

# SANCL: Multimodal Review Helpfulness Prediction with Selective Attention and Natural Contrastive Learning

Anonymous ACL submission

## Abstract

With the boom of e-commerce, Multimodal Review Helpfulness Prediction (MRHP) that identifies the helpfulness score of multimodal product reviews has become a research hotspot. Previous work on this task focuses on attention-based modality fusion, information integration, and relation modeling, which primarily exposes the following drawbacks: 1) the model may fail to capture the really essential information due to its indiscriminate attention formulation; 2) lack appropriate modeling methods that takes full advantage of correlation among provided data. In this paper, we propose SANCL: Selective Attention and Natural Contrastive Learning for MRHP. SANCL adopts a probe-based strategy to enforce high attention weights on the regions of greater significance. It also constructs a contrastive learning framework based on natural matching properties in the dataset. Experimental results on two benchmark datasets with three categories show that SANCL achieves state-of-the-art baseline performance with lower memory consumption.

## 1 Introduction

We have witnessed an acceleration towards an e-commerce boom that has transpired over the past decades (Vulkan, 2020). In the virtual bazaar, countless deals are made between mutually invisible sellers and customers from time to time. For customers, it may be their biggest headache to determine whether they should pay for a good when being overwhelmed by tempting advertisements, as they can hardly learn about the true information about a product in face of the seller’s meticulous promotion without any references. In this situation, reviews in e-shops that can provide justification information, are thus of great value to customers. However, the quality of reviews under a certain product page can be disparate—many customers are willing to leave informative feedback on the product, while many others arbitrarily write a few

words and even paste irrelevant messages in their comments. Therefore, from the perspective of online shopping platforms, they would be welcome and attractive to customers if they provide a service that can intelligently filter and place the most helpful reviews at the top position. The task in the machine learning field to solve this problem is Review Helpfulness Prediction (RHP) (Tang et al., 2013).

As the thriving of multimodal learning research and the handy accessibility of multimodal data in this Internet era, the latest progress incorporated image (vision modality) information into the review helpfulness prediction (RHP) (Liu et al., 2021) as Multimodal RHP (MRHP). Although previous work attained excellent results in MRHP, there are still some drawbacks. First, the attention mechanism in these works for representation learning follows the most basic setting—it directly computes out the attention scores based on the representation vectors of tokens or sentences, without any further intervention on the obtained weights (Fan et al., 2019). Generally, the amount of task-related information in each sentence in a given piece of review may vary greatly—since customers usually casually write these reviews and may insert some meaningless words, such as emotional appreciation or complaint that can not benefit the viewers. We observed that due to dataset characteristics in the MRHP task, there are cues to help locate those key sentences in the review text. Therefore, we proposed a probe-based selective attention mechanism to employ them for better attention results.

Secondly, it has been revealed that the correlation, e.g., the similarity of feature vectors, among multimodal and multi-domain data is an essential factor in modeling (Xu et al., 2020; Chen et al., 2019). Nevertheless, existing studies (Xu et al., 2020; Liu et al., 2021) simply quantified them in similarity metrics, such as cosine value, for direct classification use. Though gained appreciative re-

042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082

sults, we believe they can be better utilized through the contrastive learning scheme to enhance the quality of learned representations.

In this paper, we propose a novel framework, SANCL, to incorporate these two points. In SANCL we first generate a special “probe” mask that highlights the key sentences from the product and review text. The mask then attends the computation attention modules to help focus more on those task-related sentences. Then we construct a contrastive learning framework to learn better modality representations with internal correlations of data. Based on contrastive predictive coding (CPC) (Oord et al., 2018), the framework is composed of two feature spaces (domains). Each domain takes specific combinations of projected representations as input, according to their relation types from our analysis. Through optimization over the contrastive score, the multimodal and multi-domain representations can learn from the inherent relations. Our contribution can be summarized as follows:

- We design a selective attention approach, including the probe mask generation and mask-based attention computation, for the information aggregation in MRHP tasks.
- We analyze the characteristics and relations in multimodal reviews and formulate a contrastive learning framework to refine the learned representations.
- Extensive experiments on three publicly available datasets show our approach achieves state-of-the-art performance with lower memory consumption.

## 2 Related Work

In this section, we briefly recap some relevant work in the field of review helpfulness prediction and multimodal contrastive learning.

**Review Helpfulness Prediction** Customer reviews play an important role in helping customers investigate products before determining whether to purchase. (Zhu and Zhang, 2010; Diaz and Ng, 2018; Gamzu et al., 2021). Support vector regression (SVM) was first employed to automatically judge the review helpfulness (Kim et al., 2006; Zhang and Varadarajan, 2006; Tsur and Rappoport, 2009). Later, linear regression (Lu et al., 2010;

Ghose and Ipeirotis, 2010), extended tensor factorization (Moghaddam et al., 2012), and probabilistic matrix factorization models (Tang et al., 2013) have been applied to integrate complicated constraints into the learning process. With the development of deep learning, deep neural networks (Lee and Choeh, 2014; Fan et al., 2018; Chen et al., 2018) have been utilized to model the sophisticated elements in this task. Recently, Qu et al. (2020) proposed a graph neural network to capture the intrinsic relationship between the products and their reviews. However, most existing studies only focus on the text of reviews, neglecting the images that usually exist in online reviews. This paper takes advantage of the images and proposes a novel contrastive learning framework with a selective attention mechanism to learn expressive multimodal features.

**Multimodal Representation Learning** The foremost problem of multimodal tasks lies in multimodal representation learning (Baltrušaitis et al., 2018). The concept of multimodal representation learning covers many techniques, such as multimodal fusion (Vielzeuf et al., 2018; Wang et al., 2020; Mai et al., 2020; Han et al., 2021a), multimodal contrastive learning (Yuan et al., 2021; Han et al., 2021b), etc. Attention-based architectures are the basic routine in multimodal fusion, but the formulations are similar. In this paper, knowing about the particularity of MRHP and its dataset, we devise a novel attention mechanism to better aggregate information in textual data. Additionally, we also upgrade the application of contrastive learning. Unlike the ordinary treatment that divides samples into positive and negative groups according to “from myself” or “not from myself” (Cui et al., 2020; Liang et al., 2020), we extract contrastive pairs according to the natural correlation in the dataset and construct the framework of two feature spaces termed as domains.

## 3 Method

In this section, we first introduce the problem definition of Multimodal Review Helpfulness Prediction (MRHP). Then we elaborate on the model architecture and processing pipeline of our method.

### 3.1 Problem Definition

Given a collection of product descriptions  $\mathcal{P} = \{P_1, P_2, \dots, P_N\}$  and associated reviews  $\mathcal{R} = \{R_1, R_2, \dots, R_N\}$  gleaned from an e-shopping

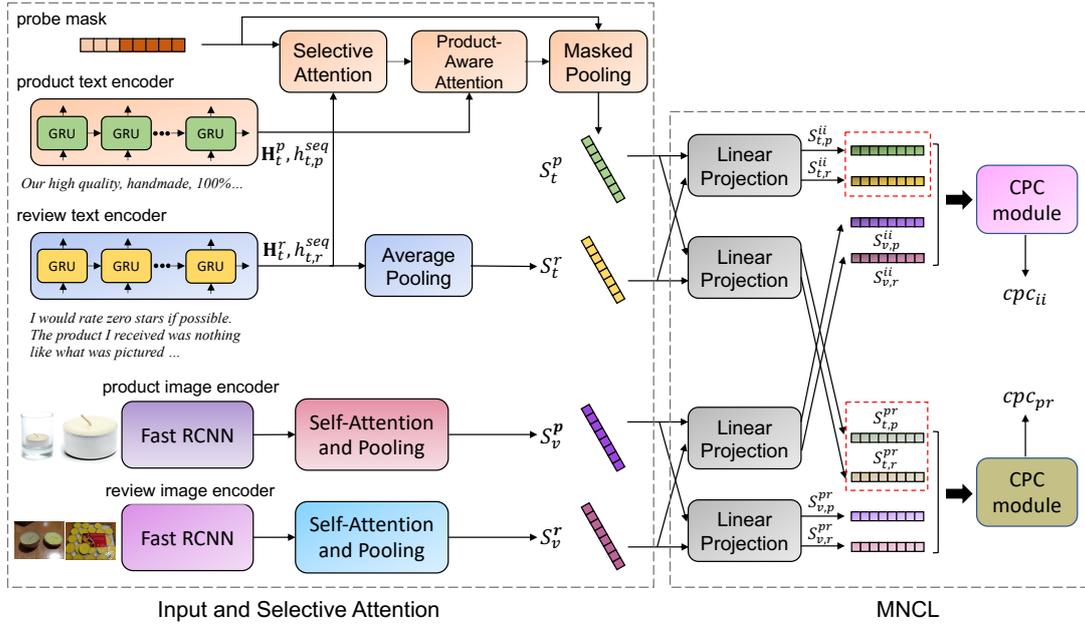


Figure 1: The overview of SANCL. The output layer is omitted. Features in red boxes ( $S_{v,r}^{ii}, S_{t,r}^{ii}, S_{v,r}^{pr}, S_{t,r}^{pr}$ ) are used in final helpfulness score prediction.

website. Each product description  $P_i \in \mathcal{P}$  contains the product name  $n_{p_i}$  plus the text and image descriptions  $T_{p_i}$  and  $I_{p_i}$ . The underlying review collection  $R_i$  associated with product  $i$  contains  $m$  review pieces  $R_i = \{r_{i1}, r_{i2}, \dots, r_{im}\}$ . Each review data frame is composed of images  $I_{r_{ij}}$  and text  $T_{r_{ij}}$  as well. We exhibit an example of input data at the model’s input position in Figure 1. All review pieces are annotated with helpfulness scores  $s_{ij} \in \{0, 1, 2, 3, 4\}$ . Multimodal review helpfulness prediction can be formulated as a regression task that aims to predict the helpfulness score of each review piece, and a ranking task to sort these reviews by their scores in descending order.

### 3.2 Overview

The overall architecture of SANCL is depicted in Figure 1. We first generate a probe mask for each review according to the corresponding product name and review text as shown in Figure 2. The probe mask highlights the sentences that mention the product, which then participates in the computation of selective attention to produce text representations. For images, we feed the features extracted by pre-trained visual neural networks to two self-attention modules to produce image representations. Then we project these representations of each modality in both product description and customer review into two shared spaces (domains). We finally develop a contrastive learning module to compute the cross-modality and review-product contrastive scores, which further improves the qual-

ity of representations output from attention modules.

### 3.3 Input Encoding

**Context-aware Textual Representation** For both review and product text, we initialize the token representations with GloVe (Pennington et al., 2014)<sup>1</sup> or pre-trained models as  $\mathbf{E}_t = \{e_1^t, e_2^t, \dots, e_l^t\} \in \mathbb{R}^{l \times d_e^t}$ , where  $l$  is the length (number of tokens) of a given sentence and  $d_e^t$  is the embedding dimension. We then send these embeddings to a uni-directional Gated Recurrent Unit (GRU) (Cho et al., 2014), yielding token-wise and sequence representations  $\mathbf{H}_t = \{h_1^t, \dots, h_l^t\}$  and  $h_t^{seq}$ :

$$\mathbf{H}_t, h_t^{seq} = \text{GRU}(\mathbf{E}; \theta_t). \quad (1)$$

where  $\theta_t$  is the parameters in GRU.

**Visual Feature Extraction** We apply Faster R-CNN (Ren et al., 2015) on raw images and yield the hidden representations  $\mathbf{E}_v = \{e_1^v, e_2^v, \dots, e_n^v\} \in \mathbb{R}^{n \times d_e^v}$  in the last layer before the classifier to map the Regions of Interest (RoI) in an image to a hidden space, where  $n$  is the number of hot regions detected in the image and  $d_e^v$  is the vector lengths of hidden representations. Then same as Liu et al. (2021), we feed them into a self-attention module that outputs the encoded image representations  $\mathbf{H}_v = \{h_1^v, h_2^v, \dots, h_n^v\}$ .

<sup>1</sup>We used glove.840B.300d in our experiments.

### 3.4 Probe-based Selective Attention (PSA)

Having gained elementary encoded multimodal representations, interactions between parallel review pieces and corresponding product descriptions are required to form the product-aware review representations. Previous work primarily formulated these interactions as token-wise description–review attention (Fan et al., 2019; Qu et al., 2020). Though succinct and effective, this token-by-token or sentence-by-sentence computation scheme may neglect distinct the relative importance among sentences and missing really task-related information. Because only through loss back-propagation without any re-weighting operation, it can not always be ensured that larger weights will be put on those key sentences. To mitigate this issue, we propose the selective attention approach. It first generates a special “probe” mask then performs discriminate attention based on that.

**Probe Mask Generation** The probe mask should reflect the position (i.e., in which sentence) where the product is mentioned in a review. An example of the generation process is displayed in Figure 2. We first retrieve the core words from the product name by looking up its dependency tree and picking the lemmatized form of the words around the root. Next, we use coreference resolution to identify all coreference clusters in the review.

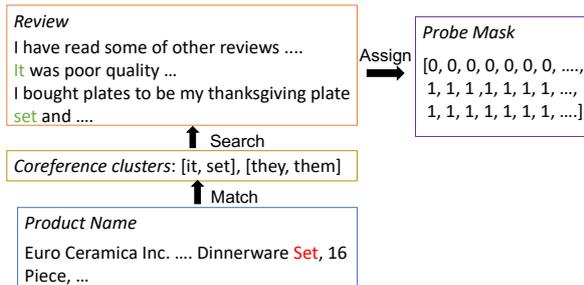


Figure 2: An example of mask generation

There are three possible resolution results: (1) A cluster containing the core word of the product name; (2) At least one cluster exists but the core word is missing in all clusters; (3) No coreference cluster exists. For (1) and (3), we do not require extra steps as the existence of entity clusters can be confirmed. For (2), we are still uncertain whether an entity cluster is in the text. We devise a simple rule to tackle this situation—we regard the first

cluster as product name mention cluster, based on our observation that the first repeatedly mentioned pronouns in a review are more likely to refer to the product. After locating these product name mentions, we create the probe mask  $M \in \mathbb{R}^{1 \times l}$  by assigning 1 to the positions of those mentioned sentences and 0 to others. The process is summarized in Algorithm 1.

---

#### Algorithm 1: Probe Mask Generation

---

**Input:** Review sentences  $R$ , product name  $P$   
**Output:** Probe mask  $M$

*# core words and coreference clusters extraction:*  
 $\hat{R} \leftarrow \text{Lemmatise}(R), \hat{P} \leftarrow \text{Lemmatise}(P);$   
*# core words extraction:*  
 $T \leftarrow \text{DependencyParse}(\hat{P});$   
 $W \leftarrow \text{FindWordsNearRoot}(T);$   
 $\text{clusters} \leftarrow \text{FindCoreferenceCluster}(\hat{R})$

*# mask generation:*  
 $M \leftarrow \text{ZeroInit}(R.\text{size})$   
**if**  $C = \emptyset$  **then**  
  | **return**  $M$   
**end**  
**foreach**  $c$  **in**  $\text{clusters}$  **do**  
  | **if any**  $w \in W$  **in**  $c$  **then**  
  | |  $\text{gold\_cluster} = c$   
  | **end**  
**end**  
**if**  $\text{gold\_cluster} = \emptyset$  **then**  
  |  $\text{gold\_cluster} = \text{clusters}[0]$   
**end**  
**foreach**  $\text{sent} \in \hat{R}$  **do**  
  | **if any**  $w \in \text{gold\_cluster}$  **in**  $\text{sent}$  **then**  
  | |  $M[\text{sent.start} : \text{sent.end}] \leftarrow \text{True}$   
  | **end**  
**end**  
**return**  $M$

---

**Selective Attention with Probe Mask** There are three steps to acquire product-aware review representations—self-attention, cross-text attention, and pooling, among which the first and last steps take advantage of probe masks generated. We first transform the probe mask to a new form:

$$M' = \alpha M + \beta(1 - M), \quad (2)$$

where  $\alpha > \beta > 0$  since we expect the mask could help focus more on the sentences where the product is mentioned. This effect embodies in the self attention computation of the review text  $\mathbf{H}_t^r$ , where the fundamental attention weights are computed as:

$$\mathbf{A} = \text{softmax}(\mathbf{W}\mathbf{H}_t^r), \quad (3)$$

We renew the original attention matrix  $A \in \mathbb{R}^{l \times l}$ :

$$\mathbf{A}' = (M')^T M' \odot \mathbf{A}, \quad (4)$$

In this process, the attention weights are actually re-weighted as

$$a'_{ij} = \begin{cases} \alpha^2 a_{ij}, & \text{if } m_i = m_j = 1 \\ \alpha\beta a_{ij}, & \text{if } m_i = 1, m_j = 0 \\ \beta^2 a_{ij}, & \text{if } m_i = m_j = 0 \end{cases}, \quad (5)$$

An intuitive explanation to this would be that a token receives more information from hot regions (whose mask value is 1) than non-hot regions, and the strength ratio of these two regions is  $\alpha/\beta$ . Naturally we set  $\alpha$  to 1.0 while generating  $\beta$  individually for each review from its sequence representation:

$$\beta = \text{sigmoid}(\mathbf{W}_{gen} h_{t,r}^{seq} + b_{gen}), \quad (6)$$

where  $\mathbf{W} \in \mathbb{R}^{d_h \times 1}$  is the weight matrix and  $b \in \mathbb{R}$  is the bias. The sigmoid function ensures that  $\beta \in (0, 1)$ . Then we acquire the self-attention results as in common practice:

$$\mathbf{H}_t^{r'} = \mathbf{H}_t^r + \mathbf{A}' \mathbf{H}_t^r, \quad (7)$$

In cross-text attention, since weighted-sum is performed on the product text, we do not utilize probe masks in this stage and obtain  $\mathbf{H}_t^{r''}$ . Finally, we average the result with the probe mask by a weighted sum to aggregate these sentence representations:

$$S_t^r = \text{weighted\_sum}(\mathbf{H}_t^{r''}, M^r). \quad (8)$$

Note that for review image representations there are only cross-image attention and average pooling to yield  $S_v^r$ .

### 3.5 Multi-domain Natural Contrastive Learning (MNCL)

From the theory of mutual information, training to distinguish positive samples from negative ones in terms of their similarity can enrich the learned representations and enhance downstream tasks' performance. In our work, we are concerned about the natural relations and split them into two domains: the inner-instance domain (*ii*) and product-review (*pr*) domain. Before forwarding the input representations into the MNCL module, all pooled representations are projected to the domain-specific shared spaces through two linear layers and an activation layer in between. We denote them as  $S_{m,f}^d$ , where  $m \in \{t, v\}$  is the modality type,  $f \in \{r, p\}$  is the field (review or product description) and  $d \in \{ii, pr\}$  is the domain name:

$$S_{m,f}^d = W_{m,f,2}^d \text{Tanh}(W_{m,f,1}^d S_m^d + b_{m,f,1}^d) + b_{m,f,2}^d \quad (9)$$

where  $W_{m,*,i}^d$  and  $b_{m,*,i}^d$  are weights and biases in the  $i$ -th layer of the projection network. Note that the data in the same modality and domain share the same network parameters. In the succeeding content, we are going to describe details of the two contrastive-learning domains, mainly concerning how to pick positive and negative samples for contrastive learning and training.

**Inner Instance (II) Domain** In the inner instance domain, we separate positive and negative pairs according to how similar the representations between image and text are in a single training instance. First, from the sellers' perspective, the text and image of a product should match well so as to attract customers. Thus we mark text-image pairs of product descriptions as positive ones (the set of these pairs is denoted as  $\mathbb{S}_{ii}^p$ ). Therefore, we mark the former as positive (the set is denoted as  $\mathbb{S}_{ii}^+$ ) and the latter as negative ones (the collection is denoted as  $\mathbb{S}_{ii}^-$ ). Besides, from our observation, reviews that achieve high helpfulness scores possess a high similarity between its text and the attached image.

**Product-Review (PR) Domain** The semantic matching property also exists between product descriptions and their associated reviews. As helpfulness is dependent on how well a review is pertinent to the theme of the product, we argue that review pieces of high helpfulness scores ( $\mathbb{S}_{pr}^+$ ) should match the product introduction both visually and literally, while those low-score pieces ( $\mathbb{S}_{pr}^-$ ) match the introduction poorly in both modalities.

**Multi-domain Contrastive Predictive Coding (MCPC)** In contrastive predictive coding (Oord et al., 2018), we need to compute contrastive scores for every sample pair. According to the common approach (Yuan et al., 2021; Han et al., 2021b), exponential function is chosen as the score function:

$$\varphi(\mathbf{A}, \mathbf{B}) = \exp\left(\frac{\text{norm}(\mathbf{A}^T) \text{norm}(\mathbf{B})}{\tau}\right), \quad (10)$$

where  $\text{norm}(\ast)$  is the l2-norm function,  $\tau$  is the temperature hyper-parameter, for simplicity we keep its value 1.0 in our experiments. By noise contrastive estimation (Gutmann and Hyvärinen, 2010), in the inner instance domain the score is

384 computed as:

$$385 \quad cpc_{ii} = - \sum_{(S_{t,j}, S_{v,j}) \in (\mathbb{S}_{ii}^+ \cup \mathbb{S}_{ii}^P)} \log \frac{\varphi(S_{t,j}, S_{v,j})}{\sum_{S_k \in (\mathbb{S}_{ii}^+ \cup \mathbb{S}_{ii}^- \cup \mathbb{S}_{ii}^P)} \varphi(S_{t,k}, S_{v,k})}, \quad (11)$$

388 where  $(S_{t,j}, S_{v,j})$  are the text-image pair from the  
 389 instance, i.e., a review piece or product descrip-  
 390 tion. The summation is over  $\mathbb{S}_{ii}^+$  and  $\mathbb{S}_{ii}^P$  because  
 391 instances counted here are from both product de-  
 392 scriptions and review pieces. Similarly in product-  
 393 review domain the score is:

$$394 \quad cpc_{pr}^m = - \sum_{S_{m,j}^r \in \mathbb{S}_{pr}^+} \log \frac{\varphi(S_{m,j}^r, S_{m,j}^P)}{\sum_{S_k \in (\mathbb{S}_{pr}^+ \cup \mathbb{S}_{pr}^-)} \varphi(S_{m,k}^r, S_{m,k}^P)}. \quad (12)$$

$$395 \quad cpc_{pr} = cpc_{pr}^t + cpc_{pr}^v \quad (13)$$

396 where  $S_{m,j}^r$  is the representation of modality  $m$  in  
 397 review  $r$  from the positive review set  $\mathbb{S}_{pr}^+$  and  $S_{m,j}^P$   
 398 is the counterpart of the corresponding product.

### 402 3.6 Prediction and Training

403 We select all review-related representations  
 404 in the common spaces of two domains  
 405  $(S_{t,r}^{ii}, S_{t,r}^{pr}, S_{v,r}^{ii}, S_{v,r}^{pr})$  and concatenate them  
 406 as features for final prediction ( $\mathbf{F}$ ). After concate-  
 407 nating these features, a linear layer takes them as  
 408 input and outputs the helpfulness score predictions  
 409  $\xi_r$ :

$$410 \quad \mathbf{F} = \text{concat}([S_{t,r}^{ii}, S_{t,r}^{pr}, S_{v,r}^{ii}, S_{v,r}^{pr}]) \quad (14)$$

$$411 \quad \xi_r = \mathbf{W}_o \mathbf{F} + b_o, \quad (15)$$

412 where  $\mathbf{W}_o$  and  $b_o$  are the weight matrix and bias  
 413 in the output layer. Same as Liu et al. (2021), we  
 414 adopt the standard pairwise ranking loss as the task  
 415 loss:

$$416 \quad \mathcal{L}_{task} = \sum_i \max(0, \gamma - \xi_{r^+,i} + \xi_{r^-,i}), \quad (16)$$

417 where  $r^+, r^-$  are an arbitrary pair of review pieces  
 418 under product  $P_i$ ,  $\gamma$  is a scaling factor. Contrastive  
 419 losses make up the auxiliary loss:

$$420 \quad \mathcal{L}_{aux} = cpc_{ii} + cpc_{pr} \quad (17)$$

421 Hence the total loss for training is ( $\kappa$  is a hyper-  
 422 parameter to adjust the effect of auxiliary loss):

$$423 \quad \mathcal{L} = \mathcal{L}_{task} + \kappa \mathcal{L}_{aux} \quad (18)$$

## 426 4 Experimental Settings

427 This section presents the datasets used in the exper-  
 428 iments, baseline models, and metrics.

### 429 4.1 Datasets

430 We conduct experiments on three MRHP datasets  
 431 (Liu et al., 2021) in different categories: *Cloth-*  
 432 *ing, Shoes & Jewelry, Home & Kitchen* and *Elec-*  
 433 *tronics*. The text and images in these datasets are  
 434 crawled from Amazon online shops through 2018  
 435 until 2019. The helpfulness scores are annotated  
 436 according to the number of helpfulness votes to  
 437 the reviews. More specifically the scores equal to  
 438  $\lfloor \log_2 n_{votes} \rfloor$  and are clipped into  $[0, 4]$ . Details of  
 439 datasets are provided in Appendix.

### 440 4.2 Baseline Models

441 Following previous work, we first compare our  
 442 model with a bunch of baselines in the text-  
 443 only setting, which examines and signifies the  
 444 enhancement by our selective attention mecha-  
 445 nism and text-related contrastive learning modules.  
 446 The baseline candidates contain Multi-Perspective  
 447 Matching (BiMPM) network (Wang et al., 2017),  
 448 Embedding-gated CNN (EG-CNN) (Chen et al.,  
 449 2018), Convolutional Kernel-based Neural Rank-  
 450 ing Model (Conv-KNRM) (Dai et al., 2018) and  
 451 Product-aware Helpfulness Prediction Network  
 452 (PRHNet) (Fan et al., 2019). In multimodal set-  
 453 tings, we pick a collection of state-of-the-art multi-  
 454 modal helpfulness prediction models for compari-  
 455 son:

- 456 • **SSE-Cross** (Abavisani et al., 2020): The  
 457 Stochastic Shared Embeddings (SSE) Cross-  
 458 modal Attention Network introduces a novel  
 459 cross-attention mechanism that can filter noise  
 460 components from weak modalities which may  
 461 mislead the model to make wrong predictions  
 462 on a sample. SSE is adopted as the regulariza-  
 463 tion technique to alleviate over-fitting in the  
 464 fusion process to further prompt the prediction  
 465 accuracy.
- 466 • **D&R Net** (Xu et al., 2020): The Decomposi-  
 467 tion and Relation Network learns the common-  
 468 ality and discrepancy between image and text  
 469 in decomposition network and the multi-view  
 470 semantic association information in relation  
 471 network.
- 472 • **MCR** (Liu et al., 2021): The Multi-  
 473 perspective Coherent reasoning method in-

| Setting           | Model                               | Cloth. & Jew             |                          |                          | Electronics              |                          |                          | Home & Kitchen           |                          |                          |
|-------------------|-------------------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
|                   |                                     | MAP                      | N-3                      | N-5                      | MAP                      | N-3                      | N-5                      | MAP                      | N-3                      | N-5                      |
| Text-only         | BiMPM* (Wang et al., 2017)          | 57.7                     | 41.8                     | 46.0                     | 52.3                     | 40.5                     | 44.1                     | 56.6                     | 43.6                     | 47.6                     |
|                   | EG-CNN* (Chen et al., 2018)         | 56.4                     | 40.6                     | 44.7                     | 51.5                     | 39.4                     | 42.1                     | 55.3                     | 42.4                     | 46.7                     |
|                   | Conv-KNRM* (Dai et al., 2018)       | 57.2                     | 41.2                     | 45.6                     | 52.6                     | 40.5                     | 44.2                     | 57.4                     | 44.5                     | 48.4                     |
|                   | PRHNet† (Fan et al., 2019)          | 58.23                    | 43.36                    | 47.21                    | 52.31                    | 40.43                    | 43.88                    | 57.11                    | 44.46                    | 48.27                    |
|                   | SANCL (Ours)                        | <b>58.98<sup>‡</sup></b> | <b>44.75<sup>‡</sup></b> | <b>48.57<sup>‡</sup></b> | <b>53.03<sup>‡</sup></b> | <b>41.03<sup>‡</sup></b> | <b>44.77<sup>‡</sup></b> | <b>58.03<sup>‡</sup></b> | <b>45.59<sup>‡</sup></b> | <b>49.31<sup>‡</sup></b> |
|                   | BERT (Devlin et al., 2018)          | 56.47                    | 42.98                    | 46.84                    | 51.95                    | 39.77                    | 43.11                    | 56.62                    | 42.12                    | 46.87                    |
|                   | PRHNet+BERT† (Fan et al., 2019)     | 57.51                    | 43.65                    | 47.74                    | 52.28                    | 40.66                    | 44.02                    | 57.32                    | 44.74                    | 48.42                    |
| SANCL+BERT (Ours) | <b>58.49<sup>‡</sup></b>            | <b>44.91<sup>‡</sup></b> | <b>48.69<sup>‡</sup></b> | <b>53.13<sup>‡</sup></b> | <b>41.77<sup>‡</sup></b> | <b>45.01<sup>‡</sup></b> | <b>58.20<sup>‡</sup></b> | <b>45.83<sup>‡</sup></b> | <b>49.65<sup>‡</sup></b> |                          |
| Multimodal        | SSE-Cross* (Abavisani et al., 2020) | 65.0                     | 56.0                     | 59.1                     | 53.7                     | 43.8                     | 47.2                     | 60.8                     | 51.0                     | 54.0                     |
|                   | D&R Net* (Xu et al., 2020)          | 65.2                     | 56.1                     | 59.2                     | 53.9                     | 44.2                     | 47.5                     | 61.2                     | 51.8                     | 54.6                     |
|                   | MCR† (Liu et al., 2021)             | 66.96                    | 58.03                    | 61.06                    | 55.86                    | 46.32                    | 49.45                    | 63.17                    | 53.85                    | 57.14                    |
|                   | SANCL (Ours)                        | <b>67.26</b>             | <b>58.61<sup>‡</sup></b> | <b>61.48<sup>‡</sup></b> | <b>56.19</b>             | <b>46.98<sup>‡</sup></b> | <b>49.92<sup>‡</sup></b> | <b>63.35</b>             | <b>54.28<sup>‡</sup></b> | <b>57.40</b>             |
|                   | MCR+BERT (Liu et al., 2021)         | 65.81                    | 55.94                    | 58.75                    | 55.15                    | 45.67                    | 48.62                    | 62.39                    | 52.91                    | 56.09                    |
|                   | SANCL+BERT (Ours)                   | <b>66.52<sup>‡</sup></b> | <b>56.73<sup>‡</sup></b> | <b>59.90<sup>‡</sup></b> | <b>56.04<sup>‡</sup></b> | <b>46.77<sup>‡</sup></b> | <b>49.95<sup>‡</sup></b> | <b>62.74</b>             | <b>53.65<sup>‡</sup></b> | <b>56.91<sup>‡</sup></b> |

Table 1: Results on three datasets; all reported metrics are the average of five runs; “\*” are from Liu et al. (2021) and “†” are from the open-source code in Liu et al. (2021); “‡” represent the results significantly outperforms PRHNet and MCR with p-value < 0.05 based on paired t-test.

474 corporates the joint reasoning across textual  
475 and visual modalities from both the product  
476 and the review. Three types of coherence are  
477 modeled to learn effective modality represen-  
478 tations for the helpfulness prediction.

479 In both settings, we also test our method with  
480 BERT (Devlin et al., 2018) encoder and compare  
481 that to the respective SOTA models on BERT. In ad-  
482 dition, we test and record basic BERT performance  
483 (BERT+a double linear layers).

### 484 4.3 Metrics

485 As MRHP is a ranking task, the metrics for com-  
486 parison are as well ranking-customized. After sort-  
487 ing all prediction-truth scores in descending order,  
488 the Mean Average Precision (MAP) computes the  
489 mean precision of top-1 to top-K samples. K is  
490 usually large enough to ensure top-K can encom-  
491 pass the entire collection of reviews under every  
492 product. The Normalized Discounted Cumulative  
493 Gain (NDCG-N) (Järvelin and Kekäläinen, 2017;  
494 Diaz and Ng, 2018) purely reckons the gain value  
495 over top-N predictions (N is 3 and 5 in our exper-  
496 iments), which simulates the real circumstances  
497 of a typical customer who would always read the  
498 topmost reviews.

## 499 5 Results and Analysis

500 In this section, we will compare our approach with  
501 several advanced baselines and explore how it im-  
502 proves in the multimodal helpfulness prediction  
503 task.

## 504 5.1 Performance Comparison

505 We list the performance of our model and base-  
506 lines in Table 1. Notably, SANCL consistently  
507 outperforms all the baselines in both text-only and  
508 multimodal, BERT and Glove initialization settings.  
509 These outcomes initially demonstrate the efficacy  
510 of our method in MRHP tasks. It is surprising that  
511 we cannot gain significant performance boost by  
512 replacing Glove with BERT as the text encoder. We  
513 speculate the reason is that Glove embeddings are  
514 expressive enough for this task.

515 Moreover, it can be claimed that SANCL is a  
516 lightweight model compared to the multimodal  
517 SOTA, since the model size and GPU memory  
518 consumption of SANCL are much lower than  
519 MCR. The total number of parameters is 2.63M  
520 in MCR and 1.41M (exclude the embedding layer)  
521 in SANCL respectively, which indicates a double  
522 efficiency. The average GPU memory usage of  
523 SANCL during the training on Amazon-MRHP  
524 Home & Kitchen is around 2.4G, while MCR oc-  
525 cupies an average of 13.7G GPU memory during  
526 training, which is 4.7 times higher than SANCL.

## 527 5.2 Ablation Study

528 To verify the benefits of our proposed method, we  
529 carry out comprehensive ablation experiments on  
530 the Amazon electronics dataset, including the selec-  
531 tive attention and contrastive learning components.  
532 In selective attention, we first replace learned  $\beta$   
533 with a fixed value of 0.5, since we find most  $\beta$   
534 values in our experiments are around 0.5. Next, we re-  
535 move the entire selective attention module and only  
536 preserve the primitive attention computation. The

| Description                        | MAP          | N-3          | N-5          |
|------------------------------------|--------------|--------------|--------------|
| SANCL                              | <b>56.19</b> | <b>46.98</b> | <b>49.92</b> |
| Attention                          |              |              |              |
| w/o learned $\beta$ (fixed at 0.5) | 55.61        | 46.37        | 49.58        |
| w/o probe mask                     | 55.43        | 46.11        | 49.45        |
| Contrastive learning               |              |              |              |
| w/o $cpc_{ii}$                     | 55.54        | 46.29        | 49.23        |
| w/o $cpc_{pr}$                     | 55.81        | 46.40        | 49.47        |
| w/o $cpc_{ii}$ and $cpc_{pr}$      | 55.35        | 46.28        | 49.09        |

Table 2: Ablation study of SANCL on the Electronics dataset.

decline in the outcome of both situations manifests that the probe-based selective attention amends the cross-text information exchange between text fields. For multi-domain contrastive learning, we delete the CPC losses of a single or both domains in training. The results indicate that both domains have a positive impact on performance. Moreover, the effect of the two domains does not counteract their collaboration, as we observe accumulated benefits when they operate together.

### 5.3 Case Study

To understand how our model deals with samples in-depth, we randomly picks up a test product-review instance from the test set of Amazon Home & Kitchen to explain how SANCL works and, as shown in Table 3.

In this example, the customer bought the pins to fix the edge of his sofa. Instead of photoing pins themselves, the customer only presented the tidy sofa after installing the pins. We first visualize the attention weights in test time, as shown in Fig. 3. Note that only the first sentence in the review contains the elements in the coreference clusters, which we have emphasized with italics and underline in Table 3. Consistently, we observe the significant larger weights in the region of first sentence (row/column 1-19) while the rest region’s weights are much smaller. We also ran MCR and collect its prediction on this example, and it is clear that MCR commits a severe error here, probably caused by the direct classification on the unimodal cosine similarity. In our approach, as we carefully analyze and classify the positive and negative pairs in the multi-domain contrastive learning framework, the huge semantic similarity between review text and image and between product description and review text, indicated by the high CPC scores  $S_{v,r}^{ii}$ , assists the model to correctly predict the score.

**Product Name:** Twisty Pins for Upholstery, Slipcovers and Bedskirts 50/pkg

**Product description:** Package of 50 Clear Twisty Pins for securing fabrics and accent trims. Nickel plated steel pin 1/2" in diameter clear top, wire twist 3/8" long. Perfect for Medium to light weight fabrics, bed skirts, bed ruffles, slipcovers and upholstery.



**Review (Helpfulness Score: 4):** I bought *these* to pin the loose material on a sofa cover and they worked like a charm. The sofa cover definitely looks form fitting now.



**Predictions:** SANCL: 4.5291 MCR: -1.0832

**CPC score:**  $cpc_{ii} = 0.82$ ,  $cpc_{pr}^t = 0.76$ ,  $cpc_{pr}^v = 0.21$

Table 3: Examples from the Amazon Home & Kitchen test set.

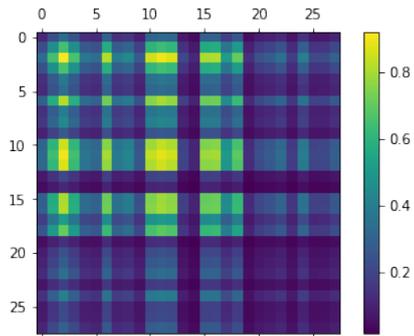


Figure 3: The self-attention weights of the review text in the given example at test time ( $\beta=0.57$ )

## 6 Conclusion

We propose a novel framework, SANCL, for the task of multimodal review helpfulness prediction (MRHP) in this paper. We first present a selective attention mechanism, which purposefully aggregates information from these crucial sentences in the review text by generating the probe mask that exerts re-normalization on the attention weights and pooling stage. We then build up a multi-domain natural contrastive learning framework in our model. It exploits the natural relations among the data from different fields and modalities in the dataset to enhance the model’s capacity of multimodal representation learning. Results of comprehensive experiments and analyses demonstrate the superiority of our model over the comparable baselines and the efficacy of the novel components.

## References

- 592 Mahdi Abavisani, Liwei Wu, Shengli Hu, Joel Tetreault, 647  
593 and Alejandro Jaimes. 2020. [Multimodal categoriza- 648](#)  
594 [tion of crisis events in social media](#). In *Proceedings 649*  
595 [of the IEEE/CVF Conference on Computer Vision 650](#)  
596 [and Pattern Recognition](#), pages 14679–14689.
- 598 Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe 651  
599 Morency. 2018. [Multimodal machine learning: A 652](#)  
600 [survey and taxonomy](#). *IEEE transactions on pattern 653*  
601 [analysis and machine intelligence](#), 41(2):423–443.
- 602 Cen Chen, Minghui Qiu, Yinfei Yang, Jun Zhou, Jun 654  
603 Huang, Xiaolong Li, and Forrest Sheng Bao. 2019. [Multi-domain gated cnn for review helpfulness pre- 655](#)  
604 [diction](#). In *The World Wide Web Conference*, pages 656  
605 2630–2636. 657
- 607 Cen Chen, Yinfei Yang, Jun Zhou, Xiaolong Li, and 658  
608 Forrest Bao. 2018. [Cross-domain review helpfulness 659](#)  
609 [prediction based on convolutional neural networks 660](#)  
610 [with auxiliary domain discriminators](#). In *Proceed- 661*  
611 [ings of the 2018 Conference of the North American 662](#)  
612 [Chapter of the Association for Computational Lin- 663](#)  
613 [guistics: Human Language Technologies, Volume 2 664](#)  
614 [\(Short Papers\)](#), pages 602–607.
- 615 Kyunghyun Cho, Bart Van Merriënboer, Caglar Gul- 665  
616 cehre, Dzmitry Bahdanau, Fethi Bougares, Holger 666  
617 Schwenk, and Yoshua Bengio. 2014. [Learning 667](#)  
618 [phrase representations using rnn encoder-decoder 668](#)  
619 [for statistical machine translation](#). *arXiv preprint 669*  
620 *arXiv:1406.1078*.
- 621 Wanyun Cui, Guangyu Zheng, and Wei Wang. 2020. 670  
622 [Unsupervised natural language inference via decou- 671](#)  
623 [pled multimodal contrastive learning](#). In *Proceed- 672*  
624 [ings of the 2020 Conference on Empirical Methods 673](#)  
625 [in Natural Language Processing \(EMNLP\)](#), pages 674  
626 5511–5520.
- 627 Zhuyun Dai, Chenyan Xiong, Jamie Callan, and 675  
628 Zhiyuan Liu. 2018. [Convolutional neural networks 676](#)  
629 [for soft-matching n-grams in ad-hoc search](#). In *Pro- 677*  
630 [ceedings of the eleventh ACM international confer- 678](#)  
631 [ence on web search and data mining](#), pages 126–134. 679
- 632 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 680  
633 Kristina Toutanova. 2018. Bert: Pre-training of deep 681  
634 bidirectional transformers for language understand- 682  
635 ing. *arXiv preprint arXiv:1810.04805*.
- 636 Gerardo Ocampo Diaz and Vincent Ng. 2018. [Modeling 683](#)  
637 [and prediction of online product review helpfulness: 684](#)  
638 [a survey](#). In *Proceedings of the 56th Annual Meet- 685*  
639 [ing of the Association for Computational Linguistics 686](#)  
640 [\(Volume 1: Long Papers\)](#), pages 698–708.
- 641 Miao Fan, Chao Feng, Lin Guo, Mingming Sun, and 687  
642 Ping Li. 2019. [Product-aware helpfulness prediction 688](#)  
643 [of online reviews](#). In *The World Wide Web Confer- 689*  
644 [ence](#), pages 2715–2721.
- 645 Miao Fan, Yue Feng, Mingming Sun, Ping Li, Haifeng 690  
646 Wang, and Jianmin Wang. 2018. [Multi-task neural 691](#)  
[learning architecture for end-to-end identification of 692](#)  
[helpful reviews](#). In *2018 IEEE/ACM International 693*  
*Conference on Advances in Social Networks Analysis 694*  
*and Mining (ASONAM)*, pages 343–350. IEEE. 695
- Iftah Gamzu, Hila Gonen, Gilad Kutiel, Ran Levy, and 696  
Eugene Agichtein. 2021. [Identifying helpful sen- 697](#)  
[tences in product reviews](#). In *Proceedings of the 698*  
*2021 Conference of the North American Chapter of 699*  
*the Association for Computational Linguistics: Hu- 700*  
*man Language Technologies*, pages 678–691, Online. 701  
Association for Computational Linguistics. 702
- Anindya Ghose and Panagiotis G Ipeirotis. 2010. [Esti- 703](#)  
[mating the helpfulness and economic impact of prod- 704](#)  
[uct reviews: Mining text and reviewer characteristics](#). *IEEE transactions on knowledge and data engineer- 705*  
*ing*, 23(10):1498–1512. 706
- Michael Gutmann and Aapo Hyvärinen. 2010. Noise- 707  
contrastive estimation: A new estimation principle 708  
for unnormalized statistical models. In *Proceedings 709*  
of the thirteenth international conference on artificial 710  
intelligence and statistics, pages 297–304. JMLR 711  
Workshop and Conference Proceedings. 712
- Wei Han, Hui Chen, Alexander Gelbukh, Amir 713  
Zadeh, Louis-philippe Morency, and Soujanya Poria. 714  
2021a. [Bi-bimodal modality fusion for correlation- 715](#)  
[controlled multimodal sentiment analysis](#). In *Pro- 716*  
*ceedings of the 2021 International Conference on 717*  
*Multimodal Interaction*, pages 6–15. 718
- Wei Han, Hui Chen, and Soujanya Poria. 2021b. [Im- 719](#)  
[proving multimodal fusion with hierarchical mutual 720](#)  
[information maximization for multimodal sentiment 721](#)  
[analysis](#). In *Proceedings of the 2021 Conference on 722*  
*Empirical Methods in Natural Language Processing*, 723  
pages 9180–9192. 724
- Kalervo Järvelin and Jaana Kekäläinen. 2017. [Ir eval- 725](#)  
[uation methods for retrieving highly relevant doc- 726](#)  
[uments](#). In *ACM SIGIR Forum*, volume 51, pages 727  
243–250. ACM New York, NY, USA. 728
- Soo-Min Kim, Patrick Pantel, Timothy Chklovski, and 729  
Marco Pennacchiotti. 2006. [Automatically assessing 730](#)  
[review helpfulness](#). In *Proceedings of the 2006 Con- 731*  
*ference on empirical methods in natural language 732*  
*processing*, pages 423–430. 733
- Sangjae Lee and Joon Yeon Choeh. 2014. [Predicting 734](#)  
[the helpfulness of online reviews using multilayer 735](#)  
[perceptron neural networks](#). *Expert Systems with 736*  
*Applications*, 41(6):3041–3046. 737
- Zujie Liang, Weitao Jiang, Haifeng Hu, and Jiaying 738  
Zhu. 2020. [Learning to contrast the counterfactual 739](#)  
[samples for robust visual question answering](#). In *Proceedings of the 2020 Conference on Empirical 740*  
*Methods in Natural Language Processing (EMNLP)*, 741  
pages 3285–3292. 742

|     |   |   |     |
|-----|---|---|-----|
| 700 | Junhao Liu, Zhen Hai, Min Yang, and Lidong Bing.                              | Valentin Vielzeuf, Alexis Lechery, Stéphane Pateux,                         | 756 |
| 701 | 2021. <a href="#">Multi-perspective coherent reasoning for help-</a>          | and Frédéric Jurie. 2018. <a href="#">Centralnet: a multilayer ap-</a>      | 757 |
| 702 | <a href="#">fulness prediction of multimodal reviews</a> . In <i>Pro-</i>     | <a href="#">proach for multimodal fusion</a> . In <i>Proceedings of the</i> | 758 |
| 703 | <i>ceedings of the 59th Annual Meeting of the Asso-</i>                       | <i>European Conference on Computer Vision (ECCV)</i>                        | 759 |
| 704 | <i>ciation for Computational Linguistics and the 11th</i>                     | <i>Workshops</i> , pages 0–0.   | 760 |
| 705 | <i>International Joint Conference on Natural Language</i>                     |   |     |
| 706 | <i>Processing (Volume 1: Long Papers)</i> , pages 5927–                       | Nir Vulkan. 2020. <i>The Economics of E-commerce</i> .                      | 761 |
| 707 | 5936.   | Princeton University Press.   | 762 |
| 708 | Yue Lu, Panayiotis Tsaparas, Alexandros Ntoulas, and                          | Yikai Wang, Wenbing Huang, Fuchun Sun, Tingyang                             | 763 |
| 709 | Livia Polanyi. 2010. <a href="#">Exploiting social context for</a>            | Xu, Yu Rong, and Junzhou Huang. 2020. <a href="#">Deep mul-</a>             | 764 |
| 710 | <a href="#">review quality prediction</a> . In <i>Proceedings of the 19th</i> | <a href="#">timodal fusion by channel exchanging</a> . <i>Advances in</i>   | 765 |
| 711 | <i>international conference on World wide web</i> , pages                     | <i>Neural Information Processing Systems</i> , 33.                          | 766 |
| 712 | 691–700.  |   |     |
| 713 | Sijie Mai, Haifeng Hu, and Songlong Xing. 2020.                               | Zhiguo Wang, Wael Hamza, and Radu Florian. 2017.                            | 767 |
| 714 | <a href="#">Modality to modality translation: An adversarial rep-</a>         | <a href="#">Bilateral multi-perspective matching for natural lan-</a>       | 768 |
| 715 | <a href="#">resentation learning and graph fusion network for</a>             | <a href="#">guage sentences</a> . In <i>IJCAI</i> .                         | 769 |
| 716 | <a href="#">multimodal fusion</a> . In <i>Proceedings of the AAAI Con-</i>    | Nan Xu, Zhixiong Zeng, and Wenji Mao. 2020. <a href="#">Reason-</a>         | 770 |
| 717 | <i>ference on Artificial Intelligence</i> , volume 34, pages                  | <a href="#">ing with multimodal sarcastic tweets via modeling</a>           | 771 |
| 718 | 164–172.  | <a href="#">cross-modality contrast and semantic association</a> . In       | 772 |
| 719 | Samaneh Moghaddam, Mohsen Jamali, and Martin Es-                              | <i>Proceedings of the 58th Annual Meeting of the Asso-</i>                  | 773 |
| 720 | ter. 2012. <a href="#">Etf: extended tensor factorization model</a>           | <i>ciation for Computational Linguistics</i> , pages 3777–                  | 774 |
| 721 | <a href="#">for personalizing prediction of review helpfulness</a> .          | 3786.   | 775 |
| 722 | In <i>Proceedings of the fifth ACM international confer-</i>                  | Xin Yuan, Zhe Lin, Jason Kuen, Jianming Zhang, Yilin                        | 776 |
| 723 | <i>ence on Web search and data mining</i> , pages 163–172.                    | Wang, Michael Maire, Ajinkya Kale, and Baldo Fai-                           | 777 |
| 724 | Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018.                        | eta. 2021. <a href="#">Multimodal contrastive training for vi-</a>          | 778 |
| 725 | <a href="#">Representation learning with contrastive predictive</a>           | <a href="#">sual representation learning</a> . In <i>Proceedings of the</i> | 779 |
| 726 | <a href="#">coding</a> . <i>arXiv preprint arXiv:1807.03748</i> .             | <i>IEEE/CVF Conference on Computer Vision and Pat-</i>                      | 780 |
| 727 | Jeffrey Pennington, Richard Socher, and Christopher D                         | <i>tern Recognition</i> , pages 6995–7004.                                  | 781 |
| 728 | Manning. 2014. <a href="#">Glove: Global vectors for word rep-</a>            | Zhu Zhang and Balaji Varadarajan. 2006. <a href="#">Utility scor-</a>       | 782 |
| 729 | <a href="#">resentation</a> . In <i>Proceedings of the 2014 conference</i>    | <a href="#">ing of product reviews</a> . In <i>Proceedings of the 15th</i>  | 783 |
| 730 | <i>on empirical methods in natural language processing</i>                    | <i>ACM international conference on Information and</i>                      | 784 |
| 731 | <i>(EMNLP)</i> , pages 1532–1543.   | <i>knowledge management</i> , pages 51–57.                                  | 785 |
| 732 | Xiaoru Qu, Zhao Li, Jialin Wang, Zhipeng Zhang,                               | Feng Zhu and Xiaoquan Zhang. 2010. <a href="#">Impact of online</a>         | 786 |
| 733 | Pengcheng Zou, Junxiao Jiang, Jiaming Huang, Rong                             | <a href="#">consumer reviews on sales: The moderating role of</a>           | 787 |
| 734 | Xiao, Ji Zhang, and Jun Gao. 2020. <a href="#">Category-aware</a>             | <a href="#">product and consumer characteristics</a> . <i>Journal of</i>    | 788 |
| 735 | <a href="#">graph neural networks for improving e-commerce</a>                | <i>marketing</i> , 74(2):133–148.   | 789 |
| 736 | <a href="#">review helpfulness prediction</a> . In <i>Proceedings of the</i>  |   |     |
| 737 | <i>29th ACM International Conference on Information</i>                       |   |     |
| 738 | <i>&amp; Knowledge Management</i> , pages 2693–2700.                          |   |     |
| 739 | Shaoqing Ren, Kaiming He, Ross Girshick, and Jian                             |   |     |
| 740 | Sun. 2015. <a href="#">Faster r-cnn: Towards real-time object</a>             |   |     |
| 741 | <a href="#">detection with region proposal networks</a> . <i>Advances</i>     |   |     |
| 742 | <i>in neural information processing systems</i> , 28:91–99.                   |   |     |
| 743 | Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang.                             |   |     |
| 744 | 2019. How to fine-tune bert for text classification?                          |   |     |
| 745 | In <i>China National Conference on Chinese Computa-</i>                       |   |     |
| 746 | <i>tional Linguistics</i> , pages 194–206. Springer.                          |   |     |
| 747 | Jiliang Tang, Huiji Gao, Xia Hu, and Huan Liu. 2013.                          |   |     |
| 748 | <a href="#">Context-aware review helpfulness rating prediction</a> .          |   |     |
| 749 | In <i>Proceedings of the 7th ACM Conference on Rec-</i>                       |   |     |
| 750 | <i>ommender Systems</i> , pages 1–8.  |   |     |
| 751 | Oren Tsur and Ari Rappoport. 2009. <a href="#">Revrnk: A fully</a>            |   |     |
| 752 | <a href="#">unsupervised algorithm for selecting the most help-</a>           |   |     |
| 753 | <a href="#">ful book reviews</a> . In <i>Proceedings of the Interna-</i>      |   |     |
| 754 | <i>tional AAAI Conference on Web and Social Media</i> ,                       |   |     |
| 755 | volume 3.   |   |     |

## A Dataset Specification

Specifications of the two datasets are in Table 4 and 5 below.

| Amazon-MRHP (Products/Reviews) |               |              |               |
|--------------------------------|---------------|--------------|---------------|
| Category                       | Cloth. & Jew. | Elec.        | Home & Kitch. |
| Train                          | 12074/277308  | 10564/240505 | 14570/369518  |
| Dev                            | 3019/122148   | 2641/84402   | 3616/92707    |
| Test                           | 3966/87492    | 3327/79750   | 4529/111593   |

Table 4: Statistics of the Amazon-MRHP dataset.

| Lazada-MRHP (Products/Reviews) |               |            |               |
|--------------------------------|---------------|------------|---------------|
| Category                       | Cloth. & Jew. | Elec.      | Home & Kitch. |
| Train                          | 6596/104093   | 3848/41828 | 2939/36991    |
| Dev                            | 1649/26139    | 963/10565  | 736/9611      |
| Test                           | 2062/32274    | 1204/12661 | 920/12551     |

Table 5: Statistics of the Lazada-MRHP dataset.

## B Hyperparameter Search

The optimal hyperparameter settings are provided in Table 6 and 7.

| Glove Hyperparameters  |               |           |               |
|------------------------|---------------|-----------|---------------|
|                        | Cloth. & Jew. | Elec.     | Home & Kitch. |
| learning rate          | $1e^{-4}$     | $5e^{-5}$ | $1e^{-4}$     |
| text embedding dim     | 300           | 300       | 300           |
| text embedding dropout | 0.5           | 0.5       | 0.2           |
| image embedding dim    | 128           | 128       | 128           |
| LSTM hidden dim        | 128           | 128       | 128           |
| shared space hidden    | 64            | 64        | 64            |
| $\kappa$               | 0.25          | 0.1       | 0.1           |
| batch size             | 32            | 32        | 32            |

Table 6: Hyperparameters for all categories using glove-300d embeddings.

| Amazon-MRHP Hyperparameters |               |           |               |
|-----------------------------|---------------|-----------|---------------|
|                             | Cloth. & Jew. | Elec.     | Home & Kitch. |
| learning rate               | $2e^{-5}$     | $2e^{-5}$ | $2e^{-5}$     |
| text embedding dim          | 768           | 768       | 768           |
| text embedding dropout      | 0.5           | 0.5       | 0.5           |
| LSTM hidden dim             | 128           | 128       | 128           |
| image embedding dim         | 128           | 128       | 128           |
| shared space hidden         | 64            | 64        | 64            |
| $\kappa$                    | 0.3           | 0.25      | 0.25          |
| batch size                  | 32            | 32        | 32            |

Table 7: Hyperparameters for all categories using BERT as encoder

We use the same set of settings for text-only and multimodal modes for the same category dataset. The search space of these hyperparameters are:

learning rate in  $\{1e^{-4}, 2e^{-5}\}$ , text embedding dropout in  $\{0.2, 0.5\}$ ,  $\kappa$  in  $\{0.1, 0.25, 0.3, 0.5\}$ , shared space hidden dimension in  $\{64, 128\}$ . We train and test each dataset on a single Tesla V100 GPU. In BERT experiments, we use shared a BERT encoder for both product description and review text. To balance the computation cost and model performance, following Sun et al. (2019), we fine-tune the last four layers of the BERT encoder.

## C Language Tools

For coreference resolution, we use neuralcoref, an extension that can be placed on SpaCy processors. For BERT model, we use the huggingface transformers package to load.