

Wrocław University of Science and Technology

ARTIFICIAL INTELLIGENCE AND COMPUTER VISION

Faculty of Electronics, Photonics and Microsystems
AIR DRAWING CALCULATOR (AI AND CV)

Theme of class: Final Report

Students:

1. Stanislav Kustov 275512

Date of class: 2025-01-30

Submission Date: 2025-01-30

Lab assistant: Mateusz Cholewiski

Contents

1	Introduction	2
1.1	Purpose and Applications	2
1.2	How It Works	2
1.3	Technologies Used	2
1.4	Overview	2
2	Assumptions	3
2.1	Functional Assumptions	3
2.2	Design Assumptions	3
2.3	Hardware Components	3
2.4	Software Assumptions	3
3	Description of Software	4
3.1	Hand Tracking	4
3.2	Fingertip Tracking	6
3.3	Air Drawing and Erasing	8
3.4	Digit and Operator Recognition	9
3.5	Equation Processing and Solving	10
3.5.1	Performing the Calculation	11
3.5.2	Examples of Different Group Counts	11
3.5.3	Visual Indicators for Operators	12
3.5.4	Summary	13
3.6	Training the Digit Recognition Model: Code, Methods, Results, and Discussion	14
3.6.1	Model Architecture and Training Code	14
3.6.2	Training Methods and Improvements	17
3.6.3	Results and Interpretation	18
3.6.4	Conclusion	19
3.7	Training the Operator Recognition Model: Code, Methods, Results, and Discussion	20
3.7.1	Model Architecture and Training Code	20
3.7.2	Results and Interpretation	23
3.7.3	Conclusion	24
4	Conclusion	25

1 Introduction

The **Air Drawing Calculator** is an interactive application that combines **Computer Vision (CV)** and **Artificial Intelligence (AI)** to recognize and process handwritten mathematical expressions in real-time. Using **hand tracking, gesture recognition, and deep learning**, this system allows users to draw numbers and mathematical operators in the air and instantly compute results.

1.1 Purpose and Applications

The purpose of this project is to develop an interactive application that combines hand tracking, gesture recognition, and optical character recognition (OCR) to create a virtual environment for drawing and recognizing handwritten symbols and equations.

1.2 How It Works

The system consists of two key components:

- **Hand Tracking & Gesture Recognition:**
 - Uses **Mediapipe** to track the user's hand.
 - The fingertip acts as a digital pen for writing numbers and operators.
 - A palm gesture allows erasing content effortlessly.
- **AI-Based Symbol Recognition:**
 - A **CNN model** classifies handwritten digits (0-9).
 - A separate model detects operators (+, -, *, /).
 - Small unintended marks are ignored, while deliberate dots are recognized as decimal points.

1.3 Technologies Used

This project is built using:

- **OpenCV** For image processing and video capture.
- **Mediapipe** For hand tracking and gesture detection.
- **TensorFlow/Keras** For deep learning-based handwriting recognition.
- **NumPy** For numerical operations and data processing.

1.4 Overview

This report covers:

- Implementation of hand tracking and gesture recognition.
- AI models for handwritten digit and operator recognition.
- Testing, accuracy improvements, and future enhancements.

By combining AI with gesture-based input, the Air Drawing Calculator offers a new, interactive way to solve mathematical problems without the need for a physical interface.

2 Assumptions

2.1 Functional Assumptions

- The user will write mathematical expressions clearly and legibly.
- The system assumes that only one equation or a set of numbers is written at a time.
- Small accidental marks will be ignored unless they resemble meaningful digits or operators.
- Hand gestures will be used for writing, erasing, and confirming calculations.
- The system will automatically detect and process digits, operators, and decimal points.

2.2 Design Assumptions

- The application is designed for real-time use with minimal processing delay.
- Users will interact with the system via hand tracking and gesture recognition.
- The drawing area is a virtual canvas displayed over the camera feed.
- The system prioritizes single-line input but supports multi-line numbers.
- Operators and digits are written separately to ensure accurate recognition.

2.3 Hardware Components

- A camera with a reasonable resolution (e.g., 720p or higher) is required for accurate hand tracking.
- A computer with a CPU or GPU capable of running deep learning models in real-time.
- Adequate lighting conditions are assumed to improve gesture recognition accuracy.
- Users should maintain a steady hand while writing to enhance recognition reliability.

2.4 Software Assumptions

- The system is built using Python and requires OpenCV, Mediapipe, and TensorFlow/Keras.
- The machine learning models are pre-trained and loaded during execution.
- The software runs on a local machine and does not require an internet connection.
- The program assumes that all libraries and dependencies are installed and configured properly.
- The application will run smoothly on modern operating systems (Windows, Linux, macOS).

3 Description of Software

3.1 Hand Tracking

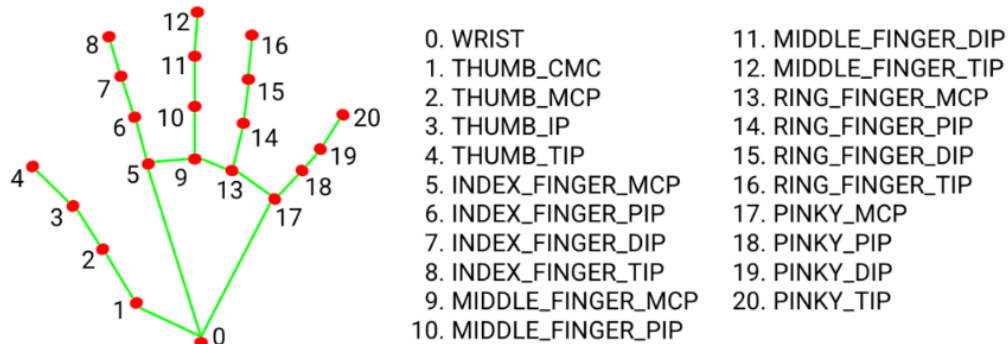
The development of the Air Calculator began with the implementation of real-time hand tracking and gesture recognition using **Mediapipe**, an advanced framework for detecting and analyzing hand landmarks.

The **HandLandmarker** module in Mediapipe accurately identifies the position of the user's hands and fingertips, establishing the foundation for gesture-based interactions such as drawing and erasing.



The implementation was based on the official Mediapipe documentation and GitHub repository. Below is a simplified example demonstrating how Mediapipe detects hand landmarks in an image:

The hand landmark model bundle detects the keypoint localization of 21 hand-knuckle coordinates within the detected hand regions. The model was trained on approximately 30K real-world images, as well as several rendered synthetic hand models imposed over various backgrounds.



The hand landmarker model bundle contains a palm detection model and a hand landmarks detection model. The Palm detection model locates hands within the input image, and the hand landmarks detection model identifies specific hand landmarks on the cropped hand image defined by the palm detection model.

```
1 # Step 1: Import necessary modules
2 import mediapipe as mp
3 from mediapipe.tasks import python
4 from mediapipe.tasks.python import vision
5
6 # Step 2: Create a HandLandmarker object
7 base_options = python.BaseOptions(model_asset_path='hand_landmarker.task')
8 options = vision.HandLandmarkerOptions(base_options=base_options, num_hands=2)
9 detector = vision.HandLandmarker.create_from_options(options)
10
11 # Step 3: Load an input image
12 image = mp.Image.create_from_file("image.jpg")
13
14 # Step 4: Detect hand landmarks
15 detection_result = detector.detect(image)
16
17 # Step 5: Visualize the detected hand landmarks
18 annotated_image = draw_landmarks_on_image(numpy_view(), detection_result)
19 cv2.imshow(cv2.cvtColor(annotated_image, cv2.COLOR_RGB2BGR))
```

This code initializes the Mediapipe hand tracking system, loads an image, and detects hand landmarks, which can then be used for further processing.

3.2 Fingertip Tracking

With hand tracking established, the next phase involved real-time tracking of the **fingertip position**, which is essential for the air drawing functionality. By focusing on the index fingertip (landmark 8), the system can precisely capture the users drawing movements. Additionally, a **thumb gesture** serves as a control mechanism, allowing users to switch between drawing and idle modes effortlessly.

```
1 fingertip_x = hand_landmarks.landmark[8].x * frame.shape[1]
2 fingertip_y = hand_landmarks.landmark[8].y * frame.shape[0]
3 fingertip_position = (int(fingertip_x), int(fingertip_y))
4
5 if is_thumb_extended(hand_landmarks):
6     drawing_mode = True
7 else:
8     drawing_mode = False
9
10 if drawing_mode:
11     if previous_position is not None:
12         cv2.line(canvas, previous_position, fingertip_position, (255, 255, 255), 5)
13         previous_position = fingertip_position
14 else:
15     previous_position = None
```

Step-by-Step Explanation:

- **Fingertip Position Tracking:** The index fingertip (landmark 8) is detected and converted into pixel coordinates:

```
fingertip_x = hand_landmarks.landmark[8].x * frame.shape[1]
fingertip_y = hand_landmarks.landmark[8].y * frame.shape[0]
fingertip_position = (int(fingertip_x), int(fingertip_y))
```

- **Thumb Gesture Detection:** The system determines whether the thumb is extended to activate drawing mode:

```
if is_thumb_extended(hand_landmarks):
    drawing_mode = True
else:
    drawing_mode = False
```

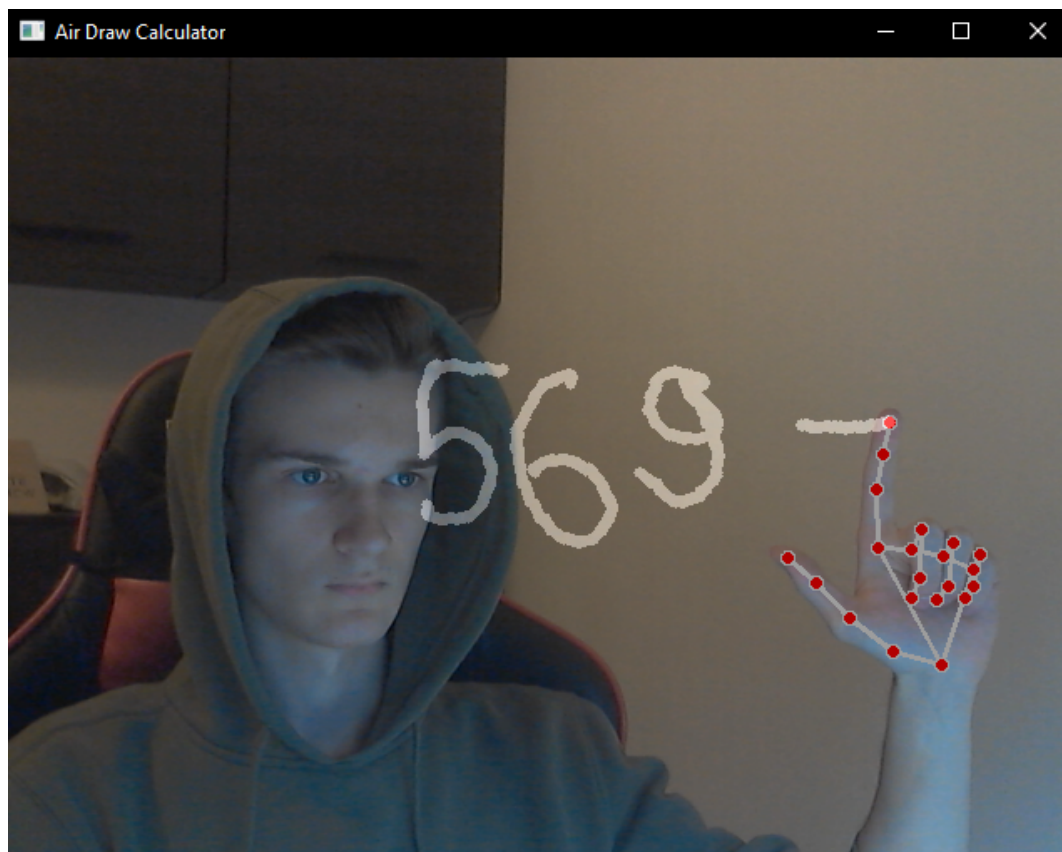
- **Drawing Mode:** When drawing is enabled, lines are drawn between consecutive fingertip positions:

```
if drawing_mode:
    if previous_position is not None:
        cv2.line(canvas, previous_position, fingertip_position, (255, 255, 255), 5)
        previous_position = fingertip_position
    else:
        previous_position = None
```

- **Canvas Initialization:** A black canvas matching the frame dimensions is used as the drawing surface:

```
if canvas is None:
    canvas = np.zeros_like(frame)
```

This implementation allows users to **draw lines in real-time** by moving their fingertip while extending the thumb. When the thumb is not extended, drawing is disabled.



3.3 Air Drawing and Erasing

Once the fingertip position is tracked, the system allows users to draw in the air. The path of the fingertip is visualized as white strokes on a black canvas.

Drawing Mode:

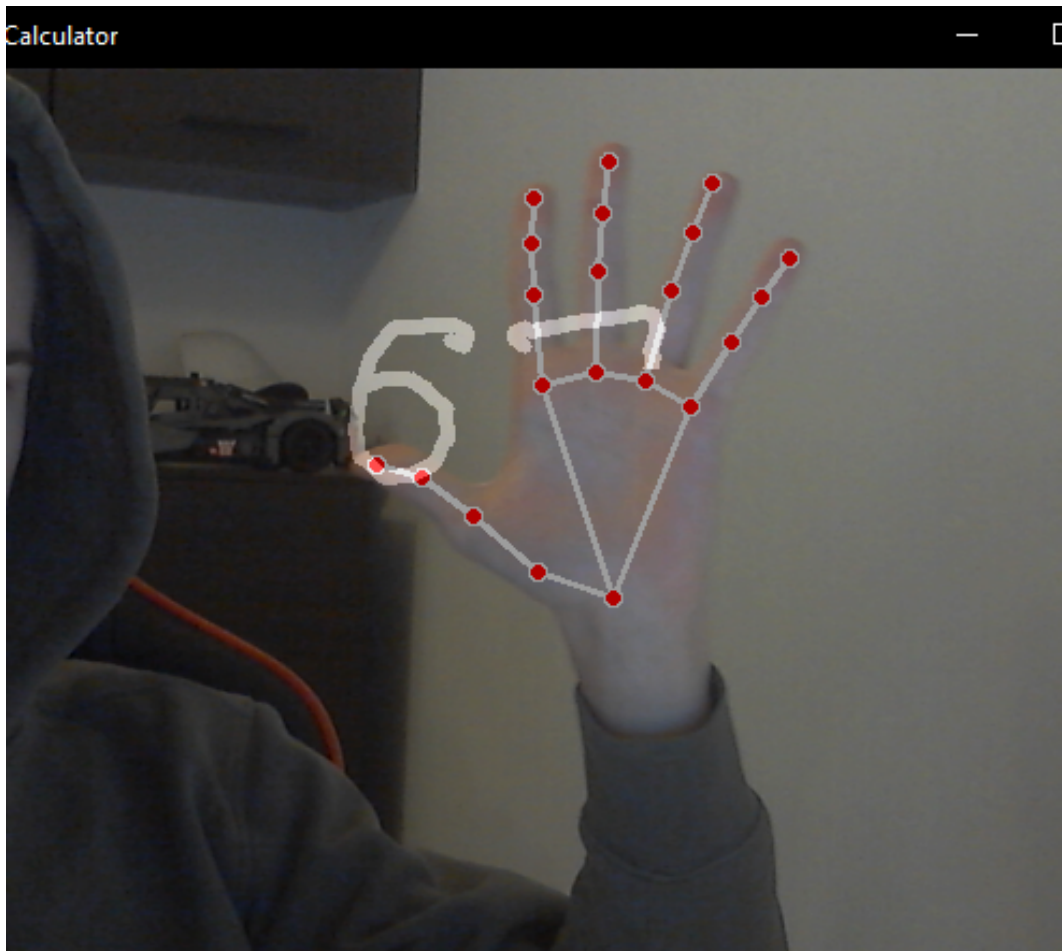
- The user extends their thumb to activate drawing mode.
- Movement of the index fingertip is recorded and displayed on the screen.

Erasing Mode:

- If the user's palm is detected, an erasing operation is triggered.
- A black rectangle removes drawn elements.
- If the palm detection fails, the system erases based on finger position.

Erasing Implementation:

```
57 if is_hand_open(hand_landmarks):  
58     erase_area(hand_landmarks, canvas, frame.shape)
```



3.4 Digit and Operator Recognition

When the user finishes writing an expression, the system segments the written input into separate components:

- **Digits (0-9):** Recognized using a CNN trained on handwritten digits.
- **Operators (+, -, *, /):** Recognized using a separate CNN for mathematical symbols.

Digit Recognition:

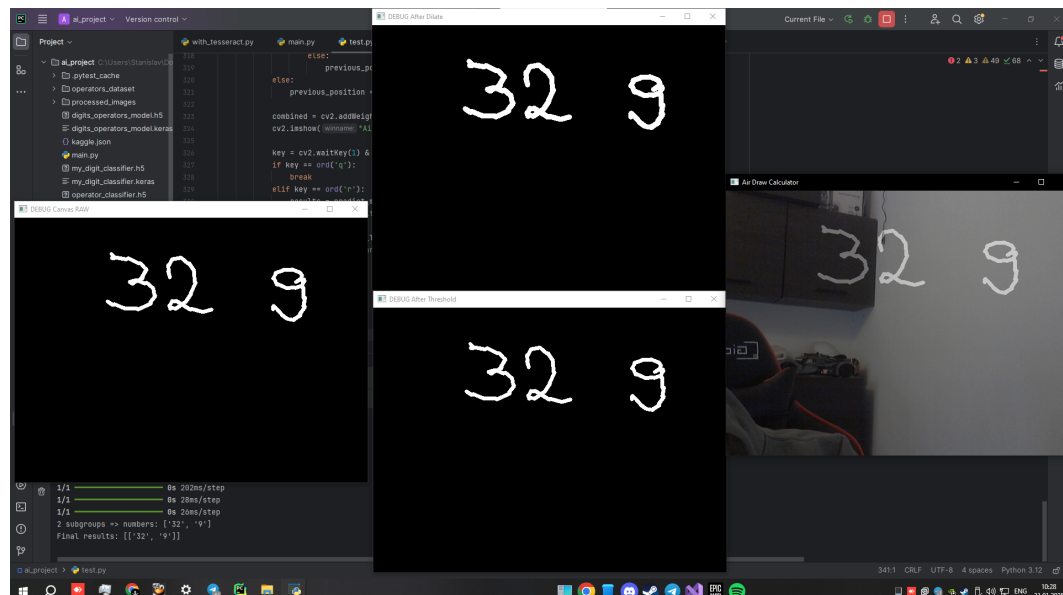
- The digit is extracted and resized to **28x28** pixels.
- The image is normalized and passed into the trained CNN model.
- The output is a classification result in the range **0-9**.

Operator Recognition:

- The extracted operator is inverted (black to white).
- It is resized and normalized before being classified by the CNN model.
- The output is mapped to one of the four basic operators.

Digit Prediction Code:

```
82 final_28 = cropped.astype('float32') / 255.0
83 input_data = final_28.reshape(1, 28, 28, 1)
84 preds = digits_model.predict(input_data, verbose=1)
85 digit = np.argmax(preds)
```



3.5 Equation Processing and Solving

When the user presses the **R** key, the program runs a function called `predict_all_in_one_line(canvas)` that does the following:

1. **Preprocessing the Canvas:**

- Converts the black-and-white “drawing” (the canvas) to grayscale.
- Applies a binary threshold (e.g., value of 127) so that only solid white contours remain.
- Performs a small dilation to ensure the contours are connected clearly.
- Uses `cv2.findContours` to detect all white patches (each patch is a contour).

2. **Extracting Bounding Boxes:**

- For each contour, the program computes a bounding box (x, y, w, h) .
- These boxes are stored in a list.

3. **Grouping Bounding Boxes (Sublines):**

- The list of boxes is sorted from left to right based on the **x**-coordinate.
- The function `split_line_into_groups` checks the *horizontal gap* between consecutive boxes. For two boxes with coordinates (x, y, w, h) and (px, py, pw, ph) , the gap is

$$\text{gap} = x - (px + pw).$$

- If $\text{gap} \leq \text{space_x_thresh}$ (default = 15 pixels), the boxes are merged into the same group, which usually indicates consecutive parts of a single number (like “1” next to “5” forming “15”).
- If $\text{gap} > 15$, the program starts a new group, indicating a larger separation (for example, a gap between a digit and an operator).

4. **Determining if it is an Equation:**

- The program then counts how many **groups** were formed.
- If there are exactly 3 groups, the code assumes the user wrote:

(Group 1: Number) (Group 2: Operator) (Group 3: Number).

- If there are more or fewer than 3 groups, the program simply treats them as a *list of recognized items* (most likely just digits or decimal points) and prints them separated by semicolons without trying to do any computation.

Why three groups means “Equation”: By design, a basic arithmetic expression like “3 + 5” or “15 * 24” naturally splits into **three** groups:

- *Left operand* (Group 1) – e.g., “3” or “15”.
- *Operator* (Group 2) – e.g., “+” or “*”.
- *Right operand* (Group 3) – e.g., “5” or “24”.

If the user drew more symbols (e.g., “3”, “2”, “15”, “46”) or fewer symbols (e.g., “15”, “24”), the system does not recognize a standard “number – operator – number” pattern.

3.5.1 Performing the Calculation

Once the program determines that there are exactly three groups, it takes the following steps to compute the result. Below is the relevant portion of the code:

```
110 if len(sublines) == 3:
111     # Convert recognized strings to float
112     left_val = float(left_str)
113     right_val = float(right_str)
114
115     # Operator check
116     if op_str == '+':
117         res = left_val + right_val
118     elif op_str == '-':
119         res = left_val - right_val
120     elif op_str == '*':
121         res = left_val * right_val
122     elif op_str == '/':
123         res = left_val / right_val if right_val != 0 else "ERR(div0)"
```

- **String to Number Conversion:** The left and right groups (e.g., "3" or "15.2") are transformed into floating-point values.
- **Recognizing the Operator:** The middle group is passed to the operator model, which returns '+', '-', '*', or '/'.
- **Applying Arithmetic:** Depending on the operator, the program uses a basic `if-elif` block to do addition, subtraction, multiplication, or division.
- **Division by Zero:** If `right_val` is zero and the operator is '/', the program returns "ERR(div0)".
- **Rounding to an Integer (if needed):** After the calculation, if the result is very close to an integer (for example, 8.000000001), the code converts it to an `int` to avoid printing 8.0.

3.5.2 Examples of Different Group Counts

Example 1: Three Groups, Equation ("3 + 5")

- First group is recognized as "3" → converts to 3.0.
- Second group is recognized as "+".
- Third group is recognized as "5" → converts to 5.0.
- The system calculates $3.0 + 5.0 = 8.0$ and then rounds it to 8.
- It prints: "3 + 5 = 8".

Example 2: Four Groups (“3”, “2”, “15”, “46”)

- The user draws four distinct symbols far enough apart that each becomes its own group.
- There is no operator recognized in a separate middle group.
- The system outputs them as: "3; 2; 15; 46" with no attempt at arithmetic, because `len(sublines) = 4`.

Example 3: Two Groups (“15”, “24”)

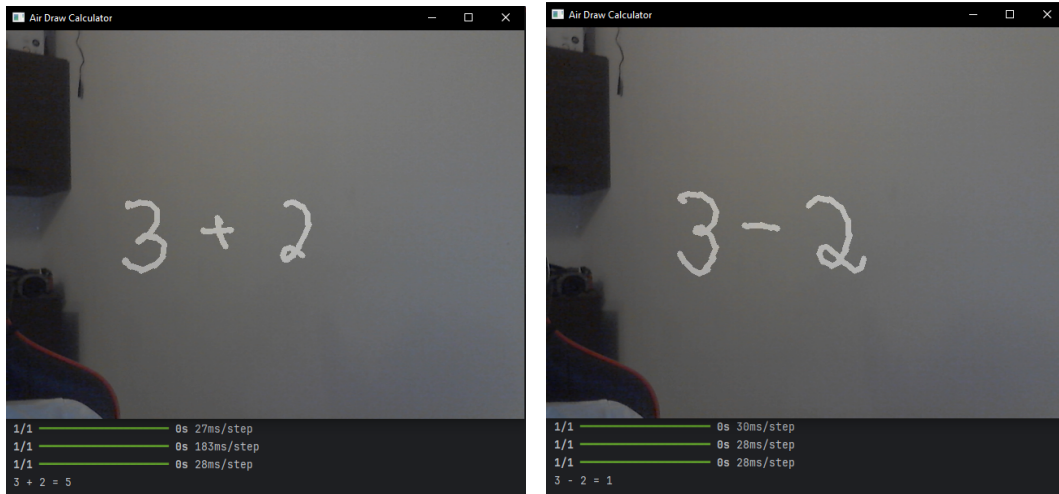
- Again, no operator group exists.
- The program will simply output: "15; 24".

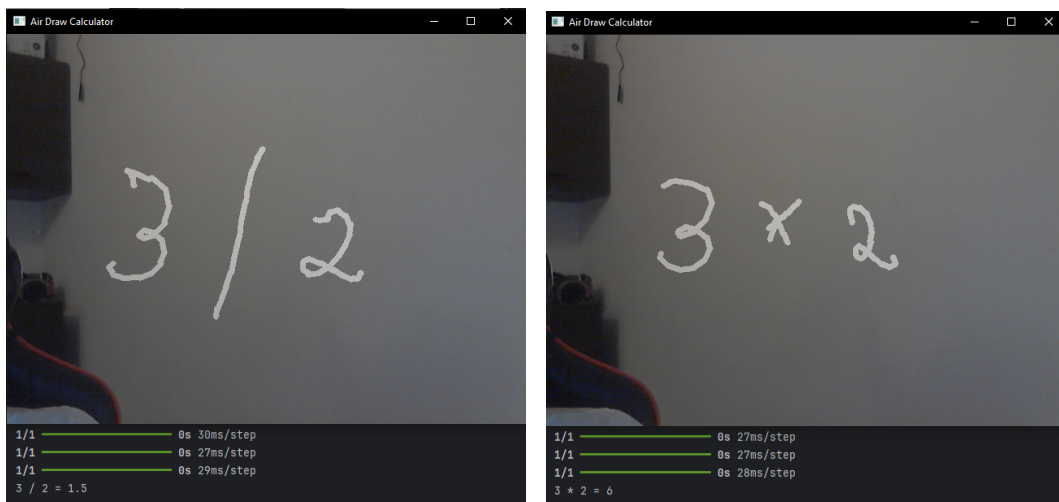
Example 4: Three Groups with an Operator (“15 * 24”)

- Group 1: "15" → `float(15.0)`
- Group 2: "*" → `*`
- Group 3: "24" → `float(24.0)`
- The result is $15 * 24 = 360$, so the console prints "15 * 24 = 360".

3.5.3 Visual Indicators for Operators

Below are example images (plus, minus, divide, times) that the model recognizes:





3.5.4 Summary

In summary, the program classifies written symbols into digits (possibly forming multi-digit numbers) and exactly one operator if three groups are formed. It then converts the digit strings to floats, applies the recognized operator, and handles errors such as division by zero. If there are fewer or more groups, the program does not perform any calculation, instead simply reporting the recognized digits as separate items.

3.6 Training the Digit Recognition Model: Code, Methods, Results, and Discussion

Overview: In this section, I describe how I trained a Convolutional Neural Network (CNN) on the MNIST dataset. I explain the key methods I used to boost performance, including data augmentation, specific CNN blocks, and additional layers that improved generalization. I also show the training and validation accuracy over 20 epochs and interpret the final results.

3.6.1 Model Architecture and Training Code

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import tensorflow as tf
4 from tensorflow.keras import layers, models
5 from tensorflow.keras.datasets import mnist
6 from tensorflow.keras.preprocessing.image import ImageDataGenerator
7
8 # 1. Load MNIST
9 (x_train, y_train), (x_test, y_test) = mnist.load_data()
10
11 x_train = x_train.astype('float32') / 255.0
12 x_test = x_test.astype('float32') / 255.0
13
14 # Add a channel dimension: from (28,28) to (28,28,1)
15 x_train = np.expand_dims(x_train, axis=-1)
16 x_test = np.expand_dims(x_test, axis=-1)
17
18 num_classes = 10
19 y_train = tf.keras.utils.to_categorical(y_train, num_classes)
20 y_test = tf.keras.utils.to_categorical(y_test, num_classes)
21
22 # 2. Data Augmentation
23 datagen = ImageDataGenerator(
24     rotation_range=25,
25     width_shift_range=0.2,
26     height_shift_range=0.2,
27     zoom_range=0.2,
28     shear_range=10.0,
29     fill_mode='constant',
30     cval=0.0
31 )
32 datagen.fit(x_train)
33
34 # 3. Define the CNN Model
35 model = models.Sequential()
36
37 # Block 1
38 model.add(layers.Conv2D(32, (3,3), padding='same', input_shape=(28,28,1)))
39 model.add(layers.BatchNormalization())
40 model.add(layers.Activation('relu'))
```

```
41
42 model.add(layers.Conv2D(32, (3,3), padding='same'))
43 model.add(layers.BatchNormalization())
44 model.add(layers.Activation('relu'))
45 model.add(layers.MaxPooling2D((2,2))) # Down to 14x14
46
47 # Block 2
48 model.add(layers.Conv2D(64, (3,3), padding='same'))
49 model.add(layers.BatchNormalization())
50 model.add(layers.Activation('relu'))
51
52 model.add(layers.Conv2D(64, (3,3), padding='same'))
53 model.add(layers.BatchNormalization())
54 model.add(layers.Activation('relu'))
55 model.add(layers.MaxPooling2D((2,2))) # Down to 7x7
56
57 # Block 3
58 model.add(layers.Conv2D(128, (3,3), padding='same'))
59 model.add(layers.BatchNormalization())
60 model.add(layers.Activation('relu'))
61
62 model.add(layers.Conv2D(128, (3,3), padding='same'))
63 model.add(layers.BatchNormalization())
64 model.add(layers.Activation('relu'))
65 model.add(layers.MaxPooling2D((2,2))) # Down to 3x3
66
67 # Dense layers
68 model.add(layers.Flatten())
69 model.add(layers.Dense(256))
70 model.add(layers.BatchNormalization())
71 model.add(layers.Activation('relu'))
72 model.add(layers.Dropout(0.4))
73
74 model.add(layers.Dense(128))
75 model.add(layers.BatchNormalization())
76 model.add(layers.Activation('relu'))
77 model.add(layers.Dropout(0.4))
78
79 model.add(layers.Dense(num_classes, activation='softmax'))
80
81 model.compile(
82     optimizer='adam',
83     loss='categorical_crossentropy',
84     metrics=['accuracy']
85 )
86
87 model.summary()
88
89 # 4. Training
90 batch_size = 64
91 epochs = 20
92 history = model.fit(
```



```
93     datagen.flow(x_train, y_train, batch_size=batch_size),
94     steps_per_epoch=len(x_train)//batch_size,
95     epochs=epochs,
96     validation_data=(x_test, y_test)
97 )
98
99 # 5. Plot Accuracy
100 plt.plot(history.history['accuracy'], label='Train acc')
101 plt.plot(history.history['val_accuracy'], label='Val acc')
102 plt.legend()
103 plt.title("Training vs. Validation Accuracy")
104 plt.show()
105
106 # 6. Evaluate on Test Set
107 test_loss, test_acc = model.evaluate(x_test, y_test, verbose=0)
108 print("Accuracy on test:", test_acc)
109
110 # 7. Save the Model
111 model.save("super_digit_classifier.keras", save_format="keras")
```

Key Points:

- **Data Augmentation:** I use rotations, shifts, zooming, and shearing to mimic a wide range of handwriting styles and random distortions.
- **CNN Blocks:** Each “block” contains two Conv2D layers and a MaxPooling2D layer. BatchNormalization follows each convolution to stabilize gradients, while ReLU provides non-linearity.
- **Layer Separation:** After three convolution blocks, I transition from the learned feature maps to fully connected layers (Dense layers). This separation is key: the lower blocks focus on extracting spatial patterns, while the dense layers specialize in the final classification.
- **Dropout & BatchNormalization in Dense Layers:** Including Dropout(0.4) in two successive dense layers helps prevent overfitting by randomly “dropping out” neurons. The BatchNormalization layers help reduce internal covariate shift, accelerating convergence.
- **Final Layer:** A softmax layer with 10 outputs corresponds to the ten MNIST digit classes.

3.6.2 Training Methods and Improvements

1. Data Augmentation with ImageDataGenerator:

- *Rotation (up to 25 degrees)*: Creates tilted versions of digits, forcing the network to learn rotation invariance.
- *Width/Height Shifts (up to 20%)*: Moves digits around the canvas, helping the model detect digits in varied positions.
- *Zoom and Shear*: Simulates zoomed-in or slanted digits, improving robustness against shape distortions.

2. Multiple Convolution Blocks:

- *Block 1 (32 filters), Block 2 (64 filters), Block 3 (128 filters)*: Each block increases the network depth, allowing it to learn progressively abstract features (edges, corners, strokes, etc.).
- *MaxPooling2D*: Reduces the spatial dimension after each pair of Conv2D layers, speeding up training and preventing oversized feature maps.

3. Separated Dense Layers:

- After flattening, I feed the features into two Dense layers (256 and 128 units). This separation lets me apply dropout and further normalization specifically in the classification part of the network.
- A dropout rate of 0.4 is relatively high but proves effective for managing overfitting on MNIST.

4. Adam Optimizer: I train using the Adam optimizer, which adapts the learning rate automatically. This optimizer usually converges quickly on MNIST.

3.6.3 Results and Interpretation

```
1 Epoch 1/20
2 937/937 | 47s 46ms/step - accuracy: 0.7315 - val_accuracy: 0.9515
3 Epoch 2/20
4 937/937 | 43s 47ms/step - accuracy: 0.9375 - val_accuracy: 0.9555
5 ...
6 Epoch 15/20
7 937/937 | accuracy: 0.9811 - val_accuracy: 0.9931
8 Epoch 20/20
9 937/937 | accuracy: 0.9844 - val_accuracy: 0.9872
10
11 Accuracy on test: 0.9872000217437744
```

Training Log Snippet

- The model starts around 73% training accuracy during the first epoch, surging to over 95% on the validation set by its end.
- By later epochs, the training accuracy occasionally hits 100% on some passes, while validation remains mostly in the 95–99% range, sometimes exceeding 99%.
- The final reported test accuracy is **98.72%**, a strong result for MNIST, especially considering the extra distortions introduced by augmentation.

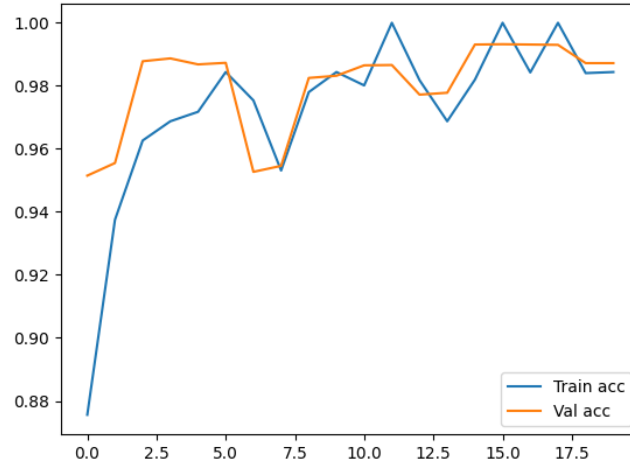


Figure 3: Training vs. Validation Accuracy over 20 epochs

Accuracy Plot Analysis

- **Rapid Early Improvement:** The model quickly learns essential digit shapes, rising from $\sim 88\%$ to $\sim 96\%$ in just a few epochs.
- **Fluctuations from Augmentation:** I observe occasional dips (e.g., around Epochs 7–8), reflecting the models adaptation to heavily augmented batches. It recovers swiftly, illustrating robustness.
- **High Final Accuracy:** After approximately 15–20 epochs, both curves stabilize above 97–98%. The training curve occasionally reaches 100%, while validation accuracy hovers around 98–99%.
- **Modest Overfitting:** Some minor gap appears between training and validation, but dropout and batch normalization help keep that gap in check.

3.6.4 Conclusion

- **Data Augmentation + CNN Blocks:** Combining augmentation with progressively deeper CNN blocks has proven effective, providing strong invariance to common handwriting distortions.
- **Separated Dense Layers with Dropout:** Including two dense layers (256 and 128 neurons) with a 40% dropout significantly reduces overfitting and improves generalization.
- **High Accuracy on MNIST (98.72%):** This confirms that my model is well-tuned for digit recognition tasks, thanks to the chosen methods (augmentation, block architecture, batch normalization, and dropout).
- **Practical Implications:** I saved the final network as `super_digit_classifier.keras`. This model can be easily loaded for inference or integrated into other applications (e.g., a digit recognition pipeline or an air-writing calculator).

3.7 Training the Operator Recognition Model: Code, Methods, Results, and Discussion

Overview: In this section, I explain how I trained a Convolutional Neural Network (CNN) to recognize four mathematical operators: addition (`add`), division (`divide`), multiplication (`multiply`), and subtraction (`subtract`). I used a folder-based dataset with separate directories for train, validation, and test splits. Below is the full code and a discussion of the training process, data augmentation, and final accuracy results.

3.7.1 Model Architecture and Training Code

```
1 import os
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import tensorflow as tf
5 from tensorflow.keras import layers, models
6 from tensorflow.keras.preprocessing.image import ImageDataGenerator
7
8 # Paths to train, val, test folders
9 base_dir = r"C:\Users\Stanislav\Downloads\archive (5)\final_symbols_split_ttv"
10 train_dir = os.path.join(base_dir, "train")
11 val_dir = os.path.join(base_dir, "val")
12 test_dir = os.path.join(base_dir, "test")
13
14 # 1) Generators
15 train_datagen = ImageDataGenerator(
16     rescale=1./255,
17     rotation_range=25,
18     width_shift_range=0.2,
19     height_shift_range=0.2,
20     zoom_range=0.2,
21     shear_range=10.0,
22     fill_mode='constant',
23     cval=0.0
24 )
25 val_datagen = ImageDataGenerator(rescale=1./255)
26 test_datagen = ImageDataGenerator(rescale=1./255)
27
28 batch_size = 32
29 img_size = (28, 28) # 28x28, grayscale
30
31 train_generator = train_datagen.flow_from_directory(
32     directory=train_dir,
33     target_size=img_size,
34     color_mode='grayscale',
35     class_mode='categorical',
36     batch_size=batch_size,
37     shuffle=True
38 )
39
40 val_generator = val_datagen.flow_from_directory(
```

```
41     directory=val_dir,
42     target_size=img_size,
43     color_mode='grayscale',
44     class_mode='categorical',
45     batch_size=batch_size,
46     shuffle=False
47 )
48
49 test_generator = test_datagen.flow_from_directory(
50     directory=test_dir,
51     target_size=img_size,
52     color_mode='grayscale',
53     class_mode='categorical',
54     batch_size=batch_size,
55     shuffle=False
56 )
57
58 num_classes = train_generator.num_classes
59 print("Number of classes:", num_classes)
60 print("class_indices (train):", train_generator.class_indices)
61
62 # 2) CNN Model
63 model = models.Sequential()
64
65 # Block 1
66 model.add(layers.Conv2D(32, (3,3), padding='same', input_shape=(28,28,1)))
67 model.add(layers.BatchNormalization())
68 model.add(layers.Activation('relu'))
69
70 model.add(layers.Conv2D(32, (3,3), padding='same'))
71 model.add(layers.BatchNormalization())
72 model.add(layers.Activation('relu'))
73 model.add(layers.MaxPooling2D((2,2))) # -> 14x14
74
75 # Block 2
76 model.add(layers.Conv2D(64, (3,3), padding='same'))
77 model.add(layers.BatchNormalization())
78 model.add(layers.Activation('relu'))
79
80 model.add(layers.Conv2D(64, (3,3), padding='same'))
81 model.add(layers.BatchNormalization())
82 model.add(layers.Activation('relu'))
83 model.add(layers.MaxPooling2D((2,2))) # -> 7x7
84
85 # Block 3
86 model.add(layers.Conv2D(128, (3,3), padding='same'))
87 model.add(layers.BatchNormalization())
88 model.add(layers.Activation('relu'))
89
90 model.add(layers.Conv2D(128, (3,3), padding='same'))
91 model.add(layers.BatchNormalization())
92 model.add(layers.Activation('relu'))
```

```
93 model.add(layers.MaxPooling2D((2,2))) # -> 3x3
94
95 # Dense part
96 model.add(layers.Flatten())
97 model.add(layers.Dense(256))
98 model.add(layers.BatchNormalization())
99 model.add(layers.Activation('relu'))
100 model.add(layers.Dropout(0.4))
101
102 model.add(layers.Dense(128))
103 model.add(layers.BatchNormalization())
104 model.add(layers.Activation('relu'))
105 model.add(layers.Dropout(0.4))
106
107 model.add(layers.Dense(num_classes, activation='softmax'))
108
109 model.compile(
110     optimizer='adam',
111     loss='categorical_crossentropy',
112     metrics=['accuracy']
113 )
114
115 model.summary()
116
117 # 3) Training
118 epochs = 20
119 history = model.fit(
120     train_generator,
121     epochs=epochs,
122     validation_data=val_generator
123 )
124
125 # 4) Plot Accuracy
126 plt.plot(history.history['accuracy'], label='Train acc')
127 plt.plot(history.history['val_accuracy'], label='Val acc')
128 plt.legend()
129 plt.title("Operator Model Accuracy")
130 plt.show()
131
132 # 5) Test Evaluation
133 test_loss, test_acc = model.evaluate(test_generator, verbose=0)
134 print("Test accuracy:", test_acc)
135
136 # 6) Save the Model
137 model.save("operators_model_aiaiai.keras")
138 print("Model saved: operators_model.keras")
139
140 # 7) Inverse Class Dictionary
141 class_indices = train_generator.class_indices
142 inv_map = {v: k for k,v in class_indices.items()}
143 print("Inverse map:", inv_map)
```

Key Points:

- **Dataset:** The dataset is split into `train`, `val`, and `test` folders, each containing four subdirectories: `add`, `divide`, `multiply`, and `subtract`.
- **Data Augmentation for Training:**
 - `rotation_range=25`, `width_shift_range=0.2`, `height_shift_range=0.2`, `zoom_range=0.2`, `shear_range=10.0`
 - This helps the network handle variations in operator drawings, scale, and position.
- **Validation and Test Generators:** Only rescaling is applied (no augmentation). This ensures the validation and test sets remain consistent for accurate performance checks.
- **CNN Blocks:** Like before, I employ three convolution blocks, each containing two `Conv2D` + `BatchNormalization` layers and a `MaxPooling2D` layer.
- **Dense Layers:** I flatten the feature maps, then use two fully connected layers (256, 128) with `Dropout(0.4)`, `BatchNormalization`, and `ReLU`.
- **Final Layer:** A softmax output for the 4 classes (`add`, `divide`, `multiply`, `subtract`).

3.7.2 Results and Interpretation

```
1 Epoch 1/20
2 364/364 ... accuracy: 0.7753 - val_accuracy: 0.2045
3 Epoch 2/20
4 364/364 ... accuracy: 0.9501 - val_accuracy: 0.9193
5 Epoch 7/20
6 364/364 ... accuracy: 0.9833 - val_accuracy: 0.9767
7 Epoch 14/20
8 364/364 ... accuracy: 0.9872 - val_accuracy: 0.2604
9 Epoch 17/20
10 364/364 ... accuracy: 0.9906 - val_accuracy: 0.9974
11 Epoch 20/20
12 364/364 ... accuracy: 0.9933 - val_accuracy: 0.9979
13
14 Test accuracy: 0.9994066953659058
```

Training Log Snippet

- The training accuracy quickly escalates above 95% by epoch 2, showing that the model is capable of distinguishing most operator images with few errors.
- Validation accuracy fluctuates (e.g., it drops to around 20% at epoch 1, 26% at epoch 14, then jumps to above 99%). These strong oscillations may be due to differences in how augmented data aligns with the small-ish validation set, or potentially class imbalance in certain subsets.
- By the end of training (epoch 20), validation accuracy sits near 99.79%, and **test accuracy is about 99.94%**.

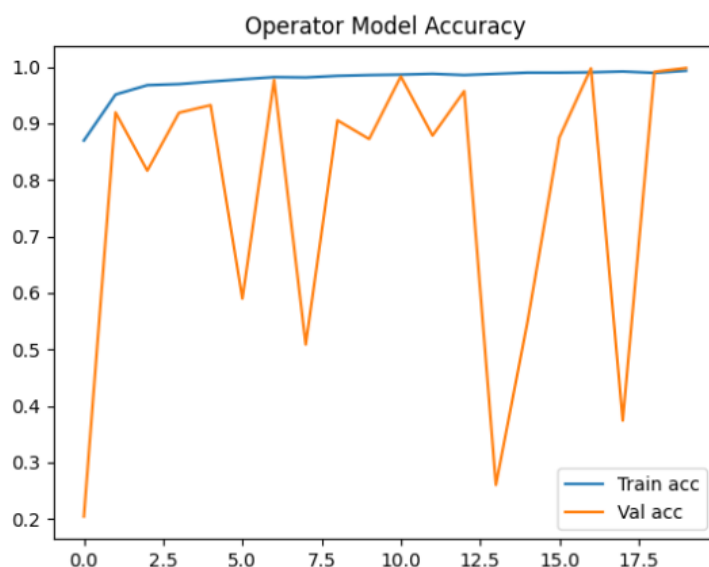


Figure 4: Train vs. Validation Accuracy for the Operator Model

Accuracy Plot Analysis

- **High and Steady Training Accuracy:** The blue line (train accuracy) remains above 97–98% from early in the training, reaching nearly 99–100% by later epochs.
- **Significant Validation Variations:** The orange line (val accuracy) shows big swings (from around 20% to over 90–95% or higher). These fluctuations might occur if certain batches in the validation set are more challenging or if there’s a class imbalance that occasionally confuses the model.
- **Final Convergence:** Despite the dips, the validation accuracy ends extremely high (around 99–100%), indicating the network eventually generalizes well.

3.7.3 Conclusion

- **Data Augmentation for Operators:** The strategy of random rotations, shifts, zooming, and shearing helps the model handle a diverse set of handwritten operator images.
- **Robust CNN Blocks + Dense Layers:** Repeating convolution blocks and adding two dense layers with dropout effectively reduce overfitting, even though occasional validation drops appear
- **Excellent Final Performance:** A test accuracy of **99.94%** confirms that the model distinguishes `add`, `divide`, `multiply`, and `subtract` with near-perfect accuracy on the test set.

“\latex

4 Conclusion

In this report, I presented the development of the "Air Calculator" system, which enables the recognition of handwritten mathematical expressions through real-time hand gesture tracking and convolutional neural networks (CNNs). I detailed every aspect of the project from the foundational hardware and functional assumptions to the comprehensive software design and the training of both digit and operator recognition models. My key conclusions are summarized below:

- **Innovative User Interaction:** I developed a system that allows users to write mathematical expressions in mid-air using hand gestures, leveraging technologies such as Mediapipe for hand tracking and OpenCV for image processing. This approach results in an intuitive and interactive user experience.
- **Accurate Hand and Finger Tracking:** By focusing on precise tracking of the index fingertip and incorporating a thumb gesture for mode switching, I ensured that the system reliably captures drawing inputs and manages actions like drawing and erasing on a virtual canvas.
- **Effective Input Segmentation and Processing:** I implemented an efficient segmentation process that divides the drawn symbols into distinct groups, identifying operators and operands. For standard arithmetic expressions (e.g., "**number -- operator -- number**"), the system accurately performs the corresponding calculations while gracefully handling errors (such as division by zero).
- **High-Performing Recognition Models:** The digit recognition model, which I trained on the MNIST dataset using data augmentation and a robust CNN architecture, achieved an accuracy of approximately 98.72%. Similarly, the operator recognition model demonstrated an exceptional test accuracy of around 99.94%. These results confirm the reliability and robustness of the recognition models under various handwriting distortions.
- **Comprehensive System Design:** I provided detailed explanations of the CNN architectures, data augmentation techniques, and modular software components. This comprehensive approach not only clarifies the systems operation but also facilitates future improvements and integrations into other applications.

In summary, I successfully demonstrated the potential of modern computer vision and deep learning technologies to create an interactive tool for recognizing handwritten mathematical expressions. The high accuracy of the models and the intuitive gesture-based control of the system open up promising avenues for further development in educational, commercial, and research settings.

Appendix

The full Python code is listed below:

```
1 import cv2
2 import mediapipe as mp
3 import numpy as np
4 from tensorflow.keras.models import load_model
5
6 print("Starting Air Calculator...")
7
8 digits_model = load_model("super_digit_classifier.keras")
9 operators_model = load_model("operators_model.keras")
10
11 mp_hands = mp.solutions.hands
12 mp_drawing = mp.solutions.drawing_utils
13
14 canvas = None
15 previous_position = None
16
17 DOT_IGNORE_AREA = 10
18 DOT_AREA_THRESHOLD = 40
19
20 def crop_and_predict_digit_or_dot(cropped):
21     """Predict digit 0..9 or '.' or None if area too small. Uses verbose=1."""
22     h, w = cropped.shape[:2]
23     area = w*h
24     if area < DOT_IGNORE_AREA:
25         return None
26     if area < DOT_AREA_THRESHOLD:
27         return '.'
28     max_side = max(w, h)
29     scale = 20.0 / max_side
30     new_w = int(w * scale)
31     new_h = int(h * scale)
32     resized = cv2.resize(cropped, (new_w, new_h), interpolation=cv2.INTER_AREA)
33     out28 = np.zeros((28, 28), dtype=np.uint8)
34     sx = (28 - new_w) // 2
35     sy = (28 - new_h) // 2
36     out28[sy:sy+new_h, sx:sx+new_w] = resized
37     final_28 = out28.astype('float32') / 255.0
38     x = final_28.reshape(1, 28, 28, 1)
39     preds = digits_model.predict(x, verbose=1) # verbose=1
40     d = np.argmax(preds)
41     return str(d)
42
43 def crop_and_predict_operator(cropped):
44     """Predict + / * - with inversion. Uses verbose=1 for progress."""
45     h, w = cropped.shape[:2]
46     max_side = max(w, h)
47     scale = 20.0 / max_side
48     new_w = int(w * scale)
```

```
49     new_h = int(h * scale)
50     resized = cv2.resize(cropped, (new_w, new_h), interpolation=cv2.INTER_AREA)
51     out28 = np.zeros((28, 28), dtype=np.uint8)
52     sx = (28 - new_w) // 2
53     sy = (28 - new_h) // 2
54     out28[sy:sy+new_h, sx:sx+new_w] = resized
55     out28 = 255 - out28
56     final_28 = out28.astype('float32') / 255.0
57     x = final_28.reshape(1, 28, 28, 1)
58     preds = operators_model.predict(x, verbose=1)
59     idx = np.argmax(preds)
60     op_map = {0: '+', 1: '/', 2: '*', 3: '-'}
61     return op_map.get(idx, '?')
62
63 def preprocess_and_find_contours(img):
64     """Threshold + dilate => contours."""
65     gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
66     _, bin_canvas = cv2.threshold(gray, 127, 255, cv2.THRESH_BINARY)
67     kernel = np.ones((3,3), np.uint8)
68     bin_dil = cv2.dilate(bin_canvas, kernel, iterations=1)
69     contours, _ = cv2.findContours(bin_dil, cv2.RETR_EXTERNAL, cv2.CHAIN_APPROX_SIMPLE)
70     return bin_dil, contours
71
72 def split_line_into_groups(boxes, space_x_thresh=15):
73     """Sort by x, group bounding boxes if gap <= space_x_thresh."""
74     boxes.sort(key=lambda b: b[0])
75     sublines = []
76     current = [boxes[0]]
77     for i in range(1, len(boxes)):
78         x,y,w,h = boxes[i]
79         px,py,pw,ph = boxes[i-1]
80         gap = x - (px+pw)
81         if gap > space_x_thresh:
82             sublines.append(current)
83             current = [boxes[i]]
84         else:
85             current.append(boxes[i])
86     if current:
87         sublines.append(current)
88     return sublines
89
90 def process_subgroup(bin_dil, boxes):
91     """Merge all bounding boxes => string of digits or '.'."""
92     boxes.sort(key=lambda b: b[0])
93     s = ""
94     for (x,y,w,h) in boxes:
95         cropped = bin_dil[y:y+h, x:x+w]
96         c = crop_and_predict_digit_or_dot(cropped)
97         if c is not None:
98             s += c
99     if s.endswith('.'):
100         s = s[:-1]
```

```
101     if s == "":
102         s = "0"
103     return s
104
105 def process_operator_subgroup(bin_dil, boxes):
106     """Assume first bounding box is operator."""
107     boxes.sort(key=lambda b: b[0])
108     x,y,w,h = boxes[0]
109     cropped = bin_dil[y:y+h, x:x+w]
110     return crop_and_predict_operator(cropped)
111
112 def predict_all_in_one_line(canvas_img):
113     """We read all bounding boxes as if they are in one single line."""
114     bin_dil, contours = preprocess_and_find_contours(canvas_img)
115     boxes = []
116     for c in contours:
117         x,y,w,h = cv2.boundingRect(c)
118         boxes.append((x,y,w,h))
119     if not boxes:
120         return
121     sublines = split_line_into_groups(boxes, space_x_thresh=15)
122     if len(sublines) == 3:
123         # interpret as equation
124         left_str = process_subgroup(bin_dil, sublines[0])
125         op_str = process_operator_subgroup(bin_dil, sublines[1])
126         right_str = process_subgroup(bin_dil, sublines[2])
127         try:
128             lv = float(left_str)
129             rv = float(right_str)
130             res = None
131             if op_str == '+': res = lv + rv
132             elif op_str == '-': res = lv - rv
133             elif op_str == '*': res = lv * rv
134             elif op_str == '/':
135                 if rv == 0: res = "ERR(div0)"
136             else: res = lv / rv
137             if isinstance(res, float) and not isinstance(res, str):
138                 if abs(res - round(res)) < 1e-9:
139                     res = int(round(res))
140             print(f"{left_str} {op_str} {right_str} = {res}")
141         except:
142             print(f"{left_str} {op_str} {right_str} = ERR")
143     else:
144         # multiple subgroups => print them with semicolon
145         results = []
146         for grp in sublines:
147             val = process_subgroup(bin_dil, grp)
148             results.append(val)
149         print("; ".join(results))
150
151 def is_thumb_extended(landmarks):
152     thumb_tip = landmarks.landmark[4]
```

```
153     thumb_base = landmarks.landmark[2]
154     return abs(thumb_tip.x - thumb_base.x) > 0.05
155
156 def is_hand_open(landmarks):
157     tip_ids = [8,12,16,20]
158     base_ids = [6,10,14,18]
159     for tip_id, base_id in zip(tip_ids, base_ids):
160         if landmarks.landmark[tip_id].y > landmarks.landmark[base_id].y:
161             return False
162     return True
163
164 def erase_area(landmarks, canvas_img, shape):
165     try:
166         palm_ids = [0,1,5,9,13,17]
167         pls = [landmarks.landmark[i] for i in palm_ids]
168         xs = [int(lm.x*shape[1]) for lm in pls]
169         ys = [int(lm.y*shape[0]) for lm in pls]
170         x_min, x_max = min(xs), max(xs)
171         y_min, y_max = min(ys), max(ys)
172         cv2.rectangle(canvas_img, (x_min, y_min), (x_max, y_max), (0,0,0), -1)
173     except:
174         for fid in [8,12,16,20]:
175             fx = int(landmarks.landmark[fid].x*shape[1])
176             fy = int(landmarks.landmark[fid].y*shape[0])
177             cv2.circle(canvas_img, (fx, fy), 20, (0,0,0), -1)
178
179 def main():
180     global canvas, previous_position
181     mp_hands = mp.solutions.hands
182     mp_drawing = mp.solutions.drawing_utils
183     cap = cv2.VideoCapture(0)
184     with mp_hands.Hands(model_complexity=0, min_detection_confidence=0.8,
185                         min_tracking_confidence=0.5) as hands:
186         while True:
187             ret, frame = cap.read()
188             if not ret:
189                 break
190             frame = cv2.flip(frame, 1)
191             rgb = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
192             if canvas is None:
193                 canvas = np.zeros_like(frame)
194             results = hands.process(rgb)
195             if results.multi_hand_landmarks:
196                 for lm in results.multi_hand_landmarks:
197                     mp_drawing.draw_landmarks(frame, lm, mp_hands.HAND_CONNECTIONS)
198                     fx = int(lm.landmark[8].x * frame.shape[1])
199                     fy = int(lm.landmark[8].y * frame.shape[0])
200                     if is_hand_open(lm):
201                         erase_area(lm, canvas, frame.shape)
202                     continue
203                     if is_thumb_extended(lm):
204                         if previous_position is not None:
```

```
205         cv2.line(canvas, previous_position, (fx, fy), (255,255,255), 5)
206         previous_position = (fx, fy)
207     else:
208         previous_position = None
209     else:
210         previous_position = None
211     combined = cv2.addWeighted(frame, 0.7, canvas, 0.3, 0)
212     cv2.imshow("Air Draw Calculator", combined)
213     key = cv2.waitKey(1) & 0xFF
214     if key == ord('q'):
215         break
216     elif key == ord('r'):
217         predict_all_in_one_line(canvas)
218     elif key == ord('c'):
219         canvas = np.zeros_like(frame)
220     cap.release()
221     cv2.destroyAllWindows()
222
223 if __name__ == "__main__":
224     main()
225
```