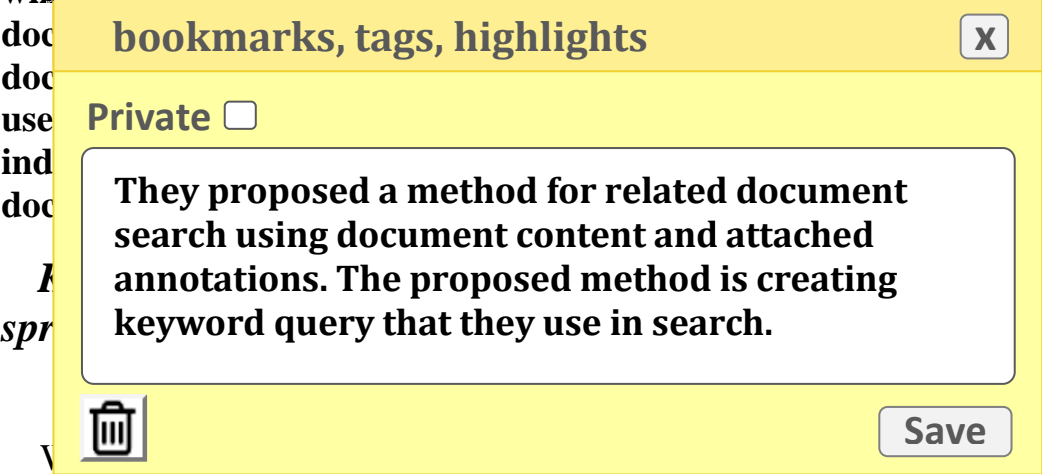# Related Document Search Using User Created Annotations

Jakub Ševcech

supervised by prof. Maria Bielikova
Slovak University of Technology in Bratislava
Faculty of Informatics and Information Technologies
Ilkovičova, 842 16 Bratislava, Slovakia
E-mail: sevo.jakub@gmail.com

*Abstract*—**We often use various services for creating bookmarks, tags, highlights and other types of annotations while surfing the Internet or when reading electronic documents... [obscured]**

bookmarks, tags, highlights

**Private** ☐

**They proposed a method for related document search using document content and attached annotations. The proposed method is creating keyword query that they use in search.**

Save

to insert various types of annotations into web pages and PDF documents displayed in browser. We have developed an extension for Firefox, which allows users highlight text bookmark pages and organizing them by tags, attach comments to text selections and attach notes to documents.
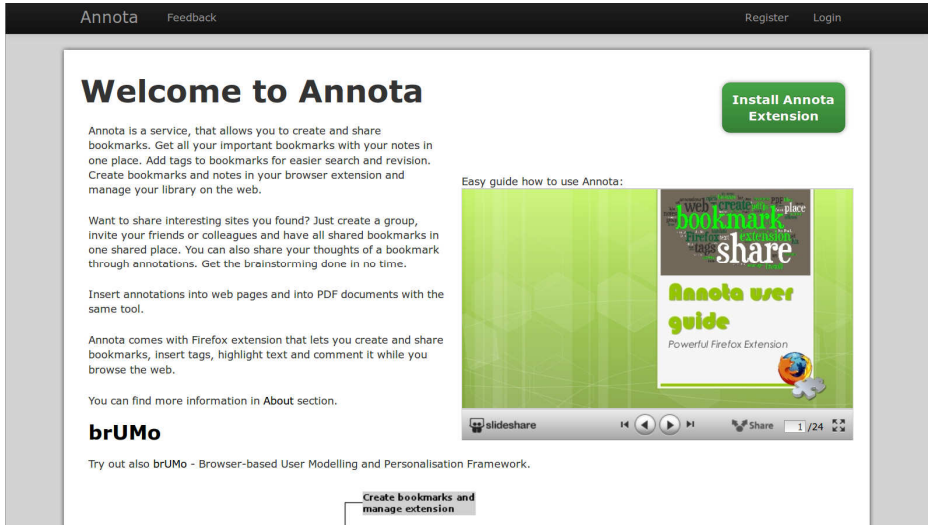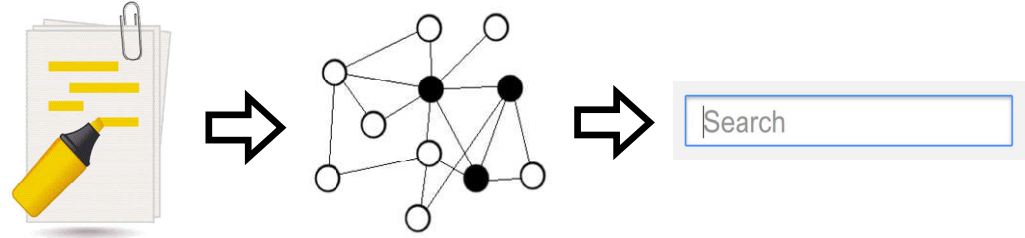


Figure 1 Annota, insert annotations into documents

## II. METHOD

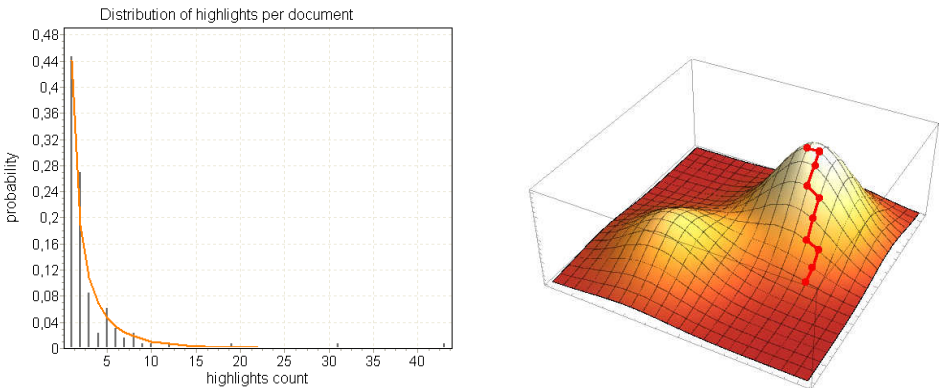Query construction for related document search:

1. Text to graph transformation
2. Spreading activation

Text words transformed to graph nodes. Edges created using neighborhood of words. Initial activation inserted to nodes with attached annotations. Spreading activation in created graph. Nodes with concentrated activation used as query words. We are using different weights for various annotation types. Important is number of iteration.



## III. SIMULATION

We have preformed a simulation to determine best weights for different types of annotations and to compare with commonly used TF-IDF based method. We analyzed properties of annotations from Annota.



We have generated annotations into documents from Wikipedia using parameters of annotations from Annota. Using proposed method, we created a query from document content and annotations. We have searched for similar documents. Relevant documents were those from the same category. We used hill climbing algorithm to optimize query construction parameters. Comparison of spreading activation based and TF-IDF based methods using simulation is in Table 1.

Table 1 Results of simulation with query construction using generated annotations

| Method | Precision |
|---|---|
| TF-IDF, no annotations | 21.32% |
| Proposed, no annotations | 21.96% |
| TF-IDF, generated annotations | 33.64% |
| Proposed, generated annotations | 37.07% |
| TF-IDF, whole fragment annotated | 43.20% |
| Proposed, whole fragment annotated | 53.34% |

Annotations are improving query construction process for different keyword extraction methods. Proposed method outperformed TF-IDF based method. Query construction process is independent from other documents but can be easily extended.

## IV. REFERENCES

[1] Faculty of Informatics and Information Technologies, Retrieved from http://www.fiit.stuba.sk/

[2] PeWe (Personalized Web) Group, Retrieved from http://pewe.fiit.stuba.sk/

[3] Annota, Retrieved from http://annota.fiit.stuba.sk/