



NÁVRH CHATBOTA PRE POUŽITIE PROTI SCAM KOMUNIKÁCI

BC. DIANA FORTUNOVÁ

ÚINF – ÚSTAV INFORMATIKY

VEDÚCI: RNDR. ĽUBOMÍR ANTONI, PHD.

KONZULTANT: DOC. RNDR. JUDR. PAVOL SOKOL, PHD.



john888mrmichael@gmail.com

komu: ▼

so 12. 11. 12:16

Dobrý deň,

Prijmite, prosím, moje ospravedlnenie, nechcem zasahovať do vášho súkromia, napísal som vám predchádzajúci e-mail, ale bez odpovede, v prvom e-maile som sa vám zmienil o mojom zosnulom klientovi, od jeho smrti som dostal niekoľko posledných od jeho banke, kde pred smrťou zložil zálohu vo výške 5,4 milióna dolárov, ma banka požiadala, aby som poskytol jeho najbližším príbuzným alebo niektorým z jeho príbuzných, ktorí si môžu uplatniť nárok na jeho finančné prostriedky, inak budú bankou skonfiškované, keďže som ich nevedel nájsť ktorýkoľvek z jeho príbuzných, preto som vás kontaktoval s týmto nárokom.

Veľmi si vážime vašu naliehavú odpoveď a spoluprácu a získajte ďalšie podrobnosti o tejto transakcii.

Vďaka

S úctou,

Advokát John Michael.

CONSUMER SENTINEL

By FTC (2022)

CONSUMER
SENTINEL
NETWORK
DATA BOOK 2022



SNAPSHOT

5.3
MILLION
REPORTS

TOP THREE CATEGORIES

- 1 Identity Theft
- 2 Imposter Scams
- 3 Credit Bureaus, Info Furnishers and Report Users

2.5 million fraud reports

25% reported a loss

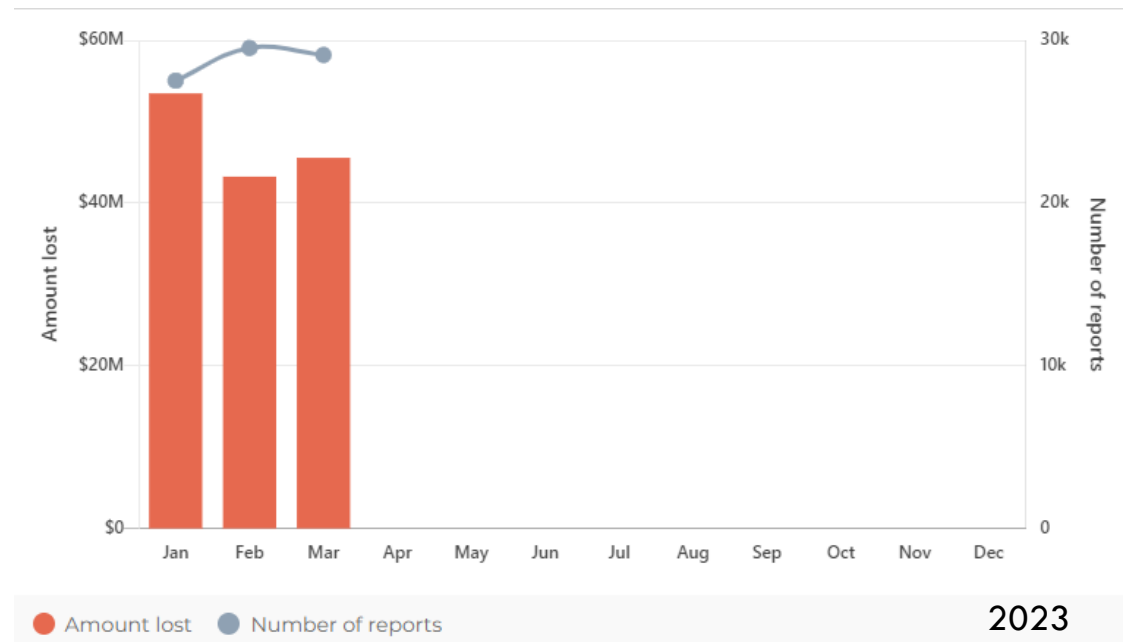
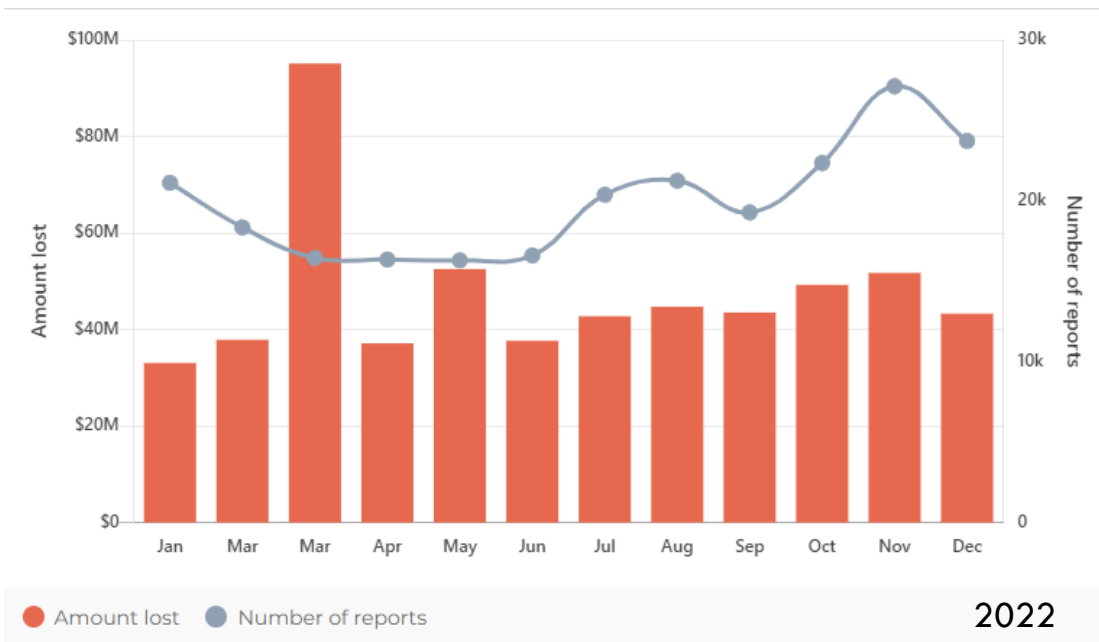


\$9.0 billion
total fraud losses

\$650
median loss

SCAMWATCH

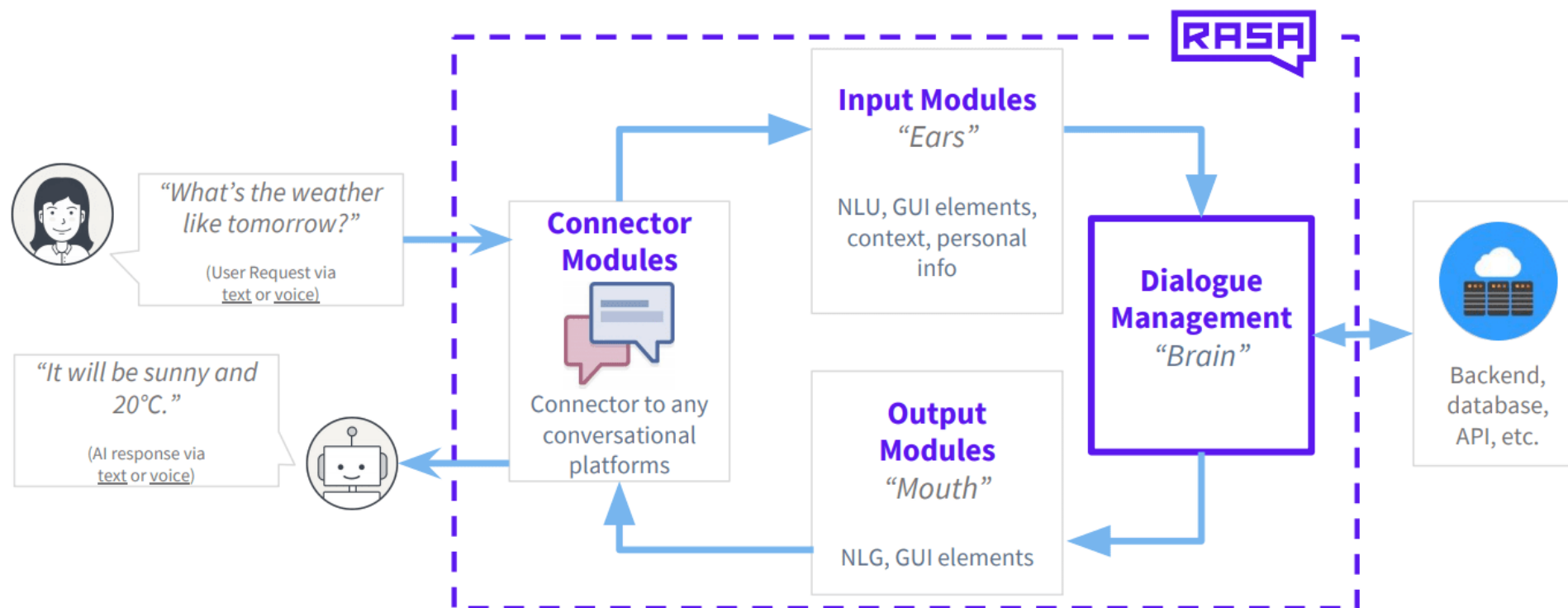
By ACCC



MODEL SYSTÉMU KONVERZAČNÉHO ROBOTY

Introduction

Rasa the OSS to build conversational software with ML



Alternatives:  Dialogflow  wit.ai 

RASA

VSTUPNÝ MODUL KONVERZAČNÉHO ROBOTY

- Named Entity Recognition (NER)
- Textová analýza pre klasifikáciu podvodných správ

NAMED ENTITY RECOGNITION (NER)

- Prvý krok k extrakcii informácií
- Hľadá a klasifikuje pomenované entity v texte
- Používaný v mnohých oblastiach NLP (Spracovanie prirodzeného jazyka)



Named Entity Recognizer i

Disease (Perceptron) ▼

Number Of Terms i

25

Number Of Concepts i

25

Extract

Reset

Extraction Results

Named Entities

Concepts

Important Terms

Document Classification

Disease i

tumor

cancer

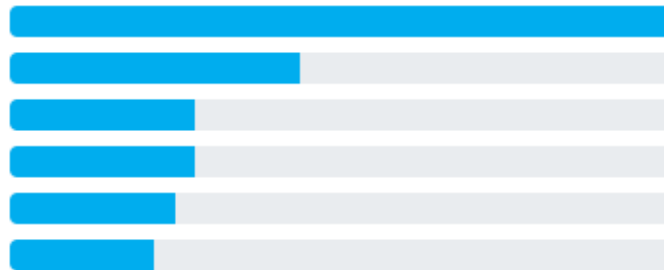
tumors

myelodysplastic syndromes

breast cancer

metastasis

Score



Frequency

3

1

1

1

2

1

AKO NA IMPLEMENTÁCIU ?

- Platforma NLTK

Identify named entities:

```
>>> entities = nltk.chunk.ne_chunk(tagged)
>>> entities
Tree('S', [(('At', 'IN'), ('eight', 'CD'), ("o'clock", 'JJ'),
            ('on', 'IN'), ('Thursday', 'NNP'), ('morning', 'NN')),
            Tree('PERSON', [(('Arthur', 'NNP')]),
                ('did', 'VBD'), ("n't", 'RB'), ('feel', 'VB'),
                ('very', 'RB'), ('good', 'JJ'), ('.', '.')])])
```

Tokenize and tag some text:

```
>>> import nltk
>>> sentence = """At eight o'clock on Thursday morning
... Arthur didn't feel very good."""
>>> tokens = nltk.word_tokenize(sentence)
>>> tokens
['At', 'eight', "o'clock", 'on', 'Thursday', 'morning',
'Arthur', 'did', "n't", 'feel', 'very', 'good', '.']
>>> tagged = nltk.pos_tag(tokens)
>>> tagged[0:6]
[(('At', 'IN'), ('eight', 'CD'), ("o'clock", 'JJ'), ('on', 'IN'),
('Thursday', 'NNP'), ('morning', 'NN'))]
```

PERSON
NORP
FAC
ORG
GPE
LOC
PRODUCT
EVENT
WORK OF ART
LAW
LANGUAGE
DATE
TIME
PERCENT
MONEY
QUANTITY
ORDINAL
CARDINAL

spaCy

Alan Turing PERSON was born in 1912 DATE at Paddington GPE , London GPE .

```
[('Alan Turing', 'PERSON'), ('1912', 'DATE'), ('Paddington', 'GPE'), ('London', 'GPE')]
```

- Alan Turing - 0 : 11 - PERSON
- 1912 - 24 : 28 - DATE
- Paddington - 32 : 42 - GPE
- London - 44 : 50 - GPE

Apple ORG sold nearly 20 thousand CARDINAL iPods PRODUCT for a profit \$6 million MONEY .

```
[('Apple', 'ORG'), ('nearly 20 thousand', 'CARDINAL'), ('iPods', 'PRODUCT'), ('$6 million', 'MONEY')]
```

- Apple - 0 : 5 - ORG
- nearly 20 thousand - 11 : 29 - CARDINAL
- iPods - 30 : 35 - PRODUCT
- \$6 million - 49 : 59 - MONEY

ENTITY POUŽITEĽNÉ V BEZPEČNOSTI

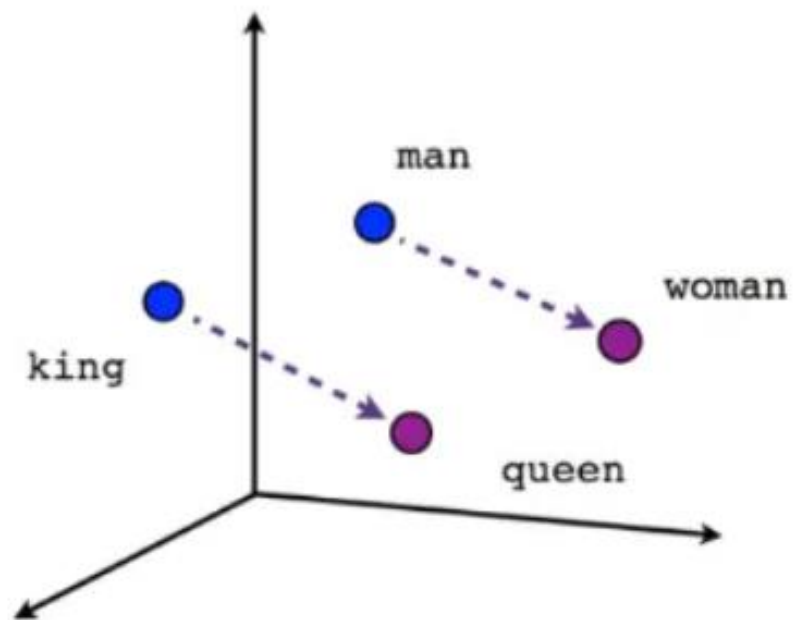
- Osoby – mená, zamestnanie, popis, ... (Trump, teacher, child, ex husband)
- Organizácie – Google, Facebook, ...
- Geografická poloha – štáty, mestá, obce, ...
- Produkty – iPod, tablet, smartphone, ...
- Právne dokumenty(LAW)
- Dátum a čas
- Sumy peňazí

niečo navyiac ?



TEXTOVÁ ANALÝZA

- Word2vec



Male-Female

fig 1: king to queen is like man to woman. it is illegal to write about **word2vec** without attaching this plot

AKO NA TO ?

- (CBOW)
- Skip-Gram model

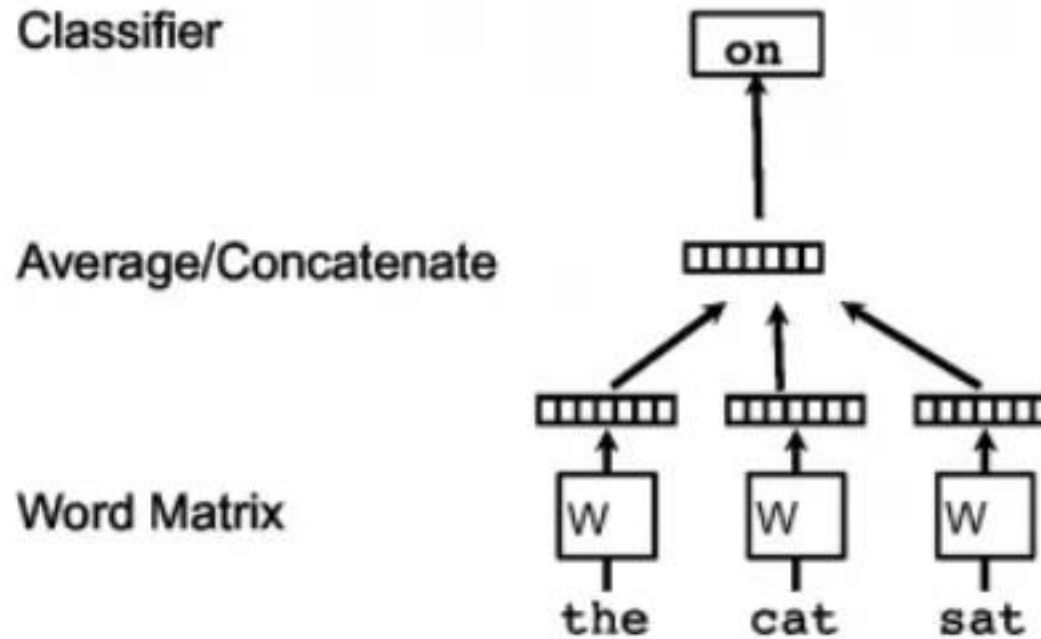


fig 2: CBOW algorithm sketch: the words "the" "cat" "sat" are used to predict the word "on"

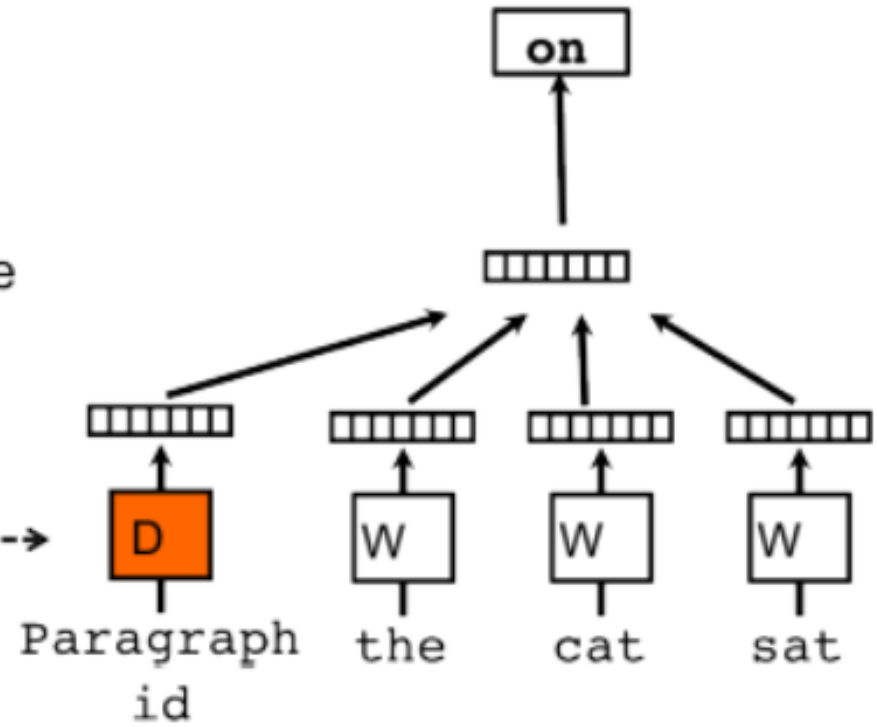
AKO NA TO ?

- Doc2vec
PV-DM

Classifier

Average/Concatenate

Paragraph Matrix



AKO NA TO ?

- Doc2vec
PV-DBOW

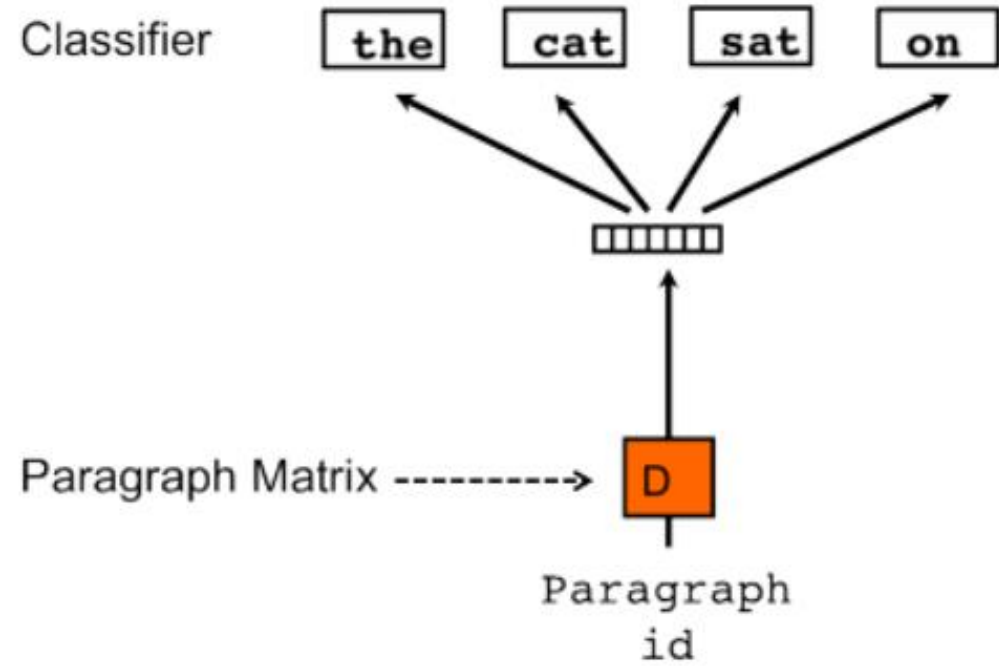


fig 4: PV-DBOW model

4 HLAVNÉ KATEGÓRIE SPRÁV

- Dedičstvo
- Zoznamka
- Predvolanie na políciu
- Výhra



CIELE PRE NASLEDUJÚCE MESIACE

- Pomocou nástrojov predstavených v prezentácii vytvoriť program pre klasifikáciu entít z oblasti bezpečnosti v podvodných správach
- Klasifikovať doposiaľ prijaté správy do zvolených kategórií
- Pustiť sa do návrhu riešenia pre nasledujúci modul

The background features a series of concentric, light blue circles centered on the page. In the four corners, there are stylized circuit board traces in a darker blue color, with small circles at the end of the lines, resembling electronic components or data paths.

ĎAKUJEM ZA POZORNOST