

Supervised Keyphrase Extraction as Positive Unlabeled Learning

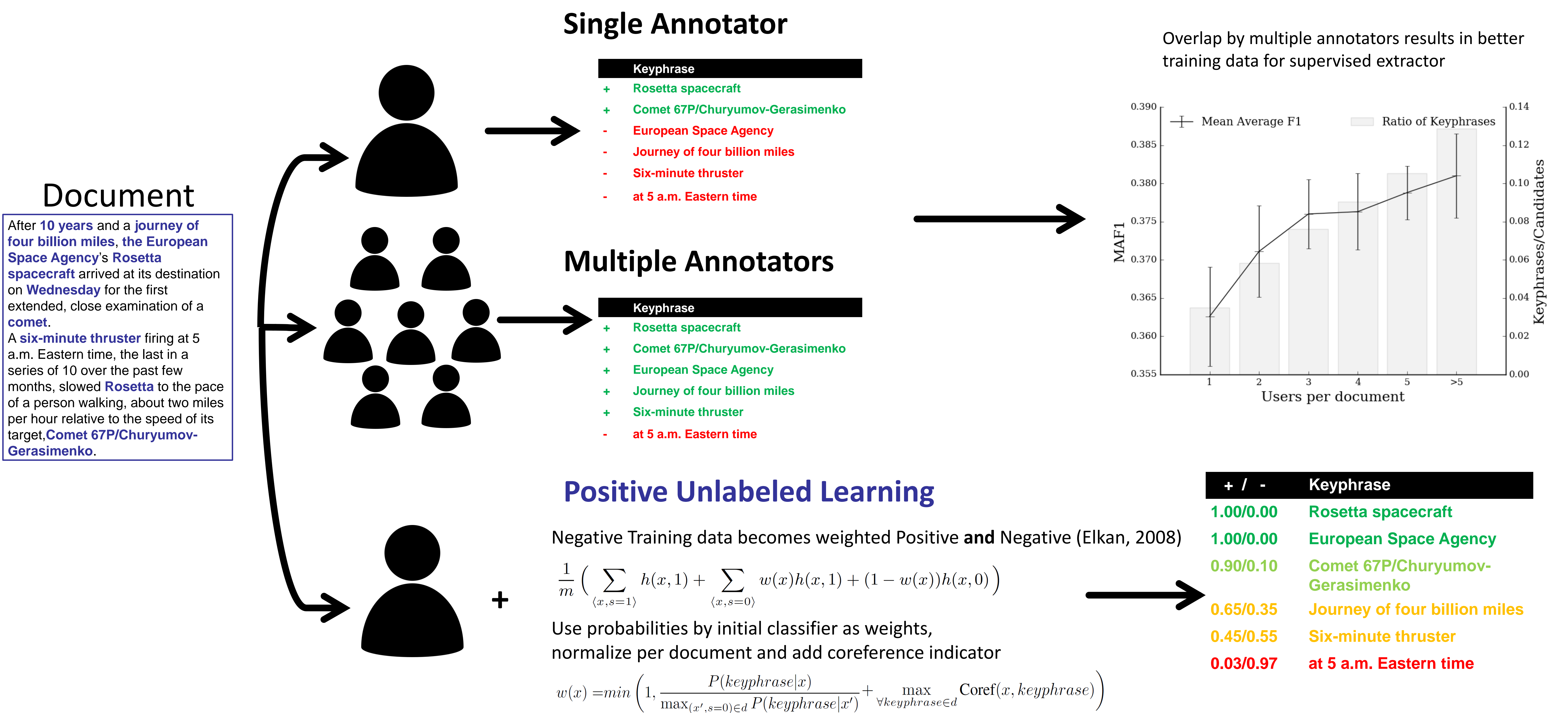
Lucas Sterckx*, Cornelia Caragea†, Thomas Demeester*, Chris Develder*
 *Ghent University – imec, †University of North Texas

Training Data for Supervised Keyphrase Extraction

- Supervised keyphrase extraction = binary classification of keyphrase candidates
 - State-of-the-art but requires **training data**
- Problems with keyphrase annotations and resulting training data
 - Noisy : keyphrases are highly subjective
 - Unbalanced : negatives outweigh positives many times

We create large test collections of articles with many different opinions per document, evaluate the effect on extraction performance, and present a procedure for supervised keyphrase extraction with noisy labels.

Keyphrase Extraction as Positive Unlabeled Learning



Experiments										
Method	Online News		Lifestyle Magazines		WWW		KDD		Inspec	
	MAF ₁	P@5	MAF ₁	P@5	MAF ₁	P@5	MAF ₁	P@5	MAF ₁	P@5
Single Annotator	.364	.416	.294	.315	.230	.189	.266	.200	.397	.432
Multiple Annotators	<u>.381</u>	<u>.426</u>	.303	<u>.327</u>	/	/	/	/	/	/
Self Training	.366	.417	.301	.317	.236	.190	.269	.196	.401	.434
Reweighting	.364	.417	.297	.313	.238	.189	.275	.201	401	.429
Reweighting + Norm. + Coref	.374	.419	<u>.305</u>	.322	.245	.194	.275	.200	.402	.434

Conclusions

- Keyphrase datasets with multiple annotations per document
- Supervised keyphrase extraction as **Positive Unlabeled Learning**
 - Treat non-selected phrases as **unlabeled instead of negative**
 - Reweigh keyphrases using noisy classifier prediction, **normalization and coreference**
- Future work : **Evaluation** measures for keyphrases