

# Hľadanie podobností v textoch ľudových

## ▶ piesní

Michal Mižák

Vedúci práce: doc. RNDr. Stanislav Krajči PhD

# Výstup práce

- ▶ **Vložím pieseň alebo úryvok textu**
- ▶ **Dostanem zoznam podobných piesní alebo pieseň, ktorá úryvok v nejakej forme obsahuje**
- ▶ **Dostanem štatistické vyhodnotenie podobnosti**
- ▶ **Porovnanie rôznych konfigurácií nášho „Vector space“ algoritmu**
  
- ▶ **Bonus:**
  - ▶ „Mergovací“ systém
  - ▶ API jednoducho využiteľné v iných aplikáciách

# Využitie

- ▶ Od dediny ku dedine, od domu k domu **pretvorený/ skomolený text**
- ▶ Pri folklórnom výskume sa väčšinou **piesne získavajú od starších ľudí**, ktorí spievajú tak, ako si spomenú
  - ▶ Síce **autentické**, ale odborníci to môžu **opraviť**
- ▶ Pri rozumnej databáze porovnávanie **regionálnych rozdielov** v piesňach
- ▶ Priestor pre odborníkov a historikov na „skonzistentnenie“ záznamov o piesňach
  - ▶ Nesprávny zápis hlások v nárečí

# Využitie

- ▶ **Dzifče** počarovne, **ňezal'up** še do mňa  
ja chlapec vandrovni, co ci budze zo mňa  
  
Ja chlapec **vandrovni**, zo šireho poľa  
co ci budze zo mňa, frajirečko moja?
- ▶ **Dzifče** počarovne, **nezal'ub** se do mne  
ja chlapec **vandrovny**, co ci budze zo mňa  
  
Ja chlapec **vandrovny**, zo šireho poľa  
co ci budze zo mne, frajirečko moja?
- ▶ **Dzivče** počarovne, **nepopatraj** na mne  
ja chlapec vandrovny, co ci budze zo mne  
  
Ja chlapec **vandrovný** zo šireho poľa  
co ci budze zo mne draha duša moja?

# Technológie

- ▶ Programovací jazyk Java
  - ▶ Knižnica JOOQ
- ▶ Databáza
  - ▶ MySQL
- ▶ User interface
  - ▶ JSF

# JOOQ v skratke

```
Long id = create.select()
    .from(T_SONG)
    .where(T_SONG.LYRICS.eq(song.getLyrics()))
    .fetchOne(T_SONG.SONGID);
```


```
SongRecord songRecord = create
    .insertInto(T_SONG, T_SONG.TITLE, T_SONG.LYRICS, T_SONG.CLEANLYRICS,
        T_SONG.SONGSTYLE, T_SONG.REGION, T_SONG.SOURCE)
    .values(s.getTitle(), s.getLyrics(), s.getCleanLyrics(),
        s.getSongStyle(), s.getRegion(), s.getSource())
    .onDuplicateKeyIgnore()
    .returning(T_SONG.SONGID)
    .fetchOne();
```

DONE!

1. Eventuálne  
appka chce  
bežať na serveri

2. Vraj to je  
celkom dobré

JSF



To Do...

# Príprava databázy

- ▶ Scraping a rôzne formy skriptovania

- ▶ Dáta bolo treba niekde zohnať

- ▶ Viktor Gliganič - primáš (1. huslista) Ľudovej hudby FS Zemplín

- ▶ Piesne372

- ▶ Iné zdroje v stave ToDo

☑ DONE!



# „Čistenie“ textov

- ▶ Opakovačky
  - ▶ Štandardná štruktúra ľudovky
    - ▶ [: Sloha :] [: Refrén :]
  - ▶ Pre samotnú pieseň nie je opakovaný text refrénu/slohy zaujímavý
    - ▶ Zátvorky a špeciálne znaky (teda {[/:.,““]}) tým pádom môžeme ignorovať
- ▶ Čísla ignorujeme tiež

☑ DONE!

# Čistenie a parsovanie textov

- ▶ „ ... “
  - ▶ Ignorovaný nedopísaný text
- ▶ Čiastočne zatriedenie piesní podľa regiónov, tónin (dur, mol) a charakteru (valčík, polka...)

DONE!

# Algoritmus

- ▶ Potrebujeme algoritmicky porovnať dokumenty
  - ▶ Information retrieval
- ▶ Myšlienka:
  - ▶ Vezmem slová ako názvy stĺpcov
  - ▶ Riadok = pieseň, hodnota = početnosť
  - ▶ Získam vektor
  - ▶ Vektory algebraicky porovnáam a získam vzdialenosť

# Vector space algorithm v rychlosti

DONE!

- ▶ „raz raz dva dva raz raz“

| raz | dva |
|-----|-----|
| 4   | 2   |

- ▶ [4, 2]

- ▶ „raz raz dva raz dva raz raz“

| raz | dva |
|-----|-----|
| 5   | 2   |

- ▶ [5, 2]

Lepši dva raz po raz jak raz raz-dva

# TermScheme - N-gramy

- ▶ Poradie?
  - ▶ Porovnávanie n-tíc slov = n-gramov

DONE!

# Vector space algorithm v rychlosti

DONE!

- ▶ „raz raz dva dva raz raz“

|     |     |
|-----|-----|
| raz | dva |
| 4   | 2   |

- ▶ 2-gramy

|         |         |         |         |
|---------|---------|---------|---------|
| raz raz | raz dva | dva dva | dva raz |
| 2       | 1       | 1       | 1       |

- ▶ 3-gramy

|             |             |             |             |
|-------------|-------------|-------------|-------------|
| raz raz dva | raz dva dva | dva dva raz | dva raz raz |
| 1           | 1           | 1           | 1           |

Lepši dva raz po raz jak raz raz-dva

# Vector space

- ▶ Vzdialenosť dvoch vektorov
  - ▶ Kosínusová miera

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|_2 \|\mathbf{B}\|_2} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}, w$$

☑ DONE!

But wait, there's more...



# Vector inclusion

Poznámka - pracujeme s  
dimenziami vektora

# Vector inclusion

- ▶ Vezmime vektory s dimenziami
  - ▶ <šírka, výška, hĺbka>
  - ▶ <výška, šírka, hĺbka, čas>
  - ▶ <šírka, výška>
- ▶ Chcem nalepiť plagát na skriňu...
- ▶ Musíme vektory „naštelovať“ podľa dimenzií



„Slová“  
textu A

„Slová“  
textu B

# Vector inclusion

- ▶ Vektor  $v1 = \langle A, B, C \rangle$  |  $\langle a1, b1, c \rangle$ 
  - ▶ Reprezentuje text A
- ▶ Vektor  $v2 = \langle B, A, E \rangle$  |  $\langle b2, a2, e \rangle$ 
  - ▶ Reprezentuje text B

# Vector inclusion

- ▶ A forma - Slovu **C** vo  $v_2$  priradíme 0
  - ▶  $\langle A, B, C \rangle$
  - ▶  $v_1 = \langle a_1, b_1, c \rangle$
  - ▶  $v_2 = \langle a_2, b_2, 0 \rangle$

DONE!

# Vector inclusion

- ▶ B forma - Slovu **E** vo  $v_1$  priradíme 0
  - ▶  $\langle A, B, E \rangle$
  - ▶  $v_1 = \langle a_1, b_1, 0 \rangle$
  - ▶  $v_2 = \langle a_2, b_2, e \rangle$

To Be Continued...

Všetky „slová“

A Venn diagram consisting of a large light blue outer circle and two smaller overlapping circles inside it. The left inner circle is a darker blue and contains the text „Slová“ textu A. The right inner circle is a teal color and contains the text „Slová“ textu B. The intersection of the two inner circles is shaded with a lighter teal color. The text „Všetky „slová““ is positioned at the top of the large outer circle.

„Slová“  
textu A

„Slová“  
textu B

# Vector inclusion

- ▶ Zjednotenie

- ▶  $\langle A, B, C, E \rangle$

- ▶ Prienik

- ▶  $\langle A, B \rangle$

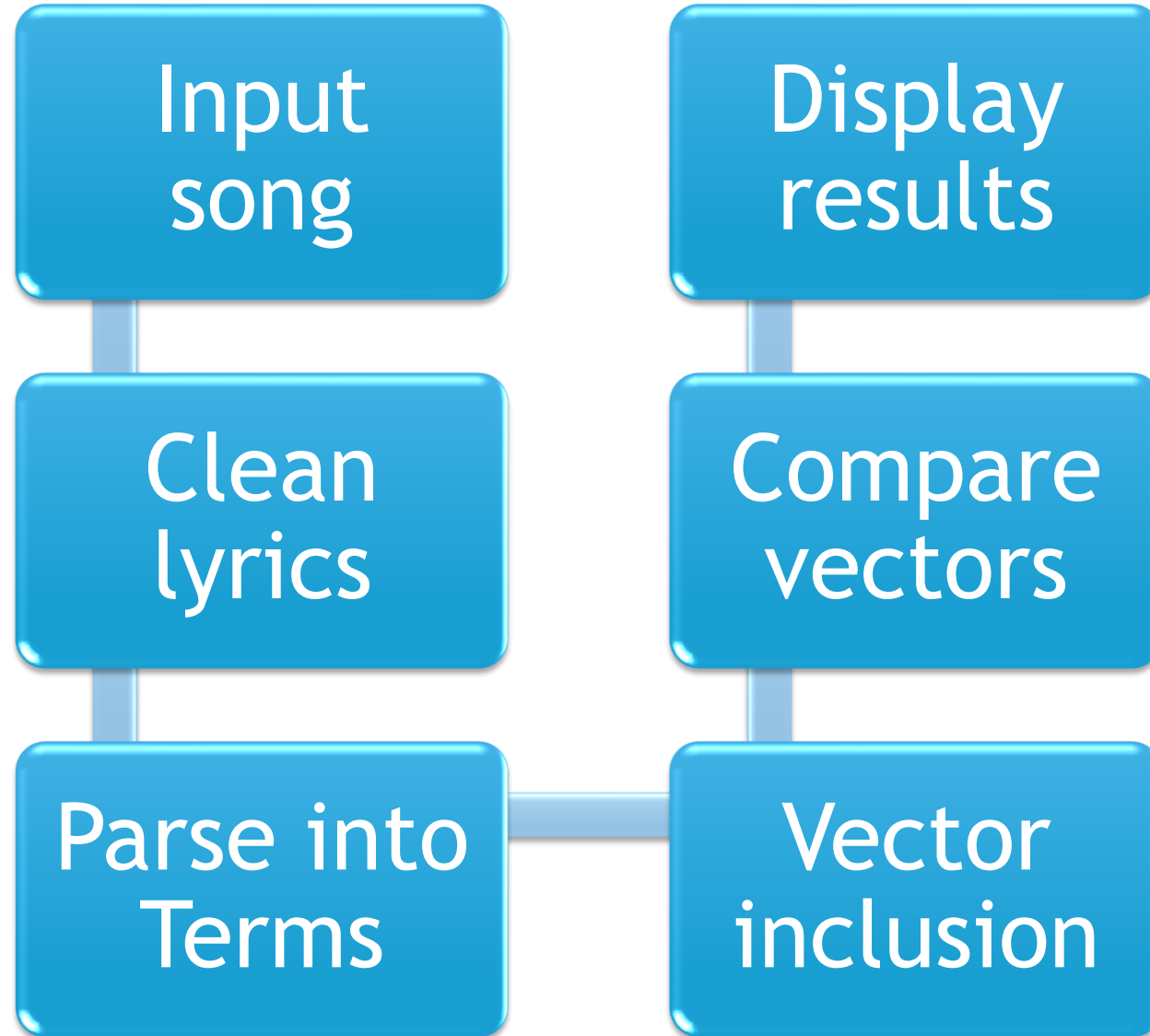
- ▶ Všetky slová databázy

- ▶  $\langle A, B, C, D, E, F, G, H, I, J, K... \rangle$



To Be Continued...





But wait, there's more...

# Term weighting

- ▶ „Slová“ v texte sú pre nás rôzne zaujímavé
  - ▶ Tf - term frequency
  - ▶ Idf - inverse frequency
  - ▶ Tf-idf (term frequency-inverse document frequency)
    - ▶ „A ja taka **čarna** jak **čarna čarnica**...“
    - ▶ „A ja taka **dzivočka cingi lingi bom**...“
  - ▶ Mnoho ďalších...



To Be Continued...

# TermScheme - zovšeobecnenie

- ▶ „Ja chlapec vandrovny, zo šireho poľa“
  - ▶ Dvojice
    - ▶ <Ja, vandrovny>, <chlapec, zo>, <vandrovny, šireho>
    - ▶ <Ja, poľa>
  - ▶ Trojice
    - ▶ <Ja, chlapec, zo>, <chlapec, vandrovny, šireho>
  - ▶ Mnoho ďalších...

DONE!

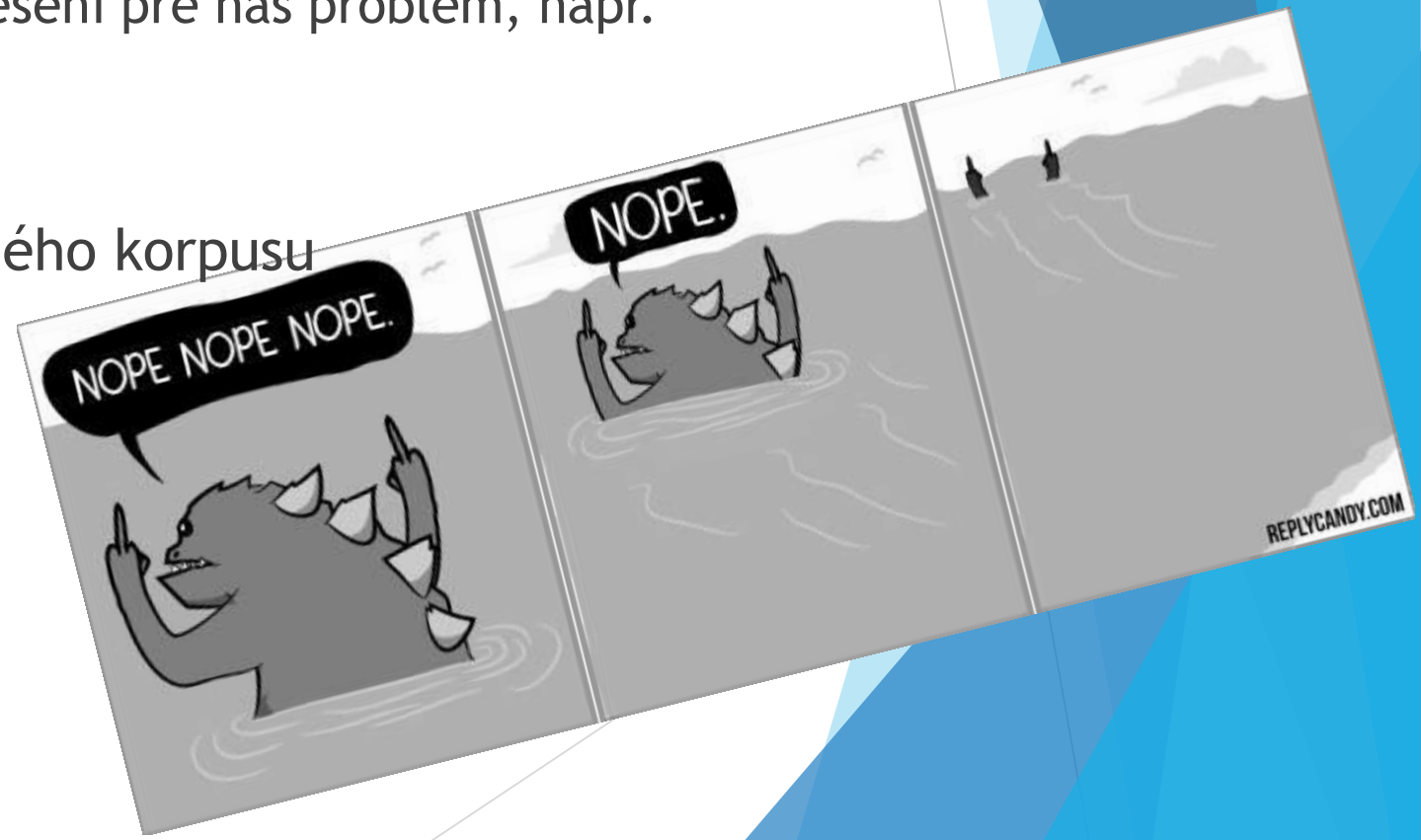
# Term comparison algorithm

- ▶ Keď už porovnávame vektory, môžeme porovnávať aj „slová“

DONE!

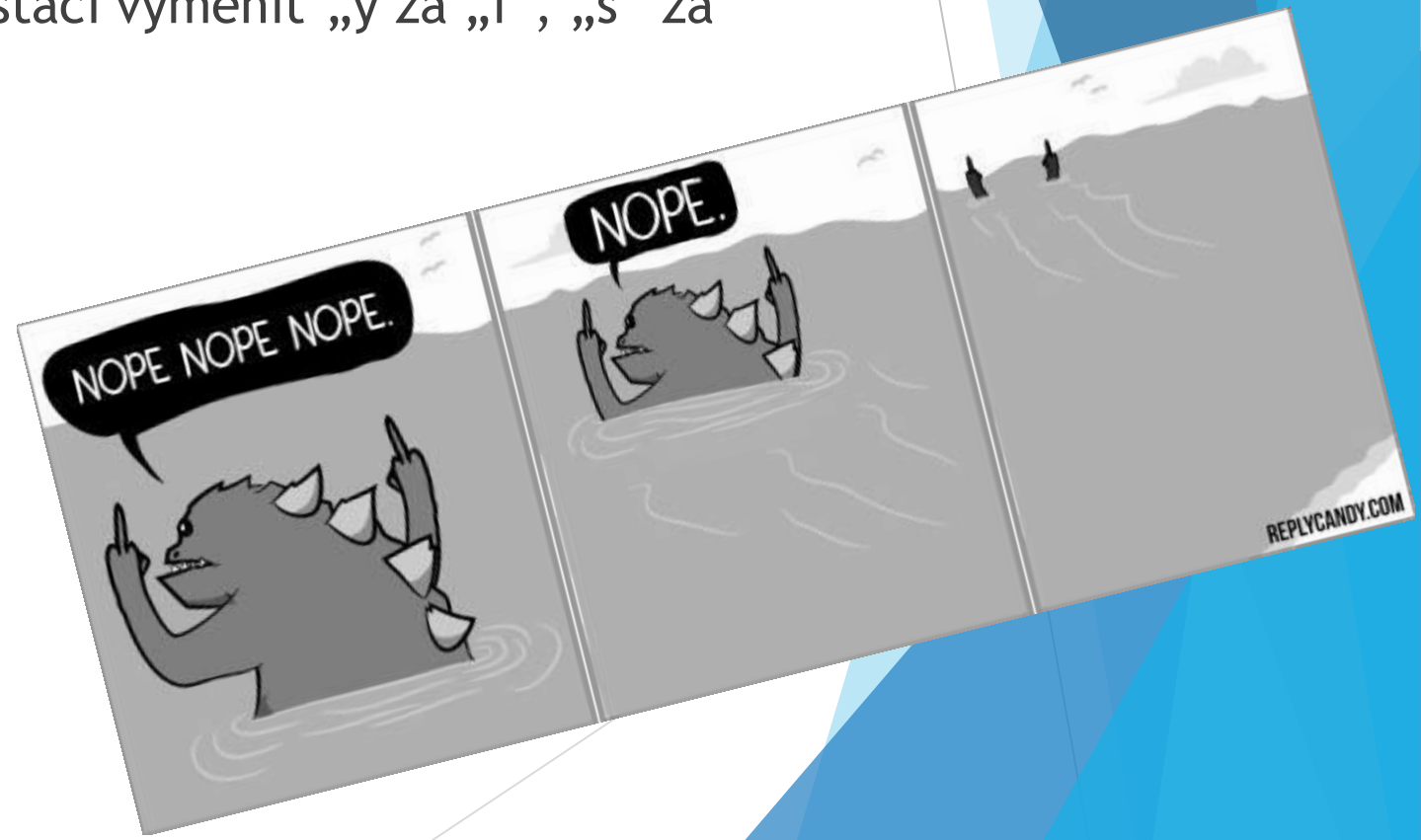
# Jazykovedný prístup

- ▶ Lemmatization, stemming - porterov algoritmus
  - ▶ Analýza efektivity existujúcich riešení pre náš problém, napr. STUBA
  - ▶ Ohýbanie riešení na ľudový jazyk
- ▶ Korpus nárečí Slovenského národného korpusu



# Tvaroslovník

- ▶ Tvaroslovník - vytváranie vlastného ľudového „tvaroslovníka“
  - ▶ Podobnosti sú často veľmi veľké, stačí vymeniť „y“ za „i“, „š“ za „ś“, „h“ za „g“ a pod.
  - ▶ Generovanie koncoviek



# Term similarity

- ▶ Levenshteinova vzdialenosť
- ▶ Mnoho ďalších...

1. **k**itten → **s**itten (substitution of "s" for "k")
2. **s**itten → **s**ittin (substitution of "i" for "e")
3. **s**ittin → **s**ittin**g** (insertion of "g" at the end).



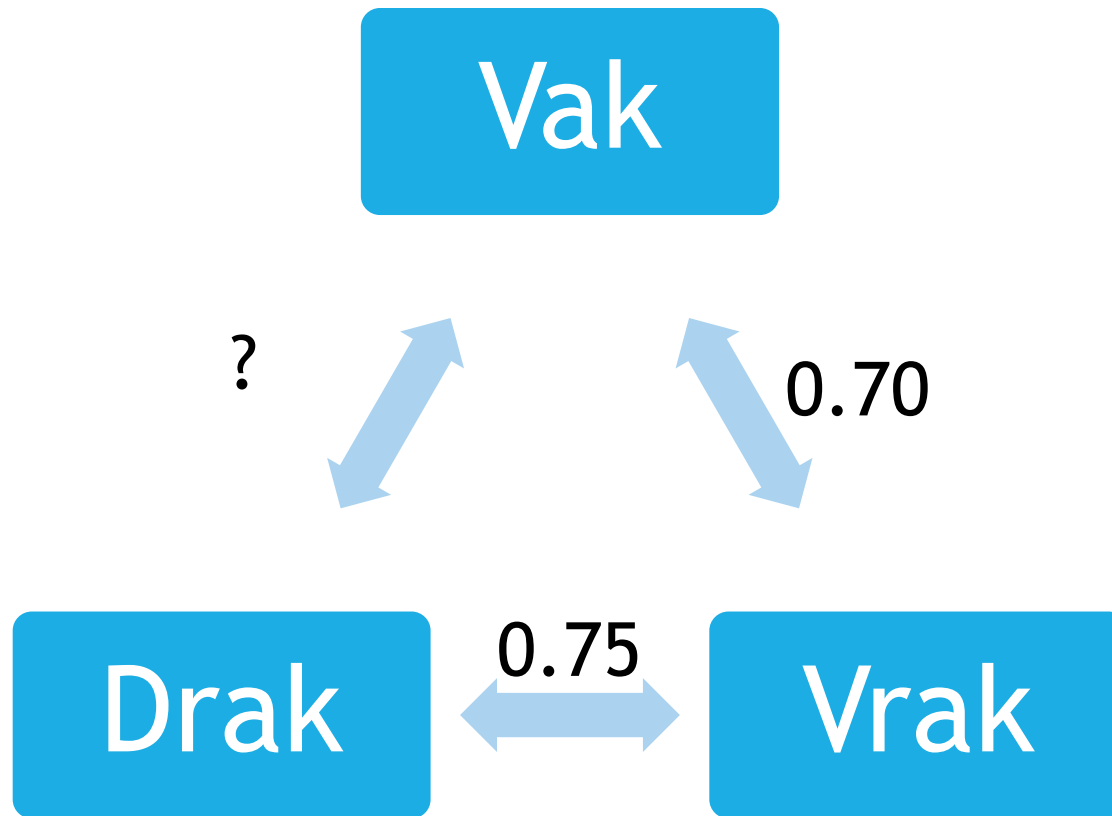
To Be Continued...



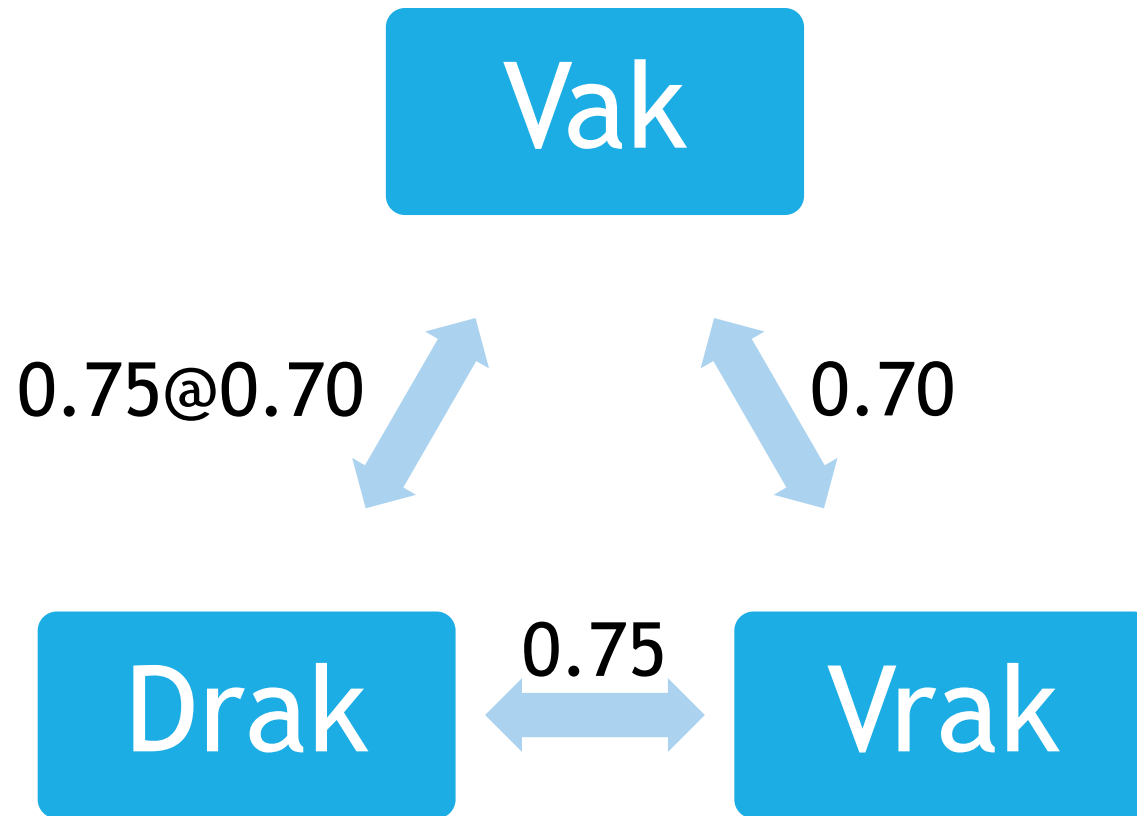
# Indirect term similarity

- ▶ Idea - „slová“ si nemusia byť podobné priamo, ale cez nejakého suseda
- ▶ Ilustrácia - číslo naznačuje, nakoľko sú si „slová“ podobné v danom porovnávacom algoritme

# Indirect term similarity

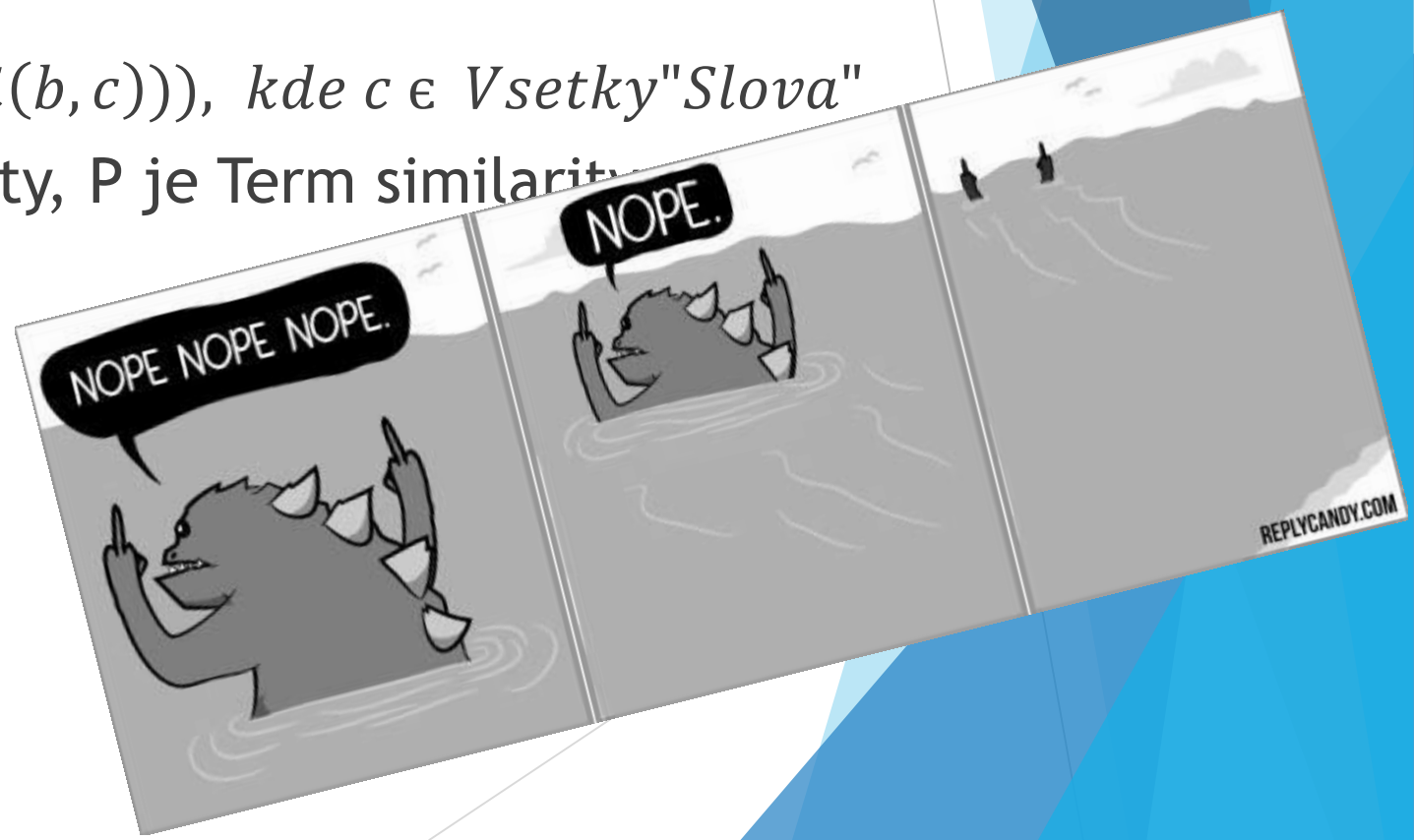


# Indirect term similarity



# Indirect term similarity

- ▶ Dostaneme teda vzorec
- ▶  $P_i(a, b) = \max(P(a, b), \max(P_i(a, c) @ P_i(b, c)))$ , kde  $c \in V$  všetky "Slova"
- ▶ Kde  $P_i$  je Indirect term similarity,  $P$  je Term similarity  
 $a, b$  sú „slová“



Praktickejšie riešenie:

# Term comparison + Tolerance

- ▶ Ak je podobnosť dvoch „slov“ v tolerancii, označíme ich za totožné
- ▶ Dáme ich do rovnakej **triedy ekvivalencie**

# Tolerancia

```
public enum Tolerance {  
    NONE, LOW, MEDIUM, HIGH  
}
```

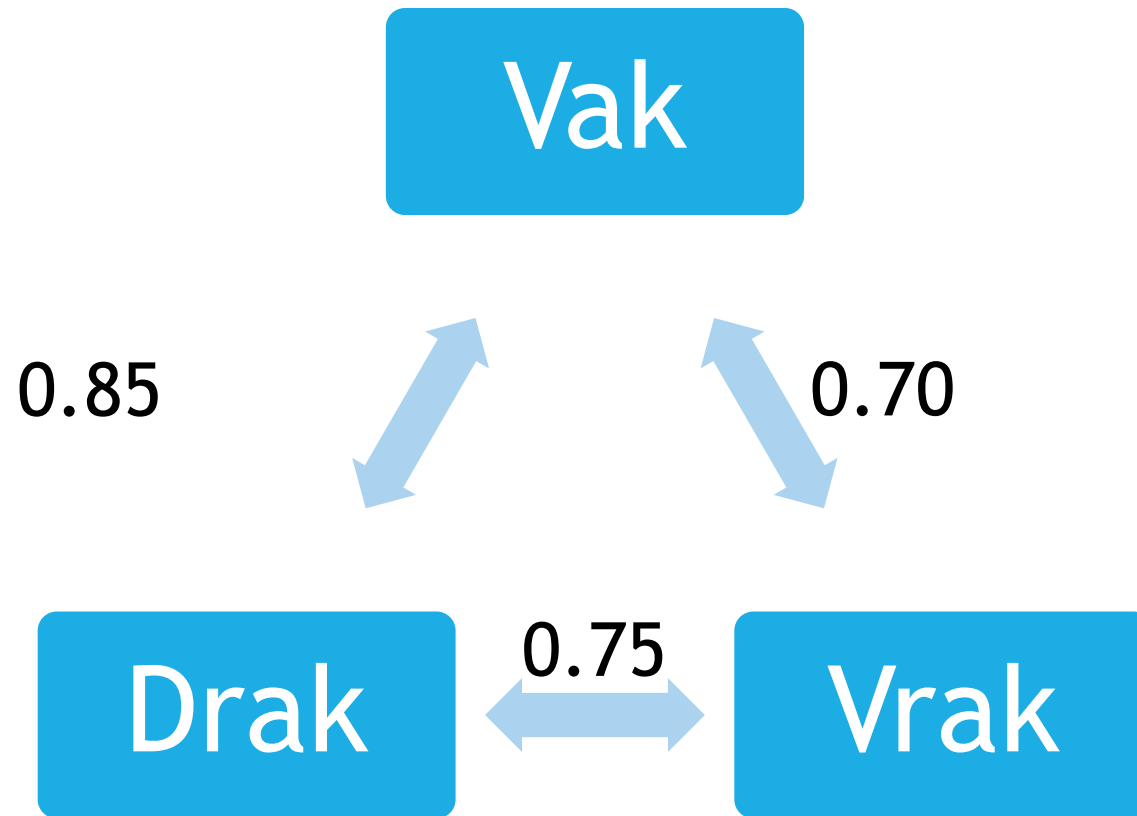
```
public interface IToleranceCalculator {  
    Double calculateTolerance(Tolerance tolerance,  
                             TermComparisonAlgorithm termComparisonAlgorithm);  
}
```

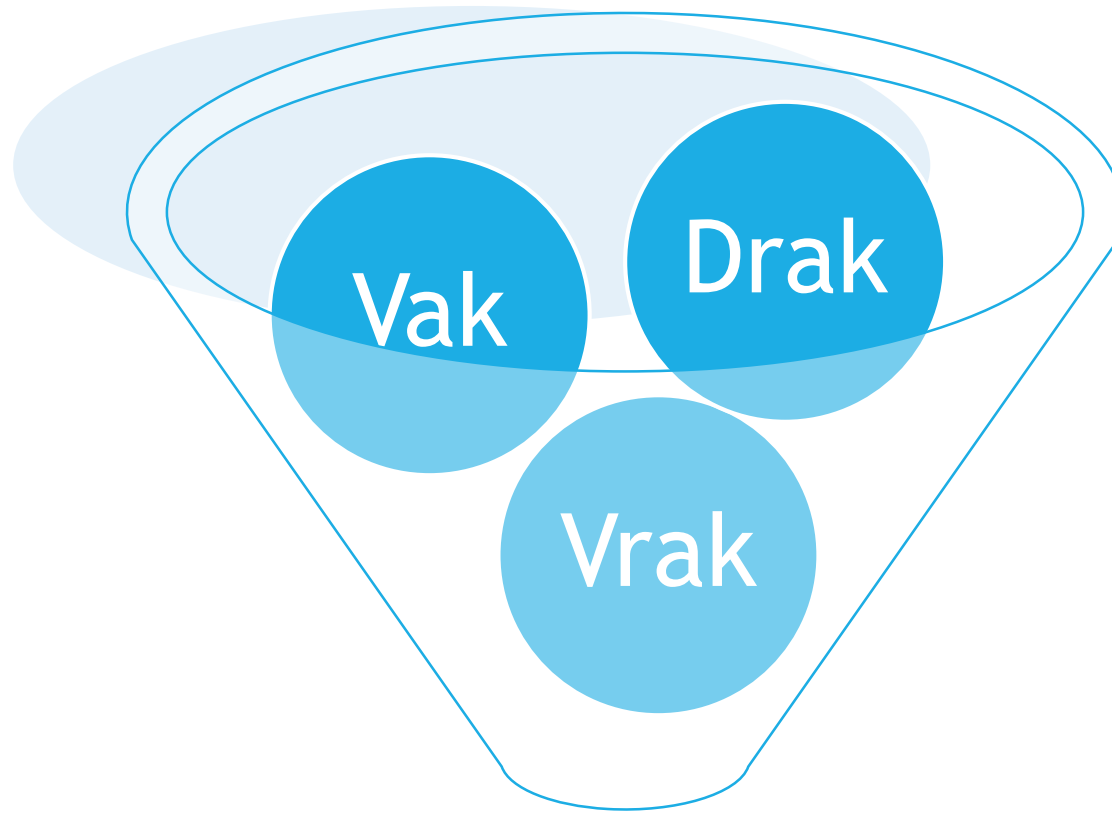
# Term grouping

- ▶ Tolerance = 0.8
- ▶ TermComparisonAlgorithm = NAIVE



# Term grouping





TermGroupXYZ

DONE!

# Term grouping

- ▶ Teda {Vak, Drak, Vrak} označíme ako triedu ekvivalencie pri tolerancii = 0.8 a naivnom porovnávacom algoritme

DONE!

# Rôzne prístupy a taktiky

```
public enum TermGroupMatchingStrategy {  
    MATCH_ALL, MATCH_ONE  
}
```

```
public enum TermGroupMergingStrategy {  
    MERGE_ANY, MERGE_ALL  
}
```

# Vector comparison algorithm

$$|Q \cap D|$$

Simple matching (coordination level match)

$$2 \frac{|Q \cap D|}{|Q| + |D|}$$

Dice's Coefficient

$$\frac{|Q \cap D|}{|Q \cup D|}$$

Jaccard's Coefficient

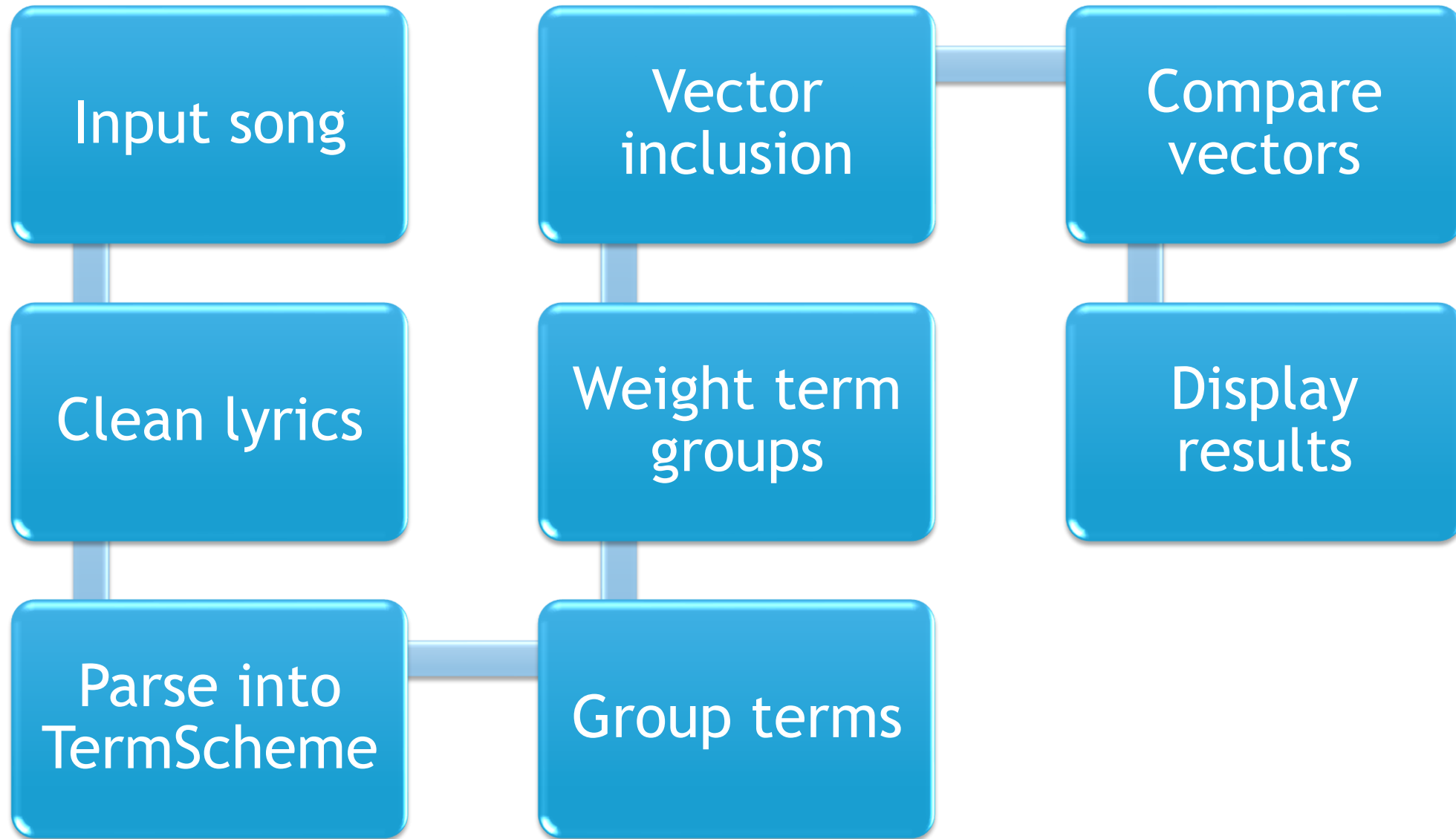
$$\frac{|Q \cap D|}{|Q|^{1/2} \times |D|^{1/2}}$$

Cosine Coefficient (what we studied)

$$\frac{|Q \cap D|}{\min(|Q|, |D|)}$$

Overlap Coefficient

To Be Continued...



```
switch (vectorInclusion) {
    case A:
        return aFormation(a, b, termComparator, termComparisonAlgorithm, tolerance);
    case B:
        return bFormation(a, b, termComparator, termComparisonAlgorithm, tolerance);
    case UNIFICATION:
        return unificationFormation(a, b, termComparator, termComparisonAlgorithm, tolerance);
    case INTERSECTION:
        return intersectionFormation(a, b, termComparator, termComparisonAlgorithm, tolerance);
    case ALL:
        return allFormation(a, b, termComparator, termComparisonAlgorithm, tolerance);
    default:
        throw new RuntimeException("Unimplemented vector inclusion");
}
```

```
/**
 * N-gram...
 */
private TermScheme termScheme;
private Integer termDimension;

/**
 * TF, IDF, TF-IDF
 */
private TermWeightType termWeightType;

/**
 * NAIVE, LEVENSHTEIN_DISTANCE
 */
private TermComparisonAlgorithm termComparisonAlgorithm;
private Tolerance tolerance;

private TermGroupMatchingStrategy termGroupMatchingStrategy;
private TermGroupMergingStrategy termGroupMergingStrategy;

/**
 * A, B, INTERSECTION, UNIFICATION, ALL
 */
private VectorInclusion vectorInclusion;

/**
 * SIMPLE_MATCHING, COS_COEFFICIENT, DICE'S_COEFFICIENT, JACCARD'S_COEFFICIENT...
 */
private VectorComparisonAlgorithm vectorComparisonAlgorithm;
```



# Zdroje

- ▶ Texty ľudových piesní
  - ▶ Rôzne stránky
  - ▶ Kontakty u nadšencov
- ▶ <https://en.wikipedia.org/wiki/N-gram>
- ▶ [https://en.wikipedia.org/wiki/Levenshtein\\_distance](https://en.wikipedia.org/wiki/Levenshtein_distance)
- ▶ <https://courses.cs.washington.edu/courses/cse573/12sp/lectures/17-ir.pdf>
- ▶ <https://nlp.stanford.edu/IR-book/>



To Be Continued...

Ďakujem za pozornosť