

Research Paper Recommendation System

Using Document Level Embedding

Akshay Gupta | Anant Dev | Muhammed Shariq Nawaz | Gaurang Raje

Flow of our Talk

- Motivation
- Main Idea
- Existing Research
- Our Method
- Our Progress
- Our Hypothesis



Motivation

- Exponential rise in research papers published
- Accessing related documents is monotonous
- Simplify the literature review process
- Get relevant citations of related papers



Process Steps

- Language Transformer/ Sentence Transformer
- Document Embedding
- Output - The N-dimensional embeddings
- Cosine Similarity
- Final output - Top-K results

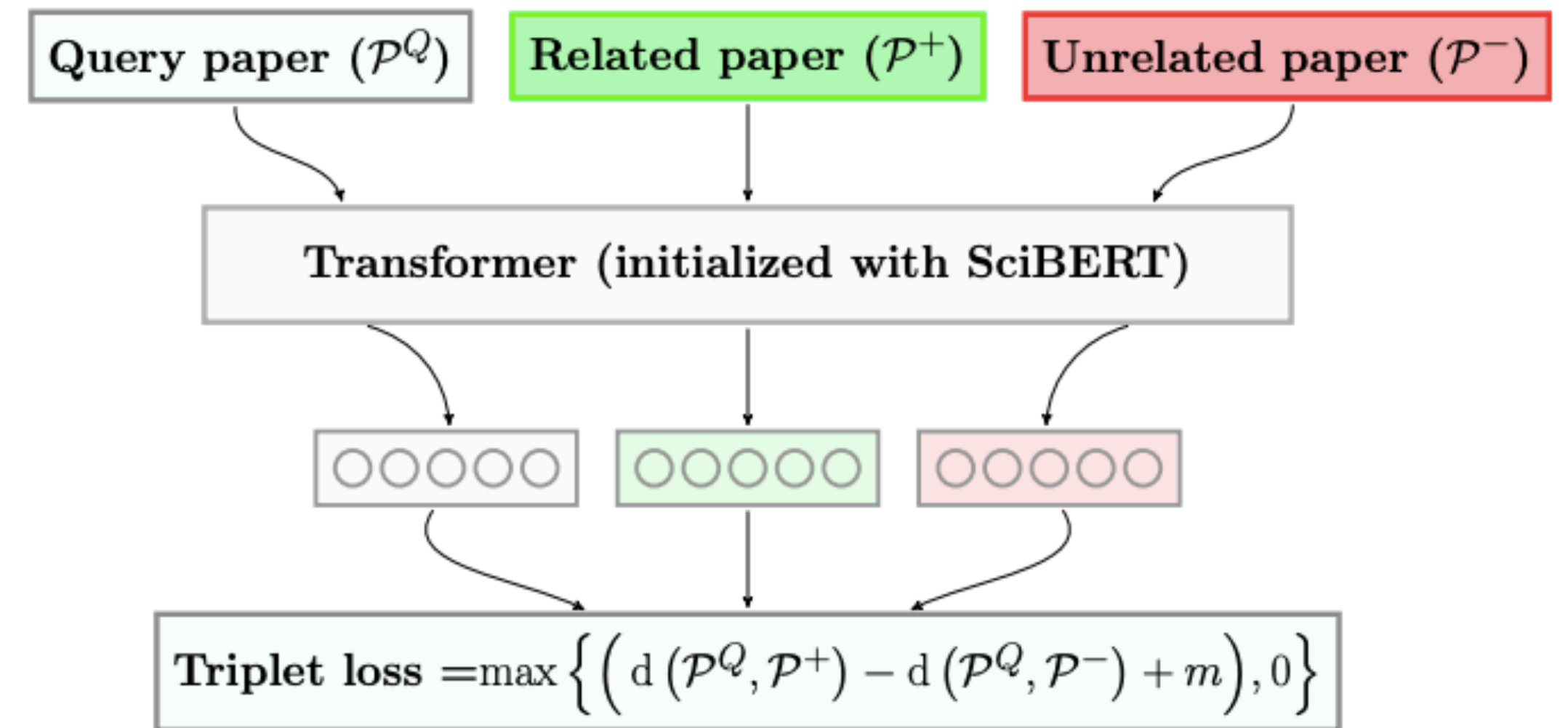


Figure 1: Overview of SPECTER.

Existing Research

Work by Haruna et al.

- Based on contextual metadata
- Heuristic algorithm to score related papers
- Normalised score returns top-N recommendations
- Evaluation metrics:

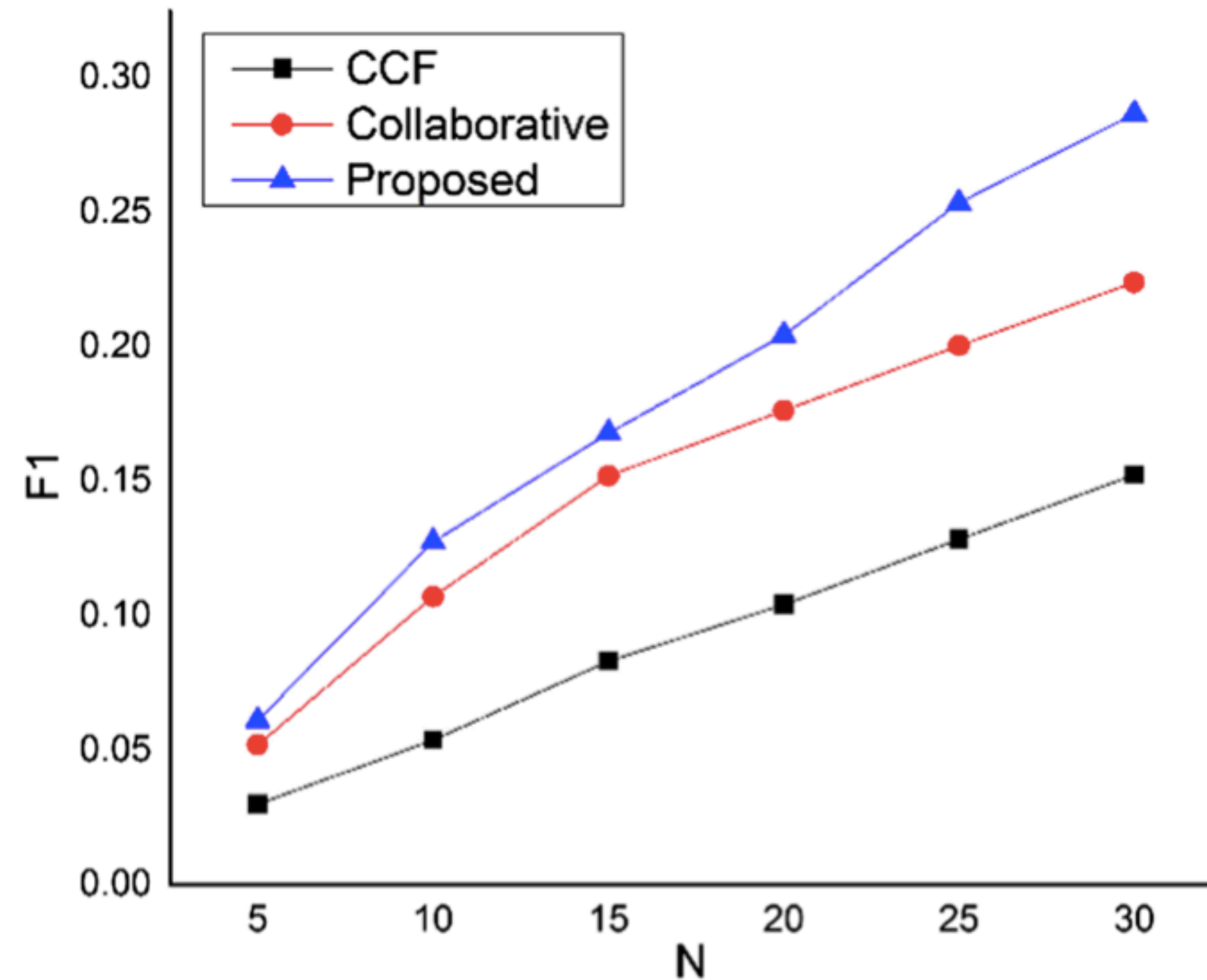
$$precision = \frac{\sum (relevant_papers) \cap \sum (retrieved_papers)}{\sum (retrieved_papers)}$$

$$recall = \frac{\sum (relevant_papers) \cap \sum (retrieved_papers)}{\sum (retrieved_papers)}$$

$$F1 = \frac{2 \times precision \times recall}{precision + recall}$$

Work by Haruna et al.

- Results:



- Drawbacks:

- Divergent Techniques papers
- Relevance/preference score

Arxiv Neural Search Repository

- Uses categories and abstracts for similarity search
- Sources above elements from Kaggle Arxiv Metadata dataset
- FAISS for indexing
- DistillBERT for generation of query embeddings

Arxiv Neural Search Repository

localhost:8501

Apps SID CRUD Artificial Intelligenc... Improve HCP Email... Storyboard your ide... 14 Most Important... Average Click-Thro... 8 Ways Machine Le... Epocrates Versus M... Other bookmarks Reading list

RUNNING... Stop

Arxiv AI/ML Vector search with Sentence Transformers and Faiss

Search box

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder-decoder configuration. The best performing models also connect the encoder and decoder through an attention mechanism. We

Filters

Number of search results

10

10

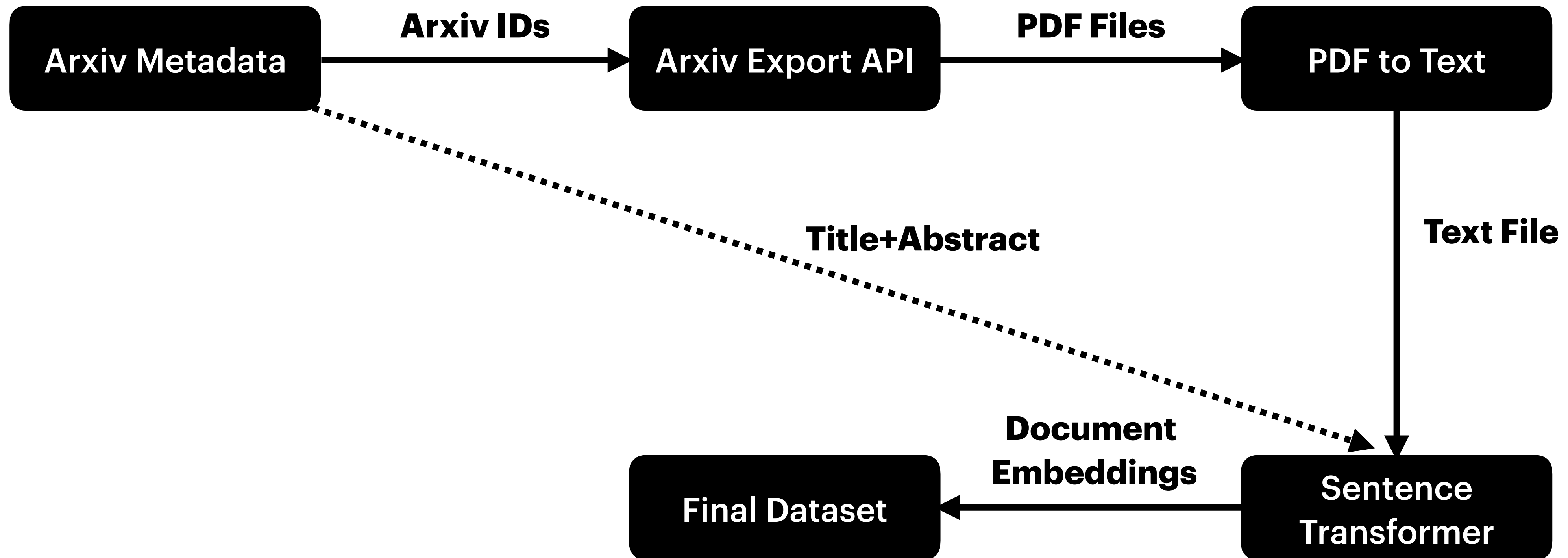
50

Select the AI/ML Categories

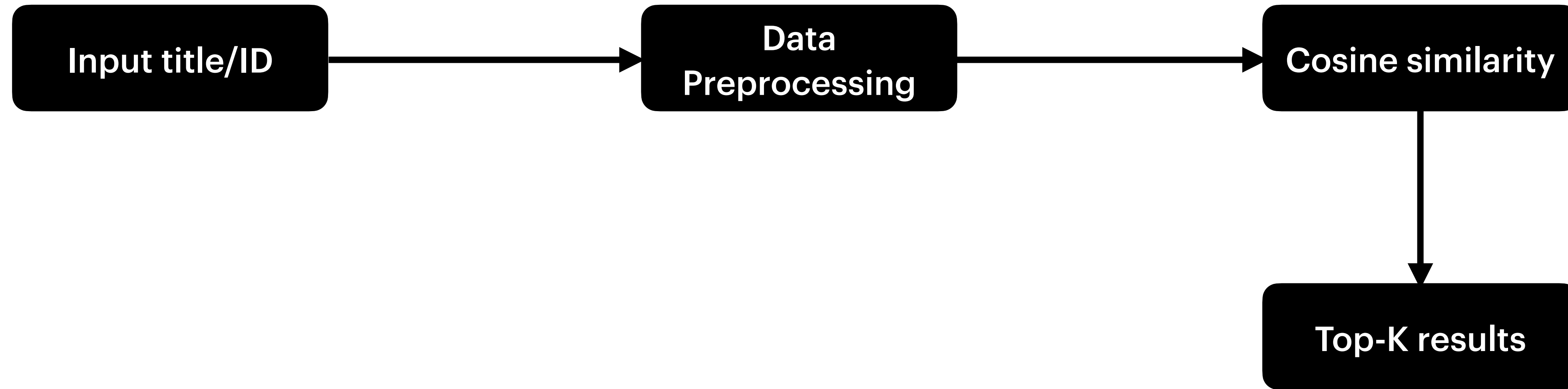
Choose an option

Our Method

Dataset Preprocessing



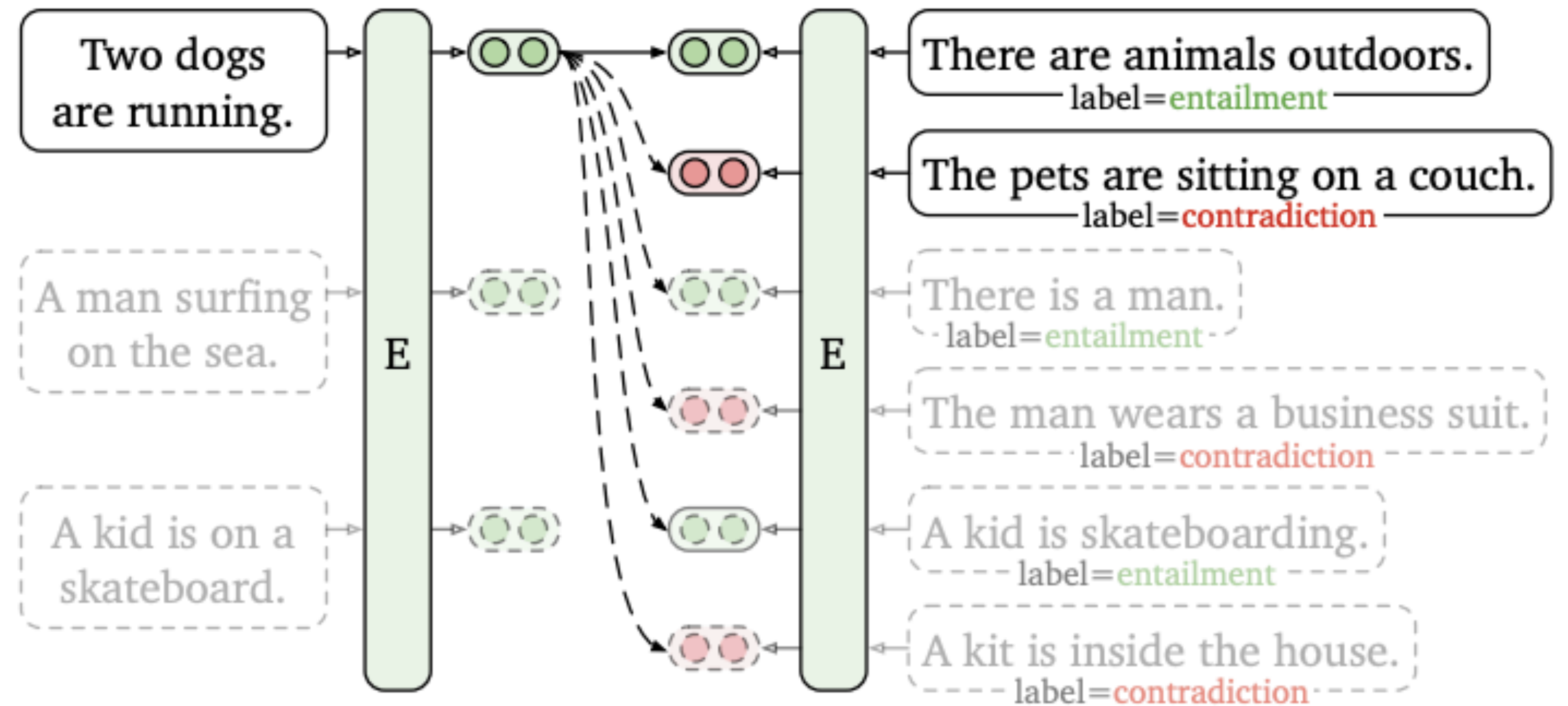
Testing on New Inputs



Pre-trained Models

- Contrastive learning:
 - SimCSE
 - DeCLUTTR
- SPECTER

(b) Supervised SimCSE



Our Progress so far...

[11]:

```
# This paper was a EMNLP 2020 Honourable Mention Papers
search_papers(title='If beam search is the answer, what was the question?',
              abstract='Quite surprisingly, exact maximum a posteriori (MAP) decoding of neural language generato
```

Batches: 100%  1/1 [00:00<00:00, 20.50it/s]

	title	authors	arxiv_id	score
0	Infinite Viterbi alignments in the two state hidden Markov models	J. Lember, A. Koloydenko	0711.0928	0.819104
1	On the Vocabulary of Grammar-Based Codes and the Logical Consistency of Texts	{\L}ukasz D{\k(e)}bowski	0810.3125	0.814937
2	CLAIRLIB Documentation v1.03	Dragomir Radev, Mark Hodges, Anthony Fader, Mark Joseph, Joshua Gerrish, Mark Schaller, Jonathan dePeri, Bryan Gibson	0712.3298	0.806504
3	Questions & Answers for TEI Newcomers	Laurent Romary (LORIA)	0812.3563	0.803840
4	The emerging field of language dynamics	S. Wichmann	0801.1415	0.799843
5	In memoriam Maurice Gross	Eric Laporte (IGM-LabInfo)	0711.3452	0.799252
6	A constructive proof of the existence of Viterbi processes	J. Lember, A. Koloydenko	0804.2138	0.798086
7	Morphic and Automatic Words: Maximal Blocks and Diophantine Approximation	Yann Bugeaud, Dalia Krieger, Jeffrey Shallit	0808.2544	0.795205
8	Three Lectures on Automatic Structures	Bakhadyr Khoussainov, Mia Minnes	0809.3430	0.789193
9	UNL-French deconversion as transfer & generation from an interlingua with possible quality enhancement through offline human interaction	Gilles s'erasset (IMAG, Clips - Imag, Lig), Christian Boitet (IMAG, Clips - Imag, Lig)	0811.0579	0.785540

Our Progress so far...

[10]:

```
this paper was the EMNLP 2020 Best Paper
ch_papers(title='Digital Voicing of Silent Speech',
          abstract='In this paper, we consider the task of digitally voicing silent speech, where silently mouthe
```

Batches: 100%  1/1 [00:00<00:00, 18.10it/s]

	title	authors	arxiv_id	score
0	EasyVoice: Integrating voice synthesis with Skype	Paulo A. Condado and Fernando G. Lobo	0706.3132	0.818029
1	Amélioration des Performances des Systèmes Automatiques de Reconnaissance de la Parole pour la Parole Non Native	Ghazi Bouselmi (INRIA Lorraine - LORIA), Dominique Fohr (INRIA Lorraine - LORIA), Irina Illina (INRIA Lorraine - LORIA), Jean-Paul Haton (INRIA Lorraine - LORIA)	0711.1038	0.807498
2	Adjusted Viterbi training for hidden Markov models	J. Lember, A. Koloydenko	0709.2317	0.770014
3	Exploiting Nonlinear Recurrence and Fractal Scaling Properties for Voice Disorder Detection	Max A Little, Patrick E McSharry, Stephen J Roberts, Declan AE Costello and Irene M Moroz	0707.0086	0.766409
4	Probabilistic SVM/GMM Classifier for Speaker-Independent Vowel Recognition in Continues Speech	Mohammad Nazari, Abolghasem Sayadiyan, SeyedMajid Valiollahzadeh	0812.2411	0.761897
5	Acoustic Features and Perceptive Cues of Songs and Dialogues in Whistled Speech: Convergences with Sung Speech	Julien Meyer (LAB-Upc)	0711.3704	0.760309
6	Combined Acoustic and Pronunciation Modelling for Non-Native Speech Recognition	Ghazi Bouselmi (INRIA Lorraine - LORIA), Dominique Fohr (INRIA Lorraine - LORIA), Irina Illina (INRIA Lorraine - LORIA)	0711.0811	0.756685
7	Phoneme recognition in TIMIT with BLSTM-CTC	Santiago Fernandez, Alex Graves, Juergen Schmidhuber	0804.3269	0.753250
8	A biomechanical model of the face including muscles for the prediction of deformations during speech production	Julie Groleau (TIMC), Matthieu Chabanas (TIMC), Christophe Marecaux (TIMC), Natacha Payrard, Brice Segaud, Michel Rochette, Pascal Perrier (ICP), Yohan Payan (TIMC)	0803.3924	0.750757
9	Suppléance perceptive par électro-stimulation linguale embarquée : perspectives pour la prévention des escarres chez le blessé médullaire	Olivier Chenu (TIMC), Nicolas Vuillerme (TIMC), Alexandre Moreau-Gaudry (TIMC), Anthony Fleury (TIMC), Jacques Demongeot (TIMC), Yohan Payan (TIMC)	0711.3786	0.746030

Our Hypothesis

Evaluation

- Use overlapping citations in target and recommended papers
- Count common citations
- Normalise scores with total number of citations
- Determine a threshold for correctness of result
- Calculate accuracy of the recommender
- Aim is to achieve better results than arxiv neural search

References

- Haruna, Khalid; Ismail, Maizatul Akmar; Qazi, Atika; Kakudi, Habeebah Adamu; Hassan, Mohammed; Muaz, Sanah Abdullahi; Chiroma, Haruna (2020).
- <https://github.com/karndeb/Arxiv-Neural-Search>
- DeCLUTR: Deep Contrastive Learning for Unsupervised Textual Representations
- [2104.08821] SimCSE: Simple Contrastive Learning of Sentence Embeddings
- Cohan, Arman, et al. "Specter: Document-level representation learning using citation-informed transformers." arXiv preprint arXiv:2004.07180 (2020).

Thoughts?
Questions?

Thank you.