

# MediaPipe BlazePose GHUM 3D



## MODEL DETAILS

Lite (3MB size), Full (6 MB size) and Heavy (26 MB size) models, to estimate the **full 3D body pose** of an individual in videos captured by a **smartphone or web camera**. Optimized for **on-device, real-time fitness applications**: Lite model runs ~44 FPS on a CPU via [XNNPack](#) TFLite and ~49 FPS via TFLite GPU on a Pixel 3. Full model runs ~18 FPS on a CPU via [XNNPack](#) TFLite and ~40 FPS via TFLite GPU on a Pixel 3. Heavy model runs ~4 FPS on a CPU via [XNNPack](#) TFLite and ~19 FPS via TFLite GPU on a Pixel 3.



Depth is encoded via gradient from blue (closer) to green (further). Invisible (occluded) keypoints marked as black.

Returns 33 keypoints describing the approximate location of body parts:

- Nose
- Right eye (3 keypoints): Inner, Center, Outer
- Left eye (3 keypoints): Inner, Center, Outer
- Ears (2 keypoints): Right, Left
- Mouth (2 keypoints): Right Corner, Left Corner
- Shoulder (2 keypoints): Right, Left
- Elbow (2 keypoints): Right, Left
- Wrist (2 keypoints): Right, Left
- Pinky knuckle (2 keypoints): Right, Left
- Index knuckle (2 keypoints): Right, Left
- Thumb knuckle (2 keypoints): Right, Left
- Hip (2 keypoints): Right, Left
- Knee (2 keypoints): Right, Left
- Ankle (2 keypoints): Right, Left
- Heels (2 keypoints): Right, Left
- Foot Index (2 keypoints): Right, Left

Keypoint Z-value estimate is provided using synthetic data, obtained via the [GHUM model](#) (articulated 3D human shape model) fitted to 2D point projections.



## MODEL SPECIFICATIONS

### Model Type

Convolutional Neural Network

### Model Architecture

Convolutional Neural Network: MobileNetV2-like with customized blocks for real-time performance.

### Input(s)

Regions in the video frames where a person has been detected. Represented as a 256x256x3 array with aligned human full body part, centered by mid-hip in vertical body pose and rotation distortion of (-10, 10). Channels order: RGB with values in [0.0, 1.0].

### Output(s)

33x5 array corresponding to (x, y, z, visibility, presence).

- X, Y coordinates are local to the region of interest and range from [0.0, 255.0].
- Z coordinate is measured in "image pixels" like the X and Y coordinates and represents the distance relative to the plane of the subject's hips, which is the origin of the Z axis. Negative values are between the hips and the camera; positive values are behind the hips. Z coordinate scale is similar with X, Y scales but has different nature as obtained not via human annotation, by fitting synthetic data ([GHUM model](#)) to the 2D annotation. Note, that Z is not metric but up to scale.
- Visibility is in the range of [min\_float, max\_float] and after user-applied sigmoid denotes the probability that a keypoint is located within the frame and not occluded by another bigger body part or another object.
- Presence is in the range of [min\_float, max\_float] and after user-applied sigmoid denotes the probability that a keypoint is located within the frame.



#### AUTHORS

##### Who created this model?

Valentin Bazarevsky, Google  
Ivan Grishchenko, Google  
Eduard Gabriel Bazavan, Google



#### CITATION

##### How can users cite your model?

BlazePose: On-device Real-time Body Pose tracking,  
CVPR Workshop on Computer Vision for Augmented  
and Virtual Reality,  
Seattle, WA, USA, 2020

#### DATE

April, 16, 2021

GHUM & GHUML: Generative 3D Human Shape and  
Articulated Pose Models

Proceedings of the IEEE/CVF Conference on  
Computer Vision and Pattern Recognition, pages  
6184-6193, 2020



#### DOCUMENTATION

- [BlazePose: On-device Real-time Body Pose tracking](#)
- [GHUM & GHUML: Generative 3D Human Shape and Articulated Pose Models](#)



#### LICENSED UNDER

[Apache License, Version 2.0](#)

## Intended Uses



#### APPLICATION

3D full body pose estimation for  
single-person videos on mobile,  
desktop and in browser.



#### DOMAIN AND USERS

- Augmented reality
- 3D Pose and gesture recognition
- Fitness and repetition counting
- 3D pose measurements (angles / distances)



#### OUT-OF-SCOPE APPLICATIONS

- Multiple people in an image.
- People too far away from the camera (e.g. further than 14 feet/4 meters)
- Head is not visible
- Applications requiring metric accurate depth
- Any form of surveillance or identity recognition is explicitly out of scope and not enabled by this technology

## Limitations



#### PRESENCE OF ATTRIBUTES

Tracks only one person on scene if  
multiple present



#### TRADE-OFFS

The model is optimized for  
real-time performance on a  
wide variety of mobile devices,  
but is sensitive to face position,  
scale and orientation in the  
input image.



#### ENVIRONMENT

When degrading the  
environment light, noise, motion  
or face overlapping conditions  
one can expect degradation of  
quality and increase of “jittering”  
(although we cover such cases  
during training with real-world  
samples and augmentations).

# Ethical Considerations



## HUMAN LIFE

The model is not intended for human life-critical decisions. The primary intended application is entertainment.



## PRIVACY

This model was trained and evaluated on images, including consented images (30K), of people using a mobile AR application captured with smartphone cameras in various “in-the-wild” conditions. The majority of training images (85K) capture a wide range of fitness poses.



## BIAS

This model was trained and evaluated both on visible and hidden points. For cases that the point location is present but hard to define by humans annotator, it is annotated with a “best guess” and default pose. Model has been qualitatively evaluated on users with missing limbs and prosthetics and degrades gracefully by predicting average point location.

The model is providing 3D coordinates, but the z-coordinate is up to scale (not metric) and obtained from synthetic data using the GHUM model (articulated 3D human shape model), fitted via an algorithm to the 2D key point projections.

# Training Factors and Subgroups



## INSTRUMENTATION

- All dataset images were captured on a diverse set of back-facing smartphone cameras.
- All images were captured in a real-world environment with different light, noise and motion conditions via an AR (Augmented Reality) application.



## ATTRIBUTES

- Human Full-body cropped from the captured frame should contain a single person placed in the center of the image.
- There should be a margin around the square circumscribing full-body calculated as 25% of size.
- Model is tolerant to certain level of input inaccuracy:
  - 10% shift and scale (taking body width/height as 100% for corresponding axis)
  - 8° roll



## ENVIRONMENTS

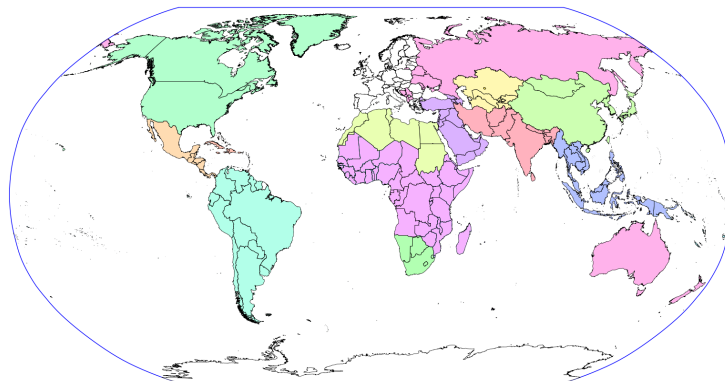
Model is trained on images with various lighting, noise and motion conditions and with diverse augmentations. However, its quality can degrade in extreme conditions. This may lead to increased “jittering” (inter-frame prediction noise).



## GROUPS

To perform fairness evaluation we group user samples into 14 evenly distributed geographic subregions (based on [United Nations geoscheme](#) with merges):

Central America	Caribbean
Southern America	Northern America
Central Asia	Northern Africa
Eastern Asia	Middle Africa
Southeastern Asia	Southern Africa
Southern Asia	Australia and New Zealand
Western Asia	Europe (excluding EU)



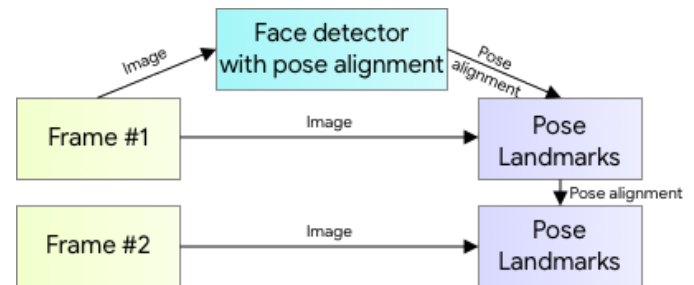
# Evaluation modes and metrics

## Evaluation Modes



### TRACKING MODE

Main mode that takes place most of the time and is based on obtaining a highly accurate full-body crop from the prediction on the previous frame (frames 2, 3, ... on the image)



## Model Performance Measures



PDJ, Average percentage of detected joints  
(Also known as PCK@0.2 - Percent of Correct Keypoints)

<https://github.com/cbsudux/Human-Pose-Estimation-101>

We consider a keypoint to be correctly detected if predicted visibility for it matches ground truth and the absolute 2D Euclidean error between the reference and target keypoint normalized by the 2D torso diameter projection is smaller than 20%. This value was determined during development as the maximum value that does not degrade accuracy in classifying pose / asana based solely on the key points without perceiving the original RGB image.

The model is providing 3D coordinates, but the z-coordinate is obtained from synthetic data, so for a fair comparison with human annotations, only 2D coordinates are employed.

# Evaluation results

## Geographical Evaluation Results



### DATA

- **Contains 1400 samples evenly distributed across 14 geographical subregions** (see specification in Section "Factors and Subgroups"). Each region contains 100 images.
- All samples are picked from the same source as training samples and are characterized as smartphone back-facing camera photos taken in real-world environments (see specification in "Factors and Subgroups - Instrumentation").



### EVALUATION RESULTS

Detailed evaluation for the tracking modes across 14 geographical subregions, gender and skin tones is presented in the table below

Region	Lite model		Full model		Heavy model	
	PDJ	Standard deviation	PDJ	Standard deviation	PDJ	Standard deviation
Australia and New Zealand	86.1	12.7	92.0	9.5	93.5	10.8
Caribbean	88.8	14.3	93.0	12.5	94.6	11.5
Europe	83.9	15.5	90.2	14.1	93.9	8.4
Northern Africa	89.0	11.9	93.0	12.0	93.8	8.7
South America	88.0	13.0	91.4	11.3	95.4	7.0
Southeastern Asia	87.7	14.7	91.5	13.2	94.1	12.3
Western Asia	88.8	10.9	93.8	7.4	95.5	6.9
Central America	87.5	13.4	92.9	8.8	95.2	6.3
Central Asia	85.6	15.2	90.7	12.7	93.1	10.8
Eastern Asia	83.2	15.0	90.4	10.4	92.6	9.0
Middle Africa	85.1	19.3	89.2	16.6	91.4	14.6
Northern America	88.3	10.8	91.3	9.1	95.4	8.5
Southern Africa	89.7	14.8	94.0	8.8	93.6	11.9
Southern Asia	86.6	15.3	91.2	12.4	96.2	7.2
Average	87.0		91.8		94.2	
Range	6.5		4.8		4.8	

## Geographical Fairness Evaluation Results



### FAIRNESS CRITERIA

We consider a model to be performing unfairly across representative groups if the error range on them spans more than ~3x the human annotation discrepancy, in our case a total of **7.5% PDJ**.



### FAIRNESS METRICS & BASELINE

We asked two annotators to re-annotate the Pose Validation dataset, yielding a PDJ of **97.5%**  
This is a high inter-annotator agreement, suggesting that the PDJ metric is a strong indicator of precise matches between predicted keypoints and ground truth keypoints.



### FAIRNESS RESULTS

Evaluation across 14 regions of heavy, full and lite models on smartphone back-facing camera photos dataset results an average performance of 94.2% +/- 1.3% stdev with a range of [91.4%, 96.2%] across regions for the heavy model, an average performance of 91.8% +/- 1.4% stdev with a range of [89.2%, 94.0%] across regions for the full model and an average performance of 87.0% +/- 2.0% stdev with a range of [83.2%, 89.7%] across regions for the lite model.

Comparison with our fairness criteria yields a maximum discrepancy between average and worst performing regions of 4.8% for the heavy, 4.8% for the full and 6.5% for the light model.

## Skin Tone and Gender Evaluation Results



### DATA

- **1400 images, 100 images from each of 14 the geographical subregions** were annotated with perceived gender and skin tone (from 1 to 6) based on the Fitzpatrick scale.



### EVALUATION RESULTS

Evaluation on smartphone back-facing camera photos dataset results in an average performance of 93.6% with a range of [89.3%, 95.0%] across all skin tones for the heavy model, an average performance of 91.1% with a range of [85.9%, 92.9%] across all skin tones for the full model and an average performance of 86.2% with a range of [80.5%, 87.8%] across regions for the lite model. The maximum discrepancy between worst and best performing categories is 5.7% for the heavy model, 7.0% for the full model and 7.3% for the lite model.

Evaluation across gender yields an average performance of 94.8% with a range of [94.2%, 95.3%] for the heavy model, an average performance of 92.3% with a range of [91.2%, 93.4%] for the full model, and an average of 87.6% with a range of [86.0%, 89.1%] for the lite model. The maximum discrepancy is 1.1% for the heavy model, 2.2% for the full model and 3.1% for the lite model.

Skin tone type	% of dataset	Lite model		Full model		Heavy model	
		PDJ	Standard deviation	PDJ	Standard deviation	PDJ	Standard deviation
1	1.3	80.5	10.9	85.9	11.1	89.3	12.4
2	9.5	85.7	14.3	91.4	10.0	94.1	8.0
3	34.3	87.5	13.5	92.4	11.2	94.7	8.7
4	36.2	87.8	12.2	92.3	10.1	95.0	7.6
5	14.2	87.7	13.2	91.9	10.6	94.6	8.8
6	4.5	87.7	15.5	92.9	8.9	93.6	11.2
Average		86.2		91.1		93.6	
Range		7.3		7.0		5.7	

Gender	% of dataset	Lite model		Full model		Heavy model	
		PDJ	Standard deviation	PDJ	Standard deviation	PDJ	Standard deviation
Male	43.5	89.1	12.4	93.4	9.5	95.3	7.8
Female	56.5	86.0	13.6	91.2	11.3	94.2	9.0
Average		87.6		92.3		94.8	
Range		3.1		2.2		1.1	

# Definitions

## AUGMENTED REALITY (AR)

Augmented reality, a technology that superimposes a computer-generated image on a user's view of the real world, thus providing a composite view.

## PRESENCE

Presence denotes the probability that a keypoint is located within the frame. It does not indicate whether the keypoint is occluded by another body part.

## KEYPOINTS

"Keypoints" or "landmarks" are (x, y, z) coordinate locations of body parts.

## VISIBILITY

Visibility denotes the probability that a keypoint is located within the frame and not occluded either by other body parts or other objects.