

# Up-convolutional networks and their applications

---

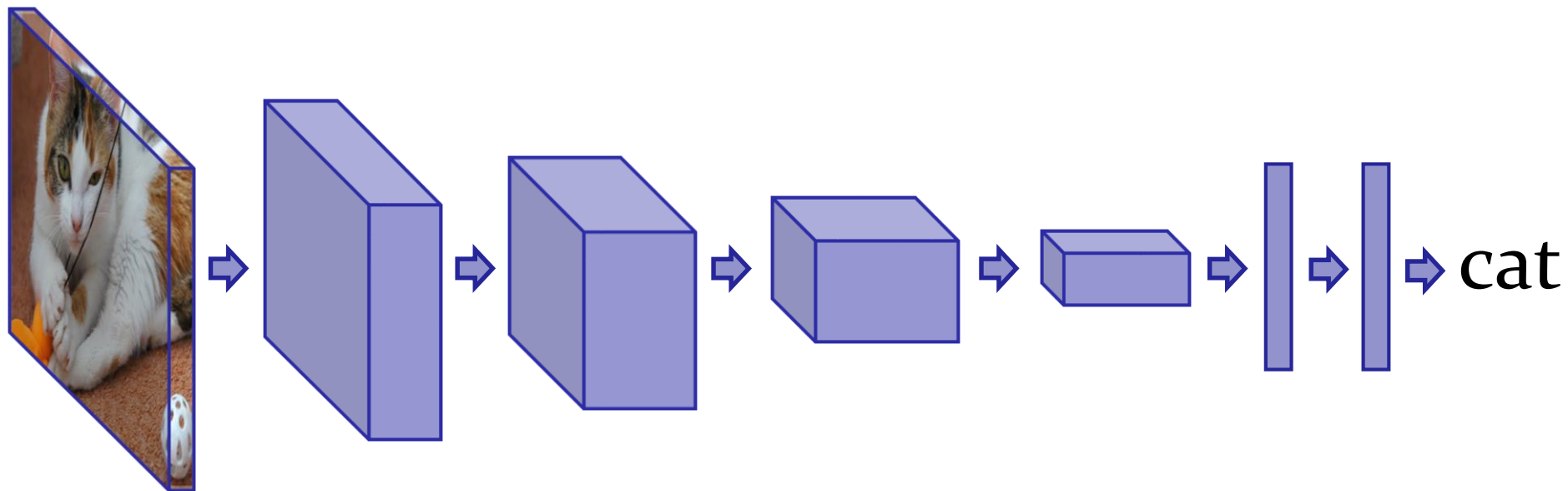
Alexey Dosovitskiy

University of Freiburg / Intel Labs

28.12.2016

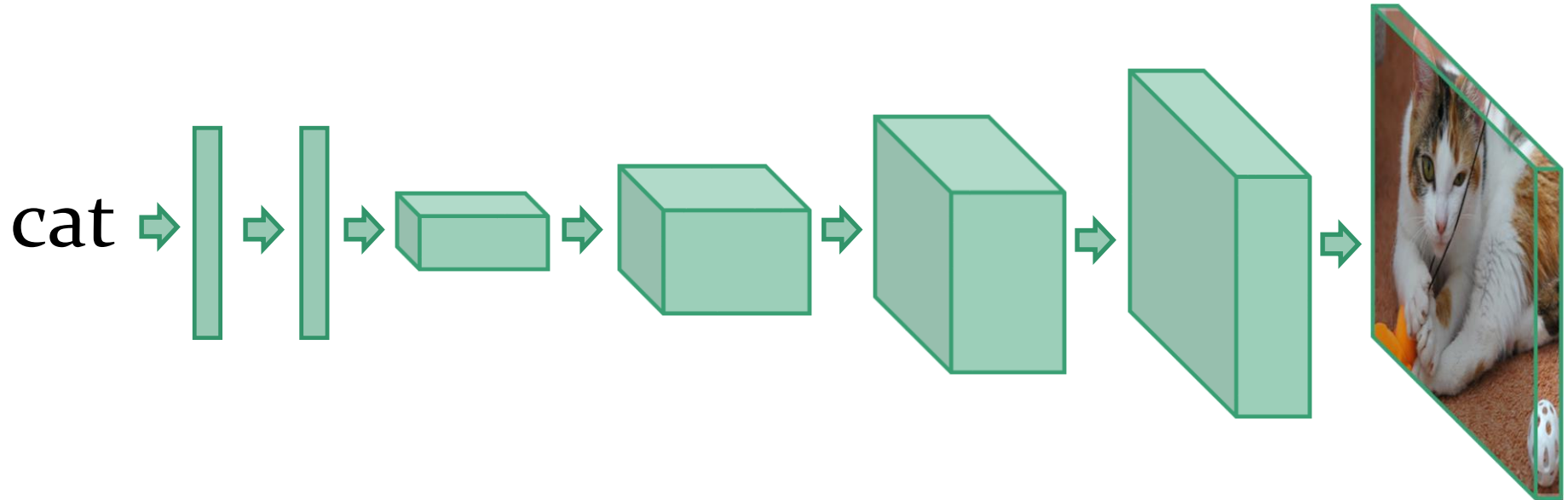


# Convolutional network



Convolutional network

# Up-convolutional network



Up-convolutional network  
(a.k.a. “deconvolutional”)

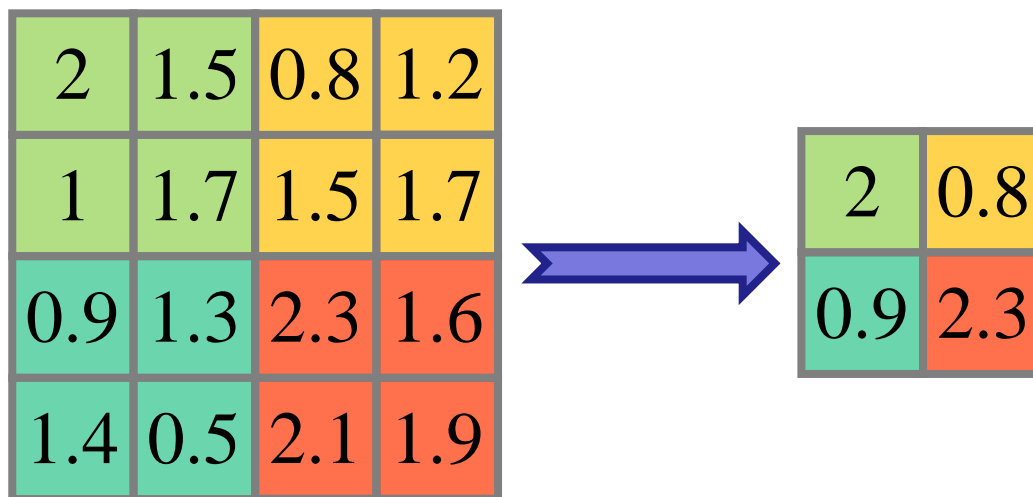
# This talk

- Up-convolutional networks
- End-to-end estimation of motion and depth
- Inverting ConvNets with perceptual metrics
- Visualizing neurons and generating images
- Learning to play Doom

# Pooling and unpooling

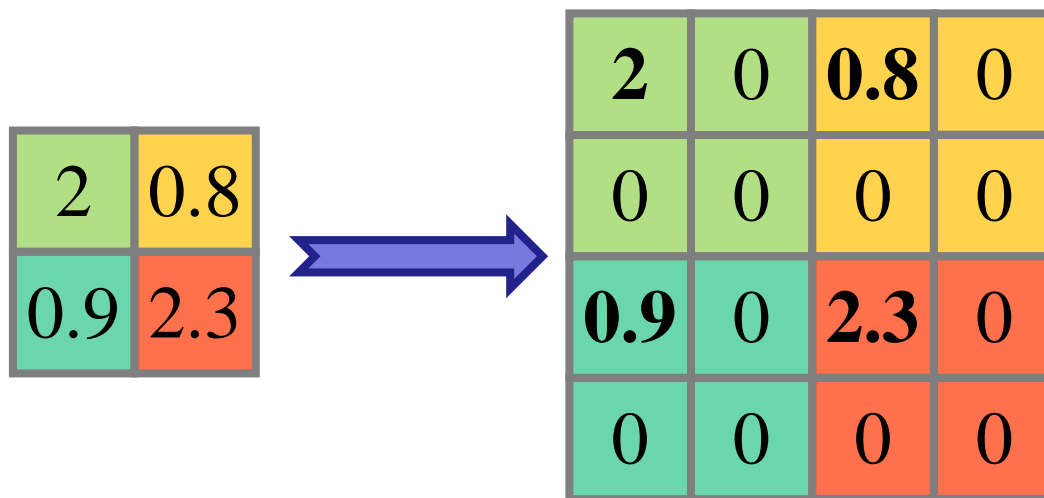
- Pooling = shrinking the feature maps

Convolution  
with stride

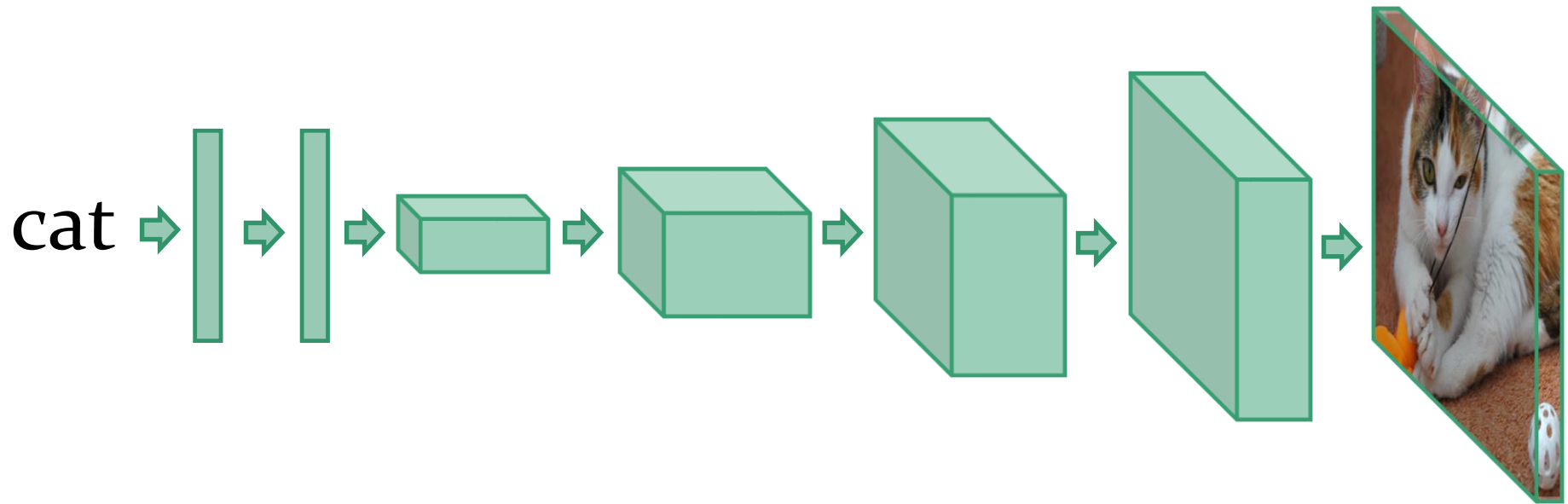


- Unpooling = expanding the feature maps (upsampling)

“Bed of nails”  
upsampling



# Up-convolutional network



Up-convolutional network  
(a.k.a. “deconvolutional”)

What can we do with this thing?

# End-to-end estimation of motion and depth

Joint work with the group of Daniel Cremers



Philipp  
Fischer



Eddy  
Ilg



Nikolaus  
Mayer



Philip  
Häusser



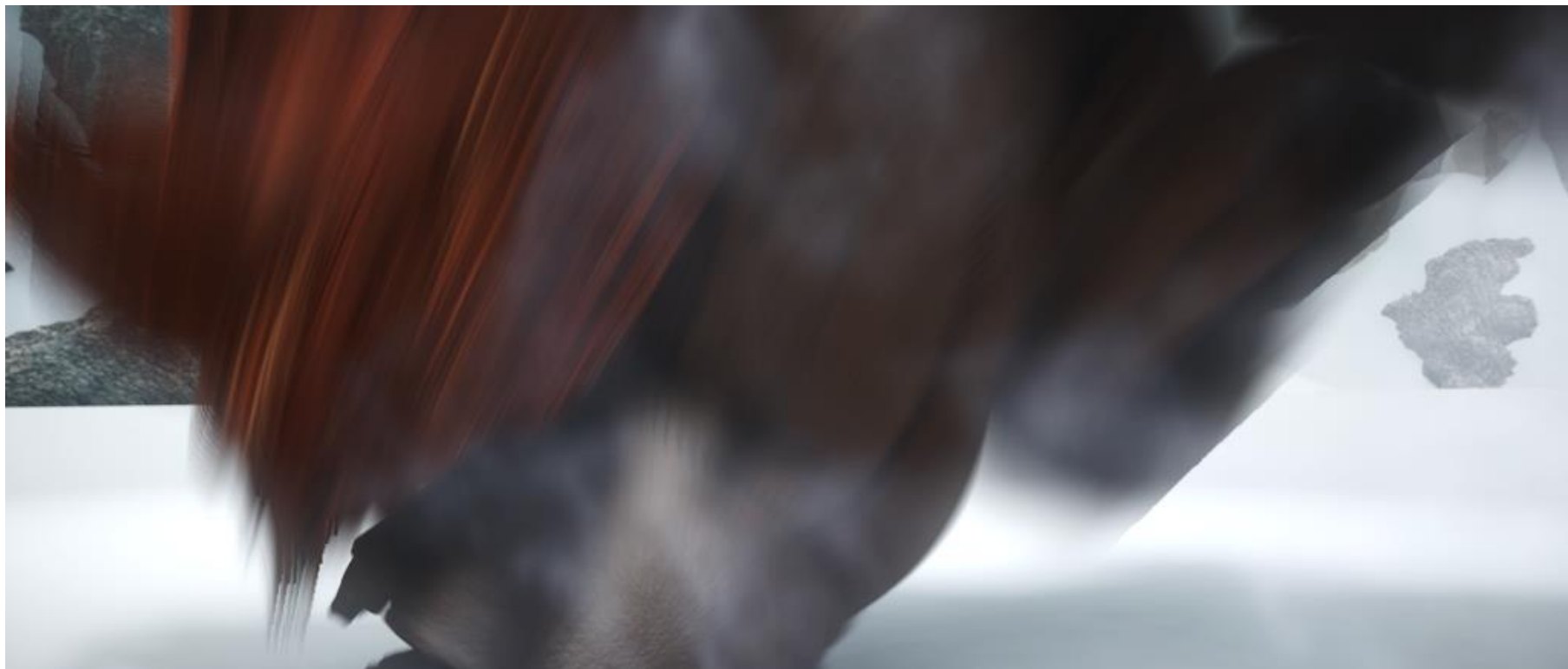
Caner  
Hazirbas



Vladimir  
Golkov

ICCV 2015, CVPR 2016, arxiv 2016

# Optical flow estimation is difficult!



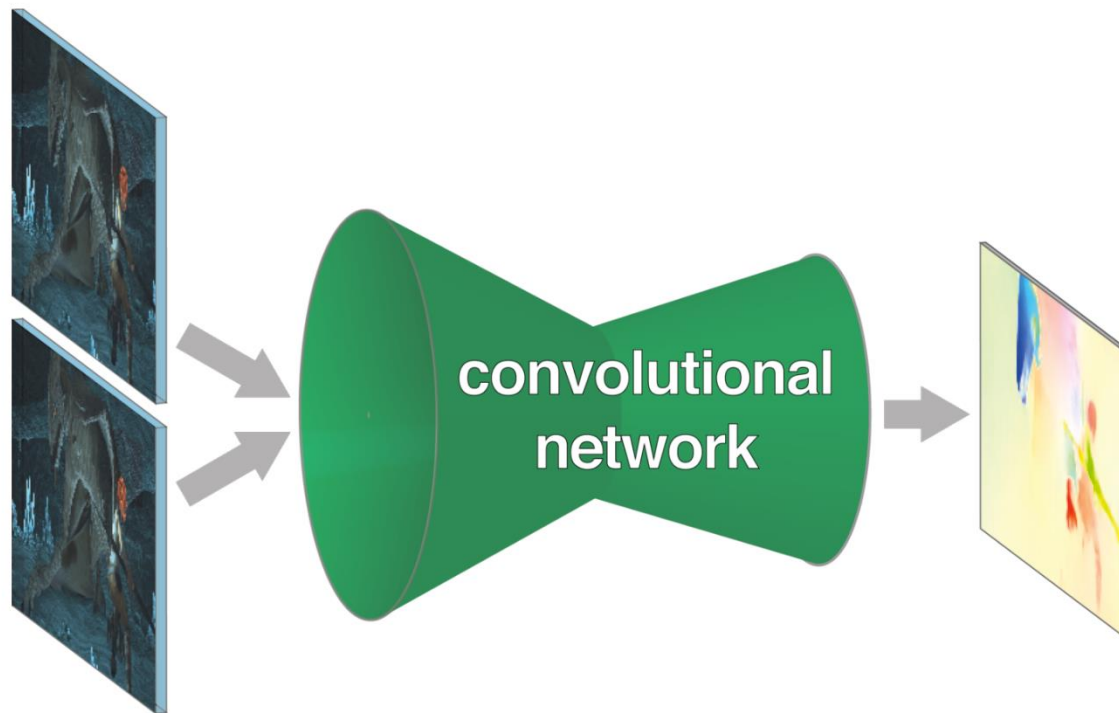


# Optical flow estimation is difficult!



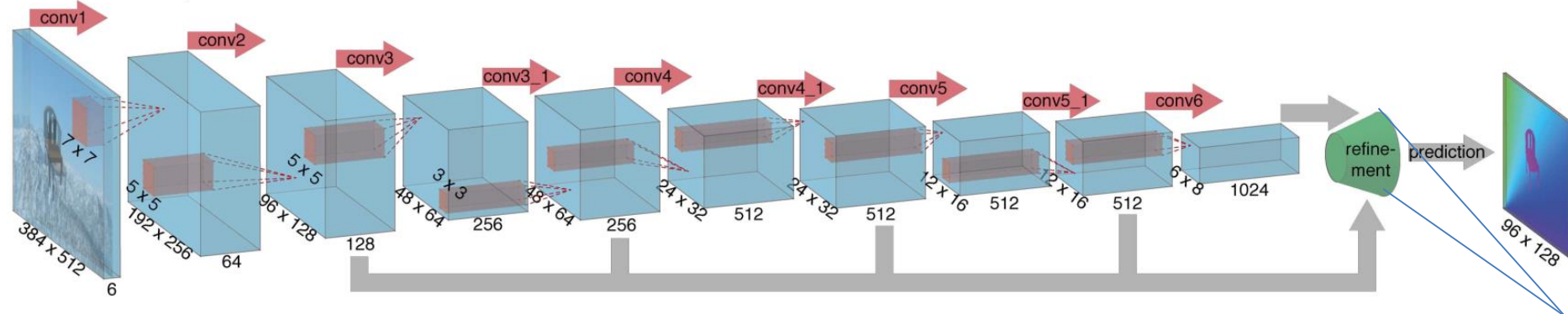
# End-to-end optical flow

- Standard approach: matching + aggregation
- End to end: two frames in, flow out

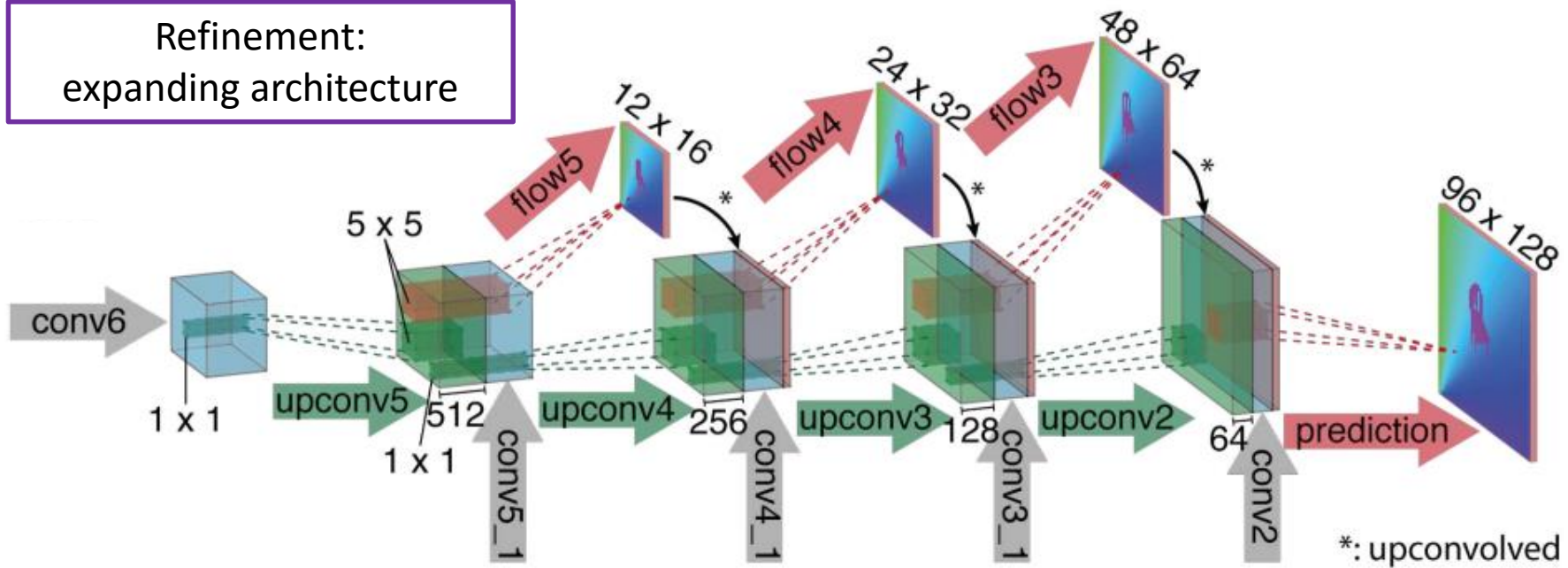


# Architecture: FlowNetSimple

FlowNetSimple

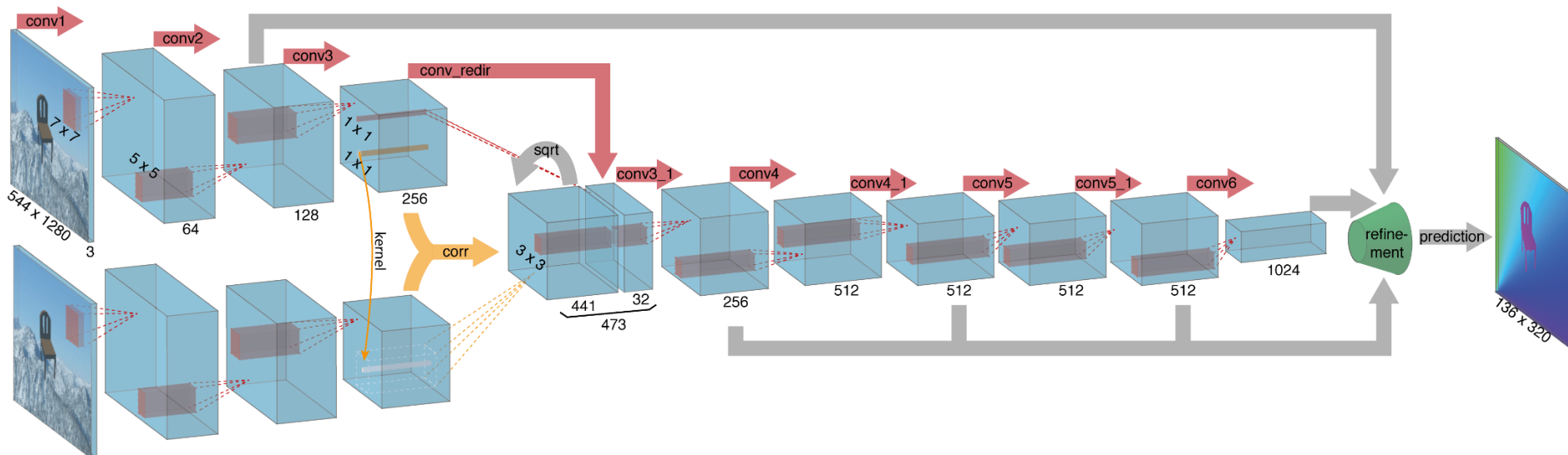


Refinement:  
expanding architecture



# Architecture: FlowNetCorr

FlowNetCorr

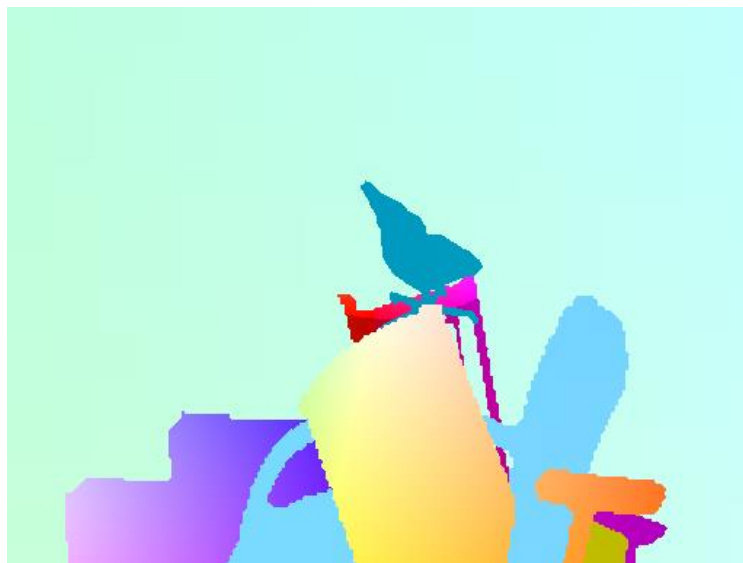


# How to train?

- Of course, supervised!
- Where do we get training data?



# The “Flying chairs” dataset



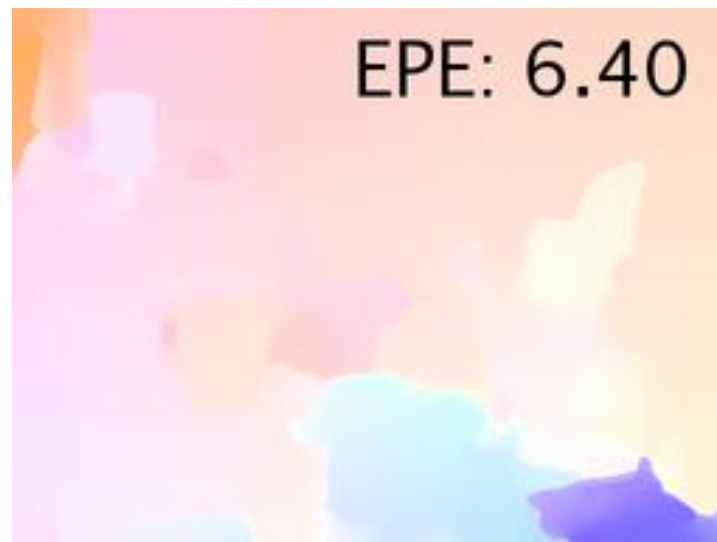
Rendered image

Optical flow

# It works on “Flying chairs” !



Input images



EpicFlow (Revaud et al. 2015)



Ground truth



FlowNetCorr

# And it generalizes!



Input images



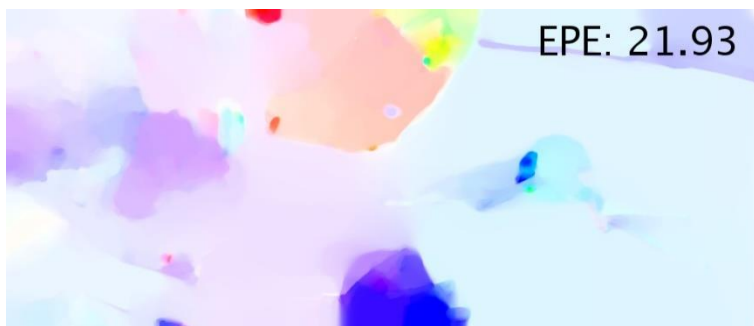
FlowNetSimple



Ground truth



FlowNetCorr



LDOF (Brox-Malik 2011)



EpicFlow (Revaud et al. 2015)

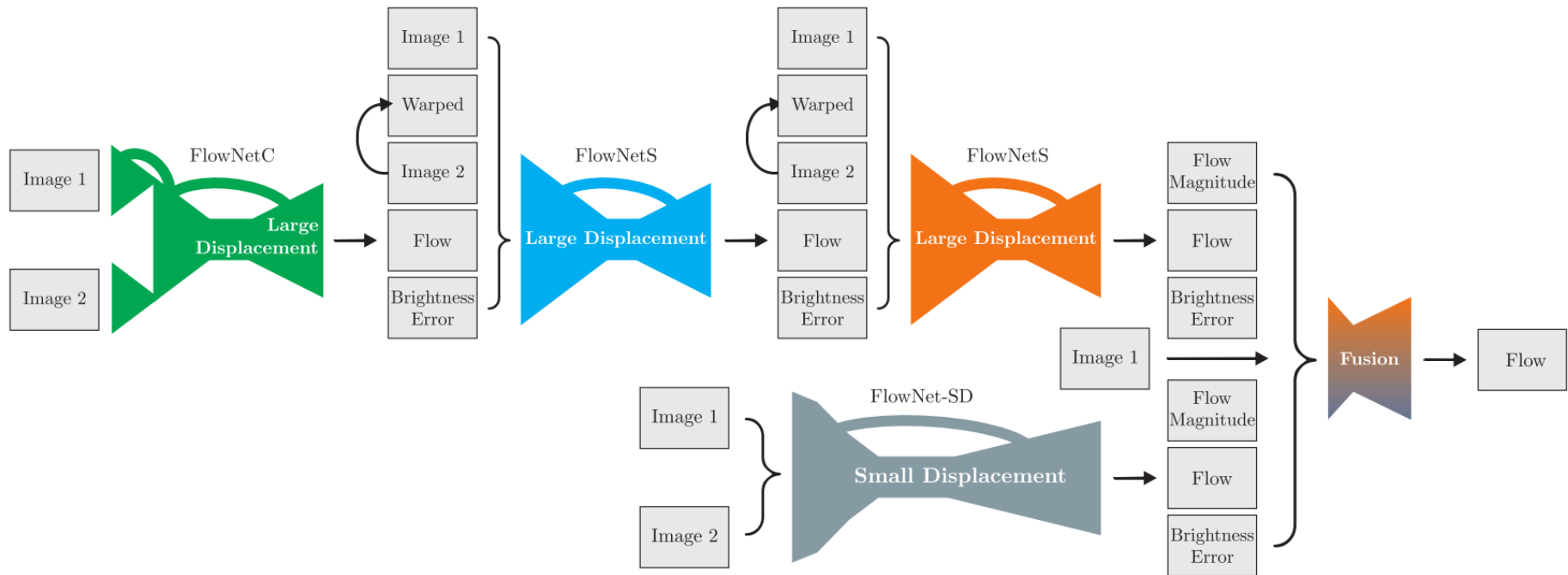


# Better dataset: FlyingThings3D

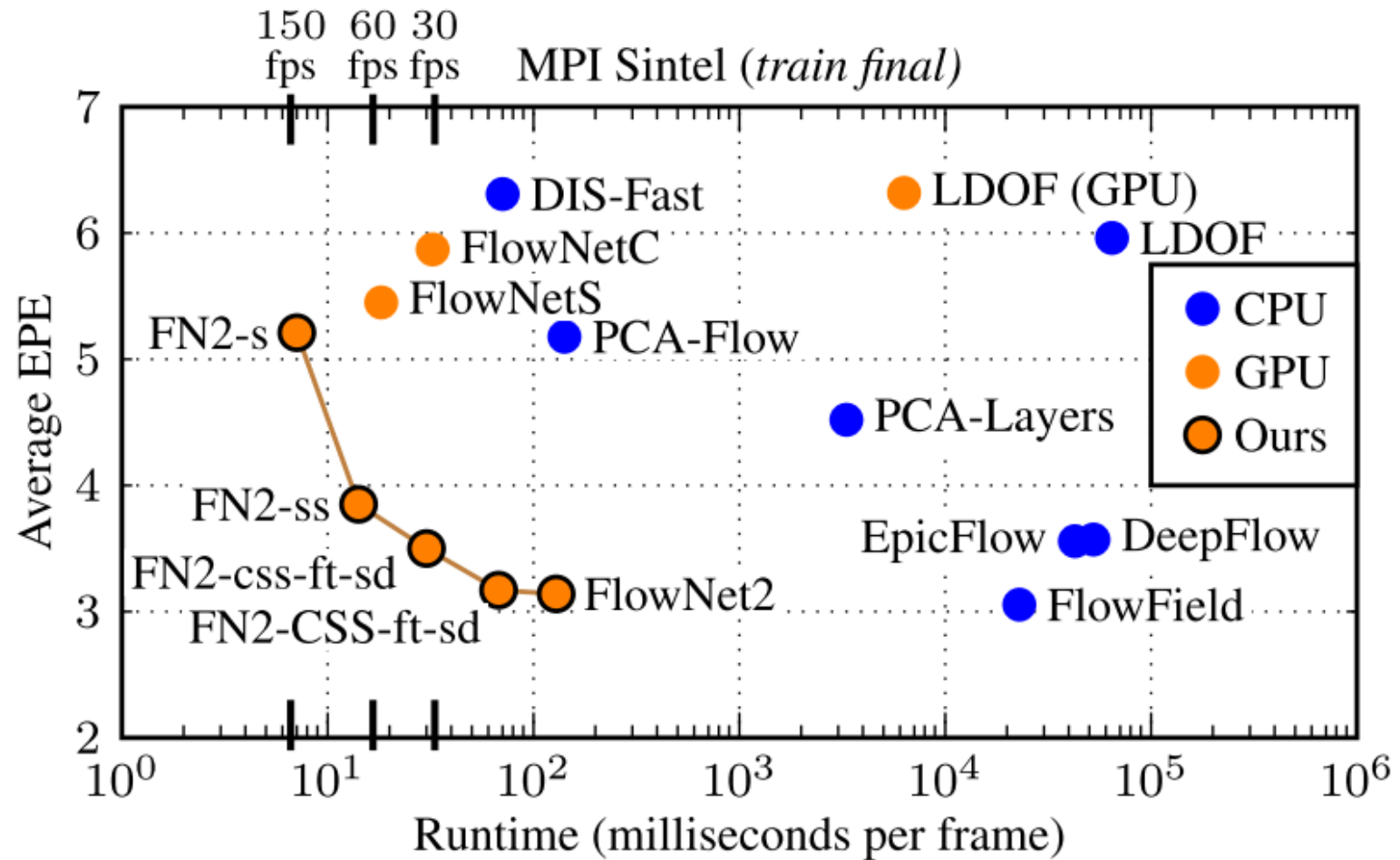


# Better architecture: FlowNet 2.0

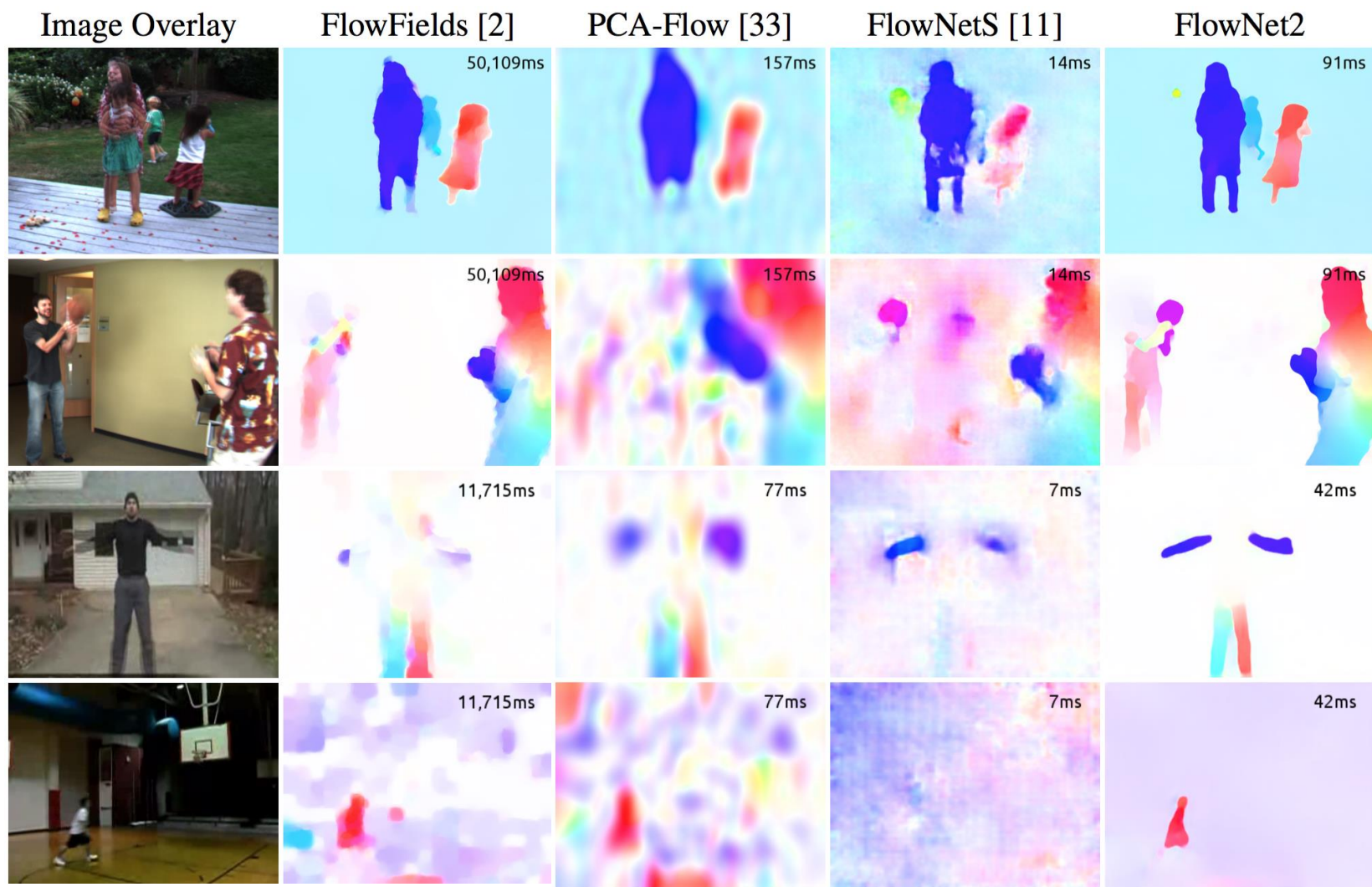
- Stacking and warping



# FlowNet 2.0: Sintel



# FlowNet 2.0: Real data



↑  
0.02-0.1 FPS

↖ ↗  
10-20 FPS



# FlowNet 2.0: KITTI

Image Overlay



PCA-Flow [33]



Ground Truth



FlowNetS [11]



FlowFields [2]



FlowNet2-kitti



Image Overlay



PCA-Flow [33]



Ground Truth



FlowNetS [11]













FlowFields [2]



FlowNet2-kitti



# FlowNet 2.0: KITTI

	Method	Setting	Code	Fl-bg	Fl-fg	<b>Fl-all</b>	Density	Runtime
1	<a href="#">PRSM</a>	 	<a href="#">code</a>	5.33 %	17.02 %	7.28 %	100.00 %	300 s
C. Vogel, K. Schindler and S. Roth: <a href="#">3D Scene Flow Estimation with a Piecewise Rigid Scene Model</a> . ijcv 2015.								
2	<a href="#">OSF+TC</a>	 		5.76 %	16.61 %	7.57 %	100.00 %	50 min
3	<a href="#">OSF</a>		<a href="#">code</a>	5.62 %	22.17 %	8.37 %	100.00 %	50 min
M. Menze and A. Geiger: <a href="#">Object Scene Flow for Autonomous Vehicles</a> . Conference on Computer Vision and Pattern Recognition								
4	<a href="#">SSFAV</a>			7.10 %	21.22 %	9.45 %	100.00 %	5 min
5	<a href="#">FlowNet2</a>			10.75 %	15.14 %	11.48 %	100.00 %	0.12 s
6	<a href="#">SDF</a>			8.61 %	26.69 %	11.62 %	100.00 %	TBA
M. Bai*, W. Luo*, K. Kundu and R. Urtasun: <a href="#">Exploiting Semantic Information and Deep Matching for Optical Flow</a> . ECCV 2016.								
7	<a href="#">FSF+MS</a>	  		8.48 %	29.62 %	12.00 %	100.00 %	2.7 s
8	<a href="#">CNNF+PMBP</a>			10.08 %	23.18 %	12.26 %	100.00 %	45 min
9	<a href="#">CSF</a>			10.40 %	30.33 %	13.71 %	100.00 %	80 s
Z. Lv, C. Beall, P. Alcantarilla, F. Li, Z. Kira and F. Dellaert: <a href="#">A Continuous Optimization Approach for Efficient and Accurate Scene Flow Estimation</a> . CVPR 2017.								

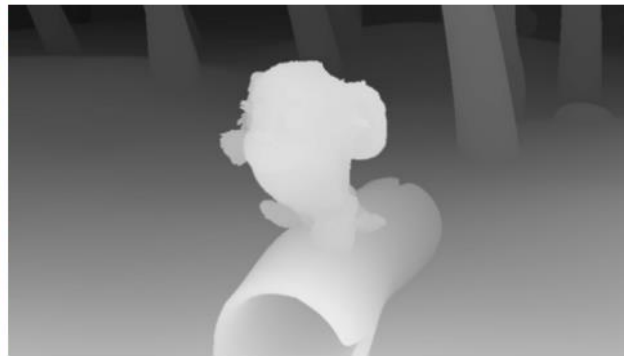
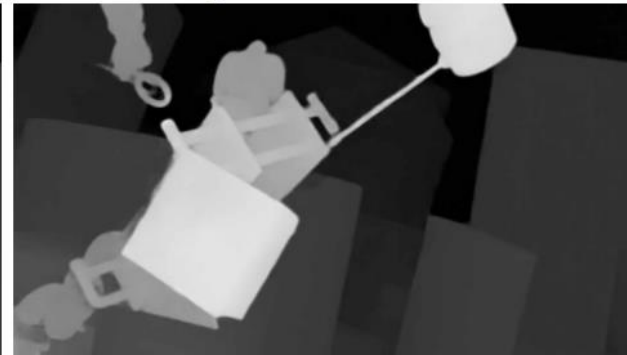
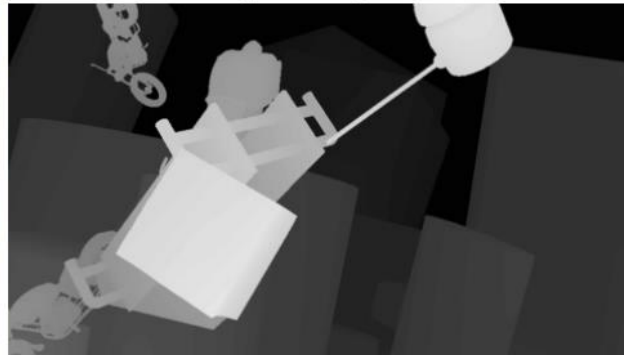
Depth estimation!

# Disparity estimation

RGB image (L)

Disparity GT

DispNetCorr1D





# Disparity: KITTI

RGB image (L)



DispNetCorr1D-K



SGM prediction

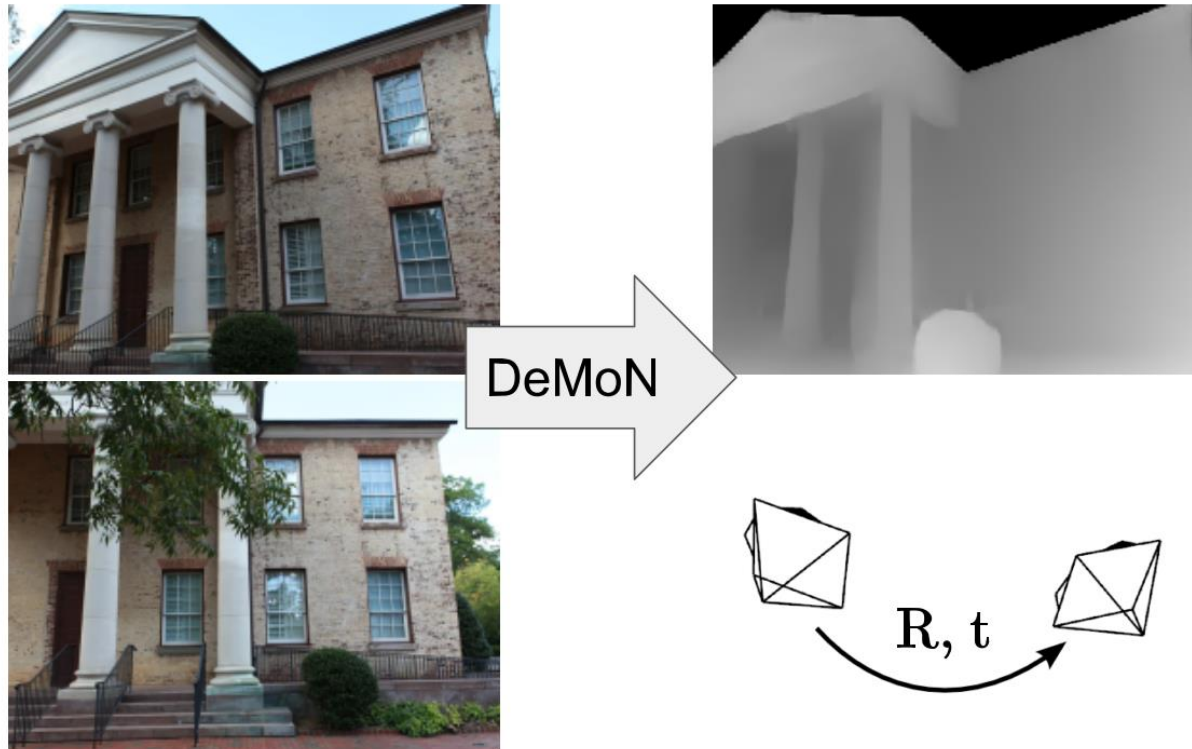


MC-CNN prediction

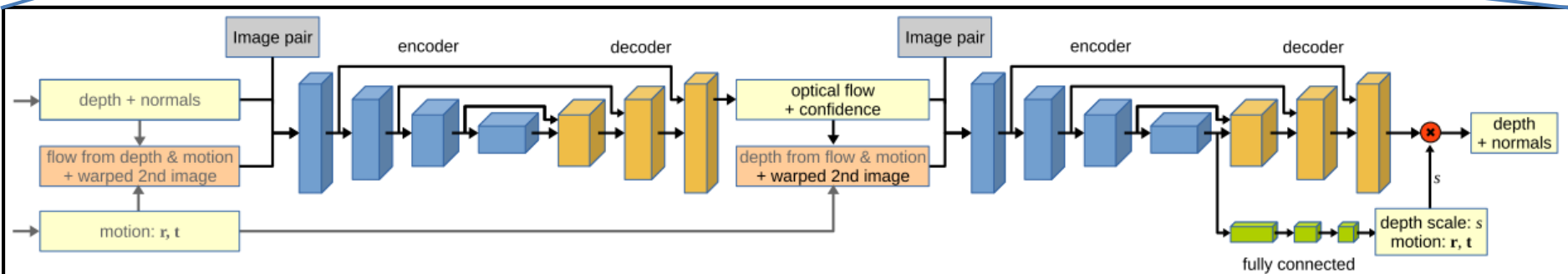
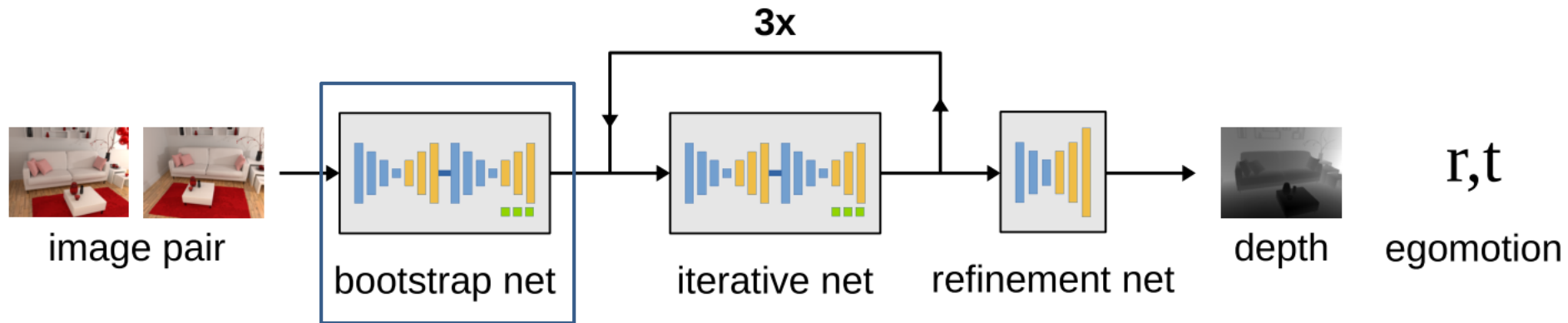


# DeMoN: monocular stereo

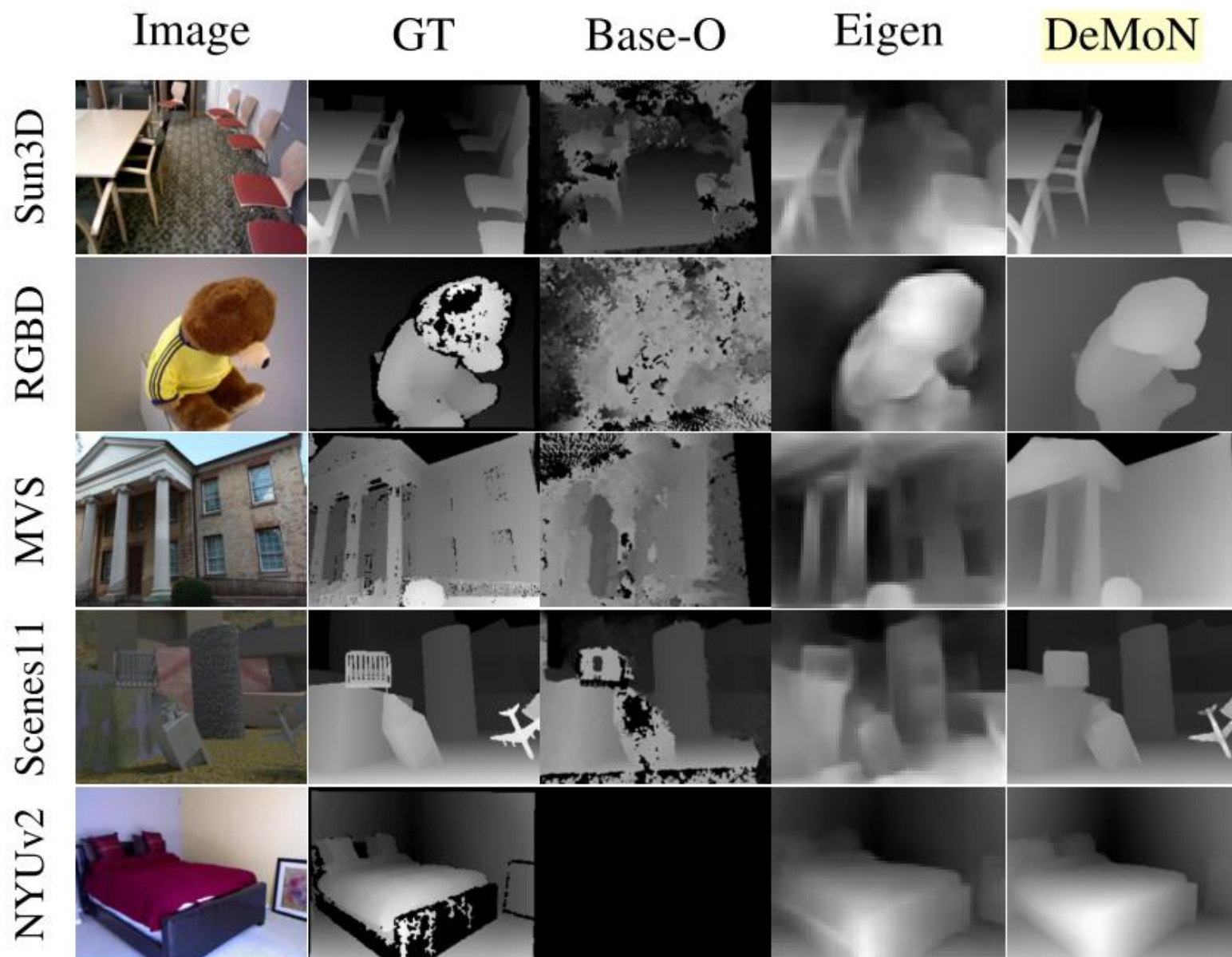
- **D**epth and **M**otion **N**etwork
  - Two frames in
  - Depth and camera motion out



# DeMoN: architecture



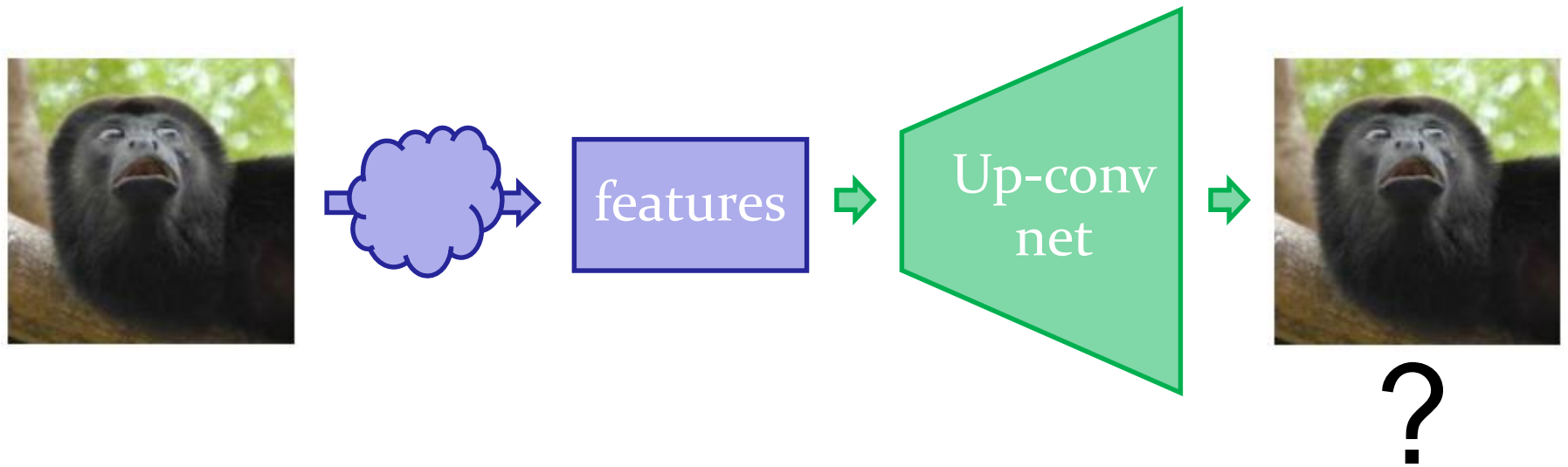
# DeMoN: architecture



- (Up)ConvNets can estimate motion and depth end to end
- State-of-the-art performance at interactive framerates

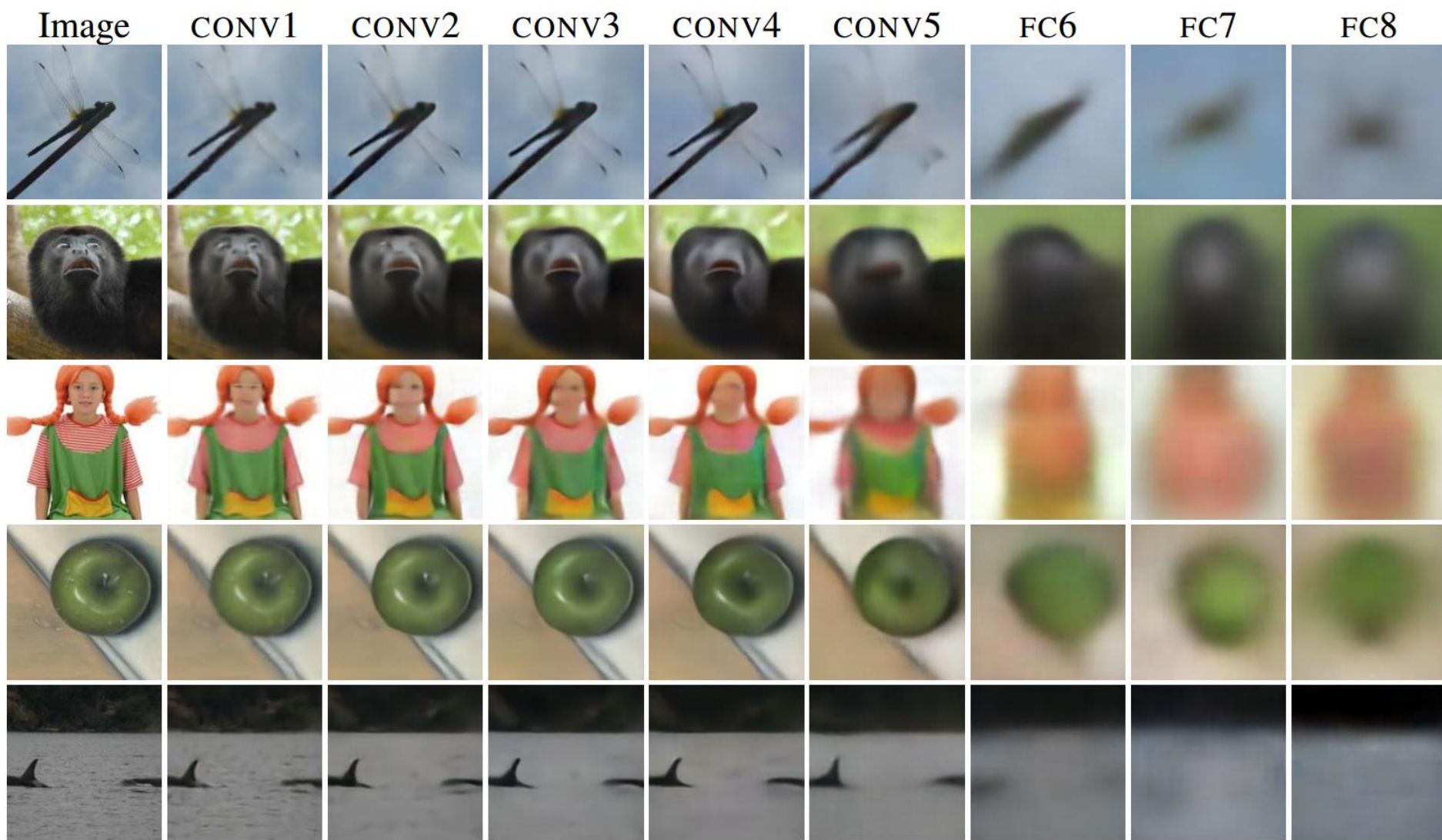
# Inverting ConvNets with perceptual metrics

# Inverting representations with ConvNets



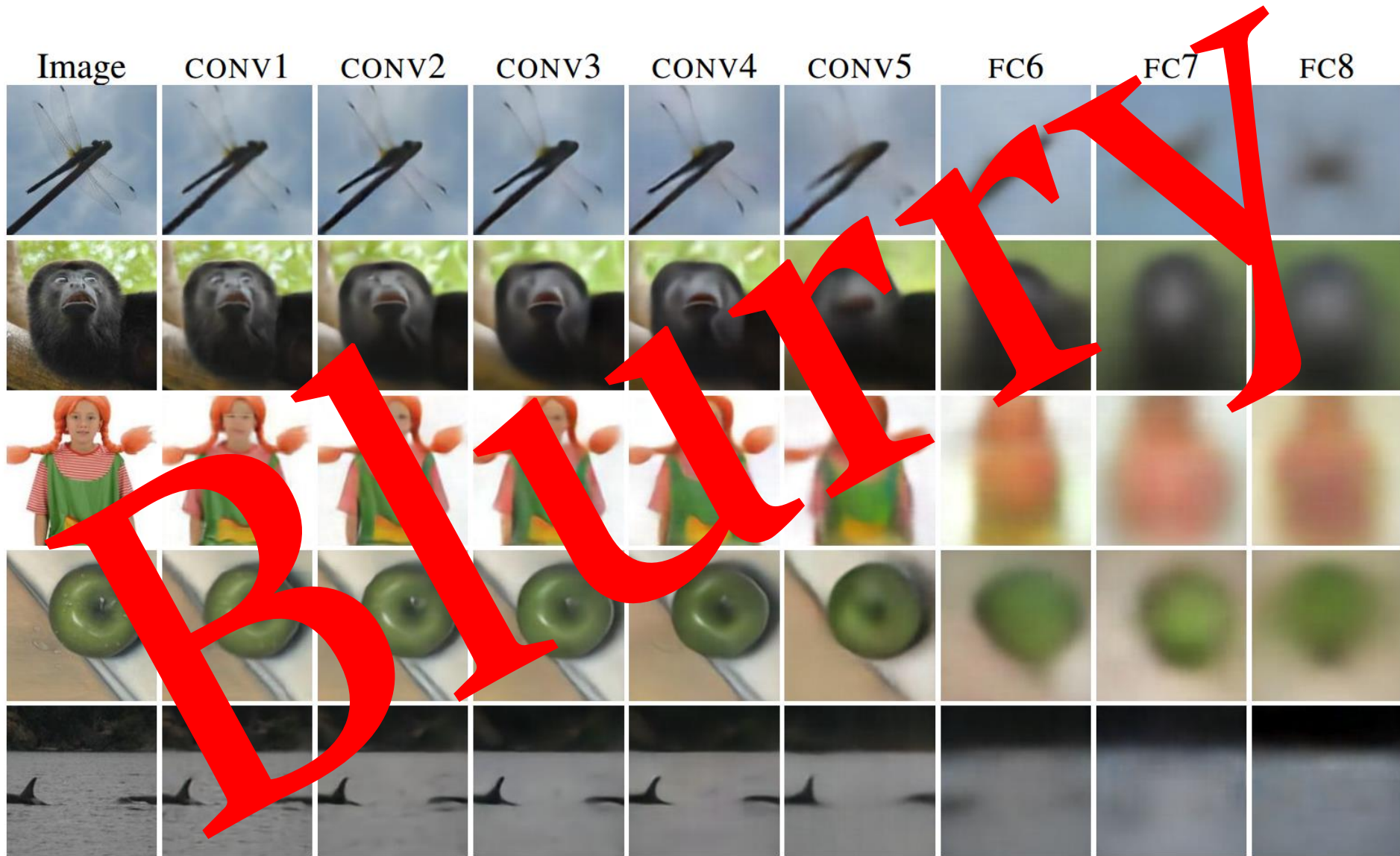


# Inverting AlexNet





# Inverting AlexNet



# Why so blurry?

- Problem: the feature vector does not contain the precise locations of all details
- Solution: with appropriate loss function it need not!

# Deep perceptual similarity metric

- Want to be sensitive to important properties, but invariant to irrelevant deformations
- Instead of the image space, measure image similarity in the feature space
- Add adversarial loss as a natural image prior

Original      Img loss

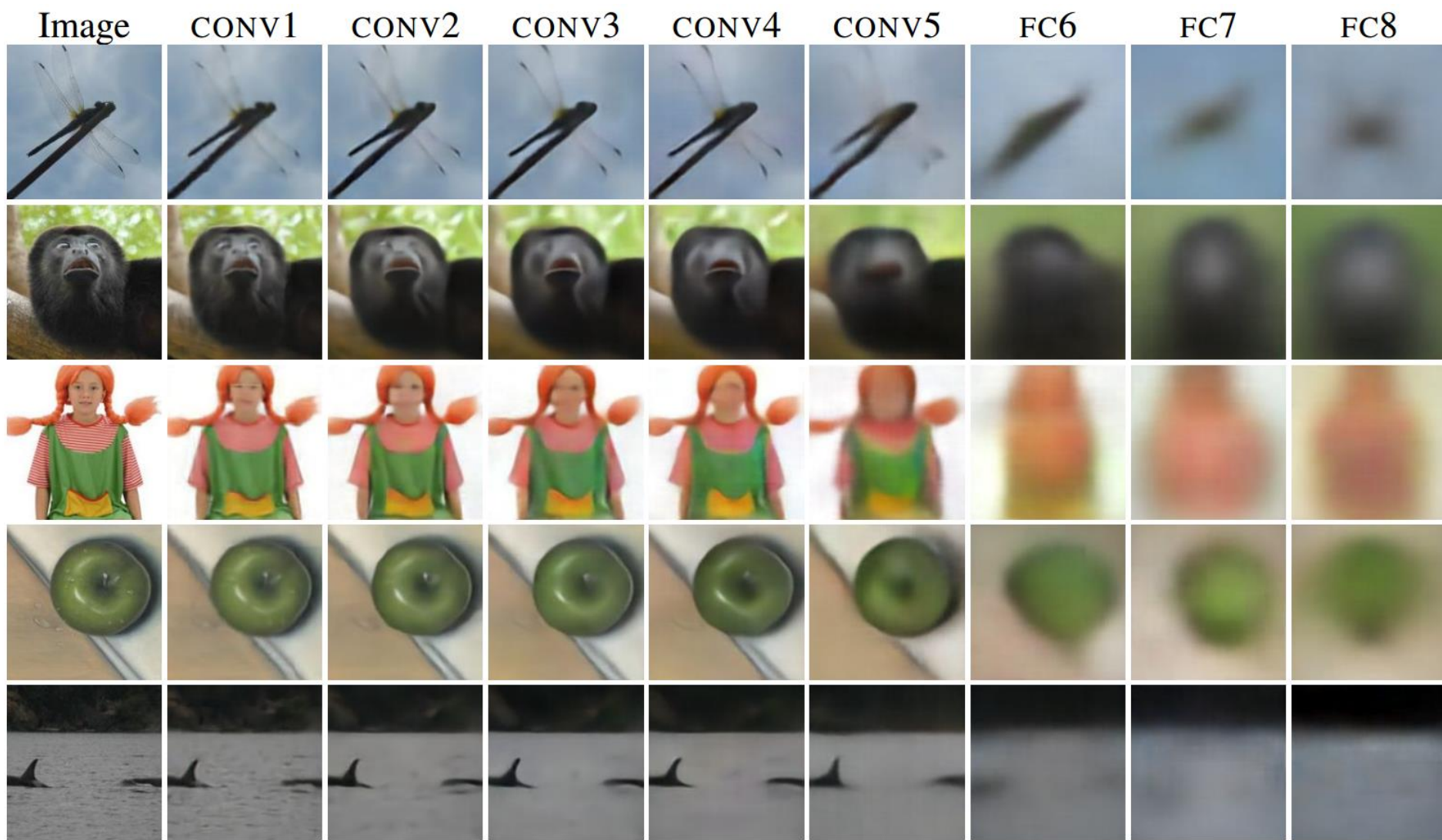


# Deep perceptual similarity metric

- Want to be sensitive to important properties, but invariant to irrelevant deformations
- Instead of the image space, measure image similarity in the feature space
- Add adversarial loss as a natural image prior

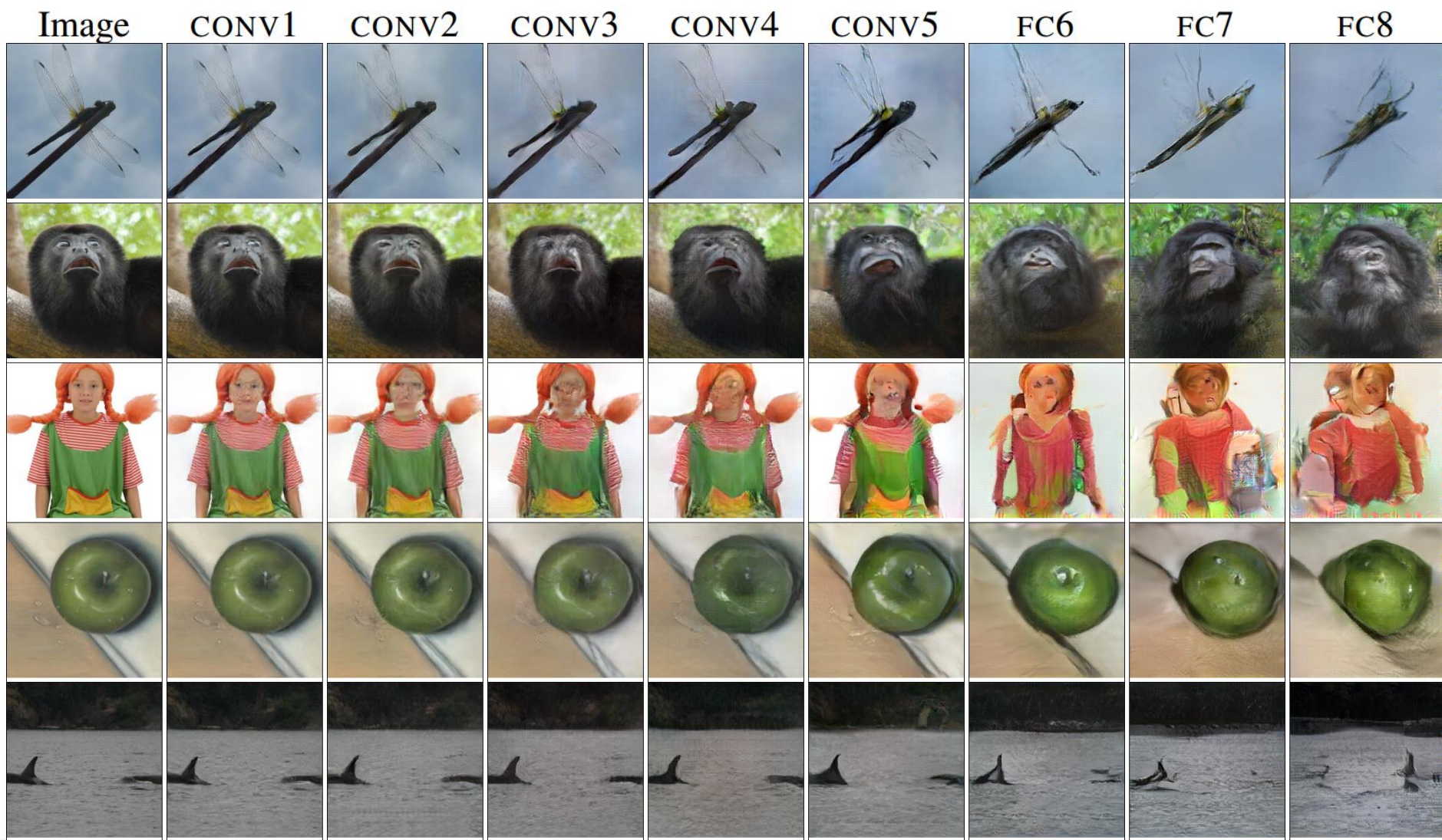


# Inverting AlexNet: Euclidean loss



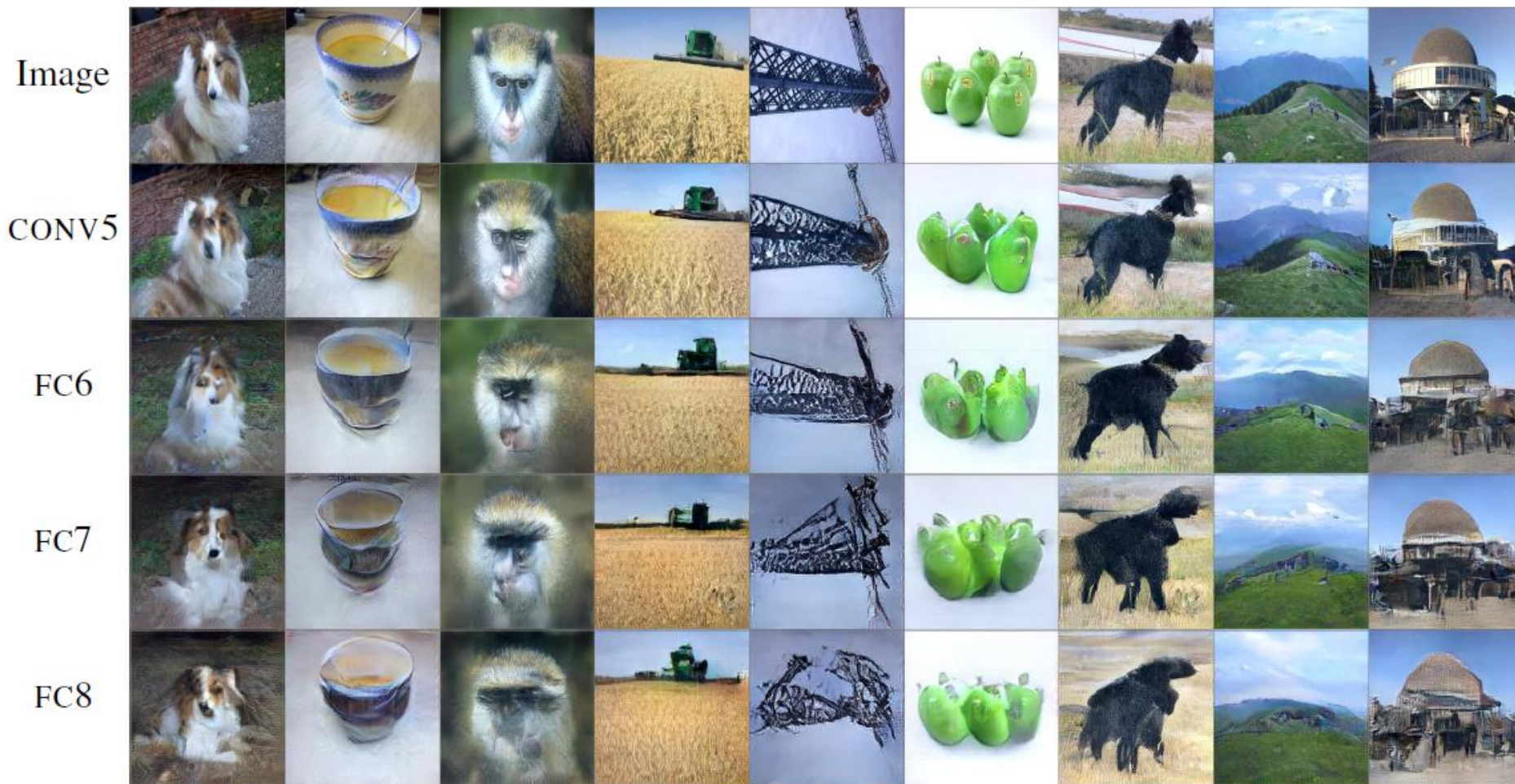


# Inverting AlexNet: DeePSiM loss





# Inverting AlexNet: more results



- Superresolution [Johnson et al. 2016], [Ledig et al. 2016]
- Image compression
- Denoising
- Analysis of deep networks
- Generative models



# Visualizing neurons and generating images



Anh  
Nguyen

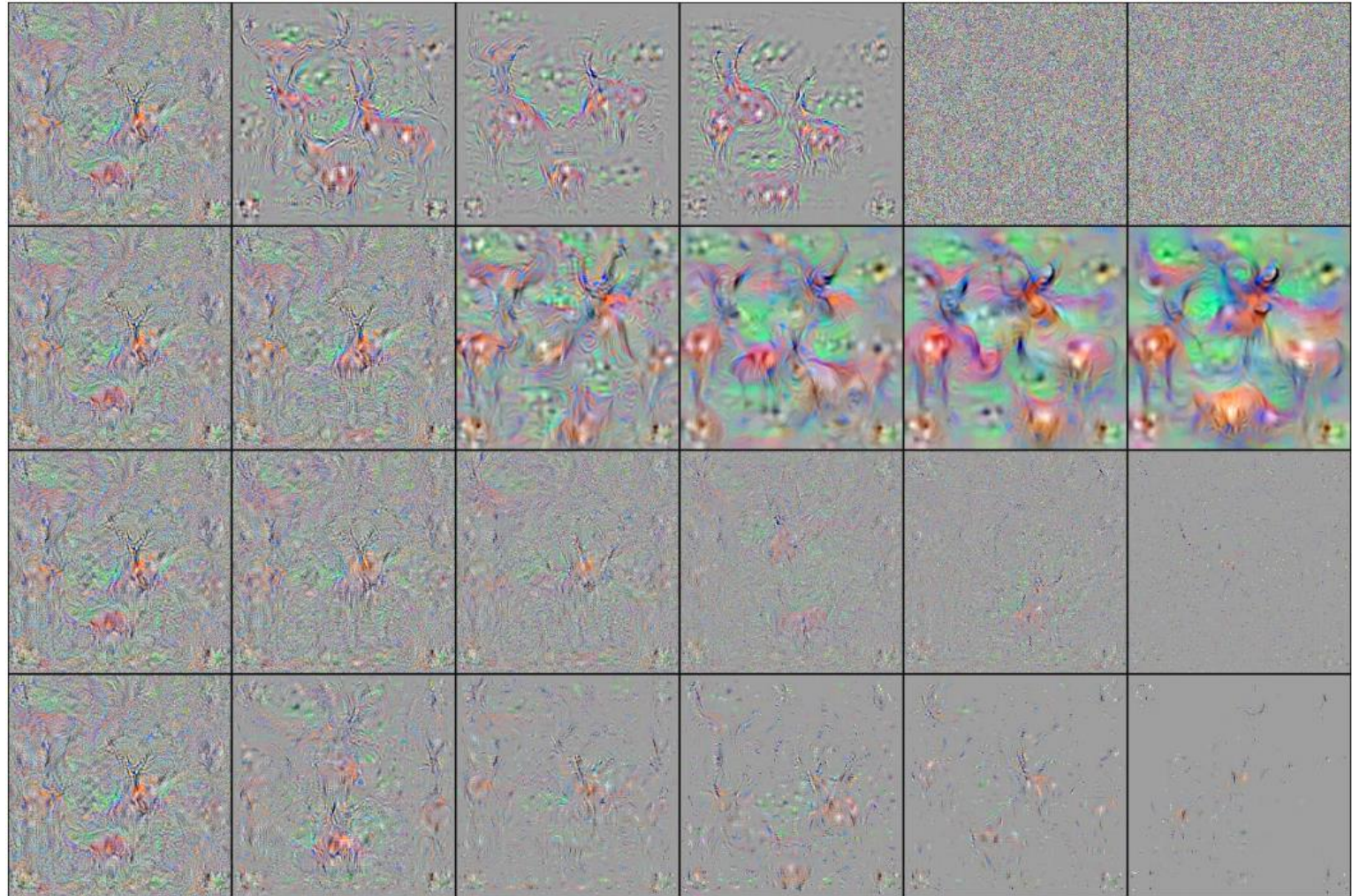


Jason  
Yosinski



Jeff  
Clune

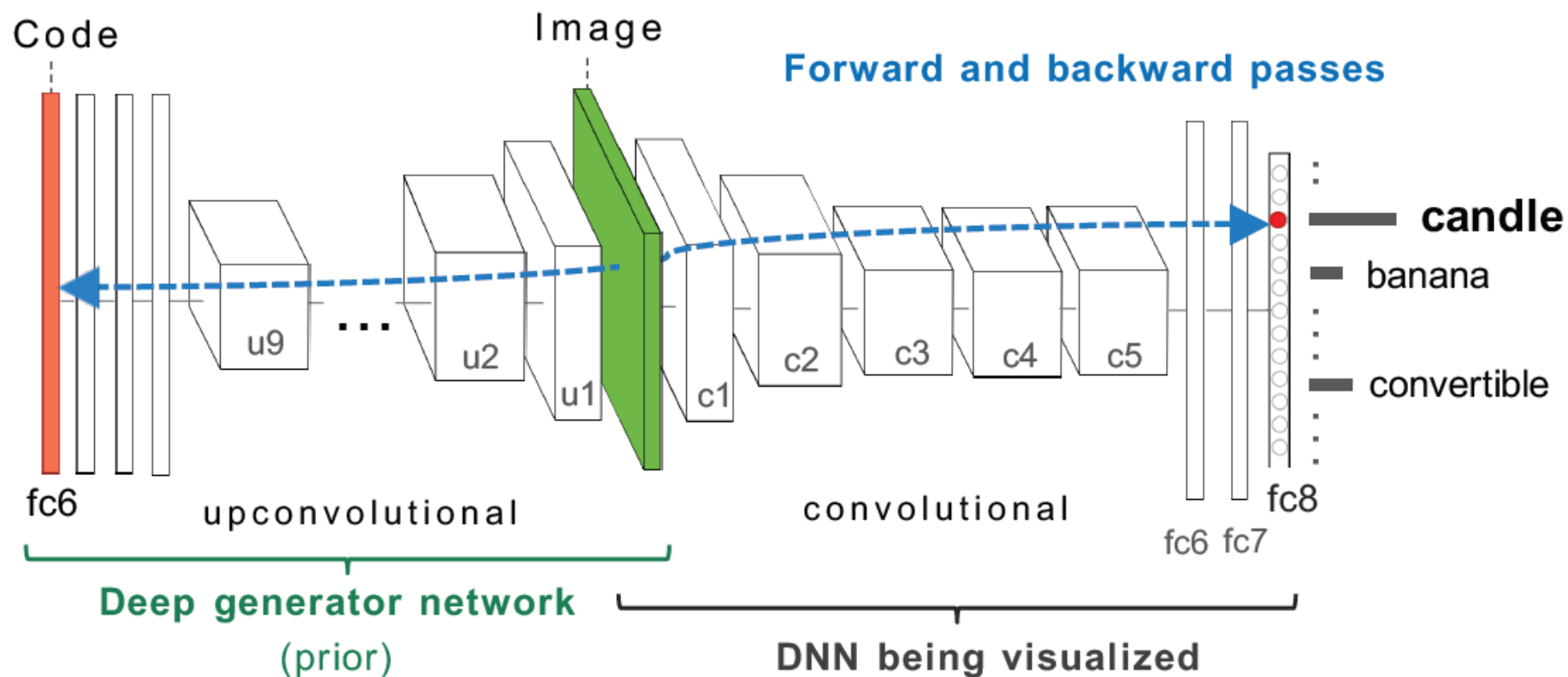
# Activation maximization



Yosinski et al. 2015

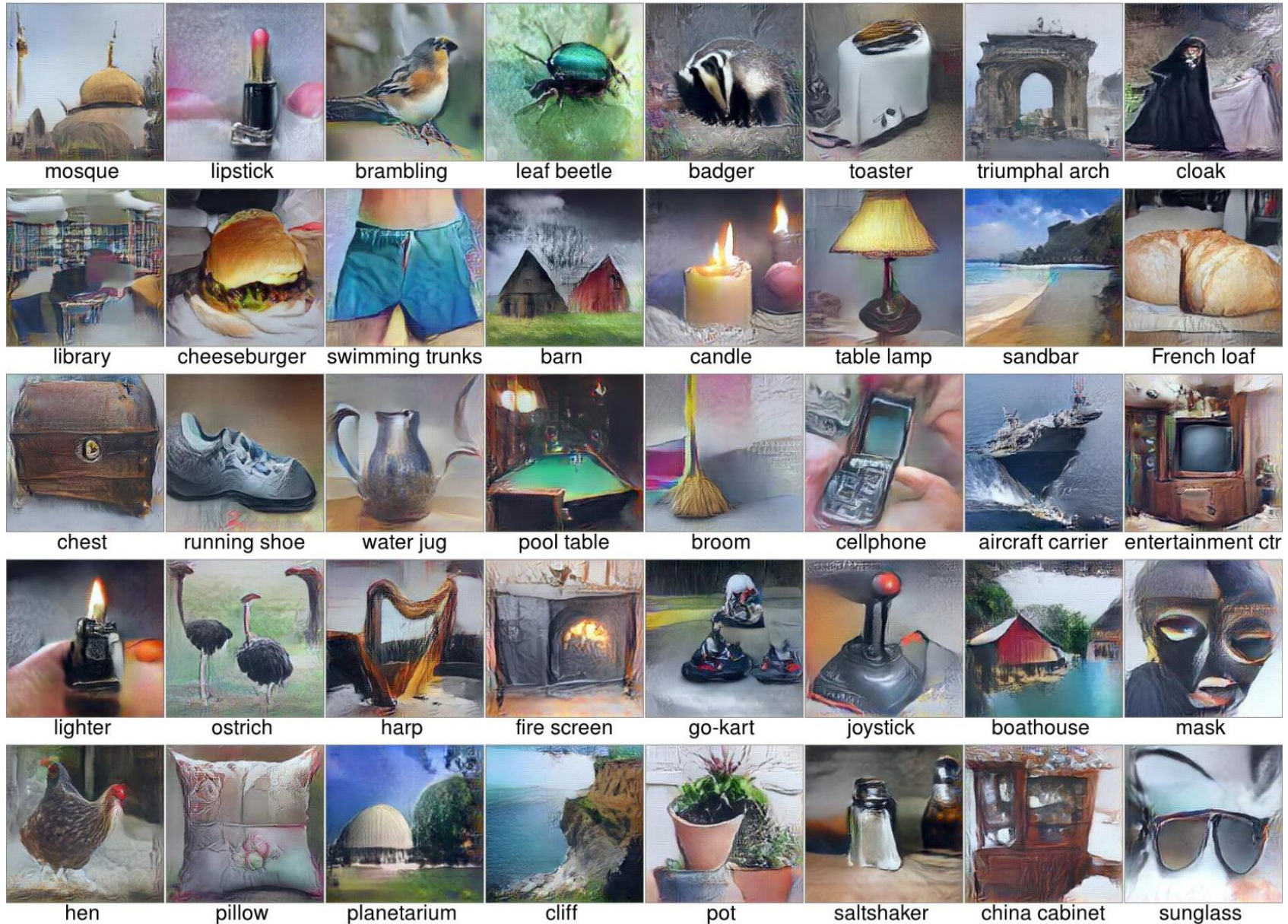
Hand-crafted priors are not good enough

# AM with an UpConvNet prior





# Activating FC8 neurons: ImageNet





# Activating FC8 neurons: Places

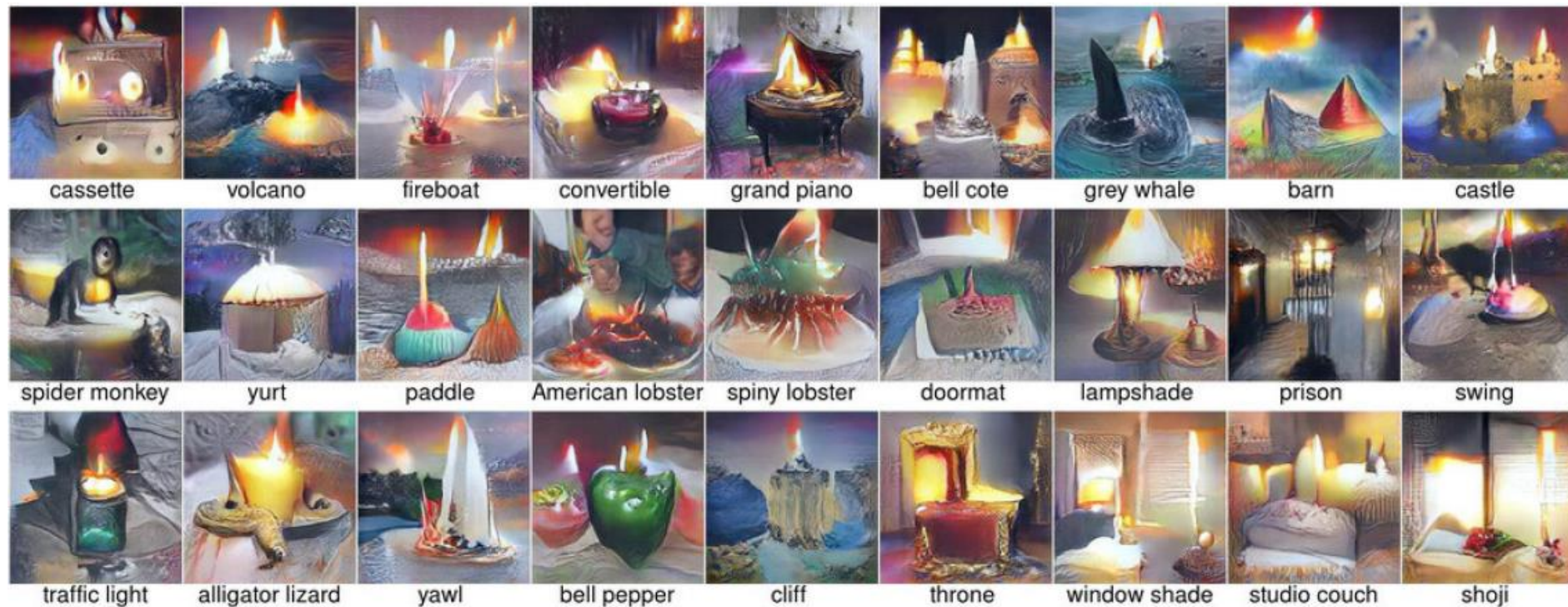


# Activating FC8 neurons: 2 classes





# Activating FC8 neurons: 2 classes





# Activating neurons from different layers

Layer 8



auditorium

food court

doorway/outdoor

conference room

igloo

Layer 7



Layer 6



# Generative model?

- Pictures are nice, but is it a real generative model?
- If we add some noise during optimization, it is!



# Plug-and-play generative networks

(a) Real: top 9



(c) Real: random 9



# Plug-and-play generative networks

(a) Real: top 9



(b) DGN-AM [36]



(c) Real: random 9



(d) PPGN (this)





# Plug-and-play generative networks



redshank

ant

monastery



volcano



# PPGN: sentence to image



a red car parked on the side of a road

a blue car parked on the side of a road



a pizza on a plate at a restaurant

someone is just about to cut the pizza



oranges on a table next to a liquor bottle

a pile of oranges sitting in a wooden crate

- Perceptual metrics for better image generation
- ConvNets are surprisingly invertible
- Plug-and-play generative networks produce great high resolution images

# Sensorimotor control (learning to play Doom)



Vladlen  
Koltun

arxiv 2017

# Motivation

Reinforcement learning

Single goal

Scalar reward

Maximize returns

Real life

Multitude of goals

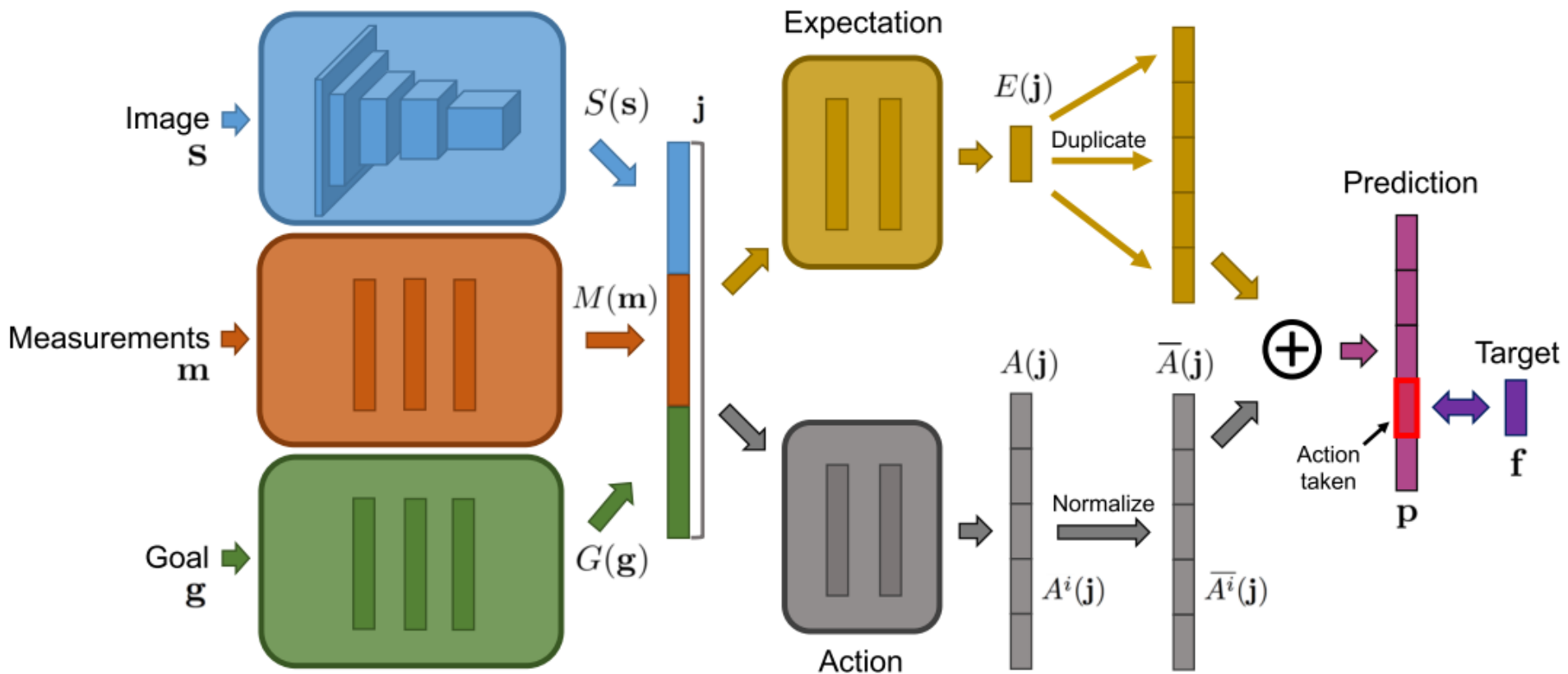
Rich sensory stream

Explore the world

# Reinforcement learning vs real life

- Idea: predict future *measurements*
  - Hunger, pain, cold, sleep
  - Health, ammo, frags
- Formulate goals in terms of these
  - Minimize hunger, pain, cold, sleepiness
  - Maximize health, ammo, frags
- Predict with simple supervised training

# Architecture





# ViZDoom tasks



D1: Basic



D2: Navigation

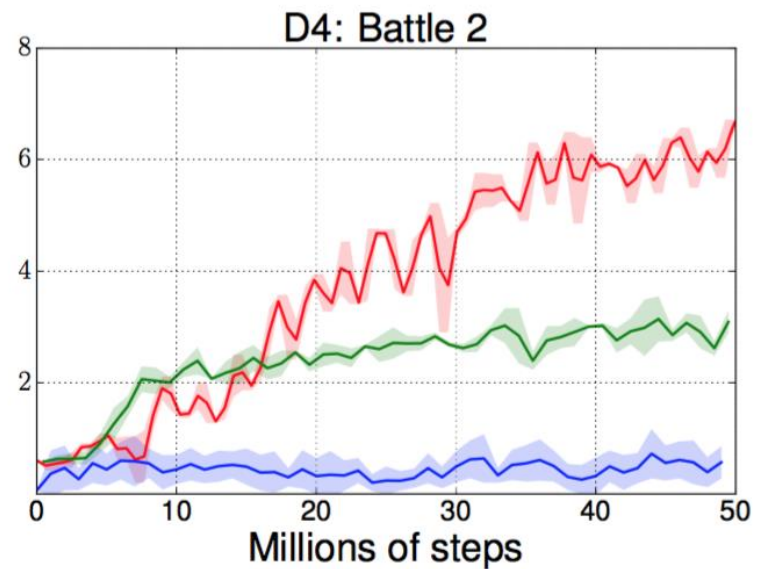
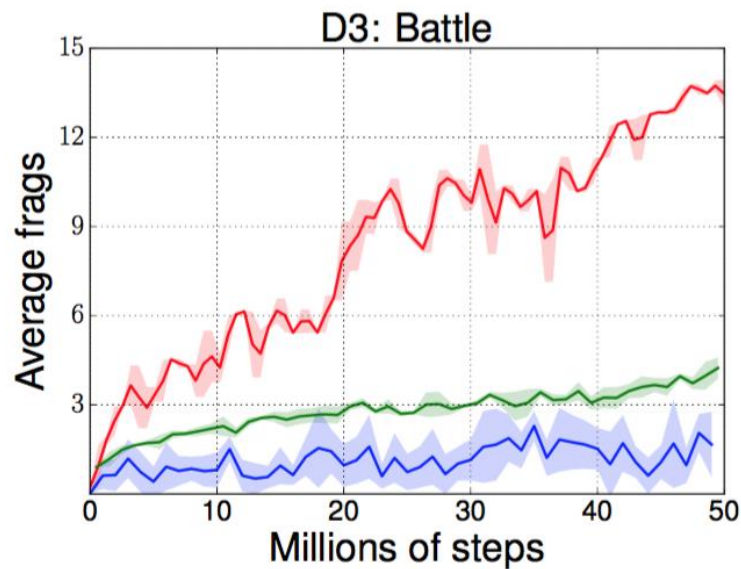
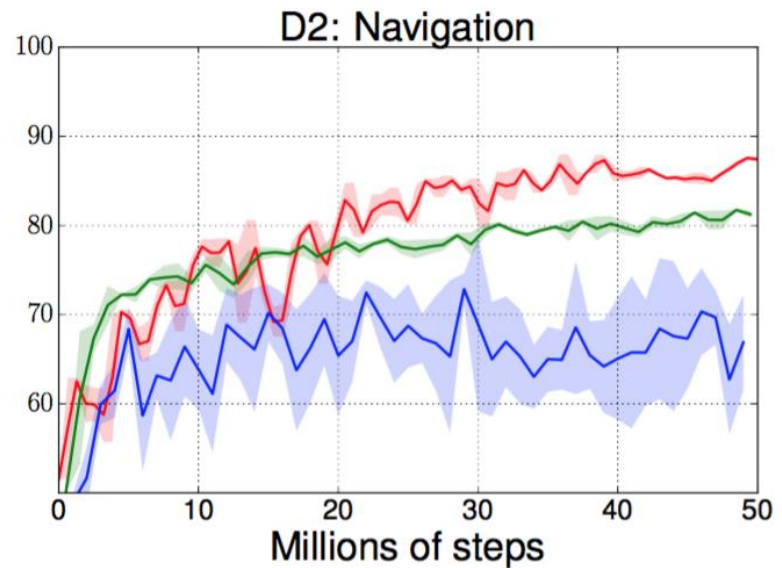
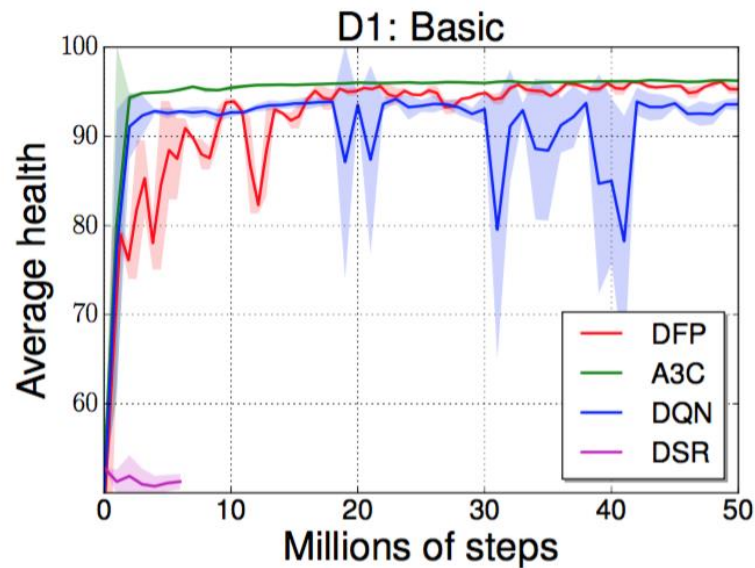


D3: Battle



D4: Battle 2

# ViZDoom results



## Learning to Act by Predicting the Future

Alexey Dosovitskiy    Vladlen Koltun

# ViZDoom competition

Place	Team	1	2	3	4	5	6	7	8	9	10	11	12	Total
1	IntelAct	29	21	23	21	6	11	9	6	30	32	33	35	256
2	The Terminators	22	17	21	15	13	12	6	5	14	13	13	13	164
3	TUHO	8	11	13	12	0	-1	-1	-4	2	2	6	3	51
4	ColbyCS	2	4	0	1	-1	0	-1	0	3	3	4	3	18
5	5vision	3	0	4	2	1	0	1	0	0	-1	1	1	12
6	Ivomi	3	0	1	0	1	-1	-4	-4	1	1	0	0	-2
7	PotatoesArePrettyOk	0	0	2	0	-1	-3	-1	0	-2	-1	-1	-2	-9

# What next?

- Deep learning and simulation
- Learning models of environments
  - Future prediction
  - Planning
  - Analysis by synthesis
- Coupled perception and control

Looking for interns! (Munich, Santa Clara)



# Summary



End-to-end motion and depth estimation



Inverting ConvNets and perceptual metrics



Visualizing neurons and generating images



Sensorimotor control

Looking for interns! (Munich, Santa Clara)