

# Deep 3D Geometry Computer Vision - Selected Recent Research

Yang Li, Richard Xu

January 28, 2020

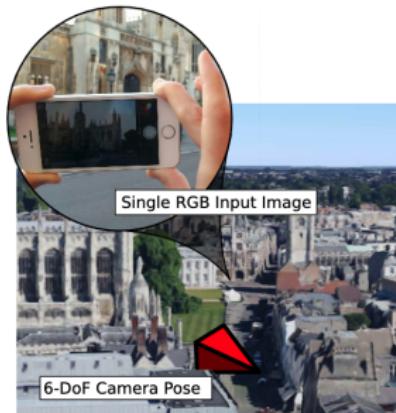
Selected topics on cutting-edge research using **Deep Learning** for solving 3D computer vision problems:

- ▶ Single image → Camera Model estimation
- ▶ Multi-Person 3D pose estimation using multi-view
- ▶ GAN-based 3D pose estimation
- ▶ Deep Learning based Structure from Motion
- ▶ Deep Learning based Depth Estimation

Take home message: Deep Learning is seemingly able to learn how to transform every forms of data to a supervised target

- ▶ but of course it's not true - in 3D, let's get the Geometry and CNN to work together!

- ▶ **PoseNet** is a deep learning framework
- ▶ regress: “single image” → “6-DoF camera pose” using a convolutional neural network



---

<sup>1</sup>Alex Kendall, Matthew Grimes, and Roberto Cipolla. PoseNet: A convolutional network for real-time 6-DoF camera relocalization. In: ICCV2015

- ▶ **Input:** image  $I$ :
  - ▶ **Output:** a pose vector  $[t, q]$  note we did not use Euler angle  $r$  here
    - 1. translation  $t$
    - 2. rotation  $q$  (but in quaternion form) 4DoF
  - ▶ in interpolation: Cartesian coordinates are independent of one another, but Euler angles are not, and Gimbal Lock problem
  - ▶ Quaternion solves it

Intuitively, loss function is:

$$\mathcal{L}(\hat{\mathbf{t}}, \hat{\mathbf{q}}) = \|\hat{\mathbf{t}} - \mathbf{t}\|_2 + \beta \left\| \hat{\mathbf{q}} - \frac{\mathbf{q}}{\|\mathbf{q}\|} \right\|_2$$

- ▶  $\beta$  is the weight between  $\mathbf{t}$  and  $\mathbf{q}$

## Drawback of PoseNet

$$\mathcal{L}(\hat{\mathbf{t}}, \hat{\mathbf{q}}) = \|\hat{\mathbf{t}} - \mathbf{t}\|_2 + \beta \left\| \hat{\mathbf{q}} - \frac{\mathbf{q}}{\|\mathbf{q}\|} \right\|_2$$

- ▶ Experiments show jointly regress  $\mathbf{t}$  and  $\mathbf{q}$  is better than training them separately
- ▶ hyperparameter  $\beta$  required significant tuning to get reasonable results. (two different scale)

## Geometric Reprojection Error<sup>2</sup>

- ▶ instead of measure loss in  $[\mathbf{t}, \mathbf{q}]$  space which is unnatural and difficult to determine the relative weights
- ▶ we measure loss in terms of projection error:

$$\mathcal{L}_g(I) = \frac{1}{|\mathcal{G}'|} \sum_{g_i \in \mathcal{G}'} \|\pi(\mathbf{x}, \mathbf{q}, \mathbf{g}_i) - \pi(\hat{\mathbf{x}}, \hat{\mathbf{q}}, \mathbf{g}_i)\|_\gamma$$

- ▶ Function  $\pi$  just a regular  $3D \rightarrow 2D$  projection, for some **test point  $\mathbf{g}$** :

$$\begin{bmatrix} u' \\ v' \\ w' \end{bmatrix} = \mathbf{K}(\mathbf{R}\mathbf{g} + \mathbf{t})$$

$$\pi(\mathbf{t}, \mathbf{q}, \mathbf{g}) = \begin{bmatrix} u'/w' \\ v'/w' \end{bmatrix}$$

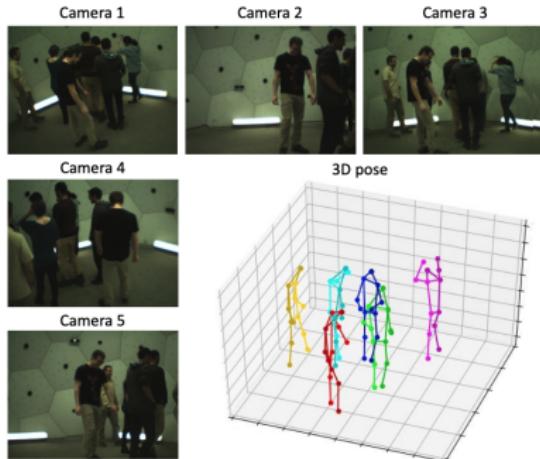
- ▶ instead of regress image  $\rightarrow$  parameter directly
- ▶ now it regresses **image  $\rightarrow$  image location (as the result of the parameter)**

---

<sup>2</sup>Kendall, Alex and Cipolla, and Roberto. Geometric loss functions for camera pose regression with deep learning. In: CVPR2017

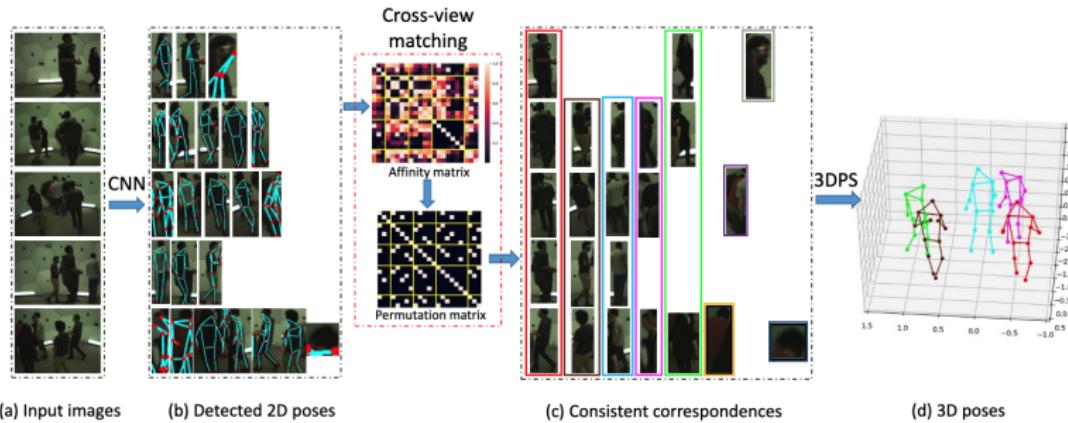
# Multi-person 3D pose estimation from multiple views

- ▶ Reconstruct 3D human pose for multiple persons in a few calibrated camera views.
- ▶ Challenge: Find the cross-view correspondences among noisy and incomplete 2D pose predictions.



# Overall Architecture

Overview of the approach.<sup>3</sup>



<sup>3</sup>Dong J, Jiang W, Huang Q, et al. Fast and robust multi-person 3d pose estimation from multiple views. In: CVPR2019

# Affinity Matrix

Use Cascaded Pyramid Network (CPN) trained on MSCOCO dataset for 2D pose detection.

- ▶  $V$  cameras view
- ▶ affinity matrix  $\mathbf{A}_{ij} \in \mathbb{R}^{p_i \times p_j}$ : whose elements represent affinity scores between two sets of bounding boxes in view  $i$  and view  $j$  of all its bounding boxes
- ▶ how to calculate affinity matrix between two bounding boxes:
  - ▶ Appearance similarity: pre-trained person re-ID network.
  - ▶ Geometric compatibility: satisfy the epipolar constraint.

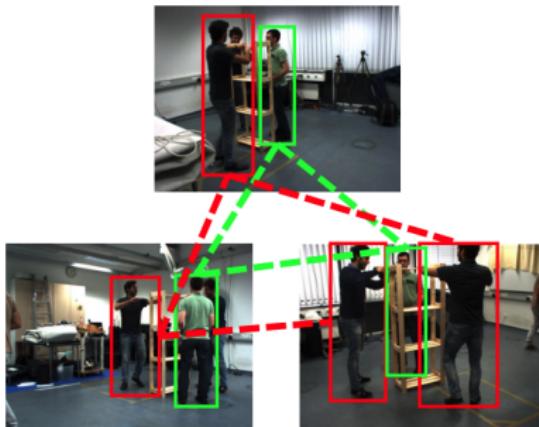
$$D_g(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{2N} \sum_{n=1}^N d_g\left(\mathbf{x}_i^n, \mathcal{L}_{ij}\left(\mathbf{x}_j^n\right)\right) + d_g\left(\mathbf{x}_j^n, \mathcal{L}_{ji}\left(\mathbf{x}_i^n\right)\right)$$

- ▶  $\mathcal{L}_{ij}\left(\mathbf{x}_j^n\right)$  returns the epipolar line associated with  $x_j^n$  from the other view using Fundamental matrix:  $\mathcal{L}_{ij}\left(\mathbf{x}_j^n\right) \equiv \mathbf{F}_{ij}\mathbf{x}_j^n$
- ▶  $d_g(\cdot, l)$  is the point-to-line distance for  $l$ .

- ▶  $P_{ij} \in \{0, 1\}^{p_i \times p_j}$  is the variable to be optimized
- ▶  $P_{ij} = 1$  means correspondence, 0 otherwise
- ▶  $P_{ij}$  is a partial permutation matrix

# Correspondence Matrix

- ▶ Under multiple views, just matching separately for each pair of views, and ignores the **cycle-consistency constraint** can lead to inconsistent results



- ▶ **solution:**

$$\text{rank}(\mathbf{P}) \leq s, \mathbf{P} \succeq 0$$

where  $s$  is the underlying number of people in the scene

- ▶ Since  $s$  is unknown in advance, we should minimize  $\text{rank}(\mathbf{P})$

# Overall Objective

- ▶ real objective function is:

$$\begin{aligned}f(\mathbf{P}) &= - \sum_{i=1}^n \sum_{j=1}^n \langle \mathbf{A}_{ij}, \mathbf{P}_{ij} \rangle + \lambda \cdot \text{rank}(\mathbf{P}) \\&= -\langle \mathbf{A}, \mathbf{P} \rangle + \lambda \cdot \text{rank}(\mathbf{P})\end{aligned}$$

- ▶ approximated by:

$$\min_{\mathbf{P}} \left( -\langle \mathbf{A}, \mathbf{P} \rangle + \lambda \underbrace{\|\mathbf{P}\|_*}_{\text{nuclear norm}} \right)$$

- ▶ Schatten p-norms, for matrix of size  $m \times n$ , this kind of norm “operates on singular values”:

$$\|A\|_p = \left( \sum_{i=1}^{\min\{m,n\}} \sigma_i^p(A) \right)^{\frac{1}{p}}$$

- ▶ when  $p = 1$ :

$$\|A\|_1 = \left( \sum_{i=1}^{\min\{m,n\}} \sigma_i(A) \right) = \text{trace} \left( \sqrt{A^* A} \right)$$

note  $\sqrt{A^* A}$  is a matrix!

$$B = \sqrt{A^* A} \implies B B = A^* A$$

- ▶ we give a special name  $\|A\|_* = \|A\|_1$

## lastly: 3D Pose Generation (pictorial structure model)

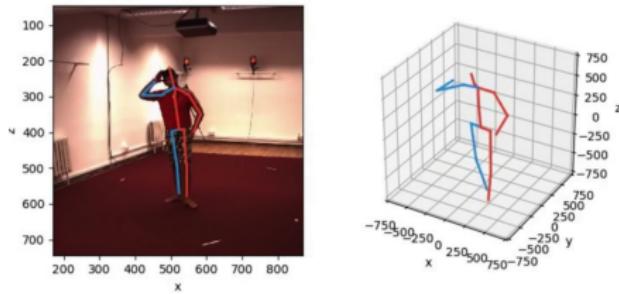
- $\mathbf{T} \equiv (T_1, \dots, T_N)$  joint pose of a person in all  $V$  views:

$$p(\mathbf{T}|I) \propto \underbrace{\prod_{v=1}^V \prod_{i=1}^N p(l_v | \pi_v(t_i))}_{p(l_1, \dots, l_V | \mathbf{T})} \underbrace{\prod_{(i,j) \in \varepsilon} p(t_i, t_j)}_{p(\mathbf{T})}$$

- $\pi_v(t_i)$  2D projection of  $t_i$  on  $v$ -th view
- $p(l_v | \pi_v(t_i))$  is likelihood of heatmap  $l_v$
- $p(t_i, t_j)$  structural dependency between joint  $t_i$  and  $t_j$

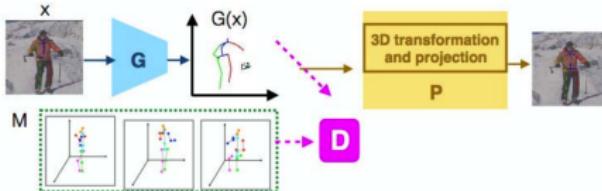
# New topic: Monocular 3D Human Pose Estimation

- ▶ Estimate human body (pose) configuration in 3D space from single image or video
- ▶ focus on the methods combining **geometric priors + neural networks**



# Adversarial Inverse Graphics Network <sup>4</sup>

- ▶ Adversarial Inverse Graphics Network (AIGN) for 3D human pose estimation



- ▶ Basically a conditional GAN + reconstruction error
- ▶ GAN is generally not reversible, i.e.,  $z \rightarrow \mathbf{x}$  is not possible
- ▶ however here  $\mathbf{x}$  is 3D and  $z$  is 2D!
- ▶ Question, can you think of other application may work as well? for example, super-resolution

<sup>4</sup>Tung H Y F, Harley A W, Seto W, et al. Adversarial inverse graphics networks: Learning 2d-to-3d lifting and image-to-image translation from unpaired supervision. In: ICCV2017

- ▶ Decompose 3D human shape into a view-aligned 3D linear basis:

$$S_{3D} = \sum_{j=1}^{|B|} w_j B_j$$

- ▶ model and a rotation matrix

$$\mathbf{x}_{3D} = \mathbf{R} S_{3D} \quad \text{where} \quad \mathbf{R} = \mathbf{R}^x(\alpha) \mathbf{R}^y(\beta) \mathbf{R}_t^z(\gamma)$$

- ▶ Shape basis  $\{B_j\}$  is obtained using PCA on orientation-aligned 3D poses in priorly
- ▶ in here, Generator generates weights  $\mathbf{w} \equiv \{w_1, \dots, w_{60}\}$  to get  $\rightarrow \mathbf{x}_{3D}$ , instead of to generate  $\mathbf{x}_{3D}$  directly

# Loss Function

► “re-projected” 2D points:  $\tilde{\mathbf{x}}_{2D}^{\text{proj}}$  depends on:

1.  $\mathbf{x}_{3D}$  via “weights”  $\mathbf{w} \in \mathbb{R}^{60}$
2. focal length  $f$
3. Euler rotation angles  $\alpha, \beta, \gamma$
4. shift  $[c_x, c_y]$

$$\tilde{\mathbf{x}}_{2D}^{\text{proj}} = \mathbf{K}_f \left( \mathbf{R} \mathbf{x}_{3D} + \begin{bmatrix} c_x \\ c_y \\ 0 \end{bmatrix} \right)$$

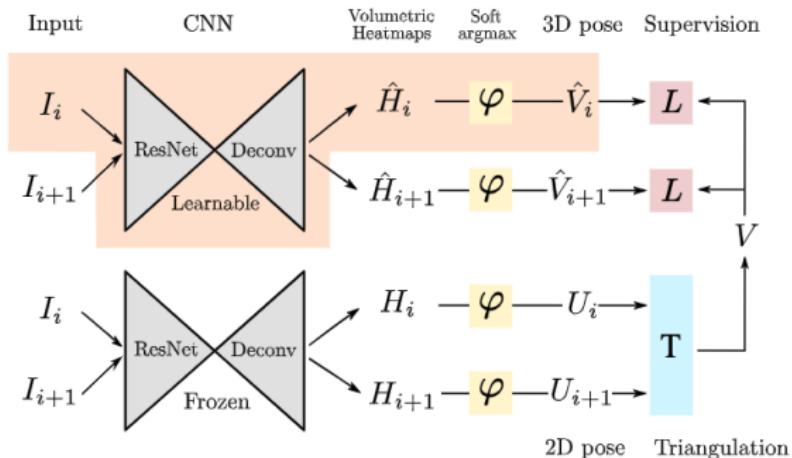
# Loss Function

$$\min_{\mathcal{G}} \max_{\mathcal{D}} \mathbb{E}_{\mathbf{x} \in \Omega_{\mathbf{x}}} \underbrace{\|\hat{\mathbf{x}}_{2D}^{\text{proj}} - \mathbf{x}\|_2}_\text{reconstruction loss} + \\ \beta \sum_{i=1}^K \underbrace{\log D_i(M_i) + \log \left(1 - D_i(\hat{\mathcal{G}}_i(x))\right)}_\text{adversarial loss}$$

- ▶  $\hat{\mathbf{x}}_{2D}^{\text{proj}}$ : “re-projected” 2D
- ▶  $\mathbf{x}$ : “argmax of predicted 2D heatmaps”
- ▶ Discriminator network discriminates between generated keypoints  $\mathbf{x}_{3D}$  and a database of 3D human skeletons

# EpipolarPose Network 5

## ► EpipolarPose network:



<sup>5</sup>Kocabas M, Karagoz S, Akbas E. Self-supervised learning of 3d human pose using multi-view geometry. In: CVPR2019

## training

- ▶ a pair of images ( $I_i, I_{i+1}$ ) from two cameras is fed into the CNN pose estimators
- ▶ 3D pose ( $V$ ) generated by the lower branch using triangulation is used as a training signal for CNN in upper branch

## inference

- ▶ EpipolarPose is a monocular method: it takes a single image ( $I_i$ ) as input and estimates the corresponding 3D pose.

# Obtain 3D Coordinate through Triangulation

- ▶ projection is use traditional stuff:

$$\begin{bmatrix} X_{i,j} \\ Y_{i,j} \\ Z_{i,j} \\ 1 \end{bmatrix} = \mathbf{K}[\mathbf{R} \mid \mathbf{t}] \begin{bmatrix} X_j \\ Y_j \\ Z_j \\ 1 \end{bmatrix}, \mathbf{K} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}, \mathbf{t} = \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix}$$

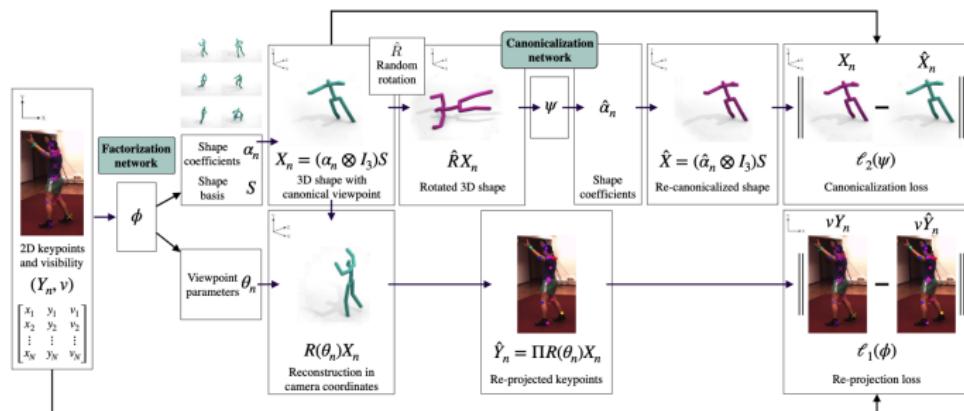
- ▶ When  $[\mathbf{R}|\mathbf{t}]$  are available, using triangulation to obtain 3D pose

## When extrinsic parameters are **unavailable**:

- ▶ Assume first camera is at center of coordinate system: means  $\mathbf{R}$  of first camera is identity matrix
- ▶ For corresponding joints in view  $(i, i + 1)$ :  $U_i$  and  $U_{i+1}$ :
  1. estimate fundamental matrix  $\mathbf{F}$  :  $U_{i,j}\mathbf{F}U_{i+1,j} = 0$  for  $\forall j$  using RANSAC
  2. from  $\mathbf{F}$  (acts like data), calculate Essential matrix  $\mathbf{E} = \mathbf{K}^\top \mathbf{F} \mathbf{K}$
  3. by decomposing  $\mathbf{E}$  (acts like data now) with SVD, obtain  $\mathbf{E} = \mathbf{R}[\mathbf{t}_\times]$

# Canonical 3D Pose Networks

- ▶ Overview architecture of the Canonical 3D Pose Networks (C3DPO)
- ▶ aims at extracting 3D models of deformable objects from 2D keypoint annotations in unconstrained images
- ▶ notice the unique mechanism in disentangle “shape” and “rotation”



6

<sup>6</sup>Novotny D, Ravi N, Graham B, et al. C3DPO: Canonical 3D Pose Networks for Non-Rigid Structure From Motion. In: ICCV2019

- ▶ **input**

$$Y_n = (y_{n,1}, \dots, y_{n,K}) \in \mathbb{R}^{2 \times K}$$

representing  $N$  views  $y_1, \dots, y_N$  of a rigid objects of  $K$  2D keypoints each

- ▶ views are generated from a **latent** single set of 3D keypoints

$$\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_K) \in \mathbb{R}^{3 \times K}$$

- ▶  $N$  rigid motion  $(\mathbf{R}_n, \mathbf{t}_n)$ .
- ▶ standard projection matrix:

$$Y_n = [\mathbf{R}_n \mid \mathbf{t}] \begin{bmatrix} \mathbf{X} \\ 1 \end{bmatrix} = \underbrace{\mathbf{R}_n \mathbf{X} + \mathbf{t}_n}_{\mathbf{P}} \\ \implies Y_{n,k} = \mathbf{R}_n \mathbf{X}_k + \mathbf{t}_n$$

where  $\mathbf{P} : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  is the camera projection function.

- ▶ the paper, “for simplicity” assumes  $\mathbf{P} = [\mathbf{I}_2 \quad \mathbf{0}]$ , and all keypoints are centered together.
- ▶ yields simplified system of equations

$$y_{nk} = \mathbf{M}_n \mathbf{X}_k \quad \text{where } \mathbf{M}_n = \mathbf{R}_n$$

- ▶ Structure from Motion can be formulated as factoring the views  $Y$  into viewpoints  $M$  and structure  $X$ .

$$Y = \begin{bmatrix} y_{1,1} & \cdots & y_{1,K} \\ \vdots & \ddots & \vdots \\ y_{N,1} & \cdots & y_{N,K} \end{bmatrix}, \mathbf{M} = \begin{bmatrix} \mathbf{M}_1 \\ \vdots \\ \mathbf{M}_N \end{bmatrix}, \underbrace{Y}_{2N \times K} = \underbrace{\mathbf{M}}_{2N \times 3} \underbrace{\mathbf{X}}_{3 \times K}$$

- ▶ In Non-Rigid SFM, structure  $\mathbf{X}_n$  is allowed to deform from one view to next
- ▶ Constrain deformation using simple linear model
- ▶ Express  $\mathbf{X}_n$  as:

$$\mathbf{X}_n = \mathbf{X}(\alpha_n, S)$$

1.  $\alpha_n \in \mathbb{R}^D$ : small vector of view-specific pose parameter
2.  $S \in \mathbb{R}^{3D \times K}$ : view-invariant shape basis

- ▶ Non-Rigid Structure from Motion can be formulated as:

$$\underbrace{\mathbf{Y}}_{2N \times K} = \underbrace{\bar{\mathbf{M}}}_{2N \times 3N (\text{ sparse })} \underbrace{(\underbrace{\alpha}_{N \times D} \otimes \mathbf{I}_3)}_{3N \times 3D} \underbrace{\mathbf{S}}_{3D \times K}$$

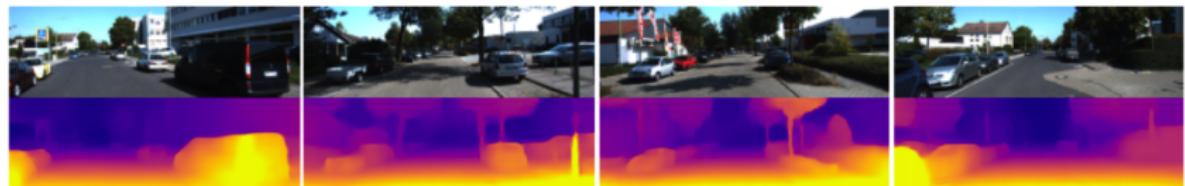
- ▶ The network is learned by minimizing the reprojection loss:

$$\ell_1(Y, v; \Phi, S) = \frac{1}{K} \sum_{k=1}^K v_k \cdot \|Y_k - M(\theta)(\alpha \otimes I_3)S_{:,k}\|_\epsilon$$

$v_k$ : binary indicating whether  $k^{\text{th}}$  keypoint is visible in that particular view

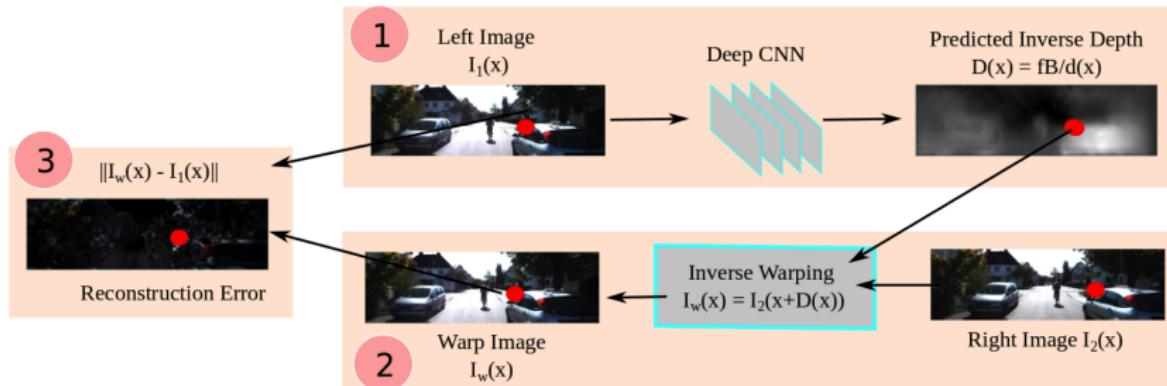
# Depth Estimation

- ▶ Depth estimation: designed to estimate depth from 2D images
- ▶ some works used supervised method: drawback: large amount of manually labelled data (image and corresponding depth)



# Learning from Left-Right Correspondence

Overall architecture of a stereopsis based auto-encoder. <sup>7</sup>



## ► input

1. pair of rectified stereo images
2. known camera motion between the two

<sup>7</sup>Garg R, BG V K, Carneiro G, et al. Unsupervised CNN for single view depth estimation: Geometry to the rescue. In: ECCV2016

# Overall Architecture

- ▶ **encoder** ① CNN map left image ( $I_1$ ) to depth map  $d(x)$ ,
- ▶ then  $d(x)$  allows to estimate disparity  $D(x) = \frac{Bf}{Z}$
- ▶ **decoder** ② knowing disparity  $D(x)$ , one obtains synthesized version of  $I_1$ , i.e.,  $I_1^w$  by backward warp  $I_1^w(x) = I_2(x + D(x))$
- ▶ reconstructed loss between  $I_1$  and  $I_1^w$  ③ via a simple loss.

It's an interesting auto-encoder

1. **Encoder:**  $I_1 \rightarrow d(x)$ , simple! just CNN
2. **Decoder:**  $d(x) \rightarrow I_1^w$ , no neural network but geometry  
however this step needs the help of  $I_2$  and camera parameter ( $B, f, Z$ )

- ▶ **Photometric error** basically just the reconstruction error:

$$\begin{aligned}\mathcal{L}_{\text{reconstruction}}^i &= \int_{\Omega} \left\| \mathbf{l}_w^i(x) - \mathbf{l}_1^i(x) \right\|^2 dx \\ &= \int_{\Omega} \left\| \mathbf{l}_2^i(x + \underbrace{D^i(x)}_{fB/d^i(x)}) - \mathbf{l}_1^i(x) \right\|^2 dx\end{aligned}$$

- ▶ **Smooth error** Use the L2 regularization on the disparity discontinuities as our prior to deal with the aperture problem.

$$\mathcal{L}_{\text{smooth}}^i = \|\nabla D^i(x)\|^2$$

- ▶ Jointly estimate **monocular depth** and **camera motion** from unstructured video sequence.
- ▶ End-to-end model mapping input pixels to:
  1. ego-motion (parameterized as 6-DoF transformation matrices)
  2. underlying scene structure (parameterized as per-pixel depth maps under a reference view)

---

<sup>8</sup>Zhou T, Brown M, Snavely N, et al. Unsupervised learning of depth and ego-motion from video. In: CVPR2017



# View Synthesis as Supervision

- ▶ overall loss is its reconstruction between:
  - ▶ input view of a scene  $\mathbf{I}_t$
  - ▶ a synthesize scene  $\hat{\mathbf{I}}_s$  seen from a different camera pose  $\mathbf{I}_{s \neq t}$

$$\mathcal{L}_{vs} = \sum_s \sum_p |\mathbf{I}_t(p) - \hat{\mathbf{I}}_s(p)|$$

- ▶ in order to obtain a  $\hat{\mathbf{I}}_s$  seen from a different camera pose  $\mathbf{I}_{s \neq t}$ , it trains
  1. **monocular depth** generation and
  2. **camera pose** between views

Together, they learned the relationship of how a corresponding pixel transform from one view to another

## How a corresponding pixel transform from one view to another

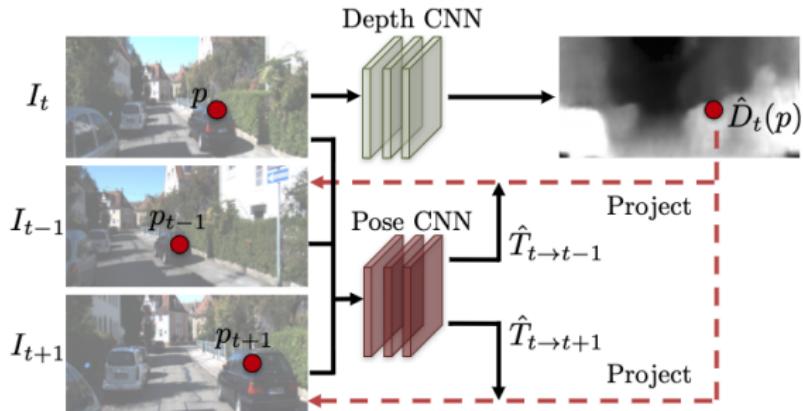
- ▶ reconstruct target view  $\mathbf{I}_t$  by sampling pixels from a source view  $\mathbf{I}_s$  based on the predicted depth map  $\hat{D}_t$  and relative pose  $\hat{T}_{t \rightarrow s}$ .

$$p_s \propto \mathbf{K} \hat{T}_{t \rightarrow s} \underbrace{\hat{D}_t(p_t) \mathbf{K}^{-1} p_t}_{\text{X in } t \text{ coordinate}}$$

- ▶ learned both  $\hat{D}_t$  and  $\hat{T}_{t \rightarrow s}$  help to obtain the reverse mapping:  $s \rightarrow t$

# Overall Architecture

Overview of the supervision pipeline based on view synthesis.



- ▶ **Depth CNN** takes only target view as input, outputs per-pixel depth map  $\hat{D}_t(p_t)$
- ▶ **Pose CNN** takes both the target view ( $I_t$ ) and the nearby/source views (e.g.,  $I_{t-1}$  and  $I_{t+1}$ ) as input, and outputs the relative camera poses ( $\hat{T}_{t \rightarrow t-1}$ ,  $\hat{T}_{t \rightarrow t+1}$ ).

**inverse warp** the source views to reconstruct the target view  $\hat{\mathbf{I}}(s)$ , and the photometric reconstruction loss is used:

$$\mathcal{L}_{vs} = \sum_s \sum_p \left| \mathbf{I}_t(p) - \hat{\mathbf{I}}_s(p) \right|$$