

# AUTOMATIC GENERATION OF TOPIC LABELS

---

Areej Alokaili\*  
areej.okaili@sheffield.ac.uk  
University of  
Sheffield

Nikolaos Aletras  
n.aletras@sheffield.ac.uk  
University of  
Sheffield

Mark Stevenson  
mark.stevenson@sheffield.ac.uk  
University of  
Sheffield

SIGIR 2020

M11115020 吳宇祥 M11015084 魏向晨  
M11115079 張家維 M11115099 郭建鴻  
M11115073 張晉嘉 M11115103 施昱民

# TABLE OF CONTENTS

---

01 INTRODUCTION

02 MODEL

03 DATA

04 EXPERIMENT & CONCLUSION



# 01

## INTRODUCTION

---

# TOPIC MODEL

---

- Identifying the underlying themes in document collections
- An unsupervised machine learning technique
- A quick and easy way to start analyzing data
- The topic represented by a list of terms ranked by their probability can be difficult to interpret
  - Various approaches have been developed to assign descriptive labels to topics

# DIFFICULTIES

---

- Previous work on the automatic assignment of labels to topics has relied on a two-stage approach:
  - Candidate labels are retrieved from a large pool
  - Re-ranked based on semantic similarity to the topic terms
- These extractive approaches can only assign candidate labels from a restricted set that may not include any suitable ones

# SOLUTION

---

- Sequence-to-sequence neural-based approach
  - The model is trained over a new large synthetic dataset created using [distant supervision](#)
  - The method is evaluated by comparing the labels which generates to ones rated by humans.

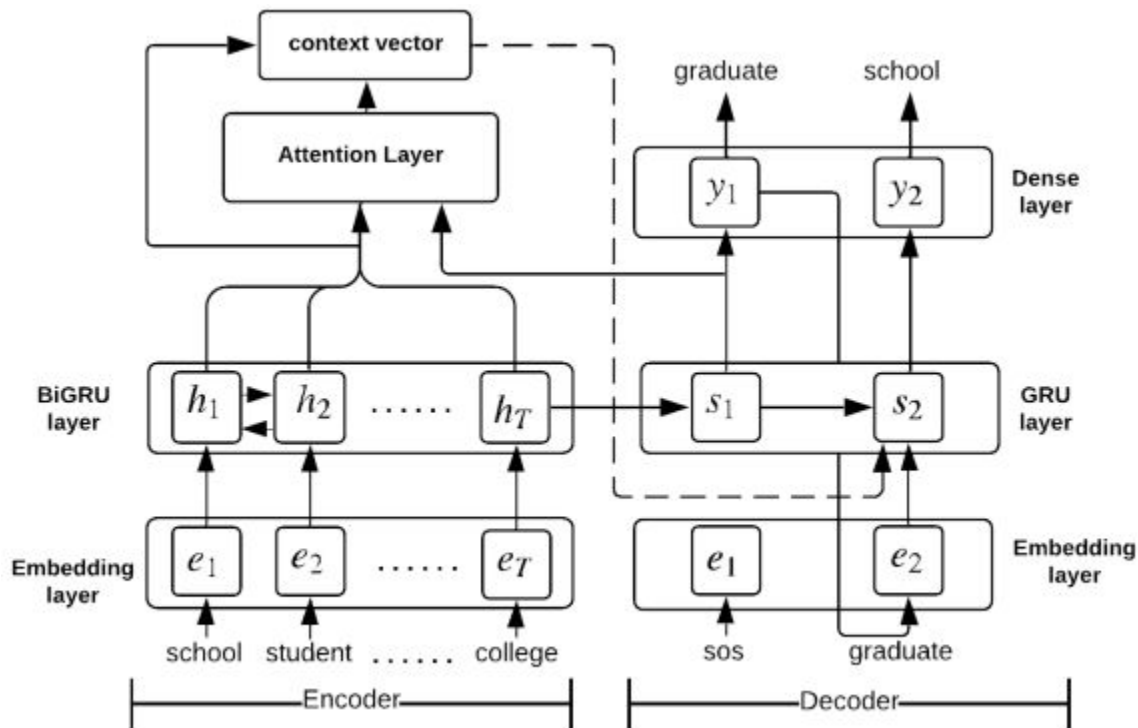


# 02

## MODEL

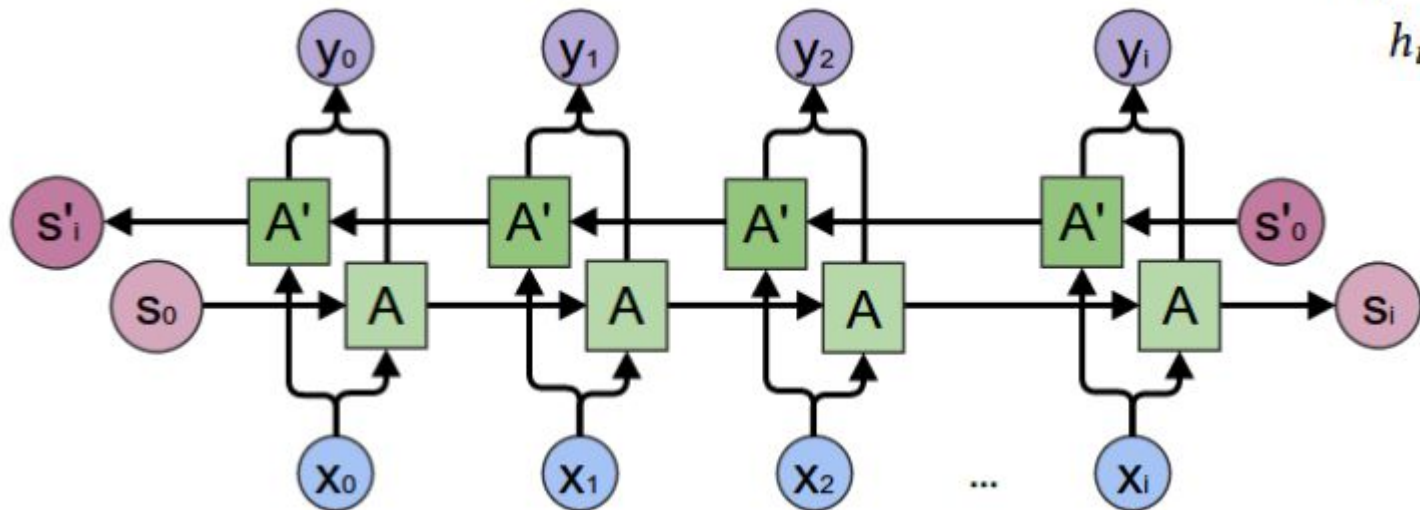
---

# ARCHITECTURE





# BIDIRECTIONAL GRU



$$hf_t = GRU(x_t, h_{t-1})$$

$$hb_t = GRU(x_t, h_{t-1})$$

$$h_t = [hf_t; hb_t]$$

# DECODER

- GRU Layer :

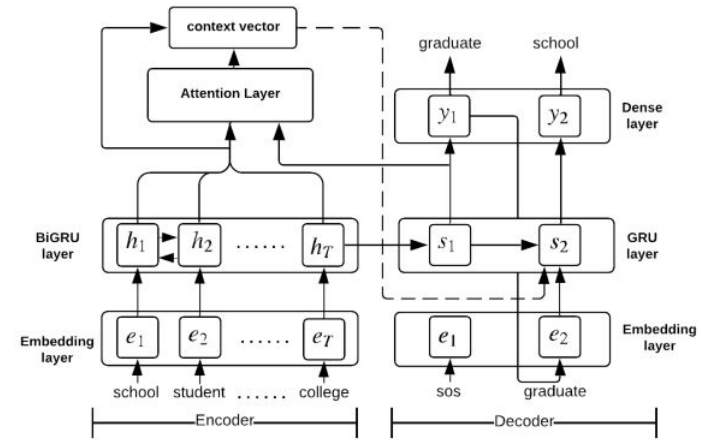
$$s_t = GRU(y_{t-1}, s_{t-1}, c_t)$$

$$c_t = \sum_{j=1}^{T_x} \alpha_{tj} h_j$$

- Dense Layer :

$$P(y_t | \{y_1, \dots, y_{t-1}\} X) = Dense(s_t)$$

$$y_t = \operatorname{argmax}(P(y_t | \{y_1, \dots, y_{t-1}\}, X))$$





# 03

## Dataset

---

# DATASET

---

Total 300,000 pairs of topics and labels from **ds\_wiki\_tfidf** & **ds\_wiki\_sent**

After standard preprocessing steps, which were applied to clean datasets such as removal of numbers, special characters and stop words

Total 250,506 pairs of topics and labels left

- training data:226,282(91%)
- validate data:12,424(4%)
- testing data:11,800(5%)

# DATASET

Training data :

- ds\_wiki\_tfidf & ds\_wiki\_sent

Testing data :

- ds\_wiki\_tfidf & ds\_wiki\_sent
- topics\_bhatia & topics\_bhatia\_tfidf

**Automatic Labelling of Topics with Neural Embeddings**

Additional testing data :

- wiki-dataset:extract from enwiki-20221120-pages-articles

<a href="#">zhwiki-latest-linktarget.sql.gz-rss.xml</a>	21-Nov-2022 18:06	766
<a href="#">zhwiki-latest-md5sums.txt</a>	21-Nov-2022 21:59	9278
<a href="#">zhwiki-latest-page.sql.gz</a>	20-Nov-2022 16:13	243472618
<a href="#">zhwiki-latest-page.sql.gz-rss.xml</a>	21-Nov-2022 18:07	748
<a href="#">zhwiki-latest-page_props.sql.gz</a>	20-Nov-2022 16:14	40657538
<a href="#">zhwiki-latest-page_props.sql.gz-rss.xml</a>	21-Nov-2022 18:07	766
<a href="#">zhwiki-latest-page_restrictions.sql.gz</a>	20-Nov-2022 16:14	519176
<a href="#">zhwiki-latest-page_restrictions.sql.gz-rss.xml</a>	21-Nov-2022 18:07	787
<a href="#">zhwiki-latest-pagelinks.sql.gz</a>	20-Nov-2022 16:03	1121341240
<a href="#">zhwiki-latest-pagelinks.sql.gz-rss.xml</a>	21-Nov-2022 18:06	763
<a href="#">zhwiki-latest-pages-articles-multistream-index...</a>	21-Nov-2022 02:47	35378851
<a href="#">zhwiki-latest-pages-articles-multistream-index...</a>	21-Nov-2022 21:59	835
<a href="#">zhwiki-latest-pages-articles-multistream-index1...</a>	21-Nov-2022 02:20	702484
<a href="#">zhwiki-latest-pages-articles-multistream-index1...</a>	21-Nov-2022 21:59	868
<a href="#">zhwiki-latest-pages-articles-multistream-index2...</a>	21-Nov-2022 02:21	1913845
<a href="#">zhwiki-latest-pages-articles-multistream-index2...</a>	21-Nov-2022 21:59	883
<a href="#">zhwiki-latest-pages-articles-multistream-index3...</a>	21-Nov-2022 02:23	3867086
<a href="#">zhwiki-latest-pages-articles-multistream-index3...</a>	21-Nov-2022 21:59	886
<a href="#">zhwiki-latest-pages-articles-multistream-index4...</a>	21-Nov-2022 02:26	6279169
<a href="#">zhwiki-latest-pages-articles-multistream-index4...</a>	21-Nov-2022 21:59	889
<a href="#">zhwiki-latest-pages-articles-multistream-index4...</a>	21-Nov-2022 02:21	2423050
<a href="#">zhwiki-latest-pages-articles-multistream-index4...</a>	21-Nov-2022 21:59	889
<a href="#">zhwiki-latest-pages-articles-multistream-index5...</a>	21-Nov-2022 02:23	4631468
<a href="#">zhwiki-latest-pages-articles-multistream-index5...</a>	21-Nov-2022 21:59	889
<a href="#">zhwiki-latest-pages-articles-multistream-index5...</a>	21-Nov-2022 02:29	2474806
<a href="#">zhwiki-latest-pages-articles-multistream-index5...</a>	21-Nov-2022 21:59	889
<a href="#">zhwiki-latest-pages-articles-multistream-index6...</a>	21-Nov-2022 02:35	6449017
<a href="#">zhwiki-latest-pages-articles-multistream-index6...</a>	21-Nov-2022 21:59	889
<a href="#">zhwiki-latest-pages-articles-multistream-index6...</a>	04-Nov-2022 11:20	6491030
<a href="#">zhwiki-latest-pages-articles-multistream-index6...</a>	07-Nov-2022 02:12	889
<a href="#">zhwiki-latest-pages-articles-multistream-index6...</a>	21-Nov-2022 02:34	6598255
<a href="#">zhwiki-latest-pages-articles-multistream-index6...</a>	21-Nov-2022 21:59	889
<a href="#">zhwiki-latest-pages-articles-multistream.xml.bz2</a>	21-Nov-2022 02:47	2694577747
<a href="#">zhwiki-latest-pages-articles-multistream.xml.bz...</a>	21-Nov-2022 21:59	817
<a href="#">zhwiki-latest-pages-articles-multistream1.xml.p...</a>	21-Nov-2022 02:20	222027702
<a href="#">zhwiki-latest-pages-articles-multistream1.xml.p...</a>	21-Nov-2022 21:59	850
<a href="#">zhwiki-latest-pages-articles-multistream2.xml.p...</a>	21-Nov-2022 02:21	281203796
<a href="#">zhwiki-latest-pages-articles-multistream2.xml.p...</a>	21-Nov-2022 21:59	865
<a href="#">zhwiki-latest-pages-articles-multistream3.xml.p...</a>	21-Nov-2022 02:23	366642855
<a href="#">zhwiki-latest-pages-articles-multistream3.xml.p...</a>	21-Nov-2022 21:59	868
<a href="#">zhwiki-latest-pages-articles-multistream4.xml.p...</a>	21-Nov-2022 02:26	291756311
<a href="#">zhwiki-latest-pages-articles-multistream4.xml.p...</a>	21-Nov-2022 21:59	871
<a href="#">zhwiki-latest-pages-articles-multistream4.xml.p...</a>	21-Nov-2022 02:21	138309890
<a href="#">zhwiki-latest-pages-articles-multistream4.xml.p...</a>	21-Nov-2022 21:59	871
<a href="#">zhwiki-latest-pages-articles-multistream5.xml.p...</a>	21-Nov-2022 02:23	318977219
<a href="#">zhwiki-latest-pages-articles-multistream5.xml.p...</a>	21-Nov-2022 21:59	871
<a href="#">zhwiki-latest-pages-articles-multistream5.xml.p...</a>	21-Nov-2022 02:29	198039574
<a href="#">zhwiki-latest-pages-articles-multistream5.xml.p...</a>	21-Nov-2022 21:59	871
<a href="#">zhwiki-latest-pages-articles-multistream6.xml.p...</a>	21-Nov-2022 02:35	473913869
<a href="#">zhwiki-latest-pages-articles-multistream6.xml.p...</a>	21-Nov-2022 21:59	871



# 04

## Experiment

---

# EVALUATION

---

$$score\_topic_t = \max_{i=[1,...,n]} BERTScore(l_t, gold\_l_{ti})$$

$$score\_model = \frac{1}{T} \sum_{t=1}^T score\_topic_t$$

BERTScore: Evaluating Text Generation with BERT (ArXiv 2019)

# EXPERIMENT

			BERTScore			
			P	R	F	
Baselines		Top-2 label	0.902	0.912	0.902	
		Top-3 label	0.870	0.903	0.882	
Train data	ds_wiki_tfidf	Test data	topics_bhatia	0.922* <sup>†</sup>	0.928* <sup>†</sup>	0.922* <sup>†</sup>
			topics_bhatia_tfidf	0.926* <sup>†</sup>	0.930* <sup>†</sup>	0.925* <sup>†</sup>
	ds_wiki_sent	topics_bhatia	0.919 <sup>†</sup>	0.926 <sup>†</sup>	0.919 <sup>†</sup>	
		topics_bhatia_tfidf	<b>0.930*<sup>†</sup></b>	<b>0.933*<sup>†</sup></b>	<b>0.929*<sup>†</sup></b>	

\* and † denote statistically significant difference ( $p < 0.001$ ) compared to Top-2 label and Top-3 label, respectively.

Our result

ds_wiki_tfidf	Another wiki-data	0.916	0.868	0.889
---------------	-------------------	-------	-------	-------



# DATASET EXTRACT

<https://dumps.wikimedia.org/zhwiki/latest/>

```
from gensim.corpora.wikicorpus import extract_pages, filter_wiki
import bz2file
import re
import opencc
from tqdm import tqdm
import codecs
wiki = extract_pages(bz2file.open('enwiki-20221120-pages-articles-multistream19.xml'))

def wiki_replace(d):
    s = d[i]
    s = re.sub('.*?[\[\s\S]*?|)', '', s)
    s = re.sub('<gallery>[\s\S]*?</gallery>', '', s)
    s = re.sub('.*?{{{[^\}]*?|}}}', '\[\[\2\]', s)
    s = filter_wiki(s)
    s = re.sub('.*?*\n\{2}', '', s)
    s = re.sub('\n+', '\n', s)
    s = re.sub('\n[:;]\n+', '\n', s)
    s = re.sub('\n==', '\n\n==', s)
    s = u[' ' + d[0] + u'] \n' + s
    print(s)
    return s.strip()

i = 0
f = codecs.open('wiki.txt', 'w', encoding='utf-8')
w = tqdm(wiki, desc='已獲取篇文章')
for d in w:
    if not re.findall('[a-zA-Z]+:', d[0]) and d[0] and not re.findall(u'^(#)', d[1]):
        s = wiki_replace(d)
        f.write(s + '\n\n\n')
        i += 1
        if i % 100 == 0:
            w.set_description(u'已獲取篇文章%i')

f.close()
```

zhwiki-latest-linktarget.sql.gz-rss.xml	21-Nov-2022 18:06	766
zhwiki-latest-mdsums.txt	21-Nov-2022 21:59	9278
zhwiki-latest-page.sql.gz	20-Nov-2022 16:13	243472618
zhwiki-latest-page.sql.gz-rss.xml	21-Nov-2022 18:07	748
zhwiki-latest-page_props.sql.gz	20-Nov-2022 16:14	40657538
zhwiki-latest-page_props.sql.gz-rss.xml	21-Nov-2022 18:07	766
zhwiki-latest-page_restrictions.sql.gz	20-Nov-2022 16:14	519176
zhwiki-latest-page_restrictions.sql.gz-rss.xml	21-Nov-2022 18:07	787
zhwiki-latest-pagelinks.sql.gz	20-Nov-2022 16:03	1121341240
zhwiki-latest-pagelinks.sql.gz-rss.xml	21-Nov-2022 18:06	763
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:47	35378851
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	835
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:20	702484
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	868
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:21	1913845
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	883
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:23	3867086
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	886
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:26	6279169
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	889
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:21	24230950
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	889
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:23	4631468
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	889
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:29	2474886
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	889
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:35	6449017
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	889
zhwiki-latest-pages-articles-multistream-index...	04-Nov-2022 11:20	6491030
zhwiki-latest-pages-articles-multistream-index...	07-Nov-2022 02:12	889
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 02:34	6598255
zhwiki-latest-pages-articles-multistream-index...	21-Nov-2022 21:59	889
zhwiki-latest-pages-articles-multistream.xml.bz2	21-Nov-2022 02:47	269457747
zhwiki-latest-pages-articles-multistream.xml.bz2	21-Nov-2022 21:59	817
zhwiki-latest-pages-articles-multistream.xml.p...	21-Nov-2022 02:20	222027702
zhwiki-latest-pages-articles-multistream.xml.p...	21-Nov-2022 21:59	850
zhwiki-latest-pages-articles-multistream2.xml.p...	21-Nov-2022 02:21	281203796
zhwiki-latest-pages-articles-multistream2.xml.p...	21-Nov-2022 21:59	865
zhwiki-latest-pages-articles-multistream3.xml.p...	21-Nov-2022 02:23	366642855
zhwiki-latest-pages-articles-multistream3.xml.p...	21-Nov-2022 21:59	868
zhwiki-latest-pages-articles-multistream4.xml.p...	21-Nov-2022 02:26	291756311
zhwiki-latest-pages-articles-multistream4.xml.p...	21-Nov-2022 21:59	871
zhwiki-latest-pages-articles-multistream4.xml.p...	21-Nov-2022 02:21	138309890
zhwiki-latest-pages-articles-multistream4.xml.p...	21-Nov-2022 21:59	871
zhwiki-latest-pages-articles-multistream5.xml.p...	21-Nov-2022 02:23	318977219
zhwiki-latest-pages-articles-multistream5.xml.p...	21-Nov-2022 21:59	871
zhwiki-latest-pages-articles-multistream5.xml.p...	21-Nov-2022 02:29	198039574
zhwiki-latest-pages-articles-multistream5.xml.p...	21-Nov-2022 21:59	871
zhwiki-latest-pages-articles-multistream6.xml.p...	21-Nov-2022 02:35	473913869
zhwiki-latest-pages-articles-multistream6.xml.p...	21-Nov-2022 21:59	871

# CASE STUDY ANOTHER WIKI DATA

---

<https://dumps.wikimedia.org/enwiki/latest/>

Original data(term) : vmware server virtual oracle update virtualization application infrastructure management microsoft plesk web hosting dns windows linux accounts reseller software aps sql administrator panel cloudlinux automation swsoft platform packaging designed versions

Ground truth : cloud computing, microsoft exchange server, vmware, web application, virtualization, operating system

Predict : vmware server

# CONCLUSION

---

- Present the first seq2seq model to generate textual labels for automatically generated topics.
- Presented a dataset built from Wikipedia and use BERTScore to measure the similarities between the generated labels and gold standard labels.



# THANKS

CREDITS: This presentation template was created by **Slidesgo**, including icons by **Flaticon**, infographics & images by **Freepik**

# 實作心得

---

- Present the first seq2seq model to generate textual labels for automatically generated topics.
- Presented a dataset built from Wikipedia and use BERTScore to measure the similarities between the generated labels and gold standard labels.

# 實作心得

---

- Present the first seq2seq model to generate textual labels for automatically generated topics.
- Presented a dataset built from Wikipedia and use BERTScore to measure the similarities between the generated labels and gold standard labels.