# CMA-CLIP: Cross-Modality Attention CLIP for Image-Text Classification

Huidong Liu*
huidliu@cs.stonybrook.edu
Stony Brook University
Stony Brook, NY, USA

Shaoyuan Xu*
shaoyux@amazon.com
Amazon Inc.
Seattle, WA, USA

Jinmiao Fu*
jinmiaof@amazon.com
Amazon Inc.
Seattle, WA, USA

Yang Liu
Ning Xie
Chien-Chih Wang
yliuu@amazon.com
xining@amazon.com
ccwang@amazon.com
Amazon Inc.
Seattle, WA, USA

Bryan Wang
Yi Sun
brywan@amazon.com
yisun@amazon.com
Amazon Inc.
Seattle, WA, USA

## ABSTRACT

Modern Web systems such as social media and e-commerce contain rich contents expressed in images and text. Leveraging information from multi-modalities can improve the performance of machine learning tasks such as classification and recommendation. In this paper, we propose the Cross-Modality Attention Contrastive Language-Image Pre-training (CMA-CLIP), a new framework which unifies two types of cross-modality attentions, sequence-wise attention and modality-wise attention, to effectively fuse information from image and text pairs. The sequence-wise attention enables the framework to capture the fine-grained relationship between image patches and text tokens, while the modality-wise attention weighs each modality by its relevance to the downstream tasks. In addition, by adding task specific modality-wise attentions and multilayer perceptrons, our proposed framework is capable of performing multi-task classification with multi-modalities.

We conduct experiments on a Major Retail Website Product Attribute (MRWPA) dataset and two public datasets, Food101 and Fashion-Gen. The results show that CMA-CLIP outperforms the pre-trained and fine-tuned CLIP by an average of 11.9% in recall at the same level of precision on the MRWPA dataset for multi-task classification. It also surpasses the state-of-the-art method on Fashion-Gen Dataset by 5.5% in accuracy and achieves competitive performance on Food101 Dataset. Through detailed ablation studies, we further demonstrate the effectiveness of both cross-modality attention modules and our method's robustness against noise in image and text inputs, which is a common challenge in practice.

## 1 INTRODUCTION

Inspired by the recent rise of the pre-trained NLP models such as BERT [11], learning to classify image-text pairs for vision-language ($VL$) tasks using Transformer [37] based encoders has received much attention as both modalities can be informative and beneficial to each other. Current methods can be classified into two main categories: one-stream methods and two-stream methods. One-stream methods capture the cross-modality attention across

image and text by concatenating them at early stage and input the concatenated feature into one unified transformer encoder. Two-stream methods first extract the image and text features using two separate encoders and then learn their cross-modal relationship through various methods such as contrastive learning [3], etc.

Among those one-stream and two-stream methods, the Contrastive Language-Image Pre-Training (CLIP) [30] has achieved great success recently. CLIP is trained on the WebImageText (WIT) Dataset which consists of 400 million image-text pairs collected from a variety of publicly available sources on the Web and achieves many state-of-the-art results in zero-shot learning tasks, pre-training tasks, and supervised classification tasks when a linear probe is added on top of it.

Despite of CLIP's [30] strength, it is mainly designed for zero-shot image classification, resulting in the limitation of its ability to leverage both image and text input when available. Since the training of CLIP only involves global image and text features, thus the fine-grained relationship between image patches and text tokens are not modeled. Such relationship is useful in fine-trained classification tasks, especially in the situations where only a small proportion of image patches or text tokens are related to the classification tasks. Moreover, since it chooses the user-defined textual description called "prompts" with class value that matches the most with the image as the classification result, significant efforts to engineer the prompts for optimizing downstream tasks are required. Last but not least, in practice, it is quite common for the image-text pairs to contain noise. For instance, on E-commerce websites, some images or text could be irrelevant to the product due to catalog errors. In social media apps, users might enter irrelevant textual comments or upload unrelated images. Treating input from both modalities equally in such situation may lead to poor classification performances as one of the modalities could be pure noise. To address the aforementioned issues, in this paper we propose the Cross-Modality Attention CLIP (CMA-CLIP). Our contributions include:

- We combine CLIP with a sequence-wise attention module, which refines the CLIP-generated image and text embeddings

---

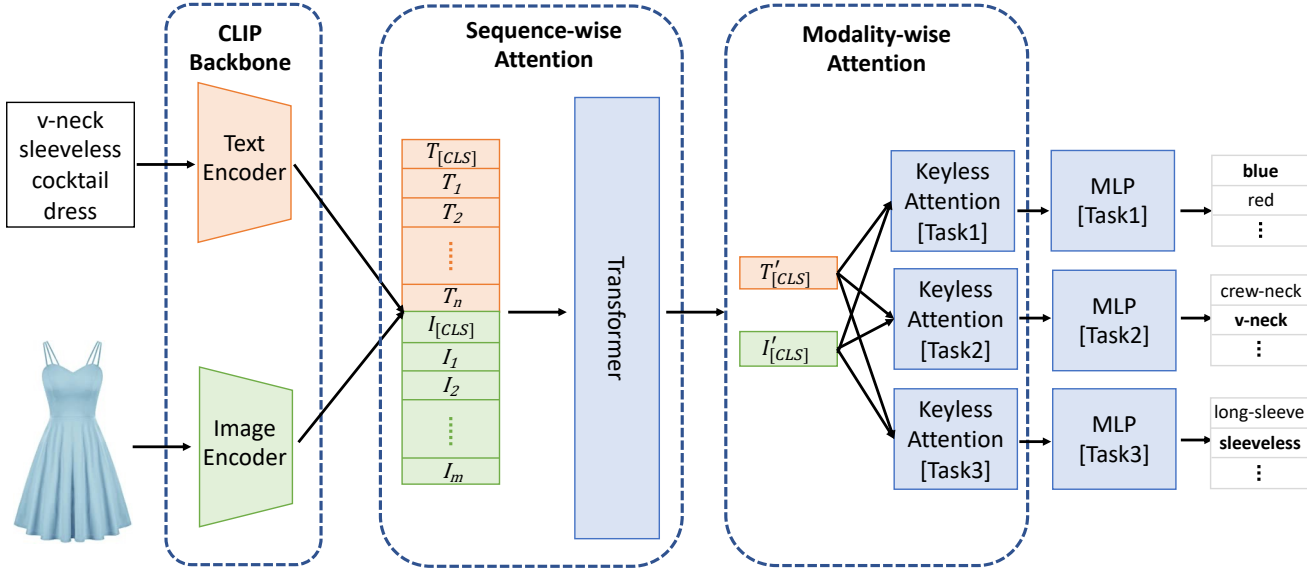*All three authors contributed equally to this research.

**Figure 1: The pipeline of our proposed CMA-CLIP.**

by modeling the relationship among the embedding of the sequence of image patches and text tokens. This transformer-based module makes the embedding more context-aware. *E.g.*, the embedding of the black image patches are more correlated with the '*black*' token in text. We experimentally prove that such refinement can improve the performance of classification tasks.

- We adopt a modality-wise attention module to assign learnable weights to each modality that measures its relevance to the classification task. The impact of the irrelevant modality will be dampened, and therefore our network is robust against noisy image or text inputs, which is a common challenge in practice.
- We add task specific modality-wise attentions and MLP heads on top of the sequence-wise attention module, so that the same network can be leveraged for multi-task classification. Moreover, compared with CLIP, this architecture enables the network to leverage both image and text inputs in both training and inference stage.
- On the MRWPA, CMA-CLIP outperforms the raw CLIP and the fine-tuned CLIP (fine-tuned using image-text pairs from a major retail website) on the classification of three product attributes by 10.9% and 12.9% in recall at the same level of precision. It also improves the state-of-the-art performance [44] on Fashion-Gen Dataset from 88.1% to 93.6% in accuracy and achieves competitive performance on Food101 Dataset against [21].

## 2 RELATED WORK

Current multi-modality learning methods are mainly one-stream and two-steam where one-stream methods use a single Transformer encoder to process the concatenated image and text embedding, while two-stream methods use both image encoder and text encoder to extract image and text embeddings at early stage and then learn their cross-modal relationship.

The one-stream methods, such as ViLBERT [28], VisualBERT [26], VL-BERT [35], Unicoder-VL [25], ImageBERT [29] and Unified VLP [43], concatenate the image's Region-Of-Interest (ROI) patches and text tokens as the input tokens for BERT [11]. These models are typically pre-trained using tasks including Masked Language Modeling (MLM), Masked Region Modeling (MRM), Multi-Model Alignment Prediction (MMAP). The UNITER [8] and OSCAR [27] incorporate additional pre-training tasks. The UNITER uses the Optimal Transport (OT) [38] to model the relationship between the image patches and the text tokens. The OSCAR uses the object categories detected by the Faster-RCNN [31] and encodes the category text as additional input tokens to BERT. Instead of using the Faster-RCNN to detect ROIs, methods such as ICMLM [33], Pixel-BERT [18], and SOHO [17] use a CNN to extract the feature maps of an image and use the depth vectors in feature maps as image tokens. Such configuration is able to capture the semantic connection between image pixels and text tokens, which is overlooked by region based image features extracted by Faster-RCNN. Similar to ICMLM, the VirTex [10] also uses the depth vectors in feature maps as image tokens and input the image and text tokens into a forward transformer decoder and a backward transformer decoder. Instead of feeding the whole image or ROI into a CNN, methods like FashionBERT [13], KaleidoBERT [44], and ViLT [22] cut an image into patches and treat each patch as an "image token". For FashionBERT, it uses a pre-trained image model such as InceptionV3 [36] or ResNeXt-101 [40] to extract image features. Different from FashionBERT, KaleidoBERT adopts the SAT [41]

network to generate description of salient image patches aiming to find an approximate correlation between image patches and text tokens to serve their pre-training tasks, *i.e.*, Aligned Masked Language Modeling (AMLM), Image and Text Matching (ITM), and Aligned Kaleido Patch Modeling (AKPM). ViLT differs from all above-mentioned methods by simply applying linear projection on flattened image patches which greatly reduce the model size, thus leading to significant runtime and parameter efficiency.

The two-stream methods are mainly motivated by self-supervised learning methods [2–7, 14, 15]. In self-supervised learning, two views (*e.g.*, two augmentations), of a single image are forwarded into one network respectively. Their outputs are compared using the contrastive loss [3], so that the two views of the same image are much similar than the two views from two different images. The ConVIRT [42] adopts this idea on the self-supervised learning of image-text pairs. Two networks are used to extract the image and text features respectively. The image and text features from the paired image-text input is trained to be much similar than the unpaired ones. CLIP [30] is a simplified version of ConVIRT, where the text in each image-text pair is a single sentence instead of a pool of sentences as in ConVIRT. Similar as CLIP, BriVL [19] uses MoCo [15] which is a more advanced cross-modal contrastive learning algorithm to help train the network with limited GPU memory by leveraging more negative samples. The ALIGN [20] collects 1.8 billion image-text pairs, and adopts a similar network architecture as CLIP. The performance of ALIGN is comparable to CLIP on the ImageNet dataset for the classification task, Flickr30K and MSCOCO datasets for image-text retrieval task.

In order to perform image-text classification tasks, a classification layer needs to be added on top of the image-text embeddings of the pre-trained models such as VL-BERT [35], UNITER [8], et al. The MMBT [21] is specifically designed for image-text classification tasks. Unlike VL-BERT or UNITER, MMBT directly loads weights from BERT which does not require pre-training. In MMBT, ResNet [16] is used to extract image features. The image features are projected into the same space as the text tokens, used as image tokens and feed into a BERT together with text tokens. A linear layer is added on the classification embedding to perform supervised tasks.

However, both one-stream and two-stream methods have their own inadequacies. One-stream methods heavily rely on pre-trained Faster-RCNN [31] or ResNet [16] to extract image feature which does not support end-to-end training of the whole network, and therefore, the extracted image and text features are not optimized to model the image-text relationship. Two-stream methods only focus on learning global image and text features, and cannot capture the fine-grained relationship between image patches and text tokens.

In order to overcome the aforementioned disadvantages, our proposed method fuses both one-stream and two-stream architectures which complements each other's inadequacies. It leverages the pre-trained CLIP, a two-stream architecture, to capture the overall alignment between image and text. Subsequently, we add a sequence-wise attention module, which is a transformer based cross-modality attention module used in most one-stream architectures, to capture the fine-grained relationship between image patches and text tokens. SemVLP [24] applies similar fusion logic. The difference between our method and SemVLP is that, SemVLP leverages the same cross-modality attention module to capture both

high-level and fine-grained relationship between image and text. It was pre-trained on tasks such as MLM and MRM, whereas CLIP was directly pre-trained to maximize the overall image-text alignment through contrastive learning. More importantly, SemVLP does not consider the situation where one of the modalities is irrelevant to the downstream classification tasks due to input noises, which is a common challenge in practice. To handle such situation, we add a modality-wise attention module, which learns the importance of both modalities so that the irrelevant modality can be dampened for the classification tasks. At last, by adding task specific modality-wise attentions and MLPs, our model is able to perform multi-task classifications.

## 3 METHOD

The rest of the paper is arranged as follows: In Section 3, we first give a brief review of CLIP, and then we introduce our proposed CMA-CLIP with detailed explanation of each component. In Section 4, we introduce the datasets that we use, the corresponding experimental results, the visualization of the sequence-wise attention module, and the ablation study to prove the effectiveness of modality- and sequence-wise attention modules. In Section 5, we conclude this paper and elaborate our future work.

### 3.1 Contrastive Language-Image Pre-Training (CLIP)

The Contrastive Language-Image Pre-Training (CLIP) consists of an image encoder and a text encoder. For each image-text pair, the image and text encoders project the pair into an image and text embedding in the same multi-modal space. Given $N$ image-text pairs, the training objective of CLIP is to maximize the cosine similarity of the paired image and text embedding while minimize the cosine similarity of the unpaired ones.

During inference, for a classification task with $K$ classes, it first uses the $K$ class values to construct $K$ prompts such as 'A photo of **{class value}**'. These $K$ prompts are then projected to $K$ text embeddings by the text encoder. For any given image, it is projected to an image embedding by the image encoder, then CLIP computes the cosine similarities between the image embedding and those $K$ text embeddings. The class value with the largest similarity is then considered as the class prediction.

CLIP is trained using WIT Dataset which contains 400 million image-text pairs collected from the Web. According to the results reported in [30], its zero-shot classification performance surpasses the supervised linear classifier fitted on ResNet50 [16] features on datasets such as StanfordCars [23], Country211 [30], Food101 [1], and UCF101 [34] etc.

### 3.2 The Cross-Modality Attention CLIP (CMA-CLIP)

CLIP focuses on the learning of the global image and text features. In CMA-CLIP, we build a sequence-wise attention module to capture the fine-grained relationship between the image patches and the text tokens such as the black image patches and the 'black' tokens in text. This module leverages the transformer architecture [37]. It takes the sequence of embeddings corresponding to all the image patches and text tokens generated by CLIP as input. The

| Dataset | Image | Label | Text |
|---|---|---|---|
| MRWPA | | Striped | Women Multi Striped Sheath Dress |
| Food101 Dataset | | Pancakes | Banana Bread Pancakes with Cinnamon Cream Cheese Syrup - Cooking Classy |
| Fashion-Gen Dataset | | Jeans | Slim-fit jeans in dark blue. Distressing throughout. Fading at front. Textured black leather logo patch at back waist. Silver-tone metal logo plaque at back pocket. Contrast stitching in tan. Red logo tab at button-fly. |

**Table 1: Data example for each of the datasets**

module outputs two embeddings incorporating the aggregated image and text information. Instead of directly leveraging these two embeddings for classification, we add a modality-wise attention module to handle the situation where a certain modality (image or text) is irrelevant to the classification task. This is because, in practice, it is common for the image-text pairs to contain noise. *E.g.*, a retailer might upload wrong product images to E-commerce website, or a user might enter random textual comments on social media apps. To handle such situations, we leverage the similar architecture as in [39] to learn the importance of each modality to the classification tasks. The sum of the two embeddings weighted by their importances is followed by a MLP head for the classification. To leverage the network for multiple classification tasks, we configure task specific modality-wise attention modules and MLP heads. The complete architecture of CMA-CLIP is shown in Fig. 1.

*3.2.1 Sequence-wise Attention.* In our implementation, the sequence-wise attention module is a transformer encoder [37]. Let $X \in \mathbb{R}^{s \times d}$ be the matrix of the sequence of embedding of all the image patches and text tokens generated by CLIP, where $s$ is the length of the sequence and $d$ is the dimension of the embedding. Let $W_K \in \mathbb{R}^{s \times d}$, $W_Q \in \mathbb{R}^{s \times d}$ and $W_V \in \mathbb{R}^{s \times d}$ be the projection matrices which project each embedding in $X$ to key space, query space and value space respectively:

$$K = XW_K, \quad Q = XW_Q, \quad V = XW_V \quad (1)$$

The embedding matrix $X$ is updated as

$$\text{Attention}(K, Q, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V \quad (2)$$

The self-attention block learns a similarity matrix $QK^T$ between each pair of embeddings in $X$. Each embedding in the sequence is

then updated as the average of the projected embedding across all the embeddings in the value space weighted by their similarities. This sequence-wise attention module captures the fine-grained relationship between each image patch and text token.

*3.2.2 Modality-wise Attention.* After the image feature and the text feature are generated, they need to be aggregated to form the final feature for the classification tasks. In order to dampen the irrelevant modality, we leverage the keyless attention module proposed in [39]. Given an image-text pair, the sequence-wise attention module will output $I'_{[CLS]}$ and $T'_{[CLS]}$ as the global image and text embedding, respectively. The aggregated embedding is the weighted average of $I'_{[CLS]}$ and $T'_{[CLS]}$:

$$c = \lambda I'_{[CLS]} + (1 - \lambda)T'_{[CLS]} \quad (3)$$

The weight $\lambda$ is computed by:

$$e_I = w^T I'_{[CLS]} \quad (4)$$

$$e_T = w^T T'_{[CLS]} \quad (5)$$

$$\lambda = \frac{\exp(e_I)}{\exp(e_I) + \exp(e_T)} \quad (6)$$

where $w$ is a learnable parameter vector that is of the same dimension as $I'_{[CLS]}$ and $T'_{[CLS]}$.

**Algorithm 1** Training Process of CMA-CLIP. Line 4-16: Warm-up Stage, Line 18-21: End-to-End Training Stage, Line 23-26: Tuning Stage

**Input:** Image-text pairs and their labels of K tasks $(\mathcal{X}, \mathcal{Y}, \mathcal{L}_1, \mathcal{L}_2, ..., \mathcal{L}_K)$.

1: *CLIP* denotes the CLIP model including the image and text encoders. $CMA_{SA}$ denotes the sequence-wise attention model. $CMA_{MA}$ denotes the modality-wise attention model. *MLP* denotes the multi-layer perception head for classification.

2: **Training:**

3: Use the image encoder and the text encoder from CLIP, and freeze their weights.

4: Set $lr_{CLIP} = 0$, $lr_{CMA} = 1e-5$, $lr_{MLP} = 1e-5$            ▷ since the CLIP module is frozen, we set $lr_{CLIP}$ to be 0

5: **while** $i < MaxEpoch1$ **do**

6:      $(\hat{\mathcal{X}}_{CLIP}, \hat{\mathcal{Y}}_{CLIP}) = CLIP(\mathcal{X}, \mathcal{Y})$            ▷ apply CLIP's image and text encoders on $\mathcal{X}$ and $\mathcal{Y}$

7:      $\hat{\mathcal{X}}\hat{\mathcal{Y}}_{concat} = concat(\hat{\mathcal{X}}_{CLIP}, \hat{\mathcal{Y}}_{CLIP})$            ▷ concatenate the image feature $\hat{\mathcal{X}}_{CLIP}$ and the text feature $\hat{\mathcal{Y}}_{CLIP}$

8:      $(\hat{\mathcal{X}}_{SA}, \hat{\mathcal{Y}}_{SA}) = CMA_{SA}(\hat{\mathcal{X}}\hat{\mathcal{Y}}_{concat})$            ▷ apply the sequence-wise attention module to the concatenated feature $\hat{\mathcal{X}}\hat{\mathcal{Y}}_{concat}$

9:      **for** k = 1 to K **do**

10:          $\mathcal{U}_{MA}^k = CMA_{MA}^k(\hat{\mathcal{X}}_{SA}, \hat{\mathcal{Y}}_{SA})$            ▷ apply the modality-wise attention module to get the weighted average feature $\mathcal{U}_{MA}^k$

11:          $\hat{\mathcal{L}}_k = MLP(\mathcal{U}_{MA}^k)$            ▷ use the MLP head for classification task to get prediction $\hat{\mathcal{L}}_k$

12:          $Loss_k = Softmax(\mathcal{L}_k, \hat{\mathcal{L}}_k)$            ▷ compute Softmax Loss between ground truth $\mathcal{L}_k$ and prediction $\hat{\mathcal{L}}_k$

13:      **end for**

14:      $Loss = Loss_1 + Loss_2 + ... + Loss_K$

15:      $\theta^{i+1} = Adam(Loss, \theta^i)$            ▷ update model parameters using Adam

16: **end while**

17: Release all modules. Fine-tune on the best check-points from the previous stage.

18: Set $lr_{CLIP} = 1e-5$, $lr_{CMA} = 1e-5$, $lr_{MLP} = 1e-5$            ▷ since all modules are released, all learning rates are set to $1e-5$

19: **while** $i < MaxEpoch2$ **do**

20:      Repeat Line 6-15.

21: **end while**

22: Freeze the image encoder, the text encoder, the sequence-wise attention transformer. Fine-tune on the best check-points from the previous stage.

23: Set $lr_{CLIP} = 0$, $lr_{CMA} = 0$, $lr_{MLP} = 1e-5$            ▷ since both CLIP and CMA modules are frozen, we set $lr_{CLIP}$ and $lr_{CMA}$ to be 0

24: **while** $i < MaxEpoch3$ **do**

25:      Repeat Line 6-15.

26: **end while**

**Output:** *CLIP* model, $CMA_{SA}$ model, $CMA_{MA}$ model, and *MLP* model.

| Method | Data | Color | Pattern | Style | Avg. |
|---|---|---|---|---|---|
| Raw CLIP | WIT | 47.3 | 58.0 | 22.1 | 42.5 |
| Fine-tuned CLIP | MRWPA | 53.4 | 56.5 | 11.7 | 40.5 |
| CMA-CLIP | MRWPA | **61.1** | **76.3** | **22.9** | **53.4** |

**Table 2: Recall (%) at 90% precision on the MRWPA dataset.**

| Method | Fashion-Gen |
|---|---|
| FashionBERT | 85.3 |
| ImageBERT | 80.1 |
| OSCAR | 84.2 |
| KaleidoBERT | 88.1 |
| CMA-CLIP | **93.6** |

**Table 4: Accuracies (%) on the Fashion-Gen dataset.**

| Method | Food101 |
|---|---|
| ViT | 81.8 |
| BERT | 87.2 |
| CLIP | 88.8 |
| MMBT | 92.1 |
| CMA-CLIP | **93.1** |

**Table 3: Accuracies (%) on the Food101 dataset.**

is used to compute the loss for this classification task. For multi-task classification, we add task-specific modality-wise attention and MLP for each task separately. This is because the relevance of modality is dependent on the task, hence we need task-specific modality-wise attentions and multiple MLPs as the classification heads.

*3.2.3 MLP Heads for Classification.* For any classification task, an Multi-Layer Perception (MLP) head is added on top of the final feature outputted by the modality-wise attention. The cross entropy

| Attribute | Image | Label | Title |
|-----------|-------|-------|-------|
| Color | | Black | Portland [Black] T-Shirt Dress |
| Pattern | | Plain | Women's Sexy V Neck Crisscross Backless Cocktail Party Bodycon Peplum [Plain] Dress, White, L |
| Pattern | | Plain | Women's [Plain] Mini Dungaree |

**Table 5: Examples where CMA-CLIP is able to give the correct attribute classification while CMA-CLIP w/o the modality-wise attention cannot. Noting that, text tokens which are shown in red are the attribute label keywords that do not exist in the original titles. We add them in and re-do the prediction with both methods to further validate that, without the modality-wise attention, the model is not able to manage noise properly.**

| Method | Color | Pattern | Style | Avg. |
|--------|-------|---------|-------|------|
| CMA-CLIP | **61.1** | **76.3** | **22.9** | **53.4** |
| CMA-CLIP w/o $MA$ | 60.0 | 67.9 | 14.5 | 47.5 |
| CMA-CLIP w/o ($MA + SA$) | 57.3 | 60.3 | 19.8 | 45.8 |

**Table 6: Ablation study of CMA-CLIP. Recall (%) at 90% precision on the MRWPA with Color, Pattern and Style attributes. $MA$ denotes modality-wise attention and $SA$ denotes to sequence-wise attention.**

# 4 EXPERIMENTS

## 4.1 Datasets

We perform experiments on three datasets, the MRWPA Dataset, the Food101 Dataset [1] and the Fashion-Gen [32] Dataset. All three datasets consist of image-text pairs. Data samples from the three datasets are shown in Table 1.

*4.1.1 The MRWPA Dataset.* This dataset includes the product image and title pairs of dress products from a major retail website. The goal is to classify three dress related product attributes, color, pattern, and style. Color attribute has 17 classes such as black and white, pattern has 12 classes such as graphic and plain, and style has 21 classes such as pencil and a-line. The training data consists of 5.8 Million product image-title pairs. We also prepare 310 and 132 image-title pairs as the validation and test set, which are used for hyper-parameter tuning and performance evaluation respectively.

*4.1.2 The Food101 Dataset.* This dataset contains 101 food categories. The goal is to classify each image-text pair to a food category. We download the preprocessed images and texts from the Kaggle competition[1]. In the processed data, 67971 images are in the training set, and 22715 images are in the testing set. During training, we randomly split 80% of the data in the training set for training and the rest 20% data for validation.

*4.1.3 The Fashion-Gen Dataset.* This dataset contains 293,008 fashion images. Each image is paired with a text describing the image. This dataset contains 48 main categories, such as "DRESSES", "JEANS", "SKIRTS", "SHIRTS", etc., and 121 sub-categories, such as "SHORT DRESSES", "LEATHER JACKETS", "MID LENGTH SKIRTS", "T-SHIRTS" and so on. In our experiments, we perform 121 sub-category classification. We use the same data as used in [44] for training and testing. The number of training data is 260480, and the number of testing data is 32528.

## 4.2 Implementation and Settings

*4.2.1 Experiment Settings.* Same as CLIP, the image encoder of CMA-CLIP is a 12-layer 768-width ViT-B/32 [12] with 12 attention heads, and the text encoder of CMA-CLIP is a 12-layer 512-width Transformer with 8 heads used in [37]. The sequence-wise attention transformer is also a 12-layer 512-width model with 8 attention heads. In all the experiments, the batch size is set to 1024, weight decay of Adam is set to $1e-4$, and the learning rate is set to $1e-5$.

---

[1]https://www.kaggle.com/gianmarco96/upmcfood101

(a) MRWPA: Fashion Women Graphic Print Round Neck Ringer T-Dress Yellow

(b) MRWPA: Women's Summer Long Sleeve Casual Loose T-Shirt Dress

(c) MRWPA: Fashion Women's Oversized Short Sleeves Floral Print Mid-Long Dress Yellow Size UK 16

(d) Food101: The Brewer & The Baker: Lobster Tail Ale & Lobster Rolls

(e) Food101: Scallop Recipe for Beginners | Pop Sugar Food

(f) Food101: Shrimp and Grits - Picture-Perfect Meals \xc2 \xae Picture-Perfect Meals

(g) Fashion-Gen: Grained calfskin shoulder bag in 'plaid' red. Curb chain shoulder strap. Logo stamp gold-tone...

(h) Fashion-Gen: kinny-fit jeans in mid blue wash. White paint at leg. Low waist. Fading and distressing...

(i) Fashion-Gen: Ankle-high grained leather boots in black. Pointed toe. Zip closure at heel. Tonal leather...

**Figure 2: Visualization examples. Each image is highlighted using the attention map between the image embedding and the embedding of the most relevant text token.**

*4.2.2 Training Strategy.* We use the pre-trained weights of CLIP as the initial weights of the image encoder and text encoder in CMA-CLIP. We randomly initialize the weights in the sequence-wise attention module, modality-wise attention module and MLP. As CMA-CLIP contains a mixture of pre-trained weights and randomly initialized weights, instead of training the model end-to-end which may cause under- or over-fitting of certain modules, we adopt a multi-stage training strategy to train CMA-CLIP. The training stages are listed below:

- **Warm-up stage.** In this stage, the weights of the image encoder and the text encoder are frozen. We train the sequence-wise attention, the modality-wise attention and the MLP modules.

- **End-to-end training stage.** In this stage, we unfreeze the weights of the image encoder and the text encoder, and train all the components together.

- **Tuning stage.** This stage is for multi-task training. The weights of the image encoder, the text encoder and the sequence-wise attention are frozen. We train the modality-wise attentions and MLPs for all the tasks.

*4.2.3 Implementation.* Detailed training process of CMA-CLIP is summarized in Algorithm 1. For the MRWPA Dataset, all three stages are trained for 20 epochs. For Food101 and Fashion-Gen Datasets, Warm-up stage is trained for 100 epochs and End-to-end training stage is trained for 300 epochs. Since for the two public datasets, they are both single-task classification so the Tuning stage is not needed. During the Warm-up stage, due to the freeze of the CLIP module, only the check-points of the sequence-wise attention, the modality-wise attention and the MLP modules are updated. During the End-to-end training stage, check-points of all three modules are updated. And during the Tuning stage, only the check-points of the modality-wise attention the MLP modules are updated. At the end of each training stage, the best check-points with the lowest validation accuracy are used for either next stage's fine-tuning or inference.

### 4.3 Experimental Results

*4.3.1 The MRWPA Dataset.* We compare CMA-CLIP with the zero-shot performance of raw CLIP and fine-tuned CLIP (fine-tuned using image-title pairs in MRWPA dataset) in terms of the recall at 90% precision for the color, pattern, and style attributes. The results are included in Table 2. We observe that CMA-CLIP consistently outperforms both raw CLIP and fine-tuned CLIP by a large margin across all three attributes.

*4.3.2 The Food101 Dataset.* On the Food101 dateset, we compare CMA-CLIP with two single-modality baseline methods including BERT [11] and ViT [12], and two multi-modality baseline methods including raw CLIP [30] with same ViT-B/32 and MMBT [21]. Results are included in Table 3. CMA-CLIP achieves the best accuracy of 93.1%, which improves 1% over the a current strong baseline method MMBT. Using only image features achieves 81.8% by ViT, and using only text features achieves 87.2% by BERT.

*4.3.3 The Fashion-Gen Dataset.* On the Fashion-Gen dataset, we compare CMA-CLIP with multiple SOTA methods including Fash-ionBERT [13], ImageBERT [29], OSCAR [27] and KaleidoBERT [44]. CMA-CLIP achieves the highest accuracy of 93.6%, which improves over KaleidoBERT [44], the previous SOTA method, by 5.5%.

## 4.4 Ablation Study

We conduct systematic ablation study to validate the effectiveness of modality-wise attention module and sequence-wise attention module by removing them sequentially and comparing the performance with CMA-CLIP. Detailed results are shown in Table 6.

On MRWPA, the average recall at 90% precision across the 3 attributes drops from 53.4% to 47.5% after removing the modality-wise attention module. This is because the proportions of titles that contain tokens related to color, pattern, and style are 67%, 25% and 15% respectively. When a title does not contain any tokens related to an attribute, it becomes irrelevant for the classification of that attribute. The performance drop indicates that the modality-wise attention module significantly improves CMA-CLIP's robustness against noisy inputs.

To illustrate our model's robustness to input noise, in Table 5 we randomly pick some product image-title examples that CMA-CLIP is able to give correct classification whereas CMA-CLIP without modality-wise attention module cannot. We can clearly observe that in those examples, the product titles do not contain any tokens related to the attribute labels. Furthermore, for these examples, we complete the titles by adding the label related keywords and re-test them. This time, both methods can provide correct classification results which further proves that the modality-wise attention has the ability of filtering out irrelevant information (text without label information is considered as noise). We are not able to select similar examples in the Food101 and Fashion-Gen datasets, because in these two public datasets, there are no images or text that are irrelevant to the classification task.

The average recall drops from 47.5% to 45.8% after further removing the sequence-wise attention module. The sequence-wise attention module enhances the context-awareness of the image and text embedding by capturing the fine-grained correlation among image patches and text tokens, and the resulting embedding is expected to yield better results for classifications. The performance drop supports this conclusion.

We also visualize the result of sequence-wise attention for MR-WPA, Food101 and Fashion-Gen datasets in Figure 2. For each text input, we locate the token that is related to the classification task, and visualize the image patches that are most correlated to it by checking the inner product between the query embedding of the text token and the key embeddings of the image patches. In Figure 2, red regions are where the correlation is high. We observe that the sequence-wise attention is able to identify the highly correlated image patches and text tokens across all three datasets.

## 5 CONCLUSION

In this paper, we propose the CMA-CLIP, which unifies two types of cross-modality attentions: sequence-wise attention, a transformer based attention module that captures the fine-grained relationship between image patches and text tokens, and modality-wise attention, which learns the importance of image and text modalities in order to filter out the irrelevant modality for the classification task. We also design task specific modality-wise attentions and MLPs so that we can leverage a unified network for multi-task classifications. We evaluate our method on the MRWPA Dataset, the Food101 dataset and the Fashion-Gen dataset. CMA-CLIP outperforms the pre-trained and fine-tuned CLIP by an average of 11.9% in recall at the same level of precision on the MRWPA Dataset for the classifications for color, pattern, and style attributes. It also surpasses the state-of-the-art method on the Fashion-Gen Dataset by 5.5% in accuracy and achieves competitive performance on the Food101 Dataset. For the future work, we are interested in training CMA-CLIP with other datasets to enable the contrastive loss, improving CMA-CLIP's robustness against noisy labels, and also, exploring semi-supervised learning methods so that unlabeled image-text pairs can be leveraged in the training process to improve model generalizability.

## REFERENCES

[1] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. 2014. Food-101 – Mining Discriminative Components with Random Forests. In *European Conference on Computer Vision.* 446–461.

[2] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. 2020. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882* (2020).

[3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning.* PMLR, 1597–1607.

[4] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. 2020. Big Self-Supervised Models are Strong Semi-Supervised Learners. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 22243–22255. https://proceedings.neurips.cc/paper/2020/file/fcbc95ccdd551da181207c0c1400c655-Paper.pdf

[5] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. 2020. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297* (2020).

[6] Xinlei Chen and Kaiming He. 2021. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 15750–15758.

[7] Xinlei Chen, Saining Xie, and Kaiming He. 2021. An empirical study of training self-supervised vision transformers. *arXiv preprint arXiv:2104.02057* (2021).

[8] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Uniter: Universal image-text representation learning. In *European conference on computer vision.* Springer, 104–120.

[9] François Chollet. 2017. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 1251–1258.

[10] Karan Desai and Justin Johnson. 2021. Virtex: Learning visual representations from textual annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 11162–11173.

[11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers).* 4171–4186. https://aclanthology.org/N19-1423

[12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is

Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*. https://openreview.net/forum?id=YicbFdNTTy

[13] Dehong Gao, Linbo Jin, Ben Chen, Minghui Qiu, Peng Li, Yi Wei, Yi Hu, and Hao Wang. 2020. Fashionbert: Text and image matching with adaptive loss for cross-modal retrieval. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2251–2260.

[14] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, koray kavukcuoglu, Remi Munos, and Michal Valko. 2020. Bootstrap Your Own Latent - A New Approach to Self-Supervised Learning. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 21271–21284. https://proceedings.neurips.cc/paper/2020/file/f3ada80d5c4ee70142b17b8192b2958e-Paper.pdf

[15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9729–9738.

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[17] Zhicheng Huang, Zhaoyang Zeng, Yupan Huang, Bei Liu, Dongmei Fu, and Jianlong Fu. 2021. Seeing Out of tHe bOx: End-to-End Pre-training for Vision-Language Representation Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12976–12985.

[18] Zhicheng Huang, Zhaoyang Zeng, Bei Liu, Dongmei Fu, and Jianlong Fu. 2020. Pixel-BERT: Aligning Image Pixels with Text by Deep Multi-Modal Transformers. *CoRR* abs/2004.00849 (2020).

[19] Yuqi Huo, Manli Zhang, Guangzhen Liu, Haoyu Lu, Yizhao Gao, Guoxing Yang, Jingyuan Wen, Heng Zhang, Baogui Xu, Weihao Zheng, et al. 2021. WenLan: Bridging vision and language by large-scale multi-modal pre-training. *arXiv* abs/2103.06561 (2021).

[20] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. *arXiv* abs/2102.05918 (2021).

[21] Douwe Kiela, Suvrat Bhooshan, Hamed Firooz, Ethan Perez, and Davide Testuggine. 2019. Supervised multimodal bitransformers for classifying images and text. *arXiv* abs/1909.02950 (2019).

[22] Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. ViLT: Vision-and-Language Transformer Without Convolution or Region Supervision. In *Proceedings of the 38th International Conference on Machine Learning*, Vol. 139. 5583–5594.

[23] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 2013. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*. 554–561.

[24] Chenliang Li, Ming Yan, Haiyang Xu, Fuli Luo, Wei Wang, Bin Bi, and Songfang Huang. 2021. SemVLP: Vision-Language Pre-training by Aligning Semantics at Multiple Levels. *ArXiv* abs/2103.07829 (2021).

[25] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. 2020. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 11336–11344.

[26] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv* abs/1908.03557 (2019).

[27] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*. Springer, 121–137.

[28] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. ViLBERT: Pre-training Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2019/file/c74d97b01eae257e44aa9d5bade97baf-Paper.pdf

[29] Di Qi, Lin Su, Jia Song, Edward Cui, Taroon Bharti, and Arun Sacheti. 2020. Imagebert: Cross-modal pre-training with large-scale weak-supervised image-text data. *arXiv preprint arXiv:2001.07966* (2020).

[30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, Vol. 139. 8748–8763. http://proceedings.mlr.press/v139/radford21a.html

[31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems* 28 (2015), 91–99.

[32] Negar Rostamzadeh, Seyedarian Hosseini, Thomas Boquet, Wojciech Stokowiec, Ying Zhang, Christian Jauvin, and Chris Pal. 2018. Fashion-gen: The generative fashion dataset and challenge. *arXiv* abs/1806.08317 (2018).

[33] Mert Bulent Sariyildiz, Julien Perez, and Diane Larlus. 2020. Learning visual representations with caption annotations. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*. Springer, 153–170.

[34] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. 2012. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402* (2012).

[35] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. VL-BERT: Pre-training of Generic Visual-Linguistic Representations. In *International Conference on Learning Representations*. https://openreview.net/forum?id=SygXPaEYvH

[36] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, Los Alamitos, CA, USA, 2818–2826.

[37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc., 5998–6008.

[38] Cédric Villani. 2009. *Optimal transport: old and new*. Vol. 338. Springer.

[39] Long Xiang, Gan Chuang, Melo Gerard d, Liu Xiao, Li Yandong, Li Fu, and Wen Shilei. 2018. Multimodal Keyless Attention Fusion for Video Classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*.

[40] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. 2017. Aggregated Residual Transformations for Deep Neural Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 5987–5995.

[41] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *Proceedings of the 32nd International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 37)*, Francis Bach and David Blei (Eds.). PMLR, Lille, France, 2048–2057. https://proceedings.mlr.press/v37/xuc15.html

[42] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. 2020. Contrastive learning of medical visual representations from paired images and text. *arXiv* abs/2010.00747 (2020).

[43] Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. 2020. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 13041–13049.

[44] Mingchen Zhuge, Dehong Gao, Deng-Ping Fan, Linbo Jin, Ben Chen, Haoming Zhou, Minghui Qiu, and Ling Shao. 2021. Kaleido-BERT: Vision-Language Pre-training on Fashion Domain. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12647–12657.