# Embedded Digital Image Processing

# EE4065

# Homework 6

Mustafa Öztürk

150722016

Emirhan Şener

150721069

Faculty of Engineering
Electrical Electronic Engineering

Fall 2025

# ACRONYMS

**CNN** Convolutional Neural Network

**EE4065** Embedded Digital Image Processing

**Hu Moments** Hu Invariant Moments

**MCU** Microcontroller Unit

**MNIST** Modified National Institute of Standards and Technology

**RAM** Random Access Memory

**UART** Universal Asynchronous Receiver Transmitter

# CONTENTS

# LIST OF FIGURES

# 1. INTRODUCTION

This report presents the implementation and results of the $6^{th}$ homework for the Embedded Digital Image Processing (EE4065) course. Building upon previous assignments where we explored basic neural networks and feature extraction, this work focuses on deploying deep Convolutional Neural Network (CNN) architectures for image classification. The primary objective is to implement handwritten digit recognition on the STM32F446RE Microcontroller Unit (MCU), balancing model accuracy with the strict resource constraints of an embedded system.

The tasks in this assignment are based on Section 13.7 of the course textbook [1]. We utilized the Modified National Institute of Standards and Technology (MNIST) dataset to train and evaluate several standard architecture families, including EfficientNet, ResNet, MobileNet, and SqueezeNet. The goal was to compare these models not only on classification performance but also on their deployability regarding inference speed and memory footprint.

We used the **ST Edge AI Developer Cloud** [2] to optimize and convert our trained TensorFlow models into C code executable on the STM32. The workflow involved converting models from .h5 to .tflite format and applying INT8 quantization to minimize size. During this process, we observed that the optimized EfficientNet and SqueezeNet models exceeded the 512kB flash memory limit of the target MCU, rendering them unsuitable for this specific hardware despite sufficient Random Access Memory (RAM). Consequently, the on-board implementation phase focuses on ResNetV1, MobileNetV1, FD MobileNetV1, and ST FD MobileNetV1.

All models were initially trained and validated on a PC using Python and TensorFlow [3]. We analyzed the performance using accuracy metrics and confusion matrices. The compatible models were then generated as optimized projects via the cloud service and deployed to the MCU to demonstrate real-time inference capabilities.

# 2.  PROBLEMS

## 2.1.  Q-1) SECTION 13.7 APPLICATION: HANDWRITTEN DIGIT RECOGNITION FROM DIGITAL IMAGES

**Implement the end-of-chapter application from the course textbook [1]. The goal is to perform Handwritten Digit Recognition from Digital Images using CNN models. You are required to use all the mentioned CNN models (SqueezeNet, EfficientNet, etc.) in the course and implement the solution on the target hardware.**

### Theory

### CNNs and Quantization on Embedded Systems

Deploying CNNs on the STM32F446RE requires balancing model complexity with limited Flash (512kB) and RAM (128kB). A critical optimization is **Quantization**, which converts 32-bit floating-point weights to 8-bit integers. This reduces model size by 75% and accelerates inference on MCUs without floating-point units, utilizing integer arithmetic for matrix operations.

### Efficient Architectures

We evaluated several architectures optimized for efficiency:

- **ResNetV1:** Utilizes residual skip connections to allow deeper networks while maintaining training stability.

- **MobileNetV1:** Designed for embedded vision, it uses *Depthwise Separable Convolutions* to drastically reduce parameters and computation compared to standard convolutions.

- **SqueezeNet & EfficientNet:** architectures that prioritize high accuracy-per-parameter. However, their complex graph structures often result in binary sizes that exceed the Flash memory of smaller MCUs like the F446RE.

## ST Edge AI Developer Cloud

The **ST Edge AI Developer Cloud** serves as the deployment bridge. It optimizes the '.tflite' models, validates memory constraints (Flash/RAM), and generates the optimized C library (X-CUBE-AI) required to run inference on the STM32 hardware.

## Procedure

## Step 1: Training the CNN Models

We trained a variety of CNN architectures to compare their performance and suitability for embedded deployment. The models include ResNetV1, SqueezeNetV1.1, EfficientNetV2-B0, MobileNetV1/V2, FDMobileNet, and ST-optimized variants.

**Training Process:** The models were trained using the MNIST dataset. Since MNIST images are 28x28 grayscale, we preprocessed them by resizing to 32x32 and converting them to 3-channel (RGB) format to match the input requirements of standard CNN backbones. Pixel values were normalized to the [0, 1] range.

Training was performed using the Adam optimizer with a learning rate of 0.001. We employed callbacks for model checkpointing (saving the best model based on validation accuracy), early stopping (patience of 5 epochs), and learning rate reduction on plateaus.

**Listing 2.1:** Training script for all CNN models (train_all_models.py)

```
1  import os
2  import sys
3  import numpy as np
4  import tensorflow as tf
5  from tensorflow import keras
6  import matplotlib
7  matplotlib.use('Agg')
8  import matplotlib.pyplot as plt
9  from sklearn.metrics import confusion_matrix, classification_report,
       accuracy_score, precision_score
10
11 sys.path.insert(0, os.path.join(os.path.dirname(__file__), 'suplementary
       '))
12
13 from resnetv1 import get_resnetv1
14 from squeezenetv11 import get_squeezenetv11
```

```
15  from efficientnetv2 import get_efficientnetv2
16  from mobilenetv1 import get_mobilenetv1
17  from mobilenetv2 import get_mobilenetv2
18  from fdmobilenet import get_fdmobilenet
19  from st_efficientnet_lc_v1 import get_st_efficientnet_lc_v1
20  from st_fdmobilenet_v1 import get_st_fdmobilenet_v1
21
22
23  # Configuration
24  NUM_CLASSES = 10
25  DATA_SHAPE = (32, 32, 3)
26  EPOCHS = 20
27  BATCH_SIZE = 64
28  PATIENCE = 5
29  MODELS_DIR = os.path.join(os.path.dirname(__file__), 'models')
30  RESULTS_DIR = os.path.join(os.path.dirname(__file__), 'results')
31
32
33  def prepare_tensor(images, out_shape):
34      images = tf.expand_dims(images, axis=-1)
35      images = tf.repeat(images, 3, axis=-1)
36      images = tf.image.resize(images, out_shape[:2])
37      images = images / 255.0
38      return images
39
40
41  def load_mnist_data():
42      print("Loading MNIST dataset...")
43      (train_images, train_labels), (val_images, val_labels) = tf.keras.
            datasets.mnist.load_data()
44
45      print("Preprocessing data...")
46      train_images = prepare_tensor(train_images, DATA_SHAPE)
47      val_images = prepare_tensor(val_images, DATA_SHAPE)
48
49      train_labels = tf.keras.utils.to_categorical(train_labels,
            NUM_CLASSES)
50      val_labels = tf.keras.utils.to_categorical(val_labels, NUM_CLASSES)
51
52      print(f"Training samples: {len(train_images)}")
53      print(f"Validation samples: {len(val_images)}")
54      print(f"Image shape: {train_images.shape[1:]}")
55
56      return (train_images, train_labels), (val_images, val_labels)
57
58
59  def create_callbacks(model_name):
60      model_path = os.path.join(MODELS_DIR, f'{model_name}.h5')
61
62      callbacks = [
63          keras.callbacks.ModelCheckpoint(
64              model_path,
65              save_best_only=True,
66              monitor='val_accuracy',
67              verbose=1
68          ),
69          keras.callbacks.EarlyStopping(
70              monitor='val_accuracy',
71              patience=PATIENCE,
72              verbose=1,
73              restore_best_weights=True
74          ),
75          keras.callbacks.ReduceLROnPlateau(
76              monitor='val_loss',
```

```
77              factor=0.5,
78              patience=3,
79              verbose=1,
80              min_lr=1e-6
81          )
82      ]
83      return callbacks
84
85
86  def train_model(model, model_name, train_data, val_data):
87      print(f"\n{'='*60}")
88      print(f"Training {model_name}")
89      print(f"{'='*60}")
90
91      # Compile model
92      model.compile(
93          optimizer=keras.optimizers.Adam(learning_rate=1e-3),
94          loss='categorical_crossentropy',
95          metrics=['accuracy']
96      )
97
98      model.summary()
99
100     # Train
101     train_images, train_labels = train_data
102     val_images, val_labels = val_data
103
104     history = model.fit(
105         x=train_images,
106         y=train_labels,
107         epochs=EPOCHS,
108         batch_size=BATCH_SIZE,
109         validation_data=(val_images, val_labels),
110         callbacks=create_callbacks(model_name),
111         verbose=1
112     )
113
114     # Evaluate
115     loss, accuracy = model.evaluate(val_images, val_labels, verbose=0)
116     print(f"\n{model_name} - Final Validation Accuracy: {accuracy*100:.2
            f}%")
117
118     return model, history, accuracy
119
120
121 def evaluate_and_plot_metrics(model, model_name, val_images, val_labels)
        :
122     os.makedirs(RESULTS_DIR, exist_ok=True)
123
124     # Get predictions
125     y_pred_prob = model.predict(val_images, verbose=0)
126     y_pred = np.argmax(y_pred_prob, axis=1)
127     y_true = np.argmax(val_labels, axis=1)
128
129     # Compute metrics
130     acc = accuracy_score(y_true, y_pred)
131     prec = precision_score(y_true, y_pred, average='weighted')
132
133     print(f"\n{model_name} Metrics:")
134     print(f"  Accuracy:  {acc*100:.2f}%")
135     print(f"  Precision: {prec*100:.2f}%")
136
137     # Confusion matrix
138     cm = confusion_matrix(y_true, y_pred)
```

```
139
140      # Plot confusion matrix
141      plt.figure(figsize=(8, 6))
142      plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
143      plt.title(f'{model_name} Confusion Matrix')
144      plt.colorbar()
145      plt.xlabel('Predicted')
146      plt.ylabel('True')
147      plt.xticks(range(10))
148      plt.yticks(range(10))
149
150      # Add text annotations
151      for i in range(10):
152          for j in range(10):
153              plt.text(j, i, str(cm[i, j]), ha='center', va='center',
154                       color='white' if cm[i, j] > cm.max()/2 else 'black'
                         )
155
156      plt.tight_layout()
157      plt.savefig(os.path.join(RESULTS_DIR, f'{model_name}
             _confusion_matrix.png'), dpi=150)
158      plt.close()
159
160      print(f"  Saved: {model_name}_confusion_matrix.png")
161
162      return acc, prec
163
164
165  def plot_training_history(history, model_name):
166      os.makedirs(RESULTS_DIR, exist_ok=True)
167
168      # Accuracy plot
169      plt.figure(figsize=(10, 4))
170
171      plt.subplot(1, 2, 1)
172      plt.plot(history.history['accuracy'], label='Train')
173      plt.plot(history.history['val_accuracy'], label='Validation')
174      plt.title(f'{model_name} - Accuracy')
175      plt.xlabel('Epoch')
176      plt.ylabel('Accuracy')
177      plt.legend()
178      plt.grid(True)
179
180      # Loss plot
181      plt.subplot(1, 2, 2)
182      plt.plot(history.history['loss'], label='Train')
183      plt.plot(history.history['val_loss'], label='Validation')
184      plt.title(f'{model_name} - Loss')
185      plt.xlabel('Epoch')
186      plt.ylabel('Loss')
187      plt.legend()
188      plt.grid(True)
189
190      plt.tight_layout()
191      plt.savefig(os.path.join(RESULTS_DIR, f'{model_name}_training_curves
             .png'), dpi=150)
192      plt.close()
193
194      print(f"  Saved: {model_name}_training_curves.png")
195
196
197  def get_all_models():
198      models = {}
199
```

```python
200        # ResNet V1
201        print("Creating ResNet V1...")
202        try:
203            models['resnet_v1'] = get_resnetv1(
204                num_classes=NUM_CLASSES,
205                input_shape=DATA_SHAPE,
206                depth=8,
207                dropout=0.2
208            )
209        except Exception as e:
210            print(f"Error creating ResNet: {e}")
211
212        # SqueezeNet V1.1
213        print("Creating SqueezeNet V1.1...")
214        try:
215            models['squeezenet_v11'] = get_squeezenetv11(
216                num_classes=NUM_CLASSES,
217                input_shape=DATA_SHAPE,
218                dropout=0.5
219            )
220        except Exception as e:
221            print(f"Error creating SqueezeNet: {e}")
222
223        # EfficientNet V2-B0
224        print("Creating EfficientNet V2-B0...")
225        try:
226            models['efficientnet_v2_b0'] = get_efficientnetv2(
227                input_shape=DATA_SHAPE,
228                model_type='B0',
229                num_classes=NUM_CLASSES,
230                dropout=0.2,
231                pretrained_weights=None  # Train from scratch
232            )
233        except Exception as e:
234            print(f"Error creating EfficientNet: {e}")
235
236        # MobileNetV1
237        print("Creating MobileNetV1...")
238        try:
239            models['mobilenet_v1'] = get_mobilenetv1(
240                input_shape=DATA_SHAPE,
241                alpha=0.25,
242                num_classes=NUM_CLASSES,
243                dropout=0.2,
244                pretrained_weights=None
245            )
246        except Exception as e:
247            print(f"Error creating MobileNetV1: {e}")
248
249        # MobileNetV2
250        print("Creating MobileNetV2...")
251        try:
252            models['mobilenet_v2'] = get_mobilenetv2(
253                input_shape=DATA_SHAPE,
254                alpha=0.35,
255                num_classes=NUM_CLASSES,
256                dropout=0.2,
257                pretrained_weights=None
258            )
259        except Exception as e:
260            print(f"Error creating MobileNetV2: {e}")
261
262        # FDMobileNet
263        print("Creating FDMobileNet...")
```

```python
264        try:
265            models['fdmobilenet'] = get_fdmobilenet(
266                input_shape=DATA_SHAPE,
267                num_classes=NUM_CLASSES,
268                alpha=0.25,
269                dropout=0.2
270            )
271        except Exception as e:
272            print(f"Error creating FDMobileNet: {e}")
273
274        # ST EfficientNet LC V1
275        print("Creating ST EfficientNet LC V1...")
276        try:
277            models['st_efficientnet_lc_v1'] = get_st_efficientnet_lc_v1(
278                input_shape=DATA_SHAPE,
279                num_classes=NUM_CLASSES,
280                dropout=0.2
281            )
282        except Exception as e:
283            print(f"Error creating ST EfficientNet LC: {e}")
284
285        # ST FDMobileNet V1
286        print("Creating ST FDMobileNet V1...")
287        try:
288            models['st_fdmobilenet_v1'] = get_st_fdmobilenet_v1(
289                input_shape=DATA_SHAPE,
290                num_classes=NUM_CLASSES,
291                dropout=0.2
292            )
293        except Exception as e:
294            print(f"Error creating ST FDMobileNet: {e}")
295
296        return models
297
298
299    def main():
300        os.makedirs(MODELS_DIR, exist_ok=True)
301        train_data, val_data = load_mnist_data()
302
303        # Create models
304        print("\n" + "="*60)
305        print("Creating CNN Models")
306        print("="*60)
307        models = get_all_models()
308
309        # Train each model
310        results = {}
311        val_images, val_labels = val_data
312
313        for model_name, model in models.items():
314            try:
315                trained_model, history, accuracy = train_model(model,
                    model_name, train_data, val_data)
316
317                # Save training curves
318                plot_training_history(history, model_name)
319
320                # Compute metrics and save confusion matrix
321                acc, prec = evaluate_and_plot_metrics(trained_model,
                    model_name, val_images, val_labels)
322                results[model_name] = {'accuracy': acc, 'precision': prec}
323
324                # Clear session to free memory
325                keras.backend.clear_session()
```

```
326
327        except Exception as e:
328            print(f"Error training {model_name}: {e}")
329            results[model_name] = {'accuracy': 0.0, 'precision': 0.0}
330
331    # Print final results
332    print("\n" + "="*60)
333    print("TRAINING RESULTS SUMMARY")
334    print("="*60)
335    print(f"{'Model':<25} {'Accuracy':>10} {'Precision':>10}")
336    print("-"*47)
337    for model_name, metrics in sorted(results.items(), key=lambda x: x
           [1]['accuracy'], reverse=True):
338        print(f"{model_name:<25} {metrics['accuracy']*100:>9.2f}% {
               metrics['precision']*100:>9.2f}%")
339
340    print(f"\nModels saved to: {MODELS_DIR}")
341    return results
342
343 if __name__ == "__main__":
344    main()
```

**Quantization and Conversion:** Initially, we attempted to convert the models directly using y5 to tflite workflow. However, we found that without quantization, the model sizes were too large to fit into the limited Flash memory of the STM32F446RE. Therefore, we implemented a quantization script that converts the Keras models (.h5) to TensorFlow Lite (.tflite) format while applying full integer quantization ().

The quantization process uses a representative dataset drawn from the MNIST training data to estimate the dynamic range of activations. This ensures that the conversion from floating-point to integer does not significantly degrade accuracy.

**Listing 2.2:** Quantization script (quantize_models.py)

```
1  import os
2  import sys
3  import numpy as np
4  import tensorflow as tf
5  from tensorflow import keras
6
7  # Configuration
8  MODELS_DIR = os.path.join(os.path.dirname(__file__), 'models')
9  DATA_SHAPE = (32, 32, 3)
10
11 def prepare_tensor(images, out_shape):
12     images = tf.expand_dims(images, axis=-1)
13     images = tf.repeat(images, 3, axis=-1)
14     images = tf.image.resize(images, out_shape[:2])
15     images = images / 255.0
16     return images
17
18 def get_representative_dataset():
19     print("Loading MNIST for representative dataset...")
20     (train_images, _), _ = tf.keras.datasets.mnist.load_data()
21
22     # Preprocess a subset of images
```

```python
23      # We use 100 samples as recommended for representative datasets
24      num_calibration_steps = 100
25
26      # Shuffle and pick a subset
27      indices = np.random.choice(len(train_images), num_calibration_steps,
            replace=False)
28      calibration_images = train_images[indices]
29
30      # Preprocess
31      print("Preprocessing representative data...")
32      calibration_images = prepare_tensor(calibration_images, DATA_SHAPE)
33
34      def representative_data_gen():
35          for input_value in calibration_images:
36              # Model expects [1, 32, 32, 3]
37              input_value = tf.expand_dims(input_value, axis=0)
38              yield [input_value]
39
40      return representative_data_gen
41
42  def quantize_model(h5_path):
43      base_name = os.path.splitext(os.path.basename(h5_path))[0]
44      output_path = os.path.join(os.path.dirname(h5_path), f'{base_name}
            _quant.tflite')
45
46      try:
47          print(f"\nProcessing: {base_name}")
48
49          # Load Keras model
50          model = keras.models.load_model(h5_path)
51
52          # Create converter
53          converter = tf.lite.TFLiteConverter.from_keras_model(model)
54
55          # Set optimization flags
56          converter.optimizations = [tf.lite.Optimize.DEFAULT]
57
58          # Set representative dataset
59          converter.representative_dataset = get_representative_dataset()
60
61          # Ensure full integer quantization
62          converter.target_spec.supported_ops = [tf.lite.OpsSet.
            TFLITE_BUILTINS_INT8]
63
64          # Convert
65          print("Converting with quantization...")
66          tflite_model = converter.convert()
67
68          # Save
69          with open(output_path, 'wb') as f:
70              f.write(tflite_model)
71
72          original_size = os.path.getsize(h5_path) / 1024
73          quant_size = len(tflite_model) / 1024
74
75          print(f"Saved: {output_path}")
76          print(f"Original H5 size: {original_size:.1f} KB")
77          print(f"Quantized size:   {quant_size:.1f} KB")
78          print(f"Reduction ratio:  {original_size/quant_size:.2f}x")
79
80          return {
81              'name': base_name,
82              'original_size': original_size,
83              'quant_size': quant_size
```

```
84                }
85
86        except Exception as e:
87            print(f"Error converting {base_name}: {e}")
88            return None
89
90    def main():
91        if not os.path.exists(MODELS_DIR):
92            print(f"Models directory not found: {MODELS_DIR}")
93            return
94
95        h5_files = [f for f in os.listdir(MODELS_DIR) if f.endswith('.h5')]
96
97        if not h5_files:
98            print("No .h5 files found to quantize.")
99            return
100
101       results = []
102       for f in h5_files:
103           h5_path = os.path.join(MODELS_DIR, f)
104           res = quantize_model(h5_path)
105           if res:
106               results.append(res)
107
108       # Summary
109       print("\n" + "="*60)
110       print("QUANTIZATION SUMMARY")
111       print("="*60)
112       print(f"{'Model':<25} {'Orig (KB)':>10} {'Quant (KB)':>12} {'Ratio':>8}")
113       print("-"*57)
114
115       for r in results:
116           ratio = r['original_size'] / r['quant_size']
117           print(f"{r['name']:<25} {r['original_size']:>10.1f} {r['quant_size']:>12.1f} {ratio:>8.2f}x")
118
119   if __name__ == "__main__":
120       main()
```

## Step 2: ST Edge AI Developer Cloud Optimization

After generating the quantized .tflite models, we utilized the **ST Edge AI Developer Cloud** [2] to validate their deployability on the STM32F446RE. This platform automates the complex task of mapping neural network operators to the specific hardware accelerator and memory architecture of the target MCU.

**Model Analysis and Validation:** We uploaded the quantized models to the cloud platform and selected the **NUCLEO-F446RE** as the target board. The service analyzes the model's structure to determine the required Flash memory (for weights) and RAM (for activation buffers).
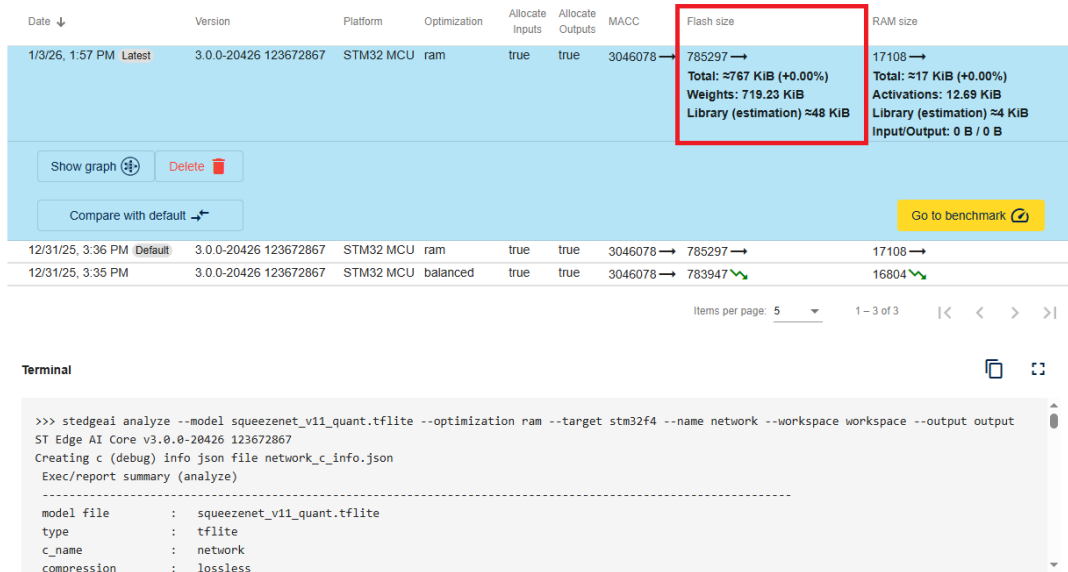
14

**Figure 2.1:** ST Edge AI Developer Cloud analysis dashboard showing SqueezenetV11 memory usage exceeds in flash size.

The analysis revealed a critical hardware limitation. Despite having sufficient RAM for the activation buffers, the binary size for both **EfficientNet** and **SqueezeNet** exceeded the 512kB Flash memory limit of the STM32F446RE. As a result, these models could not be deployed.

**Project Generation:** The remaining models **ResNetV1**, **MobileNetV1**, **FD-MobileNetV1**, and **ST-FD-MobileNetV1** successfully passed the memory validation checks. For these models, we ran the cloud benchmark to estimate inference latency and then generated the optimized C-code. The platform provided a downloadable STM32CubeIDE project ('.zip') containing the X-CUBE-AI library pre-configured with the network weights and runtime environment.
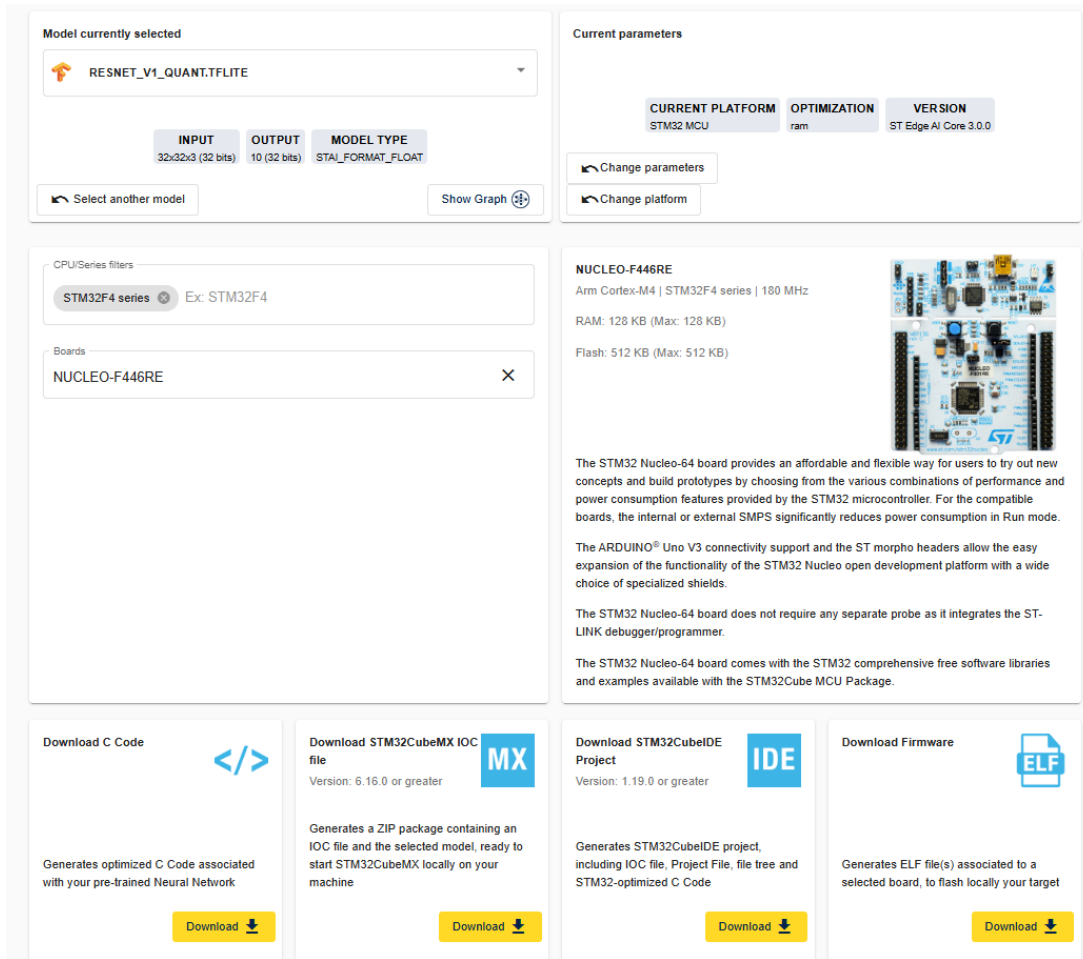
**Figure 2.2:** Successful generation of the STM32CubeIDE project for ResNetV1.

## Step 3: STM32 Firmware Implementation

The generated STM32CubeIDE projects contain the X-CUBE-AI library, which provides the runtime environment for the neural networks. However, the default `main.c` file requires modification to handle real-time data ingestion and specific input preprocessing requirements.

We implemented a unified firmware structure that works across all deployed models (ResNetV1, MobileNetV1, etc.). The core logic involves three key components: Preprocessing, Initialization, and the Inference Loop.

**Input Preprocessing:** The models were trained on 32x32 RGB images, but the incoming data from the UART interface is raw 28x28 grayscale (MNIST format). We implemented a `PrepareInput` function to bridge this gap. This function performs:

1. **Resizing:** Bilinear interpolation from 28x28 to 32x32 pixels.

16

2. **Normalization:** Scaling pixel values from [0, 255] to [0.0, 1.0].

3. **Channel Expansion:** Replicating the single grayscale channel into 3 RGB channels to match the network's input tensor shape (32, 32, 3).

**Listing 2.3:** Preprocessing function handling resize and normalization

```
/* Bilinear interpolation resize 28x28 -> 32x32, grayscale -> RGB,
    normalize */
static void PrepareInput(uint8_t *img28, float *in32) {
  int x32, y32, c;
  float scale = 28.0f / 32.0f;

  for (y32 = 0; y32 < IMG_SIZE_32; y32++) {
    for (x32 = 0; x32 < IMG_SIZE_32; x32++) {
      /* Map 32x32 coord to 28x28 */
      float src_x = x32 * scale;
      float src_y = y32 * scale;

      int x0 = (int)src_x;
      int y0 = (int)src_y;
      int x1 = (x0 < IMG_SIZE_28 - 1) ? x0 + 1 : x0;
      int y1 = (y0 < IMG_SIZE_28 - 1) ? y0 + 1 : y0;

      float dx = src_x - x0;
      float dy = src_y - y0;

      /* Bilinear interpolation */
      float p00 = img28[y0 * IMG_SIZE_28 + x0];
      float p10 = img28[y0 * IMG_SIZE_28 + x1];
      float p01 = img28[y1 * IMG_SIZE_28 + x0];
      float p11 = img28[y1 * IMG_SIZE_28 + x1];

      float val = p00 * (1 - dx) * (1 - dy) +
                  p10 * dx * (1 - dy) +
                  p01 * (1 - dx) * dy +
                  p11 * dx * dy;

      /* Normalize to [0,1] */
      float norm_val = val / 255.0f;

      /* Replicate to RGB channels (HWC format: 32x32x3) */
      int idx = (y32 * IMG_SIZE_32 + x32) * 3;
      for (c = 0; c < 3; c++) {
        in32[idx + c] = norm_val;
      }
    }
  }
}
```

**Network Initialization:**   The `AI_Init` function configures the neural network context. It links the model weights (stored in Flash) and the activation buffers (allocated in RAM) to the inference engine.

**Listing 2.4:** Network initialization using STAI API

```
static int AI_Init(void) {
```

```
 2    if (stai_network_init((stai_network *)net_ctx) != STAI_SUCCESS)
 3      return -1;
 4
 5    /* Set memory for activations */
 6    if (stai_network_set_activations((stai_network *)net_ctx, act_ptrs, 1)
          != STAI_SUCCESS)
 7      return -1;
 8
 9    /* Link input and output buffers */
10    if (stai_network_set_inputs((stai_network *)net_ctx, in_ptrs, 1) !=
          STAI_SUCCESS)
11      return -1;
12
13    if (stai_network_set_outputs((stai_network *)net_ctx, out_ptrs, 1) !=
          STAI_SUCCESS)
14      return -1;
15
16    return 0;
17  }
```

**Main Inference Loop:**   The main loop utilizes a blocking UART receive call to wait for a synchronization byte (0xBB) and the image payload (784 bytes). Upon reception, the image is processed, inference is executed synchronously, and the predicted class with the highest probability (ArgMax) is transmitted back to the PC.

**Listing 2.5:** Main execution loop

```
 1    /* ... Initialization code ... */
 2
 3    while (1) {
 4      /* Wait for Sync Byte */
 5      if (HAL_UART_Receive(&huart2, &sync, 1, HAL_MAX_DELAY) == HAL_OK &&
            sync == SYNC_BYTE) {
 6        HAL_GPIO_WritePin(LD2_GPIO_Port, LD2_Pin, GPIO_PIN_RESET); // LED
              OFF = Busy
 7
 8        /* Receive Raw Image */
 9        if (HAL_UART_Receive(&huart2, img_buf, IMG_BYTES, 5000) == HAL_OK)
              {
10
11          /* Preprocess: 28x28 -> 32x32x3 float */
12          PrepareInput(img_buf, in_buf);
13
14          /* Run Inference */
15          if (stai_network_run((stai_network *)net_ctx, STAI_MODE_SYNC) ==
                STAI_SUCCESS) {
16            res = ArgMax(out_buf, 10); /* Get Class Index */
17          } else {
18            res = 0xFF; /* Error */
19          }
20
21          /* Transmit Result */
22          HAL_UART_Transmit(&huart2, &res, 1, 100);
23        }
24        HAL_GPIO_WritePin(LD2_GPIO_Port, LD2_Pin, GPIO_PIN_SET); // LED ON
              = Ready
25      }
26    }
```

## Step 4: Python UART Client

To validate the deployed model, we developed a Python script ('mnist_uart.py') that acts as the interface between the PC and the STM32F446RE. This script handles image loading, preprocessing (resizing to 28x28, inverting colors if necessary for white-on-black digits), and UART communication.

The client sends a synchronization byte ('0xBB') followed by the raw pixel data (784 bytes) to the board. It then waits for the predicted digit index to be returned.

**Listing 2.6:** Python UART client for MNIST inference (mnist_uart.py)

```python
import argparse
import numpy as np
import cv2
import serial

SYNC_BYTE = 0xBB

def load_image(path):
    img = cv2.imread(path, cv2.IMREAD_GRAYSCALE)
    if img is None:
        raise ValueError(f"Cannot load: {path}")
    img = cv2.resize(img, (28, 28))
    # Invert colors if the image is black-on-white (like standard
        writing)
    # MNIST requires white-on-black
    if np.mean(img) > 127:
        img = 255 - img
    return img

def send_and_receive(img, port='COM9', baudrate=115200):
    try:
        with serial.Serial(port, baudrate, timeout=5) as ser:
            ser.write(bytes([SYNC_BYTE]))
            ser.write(img.flatten().tobytes())
            print(f"Sent 784 bytes")
            result = ser.read(1)
            return result[0] if len(result) == 1 else -1
    except Exception as e:
        print(f"Error: {e}")
        return -1

def display(digit, img=None):
    canvas = np.zeros((600, 800, 3), dtype=np.uint8)
    cv2.putText(canvas, str(digit), (280, 400), cv2.FONT_HERSHEY_SIMPLEX
        , 15, (0, 255, 0), 30)
    cv2.putText(canvas, "Recognized Digit", (150, 60), cv2.
        FONT_HERSHEY_SIMPLEX, 1.5, (255,255,255), 3)
    if img is not None:
        # Show the input image in the corner
        canvas[80:220, 20:160] = cv2.cvtColor(
            cv2.resize(img, (140, 140), interpolation=cv2.INTER_NEAREST)
                ,
            cv2.COLOR_GRAY2BGR
        )
    cv2.imshow('Result', canvas)
    cv2.waitKey(0)
    cv2.destroyAllWindows()
```

```
44
45  def main():
46      parser = argparse.ArgumentParser()
47      parser.add_argument('image', help='Path to the image file')
48      parser.add_argument('--port', default='COM9', help='Serial port (
            default: COM9)')
49      args = parser.parse_args()
50
51      img = load_image(args.image)
52      print(f"Sending 28x28 image to MCU...")
53      digit = send_and_receive(img, args.port)
54
55      if 0 <= digit <= 9:
56          print(f"*** Recognized: {digit} ***")
57          display(digit, img)
58      else:
59          print("Failed or Timed Out")
60
61  if __name__ == '__main__':
62      main()
```

**Usage:** The script is executed from the command line. It requires the path to the image file and optionally accepts a specific COM port. The recognized digit is displayed in an OpenCV window along with a visual confirmation.

Example command:

```
python mnist_uart.py two.png --port COM9
```

## Results

## Model Performance Analysis

All six trained models achieved high convergence on the MNIST dataset, demonstrating validation accuracies exceeding 95%. Figure 2.3 illustrates the training progression over 20 epochs. ResNetV1 and MobileNetV1 showed the most stable convergence, reaching optimal accuracy within the first 5 epochs.
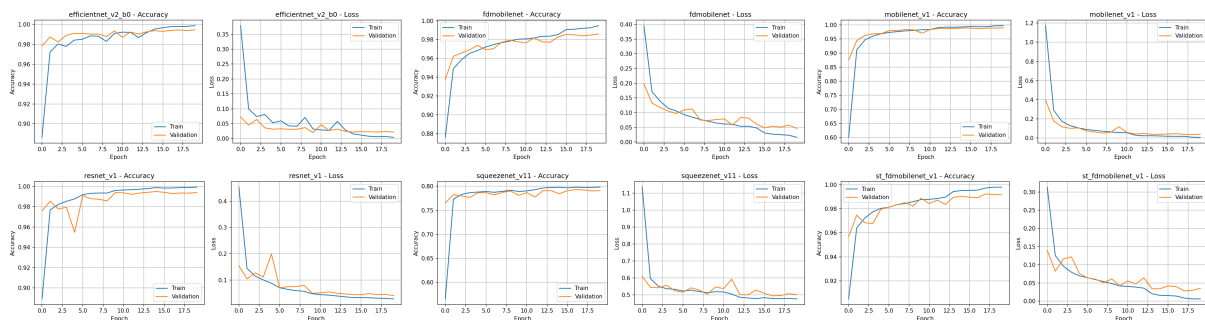


**Figure 2.3:** Training accuracy and loss curves for all evaluated models.

The confusion matrices in Figure 2.4 provide a deeper insight into classification performance. The diagonal dominance across all models confirms robust feature extraction. **ResNetV1** and **EfficientNetV2** exhibited the fewest off-diagonal errors, indicating superior precision. But the thing is SqueezeNet has some off labels as can be seen in the figure. That may be why we couldn't be able to fit into MCU because of some computational error that happened during the training.
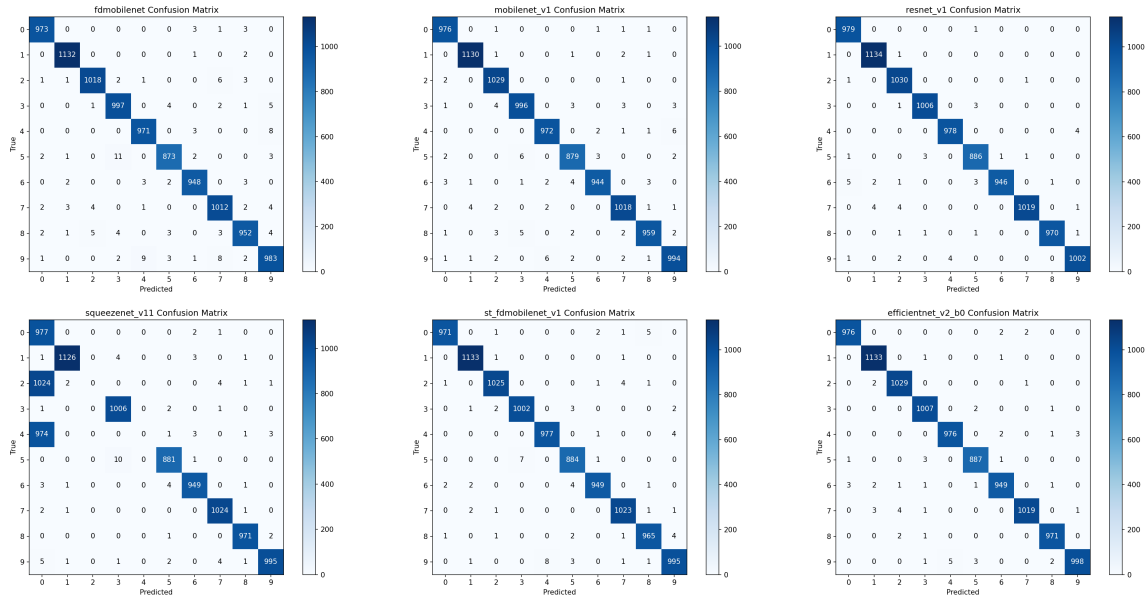


**Figure 2.4:** Confusion matrices comparing the classification performance of the six architectures.

## Hardware Resource Utilization

The critical constraint for this project was the 512kB Flash memory of the STM32F446RE. Table 2.1 summarizes the memory footprint of each model after quantization and optimization by the ST Edge AI Cloud.

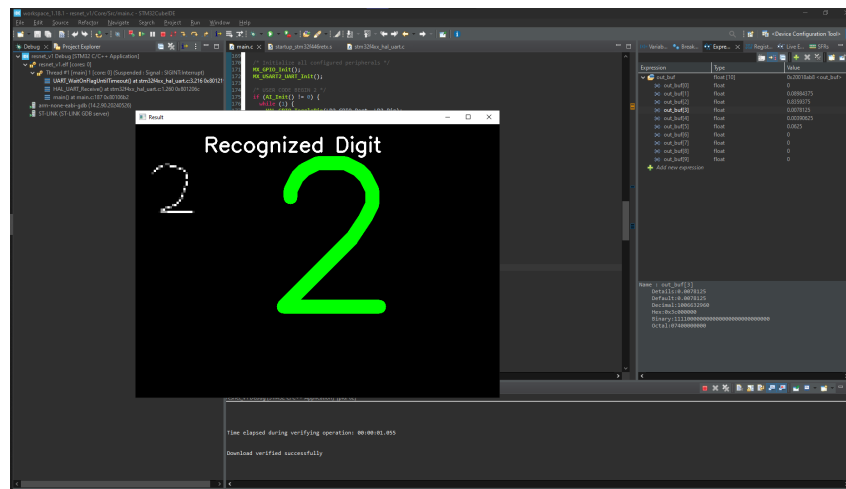**Table 2.1:** Resource Utilization and Deployment Status on STM32F446RE

| Model Architecture | Flash Usage (kB) | RAM Usage (kB) | Deployment Status |
|---|---|---|---|
| **ResNetV1** | **112** | **44** | **Success** |
| FDMobileNetV1 | 160 | 13 | Success |
| ST-FDMobileNetV1 | 178 | 13 | Success |
| MobileNetV1 | 252 | 13 | Success |
| EfficientNetV2-B0 | 6,000 | 84 | Failed (Flash Overflow) |
| SqueezeNetV1.1 | 767 | 17 | Failed (Flash Overflow) |

Despite their high accuracy, both EfficientNetV2 (6MB) and SqueezeNetV1.1 (767kB) exceeded the available Flash memory. While SqueezeNet is often cited as a "lightweight"
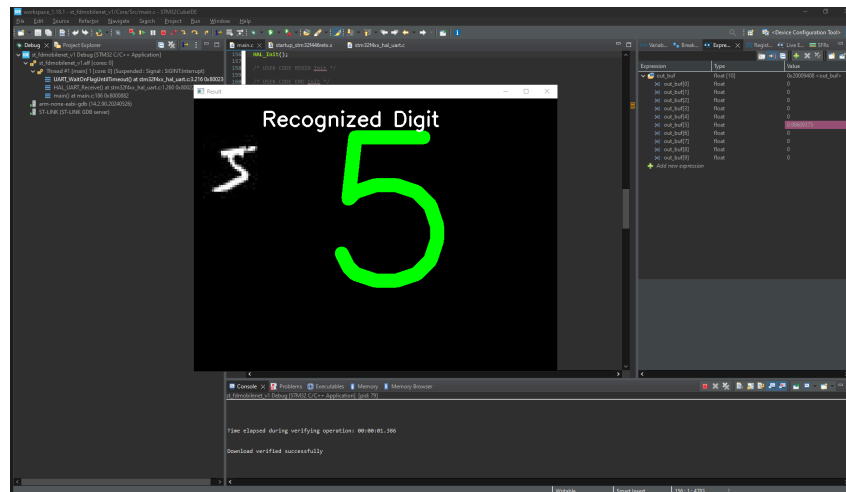
architecture, its parameter density proved too high for this specific microcontroller class without external memory. Conversely, ResNetV1 was the most balanced candidate, consuming only 112kB of Flash while maintaining high accuracy.

## Real-Time Inference Verification

The four compatible models were successfully deployed to the MCU. We verified real-time performance using the Python UART client. Figure 2.5 demonstrates the successful classification of handwritten digits sent from the PC to the STM32F446RE.

**(a)** ResNetV1 correctly classifying digit '2'.

**(b)** ST-FDMobileNetV1 correctly classifying digit '5'.

**Figure 2.5:** Real-time inference results from the STM32F446RE visualized via the Python client and showing computed probabilities of the label shown in STM32CubeIDE

# 3. CONCLUSION

In this homework, we successfully deployed deep CNN architectures onto the STM32F446RE MCU for handwritten digit recognition. The main goal was to evaluate different neural network topologies and implement a complete edge AI pipeline, navigating the strict memory constraints of a mid-range microcontroller.

We trained multiple architectures, including ResNetV1, MobileNetV1, SqueezeNet, and EfficientNet, using the MNIST dataset. Unlike previous assignments that relied on manual feature extraction like Hu Invariant Moments (Hu Moments), this work utilized deep learning models that process raw image data directly. We developed a Python client to transmit images via Universal Asynchronous Receiver Transmitter (UART), and implemented firmware on the MCU to perform real-time preprocessing (resizing and normalization) and inference.

A critical insight from this work was the impact of model complexity on hardware feasibility. While models like EfficientNetV2 and SqueezeNet offer high theoretical performance, we observed that they exceeded the 512kB Flash memory limit of the STM32F446RE, rendering them undeployable without external memory. Conversely, **ResNetV1** and **MobileNetV1**, combined with INT8 quantization, provided an optimal balance, fitting comfortably within the hardware resources while maintaining high classification accuracy.

We used the **ST Edge AI Developer Cloud** to bridge the gap between TensorFlow and the embedded C environment. This tool, along with our custom quantization scripts, allowed us to optimize model footprints significantly. This assignment highlighted the importance of hardware aware design, shows that successful edge AI deployment requires not just training accurate models but strictly adjusting to the physical constraints of the target device.

# BIBLIOGRAPHY

[1] C. Ünsalan, B. Höke, E. Atmaca, Embedded Machine Learning with Microcontrollers: Applications on STM32 Boards, Springer Nature, 2025.

[2] STMicroelectronics, St edge ai developer cloud, `https://stedgeai-dc.st.com` (2024).

[3] M. Abadi, et al., Tensorflow: Large-scale machine learning on heterogeneous systems, `https://www.tensorflow.org` (2015).