

---

# Supplementary Material: Toward Robustness against Label Noise in Training Deep Discriminative Neural Networks

---

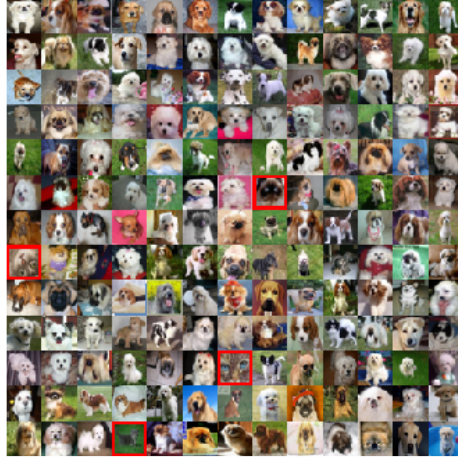
Arash Vahdat  
D-Wave Systems Inc.  
Burnaby, BC, Canada  
avahdat@dwavesys.com

## 1 Visualization

As shown in the E step in Sec. 3.3, the variational distribution  $q$  infers latent clean labels by combining information from both the image-based CRF-CNN model  $p_{\theta}(\hat{\mathbf{y}}, \mathbf{h}|\mathbf{y}, \mathbf{x})$  and the label-based auxiliary distribution  $p_{aux}^{\alpha}(\hat{\mathbf{y}}, \mathbf{h}|\mathbf{y})$ . In our experiments, we observe that in general  $q$  proposes clean labels more accurately than the auxiliary distribution. Fig. 1 compares  $q$  against  $p_{aux}$  in terms of its ability to infer clean labels for a few instances in the noisy training set ( $D_N$ ) for the COCO experiment with actual Flickr tags. In Fig. 2, examples of the recovered clean labels are visualized for the CIFAR-10 experiment.

	<b>Flickr</b> $\emptyset$ <b><math>p_{aux}</math></b> person <b><math>q</math></b> skateboard, person <b>clean</b> skateboard, person
	<b>Flickr</b> $\emptyset$ <b><math>p_{aux}</math></b> person <b><math>q</math></b> person, baseball glove, baseball bat <b>clean</b> person, baseball glove, baseball bat, sports ball, chair, bench
	<b>Flickr</b> 2009, miami <b><math>p_{aux}</math></b> person <b><math>q</math></b> person, tennis racket <b>clean</b> person, tennis racket
	<b>Flickr</b> uploaded:by=flickr_mobile, flickriosapp:filter=NoFilter <b><math>p_{aux}</math></b> person <b><math>q</math></b> person, surfboard <b>clean</b> person, surfboard
	<b>Flickr</b> computer <b><math>p_{aux}</math></b> $\emptyset$ <b><math>q</math></b> laptop, mouse, tv, keyboard <b>clean</b> laptop, mouse, tv, keyboard
	<b>Flickr</b> square, iphoneography, square format, instagram app, uploaded:by=instagram <b><math>p_{aux}</math></b> $\emptyset$ <b><math>q</math></b> cup, dining table, bottle, bowl <b>clean</b> cup, dining table, bottle, bowl, spoon, hot dog
	<b>Flickr</b> food, square, square format, nikon, white, orange, fruit, color, India, photography, table, project365, bowl, colour, wood, 50mm, nikkor, bokeh <b><math>p_{aux}</math></b> orange, apple, banana <b><math>q</math></b> orange, bowl, dining table <b>clean</b> orange, bowl
	<b>Flickr</b> home, light, photo, art, chair, room, table, architecture, apartment, interior, couch, decor, beauty, design, live, lamp, indoor, furniture, relaxed, sofa, flooring, modern <b><math>p_{aux}</math></b> chair, couch, vase, book, dining table, sink, clock, bed, potted plant <b><math>q</math></b> chair, couch, vase, book <b>clean</b> chair, dining table, tv

Figure 1: Visualization of inferred labels for a few instances in the noisy training set ( $D_N$ ) of the COCO dataset. Flickr labels represent the noisy labels extracted from Flickr tags, whereas clean labels are the true labels ignored during training.  $p_{aux}$  and  $q$  correspond to the labels that are extracted using these distribution by thresholding at 0.5. The auxiliary distribution  $p_{aux}$  tends to assign the label “person” to the images with no tag while  $q$  adds more clean labels. In the last two images,  $q$  removes a few unrelated labels.



(a) cat  $\rightarrow$  dog



(b) dog  $\rightarrow$  cat



(c) automobile  $\rightarrow$  truck



(d) horse  $\rightarrow$  deer

Figure 2: Our proposed model can recover clean labels in the noisy training dataset. Here, corrupted instances are visualized for different categories in the CIFAR-10 training dataset. Sub-figures (a) through (d), captioned with *annotated label*  $\rightarrow$  *inferred label*, represents the instances that are labeled with the annotated label but have been assigned to the inferred label by our proposed variational distribution  $q$ . In this visualization, images are sorted based on the confidence of  $q$  for the inferred label from left to right and top to bottom, and the mistaken instances are marked with the red frame. The probability that  $q$  assigns for the inferred label is typically very high ( $> 0.9$ ) for these images, which indicates that  $q$  is confident in changing the noisy labels.