Supervised convolutional models for natural image understanding in road traffic

Siniša Šegvić, Ivan Krešo, Josip Krapac Faculty of Electrical Engineering and Computing University of Zagreb

INTRODUCTION: ABOUT

Semantic segmentation: image understanding at the pixel-level

- pixel-level: associate each pixel with a class
- image understanding: classes have a high-level meaning traffic participants: person (red), car (blue), bicycle (dark red)
 objects: pole (light grey), traffic sign (yellow), traffic light (orange)
 landscape: road (purple), sidewalk (pink), building (dark grey),
 vegetation (dark green), terrain (light green), sky (light blue)



INTRODUCTION: AGENDA

- 1. About semantic segmentation
 - overview, fully convolutional approach, problems
- 2. Achieving scale invariance by depth-driven selection
 - improve recognition with reconstruction
- 3. Restoring the resolution with ladder-style upsampling
 - □ blend semantic information with spatial accuracy
 - recover semantic information with the DenseNet architecture
- 4. Experiments
 - □ performance criteria, datasets, results
- 5. Conclusion

OVERVIEW: RECOGNITION

Example: discriminate bison from oxen



- 1. express the program with many free parameters
 - the parameters determine a transformation which we call the model
- 2. fit parameters on the training set
- 3. evaluate performance on the test set

Success depends on the model, training set and processing power.

OVERVIEW: ARCHITECTURE

Deep convolutional model for image classification [krizhevsky12nips]

- input: image; output: distribution over 1000 classes
- □ fitness criterion: average log probability of the correct class
- structure: a succession of convolutions and poolings
 - gradual decrease of resolution and increase of the semantic depth
- recent architectures: O(10²) layers, O(10⁶) parameters, O(10⁹) multiplications for a 224x224 image!



OVERVIEW: SLIDING WINDOW

- A classification model can be applied to the segmentation task:
 - analyze the image in the **sliding window** fashion
 - $\hfill\square$ each patch produces one pixel of the semantic map
 - segmentation groundtruth allows end-to end training
 - □ each pixel becomes one component of the fitness criterion
 - optimized implementation required in practice
 - \square 10⁶ pixels \times 10⁹ multiplications?



OVERVIEW: GOING FULLY CONVOLUTIONAL

Luckily, the processing of neighbouring patches involves calculating many common latent activations

More efficient: perform the classification layer-wise [long15cvpr]:

- the resulting semantic map is subsampled due to pooling
- □ this can be relaxed to some extent with dilated filtering [yi16iclr]



[long15cvpr]

PROBLEMS: LARGE OBJECTS

Classifying pixels at large objects may require a huge receptive field.

- many local neighbourhoods are not discriminative enough
- □ their center pixel can only be recognized in a larger context
- problem arises when the context is larger than receptive field



[kreso16gcpr]

PROBLEMS: SMALL OBJECTS

Detecting small objects with a huge receptive field wastes resources:

- small and simple objects can be recognized with few layers
- □ latter layers forward their activations over and over again
- that leads to loss of the representational power



[kreso16gcpr]

PROBLEMS: MEMORY REQUIREMENTS

Successful segmentation architectures rely on pretrained models designed for ImageNet: small input resolution, single pixel output

- □ in segmentation we have large resolution both on input and output
- $\hfill\square$ brute force output restoration (dilated filtering) is feasible up to 8x \uparrow



Mainly relevant for training, where all activations must be cached

□ however, lean memory requirements favour short evaluation times

SCALE INVARIANCE: IDEA

Use stereo reconstruction to disentangle appearance from scale

- □ independently extract features from all levels of image pyramid
- analyze each pixel at the pyramid level determined by its distance from the camera
- effect: image objects are perceived at the common scale regardless of the camera distance



[kreso16gcpr]

SCALE INVARIANCE: ARCHITECTURE

We introduce a new layer: the scale selection multiplexer

- □ the multiplexer assembles pieces from appropriate pyramid levels
- the back-end receives a scale-invariant feature mosaic
- the scale multiplexer is compatible with end-to-end training



SCALE INVARIANCE: RESULTS

Scale selection solves the problems of large and small objects:

- nearby objects are recognized in diminished images
- □ far objects are recognized at the original resolution



Effects of scale selection (mIoU): 56.4 \rightarrow 64.4

Predictions are still subsampled, the memory problem remains

SMART UPSAMPLING: IDEA

Using a full-fledged ImageNet-class model in each pixel is wasteful:

- boundary refinement should be easier than recognition
- upsample a deep representation by blending it with a higher resolution earlier layer [valpola14arxiv,ronneberger15arxiv]

The resulting architecture operates as follows:

- the downsampling datapath infers the semantic information
- the upsampling datapath refines the boundaries
- □ the lateral connections ensure the blending: $\hat{F}_t = g_t(F_t, \hat{F}_{t+1})$



SMART UPSAMPLING: DENSENET

The downsampling datapath can be any classification architecture

we compare ResNet [he16cvpr] and DenseNet [huang17cvpr]



both architectures favour the gradient flow towards the early layers

hypothesis: DenseNet is better when the class complexities vary

 $CCVW2017 \rightarrow Smart upsampling 15/28$

SMART UPSAMPLING: LADDER-STYLE

Hypothesis: recognizing objects is harder than boundary refinement

use a lean representation in the upsampling datapath

this results in huge memory savings, unlike [jegou16arxiv]



We split DB4 into DB4a and DB4b to increase receptive field

- □ this results in 64× downsampling (appropriate for large objects)
- □ this compromises ImageNet initialization for DB4b
- □ however, fine-tuning succeds to recover

SMART UPSAMPLING: RESULTS

We train on full Cityscapes resolution with bs=2 on two GTX1070

We recover fine details lost due to $64 \times$ downsampling

- □ middle: upsampling 64× by interpolation
- □ bottom: upsampling $16 \times$ by blending and $4 \times$ by interpolation



Effect of blending (mIoU): 62.5 \rightarrow 72.8

EVALUATION: DATASETS

Pascal VOC 2012x [everingham10ijcv]:

- generic photographs
- 10 indoor, 10 outdoor classes
- □ 12000 images, <.25 MPixel
- KITTI [geiger13ijrr]:
 - driver's perspective, 11 classes
 - 450 stereo images, 0.5 MPixel
 - reconstruction groundtruth
 - odometry groundtruth
 - Karlsruhe, fine weather





EVALUATION: DATASETS (2)

- □ Cityscapes [cordts16cvpr]:
 - □ driver's perspective, 19 classes
 - □ 5000 stereo images, 2MPixel
 - 20000 coarsely annotated images
 - instance level annotations
 - □ 50 cities, spring to autumn
- Vistas [neuhold17iccv]:
 - driver's perspective, 100 classes
 - 25000 images, 2-8 MPixel
 - instance level annotations
 - worldwide, various weather







EVALUATION: CITYSCAPES [CORDTS16CVPR]



Pros: fine annotation, many classes, well-chosen categories, complex cluttered scenes, variety of scale



 $CCVW2017 \rightarrow Evaluation$ (2) 20/28

EVALUATION: PERFORMANCE CHARACTERIZATION



Typically, the performance is expressed as mean IoU over all classes

 $\square \ \mathsf{mloU} = \frac{\sum_c \mathrm{IoU}_c}{C}$

- □ this increases the influence of rare classes with few training pixels
- □ examples: wall, fence, pole, bottle, potted plant

To ensure integrity, labels of test subsets are withheld from public datasets

The performance on the test set is determined by submitting results to the evaluation server

EVALUATION: TEST ERRORS



bus/truck road/sidewalk





bus/car/building sidewalk road pedestrian/rider

EVALUATION: TEST ERRORS (2)





truck/car motorcycle/bicycle building/wall building/sign



EVALUATION: TEST ERRORS (3)





bus/tram sidewalk/road



bus/tram/car bus/truck/wall/fence





car/truck/tram/ car/building/sign

EVALUATION: WORST PERFORMANCE ON VAL





















Problems: i) ambiguous sidewalk, ii) unusual fence, iii) large truck, iv) fence vs wall, occlusions.

EVALUATION: PERFORMANCE



[zhao17arxiv]

CONCLUSION

Scale selection: Cityscapes mIoU 56.4 \rightarrow 64.4

Ladder-style upsampling (mIoU): Cityscapes 62.5 \rightarrow 72.8

ResNet vs DenseNet: Cityscapes 69.5 \rightarrow 72.8; VOC12 63.0 \rightarrow 70.2

Performance on test (mIoU): Cityscapes 74.6; VOC 2012 AUG 78.0

Able to train at full Cityscapes resolution with bs=2 on $2 \times GTX1070$ Able to process 1024×448 images at 31Hz with a 74.6 mIoU model Able to recover from $64 \times$ subsampling





CONCLUSION: DISCUSSION

Thank you for your attention!

Questions?



This presentation has been fully supported by Croatian Science Foundation under the project I-2433-2014. http://multiclod.zemris.fer.hr

The Titan X used in experiments was donated by NVIDIA Corporation.

APPENDIX: ADVERSARIAL EXAMPLES



