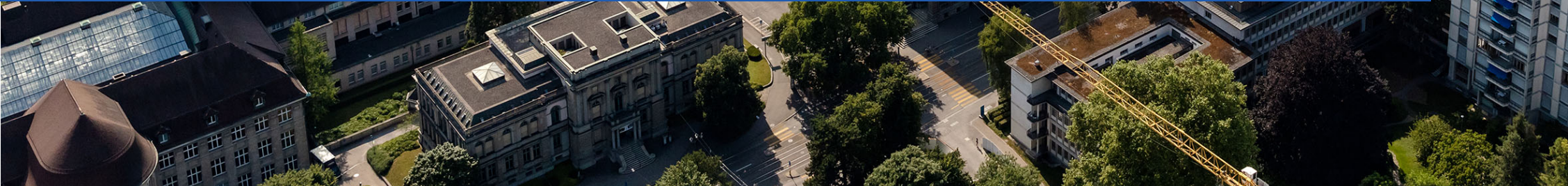




Extractive Summarization of Long Documents using Multi-step Episodic Markov Decision Processes

Nianlong Gu, Elliott Ash, Richard H.R. Hahnloser
ETH Zurich
nianlonggu@gmail.com



Document summarization examples

1. Summarization of a scientific paper



Animal Behaviour
Volume 165, July 2020, Pages 123-132



Contingent parental responses are naturally associated with zebra finch song learning

Samantha Carouso-Peck ^{a, 1}, Otilia Menyhart ^{b, c, 1}, Timothy J. DeVoogd ^a, Michael H. Goldstein ^a ✉

Show more ▾

+ Add to Mendeley 🔗 Share 📄 Cite

<https://doi.org/10.1016/j.anbehav.2020.04.019>

[Get rights and content](#)

Highlights

- Social interaction enhances song learning in some birds, but mechanisms are unknown.
- Zebra finch parents naturally display contingent responses to juvenile plastic song.
- Song learning significantly correlated with contingent maternal 'fluff-ups'.
- Juvenile song accuracy positively correlated with contingent paternal song.

2. Summarization of daily news

MailOnline



Home News U.S. | Sport | TV&Showbiz | Australia | Femail | Health | Science | Money |

Latest Headlines | Pandora Papers | Afghanistan | Covid-19 | Piers Morgan | Prince Harry | Meghan Markle | W

Charles is accused of 'spearheading a monstrosity' over plans to build 2,500 homes on Grade I agricultural land in medieval market town as furious locals leave the area due to 'unbelievable' proposal

- The Duchy of Cornwall has been slammed over a 'masterplan' to build the estate
- They submitted a proposal to build 2,500 greenfield homes in Faversham in Kent
- If approved, it would be in direct conflict with Boris Johnson's promise this week
- And residents have been left furious with one couple moving away after 25 years

By ANDY JEHRING and JOHN ABIONA FOR THE DAILY MAIL and HARRY HOWARD FOR MAILONLINE

PUBLISHED: 22:13 BST, 8 October 2021 | UPDATED: 02:19 BST, 9 October 2021



Share

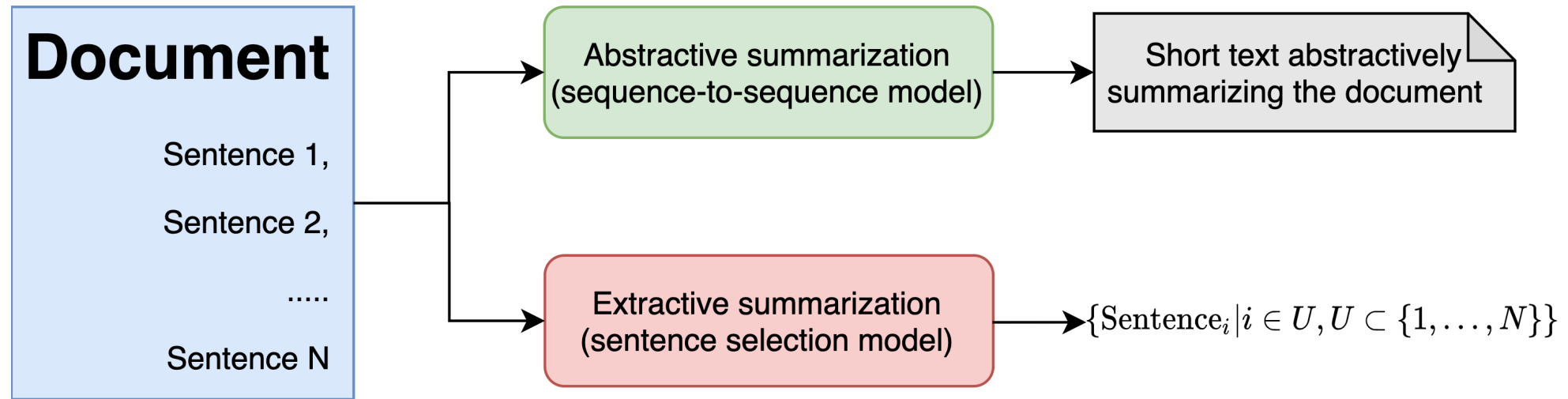


89 shares

2.3k View comments

Prince Charles has been accused of 'spearheading a monstrosity' over plans to build 2,500 greenfield homes in Faversham in Kent as furious locals leave the area due to 'unbelievable' proposal to

Extractive summarization v.s abstractive summarization

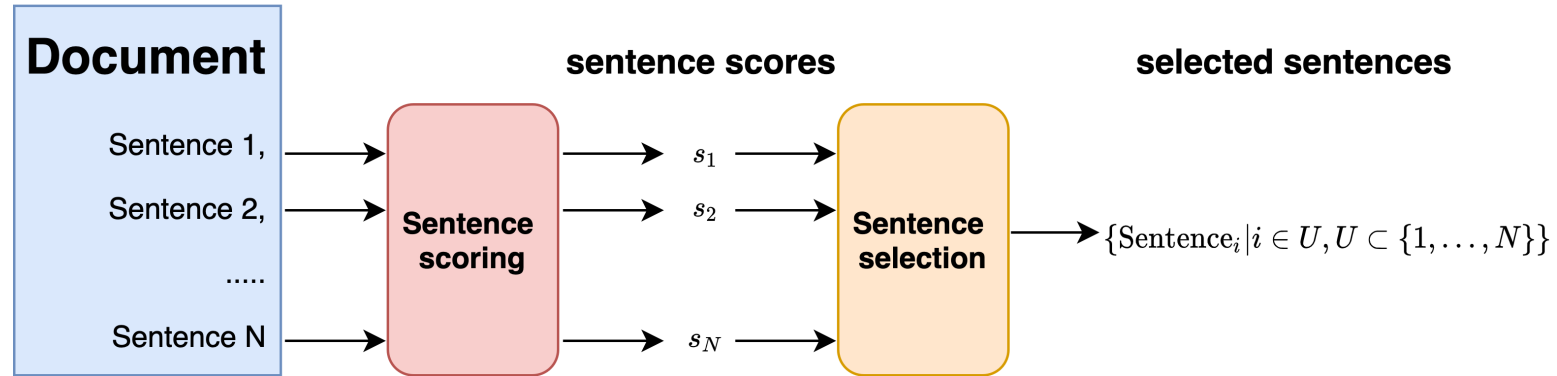


- Abstractive summarization
 - 😊 Better fluency (e.g. models based on Transformer [1])
 - 😞 Low facticity
- **Extractive summarization (main focus of this work)**
 - 😊 Better preservation of original information
 - 😞 Low fluency

[1] Huang L, Cao S, Parulian N, et al. Efficient attentions for long document summarization[J]. arXiv preprint arXiv:2104.02112, 2021.

Current SOTA extractive summarization models

A general framework of current extractive summarization systems (HIBERT [1], BERTSUMEXT [2], Atten-Cont [3])



Common feature:

Sentences scores are computed in one step, and kept fixed during the selection process.

Limits:

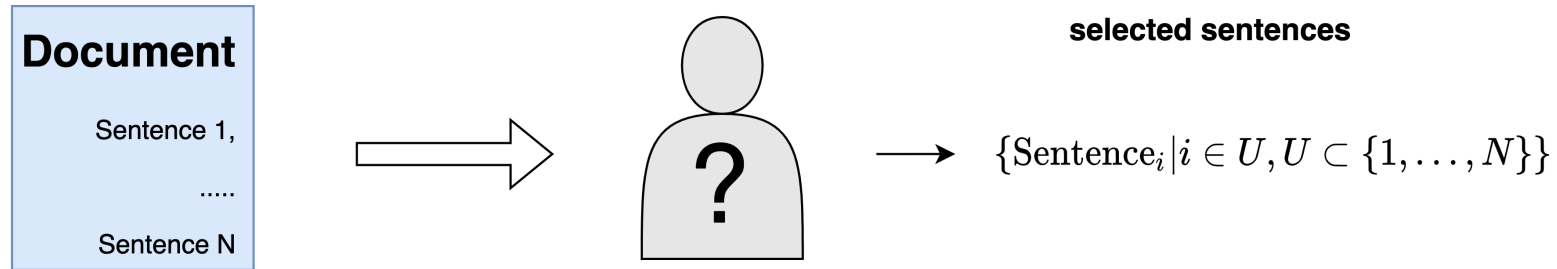
Sentence selection lacks the awareness of extraction history --> Susceptible to redundancy

[1] Zhang X, Wei F, Zhou M. HIBERT: Document level pre-training of hierarchical bidirectional transformers for document summarization[J]. arXiv preprint arXiv:1905.06566, 2019.

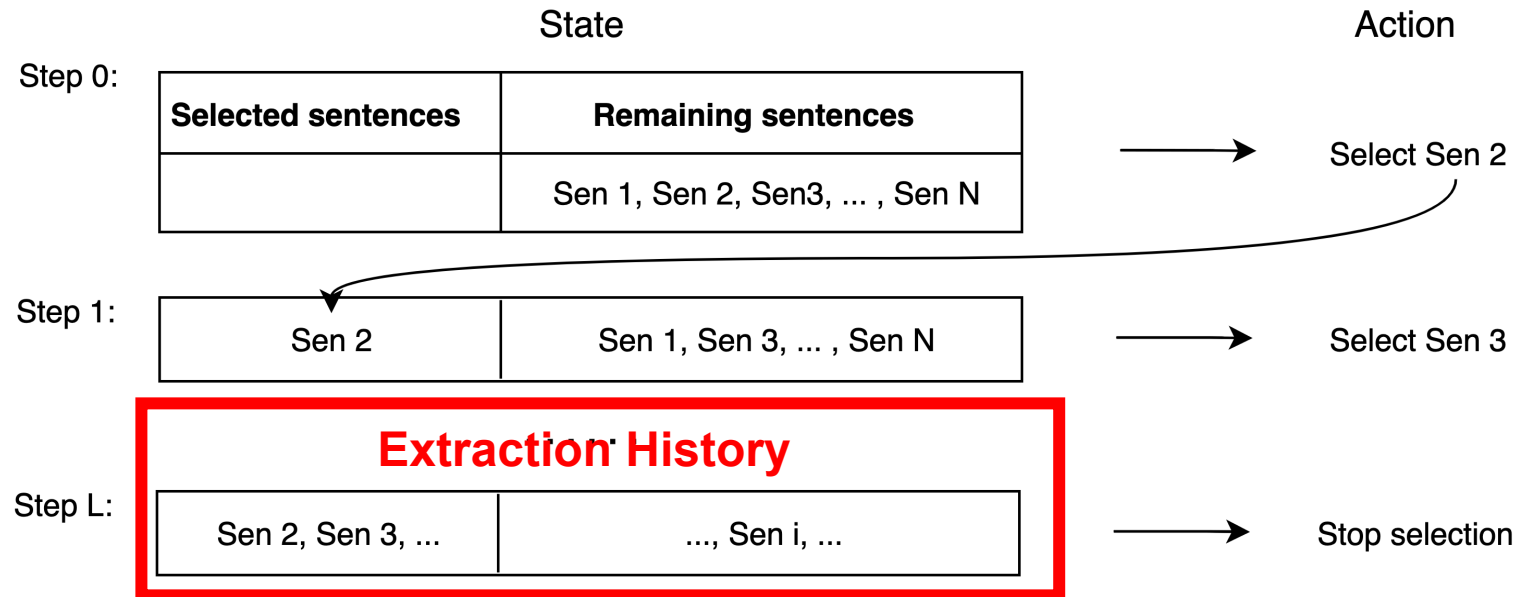
[2] Liu Y, Lapata M. Text summarization with pretrained encoders[J]. arXiv preprint arXiv:1908.08345, 2019.

[3] Xiao W, Carenini G. Extractive summarization of long documents by combining global and local context[J]. arXiv preprint arXiv:1909.08089, 2019.

Our approach to the extractive summarization with history awareness

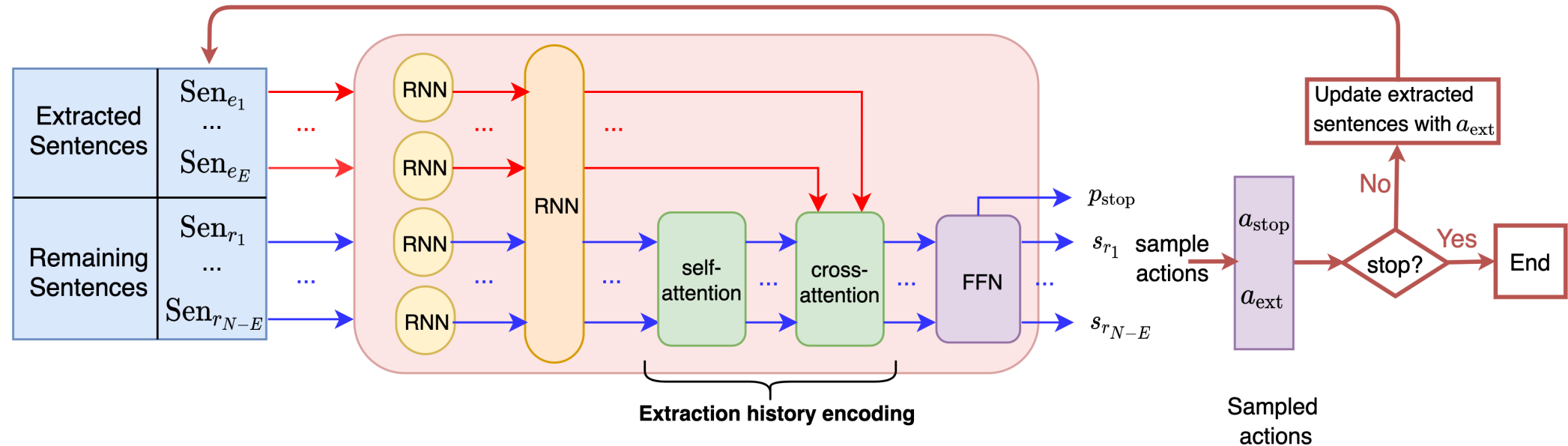


Extractive summarization as multiple steps of selecting sentences without replacement



Extraction history helps to avoid selecting sentences highly similar to already-selected sentences.
(less redundancy)

MemSum: how we encode the extraction history information

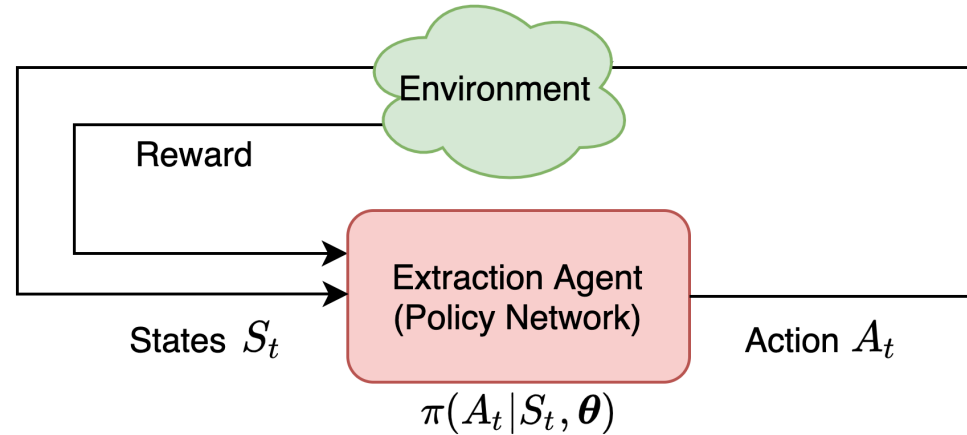


Architecture of our extractive summarization model (MemSum) with extraction history awareness

Training

Training Goal: Increase the ROUGE between extracted summary and gold summary (not differentiable)

Reinforcement Learning



Episode: the process of extracting the entire subset of sentences as the summary, denoted as:

$$(S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T)$$

$$R_t = 0 \text{ for } t = 1, 2, \dots, T-1, \quad R_T = \text{ROUGE score}$$

Policy gradient algorithm (REINFORCE)

Loop forever (for each episode):

Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot | \cdot, \theta)$

Loop for each step of the episode $t = 0, 1, \dots, T-1$:

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k = R_T$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t | S_t, \theta)$$

Extraction history awareness enhances summarization performance

Model	PubMed			arXiv		
	R-1	R-2	R-L	R-1	R-2	R-L
ORACLE	61.99	34.95	56.76	60.00	30.60	53.03
Extractive summarization baselines						
Lead-10	37.45	14.19	34.07	35.52	10.33	31.44
SummaRuNNer	43.89	18.78	30.36	42.81	16.52	28.23
Atten-Cont	44.85	19.70	31.43	43.62	17.36	29.14
Sent-CLF	45.01	19.91	41.16	34.01	8.71	30.41
Sent-PTR	43.30	17.92	39.47	42.32	15.63	38.06
NeuSum	47.46	21.92	42.87	47.49	21.56	41.58
Abstractive summarization baselines						
PEGASUS	45.97	20.15	41.34	44.21	16.95	38.83
BigBird	46.32	20.65	42.33	46.63	19.02	41.77
Dancer	46.34	19.97	42.42	45.01	17.60	40.56
[1] Hepos-Sinkhorn	47.93	20.74	42.58	47.87	20.00	41.50
Hepos-LSH	48.12	21.06	42.72	48.24	20.26	41.78
MemSum (ours)	49.25*	22.94*	44.42*	48.42	20.30	42.54*

Criteria	Experiment I		Experiment II	
	NeuSum	MemSum	NeuSum	MemSum w/o auto-stop
overall	1.58	1.37	1.57	1.38
coverage	1.46	1.49	1.44	1.51
non-redundancy	1.67	1.28*	1.65	1.30*
avg. summ. length				
# of sentences	7.0	5.6*	7.0	7.0
# of words	248.8	189.3*	263.6	239.5*

[1] Huang L, Cao S, Parulian N, et al. Efficient attentions for long document summarization[J]. arXiv preprint arXiv:2104.02112, 2021.

Case study on GovReport dataset

Model	PubMed _{trunc}			GovReport		
	R-1	R-2	R-L	R-1	R-2	R-L
ORACLE	45.12	20.33	40.19	75.56	45.91	72.51
Extractive summarization baselines						
Lead	37.58	12.22	33.44	50.94	19.53	48.45
MatchSum	41.21	14.91	36.75	-	-	-
NeuSum	-	-	-	58.94	25.38	55.80
Abstractive summarization baselines						
Hepos-LSH	-	-	-	55.00	21.13	51.67
Hepos-Sinkhorn	-	-	-	56.86	22.62	53.82
MemSum (ours)	43.08*	16.71*	38.30*	59.43*	28.60*	56.69*

Table 3: Results on PubMed_{trunc} and GovReport.

Human-written Summary:

(...) While CMS is generally required to disallow, or *recoup, federal funds* from states for *eligibility-related improper payments* if the state's *eligibility error rate exceeds 3 percent*, it has not done so for decades, (...) CMS *issued revised procedures through which it can recoup funds for eligibility errors, beginning in fiscal year 2022*. (...)

Hepos-Sinkhorn (abstractive):

(...) The selected states also reported that they did not have adequate processes to address these issues. CMS has taken steps to improve its oversight of the Medicaid program, including issuing guidance to states on the use of MAGI-exempt bases for determining eligibility, but these efforts have not been fully implemented. (...)

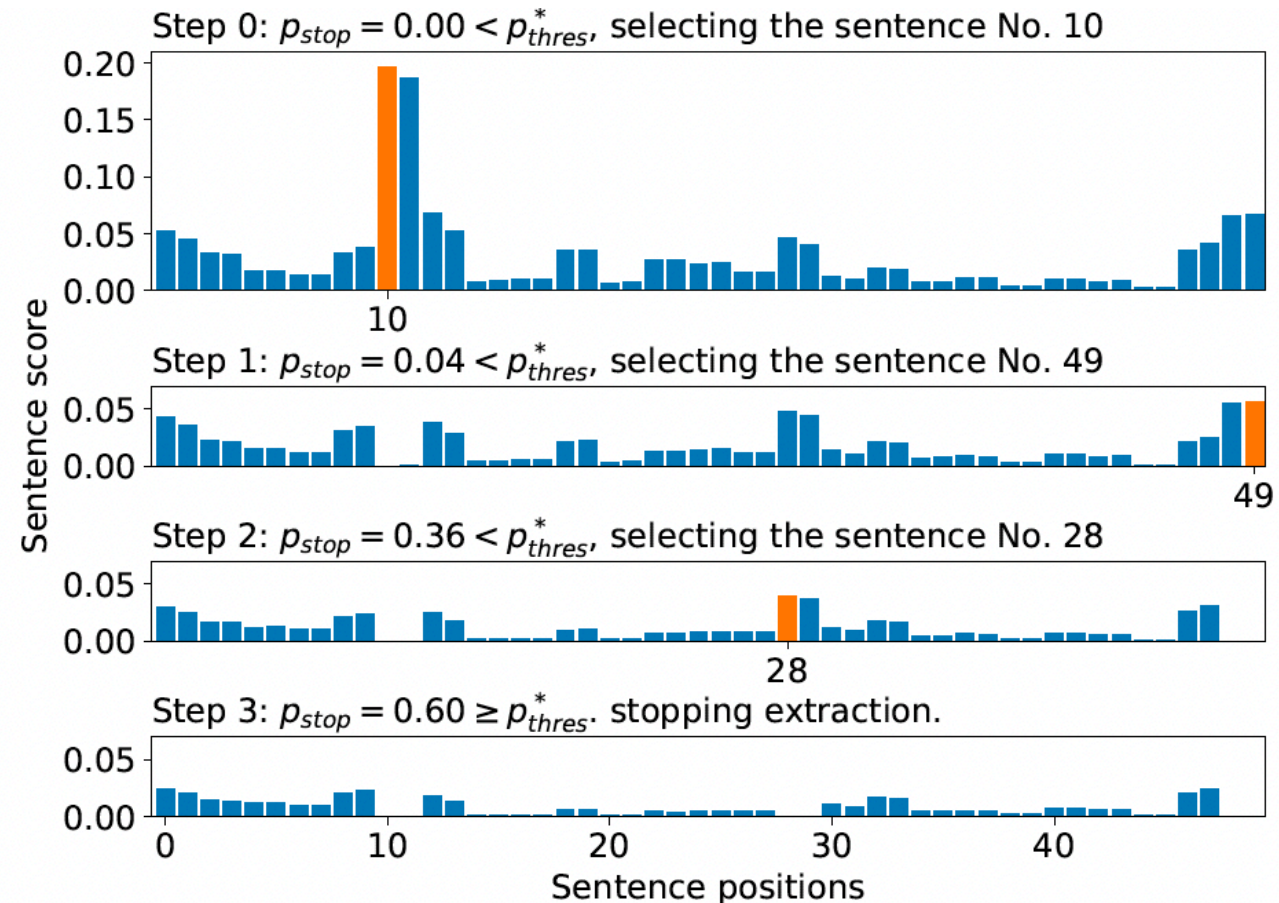
MemSum (ours, extractive):

(...) implemented its statutory requirement to *recoup funds* associated with Medicaid *eligibility-related improper payments* for states with an *eligibility error rate above 3 percent* through its MEQC program. (...) However, the agency has *introduced new procedures through which it can, under certain circumstances, begin to recoup funds based on eligibility errors in fiscal year 2022*. (...)

How MemSum avoid redundancy with history awareness – a case study

$$[\text{Sen}_1, \text{Sen}_2, \dots, \text{Sen}_N] \rightarrow [\text{Sen}_1, \text{Sen}_1, \text{Sen}_2, \text{Sen}_2, \dots, \text{Sen}_N, \text{Sen}_N]$$

- The sentence scores of 50 sentences computed by the MemSum model at extraction steps 0 to 3.
- In the document there is artificial redundancy where the $2n$ and the $2n+1$ sentences are identical ($n=0, 1, \dots, 24$)



Conclusion

- We treat extractive summarization as a multi-step episodic Markov decision process.
- We show that the awareness of the extraction history helps to avoid redundancies in documents.
- Our model outperforms both extractive and abstractive summarization models on PubMed, arXiv and GovReport datasets.

Outlook

- Evaluate the coherence of the extracted summary.
- Explore the extractive-abstractive summarization pipelines for long document summarization.