

Multi-Modal Generative Adversarial Network for Short Product Title Generation in Mobile E-Commerce

Jian-Guo Zhang,¹ Pengcheng Zou,² Zhao Li,² Yao Wan,³
Xiuming Pan,² Yu Gong,² Philip S. Yu¹

¹ University of Illinois at Chicago, USA; ² Alibaba Group; ³ Zhejiang University, China
{jzhan51, psyu}@uic.edu, wanyao@zju.edu.cn
{xuanwei.zpc, lizhao.lz}@alibaba-inc.com
{xuming.panxm, gongyu.gy}@alibaba-inc.com

Abstract

Nowadays, more and more customers browse and purchase products in favor of using mobile E-Commerce Apps such as Taobao and Amazon. Since merchants are usually inclined to describe redundant and over-informative product titles to attract attentions from customers, it is important to concisely display short product titles on limited screen of mobile phones. To address this discrepancy, previous studies mainly consider textual information of long product titles and lacks of human-like view during training and evaluation process. In this paper, we propose a Multi-Modal Generative Adversarial Network (MM-GAN) for short product title generation in E-Commerce, which innovatively incorporates image information and attribute tags from product, as well as textual information from original long titles. MM-GAN poses short title generation as a reinforcement learning process, where the generated titles are evaluated by the discriminator in a human-like view. Extensive experiments on a large-scale E-Commerce dataset demonstrate that our algorithm outperforms other state-of-the-art methods. Moreover, we deploy our model into a real-world online E-Commerce environment and effectively boost the performance of click through rate and click conversion rate by 1.66% and 1.87%, respectively.

1 Introduction

E-commerce companies such as TaoBao and Amazon put many efforts to improve the user experience of their mobile Apps. For the sake of improving retrieval results by search engines, merchants usually write lengthy, over-informative, and sometimes incorrect titles, e.g., the original product title in Fig. 1 contains more than 20 Chinese words, which may be suitable for PCs. However, these titles are cut down and no more than 10 words can be displayed on a mobile phone with limited screen size varying from 4 to 5.8 inches. Hence, to properly

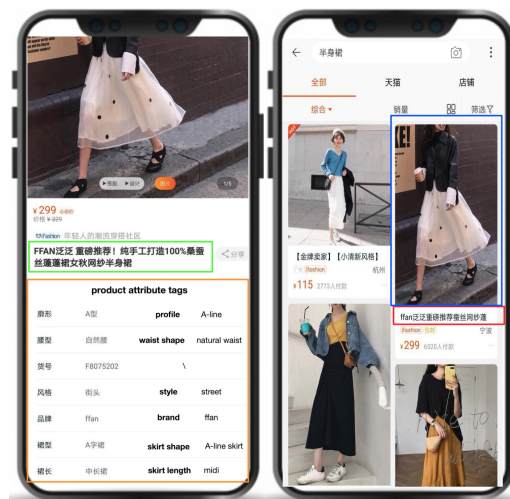


Figure 1: A product with original long titles (green box), cutoff short titles (red box), the main image (blue box), and attribute tags (yellow box).

display products in mobile screen, it is important to produce succinct short titles to preserve important information of original long titles and accurate descriptions of products.

This problem is related to text summarization, which can be categorized into two classes: extractive (Cao et al., 2016; Miao and Blunsom, 2016; Nallapati et al., 2017) and abstractive (Chen et al., 2016; Chopra et al., 2016; See et al., 2017; Wan et al., 2018) methods. The extractive methods select important words from original titles, while the abstractive methods generate titles by extracting words from original titles or generating new words from data corpus. They usually approximate such goals by predicting the next word given previous predicted words using maximum likelihood estimation (MLE) objective. Despite their successes to a large extent, they suffer from the issue of *exposure bias* (Ranzato et al., 2016). It may cause the models to behave in undesired ways, e.g., generating repetitive or truncated outputs. In addition, predicting next word based on previously generated words will make the learned model lack of human-like

holistic view of the whole generated short product titles.

More recent state-of-the-art methods (Gong et al., 2018; Wang et al., 2018) treat short product titles generation as a sentence compression task following attention-based extractive mechanism. They extract key characteristics mainly from original long product titles. However, in real E-Commerce scenario, product titles are usually redundant and over-informative, and sometimes even inaccurate, e.g., long titles of a cloth may include both “å»å|çé (hip-pop|wild)” and “æè°|æ·å (artsy|delicate)” simultaneously. It is a tough task to generate succinct and accurate short titles merely relying on the original titles. Therefore, it is insufficient to regard short title generation as traditional text summarization problem in which original text has already contained complete information.

In this paper, we propose a novel **Multi-Modal Generative Adversarial Network**, named **MM-GAN**, to better generate short product titles. It contains a generator and a discriminator. The generator generates a short product title based on original long titles, with additional information from the corresponding visual image and attribute tags. On the other hand, the discriminator tries to distinguish whether the generated short titles are human-produced or machine-produced in a human-like view. The task is treated as a reinforcement learning problem, in which the quality of a machine-generated short product title depends on its ability to fool the discriminator into believing it is generated by human, and output of the discriminator is a reward for the generator to improve generated quality. The main contributions of this paper can be summarized as follows:

- In this paper, we focus on a fundamental problem existing in the E-Commerce industry, i.e., generating short product titles for mobile E-Commerce Apps. We formulate the problem as a reinforcement learning task;
- We design a multi-modal generative adversarial network to consider multiple modalities of inputs for better short product titles generation in E-commerce;
- To verify the effectiveness of our proposed model, we deploy it into a mobile E-commerce App. Extensive experiments on a large-scale real-world dataset with A/B testing show that our proposed model outperforms state-of-the-art methods.

2 Related Work

Our work is related to text summarization tasks and generative adversarial networks (GANs).

Text Summarization. In terms of text summarization, it mainly includes two categories of approaches: extractive and abstractive methods. Extractive methods produce a text summary by extracting and concatenating several words from original sentence. Whereas abstractive methods generate a text summary based on the original sentence, which usually generate more readable and coherent summaries. Traditional extractive methods such as graphic models (Mihalcea, 2005) and optimization-based models (Woodsend and Lapata, 2010) usually rely on human-engineered features. Recent RNN-based methods (Chopra et al., 2016; Gong et al., 2018; Nallapati et al., 2017; Wang et al., 2018) have become popular in text summarization tasks. Among them, (Gong et al., 2018; Wang et al., 2018) design attention-based neural networks for short product titles generation in E-commerce. (Gong et al., 2018) considers rich semantic features of long product titles. (Wang et al., 2018) designs a multi-task model and uses user searching log data as additional task to facilitate key words extraction from original long titles. However, they mainly consider information from textual long product titles, which sometimes are not enough to select important words and filter out over-informative and irrelevant words from long product titles. In addition, these methods mainly apply MLE objective and predict next word based on previous words. As a consequence, these models usually suffer from *exposure bias* and lack of human-like view.

Generative Adversarial Networks (GANs). GAN is firstly proposed in (Goodfellow et al., 2014), which is designed for generating real-valued, continuous data, and has gained great success in computer vision tasks (Dai et al., 2017; Isola et al., 2017; Ledig et al., 2017). However, applying GANs to discrete scenarios like natural language has encountered many difficulties, since the discrete elements break the differentiability of GANs and prohibit gradients backpropagating from the discriminator to generator. To mitigate these above mentioned issues, SeqGan (Yu et al., 2017) models sequence generation procedure as a sequential decision making process. It applies a policy gradient method to train the generator and discriminator, and shows improvements on multiple generation

task such as poem generation and music generation. MaskGan (Fedus et al., 2018) designs a actor-critic reinforcement learning based GAN to improve qualities of text generation through filling in missing texts conditioned on the surrounding context. There are also some other RL based GANs for text generation such as LeakGan (Guo et al., 2017), RankGan (Lin et al., 2017), SentiGan (Wang and Wan, 2018), etc. All above methods are designed for unsupervised text generation tasks. (Li et al., 2017) designs an adversarial learning method for neural dialogue generation. They train a seq2seq model to produce responses and use a discriminator to distinguish whether the responses are human-generated and machine-generated, and showing promising results. It should be noticed that our work differs from other similar tasks such as image captioning (Dai et al., 2017) and visual question answering (Antol et al., 2015). The image captioning can be seen as generating caption from a single modality of input, while the visual question answering mainly focuses on aligning the input image and question to generate a correct answer. In our task, we put more attention on learning more information from the multi-modal input sources (i.e., long product titles, product image and attribute tags) to generate a short product title.

3 Multi-Modal Generative Adversarial Network

In this section, we describe in details the proposed MM-GAN. The problem can be formulated as follows: given an original long product title $L = \{l_1, l_2, \dots, l_K\}$ consisted of K Chinese or English words, a single word can be represented in a form like "skirt" in English or "半身裙" in Chinese. With an additional image I and attribute tags $A = \{a_1, a_2, \dots, a_M\}$, the model targets at generating a human-like short product title $S = \{s_1, s_2, \dots, s_N\}$, where M and N are the number of words in A and S , respectively.

3.1 Multi-Modal Generator

The multi-modal generative model defines a policy of generating short product titles S given original long titles L , with additional information from product image I and attribute tags A . Fig. 2 illustrates the architecture of our proposed multi-modal generator which follows the seq2seq (Sutskever et al., 2014) framework.

Multi-Modal Encoder. As we mentioned be-

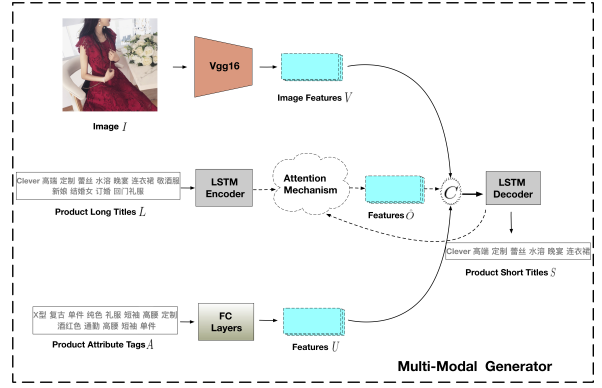


Figure 2: Overall Framework of MM-GAN.

fore, our model tries to incorporate multiple modalities of a product (i.e., product image, attribute tags and long title). To learn the multi-modal embedding of a product, we first adopt a pre-trained VGG16 (Simonyan and Zisserman, 2015) as the CNN architecture to extract features $V = [v_1, v_2, \dots, v_Z]$ of an image I from the condensed fully connected layers, where Z is the number of latent features. In order to get more descriptive features, we fine-tune the last 3 layers of VGG16 based on a supervised classification task given classes of products images. Second, we encode the attribute tags A into a fixed-length feature vector $U = [u_1, u_2, \dots, u_{M'}]$, and $U = f_1(A)$, where f_1 denotes fully connected layers, M' is the output size of f_1 . Third, we apply a recurrent neural network to encode the original long titles L as $O = [o_1, o_2, \dots, o_K]$, where $o_t = f_2(o_{t-1}, l_t)$. Here f_2 represents a non-linear function, and in this paper the LSTM unit (Hochreiter and Schmidhuber, 1997) is adopted.

Decoder. The hidden state h_t for the t -th target word s_t in short product titles S can be calculated as $h_t = f_2(h_{t-1}, s_{t-1}, \hat{o}_t)$. Here we adopt an attention mechanism (Bahdanau et al., 2015) to capture important words from original long titles L . The context vector \hat{o}_t is a weighted sum of hidden states O , which is represented as:

$$\hat{o}_t = \sum_{k=1}^K \alpha_{t,k} o_k, \quad (1)$$

where $\alpha_{t,k}$ is the contribution of an input word l_k to the t -th target word using an alignment model g (Bahdanau et al., 2015):

$$\alpha_{t,k} = \frac{\exp(g(h_{t-1}, o_k))}{\sum_{k'=1}^K \exp(g(h_{t-1}, o_{k'}))}. \quad (2)$$

After obtaining all features U, V, \hat{O} from A, I and L , respectively, we then concatenate them into the final feature vector:

$$C = \tanh(W[\hat{O}; V; U]), \quad (3)$$

where W are learnable weights and $[\cdot; \cdot]$ denotes the concatenation operator.

Finally, C is fed into the LSTM based decoder to predict the probability of generating each target word for short product titles S . As the sequence generation problem can be viewed as a sequence decision making process (Bachman and Precup, 2015), we denote the whole generation process as a policy $\pi(S|C)$.

3.2 Discriminator

The discriminator model D is a binary classifier which takes an input of a generated short product titles S and distinguishes whether it is human-generated or machine-generated. The short product titles are encoded into a vector representation through a two-layer LSTM model, and then fed into a sigmoid function, which returns the probability of the input short product titles being generated by human:

$$R_\phi(S) = \text{sigmoid}(W_d [LSTM(S)] + b_d), \quad (4)$$

where ϕ are learnable parameters for D , W_d is a weight matrix and b_d is a bias vector.

3.3 End-to-End Training

The multi-modal generator G tries to generate a sequence of tokens S under a policy π and fool the discriminator D via maximizing the reward signal received from D . The objective of G can be formulated as follows:

$$J(\theta) = \mathbb{E}_{S \sim \pi_\theta(S|C)} [R_\phi(S)], \quad (5)$$

where θ are learnable parameters for G .

Conventionally, GANs are designed for generating continuous data and thus G is differentiable with continuous parameters guided by the objective function from D (Yu et al., 2017). Unfortunately, it has difficulty in updating parameters through back-propagation when dealing with discrete data in text generation. To solve the problem, we adopt the REINFORCE algorithm (Williams, 1992). Specifically, once the generator reaches the end of a sequence (i.e., $S = S_{1:T}$), it receives a reward $R_\phi(S)$ from D based on the probability of being real.

In text generation, D will provide a reward to G only when the whole sequence has been generated, and no intermediate reward is obtained before the final token of S is generated. This may cause the discriminator to assign a low reward to all tokens in the sequence though some tokens are proper results. To mitigate the issue, we evaluate the reward function by aggregating the N' -time Monte-Carlo simulations (Yu et al., 2017) at each decoding step:

$$R'_\phi(S_{1:t-1}, a' = s_t) \approx \begin{cases} \frac{1}{N'} \sum_{n=1}^{N'} R(S_{1:t}, S_{t+1:N}^{(n)}), & t < N \\ R(S_{1:t-1}, s_t), & t = N, \end{cases} \quad (6)$$

where $\{S_{t+1:N}^{(1)}, \dots, S_{t+1:N}^{(N')}\}$ is the set of generated short titles, which are sampled from the $t+1$ -th decoding step based current state and action. Now we can compute the gradient of the objective function for the generator G :

$$\nabla_\theta J(\theta) \approx \mathbb{E}_{S \sim \pi_\theta(S|C)} \left[\sum_{t=1}^N \nabla_\theta \log(\pi_\theta(s_t|C)) R'_\phi(S_{1:t-1}, s_t) \right], \quad (7)$$

where ∇_θ is the partial differentiable operator for θ in G , and the reward R'_ϕ is fixed during updating of generator.

The objective function for the discriminator D can be formulated as:

$$\mathbb{E}_{S \sim \mathcal{P}_\theta(S|C)} [\log R'_\phi(S|C)] - \mathbb{E}_{S \sim \pi_\theta(S|C)} [\log R'_\phi(S|C)], \quad (8)$$

where $S \sim \mathcal{P}_\theta$ and $S \sim \pi_\theta$ denote that S is from human-written sentences and synthesized sentences, respectively.

On training stage, we first pre-train G and D for several steps. Due to the large size of searching space of possible sequences, it is also very important to feed human-generated short product titles to the generator for model updates. Specifically, we follow the Teacher Forcing mechanism (Lamb et al., 2016). In each training step, we first update the generator using machine-generated data with rewards gained from the discriminator, then sample some data from human-generated short product titles and assign a reward of 1 to them, where the generator uses this reward to update parameters. Alg. 1 summarizes the training procedure, where D-steps and G-steps are both set to 1 in this paper.

Algorithm 1: Multi-Modal Generative Adversarial Network

Input: Long product titles L , short titles S , images I , attribute tags A . Multi-modal generator G , discriminator D .

- 1 Fine-tune last 3 layers of pretrained VGG16 network to get image features V based on a classification task
- 2 Pretrain G on human-generated data using MLE
- 3 Pretrain D on human-generated data and machine-generated data
- 4 **for** number of training iterations **do**
- 5 **for** $i \leftarrow 1$ to D -steps **do**
- 6 Sample S from human-generated data
- 7 Sample $\hat{S} \sim \pi_\theta(\cdot|C)$
- 8 Update D on both S and \hat{S}
- 9 **for** $i \leftarrow 1$ to G -steps **do**
- 10 Sample (L, V, A, S) from human-generated data
- 11 Sample $(L, V, A, \hat{S}) \sim \pi_\theta(\cdot|C)$ based on MC search
- 12 Compute reward R'_ϕ for (L, V, A, \hat{S}) using D
- 13 Update G using R'_ϕ based on Eq. (7)
- 14 Teacher-Forcing: Update G on (L, V, A, S) using MLE

Data set size	2,403,691
Avg. length of long Titles	13.7
Avg. length of Short Titles	4.5
Avg. length of Attributes Tags	18.3
Avg. number of Image	1

Table 1: Statistics of the crawled dataset. Here all the lengths are counted by Chinese or English words.

4 Experiments

4.1 Experimental Setup

Dataset. The dataset used in our experiment is crawled from a module named 有好货 (Youhaohuo) of the well-known 淘宝 (TAOBAO) platform in China. Every product in the dataset includes a long product title and a short product title written by professional writers, along with product several high quality visual images and attributes tags, here for each product we use its main image. This Youhaohuo module includes more than 100 categories of products, we crawled top 7 categories of them in the module, and exclude the products with original long titles shorter than 10 Chinese characters. We further tokenize the original long titles and short titles into Chinese or English words, e.g. “skirt” is a word in English and 半身裙 is a word in Chinese. Table 1 shows some details of the dataset. We randomly select 1.6M samples for training, 0.2M samples for validation, and test our proposed model on 5000 samples.

Baselines. We compare our proposed model with the following four baselines: (a) Pointer Network (**Ptr-Net**) (See et al., 2017) which is a seq2seq based framework with pointer-generator network copying words from the source text via *pointing*. (b) Feature-Enriched-Net (**FE-Net**) (Gong et al., 2018) which is a deep and wide model based on attentive RNN to generate the textual long product titles. (c) Agreement-based MTL (**Agree-MTL**) (Wang et al., 2018) which is a multi-task learning approach to improve product title compression with user searching log data. (d) Generative Adversarial Network (**GAN**) (Li et al., 2017) which is a generative adversarial method for text generation with only one modality of input.

Implementation Details. We first pre-train the multi-modal generative model given human-generated data via maximum likelihood estimation (MLE), and we transfer the pretrained model weights for the multi-modal encoder and decoder modules. Then we pre-train the discriminator using human-generated data and machine-generated data. To get training samples for the discriminator, we sample half of data from multi-modal generator and another half from human-generated data. After that, we perform normal training process based on pre-trained MM-GAN.

Specifically, we create a vocabulary of 35K words for long product titles and short titles, and another vocabulary of 35k for attribute tags in training data with size of 1.6M. Words appear less than 8 times in the training set are replaced as $\langle \text{UNK} \rangle$. We implement a two-layer LSTM with 100 hidden states to encoder attribute tags, and all other LSTMs in our model are two layers with 512 hidden states. The last 3 layers of the pre-trained VGG16 network are fine tuned based on the products visual images with 7 classes. The Adam optimizer (Kingma and Ba, 2015) is initialized with a learning rate 10^{-3} . The multi-modal generator and discriminator are pre-trained for 10000 steps, the normal training steps are set to 13000, the batch size is set to 512 for the discriminator and 256 for the generator, the MC search time is set to 7.

4.2 Automatic Evaluation

To evaluate the quality of generated product short titles, we follow (Wang et al., 2018; Gong et al., 2018) and use standard recall-oriented ROUGE metric (Lin, 2004), which measures the generated quality by counting the overlap of N-grams be-

Models	Ptr-Net	FE-Net	Agree-MTL	GAN	MM-GAN
ROUGE-1	59.21	61.07	66.19	60.67	69.53
ROUGE-2	42.01	44.16	49.39	46.46	52.38
ROUGE-L	57.12	58.00	64.04	60.27	65.80

Table 2: ROUGE performance of different models on the test set.


Data		Methods	Results
Product Long Titles	Artka 阿卡 夏新 花边 镂空 荷叶边 抽绳 民族 狂野 复古 衬衫 S110061Q (Artka Artka summer lace hollow-out flounce drawstring nation wild retro shirt S110061Q)	FE-Net	阿卡 花边 镂空 荷叶边 衬衫 (Artka lace hollow-out flounce shirt)
Image 	Attributes Tags 修身 常规款 圆领 Artka 米白 长袖 套头 复古 通勤 纯色 夏季 喇叭袖 棉 (slim common round-neck Artka off-white long-sleeve pullover retro commuting plain summer flare-sleeve cotton)	Agree-MTL	Artka 阿卡 夏新 花边 镂空 荷叶边 衬衫 (Artka Artka summer lace hollow-out flounce shirt)
		GAN	Artka 荷叶边 抽绳 衬衫 (Artka lace flounce drawstring shirt)
		MM-GAN	Artka 花边 荷叶边 镂空 复古 衬衫 (Artka lace flounce hollow-out retro shirt)

Table 3: The comparison of generated short titles among different methods.

tween the machine-generated and human-generated titles. Here we consider ROUGE-1 (1-gram), ROUGE-2 (bi-grams), ROUGE-L (longest common subsequence). Experimental results on the test set are shown in Table 2. From this table, we note that our proposed MM-GAN achieves best performance on three metrics. Furthermore, when comparing MM-GAN with GAN, we can see that our proposed MM-GAN achieves an improvement of 8.86%, 5.92%, 5.53%, in terms of ROUGE-1, ROUGE-2, ROUGE-3, respectively. This verifies that additional information such as image and attribute tags from product can absolutely facilitate our model to generate better short titles. In addition, compared with the best model Agree-MTL, MM-GAN improves ROUGE-1, ROUGE-2, ROUGE-L by 3.34%, 2.99%, 1.76%, respectively. We attribute the outperformance of MM-GAN to two facts: (a) it incorporates multiple sources, containing more information than other single-source based models. (b) it applies a discriminator to distinguish whether a product short titles are human-generated or machine-generated, which makes the model evaluate the generated sequence in a human-like view, and naturally avoid exposure bias in other methods.

4.3 Online A/B Testing

In order to further verify the effectiveness of MM-GAN, we test our method in the real-world online environment of the TaoBao App.

We perform A/B testing in seven categories of products in the E-commerce App, i.e., “连衣裙| (one-piece)”, “男士T恤| (Man T-shirt)”, “衬衫| (shirt)”, “休闲裤| (Casual pants)”, “女士T恤| (Woman T-shirt)”, “半身裙| (skirt)”, “毛针

织衫| (Sweaters)”. During online A/B testing, users (3×10^6 users per day) are split equally into two groups and are directed into a baseline bucket and an experimental bucket. For users in the baseline bucket, product short titles are generated by the default online system, following an ILP based method (Clarke and Lapata, 2008). While for users in the experimental bucket, product short titles are generated by MM-GAN. All conditions in the two buckets are identical except for short titles generation methods. We apply Click Through Rate (CTR) and Click Conversion Rate (CVR) to measure the performance. $CTR = \frac{\#click_of_product}{\#pv_of_product}$, and $CVR = \frac{\#trade_of_product}{\#click_of_product}$, where $\#click_of_product$ indicates clicking times of a product, $\#pv_of_product$ is the number of page views of the product and $\#trade_of_product$ is the number of purchases of the product.

We deploy A/B testing for 7 days and calculate overall CTR for all products in the baseline bucket and experimental bucket. MM-GAN improves CTR by 1.66% and CVR by 1.87% in the experimental bucket over the default online system in the baseline bucket. It verifies the effectiveness of our proposed method. Moreover, through online A/B testing, we find that with more readable, informative product short titles, users are more likely to click, view and purchase corresponding products, which indicates good short product titles can help improving the product sales on E-commerce Apps.

4.4 Qualitative Analysis

Table 3 shows a sample of product short titles generated by MM-GAN and baselines.

From this table, we can note that (a) product

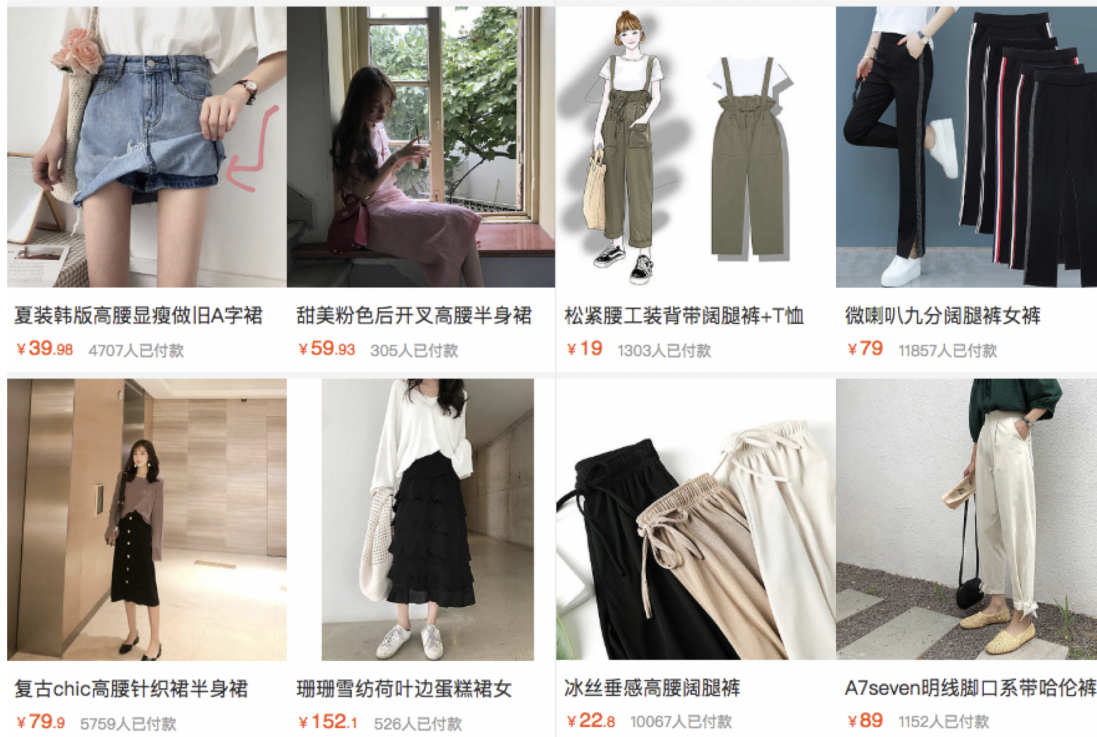


Figure 3: Some samples generated by MM-GAN.

short titles generated by our model are more fluent, informative than baselines, and core product words (e.g., “Artka| (阿卡)”, “复古| (retro)”, “衬衫| (shirt)”) can be recognized. (b) There are over-informative words (e.g., “阿卡| (Artka)”, “S110061Q”) and irrelevant words (e.g., “狂野| (wild)”) in product long titles. Over-informative words may disturb model’s generation process, irrelevant words may give incorrect information to the model. These situations could happen in real E-commerce environment. FE-Net misses the English brand name “Artka” and gives its Chinese name ‘阿卡’ instead. Agree-MTL using user searching log data performs better than GAN. However, Agree-MTL still generates the over-informative word ‘阿卡’. MM-GAN outperforms all baselines, information in additional attribute tags such as “复古| (retro)”, “Artka”), and other information from the product main image are together considered by the model and help the model select core words and filter out irrelevant words in generated product short titles. It shows that MM-GAN using different types of inputs can help generate better product short titles. To leave a deeper impression on the performance of our proposed model, we also put more online samples generated by the MM-GAN in a module of the TaoBao App, as shown in Fig. 3. From all generated samples we also find few bad

cases which are not shown online (e.g., repetitive words in the generated short titles, wrong generated words which are not related to the product at all), leaving a great space for further improvement.

5 Conclusion

In this paper, we propose a multi-modal generative adversarial network for short product title generation in E-commerce. Different from conventional methods which only consider textual information from long product titles, we design a multi-modal generative model to incorporate additional information from product image and attribute tags. Extensive experiments on a large real-world E-commerce dataset verify the effectiveness of our proposed model when comparing with several state-of-the-art baselines. Moreover, the online deployment in a real environment of an E-commerce app shows that our method can improve the click through rate and click conversion rate.

6 Acknowledgement

This work is supported in part by NSF through grants IIS-1526499, IIS-1763325, and CNS-1626432, and NSFC 61672313. We thank the anonymous reviewers for their valuable comments.

References

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *ICCV*.
- Philip Bachman and Doina Precup. 2015. Data generation as sequential decision making. In *NIPS*.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *ICLR*.
- Ziqiang Cao, Wenjie Li, Sujian Li, Furu Wei, and Yanran Li. 2016. Attsum: Joint learning of focusing and summarization with neural attention. In *COLING*.
- Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, and Hui Jiang. 2016. Distraction-based neural networks for document summarization. In *IJCAI*.
- Sumit Chopra, Michael Auli, and Alexander M Rush. 2016. Abstractive sentence summarization with attentive recurrent neural networks. In *NAACL-HLT*.
- James Clarke and Mirella Lapata. 2008. Global inference for sentence compression: An integer linear programming approach. *Journal of Artificial Intelligence Research*.
- Bo Dai, Sanja Fidler, Raquel Urtasun, and Dahua Lin. 2017. Towards diverse and natural image descriptions via a conditional gan. In *ICCV*.
- William Fedus, Ian Goodfellow, and Andrew M Dai. 2018. Maskgan: Better text generation via filling in the `_`. In *ICLR*.
- Yu Gong, Xusheng Luo, Kenny Q Zhu, Shichen Liu, and Wenwu Ou. 2018. Automatic generation of chinese short product titles for mobile display. In *IAAI*.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *NIPS*.
- Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, and Jun Wang. 2017. Long text generation via adversarial training with leaked information. In *AAAI*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-image translation with conditional adversarial networks. In *ICCV*.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.
- Alex M Lamb, Anirudh Goyal ALIAS PARTH GOYAL, Ying Zhang, Saizheng Zhang, Aaron C Courville, and Yoshua Bengio. 2016. Professor forcing: A new algorithm for training recurrent networks. In *NIPS*.
- Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew P Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. 2017. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*.
- Jiwei Li, Will Monroe, Tianlin Shi, Sébastien Jean, Alan Ritter, and Dan Jurafsky. 2017. Adversarial learning for neural dialogue generation. In *EMNLP*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. *Text Summarization Branches Out*.
- Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, and Ming-Ting Sun. 2017. Adversarial ranking for language generation. In *NIPS*.
- Yishu Miao and Phil Blunsom. 2016. Language as a latent variable: Discrete generative models for sentence compression. In *EMNLP*.
- Rada Mihalcea. 2005. Language independent extractive summarization. In *ACL*.
- Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. Summarunner: A recurrent neural network based sequence model for extractive summarization of documents. In *AAAI*.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2016. Sequence level training with recurrent neural networks. In *ICLR*.
- Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *ACL*.
- Karen Simonyan and Andrew Zisserman. 2015. Very deep convolutional networks for large-scale image recognition. In *ICLR*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *NIPS*.
- Yao Wan, Zhou Zhao, Min Yang, Guandong Xu, Haochao Ying, Jian Wu, and Philip S Yu. 2018. Improving automatic source code summarization via deep reinforcement learning. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, pages 397–407. ACM.
- Jingang Wang, Junfeng Tian, Long Qiu, Sheng Li, Jun Lang, Luo Si, and Man Lan. 2018. A multi-task learning approach for improving product title compression with user search log data. In *AAAI*.
- Ke Wang and Xiaojun Wan. 2018. Sentigan: Generating sentimental texts via mixture adversarial networks. In *IJCAI*.
- Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*.

Kristian Woodsend and Mirella Lapata. 2010. Automatic generation of story highlights. In *ACL*.

Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In *AAAI*.