

6. Appendix

6.1. Architecture details

As a further explanation of the model architecture, we plot the conceptual attention module in Figure 5.

6.2. Experiments

All experiments are implemented by Python 3.8 and Pytorch 1.5.1. We run each of the experiments on an NVIDIA GeForce RTX 2080Ti GPU. Primarily we use the default Adam optimizer to train our model. The other training details are revealed in Table 5.

6.3. Comparison experiments

In Table 6 and 7, we reveal the settings for the baseline experiments. For the ablation study in Section 4.5, we use the same settings as in the previous section while altering the considered factors.

6.4. Dataset statistics extended

As extended from Table 1, in this section we provide some further details about the datasets and their label graphs.

1. **Pet datasets.** We construct 14 augmented nodes in addition to the 39 label nodes extracted from both datasets. The total edge number is 119.
2. **Flower datasets.** We augment the label graph by 18 intermediate nodes with a total edge number 129.
3. **Arxiv.** We adopt the original Arxiv label hierarchy graph structuring the website⁴.

Due to the size issues of these graphs, we only visualize the moderate-size graph of the Pet datasets in Figure 6. We will release the other label graphs in the format of metafiles upon publication.

⁴<http://arxiv.org/>

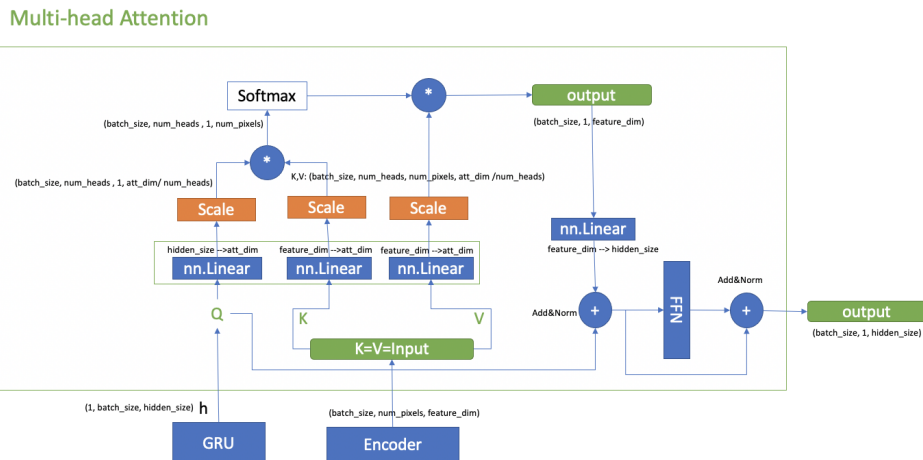


Figure 5. Multi-head attention module architecture.

Table 5. Training settings for the classification experiments of our label space augmentation approach (*label-aug*). *LRE* and *LR* indicate the learning rate adopted for the encoder backbone (*EfficientNet-B4* and *BERT/LSTM* respectively for image and text experiments) and the rest of model parameters. *Dynamic-Reduce(n)* means we lower the learning rate by half when the dev set performance does not improve for a consecutive n epochs.

DATASET	Batch Size	Initial LRE	Initial LR	Decay Period(epoch)	Attention Size	Hidden Size	Input
<i>Oxford-IIIT Pet</i>	8	0.00005	0.0004	10	1024	1024	380*380
<i>PetFusion</i>	24	0.0001	0.0004	10	1024	1024	380*380
<i>102 Category Flower</i>	12	0.0001	0.0005	<i>Dynamic Reduce(5)</i>	1024	1024	380*380
<i>FlowerFusion</i>	12	0.0002	0.0006	<i>Dynamic Reduce(5)</i>	1024	1024	380*380
<i>ArxivOriginal</i>	128	0.000015	0.000015	<i>Dynamic Reduce(5)</i>	768	768	≤ 512
<i>ArxivFusion</i>	128	0.000015	0.000015	<i>Dynamic Reduce(5)</i>	768	768	≤ 512
<i>ArxivOriginal-LSTM</i>	128	0.0007	0.0007	<i>Dynamic Reduce(5)</i>	768	768	≤ 512

Table 6. Training settings for baseline experiments. Namely the *EfficientNet-B4 + FFN* and *BERT/LSTM + FFN* experiments for image and text domain respectively.

DATASET	Epoch	Batch Size	Input	Initial LR	Decay Period(epoch)
<i>Oxford-IIIT Pet</i>	60	16	380*380	0.00005	10
<i>102 Category Flower</i>	60	16	380*380	0.0001	<i>Dynamic Reduce(5)</i>
<i>ArxivOriginal</i>	40	128	≤ 512	0.000015	<i>Dynamic Reduce(5)</i>
<i>ArxivFusion</i>	40	128	≤ 512	0.000015	<i>Dynamic Reduce(5)</i>
<i>ArxivOriginal-LSTM</i>	80	128	≤ 512	0.0007	<i>Dynamic Reduce(5)</i>

Table 7. Training settings for *Pseudo Label* method.

DATASET	Epoch	Batch Size	Input	Initial LR	Decay Period(epoch)
<i>Oxford-IIIT Pet</i>	60	16	380*380	0.0001	10
<i>PetFusion</i>	60	16	380*380	0.0001	10
<i>102 Category Flower</i>	60	16	380*380	0.0004	<i>Dynamic Reduce(5)</i>
<i>FlowerFusion</i>	60	16	380*380	0.0004	<i>Dynamic Reduce(5)</i>

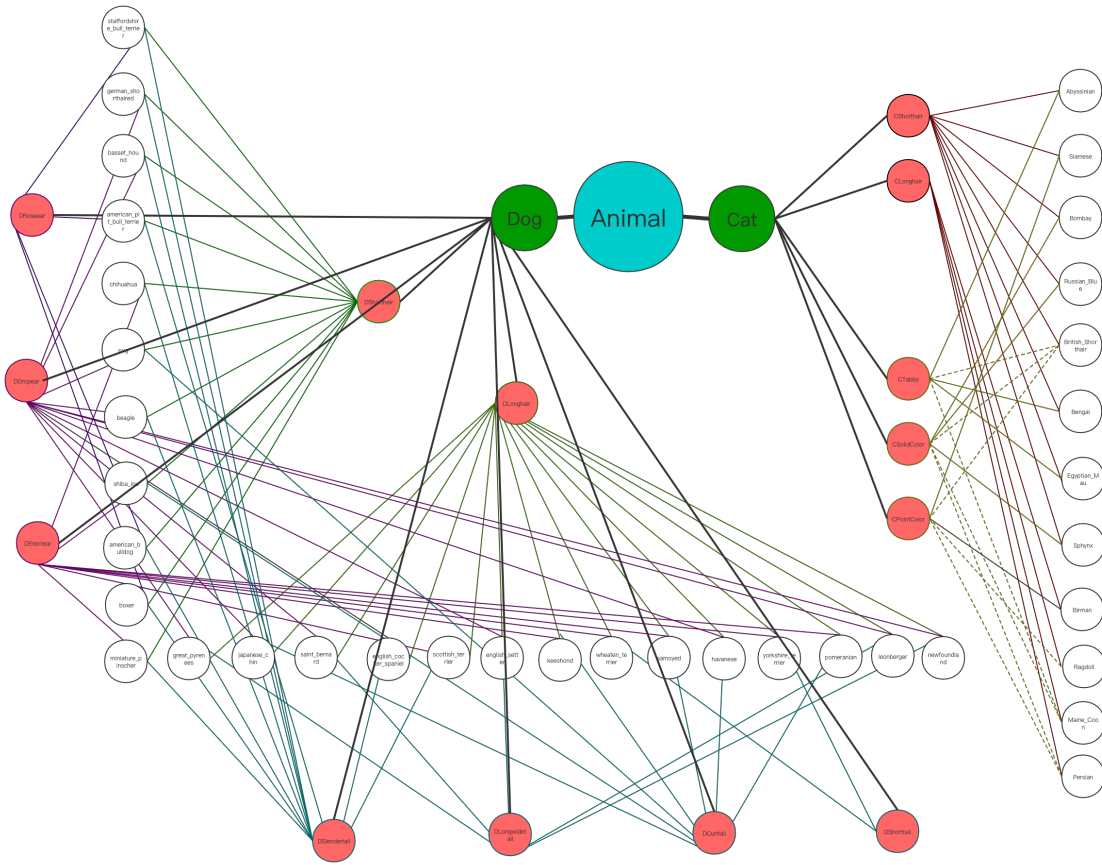


Figure 6. Complete label graph for *Pet datasets*. The red and green color denote augmented nodes at different hierarchies, and white color indicates label nodes. The dotted lines indicate nondeterministic paths.