

Získavanie metadát o vzťahoch a obsahu na webe

Tomáš Uherčík
Ing. Marián Šimko

Motivácia

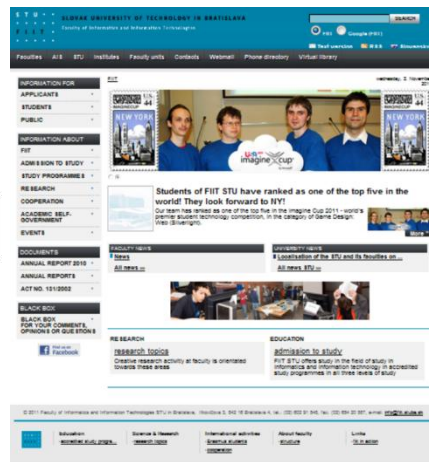
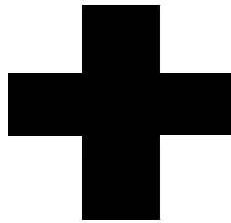
- ▶ Problém – Automatické získavanie anotácií o internetových zdrojoch (prepojeniach) na webe
- ▶ **Kde sú metadáta ?**
 - Mikroblogy
 - Priamo v obsahu

Motivácia

- ▶ Informácie z mikrobloggerov:
 - dezinformácie
 - zavádzajúce informácie
 - reklamní boti
 - nálady, pocity, subjektívne názory používateľov
- ▶ Z obsahu nezískame niektoré metadáta
 - obrázky
 - videá
 - kvalita (hodnotenie) dobrý? zlý?

Návrh metódy

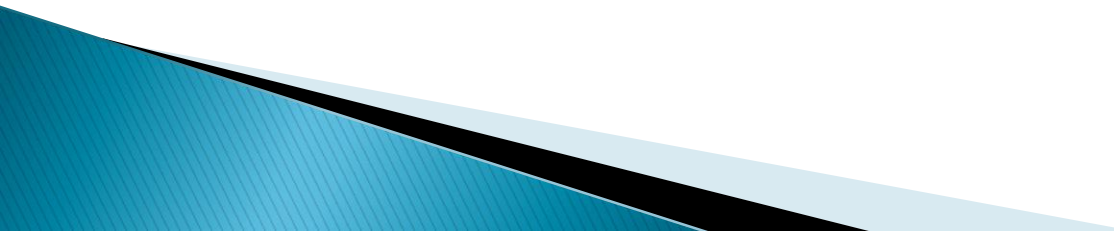
- ▶ Získať kľúčové slová z mikrobloggerov
- ▶ Získať kľúčové slová z analýzy zdroja
- ▶ Vybrať z nich tie najlepšie



Klíčové slová z mikroblogov

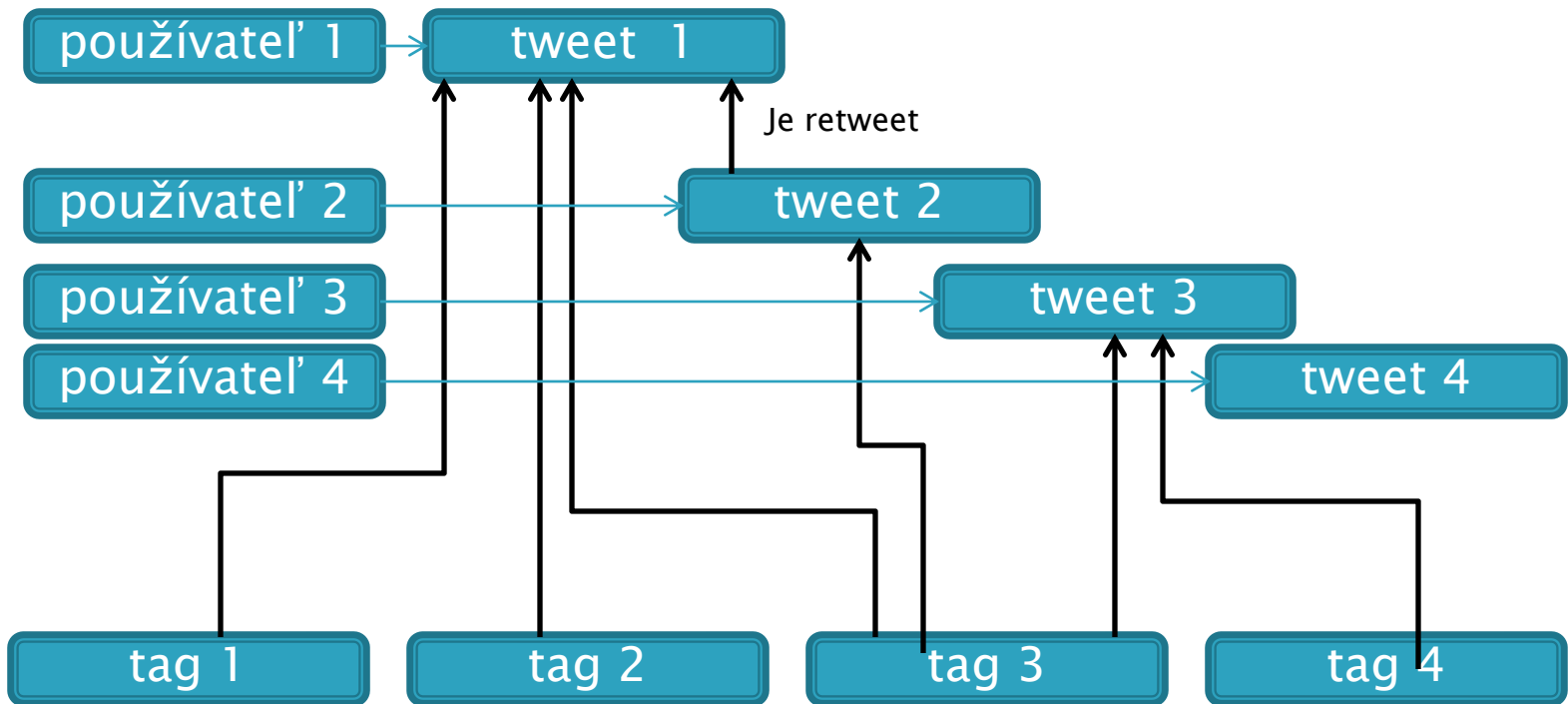
- ▶ Relevancia tagu:
 - Ohodnotenie používateľa
 - Tweet, z ktorého pochádza
 - Hashtag
 - Počet výskytov slova, TF-IDF, TextRank

Používatelia mikrobloggerov

- ▶ PageRank
 - ▶ HITS
 - ▶ TrustRank
 - ▶ TunkRank
 - ▶ NodeRank
- 

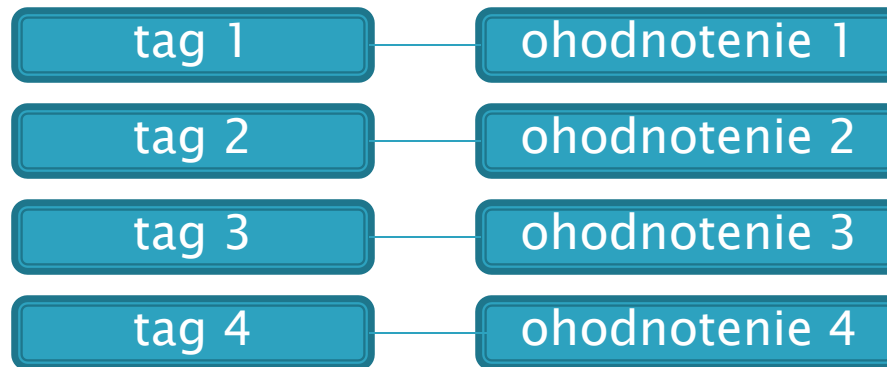
Kľúčové slová z mikrobloggerov

- ▶ Príklad pre vyhľadávanie zdroja z



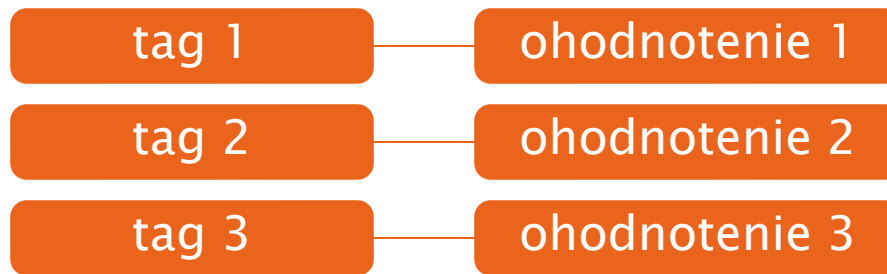
Kľúčové slová z mikroblov

- ▶ Množina tagov spolu s ich ohodnotením



Analýza textu

- ▶ Získanie množiny tagov spolu s ich rankingom



Výsledná množina klíčových slov

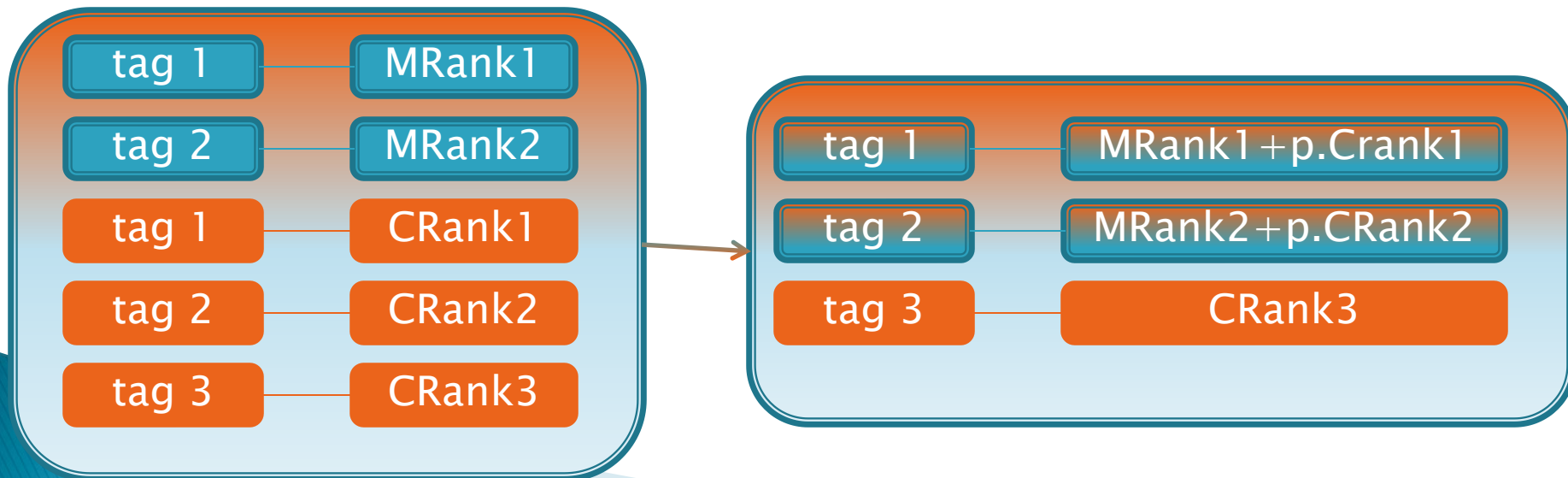
- ▶ Zoradená množina z oboch zdrojov
- ▶ Zorad'ovacie funkcie:

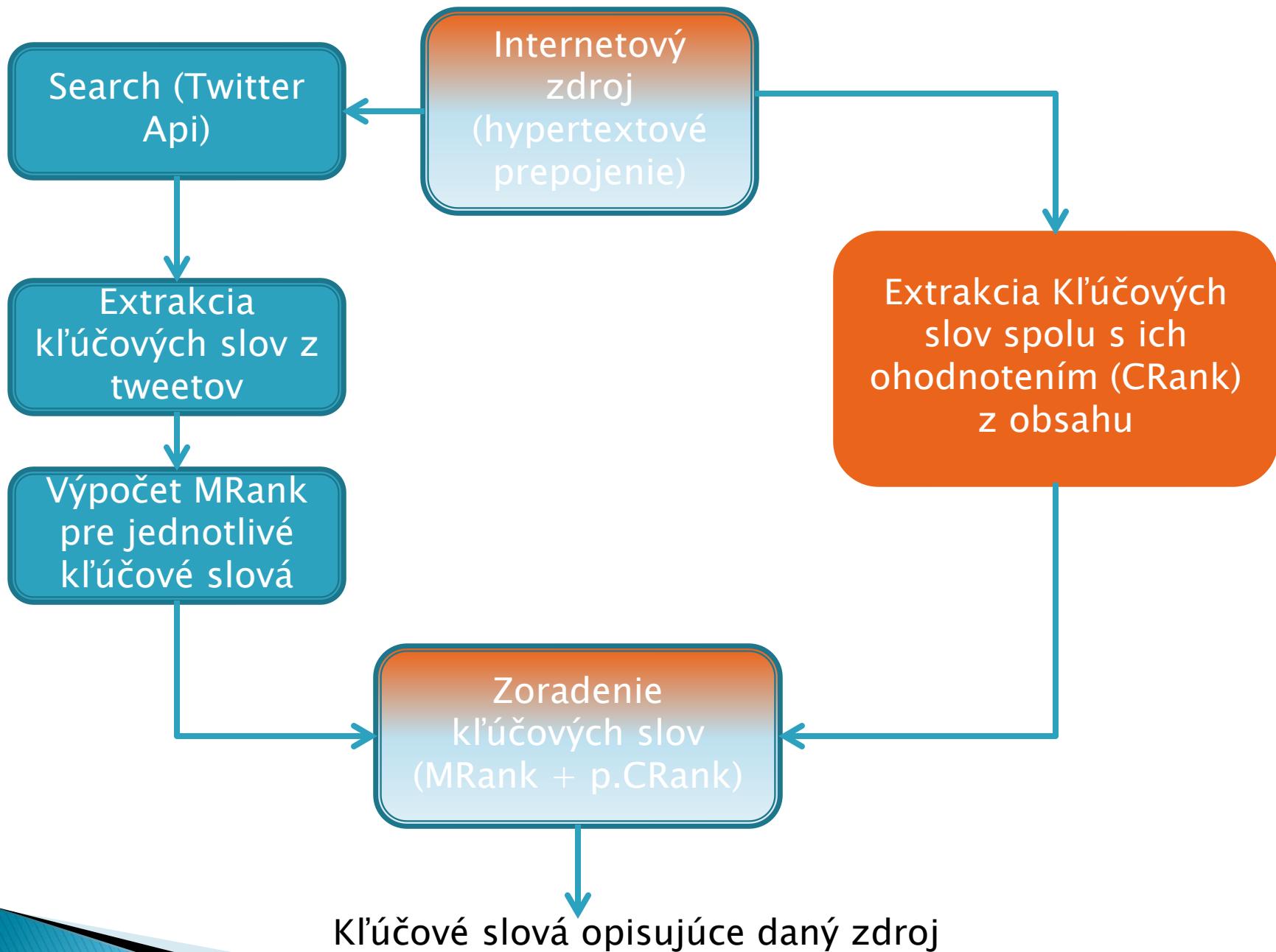
- Pre mikroblogy
- Pre obsah(content)

MRank(t)

CRank(t)

$$\text{Rank}(t) = \text{MRank}(t) + p \cdot \text{CRank}(t)$$





Overenie

- ▶ Použitie recent liniek z delicious.com
- ▶ Porovnanie kľúčových slov

