

Descriptions from the Customers: Comparative Analysis of Review-based Product Description Generation Methods

SLAVA NOVGORODOV, eBay Research, Israel

IDO GUY, eBay Research, Israel

GUY ELAD, Technion, Israel Institute of Technology, Israel

KIRA RADINSKY, Technion, Israel Institute of Technology, Israel

Product descriptions play an important role in the e-commerce ecosystem. Yet, on leading e-commerce websites product descriptions are often lacking or missing. In this work, we suggest to overcome these issues by generating product descriptions from user reviews. We identify the set of candidates using a supervised approach that extracts review sentences in their original form, diversifies them, and selects the top candidates. We present extensive analyses of the generated descriptions, including a comparison to the original descriptions and examination of review coverage. We also perform an A/B test that demonstrates the impact of presenting our descriptions on user traffic.

CCS Concepts: • **Information systems** → **Electronic commerce**; **Online shopping**; • **Applied computing** → **Electronic commerce**; • **Computing methodologies** → *Natural language generation*; *Multi-task learning*.

Additional Key Words and Phrases: Deep multi-task learning; electronic commerce; language generation; user-generated content.

ACM Reference Format:

Slava Novgorodov, Ido Guy, Guy Elad, and Kira Radinsky. 2020. Descriptions from the Customers: Comparative Analysis of Review-based Product Description Generation Methods. *ACM Trans. Internet Technol.* 20, 4, Article 44 (October 2020), 30 pages. <https://doi.org/10.1145/3418202>

1 INTRODUCTION

The importance of content on e-commerce websites has been widely recognized. High-quality and trusted product content has been empirically shown to have a substantial influence on user behavior, which is manifested in conversion rates and the volume of sales [31, 39, 43]. Product descriptions are an important element of the content displayed on product pages, alongside the product's title, image, and key attributes, such as model name, color, or size. Yet, such descriptions are often lacking or missing; for example, the majority of the Fashion products on eBay have no description at all. Even when available, product descriptions are often long and tedious to read, containing a lot of information that is insignificant for potential buyers. Our own analysis indicates that substantial portions of the product description sentences include details specific to a single listing or seller, information about the brand as a whole, and pure marketing statements.

Authors' addresses: Slava Novgorodov, eBay Research, Netanya, Israel, snovgorodov@ebay.com; Ido Guy, eBay Research, Netanya, Israel, idoguy@acm.org; Guy Elad, Technion, Israel Institute of Technology, Haifa, Israel, sguyelad@cs.technion.ac.il; Kira Radinsky, Technion, Israel Institute of Technology, Haifa, Israel, kirar@cs.technion.ac.il.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2020 Association for Computing Machinery.

1533-5399/2020/10-ART44 \$15.00

<https://doi.org/10.1145/3418202>

We refer to a *product description* as a written (textual) presentation of what the product is, how it can be used, and why it is worth purchasing. The purpose of a product description is to provide customers with details about the features and benefits of the product so they are compelled to buy.¹ In line with previous research on e-commerce content, we expect a good description to be informative, readable, objective, and relevant to the product (e.g., as opposed to a specific listing or a whole brand) [64]. We focus on concise descriptions of several sentences, which can be quickly consumed in their entirety and are especially suitable for small-screen devices, as e-commerce mobile applications have seen a remarkable growth and account for a major portion of the overall e-commerce traffic [32, 34]. Like other types of product content, credible descriptions have been shown to increase sales, while lacking descriptions withhold users from reaching a purchase decision or effectively searching for products [27, 40].

In light of the aforementioned challenges, we propose to use the “crowd” to generate trustworthy product descriptions, by leveraging the products’ user reviews [17, 30, 42]. Such reviews are often abundant on e-commerce websites and reflect the perspective of those who have already purchased the product. Therefore, potential buyers tend to trust reviews much more than they trust seller-provided content [33, 65].

User reviews primarily aim at reflecting a buyer’s subjective perspective and include personal opinions, stories, experiences, and complaints, which are not suitable to include in a product description (e.g., “*My old shoes wore down and I needed a new pair*” or “*I can plug it to each of my three cars*”). The large volumes of user reviews accumulated for popular products², with each review typically containing multiple sentences, makes them practically impossible to consume. As a result, users often read only a few reviews and may miss helpful information that appears in others. We observe that some portion of the review sentences are descriptive of the product [17], and suggest an extractive approach to generate *crowd-based* descriptions by combining original review sentences. While our approach require a product to be mature enough to have a certain amount of reviews, it can help boost the product’s exposure after gaining these reviews and, on the other hand, enhance the buyer experience with crowd-based descriptions for popular products.

The transformation from reviews to descriptions is a challenging task, which, to the best of our knowledge, is novel. While reviews aim to reflect the buyer’s perspective, descriptions typically reflect the viewpoint of the seller. Moreover, while reviews are meant to reflect a variety of subjective opinions, descriptions are expected to provide objective fact-based information. Several prior studies have examined review summarization (e.g., [29, 49, 79]), however such summaries do not necessarily contain descriptions of the product. In other words, a sentence that may be pivotal to the set of reviews, and therefore for its summary, might not be appropriate for a description. For example, the sentence “my girlfriend liked this dress a lot” may be included in a summary, but poses low value for a product description, as it is subjective and provides little information about the product.

Our extraction of candidate sentences from user reviews, to be included in the product depiction, is primarily supervised. We examine both classic machine learning models and deep learning approaches for the classification task, trained over thousands of sentences in two key e-commerce domains: Fashion and Motors. We also analyze the key reasons making review sentences unsuitable for a description. A deep multi-task learning classifier, which is based on mapping the top reasons to auxiliary tasks, is found to yield the best performance for the candidate identification task.

Following, we select the top sentences out of the candidate set for the final product description. To this end, we use a sentence similarity measure that helps diversify and avoid redundancies. Semantic similarity based on word embedding is found to be more effective than a bag-of-words approach. We experiment with several basic methods to produce the final description, which rely on the classification score from the candidate extraction

¹This definition is largely based on that by the Shopify e-commerce platform:
<https://www.shopify.com>

²For instance, a 200K Fire TV Stick with Alexa Voice Remote has over 200,000 reviews on Amazon.

process and the similarity measure. We perform a large-scale evaluation of the descriptions based on thousands of ratings from professional annotators, comparing the different methods and inspecting three description lengths: 3, 5, and 7 sentences. In addition, we provide a coverage analysis of the generated descriptions compared to the original seller-provided description and the complete reviews. We conclude our analysis by running an A/B test in a production environment that examines the impact on user traffic.

The main contributions of this work can be summarized as follows:

- To the best of our knowledge, this is the first work to suggest the extraction of product descriptions from reviews.
- We provide analysis and examples of how review sentences can be used for descriptions (what portion, which kind of sentences) across two principal yet very different e-commerce domains.
- We develop a classifier for identifying review sentences suitable for product descriptions, reaching an AUC of over 0.92.
- We present an end-to-end system for description generation from reviews, comparing different approaches for sentence selection, reaching an average rating of 4.3 (out of 5) per description.
- We present an extensive analysis of our generated descriptions, including their review coverage, comparison to the seller-provided descriptions, and their applicability across different e-commerce domains.
- We demonstrate results of an A/B test performed in our production environment and discuss its impact on user traffic.

2 RELATED WORK

Textual product descriptions have been explored in the e-commerce literature along with other seller-provided product content types, such as titles [21, 56], images [7, 22], and attributes [56, 61, 73]. Some of the studies refer to the product description in a broad sense, which encompasses the other content types, e.g., the list of structured attributes [2], title [74], or product image [52]. Other studies refer to a product description similarly to us, as the textual writeup that extends the title and attributes. Probst et al. [63] studied the extraction of attribute-value pairs from such product descriptions, in order to enrich the product's structured representation, used for tasks such as recommendation and matching. Shinzato and Sekine [69] proposed an unsupervised approach for the same task. Dumitru et al. [12] applied text mining and clustering techniques over product descriptions in order to recommend product features for a given domain. In a recent study, Pryzant et al. [64] showed that product descriptions of the type we study can help predict the product's business outcome. Experimentation was based on product descriptions and sales records from the Rakuten Japanese e-commerce website. None of these studies, however, provided a definition of a product description as presented in this work. Our previous work [59] addresses the same problem of description generation from reviews, however the current paper provides additional analysis, extended set of experiments and initial A/B testing results.

As mentioned in the Introduction, a related task to the one we explore is review summarization. Different from standard text summarization [20], where the goal is to generate a concise summary for a single [71] or multi-document [44], review summarization aims at extracting and summarizing opinions about a product from multiple reviews [49, 78]. Most studies focus on identifying the key attributes of an entity, such as a product, a movie, or a hotel, and then extracting key phrases that describe these attributes or the sentiment towards them (e.g., [29, 41, 62, 81]). Techniques used for this type of summarization include rule-based mining [46], topic modeling [55, 75], and neural networks [50, 76]. Ganesan et al. [15] proposed a graph-based summarization framework that generated concise abstractive opinion summaries of products. They represented the text opinions as a graph and used predefined rules to extract sub-paths from the graph and turn these into sentences. Although the sentences were readable, they missed crucial information and aspects of the product [35]. Gernai et al. [16] generated abstractive summarization of product reviews using discourse structure. They used templates to

Table 1. Fashion and Motors dataset characteristics.

	Fashion				Motors			
	Avg	Std	Median	Max	Avg	Std	Median	Max
Reviews per product	1118	1371	596	8271	883	985	378	7352
Sentences per review	3.71	3.09	3	103	4.05	3.73	3	98
Words per sentence	8.04	6.43	6	58	9.24	7.26	7	63
Number of products	892				807			
Number of reviews	997,274				712,904			

generate natural language summaries. The summary created was a statistical overview of the product with no detailed product information. An overview of review summarization techniques can be found in several surveys [35, 45, 47, 60].

While summarization seeks a good coverage of the main topics within the set of reviews, sometimes revolving around key product attributes, we aim to identify a unique subset of the reviews' content that is descriptive of the product. A sentence that may be pivotal to the set of reviews, and thereby for its summary, might not be appropriate for a description. For example, the sentence "*I like this hat very much*" may be included in a review summary, but poses low value for a product description, as it is subjective and provides little information about the product. In our experiments, described later in this paper, we found that fewer than 10% of the review sentences were suitable "as is" to take part in the product description.

Another related body of research has focused on extracting experiences [54, 58] and tips [23, 80] from user reviews. Somewhat similarly to the motivation presented in this work, these studies aim at helping users sift through the large volumes of reviews by identifying a more specific type of information within the reviews. Nonetheless, extracting experiences and tips is each inherently different than extracting descriptive sentences: experiences are subjective in nature, reflecting the unique viewpoint of an individual user (or group) and are thus not suitable, almost by definition, to be part of a product description. Tips, on the other hand, are defined as concrete and typically actionable pieces of advice. Therefore, their extraction actually aims at excluding purely descriptive sentences, of the type pursued in this work.

3 DATASETS AND CHARACTERISTICS

In this section, we describe the datasets used for our analysis and experimentation and their characteristics.

Datasets. Our research is based on two product datasets from two principal yet very different e-commerce domains: Fashion (clothing, shoes, and jewelry) and Motors (automotive parts and accessories). Both datasets were obtained from a large e-commerce website in the United States, representing best-selling products in each of the two domains.³ The datasets contain, per product, both its description and user reviews. Table 1 presents the characteristics of the two datasets. The number of products is rather similar in both datasets, while for Fashion there are more reviews per product and for Motors the number of sentences per review and their length in words is slightly higher. In addition, we used two larger datasets with over 10 million reviews of over 10,000 best-selling products in each domain (Fashion and Motors), for pre-training word embeddings, as will be described later in this paper.

Data Annotation. Labeling for training and evaluation in this work was performed by in-house professional editors (annotators), with domain expertise in both Fashion and Motors. The pool included a total of 20 editors,

³As of July 2018.

of whom different subsets were selected for different tasks, proportionally to the task's size. Unless otherwise stated, each evaluation was performed by a single editor.

Description Characteristics. The descriptions in our dataset are substantially longer than those we aim to generate. The average number of sentences per description in the Fashion dataset is 28.2 (std: 12.9, median: 27, min: 11, max: 65) and for motors it is 26.8 (std: 12.1, median: 29, min: 9, max: 68). To get a preliminary sense of the content of these descriptions, 50 descriptions from each dataset were annotated by two professional annotators. Only 45% of the sentences were labeled as suitable for a product description, with the key reasons for the sentences not being adequate including purely subjective marketing statements, accounting for slightly over 20% of the sentences (e.g., “*Give your clothes the luxury they deserve with these wonderful hangers!*”); information specific to a seller or a listing (18%; “*1 year limited warranty*”), and description of the brand as a whole (15%; “*For over 35 years, we have been one of the largest sunglasses brands.*”)

Reviews vs. Descriptions. In essence, reviews and descriptions hold fundamentally contrasting characteristics: reviews reflect a subjective opinion based on an individual experience, while descriptions are expected to be “dryer”, explaining what the product is and why it is worth purchasing. As a first step, we set out to examine the most prominent language differences between reviews and descriptions. To this end, we used Kullback-Leibler (KL) divergence, which is a non-symmetric distance measure between two given distributions [4]. Specifically, we calculated the terms, per each of the two domains, which contribute the most to the KL divergence between the language model of the reviews versus the language model of the descriptions and vice versa [24].

Table 2 presents the most distinctive unigrams. It can be seen that the unigram list most characterizing reviews (relative to descriptions) is rather similar between Fashion and Motors. The first-person pronouns ‘i’ and ‘my’ are at the top of the two lists, both common on reviews yet hardly occurring on descriptions. For example, ‘i’ occurs on 3.43% of the Fashion review sentences and 3.11% of the Motors review sentences, whereas for descriptions it occurs on 0.03% and none of the sentences for both datasets, respectively. Other prominent unigrams on the review lists include ‘was’, which typically reflects a past-tense experience (‘bought’ can also be observed on the Fashion list); the third-person pronouns ‘it’ and ‘they’; the adjectives ‘good’ and ‘great’, which often reflect a subjective opinion; and the emphasizing adverb ‘very’. Further down the list, beyond Table 2, we also encountered unigrams that refer to aspects of a specific listing of the product, which may vary from one seller to another, such as ‘price’, ‘cheap’, and ‘shipping’; other third-person references such as ‘she’ or ‘son’; and verbs that reflect subjectivity (‘recommend’, ‘like’, ‘love’) or past tense (‘arrived’, ‘purchased’, ‘got’, ‘ordered’). Inspecting the bigram lists, pairs such as ‘i have’, ‘i love’, and ‘for my’ were at the top of both the Fashion and Motors lists.

The most characterizing unigrams for descriptions relative to reviews are more dissimilar between Fashion and Motors. As can be seen in Table 2, each list includes its own domain-specific descriptive words, such as ‘jewelry’, and ‘cotton’ for Fashion, or ‘power’ and ‘fuel’ for Motors. The second-person pronoun ‘your’ ranks high on both lists, indicating that while first- and third-person language is used almost exclusively on reviews, second-person language is more characteristic of descriptions (e.g., ‘your’ occurs on 1.34% of the Fashion description sentences versus only 0.18% of the review sentences). The only other common words between the two lists are ‘provides’ and ‘designed’. Inspecting the bigram lists characterizing descriptions versus reviews, we observed similar trends, with domain-specific words, such as ‘fuel economy’, ‘power generator’, ‘machine wash’, and ‘your clothes’, as well as descriptive phrases such as ‘operates at’, ‘backed by’, ‘designed to’, and ‘suitable for’.

Overall, the above analysis gives an indication of the key differences in the language of reviews versus descriptions. Following, review sentences to be used as part of a description need to be carefully selected. We elaborate on this process in the next section.

Table 2. Most distinctive unigrams for reviews vs. descriptions ('Reviews') and vice versa ('Descriptions') in Fashion and Motors.

Fashion		Motors	
Reviews	Descriptions	Reviews	Descriptions
i	your	i	power
my	designed	my	generator
it	you	it	protection
they	features	was	your
was	comfort	great	fuel
but	jewelry	but	torqx
very	cotton	have	portable
great	imported	so	provides
good	inches	works	uego
these	fashion	this	economy
them	polyester	good	watts
bought	technology	they	designed
me	provides	very	advanced

4 CANDIDATE SENTENCE EXTRACTION

In this section, we describe a key component of our description generation method: the extraction of candidate review sentences that can be used for the product's description. Given the set of all user reviews for the product, our goal is to identify a set of sentences that can be used in their original form for a description of the same product. We first apply rule-based filtering based on the analysis presented in the previous section. We then apply a supervised approach that learns to identify review sentences suitable for a description. We examine different types of classifiers for this task and compare their performance based on a large labeled set of review sentences.

4.1 Rule-Based Filtering

Considering our analysis of linguistic differences between descriptions and reviews, we established several simple rules to identify review sentences that cannot be used as part of a description:

(1) **Short:** sentences of 3 words or fewer generally introduce little information and do not flow well as part of a product description. For example, "*Recommended*", "*Very good quality*", "*no complaints*", "*worth every penny*", or "*great product*". as well as "*good fit*" or "*very soft*" for Fashion and "*easy to install*" or "*works as expected*" for Motors, were among the most common short review sentences in our datasets. We encountered very few exceptions for short sentences that can be used in descriptions, e.g., "*Includes extra batteries*", "*Provides 180° view*", and "*Made in Japan*". Yet, these accounted for less than 0.3% of the short sentences. Moreover, they are only informative to a limited degree. Overall, short sentences accounted for 17.1% of all review sentences in the Fashion domain and 17.9% in Motors.

(2) **Personal:** sentences with a first-person pronoun, such as 'i', 'my', 'our', or 'us', or a 3rd-person *personal* pronoun, such 'she', 'his', 'hers', but not 'it' or 'them'. As demonstrated in the previous section, such pronouns hardly ever occur on a product description. Examples include "*I like the color of these jeans*"; "*Perfect fit for our car*"; "*My husband makes good use of them*"; and "*best gift for his birthday*". Overall, 35.9% of the review sentences matched this filtering criterion in the Fashion domain and 37.8% in Motors.

(3) **Listing-specific:** Some review sentences refer to listing-specific aspects, as observed in the previous section. Examples include "*Great value for a fair price*", "*Delivery was smooth and fast*", or "*The seller was responsive and*

Table 3. Reasons for review sentences labeled ‘*bad*’ and their distribution (portion of all sentences marked *bad*) for Fashion and Motors.

Reason	% Fashion	% Motors	Example
Subjective	52.50%	52.43%	It was the easiest jumpstart ever.
Missing context	16.86%	16.82%	Otherwise it remains idle.
Refers to a listing’s aspect	8.40%	6.73%	10 bucks for 3 pairs is a great deal.
Non-informative	7.95%	6.42%	This shirt is great.
Poor language and spelling	4.91%	5.17%	Extremely easy setup let’s you pull you vehecle’s code fast.
Negative sentence	3.90%	4.25%	Only issue is the pretty thin material.
Expresses explicit doubt	2.40%	2.30%	Probably good also for bicycle tires.
Refers to the description	1.83%	1.74%	The hat is exactly as described.
Other	1.49%	1.24%	Like others here have said, this gas can has a long rotating nozzle.
Too specific/detailed	0.64%	2.52%	Great for Honda 2003 2.0L.
Offensive language	0.12%	0.19%	Fantastic product, bright as sh*t.

helpful”. Since our goal is to produce a description at the product level rather than the listing (item) level, such aspects are not suitable for referencing as part of the description, since they may vary according to the seller. Our blacklist for this rule included the unigrams ‘price’, ‘cheap’, ‘expensive’, ‘delivery’, ‘shipping’, ‘seller’, and ‘warranty’. Overall, 5.7% of the review sentences matched this rule in the Fashion domain and 6.1% in Motors.

Our rules aim to filter out sentences that are not suitable for a description with a very high likelihood, almost by definition. We therefore did not filter out other potential candidates, such as sentences in past tense, since these can sometimes be appropriate (e.g., “Tested on several cars”). We also did not automatically filter out sentences from reviews with low ratings, because the vast majority of the reviews in our dataset had a positive rating, in line with past work that indicated online user reviews tend to the positive [9]. Overall, 53.7% and 55.2% of the review sentences were filtered out using the three rules above, for Fashion and Motors, respectively.⁴

4.2 Automatic Classification

After the initial rule-based filtering, we set out to explore a supervised approach, where we trained a classifier to predict whether a product review sentence is suitable as a description sentence for the same product. To this end, we sampled uniformly at random, out of all review sentences that were not filtered out by the rules, 25K sentences for each of the two datasets, Fashion and Motors. Each of the 25K sentences were then labeled by a group of 10 annotators, who were asked to mark them as either ‘*good*’ or ‘*bad*’, i.e., suitable to be part of a product description or not. In case the sentence was labeled *bad*, the annotators also selected a reason. The set of reasons was identified in an earlier round of labeling, and included ‘other’ in case none of the 10 reasons was appropriate. The annotators received detailed guidelines, explaining what makes a sentence suitable versus unsuitable for a description, with examples of *good* and *bad* sentences, as well as examples for each of the possible reasons for *bad*. They also performed qualification tests, i.e., an iterative process of labeling, followed by feedback from other annotators, until the quality was aligned among all. At the end of the process, the inter-annotator agreement for the task of *good* versus *bad* labeling, measured by Cohen’s kappa [10], was 0.89 for Fashion and 0.9 for Motors, calculated over a set of 300 sentences labeled by two different annotators. Overall, 8.55% of the Fashion and 7.97% of the Motors sentences were labeled *good*.⁵

⁴The portions of all three rules do not sum up to the total number of filtered sentences, since some sentences matched more than one rule.

⁵The dataset is available at https://tdk.cs.technion.ac.il/research-files/description_generation_from_reviews.zip

Table 4. Most distinctive unigrams and bigrams for *good* sentences vs. *bad* sentences (‘*Good*’) and vice versa (‘*Bad*’) in Fashion and Motors.

Fashion				Motors			
<i>Good</i>	<i>Bad</i>	<i>Good</i>	<i>Bad</i>	<i>Good</i>	<i>Bad</i>	<i>Good</i>	<i>Bad</i>
easy	was	easy to	a little	easy	was	easy to	it was
very	but	good quality	it was	very	but	to use	so far
sturdy	not	well made	so far	great	it	very easy	than the
quality	it	very sturdy	a bit	use	would	to install	better than
comfortable	love	perfect for	but the	install	not	works great	a little
nice	than	are very	as expected	quality	than	good quality	a bit
well	as	they fit	they were	works	had	well made	seems to
great	buy	great for	but it	well	as	to apply	as advertised
durable	had	high quality	happy with	nice	buy	very well	as described
is	would	to assemble	but they	your	this	comes with	happy with

Table 3 lists the different reasons for *bad* sentences. It can be seen that the distribution is similar for Fashion and Motors, with subjective sentences accounting for a little over half of the *bad* sentences in both, followed by sentences with a missing context. The only noticeable difference between the domains is for the too specific/narrow reason, which is generally uncommon, but occurred substantially more frequently in Motors.

4.2.1 Good vs. Bad Review Sentences. Before building our classifier, we performed a statistical analysis comparing the review sentences labeled *good* by our annotators with review sentences labeled *bad*.

Table 4 presents the most distinctive unigrams and bigrams for *good* review sentences versus *bad* review sentences and vice versa, for Fashion and Motors. Distinctive terms were calculated using KL divergence as described in Section 3. It can be seen that *good* sentences in both domains include positive adjectives and adverbs, such as ‘easy’ and ‘very’, which are at the top of both unigram lists, as well as ‘great’, ‘sturdy’, ‘nice’, and ‘well’, which are used to describe products as easy to use/install/assemble; being of good/high quality or well made; or being perfect/excellent/great for a specific use. While many of the terms are common for both domains, in Fashion traits such as comfortable, durable, or fitting well are salient, while in Motors aspects of installation, application, and work are more dominant. For *bad* sentences, there is also a substantial overlap between the domains, implying the language differences between *good* and *bad* sentences can be generalized. The list of terms reflects some of the key reasons for sentences not fitting into a description, such as subjective viewpoints (e.g., ‘better than’, ‘happy with’, ‘love’), personal experiences in past tense (‘it was’), expression of doubt (‘seems to’), negative sense (‘but’, ‘not’) and a reference to the description itself (‘as described’, ‘as advertised’).

Table 5 shows the portion of sentences labeled as *good* according to various characteristics of the sentence and the review it originated from, for Fashion and Motors. Inspecting the sentence length, starting 4 words, it can generally be seen that shorter sentences are somewhat more likely to be labeled *good*. For Fashion, the “optimal” sentence length is 6-8 words, while for Motors the percentage of *good* sentences consistently decreases with the length. Long sentences are more likely to include expressions that would make them unsuitable for a description based on some of the reasons detailed in Table 3, particularly subjectivity, which is the most common. Looking at the length of the review, it can be observed that sentences originating from long reviews are less likely to be *good*. For both Fashion and Motors, the review length that yields the highest portion of *good* is 3 sentences. Finally, inspecting the position of the sentence within the review (while controlling for the review length), there is a

Table 5. Percentage of *good* sentences in Fashion and Motors distributed by sentence length, review length, and sentence position within the review.

Sentence length (words)		4	5	6	7 – 8	9 – 10	11 – 12	13 – 15	16+
	Fashion	9.22%	9.96%	10.21%	10.38%	8.32%	7.81%	6.69%	6.05%
	Motors	11.07%	9.91%	9.43%	8.86%	7.42%	6.86%	6.60%	6.63%
Review length (sentences)		1	2	3	4	5 – 6	7 – 9	10+	
	Fashion	7.49%	9.19%	10.27%	9.46%	8.33%	8.11%	6.81%	
	Motors	9.15%	10.04%	10.30%	9.88%	8.24%	6.69%	5.66%	
Sentence position in review	Fashion		1 st	2 nd	3 rd	4 th	5 th		
		Review length 3	11.48%	9.24%	8.90%				
		Review length 4	10.67%	11.16%	7.02%	5.27%			
	Motors	Review length 5	8.92%	10.99%	8.72%	6.64%	4.35%		
		Review length 3	12.25%	8.78%	7.56%				
		Review length 4	11.31%	10.76%	7.97%	6.82%			
		Review length 5	9.89%	11.50%	8.88%	5.59%	3.86%		

consistent trend indicating that sentences that occur towards the beginning of the review (not necessarily the first sentence) have higher likelihood to be *good*. The last sentence of the review has a particularly low likelihood to fit a description. This trend is consistent for Fashion and Motors and persists for higher review lengths not shown in Table 5. Overall, it suggests that users are more likely to include descriptive details earlier in the review.

4.2.2 Supervised Approaches. For the *good* versus *bad* classification task, we experimented with both traditional machine learning models and deep learning approaches. For the latter, we examined the effect of an attention mechanism and the use of a multi-task learning approach that tries to predict a reason for a sentence being inadequate for a description. We used 5-fold cross-validation to tune the hyper-parameters and evaluate the classifiers. As evaluation metric, we used the area under the ROC curve (AUC).

Naïve Bayes and XGBoost. We examined two common models for text classification: Naïve Bayes [66] and XGBoost [8]. Our features included textual features, specifically the unigrams, bigrams, and trigrams of the review sentence, and statistical features of the sentence and originating review, based on Table 5. For Naïve Bayes, we tuned the type of n-grams (unigrams, bigrams, or trigrams) and for XGBoost, in addition, the maximum depth of a tree, minimum child weight, as well as the learning rate and number of rounds (trees). Results for both classifiers are presented at the top rows of Table 6, indicating XGBoost achieved a higher AUC.

To better understand the contribution of non-textual features, which are based on the characteristics presented in Table 5, we trained the XGBoost classifier with different feature subsets (while always using textual features). Feature types included sentence length (in words and characters), review length (in words and sentences), and position (absolute position and normalized by the total number of review sentences). Results, presented in Table 7, indicate that position features yielded a slightly higher performance gain than sentence and review length. Overall, however, as can be observed in the ‘All’ row, the contribution of these features on top of the textual features was minor (+0.97% for Fashion and +0.96% for Motors). We therefore did not use these features in the remainder of our experiments and worked with the review text only.

Table 6. AUC performance results for classifying review sentences as *good* or *bad* for the product’s description.

Classifier	Fashion	Motors
Naïve Bayes	0.779	0.798
XGBoost	0.831	0.839
LSTM	0.914	0.916
LSTM-Attention	0.915	0.916
LSTM-MTL	0.924	0.924

Table 7. AUC performance of the XGBoost classifier with textual features when using (‘Only’) or disregarding (‘Exclude’) additional types of features. The ‘All’ row refers to all additional types jointly.

Additional Features	Fashion		Motors	
	Only	Exclude	Only	Exclude
Sentence Length	0.824	0.828	0.833	0.835
Review Length	0.828	0.826	0.835	0.834
Sentence Position	0.829	0.828	0.837	0.836
All	0.831	0.823	0.839	0.831

While the XGBoost classifier demonstrated reasonable performance, we set out to explore how it can be further enhanced. In recent years, various applications in text classification have shown significant improvement by the use of deep learning models [19, 38]. We therefore explored several deep learning approaches for our task:

LSTM. A recurrent neural network based on long short-term memory (LSTM) [28] architecture, with pre-trained word2vec embeddings using continuous bag of words (CBOW) with negative sampling [53]. We experimented with pre-training over Wikipedia and over our own datasets of more than 10 million product reviews for Fashion and Motors, respectively, as described in Section 3. Pre-training using our own data, with separate models for Fashion and Motors, consistently achieved a slightly better performance. We henceforth only report the results when using word2vec pre-trained based on our own data.

LSTM with Attention. Attention mechanisms enable the network to focus on relevant parts of the input [77]. The overall architecture of the “attention network” consists of two components: an LSTM-based word sequence encoder and a word-level attention layer. Given a review with the words $w_i, i \in [1, N]$, we first embed the words using pre-trained word2vec, as previously described. We then use the LSTM network to produce the hidden states $h_i, i \in [1, N]$. The attention mechanism is subsequently used to put more focus on certain words in the review sentence. For example, in the sentence “*this is a high quality product*” the words “*high quality*” should receive higher weight. To this end, we use the attention mechanism as follows:

$$u_i = \tanh(W_w h_i + b_w) \quad (1)$$

$$\alpha_i = \frac{\exp(u_i u_s)}{\sum_j \exp(u_j u_s)} \quad (2)$$

$$s = \sum_i^N \alpha_i h_i \quad (3)$$

We feed the word annotation h_i through a single-layer perceptron network to receive u_i , a latent representation of h_i . Then, we calculate the similarity of u_i with a word-level context vector, normalized by a softmax function,

to produce the word's importance weight. We then construct the review vector as a weighted sum of the word annotations based on each word's weight.

Deep Multi-Task Learning. Multi-task learning (MTL) is based on the idea that features trained for one task can be useful for related tasks. Models for all tasks of interest are jointly trained with an additional linkage between their trainable parameters, aiming at improving the generalization error [6]. Multi-task learning can be viewed as a form of inductive transfer learning, which can help improve the model by introducing an inductive bias. In the case of MTL, the inductive bias is provided by the *auxiliary tasks*, which lead the model to prefer hypotheses that explain more than one task. MTL has been widely used for deep learning tasks, e.g., in computer vision [18] and natural language processing [25, 70].

Our network is based on an MTL hard parameter sharing architecture [67], which includes an LSTM layer shared among all tasks and separate feed-forward networks per task. We learned the task of predicting if a review sentence could be included as part of a description (*good/bad*) with the additional auxiliary tasks that specialize on a specific reason, particularly each of the top four reasons listed in Table 3, which cover over 80% of the *bad* sentences in both Fashion and Motors. Each specialized classifier is trained to predict if the sentence falls under the specific reason for not fitting in a description. For example, the 'subjectivity' classifier learns to predict subjective sentences. Overall, the specialized classifiers may help the main task of predicting *good* or *bad* sentences. Similarly to the LSTM with attention, given a review sentence, we first embed the words using our pre-trained word2vec, and then use the shared LSTM encoder to produce the latent representation of the review sentence. This representation is then fed in parallel into 2-5 fully-connected feed-forward networks (one per task), each with two hidden layers, trying to predict the output for that task independently. The ultimate loss function to be optimized assigns a weight to each of the task-specific loss functions. Notice the weighted loss function serves for optimization at training time, while testing remains for the *good* versus *bad* task only. In Figure 1, we demonstrate our MTL hard parameter sharing network, with a shared LSTM layer and 5 separate feed-forward layers, one for each task.

For each of the three deep learning methods, hyper-parameter tuning included the batch size, the dropout rates, and the number of hidden units in the LSTM layers. In addition, we experimented with both Adam [36] and RMSProp [72] optimizers. For the MTL architecture, we also tuned the subset of tasks to be included and their weight in the loss function, as well as the dropout rates and number of hidden units of the feed-forward layers.

4.2.3 Performance Results. Table 6 presents the AUC results for the *good* versus *bad* task. Evidently, the deep learning classifiers achieved substantially better results than the Naïve Bayes and XGBoost classifiers. The attention mechanism did not lead to any performance gain, possibly due to the sentences' length: as shown in Table 1, the median number of words per review sentence is 6 and 7, for Fashion and Motors, respectively. The effectiveness of the attention mechanism, however, typically lies in longer sentences [48]. The LSTM-based MTL model achieved the best AUC results for both datasets, Motors and Fashion. Table 8 presents the auxiliary task subsets that yielded the best performance. For Motors, the combination of all four auxiliary tasks led to the best performance, while for Fashion it was the combination of the subjective and missing context tasks. Indeed, these are the two most common reasons for sentences being labeled *bad*, as shown in Table 3. Other combinations listed in Table 8 yielded close performance to the top ones. Noticeably, using only one auxiliary task – for identifying subjective sentences – produced a substantial part of the performance gain compared to the vanilla LSTM, especially for Fashion. As previously shown (Table 3), 'subjective' is the most common reason, covering over half of the *bad* sentences for each of the two domains.

As mentioned before, we compared Adam and RMSProp optimizers. The latter was found to perform better, slightly but consistently across different number of tasks, for both Fashion and Motors.

4.2.4 Effect of Labeled Data Size on Performance. For the *good* versus *bad* learning task, we had a large dataset of training data at our disposal, with 25,000 example for each of the two domain, used in a 5-fold cross validation

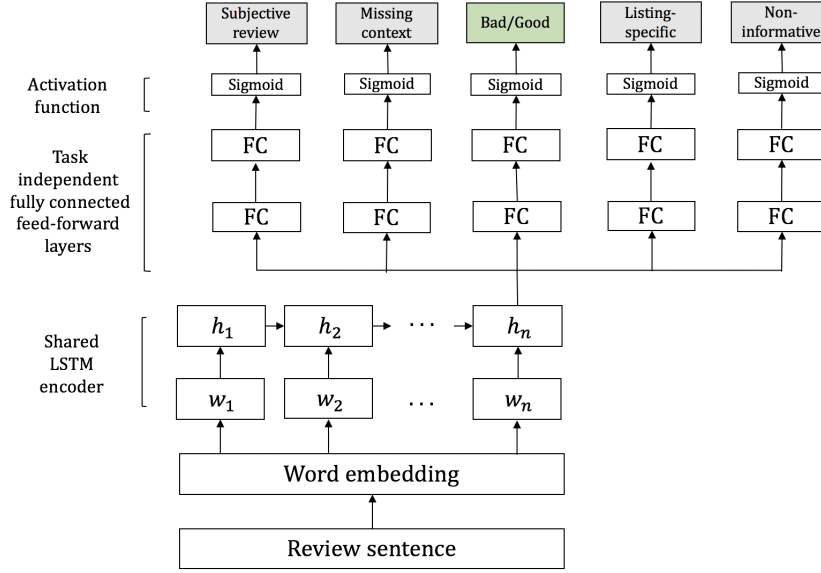


Fig. 1. Hard parameter sharing for deep multi-task learning architecture.

Table 8. AUC performance of the deep MTL classifier with different subsets of the four auxiliary tasks: subjective ('subj'), missing context ('MC'), listing-specific ('LS'), and non-informative ('NI').

Fashion		Motors	
Subj, MC, LS, NI	0.921	Subj	0.920
LS, NI	0.922	MC, LI	0.922
Subj	0.922	MC, LS, NI	0.923
Subj, LS, NI	0.923	SUBJ, MC	0.923
Subj, MC	0.924	Subj, MC, LS, NI	0.924

setting as described above. Since labeled data is typically expensive to obtain, we set out to explore the effect of the labeled data size on the performance of our model. To this end, we experimented with smaller labeled sets by sampling uniformly at random subsets from the full labeled dataset in each of the two domains. Figure 2 depicts the AUC results obtained using 5-fold cross validation with the LSTM classifier, according to the size of the labeled data used for cross validation. While there is indeed a trade off between the labeled data size and the model's performance, it can be seen that fairly good results can be obtained using a smaller number of labeled instances. For example, using 10,000 examples yielded an AUC of 0.897 in Motors, down 2.1% compared to the model using 25,000 instances (for Fashion the parallel decrease was 2.3%). Using only a tenth of the original labeled set, i.e., 2500 instances, yielded an AUC of 0.839 for Motors, down 8.4% compared to the "full" model (8.5% for Fashion). Across all data sizes we experimented with, the trade off between the AUC and size of the labeled data was similar between Fashion and Motors. Using semi-supervised approaches, leveraging large volumes of unlabeled review data that is publicly available in many e-commerce platforms, can further improve the trade-offs reported here between training data and performance of identifying *good* sentences.

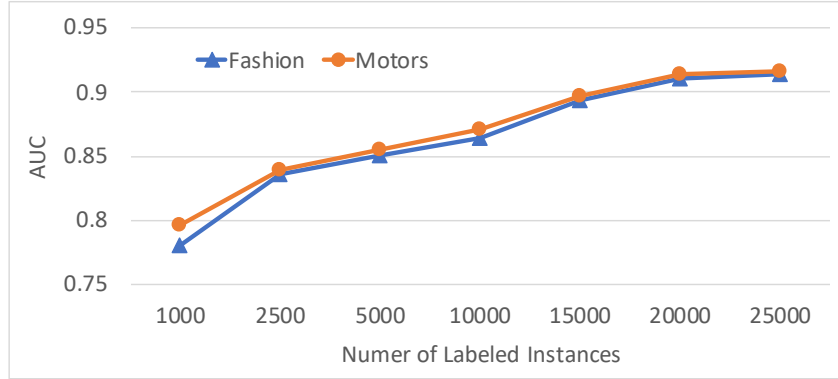


Fig. 2. AUC performance results of the LSTM classifier by the number of labeled instances used for 5-fold cross validation.

5 DIVERSIFICATION

In order to create coherent and concise descriptions, we wished to avoid the inclusion of very similar sentences in the same description. After identifying candidate sentences for the description, our next step was therefore the identification of redundant sentences in order to increase the diversity in the content of the final description. The diversification phase was based on the computation of similarity between candidate sentences. To this end, we used the common cosine similarity measure, while experimenting with three sentence representation methods:

Weighted Bag-of-Words (BOW). Each sentence is represented as the TD-IDF weighted vector of its words [1].

Average Embedding. The weighted BOW approach measures similarity based on actual word overlap. In order to capture deeper semantic similarity, an embedded representation of the sentence can be used. Specifically, we apply word embedding using word2vec and average over all the sentence's words to produce the final sentence representation. The word2vec model is trained separately for Fashion and Motors, based on the dataset of over 10 million reviews per domain, as described in Section 3.

Weighted Embedding. This method applies the same word2vec embedding as described above, but instead of averaging the embeddings to generate the sentence representation, a common word weighting approach is applied [3]. This weighting is performed using TF-IDF scores, as it has been empirically shown to achieve similar results to learning the language model of the sentences.

To decide whether two sentences are different enough to be included in the description, a similarity threshold θ for each method had to be determined. To avoid using a hard-coded threshold, we learned θ in an unsupervised manner from the data, for each of the two domains (Fashion and Motors). To this end, we considered the product descriptions in both of our datasets, which were manually curated by domain experts. We measured the similarity between each pair of sentences within each description and considered the 90-th percentile as the threshold θ . The resultant values of θ were 0.73 for Fashion and 0.76 for Motors. Intuitively, we set our threshold to allow a similarity level that is aligned with the typical degree of similarity that exists in professionally-written descriptions.

In order to compare the three representations, we conducted a small-scale experiment. For each domain (Fashion and Motors), we sampled 1000 pairs of sentences in the following way: we sampled 200 products from the respective dataset uniformly at random. For each product, we sampled 10 pairs of review sentences, uniformly at random out of all (unordered) pairs of sentences labeled *good*. We evaluated only *good* sentences, since the *bad* sentences were not considered for the final description. The sampled pairs were labeled by 3 professional annotators as either similar or not. The annotators were instructed to label a pair of sentences as similar in case

Table 9. Accuracy for classifying sentence pairs as similar or not using different representations.

Sentence Representation	Fashion	Motors
Weighted bag of words	76.6	78.3
Average embedding	88.3	88.6
Weighted embedding	92.2	91.7

Table 10. Example sentence pairs classified as ‘similar’ by average and weighted embedding and ‘not similar’ by BOW.

Good looking and sturdy hangers.	Hangers look nice and very solid.
Flexible construction and good fit.	Very well made and great fit.
Fantastic way to maximize space.	Great for minimizing storage space.

one of them did not add any substantial piece of information to the reader on top of the other. The agreement between annotators, measured by Cohen’s Kappa, was 0.84 for Fashion and 0.86 for Motors, calculated over 100 pairs for each domain, which were evaluated by two different annotators.

Overall, 35.2% of the sentences were labeled as similar, indicating there was indeed a high level of redundancy among *good* sentences. We tested each of the three methods against this labeled set by calculating its accuracy, i.e., the portion of pairs it correctly classified out of the 1000 pairs. Results are depicted in Table 9. It can be seen that both embedding-based methods reached a substantially higher accuracy than the weighted BOW method, for both domains. Weighted embedding slightly outperformed average embedding at over 90% for both domains, and was therefore our choice as sentence similarity method for the rest of our experiments.

Table 10 presents a few examples deemed ‘not similar’ by the weighted BOW method and ‘similar’ by both average and weighted embedding, as well as the human annotators. These demonstrate why semantic representations capture similarities that may be overlooked by a purely-lexical method. For example, for the first pair, the tokens ‘and’, ‘are’, ‘nice’ and ‘hangers’ appear in both sentences. As the first two are generic words, the Bag-of-Words model gives higher weight to the words ‘hangers’ and ‘nice’ tokens, and outputs a low similarity between the two sentences. On other hand, the embedding algorithm observes the similarity between the words ‘looking’ and ‘look’, and also ‘strudy’ and ‘solid’. Hence, the embedding algorithm marked the pair as “similar” (exactly as the human expert).

6 FINAL DESCRIPTION GENERATION

In this section, we describe our methods for producing the final product description. Following, we present a detailed evaluation of the generated descriptions. Our evaluation compares the different methods for producing the descriptions and different description lengths by their overall quality ratings, as well as by ratings of more specific description aspects. In addition, we examine the coverage of our descriptions both with respect to the product’s set of reviews and the product’s original description. We finally present results over publicly available data and experiment with cross-domain description generation.

6.1 Sentence Selection

The final step of our product description generation process is the selection of the final sentence set. From the previous steps, we have a list of candidate sentences, with their classification score reflecting their likelihood to be suitable to a product description, a similarity metric that can be applied to a pair of candidate sentences to measure their closeness, and a similarity threshold θ reflecting the desired similarity between sentences in a

description. With these in hand, we examine four methods to generate the final description, given an integer K that determines the desired number of sentences in the description.

Greedy approach. This method traverses the list of candidates according to their *good/bad* classification score, from highest to lowest, and adds a candidate to the description if (and only if) it is not similar (i.e., has a similarity score of θ or higher) to any of the candidates already selected for the description. The process stops when K sentences have been added.

LexRank. This method uses LexRank, a common extractive summarization method that yielded high performance results for several text summarization tasks [14]. Specifically, LexRank assigns an importance score per sentence, using random walks and eigenvector centrality. We apply LexRank on top of the list of candidate sentences and use its score to rank the candidate sentences. Intuitively, the LexRank score indicates how well the candidate covers the information included in the other candidates. The method then operates similarly to the greedy approach: after ranking the candidate sentences by their LexRank score, the list is traversed by descending score and a sentence is added in case it is not similar to any of those previously selected for the description.

K-means classification score. This method partitions the candidate sentences into K clusters using the k -means algorithm [51] with $k=K$. The distance between sentences is calculated as $(1-s)$, where s is our sentence similarity measure described in Section 5. For the final product description, the clusters are traversed from the largest (representing the highest number of review sentences) to the smallest and from each cluster, the sentence with the highest classification score is added to the description.

K-means centroid. The method works as the previous one, but instead of selecting the candidate with the highest classification score from each of the K clusters, the candidate closest to the centroid is selected, assuming it best represents the cluster.

Note that smart ordering of the sentences within the final description is beyond the scope of this work. Currently, the sentences are ordered by classification score (greedy), summarization score (LexRank), or cluster size (both k -means methods).

6.2 End-to-End Description Evaluation

For an end-to-end evaluation of our approach, we set out to examine full descriptions generated using one of the four methods described above. For candidate sentence extraction, we used the LSTM-MTL classifier and for similarity measure we used weighted sentence embedding, since both were found to yield the best performance for their respective tasks, as previously reported. We experimented with descriptions of three lengths: 3, 5, and 7 sentences. These represent relatively concise content, which can be swiftly consumed in its entirety, and is especially suitable for small-screen devices, such as mobile phones. We generated descriptions for all products in both the Fashion and Motors datasets. Using our candidate extraction method (Section 4), we identified review sentences that were suitable for a description for each product in the sample. On average, for each Fashion product we identified 46.1 *good* sentences (std: 9.3, median: 43, min: 24, max: 97) and for Motors 47.3 (std: 11.7, median: 45, min: 19, max: 103). For each such product, we then generated descriptions of length 3, 5, and 7, using each of the four methods described above (a total of 12 description versions). These descriptions were evaluated by 15 professional annotators. We ensured each annotator evaluated no more than one description per product.

Annotators were presented with the product's title, image, and generated description and were asked to rate the quality of the description's text for serving as the product's description, on a 5-point Likert scale, from 'very bad' to 'very good'. In addition, to examine finer-grained aspects of the description's quality, along the lines presented in the Introduction, annotators were asked to assess, on the same scale, to what extent the description's text was readable, informative, objective, and relevant to the product. They were provided with good and bad examples for each of the questions. Figure 3 demonstrates the user interface developed to collect annotators' input.

Champion 3500-Watt Portable Generator



Description: Super dependable and really quiet. Very easy to use. Runs without issues at constant voltage.

Overall description quality:

☐ Very Bad ☐ Bad ☐ Average ☐ Good ☐ Very Good

	Very Bad		Average		Very Good	
Readable:	<input type="radio"/> 1	<input type="radio"/> 2	<input type="radio"/> 3	<input type="radio"/> 4	<input type="radio"/> 5	
Informative:	<input type="radio"/> 1	<input type="radio"/> 2	<input type="radio"/> 3	<input type="radio"/> 4	<input type="radio"/> 5	
Objective:	<input type="radio"/> 1	<input type="radio"/> 2	<input type="radio"/> 3	<input type="radio"/> 4	<input type="radio"/> 5	
Relevant to the product:	<input type="radio"/> 1	<input type="radio"/> 2	<input type="radio"/> 3	<input type="radio"/> 4	<input type="radio"/> 5	

Comment:

Fig. 3. Product description evaluation interface.

The inter-annotator agreement for the main quality evaluation question, measured using weighted Kappa [11], was 0.83 for Fashion and 0.84 for Motors, calculated over a set of 200 generated descriptions evaluated by two different annotators. Table 11 presents a few examples of our generated descriptions, which were rated 5.

6.2.1 Rating Results. Table 12 shows the ratings of the main question (overall description quality) in our end-to-end evaluation. Generally, the generated descriptions received high ratings, with an average of above 4 for all methods except for the greedy approach. Ratings for Fashion and Motors were similar, giving some indication that our method may be generalized to other e-commerce domains. Among the four methods, LexRank achieved the highest ratings, consistently for all three description lengths for both Fashion and Motors. At the other extreme, the greedy approach yielded substantially lower ratings than the other methods. K-means with the closest-to-centroid selection was consistently the second best, higher than k-means with highest-score selection. This suggests that selecting the sentence based on representation of the whole cluster rather than the individual prediction score for description suitability is preferable. The superiority of the LexRank approach implies that good representation of the whole set of candidates for final sentence selection works well. All differences between the methods were statistically significant, except between LexRank and k-means-centroid with $K=3$.⁶ As for the length of the description, there was no consistent trend, but the highest ratings for LexRank were achieved for $K=5$ sentences by an insignificant difference from both $K=3$ and $K=7$ sentences.

Figure 4 presents the rating results for the additional questions in our evaluation, relating to the different quality aspects. Across all questions, LexRank consistently achieved the best ratings, with the rest of the methods ranked similarly to the general quality question. Among the four questions, ratings were highest for product relevance, with LexRank achieving an average of over 4.5 for all three description lengths. This indicates that as could be expected, the reviews serve as a relevant source of information about the product. Objectivity was also rated high in general, indicating our method for filtering the objective parts out of the subjective reviews (step 1) works well. Readability received the lowest rating, with the greedy approach (and for $K=5$ and $K=7$ also k-means with highest-score selection) performing especially poorly. Inspecting the ratings by description length, objectivity and readability tend to decrease with description length, while informativeness increases. Product-relevance remains stable regardless of the description length. Overall, it is expected that as the description is longer, it is more likely to contain non-objective parts and become less readable due to connectivity or repetition issues, while it is also likely to contain more details and therefore become more informative. Inspecting the different ratings in Figure 4, descriptions of 5 sentences seem to yield the best trade-off between these two trends. For

⁶Statistical significance was measured using one-way ANOVA with Tukey post-hoc comparisons for $p < 0.01$.

Table 11. Example descriptions rated 5.

Domain	Product	Description
Fashion	Dri-fit shirt	Perfect quality of the material. Great for colder weather. Holds up to multiple washings. It does keep your body heat in. Definitely fits well without restricting movement.
Fashion	Socks	Perfect thickness for shoes or boots. Extra padding at toes. The quality is excellent. Easy to handle and very comfortable. No shrinkage in washer or dryer.
Fashion	Hat	Good hat for warm weather. Breathable and protects from the sun. Lightweight material and great fit. Very easy to wash. This hat can fit both males and females perfectly with the adjustable back.
Fashion	Jacket	The cuffs on the sleeves are adjustable, which is perfect for keeping wind out when biking or for just a tighter fit around your wrists. It's very soft on the skin. The zipper works easily. The material is well made and the hood tucks away easily and hides well when not needed. Folds up to fit in bag.
Motors	Scratch removal system	Great for minor to moderate damage. Use it for removing scratches from your car. There is enough for several repairs. Coat twice for extra protection. Very recommended after wax.
Motors	Lights	This is a great led bulb. Very good at night. Way better than halogen. Very easy too install. Good wide beam and bright white color.
Motors	Motor oil	It works well and lasts for a reasonable period of time. The oil is still very clean when changed. Engine runs smoother and pulls better. Good oil for any motorcycle. This oil is made for high temp engines.
Motors	Mirrors	The mirrors are definitely helpful when backing into a striped parking spot because the stripes are easy to see in the mirror. The advantage is that you can better see the area of interest. It definitely serves its purpose. Every new driver should have these. They cover the blind spots perfectly.
Toys	Magnetic cubes	Works for all ages – toddlers to teens. The letters are a great addition and a fun way for young kids to learn their alphabet. Very colorful and high quality magnetic toy! Very easy to build. Perfect for family time.
Toys	Puzzle	The pieces are thick and lock together well, even on carpet. It comes in a nice box for easy storage. The pieces are big so no one can choke. The picture has great colors and is very bright. It gives beginning skills, and gives a chance for social interaction.
Electronics	TV	This TV is very capable for anything. The antenna integration and USB storage for recording live programming are great. The remote and menu's are easy to use. The sound is also surprisingly good. Good screen clarity nice picture.
Electronics	Mobile phone	This phone has a very good battery. The screen is big enough for watching movies. Includes power charger and headphones. The quality of photos is excellent. The phone is very smart and connects well with other devices.

instance, consider the first example in Table 11. When the generated description with $K=3$ included the first 3

Table 12. Average rating of description quality of 3, 5, and 7 sentences for products in Fashion and Motors.

	K=3		K=5		K=7	
	Fashion	Motors	Fashion	Motors	Fashion	Motors
Greedy	3.95	3.92	3.81	3.81	3.80	3.77
K-means score	4.10	4.09	4.08	4.06	4.04	4.01
K-means centroid	4.29	4.28	4.17	4.22	4.26	4.23
LexRank	4.36	4.30	4.38	4.35	4.36	4.33

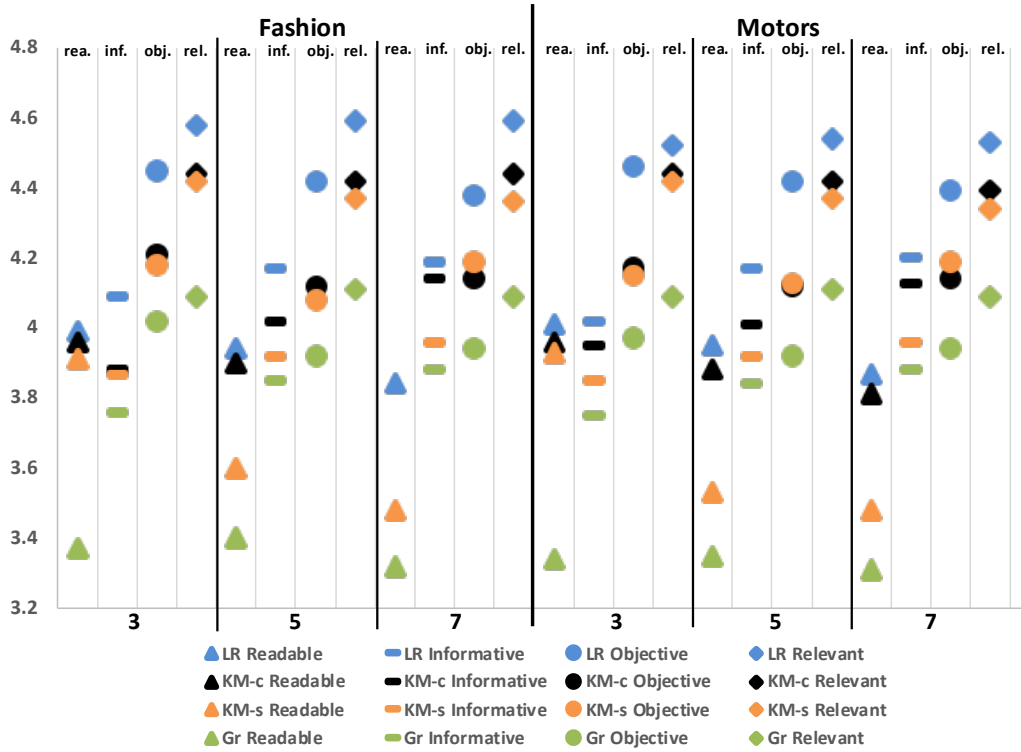


Fig. 4. Average rating of four aspects (readable, informative, objective, relevant) for descriptions of length 3, 5, and 7 sentences, generated by four methods: LexRank (LR), K-Means centroid (KM-c), K-Means score (KM-s), and Greedy (Gr).

sentences, its ‘informative’ rating reduced. On the other hand, when the $K=7$ version included the extra sentence “Best socks for sports, especially running,” the ‘objective’ rating reduced.

We also examined the tie between the description length, for a fixed number of sentences (K), and its overall rating. We conjectured that longer sentences may yield higher ratings. However, we did not find any correlation. The ratings for all questions remained stable regardless of description length for the different values of K and both domains.

Table 13. Average number of reviews (sentences) needed to produce a deception of 3, 5, and 7 sentences, as measured for 100 products from Fashion and 100 products from Motors.

Domain	3-sentence description	5-sentence description	7-sentence description
Fashion	15.6 (64.7)	18.5 (90.9)	23.2 (109.3)
Motors	16.2 (67.2)	19.5 (94.2)	25.1 (112.4)

6.2.2 Evaluation over public data. The datasets used in this work are proprietary and cannot be shared due to business sensitivity. We therefore set out to get a sense how well our method performs on publicly available data. As a first step, we observe that 13.5% (120 in total) of the Fashion products and 16.7% (135) of the Motors products sent for our end-to-end evaluation are included in a large publicly-released e-commerce dataset [26]. Inspecting the ratings for these product subsets only, the results were almost completely identical to those reported for all products (Table 12 and Figure 4) and are thus not separately reported.

As a second (and final) step, we conducted another evaluation, in which all products in the public dataset [26] with at least 10 reviews between January and July 2014 (the most recent months in the dataset) were evaluated end-to-end using the LexRank method with $k=5$ sentences, by 3 professional annotators. This set included 140 Fashion products and 120 Motors products. Results, depicted in the upper section of Table 16, are very similar ($\pm 1.5\%$) to those reported for our own datasets in Table 12 and Figure 4.

6.2.3 Required number of reviews. Our end-to-end approach requires a given product to have a certain number of reviews, so that we can produce a description. In this section, we examine the number of reviews and review sentences required to produce a description using our method. To this end, we sampled uniformly at random 300 products from the Fashion domain and 300 products from the Motors domain. For each product, we considered its reviews in a random order and inspected how many reviews and review sentences were required to produce a description of $K \in \{3, 5, 7\}$ sentences. Recall that in order to produce a description of K sentences, we need to identify enough good sentences, such that K of them will be diverse enough according to the similarity threshold θ , as detailed in Section 5. The values selected for θ , as described, are 0.73 for Fashion and 0.76 for Motors. Table 13 presents the average number of reviews and review sentences, per domain, which were required to produce a description of length $K \in \{3, 5, 7\}$.

It can be seen that in order to generate a description of 3 sentences, the average number of reviews needed was around 16, while for generating a description of 7 sentences it was between 23 reviews for Fashion and 25 for Motors. For popular products, as in our datasets, there are enough reviews to generate such descriptions. However, for long-tail or brand new products, there might not be enough reviews to produce an appropriate description from the crowd. Our approach is therefore suitable for boosting products that already gained some popularity, by improving their user experience and visibility on search engines.

6.3 Coverage Analysis

6.3.1 Review Coverage. In this section, we set out to explore how well the selected sentences span the existing information in the reviews. Given a generated description (for $K \in \{3, 5, 7\}$), we wish to assess how distant from it are the *good* review sentences not selected for the description. To this end, for each *good* candidate not included in the generated description, we measured its minimal distance to a sentence in the description, i.e., the distance to the most similar sentence in the description. The distance $d=1-s$ is calculated as previously explained, based on the similarity between the weighted embedding representations of the sentences. Table 14 presents the distribution of this minimal distance, for the LexRank method, across all the sentences labeled *good*, but not included in the description, for both Fashion and Motors, across three bins. The first bin ranges from 0 to $\gamma=1-\theta$, where θ is the

Table 14. Distribution of *good* sentences not included in the description according to their minimum distance from a sentence in the description. Distribution is partitioned into three bins according to the distance threshold γ .

K	Fashion			Motors		
	$0 - \gamma$	$\gamma - 2\gamma$	$2\gamma - 1$	$0 - \gamma$	$\gamma - 2\gamma$	$2\gamma - 1$
3	40.2%	41.6%	18.2%	37.1%	42.8%	20.1%
5	55.3%	37.3%	7.4%	52.2%	40.2%	7.6%
7	62.3%	32.4%	5.3%	62.1%	31.7%	6.2%

similarity threshold calculated for Fashion and Motors, as explained in Section 5. Intuitively, review sentences in this bin are similar to at least one sentence in the description, and are thus represented in the description. The second bin ranges from γ to 2γ , while the third ranges from 2γ to 1 (the maximum distance value). Sentences in the third bin are particularly distant from all description sentences.

As can be seen in Table 14, when $K=3$, nearly 40% of the sentences not included in the description fall into the $[0-\gamma]$ bin, indicating they are “covered” by the description. Naturally, as the number of sentences in the description grows, this portion increases, up to about 62% for both Fashion and Motors for $K=7$. On the other hand, the portion of sentences within the $[2\gamma-1]$ bin decreases from around 20% for $K=3$ to 5%–6% for $K=7$. These changes are noticeably more substantial, however, from $K=3$ to $K=5$ than from $K=5$ to $K=7$, both in terms of the increase in the $[0-\gamma]$ bin and the decrease in the $[2\gamma-1]$ bin. It can also be observed that the coverage is similar for Fashion and Motors, with a slight advantage to the former. Finally, we note that calculating the same distribution for the three other methods indicated that LexRank had the best coverage out of all four method, with the highest portion of the sentences under the first bin and lowest under the third bin. At close second across all K values was the k -means centroid method, while both the greedy and k -means score method reaching a substantially lower coverage. This gives another indication to the superiority of the LexRank method followed by the k -means centroid method, as was also reflected by the ratings.

6.3.2 Comparison with the Original Description. We used the same type of coverage analysis to compare the generated descriptions of the products in each of the two domains, with the original descriptions as included in the respective datasets. Recall that the original descriptions are substantially longer than those we generate, as described in Section 3. The upper section of Table 15 presents the distribution of original description sentences relative to the generated descriptions. It can be seen that only 12.5%–13% of the original description sentences fall into the $[0-\gamma]$ bin for $K=3$ up to around 20% for $K=7$. This can be explained by the length of the original description, spanning a large number of sentences. In addition, as already mentioned, the original descriptions include content that does not appear in the generated descriptions, such as marketing statements (e.g. “*Our genius designers did all they could, making this trifold wallet slim & compact enough*”) or seller-specific notes (e.g., “*This item is not eligible for international shipping*”).

The lower section of Table 15 shows the distribution of generated sentences relative to the original description sentences. It can be seen that a substantial portion, from over a third for $K=3$ to over a half for $K=7$, are similar to at least one sentence in the original description ($[0-\gamma]$ bin). For example, the sentence “*Steering wheel covers fit middle size steering wheels*” appeared in our generated description and was similar to the sentence “*the covers fit in medium-sized steering wheel*” in the original description. On the other hand, at least one sentence in the generated description falls within the $[2\gamma-1]$ bin, indicating it is very different than any sentence in the original description. This indicates that our crowd-based descriptions capture content that is similar to the seller-provided content, but also sentences that do not tend to appear in classic descriptions. Examples for the latter include

Table 15. Distribution of sentences according to their minimum distance from a sentence in the generated description. Distribution is partitioned into three bins according to the distance threshold γ .

Sentence set	K	Fashion			Motors		
		$0 - \gamma$	$\gamma - 2\gamma$	$2\gamma - 1$	$0 - \gamma$	$\gamma - 2\gamma$	$2\gamma - 1$
Original description sentences	3	13.2%	68.2%	18.6%	12.5%	67.2%	20.3%
	5	17.2%	72.2%	10.6%	16.1%	70.9%	13.0%
	7	20.1%	72.4%	7.5%	19.8%	70.1%	10.1%
Generated description sentences	3	36.3%	28.9%	34.8%	37.7%	30.2%	32.1%
	5	48.2%	30.2%	21.6%	46.3%	31.2%	22.5%
	7	52.5%	31.0%	16.5%	51.2%	32.1%	16.7%

Table 16. Average rating of 5-sentence descriptions generated using the LexRank method over public data for Fashion and Motors (upper section) and for Toys and Electronics using cross-domain learning (lower section).

	Overall	Readable	Informative	Objective	Relevant
Fashion	4.34	3.88	4.20	4.40	4.53
Motors	4.32	3.95	4.22	4.43	4.50
Toys	4.05	3.62	4.09	4.21	4.27
Electronics	4.11	3.60	4.05	4.18	4.21

“Definitely fits well without restricting movement”, *“You can fit 2 full standard size storage bins full of clothes in one bag”*, or *“No shrinkage in washer or dryer”*.

6.4 Cross-Domain Generation

We also set out to explore if our approach can be used across domains. To this end, we trained a model based on both the Fashion and Motors data and generated descriptions for two additional e-commerce domains: Toys and Electronics. As in the previous experiment, we generated descriptions of 5 sentences using the LexRank method for all products in these domains with at least 10 reviews between January and July 2014, resulting in a total of 55 products for Toys and 75 for Electronics [26]. These descriptions were evaluated by the same 3 annotators and the lower section of Table 16 presents their ratings. It can be seen that while the overall quality ratings are somewhat lower than for Fashion and Motors, they are still just above 4 on average. There is a rather similar rating decrease across all four quality aspects relative to Fashion and Motors, with the sharpest decline being for readability. The last two examples in Table 11 demonstrate generated descriptions for Toys and Electronics.

Overall, the cross-domain description generation results indicate that applying a model trained jointly on two domains to other domains works reasonably well. This implies that characterizing qualities of review sentences suitable to be used as part of a description can be generalized across domains. To further explore this, we set out to examine what portion of the review sentences in the ‘*target*’ domains (Toys and Electronics) in our dataset have a similar sentence in the review set of the ‘*base*’ domains used for training (Fashion and Motors). To this end, for each target and base domain pairs, we examined the distribution of all review sentences in the target domain according to their minimum distance from a sentence out of all review sentences in the base domain. Table 17 presents the distribution of the minimum distance of target domain review sentences from a sentence in the base

Table 17. Distribution of all review sentences in Toys and Electronics according to their minimum distance from a review sentence in Fashion and Motors. Distribution is partitioned into three bins according to the distance threshold γ

	$0 - \gamma$	$l - 2\gamma$	$2\gamma - 1$
Toys vs. Fashion	46.2%	44.5%	9.3%
Toys vs. Motors	44.2%	51.4%	4.4%
Electronics vs. Fashion	35.9%	57.8%	6.3%
Electronics vs. Motors	52.4%	41.9%	5.7%

Table 18. Example of similar sentences in different domain pairs

	Sentence 1	Sentence 2
Toys vs. Fashion	It is also very durable and colorful Perfect for sports and outdoor activity Holds up to multiple washings	Very colorful and durable product Great for football and other sports You can wash it many times
Toys vs. Motors	Works great on all batteries Very stable and super silent Comes with a nice pouch to keep it clean	Accepts different types of batteries Super dependable and really quiet Easy for clean up and storage
Electronics vs. Fashion	The material is very durable Not too tight but also not baggy The lenses are crystal clear	Made of high quality material Fits without feeling too tight Clear lenses, and easy to focus
Electronics vs. Motors	It is really easy to use and charges pretty fast Includes many types of cables Lightweight and very compact compressor	Very useful and charges fast Many connections, such as USB, VGA and DVI The projector is very portable

domain reviews across three buckets, as described in Section 6.3.1, across all four target-base pairs. It can be seen that a substantial portion of the review sentences in Toys and Electronics have a similar sentence (distance of γ or lower or, in other words, similarity of θ or higher, as described in Section 6.3.1) in the base domain, Fashion or Motors. The highest portion in the $[0 - \gamma]$ bucket is for the Electronics domain, when matched against Motors. In this case, over half (52.4%) of the Electronics review sentences have at least one similar sentence in the Motors review. The lowest portion is for the Fashion domain against Electronics, at 35.9%. Only a small portion of the review sentences “fall” in the third bucket, with all review sentences in the base domain being very different (a distance of at least 2γ from each review sentence in the base domain). In Table 18, we listed a few examples for similar sentences (distance between 0 and γ) per each of the four domain pairs. Overall, our analysis demonstrates that review language often poses similar patterns across domains, which can be leveraged, as we have shown, for applying our description generation model across domains.

6.4.1 Transfer Learning for Cross-Domain Descriptions Generation. Transfer learning is a method used to share “knowledge” acquired while solving one problem and re-applying it to a different but somewhat related problem. Recent research by Mou et al. [57] studied the transfer learning quality for various NLP applications. The authors of the research concluded that the ability of a neural network to be effectively transferable depends on how semantically close the source and the target tasks are. In our case, as discussed above, the domains share some similarity hence we tried to apply a transfer learning technique. To this end, since we had training data only in

Table 19. Average rating of 5-sentence descriptions generated using the transfer learning approach for Toys and Electronics from Fashion and Motors domains, respectively.

	Overall	Readable	Informative	Objective	Relevant
Toys (from Fashion)	4.10	3.72	4.14	4.25	4.30
Electronics (from Motors)	4.22	3.70	4.11	4.24	4.25

Table 20. A/B test results - change in Views

	Popular Products		Random Products	
	Treatment	Control	Treatment	Control
Fashion	+23.8%	+2.4%	+11.7%	+4.1%
Motors	+12.3%	+5.3%	+9.4%	+3.9%

Fashion and Motors domains, we first performed low-resource tagging of review sentences in the target domains: 250 sentences in Toys and 250 in Electronics (which account for 1% of the labeled data we had in the source domain). We then used the previously-trained models in the Fashion and Motors domains (25K labeled examples in each of the domains), and preformed low-resource domain adaption [82]. This is the most suitable solution for domain adaptation in cases where there are two similar tasks and little training data in the target domain [5]. Concretely, we reused the trained Fashion (Motors) model, and performed additional training over the 250 labeled sentences in the Toys (Electronics) domain. This additional training allowed us to fine-tune the model weights and capture language nuances in the target domain. After building the models, we performed the end-to-end description generation process, and asked our experts to manually evaluate the results over the 55 products in Toys and 75 products in Electronics. Table 19 presents the results.

It can be seen that the results using transfer learning are slightly better than applying the model trained on both Fashion and Motors without any transfer learning (as depicted in the lower part of Table 16) and very close to the results of the domains that were trained using much more in-domain data (as depicted in the upper part of Table 16). Overall, we observe that using labeled data across domains is quite effective and can save substantial labeling efforts while leveraging already-labeled data in specific domains.

7 A/B TESTING IN PRODUCTION

To examine the effect of our generated crowd-based description “in vivo”, we performed a small-scale A/B test [37] in production for both the Fashion and Motors domains on the eBay US website. Our goal was to measure changes in traffic, as product page views are highly correlated with purchases and ultimate revenue. To this end, we selected 100 products for evaluation in each domain. Half of the products in each domain were *popular* products, among the top 0.5% according to daily page views, while the other half were random products, sampled uniformly at random from all products in the domain. For each product in these product groups, we presented the added the generated description to the “About this product” section in the product page titled as “Description from our customers”. Figure 5 illustrates a sample product page enriched by a generated description. For each of the sampled products, we set out to explore the change in traffic as a result of the description addition. To this end, we compared the traffic, in terms of product page views, in the 30 days preceding the addition of the crowd-based description to the page, with the number of page views in the following 30 days. For the control group, we

Shop by category

Search for anything

All Categories

Search

Advanced

eBay > Clothing, Shoes & Accessories > Men's Clothing > Activewear > Activewear Tops
Share

Under Armour Cold Gear Compression Mock (1265648) XL 100white

★★★★★ 2 product ratings | About this product

Brand new
ILS 240.79

New (other)
ILS 199.32

Brand new: lowest price ⓘ
Approx.
ILS 148.22
+ ILS 92.57 Shipping
US \$39.95

eBay MONEY BACK GUARANTEE

Buy It Now

Add to cart

Watch

About this product

Description from customers reviews:

Great for colder weather. Definitely fits well without restricting movement. Holds up to multiple washings. It does keep your body heat in. The quality of under armour is wonderful.

Product Identifiers	
GTIN	0888376860592
BRAND	Under armour
MPN	1265648

Show More ▾

Fig. 5. Product description evaluation interface.

sampled a similar number of products (100 in each domain, half popular, half random) of the same type from the same leaf categories, with a similar number of page views (+/−5%) in the 30 days preceding our test.

Table 20 presents the average change in product page views along the two 30-day periods for the “treatment” group (the products whose pages were enriched with the generated descriptions) and the control group. For the control group, there is a slight uplift in page views, indicating the period preceding the test’s start date generally enjoyed a slightly higher traffic. For both popular and random products, in both the Fashion and Motors domains,

the uplift across the treatment group's products was significantly higher than for the control group.⁷ The uplift boost for the treatment group was stronger for Fashion products, implying that textual descriptions are more important for products in this domain. The uplift boost was also stronger for popular products, with a high +23.8% for popular Fashion products (compared to +2.4% for the control group). Overall these results show a clear impact of the crowd-based description on product page traffic. This can stem from a variety of sources, one of them is the improvement of product occurrence on both internal and external search, due to the additional text, which adds new information relevant to the product. In addition, the better user experience enabled by the addition of the descriptions can lead users to visit product pages more frequently. These results complement the description quality evaluation, showing that our generated descriptions are not only of high quality, but can actually lead to a significant impact on user behaviour within an e-commerce website.

8 DISCUSSION AND LIMITATIONS

Our approach for extracting product descriptions requires an initial set of reviews for the product, as detailed in Section 6.2.3. It follows that for products with few reviews or none at all, our method suffers from a cold start problem [68]. That is, while our goal is to promote products by providing them with better descriptions, which, in turn, can enhance the user experience for the product page and its visibility on search engines, it will not be applicable for products that have received little exposure. Indeed, our approach aims to boost products after they have attracted some attention reflected in reviews from buyers. As discussed in Section 6.2.3, 23-25 reviews or 109-112 sentences were needed, on average, to generate a description of 7 suitable diverse sentences. One could think of mitigating this cold start problem by providing crowd-based descriptions based on reviews from similar products; however, we believe this approach carries risk, as similar products may pose different qualities (e.g., a different camera for similar cellphones), which may render misleading descriptions. Future work can examine how to identify descriptive sentences that can be shared between similar products. Another more practical approach would be to promote a fresh product by other means, such as targeted ads or homepage recommendation; after gaining some momentum and reviews, our approach can further boost the product, improving its "findability" and page experience. Finally, as we will demonstrate later in this section, simple rules can be applied to modify review sentences so they can fit a description and, as a result, reduce the total number of reviews and sentences needed to produce a proper product description.

For popular products, with an established number of reviews, our method enhances the buyer experience, by providing a description from a past buyers' viewpoint. The impact of this enhanced experience can be substantial, since popular products attract high numbers of potential buyers. Indeed, our A/B testing indicates a particularly significant effect on page views for popular products that included our descriptions. While we use sentences that already exist in reviews, we surface information that may require traversing a vast number of reviews, with many other types of content, such as personal experiences and subjective opinions. Moreover, the description part of the product page is indexed by external search engines, while the reviews might be excluded due to their large volume and content quality. Since we work at the sentence level, rather than the review level, common approaches for review triage, such as by the review's "helpfulness" votes, are not applicable for our task. We distill the descriptive sections hidden within reviews to provide a concise and diverse crowd-based description. Our approach is therefore different than methods commonly applied for review summarization, which aim to identify the most common aspects in a product's set of reviews [29, 49, 79].

Our description generation method is extractive: it is based on identifying the relevant sentences in reviews and drawing them out to create a crowd-based description. Beyond extraction, other techniques can be used to generate the description. For example, rules can be used to remove subjective parts of sentences that include viable descriptive information, to avoid their disqualification. Language adaptations can be applied to alter sentences

⁷statistical significance measured using one-tailed paired t-test.

Table 21. Example review sentences from our dataset that can be altered using rules to fit a description. ‘<BOS>’ represents the beginning of a sentence.

Original sentence	Modified sentence	Rule scheme
I recommend them to anyone for hiking and hitting the hills, certainly helps save the feet.	Recommended to anyone for hiking and hitting the hills, certainly helps save the feet.	<BOS>I recommend it/this/them/these → Recommended
The lenses are crisp and clear, and I like how the arms “lock” in place.	The lenses are crisp and clear.	Remove the part of sentence starting with “and I”
Also they are be easy to launder and dry.	They are be easy to launder and dry.	Remove transitional or conjunctive adverbs in the beginning of a sentence
In addition, they are mostly made of cotton but Gold Toe adds stretch nylon and a little bit of spandex.	They are mostly made of cotton but Gold Toe adds stretch nylon and a little bit of spandex.	
It will remove the dust from hard water and rain deposits.	This product will remove the dust from hard water and rain deposits.	<BOS>It → This product
I use it on exhaust sensor threads and spark plugs with no issue.	Can be used on exhaust sensor threads and spark plugs with no issue.	<BOS>I use it → Can be used

that are not suitable to use “as is” in a description, e.g., to avoid a missing context or to move from a subjective to objective phrasing. Table 21 lists a few examples for such rules that can transform sentences from our dataset into variants that can be used in a product description. More advanced natural language processing techniques, such as abstractive summarization and language generation, can also be applied to create a product description with review sentences in a different form than the original one that appears in the product’s review [15, 16].

Our supervised deep multi-task classifier was trained using a plethora of labeled sentences, annotated as suitable or not to be part of the product’s description. Our experiments indicated that up to 90% of the data can be spared, for a rather minor trade-off in performance. While the annotation task is simple and can be completed in a relatively short time (annotators completed 60 sentences per hour, on average), with a large inter-annotator agreement, it is always desirable to spare labeling efforts. We examined a cross-domain approach that requires no additional domain-specific data and demonstrated good results. Our initial experimentation with transfer learning from one domain, with a large amount of training data, to another, showed further improvement in the cross-domain performance, while using only little amount of labeled examples in the target domains. These results are especially promising for large e-commerce platforms, which span many different domains and categories. As our analysis shows, much of the language that distinguishes review sentences that can be used for a description from all other sentences is common across two very different domains: Fashion and Motors. Other semi-supervised learning and/or distant supervision techniques may be used to further reduce labeling efforts for this task.

Our approach selects the top sentences to include in the product’s description, but does not address the optimization of their order. As can be seen in the examples provided in Table 11, the descriptions often flow fairly well without special ordering. This could be due to the fact that standalone sentences were selected (missing context was the second most common reason for *bad* sentences, as detailed in Table 3) and also as we excluded similar sentences to one another. On the other hand, readability was the lowest-rated aspect of the generated

description, implying that there is still room for improving the description's flow, by optimizing the sentence order and connectivity.

9 CONCLUSION AND FUTURE WORK

In this paper, we presented a method for automatic generation of product descriptions based on online user reviews. Product descriptions provided by sellers are often missing or lacking. Hence, customers turn to product reviews and often spend a significant amount of time sifting through them before making a purchase. However, users might still find it hard to extract the relevant information from thousands of reviews available online. Moreover, many reviews include personal and subjective opinions, while the users are sometimes only interested in the key product details before purchasing [17].

We first studied the main differences between the reviews and descriptions and adopted a supervised deep multi-task learning approach to identify appropriate review sentences. Afterwards, we introduced a similarity measure between review sentences that helped us eliminate redundancies when creating the descriptions. Finally, we used extractive text summarization to create a coherent and concise product description. We provided an extensive set of experiments that demonstrated our approach is productive, including an A/B test in our production environment that showed increase in traffic.

Our experiments inspected two principally different e-commerce domains: Fashion and Motors. The similar results throughout our evaluation for both domains, as well as the analysis revealing common review language characterizing sentences that can be used for a description, suggest that our approach can be generalized to additional domains. Indeed, cross-domain description generation, trained over a combination of Fashion and Motors reviews, showed promising results when applied to two other primary e-commerce domains: Toys and Electronics. Our evaluation and examples also illustrate that our descriptions are somewhat different than seller-provided descriptions. A presentation of such descriptions on product pages may include an indication they are crowd-based or stem from user reviews. In addition, the descriptions can be presented as a bulleted list, as currently done on some e-commerce websites, to provide an alternative user experience and mitigate readability issues. In our A/B test, we presented a user interface that can be used to present review-based descriptions.

For future work, we plan to focus on four main directions. Currently, we select K sentences for the description without any special ordering. We believe that proper ordering may increase the readability of the descriptions. Second, automatically deriving the value of K based on the product and review content can further improve description quality. Third, abstractive approaches [15, 16], which adapt the original review text and combine content from different sentences, are also worth exploring, and can help generate descriptions of even higher quality based on fewer reviews. Finally, personalizing the generated descriptions for the individual consumer can help make them even more compelling [13].

REFERENCES

- [1] Palakorn Achananuparp, Xiaohua Hu, and Xiajiong Shen. 2008. The evaluation of sentence similarity measures. In *Proc. of DaWaK*. 305–316.
- [2] Mathieu Acher, Anthony Cleve, Gilles Perrouin, Patrick Heymans, Charles Vanbeneden, Philippe Collet, and Philippe Lahire. 2012. On extracting feature models from product descriptions. In *Proc. of VaMoS*. 45–54.
- [3] Sanjeev Arora, Yingyu Liang, and Tengyu Ma. 2017. A simple but tough-to-beat baseline for sentence embeddings. In *Proc. of ICLR*.
- [4] Adam Berger and John Lafferty. 1999. Information retrieval as statistical translation. In *Proc. of SIGIR*. 222–229.
- [5] Avi Bleiweiss. 2019. LSTM Neural Networks for Transfer Learning in Online Moderation of Abuse Context. In *ICAART*. SciTePress, 112–122.
- [6] Rich Caruana. 1998. Multitask learning. In *Learning to learn*. Springer, 95–133.
- [7] Y. Chae, M. Nakazawa, and B. Stenger. 2018. Enhancing product images for click-through rate improvement. In *Proc. of ICIP*. 1428–1432.
- [8] Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. In *Proc. of KDD*. 785–794.
- [9] Judith A. Chevalier and Dina Mayzlin. 2006. The effect of word of mouth on sales: Online book reviews. *Journal of Marketing Research* 43, 3 (2006), 345–354.

- [10] Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement* 20, 1 (1960), 37–46.
- [11] Jacob Cohen. 1968. Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial credit. *Psychological bulletin* 70, 4 (1968), 213–220.
- [12] Horatiu Dumitru, Marek Gibiec, Negar Hariri, Jane Cleland-Huang, Bamshad Mobasher, Carlos Castro-Herrera, and Mehdi Mirakhorli. 2011. On-demand feature recommendations derived from mining public product descriptions. In *Proc. of ICSE*. 181–190.
- [13] Guy Elad, Ido Guy, Slava Novgorodov, Benny Kimelfeld, and Kira Radinsky. 2019. Learning to Generate Personalized Product Descriptions. In *Proc. of CIKM*. 389–398.
- [14] Günes Erkan and Dragomir R Radev. 2004. Lexrank: Graph-based lexical centrality as salience in text summarization. *Journal of artificial intelligence research* 22 (2004), 457–479.
- [15] Kavita Ganesan, ChengXiang Zhai, and Jiawei Han. 2010. Opinosis: A graph-based approach to abstractive summarization of highly redundant opinions. In *Proc. of COLING*. 340–348.
- [16] Shima Gerani, Yashar Mehdad, Giuseppe Carenini, Raymond T. Ng, and Bitu Nejat. 2014. Abstractive summarization of product reviews using discourse structure. In *Proc. of EMNLP*. 1602–1613.
- [17] Anindya Ghose and Panagiotis G. Ipeirotis. 2011. Estimating the helpfulness and economic impact of product reviews: Mining text and reviewer characteristics. *IEEE Transactions on Knowledge and Data Engineering* 23, 10 (2011), 1498–1512.
- [18] Ross Girshick. 2015. Fast r-cnn. In *Proc. of ICCV*. 1440–1448.
- [19] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. Domain adaptation for large-scale sentiment classification: A deep learning approach. In *Proc. of ICML*. 513–520.
- [20] Jade Goldstein, Mark Kantrowitz, Vibhu Mittal, and Jaime Carbonell. 1999. Summarizing text documents: sentence selection and evaluation metrics. In *Proc. of SIGIR*. 121–128.
- [21] Vishrawas Gopalakrishnan, Suresh Parthasarathy Iyengar, Amit Madaan, Rajeev Rastogi, and Srinivasan Sengamedu. 2012. Matching product titles using web-based enrichment. In *Proc. of CIKM*. 605–614.
- [22] Anjan Goswami, Naren Chittar, and Chung H. Sung. 2011. A study on the impact of product images on user clicks for online shopping. In *Proc. of WWW*. 45–46.
- [23] Ido Guy, Avihai Mejer, Alexander Nus, and Fiana Raiber. 2017. Extracting and ranking travel tips from user-generated reviews. In *Proc. of WWW*. 987–996.
- [24] Ido Guy and Bracha Shapira. 2018. From Royals to Vegans: Characterizing Question Trolling on a Community Question Answering Website. In *Proc. of SIGIR*. 835–844.
- [25] Kazuma Hashimoto, Caiming Xiong, Yoshimasa Tsuruoka, and Richard Socher. 2016. A joint many-task model: Growing a neural network for multiple NLP tasks. *arXiv preprint abs/1611.01587* (2016).
- [26] Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *Proc. of WWW*. 507–517.
- [27] Sharon Hirsch, Ido Guy, Alexander Nus, Arnon Dagan, and Oren Kurland. 2020. Query Reformulation in E-Commerce Search. In *Proc. of SIGIR*. 1319–1328.
- [28] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [29] Mingqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proc. of KDD*. 168–177.
- [30] Nan Hu, Paul A. Pavlou, and Jennifer Zhang. 2006. Can online reviews reveal a product’s true quality?: Empirical findings and analytical modeling of online word-of-mouth communication. In *Proc. of EC*. 324–330.
- [31] Alice Jiang, Zhilin Yang, and Minjoon Jun. 2013. Measuring consumer perceptions of online shopping convenience. *Journal of Service Management* 24, 2 (2013), 191–214.
- [32] Gagandeep Kaur and Gagandeep Kaur. 2016. Mobile applications are major players in the world of e-commerce. *International Journal of Advanced Research in IT and Engineering* 5, 2 (2016), 13–21.
- [33] Zehra Kavasoglu and Şule Gündüz Ögüdücü. 2013. Personalized summarization of customer reviews based on user’s browsing history. *IADIS International Journal on Computer Science & Information Systems* 8, 2 (2013), 147–158.
- [34] H. Khalid, E. Shihab, M. Nagappan, and A.E. Hassan. 2015. What do mobile app users complain about? *IEEE Software* 32, 3 (2015), 70–77.
- [35] Hyun Duk Kim, Kavita Ganesan, Parikshit Sondhi, and Chengxiang Zhai. 2011. Comprehensive review of opinion summarization. *UIUC Technical Report* (2011).
- [36] Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint abs/1412.6980* (2014).
- [37] Ron Kohavi and Roger Longbotham. 2017. Online Controlled Experiments and A/B Testing. *Encyclopedia of machine learning and data mining* 7, 8 (2017), 922–929.
- [38] Siwei Lai, Liheng Xu, Kang Liu, and Jun Zhao. 2015. Recurrent convolutional neural networks for text classification.. In *Proc. of AAAI*. 2267–2273.
- [39] Eun-Ju Lee and Soo Yun Shin. 2014. When do consumers buy online product reviews? Effects of review quality, product type, and reviewer’s photo. *Comput. Hum. Behav.* 31 (2014), 356–366.
- [40] Beibei Li, Anindya Ghose, and Panagiotis G Ipeirotis. 2011. Towards a theory model for product search. In *Proc. of WWW*. 327–336.

- [41] Fangtao Li, Chao Han, Minlie Huang, Xiaoyan Zhu, Ying-Ju Xia, Shu Zhang, and Hao Yu. 2010. Structure-aware review mining and summarization. In *Proc. of COLING*. 653–661.
- [42] Xinxin Li and Lorin M Hitt. 2008. Self-selection and information role of online product reviews. *Information Systems Research* 19, 4 (2008), 456–474.
- [43] Moez Limayem, Mohamed Khalifa, and A. Frini. 2000. What makes consumers buy from Internet? A longitudinal study of online shopping. *IEEE Trans. Systems, Man, and Cybernetics, Part A* 30, 4 (2000), 421–432.
- [44] Chin-Yew Lin and Eduard Hovy. 2002. From single to multi-document summarization: A prototype system and its evaluation. In *Proc. of ACL*. 457–464.
- [45] Bing Liu. 2012. *Sentiment analysis and opinion mining*. Morgan & Claypool Publishers.
- [46] Qian Liu, Zhiqiang Gao, Bing Liu, and Yuanlin Zhang. 2015. Automated rule selection for aspect extraction in opinion mining. In *Proc. of IJCAI*. 1291–1297.
- [47] Roque Enrique Lpez Condori and Thiago Alexandre Salgueiro Pardo. 2017. Opinion summarization methods. *Expert Syst. Appl.* 78, C (July 2017), 124–134.
- [48] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. *arXiv preprint abs/1508.04025* (2015).
- [49] Duy Khang Ly, Kazunari Sugiyama, Ziheng Lin, and Min-Yen Kan. 2011. Product review summarization from a deeper perspective. In *Proc. of JCDL*. 311–314.
- [50] Dehong Ma, Sujian Li, Xiaodong Zhang, and Houfeng Wang. 2017. Interactive attention networks for aspect-level sentiment classification. In *Proc. of IJCAI*. 4068–4074.
- [51] James MacQueen et al. 1967. Some methods for classification and analysis of multivariate observations. In *Proc. of Berkeley symposium on mathematical statistics and probability*, Vol. 1. 281–297.
- [52] Deborah Brown McCabe and Stephen M Nowlis. 2003. The effect of examining actual products or product descriptions on consumer preference. *Journal of Consumer Psychology* 13, 4 (2003), 431–439.
- [53] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint abs/1301.37810* (2013).
- [54] Hye-Jin Min and Jong C Park. 2012. Identifying helpful reviews based on customer’s mentions about experiences. *Expert Systems with Applications* 39, 15 (2012), 11830–11838.
- [55] Samaneh Moghaddam and Martin Ester. 2012. On the design of LDA models for aspect-based opinion mining. In *Proc. of CIKM*. 803–812.
- [56] Ajinkya More. 2016. Attribute extraction from product titles in eCommerce. *arXiv preprint abs/1608.04670* (2016).
- [57] Lili Mou, Zhao Meng, Rui Yan, Ge Li, Yan Xu, Lu Zhang, and Zhi Jin. 2016. How transferable are neural networks in nlp applications? *arXiv preprint arXiv:1603.06111* (2016).
- [58] Quang Nguyen. 2012. *Detecting experience revealing sentences in product reviews*. Ph.D. Dissertation. University of Amsterdam.
- [59] Slava Novgorodov, Ido Guy, Guy Elad, and Kira Radinsky. 2019. Generating product descriptions from user reviews. In *Proc. of WWW*. 1354–1364.
- [60] Bo Pang and Lillian Lee. 2008. Opinion mining and sentiment analysis. *Found. Trends Inf. Retr.* 2, 1-2 (Jan. 2008), 1–135.
- [61] Eun Joo Park, Eun Young Kim, Venessa Martin Funches, and William Foxx. 2012. Apparel product attributes, web browsing, and e-impulse buying on shopping websites. *Journal of Business Research* 65, 11 (2012), 1583–1589.
- [62] A.M. Popescu and Oren Etzioni. 2005. Extracting product features and opinions from reviews. In *Proc. of HLT*. 339–346.
- [63] Katharina Probst, Rayid Ghani, Marko Krema, Andrew Fano, and Yan Liu. 2007. Semi-supervised learning of attribute-value pairs from product descriptions. In *Proc. of IJCAI*. 2838–2843.
- [64] Reid Pryzant, Young-Joo Chung, and Dan Jurafsky. 2017. Predicting sales from the language of product descriptions. In *Proc. of ECOM (SIGIR Workshops)*.
- [65] Pradeep Racherla, Munir Mandviwalla, and Daniel J Connolly. 2012. Factors affecting consumers’ trust in online product reviews. *Journal of Consumer Behaviour* 11, 2 (2012), 94–104.
- [66] Irina Rish. 2001. An empirical study of the naive Bayes classifier. In *IJCAI 2001 workshop on empirical methods in artificial intelligence*, Vol. 3. 41–46.
- [67] Sebastian Ruder. 2017. An overview of multi-task learning in deep neural networks. *arXiv preprint abs/1706.05098* (2017).
- [68] Andrew I. Schein, Alexandrin Popescul, Lyle H. Ungar, and David M. Pennock. 2002. Methods and Metrics for Cold-Start Recommendations. In *Proc. of SIGIR*. 253–260.
- [69] Keiji Shinzato and Satoshi Sekine. 2013. Unsupervised extraction of attributes and their values from product description. In *Proc. of ACL*. 1339–1347.
- [70] Anders Søgaard and Yoav Goldberg. 2016. Deep multi-task learning with low level tasks supervised at lower layers. In *Proc. of ACL*, Vol. 2. 231–235.
- [71] Krysta M. Svore, Lucy Vanderwende, and Chris J.C. Burges. 2007. Enhancing single-document summarization by combining RankNet and third-party sources. In *Proc. of EMNLP-CoNLL*.

- [72] Tijmen Tieleman and Geoffrey Hinton. 2012. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning* 4, 2 (2012), 26–31.
- [73] Hen Tzaban, Ido Guy, Asnat Greenstein-Messica, Arnon Dagan, Lior Rokach, and Bracha Shapira. 2020. Product Bundle Identification Using Semi-Supervised Learning. In *Proc. of SIGIR*. 791–800.
- [74] Damir Vandic, Flavius Frasincar, and Uzay Kaymak. 2018. A framework for product description classification in e-commerce. *Journal of Web Engineering* 17, 1&2 (2018), 001–027.
- [75] Shuai Wang, Zhiyuan Chen, and Bing Liu. 2016. Mining aspect-specific opinion using a holistic lifelong topic model. In *Proc. of WWW*. 167–176.
- [76] Wenya Wang, Sinno Jialin Pan, Daniel Dahlmeier, and Xiaokui Xiao. 2017. Coupled multi-layer attentions for co-extraction of aspect and opinion terms. In *Proc. of AAAI*. 3316–3322.
- [77] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. Hierarchical attention networks for document classification. In *Proc. of NAACL*. 1480–1489.
- [78] Koji Yatani, Michael Novati, Andrew Trusty, and Khai N. Truong. 2011. Analysis of adjective-noun word pair extraction methods for online review summarization. In *Proc. of IJACI*. 2771–2776.
- [79] Naitong Yu, Minlie Huang, Yuanyuan Shi, et al. 2016. Product review summarization by exploiting phrase properties. In *Proc. of COLING*. 1113–1124.
- [80] Di Zhu, Theodoros Lappas, and Juheng Zhang. 2018. Unsupervised tip-mining from customer reviews. *Decision Support Systems* 107 (2018), 116–124.
- [81] Li Zhuang, Feng Jing, and Xiao-Yan Zhu. 2006. Movie review mining and summarization. In *Proc. of CIKM*. 43–50.
- [82] Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation. *arXiv preprint arXiv:1604.02201* (2016).