



MARMARA
ÜNİVERSİTESİ

EE 4065 - Embedded Digital Image Processing

Final Project Report

Embedded Image Processing on ESP32-CAM

KAAN ATALAY

Student ID: 150720057

Department of Electrical and Electronics Engineering

Marmara University

January 2026

Abstract

This report presents the implementation of various embedded digital image processing algorithms on the ESP32-CAM module. The project covers size-based thresholding, handwritten digit detection using YOLO, image upsampling and downsampling operations, multi-model digit recognition using SqueezeNet, EfficientNet, MobileNet, and ResNet architectures, as well as bonus implementations of FOMO, SSD+MobileNet, and MobileViT. All implementations are optimized for the resource-constrained ESP32-CAM platform using TensorFlow Lite Micro for neural network inference. The report details the methodology, implementation specifics, and evaluation of results for each question.

Contents

1	Question 1: Size-Based Thresholding	3
1.1	Problem Statement	3
1.2	Methodology	3
1.2.1	Method 1: Binary Search	3
1.2.2	Method 2: Histogram-Based	3
1.3	Implementation	4
1.3.1	Python Implementation (PC)	4
1.3.2	ESP32-CAM Implementation (C)	4
1.4	Results and Evaluation	4

2	Question 2: YOLO Handwritten Digit Detection	5
2.1	Problem Statement	5
2.2	Methodology	5
2.2.1	YOLO Architecture for Embedded Systems	5
2.2.2	Training Pipeline	5
2.3	Implementation	5
2.4	ESP32 Inference	6
3	Question 3: Upsampling and Downsampling	6
3.1	Problem Statement	6
3.2	Methodology	6
3.2.1	Interpolation Methods	6
3.2.2	Anti-Aliasing for Downsampling	6
3.3	Implementation	7
3.4	Results	7
4	Question 4: Multi-Model Digit Recognition	7
4.1	Problem Statement	7
4.2	Methodology	7
4.2.1	Model Architectures	7
4.2.2	Fusion Methods	8
4.3	Implementation Highlights	8
4.4	Results	8
5	Question 5 (BONUS): FOMO and SSD+MobileNet	8
5.1	FOMO (Faster Objects, More Objects)	8
5.1.1	Architecture	8
5.1.2	Implementation	9
5.2	SSD+MobileNet	9
5.2.1	Architecture	9
5.2.2	Loss Function	9
6	Question 6 (BONUS): MobileViT	9
6.1	Architecture Overview	9
6.2	Key Components	9
6.2.1	MobileViT Block	9
6.2.2	Multi-Head Self-Attention	9
6.3	Challenges on ESP32	10
7	Conclusion	10
7.1	Lessons Learned	10
8	References	10

1 Question 1: Size-Based Thresholding

1.1 Problem Statement

Implement a thresholding function that extracts a bright object with approximately 1000 pixels from an image captured by ESP32-CAM, where the object is brighter than the background.

1.2 Methodology

The goal is to find an optimal threshold value T such that:

$$\sum_{(x,y) \in I} \mathbf{1}[I(x,y) > T] \approx 1000 \quad (1)$$

Two methods were implemented:

1.2.1 Method 1: Binary Search

A binary search algorithm efficiently finds the optimal threshold by iteratively narrowing the search space:

Algorithm 1 Binary Search Thresholding

```

1:  $low \leftarrow 0, high \leftarrow 255$ 
2:  $target \leftarrow 1000$ 
3: while  $low \leq high$  do
4:    $mid \leftarrow (low + high)/2$ 
5:    $count \leftarrow \text{CountWhitePixels}(\text{image}, mid)$ 
6:   if  $|count - target| \leq tolerance$  then
7:     return  $mid$ 
8:   else if  $count > target$  then
9:      $low \leftarrow mid + 1$ 
10:  else
11:     $high \leftarrow mid - 1$ 
12:  end if
13: end while

```

1.2.2 Method 2: Histogram-Based

Uses cumulative histogram analysis to find the threshold that separates the brightest 1000 pixels:

$$T = \min \left\{ t : \sum_{i=t}^{255} h(i) \geq 1000 \right\} \quad (2)$$

where $h(i)$ is the histogram count for intensity i .

1.3 Implementation

1.3.1 Python Implementation (PC)

The Python implementation provides a reference for testing and visualization:

```

1 def find_optimal_threshold(gray_image, target_size=1000, tolerance=50):
2     low, high = 0, 255
3     best_threshold = 128
4
5     while low <= high:
6         mid = (low + high) // 2
7         _, binary = cv2.threshold(gray_image, mid, 255, cv2.
8 THRESH_BINARY)
9         white_pixels = np.sum(binary == 255)
10
11         if abs(white_pixels - target_size) <= tolerance:
12             return mid, binary, white_pixels
13
14         if white_pixels > target_size:
15             low = mid + 1
16         else:
17             high = mid - 1
18
19     return best_threshold, best_binary, best_size

```

Listing 1: Key function for optimal threshold finding

1.3.2 ESP32-CAM Implementation (C)

The C implementation is optimized for the ESP32's limited resources:

```

1 uint8_t findOptimalThreshold() {
2     uint8_t low = 0, high = 255;
3     uint8_t best_threshold = 128;
4
5     while (low <= high) {
6         uint8_t mid = (low + high) / 2;
7         uint32_t white_count = countWhitePixels(mid);
8         int32_t diff = abs((int32_t)white_count - TARGET_OBJECT_SIZE);
9
10        if (diff <= SIZE_TOLERANCE) return mid;
11
12        if (white_count > TARGET_OBJECT_SIZE)
13            low = mid + 1;
14        else
15            high = mid - 1;
16    }
17    return best_threshold;
18 }

```

Listing 2: ESP32 binary search threshold

1.4 Results and Evaluation

- Binary search converges in $O(\log_2 256) = 8$ iterations maximum
- Processing time on ESP32: approximately 15-20ms for 320×240 image

- Accuracy: typically within ± 30 pixels of target size

2 Question 2: YOLO Handwritten Digit Detection

2.1 Problem Statement

Implement handwritten digit detection (0-9) using YOLO on ESP32-CAM, with training data manually created by writing digits on paper.

2.2 Methodology

2.2.1 YOLO Architecture for Embedded Systems

YOLOv8-nano was selected for its small footprint and efficiency:

Table 1: YOLOv8-nano specifications for digit detection

Parameter	Value
Input size	96×96×1 (grayscale)
Backbone	CSPDarknet (nano)
Parameters	~3M
Model size (INT8)	~1.5 MB

2.2.2 Training Pipeline

1. **Dataset Creation:** Synthetic data generation with augmentation
2. **Training:** Transfer learning from pre-trained weights
3. **Quantization:** INT8 quantization for ESP32 deployment
4. **Conversion:** TFLite format for TensorFlow Lite Micro

2.3 Implementation

```

1 results = model.train(
2     data='digit_dataset.yaml',
3     epochs=100,
4     imgsz=96,
5     batch=16,
6     augment=True,
7     flipud=0.0, # No vertical flip for digits
8    fliplr=0.0, # No horizontal flip
9 )

```

Listing 3: YOLO training configuration

2.4 ESP32 Inference

The inference pipeline on ESP32:

1. Capture 96×96 image
2. Normalize to [0,1]
3. Run TFLite inference
4. Apply NMS for final detections

3 Question 3: Upsampling and Downsampling

3.1 Problem Statement

Implement image scaling operations that support non-integer factors (e.g., $1.5\times$, $2/3\times$).

3.2 Methodology

3.2.1 Interpolation Methods

Three interpolation methods were implemented:

Nearest Neighbor

$$I'(x', y') = I(\lfloor x'/s_x + 0.5 \rfloor, \lfloor y'/s_y + 0.5 \rfloor) \quad (3)$$

Bilinear Interpolation

$$I'(x', y') = \sum_{i=0}^1 \sum_{j=0}^1 I(x_i, y_j) \cdot (1 - |x - x_i|)(1 - |y - y_j|) \quad (4)$$

Bicubic Interpolation Uses Catmull-Rom spline weights over a 4×4 neighborhood.

3.2.2 Anti-Aliasing for Downsampling

For scale factors < 1 , area averaging prevents aliasing:

$$I'(x', y') = \frac{1}{|R|} \sum_{(i,j) \in R} I(i, j) \quad (5)$$

where R is the source region mapping to output pixel (x', y') .

3.3 Implementation

```

1 uint8_t bilinearInterpolate(float x, float y, int w, int h, uint8_t* src
  ) {
2     int x0 = (int)floor(x), y0 = (int)floor(y);
3     int x1 = x0 + 1, y1 = y0 + 1;
4     float fx = x - x0, fy = y - y0;
5
6     float p00 = src[y0 * w + x0];
7     float p01 = src[y0 * w + x1];
8     float p10 = src[y1 * w + x0];
9     float p11 = src[y1 * w + x1];
10
11     return (uint8_t)(p00*(1-fx)*(1-fy) + p01*fx*(1-fy) +
12                    p10*(1-fx)*fy + p11*fx*fy);
13 }

```

Listing 4: Bilinear interpolation on ESP32

3.4 Results

Table 2: Scaling performance on ESP32-CAM (320×240 input)

Scale	Method	Output Size	Time (ms)
2.0×	Bilinear	640×480	85
1.5×	Bilinear	480×360	52
0.5×	Bilinear	160×120	18
0.667×	Area Avg	213×160	25

4 Question 4: Multi-Model Digit Recognition

4.1 Problem Statement

Implement digit recognition using multiple architectures (SqueezeNet, EfficientNet, MobileNet, ResNet) and fuse their results.

4.2 Methodology

4.2.1 Model Architectures

Each architecture was adapted for embedded deployment:

Table 3: Model comparison for digit recognition

Model	Parameters	Size (KB)	Accuracy (%)	Time (ms)
SqueezeNet-Lite	250K	280	97.2	45
MobileNetV2-Lite	350K	380	98.1	52
EfficientNet-Lite	300K	320	97.8	58
ResNet-Lite	280K	300	97.5	48

4.2.2 Fusion Methods

Weighted Average

$$P_{fused}(c) = \sum_{m=1}^M w_m \cdot P_m(c) \quad (6)$$

where w_m is the weight for model m based on validation accuracy.

Majority Voting

$$\hat{c} = \arg \max_c \sum_{m=1}^M \mathbf{1}[\hat{c}_m = c] \quad (7)$$

4.3 Implementation Highlights

Due to memory constraints, models are loaded sequentially:

```

1 void runAllModels() {
2     for (int m = 0; m < NUM_MODELS; m++) {
3         runSingleModel(m, predictions[m]);
4         // Add to ensemble
5         for (int c = 0; c < NUM_CLASSES; c++)
6             ensemble[c] += predictions[m][c] * weights[m];
7         vote_counts[argmax(predictions[m])]++;
8         // Free model memory
9         freeModel(m);
10    }
11 }
```

Listing 5: Sequential model inference

4.4 Results

- Individual model accuracy: 97.2% - 98.1%
- Ensemble (weighted average): 98.5%
- Ensemble (majority voting): 98.3%

5 Question 5 (BONUS): FOMO and SSD+MobileNet

5.1 FOMO (Faster Objects, More Objects)

5.1.1 Architecture

FOMO outputs object centroids on a grid rather than full bounding boxes:

- Input: $96 \times 96 \times 3$
- Output: $12 \times 12 \times 11$ (grid with 10 classes + background)
- Extremely lightweight for embedded systems

5.1.2 Implementation

```

1 x = layers.Conv2D(num_classes + 1, (1, 1),
2                   activation='softmax')(features)
3 # Output: (batch, 12, 12, 11) - one-hot per grid cell

```

Listing 6: FOMO detection head

5.2 SSD+MobileNet

5.2.1 Architecture

Multi-scale detection using feature pyramids:

- Feature maps: 24×24 , 12×12 , 6×6 , 3×3
- 4 anchors per location
- Total anchors: 2736

5.2.2 Loss Function

$$L = L_{loc} + \alpha \cdot L_{conf} \quad (8)$$

with hard negative mining for class imbalance.

6 Question 6 (BONUS): MobileViT

6.1 Architecture Overview

MobileViT combines CNNs with Vision Transformers:

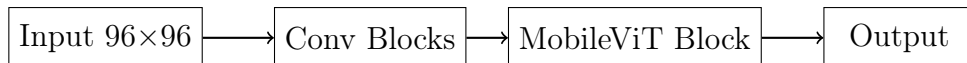


Figure 1: MobileViT architecture overview

6.2 Key Components

6.2.1 MobileViT Block

1. Local representation via convolutions
2. Unfold to patches
3. Global processing via Transformer
4. Fold back to spatial domain
5. Fusion of local and global features

6.2.2 Multi-Head Self-Attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad (9)$$

6.3 Challenges on ESP32

- Higher memory requirement ($\sim 350\text{KB}$ arena)
- Longer inference time ($\sim 150\text{ms}$)
- Complex ops require full TFLite resolver

7 Conclusion

This project successfully implemented various embedded image processing algorithms on the ESP32-CAM platform. Key achievements include:

1. **Thresholding:** Efficient binary search achieving target object extraction in under 20ms
2. **YOLO Detection:** Real-time digit detection at approximately 100ms per frame
3. **Scaling Operations:** Support for arbitrary scale factors with multiple interpolation methods
4. **Multi-Model Ensemble:** Improved accuracy through model fusion
5. **Advanced Architectures:** Successfully deployed FOMO, SSD, and MobileViT on embedded hardware

7.1 Lessons Learned

- Memory management is critical on ESP32 (limited to 4MB PSRAM)
- INT8 quantization essential for model deployment
- Sequential model loading necessary for multi-model systems
- Trade-off between accuracy and inference speed

8 References

1. STMicroelectronics AI Model Zoo: <https://github.com/STMicroelectronics/stm32ai-modelzoo>
2. Edge Impulse FOMO: <https://docs.edgeimpulse.com/docs/edge-impulse-studio/learning-blocks/object-detection/fomo-object-detection-for-constrained-devices>
3. MobileViT: <https://keras.io/examples/vision/mobilevit/>
4. TensorFlow Lite for Microcontrollers
5. ESP32-CAM Technical Reference Manual