CAVSS: A Commonsense Augmented Variational Sequence to Sequence Model for Language Generation

Anonymous ACL submission

Abstract

1

Commonsense knowledge as external 2 knowledge enhances the semantic under-3 standing of the input sequences of the 4 model and is of guidance to text generation 5 models. In this paper, we propose a novel 6 approach of incorporating commonsense 7 knowledge for enhancing the performance 8 of end-to-end text generation models. 9 Firstly, given an input sequence and retriev-10 ing the relevant knowledge triples, the em-11 bedding of the commonsense knowledge 12 and the context vector encoded in the en-13 coder part are spliced for sampling, allow-14 ing the prior distribution to approximately 15 fit the posterior distribution to achieve the 16 selection of appropriate knowledge even 17 without posterior information. Then an au-18 toregressive transformation is applied to the 19 sampling to prevent the problem of too 20 slow fitting of simple Gaussian distribution, 21 and a new learning objective is designed in 22 the training phase to make this transformed 23 distribution fit the posterior distribution. In 24 addition, we perform variational operations 25 on the decoding part of the attention mech-26 anism to weaken the attention strength and 27 prevent reconstruction from playing a deci-28 sive role in generation while ignoring other 29 modules. Experiments show that our pro-30 posed model can generate more fluent and 31 significantly more diverse sentences, and 32 the contributions of each module to the 33 model are analyzed, achieving satisfactory 34 results. 35

36 1 Introduction

³⁷ Natural language generation (NLG) is one of the ⁷² from the Reconstruction term, and the mainstream
³⁸ main problems in natural language processing. ⁷³ method is mainly to add loss to let the latent vari³⁹ Currently, seq-to-seq (Sequence to Sequence) ⁷⁴ ables participate in the prediction directly, which
⁴⁰ model is the main method of text generation. ⁷⁵ alleviates the posterior collapse problem, but the
⁴¹ However, due to the flaws of model and the lack ⁷⁶ defects of model itself have not been solved (Zhao

⁴² of commonsense reasoning, seq-to-seq model
⁴³ may generates a lot of low-quality text with sim⁴⁴ ple repetition or factual errors (Brown et al.,
⁴⁵ 2020)(Bi et al., 2019).

In order to solve these problems, related work 46 47 firstly adopted VAE (Variational Autoencoder) 48 model as encoder to improve the diversity of the 49 generated text. However, posterior collapse, ⁵⁰ which is generally caused by the disappearance of 51 KL divergence, is a flaw in VAE. Posterior col-52 lapse is mainly prevented from two aspects, that 53 is, KL term and reconstruction. One resolution is 54 KL cost annealing (Bowman et al., 2016), which 55 is done by improving KL item by multiplying a ⁵⁶ weight coefficient on the KL which is allocated to 57 0 at the beginning and gradually increased to 1 58 during training, optimizing Reconstruction part ⁵⁹ with high priority, and gradually paying attention 60 to KL, can be seen as a gradient from AutoEn-61 coder to Variational AutoEncoder. Another resolu-62 tion is making the latent variables more flexibility 63 on the strength of reversible transformation (Chen 64 et al. 2016). Unfortunately, a pure Gaussian distri-65 bution is difficult to fit the distribution of real data 66 (Papamakarios et al., 2021), so some transfor-67 mations are needed to map the original Gaussian 68 space to a suitable feature representation space. 69 Mathematical transformations are difficult to sim-70 ulate or even cannot fit for data in different appli-71 cation scenarios. The other solution is to improve 72 from the Reconstruction term, and the mainstream 73 method is mainly to add loss to let the latent vari-74 ables participate in the prediction directly, which 75 alleviates the posterior collapse problem, but the

77 et al., 2017).

Secondly, related work introduces external 115 78 79 knowledge into the model to improve the factual 116 ⁸⁰ accuracy of the generated text and have made a lot 117 81 of progress (Zhu et al., 2017; Li et al., 2021; Yu et $_{\rm 82}$ al., 2020). However, the shortcomings of these $^{\rm 118}$ 2 83 works are that most of the external knowledge ⁸⁴ used is domain-specific rather than general com-85 monsense knowledge. At the same time, external ⁸⁶ knowledge is directly embedded as a whole in the 87 process of knowledge embedding and result in the 88 loss of semantics because relationship between 89 edges and nodes in the graph structure cannot be 90 established (Qiao et al., 2020).

In response to the above problems, this paper 91 92 proposes a new Commonsense knowledge Aug-93 mented Variational Seq-to-Seq generation model 94 CAVSS. The main contributions are:

To solve the problem of posterior collapse, we 95 1. propose an autoregressive sampling method 96 from the KL perspective by defining a new 97 conversion function with the help of asymp-98 totic property of the fully connected layer's 99 parameter approximation. It can train the sam-100 pling space of the fitted data, and a new learn-101 ing objective is designed to guide the direc-102 tion of model training. From the Reconstruc-103 tion perspective, the attention mechanism in 104 the Seq-to-Seq model is weakened to prevent 105 the model from bypassing other modules. 106

Using commonsense knowledge to enhance 107 2. 108 and designing a new attention mechanism 109 when embedding to make full use of the re-110 trieved knowledge graph structure. According 111 to the graph structure, nodes (entities) are con-112 nected by edges (relations) to form triples 113

such as (thunder, RelatedTo, shocking). Word "thunder" is related to "shocking", which helps the model to generate commonsense and diverse sequences.

Related Work

114

119 Sequence-to-Sequence (Seq2Seq) general model 120 has been successfully applied in the Generation 121 tasks (Sutskever et al., 2014), and the attention-122 fused Seq2Seq model greatly improves the quality 123 of text generation (Bahdanau et al.2015). Varia-124 tional AutoEncoder (VAE) is widely used in vari-125 ous tasks (generating summaries, dialogue sys-126 tems, etc.), and produces many improved versions. ¹²⁷ β-VAE can increase the quality of Reconstruction 128 by adding constraints in latent representation 129 space (Irina et al., 2016), and VAE guided by in-130 ternal knowledge (topic) and Householder Trans-131 formation can make the approximate posterior of 132 latent code highly flexible (Wang et al., 2019).

For those knowledge graphs that are con-134 structed based on data outside the input text (e.g., 135 ConceptNet), we refer to them as external 136 knowledge. Internal knowledge often refers to 137 keywords, subject word and other ways to pro-138 mote the generated text closely to the topic (Wei 139 et al., 2019; Wang et al., 2019, Li et al., 2020), in-140 ternal knowledge plays an active role in under-141 standing the input sequence, while texts generated 142 with external knowledge are more diverse. The 143 method of juxtaposing knowledge graph construcmodel's understanding of entities in sentences, 144 tion, data preprocessing and generating sequences 145 forms a way to generate sequences in an end-to-146 end manner. Although the generated sequences 147 may be diverse, it may be caused by error propa-148 gation (Liu et al., 2021).

149 150



Figure 1: Overview Of CAVSS.

151 3 Model

152 3.1 Variational Sequence to Sequenc
153 model (VSS)

154 3.1.1 Variational Auto-Encoder (VAE)

¹⁵⁵ VAE (Variational Autoencoder)(Kingma and ¹⁵⁶ Welling, 2013) has been successfully applied in ¹⁵⁷ various applications. VAE contains the framework ¹⁵⁸ of AE (Auto-Encoding), VAE directly encodes the ¹⁵⁹ input sequence as latent variable z. In the decod-¹⁶⁰ ing stage, the features are sampled on top of z:

161 $z = m + \exp(\sigma) * \varepsilon$ (1) 162 Where *m* and σ are the mean and variance 163 calculated from the input sequence *X* after pass-164 ing through the neural network, and ε is the noise 165 obtained by sampling from the σ reparameter 166 trick. The desired sequence will be decoded on 167 this feature to obtain.

VAE uses variational inference to continuously approximate the probability of the posterior by learning all observed parameters. The learning objective function is a variational lower bound on the log-likelihood of the edges of the data:

173
$$logp_{\theta}(y) \ge ELBO = \mathbb{E}_{z \sim q_{\varphi}(z|\chi)}(logp_{\theta}(y|z)) - KL(q_{\varphi}(z|y)||p(z))$$
(2)

175 3.1.2 Motivation

Sequence ¹⁷⁶ The KL divergence of VAE makes up for the de-177 fect that the Auto-Encoding model can only re-178 construct on the input sentences, and AE cannot 179 generate new samples. In the process of encoding, 180 VAE adds the sampling of Gaussian distribution 181 to form latent variable z, and then decodes from $_{182}$ z to generate new samples. Where z contains the 183 noise ε obtained from σ reparameterization 184 trick. The training goal of the VAE is to generate 185 samples from z that are similar but not identical 186 to the input, and reconstruct x from the distribu-187 tion of z. However, for z, if $z \subseteq x$, that means z, x are not independent of each other, the model ¹⁸⁹ will ignore z, $q_{\varphi}(z|x) = Dirac_Distribution(z_0)$, 190 the first term of equation (2) in the model degen-191 erates to $logp_{\theta}(x|z_0)$, at which point the model ¹⁹² bypasses the latent variable z, making the recon-193 struction of x independent of the variational pro-194 cess, so that z is added to this with a reparame-195 terization trick containing ε to decouple the ran-196 domness in the latent variable z from the formal 197 information of the data. If $z \not\subseteq x$, and z, x are 198 completely independent of each other, that is, the 2) ¹⁹⁹ noise of the variational distribution is too large, 200 resulting in the phenomenon of "posterior col-²⁰¹ lapse". $q_{\varphi}(z|x) = p_{\theta}(z)$, equation (2) in the model ²⁰² degenerates to $ELBO = \mathbb{E}_{z \sim q_{\varphi}(Z|X)}(logp_{\theta}(x|z))$, ²⁴⁵ the attention to the semantic information of the 203 the KL divergence vanishes, the variational distri- 246 original input hidden vector. Therefore, in addi-²⁰⁴ bution Regardless of x, it is very difficult to re- ²⁴⁷ tion to the sampling of the input sequence in the 205 construct output sequence from z. Therefore, the 248 Encoder part, VSS also performs a sampling pro- $_{206}$ selection of latent variable z is particularly im- $_{249}$ cess in the Decoder stage to obtain variational at-207 portant. The variational distribution in VAE often 250 tention. Assuming that the hidden state of decoder 208 uses Gaussian distribution to sample z. The opti- 251 at all moments is s_1, s_2, \dots, s_m , then in seq2seq we 209 mal parameters within the family cannot make 252 have: 210 $q_{\varphi}(z|x)$ and $p_{\theta}(x|z)$ equal, and ELBO cannot 253 ²¹¹ reach the upper bound $logp_{\theta}(x)$. The approach in ²⁵⁴ ²¹² VAE is to add auxiliary variables $\varepsilon \sim N(0, I)$, z = 255213 $m + \exp(\sigma) * \varepsilon$ So we design a flow model with 256 214 the addition of the non-affine transformation T_{257} rates the semantic information of the hidden vec-215 that convert the standard Gaussian distribution 258 tor of Encoder is obtained by calculating the atten-216 into a complex distribution, the distribution is 259 tion weights. The latent variable z_{att} is sampled 217 gradually fitted thanks to the approximation abil- 260 from the hidden vector context_i, and VSS uses a 218 ity of MLP.

219 3.1.3 Sequence to Sequence with Attention

²²⁰ The model has input $X = x_1, x_2 \dots x_n$, and $e(x_t)$ ²⁶⁴ by the encoder and the decoder input at the next is obtained through the embedding layer, the en- 265 moment, that is $e(y_{t-1}) = (embed(y_{t-1}): z_{att})$, 222 coder receives the input text, after encoding to get $_{266} y_{t-1}$ is the groundtruth in the training phase and ²²³ the hidden vector $h_{enc} = h_1, h_2 \dots h_n$. After getting ²⁶⁷ the predicted value from the previous moment is 224 each hidden vector, the semantic vector c = 268 input to decoder in the testing phase. In order to 225 $RNN(h_{enc})$ is generated. In the decoding phase, 269 prevent the decoder from being limited in obtain-226 the hidden vector of decoder and the semantic 270 ing information from the original hidden vector $_{227}$ vector c are firstly used to calculate the attention $_{271}$ space, we designed a variational attention mechaweights to get the new semantic vector c', and fi- $_{272}$ nism ²²⁹ nally the new semantic information and output of ²⁷³ LSTM($e(y_{t-1}), c, s_{t-1})$, and decoder generates a ²³⁰ the previous step to generate the next word $y_t = 274$ token by sampling from the output probability dis-²³¹ $\prod_{t=1}^{t} p(y_t|y_1, y_2, \dots, y_{t-1}, c').$

232 3.1.4 Sampling in Decoder

233 As demonstrated by Zheng et al. (Zheng et 278 ²³⁴ al.2018), it is shown that the attention mechanism ²⁷⁹ 235 is so powerful that removing other connections 280 236 between the encoder and decoder has little effect 237 on the BLEU score of the generated sequence. ²⁸¹ 3.2 CAVSS: Commonsense Augmented VSS 238 Therefore, a Sequence-to-Sequence with deter- 282 In the field of generation, there may be multiple 239 ministic attention may learn reconstructions 283 choices of candidate words, the embedding of the 240 mainly from attention, while the posterior of the 284 knowledge graph is incorporated into the process 241 latent space can fit its prior to minimize the KL 285 of sampling the latent variable z. Before this, the 242 term. In our model, this strong deterministic atten- 286 encoding process in the encoder is improved. As 243 tion may cause the model to ignore other modules, 287 shown in Figure 1, we equip encoder with a com-

$$\alpha_{ij} = softmax(e_{ij}) \tag{3}$$

$$e_{ij} = h_i^T W_e s_j \tag{4}$$

$$context_j = \sum_{i=1}^n \alpha_{ij} h_i \tag{5}$$

The context vector $context_i$ that incorpo-²⁶¹ bidirectional LSTM in the encoding stage to ob- $_{262}$ tain the memory unit c and hidden vector h of $_{263}$ the entire sentence. VSS splices the z sampled the in decoder, where $s_t =$ 275 tribution, which can be calculated as follows,

$$y_t = \prod_{t=1}^t p(y_t | y_1, y_2, \dots y_{t-1}, s_t, z_{att}) \quad (6)$$

The final loss is:

$$\mathcal{L}_{VSS} = \mathbb{E}_{z \sim q_{\varphi}(z|x), z_{att} \sim q_{\varphi}(z_{att}|x)} (logp_{\theta}(y|z, z_{att})) -KL(q_{\varphi}(z_{att}|y)||p(z_{att})) - KL(q_{\varphi}(z|y)||p(z))$$
(7)

244 and we believe that the decoder needs to weaken 288 monsense knowledge graph attention mechanism

276

277

289 KGA to incorporate commonsense knowledge 333 Empty (gray node). Then, the knowledge in the ²⁹⁰ from ConceptNet. In ConceptNet semantic net- ³³⁴ interpreter computes the graph vector G_{ent_l} of the 291 work, there is a knowledge triple tri = (h, r, t), 335 retrieved graph using a static graph attention 292 and the *h*-head node and the *t*-tail node have the $_{336}$ mechanism.

²⁹³ relation r. Assuming that there are l entities in ³³⁷ ²⁹⁴ the input sequence x_k , then we have $entity_{x_k} = 338$ and then encoding more structured semantic infor-295 $\{ent_1, ent_2, ..., ent_l\}$, where each entity corre- 339 mation, we design the following attention, each sponds to a different number of triples group, then 340 entity subgraph $g_{ent_l} = \{tri_1, tri_2, ..., tri_{N_l}\}$ as in-²⁹⁷ we have $g_{ent_l} = \{tri_1, tri_2, ..., tri_{N_l}\}$, for entities ³⁴¹ put to construct the graph vector G_{ent_l} : 298 not in the commonsense database, We assign them 342 299 the value of *empty*.

When given a set of input and output se- 344 300 301 quences, for external commonsense knowledge, 345 302 the model can only select valid commonsense 346 tionship nodes, head entities and tail entities, and 303 knowledge triples based on prior distribution 347 the attention weight measures the degree of asso-304 learning, and it is difficult to obtain the correct 348 ciation between head and tail and relation. The 305 posterior distribution in the inference stage. Our 349 graph vector G_{ent_1} is a weighted sum that com-306 solution is to splice the context vectors of the em- 350 bines the semantic relationships of head and tail 307 bedding and encoder parts of the commonsense 351 nodes. ³⁰⁸ knowledge before sampling, so that the prior dis-³⁰⁹ tribution approximates the posterior distribution, ³⁵² **3.2.2 Graph Attention** 310 so that appropriate knowledge can be selected 353 As shown in Figure 1, the graph in the model is 311 even without the posterior information. We intro- 354 embedded after the encoder. As described in sec- $_{312}$ duce an auxiliary loss(section 3.2.4), called align $_{355}$ tion 3.2.1, the vector c with the semantic inforloss, to measure the distance between the prior and 356 mation of the input sequence enters the KGA posterior distributions. 314

The graph attention mechanism needs to gen- 358 the probability of using each triplet: 315 316 erate vector representations for the retrieved sub- 359 317 graphs, but the relationship between entities is of- 360 318 ten not negligible, so the KGA in the model is di- 361 319 vided into two modules: the graph embedding at- 362 ³²⁰ tention module and graph attention module.

321 3.2.1 Attention in Graph Embedding

 $_{323}$ of relation vectors and head-tail entities by re- $_{367}$ and the latent variable z is sampled on this basis 324 trieving the entire 325 knowledge base using each word in Input (red ³²⁶ node) as a key entity. The retrieved graph consists ³⁶⁹ **3.2.3 Autoregressive Transformer** ³²⁷ of a key entity (red node), its neighboring entities ³⁷⁰ According to section 3.1.2, the autoregressive 328 (blue nodes are the head nodes and green nodes 371 sampling transformation is proposed to prevent $_{329}$ are the tail nodes), and relationships (directed $_{372}$ the KL vanishing phenomenon. The mean m and $_{330}$ edges). For common words that do not match en- $_{373}$ variance σ sampling in VAE are both imple-331 titles in the commonsense knowledge base (such 374 mented through the fully connected feedforward

Considering the relationship between nodes

$$\tau_n = W_r r_n \tanh \left(W_h h_n + W_t t_n \right) \qquad (8)$$

$$\alpha_n^{EKG} = softmax(\tau_n) \tag{9}$$

$$G_{ent_l} = \sum_{n=1}^{N_l} \alpha_n^{EKG} (h_n; t_n)$$
(10)

 W_r, W_h, W_t are the weight matrices of rela-

³⁵⁷ module to query the target triplet and calculates

$$\alpha_{lt}^{KGA} = c_t W_c G_{ent_l} \tag{11}$$

$$\beta_{l_t}^{KGA} = softmax(\alpha_{l_t}^{KGA})$$
(12)

$$k_l = \sum_{t=1}^{N_l} \beta_{l_t}^{KGA} G_{ent_l} \tag{13}$$

 W_c is a trainable parameter, the graph vector $_{363} k_l$ is the weighted sum of the target graph embed- $_{364}$ ding, and the graph vector k_l is spliced with the $_{365}$ context vector c output by the encoder to obtain ³²² This module aims to facilitate the semantic fusion ³⁶⁶ $(k_i; c)$. The $(k_i; c)$ is input to the sampling layer, general commonsense ₃₆₈ to yield $z = GaussianSampling((k_1:c))$.

332 as in), they are represented by the special node 375 neural network structure, so the sampling can be

376 trained. By continuously iterating the reversible 417 $_{377}$ transformation T, the latent variable z is made $_{418}$ 378 smoother and more flexible, and the direct use of 419 379 Gaussian sampling is not accurate enough.

$$T(z) = \frac{1}{\kappa} \sum_{j=1}^{K} w_j f(w_j' z + b_j)$$
(14)

381 $_{382}$ ter K transformations. w_i, w'_i, b_i are the training $_{424}$ we want to favor the sampling of the knowledge ₃₈₃ parameters of the neural network. Here $f(\cdot)$ uses 425 graph to generate more diverse texts, we set the ³⁸⁴ the PRelu function. Since the model incorporates ⁴²⁶ hyperparameter β_{att} between the two KLs. 385 the semantics of the knowledge graph into the ³⁸⁶ context vector for sampling as well, the sampling ⁴²⁷ **4** Experiment $_{387}$ space of z will be slightly larger than the sam-³⁸⁸ pling space of the original input. Therefore, we ³⁸⁹ add autoregressive transformation hoping to grad- ⁴²⁹ Commonsense External Knowledge ³⁹⁰ ually fit the distribution of P(x).

391 3.2.4 Training Target

³⁹³ duced in section 3.2.3. The parameters of this ⁴³⁴ and understand human intent. It consists of nodes 394 transformation are trainable, but this training 435 that represent concepts expressed as words or 395 lacks objects that need to be aligned. In other 436 phrases in natural language, and in which the re-³⁹⁶ words, these parameters require a training target. ⁴³⁷ lationships of these concepts are labeled, e.g. ³⁹⁷ We introduce the alignment vector z_{align} , where ⁴³⁸ (London, AtLocation, American), these features ³⁹⁸ the latent variable z sampled on the Gaussian ⁴³⁹ are important in the learning of the model. ³⁹⁹ distribution tends to the distribution of P(x) after ⁴⁴⁰ **Dataset** 400 transformation T. Therefore, in the training phase, 441 We applied the model to a question generation 401 y is used as the existing data to be input into the 442 task using the Stanford Q&A dataset (Rajpurkar et ⁴⁰² model together with x to obtain the z_{align} vec- ⁴⁴³ al., 2016). The attention mechanism is particularly 403 tor, and the distance between z and z_{align} distri- 444 important when generating questions based on 404 bution is approximated by KL divergence. Note: 445 sentences and hopefully open-ended questions. 405 This part is only used during the training phase. 446 The integration of commonsense knowledge al-Add the loss of z_{align} sampling to the original 447 lows the generated questions to tend to be diverse. 407 loss:

$$uos align = -KL((q_{\varphi}(z_{align}|y, x, g)||p(z_{align})) + \mathbb{E}_{z_{align} \sim q_{\varphi}(z_{align}|x)}(logp_{\theta}(y|z, z_{align}))$$

410 3.3 Loss

412 focusing only on text generation based on recon-413 struction, we add a hyperparameter γ_{KL} to the 455 alogue and other text generation tasks. 414 loss function to balance the reconstruction loss 456 Dist: We use Dist-1, Dist-2 (Li et al.2016) to 415 and KL loss. The new loss function is obtained as 457 measure the diversity of generation. We count the 416 follows:

$$\mathcal{L} = \mathbb{E}_{z \sim q_{\varphi}(z|x, y, g)}(logp_{\theta}(y|z, z_{att}, z_{align})) - \gamma_{KL}(KL(q_{\varphi}(z_{att}|y)||p(z_{att})) + KL((q_{\varphi}(z_{align}|y, x, g)||p(z_{align})) + \beta_{att}KL(q_{\varphi}(z|y)||p(z))$$
(16)

Since the sampling of the knowledge graph 422 and the sampling of the decoder layer are inte-The final latent variable T(z) is obtained af- 423 grated into the sampling of the knowledge graph,

428 4.1 Dataset

421

430 A commonsense knowledge base is built using 431 ConceptNet, a semantic network that contains a 432 large amount of information a computer should ³⁹² The autoregressive transformation T is intro-⁴³³ know about the world to help it do better searches

448 4.2 Evaluation

449 Automatic Evaluation

¹⁵⁾₄₅₀ **BLEU:** We use BLEU-1 to BLEU-4 scores (Pap-451 ineni et al.2002) as a criterion for evaluating the 452 accuracy of generated sentences by measuring 411 To prevent the model from ignoring sampling and 453 word overlap between ground truth and generated 454 sentences, widely used in machine translation, di-

458 proportion of distinct 1-grams and 2-grams in the

459 generated sentences to evaluate the diversity of 469 averaged as the scores for each indicator. We de-470 fine the following four metrics: Fluency (whether 460 the output.

Manual Evalution 461

462 ⁴⁶³ accurately evaluate the quality of the generated se- ⁴⁷³ generated sequence to the topic), **Diversity** 464 quences. We ask evaluators to assess each group 474 (whether the sequence includes new information 465 of 100 generated sequences, and for each gener- 475 or knowledge in addition to the original input con-466 ated sequence, three evaluators are hired to give a 476 tent), Commonsense (whether the generated se-467 score from 1 to 5 based on the following three 477 quence is incorrect in terms of commonsense). 468 metrics. The scores of the three raters were

471 the sequence is appropriate in terms of grammar We also perform human evaluation to more 472 and logic), Topic (the degree of relevance of the

			Automatic Ev	valuation		
Models	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Dist-1	Dist-2
VAE	29.31	12.42	6.55	3.61	-	-
Seq2Seq	31.34	13.79	7.36	4.26	-	-
VSS	32.71	16.24	9.63	5.91	0.140	0.211
CAVSS	33.09	16.4	9.81	6.1	0.144	0.219
CAVSS-T	33.08	16.52	9.91	6.2	0.158	0.229
CAVSS-a	32.95	16.4	9.77	6.03	0.136	0.203
Full Model	33.46	16.82	10.13	6.32	0.142	0.214

Table 1:	Automatic	Evaluation	with BLEU	J and Dist.
----------	-----------	------------	-----------	-------------

	Ν	/Ianual Evaluat	ion	
Models	Fluency	Topic	Diversity	Commonsense
VSS	3.93	3.96	3.93	3.92
CAVSS	3.91	4.06	4.08	4.03
CAVSS-T	4.01	3.97	4.12	4.08
CAVSS-a	4.03	4.09	3.98	3.97
Full Model	4.07	4.11	4.09	4.14

Table 2: Manual Evaluation with Fluency, Topic and Diversity.



Figure 2: BLEU1~4 with different γ_{KL} values.

478 5 **Result and Analysis**

479 An ablation study of text quality. To understand 480 the contribution of each component of our model to ⁴⁸¹ the task, we train two ablation versions of the model: 482 with or without commonsense knowledge for VSS 483 ("w/o KG") and with or without autoregressive

485 without align sampling ("w/o a"). Table 1 and Table 486 2 show the automatic evaluation scores and human 487 evaluation results for the ablation study.

Experimental results. Table 1 and Table 2 489 show the performance of our model and the baseline 490 model, where we construct the traditional Seq2Seq 491 model and the VAE model. By comparing VSS and 484 transformation for CAVSS ("w/o T"), with or 492 CAVSS, we find that without commonsense 493 knowledge, model performance degrades in all 494 metrics. The improved Topic scores indicates that by 528 interesting that CAVSS-T achieves the highest diver-495 learning the correlation between an entity and its 529 sity score when align-aligned sampling is removed $_{496}$ neighboring concepts, concepts that are more $_{530}$ and an autoregressive sampling transformation T is 497 closely related to an entity receive higher attention 531 added. We believe that the lack of z_{align} alignment 498 during the generation process. The reason for the im- 532 sampling constraints makes the text generation diproved Diversity is that the expansion of external 533 versity increase, and the quality of generated se-499 500 knowledge graph information makes the output text 534 quences is not significantly improved. more novel. The improvement in Fluency is because 535 Loss factor experiment. We adjust the hy-501 so by learning the relationship between entities, it can 536 perparameter γ_{KL} in Equation 16 with the BLEU ⁵⁰³ help the model to better fit the data. All versions of ⁵³⁷ index as the criterion, and we apply the Full Model our CAVSS models outperform the baselines in all 538 for experiments. As shown in the Figure 2, γ_{KL} evaluation metrics. In the manual evaluation, 539 takes values in 0.3 and 2.0 both have higher scores. 505 CAVSS-T obtained the highest score (4.12) for the 540 When γ_{KL} =0.3, that is, the reconstruction part is Diversity index and Full Model obtained the highest 541 strong, the model does not need to extract features score (4.07,4.11,4.14) for the Fluency, Topic, and 542 from latent variables, and the model construction 508 Commonsense indices, and similar conclusions can 543 fails, although various indices are improved, it debe drawn from the automatic evaluation. Similar 544 feats the original intent of the model. When 510 conclusions can also be drawn from the automatic 545 γ_{KL} =2.0, the model focuses on optimizing the KL 512 evaluation. The improvement of CAVSS-T in gener- 546 term. The model generates sequences through latent 513 ating text Diversity and Commonsense is significant, 547 variables and achieves better results in terms of gens14 and this improvement comes from our external com- 548 eration quality. γ_{KL} serves as a regularization factor 515 monsense knowledge, as our sentence representa- 549 that aims to constrain the capacity of latent variables 516 tions are generated by an autoregressive transfor- 550 and find the right balance between the Reconstruc-517 mation of samples of continuous latent variables. 551 tion part and KL, which is consistent with the find-⁵¹⁸ Compared to the baseline, this step introduces more ⁵⁵² ings of Higgins et al. (Higgins et al., 2016). ⁵¹⁹ randomness. When using z_{align} alignment sam- ⁵⁵³ Case study. As shown in Table 3, for the origi-520 pling on the basis of CAVSS-T, i.e. (Full Model), 554 nal input, there are triples "magnitude Related To 521 Full Model achieves the best performance in BLEU 555 earthquake", "scale RelatedTo deep", "estimated Re-522 (33.46). However, we found that the Full Model did 556 lated To model", and the model performs KGA on the 523 not significantly outperform the CAVSS model on 557 knowledge triples. By learning relations and entities, 524 diversity metrics (Dist-1, Dist-2). The results show 558 we pay different attention to the neighboring entities s25 that aligning the sampling is beneficial for the model 559 of different entities in the original sentence and use 526 to better fit the test set, but it does not significantly 560 this structured information to encode and decode 527 help other important metrics such as Dist. We find it 561 them for the purpose of generating diverse sentences. Input: In a united states geological survey uses study preliminary runture models of the earthquake indicated dis

input: In a united states geological survey usgs study preliminary rupture models of the earthquake indicated dis-		
placement of up to 9 meters along a fault 240 km long by 20 km deep.		
Output: How large was the displacement?		
VSS	What percentage of the earthquake was conducted by the earthquake?	
CAVSS	What was the magnitude of the earthquake?	
CAVSS-T	What was the scale of the earthquake?	
CAVSS-a	How deep was the earthquake that damaged the US?	
Full Model	What was the estimated displacement of the earthquake in the US?	

Table 3: Sample questions generated by all the models.

562 References

563 (Kingma and Welling, 2013) Kingma, Diederik P., and 608 3301.

564 Max Welling. "Auto-encoding variational bayes." 609 (Higgins et al., 2016) Irina Higgins, Loic Matthey,

565 arXiv preprint arXiv:1312.6114 (2013).

- 567 Jiangnan Xia, and Chenliang Li. 2019. Incorporating 612 Lerchner (2016). beta-vae: Learning basic visual con-

External Knowledge into Machine Reading for Gener- 613 cepts with a constrained variational framework.

569 ative Question Answering. In Conference on Empirical 614 (Sutskever et al., 2014) Ilya Sutskever, Oriol Vinyals,

572 cessing (EMNLP-IJCNLP).

575 and Yoshua Bengio. Graph attention networks. CoRR, 620 ral machine translation. Transactions of the Association 576 abs/1710.10903, 2017.

581 Conference on Computational Natural Language 626 Representations.

582 Learning, pages 10–21, Berlin, Germany. Association 627 (Wang et al., 2019) Wenlin Wang, Zhe Gan, Hongteng

583 for Computational Linguistics.

585 Salimans, Yan Duan, Prafulla Dhariwal, John Schul- 630 guided variational autoencoders for text generation. 586 man, Ilya Sutskever and Pieter Abbeel. Variational 631 arXiv preprint arXiv:1903.07137.

⁵⁸⁷ lossy autoencoder(J). arXiv preprint arXiv:1611.02731, 632 (Wang et al., 2018) Wenlin Wang, Yunchen Pu, Vinay 588 2016.

589 (Yu et al., 2020) Wenhao Yu, Chenguang Zhu, Zaitang 634 Piyush Rai, and Lawrence Carin. 2018b. Zero-shot 590 Li, Zhiting Hu, Qingyun Wang, Heng Ji and Meng 635 learning via class-conditioned deep generative models.

- ⁵⁹¹ Jiang (2020). A survey of knowledge-enhanced text 636 In AAAI.
- ⁵⁹² generation. arXiv preprint arXiv:2010.04389.

⁵⁹⁶ print arXiv:2010.05511.

- 600 bilistic modeling and inference. Journal of Machine 645 (Wei et al., 2019) Xiangpeng Wei, Yue Hu, Luxi Xing,
- 601 Learning Research, 22(57), 1-64.

- 604 Aaron C Courville, and Yoshua Bengio. 2017. A hier- 649 (AAAI).

605 archical latent variable encoder-decoder model for 650 (Li et al., 2020) Haoran Li, Junnan Zhu, Jiajun Zhang,

606 generating dialogues. In Proceedings of the 31st AAAI

607 Conference on Artificial Intelligence, pages 3295-

610 Arka Pal, Christopher Burgess, Xavier Glorot, Mat-(Bi et al., 2019) Bin Bi, Chen Wu, Ming Yan, Wei Wang, 611 thew Botvinick, Shakir Mohamed and Alexander

570 Methods in Natural Language Processing and Interna- 615 and Quoc V Le. Sequence to sequence learning with 571 tional Joint Conference on Natural Language Pro- 616 neural networks. In NIPS, pages 3104-3112, 2014.

617 (Zheng et al. 2018) Zaixiang Zheng, Hao Zhou, Shujian 573 (Velickovic et al., 2017) Petar Velickovic, Guillem Cu- 618 Huang, Lili Mou, Xinyu Dai, Jiajun Chen, and 574 curull, Arantxa Casanova, Adriana Romero, Pietro Lio, 619 Zhaopeng Tu. 2018. Modeling past and future for neu-621 for Computational Linguistics, pages 145–157

577 (Bowman et al., 2016) Samuel R. Bowman, Luke Vil- 622 (Bahdanau et al.2015) Dzmitry Bahdanau, Kyunghyun 578 nis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and 623 Cho, and Yoshua Bengio. 2015. Neural machine trans-579 Samy Bengio. 2016. Generating Sentences from a Con- 624 lation by jointly learning to align and translate. In Pro-580 tinuous Space. In Proceedings of The 20th SIGNLL 625 ceedings of the International Conference on Learning

628 Xu, Ruiyi Zhang, Guoyin Wang, Dinghan Shen, 584 (Chen et al., 2016) Xi Chen, Diederik P. Kingma, Tim 629 Changyou Chen and Lawrence Carin (2019). Topic-

633 Kumar Verma, Kai Fan, Yizhe Zhang, Changyou Chen,

637 (Li et al., 2021) Jing Li, Qingbao Huang, Yi Cai, 593 (Qiao et al., 2020) Qiao, L., Yan, J., Meng, F., Yang, Z., 638 Yongkang Liu, Mingyi Fu and Qing Li(2021). Topic-594 & Zhou, J. (2020). A sentiment-controllable topic-to- 639 level knowledge sub-graphs for multi-turn dialogue ⁵⁹⁵ essay generator with topic knowledge graph. arXiv pre- ⁶⁴⁰ generation. Knowledge-Based Systems, 234, 107499.

641 (Zhu et al., 2017) Wenya Zhu, Kaixiang Mo, Yu Zhang, 597 (Papamakarios et al., 2021) Papamakarios, G., 642 Zhangbin Zhu, Xuezheng Peng, and Qiang Yang. Flex-598 Nalisnick, E., Rezende, D. J., Mohamed, S., & Laksh- 643 ible end-to-end dialogue system for knowledge ⁵⁹⁹ minarayanan, B. (2021). Normalizing flows for proba- 644 grounded conversation. CoRR, abs/1709.04264, 2017 646 Yipeng Wang, and Li Gao. 2019. Translating with Bi-602 (Serban et al., 2017) Iulian Vlad Serban, Alessandro 647 lingual Topic Knowledge for Neural Machine Transla-603 Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, 648 tion. In AAAI Conference on Artificial Intelligence

- 651 Chengqing Zong, and Xiaodong He. 2020. Keywords- 696 (Liu et al., 2021) Shuang Liu, Nannan Tan, Yaqian Ge
- Conference on Artificial Intelligence (AAAI). 653
- 654 (Zhang et al., 2018) Saizheng Zhang, Emily Dinan, 699 on Pointer Network. Information, 12(3), 136.

655 Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason 700 (Zhao et al., 2017) Tiancheng Zhao, Ran Zhao, Maxine

656

- 658 ation Computational Linguistics (ACL).
- 659 (Zhou et al., 2018) Hao Zhou, Minlie Huang, Tianyang

660 Zhang, Xiaoyan Zhu, and Bing Liu. 2018. Emotional

chatting machine: Emotional conversation generation

662 with internal and external memory. In AAAI Confer-

663 ence on Artificial Intelligence (AAAI).

664 (Brown et al., 2020) Tom Brown, Benjamin Mann,

665 Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla

666 Dhariwal., ... & Dario Amodei (2020). Language mod-

667 els are few-shot learners. Advances in neural information processing systems, 33, 1877-1901.

669 (Rajpurkar et al.2016) Pranav Rajpurkar, Jian Zhang,

670 Konstantin Lopyrev, and Percy Liang. 2016. SQuAD:

⁶⁷¹ 100,000+ questions for machine comprehension of text.

672 In Proceedings of the 2016 Conference on Empirical

673 Methods in Natural Language Processing, pages 2383-674 2392.

675 (Mikolov et al.2013) Tomas Mikolov, Ilya Sutskever, 676 Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Dis-

677 tributed representations of words and phrases and their 678 compositionality. In Advances in Neural Information

Processing Systems, pages 3111–3119.

680 (Kingma and Ba 2015) Diederik Kingma and Jimmy

681 Ba. 2015. Adam: A method for stochastic optimization.

682 In Proceedings of the International Conference on 683 Learning Representations.

(Papineni et al.2002) Kishore Papineni, Salim Roukos, 684

685 Todd Ward, and Wei-Jing Zhu. 2002. BLEU: A method

for automatic evaluation of machine translation. In Pro-687 ceedings of the 40th Annual Meeting on Association

for Computational Linguistics, pages 311–318.

689 (Li et al.2016) Jiwei Li, Michel Galley, Chris Brockett,

690 Jianfeng Gao, and Bill Dolan. 2016. A diversity-pro-

⁶⁹¹ moting objective function for neural conversation mod-

692 els. In Proceedings of the 2016 Conference of the North

693 American Chapter of the Association for Computa-

694 tional Linguistics: Human Language Technologies,

695 pages 110-119.

652 guided abstractive sentence summarization. In AAAI 697 and Niko Lukač. (2021). Research on Automatic Ques-698 tion Answering of Generative Knowledge Graph Based

Weston. 2018. Personalizing Dialogue Agents: I have a 701 Eskenazi (2017). Learning discourse-level diversity for 657 dog, do you have pets?. In Annual Meeting of Associ- 702 neural dialog models using conditional variational au-703 toencoders. arXiv preprint arXiv:1703.10960.