



# Hľadanie vzťahov medzi kľúčovými slovami

**Peter Kajan**

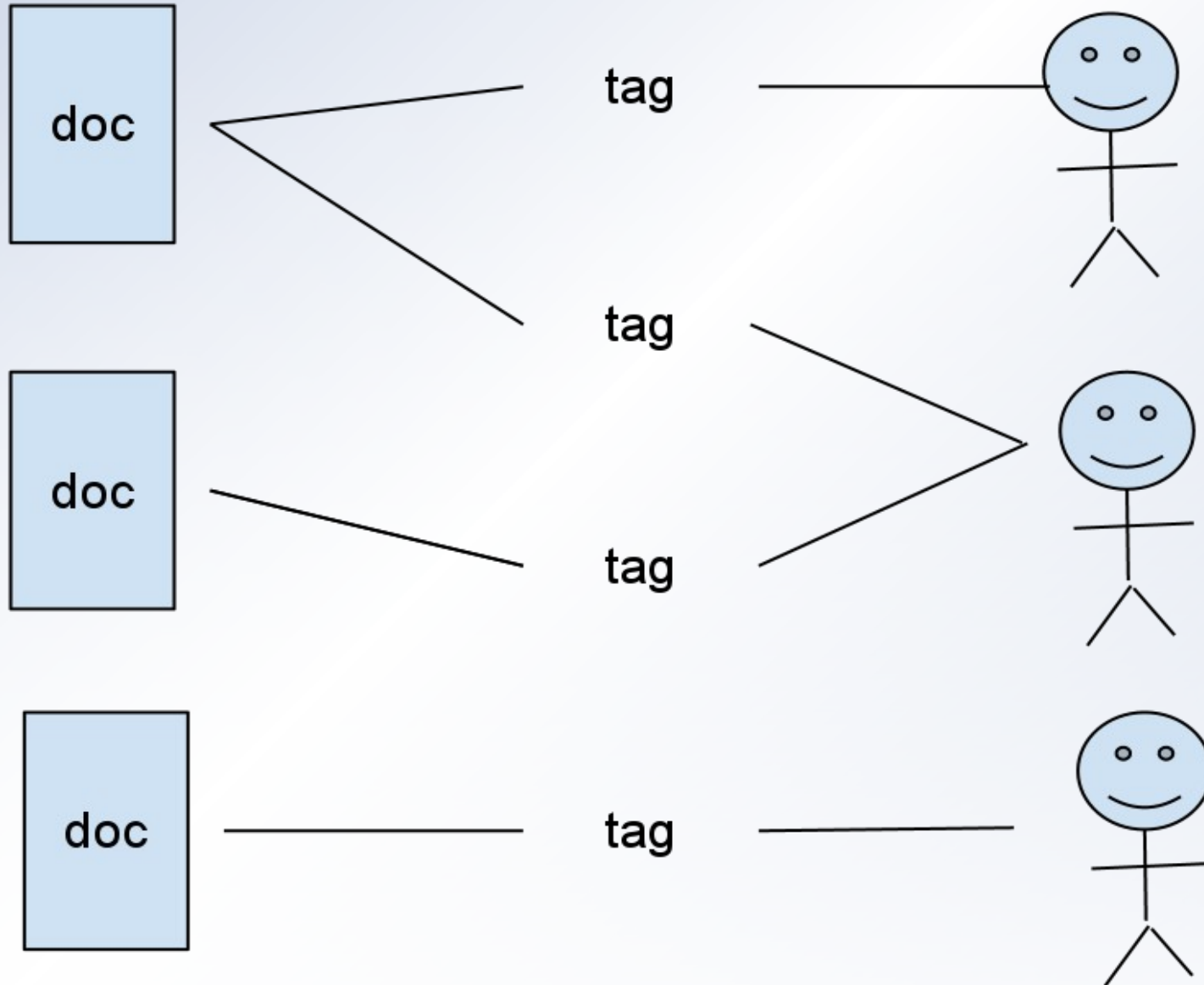
**Vedúci práce: Ing. Michal Barla, PhD**

Fakulta informatiky a informačných technológií  
Slovenská technická univerzita

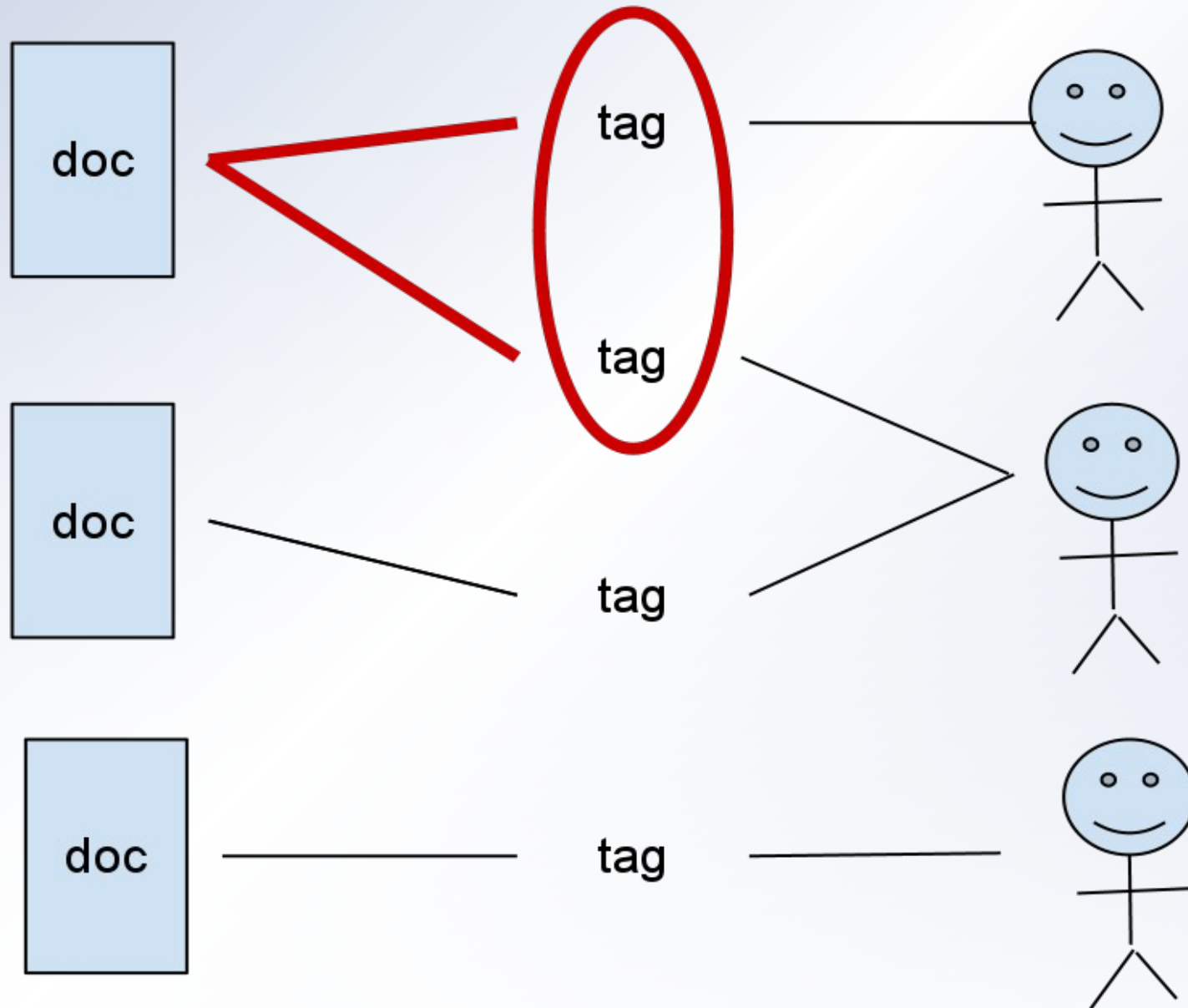
# Motivácia

- Príbuznosť objektov
  - Dokumentov (stránok)
  - Používateľov
  - Komunit ...
- Klasifikácia
  - ...
- Odporúčanie tagov
- Znalostná báza

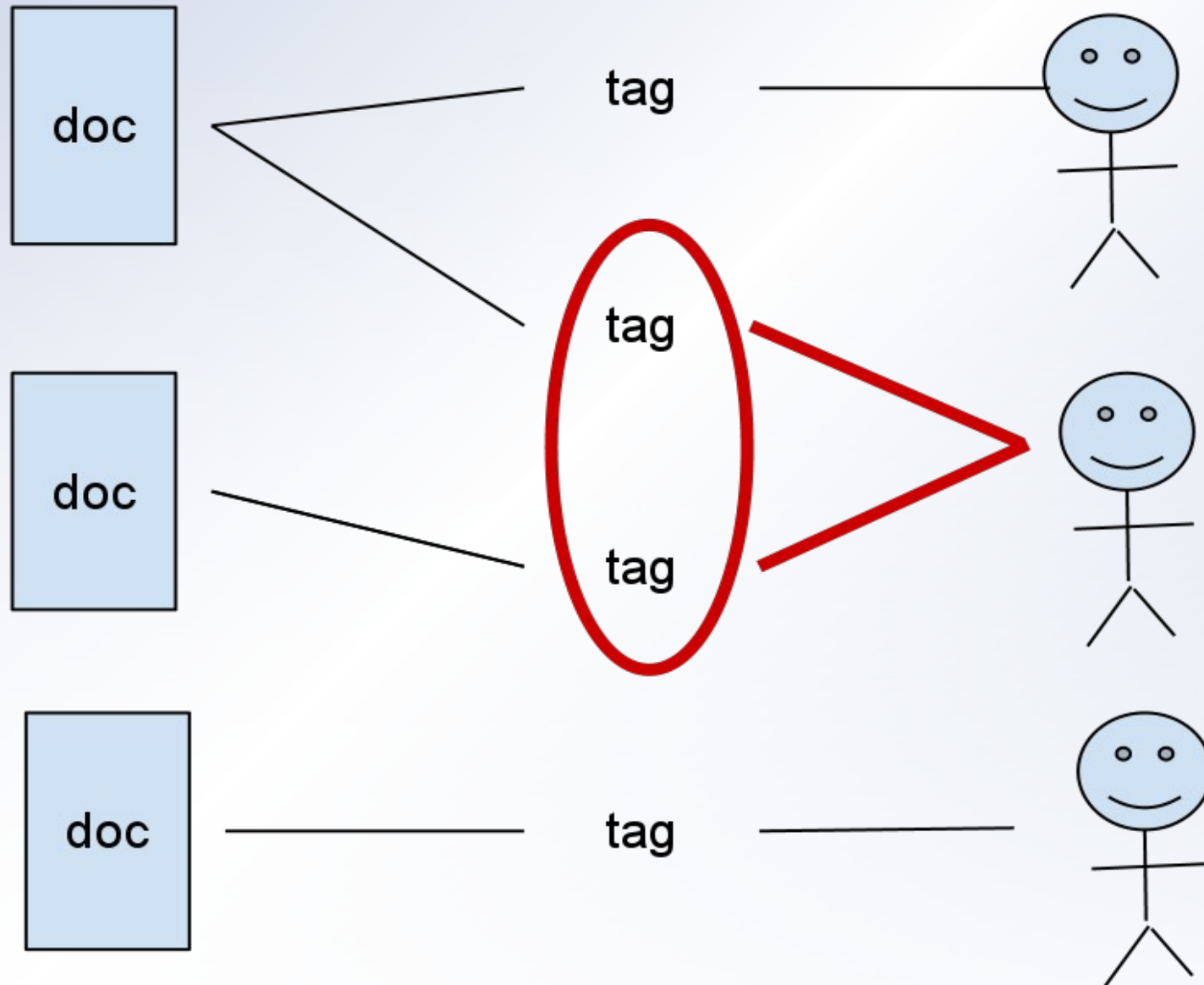
# Predošlé prístupy



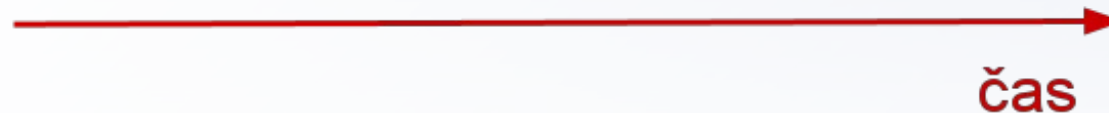
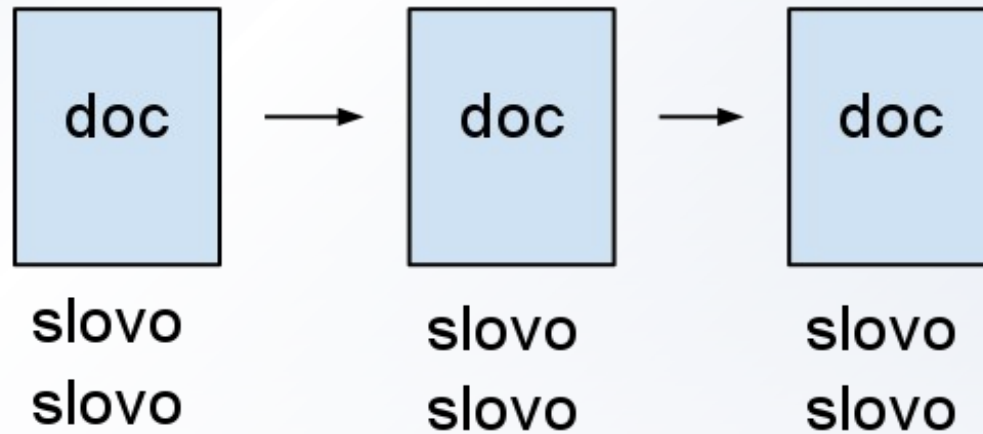
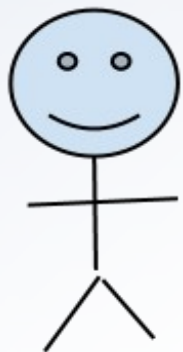
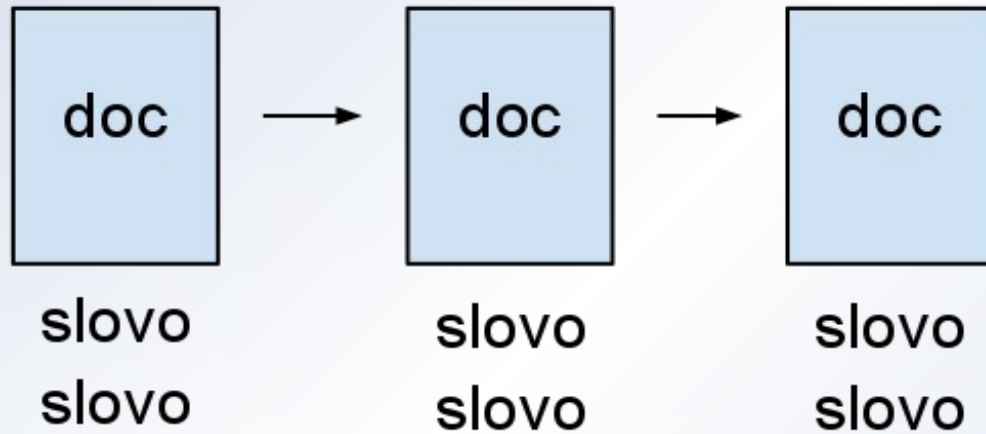
# Predošlé prístupy



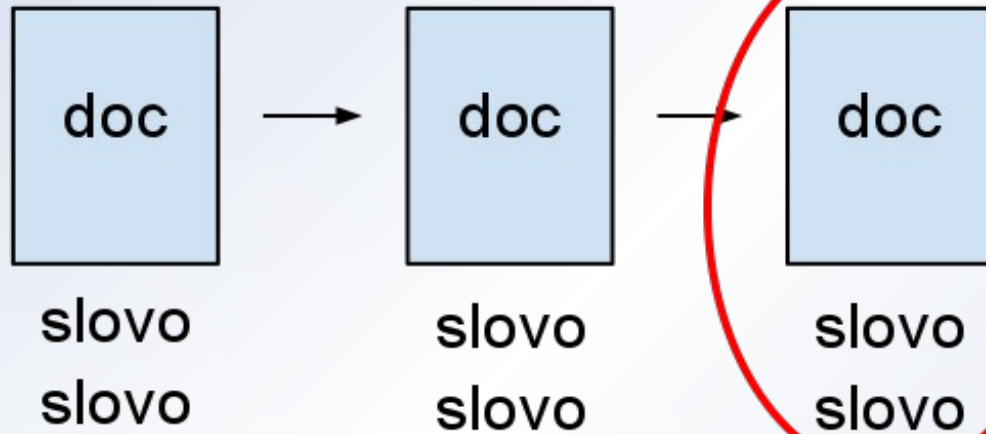
# Predošlé prístupy



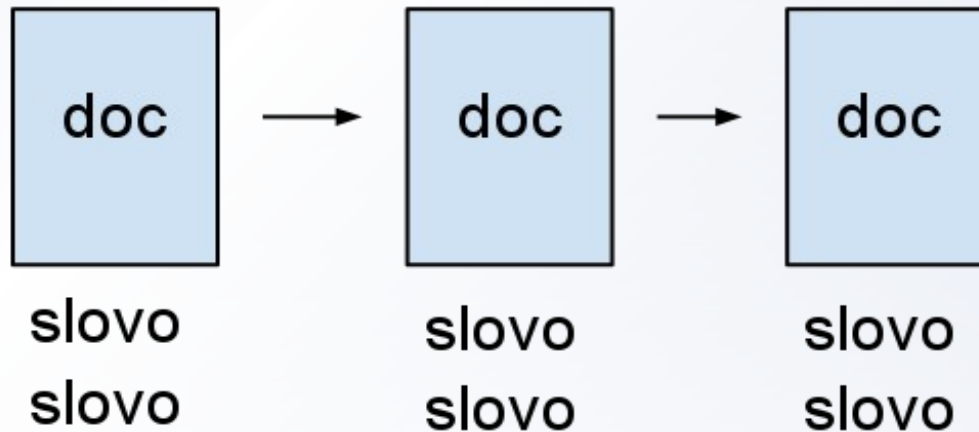
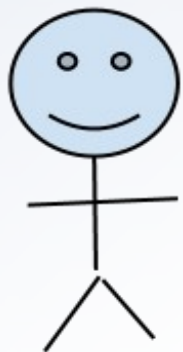
# Náš prístup



# Náš prístup

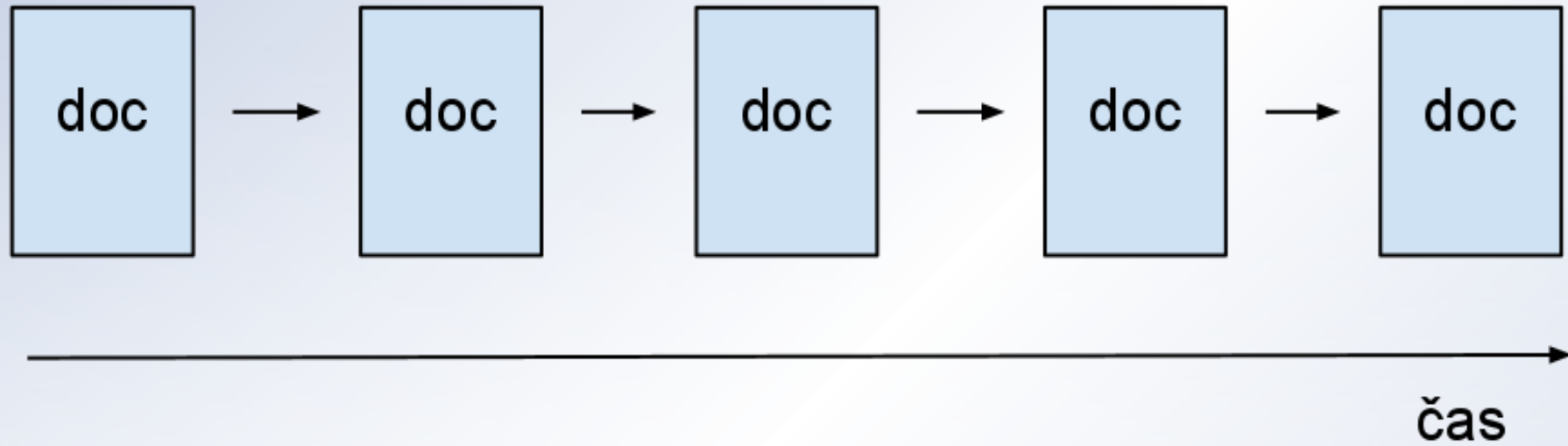


Stačí  
podobnosť  
stránok

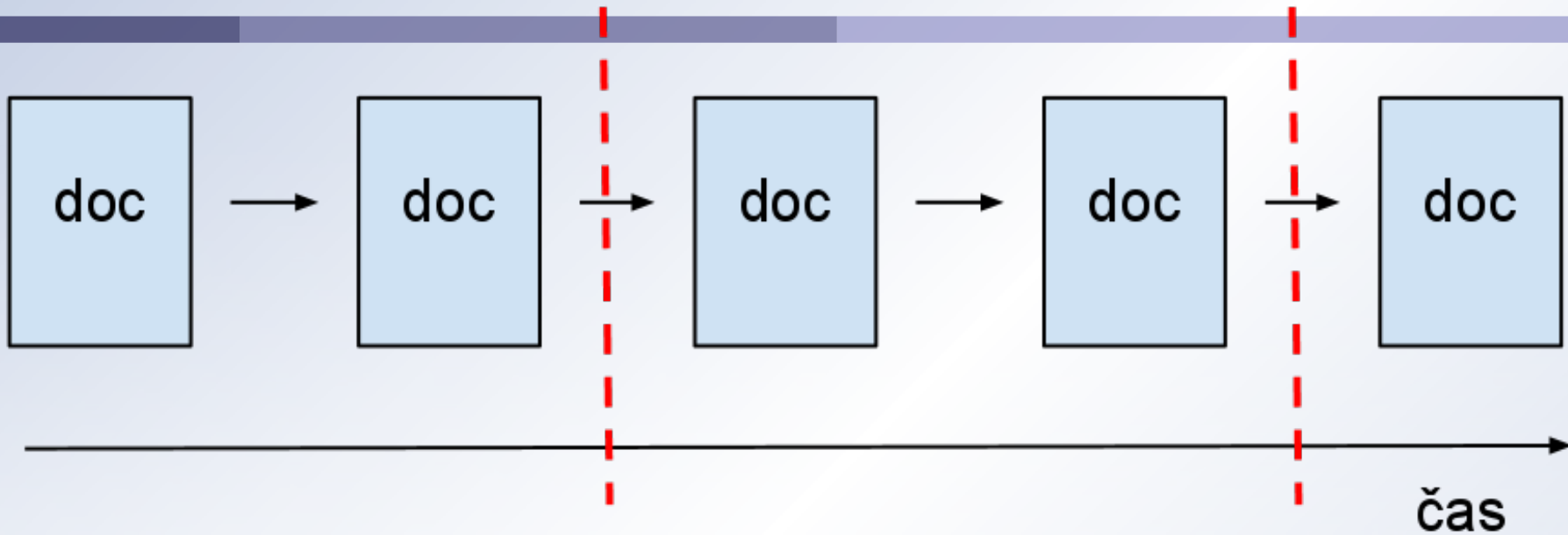


čas

# Sessions

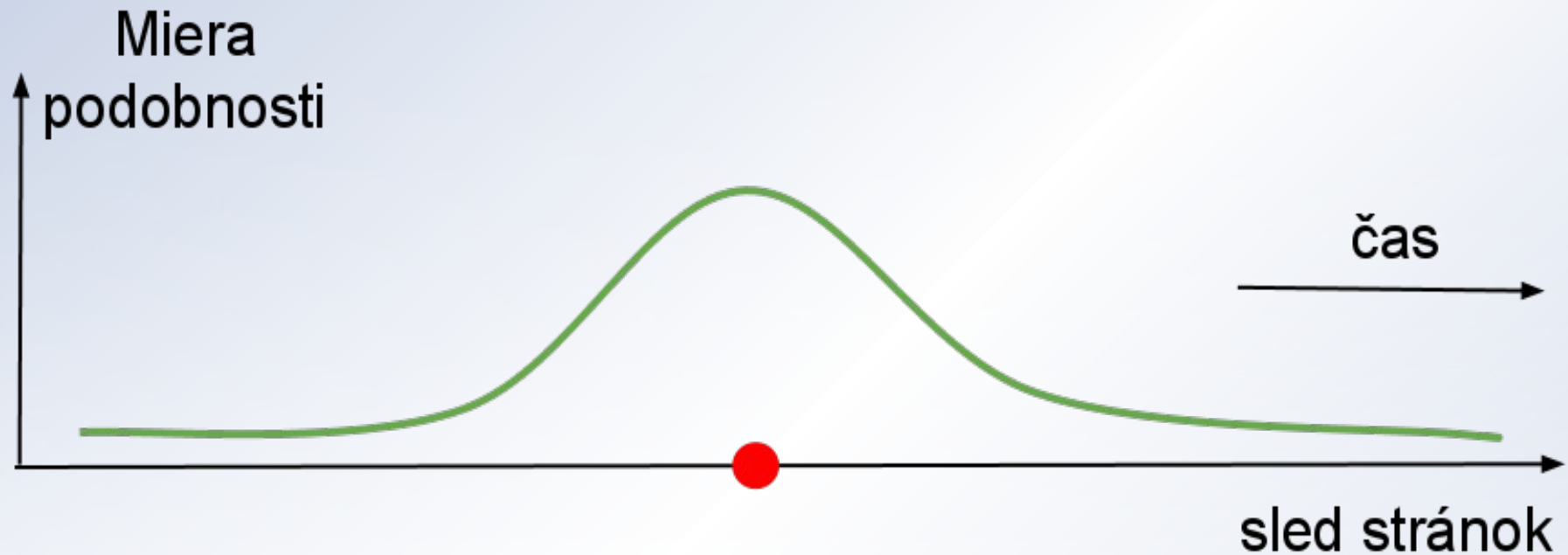


# Sessions



- Časové okno
- Levensteinova vzdialenosť klúč. slov
- Taxonómia (existujúca, Linked data)
- Naša taxonómia

## Podobnosť stránok vrámci session



- Zohľadniť čas strávený na stránke
- Penalizovať časté stránky (TF-IDF)

# Podobnost' klíčových slov

- Podobnost'

$$termSim(t_1, t_2) += pageSim(p_1, p_2) \quad \forall t_1 \in p_1, t_2 \in p_2, pageSim(p_1, p_2)$$

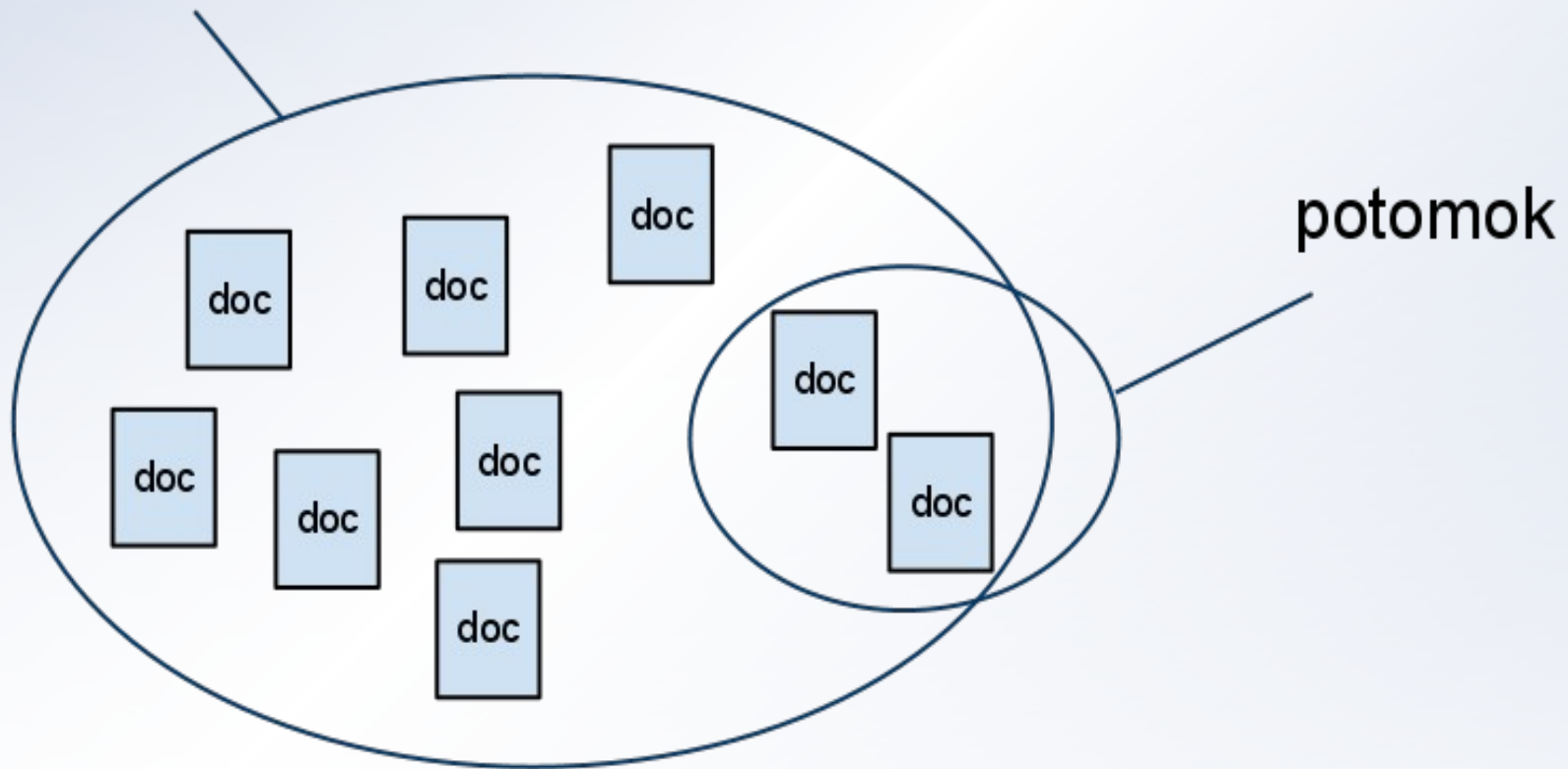
- Příslušnost'

$$termSim(t_1, t_2) += tfIdf(t_1) * pageSim(p_1, p_2)$$

- Podobnost' z příslušností
- Rovnako pre stránky

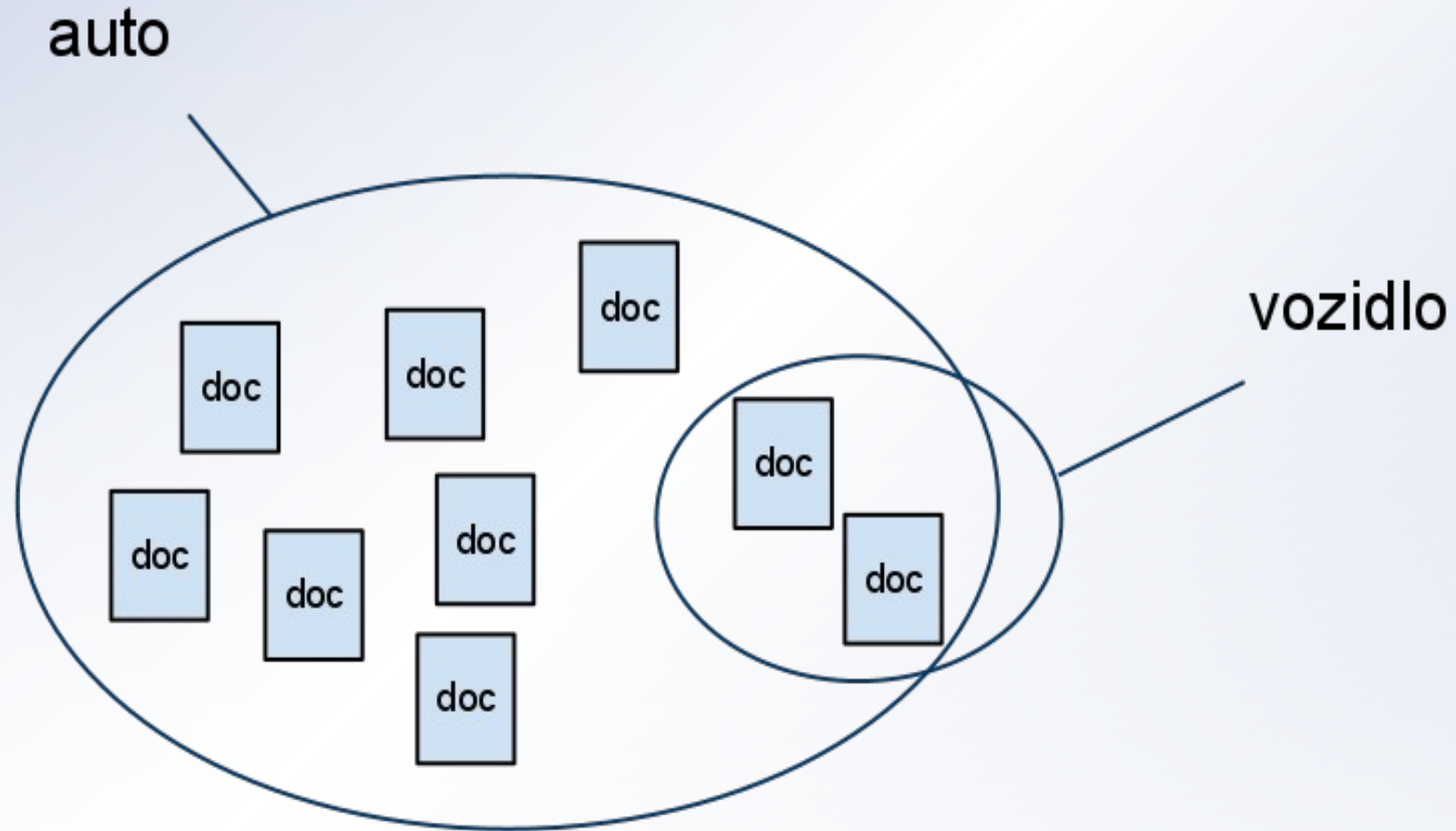
# Hierarchia - predošlé prístupy

predok

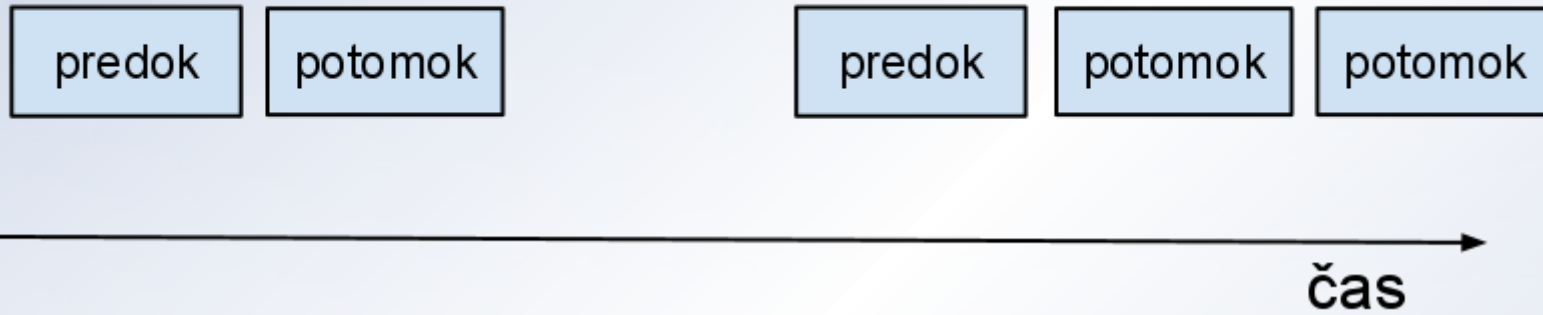


potomok

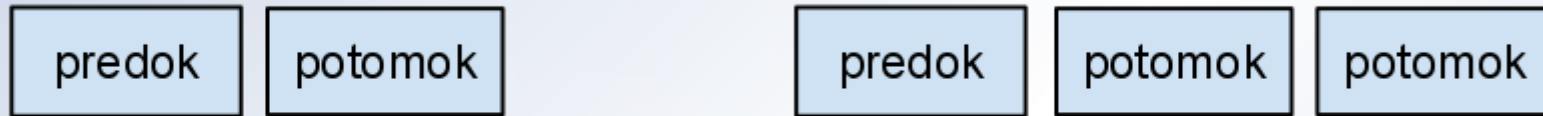
# Hierarchia - predošlé prístupy



# Hierarchia – náš nápad



## Hierarchia – náš nápad



## Výstup - klastre

session, google, field, app, name, type, Fabulous, Browse, tab  
session, firefox 5.0 ca, Fabulous Tabs, awesome tab, Awesome  
Bar, open tabs, panorama, Awesome Bar, tag

Javascript, Firefox, input, menu bar, web development platform, Vim  
text editor, Firefox extension, default configuration, error messages,  
syntax error, XPI

proxy servers, apache server

table, script, hash, cells

Repository, repositories, code, flash player, GitHub, Metal history,  
interactive overview, genres, Synchronize, sync button

## Overenie

- Príbuznosť dokumentov
- Rozširovanie dopytov
- Odporúčanie tagov

## Zhrnutie

- Vzťahy kľúčových slov z postupností
  - Viac dát, netreba tagovať
  - Časová informácia
- Detekcia sessions z taxonómie (iteratívne vylepšovanie)
- Využitie TF-IDF a asymetrickej podobnosti slov, stránok
- Identifikácia hierarchie z postupností
- Overenie riešenia