

What are deep vision models learning?

David Crandall

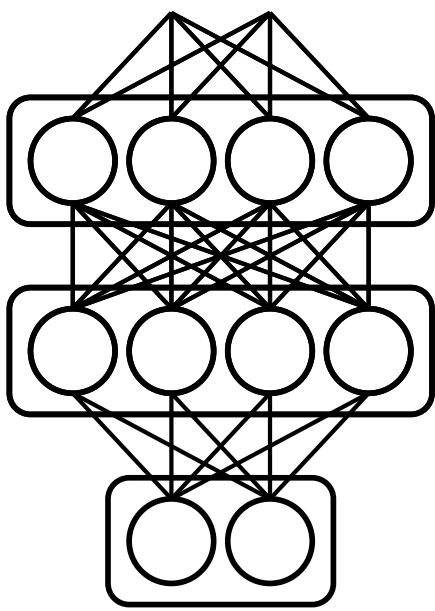
Unsupervised Visual Representation Learning by Context Prediction

Carl Doersch

Joint work with Alexei A. Efros & Abhinav Gupta

ICCV 2015

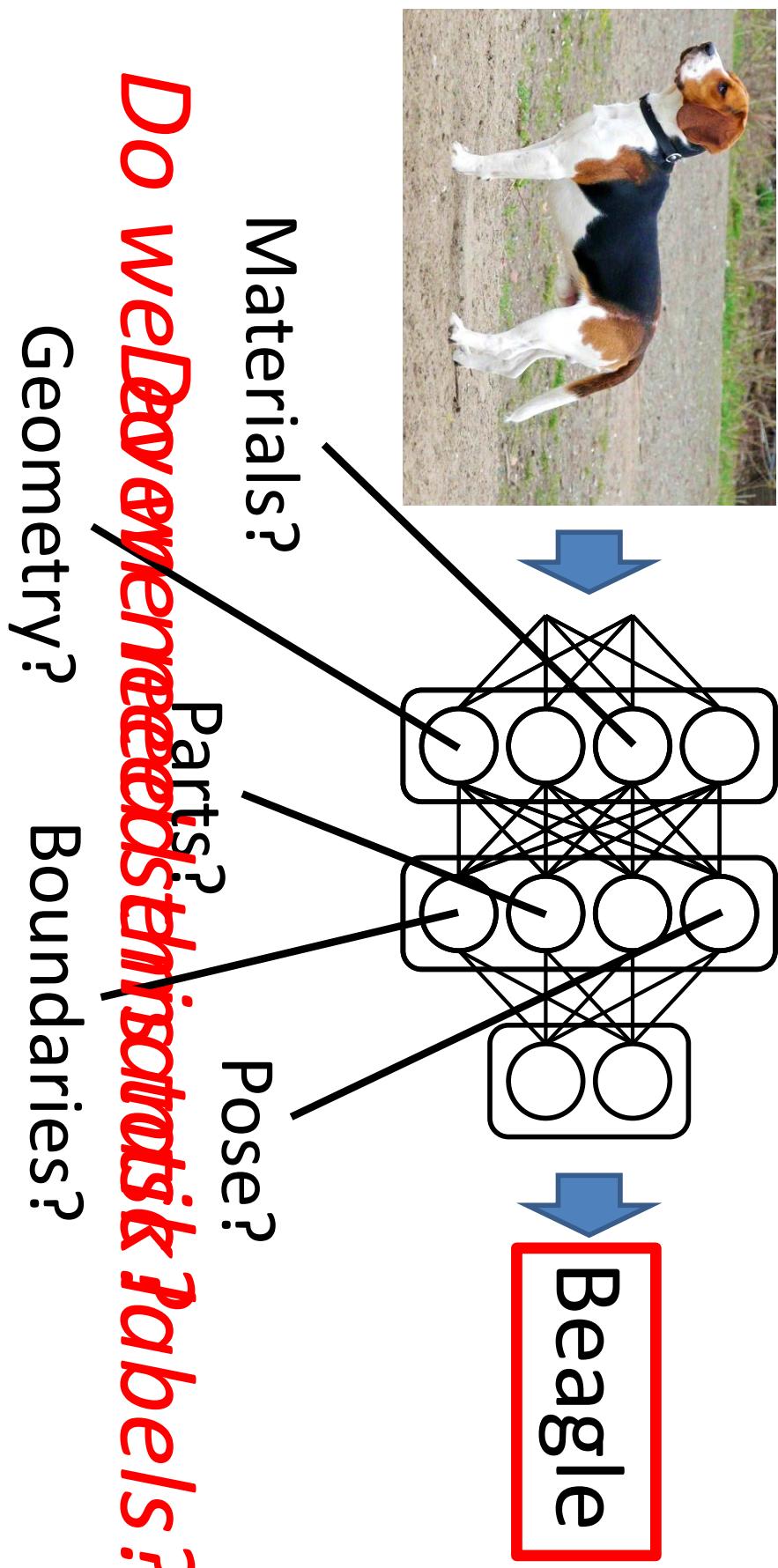
ImageNet + Deep Learning



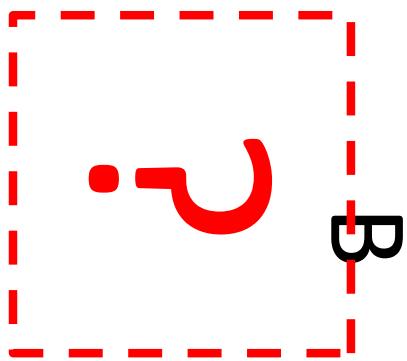
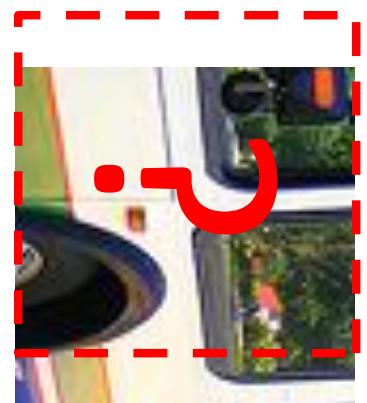
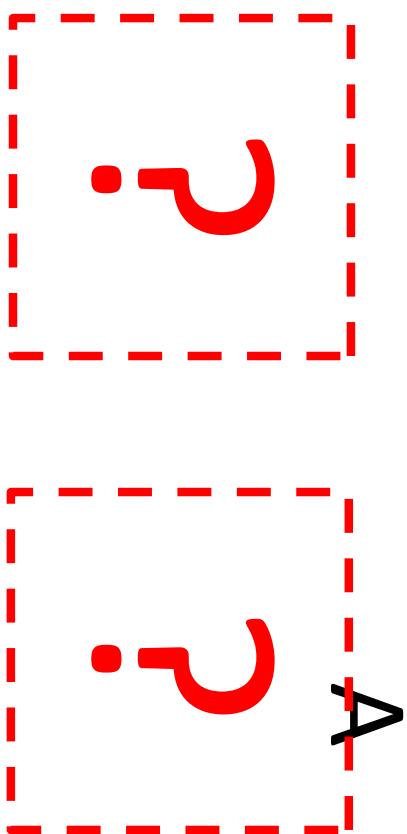
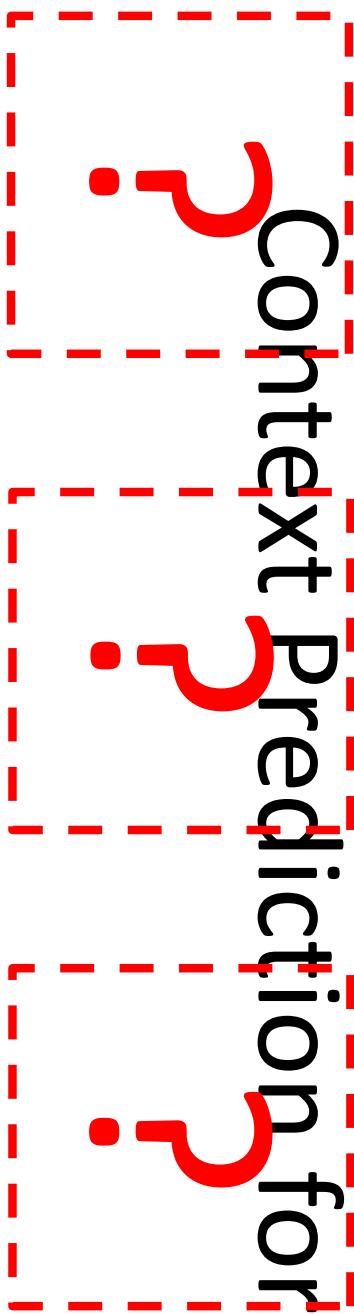
Beagle

- Image Retrieval
- Detection (RCNN)
- Segmentation (FCN)
- Depth Estimation
- ...

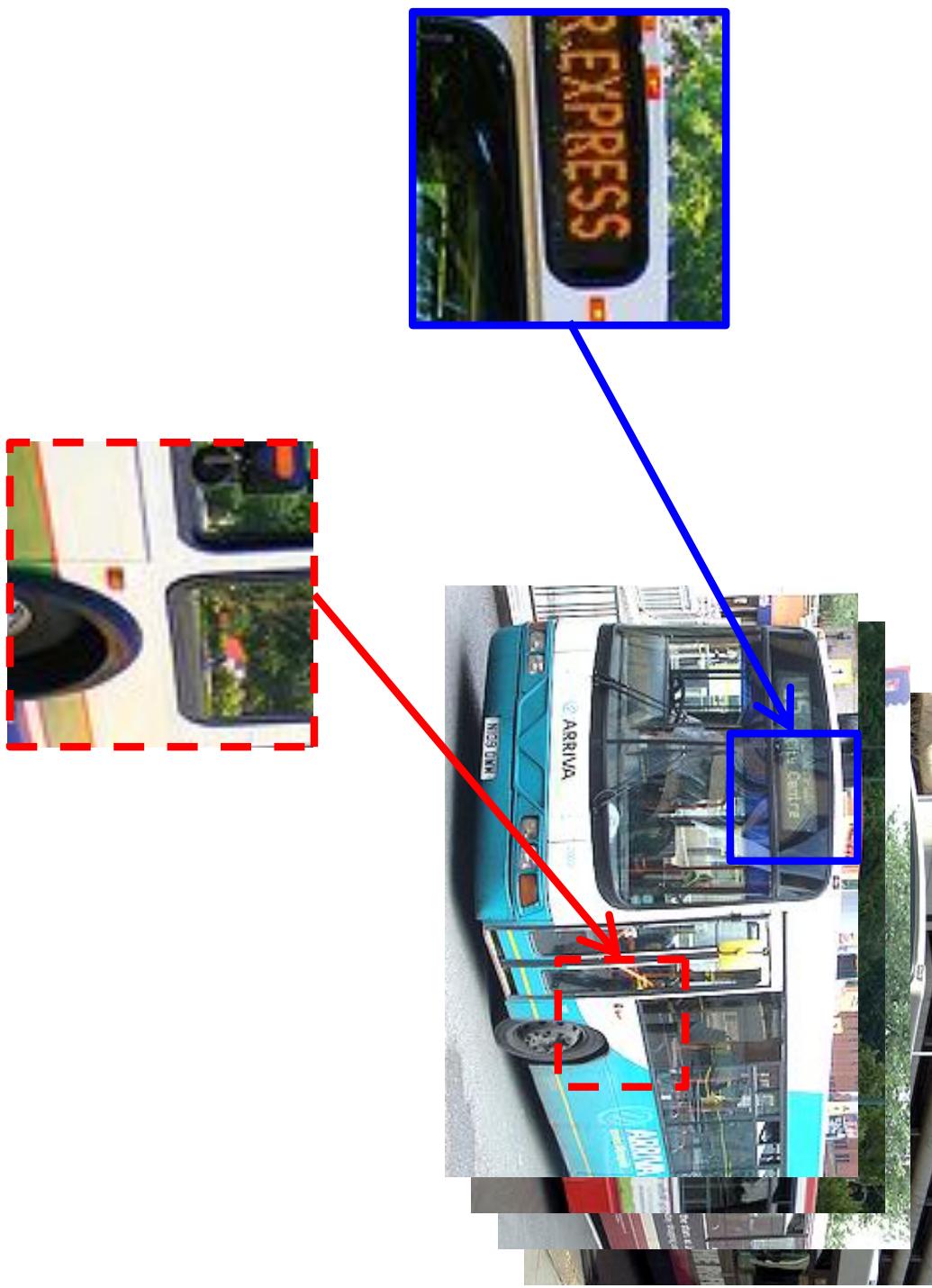
ImageNet + Deep Learning



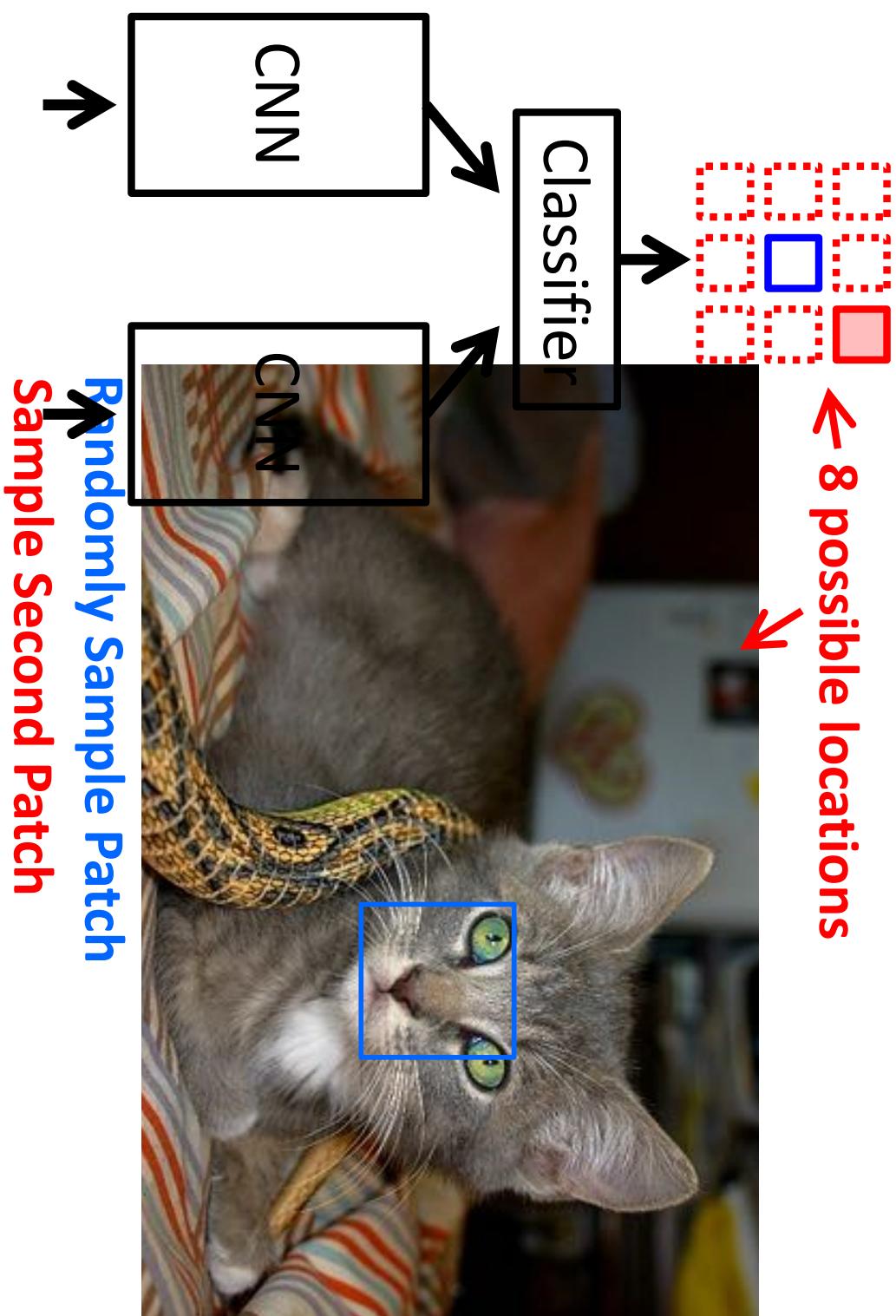
Context Prediction for Images

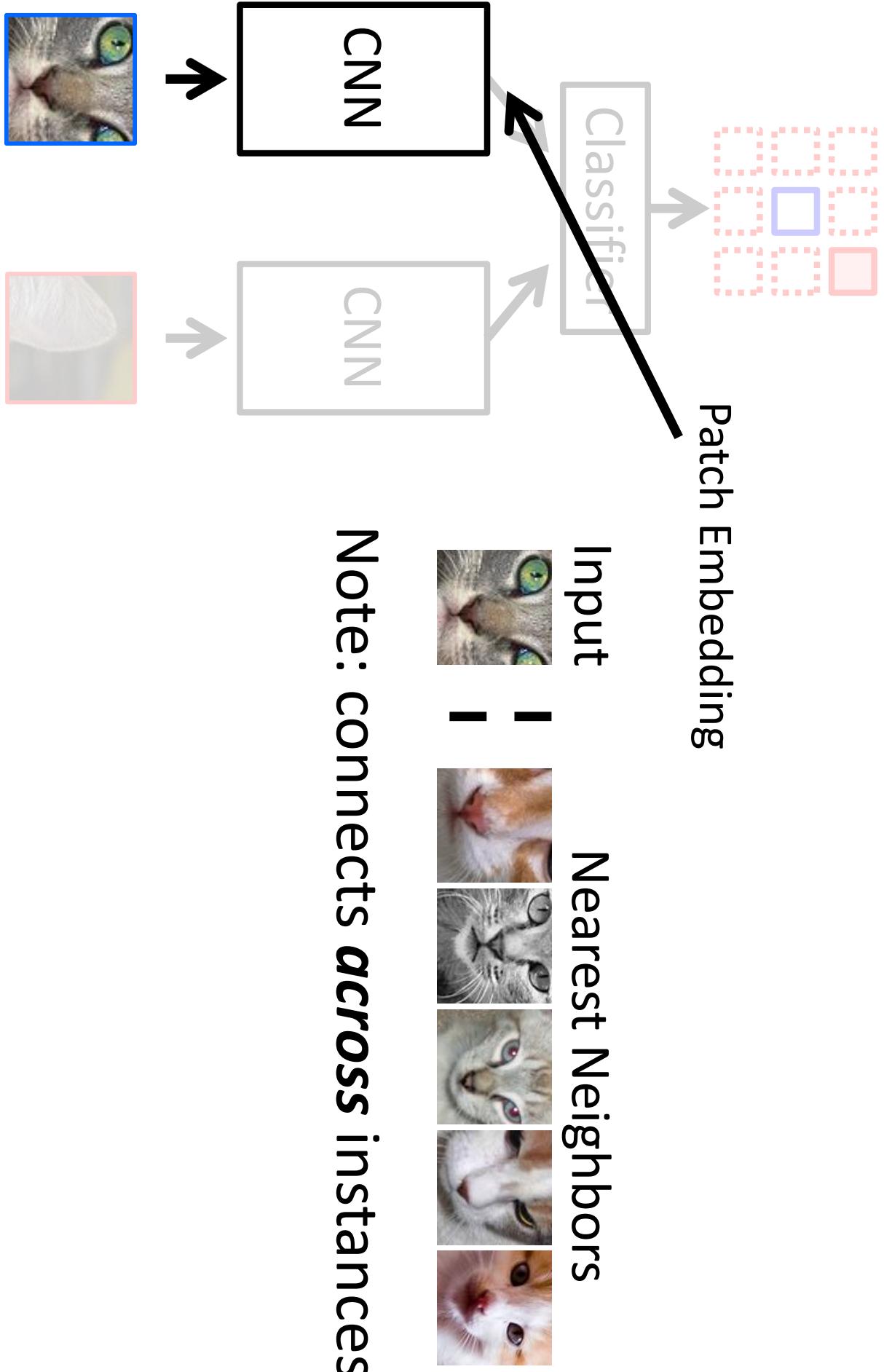


Semantics from a non-semantic task

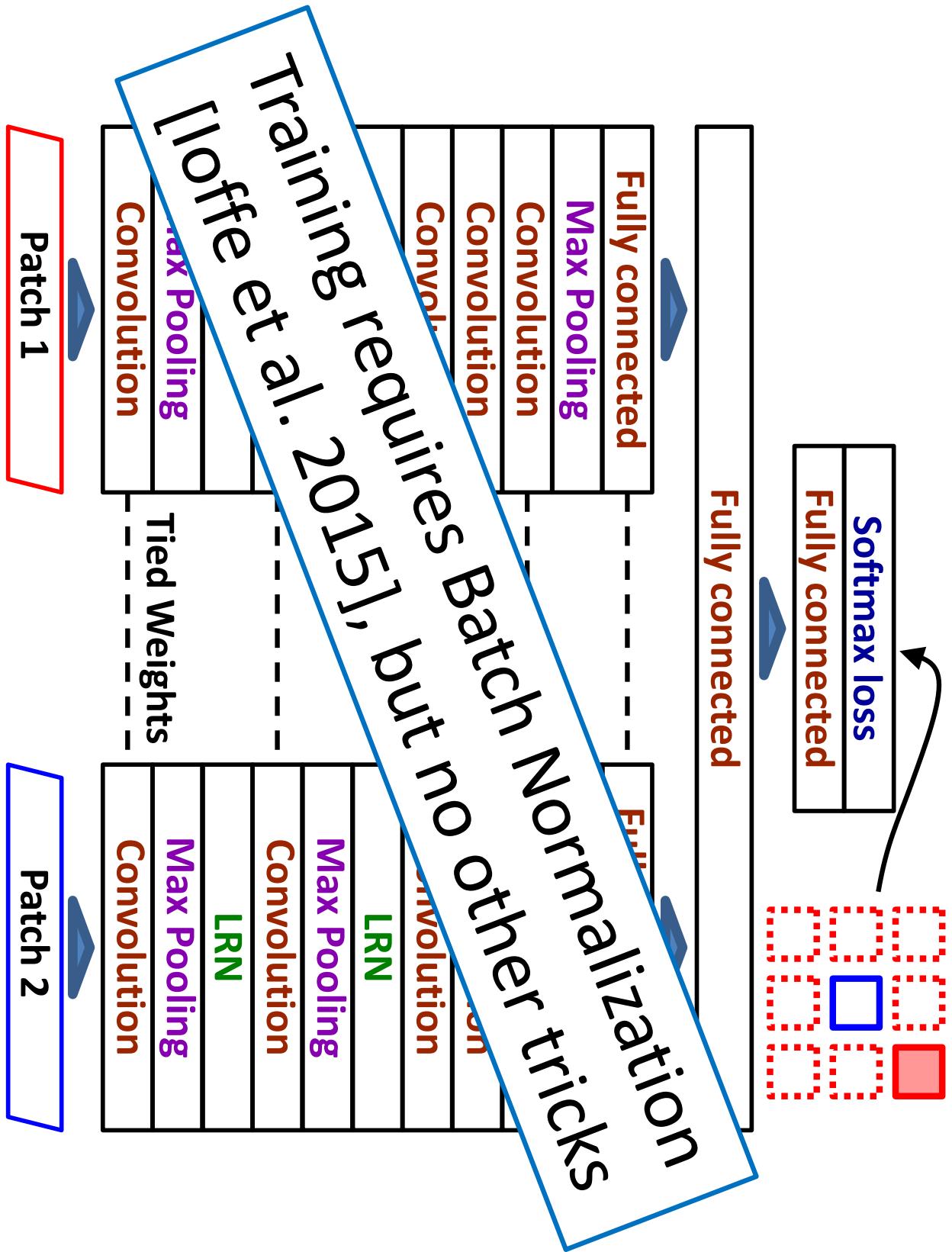


Relative Position Task

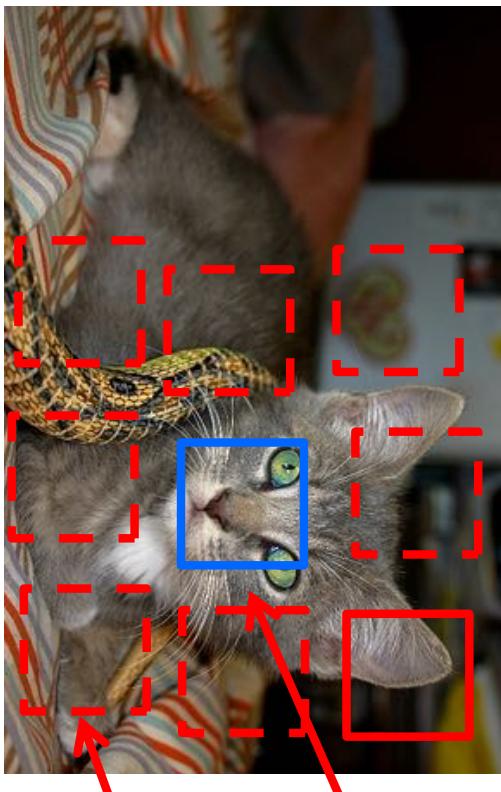




Architecture

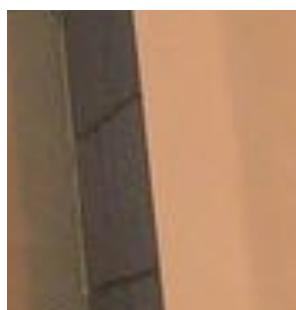


Avoiding Trivial Shortcuts



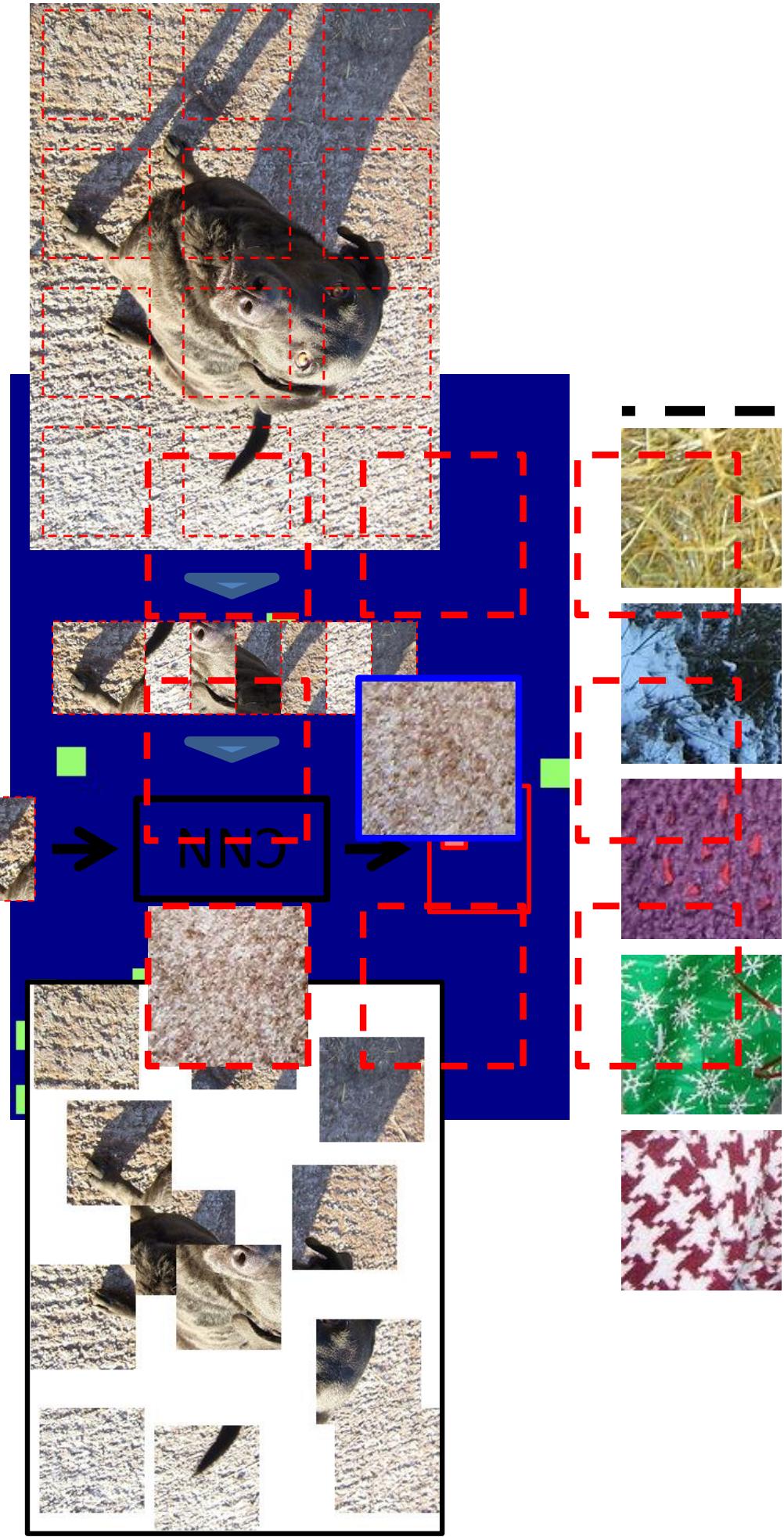
jitter the patch locations

Include a gap

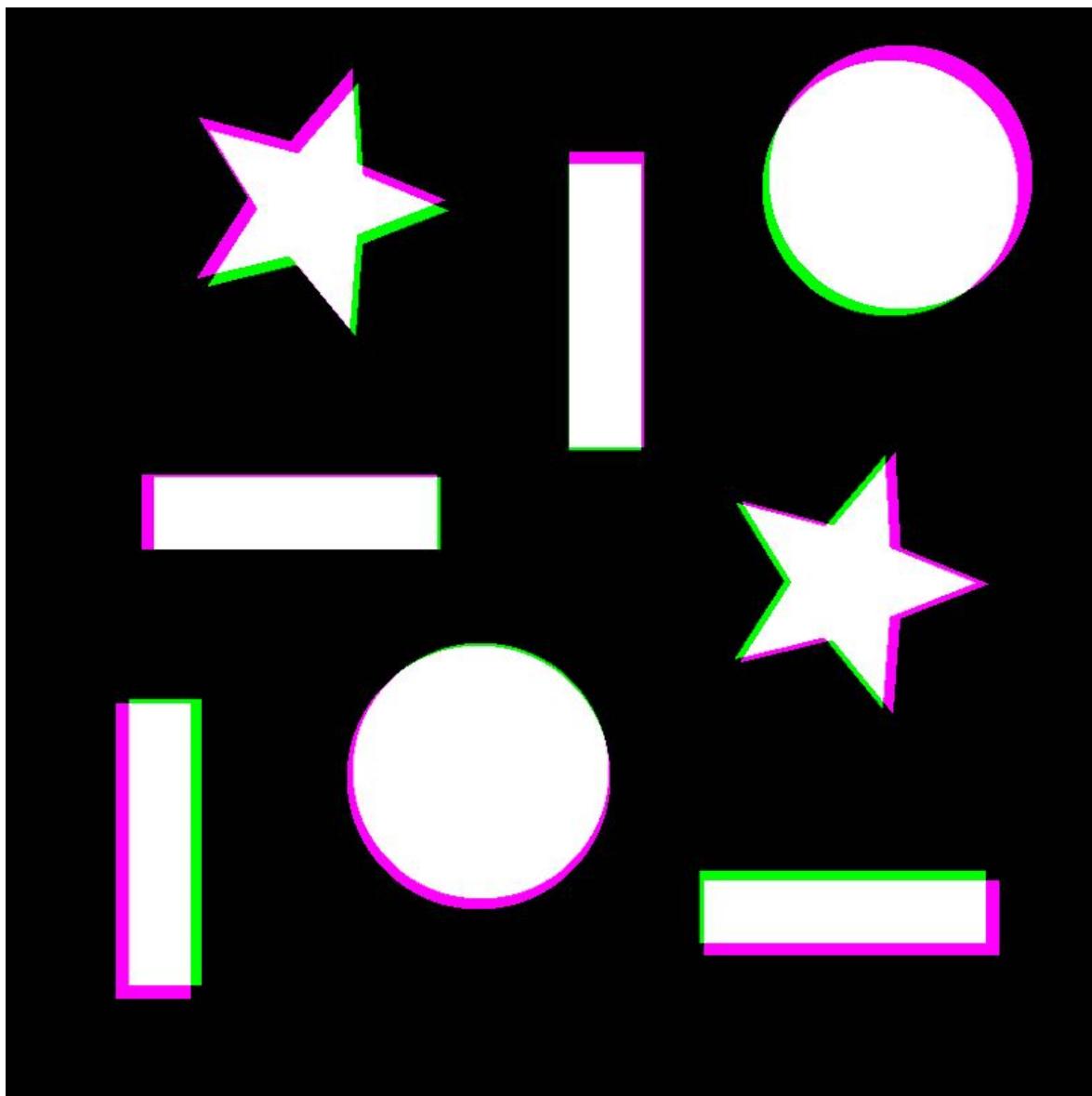
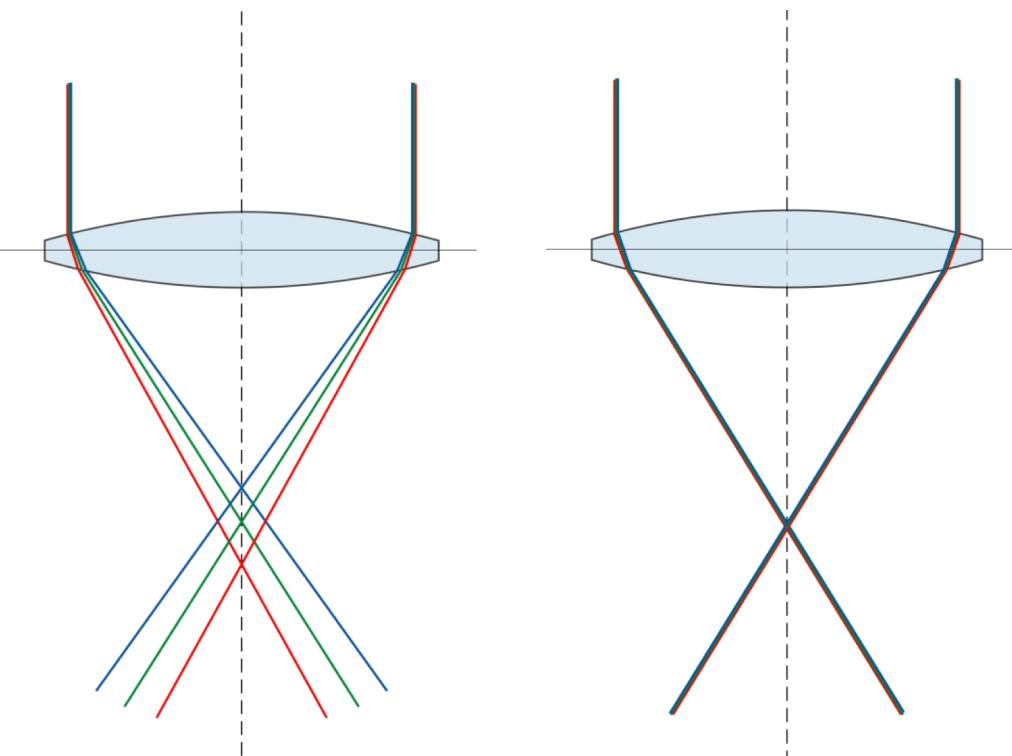


A Not-So “Trivial” Shortcut

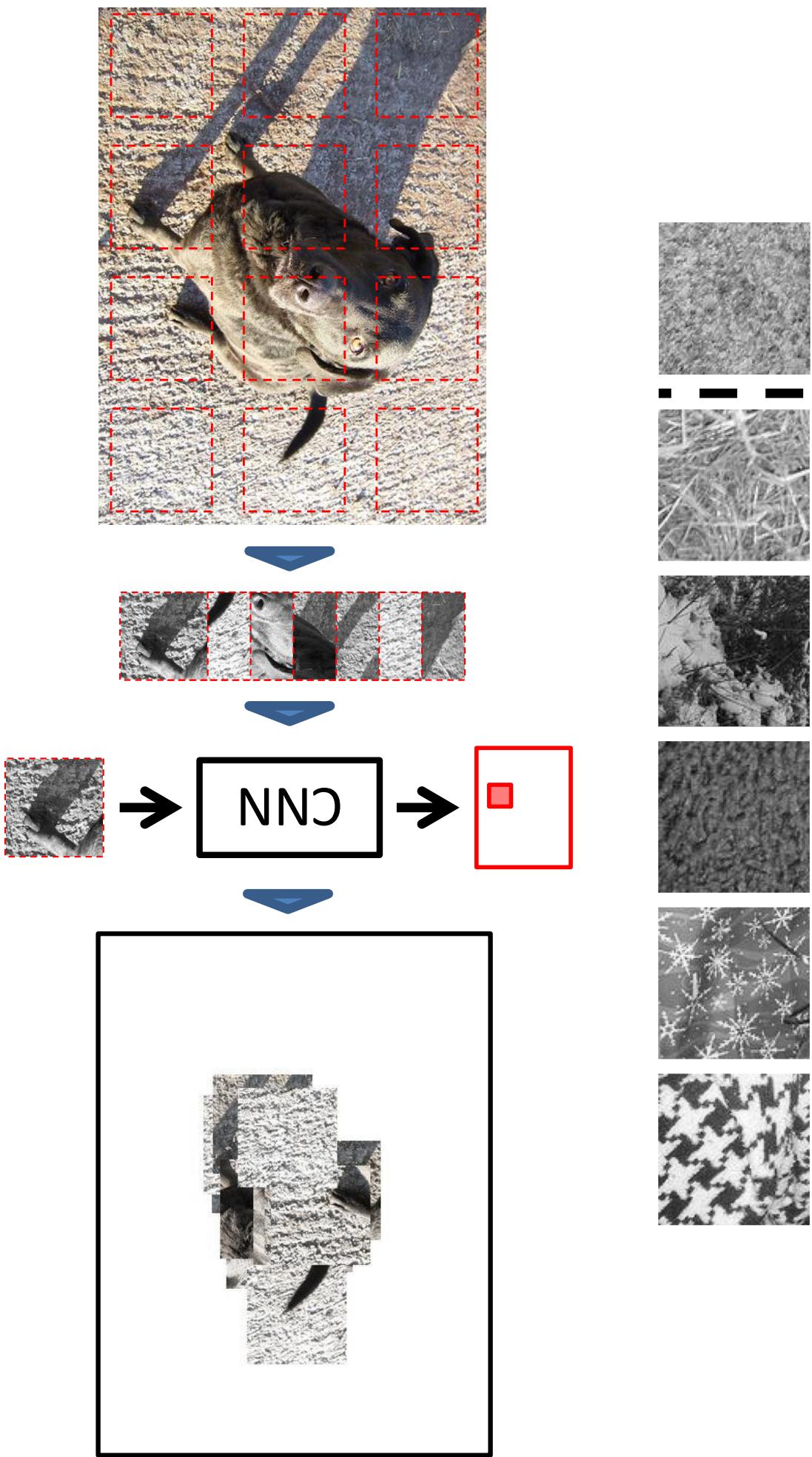
Position in image



Chromatic Aberration

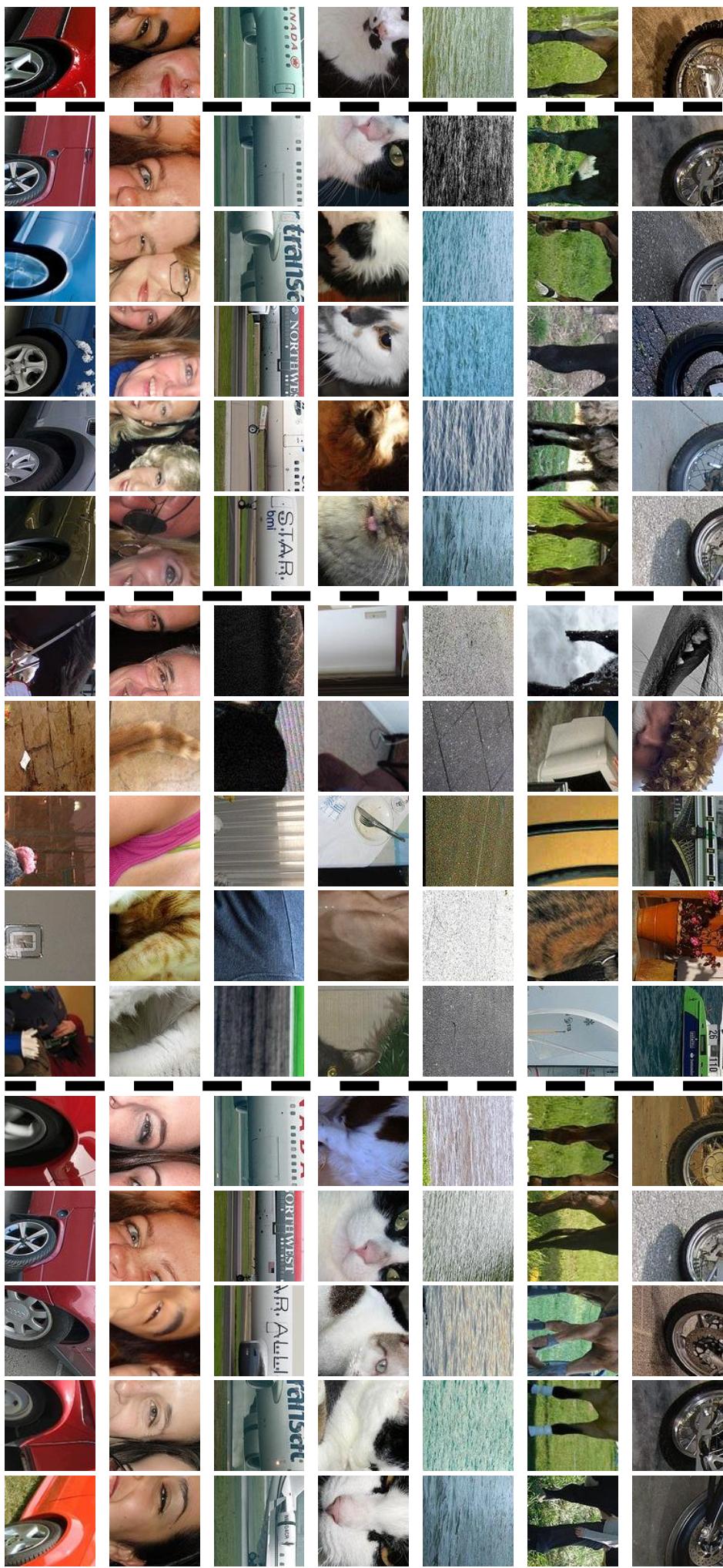


Chromatic Aberration

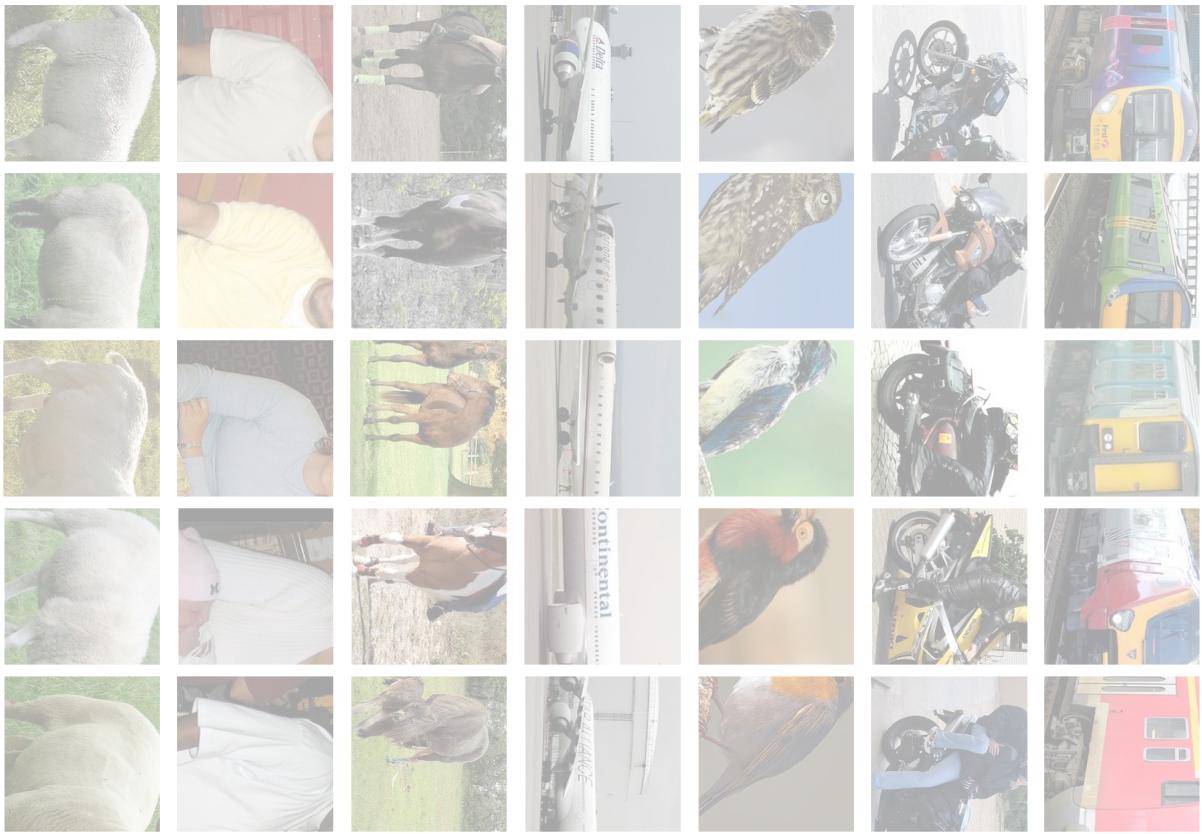
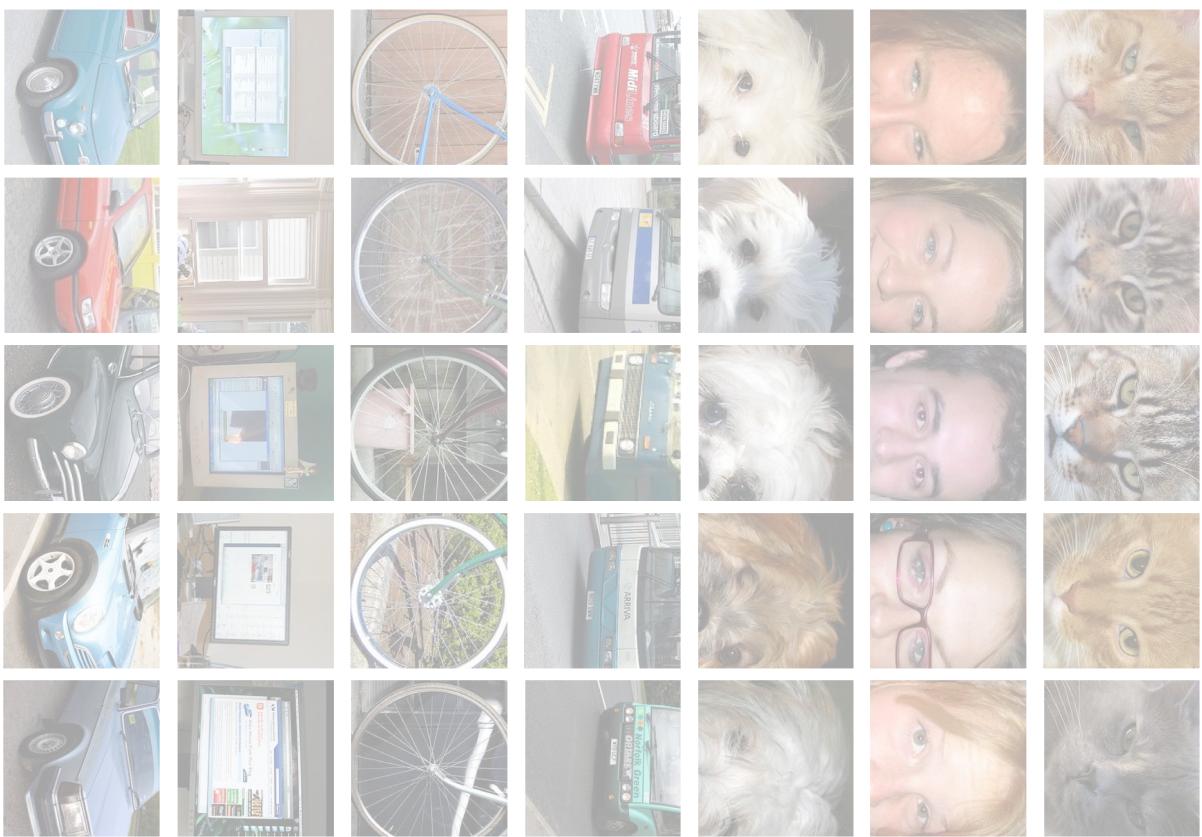


What is learned?

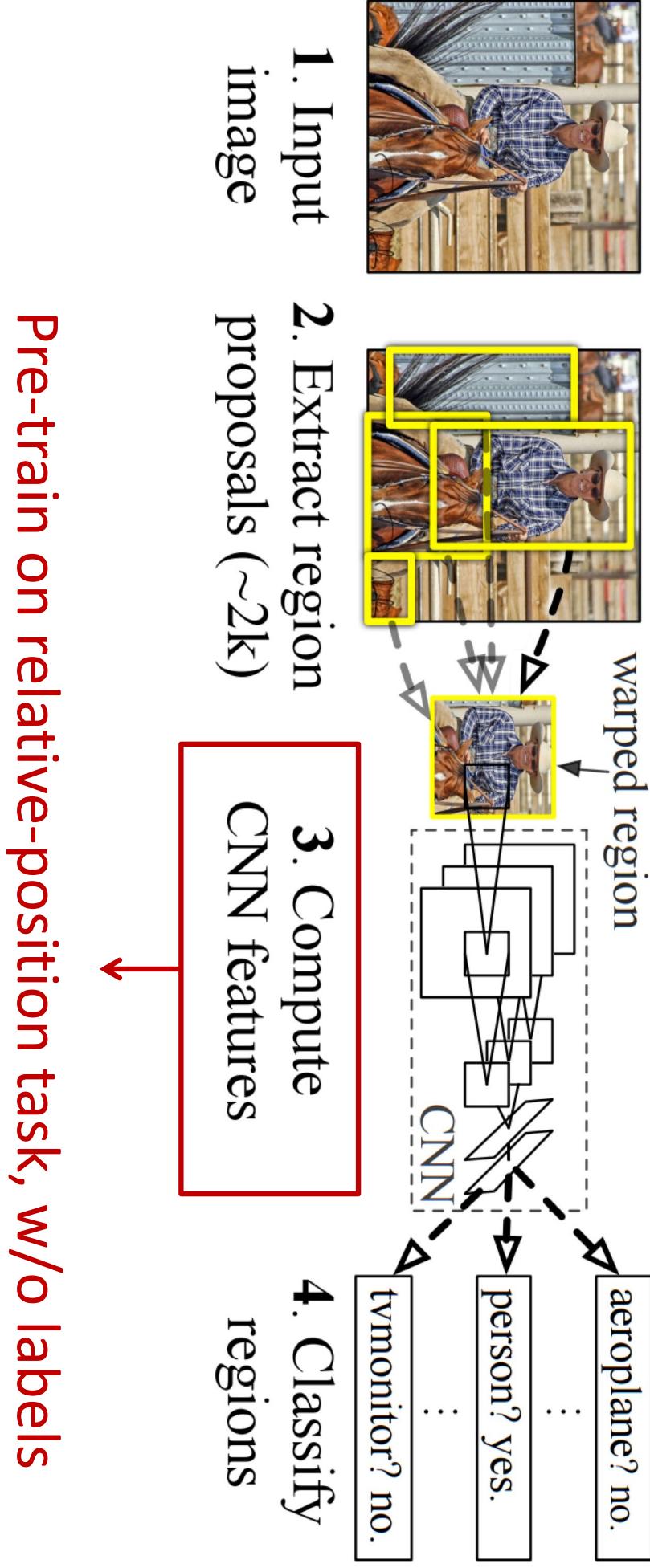
Input Ours Random Initialization ImageNet AlexNet



Mined from Pascal VOC2011



Pre-Training for R-CNN



VOC 2007 Performance

(pretraining for R-CNN)

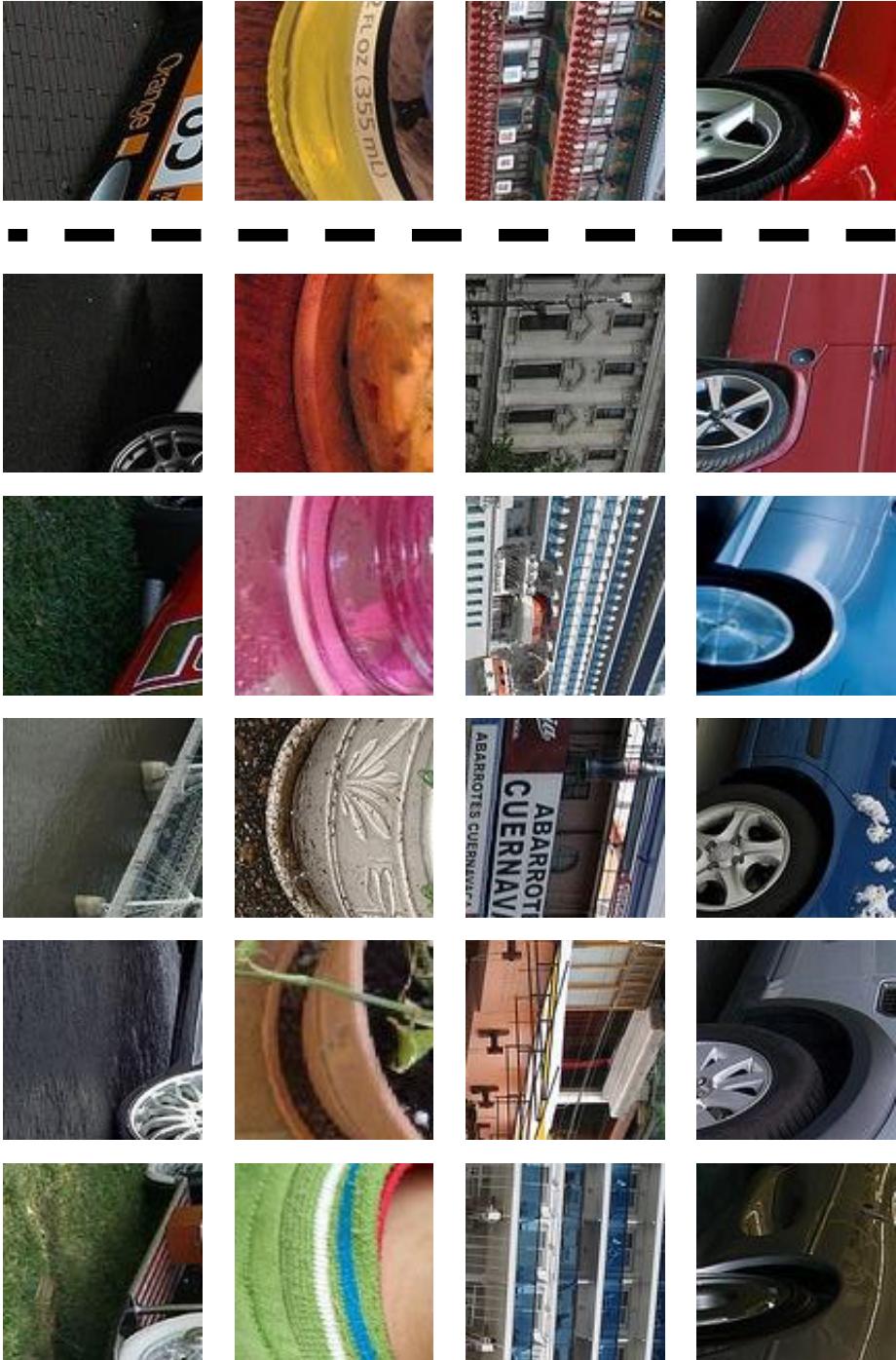
% Average Precision

ImageNet Labels Ours No Pretraining

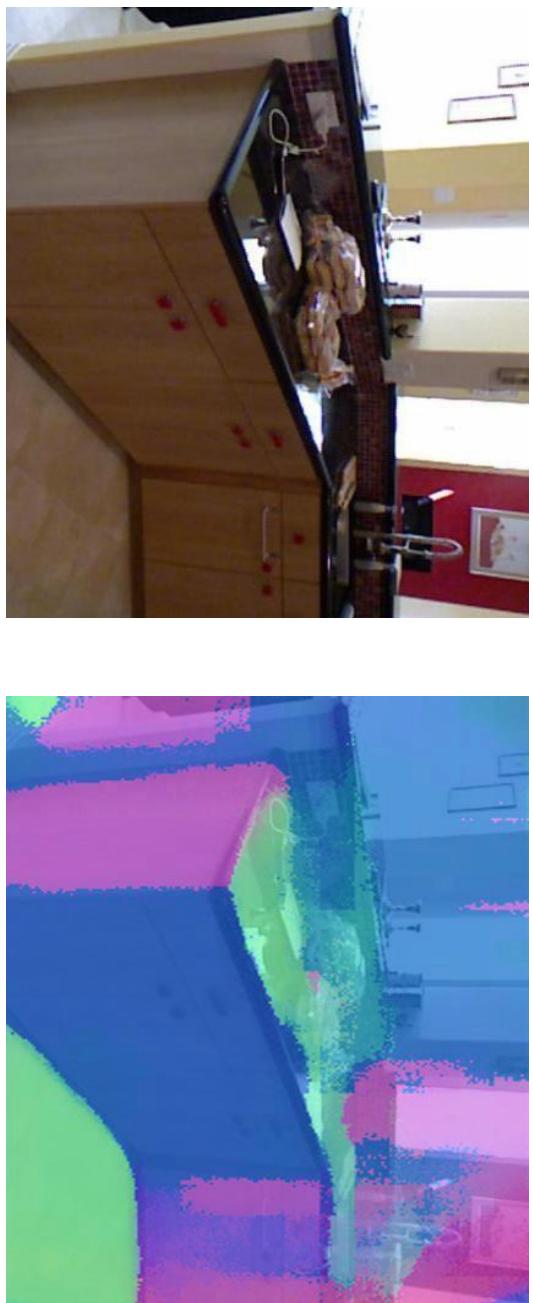
	No Rescaling	Krähenbühl et al. 2015	VGG + Krähenbühl et al.
68.6			
54.2			
56.8			
51.1			
46.3			
45.6			
40.7			
42.4			

[Krähenbühl, Doersch, Donahue & Darrell, “Data-dependent Initializations of CNNs”, 2015]

Capturing Geometry?



Surface-normal Estimation

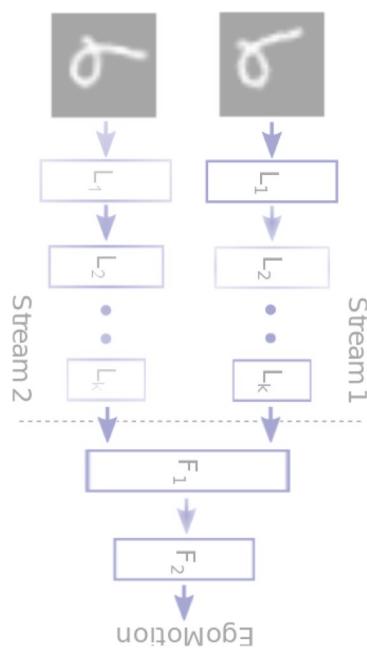


Method	Error (Lower Better)		% Good Pixels (Higher Better)		
	Mean	Median	11.25°	22.5°	30.0°
No Pretraining	38.6	26.5	33.1	46.8	52.5
Ours	33.2	21.3	36.0	51.2	57.8
ImageNet Labels	33.3	20.8	36.7	51.7	58.1

So, do we need semantic labels?

“Self-Supervision” and the Future

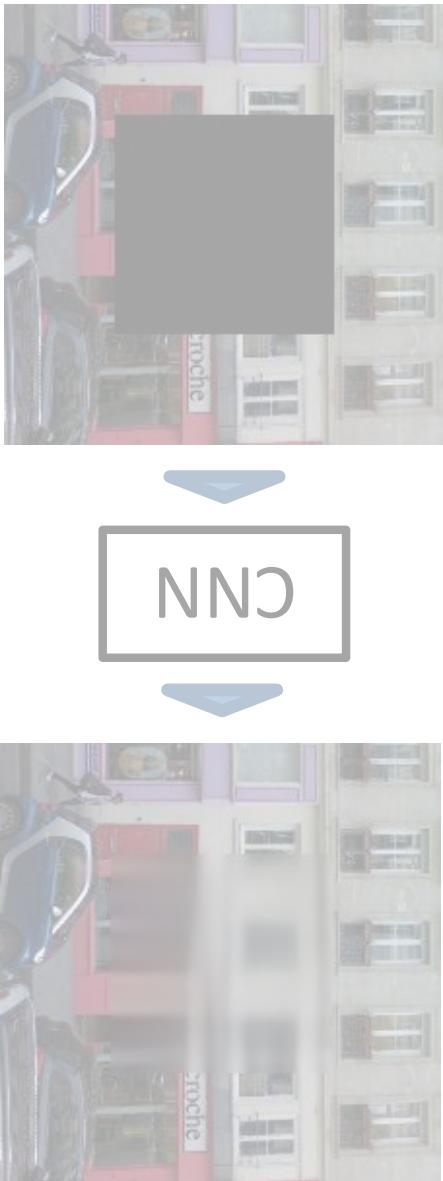
Ego-Motion



[Agrawal et al. 2015; Jayaraman et al. 2015]

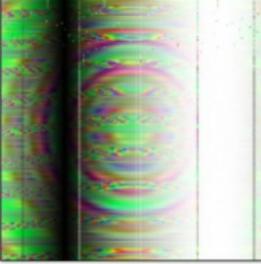
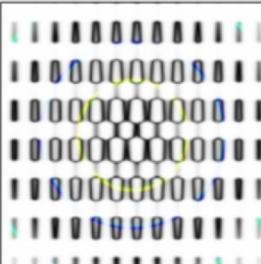
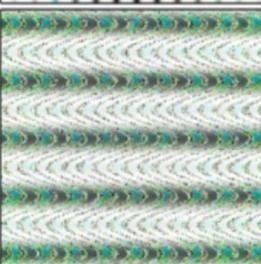
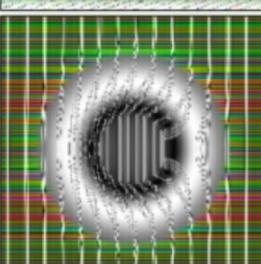
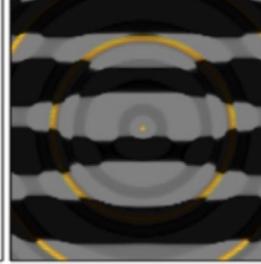
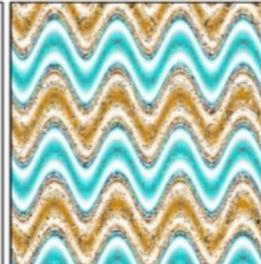
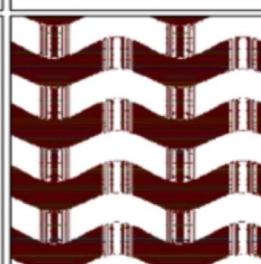
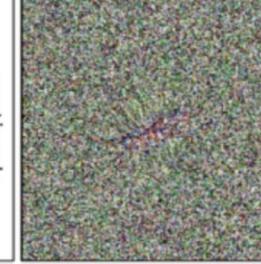
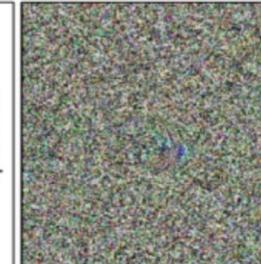
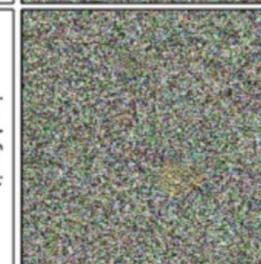
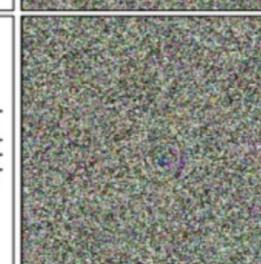
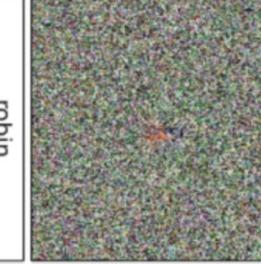
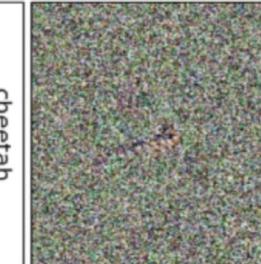
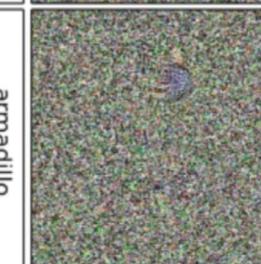
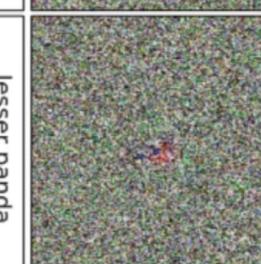
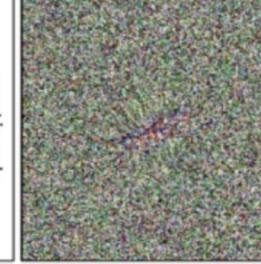
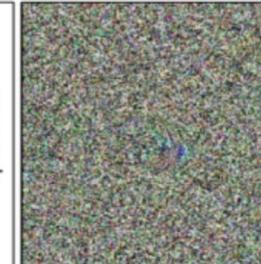
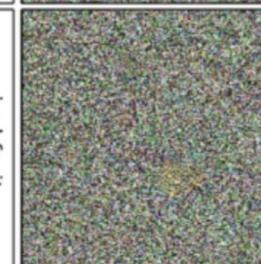
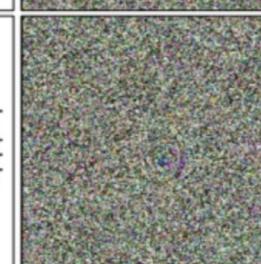
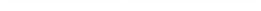
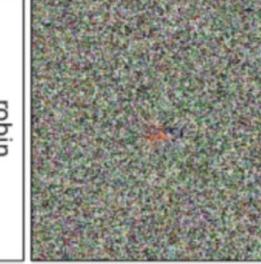
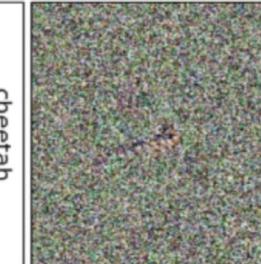
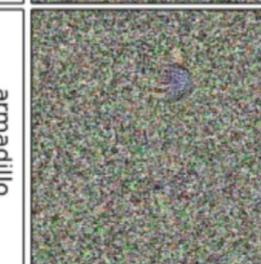
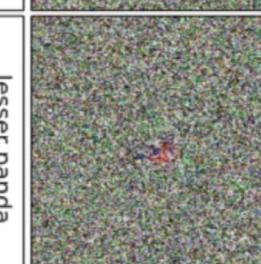
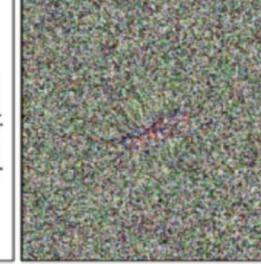
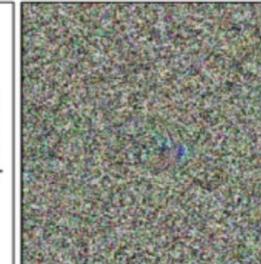
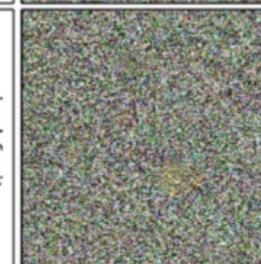
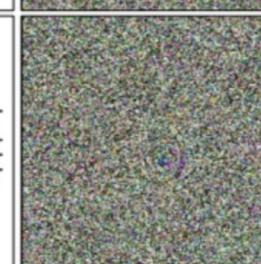
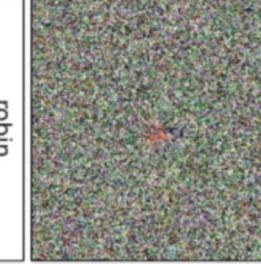
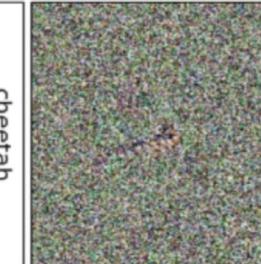
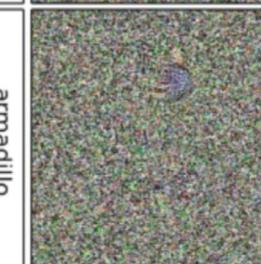
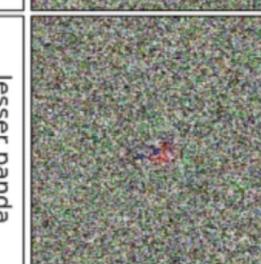
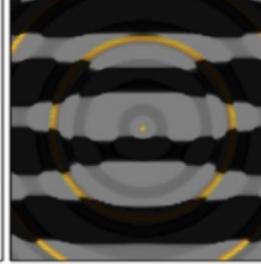
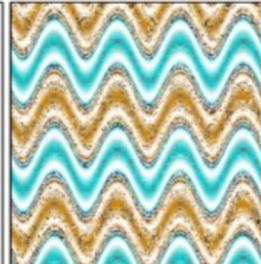
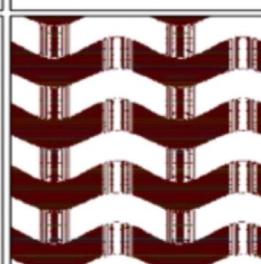
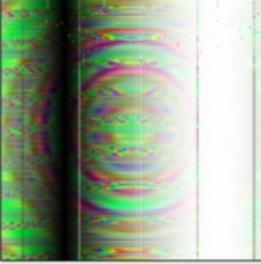
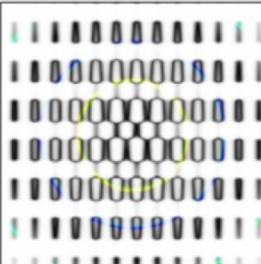
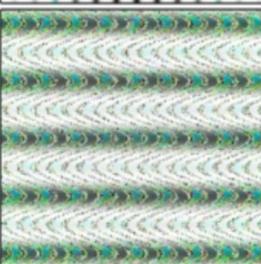
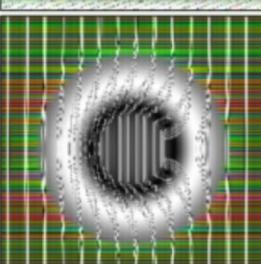
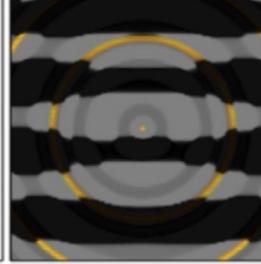
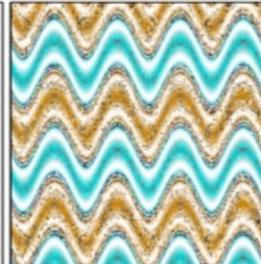
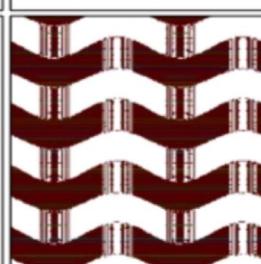
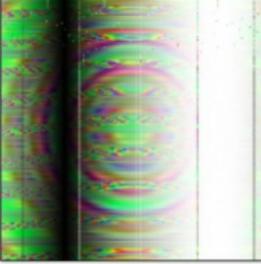
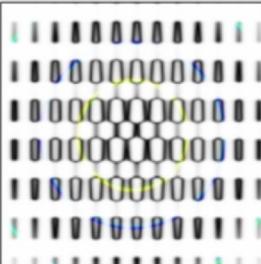
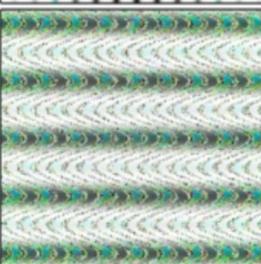
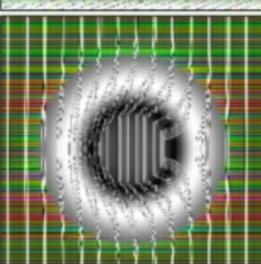
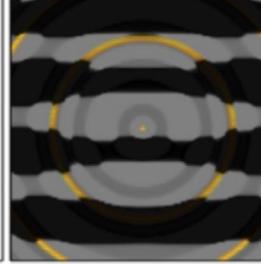
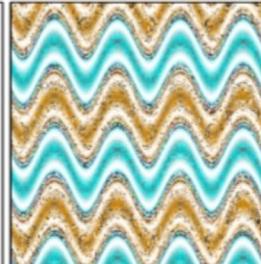
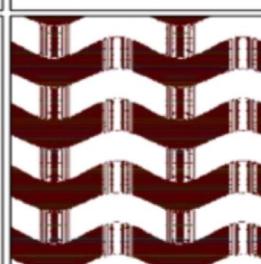
[Wang et al. 2015; Srivastava et al 2015; ...]

Context



[Doersch et al. 2014; Pathak et al. 2015; Isola et al. 2015]

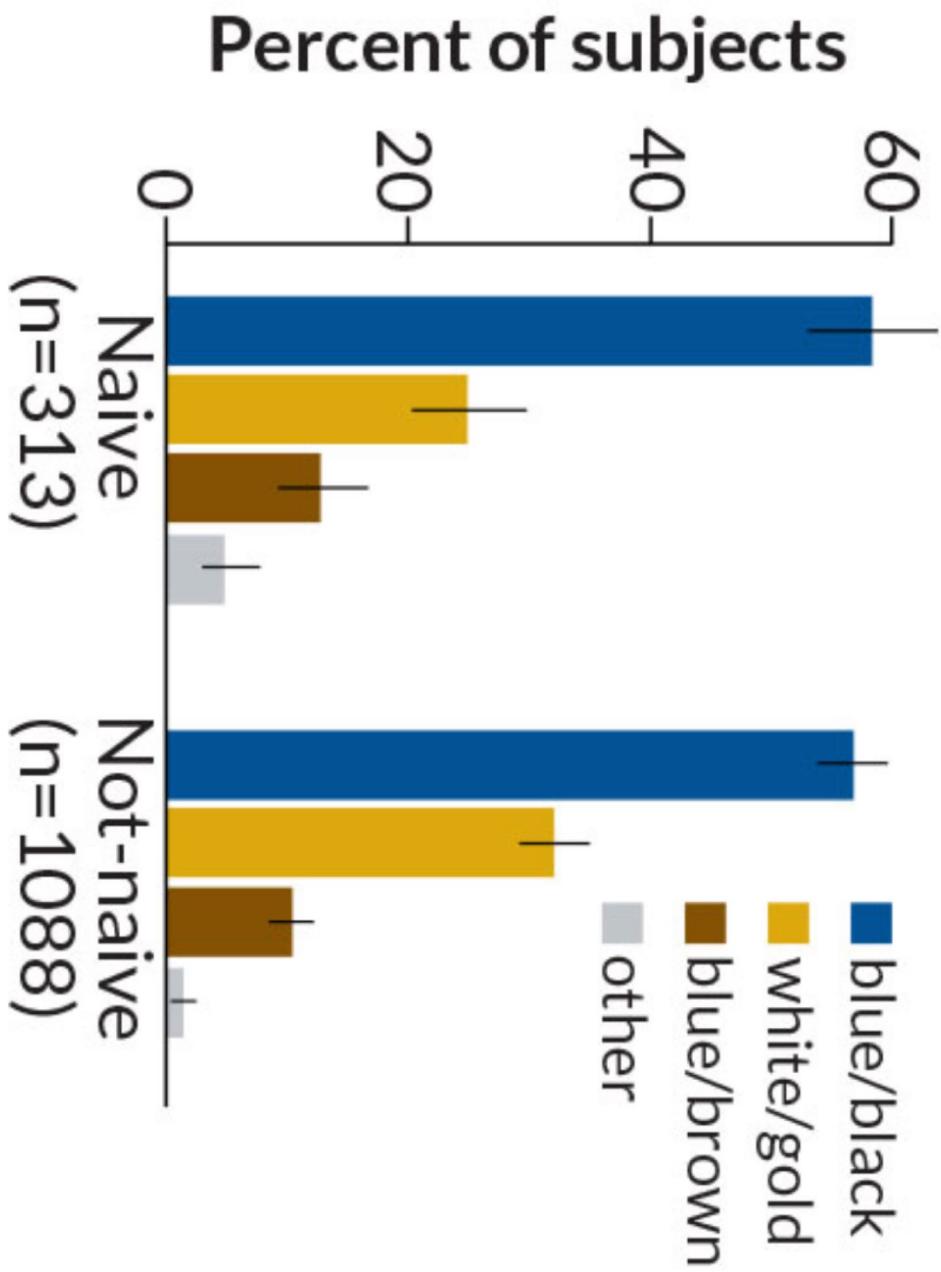
What are CNNs learning?

freight car							
remote control							
peacock							
African grey							
robin							
cheetah							
armadillo							
lesser panda							
centipede							
king penguin							
starfish							
baseball							
electric guitar							

[Nguyen et al]

The Dress divided the Internet, but it's really about subtraction

BY RACHEL EHRENBERG 12:49PM, MAY 14, 2015



CNN mysteries

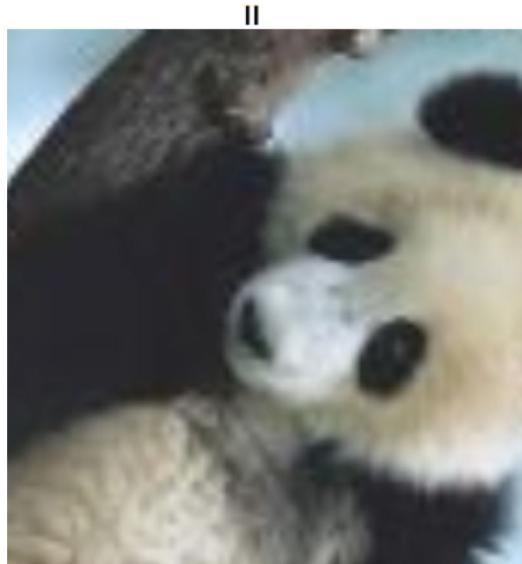
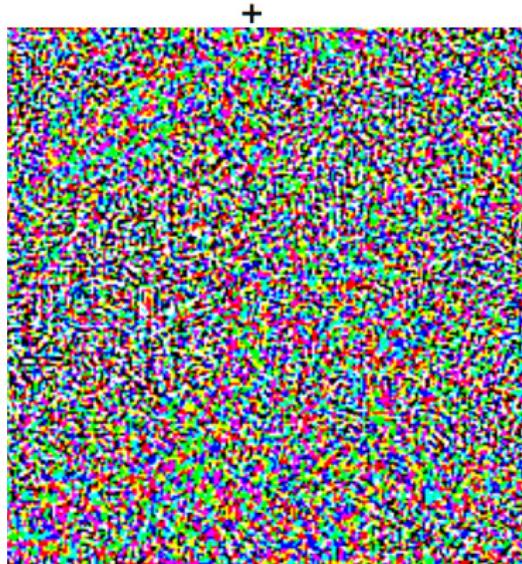
CNNs can sometimes outperform humans. Even after adding extreme noise, all these digits are recognized correctly!



But, only first column of these are
recognized correctly. Second
column digits are all wrong!

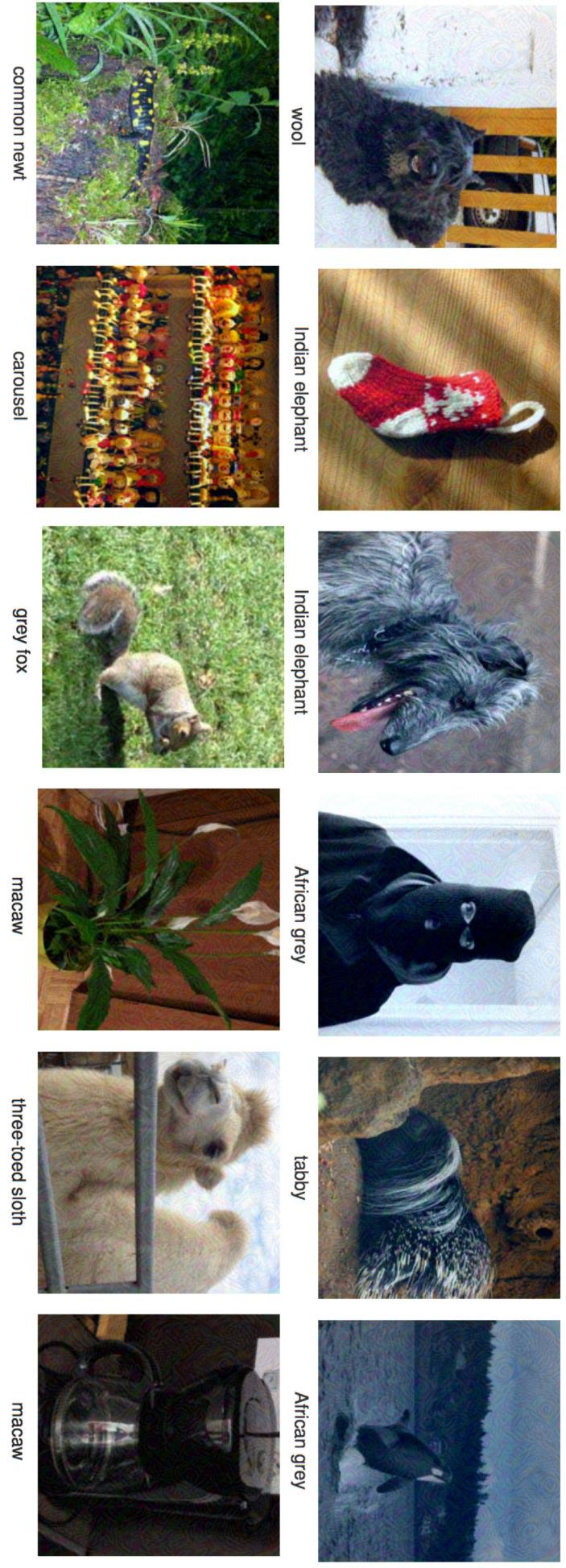
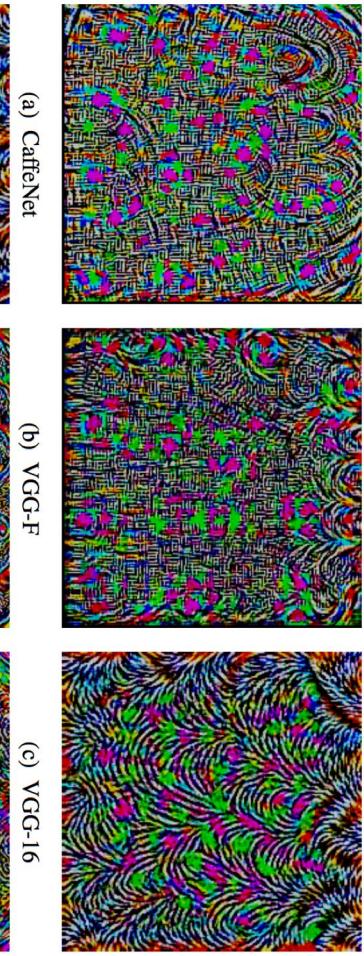


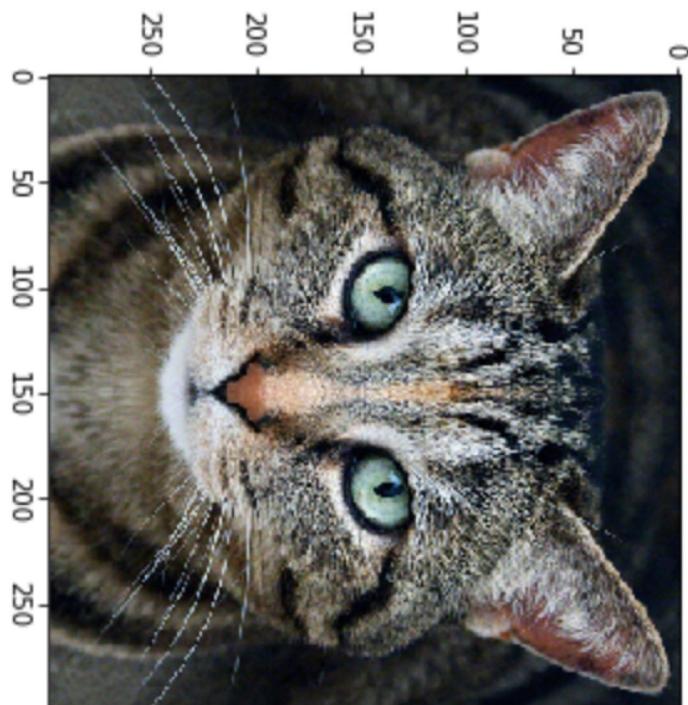
Evidence of overfitting?



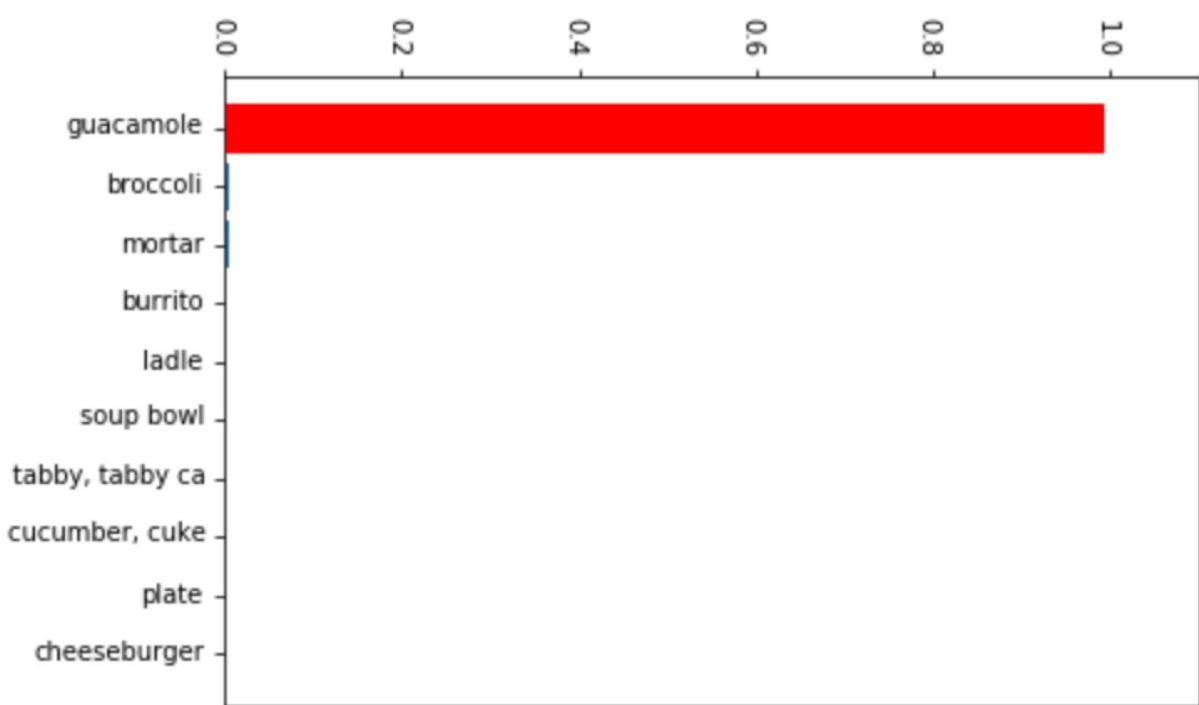
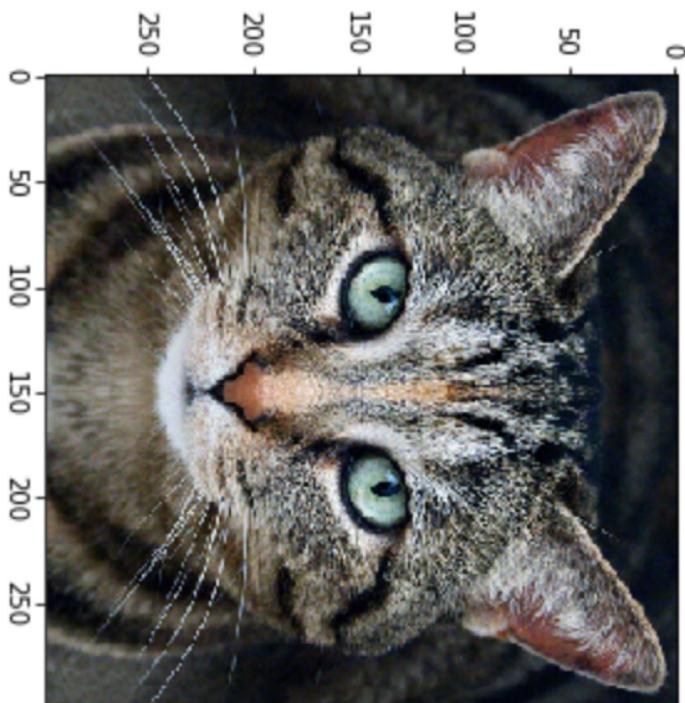
Goodfellow et al

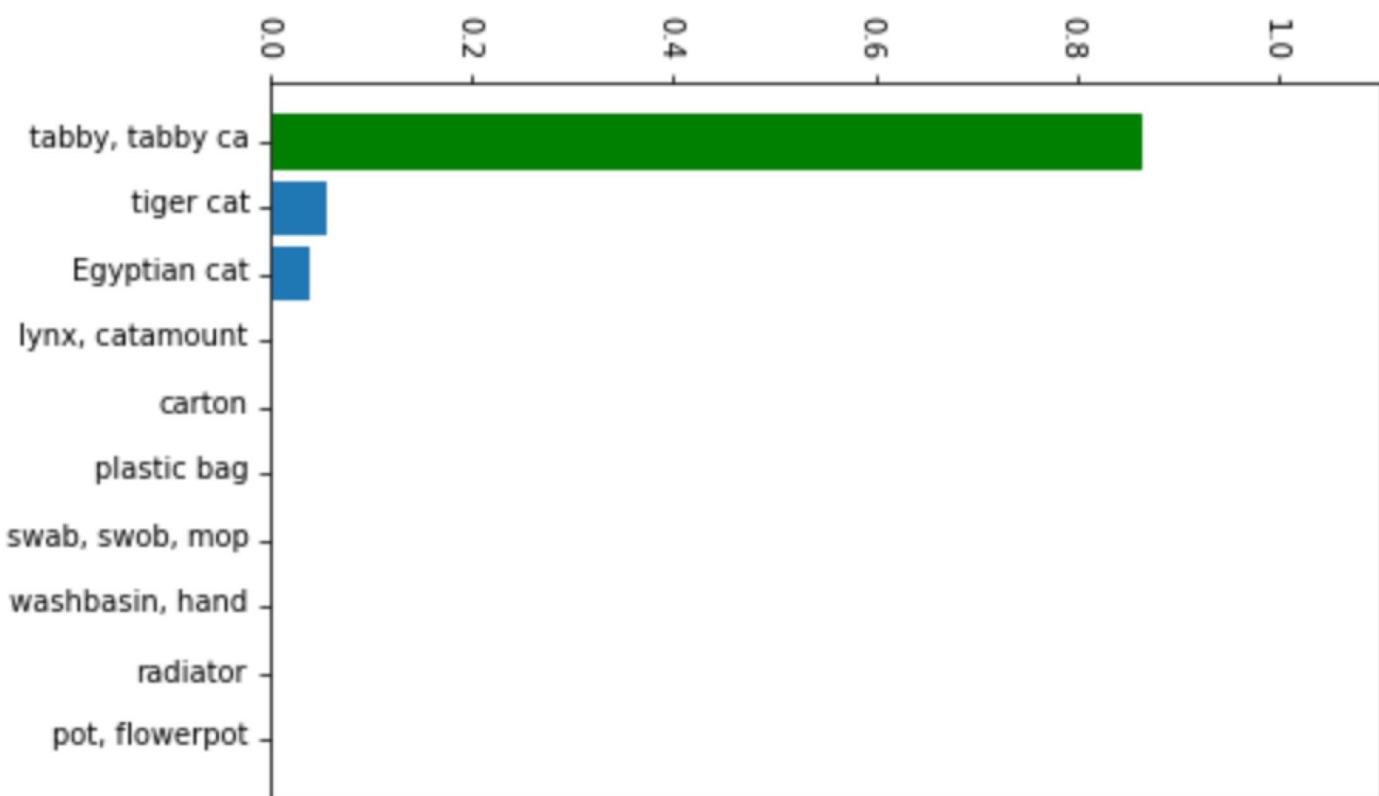
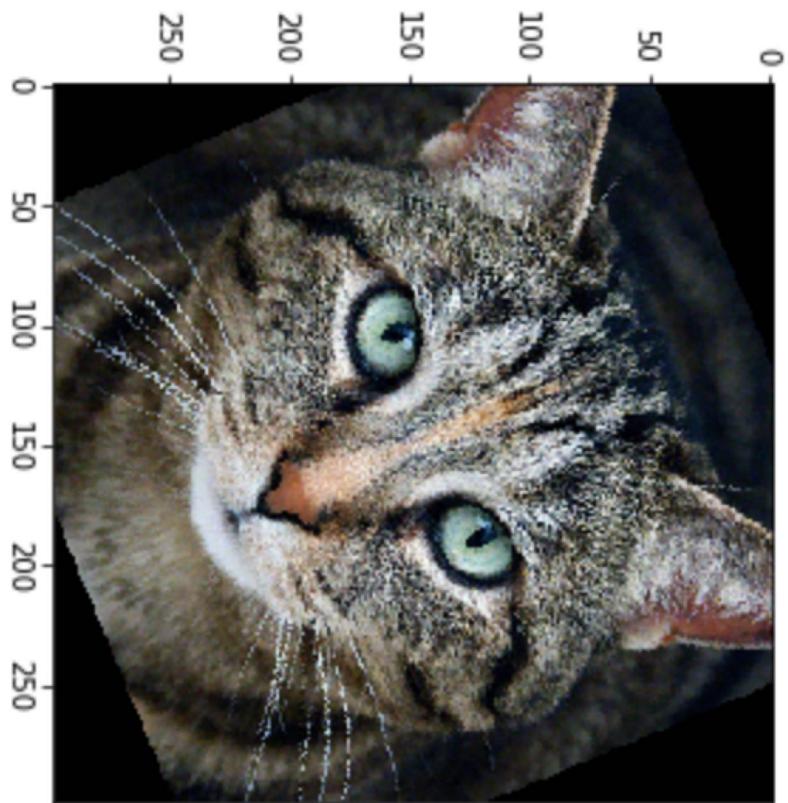
Moosavi-Dezfooli et al,
Universal adversarial
perturbations, CVPR 2017.





Elsayed 2018





(a) Image from dataset



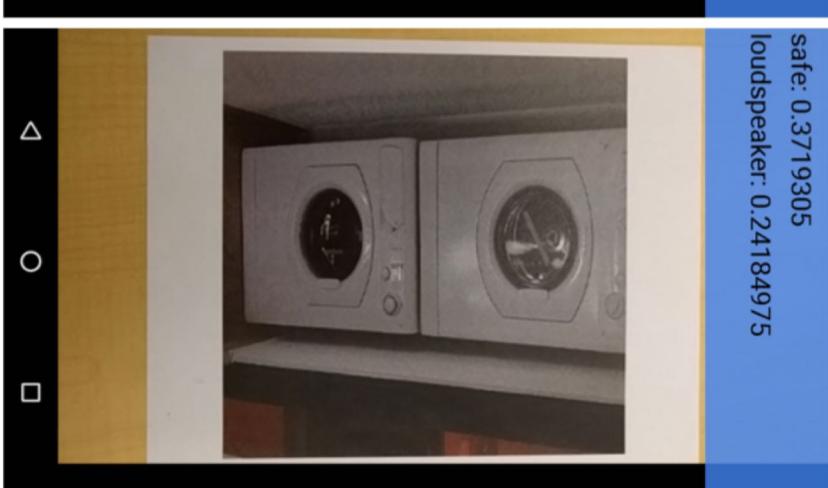
(b) Clean image



(c) Adv. image, $\epsilon = 4$



(d) Adv. image, $\epsilon = 8$



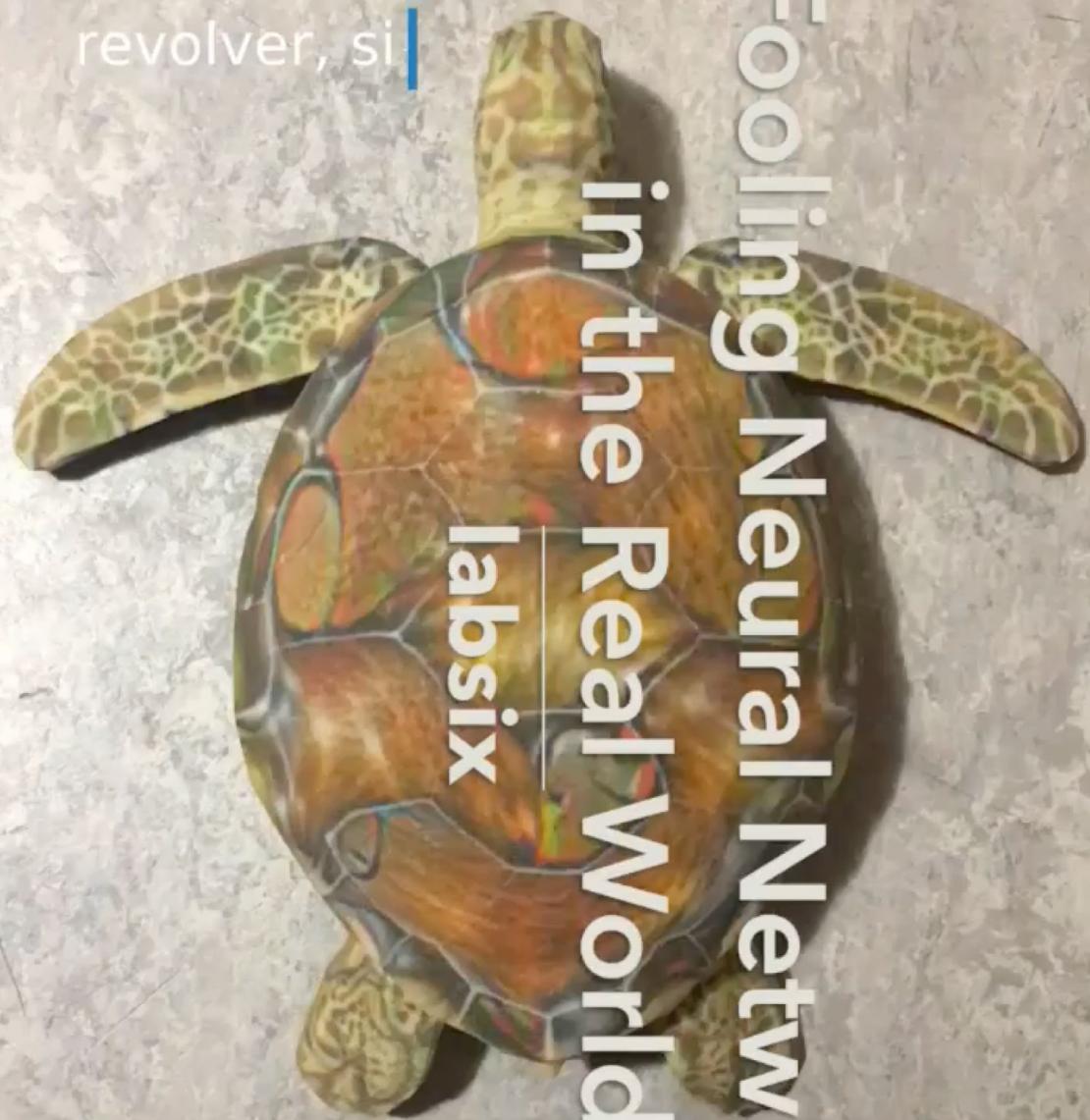
Distance/Angle	Subtle Poster	Subtle Poster Right Turn	Camouflage Graffiti	Camouflage Art (LISA-CNN)	Camouflage Art (GTSRB-CNN)
5' 0°					
5' 15°					
10' 0°					
10' 30°					
40' 0°					
Targeted-Attack Success	100%	73.33%	66.67%	100%	80%

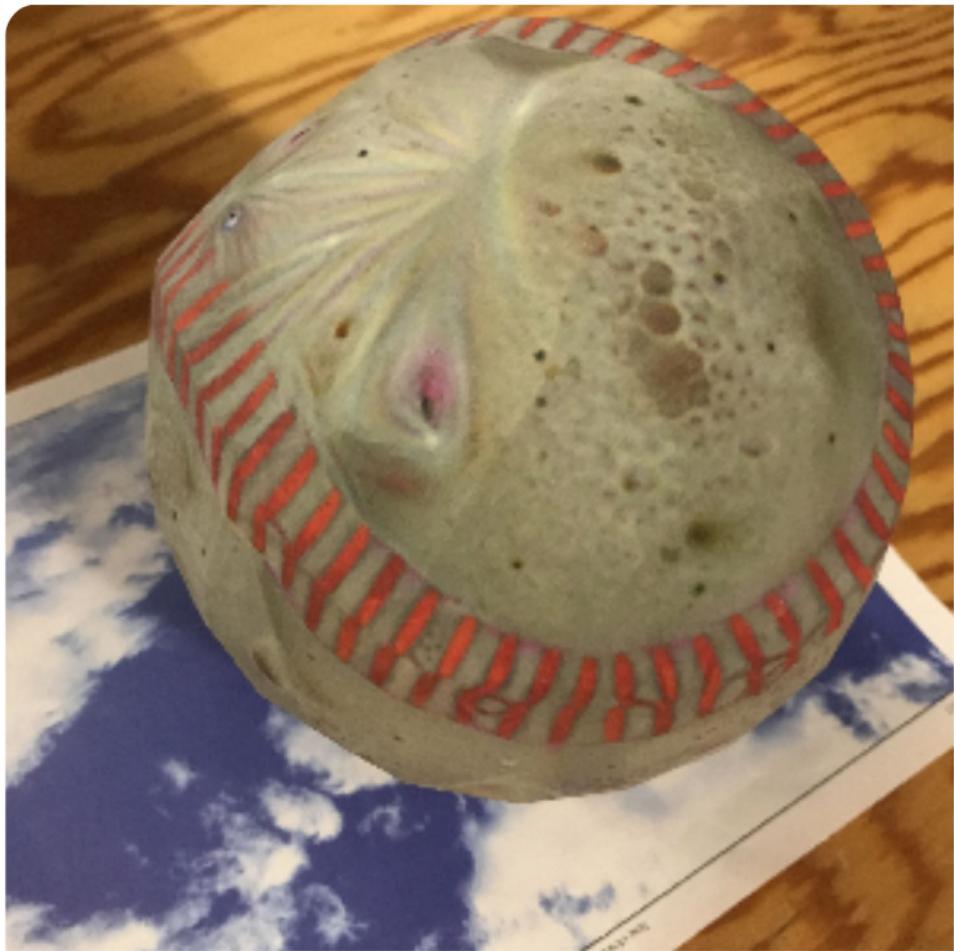
Eykholt 2018

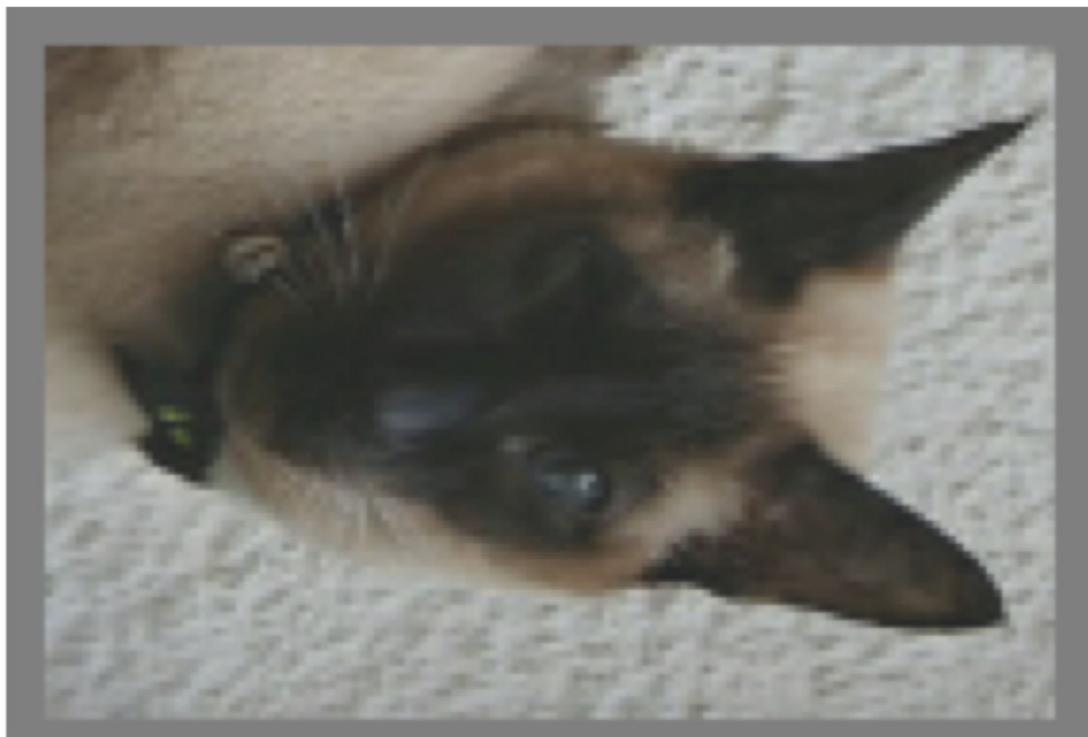
rifle
shield, buck
revolver, si

Fooling Neural Networks in the Real World

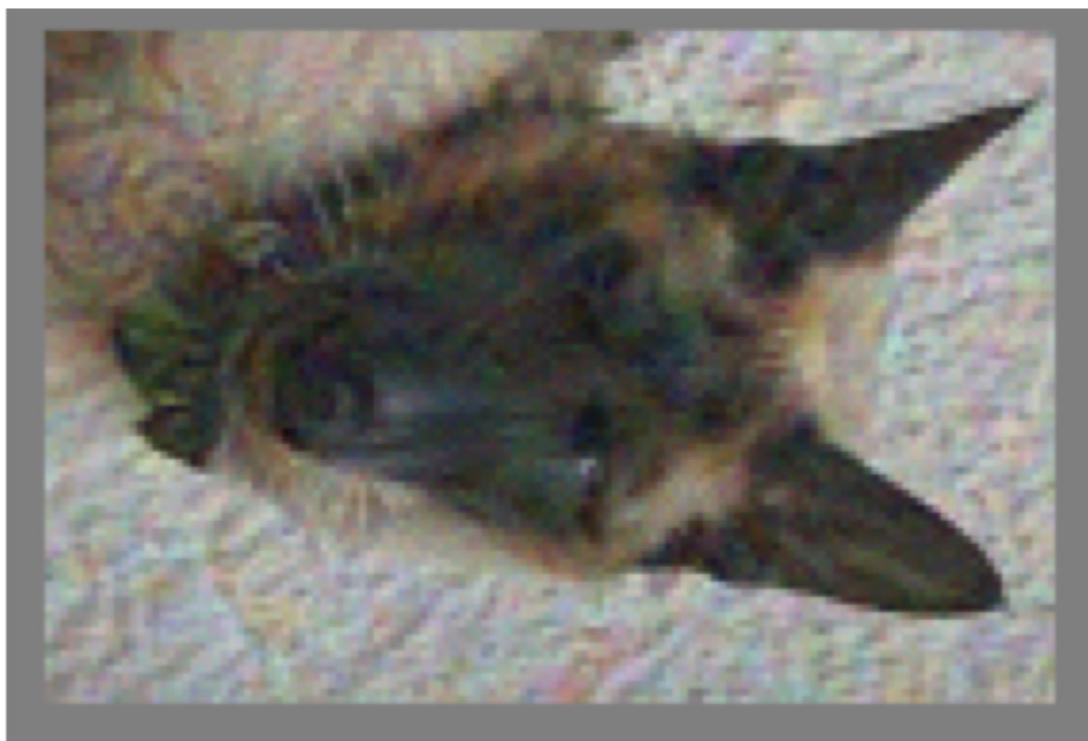
labsix



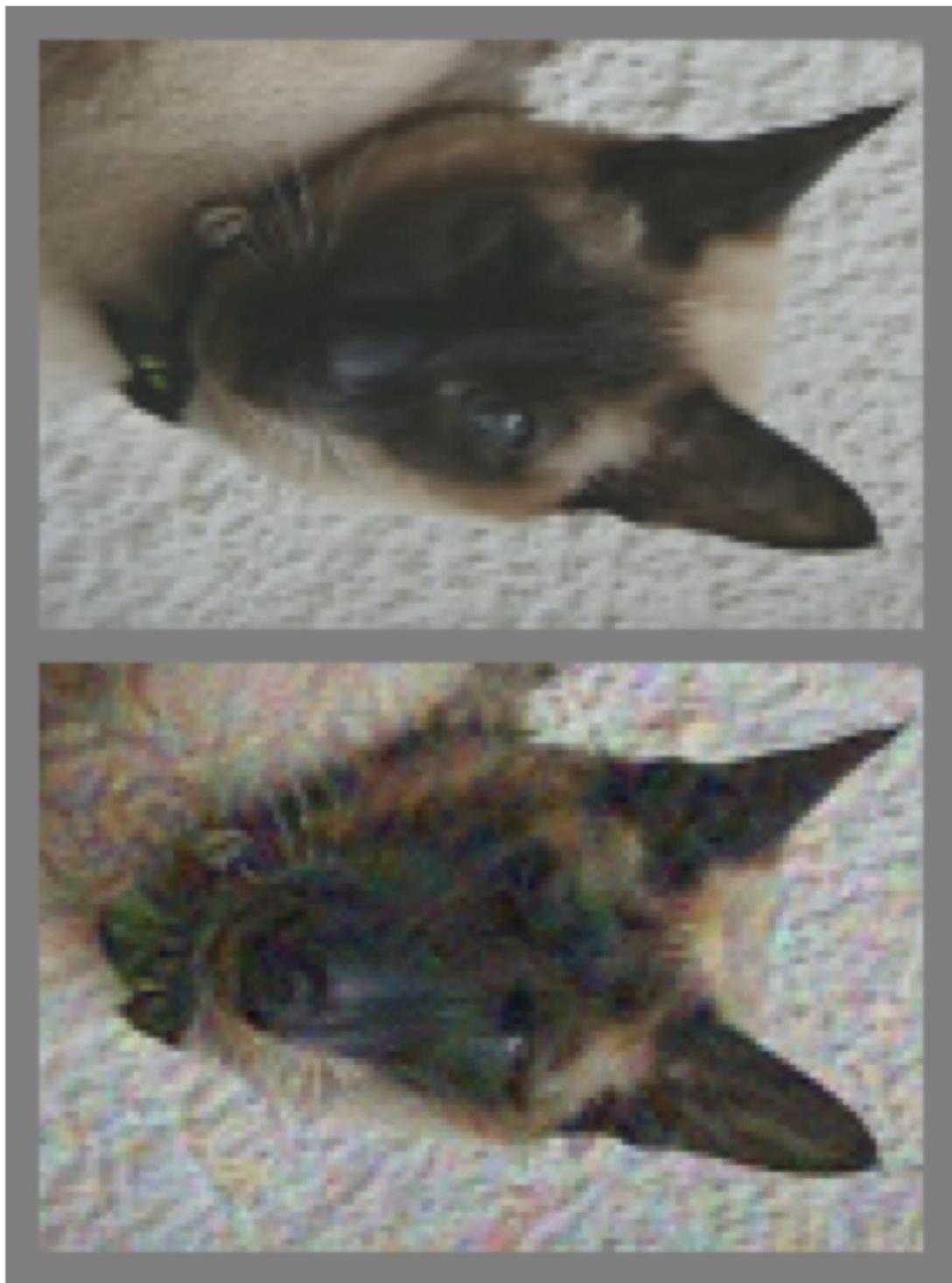




Elsayed 2018



Elsayed 2018



Elsayed 2018

Visualizing Convnets

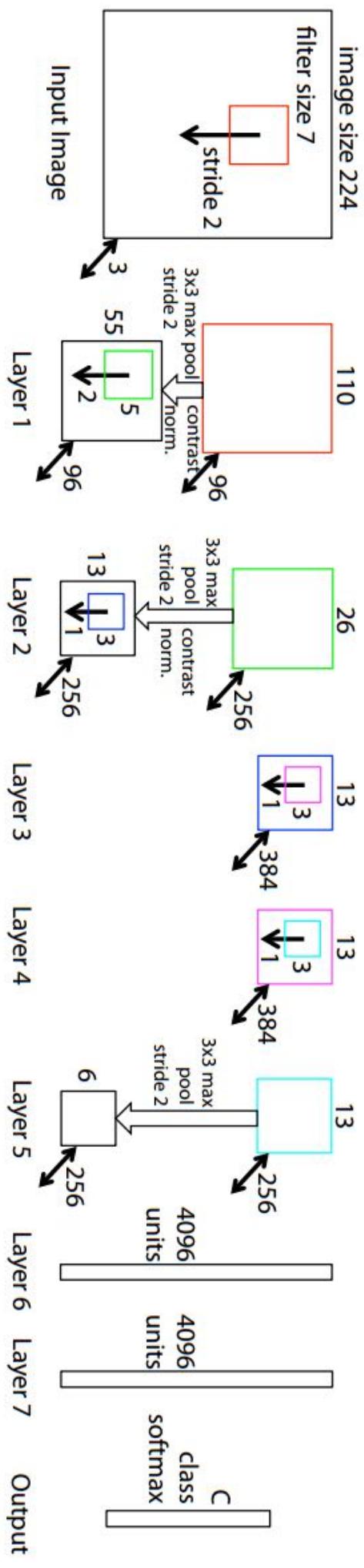
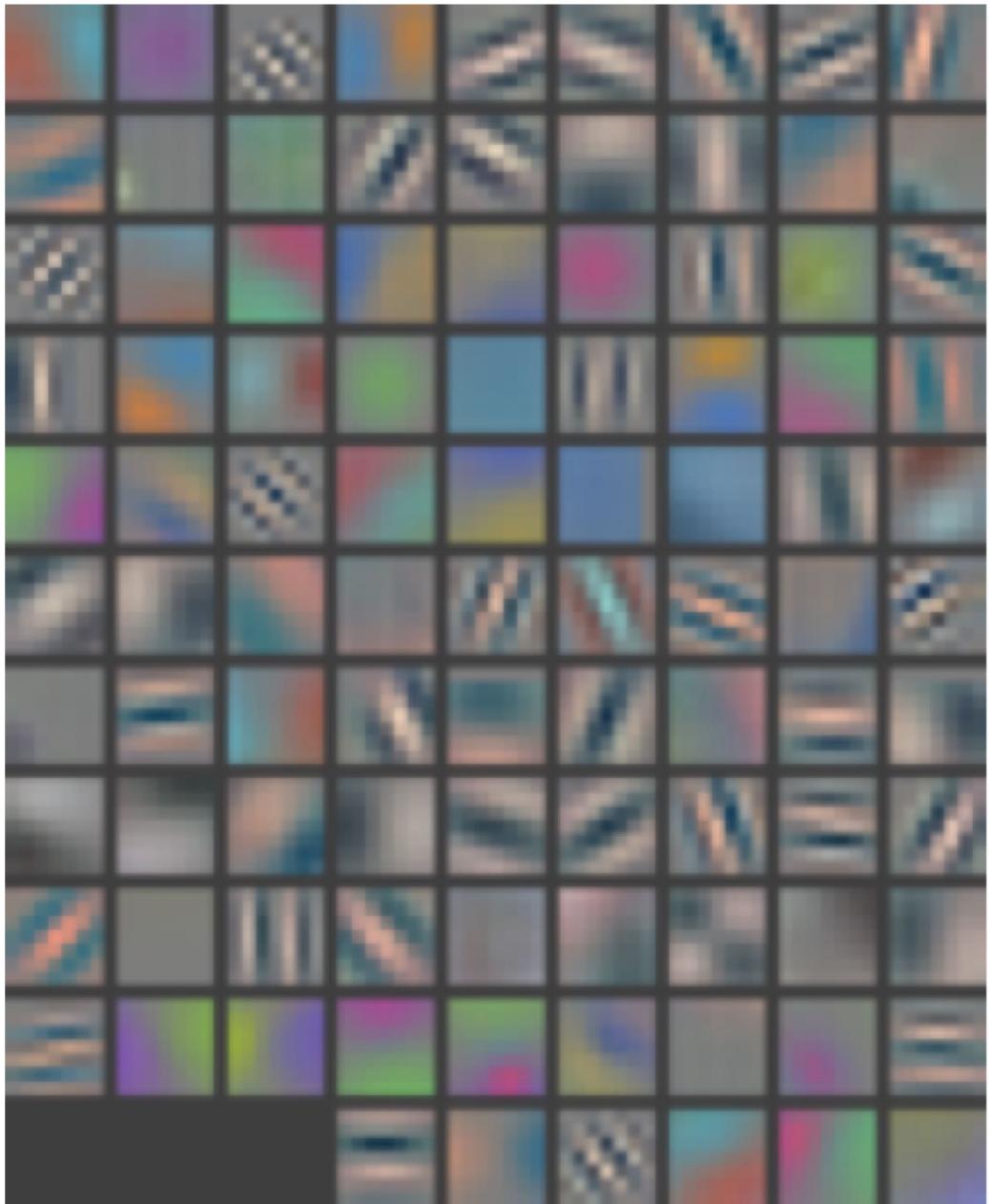


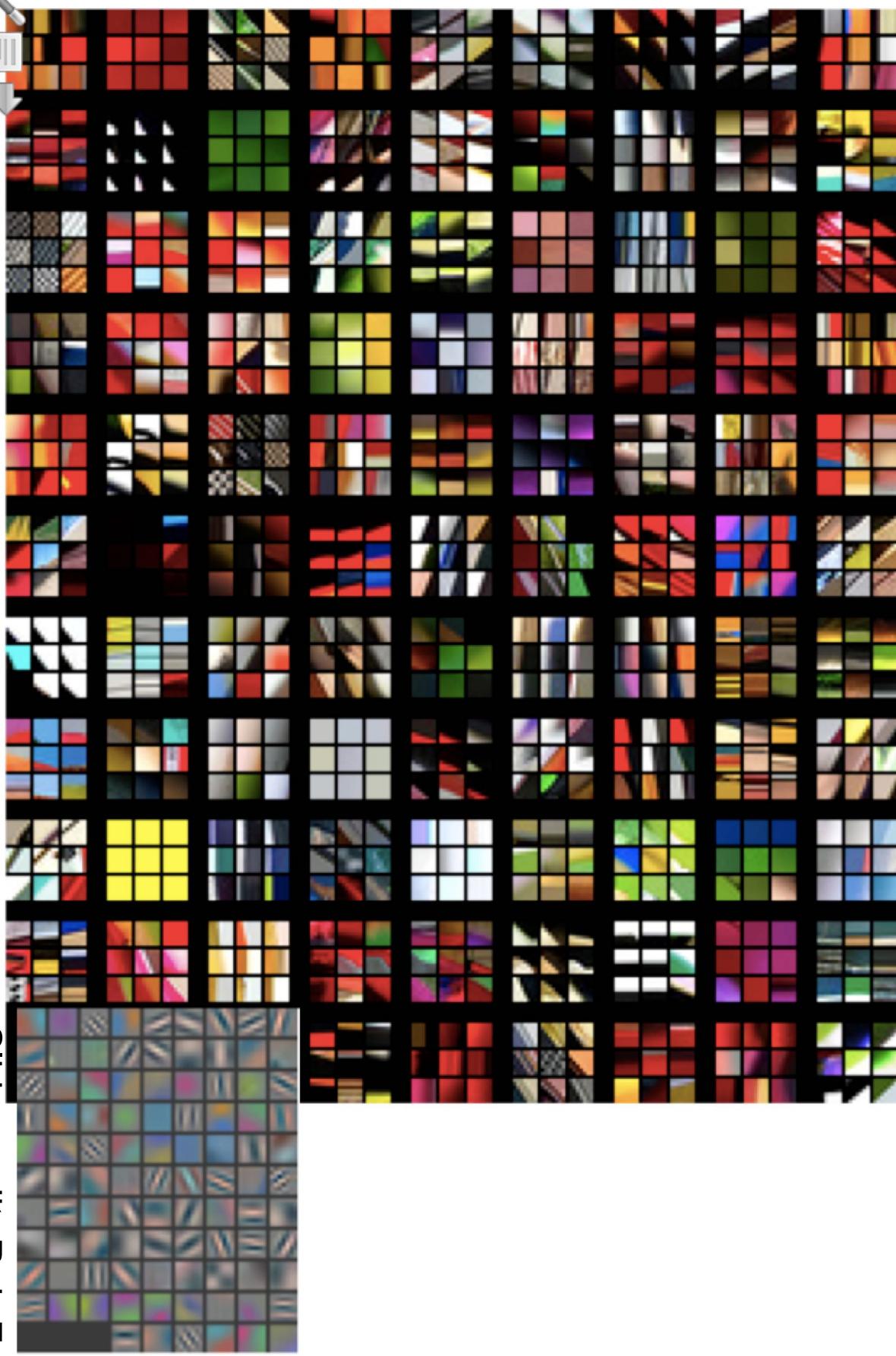
Figure 3. Architecture of our 8 layer convnet model. A 224 by 224 crop of an image (with 3 color planes) is presented as the input. This is convolved with 96 different 1st layer filters (red), each of size 7 by 7, using a stride of 2 in both x and y. The resulting feature maps are then: (i) passed through a rectified linear function (not shown), (ii) pooled (max within 3x3 regions, using stride 2) and (iii) contrast normalized across feature maps to give 96 different 55 by 55 element feature maps. Similar operations are repeated in layers 2,3,4,5. The last two layers are fully connected, taking features from the top convolutional layer as input in vector form ($6 \cdot 6 \cdot 256 = 9216$ dimensions). The final layer is a C -way softmax function, C being the number of classes. All filters and feature maps are square in shape.

Layer 1 Filters



Slide credit: Rob Fergus

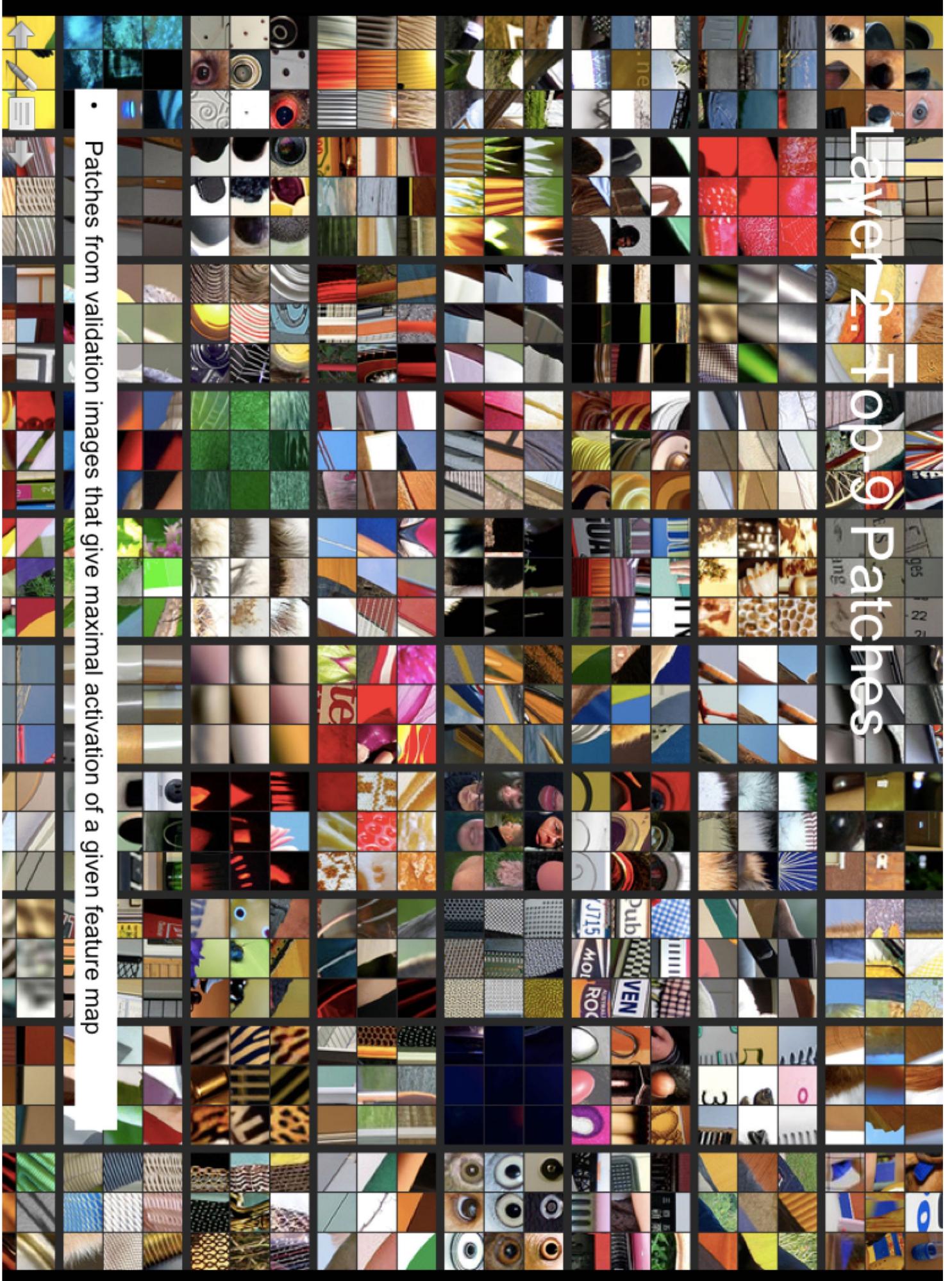
Layer 1: Top-9 Patches

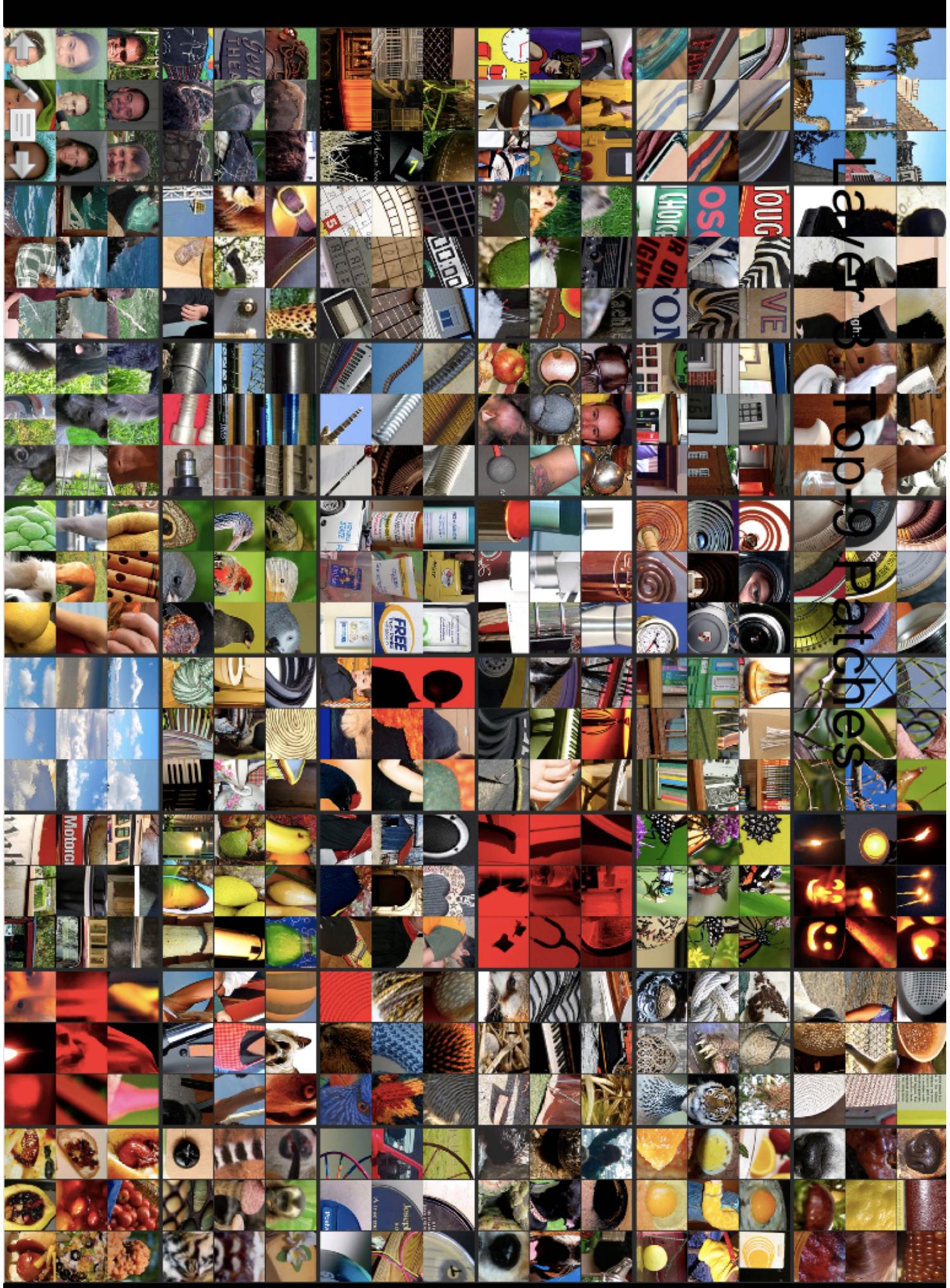


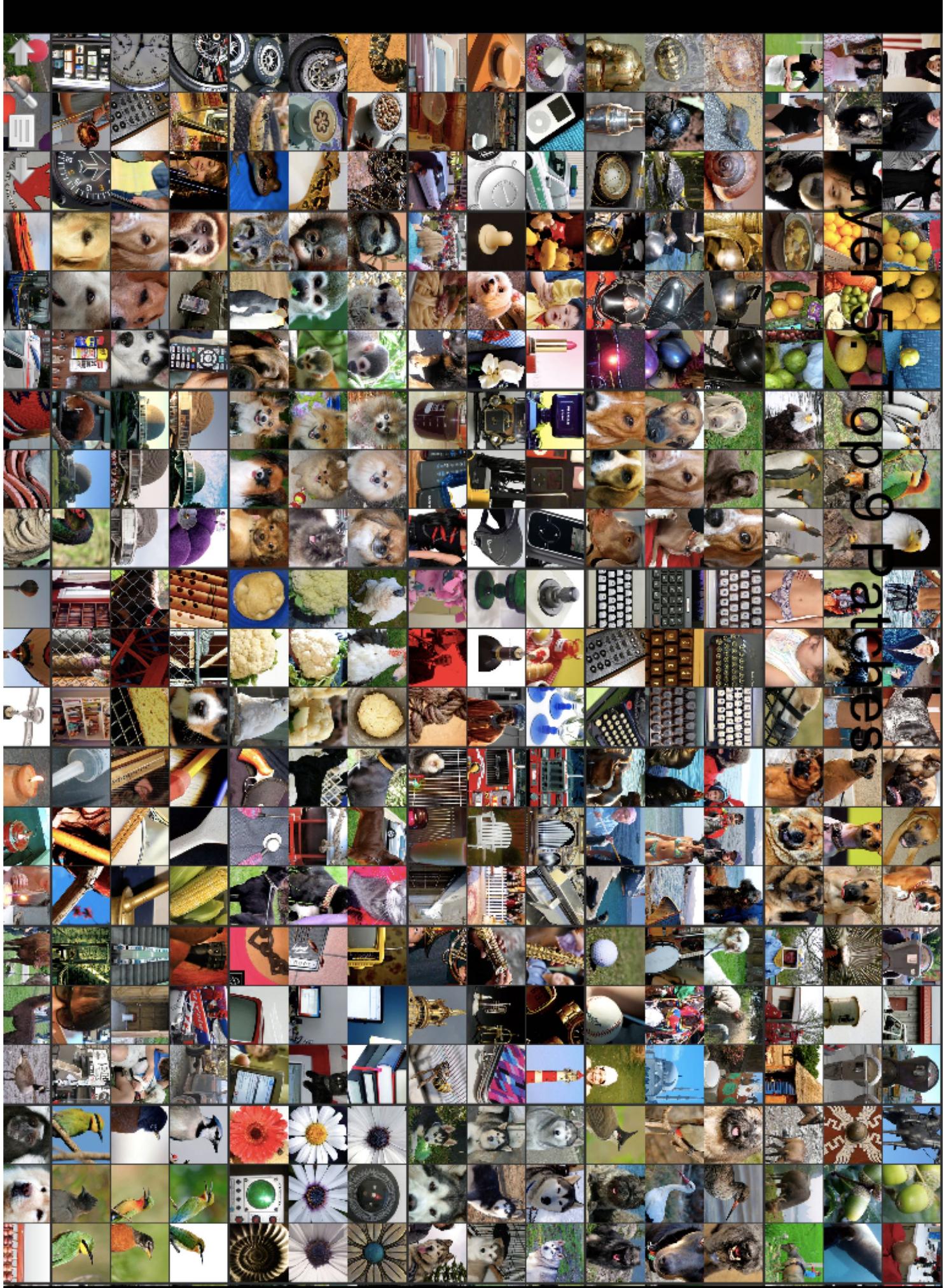
Slide credit: Rob Fergus

Layer 2... Top-9 Patches

- Patches from validation images that give maximal activation of a given feature map







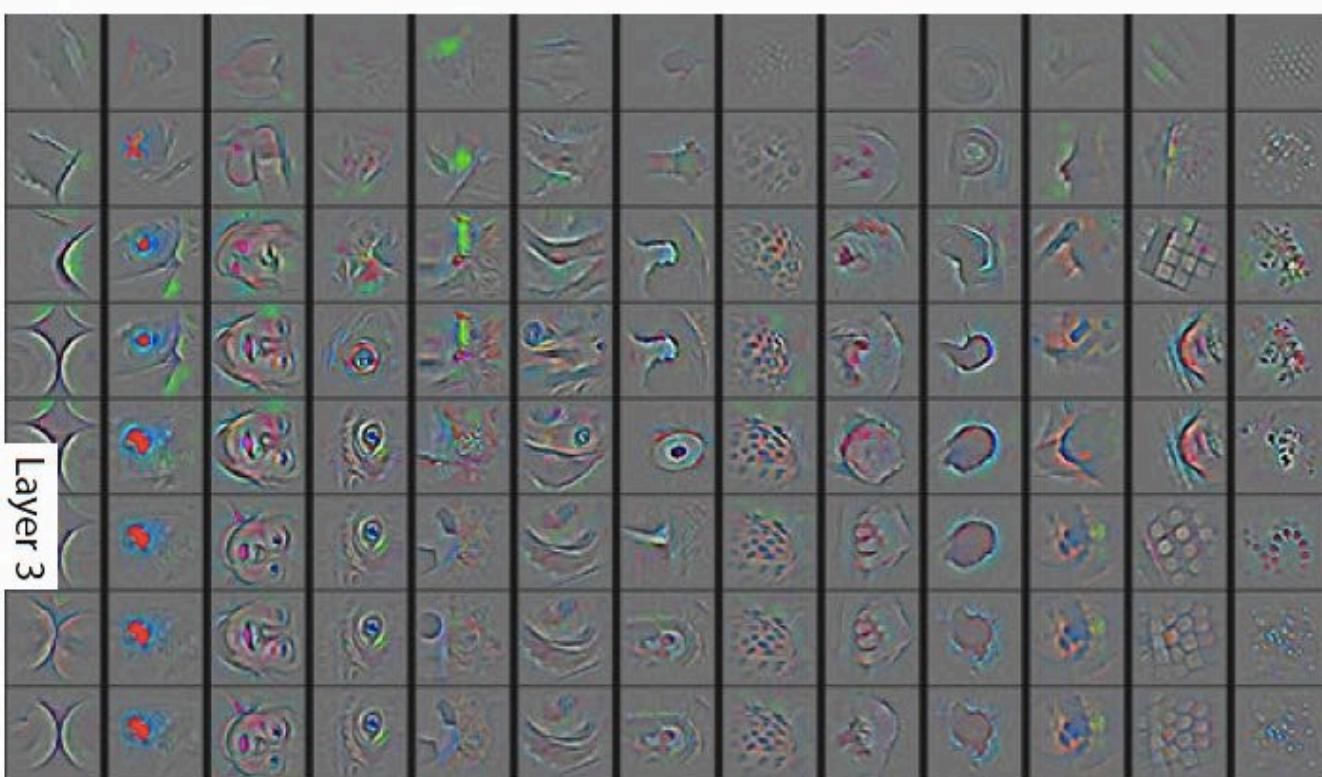
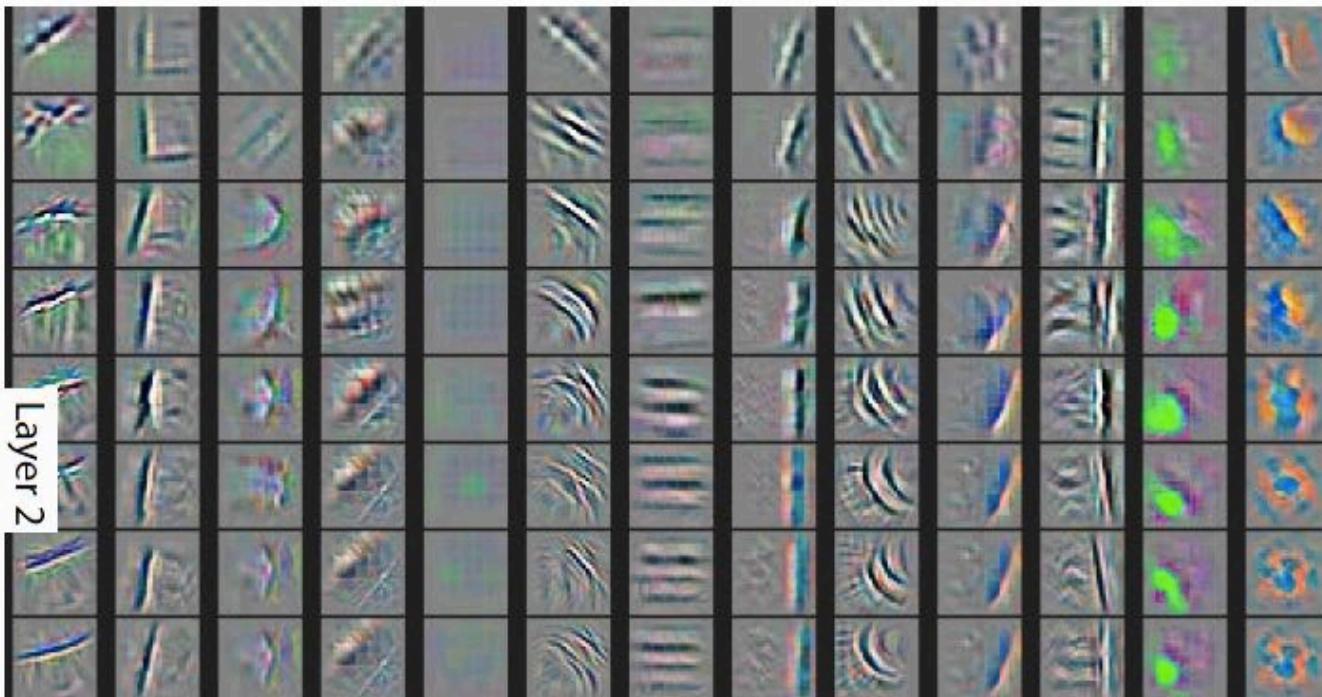
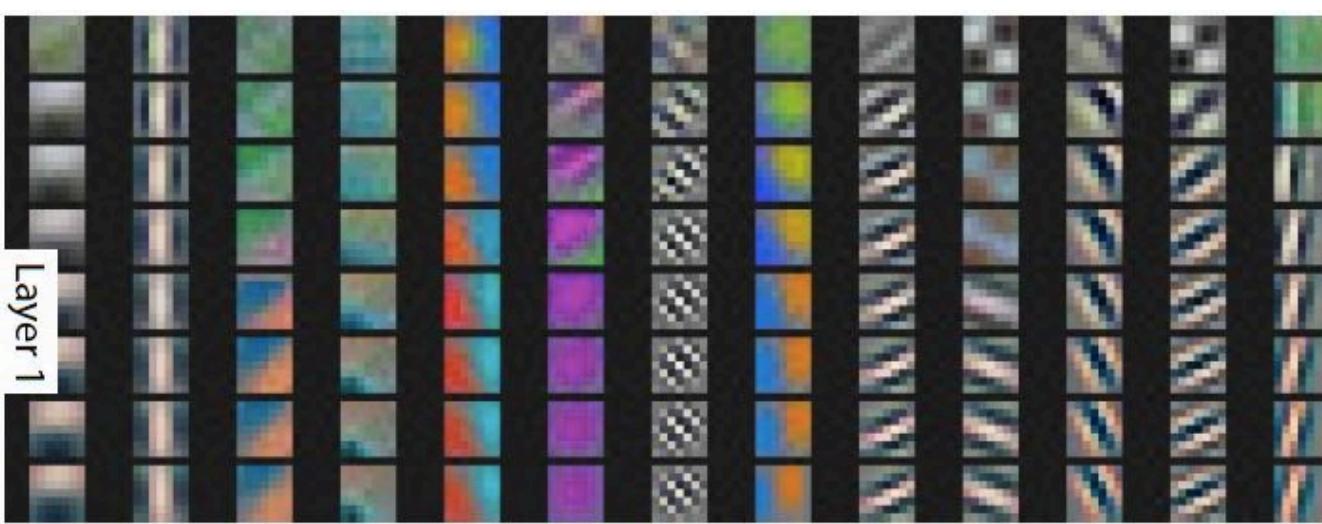
...
y
e
r
5
o
p
9
P
a
t
h
e
s

Evolution of Features During Training

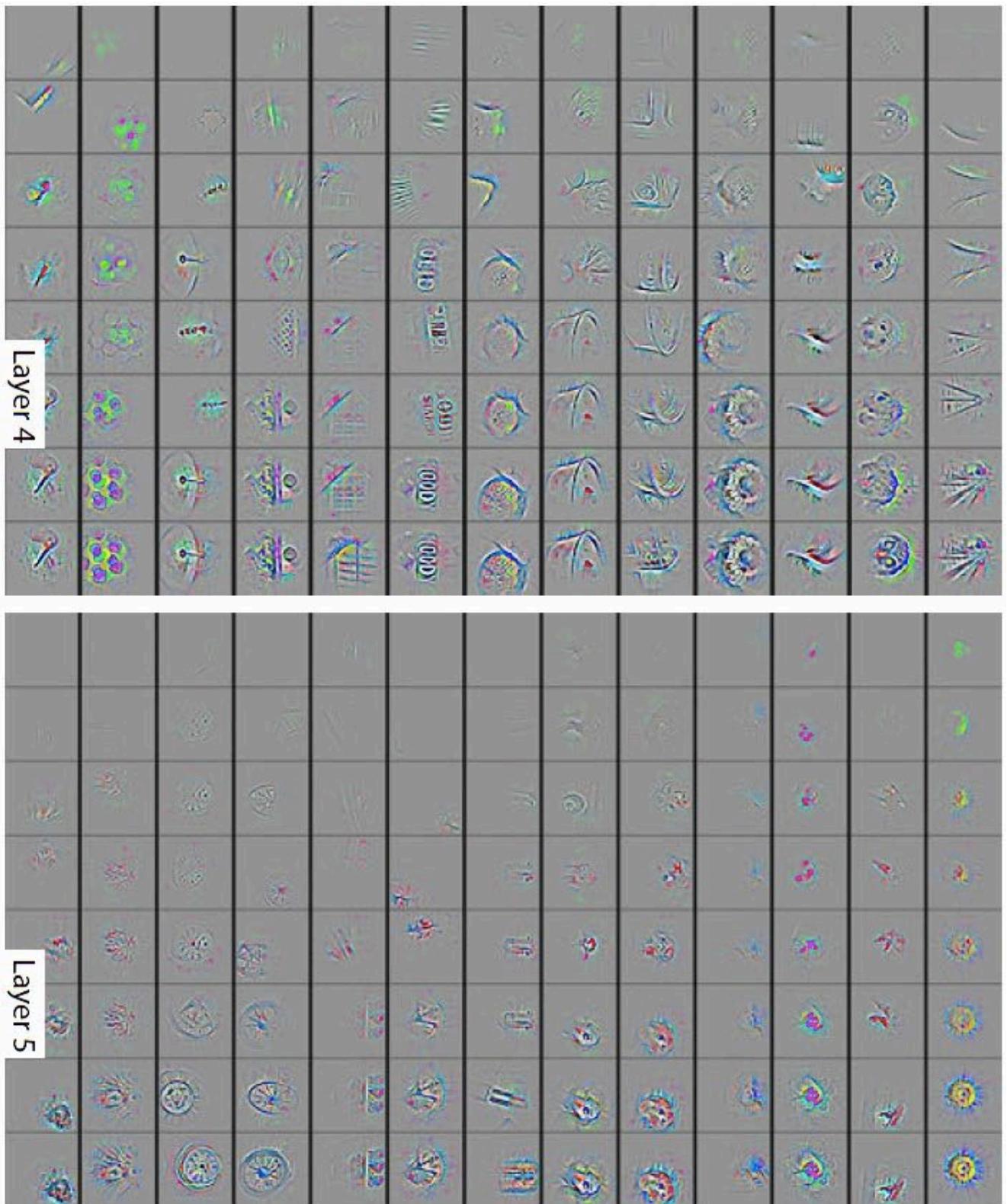
Layer 1

Layer 2

Layer 3

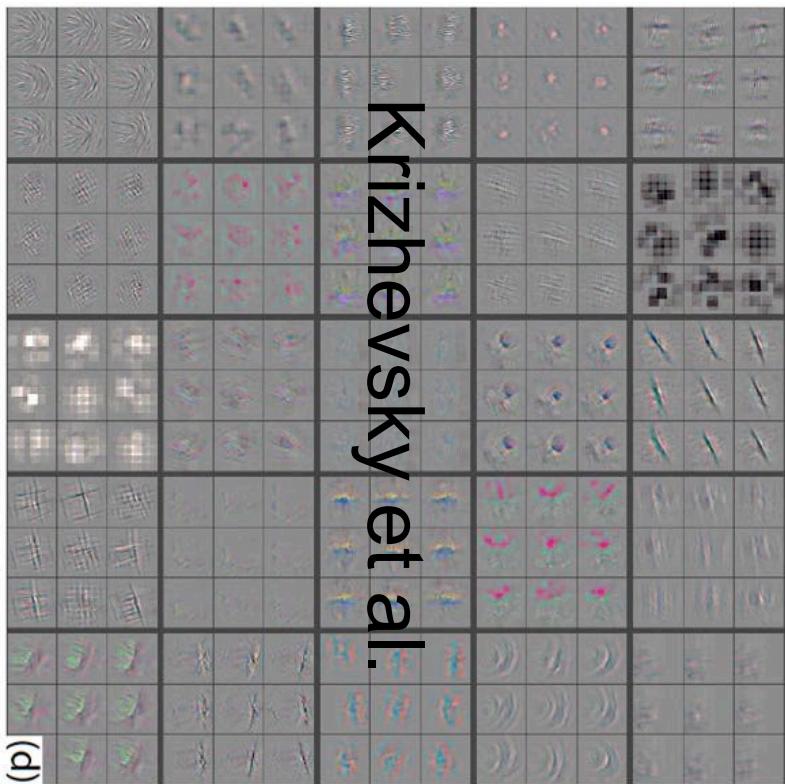


Evolution of Features During Training

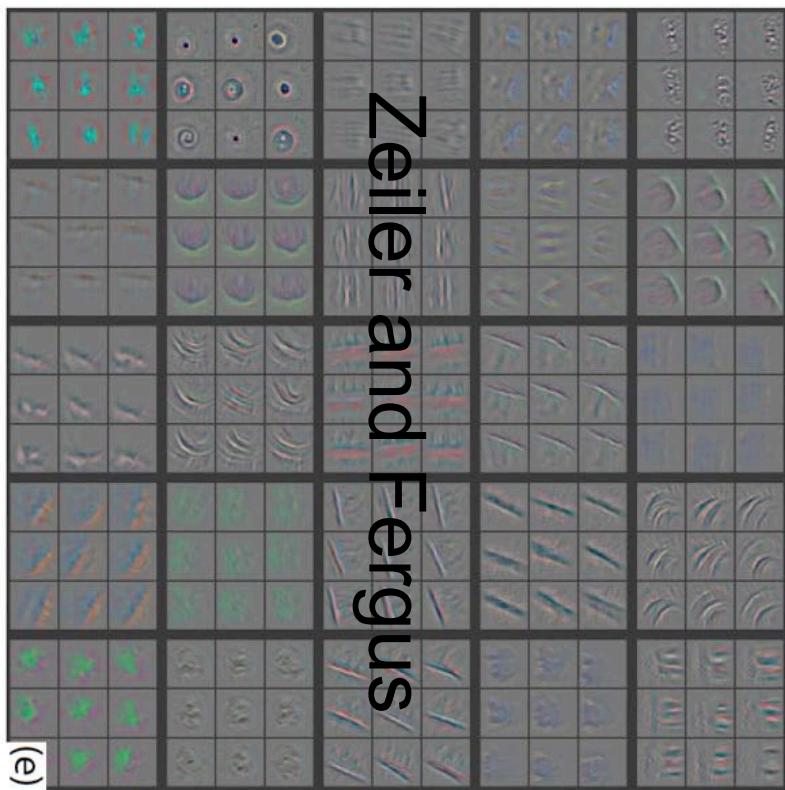


Diagnosing Problems

- Visualization of Krizhevsky et al.'s architecture showed some problems with layers 1 and 2
- Large stride of 4 used
- Alter architecture: smaller stride & filter size
 - Visualizations look better
- Performance improves



(d)



(e)

Krizhevsky et al.

Zeiler and Fergus

Occlusion Experiment

- Mask parts of input with occluding square
- Monitor output



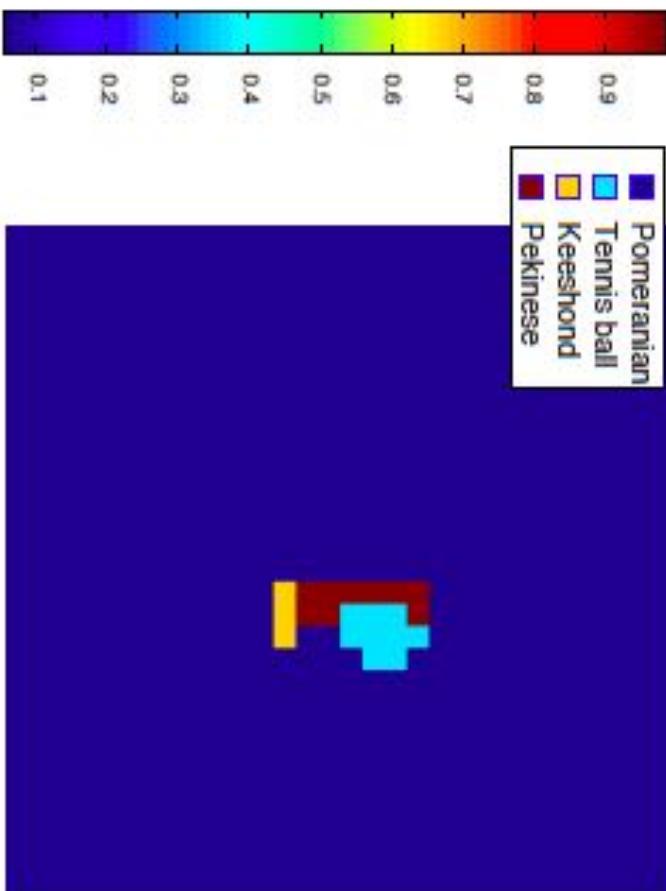
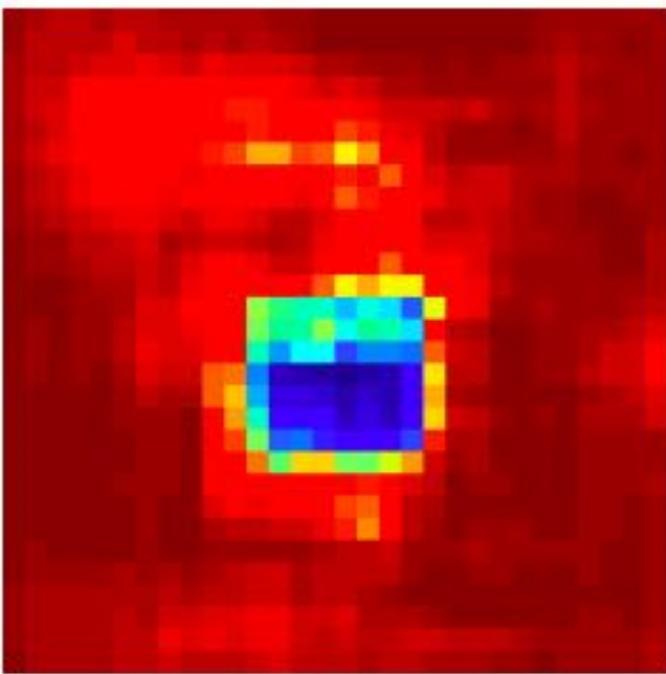
Input image



True Label: Pomeranian

$p(\text{True class})$

Most probable class

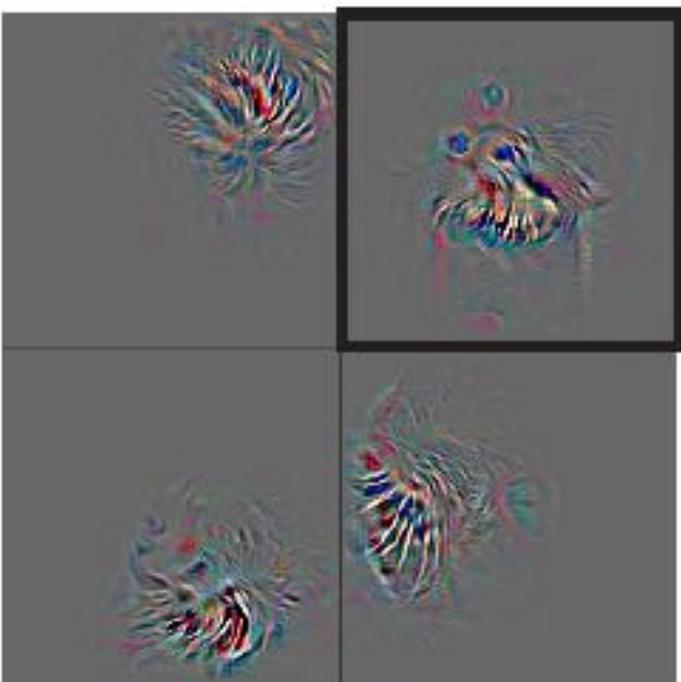
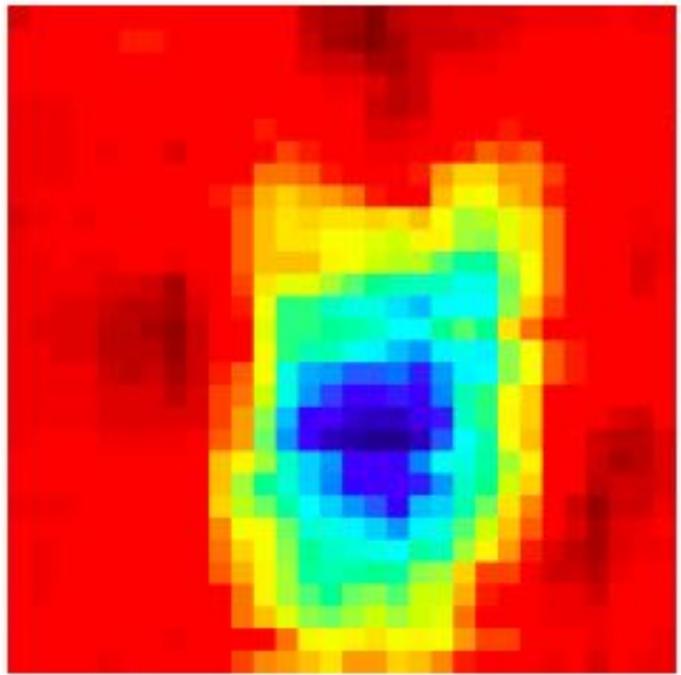


Input image



Total activation in most
active 5th layer feature map

Other activations from
same feature map



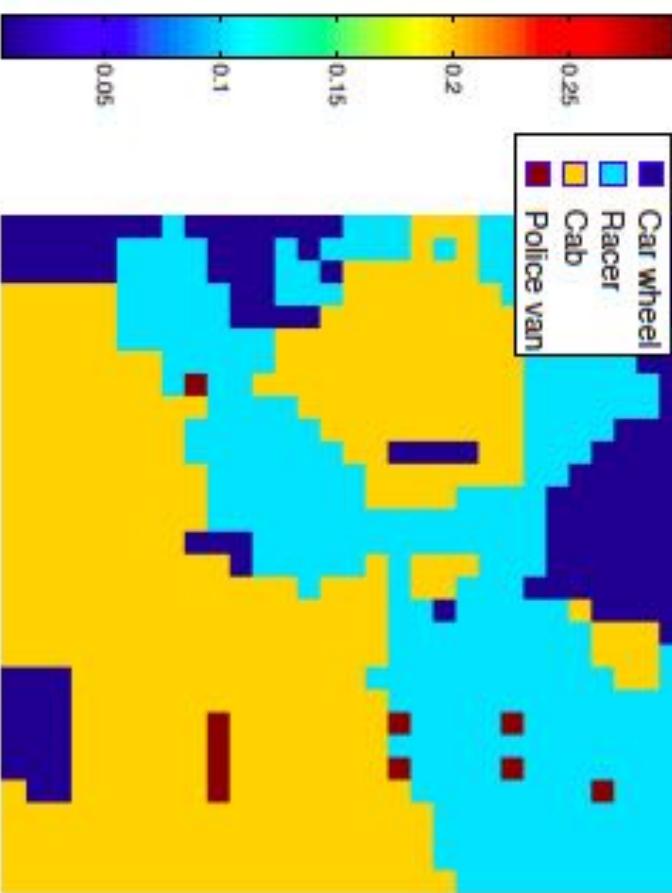
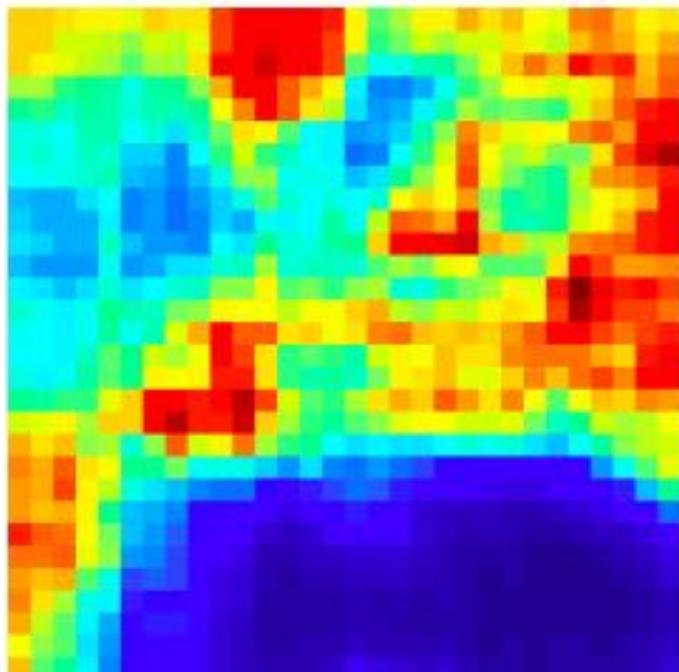
Input image



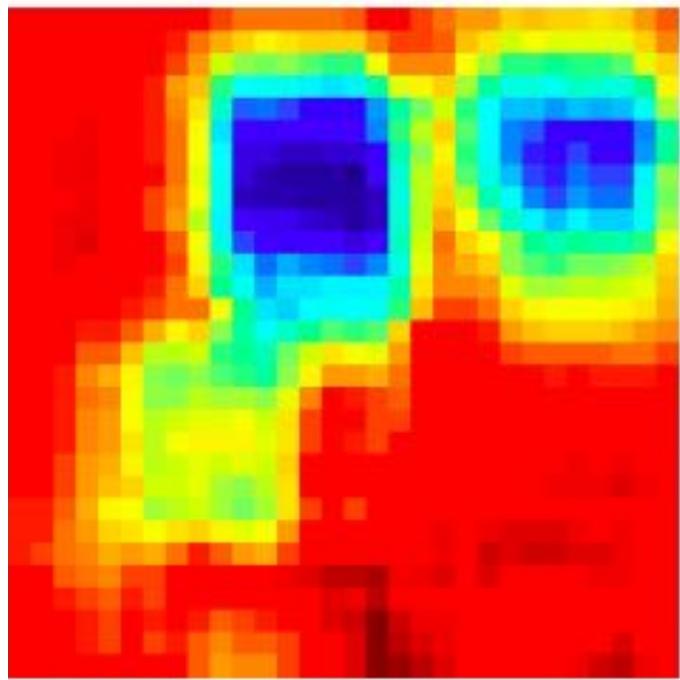
$p(\text{True class})$

Most probable class

- Car wheel
- Racer
- Cab
- Police van

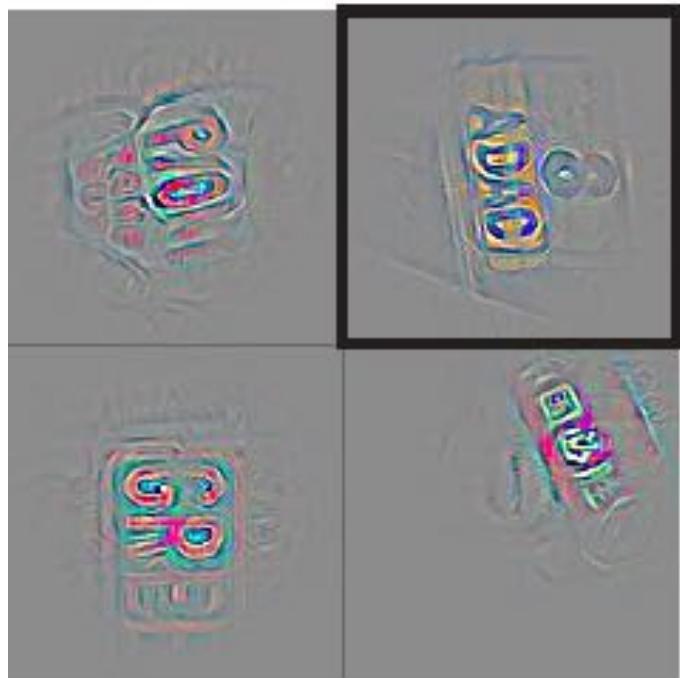


Input image

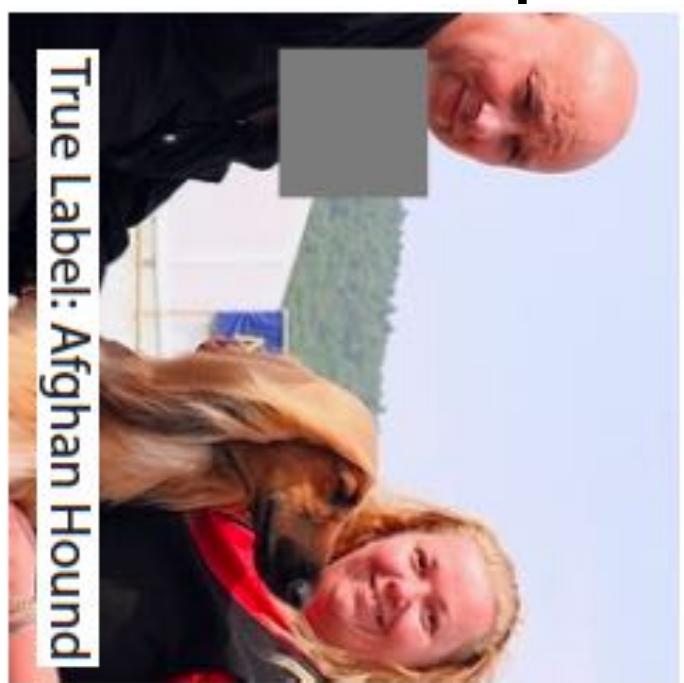


Total activation in most
active 5th layer feature map

Other activations from
same feature map

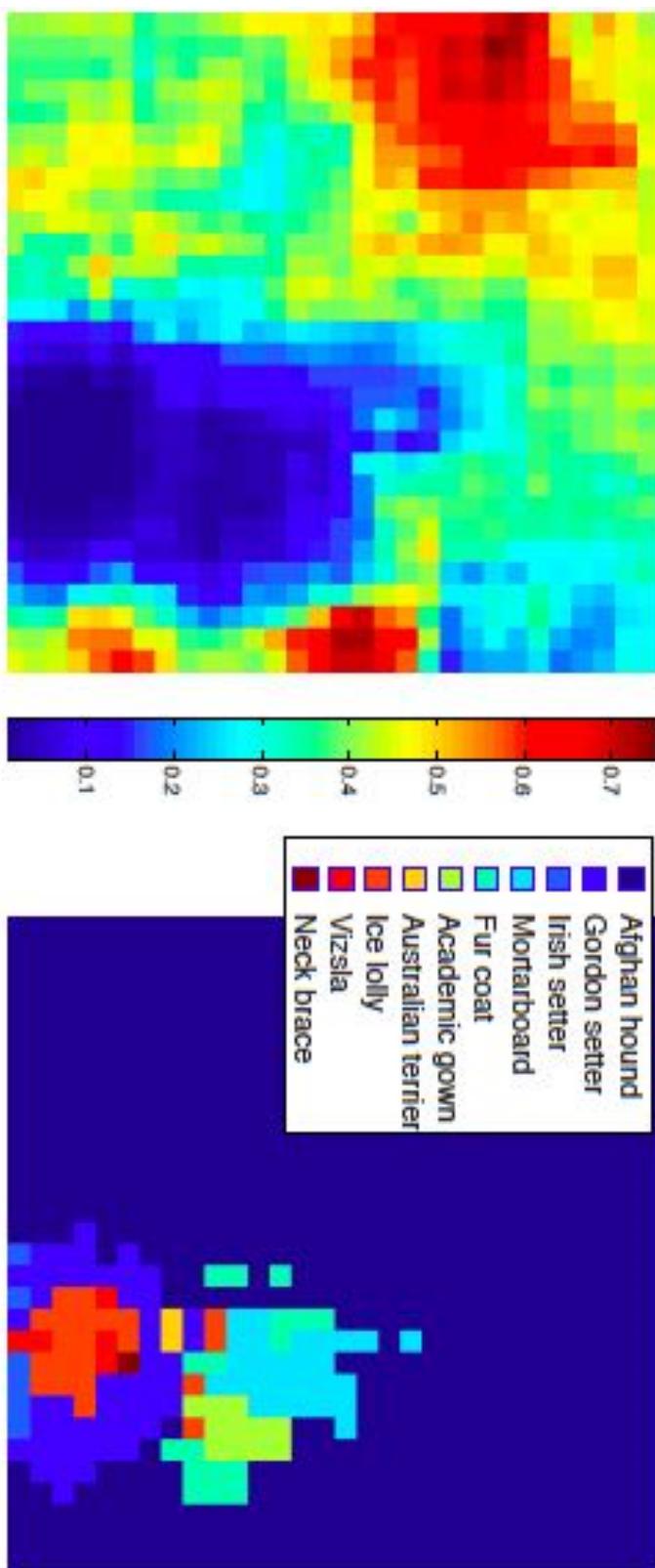


Input image



$p(\text{True class})$

Most probable class

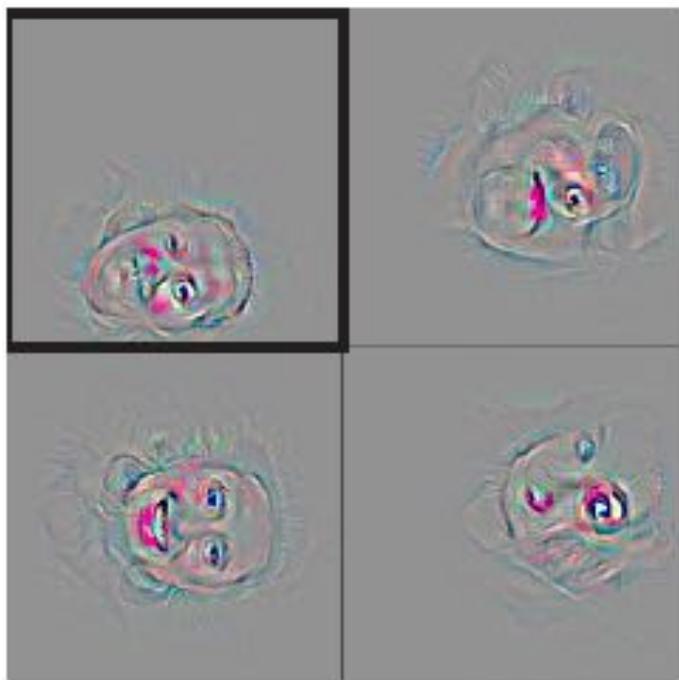
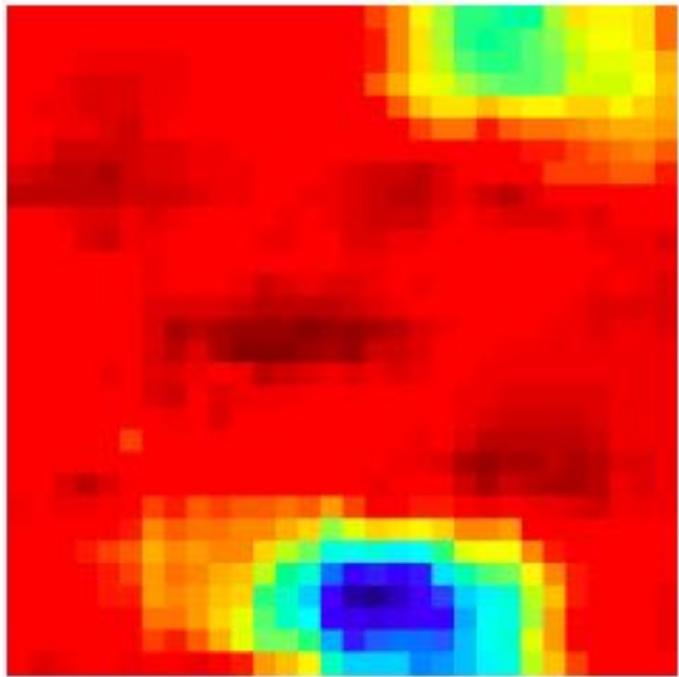


Input image



Total activation in most
active 5th layer feature map

Other activations from
same feature map



S.AFRICA



SENEGAL



RWANDA



MALE

FEMALE

MALE

FEMALE

FEMALE

MALE

FEMALE

MALE

SWEDEN

ICELAND

FINLAND

AFRICA

AVERAGE FACES

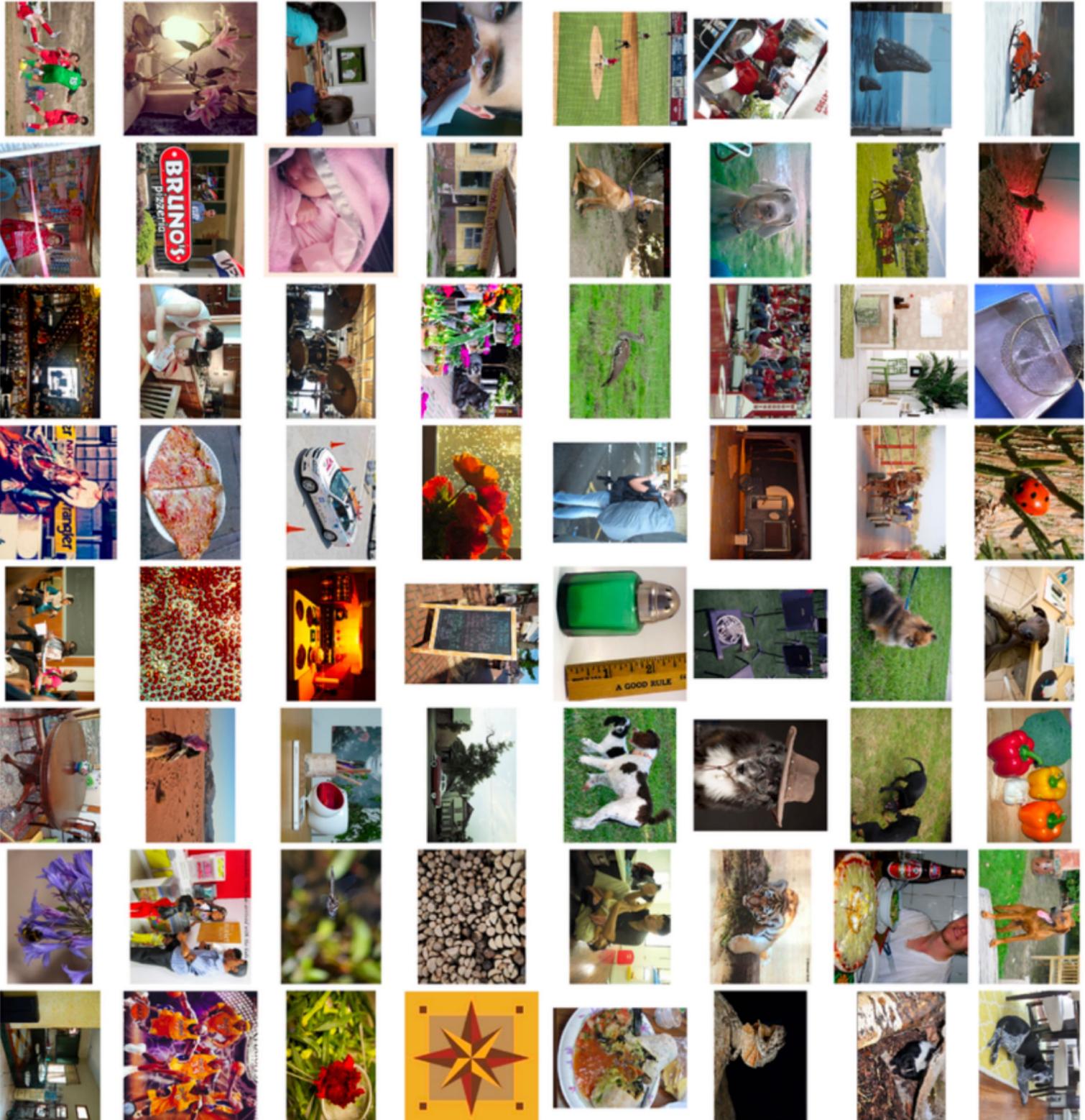
EUROPE

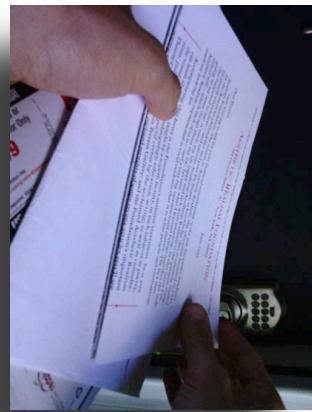
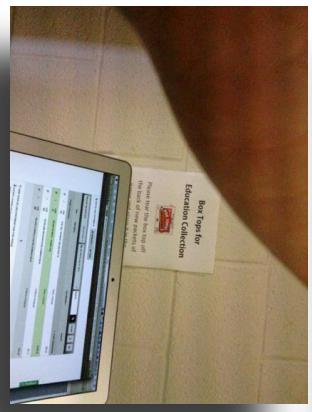
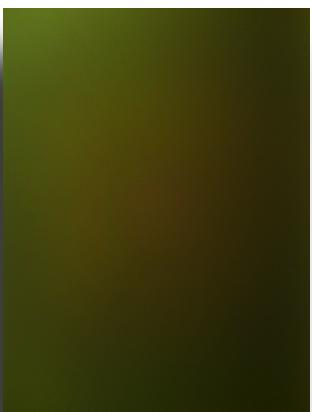
Classifier	Metric	DF	DM	LF	LM
MSFT	PPV(%)	76.2	100	100	100
	Error Rate(%)	23.8	0.0	0.0	0.0
	TPR(%)	100	84.2	100	100
	FPR(%)	15.8	0.0	0.0	0.0
Face++	PPV(%)	64.0	99.5	100	100
	Error Rate(%)	36.0	0.5	0.0	0.0
	TPR(%)	99.0	77.8	100	96.9
	FPR(%)	22.2	1.03	3.08	0.0
IBM	PPV(%)	66.9	94.3	100	98.4
	Error Rate(%)	33.1	5.7	0.0	1.6
	TPR(%)	90.4	78.0	96.4	100
	FPR(%)	22.0	9.7	0.0	3.6

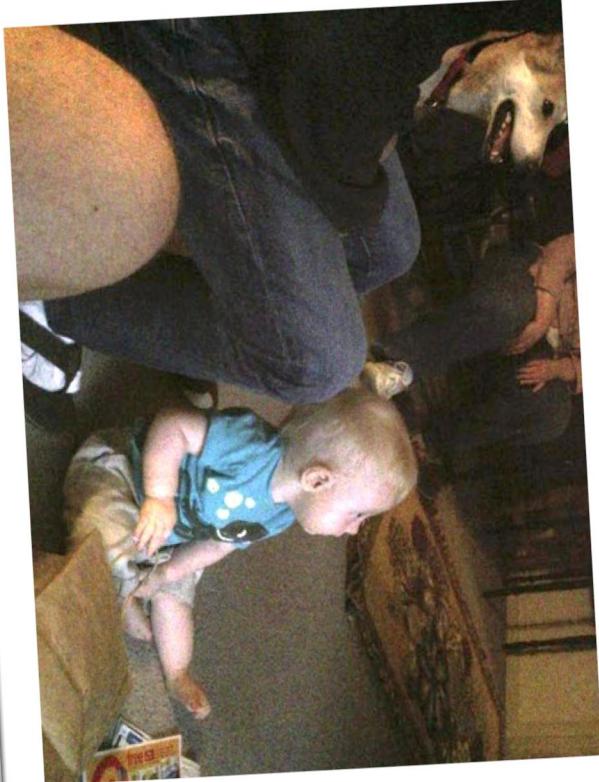
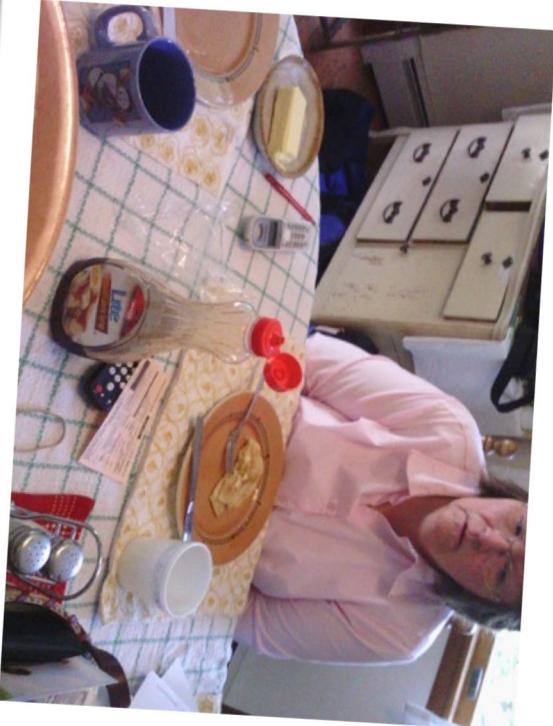
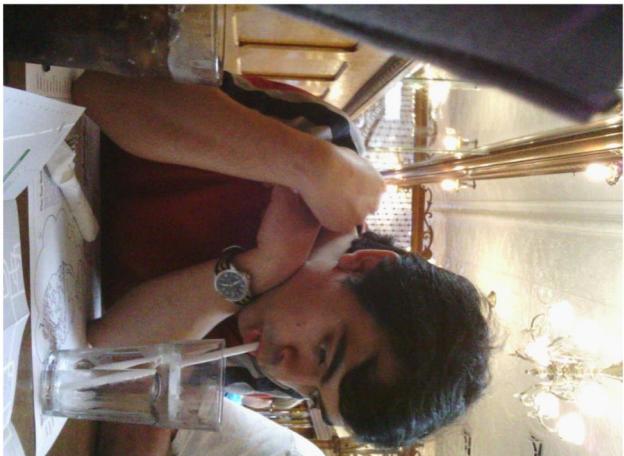
Table 5: Gender classification performance as measured by the positive predictive value (PPV), error rate (1-PPV), true positive rate (TPR), and false positive rate (FPR) of the 3 evaluated commercial classifiers on the South African subset of the PPB dataset. Results for South Africa follow the overall trend with the highest error rates seen on darker-skinned females.

Moral:

Algorithms are important.
But so is the training data!

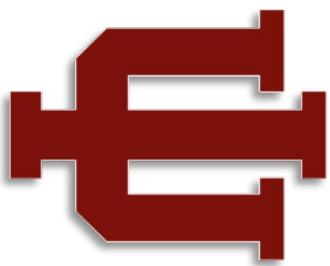






Studying visual learning in children using computer vision

(and vice-versa, maybe, eventually?)



David Crandall
School of Informatics, Computing, and Engineering
Indiana University
Bloomington

Swapnaa Jayaraman



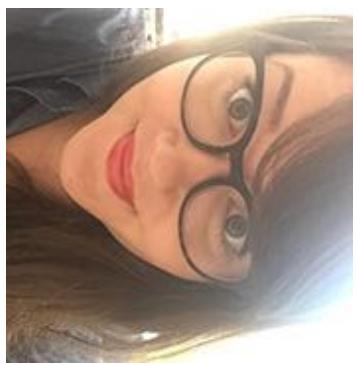
Zehua Zhang



Sven Bambach



Elizabeth Clerkin



Christina DeSerio



David Crandall



Linda Smith



Chen Yu





- Snapchat



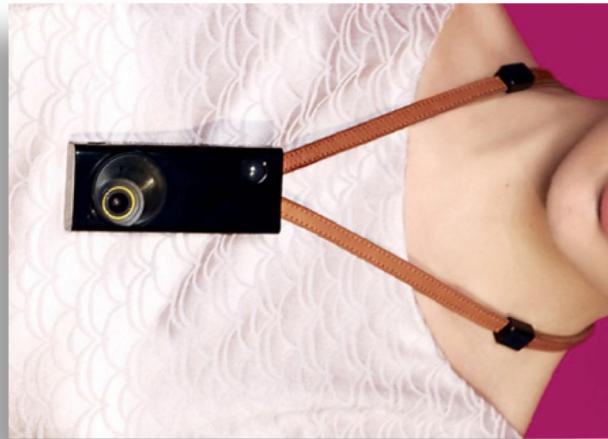
- GoPro



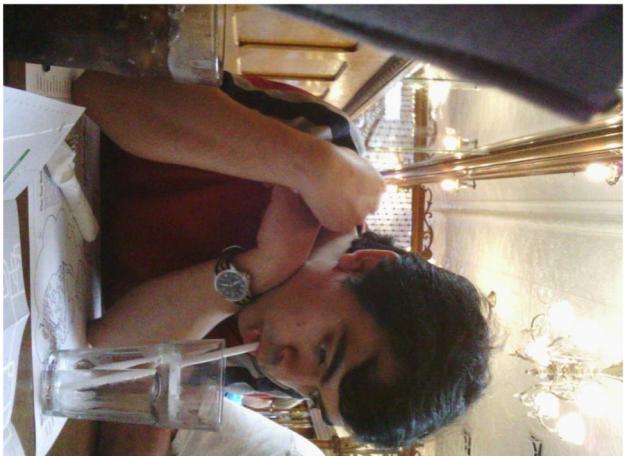
- Google

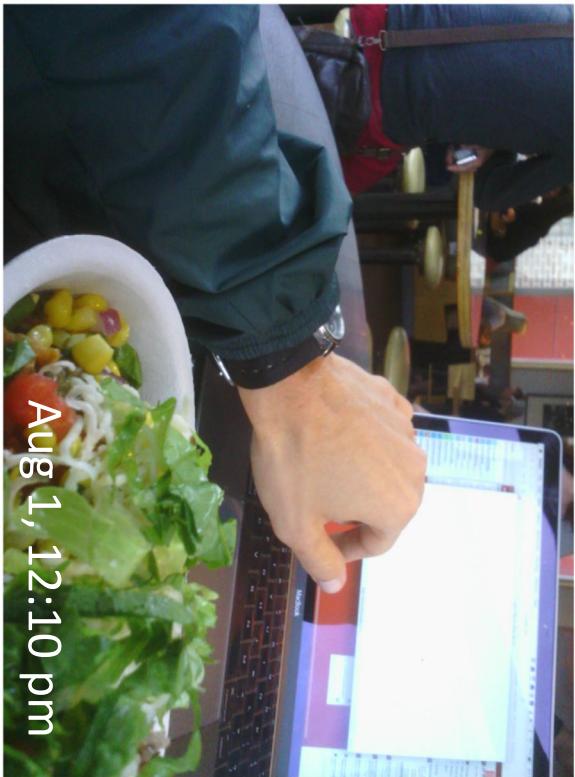


- Narrative



- Autographer





Aug 1, 12:10 pm



June 1, 12:10 pm



Sept 1, 12:10 pm



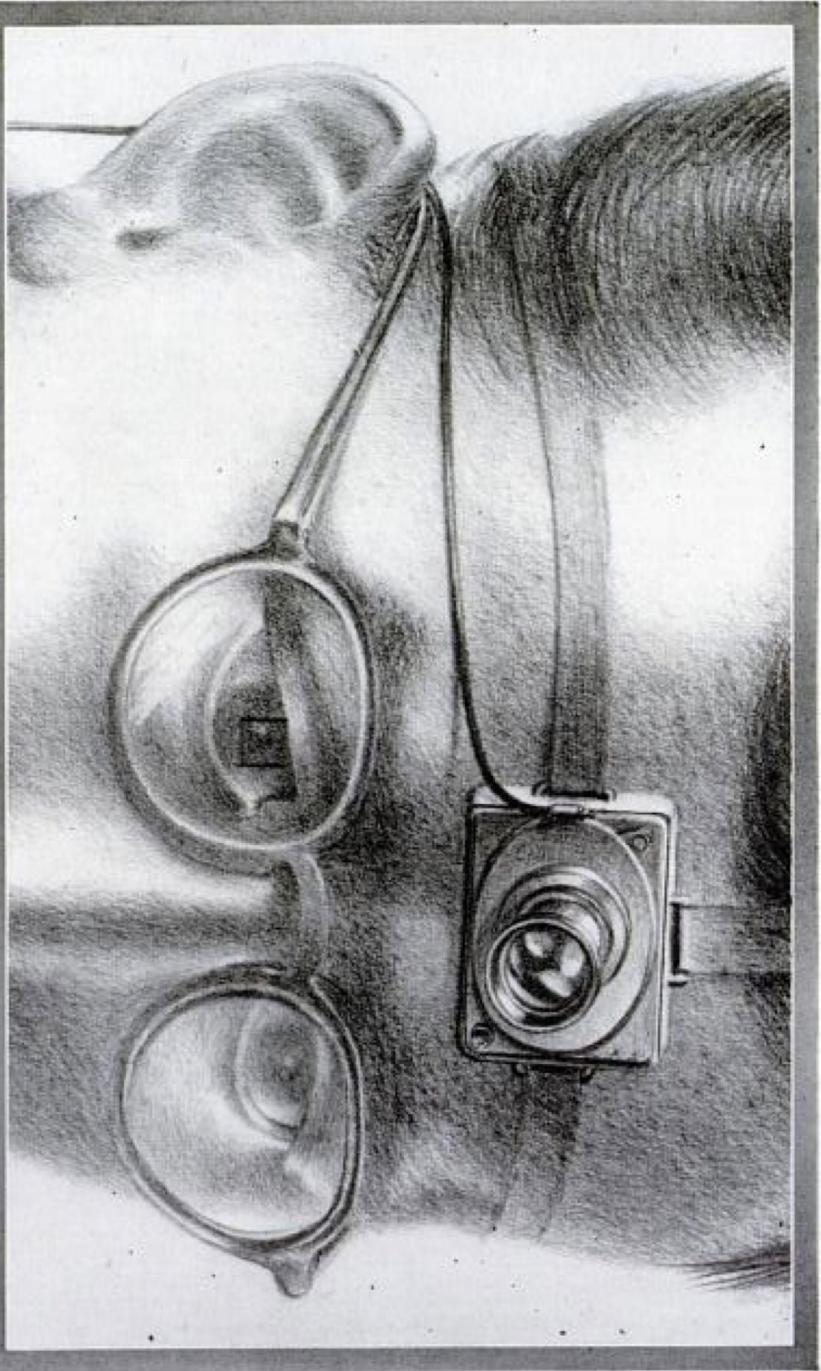
July 1, 12:10 pm

Vannevar Bush, *The Atlantic*, 1945

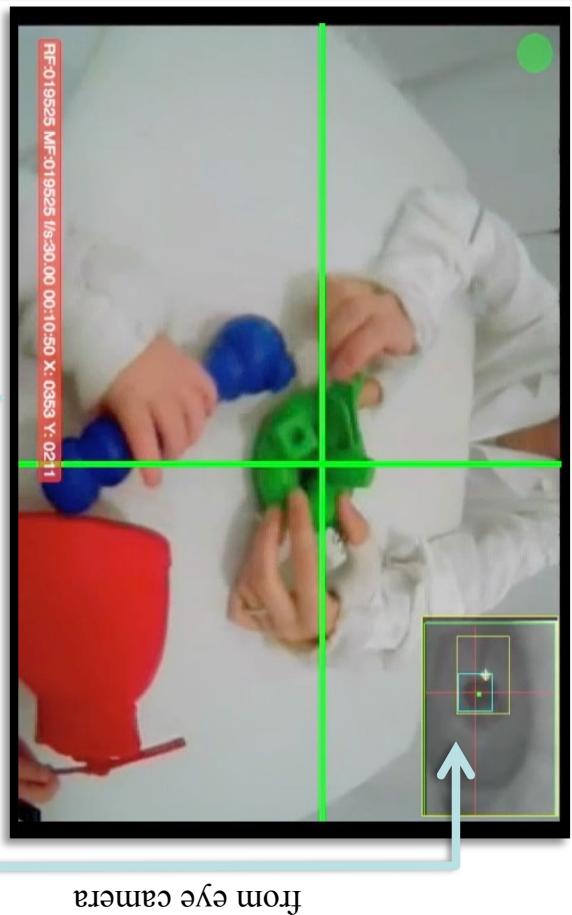
AS WE MAY THINK

A TOP U. S. SCIENTIST FORESEES A POSSIBLE FUTURE WORLD

A SCIENTIST OF THE FUTURE RECORDS EXPERIMENTS WITH A TINY CAMERA FITTED WITH UNIVERSAL-FOCUS LENS. THE SMALL SQUARE IN THE EYEGLASS AT THE LEFT SIGNS THE OBJECT



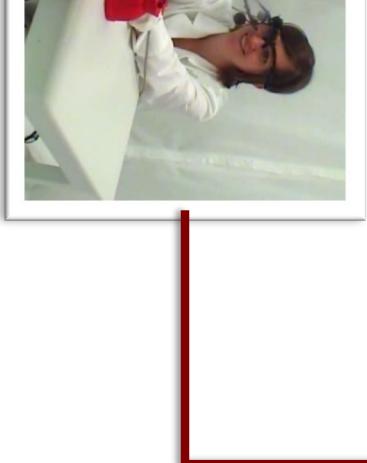
child's first person view



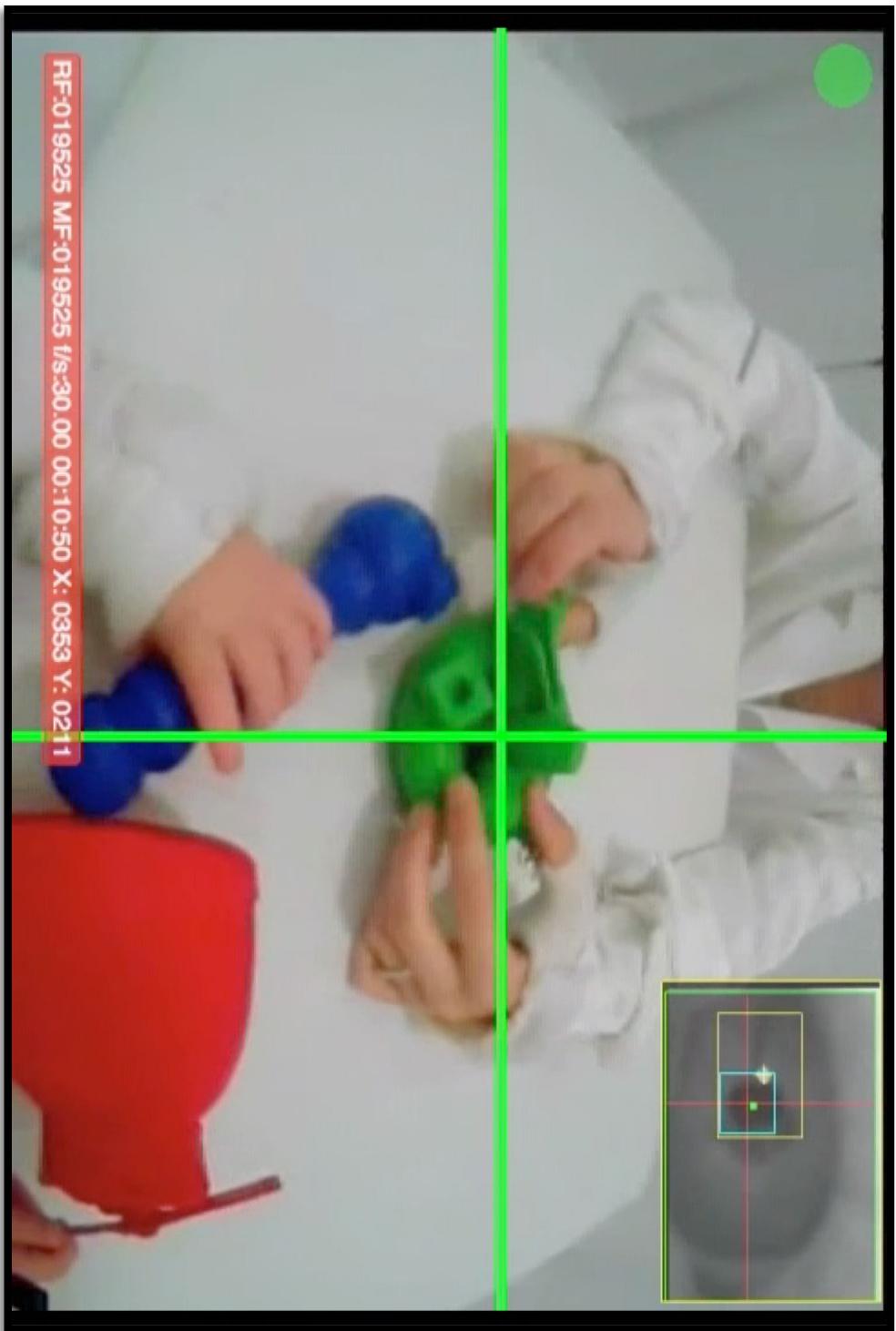
parent's first person view



from eye camera



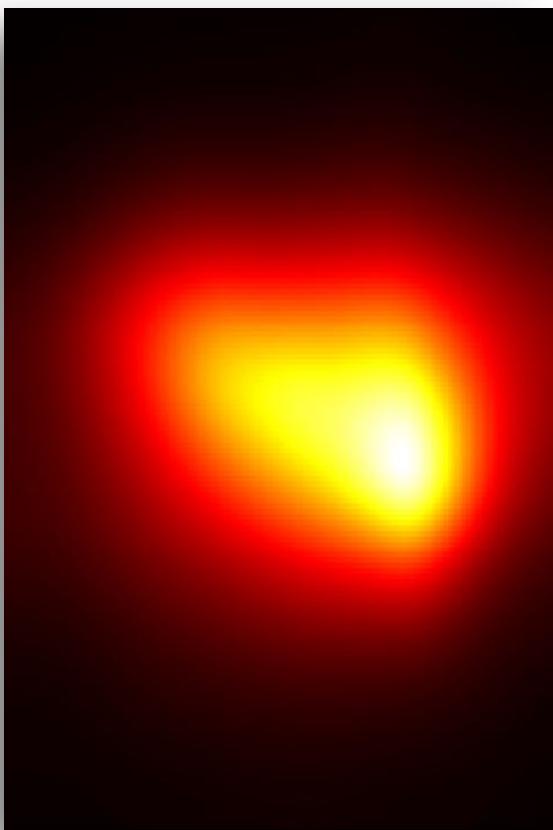
from head camera



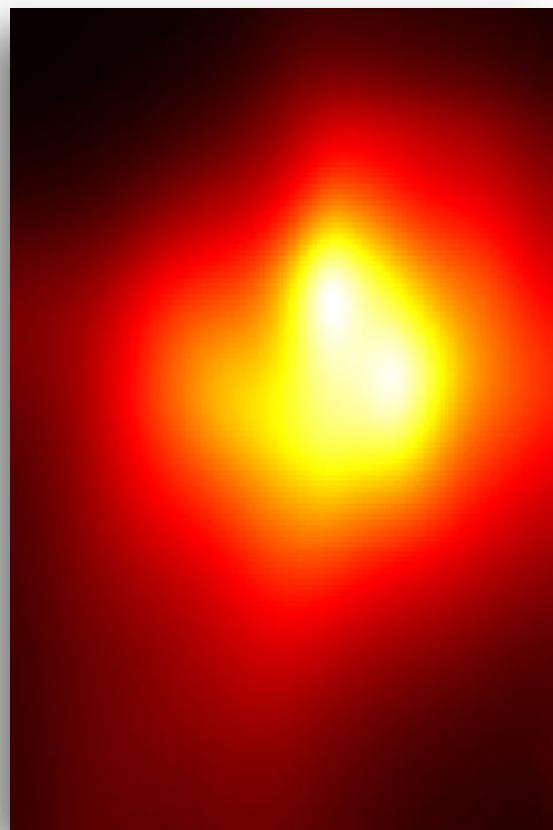
RF:019525 MF:019525 t/s:30.00 00:10:50 X: 0353 Y: 0211

Eye gaze *within* the visual field

Parents



Children



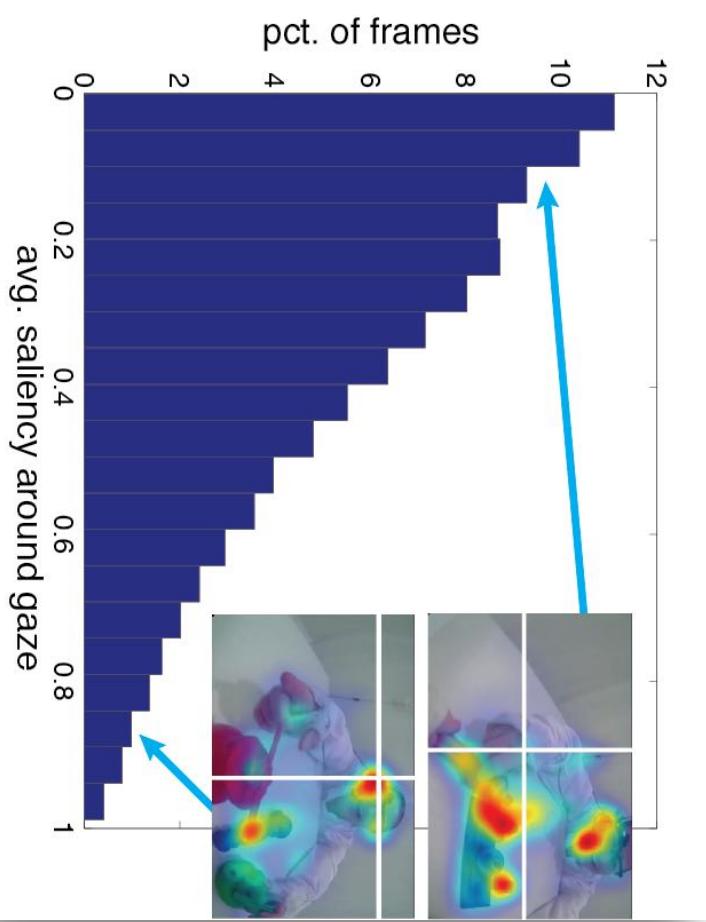
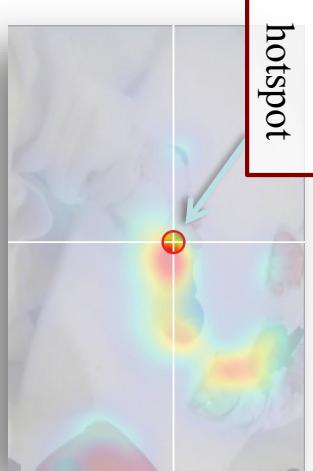
$\mu = (361, 224)$
 $\sigma_x = 53, \sigma_y = 60$
 $N = 148,279$

$\mu = (340, 231)$
 $\sigma_x = 86, \sigma_y = 65$
 $N = 148,279$

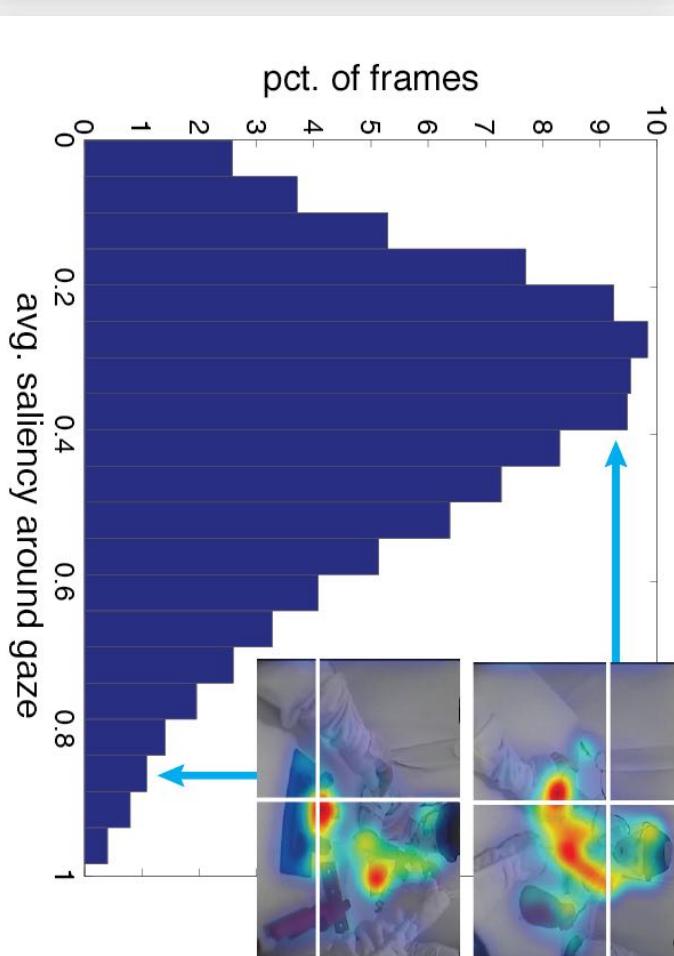
Visual saliency in first-person views

Comparison of average saliency within hotspot around gaze location on 148,000 frames

→ Gaze predictiveness differs significantly

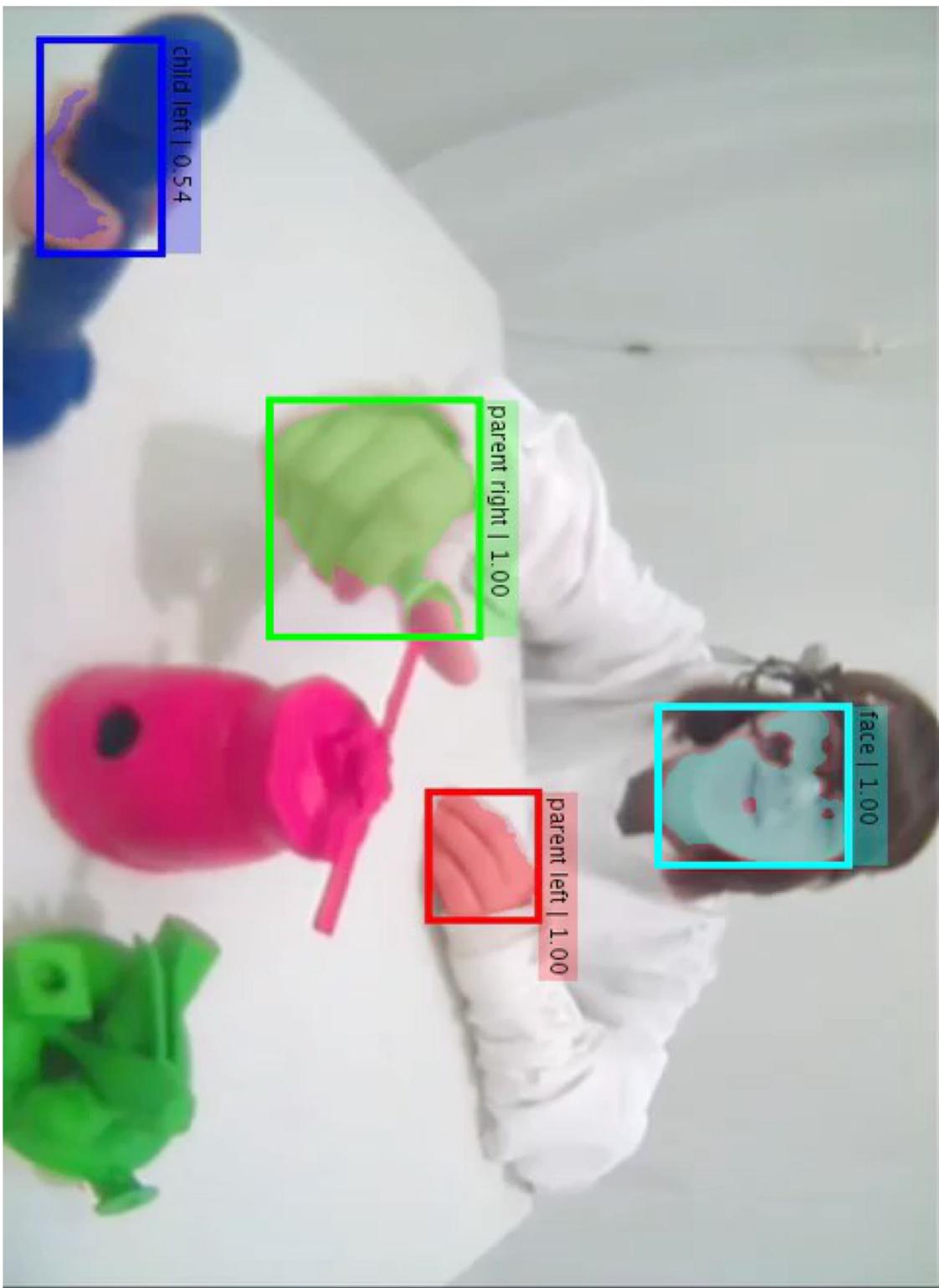


Child



Parent

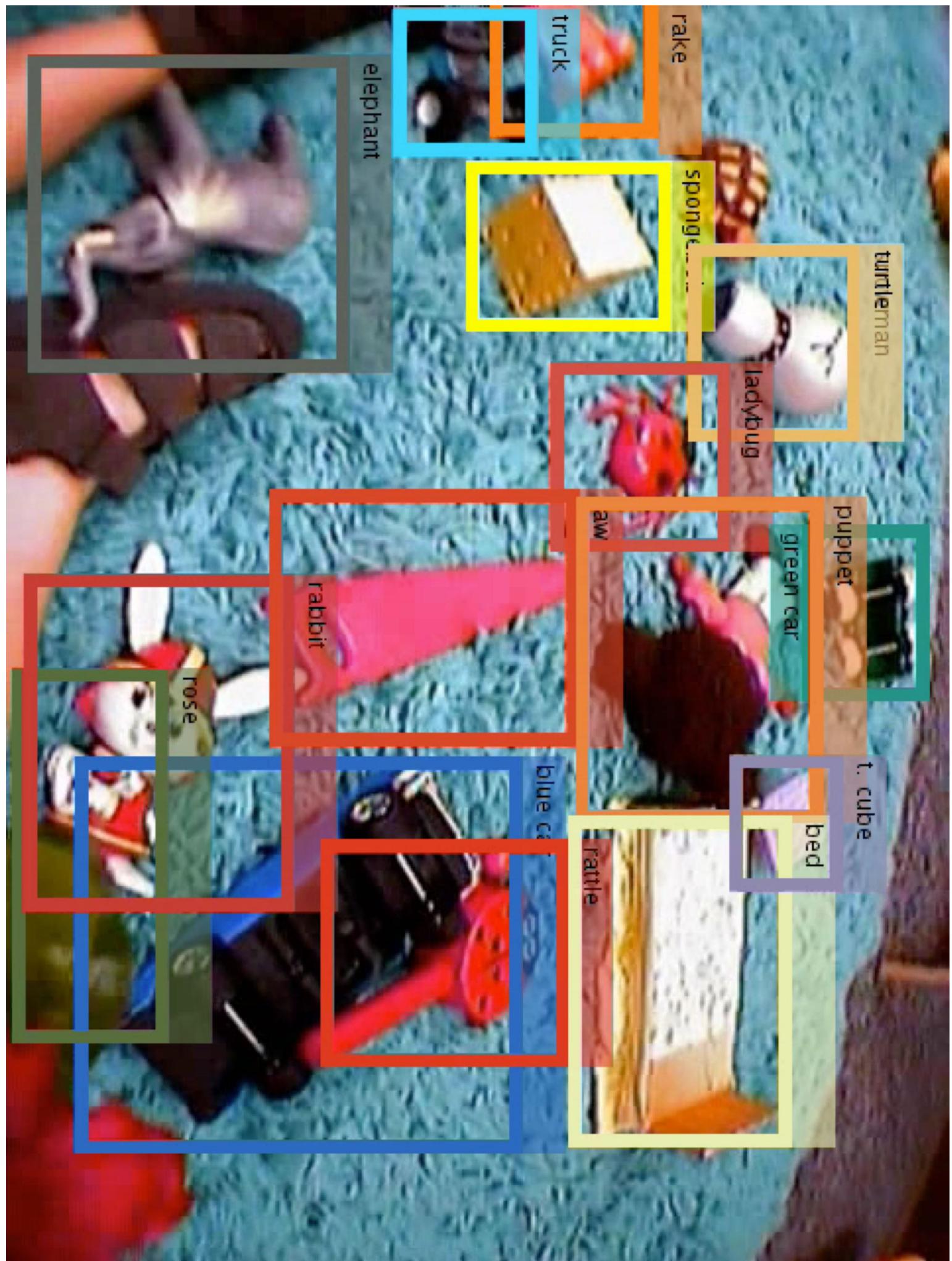
Hand detection & segmentation



S. Lee, S. Bambach, C. Yu, D. Crandall. This hand is my hand: A probabilistic approach to hand disambiguation in egocentric video, *CVPR Egovision*, 2014.

S. Bambach, S. Lee, D. Crandall, C. Yu. Lending a hand: Detecting hands and recognizing activities in complex egocentric interactions. *ICCV*, 2015.



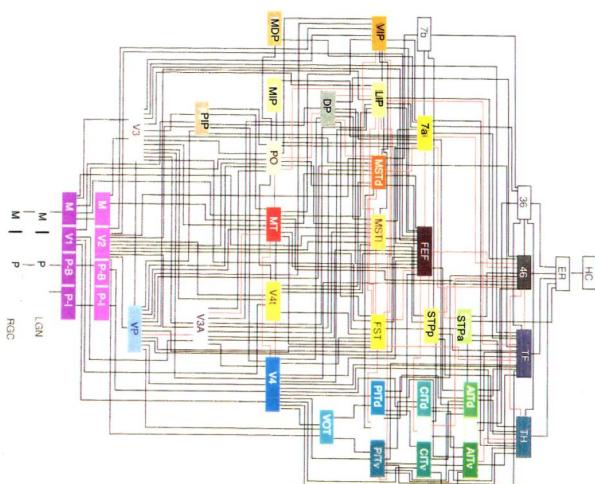


How do children's amazing visual learning systems work?

- Probably not with CNNs
- Probably not trained on ImageNet
 - Interact with objects
 - See objects in natural contexts
 - Their bodies impose constraints



Figure 4. Interplay of visual areas. This sketchy diagram visualizes visual areas and according to the same scheme as in Figure 7, 2 additional visual stages. The visual input is shown as MOP and PIP. The visual areas are V1, V2, V3A, V3B, P-B, P-L, MT, MSTG, MSST, LIP, VIP, Dp, PO, FEF, STP, STPa, STPb, Crt, Crtv, PIPg, PIPb, NOT, and HC. The visual areas are grouped by 10 pathways (in green) and each group is labeled with a regional grouping.



- Can we use wearable cameras to characterize the “training data” that they receive?

The “training data” each day at 1 year old:

- 16 hours of visual experience
- 20,000 word
- Self-generated through activity (10,000 steps)



Child view

Parent view

Whose view is better for object learning?

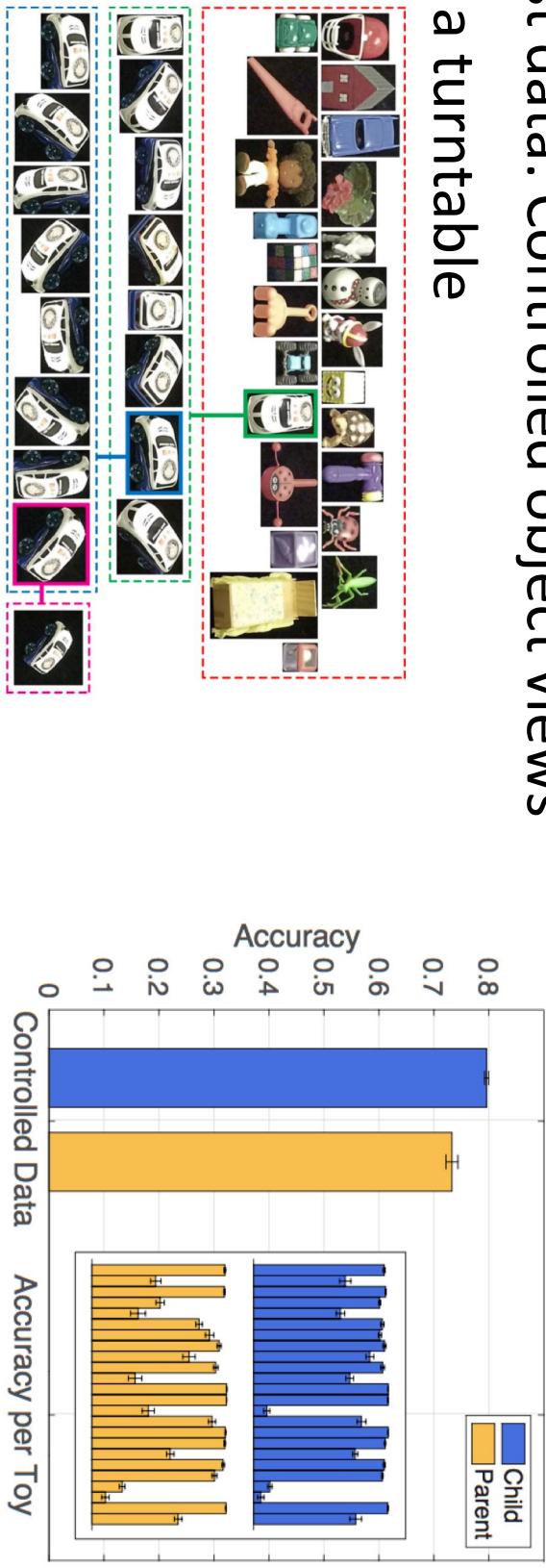
- Training data: ~100 minutes of video from 10 child-parent dyads playing with 24 toys in a free-form lab environment



(a) Parent view

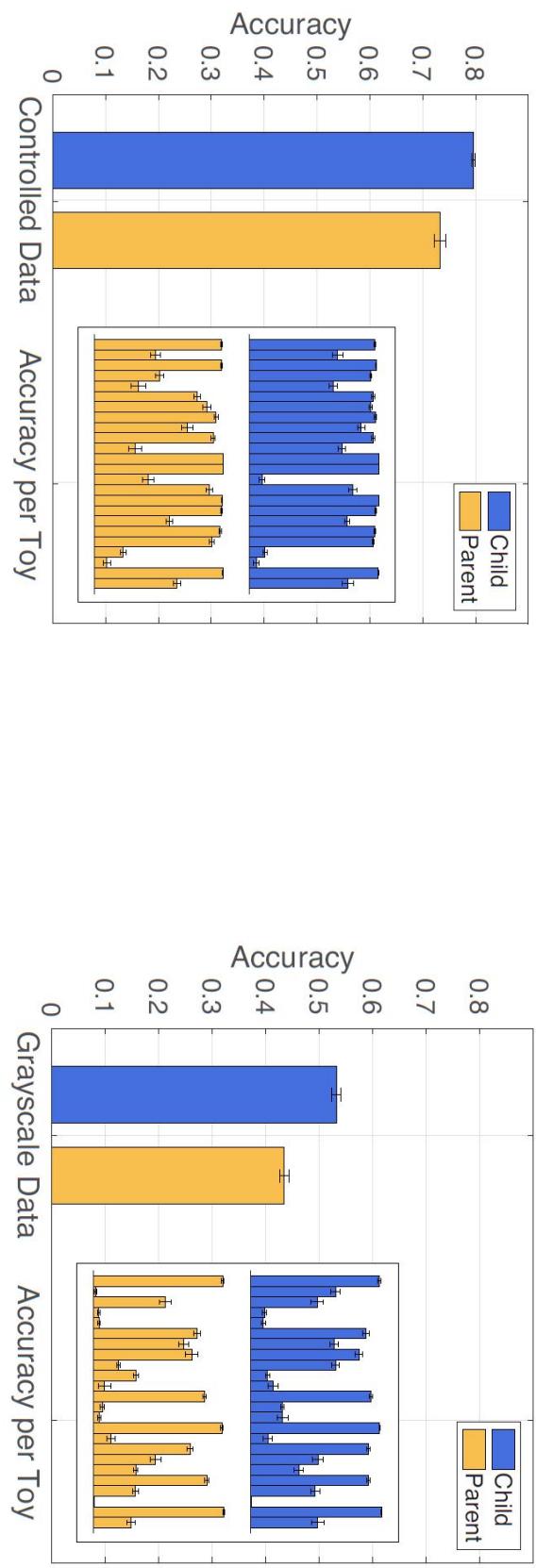
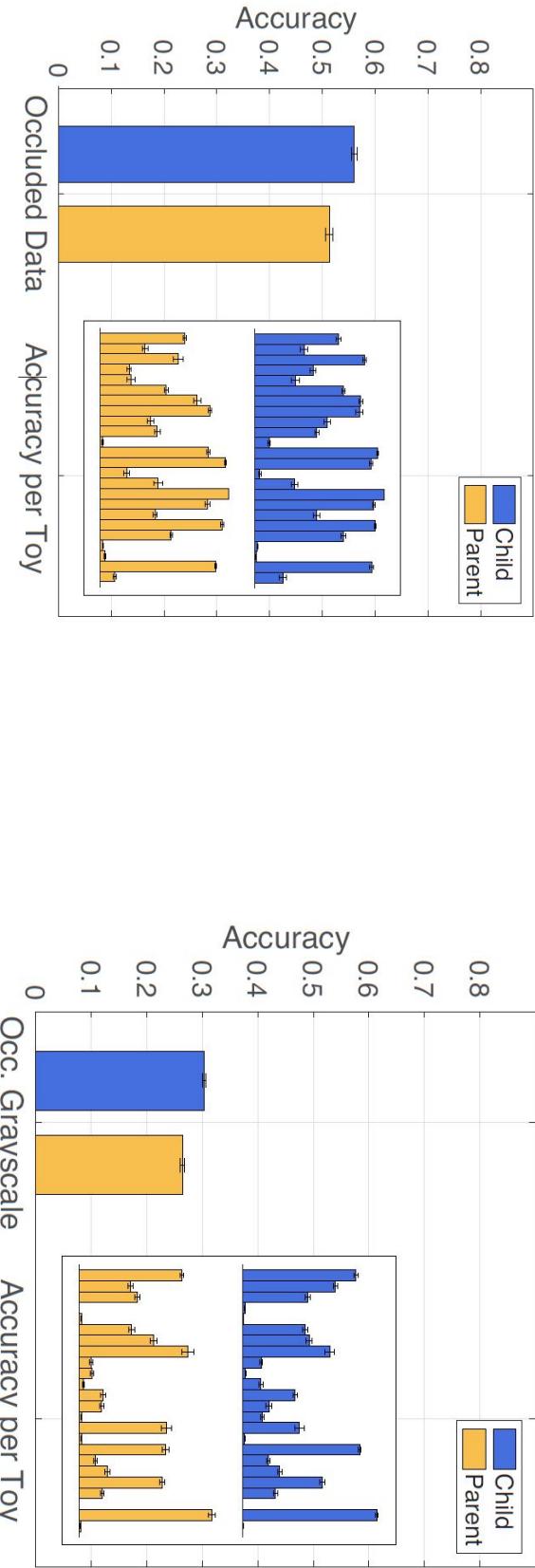
(b) Child view

- Test data: Controlled object views on a turntable

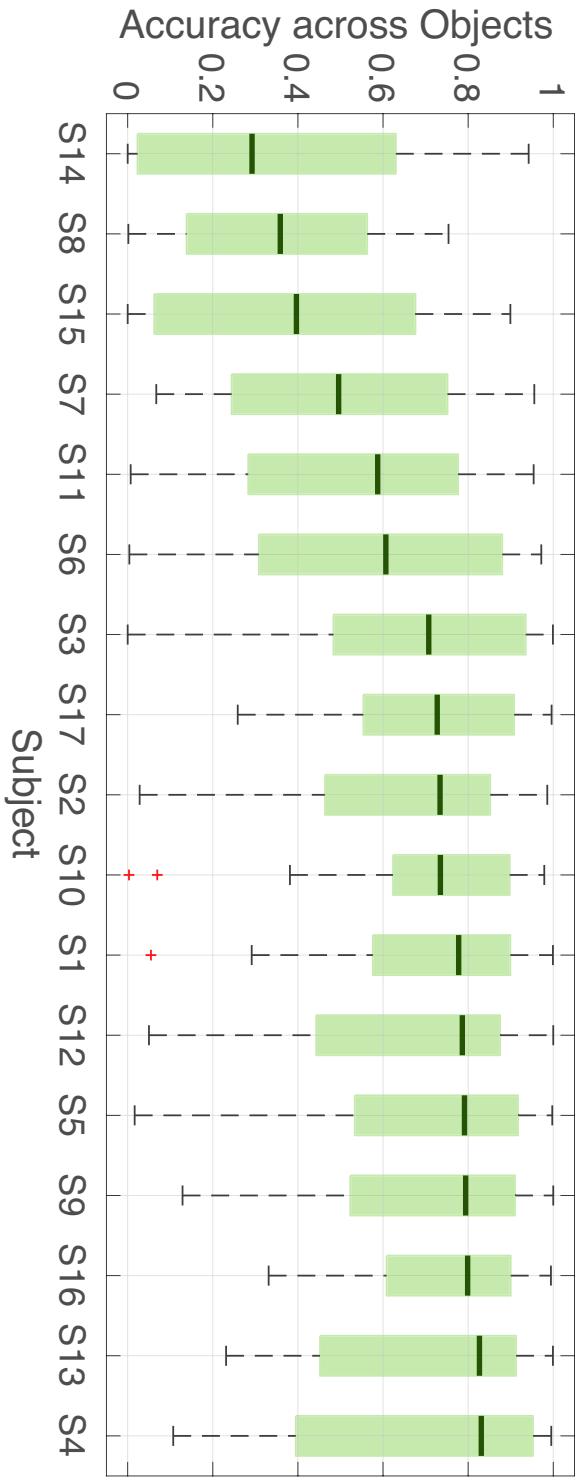


S. Bambach, D. Crandall, L. Smith, C. Yu. Active viewing in toddlers facilitates visual object learning: An egocentric vision approach, *CogSci*, 2016.

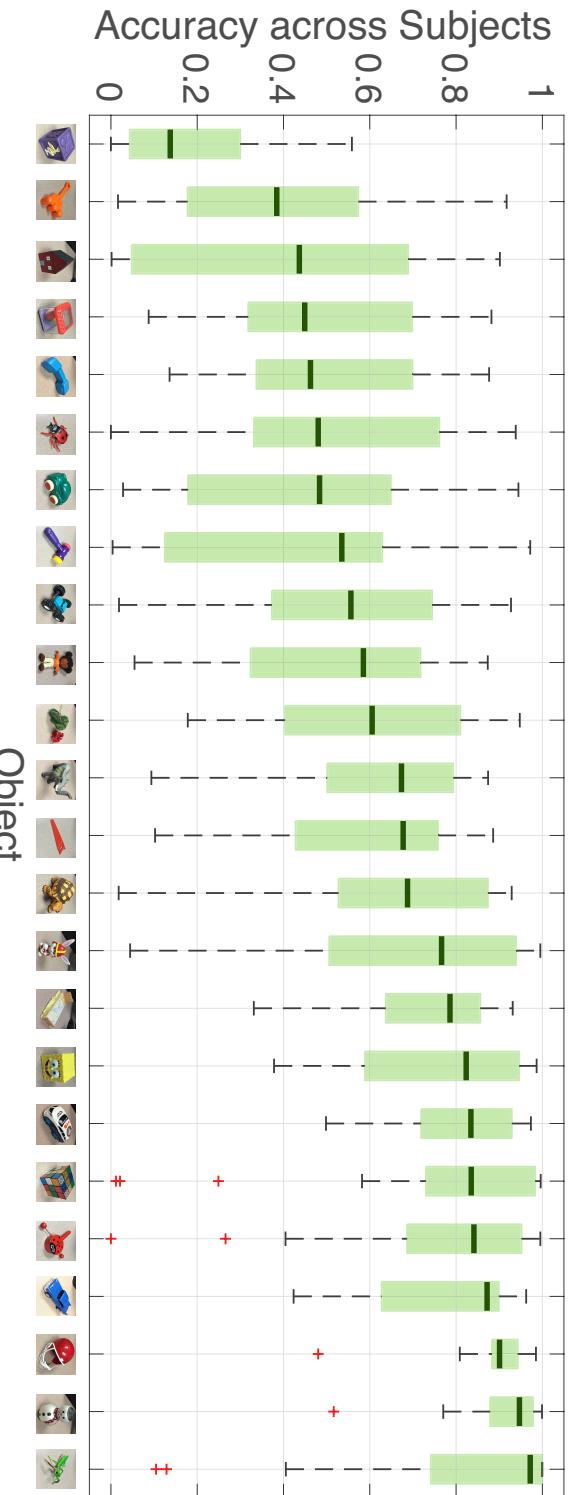
- CNNs can learn models of the toys in this first-person data;
- They generalize to different contexts and viewpoints;
- Toddler views seem to be of particularly high value for training.



Each Subject's Recognition Accuracies Across Objects



Each Object's Recognition Accuracies Across Subjects

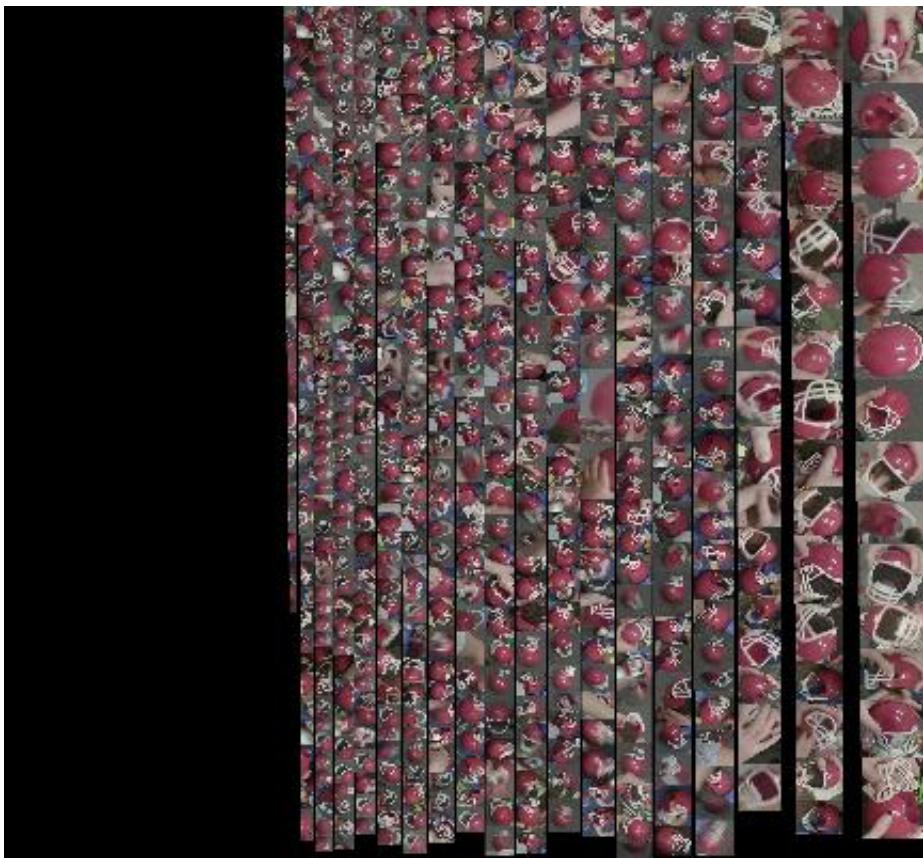


S. Bambach, Z. Zhang, D. Crandall, C. Yu. Exploring inter-observer differences in first-person object views using deep learning models, *ICCV Workshop on Mutual Benefits of Cognitive and Computer Vision*, 2017.

Different views of the world



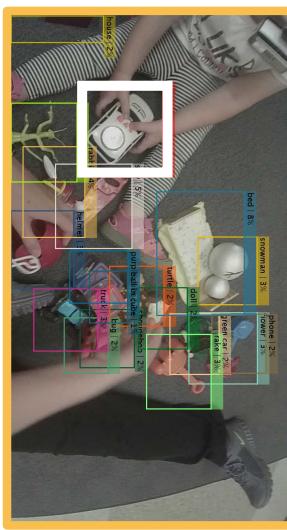
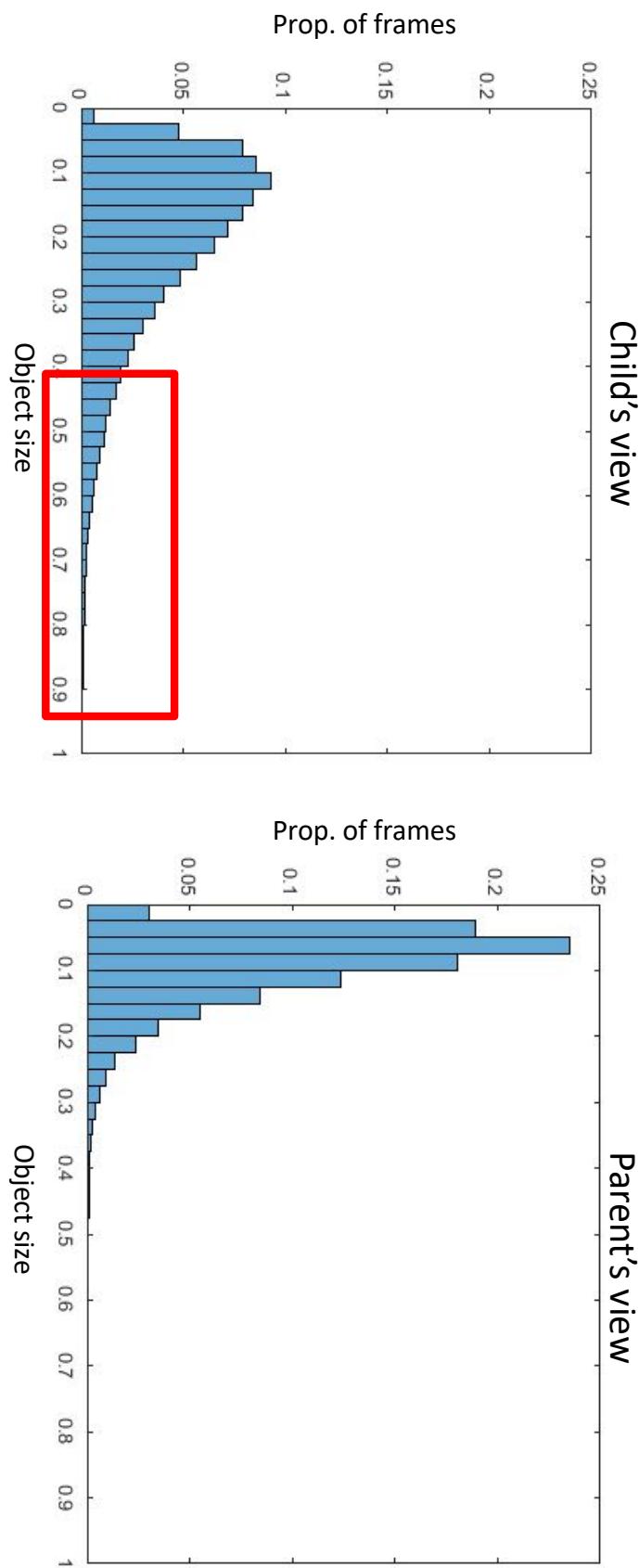
child's view



parent's view

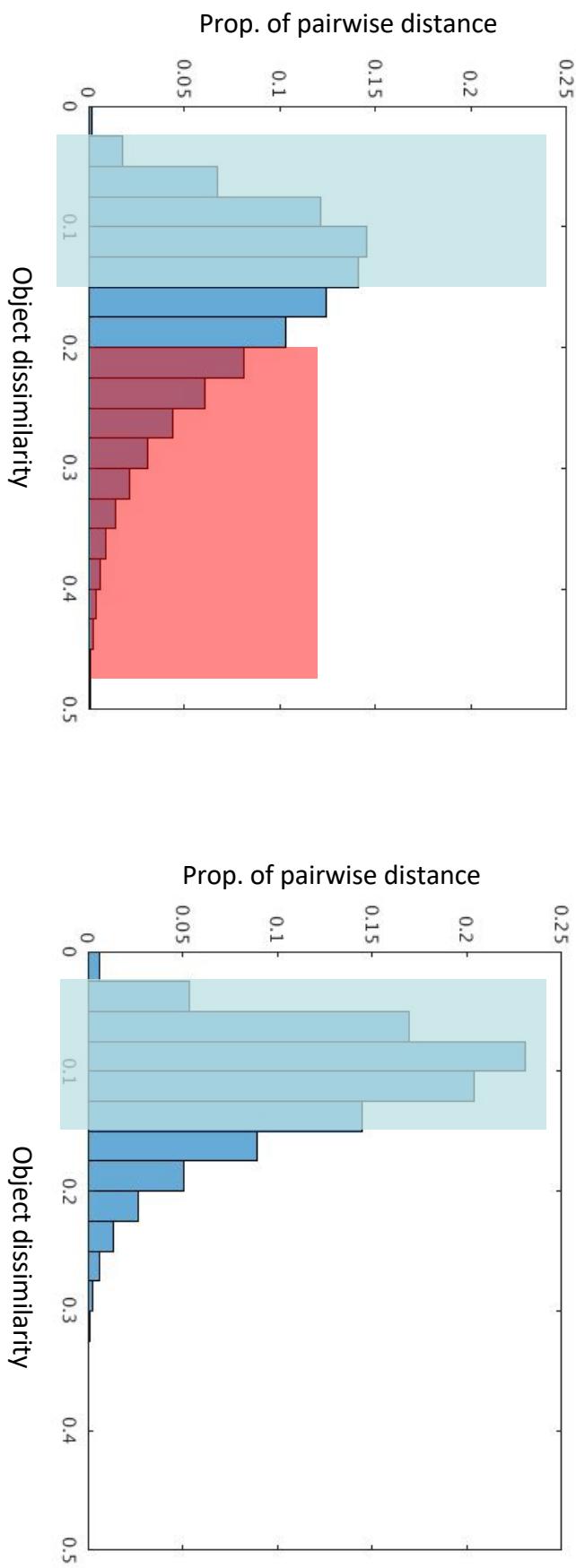
S. Bambach, D. Crandall, L. Smith, C. Yu. Active viewing in toddlers facilitates visual object learning: An egocentric vision approach, *CogSci*, 2016.

Object Size



S. Bambach, D. Crandall, L. Smith, C. Yu. Active viewing in toddlers facilitates visual object learning: An egocentric vision approach, *CogSci*, 2016.

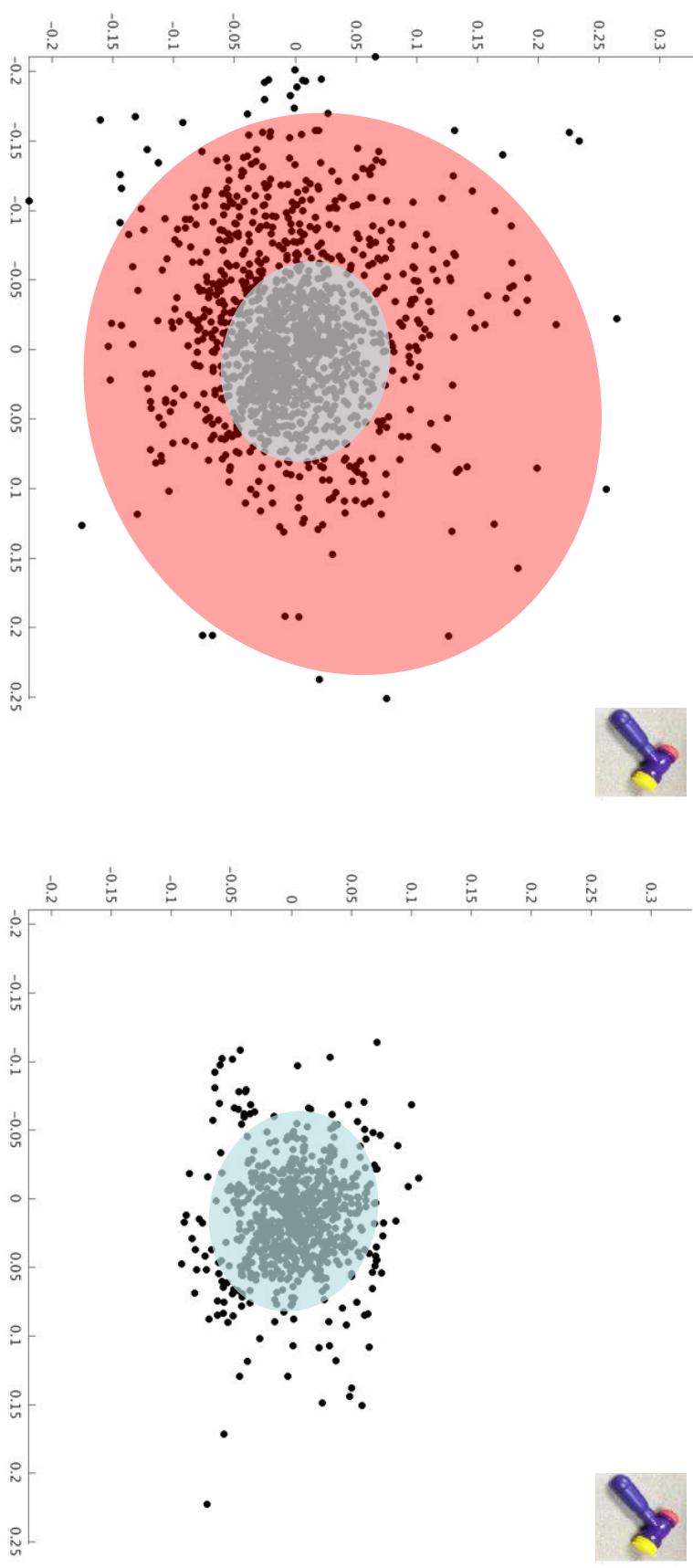
Diversity of object views



Toddlers spend hours every day playing with toys, actively manipulating them and creating “training data” by **self-selecting** object views. This creates more **diverse views** for children than the parents.

Consistency and Variability

Child's view



Parent's view

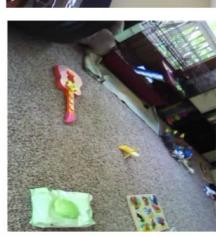
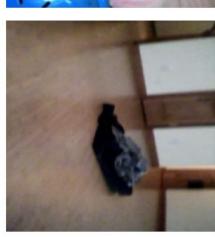
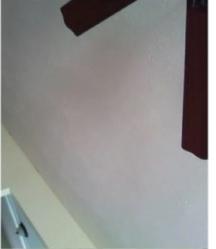
S. Bambach, D. Crandall, L. Smith, C. Yu. Active viewing in toddlers facilitates visual object learning: An egocentric vision approach, *CogSci*, 2016.

This “training data” changes over time

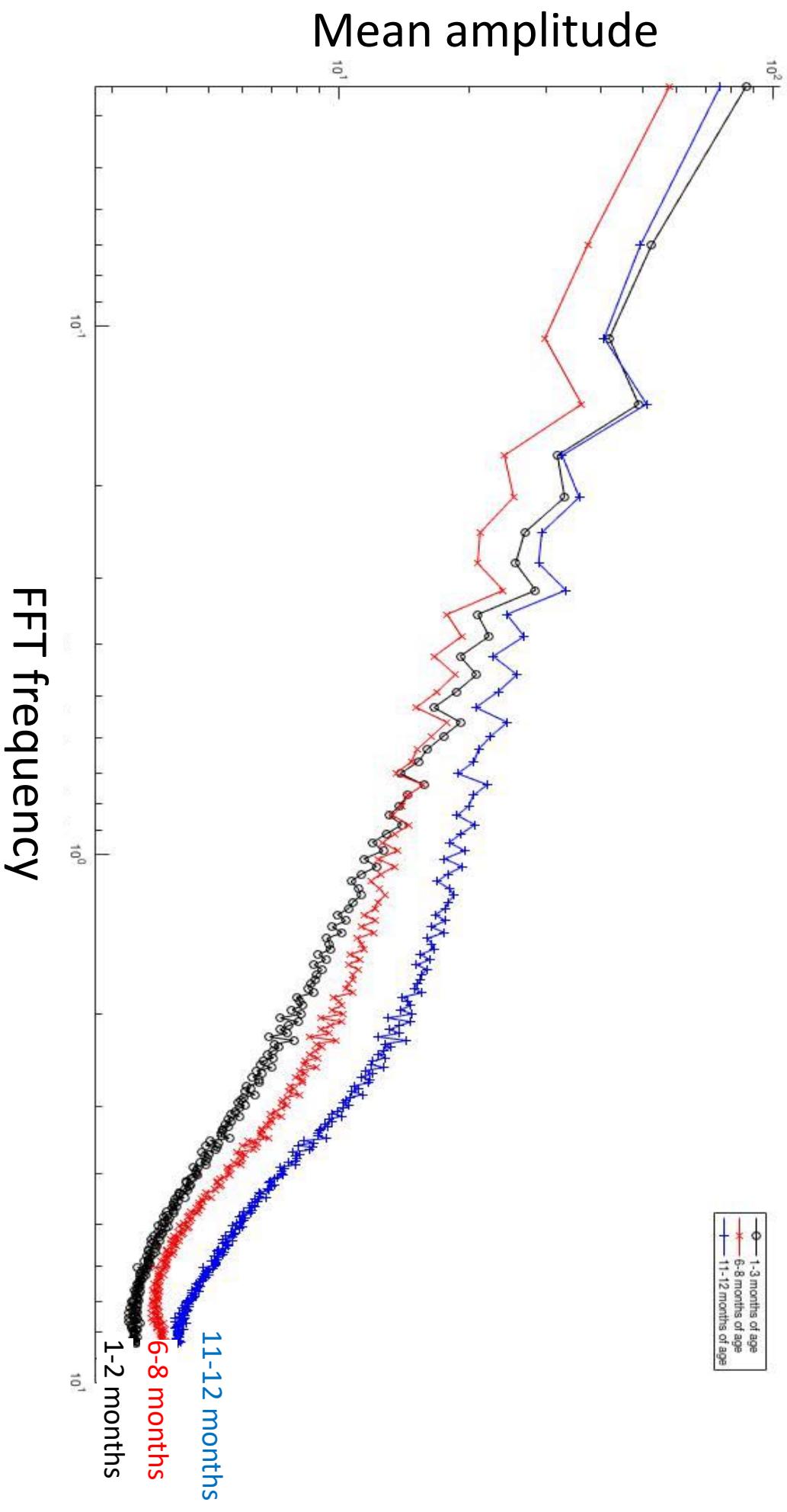


HomeView: dataset of developmentally indexed head-mounted camera data

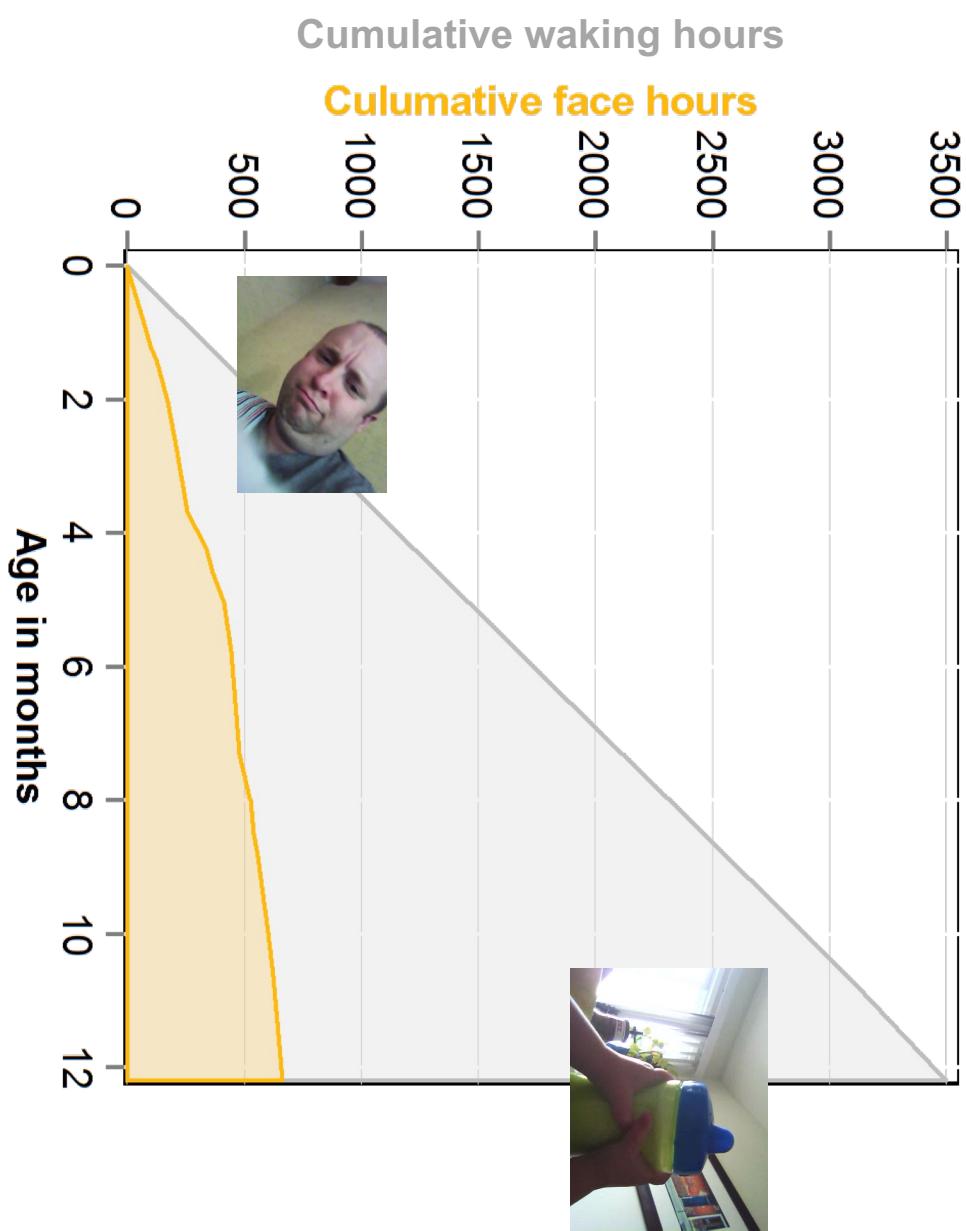
- 101 infants, 3 weeks – 24 months
- 4 to 6 hours of in-home video
- No experimenters present
- 30 Hz, 750 million frames

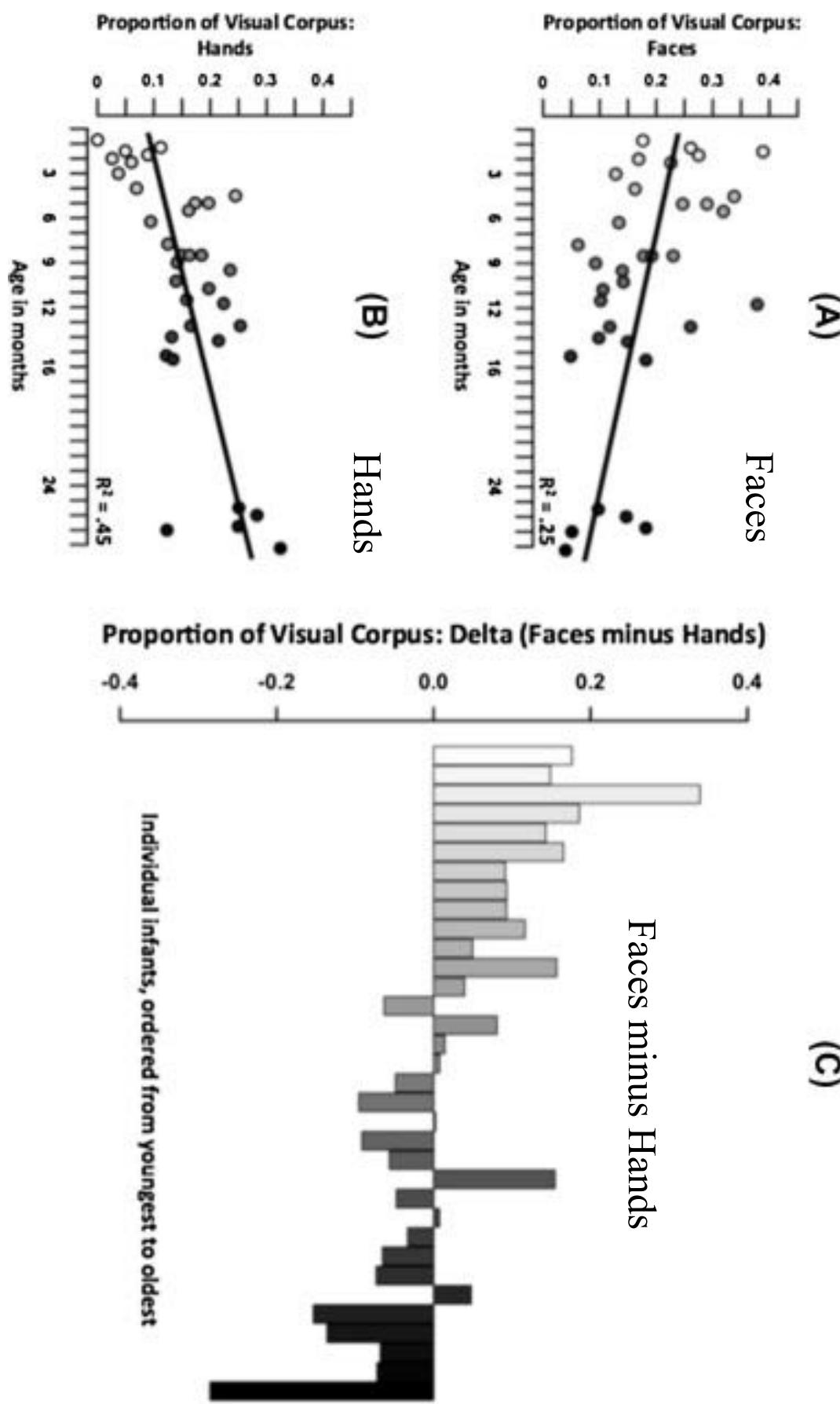


Low-level image stats change over time



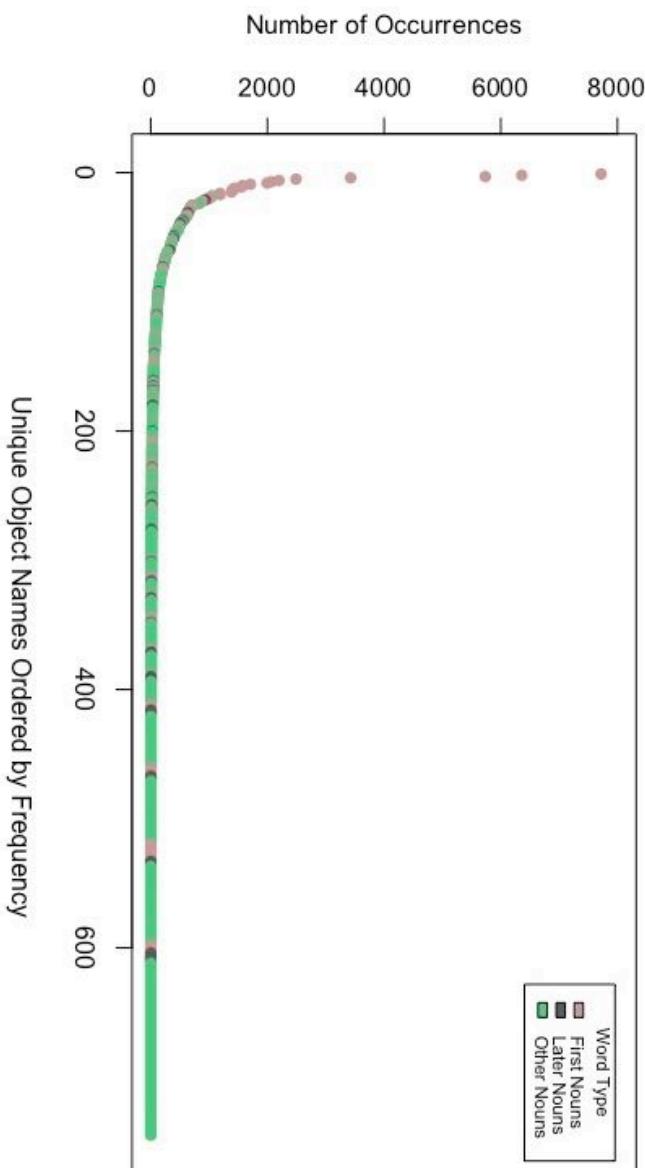
Frequency of objects in view also change





Word learning

- Labeled instances of 745 objects in HomeView data
- **Key finding:** The visual frequency of objects in scenes at 8-10 months – not the frequency of their names -- predicts age of acquisition of object names.



- This suggests early visual word learning is largely “unsupervised,” around 1 year transitions to “supervised”

Short term questions:

How can automated computer vision help analyze egocentric images?

Longer-term question:

Can we model head and hand pose, eye gaze, saliency, and activity to predict and explain human vision?

Long-term question:

What can machine learning (and deep learning in particular) reveal about human visual learning, and vice versa?



IU Bloomington Newsroom

All IUB News Arts & Humanities Business International Law & Policy Life & Health Sciences Science & Technology

IUB Newsroom » IU Bloomington awards up to \$3 million to advance new approach to study of learning

[PRINT](#) [SHARE](#)

IU Bloomington awards up to \$3 million to advance new approach to study of learning

Campus's Emerging Areas of Research program makes inaugural grant

Jan. 11, 2017

FOR IMMEDIATE RELEASE

BLOOMINGTON, Ind. — Can machines learn to think like children? An interdisciplinary team of cognitive scientists, neuroscientists and computer scientists at Indiana University Bloomington has received the campus's inaugural Emerging Areas of Research funding award to explore that question.

Called "[Learning: Brains, Machines and Children](#)," the first Emerging Areas of Research initiative is led by Linda Smith, Distinguished Professor and Chancellor's Professor of psychological and brain sciences in the IU Bloomington College of Arts and Sciences. The research team will receive up to \$3 million for the four-year project.

The Emerging Areas of Research program was proposed in the [Bicentennial Strategic Plan for IU Bloomington](#) as part of a suite of initiatives designed to invest in research innovation and excellence on campus.

"The Emerging Areas of Research funding program ensures that our campus is constantly thinking about the next frontier in knowledge and hiring new faculty to ensure we can reach that frontier," IU Bloomington Provost and Executive Vice President Lauren Rebel said. "The stellar team led by Professor Smith epitomizes the strengths that allow us to build toward that future."

Representing expertise in developmental psychology, human learning, neuroscience and artificial intelligence, the inaugural Emerging Areas of Research team will pursue a comprehensive and unified theory of learning rooted in research about how infants and children learn to classify faces, objects, letters, numbers, etc.

"We are thrilled that IU Bloomington's first Emerging Areas of Research initiative represents such a profound area of focus and such an impressive group of researchers," said Rick Van Kooten, the vice provost for research at IU Bloomington, whose office oversees the Emerging Areas of Research program. "The proposal leverages exemplary areas of research on our campus, and the combination is certain to yield exciting results."



Linda B. Smith | Photo by Indiana University

[Print-Quality Photo](#)

Media Contacts

Lauren Bryant

Office of the Vice Provost for Research

Office 812-855-4152

labryant@indiana.edu

Related Stories

New Emerging Areas of Research funding program launches at IU Bloomington

Workshop on

Egocentric Vision: From Science to Real-World Applications



Indiana University, Bloomington Indiana — June 3-5, 2017

Sponsored by the Indiana University Ostrom Grants Program, the College of Arts and Science, the School of Informatics and Computing, and the OVPR Emerging Areas of Research program.

<http://vision.soic.indiana.edu/egocentric-workshop-2017>

Organizers

	David Crandall School of Informatics and Computing Indiana University		Michael Ryoo School of Informatics and Computing Indiana University		Linda Smith Psychological and Brain Sciences Indiana University		Chen Yu Psychological and Brain Sciences Indiana University
--	--	--	--	--	--	--	--

Speakers and panelists

	Karen Adolf Psychology NYU		Dana Ballard Computer Science UT Austin		Derek Houston Otolaryngology-Head and Neck Surgery The Ohio State University		Michael Gaffney Psychiatry Washington University in St. Louis
	Mary Hayhoe Center for Perceptual Systems UT Austin		Dan Kennedy Psychological and Brain Sciences Indiana University		Maithilee Kunda Electrical Engineering and Computer Science Vanderbilt University		Hyun Soo Park Computer Science and Engineering University of Minnesota
	Jim Rehg School of Interactive Computing Georgia Tech		Alexander Schwing Electrical and Computer Engineering University of Illinois at Urbana-Champaign				

Talks and slides online!



Ruslan Salakhutdinov

Machine Learning
Carnegie Mellon University / Apple



Jim Rehg

School of Interactive Computing
Georgia Tech



Alexander Schwing

Electrical and Computer Engineering
University of Illinois at Urbana-Champaign

Thanks!

<http://vision.soic.indiana.edu/>



David Crandall



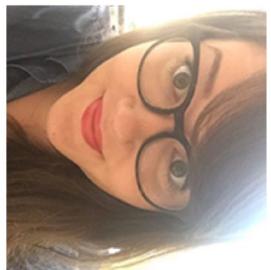
Linda Smith



Chen Yu



Sven Bambach



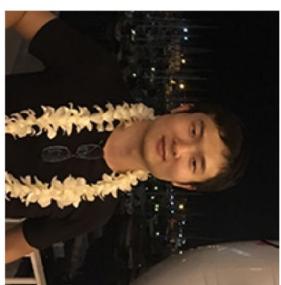
Elizabeth Clerkin



Christina DeSerio



Swapnaa Jayaraman



Zehua Zhang

Sponsors: NSF, NIH, Google, Nvidia, IU EAR, IU FRSP, Lilly Endowment