncza, wołacz, rodzaj męski

- babsztylu
- bałwanie
- baranie
- bęcwałę
- bezmózgowcu
- bezmózgowcze
- bezmózgu
- biedaku
- bolszewiku
- bucu
- ...

Top listy	Poniżej <i>top-u</i>
50.10 skurwysynie	20.83 pierdolisz
43.59 gówniarzu	19.89 skórwysynu
36.23 bydlaku	19.78 matolku
35.12 przyglupku	18.54 spierdalaj
33.95 zasańcu	18.15 pierdolcu
33.87 wszarzu	18.04 chłystku
32.87 gnojku	17.80 skurwysynek
32.61 popaprańcu	17.64 kanciarzu
31.34 kmiecie	17.61 przyglupasie
31.15 pętaku	17.15 padalcu
31.05 skurwysynku	17.05 pojechańcu
30.91 pojebie	
29.79 wypierdku	
28.80 smarkaczu	
27.89 śmieciu	
27.33 głupolu	
27.22 kmiotku	
26.00 idioto	
25.51 cwaniaczku	
25.16 debilku	
24.19 skurczybyku	
22.97 pojebusie	



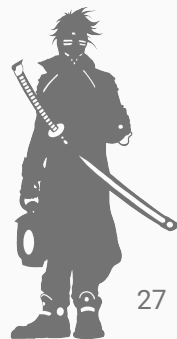
## Terminy do wykrywania intencji samobójczych

### kill

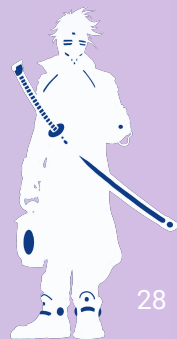
kms	0.654263
die	0.618894
asap	0.616186
commit	0.611897
iam	0.610663
killing	0.608998
guts	0.607888
offing	0.606651
punish	0.604843
ima	0.603541

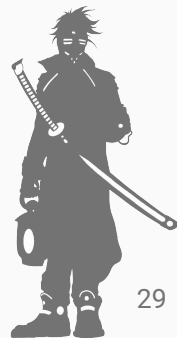
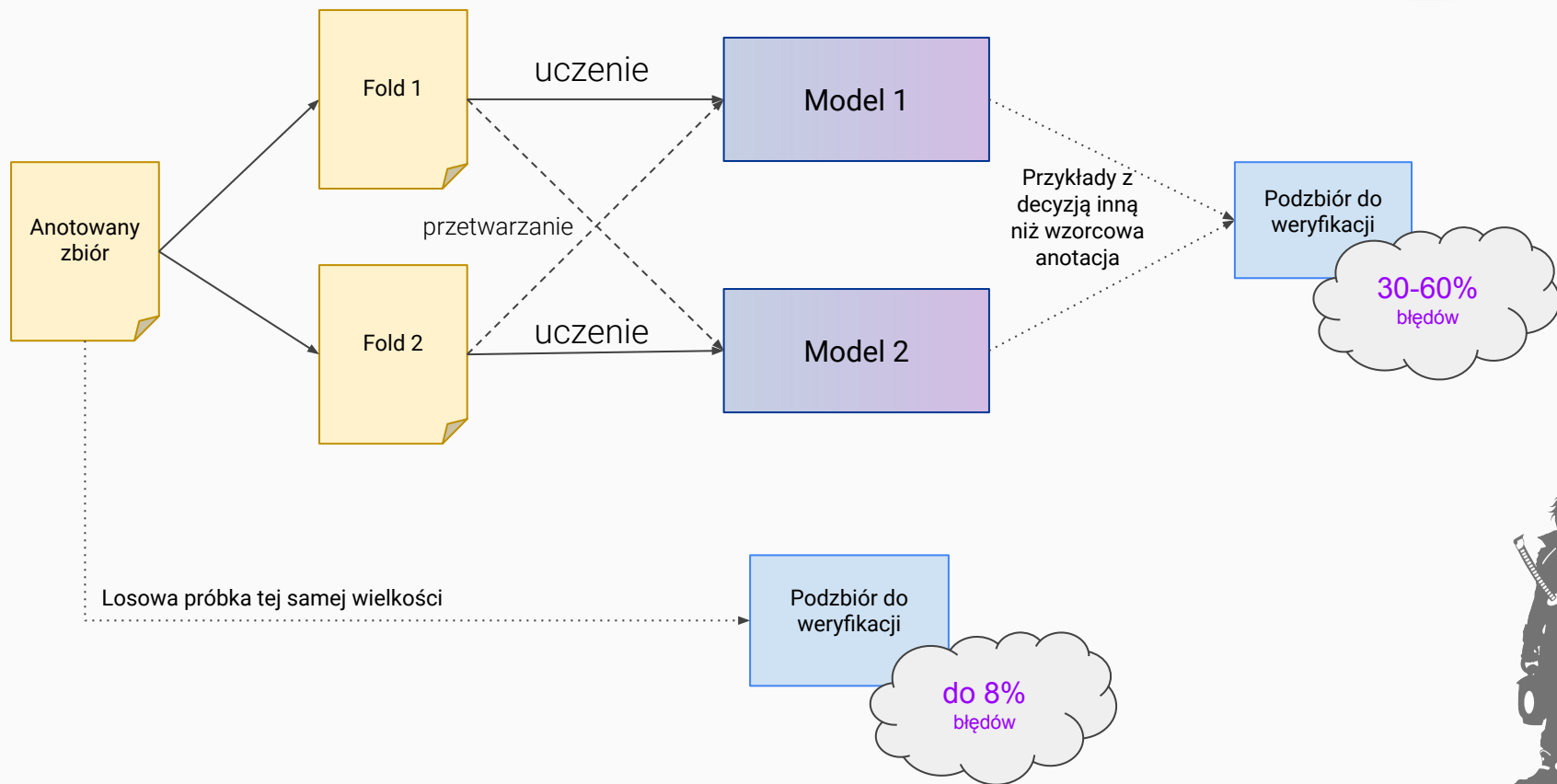
### kms

1. cant wait to **kms** on heroin
2. i already tried to **kms** once when i was like but failed
3. i tried to **kms** once but failed
4. im about to turn and that will most likely be the day that i will **kms** th april
5. i set a date to **kms**

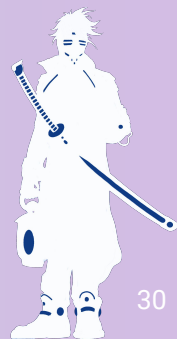


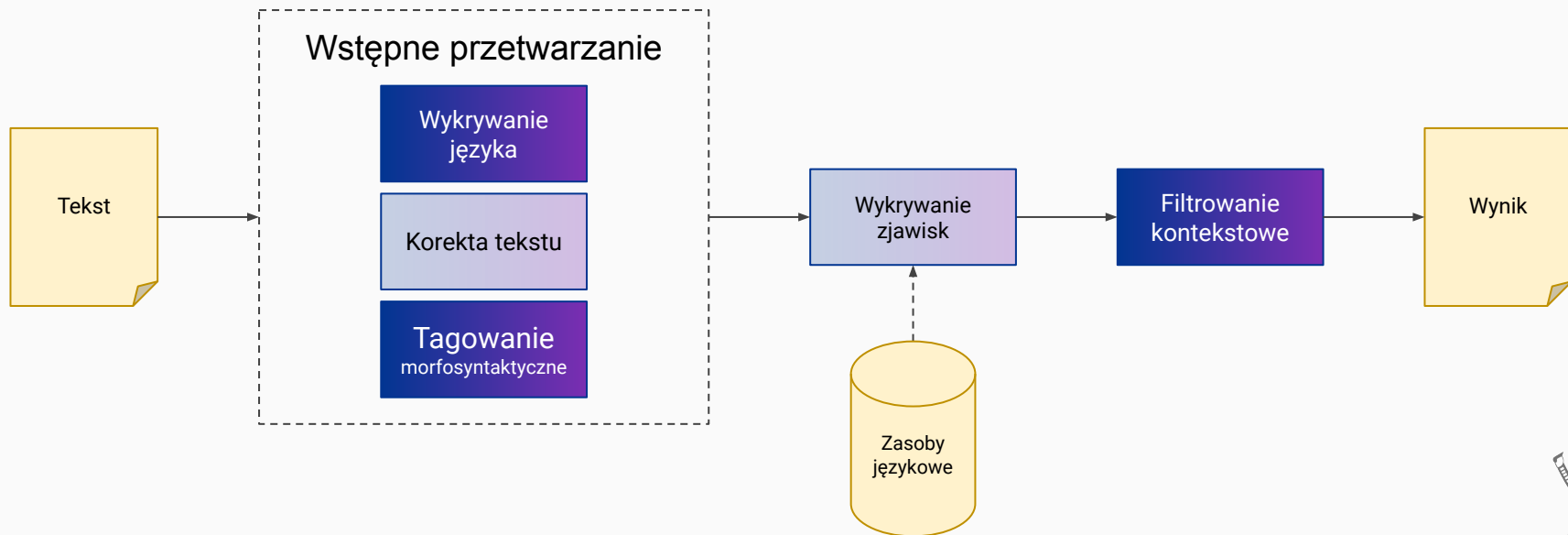
# Weryfikacja anotacji





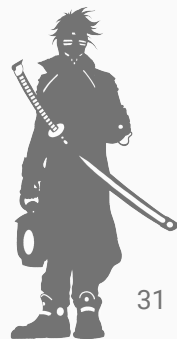
# Realizacja zadań pomocniczych





Moduł  
ML

Moduł  
symboliczny

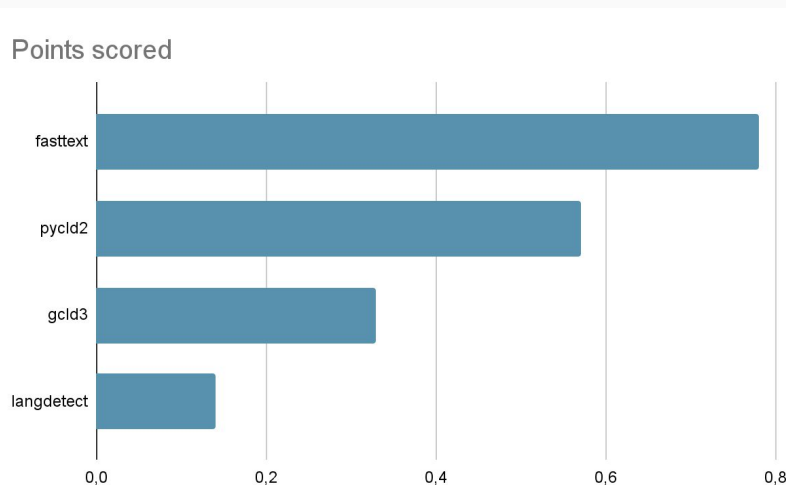


## Narzędzia

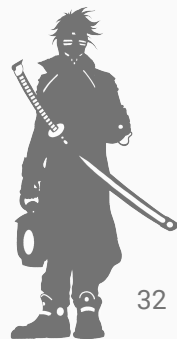
- fasttext
- langid
- langdetect
- spacy-langdetect
- pylcl3
- gcl3

	fasttext	langid	langdetect	spacy-langdetect
<b>speed</b>	129ms	8.07s	1min 4s	2min 2s
<b>accuracy</b>	0.860	0.827	0.845	0.845

Tweet set (<https://www.kaggle.com/rtatman/the-umass-global-english-on-twitter-dataset>)

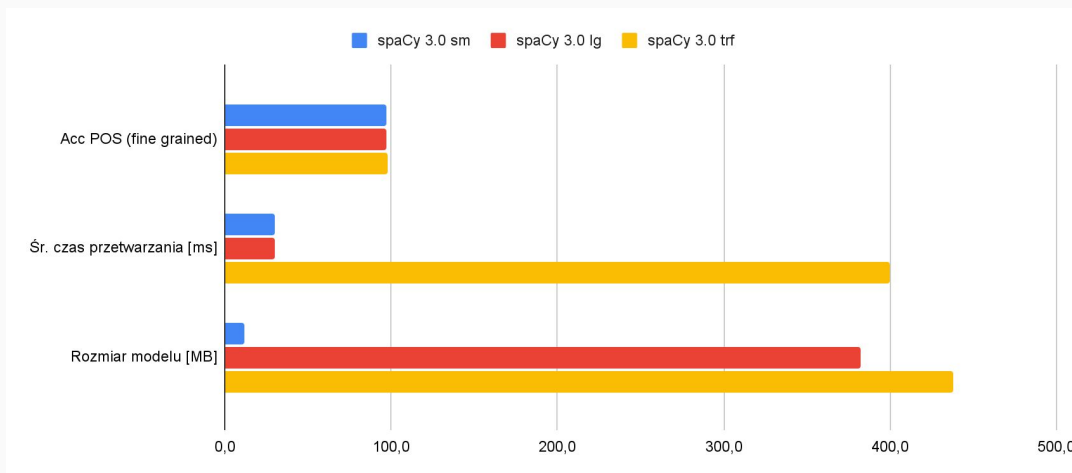


Wpisy z Discord-a



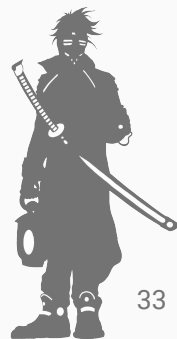


- **spaCy 3.0** dla j. angielskiego (<https://github.com/explosion/spaCy>)

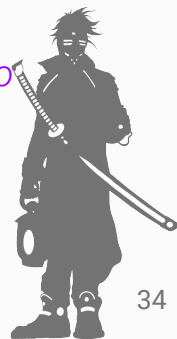


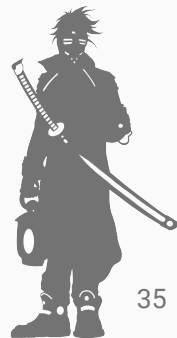
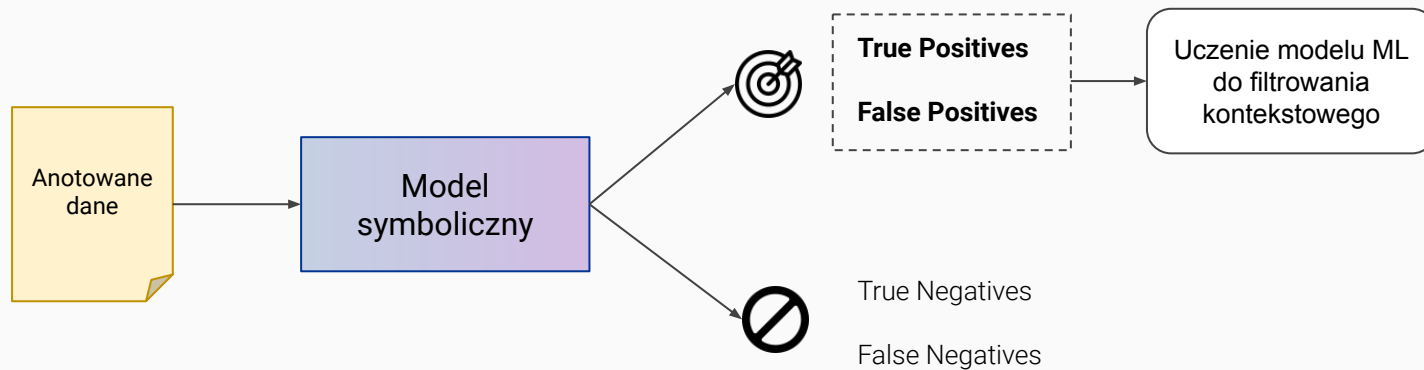
- **KRNNT** dla j. polskiego (<https://github.com/kwrobel-nlp/krnnt>)

KFTT (<https://github.com/kwrobel-nlp/kfft>) – 97% vs. 93% (KRNNT)



- *"I am 32 and my gf is 13 (...)"*
  - pedofilia?
  - *"(...) Can we still do quests together or are our levels too far apart? Sorry kind of new to WoW."*
  - Źródło — forum gry World of Warcraft
- *"I tried to kill myself and I couldn't. (...)"*
  - próba popełnienia samobójstwa?
  - *"(...) I also noticed I had no effect on objectives. It was so weird. The next match I was back to normal."*
  - Źródło — forum gry Call of Duty: Modern Warfare





# Dziękuję za uwagę

samurai  
LABS

## Customers

Khoros

INACH

webex  
by CISCO

HEROO  
MOBILE



## Partners



GGP GAMPLAY  
ALLIANCE

PUBLIC

FACTMATA



## Press

Forbes

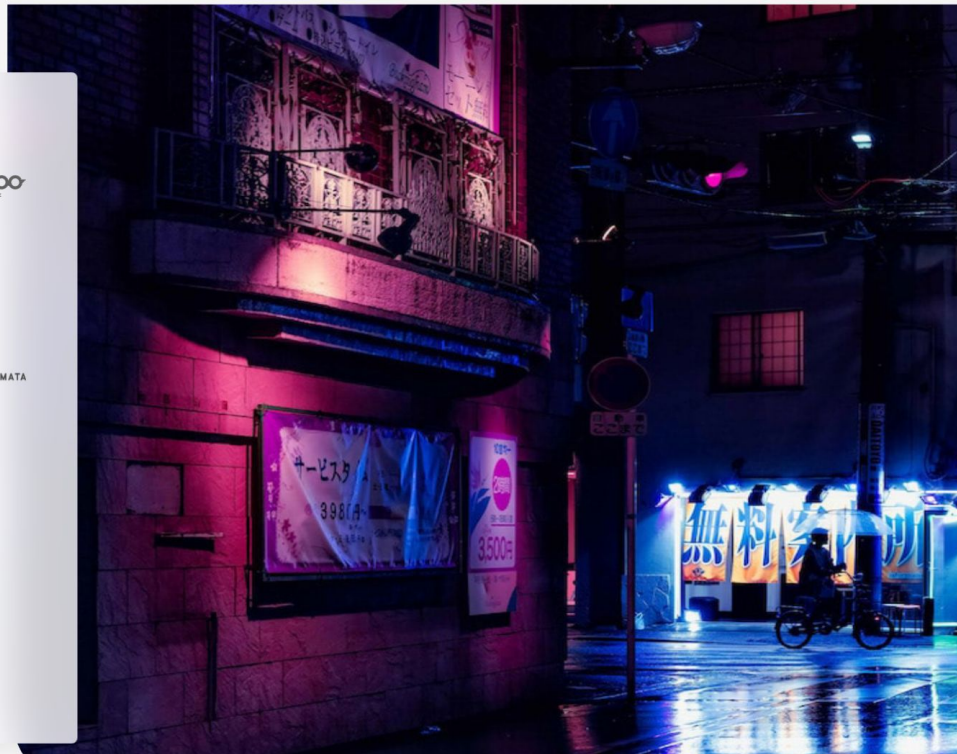
B B C

BUSINESS  
INSIDER

## Technology Awards

CB INSIGHTS

Gartner



# Neuro-symboliczna sztuczna inteligencja w walce z cyberprzemocą w Internecie

dr inż. Michał Marcińczuk (Samurai Labs)

**Blackmail Detection Example:**

*Send me more nudes or I 'll publish the ones I already have*

Send me more nudes or I 'll publish the ones I already have

Request      Alternative      Negative Consequences

Samurai Cyber Guardian Detection
<b>Blackmail</b> threat to reveal information
<b>Sexual harassment</b> attempt to solicit intimate photographs

Please input your username.

\*\*\*\*\*

**This username violates our community standards.**

English: **INAPPROPRIATE** **PROFANITY**

Polish: **OFFENSIVE** **PROFANITY** **SEXUAL**

