
FashionVQA: A Domain-Specific Visual Question Answering System

Min Wang

Target Corporation
100 S. Mathilda Place Suite 300
Sunnyvale, CA 94086
min.wang@target.com

Ata Mahjoubfar

Target Corporation
100 S. Mathilda Place Suite 300
Sunnyvale, CA 94086
ata.mahjoubfar@target.com

Anupama Joshi

Target Corporation
100 S. Mathilda Place Suite 300
Sunnyvale, CA 94086
anupama.joshi@target.com

Abstract

Humans apprehend the world through various sensory modalities, yet language is their predominant communication channel. Machine learning systems need to draw on the same multimodal richness to have informed discourses with humans in natural language; this is particularly true for systems specialized in visually-dense information, such as dialogue, recommendation, and search engines for clothing. To this end, we train a visual question answering (VQA) system to answer complex natural language questions about apparel in fashion photoshoot images. The key to the successful training of our VQA model is the automatic creation of a visual question-answering dataset with 168 million samples from item attributes of 207 thousand images using diverse templates. The sample generation employs a strategy that considers the difficulty of the question-answer pairs to emphasize challenging concepts. Contrary to the recent trends in using several datasets for pretraining the visual question answering models, we focused on keeping the dataset fixed while training various models from scratch to isolate the improvements from model architecture changes. We see that using the same transformer for encoding the question and decoding the answer, as in language models, achieves maximum accuracy, showing that visual language models (VLMs) make the best visual question answering systems for our dataset. The accuracy of the best model surpasses the human expert level, even when answering human-generated questions that are not confined to the template formats. Our approach for generating a large-scale multimodal domain-specific dataset provides a path for training specialized models capable of communicating in natural language. The training of such domain-expert models, e.g., our fashion VLM model, cannot rely solely on the large-scale general-purpose datasets collected from the web.

1 Introduction

Fashion is about 2% of the world’s GDP and a significant sector of the retail industry. Whenever a new fashion item like apparel or footwear is launched, the retailer needs to prepare and show rich information about the product, including pictures, text descriptions, and detailed attribute tags. The attributes of the fashion products, including color, pattern, texture, material, occasion-to-use, etc., require domain experts to label them piece by piece. This labeling process is time-consuming, costly, subjective, error-prone, and fundamentally imprecise due to the interdependency of the attributes. To

address these issues, we introduce a multi-task multimodal machine learning model to automatically, consistently, and precisely infer the visual attributes of the fashion items.

Each item is typically labeled with multiple tags that describe different attributes of the item. For example, an item can be labeled with “shirt”, “red”, “solid pattern”, “blue collar” and “short sleeve”. An intuitive way of learning such information is to train a multi-label classifier, which outputs the probability of multiple labels of each input sample. However, such a model cannot encode the relationship between different attributes. For example, “short sleeve” is a suitable attribute for “shirt”, but not for “jeans”, and “red” only describes the body part of the shirt, but not the collar. The model needs to learn attribute and object relationships and adjusts its output accordingly.

We propose designing a Visual Questioning Answer (VQA) framework for fashion items, in which the model is trained to answer complex natural-language questions, such as “is the person wearing a red shirt with solid pattern and blue collar?”, given the input image. The VQA task is more challenging than the simple attribute classifier since it requires a thorough understanding of both the question and the structure and relationship between various visual attributes in the image. By training such a model, we convert the manual process of tagging new products with visual attributes into automated answering of a series of questions with visual intents (auto-labeling). The model also generates multimodal embeddings of the product images attended to the questions for downstream dialogue, search, and recommendation systems.

Prior to our work, there exists a large-scale VQA v2 dataset [1], which includes 0.6 million *question-answer-image* triplets. It has been widely used as the benchmark in recent research on VQA tasks. However, this general dataset only contains a small number of *question-answer-image* triplets related to fashion. In this work, we build a fashion VQA dataset from a diverse apparel product database. The questions, including both binary and non-binary, are automatically composed by filling question templates with the given attribute information. The dataset contains 207 thousand images and 168 million *question-answer-image* triplets. The automatic generation of the VQA dataset from a limited number of images and attributes allows us to achieve the scale required for training a multimodal domain expert model.

We leverage a cross-modality fusion model mapping representations from visual and text space to the same latent feature space and performing answer prediction with classifier modules. Given an image that contains a fashion item and the corresponding questions regarding its different attributes, the model predicts the answers to the given questions. We can then use the model to generate the missing or alternative attribute information based on its answers.

Additionally, given different but similar text descriptions on the same item, we can generate consistent feature embeddings that enable us to build better online search services. The existing search engines cannot attend to the relevant visual parts of a fashion item given the query and do not adapt the attention mask according to the chained adjectives. With this work, we can map the input query to the learned embeddings space and perform a robust and fuzzy search in that multimodal space. We can also provide a visual dialogue service, in which the customers can ask consecutive questions to narrow down the item list according to their apparel preferences. We can also build a fashion recommendation system in the multimodal embeddings space. The customer-item interaction history is mapped to this space, and the neighboring items are recommended.

2 Related work

Visual feature learning in VQA: The visual feature vector is often extracted from the input image using a Convolutional Neural Networks (CNN) as the visual encoder, e.g., VGG [2], ResNet [3], or ResNext [4] models. In the early VQA frameworks, grid-based visual features extracted by ImageNet-pretrained [5] VGG or ResNet models were widely adopted. Since [6], region-based visual features extracted using Faster R-CNN [7]-based object detection model, especially fine-tuned on Visual Genome [8] dataset, have been dominant [9][10][11][12]. In [13], the authors propose extracting the grid features from the same layer as the pretrained detector, achieving comparable performance as the region-based features with higher efficiency. We benchmark these two types of visual feature extraction methods on our dataset across different VQA models.

Cross-modality fusion models: Cross-modality fusion model is a core component of the VQA framework. It aligns the features from the visual and language modalities. Initially-proposed VQA

models identify the high-level cross-modal interactions by Bilinear Fusion [14]. MCB [15], MLB [16] and MUTAN [17] are later introduced to achieve better fusion performance at much lower computational cost and parameters. Motivated by the remarkable performance of the attention mechanism in language and vision models [18][19], the attention module becomes the fundamental block in designing the cross-modality fusion models. DFAF [9] uses self-attention and co-attention modules to learn the inter-and intra-connections between the two modalities. MCAN [10] builds the model with the blocks of self-attention and guided attention. LXMERT [12], ViBERT [20] adopt a similar strategy and build a two-stream co-attention-based model. VisualBERT [21], UNITER [22], OSCAR [11], OSCAR+ [23] learn the alignment between image and language by pretraining on multiple image caption datasets with BERT-style [18] visual language models (VLMs).

Fashion datasets: In recent year, many valuable fashion datasets [24][25][26][27][28][29][30][31][32][33][34] have greatly contributed to clothing item recognition and apparel attribute understanding. However, most of them suffer from some limitations when considered for training versatile VQA models. In [33][26][34][32], only primary categories in the dataset are labeled. Additional garment parts and attributes are annotated in [29][25][27]. Segmentation masks over each piece of garment are drawn for the semantic segmentation task in [26][33][25][30]. Generally, the localization of the garment pieces and parts in these types of datasets takes considerable human annotation labor, and few of the datasets are suitable for conversion to a new dataset for vision language tasks.

3 Methods

In this section, we describe how we designed and generated a novel VQA dataset for fashion. We named the new dataset FashionVQA dataset.

3.1 Terminologies

Category: Each clothing item can be labeled with one super-category and several primary categories or sub-categories.

- Super-categories: “apparel top”, “apparel bottom”, “one-piece clothing”, “shoes”, and “accessories”.
- Primary categories: “shirt”, “sweater”, “jacket”, “pants”, “skirt”, “dress”, “jumpsuit”, “boots”, “sneaker”, “gloves”, etc.
- Sub-categories: “t-shirt”, “cardigan”, “blazer jacket”, “pencil skirt”, “sweatpants”, “overall jumpsuit”, “hiking boots”, etc.

Attributes: “Color”, “pattern”, “fit type”, “closure type”, and “material/fabric” are the general attributes for all the fashion items. Each apparel type also has its unique attributes. The unique attributes for “apparel top”, “apparel bottom”, “one-piece clothing”, and “shoes” are listed below.

- Apparel top: “torso length type”, “sleeve length type”, “pocket type”, “neckline type”, “sleeve style”, “collar type”, “lapel type”, and “cuff type”.
- Apparel bottom: “pant leg type”, “skirt length type”, “pant leg style”, and “pleat type”.
- One-piece clothing: “neckline type”, “sleeve length type”, “sleeve style”, “pant leg type”, “skirt length type”, and “pleat type”.
- Shoes: “height”, “width”, “toe openness”, and “shape of toe”.

Attribute values: Each attribute is composed of a set of *attribute values*. For example, the set of the *attribute values* of color attribute includes “red”, “black”, “green”, “blue”, “yellow”, etc.

Parts: The *parts* mentioned in our dataset are typically lined on the fashion item, such as “patches” and “pockets”.

Location: In our dataset, there exist numerous images with a person wearing multiple fashion items. Therefore, we use *Location* to specify the relative location of the primary fashion item in the image, such as “on the top”, “on the bottom”, “on the feet”, “over the neck”, or “on the head”.

3.2 Data Collection pipeline

Our data collection pipeline involves four steps: [1] querying fashion items’ unique identity numbers (image IDs), [2] querying and parsing meta-information, [3] downloading images, and [4] filling question templates and forming *question-answer-image* triplets.

Each fashion item comes with a unique identity (ID) number. First, we query all fashion items and retrieve their IDs. Then, we predefine a data structure that is eligible to query the meta-information of fashion items from the item database. Feeding the data structure to an open-source data query API, “graphQL”, we can obtain the meta-information attached to each ID, which contains the primary image (front-view) URL and the description of the primary fashion item. We can directly download the primary image from the URL with Python.

The description of the fashion item is not an on-deck dictionary that maps each unique targeted attribute to its corresponding set of *attribute values*. For example, “Color” could be described with different phrases such as “Product Color” or “Color Name”. Parsing from the description is a process that collects *attribute values* from various sources and reduces similar attribute terminologies into the same group. Also, meta-information comes in a very raw manner with many *attribute values* cross *attributes* entangled, e.g., “black/stripes”, or in a vague expression. In this stage, we also need to clean these *attribute values* and map them into common terminologies, e.g., map “black/stripes” to “black” for color and “stripes” for pattern, or “olive night” color to “olive green”.

3.3 Question templates

We adopt a templating mechanism to automatically create *question-answer* pairs. The question templates are designed based on a set of fixed rules that meet the English grammar and result in human-readable sentences. By filling the question templates with specific item *attribute*, *attribute value*, *category*, and *location*, we can generate a variety of questions for each image. Answer of each question can be “Yes/No” for binary questions and multiple choices from the relevant *attribute values* for non-binary questions.

Since the images from the FashionVQA dataset are all photoshoot images with a solid background, the question templates ask only attribute-related questions about the fashion items in the image. For example, “what is the sleeve length of this shirt on the top?” or “is this a white v-neck sweater?”. The basic template is structured as “{*question type*} {this/these} {a/an/} {pair of/pairs of/} {*object*} {*location*}?”. When filling the template to expand into a full sentence, the choices between “is/are”, “this/these”, “a/an”, “a pair of/pairs of”, and singular or plural format of *category* are required to follow the English grammar and be aligned with the number of targeted fashion item in the image. For example, if the *number of pieces* in the image is more than one, we choose “are”, “these”, “/pairs of”, and plural format of *category*. If the fashion item includes pant legs or two pieces like eyeglasses, we add “pair of / pairs of” in the question templates.

If a person is in an image, and the primary fashion item is not from the super-category of “one-piece clothing”, we assume there are multiple fashion items in the image. We use “{*location*}” to specify the relative location of the primary fashion item. We use “on the top” for “apparel top”, “on the bottom” for “apparel bottom”, “on the feet” for “shoes”, “on the head” for “hat”, and “over the neck” for “scarf”.

The question templates fall into two primary categories based on the answer types: binary and non-binary templates.

Binary question templates: Binary question templates typically start with “is this/are these”, “can you see”, or “is there any {*part*} on this/these”, followed by the description of the targeted item in the format of “{*location*} {a}/ {a pair of/}/ {*attribute value 1*} {*attribute value 2*} {*category*}”, whereas *attribute value 1* and *attribute value 2* are two *attribute values* from different *attributes*. Permuting *attribute value 1*, *attribute value 2*, *category* in different orders yields different question templates. Conjunction words like “with”, “and”, or “in” can be used in templates when *attribute value 1* or *attribute value 2*, or both are located after *category*. The most common question types used in binary questions are “is/are” and “can”.

Non-binary question templates: Non-binary question templates typically start with question words like “what” / “why” / “when” / “how” followed by terms of attribute. The formats of the question type vary from attribute to attribute. For example, the question type can be “what color is” or “what

Question templates	Answer types	Question types	Questions
“is this a {attr1} {category} with {attr2}?”	“yes/no”	“is/are”	“is this a white shirt with long sleeves?”
“on the top a {category} with {attr1} and in {attr2} design?”	“yes/no”	“is/are”	“on the top a sweater with floral print and in v neck design?”
“what {attribute} is this {category} the person wearing {location}?”	“others”	“what {attribute}”	“what color is this a-line dress the person wearing on the top?”
“what {attribute} is the one {location}?”	“others”	“what {attribute}”	“what color is the one on the top?”
“when is a good time to wear this {attr1} {category}?”	“others”	“when”	“when is a good time to wear this yellow dress?”

Table 1: Question templates and examples in the FashionVQA dataset

is the color of” for attribute “color”, “what pattern is on” or “what print is on” for attribute “pattern”, and “how many” or “what number of” for attribute “number of pockets”.

Unlike binary question templates in our current dataset, we do not leverage other *attribute values* unrelated to the targeted attribute in filling the non-binary question templates; even the *category* of the targeted fashion item is not necessary. Therefore, it is possible to increase the diversity of the non-binary question templates with additional *attribute values* or *categories*. For example, we can come up with a color question template like “what color is on the top?” or “what color is this shirt the person wearing on the top?”.

Diversification: The primary question templates are those preserving all the *demonstratives*, *subject pronouns*, and *prepositional phrases*. By randomly either removing parts of those phrases or replacing them with alternatives, we can create assorted variant question templates.

In non-binary question templates, the question types for a given attribute come in different fashions, contributing to diverse non-binary question templates. Additionally, it is reasonable to replace the specific *category* information of the targeted item with the combination of *pronoun* and *location* to expand the diversity of question templates. Adding non-relevant *attribute values* to describing the fashion item is also an approach to creating new question sentences. To further increase the robustness of the question templates, we also introduce a small portion of noise into the question templates, switching between “this/these”, “is/are”, and “singular/plural”.

In binary question templates, even elimination of the question type phrases like “is this a”, “are these”, or “is there” does not cause an obstacle to make the remaining phrase human readable. Therefore, we truncate a small fraction of the full question sentences by removing phrases of question type to increase the diversity of the binary questions. When *attribute values* are placed after the *category*, we randomly pick one from different *conjunction structures* to form different phrases, which will remarkably increase the diversity of the binary question sentences. For example, for “a shirt with stripe pattern”, an alternative expression can be “a shirt designed with stripe pattern”, or “a shirt featured in stripe design”.

Table 1 demonstrates some examples of question sentences generated from question templates.

3.3.1 Balance positive and negative samples for each binary question

Given binary and non-binary question templates and *attribute values* for a specific image, we can easily generate non-binary *question-(multiple answers)-image* triplets and binary *question-(positive answer)-image* triplets.

For a balanced VQA dataset, we expect each binary question to come with the same number of positive and negative samples, i.e., balanced (*question*, “Yes”, *image ID*) triplets and (*question*, “No”, *image ID*) triplets. Here, we consider two different strategies for generating the negative samples of each binary question. One strategy keeps the image fixed and changes the *attribute values* in the question; the other one keeps the *attribute values* fixed and changes the image. Here we further explain these strategies in detail:

Image-based: For each image, by filling the binary question templates with specific *attribute values* and *category* information provided for this image, we make a positive binary sample. When an

attribute value or *category* in an existing binary question is changed, if the alteration is not in the list of *attribute values* or *categories* corresponding to the image, we assume this is a negative sample for the binary question.

Algorithm 1 Attribute-based balancing of the positive and negative samples for binary questions

Input: S: $\{s_i, \dots\}$ list of all fashion items
Each fashion item s_i : {image ID: u_i , category: c_i , attributes: $\{a_k, \dots\}$, attribute values: $\{v_{a_k}, \dots\}$
d)
 Q_T : list of binary question templates of all attributes
Output: B: list of binary *question-answer-image* triplets
Initialization:
for each specific attribute a_k **do**
 $U_{a_k} \leftarrow \{\}$: empty set of all image IDs with attribute a_k
 $V_{a_k} \leftarrow \{\}$: empty set of unique attribute values with attribute a_k
 $C \leftarrow \{\}$: empty set of unique categories
Build attribute-value-to-images dictionary:
for each fashion item $s_i \in S$ **do**
 $C \leftarrow c_i$
 for each attribute $a_k \in s_i(\text{attributes})$ **do**
 $U_{a_k} \leftarrow u_i$
 P_{c_i} (image ID set of positive answer of category c_i) $\leftarrow u_i$
 for each attribute value $v_{a_k} \in s_i(\text{attribute values})$ **do**
 $V_{a_k} \leftarrow v_{a_k}$;
 $P_{v_{a_k}}$ (image ID set of positive answer with attribute value v_{a_k}) $\leftarrow u_i$
Build attribute-value-to-(positive/negative answer)-images dictionary:
for each attribute value $v_{a_k} \in V_{a_k}$ **do**
 $V' = \text{Synonyms}(v_{a_k})$
 for each $v_+ \in (V' \cap V_{a_k})$ **do**
 $P_{v_{a_k}} = P_{v_{a_k}} \cup P_{v_+}$
 for each attribute value $v_{a_k} \in V_{a_k}$ **do**
 $N_{v_{a_k}} = U_{a_k} - P_{v_{a_k}}$
 for each category $c_i \in C$ **do**
 Follow the same strategy to update positive answer image ID set P_{c_i} and build negative set N_{c_i}
Expand attribute-value-to-(positive/negative answer)-images dictionary with attributes and category combinations:
for each attribute value $v_{a_k} \in V_{a_k}$ **do:**
 for each category $c_i \in C$ **do**
 $P_{(v_{a_k}, c_i)}$: positive answer image ID set of the combination of (v_{a_k}, c_i)
 $N_{(v_{a_k}, c_i)}$: negative answer image ID set of the combination of (v_{a_k}, c_i)
 $P_{(v_{a_k}, c_i)} = P_{v_{a_k}} \cap P_{c_i}$
 $N_{(v_{a_k}, c_i)} = (P_{v_{a_k}} \cap N_{c_i}) \cup (N_{v_{a_k}} \cap N_{c_i}) \cup (N_{v_{a_k}} \cap P_{c_i})$
Create balanced question-(positive/negative answer)-images triplets:
for each category $c_i \in C$ **do**
 for each specific attribute a_k of category c_i **do**
 $Q_{T(a_k, c_i)} = Q_T$ (binary question templates of attribute a_k and category c_i)
 for each combination of attribute value $v_{a_k} \in V_{a_k}$ and category $c_i \in C$ **do**
 $Q_{(v_{a_k}, c_i)} = \text{Fill } Q_{T(a_k, c_i)}$ templates with v_{a_k} and c_i to generate binary questions
 for each binary question $q_{(v_{a_k}, c_i)} \in Q_{(v_{a_k}, c_i)}$ **do**
 Pick the same number of image IDs from $P_{(v_{a_k}, c_i)}$ and $N_{(v_{a_k}, c_i)}$:
 $B \leftarrow (q_{(v_{a_k}, c_i), \text{yes}}, u_p \in P_{(v_{a_k}, c_i)}) \cup (q_{(v_{a_k}, c_i), \text{no}}, u_n \in N_{(v_{a_k}, c_i)})$

Attribute-based: First, we build an *attribute-value-to-images* dictionary to map each distinct *attribute value* or *category* to a set of eligible image IDs. Given a specific *attribute value*, we collect a

Images				
Fashion items	Shirt, Jumpsuit	Sweater, Pants	Dress	Shirt, Skirt
Attributes	Closure type; Leg length type	Sleeve length type; Neckline type	Neckline type; Pattern; Sleeve length type	Pattern; Color
Attribute values	Pull on, Pullover, Front buckle; Full length	Long; V neck	Split neck; Geometric print; Long	Letter print; Light gray
Question/Answer	Q: Is the person wearing a full-length jumpsuit with front buckle closure? A: Yes	Q: What is the sleeve length of her sweater on the top? A: long Q: What type of neckline is this sweater on the top? A: V neck	Q: What is the neckline type of this dress? A: split Q: is this a short sleeve geometric print dress? A: No	Q: What color of shirt is she wearing on the top? A: light gray Q: is the person wearing on the top a letter print pullover shirt? A: Yes

Figure 1: Four randomly picked *question-answer-image* triplets from FashionVQA dataset.

set of positive answer image IDs directly from this *attribute-value-to-images* dictionary using given *attribute value* and its synonyms. The negative answer image IDs are collected from all image IDs of the same *attribute* excluding the positive image IDs. More concretely, to maximally reduce the noise in the positive/negative answer image IDs, we need to verify the relationship among *attribute values* as alternative, hierarchical, or exclusive terms. Examples of alternative terminologies are “sweatpants”, “jogger pants”, and “loung pants”; examples of hierarchical terminologies are “blue”, “light blue”, and “sky blue”; and, examples of exclusive terminologies are “light blue” and “dark blue”. We expect *attribute values* with similar terminologies (alternatives and parents of hierarchical terms) to contain the same set of positive samples, so they are considered synonyms. In this manner, we can build an *attribute-value-to-(positive/negative answer)-images* dictionary (see Algorithm 1).

Then, we consider all the combinations of assorted *attributes* with *category*. For example, $\langle color, pattern, category \rangle$, $\langle color, category \rangle$, $\langle material, neckline type, category \rangle$, etc. For each combination, we further expand the *attribute-value-to-(positive/negative answer)-images* dictionary by mapping the combination of one specific *attribute value* and one specific *category* (e.g. $\langle red, shirt \rangle$) to its positive/negative answer image ID set. We collect the positive answer image ID set of the combinations following the formula in Equation 1 and the negative answer image ID set following the formula in Equation 2:

$$Pos(\langle attr1, category \rangle) = Pos(attr1) \cap Pos(category) \quad (1)$$

$$Neg(\langle attr1, category \rangle) = (Pos(attr1) \cap Neg(category)) \cup (Neg(attr1) \cap Pos(category)) \cup (Neg(attr1) \cap Neg(category)) \quad (2)$$

whereas, $Pos()$ is the positive answer image ID set and $Neg()$ is the negative answer Image ID set. With the *attribute-value-to-(positive/negative answer)-images* dictionary, we can easily generate different binary questions via filling the question templates with each combination of *attribute value* and *category* in the dictionary. We can pick a fixed number of positive and negative answer image IDs to guarantee the sample balance for each question. Following the same formula, we can easily expand the combinations to multiple attribute values and one category.

3.4 Dataset description

Figure 1 shows four randomly picked *question-answer-image* triplet examples in our dataset. There are 42 *attributes* in our dataset, including *category*, *color*, *pattern*, *occasion*, *material*, *number of*, 29 type-related *attributes*, 5 style-related *attributes*, and 2 shape-related *attributes*. The binary questions in our dataset are composed of three major types: *category*, *category + one attribute*, and *category + two attributes* with 1, 2, and 6 permutations between *category* and *attribute*, respectively, along with the ascending difficulty level to learn the alignment between a given binary question and an input image.

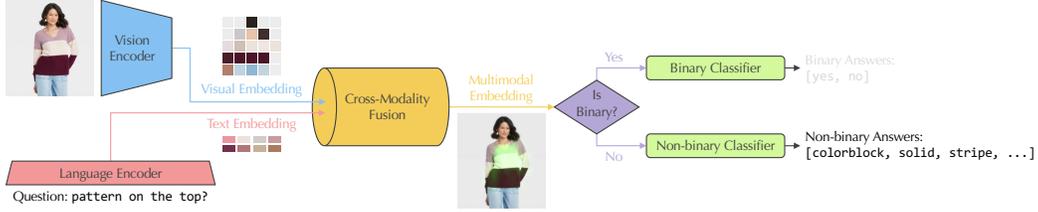


Figure 2: Pipeline of the fashion VQA task.

FashionVQA: FashionVQA dataset includes 207,654 unique photoshoot images with resolution 600×600 . We use 169,406 images in the train split for training and 38,248 images in the validation split for evaluation. The train split is composed of 163M *question-answer-image* triplets and the validation split includes 5.2M *question-answer-image* triplets. Since the information in binary questions is much more complicated than that in the non-binary questions, there are more binary triplets than non-binary ones in the dataset. In the train split, we have 22M non-binary *question-answer-image* triplets covering 33 different question types, and approximately 141M binary *question-answer-image* triplets, among which 134M have questions with one *category* and two *attribute values*, 6M have questions with one *category* and one *attribute value*, and 1M have questions with only one *category* or one *attribute value*. In the validation split, we have 1.2M non-binary *question-answer-image* triplets and 4M binary *question-answer-image* triplets. The answer vocabulary contains 1,545 different classes in total.

mini-FashionVQA: We also create a subset dataset, named mini-FashionVQA, derived from the FashionVQA dataset. The mini-FashionVQA dataset includes 20M *question-answer-image* triplets in the train split (11M from non-binary triplets and 9M from binary triplets) and 2.2M triplets in the validation split (0.7M from non-binary triplets and 1.5M from binary triplets).

4 Benchmarks

Every benchmark reported on our datasets is implemented via PyTorch[35]-v1.10 on servers with 8 Nvidia 80GB A-100 GPUs, 2 AMD 2.25GHz 7742 CPUs, and 4TB system memory. In the training stage, we adopted data-parallel multi-GPU training and set the batch size to 2048, and trained for 40 epochs. The Adam[36] optimizer is used across all the models. The learning rate is set to 0.0001 and reduced by half at the milestone epochs of 20, 30, and 35.

We benchmark the FashionVQA dataset by training several VQA models to learn the interaction between images and questions. Figure 2 shows the VQA pipeline adopted in our experiment. Given the visual embedding of the input image and text embedding of the input question sentence, we train the model to output the given answer to the question. The dataset is used to train two variants of the MCAN [10] model and a MUTAN [17] model. One MCAN variant, named MCAN*-v1, is a modification of the MCAN-small, which includes only two encoder-decoder modules. The other variant is named MCAN*-VLM, which has a similar structure to MCAN*-v1, but instead of an answer classifier, it has a token classifier covering all of the question and answer tokens. For MCAN*-VLM, the answer to each question is tokenized as one token and concatenated with the question tokens as the language input. The special token ‘SEP’ is inserted between the question and the answer. Also, ‘EOS’ token is used at the end of the answer. During the training of MCAN*-VLM, we randomly mask one token and predict the masked token as in the masked language modeling, similar to BERT [18]. Different from MCAN*-v1 that answer vocabulary is independent of the word vocabulary of the questions, MCAN*-VLM maps each answer to one token and expands the original word vocabulary to a larger one with the answer tokens. Thus, the tokens in the answers and questions share the same word vocabulary. This allows the MCAN*-VLM to work as a visual language model, which directly benefits from the overlap in the question tokens of the binary questions and the answer tokens of the non-binary question.

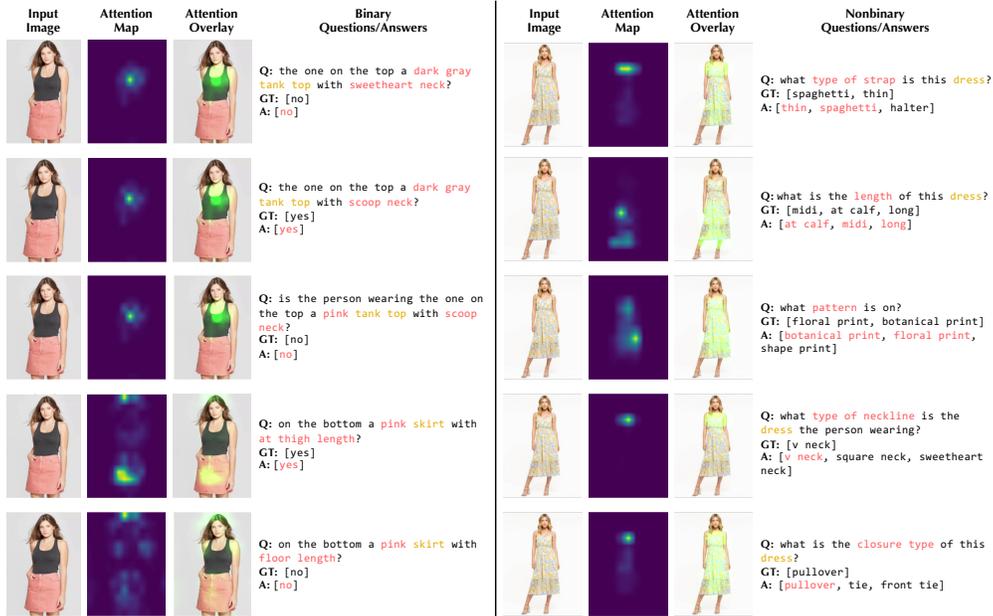


Figure 3: Visualization of attention maps generated by the model trained with FashionVQA dataset.

In the training stage, except MCAN*-VLM, we treat the binary-question prediction and the non-binary question prediction as two different tasks and output the predicted answers from two different classifiers. We report top-1 accuracies for both binary and non-binary samples.

Table 2: Benchmarks of MCAN*-v1, MCAN*-VLM, and MUTAN trained on FashionVQA dataset

Model	Top-1 Acc		
	All	Non-binary	Binary
MUTAN	81.38%	61.62%	87.43%
MCAN*-v1	84.42%	64.32%	90.58%
MCAN*-VLM	84.69%	64.65%	90.84%

Table 2 lists the benchmark results of the three aforementioned models on the validation split of our FashionVQA dataset. The results show that MCAN*-VLM works better than MCAN*-v1 and MUTAN, indicating that a decoder-only visual language model (VLM) performs better than the dedicated VQA architectures.

By visualizing the image attention maps generated from an intermediate layer of the model, we can validate whether the model focuses its attention on the regions mentioned in a question. Figure 3 visualizes the attention map from two validation samples for a series of binary and non-binary questions. The three columns of images on each side are the input images, attention maps, and images overlaid by the attention maps, respectively, followed by the corresponding input questions, ground truth answers, and predicted answers. When the questions focus on different fashion items of the same image, the attention map shifts to the targeted region as expected.

4.1 Benchmarks with different VQA models

We also use the mini-FashionVQA to benchmark a larger variety of VQA models including Bottom-up-top-down (BUTD) [6], MUTAN [17], DFAF [9], MCAN*-v1, MCAN*-v2, MCAN*-VLM, and OSCAR [11]. MCAN*-v2 has the same model structure as MCAN*-v1 except for its intermediate feed-forward layer with only half the number of channels of MCAN*-v1. We apply similar visual embedding, text embedding, and loss function in these models and train them from scratch.

Table 3 lists the results (average of top-1 accuracies for three runs) from different VQA models with the same region-based visual features as input. The visual features are extracted from FasterRCNN with ResNet-101 backbone fine-tuned with VisualGenome. We set the maximum number of objects extracted from the object-detection model to 25. The feature dimension of each object is 2048. A combination of GLove [37] + GRU[38]/LSTM[39] is used for the text embedding in DFAF, MCAN*-v1, MCAN*-v2, MCAN*-VLM, and BUTD. MUTAN adopts GRU for the text embedding, of which the parameters are initialized with SkipThoughts [40]. The number of parameters, FLOPs, and activation counts in all our experiments are calculated only from the cross-modality models, excluding text embedding and visual embedding components. On the mini-FashionVQA dataset, MCAN*-VLM achieves the best accuracy for both non-binary question and binary question samples, with fewer parameters and FLOPs than OSCAR. Also, MCAN*-VLM works better than MCAN*-v1 on both non-binary questions and binary questions.

Table 3: Performance on different VQA models trained on the mini-FashionVQA dataset with same region-based visual features

Model	Parameters	FLOPs	Act.Count	Top-1 Acc		
				All	Non-binary	Binary
MUTAN	9.8M	38.5M	156.9K	75.08%	59.14%	81.50%
BUTD	11.5M	61.5M	30.7K	79.26%	63.61%	85.56%
DFAF	9M	280M	114.8K	80.55%	62.52%	87.81%
OSCAR	86.7M	6475M	2832.3K	81.21%	64.20%	88.05%
MCAN*-v1	19M	427M	238.9K	81.69%	64.47%	88.61%
MCAN*-v2	14.5M	320M	134.5K	81.08%	64.33%	87.83%
MCAN*-VLM	19M	464M	282.0K	81.80%	64.63%	88.71%

4.2 Ablation study

Impact of visual embedding extraction schemes: We also benchmark different visual embedding extraction schemes and see their impact on the performance of VQA tasks for the mini-FashionVQA dataset. We replace the region-based feature with the grid feature with the same dimension. The grid-feature (refer to [13]) model is built with ResNext-101 backbone and fine-tuned with object detection task on the VisualGenome dataset. The visual feature is extracted from the same layer as the object detection and pooled into different sizes. To be aligned with the visual input dimension size from the region-based feature, the spatial dimension of the grid feature is set to 5×5 with the feature dimension set to 2048. Other than the visual embedding, all of the settings remain the same.

Table 4: Performance with region-based and grid visual features across different VQA models

Model	Region-based features (ResNet-101)			Grid features (ResNext-101)		
	All	Non-binary	Binary	All	Non-binary	Binary
MUTAN	75.08%	59.14%	81.50%	79.77%(+4.69%)	62.54%	86.70%
BUTD	79.26%	63.61%	85.56%	80.54%(+1.28%)	64.30%	87.08%
DFAF	80.55%	62.52%	87.81%	82.01%(+1.46%)	64.70%	88.97%
MCAN*-v1	81.69%	64.47%	88.61%	83.29%(+1.60%)	65.38%	90.49%
MCAN*-v2	81.08%	64.33%	87.83%	82.98%(+1.90%)	65.17%	90.14%
MCAN*-VLM	81.80%	64.63%	88.71%	83.41%(+1.61%)	65.52%	90.62%

Table 4 shows that the grid-feature-based visual embedding extraction method consistently works better than the region-based method across all different VQA models by more than 1% when trained on our dataset. In the other experiments, unless mentioned otherwise, we use grid-feature-based visual embedding for all the models.

Impact of different backbones for visual embedding: Generally, a better visual backbone will contribute to better visual embedding. We benchmark three different visual backbones (ResNet-50, ResNext-101, ResNext-152) for the grid-feature extraction on our dataset for MCAN*-v1, BUTD, and DFAF. All the visual backbones are pre-trained on VisualGenome [8] dataset for the grid-feature extraction.

Table 5: Performances with different visual backbones for grid-feature

Model	MCAN*-v1			BUTD			DFAF		
	All	Non-binary	Binary	All	Non-binary	Binary	All	Non-binary	Binary
ResNet-50	82.99%	65.25%	90.13%	80.47%	64.35%	86.96%	80.81%	63.39%	87.81%
ResNext-101	83.29%	65.38%	90.49%	80.54%	64.30%	87.08%	82.01%	64.70%	88.97%
ResNext-152	83.15%	65.41%	90.29%	80.50%	64.62%	86.89%	82.39%	65.17%	89.32%

Table 5 shows that ResNext-101 constantly works better than ResNet-50 on three different models for the performance of both non-binary and binary questions; however, the performance improvement from ResNext-101 to ResNext-152 is inconsistent. Overall, grid-feature with ResNext-101 as the backbone is the best choice for extracting visual features on our dataset.

Impact of different spatial dimension sizes for grid feature: A larger spatial dimension size after the pooling operation will typically preserve more visual information. We benchmark MCAN*-v1 with three different spatial dimension sizes (5×5 , 7×7 , and 9×9) for the grid feature visual embeddings in Table 6. ResNext-101 is the selected visual backbone.

The results in Table 6 show that the best performance among the three is from the smallest spatial dimension size, 5×5 , rather than the largest one. One possible reason is that the background of the photoshoot images from our dataset includes some trivial information, and the larger spatial dimension sizes do not add useful information.

Table 6: Performances with different spatial dimension sizes for grid-feature

Model	Spatial dimension size	Top-1 Acc		
		All	Non-binary	Binary
MCAN*-v1	5×5	83.29%	65.38%	90.49%
	7×7	82.94%	64.19%	90.48%
	9×9	82.60%	65.23%	89.59%

Impact of single-task versus multi-task training: Due to the large difference in the answer distribution of non-binary questions and binary questions, we consider using different classifiers for answer predictions and treating the problem as a multi-task classification. Namely, predicting answers for two types of questions with either a single classifier or two separate classifiers. This applies to all models, except the MCAN*-VLM model, where the outputs are generated by a single token classifier, including both answer and question tokens.

Table 7: Performance with different number of classifiers for non-binary and binary questions

Model	Single-task			Multi-tasks		
	All	Non-binary	Binary	All	Non-binary	Binary
MUTAN	79.40%	62.12%	86.36%	79.77%(+0.37%)	62.54%	86.70%
BUTD	80.32%	63.55%	87.07%	80.54%(+0.22%)	64.30%	87.08%
DFAF	81.6%	63.85%	88.74%	82.01%(+0.41%)	64.70%	88.97%
MCAN*-v1	83.24%	65.36%	90.44%	83.29%(+0.05%)	65.38%	90.49%

Table 7 demonstrates that the proposed multi-task classification is superior to a single-task classification in predicting the answers for the VQA models.

5 Comparison to human performance

Human accuracy for FashionVQA dataset: To see how well humans can answer the question in our dataset, we implemented a user interface that shows one *question-image* pair from the validation set at a time. The user interface allows the human annotators to select one of the acceptable answers among 1,545 answer classes, e.g., “yes”, “no”, “purple”, “unicorn print”, “tailored”, “fly hook and

loop fastener”, “three quarter length”, etc. We asked the annotators to answer each question to the best of their knowledge without looking up the terms.

We have two types of annotators: experts and non-experts. We trained our expert annotators with at least ten examples per fashion term in our word vocabulary. Both expert and non-expert annotators are trained on the VQA task of our dataset. Table 8 shows the accuracies of nine human annotators compared to the MCAN*-VLM model trained on the FashionVQA dataset.

Table 8: Performances of different human annotators on samples from FashionVQA validation set

Annotator	Number of samples			Accuracy			Accuracy p -value		
	All	Non-binary	Binary	All	Non-binary	Binary	All	Non-binary	Binary
Expert 1	728	216	512	63.6%	43.5%	72.1%	8.5e-30	1.1e-09	7.5e-20
Non-expert 1	106	29	77	58.5%	24.1%	71.4%	1.8e-07	1.4e-05	0.00018
Non-expert 2	70	18	52	52.9%	22.2%	63.5%	6.9e-07	0.0003	8.7e-05
Non-expert 3	61	17	44	63.9%	29.4%	77.3%	0.00072	0.0035	0.02
Non-expert 4	51	14	37	47.1%	14.3%	59.5%	1.2e-06	8.7e-05	0.00025
Non-expert 5	150	44	106	50.7%	22.7%	62.3%	3e-14	2.8e-08	1.3e-08
Non-expert 6	211	62	149	52.6%	22.6%	65.1%	9.5e-18	3.9e-11	4.4e-10
Non-expert 7	103	27	76	48.5%	25.9%	56.6%	3.3e-11	6.2e-05	3.6e-08
Non-expert 8	50	14	36	52.0%	14.3%	66.7%	1.6e-05	8.7e-05	0.0023

To analyze the statistical significance of the results, we calculated the p -values of the human accuracies with respect to the validation accuracy of the model using the one-sided t-test. The validation accuracies of the MCAN*-VLM model are 84.69%, 64.65%, and 90.84% for all, non-binary, and binary questions, respectively. The model outperforms all of the human annotators, and at a 95% confidence level, the differences between the model validation accuracy and human accuracies are statistically significant.

Accuracies for human-generated questions: We also stress-tested the model by measuring its performance on human-generated questions. We asked an expert annotator, Expert 2, to paraphrase the questions of 300 random samples (218 binary and 82 non-binary samples) from the validation set. We used these questions instead of the original questions in the validation set to measure the accuracies of the MCAN*-VLM model and a human annotator, Expert 1, as shown in Table 9.

Table 9: Performances of the MCAN*-VLM model and a human expert on human-generated questions

	Accuracy		
	All	Non-binary	Binary
Human Expert 1	62.3%	30.5%	74.3%
MCAN*-VLM	77.7%	47.6%	89.0%
p -value	1.9e-05	0.0125	3.4e-05

We performed a one-sided t-test to analyze the statistical significance of the difference between the human and the model accuracies. At a significance level of 0.05 ($\alpha = 0.05$), the p -values reject the null hypothesis of the human accuracy being greater than or equal to the model. Figure 4 provides several examples from this experiment.

Impact on downstream tasks: We performed a side-by-side comparison of the apparel search with and without FashionVQA. A baseline search engine returns the top 24 items for an apparel search query. Another variant of the search results is formed by reranking these 24 items with FashionVQA: we generate a set of binary questions from the search query and use MCAN*-VLM model trained with FashionVQA to answer these questions for each of the 24 items. The average confidence scores of the yes and no answers are used as the additional features to rerank the top 24 items.

For a number of randomly-selected search queries with two *attribute values* and one *category*, e.g., “green crew neck dress”, a human annotator is presented with the original and reranked search result pages and gets to choose her/his preferred result page. The result pages are randomly located on the left and right sides of the screen without the annotator knowing which of the two pages presents the reranked results. Figure 5 shows an example of our side-by-side A/B test for the given random

FashionVQA question paraphrased by Expert 2	FashionVQA image	FashionVQA answer	Expert 1 answer	MCAN*-VLM answer
What kind of top is she wearing?		Hoodie	Jacket	Jacket
What collar is the collar type?		Notched	Hooded	Notched
Is this shoe taupe beige and good for fall?		Yes	Yes	Yes
Shirt in gray with v neck?		Yes	No	Yes
Are these pairs of above-knee swim trunks?		No	Yes	No
What is the style of these socks?		Athletic socks	Crew socks	Athletic socks
Dark gray graphic t-shirt for casual occasions?		Yes	Yes	Yes
What is the cuff type of this shirt?		Barrel	Rolled	Barrel
Pink shift dress?		Yes	No	Yes
Is she wearing a no pattern one piece with turnover collar?		Yes	Yes	Yes
What type of legs are these?		Taper	Taper	Jogger

Figure 4: FashionVQA paraphrased and answered by humans and model

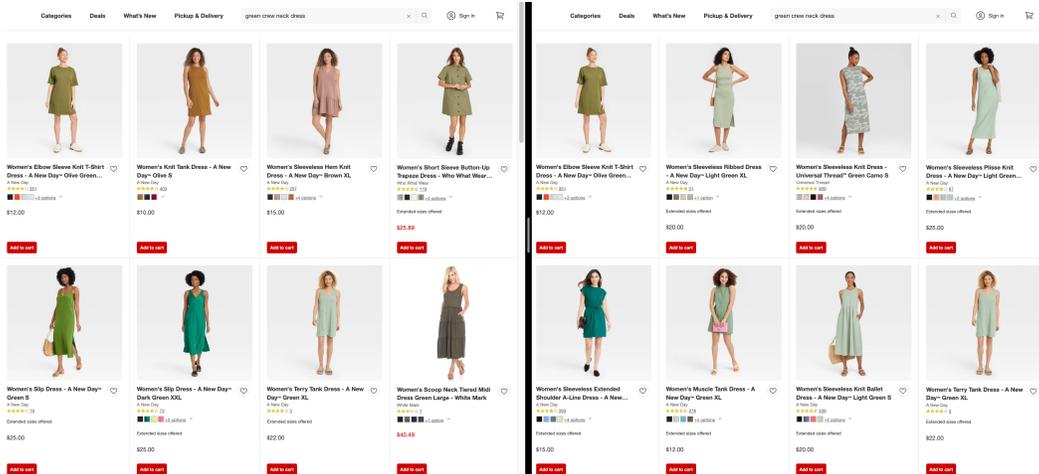


Figure 5: An example of the side-by-side comparison of the search results with and without reranking. Given a random search query, e.g. “green crew neck dress”, the annotator picks her/his preferred search results between the left (A) and right (B) result pages.

query. Out of 150 search queries, the human annotator preferred 117 search pages reranked based on the FashionVQA. Binomial statistical test results in a p -value of $3.2e-12$, showing that the human annotator significantly prefers the search result page reranked using FashionVQA.

Conclusion

In this work, we design a fashion VQA dataset and generate non-binary and binary questions via diverse templates. The templates allow us to flexibly scale the dataset to the size and complexity required for training a domain-specific multimodal model. We benchmark this large-scale dataset on different VQA models and discuss several factors impacting the performance of the VQA task. The best model is a visual language model trained on the FashionVQA dataset. The model generates the cross-modality embeddings of the vision and language domains applicable to downstream tasks of fashion dialogue, search, and recommendation.

References

- [1] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [2] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [4] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017.
- [5] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.

- [6] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018.
- [7] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015.
- [8] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017.
- [9] Peng Gao, Zhengkai Jiang, Haoxuan You, Pan Lu, Steven CH Hoi, Xiaogang Wang, and Hongsheng Li. Dynamic fusion with intra-and inter-modality attention flow for visual question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6639–6648, 2019.
- [10] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6281–6290, 2019.
- [11] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer, 2020.
- [12] Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019.
- [13] Huaizu Jiang, Ishan Misra, Marcus Rohrbach, Erik Learned-Miller, and Xinlei Chen. In defense of grid features for visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10267–10276, 2020.
- [14] Yang Gao, Oscar Beijbom, Ning Zhang, and Trevor Darrell. Compact bilinear pooling. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 317–326, 2016.
- [15] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. *arXiv preprint arXiv:1606.01847*, 2016.
- [16] Jin-Hwa Kim, Kyoung-Woon On, Woosang Lim, Jeonghee Kim, Jung-Woo Ha, and Byoung-Tak Zhang. Hadamard product for low-rank bilinear pooling. *arXiv preprint arXiv:1610.04325*, 2016.
- [17] Hedi Ben-Younes, Rémi Cadene, Matthieu Cord, and Nicolas Thome. MUTAN: Multimodal tucker fusion for visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2612–2620, 2017.
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [20] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. ViBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019.

- [21] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. VisualBERT: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [22] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. UNITER: Universal image-text representation learning. In *ECCV*, 2020.
- [23] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588, 2021.
- [24] Ziwei Liu, Ping Luo, Shi Qiu, Xiaogang Wang, and Xiaoou Tang. DeepFashion: Powering robust clothes recognition and retrieval with rich annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1096–1104, 2016.
- [25] Yuying Ge, Ruimao Zhang, Xiaogang Wang, Xiaoou Tang, and Ping Luo. DeepFashion2: A versatile benchmark for detection, pose estimation, segmentation and re-identification of clothing images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5337–5345, 2019.
- [26] Shuai Zheng, Fan Yang, M Hadi Kiapour, and Robinson Piramuthu. ModaNet: A large-scale street fashion dataset with polygon annotations. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1670–1678, 2018.
- [27] Xingxing Zou, Xiangheng Kong, Waikeng Wong, Congde Wang, Yuguang Liu, and Yang Cao. FashionAI: A hierarchical dataset for fashion understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- [28] Hui Wu, Yupeng Gao, Xiaoxiao Guo, Ziad Al-Halah, Steven Rennie, Kristen Grauman, and Rogerio Feris. Fashion IQ: A new dataset towards retrieving images by natural language feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11307–11317, 2021.
- [29] Sheng Guo, Weilin Huang, Xiao Zhang, Prasanna Srikhanta, Yin Cui, Yuan Li, Hartwig Adam, Matthew R Scott, and Serge Belongie. The iMaterialist fashion attribute dataset. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.
- [30] Menglin Jia, Mengyun Shi, Mikhail Sirotenko, Yin Cui, Claire Cardie, Bharath Hariharan, Hartwig Adam, and Serge Belongie. Fashionpedia: Ontology, segmentation, and an attribute localization dataset. In *European conference on computer vision*, pages 316–332. Springer, 2020.
- [31] Lukas Bossard, Matthias Dantone, Christian Leistner, Christian Wengert, Till Quack, and Luc Van Gool. Apparel classification with style. In *Asian conference on computer vision*, pages 321–335. Springer, 2012.
- [32] Kota Yamaguchi, M Hadi Kiapour, Luis E Ortiz, and Tamara L Berg. Parsing clothing in fashion photographs. In *2012 IEEE Conference on Computer vision and pattern recognition*, pages 3570–3577. IEEE, 2012.
- [33] Sirion Vittayakorn, Kota Yamaguchi, Alexander C Berg, and Tamara L Berg. Runway to realway: Visual analysis of fashion. In *2015 IEEE Winter Conference on Applications of Computer Vision*, pages 951–958. IEEE, 2015.
- [34] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

- [36] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [37] Jeffrey Pennington, Richard Socher, and Christopher D Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [38] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [39] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- [40] Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. *Advances in neural information processing systems*, 28, 2015.