

# Hierarchical Classifier

*intern@mmlab, Sep.13*

---

## Abstract

**1).** Explore Dataset **2).** Ablation study on Baseline model. **3).** A unified method combining coarse-grained and fine-grained classification task.

---

## Contents

<b>1</b>	<b>Related Work</b>	<b>1</b>
1.1	Hierarchical Softmax Classifier ( <a href="#">Mikolov et al., 2013</a> ) . . . . .	1
1.2	Multiple Granularity Descriptors ( <a href="#">Wang et al., 2015</a> ) . . . . .	2
1.3	ImageNet ( <a href="#">Deng et al., 2010</a> ) . . . . .	2
1.4	Structure prediction . . . . .	3
<b>2</b>	<b>ImageNet</b>	<b>3</b>
2.1	Observation: . . . . .	3
<b>3</b>	<b>Training on ImageNet10k</b>	<b>5</b>
<b>4</b>	<b>Train on Cifar100</b>	<b>6</b>

## 1. Related Work

### 1.1. Hierarchical Softmax Classifier ([Mikolov et al., 2013](#))

The large number of classes of ImageNet, *i.e.* 22K categories, can be expensive to train, if use basic full softmax classifier.  $w_O$  need to travel through all possible classes, *e.g.*  $10^7$  words in dictionary for NLP task.

$$p(w_O|w_I) = \frac{\exp(v'_{w_O} \top v_{w_I})}{\sum_{w=1}^W \exp(v'_w \top v_{w_I})}$$

The idea of hierarchical softmax is **1).** view 20K classes as a tree generated by Huffman Coding (If just want to reduce training cost. The cost is  $O(L(w_O))$ ). **2).** train a multi-label classifier at each node on the tree.

$$p(w|w_I) = \prod_{j=1}^{L(w)-1} \sigma \left( [n(w, j+1) = \text{ch}(n(w, j))] \cdot v'_{n(w, j)}{}^\top v_{w_I} \right)$$

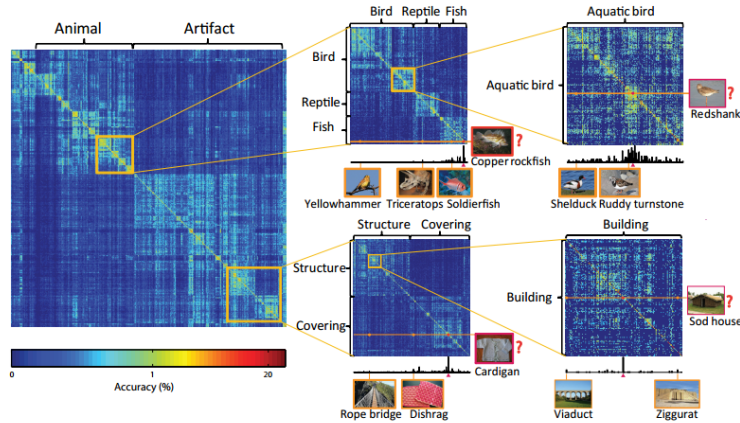
But it is different with our work:

- The learned representation is flatten. The dilemma between different granularity still exists.
- If trained on ImageNet 21K (A more practical case), we can use *Flattened Hit Rate@Top k* or other metric to evaluate it. Given a new species ‘Leopard-tiger’ it can only say ‘It is animal  $\succ$  leopard  $\succ$  tiger’, but it can not say ‘It is animal, to be specific, may be leopard or tiger’.

### 1.2. Multiple Granularity Descriptors ([Wang et al., 2015](#))

### 1.3. ImageNet ([Deng et al., 2010](#))

- Flattened Classifier: Train on ImageNet7K, equivalent to our flattened baseline. The hierarchical structure of classes is learned.



- Measure Grainuity of Dataset:  $h(i, j)$  between categories  $i$  and  $j$ , as the height of their lowest common ancestor. Grainuity =  $\text{mean}_{i,j \in c} \{h(i, j)\}$ .
- Loss for Hierarchical Dataset: ‘redshank’ is a ‘bird’, classiy it as ‘bird’ is less informative, but at least right. Rather than using 0-1 Loss, (Deng et al., 2010) use Hierarchical Loss.  $L(x) = \sum_{j=1}^K C_{i,j} p_j(x)$ , where  $C_{i,j} = h(i, j)$ .

#### 1.4. Structure prediction

## 2. ImageNet

- ImageNet22K: Whole Dataset. 22K cateories from the Fall 2011 release of ImageNet, only including leaf categories.
- ImageNet10K(Deng et al., 2010): 10184 categories from the Fall 2009 release of ImageNet, including both internal and leaf nodes with more than 200 images in each category.
- ImageNet7K: 7404 leaf categories from ImageNet10K. Only leaf.
- ImageNet1K(ILSCV 2010 – 2012 Classification Task): 1000 leaf categories randomly sampled from ImageNet7K.

#### 2.1. Observation:

- Approximately Tree Structure

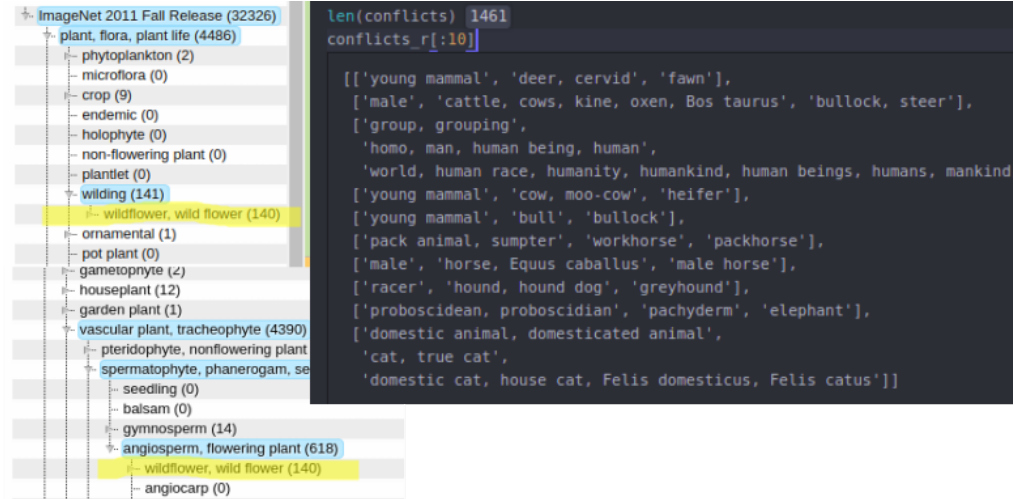


Figure 1: A problem in ImageNet: It is true that fawn is both yong manmal and deer,thus wordnet is a DAG. But the conflict is few only 1461, compared to 82K<sup>2</sup>, so can be approximately viewed as a tree.

We will use DFS to force it to become a tree.

- Analyze Hierarchical Structure of Classes

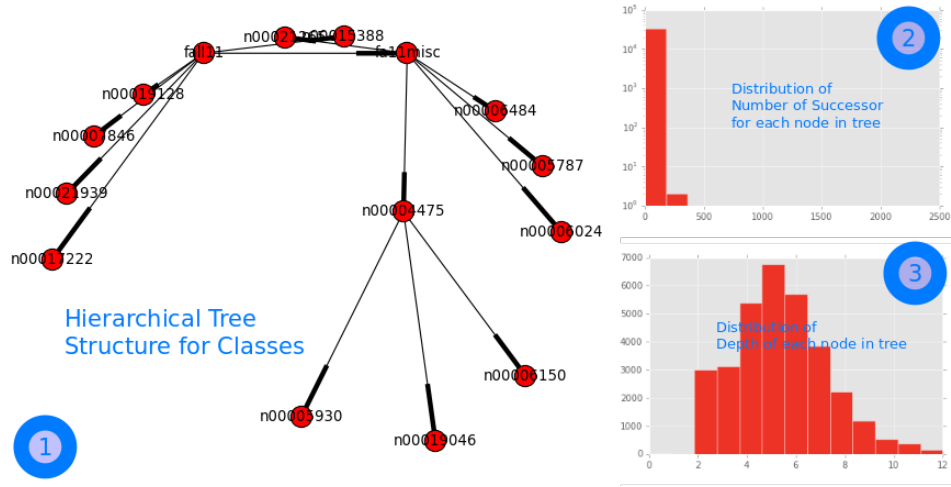


Figure 2: Note that in 2)., the Y-axis is log-scaled

- Log-tail Distribution of Number of Instance in each Class

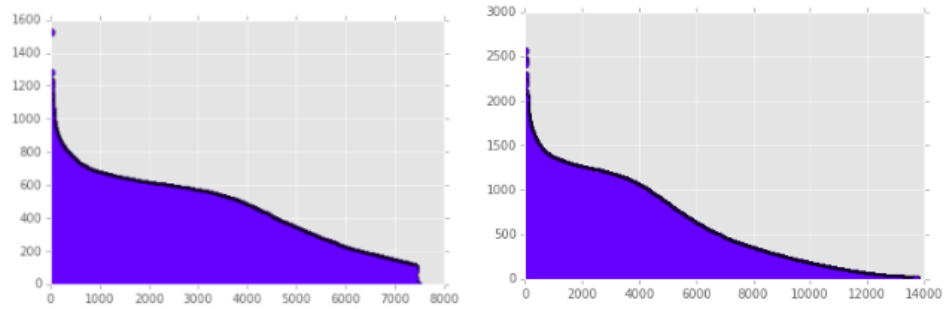


Figure 3: **Left:** Number of instance in each class for classes in leaves of ImageNet10K (At least 200 images per classes), **Right:** For classes in ImageNet22K (Partially downloaded). update.

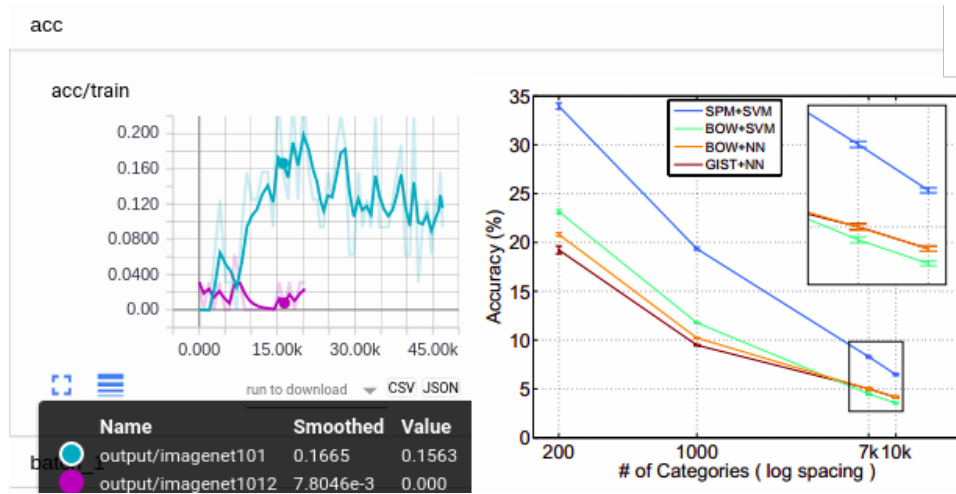
<sup>2</sup>82K contains all classes from coarse to fine, 22K-ImageNet refer to the most fine-grained classes (leaf node) is 22K.

- We will use ImageNet10K for our experiments
  - Categories are over-lapping between internal-node and leaves.
  - Contains ImageNet1K, where pretrained model trained from. But we may ignore this problem.
  - The benefit is number of images in each class is more than 200, although training data still unbalanced.
- No standard train/valid/test split for the task
  - Some people use 50/50, which is a bit crazy, but it is the convention of the PASCAL VOC Challenge, and (Deng et al., 2010) use this split protocol.
  - Some use 90/10, which can be quite different. It is split protocol for ILSVR Competition.

### 3. Training on ImageNet10k

After removing internal node, ImageNet10k contains 7461 categories. Thus in fact we are using ImageNet7k proposed by (Deng et al., 2010) but we can still give prediction for all classes in ontology tree for ImageNet10k. Train baseline on this dataset:

- Resnet101  $\rightarrow$  7x7x1024  $\rightarrow$  global ave pooling  $\rightarrow$  1024  $\rightarrow$  FC  $\rightarrow$  7461
  - Step 020660 - Loss 9.43 - acc 0.00% (0.37sec/step)
- Resnet101  $\rightarrow$  7x7x1024  $\rightarrow$  conv  $\rightarrow$  2x2x1024  $\rightarrow$  Flatten + FC  $\rightarrow$  7461
  - Step 042420 - Loss: 5.86 acc: 34.38% (0.870 sec/step)
  - Step 045160 - Loss: 5.48 acc: 3.12% (0.697 sec/step)
- Training expense still comes from I/O, rather than normalization computation in softmax layer.
- Due to imbalance image number/classes in one batch, training accuracy is fluctuating. The test accuracy(average on each instance) reported in (Deng et al., 2010) is 7%. Our train averaging accuracy is 16%.



**Left:** Learning curve of our model, **Right:** Accuracy reported in (Deng et al., 2010)

#### 4. Train on Cifar100

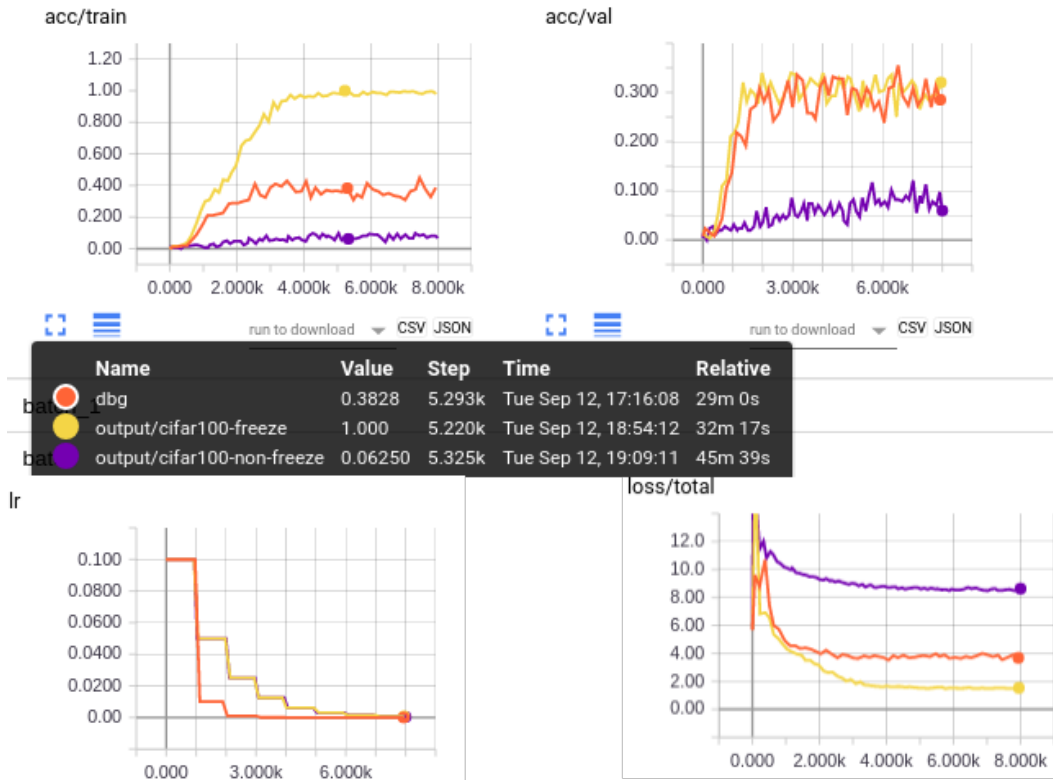
The Structure of Cifar100(20 super classes, 100 fine classes) is similar to our setting, except that instance in each class is balance(100 imgs/class) and batch size  $256 > 100$ .

Note that for multiple loss experiments, we will reorder the label for each class, the class label for semantically similar class will be similar, *e.g.* mapping 4, 30, 55, 72, 95 to 0, 1, 2, 3, 4. Using this class label, mapping from coarse label to fine label can be calculated directly rather than using HashTable.

```
In [17]: cifar100.mapping_human
Out[17]: {'aquatic_mammals': {'beaver', 'dolphin', 'otter', 'seal', 'whale'},
          'fish': {'aquarium_fish', 'flatfish', 'ray', 'shark', 'trout'},
          'flowers': {'orchid', 'poppy', 'rose', 'sunflower', 'tulip'},
          'food_containers': {'bottle', 'bowl', 'can', 'cup', 'plate'},
In [18]: cifar100.mapping
Out[18]: {0: {4, 30, 55, 72, 95},
          1: {1, 32, 67, 73, 91},
          2: {54, 62, 70, 82, 92},
          3: {9, 10, 16, 28, 61},
```

- How to train fast:

- Finetune = Initialize from pretrained model + Only train last parameter block
- But observe overfitting if learning rate is large for a long time. High val acc on cifar100 can be 75.72%, but we will first use this model to iterate ideas.



- Combine multiple loss

Given the feature of a sample  $\mathbf{x}_i$ , whose ground truth is  $\mathbf{y}_i = c_i$ , *i.e.*  $y_i^{c_i}$ , we want to know the propabability of each possible category  $y_i$  to form a prediction vector for this sample  $\mathbf{P}(\mathbf{y}_i|\mathbf{x}_i) = [P(y_i^{c_0}|\mathbf{x}_i), \dots, P(y_i^{c_n}|\mathbf{x}_i)]^T$ .

If we have semantic prior from wordnet tree, then

$$\begin{aligned}
P(y_i^{c_i} | \mathbf{x}_i) &= \sum_{c_k \in \text{Par}(c_i)} P(y_i^{c_i} | y_i^{c_k}, \mathbf{x}_i) P(y_i^{c_k} | \mathbf{x}_i) \\
&= P(y_i^{c_i} | y_i^{\text{Par}(c_i)}, \mathbf{x}_i) P(y_i^{\text{Par}(c_i)} | \mathbf{x}_i)
\end{aligned}$$

*E.g.*, let  $\mathbf{y}_i = c_i = \text{tiger}$ ,  $\text{Par}(c_i) = \text{carnivore}$  and  $\text{Par}(c_j) = \text{herbivore}$ , then  $P(y_i^{c_i} | y_i^{\text{Par}(c_j)}, \mathbf{x}_i) = 0$ , since we infer  $y_i^{c_i} \perp y_i^{\text{Par}(c_j)} | \mathbf{x}_i$  from wordnet tree.

Hierachy softmax force wordnet tree into *binary* tree using hierachuy clustering (Morin and Bengio, 2005), due to computation cost consideration. We aim to fully utilize the semantic prior: **1).** use geometry structure between classes to boost performance, **2).** use taxonomy tree to handle unbalanced training examples.

Cross entroy loss is widely used in multi-class classification. Let  $\mathbf{f} = [f_0, \dots, f_n]^T$  be the logits, then cross entroy loss is defined as  $L_i = -\log(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}})$ . Based on it, We design a hierachical loss incorprating semantic prior from wordnet.

To begin with a simple case, we consider cifar100, Ref. to Fig. 4:

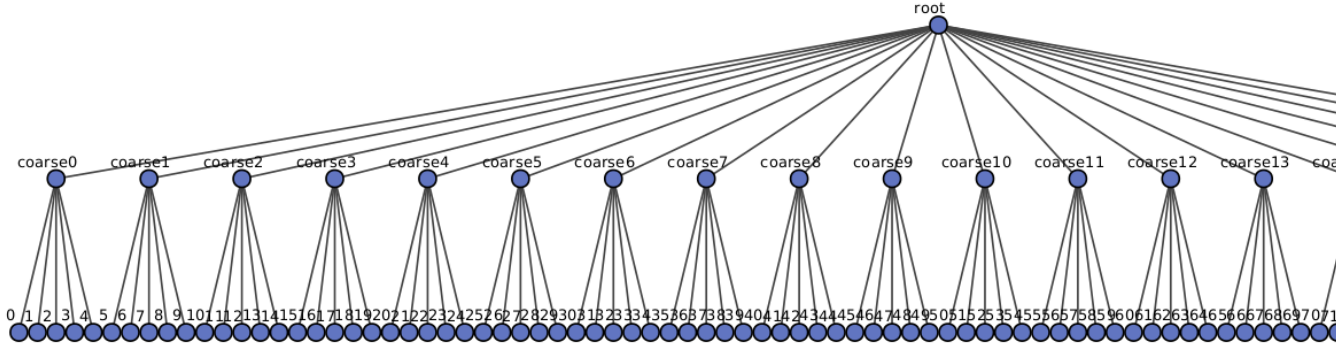


Figure 4: Coarse and Fine Categories of Cifar100

Assume we have Function  $\text{PAR}(\text{node}, \text{depth})$ , Ref. to Alg. 1 and  $\text{CH}(\text{node}, \text{depth})$ , Ref. to Alg. 2 and we define  $\text{FIND} = \text{lambda } y_i, \text{start}, \text{end} : \text{CH}(\text{PAR}(y_i, \text{start}), \text{end})$ .

$$\begin{aligned}
L_i^{100} &= -\log \left( \frac{\sum_{j \in \text{FIND}(y_i, 2, 2)} e^{f_j}}{\sum_{j \in \text{FIND}(y_i, 0, 2)} e^{f_j}} \right) \\
L_i^{20} &= -\log \left( \frac{\sum_{j \in \text{FIND}(y_i, 1, 2)} e^{f_j}}{\sum_{j \in \text{FIND}(y_i, 0, 2)} e^{f_j}} \right) \\
L_i^{\text{group}} &= -\log \left( \frac{\sum_{j \in \text{FIND}(y_i, 2, 2)} e^{f_j}}{\sum_{j \in \text{FIND}(y_i, 1, 2)} e^{f_j}} \right)
\end{aligned}$$



---

**Algorithm 1** Find Parents of *node* at *depth*

---

```
1: function PAR(node,depth)
2: % Input: node: the ground truth label to query;
3: %      depth: parent node's depth.
4: % Output: parent node at depth
5:   if node.depth  $\geq$  depth then
6:     return node
7:   else
8:     return PAR(node.parent,depth)
```

---

---

**Algorithm 2** Find Child Nodes of *root* at *depth*

---

```
1: function CH(root,depth)
2: % Input: root: the root node to query;
3: %      depth: child nodes' depth.
4: % Output: list of child nodes at depth
5:   ChildNodes  $\leftarrow$  list()
6:   for node in root.childs do
7:     if node.depth  $\geq$  depth then
8:       ChildNodes.extend(list([node]))
9:     else
10:      ChildNodes.extend(CH(node, depth))
11:   return ChildNodes
```

---

Deng, J., Berg, A. C., Li, K., Fei-Fei, L., 2010. What does classifying more than 10,000 image categories tell us? In: European conference on computer vision. Springer, pp. 71–84.

Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., Dean, J., 2013. Distributed representations of words and phrases and their compositionality. In: Advances in neural information processing systems. pp. 3111–3119.

Morin, F., Bengio, Y., 2005. Hierarchical probabilistic neural network language model. In: Aistats. Vol. 5. pp. 246–252.

Wang, D., Shen, Z., Shao, J., Zhang, W., Xue, X., Zhang, Z., 2015. Multiple granularity descriptors for fine-grained categorization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2399–2406.