

Visual Transformers with Primal Object Queries for Multi-Label Image Classification

Vacit Oguz Yazici

Computer Vision Center/Wide Eyes Technologies

Barcelona, Spain

Email: voyazici@cvc.uab.es

Joost van de Weijer

Computer Vision Center,

Universitat Autònoma de Barcelona

Email: joost@cvc.uab.es

Longlong Yu

Wide-Eyes Technologies

Barcelona, Spain

Email: longyu@wide-eyes.it

Abstract—Multi-label image classification is about predicting a set of class labels that can be considered as orderless sequential data. Transformers process the sequential data as a whole, therefore they are inherently good at set prediction. The first vision-based transformer model, which was proposed for the object detection task introduced the concept of *object queries*. Object queries are learnable positional encodings that are used by attention modules in decoder layers to decode the object classes or bounding boxes using the region of interests in an image. However, inputting the same set of object queries to different decoder layers hinders the training: it results in lower performance and delays convergence. In this paper, we propose the usage of *primal object queries* that are only provided at the start of the transformer decoder stack. In addition, we improve the mixup technique proposed for multi-label classification. The proposed transformer model with *primal object queries* improves the state-of-the-art class wise F1 metric by 2.1% and 1.8%; and speeds up the convergence by 79.0% and 38.6% on MS-COCO and NUS-WIDE datasets respectively.

I. INTRODUCTION

The task of predicting the presence of visual concepts in images is known as multi-label classification. The visual concepts in general refer to a set of objects, but could also refer to other visual concepts such as attributes. Multi-label classification is difficult because of the wide range of classes that is typically considered, the wide variety of scales in which these classes can occur, and the complex inter-dependencies between classes [1], [2].

The field of multi-label image classification has seen much progress in recent years. Earlier works exploited graphical models to model label relations explicitly [3], [4]. Then, CNN-RNN models were proposed [5]–[7] to capture label correlations. Later, to learn label dependencies explicitly, graph convolutional networks were proposed [8]. Even though significant progress has been made in multi-label image classification in recent years, systems still suffer from problems common to recursive methods (such as modeling long-term dependencies) or fail to capture the complex relations between the many visual concepts involved in multi-label classification.

Transformers were first proposed in [9]. Unlike recurrent models, transformers process data simultaneously. In case of recurrent models, if the decoding process is interrupted the decoder is forced to output a *termination* token, which will cause the not-yet attended classes to be missed. Due to the non-sequential nature of transformers, they do not suffer from

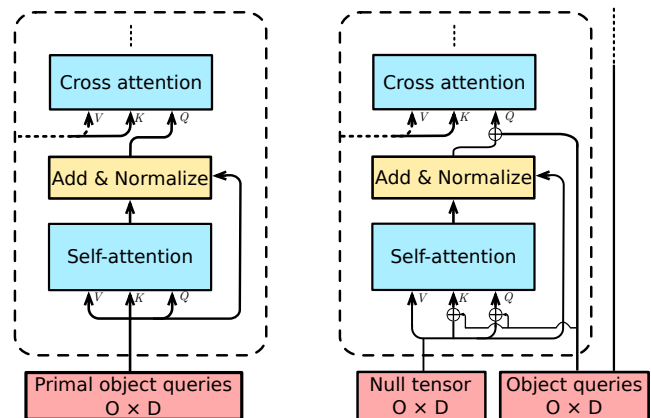


Fig. 1. Comparison between the primal object queries and object queries. Primal object queries are given directly to the first decoder layer instead of being added as positional encodings in every decoder layer. This leads to faster convergence and better performance.

this problem. Therefore, although they were firstly used for various NLP tasks [10]–[12], they became also widely used for other tasks. The first visual CNN-transformer model for a vision task was proposed for object detection [13]. Recently, a novel and pure transformer model that uses a sequence of image patches as input data was proposed for image classification [14]. Lanchantin et al. [15] proposed the first transformer model for multi-label image classification.

Carion et al. [13] introduced the concept of *object queries*. As a set of learnable positional encodings, they are added to query and key tensors in attention modules of every decoder layer. However, using the same set of object queries in different decoder layers is not optimal for training due to the different set of relations learned by each of the layers. As we empirically demonstrate in the experimental section, it leads to lower performance and slower convergence. Therefore, in this paper, we introduce *primal object queries* that differ from standard object queries in the way that they are input to a transformer decoder stack (see Figure 1). Moreover, we improve the mixup technique for multi-label classification to achieve significantly better results.

The main contributions in this paper are:

- We introduce *primal object queries* that obtain better results and yield faster convergence than standard object

queries by 79.0% and 38.6% on MS-COCO and NUS-WIDE datasets respectively.

- We improve the mixup for multi-label classification and show that it results in a significant improvement.
- Evaluation shows that we improve the state-of-the-art class wise F1 score by 2.1% and 1.8% on the MS-COCO and NUS-WIDE datasets respectively.

II. RELATED WORK

Transformers. Transformers were first proposed in [9], and since became the state-of-the-art approach for sequence-to-sequence tasks such as machine translation and visual question answering. The transformer self-attention module attends all the sequential data at once, hence handling the long sequences better than RNNs, which struggle with long-term dependencies. A challenging aspect of transformers is the large number of parameters they require, which require large amounts of data to fit properly. Since such large-scale training datasets are scarce, and an important limitation to use transformer models in many practical applications, Devlin et al. [10] proposed a way of training transformers in an unsupervised manner on readily available large unsupervised text corpus, and showed that with a simple fine-tuning procedure state-of-the-art results could be achieved on various tasks. This led to a wider popularization of transformer models, and usage in areas different than the originally proposed ones, such as image captioning [11], [12] and object detection [13]. Carion et al. [13] introduced the new concept of object queries to be given input for the transformer through the training. These object queries are learned through the training, and according to the analysis done by the authors, each object query learns to focus on different areas of images and box sizes. Dosovitskiy et al. [14] proposed the first pure transformer model for image classification. They split an image into smaller patches and convert it to a sequential data (which also consists of tokens) to be processed by transformer encoder layers. Although they showed that transformers can give promising results for pure vision tasks, they still fell behind other convolutional models. Yuan et al. [16] proposed a new tokenization system that reduced the number of parameters and achieved comparable results with other convolutional models. Recently, Lanchantin et al. [15] introduced the first transformer model for multi-label image classification. They exploited self-attention modules to learn label dependencies for the purpose of predicting a set of labels given a set of masked label embeddings and image features.

Multi-label Classification. Multi-label classification seeks to predict a variable number of labels for every single image, ideally capturing all the relevant visual concepts, such as objects or attributes, that appear on the image. Traditional methods for multi-label classification ignore label correlation, which can help boost the performance of certain under-represented classes. The few early works that tried to leverage label correlations for multi-label classification exploited graphical models such as Conditional Random Fields (CRFs) [3] or Dependency

Networks [4]. More recently, the idea to exploit RNN models to capture label correlations was proposed in [17] and [5], where the low dimensional internal state of the network was used to model label dependencies. Wang et al. [5] combined CNN and RNN architectures, and learned a joint image-label embedding space to learn label dependencies. However, since LSTMs produce sequential outputs, a fixed order was imposed to the labels, and the model learned to predict the output in the same order. This led to problems such as skipping predictions for classes that appear earlier in the sequence but less relevant in the image. Yazici et al. [7], proposed a CNN-RNN model which they trained with an orderless loss function to avoid the drawbacks of imposing a fixed label order in RNN. Finally, ML-GCN [18], exploited graph convolutional networks to capture label dependencies.

Mixup. Zhang et al. [19] proposed to blend images and their associated labels randomly to improve generalization of the models. It was shown that the mixup was beneficial to avoid overconfident predictions in several tasks such as image classification [19], [20], object detection [21], text classification [20] and semantic segmentation [22]. Verma et al. [23] proposed to combine hidden states of paired samples in addition to their images and labels. Islam et al. [22] conditioned the mixup on the labels of paired samples. If the mixup is done without any constraints, then the majority of the pairs will mostly include the most frequent classes. To avoid the model to have a strong bias for frequent classes, they combined images based on a uniform distribution across categories. Wang et al. [24] is the first work that exploited the mixup technique for multi-label image classification. Although the authors did not report any improvement over the baselines in case of single models, they noted that an ensemble of models trained with mixup achieved better results.

III. METHOD

Our method for the multi-label classification is based on a transformer architecture. We try different strategies to assign labels to the decoder output. In this section, we briefly introduce some transformer concepts, explain our proposed architecture with different losses and, finally, present our different adaptations of the mixup technique.

A. Transformers

The main component of transformers is a self-attention module [9] which uses a set of weights W to compute the *query* (Q), *key* (K) and *value* (V) vectors with the input vector:

$$Q = W_Q X, \quad K = W_K X, \quad V = W_V X \quad (1)$$

Then, these three vectors are combined to compute the output of the self-attention layer:

$$A = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V \quad (2)$$

The result of the dot product between the *query* and the *key* is divided by the square root of the dimensionality D to have smoother softmax values. Many encoder layers, consisting of

TABLE I
COMPARISON OF THE PERFORMANCES OF DETR, DETR* AND T-POQ MODELS WITH DIFFERENT NUMBER OF DECODER LAYERS ON MS-COCO.

	# of decoder layers								
	2			3			4		
	C-P	C-R	C-F1	C-P	C-R	C-F1	C-P	C-R	C-F1
DETR	75.9	64.7	69.9	74.9	64.6	69.4	76.8	64.2	69.9
DETR*	75.9	65.9	70.5	75.9	66.1	70.7	76.5	65.4	70.5
T-POQ	76.6	66.0	70.9	76.5	66.1	70.9	73.8	67.1	70.3

a self-attention and a feed-forward neural network are stacked in the encoder part of the network, and the output of the final encoder layer (A) is passed to the *cross attention module* in every decoder layer. The cross attention module has the same structure as the self attention module, with the only difference that the output of the final encoder layer is used as the key and value vectors for every decoder layer during the forward pass, while the query is obtained from the input of the decoder.

The input of the transformer decoder stack depends on the task and design of the architecture. In case of an NLP task, it might be class embeddings [12] or masked output embeddings [9]. In case of a vision-based task, it might be a set of object queries [13]. In [13], the object queries are added to the query input of each decoder layer (the query input of the first decoder layer consists of zeros as can be seen in Figure 1). The fact that object queries are provided to all decoder layers complicates their training: we conjecture that the requirement to be useful at multiple hierarchical levels (and therefore at different semantic levels) is hard to fulfill. In order to verify this, we conduct several experiments where we compare the approach in [13] (denominated as DETR) with inputting unique sets of object queries to each decoder layer (denominated as DETR*). The first two rows of Table I, show the performance of the two models with different number of decoder layers. DETR* obtains significantly better results than DETR which confirms our assumption that inputting the same set of object queries to different decoder layers is not optimal. However, although learning a separate set of object queries for each decoder layer improves the results, it might be computationally redundant and not feasible when the number of decoder layers increases.

B. Overview of architecture

Our architecture consists of two parts: a backbone that processes the input image, and a transformer, which can be further divided into encoder and decoder (see Figure 2).

a) Backbone: We use a convolutional neural network to obtain a feature representation of the input image, but since the transformer requires sequential data, we linearize the feature map of the last convolutional layer of the backbone network along the spatial dimensions, and use them as the input sequence for the encoder. More precisely, let the last convolutional layer of the CNN output a feature map of size $H \times W \times C$, where H and W are the spatial dimensions, and C the number of channels, we reshape it to $HW \times C$ to obtain a sequence of feature vectors with dimension C .

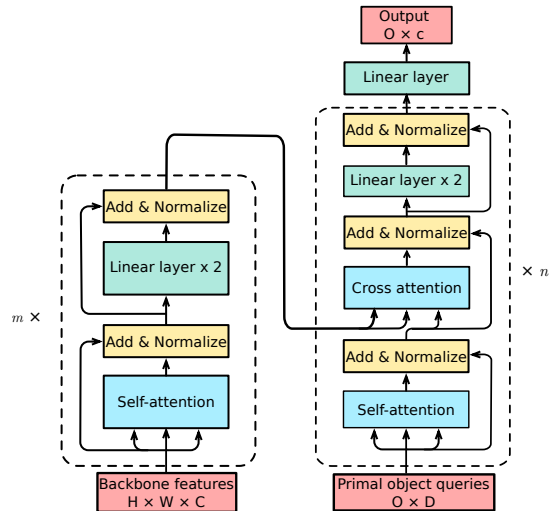


Fig. 2. The overall transformer architecture used in this paper.

b) Transformer encoder and decoder: In the transformer encoder, the linearized feature map is passed through the self-attention and linear layers, and then the result is passed to the cross attention module of every decoder layer. The initial linear layer reduces the dimension to D . We do not add learnable or fixed positional encodings to encoder features since it does not give any improvements. We attribute this to using semantically strong features that do not require additional spatial information for the task of multi-label classification. We quantitatively verified that adding positional encodings is only beneficial when features from earlier layers are used, however, that also leads to lower results.

The input of the first decoder layer are *primal object queries* which differ from the object queries proposed in [13] in the way that they are given to transformer decoder stack. In our model, primal object queries are given directly to the first decoder layer instead of being added as positional encodings in every decoder layer. In this way, we avoid the drawbacks of the standard object queries which were mentioned in the previous section. In Table I, we also added the results of our method (T-POQ) and it can be seen that without the overhead of learning separate object queries for each decoder layer (as does DETR*) our method obtains equal or better results. Interestingly, in the experimental section we will also show that T-POQ obtains a significant speed-up in training when compared to the DETR model. Finally, the number of primal object queries is O while the dimension of the queries is D . Then, a linear layer outputs a tensor whose shape is $O \times c$ where c is the number of classes.

c) Training losses: Using object queries for multi-label classification is a new concept. Previously, e.g. in CNN-RNN models, the labels are ordered either dynamically [7] or an imposed fixed-order is applied [5]. The forward pass is done in a recursive way, therefore the length of the output tensor of the RNN is bounded by the number of labels of an image that has the most labels. However, in case of transformers, the

forward pass is not done recursively which gives the flexibility of determining the size of the output tensor. We will consider two ways to train transformers for multi-label classification.

Firstly, we consider the case where we align each object query with a specific class. In this case $O = c$ and each object query specializes in detecting a single object class. In case of non-existent labels, *empty* tokens are assigned to the object queries that are in charge of these labels. We call this the *exhaustive model*. Here is the loss equation for the exhaustive model:

$$S_{ji} = \frac{e^{x_{ji}}}{\sum_{i=1}^c e^{x_{ji}}}, \quad \mathcal{L}_{exh} = -\frac{1}{O} \sum_{j=1}^O \sum_{i=1}^c y_{ji} \log S_{ji} \quad (3)$$

$$\text{subject to } y_{ji} \in \{0, 1\}, \\ \sum_j y_{ji} = 1 \quad \forall i \in L \quad \text{and} \quad \sum_j y_{ji} = 0 \quad \forall i \notin L$$

where $y_{ji} \in \mathbb{N}^{O \times c}$, x_{ji} and L are the class labels, output tensor and set of class labels, respectively. A drawback of the exhaustive approach is that it scales linearly with the number of classes, and might be infeasible for datasets with many labels. However, reducing the number of object queries requires one object query to be in charge of multiple labels, which results in an assignment problem. Therefore, we employ an orderless loss inspired by the one proposed in [7]. We call this model the *aligned model* and it no longer requires to scale linearly with the number of classes in the dataset:

$$\mathcal{L}_{align} = \min_y \sum_{i=1}^c y_{ji} \log S_{ji} \quad (4)$$

$$\text{subject to } y_{ji} \in \{0, 1\}, y_{ji} = 1 \text{ if } \hat{l}_j \in L \text{ and } i = \hat{l}_j, \\ \sum_j y_{ji} \geq 1 \quad \forall i \in L \quad \sum_j y_{ji} = 0 \quad \forall i \notin L$$

where \hat{l}_j is the class predicted by the model at object query j . The order of the labels in y is chosen in such a way that it gives the minimum cross entropy loss. This label assignment problem can be solved with the Hungarian algorithm. In addition, we impose the constraint, which was proposed in [7], that assigns the class \hat{l}_j to object query j if \hat{l}_j belongs to L . Therefore, the same class may be assigned to several object queries.

C. Mixup

Mixup was proposed in [19] and was found to significantly improve results for image classification, object detection and NLP tasks [19]–[21]. We consider three ways of adapting the mixup technique to the training: soft, hard and restricted hard.

To train the model, a dataset with pairs of images and sets of labels is used. Let (I, T) be the pairs in a batch containing N images $I = \{i_1, i_2, \dots, i_N\}$ and labels $T = \{t_1, t_2, \dots, t_N\}$, $t_i \in \{0, 1\}^c$ where c is the number of classes in the dataset. For the *soft mixup*, which is the original mixup as proposed in [19], we sample random weights from a beta distribution $\lambda \sim \text{Beta}(\alpha, \alpha)$, $\alpha \in (0, \infty)$ to use them to mixup images and their associated labels:

$$i_m = \lambda i_i + (1 - \lambda) i_j \\ t_m = \lambda t_i + (1 - \lambda) t_j \quad (5)$$

where i_i and i_j are randomly selected images, and i_m is the mixed up version.

The soft mixup makes sense for single class image classification where the last layer is typically a softmax. However, for multi-label image classification multiple labels can be present in the image. Therefore, we use the mixup proposed in [24] and denominate it as *hard mixup* where the union of labels is taken instead of the average. As proposed by the author, we alternatively enable and disable the application of mixup in every epoch and we use the ratio of 0.5 : 0.5 for images. In order to apply the mixup in all epochs and have both mixed and non-mixed images in a batch we consider the *restricted hard mixup* as a final setup, where we apply the mixup in every epoch and restrict it by applying it only to half of the images in the batch. We use the last half of the batch to mix with the first half. For instance, if we set the batch size to 4, after the mixup the images and labels become $I = \{(i_1 + i_3)/2, (i_2 + i_4)/2, i_3, i_4\}$ and $T = \{t_1 \cup t_3, t_2 \cup t_4, t_3, t_4\}$ respectively.

IV. EXPERIMENTS

Datasets and setting. We evaluate our models on MS-COCO [25] and NUS-WIDE [26] datasets. **MS-COCO** consists of 82,081 training and 40,137 test images for 80 object categories. **NUS-WIDE** consists of 269,648 images with a total number of 5,018 unique labels. However, annotations for 81 labels are more trustworthy and used for evaluation. After removing images that do not belong to the 81 labels, 209,347 images remain.

Evaluation metrics. We use *per-class* and *overall* precision, recall and F1 scores. For the MS-COCO, we also report the mean average precision (mAP) score.

Network training. The number of object queries is set to 25 for all datasets. The internal dimension of the transformer is 512. The backbone network, which is pre-trained on ImageNet [27], uses an SGD optimizer with learning rate 0.001 and momentum 0.9. The rest of the model is trained with the ADAM optimizer with a learning rate 0.0001 for 40 epochs. The batch size is 32, and random affine transformations and contrast changes are applied as data augmentation. The batch norm layers in the backbone are frozen during the training¹.

A. Ablation

All the ablation studies are done on MS-COCO dataset. and are repeated three times. The final result is the average of the three experiments. For validation we use 25% of training data and the rest as train set. The snapshot that gives the highest score on the validation set is evaluated on the test set. The input image size is 288. Unless stated otherwise, the used transformer model is the aligned model.

In Table II, LSTM and transformer models with different encoder and decoder layers are compared. The LSTM model is trained with an orderless loss [7] and has the same backbone and hidden size as the transformer model. From the results it can be seen that adding encoder layers does not yield

¹<https://github.com/voyazici/visual-transformers-classification>

any significant improvement. When we evaluate the attention maps generated by the self-attention module in the encoder, unlike in [13], the encoder does not separate instances which is not crucial for multi-label classification unlike for object-detection. On the other hand, more decoder layers do lead to improvement in the recall metric. Moreover, our transformer model slightly increases the total number of parameters and computational cost (55.2M and 27.1 GFLOPS) compared to the backbone model (44.6M and 25.9 GFLOPS).

TABLE II
COMPARISON OF LSTM AND DIFFERENT TRANSFORMER MODELS.

Model	# of enc.	# of dec.	mAP	C-P	C-R	C-F1	O-P	O-R	O-F1
LSTM	-	-	75.3	76.2	66.7	71.1	78.8	71.3	74.9
Trans.	0	1	78.0	79.6	67.9	73.3	81.8	71.9	76.5
Trans.	1	1	77.6	77.9	68.8	73.1	80.2	72.6	76.2
Trans.	1	2	78.3	79.6	68.2	73.5	80.8	72.4	76.3
Trans.	2	2	78.2	78.2	69.3	73.5	80.2	73.1	76.5
Trans.	2	3	77.9	77.4	69.6	73.3	79.3	73.6	76.4

TABLE III
SUBTRACTION OF LSTM SCORES FROM TRANSFORMER SCORES ON DIFFERENT NUMBER OF LABELS PER IMAGE.

	Number of labels per image										Avg
	1	2	3	4	5	6	7	8	9	+10	
precision	2.1	1.5	1.5	2.4	1.3	1.8	1	-0.5	0.7	0.0	0.8
recall	0.5	1	1.1	1.7	1.2	1.5	1.2	2	2.5	2.1	1.7
F1	1.5	1.3	1.3	2	1.3	1.6	1.1	1.1	1.9	1.5	1.5

In Table III, we compare the performance of the LSTM and transformer model (the one with one encoder and two decoder layers) on images that have different number of labels. The values are the subtraction of average LSTM scores from average transformer scores. When the number of labels increases, the transformer model misses fewer classes that leads to fewer false negatives and higher recall.

In Table IV, aligned and exhaustive models (one encoder and two decoder layers) are compared. The results of the models are comparable. However, for the comparison with the state-of-the-art models, we will use the aligned model. It is a more compact model and the overhead cost of the alignment step (0.9 ms per image) is negligible. Also in Table IV, different mixup setups are compared. The α value for the soft mixup is 0.4. Mixup improves the results considerably. The best result is obtained with the restricted hard mixup. Higher precision and mAP scores show that the restricted hard mixup model has more confident predictions. We attribute the superiority of the restricted hard mixup over the soft mixup to soft labels being detrimental for modelling label correlations; and the superiority over the hard mixup to lower variance

TABLE IV
COMPARISON OF EXHAUSTIVE AND ALIGNED MODELS AND MIXUP METHODS.

	mAP	C-P	C-R	C-F1	O-P	O-R	O-F1
Exhaustive	78.3	77.5	69.5	73.2	79.8	73.3	76.4
Aligned	78.1	77.6	69.4	73.3	79.9	73.3	76.4
Aligned + soft mixup	78.6	79.1	69.8	74.2	80.9	73.7	77.1
Aligned + hard mixup	79.0	79.7	69.4	74.2	81.4	73.3	77.1
Aligned + restr. hard mixup	79.6	80.2	69.7	74.6	82.1	73.5	77.5

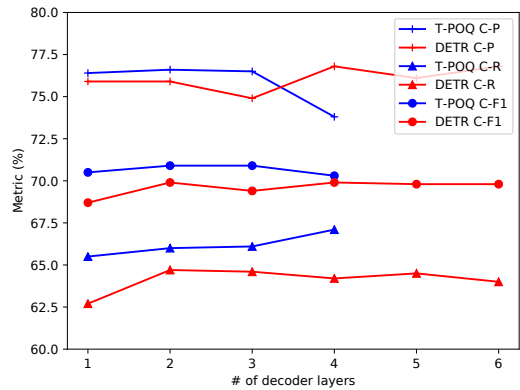


Fig. 3. Comparison of the proposed and DETR approaches with different number of decoder layers on MS-COCO.

during the gradient update due to having mixed and non-mixed samples in the same batch.

B. Object Queries

In Figure 3, we compare our primal object queries with the object queries in DETR [13]. We show the change of C-P, C-R and C-F1 metrics with different number of decoder layers. For simplicity, we do not employ the mixup and the image size is 224×224 . We set the learning rate of the backbone to 0.0001 for the DETR model to make the model converge during training. It can be seen that our approach yields significantly higher C-F1 scores in every setup. Also, with more decoder layers it achieves higher recall. In addition, the model with one decoder layer already performs comparable with the models that have more decoder layers. On the other hand, the model with the DETR approach requires at least two decoder layers to achieve comparable performance. We attribute this superiority to the residual connection after the cross attention module which enables the propagation of our primal object queries to the next layer. We empirically confirm this by disabling the residual connection and obtaining the same results as the DETR model in the case that the number of decoder layers is one. The same fact causes 79.0% and 38.6% faster convergence on MS-COCO and NUS-WIDE datasets respectively. In order to compare the convergence between the two models, we determine the epoch number at which C-F1 stops improving for both models. Then, we calculate the percentage change of the epoch number for setups that have up to three decoder layers. Finally, we average the percentages from different setups to get the relative change in the convergence speed. In Figure 4, the DETR model starts to converge much later than our model (the setup with three decoder layers). Moreover, due to the lack of the propagation of the object queries, the DETR model starts from a significantly lower point compared to our model. Consequently, when we compare the best models, we achieve improvements over the DETR approach by 1.2%-1.5% in all metrics.

Next, we analyze what the primal object queries of the aligned model learn when they must account for multiple

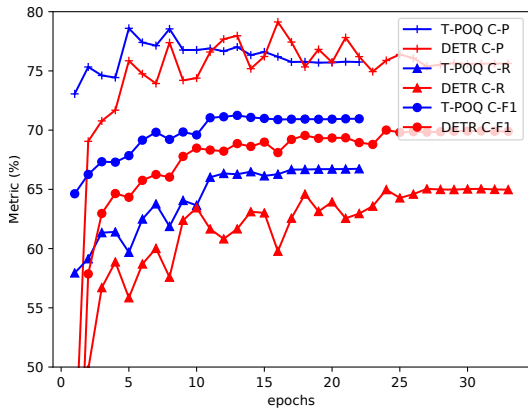


Fig. 4. Convergence of the proposed and DETR models on MS-COCO.

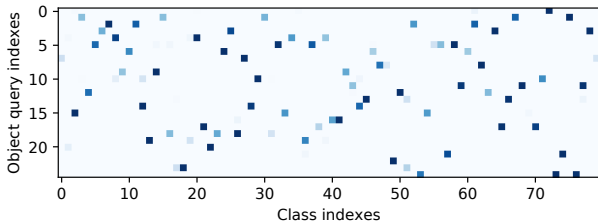


Fig. 5. Normalized counts of predicted classes. Each object query recognizes a subset of classes.

classes with the orderless loss. In Figure 5 (darker shades indicate higher values), we display the normalized counts of the predicted classes for each object query. Each object query learns to recognize a subset of classes. Since each of them can only make one prediction during one forward pass, the subsets consist of classes that are not likely to exist together. For example, the fourth object query learns to recognize *bear*, *cow*, *suitcase*, and *wine glass*.

C. Comparison with the SOTA

We compare our results with several recent models: CNN-RNN [5], SR CNN-RNN [6], Chen et al. [28], Li et al. [29] and PLA [7] are models that include RNNs either to model label relations or recursively generate attention maps. SRN [30] proposes a Spatial Regularization Network that generates attention maps for all labels. ACfs [8], proposes a two-branch network with an original image and its transformed image as inputs and imposes an additional loss to ensure the consistency between attention maps of both versions. DER [31] proposes to train a detection model in three steps in which it learns class-aware attention maps, models label correlations explicitly and measures similarity between label embeddings. ML-GCN [18], exploits graph convolutional networks to capture label dependencies. C-Tran [15] is the only transformer model that we compare with. It exploits self-attention layers in a transformer encoder to learn label correlations given image features and a set of masked label embeddings and does not include any transformer decoder layer unlike our model. For MS-COCO

TABLE V
COMPARISON WITH STATE-OF-THE-ART ON MS-COCO.

Methods	Image size	mAP	C-P	C-R	C-F1	O-P	O-R	O-F1
SRN [30]	224	77.1	81.6	65.4	71.2	82.7	69.9	75.8
ACFs [8]	288	77.5	77.4	68.3	72.2	79.8	73.1	76.3
PLA [7]	288	-	80.4	68.9	74.2	81.5	73.3	77.1
ML-GCN [18]	448	83.0	85.1	72.0	78.0	85.8	75.4	80.3
DER [31]	448	82.9	84.7	71.6	77.6	86.0	74.9	80.0
C-Tran [15]	576	85.1	86.3	74.3	79.9	87.7	76.5	81.7
T-POQ	224	77.9	79.5	67.4	73.0	81.5	71.3	76.1
T-POQ	288	80.6	80.9	70.9	75.6	82.5	74.4	78.2
T-POQ	448	84.5	82.9	75.8	79.2	84.4	78.4	81.3
T-POQ	576	86.2	84.1	77.9	80.9	85.0	80.6	82.8

TABLE VI
COMPARISON WITH STATE-OF-THE-ART ON NUS-WIDE.

Methods	C-P	C-R	C-F1	O-P	O-R	O-F1
CNN-RNN [5]	40.5	30.4	34.7	49.9	61.7	55.2
Chen et al. [28]	59.4	50.7	54.7	69.0	71.4	70.2
SR CNN-RNN [6]	55.7	50.2	52.8	70.6	71.4	71.0
Li et al. [29]	44.2	49.3	46.6	53.9	68.7	60.4
LSEP [33]	66.7	45.9	54.4	76.8	65.7	70.8
PLA [7]	60.7	52.4	56.2	72.0	72.8	72.4
T-POQ	66.0	52.7	58.6	74.7	71.8	73.2
SRN [30]	65.2	55.8	58.5	75.5	71.5	73.4
DER [31]	64.2	57.9	60.9	75.5	73.0	74.2
T-POQ	66.5	56.0	60.8	75.1	73.2	74.1
T-POQ*	66.8	57.8	62.0	75.6	74.3	74.9

experiments ResNet-101 architecture is used for the backbone, and for NUS-WIDE we run experiment with both ResNet101 and VGG16. For the comparison with SOTA models, we use the aligned transformer model with one encoder and two decoder layers. We also employ the restricted hard mixup.

We outperform all the state-of-art on MS-COCO (see Table V). The performance superiority is more apparent in the recall metrics, since the transformer model is less likely to miss classes. The results on the NUS-WIDE dataset can be seen in Table VI. The results on the top part of the table use the split proposed by [32] with a VGG16 backbone, while the ones on the lower part use the original split and a ResNet-101 backbone. For the T-POQ* model, we resize the input image to 448×448 to make the comparison fair with the DER [31]. For the T-POQ model, we resize it to 224×224 to make the comparison fair with SRN [30]. We surpass all other models, especially in the class-wise metrics which are more relevant, since NUS-WIDE is an unbalanced dataset.

V. CONCLUSIONS

We introduced the *primal object queries* that achieved significantly better results and a large speed-up of training convergence for both MS-COCO and NUS-WIDE datasets. Our model is unique in that it achieves to learn long-term dependencies, adapts and integrates the mixup technique for multi-label classification successfully, and obtains state-of-the-art results for multi-label classification on the MS-COCO and NUS-WIDE datasets by a large margin.

Acknowledgements: We acknowledge the Spanish projects PID2019-104174GB-I00 and the Industrial Doctorate Grant 2016 DI 039 of the Ministry of Economy and Knowledge of the Generalitat de Catalunya, and its CERCA Program.

REFERENCES

- [1] W. Bi and J. Kwok, "Efficient multi-label classification with many labels," in *International Conference on Machine Learning*. PMLR, 2013, pp. 405–413.
- [2] J. Wehrmann, R. Cerri, and R. Barros, "Hierarchical multi-label classification networks," in *International Conference on Machine Learning*. PMLR, 2018, pp. 5075–5084.
- [3] N. Ghamrawi and A. McCallum, "Collective multi-label classification," in *Proceedings of the 14th ACM International Conference on Information and Knowledge Management*. ACM, 2005, pp. 195–200.
- [4] Y. Guo and S. Gu, "Multi-label classification using conditional dependency networks," in *Twenty-Second International Joint Conference on Artificial Intelligence*, 2011.
- [5] J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, and W. Xu, "Cnn-rnn: A unified framework for multi-label image classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2285–2294.
- [6] F. Liu, T. Xiang, T. M. Hospedales, W. Yang, and C. Sun, "Semantic regularisation for recurrent image annotation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2872–2880.
- [7] V. O. Yazici, A. Gonzalez-Garcia, A. Ramisa, B. Twardowski, and J. v. d. Weijer, "Orderless recurrent models for multi-label classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 440–13 449.
- [8] H. Guo, K. Zheng, X. Fan, H. Yu, and S. Wang, "Visual attention consistency under image transforms for multi-label image classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 729–739.
- [9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "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)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186.
- [11] M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara, "Meshed-memory transformer for image captioning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 578–10 587.
- [12] X. Zhu, L. Li, J. Liu, H. Peng, and X. Niu, "Captioning transformer with stacked attention modules," *Applied Sciences*, vol. 8, no. 5, p. 739, 2018.
- [13] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision*. Springer, 2020, pp. 213–229.
- [14] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2021.
- [15] J. Lanchantin, T. Wang, V. Ordonez, and Y. Qi, "General multi-label image classification with transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 16 478–16 488.
- [16] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z. Jiang, F. E. Tay, J. Feng, and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," *arXiv preprint arXiv:2101.11986*, 2021.
- [17] J. Jin and H. Nakayama, "Annotation order matters: Recurrent image annotator for arbitrary length image tagging," in *International Conference on Pattern Recognition*. IEEE, 2016, pp. 2452–2457.
- [18] Z.-M. Chen, X.-S. Wei, P. Wang, and Y. Guo, "Multi-label image recognition with graph convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5177–5186.
- [19] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," in *International Conference on Learning Representations*, 2018.
- [20] S. Thulasidasan, G. Chennupati, J. Bilmes, T. Bhattacharya, and S. Michalak, "On mixup training: Improved calibration and predictive uncertainty for deep neural networks," in *Advances in Neural Information Processing Systems*, 2019.
- [21] Z. Zhang, T. He, H. Zhang, Z. Zhang, J. Xie, and M. Li, "Bag of freebies for training object detection neural networks," *arXiv preprint arXiv:1902.04103*, 2019.
- [22] M. A. Islam, M. Kowal, K. G. Derpanis, and N. D. Bruce, "Feature binding with category-dependant mixup for semantic segmentation and adversarial robustness," in *British Machine Vision Conference*, 2020.
- [23] V. Verma, A. Lamb, C. Beckham, A. Najafi, I. Mitliagkas, D. Lopez-Paz, and Y. Bengio, "Manifold mixup: Better representations by interpolating hidden states," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6438–6447.
- [24] Q. Wang, N. Jia, and T. P. Breckon, "A baseline for multi-label image classification using an ensemble of deep convolutional neural networks," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 644–648.
- [25] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision*. Springer, 2014, pp. 740–755.
- [26] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng, "Nus-wide: a real-world web image database from national university of singapore," in *Proceedings of the ACM International Conference on Image and Video Retrieval*, 2009, pp. 1–9.
- [27] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [28] S.-F. Chen, Y.-C. Chen, C.-K. Yeh, and Y.-C. F. Wang, "Order-free rnn with visual attention for multi-label classification," in *AAAI Conference on Artificial Intelligence*, 2018.
- [29] L. Li, S. Wang, S. Jiang, and Q. Huang, "Attentive recurrent neural network for weak-supervised multi-label image classification," in *2018 ACM Multimedia Conference on Multimedia Conference*. ACM, 2018, pp. 1092–1100.
- [30] F. Zhu, H. Li, W. Ouyang, N. Yu, and X. Wang, "Learning spatial regularization with image-level supervisions for multi-label image classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5513–5522.
- [31] Z. Chen, Q. Cui, X.-S. Wei, X. Jin, and Y. Guo, "Disentangling, embedding and ranking label cues for multi-label image recognition," *IEEE Transactions on Multimedia*, 2020.
- [32] J. Johnson, L. Ballan, and L. Fei-Fei, "Love thy neighbors: Image annotation by exploiting image metadata," in *International Conference on Computer Vision*, 2015, pp. 4624–4632.
- [33] Y. Li, Y. Song, and J. Luo, "Improving pairwise ranking for multi-label image classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3617–3625.