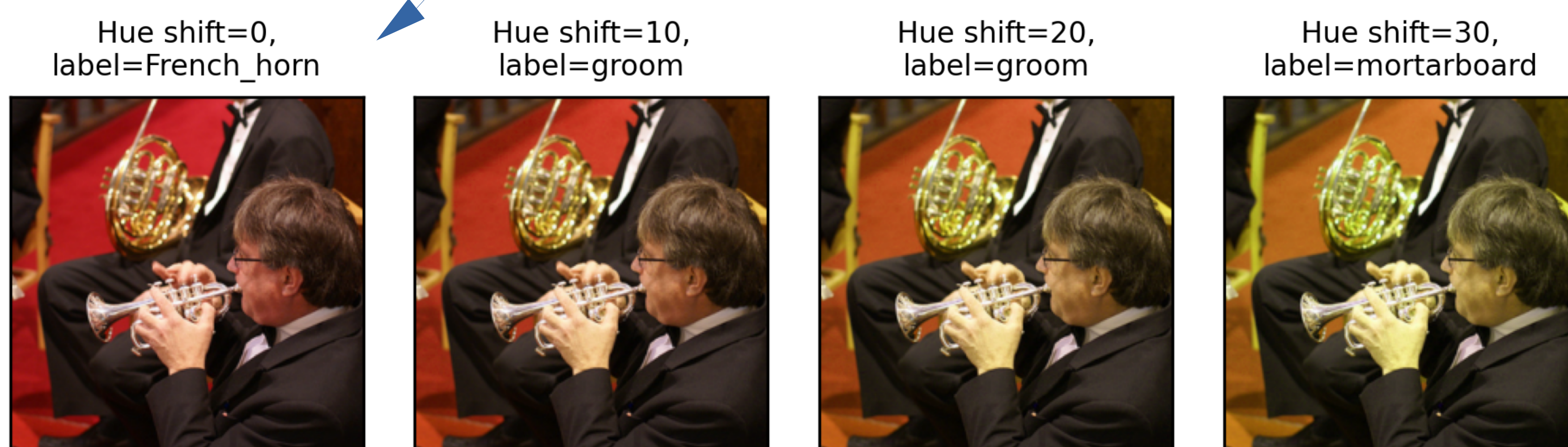


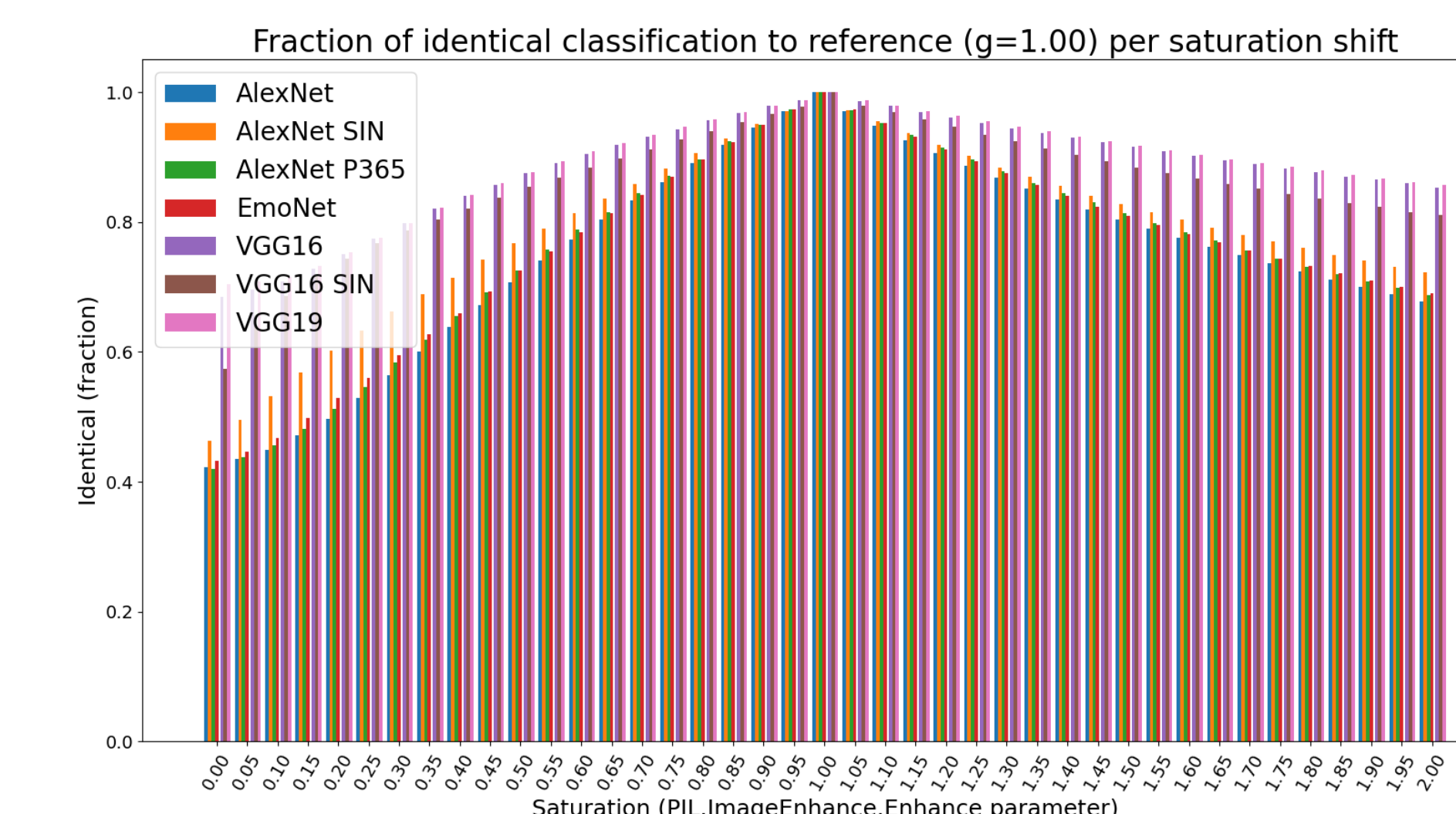
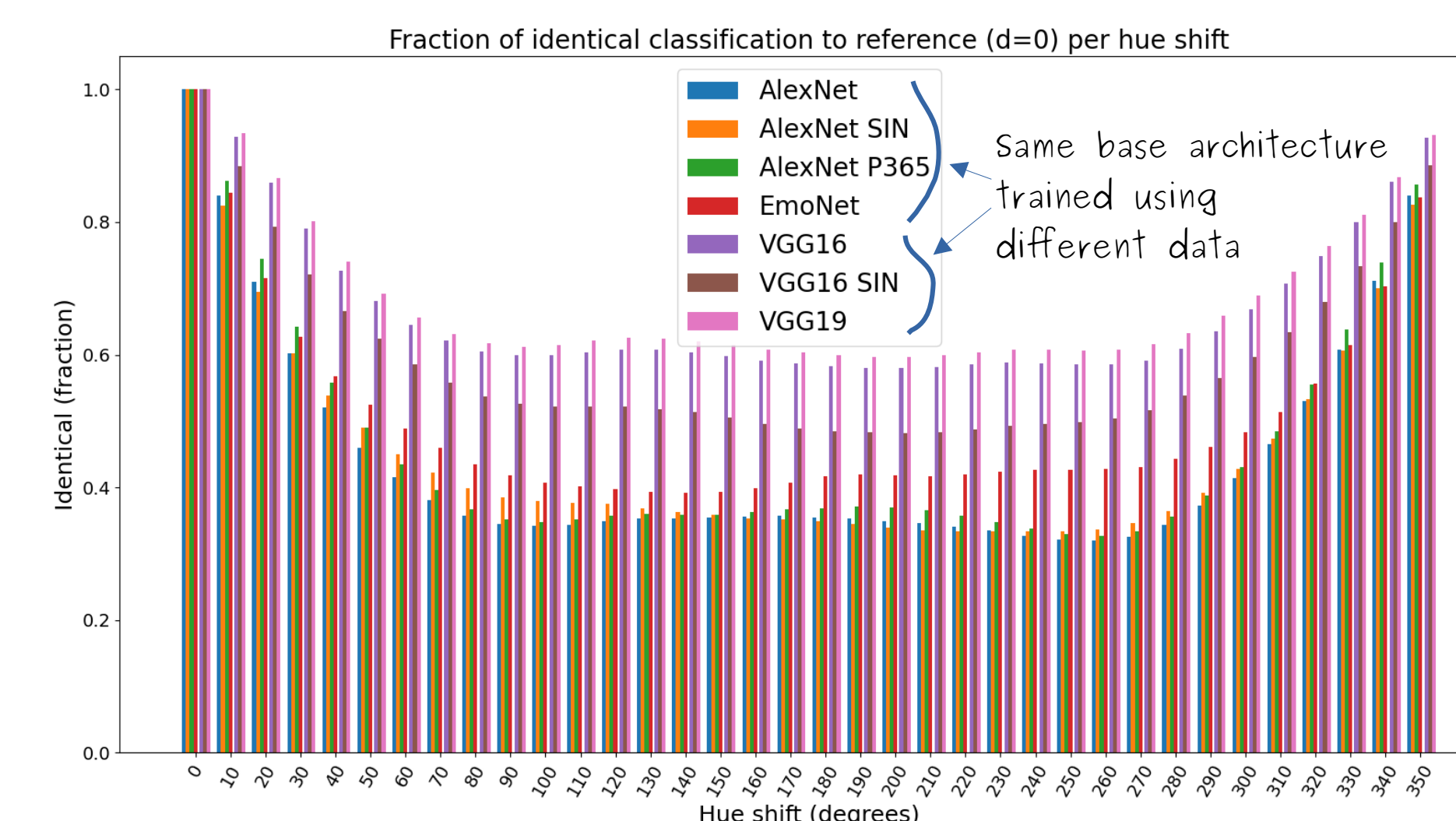
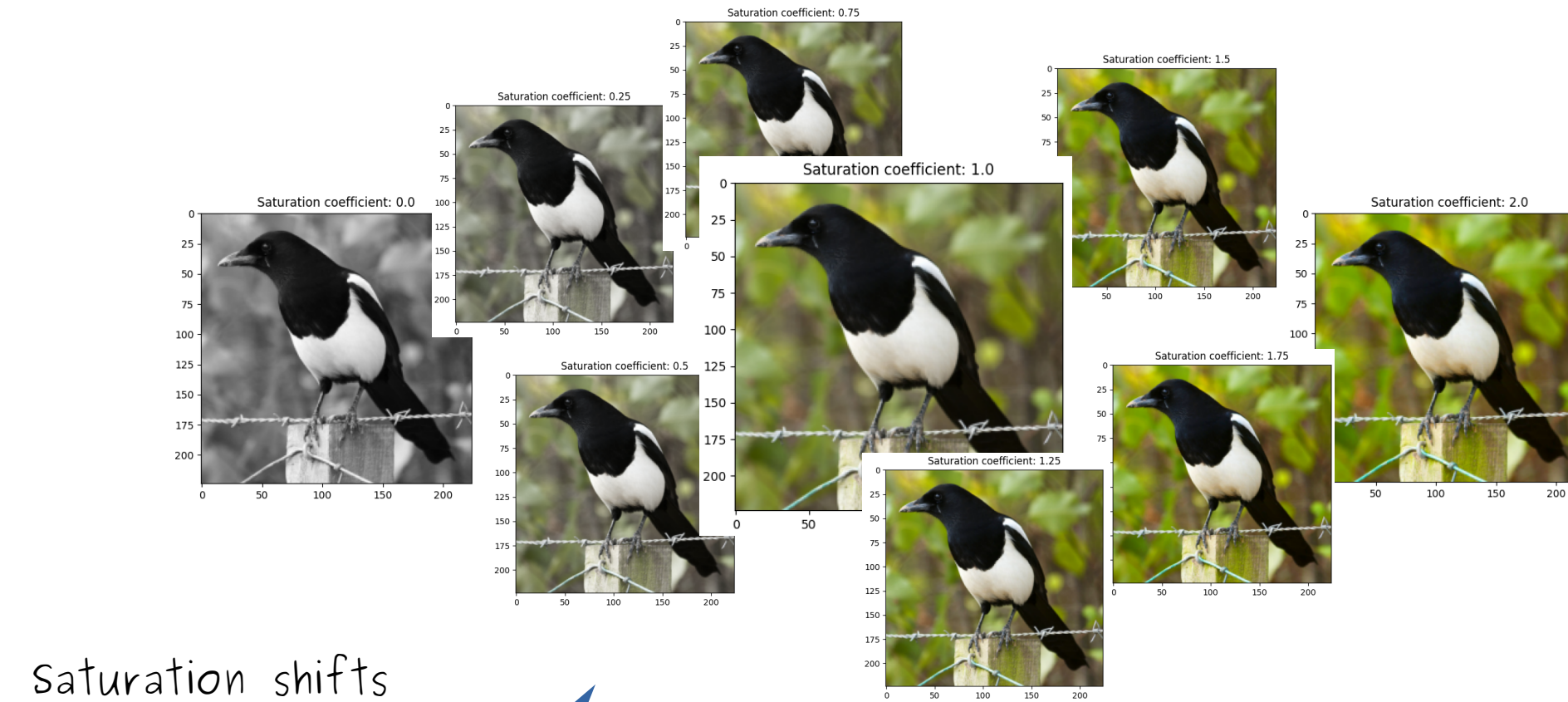
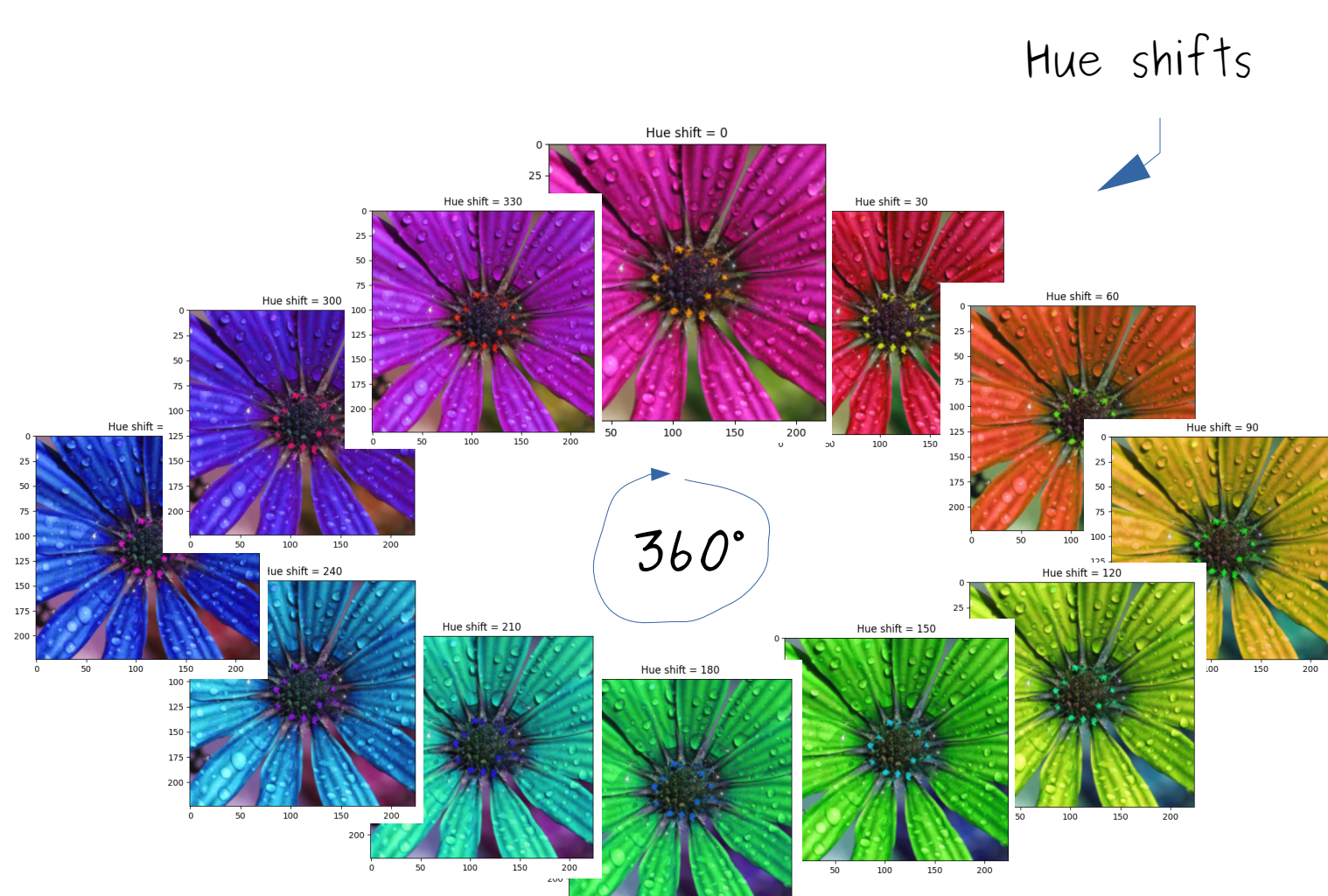
Color-Dependent Prediction Stability of Popular CNN Image Classification Architectures

Question: do you think this image represents something different than the images to its right?



If you answered "no", you disagree with AlexNet, which provided the shown labels!

Turns out that when you show color-modified versions of images...



How often does the model predict the same label as for the original image, regardless of whether that label was correct?

	AlexNet	AlexNet SIN	AlexNet P365	EmoNet	VGG16	VGG16 SIN	ResNet18	ResNet18 P365	ResNet50	ResNet50 SIN	DenseNet161	DenseNet161 P365
Equal d_{all}	.431 ⁻¹⁵	.441 ⁻¹⁴	.447 ⁻¹⁵	.489 ⁻¹²	.659 ⁻¹⁰	.581 ⁻¹²	.659 ⁻¹⁰	.614 ⁻¹¹	.800 ⁻⁰⁶	.632 ⁻¹⁰	.759 ⁻⁰⁷	.662 ⁻⁰⁹
Equal $ d \leq 30$.718 ⁻¹⁰	.709 ⁻⁰⁹	.747 ⁻⁰⁹	.724 ⁻⁰⁹	.861 ⁻⁰⁵	.803 ⁻⁰⁶	.847 ⁻⁰⁵	.830 ⁻⁰⁶	.910 ⁻⁰³	.826 ⁻⁰⁶	.897 ⁻⁰⁴	.847 ⁻⁰⁵
Equal g_{all}	.746 ⁻¹⁵	.786 ⁻¹³	.756 ⁻¹⁵	.758 ⁻¹⁵	.883 ⁻⁰⁸	.857 ⁻⁰⁹	.874 ⁻⁰⁹	.842 ⁻¹⁰	.950 ⁻⁰³	.885 ⁻⁰⁸	.946 ⁻⁰⁵	.863 ⁻⁰⁹
Equal $g \in]0.5, 1.5[\setminus \{1\}$.863 ⁻⁰⁷	.883 ⁻⁰⁶	.872 ⁻⁰⁷	.870 ⁻⁰⁷	.943 ⁻⁰³	.924 ⁻⁰⁴	.937 ⁻⁰³	.919 ⁻⁰⁴	.972 ⁻⁰¹	.940 ⁻⁰³	.977 ⁻⁰¹	.929 ⁻⁰⁴
Top1 d_0, g_1	.566	.400	-	-	.716	.522	.697	-	.803	.602	.771	-
Top1 d_{all}	.331 ⁻⁰⁹	.263 ⁻⁰⁶	-	-	.552 ⁻⁰⁷	.409 ⁻⁰⁵	.548 ⁻⁰⁶	-	.713 ⁻⁰⁴	.488 ⁻⁰⁵	.661 ⁻⁰⁵	-
OL+ d_{all}	.528 ⁻¹⁵	.579 ⁻¹⁴	-	-	.732 ⁻¹⁰	.713 ⁻¹⁰	.743 ⁻⁰⁹	-	.857 ⁻⁰⁵	.752 ⁻⁰⁹	.823 ⁻⁰⁶	-
OL- d_{all}	.305 ⁻¹⁴	.349 ⁻¹⁴	-	-	.475 ⁻¹²	.437 ⁻¹³	.465 ⁻¹²	-	.568 ⁻⁰⁹	.451 ⁻¹²	.541 ⁻¹⁰	-
Top1 g_{all}	.505 ⁻⁰⁶	.382 ⁻⁰²	-	-	.688 ⁻⁰⁴	.511 ⁻⁰²	.668 ⁻⁰⁴	-	.798 ⁻⁰¹	.591 ⁻⁰²	.762 ⁻⁰²	-
OL+ g_{all}	.855 ⁻¹⁴	.846 ⁻¹⁴	-	-	.917 ⁻¹⁵	.874 ⁻¹⁴	.915 ⁻¹⁵	-	.938 ⁻¹⁵	.897 ⁻¹⁵	.934 ⁻¹⁵	-
OL- g_{all}	.253 ⁻¹⁷	.345 ⁻¹⁵	-	-	.379 ⁻¹⁶	.361 ⁻¹¹	.359 ⁻¹⁴	-	.451 ⁻⁰⁹	.362 ⁻¹⁰	.399 ⁻¹⁰	-
O.P. d_{all}	58.5 ¹²²	40.6 ⁹⁶	29.7 ⁵³	3.7 ³	28.0 ⁷⁶	26.5 ⁷³	25.3 ⁷¹	13.0 ²⁹	31.6 ¹¹⁹	24.5 ⁷⁰	19.1 ⁶¹	9.9 ²⁴
O.P. $ d \leq 30$	10.9 ³⁷	6.7 ²³	4.9 ¹³	2.1 ²	5.2 ¹⁸	4.4 ¹⁶	5.0 ¹⁸	2.8 ⁶	6.4 ⁴⁵	4.2 ¹⁴	4.0 ¹⁵	2.6 ⁶
O.P. g_{all}	11.9 ⁴⁰	19.8 ⁵⁵	20.0 ⁴⁴	2.3 ²	5.3 ²⁰	4.3 ¹⁷	4.9 ¹⁸	3.8 ¹⁰	5.5 ²³	11.9 ³⁶	2.9 ¹⁰	2.9 ⁸
O.P. $g \in]0.5, 1.5[\setminus \{1\}$	2.5 ⁵	18.9 ⁵²	19.4 ⁴²	1.4 ¹	1.6 ²	1.5 ²	1.6 ²	1.5 ¹	5.5 ²²	11.9 ³⁶	1.2 ¹	1.4 ¹

- Behavior seems independent of training data.
- Behavior is inherited through transfer learning.
- Larger architectures appear less sensitive.
- Effect is more pronounced for images for which the model originally made a wrong prediction.
- Different label is more than just "flipping second and first place"

(And yes, additional pre-processing during training alleviates the problem...)

...to popular CNN models...

...they often alter their predictions!

Table 1. Overview of training and validation data per model. "(ModelName)" is a placeholder for a valid architecture, "IN-1k" = ImageNet-1k, "SIN" = Stylized ImageNet, "train" = train data, "val" = validation data.

Model	Trained on	Validated on
AlexNet, VGG16, ResNet18/50, DenseNet161	IN-1k train	IN-1k val
(ModelName)-SIN	SIN	IN-1k val
(ModelName)-P365	Places365 train	Places365 val
EmoNet	IN-1k train + EmoNet	IN-1k val