

G-RMI Keypoints Detection

COCO Visual Recognition Challenges Workshop @ ECCV 2016

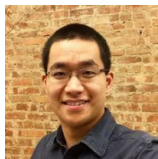
Presenter: George Papandreou (gpapan@google.com)

Team: Tyler Zhu, Nori Kanazawa, George Papandreou, Alex Toshev,
Hartwig Adam, Chris Bregler, Kevin Murphy, Jonathan Tompson

Google Research and Machine Intelligence



Team Members



Tyler
Zhu



Nori
Kanazawa



George
Papandreou



Alex
Toshev



Hartwig
Adam



Chris
Bregler



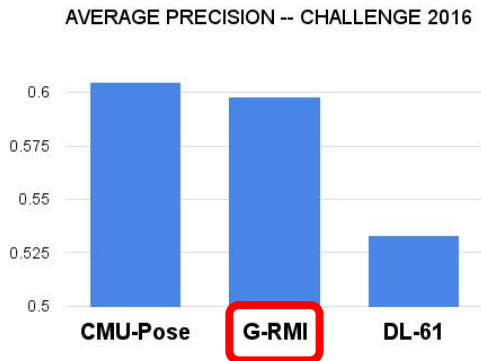
Kevin
Murphy



Jonathan
Tompson

Summary: G-RMI Keypoint Detection System

- Two-stage system:
 - Box person detector
 - Human pose estimator
- Ranked #2. AP during competition:
 - **CMU-Pose**: 0.605 (challenge), 0.618 (testdev)
 - **G-RMI**: 0.598 (challenge), 0.605 (testdev)
 - **DL-61**: 0.533 (challenge), 0.544 (testdev)
 - ...
- AP post-competition:
 - **G-RMI**: 0.668 (testdev)
- Key technical aspects:
 - State-of-art person box detector
 - Pose estimator featuring highly localized keypoint activation maps
 - Effective box proposal rescoring by the pose estimator



● Comparison of using COCO-only as well as COCO + in-house data for training

System Overview



(1) Person detection



(2) Pose estimation

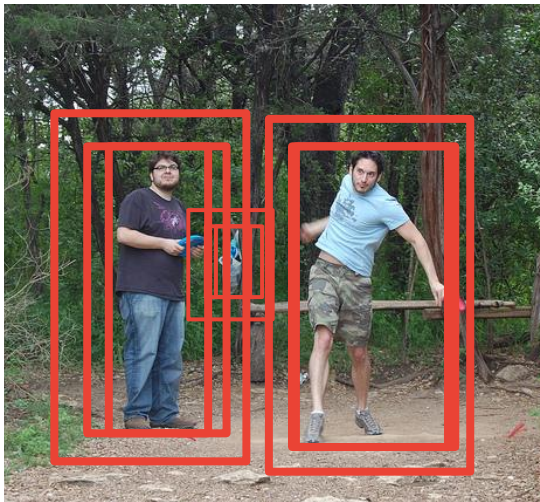
Person Detection

1. **G-RMI** Box Detection entry trained on COCO data (Inception-ResNet Faster-RCNN model ensemble)
 - **0.584** person keypoint AP.
2. Person-specific detector trained on COCO + ImageNet + in-house dataset (single model, multi-crop)
 - **0.592** person keypoint AP.
3. Our best result (before the deadline): Union of (1) + (2)
 - **0.605** person keypoint AP.



More info: G-RMI detection team presentation

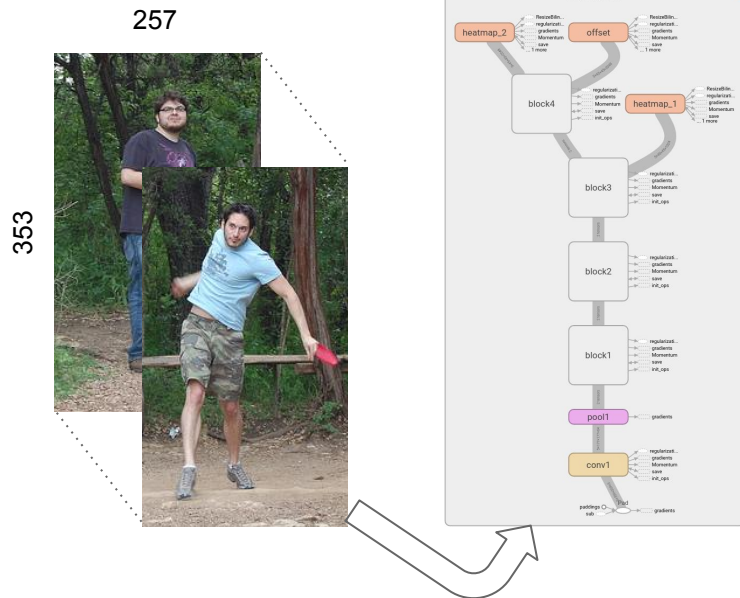
From Box to Pose Proposals



- Aspect ratio normalization
 - $\text{height/width} = 353/257$
- Box enlargement
 - Train scale factor in $[1.0, 1.5]$
 - Eval scale factor 1.25
- Crop extraction
 - height = 353
 - width = 257



Pose Estimation Network



- Single ImageNet-pretrained Resnet-101 producing heatmaps and offsets fully-convolutionally¹
- Dense (stride=8) feature extraction via atrous convolution², followed by bilinear interpolation
 - 353x257 crop → 45x33 feature maps → 353x257 feature maps
- Intermediate supervision, similar to MPII's DeeperCut system³

1. He, K., Zhang, X., Ren, S., & Sun, J. (2015). Deep residual learning for image recognition. CVPR 2016
2. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., & Yuille, A. L. (2016). DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs. arXiv:1606.00915.
3. Insafutdinov, E., Pishchulin, L., Andres, B., Andriluka, M., & Schiele, B. (2016). DeeperCut: A Deeper, Stronger, and Faster Multi-Person Pose Estimation Model. arXiv:1605.03170.

Powerhorse: Tensorflow and TF-Slim



- Scalable distributed infrastructure
- Multi-machine + Multi-GPU training
- Rich library of SoA vision models:
 - Inception
 - ResNet
 - Inception-Resnet
 - ...



Your laptop



Datacenters



Mobile



Raspberry
Pi



Tensor
Processing Unit

Pose Estimation Net: Heatmap Output



- Heatmap field for each keypoint
 - 17 channels (1 within a disk around each keypoint, 0 outside)
 - Sigmoid cross entropy loss
- CNN layers 52 (intermediate) and 101 (final)

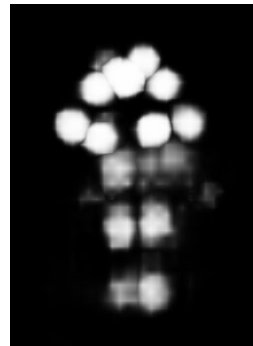
Crop



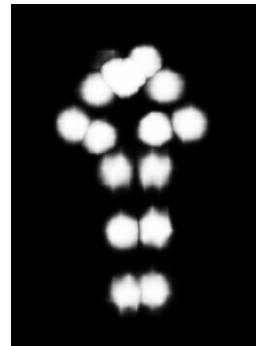
Target



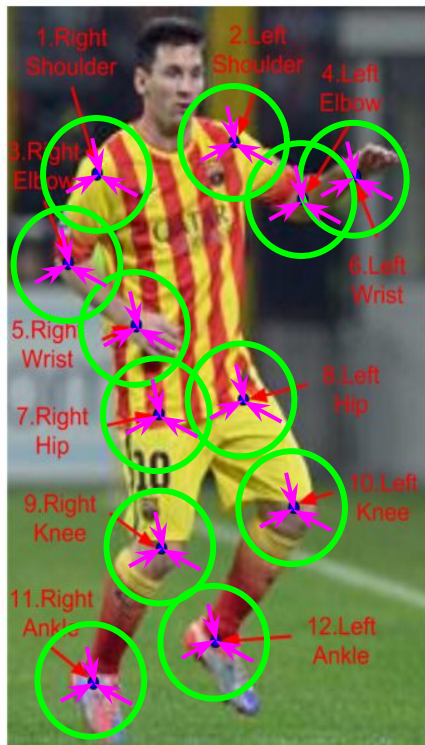
Net-Layer52



Net-Layer101



Pose Estimation Net: Offset Output

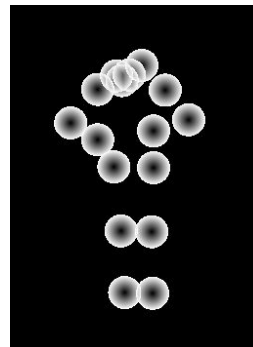


- Offset field towards the center of the disk
 - 34 channels for x- and y- offsets
 - Huber loss, only active within disks
- Only at CNN layer 101 (final)

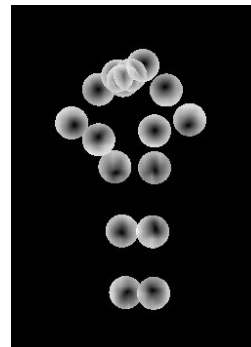
Crop



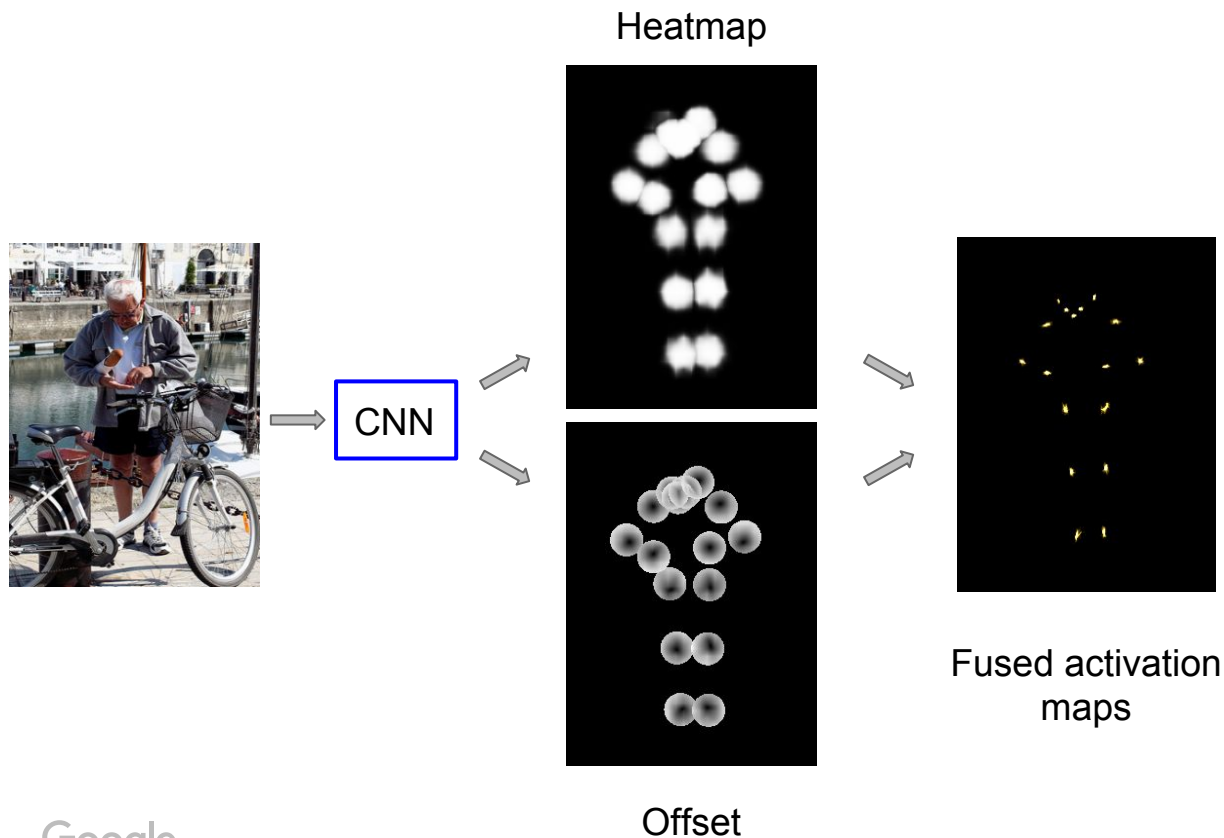
Target



Net-Layer101



Fusing Heatmap and Offset Outputs

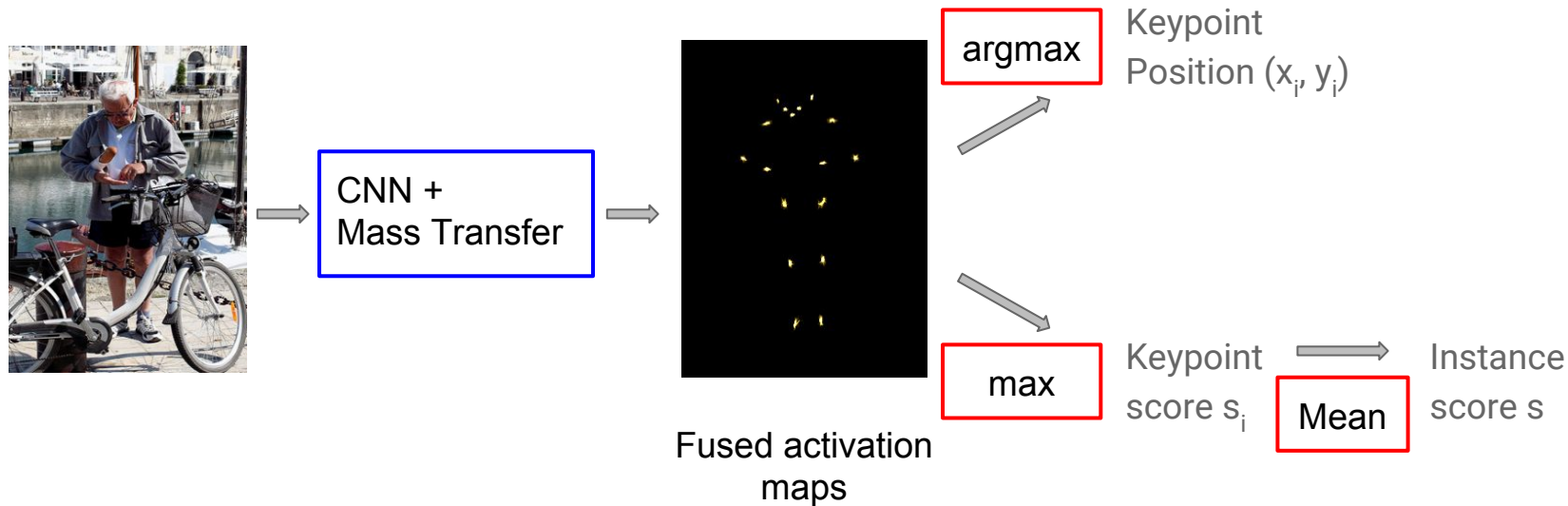


Algo: *Offset-guided mass transfer*

For each point in the heatmap:

- (1) Transfer its mass by the corresponding offset.
- (2) Accumulate into fused activation maps.

Final Pose Prediction: Keypoint Position and Score



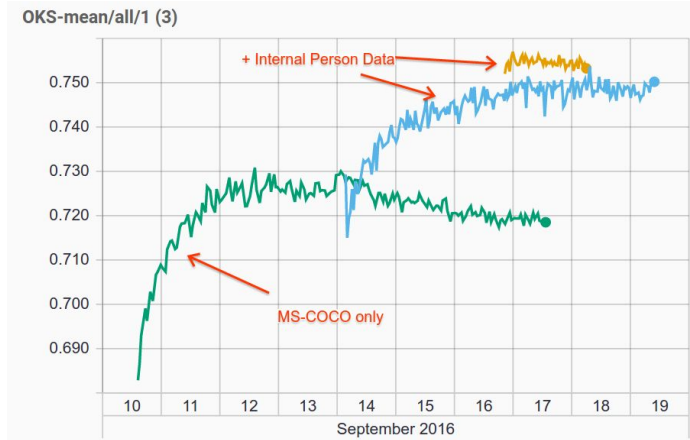
Pose rescoring is *crucial*:

0.05 AP boost compared to using Faster-RCNN box scores.

Pose Estimation-Only Results

COCO pose estimation with oracle boxes

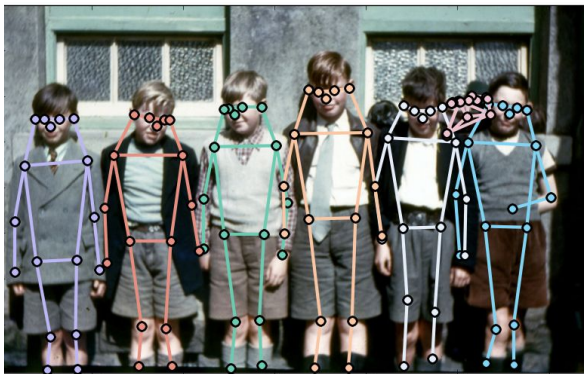
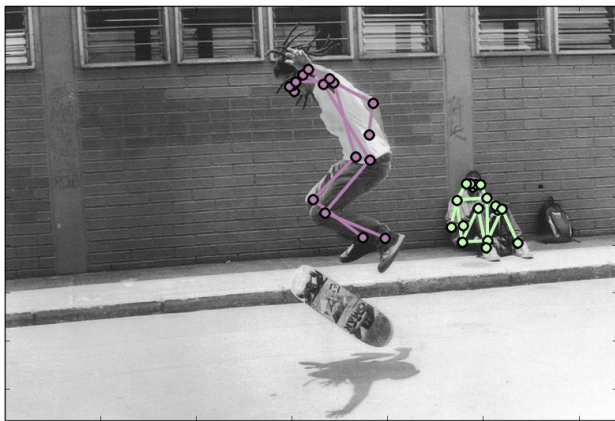
- How well can we predict the pose, given ground truth person boxes?
- Metric: AR in mini-val given ground truth box
- 0.730 using COCO only annotations.
- 0.756 also using in-house person keypoint annotations.



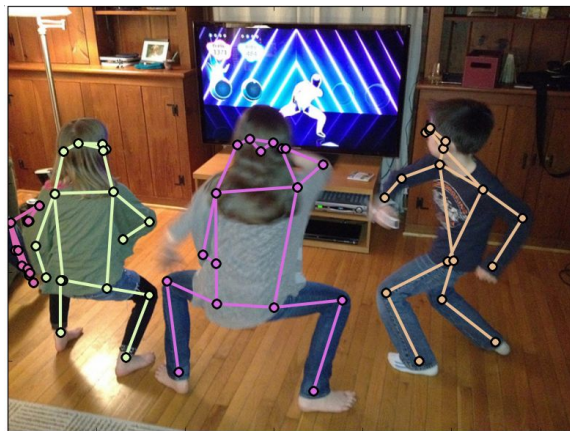
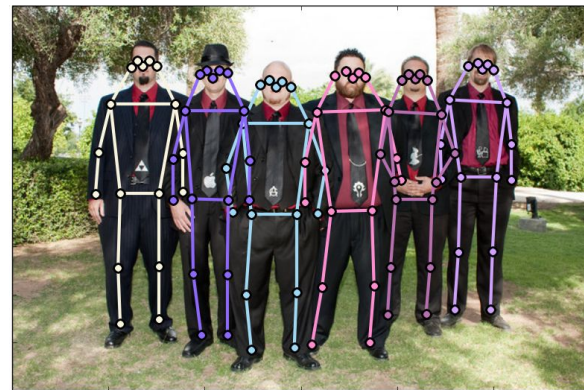
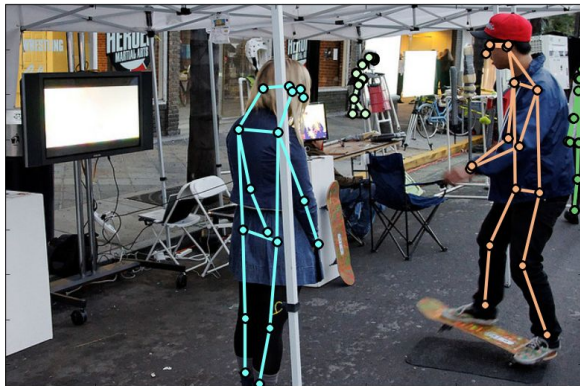
MPII Single-Person task

- Sanity check for our pose estimation network:
 - Training on MPII only data
 - 89.0% PCKh@0.5 (close to state-of-art)

Full System Results: COCO Images



Full System Results: COCO Images



COCO Quantitative Results (Competition)

- Competition results on “*Challenge*” split.

TEAM	AP	AP@.5	AP@.75	AP (M)	AP (L)	AR	AR@.5	AR@.75	AR (M)	AR (L)
CMU-Pose	0.605	0.834	0.664	0.551	0.681	0.659	0.864	0.713	0.594	0.748
G-RMI	0.598	0.81	0.651	0.567	0.667	0.664	0.865	0.712	0.618	0.726
DL-61	0.533	0.751	0.485	0.555	0.548	0.708	0.828	0.688	0.74	0.782
R4D	0.497	0.743	0.545	0.456	0.556	0.556	0.773	0.603	0.491	0.644
umich_vl7	0.434	0.722	0.449	0.364	0.534	0.499	0.758	0.52	0.387	0.652

COCO Quantitative Results (Post Competition)

- Latest results on “*testdev*” split.

TEAM	AP	AP@.5	AP@.75	AP (M)	AP (L)	AR	AR@.5	AR@.75	AR (M)	AR (L)
G-RMI (competition entry)	0.605	0.822	0.662	0.576	0.666	0.662	0.866	0.714	0.619	0.722
G-RMI (post competition)	0.668	0.863	0.734	0.630	0.733	0.716	0.896	0.776	0.669	0.782

- Bug discovered: Aspect ratio mismatch between train and eval code.
- Added OKS-based non-maximum suppression on the pose result.

COCO Quantitative Results (Post Competition)

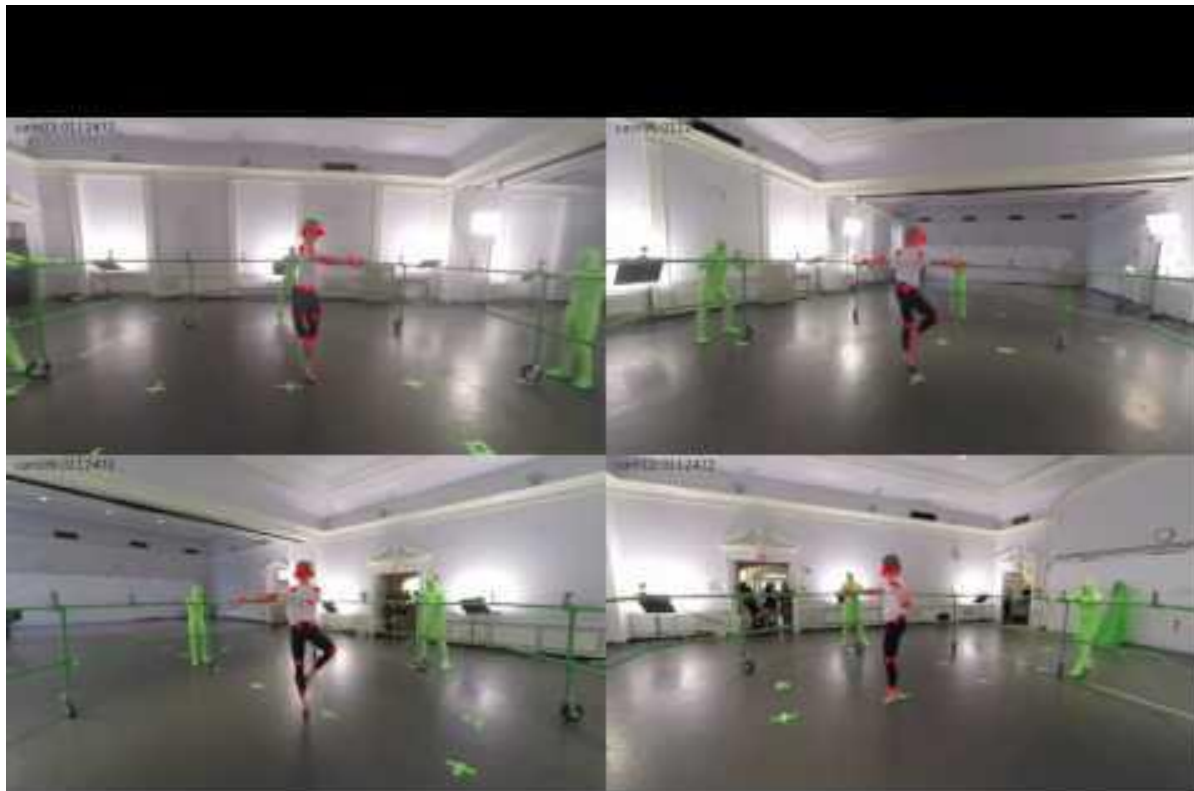
- Latest results on “*testdev*” split.
- AP using the G-RMI Object Detection team’s boxes trained only on COCO box annotations.
- Effect of training the pose estimator on COCO-only pose annotations vs. COCO+in-house pose annotations.
- Effect of OKS-based Non-Maximum Suppression.

AP (testdev)	No OKS-based NMS	With OKS-based NMS
COCO-only	0.601	0.636
COCO+in-house	0.628	0.668

Lessons from COCO Keypoints Challenge

- New dataset, new metric, first year the challenge runs
- Challenging problem:
 - Person detection and pose estimation
- Important differences compared to previous pose datasets:
 - Many example persons with severe occlusion
 - Large scale variations (and scale is not considered known)
- Lots of room for improvement!
- Interesting research problems:
 - Example: Two-stage or single-stage system?
- Pose estimation is becoming mature for in-the-wild deployment!

David Hallberg Dance Sequence Results



David Hallberg 3D Trajectory Reconstruction



Thanks!

- Person detection and pose estimation team
 - Tyler Zhu, Nori Kanazawa, George Papandreou, Alex Toshev, Hartwig Adam, Chris Bregler, Kevin Murphy, Jonathan Tompson
- G-RMI object detection team
 - Alireza Fathi, Ian Fischer, Sergio Guadarrama, Jonathan Huang, Anoop Korattikara, Kevin Murphy, Vivek Rathod, Yang Song, Chen Sun, Zbigniew Wojna, Menglong Zhu
- TF-slim
 - Sergio Guadarrama, Nathan Silberman
- Person annotation team
 - Akshay Gogia, Gursheesh Kour, Manish Arora
- Special thanks:
 - Georgia Gkioxari, Dumitru Erhan, Jitendra Malik, Chuck Rosenberg, Rahul Sukthankar, Jay Yagnik