



RESUME ANALYZER USING LLM

Chandhana H, Gowsika Sree S M, Asvitha VE

¹Studuent, Dept. of Artificial Intelligence and Data Science, Anna University, IN ^{x2}Studuent, Dept. of Artificial Intelligence and Data Science, Anna University, IN ³Studuent, Dept. of Artificial Intelligence and Data Science, Anna University, IN

_____***_________***

Abstract - — This study introduces a resume analyzer using large language models (LLMs) to automate and enhance the classification of resumes in hiring databases. Traditional resume classification, typically done manually, is time consuming and prone to inconsistency, especially as online recruitment continues to grow. Conventional machine learning approaches to resume categorization face limitations due to sparse labeled training data, scalability issues, and variable data quality. Our proposed method leverages a graph multi headed attention network (MGAT) model in a domain adaptation framework, trained on structured job post data to classify unstructured resume data. By treating the job post dataset as the source domain and the resume dataset as the target domain, this approach enhances classification accuracy and reduces reliance on extensive resume training data. The MGAT-based solution also addresses the challenges posed by long, variably formatted resumes, improving both classification efficiency and semantic alignment between resumes and job postings. This work represents a pioneering application of graph neural networks in the domain adaptation context for resume classification, presenting a scalable and effective solution for recruitment automation.

Key Words— Resume Classification, Large Language Models (LLMs), Recruitment Automation, Semantic matching, Machine Learning, Graph Multi-Headed Attention Network(MGAT)

1.INTRODUCTION

Classifying resumes is crucial for maintaining the hiring database and expediting the hiring procedure. Resume categorization aids the process by restricting the resume search range inside its category by classifying resumes according to occupation kinds. As a standard procedure, recruiters classify resumes according to their education, experiences, and other details after receiving them. They keep the resumes in databases. Screening resumes is generally a tedious and time-consuming process. An additional issue is that recruiters often classify resumes according to the industries of the applicants, even though some resumes may encompass multiple industries. The manual sorting of resumes has faced challenges related to database rearrangement and overwhelming data volume as online recruitment has grown. Supervised machine learning-based algorithms have demonstrated the effectiveness of machine learning applications in recent resume categorization studies. Nevertheless, labeled training data is absent from the supervised learning models [1]. Additionally, data labeling is done by hand, which takes a lot of time and has issues with scalability © 2024, IRJEdT

and quality stability [2]. Classification becomes considerably more challenging for resumes that primarily lack labeled data [3,4]. The effectiveness of the supervised model was limited because the majority of resume classification experiments utilized tiny training and testing

datasets [5, 6]. Resumes and job postings are matched for job recommendations in the recruitment process; this keyword based method is predicated on the idea that the contents of resumes and job postings are semantically comparable. The applicant's biographical information, employment history, educational background, and additional qualifications, such projects and certificates, are among the items that may overlap between a resume and a job posting. Therefore, it is possible to classify resumes using transfer learning on the features listed in job post data. Conversely, job postings are more structured to outline requirements for possible applicants, however resume styles differ slightly, making it more difficult to classify resume data. In addition to the described lack of resume data, another issue is that resumes are typically long documents. Due to their significantly increased memory and time requirements, which result in a model that is not stable, the existing contextual embedding techniques, such as BERT-based approaches, are not appropriate for handling lengthy texts [7-9]. Even though a resume's basic format includes sections like education, employment history, and personal information, the representations vary depending on the profession. While curriculum vitae in the software business may contain a project part, resumes in the creative industry should highlight personal qualities. Furthermore, the parts are flexible and can be tailored to each individual.

This study suggests a graph multi-headed attention neural network (MGAT) based on a domain adaptation technique to mitigate the issues raised. The MGAT model is trained on the labeled job post dataset before being used to categorize the resumes. As a result, the resume dataset is the target domain and the job post dataset is the source domain. These terms were used interchangeably in this investigation. The suggested method sought to improve classification efficiency using the domain adaptation strategy while lowering the dependency on resuming training data. This study is the first to use the graph neural networks variation in the domain adaption strategy to enhance resume classification, thanks to a comprehensive search of the pertinent literature. This is





findings

revealed a 34%

how the remainder of the study is structured. The literature review and associated investigations are covered in Section

2. The problem formulation is shown in Section 3. In Section 4, the methodology, structure, and characteristics of the suggested strategy, while the information, analysis, and outcomes of the experiments and Section 5 presents the discussions. Lastly, Section 6 offers the findings and further investigation.

2. LITERATURE REVIEW

2.1 Resume Classification and Domain Adaptation

In the early resume classification method, the resume was first segmented using the hidden Markov model (HMM), and then the crucial information was extracted from the labelled segmentations using a support vector machine (SVM) [10]. The precision was impacted as data moved through two poorly related models. The unstructured resumes were transformed into a set of concepts by Zaroor et al. [11], who then used knowledge-based assisted classification and conceptual matching to classify 10,000 job postings and 2000 resumes. The occupational categories incorporated from the DICE1 and O*NET2 categorization methods were compared with the resumes and job postings. Only the terms in the predefined list created by NER were eliminated, despite the fact that additional proper names and personal information were cleaned. Another study used five machine learning models-Naïve Bayes, Multinomial Naïve Bayes, Linear SVC, Bernoulli Naïve Bayes, and Logistic Regression-to build a voting classifier that predicted the resume category [12]. Only tokens containing tags as a noun and proper nouns were selected for additional processing, according to the results. The majority of votes utilizing the individual classifier influenced the final selection. The time required to handle the massive volume of data presented another difficulty for their method, since each classifier had to train the model and generate predictions independently. In order to classify 2000 resumes, Nasser et al. [3] used convolutional neural networks (CNN) in conjunction with the Glove embedding model.

words, and punctuation were used to clean the dataset. Additionally, they classified resumes in particular areas into distinct groups using the hierarchy classification, with each level representing a binary classification of resumes that were either technology-related or nontechnical. On the resume dataset, El Mohadab et al. [13] used the decision tree classifier and vector model. The attention mechanism is crucial to natural language processing [14]. Only [15] has employed the attention mechanism for resume categorization thus far; they applied Bi-LSTM in conjunction with a bespoke attention layer and created word embeddings using the Word2Vec model. In conclusion, several studies faced the issues of tiny resume datasets and limited categories as a result of inadequate labelled data. Another subject included in the study of resume classification is the domain adaption technique. Sayfullina et al. [4] trained the CNN model on a dataset of 80,000 job postings and classified 523 resumes in terms of training data labelling and quality. The © 2024, IRJEdT

decline in the model's ability to forecast job post to resume categories. Additionally, each sample's text length was restricted to 100 words. In a similar vein, Ramraj et al. [6]used CNN to extract data characteristics for classification; nevertheless, their approach did not considerably increase classification efficiency.

2.2 Graph-based Neural Networks in Text Classification

Over the past ten years, graph-based neural networks have grown in popularity and found several uses in a variety of industries. To capture word relations in the many graph-based methods for corpus, text classification were put forth. In order to enhance the existing text classification techniques, Wang et al. [16] worked on entity characteristics using hierarchical graph learning. Prior research focused on graph representation learning techniques, like gated graph neural networks (GGNN) and graph convolutional networks (GCN), for node representation in order to attain semantic correlations in the corpus. Using GCN, Yao et al. [18] constructed the graph using the number of documents in the corpus and the unique words. As a result, these approaches only addressed global interaction. Another study by Huang et al. [19] took into account the word level when creating the graph, but they only looked at a single pair of words with a fixed edge between them. The incapacity to produce embeddings for invisible nodes was the constraint on graph creation for the entire corpus. However, the previously noted problem was resolved by using a single graph for every document. In order to capture the semantic relationship of word-to-word relations, Zhang et al. [20] used GGNN and the cooccurrence relation within a fixed-sized context frame. By using this technique, the model was able to predict the new node that was absent from the original training graph. In a similar vein, by creating a homogenous graph with word nodes and cooccurrence statistics, graph fusion networks

[21] also enhanced the reconstruction problem of transductive approaches. The co-occurrence statistic in this study was based on word pairings and point-wise mutual information to extract the features of joint and marginal probabilities, rather than the slide-window method. However, even though word cooccurrence varied depending on the context, the final document representation graph did not take edge weights and attention weights into consideration.

The effectiveness of the attention mechanism in handling text-based data has been demonstrated by the developments of modern NLP techniques, such as transformers and BERTbased approaches. In order to address the brief text categorization problem, Linmei et al. [22] took advantage of dual-level graph attention networks (GAT). To address sentiment classification, Gan and Tang [23] combined dependency parsing with GAT. After BERT (Huggingface) inserted the input text





four primary

was used to retain the word's context within the sentence. These methods, however, were incompatible with the issues with lengthy documents. In order to categorize the resumes, this work used the MGAT experimentally to carefully learn the latent features from the job post dataset. The suggested combination was verified using a sizable collection of real resumes.

The scope of the proposed work involves creating a robust

ISL recognition system using Convolutional Neural Networks (CNNs). This system is trained on a diverse dataset that captures various ISL hand gestures under different conditions, enabling accurate recognition and interpretation.

METHODOLOGY

Resumes and job postings were typically pre-processed independently. To create the matrix A of word characteristics and their co-occurrence, the generated unique words were merged to create a set of unique words, which were then embedded via Glove embedding. The GAT classifier was then trained using the adjacency matrix of training graphs that was taken out of matrix A. Without retraining the classifier, resumes in the target domain were subsequently classified using the taught model.

2.3 Data Pre-Processing and Word Representation

The résumé and job post's raw datasets were preprocessed independently using standard methods, which included removing non-ASCII strings, common stop words, punctuation, and anomalous characters fig 3.1. Proper names, addresses, and contact details are examples of personal information that does not significantly improve resume categorization performance but cannot be removed using regex patterns or conventional preprocessing stop word removal. As a result, the POS tagging model was used to preprocess the resume noises. POS tagging assigned a proper POS to each token on the resume. The only tokens that were experimentally distilled were those that had the tags "NOUN," "ADJ," "VERB," and "ADV," which stand for noun, verb, adjective, and adverb, respectively. For instance, the token "Williamson," which is a person's name, was marked as "PROPN" (a pronoun) in Figure 3. Additionally, the misspelled and irregular terms were eliminated. Thus, the example statement was shortened to "software development consultant system architect." All of these basic keywords, along with the remaining general text, can offer distinguishing features for the classification process. The pre-trained Glove embedding was chosen for the word representation selection in order to improve graph creation and create a reliable model for a lot of data. First, using the transformer approach for embedding would take a lot of time because resumes and job postings are long text documents. Second, the corpus itself is discrete data since the preprocessing method keeps the POS words. Last but not least, the Glove embedding theory applies to the weights of edges in graphs created using the co-occurrence weights of the words. Since hi is a feature of the unique word vi, each CV or job posting was transformed into a set of features $h = {h1,h2,...,hu},hi \in Rv*d$.

2.4 Domain Adaptation Approach

Semantic similarity between the resume and the job posting served as the foundation for the domain adaption approach's classification of resumes. Word characteristics and word correlation were kept in the shared adjacency matrix A, which was created by processing and assembling the important data from both datasets for this investigation. The target domain is typically solely utilized for evaluation; its content is respected but disregarded. In this suggested method, the shared representation was created by combining the knowledge of resume content with the information of job post material. In particular, Eq. expresses the entire amount of unique words, which equals the knowledge of all papers. The adjacency matrix $A \in Ru*u$, a shared representation for both domains, was then constructed using the set of unique words V embedded in a set of words featured with 300 dimensions. The word co- occurrence in slide window 3 displayed the correlation of words in the matrix. The attention weights from matrix A were passed down to the edge weights since the word co- occurrences weights were set as the edge weights between nodes. The suggested approach was created to closely monitor the most prevalent patterns in the source domain during the training phase, both with regard to the correlation between the keywords and their patterns.





embedding



3. EXPERIMENTS

In this section, the effectiveness of the suggested method was assessed through a case study using real datasets. Regarding the target domain, the labelled resume dataset from a headhunting firm was utilized. When conducting the evaluation, the target domain was excluded from the collection of occupational categories, which are listed in Table 1 and are based on the present set of enterprises. The job post datasets from public datasets (10,000 Data ScientistJob Postings from the USA — Dataset by Jobspikr, [27]; 30,000 Job Postings from SEEK Australia — Dataset by Prompt cloud, [28]; Jobs on Naukri.Com, [29]) were combined to form the source domain. For consistency, the job posts from the same categories were grouped to form a shared common set with the target domain

Fig 3.1 Flow

chart

3.1 Baseline Models Selection

The résumé and job post's raw datasets were preprocessed independently using standard methods, which included removing non-ASCII strings, common stop words, punctuation, and anomalous characters. Proper names, addresses, and contact details are examples of personal information that does not significantly improve resume categorization performance but cannot be removed using regex patterns or conventional preprocessing stop word removal. As a result, the POS tagging model was used to preprocess the resume noises. POS tagging assigned a proper POS to each token on the resume. The only tokens that were experimentally distilled were those that had the tags "NOUN," "ADJ," "VERB," and "ADV," which stand for noun, verb, adjective, and adverb, respectively. For instance, the token "Williamson," which is a person's name, was marked as "PROPN" (a pronoun) in Figure 3. Additionally, the misspelled and irregular terms were eliminated. Thus, the example statement was shortened to "software development consultant system architect." All of these basic keywords, along with the remaining general text, can offer distinguishing features for the classification process. The pre-trained Glove embedding was chosen for the word representation selection in order to improve graph creation and create a reliable model for a lot of data. First, using the transformer approach for embedding would take a lot of time because resumes and job postings are long text documents. Second, the corpus itself is discrete data since the preprocessing method keeps the four primary POS words. Last but not least, the Glove theory applies to the weights of edges in graphs created using the co-occurrence weights of the words. Since hi is a feature of the unique word vi, each CV or job posting was transformed into a set of features $h = {h1,h2,...,hu},hi \in Rv*d$

3.2 Experiment Setting

The input features for the suggested and baseline models were identical for training, validation, and testing data in order to maintain consistency. The training data was experimentally reorganized to three folds using the stratified k-folds cross-validation method with k = 3, a variant of Kfold cross-validation. This approach ensured that each fold was a good representation of the dataset because both datasets were unbalanced, maintaining the same percentage of observations to a given label in each





more

fold. . During the preprocessing phase, the POS tagging model assigned each token a suitable POS tag by applying the "en_core_web_lg" package created by spaCy (English: spaCy Models Documentation [34]). GloVe: Global Vectors for Word Representation [35] employed 300 dimensions in the pre- trained glove embedding. During the trials and fine-tuning processes, the relevant parameters were established. After that, the model was trained and validated using data from the source domain over ten epochs, with a batch size of 64 in the source domain for each fold. Two-thirds of the data was designated as training data, while one-third was designated as validation data. All of the data was used in the target domain testing to confirm how well the trained model performed using the domain adaption technique.

The dropout rate was set at $0.3\ \text{and}\ the\ learning\ rate at$

0.001. The goal of every parameter choice was to create a reliable and effective model. The deep graph library [36], the Python programming language, and the PyTorch backend API were used to create and train the models based on graph neural networks. Every experiment is conducted on a PC running Windows 10 that has a 12th generation Intel R core i7 20 processor running at 3.61 GHz with 32 GB of RAM. The average computational time with this specification was roughly two hours for the training and validation of the BiLSTM model and two and a half hours for graph-based models. 5.4. Results and discussion The performance of the suggested and baseline models is shown and examined in this subsection. The primary goal of the studies was to assess the trained model's effectiveness in the intended field. There was a distribution shift between the source and target domains, and the data were unbalanced. As a result, the F1 score was a crucial metric. All models underwent five iterations of the experiments, and the average outcomes were published. In both the source and target domains, the MGAT strategy fared better overall than the other models. For the target domain, the strategy utilizing the Bi-LSTM model obtained the lowest score, 0.68. The primary cause was the model's input feature characteristics. Although the Bi-LSTM model's input required an LSTM—PyTorch 1.11.0 documentation [37] input sequence, the input features in this method were a collection of distinct, one-of-a-kind words. There was a lot of noisy data in the target domain prior to the preprocessing step, necessitating a trade-off between text representation and text preprocessing. The context was not reserved and was not in sequence form once the stopwords and biased information were eliminated.Compared to other approaches, the GCN model's efficiency was more significant; its F1 score in the target domain was 0.75, but the performance in the source domain was the lowest.

This outcome demonstrated the effectiveness of graph convolutional networks and provided an explanation for their recent surge in popularity. GCN, on the other hand, assigns equal weight to each node. Since this study focused on text type data, the attention method would be more appropriate because some nodes should be © 2024, IRJEdT Volume

significant than others. In specifics, instead of computing the GCN explicitly, the GAT calculated the coefficient implicitly. The MGAT learns to reweight the propagation weight depending on the hidden features using a learnable attention function, rather than relying solely on the graph structure to guide the propagation. This function was calculated as the concatenation of the affine converted hidden features of two nodes and the SoftMax normalized inner product between a learnable vector. One of the two most popular attention mechanisms, the dot-product attention technique, was employed in the dotGAT model. On the other hand, the suggested strategy employed additive attention. In contrast to the MGAT, the dotGAT model's performance in the experiments tended to be unstable. The findings of the F1 score standard deviation on the

The findings of the F1 score standard deviation on the targetdomain are displayed in Table 3. High accuracy

variability was shown by the larger dotGAT standard deviation findings. Initially, the attention weights between a random word feature and the collection of other word features in the phrase were determined using the dot- product. Therefore, if the learnable weights are not randomly initialized correctly, the remaining process may become unstable. The graph multi-headed attention utilized in MGAT concurrently computes eight heads of attention, each with a unique learnable weight, in order to get over this restriction, particularly for lengthy documents. The new feature vector is then obtained by concatenating the features. Consequently, the usage of the MGAT decreased the instability in dot-product attention and produced the highest score by examining various facets of the features from the preceding layer.

3.3 Statistical Hypothesis

The performance of the suggested model is evaluated in this section against the comparative models on the target domain using the two-sample t-test. The difference between the proposed model's mean and the compared model is the hypothesis h1. As a result, the MGAT model performs better than the comparative model if hypothesis h1 is adopted. The outcomes of statistical hypothesis testing are shown in Table 6. Hypothesis h1 is accepted since all of the p-value scores are less than 0.05, indicating that the suggested model outperforms the other models.

3.4 POS Tagging Filter Analysis

The goal of POS tagging is to increase the model's performance by lowering the noise in user-generated data. To assess the effectiveness of the POS tagging in the resume dataset, more tests were carried out. The effect of POS tagging on improving the cover proportion of word embedding is seen in Table 4. The performance of the downstream tasks depends on the percentages of vocabulary and all text embedding since the pre-trained word vector embedded the words into word features. More information is retained the higher the proportion of embedding that is covered. Additionally, the percentage of





can be seen

embedding increased to 35.12% and the percentage of total text embedding increased from 52.97% to 98.05% with the addition of the POS filter. The most common bigrams in a banking and finance group example are examined in more detail. Although the bi-grams "New York," "Microsoft office," and "MS office" were most frequently used, they weren't the keywords that set the banking and financial group apart from the others. However, the top bi-grams in resumes that had the POS filter applied further were more specific to the banking and finance industry.

Furthermore, the use of POS tagging allowed for the selective distillation of the POS, allowing for flexible application to match the challenges and data characteristics. The percentage of nouns, verbs, adjectives, and adverbs in each group is displayed in Fig. 8a. Each group had a higher proportion of nouns than the other POSs, although adverbs made up a smaller percentage. The most common bi-grams in each POS in the banking and finance group, for instance, are displayed in Fig. 8b.

Lastly, in-depth tests were carried out to confirm how the POS choices affected the model's performance. The three methods are listed in Table 5: Tokens with tags of "NOUN," "ADJ," "VERB," and "ADV" are retained in MGAT (noun and verb), which is used in this study; MGAT (noun and verb) — which uses the pre-processed text without applying a POS filter; and MGAT (noun and verb) - which only filters the tokens with tags of "NOUN" and "VERB." Due to the low percentage of embedding cover, the raw text dataset was omitted. The increase in F1 score from 0.75 to 0.80 confirmed the efficacy of the POS tagging. The corpus had a low embedding cover proportion and a high percentage of noise information when the POS tagging was not used, which appears to have hurt the model's performance on the target domain. The case of the remaining noun and verb increased by 0.75 in the F1 score, as indicated by the data in Table 5. The number of significant traits decreased along with the substantial reduction in vocabulary. Additionally, cutting out a lot of words had a direct impact on how long many short resumes were.

3.5 Target Source Misprediction Analysis

Even though the suggested approach demonstrated effectiveness when the target domain's F1 score was 0.80, more research was done on the collection of misprediction resumes to obtain a deeper understanding. The goal domain is shown in Fig. 9, where both true and incorrect predictions are included for each category group. Series 2 in orange indicates the instances in which the model was unable to forecast the ground truth categories, whereas series 1 in blue indicates the resumes that the algorithm was able to correctly predict. The predetermined number on the horizontal axis displays the categories. All things considered, the primary cause of the target domain misprediction was the possibility that a single resume could fall under multiple closely related categories. As

from the graphic, the sales and technical groups had the most incorrect forecasts when compared to the orange columns. In these two categories, the trained model's prediction accuracy was only about 60% of the ground truth. In order to further understand the misclassified group, additional analysis was conducted on these two groups using the misprediction dataset. Figure 10 displays the sales group's performance, whereas Figure 11 displays the technical group's. the input features' similarity. The remaining misprediction may have resulted from the close relationship between sales group characteristics and the other categories, since candidates are presumed to possess sufficient expertise to sell a wide range of products from all industries. The model's input features include the distinct words found in the resume and job posting; as a result, during training, the model became confused and incorrectly predicted the resume category. The situation described above also applies to the technical group. The IT group made the majority of the poor choices. First, software and firmware engineers may have talents in common with hardware engineers. Electronic engineering, for instance, may fall within the software and firmware category. Second, both the IT and technical sectors value common talents like MATLAB coding, programming languages like C++, and soft skills. Field engineering and the industrial management group make up a minor fraction of the total. The maintenance engineer, the industrial management group, and certain technical positions involving automation, machine operation, and warehouse technical operations provide comparable examples. Similar abilities, such as AutoCAD drawing and simulation, are shared by the field group and the technical group. However, the model nearly effectively separated the technical group from the hospitality group and the banking and finance group; as it is, these groupings do not share any traits.

4.5 CONCLUSION AND FUTURE WORK

- The excessive collection of data has gotten increasingly difficult due to the rapid advancement of information and communication technology. The recruitment process can be streamlined and the classification process made more effective with an automated resume classification approach. In order to improve overall performance, this study used graph multiheaded attention networks in the domain adaption approach to investigate the problem of categorizing resume data for the training model in aclassification task. Initially, the MGAT model was trainedusing the job post data. The trained model was then used to forecast resume categories based on the semantic similarity of the relevant sections of resumes and job postings.
- POS tagging was subsequently used to extract the essential characteristics and significantly cut down on superfluous data in the resume dataset in order to increase the model's efficiency. The pre-processed unique terms and their co- occurrence weights were then used to create the document graphs. Graph multi-





Stat].

headed attention networks were utilized to construct the document embedding for classification by locally and attentively learning the node features in the document graphs.

With an F1 score of 80%, the testing findings demonstrate that the suggested strategy performed noticeably better on the target domain. This study's primary contribution was to suggest a unique method for classifying resumes without the requirement to label the resume data by utilizing graph multi-headed attention networks with domain adaptability. Lastly, by reducing the scope of the applicant search, the suggested approach successfully optimized the hiring process bylowering manual labour, bias, and domain knowledge constraints.

REFERENCES

- L. Sayfullina, E. Malmi, Y. Liao, A. Jung, Domain adaptation for resume classification using convolutional neural networks, 2017, ArXiv:1707.05576 [Cs].
- A. Jha, V. Rakesh, J. Chandrashekar, A. Samavedhi, C.K. Reddy, Supervised contrastive learning for interpretable long-form document matching, 2022, ArXiv:2108.09190arXiv.
- Huggingface/Transformers Transformers State-of-the-Art Machine Learning for Pytorch, Tensorflow, and JAX, 2022, Retrieved February 27, 2022, from https://github.com/huggingface/transformers.
- T.N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, 2017, ArXiv:1609.02907 [Cs,

[1] M. El Mohadab, B. Bouikhalene, S. Safi, Automatic CV processing for scientific research using data mining algorithm, J. King Univ. Comput. Inform. Sci. 32 (5) (2020) 561–567.

[2] F. Rousseau, E. Kiagias, M. Vazirgiannis, Text categorization as a graph classification problem, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), 2015, pp. 1702–1712.

[3] M.-T. Luong, H. Pham, C.D. Manning, Effective approaches to attention based neural machine translation, 2015, ArXiv:1508.04025 [Cs]