

0. paper

Hand Keypoint Detection in Single Images using Multiview Bootstrapping

1. motivation

Due to heavy occlusions, even manual hand keypoint annotations are difficult to get right.

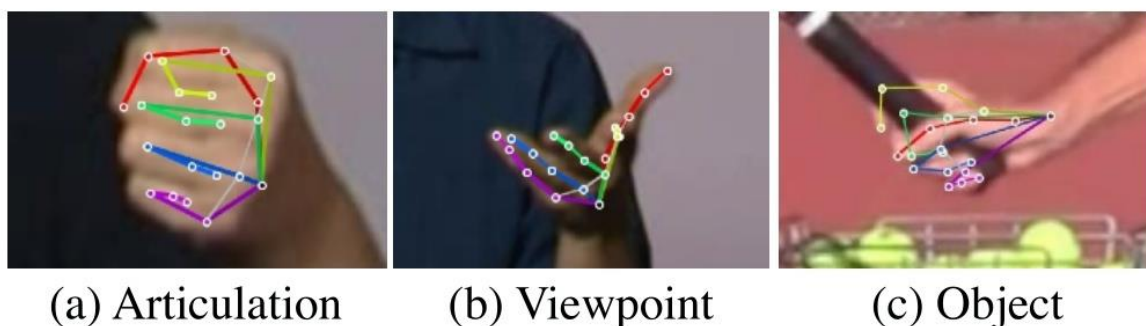


Figure 2: Hand annotation is difficult in single images because joints are often occluded due to (a) articulations of other parts of the hand, (b) a particular viewing angle, or (c) objects that the hand is grasping.

based on the observation: even if a particular image of the hand has significant occlusion, there often exists an **unoccluded view**.

2. multiview bootstrapped training

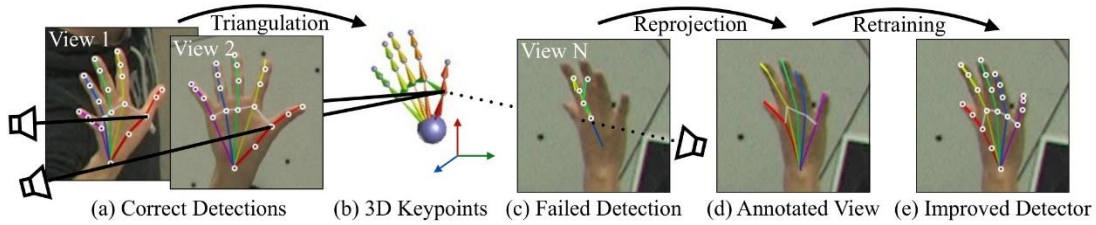


Figure 3: Multiview Bootstrapping. (a) A multiview system provides views of the hand where keypoint detection is easy, which are used to triangulate (b) the 3D position of the keypoints. Difficult views with no keypoint detections can be annotated using the reprojected 3D keypoints, and used to retrain (e) an improved detector that now works on difficult views.

Given the **initial keypoint detector d_0** and a dataset of **unlabeled multiview images**, our objective is to use the detector to generate a set of labeled images, which can be used to train an improved detector d_1 .

$$d(\mathbf{I}) \mapsto \{(\mathbf{x}_p, c_p) \text{ for } p \in [1 \dots P]\} . \quad (1)$$

detector $d(\mathbf{I})$;

Image: \mathbf{I} ;

P : hand keypoint number;

\mathbf{x}_p : keypoint **2d location**;

c_p : **detection confidence** of \mathbf{x}_p .

Training set with annotations:

$$\mathcal{T}_0 := \{(\mathbf{I}^f, \{\mathbf{y}_p^f\}) \text{ for } f \in [1 \dots N_0]\} , \quad (2)$$

for frame f , y is the labeled keypoint, total N_0 pairs.

Algorithm 1 Multiview Bootstrapping

Inputs:

- Unlabeled images: $\{\mathbf{I}_v^f \text{ for } v \in \text{views}, f \in \text{frames}\}$
- Keypoint detector: $d_0(\mathbf{I}) \mapsto \{(\mathbf{x}_p, c_p) \text{ for } p \in \text{points}\}$
- Labeled training data: \mathcal{T}_0

for iteration i in 0 to K :

1. Triangulate keypoints from weak detections

for every frame f :

 (a) Run detector $d_i(\mathbf{I}_v^f)$ on all views v (Eq. (5))

 (b) Robustly triangulate keypoints (Eq. (6))

2. Score and sort triangulated frames (Eq. (7))

3. Retrain with N -best reprojections (Eq. (8))

$d_{i+1} \leftarrow \text{train}(\mathcal{T}_0 \cup \mathcal{T}_{i+1})$

Outputs: Improved detector $d_K(\cdot)$ and training set \mathcal{T}_K

Given \mathbf{V} views of an object in a particular frame \mathbf{f} :

$$\mathbf{X}_p^f = \arg \min_{\mathbf{X}} \sum_{v \in \mathcal{I}_p^f} \|\mathcal{P}_v(\mathbf{X}) - \mathbf{x}_p^v\|_2^2, \quad (6)$$

Try to **minimize the reprojection error** to obtain the final triangulated position.

The triangulation is RANSAC approach.

Then order the frames by P keypoints` score:

$$\text{score}(\{\mathbf{X}_p^f\}) = \sum_{p \in [1 \dots P]} \sum_{v \in \mathcal{I}_p^f} c_p^v. \quad (7)$$

知乎 @banana16314

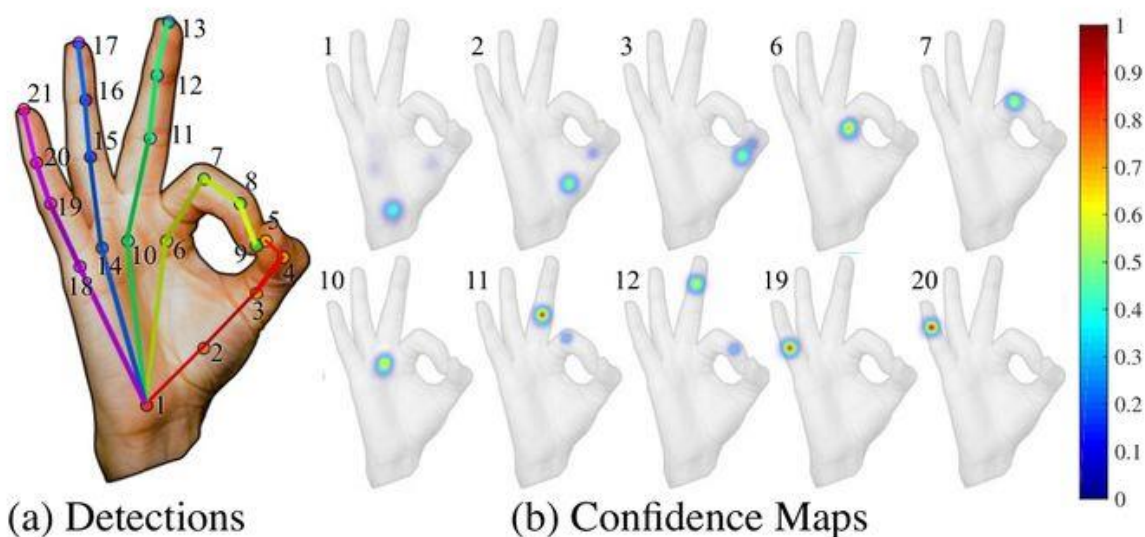


Figure 4: (a) Input image with 21 detected keypoints. (b) Selected confidence maps produced by our detector, visualized as a “jet” colormap overlaid on the input.

知乎 @banana16314

finally, Retraining with N-best Reprojections.