

SIZER: A Dataset and Model for Parsing  
3D Clothing and Learning Size Sensitive  
3D Clothing

Garvita Tiwari<sup>1(B)</sup>, Bharat Lal Bhatnagar<sup>1</sup>, Tony Tung<sup>2</sup>,  
and Gerard Pons-Moll<sup>1</sup>

<sup>1</sup>MPI for Informatics, Saarland Informatics Campus, Saarbrücken, Germany

{gtiwari,bbhatnag,gpons}@mpi-inf.mpg.de

<sup>2</sup>Facebook Reality Labs, Sausalito, USA

tony.tung@fb.com

Fig. 1. SIZER dataset of people with clothing size variation. (Left): 3D scans  
of people

captured in different clothing styles and sizes.(Right): T-shirt and short pants  
for sizes

small and large, which are registered to a common template.

Abstract. While models of 3D clothing learned from real data exist, no  
method can predict clothing deformation as a function of garment size.

In this paper, we introduce SizerNet to predict 3D clothing conditioned  
on human body shape and garment size parameters, and ParserNet to infer garment m  
eshes and shape under clothing with personal details in a

single pass from an input mesh. SizerNet allows to estimate and visualize  
the dressing effect of a garment in various sizes, and ParserNet allowsto edit cl  
othing of an input mesh directly, removing the need for scan

segmentation, which is a challenging problem in itself. To learn these  
models, we introduce the SIZER dataset of clothing size variation which  
includes 100 different subjects wearing casual clothing items in various

sizes, totaling to approximately 2000 scans. This dataset includes the  
scans, registrations to the SMPL model, scans segmented in clothing  
parts, garment category and size labels. Our experiments show better

Electronic supplementary material The online version of this chapter ([https://doi.org/10.1007/978-3-030-58580-8\\_1](https://doi.org/10.1007/978-3-030-58580-8_1)) contains supplementary material, which is a  
vail-

able to authorized users.

©/circlecopyrtSpringer Nature Switzerland AG 2020

A. Vedaldi et al. (Eds.): ECCV 2020, LNCS 12348, pp. 1–18, 2020.<https://doi.org/10.1007/978-3-030-58580-8>

\*\*\*\*\*

LIMP: Learning Latent Shape  
Representations with Metric  
Preservation Priors

Luca Cosmo<sup>1,2(B)</sup>, Antonio Norelli<sup>1</sup>, Oshri Halimi<sup>3</sup>, Ron Kimmel<sup>3</sup>,  
and Emanuele Rodolà<sup>1</sup>

<sup>1</sup>Sapienza University of Rome, Rome, Italy

cosmo@di.uniroma1.it

<sup>2</sup>University of Lugano, Lugano, Switzerland

<sup>3</sup>Technion - Israel Institute of Technology, Haifa, Israel

Abstract. In this paper, we advocate the adoption of metric preserva-  
tion as a powerful prior for learning latent representations of deformable  
3D shapes. Key to our construction is the introduction of a geometricdistortion  
criterion, defined directly on the decoded shapes, translating  
the preservation of the metric on the decoding to the formation of linear  
paths in the underlying latent space. Our rationale lies in the observa-tion tha  
t training samples alone are often insufficient to endow generative  
models with high fidelity, motivating the need for large training datasets.

In contrast, metric preservation provides a rigorous way to control theamount of  
geometric distortion incurring in the construction of the latent  
space, leading in turn to synthetic samples of higher quality. We further  
demonstrate, for the first time, the adoption of differentiable intrinsicdistances  
in the backpropagation of a geodesic loss. Our geometric pri-  
ors are particularly relevant in the presence of scarce training data, where

learning any meaningful latent structure can be especially challenging. The effectiveness and potential of our generative model is showcased in applications of style transfer, content generation, and shape completion.

Keywords: Learning shapes

•Generative model •Metric distortion

1

\*\*\*\*\*

Unsupervised Sketch to Photo Synthesis

Runtao Li<sup>1</sup>, Qian Yu<sup>1,2(B)</sup>, and Stella X. Yu<sup>1</sup>

<sup>1</sup>UC Berkeley/ICSI, Berkeley, USA

qianyu@buaa.edu.cn

<sup>2</sup>Beihang University,

Xueyuan Rd. No. 37, Haidian District, Beijing, China

**Abstract.** Humans can envision a realistic photo given a free-hand sketch that is not only spatially imprecise and geometrically distorted but also without colors and visual details. We study unsupervised sketch to photo synthesis for the first time, learning from unpaired sketch and photo data where the target photo for a sketch is unknown during training. Existing works only deal with either style difference or spatial deformation alone, synthesizing photos from edge-aligned line drawings or transforming shapes within the same modality, e.g., color images. Our insight is to decompose the unsupervised sketch to photo synthesis task into two stages of translation: First shape translation from sketches to grayscale photos and then content enrichment from grayscale to color photos. We also incorporate a self-supervised denoising objective and an attention module to handle abstraction and style variations that are specific to sketches. Our synthesis is sketch-faithful and photo-realistic, enabling sketch-based image retrieval and automatic sketch generation that captures human visual perception beyond the edge map of a photo.

1

\*\*\*\*\*

Unsupervised Sketch to Photo Synthesis 47

(a) (b) (c) (d) (a) (b) (c) (d) (e) (f)

Fig. 7. Left: Synthesized results when the edge map is used as the intermediate goal

instead of the grayscale photo. (a) Input sketch; (b) Synthesized edge map, (c) Synthesized RGB photo using the edge map; (d) Synthesized RGB photo using grayscale (Ours).

Right: Our model can successfully deal with noise sketches, which are not well handled

by another attention-based model, UGATIT. For an input sketch (a), our model produce an attention mask (b); (c) and (d) are grayscale images produced by vanilla and

our model. (e) and (f) compare ours with the result of UGATIT. (Color figure online)

Fig. 8. Comparisons of paired and unpaired training for shape translation. There are

four examples. For each example, the 1st one is the input sketch, the 2nd and the 3rd are grayscale images synthesized by Pix2Pix and our model respectively. Note that for

each example, although the input sketches are different visually, Pix2Pix produces a

similar-looking grayscale image. Our results are more faithful to the sketch.

#### 4.3 Ablation Study

**Two-Stage Architecture.** Two-stage architecture is the key to the success of our model. This strategy can be easily adapted by other models such as cycleGAN. Table 2 compares the performance of the original cycleGAN and its two-stage version (i.e., cycleGAN is used only for shape translation while the content

enrichment network is the same as ours). The two-stage version outperforms the original cycleGAN by 27.55 (on ShoeV2) and 68.33 (on ChairV2), indicating the significant benefits brought by this architectural design.

Edge Map vs. Grayscale as the Intermediate Goal. We choose grayscale as our intermediate goal of translation. As shown in Fig. 1, edge maps could be an alternative since it does not have shape deformation either. We can first translate sketch to an edge map, and then fill the edge map with colorful details

Table 2 and Fig. 7 show that using the edge map is worse than using the grayscale. Our explanations are: 1) Grayscale images contain more visual details thus can provide more learning signals for training shape translation network; 2) Content enrichment is easier for grayscale as they are closer to color photos than edge maps. The grayscale is also easier to obtain in practice.

Deal with Abstraction and Style Variations. We have discussed the problem encountered during shape translation in Sect. 3.1, and further introduced 1) a self-supervised objective along with noise sketch composition strategies and

\*\*\*\*\*

A Simple Way to Make Neural Networks

Robust Against Diverse Image

Corruptions

Evgenia Rusak<sup>1,2(B)</sup>, Lukas Schott<sup>1,2</sup>, Roland S. Zimmermann<sup>1,2</sup>,  
Julian Bitterwolf<sup>2</sup>, Oliver Bringmann<sup>1</sup>, Matthias Bethge<sup>1,2</sup>,  
and Wieland Brendel<sup>1,2</sup>

<sup>1</sup>University of Tübingen, Tübingen, Germany

{evgenia.rusak, lukas.schott, roland.zimmermann,  
oliver.bringmann, matthias.bethge, wieland.brendel}@uni-tuebingen.de

<sup>2</sup>International Max Planck Research School for Intelligent Systems,  
Tübingen, Germany

julian.bitterwolf@uni-tuebingen.de

**Abstract.** The human visual system is remarkably robust against a wide range of naturally occurring variations and corruptions like rain or snow. In contrast, the performance of modern image recognition models strongly degrades when evaluated on previously unseen corruptions. Here, we demonstrate that a simple but properly tuned training with additive Gaussian and Speckle noise generalizes surprisingly well to unseen corruptions, easily reaching the state of the art on the corruption benchmark ImageNet-C (with ResNet50) and on MNIST-C. We build on top of these strong baseline results and show that an adversarial training of the recognition model against locally correlated worst-case noise distributions leads to an additional increase in performance. This regularization can be combined with previously proposed defense methods for further improvement.

**Keywords:** Image corruptions

· Robustness · Generalization ·

Adversarial training

1

\*\*\*\*\*

SoftPoolNet: Shape Descriptor for Point

Cloud Completion and Classification

Yida Wang<sup>1(B)</sup>, David Joseph Tan<sup>2</sup>, Nassir Navab<sup>1</sup>, and Federico Tombari<sup>1,2</sup>

<sup>1</sup>Technische Universität München, München, Germany

yida.wang@tum.de

<sup>2</sup>Google Inc., Menlo Park, USA

**Abstract.** Point clouds are often the default choice for many applications as they exhibit more flexibility and efficiency than volumetric data. Nevertheless, their unorganized nature – points are stored in an unordered way – makes them less suited to be processed by deep learning pipelines. In this paper, we propose a method for 3D object completion and classification based on point clouds. We introduce a new way of organizing the extracted features based on their activations, which we name

soft pooling. For the decoder stage, we propose regional convolutions, a novel operator aimed at maximizing the global activation entropy. Furthermore, inspired by the local re $\square$ ning procedure in Point Completion Network (PCN), we also propose a patch-deforming operation to simulate deconvolutional operations for point clouds. This paper proves that our regional activation can be incorporated in many point cloud architectures like AtlasNet and PCN, leading to better performance for geometric completion. We evaluate our approach on different 3D tasks such as object completion and classification, achieving state-of-the-art accuracy.

1

\*\*\*\*\*

Hierarchical Face Aging Through

Disentangled Latent Characteristics

Peipei Li<sup>1,3</sup>, Huaibo Huang<sup>1,3,4</sup>, Yibo Hu<sup>1</sup>, Xiang Wu<sup>1</sup>, Ran He<sup>1,2,3(B)</sup>, and Zhenan Sun<sup>1,2,3</sup>

<sup>1</sup>Center for Research on Intelligent Perception and Computing, NLPR, CASIA, Beijing, China

{peipei.li, huaibo.huang}@cripac.ia.ac.cn, huyibo871079699@gmail.com, alfredxiangwu@gmail.com, {rhe, znsun}@nlpr.ia.ac.cn

<sup>2</sup>Center for Excellence in Brain Science and Intelligence Technology, CAS, Beijing, China

<sup>3</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China

<sup>4</sup>Artificial Intelligence Research, CAS, Jiaozhou, Qingdao, China

**Abstract.** Current age datasets lie in a long-tailed distribution, which brings difficulties to describe the aging mechanism for the imbalance ages. To alleviate it, we design a novel facial age prior to guide the aging mechanism modeling. To explore the age effects on facial images, we propose a Disentangled Adversarial Autoencoder (DAAE) to disentangle the facial images into three independent factors: age, identity and extraneous information. To avoid the "wash away" of age and identity information in face aging process, we propose a hierarchical conditional generator bypassing the disentangled identity and age embeddings to the high-level and low-level layers with class-conditional BatchNorm. Finally, a disentangled adversarial learning mechanism is introduced to boost the image quality for face aging. In this way, when manipulating the age distribution, DAAE can achieve face aging with arbitrary ages. Further, given an input face image, the mean value of the learned age posterior distribution can be treated as an age estimator. These indicate that DAAE can efficiently and accurately estimate the age distribution in a disentangling manner. DAAE is the first attempt to achieve facial age analysis tasks, including face aging with arbitrary ages, exemplar-based face aging and age estimation, in a universal framework. The qualitative and quantitative experiments demonstrate the superiority of DAAE on five popular datasets, including CACD2000, Morph, UTKFace, FG-NET and AgeDB.

**Keywords:** Facial age analysis

• Variational autoencoder

P. Li, H. Huang and R. He—Equal contribution.

Electronic supplementary material The online version of this chapter ([https://doi.org/10.1007/978-3-030-58580-8\\_6](https://doi.org/10.1007/978-3-030-58580-8_6)) contains supplementary material, which is available to authorized users.

©/circlecopyrt Springer Nature Switzerland AG 2020

A. Vedaldi et al. (Eds.): ECCV 2020, LNCS 12348, pp. 86–101, 2020.

<https://doi.org/10.1007/978-3-030-58580-8>

—

\*\*\*\*\*

Hybrid Models for Open Set Recognition

Hongjie Zhang<sup>1</sup>, Ang Li<sup>2</sup>, Jie Guo<sup>1</sup>, and Yanwen Guo<sup>1(B)</sup>

<sup>1</sup>State Key Laboratory for Novel Software Technology, Nanjing University,

Nanjing 210023, China  
hjzhang@smail.nju.edu.cn, {guojie,ywguo}@nju.edu.cn  
2DeepMind, Mountain View, CA, USA  
anglilli@google.com

**Abstract.** Open set recognition requires a classifier to detect samples not belonging to any of the classes in its training set. Existing methods fit a probability distribution to the training samples on their embedding space and detect outliers according to this distribution. The embeddingspace is often obtained from a discriminative classifier. However, such discriminative representation focuses only on known classes, which may not be critical for distinguishing the unknown classes. We argue that the representation space should be jointly learned from the inlier classifier and the density estimator (served as an outlier detector). We propose theOpenHybrid framework, which is composed of an encoder to encode the input data into a joint embedding space, a classifier to classify samples to inlier classes, and a flow-based density estimator to detect whether a sample belongs to the unknown category. A typical problem of existing flow-based models is that they may assign a higher likelihood to outliers. However, we empirically observe that such an issue does not occur in our experiments when learning a joint representation for discriminative and generative components. Experiments on standard open set benchmarks also reveal that an end-to-end trained OpenHybrid model significantly outperforms state-of-the-art methods and flow-based baselines.

**Keywords:** Flow-based model  
·Density estimation ·Image classification

1

\*\*\*\*\*  
**TopoGAN: A Topology-Aware Generative Adversarial Network**

Fan Wang(B), Huidong Liu, Dimitris Samaras, and Chao Chen  
Stony Brook University, Stony Brook, NY 11794, USA  
{fanwang1,huidliu,samaras}@cs.stonybrook.edu, chao.chen.1@stonybrook.edu

**Abstract.** Existing generative adversarial networks (GANs) focus on generating realistic images based on CNN-derived image features, but fail to preserve the structural properties of real images. This can be fatal in applications where the underlying structure (e.g., neurons, vessels, membranes, and road networks) of the image carries crucial semantic meaning. In this paper, we propose a novel GAN model that learns the topology of real images, i.e., connectedness and loopiness. In particular, we introduce a new loss that bridges the gap between synthetic image distribution and real image distribution in the topological feature space. By optimizing this loss, the generator produces images with the same structural topology as real images. We also propose new GAN evaluation metrics that measure the topological realism of the synthetic images. We show in experiments that our method generates synthetic images with realistic topology. We also highlight the increased performance that our method brings to downstream tasks such as segmentation.

**Keywords:** Topology  
·Persistent homology ·Generative Adversarial Network

1

\*\*\*\*\*  
**Learning to Localize Actions from Moments**

Fuchen Long<sup>1</sup>, Ting Yao<sup>2(B)</sup>, Zhaofan Qiu<sup>1</sup>, Xinmei Tian<sup>1</sup>, Jiabo Lu<sup>3</sup>, and Tao Mei<sup>2</sup>

<sup>1</sup>University of Science and Technology of China, Hefei, China  
longfc.ustc@gmail.com, zhaofanqiu@gmail.com, xinmei@ustc.edu.cn  
<sup>2</sup>JD AI Research, Beijing, China

tingyao.ustc@gmail.com, tmei@jd.com

3University of Rochester, Rochester, NY, USA

jluo@cs.rochester.edu

**Abstract.** With the knowledge of action moments (i.e., trimmed video clips that each contains an action instance), humans could routinely localize an action temporally in an untrimmed video. Nevertheless, most practical methods still require all training videos to be labeled with temporal annotations (action category and temporal boundary) and develop the models in a fully-supervised manner, despite expensive labeling efforts and inapplicable to new categories. In this paper, we introduce a new design of transfer learning type to learn action localization for a large set of action categories, but only on action moments from the categories of interest and temporal annotations of untrimmed videos from a small set of action classes. Specifically, we present Action Herald Networks (AherNet) that integrate such design into an one-stage action localization framework. Technically, a weight transfer function is uniquely devised to build the transformation between classification of action moments or foreground video segments and action localization in synthetic contextual moments or untrimmed videos. The context of each moment is learnt through the adversarial mechanism to differentiate the generated features from those of background in untrimmed videos. Extensive experiments are conducted on the learning both across the splits of ActivityNet v1.3 and from THUMOS14 to ActivityNet v1.3. Our AherNet demonstrates the superiority even comparing to most fully-supervised action localization methods. More remarkably, we train AherNet to localize actions from 600 categories on the leverage of action moments in Kinetics-600 and temporal annotations from 200 classes in ActivityNet v1.3.

This work was performed at JD AI Research.

Electronic supplementary material The online version of this chapter ([https://doi.org/10.1007/978-3-030-58580-8\\_9](https://doi.org/10.1007/978-3-030-58580-8_9)) contains supplementary material, which is available to authorized users.

©/circlecopyrt Springer Nature Switzerland AG 2020

A. Vedaldi et al. (Eds.): ECCV 2020, LNCS 12348, pp. 137–154, 2020. <https://doi.org/10.1007/978-3-030-58580-8>

\*\*\*\*\*

ForkGAN: Seeing into the Rainy Night

Ziqiang Zheng<sup>1</sup>, Yan Gao<sup>2</sup>(B), Xinran Han<sup>3</sup>, and Jianbo Shi<sup>3</sup>

<sup>1</sup>UISEE Technology (Beijing) Co., Ltd., Beijing, China

zhengziqiang1@gmail.com

<sup>2</sup>Kyoto University, Kyoto, Japan

wu.yang.8c@kyoto-u.ac.jp

<sup>3</sup>University of Pennsylvania, Philadelphia, USA

{hxinran, jshi}@seas.upenn.edu

**Abstract.** We present a ForkGAN for task-agnostic image translation that can boost multiple vision tasks in adverse weather conditions. Three tasks of image localization/retrieval, semantic image segmentation, and object detection are evaluated. The key challenge is achieving high-quality image translation without any explicit supervision, or task awareness. Our innovation is a fork-shape generator with one encoder and two decoders that disentangles the domain-specific and domain-invariant information. We force the cyclic translation between the weather conditions to go through a common encoding space, and make sure the encoding features reveal no information about the domains. Experimental results show our algorithm produces state-of-the-art image synthesis results and boost three vision tasks' performances in adverse weathers.

**Keywords:** Light illumination

·Image-to-image translation ·Image

synthesis ·Generative adversarial networks

\*\*\*\*\*

TCGM: An Information-Theoretic  
Framework for Semi-supervised  
Multi-modality Learning

Xinwei Sun<sup>1</sup>, Yilun Xu<sup>2</sup>, Peng Cao<sup>2</sup>, Yuqing Kong<sup>2(B)</sup>, Lingjing Hu<sup>3</sup>,  
Shanghang Zhang<sup>4(B)</sup>, and Yizhou Wang<sup>2,5</sup>

<sup>1</sup>Microsoft Research-Asia, Beijing, China

xinsun@microsoft.com

<sup>2</sup>Center on Frontiers of Computing Studies, Advanced Institute of Information  
Technology, Department of Computer Science, Peking University, Beijing, China  
{xuyilun, caopeng2016, yuqing.kong, Yizhou.Wang}@pku.edu.cn

<sup>3</sup>Yanjiang Medical College, Capital Medical University, Beijing, China

hulj@ccmu.edu.cn

<sup>4</sup>UC Berkeley, Berkeley, USA

shz@eecs.berkeley.edu

<sup>5</sup>Deepwise AI Lab, Beijing, China

**Abstract.** Fusing data from multiple modalities provides more information to train machine learning systems. However, it is prohibitively expensive and time-consuming to label each modality with a large amount of data, which leads to a crucial problem of semi-supervised multi-modal learning. Existing methods suffer from either ineffective fusion across modalities or lack of theoretical guarantees under proper assumptions. In this paper, we propose a novel information-theoretic approach - namely, Total Correlation Gain Maximization (TCGM) - for semi-supervised multi-modal learning, which is endowed with promising properties: (i) it can utilize effectively the information across different modalities of unlabeled data points to facilitate training classifiers of each modality (ii) it has theoretical guarantee to identify Bayesian classifiers, i.e., the ground truth posteriors of all modalities. Specifically, by maximizing TC-induced loss (namely TC gain) over classifiers of all modalities, these classifiers can cooperatively discover the equivalent class of ground-truth classifiers; and identify the unique ones by leveraging limited percentage of labeled data. We apply our method to various tasks and achieve state-of-the-art results, including the news classification (Newsgroup dataset), emotion recognition (IEMOCAP and MOSI datasets), and disease prediction (Alzheimer's Disease Neuroimaging Initiative dataset).

X. Sun and Y. Xu—Equal Contribution.

Electronic supplementary material The online version of this chapter ([https://doi.org/10.1007/978-3-030-58580-8\\_11](https://doi.org/10.1007/978-3-030-58580-8_11)) contains supplementary material, which is available to authorized users.

© Springer Nature Switzerland AG 2020

A. Vedaldi et al. (Eds.): ECCV 2020, LNCS 12348, pp. 171–188, 2020.

[https://doi.org/10.1007/978-3-030-58580-8\\_11](https://doi.org/10.1007/978-3-030-58580-8_11)

\*\*\*\*\*

ExchNet: A Unified Hashing Network  
for Large-Scale Fine-Grained Image  
Retrieval

Quan Cui<sup>1,3</sup>, Qing-Yuan Jiang<sup>2</sup>, Xiu-Shen Wei<sup>3(B)</sup>, Wu-Jun Li<sup>2</sup>,  
and Osamu Yoshie<sup>1</sup>

<sup>1</sup>Graduate School of IPS, Waseda University, Fukuoka, Japan

cui-quan@toki.waseda.jp, yoshie@waseda.jp

<sup>2</sup>National Key Laboratory for Novel Software Technology, Department of Computer  
Science and Technology, Nanjing University, Nanjing, China

qyjiang24@gmail.com, liwujun@nju.edu.cn

<sup>3</sup>Megvii Research Nanjing, Megvii Technology, Nanjing, China

weixs.gm@gmail.com

**Abstract.** Retrieving content relevant images from a large-scale fine-grained dataset could suffer from intolerably slow query speed and highly redundant storage cost, due to high-dimensional real-valued embeddings

which aim to distinguish subtle visual differences of fine-grained objects. In this paper, we study the novel fine-grained hashing topic to generate compact binary codes for fine-grained images, leveraging the search and storage efficiency of hash learning to alleviate the aforementioned problems. Specifically, we propose a unified end-to-end trainable network, termed as ExchNet. Based on attention mechanisms and proposed attention constraints, ExchNet can firstly obtain both local and global features to represent object parts and the whole fine-grained objects, respectively. Furthermore, to ensure the discriminative ability and semantic meaning's consistency of these part-level features across images, we design a local feature alignment approach by performing a feature exchanging operation. Later, an alternating learning algorithm is employed to optimize the whole ExchNet and then generate the final binary hash codes. Validated by extensive experiments, our ExchNet consistently outperforms state-of-the-art generic hashing methods on five fine-grained datasets. Moreover, compared with other approximate nearest neighbor methods, ExchNet achieves the best speed-up and storage reduction, revealing its efficiency and practicality.

Keywords: Fine-Grained Image Retrieval

·Learning to hash ·

Feature alignment ·Large-scale image search

Q. Cui, Q.-Y. Jiang—Equal contribution.

Electronic supplementary material The online version of this chapter ([https://doi.org/10.1007/978-3-030-58580-8\\_12](https://doi.org/10.1007/978-3-030-58580-8_12)) contains supplementary material, which is available to authorized users.

© Springer Nature Switzerland AG 2020

A. Vedaldi et al. (Eds.): ECCV 2020, LNCS 12348, pp. 189–205, 2020. <https://doi.org/10.1007/978-3-030-58580-8>

\_1

\*\*\*\*\*

TSIT: A Simple and Versatile Framework  
for Image-to-Image Translation

Liming Jiang<sup>1</sup>, Changxu Zhang<sup>2</sup>, Mingyang Huang<sup>3</sup>, Chunxiao Liu<sup>3</sup>,  
Jianping Shi<sup>3</sup>, and Chen Change Loy<sup>1(B)</sup>

<sup>1</sup>Nanyang Technological University, Singapore, Singapore

{liming002, ccloy}@ntu.edu.sg

<sup>2</sup>University of California, Berkeley, CA, USA

zhangcx@berkeley.edu

<sup>3</sup>SenseTime Research, Beijing, China

{huangmingyang, liuchunxiao, shijianping}@sensetime.com

**Abstract.** We introduce a simple and versatile framework for image-to-image translation. We unearth the importance of normalization layers, and provide a carefully designed two-stream generative model with newly proposed feature transformations in a coarse-to-fine fashion. This allows multi-scale semantic structure information and style representation to be effectively captured and fused by the network, permitting our method to scale to various tasks in both unsupervised and supervised settings. No additional constraints (e.g., cycle consistency) are needed, contributing to a very clean and simple method. Multi-modal image synthesis with arbitrary style control is made possible. A systematic study compares the proposed method with several state-of-the-art task-specific baselines, verifying its effectiveness in both perceptual quality and quantitative evaluation. GitHub: <https://github.com/EndlessSora/TSIT>.

1

\*\*\*\*\*

ProxyBNN: Learning Binarized Neural  
Networks via Proxy Matrices

Xiangyu He<sup>1,2</sup>, Zitao Mo<sup>1</sup>, Ke Cheng<sup>1,2</sup>, Weixiang Xu<sup>1,2</sup>, Qinghao Hu<sup>1</sup>,  
Peisong Wang<sup>1</sup>, Qingshan Liu<sup>4</sup>, and Jian Cheng<sup>1,2,3(B)</sup>

<sup>1</sup>NLPR, Institute of Automation, Chinese Academy of Sciences, Beijing, China



{xiangyu.he, qinghao.hu, peisong.wang, jcheng}@nlpr.ia.ac.cn

2School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China

3Center for Excellence in Brain Science and Intelligence Technology, Beijing, China

4Nanjing University of Information Science and Technology, Nanjing, China

Abstract. Training Binarized Neural Networks (BNNs) is challenging due to the discreteness. In order to efficiently optimize BNNs through backward propagations, real-valued auxiliary variables are commonly used to accumulate gradient updates. Those auxiliary variables are then directly quantized to binary weights in the forward pass, which brings about large quantization errors. In this paper, by introducing an appropriate proxy matrix, we reduce the weights quantization error while circumventing explicit binary regularizations on the full-precision auxiliary variables. Specifically, we regard pre-binarization weights as a linear combination of the basis vectors. The matrix composed of basis vectors is referred to as the proxy matrix, and auxiliary variables serve as the coefficients of this linear combination. We are the first to empirically identify and study the effectiveness of learning both basis and coefficients to construct the pre-binarization weights. This new proxy learning contributes to new leading performances on benchmark datasets.

Keywords: Binarized Neural Networks

•Proxy matrix

1

\*\*\*\*\*

HMOR: Hierarchical Multi-person

Ordinal Relations for Monocular

Multi-person 3D Pose Estimation

Can Wang<sup>2</sup>, Jiefeng Li<sup>1</sup>, Wentao Liu<sup>2</sup>, Chen Qian<sup>2</sup>, and Cewu Lu<sup>1(B)</sup>

<sup>1</sup>Shanghai Jiao Tong University, Shanghai, China

{ljflikit, lucewu}@sjtu.edu.cn

<sup>2</sup>SenseTime Research, Beijing, China

{wangcan, liuwentao, qianchen}@sensetime.com

Abstract. Remarkable progress has been made in 3D human pose estimation from a monocular RGB camera. However, only a few studies explored 3D multi-person cases. In this paper, we attempt to address the lack of a global perspective of the top-down approaches by introducing a novel form of supervision - Hierarchical Multi-person Ordinal Relations (HMOR). The HMOR encodes interaction information as the ordinal relations of depths and angles hierarchically, which captures the body-part and joint level semantic and maintains global consistency at the same time. In our approach, an integrated top-down model is designed to leverage these ordinal relations in the learning process. The integrated model estimates human bounding boxes, human depths, and root-relative 3D poses simultaneously, with a coarse-to-fine architecture to improve the accuracy of depth estimation. The proposed method significantly outperforms state-of-the-art methods on publicly available multi-person 3D pose datasets. In addition to superior performance, our method costs lower computation complexity and fewer model parameters.

Keywords: 3D human pose

•Ordinal relations •Integrated model

1

\*\*\*\*\*

Mask2CAD: 3D Shape Prediction by

Learning to Segment and Retrieve

Weicheng Kuol, <sup>2(B)</sup>, Anelia Angelova<sup>1,2</sup>, Tsung-Yi Lin<sup>1,2</sup>, and Angela Dai<sup>3</sup>

<sup>1</sup>Google AI, Mountain View, USA

weicheng@google.com, anelia@google.com, tsungyi@google.com

<sup>2</sup>Robotics at Google, Munich, Germany

<sup>3</sup>Technical University of Munich, Munich, Germany

angela.dai@tum.de

**Abstract.** Object recognition has seen significant progress in the image domain, with focus primarily on 2D perception. We propose to leverage existing large-scale datasets of 3D models to understand the underlying 3D structure of objects seen in an image by constructing a CAD-based representation of the objects and their poses. We present Mask2CAD, which jointly detects objects in real-world images and for each detected object, optimizes for the most similar CAD model and its pose. We construct a joint embedding space between the detected regions of an image corresponding to an object and 3D CAD models, enabling retrieval of CAD models for an input RGB image. This produces a clean, lightweight representation of the objects in an image; this CAD-based representation ensures a valid, efficient shape representation for applications such as content creation or interactive scenarios, and makes a step towards understanding the transformation of real-world imagery to a synthetic domain. Experiments on real-world images from Pix3D demonstrate the advantage of our approach in comparison to state of the art. To facilitate future research, we additionally propose a new image-to-3D baseline on ScanNet which features larger shape diversity, real-world occlusions, and challenging image views.

1

\*\*\*\*\*

**A Unified Framework of Surrogate Loss  
by Refactoring and Interpolation**

Lanlan Li<sup>1,2</sup>, Mingzhe Wang<sup>2</sup>, and Jia Deng<sup>2(B)</sup>

<sup>1</sup>University of Michigan, Ann Arbor, MI 48105, USA

llanlan@umich.edu

<sup>2</sup>Princeton University, Princeton, NJ 08544, USA

{mingzhew, jiadeng}@cs.princeton.edu

**Abstract.** We introduce UniLoss, a unified framework to generate surrogate losses for training deep networks with gradient descent, reducing the amount of manual design of task-specific surrogate losses. Our key observation is that in many cases, evaluating a model with a performance metric on a batch of examples can be refactored into four steps: from input to real-valued scores, from scores to comparisons of pairs of scores, from comparisons to binary variables, and from binary variables to the final performance metric. Using this refactoring we generate differentiable approximations for each non-differentiable step through interpolation. Using UniLoss, we can optimize for different tasks and metrics using one unified framework, achieving comparable performance compared with task-specific losses. We validate the effectiveness of UniLoss on three tasks and four datasets. Code is available at <https://github.com/princeton-vl/uniloss>.

**Keywords:** Loss design

·Image classification ·Pose estimation

1

\*\*\*\*\*

**Deep Reflectance Volumes: Relightable  
Reconstructions from Multi-view**

Photometric Images

Sai Bilal<sup>(B)</sup>, Zexiang Xu<sup>1,2</sup>, Kalyan Sunkavalli<sup>2</sup>, Miloš Škafraňský<sup>2</sup>,

Yannick Hold-Georroy<sup>2</sup>, David Kriegman<sup>1</sup>, and Ravi Ramamoorthi<sup>1</sup>

<sup>1</sup>University of California, San Diego, USA

bisai@cs.ucsd.edu

<sup>2</sup>Adobe Research, San Jose, USA

**Abstract.** We present a deep learning approach to reconstruct scene appearance from unstructured images captured under collocated pointlighting. At the heart of Deep Reflectance Volumes is a novel volumetric scene representation consisting of opacity, surface normal and reflectance voxel grids. We present a novel physically-based differentiable volumetric ray marching

g framework to render these scene volumes under arbitrary viewpoint and lighting. This allows us to optimize the scene volumes to minimize the error between their rendered images and the captured images. Our method is able to reconstruct real scenes with challenging non-Lambertian reflectance and complex geometry with occlusions and shadowing. Moreover, it accurately generalizes to novel viewpoints and lighting, including non-collocated lighting, rendering photorealistic images that are significantly better than state-of-the-art mesh-based methods. We also show that our learned reflectance volumes are editable, allowing for modifying the materials of the captured scenes.

Keywords: View synthesis

·Relighting ·Appearance acquisition ·

Neural rendering

1

\*\*\*\*\*

Memory-Augmented Dense Predictive

Coding for Video Representation

Learning

Tengda Han(B), Weidi Xie , and Andrew Zisserman

Visual Geometry Group, Department of Engineering Science,

University of Oxford, Oxford, UK

{htd,weidi,az }@robots.ox.ac.uk

Abstract. The objective of this paper is self-supervised learning from video, in particular for representations for action recognition. We make the following contributions: (i) We propose a new architecture and learning framework Memory-augmented Dense Predictive Coding (MemDPC ) for the task. It is trained with a predictive attention mechanism over the set of compressed memories , such that any future states can always be constructed by a convex combination of the condensed representations, allowing to make multiple hypotheses efficiently. (ii) We investigate visual-only self-supervised video representation learning from RGB frames, or from unsupervised optical flow, or both. (iii) We thoroughly evaluate the quality of the learnt representation on four different downstream tasks: action recognition, video retrieval, learning with scarce annotations, and unintentional action classification. In all cases, we demonstrate state-of-the-art or comparable performance over other approaches with orders of magnitude fewer training data.

1

\*\*\*\*\*

PointMixup: Augmentation

for Point Clouds

Yunlu Chen<sup>1</sup>(B), Vincent Tao Hu<sup>1</sup>(B), Efstratios Gavves<sup>1</sup>, Thomas Mensink<sup>1,2</sup>,

Pascal Mettes<sup>1</sup>, Pengwan Yang<sup>1,3</sup>, and Cees G. M. Snoek<sup>1</sup>

<sup>1</sup>University of Amsterdam, Amsterdam, The Netherlands

{y.chen<sup>3</sup>,t.hu }@uva.nl

<sup>2</sup>Google Research, Amsterdam, The Netherlands

<sup>3</sup>Peking University, Beijing, China

Abstract. This paper introduces data augmentation for point clouds by interpolation between examples. Data augmentation by interpolation has shown to be a simple and effective approach in the image domain. Such a mixup is however not directly transferable to point clouds, as we do not have a one-to-one correspondence between the points of two different objects. In this paper, we define data augmentation between point clouds as a shortest path linear interpolation. To that end, we introduce PointMixup, an interpolation method that generates new examples through an optimal assignment of the path function between two point clouds. We prove that our PointMixup finds the shortest path between two point clouds and that the interpolation is assignment invariant and linear. With the definition of interpolation, PointMixup allows to introduce strong interpolation-based regularizers such as mixup and manifold

mixup to the point cloud domain. Experimentally, we show the potential of PointMixup for point cloud classification, especially when examples are scarce, as well as increased robustness to noise and geometric transformations to points. The code for PointMixup and the experimental details are publicly available (Code is available at: <https://github.com/yunlu-chen/PointMixup/> ).

Keywords: Interpolation

•Point cloud classification •Data augmentation

1

\*\*\*\*\*

Identity-Guided Human Semantic Parsing  
for Person Re-identification

Kuan Zhul,2(B), Haiyun Guo1,Zhiwei Liul,2, Ming Tang1,3,  
and Jinqiao Wang1,2

1National Laboratory of Pattern Recognition, Institute of Automation,  
Chinese Academy of Sciences, Beijing, China

{kuan.zhu,haiyun.guo,zhiwei.liu,tangm,jqwang}@nlpr.ia.ac.cn

2School of Artificial Intelligence, University of Chinese Academy of Sciences,  
Beijing, China

3Shenzhen Inovance Limited, Shenzhen, China

Abstract. Existing alignment-based methods have to employ the pre-trained human parsing models to achieve the pixel-level alignment, and cannot identify the personal belongings (e.g., backpacks and reticule) which are crucial to person re-ID. In this paper, we propose the identity-guided human semantic parsing approach (ISP) to locate both the human body parts and personal belongings at pixel-level for aligned person re-ID only with person identity labels. We design the cascaded clustering on feature maps to generate the pseudo-labels of human parts. Specifically, for the pixels of all images of a person, we first group them to foreground or background and then group the foreground pixels to human parts. The cluster assignments are subsequently used as pseudo-labels of human parts to supervise the part estimation and ISP iteratively learns the feature maps and groups them. Finally, local features of both human body parts and personal belongings are obtained according to the self-learned part estimation, and only features of visible parts are utilized for the retrieval. Extensive experiments on three widely used datasets validate the superiority of ISP over lots of state-of-the-art methods. Our code is available at <https://github.com/CASIA-IVA-Lab/ISP-reID>.

Keywords: Person re-ID

•Weakly-supervised human parsing •  
Aligned representation learning

1

\*\*\*\*\*

Learning Gradient Fields for Shape  
Generation

Ruojin Cai(B), Guandao Yang, Hadar Averbuch-Elor, Zekun Hao,  
Serge Belongie, Noah Snavely, and Bharath Hariharan

Cornell University, Ithaca, USA

rc844@cornell.edu

Abstract. In this work, we propose a novel technique to generate shapes from point cloud data. A point cloud can be viewed as samples from a distribution of 3D points whose density is concentrated near the surface of the shape. Point cloud generation thus amounts to moving randomly sampled points to high-density areas. We generate point clouds by performing stochastic gradient ascent on an unnormalized probability density, thereby moving sampled points toward the high-likelihood regions. Our model directly predicts the gradient of the log density field and can be trained with a simple objective adapted from score-based generative models. We show that our method can reach state-of-the-art performance

for point cloud auto-encoding and generation, while also allowing for extraction of a high-quality implicit surface. Code is available at <https://github.com/RuojinCai/ShapeGF> .

Keywords: 3D generation

•Generative models

1

\*\*\*\*\*

COCO-FUNIT: Few-Shot Unsupervised

Image Translation with a Content

Conditioned Style Encoder

Kuniaki Saito<sup>1,2(B)</sup>, Kate Saenko<sup>1</sup>, and Ming-Yu Liu<sup>2</sup>

<sup>1</sup>Boston University, Boston, USA

{keisaito,saenko}@bu.edu

<sup>2</sup>NVIDIA, Santa Clara, USA

mingyul@nvidia.com

Abstract. Unsupervised image-to-image translation intends to learn a mapping of an image in a given domain to an analogous image in a different domain, without explicit supervision of the mapping. Few-shot unsupervised image-to-image translation further attempts to generalize the model to an unseen domain by leveraging example images of the unseen domain provided at inference time. While remarkably successful, existing few-shot image-to-image translation models find it difficult to preserve the structure of the input image while emulating the appearance of the unseen domain, which we refer to as the content loss problem.

This is particularly severe when the poses of the objects in the input and example images are very different. To address the issue, we propose a new few-shot image translation model, COCO-FUNIT, which computes the style embedding of the example images conditioned on the input image and a new module called the constant style bias. Through extensive experimental validations with comparison to the state-of-the-art, our model shows effectiveness in addressing the content loss problem. Code and pretrained models are available at <https://nvlabs.github.io/COCO-FUNIT/> .

Keywords: Image-to-image translation

•Generative Adversarial

Networks

1

\*\*\*\*\*

Corner Proposal Network for

Anchor-Free, Two-Stage Object Detection

Kaiwen Duan<sup>1</sup>, Lingxi Xie<sup>2</sup>, Honggang Qil, Song Bai<sup>3</sup>, Qingming Huang<sup>1,4(B)</sup>, and Qi Tian<sup>2(B)</sup>

<sup>1</sup>University of Chinese Academy of Sciences, Beijing, China

duankaiwen17@mails.ucas.ac.cn, {hgqi,qmhuang}@ucas.ac.cn

<sup>2</sup>Huawei Inc., Shenzhen, China

198808xc@gmail.com, tian.qil@huawei.com

<sup>3</sup>Huazhong University of Science and Technology, Wuhan, China

songbai.site@gmail.com

<sup>4</sup>Peng Cheng Laboratory, Shenzhen, China

Abstract. The goal of object detection is to determine the class and location of objects in an image. This paper proposes a novel anchor-free, two-stage framework which first extracts a number of object proposals by finding potential corner keypoint combinations and then assigns a class label to each proposal by a standalone classification stage. We demonstrate that these two stages are effective solutions for improving recall and precision, respectively, and they can be integrated into an end-to-end network. Our approach, dubbed Corner Proposal Network (CPN), enjoys the ability to detect objects of various scales and also avoids being confused by a large number of false-positive proposals. On the MS-COCO dataset, CPN achieves an AP of 49.2% which is competitive among state-

of-the-art object detection methods. CPN also sets the scenario of computational efficiency, which achieves an AP of 41.6%/39.7% at 26.2/43.3 FPS, surpassing most competitors with the same inference speed. Code is available at <https://github.com/Duankaiwen/CPNDet>.

Keywords: Object detection

•Anchor-free detector •Two-stage detector •Corner keypoints •Object proposals

1

\*\*\*\*\*

PhraseClick: Toward Achieving Flexible

Interactive Segmentation by Phrase  
and Click

Henghui Ding<sup>1(B)</sup>, Scott Cohen<sup>2</sup>, Brian Price<sup>2</sup>, and Xudong Jiang<sup>1</sup>

<sup>1</sup>Nanyang Technological University, Singapore, Singapore

{ding0093,exdjiang}@ntu.edu.sg

<sup>2</sup>Adobe Research, San Jose, USA

{scohen,bprice}@adobe.com

**Abstract.** Existing interactive object segmentation methods mainly take spatial interactions such as bounding boxes or clicks as input. However, these interactions do not contain information about explicit attributes of the target-of-interest and thus cannot quickly specify what the selected object exactly is, especially when there are diverse scales of candidate objects or the target-of-interest contains multiple objects. Therefore, excessive user interactions are often required to reach desirable results. On the other hand, in existing approaches attribute information of objects is often not well utilized in interactive segmentation. We propose to employ phrase expressions as another interaction input to infer the attributes of target object. In this way, we can 1) leverage spatial clicks to locate the target object and 2) utilize semantic phrases to qualify the attributes of the target object. Specifically, the phrase expressions focus on “what” the target object is and the spatial clicks are in charge of “where” the target object is, which together help to accurately segment the target-of-interest with smaller number of interactions. Moreover, the proposed approach is flexible in terms of interaction modes and can efficiently handle complex scenarios by leveraging the strengths of each type of input. Our multi-modal phrase+click approach achieves new state-of-the-art performance on interactive segmentation. To the best of our knowledge, this is the first work to leverage both clicks and phrases for interactive segmentation.

Keywords: Interactive segmentation

•Click •Phrase •Flexible •

Attribute

1

\*\*\*\*\*

Unified Multisensory Perception:

Weakly-Supervised Audio-Visual

Video Parsing

Yapeng Tian<sup>1(B)</sup>, Dingzeyu Li<sup>2</sup>, and Chenliang Xu<sup>1</sup>

<sup>1</sup>University of Rochester, Rochester, USA

{yapengtian,chenliang.xu}@rochester.edu

<sup>2</sup>Adobe Research, Seattle, USA

dinli@adobe.com

**Abstract.** In this paper, we introduce a new problem, named audio-visual video parsing, which aims to parse a video into temporal event segments and label them as either audible, visible, or both. Such a problem is essential for a complete understanding of the scene depicted inside a video. To facilitate exploration, we collect a Look, Listen, and Parse (LLP) dataset to investigate audio-visual video parsing in a weakly-supervised manner. This task can be naturally formulated as a Multimodal Multiple Instance Learning (MMIL) problem. Concretely, we propose a novel

hybrid attention network to explore unimodal and cross-modal temporal contexts simultaneously. We develop an attentive MMIL pooling method to adaptively explore useful audio and visual content from different temporal extent and modalities. Furthermore, we discover and mitigate modality bias and noisy label issues with an individual-guided learning mechanism and label smoothing technique, respectively. Experimental results show that the challenging audio-visual video parsing can be achieved even with only video-level weak labels. Our proposed framework can effectively leverage unimodal and cross-modal temporal contexts and alleviate modality bias and noisy labels problems.

Keywords: Audio-visual video parsing  
·Weakly-supervised ·LLP  
dataset

1  
\*\*\*\*\*  
Learning Delicate Local Representations  
for Multi-person Pose Estimation

Yuanhao Cai<sup>1,2</sup>, Zhicheng Wang<sup>1(B)</sup>, Zhengxiong Luo<sup>1,3</sup>, Binyi Yin<sup>1,4</sup>,  
Angang Du<sup>1,5</sup>, Haoqian Wang<sup>2</sup>, Xiangyu Zhang<sup>1</sup>, Xinyu Zhou<sup>1</sup>, Erjin Zhou<sup>1</sup>,  
and Jian Sun<sup>1</sup>

<sup>1</sup>Megvii Inc., Beijing, China  
{caiyuanhao,wangzhicheng,zxy,zex,j,zhangxiangyu,sunjian}@megvii.com  
<sup>2</sup>Tsinghua University, Beijing, China  
wanghaoqian@tsinghua.edu.cn

<sup>3</sup>Chinese Academy of Sciences, Beijing, China  
<sup>4</sup>Beihang University, Beijing, China  
<sup>5</sup>Ocean University of China, Qingdao, China

Abstract. In this paper, we propose a novel method called Residual Steps Network (RSN). RSN aggregates features with the same spatial size (Intra-level features) efficiently to obtain delicate local representations, which retain rich low-level spatial information and result in precise keypoint localization. Additionally, we observe the output features contribute differently to final performance. To tackle this problem, we propose an efficient attention mechanism - Pose Refine Machine (PRM) to make a trade-off between local and global representations in output features and further refine the keypoint locations. Our approach won the 1st place of COCO Keypoint Challenge 2019 and achieves state-of-the-art results on both COCO and MPII benchmarks, without using extra training data and pretrained model. Our single model achieves 78.6 on COCO test-dev, 93.0 on MPII test dataset. Ensembled models achieve 79.2 on COCO test-dev, 77.1 on COCO test-challenge dataset. The source code is publicly available for further research at <https://github.com/caiyuanhao1998/RSN/>.

Keywords: Human pose estimation  
·COCO ·MPII·Feature  
aggregation ·Attention mechanism

1  
\*\*\*\*\*  
Learning to Plan with Uncertain  
Topological Maps

Edward Beeching<sup>1(B)</sup>, Jilles Dibangoye<sup>1</sup>, Olivier Simonin<sup>1</sup>,  
and Christian Wolf<sup>2</sup>  
<sup>1</sup>INRIA Chroma team, CITI Lab. INSA Lyon, Villeurbanne, France  
{edward.beeching,jilles.dibangoye,olivier.simonin}@insa-lyon.fr  
<sup>2</sup>Université de Lyon, INSA-Lyon, LIRIS, CNRS, Lyon, France  
christian.wolf@insa-lyon.fr  
<https://team.inria.fr/chroma/en/>

Abstract. We train an agent to navigate in 3D environments using a hierarchical strategy including a high-level graph based planner and a local policy. Our main contribution is a data driven learning based approach for planning under uncertainty in topological maps, requiring an

estimate of shortest paths in valued graphs with a probabilistic structure. Whereas classical symbolic algorithms achieve optimal results on noise-less topologies, or optimal results in a probabilistic sense on graphs with probabilistic structure, we aim to show that machine learning can overcome missing information in the graph by taking into account rich high-dimensional node features, for instance visual information available at each location of the map. Compared to purely learned neural whitebox algorithms, we structure our neural model with an inductive bias for dynamic programming based shortest path algorithms, and we show that a particular parameterization of our neural model corresponds to the Bellman-Ford algorithm. By performing an empirical analysis of our method in simulated photo-realistic 3D environments, we demonstrate that the inclusion of visual features in the learned neural planner outperforms classical symbolic solutions for graph based planning.

Keywords: Visual navigation

·Topological maps ·Graph neural networks

1

\*\*\*\*\*

Neural Design Network: Graphic Layout Generation with Constraints

Hsin-Ying Lee<sup>2</sup>, Lu Jiang<sup>1</sup>, Irfan Essa<sup>1,4</sup>, Phuong B. Le<sup>1</sup>, Haifeng Gong<sup>1</sup>, Ming-Hsuan Yang<sup>1,2,3</sup>, and Weilong Yang<sup>1(B)</sup>

<sup>1</sup>Google Research, Mountain View, USA

<sup>2</sup>University of California, Merced, Merced, USA

<sup>3</sup>Yonsei University, Seoul, South Korea

<sup>4</sup>Georgia Institute of Technology, Atlanta, USA

weilongyang@google.com

**Abstract.** Graphic design is essential for visual communication with layouts being fundamental to composing attractive designs. Layout generation differs from pixel-level image synthesis and is unique in terms of the requirement of mutual relations among the desired components. We propose a method for design layout generation that can satisfy user-specified constraints. The proposed neural design network (NDN) consists of three modules. The first module predicts a graph with complete relations from a graph with user-specified relations. The second module generates a layout from the predicted graph. Finally, the third module fine-tunes the predicted layout. Quantitative and qualitative experiments demonstrate that the generated layouts are visually similar to real design layouts. We also construct real designs based on predicted layouts for a better understanding of the visual quality. Finally, we demonstrate a practical application on layout recommendation.

1

\*\*\*\*\*

Learning Open Set Network with Discriminative Reciprocal Points

Guangyao Chen<sup>1</sup>, Limeng Qiao<sup>1</sup>, Yemin Shi<sup>1</sup>, Peixi Peng<sup>1(B)</sup>, Jialin Li<sup>2,3</sup>, Tiejun Huang<sup>1,3</sup>, Shiliang Pu<sup>4</sup>, and Yonghong Tian<sup>1,3(B)</sup>

<sup>1</sup>Department of Computer Science and Technology, Peking University, Beijing, China

{gy.chen,qiaolm,pxpeng,yhtian}@pku.edu.cn

<sup>2</sup>State Key Laboratory of Virtual Reality Technology and Systems, SCSE, Beihang University, Beijing, China

<sup>3</sup>Peng Cheng Laboratory, Shenzhen, China

<sup>4</sup>Hikvision Research Institute, Hangzhou, China

**Abstract.** Open set recognition is an emerging research area that aims to simultaneously classify samples from predefined classes and identify the rest as 'unknown'. In this process, one of the key challenges is to reduce the risk of generalizing the inherent characteristics of numerous unknown samples learned from a small amount of known data. In this



paper, we propose a new concept, Reciprocal Point, which is the potential representation of the extra-class space corresponding to each known category. The sample can be classified to known or unknown by the otherness with reciprocal points. To tackle the open set problem, we offer a novel open space risk regularization term. Based on the bounded space constructed by reciprocal points, the risk of unknown is reduced through multi-category interaction. The novel learning framework called Reciprocal Point Learning (RPL), which can indirectly introduce the unknown information into the learner with only known classes, so as to learn more compact and discriminative representations. Moreover, we further construct a new large-scale challenging aircraft dataset for open set recognition: Aircraft 300 (Air-300). Extensive experiments on multiple benchmark datasets indicate that our framework is significantly superior to other existing approaches and achieves state-of-the-art performance on standard open set benchmarks.

1

\*\*\*\*\*

Convolutional Occupancy Networks

Songyou Peng<sup>1,2(B)</sup>, Michael Niemeyer<sup>2,3</sup>, Lars Mescheder<sup>2,4</sup>,  
Marc Pollefeys<sup>1,5</sup>, and Andreas Geiger<sup>2,3</sup>

<sup>1</sup>ETH Zurich, Zurich, Switzerland

songyou.peng@inf.ethz.ch

<sup>2</sup>Max Planck Institute for Intelligent Systems, Tübingen, Germany

<sup>3</sup>University of Tübingen, Tübingen, Germany

<sup>4</sup>Amazon, Tübingen, Germany

<sup>5</sup>Microsoft, Zurich, Switzerland

**Abstract.** Recently, implicit neural representations have gained popularity for learning-based 3D reconstruction. While demonstrating promising results, most implicit approaches are limited to comparably simple geometry of single objects and do not scale to more complicated or large-scale scenes. The key limiting factor of implicit methods is their simple fully-connected network architecture which does not allow for integrating local information in the observations or incorporating inductive biases such as translational equivariance. In this paper, we propose Convolutional Occupancy Networks, a more flexible implicit representation for detailed reconstruction of objects and 3D scenes. By combining convolutional encoders with implicit occupancy decoders, our model incorporates inductive biases, enabling structured reasoning in 3D space. We investigate the effectiveness of the proposed representation by reconstructing complex geometry from noisy point clouds and low-resolution voxel representations. We empirically find that our method enables the fine-grained implicit 3D reconstruction of single objects, scales to large indoor scenes, and generalizes well from synthetic to real data.

1

\*\*\*\*\*

Multi-person 3D Pose Estimation

in Crowded Scenes Based

on Multi-view Geometry

He Chen<sup>1(B)</sup>, Pengfei Guo<sup>1</sup>, Pengfei Li<sup>1</sup>, Gim Hee Lee<sup>2</sup>,  
and Gregory Chirikjian<sup>1,2</sup>

<sup>1</sup>The Johns Hopkins University, Baltimore, USA

{hchen136, pguo4, pli32, gchirik1}@jhu.edu

<sup>2</sup>National University of Singapore, Singapore, Singapore

gimhee.lee@comp.nus.edu.sg, mpegre@nus.edu.sg

**Abstract.** Epipolar constraints are at the core of feature matching and depth estimation in current multi-person multi-camera 3D human pose estimation methods. Despite the satisfactory performance of this formulation in sparser crowd scenes, its effectiveness is frequently challenged under denser crowd circumstances mainly due to two sources of ambiguity. The first is the mismatch of human joints resulting from the simple cues

provided by the Euclidean distances between joints and epipolar lines. The second is the lack of robustness from the naive formulation of the problem as a least squares minimization. In this paper, we depart from the multi-person 3D pose estimation formulation, and instead reformulate it as crowd pose estimation. Our method consists of two key components: a graph model for fast cross-view matching, and a maximum a posteriori (MAP) estimator for the reconstruction of the 3D human poses. We demonstrate the effectiveness and superiority of our proposed method on four benchmark datasets. Our code is available at: <https://github.com/HeCraneChen/3D-Crowd-Pose-Estimation-Based-on-MVG>.

Keywords: 3D pose estimation  
·Occlusion ·Correspondence problem

1

\*\*\*\*\*

TIDE: A General Toolbox for Identifying

Object Detection Errors

Daniel Bolya(B), Sean Foley, James Hays, and Judy Ho■man

Georgia Institute of Technology, Atlanta, USA

dbolya@gatech.edu

Abstract. We introduce TIDE, a framework and associated toolbox (<https://dbolya.github.io/tide/>) for analyzing the sources of error in object detection and instance segmentation algorithms. Importantly, our framework is applicable across datasets and can be applied directly to output prediction ■les without required knowledge of the underlying prediction system. Thus, our framework can be used as a drop-in replacement for the standard mAP computation while providing a comprehensive analysis of each model's strengths and weaknesses. We segment errors into six types and, crucially, are the ■rst to introduce a technique for measuring the contribution of each error in a way that isolates its effect on overall performance. We show that such a representation is critical for drawing accurate, comprehensive conclusions through in-depth analysis across 4 datasets and 7 recognition models.

Keywords: Error diagnosis  
·Object detection ·Instance segmentation

1

\*\*\*\*\*

PointContrast: Unsupervised Pre-training

for 3D Point Cloud Understanding

Saining Xiel(B), Jiatao Gul, Demi Guol, Charles R. Qil, Leonidas Guibas2, and Or Litany2

1Facebook AI Research, Menlo Park, USA

xiesaining@gmail.com

2Stanford University, Stanford, USA

Abstract. Arguably one of the top success stories of deep learning is transfer learning. The ■nding that pre-training a network on a rich source set (e.g., ImageNet) can help boost performance once ■ne-tuned on a usually much smaller target set, has been instrumental to many applications in language and vision. Yet, very little is known about its usefulness in 3D point cloud understanding. We see this as an opportunity considering the effort required for annotating data in 3D. In this work, we aim at facilitating research on 3D representation learning. Different from previous works, we focus on high-level scene understanding tasks. To this end, we select a suit of diverse datasets and tasks to measure the effect of unsupervised pre-training on a large source set of 3D scenes. Our ■ndings are extremely encouraging: using a unified triplet of architecture, source dataset, and contrastive loss for pre-training, we achieve improvement over recent best results in segmentation and detection across 6 different benchmarks for indoor and outdoor, real and synthetic datasets - demonstrating that the learned representation can generalize across domains. Furthermo

re, the improvement was similar to supervised pre-training, suggesting that future efforts should favor scaling data collection over more detailed annotation. We hope these findings will encourage more research on unsupervised pretext task design for 3D deep learning.

Keywords: Unsupervised learning

•Point cloud recognition •

Representation learning •3D scene understanding

1

\*\*\*\*\*

DSA: More Efficient Budgeted Pruning  
via Differentiable Sparsity Allocation

Xuefei Ning<sup>1</sup>, Tianchen Zhao<sup>2</sup>, Wen Shu<sup>1</sup>, Peng Lei<sup>2</sup>,  
Yu Wang<sup>1(B)</sup>, and Huazhong Yang<sup>2</sup>

<sup>1</sup>Department of Electronic Engineering, Tsinghua University, Beijing, China  
foxdoraame@gmail.com, yu-wang@tsinghua.edu.cn

<sup>2</sup>Department of Electronic Engineering, Beihang University, Beijing, China  
ztc16@buaa.edu.cn

**Abstract.** Budgeted pruning is the problem of pruning under resource constraints. In budgeted pruning, how to distribute the resources across layers (i.e., sparsity allocation) is the key problem. Traditional methods solve it by discretely searching for the layer-wise pruning ratios, which lacks efficiency. In this paper, we propose Differentiable Sparsity Allocation (DSA), an efficient end-to-end budgeted pruning flow. Utilizing an overall differentiable pruning process, DSA finds the layer-wise pruning ratios with gradient-based optimization. It allocates sparsity in continuous space, which is more efficient than methods based on discrete evaluation and search. Furthermore, DSA could work in a pruning-from-scratch manner, whereas traditional budgeted pruning methods are applied to pre-trained models. Experimental results on CIFAR-10 and ImageNet show that DSA could achieve superior performance than current iterative budgeted pruning methods, and shorten the time cost of the overall pruning process by at least 1.5 × in the meantime.

Keywords: Budgeted pruning

•Structured pruning •Model  
compression

1

\*\*\*\*\*

Circumventing Outliers of AutoAugment  
with Knowledge Distillation

Longhui Wei<sup>1,2(B)</sup>, An Xiao<sup>1</sup>, Lingxi Xie<sup>1</sup>, Xiaopeng Zhang<sup>1</sup>, Xin Chen<sup>1,3</sup>,  
and Qi Tian<sup>1</sup>

<sup>1</sup>Huawei Inc., Shenzhen, China  
weilh2568@gmail.com, {xiaolan1,tian.qi1}@huawei.com, 198808xc@gmail.com,  
zxphistory@gmail.com

<sup>2</sup>University of Science and Technology of China, Hefei, China

<sup>3</sup>Tongji University, Shanghai, China

1410452@tongji.edu.cn

**Abstract.** AutoAugment has been a powerful algorithm that improves the accuracy of many vision tasks, yet it is sensitive to the operator space as well as hyper-parameters, and an improper setting may degenerate network optimization. This paper delves deep into the working mechanism, and reveals that AutoAugment may remove part of discriminative information from the training image and so insisting on the ground-truth label is no longer the best option. To relieve the inaccuracy of supervision, we make use of knowledge distillation that refers to the output of a teacher model to guide network training. Experiments are performed in standard image classification benchmarks, and demonstrate the effectiveness of our approach in suppressing noise of data augmentation and stabilizing training. Upon the cooperation of knowledge distillation and AutoAugment, we claim the new state-of-the-art on ImageNet classi-

■cation with a top-1 accuracy of 85.8%.

Keywords: AutoML

·AutoAugment ·Knowledge distillation

1

\*\*\*\*\*

S2DNet: Learning Image Features for  
Accurate Sparse-to-Dense Matching

Hugo Germain<sup>1(B)</sup>, Guillaume Bourmaud<sup>2</sup>, and Vincent Lepetit<sup>1,2</sup>

<sup>1</sup>LIGM, 'Ecole des Ponts, Univ Gustave Eiffel, CNRS, Marne-la-vall' ee, France

{hugo.germain,vincent.lepetit }@enpc.fr

<sup>2</sup>Laboratoire IMS, Universit' e de Bordeaux, Bordeaux, France

guillaume.bourmaud@u-bordeaux.fr

Abstract. Establishing robust and accurate correspondences is a fundamental backbone to many computer vision algorithms. While recent learning-based feature matching methods have shown promising results in providing robust correspondences under challenging conditions, they are often limited in terms of precision. In this paper, we introduce S2DNet, a novel feature matching pipeline, designed and trained to efficiently establish both robust and accurate correspondences. By leveraging a sparse-to-dense matching paradigm, we cast the correspondence learning problem as a supervised classification task to learn to output highly peaked correspondence maps. We show that S2DNet achieves state-of-the-art results on the HPatches benchmark, as well as on several long-term visual localization datasets.

Keywords: Feature matching

·Classification ·Visual localization

1

\*\*\*\*\*

RTM3D: Real-Time Monocular

3D Detection from Object Keypoints

for Autonomous Driving

Peixuan Li<sup>1,2,3,4,5</sup>, Huaici Zhao<sup>1,2,4,5(B)</sup>, Pengfei Liu<sup>1,2,3,4,5</sup>,

and Feidao Cao<sup>1,2,3,4,5</sup>

<sup>1</sup>Shenyang Institute of Automation, Chinese Academy of Sciences, Shenyang, China

hczhao@sia.cn

<sup>2</sup>Institutes for Robotics and Intelligent Manufacturing,  
Chinese Academy of Sciences, Shenyang, China

<sup>3</sup>University of Chinese Academy of Sciences, Beijing, China

<sup>4</sup>Key Laboratory of Opto-Electronic Information Processing,  
Chinese Academy of Sciences, Shenyang, China

<sup>5</sup>Key Lab of Image Understanding and Computer Vision, Shenyang, Liaoning, China

Abstract. In this work, we propose an efficient and accurate monocular 3D detection framework in single shot. Most successful 3D detectors take the projection constraint from the 3D bounding box to the 2D box as an important component. Four edges of a 2D box provide only four constraints and the performance deteriorates dramatically with the small error of the 2D detector. Different from these approaches, our method predicts the nine perspective keypoints of a 3D bounding box in image space, and then utilize the geometric relationship of 3D and 2D perspectives to recover the dimension, location, and orientation in 3D space. In this method, the properties of the object can be predicted stably even when the estimation of keypoints is very noisy, which enables us to obtain fast detection speed with a small architecture. Training our method only uses the 3D properties of the object without any extra annotations, category-specific 3D shape priors, or depth maps. Our method is the first real-time system (FPS >24) for monocular image 3D detection while achieves state-of-the-art performance on the KITTI benchmark.

Keywords: Real-time monocular 3D detection

·Autonomous

driving ·Keypoint detection

1

\*\*\*\*\*

## Video Object Segmentation with Episodic Graph Memory Networks

Xiankai Lu<sup>1</sup>, Wenguan Wang<sup>2(B)</sup>, Martin Danelljan<sup>2</sup>, Tianfei Zhou<sup>1</sup>,  
Jianbing Shen<sup>1</sup>, and Luc Van Gool<sup>2</sup>

<sup>1</sup>Inception Institute of Artificial Intelligence, Abu Dhabi, UAE  
carrierlxk@gmail.com

<sup>2</sup>ETH Zurich, Zurich, Switzerland  
wenguanwang.ai@gmail.com

<https://github.com/carrierlxk/GraphMemVOS>

**Abstract.** How to make a segmentation model efficiently adapt to a specific video as well as online target appearance variations is a fundamental issue in the field of video object segmentation. In this work, a graph memory network is developed to address the novel idea of "learning to update the segmentation model". Specifically, we exploit an episodic memory network, organized as a fully connected graph, to store frames as nodes and capture cross-frame correlations by edges. Further, learnable controllers are embedded to ease memory reading and writing, as well as maintain a fixed memory scale. The structured, external memory design enables our model to comprehensively mine and quickly store new knowledge, even with limited visual information, and the differentiable memory controllers slowly learn an abstract method for storing useful representations in the memory and how to later use these representations for prediction, via gradient descent. In addition, the proposed graph memory network yields a neat yet principled framework, which can generalize well to both one-shot and zero-shot video object segmentation tasks. Extensive experiments on four challenging benchmark datasets verify that our graph memory network is able to facilitate the adaptation of the segmentation network for case-by-case video object segmentation.

**Keywords:** Video segmentation

• Episodic graph memory • Learn to  
update

1

\*\*\*\*\*

## Rethinking Bottleneck Structure for Efficient Mobile Network Design

Daquan Zhou<sup>1,2,3(B)</sup>, Qibin Hou<sup>1(B)</sup>, Yunpeng Chen<sup>2</sup>, Jiashi Feng<sup>1</sup>,  
and Shuicheng Yan<sup>2</sup>

<sup>1</sup>National University of Singapore, Singapore, Singapore  
zhoudaquan21@gmail.com, andrewhoux@gmail.com, elefjia@nus.edu.sg

<sup>2</sup>Yitu Technology, Singapore, Singapore  
{yunpeng.chen, shuicheng.yan}@yitu-inc.com

<sup>3</sup>Institute of Data Science, NUS, Singapore, Singapore

**Abstract.** The inverted residual block is dominating architecture design for mobile networks recently. It changes the classic residual bottleneck by introducing two design rules: learning inverted residuals and using linear bottlenecks. In this paper, we rethink the necessity of such design changes and find it may bring risks of information loss and gradient confusion. We thus propose to flip the structure and present a novel bottleneck design, called the sandglass block, that performs identity mapping and spatial transformation at higher dimensions and thus alleviates information loss and gradient confusion effectively. Extensive experiments demonstrate that, different from the common belief, such bottleneck structure is more beneficial than the inverted ones for mobile networks. In ImageNet classification, by simply replacing the inverted residual block with our sandglass block without increasing parameters and computation, the classification accuracy can be improved by more than 1.7% over MobileNetV2. On Pascal VOC 2007 test set, we observe that there is

also 0.9% mAP improvement in object detection. We further verify the effectiveness of the sandglass block by adding it into the search space of neural architecture search method DARTS. With 25% parameter reduction, the classification accuracy is improved by 0.13% over previous DARTS models. Code can be found at: [https://github.com/zhoudaquan/rethinking\\_bottleneck\\_design](https://github.com/zhoudaquan/rethinking_bottleneck_design).

Keywords: Sandglass block · Residual block · Efficient architecture design · Image classification

D. Zhou and Q. Hou—Authors contributed equally. D. Zhou—Work done during an internship at Yitu Tech.

Electronic supplementary material The online version of this chapter ([https://doi.org/10.1007/978-3-030-58580-8\\_40](https://doi.org/10.1007/978-3-030-58580-8_40)) contains supplementary material, which is available to authorized users.

© Springer Nature Switzerland AG 2020

A. Vedaldi et al. (Eds.): ECCV 2020, LNCS 12348, pp. 680–697, 2020. [https://doi.org/10.1007/978-3-030-58580-8\\_40](https://doi.org/10.1007/978-3-030-58580-8_40)

\_4

\*\*\*\*\*

Side-Tuning: A Baseline for Network Adaptation via Additive Side Networks

Jeffrey O. Zhang<sup>1(B)</sup>, Alexander Saxl<sup>1</sup>, Amir Zamir<sup>3</sup>, Leonidas Guibas<sup>2</sup>, and Jitendra Malik<sup>1</sup>

<sup>1</sup>UC Berkeley, Berkeley, USA

[jozhang@berkeley.edu](mailto:jozhang@berkeley.edu)

<sup>2</sup>Stanford University, Stanford, USA

<sup>3</sup>Swiss Federal Institute of Technology (EPFL), Lausanne, Switzerland

<http://sidetuning.berkeley.edu>

**Abstract.** When training a neural network for a desired task, one may prefer to adapt a pre-trained network rather than starting from randomly initialized weights. Adaptation can be useful in cases when training data is scarce, when a single learner needs to perform multiple tasks, or when one wishes to encode priors in the network. The most commonly employed approaches for network adaptation are fine-tuning and using the pre-trained network as a fixed feature extractor, among others. In this paper, we propose a straightforward alternative: side-tuning. Side-tuning adapts a pre-trained network by training a lightweight “side” network that is fused with the (unchanged) pre-trained network via summation. This simple method works as well as or better than existing solutions and it resolves some of the basic issues with fine-tuning, fixed features, and other common approaches. In particular, side-tuning is less prone to overfitting, is asymptotically consistent, and does not suffer from catastrophic forgetting in incremental learning. We demonstrate the performance of side-tuning under a diverse set of scenarios, including incremental learning (iCIFAR, iTaskonomy), reinforcement learning, imitation learning (visual navigation in Habitat), NLP question-answering (SQuAD v2), and single-task transfer learning (Taskonomy), with consistently promising results.

**Keywords:** Sidetuning

· Finetuning · Transfer learning ·

Representation learning · Lifelong learning · Incremental learning ·

Continual learning

1

\*\*\*\*\*

Towards Part-Aware Monocular 3D

Human Pose Estimation: An Architecture

Search Approach

Zerui Chen<sup>1,3(B)</sup>, Yan Huang<sup>1</sup>, Hongyuan Yu<sup>1,3</sup>, Bin Xu<sup>3</sup>, Kehua Ni<sup>1</sup>, Yiru Guo<sup>5</sup>, and Liang Wang<sup>1,2,4</sup>

<sup>1</sup>Center for Research on Intelligent Perception and Computing, NLPR, CASIA,

Beijing, China

{zerui.chen,hongyuan.yu,ke.han }@cripac.ia.ac.cn,

{yhuang,wangliang }@nlpr.ia.ac.cn

2Center for Excellence in Brain Science and Intelligence Technology, CAS,

Beijing, China

3School of Artificial Intelligence, University of Chinese Academy of Sciences,

Beijing, China

xuebin2018@ia.ac.cn

4Chinese Academy of Sciences, Artificial Intelligence Research (CAS-AIR),

Beijing, China

5School of Astronautics, Beihang University, Beijing, China

guoyiru@buaa.edu.cn

**Abstract.** Even though most existing monocular 3D pose estimation approaches achieve very competitive results, they ignore the heterogeneity among human body parts by estimating them with the same network architecture. To accurately estimate 3D poses of different body parts, we attempt to build a part-aware 3D pose estimator by searching a set of network architectures. Consequently, our model automatically learns to select a suitable architecture to estimate each body part. Compared to models built on the commonly used ResNet-50 backbone, it reduces 62% parameters and achieves better performance. With roughly the same computational complexity as previous models, our approach achieves state-of-the-art results on both the single-person and multi-person 3D pose estimation benchmarks.

**Keywords:** 3D pose estimation

·Body parts ·Neural architecture search

1

\*\*\*\*\*

**REVISE: A Tool for Measuring and Mitigating Bias in Visual Datasets**

Angelina Wang(B), Arvind Narayanan , and Olga Russakovsky

Princeton University, Princeton, USA

angelina.wang@princeton.edu

**Abstract.** Machine learning models are known to perpetuate and even amplify the biases present in the data. However, these data biases frequently do not become apparent until after the models are deployed. To tackle this issue and to enable the preemptive analysis of large-scale dataset, we present our tool. REVISE (REvealing VISual biasSEs) is a tool that assists in the investigation of a visual dataset, surfacing potential biases currently along three dimensions: (1) object-based, (2) gender-based, and (3) geography-based. Object-based biases relate to size, context, or diversity of object representation. Gender-based metrics aim to reveal the stereotypical portrayal of people of different genders. Geography-based analyses consider the representation of different geo-graphic locations. REVISE sheds light on the dataset along these dimensions; the responsibility then lies with the user to consider the cultural and historical context, and to determine which of the revealed biases may be problematic. The tool then further assists the user by suggesting actionable steps that may be taken to mitigate the revealed biases. Overall, the key aim of our work is to tackle the machine learning bias problem early in the pipeline. REVISE is available at <https://github.com/princetonvisualai/revise-tool>.

**Keywords:** Dataset bias

·Dataset analysis ·Computer vision fairness

1

\*\*\*\*\*

**Contrastive Learning for Weakly Supervised Phrase Grounding**

Tanmay Gupta(B), Arash Vahdat<sup>3</sup>, Gal Chechik<sup>2,3</sup>, Xiaodong Yang<sup>3</sup>,

Jan Kautz<sup>3</sup>, and Derek Hoiem<sup>1</sup>

<sup>1</sup>University of Illinois Urbana-Champaign, Champaign, USA

tgupta6@illinois.edu

<sup>2</sup>Bar Ilan University, Ramat Gan, Israel

<sup>3</sup>NVIDIA, Santa Clara, USA

**Abstract.** Phrase grounding, the problem of associating image regions to caption words, is a crucial component of vision-language tasks. We show that phrase grounding can be learned by optimizing word-region attention to maximize a lower bound on mutual information between images and caption words. Given pairs of images and captions, we maximize compatibility of the attention-weighted regions and the words in the corresponding caption, compared to non-corresponding pairs of images and captions. A key idea is to construct effective negative captions for learning through language model guided word substitutions. Training with our negatives yields a ~10% absolute gain in accuracy over randomly-sampled negatives from the training data. Our weakly-supervised phrase grounding model trained on COCO-Captions shows a healthy gain of 5.7% to achieve 76.7% accuracy on Flickr30K Entities benchmark. Our code and project material will be available at <http://tanmaygupta.info/info-ground>.

**Keywords:** Mutual information

·InfoNCE ·Grounding ·Attention

1

\*\*\*\*\*

Collaborative Learning of Gesture

Recognition and 3D Hand Pose

Estimation with Multi-order Feature

Analysis

Siyuan Yang<sup>1,2</sup>, Jun Liu<sup>3(B)</sup>, Shijian Lu<sup>4</sup>, Meng Hwa Er<sup>2</sup>, and Alex C. Kot<sup>2</sup>

<sup>1</sup>Rapid-Rich Object Search (ROSE) Lab, Interdisciplinary Graduate Programme, Nanyang Technological University, Singapore, Singapore

siyuan005@e.ntu.edu.sg

<sup>2</sup>School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore, Singapore

{emher,eackot}@ntu.edu.sg

<sup>3</sup>Information Systems Technology and Design Pillar, Singapore University of Technology and Design, Singapore, Singapore

junliu@sutd.edu.sg

<sup>4</sup>School of Computer Science and Engineering, Nanyang Technological University, Singapore, Singapore

shijian.Lu@ntu.edu.sg

**Abstract.** Gesture recognition and 3D hand pose estimation are two highly correlated tasks, yet they are often handled separately. In this paper, we present a novel collaborative learning network for joint gesture recognition and 3D hand pose estimation. The proposed network exploits joint-aware features that are crucial for both tasks, with which gesture recognition and 3D hand pose estimation boost each other to learn highly discriminative features. In addition, a novel multi-order multi-stream feature analysis method is introduced which learns posture and multi-order motion information from the intermediate feature maps of videos effectively and efficiently. Due to the exploitation of joint-aware features uncommon, the proposed technique is capable of learning gesture recognition and 3D hand pose estimation even when only gesture or pose labels are available, and this enables weakly supervised network learning with much reduced data labeling efforts. Extensive experiments show that our proposed method achieves superior gesture recognition and 3D hand pose estimation performance as compared with the state-of-the-art.

**Keywords:** Gesture recognition

·3D hand pose estimation ·



Multi-order multi-stream feature analysis ·Slow-fast feature analysis ·  
Multi-scale relation

Electronic supplementary material The online version of this chapter ( [https://doi.org/10.1007/978-3-030-58580-8\\_45](https://doi.org/10.1007/978-3-030-58580-8_45)) contains supplementary material, which is available to authorized users.

©/circlecopyrtSpringer Nature Switzerland AG 2020

A. Vedaldi et al. (Eds.): ECCV 2020, LNCS 12348, pp. 769–786, 2020.<https://doi.org/10.1007/978-3-030-58580-8>

\_4

\*\*\*\*\*