

5. MediaPipe development

5. MediaPipe development

5.1. Introduction

5.2. Use

5.3. MediaPipe Hands

5.4. MediaPipe Pose

5.5. dlib

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>

5.1. Introduction

MediaPipe is a data stream processing machine learning application development framework developed and open sourced by Google. It is a graph-based data processing pipeline for building data sources that use 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.

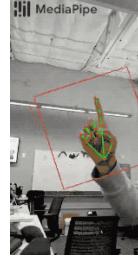
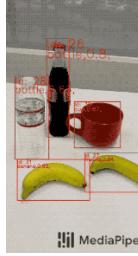
The core framework of MediaPipe is implemented in C++ and supports languages such as Java and Objective C. The main concepts of MediaPipe include Packet, Stream, Calculator, Graph and Subgraph.

Features of MediaPipe:

- End-to-end acceleration: Built-in fast ML inference and processing accelerates even on common hardware.
- Build Once, Deploy Anywhere: Unified solution for Android, iOS, desktop/cloud, web and IoT.
- Ready-to-use solutions: Cutting-edge ML solutions that demonstrate the full capabilities of the framework.
- Free and open source: frameworks and solutions under Apache2.0, fully extensible and customizable.

Deep Learning Solutions in MediaPipe

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

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			✓			

	Android	iOS	C++	Python	JS	Coral
MediaSequence			<input checked="" type="checkbox"/>			
YouTube 8M			<input checked="" type="checkbox"/>			

5.2. Use

----- ROS -----

```
roslaunch yahboomcar_mediapipe cloud_viewer.launch          # Point cloud
view: support 01~04
roslaunch yahboomcar_mediapipe 01_HandDetector.launch      # hand detection
roslaunch yahboomcar_mediapipe 02_PoseDetector.launch      # pose detection
roslaunch yahboomcar_mediapipe 03_Holistic.launch          # overall
detection
roslaunch yahboomcar_mediapipe 04_FaceMesh.launch #         face detection
roslaunch yahboomcar_mediapipe 05_FaceEyeDetection.launch # face
recognition
```

If using a monocular camera or a Raspberry PI CSI camera, change the following file names:

01_HandDetector_usb.py、02_PoseDetector_usb.py、03_Holistic_usb.py、
 04_FaceMesh_usb.py、05_FaceEyeDetection_usb.py--> 01_HandDetector.py、
 02_PoseDetector.py、03_Holistic.py、04_FaceMesh.py、05_FaceEyeDetection.py

If using a jetson CSI camera, change the following file names:

01_HandDetector_csi.py、02_PoseDetector_csi.py、03_Holistic_csi.py、
 04_FaceMesh_csi.py、05_FaceEyeDetection_csi.py-->
 01_HandDetector.py、02_PoseDetector.py、03_Holistic.py、04_FaceMesh.py、
 05_FaceEyeDetection.py

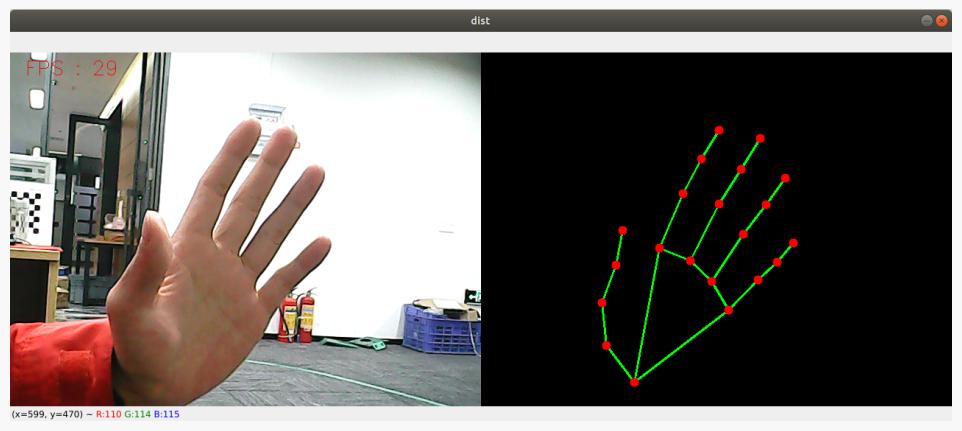
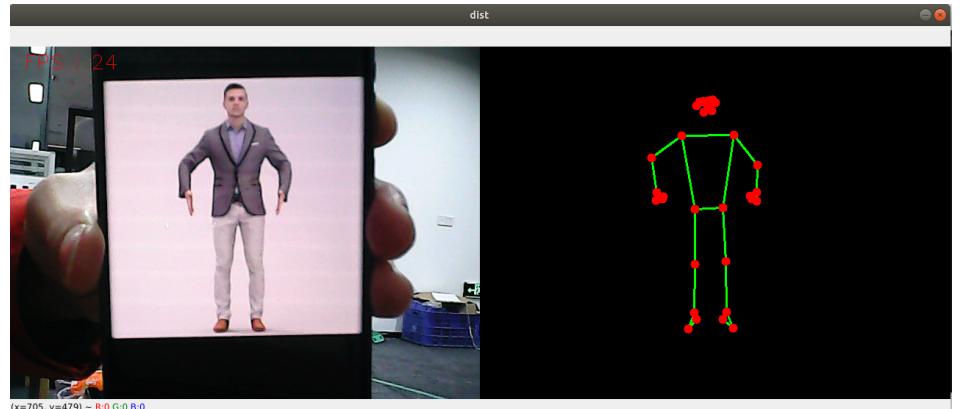
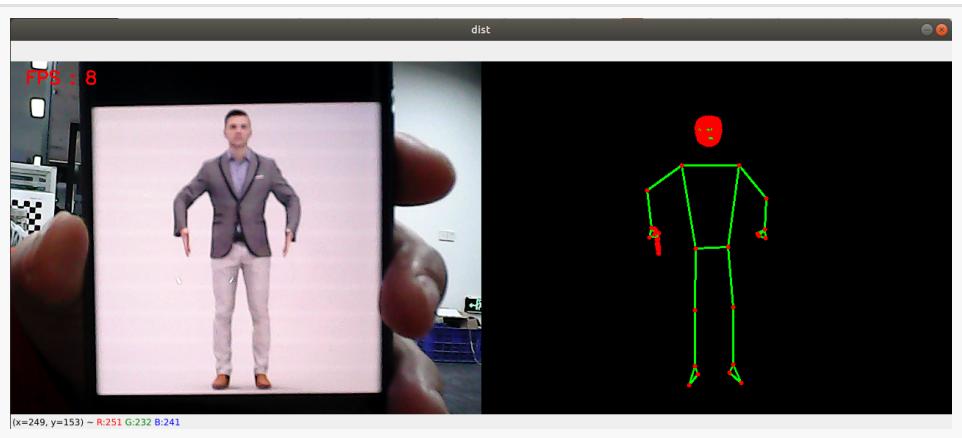
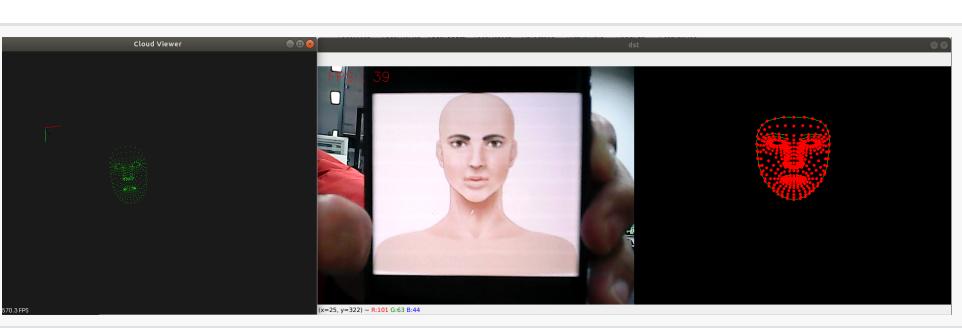
----- not ROS -----

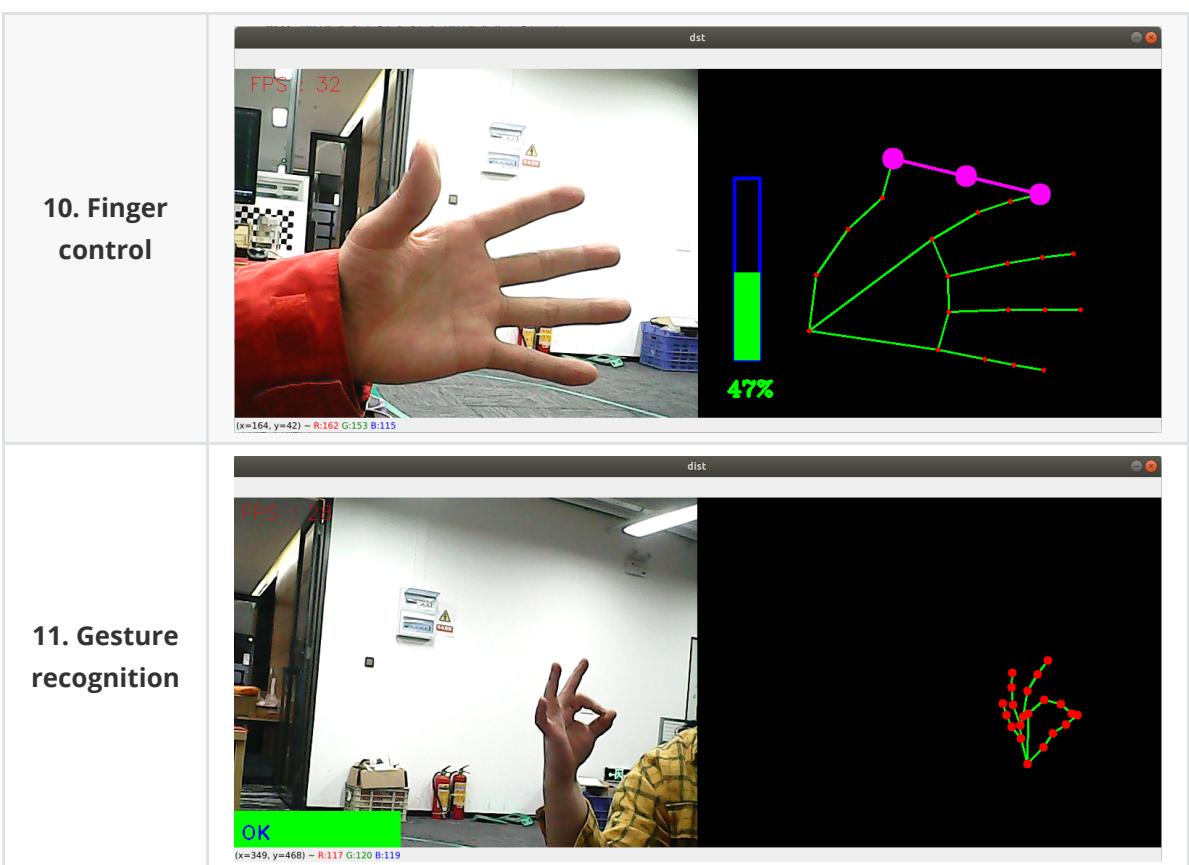
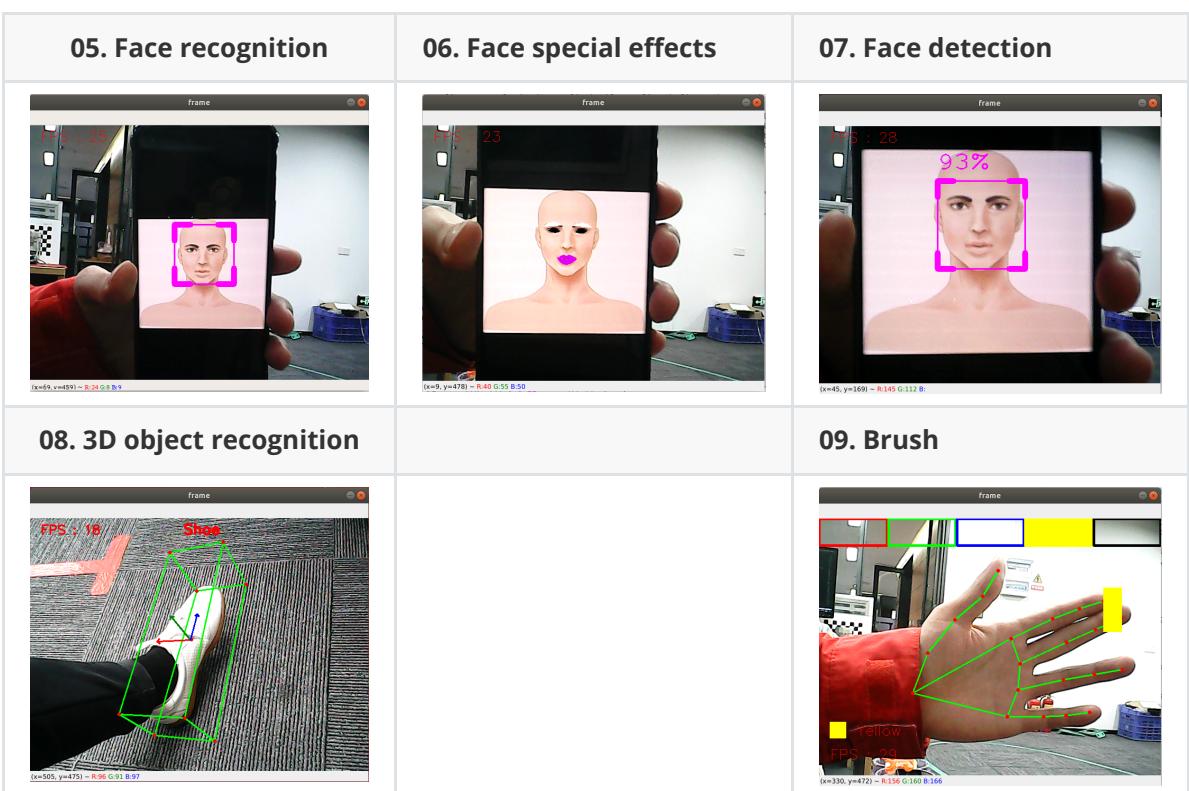
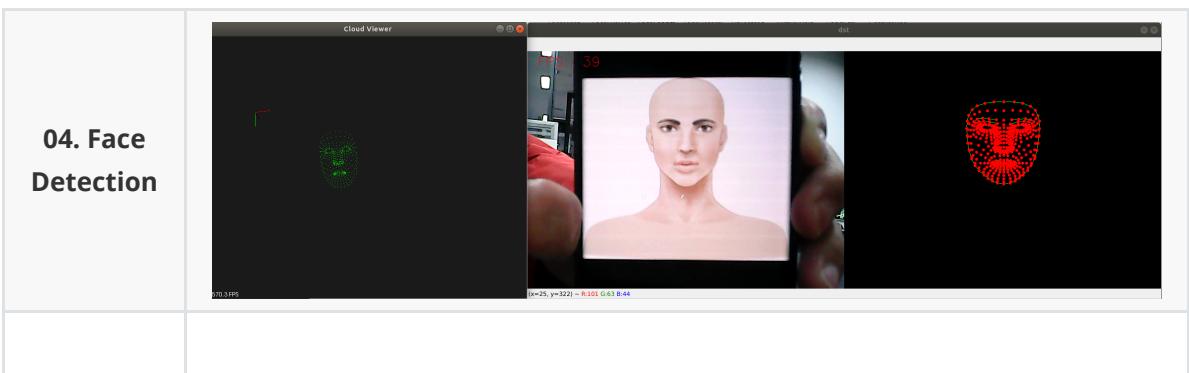
```
cd ~/yahboomcar_ws/src/yahboomcar_mediapipe/scripts      # Enter the directory
where the source code is located
python3 06_FaceLandmarks.py      # face effects
python3 07_FaceDetection.py     # face detection
python3 08_Objectron.py        # 3D object recognition
python3 09_VirtualPaint.py    # brushes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

In the process of use, it should be noted that

- Hand detection, posture detection, overall detection, and face detection all have point cloud viewing functions. Take face detection as an example.
- All functions 【q key】 to exit.
- Overall detection: including hand, face, body pose detection.
- 3D Object Recognition: Recognized objects are: ['Shoe', 'Chair', 'Cup', 'Camera'], a total of 4 categories; click [f key] to switch the recognition object; jetson series cannot use keyboard keys to switch recognition The object needs to change the [self.index] parameter in the source code.

- Brush: When the index finger and middle finger of the right hand are combined, it is in the selected state, and the color selection box will pop up at the same time. When the two fingertips move to the corresponding color position, select the color (black is the eraser); the index finger and the middle finger start to be in the drawing state, which can be displayed on the drawing board. Draw arbitrarily.
- Finger control: Click [f key] to switch the recognition effect.
- Finger recognition: gesture recognition designed with the right hand as the criterion, can be accurately recognized when certain conditions are met. Recognized gestures are: [Zero, One, Two, Three, Four, Five, Six, Seven, Eight, Ok, Rock, Thumb_up (like), Thumb_down (thumb down), Heart_single (one-handed heart)] , a total of 14 categories.

01. Hand detection	 <p>This interface shows hand detection. The left window displays a live video feed of a person's hand, with the FPS (Frames Per Second) value shown as 29. The right window shows a skeleton diagram of the hand with red dots at joints and green lines connecting them.</p>
02. Attitude detection	 <p>This interface shows attitude detection. The left window displays a live video feed of a person standing, with the FPS value shown as 24. The right window shows a skeleton diagram of the person with red dots at joints and green lines connecting them.</p>
03. Overall inspection	 <p>This interface shows overall inspection. The left window displays a live video feed of a person standing, with the FPS value shown as 8. The right window shows a skeleton diagram of the person with red dots at joints and green lines connecting them.</p>
04. Face Detection	 <p>This interface shows face detection. The left window displays a live video feed of a face, with the FPS value shown as 39. The middle window shows a 3D mesh of the face. The right window shows a skeleton diagram of the face with red dots at joints and green lines connecting them.</p>

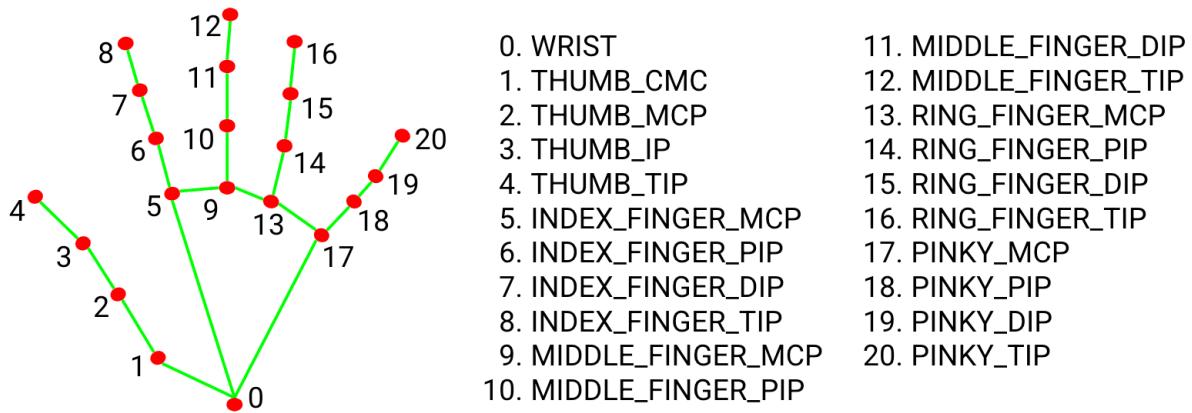


5.3. MediaPipe Hands

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 a single frame.

After palm detection is performed on the entire image, accurate key point positioning is performed on the 21 3D hand joint coordinates in the detected hand region by regression according to the hand marker model, that is, direct coordinate prediction. The model learns consistent internal hand pose representations and is robust to even partially visible hands and self-occlusion.

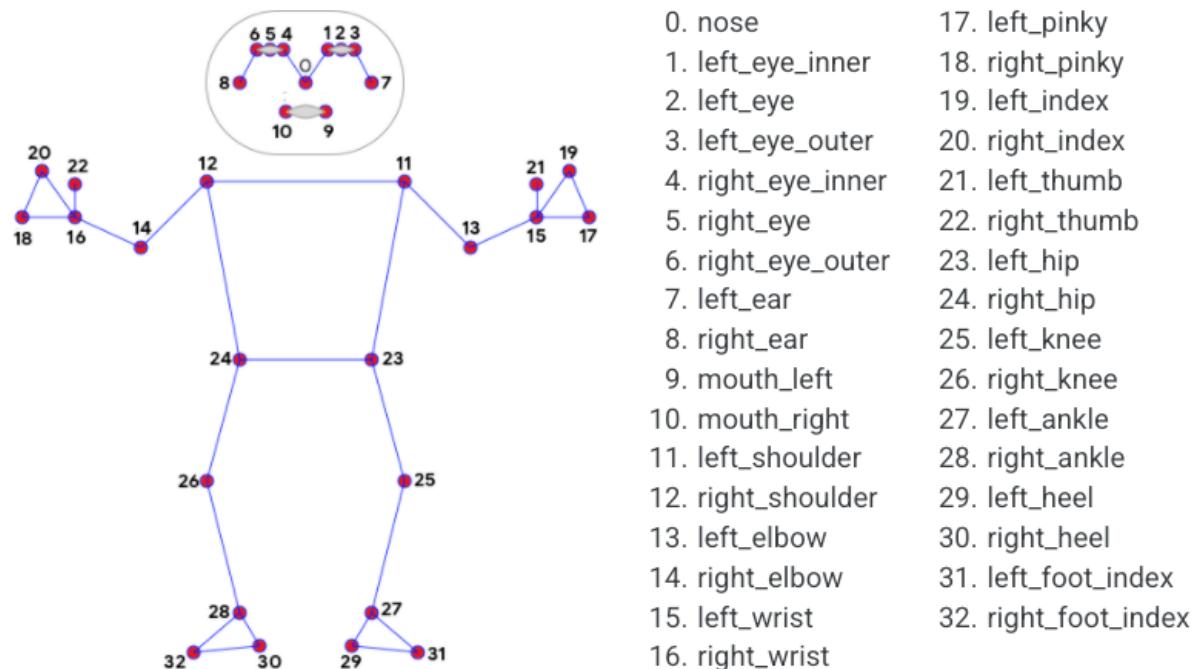
To obtain ground truth data, about 30K real-world images were manually annotated with 21 3D coordinates as shown below (Z-values are obtained from the image depth map, if there is a Z-value for each corresponding coordinate). To better cover the possible hand poses and provide additional supervision on the nature of the hand geometry, high-quality synthetic hand models in various contexts are also drawn and mapped to the corresponding 3D coordinates.



5.4. MediaPipe Pose

MediaPipe Pose is an ML solution for high-fidelity body pose tracking that infers 33 3D coordinates and a full-body background segmentation mask from RGB video frames using the BlazePose research that also powers the ML Kit pose detection API.

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



5.5. dlib

The corresponding case is the face effect.

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 points 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

