# Convolutional models for image understanding

principles, challenges and research outlook

Siniša Šegvić
UniZg-FER

# AGENDA

- Introduction: image understanding
- Convolutional models: principles performance
- Challenges towards intelligent perception
- Overview of our recent research
- Conclusions

# INTRODUCTION: IMAGE CLASSIFICATION

Story of natural image understanding starts with image classification

Consider the problem of discriminating bison from oxen:



[image-net.org]

Hard because we do not know where is the defining object

Especially hard when intra-class variance is large and inter-class variance is small

# INTRODUCTION: A RULE-BASED APPROACH?

We could try to solve image classification by a rule-based system:

1. concentrate on image regions projected from four legged animals
2. oxen have longer horns or no horns at all
3. bison are dark brown, etc, etc

It's explainable and neat, but nobody succeeded to make that work!



[image-net.org]

# INTRODUCTION: LEARNING

Hence, we look up approaches which are able to learn functionality from the data

A machine learning approach would roughly follow the following steps:

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

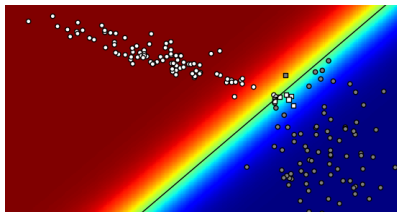Success depends on the model, training set and processing power

1. model may have insuficient (or excessive) capacity

2. the training set may be too small or not representative enough

3. insufficient processing power ⇒ the training may not converge

# INTRODUCTION: SHALLOW MODELS

Transformation: data → decision

1. one linear projection
2. optional squashing non-linearity

An example in 2D: $y_i = \sigma(\mathbf{w}^\top \mathbf{x} + b)$



Advantages of shallow models:

☐ the best solution is guaranteed and fast

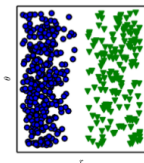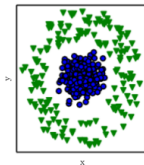☐ this is the best approach when classes are linearly separable

Unfortunately, shallow models are a poor fit for image classification:

☐ the model can only learn a lookup table

☐ applicable only for very simple tasks (eg. 88% on MNIST)

☐ insufficient capacity, tendency to underfit

# INTRODUCTION: DATA REPRESENTATION

Classification can profit from good representation:

- □ it would be easy to discriminate bison from oxen if some magical algorithm converted the input image into a binary vector:
  [fur?, small horn?, wilderness?, ...]

- □ most bison would be: [1, 1, 1, ...]

- □ most oxen would be: [0, 0, 0, ...]



[goodfellow16]

Hand crafting quality representation is hard:
- □ Greek and Romans did not invent 0 in 1000 years of civilization
- □ that hindered development of science and engineering
  - □ MCMLXXI + XXIX = ?
  - □ MXXIV : LXIV = ?

A lot of hard work of smart people went into feature engineering

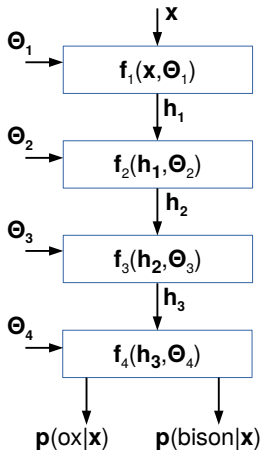# INTRODUCTION: COMPOSITIONAL PARADIGM

Deep model: a sequence of learned non-linear transformations

Why deep learning was not successful?

- □ no guarantee of the learning success
- □ non-competitive performance

Why deep learning became popular?

- □ better modeling and training
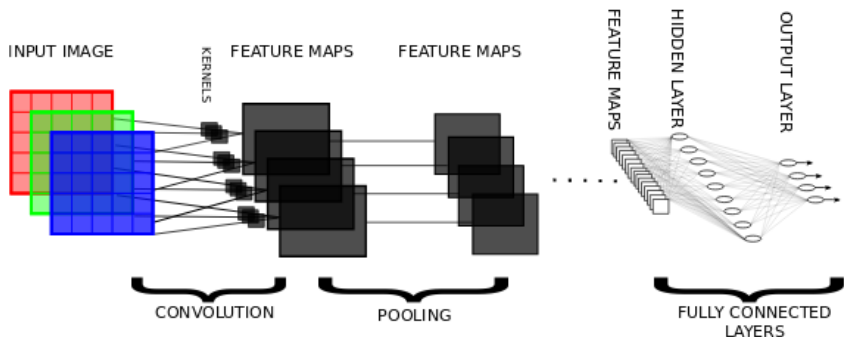- □ large datasets ($n=10^6$)
- □ processing power (TFLOPS)

Suitable for **articulated** data: images, language, speech, bioinformatics...



$\Theta_1 \rightarrow$ $\mathbf{f}_1(\mathbf{x}, \Theta_1)$

$\mathbf{x}$

$\mathbf{h}_1$

$\Theta_2 \rightarrow$ $\mathbf{f}_2(\mathbf{h}_1, \Theta_2)$

$\mathbf{h}_2$

$\Theta_3 \rightarrow$ $\mathbf{f}_3(\mathbf{h}_2, \Theta_3)$

$\mathbf{h}_3$

$\Theta_4 \rightarrow$ $\mathbf{f}_4(\mathbf{h}_3, \Theta_4)$

$\mathbf{p}(\text{ox}|\mathbf{x})$    $\mathbf{p}(\text{bison}|\mathbf{x})$
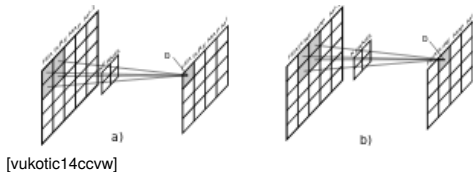
# CONVOLUTIONAL MODELS: ARCHITECTURE

Deep models for image classification typically consist of:

- ☐ **convolutions** (linear, recognize object parts)
- ☐ poolings (reduce representation dimensionality)
- ☐ projections (linear, recognize image as a whole)
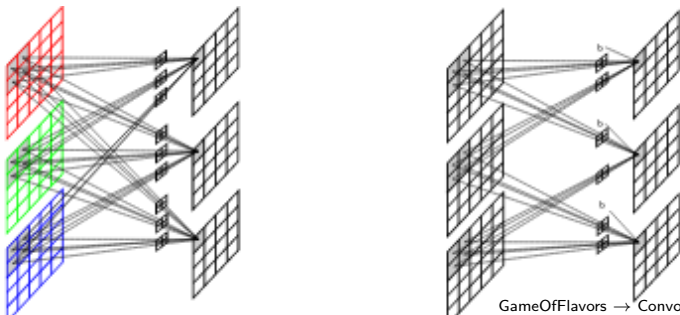- ☐ elementwise non-linearities, eg. max(0,x)

# CONVOLUTIONAL MODELS: CONVOLUTIONS

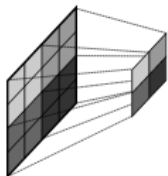Task: convolve the previous feature map with a kernel



[vukotic14ccvw]

We typically have multiple feature maps on output ⇒ multiple kernels

We typically have several feature maps on input ⇒ kernels are 3D
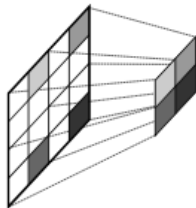
# Convolutional models: poolings

Task: reduce dimensionality to relax memory requirements



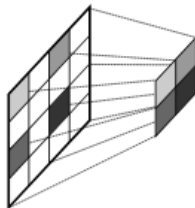[vukotic14ccvw]

Most often implemented as average pooling or max pooling
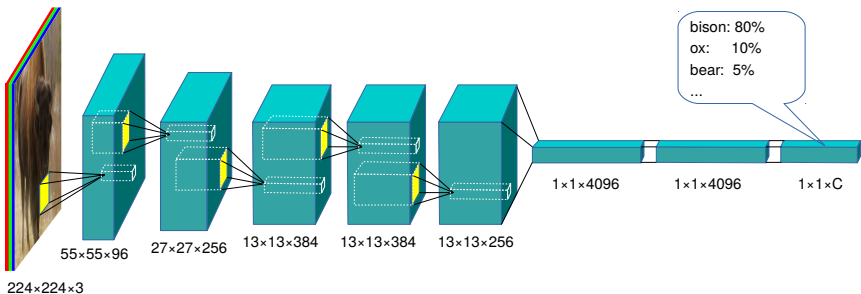
It also increases translation invariance:



[vukotic14ccvw]

# CONVOLUTIONAL MODELS: LARGE-SCALE

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: $O(10^2)$ layers, $O(10^6)$ parameters; $O(10^8)$ bytes and $O(10^9)$ operations for a $224 \times 224$ image!



bison: 80%
ox:    10%
bear:  5%
...

55×55×96   27×27×256   13×13×384   13×13×384   13×13×256   1×1×4096   1×1×4096   1×1×C

224×224×3

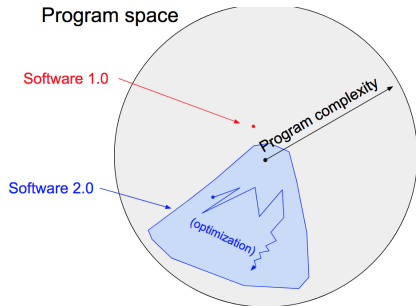# Convolutional models: Two views



[xkcd 1838]

Neural networks →
~~Deep learning~~ →
Differential programming
(software 2.0)
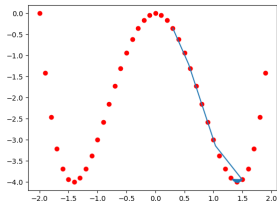


[karpathy17medium]

# Convolutional models: Diff. programming

```
import torch

step=0.12

x=torch.rand(1, requires_grad=True)
#x=torch.tensor(1.9, requires_grad=True)

for i in range(100):
  y = x**4 - 4*x**2
  y.backward()
  print(x, y, x.grad)
  x.data = x - step*x.grad
  x.grad.zero_()
```



The workhorse algorithm: reverse-mode automatic differentation

# Convolutional models: ImageNet

One of the most popular vision datasets [russakovsky15ijcv]

- □ annual challenges: classification, localization, detection in video
- □ we focus on the classification challenge: $10^6$ images, $10^3$ classes
- □ fine-grained animals, objects, materials, sports, dishes...
- □ evaluation metric: top-five prediction error (trained human: 5%)

red fox (100)  hen-of-the-woods (100)  ibex (100)  goldfinch (100) flat-coated retriever (100)
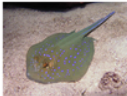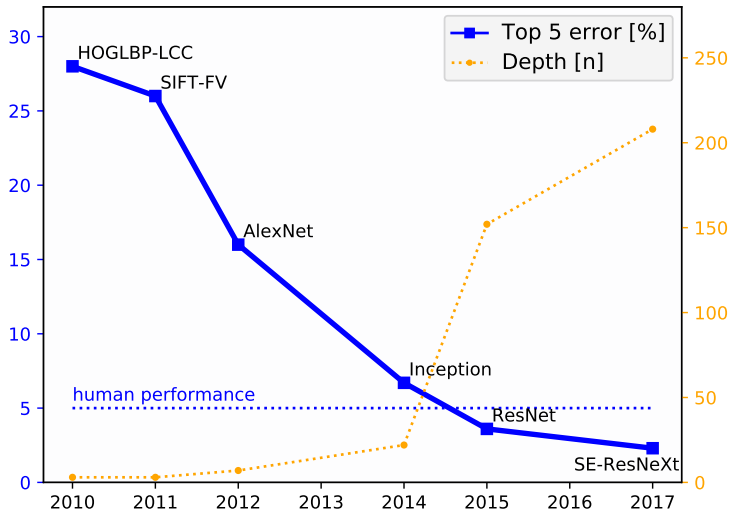
tiger (100)  hamster (100)  porcupine (100)  stingray (100)  Blenheim spaniel (100)

# CONVOLUTIONAL MODELS: IMAGENET HARD EXAMPLES

Tasks which are difficult for **humans** [russakovsky15ijcv]:

- fine-grained classification (e.g. 120 breeds of dogs!)
- exotic classes (pulley, spotlight, maypole)

Tasks which are difficult for **GoogleNet** (2014, 6.7%):

- little and thin objects, filtered and atypical images
- abstraction (a toy hatchet, images with text)
- large intra-class variance, small between-class variance



muzzle (71)  hatchet (68)  water bottle (68)  velvet (68)  loupe (66)

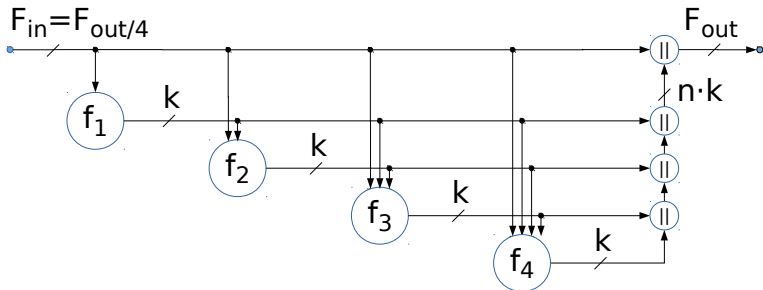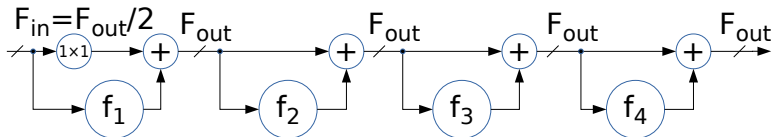hook (66)  spotlight (66)  ladle (65)  restaurant (64)  letter opener (59)

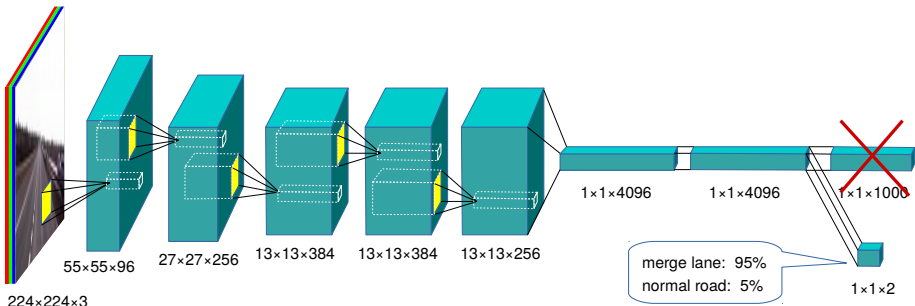# CONVOLUTIONAL MODELS: FURTHER IMPROVEMENT

Important idea: enhance the gradient flow towards early layers
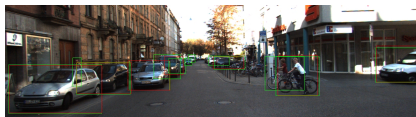
# CONVOLUTIONAL MODELS: KNOWLEDGE TRANSFER

A deep classification model can be fine-tuned for another (easier) task:

- □ cut-off the last few layers

- □ connect the remaining layers with the back end for the new task

- □ train the resulting model on new images

- □ inherited layers are well-trained so we can train with less data (few thousands images)



224×224×3

55×55×96   27×27×256   13×13×384   13×13×384   13×13×256

1×1×4096      1×1×4096      1×1×1000

merge lane: 95%
normal road: 5%

1×1×2

# CONVOLUTIONAL MODELS: TASKS

- Object localization:
  - detect objects...
  - ...and indicate their location



- Semantic segmentation:
  - label all image pixels...
  - ...with semantic classes



- Stereoscopic reconstruction:
  - label all image pixels...
  - ... with metric distances



- Recognition in video, object tracking, styling and generating images, ...

# HARDWARE: CPU vs GPU

CPU core (AVX) can dispatch up to 2 FMA instructions / cycle

- ☐ 2·2·8 operations @3GHz → 100 GFLOPS
- ☐ if a CPU has 10 cores - that's 1 TFLOPS

Modern GPUs achieve 10+ TFLOPS in matrix multiplication

- ☐ in practice, the advantage is ×50 due to better memory bandwidth.

Price of training a simple ImageNet model:

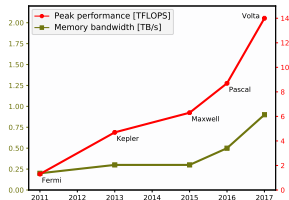| GPU | TFLOPS | price | time |
|---|---|---|---|
| GTX 1070 | 6.5 | 4000 kn | 52 h |
| GTX 1080 Ti | 11.3 | 7000 kn | 31 h |
| P100 FP16 | 25.0 | 47000 kn | 13 h |
| DGX-1 FP16 | 170.0 | $129000 | 2 h |

# HARDWARE: GPU

Best performance/price is obtained on gaming GPUs:

- □ NVIDIA Titan X: 250W, 12GB, 11 TFLOPS, 3000 GPU cores
- □ Radeon Vega FE: 300W, 16GB, 13 TFLOPS, 4000 GPU cores
- □ NVIDIA Titan V: 250W, 16GB, 110 TFLOPS (?), 12nm, $8cm^2$, 21Gt



Actual non-representative measurement:

- □ GTX1070 (3500 kn) $60\times$ faster than E3-1220 v3 (1700kn)

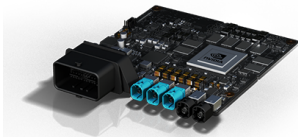Additional challenge: deliver such performance with low power

# HARDWARE: EMBEDDED

Novel hardware concept: processing units for artificial intelligence

- ☐ fast matrix multiplication (10 TFLOPS)
- ☐ low power, low precision (8-32 bit)

Main players:

- ☐ NVIDIA TX2: 1.5 TOPS, 15W, 8GB RAM
- ☐ NVIDIA Xavier: 20 TOPS, 20W, automotive certificate (ISO 26262)
- ☐ Google TPU2: 45 TFLOPS
- ☐ Microsoft + Intel DPU (Stratix-10): 40 TFLOPS

# CHALLENGES: FALSE POSITIVES DUE TO CONTEXT

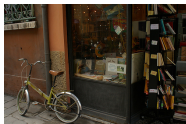Current models tend to produce false positive detections due to context

- good performance likely due to recognition of **easy context**

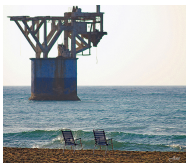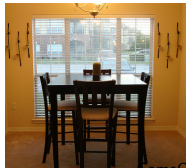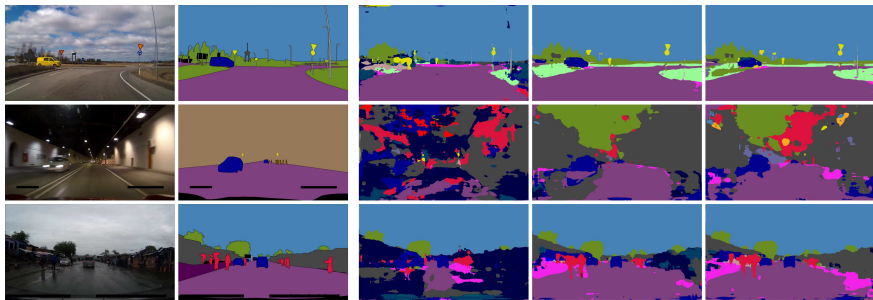Results on Pascal VOC 2012 (confident TP, high FP, low FN):

boat

bottle

bus

car

# CHALLENGES: CROSS-DATASET GENERALIZATION

Often our models generalize well only within dataset

For instance, images from the novel WildDash dataset (left) fool most models trained on popular datasets such as Cityscapes or Vistas (right)
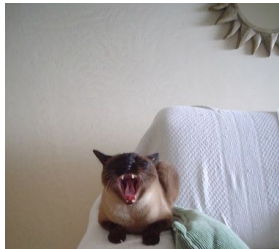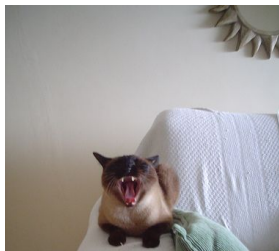


[zendel18eccv]

Conclusion: models tend to <u>overfit</u> to **dataset specifics**

   □ camera, weather, environment, climate, ...

# CHALLENGES: ADVERSARIAL EXAMPLES

Imperceptible perturbations may invalidate prediction [szegedy14iclr]:



[kreso18ep]

# CHALLENGES: ADVERSARIAL EXAMPLES (2)

Adversarial perturbation:

$$\delta = \arg \min_{\delta} p(Y = y_i | \mathbf{x}_i + \delta, \Theta)$$

Adversarial example: $\mathbf{x}_i + \delta$

Existence of adversarial examples suggests that current vision systems are free-riding on **easy features** while ignoring the gist of the scene



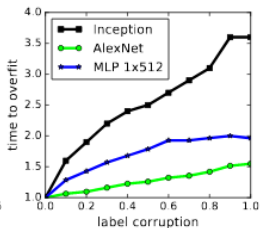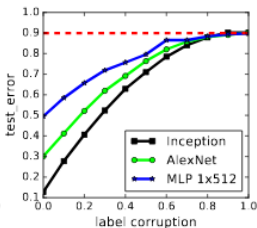| $\boldsymbol{x}$ | $\mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$ | $\boldsymbol{x} + \epsilon \mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$ |
| :---: | :---: | :---: |
| "panda" | "nematode" | "gibbon" |
| 57.7% confidence | 8.2% confidence | 99.3 % confidence |

# CHALLENGES: THEORY

Effective capacity of deep models is large enough to shatter popular image classification datasets:



(a) learning curves    (b) convergence slowdown    (c) generalization error growth
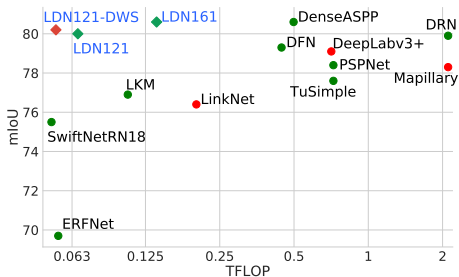
[zhang17iclr]

In simple words, the model is able to memorize the entire training data

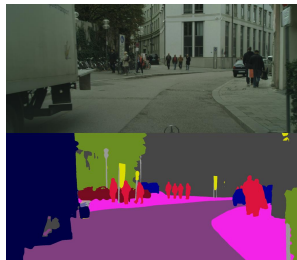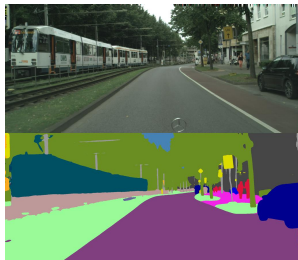Yet, the models generalize well when trained on well sorted data

A theory to explain this behaviour is missing.

# Results: Efficient semantic segmentation

- Goal: pixel-level classification
- DenseNet-121 backbone
- Lean upsampling path
- Memory-efficient training
- SotA results with very efficient model



[kreso19review]

# RESULTS: ROBUST VISION CHALLENGE 2018

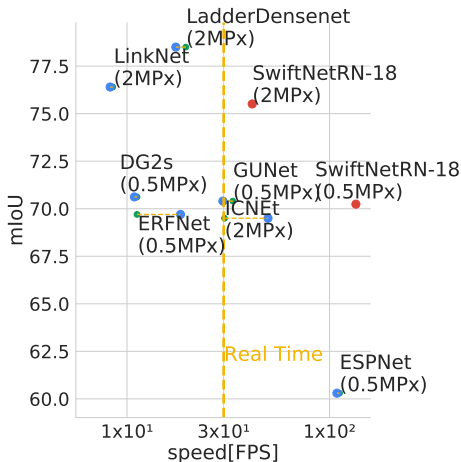☐ 2nd place, 4 datasets, mutiple domains: driving and indoor

| 🏆 | Method | KITTI (Detailed subrankings) | ScanNet (Detailed subrankings) | Cityscapes (Detailed subrankings) | WildDash (Detailed subrankings) |
|---|---|---|---|---|---|
| 1 | MapillaryAI_ROB | 1 | 1 | 1 | 1 |
| | In-Place Activated BatchNorm for Memory-Optimized Training of DNNs [Project page] - Submitted by Peter Kontschieder (Mapillary Research) | | | | |
| 2 | LDN2_ROB | 3 | 2 | 2 | 3 |
| | Ladder-style DenseNets for Semantic Segmentation of Large Natural Images [Project page] - Submitted by Ivan Krešo (University of Zagreb, Faculty of Electrical Engineering and Computing) | | | | |
| 3 | IBN-PSP-SA_ROB | 2 | 3 | 3 | 4 |
| | | | | | Submitted by Anonymous |
| 4 | AHiSS_ROB | 5 | 8 | 5 | 2 |
| | Training of Convolutional Networks on Multiple Heterogeneous Datasets for Street Scene Semantic Segmentation [Project page] - Submitted by Panagiotis Meletis (Eindhoven University of Technology) | | | | |
| 5 | VENUS_ROB | 4 | 4 | 4 | 8 |
| | | | | VENUS-Net for RobustVision - Submitted by Anonymous | |
| 6 | AdapNetv2_ROB | 5 | 5 | 6 | 7 |
| | | | | | Submitted by Anonymous |
| 7 | VlocNet++_ROB | 7 | 5 | 9 | 5 |
| | | | | | Submitted by Anonymous |
| 8 | APMoE_seg_ROB | 8 | 7 | 7 | 5 |
| | Pixel-wise Attentional Gating for Parsimonious Pixel Labeling [Project page] - Submitted by Shu Kong (University of California at Irvine) | | | | |
| 9 | BatMAN_ROB | 8 | 10 | 8 | 6 |
| | | | | | Submitted by Robert Peharz (University of Cambridge) |
| 10 | FCN101_ROB | 10 | 9 | 10 | 10 |
| | | | | | Submitted by Anonymous |

# RESULTS: REAL-TIME SEMANTIC SEGMENTATION

- Increased receptive field can compensate for model capacity

- Novel approach for increasing receptive field: pyramid fusion

- SotA real-time segmentation: lightweight ImageNet-pretrained models with increased receptive field



[orsic19cvpr]



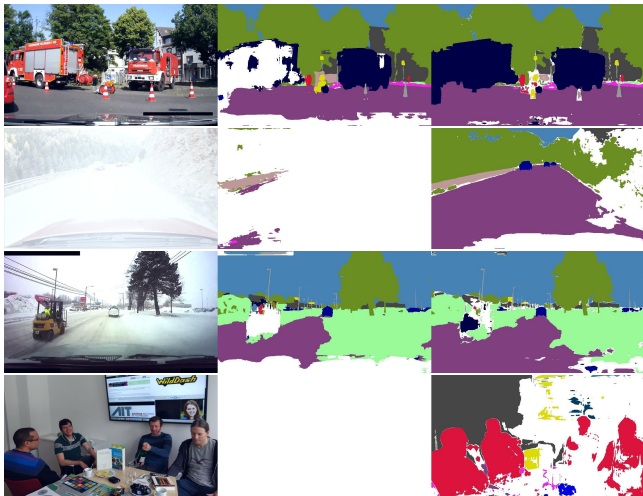[orsic19cvpr]

# RESULTS: PEDESTRIAN TRACKING

□ Deep learning for pedestrian detection and appearance representation

□ Detection: Mask R-CNN ResNet-50 pre-trained on COCO and Cityscapes

□ Appearance: correspondence embedding based on ResNet-18 trained using angular loss

□ Tracking: fusion of probabilistic data association and appearance descriptors

□ Best result on MOT2015-3D



[bicanic19review]

- □ Outlier detection = locating unknown concepts in images

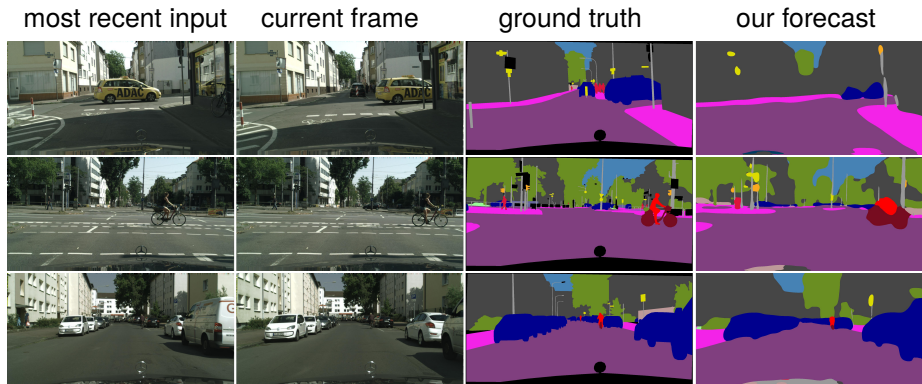- □ Fast evaluation, concurrent with semantic segmentation

# RESULTS: SEMANTIC SEGMENTATION AND OUTLIER DETECTION

☐ First place on WildDash benchmark (very difficult images)

| Model | Meta Avg | Classic | | | | Negative |
|---|---|---|---|---|---|---|
| | mIoU cla | mIoU cla | iIoU cla | mIoU cat | iIoU cat | mIoU cla |
| APMoE_seg_ROB [kong18] | 22.2 | 22.5 | 12.6 | 48.1 | 35.2 | 22.8 |
| DRN_MPC [yu17] | 28.3 | 29.1 | 13.9 | 49.2 | 29.2 | 15.9 |
| DeepLabv3+_CS [chen18] | 30.6 | 34.2 | 24.6 | 49.0 | 38.6 | 15.7 |
| LDN2_ROB [kreso18] | 32.1 | 34.4 | 30.7 | 56.6 | 47.6 | 29.9 |
| MapillaryAI_ROB [bulo17] | 38.9 | 41.3 | **38.0** | 60.5 | **57.6** | 25.0 |
| AHiSS_ROB [meletis18] | 39.0 | 41.0 | 32.2 | 53.9 | 39.3 | 43.6 |
| LDN169, two heads (ours) | 41.8 | **43.8** | 37.3 | 58.6 | 53.3 | **54.3** |
| LDN169, single OE head (ours) | **42.7** | 43.3 | 31.9 | **60.7** | 50.3 | 52.8 |

[bevandic19review]

most recent input     current frame     ground truth     our forecast

[saric19review]

# RESULTS: SEMANTIC FORECASTING

|  | mIoU | mIoU-MO |
|---|---|---|
| Oracle | 72.5 | 71.5 |
| Copy last segmentation | 38.6 | 29.6 |
| Luc Dil10-S2S | 47.8 | 40.8 |
| Luc Mask-S2S | / | 42.4 |
| Luc Mask-F2F | / | 41.2 |
| Terwilliger | 51.5 | 46.3 |
| Bhattacharyya | 51.2 | / |
| Luc F2F (our implementation) | 45.6 | 39.0 |
| our | **52.4** | **48.3** |
| our (more data) | **53.6** | **49.9** |

[saric19review]

Thank you for your attention!

Questions?