

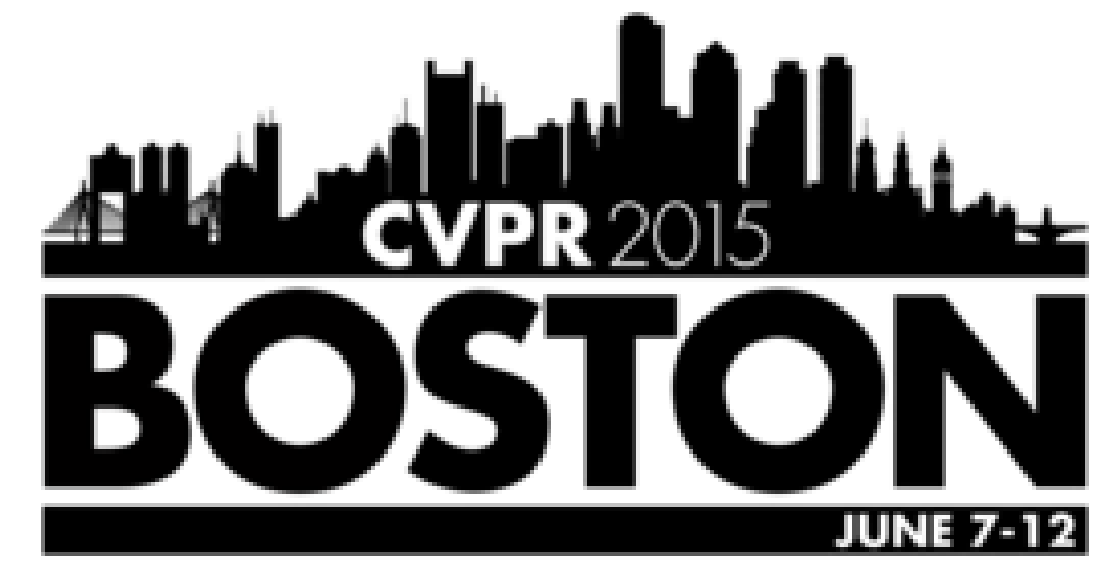


ROYAL INSTITUTE  
OF TECHNOLOGY

# From Generic to Specific Deep Representations for Visual Recognition

Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, Stefan Carlsson

KTH (Royal Institute of Technology), Computer Vision and Active Perception Lab. (CVAP), Stockholm, Sweden



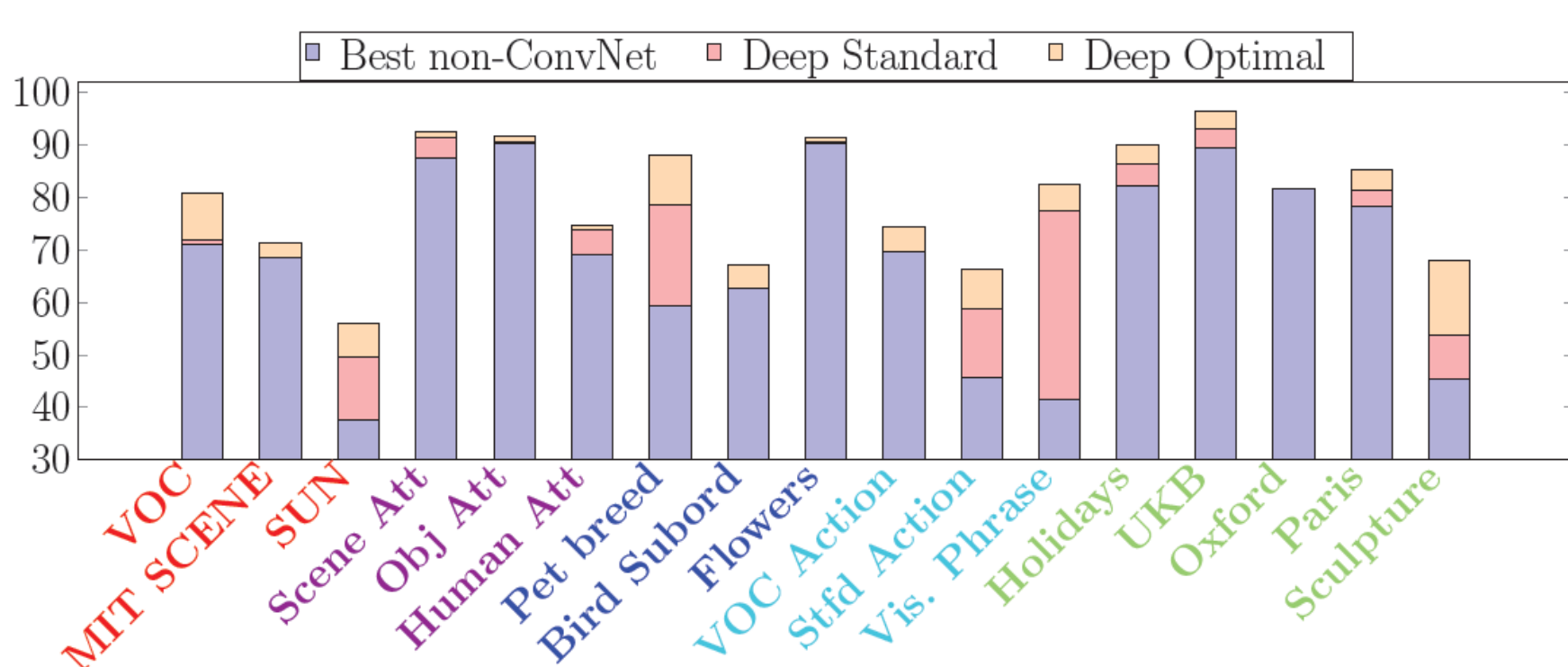
DeepVision  
DeepVision: Deep Learning in Computer Vision 2015

## Goal

- Understanding factors for transferability of a Generic ConvNet representation to different target tasks
- Analyze the correlation of the transferability factors and source to target tasks distance

## Motivation

- Numerous computer vision tasks are affected by deep learning
  - Object detection
  - Scene recognition
  - Pose Estimation
  - Semantic Segmentation
- Better ConvNet representation often beats more complicated reasoning/modeling
  - e.g. Deeper networks rep. + SVM often outperforms shallow networks representation + complicated model



## Contributions

- We propose a set of pre-training and post-training factors when transferring a generic ConvNet representation
- Categorically organize the different computer vision tasks
- Extensively study the proposed factors and their correlation with distance of source to target tasks
- state of the art performance on various (=16) recognition tasks

## Transferability Factors

- We divide the transferability factors into two groups
- We call the decisions involved before learning the generic ConvNet representation on the source task, **learning factors**
- We further identify factors which are relevant after optimizing the ConvNet on the source task: **post-learning factors**

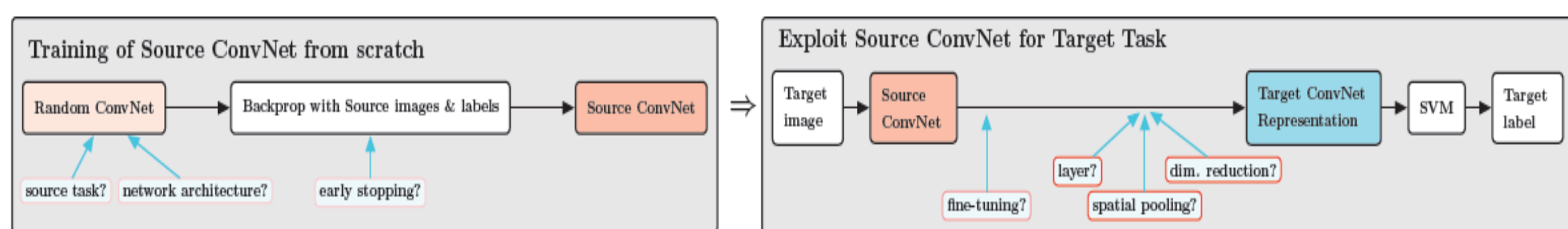
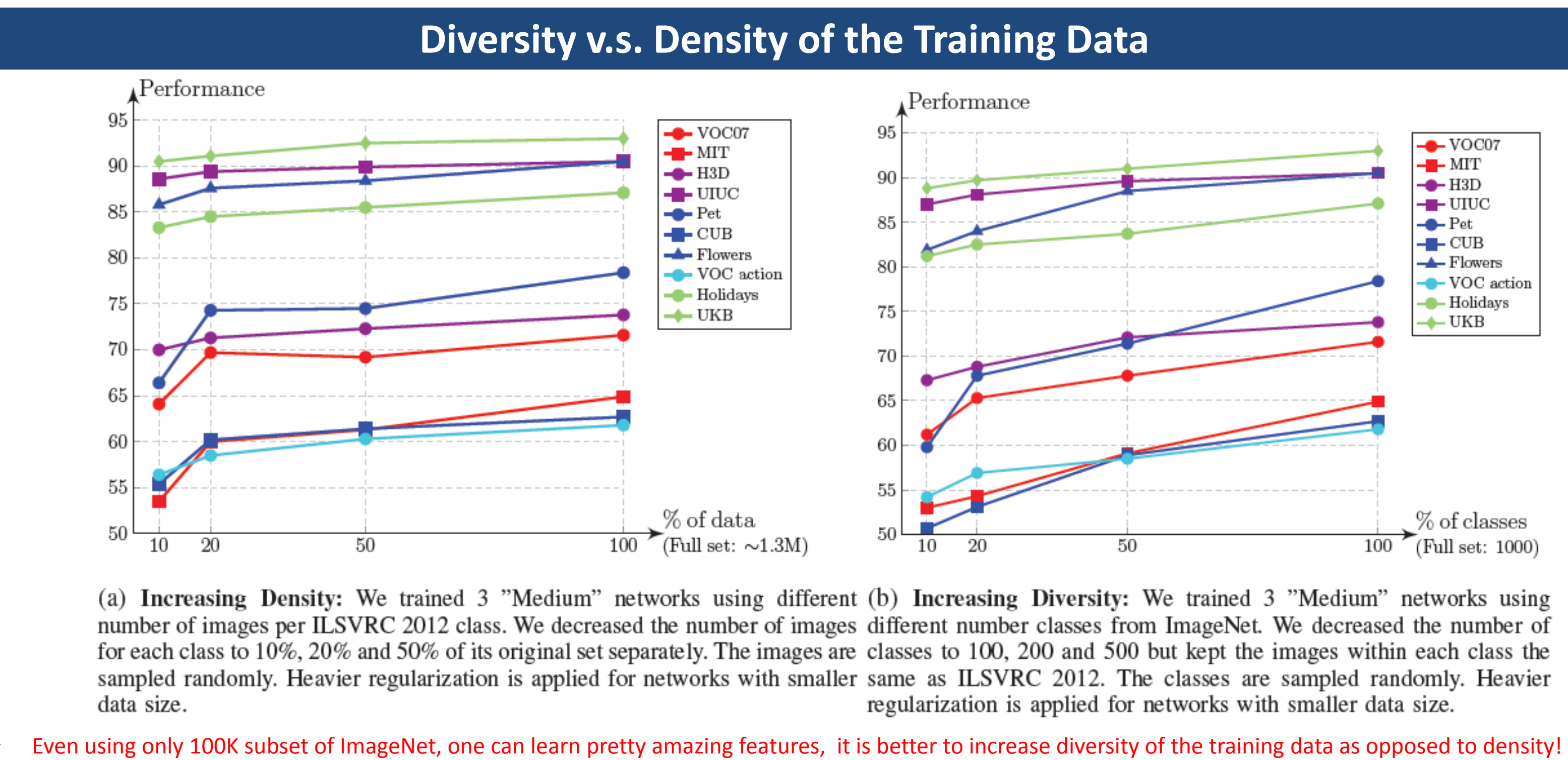
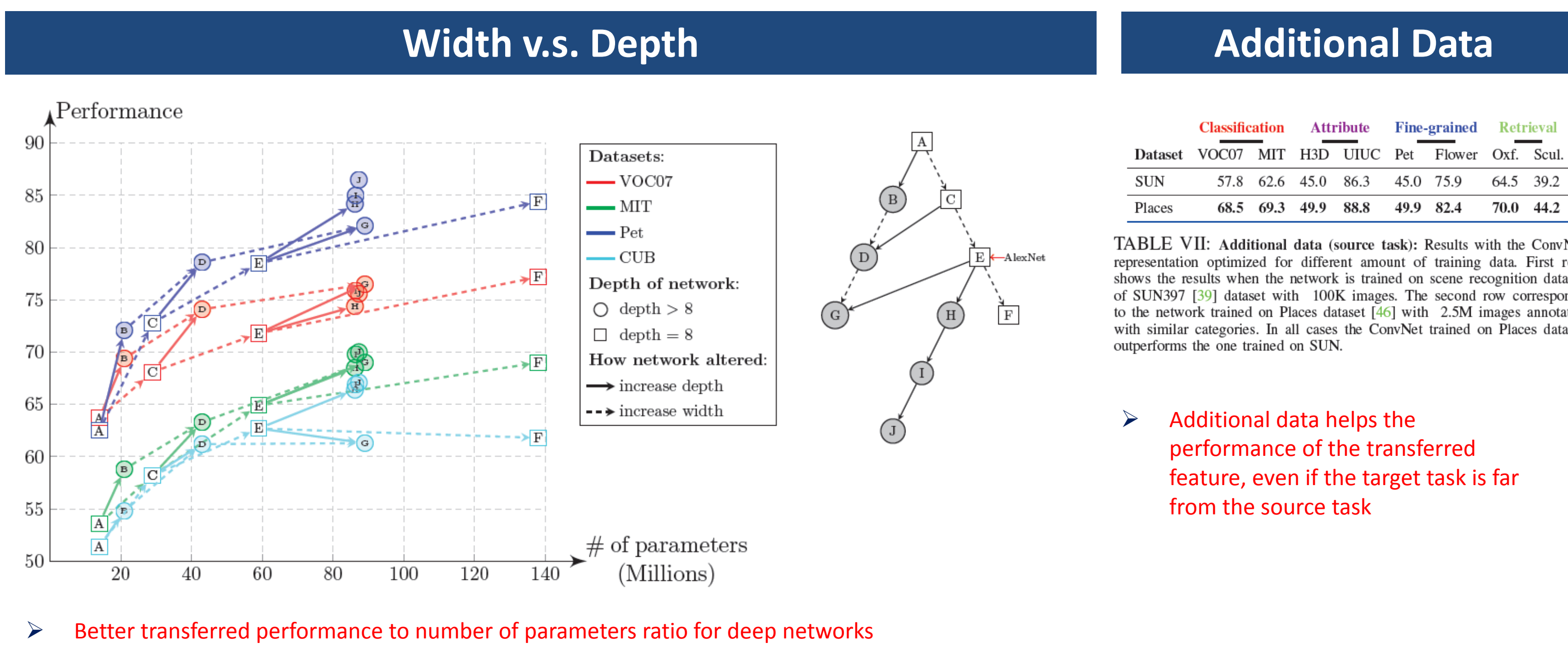
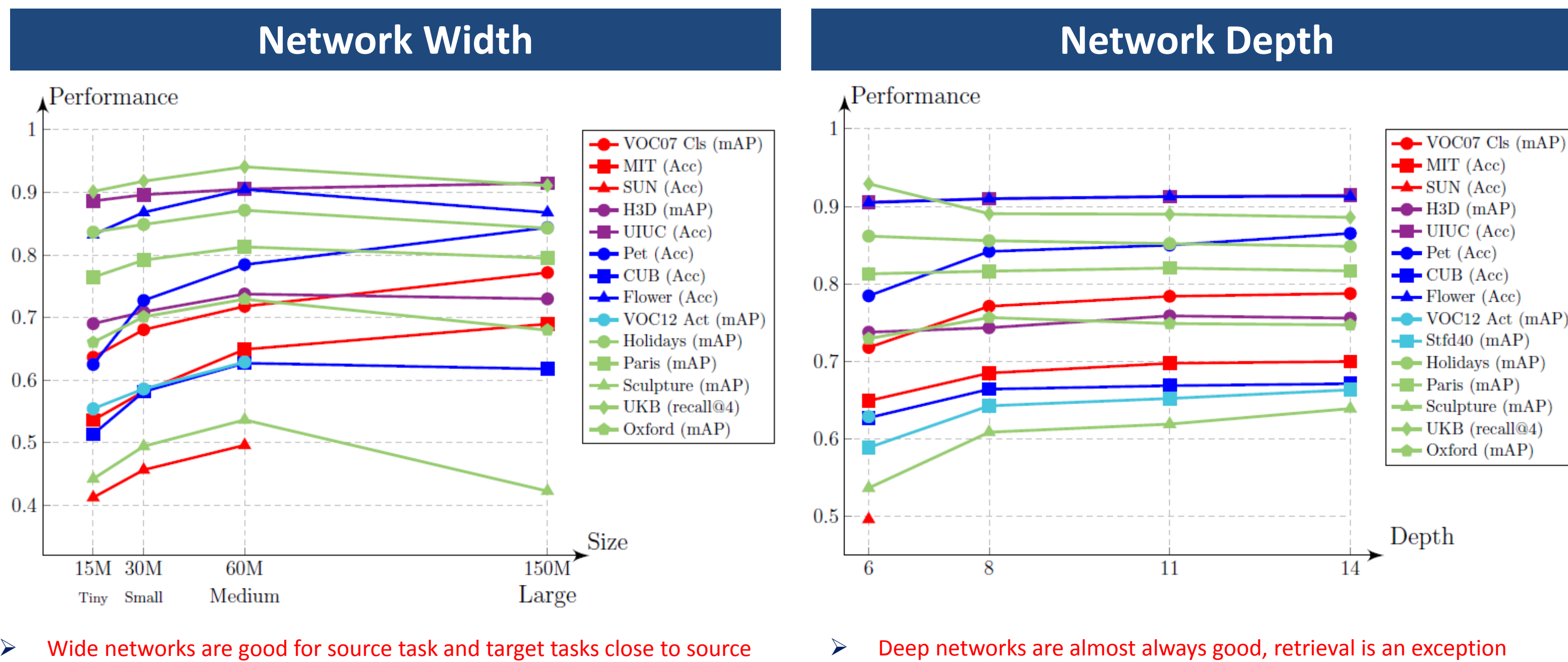
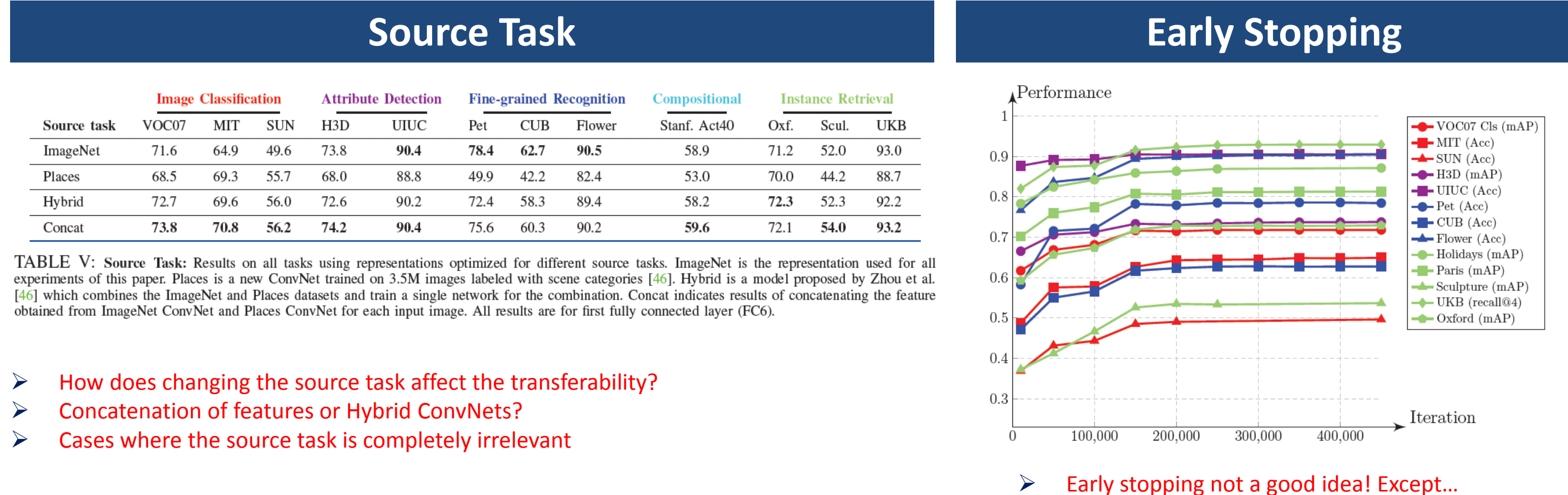


Fig. 2: Transferring a ConvNet Representation ConvNet representations are effective for visual recognition. The picture above shows the pipeline of transferring a source ConvNet representation to a target task of interest. We define several factors which control the transferability of such representations to different tasks (questions with blue arrow). These factors come into play at different stages of transfer. Optimizing these factors is crucial if one wants to

## Learning Factors



## Range of Tasks

- Distance of source-target task can be analyzed from different viewpoints and can become ambiguous. We take the following parameters into consideration:
- Target classes being super-category or sub-category of the source task
- Learning based or metric based tasks
- Explanatory classes
- Image acquisition
- ...

Image Classification	Attribute Detection	Fine-grained Recognition	Compositional	Instance Retrieval
PASCAL VOC Object [9] MIT 67 Indoor Scenes [29] SUN 397 Scene [40]	H3D human attributes [6] Object attributes [10] SUN scene attributes [26]	Cat&Dog breeds [25] Stanford 40 Actions [41] 102 Flowers [23]	VOC Human Action [9] Bird subordinate [38] Visual Phrases [30]	Holiday scenes [16] Paris buildings [27] Sculptures [4]

## Post-Learning Factors

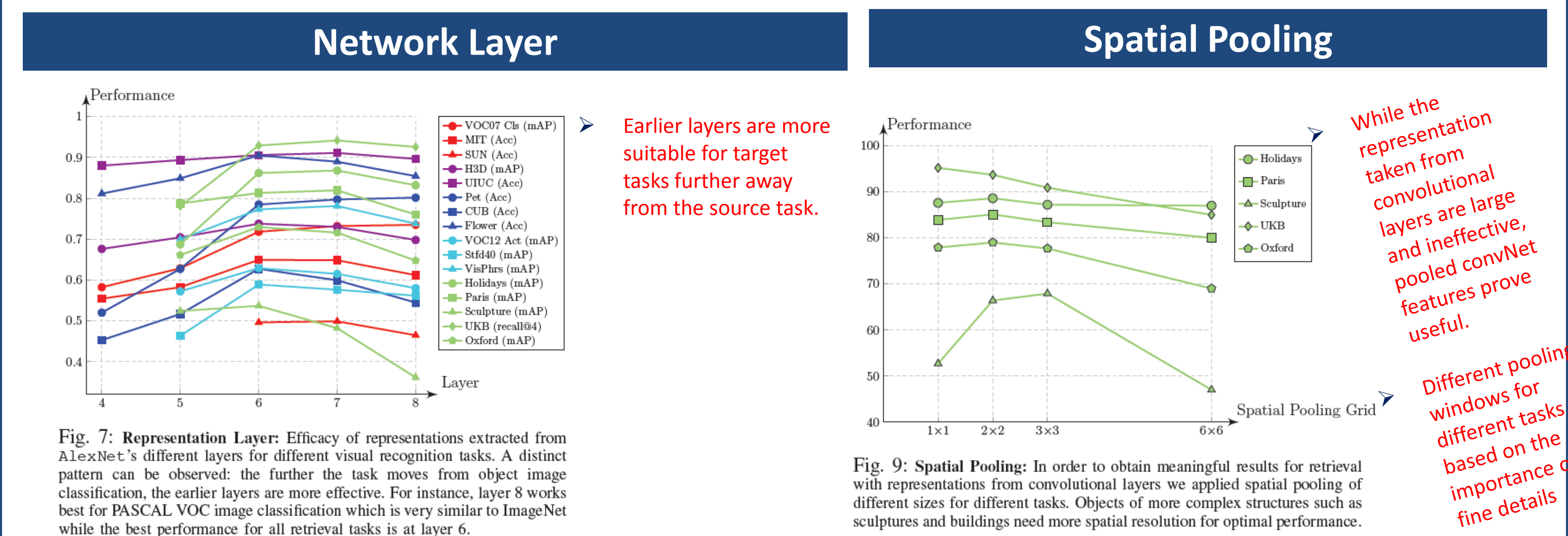
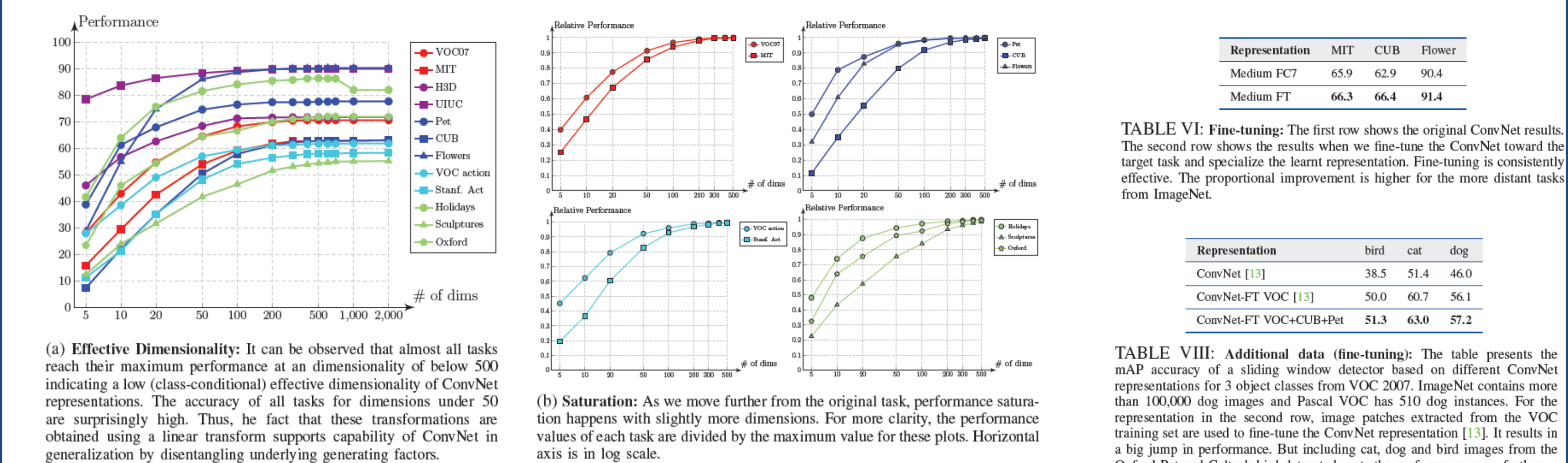


Fig. 7: Representation Layer: Efficacy of representations extracted from AlexNet's different layers for different visual recognition tasks. A distinct pattern can be observed: the further the task moves from object image classification, the earlier layers are more effective. For instance, layer 8 works best for PASCAL VOC image classification which is very similar to ImageNet while the best performance for all retrieval tasks is at layer 6.

Fig. 9: Spatial Pooling: In order to obtain meaningful results for retrieval with representations from convolutional layers we applied spatial pooling of different sizes for different tasks. Objects of more complex structures such as sculptures and buildings need more spatial resolution for optimal performance.

## Dimensionality Reduction



(a) Effective Dimensionality: It can be observed that almost all tasks reach their maximum performance at a dimensionality of below 500 indicating a low (class-conditional) effective dimensionality of ConvNet representations. The accuracy of all tasks for dimensions under 50 are surprisingly high. Thus, the fact that these representations are obtained using a linear transform supports capability of ConvNet in generalization by disentangling underlying generating factors.

(b) Saturation: As we move further from the original task, performance saturation happens with slightly more dimensions. For more clarity, the performance values of each task are divided by the maximum value for these plots. Horizontal axis is in log scale.

TABLE VI: Fine-tuning: The first row shows the original ConvNet results. The second row shows the results when we fine-tune the ConvNet toward the target task and specialize the learnt representation. Fine-tuning is consistently effective. The proportional improvement is higher for the more distant tasks from ImageNet.

Representation	MIT	CUB	Flower
ConvNet [13]	65.9	62.9	90.4
Medium FT	66.3	66.4	91.4

TABLE VIII: Additional data (fine-tuning): The table presents the mAP accuracy of a sliding window decoder based on different ConvNet representations for 3 object classes from VOC 2007. ImageNet contains more than 100,000 dog images and PASCAL VOC has 510 dog instances. For the representation in the second row, image patches extracted from the VOC training set are used to fine-tune the ConvNet representation [13]. It results in a big jump in performance. But including cat, dog and bird images from the Oxford Pet and Caltech bird datasets boosts the performance even further.

- Effective dimensionality of a ConvNet representation for various target task is between 200 to 500, closer tasks have slightly lower effective dimensionality

## Final Results Table

	Image Classification				Attribute Detection				Fine-grained Recognition				Compositional				Instance Retrieval			
	VOC07	MIT	SUN	SUNAct	UIUC	H3D	Pet	Flower	VOCa	Act40	Phrase	Holid.	UKB	Oxf.	Paris	Scul.				
non-ConvNet	[34]	[22]	[39]	[26]	[37]	[44]	[25]	[12]	[18]	[24]	[41]	[30]	[36]	[45]	[36]	[4]				
Deep Standard	71.8	64.9	49.6	91.4	90.6	73.8	78.5	62.8	90.5	69.2	58.9	77.3	86.2	93.0	73.0	81.3	53.7			
Deep Optimized <sup>†</sup>	80.7	71.3	56.0	92.5	91.5	74.6	88.1	67.1	91.3	74.3	66.4	82.3	90.0	96.3	79.0	85.1	67.9			
Err. Reduction	32%	18%	13%	13%	10%	4%	45%	12%	8%	17%	18%	22%	28%	47%	22%	20%	31%			
Source Task	ImgNet	Hybrid	Hybrid	Hybrid	Hybrid	ImgNet	ImgNet	ImgNet	ImgNet	ImgNet	ImgNet	Hybrid	ImgNet	ImgNet	ImgNet	ImgNet				
Network Depth	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium	Medium				
Rep. Layer	last	last	last	last	2nd last	2nd last	2nd last	2nd last	3rd last	3rd last	3rd last	3rd last	4th last	4th last	4th last	4th last				
Pooling	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×				
Deep Optimized MS	80.7	71.3	56.0	92.5	91.5	74.6	88.1	67.1	91.3	74.3	66.4	82.3	90.0	96.3	79.0	85.1	67.9			
Deep Optimized MS+ML	80.7	71.3	56.0	92.5	91.5	74.6	88.1	67.1	91.3	74.3	66.4	82.3	90.0	96.3	79.0	85.1	67.9			

TABLE IX: Final Results: Final results of the deep representation with optimized factors along with a linear SVM compared to the non-ConvNet state of the art. In the bottom half of the table the factors used for each task are noted. We achieve up to a 50% reduction of error by optimizing transferability factors. Relative error reductions refer to how much of the remaining error (from Deep Standard) is decreased. "Deep Standard" is the common choice of parameters - a Medium sized network of depth 8 trained on ImageNet with representation taken from layer 6 (FC6).

- These factors are important! Taking different factors into account, we achieved up to 40% reduction of classification error.