

# Heterogeneous Graph Neural Networks for Keyphrase Generation

**Jiacheng Ye<sup>1\*</sup>, Ruijian Cai<sup>1\*</sup>, Tao Gui<sup>2†</sup> and Qi Zhang<sup>1†</sup>**

<sup>1</sup>School of Computer Science, Shanghai Key Laboratory of Intelligent Information Processing,  
Fudan University, Shanghai, China

<sup>2</sup>Institute of Modern Languages and Linguistics, Fudan University  
{yejc19, 19210240253, tgui, qz}@fudan.edu.cn

EMNLP 2021

[https://github.com/jiacheng-ye/kg\\_gater](https://github.com/jiacheng-ye/kg_gater)



**gesis**  
Leibniz-Institut  
für Sozialwissenschaften



Reported by Dongdong Hu

# Introduction

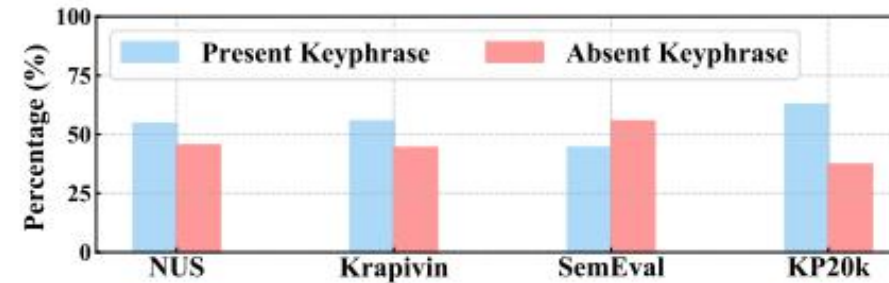


Figure 1: Proportion of present and absent keyphrases among four datasets. Although the previous methods for keyphrase generation have shown promising results on present keyphrase predictions, they are not yet satisfactory on the absent keyphrase predictions, which also occupy a large proportion.

Relying solely on the source document can result in generating uncontrollable and inaccurate absent keyphrases.

# Method

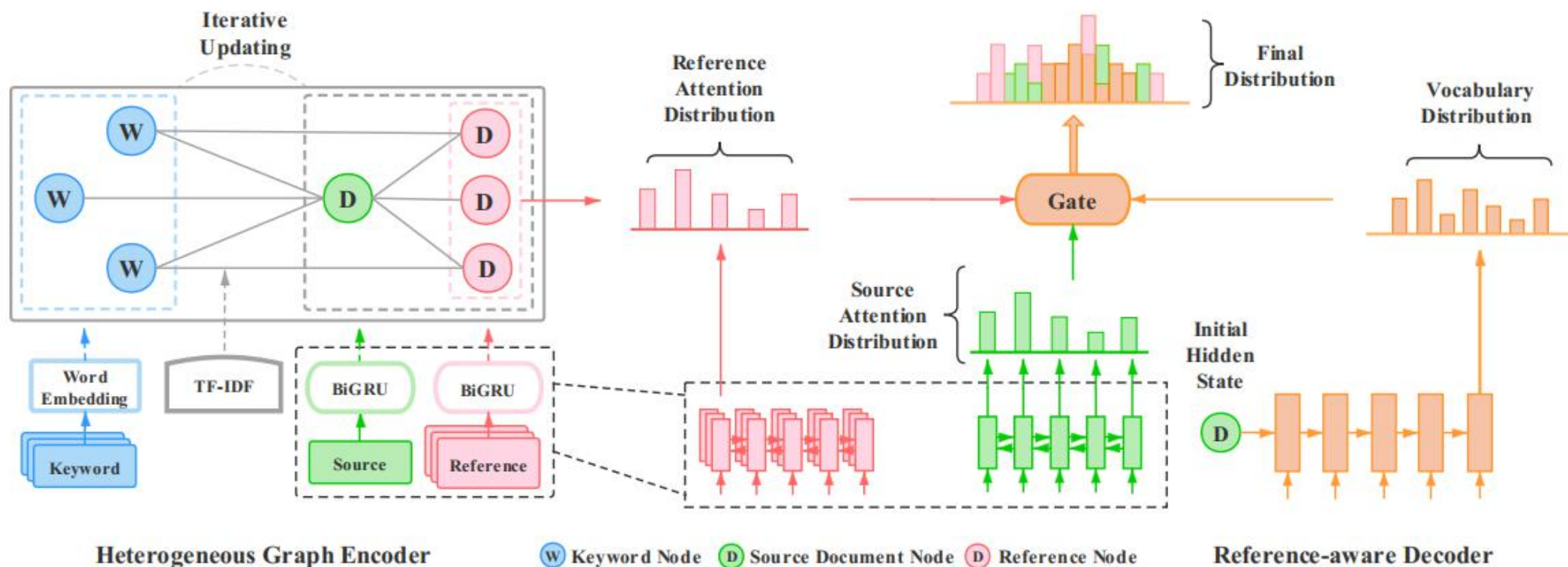


Figure 2: Graphical illustration of our proposed GATER. We first retrieve references using the source document, where each reference is the concatenation of document and keyphrases pair from the training set. Then we construct a heterogeneous graph and perform iterative updating. Finally, the source document node is extracted to decode the keyphrase sequence with a hierarchical attention and copy mechanism.



we first represent the source document and all the reference candidates as TF-IDF weighted uni/bi-gram vectors. Then, the most similar  $K$  references  $\mathcal{X}^r = \{\mathbf{x}^{r_i}\}_{i=1,\dots,K}$  are retrieved by comparing the cosine similarities of the vectors of the source document and all the references.

# Method

## Graph Construction



$$G = \{V, E\}$$

$$V = V_w \cup V_d$$

$$V_w = \{w_i\} (i \in \{1, \dots, m\})$$

$$V_d = \mathbf{x} \cup \mathcal{X}^r$$

$$E = E_{d2d} \cup E_{w2d}$$

$$E_{d2d} = \{e_k\} (k \in \{1, \dots, K\})$$

$$E_{w2d} = \{e_{i,j}\} (i \in \{1, \dots, m\}, j \in \{1, \dots, K+1\})$$

Similarly, we also infuse TF-IDF values in the edge weights of  $E_{d2d}$  as a prior statistical  $n$ -gram similarity between documents.

# Method

## Graph Initializers

document node

$$e^w$$

$$\mathbf{d} = [\vec{m}_1; \overleftarrow{m}_{L_x}] \text{ and } m_i = [\vec{m}_i; \overleftarrow{m}_i]$$

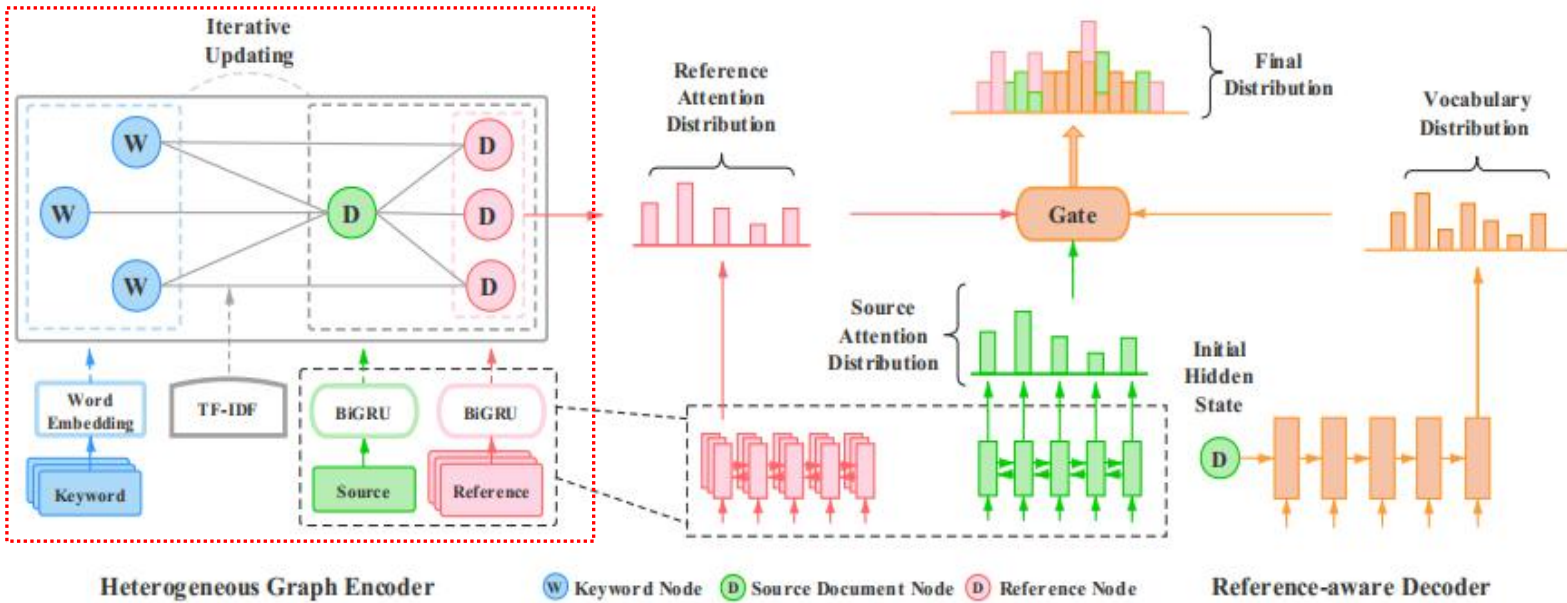
keyword node

$$\mathbf{w}_i = e^w(w_i).$$

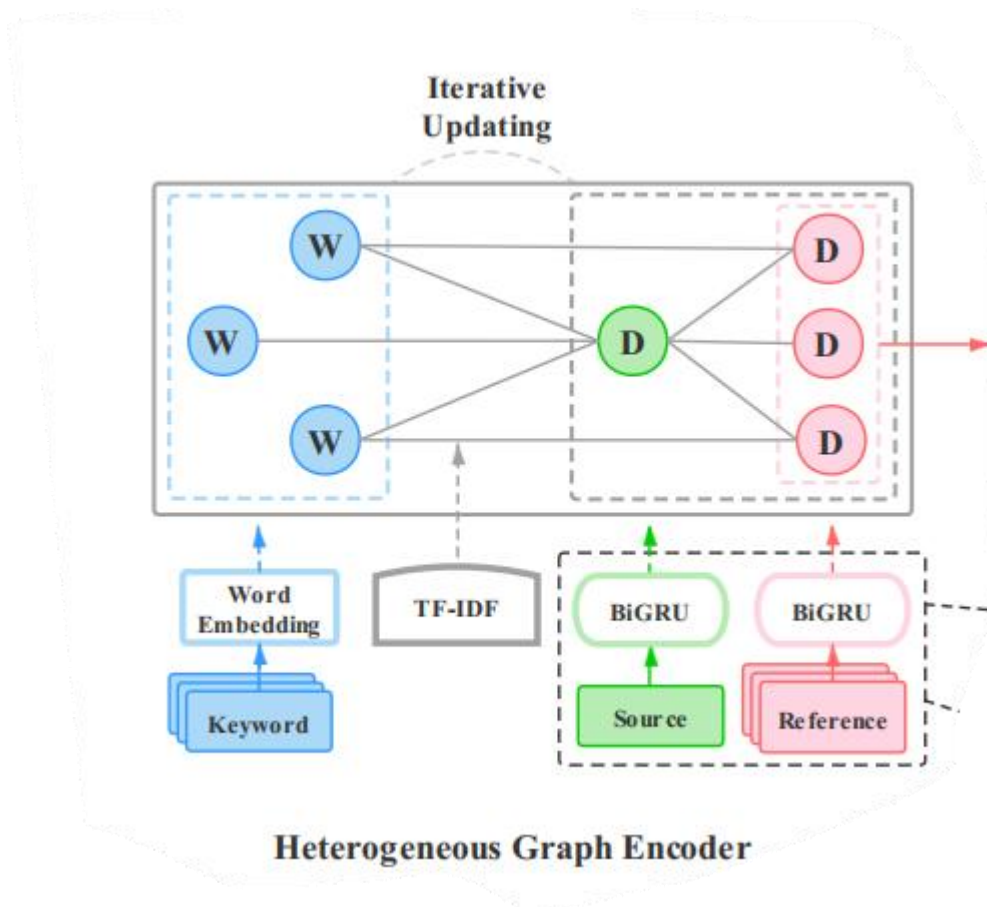
edge initializers

$$E_{d2d} \quad E_{d2w}$$

$$e^{d2d} \text{ and } e^{w2d}$$



# Method



$$z_{ij} = \text{LeakyReLU}(\mathbf{w}_a^T [\mathbf{W}_q \mathbf{h}_i; \mathbf{W}_k \mathbf{h}_j; \mathbf{e}_{ij}])$$

$$\alpha_{ij} = \text{softmax}_j(z_{ij}) = \frac{\exp(z_{ij})}{\sum_{k \in \mathcal{N}_i} \exp(z_{ik})}$$

$$\mathbf{u}_i = \sigma\left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \mathbf{W}_v \mathbf{h}_j\right),$$

GAT ( $\mathbf{H}, \mathbf{H}, \mathbf{H}, \mathbf{E}$ ) to denote the GAT aggregating layer

$\mathbf{H}$  is used for query, key, and value

$$\mathbf{H}_w^1 = \text{FFN}(\text{GAT}(\mathbf{H}_w^0, \mathbf{H}_d^0, \mathbf{H}_d^0, \mathbf{E}_{w2d}) + \mathbf{H}_w^0)$$

$$\mathbf{H}_d^1 = \text{FFN}(\text{GAT}(\mathbf{H}_d^0, \mathbf{H}_w^1, \mathbf{H}_w^1, \mathbf{E}_{w2d}) + \mathbf{H}_d^0)$$

$$\mathbf{H}_d^1 = \text{FFN}(\text{GAT}(\mathbf{H}_d^1, \mathbf{H}_d^1, \mathbf{H}_d^1, \mathbf{E}_{d2d}) + \mathbf{H}_d^1), \quad (2)$$

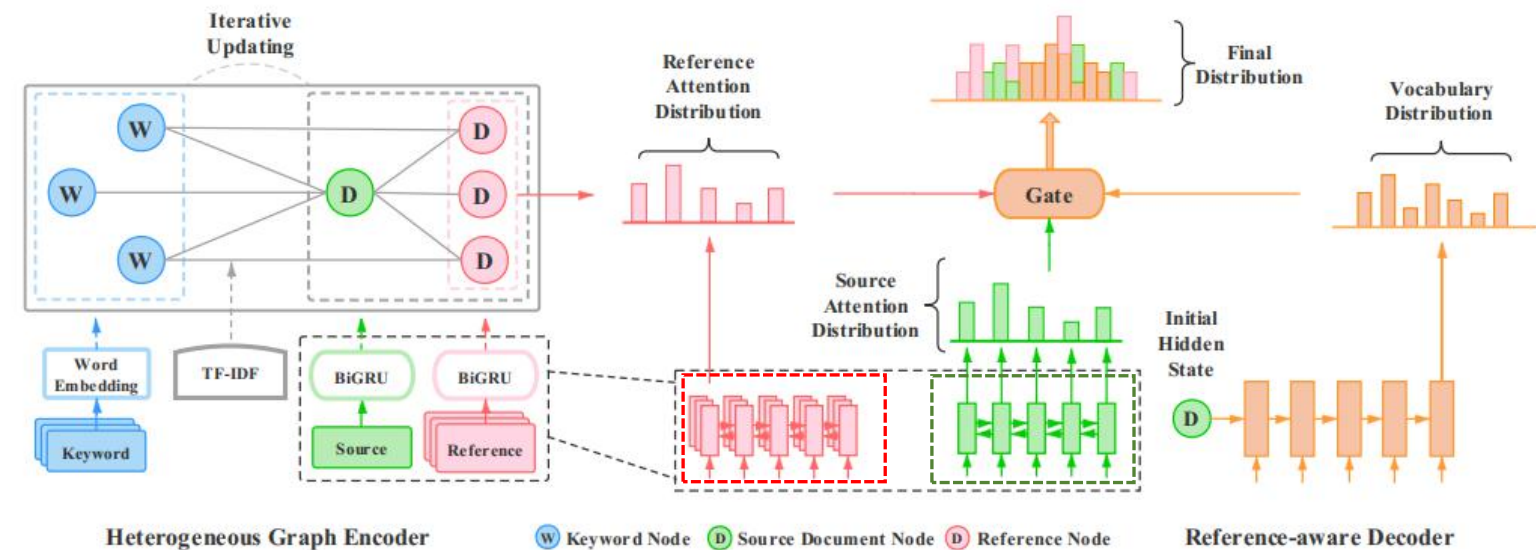
we introduce a residual connection and position-wise feed-forward (FFN) layer consisting of two linear transformations.

seperate  $\mathbf{H}_d^1$  into  $\mathbf{d}^s$  and  $\mathbf{D}^r = \{\mathbf{d}^{r_i}\}_{i=1, \dots, K}$



# Experiments

$H_d^I$  into  $d^s$  and  $D^r = \{d^{r_i}\}_{i=1, \dots, K}$



$$M^r = \{M^{r_i}\}_{i=1, \dots, K}$$

$$M^{r_i} = \{m_j^{r_i}\}_{j=1, \dots, L_{r_i}}$$

$$M^s = \{m_i^s\}_{i=1, \dots, L_x}$$

$$\mathbf{h}_t = \text{GRU}(\mathbf{e}^w(y_{t-1}), \mathbf{h}_{t-1})$$

$$\mathbf{c}_t = \text{hier\_attn}(\mathbf{h}_t, M^s, M^r, D^r) \quad (3)$$

$$\tilde{\mathbf{h}}_t = \tanh(W_c[\mathbf{c}_t; \mathbf{h}_t]),$$

hier\_attn

$$\mathbf{c}_t^s = \sum_{i=1}^{L_x} a_{t,i}^s \mathbf{m}_i^s; \mathbf{c}_t^r = \sum_{i=1}^K \sum_{j=1}^{L_{r_i}} a_{t,i}^r a_{t,j}^{r_i} \mathbf{m}_j^{r_i} \quad (4)$$

$$\mathbf{c}_t = g_{ref} \cdot \mathbf{c}_t^s + (1 - g_{ref}) \cdot \mathbf{c}_t^r,$$

where  $\mathbf{a}_t^s$  is a word-level attention distribution over words from the source document using  $M^s$

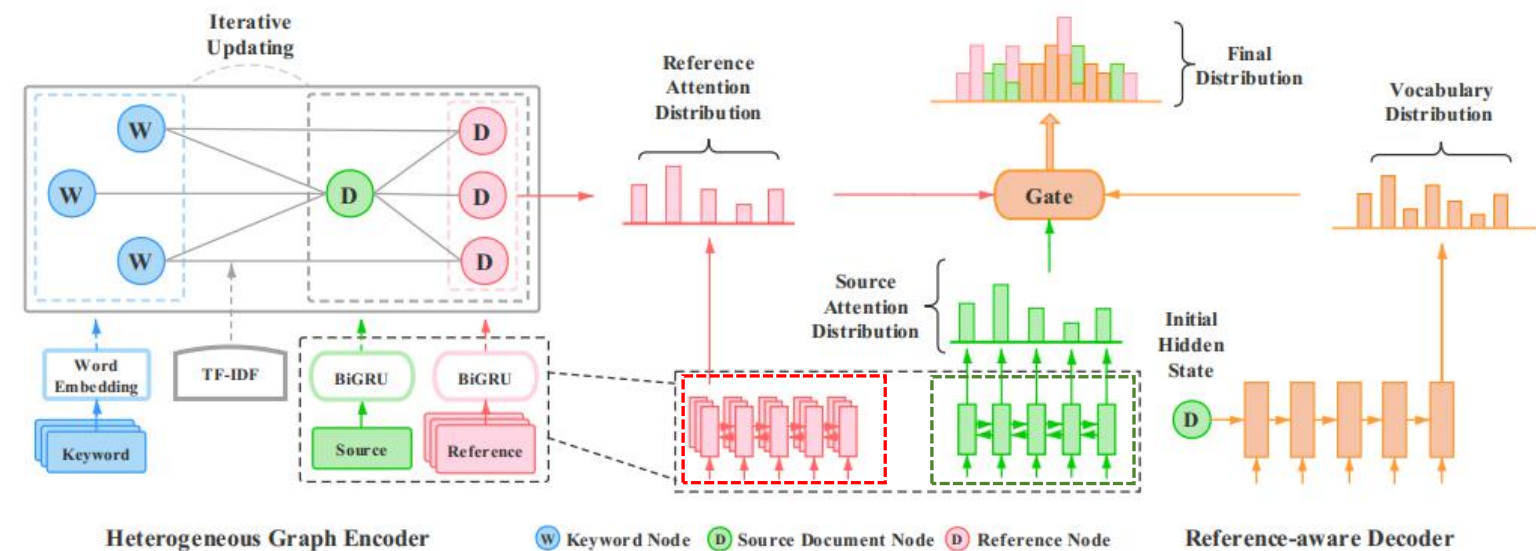
$\mathbf{a}_t^r$  is an attention distribution over references using  $D^r$ ,

$\mathbf{a}_t^{r_i}$  is a word-level attention distribution over words from  $i$ -th reference using  $M^{r_i}$ .

$$g_{ref} = \text{sigmoid}(\mathbf{w}_{ref}[\mathbf{c}_t^s; \mathbf{c}_t^r])$$

# Method

$H_d^I$  into  $d^s$  and  $D^r = \{d^{r_i}\}_{i=1,\dots,K}$



$$M^r = \{M^{r_i}\}_{i=1,\dots,K}$$

$$M^{r_i} = \{m_j^{r_i}\}_{j=1,\dots,L_{r_i}}$$

$$M^s = \{m_i^s\}_{i=1,\dots,L_x}$$

hierarchical copy mechanism

$$P_{V'}(y_t) = p_1 P_V(y_t) + p_2 P_{V_x}(y_t) + p_3 P_{V_{x^r}}(y_t)$$

$$P_V(y_t) = \text{softmax}(\text{MLP}([\mathbf{h}_t; \tilde{\mathbf{h}}_t]))$$

$$P_{V_x}(y_t) = \sum_{i:x_i=y_t} a_{t,i}^s$$

$$P_{V_{x^r}}(y_t) = \sum_i \sum_{j:x_j^r=y_t} a_{t,j}^{r_i}$$

$$\mathbf{p} = \text{softmax}(\mathbf{W}_p[\mathbf{h}_t; \mathbf{h}_t; \mathbf{e}^w(y_{t-1})]) \in \mathbb{R}^3$$

loss

$$\mathcal{L}_{\text{ONE2ONE}}(\theta) = - \sum_{i=1}^{|\mathcal{Y}|} \sum_{t=1}^{L_{y_i}} \log P_{V'}(y_{i,t} | \mathbf{y}_{i,1:t-1}, \mathbf{x}; \theta)$$

$$\mathcal{L}_{\text{ONE2SEQ}}(\theta) = - \sum_{t=1}^{L_{y^*}} \log P_{V'}(y_t^* | \mathbf{y}_{1:t-1}^*, \mathbf{x}; \theta)$$



# Experiments

| Model                           | NUS                       |                           |                           |                           | SemEval                   |                           |                           |                           | KP20k                     |                    |                           |                           |
|---------------------------------|---------------------------|---------------------------|---------------------------|---------------------------|---------------------------|---------------------------|---------------------------|---------------------------|---------------------------|--------------------|---------------------------|---------------------------|
|                                 | Present                   |                           | Absent                    |                           | Present                   |                           | Absent                    |                           | Present                   |                    | Absent                    |                           |
|                                 | $F1@5$                    | $F1@10$                   | $R@10$                    | $R@50$                    | $F1@5$                    | $F1@10$                   | $R@10$                    | $R@50$                    | $F1@5$                    | $F1@10$            | $R@10$                    | $R@50$                    |
| CopyRNN (Meng et al., 2017)     | 0.311                     | 0.266                     | 0.058                     | 0.116                     | 0.293                     | 0.304                     | 0.043                     | 0.067                     | 0.333                     | 0.262              | 0.125                     | 0.211                     |
| CorrRNN (Chen et al., 2018)     | 0.318                     | 0.278                     | 0.059                     | -                         | 0.320                     | 0.320                     | 0.041                     | -                         | -                         | -                  | -                         | -                         |
| TG-Net (Chen et al., 2019b)     | 0.349                     | 0.295                     | 0.075                     | 0.137                     | 0.318                     | 0.322                     | 0.045                     | 0.076                     | 0.372                     | 0.315              | 0.156                     | 0.268                     |
| KG-KE-KR-M (Chen et al., 2019a) | 0.344                     | 0.287                     | 0.123                     | <b>0.193</b>              | 0.329                     | 0.327                     | 0.049                     | 0.090                     | 0.400                     | <b>0.327</b>       | 0.177                     | 0.278                     |
| CopyRNN-GATER (Ours)            | <b>0.374</b> <sub>4</sub> | <b>0.304</b> <sub>4</sub> | <b>0.126</b> <sub>3</sub> | <b>0.193</b> <sub>2</sub> | <b>0.366</b> <sub>3</sub> | <b>0.340</b> <sub>4</sub> | <b>0.056</b> <sub>1</sub> | <b>0.092</b> <sub>2</sub> | <b>0.402</b> <sub>1</sub> | 0.324 <sub>1</sub> | <b>0.186</b> <sub>0</sub> | <b>0.285</b> <sub>1</sub> |

Table 1: Keyphrase prediction results of all the models trained under ONE2ONE paradigm. The best results are bold. The subscript are corresponding standard deviation (e.g., 0.285<sub>1</sub> means 0.285±0.001).

# Experiments

| Model                          | NUS                |                          |                          |                          | SemEval                  |                          |                          |                          | KP20k              |                          |                          |                          |
|--------------------------------|--------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------|--------------------------|--------------------------|--------------------------|
|                                | Present            |                          | Absent                   |                          | Present                  |                          | Absent                   |                          | Present            |                          | Absent                   |                          |
|                                | <i>F1@5</i>        | <i>F1@M</i>              | <i>F1@5</i>              | <i>F1@M</i>              | <i>F1@5</i>              | <i>F1@M</i>              | <i>F1@5</i>              | <i>F1@M</i>              | <i>F1@5</i>        | <i>F1@M</i>              | <i>F1@5</i>              | <i>F1@M</i>              |
| catSeq (Yuan et al., 2020)     | 0.323              | 0.397                    | 0.016                    | 0.028                    | 0.242                    | 0.283                    | 0.020                    | 0.028                    | 0.291              | 0.367                    | 0.015                    | 0.032                    |
| catSeqD (Yuan et al., 2020)    | 0.321              | 0.394                    | 0.014                    | 0.024                    | 0.233                    | 0.274                    | 0.016                    | 0.024                    | 0.285              | 0.363                    | 0.015                    | 0.031                    |
| catSeqCorr (Chan et al., 2019) | 0.319              | 0.390                    | 0.014                    | 0.024                    | 0.246                    | 0.290                    | 0.018                    | 0.026                    | 0.289              | 0.365                    | 0.015                    | 0.032                    |
| catSeqTG (Chan et al., 2019)   | 0.325              | 0.393                    | 0.011                    | 0.018                    | 0.246                    | 0.290                    | 0.019                    | 0.027                    | 0.292              | 0.366                    | 0.015                    | 0.032                    |
| SenSeNet (Luo et al., 2020)    | <b>0.348</b>       | 0.403                    | 0.018                    | 0.032                    | 0.255                    | 0.299                    | 0.024                    | 0.032                    | <b>0.296</b>       | 0.370                    | 0.017                    | 0.036                    |
| catSeq-GATER (Ours)            | 0.337 <sub>4</sub> | <b>0.418<sub>4</sub></b> | <b>0.033<sub>3</sub></b> | <b>0.054<sub>4</sub></b> | <b>0.257<sub>3</sub></b> | <b>0.309<sub>4</sub></b> | <b>0.026<sub>4</sub></b> | <b>0.035<sub>5</sub></b> | 0.295 <sub>2</sub> | <b>0.384<sub>1</sub></b> | <b>0.030<sub>1</sub></b> | <b>0.060<sub>2</sub></b> |

Table 2: Keyphrase prediction results of all the models trained under ONE2SEQ paradigm. The best results are bold. The subscript are corresponding standard deviation (e.g., 0.060<sub>2</sub> means 0.060±0.002).

# Experiments

| Model                              | Present      |              | Absent       |              |
|------------------------------------|--------------|--------------|--------------|--------------|
|                                    | $F_1@5$      | $F_1@M$      | $F_1@5$      | $F_1@M$      |
| catSeq-GATER                       | <b>0.295</b> | <b>0.384</b> | <b>0.030</b> | <b>0.060</b> |
| <i>Input Reference</i>             |              |              |              |              |
| - retrieved documents              | 0.293        | 0.377        | 0.026        | 0.052        |
| - retrieved keyphrases             | 0.291        | 0.369        | 0.018        | 0.037        |
| - both                             | 0.291        | 0.367        | 0.015        | 0.032        |
| <i>Heterogeneous Graph Encoder</i> |              |              |              |              |
| - $d2d$ edge                       | 0.294        | 0.379        | 0.024        | 0.049        |
| - $w2d$ edge                       | 0.294        | 0.379        | 0.026        | 0.052        |
| - both                             | 0.293        | 0.371        | 0.020        | 0.041        |
| <i>Reference-aware Decoder</i>     |              |              |              |              |
| - hierarchical copy                | 0.293        | 0.373        | 0.022        | 0.042        |
| - hierarchical attention           | 0.291        | 0.368        | 0.018        | 0.036        |

Table 3: Ablation study of catSeq-GATER on **KP20k** dataset. All references are ignored in graph encoder when removing  $d2d$  edge and the heterogeneous graph becomes homogeneous graph when removing  $w2d$  edge.



# Experiments

| Model      | Present      |              | Absent       |              |
|------------|--------------|--------------|--------------|--------------|
|            | $F1@5$       | $F1@M$       | $F1@5$       | $F1@M$       |
| catSeqD    | 0.285        | 0.363        | 0.015        | 0.031        |
| + GATER    | <b>0.294</b> | <b>0.381</b> | <b>0.025</b> | <b>0.051</b> |
| catSeqCorr | 0.289        | 0.365        | 0.015        | 0.032        |
| + GATER    | <b>0.296</b> | <b>0.384</b> | <b>0.030</b> | <b>0.060</b> |
| catSeqTG   | 0.292        | 0.366        | 0.015        | 0.032        |
| + GATER    | <b>0.293</b> | <b>0.380</b> | <b>0.025</b> | <b>0.052</b> |

Table 4: Results of applying our GATER to other baseline models on **KP20k** test set. The best results are bold.

# Experiments

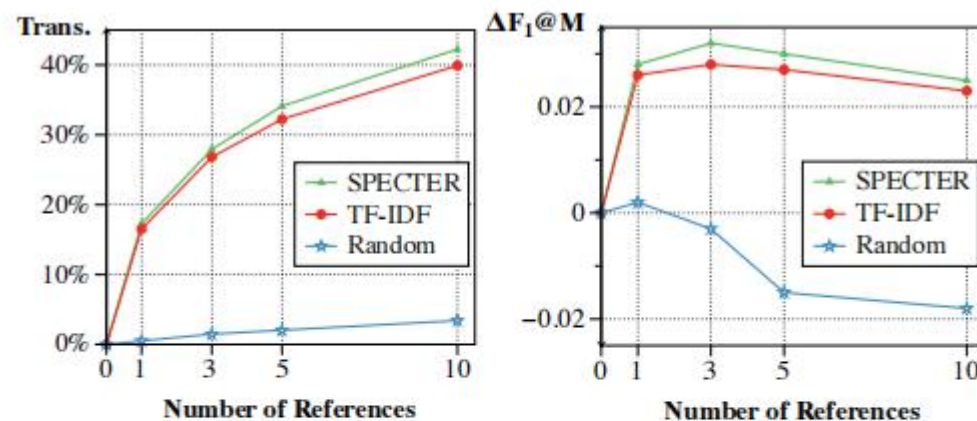


Figure 3: Transforming rate and  $\Delta F_1@M$  for absent keyphrases under different types of retrievers on **KP20k** dataset for catSeq-GATER. We study a random retriever, a sparse retriever based on TF-IDF and a dense retriever based on SPECTER.

# Experiments

---

**Document:** **natural convection** in porous annular domains **mimetic scheme** and **family of steady states**. **natural convection** of the incompressible fluid in the porous media based on the darcy hypothesis (lapwood convection) gives an intriguing branching off of one parameter family of steady patterns. this scenario may be suppressed in computations when governing equations are approximated by schemes which do not preserve the **cosymmetry** property ...

---

|                                        |                                                                                                                                                                                                                                          |
|----------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Present Keyphrases</b>              | natural convection; mimetic scheme; family of steady states; cosymmetry                                                                                                                                                                  |
| <b>CopyRNN</b>                         | <b>natural convection</b> ; porous media; polar coordinates; mimetic; annular porous domain; porous domain; finite difference; steady states; <b>mimetic scheme</b> ; ...                                                                |
| <b>KG-KE-KR-M</b>                      | <b>natural convection</b> ; porous media; <b>mimetic scheme</b> ; mimetic; polar coordinates; ...                                                                                                                                        |
| <b>CopyRNN-G<sub>ATER</sub> (Ours)</b> | <b>natural convection</b> ; porous media; <b>mimetic scheme</b> ; <b>cosymmetry</b> ; mimetic; darcy hypothesis; finite difference; polar coordinates; ...                                                                               |
| <b>Absent Keyphrases</b>               | darcy law; porous medium; finite difference method                                                                                                                                                                                       |
| <b>CopyRNN</b>                         | convective convection; annular porous media; mimetic method; <b>finite difference method</b> ; ...                                                                                                                                       |
| <b>KG-KE-KR-M</b>                      | cosymmetry convection; mimetic method; <b>darcy law</b> ; <b>convective patterns</b> ; lapwood property; annular porous media; <b>finite difference method</b> ; ...                                                                     |
| <b>CopyRNN-G<sub>ATER</sub> (Ours)</b> | <b>darcy law</b> ; <b>convective patterns</b> ; <b>porous medium</b> ; <b>multicomponent fluid</b> ; <b>finite difference method</b> ; family convection; cosymmetry convection; <b>staggered grids</b> ; <b>darcy formulation</b> ; ... |

---

Figure 4: Example of generated keyphrases by different models. The top 10 predictions are compared and some incorrect predictions are omitted for simplicity. The correct predictions are in bold blue and bold red for present and absent keyphrase, respectively. The absent predictions that appear in the references are highlighted in yellow, where only the keyphrases of retrieved documents are considered as references for KG-KE-KR-M.





# Thanks