

# Contrastive Multi-document Question Generation

Woon Sang Cho\* Yizhe Zhang† Sudha Rao† Asli Celikyilmaz†  
Chenyan Xiong† Jianfeng Gao† Mengdi Wang\* Bill Dolan†

\*Princeton University

†Microsoft Research AI

\*{woonsang, mengdiw}@princeton.edu

†{yizhzhang, sudhra, aslicel, cxiong, jfgao, billdol}@microsoft.com

## Abstract

Multi-document question generation focuses on generating a question that covers the common aspect of multiple documents. However, models trained only using the targeted (“positive”) document set generate questions that are generic i.e. they cover a larger scope than delineated by the document set. To address this challenge, we introduce the contrastive learning strategy where given “positive” and “negative” sets of documents, we generate a question that is closely related to the “positive” set but is far away from the “negative” set. This setting allows generated questions to be more specific and related to the target document set. To generate such specific questions, we propose Multi-Source Coordinated Question Generator (MSCQG), a novel framework that includes a supervised learning (SL) stage and a reinforcement learning (RL) stage. In the SL stage, a single-document question generator is trained. In the RL stage, a coordinator model is trained to find optimal attention weights among multiple single-document generator instances, by optimizing a reward designed to promote specificity of generated questions. We also develop an effective auxiliary objective, named Set-induced Contrastive Regularization (SCR) that improves the coordinator’s contrastive learning during the RL stage. We show that our model significantly outperforms several strong baselines based on retrieval and neural generation, as measured by automatic metrics and human evaluation.

## 1 Introduction

User queries on web search engines can sometimes be vague. Search engines may resolve this ambiguity by suggesting clarification options back to the user in the form of questions (Braslavski et al., 2017; Aliannejadi et al., 2019b; Zamani et al., 2020). Following Cho et al. (2019b), this approach

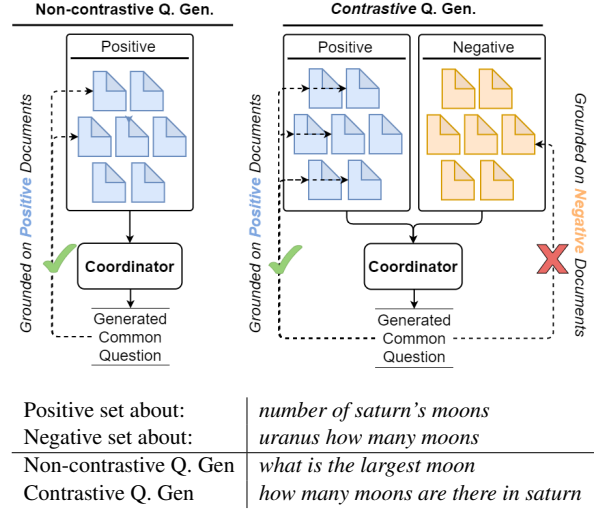


Figure 1: Non-contrastive and contrastive method for multidocument question generation. Left: non-contrastive modeling that takes input as a set of positive documents. However, model-generated questions from this method are rather generic and not specific to the input documents. Right: *contrastive* modeling, which considers both positive and negative document sets, and learns to generate common questions that are more grounded on the positive document set.

may involve three phases: *i) retrieval*: gather the initial return documents by the search engine; *ii) partition*: partition the documents into semantically similar clusters in an unsupervised manner; *iii) multi-document question generation*: generate a clarification question by finding an “overlap” among documents in each cluster. In principle, the clarification questions should be “mutually exclusive” rather than generic and bland, otherwise it is counter to the objective of clarification (Radlinski and Craswell, 2017; Rao and Daumé III, 2018, 2019). In this work, we focus on the last step.

We address this challenge by leveraging contrastive learning. Given a set of positive documents  $\mathcal{D}^+$  and a set of negative documents  $\mathcal{D}^-$  (where  $\mathcal{D}^-$  is semantically close to  $\mathcal{D}^+$ ), we propose a new

strategy to generate a question that is relevant to  $\mathcal{D}^+$  and far away from  $\mathcal{D}^-$ . The comparison between the contrastive and non-contrastive multi-document question generation is illustrated in Figure 1.

Ideally, the model that leverages the challenging “negative” documents in addition to the “positive” documents, to push the model to distinguish these two sets and constrain the generation to be specific to the positive set. The similarity between the  $\mathcal{D}^+$  and  $\mathcal{D}^-$  makes the generation more challenging and forces the model to be as specific as possible in order to distinguish between the two sets.

This task is particularly challenging because *i)* there does not exist direct supervised *ground-truth* multi-document question given positive and negative sets of documents. *ii)* The whole procedure involves multiple aspects including language understanding, inter-document information aggregation, coordinative planning and language generation. To address these, we employ a multi-step strategy using reinforcement learning.

We also propose a novel auxiliary objective, Set-induced Contrastive Regularization (SCR) (Section 4), which heuristically drives the coordinator’s generation closer towards  $\mathcal{D}^+$  by minimizing KL divergence between aggregated contrastive word distributions and distributions induced by  $\mathcal{D}^+$ . Likewise drives it away from  $\mathcal{D}^-$  by maximizing KL divergence but limiting this effect by monitoring how similar the two sets of distributions are.

We evaluate a generated question from each model by assessing how many of the input documents it can reverse-retrieve, using a publicly available state-of-the-art pre-trained ranker (Section 4), and by crowd-sourced human judgments.

Our contributions are summarized in below: *i)* We develop a novel Multi-Source Coordinated Question Generator (MSCQG) model that is trained using a hierarchical generation scheme. The document-specific generation is fine-tuned from GPT-2 and the inter-document coordinator is trained using reinforcement learning. *ii)* We introduce Set-induced Contrastive Regularization (SCR), an auxiliary regularizer that pushes MSCQG toward  $\mathcal{D}^+$  relative to  $\mathcal{D}^-$  while limiting the effect of  $\mathcal{D}^-$  in a principled manner. *iii)* Empirical results show that our model is able to generate more grounded and specific questions, significantly outperforming existing baseline models in automatic measures and human evaluation.

## 2 Related Work

**Multi-Source Encoder-Decoder:** Ensemble set induction mechanism (Rokach, 2010) that has been widely applied to neural machine translation tasks (NMT) (Bojar et al., 2014). Firat et al. (2016) introduced a new type of ensemble of NMT systems which take inputs as multiple sentences in different languages and output a translation into a single language. Each NMT system is trained on a mono-lingual source to target language translation dataset. Garmash and Monz (2016) further developed the multi-source encoder-decoder framework for the multi-lingual neural machine translation task, by learning to assign uneven attention weights, called *expert combination weights* among multi-lingual NMT systems. For such multi-lingual translation tasks, the target translation is available. However, in this task of generating multi-document questions, the target does not exist which makes it more challenging. To handle multi-source input, we take a similar multi-source encoder-decoder approach for our coordinator model, which is trained via RL, rather than supervised learning.

**Question Generation:** Most prior work on question generation has been on single document i.e. given a document and an answer phrase in the document, generate a question that is answered by the answer phrase (Heilman, 2011; Rus et al., 2010). However, in our work, we aim to generate a multi-document question is answerable by multiple input documents. Recently, sequence-to-sequence based neural network models have defined the state-of-the-art for question generation (Du et al., 2017a; Duan et al., 2017a). For a survey, see Pan et al. (2019). Our generator model, on the other hand, is based on the more recent GPT-2 (Radford et al., 2019) generation model, and this forms the underlying component of our question generating system. Fan et al. (2018) propose a visual question generation model to generate natural questions about images using reinforcement learning where they use naturalness and human-like as reward signals. In our work, we use retrieval statistics, similar to Nogueira and Cho (2017), derived from a document-question ranker as the reward for training our coordinator model in isolation, rather than the entire generating pipeline.

**Contrastive learning in NLP:** Contrastive learn-

ing has been widely used in NLP (Smith and Eisner, 2005; Collobert et al., 2011; Bordes et al., 2013; Hjelm et al., 2019; Deng et al., 2020). Broadly, contrastive learning methods differentiate observed data from artificial negative examples. Gutmann and Hyvärinen (2010) leverages the Noise Contrastive Estimation (NCE) metric to differentiate the target sample from noise samples. Negative Sampling proposed by Mikolov et al. (2013) is a simplified variation of NCE loss. Recently, contrastive learning has also been employed in learning sentence representation (Clark et al., 2020). Our contrastive learning approach is fundamentally different from above, which learn representations. To our best knowledge, we are the first to leverage contrastive learning and establish set-induced penalization in the context of question generation.

### 3 Preliminary

Cho et al. (2019b) introduced a task of generating common questions that can be answered by multiple documents. The authors employ an information aggregation mechanism by linearly interpolating the output distributions of each decoder. Initially, a recurrent neural network (Rumelhart et al., 1988; Werbos, 1990) sequence-to-sequence (Seq2Seq) model (Sutskever et al., 2014; Bahdanau et al., 2014) is trained from a single document (input) and a single question (output) that is *answerable* by the input document. Using multiple instances of the Seq2Seq model, individual inferences are made from  $N$  input documents after which decoding distributions  $\pi_{i,t}$  ( $i \in \{1, \dots, N\}$ ) are averaged to generate a common word for the question at time  $t$ .

$$\pi_t = \frac{1}{N}(\pi_{1,t} + \pi_{2,t} + \dots + \pi_{N,t}) \quad (1)$$

where  $\pi_t$  is the common decoding distribution at time  $t$ . Further heuristics are applied on  $\pi_t$  to improve their performance. As shown in their evaluation using automatic metrics and human judgments, the generated questions are, by and large, relevant to the document set, but can be relatively generic. This is because the training objective only requires the generation to be answerable by all documents, thus a generic question that covers broad topics will not be penalized.

### 4 Method

The overview of our model is illustrated in Figure 2. The method concerns training two ma-

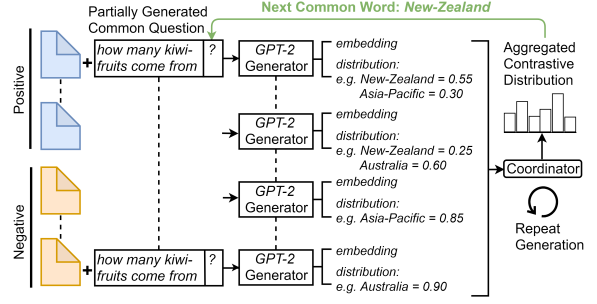


Figure 2: System overview. The example is an illustration using fictitious tokens for ease of understanding. Our MSCQG model learns to attend different weights and form a final aggregated distribution at each decoding time, given the input embeddings and distributions. The decision to enforce or penalize the negative set distributions to the aggregated distribution is controlled in a principled manner.

major components: *document-specific generator* and *inter-document coordinator*.

The *document-specific generator*, which generates a question from a single document, is trained by fine-tuning the pre-trained GPT-2 model. The *inter-document coordinator* is a novel hierarchical generation component that integrates multiple-document information from the *document-specific generator* instances. This coordinator is trained using reinforcement learning and an auxiliary regularization to encourage (or discourage) the generation to be more relevant to the positive (or negative) set. The source code and dataset will be released upon publication.

**Document-specific GPT-2 Generator:** At the first pre-training stage, we load the publicly available GPT-2 (Radford et al., 2019) model as our underlying document-specific generator. The GPT-2 model leverages massive out-domain data and serves as good initialization to generate grammatical and informative question. Then, we further fine-tune the language model on MS-MARCO (Nguyen et al., 2016) *selected* document as an input followed by a special separator and the corresponding question as an output.

At every generation step, the previously generated words are concatenated to all documents as inputs to the generator instances, which yield updated hidden states and output distributions for the next common word.

**Ranking-based Rewards Computation:** Before learning a policy for the inter-document coordi-

nator, we first describe the calculating the reward signal based on retrieval statistics from a BERT-based document ranker.

The BERT-based (Devlin et al., 2018a) ranker (Nogueira and Cho, 2019) (*Ranker*), a state-of-the-art model <sup>1</sup> in the MS-MARCO Passage Retrieval task (Nguyen et al., 2016), is trained to rank (document, question) pairs. This ranker assigns high scores for *true positive* document and question pairs. We assume the ranker delivers accurate reward signal since it achieves good performance on the challenging retrieval task. Let  $\tilde{q}$  be the generated question from the underlying generator block and coordinator with the positive and negative document sets ( $\mathcal{D}^+$  and  $\mathcal{D}^-$ ) as the input.

$$\text{Ranker}(d, \tilde{q}) = \text{score} \in (0, 1) \quad (2)$$

$$\forall d \in \mathcal{D}^+, \mathcal{D}^-$$

We pair  $\tilde{q}$  with each of the documents in the positive and negative set, and evaluate the question-document pairs through the ranker for answer-relevancy. These retrieval scores for each document that lie in  $(0, 1)$  are sorted in descending order. Using score memberships in  $\mathcal{D}^+$  or  $\mathcal{D}^-$ , we compute retrieval statistics, such as *Precision@10* and *mean-Average-Precision (mAP)* (Zhu, 2004) which are candidate non-differentiable rewards  $R$ .

**Inter-generator Coordinator:** Next, we train a coordinator system using policy gradient to optimize the reward described above. The separation between the generator and the coordinator aims to ease the RL training by significantly reducing the action space.

Note that the generator model is fixed during this stage. We find that using RL to train the entire generating pipeline yields large variance since the action space is large and the auto-regressive nature of the generation process further amplifies such variance. Therefore, we fix the underlying generator component and then on top of multiple instances of the underlying generator, we stack our transformer-based (Vaswani et al., 2017) coordinator model which is trained using RL in isolation. Instead of training both token-level GPT-2 and document-level coordinator over multiple GPT-2 instances using RL, only the coordinator is trained using RL, which structure dramatically reduces variance.

The coordinator is a transformer-based (Vaswani et al., 2017) model with damped contrastive dis-

tribution. Unlike Transformer Decoder (Liu et al., 2018), there is no causal mask. Instead, the coordinator model uses the hidden states updated every decoding time from the underlying fine-tuned GPT-2 (Radford et al., 2019) language model generators.

We add learned *cluster embedding*  $c_i$  to the input document hidden states  $h_i$ , similar to learned positional embedding (Devlin et al., 2018a), to indicate whether the source document  $i$  is in  $\mathcal{D}^+$  or  $\mathcal{D}^-$ .

$$x_i^0 = h_i + c_i \quad (3)$$

The coordinator model consists of  $n$  recurrent transformers blocks (Vaswani et al., 2017), followed by three different feed-forward layers ( $\text{FF}_w$ ,  $\text{FF}_v$ , and  $\text{FF}_z$ ) to output  $w$ ,  $v$ , and  $z$ .

$$x^k = \text{Add-Norm}(u, \text{FF}_x(u)) \quad (4)$$

$$u = \text{Add-Norm}(x^{k-1}, \text{MultiHead}(x^{k-1})) \quad (5)$$

for  $k = 1, \dots, n$

$$w = \text{FF}_w(x^n) \quad (6)$$

$$v = \text{FF}_v(x^n) \quad (7)$$

$$z = \text{FF}_z(x^n) \quad (8)$$

$w$  and  $v$  are the attention weights among the positive documents  $\mathcal{D}^+$ , and negative documents  $\mathcal{D}^-$ . We weight the documents unevenly because often times not all the top-10 documents in the set share the same content. Thus, we leave to the model to learn the optimal attention weights among positive and negative sets that produce a more grounded common question.

$z$  parametrizes  $\eta$  in how much the coordinator model penalizes, or sometimes reinforces, weighted average of decoding distributions from the negative set  $\mathcal{D}^-$ .  $\eta$  is a simple heuristic variation of  $\tanh$  such that the image lies in  $(-1, 0.5)$  for all real numbers  $\mathbb{R}$ . Thus  $\eta$  is a damped penalization coefficient.

$$\eta(z) = -\frac{e^{2z} - 0.5}{e^{2z} + 1} \in (-1, 0.5) \quad \forall z \in \mathbb{R} \quad (9)$$

Given  $w$ ,  $v$ , and  $z$ , we obtain the final common question decoding distribution at test time  $t$ .

$$\pi_\theta^t = \frac{1}{C} \left[ \sum_{i \in \mathcal{D}^+} w_{i,\theta}^t \pi_i^t - \eta(z_\theta^t) \cdot \sum_{i \in \mathcal{D}^-} v_{i,\theta}^t \pi_i^t \right]_+ \quad (10)$$

where  $\theta$  is the coordinator's parameters, the subscript  $+$  is an operator that selects non-negative weighted tokens, and  $C$  is the normalizing factor into a distribution. The sequentially decoded

<sup>1</sup><http://www.ms-marco.org/leaders.aspx>



common question words (partial sequence) are concatenated to all input documents in  $\mathcal{D}^+$  and  $\mathcal{D}^-$  followed by EOS token, to obtain new hidden states and decoding distributions. The decoding process is repeated until the generation is complete.

**Policy Gradient Loss:** The policy gradient loss is defined as follows:

$$\mathcal{L}_{PG}(\theta) = -\mathbb{E} \left[ (R(\tilde{q}|\mathcal{D}^+, \mathcal{D}^-) - R_{\text{baseline}}) \cdot \sum_t \log \pi_{\theta}^t(o_t | \tilde{q}_{<t}, G, \mathcal{D}^+, \mathcal{D}^-) \right] \quad (11)$$

With a complete generation  $\tilde{q}$ , a terminal retrieval statistics reward is computed from the *Ranker* scores and score memberships, noted as  $R(\tilde{q}|\mathcal{D}^+, \mathcal{D}^-)$ . This reward weights the sum of log-likelihoods of generating the observed words  $o_t$  given the generation so far  $\tilde{q}_{<t}$ , from the underlying generator  $G$ , and the two document sets.

**Set-induced Contrastive Regularization:** We further propose an auxiliary to provide richer signals when optimizing the coordinator model. The intuition is that we would like to encourage the coordinator model to generate common questions *toward* the positive set  $\mathcal{D}^+$  relative to the negative set  $\mathcal{D}^-$ . We name the regularizer as *Set-induced Contrastive Regularization* (SCR) because the decoding distributions from  $\mathcal{D}^+$  and  $\mathcal{D}^-$  guide the coordinator to learn to make contrasts between the two sets. Although the decoding distributions from  $\mathcal{D}^+$  and  $\mathcal{D}^-$  are not gold supervision signals, modifying distributional distance toward or away from them helps regulate specificity to  $\mathcal{D}^+$ . The former idea can be formulated as *minimizing* the KL-divergence, evaluated at timestep  $t$ :

$$\min_{\theta} \mathcal{L}_{KL,t}^{\text{pos}}(\theta) = \min_{\theta} \sum_{i \in \mathcal{D}^+} \left[ D_{KL}(\pi_{\theta}^t || \pi_i^t) + D_{KL}(\pi_i^t || \pi_{\theta}^t) \right] \quad (12)$$

We minimize both the forward and the reverse KL divergence since the forward KL does not penalize high mass of  $\pi_{\theta}$  where  $\pi_i$  does not. Likewise for the reverse KL. On the other hand, the latter idea can be formulated as *maximizing* the KL-divergence against the negative set, evaluated at time step  $t$ :

$$\max_{\theta} \mathcal{L}_{KL,t}^{\text{neg}}(\theta) = \max_{\theta} \sum_{i \in \mathcal{D}^-} \left[ D_{KL}(\pi_{\theta}^t || \pi_i^t) + D_{KL}(\pi_i^t || \pi_{\theta}^t) \right] \quad (13)$$

However, we need to *cap* the negative set penalty rather than naively maximizing it, more restrictively if the positive set and the negative sets are semantically close. Intuition is that if the KL divergence against the negative set is too large, then we do not penalize further. Therefore, we define our contrastive regularization function as follows:

$$\mathcal{L}_{SCR}(\theta) = \frac{1}{T} \sum_{t=1}^T \left[ \mathcal{L}_{KL,t}^{\text{pos}}(\theta) - \mathcal{L}_{KL,t}^{\text{neg}}(\theta) \cdot \mathbb{1}_{\nu_t \cdot \mathcal{L}_{KL,t}^{\text{neg}}(\theta) > \mathcal{L}_{KL,t}^{\text{pos}}(\theta)} \right] \quad (14)$$

where  $T$  is the length of the completed generation, and  $\nu_t$  is the similarity measure between positive and negative sets at decoding time  $t$ . Specifically,

$$\nu_t = \cos \text{sim} \left( \frac{1}{|\mathcal{D}^+|} \sum_{i \in \mathcal{D}^+} \pi_i^t, \frac{1}{|\mathcal{D}^-|} \sum_{i \in \mathcal{D}^-} \pi_i^t \right) \quad (15)$$

**Entropy Loss:** We add negative entropy loss  $\mathcal{L}_H$  across the attention weights  $w$  and  $v$ , averaged over  $T$  to encourage the model attend to all the documents rather than attend to a small subset of the documents and risk losing positive and negative set representational information.

$$\mathcal{L}_H(\theta) = \frac{1}{T} \sum_{t=1}^T \left[ \sum_{i \in \mathcal{D}^+} w_{i,\theta}^t \log w_{i,\theta}^t + \sum_{i \in \mathcal{D}^-} v_{i,\theta}^t \log v_{i,\theta}^t \right] \quad (16)$$

We finally optimize for the following loss:

$$\mathcal{L}(\theta) = \lambda_1 \mathcal{L}_{PG}(\theta) + \lambda_2 \mathcal{L}_{SCR}(\theta) + \lambda_3 \mathcal{L}_H(\theta) \quad (17)$$

where  $\lambda_{1,2,3}$  are the scaling hyper-parameters.

## 5 Experiments

**Dataset:** We use the MS-MARCO Q&A dataset (Nguyen et al., 2016) where for the Bing query  $q$ , we consider the top-10 retrieved documents as our positive set  $\mathcal{D}^+$ . To get our negative set  $\mathcal{D}^-$ , we use the MS-MARCO-Conversational Search<sup>2</sup> dataset to first find a query  $q'$  that is similar to  $q$  yet not a paraphrase and consider the top-10 documents retrieved for  $q'$  as our negative set  $\mathcal{D}^-$ . In total, we gather 100K/10K/10K training,

<sup>2</sup><https://github.com/microsoft/MSMARCO-Conversational-Search>

Table 1: Retrieval performance. “Out-Sample IR” refers to the evaluation data sample that consists of 10+10 documents  $\mathcal{D}^+$  and  $\mathcal{D}^-$ . “Search-Engine Augmented IR” refers to augmenting the out-sample into 100 documents in total through Lucene.

Model	Out-Sample IR				Search-Engine Augmented IR				
	mAP	RPrec	MRR (=MRR@10)	nDCG	mAP	RPrec	MRR	MRR@10	nDCG
Top-TFIDF @100	0.416	0.533	0.696	0.545	0.113	0.0588	0.0260	0.0050	0.181
Top-Frequent @100	0.680	0.742	0.921	0.779	0.171	0.129	0.0404	0.0119	0.204
MSQG (Cho et al. '19)	-	-	-	-	-	-	0.0704	<b>0.0441</b>	0.234
MSQG <sub>GPT2</sub>	0.713	0.763	0.945	0.804	0.245	0.217	0.0714	0.0400	0.240
MSCQG <sub>SCR</sub>	0.751	0.790	0.974	0.836	0.258	0.234	0.0745	0.0420	0.245
MSCQG <sub>PG</sub>	0.753	0.791	0.978	0.838	0.256	0.232	0.0742	0.0421	0.244
MSCQG <sub>PG+SCR</sub>	<b>0.767</b>	<b>0.803</b>	<b>0.981</b>	<b>0.849</b>	<b>0.265</b>	<b>0.242</b>	0.0748	0.0420	0.245
MSCQG <sub>PG+SCR+H</sub>	0.765	0.800	0.976	0.847	0.262	0.239	<b>0.0759</b>	0.0434	<b>0.246</b>
Oracle Questions for $\mathcal{D}^+$	0.759	0.797	0.976	0.842	0.292	0.273	0.0846	0.0495	0.256

Table 2: Comparison against the *oracle* MARCO questions for  $\mathcal{D}^+$ . Since retrieval scores cannot give a complete picture of the generation, we aim to understand how close the generations are in terms of various metrics. The numbers show that our proposed model generates common questions similar to the *oracle* MARCO questions. **Notations:** BL for BLEU; ST for Skip-Thought similarity; EM for Embedding Mean similarity; VE for Vector Extrema similarity; and GM for Greedy Matching.

	BL-1	BL-2	BL-3	BL-4	METEOR	ROUGE_L	CIDEr	ST	EM	VE	GM
Oracle Question for $\mathcal{D}^-$	0.449	0.291	0.177	0.100	0.215	0.428	1.076	0.547	0.766	0.617	0.697
Top-TFIDF @100	0.253	0.157	0.104	0.075	0.195	0.339	1.174	0.470	0.747	0.575	0.671
Top-Frequent @100	0.438	0.328	0.260	0.217	0.281	0.476	2.684	0.573	0.799	0.682	0.735
MSQG <sub>GPT2</sub>	0.457	0.313	0.207	0.139	0.282	0.494	1.993	0.563	0.814	0.705	0.768
MSCQG <sub>SCR</sub>	0.501	0.363	0.260	0.193	0.303	0.535	2.533	0.604	0.829	0.729	0.786
MSCQG <sub>PG</sub>	0.562	0.418	0.310	0.234	0.304	0.565	2.702	0.630	0.844	0.734	0.798
MSCQG <sub>PG+SCR</sub>	<b>0.589</b>	<b>0.449</b>	<b>0.339</b>	<b>0.262</b>	<b>0.323</b>	<b>0.591</b>	<b>2.994</b>	<b>0.647</b>	<b>0.858</b>	<b>0.759</b>	<b>0.815</b>
MSCQG <sub>PG+SCR+H</sub>	0.573	0.436	0.330	0.255	0.321	0.583	2.946	0.641	0.851	0.752	0.808

development, and evaluation data points. Details of the pre-processing and experimental configuration are in the Supplementary Materials (SM).

**Automatic evaluation:** The generated common questions are evaluated through standard retrieval-based metrics: MRR and MRR10 (Voorhees, 1999; Radev et al., 2002a), nDCG (Järvelin and Kekäläinen, 2002), precision, mAP. These metrics are computed from the 10 positive and 10 negative document sets (Out-Sample IR). In addition, for each generated question, we use Lucene<sup>3</sup> to retrieve the most challenging 100 MARCO documents via BM25 (Robertson and Zaragoza, 2009), and compute the retrieval statistics (Search-Engine Augmented IR), done in MS-MARCO Retrieval task.

The generated questions are also evaluated in terms of BLEU (Papineni et al., 2002), ROUGE (Lin and Hovy, 2003), METEOR (Banerjee and Lavie, 2005), CIDEr (Vedantam et al., 2015), Greedy Matching (Rus and Lintean, 2012), Skip-Thought (Kiros et al., 2015), Embedding Average (Kenter et al., 2016) and Vector Extrema

Table 3: Pairwise comparison and  $d^+/d^-$  comparison of human evaluation. M=MSCQG, B=MSQG<sub>GPT2</sub>, O=Oracle. Preferences are expressed in percentage (%). Comparison results are statistically significant ( $p < 0.01$ ) unless indicated \*. Ans., Rel., Flu. and Ovr. denotes Answerability, Relevancy and Fluency and Overall, respectively.

Criteria	Pair (M vs. B)			Pair (M vs. O)			Pair (B vs. O)		
	M	B	=	M	O	=	B	O	=
Ans.	<b>52.2</b>	17.5	30.3	39.2	19.8	<b>41.0</b>	32.3	<b>42.5</b>	25.2
Rel.	<b>53.3</b>	18.7	28.0	35.2	22.2	<b>42.7</b>	31.7	<b>44.2</b>	24.2
Flu.	<b>49.3</b>	22.3	28.3	<b>50.8</b>	24.7	24.5	<b>43.7</b>	32.7	23.7
Ovr.	<b>57.5</b>	21.3	21.2	<b>49.5</b>	27.0	23.5	38.3	<b>42.8*</b>	18.8

Criteria	M			B			O		
	$d^+$	$d^-$	=	$d^+$	$d^-$	=	$d^+$	$d^-$	=
Ans.	<b>70.7</b>	10.2	19.2	<b>61.2</b>	14.0	24.8	<b>67.2</b>	13.7	19.2
Rel.	<b>72.2</b>	11.3	16.5	<b>62.2</b>	16.3	21.5	<b>70.3</b>	14.3	15.3
Ovr.	<b>72.0</b>	10.5	17.5	<b>63.2</b>	14.2	22.7	<b>69.0</b>	14.0	17.0

(Forgues et al., 2014) cosine similarities.

**Human evaluation:** We conduct human evaluation through Amazon Mechanical Turk where we evaluate questions generated by MSCQG, MSQG<sub>GPT2</sub>, and the *oracle* question in four criteria: *fluency*, *relevancy*, *answerability*, and *overall*. We randomly select 600 ( $d, q_A, q_B$ ) tuples where

<sup>3</sup><https://lucene.apache.org/>

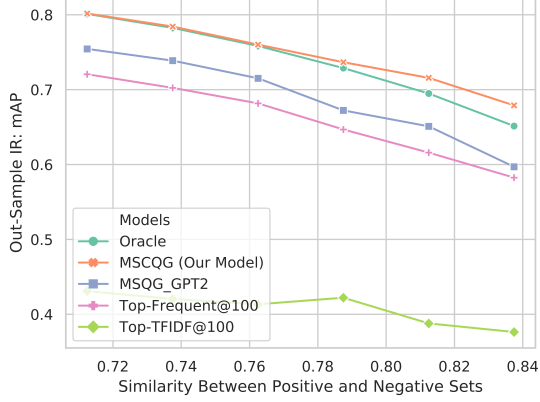


Figure 3: Out-Sample IR: mAP among  $\mathcal{Q}^+$  and  $\mathcal{Q}^-$

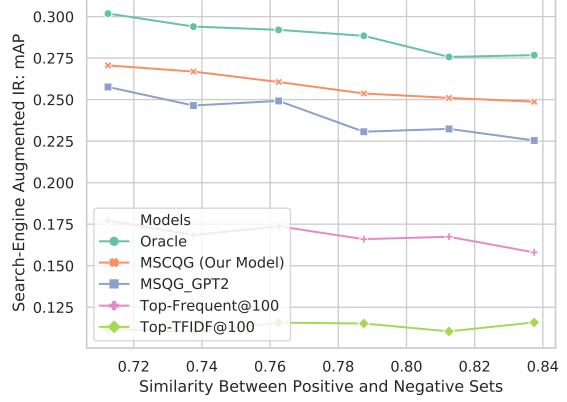


Figure 4: Search-Engine Augmented IR: mAP

Figure 3 shows that our model  $\text{MSCQG}_{PG+SCR+H}$  outperforms the *oracle* questions by a small margin on the Out-Sample IR. In the larger retrieval evaluation using Lucene, it performs subpar against the *oracle* questions, but performs significantly better than all the considered baseline models, shown in Figure 4.

the  $d$  is any from  $\mathcal{Q}^+$  and  $q_A, q_B$  from the three questions. In addition, we evaluate 600  $(d^+, d^-, q)$  tuples where given a question,  $d^+, d^-$  are randomly chosen from  $\mathcal{Q}^+$  and  $\mathcal{Q}^-$ . Each sample is judged by 3 crowd-sourced workers who passed a spam-detection screening, totalling 3,600 samples to obtain reliable results. For details, see SM.

### 5.1 Baseline models

**Multi-Source Question Generator:** This model  $\text{MSQG}_{GPT2}$  is similar to MSQG in Cho et al. (2019b). It processes individual documents in parallel through the fine-tuned GPT-2 generator, rather than RNN-based Seq2Seq model in MSQG, and averages the decoding distributions at test time  $t$ .

$$\pi_{\text{MSQG}_{GPT2}}^t = \frac{1}{|\text{pos}|} \sum_i^{\text{pos}} \pi_i^t \quad (18)$$

Unlike MSQG in Cho et al. (2019b), no further heuristic modifications are made to the model.

**Top-TFIDF@K:** Why do we not simply retrieve the top question implied by the 10 positive documents? To this end, we design a retrieval baseline using the learned TF-IDF (Luhn, 1957; Jones, 1972; Salton and McGill, 1983) weights. This baseline re-evaluates the collection of retrieved questions against all documents in  $\mathcal{Q}^+$  using TF-IDF, and retrieves the most relevant question. For pseudo-code details, see SM.

**Top-Frequent@K:** Another retrieval model is to find an intersecting subset among all the 10 top- $k$  question sets. For pseudo-code details, see SM.

### 5.2 Results and Analysis

**Model comparison and ablation study:** For simplicity, we abuse the term *oracle* by calling the ground-truth question that retrieves  $\mathcal{Q}^+$  when constructing the dataset as the *oracle* question. However, these questions are not gold questions as they might not be the most relevant and specific questions to the given positive and negative sets.

Table 1 shows that our proposed model is effective at generating common questions given multiple documents. In particular, it shows that policy gradient or set-induced contrastive regularization alone is effective in improving performance. The coordinator performs better when optimized for both policy gradient and regularization objectives.

The retrieval results for the questions that initially clustered  $\mathcal{Q}^+$  sets are presented. Note that these are not *gold* questions because in most cases not all the retrieved documents in  $\mathcal{Q}^+$  answer the questions. For clarity of our presentation, we abuse the term and name them as *oracle* questions. Automatic metric results show that our methods are upper-bounded by the *oracle* MARCO questions.

Entropy regularization improves the search-engine augmented IR scores, in particular, MRR. However, it is not crucial as supplemented by Table 2. For additional results, see Table 4 in SM.

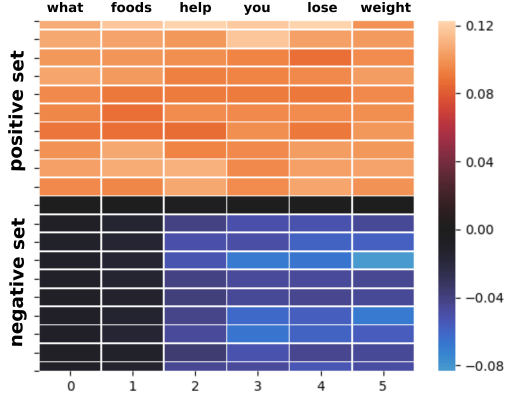


Figure 5: Example attention weights #1

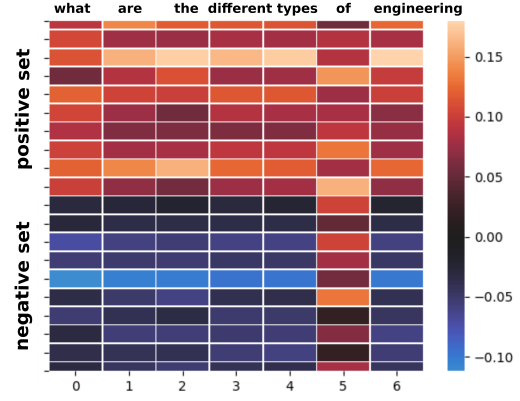


Figure 6: Example attention weights #2

Visualization of sequential attention weights. In the vertical axis, 0-9 indices indicate documents in  $\mathcal{D}^+$ , and 10-19 in  $\mathcal{D}^-$ . The displayed negative weights  $v$  are adjusted by  $\eta(z)$ , see equation 10. Figure 5 shows that the model learns to push the sequential generation semantics more toward  $\mathcal{D}^+$  by gradually penalizing  $\mathcal{D}^-$ . Figure 6 shows that frequent and semantically less distinguishing words such as ‘*of*’ are encouraged even by  $\mathcal{D}^-$ , which empirically aligns with our intuition for TF-IDF.

**Model performance v.s. similarities between  $\mathcal{D}^+$  and  $\mathcal{D}^-$ :**  $\cos \text{sim}(\mathcal{D}^+, \mathcal{D}^-)$  is approximated using the *oracle* questions that are available in the dataset. The similarity is computed by the cosine similarity of the two GEN-Encoder (Zhang et al., 2019) representations. Figures 3 and 4 show that our model generated common questions are more grounded on the positive documents than the baseline model generations. The more similar the two sets  $\mathcal{D}^+$  and  $\mathcal{D}^-$ , the more difficult for the models, even humans, to distinguish which document is more relevant, if not answerable, given the generated question. The model outperforms the baseline model uniformly across different similarities between  $\mathcal{D}^+$  and  $\mathcal{D}^-$ .

**Role of  $\mathcal{D}^-$  by visualizing  $w$ ,  $v$ , and  $z$ :** Figures 5 and 6 show that our model MSCQG learns to gradually penalize  $\mathcal{D}^-$  as it sequentially generates words that are more grounded on  $\mathcal{D}^+$ . Notice the roughly uniform weights across  $\mathcal{D}^+$  but increasing penalization weights across  $\mathcal{D}^-$ , in decoding time.

$\eta$ , which is controlled by the  $z$ , is learned to encourage, rather than discourage, certain words during decoding. We observe that words that are not semantically distinguishing between  $\mathcal{D}^+$  and  $\mathcal{D}^-$ , are encouraged by the coordinator to maintain readability. For example, the weights of the word *of* is mostly non-negative, whereas weights for other words are mostly negative. This indicates that the coordinator learns to selectively

activate/suppress decoding of certain words by coordinating information from  $\mathcal{D}^+$  and  $\mathcal{D}^-$ .

**Human judgments:** Table 3 shows that our model significantly outperforms the strong baseline in every aspect. The results are statistically significant as marked, drawn from the large number of evaluations. Furthermore, we draw a more favorable conclusion toward our model-generated questions when compared against the *oracle* questions than from the automatic metrics, which are approximate yet reasonable metrics. The pairwise agreement between judges is  $54\% \pm 1\%$ . The Cohen’s Kappa score is  $0.19 \pm 0.01$ . Note that this is a reasonable number given the “*same*” or ambiguous option in pairwise comparisons. *Ranker* achieves a relatively high Pearson correlation of 0.6 with respect to human evaluation. For details, see SM.

## 6 Conclusion

We proposed a novel coordinator model that can generate common questions that are more grounded on documents of interest. This coordinator model consists of transformer blocks, and is trained through reinforcement learning and an effective auxiliary loss: Set-induced Contrastive Regularization (SCR). Experiment results show that our model significantly outperforms the previous neural generation model as well as strong retrieval baselines in automatic and human metrics.



## References

- Mohammad Aliannejadi, Hamed Zamani, Fabio Crestani, and Bruce Croft. 2019a. Asking clarifying questions in open-domain information-seeking conversations. In *SIGIR '19*.
- Mohammad Aliannejadi, Hamed Zamani, Fabio Crestani, and W. Bruce Croft. 2019b. [Asking clarifying questions in open-domain information-seeking conversations](#). In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR'19, page 475–484, New York, NY, USA. Association for Computing Machinery.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. [Neural machine translation by jointly learning to align and translate](#). *CoRR*, abs/1409.0473.
- Satanjeev 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*, pages 65–72.
- Delphine Bernhard. 2010. [Query expansion based on pseudo relevance feedback from definition clusters](#). In *Coling 2010: Posters*, pages 54–62, Beijing, China. Coling 2010 Organizing Committee.
- Ondrej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, et al. 2014. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the ninth workshop on statistical machine translation*, pages 12–58.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. [Translating embeddings for modeling multi-relational data](#). In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2787–2795. Curran Associates, Inc.
- Pavel Braslavski, Denis Savenkov, Eugene Agichtein, and Alina Dubatovka. 2017. [What do you mean exactly? analyzing clarification questions in cqa](#). In *Proceedings of the 2017 Conference on Conference Human Information Interaction and Retrieval*, CHIIR '17, page 345–348, New York, NY, USA. Association for Computing Machinery.
- Christian Buck, Jannis Bulian, Massimiliano Ciaramita, Wojciech Gajewski, Andrea Gesmundo, Neil Houlsby, and Wei Wang. 2017. [Ask the right questions: Active question reformulation with reinforcement learning](#).
- Guihong Cao, Jian-Yun Nie, Jianfeng Gao, and Stephen Robertson. 2008. [Selecting good expansion terms for pseudo-relevance feedback](#). In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '08, pages 243–250, New York, NY, USA. ACM.
- Asli Celikyilmaz, Antoine Bosselut, Xiaodong He, and Yejin Choi. 2018. [Deep communicating agents for abstractive summarization](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1662–1675, New Orleans, Louisiana. Association for Computational Linguistics.
- Woon Sang Cho, Pengchuan Zhang, Yizhe Zhang, Xijun Li, Michel Galley, Chris Brockett, Mengdi Wang, and Jianfeng Gao. 2019a. Towards coherent and cohesive long-form text generation. In *Proceedings of the First Workshop on Narrative Understanding*, pages 1–11.
- Woon Sang Cho, Yizhe Zhang, Sudha Rao, Chris Brockett, and Sungjin Lee. 2019b. Generating a common question from multiple documents using multi-source encoder-decoder models. In *The 3rd Workshop on Neural Generation and Translation*.
- Eric Chu and Peter J. Liu. 2018. [Unsupervised neural multi-document abstractive summarization](#). *CoRR*, abs/1810.05739.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. [ELECTRA: Pre-training text encoders as discriminators rather than generators](#). In *ICLR*.
- Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. 2011. [Natural language processing \(almost\) from scratch](#). *Journal of Machine Learning Research*, 12(76):2493–2537.
- Hal Daumé, John Langford, and Daniel Marcu. 2009. [Search-based structured prediction](#). *Machine Learning*, 75(3):297–325.
- Yuntian Deng, Anton Bakhtin, Myle Ott, Arthur Szlam, and Marc'Aurelio Ranzato. 2020. [Residual energy-based models for text generation](#). In *International Conference on Learning Representations*.
- Nina Dethlefs and Heriberto Cuayáhuitl. 2010. Hierarchical reinforcement learning for adaptive text generation. In *Proceedings of the 6th International Natural Language Generation Conference*, pages 37–45. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018a. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018b. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.

- Rui Dong and David Smith. 2018. [Multi-input attention for unsupervised OCR correction](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2363–2372, Melbourne, Australia. Association for Computational Linguistics.
- Marco Dorigo and Marco Colombetti. 1994. Robot shaping: Developing autonomous agents through learning. *Artificial intelligence*, 71(2):321–370.
- Xinya Du, Junru Shao, and Claire Cardie. 2017a. Learning to ask: Neural question generation for reading comprehension. *arXiv preprint arXiv:1705.00106*.
- Xinya Du, Junru Shao, and Claire Cardie. 2017b. [Learning to ask: Neural question generation for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1342–1352, Vancouver, Canada. Association for Computational Linguistics.
- Nan Duan, Duyu Tang, Peng Chen, and Ming Zhou. 2017a. Question generation for question answering. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 866–874.
- Nan Duan, Duyu Tang, Peng Chen, and Ming Zhou. 2017b. [Question generation for question answering](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 866–874, Copenhagen, Denmark. Association for Computational Linguistics.
- Zhihao Fan, Zhongyu Wei, Siyuan Wang, Yang Liu, and Xuanjing Huang. 2018. [A reinforcement learning framework for natural question generation using bi-discriminators](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1763–1774, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Orhan Firat, Baskaran Sankaran, Yaser Al-Onaizan, Fatos T. Yarman-Vural, and Kyunghyun Cho. 2016. [Zero-resource translation with multi-lingual neural machine translation](#). *CoRR*, abs/1606.04164.
- Gabriel Forgues, Joelle Pineau, Jean-Marie Larchevêque, and Réal Tremblay. 2014. Bootstrapping dialog systems with word embeddings. In *Nips, modern machine learning and natural language processing workshop*, volume 2.
- Jianfeng Gao, Michel Galley, and Lihong Li. 2018. Neural approaches to conversational AI. *arXiv preprint arXiv:1809.08267*.
- Ekaterina Garmash and Christof Monz. 2016. [Ensemble learning for multi-source neural machine translation](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1409–1418, Osaka, Japan. The COLING 2016 Organizing Committee.
- Michael Gutmann and Aapo Hyvärinen. 2010. [Noise-contrastive estimation: A new estimation principle for unnormalized statistical models](#). In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 297–304, Chia Laguna Resort, Sardinia, Italy. PMLR.
- Vrindavan Harrison and Marilyn Walker. 2018. [Neural generation of diverse questions using answer focus, contextual and linguistic features](#).
- Michael Heilman. 2011. *Automatic factual question generation from text*. Ph.D. thesis, Carnegie Mellon University.
- Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Philip Bachman, Adam Trischler, and Yoshua Bengio. 2019. [Learning deep representations by mutual information estimation and maximization](#). In *ICLR 2019*. ICLR.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural Comput.*, 9(8):1735–1780.
- Srinivasan Janarthanam and Oliver Lemon. 2009. [Learning lexical alignment policies for generating referring expressions for spoken dialogue systems](#). In *Proceedings of the 12th European Workshop on Natural Language Generation (ENLG 2009)*, pages 74–81, Athens, Greece. Association for Computational Linguistics.
- Natasha Jaques, Shixiang Gu, Dzmitry Bahdanau, José Miguel Hernández-Lobato, Richard E Turner, and Douglas Eck. 2017. Sequence tutor: Conservative fine-tuning of sequence generation models with kl-control. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1645–1654. JMLR. org.
- Kalervo Järvelin and Jaana Kekäläinen. 2002. [Cumulated gain-based evaluation of ir techniques](#). *ACM Trans. Inf. Syst.*, 20(4):422–446.
- Karen Spärck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation*, 28:11–21.
- Nal Kalchbrenner and Phil Blunsom. 2013. Recurrent continuous translation models. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1700–1709.
- Tom Kenter, Alexey Borisov, and Maarten De Rijke. 2016. Siamese cbow: Optimizing word embeddings for sentence representations. *arXiv preprint arXiv:1606.04640*.
- Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980.

- Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S Zemel, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2015. Skip-thought vectors. *arXiv preprint arXiv:1506.06726*.
- Igor Labutov, Sumit Basu, and Lucy Vanderwende. 2015. [Deep questions without deep understanding](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 889–898, Beijing, China. Association for Computational Linguistics.
- Logan Lebanoff, Kaiqiang Song, and Fei Liu. 2018. [Adapting the neural encoder-decoder framework from single to multi-document summarization](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4131–4141, Brussels, Belgium. Association for Computational Linguistics.
- Oliver Lemon. 2008. Adaptive natural language generation in dialogue using reinforcement learning. In *Proceedings of the 12th SEMdial Workshop on the Semantics and Pragmatics of Dialogues*, pages 149–156.
- Canjia Li, Yingfei Sun, Ben He, Le Wang, Kai Hui, Andrew Yates, Le Sun, and Jungang Xu. 2018a. [NPRF: A neural pseudo relevance feedback framework for ad-hoc information retrieval](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4482–4491, Brussels, Belgium. Association for Computational Linguistics.
- Ruizhi Li, Xiaofei Wang, Sri Harish Reddy Mallidi, Takaaki Hori, Shinji Watanabe, and Hynek Herman sky. 2018b. [Multi-encoder multi-resolution framework for end-to-end speech recognition](#). *CoRR*, abs/1811.04897.
- Jindřich Libovický and Jindřich Helcl. 2017. [Attention strategies for multi-source sequence-to-sequence learning](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 196–202, Vancouver, Canada. Association for Computational Linguistics.
- Jindřich Libovický, Jindřich Helcl, and David Mareček. 2018. [Input combination strategies for multi-source transformer decoder](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 253–260, Belgium, Brussels. Association for Computational Linguistics.
- Chin-Yew Lin and Eduard Hovy. 2003. Automatic evaluation of summaries using n-gram co-occurrence statistics. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1, NAACL '03*, pages 71–78, Stroudsburg, PA, USA.
- Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. *International Conference on Learning Representations*.
- Yang Liu and Mirala Lapata. 2019. Hierarchical transformers for multidocument summarization. In *ACL*.
- Yi Luan, Yangfeng Ji, Hannaneh Hajishirzi, and Boyang Li. 2016. Multiplicative representations for unsupervised semantic role induction. In *ACL*.
- Hans Peter Luhn. 1957. A statistical approach to mechanized encoding and searching of literary information. *IBM Journal of research and development*, 1(4):309–317.
- Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. [Effective approaches to attention-based neural machine translation](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1412–1421, Lisbon, Portugal. Association for Computational Linguistics.
- Donald Metzler and W. Bruce Croft. 2007. [Latent concept expansion using markov random fields](#). In *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '07*, pages 311–318, New York, NY, USA. ACM.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositional-ity. In *Advances in neural information processing systems*, pages 3111–3119.
- Nasrin Mostafazadeh, Chris Brockett, Bill Dolan, Michel Galley, Jianfeng Gao, Georgios Spithourakis, and Lucy Vanderwende. 2017. [Image-grounded conversations: Multimodal context for natural question and response generation](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 462–472, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Vinod Nair and Geoffrey E. Hinton. 2010. [Rectified linear units improve restricted boltzmann machines](#). In *Proceedings of the 27th International Conference on Machine Learning, ICML'10*, pages 807–814, USA. Omnipress.
- Andrew Y Ng, Daishi Harada, and Stuart Russell. 1999. Policy invariance under reward transformations: Theory and application to reward shaping. In *ICML*, volume 99, pages 278–287.
- 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 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, December 9, 2016.
- Kyosuke Nishida, Itsumi Saito, Kosuke Nishida, Kazutoshi Shinoda, Atsushi Otsuka, Hisako Asano, and Junji Tomita. 2019. [Multi-style generative reading comprehension](#). *CoRR*, abs/1901.02262.
- Yuta Nishimura, Katsuhito Sudoh, Graham Neubig, and Satoshi Nakamura. 2018. [Multi-source neural machine translation with missing data](#). In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 92–99, Melbourne, Australia. Association for Computational Linguistics.
- Rodrigo Nogueira and Kyunghyun Cho. 2017. [Task-oriented query reformulation with reinforcement learning](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 574–583, Copenhagen, Denmark. Association for Computational Linguistics.
- Rodrigo Nogueira and Kyunghyun Cho. 2019. [Passage re-ranking with BERT](#). *CoRR*, abs/1901.04085.
- Franz Josef Och and Hermann Ney. 2001. Statistical multi-source translation. In *In MT Summit 2001*, pages 253–258.
- Liangming Pan, Wenqiang Lei, Tat-Seng Chua, and Min-Yen Kan. 2019. [Recent advances in neural question generation](#). *CoRR*, abs/1905.08949.
- 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 on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Daraksha Parveen and Michael Strube. 2014. Multi-document summarization using bipartite graphs. In *TextGraphs-9: the workshop on Graph-based Methods for Natural Language Processing*.
- Ramakanth Pasunuru and Mohit Bansal. 2017. Reinforced video captioning with entailment rewards. *arXiv preprint arXiv:1708.02300*.
- Romain Paulus, Caiming Xiong, and Richard Socher. 2017. [A deep reinforced model for abstractive summarization](#). *CoRR*, abs/1705.04304.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, pages 1532–1543.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL*.
- Dragomir Radev. 2000. A common theory of information fusion from multiple text sources step one: Cross-document structure. In *1st SIGdial Workshop on Discourse and Dialogue*, pages 78–83.
- Dragomir R Radev, Hong Qi, Harris Wu, and Weiguo Fan. 2002a. Evaluating web-based question answering systems. In *LREC*.
- Dragomir R. Radev, Hong Qi, Harris Wu, and Weiguo Fan. 2002b. [Evaluating web-based question answering systems](#). In *Proceedings of the Third International Conference on Language Resources and Evaluation (LREC’02)*, Las Palmas, Canary Islands - Spain. European Language Resources Association (ELRA).
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Filip Radlinski and Nick Craswell. 2017. [A theoretical framework for conversational search](#). In *Proceedings of the 2017 Conference on Conference Human Information Interaction and Retrieval, CHIIR ’17*, page 117–126, New York, NY, USA. Association for Computing Machinery.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv e-prints*.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. [Sequence level training with recurrent neural networks](#). *CoRR*, abs/1511.06732.
- Sudha Rao and Hal Daumé III. 2018. [Learning to ask good questions: Ranking clarification questions using neural expected value of perfect information](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2737–2746, Melbourne, Australia. Association for Computational Linguistics.
- Sudha Rao and Hal Daumé III. 2019. [Answer-based Adversarial Training for Generating Clarification Questions](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 143–155, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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*, pages 7008–7024.
- Verena Rieser and Oliver Lemon. 2009. Natural language generation as planning under uncertainty for spoken dialogue systems. In *Empirical methods in natural language generation*, pages 105–120. Springer.



- Stephen Robertson and Hugo Zaragoza. 2009. [The probabilistic relevance framework: Bm25 and beyond](#). *Found. Trends Inf. Retr.*, 3(4):333–389.
- Lior Rokach. 2010. [Ensemble-based classifiers](#). *Artificial Intelligence Review*, 33(1):1–39.
- A Rothe, Brenden Lake, and Todd Gureckis. 2016. Asking and evaluating natural language questions. In *Proceedings of the 38th Annual Conference of the Cognitive Science Society*.
- David E Rumelhart, Geoffrey E Hinton, Ronald J Williams, et al. 1988. Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1.
- Vasile Rus and Mihai Lintean. 2012. A comparison of greedy and optimal assessment of natural language student input using word-to-word similarity metrics. In *Proceedings of the Seventh Workshop on Building Educational Applications Using NLP*, pages 157–162. Association for Computational Linguistics.
- Vasile Rus, Brendan Wyse, Paul Piwek, Mihai Lintean, Svetlana Stoyanchev, and Cristian Moldovan. 2010. The first question generation shared task evaluation challenge. In *Proceedings of the 6th International Natural Language Generation Conference*, pages 251–257. Association for Computational Linguistics.
- G. Salton. 1971. *The SMART Retrieval System—Experiments in Automatic Document Processing*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- Gerard Salton and Michael J McGill. 1983. *Introduction to modern information retrieval*. mcgraw-hill.
- 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 (Volume 1: Long Papers)*, pages 1073–1083, Vancouver, Canada. Association for Computational Linguistics.
- 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 (Volume 1: Long Papers)*, pages 588–598, Berlin, Germany. Association for Computational Linguistics.
- Shikhar Sharma, Layla El Asri, Hannes Schulz, and Jeremie Zumer. 2017. [Relevance of unsupervised metrics in task-oriented dialogue for evaluating natural language generation](#). *CoRR*, abs/1706.09799.
- Noah A. Smith and Jason Eisner. 2005. [Contrastive estimation: Training log-linear models on unlabeled data](#). In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pages 354–362, Ann Arbor, Michigan. Association for Computational Linguistics.
- Lin Feng Song, Zhiguo Wang, and Wael Hamza. 2017. [A unified query-based generative model for question generation and question answering](#).
- Lin Feng Song, Zhiguo Wang, Wael Hamza, Yue Zhang, and Daniel Gildea. 2018. [Leveraging context information for natural question generation](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 569–574, New Orleans, Louisiana. Association for Computational Linguistics.
- 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 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3930–3939, Brussels, Belgium. Association for Computational Linguistics.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *NIPS*.
- Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. Policy gradient methods for reinforcement learning with function approximation. In *Proceedings of the 12th International Conference on Neural Information Processing Systems*, NIPS’99, pages 1057–1063. MIT Press.
- Zhaopeng Tu, Zhengdong Lu, Yang Liu, Xiaohua Liu, and Hang Li. 2016. [Modeling coverage for neural machine translation](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 76–85, Berlin, Germany. Association for Computational Linguistics.
- 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 *Advances in neural information processing systems*, pages 5998–6008.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.
- Ellen M. Voorhees. 2001. [The trec question answering track](#). *Nat. Lang. Eng.*, 7(4):361–378.
- EM Voorhees. 1999. Proceedings of the 8th text retrieval conference. *TREC-8 Question Answering Track Report*, pages 77–82.

- Xiaojun Wan. 2008. An exploration of document impact on graph-based multi-document summarization. In *Conference on Empirical Methods in Natural Language Processing*.
- Yizhong Wang, Kai Liu, Jing Liu, Wei He, Yajuan Lyu, Hua Wu, Sujian Li, and Haifeng Wang. 2018. [Multi-passage machine reading comprehension with cross-passage answer verification](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1918–1927, Melbourne, Australia. Association for Computational Linguistics.
- P. J. Werbos. 1990. [Backpropagation through time: what it does and how to do it](#). *Proceedings of the IEEE*, 78(10):1550–1560.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Jinxi Xu and W. Bruce Croft. 1996. [Query expansion using local and global document analysis](#). In *Proceedings of the 19th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’96*, pages 4–11, New York, NY, USA. ACM.
- Ming Yan, Jiangnan Xia, Chen Wu, Bin Bi, Zhongzhou Zhao, Ji Zhang, Luo Si, Rui Wang, Wei Wang, and Haiqing Chen. 2018. [A deep cascade model for multi-document reading comprehension](#). *CoRR*, abs/1811.11374.
- Hamed Zamani, Susan Dumais, Nick Craswell, Paul Bennett, and Gord Lueck. 2020. [Generating clarifying questions for information retrieval](#). In *The Web Conference 2020 (formerly WWW conference)*.
- Chengxiang Zhai and John Lafferty. 2001. [Model-based feedback in the language modeling approach to information retrieval](#). In *Proceedings of the Tenth International Conference on Information and Knowledge Management, CIKM ’01*, pages 403–410, New York, NY, USA. ACM.
- Hongfei Zhang, Xia Song, Chenyan Xiong, Corby Rosset, Paul N. Bennett, Nick Craswell, and Saurabh Tiwary. 2019. [Generic intent representation in web search](#). In *Proceedings of the 42Nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR’19*, pages 65–74, New York, NY, USA. ACM.
- Jianmin Zhang, Jiwei Tan, and Xiaojun Wan. 2018. Adapting neural single-document summarization model for abstractive multi-document summarization: A pilot study. In *International Conference on Natural Language Generation*.
- Mu Zhu. 2004. Recall, precision and average precision. *Department of Statistics and Actuarial Science, University of Waterloo, Waterloo*, 2:30.
- Barret Zoph and Kevin Knight. 2016. [Multi-source neural translation](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 30–34, San Diego, California. Association for Computational Linguistics.

## Supplementary Materials

### A Data Pre-processing Details

**Data Pre-processing:** MS-MARCO Q&A dataset (Nguyen et al., 2016) contains 1,010,916 questions, in which each question is associated with top-10 documents. Each data point contains a question and its top-10 returned documents from the Bing search engine<sup>4</sup>. This question is not a target itself since not all the top-10 retrieved documents answer the question. However, it can give relative evaluation against a model-generated common question based on the top-10 retrieved documents. We target a broader class of problems where only document groups are available but no such group-inducing or *oracle* questions.

In fact, among the top-10 retrieved documents, often one document is labeled ‘*selected*’ by human annotators to indicate that the document answers the question (*true positive*), and left unknown or unlabeled for the rest of the documents, implying they may or may not answer the question (*true negative* or *false negative*). This label information is used to train the underlying generator block of their MSQG model (Cho et al., 2019b). A single *selected* MS-MARCO (Nguyen et al., 2016) document is fed into a long short-term memory-based sequence-to-sequence model to output the corresponding question. An example of the input *selected* document is: *The House of Representatives shall be composed of Members chosen every second Year by the People of the several States.... Article I, Section 2, Clause 1*, and the corresponding question is: *how long is a term for a member of the house of representatives*. We chose this dataset since the question that retrieve the top-10 documents can shed light to relative performance of our model.

To find two 10-document sets  $\mathcal{D}^+$  and  $\mathcal{D}^-$  that are similar, we find a pair of questions that are semantically similar. However, computing pair-wise similarities among roughly 1 million questions is computationally intractable. Therefore, we leverage another dataset: MS-MARCO-Conversational Search<sup>5</sup>: an artificially constructed public dataset that simulate user search sequences.

Each data point or session is an artificial sequence of similar questions grounded on true user behavior. Since many similar questions are

grouped together, we can reduce the search space for finding pairs of similar questions. Then we take pairs of high semantic similarity ( $\geq 0.7$ ) yet not a paraphrase ( $\leq 0.85$  following their classification criteria) using GEN-Encoder (Zhang et al., 2019) which two associated 10-document sets do not have overlaps, primarily for prototype evaluation convenience. For deployment models, one may choose to allow overlaps between two sets for more challenging learning. From the two similar 10-document sets, either one is set to positive  $\mathcal{D}^+$  or negative  $\mathcal{D}^-$ , yielding two data points for the our derived dataset.

These pre-processing steps yield 346,215 data points, each of which contains a pair of positive and negative questions, and positive and negative 10-document sets. Training MSCQG on the entire dataset requires processing about 7 million MARCO documents. This is computationally intensive and takes about two days on 8 Nvidia Tesla V100 GPU cards for a single epoch. Therefore, for building small research prototypes and benchmarks, we will also release a subset of the data, that consists of 100K/10K/10K training, development, and evaluation data points.

#### Data Example:

##### Oracle question for $\mathcal{D}^+$ :

number of saturn’s moons

##### Oracle question for $\mathcal{D}^-$ :

uranus how many moons

##### Positive Set $\mathcal{D}^+$ :

1. *moons of saturn. there are 62 moons orbiting saturn. the moons of saturn vary not only in size but also in composition and shape. the largest of the moons of saturn is the aptly named titan, more than 5,000 km across and is bigger than mercury. there are 7 major moons of saturn and the rest are grouped based on the mythology from which it is taken.*

2. *iapetus with a diameter of 1,470 km, it is the 3rd largest moon of saturn. it was discovered by giovanni cassini in 1671. it has a distinct feature of having a bright and dark hemisphere. dione the 4th largest moon of saturn named after a vague character in greek mythology.*

<sup>4</sup><https://www.bing.com>

<sup>5</sup><https://github.com/microsoft/MSMARCO-Conversational-Search>

3. titan is the largest of saturn's moons and the first to be discovered. titan is the only moon in the solar system known to have a significant atmosphere. nitrogen and methane extend around the moon 10 times as far into space as earth's atmosphere, sometimes falling to the surface in the form of methane rain.

4. saturn has at least 150 moons and moonlets, 53 of which have formal names. titan, the largest, comprises more than 90% of the mass in orbit around saturn, including the rings. saturn's second-largest moon, rhea, may have a tenuous ring system of its own, along with a tenuous atmosphere.

5. their journeys around the ringed planet average from half an earth day to just over four earth years. saturn's moons formed early in the history of the solar system. one of the moons, titan, makes up 96 percent of the mass orbiting the planet. scientists think that the system may have originally housed two such moons, but the second broke up, creating the debris that formed the rings and smaller, inner moons.

6. saturn has a prominent ring system that consists of nine continuous main rings and three discontinuous arcs and that is composed mostly of ice particles with a smaller amount of rocky debris and dust. sixty-two moons are known to orbit saturn, of which fifty-three are officially named.

7. sixteen of the moons are tidally locked, with one face permanently turned toward saturn. the first moon was discovered in 1655. over the next 200 years, the other seven major satellites were spotted. by 1997, astronomers on earth had found 18 moons in orbit around the planet.

8. saturn is the sixth planet from the sun and the second-largest in the solar system, after jupiter. it is a gas giant with an average radius about nine times that of earth. although only one-eighth the average density of earth, with its larger volume saturn is just over 95 times more massive.

9. this temporary name usually consists of the year of discovery and a number indicating the order of discovery in that year. in the case of saturn's moons, these provisory names follow the

format s/2005-s1, s/2005-s2 etc. the first s (before the slash) is for saturn. the second s (after the dash) is for satellite.

10. this does not include the hundreds of moonlets comprising the rings. titan, saturn's largest moon, and the second-largest in the solar system, is larger than the planet mercury, although less massive, and is the only moon in the solar system to have a substantial atmosphere.

#### **Negative Set $\mathfrak{D}^-$ :**

11. uranus has 27 moons that we know of. five of the moons are large and the rest are much smaller. the five large moons are called miranda, ariel, umbriel, titania, and oberon. titania is the largest moon of uranus and it is covered with small craters, a few large craters, and very rough rocks. ariel is the brightest moon of uranus and has canyons and valleys as well as a lot of craters. umbriel is very dark.

12. uranus can't seem to catch a break these days. besides spinning on its side like the drunkard of the solar system and being the butt of everyone's jokes, new research suggests several of its tiny moons will collide in a million years. uranus can't seem to catch a break these days.

13. the gas giant uranus is the third largest planet in our solar system, has many moons, a ring system, and composed of gases and ices. universe today space and astronomy news login

14. the researchers used cressida's mass and orbit to determine its possible doom. since uranus' 27 moons are tightly packed together, the team posits that in a million years, cressida will likely have a deadly encounter with one of its neighboring moons, called desdemona. previous research and simulations suggest cupid and belinda will also probably smack into each other some time between 1,000 and 10 million years from now.

15. puck, at 162 km, is the largest of the inner moons of uranus and the only one imaged by voyager 2 in any detail while puck and mab are the two outermost inner satellites of uranus. all inner moons are dark objects.

16. uranus, which takes its name from the greek



god of the sky, is a gas giant and the seventh planet from our sun. it is also the third largest planet in our solar system, ranking behind jupiter and saturn. like its fellow gas giants, it has many moons, a ring system, and is primarily composed of gases that are believed to surround a solid core.

17. in 1986, the voyager 2 spacecraft hit the jackpot while studying uranus and discovered 10 other moons, including desdemona and cressida. since then, hubble observations have helped bring that number up to 27 for now.

18. at an average distance of 3 billion km from the sun, it takes uranus roughly 84 years (or 30,687 days) to complete a single orbit of the sun. 1 the rotational period of the interior of uranus is 17 hours, 14 minutes. as with all giant planets, its upper atmosphere experiences strong winds in the direction of rotation.

19. uranus' size, mass and orbit: with a mean radius of approximately 25,360 km, a volume of 6.833—1013 km<sup>3</sup>, and a mass of 8.68 — 1025 kg, uranus is approximately 4 times the sizes of earth and 63 times its volume.

20. uranus has 27 known satellites, which are divided into the categories of larger moons, inner moons, and irregular moons (similar to other gas giants). the largest moons of uranus are, in order of size, miranda, ariel, umbriel, oregon and titania.

## B Retrieval Baselines

### Top-TFIDF@K and Top-Frequent@K

The retrieval baselines are designed to give a relative insight into the performance between MSQG in Cho et al. (2019b) and our novel coordinator model. We use Lucene to retrieve *questions* instead of documents from a corpus composed of the 1,010,916 MS-MARCO questions. The retrieved questions from Top-TFIDF@K and Top-Frequent@K baselines are evaluated in the same manner as the generated ones.

For the intersection to be non-empty,  $k$  should be sufficiently large. However, even for  $k = 1000$ , there were no intersecting subset questions for almost all cases. Therefore, we relax the intersection among all 10 retrieved sets, into finding the most frequently occurring question among the 10 top- $k$

retrieved sets.  $k = 100$  was an appropriate value that is not too large to retrieve remotely relevant questions, and not too small to yield vastly different retrieval sets. If there are multiple questions with the same count, we randomly choose one.

---

### Algorithm 1 Top-TFIDF@K

---

**Input:**  $\mathcal{D}^+$ , Corpus  $\mathbb{C}$

For each  $d \in \mathcal{D}^+$ , retrieve top-K questions in  $\mathbb{C}$ ;  
Using all unique questions  $\mathbb{Q}$ , compute TF-IDF;  
Let  $\Psi$  be the TF-IDF transform operator;  
 $q^* = \arg \max_{q \in \mathbb{Q}} \sum_{d \in \mathcal{D}^+} \cos \text{sim}(\Psi_q, \Psi_d)$ ;

**Output:**  $q^*$

---



---

### Algorithm 2 Top-Frequent@K

---

**Input:**  $\mathcal{D}^+$ , Corpus  $\mathbb{C}$

For each  $d \in \mathcal{D}^+$ , retrieve top-K questions in  $\mathbb{C}$ ;  
Let  $\mathbb{S}_d$  be the retrieved set for each  $d$ ;  
 $q^* = \arg \max_{q \in \mathbb{Q}} \sum_{d \in \mathcal{D}^+} \mathbb{1}_{q \in \mathbb{S}_d}$ ;

**Output:**  $q^*$

---

## C Experiment Configurations

**Document-specific GPT-2 Generator:** From each document  $i$ , the generator yields its final layer hidden state  $h_i \in \mathbb{R}^H$  ( $H = 768$  is the hidden dimension) and a document-specific discrete output distribution  $\pi_i \in \mathbb{R}^V$  ( $V = 50257$  is the vocabulary dimension) from the learned language model head.

**Coordinator:** The input size is 20 with the dimensionality of the embeddings and hidden states as 768. The number of recurrent layers is 2, with 4 attention heads in each layer. The epsilon value used in the layer normalization is set to  $1e-5$ . The number of cluster embeddings is 2 (positive or negative). The standard deviation of the truncated normal initializer for weight matrices is 0.02.  $\lambda_1, \lambda_2, \lambda_3 = 1.0, 100.0, 0.1$ . Maximum generation length is 20 tokens. We use the BERT (Devlin et al., 2018a) version of Adam optimizer (Kingma and Ba, 2014) with weight decay of 0.01 and learning rate of  $1e-5$ . We trained the coordinator model by maximizing *Precision@10* with *oracle* questions as the policy gradient baseline. Additional experiment results using a different baseline - self-critic (Rennie et al., 2017) - is shown in Table 4. This shows that our proposed model framework is effective even with any of the two policy gradient

baselines. Conceptually, the coordinator model would generate a common question that can better retrieve the documents from the positive set, aided by the negative set.

## D Human Evaluation Details

We performed two human evaluations: In the first experiment, we showed judges one randomly selected positive document followed by a pair of questions from the three sources. Judges were asked to evaluate which one of the two questions is preferred based on four criteria. For each pair of three sources, we evaluated 200 same random samples for each judge (or 600 samples for the 3 judges), totalling 1,800 samples.

In the second experiment, human annotators evaluated contrastive ability from 1,800 samples of one question, followed by two documents each from the positive and negative sets. Note that our model is trained to generate questions, accounting for the negative set.

The results were averaged across all samples and judges.

For computing the Pearson correlation between the ranker and human evaluation results, we map *Option A preferred*  $\rightarrow 0$ , *Same*  $\rightarrow 0.5$ , *Option B preferred*  $\rightarrow 1$ , accounting for the random assignments between *A* and *B*. This projection ensures that image of two metrics are the same (between 0 and 1). Then we compute the correlation value between two results.

Table 4: Additional retrieval performance using self-critic (Rennie et al., 2017) baseline in the policy gradient, applicable to datasets with no *oracle* questions. It shows that our framework is also effective using a different baseline. The superscript *null-neg* denotes models that do not use negative attentions when generating questions. This shows the importance of the negative set in promoting specificity in the generated question. It further corroborates that the non-uniform weighted-sum scheme among  $\mathcal{D}^+$  improves performance because not all documents in  $\mathcal{D}^+$  revolve around the same topic, and the model learns to address this nature of the dataset through unequal weights and generate a more representative question.

Model	Out-Sample IR				Search-Engine Augmented IR				
	mAP	RPrec	MRR (=MRR@10)	nDCG	mAP	RPrec	MRR	MRR@10	nDCG
Top-TFIDF @100	0.416	0.533	0.696	0.545	0.113	0.0588	0.0260	0.0050	0.181
Top-Frequent @100	0.680	0.742	0.921	0.779	0.171	0.129	0.0404	0.0119	0.204
MSQG (Cho et al. '19)	-	-	-	-	-	-	0.0704	0.0441	0.234
MSQG <sub>GPT2</sub>	0.713	0.763	0.945	0.804	0.245	0.217	0.0714	0.0400	0.240
MSCQG <sub>SCR</sub>	0.751	0.790	0.974	0.836	0.258	0.234	0.0745	0.0420	0.245
MSCQG <sub>PG+SCR+H</sub> <sup>self-critic,null-neg</sup>	0.714	0.764	0.945	0.805	0.247	0.220	0.0724	0.0407	0.241
MSCQG <sub>PG</sub> <sup>self-critic</sup>	<b>0.762</b>	<b>0.798</b>	<b>0.982</b>	<b>0.845</b>	0.259	0.237	0.0746	0.0420	0.244
MSCQG <sub>PG+SCR</sub> <sup>self-critic</sup>	0.760	0.797	0.977	0.843	0.260	0.236	0.0744	0.0416	0.245
MSCQG <sub>PG+SCR+H</sub> <sup>self-critic</sup>	0.760	0.797	0.977	0.843	<b>0.262</b>	<b>0.238</b>	<b>0.0771</b>	<b>0.0444</b>	<b>0.247</b>
MSCQG <sub>PG+SCR+H</sub> <sup>orcl-critic,null-neg</sup>	0.717	0.766	0.950	0.808	0.246	0.220	0.0722	0.0404	0.241
MSCQG <sub>PG</sub> <sup>orcl-critic</sup>	0.753	0.791	0.978	0.838	0.256	0.232	0.0742	0.0421	0.244
MSCQG <sub>PG+SCR</sub> <sup>orcl-critic</sup>	<b>0.767</b>	<b>0.803</b>	<b>0.981</b>	<b>0.849</b>	<b>0.265</b>	<b>0.242</b>	0.0748	0.0420	0.245
MSCQG <sub>PG+SCR+H</sub> <sup>orcl-critic</sup>	0.765	0.800	0.976	0.847	0.262	0.239	<b>0.0759</b>	<b>0.0434</b>	<b>0.246</b>
Oracle Questions for $\mathcal{D}^+$	0.759	0.797	0.976	0.842	0.292	0.273	0.0846	0.0495	0.256