# V2P: Vision-to-Prompt based Multi-Modal Product Summary Generation

Xuemeng Song†, Liqiang Jing†, Dengtian Lin†,
Zhongzhou Zhao§, Haiqing Chen§, Liqiang Nie†

†Shandong University, Shandong, China,
§Alibaba Group, Hangzhou, China

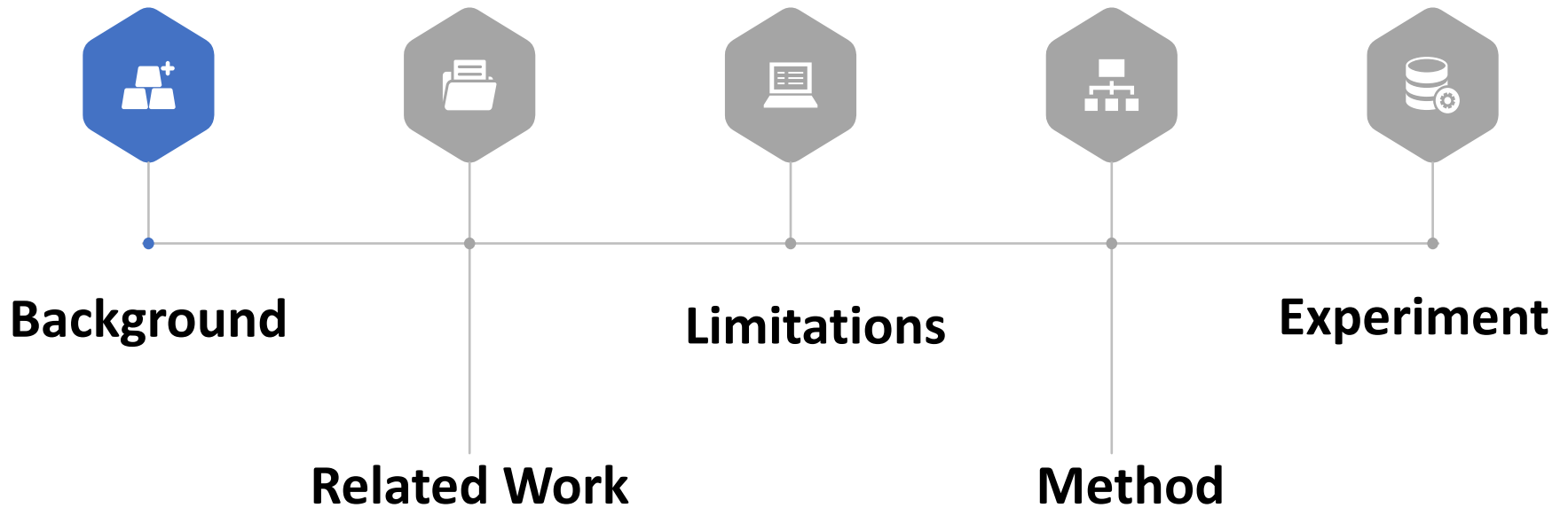Presenter: Liqiang Jing

jingliqiang6@gmail.com

# Outline



Background

Related Work

Limitations

Method

Experiment

# Background

Online shopping platform always shows various information about the product.



https://www.oberlo.com/blog/online-shopping-statistics

# Background

➤ Multi-Modal Product Summarization Generation (MPSG)

# Outline



Background

Related Work

Limitations

Method

Experiment

# Related Work

> ## Traditional Product Summarization



Jiahui Liang et al. CUSTOM: Aspect-Oriented Product Summarization for E-Commerce. NLPCC. 2021.

# Related Work

➢ Multi-Modal Product Summarization





Haoran Li et al. Aspect-Aware Multimodal Summarization for Chinese E-Commerce Products. AAAI. 2020.

# Outline



Background

Related Work

Limitations

Method

Experiment

# Limitations

➢ L1: Overlook the benefit of pre-training.

➢ L2: Lack the representation-level supervision.

➢ L3: Ignore the diversity of the seller-generated data.

# Outline

Background

Related Work

Limitations

Method

Experiment

# Method



Fig.1: The proposed V2P scheme.

# Method



**Vision-based Attribute Prediction**

➢ Attribute Prediction

$$\hat{s} = Swin(V),$$

$$\widehat{\mathcal{A}} = \{a^k | \hat{s}^k \geq \theta, k = 1, \cdots, S\},$$

$$L_{BCE} = \min_{\Theta_a} -[s \ln(\sigma(\hat{s})) + (1-s) \ln(1 - \sigma(\hat{s}))],$$

# Method



**Attribute Prompt-guided Summary Generation**

**Text and Attribute Prompt Embedding**

$$\widehat{X} = [X, a^{k_1}, a^{k_2}, \cdots, a^{k_Q}]$$

$$\mathbf{e}^j = \mathbf{W}^T \mathbf{g}^j, \quad j = 1, \cdots, U,$$

$$\mathbf{E} = [\mathbf{e}^1; \mathbf{e}^2; \cdots; \mathbf{e}^U] + \mathbf{E}_{pos},$$

**BART-based Summary Generation**

$$\mathbf{Z} = \mathcal{E}(\mathbf{E}),$$

$$\hat{\mathbf{p}}_j = \mathcal{D}(\mathbf{Z}, \hat{y}_1, \hat{y}_2, \cdots, \hat{y}_{j-1}),$$

13

# Method



**Attribute Prompt-guided Summary Generation**

**Output-level and Representation-level Supervision**

$$\mathcal{L}_{CE} = -\frac{1}{L}\sum_{j=1}^{L} log(\hat{\mathbf{p}}_j[t*]),$$

$$\mathcal{L}_{ReS} = \mathbb{E}_{p(z_{1:B})p(\tilde{z}|z_1)}\Big[ -log\frac{e^{f(z^y,z_1)}}{\sum_{i=1}^{B} e^{f(z^y,z_i)}}\Big],$$

**Data Augmentation-based Robustness Regularization**

$$\mathcal{L}'_{CE} = -\frac{1}{L}\sum_{i=1}^{L} log(\hat{\mathbf{p}}'_j[t*]),$$

$$\mathcal{L}_{KL} = \sum_{j=1}^{L} D_{KL}(\hat{\mathbf{p}}_j\|\hat{\mathbf{p}}'_j) = \sum_{j=1}^{L}\sum_{k=1}^{L} \hat{p}_{jk} log(\frac{\hat{p}_{jk}}{\hat{p}'_{jk}}).$$

# Outline



Background

Related Work

Limitations

Method

Experiment

# Experiment

➤ Dataset

Table 1: Dataset statistics of CEPSUM.

| Category | Home Appliances | Clothing | Cases&Bags |
|---|---|---|---|
| #Train Sample | 437,646 | 790,297 | 97,510 |
| #Valid Sample | 10,000 | 10,000 | 5,000 |
| #Test Sample | 10,000 | 10,000 | 5,000 |
| Avg Input Len | 335 | 286 | 299 |
| Avg Output Len | 79 | 78 | 79 |

➤ Evaluation Metric

Rouge-1, Rouge-2, Rouge-L.

# Experiment

## ➢ On Model Comparison

Table 2: Performance comparison among different methods.

| Dataset/ Model | Home Appliances | | | Clothing | | | Cases&Bags | | |
|---|---|---|---|---|---|---|---|---|---|
| | Rouge-1 | Rouge-2 | Rouge-L | Rouge-1 | Rouge-2 | Rouge-L | Rouge-1 | Rouge-2 | Rouge-L |
| Lead | 21.97 | 9.54 | 12.79 | 19.83 | 8.39 | 13.56 | 21.49 | 9.37 | 14.19 |
| LexBank | 24.06 | 10.01 | 18.19 | 26.87 | 9.01 | 17.76 | 27.09 | 9.87 | 18.03 |
| Seq2Seq | 21.57 | 7.18 | 17.61 | 23.05 | 6.84 | 16.82 | 23.18 | 6.94 | 17.29 |
| MASS | 28.19 | 8.02 | 18.73 | 26.73 | 8.03 | 17.72 | 27.19 | 9.03 | 18.17 |
| PG | 31.31 | 10.93 | 21.11 | 29.11 | 9.24 | 19.92 | 31.11 | 10.27 | 21.79 |
| MMPG | 32.88 | 11.88 | 21.96 | 30.73 | 10.29 | 21.25 | 32.69 | 11.78 | 22.27 |
| VG-Bart-1 | 32.71 | 11.46 | 22.87 | 31.36 | 9.94 | 21.34 | 32.73 | 10.26 | 22.44 |
| VG-Bart-2 | 32.73 | 11.74 | 23.61 | 31.63 | 10.08 | 21.64 | 33.30 | 11.31 | 23.13 |
| **V2P** | **34.47*** | **12.63*** | **25.09*** | **35.05*** | **11.98*** | **22.62*** | **34.65*** | **11.89*** | **24.53*** |
| Improvement . ↑ | 4.84% | 6.32% | 6.27% | 10.81% | 16.42% | 5.06% | 4.05% | 0.93% | 6.05% |

V2P consistently surpasses all the baselines, exhibiting the effectiveness of the proposed scheme.

# Experiment

## ➢ On Ablation Study

Table 3: Performance comparison between V2P and its derivatives.

| Dataset/ Model | Home Appliances | | | Clothing | | | Cases&Bags | | |
|---|---|---|---|---|---|---|---|---|---|
| | Rouge-1 | Rouge-2 | Rouge-L | Rouge-1 | Rouge-2 | Rouge-L | Rouge-1 | Rouge-2 | Rouge-L |
| **V2P** | **34.47** | **12.63** | **25.09** | **35.05** | **11.98** | **22.62** | **34.65** | **11.89** | **24.53** |
| V2P-w/o-Image | 32.78 | 11.38 | 24.26 | 30.71 | 9.67 | 20.96 | 32.95 | 10.53 | 23.17 |
| V2P-w/o-ReS | 33.74 | 11.73 | 24.64 | 34.81 | 11.67 | 21.76 | 33.77 | 11.11 | 23.48 |
| V2P-w/o-Robust | 34.22 | 12.39 | 24.76 | 34.62 | 11.47 | 21.10 | 33.19 | 10.02 | 22.56 |
| V2P-w-VGG | 32.90 | 11.67 | 23.92 | 34.31 | 11.60 | 21.07 | 33.46 | 10.99 | 23.72 |
| V2P-w-Res | 33.15 | 11.90 | 24.12 | 34.49 | 11.67 | 21.13 | 19.44 | 4.58 | 13.34 |

V2P obtains the best performance, which verifies these components are significant in our model.

# Experiment

➢ On Case Study



**Product Long Text Description:**
Off season women's classic Pima cotton cardigan, navy blue, cotton, slim fit, long sleeve, round neck, suitable for different seasons. This cardigan is made of ring spun Pima cotton with high softness, slim scissors, slim cutting, long sleeve, rib knitted cuffs. Rib knit hem. Button placket.

**Product Image:**

**GT Summary:** This cardigan is made of ring spun Pima cotton with high softness. The slim fit is suitable for different seasons. Thread knitted hem effectively modifies the waist line, which is quite thin.

**V2P:** This cardigan is made of ring spun leather horse cotton with high softness. It is exquisite and fashionable. The fabric is neat and stylish. The slim version shows more temperament. The cardigan design modifies the facial lines, and the button cardigan is easy to wear and take off.
(Rouge-1: 56.20%, Rouge-2: 42.02%, Rouge-L: 54.00%)

**V2P-w/o-Image:** Made of high-quality Pima cotton fabric, it feels soft and delicate, has good air permeability and brings a comfortable wearing experience. The classic round neck design naturally fits the neck and is beautiful.
(Rouge-1: 21.90%, Rouge-2: 4.44%, Rouge-L: 15.93%)

Figure 2: Comparison between the summaries generated by our model and its variant V2P-w/o-image for a clothing product.

# Conclusion

➢ We design a **vision-to-prompt** based multi-modal product summary generation scheme, where the **heterogeneous multi-modal data** are unified in the same space by converting the vision modality into **semantic attribute prompts**.

➢ We are the first to perform both **output-level** and **representation-level** supervision for MPSG simultaneously and introduce the **data augmentation-based** robustness regularization.

➢ We are the first to adapt the **generative pre-trained language model** to solve the MPSG task and achieve superior performance on a real-world dataset.

# Thanks for your listening.



**Codes are available!**