

High-Recall Document Retrieval from Large-Scale Noisy Documents via Visual Analytics based on Targeted Topic Modeling

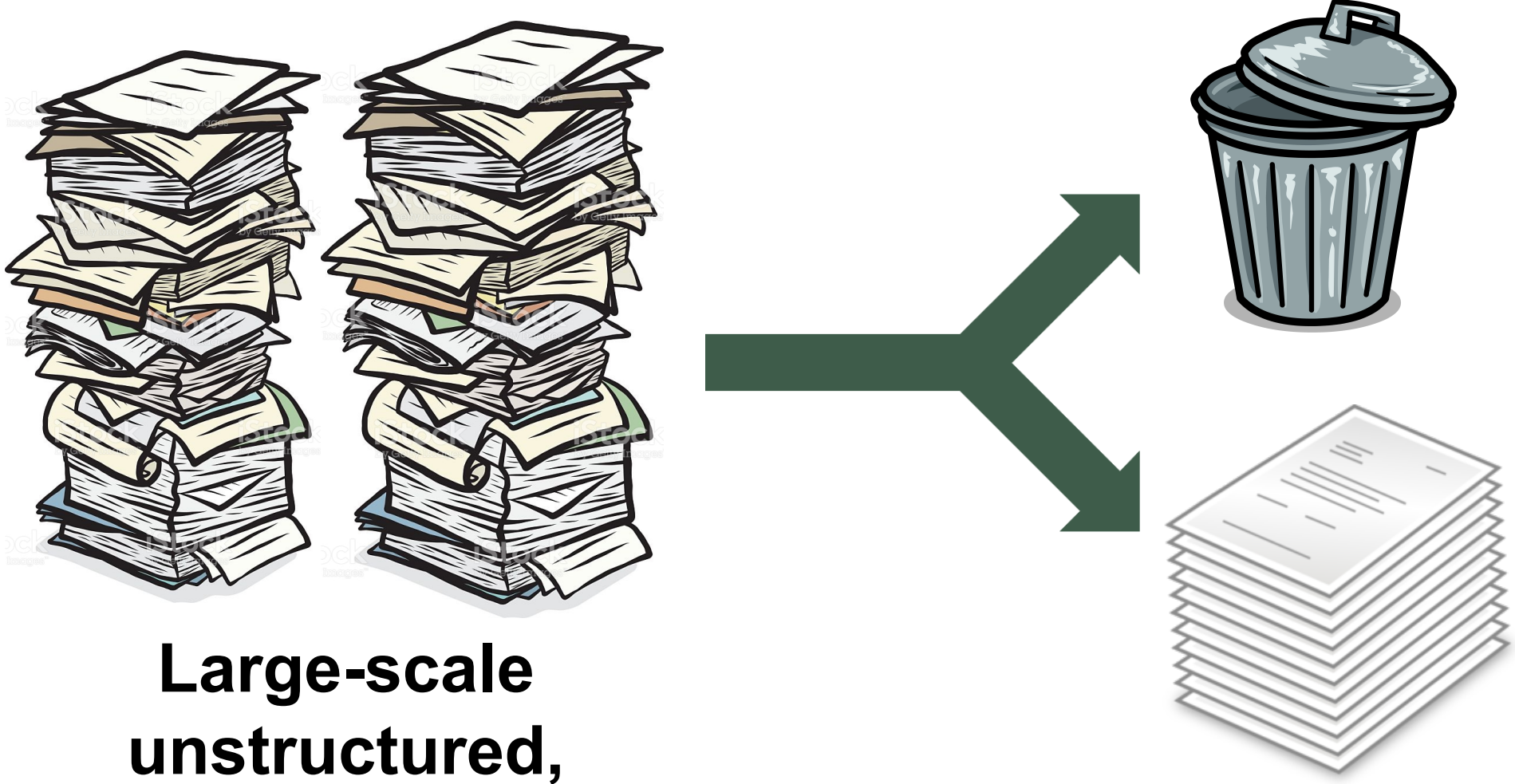
Hannah Kim¹, Jaegul Choo², Alex Endert¹, and Haesun Park¹

¹Georgia Tech, ²Korea University

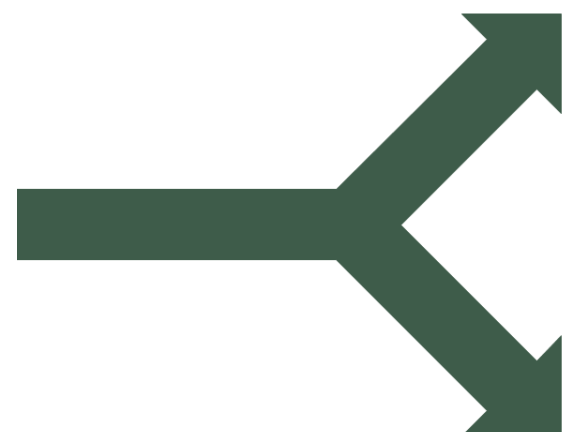
hannahkim@gatech.edu



Information Retrieval

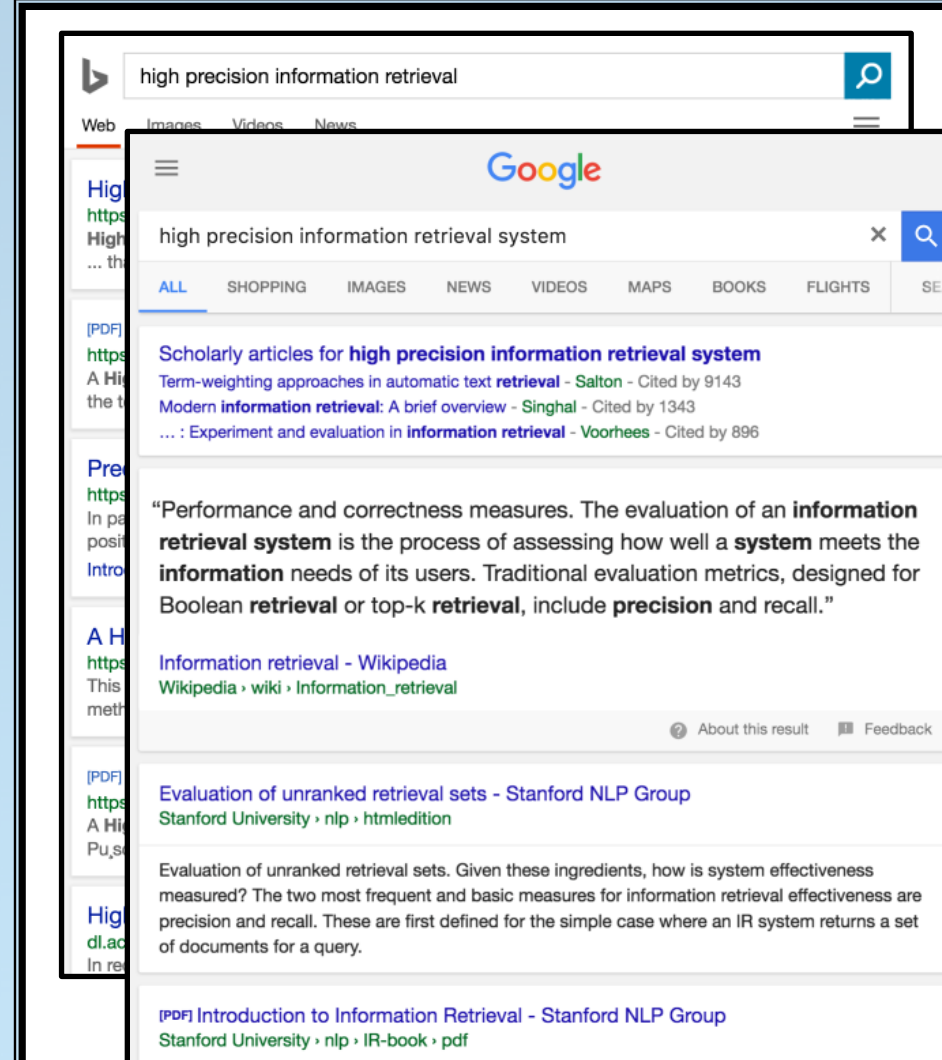


Large-scale unstructured, noisy text data



Relevant documents (about a particular event, subject, or product)

High Precision vs. High Recall



Traditional retrieval systems:

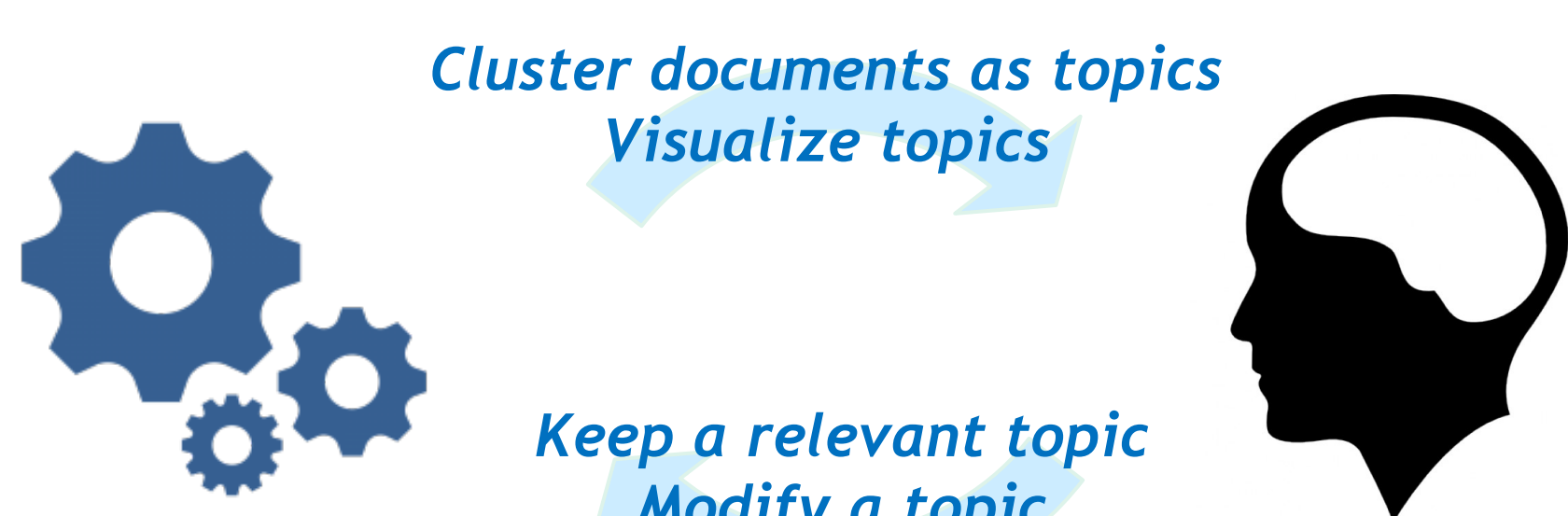
- Focus only on high precision
- Retrieve **a number of most relevant** documents
- Ex) Google, Bing, Twitter, PubMed

Our system:

- Focuses on high precision and high recall
- Retrieves **ALL relevant** documents
- Is suited when missing any relevant item is critical
- Ex) marketing, social media, legal cases, medical cases, literature review, etc

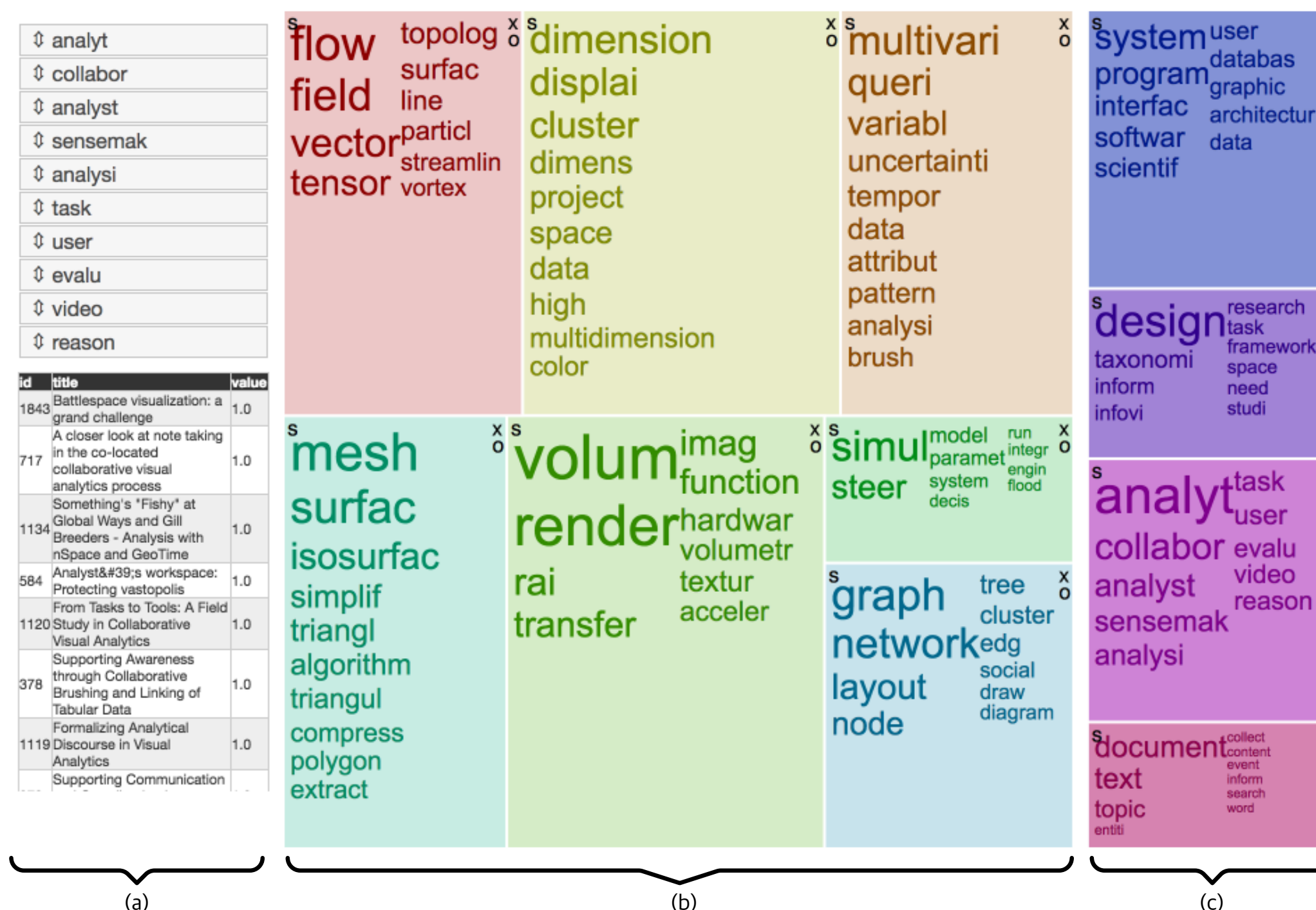
Interactive Visual Retrieval System

Workflow



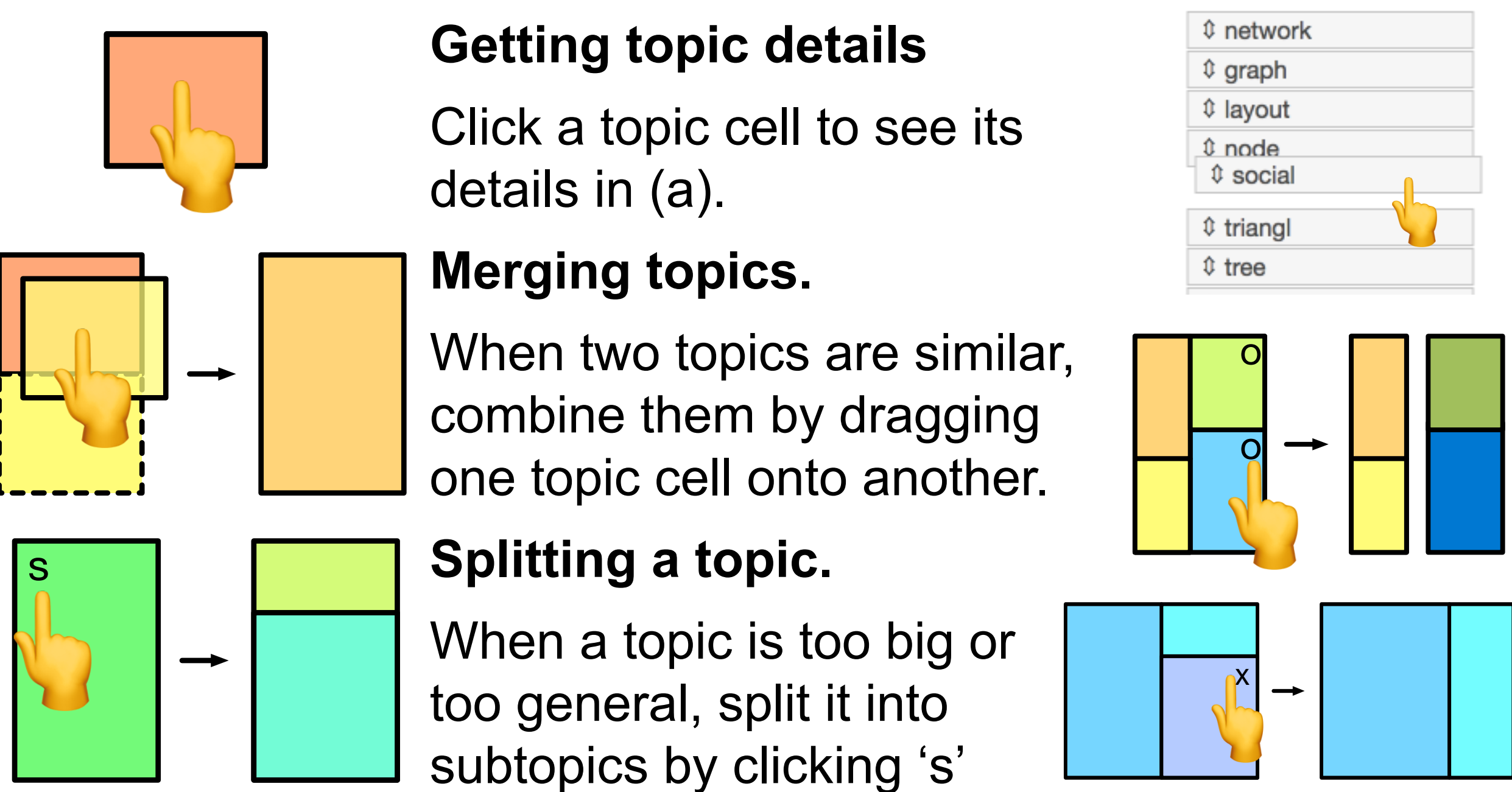
1. Documents are clustered w.r.t. their topics.
2. Topics are visualized and used as exploration units.
3. Users can inspect a topic and keep/modify/remove it iteratively.

System Design

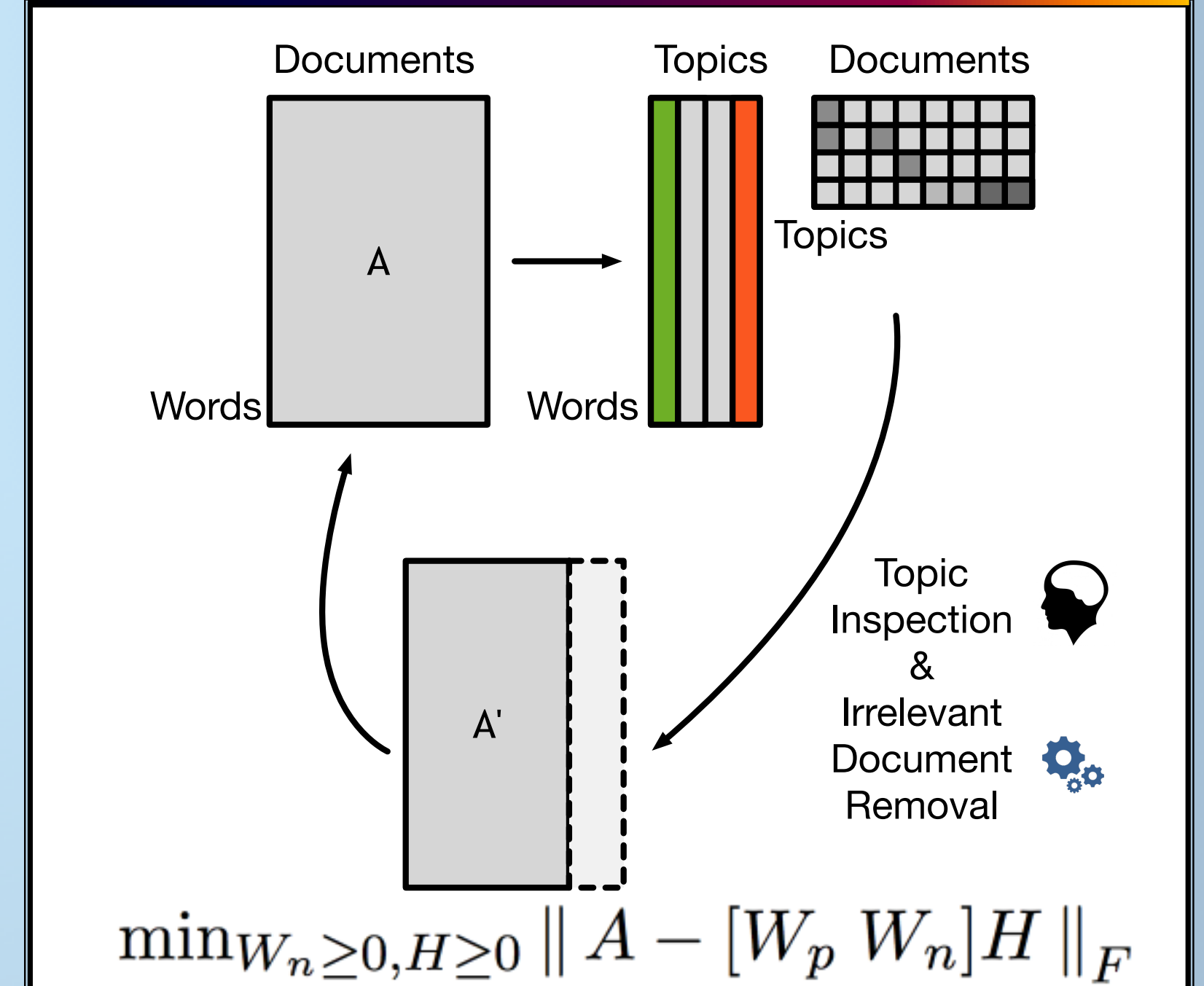


- (a) The topic detail panel with an interactive list of keywords and a document table
 - (b) Main topic treemap visualization for interactive topic exploration where each cell represents a topic
 - (c) Confirmed topic treemap visualization which shows relevant topics that are confirmed by users
- In (b) and (c), semantically similar topics are placed closer.

Supported Interaction



Targeted Topic Modeling



Usage Scenario with IEEE VIS publication dataset

