# 4.MediaPipe Development

mediapipe github：https://github.com/google/mediapipe

mediapipe official website： https://google.github.io/mediapipe/

dlib official website： http://dlib.net/

dlib github： https://github.com/davisking/dlib

# 4.Introduction

MediaPipe is an open source data stream processing machine learning application development framework developed by Google. It is a graph-based data processing pipeline for building data sources using many forms, such as video, audio, sensor data, and any time series data. MediaPipe is cross-platform and can run on embedded platforms (Raspberry Pi, etc.), mobile devices (iOS and Android), workstations and servers, and supports mobile GPU acceleration. MediaPipe provides cross-platform, customizable ML solutions for live and streaming media.

The core framework of MediaPipe is implemented in C++ and provides support for languages such as Java and Objective C. The main concepts of MediaPipe include Packet, Stream, Calculator, Graph and Subgraph.

Features of MediaPipe：

- End-to-end acceleration: Built-in fast ML inference and processing accelerates even on commodity hardware.
- Build once, deploy anywhere: Unified solution for Android, iOS, desktop/cloud, web and IoT.
- Ready-to-use solutions: Cutting-edge ML solutions that showcase the full capabilities of the framework.
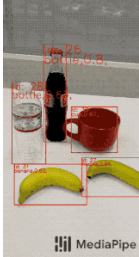- Free and open source: framework and solutions under Apache2.0, fully scalable and customizable.

Deep Learning Solutions in MediaPipe

| Face Detection | Face Mesh | Iris | Hands | Pose | Holistic |
|---|---|---|---|---|---|
|  |  |  |  |  |  |

| Hair Segmentation | Object Detection | Box Tracking | Instant Motion Tracking | Objectron | KNIFT |
|---|---|---|---|---|---|
|  |  |  |  |  |  |

|  | Android | iOS | C++ | Python | JS | Coral |
|---|---|---|---|---|---|---|
| Face Detection | ☑ | ☑ | ☑ | ☑ | ☑ | ☑ |
| Face Mesh | ☑ | ☑ | ☑ | ☑ | ☑ | |
| Iris | ☑ | ☑ | ☑ | | | |
| Hands | ☑ | ☑ | ☑ | ☑ | ☑ | |
| Pose | ☑ | ☑ | ☑ | ☑ | ☑ | |
| Holistic | ☑ | ☑ | ☑ | ☑ | ☑ | |
| Selfie Segmentation | ☑ | ☑ | ☑ | ☑ | ☑ | |
| Hair Segmentation | ☑ | | ☑ | | | |
| Object Detection | ☑ | ☑ | ☑ | | | ☑ |
| Box Tracking | ☑ | ☑ | ☑ | | | |
| Instant Motion Tracking | ☑ | | | | | |
| Objectron | ☑ | | ☑ | ☑ | ☑ | |
| KNIFT | ☑ | | | | | |
| AutoFlip | | | ☑ | | | |
| MediaSequence | | | ☑ | | | |
| YouTube 8M | | | ☑ | | | |

## 4.2. Use

```
cd ~/yahboomcar_ws/src/yahboomcar_mediapipe/scripts      # Enter the directory
where the source code is located
python3 06_FaceLandmarks.py                              # face effects
python3 07_FaceDetection.py                              # Face Detection
python3 08_Objectron.py                                  # 3D object recognition
python3 09_VirtualPaint.py                               # Paintbrush
python3 10_HandCtrl.py                                   # finger control
python3 11_GestureRecognition.py                         # Gesture Recognition
```

During use, you need to pay attention to the following:

- All functions [q key] are for exit.
- 3D object recognition: Recognizable objects are: ['Shoe', 'Chair', 'Cup', 'Camera'], a total of 4 categories; click the [f key] to switch to recognize objects; jetson series cannot use keyboard keys to switch recognition The object needs to change the [self.index] parameter in the source code.
- Paintbrush: When the index finger and middle finger of the right hand are combined, it is in the selection state, and the color selection box will pop up at the same time. When the two fingertips move to the corresponding color position, the color will be selected (black is the eraser); Draw freely on.
- Finger control: Click [f key] to switch the recognition effect.

- Gesture recognition: Gesture recognition designed for the right hand can be accurately recognized when certain conditions are met. Recognizable gestures are: [Zero, One, Two, Three, Four, Five, Six, Seven, Eight, Ok, Rock, Thumb_up (like), Thumb_down (thumb down), Heart_single (one-handed comparison)] , a total of 14 categories. |
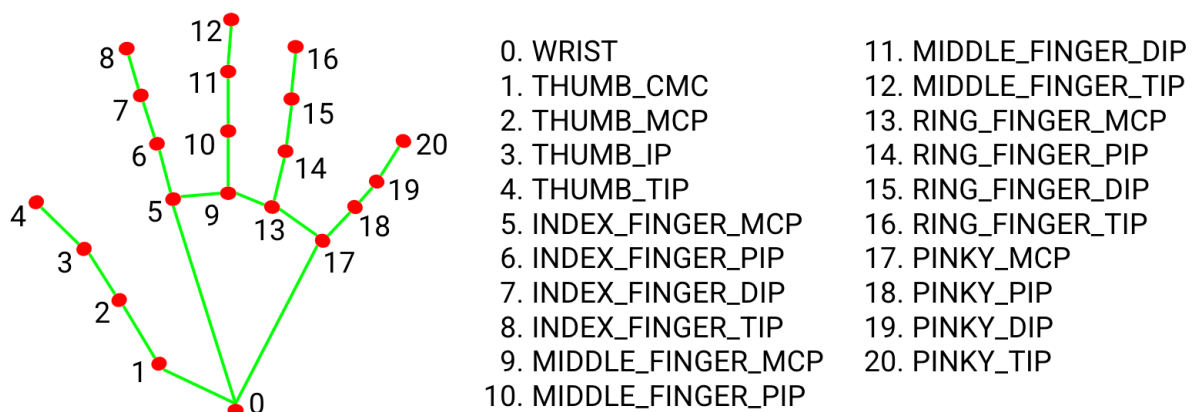
| :-------------: | ---------------------------------------------------------- |

| 06.face effects | 07. Face detection |
| :-- | :-- |
| | |
| 08. Three-dimensional object recognition | 09. Paintbrush |
| | |

| 10. Finger control | |
| :--: | :-- |
| 11. Gesture recognition | |

## 4.3.MediaPipe Hands

MediaPipe Hands is a high fidelity hand and finger tracking solution. It uses machine learning (ML) to infer the 3D coordinates of 21 hands from one frame.

After palm detection on the entire image, the 21 3D hand joint coordinates in the detected hand area are accurately positioned by regression according to the hand marking model, that is, direct coordinate prediction. The model learns a consistent internal hand pose representation that is robust even to partially visible hands and self-occlusions.
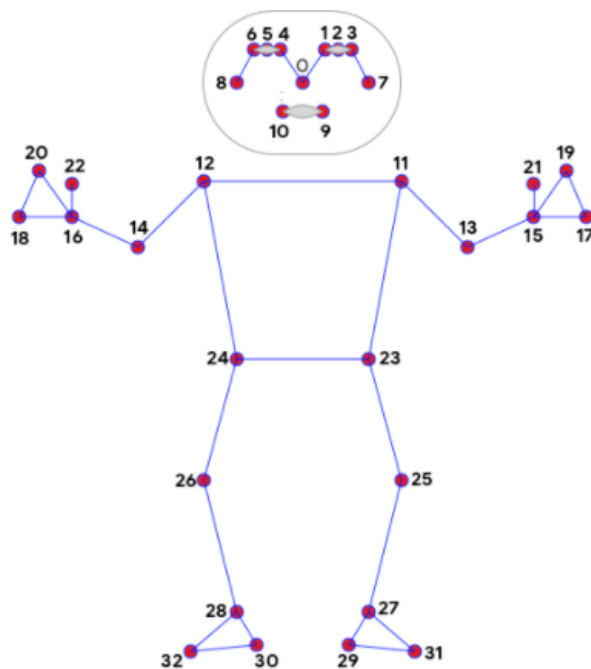
In order to obtain ground truth data, about 30K real-world images were manually annotated with 21 3D coordinates, as shown below (get the Z value from the image depth map, if there is a Z value for each corresponding coordinate). To better cover possible hand poses and provide additional supervision over the nature of the hand geometry, high-quality synthetic hand models in various backgrounds are also drawn and mapped to corresponding 3D coordinates.



```
0. WRIST                    11. MIDDLE_FINGER_DIP
1. THUMB_CMC                12. MIDDLE_FINGER_TIP
2. THUMB_MCP                13. RING_FINGER_MCP
3. THUMB_IP                 14. RING_FINGER_PIP
4. THUMB_TIP                15. RING_FINGER_DIP
5. INDEX_FINGER_MCP         16. RING_FINGER_TIP
6. INDEX_FINGER_PIP         17. PINKY_MCP
7. INDEX_FINGER_DIP         18. PINKY_PIP
8. INDEX_FINGER_TIP         19. PINKY_DIP
9. MIDDLE_FINGER_MCP        20. PINKY_TIP
10. MIDDLE_FINGER_PIP
```

## 4.4.MediaPipe Pose

MediaPipe Pose, an ML solution for high-fidelity body pose tracking, leverages BlazePose research to infer 33 3D coordinates and a full-body background segmentation mask from RGB video frames, which also powers the ML Kit pose detection API.

The landmark model in MediaPipe poses predicts the location of 33 pose coordinates (see image below).

| | |
|---|---|
| 0. nose | 17. left_pinky |
| 1. left_eye_inner | 18. right_pinky |
| 2. left_eye | 19. left_index |
| 3. left_eye_outer | 20. right_index |
| 4. right_eye_inner | 21. left_thumb |
| 5. right_eye | 22. right_thumb |
| 6. right_eye_outer | 23. left_hip |
| 7. left_ear | 24. right_hip |
| 8. right_ear | 25. left_knee |
| 9. mouth_left | 26. right_knee |
| 10. mouth_right | 27. left_ankle |
| 11. left_shoulder | 28. right_ankle |
| 12. right_shoulder | 29. left_heel |
| 13. left_elbow | 30. right_heel |
| 14. right_elbow | 31. left_foot_index |
| 15. left_wrist | 32. right_foot_index |
| 16. right_wrist | |

## 4.5.dlib

The corresponding case is face special effects.

DLIB is a modern C++ toolkit containing machine learning algorithms and tools for creating complex software in C++ to solve real-world problems. It is widely used by industry and academia in fields such as robotics, embedded devices, mobile phones, and large-scale high-performance computing environments.

The dlib library uses 68 points to mark important parts of the face, such as 18-22 points to mark the right eyebrow, and 51-68 to mark the mouth. Use the get_frontal_face_detector module of the dlib library to detect the face, and use the shape_predictor_68_face_landmarks.dat feature data to predict the face feature value