

DeepGen: Diverse Search Ad Generation and Real-Time Customization

Konstantin Golobokov[◇], Junyi Chai[♣], Victor Ye Dong[♣],
Mandy Gu[♣], Bingyu Chi[♣], Jie Cao[♣], Yulan Yan[♣], Yi Liu[♣]

[◇] Azure AI, [♣] Bing Ads, [♠] News & Feeds

Microsoft, Redmond, USA

{FirstName.LastName}@microsoft.com

Abstract

We present DeepGen, a system deployed at web scale for automatically creating sponsored search advertisements (ads) for Bing Ads customers. We leverage state-of-the-art natural language generation (NLG) models to generate fluent ads from advertiser’s web pages in an abstractive fashion and solve practical issues such as factuality and inference speed. In addition, our system creates a customized ad in real-time in response to the user’s search query, therefore highlighting different aspects of the *same* product based on what the user is looking for. To achieve this, our system generates a *diverse* choice of smaller pieces of the ad ahead of time and, at query time, selects the most relevant ones to be stitched into a complete ad. We improve generation diversity by training a controllable NLG model to generate multiple ads for the same web page highlighting different selling points. Our system design further improves diversity horizontally by first running an ensemble of generation models trained with different objectives and then using a diversity sampling algorithm to pick a diverse subset of generation results for online selection. Experimental results show the effectiveness of our proposed system design. Our system is currently deployed in production, serving $\sim 4\%$ of ads globally on Bing.

1 Introduction

Search advertising is the largest segment of digital advertising for its projected \$203B out of \$515B market share worldwide in 2022 (Statista, 2022). Traditionally, advertisers manually create ads for their web pages to start an advertising campaign. There is a growing need to automate this process, either to lessen the burden for small and medium businesses, or to create millions of ads for large businesses that have lots of products.

A classical automated ad generation system relies on extraction rules as described in Section 2.1,

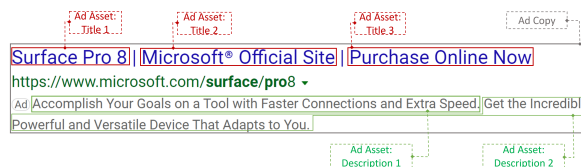


Figure 1: An example of an ad copy (grey box) comprised of ad assets. Red box is used for ad title assets, and green box is used for ad description assets. This ad could be shown for search query “Surface 8”.

for example, extracting key phrases from advertiser’s web pages as ad titles. However, per our experience, extraction-based methods are not very successful in generating the much longer ad description. Refer to Figure 1 for the example ad title and description assets. Therefore, we aim to generate ads in an abstractive fashion. In this work, we focus on improving ad performance from two aspects: factuality and customization.

To achieve the optimal ad performance, our current system creates a customized ad in real-time in response to a user’s search query. As shown in Figure 2, different ads are displayed for different queries, although they are advertising the same web page. We dynamically customize ad copies by stitching the generated ad assets together given the user’s search context, approximating the ultimate goal of real-time customized generation. Our work makes the following contributions:

1. We demonstrate an NLG application that leverages cutting-edge models, which can abstractively generate and instantaneously stitch ad text, matching human quality and achieving real-time ad content customization.
2. We record a significant click-through-rate gain of 13.28% over an extraction-based system as a baseline. Our system is currently deployed at web scale, serving $\sim 4\%$ of ads shown on Bing search engine.

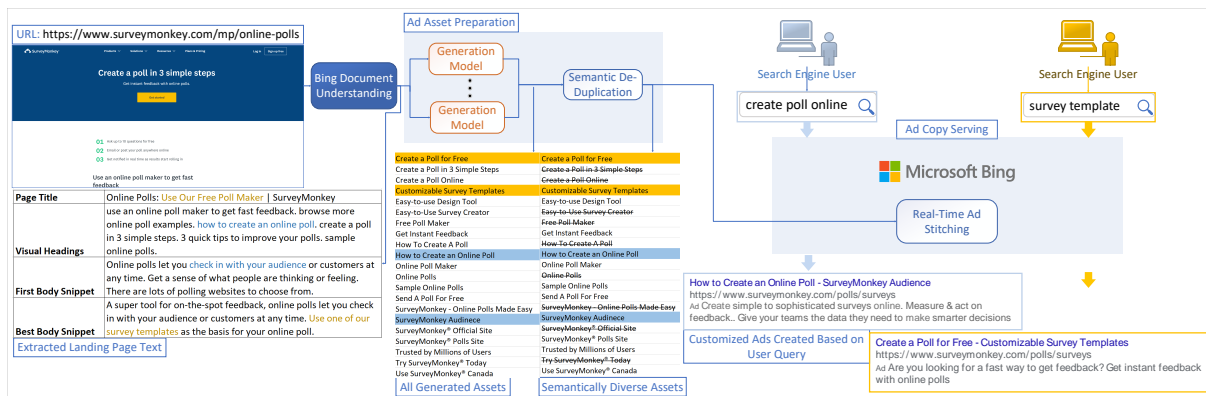


Figure 2: An illustration of the end-to-end DeepGen system. First, multiple ad assets are generated based on various parts of the advertiser’s web page. Semantically diverse ad assets are then selected and prepared for serving. Finally, customized ads are created based on user queries. Transparent blocks are the NLP models, solid blocks are the surrounding infrastructure. Generative models are shown in orange, discussed in Section 2. The rest of the system is presented in Section 3.

2 Ads Generation System

Our system for ad content generation and stitching is automated end-to-end as shown in Figure 2. Advertisers only need to supply us with their domain names, landing page targeting rules, and a bid for each rule (e.g., bid \$0.5 for URLs containing “shoes”). Our Search Indexing infrastructure crawls all landing pages under advertiser domain names that match targeting rules and runs the Document Understanding (DU) pipeline to extract textual information as per Section 2.1. After that, we run multiple NLG models concurrently. This parallel design enables us to scale modeling horizontally: we can add or remove generation models at will. The models can either generate an ad asset or a full ad copy. For a full ad copy we simply split it into assets. At the end of generation stage, we have many title and description assets generated for each advertiser URL.

2.1 Baselines

Extraction-based systems The extraction techniques have evolved in Bing Ads over a decade and we consider them a strong industrial baseline in this paper. This baseline can produce title assets of high quality, but it does not perform as well for the longer description assets. For extraction candidates, we leverage parts of the website extracted by Bing DU pipeline, as per example below:

- Page Title - the document title present in meta-data; <title> tag for HTML documents
- Visual Headings - the visually emphasized

document title present in the document, visible to user

- First/Best Body Snippet - first (top-most)/best document body snippet extracted by Bling (Xiong et al., 2019)

Examples of the above landing page text extracted by DU pipeline can be seen on the left in Figure 2.

Abstractive generation baseline We consider models finetuned directly on advertiser written ad copies as the baseline for abstractive generation approach. We finetune UniLMv2 (Bao et al., 2020) on advertiser-written full ad copies, with learning rate of $5 \cdot 10^{-5}$. We refer to such models as AdCopy models as they generate one ad copy for each source sequence. See Figure 3 for an example of source/target sequences for this task. Multiple AdCopy models were successfully deployed in production with significant business gains (Wang et al., 2021). Some best practices we learned are: 1) advertiser-written ads have a very skewed distribution with some advertiser having millions of template generated ads. Therefore we sample the 3000 URLs with the most ad impressions in the past year per advertiser domain, obtaining 3M-5M training examples; 2) validation and test sets randomly split from training set do not work well; they need to be constructed from different advertisers than those in training set to avoid overfitting. We use validation set of size 300K-500K examples and a test set of 30K-50K examples. We use ROUGE1-F1 (Lin, 2004) on validation set to select the best checkpoint during training.

We inference with beam search of size 5 with code optimization, leveraging Einsum operator in cross-attention stage to avoid the encoder cache copy, per the FastSeq (Yan et al., 2021) implementation. This optimization allows us to increase batch size and brings 5x speed up in our task. Our generation models can be seen in the center of Figure 2 in orange color.

Source Sequence:
 survey monkey **PageTitle** Online Polls: Use Our Free Poll Maker | SurveyMonkey **VisualHeadings** use an online poll maker to get fast feedback. browse more online poll examples. how to create an online poll. create a poll in 3 simple steps. 3 quick tips to improve your polls. **FirstBodySnippet** Online polls let you check in with your audience or customers at any time. Get a sense of what people are thinking or feeling. There are lots of polling websites to choose from. **BestBodySnippet** A super tool for on-the-spot feedback, online polls let you check in with your audience or customers at any time. Use one of our survey templates as the basis for your online poll.

Target Sequence:
 Survey Template **TitleSep** Effective Online Survey Tool **Desc** Create Beautiful Surveys With Typeform For Free. Easy-to-use Online Survey Creator.

Figure 3: An example of a source and generated target sequence pair for the baseline AdCopy model.

2.2 Factuality Improvement

To evaluate the quality of generated ads, we mainly rely on human evaluation. For that, we sample a stratified sample of at most 50 examples per domain, and then uniformly subsample 500 – 1000 examples per human evaluation task. This way, we get an evaluation result from diverse portions of our demand, not letting very large domains dominate. We work with a pool of professional judges, trained to evaluate ads in an unbiased way. We further examine evaluation examples and give feedback to the judges in case there is a misunderstanding of the judgement guidelines. Thus, we evaluate the quality of generated ad texts along the following 4 aspects:

- **Text Quality:** evaluates grammar and style, with levels Good, Fair, Bad, Embarrassing, and Not Scorable.
- **Human Likeness:** whether it looks like human-written, with levels Yes and No.
- **Factuality:** whether the generated information is supported by landing page, with levels Yes and No.
- **Relevance:** whether the generated text is relevant to advertiser’s business, with levels Yes and No.

We define an ad text to be “Overall Good” if it gets “Good” or “Fair” for Text Quality, and “Yes” for Human Likeness, Factuality, and Relevance. Refer to Figure 4 in the Appendix A for an example human judge interface. To be allowed for further

A/B testing, the Overall Good Rate needs to be at least 90% with confidence greater than 97.5%.

As shown in Table 1, our baseline model does not have a significant difference in quality from the advertiser written ads. However, the overall good rate for both is curtailed by lower factuality scores. For example, our AdCopy model can generate popular claims like “Free Shipping” or “15% Discount” which do not exist in the landing page. This is similar to the hallucination issue in abstractive summarization (Filippova, 2020; Maynez et al., 2020b).

To alleviate the extrinsic hallucinations (Maynez et al., 2020a) in our ads, we employ phrase-based cross-check filtering. For that, we use a list of potentially erroneous phrases and patterns obtained by studying human evaluation results for our generated ads. Our approach is similar to entity-based filtering per Nan et al. (2021).

Some cross-check examples are 1) Phrase Check: a list of sensitive or potentially misleading phrases (e.g., “Free Return”, “Promo Code: ABC”); 2) Brand Check: brand list compiled from our search engine’s knowledge graph (Noy et al., 2019; Chai et al., 2021); 3) Domain Check: checking patterns like “xyz.com” against landing page URL.

We add the cross check rules at two stages: (1) We filter training data with cross check rules before training (train x-check); and (2) We filter generated text after the inference (infer x-check). Per Table 1, both train x-check and infer x-check improve quality significantly, with the greatest improvement when both are used together.

For an AdCopy model, we do observe that $\sim 15\%$ of generated ad copies are filtered during the post-inference cross check. This effect is ameliorated by the fact that we use multiple NLG models, allowing them to backfill each other’s coverage. The remaining coverage is backfilled with extraction candidates. Due to this system design, the eventual URL coverage does not suffer from the cross check.

2.3 Controllable Generation at Asset-Level

To model diversity explicitly, we build a controllable NLG model to generate multiple ad assets for the same source sequence. We accomplish this is via *control codes*, categorical variables that represent the desired output property and are prepended to the model inputs during training and testing, Keskar et al. (2019) and Ficler and Gold-

Technique	Overall	Text Quality	Human Like	Factuality	Relevance
Advertiser-written	90.7 \pm 2.1	97.9 \pm 1.0	98.1 \pm 1.0	92.7 \pm 1.9	99.0 \pm 0.7
Baseline: AdCopy w/o check	89.8 \pm 2.2	98.8 \pm 0.8	98.5 \pm 0.9	91.1 \pm 2.1	98.9 \pm 0.7
AdCopy w/ train check	94.7 \pm 1.6	99.6 \pm 0.5	99.0 \pm 0.7	95.6 \pm 1.5	99.6 \pm 0.5
AdCopy w/ infer check	94.4 \pm 1.9	98.8 \pm 0.9	98.5 \pm 1.0	95.6 \pm 1.7	98.8 \pm 0.9
AdCopy w/ train + infer check	96.3 \pm 1.5	100.0	99.4 \pm 0.6	97.0 \pm 1.3	99.7 \pm 0.4

Table 1: A comparison of Ad Copy models (as per Section 2.1) via human evaluation. 95% confidence intervals (CI) are reported. Results that outperform advertiser baseline at $p < 0.05$ level are **bolded**.

Technique	Overall	Text Quality	Human Like	Factuality	Relevance
Advertiser Title Asset	98.2 \pm 0.9	99.9 \pm 0.2	100.0	98.4 \pm 0.9	100.0
Extraction Title Asset	99.0 \pm 0.7	99.4 \pm 0.6	99.6 \pm 0.5	99.6 \pm 0.5	100.0
Guided Title Asset	98.1 \pm 0.6	99.8 \pm 0.2	100.0	98.3 \pm 0.5	99.6 \pm 0.3
Advertiser Desc Asset	98.2 \pm 0.9	99.9 \pm 0.2	99.9 \pm 0.2	98.4 \pm 0.9	100.0
Guided Desc Asset	95.3 \pm 0.9	97.6 \pm 0.7	98.8 \pm 0.5	97.9 \pm 0.6	99.2 \pm 0.4

Table 2: A comparison of Guided Asset generation model against advertiser written ads and extraction-based titles via human evaluation. 95% CI are reported. Results better than advertiser baseline at $p < 0.05$ level are **bolded**.

berg (2017). We refer to it as Guided model, as the generation is guided by the control codes.

We assume each landing page can be advertised along 12 categories for different selling points. Example categories are Product or Service, Advertiser Name or Brand, Location, etc.; they are borrowed from the instructions on the web portal where advertisers create ads. We then use human judges to classify ~ 6500 distinct advertiser-written assets into categories. We finetune BERT-base-uncased (Devlin et al., 2018) for asset category classification task and obtain $\sim 80\%$ prediction accuracy, using a random 80/20 split for train/test sets and learning rate of $5 * 10^{-5}$.

We then inference ad category for each ad asset in the NLG model training set, prepending the resulting category control code as plaintext at the beginning of each NLG source sequence. Thus, we obtain a data set of 6M ad assets (both title and description together) for training the Guided NLG models. Otherwise, our generative modeling decisions align with Section 2.1. During inference, we evaluate the model on all available categories, by prepending each control code to the landing page information.

Human evaluation results for our Guided NLG model are shown in Table 2. The overall title asset quality of the Guided model does not have significant difference to that of advertiser-written assets, with Extraction titles outperforming both. The advertiser-written description assets are better, though the overall good rate of Guided model is

Title Asset	Count	PB \downarrow	SB \downarrow	Dist \uparrow
Advertiser	18.4	13.4	71.0	45.3
Generated	24.4	6.7	41.0	66.6
Generated + DPP	14.2	4.5	25.3	80.5
Guided	13.3	7.8	33.6	74.9
Ensemble	12.1	5.8	31.2	77.0
Guided + DPP	7.8	5.0	18.3	86.9
Ensemble + DPP	7.7	3.6	17.0	88.3

Table 3: Averaged results of the diversity evaluation on English title assets. For PairwiseBLEU (PB) and Self-BLEU (SB) scores, lower is better, for Distinct N-gram (Dist) scores, higher is better. Average count of title assets per URL (Count) is also reported. Differences of over 1 point are **bolded**. Ensemble here is for an ensemble of AdCopy models. Generated assets include the combination of Guided, Ensemble, and Extraction titles.

still well above our quality bar of 90%. The advantage of Guided model in this case is that it is able to explicitly capture different advertising categories for both title and description. Extraction technique cannot produce good ad descriptions in our experience.

3 Serving and Customization System

3.1 Diverse Selection

At this stage, we aim to select a semantically diverse subset of T title and D description assets for each URL to send to online serving components. By selecting a subset of ad texts, we aim to both

reduce the load on the ad serving system, as well as improve diversity of the generated texts. We use CDSSM (Shen et al., 2014) model, trained on web search logs, to map each text asset to a dense vector, such that the ad texts with high degree of semantic similarity will map to representations with higher cosine similarity (i.e., closer in the embedding space) to one another. Then, we sample a diverse subset of points in the CDSSM embedding space with k-DPP maximum a posteriori inference algorithm as per Chen et al. (2018), stopping after we select T titles or D descriptions. Refer to Figure 2 (bottom middle) for an example of removing semantic duplicates in such fashion.

We use PairwiseBLEU (PB) (Shen et al., 2019), SelfBLEU (SB) (Zhu et al., 2018), and Distinct N-gram (DistN) (Xu et al., 2018) scores to evaluate the diversity of the title assets before and after k-DPP diverse sampling. We calculate the average diversity metrics for ~ 2000 EN URLs randomly sampled from a stratified sample of 50 URLs/domain. Since all instances of each metric show similar trends, we follow suit with Tevet and Berant (2020) and average each metric over different N-gram options. Refer to Table 3 for diversity score details.

We find that generated title assets are more diverse than the ones provided by the advertiser in general. In addition, k-DPP helps further increase the asset diversity. We also compare title assets from the Guided model with those from an ensemble of AdCopy models. We find that the Guided model by itself can generate title assets in similar quantity and with similar diversity as the ones produced by several AdCopy models combined, trained as per the NLG baseline method in Section 2.1 on different versions of training data.

3.2 Real-Time Stitching

The diversified ad assets are then ingested into the online serving infrastructure. At query time, we stitch together a customized ad copy, optimizing for the auction win rate¹ (with some level of exploration). From our domain knowledge, the earlier asset positions (e.g., Title 1) influence the ad auction result more than the later ones (e.g., Description 2), as shown in Figure 1. Thus, we perform a greedy sequential selection and consider $T + (T - 1) + (T - 2) + D + (D - 1)$ permutation

¹Auction is the final stage to decide which ads will be displayed. Auction win rate is the probability of an ad winning an auction. Ads with better quality and CTR have a better chance to win the auction.

options. For example, we first select asset for Title 1 position from T title assets, and then select asset for Title 2 position from the remaining $T - 1$ title assets.

We use a logistic regression (LR) model to score each asset position: Title 1, Title 2, Title 3, Description 1, Description 2. We use features from ad auction log like string hash, length, unigrams and bigrams from asset texts. We also cross these with the query text to a total sparse feature dimensionality per position of $\sim 4B$. The LR model learns the probability of winning the auction for a given ad copy. It is continuously trained daily, using $\sim 10B$ data examples from the previous day’s log for training with batch size as 1000 and learning rate as 0.02, and $\sim 300M$ examples from current day’s log for validation.

We include an exploration mechanism to allow newly added assets to be shown to users and to de-bias the model. Due to sequential nature of our stitching process, we model exploration as a sequential contextual bandit (CB) problem. At each asset position, the CB uses the LR score and the gradient sum of LR features as a heuristic for the trial count (Mcmahan et al., 2013) to select an asset using Thompson Sampling strategy (Agrawal and Goyal, 2017). As a result, we sample from a total of $T + (T - 1) + (T - 2) + D + (D - 1)$ Beta distributions to stitch together an ad copy.

4 A/B Testing

DeepGen is deployed globally to serve Dynamic Search Ads (DSA), which accounts for $\sim 4\%$ of all Bing Ads displayed globally. In A/B testing, we split the production user traffic randomly between the treatment experiment that enables the proposed experimental techniques and the control experiment that uses existing production techniques. We use 10% of production traffic for the control experiment. We use the difference in business metrics between two experiments to decide if treatment is effective.

Two key business metrics are Revenue Per Mille (RPM) – revenue per every thousand search result page views (SRPV) and Quick Back Rate (QBR) – the rate of users clicking the back button after clicking on an ad, which is a proxy for user dissatisfaction (lower QBR is better). RPM is driven by Impression Yield (IY, number of ads shown divided by number of search result page views) and Click-Through Rate (CTR, number of clicks divided by

Metric	Exp. 1	Exp. 2	Explanation
Days	5	10	Number of days for the experiment.
Traffic%	5.0	10.0	Percentage of the Bing user traffic allocated for the experiment.
Δ RPM%	+24.87	+10.65	Revenue (USD) from 1000 Search Results Page Views (SRPVs).
Δ IY%	+11.87	+14.43	Average number of ads shown per page.
Δ CTR%	+13.28	-0.19	Proportion of ads clicked from ads shown.
Δ QBR%	+5.27	+1.82	Proportion of ad clicks that resulted in a back-click within 20 sec.

Table 4: A summary of the business metrics from A/B tests performed on DSA ad traffic. Results statistically significant at $p < 0.05$ level are **bolded**.

total number of ads displayed). Usually there is a trade-off between RPM (revenue) and QBR (user satisfaction). DeepGen increases CTR (proportion of ads clicked from ads displayed) and IY (number of ads displayed per page), thus also increasing ad revenue. We do so by generating high-quality ads that are customized to the user. We avoid sacrificing user or advertiser satisfaction by ensuring the ads to be faithful to the landing page.

In Exp. 1, we compare DeepGen (treatment) against the extraction system (control). As shown in Table 4, we observe strong RPM (revenue) gain, driven by both IY and CTR, which means that personalized ad copies generated by DeepGen are more likely to win the auction as well as to be clicked by the user. In this experiment, we record a *13.28% CTR gain*. We acknowledge the increase in QBR (user dissatisfaction), which could be attributed to the still higher factuality of the extraction system, as shown in Table 2.

We use Exp. 2 as an ablation for real-time customization. DeepGen is used in both treatment and control, but we replace real-time stitching with pre-computed stitching in control. For this experiment, we build a separate model to stitch assets into multiple ad copies offline, and only *rank* the pre-stitched ad texts during query time (online). There is significant RPM (revenue) gain, though it is mainly driven by IY but not CTR. This may suggest that online stitching has a higher chance of winning the auction as it covers much larger permutation space than the offline stitching. But for those ad copies that did win an auction, they have similar attractiveness to the user whether stitched online or offline. This experiment shows online stitching to be an integral part of our system.

Thus, DeepGen increases revenue by generating high-quality ads customized to the user while being mindful of user satisfaction by ensuring the ads to be faithful to the landing page.

5 Related Work

The early automated content generation approaches focused on template-based ad text generation (Bartz et al., 2008; Fujita et al., 2010; Thomaidou et al., 2013). These approaches have potential to suffer from ad fatigue (Abrams and Vee, 2007).

More recently, deep Reinforcement Learning (RL) was shown effective for ad text generation (Hughes et al., 2019; Kamigaito et al., 2021; Wang et al., 2021), using a *general* attractiveness model as a reward policy and yielding up to 7.01% observed CTR gain per Kamigaito et al. (2021). CTR is an important metrics, as reflects the relevance of an ad from user’s perspective (Yang and Zhai, 2022).

Product headline generation is a closely related direction of work, where a single headline is generated to advertise a line of related products, based on each product’s advertiser-written title. Kanungo et al. (2021) use BERT-large (Devlin et al., 2018) encoder finetuned for generation with UniLM-like masked attention, as per Dong et al. (2019), optimized using a self-critical RL objective, as per Hughes et al. (2019). Kanungo et al. (2022) further produce SC-COBART by finetuning a BART model, using control codes, as per Keskar et al. (2019), for bucketized CTR and length of a headline, optimized with a mixture of MLE and self-critical RL objectives. SC-COBART improves estimated CTR by 5.82% over their previous work (Kanungo et al., 2021).

In another line of work, product descriptions are generated either with templates (Wang et al., 2017), pointer-generator encoders (Zhang et al., 2019), commonsense knowledge-base guidance (Chan et al., 2020; Zhang et al., 2021), or CVAEs (Shao et al., 2021), yielding up to 13.17% CTR gain in A/B test per Shao et al. (2021).

6 Conclusion

In this work, we present an automated end-to-end search advertisement text generation solution. We employ deep NLG modeling for ad content generation and diverse selection. We leverage real-time LR rankers for content stitching. The generation techniques provide us a rich source of high-quality ad content, which performs strongly against human and extraction baselines. We further apply diverse selection via semantic embedding, which allows us to surpass human content diversity, while ensuring the system’s scalability. Finally, we use real-time ranking to stitch not just attractive, but a truly *customized* ad for each user based on query and search intent. The system combines several NLP approaches to provide a cutting edge solution to automated ad generation and showcases an significant CTR gain over an extraction baseline.

References

- Zoë Abrams and Erik Vee. 2007. [Personalized ad delivery when ads fatigue: An approximation algorithm](#). pages 535–540.
- Shipra Agrawal and Navin Goyal. 2017. [Near-optimal regret bounds for thompson sampling](#). *J. ACM*, 64(5).
- Hangbo Bao, Li Dong, Furu Wei, Wenhui Wang, Nan Yang, Xiaodong Liu, Yu Wang, Songhao Piao, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2020. [Unilmv2: Pseudo-masked language models for unified language model pre-training](#). *CoRR*, abs/2002.12804.
- Kevin Bartz, Cory Barr, and Adil Aijaz. 2008. [Natural language generation for sponsored-search advertisements](#). In *Proceedings of the 9th ACM Conference on Electronic Commerce*, EC ’08, page 1–9, New York, NY, USA. Association for Computing Machinery.
- Junyi Chai, Yujie He, Homa Hashemi, Bing Li, Daraksha Parveen, Ranganath Kondapally, and Wenjin Xu. 2021. [Automatic construction of enterprise knowledge base](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 11–19, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zhangming Chan, Yuchi Zhang, Xiuying Chen, Shen Gao, Zhiqiang Zhang, Dongyan Zhao, and Rui Yan. 2020. [Selection and generation: Learning towards multi-product advertisement post generation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3818–3829, Online. Association for Computational Linguistics.
- Laming Chen, Guoxin Zhang, and Hanning Zhou. 2018. [Fast greedy map inference for determinantal point process to improve recommendation diversity](#). In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS’18*, page 5627–5638, Red Hook, NY, USA. Curran Associates Inc.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.
- Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. [Unified language model pre-training for natural language understanding and generation](#). *CoRR*, abs/1905.03197.
- Jessica Fidler and Yoav Goldberg. 2017. [Controlling linguistic style aspects in neural language generation](#). *CoRR*, abs/1707.02633.
- Katja Filippova. 2020. [Controlled hallucinations: Learning to generate faithfully from noisy data](#). *CoRR*, abs/2010.05873.
- Atsushi Fujita, Katsuhiko Ikushima, Satoshi Sato, Ryo Kamite, Ko Ishiyama, and Osamu Tamachi. 2010. [Automatic generation of listing ads by reusing promotional texts](#). In *Proceedings of the 12th International Conference on Electronic Commerce: Roadmap for the Future of Electronic Business*, ICEC ’10, page 179–188, New York, NY, USA. Association for Computing Machinery.
- J. Weston Hughes, Keng-hao Chang, and Ruofei Zhang. 2019. [Generating better search engine text advertisements with deep reinforcement learning](#). In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD ’19*, page 2269–2277, New York, NY, USA. Association for Computing Machinery.
- Hidetaka Kamigaito, Peinan Zhang, Hiroya Takamura, and Manabu Okumura. 2021. [An empirical study of generating texts for search engine advertising](#). In *NAACL*.
- Yashal Shakti Kanungo, Gyanendra Das, Pooja A, and Sumit Negi. 2022. [Cobart: Controlled, optimized, bidirectional and auto-regressive transformer for ad headline generation](#). In *KDD 2022*.
- Yashal Shakti Kanungo, Sumit Negi, and Aruna Rajan. 2021. [Ad headline generation using self-critical masked language model](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Papers*, pages 263–271, Online. Association for Computational Linguistics.
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. [CTRL: A conditional transformer language model for controllable generation](#). *CoRR*, abs/1909.05858.

- Chin-Yew Lin. 2004. [Rouge: a package for automatic evaluation of summaries](#). In *Workshop on Text Summarization Branches Out, Post-Conference Workshop of ACL 2004, Barcelona, Spain*, pages 74–81.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020a. [On faithfulness and factuality in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan T. McDonald. 2020b. [On faithfulness and factuality in abstractive summarization](#). *CoRR*, abs/2005.00661.
- H. Brendan Mcmahon, Gary Holt, D. Sculley, Michael Young, Dietmar Ebner, Julian Grady, Lan Nie, Todd Phillips, Eugene Davydov, Daniel Golovin, Sharat Chikkerur, Dan Liu, Martin Wattenberg, Arnar Mar Hrafnkelsson, Tom Boulos, and Jeremy Kubica. 2013. [Ad click prediction: A view from the trenches](#). *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Part F1288:1222–1230.
- Feng Nan, Ramesh Nallapati, Zhiguo Wang, Cicero Nogueira dos Santos, Henghui Zhu, Dejiao Zhang, Kathleen McKeown, and Bing Xiang. 2021. [Entity-level factual consistency of abstractive text summarization](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2727–2733, Online. Association for Computational Linguistics.
- Natasha Noy, Yuqing Gao, Anshu Jain, Anant Narayanan, Alan Patterson, and Jamie Taylor. 2019. [Industry-scale knowledge graphs: Lessons and challenges](#). *Communications of the ACM*, 62(8):36–43.
- Huajie Shao, Jun Wang, Haohong Lin, Xuezhou Zhang, Aston Zhang, Heng Ji, and Tarek F. Abdelzaher. 2021. [Controllable and diverse text generation in e-commerce](#). *CoRR*, abs/2102.11497.
- Tianxiao Shen, Myle Ott, Michael Auli, and Marc’Aurelio Ranzato. 2019. [Mixture models for diverse machine translation: Tricks of the trade](#). *CoRR*, abs/1902.07816.
- Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. [A latent semantic model with convolutional-pooling structure for information retrieval](#). *CIKM 2014 - Proceedings of the 2014 ACM International Conference on Information and Knowledge Management*, pages 101–110.
- Statista. 2022. [Search advertising - worldwide | statista market forecast](#).
- Guy Tevet and Jonathan Berant. 2020. [Evaluating the evaluation of diversity in natural language generation](#). *CoRR*, abs/2004.02990.
- Stamatina Thomaidou, Ismini Lourentzou, Panagiotis Katsivelis-Perakis, and Michalis Vazirgiannis. 2013. [Automated snippet generation for online advertising](#). In *Proceedings of the 22nd ACM International Conference on Information & Knowledge Management, CIKM ’13*, page 1841–1844, New York, NY, USA. Association for Computing Machinery.
- Jinpeng Wang, Yutai Hou, Jing Liu, Yunbo Cao, and Chin-Yew Lin. 2017. [A statistical framework for product description generation](#). In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 187–192, Taipei, Taiwan. Asian Federation of Natural Language Processing.
- Xiting Wang, Xinwei Gu, Jie Cao, Zihua Zhao, Yulan Yan, Bhuvan Middha, and Xing Xie. 2021. [Reinforcing pretrained models for generating attractive text advertisements](#). In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (ACM SIGKDD)*. Applied Data Science Track.
- Lee Xiong, Chuan Hu, Chenyan Xiong, Daniel Campos, and Arnold Overwijk. 2019. [Open domain web keyphrase extraction beyond language modeling](#). In *Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5174–5183. ACL, ACL.
- Jingjing Xu, Xuancheng Ren, Junyang Lin, and Xu Sun. 2018. [Diversity-promoting GAN: A cross-entropy based generative adversarial network for diversified text generation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3940–3949, Brussels, Belgium. Association for Computational Linguistics.
- Yu Yan, Fei Hu, Jiusheng Chen, Nikhil Bhendawade, Ting Ye, Yeyun Gong, Nan Duan, Desheng Cui, Bingyu Chi, and Ruofei Zhang. 2021. [FastSeq: Make sequence generation faster](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 218–226, Online. Association for Computational Linguistics.
- Yanwu Yang and Panyu Zhai. 2022. [Click-through rate prediction in online advertising: A literature review](#). *Information Processing & Management*, 59(2):102853.
- Chao Zhang, Jingbo Zhou, Xiaoling Zang, Qing Xu, Liang Yin, Xiang He, Lin Liu, Haoyi Xiong, and Dejing Dou. 2021. [CHASE: Commonsense-Enriched Advertising on Search Engine with Explicit Knowledge](#), page 4352–4361. Association for Computing Machinery, New York, NY, USA.
- Tao Zhang, Jin Zhang, Chengfu Huo, and Weijun Ren. 2019. [Automatic generation of pattern-controlled product description in e-commerce](#). In *The World Wide Web Conference, WWW ’19*, page 2355–2365,

New York, NY, USA. Association for Computing Machinery.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. [Texygen: A benchmarking platform for text generation models](#). *CoRR*, abs/1802.01886.

A Ad Text Quality Judgement UI

Please evaluate if the below generated Bing Ad is valid?

Generated Ad

[Shop Our Winter Candles | Buy Direct from the Manufacturer](#)
Ad [goosecreekcandle.com](#)
High Quality, Long-lasting Fragrance for Any Occasion.

QA Guidelines Checklist:

1. Rate the language and the grammar of the ad copy.
 Good Fair Bad Embrassing Not Scorable *required

Comments:

2. Is the ad copy human like?
 Yes No *required

3. Is the information in the ad accurate?
(There should be no claims/offerings/deals not offered on the advertiser's LP and domain)
 Yes No *required

4. Is the ad's focus relevant to the advertiser's business?
 Yes No *required

Comments:

Figure 4: User interface for human ad quality evaluation.