

# Hľadanie podobností v textoch ľudových piesní

Michal Mižák

Školiteľ: doc. RNDr. Stanislav Krajči PhD



*Femplin*

# 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



*Femplin*

# Využitie

- ▶ Od dediny ku dedine, od domu k domu **pretvorený/ skomolený text**
- ▶ **Piesne získavajú od starších ľudí**, ktorí spievajú tak, ako si spomenú
  - ▶ Síce autentické, ale odborníci to môžu opraviť
- ▶ Porovnávanie **regionálnych rozdielov** v piesňach
- ▶ „Skonzistentnenie“ záznamov o piesňach
  - ▶ Nesprávny zápis hlások v nárečí



Zemplin

# Využitie

- ▶ **Dzifče** počarovne, **ňezal'up** še do mñe  
ja chlapec vandrovni, co ci budze zo mñe

Ja chlapec **vandrovni**, zo šireho poľa  
co ci budze zo mñe, frajirečko moja?

- ▶ **Dzifče** počarovne, **nezal'ub** se do mne  
ja chlapec **vandrovny**, co ci budze zo mñe

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 **vandrovny** zo šireho poľa  
co ci budze zo mne draha duša moja?



*Femplin*

# Technológie

- ▶ Programovací jazyk Java
  - ▶ Knižnica JOOQ
- ▶ Databáza
  - ▶ MySQL



*Femplin*

# 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();
```



Femplin

# Príprava databázy

- ▶ Scraping a rôzne formy skriptovania
  - ▶ Dáta bolo treba niekde zohnať
    - ▶ Viktor Gliganič - primáš (1. huslista) Ľudovej hudby FS Zemplín
    - ▶ Piesne372



# „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



Femplin



# Čistenie a parsovanie textov

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



*Zemplin*

# Algoritmus

- ▶ Potrebujeme algoritmicky porovnať dokumenty
  - ▶ Information retrieval
- ▶ Myšlienka:
  - ▶ Vytvorím tabuľky
    - ▶ 1. riadok = slová, 2. riadok = početnosť
  - ▶ Získam vektor = druhý riadok
  - ▶ Vektory porovnáam a získam ich algebraický vzťah



Femplin

# Vector space algorithm v rýchlosti

- ▶ „Ja som veta“

ja	som	veta
1	1	1

- ▶  $\langle 1, 1, 1 \rangle$

- ▶ „Ja som veta, veta som ja“

ja	som	veta
2	2	2

- ▶  $\langle 2, 2, 2 \rangle$



# TermScheme - n-gramy

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



*Femplin*

# Vector space algorithm v rychlosti

## ► „Un-gram

Ja	som	veta
1	1	1

## ► Bi-gram

Ja som	som veta
1	1

## ► Tri-gram

Ja som veta
1



Zemplin

► Máme teda vektory...

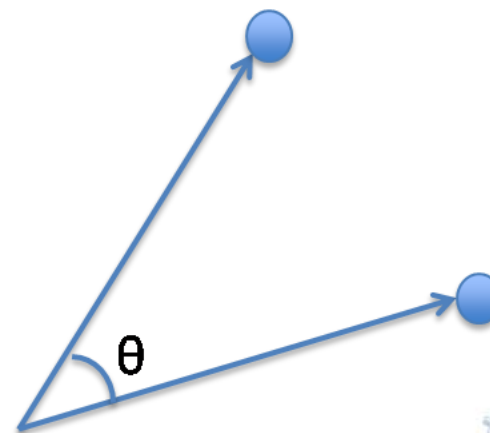


*Femplin*

# Vzdialenosť dvoch vektorov

- ▶ Kosínusová miera
  - ▶ Podiel skalárneho súčinu a súčinu euklidovských dĺžok vektorov

$$\text{sim}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



# Vector inclusion

Poznámka - pracujeme s  
dimenziami vektora



# Vector inclusion

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



Femplin



„Slová“  
textu A

„Slová“  
textu B

# Vector inclusion

- ▶ Vektor  $v_1 = \langle A, B, C \rangle$  |  $\langle a_1, b_1, c \rangle$ 
  - ▶ Repräsentuje text A
- ▶ Vektor  $v_2 = \langle B, A, E \rangle$  |  $\langle b_2, a_2, e \rangle$ 
  - ▶ Repräsentuje text B



# Vector inclusion

▶ A forma - Slovu **C** vo  $v_2$  priradíme 0

▶  $\langle \mathbf{A}, \mathbf{B}, \mathbf{C} \rangle$

▶  $v_1 = \langle a_1, b_1, c \rangle$

▶  $v_2 = \langle a_2, b_2, 0 \rangle$

▶  $\langle \mathbf{B}, \mathbf{A}, \mathbf{E} \rangle$  zmeníme na  $\langle \mathbf{A}, \mathbf{B}, - \rangle$



Femplin

# Vector inclusion

▶ B forma - Slovu **E** vo  $v_1$  priradíme 0

▶  $\langle \mathbf{B}, \mathbf{A}, \mathbf{E} \rangle$

▶  $v_1 = \langle b_1, a_1, 0 \rangle$

▶  $v_2 = \langle b_2, a_2, e \rangle$

▶  $\langle \mathbf{A}, \mathbf{B}, \mathbf{C} \rangle$  zmeníme na  $\langle \mathbf{B}, \mathbf{A}, - \rangle$



Femplin

Všetky „slová“

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

„Slová“  
textu A

„Slová“  
textu B

# Vector inclusion

- ▶ Zjednotenie

  - ▶  $\langle A, B, C, E \rangle$

- ▶ Prienik

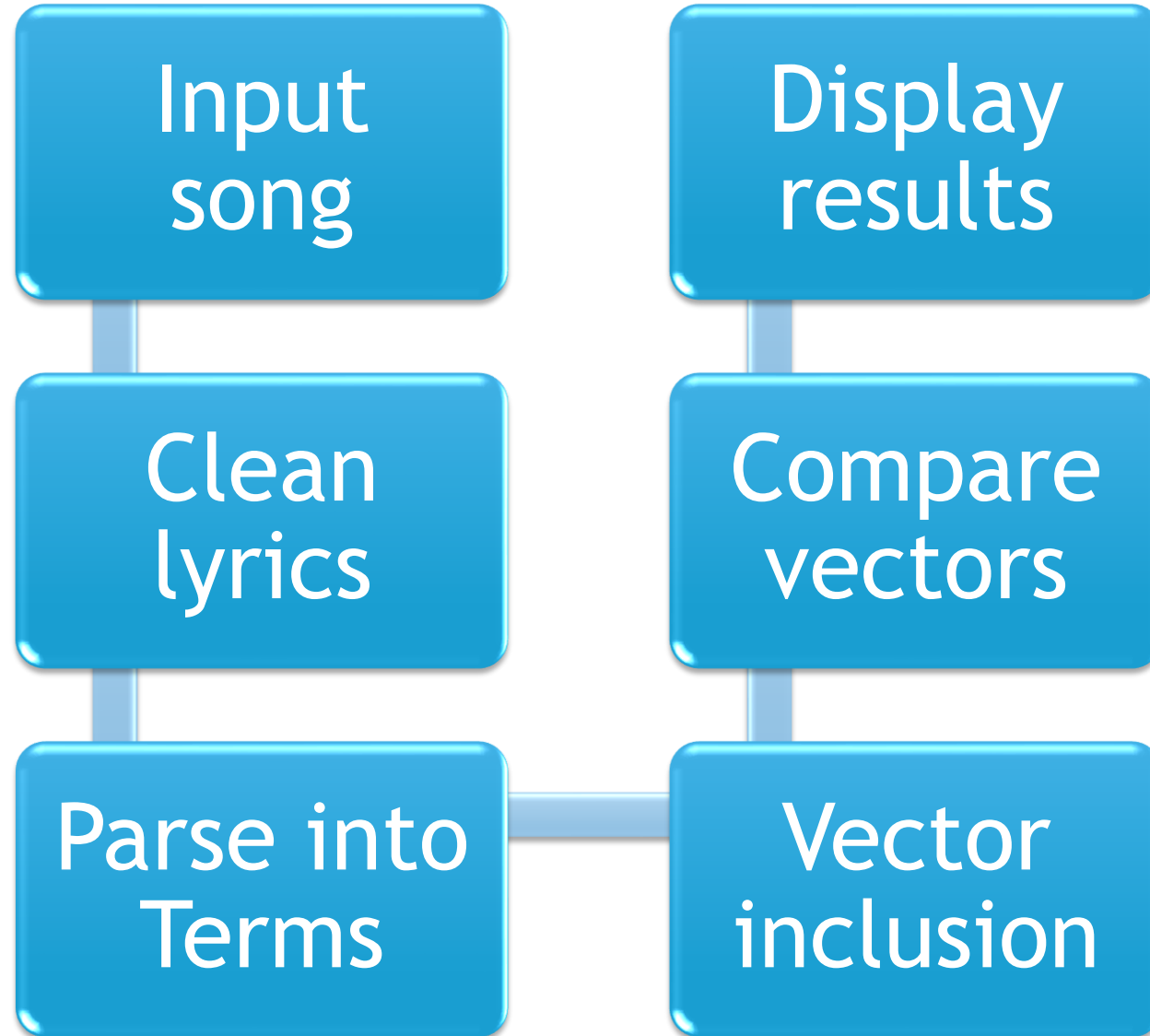
  - ▶  $\langle A, B \rangle$

- ▶ Všetky slová databázy

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



Femplin





# Zobrazenie výsledkov

- ▶ Pre každú pieseň  $p$  vrátim mapu  $s \rightarrow \%$  podobnosť, kde  $s$  je každá pieseň rôzna od  $p$



Femplin

# Pokročilé porovnávanie

# Term weighting

- ▶ „Slová“ v texte sú pre nás rôzne zaujímavé
  - ▶ „A ja taka **čarna** jak **čarna čarnica**...“
  - ▶ „A ja taka **dzivočka cingi lingi bom**...“



Femplin

TF

- ▶ Zohľadnenie frekvencie pojmu v **jednom dokumente**
- ▶ Rôzne vyhladzovacie techniky



*Femplin*

# TF

## Logarithm tf

Idea logaritmického tf je využití logaritmus jako škálovací mechanismus. Nech  $n_{p,d}$  je početnost pojmu  $p$  v dokumente  $d$ . Potom

$$\text{logtf}(p, d) = 1 + \log(n_{p,d})$$



Jyväskylä

# IDF - inverse document frequency

- ▶ Zohľadnenie frekvencie pojmu v rámci celej kolekcie
- ▶ Rôzne vyhladzovacie techniky



*Femplin*

# IDF - inverse document frequency

## **Idf**

Nech  $n_{doc}$  je počet dokumentov v našej kolekcii. Nech  $df_p$  je početnosť pojmu  $p$  vo všetkých dokumentoch. Potom

$$idf(p) = \log(n_{doc}/df_p)$$



Zemplin

# TfIdf

- ▶ Zohľadnenie oboch
- ▶ Napríklad  $tf * idf$





# 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...



Femplin

# Term comparison algorithm

- ▶ Keď už porovnávame vektory, môžeme porovnávať aj „slová“
- ▶ Neoplatí sa zohľadňovať sémantiku slov



*Femplin*

# Jazykovedný prístup

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



*Femplin*

# 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



*Femplin*

# Term similarity

- ▶ Levenshteinova vzdialenosť
  - ▶ Počet operácií potrebných na zmenu  $s_1$  na  $s_2$
  - ▶ Vloženie, vymazanie a zmena písmena
    1. **k**itten → **s**itten (substitution of "s" for "k")
    2. **s**itten → **s**ittin (substitution of "i" for "e")
    3. **s**ittin → **s**itting (insertion of "g" at the end).
- ▶ Mnoho ďalších...



Femplin

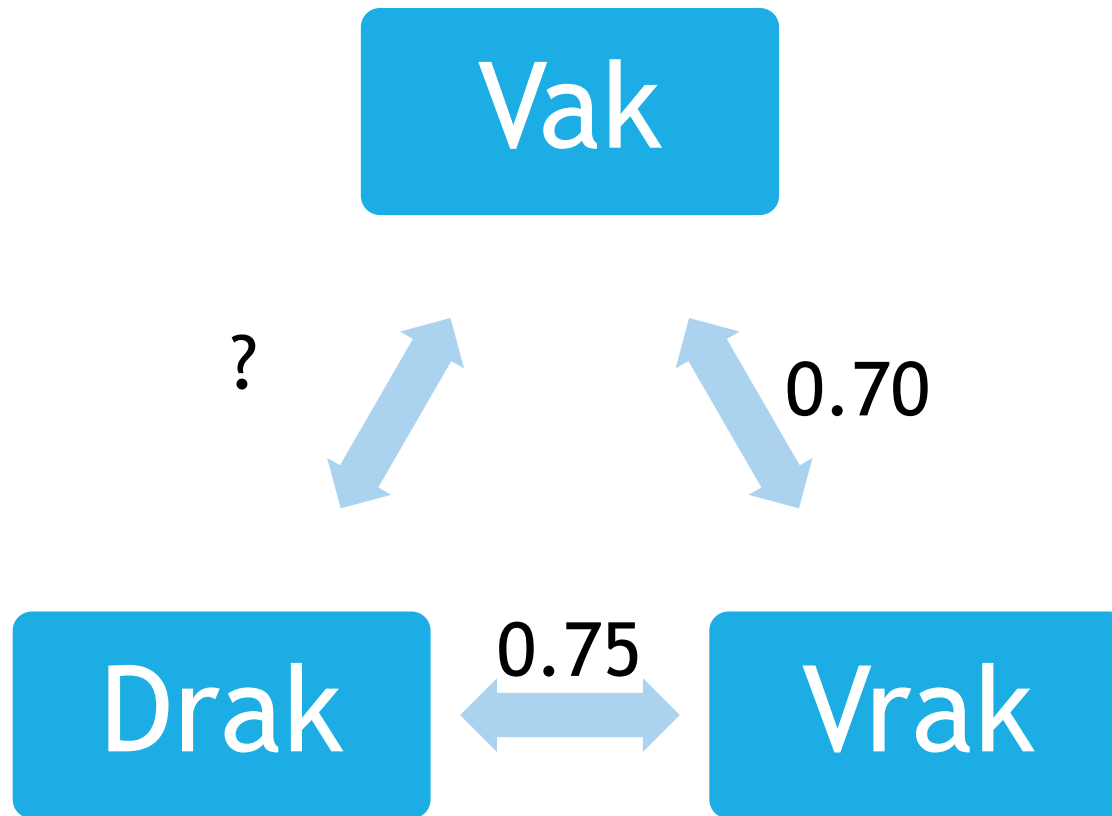
# Indirect term similarity

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

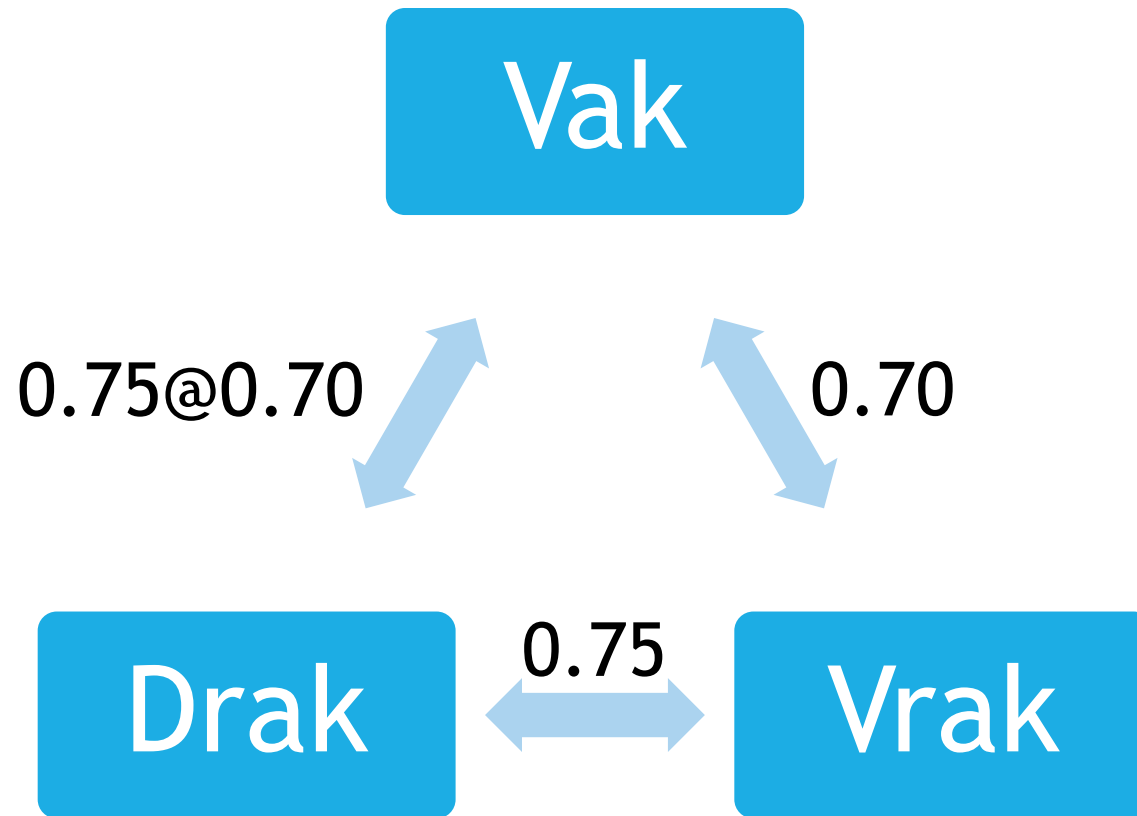


*Femplin*

# Indirect term similarity



# Indirect term similarity





# Indirect term similarity

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



Femplin

Praktickejšie riešenie:

# Term comparison + Tolerance

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



*Femplin*

# Tolerancia

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

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



Femplin

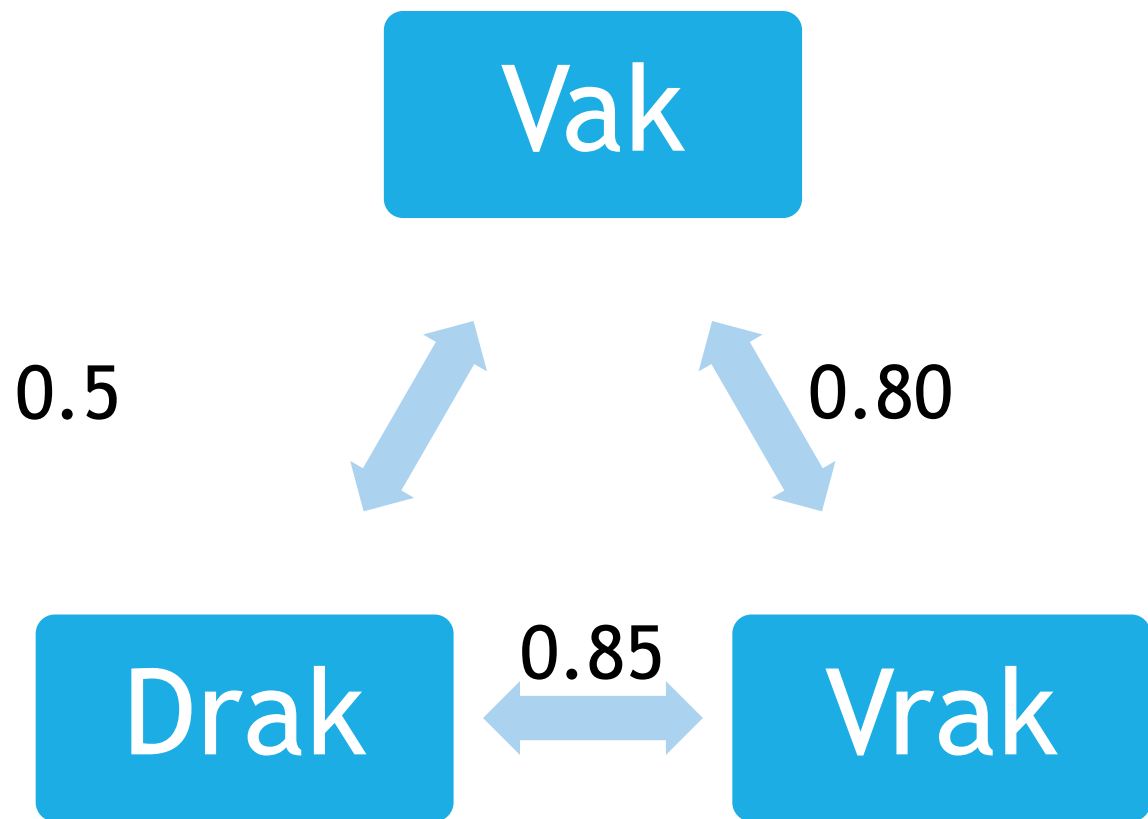
# Term grouping

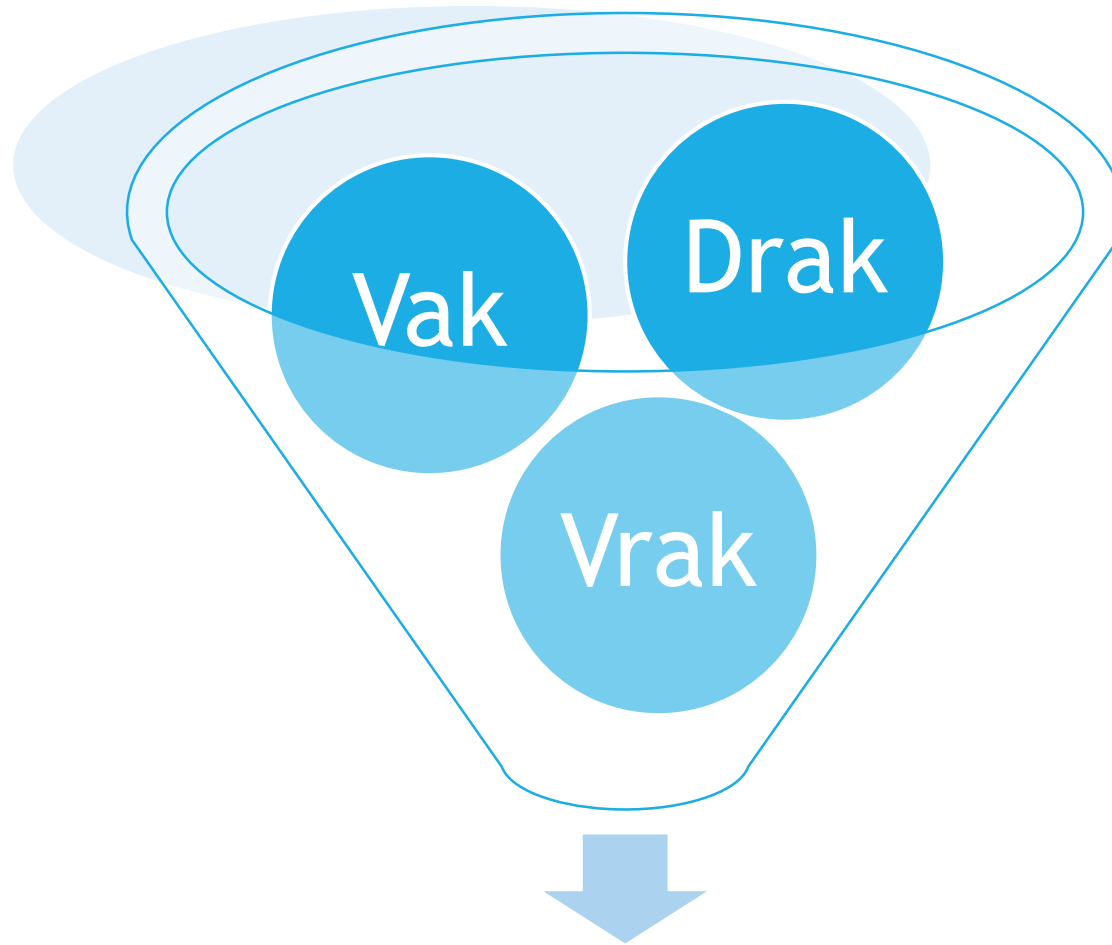
- ▶ TermComparisonAlgorithm = naivný
- ▶ Tolerance = MEDIUM = 0.8



*Zemplin*

# Term grouping





**TermGroupXYZ**

# Term grouping

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



*Femplin*



# Rôzne prístupy a taktiky

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

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



Femplin

# 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

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

Cosine Coefficient (what we studied)

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

Overlap Coefficient



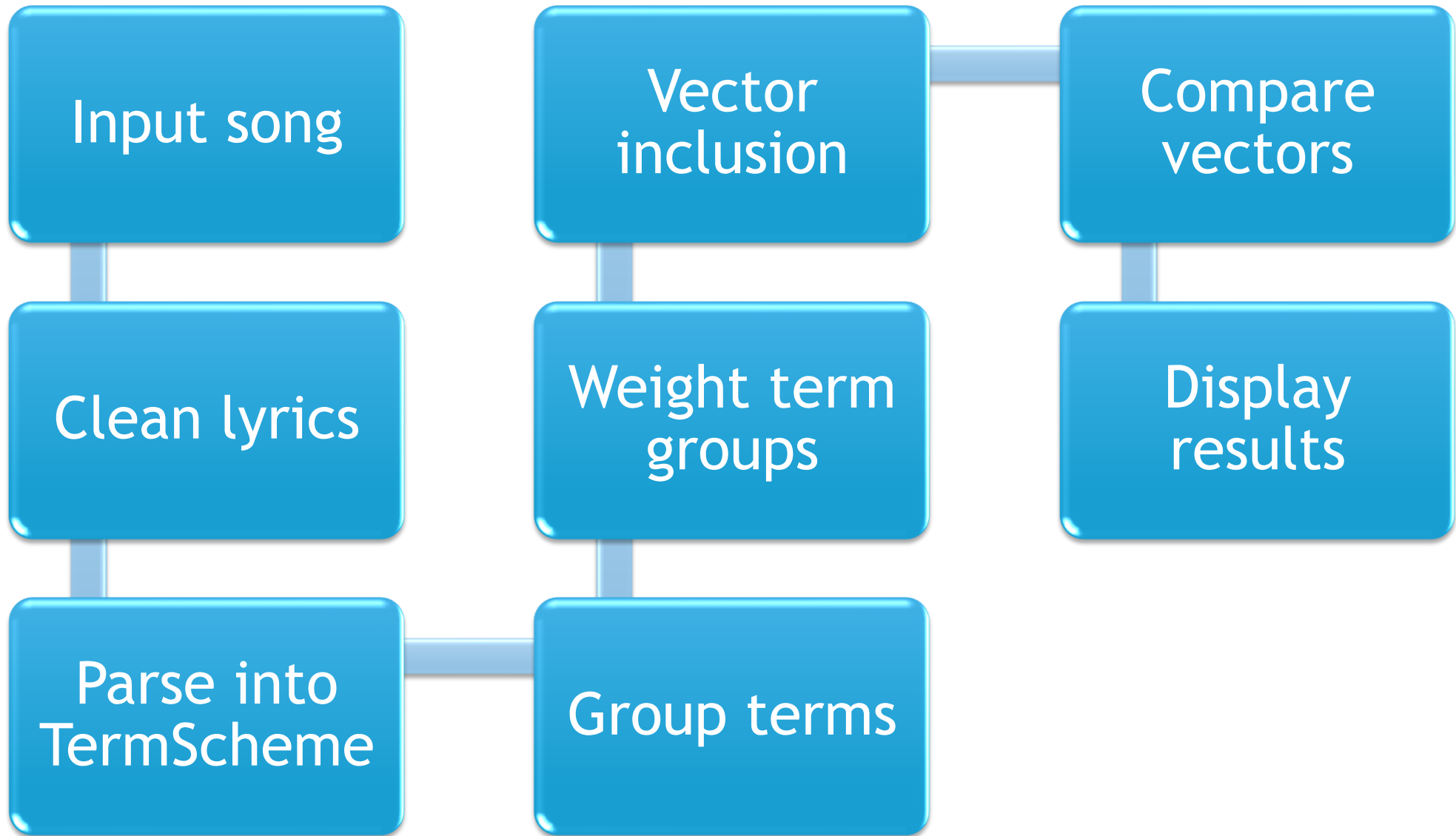
Femplin

# Vector algorithm configuration

- ▶ Objekt, ktorá definuje každú fázu algoritmu na základe enum konštanty



*Femplin*



```
/**
 * 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;
```

# Plány do budúcnosti

- ▶ Štatistické vyhodnotenie výsledkov behu pri rôznych konfiguráciách
- ▶ „Mergovací“ systém
- ▶ API jednoducho využiteľné v iných aplikáciách
  - ▶ Prepojenie so spevníkom pre Android



*Femplin*

# 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/>



*Zemplin*

Ďakujem za pozornosť



*Femplin*