

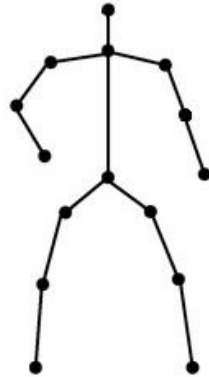
Human Pose Lecture

Dr. Sergio Escalera
University of Barcelona

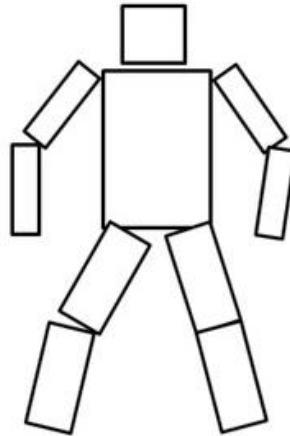


What is articulated human pose?

- Given body kinematic tree, human pose is defined as the vector of joints locations either in 2D (image coordinate in pixels) or 3D (world coordinate in millimeters).
- Human body pose estimation is a way of representing humans in the images,
- It can be represented and estimated from coarse to fine models.
- Classically addressed by pictorial structures and graphical models



Kinematic model



Cardboard model



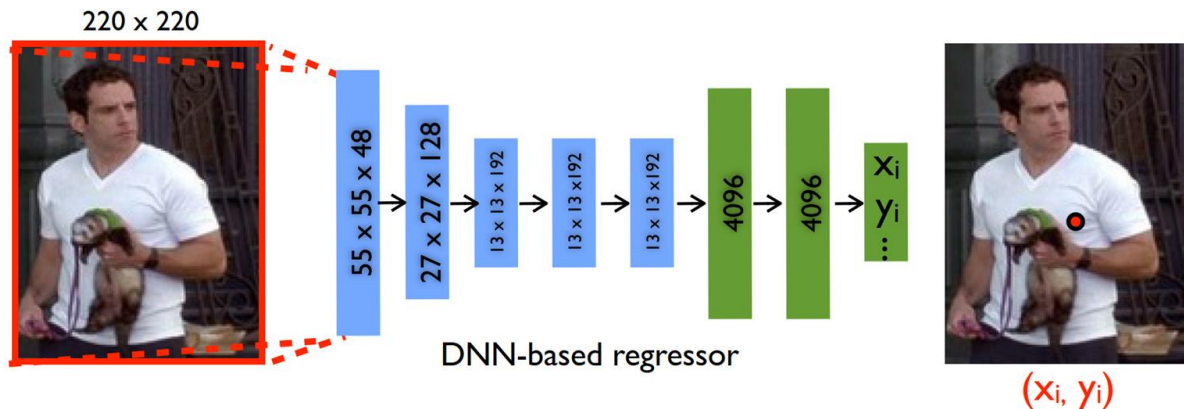
Volumetric model

Challenges



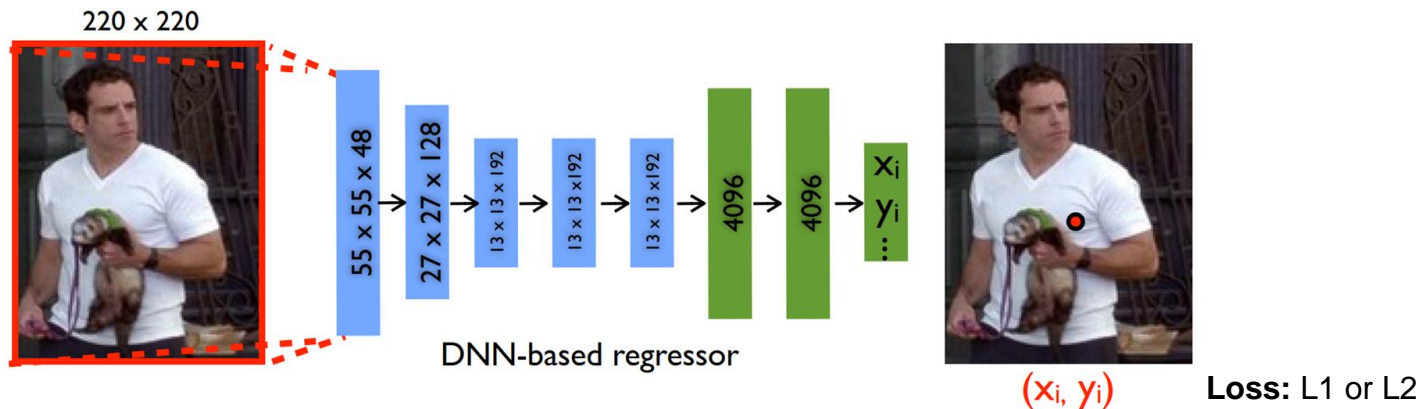
Human pose regression

- The output of the network is directly a vector of (x,y) joints locations,
- This is a difficult task for the network because
 - Pose vector is a highly nonlinear variable,
 - Network must deal with scale and translation as well.



Human pose regression

- The output of the network is directly a vector of (x,y) joints locations,
- This is a difficult task for the network because \leftarrow **A better alternative:** a dense probability map for each joint.
 - Pose vector is a highly nonlinear variable, \leftarrow **Solution:** cascade of pose regressors.
 - Network must deal with scale and translation as well. \leftarrow **Solution:** crop person bounding box as pre-processing.

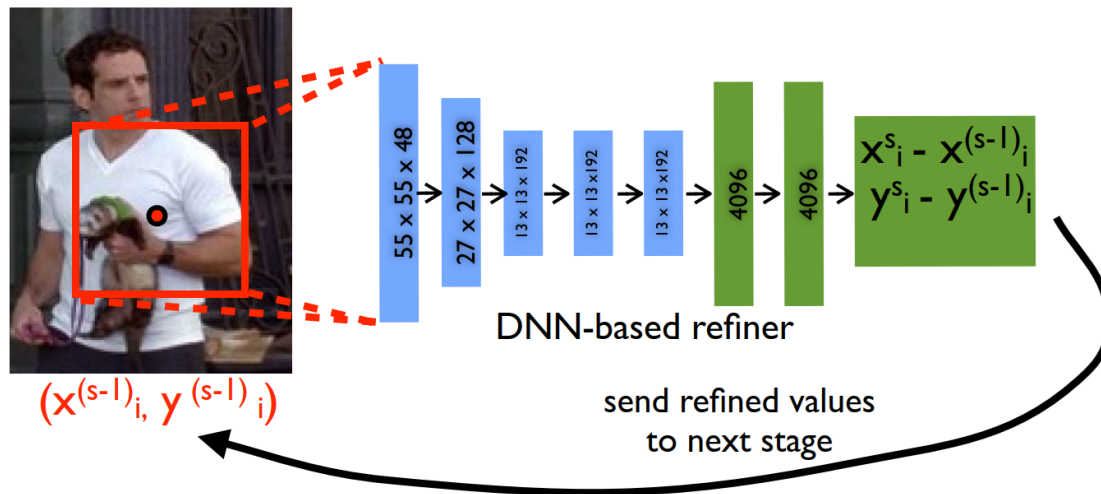


Architecture: AlexNet

Cascade of pose regressors

- Given an initial estimation of the joints (in holistic view), iteratively refines the error (in local view).

Stage s

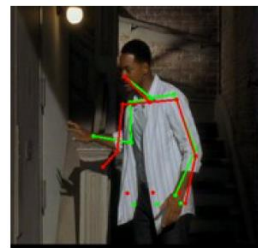


Cascade of pose regressors

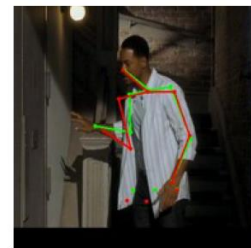
Limitations:

- Sensitive to initial estimation,
- Local solution, easily can stick to local minima,
- Lack of structured predictions.

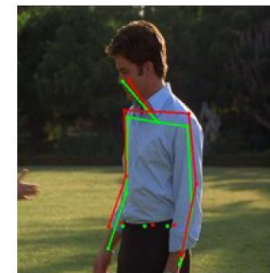
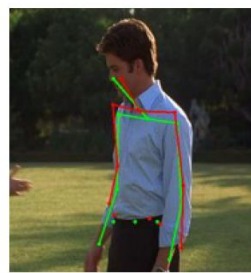
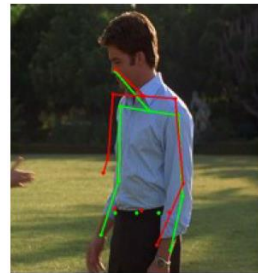
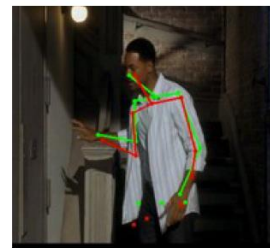
Initial stage 1



stage 2



stage 3

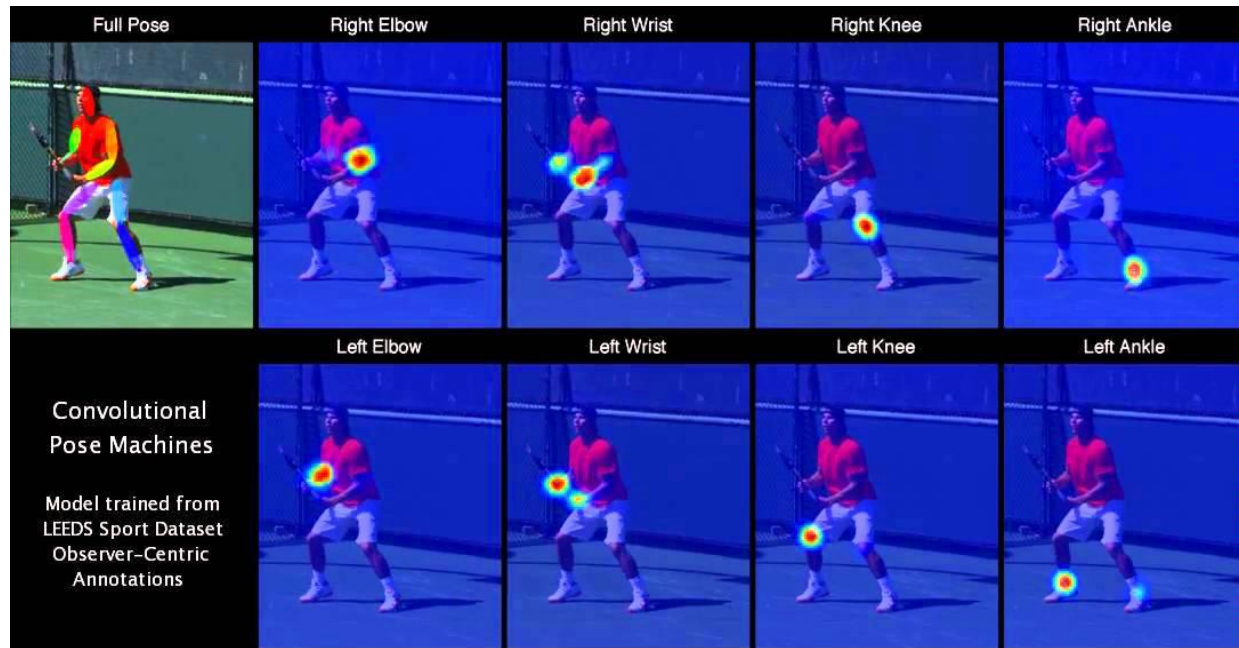


Results on LSP dataset



Joints heatmaps instead of pose vector

- It is an easier problem for the network,
- Still can be combined with pose regressors,
- The map can be fed into graphical models to learn higher order joint relationships.



Context is important

Which part belongs to a human body?



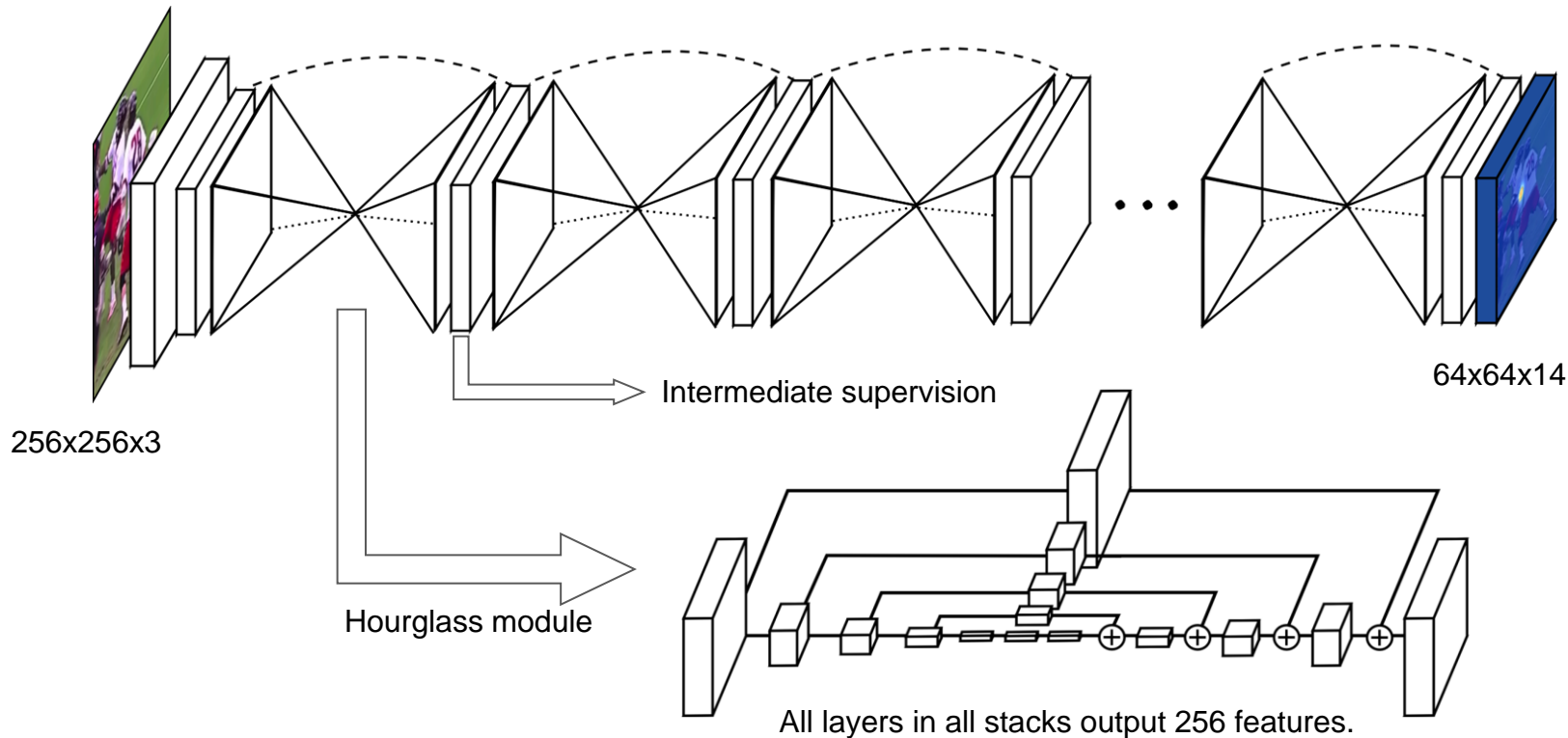
Context is important

Which part belongs to a human body?

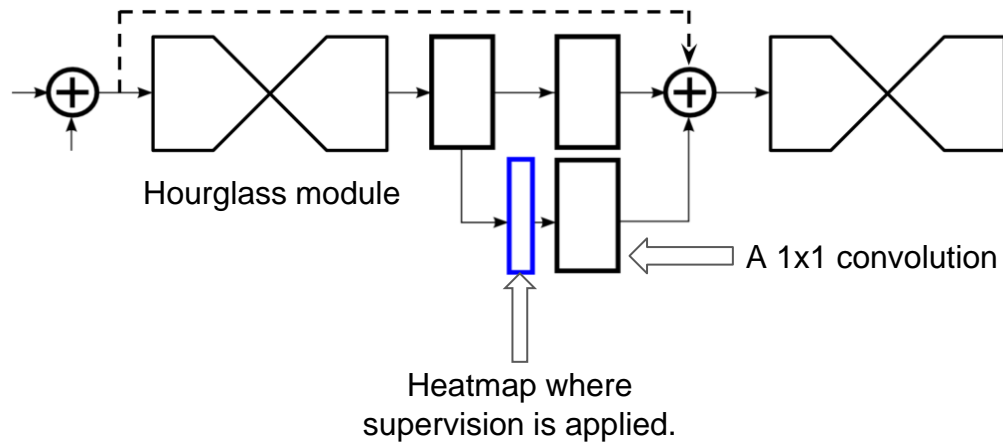
- ☐ Local evidence is weak,
- ☐ Larger receptive field = more context,
- ☐ Recover from failures by cascading.



Stacked Hourglass Network



Stacked Hourglass Network: intermediate supervision

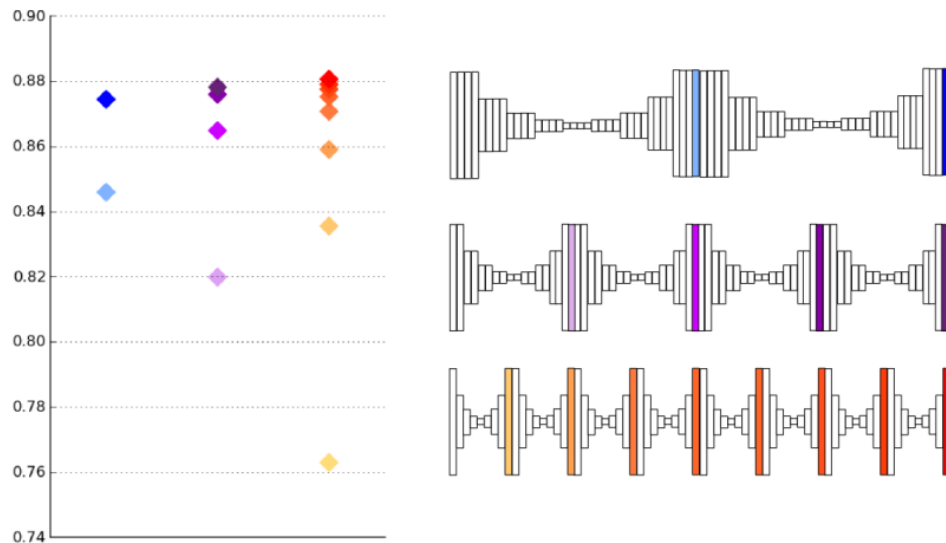


Stacked Hourglass Network

Cons:

- Quite heavy in GPU,
- Joints may be confused with background,
- Still does not explicitly deal with structure.

Intermediate Prediction Accuracy (Validation, PCKh@0.5)



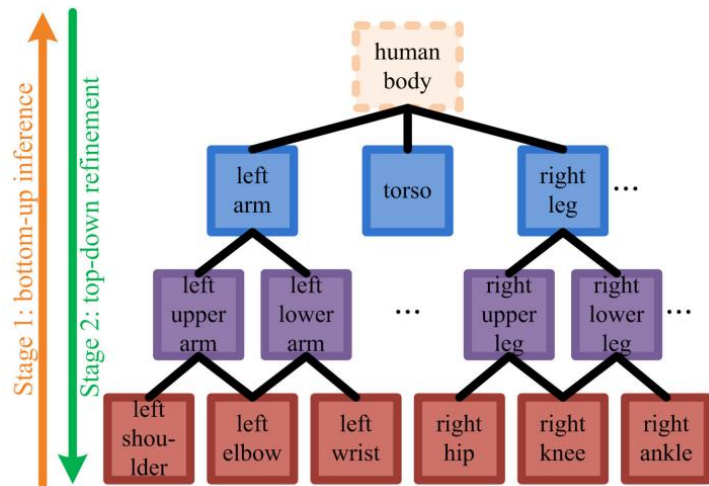
Some results on MPII dataset



Deeply learned compositional models

1. We have a compositional model of part-subpart relationships,
2. Traditional solutions by modelling the graph (or tree) with Gibbs formulation,
3. We want a bottom-up and top-down inference,

Question: How to model it with CNN models?



$$p(\Omega|\mathbf{I}) = \frac{1}{Z} \exp\{-E(\Omega, \mathbf{I})\}$$

$$S(\Omega) \equiv -E(\Omega, \mathbf{I}) = \sum_{u \in \mathcal{V}^{leaf}} \phi_u^{leaf}(w_u, \mathbf{I}) + \sum_{u \in \mathcal{V}^{and}} \sum_{v \in ch(u)} \phi_{u,v}^{and}(w_u, w_v)$$

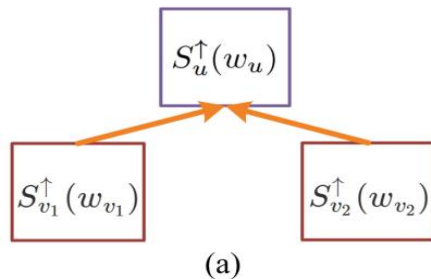
↑
Score to be
minimized

↑
Unary term

↑
Pairwise term

First define updating rule

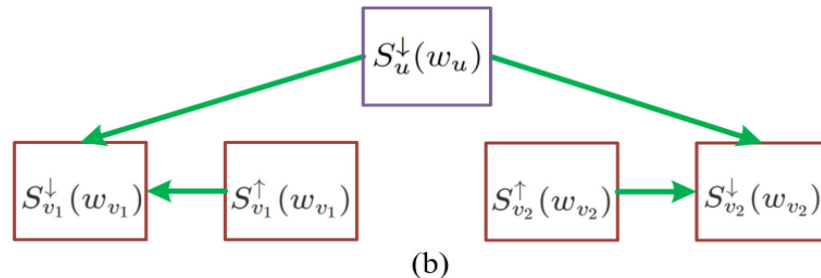
Bottom-up inference:



$$(\text{Leaf}) \quad S_u^\uparrow(w_u) = \phi_u^{leaf}(w_u, \mathbf{I})$$

$$(\text{And}) \quad S_u^\uparrow(w_u) = \sum_{v \in ch(u)} \max_{w_v} [\phi_{u,v}^{and}(w_u, w_v) + S_v^\uparrow(w_v)]$$

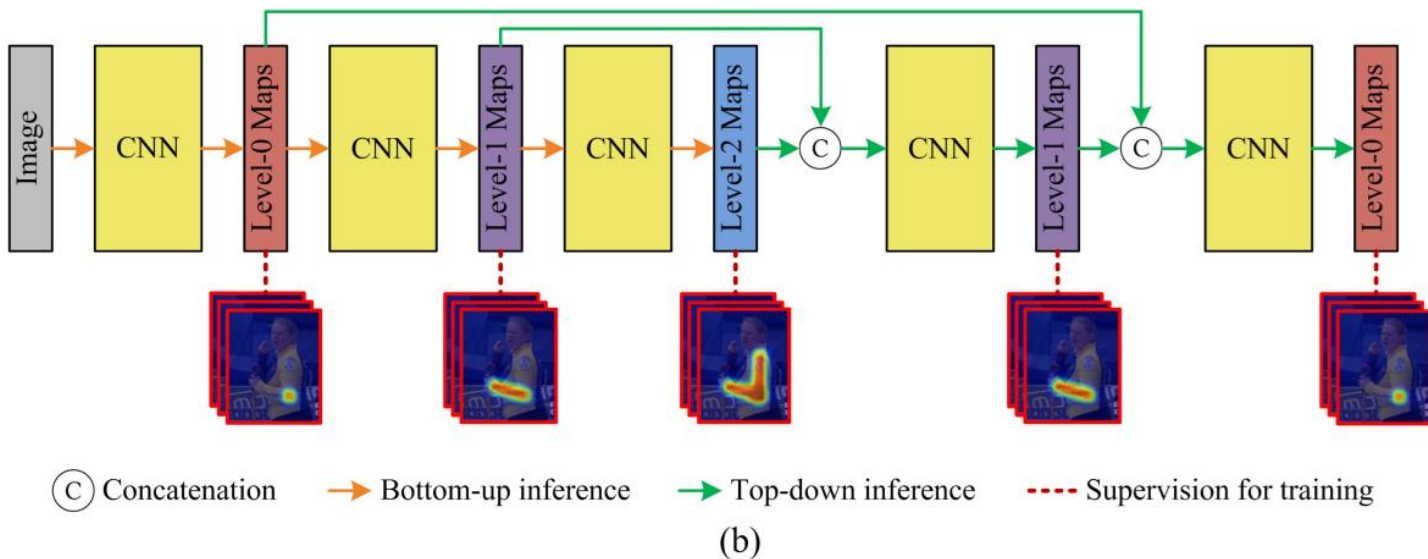
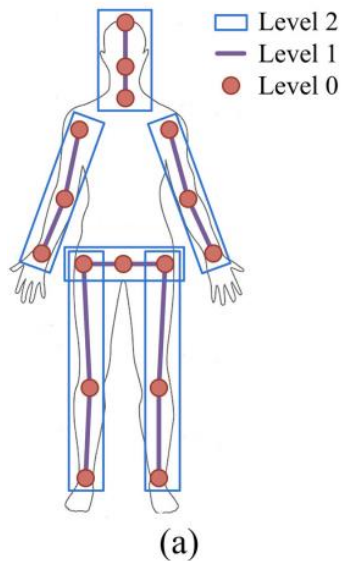
Top-down inference:



$$(\text{Root}) \quad w_u^* = \operatorname{argmax}_{w_u} S_u^\downarrow(w_u) \equiv \operatorname{argmax}_{w_u} S_u^\uparrow(w_u)$$

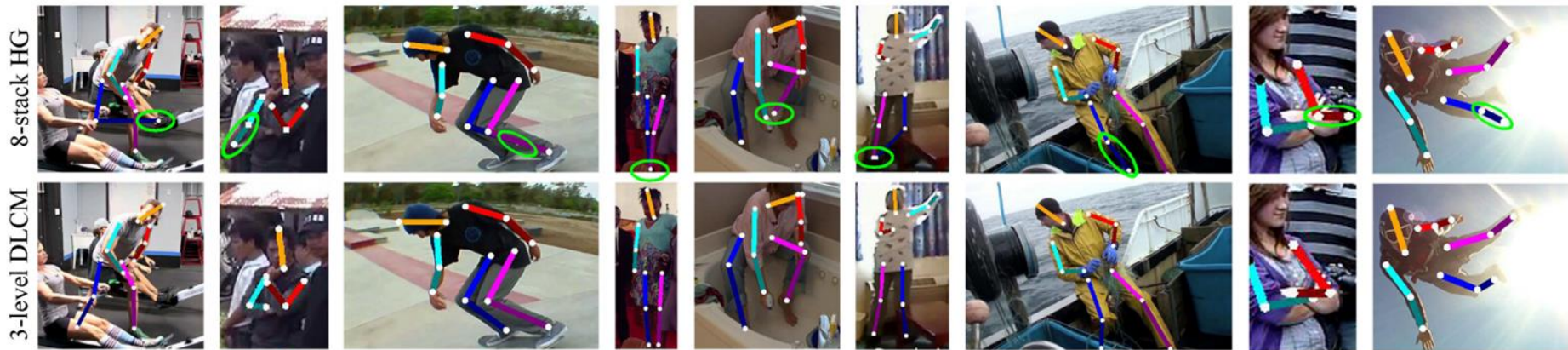
$$(\text{Non-root}) \quad w_v^* = \operatorname{argmax}_{w_v} S_v^\downarrow(w_v) \equiv \operatorname{argmax}_{w_v} [\phi_{u,v}^{and}(w_u^*, w_v) + S_v^\uparrow(w_v)]$$

Revisit stacked hourglass network



Heatmaps of level i ($i > 0$) are composed of heatmaps of level 0 and level i .
CNN \rightarrow Hourglass Module

Some results comparing to 8-stack hourglass

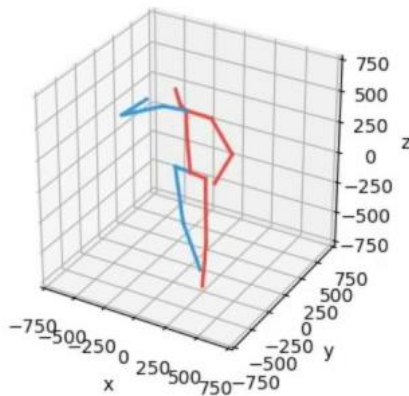


3D human body pose

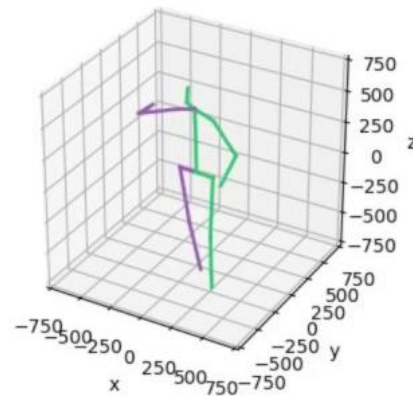
- 3D pose is the vector of body joints in 3D space,
- 2D pose is the projection of 3D pose to image plane,
- 3D pose is defined with the same models as 2D pose,
- 3D pose estimation is a challenging task since depth dimension is lost during projection to RGB image.



2d observation



3d ground truth



3d prediction


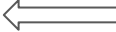
3D human body pose estimation

Solutions:

- Directly regress 3D pose from RGB image,
- Estimate 2D pose and then regress 3D pose from 2D,
- Apply volumetric heatmaps and estimate 3D pose,
- Apply multi-task learning by the use of different modalities.

3D human body pose estimation

Solutions:

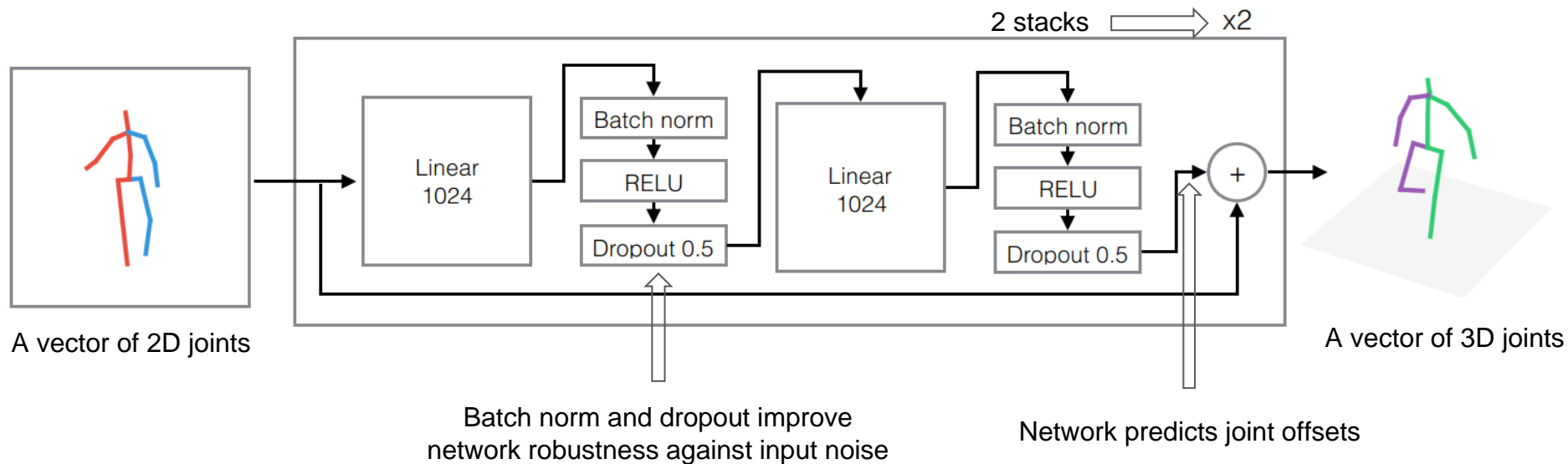
- Directly regress 3D pose from RGB image,  A difficult problem to the network similar to 2D pose.
- Estimate 2D pose and then regress 3D pose from 2D,
- Apply volumetric heatmaps and estimate 3D pose,
- Apply multi-task learning by the use of different modalities.  Annotating data for all modalities is not a trivial task for many datasets.

3D human body pose estimation

Solutions:

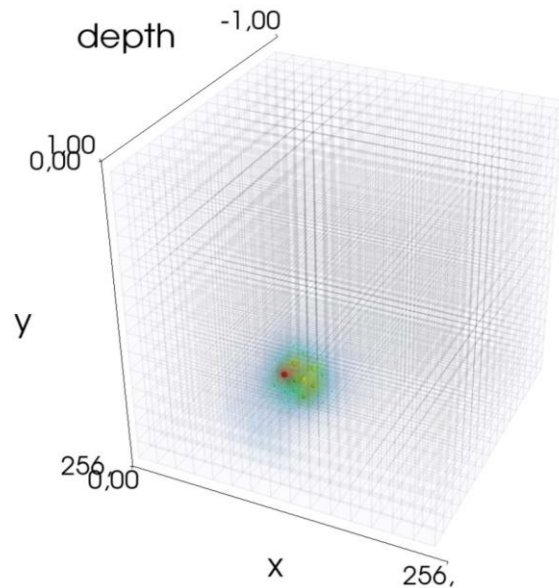
- Directly regress 3D pose from RGB image,
- Estimate 2D pose and then regress 3D pose from 2D,
- Apply volumetric heatmaps and estimate 3D pose,
- Apply multi-task learning by the use of different modalities.

Lifting 2D joints to 3D: A simple solution

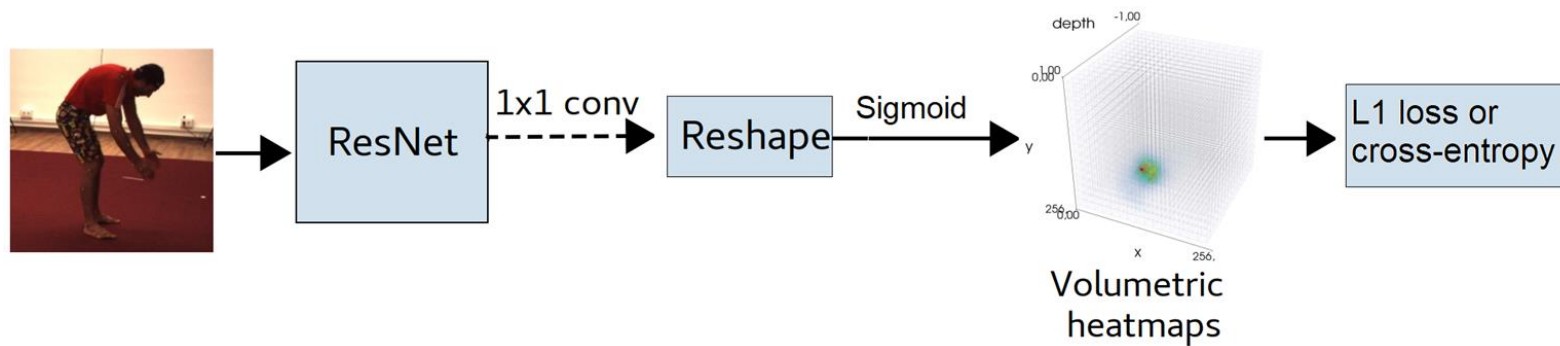


Volumetric heatmaps for 3D pose estimation

- A volumetric heatmap is defined as a tensor in the form of ($\#Width \times \#Height \times \#Depth \times \#Joints$),
- A joint's depth is a continuous value which is discretized into several bins, i.e. $\#Depth$,
- Ground truth heatmap can be defined by a Gaussian.



A simple solution



A simple solution - data augmentation



2638 occluder objects from Pascal VOC

Filter out 'person', 'truncated', 'difficult' and small object segments



Augmented inputs with pasted occluders

Applied with 50% probability, 1–8 objects,
at random scale, at random position

Applications



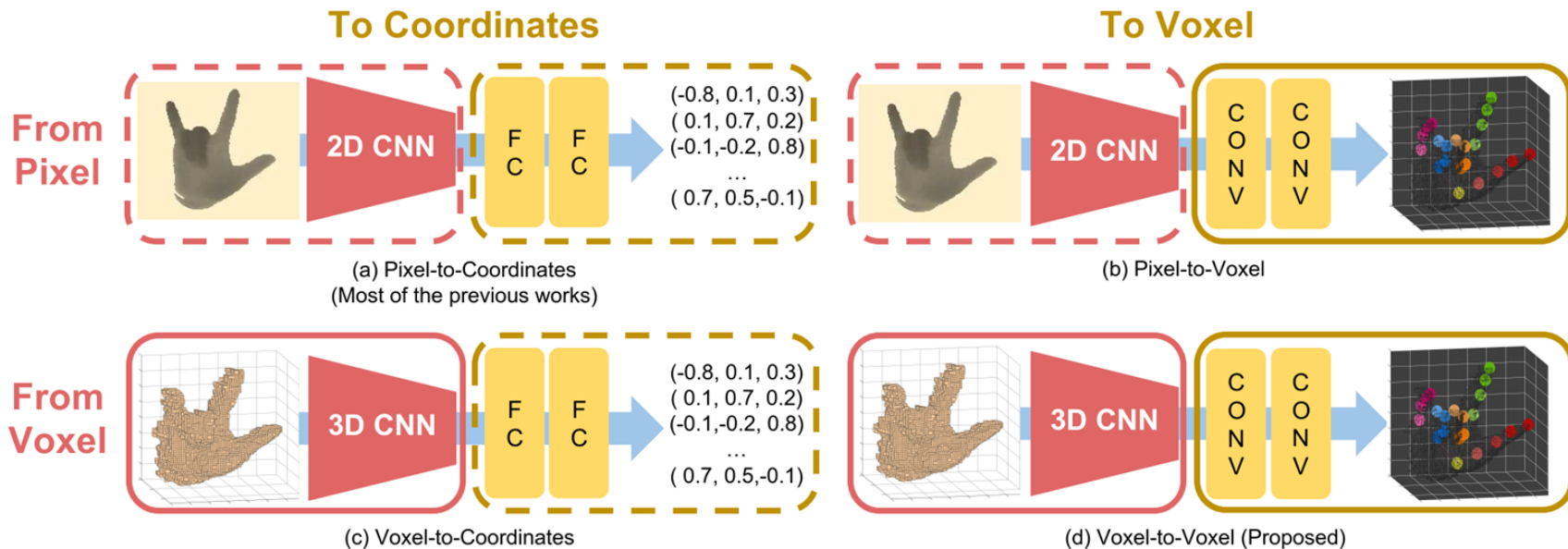
3D hand pose estimation based on depth images (similar strategies applied on RGB images)



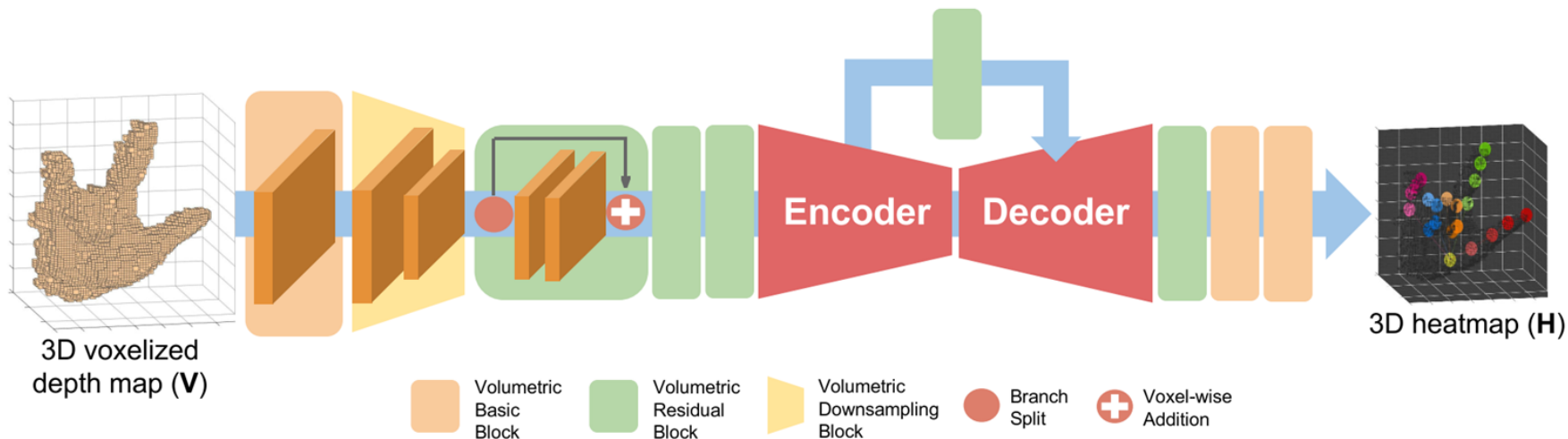
Applications

- Human-computer interaction,
- Virtual reality,
- Training robot hands,
- Sign language and gesture recognition.

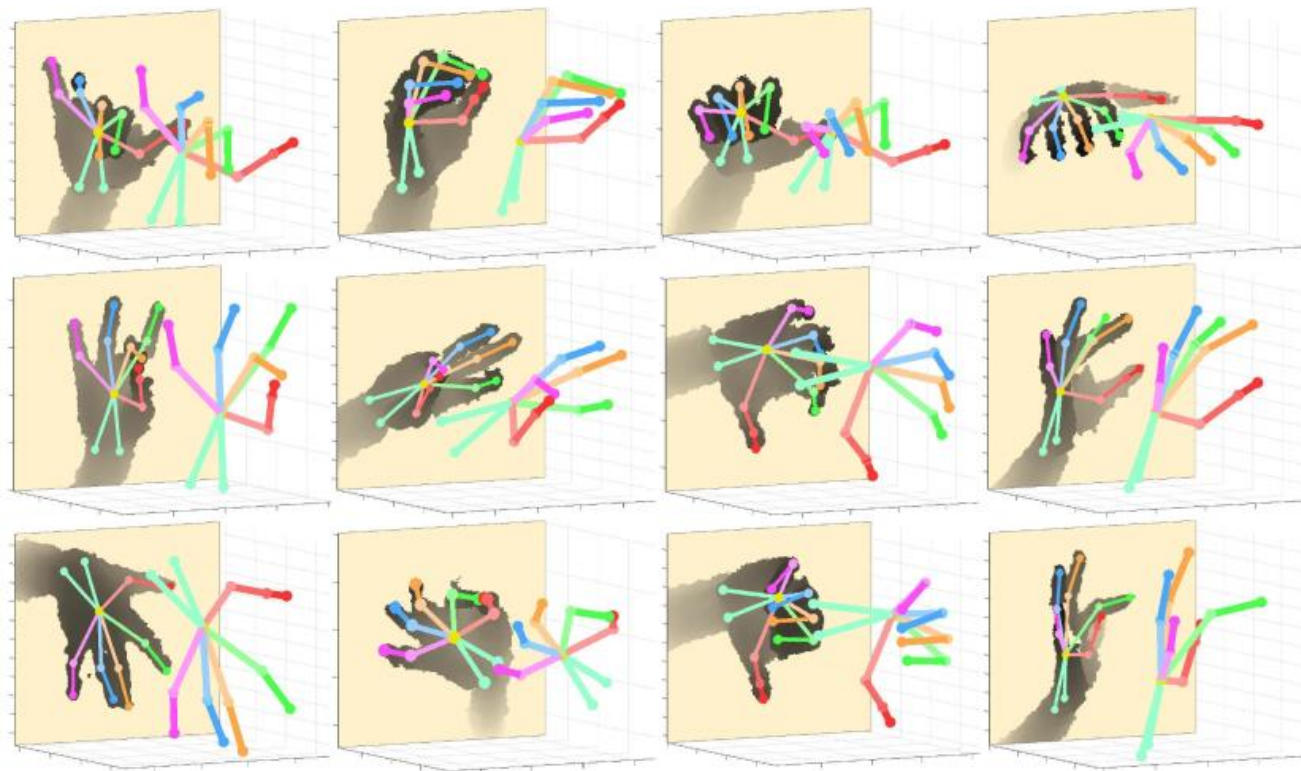
Voxel to voxel predictions



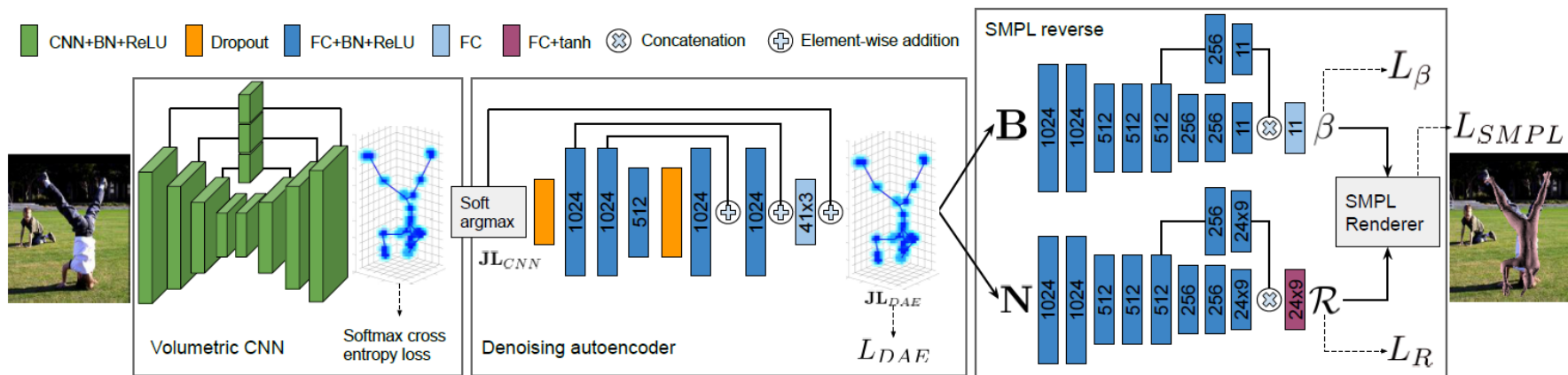
Voxel to voxel predictions



Voxel to voxel predictions - results on NYU dataset



Bonus 1: Volumetric models from 2D (even 3D clothes)



SMPLR: Deep SMPL reverse for 3D human pose and shape recovery
Meysam Madadi, Hugo Bertiche, Sergio Escalera

Bonus 2: Faces

Automatic human face and body analysis



(from Martinez (2019). Context may reveal how you feel. PNAS)

Face

- Identity of the person
- *Perceived* age
- *Perceived* gender
- *Perceived* attractiveness
- Ethnicity/race/skin color
- Head pose
- Gaze direction
- Facial expression of *emotion*: crying (of sadness)
- Other facial features

+ Body

- Body pose
- Hand pose
- Bodily expression of *emotion*: crying (of despair)
- Other body features

+ Context

- Activity/actions/intentions: concert
- Experienced emotion: crying of joy
- Interaction with other people/objects

Bonus 2: Faces

Automatic human face and body analysis

Applications for good

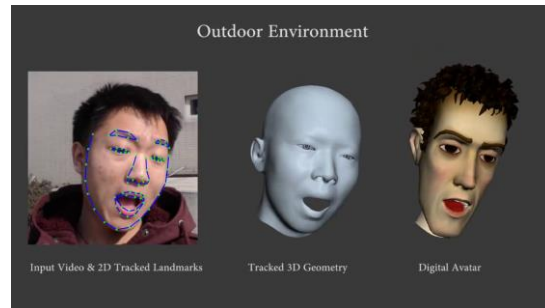
- Human-Computer/Robot interaction
- Early intervention and medical diagnosis (e.g., ASD, depression)
- Augmented reality
- Image synthesis (e.g., digital avatars)
- Education (measuring engagement, effectiveness)
- Gaming
- Driving assistance
- Market research
- Assistive living (e.g., image captioning, fall detection)

Personalized empathic virtual agents



(from Justo et al (2020) Analysis of the interaction between elderly people and a simulated virtual coach. JAIHC)

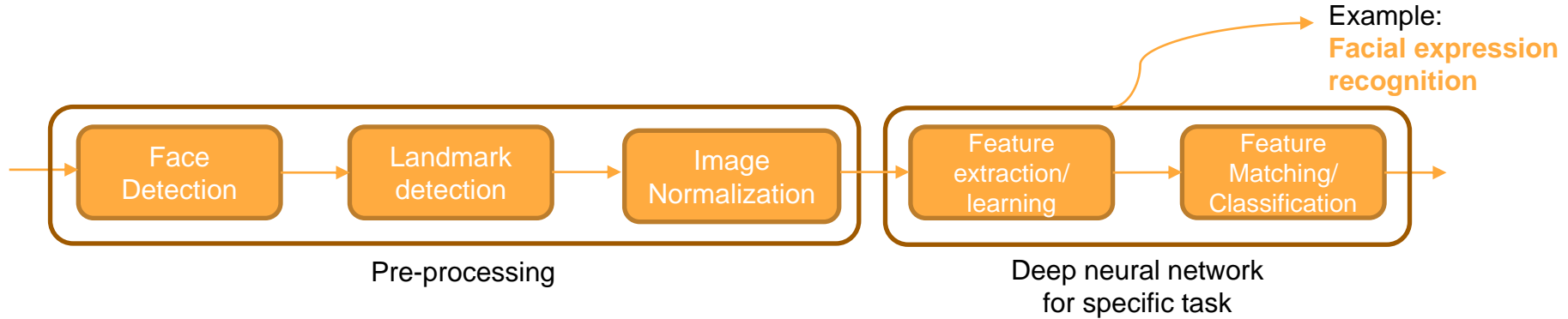
Digital avatars



(from https://gigazine.net/gsc_news/en/20140813-real-time-facial-tracking/) 36

Bonus 2: Faces

Face analysis pipeline



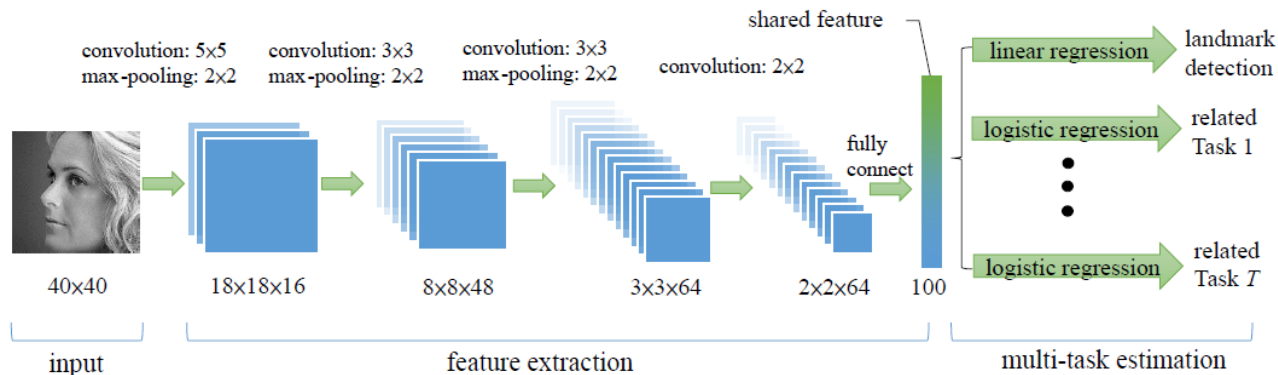
Pre-processing:

1. **Face detection:** detect bounding box surrounding the face.
→ Bounding box coordinates regression (xmin, xmax, ymin, ymax).
2. **Landmark detection:** detect facial landmarks (fiducial points or keypoints).
→ Keypoint coordinates regression (x,y)
3. **Image normalization:** align and normalize the image to reduce variation in face scale and in-plane rotation.

Bonus 2: Faces

Multi-task learning: *optimizing facial landmark detection (the main task) with related/auxiliary tasks.*

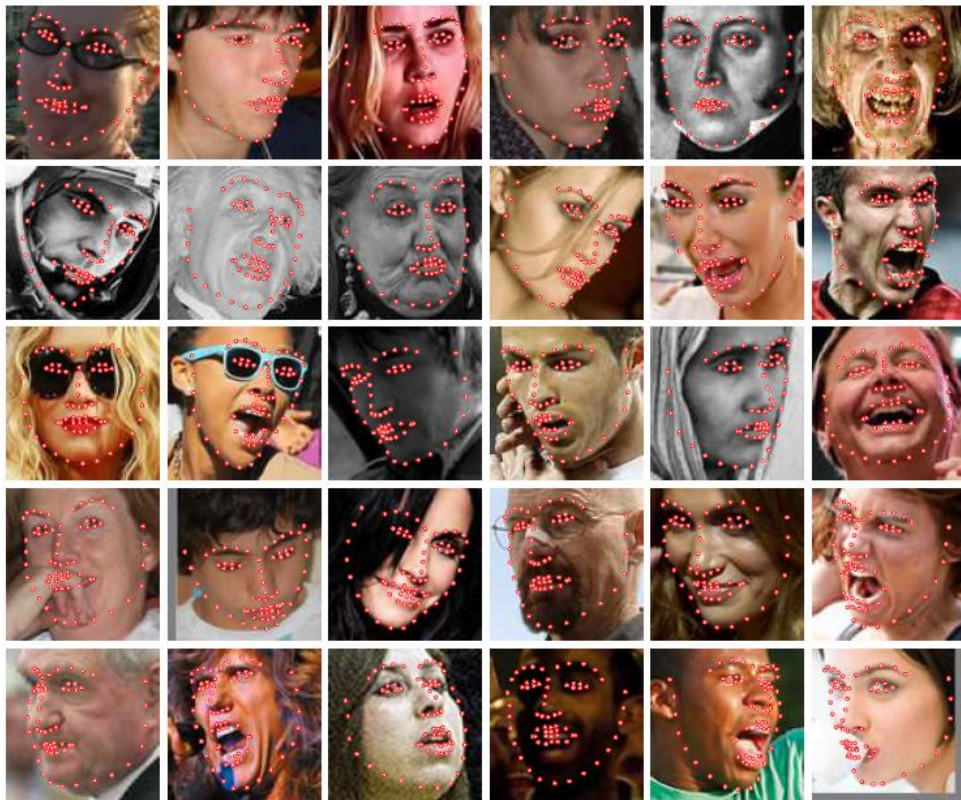
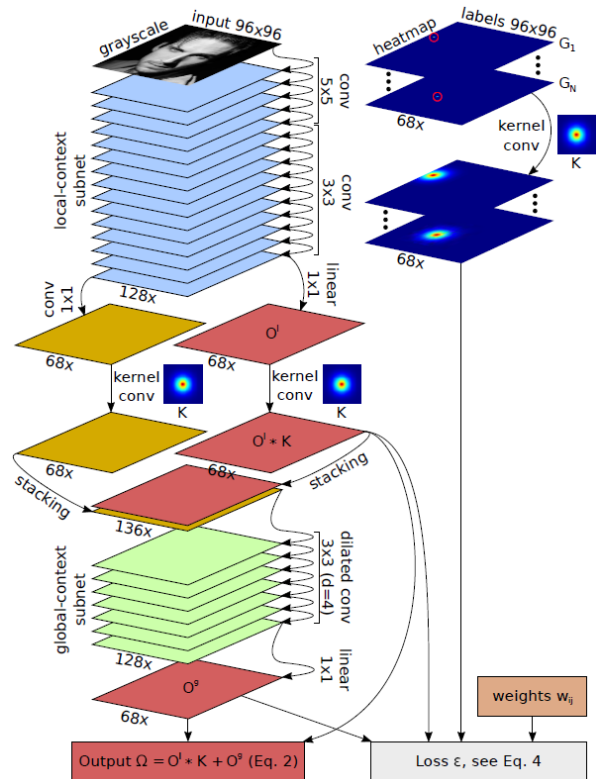
2014 – TDCN*



TDCN								
Auxiliary Tasks	wearing glasses	×	×	✓	×	✓	×	×
	smiling	×	✓	×	×	×	×	×
	gender	female	male	female	female	male	male	female
	pose	right profile	frontal	frontal	left	frontal	frontal	right profile

*Zhang et al. (2014) Facial landmark detection by Deep Multi-task Learning. ECCV

Bonus 2: Faces



Robust Facial Landmark Detection via a Fully-Convolutional Local-Global Context Network, 2019

Human Pose Lecture

Dr. Sergio Escalera
University of Barcelona