

#### WHY DO WE NEED TO INTEGRATE PRIOR SEMANTIC KNOWLEDGE?

The New York Times		Semanti
GOOGLE Google Ph BY CONOR DOUGHE	otos Mistakenly Labels Black People 'Gorillas'	always o similarity
<b>Email</b>	Google continued to apologize Wednesday for a flaw in Google Photos, which was released to <u>great fanfare</u> in May, that led the new	
<b>f</b> Share	application to mistakenly label photos of black people as "gorillas." The company said it had fixed the problem and was working to	
🔰 Tweet	figure out exactly how it happened.	
Save	"We're appalled and genuinely sorry that this happened," said a Google representative in an emailed statement. "We are taking immediate action to prevent this type of result from appearing."	
More		

- Simply removing the "gorilla" class is a sub-optimal "fix" and only works for classification, not for content-based image retrieval.
- Better: learning a feature representation that carries semantic information.

#### **MEASURING SEMANTIC SIMILARITY**

- Previous works focused on learning the semantic similarity between classes either from text [1] or the images themselves [2, 3].
- However, semantic knowledge is already available in the form of class taxonomies for almost all concepts of the world (WordNet [4], Wikispecies, Open Tree of Life, ...).
- We use a semantic similarity measure derived from this graph of concepts [5].



# **Hierarchy-based Image Embeddings for Semantic Image Retrieval**

Björn Barz and Joachim Denzler Computer Vision Group, Friedrich Schiller University Jena, Germany

ic similarity does not correlate with visual



visually simila

semantically similar

### HIERARCHY-BASED SEMANTIC EMBEDDINGS





Map images onto their target class embeddings in this semantic space using a simple loss: with additional classification objective **Cosine Distance** cross-entropy loss  $\mathcal{L}_{\text{CORR+CLS}}(x, y)$ 

 $\mathcal{L}_{\text{CORR}}(x, y) = 1 - \psi(x)^T \varphi(c_y)$ 

L2-normalized CNN output

Code available! github.com/cvjena/semantic-embeddings

Goal: learn image features whose dot product resembles the semantic similarity of their classes

semantically similar, though visually dissimilar





Class embeddings on the unit hyper-sphere do not need to be learned but can be computed explicitly by solving:

 $\forall_{1 \leq i,j \leq n}: \varphi(c_i)^T \varphi(c_j) = \operatorname{sim}(c_i, c_j)$  $\forall_{1 \leq i \leq n} : \|\varphi(c_i)\| = 1$ known classes class embedding function

 $= \mathcal{L}_{\text{CORR}}(x, y) + \lambda \cdot \mathcal{L}_{\text{XENT}}(f(\psi(x)), y)$ 

additional FC layer with softmax —





IFAR-100: chimpanze







- Straightforward training.





## Semantically more consistent retrieval results. Semantically meaningful feature space based on prior knowledge.

[1] Frome, Corrado, Shlens, Bengio, Dean, Ranzato, Mikolov. "DeViSE: A Deep Visual-Semantic I

2] Wen, Zhang, Li, Qiao. "A Discriminative Feature Learning Approach for Deep Face Recognition

Sun, Wei, Ren, Ma. "Label Embedding Network: Learning Label Representation for Soft Training o Deep Networks." arXiv:1710.10393, 2017.

Fellbaum. "WordNet." Wiley Online Library, 1998.

5] Deng, Berg, Fei-Fei. "Hierarchical Semantic Indexing for Large-Scale Image Retrieval." CVPR 200