

David Novotny^{1,2}, Diane Larlus², Andrea Vedaldi¹

¹ Visual Geometry Group, University of Oxford, UK

² Naver Labs Europe, Grenoble, France

Introduction

The task:

Semantic matching

Given a pair of semantically related objects
=> estimate matches between corresponding parts

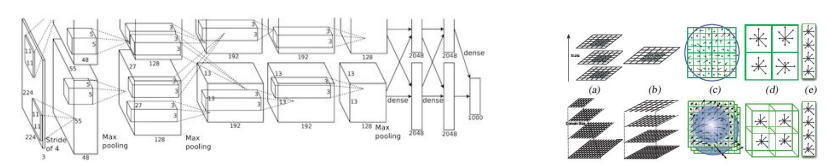
Motivation:

Fully supervised approaches [4,6] require expensive annotations / synthetic datasets => we target **weak supervision**

Weakly supervised approaches:

Step 1. Extract pixel-wise descriptors

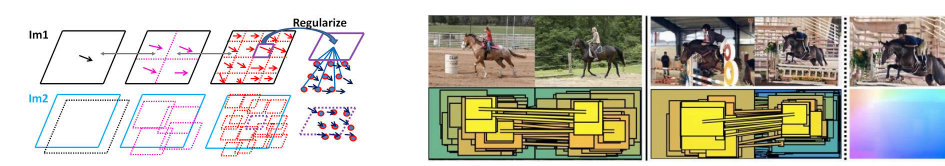
> Pretrained deep features, HoG, ...



The main focus of the paper

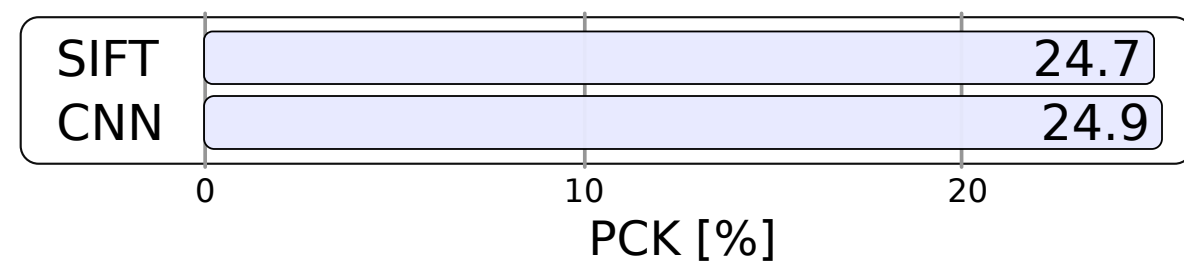
Step 2. Use a matching algorithm

> DSP [1], Proposal Flow [2], SIFT Flow [3], ...



Pixel-wise descriptors - deep vs. engineered features:

Semantic matching accuracy on Pascal VOC [5]



=> deep features on par with engineered ones

Deep features trained with a **global classification loss**

=> attention to most discriminative regions
=> invariance to geometry of the input

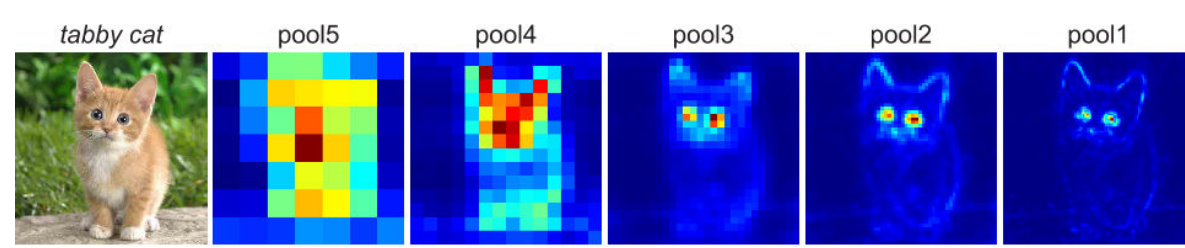
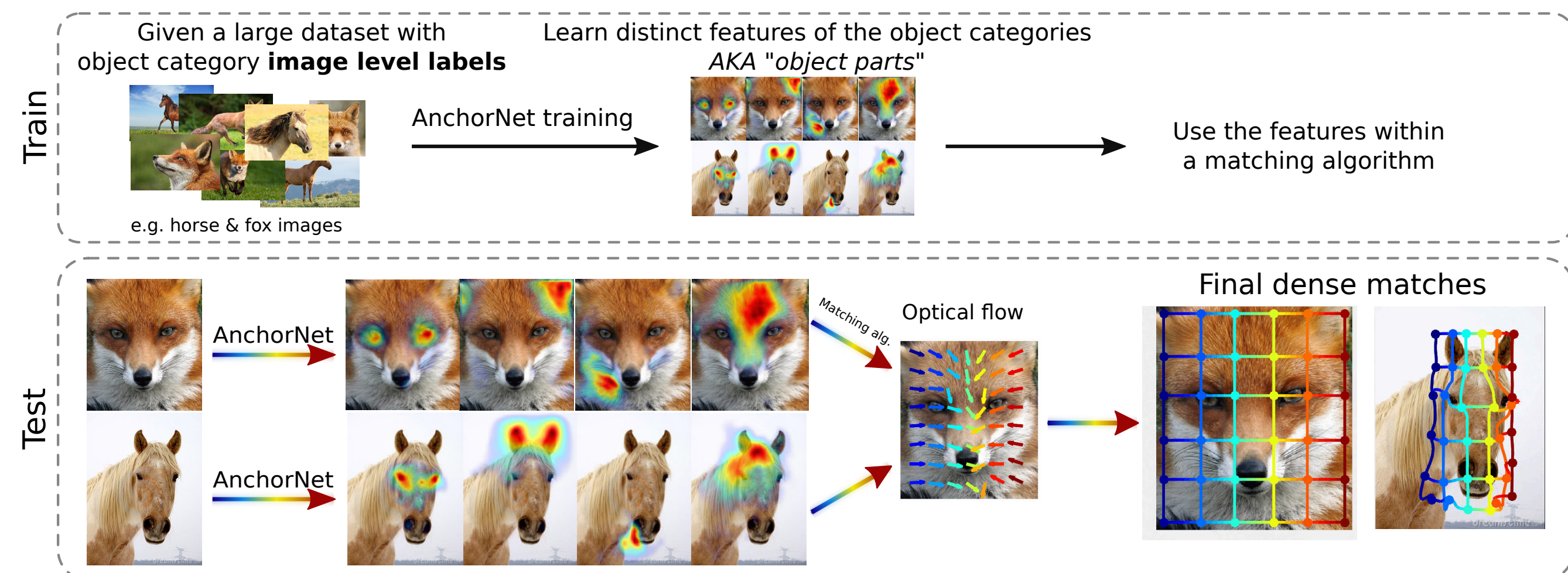


Image courtesy of [10]

=> **Main challenge:** Design the training s.t. features are not invariant to geometry

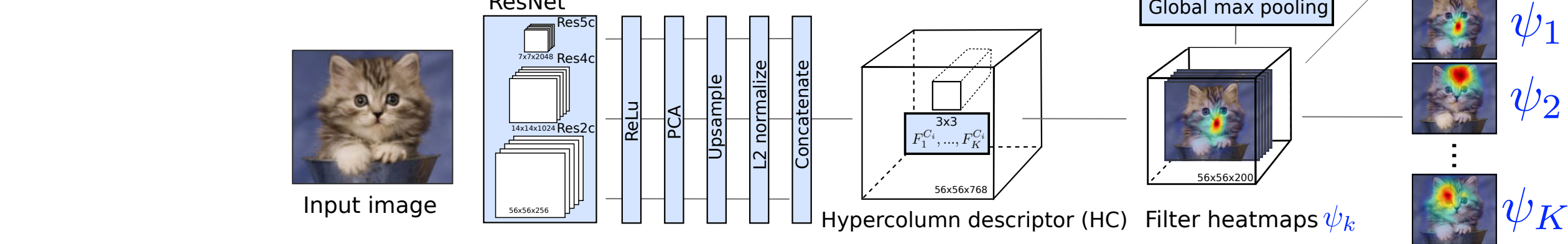
Proposed approach - overview:



AnchorNet architecture

Class specific features:

AnchorNet automatically discovers **discriminative** & **diverse** class specific features
=> discovered features = "Anchors"

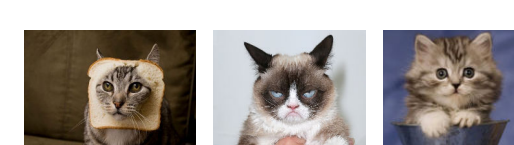


Discriminability loss

$$\mathcal{L}_{Discr} = -y_I \sum_{k=1}^K \max \psi_k$$

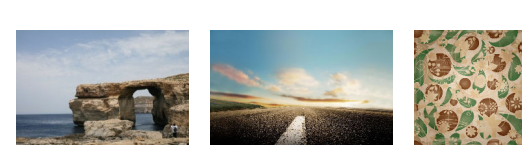
> identifies the representative of each filter by global max pooling
> increases response in positive images & decreases for negatives

Positive images



$$y_I = 1$$

Negative images



$$y_I = -1$$

Diversity loss

$$\mathcal{L}_{Div} = \sum_{i \neq j} \left\| \frac{\langle \psi_i, \psi_j \rangle}{\|\psi_i\|_F \|\psi_j\|_F} \right\|^2$$

> enforces orthogonal responses of the filters

Training **only** with Discriminability loss



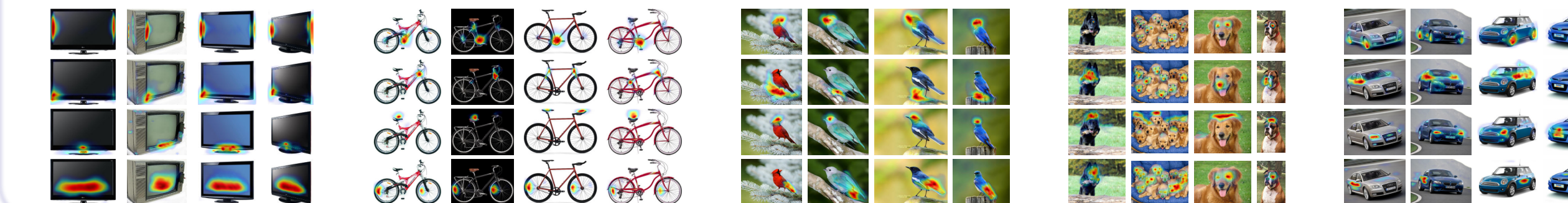
NO diversity loss => redundant features

Training with Discriminability & Diversity losses



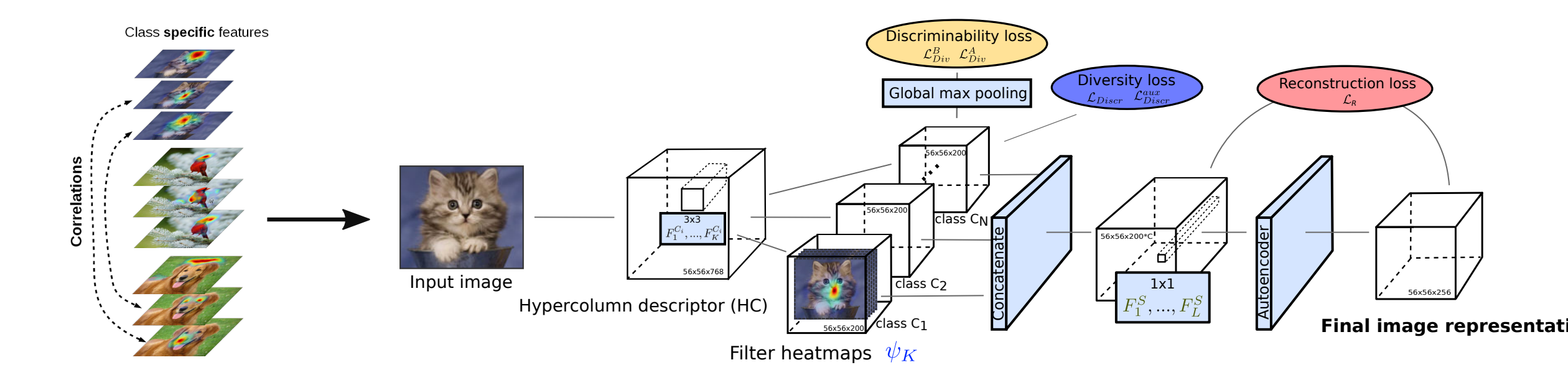
WITH diversity loss => full object coverage

Example learned class specific filters



Class agnostic features:

Class specific filters exhibit correlations => use an **autoencoder** to compress the features shared between classes



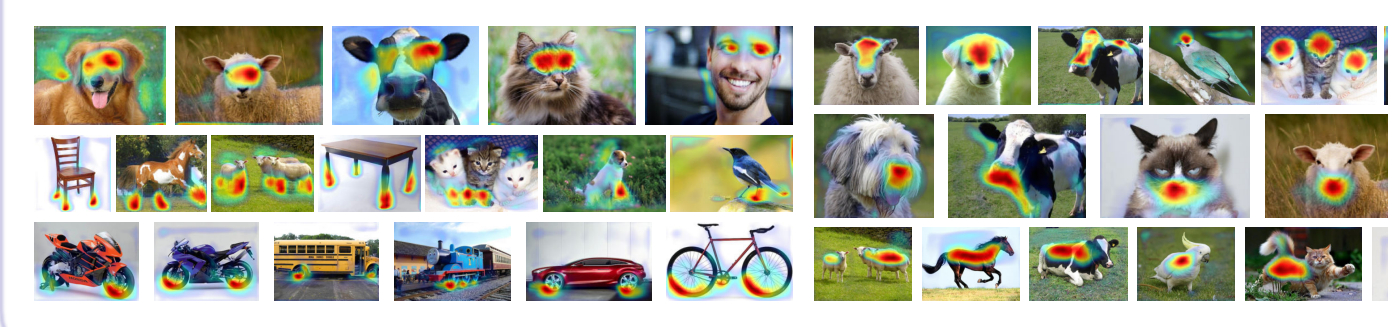
Reconstruction loss

$$\mathcal{L}_R = D(\Gamma, (F^S)^T F^S \Gamma)$$

$\Gamma = \text{stack}(\psi_1^S, \dots, \psi_K^S)$ - concatenation of the class-spec. features ψ_k

> compresses the correlated class specific features
> produces a class agnostic representation

Example learned class agnostic filters



Experiments

Semantic matching:

Given a pair of images of **the same object category**
=> estimate matches between corresponding parts

Evaluated approach

Step 1. Extract pixelwise descriptors

> **AnchorNet features**
> **ANet-class** ... class specific features
> **ANet** ... class agnostic features
> SIFT, HoG, Hypercolumns

Step 2. Match descriptors using a matching alg.

> DSP [1]
> Proposal Flow [2]

Benchmarks

Pascal Parts [7]

> part segmentation masks + keypoint annotations
> the 20 Pascal VOC classes

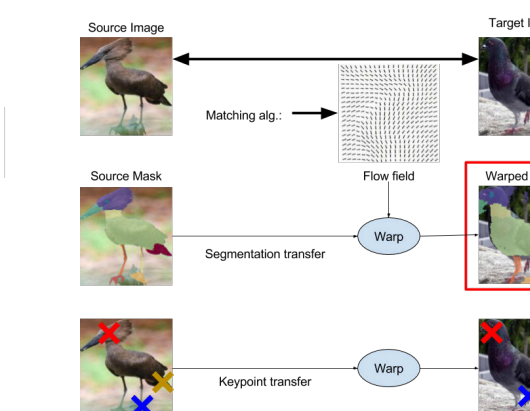


PF Dataset [2]

> dense correspondence annotations
> 6 different object classes



Evaluation procedure



Pascal Parts - segmentation transfer - IoU

362 | [ImageNet-Cow to ImageNet-Car](#) | [ImageNet-Car to ImageNet-Cow](#)

	mean	air	bird	boat	bottle	bus	car	chair	cow	dog	horse	motor	person	sheep	sofa	table	train	tv
DSP - ANet-class	0.46	0.31	0.49	0.32	0.33	0.75	0.31	0.47	0.23	0.33	0.30	0.30	0.41	0.22	0.46	0.47	0.45	0.48
DSP - ANet	0.46	0.30	0.47	0.29	0.30	0.73	0.30	0.46	0.25	0.33	0.31	0.30	0.39	0.20	0.44	0.47	0.45	0.48
DSP - SIFT [1]	0.38	0.28	0.38	0.21	0.46	0.80	0.30	0.45	0.16	0.30	0.31	0.30	0.30	0.12	0.40	0.37	0.38	0.46
Proposal Flow - ANet-class	0.43	0.30	0.43	0.28	0.34	0.71	0.30	0.45	0.24	0.32	0.31	0.30	0.35	0.21	0.45	0.46	0.44	0.50
Proposal Flow - ANet	0.42	0.29	0.41	0.30	0.33	0.70	0.29	0.45	0.25	0.31	0.31	0.30	0.31	0.17	0.43	0.39	0.44	0.50
Proposal Flow - HoG [2]	0.41	0.25	0.45	0.23	0.34	0.70	0.29	0.44	0.19	0.30	0.31	0.30	0.35	0.16	0.41	0.35	0.44	0.50

Pascal Parts - keypoint transfer - PCK@0.05 [%]

percentage of correct keypoints

	mean	air	bird	boat	bottle	bus	car	chair	cow	dog	horse	motor	person	sheep	sofa	table	train	tv
DSP - ANet-class	0.34	0.20	0.38	0.06	0.38	0.44	0.39	0.14	0.18	0.16	0.11	0.11	0.41	0.11	0.37	0.37	0.37	
DSP - ANet	0.23	0.20	0.26	0.06	0.38	0.42	0.34	0.14	0.17	0.17	0.13	0.13	0.38	0.11	0.37	0.37	0.37	
DSP - SIFT [1]	0.19	0.17	0.20	0.05	0.19	0.33	0.34	0.09	0.17	0.12	0.09	0.10	0.35	0.11	0.37	0.37	0.37	
Proposal Flow - ANet-class	0.31	0.19	0.35	0.22	0.30	0.42	0.10	0.14	0.11	0.07	0.10	0.09	0.38	0.11	0.37	0.37	0.37	
Proposal Flow - ANet	0.30	0.19	0.30	0.22	0.28	0.40	0.10	0.14	0.11	0.09	0.10	0.09	0.38	0.11	0.37	0.37	0.37	
Proposal Flow - HoG [2]	0.27	0.20	0.26	0.05	0.20	0.31	0.29	0.10	0.17	0.13	0.05	0.13	0.21	0.11	0.37	0.37	0.37	

PF Dataset

Accuracy for PFH

	mean	air	bird	boat	bottle	bus	car	chair	cow	dog	horse	motor	person	sheep	sofa	table	train	tv
DSP - ANet-class	0.34	0.20	0.38	0.06	0.38	0.44	0.39	0.14	0.18	0.16	0.11	0.11	0.41	0.11	0.37	0.37	0.37	
DSP - ANet	0.23	0.20	0.26	0.06	0.38	0.42	0.34	0.14	0.17	0.17	0.13	0.13	0.38	0.11	0.37	0.37	0.37	
DSP - SIFT [1]	0.19	0.17	0.20	0.05	0.19	0.33	0.34	0.09	0.17	0.12	0.09	0.10	0.35	0.11	0.37	0.37	0.37	
Proposal Flow - ANet-class	0.31	0.19	0.35	0.22	0.30	0.42	0.10	0.14	0.11	0.07	0.10	0.09	0.38	0.11	0.37	0.37	0.37	
Proposal Flow - ANet	0.30	0.19	0.30	0.22	0.28	0.40	0.10	0.14	0.11	0.09	0.10	0.09	0.38	0.11	0.37	0.37	0.37	
Proposal Flow - HoG [2]	0.27	0.20	0.26	0.05	0.20	0.31	0.29	0.10	0.17	0.13	0.05	0.13	0.21	0.11	0.37	0.37	0.37	

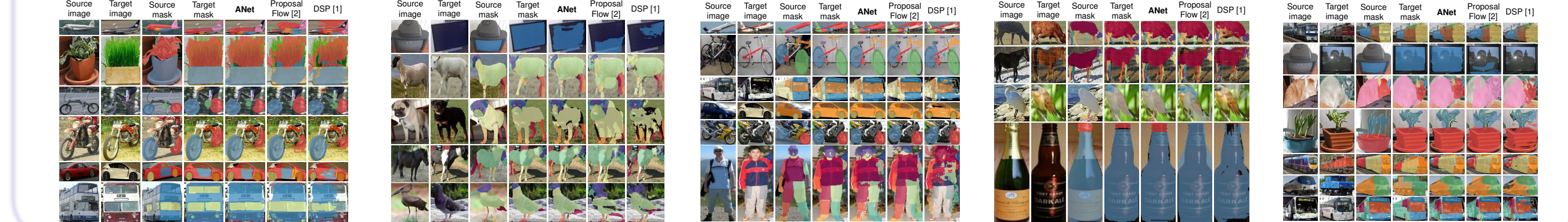
→ State-of-the-art performance on the segmentation transfer task

-> AnchorNet features improve performance of matching algorithms

-> AnchorNet vs AnchorNet-class perform on par => successful conversion from class specific to class agnostic features

-> PF Dataset - naive matching of AnchorNet features similar to matching engineered features with a sophisticated algorithm

Pascal VOC - qualitative results



Cross-class semantic matching:

Given a pair of images of **related object categories**
=> estimate matches between corresponding parts

Benchmarks

Pascal Parts [7]

> related object classes share part segmentations
> e.g. "car" and "bus" classes share "wheel", "door", "window", ... parts
> evaluation of segmentation transfer between images of meaningful classes
> low number of categories, large number of shared parts

Animal Parts [8]

> images from the ImageNet dataset annotated with "eye" and "foot" keypoints
> evaluation of keypoint transfer between images from different animal domains
> large number of categories, low number of parts



Animal Parts - Cross class keypoint matching accuracy - PCK

Transfer accuracy between animal domains

Heatmap showing PCK transfer accuracy between animal domains. The color scale ranges from 0.00 (blue) to 0.40 (red). The rows represent the source domain and the columns represent the target domain. The diagonal elements are 1.00, indicating perfect self-transfer.

	air	bird	boat	bottle	bus	car	chair	cow	dog	horse	motor	person	sheep	sofa	table	train	tv
AnchorNet	0.34	0.20	0.38	0.06	0.38	0.44	0.39	0.14	0.18	0.16	0.11	0.11	0.41	0.11	0.37	0.37	0.37
SIFT [1]	0.23	0.20	0.26	0.06	0.38	0.42	0.34	0.14	0.17	0.17	0.13	0.13	0.38	0.11	0.37	0.37	0.37
Proposal Flow	0.31	0.19	0.35	0.22	0.30	0.42	0.10	0.14	0.11	0.07	0.10	0.09	0.38	0.11	0.37	0.37	0.37

Mean PCK over all matches

Matching Map	DSP	Proposal Flow		
Feature	Anchor	SIFT [1]	Anchor	SIFT [1]
PCK (0 - 0.5)	0.31	0.26	0.33	0.28
PCK (0 - 0.1)	0.24	0.17	0.32	0.18

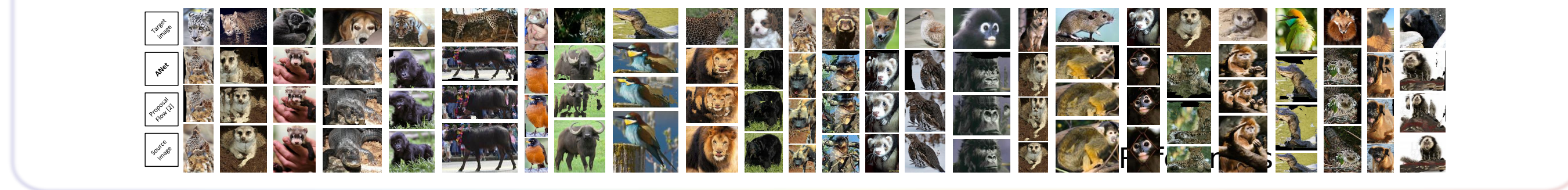
Pascal Parts - Cross-class segmentation transfer - IoU

> Mean IoU over transferred segmentations of shared parts

	air	bird	boat	bottle	bus	car	chair	cow	dog	horse	motor	person	sheep	sofa	table	train	tv
AnchorNet	0.45	0.31	0.49	0.32	0.33	0.75	0.31	0.47	0.23	0.33	0.30	0.30	0.41	0.22	0.46	0.47	0.45
SIFT [1]	0.38	0.28	0.38	0.21	0.46	0.80	0.30	0.45	0.16	0.30	0.31	0.30	0.30	0.12	0.40	0.37	0.38
Proposal Flow	0.43	0.30	0.43	0.28	0.34	0.71	0.30	0.45	0.24	0.32	0.31	0.30	0.35	0.21	0.45	0.46	0.44
Proposal Flow - ANet-class	0.42	0.29	0.41	0.30	0.33	0.70	0.29	0.45	0.25	0.31	0.31	0.30	0.31	0.17	0.43	0.39	0.44
Proposal Flow - HoG [2]	0.41	0.25	0.45	0.23	0.34	0.70	0.29	0.44	0.19	0.30	0.31	0.30	0.35	0.16	0.41	0.35	0.44

> AnchorNet features bring significant improvement over considered baseline features
> State-of-the-art performance on both datasets

Animal Parts - qualitative results



Conclusions

- > Proposed **AnchorNet** - a weakly supervised architecture for learning geometry-sensitive features
- > The learned features are invariant to appearance making them suitable for semantic matching tasks
- > Experimentally verified that the features improve performance of existing matching algorithms
- > State-of-the-art performance on semantic matching and on novel cross-class semantic matching task

References

- [1] Kim et al. "Deformable spatial pyramid pooling for fast dense correspondences." CVPR 2018.
- [2] Hartmann et al. "Proposal flow: CVPR 2018.
- [3] Liu et al. "Sift flow: Dense correspondence across scenes and its applications." CVPR 2011.
- [4] Long et al. "Universal correspondence network." NIPS 2016.
- [5] Long et al. "Do convnets learn correspondence?" NIPS 2014.
- [6] Zhou et al. "Learning dense correspondence via 3d-guided cycle consistency." CVPR 2016.
- [7] Chen et al. "Detect what you can: Detecting and representing objects using holistic models." CVPR 2015.
- [8] Hartmann et al. "Hypercolumns for object segmentation and fine-grained localization." CVPR 2015.
- [9] Zhang et al. "Region-wise neural attention by location background." ECCV 2016.