# Geometry in Computer Vision

Natalia Neverova
Research Scientist, Facebook AI

I. Learning correspondences
II. 3D reconstruction
III. Generation
Application: human-centered tasks
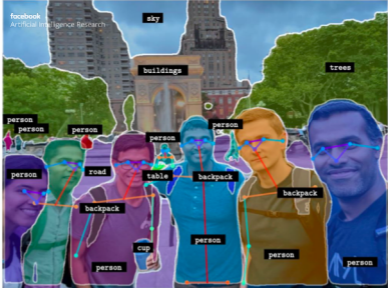
# 3D-fying panoptic perception in the wild

**Object boxes**

**Masks**

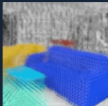[He et al. Mask-RCNN. CVPR. 2017]

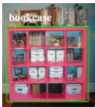# 3D-fying panoptic perception in the wild
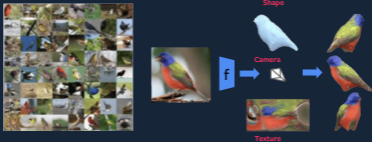
**Object boxes**

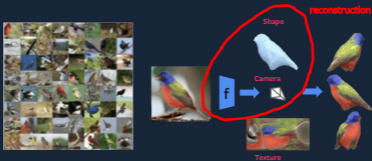**Masks**

**3D geometry**

[Gkioxari et al. Mesh-RCNN. ICCV, 2019]

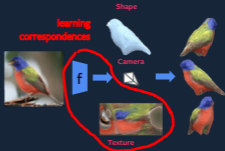# Objects: shape and appearance decomposition



[Kanazawa et al. Learning Category-Specific Mesh Reconstruction from Image Collections. ECCV, 2018]

# Objects: shape and appearance decomposition



[Kanazawa et al. Learning Category-Specific Mesh Reconstruction from Image Collections. ECCV, 2018]

# Objects: shape and appearance decomposition



[Kanazawa et al. Learning Category-Specific Mesh Reconstruction from Image Collections. ECCV, 2018]

# I. Learning correspondences

# Learning correspondences: image to image



image A                    image B

Let g be the **correspondence field** between images A and B
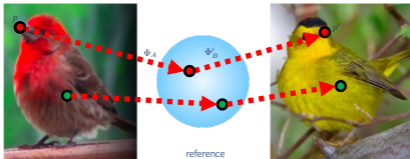
# Learning correspondences



image A · reference · image B

**Factorized** correspondence field: $g = \Phi_B^{-1} \circ \Phi_A$

3D geometry is **irrelevant**, we only need **an index set** over the object surface

# Learning correspondences: image to 3D model



image A

3D model

# Learning correspondences

**Supervised**
- model-driven;
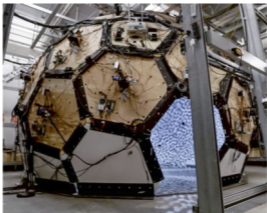- data-driven.

**Un/self-supervised**
- equivariance;
- cycle consistency.

# Model-driven: synthetic data



[Zhou et al. Learning Dense Correspondence via 3D-guided Cycle Consistency. ECCV, 2016]
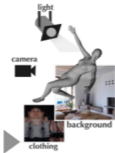
# Synthetic data: articulated objects?



[Joo et al. Panoptic studio: A massively multiview system for social interaction capture. PAMI, 2016]
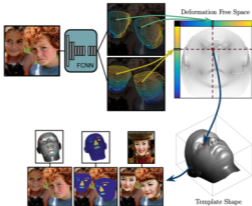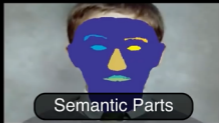
SMPL Model

# Synthetic data: articulated objects?
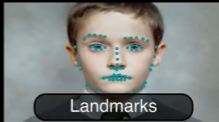


[Varol et al. "Learning from Synthetic Humans". CVPR, 2017]

**Very different image statistics!**

# Model-driven: sparse annotations + fitting



[Güler et al. "DenseReg: Fully Convolutional Dense Shape Regression In-the-Wild". CVPR, 2017]

Input Image

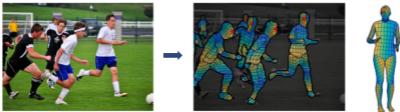Landmarks

Dense Coordinates

Semantic Parts

# Full body articulation?



[Lassner et al. "Unite the People: Closing the Loop Between 3D and 2D Human Representations". CVPR, 2017]

**Poor approximation of real data!**
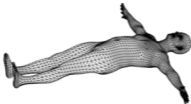
# Data-driven approach: DensePose



[Güler et al. "DensePose: learning dense correspondences in the wild". CVPR, 2018]

Eliminates dependency on a **specific 3D model** and its **expressivity** (as long as **semantics is preserved**)
**No domain gap**, annotations are **easier to obtain**
Human annotation **errors can be significant** due too ambiguities

# Dense correspondence task

# COCO-DensePose dataset

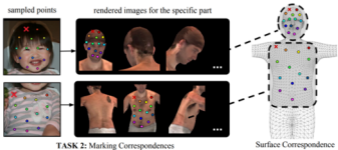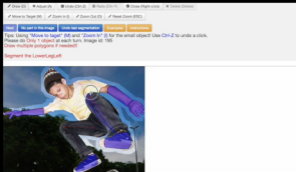Annotation task 1: body part segmentation

# COCO-DensePose dataset

Annotation task 2: marking sparse correspondences



sampled points    rendered images for the specific part

**TASK 2: Marking Correspondences**

Surface Correspondence

# COCO DensePose:
## Collecting Data



## Task - 1
### Part Segmentation

# COCO-DensePose dataset

50 annotated instances, 5 million correspondences (~100 points/image)

# Evaluation metric: geodesic distance



For instance based frameworks:

$$\text{GPS}_j = \frac{1}{|P_j|} \sum_{p \in P_j} \exp\left( \frac{-g(i_p, \hat{i}_p)^2}{2\kappa^2} \right)$$

**geodesic point similarity (GPS)**

# Architecture: DensePose-RCNN



[He, Gkioxari, Dollar, Girshick. Mask-RCNN. ICCV, 2017]

github.com/facebookresearch/DensePose

facebook
Artificial Intelligence Research

Textures taken from SURREAL dataset.
Varol, Gül, et al. "Learning from synthetic humans." CVPR 2017.

# Real-time demos on desktop & mobile

# DensePoseTrack dataset



[Andriluka et al. PoseTrack: A Benchmark for Human Pose Estimation and Tracking, CVPR, 2018]

# DensePoseTrack dataset



[Neverova et al. Slim DensePose: Thrifty Learning with Motion Cues. CVPR, 2019]

Labeled images: 1680 / 782 (training / validation)
Instances: 8274 / 4753
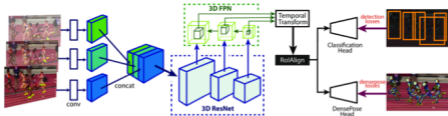Correspondences: 800142 / 459348
Every 2nd frame for 4 frames, every 8th frame otherwise
Ignored: instances with <6 keypoints, severe motion blur
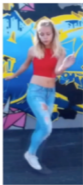
# Flow-guided 3D DensePose-RCNN



[Khalidov et al., 2019]

baseline

ours

input

[Neverova et al. Slim DensePose: Thrifty Learning with Motion Cues. CVPR, 2019]

# Learning correspondences

**Supervised**
- model-driven;
- data-driven.

**Un/self-supervised**
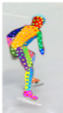- equivariance;
- cycle consistency.

Learning with less supervision?

# Cost efficient annotation process



full annotations      full dense annotations      sparse annotations      keypoints

[Neverova et al. Slim DensePose: Thrifty Learning with Motion Cues. CVPR, 2019]

# Cost efficient annotation process



Points on the SMPL model

points with UV, %

0   1   5   10   20   50   100

# Correspondences by self-supervision



GT propagation & equivariance

[Neverova et al. Slim DensePose:
Thrifty Learning with Motion Cues. CVPR, 2019]



cycle-consistency

[Kulkarni et al. Canonical Surface Mapping
via Geometric Cycle Consistency. ICCV, 2019]

# GT propagation vs equivariance



Transfer a given label to a new frame

GT propagation

Constrain unknown labels to be consistent

equivariance

# Synthetic equivariance: thin-plane splines (TPS)

# Synthetic equivariance: thin-plane splines (TPS)



The **known mapping between points in a pair of original-deformed frames** is used both for **data augmentation (sparsely)** and enforcing **equivariance (densely).**

# Flow-guided temporal equivariance



frames t, t+1...3
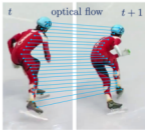
flow

propagated annotations

**Optical flow** is used both for **data augmentation (sparse points)** and enforcing **inter-frame temporal equivariance (densely)**

# Flow-guided temporal equivariance

Real transform >> synthetic

GT propagation >> equivariance

Combination > individual

Cycle-consistency

# II. Reconstruction

# 3D reconstruction

## 3D-supervised

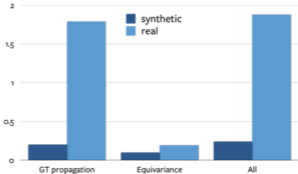- synthetic data;
- multi-view data / video;
- manual annotations (?).

## 2D-supervised

# Model-based 3D reconstruction

Predicting **parameters of the SMPL model** based on DensePose representation rather than RGB

Synthetic data helps due to smaller domain gap in this space

[DenseRac: Joint 3D Pose and Shape Estimation by Dense Render-and-Compare. ECCV, 2018]

# Synthetic Data from Virtual World

facebook
Reality Labs

**Mixamo**
(www.mixamo.com)

Offering free animated
3D characters

Thousands of
customizable 3D
animations

# Differentiable rendering

facebook
Reality Labs

Learning with less supervision?

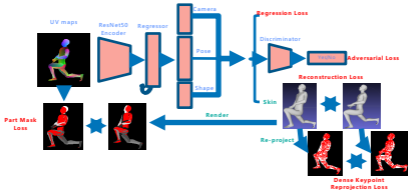# 3D reconstruction with 2D supervision



[Novotny, Ravi et al. C3DPO: Canonical 3D Pose Networks for Non-Rigid Structure From Motion. ICCV, 2019]

# Canonicalization network

# From sparse landmarks to dense annotations



Qualitative results on **synthetic renderings** using the SMPL model

# III. Generation

# Generation

AR

VR

Creativity / Gaming

# Mapping from image to texture space

Textures taken from SURREAL dataset.
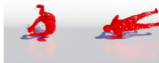Varol, Gül, et al. "Learning from synthetic humans." CVPR 2017.

facebook
Artificial Intelligence Research

Input Image

Target Image

DensePose Texture

Inpainted Texture
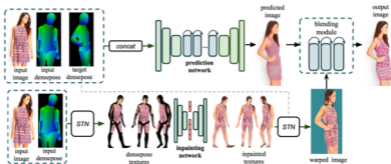
Inpainted Texture Transfer

[Neverova et al. Dense Pose Transfer [CVPR, 2019]]

# Texture inpainting in UV space



The inpainting network learns to reconstruct **full body texture** from **partial observations** by autoencoding in a normalized texture space

# Two stream model



The inpainting network introduces **generalization over the pose** space for free

# DensePose vs sparse keypoints conditioning



Keypoint-based

DensePose-based

Conditioning on DensePose resolves **ambiguity in z-ordering** and encourages **anatomical plausibility**

From fits-them-all to personalized models

input　　poses　　output　　　　input　　poses　　output
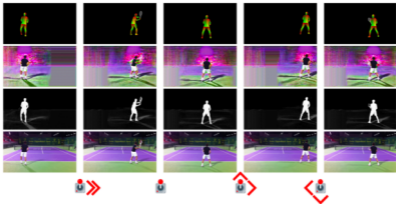
input　　poses　　output　　　　input　　poses　　output

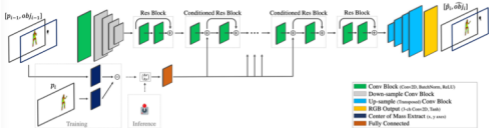[Wang, Liu, Zhu, Liu, Tao, Kautz, Catanzaro. Video-to-Video Synthesis. NeurIPS, 2018]

# Vid2game: creating controllable characters



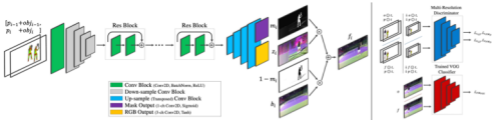[Gafni, Wolf, Taigman. Vid2Game: Controllable Characters Extracted from Real-World Videos. arXiv:1904.08379, 2019]

# Vid2game: creating controllable characters



DensePose representation of the next frame is predicted conditioned
on a current DensePose and an instruction

# Vid2game: creating controllable characters



The new frame is rendered through decomposition: **active character** (predicted) + **background** (copied)
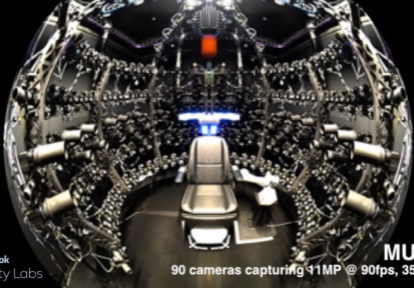
# Evaluation - Walking (Controllable Results)

facebook
Artificial Intelligence Research

... and more personalization

facebook Reality Labs

**MUGSY**
90 cameras capturing 11MP @ 90fps, 350 lights
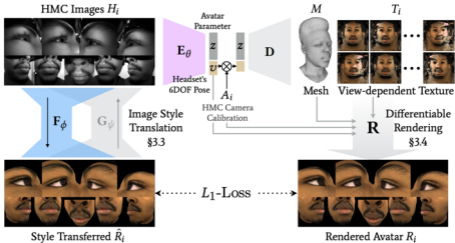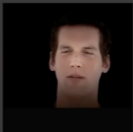
Deep Appearance Variational Autoencoder

Average texture    Encoder    Decoder    View-specific texture

Z
V

Mesh    Mesh

Ground Truth          Rendered Avatar
(novel view)           (novel view)

HMC Images $H_i$

Avatar Parameter

$M$

$T_i$

$\mathbf{E}_\theta$

$z$  $z$

$\mathbf{D}$

Headset's 6DOF Pose

$v$

$A_i$

HMC Camera Calibration

Mesh

View-dependent Texture

$\mathbf{F}_\phi$  $\mathbf{G}_\psi$

Image Style Translation §3.3

$\mathbf{R}$

Differentiable Rendering §3.4

Style Transferred $\hat{R}_i$

$\cdots\cdots$ $L_1$-Loss $\cdots\cdots$

Rendered Avatar $R_i$

# Questions?

**facebook**
Artificial Intelligence Research