

# Získavanie metadát o vzťahoch a obsahu na webe

Tomáš Uherčík  
Ing. Marián Šimko

# Motivácia

## ▶ Problém

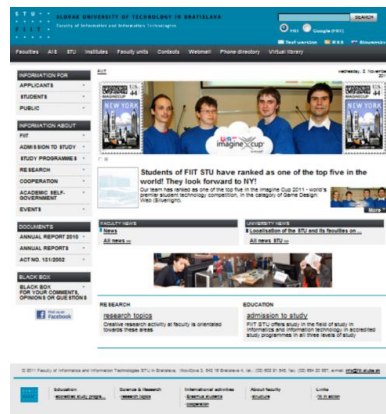
automatické získavanie metadát o internetových zdrojoch (prepojeniach) na webe

## ▶ Kde hľadať metadáta ?

- Obsah
- Prepojenia
- Záznamy o používaní webu
- Externé zdroje, ktoré obsahujú referenciu na *URL*
  - mikroblogy

# Návrh metódy

- ▶ Získanie kľúčových slov z mikroblogov
- ▶ Získanie kľúčových slov z textovej analýzy zdroja
- ▶ Výber tých najlepších z nich



## Metadáta

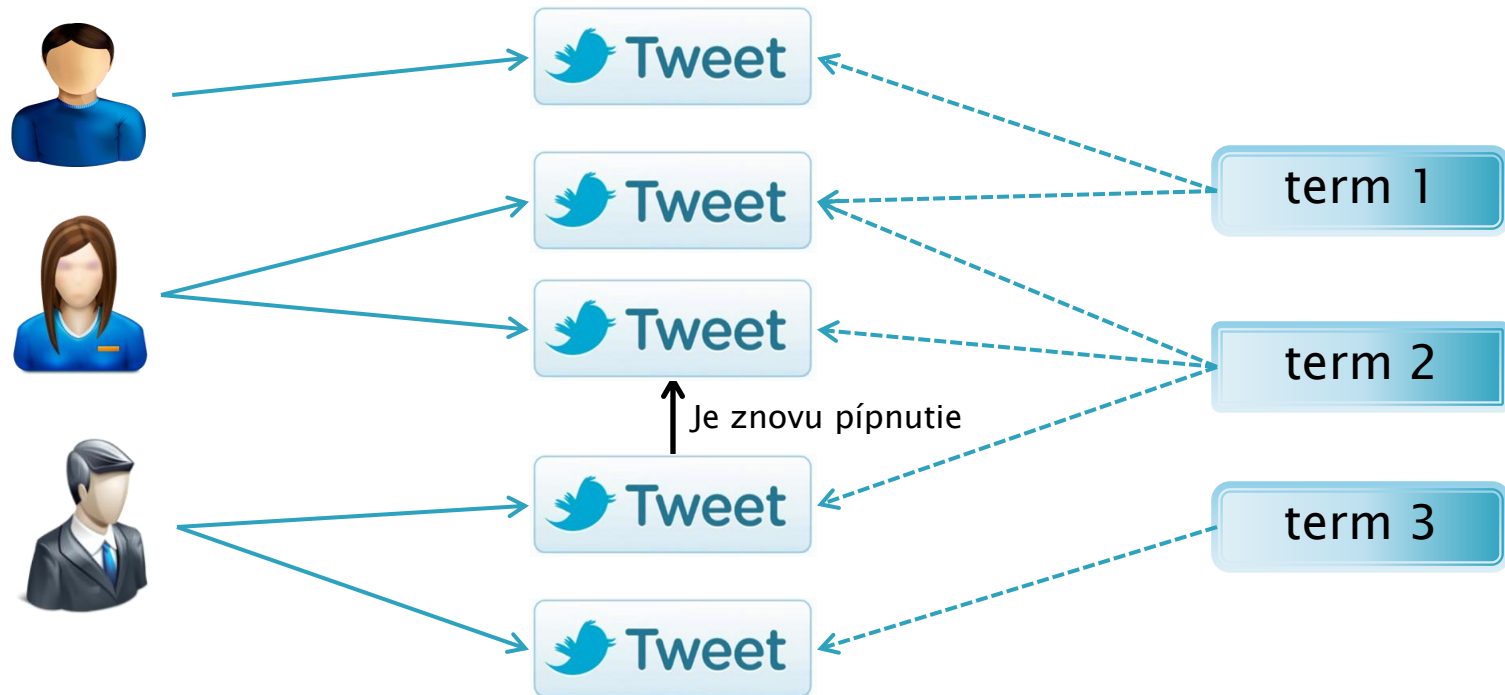
Kľúčové slovo 1

Kľúčové slovo 2

Kľúčové slovo 3

# Kľúčové slová z mikrobloggerov

- ▶ Vyhľadanie zadanej *URL* v mikrobloggeru:



# Kľúčové slová z mikrobloggerov

- ▶ Relevancia kľúčového slova:
  - Relevancia v texte pípnutí
    - Text Rank
    - Alchemy API
  - Ohodnotenie používateľa, ktorý pípnutie napísal
    - Tunk Rank

$$TunkRank(X) = \sum_{Y=\text{nasledovníci}(X)} \frac{1 + p \cdot TunkRank(Y)}{|\text{nasledovníci}(Y)|}$$

kde  $p$  je pravdepodobnosť znovu pípnutia



# Úpravený Tunk Rank – ARank

- ▶ Tunk Rank obohatený o závislosť od periódy publikovania pípnutí:

$$ARank(X) = \sum_{Y=\text{nasledovníci}(X)} \frac{1 + \frac{p}{\log(T)} \cdot ARank(Y)}{|\text{nasledovníci}(Y)|}$$

kde:

$p$  – pravdepodobnosť znovu pípnutia

$T$  – perióda ako často používateľ publikuje pípnutia

$X$  – používateľ

# Výsledná množina klíčových slov

- ▶ Zoradená množina z oboch zdrojov:

- Pre *mikroblogy*

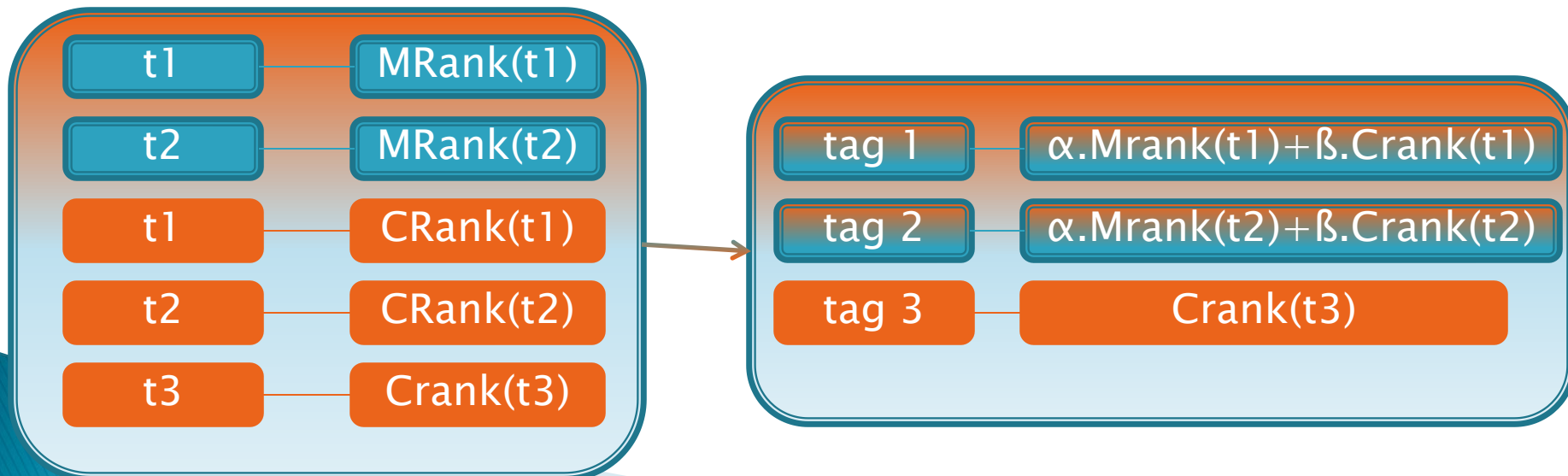
**MRank(t, ARank(t))**

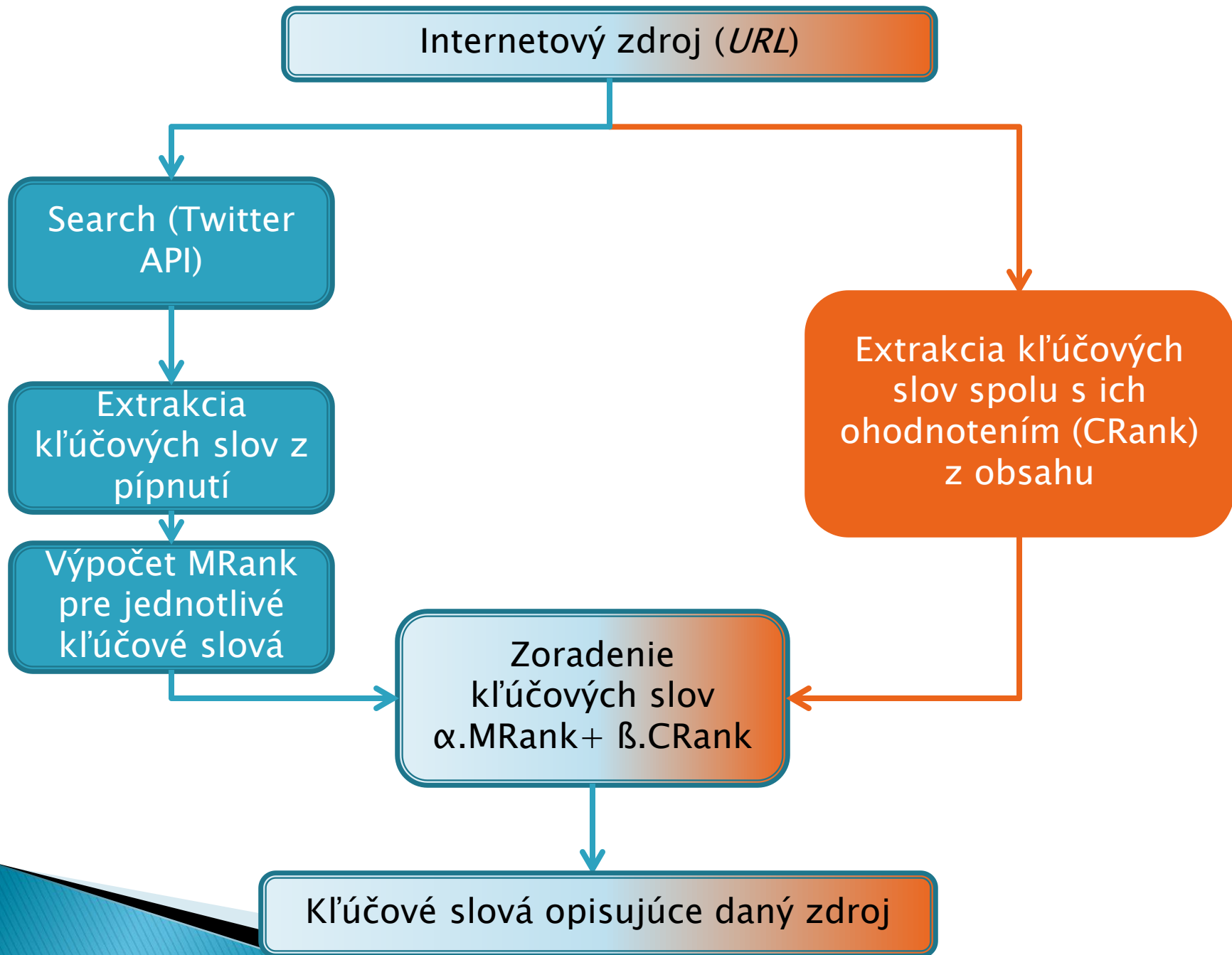
- Pre *obsah(content)*

**CRank(t)**

$$\text{Rank}(t) = \alpha \cdot \text{MRank}(t, \text{ARank}(t)) + \beta \cdot \text{CRank}(t)$$

kde  $t$  - klúčové slovo,  $\alpha$ ,  $\beta$  - normalizačné konštanty





# Overenie

- ▶ Vyhodnotiť prínos obohatenia kľúčovými slovami z mikrobloggerov:
  - ľudskí experti
  - porovnanie s kľúčovými slovami z tagovacích systémov
  - porovnanie s kľúčovými slovami z analýzy textu
- ▶ Overiť na stálom datasete

# Ďalšia práca

- ▶ Vystavenie webovej služby na extrakciu kľúčových slov o *URL*
  - Možnosť rozšírenia služby *metall*
- ▶ Vylepšenie už implementovaných častí:
  - Sofistikovanejší spôsob porovnávania kľúč. slov
  - Vylepšenie implementácie Text Rank
  - Vyhodnotenie typu obsahu *URL*
    - Automatické nastavenie hodnoty norm. konštánt  $\alpha$  a  $\beta$

$$\text{Rank}(t) = \alpha \cdot \text{MRank}(t, \text{ARank}(t)) + \beta \cdot \text{CRank}(t)$$

# Príklad extrakcie kľúčových slov

- ▶ Pre video: *President Obama Speaks about Insourcing American Jobs*

Twitter	Obsah URL	Zjednotené
american jobs	House business leaders	President Obama
insourcing	President Obama	insourcing
President Obama	America	american jobs
		house business leaders
		America

# Príklad extrakcie kľúčových slov

- ▶ Pre video: Santigold L.E.S. Artistes

Twitter	Obsah <i>URL</i>	Zjednotené
santigold	Nima Nourizadeh	santigold
Youtube video		Youtube video
L.E.S. Artistes		L.E.S. Artistes
		Nima Nourizadeh

