# Project 4: Video Search
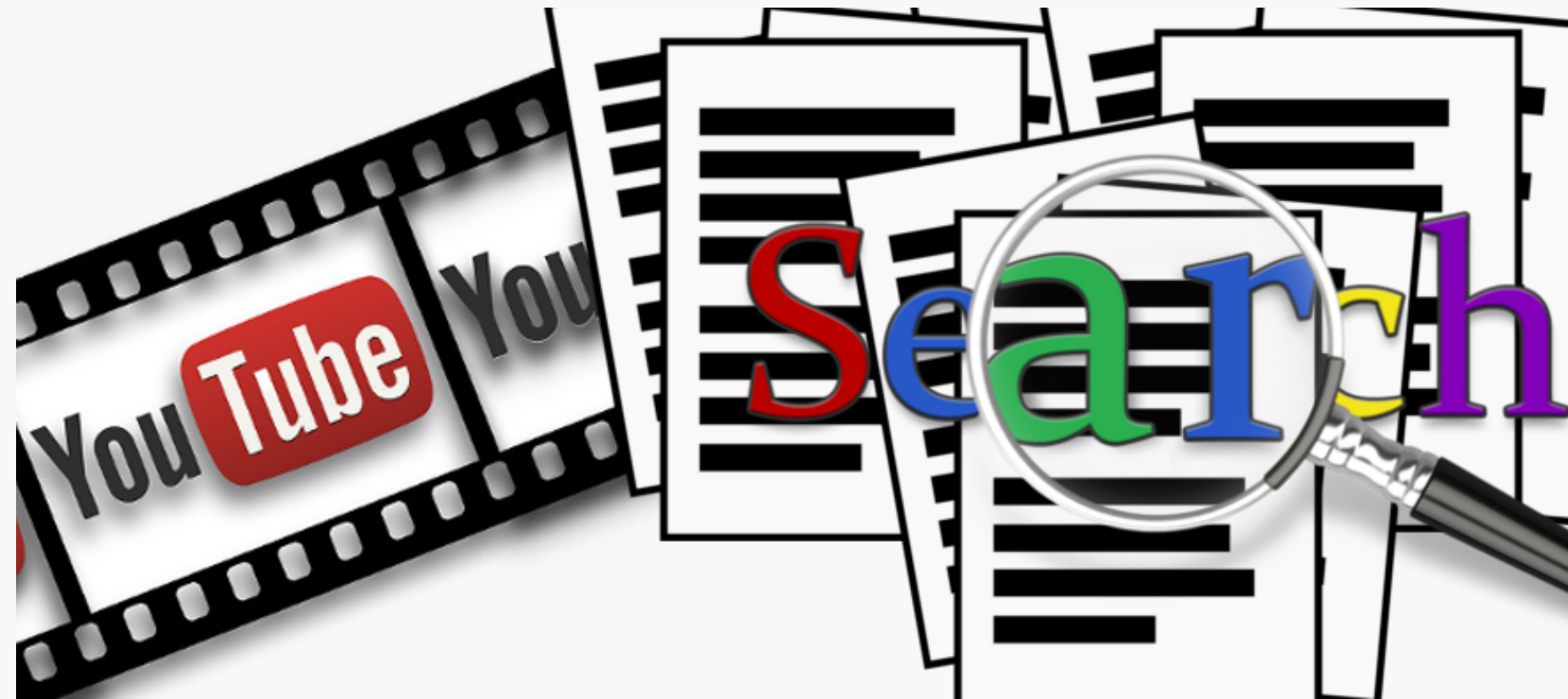
## Problem description

► Video search engine
► YouTube does not use visual content
► Input: visual query as text (innovative)
► Output: List of links with specified time



## Method

► First step: Mapping video to text.
  **Metadata** Collect data around the videos e.g. title, description, likes, etc.
  **Captions** Neural description of video frames - description of visual content.
  **Expansion** Expanding human-like neural description with WordNet hypernyms and hyponyms.
► Second step: Getting a "good" search engine on those special text documents
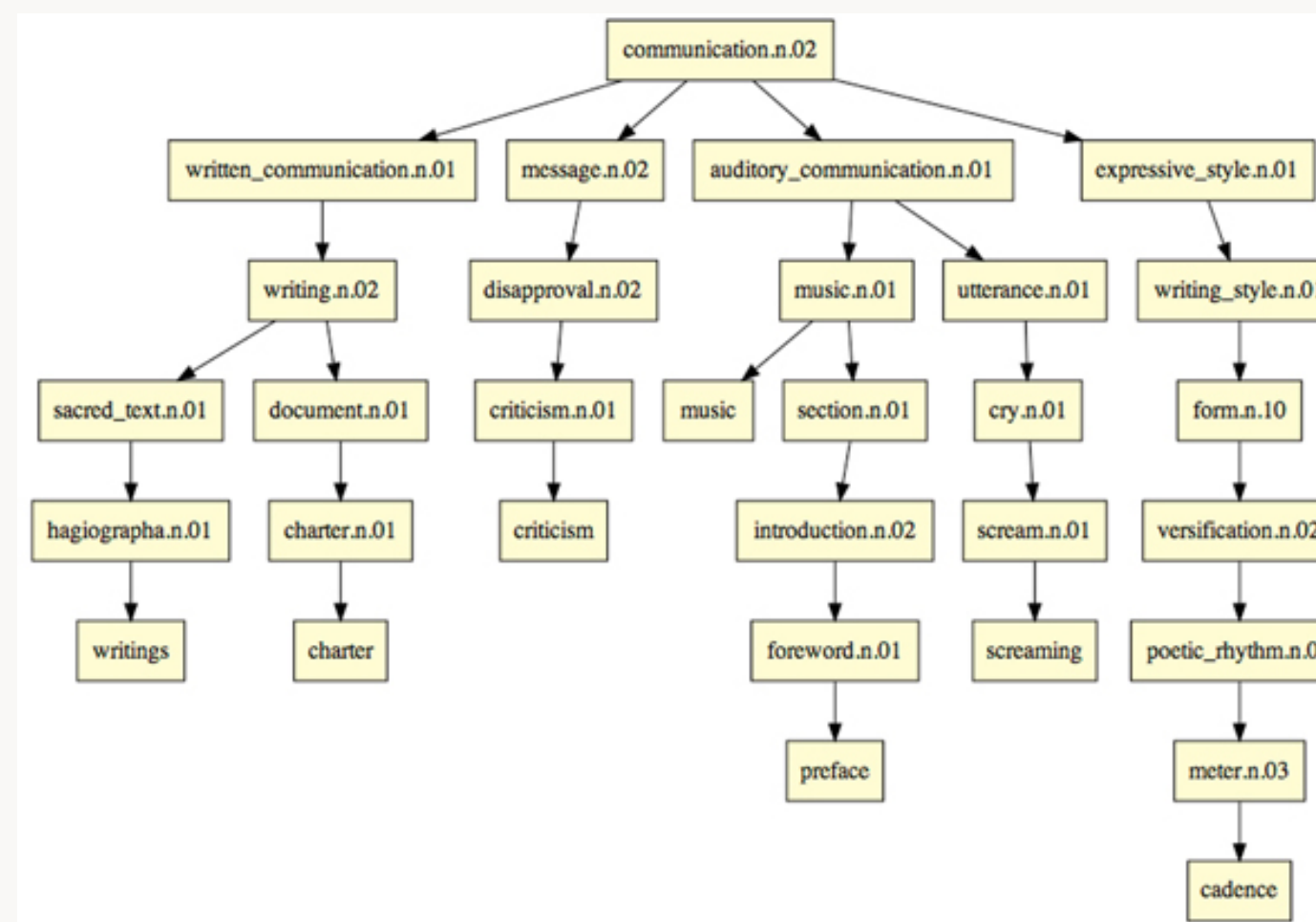  **Requirements:** high precision and recall, fast and ergonomic

## Data

► Generated Video-captions
► Meta-data
  • Title, Description
  • Views, Favorites, Votes, Comments

## Visual content to text mapping

► Generated human-like neural descriptions.



► Expanded descriptions.



## Experiments

Evaluating search engine with three queries:
► Large airplane
► Playing with a dog
► Man on a bench

For each query:
► Precision and recall at 1, 3, 5, 10, 20, 30
► Area under the precision-recall graph

## Results

Choice of:
– Language model for text mapping

| Search | CNN_S | ILSVRC_16 |
|--------|-------|-----------|
| Q1 | 0.150 | 0.404 |
| Q2 | 0.500 | 0.092 |
| Q3 | 0.185 | 0.45 |
| Average | 0.278 | 0.315 |

– Retrieval model → ranked retrieval using tf-idf and popularity.

| Search | Boolean | Ranked |
|--------|---------|--------|
| Q1 | 0.178 | 0.404 |
| Q2 | 0.23 | 0.092 |
| Q3 | 0.323 | 0.45 |
| Average | 0.246 | 0.315 |

Table: Size of the area under the curve according to the retrieval model and the query considered
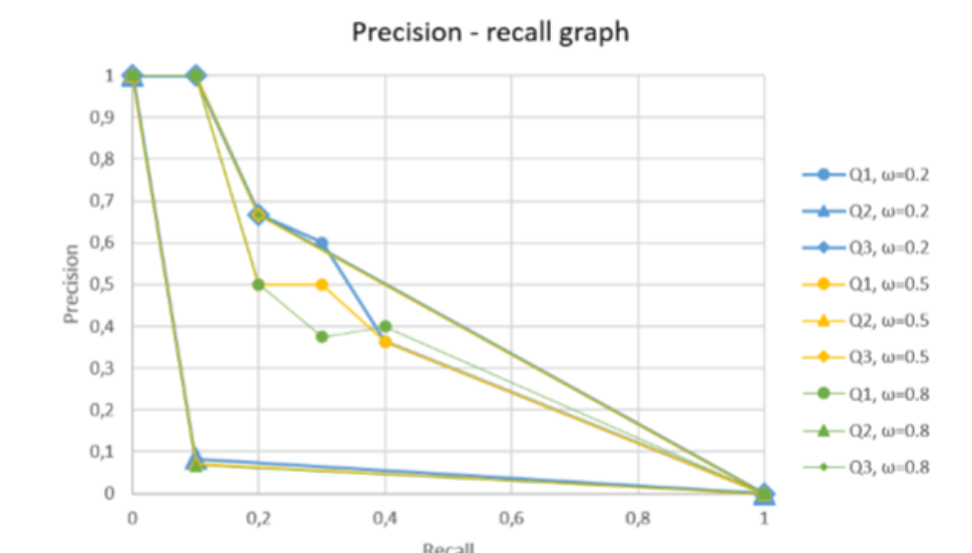


Figure: Precision - recall graph using different weights for the popularity compared to the tf-idf.

## Conclusions

We provide users a video search engine available online:
► Fast and ergonomic.
► Precision: generally high for top 3 results, then very low

**Marc Beillevaire, Mateusz Buda, Ted Cassirer, Laura Jacquemod, Christoph Kaiser**
marcbei@kth.se, buda@kth.se, ted.cassirer@gmail.com, lauraja@kth.se, ckai@kth.se
KTH Royal Institute of Technology, Department of Computer Science

KTH VETENSKAP OCH KONST