



Politecnico
di Torino

SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation

Vijay Badrinarayanan, Alex Kendall, Roberto Cipolla, Senior Member, IEEE,
Conference

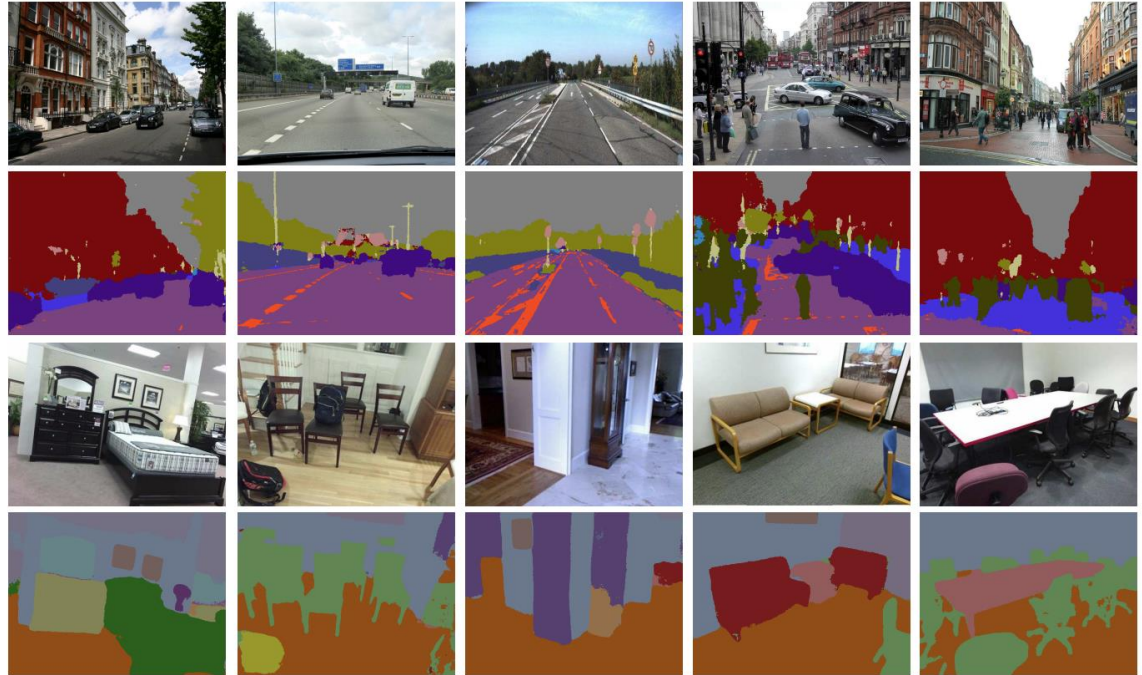
Davide Bartoletti
Domenico Bulfamante
Giovanni Sciortino

MLDL
A.A 2021/2022

Introduction

Semantic Segmentation:

Process of linking each pixel in an image to a class label.

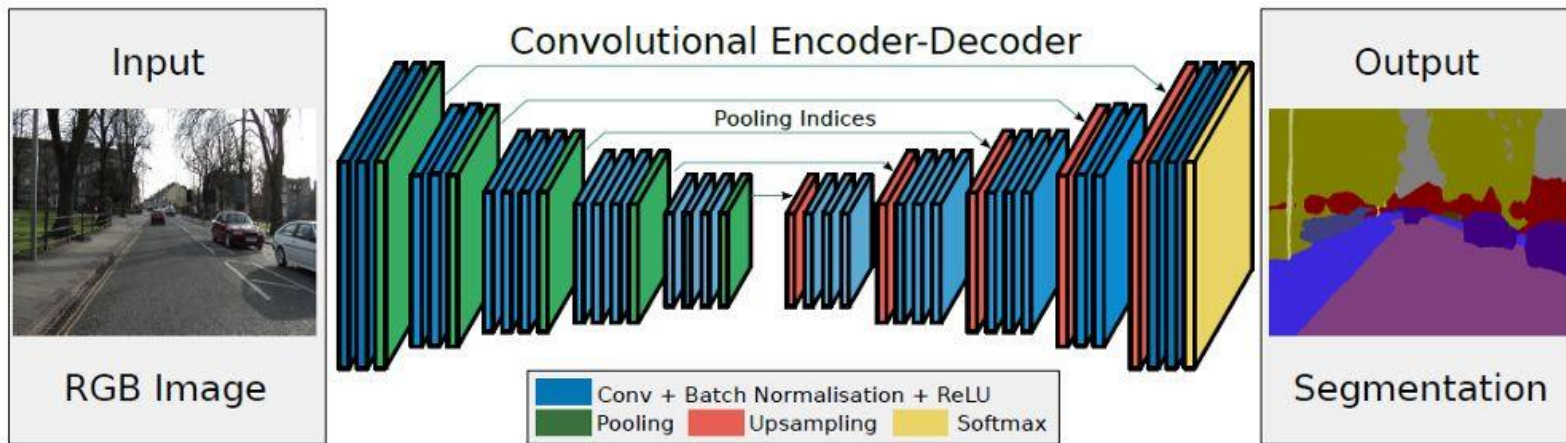


SegNet predictions on road scenes and indoor scenes

SegNet

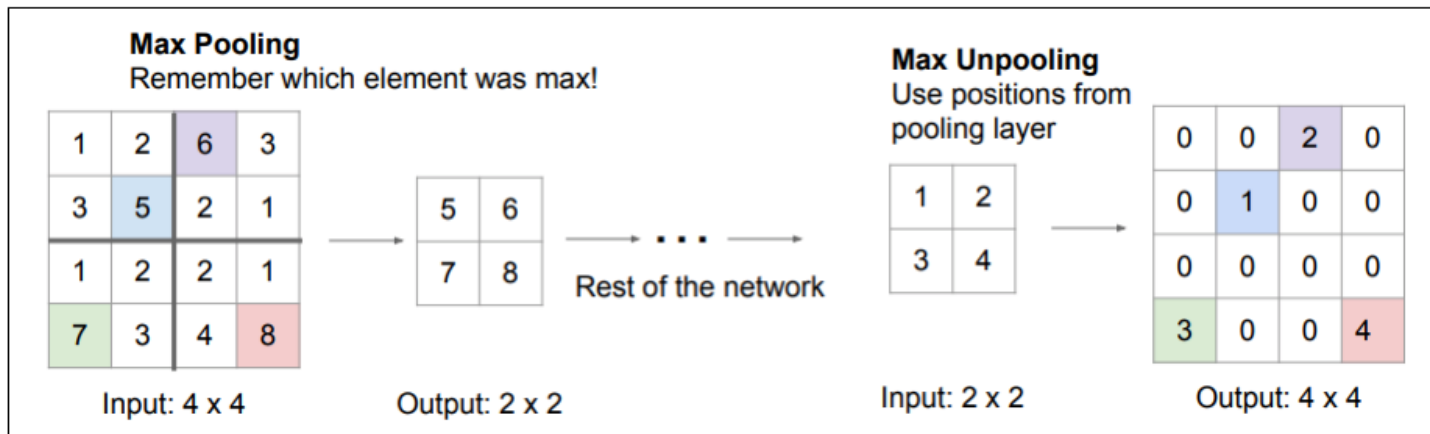
- Encoder Network : VGG16 with no fully connected layers
- Decoder Network : hierarchy of decoders
- Trade off between memory and accuracy

*Max-pooling
Indices*

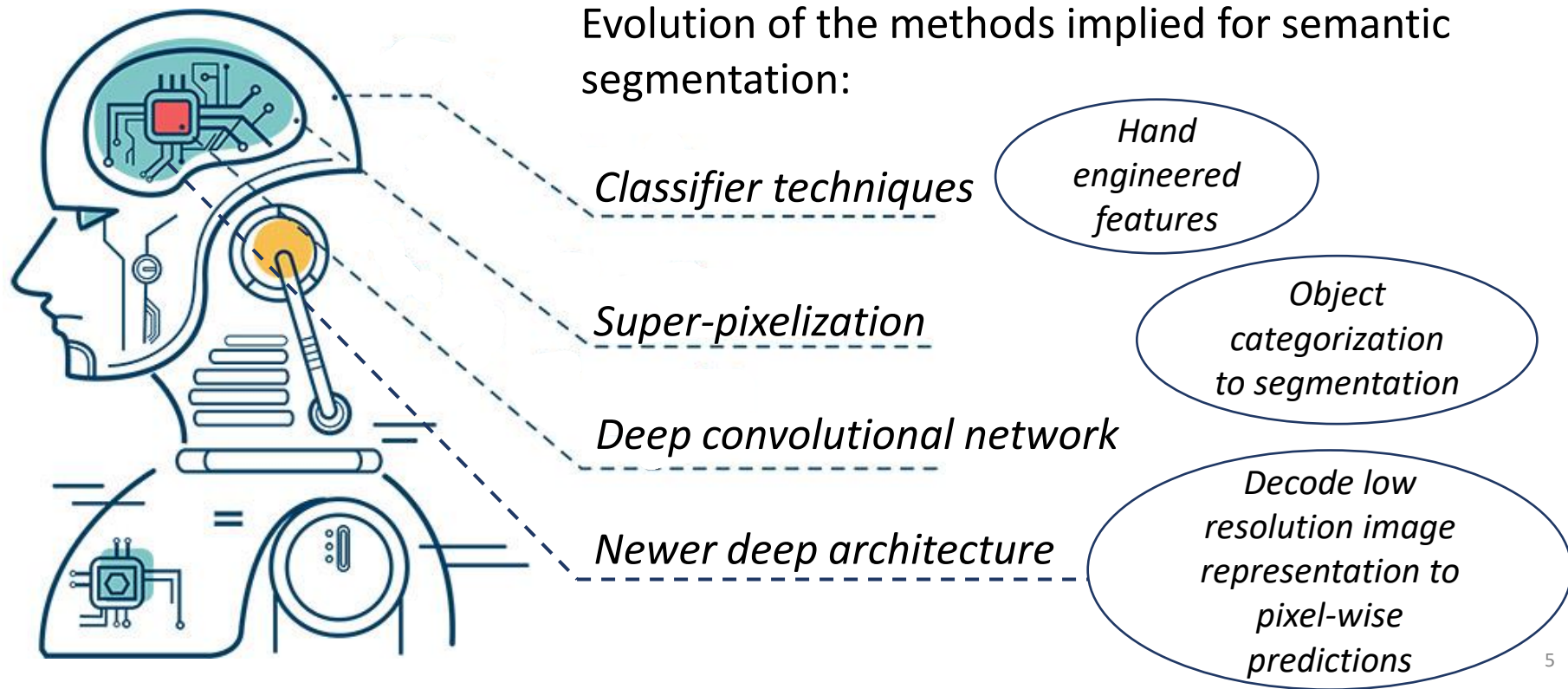


Why the usage of max-indices in the decoder phase?

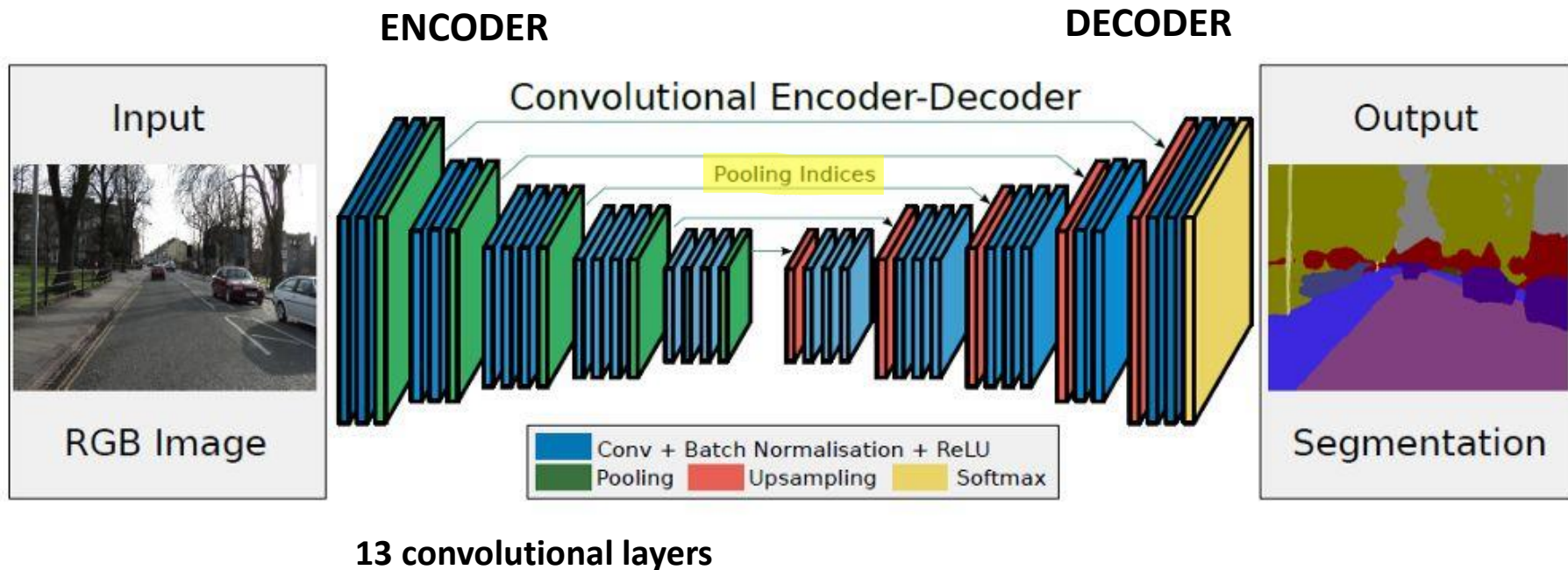
- 1) Improves boundary delineation
- 2) Reduces the number of parameters enabling end-to-end training
- 3) This form of upsampling can be incorporated into any encoder-decoder architecture



Related Work



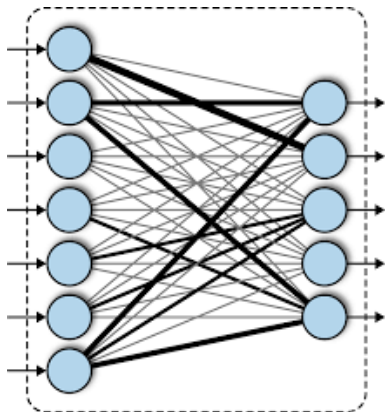
Architecture



Two similar architectures:

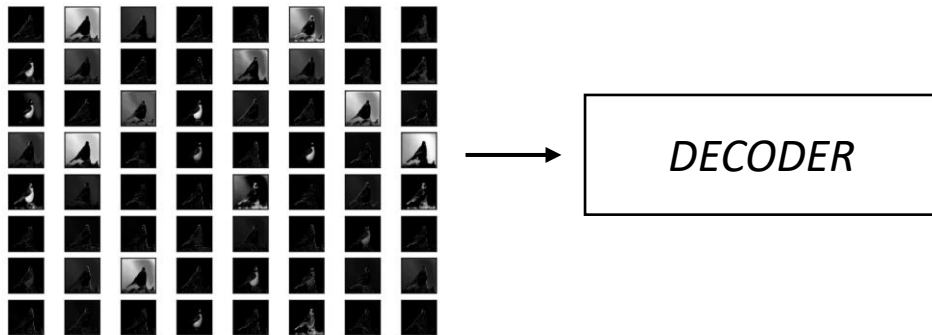
DeconvNet:

- Larger parameterization and computational resources
- Usage of fully connected layers



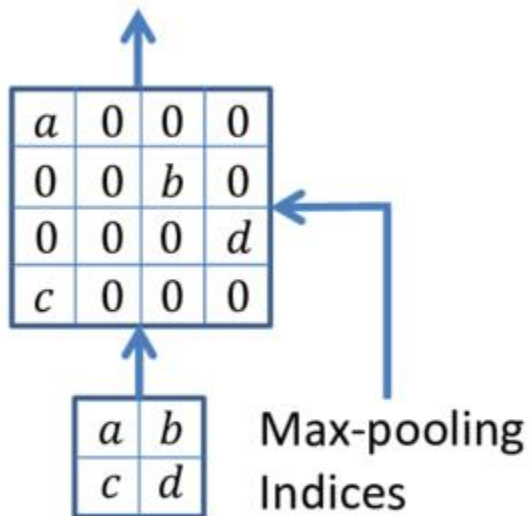
U-Net:

- No reuse of pooling indices
- Transferring of the entire feature map to the decoders



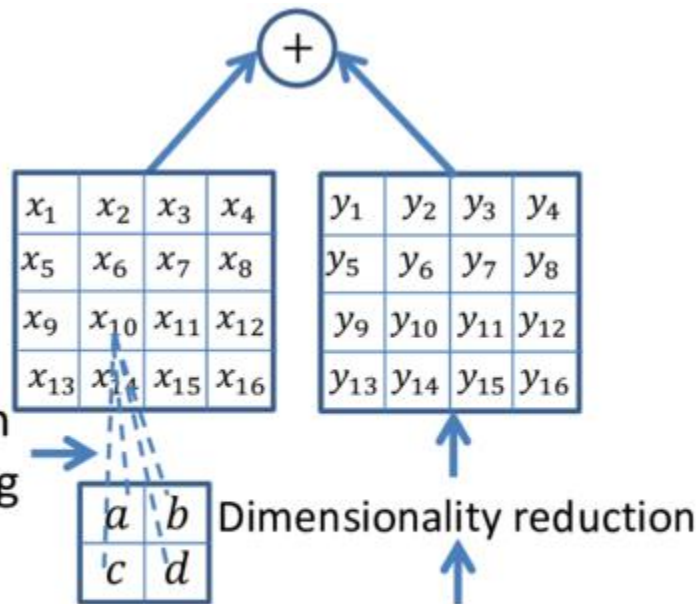
SegNet vs Fully Convolutional Network

Convolution with trainable decoder filters



SegNet

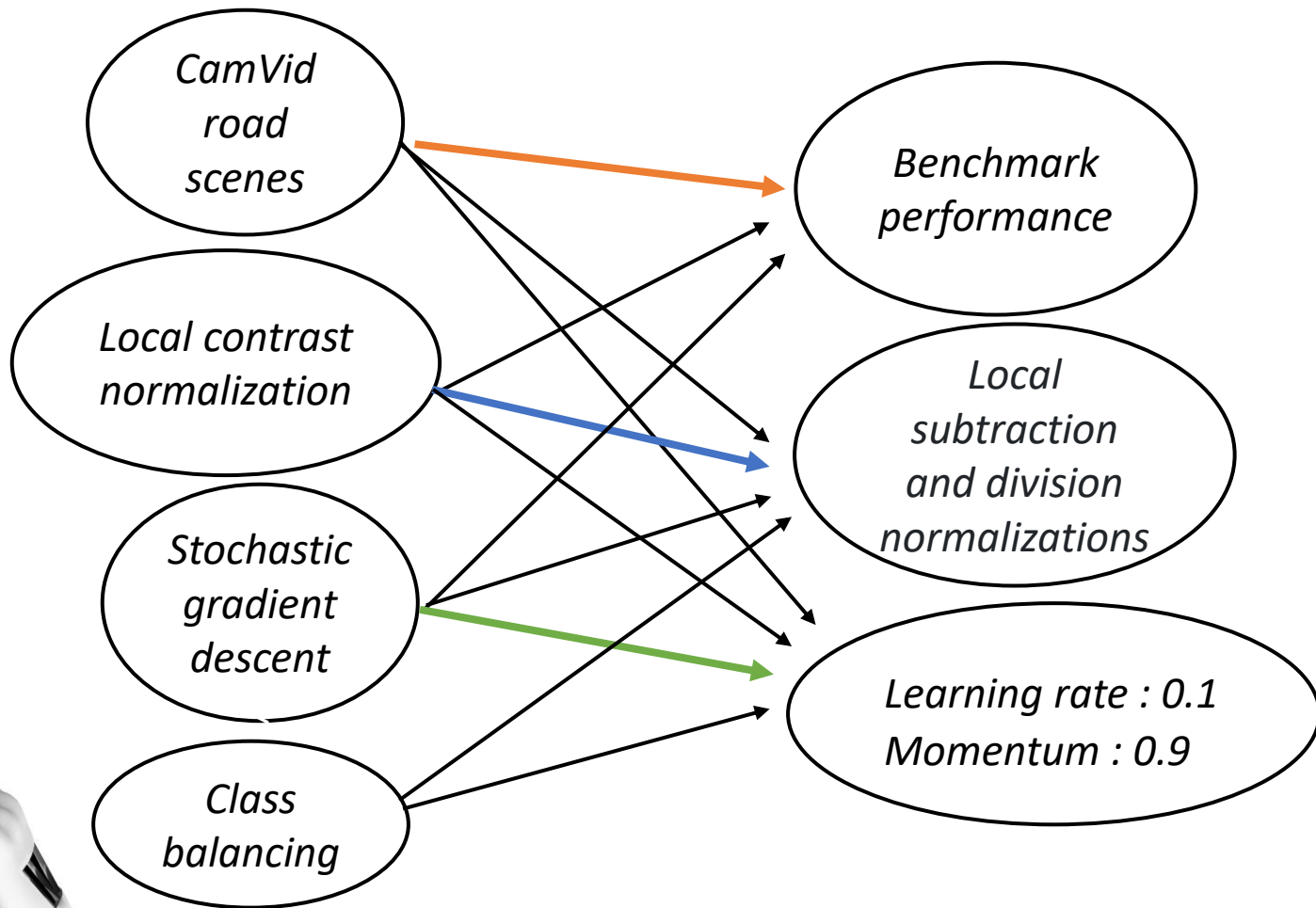
Deconvolution
for upsampling



Encoder feature map

FCN

Training



Analysis

- Global accuracy
- Class average accuracy
- Mean intersection over union – mIoU
- Semantic contour measure - BF

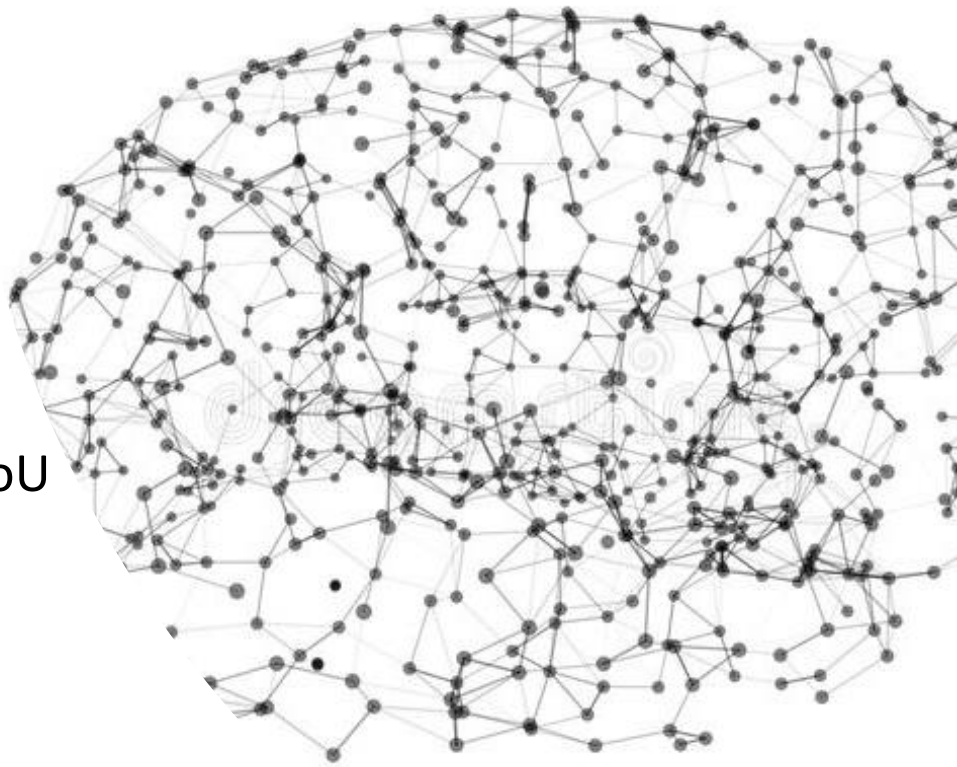


Table of comparison of decoder variants

Variant	Params (M)	Storage multiplier	Infer time (ms)	Median frequency balancing								Natural frequency balancing							
				Test				Train				Test				Train			
				G	C	mIoU	BF	G	C	mIoU		G	C	mIoU	BF	G	C	mIoU	
Fixed upsampling																			
Bilinear-Interpolation	0.625	0	24.2	77.9	61.1	43.3	20.83	89.1	90.2	82.7		82.7	52.5	43.8	23.08	93.5	74.1	59.9	
Upsampling using max-pooling indices																			
SegNet-Basic	1.425	1	52.6	82.7	62.0	47.7	35.78	94.7	96.2	92.7		84.0	54.6	46.3	36.67	96.1	83.9	73.3	
SegNet-Basic-EncoderAddition	1.425	64	53.0	83.4	63.6	48.5	35.92	94.3	95.8	92.0		84.2	56.5	47.7	36.27	95.3	80.9	68.9	
SegNet-Basic-SingleChannelDecoder	0.625	1	33.1	81.2	60.7	46.1	31.62	93.2	94.8	90.3		83.5	53.9	45.2	32.45	92.6	68.4	52.8	
Learning to upsample (bilinear initialisation)																			
FCN-Basic	0.65	11	24.2	81.7	62.4	47.3	38.11	92.8	93.6	88.1		83.9	55.6	45.0	37.33	92.0	66.8	50.7	
FCN-Basic-NoAddition	0.65	n/a	23.8	80.5	58.6	44.1	31.96	92.5	93.0	87.2		82.3	53.9	44.2	29.43	93.1	72.8	57.6	
FCN-Basic-NoDimReduction	1.625	64	44.8	84.1	63.4	50.1	37.37	95.1	96.5	93.2		83.5	57.3	47.0	37.13	97.2	91.7	84.8	
FCN-Basic-NoAddition-NoDimReduction	1.625	0	43.9	80.5	61.6	45.9	30.47	92.5	94.6	89.9		83.7	54.8	45.5	33.17	95.0	80.2	67.8	

G : global average

C : class average

mIoU : mean of intersection over union

BF : semantic contour measure

Analysis Summary



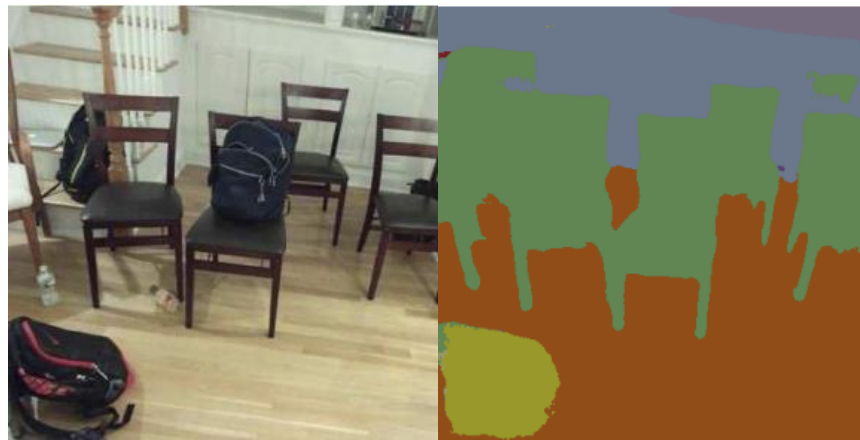
- The best performance is achieved when encoder feature maps are stored in full. This is reflected in the semantic contour delineation metric (BF) most clearly.
- When memory during inference is constrained, then compressed forms of encoder feature maps can be stored and used with an appropriate decoder to improve performance.
- Larger decoders increase performance for a given encoder network.

Benchmarking

We quantify the performance of SegNet on two scene segmentation benchmarks:



Road scene segmentation



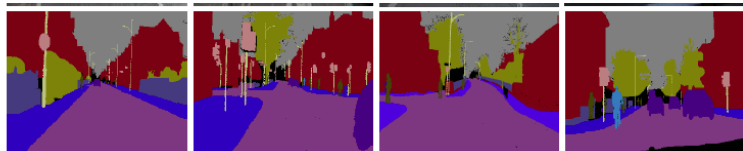
Indoor scene segmentation

Road Scene Segmentation

Test samples

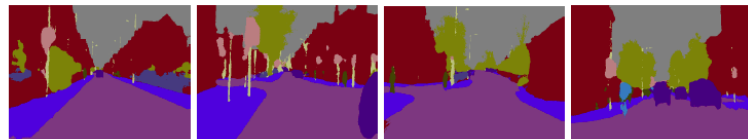


Ground Truth

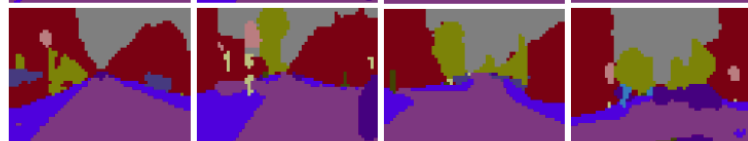


Results on CamVid day and dusk test samples

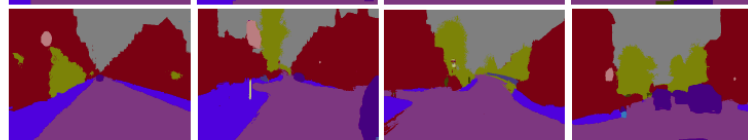
SegNet



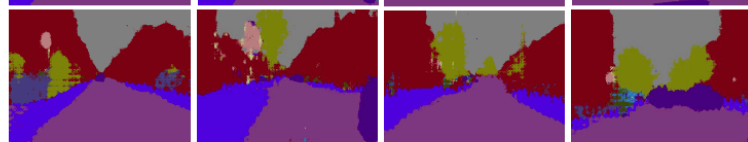
DeepLab-LargeFOV



DeepLab-LargeFOV-denseCRF



FCN



FCN (learn deconv)



DeconvNet



Method	Building	Tree	Sky	Car	Sign-Symbol	Road	Pedestrian	Fence	Column-Pole	Side-walk	Bicyclist	Class avg.	Global avg.	mIoU	BF
SfM+Appearance [28]	46.2	61.9	89.7	68.6	42.9	89.5	53.6	46.6	0.7	60.5	22.5	53.0	69.1	n/a*	
Boosting [29]	61.9	67.3	91.1	71.1	58.5	92.9	49.5	37.6	25.8	77.8	24.7	59.8	76.4	n/a*	
Dense Depth Maps [32]	85.3	57.3	95.4	69.2	46.5	98.5	23.8	44.3	22.0	38.1	28.7	55.4	82.1	n/a*	
Structured Random Forests [31]	n/a											51.4	72.5	n/a*	
Neural Decision Forests [64]	n/a											56.1	82.1	n/a*	
Local Label Descriptors [65]	80.7	61.5	88.8	16.4	n/a	98.0	1.09	0.05	4.13	12.4	0.07	36.3	73.6	n/a*	
Super Parsing [33]	87.0	67.1	96.9	62.7	30.1	95.9	14.7	17.9	1.7	70.0	19.4	51.2	83.3	n/a*	
SegNet (3.5K dataset training - 140K)	89.6	83.4	96.1	87.7	52.7	96.4	62.2	53.45	32.1	93.3	36.5	71.20	90.40	60.10	46.84
CRF based approaches															
Boosting + pairwise CRF [29]	70.7	70.8	94.7	74.4	55.9	94.1	45.7	37.2	13.0	79.3	23.1	59.9	79.8	n/a*	
Boosting+Higher order [29]	84.5	72.6	97.5	72.7	34.1	95.3	34.2	45.7	8.1	77.6	28.5	59.2	83.8	n/a*	
Boosting+Detectors+CRF [30]	81.5	76.6	96.2	78.7	40.2	93.9	43.0	47.6	14.3	81.5	33.9	62.5	83.8	n/a*	

Quantitative comparisons of SegNet with traditional methods on the CamVid 11 road class segmentation problem.

SUN RGB-D Indoor Scenes

Test samples



Ground Truth



Qualitative assessment of SegNet predictions on RGB indoor test scenes from the recently released SUN RGB-D dataset

SegNet



DeepLab-LargeFOV



DeepLab-LargeFOV-denseCRF



*FCN
(learn deconv)*



DeconvNet



Wall	Floor	Cabinet	Bed	Chair	Sofa	Table	Door	Window	Bookshelf	Picture	Counter	Blinds
83.42	93.43	63.37	73.18	75.92	59.57	64.18	52.50	57.51	42.05	56.17	37.66	40.29
Desk	Shelves	Curtain	Dresser	Pillow	Mirror	Floor mat	Clothes	Ceiling	Books	Fridge	TV	Paper
11.92	11.45	66.56	52.73	43.80	26.30	0.00	34.31	74.11	53.77	29.85	33.76	22.73
Towel	Shower curtain	Box	Whiteboard	Person	Night stand	Toilet	Sink	Lamp	Bathtub	Bag		
19.83	0.03	23.14	60.25	27.27	29.88	76.00	58.10	35.27	48.86	16.76		

Class average accuracies of SegNet predictions for the 37 indoor scene classes in the SUN RGB-D benchmark dataset

CamVid segmentation vs SUNRGB-D segmentation

Network/Iterations	40K				80K				>80K				Max iter
	G	C	mIoU	BF	G	C	mIoU	BF	G	C	mIoU	BF	
SegNet	88.81	59.93	50.02	35.78	89.68	69.82	57.18	42.08	90.40	71.20	60.10	46.84	140K
DeepLab-LargeFOV [3]	85.95	60.41	50.18	26.25	87.76	62.57	53.34	32.04	88.20	62.53	53.88	32.77	140K
DeepLab-LargeFOV-denseCRF [3]	not computed								89.71	60.67	54.74	40.79	140K
FCN	81.97	54.38	46.59	22.86	82.71	56.22	47.95	24.76	83.27	59.56	49.83	27.99	200K
FCN (learnt deconv) [2]	83.21	56.05	48.68	27.40	83.71	59.64	50.80	31.01	83.14	64.21	51.96	33.18	160K
DeconvNet [4]	85.26	46.40	39.69	27.36	85.19	54.08	43.74	29.33	89.58	70.24	59.77	52.23	260K

Quantitative comparison of deep networks for semantic segmentation on the CamVid test set when trained on a corpus of 3433 road scenes without class balancing.

Network/Iterations	80K				140K				>140K				Max iter
	G	C	mIoU	BF	G	C	mIoU	BF	G	C	mIoU	BF	
SegNet	70.73	30.82	22.52	9.16	71.66	37.60	27.46	11.33	72.63	44.76	31.84	12.66	240K
DeepLab-LargeFOV [3]	70.70	41.75	30.67	7.28	71.16	42.71	31.29	7.57	71.90	42.21	32.08	8.26	240K
DeepLab-LargeFOV-denseCRF [3]	not computed								66.96	33.06	24.13	9.41	240K
FCN (learnt deconv) [2]	67.31	34.32	24.05	7.88	68.04	37.2	26.33	9.0	68.18	38.41	27.39	9.68	200K
DeconvNet [4]	59.62	12.93	8.35	6.50	63.28	22.53	15.14	7.86	66.13	32.28	22.57	10.47	380K

Quantitative comparison of deep architectures on the SUNRGB-D dataset when trained on a corpus of 5250 indoor scenes

Conclusion

SegNet :

- Born from the need to design an efficient architecture for road and indoor scene understanding which is efficient both in terms of memory and computational time
- Stores the max-pooling indices of the feature maps and uses them in its decoder network to achieve good performance.

