



# A Survey of Natural Language Generation

CHENHE DONG, Sun Yat-Sen University, China

YINGHUI LI, Tsinghua University, China

HAIFAN GONG, Sun Yat-Sen University, China

MIAOXIN CHEN and JUNXIN LI, Tsinghua University, China

YING SHEN, Sun Yat-Sen University, China

MIN YANG, Chinese Academy of Science, China

This article offers a comprehensive review of the research on Natural Language Generation (NLG) over the past two decades, especially in relation to data-to-text generation and text-to-text generation deep learning methods, as well as new applications of NLG technology. This survey aims to (a) give the latest synthesis of deep learning research on the NLG core tasks, as well as the architectures adopted in the field; (b) detail meticulously and comprehensively various NLG tasks and datasets, and draw attention to the challenges in NLG evaluation, focusing on different evaluation methods and their relationships; (c) highlight some future emphasis and relatively recent research issues that arise due to the increasing synergy between NLG and other artificial intelligence areas, such as computer vision, text, and computational creativity.

CCS Concepts: • **Computing methodologies** → **Natural language generation**;

Additional Key Words and Phrases: Natural language generation, data-to-text generation, text-to-text generation, deep learning, evaluation

## ACM Reference format:

Chenhe Dong, Yinghui Li, Haifan Gong, Miaoxin Chen, Junxin Li, Ying Shen, and Min Yang. 2022. A Survey of Natural Language Generation. *ACM Comput. Surv.* 55, 8, Article 173 (December 2022), 38 pages.

<https://doi.org/10.1145/3554727>

173

## 1 INTRODUCTION

This article surveys the current state-of-the-art in **Natural Language Generation (NLG)**, defined as the task of generating text from underlying non-linguistic representation of information [102]. NLG has been receiving more and more attention from researchers because of its extremely challenging and promising application prospects.

Chenhe Dong and Yinghui Li Equal contribution.

This work was supported in part by the 173 program No. 2021-JCJQ-JJ-0029, the Shenzhen General Research Project under Grant JCYJ20190808182805919 and in part by the National Natural Science Foundation of China under Grant 61602013.

Authors' addresses: C. Dong, H. Gong, and Y. Shen (corresponding author), Sun Yat-Sen University, Guangzhou, China, 510275; emails: {dongchh, gonghf}@mail2.sysu.edu.cn, sheny76@mail.sysu.edu.cn; Y. Li, M. Chen, and J. Li, Tsinghua University, Shenzhen, China, 518055; emails: {liyinghu20, cmx20, ljx20}@mails.tsinghua.edu.cn; M. Yang (corresponding author), Chinese Academy of Science, Shenzhen, China, 510100; email: min.yang@siat.ac.cn.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2022 Association for Computing Machinery.

0360-0300/2022/12-ART173 \$15.00

<https://doi.org/10.1145/3554727>

### 1.1 What Is Natural Language Generation?

**Natural Language Generation (NLG)** is the process of producing a natural language text to meet specified communicative goals. The texts that are generated may range from a single phrase given in answer to a question, through multi-sentence remarks and questions within a dialog, to full-page explanations.

In contrast with the organization of the **Natural Language Understanding (NLU)** process—which can follow the traditional stages of a linguistic analysis: morphology, syntax, semantics, pragmatics/discourse—the generation process has a fundamentally different character. Generation proceeds involve the content planning, determination, and realization from content to form, from intentions and perspectives to linearly arrayed words and syntactic markers. Coupled with its application, the situation, and the discourse, they provide the basis for making choices among the alternative wordings and constructions that the language provides, which is the primary effort in constructing a text deliberately [78]. With its opposite flow of information, one might assume that a generation process could be organized like an understanding process but with the stages in opposite order.

Both data-to-text generation and text-to-text generation are instances of NLG. Generating text from image is an application of data-to-text generation. A further text-to-text generation complication is dividing NLG tasks into three categories, i.e., text abbreviation, text expansion, text rewriting and reasoning. The text abbreviation task is formulated to condense information from long texts to short ones, typically including research on text summarization [8, 9, 20, 24, 53, 96, 116], question generation [5, 26, 44, 46, 66, 112, 122, 133, 134, 153, 158], and distractor generation [30, 61, 76, 89, 98, 103, 117, 118]. The text expansion tasks, such as short text expansion [6, 106, 113, 124] and topic-to-essay generation [27, 97, 135, 146, 152], generate complete sentences or even texts from some meaningful words by considering and adding elements such as conjunctions and prepositions to transform the input words into linguistically correct outputs. The goal of text rewriting and reasoning task is to rewrite the text into another style or applying reasoning methods to create responses. There are two sub-tasks: text style transfer [13, 28, 43, 63, 72, 80, 85, 95, 143, 157] and dialogue generation [4, 45, 57, 60, 74, 136, 145, 159]. The task of visual-based text generation targets at generate the explanation or summarization of the given image or video, involving the study of image caption [1, 71, 104, 129, 148, 149], video caption [22, 49, 52, 88, 114, 121, 128, 131, 140, 142, 160, 161], and visual storytelling [40, 55, 147].

### 1.2 Why a Survey on Natural Language Generation?

Here, we will explain the reasons and motivations why the natural language generation is worth reviewing and investigating.

Reiter et al. [102] provided the most classical survey of NLG. However, the field of NLG has changed drastically in the past 20 years, with the emergence of successful deep learning methods. For example, since 2014, various neural encoder-decoder models pioneered by **sequence-to-sequence (Seq2Seq)** have been proposed to achieve the goal by learning to map input text to output text. In addition, the evaluation of NLG output should start to receive systematic attention.

Since Reiter et al. [102] published their book, various other NLG overview texts have also appeared. Gatter et al. [31] introduce the core tasks, applications, and evaluation metrics of natural language generation. While useful, this survey is not highly timely and does not include the state-of-the-art research on the novel deep learning models such as graph neural networks. Perera et al. [92] cover some tasks or architectures of NLG. Santhanam et al. [108] review the NLG research progress of dialogue systems. Mogandala et al. [82] study the integration of vision and language in multimodal NLG, such as image dialogue and video storytelling. Otter et al. [86] conclude the

progress of deep learning for NLP, display some classic text generation methods, but barely discuss the research progress of NLG. Yu et al. [151] offer a survey of the knowledge-enhanced text generation methods.

The goal of the our survey is to present a highly timely overview of NLG developments from the aspect of data-to-text generation and text-to-text generation. Though NLG has been a part of AI and Natural language processing for long time, it has only recently begun to take full advantage of recent advances in data-driven, machine learning, and deep learning approaches. Therefore, this survey will focus on introducing the latest development and future directions of deep learning methods in field of NLG. This survey can have broad audiences, researchers and practitioners, in academia and industry.

### 1.3 Contribution of This Survey

In this article, we provide a thorough review of different natural language generation tasks as well as its corresponding datasets and methods. To summarize, this article presents an extensive survey of natural language generations with the following contributions:

- (1) To give an up-to-date synthesis of deep learning research on the core tasks in NLG, as well as the architectures adopted in the field;
- (2) To detail meticulously and comprehensively various NLG tasks and datasets and draw attention to the challenges in NLG evaluation, focusing on different evaluation methods and their relationships.
- (3) To highlight some future emphasis and relatively recent research issues that arise due to the increasing synergy between NLG and other artificial intelligence areas, such as computer vision, text, and computational creativity.

The rest of this survey is organized as follows: In Section 2, we introduce the general methods of NLG to give a comprehensive understanding. From Sections 3 to 6, we will give a comprehensive introduction to the four main areas of NLG from the perspectives of task, data, and methods. In Section 7, we present the important evaluation metrics used in various aforementioned NLG tasks. Besides, we propose some problems and challenges of NLG as well as several future research directions in Section 8. We conclude our survey in Section 9.

## 2 GENERAL METHODS OF NLG

In general, the task of **natural language generation (NLG)** targets at finding an optimal sequence  $y_{<T+1} = (y_1, y_2, \dots, y_T)$  that satisfies:

$$y_{<T+1} = \arg \max_{y_{<T+1} \in \mathcal{Y}} \log P_{\theta}(y_{<T+1}|x) = \arg \max_{y_{<T+1} \in \mathcal{Y}} \sum_{t=1}^T \log P_{\theta}(y_t|y_{<t}, x), \quad (1)$$

where  $T$  represents the number of tokens of the generated sequence,  $\mathcal{Y}$  represents a set containing all possible sequences, and  $P_{\theta}(y_t|y_{<t}, x)$  is the conditional probability of the next token  $y_t$  based on its previous tokens  $y_{<t} = (y_1, y_2, \dots, y_{t-1})$  and the source sequence  $x$  with model parameters  $\theta$ .

The general methods to deal with the tasks of NLG mainly contain: Recurrent Neural Network, Transformer, Attention Mechanism, Copy and Pointing Mechanisms, Generative Adversarial Network, Memory Network, Graph Neural Network, and Pre-trained Model.

### 2.1 Recurrent Neural Network

As proposed by Reference [123], the encoder of the **sequence-to-sequence (Seq2Seq)** framework is a **Recurrent Neural Network (RNN)**, it will traverse every token (word) of the input, the input

of each time is the hidden state and input of the previous time, and then there will be an output and a new hidden state. The new hidden state will be used as the input hidden state of the next time. We usually only keep the hidden state of the last time, which encodes the semantics of the whole sentence. After the encoder processing, the last hidden state will be regarded as the initial hidden state of the decoder. The Decoder is also an RNN, which outputs one word at a time. The input of each time is the hidden state of the previous time and the output of the previous time. The initial hidden state is the last hidden state of the encoder, and the input is special. Then, RNN is used to calculate the new hidden state and output the first word, and then the new hidden state and the first word are used to calculate the second word. Until EOS is encountered and the output is finished. A standard RNN computes a sequence of outputs  $(y_1, \dots, y_T)$  given a sequence of inputs  $(x_1, \dots, x_T)$  by iterating the following equation:

$$y_t = W^{yh} \cdot h_t = W^{yh} \cdot \sigma \left( W^{hx} x_t + W^{hh} h_{t-1} \right), \quad (2)$$

where  $\sigma$  is activation function,  $W^{hx}$ ,  $W^{hh}$ ,  $W^{yh}$  are learnable parameters, and  $h_t$  is the hidden state at  $t$ th timestep.

## 2.2 Transformer

Transformer [126] is based on the encoder-decoder framework, and both encoder and decoder are composed of stacked identified layers. The encoder is used to map an input sequence of symbol representations to another sequence of continuous representations, and then the decoder auto-regressively generates an output sequence based on its previously generated symbols and the continuous representations from encoder. In the encoder, each layer contains two sub-layers, which are **multi-head self-attention mechanism (MultiHeadAttn)** and position-wise fully connected **feed-forward network (FFN)**, respectively. The multi-head self-attention can be formulated by:

$$\text{MultiHeadAttn}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O, \quad (3)$$

$$\text{head}_i = \text{Attention} \left( QW_i^Q, KW_i^K, VW_i^V \right), \quad (4)$$

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^\top}{\sqrt{d_k}} \right) V, \quad (5)$$

where  $Q, K, V$  are the query, key, and value matrices,  $d_k$  is the dimension of queries and keys. And the feed-forward network can be formulated by:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2. \quad (6)$$

Each of the stacked layers is surrounded by a residual connection followed by layer normalization. While in the decoder, each layer contains three sub-layers with an additional multi-head attention over the encoder's output. The self-attention in decoder is masked to prevent attending to subsequent tokens. In addition, to inject the position information, the position encodings are added to the input embeddings as formulated by:

$$PE_{pos, 2i} = \sin(pos/10000^{2i/d_{model}}), \quad PE_{pos, 2i+1} = \cos(pos/10000^{2i/d_{model}}), \quad (7)$$

where  $pos$  is the position and  $i$  is the dimension.

The most commonly used loss function is the conditional language modeling loss, which can be formulated as:

$$L = \sum_{t=1}^T \log P_\theta(y_t | y_{<t}, x), \quad (8)$$

where  $\log P_\theta(y_t|y_{<t}, x)$  is the log-likelihood of the  $t$ th generated token conditioned on the previously generated sequence  $y_{<t}$  and the source sequence  $x$ .

### 2.3 Attention Mechanism

Attention mechanism [15] is used to perform a mapping from a query to a series of key value pairs. There are three steps in the calculation of attention. The first step is to calculate the similarity between query and each key to get the weight. The common similarity functions are dot product, splicing, perceptron, and so on; the second step is to normalize these weights by using a softmax function; the last step is to sum the weight and the corresponding key value to get the final attention. The tokens in an article are first fed into the encoder to generate a sequence of encoder hidden states  $h_i$ , then the decoder receives the word embedding of previous step and derives the decoder state  $s_t$  at each step  $t$ . After that, the attentive context vector  $h_t^*$  can be obtained based on the attention distribution  $a^t$  as:

$$h_t^* = \sum_i a_i^t h_i, \text{ where } a_i^t = \text{softmax}(e_i^t), e_i^t = v^\top \tanh(W_h h_i + W_s s_t + b_{attn}), \quad (9)$$

in which  $v, W_h, W_s, b_{attn}$  are learnable parameters. Finally, the predicted vocabulary distribution  $P_{vocab}$  is obtained by:

$$P_{vocab} = \text{softmax}(V'(V[s_t, h_t^*] + b) + b'), \quad (10)$$

where  $V, V', b, b'$  are learnable parameters.

### 2.4 Copy and Pointing Mechanisms

The copy and pointing mechanisms proposed by PGN [110] is widely used in abstractive summarization, which is designed for alleviating the problem of inaccurate reproduced factual details via a pointer, dealing with out-of-vocabulary words and repetition via a generator and a coverage mechanism, respectively. In the pointer and generator, a generation probability  $p_{gen}$  based on the context vector  $h_t^*$ , decoder state  $s_t$ , and decoder input  $x_t$  at step  $t$  is calculated by:

$$p_{gen} = \sigma(w_h^\top h_t^* + w_s^\top s_t + w_x^\top x_t + b_{ptr}), \quad (11)$$

which serves as a soft switch to choose between generating a word based on the vocabulary probability or copying a word from the source document based on the attention distribution. The probability distribution over the extended vocabulary  $P(w)$  can be formulated as:

$$P(w) = p_{gen} P_{vocab}(w) + (1 - p_{gen}) \sum_{i: w_i = w} a_i^t. \quad (12)$$

And in the coverage mechanism, a coverage vector  $c^t$  at step  $t$  is calculated by the sum of attention distributions over previous decoder timesteps as  $c^t = \sum_{t'=0}^{t-1} a^{t'}$ , which is then added to the attention mechanism and the primary loss function as:

$$e_i^t = v^\top \tanh(W_h h_i + W_s s_t + w_c c_i^t + b_{attn}), \quad (13)$$

$$loss_t = -\log P(w_t^*) + \lambda \sum_i \min(a_i^t, c_i^t), \quad (14)$$

where  $w_t^*$  is the target word at step  $t$ .

### 2.5 Generative Adversarial Network

The **generative adversarial network (GAN)** [35] is a framework that uses an adversarial training process to estimate generative models. This framework can be regarded as a minimax two-player

game containing a generative model and a discriminative model, and these two models are simultaneously trained. The **generator (G)** captures the data distribution and tries to produce fake samples, and the **discriminator (D)** attempts to determine whether the samples come from the model distribution or data distribution. In detail, G is trained to maximize the probability identified by D for the sample coming from the data rather than G, while D is trained to maximize the probability of assigning the correct label to training samples and samples generated by G. The training process continues until the counterfeits are indistinguishable from the genuine articles. The training objective with value function  $V(D, G)$  can be formulated as:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \quad (15)$$

## 2.6 Memory Network

The end-to-end memory network [120] is based on the recurrent neural network. Before outputting a symbol, the recurrence reads from a possibly large external memory multiple times. In this framework, a discrete set of input sentences is written to the memory up to a fixed size, and a continuous representation for each memory sentence and the query (i.e., the question) is calculated and is then processed through multiple hops to output the answer. Specifically, for the single hop operation, the matching probability between each query and memory is first computed by the inner product followed by a softmax in the embedding space. The input set  $x$  and the query  $q$  are first embedded to obtain the memory vectors  $m$  and  $u$ , respectively, and the match probability  $p_i$  of the  $i$ th input sentence is calculated as:

$$p_i = \text{softmax}(u^\top m_i). \quad (16)$$

Then the output memory representation  $o$  is obtained by a weighted sum over another embedded memory sentence  $c$  based on the matching probability as:

$$o = \sum_i p_i c_i. \quad (17)$$

Finally, the final prediction  $\hat{a}$  can be obtained through a weight matrix and a softmax over the sum of output memory vector  $o$  and input embedding  $u$  as:

$$\hat{a} = \text{softmax}(W(o + u)). \quad (18)$$

To handle multiple hop operations, the memory layers are repeatedly stacked, and the input of each layer is the sum of the output memory vector and the input from its previous layer.

## 2.7 Graph Neural Network

The **graph neural network (GNN)** [109] is used to process the data in graph form, which can capture the dependency information between nodes of a graph via message passing. Compared with CNN and RNN, GNN can propagate on each node, respectively, and is able to ignore the input orders of nodes, which is more computationally efficient. The representation of each node in a graph is iteratively updated by aggregating information from its neighboring nodes and edges. So far, many propagation strategies have been proposed, such as convolution [48], RNN-based gate mechanism [59], and attention mechanism [127]. Meanwhile, many works attempt to improve the training method, such as sampling-based training [38] and unsupervised training [47]. Take the graph convolution network [48] as an example: It follows the layer-wise propagation rule as:

$$H^{(l+1)} = \sigma \left( \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)} \right), \quad (19)$$



where  $\sigma(\cdot)$  refers to an activation function,  $\tilde{A} = A + I_N$  is the adjacency matrix of the graph (including the identity matrix  $I_N$ ),  $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$  is the degree of node  $i$ ,  $W^{(l)}$  is a trainable weight matrix in the  $l$ th layer, and  $H^{(l)}$  is the matrix of node feature vectors in the  $l$ th layer.

## 2.8 Pre-trained Model

The pre-trained models can be divided into two categories: non-contextual and contextual. The non-contextual pre-trained models (e.g., Word2vec [79], GloVe [91]) can learn high quality of word and phrase representations and are widely used to initialize the embeddings and improve the model performance for generation tasks. However, since the non-contextual embeddings are static, it is unable to handle context-dependent words and out-of-vocabulary words.

To overcome these problems, contextual pre-trained models are proposed, which can dynamically change the word embeddings when the word appears in different sentences. Traditional contextual models mainly focus on the tasks of **natural language understanding (NLU)**, which are based on the architectures of LSTM (e.g., ELMo [93]) or Transformer encoder (e.g., BERT [18]).

To tackle the more challenging tasks of **natural language generation (NLG)**, the architecture of Transformer decoder is widely adopted, and the pre-training objective is further modified to adapt to the NLG task. For example, GPT [7] uses the standard language modeling as the pre-training objective. T5 [99] designs a unified text-to-text framework for both NLU and NLG. During pre-training, each corrupted token span in the input sequence is replaced with a sentinel token, and the output sequence is composed of the dropped-out spans. MASS [116] predicts the sentence fragment with input of masked sequence during pre-training. UniLM [23] propose a unified pre-training framework for both NLU and NLG tasks, which is based on a shared Transformer model with three different types of self-attention masks to switch among different tasks, including unidirectional, bidirectional, and sequence-to-sequence. BART [53] designs a reconstruction objective to restore corrupted documents, where five types of transformation strategies are proposed, including token masking, token deletion, text infilling, sentence permutation, and document rotation. PLATO [4] proposes a pre-training framework targeting the dialogue generation tasks with two reciprocal pre-training tasks, i.e., response generation and latent act recognition. The latent discrete variables are also introduced to solve the one-to-many mapping problem. ERNIE-GEN [141] is a multi-flow Seq2Seq pre-training model, which pays attention to the exposure bias problem of pre-training models on downstream NLG tasks such as question generation, and aims to make the NLG models generate more human-like and fluent texts. OFA [132] presents a unified multimodal pre-training framework to unify various vision and language tasks, such as NLU, NLG, and image classification. Transformer is applied as the backbone architecture, and the sparse coding and a unified vocabulary are utilized to represent the images and linguistic words.

## 3 TEXT ABBREVIATION

### 3.1 Task

The goal of text abbreviation is to distill key information from long texts to short ones, which consists of three subtopics: text summarization, question generation, and distractor generation.

There are two kinds of methods for text abbreviation: extractive and abstractive methods. Since the abstractive approach is more flexible and can create more human-like sentences than the extractive approach, it has been paid more and more attention in recent years and is our main focus in this article. Text summarization is the process of generating entirely new phrases and sentences to capture the meaning of the source document. Question generation concentrates on automatically generating questions from a given sentence or paragraph. Distractor generation is the automatic generation of adequate distractors for a given question answer pair generated from a given article

Table 1. Natural Language Generation Models for Text Abbreviation

Task	Model	Description
Text Summarization	MASS [116]	Transformer
	BART [53]	Transformer + Multi-task Learning
	PEGASUS [154]	Transformer + Multi-task Learning
	RCT [8]	Transformer + RNN + CNN + Long-term Dependency
	ProphetNet [96]	Transformer + Long-term Dependency
	En-Semantic-Model [20]	RNN + Long-term Dependency
	Post-Editing Factual Error Corrector [9]	Transformer + Factual Consistency
	SpanFact [24]	Transformer + Factual Consistency
Question Generation	Key-Phrase-based Question Generator [26]	Keyphrase + Template
	Dynamic Mathematical Question Generator [5]	Constraint Handling Rules
	KB-based Factoid Question Generator [112]	RNN
	Teacher Forcing and RL Based Question Generator [153]	RNN + RL
	Paragraph-level Question Generator [158]	RNN
	Answer-Position-aware Question Generator [122]	RNN + Answer-focused
	ASs2s [46]	RNN + Answer-focused
	NQG-MP [134]	RNN + Multi-task Learning
	Paraphrase Enhanced Question Generator [44]	RNN + Multi-task Learning
	CGC-QG [66]	RNN + Multi-task Learning + GNN
	PathQG [133]	RNN + Multi-task Learning + KG
	UniLM [23]	Transformer + Multi-task Learning
	ERNIE-GEN [141]	Transformer + Multi-task Learning
Distractor Generation	Educational Ontology Distractor Generator [118]	Ontology + Embedding
	Learning to Rank Based Distractor Generator [61]	Embedding + Ranking + GAN + RL
	BERT-based Distractor Generation [16]	BERT + Multi-task Learning
	Hierarchical Dual-attention Distractor Generator [30]	RNN
	EDGE [98]	RNN + Answer Interaction
	HMD-Net [76]	Transformer + RNN + Answer Interaction
	Code Compression Distractor Generator [117]	Abstract Syntax Tree
	CSG-DS [103]	LDA + KB + Ranking
	Named Entity Distractor Generator [89]	Tree + Clustering

to form an adequate multiple-choice question. We summarize the most representative methods for each subtask in Table 1.

### 3.2 Data

**3.2.1 Text Summarization.** There are mainly four datasets in the field of text summarization as shown below.

**CNN/DailyMail.** The CNN/DailyMail dataset [39] is a large-scale reading comprehension dataset. This dataset contains 93K and 220K articles collected from the CNN and Daily Mail websites, respectively, where each article has its matching abstractive summary.

**NYT.** The **New York Times (NYT)** dataset [90, 107] contains large amount of articles written and published by the *New York Times* between 1987 and 2007. In this dataset, most of the articles are manually summarized and tagged by a staff of library scientists, and there are over 650,000 article-summary pairs.

**XSum.** The **extreme summarization (XSum)** dataset [83] is an extreme summarization dataset containing BBC articles and corresponding single sentence summaries. In this dataset, 226,711 Wayback archived BBC articles are collected, which range from 2010 to 2017 and cover a wide variety of domains.

**Gigaword.** The English Gigaword dataset [36, 105] is a comprehensive collection of English newswire text data acquired by the Linguistic Data Consortium. This corpus contains four distinct international sources of English newswire and has totally 4,111,240 documents.

**3.2.2 Question Generation.** The two popular datasets for the task of question generation are shown below.



**SQuAD.** The **Stanford Question Answering Dataset (SQuAD)** [100] is a large reading comprehension dataset created by crowdworkers. The questions in this dataset are posed by crowdworkers based on a set of Wikipedia articles, and the answers are text segments from the corresponding passages. In total, SQuAD contains 107,785 question-answer pairs on 536 articles.

**MS MARCO.** The **Microsoft Machine Reading Comprehension (MS MARCO)** dataset [84] is a collection of anonymized search queries issued through Bing or Cortana for reading comprehension. The dataset contains both answerable and unanswerable questions. Each answerable question has a set of extracted passages from the retrieved response documents of Bing. There are totally 1,010,916 questions and 8,841,823 answering passages extracted from 3,563,535 web documents in this dataset.

**3.2.3 Distractor Generation.** There are three datasets widely used for distractor generation as shown below.

**SciQ.** SciQ [137] is a crowdsourced multiple choice question answering dataset, which consists of 13.7K science exam questions. The domain of this dataset covers biology, chemistry, earth science, and physics.

**MCQL.** The MCQL dataset [61] is a collection of multiple choice questions at the Cambridge O level and college level, which is crawled from the Web. This dataset totally contains 7.1K questions covering biology, physics, and chemistry.

**RACE.** RACE [51] is a reading comprehension dataset collected from the English exams in Chinese middle and high schools. This dataset contains 27,933 passages and 97,687 questions, covering all types of human articles.

### 3.3 Method

**3.3.1 Text Summarization.** Recently, the most common methods in this field are encoder-decoder-based pre-trained language models. Song et al. [116] design a novel pre-training objective to jointly pre-train the encoder and decoder, where the decoder learns to predict the masked sentence fragments in the encoder side. Given an unpaired source sentence  $x$  from the source domain  $\mathcal{X}$ , the model with parameter  $\theta$  predicts the sentence fragment  $x^{u:v}$  from position  $u$  to  $v$  with the masked sequence  $x^{\setminus u:v}$  as input, and the objective function is formulated as:

$$L(\theta; \mathcal{X}) = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \log P(x^{u:v} | x^{\setminus u:v}; \theta) = \frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \log \prod_{t=u}^v P(x_t^{u:v} | x_{<t}^{u:v}, x^{\setminus u:v}; \theta). \quad (20)$$

Lewis et al. [53] present a denoising autoencoder for pre-training, which consists of a series of noising strategies to corrupt text and a training objective to reconstruct the original sentence. In addition, Zhang et al. [154] propose a self-supervised pre-training objective specific for the text summarization task, namely, the gap-sentences generation objective, and their PEGASUS model achieves state-of-the-art performances on the mainstream datasets.

However, the pre-trained language models' ability of capturing long-term dependencies and maintaining global coherence is poor. To solve this problem, Cai et al. [8] introduce an additional encoder with a bidirectional RNN and a convolution module to simultaneously model sequential context and capture local importance to the base Transformer model. The encoder first applies a bidirectional LSTM on the source text embedding  $E$  and derives the hidden states sequence  $H$ , then uses three convolution operations with kernel sizes of 1, 3, 5 to learn  $n$ -gram features  $D$  from  $H$ , and finally a gated linear unit is applied to select features, which can be formulated by:

$$R = \sigma(W_d D + b_d) \odot (W_h H + b_h), \quad (21)$$

where  $H = \text{BiLSTM}(E)$ . Yan et al. [96] make several improvements to the traditional language models. Specifically, a new self-supervised objective called future n-gram prediction has been applied, and the n-stream self-attention mechanism is used in the decoder to be trained to predict the future n-gram of each timestep. Given a source sequence  $x$  and its target sequence  $y$ , the future n-gram prediction objective can be formulated as:

$$\mathcal{L} = - \sum_{j=0}^{n-1} \alpha_j \cdot \left( \sum_{t=1}^{T-j} \log p_{\theta}(y_{t+j} | y_{<t}, x) \right). \quad (22)$$

Ding et al. [20] propose an enhanced semantic model based on double encoders and a decoder with a Gain-Benefit gate structure, which is able to provide richer semantic information and reduce the influence of the length of generated texts on the decoding accuracy. The dual-encoder is used to capture both global and local context semantic information based on the bidirectional RNN, and the output vector of the Gain-Benefit gate module is calculated by:

$$P_t = (1 - a) \odot \mathbf{M} + a \odot \mathbf{C}_{t-1}, \text{ where } a = \text{sigmoid}(\mathbf{W}_1 \mathbf{C}_t + \mathbf{W}_2 \mathbf{S}_{t-1} + y_{t-1}), \quad (23)$$

in which  $\mathbf{C}_t$  is the contextual semantic representation with both global and local semantic information at current timestep,  $\mathbf{S}_{t-1}$  is the decoder hidden state of last timestep,  $y_{t-1}$  is the word generated at last timestep, and  $\mathbf{M}$  is the global semantic vector of the original text.

Except for the poor capability of capturing long-term dependency in traditional pre-trained language models, the problem of factual inconsistency between the generated content and source text is also severe and has not been tackled by previous works. To reduce this phenomenon, Meng et al. [9] propose an end-to-end neural corrector model with a post-editing correction strategy, which is pre-trained on artificial corrupted reference summaries. Dong et al. [24] introduce a factual correction framework containing a QA-span factual correction model and an auto-regressive one. The QA-span correction model masks and replaces one entity at a time during the iteration process. Specifically, given the source text  $x$  and a masked query  $q = (y'_1, \dots, [\text{MASK}], \dots, y'_m)$ , the correction model needs to predict the answer span via  $p(i = \text{start})$  and  $p(i = \text{end})$ , which are calculated based on the hidden states of the top layer  $h_i$  as:

$$p(i = \text{start}) = a_i^{\text{start}} = \frac{\exp(q_i^s)}{\sum_{j=0}^{H-1} \exp(q_j^s)}, \text{ where } q_i^s = \text{ReLU}(\mathbf{w}_s^\top h_i + b_s), \quad (24)$$

in which  $H$  is the number of hidden states in the encoder, and  $p(i = \text{end})$  is calculated in the similar way. The auto-regressive correction model masks all entities at the same time without iteration. Specifically, given the source text  $x$  and a masked query  $q = (y'_1, \dots, [\text{MASK}]_1, \dots, [\text{MASK}]_T, \dots, y'_m)$  from a summary with  $T$  entities, the correction model runs  $T$  steps to predict the answer span of each mask based on the corresponding masked token representation and its previously predicted entity representations. The entity representation  $\mathbf{s}_i^{\text{ent}}$  at timestep  $t$  is predicted based on the argmax and mean pooling operations, as calculated by:

$$\mathbf{s}_i^{\text{ent}} = \text{Mean-Pool}(\mathbf{h}_{p_{\text{start}}}, \mathbf{h}_{p_{\text{end}}}), \quad (25)$$

where  $p_{\text{start}} = \arg \max(a_1^{\text{start}}, \dots, a_M^{\text{start}})$ ,  $p_{\text{end}} = \arg \max(a_1^{\text{end}}, \dots, a_M^{\text{end}})$ .

**3.3.2 Question Generation.** Early methods to solve the task of question generation are mainly based on hand-crafted rules and always contain multiple procedures. Wijanarko et al. [26] propose a method that generates questions based on key-phrase that embeds Bloom's taxonomic in selecting contexts for constructing questions. The key-phrases can be selected via two methods.

The first method is based on probability of two-word sequence as formulated below:

$$score(w_i, w_j) = \frac{count(w_i w_j) - \delta}{count(w_i) \times count(w_j)}, \quad (26)$$

where  $count(w_i, w_j)$  is the number of a sequence of word  $w_i$  followed by  $w_j$ ,  $count(w_i)$  is the number of word  $w_i$  in the input document, and  $\delta$  is a constant to limit the number of phrases formed by less-frequent words. The second method is based on a Naive Bayes model, and the probability of a unique phrase is measured by:

$$Pr[key|T, D] = \frac{Pr[T|key] \times Pr[D|key] \times Pr[key]}{Pr[T, D]}, \quad (27)$$

where  $Pr[T|key]$  is the probability that a key-phrase has a  $TF \times IDF$  score  $T$ ,  $Pr[D|key]$  is the probability that it has a distance  $D$ , and  $Pr[T, D]$  is the normalization factor. Bhatia et al. [5] propose a dynamic question answering generator to generate questions and answers for the mathematical topic-quadratic equations. The randomization technique, first-order logic, and automated deduction are used for the study.

In the past few years, driven by advances in deep learning, end-to-end neural models based on Seq2Seq framework have attained more and more attention and have shown better performances. Serban et al. [112] study the factoid questions generation with RNN to transduce facts into neural language questions and provide an enormous question-answer pair corpus. Yuan et al. [153] use a Seq2Seq model with teacher forcing to improve the training and adopt policy gradient in reinforcement learning to optimize the generated results. During supervised learning, in addition to minimize the negative log-likelihood with teacher forcing, two additional signals are introduced to prevent the model from generating answer words ( $\mathcal{L}_s$ ) and encourage the output variety ( $\mathcal{L}_e$ ) as formulated below:

$$\mathcal{L}_s = \lambda_s \sum_t \sum_{\bar{a} \in \bar{\mathcal{A}}} p_\theta(y_t = \bar{a} | y_{<t}, D, A), \quad \mathcal{L}_e = \lambda_e \sum_t \mathbf{p}_t^T \log \mathbf{p}_t, \quad (28)$$

where  $\bar{\mathcal{A}}$  refers to the set of words appearing in the answer but not in the ground-truth question,  $\mathbf{p}_t$  refers to the full pointer-softmax probability of the  $t$ th word, which enables the model to interpolate between copying from the source document and generating from shortlist. And during reinforcement learning, the total reward  $R_{PPL+QA}$  is a combination of question answering reward  $R_{QA}$  (measured by F1 score) and question fluency reward  $R_{PPL}$  (measured by perplexity), as formulated below:

$$R_{PPL+QA} = \lambda_{QA} R_{QA}(\hat{Y}) + \lambda_{PPL} R_{PPL}(\hat{Y}), \quad (29)$$

$$R_{PPL}(\hat{Y}) = -2^{-\frac{1}{T} \sum_{t=1}^T \log_2 PLM(\hat{y}_t | \hat{y}_{<t})}, \quad (30)$$

$$R_{QA}(\hat{Y}) = F1(\hat{A}, A), \quad (31)$$

where  $\hat{A} = MPCM(\hat{Y})$  refers to the answer of the generated question by the **Multi-Perspective Context Matching (MPCM)** model,  $PLM$  is a language model. Zhao et al. [158] study the paragraph-level neural question generation by proposing a maxout pointer and gated self-attention networks, which mainly deals with the problem that long text (mostly paragraphs) does not perform well in the Seq2Seq model. In detail, the gated self-attention network contains two steps: (1) the encoded passage-answer representation  $\mathbf{u}_t$  is conducted matching against itself to derive the self matching representation  $\mathbf{s}_t$  at timestep  $t$ :

$$\mathbf{s}_t = \mathbf{U} \cdot \mathbf{a}_t^s = \mathbf{U} \cdot \text{softmax}(\mathbf{U}^\top \mathbf{W}^s \mathbf{u}_t). \quad (32)$$

And (2) the self matching representation  $\mathbf{s}_t$  is combined with the original passage-answer representation  $\mathbf{u}_t$ , which is then fed into a feature fusion gate to obtain the final encoded passage-answer representation  $\hat{\mathbf{u}}_t$  at timestep  $t$ :

$$\hat{\mathbf{u}} = \mathbf{g}_t \odot \mathbf{f}_t + (1 - \mathbf{g}_t) \odot \mathbf{u}_t, \text{ where } \mathbf{f}_t = \tanh(\mathbf{W}^f[\mathbf{u}_t, \mathbf{s}_t]), \mathbf{g}_t = \text{sigmoid}(\mathbf{W}^g[\mathbf{u}_t, \mathbf{s}_t]). \quad (33)$$

However, the rich information lying in the answer has not been fully explored, which leads to generating low-quality questions (e.g., having mismatched interrogative words with the answer type, copying context words far from irrelevant to the answer, including words from the target answer, and so on). To solve the problem, Sun et al. [122] propose an answer-focused and position-aware model. It incorporates the answer embedding to explicitly generate a question word matching the answer type and designs a position-aware attention mechanism by modeling the relative distance with the answer, which guides the model to copy the context words that are more close and relevant to the answer. Kim et al. [46] propose to separate the target answer from the original passage to avoid the generated question copying words from the answer. Specifically, it replaces the answer with a mask token and introduces a keyword-net to extract key information from the answer. Given the encoded answer representation  $h^a$  and the context vector  $c_t$  of current decoding step, the keyword feature of each layer of the keyword-net  $o_t^l$  can be formulated as:

$$o_t^l = \sum_j p_{tj}^l h_j^a, \text{ where } p_{tj}^l = \text{softmax}\left(\left(o_t^{l-1}\right)^\top h_j^a\right), \quad (34)$$

in which  $o_t^0$  is initialized by  $c_t$ . And the decoding hidden state  $s_t$  of current timestep is calculated by:

$$s_t = \text{LSTM}\left(y_{t-1}, s_{t-1}, c_t, o_t^L\right), \quad (35)$$

where  $y_{t-1}$  is the output token of previous timestep,  $L$  is the layer number of the keyword-net.

Moreover, many works recently attempt to conduct multi-task learning with external related tasks to further enhance the performance of question generation. Wang et al. [134] introduce a message-passing mechanism to simultaneously learn the tasks of phrase extraction and question generation, which helps the model be aware of question-worthy phrases that are worthwhile to be asked about. Jia et al. [44] conduct multi-task learning with paraphrase generation and question generation, which can diversify the question patterns of the question generation module. During training, the weights of the encoder are shared by all tasks while those of the first layer of decoder are shared with a soft sharing strategy, which is formulated by:

$$\mathcal{L}_{sf} = \sum_{d \in \mathcal{D}} \|\theta_d - \phi_d\|_2, \quad (36)$$

where  $\mathcal{D}$  is a set of shared decoder parameters,  $\theta, \phi$  refer to the parameters of the question generation task and paraphrase generation task, respectively. And a min-loss function is employed among the golden reference question and several expanded question paraphrases, which is represented by:

$$\mathcal{L}_{qg} = \min_{q \in Q} \left( -\frac{1}{T_{qg}} \sum_{t=1}^{T_{qg}} \log P\left(y_t^{qg} = \mathbf{q}_t\right) \right). \quad (37)$$

Despite the outstanding performance achieved by previous methods, the rich structure information hidden in the passage is ignored, which can be used as an auxiliary knowledge of the unstructured input text to improve the performance. Liu et al. [66] adopt a **graph convolutional network (GCN)** to identify the clue words in the input passage that should be copied into the target question. Specifically, the GCN is constructed on the syntactic dependency tree representation of each passage, and a Gumbel-Softmax layer is applied to the final representation of GCN

to sample the binary clue indicator for each word. A sample  $\mathbf{y} = (y_1, \dots, y_k)$  drawn from the Gumbel-Softmax distribution is formulated as:

$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)}, \quad (38)$$

where  $\tau$  is the temperature parameter,  $\pi_i$  is the unnormalized log probability of class  $i$ , and  $g_i$  is the Gumbel noise formulated by:

$$g_i = -\log(-\log(u_i)), \text{ where } u_i \sim \text{Uniform}(0, 1). \quad (39)$$

Wang et al. [133] construct a knowledge graph for each input sentence as the auxiliary structured knowledge and aim to generate a question based on a query path from the knowledge graph. The query representation learning is formulated as a sequence-labeling problem for identifying the involved facts to form a query, which is used to generate more relevant and informative questions.

Recently, pre-trained language models have achieved remarkable performances in question generation, far exceeding those of the previous RNN-based methods. Dong et al. [23] propose a unified model that is pre-trained with three types of natural language understanding or generation tasks, namely, unidirectional, bidirectional, and sequence-to-sequence prediction. Through the pre-training of these three tasks, the model's question-generation performance achieves significant improvements on SQuAD. Xiao et al. [141] propose an enhanced multi-flow seq2seq pre-training and fine-tuning framework to alleviate the exposure bias, which consists of an infilling generation mechanism and a noise-aware generation method, which achieves state-of-the-art performances on a wide range of datasets.

**3.3.3 Distractor Generation.** The researches mainly focus on generating multi-choice question distractors for ontologies or articles. Traditional methods primarily use hand-crafted rules or ranking method for distractor generation. Stasaski et al. [118] introduce a novel method with several ontology- and embedding-based approaches. The graph structure of the ontology is used to create complex problems linking different concepts. Liang et al. [61] introduce a ranking method with a feature-based model and a **neural net (NN)**-based model. The NN-based model consists of a generator  $G$  and a discriminator  $D$ , where  $G$  generates distractors  $d$  based on a conditional probability  $P(d|q, a)$  given question stems  $q$  and answers  $a$ , and  $D$  predicts whether a distractor sample comes from the real training data or  $G$ . The objective for  $D$  is to maximize the log-likelihood as:

$$\max_{\phi} \mathbb{E}_{d \sim P_{true}(d|q, a)} [\log(\sigma(f_{\phi}(d|q, a)))] + \mathbb{E}_{d \sim P_{\theta}(d|q, a)} [\log(1 - \sigma(f_{\phi}(d|q, a)))] \quad (40)$$

where  $f_{\phi}(d, q, a)$  is an arbitrary scoring function parameterized by  $\phi$ . And each distractor  $d_i$  sampled by  $G$  is based on another scoring function  $f_{\theta}(d, q, a)$  as formulated as:

$$p_{\theta}(d_i|q, a) = \frac{\exp(\tau \cdot f_{\theta}(d_i, q, a))}{\sum_j \exp(\tau \cdot f_{\theta}(d_j, q, a))} \quad (41)$$

where  $\tau$  is a temperature hyper-parameter. Then a cascaded learning framework is proposed to make the ranking more effective, which divides the ranking process into two stages to reduce the candidates.

Recently, deep learning-based models are widely adopted due to its overwhelming performance. For example, Chung et al. [16] utilize the BERT model to generate distractor with the autoregressive mechanism in a multi-tasking architecture. Additionally, Gao et al. [30] express the task as a sequence-to-sequence learning problem based on a hierarchical encoder-decoder network. In this model, static and dynamic attention mechanisms are adopted on the top of the hierarchical encoding structure, and a question-based initializer is used as the start point to generate distractors in the decoder. The question  $q$ , the answer  $a$ , and the word vectors in the  $i$ th sentence ( $\mathbf{w}_{i,1}, \dots, \mathbf{w}_{i,m}$ )

are first encoded via three separate bidirectional LSTM networks into  $(\mathbf{q}_1, \dots, \mathbf{q}_l)$ ,  $(\mathbf{a}_1, \dots, \mathbf{a}_k)$ , and  $(\mathbf{h}_{i,1}^e, \dots, \mathbf{h}_{i,m}^e)$ . Then another bidirectional LSTM is applied on the encoded word representations to derive the contextualized sentence representation  $(\mathbf{u}_1, \dots, \mathbf{u}_n)$ , and an average pooling layer is applied to derive their entire representations  $\mathbf{q}$ ,  $\mathbf{a}$ , and  $\mathbf{s}_i$ . After that, the static attention distribution  $\gamma_i$  can be derived through a matching layer and a normalization layer as:

$$\gamma_i = \text{softmax}(o_i/\tau), \text{ where } \tau = \text{sigmoid}(\mathbf{w}_q^\top \mathbf{q} + b_q), \quad o_i = \lambda_q \mathbf{s}_i^\top \mathbf{W}_m \mathbf{q} - \lambda_a \mathbf{s}_i^\top \mathbf{W}_m \mathbf{a} + \mathbf{b}_m. \quad (42)$$

In the decoder size, the decoder generates the hidden state  $\mathbf{h}_t^d$  at the  $t$ th timestep through an LSTM network. Then the sentence-level and word-level dynamic attention  $\beta_i$  and  $\alpha_{i,j}$  are formulated as:

$$\beta_i = \mathbf{u}_i^\top \mathbf{W}_{d_1} \mathbf{h}_t^d, \quad \alpha_{i,j} = \mathbf{h}_{i,j}^e \mathbf{W}_{d_2} \mathbf{h}_t^d. \quad (43)$$

Finally, the static and dynamic attentions are combined into  $\tilde{\alpha}_{i,j}$  to reweight the article token representations and predict the probability distribution  $P_V$  over vocabulary  $V$ :

$$P_V = \text{softmax}(\mathbf{W}_V \tilde{\mathbf{h}}_t^d + \mathbf{b}_V), \text{ where } \tilde{\mathbf{h}}_t^d = \tanh(\mathbf{W}_{\tilde{\mathbf{h}}}[\mathbf{h}_t^d; \mathbf{c}_t]), \quad \mathbf{c}_t = \sum_{i,j} \tilde{\alpha}_{i,j} \mathbf{h}_{i,j}^e, \quad \tilde{\alpha}_{i,j} = \frac{\alpha_{i,j} \beta_i \gamma_i}{\sum_{i,j} \alpha_{i,j} \beta_i \gamma_i}. \quad (44)$$

However, the answer interaction is not considered by previous works and the incorrectness of the generated distractors cannot be guaranteed. To address this problem, Qiu et al. [98] propose a framework consisting of reforming modules and an attention-based distractor generator, which is the state-of-the-art method on most widely adopted datasets (e.g., RACE). The reforming modules use the semantic distances to constrain the effect of words that are strongly related to the correct answer, and the distractor generator leverages the information of the reformed question and passage to generate the initial state and context vector, respectively. In detail, three contextual encoders are first applied to encode the passage, question, and its answer into  $\mathbf{P}$ ,  $\mathbf{Q}$ , and  $\mathbf{A}$ , then an attention mechanism and a fusion kernel are leveraged to enrich the question and answer representations into  $\tilde{\mathbf{Q}}$  and  $\tilde{\mathbf{A}}$ , where  $\tilde{\mathbf{Q}}$  is formulated by:

$$\tilde{\mathbf{Q}} = \text{Fuse}(\mathbf{Q}, \bar{\mathbf{Q}}) = \tanh([\mathbf{Q}; \bar{\mathbf{Q}}; \mathbf{Q} - \bar{\mathbf{Q}}; \mathbf{Q} \circ \bar{\mathbf{Q}}] \mathbf{W}_f + \mathbf{b}_f), \quad (45)$$

$$\bar{\mathbf{Q}} = \text{Attn}(\mathbf{Q}, \mathbf{P}) \mathbf{P} = \text{softmax}\left(\frac{\mathbf{Q} \mathbf{P}^T}{\sqrt{d}}\right) \mathbf{P}. \quad (46)$$

In the reforming question module, the reformed question  $\dot{\mathbf{Q}}_i$  is calculated through a self-attend layer and a gate layer as:

$$\dot{\mathbf{Q}} = \text{Gate}(\tilde{\mathbf{Q}}_i, \tilde{\mathbf{v}}^a) \tilde{\mathbf{Q}}_i = \left(\tilde{\mathbf{Q}}_i \mathbf{W}_g^q \tilde{\mathbf{v}}^{a^\top} + b_g^q\right) \tilde{\mathbf{Q}}_i, \quad (47)$$

$$\tilde{\mathbf{v}}_a = \text{SelfAlign}(\tilde{\mathbf{A}}) = \text{softmax}(\tilde{\mathbf{A}} \mathbf{W}_a)^\top \tilde{\mathbf{A}}. \quad (48)$$

And in the reforming passage module, the reformed passage  $\tilde{\mathbf{P}}$  is calculated by:

$$\tilde{\mathbf{P}} = \text{Fuse}(\dot{\mathbf{P}}, \bar{\mathbf{P}}), \quad \bar{\mathbf{P}} = \text{Attn}(\dot{\mathbf{P}}, \dot{\mathbf{Q}}) \dot{\mathbf{Q}}, \quad \dot{\mathbf{P}}_i = \text{Gate}(\mathbf{P}_i, \hat{\mathbf{v}}^a) \mathbf{P}_i, \quad (49)$$

$$\hat{\mathbf{v}}^a = \text{SelfAlign}(\hat{\mathbf{A}}), \quad \hat{\mathbf{A}} = \text{Fuse}(\tilde{\mathbf{A}}, \bar{\mathbf{A}}), \quad \bar{\mathbf{A}} = \text{Attn}(\tilde{\mathbf{A}}, \tilde{\mathbf{Q}}) \tilde{\mathbf{Q}}. \quad (50)$$

Maurya et al. [76] use a single encoder to encode the input triplet and three decoders to generate three distractors. The encoder employs SoftSel operation and a gated mechanism to capture the semantic relations among the elements of the input triplet.

In addition, there are many other scenarios for distractor generation. Srinivas et al. [117] develop a semi-automatic tool to help teachers quickly create multiple choice questions on code understanding. The tool first captures each code structure in the form of an abstract syntac tree and then trains a code model that maps functions to vectors. Ren et al. [103] create a model based



on a context-sensitive candidate set generator and a distractor selector for cloze-style multiple choice questions. The model first uses the correct answer as the key, combines the LDA model to mine the topic of the context, finds similar words in the semantic database, and then selects a specified number of misleading items according to a ranking model. Specifically, the probability distribution over all entities subsumed by the concepts in  $C$  is calculated based on the posterior probability  $p(c|a, q)$  as:

$$p_i = p(d_i|a, q) \propto \sum_{c \in C} p(d_i|c) p(c|a, q), \text{ where } p(c|a, q) \propto p(c|a) \sum_{k=1}^K \pi_{a,q}^{(k)} \gamma_c^{(k)}, \quad (51)$$

in which  $c$  is the concept,  $\pi_{a,q}$  is the topic distribution of complete sentence formed by the stem and key,  $\gamma_c$  is the topic distribution of concept  $c$ ,  $p(c|a)$  is the prior probability of  $a$  belonging to  $c$ ,  $K$  is the total number of topics, and  $p(d|c)$  is the typicality. Patra et al. [89] develop a system to generate named entity distractors. The system performs two types of similarity computation, namely, statistical and semantic. To speed up the distractor selection procedure, a hierarchical clustering method is proposed to represent the entities, where the entity similarities are embedded in a tree structure. The distractors are selected from the nearby entities of the correct answer of the question in the tree. Specifically, the statistical distance between the numeric attributes is calculated by:

$$\text{Sim}(P, Q) = 1 - \frac{1}{L} \sum_{i=1, \dots, L} \frac{(P_i \sim Q_i)}{\max(P_i, Q_i)}, \quad (52)$$

where  $P$  and  $Q$  represent two vectors corresponding to the target entities,  $L$  is the total number of numeric attributes. The hierarchical distance between two entities  $(x, x')$  is normalized by:

$$\text{Sim}(x, x') = \frac{d_1(x, x')}{\sqrt{d_1(x, x) \odot d_1(x', x')}}, \quad (53)$$

where  $d_1(x, x')$  is the highest tree level connecting  $x$  and  $x'$ . And the semantic similarity score between the key  $x$  and a candidate distractor  $x'$  is formulated as:

$$\text{Sim}(x, x') = \frac{|(\text{triplet}_i \in x)| \& (\text{triplet}_i \in x')|_i}{|\text{triplet}_j \in x|_j}, \quad (54)$$

where the normalization factor is the size of the triplet set corresponding to the key.

## 4 TEXT EXPANSION

### 4.1 Task

The main purpose of text expansion is to inflate the short texts to longer ones that contains more abundant information, which can be divided into two aspects: short text expansion and topic-to-essay generation. Short text expansion aims to expand a short text into a richer representation based on a set of long documents. Topic-to-essay generation aims at generating human-like diverse and topic-consistent paragraph-level text with a set of given topics. The most representative methods for each subtask are shown in Table 2.

### 4.2 Data

**4.2.1 Short Text Expansion.** There are mainly four datasets for short text expansion as listed below.

*Wikipedia.* Reference [124] constructs this dataset through a snapshot of English Wikipedia. The titles of Wikipedia articles are regarded as short texts, and the abstract of all Wikipedia articles are leveraged to construct the related long documents. This dataset contains 30,000 short texts and 4,747,988 long documents in total.

Table 2. Natural Language Generation Models for Text Expansion

Task	Model	Description
Short Text Expansion	Associated-Query-based Query Expander [6]	Ranking
	ExpaNet [124]	Retrieval + Memory Network
	FC-LSTM [113]	LDA + RNN + Ranking
	Fiction Sentence Expander [106]	RNN
Topic-to-Essay Generation	MTA-LSTM [27]	RNN + Global Coherence
	SRENN [135]	RNN + Retrieval + Global Coherence
	UD-GAN [152]	RNN + GAN + RL
	KB-based Topic-to-essay Generator [146]	CNN + RNN + GAN + RL + KB
	SCTKG [97]	CNN + RNN + GAN + RL + KG

*DBLP*. Reference [124] uses the DBLP bibliography database to construct this dataset. The titles of computer science literature represent short texts, and the abstracts of all papers are collected to construct the corresponding long documents. Statistically, this dataset consists of 81,479 short texts and 480,558 long documents.

*Programmableweb*. Reference [113] establishes a real-world dataset for service recommendation and description expansion, which is crawled from programmableweb.com in 2013 and 2016. The entire dataset contains 16,012 APIs, 7,816 Mashups, and 16,449 links between them.

*Fiction Corpus*. Reference [106] creates an English fiction corpus, which is obtained by applying sentence compression techniques on a modern fiction corpus scraped from online resources. Each story sentence has a corresponding compression. In total, the dataset contains around 17,000,000 sentences.

**4.2.2 Topic-to-essay Generation.** There are primarily five datasets for topic-to-essay generation as listed below.

*ESSAY*. The ESSAY dataset [27] is a large collection of topic compressions crawled from the Internet. The topic words are extracted by TextRank. This dataset totally contains 305,000 paragraph-level essays.

*ZhiHu*. Reference [27] constructs this dataset by crawling from a Chinese question-and-answering website called ZhiHu, which consists of 55,000 articles. The topic words of each article are specified by users in the community.

*Movie Reviews*. The **Stanford Sentiment Treebank (SST)** dataset [115] contains 11,855 single sentences of movie reviews. This dataset has two sentiment classes for each review and has a set of fully labeled parse trees.

*Beer Reviews*. Reference [77] creates this dataset by crawling from the beer review website Beer-Advocate. In total, this dataset has 1,586,259 ratings on 66,051 items that are scored by 33,387 users.

*Customer Reviews*. Reference [41] constructs a customer review dataset consisting of five electronics products, which is collected from Amazon.com and C|net.com. There are totally 1,886 items in this dataset.

## 4.3 Method

**4.3.1 Short Text Expansion.** Early works primarily use statistical similarity for short text expansion. Billerbeck et al. [6] propose a method about the query expansion using associated queries for web search engines. The main method is to associate a query that closely matches the document

in a given log containing a large number of queries. Then, the query associated with the query document is a reasonable description of the document and can be used to expand the query text.

Benefiting from the development of deep learning, many end-to-end frameworks based on neural networks have been developed. Tang et al. [124] propose an end-to-end solution based on deep memory network for short text extension. In this work, the original short text  $q$  is first used as a query to search for a set of potentially relevant long documents  $C_q$ , which will be used as the material for text expansion, from an external large collection  $C$ . In the following process, the short text  $q$  is represented as the average vector of words in it, i.e.,  $\vec{q}$ , and each document is also represented by the average vector of words belong to it, i.e.,  $\vec{d}$ . To further identify the relevant documents in  $C_q$ , both soft attention and hard attention mechanisms are utilized, and the information read from the document set can be written as:

$$\vec{o} = \sum_{i=1}^K p_i \cdot \vec{d}_i = \sum_{i=1}^K \text{softmax}\left(\frac{\vec{q}^T \vec{d}_i + g_i}{\tau}\right) \cdot \vec{d}_i, \quad (55)$$

where  $g_i$  follows the Gumbel(0,1) distribution, and  $\tau$  is the temperature hyperparameter. Then the two sources of information,  $\vec{q}$  and  $\vec{o}$ , are integrated by GRU as follows:

$$\vec{z} = \sigma(\mathbf{W}^{(z)}\vec{q} + \mathbf{U}^{(z)}\vec{o}), \quad (56)$$

$$\vec{r} = \sigma(\mathbf{W}^{(r)}\vec{q} + \mathbf{U}^{(r)}\vec{o}), \quad (57)$$

$$\vec{o}' = \tanh(\mathbf{W}\vec{q} + \vec{r} \circ \mathbf{U}\vec{o}), \quad (58)$$

$$\vec{q}' = (1 - \vec{z}) \circ \vec{q} + \vec{z} \circ \vec{o}', \quad (59)$$

where  $\circ$  means element-wise multiplication, both the Sigmoid function  $\sigma(x)$  and  $\tanh(x)$  are operated on element-wise. The output  $\vec{q}'$  is the expanded representation of the input short text  $q$ . This process can be repeated several times for the same  $q$  to continuously extend its representation to mimic the human behavior when querying a piece of short text.

Shi et al. [113] present a text expansion method for service recommendation system. The description of services is first expanded at sentence level by a probabilistic topic model; in this process, the similarity between target sentence  $u$  and another sentence  $v$  from existed corpus can be calculated as:

$$\text{Similarity}(u, v) = \mu \cdot D_{JS}(u, v) + (1 - \mu) \cdot D_{JS}(S_u, S_v), \quad (60)$$

where  $S_u$  and  $S_v$  are descriptions containing  $u$  and  $v$ , respectively.  $\mu$  is a parameter used to balance weights of sentence and description on the final similarity measurement.  $D_{JS}$  is the function of JS divergence, which is used to measure the similarity between two items; for more details, readers can refer to Reference [113]. So, for each target sentence, a collection of similar sentences can be found and ranked in descent order to select top  $N$  most similar ones for description extension.

Safovich et al. [106] design a neural sentence expander trained on a corpus of fiction sentence compressions for the task of sentence expansion and enhancement, which is the state-of-the-art method on existing datasets (e.g., Fiction Corpus). In this method, a Seq2Seq model is used to predict sentences from the original input. To tackle the problem of copying input words during generation process, the modified negative log-likelihood loss function is used to increase the significance of learning new words. Specifically, the modified cross-entropy of generating is calculated as:

$$\mathcal{L} = - \sum_t (1 + \lambda I_{T-S}(w_t)) \log p(w_t | w_1, \dots, w_{t-1}), \quad (61)$$

in which  $I$  denotes the indicator function for a set,  $T$  the ground truth,  $S$  the source tokens,  $\lambda$  is a parameter that controls the learning importance of the words in target sentence while not in original input. With this modification, the model no longer degenerates to copying its input and can be trained as desired. In addition, to increase the novelty of generated sentence, Safovich et al. [106] also develop a controlled sampling method for the decoding process so the model can generate diverse output.

**4.3.2 Topic-to-essay Generation.** With the popularity of deep learning, the RNN-based Seq2Seq framework has been widely used in this field. Many works have been proposed to improve the traditional Seq2Seq framework.

To improve the global coherence of the generated essays, Feng et al. [27] propose an LSTM model with attention mechanism for essay generation. The main idea of this work is to maintain a topic coverage vector, each dimension of which represents the degree to which a topic word needs to be expressed in future generation, to adjust the attention policy, so the model can consider more about unexpressed topic words. Specifically, the topic coverage vector is updated by parameter  $\phi_j$ , and it can be regarded as the discourse-level weight of  $topic_j$ , the word embedding of topic word  $i$ . Topic coverage vector  $C_t$  is initialized as a  $k$  dimensional vector, i.e.,  $C_0$ ,  $k$  is the number of input topic words, and each value of  $C_0$  is 1.0. At timestep  $t$  in generation process, each element  $c_{t,j}$  is updated as follows:

$$c_{t,j} = c_{t-1,j} - \frac{\alpha_{t,j}}{\phi_j}, \quad (62)$$

in which  $\alpha_{t,j} = \exp(g_{tj}) / \sum_{i=1}^K \exp(g_{ti})$  is the attention weight of topic word  $i$  at timestep  $t$ , and  $g_{tj}$  is the attention score on  $topic_j$  at timestep  $t$ , which is calculated as:

$$g_{tj} = c_{t-1,j} v_a^T \tanh(W_a h_{t-1} + U_a topic_j), \quad (63)$$

where  $h_{t-1}$  is the hidden representation of the LSTM at timestep  $t - 1$ , and  $v_a$ ,  $W_a$ ,  $U_a$  are all parameters to be optimized. Therefore, the probability of the next word  $y_t$  can be defined as:

$$P(y_t | y_{t-1}, T_t, C_t) = \text{softmax}(\text{linear}(h_t)), \quad (64)$$

where the topic representation  $T_t$  is formulated as  $T_t = \sum_{i=1}^K \alpha_{tj} topic_j$ .

Wang et al. [135] propose an enhanced neural network based on self-attention and retrieval mechanisms, and the encoder and decoder are constructed with self-attention to model longer dependence. And to alleviate the duplication problem, a retrieval process is adopted to collect topic-related sentences as an aid for essay generation. Specifically, input topic words are divided into  $m$  groups, then the material  $M = \{S_1, \dots, S_m\}$  can be collected based on cosine distance between the topic group and sentences in the corpus that is created by dividing the training set, where  $S_i$  is a sentence that corresponds to the  $i$ th topic group. For the different sequence relations between topic words and material sentences, two encoders with the same structure but different parameters are used to get the hidden representation  $H^{Topic} = \{h_1^T, \dots, h_k^T\}$  and  $H^{Material} = \{h_1^M, \dots, h_{ml}^M\}$ , based on which the hidden states of generated of essay words  $H = \{h_1, \dots, h_n\}$  are obtained by essay decoder. Finally, the output probabilities of each essay word can be computed as:

$$p(y_t | y_{t-1}, T, M) = \text{softmax}(\text{linear}(h_t)). \quad (65)$$

Meanwhile, **Generative Adversarial Net (GAN)** has shown promising results for topic-to-essay generation, which is effectively applied in Reference [152] including a GAN model and two-level discriminators. The first discriminator guides the generator to learn the paragraph-level information and sentence syntactic structure with multiple LSTMs, and the second one processes

higher-level information such as topic and sentiment for essay generation. Then the reward is calculated based on results of two discriminators, and the generator  $G_\theta$  tries to maximize expected reward from the initial state till the end state via the formulation:

$$J(\theta) = \sum_{t=1}^T E(R_t | S_{t-1}, \theta) = \sum_{t=1}^T G_\theta(y_t | Y) [\lambda(D_\phi(Y)) + (1 - \lambda)D_Y(Y)], \quad (66)$$

where  $\lambda$  is a manually set weight,  $Y$  is a complete sequence,  $R_t$  is the reward for a whole sequence,  $D_\phi$  and  $D_Y$  are the first and second discriminators, respectively.

However, the generated essays of previous works still lack novelty and topic-consistency, based on which, some works leverage the external knowledge to solve this problem. Yang et al. [146] introduce a model with the aid of commonsense knowledge in ConceptNet: In detail, each topic is used as a query to retrieve  $k$  neighboring concepts with pre-trained embeddings stored as commonsense knowledge in a memory matrix  $M_0$ ; in the decoding phase, the generator  $G_\theta$  refers to the memory matrix for text generation, and the hidden state of the decoder at timestep  $t$  is:

$$s_t = \text{LSTM}(s_{t-1}, [e(y_{t-1}); c_t; m_t]), \quad (67)$$

where  $[\cdot]$  means the concatenation of vectors,  $y_{t-1}$  is the word generated at timestep  $t - 1$ .  $c_t$  is the context vector that is computed by integrating the hidden representations of the input topic sequence, and  $m_t$  is the memory vector extracted from  $M_t$  based on the attention mechanism. As the generation progresses, the topic information that needs to be expressed keeps changing, which requires the memory matrix to be dynamically updated, so for each memory entry  $M_t^i$  in  $M_t$ , a candidate update memory  $\tilde{M}_t^i$  is computed as:

$$\tilde{M}_t^i = \tanh(U_1 M_t^i + V_1 e(y_t)), \quad (68)$$

where  $U_1$  and  $V_1$  are trainable parameters. To determine how much the  $i$ th memory entry should be updated, the adaptive gate mechanism is adopted:

$$g_t^i = \text{sigmoid}(U_2 M_t^i + V_2 e(y_t)), \quad (69)$$

where  $U_1$  and  $V_1$  are trainable parameters.  $M_t^i$  is updated by:

$$M_{t+1}^i = (\mathbf{1} - g_t^i) \odot M_t^i + g_t^i \odot \tilde{M}_t^i, \quad (70)$$

in which  $\mathbf{1}$  refers to the vector with all elements 1 and  $\odot$  denotes point-wise multiplication. Following their work, Qiao et al. [97] propose a topic-to-essay generator based on the conditional variational auto-encoder framework to control the sentiment and introduces a topic graph attention mechanism to sufficiently use the structured semantic information, which is ignored in Reference [146], so the quality of generated essays is further improved and reaches the state-of-the-art level at present on the mainstream datasets (e.g., ZhiHu).

## 5 TEXT REWRITING AND REASONING

### 5.1 Task

The target of text rewriting and reasoning is to make a reversion of the text or apply reasoning methods to generate responses, which mainly contains two subtopics: text style transfer and dialogue generation.

Text style transfer is the method to transform the attribute style of the sentence while preserving its attribute-independent content. Dialogue generation aims to automatically generate approximate answers to a series of given questions in a dialogue system. We show several classical methods for each subtask in Table 3.

Table 3. Natural Language Generation Models for Text Rewriting and Reasoning

Task	Model	Description
Text Style Transfer	Spearean Modern Language Translator [43]	RNN
	ARAE [157]	RNN + GAN + WAE
	FM-GAN [13]	RNN + GAN
	Unpaired Sentiment-to-Sentiment Translator [143]	RNN + RL + Unsupervised
	Non-Offensive Language Translator [85]	CNN + RNN + Unsupervised
	BST [95]	CNN + RNN + Unsupervised
	DualRL [72]	RNN + RL + Unsupervised
	Exploration-Evaluation-based Text Style Translator [28]	RNN + Metrics Design
	Evaluating Style Transfer for Text [80]	Embedding + CNN + RNN + Metrics Design
	Error Margins Matter in Style Transfer Evaluation [125]	VAE + RNN + Metrics Design
Dialogue Generation	PPVAE [25]	WAE + GAN
	MARM [159]	RNN
	GAN-AEL [145]	CNN + RNN + GAN
	AEM [74]	RNN + Context-response Interaction
	AGMN [57]	RNN + KG + Retrieval + Context-response Interaction
	Post-KS [60]	RNN + KB
	KADG [45]	Transformer + KB
	$p^2$ BOT [67]	Transformer + Multi-task Learning + Persona
	KnowledGPT [136]	Transformer + Knowledge Selection
	PLATO [4]	Transformer + Multi-task Learning
	TransferTransfo [138]	Transformer + Multi-task Learning

## 5.2 Data

**5.2.1 Text Style Transfer.** The following four datasets are widely applied for the task of text style transfer:

**Yelp Review.** The Yelp Review dataset is provided by the Yelp Dataset Challenge,<sup>1</sup> which contains a large amount of business review texts. This dataset contains 1.43M, 10K, and 5K pairs for training, validation, and testing, respectively.

**Amazon Food Review.** Reference [77] creates this dataset by crawling reviews from the Fine Foods category of Amazon. This dataset contains 367K, 10K, and 5K pairs for training, validation, and testing, respectively.

**EMNLP2017 WMT News.** Reference [37] picks the News section data from the EMNLP2017 WMT<sup>2</sup> Dataset, which is a large long-text corpus consists of 646,459 words and 397,726 sentences. After being preprocessed, this dataset contains 278,686 and 10,000 sentences for training and testing, respectively.

**GYAFC.** The **Grammarly's Yahoo Answers Formality Corpus (GYAFC)** [101] is a large corpus for formality stylistic transfer. The informal sentences are collected from Yahoo Answers<sup>3</sup> with domains of Entertainment & Music and Family & Relationships, and the corresponding formal sentences are created using Amazon Mechanical Turk. This dataset totally contains 106,000 sentence pairs.

**5.2.2 Dialogue Generation.** There are four popular datasets for dialogue generation, as shown below.

**DailyDialog.** The DailyDialog dataset [58] is a high-quality multi-turn dialog dataset containing daily conversations. The dialogues in the dataset are formally written by human with reasonable speaker turns, and often concentrate on a certain topic. Statistically, this dataset contains 13,118

<sup>1</sup><https://www.yelp.com/dataset/>.

<sup>2</sup><http://statmt.org/wmt17/translation-task.html>.

<sup>3</sup><https://answers.yahoo.com/answer>.



multi-turn dialogues, nearly eight average speaker turns per dialogue, and about 15 average tokens per utterance.

*UDC.* The **Ubuntu Dialogue Corpus (UDC)** [70] is created based on the two-person conversations about Ubuntu-related problems in the Ubuntu chat logs<sup>4</sup> from 2004 to 2015. This dataset contains about 1 million multi-turn dialogues, over 7 million utterances, and 100 million words.

*Persona-chat.* The Persona-chat dataset [155] is an engaging and personal chit-chat dialogue dataset collected by Amazon Mechanical Turk. Each of the paired crowdworkers condition their dialogue on a given provided profile. This dataset contains 164,356 utterances in total.

*Wizard-of-Wikipedia.* The Wizard-of-Wikipedia dataset [19] is a large crowd-sourced collection for open-domain dialogue. Each of the paired speakers conduct open-ended chit-chat, and one of the speakers needs to link the knowledge to sentences from existing Wikipedia articles connected to the topic. This dataset contains 22,311 dialogues with 201,999 turns.

### 5.3 Method

*5.3.1 Text Style Transfer.* In this field, the RNN-based Seq2Seq framework and deep latent variable model are broadly adopted. Jhamtani et al. [43] design a copy-enriched sequence-to-sequence model to transform text from modern English to Shakespearean English. The model first uses a fixed pre-trained embedding vector to represent each token and uses a bidirectional LSTM to encode sentences. In this work,  $\overrightarrow{LSTM}_{enc}$  and  $\overleftarrow{LSTM}_{enc}$  represent the forward and reverse encoder.  $h_t^{\overrightarrow{enc}}$  represents hidden state of encoder model at step t. The following equations describe the model:

$$h_t^{\overrightarrow{enc}} = \overrightarrow{LSTM}_{enc}(h_{t-1}^{enc}, E_{enc}(x_t)), \quad (71)$$

$$h_t^{\overleftarrow{enc}} = \overleftarrow{LSTM}_{enc}(h_{t+1}^{enc}, E_{enc}(x_t)), \quad (72)$$

$$h_t^{enc} = h_t^{\overrightarrow{enc}} + h_t^{\overleftarrow{enc}}. \quad (73)$$

In this work, only the forward and backward encoder states are added, and the standard connection is not used because it does not add additional parameters. Then a mixture model of RNN and pointer network are employed to transfer the text style. The pointer module provides location-based attention and output probability distribution due to pointer network module can be expressed as follows:

$$p_t^{PTR}(w) = \sum_{x_j=w} (\beta_j). \quad (74)$$

Zhao et al. [157] propose an adversarially regularized autoencoder framework to generalize the adversarial autoencoder, which combines a discrete autoencoder with a regularized latent representation of GAN and can be further formalized by the Wasserstein autoencoder. The model is trained with coordinate descent across: (1) the encoder and decoder to minimize reconstruction, (2) the critic function to approximate the  $W$  term, and (3) the encoder adversarially to the critic to minimize  $W$ :

$$(1) \min_{\phi, \psi} L_{rec}(\phi, \psi) = E_{X \sim P_*} [-\log p_{\psi}(x | enc_{\phi}(x))], \quad (75)$$

$$(2) \max_{w \in W} L_{cri}(w) = E_{X \sim P_*} [f_w(enc_{\psi}(x))] - E_{\tilde{z} \sim P_z} [f_w(\tilde{z})], \quad (76)$$

<sup>4</sup><http://irclogs.ubuntu.com/>.

$$(3) \min_{\phi} L_{enc}(\phi) = E_{X \sim P_x} [f_w(enc_{\psi}(x))] - E_{\tilde{z} \sim P_z} [f_w(\tilde{z})]. \quad (77)$$

Here,  $enc_{\psi}$  is a deterministic encoder function.  $P_z$  is the prior distribution, and  $f_w(\tilde{z})$  is the critic/discriminator.

Chen et al. [13] utilize optimal transport to improve the ability of traditional GAN in processing discrete texts with the objective of feature-mover's distance. In this work, the **feature-mover's distance (FMD)** between two sets of sentence features is then defined as:

$$D_{FMD}(P_f, P_{f'}) = \min_{T \geq 0} \sum_{i=1}^m \sum_{j=1}^n T_{ij} \cdot c(f_i, f'_j) = \min_{T \geq 0} \langle T, C \rangle, \quad (78)$$

where  $\sum_{j=1}^n T_{ij} = \frac{1}{m}$  and  $\sum_{i=1}^m T_{ij} = \frac{1}{n}$  are the constraints, and  $\langle, \rangle$  represents the Frobenius dot-product. In this work, the transport cost is defined as the cosine distance:  $c(f_i, f'_j) = 1 - \frac{f_i^T f'_j}{\|f_i\|_2 \|f'_j\|_2}$  and  $C$  is the cost matrix.

However, the problem of lacking supervised parallel data has not been well studied by the above works. To tackle this problem, many unsupervised methods have been proposed. Xu et al. [143] propose a cycled reinforcement learning approach through the cooperation between the neutralization and emotionalization modules. The neutralization module extracts the non-emotional part of the sentence with a single LSTM and a self-attention-based sentiment classifier, and the emotionalization module generates emotional words and adds them to the semantic content with a bi-decoder-based encoder-decoder framework.

Santos et al. [85] propose an unsupervised style transfer model to convert offensive language into non-offensive one. A single collaborative classifier is used to train the encoder-decoder network, and an attention mechanism with a cycle consistency loss is adopted to preserve the content. Prabhumoye et al. [95] realize the style transfer of gender, political slant, and sentiment through an unsupervised back-translation method. The transferring process is divided into two stages. The first stage learns the latent representation with back-translation of the input sentence through a language translation model, and the second stage adopts an adversarial generation technology to enable the output to match the desired style. The latent representation with back-translation of the input sentence can be described as:

$$z = Encoder(X_f; \theta_E), \quad (79)$$

where,  $x_f$  is the sentence  $x$  in language  $f$ .  $\theta_E$  represent the parameters of the encoder of language  $f \rightarrow$  language  $e$  translation system

Luo et al. [72] reformulate the traditional unsupervised style transferring task as a one-step mapping problem and propose a dual reinforcement learning framework to train the source-to-target and target-to-source mapping models. In this work, two reward methods that can evaluate style accuracy and content preservation separately are proposed.  $R_s = P(s_y|y'; \psi)$  formulate the style classifier reward where  $\psi$  is the parameter of the classifier and is fixed during the training process, and  $R_c = P(x|y'; \phi)$  represents the reward for preserving content. To encourage the model to improve both the content preservation and the style accuracy, the final reward is the harmonic mean of the above two rewards:

$$R = (1 + \beta^2) \frac{R_c \cdot R_s}{(\beta^2 \cdot R_c) + R_s}, \quad (80)$$

where  $\beta$  is a harmonic weight aiming to control the tradeoff between the two rewards.

Moreover, another challenge in this field is that there does not exist reliable evaluation metrics, which is neglected by previous works. To alleviate this issue, Fu et al. [28] propose two aspects

of evaluation metrics to measure the transfer strength and content preservation of style transfer. The transfer strength aims to evaluate whether the style is transferred through an LSTM-sigmoid classifier. The style is defined in Equation (30). This classifier is based on keras examples2. Transfer strength accuracy is defined as  $\frac{N_{right}}{N_{total}}$ ,  $N_{total}$  is the number of test data, and  $N_{right}$  is the number of correct case that is transferred to target style.

$$l_{style} = \begin{cases} paper(positive) & output \leq 0.5 \\ news(negative) & output \geq 0.5 \end{cases}. \quad (81)$$

The content preservation is used to evaluate the similarity between source and target texts and is calculated by the embedding cosine distance. Content preservation rate is defined as cosine distance (31) between source sentence embedding  $v_s$  and target sentence embedding  $v_t$ .

$$score = \frac{v_s^T v_t}{\|v_s\| \cdot \|v_t\|}. \quad (82)$$

Mir et al. [80] specify three aspects of evaluation metrics including style transfer intensity, content preservation, and naturalness. The style transfer intensity is measured by Earth Mover's Distance, the content preservation is calculated by METEOR and embedding-based metrics, and the naturalness is obtained by an adversarial evaluation method.

It is also instructive that Tikhonov et al. [125] also point out the three significant problems encountered in the evaluation metrics of style transfer. These problems mainly illustrate that the measures of style accuracy and content preservation are often different in various style transfer tasks. Therefore, they propose to take BLEU between input and human-rewritten texts into consideration to better measure the performance of style transfer models. Additionally, Duan et al. [25] propose the **Pre-train and Plug-in Variational Autoencoder (PPVAE)**, which is a model-agnostic framework towards flexible conditional text generation and consists of PretrainVAE and PluginVAE, where PretrainVAE aims to learn the original style of the sentence, and PluginVAE aims to learn the latent space of new style. The PPVAE achieves the state-of-the-art performance on the Yelp Reviews dataset.

**5.3.2 Dialogue Generation.** The RNN-based or GAN-based Seq2Seq model is widely leveraged to handle this task. Zhou et al. [159] propose a mechanism-aware neural machine based on a probabilistic RNN-based Seq2Seq framework. The model first uses latent embeddings to represent the corresponding mechanisms, then an encoder-diverter-decoder framework is leveraged to generate mechanism-aware context. In this study, there are  $M$  latent mechanisms  $M_{i=1}^M$  for response generation. Then,  $p(y|x)$  can be expanded as follows:

$$p(y|x) = \sum_{i=1}^M p(y, m_i|x) = \sum_{i=1}^M p(m_i|x)p(y|m_i, x), \quad (83)$$

where  $p(m_i|x)$  represents the probability of the mechanism  $m_i$  conditioned on  $x$ . This probability actually measures the degree that  $m_i$  can generate the response for  $x$ . The bigger of this value is, the more degree that the mechanism  $m_i$  can be used to generate the responses for  $x$ . Additionally,  $p(y|m_i, x)$  measures the probability that the response  $y$  is generated by the mechanism  $m_i$  for  $x$ . With the modeling of  $p(m_i|x)$  and  $p(y|m_i, x)$  the objective of likelihood maximization, namely,

$$\sum_{(x,y) \in D^c} \log p(y|x) = \sum_{(x,y) \in D^c} \log \sum_{i=1}^M p(m_i|x)p(y|m_i, x), \quad (84)$$

is used to learn the mechanism embeddings  $M_{i=1}^M$  and other model parameters.

Xu et al. [145] introduce a GAN framework comprising a generator, a discriminator, and an approximate embedding layer to generate informative responses. The generator uses a Seq2Seq model with GRU to generate responses, and the discriminator uses a convolutional neural network to judge the difference between human responses and machine responses. In the approximate embedding layer, the overall word embedding approximation is computed as:

$$\hat{e}_{w_i} = \sum_{j=1}^V e_j \cdot \text{softmax}(W_p(h_i + Z_i) + b_p)_j, \quad (85)$$

where  $w_p$  and  $b_p$  are the weight and bias parameters of the word projection layer, respectively, and  $h_i$  is the hidden representation of word  $w_i$ , from the decoding procedure of the generator  $G$ .

However, previous works ignore the semantic and utterance relationships between the context and response, and thus the obtained responses are not satisfying. To solve this problem, Luo et al. [74] propose an auto-encoder matching model with a mapping module. In this model, two auto-encoders are leveraged to learn the semantic representations, and the mapping module is used to learn the utterance-level dependency between the context and response. In this mapping module, for simplicity, there is only a simple feedforward network for implementation. The mapping module  $M_y$  transforms the source semantic representation  $h$  to a new representation  $t$ . To be specific, we implement a **multi-layer perceptron (MLP)**  $g(\cdot)$  for  $M_y$  and train it by minimizing the L2-norm loss  $J_3(y)$  of the transformed representation  $t = g(h)$  and the semantic representation of target response  $s$ :

$$J_3(y) = \frac{1}{2} \|t - s\|_2^2. \quad (86)$$

Li et al. [57] design a dual encoder model with an attention mechanism and a graph attention network. The attention mechanism is responsible of capturing the relationship between context and response, and the graph attention network is used to integrate the knowledge connections of domain words. The concept representation in the domain knowledge is constructed by a series of triples,  $G(x) = \{T1, T2, \dots, Tn\}$  where  $T_i$  has the same concept node  $u$  but different neighbor concept  $v$  and the graph representation of the concept  $g(x)$  can be calculated by graph attention mechanism as:

$$g(x) = \sum_{i=1}^n \alpha T_i [u_i^e; v_i^e], \text{ where } \alpha T_i = \frac{\exp(\beta_{T_i})}{\sum_{j=1}^n \exp(\beta_{T_j})}, \beta_{T_i} = \text{ReLU} \left( \left[ (u_i^e)^T W v_i^e \right] \right), \quad (87)$$

in which  $(u_i, r_i, v_i) = R_i \in G(x)$  is the  $i$ th triple in the dataset.

Some researchers try to improve the response quality by incorporating external knowledge bases and propose many strategies to select appropriate knowledge. Lian et al. [60] present a knowledge selection mechanism by separating the posterior distribution from the prior distribution. The distance between the posterior and prior distributions are minimized by the KL divergence during training, and during inference, the knowledges are selected and incorporated into the response based on the prior distribution. The **Kullback-Leibler divergence loss (KLDivLoss)**, to measure the proximity between the prior distribution and the posterior distribution, which is defined as follows:

$$L_{KL}(\theta) = \sum_{i=1}^N p(k = k_i | x, y) \log \frac{p(k = k_i | x, y)}{p(k = k_i | x)}, \quad (88)$$

where  $\theta$  denotes the model parameters.

Jiang et al. [45] propose a knowledge augmented response generation model to improve the knowledge selection and incorporation. The model consists of a divergent knowledge selector and a knowledge aware decoder, where the selector conducts a one-hop subject reasoning over facts

to reduce the subject gap in the knowledge selection, and the decoder is used to efficiently incorporate the selected fact. In addition, Liu et al. [67] concern about the importance of conversational understanding for the high-quality chit-chat systems and propose the **Persona Perception Bot** ( $P^2$  BOT). Different from other existing models,  $P^2$  BOT focuses on a important and previously overlooked concept, mutual persona perception, which is more appropriate to describe the process of information exchange that enables interlocutors to understand each other. The  $P^2$  BOT is also the current state-of-the-art model on the Persona-chat dataset.

Recently, pre-trained language models (e.g., BERT) have shown significant improvements over traditional RNN-based methods in many NLP tasks and are also applied to this task. Wang et al. [136] propose an encoder-decoder framework containing a BERT encoder and a transformer decoder. The encoder is used to learn semantic representations for both unstructured text and conversational history, and the decoder is leveraged to generate the dialogue response. In the encoder of this work, the input embedding is the sum of its token embedding, knowledge indicating embedding and position embedding:

$$I(x_i) = E(x_i) + T(x_i) + P(x_i), \quad (89)$$

where  $E(x_i)$ ,  $T(x_i)$ ,  $P(x_i)$  are word embedding, knowledge indication embedding, and position embedding, respectively. The input embeddings are then fed into BERT model to get the knowledge and dialogue history encoding representations.

Bao et al. [4] design a pre-training framework with discrete latent variables. The pre-training tasks include response generation and latent act recognition, which are jointly pre-trained through a unified network with shared parameters. Furthermore, inspired by the core idea of transfer learning, Wolf et al. [138] propose the TransferTransfo, which uses the paradigm of transfer learning to fine-tune the powerful transformer models. The specific fine-tuning tasks they select include: language modeling task, next utterance retrieval task, and generation task. Different fine-tuning tasks endow TransferTransfo with generalization performance for dialogue generation tasks in different scenarios.

## 6 FROM IMAGE TO TEXT GENERATION

### 6.1 Task

The image-based text generation aims at explaining or summarizing the visual concept of the given image, which mainly consists of three parts: image caption, video caption, and visual storytelling. The purpose of image captioning is to generate summaries from an image. Based on image captions, video caption aims to generate the summary of a series of images. Visual storytelling not only identifies the correlation between objects in a single picture but also gives the logical relationship between consecutive sequential images. It should be noted that the language generation component of VQA [1, 2, 75] model is a relatively similar to that of image caption. There exists the main distinction that current VQA systems [10, 11, 21, 33, 34, 130] are focused on reasoning process and mainly designed to choose answers from a given candidate answer set, which is not quite related to the natural language generation. Several popular methods for each subtask are shown in Table 4.

### 6.2 Data

*6.2.1 Image Caption.* The literature review of the image caption datasets is shown below.

*Flickr30k.* Reference [94] contains 31,783 images collected from Flickr. Most of these images depict humans performing various activities. Each image is paired with five crowd-sourced captions.

Table 4. From Image to Text Generation

Task	Model	Description
Image Caption	Show and Tell [129]	CNN + LSTM
	BUTD [1]	Faster-RCNN + LSTM
	Knowing When to Look [71]	CNN + LSTM
	Exploring Visual Relationship for Image Captioning [149]	Faster-RCNN + LSTM + GCN
	Self-Critical Sequence Training [104]	CNN + LSTM + RL
	Auto-Encoding Scene Graphs [148]	CNN + LSTM + GCN + RL
Video Caption	OFA [132]	Transformer + Multimodal
	LRCN [22]	CNN + LSTM
	S2VT [128]	CNN + LSTM + Knowledge
	Dense Caption Events [140]	CNN + LSTM + Daps
	Masked Transformer [160]	CNN + TCN + Transformer
	Hierarchical Reinforcement Learning [131]	CNN + LSTM + RL
	Adversarial Inference [88]	CNN + LSTM + Discriminator
	VideoBERT Pretrain [121]	CNN + BERT
	ActBERT Pretrain [161]	CNN + BERT + Multi-task Learning
	ClipBERT [52]	CNN + BERT + Clip Sampling
Visual Storytelling	UniViLM [73]	Transformer + Multimodal
	Informative Visual Storytelling [55]	CNN + GRU
	Knowledgeable Storyteller [147]	CNN + GRU + Graph
	Composite Reward [40]	CNN + RNN + MLE
	KAGS [56]	CNN + RNN + KG

*COCO*. Reference [65] is the largest image-captioning dataset, containing 82,783, 40,504, and 40,775 images for training, validation, and test, respectively. This dataset is more challenging, since most images contain multiple objects in the context of complex scenes. Each image has five human-annotated captions.

*Visual Genome*. Reference [50] is composed of dense annotations of objects, attributes, and relationships within each image to learn these models. Specifically, this dataset contains over 108K images where each image has an average of 35 objects, 26 attributes, and 21 pairwise relationships between objects.

**6.2.2 Video Caption.** There are mainly three popular datasets for video caption, as shown below.

*MSR-VTT*. Reference [142] is the most widely used video-caption dataset, which contains 7,180 videos of 20 categories. This is created by collecting 41.2 hours of 10K web video clips from a commercial video search engine.

*Charades*. Reference [114] is collected with a focus on common household activities using the Hollywood in Homes approach. This dataset contains 9,848 videos with 66,500 annotations describing 157 actions.

*ActivityNet*. Reference [49] is a large dataset that connects videos to a series of temporally annotated sentences. Each sentence describes what occurs in a unique segment of a video. Specifically, ActivityNet contains 20K videos with 100K sentence-level description.

**6.2.3 Visual Storytelling.** The following two datasets are widely used for visual storytelling:

*VIST*. Reference [42] is the most widely used dataset for visual storytelling. It contains 10,032 visual albums with 50,136 stories. Each story contains five narrative sentences, corresponding to five grounded images, respectively.



*VideoStory*. Reference [32] contains 20K videos posted publicly on a social media platform amounting to 396 hours of video with 123K sentences.

### 6.3 Method

**6.3.1 Image Captioning.** Image captioning aims to generate a description of the given image. The Show-Tell model [129] proposed an encoder-decoder-based framework that encodes images into feature vectors with **Convolution Neural Networks (CNN)**, and decodes the feature vectors into words with **Recurrent Neural Networks (RNN)**. To obtain the fine-grained visual concepts, attention-based image captioning model [71] was proposed to ground words with the corresponding part of imaging. Considering the fact that region-aware feature better fits the human visual system, Anderson et al. [1] Propose a recognized baseline called **Bottom-Up-Top-Down (BUTD)** for image caption. The BUTD is composed of two LSTM layers. The first LSTM layer is designed to capture the top-down visual attention model, while the second LSTM layer is regarded as a language model. Given the mean-pooled image feature  $\bar{\mathbf{v}}$ , a word embedding matrix  $W_e$ , and a one-hot encoding  $\Pi_t$  of the input word at time  $t$ , we could obtain the output  $\mathbf{h}_t^1$  as:

$$\mathbf{h}_t^1 = \text{LSTM} \left( \left[ \mathbf{h}_{t-1}^2, \bar{\mathbf{v}}, W_e \Pi_t \right], \mathbf{h}_{t-1}^1 \right), \quad (90)$$

and the normalized attention weight  $a_{i,t}$  can be represented by the following formulations:

$$a_{i,t} = \mathbf{w}_a^T \tanh(W_{va} \mathbf{v}_i + W_{ha} \mathbf{h}_t^1), \quad (91)$$

$$\boldsymbol{\alpha}_t = \text{softmax}(\mathbf{a}_t), \quad (92)$$

The attended image feature used as input to the language LSTM is calculated as a convex combination of all input features as  $\hat{\mathbf{v}}_t = \sum_{i=1}^K \alpha_{i,t} \mathbf{v}_i$ . The input to the language model LSTM consists of the attended image feature, concatenated with the output of the attention LSTM, given by  $\mathbf{h}_t^2 = [\hat{\mathbf{v}}_t, \mathbf{h}_t^1]$ . Using the notation  $y_{1:T}$  to refer to a sequence of words, at each timestep  $t$  the conditional distribution over possible output words is given by:

$$p(y_t | y_{1:t-1}) = \text{softmax}(W_p \mathbf{h}_t^2 + \mathbf{b}_p). \quad (93)$$

The distribution over complete output sequences is calculated as the product of conditional distributions:

$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1}). \quad (94)$$

To reduce exposure bias and metric mismatching in sequential training, notable efforts are made to optimize non-differentiable metrics using reinforcement learning [68, 104, 144]. To further boost accuracy, detected semantic concepts [29, 139, 150] are adopted in captioning framework. A more structured representation over concepts calling scene graph is further explored [148, 149] in image captioning, which can take advantage of detected objects and their relationships. Instead of using a fully detected scene graph to improve captioning accuracy, Chen et al. [14] propose to employ Abstract Scene Graph as control signal to generate intention-aware and diverse image captions.

With the progress of multimodal representation learning, Wang et al. [132] propose the OFA, a unified multimodal pretrained model that can be applied in all modalities and various tasks. The OFA is simple yet effective, and it achieves new state-of-the-art performance on the kinds of multimodal tasks, such as image captioning, text-to-image generation, and VQA.

**6.3.2 Video Caption.** The aim of video caption is to describe or summarize a video in natural language. It is a non-trivial task for computers, since it is difficult to select the useful visual features from a video clip and describe what is happening in a way that obeys the common sense of humanity.

The currently prevailing architecture for video caption is composed of a CNN-like visual encoder and an RNN-like linguistic decoder. Donahue et al. [22] design the **Long-term Recurrent Convolutional Networks (LRCNs)** that are both temporally and spatially deep. After that, Venugopalan et al. [128] introduce the S2VT, a Seq2Seq approach for video to text with the knowledge from text corpora. Krishna et al. [49] propose a captioning module that uses contextual information from past and future events to capture the dependencies between the events in a video. However, Krishna et al. [49] fail to take advantage of language to benefit event proposal with the co-training diagram. Thus, Zhou et al. [160] propose a video caption framework that produces proposal and description simultaneously. To describe a video with multiple fine-grained actions, Wang et al. [131] propose a hierarchical reinforcement learning framework that a high-level agent learns to design sub-goals and a low-level worker recognizes the primitive actions to fulfill the sub-goal. By introducing a discriminator to evaluate sentences' visual relevance to the video, language diversity & fluency, and coherence across sentences, Park et al. [88] generate more accuracy video descriptions. To resolve the dilemma that encoders of vision-language tasks are not trained end-to-end, Lei et al. [52] propose ClipBERT, a framework that applies the sparse sampling to use a few sampled clips to achieve better performance. Additionally, with the rise of multimodal learning, Luo et al. [73] propose UniVL, a unified multimodal pre-training model for video captioning. The UniVL consists of four components (i.e., two encoders for single-modal, a cross-modal encoder, and a decoder) and is pre-trained with five tasks, including language understanding and generation tasks. The highlight of UniVL is that it uses both understanding and generative tasks for cross-modal pre-training, leading to its state-of-the-art performance on video captioning.

In recent years, self-supervised learning has become increasingly important with its power to leverage the abundance of unlabeled data. Sun et al. [121] propose VideoBERT to learn bidirectional joint distributions over sequences of visual and linguistic tokens without any explicit supervision. Zhu et al. [161] propose ActNet, which models global and local visual cues for fine-grained visual and linguistic relation learning.

**6.3.3 Visual Storytelling.** Visual Storytelling not only needs to identify the correlation between objects in a single picture, but also needs to identify and learn the logical relationship between consecutive sequential images. In practice, Visual Storytelling is prone to problems such as single narrative words, rigid sentences, and content with incoherent logic. Huang et al. [42] propose to generate a coherent and reasonable story with a series of images. To deal with the issue that Visual Storytelling usually focuses on generating general description rather than the details of meaningful visual contents, Li et al. [55] propose to mine the cross-modal rules to assist the concept inference. Yang et al. [147] present a commonsense-driven generative model to introduce crucial commonsense from the external knowledge base for visual storytelling. Due to the limitation of maximum likelihood estimation on training, the majority of existing models encourage high resemblance to texts in the training database, which makes the description overly rigid and lack in diverse expressions. Therefore, Mo et al. [81] cast the task as a reinforcement learning task and propose an **Adversarial All-in-one Learning (AAL)** framework to learn a reward model, which simultaneously incorporates the information of all images in the photo stream and all texts in the paragraph and optimizes a generative model with the estimated reward. To make the Visual Storytelling model topic adaptively, Li et al. [55] introduce a gradient-based meta-learning algorithm. Conventional storytelling approaches usually focus on optimizing metrics such as BLEU, ROUGE, and CIDEr. In their paper, Hu et al. [40] revisit the issue from a different perspective by delving into what defines a natural and thematically coherent story. In addition, considering the inability of previous methods to explore latent information beyond the image and thus fail to capture consistent dependencies from the global representation, Li et al. [56] propose the KAGS, a

knowledge-enriched attention network with group-wise semantic model that achieves new state-of-the-art performance with respect to both objective and subjective evaluation metrics.

## 7 NLG EVALUATION METRICS

In the research field of Artificial Intelligence, the evaluation metrics for models of kinds of tasks have always been the focus of attention for a long time, and the same is true in the fields of NLP [119]. In this section, we mainly introduce several automatic evaluation metrics for NLG, which can be divided into two categories: untrained evaluation metrics and machine-learned evaluation metrics [12].

### 7.1 Untrained Evaluation Metrics

This category of metric is most widely used in the NLG community, since it is easy to be implemented and does not involve additional training cost, which compares machine-generated texts to human-generated ones simply based on content overlap, string distance or lexical diversity. We mainly introduce five metrics of such category, including BLEU, ROUGE, METEOR, Distinct, and Self-BLEU.

The **Bilingual Evaluation Understudy (BLEU)** metric [87] is used to calculate the co-occurrence frequency of two sentences based on the weighted average of matched n-gram phrases. BLEU was originally used to evaluate machine translation, and has been used for more and more NLG tasks, such as question generation [158], topic-to-essay generation [146], text style transfer [72], and dialogue generation [4, 64].

The **Recall-Oriented Understudy for Gisting Evaluation (ROUGE)** metric [62] is used to measure the similarity between the generated and reference texts based on the recall score. This metric is commonly used in the field of text summarization, including four types: ROUGE-n measures the n-gram co-occurrence statistics; ROUGE-l measures the longest common subsequence; ROUGE-w measures the weighted longest common subsequence; ROUGE-s measures the skip-bigram co-occurrence statistics. ROUGE has also been widely applied to other NLG tasks such as question generation [158], distractor generation [98], and dialogue generation [4].

The **Metric for Evaluation of Translation with Explicit Ordering (METEOR)** metric [3] is an improvement over BLEU to address several weaknesses including four aspects: lack of recall, use of higher order n-grams, lack of explicit word-matching between translation and reference, and use of geometric averaging of n-grams, which is calculated by the harmonic mean of the unigram precision and recall. In addition to machine translation, METEOR has been widely used in text summarization [96], question generation [158], and dialogue generation [4].

The Distinct metric [54] is used to measure the diversity of response sequences for dialogue generation. It calculates the number of distinct unigrams and bigrams in generated responses to reflect the diversity degree. To avoid preference for long sequences, the value is scaled by the total number of generated tokens.

The Self-BLEU metric [162] is also a metric to measure the diversity. Different from BLEU that only evaluates the similarity between two sentences, Self-BLEU is used to measure the resemblance degree between one sentence (hypothesis) and the rest sentences (reference) in a generated collection. It first calculates the BLEU score of every generated sentence against other sentences, then the average BLEU score is defined as the Self-BLEU score of the document, where a lower Self-BLEU score implies higher diversity.

### 7.2 Machine-learned Evaluation Metrics

This category of metric is based on machine-learned models to simulate human judges, which evaluates the similarity between machine-generated texts or between machine-generated texts and

human-generated ones. We mainly introduce three metrics of such category, containing ADEM, BLEURT, and BERTScore.

The **Automatic Dialogue Evaluation Model (ADEM)** metric [69] is used to automatically evaluate the quality of dialogue responses, where the evaluation model is trained in a semi-supervised manner with a hierarchical **recurrent neural network (RNN)** to predict the response scores. Specifically, given the dialogue context  $\mathbf{c}$ , model response  $\hat{\mathbf{r}}$ , and reference response  $\mathbf{r}$  encoded by a hierarchical RNN, the predicted score can be calculated by:

$$score = (\mathbf{c}^\top M \hat{\mathbf{r}} + \mathbf{r}^\top N \hat{\mathbf{r}} - \alpha) / \beta, \quad (95)$$

where  $M, N$  are learnable matrices initialized by the identity, and  $\alpha, \beta$  are scalar constants to initialize the predicted scores in range  $[1, 5]$ .

The **Bilingual Evaluation Understudy with Representations from Transformers (BLEURT)** metric [111] is based on BERT [18] with a novel pre-training scheme. Before fine-tuning BERT on rating data to predict human rating scores, a pre-training method is applied, where BERT is pre-trained on a large number of synthetic sentence pairs on several lexical- and semantic-level supervision signals in a multi-task manner. This pre-training process is important and can improve the robustness to quality drifts of generation systems.

The BERTScore metric [156] uses pre-trained contextual embeddings from BERT to measure the similarity between two sentences. Given the contextual embeddings of a reference sentence  $\mathbf{x}$  and a candidate sentence  $\hat{\mathbf{x}}$ , namely,  $\mathbf{x}, \hat{\mathbf{x}}$ , the recall, precision, and F1 scores are calculated by:

$$R_{\text{BERT}} = \frac{1}{|\mathbf{x}|} \sum_{x_i \in \mathbf{x}} \max_{\hat{x}_j \in \hat{\mathbf{x}}} \mathbf{x}_i^\top \hat{\mathbf{x}}_j, \quad (96)$$

$$P_{\text{BERT}} = \frac{1}{|\hat{\mathbf{x}}|} \sum_{\hat{x}_j \in \hat{\mathbf{x}}} \max_{x_i \in \mathbf{x}} \mathbf{x}_i^\top \hat{\mathbf{x}}_j, \quad (97)$$

$$F_{\text{BERT}} = 2 \frac{P_{\text{BERT}} \cdot R_{\text{BERT}}}{P_{\text{BERT}} + R_{\text{BERT}}}, \quad (98)$$

where the recall is calculated by matching each token in  $\mathbf{x}$  to a token in  $\hat{\mathbf{x}}$ , the precision is obtained by matching each token in  $\hat{\mathbf{x}}$  to a token  $\mathbf{x}$ , and greedy matching is adopted to match the most similar tokens.

### 7.3 Human Evaluation Metrics

For generation, human evaluation focuses on the explanation of two key matters: diversity and creativity, i.e., the capacity of varying their texts in form and emphasis to fit an enormous range of speaking situations, and the potential to express any object or relation as a natural language text. In further detail, human evaluation is implemented to evaluate on three aspects: Grammar (whether a generated sentence is fluent without grammatical error), Faithful (whether the output is faithful to input), and Coherent (whether a sentence is logically coherent and the order of expression is in line with human writing habits). This needs to organize the capabilities of the people who work on generation in field of computational linguistics and artificial intelligence.

## 8 PROBLEMS AND CHALLENGES

In this section, we primarily point out four problems and challenges that deserve to be tackled and investigated further, including the evaluation method, external knowledge engagement, controllable generation, and multimodal scenarios.

*Evaluation Method.* Evaluation method is still an important and open research area for the field of NLG. As pointed by Reference [17], traditional untrained evaluation metrics do not always correlate well with human judgements, while recent machine-learned metrics need a large amount of human annotations and not always have good transferability. Hence, there still exists a significant amount of challenges and improvement room in this area.

*External Knowledge Engagement.* Considering the limited information lying in the original texts and the difficulty of generating satisfying sentences [151], it is crucial to incorporate external knowledge to enhance the performance. Therefore, how to obtain useful and correlative knowledge and how to effectively incorporate the knowledge still deserve to be investigated.

*Controllable Generation.* Another challenging problem is how to generate controllable natural language as we would like it to be. Although a great body of work has been done in this area to study how to perform various kinds of controlled text generation, there is still a lack of uniform paradigms and standards about it. More importantly, how to measure the controllability of the generated text remains an open question, for different controlled contents.

*Multimodal Scenarios.* Recently, research on various applications in multimodal scenarios have gradually attracted more and more attention from NLP researchers. How to apply natural language generation methods in multimodal scenarios has been a worthy problem and promising direction. It is reasonable to believe that the utilization of rich multimodal information into natural language generation tasks will surely further advance the progress and development in this direction.

## 9 CONCLUSIONS

Over the past few years, natural language generation tasks and methods have become important and indispensable in natural language processing. This progress owes to advances in various deep learning-based methods. This article describes deep learning research on natural language generation with a historical perspective, emphasizing the special character of the problems to be solved. It begins by contrasting generation with language understanding, establishing basic concepts about the tasks, datasets, and the deep learning methods through it. A section of evaluation metrics from the output of generation systems follows, showing what kinds of performance are possible and where the difficulties are. Finally, some open problems are suggested to indicate the major challenges and future research directions of natural language generation.

## REFERENCES

- [1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. VQA: Visual question answering. In *Proceedings of the IEEE International Conference on Computer Vision*. 2425–2433.
- [3] Satyanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*.
- [4] Siqi Bao, Huang He, Fan Wang, Hua Wu, and Haifeng Wang. 2020. PLATO: Pre-trained dialogue generation model with discrete latent variable. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- [5] Rahul Bhatia, Vishakha Gautam, and Yash Garg. 2019. Dynamic question answer generator: An enhanced approach to question generation. *Int. J. Trend Scient. Res. Devel.* 3, 4 (2019).
- [6] Bodo Billerbeck, Falk Scholer, Hugh E. Williams, and Justin Zobel. 2003. Query expansion using associated queries. In *Proceedings of the 12th International Conference on Information and Knowledge Management*.
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom



- Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [8] Tian Cai, Mengjun Shen, Huailiang Peng, Lei Jiang, and Qiong Dai. 2019. Improving transformer with sequential context representations for abstractive text summarization. In *Natural Language Processing and Chinese Computing*. Association for Computational Linguistics.
- [9] Meng Cao, Yue Dong, Jiapeng Wu, and Jackie Chi Kit Cheung. 2020. Factual error correction for abstractive summarization models. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [10] Qingxing Cao, Xiaodan Liang, Bailin Li, and Liang Lin. 2019. Interpretable visual question answering by reasoning on dependency trees. *IEEE Trans. Pattern Anal. Mach. Intell.* 43, 3 (2019), 887–901.
- [11] Qingxing Cao, Wentao Wan, Keze Wang, Xiaodan Liang, and Liang Lin. 2021. Linguistically routing capsule network for out-of-distribution visual question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1614–1623.
- [12] Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. 2021. Evaluation of text generation: A survey. *arXiv preprint arXiv:2006.14799*.
- [13] Liqun Chen, Shuyang Dai, Chenyang Tao, Haichao Zhang, Zhe Gan, Dinghan Shen, Yizhe Zhang, Guoyin Wang, Ruiyi Zhang, and Lawrence Carin. 2018. Adversarial text generation via feature-mover's distance. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [14] Shizhe Chen, Qin Jin, Peng Wang, and Qi Wu. 2020. Say as you wish: Fine-grained control of image caption generation with abstract scene graphs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [15] Jan K. Chorowski, Dzmitry Bahdanau, Dzmitry Serdyuk, Kyunghyun Cho, and Yoshua Bengio. 2015. Attention-based models for speech recognition. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [16] Ho-Lam Chung, Ying-Hong Chan, and Yao-Chung Fan. 2020. A BERT-based distractor generation scheme with multi-tasking and negative answer training strategies. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, Trevor Cohn, Yulan He, and Yang Liu (Eds.). 4390–4400.
- [17] Jan Deriu, Alvaro Rodrigo, Arantxa Otegi, Guillermo Echevoyen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. 2021. Survey on evaluation methods for dialogue systems. *Artificial Intelligence Review* 54, 1 (2021), 755–810.
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- [19] Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. Wizard of Wikipedia: Knowledge-powered conversational agents. In *Proceedings of the International Conference on Learning Representations*.
- [20] Jianli Ding, Yang Li, Huiyu Ni, and Zhengquan Yang. 2020. Generative text summary based on enhanced semantic attention and gain-benefit gate. *IEEE Access* 8 (2020).
- [21] Tuong Do, Binh X. Nguyen, Erman Tjiputra, Minh Tran, Quang D. Tran, and Anh Nguyen. 2021. Multiple meta-model quantifying for medical visual question answering. In *Proceedings of the International Conference on Medical Image Computing and Computer-assisted Intervention*. 64–74.
- [22] Jeff Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Trevor Darrell, and Kate Saenko. 2015. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- [23] Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*. 13042–13054.
- [24] Yue Dong, Shuohang Wang, Zhe Gan, Yu Cheng, Jackie Chi Kit Cheung, and Jingjing Liu. 2020. Multi-fact correction in abstractive text summarization. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [25] Yu Duan, Canwen Xu, Jiaxin Pei, Jialong Han, and Chenliang Li. 2020. Pre-train and plug-in: Flexible conditional text generation with variational auto-encoders. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 253–262.
- [26] Bambang Dw, Yaya Heryadi, Hapnes Toba, and Widodo Budiharto. 2021. Question generation model based on keyphrase, context-free grammar, and Bloom's taxonomy. *Educ. Inf. Technol.* 26, 2 (2021).
- [27] Xiaocheng Feng, Ming Liu, Jiahao Liu, Bing Qin, Yibo Sun, and Ting Liu. 2018. Topic-to-essay generation with neural networks. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*.
- [28] Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, and Rui Yan. 2018. Style transfer in text: Exploration and evaluation. In *Proceedings of the AAAI Conference on Artificial Intelligence*.



- [29] Zhe Gan, Chuang Gan, Xiaodong He, Yunchen Pu, Kenneth Tran, Jianfeng Gao, Lawrence Carin, and Li Deng. 2017. Semantic compositional networks for visual captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- [30] Yifan Gao, Lidong Bing, Piji Li, Irwin King, and Michael R. Lyu. 2019. Generating distractors for reading comprehension questions from real examinations. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [31] Albert Gatt and Emiel Krahmer. 2018. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *J. Artif. Intell. Res.* 61, 1 (2018).
- [32] Spandana Gella, Mike Lewis, and Marcus Rohrbach. 2018. A dataset for telling the stories of social media videos. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [33] Haifan Gong, Guanqi Chen, Sishuo Liu, Yizhou Yu, and Guanbin Li. 2021. Cross-modal self-attention with multi-task pre-training for medical visual question answering. In *Proceedings of the International Conference on Multimedia Retrieval*. 456–460.
- [34] Haifan Gong, Guanqi Chen, Mingzhi Mao, Zhen Li, and Guanbin Li. 2022. VQAMix: Conditional triplet mixup for medical visual question answering. *IEEE Trans. Med. Imag.* (2022).
- [35] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [36] David Graff and Christopher Cieri. 2003. English Gigaword. (2003). LDC2003T05. Web Download. Philadelphia: Linguistic Data Consortium. DOI: <https://doi.org/10.35111/0z6y-q265>
- [37] Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, and Jun Wang. 2018. Long text generation via adversarial training with leaked information. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*.
- [38] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [39] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [40] Junjie Hu, Yu Cheng, Zhe Gan, Jingjing Liu, Jianfeng Gao, and Graham Neubig. 2020. What makes a good story? Designing composite rewards for visual storytelling. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [41] Mingqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [42] Ting-Hao Kenneth Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh, Lucy Vanderwende, Michel Galley, and Margaret Mitchell. 2016. Visual storytelling. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- [43] Harsh Jhamtani, Varun Gangal, Eduard Hovy, and Eric Nyberg. 2017. Shakespearizing modern language using copy-enriched sequence to sequence models. In *Proceedings of the Workshop on Stylistic Variation*.
- [44] Xin Jia, Wenjie Zhou, Xu Sun, and Yunfang Wu. 2020. How to ask good questions? Try to leverage paraphrases. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- [45] Bin Jiang, Jingxu Yang, Chao Yang, Wanyue Zhou, Liang Pang, and Xiaokang Zhou. 2020. Knowledge augmented dialogue generation with divergent facts selection. *Knowl.-based Syst.* 210 (2020).
- [46] Yanghoon Kim, Hwanhee Lee, Joongbo Shin, and Kyomin Jung. 2019. Improving neural question generation using answer separation. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [47] Thomas N. Kipf and Max Welling. 2016. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*.
- [48] Thomas N. Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *Proceedings of the International Conference on Learning Representations*.
- [49] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. 2017. Dense-captioning events in videos. In *Proceedings of the IEEE International Conference on Computer Vision*.
- [50] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *Int. J. Comput. Vis.* 123, 1 (2017).
- [51] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. RACE: Large-scale ReAding comprehension dataset from examinations. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [52] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. 2021. Less is more: ClipBERT for video-and-language learning via sparse sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [53] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language

- generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- [54] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
  - [55] Jiacheng Li, Haizhou Shi, Siliang Tang, Fei Wu, and Yueting Zhuang. 2019. Informative visual storytelling with cross-modal rules. In *Proceedings of the 27th ACM International Conference on Multimedia*.
  - [56] Tengpeng Li, Hanli Wang, Bin He, and Chang Wen Chen. 2022. Knowledge-enriched attention network with group-wise semantic for visual storytelling. arXiv preprint arXiv:2203.05346.
  - [57] Xu Li and Jinghua Zhu. 2020. Attention and graph matching network for retrieval-based dialogue system with domain knowledge. In *Proceedings of the International Joint Conference on Neural Networks*.
  - [58] Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. DailyDialog: A manually labelled multi-turn dialogue dataset. In *Proceedings of the 8th International Joint Conference on Natural Language Processing*.
  - [59] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard S. Zemel. 2016. Gated graph sequence neural networks. In *Proceedings of the International Conference on Learning Representations*.
  - [60] Rongzhong Lian, Min Xie, Fan Wang, Jinhua Peng, and Hua Wu. 2019. Learning to select knowledge for response generation in dialog systems. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*.
  - [61] Chen Liang, Xiao Yang, Neisarg Dave, Drew Wham, Bart Pursel, and C. Lee Giles. 2018. Distractor generation for multiple choice questions using learning to rank. In *Proceedings of the 13th Workshop on Innovative Use of NLP for Building Educational Applications*.
  - [62] Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*. Association for Computational Linguistics.
  - [63] Shuai Lin, Wentao Wang, Zichao Yang, Xiaodan Liang, Frank F. Xu, Eric Xing, and Zhiting Hu. 2020. Data-to-text generation with style imitation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. 1589–1598.
  - [64] Shuai Lin, Pan Zhou, Xiaodan Liang, Jianheng Tang, Ruihui Zhao, Ziliang Chen, and Liang Lin. 2021. Graph-evolving meta-learning for low-resource medical dialogue generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 13362–13370.
  - [65] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: Common objects in context. In *Proceedings of the European Conference on Computer Vision*.
  - [66] Bang Liu, Mingjun Zhao, Di Niu, Kunfeng Lai, Yancheng He, Haojie Wei, and Yu Xu. 2019. Learning to generate questions by learning what not to generate. In *Proceedings of the the World Wide Web Conference*.
  - [67] Qian Liu, Yihong Chen, Bei Chen, Jian-Guang Lou, Zixuan Chen, Bin Zhou, and Dongmei Zhang. 2020. You impress me: Dialogue generation via mutual persona perception. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 1417–1427.
  - [68] Siqi Liu, Zhenhai Zhu, Ning Ye, Sergio Guadarrama, and Kevin Murphy. 2017. Improved image captioning via policy gradient optimization of SPiDER. In *Proceedings of the IEEE International Conference on Computer Vision*.
  - [69] Ryan Lowe, Michael Noseworthy, Iulian Vlad Serban, Nicolas Angelard-Gontier, Yoshua Bengio, and Joelle Pineau. 2017. Towards an automatic Turing test: Learning to evaluate dialogue responses. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*.
  - [70] Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. The Ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems. In *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*.
  - [71] Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. 2017. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
  - [72] Fuli Luo, Peng Li, Jie Zhou, Pengcheng Yang, Baobao Chang, Xu Sun, and Zhifang Sui. 2019. A dual reinforcement learning framework for unsupervised text style transfer. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*.
  - [73] Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Xilin Chen, and Ming Zhou. 2020. UniViLM: A unified video and language pre-training model for multimodal understanding and generation. *arXiv preprint arXiv:2002.06353*.
  - [74] Liangchen Luo, Jingjing Xu, Junyang Lin, Qi Zeng, and Xu Sun. 2018. An auto-encoder matching model for learning utterance-level semantic dependency in dialogue generation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
  - [75] Sruthy Manmadhan and Binsu C. Koor. 2020. Visual question answering: A state-of-the-art review. *Artif. Intell. Rev.* 53, 8 (2020), 5705–5745.

- [76] Kaushal Kumar Maurya and Maunendra Sankar Desarkar. 2020. Learning to distract: A hierarchical multi-decoder network for automated generation of long distractors for multiple-choice questions for reading comprehension. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*.
- [77] Julian John McAuley and Jure Leskovec. 2013. From amateurs to connoisseurs: Modeling the evolution of user expertise through online reviews. In *Proceedings of the 22nd International Conference on World Wide Web*.
- [78] David D. McDonald. 2010. Natural language generation. *Handb. Nat. Lang. Process.* (2010).
- [79] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S. Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [80] Remi Mir, Bjarke Felbo, Nick Obradovich, and Iyad Rahwan. 2019. Evaluating style transfer for text. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- [81] Lingbo Mo, Chunhong Zhang, Yang Ji, and Zheng Hu. 2019. Adversarial learning for visual storytelling with sense group partition. In *Proceedings of the Asian Conference on Computer Vision*.
- [82] Aditya Mogadale, Marimuthu Kalimuthu, and Dietrich Klakow. 2020. Trends in integration of vision and language research: A survey of tasks, datasets, and methods. *arXiv preprint arXiv:1907.09358*.
- [83] Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! Topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [84] Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A human generated MACHine reading COmprehension dataset. In *Proceedings of the Workshop on Cognitive Computation: Integrating Neural and Symbolic Approaches Co-located with the 30th Annual Conference on Neural Information Processing Systems*.
- [85] Cicero Nogueira dos Santos, Igor Melnyk, and Inkit Padhi. 2018. Fighting offensive language on social media with unsupervised text style transfer. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*.
- [86] Daniel Otter, Julian Medina, and Jugal Kalita. 2020. A survey of the usages of deep learning for natural language processing. *IEEE Trans. Neural Netw. Learn. Syst.* 32, 2 (2020).
- [87] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*.
- [88] Jae Sung Park, Marcus Rohrbach, Trevor Darrell, and Anna Rohrbach. 2019. Adversarial inference for multi-sentence video description. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [89] Rakesh Patra and Sujana Kumar Saha. 2019. A hybrid approach for automatic generation of named entity distractors for multiple choice questions. *Educ. Inf. Technol.* 24, 2 (2019).
- [90] Romain Paulus, Caiming Xiong, and Richard Socher. 2018. A deep reinforced model for abstractive summarization. In *Proceedings of the International Conference on Learning Representations*.
- [91] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [92] Rivindu Perera and Parma Nand. 2017. Recent advances in natural language generation: A survey and classification of the empirical literature. *Comput. Inform.* 36, 1 (2017).
- [93] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- [94] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2015. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE International Conference on Computer Vision*.
- [95] Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, and Alan W. Black. 2018. Style transfer through back-translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*.
- [96] Weizhen Qi, Yu Yan, Yeyun Gong, Dayiheng Liu, Nan Duan, Jiusheng Chen, Ruofei Zhang, and Ming Zhou. 2020. ProphetNet: Predicting future N-gram for Sequence-to-Sequence Pre-training. In *Findings of the Association for Computational Linguistics: EMNLP 2020*.
- [97] Lin Qiao, Jianhao Yan, Fandong Meng, Zhendong Yang, and Jie Zhou. 2020. A sentiment-controllable topic-to-essay generator with topic knowledge graph. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing: Findings, EMNLP 2020*.
- [98] Zhaopeng Qiu, Xian Wu, and Wei Fan. 2020. Automatic distractor generation for multiple choice questions in standard tests. In *Proceedings of the 28th International Conference on Computational Linguistics*.

- [99] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.* 21, 1 (2020).
- [100] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [101] Sudha Rao and Joel Tetreault. 2018. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- [102] Ehud Reiter and Robert Dale. 1997. Building applied natural language generation systems. *Nat. Lang. Eng.* 3, 1 (1997).
- [103] Siyu Ren and Kenny Q. Zhu. 2021. Knowledge-driven distractor generation for cloze-style multiple choice questions. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*.
- [104] Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- [105] Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [106] Yuri Safovich and Amos Azaria. 2020. Fiction sentence expansion and enhancement via focused objective and novelty curve sampling. In *Proceedings of the IEEE 32nd International Conference on Tools with Artificial Intelligence*.
- [107] Evan Sandhaus. 2008. The New York Times Annotated Corpus. (2008). LDC2008T19. Web Download. Philadelphia: Linguistic Data Consortium. DOI : <https://doi.org/10.35111/77ba-9x74>
- [108] Sashank Santhanam and Samira Shaikh. 2019. A survey of natural language generation techniques with a focus on dialogue systems—Past, present and future directions. *arXiv preprint arXiv:1906.00500*.
- [109] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2009. The graph neural network model. *IEEE Trans. Neural Netw.* (2009).
- [110] Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*.
- [111] Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- [112] Iulian Vlad Serban, Alberto García-Durán, Caglar Gulcehre, Sungjin Ahn, Sarath Chandar, Aaron Courville, and Yoshua Bengio. 2016. Generating factoid questions with recurrent neural networks: The 30M factoid question-answer corpus. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*.
- [113] Min Shi, Yufei Tang, and Jianxun Liu. 2019. Functional and contextual attention-based LSTM for service recommendation in mashup creation. *IEEE Trans. Parallel Distrib. Syst.* 39, 9 (2019).
- [114] Gunnar A. Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. 2016. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *Proceedings of the European Conference on Computer Vision*.
- [115] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [116] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. MASS: Masked sequence to sequence pre-training for language generation. In *Proceedings of the 36th International Conference on Machine Learning*.
- [117] Milan J. Srinivas, Michelle M. Roy, and Viraj Kumar. 2019. Towards generating plausible distractors for code comprehension multiple-choice questions. In *Proceedings of the IEEE 10th International Conference on Technology for Education*.
- [118] Katherine Stasaski and Marti A. Hearst. 2017. Multiple choice question generation utilizing an ontology. In *Proceedings of the 12th Workshop on Innovative Use of NLP for Building Educational Applications*.
- [119] Keh-Yih Su, Ming-Wen Wu, and Jing-Shin Chang. 1992. A new quantitative quality measure for machine translation systems. In *Proceedings of the 14th International Conference on Computational Linguistics*. 433–439.
- [120] Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, and Rob Fergus. 2015. End-to-end memory networks. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [121] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. 2019. VideoBERT: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.
- [122] Xingwu Sun, Jing Liu, Yajuan Lyu, Wei He, Yanjun Ma, and Shi Wang. 2018. Answer-focused and position-aware neural question generation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [123] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.



- [124] Jian Tang, Yue Wang, Kai Zheng, and Qiaozhu Mei. 2017. End-to-end learning for short text expansion. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.
- [125] Alexey Tikhonov, Viacheslav Shibaev, Aleksander Nagaev, Aigul Nugmanova, and Ivan P. Yamshchikov. 2019. Style transfer for texts: To err is human, but error margins matter. *arXiv preprint arXiv:1908.06809*.
- [126] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Proceedings of the Conference on Advances in Neural Information Processing Systems*.
- [127] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph attention networks. In *Proceedings of the International Conference on Learning Representations*.
- [128] Subhashini Venugopalan, Lisa Anne Hendricks, Raymond Mooney, and Kate Saenko. 2016. Improving LSTM-based video description with linguistic knowledge mined from text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [129] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- [130] Minh H. Vu, Tommy Löfstedt, Tufve Nyholm, and Raphael Sznitman. 2020. A question-centric model for visual question answering in medical imaging. *IEEE Trans. Med. Imag.* 39, 9 (2020), 2856–2868.
- [131] Li Wang, Junlin Yao, Yunzhe Tao, Li Zhong, Wei Liu, and Qiang Du. 2018. A reinforced topic-aware convolutional sequence-to-sequence model for abstractive text summarization. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*.
- [132] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. *arXiv preprint arXiv:2202.03052*.
- [133] Siyuan Wang, Zhongyu Wei, Zhihao Fan, Zengfeng Huang, Weijian Sun, Qi Zhang, and Xuanjing Huang. 2020. PathQG: Neural question generation from facts. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [134] Siyuan Wang, Zhongyu Wei, Zhihao Fan, Yang Liu, and Xuanjing Huang. 2019. A multi-agent communication framework for question-worthy phrase extraction and question generation. *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [135] Wei Wang, Hai-Tao Zheng, and Zibo Lin. 2020. Self-attention and retrieval enhanced neural networks for essay generation. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*.
- [136] Yanmeng Wang, Wenge Rong, Jianfei Zhang, Yuanxin Ouyang, and Zhang Xiong. 2020. Knowledge grounded pre-trained model for dialogue response generation. In *Proceedings of the International Joint Conference on Neural Networks*.
- [137] Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions. In *Proceedings of the 3rd Workshop on Noisy User-generated Text*.
- [138] Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. 2019. TransferTransfo: A transfer learning approach for neural network based conversational agents. *arXiv preprint arXiv:1901.08149*.
- [139] Qi Wu, Chunhua Shen, Lingqiao Liu, Anthony Dick, and Anton Van Den Hengel. 2016. What value do explicit high level concepts have in vision to language problems? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- [140] Siying Wu, Zheng-Jun Zha, Zilei Wang, Houqiang Li, and Feng Wu. 2019. Densely supervised hierarchical policy-value network for image paragraph generation. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*.
- [141] Dongling Xiao, Han Zhang, Yu-Kun Li, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. 2020. ERNIE-GEN: An enhanced multi-flow pre-training and fine-tuning framework for natural language generation. In *Proceedings of the 29th International Joint Conference on Artificial Intelligence*. 3997–4003.
- [142] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. MSR-VTT: A large video description dataset for bridging video and language. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- [143] Jingjing Xu, Xu Sun, Qi Zeng, Xiaodong Zhang, Xuancheng Ren, Houfeng Wang, and Wenjie Li. 2018. Unpaired sentiment-to-sentiment translation: A cycled reinforcement learning approach. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*.
- [144] Ning Xu, Hanwang Zhang, An-An Liu, Weizhi Nie, Yuting Su, Jie Nie, and Yongdong Zhang. 2020. Multi-level policy and reward-based deep reinforcement learning framework for image captioning. *IEEE Trans. Multim.* 22, 5 (2020).
- [145] Zhen Xu, Bingquan Liu, Baoxun Wang, Chengjie Sun, Xiaolong Wang, Zhuoran Wang, and Chao Qi. 2017. Neural response generation via GAN with an approximate embedding layer. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.

- [146] Pengcheng Yang, Lei Li, Fuli Luo, Tianyu Liu, and Xu Sun. 2019. Enhancing topic-to-essay generation with external commonsense knowledge. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- [147] Pengcheng Yang, Fuli Luo, Peng Chen, Lei Li, Zhiyi Yin, Xiaodong He, and Xu Sun. 2019. Knowledgeable storyteller: A commonsense-driven generative model for visual storytelling. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*.
- [148] Xu Yang, Kaihua Tang, Hanwang Zhang, and Jianfei Cai. 2019. Auto-encoding scene graphs for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [149] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. 2018. Exploring visual relationship for image captioning. In *Proceedings of the European Conference on Computer Vision*.
- [150] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. 2016. Image captioning with semantic attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- [151] Wenhao Yu, Chenguang Zhu, Zaitang Li, Zhiting Hu, Qingyun Wang, Heng Ji, and Meng Jiang. 2021. A survey of knowledge-enhanced text generation. *arXiv preprint arXiv:2010.04389*.
- [152] Chenhan Yuan, Yi-Chin Huang, and Cheng-Hung Tsai. 2019. Efficient text generation of user-defined topic using generative adversarial networks. In *Proceedings of the 4th Workshop on Computational Creativity in Language Generation*.
- [153] Xingdi Yuan, Tong Wang, Caglar Gulcehre, Alessandro Sordani, Philip Bachman, Saizheng Zhang, Sandeep Subramanian, and Adam Trischler. 2017. Machine comprehension by text-to-text neural question generation. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*.
- [154] Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. PEGASUS: Pre-training with extracted gap-sentences for abstractive summarization. In *Proceedings of the 37th International Conference on Machine Learning*.
- [155] Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*.
- [156] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating text generation with BERT. In *Proceedings of the International Conference on Learning Representations*.
- [157] Junbo Zhao, Yoon Kim, Kelly Zhang, Alexander Rush, and Yann LeCun. 2018. Adversarially regularized autoencoders. In *Proceedings of the 35th International Conference on Machine Learning*.
- [158] Yao Zhao, Xiaochuan Ni, Yuanyuan Ding, and Qifa Ke. 2018. Paragraph-level neural question generation with maxout pointer and gated self-attention networks. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- [159] Ganbin Zhou, Ping Luo, Rongyu Cao, Fen Lin, Bo Chen, and Qing He. 2017. Mechanism-aware neural machine for dialogue response generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [160] Luowei Zhou, Yingbo Zhou, Jason J. Corso, Richard Socher, and Caiming Xiong. 2018. End-to-end dense video captioning with masked transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [161] Linchao Zhu and Y. Yang. 2020. ActBERT: Learning global-local video-text representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- [162] Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. TexyGen: A benchmarking platform for text generation models. In *Proceedings of the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*.

Received 10 August 2021; revised 2 July 2022; accepted 25 July 2022