

General Multi-label Image Classification with Transformers

Jack Lanchantin, Tianlu Wang, Vicente Ordonez, Yanjun Qi
University of Virginia

{jjl15sw,tianlu,vicente,yq2h}@virginia.edu

A. Appendix

A.1. Qualitative Examples

Inference with Partial Labels. In Figure 2, we show qualitative results on COCO-80 demonstrating the use of partial labels. In these examples, we first show the predictions for ResNet-101, as well as C-Tran without using partial labels. The last column shows the C-Tran predictions when using $\epsilon = 25\%$ partial labels (which is 21 labels for COCO-80) as observed, or known prior to inference. For many examples, certain labels cannot be predicted well without using partial labels.

Inference with Extra Labels. In Figure 3, we show qualitative results on CUB-312 demonstrating the use of extra labels. In the CUB-312 dataset, the extra labels are high level concepts of bird species that are not target labels. In these examples, we first show the predictions for C-Tran without using extra labels labels, and the last column shows the C-Tran predictions when using $\epsilon = 54\%$ of the extra labels (which is 60 labels for CUB-312) as observed, or known prior to inference. We can see that many bird species predictions are completely changed after using the extra labels as input to our model.

A.2. Detailed Diagram of C-Tran Settings

Figure 1 shows a detailed diagram of all possible training and inference settings used in our paper, and how C-Tran is used in each setting. By using the same random mask training, we can apply our model to any of the three inference settings.

A.3. Multi-Label Classification Metrics

$$\begin{aligned} \text{OP} &= \frac{\sum_i N_i^c}{\sum_i N_i^p} \\ \text{OR} &= \frac{\sum_i N_i^c}{\sum_i N_i^g} \\ \text{OF1} &= \frac{2 \times \text{OP} \times \text{OR}}{\text{OP} + \text{OR}} \end{aligned} \quad (1)$$

$$\begin{aligned} \text{CP} &= \frac{1}{C} \sum_i \frac{N_i^c}{N_i^p} \\ \text{CR} &= \frac{1}{C} \sum_i \frac{N_i^c}{N_i^g} \\ \text{CF1} &= \frac{2 \times \text{CP} \times \text{CR}}{\text{CP} + \text{CR}} \end{aligned} \quad (2)$$

where C is the number of labels, N_i^c is true positives for the i -th label, N_i^p is the total number of images for which the i -th label is predicted, and N_i^g is the number of ground truth images for the i -th label.

A.4. More Discussions of C-Tran

Connecting to Transformers and BERT. Our proposed method, C-Tran, draws much inspiration from works in natural language processing. The transformer model [7] proposed “self attention” for natural language translation. Self attention allows each word in the target sentence to attend to all other words (both in the source sentence and the target sentence) for translation. [2] introduced BERT for language modeling. BERT uses self attention with masked words to pretrain a language model.

Self attention and BERT are both examples of complete graphs, but on sentences rather than image features and labels. C-Tran uses the same self-attention mechanisms as [7] and [2], but instead of using only the word embeddings from a sentence, we use feature and label embeddings.

In computer vision, [1] used Transformers for object detection. Our method varies in several distinct ways. First, we are primarily interested in using partial evidence for image classification, and our unique state embeddings allow C-Tran to use such evidence. Second, we model image and label features jointly in a Transformer encoder, whereas [1] use an encoder/decoder framework. Our method allows the image features to be updated conditioned on the labels, which is a key characteristic of our model.

Connecting to Graph Based Neural Relational Learning. Another line of recent works employ object localization techniques [11, 9] or attention mechanism [8, 12] to locate semantic meaningful regions and try to identify underlying relations between regions and outputs. However, these

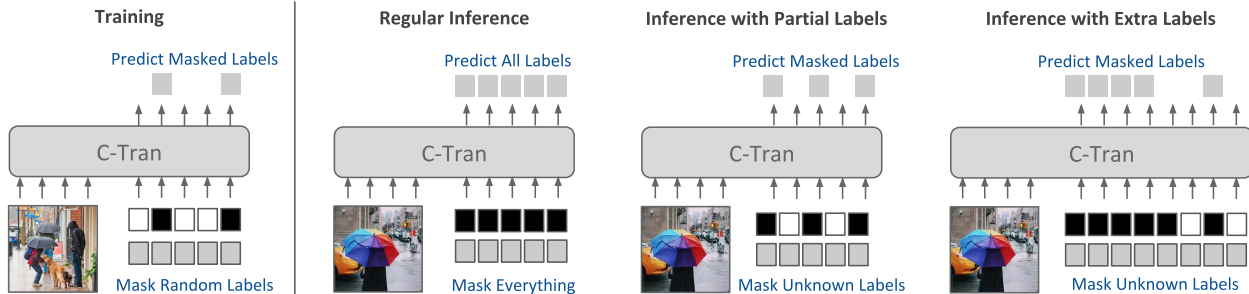


Figure 1. Detailed example of the general training method and three different inference settings where C-Tran can be applied.

Images	True Labels	ResNet-101	C-Tran	C-Tran + partial labels
 ID:000000362831	fork knife, spoon, bowl, chair, diningtable	fork, sandwich, diningtable, spoon, cup	fork, knife, diningtable, person, cake	spoon=1, trafficlight=0, bench=0, dog=0, ...
 ID:000000106216	person, car, truck, parkingmeter, horse,	person, car, truck, horse, bicycle	car, person, truck, horse, bicycle	bicycle=0, motorcycle=0, train=0, boat=0, ...
 ID:000000243213	person, bench, backpack, tennisracket, bottle, chair	person, tennisracket, chair, tie, sportsball	person, tennisracket, chair, sportsball, bench	backpack=1 parkingmeter=0, bird=0, zebra=0, ...
 ID:000000170129	airplane, train	airplane, boat, car, truck, person	airplane, boat, person, car, bird	car=0, motorcycle=0, bus=0, truck=0, ...
 ID: 000000262896	bottle, spoon, diningtable, cellphone, book	bottle, fork, diningtable, bowl, spoon	fork, spoon, bowl, book, diningtable	diningtable=1, bicycle=0, car=0, truck=0, ...

Figure 2. Qualitative examples of C-Tran + partial labels on the COCO-80 dataset. In the last column, we use = 25% partial labels, some of which are shown. Correctly predicted labels are in bold.

methods either require expensive bounding box annotations or merely get regions of interest roughly due to the lack of label supervision. One recent study by [10] also showed that modeling the associations between image feature regions and labels helps to improve multi-label performance.

In our work, C-Tran uses graph attentions and enables each target label to attend differentially to relevant parts of an input image.

For multi-label classification(MLC), [3] formulate MLC using a label graph and they introduced a conditional de-




Images	True Label	C-Tran	C-Tran + Extra Labels
 Anna_Hummingbird_0080_56366	Anna Hummingbird	Rufous Hummingbird (96%)	has_bill_shape_needle = 1, has_wing_color_green=1, has_upperparts_color=green=1, has_back_color_blue=0, has_back_color_brown=0 ... Anna Hummingbird (99%)
 Blue_Jay_0072_62944	Blue Jay	Florida Jay (99%)	has_bill_shape_all-purpose=1, has_upperparts_color_buff=1, has_upper_tail_color_grey=1, has_belly_color_red=0, has_wing_shape_broad-wings=0 ... Blue Jay (99%)
	Blue Winged Warbler	Yellow Headed Blackbird (99%)	has_upperparts_color_grey=1, has_tail_shape_rounded_tail=1, has_upper_tail_color_black=1, has_back_color_iridescent=0, has_underparts_color_purple=0 ... Blue Winged Warbler (99%)

Figure 3. Qualitative examples of C-Tran + extra labels on the CUB-312 dataset. In the last column, we use = 54% extra labels, some of which are shown.

pendency SVM where they first trained separate classifiers for each label given the input and all other true labels and used Gibbs sampling to find the optimal label set. The main drawback is that this method requires separate classifiers for each label. [6] proposes a method to label the pairwise edges of randomly generated label graphs, and requires some chosen aggregation method over all random graphs. The authors introduce the idea that variation in the graph structure shifts the inductive bias of the base learners. One recent study [5] used graph neural networks for multi-label classification on sequential inputs. The proposed method models the label-to-label dependencies using GNNs, however, does not represent input features and labels in one coherent graph. A key aspect of C-Tran is that the Transformer encoder can be viewed as a fully connected graph which is able to learn any relationships between features and labels. The Transformer attention mechanism can be regarded as a form of graph ensemble learning [4]. Above all, previous methods using graphs to model label dependencies do not allow for partial evidence information to be included in the prediction.

A.5. Label Mask Training

In Algorithm 1, we detail the label mask training (LMT) procedure. For each training sample, we select a random amount of labels to be used as “known” input labels to the model. The loss function is then computed on all unknown labels.

Algorithm 1: C-Tran Label Mask Training Procedure

```

1 loss = 0
2 for sample (x; y) in batch do
3   label_idx = range(1, `);
4   n = randint(0.25`, `);
5   unk_idx = sample(label_idx, n);
6   y_u = {y_i for i in unk_idx};
7   y_k = {y_j for j in label_idx excluding unk_idx};
8   ŷ_u = f(x; y_k);
9   for label index j in y_u do
10    | loss += -(y_j log(ŷ_j) + (1 - y_j) log(1 - ŷ_j))
11  end
12 end
13 Run backprop;
14 Update ;

```
