

# Remarkable Scientists: ExpertSearch

## Project Status Report

For this project, we worked on three different tasks:

- Topic Mining of the ExpertSearch professor bios.
- Improving Named Entity Recognition & Extraction of professor names from the bios.
- Improving recognition and extraction of professor e-mail addresses from the bios.

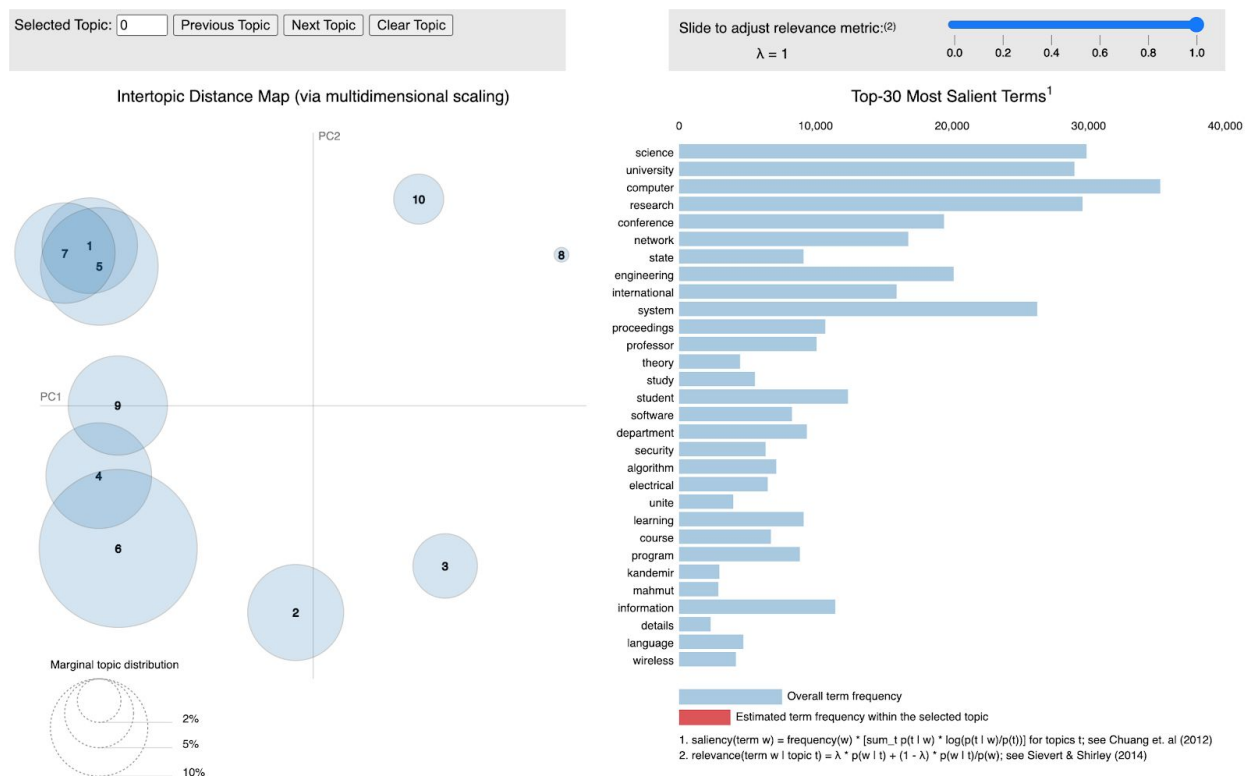
Below is the status report for each.

## Topic Mining of Professor Bios

We were able to combine [spaCy](#), [NLTK](#), and [Genism](#) to build an LDA topic model of the professor's bios, using 10 topics. We then used [pyLDavis](#) to visualize the model, and [word\\_cloud](#) to build a word cloud of the 25 highest-weighted terms in each topic. Finally, we wired the word cloud of the highest-ranked topic into the ExpertSearch search results.

### LDA Visualization

The following was created using pyLDavis, which is visualizing a 10-topic LDA model:



The following word clouds show the top 25 words in each of the 10 LDA topics:

A word cloud visualization of the SIGGRAPH 2012 program. The words are arranged in a dense, overlapping manner, with 'graphics' being the largest and most central word. Other prominent words include 'paper', 'image', 'computer', 'method', 'light', 'video', 'material', 'simulation', 'visualization', 'rendering', 'shape', 'field', 'symposium', 'base', 'imaging', 'geometry', 'modeling', 'motion', 'surface', 'present', 'optical', 'transactions', 'using', and 'visualization'.

Topic 9 Word Cloud



Topic 10 Word Cloud



## ExpertSearch Search Results

Finally, we were able to wire in the primary topic associated with the bio into the search results:

Expert Search

streaming

Q

⚙

Hulya Seferoglu

Electrical & Computer Engineering, University of Illinois at Chicago

Illinois, United States

...a. markopoulou, microcast: cooperative video **streaming** using cellular and local...a. markopoulou, microcast: cooperative video **streaming** on smartphones, in proc....le, a. markopoulou, cooperative video **streaming** on smartphones, in proc.

Steve Cole

Computer Science, Washington University in St. Louis

Missouri, United States

...include parallel computing and accelerating **streaming** applications on gpus.: ....include parallel computing and accelerating **streaming** applications on gpus..

## Next Steps

We intend to experiment with different topic counts in the LDA model to determine the best results.

We are also interested in making the ExpertSearch system browsable by allowing the user to click on a topic word cloud and see the bios ranked in descending order of relevance to that topic.

## Challenges

Metapy does not support the ability to return a list of documents based on its metadata fields, so it is unclear if we would be able to build this additional functionality in the time that we have.

## Named Entity Extraction of Professor Names

We have conducted several experiments on NER(Named Entity Extraction) solutions in the market. The best open source solution would be Stanford NER Tagger, which is also provided along with ExpertSearch project for version 2018. The provided tagger and compiled result are incomplete as it eliminates some name information during the tagging phase, part of missing information is critical to extract the main context from the text. We have done two improvements for solving this problem, firstly, we update NER tagger to the newest version which is 2020-11 online, and test on 3-class model and 4-class model separately. Secondly, we implemented a filtered mechanism from tagging results, filtering all names with more than one word, capturing all tokens tagged as 'Person' in a row, and building a two-layer tagging system(3-class model and 4-class model) to cross-validation. We use the 4-class model as the main model to capture 'PERSON' entities, and check each produced token whether it states in the 3-class model. Now the system we created is able to capture most names from context files, filtering names are irrelevant partially. For example, for the first compiled bio text, we are able to retrieve '**Tarek F. Abdelzaher Professor**' as full name rather than '**Tarek**', second one '**Sarita V. Adve**' as correct name rather than '**Sarita V. Adve Richard T. Cheng**' in the provided file

## Next Steps

We will continue working on a 2-layer tagging system to cross validation results between two models. Any results show or partial show on both models state they are 'Person/Name' entities with high possibility. Additionally, we will continue on filtering results, with introducing a counting/score system. Any name tokens appearing in high frequency should have a high score indicating its the main context of text.

## Challenges

The biggest challenge we are facing is lots of famous names and confusing information shown in the bio files. For example, 'Ann Arbor' can be a name or location. The tagging system is not able to distinguish between them. Also 'Kennedy' as a famous name, sometimes can occur multiple times in the bio page, the system even with the scoring system is not able to know 'Kennedy' is not the main context. We will introduce another non-standard model can tag 'locations' entity, then do another cross-validation between our model and the third party model result, it should lead us to a potential solution finding the main person of the context, and better filter out non-relevant names/location.

# Extraction of Professor E-Mail

We have worked on improving the existing regex based extraction of email ids from faculty bios. The original code was extracting the usual format of email ids (user@illinois.edu) and also some false positives such as cs@uiuc. Our new code is able to extract many of the alternative email id formats, for example: "**yang.r.yang at yale.edu**" "**denisew (at) uw.edu**" "**moli96 at uw.edu**". It is also able to capture email ids that were originally missing from the first output and remove some false positives. Below is a comparison of the outputs from existing code shown on the left and output from new code shown on the right.

```
amitha@ANSI-P10K70647: /mnt/c/Users/asandur2/ExpertSearch/data
4 lines: sadve@illinois.edu
alawini@illinois.edu
amato@tam.u.edu
amrgrave@illinois.edu
bphalley@illinois.edu
batesa@illinois.edu
matttox@illinois.edu

caesar@illinois.edu
rhc@illinois.edu
challen@illinois.edu

chekur@illinois.edu
rcumin@illinois.edu
davis6@illinois.edu
melkebir@illinois.edu

gcevans@illinois.edu
waf@illinois.edu

mfleck@illinois.edu
uiuc@cs
miforbes@illinois.edu
daf@illinois.edu
wfu@illinois.edu
phg@illinois.edu

egunter@illinois.edu
cs@uiuc

sariel@illinois.edu
aharris@illinois.edu
jch@illinois.edu

glherman@illinois.edu
juliahm@illinois.edu
dhoiem@illinois.edu
shj@illinois.edu

kale@illinois.edu
kkrahal@illinois.edu
cs@illinois
kirlik@illinois.edu
andreas@illinois.edu
sam@illinois.edu
wtkramer@illinois.edu

4 lines: sadve@illinois.edu
alawini@illinois.edu
amato@tam.u.edu
amrgrave@illinois.edu
bphalley@illinois.edu
batesa@illinois.edu
matttox@illinois.edu
mcacamo "at" illinois.edu
caesar@illinois.edu
rhc@illinois.edu
challen@illinois.edu
tmc "at" illinois "dot" edu
kcchang (at) illinois (dot) edu
chekur@illinois.edu
rcumin@illinois.edu
davis6@illinois.edu
melkebir@illinois.edu

gcevans@illinois.edu
waf@illinois.edu

mfleck@illinois.edu

miforbes@illinois.edu
daf@illinois.edu
wfu@illinois.edu
phg@illinois.edu
wtropp at illinois.edu

egunter@illinois.edu

sariel@illinois.edu
aharris@illinois.edu
jch@illinois.edu

glherman@illinois.edu
juliahm@illinois.edu
dhoiem@illinois.edu
shj@illinois.edu

kale@illinois.edu
kkrahal@illinois.edu
dakshita@cs.ucla.edu
kirlik@illinois.edu
andreas@illinois.edu
sam@illinois.edu
wtkramer@illinois.edu
```

As observed from the figure above, the new code is able to catch the "user at illinois dot edu" types of format email ids.

We have covered cases for the following email formats:

- user at illinois.edu
- user at illinois dot edu
- user "at" illinois "dot" edu or user "at" illinois.edu
- user (at) illinois (dot) edu
- We removed some erroneous outputs like website urls/sentences/repeated special characters

## Next Steps

We plan to convert the email ids in "user at illinois dot edu" format into the normal "user@illinois.edu" format as one of our next steps.

Although the new regex based code is a significant improvement over the existing code, we feel that this method is laborious. And even if we spend a lot more time trying to cover all the exhaustive number of cases, the gain in improvement for the tool as a whole wouldn't be much, as the number of cases are a lot. Hence we plan to look into machine learning based approaches to extract email ids.

## Challenges

The challenge here is to convert the various different email formats into the usual email id format. Although we may consider keeping it as is if we discover more new email formats along the way and it becomes too heterogeneous, because we have also seen a few websites having email ids in this alternative format.