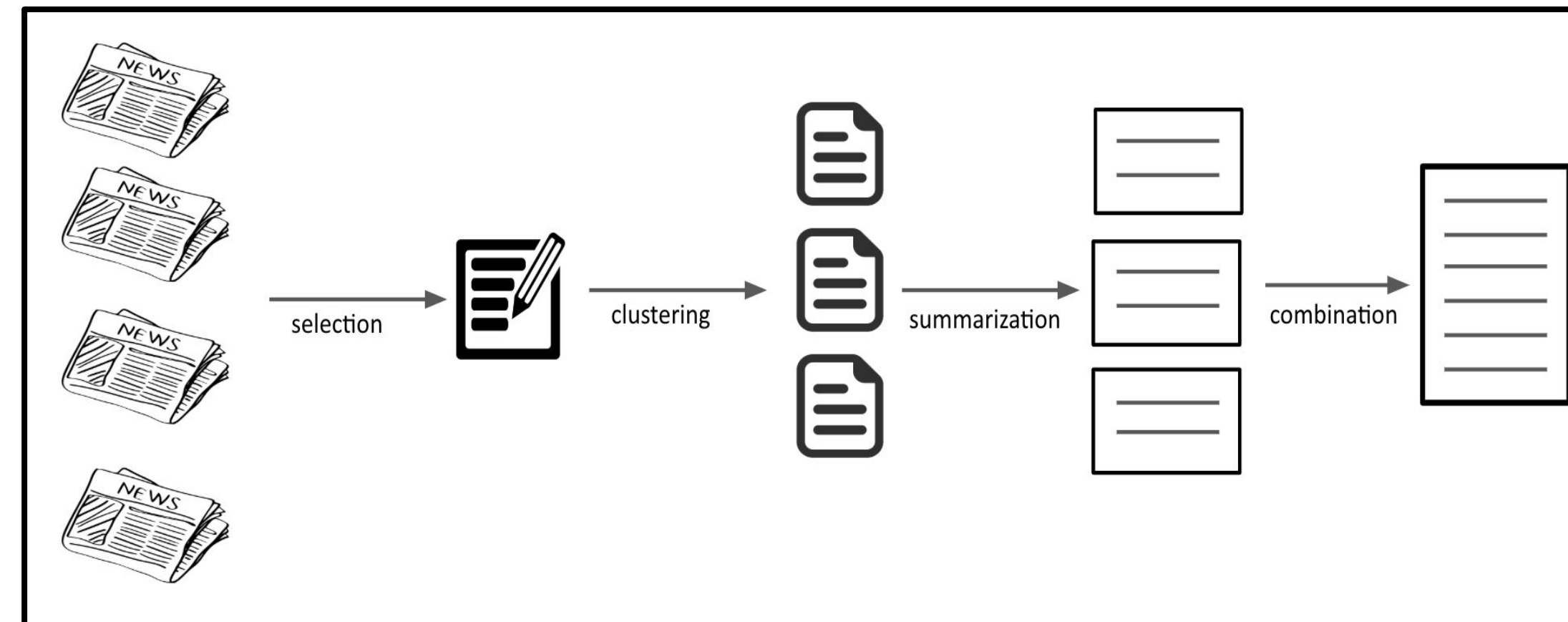


## PROBLEM

- Hundreds of sources of news abound in online and offline media, without even accounting for unofficial news spread through social media.
- People don't have the time to rigorously go through all the different news sources and form their opinion based on reading  $\geq 10$  stories on one topic. Thus, they read what they like, and get even more polarised than they were when they started reading about a topic.
- In today's highly polar political climate, it is important to have tools in place that can **summarize all the relevant news articles for readers and provide them a concise collection of arguments from different viewpoints.**

## APPROACH

- We frame the problem as one of unsupervised multi-document summarization. Our system takes a set of news documents pertaining to one topic of interest as input, and produces a concise summary of the topic based on the source documents as output.
- Unsupervised multi-document summarization is hard. Thus, we design a system with a multi-stage pipeline:
  - a) **Selection:** Select a subset of the original set of sentences from the multiple documents that are most relevant to summary generation.
  - b) **Clustering:** Identify the various aspects/viewpoints of the topic being discussed, and group sets of selected sentences discussing similar aspects together.
  - c) **Summarization:** Summarize individual clusters using extractive or abstractive methods and then concatenate all the summaries together.
- We also design an end-to-end supervised model using a **diversity loss function** to force the model to capture dissimilar viewpoints in the summary.



## PIPELINE MODEL

We optimize the various stages of our pipeline separately, and string them together to build the unsupervised multi-document summarization system.

- **Selection:** We use [1], which introduces a BERT-based model for classifying input sentences as arguments or non-arguments. During summarization, we use either just the arguments, or a combination of argumentative and non-argumentative sentences in a pre-specified ratio.
- **Clustering:** We use K-Means++ for unsupervised clustering on sentence embeddings obtained using BERT. As an alternative, we compute pairwise similarities between argumentative sentences using the argument-similarity BERT model in [1] and use agglomerative hierarchical clustering to cluster the sentences.
- **Summarization:** For K-Means clusters, we choose one sentence per cluster that lies closest to the centroid for that cluster. For agglomerative clusters, we use the sentence with maximum cumulative within-cluster similarity score.

## END-TO-END MODEL

- We modify the supervised extractive summarization loss function and add a diversity term to the existing cross entropy term.
- This ensures that if two sentences have high similarity scores, at least one of them has to be absent from the summary.
- We use the BertSum model in [2], and finetune their pretrained model with our loss function.

$$Loss = \sum_{i=1}^N BCE(v_i, \hat{v}_i) + \sum_{i=1}^N \sum_{j=1}^N v_i \hat{v}_j sim(i, j)$$

## RESULTS

**Rouge-1 Precision and Recall scores for various configurations in the pipeline model**

Config	R1-Recall	R1-Prec
All + K-means	0.2377	0.2253
Args + K-means	0.1591	<b>0.2354</b>
ArgsKM 66.7 %	0.2133	0.2219
ArgsKM 50 %	0.2262	0.2145
ArgsKM 33.3 %	<b>0.2425</b>	0.2038
Args + Agglo	0.195	0.177

We used Document Understanding Conference (DUC) 2004 dataset for evaluation of our unsupervised pipelined system. For our end-to-end model, we will be using DUC-2001-03 datasets for training and 04 for evaluation.

## FUTURE WORK

- Evaluation of the end-to-end model after finetuning with diversity loss
- Human evaluation for diversity of summaries on the DUC dataset as well as our own curated dataset.
- Experiment with abstractive summarization instead of extractive.

## REFERENCES

- [1] Nils Reimers et al. "Classification and Clustering of Arguments with Contextualized Word Embeddings". In: arXiv preprint arXiv:1906.09821(2019)
- [2] Liu, Y., & Lapata, M. (2019). Text summarization with pretrained encoders. *arXiv preprint arXiv:1908.08345*.