



THE OHIO STATE UNIVERSITY

# Hierarchical Classification with Confidence using Generalized Logits

Jim Davis<sup>1</sup>, **Tong Liang<sup>1</sup>**, James Enouen<sup>1</sup>, Roman Ilin<sup>2</sup>

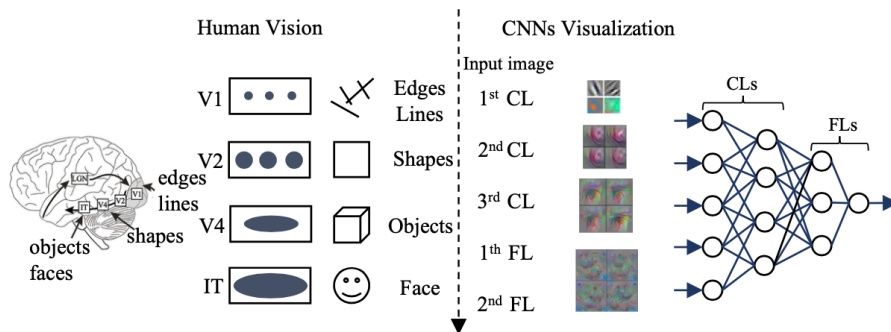
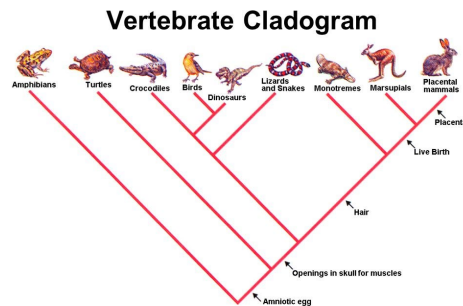
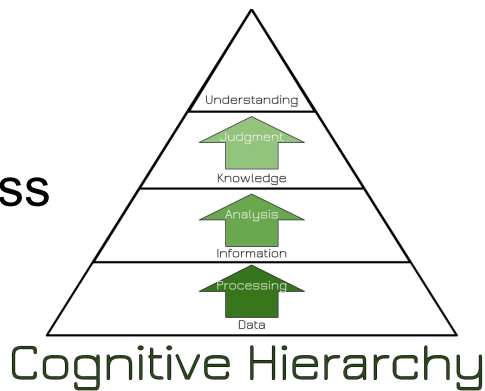
<sup>1</sup>Ohio State University, Columbus OH 43210

<sup>2</sup>AFRL/RYP, Wright-Patterson AFB OH 45433

This work is supported by the U.S. Air Force Research Laboratory

# Hierarchical Reasonings

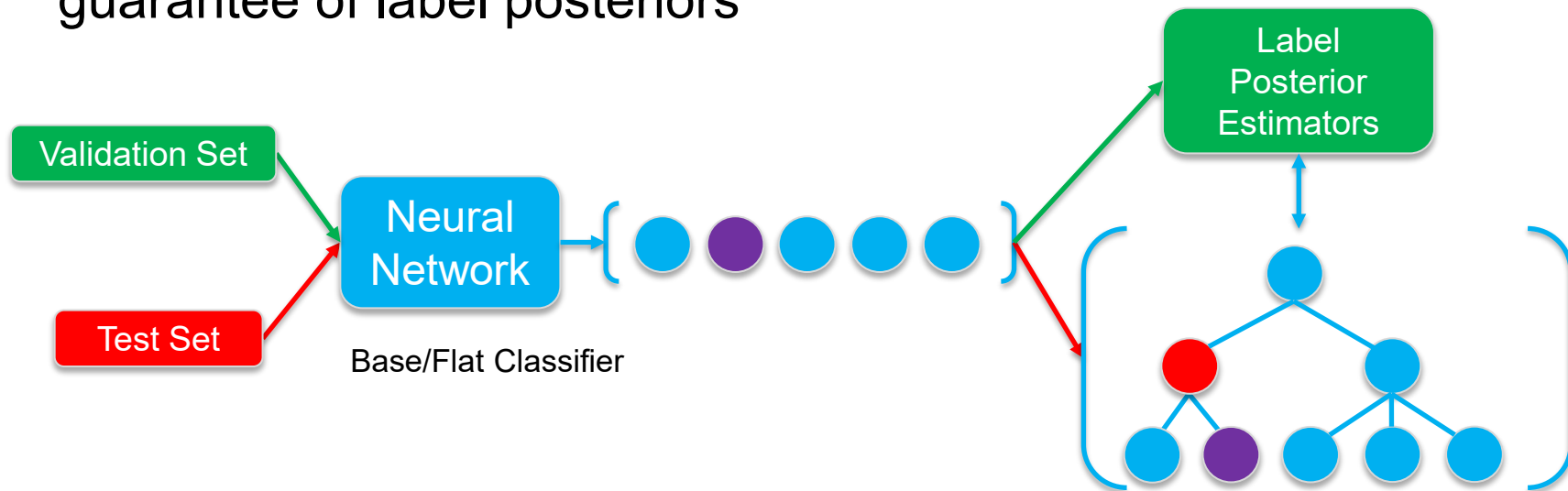
- Hierarchical reasoning prevails in many areas
  - Hierarchical cognitive process
  - Taxonomy in biology
  - Natural hierarchical representation of visual features in our brain





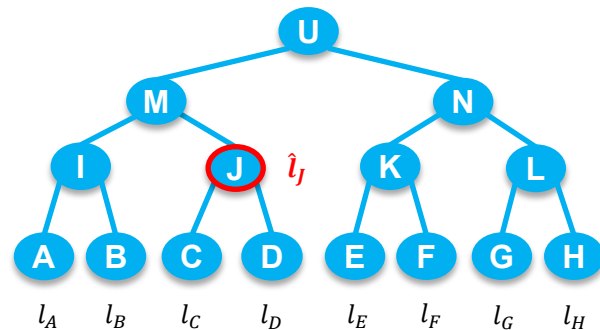
# Overview of Our Approach

- Serves as a post processing step after a given base model
- Estimates label posteriors from base model logits output of the validation set
- Produces hierarchical predictions to test set with statistical guarantee of label posteriors



# Generalized Logits

- Generalized logit = derived logit for a non-terminal class
- Computed from base classifier's softmax value of  $s_J$

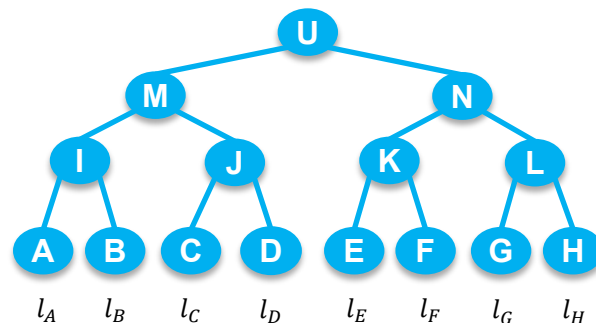


$$s_J = s_C + s_D = \frac{e^{l_C}}{e^{l_C} + \sum_{k \neq C} e^{l_k}} + \frac{e^{l_D}}{e^{l_D} + \sum_{k \neq D} e^{l_k}} = \frac{\sum_{i \in \{C, D\}} e^{l_i}}{\sum_{i \in \{C, D\}} e^{l_i} + \sum_{k \notin \{C, D\}} e^{l_k}} \triangleq \frac{e^{\hat{l}_J}}{e^{\hat{l}_J} + \sum_{k \notin J} e^{l_k}}$$

$\hat{l}_J$  is the **generalized logit** of class/node J

$$\hat{l}_J = \ln(e^{\hat{l}_J}) = \ln\left(\sum_{i \in \{C, D\}} e^{l_i}\right)$$

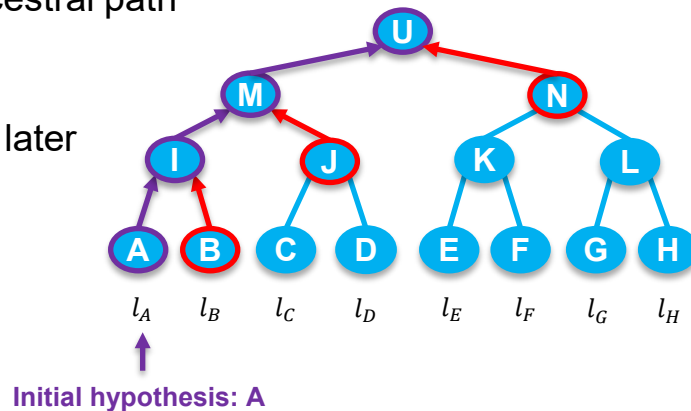
- Start with initial terminal label hypothesis
  - Option 1: Select via argmax of base classifier's logits
  - Option 2: Select via argmax of estimated terminal label posteriors
- Examine confidence of initial terminal label hypothesis
  - Is it above a given confidence threshold?
  - If YES, return that label
  - If NO, examine the remaining ancestral classes until meeting the threshold (customized)
    - Root node of the hierarchy has posterior of 1.0



# Inference: An Example

- Consider the binary tree shown on the right, evaluate the ancestral path label posteriors:

$P(A|\mathcal{L}_A)$   $\mathcal{L}_A$  is the generalized logit vector to be introduced later

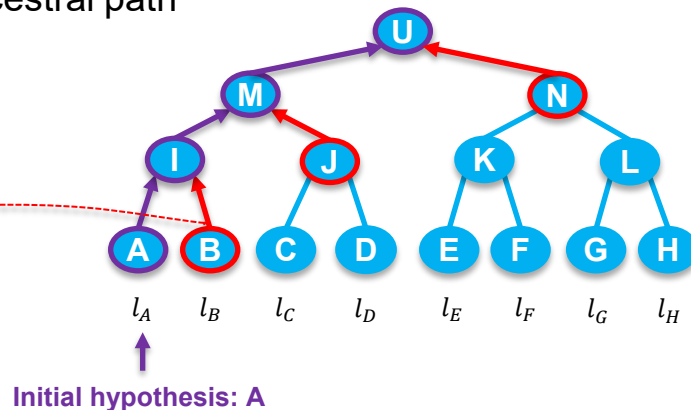


# Inference: An Example

- Consider the binary tree shown on the right, evaluate the ancestral path label posteriors:

$$P(A|\mathcal{L}_A)$$

$$P(I|\mathcal{L}_A) = P(A \cup \mathbf{B}|\mathcal{L}_A) = P(A|\mathcal{L}_A) + \mathbf{P}(\mathbf{B}|\mathcal{L}_A)$$



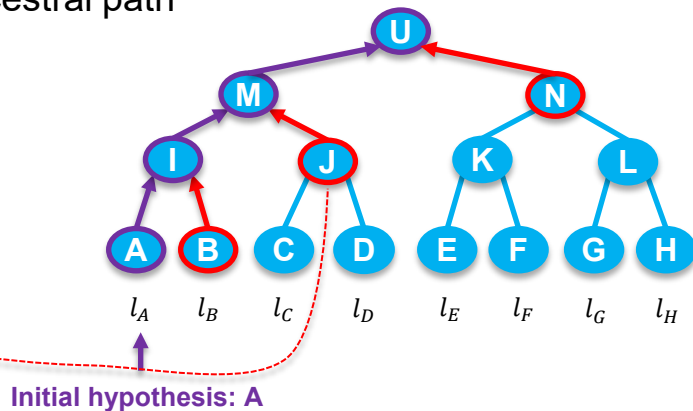
# Inference: An Example

- Consider the binary tree shown on the right, evaluate the ancestral path label posteriors:

$$P(A|\mathcal{L}_A)$$

$$P(I|\mathcal{L}_A) = P(A \cup \mathbf{B}|\mathcal{L}_A) = \mathbf{P}(A|\mathcal{L}_A) + \mathbf{P}(\mathbf{B}|\mathcal{L}_A)$$

$$P(M|\mathcal{L}_A) = P(I \cup \mathbf{J}|\mathcal{L}_A) = \mathbf{P}(I|\mathcal{L}_A) + \mathbf{P}(\mathbf{C} \cup \mathbf{D}|\mathcal{L}_A)$$



# Inference: An Example

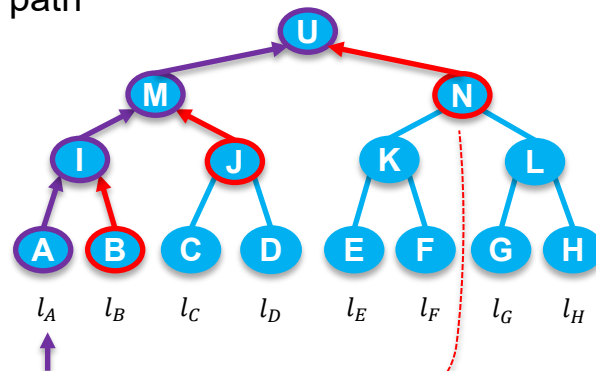
- Consider the binary tree shown on the right, evaluate the ancestral path label posteriors:

$$P(A|\mathcal{L}_A)$$

$$P(I|\mathcal{L}_A) = P(A \cup \mathbf{B}|\mathcal{L}_A) = \mathbf{P}(A|\mathcal{L}_A) + \mathbf{P}(\mathbf{B}|\mathcal{L}_A)$$

$$P(M|\mathcal{L}_A) = P(I \cup \mathbf{J}|\mathcal{L}_A) = \mathbf{P}(I|\mathcal{L}_A) + \mathbf{P}(\mathbf{C} \cup \mathbf{D}|\mathcal{L}_A)$$

$$P(U|\mathcal{L}_A) = P(M \cup \mathbf{N}|\mathcal{L}_A) = \mathbf{P}(M|\mathcal{L}_A) + \mathbf{P}(\mathbf{E} \cup \mathbf{F} \cup \mathbf{G} \cup \mathbf{H}|\mathcal{L}_A)$$



Initial hypothesis: A

# Inference: An Example

- Consider the binary tree shown on the right, evaluate the ancestral path label posteriors:

$$P(A|\mathcal{L}_A)$$

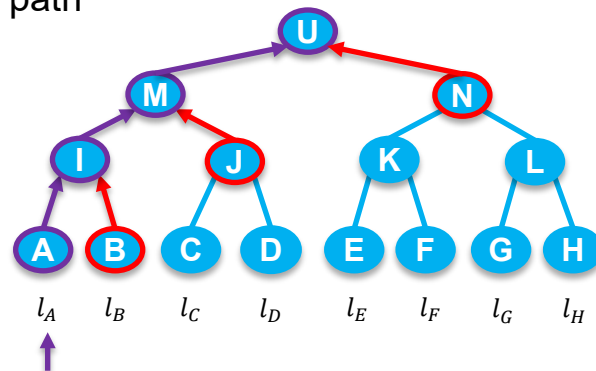
$$P(I|\mathcal{L}_A) = P(A \cup \mathbf{B}|\mathcal{L}_A) = \mathbf{P}(A|\mathcal{L}_A) + \mathbf{P}(\mathbf{B}|\mathcal{L}_A)$$

$$P(M|\mathcal{L}_A) = P(I \cup \mathbf{J}|\mathcal{L}_A) = \mathbf{P}(I|\mathcal{L}_A) + \mathbf{P}(\mathbf{C} \cup \mathbf{D}|\mathcal{L}_A)$$

$$P(U|\mathcal{L}_A) = P(M \cup \mathbf{N}|\mathcal{L}_A) = \mathbf{P}(M|\mathcal{L}_A) + \mathbf{P}(\mathbf{E} \cup \mathbf{F} \cup \mathbf{G} \cup \mathbf{H}|\mathcal{L}_A)$$

- L1 normalization before inference:

$$P(A|\mathcal{L}_A) + P(B|\mathcal{L}_A) + P(C \cup D|\mathcal{L}_A) + P(E \cup F \cup G \cup H|\mathcal{L}_A) \triangleq 1.0$$



# Inference: An Example

- Consider the binary tree shown on the right, evaluate the ancestral path label posteriors:

$$P(A|\mathcal{L}_A)$$

$$P(I|\mathcal{L}_A) = P(A \cup \mathbf{B}|\mathcal{L}_A) = P(A|\mathcal{L}_A) + \mathbf{P}(\mathbf{B}|\mathcal{L}_A)$$

$$P(M|\mathcal{L}_A) = P(I \cup \mathbf{J}|\mathcal{L}_A) = P(I|\mathcal{L}_A) + \mathbf{P}(\mathbf{C} \cup \mathbf{D}|\mathcal{L}_A)$$

$$P(U|\mathcal{L}_A) = P(M \cup \mathbf{N}|\mathcal{L}_A) = P(M|\mathcal{L}_A) + \mathbf{P}(\mathbf{E} \cup \mathbf{F} \cup \mathbf{G} \cup \mathbf{H}|\mathcal{L}_A)$$

- L1 normalization before inference:

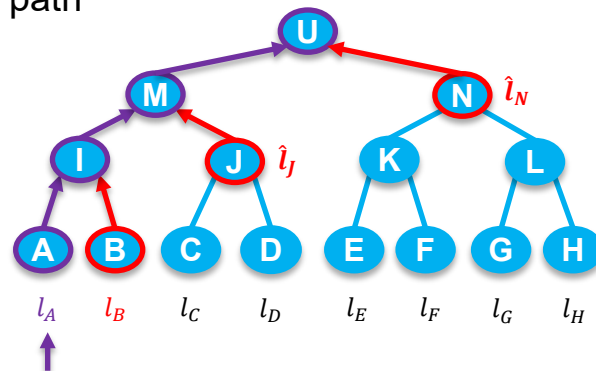
$$P(A|\mathcal{L}_A) + P(\mathbf{B}|\mathcal{L}_A) + P(\mathbf{C} \cup \mathbf{D}|\mathcal{L}_A) + P(\mathbf{E} \cup \mathbf{F} \cup \mathbf{G} \cup \mathbf{H}|\mathcal{L}_A) \triangleq 1.0$$

**A**

**B**

**J**

**N**





# Extension to Non-Binary Tree

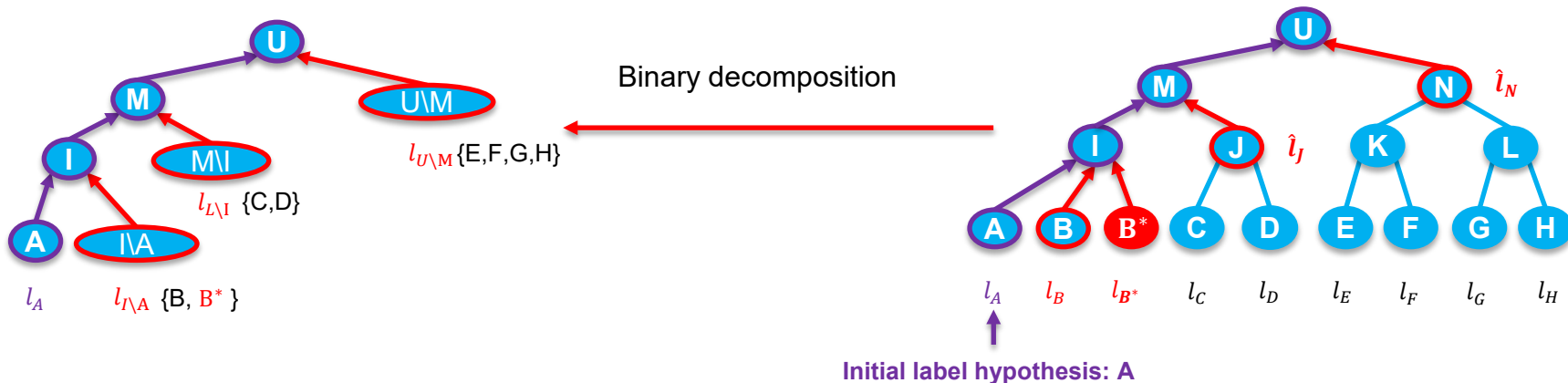
- Adding a node  $B^*$  under node  $I$

$$P(A|\mathcal{L}_A)$$

$$P(I|\mathcal{L}_A) = P(A|\mathcal{L}_A) + P(B \cup B^*|\mathcal{L}_A) = P(A|\mathcal{L}_A) + P(I \setminus A|\mathcal{L}_A)$$

$$P(M|\mathcal{L}_A) = P(I|\mathcal{L}_A) + P(C \cup D|\mathcal{L}_A) = P(I|\mathcal{L}_A) + P(M \setminus I|\mathcal{L}_A)$$

$$P(U|\mathcal{L}_A) = P(M|\mathcal{L}_A) + P(E \cup F \cup G \cup H|\mathcal{L}_A) = P(M|\mathcal{L}_A) + P(U \setminus M|\mathcal{L}_A) \triangleq 1.0$$





# Generalized Logit Feature Vector

- Adding a node  $B^*$  under node  $I$

$$P(A|\mathcal{L}_A)$$

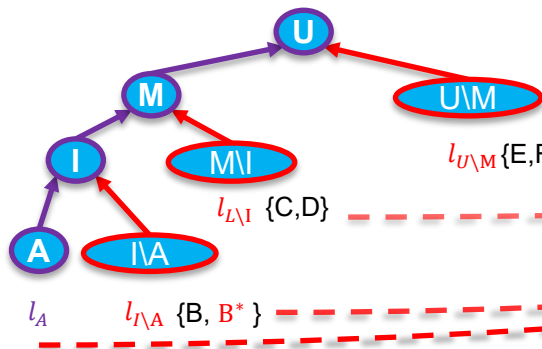
$$P(I|\mathcal{L}_A) = P(A|\mathcal{L}_A) + P(B \cup B^*|\mathcal{L}_A) = P(A|\mathcal{L}_A) + P(I \setminus A|\mathcal{L}_A)$$

$$P(M|\mathcal{L}_A) = P(I|\mathcal{L}_A) + P(C \cup D|\mathcal{L}_A) = P(I|\mathcal{L}_A) + P(M \setminus I|\mathcal{L}_A)$$

$$P(U|\mathcal{L}_A) = P(M|\mathcal{L}_A) + P(E \cup F \cup G \cup H|\mathcal{L}_A) = P(M|\mathcal{L}_A) + P(U \setminus M|\mathcal{L}_A) \triangleq 1.0$$

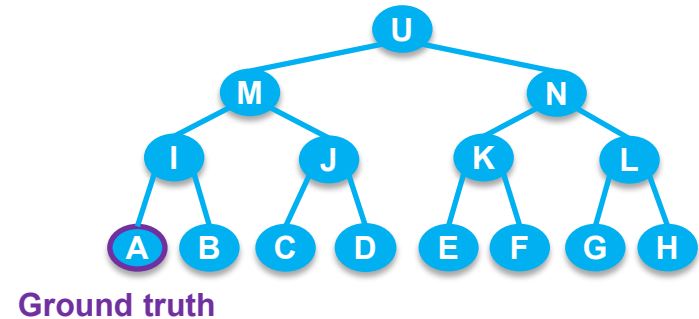
Generalized logit vector


$$\mathcal{L}_A = [l_A, \hat{l}_{I \setminus A}, \hat{l}_{M \setminus I}, \hat{l}_{U \setminus M}]$$



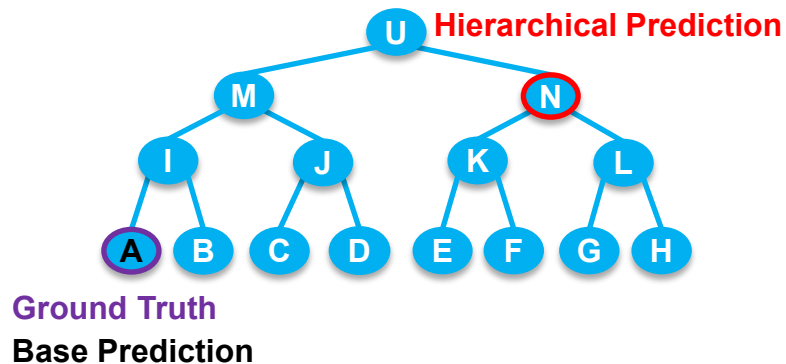
- We conducted experiments on 4 datasets with
  - ImageNet-Animal (398 terminal classes) [Davis et al., 2019]
  - CIFAR100 (100 terminal classes) [Krizhevsky, 2009]
  - CIFAR10 (10 terminal classes) [Krizhevsky, 2009]
  - Fashion-MNIST (10 terminal classes) [Xiao et al., 2017]
- Semantic hierarchy for each dataset is derived from WordNet [Davis et al., 2019]
- Compared our method with the two most related works
  - [Deng et al., 2012] employed optimization of tradeoff between accuracy and label specificity
  - [Davis et al., 2019] proposed a non-parametric histogram binning approach

- Hierarchical Classification Metrics based on originally correct (**C**) and originally incorrect (**IC**) predictions



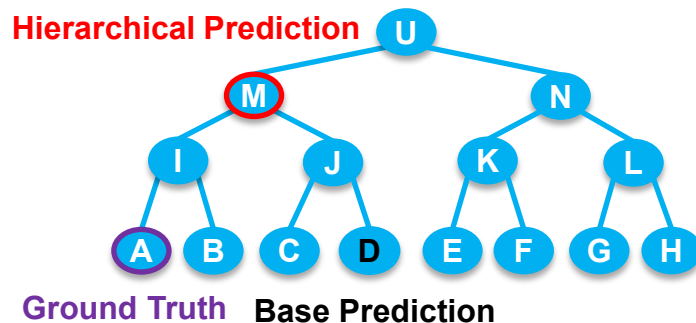
- Hierarchical Classification Metrics based on originally correct (**C**) and originally incorrect (**IC**) predictions
  -  – **C-Corrupt**: The fraction of original correct terminal predictions relabeled to incorrect labels

	Ground Truth	Base Prediction	Hierarchical Prediction
<b>C-Corrupt</b>	<b>A</b>	<b>A</b>	<b>N</b>



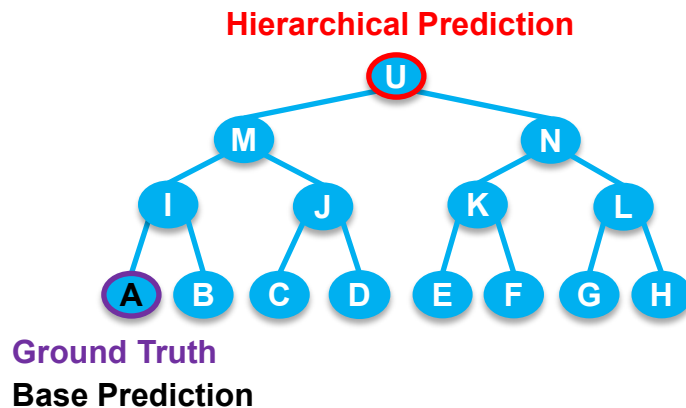
- Hierarchical Classification Metrics based on originally correct (**C**) and originally incorrect (**IC**) predictions
  - ▼ – **C-Corrupt**: The fraction of original correct terminal predictions relabeled to incorrect labels
  - ▲ – **IC-Reform**: The fraction of original incorrect terminal predictions generalized to correct labels





	Ground Truth	Base Prediction	Hierarchical Prediction
C-Corrupt	A	A	N
IC-Reform	A	D	M



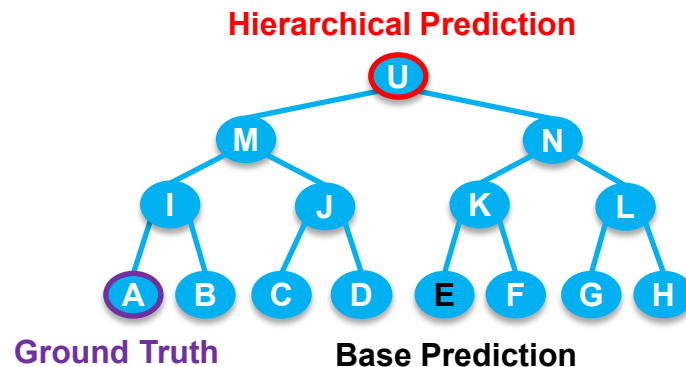
- Hierarchical Classification Metrics based on originally correct (**C**) and originally incorrect (**IC**) predictions
  - ▼ – C-Corrupt: The fraction of original correct terminal predictions relabeled to incorrect labels
  - ▲ – IC-Reform: The fraction of original incorrect terminal predictions generalized to correct labels
  - ▼ – **C-Withdrawn**: The fraction of original correct terminal predictions relabeled to root

	Ground Truth	Base Prediction	Hierarchical Prediction
C-Corrupt	A	A	N
IC-Reform	A	D	M
<b>C-Withdrawn</b>	<b>A</b>	<b>A</b>	<b>U</b>



- Hierarchical Classification Metrics based on originally correct (**C**) and originally incorrect (**IC**) predictions
  -  – C-Corrupt: The fraction of original correct terminal predictions relabeled to incorrect labels
  -  – IC-Reform: The fraction of original incorrect terminal predictions generalized to correct labels
  -  – C-Withdrawn: The fraction of original correct terminal predictions relabeled to root
  -  – **IC-Withdrawn**: The fraction of original incorrect terminal predictions relabeled to root

	Ground Truth	Base Prediction	Hierarchical Prediction
C-Corrupt	A	A	N
IC-Reform	A	D	M
C-Withdrawn	A	A	U
<b>IC-Withdrawn</b>	<b>A</b>	<b>E</b>	<b>U</b>



- Hierarchical Classification Metrics:

- ▲ – **Accuracy**: the fraction of correct hierarchical predictions (root is considered correct)
- ▲ – **Average scaled Information Gain (avg-sIG)**: corresponds to average depth of label generalizations in terms of Information Gain [Deng et al., 2012]

$$sIG(N_i) = \frac{\log_2 |T| - \log_2 (|\downarrow(N_i)|)}{\log_2 |T|}$$

$|T|$  is the total number of terminal classes

$|\downarrow(N_i)|$  is the number of terminal descendants of class  $N_i$

$$avg-sIG(N_i) = \frac{1}{M} \sum_{i=1}^M sIG(N_i)$$



# Experiments: ImageNet-Animal

- ImageNet-Animal derived from WordNet
  - Due to space limit, the lower part of the tree is omitted
  - # in (#) indicates the number of terminal classes at the branch

Unknown																													
Vertebrate																									Invertebrate				
Mammal															Reptile				Bird				Fish		...	Arthropod	...		
Placental															...	Diapsid			...	Aquatic		Oscine	...	Teleost	...	...	Insect	...	
Ungulate		Primate		Carnivore											...	Snake	Lizard	...	...	Wading	...	...	...	fish	...	...	...		
Even-toed ungulate	...	Monkey	...	Canine								Feline	...	...	...	...	...	...	bird	...	...	...	...	...	...	...			
				Dog						...	...	...	...	...															
				Hunting				Working		...																			
				Hound	Terrier	Sporting	...	Shepherd	...																				
...				...	...	...		...	...	...	...	...	...	...	...	...													
(15)	(2)	(13)	(7)	(19)	(26)	(17)	(1)	(12)	(18)	(25)	(12)	(13)	(15)	(17)	(6)	(17)	(11)	(3)	(5)	(16)	(8)	(11)	(24)	(10)	(6)	(8)	(27)	(20)	(14)



# Experiments: ImageNet-Animal




- Highest IC-Reform
- Highest overall accuracy

	Base	Proposed		Deng et al. 2012		Davis et al. 2019	
Confidence	-	90%	95%	90%	95%	90%	95%
C-Corrupt	-	.00	.00	.03	.01	.00	.00
▲ IC-Reform	-	<b>.96</b>	<b>.98</b>	.48	.72	.67	.70
▲ Accuracy	0.85	<b>.99</b>	<b>1.00</b>	.89	.95	.95	.95
avg-slG	0.85	.30	.20	.78	.69	.71	.68
C-Withdrawn	-	.07	.13	.00	.01	.01	.01
IC-Withdrawn	-	.09	.14	.00	.06	.03	.05



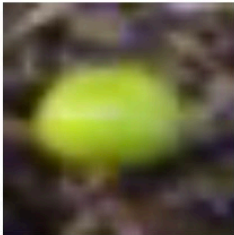
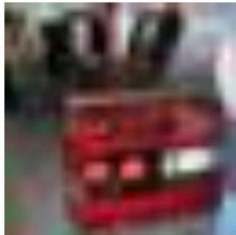

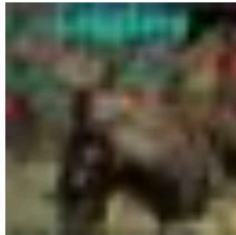


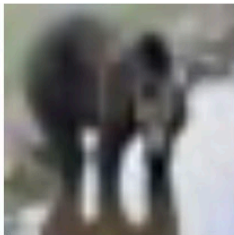
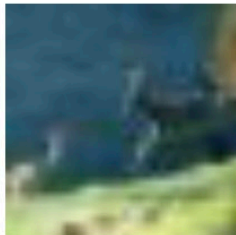
# Experiments: ImageNet-Animal

- Highest IC-Reform
- Highest overall accuracy
- Lowest avg-sIG corresponding to high withdrawns

	Base	Proposed		Deng et al. 2012		Davis et al. 2019	
Confidence	-	90%	95%	90%	95%	90%	95%
C-Corrupt	-	.00	.00	.03	.01	.00	.00
IC-Reform	-	.96	.98	.48	.72	.67	.70
Accuracy	0.85	.99	1.00	.89	.95	.95	.95
 avg-sIG	0.85	<b>.30</b>	<b>.20</b>	.78	.69	.71	.68
 C-Withdrawn	-	<b>.07</b>	<b>.13</b>	.00	.01	.01	.01
 IC-Withdrawn	-	<b>.09</b>	<b>.14</b>	.00	.06	.03	.05



# Visual Examples Across Datasets

	Generalized	C-Withdrawn	IC-Reform	IC-Withdrawn
				
<b>ground truth</b> <i>flat prediction</i> 'Hierarchical Prediction'	<b>apple</b> <i>apple</i> 'Produce'	<b>automobile</b> <i>automobile</i> 'Unknown'	<b>ankle boot</b> <i>sneaker</i> 'Footwear'	<b>lobster</b> <i>forest</i> 'Unknown'
				
<b>hammerhead</b> <i>hammerhead</i> 'Fish'	<b>bag</b> <i>bag</i> 'Unknown'	<b>bear</b> <i>elephant</i> 'Placental'	<b>deer</b> <i>airplane</i> 'Unknown'	



# Conclusion

- Estimation of label posteriors using generalized logits
  - Efficient and compact conditional vector
  - Mitigate issues of lack of validation data
- Label generalization based on semantic hierarchy
  - Bottom-up probabilistic inference framework
  - Able to correct mistakes made by the flat base classifier

Thank you!