

# E-BERT: Adapting BERT to E-commerce with Adaptive Hybrid Masking and Neighbor Product Reconstruction

Denghui Zhang<sup>1</sup>, Zixuan Yuan<sup>1</sup>, Yanchi Liu<sup>2</sup>, Zuohui Fu<sup>1</sup>, Fuzhen Zhuang<sup>3</sup>,  
Pengyang Wang<sup>4</sup>, Hui Xiong<sup>1</sup>

<sup>1</sup>Rutgers University, USA, {denghui.zhang, zy101, hxiong}@rutgers.edu

<sup>2</sup>NEC Laboratories America, Inc., USA, yanchi@nec-labs.com

<sup>3</sup>Institute of Computing Technology, Chinese Academy of Sciences, China

<sup>4</sup>University of Central Florida, USA

## Abstract

Pre-trained language models such as BERT have achieved great success in a broad range of natural language processing tasks. However, BERT cannot well support E-commerce related tasks due to the lack of two levels of domain knowledge, i.e., phrase-level and product-level. On one hand, many E-commerce tasks require accurate understanding of domain phrases, whereas such fine-grained *phrase-level* knowledge is not explicitly modeled by BERT’s training objective. On the other hand, *product-level* knowledge like product associations can enhance the language modeling of E-commerce, but they are not factual knowledge thus using them indiscriminately may introduce noise. To tackle the problem, we propose a unified pre-training framework, namely, E-BERT. Specifically, to preserve phrase-level knowledge, we introduce Adaptive Hybrid Masking, which allows the model to adaptively switch from learning preliminary word knowledge to learning complex phrases, based on the fitting progress of two modes. To utilize product-level knowledge, we introduce Neighbor Product Reconstruction, which trains E-BERT to predict a product’s associated neighbors with a denoising cross attention layer. Our investigation reveals promising results in four downstream tasks, i.e., review-based question answering, aspect extraction, aspect sentiment classification, and product classification.

## Introduction

Unsupervised pre-trained language models like BERT (Devlin et al. 2019) have greatly advanced the natural language processing research in recent years. However, these models are pre-trained on open-domain corpus and then fine-tuned for generic tasks, thus cannot well support domain-specific tasks. To this end, several domain-adaptive BERTs have been proposed recently, such as BioBERT (Lee et al. 2020), SciBERT (Beltagy, Lo, and Cohan 2019). They employ large-scale domain corpora to obtain language knowledge of specific domains, e.g., BioBERT uses 1M PubMed articles for pre-training. Despite this, they adopt the same architecture and pre-training approach as BERT does, neglecting the crucial domain knowledge which is beneficial for downstream tasks. Specifically, we find two levels of domain knowledge may not be effectively captured by BERT, i.e., *phrase-level* and *product-level* knowledge. Along this

<b>Review 1:</b> We love the <b>size of the screen</b> , although it is still light-weight and very easy to tote around.
<b>Review 2:</b> That included the extra <b>Sony Sonic Stage software</b> , the speakers and the subwoofer I got (that WAS worth the money), the <b>bluetooth mouse</b> for my supposedly bluetooth enabled computer, the <b>extended life battery</b> and the <b>docking port</b> . [...]
<b>BERT prediction:</b> 1: screen 2: software, bluetooth mouse, battery
<b>E-BERT prediction:</b> 1: size of the screen 2: Sony Sonic Stage software, bluetooth mouse, battery, docking port

Figure 1: An example of review aspect extraction, with answers marked in color. The vanilla BERT tends to make incomplete/wrong predictions, and miss aspects sometimes. But it gets improved after integrating phrase-level and product-level domain knowledge.

line, we incorporate these two-levels of domain knowledge to BERT for the E-commerce domain and perform evaluations on several related tasks.

First, many tasks in E-commerce involve understanding various *domain phrases*. However, such fine-grained *phrase-level* knowledge is not explicitly captured by BERT’s training objective, i.e., Masked Language Model (MLM). Specifically, MLM aims to predict individual masked words from incomplete input, thus being a *word-oriented* rather than *phrase-oriented* task. Although some subsequent work (Sun et al. 2020) proposes to mask phrases so that BERT can learn phrase-level knowledge, there are two major limitations: (i) they mask phrases that are simply obtained by chunking, may not be domain-specific; (ii) they discard word masking after using phrase masking, yet, we argue word-level knowledge is preliminary to phrase understanding. Figure 1 gives a motivating example from review Aspect Extraction. The task aims to extract entity aspects on which opinions have been expressed. Without phrase modeling, BERT tends to miss aspects or output incomplete aspects. It gets improved after effective phrase knowledge encoding.

On the other hand, there is rich semantic knowledge hidden in product associations, and we consider it as *product-level* knowledge. Existing models like BERT rely on co-occurrence to capture the semantics of words and phrases, which is inefficient and expensive. For instance, to teach

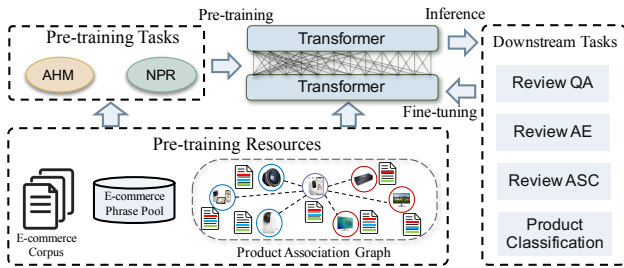


Figure 2: Overview of E-BERT.

Product	Description
Samsung Galaxy S10	OS: <b>Android</b> ; <b>5G network</b> ; <b>Dynamic AMOLED</b> ; ...
iPhone XS	<b>iOS</b> ; <b>4G signal</b> ; T-Mobile service; <b>OLED screen</b> ; ...

Figure 3: Two associated products and their descriptions.

the model that `Android` and `iOS` are semantically similar and correlated, a large number of co-occurrences of them are required in the corpus. Leveraging product associations to bridge the contents of `Samsung galaxy` and `iPhone`, we can easily enhance such semantic learning. However, it is challenging in practice because different fragments in two connected contents have different association confidence, without differentiating, it may introduce extra noise.

To enable pre-trained language models with the above two levels of domain knowledge, we propose a unified pre-training framework, E-BERT. As shown in Figure 2, we continue to use Transformer as the underlying architecture, and leverage a massive domain corpus, a high-quality E-commerce phrase pool, and a product association graph as our pre-training resources. To train E-BERT on them, we introduce two novel improvements as the pre-training tasks, i.e., Adaptive Hybrid Masking (AHM) and Neighbor Product Reconstruction (NPR):

(1) AHM extends MLM by introducing a new masking strategy. Specifically, it sets two different modes, i.e., word masking mode and phrase masking mode. The former randomly masks separate words while the latter masks domain phrases. Moreover, it can adaptively switch between the two modes based on feedback losses, enabling the model to capture word-level and phrase-level knowledge progressively.

(2) In NPR, we train E-BERT to reconstruct the neighbor products in the association graph given a central product, using its own content representation and a de-noising cross attention layer. The cross attention layer enables the model to pay more/less attention to different positions of the content according to their relevance. As a result, NPR transforms product-level knowledge into semantic knowledge without introducing too much noise.

To validate the effectiveness of the proposed approach, we fine-tune E-BERT on four downstream tasks, i.e., Review-based Question Answering (RQA), Aspect Extraction (AE), Aspect Sentiment Classification (ASC), and Product Classification. The experimental results show that E-BERT significantly outperforms BERT and several following work on these domain-specific tasks, by taking full advantage of the phrase-level and product-level knowledge.

Table 1: High-quality phrases of 6 product categories.

Category	Top-rated phrases
Automotive	jumper cables, cometic gasket, angel eyes, drink holder, static cling
Clothing, Shoes and Jewelry	high waisted jean, nike classic, removable tie, elegant victorian, vintage grey
Electronics	ipads tablets, SDHC memory card, memory bandwidth, auto switching
Office Products	decorative paper, heavy duty rubber, mailing labels, hybrid notebinder
Sports and Outdoors	basketball backboard, table tennis paddle, string oscillation, fishing tackles
Toys and Games	hulk hogan, augmented reality, teacup piggies, beam sabers, naruto uzumaki

## Methodology

In this section, we first present the pre-training resources used in E-BERT. Then, we provide an in-depth introduction about our improvements in pre-training, i.e., Adaptive Hybrid Masking and Neighbor Product Reconstruction.

### Pre-training Resources

**E-commerce Corpus** We extract millions of product titles, descriptions, and reviews from the Amazon Dataset<sup>1</sup> (Ni, Li, and McAuley 2019) to build this corpus. We divide the corpus into two sub-corpus, i.e., product corpus and review corpus. In the first corpus, each line corresponds to a product title and its description, while in the second, it corresponds to a user comment on a specific product. The corpus serves as the foundation for E-BERT to learn preliminary language knowledge.

**E-commerce Phrase Pool** To incorporate domain phrase knowledge into E-BERT, we extract plenty of domain phrases from the above corpus and build an E-commerce phrase pool in advance. Considering phrase quality, we adopt AutoPhrase<sup>2</sup>, a high efficient phrase mining method (Shang et al. 2018), can generate a quality score for each phrase based on corpus-level statistics like *popularity*, *concordance*, *informativeness*, and *completeness*. In total, we extract more than one million initial phrases. Then, we filter out phrases where  $score < 0.5$  to get quality phrases and store them in the pool. We also attach the score of each phrase in the pool, which is used for phrase sampling in AHM. Table 1 shows the top-ranked phrases from six product categories. Compared with existing work using chunking to select noun phrases for masking, our phrase pool is more diversified, fine-grained, and domain-specific. Discussion about the effects of choosing different phrase sets is shown in the experiment section.

**Product Association Graph** To enable E-BERT with product-level knowledge, we build Product Association Graph in advance. Specifically, products in our corpus are represented as nodes, and associations among products are represented as undirect edges. To obtain product associations, we extract product pairs such as substitutable and

<sup>1</sup><https://nijianmo.github.io/amazon/index.html>

<sup>2</sup><https://github.com/shangjingbo1226/AutoPhrase>

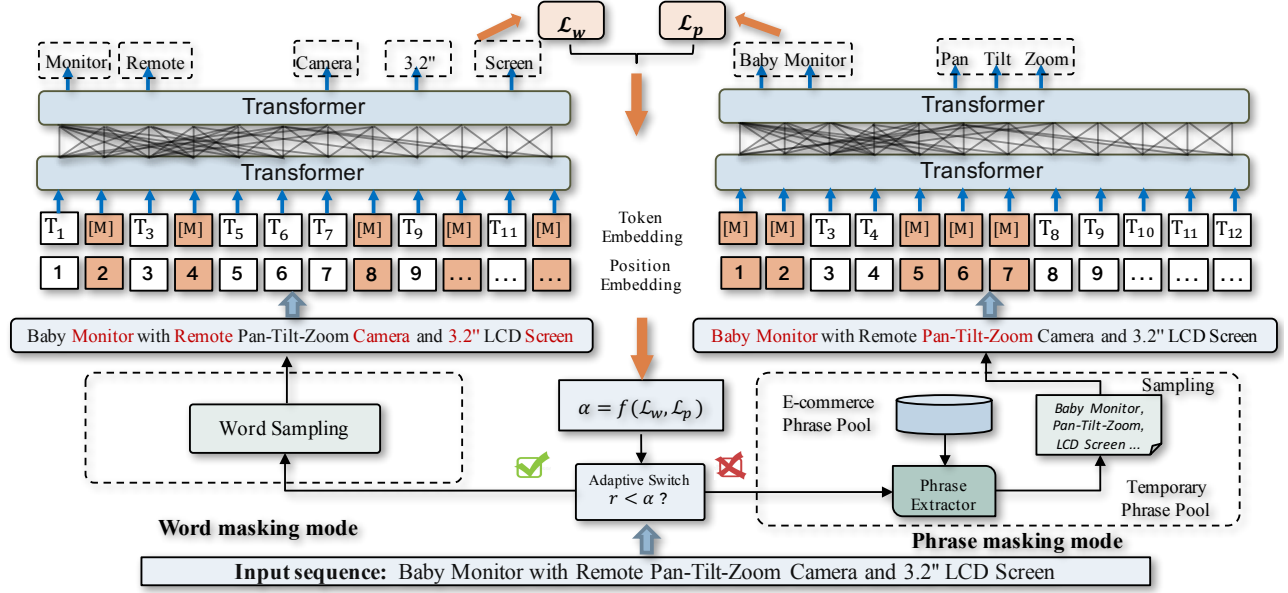


Figure 4: The illustration of Adaptive Hybrid Masking. Based on the feedback losses, it adaptively switches between two masking modes, enabling the model to learn word-level and phrase-level knowledge in a progressive manner.

complementary from the Amazon dataset, using a heuristic method based on consumer shopping statistics (McAuley, Pandey, and Leskovec 2015).

### Adaptive Hybrid Masking

To encode phrase-level knowledge effectively, we introduce a new masking strategy, namely, Adaptive Hybrid Masking (AHM), which is easy to implement by extending MLM. In AHM, we set two masking modes, i.e., *word masking* and *phrase masking* respectively. The former masks word units while the latter masks domain phrase units, resulting in inconsistent difficulty levels of reconstructing the masked tokens. To this end, we adaptively switch from predicting masked words to predicting masked phrases, enabling E-BERT to capture word-level and phrase-level knowledge in a progressive manner.

**Word Masking Mode** In this mode, we select random words from input sequence iteratively until obtain 15% tokens for masking. This scheme learns preliminary word-level semantics, which is essential for phrase understanding.

**Phrase Masking Mode** In phrase masking mode, we randomly mask consecutive tokens that can form *quality domain phrases*. Specifically, given an input sequence of tokens  $X = \{x_i\}_{i=1}^n$ , we first detect all the E-commerce phrases  $\{p_i\}_{i=1}^m$  in  $X$ , leveraging the E-commerce phrase pool  $\mathcal{P}_E$  along with a rule-based phrase matcher<sup>3</sup>. Then, we create a temporary phrase pool  $\mathcal{P}_X$  consisting of the detected phrases  $\{p_i\}_{i=1}^m$ . Some input sequences may contain too few domain phrases, therefore, to ensure we have enough and diverse phrases, we extend  $\mathcal{P}_X$  with noun phrases. That is, we extract all the noun phrases  $\{n_i\}_{i=1}^l$  in  $X$  using constituency

parsing<sup>4</sup>, and abandon the ones that have an intersection with the domain phrases  $\{p_i\}_{i=1}^m$ . Then, we add the rest “clean” noun phrases in to  $\mathcal{P}_X$ . Based on the extended  $\mathcal{P}_X$ , we sample phrases iteratively until obtain approximately 15% tokens for masking. The probability of selecting each phrase is set as the softmax of its quality score:

$$p(p_i) = \frac{\exp(s[p_i])}{\sum_{p_j \in \mathcal{P}_X} \exp(s[p_j])}, \quad (1)$$

where  $s[p_i]$  denotes the score of phrase  $p_i$ . For the supplemental phrases, it is assigned with the lowest score in  $\mathcal{P}_E$ . Phrases with higher scores are usually more E-commerce related, thus the quality-based sampling impels our model to pay more attention to unique domain phrases.

**Adaptive Switching** When learning a new language, people usually start with individual words (the vocabulary), and gradually turn to study more complex phrases and expressions. Inspired by this, we start pre-training with word masking mode, and set a time-varying parameter  $\alpha$  to adaptively switch to phrase mode.

In detail, at each iteration ( $t^{th}$ ), we calculate a “fitting index” for both modes to track their fitting progress, i.e.,  $\eta_w^t$  and  $\eta_p^t$ . The larger  $\eta_w^t$  ( $\eta_p^t$ ) is, the less sufficient the model is trained on the word (phrase) mode. Next, we calculate  $\gamma^t$ , representing the relative importance of the word mode, and rescale it to [0,1] via a non-linear unit ( $\tanh$ ) to get the probability of choosing word mode at the next iteration ( $t+1^{th}$ ), i.e.,  $\alpha^{t+1}$ :

$$\eta_w^t = \frac{\Delta_w^{t,t-1}}{\Delta_w^{t,1}} = \frac{[\mathcal{L}_w^{t-1} - \mathcal{L}_w^t]_+}{\mathcal{L}_w^1 - \mathcal{L}_w^t}, \quad (2)$$

<sup>3</sup><https://spacy.io/usage/examples#phrase-matcher>

<sup>4</sup><https://spacy.io/universe/project/self-attentive-parser>

$$\eta_p^t = \frac{\Delta_p^{t,t-1}}{\Delta_p^{t,1}} = \frac{[\mathcal{L}_p^{t-1} - \mathcal{L}_p^t]_+}{\mathcal{L}_p^1 - \mathcal{L}_p^t}, \quad (3)$$

$$\gamma^t = \frac{\eta_w^{t+1}}{\eta_p^{t+1}}, \quad (4)$$

$$\alpha^{t+1} = \tanh(\gamma^t). \quad (5)$$

where  $\Delta_w^{t,t-1}$  denotes the loss reduction of word mode between current and last iteration.  $\Delta_w^{t,1}$  denotes the total loss reduction.  $\mathcal{L}_w^t$  denotes the loss and will only be updated if word mode is selected at the  $t$ -th iteration.  $\Delta_p^{t,t-1}$ ,  $\Delta_p^{t,1}$ , denotes the counterparts in phrase mode.  $[x]_+$  is equivalent to  $\max(x, 0)$ . When  $\eta_w^{t+1} \gg \eta_p^{t+1}$ ,  $\alpha^{t+1} \approx 1$ , and the word mode becomes dominating, vice versa. In other words,  $\alpha$  controls the model to switch to the weaker mode adaptively.

Figure 4 presents an overall illustration of AHM. At each iteration, we first generate a random number  $r \in [0, 1]$ , then we select word masking mode if  $r < \alpha^t$  and select phrase mode otherwise. To be noted, for the first  $t \leq T_1$  iterations, we set  $\alpha^t = \alpha_0$  ( $\alpha_0 > 0.5$ ) to make word mode more likely to be selected. After this initial stage,  $\eta_p$  gets larger than  $\eta_w$  and  $\alpha$  decreases, consequently, the probability of selecting phrase mode gets larger. Until the end, it will switch between the two modes adaptively based on their losses, balancing word-level and phrase-level learning.

**Reconstructing Masked Tokens** For both modes after selecting tokens to be masked, following BERT to mitigate the mismatch between pre-training and fine-tuning, we replace the selected tokens with (1) the [MASK] token 80% of the time, (2) a random token 10% of the time, (3) the original token 10% of the time. Next, we predict each masked token by feeding their output embedding to a shared softmax layer (take word masking mode as example), i.e.,

$$p(X_m^t | X_{\mathcal{W}_{X^t}}^t) = \frac{\exp(\mathbf{W}_m^\top [\mathbf{E-BERT}(X_{\mathcal{W}_{X^t}}^t)]_m)}{\sum_{k \in \mathcal{V}} \exp(\mathbf{W}_k^\top [\mathbf{E-BERT}(X_{\mathcal{W}_{X^t}}^t)]_k)}, \quad (6)$$

where  $\mathbf{W}$  represents the parameters of softmax layer.  $\mathcal{V}$  denotes the vocabulary.  $X^t$  denotes the input sequence.  $\mathcal{W}_{X^t}$  denotes the set of masked tokens in word masking mode,  $\setminus$  denotes set minus,  $X_{\setminus \mathcal{W}_{X^t}}^t$  denotes the modified input where  $\mathcal{W}_{X^t}$  are masked.  $X_m^t$  denotes the masked token to be predicted,  $m \in \mathcal{W}_{X^t}$ . The overall loss function of AHM is the combined cross entropy of the two masking modes, i.e.,

$$\mathcal{L}_{\text{AHM}} = -\frac{1}{|\mathcal{D}|} \sum_{X^t \in \mathcal{D}} \alpha^t \log \prod_{m \in \mathcal{W}_{X^t}} p(X_m^t | X_{\setminus \mathcal{W}_{X^t}}^t) + (1 - \alpha^t) \log \prod_{m \in \mathcal{P}_{X^t}} p(X_m^t | X_{\setminus \mathcal{P}_{X^t}}^t), \quad (7)$$

where  $\mathcal{D}$  represents the training corpus.  $p(X_m^t | X_{\setminus \mathcal{P}_{X^t}}^t)$  denotes the prediction in phrase masking mode, calculated the same way as word masking mode.  $\mathcal{P}_{X^t}$  denotes the set of masked tokens in phrase mode.

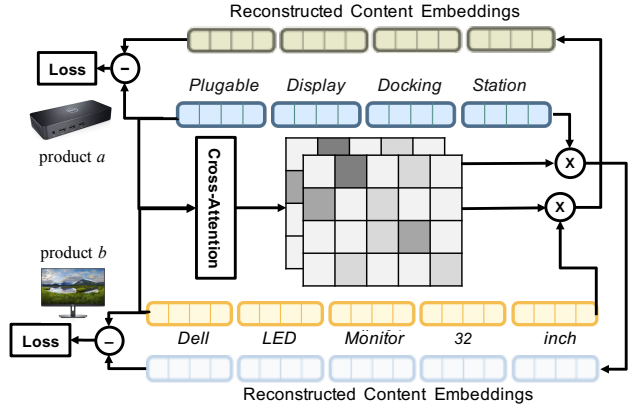


Figure 5: Illustration of Neighbor Product Reconstruction.

## Neighbor Product Reconstruction

In this task, we train E-BERT to reconstruct the neighbor (associated) product's content using the central product's content, and thereby, transform the hidden semantic knowledge into the weights of the model.

As illustrated in Figure 5, we first sample associated products from the product association graph and put their content embeddings into a pair (shown in the middle of the figure). The cross attention layer is then used to learn a correlation matrix, indicating word-level correlations between two products. Next, we multiply the correlation matrix with the central product's content embeddings to generate a set of reconstructed content embeddings for the neighbor product. They are optimized to resemble the real ones via the content reconstructing loss. Considering the central product can also be viewed as neighbor to the reconstructed product, we perform the reconstructing task in both directions.

**Content Embeddings** Given the product pair  $(a, b)$ , we feed their contents (title and description) into E-BERT respectively to get their content embeddings, i.e.,

$$\begin{aligned} \{\mathbf{w}_i\}_{i=1}^n &= \mathbf{E-BERT}(\{\mathbf{a}_i\}_{i=1}^n), \\ \{\mathbf{o}_i\}_{i=1}^n &= \mathbf{E-BERT}(\{\mathbf{b}_i\}_{i=1}^n). \end{aligned} \quad (8)$$

**Cross Attention Layer** We adopt the cross attention layer to generate two correlation matrices, i.e.,

$$\alpha_{ij} = \frac{\exp(\mathbf{w}_i \mathbf{o}_j)}{\sum_{j'} \exp(\mathbf{w}_i \mathbf{o}_{j'})}, \quad \beta_{ji} = \frac{\exp(\mathbf{w}_i \mathbf{o}_j)}{\sum_{i'} \exp(\mathbf{w}_{i'} \mathbf{o}_j)}, \quad (9)$$

where  $w_i$  indicates the  $i^{th}$  word in the product  $a$ 's content and  $o_j$  represents the  $j^{th}$  word in  $b$ . The cross attention weight  $\alpha_{ij}$  and  $\beta_{ji}$  both indicates the correlation between  $w_i$  and  $o_j$ , but using different normalizers.

**Reconstructed Embeddings** Using the attention weights, we compute a weighted average of the real content embeddings, i.e.,

$$\mathbf{w}'_i = \sum_j \alpha_{ij} \mathbf{o}_j, \quad \mathbf{o}'_j = \sum_i \beta_{ji} \mathbf{w}_i. \quad (10)$$

where  $\mathbf{w}'_i$  and  $\mathbf{o}'_j$  are the reconstructed embeddings, they are subsequently optimized to resemble the original content embeddings  $\mathbf{w}_i$  and  $\mathbf{o}_j$ . The negative effect of noise is



minimized by the cross attention as it automatically assigns smaller weights to irrelevant contents.

**Reconstructing Loss** We define the content reconstructing loss of a product pair as the Euclidean distance between their real content embeddings to the reconstructed embeddings, i.e.,

$$\langle a, b \rangle = \sum_i \|w_i - w'_i\|_2^2 + \sum_j \|o_j - o'_j\|_2^2 \quad (11)$$

We use a triplet loss as the final loss to pull relevant product-product pairs close while pushing irrelevant ones apart:

$$\mathcal{L}_{\text{NPR}} = \max(0, 1 + \langle a, b \rangle - \langle a, b^- \rangle), \quad (12)$$

where  $b^-$  is a randomly sampled negative product that is not related to  $a$ .

To be noted, we only train NPR on the product corpus where the input is formatted as content pairs  $\langle \text{content}(a), \text{content}(b) \rangle$ . We do not train NPR on the review corpus since it consists of user feedbacks, can not reflect product semantics accurately.

## Experiments

In this section, we conduct extensive experiments to answer the following research questions:

- What is the performance gain of the E-commerce corpus for each downstream task, with respect to the state-of-the-art performance?
- What is the overall performance gain of our pre-training framework incorporating two levels domain knowledge?
- What is the performance gain of each component (i.e., AHM and NPR) in E-BERT?

### Baselines

In this paper, we compare E-BERT to the following baseline methods:

- **BERT-Raw**: The vanilla BERT which is pre-trained on large-scale open-domain corpus<sup>5</sup>. We use this baseline to answer the first question.
- **BERT**: The vanilla BERT which is further post-trained on our E-commerce corpus. We compare with this baseline to answer the second question.
- **BERT-NP**: The vanilla BERT which is post-trained on our E-commerce corpus, but uses a different masking strategy, i.e., masks noun phrases instead of words. We use this to validate the effect of our domain phrase pool.
- **SpanBERT**: An variant of BERT which masks spans of tokens instead of separate tokens. We compare with it to further validate the effect of the phrase masking scheme.

For ablation studies, we further compare with the following internal baselines:

- **E-BERT-DP**: The reduced E-BERT which only uses the phrase masking mode, without word-level masking.

<sup>5</sup>We use the pre-trained model released by Huggingface.

- **E-BERT-AHM**: The reduced E-BERT which adopts AHM to adaptively change the masking mode, but not utilizes NPR to encode product-level knowledge.
- **E-BERT**: The full E-BERT, utilizing both AHM and NPR to encode two levels of domain knowledge.

### Pre-training Dataset

The dataset contains four pre-training resources:

- **Product Corpus** It contains a total of 5, 436, 547 product titles and descriptions with a size of 1.4 GB.
- **Review Corpus** It contains a total of 9, 636, 112 million product reviews with a size of 2.3 GB.
- **E-commerce Phrase Pool** It consists of 536, 332 high quality E-commerce phrases.
- **Product Association Graph** It consists of 2, 125, 352 products and 3, 484, 325 product associations.

**Phrase Overlap** Figure 6 presents the overlap between the E-commerce domain phrases and the ordinary noun phrases extracted from our corpus, divided by 5 product categories. Each entry indicates the proportion of domain phrases that also occurred in the noun phrase set. We can observe that the overlap ratio is low even for the same category, indicating our phrase pool contains more unique phrase-level knowledge.

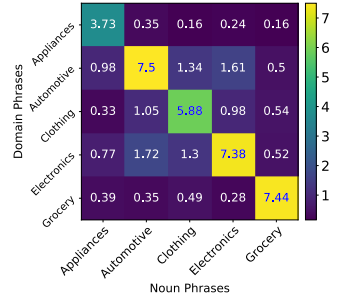


Figure 6: The overlap ratio (%) between the E-commerce phrases and ordinary noun phrases.

### Pre-training Details

All the baselines and E-BERT is initialized with the weights of the pre-trained BERT (the `bert-base-uncased` version by Huggingface, with 12 layers, 768 hidden dimensions, 12 heads, 110M parameters). We post-train all the baselines except BERT-Raw on the E-commerce corpus for 10 epochs, with batch size 32 and learning rate  $1e-5$ . For E-BERT, we adopt Continual Multi-task Learning (Sun et al. 2020) to combine AHM and NPR. To be specific, we first train AHM alone on the entire corpus for 5 epochs with the same batch size and learning rate. Then, we train AHM and NPR jointly on the product corpus for another 5 epochs. The only hyperparameter in AHM,  $T_1$ , is set to be 1 epoch.

### Downstream Tasks

**Review-based Question Answering** Given a question  $q = \{q_i\}_{i=1}^m$  and a related review  $r = \{r_i\}_{i=1}^n$ , it aims to find the span  $s = \{r_i\}_{i=s}^e$  from  $r$  that can answer  $q$ . To fine-tune this task, we adopt the standard BERT approach (Devlin et al. 2019) for span-based QA, which maximizes the sum of the log-likelihoods of the correct start and end positions.

Table 2: Results of baselines and our model on E-commerce downstream tasks (%).

Models\Tasks	Review QA				Product Classification			Review AE			Review ASC	
	<i>P.</i>	<i>R.</i>	<i>F1</i>	<i>EM</i>	<i>Acc.</i>	<i>MiF1</i>	<i>MaF1</i>	<i>P.</i>	<i>R.</i>	<i>F1</i>	<i>Acc.</i>	<i>MaF1</i>
Pre-trained on Wikipedia + BookCorpus by Huggingface.												
BERT-Raw	58.91	62.58	60.69	40.22	66.54	81.90	78.82	83.15	84.66	83.90	86.01	62.87
Further post-trained on E-commerce corpus by us.												
BERT	60.28	62.25	61.25	41.23	69.12	82.38	80.66	84.33	84.09	84.81	86.40	64.96
BERT-NP	61.39	64.57	62.94	43.35	70.28	81.72	81.34	85.23	85.71	86.11	85.79	63.21
SpanBERT	62.52	64.77	63.63	43.94	71.59	81.51	81.50	85.67	86.22	86.23	86.76	65.13
E-BERT-DP	63.76	67.02	65.77	44.63	75.07	85.84	86.28	86.80	89.47	88.11	87.84	69.02
E-BERT-AHM	65.18	68.30	66.18	<b>45.56</b>	76.61	86.35	87.32	<b>87.42</b>	<b>90.55</b>	<b>88.96</b>	<b>89.17</b>	<b>70.35</b>
E-BERT	<b>66.71</b>	<b>70.13</b>	<b>68.77</b>	45.40	<b>78.74</b>	<b>90.37</b>	<b>90.94</b>	87.35	89.61	88.42	88.43	69.32

**Review Aspect Extraction** Given a review  $r = \{r_i\}_{i=1}^n$ , the task aims to find aspects that reviewers have expressed opinions on. It is typically formalized as a sequence labeling task (Xu et al. 2019), in which each token is classified as one of  $\{B, I, O\}$ , and tokens between  $B$  and  $I$  are considered the correct aspect. We apply a dense layer and softmax layer on top of each output embedding to fine-tune.

**Review Aspect Sentiment Classification** Given an aspect  $a = \{a_i\}_{i=1}^l$  and the review sentence  $r = \{r_i\}_{i=1}^n$  where  $a$  extracted from, this task aims to classify the sentiment polarity (positive, negative, or neutral) expressed on aspect  $a$ . For fine-tuning, both  $a$  and  $r$  are input into E-BERT, and we use the [CLS] token along with a dense layer and softmax layer to predict the polarity. Training loss is the cross entropy on the polarities.

**Product Classification** Given a product title  $x = \{x_i\}_{i=1}^n$ , it aims to classify  $x$  using a predefined category hierarchy  $\mathcal{H}$ . Each product may belong to multiple categories, thus making it a multi-label classification problem. We use the [CLS] token along with a dense layer and softmax layer to perform prediction.

**Evaluation Datasets** For review QA, we evaluate on a newly released Amazon QA dataset (Miller et al. 2020), which consists of 8,967 samples. We use the laptop dataset of SemEval 2014 Task 4 (Pontiki et al. 2016) for both review AE and review ASC tasks, which contains 3,845 review sentences, 3,012 annotated aspects and the sentiment polarities on them. For product classification, we create an evaluation dataset by extracting Amazon product metadata, consisting of 10,039 product titles and 133 categories. For all the datasets, we divide them into training/validation/testing set with the ratio of 7:1:2.

**Fine-tuning Details** In each task, we adopt the standard architecture for each BERT variant. We choose the learning rate and epochs from  $\{5e-6, 1e-5, 2e-5, 5e-5\}$  and  $\{2, 3, 4, 5\}$  respectively. For each task and BERT variant, we pick the best learning rate and number of epochs on the development set and report the corresponding test results. We found the setting that works best across most tasks and models is 2 or 4 epochs and a learning rate of  $2e-5$ .

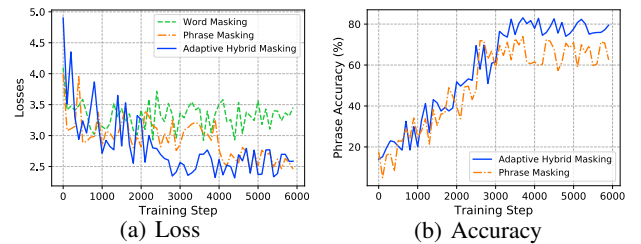


Figure 7: The convergence of different masking schemes.

**Evaluation Metrics** For review QA, we adopt the standard evaluation script from SQuAD 1.1 to report Precision, Recall, F1 scores, and Exact Match (EM). To evaluate review AE, we report Precision (P), Recall (R), and F1 score. For review ASC, we report Macro-F1 and Accuracy. Lastly, we adopt Accuracy (Acc), Micro-F1 (MiF1), and Macro-F1 (MaF1) to evaluate product classification.

## Result Analysis

Table 2 presents the results of all the baselines and E-BERT on the four tasks. First, we can see that BERT outperforms BERT-Raw on all the tasks, verifying that the E-commerce corpus can largely improve the performance on related tasks. Compared with BERT, BERT-NP and SpanBERT achieves further improvements in review QA and review AE, indicating that phrase-level knowledge is quite helpful in these extractive tasks. Comparing E-BERT-DP with BERT-NP and SpanBERT, we prove that our E-commerce phrase pool can provide more quality phrase knowledge for the downstream tasks. E-BERT outperforms all the baselines by a large margin in terms of all metrics, verifying the overall superiority of our pre-training framework in E-commerce tasks. To examine the effectiveness of product-level domain knowledge and each component of E-BERT, we present more discussions in ablation studies.

## Ablation Studies

As shown in the bottom of Table 2, E-BERT-AHM consistently outperforms E-BERT-DP in four tasks, proving that our adaptive hybrid masking strategy can utilize phrase-

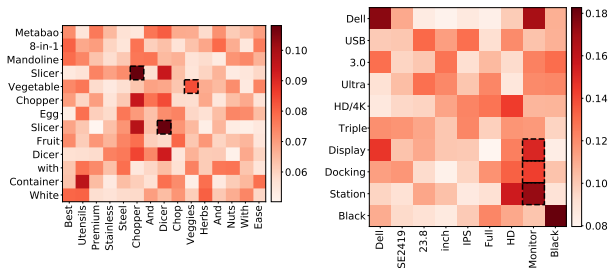


Figure 8: Visualizing the cross attention.

level knowledge in a more sufficient way. Besides, as shown in Figure 7, our masking method has a better convergence rate in terms of loss and phrase reconstruction accuracy. Compared with E-BERT-DP and E-BERT-AHM, E-BERT further encodes product-level knowledge via NPR, and it achieves significant improvements in product classification. We assume this is because there are strong category correlations between associated products, by utilizing product association knowledge, E-BERT enhances feature sharing among different instances. While the performance of review QA is also boosted, review AE and review ASC even deteriorates slightly after using NPR, indicating product-level knowledge has no significant effect on the task of review aspect analysis.

### Cross Attention Probing

Figure 8 presents the cross attention visualization of two pairs of product titles. In the first example, two similar products Mandoline Slicer and Steel Chopper are connected using the learned attention weights. The darker color indicates stronger correlations. It can be seen the cross attention automatically learn to align correlated words in two product contents, e.g., (Slicer, Chopper), (Slicer, Dicer) and (Vegetables, Veggies). In the second example, two complementary products Docking station and Dell Monitor are connected. Similarly, correlated contents such as (Display, Monitor) and (Docking station, Monitor) are aligned automatically.

### Related Work

**Pre-trained Language Models** Recent years have witnessed the great success of Pre-trained Language Models (PLMs) (Devlin et al. 2019; Peters et al. 2018; Radford et al. 2018) on a broad range of NLP tasks. Compared with traditional word embedding models (Gupta and Manning 2015), PLMs learn to represent words based on the entire input context to deal with word polysemy, thus captures semantics more accurately. Following PLMs, many endeavors have been made for further optimization. SpanBERT (Joshi et al. 2020) proposes to reconstruct randomly masked spans instead of single words. However, the span consists of random continuous words and may not form phrases, thus fails to capture phrase-level knowledge accurately. ERNIE-1.0 (Sun et al. 2019) integrates phrase-level masking and entity-level masking into BERT, which is closely related to our masking

strategy. Unlike them using simple chunking tools to get ordinary phrases, we build a high-quality E-commerce phrase pool and only mask domain phrases. Besides, we combine word-masking and phrase-masking coherently with Adaptive Hybrid Masking, accelerating the convergence without affecting performance. Due to space limit, we refer readers to the references for more work along this line (Liu et al. 2019; Lan et al. 2019; Yang et al. 2019; Brown et al. 2020; Sun et al. 2020).

**Domain-adaptive PLMs** To adapt PLMs to specific domains, several domain-adaptive BERTs have been proposed recently. BioBERT (Lee et al. 2020) and SciBERT (Beltagy, Lo, and Cohan 2019) train BERT on large-scale biomedical and scientific corpus respectively to get a pre-trained language model for biomedical and scientific NLP tasks. BERT-PT (Xu et al. 2019) propose to post-train BERT on a review corpus and obtains better performance on the task of review reading comprehension. Gururangan et al. propose to continue pre-training on domain corpus as well as task corpus and obtains more performance gains. More work along this line can be referred to (Rietzler et al. 2020; Ma et al. 2019; Jin et al. 2019; Huang, Altosaar, and Ranganath 2019). These work only leverages domain corpus for pre-training, without considering special domain knowledge like the product association graph.

**Knowledge Enhanced PLMs** Recently, to enable PLMs with world knowledge, several attempts (Wang et al. 2019; Peters et al. 2019; Zhang et al. 2019; Liu et al. 2020; Wang et al. 2020) have been made to inject knowledge into BERT leveraging Knowledge Graphs (KGs). Most of these work adopts the “BERT+entity linking” paradigm, whereas, it is not suitable for E-commerce corpus due to the lack of quality entity linkers as well as KGs in this domain. Instead, we consider utilizing the product association knowledge which is coarse-grained and may introduce noise. In E-BERT, through Neighbor Product Reconstruction and the de-noising cross attention layer, the meaning of each word in a product content is expanded to those of associated products, greatly enriches the semantic learning.

### Conclusions

In this paper, we proposed a domain-enhanced BERT for E-commerce, namely, E-BERT. We leveraged two levels of domain knowledge, i.e., phrase-level and product-level, to boost performance on related tasks. Despite the challenge of modeling phrase knowledge and reducing noise in product knowledge, we provided two technical improvements, i.e., AHM and NPR. Our investigation revealed promising results on four downstream tasks. Incorporating phrase knowledge via AHM can improve the performance significantly on all the investigated tasks. Utilizing the product-level knowledge via NPR further boosts the performance on product classification and review QA.

### References

Beltagy, I.; Lo, K.; and Cohan, A. 2019. SciBERT: A Pre-trained Language Model for Scientific Text. In *Proceed-*

- ings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), 3606–3611.
- Brown, T. B.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186.
- Gupta, S.; and Manning, C. D. 2015. Distributed representations of words to guide bootstrapped entity classifiers. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 1215–1220.
- Gururangan, S.; Marasović, A.; Swayamdipta, S.; Lo, K.; Beltagy, I.; Downey, D.; and Smith, N. A. 2020. Don’t Stop Pretraining: Adapt Language Models to Domains and Tasks. *arXiv preprint arXiv:2004.10964*.
- Huang, K.; Altosaar, J.; and Ranganath, R. 2019. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *arXiv preprint arXiv:1904.05342*.
- Jin, Q.; Dhingra, B.; Liu, Z.; Cohen, W.; and Lu, X. 2019. PubMedQA: A Dataset for Biomedical Research Question Answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2567–2577.
- Joshi, M.; Chen, D.; Liu, Y.; Weld, D. S.; Zettlemoyer, L.; and Levy, O. 2020. Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics* 8: 64–77.
- Lan, Z.; Chen, M.; Goodman, S.; Gimpel, K.; Sharma, P.; and Soricut, R. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.
- Lee, J.; Yoon, W.; Kim, S.; Kim, D.; Kim, S.; So, C. H.; and Kang, J. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 36(4): 1234–1240.
- Liu, W.; Zhou, P.; Zhao, Z.; Wang, Z.; Ju, Q.; Deng, H.; and Wang, P. 2020. K-BERT: Enabling Language Representation with Knowledge Graph. In *AAAI*, 2901–2908.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Ma, X.; Xu, P.; Wang, Z.; Nallapati, R.; and Xiang, B. 2019. Domain Adaptation with BERT-based Domain Classification and Data Selection. In *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, 76–83.
- McAuley, J.; Pandey, R.; and Leskovec, J. 2015. Inferring networks of substitutable and complementary products. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, 785–794.
- Miller, J.; Krauth, K.; Recht, B.; and Schmidt, L. 2020. The Effect of Natural Distribution Shift on Question Answering Models. *arXiv preprint arXiv:2004.14444*.
- Ni, J.; Li, J.; and McAuley, J. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 188–197.
- Peters, M.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; and Zettlemoyer, L. 2018. Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2227–2237.
- Peters, M. E.; Neumann, M.; Logan, R.; Schwartz, R.; Joshi, V.; Singh, S.; and Smith, N. A. 2019. Knowledge Enhanced Contextual Word Representations. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 43–54.
- Pontiki, M.; Galanis, D.; Papageorgiou, H.; Androutsopoulos, I.; Manandhar, S.; Al-Smadi, M.; Al-Ayyoub, M.; Zhao, Y.; Qin, B.; De Clercq, O.; et al. 2016. Semeval-2016 task 5: Aspect based sentiment analysis. In *10th International Workshop on Semantic Evaluation (SemEval 2016)*.
- Radford, A.; Narasimhan, K.; Salimans, T.; and Sutskever, I. 2018. Improving language understanding by generative pre-training.
- Rietzler, A.; Stabinger, S.; Opitz, P.; and Engl, S. 2020. Adapt or Get Left Behind: Domain Adaptation through BERT Language Model Finetuning for Aspect-Target Sentiment Classification. In *Proceedings of The 12th Language Resources and Evaluation Conference*, 4933–4941.
- Shang, J.; Liu, J.; Jiang, M.; Ren, X.; Voss, C. R.; and Han, J. 2018. Automated phrase mining from massive text corpora. *IEEE Transactions on Knowledge and Data Engineering* 30(10): 1825–1837.
- Sun, Y.; Wang, S.; Li, Y.; Feng, S.; Chen, X.; Zhang, H.; Tian, X.; Zhu, D.; Tian, H.; and Wu, H. 2019. Ernie: Enhanced representation through knowledge integration. *arXiv preprint arXiv:1904.09223*.
- Sun, Y.; Wang, S.; Li, Y.-K.; Feng, S.; Tian, H.; Wu, H.; and Wang, H. 2020. ERNIE 2.0: A Continual Pre-Training Framework for Language Understanding. In *AAAI*, 8968–8975.



Wang, R.; Tang, D.; Duan, N.; Wei, Z.; Huang, X.; Cao, C.; Jiang, D.; Zhou, M.; et al. 2020. K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808* .

Wang, X.; Gao, T.; Zhu, Z.; Liu, Z.; Li, J.; and Tang, J. 2019. KEPLER: A unified model for knowledge embedding and pre-trained language representation. *arXiv preprint arXiv:1911.06136* .

Xu, H.; Liu, B.; Shu, L.; and Philip, S. Y. 2019. BERT Post-Training for Review Reading Comprehension and Aspect-based Sentiment Analysis. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2324–2335.

Yang, Z.; Dai, Z.; Yang, Y.; Carbonell, J.; Salakhutdinov, R. R.; and Le, Q. V. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, 5753–5763.

Zhang, Z.; Han, X.; Liu, Z.; Jiang, X.; Sun, M.; and Liu, Q. 2019. ERNIE: Enhanced language representation with informative entities. *arXiv preprint arXiv:1905.07129* .