



# Embedded Systems Design

Masudul  
Imtiaz

# Embedded Systems Design

*Utilizing the Silicon Labs EFR32XG24 BLE Microcontroller*

Masudul Imtiaz, PhD

Aaron Storey, MSAI

Sigmond Kukla

& TA's of EE260

Coulson School of Engineering and Applied Sciences  
Clarkson University

# Table of contents

<b>Preface</b>	<b>i</b>
Welcome . . . . .	i
License . . . . .	i
 <b>Embedded Machine Learning</b>	 <b>ii</b>
<b>1 Introduction to Embedded Machine Learning</b>	<b>1</b>
1.1 Overview . . . . .	1
1.2 Real-World Applications of Embedded Machine Learning . . . . .	2
1.3 The EFR32MG24 for Machine Learning Applications . . . . .	3
 <b>2 The Art and Science of Machine Learning</b>	 <b>4</b>
2.1 Origins and Evolution . . . . .	4
2.1.1 Historical Context . . . . .	4
2.1.2 From Rule-Based to Learning Systems . . . . .	5
2.1.3 The Statistical Revolution . . . . .	6
2.2 Understanding Learning Systems . . . . .	7
2.2.1 The Learning Paradigm . . . . .	7
2.2.2 Learning from Data . . . . .	9
2.2.3 The Nature of Patterns . . . . .	11
2.3 The Nature of Machine Learning . . . . .	14
2.3.1 Learning as Induction . . . . .	14
2.3.2 The Role of Uncertainty . . . . .	16
2.3.3 Model Complexity and Regularization . . . . .	18
2.4 Building Learning Systems . . . . .	19
2.4.1 Data Preparation . . . . .	19
2.4.2 Model Selection and Training . . . . .	21
2.4.3 Model Evaluation and Deployment . . . . .	23
2.5 The Ethics and Governance of Machine Learning . . . . .	25
2.5.1 Fairness and Bias . . . . .	25
2.5.2 Transparency and Accountability . . . . .	25
2.5.3 Safety and Robustness . . . . .	26
2.5.4 Ethical Principles and Governance Frameworks . . . . .	27
2.6 Conclusion . . . . .	27
 <b>3 The “Hello World” of TinyML</b>	 <b>29</b>

3.1	Introduction to Microcontroller-Based Machine Learning . . . . .	29
3.2	Theoretical Foundations of TinyML for Microcontrollers . . . . .	29
3.2.1	The Computational Constraints Paradigm . . . . .	29
3.2.2	Model Compression and Quantization . . . . .	30
3.3	Building Our Sine Wave Model in Google Colab . . . . .	30
3.3.1	Generating and Processing the Dataset . . . . .	30
3.3.2	Constructing and Training the Neural Network Model . . . . .	31
3.4	Optimizing for Microcontroller Deployment . . . . .	33
3.4.1	Model Conversion and Quantization for TensorFlow Lite . . . . .	33
3.4.2	Converting to C Code for Embedded Systems . . . . .	33
3.5	Deploying with Simplicity Studio and Gecko SDK . . . . .	34
3.5.1	Creating a New Project in Simplicity Studio . . . . .	34
3.5.2	Project Structure and Important Files . . . . .	35
3.5.3	Adding Our Trained Model . . . . .	35
3.5.4	Implementing the Application Logic . . . . .	35
3.5.5	Creating the Model Integration File . . . . .	38
3.5.6	Building and Flashing the Application . . . . .	43
3.5.7	Observing the Results . . . . .	44
3.6	How it Works: Understanding the Implementation . . . . .	44
3.7	Extending the TinyML Application . . . . .	44
3.7.1	1. Adding Multiple LED Support . . . . .	44
3.7.2	2. Adding LCD Display Support . . . . .	45
3.7.3	3. Implementing Power Optimization . . . . .	46
3.7.4	4. Enhanced User Interface with Buttons . . . . .	47
3.7.5	5. Performance Profiling and Optimization . . . . .	47
3.8	Building TinyML Applications with the Gecko SDK . . . . .	48
3.8.1	Simplified Project Setup . . . . .	48
3.8.2	Streamlined Development Workflow . . . . .	48
3.8.3	Hardware-Specific Optimizations . . . . .	49
3.9	Conclusion . . . . .	49
<b>4</b>	<b>Building a TinyML Application</b> . . . . .	<b>50</b>
4.1	Understanding the Gecko SDK Approach to TinyML . . . . .	50
4.2	Setting Up Your Development Environment . . . . .	50
4.3	Creating a New Project in Simplicity Studio . . . . .	50
4.4	Exploring the Project Structure . . . . .	51
4.5	Importing the Sine Wave Model . . . . .	51
4.6	Implementing the Application . . . . .	51
4.7	Enhancing Output with PWM Control . . . . .	55
4.8	Building and Flashing the Application . . . . .	56
4.9	Debugging and Monitoring . . . . .	56
4.10	Optimizing TinyML Performance . . . . .	56
4.10.1	Memory Optimization . . . . .	56
4.10.2	Power Optimization . . . . .	57
4.10.3	Timing Performance . . . . .	58
4.11	Adding User Interaction with Buttons . . . . .	58
4.12	Enhanced Visualization with LCD (if available) . . . . .	59
4.13	Creating a Custom Component for TinyML . . . . .	60

4.14	Conclusion	64
<b>5</b>	<b>Handwriting Digit Recognition</b>	<b>65</b>
5.1	Chapter Objectives	65
5.2	Overview	65
5.3	Introduction	65
5.3.1	Challenges of Microcontroller Deployment	66
5.3.2	Chapter Objectives	66
5.4	Background & Related Work	66
5.4.1	TinyML: Machine Learning for Embedded Systems	66
5.4.2	Convolutional Neural Networks for Image Recognition	66
5.4.3	Model Optimization for Resource-Constrained Devices	67
5.5	Methodology	67
5.5.1	System Architecture	67
5.5.2	Model Design and Training	68
5.5.3	Model Optimization	69
5.5.4	Embedded Implementation	70
5.6	Implementation Details	71
5.6.1	Model Training Results	71
5.6.2	Model Quantization Effects	71
5.6.3	Embedded System Implementation	72
5.7	Results & Discussion	73
5.7.1	Classification Performance	73
5.7.2	Resource Utilization	73
5.7.3	Comparison with Cloud-Based Approaches	74
5.8	Challenges & Ethical Considerations	74
5.8.1	Technical Challenges	74
5.8.2	Ethical Considerations	75
5.9	Future Work & Conclusion	76
5.10	References	76
<b>6</b>	<b>IMU-Based Gesture Recognition</b>	<b>77</b>
6.1	Chapter Objectives	77
6.2	Introduction	77
6.3	System Architecture	77
6.4	Hardware Components	78
6.5	Model Design and Training	78
6.5.1	Dataset Preparation	78
6.5.2	CNN Architecture	79
6.5.3	Training Configuration	79
6.6	Model Optimization	80
6.6.1	Post-Training Quantization	80
6.7	Embedded Implementation	80
6.7.1	Development Environment and IMU Interface	80
6.7.2	Inference Pipeline	81
6.8	Implementation Details	82
6.8.1	Signal Processing and Sensor Fusion	83
6.8.2	Motion Detection Algorithm	83

6.9	Results & Discussion . . . . .	84
6.9.1	Classification Performance and Resource Utilization . . . . .	84
6.9.2	Comparison with Cloud-Based Approaches . . . . .	85
6.10	Technical Challenges and Solutions . . . . .	85
6.11	Future Directions . . . . .	85
6.12	Conclusion . . . . .	86
6.13	References . . . . .	86
<b>7</b>	<b>Real-Time Posture Detection Using Neural Networks</b>	<b>87</b>
7.1	Chapter Objectives . . . . .	87
7.2	Overview . . . . .	87
7.3	Introduction . . . . .	88
7.4	Hardware Configuration . . . . .	88
7.5	Development Environment . . . . .	88
7.6	Data Acquisition and Processing . . . . .	89
7.6.1	Data Collection Methodology . . . . .	89
7.6.2	Signal Processing and Feature Extraction . . . . .	89
7.7	Model Architecture and Training . . . . .	90
7.7.1	Neural Network Design . . . . .	90
7.7.2	Training Methodology . . . . .	91
7.7.3	Performance Evaluation . . . . .	91
7.8	Deployment Implementation . . . . .	92
7.8.1	Model Optimization Techniques . . . . .	92
7.8.2	Firmware Architecture . . . . .	93
7.8.3	Wireless Communication Interface . . . . .	93
7.9	Performance Analysis . . . . .	94
7.9.1	Experimental Evaluation . . . . .	94
7.9.2	Limitations and Challenges . . . . .	94
7.10	Conclusion . . . . .	95
7.11	References . . . . .	95

# Preface

## Welcome

Welcome to the “**System Design with Silicon Lab EFR32XG24 BLE Microcontroller**”. This book is designed to guide you through the process of programming, building applications, and integrating machine learning with the EFR32XG24 BLE Microcontroller. Whether you’re an engineering student or a seasoned professional, this book offers hands-on examples to make advanced concepts accessible.

You’ll learn how to: - Program the EFR32XG24 microcontroller using C. - Design and implement embedded systems applications. - Apply machine learning techniques to solve real-world problems. - Explore gesture recognition, anomaly detection, and audio-based ML solutions.

The book balances theory with practice, empowering readers to develop embedded systems that are robust, efficient, and intelligent.

If you’re interested in broader programming concepts or other machine learning platforms, we encourage you to explore additional resources and apply your learning across domains.

**i** This book was originally developed as part of the EE260 and EE513 courses at Clarkson University. The Quarto-based version serves as an example of modern technical publishing and open access education.

## License

This book is **free to use** under the [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International](#) License. You are welcome to share, adapt, and use the material for educational purposes, as long as proper attribution is given and no commercial use is made.

If you’d like to support the project or contribute, you can report issues or submit pull requests at [github.com/clarkson-edge/ee513\\_book](https://github.com/clarkson-edge/ee513_book). Thank you for helping improve this resource for the community.

# Embedded Machine Learning



## Chapter 1

# Introduction to Embedded Machine Learning

This chapter introduces the essential concepts of embedded machine learning and highlights the growing significance of TinyML in modern embedded system designs. It emphasizes the role of microcontrollers, particularly the Silicon Labs EFR32MG24, in enabling efficient, low-power machine learning inferencing for IoT applications. Whether you are a student starting your embedded ML journey or an engineer aiming to enhance your system design skills, this textbook will serve as a valuable resource to build innovative and efficient TinyML-enabled embedded solutions.

### 1.1 Overview

Embedded machine learning, often referred to as TinyML, represents a paradigm shift in computational intelligence by bringing sophisticated inferencing capabilities directly to resource-constrained embedded systems. Unlike traditional machine learning systems that rely on cloud computing or powerful edge devices, TinyML optimizes models to operate within the strict memory, processing, and power constraints of microcontrollers. This evolution enables a new class of intelligent devices that can make real-time decisions locally without requiring constant connectivity to external servers.

At the core of TinyML systems lies the microcontroller, a compact integrated circuit that combines a processor, memory, and input/output peripherals on a single chip. Modern microcontrollers like the Silicon Labs EFR32MG24 are increasingly designed with ML workloads in mind, featuring specialized hardware accelerators and optimized instruction sets that enhance neural network performance while maintaining energy efficiency.

In recent years, the demand for local intelligence in IoT devices has surged, driven by concerns about latency, privacy, bandwidth limitations, and power consumption. TinyML addresses these challenges by enabling machine learning models to run directly on microcontrollers, processing sensor data locally and making intelligent decisions without transmitting raw data to the cloud. This approach is particularly valuable for applications such as keyword spotting, gesture recognition, anomaly detection, and predictive maintenance in industrial settings.

The Silicon Labs EFR32MG24 series is one of the most advanced microcontrollers available for TinyML applications in 2024. Built on the ARM Cortex-M33 core operating at 78 MHz, it offers a powerful blend of performance and energy efficiency, with a memory footprint of 1536KB flash and 256KB RAM. The platform includes an AI/ML hardware accelerator that enhances neural network execution, making it ideal for deploying sophisticated TinyML models while maintaining battery life in portable devices.

This textbook, *Embedded Machine Learning Design with Silicon Labs EFR32MG24*, provides a comprehensive guide for students and engineers to understand and implement TinyML solutions. The book covers both theoretical foundations and practical implementations, ensuring readers gain a deep understanding of machine learning optimization for resource-constrained systems.

Throughout this book, readers will learn:

- The fundamentals of TinyML and the computational constraints paradigm
- Model compression and quantization techniques for microcontroller deployment
- Practical implementation using Google Colab for model training and Simplicity Studio for deployment
- Hands-on experience building the canonical “Hello World” of TinyML: a sine wave predictor
- Advanced techniques for power optimization, performance profiling, and model efficiency
- Real-world case studies demonstrating TinyML applications across various domains

## 1.2 Real-World Applications of Embedded Machine Learning

Embedded machine learning is transforming countless devices and technologies by enabling local intelligence in resource-constrained environments. TinyML systems execute sophisticated inferencing tasks efficiently while operating under strict constraints of power consumption, memory limitations, and processing capabilities. Examples of TinyML applications can be observed across diverse industries, showcasing the versatility and transformative potential of this technology.

In healthcare and wearables, TinyML enables continuous health monitoring without draining battery life. Smart watches and fitness trackers use embedded ML algorithms to detect irregular heartbeats, analyze sleep patterns, and recognize specific activities based on motion sensor data. These devices perform complex pattern recognition locally, only transmitting alerts or summarized insights rather than constant streams of raw data, preserving both battery life and user privacy.

Industrial IoT applications leverage TinyML for predictive maintenance and anomaly detection at the sensor level. Embedded microcontrollers equipped with ML capabilities can analyze vibration patterns from motors or machinery, detecting subtle changes that might indicate impending failure before catastrophic breakdowns occur. By processing this data directly on the device, these systems can operate in environments with limited connectivity while providing real-time insights.

Consumer electronics increasingly incorporate TinyML to enhance user experience through always-on, low-power intelligence. Voice assistants use keyword spotting models running on microcontrollers to detect wake words without sending all audio to the cloud. Smart home sensors employ ML algorithms to differentiate between routine movements and security concerns, reducing false alarms while improving response times to genuine threats.

Agricultural and environmental monitoring systems utilize TinyML to enable intelligent, autonomous operation in remote locations. Soil moisture sensors can incorporate local ML models to optimize irrigation schedules based on weather patterns, soil conditions, and crop-specific needs. Wildlife tracking devices use embedded ML to classify animal behaviors directly on the device, extending battery life from days to months by eliminating continuous data transmission.

The EFR32MG24 microcontroller is particularly well-suited for these applications due to its balance of processing power, memory resources, and energy efficiency. Its ARM Cortex-M33 core provides sufficient computational capabilities for running inference on neural networks, while its power management features enable long-term operation on battery power. The integrated ML accelerator further enhances performance for specific machine learning workloads, enabling more complex models to run efficiently.

### 1.3 The EFR32MG24 for Machine Learning Applications

The EFR32MG24 microcontroller, part of Silicon Labs' Wireless Gecko series, is specifically designed to address the growing demand for local machine learning capabilities in resource-constrained embedded systems. Built on the ARM Cortex-M33 core, it operates at a maximum frequency of 78 MHz, delivering sufficient computational power for real-time ML inferencing while maintaining energy efficiency. With 1536KB of flash memory and 256KB of RAM, it provides adequate storage for both program code and machine learning models after quantization and optimization.

A key feature that distinguishes the EFR32MG24 for ML applications is its dedicated AI/ML hardware accelerator, which enhances the execution of specific neural network operations. This accelerator enables more efficient matrix multiplications and other common ML computations, allowing for faster inference times and lower power consumption compared to software-only implementations. Combined with the DSP extensions in the Cortex-M33 architecture, this hardware support makes the EFR32MG24 an excellent platform for deploying sophisticated TinyML models.

The EFR32MG24 excels in power management, offering multiple low-power modes that are essential for battery-operated ML devices. Its Energy Management Unit (EMU) allows fine-grained control over active, sleep, and deep sleep states, enabling systems to run inferencing only when needed and remain in ultra-low-power states otherwise. This capability is critical for applications like smart sensors that may need to periodically analyze data but remain dormant most of the time.

For data acquisition and sensor integration, the EFR32MG24 provides comprehensive peripheral support, including high-precision ADCs, DACs, and various communication interfaces (UART, SPI, I2C). These peripherals enable the connection of diverse sensors for gathering the input data required by ML models. The microcontroller's wireless capabilities, particularly Bluetooth Low Energy (BLE), allow for convenient model updates, configuration changes, and the transmission of inference results when necessary.

Security features are increasingly important in ML-enabled devices, and the EFR32MG24 addresses this through hardware-based security elements including a cryptographic accelerator and secure boot mechanisms. These features help protect both the intellectual property embedded in the ML models and any sensitive data processed by the device.

The development environment for the EFR32MG24, centered around Simplicity Studio and the Gecko SDK, provides integrated support for TinyML workflows. The SDK includes optimized libraries for TensorFlow Lite Micro, enabling straightforward deployment of models trained using popular frameworks like TensorFlow. This integration streamlines the development process from model training to on-device deployment, making the platform accessible even to developers new to machine learning.

Available in the xG24-DK2601B Development Kit, the EFR32MG24 provides an ideal platform for learning and experimenting with embedded machine learning concepts, from simple inferencing tasks like our sine wave predictor to more complex applications such as sensor fusion, anomaly detection, and pattern recognition. Throughout this book, we will use this powerful yet resource-constrained platform to demonstrate the principles and practices of efficient TinyML implementation.

## Chapter 2

# The Art and Science of Machine Learning

Learning lies at the heart of intelligence, whether natural or artificial. In this chapter, we will embark on a fascinating exploration of the fundamental principles that enable machines to learn from experience. Together, we will examine both the theoretical foundations that provide a rigorous mathematical basis for machine learning, as well as the practical considerations that shape the design and implementation of modern learning systems. Our journey will take us from the historical roots of the field through to the cutting-edge research defining the current state of the art and the open challenges guiding future directions. By the end of this chapter, you will have built a comprehensive understanding of how machines can acquire, represent, and apply knowledge to solve complex problems and enhance decision-making across a wide range of domains.

To make our exploration as engaging and accessible as possible, I will aim to break down complex ideas into more easily digestible parts, building up gradually to the more advanced concepts. Along the way, I will make use of intuitive analogies, illustrative examples, and step-by-step explanations to help illuminate key points. Please feel free to ask questions or share your own insights at any point - learning is an interactive process and your contributions will only enrich our discussion!

With that in mind, let's begin our journey into the art and science of machine learning.

### 2.1 Origins and Evolution

To fully appreciate the current state and future potential of machine learning, it is helpful to understand its historical context and developmental trajectory. In this section, we will trace the origins of the field and highlight the pivotal advances that have shaped its evolution.

#### 2.1.1 Historical Context

The dream of creating intelligent machines that can learn and adapt has captivated the human imagination for centuries. In mythology and folklore around the world, we find stories of animated beings imbued with 'artificial' intelligence, from the golems of Jewish legends to the mechanical servants of ancient China. These ageless visions speak to a deep fascination with the idea of breathing life and cognizance into inanimate matter.

However, the emergence of machine learning as a scientific discipline is a more recent development, tracing its origins to the mid-20th century. In a profound sense, the birth of machine learning as we know it today arose from the convergence of several key intellectual traditions:

*Artificial intelligence* - The quest to create machines capable of intelligent behavior

*Statistics and probability theory* - The mathematical tools for quantifying and reasoning about uncertainty

*Optimization and control theory* - The principles for automated decision-making and goal-directed behavior

*Neuroscience and cognitive psychology* - The scientific study of natural learning in biological systems

Each of these tributaries contributed essential ideas and techniques that merged together to form the foundations of modern machine learning.

Some key milestones in the early history of the field:

- 1943 - Warren McCulloch and Walter Pitts publish “A Logical Calculus of the Ideas Immanent in Nervous Activity”, laying the groundwork for artificial neural networks
- 1950 - Alan Turing proposes the “Turing Test” in his seminal paper “Computing Machinery and Intelligence”, providing an operational definition of machine intelligence
- 1952 - Arthur Samuel writes the first computer learning program, which learned to play checkers better than its creator
- 1957 - Frank Rosenblatt invents the Perceptron, an early prototype of artificial neural networks capable of learning to classify visual patterns
- 1967 - Covering numbers and the Vapnik–Chervonenkis dimension (VC dimension) introduced in the groundbreaking work of Vladimir Vapnik and Alexey Chervonenkis, providing the foundations for statistical learning theory

These pioneering efforts laid the conceptual and technical groundwork for the subsequent decades of research that grew the field into the thriving discipline it is today.

### 2.1.2 From Rule-Based to Learning Systems

In its early stages, artificial intelligence research focused heavily on symbolic logic and deductive reasoning. The prevailing paradigm was that of “expert systems” - computer programs that encoded human knowledge and expertise in the form of explicit logical rules. A canonical example was MYCIN, a program developed at Stanford University in the early 1970s to assist doctors in diagnosing and treating blood infections. MYCIN’s knowledge base contained hundreds of IF-THEN rules obtained by interviewing expert physicians, such as:

```
IF (organism-1 is gram-positive) AND
   (morphology of organism-1 is coccus) AND
   (growth-conformation of organism-1 is chains)
THEN there is suggestive evidence (0.7) that
    the identity of organism-1 is streptococcus
```

By chaining together inferences based on these rules, MYCIN could arrive at diagnostic conclusions and treatment recommendations that rivaled those of human specialists in its domain.

However, the handcrafted knowledge-engineering approach of early expert systems soon ran into serious limitations:

- Knowledge acquisition bottleneck: Extracting and codifying expert knowledge proved to be extremely time-consuming and prone to inconsistencies and biases.
- Brittleness and inflexibility: Rule-based systems struggled to handle noisy data, adapt to novel situations, or keep up with changing knowledge.
- Opaque “black box” reasoning: The complex chains of inference generated by expert systems were often difficult for humans to inspect, understand, and debug.

- Inability to learn from experience: Once programmed, rule-based systems remained static and could not automatically improve their performance or acquire new knowledge.

These shortcomings highlighted the need for a fundamentally different approach - one that could overcome the rigidity and opacity of handcrafted symbolical rules and instead acquire knowledge directly from data.

### 2.1.3 The Statistical Revolution

The critical shift from rule-based to learning systems was catalyzed by two key insights: Many real-world domains are intrinsically uncertain and subject to noise, necessitating a probabilistic treatment. Expertise is often implicit and intuitive rather than explicit and axiomatic, making it more amenable to statistical extraction than symbolic codification.

Consider again the task of medical diagnosis that systems like MYCIN sought to automate. While it is possible to elicit a set of logical rules from a human expert, there are several complicating factors:

- Patients present with constellations of symptoms that are imperfectly correlated with underlying disorders.
- Diagnostic tests yield results with varying levels of accuracy and associated error rates.
- Diseases evolve over time, manifesting differently at different stages.
- Treatments have uncertain effects that depend on individual patient characteristics.
- New diseases emerge and existing ones change in their prevalence and manifestation over time. In such an environment, definitive logical rules are the exception rather than the norm. Instead, diagnosis is fundamentally a process of probabilistic reasoning under uncertainty, based on a combination of empirical observations and prior knowledge.

The key innovation that unlocked machine learning was to reframe the challenge in statistical terms:

- Instead of trying to manually encode deterministic rules, the goal became to automatically infer probabilistic relationships from observational data.
- Rather than requiring knowledge to be explicitly enumerated, learning algorithms aimed to implicitly extract latent patterns and regularities.
- In place of brittle logical chains, models learned robust statistical associations that could gracefully handle noise and uncertainty.

This shift in perspective opened up a powerful new toolbox of techniques at the intersection of probability theory and optimization. Some key formal developments:

- Maximum likelihood estimation (Ronald Fisher, 1920s): A principled framework for inferring the parameters of statistical models from observed data.
- The perceptron (Frank Rosenblatt, 1957): A simple type of artificial neural network capable of learning to classify linearly separable patterns.
- Stochastic gradient descent (Herbert Robbins & Sutton Monro, 1951): An efficient optimization procedure well-suited to large-scale machine learning problems.
- Backpropagation (multiple independent discoveries, 1970s-1980s): An algorithm for training multi-layer neural networks by propagating errors backwards through the network.
- The VC dimension (Vladimir Vapnik & Alexey Chervonenkis, 1960s-1970s): A measure of the capacity of a hypothesis space that quantifies the conditions for stable learning from finite data.
- Maximum margin classifiers and support vector machines (Vladimir Vapnik et al., 1990s): Powerful discriminative learning algorithms with strong theoretical guarantees.

Together, innovations like these provided the foundations for statistical learning systems that could effectively extract knowledge from raw data. They set the stage for the following decades of progress that would see machine learning mature into one of the most transformative technologies of our time.

## 2.2 Understanding Learning Systems

Having reviewed the historical context and key conceptual shifts behind the emergence of machine learning, we are now in a position to examine learning systems in greater depth. In this section, we will explore the fundamental principles that define the learning paradigm, the central role played by data, and the nature of the patterns that learning uncovers.

### 2.2.1 The Learning Paradigm

At its core, machine learning represents a radical departure from traditional programming approaches. To appreciate this, it is helpful to consider how we might go about solving a complex task such as object recognition using classical programming:

First, we would need to sit down and think hard about all the steps involved in identifying objects in images. We might come up with rules like:

- “an eye has a roughly circular shape”
- “a nose is usually located below the eyes and above the mouth”
- “a face is an arrangement of eyes, nose and mouth”, etc.

Next, we would translate these insights into specific programmatic instructions:

- “scan the image for circular regions”
- “check if there are two such regions in close horizontal proximity”
- “label these candidate eye regions”, etc.

We would then need to painstakingly debug and refine our program to handle all the edge cases and sources of variability we failed to consider initially.

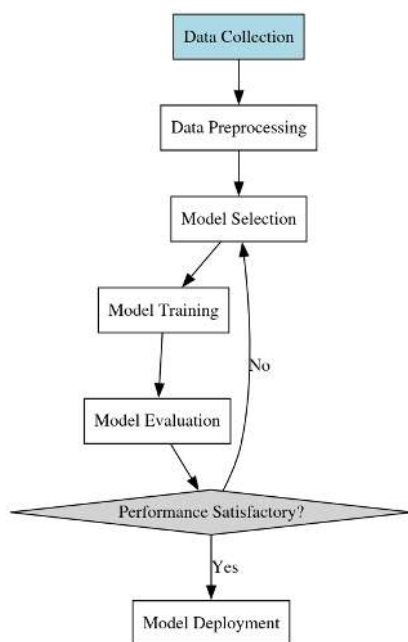
If our program needs to recognize additional object categories, we would have to return to step 1 and repeat the whole arduous process for each new class.

The classical approach places the entire explanatory burden on the human programmer - we start from a blank slate and must explicitly spell out every minute decision and edge case handling routine.

In contrast, the machine learning approach follows a very different recipe:

First, we collect a large dataset of labeled examples (e.g. images paired with the names of the objects they contain). We select a general-purpose model family that we believe has the capacity to capture the relevant patterns (e.g. deep convolutional neural networks for visual recognition). We specify a measure of success (e.g. what fraction of the images are labeled correctly) - this is our objective function.

We feed the dataset to a learning algorithm that automatically adjusts the parameters of the model so as to optimize the objective function on the provided examples. We evaluate the trained model on a separate test set to assess its ability to generalize to new cases.



Notice how the emphasis has shifted:

- Rather than having to explain “how” to solve the task, we provide examples of “what” we want and let the learning algorithm figure out the “how” for us.
- Instead of handcrafting detailed solution steps, we select a flexible model and offload the burden of tuning its parameters to an optimization procedure.
- In place of open-ended debugging, we can run controlled experiments to objectively measure generalization to unseen data.

For a more concrete illustration, consider how we might apply machine learning to the task of spam email classification:

```
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import accuracy_score

# 1. Collect labeled data
emails, labels = load_email_data()

# 2. Select a model family
vectorizer = CountVectorizer() # convert email text to word counts
classifier = MultinomialNB()   # naive Bayes with multinomial likelihood

# 3. Specify an objective function
```



```
def objective(model, X, y):  
    return accuracy_score(y, model.predict(X))  
  
# 4. Feed data to a learning algorithm  
# learn a spam classifier from 70% of the data  
train_emails, test_emails, train_labels, test_labels = train_test_split(  
    emails, labels, train_size=0.7, stratify=labels)  
X_train = vectorizer.fit_transform(train_emails)  
classifier.fit(X_train, train_labels)  
  
# 5. Evaluate generalization on held-out test set  
X_test = vectorizer.transform(test_emails)  
print("Test accuracy:", objective(classifier, X_test, test_labels))
```

This simple example illustrates the key ingredients:

- Data: A collection of example emails, along with human-provided labels indicating whether each one is spam or not.
- Model: The naive Bayes classifier, which specifies the general form of the relationship between the input features (word counts) and output labels (spam or not spam) in terms of probabilistic assumptions.
- Objective: The accuracy metric, which quantifies the quality of predictions made by the model in terms of the fraction of emails that are labeled correctly.
- Learning algorithm: The `fit` method of the classifier, which takes in the training data and finds the model parameters that maximize the likelihood of the observed labels.
- Generalization: The trained model is evaluated not on the data it was trained on, but rather on a separate test set that was held out during training. This provides an unbiased estimate of how well the model generalizes to new, unseen examples.

Of course, this is just a toy example intended to illustrate the basic flow. Real-world applications involve much larger datasets, more complex models, and more challenging prediction tasks. But the fundamental paradigm remains the same: by optimizing an objective function on a sample of training data, learning algorithms can automatically extract useful patterns and knowledge that generalize to novel situations.

## 2.2.2 Learning from Data

As the spam classification example makes clear, data plays a first-class role in machine learning. Indeed, one of the defining characteristics of the field is its focus on automatically extracting knowledge from empirical observations, rather than relying solely on human-encoded expertise. In this section, we take a closer look at how learning systems leverage data to acquire and refine their knowledge.

To begin, it is useful to clarify what we mean by “data” in a machine learning context. At the most basic level, a dataset is a collection of examples, where each example (also known as a “sample” or “instance”) provides a concrete instantiation of the task or phenomenon we wish to learn about. In the spam classification scenario, for instance, each example corresponds to an actual email message, along with a label indicating whether it is spam or not.

More formally, we can think of an example as a pair  $(x,y)$ , where:

- $x$  is a vector of input features that provide a quantitative representation of the relevant properties of the example. In the case of emails, the features might be counts of various words appearing in the message.
- $y$  is the target output variable that we would like to predict given the input features. For spam classification,  $y$  is a binary label, but in general it could be a continuous value (regression), a multi-class label (classification), or a more complex structure like a sequence or image.

A dataset, then, is a collection of  $n$  such examples:

$$D = (x_1, y_1), \dots, (x_n, y_n)$$

The goal of learning is to use the dataset  $D$  to infer a function  $f$  that maps from inputs to outputs:  $f: X \rightarrow Y$  such that  $f(x) \approx y$  for future examples  $(x, y)$  that were not seen during training.

With this formalism in mind, we can identify several key properties of data that are crucial for effective learning:

- **Representativeness:** To generalize well, the examples in  $D$  should be representative of the distribution of inputs that will be encountered in the real world. If the training data is systematically biased or skewed relative to the actual test distribution, the learned model may fail to perform well on new cases.
- **Quantity:** In general, more data is better for learning, as it provides a richer sampling of the underlying phenomena and helps the model to avoid overfitting to accidental regularities. The amount of data needed to achieve a desired level of performance depends on the complexity of the task and the expressiveness of the model class.
- **Quality:** The utility of data for learning can be undermined by issues like noise, outliers, and missing values. Careful data preprocessing, cleaning, and augmentation are often necessary to ensure that the model is able to extract meaningful signal.
- **Diversity:** For learning to succeed, the training data must contain sufficient variability along the dimensions that are relevant for the task at hand. If all the examples are highly similar, the model may fail to capture the full range of behaviors needed for robust generalization.
- **Labeling:** In supervised learning tasks, the quality and consistency of the output labels is critical. Noisy, ambiguous, or inconsistent labels can severely degrade the quality of the learned model.

To make these ideas more concrete, let's return to the spam classification example. Consider the following toy dataset:

```
train_emails = [
    "Subject: You won't believe this amazing offer!",
    "Subject: Request for project meeting",
    "Subject: URGENT: Update your information now!",
    "Hey there, just wanted to follow up on our conversation...",
    "Subject: You've been selected for a special promotion!",
]

train_labels = ["spam", "not spam", "spam", "not spam", "spam"]
```

Even without running any learning algorithms, we can identify some potential issues with this dataset:

- Small quantity: Only 5 examples is not enough to learn a robust spam classifier that covers the diversity of real-world emails. With so few examples, the model is likely to overfit to idiosyncratic patterns like the specific subject lines and fail to generalize well.
- Lack of diversity: The examples cover a very narrow range of email types (mainly short subject lines). A more representative sample would include a mix of subject lines, body text, sender information, etc. that better reflect the variability of real emails.
- Label inconsistency: On closer inspection, we might question whether the labeling is fully consistent. For instance, the 4th email seems potentially ambiguous - without more context about the content of the “conversation” it refers to, it’s unclear whether it should be classified as spam or not. Inconsistent labeling is a common source of problems in supervised learning.

To address these issues, we would want to collect a much larger and more diverse set of labeled examples. We might also need to do more careful data cleaning and preprocessing, for instance:

- Tokenizing the email text into individual words or n-grams
- Removing stop words, punctuation, and other low-information content
- Stemming or lemmatizing words to collapse related variants
- Normalizing features like word counts to avoid undue influence of message length
- Checking for and resolving inconsistencies or ambiguities in label assignments

In general, high-quality data is essential for successful learning. While it’s tempting to focus mainly on the choice of model class and learning algorithm, in practice the quality of the results is often determined by the quality of the data preparation pipeline.

*As the saying goes, “garbage in, garbage out”*\* - if the input data is full of noise, bias, and inconsistencies, no amount of algorithmic sophistication can extract meaningful patterns.\*

### 2.2.3 The Nature of Patterns

Having looked at the role of data in learning, let’s now turn our attention to the other central ingredient - the patterns that learning algorithms aim to extract. What exactly do we mean by “patterns” in the context of machine learning, and how do learning systems represent and leverage them?

In the most general sense, a pattern is any regularity or structure that exists in the data and captures some useful information for the task at hand. For instance, in spam classification, some relevant patterns might include:

- Certain words or phrases that are more common in spam messages than in normal emails (e.g. “special offer”, “free trial”, “no credit check”, etc.)
- Unusual formatting or stylistic choices that are suggestive of marketing content (e.g. excessive use of capitalization, colorful text, or images)
- Suspicious sender information, like mismatches between the stated identity and email address, or sending from known spam domains

A key insight of machine learning is that such patterns can be represented and manipulated mathematically, as operations in some formal space. For instance, the presence or absence of specific words can be encoded as a binary vector, with each dimension corresponding to a word in the vocabulary:

```
vocabulary = ["credit", "offer", "special", "trial", "won't", "believe", ...]

def email_to_vector(email):
    vector = [0] * len(vocabulary)
```

```

    for word in email.split():
        if word in vocabulary:
            index = vocabulary.index(word)
            vector[index] = 1
    return vector

# Example usage
message1 = "Subject: You won't believe this amazing offer!"
message2 = "Subject: Request for project meeting"

print(email_to_vector(message1))
# Output: [0, 1, 0, 0, 1, 1, ...]

print(email_to_vector(message2))
# Output: [0, 0, 0, 0, 0, 0, ...]

```

In this simple “bag of words” representation, each email is transformed into a vector that indicates which words from a predefined vocabulary are present in it. Already, some potentially useful patterns start to emerge - notice how the spam message gets mapped to a vector with more non-zero entries, suggesting the presence of marketing language.

Of course, this is a very crude representation that discards a lot of potentially relevant information (word order, punctuation, contextualized meanings, etc.). More sophisticated approaches attempt to preserve additional structure, for instance:

- Using counts or tf-idf weights instead of binary indicators to capture word frequencies
- Extracting  $n$ -grams (contiguous sequences of  $n$  words) to partially preserve local word order
- Applying techniques like latent semantic analysis or topic modeling to identify thematic structures
- Learning dense vector embeddings that map words and documents to points in a continuous semantic space

What these approaches all have in common is that they define a systematic mapping from the raw data (e.g. natural language text) to some mathematically tractable representation (e.g. vectors in a high-dimensional space). This mapping is where the “learning” in “machine learning” really takes place - by discovering the specific parameters of the mapping that lead to effective performance on the training examples, the learning algorithm implicitly identifies patterns that are useful for the task at hand.

To make this more concrete, let’s take a closer look at how a typical supervised learning algorithm actually goes about extracting patterns from data. Recall that the goal is to learn a function  $f : X \rightarrow Y$  that maps from input features to output labels, such that  $f(x) \approx y$  for examples  $(x, y)$  drawn from some underlying distribution.

In practice, most learning algorithms work by defining a parametrized function family  $F_\theta$  and searching for the parameter values  $\theta$  that minimize the empirical risk (i.e. the average loss) on the training examples:

$$\theta* = \operatorname{argmin}_\theta \frac{1}{n} \sum_i L(F_\theta(x_i), y_i)$$

Here  $L$  is a loss function that quantifies the discrepancy between the predicted labels  $F_\theta(x_i)$  and the true labels  $y_i$ , and the summation ranges over the  $n$  examples in the training dataset.

Different learning algorithms are characterized by the specific function families  $F$  and loss functions  $L$  that they employ, as well as the optimization procedure used to search for  $\theta^*$ . But at a high level, they all aim to find patterns - as captured by the parameters  $\theta$  - that enable the predictions  $F_\theta(x)$  to closely match the actual labels  $y$  across the training examples.

Let's make this more vivid by returning to the spam classification example. A very simple model family for this task is logistic regression, which learns a linear function of the input features:

$$F_\theta(x) = \sigma(\theta^T x)$$

Here  $x$  is the vector of word-presence features,  $\theta$  is a vector of real-valued weights, and  $\sigma$  is the logistic sigmoid function that "squashes" the linear combination  $\theta^T x$  to a value between 0 and 1 interpretable as the probability that the email is spam.

Coupled with the binary cross-entropy loss, the learning objective becomes:

$$\theta^* = \operatorname{argmin}_\theta \frac{1}{n} \sum_i [-y_i \log(F_\theta(x_i)) - (1 - y_i) \log(1 - F_\theta(x_i))]$$

where  $y_i \in \{0, 1\}$  indicates the true label (spam or not spam) for the  $i$ th training example.

Solving this optimization problem via a technique like gradient descent will yield a weight vector  $\theta^*$  such that:

- Weights for words that are more common in spam messages (like "offer" or "free") will tend to be positive, increasing the predicted probability of spam when those words are present.
- Weights for words that are more common in normal messages (like "meeting" or "project") will tend to be negative, decreasing the predicted probability of spam when those words are present.
- The magnitude of each weight corresponds to how predictive the associated word is of spam vs. non-spam - larger positive weights indicate stronger spam signals, while larger negative weights indicate stronger non-spam signals.

In this way, the learning process automatically discovers the specific patterns of word usage that are most informative for distinguishing spam from non-spam, as summarized in the weights  $\theta^*$ . Furthermore, the learned weights implicitly define a decision boundary in the high-dimensional feature space - emails that fall on one side of this boundary (as determined by the sign of  $\theta^T x$ ) are classified as spam, while those on the other side are classified as non-spam.

This simple example illustrates several key properties that are common to many learning algorithms:

The parameters  $\theta$  provide a compact summary of the patterns in the data that are relevant for the task at hand. In this case, they capture the correlations between the presence of certain words and the spam/non-spam label.

The learning process is data-driven - the specific values of the weights are determined by the empirical distribution of word frequencies in the training examples, not by any a priori assumptions or hand-coded rules.

The learned patterns are task-specific - the weights are tuned to optimize performance on the particular problem of spam classification, and may not be meaningful or useful for other tasks.

The expressiveness of the learned patterns is limited by the model family - in this case, the assumption of a simple linear relationship between word presence and spam probability. More complex model families (like deep neural networks) can capture richer, more nuanced patterns.

Of course, this is just a toy example intended to illustrate the basic principles. In practice, modern learning systems often employ much higher-dimensional feature spaces, more elaborate model families, and more sophisticated optimization procedures. But the fundamental idea remains the same -

by adjusting the parameters of a flexible model to minimize the empirical risk on a training dataset, learning algorithms can automatically discover patterns that generalize to improve performance on novel examples.

## 2.3 The Nature of Machine Learning

Having examined the fundamental components of learning systems - the data they learn from and the patterns they aim to extract - we now turn to some higher-level questions about the nature of learning itself. What does it mean for a machine to “learn” in the first place? How does this process differ from other approaches to artificial intelligence? And what challenges and opportunities does the learning paradigm present?

### 2.3.1 Learning as Induction

At a fundamental level, machine learning can be understood as a form of inductive inference - the process of drawing general conclusions from specific examples. In philosophical terms, this contrasts with deductive inference, which derives specific conclusions from general premises.

Consider a classic example of deductive reasoning:

- All men are mortal. (premise)
- Socrates is a man. (premise)
- Therefore, Socrates is mortal. (conclusion)

Here, the conclusion follows necessarily from the premises - if we accept that all men are mortal and that Socrates is a man, we must also accept that Socrates is mortal. The conclusion is guaranteed to be true if the premises are true.

Inductive reasoning, on the other hand, goes in the opposite direction:

- Socrates is a man and is mortal.
- Plato is a man and is mortal.
- Aristotle is a man and is mortal.
- Therefore, all men are mortal.

Here, the conclusion is not guaranteed to be true, even if all the premises are true - we can never be certain that the next man we encounter will be mortal, no matter how many examples of mortal men we have seen. At best, the conclusion is probable, with a degree of confidence that depends on the number and diversity of examples observed.

Machine learning can be seen as a form of algorithmic induction - instead of a human observer drawing conclusions from examples, we have a learning algorithm that discovers patterns in data and uses them to make predictions about novel cases. Just as with human induction, the conclusions of a machine learning model are never guaranteed to be true, but can be highly probable if the training data is sufficiently representative and the model family is appropriate for the task.

To make this more concrete, let's return to the spam classification example. Recall that our goal is to learn a function  $f$  that maps from email features  $x$  to spam labels  $y$ , such that  $f(x) \approx y$  for new examples  $(x, y)$  drawn from the same distribution as the training data.

In the logistic regression model we considered earlier,  $f$  takes the form:

$$f(x) = \sigma(\theta^T x)$$

where  $\theta$  is a vector of learned weights and  $\sigma$  is the logistic sigmoid function.

Now, imagine that we train this model on a dataset of 1000 labeled emails, using gradient descent to find the weights  $\theta^*$  that minimize the average cross-entropy loss on the training examples. We can then apply the learned function  $f^*$  to classify new emails as spam or not spam:

```
def predict_spam(email, weights):
    features = email_to_vector(email)
    score = weights.dot(features)
    probability = sigmoid(score)
    return probability > 0.5

# Example usage
weights = train_logistic_regression(train_emails, train_labels)

new_email = "Subject: Amazing opportunity to work from home!"
prediction = predict_spam(new_email, weights)

print(prediction) # Output: True
```

This process is fundamentally inductive:

- We start with a collection of specific examples (the training emails and their labels).
- We use these examples to learn a general rule (the weight vector  $\theta^*$ ) for mapping from inputs to outputs.
- We apply this rule to make predictions about new, unseen examples (e.g. classifying the new email as spam).

Just as with human induction, there is no guarantee that the predictions will be correct - the learned rule is only a generalization based on the limited sample of examples in the training data. If the training set is not perfectly representative of the real distribution of emails (which it almost never is), there will necessarily be some errors and edge cases that the model gets wrong.

However, if the inductive reasoning is sound - i.e. if the patterns discovered by the learning algorithm actually capture meaningful regularities in the data - then the model's predictions will be correct more often than not. Furthermore, as we train on larger and more diverse datasets, we can expect the accuracy and robustness of the learned patterns to improve, leading to better generalization performance.

Of course, spam classification is a relatively simple example as far as machine learning tasks go. In more complex domains like computer vision, natural language processing, or strategic decision-making, the input features and output labels can be much higher-dimensional and more abstract, the model families more elaborate and multilayered, and the optimization procedures more intricate and computationally intensive.

However, the fundamental inductive reasoning remains the same:

- Start with a set of training examples that (hopefully) capture relevant patterns and variations.
- Define a suitably expressive model family and objective function.
- Use an optimization algorithm to find the model parameters that minimize the objective on the training set.
- Apply the learned model to predict outputs for new, unseen inputs.
- Evaluate the quality of the predictions and iterate to improve the data, model, and optimization as needed.

The power of this paradigm lies in its generality - by framing the search for patterns as an optimization problem, learning algorithms can be applied to an extremely wide range of domains and tasks without the need for detailed domain-specific knowledge engineering. Given enough data and compute, the same basic approach can be used to learn patterns in images, text, speech, sensor readings, economic trends, user behavior, and countless other types of data.

At the same time, the generality of the paradigm also highlights some of its limitations and challenges:

- **Dependence on data quality:** The performance of a learning system is fundamentally limited by the quality and representativeness of its training data. If the data is noisy, biased, or incomplete, the learned patterns will reflect those limitations.
- **Opacity of learned representations:** The patterns discovered by learning algorithms can be highly complex and challenging to interpret. While simpler model families like linear regression produce relatively transparent representations, the internal structure of large neural networks is often inscrutable, making it difficult to understand how they arrive at their predictions.
- **Lack of explicit reasoning:** Learning systems excel at discovering statistical patterns, but struggle with the kind of explicit, logical reasoning that comes naturally to humans. Tasks that require careful deliberation, causal analysis, or manipulation of symbolic representations can be challenging to frame in purely statistical terms.
- **Potential for bias and fairness issues:** If the training data reflects societal biases or underrepresents certain groups, the learned models can perpetuate or even amplify those biases in their predictions. Careful auditing and debiasing of data and models is essential to ensure equitable outcomes.

Despite these challenges, the inductive learning paradigm has proven remarkably effective across a wide range of applications. In domains from medical diagnosis and scientific discovery to robotics and autonomous vehicles, machine learning systems are able to match or exceed human performance by discovering patterns that are too subtle or complex for manual specification. As the availability of data and computing power continues to grow, it's likely that the scope and impact of machine learning will only continue to expand.

### 2.3.2 The Role of Uncertainty

One of the most crucial things to understand about machine learning is that it is, at its core, a fundamentally probabilistic endeavor. When a learning algorithm draws conclusions from data, those conclusions are never absolutely certain, but rather statements of probability based on the patterns in the training examples.

Think back to our spam classification example. Even if our training dataset was very large and diverse, covering a wide range of both spam and legitimate emails, we can never be 100% sure that the patterns it captures will hold for every possible future email. There could always be some new type of spam that looks very different from what we've seen before, or some unusual legitimate email that happens to share many features with typical spam.

What a good machine learning model gives us, then, is not a definite classification, but a probability estimate. When we use logistic regression to predict the "spamminess" of an email, the model's output is a number between 0 and 1 that can be interpreted as the estimated probability that the email is spam, given its input features.

This might seem like a limitation compared to a deterministic rule-based system that always gives a definite yes-or-no answer. However, having a principled way to quantify uncertainty is actually



a key strength of the probabilistic approach. By explicitly representing the confidence of its predictions, a probabilistic model provides valuable information for downstream decision-making and risk assessment.

For instance, consider an email client that uses a spam filter to automatically move suspected spam messages to a separate folder. If the filter is based on a probabilistic model, we can set a confidence threshold for taking this action - say, only move messages with a 95% or higher probability of being spam. This allows us to trade off between false positives (legitimate emails moved to the spam folder) and false negatives (spam emails left in the main inbox) in a principled way.

More generally, having access to well-calibrated probability estimates opens up a range of possibilities for uncertainty-aware decision making, such as:

- Deferring to human judgment for borderline cases where the model is unsure
- Gathering additional information (e.g. asking the user for feedback) to resolve uncertainty
- Hedging decisions to balance risk and reward in the face of uncertain outcomes
- Combining predictions from multiple models to improve overall confidence

Of course, for these benefits to be realized, it's essential that the probability estimates produced by the model are actually well-calibrated - that is, they accurately reflect the true likelihood of the predicted outcomes. If a model consistently predicts 95% confidence for events that only occur 60% of the time, its uncertainty estimates are not reliable.

There are various techniques for quantifying and calibrating uncertainty in machine learning models, including:

- Explicit probability models: Some model families, like Bayesian networks or Gaussian processes, are designed to naturally produce probability distributions over outcomes. By incorporating prior knowledge and explicitly modeling sources of uncertainty, these approaches can provide principled uncertainty estimates.
- Ensemble methods: Techniques like bagging (bootstrap aggregating) and boosting involve training multiple models on different subsets or weightings of the data, then combining their predictions. The variation among the ensemble's predictions provides a measure of uncertainty.
- Calibration methods: Post-hoc calibration techniques like Platt scaling or isotonic regression can be used to adjust the raw confidence scores from a model to better align with empirical probabilities.
- Conformal prediction: A framework for providing guaranteed coverage rates for predictions, based on the assumption that the data are exchangeable. Conformal predictors accompany each prediction with a measure of confidence and a set of possible outcomes.

The importance of quantifying uncertainty goes beyond just improving decision quality - it's also crucial for building trust and promoting responsible use of machine learning systems. When a model accompanies its predictions with meaningful confidence estimates, users can make informed choices about when and how to rely on its outputs. This is especially important in high-stakes domains like healthcare or criminal justice, where the consequences of incorrect predictions can be serious.

Finally, reasoning about uncertainty is also central to more advanced machine learning paradigms like reinforcement learning and active learning:

- In reinforcement learning, an agent learns to make decisions by interacting with an environment and receiving rewards or penalties. Because the environment is often stochastic and the consequences of actions are uncertain, the agent must reason about the expected long-term value of different choices under uncertainty.

- In active learning, a model is allowed to interactively query for labels of unlabeled examples that would be most informative for improving its predictions. Selecting these examples requires estimating the expected reduction in uncertainty from obtaining their labels, based on the model's current state of knowledge.

As we continue to push the boundaries of what machine learning systems can do, the ability to properly quantify and reason about uncertainty will only become more essential. From building robust and reliable models to enabling effective human-AI collaboration, embracing uncertainty is key to unlocking the full potential of machine learning.

### 2.3.3 Model Complexity and Regularization

Another fundamental challenge in machine learning is striking the right balance between model complexity and generalization performance. On one hand, we want our models to be expressive enough to capture meaningful patterns in the data. On the other hand, we don't want them to overfit to noise or idiosyncrasies of the training set and fail to generalize to new examples.

This tradeoff is commonly known as the bias-variance dilemma:

- **Bias** refers to the error that comes from modeling assumptions and simplifications. A model with high bias makes strong assumptions about the data-generating process, which can lead to underfitting if those assumptions are wrong.
- **Variance** refers to the error that comes from sensitivity to small fluctuations in the training data. A model with high variance can fit the training set very well but may overfit to noise and fail to generalize to unseen examples.

As an analogy, think of trying to fit a curve to a set of scattered data points. A simple linear model has high bias but low variance - it makes the strong assumption that the relationship is linear, which limits its ability to fit complex patterns, but also makes it relatively stable across different subsets of the data. Conversely, a complex high-degree polynomial has low bias but high variance - it can fit the training points extremely well, but may wildly oscillate between them and make very poor predictions on new data.

In general, as we increase the complexity of a model (e.g. by adding more features, increasing the depth of a neural network, or reducing the strength of regularization), we decrease bias but increase variance. The goal is to find the sweet spot where the model is complex enough to capture relevant patterns but not so complex that it overfits to noise.

One way to control model complexity is through the choice of hypothesis space - the set of possible models that the learning algorithm can consider. A simple hypothesis space (like linear functions of the input features) will have low variance but potentially high bias, while a complex hypothesis space (like deep neural networks with millions of parameters) will have low bias but potentially high variance.

Another key tool for managing complexity is *regularization* - techniques for constraining the model's parameters or limiting its capacity to overfit. Some common regularization approaches include:

- **Parameter norm penalties:** Adding a penalty term to the loss function that encourages the model's weights to be small. L2 regularization (also known as weight decay) penalizes the squared Euclidean norm of the weights, while L1 regularization penalizes the absolute values. These penalties discourage the model from relying too heavily on any one feature.
- **Dropout:** Randomly "dropping out" (setting to zero) a fraction of the activations in a neural network during training. This prevents the network from relying too heavily on any one pathway and encourages it to learn redundant representations.

- **Early stopping:** Monitoring the model's performance on a validation set during training and stopping the optimization process when the validation error starts to increase, even if the training error is still decreasing. This prevents the model from overfitting to the training data.

The amount and type of regularization to apply is a key hyperparameter that must be tuned based on the characteristics of the data and the model. Too much regularization can lead to underfitting, while too little can lead to overfitting. Techniques like cross-validation and information criteria can help guide the selection of appropriate regularization settings.

It's worth noting that the bias-variance tradeoff and the role of regularization can vary depending on the amount of training data available. In the "classical" regime where the number of examples is small relative to the number of model parameters, regularization is essential for preventing overfitting. However, in the "modern" regime of very large datasets and overparameterized models (like deep neural networks with millions of parameters), the risk of overfitting is much lower, and the role of regularization is more subtle.

In fact, recent research has suggested that overparameterized models can exhibit "double descent" behavior, where increasing the model complexity beyond the point of interpolating the training data can actually improve generalization performance. This challenges the classical view of the bias-variance tradeoff and suggests that our understanding of model complexity and generalization is still evolving. Despite these nuances, the basic principles of managing model complexity and using regularization to promote generalization remain central to the practice of machine learning. As we train increasingly powerful models on ever-larger datasets, finding the right balance between expressiveness and constrainedness will be key to achieving robust and reliable performance.

## 2.4 Building Learning Systems

Now that we've explored some of the key theoretical principles behind machine learning, let's turn our attention to the practical considerations involved in building effective learning systems. What are the key components of a successful machine learning pipeline, and how do they fit together?

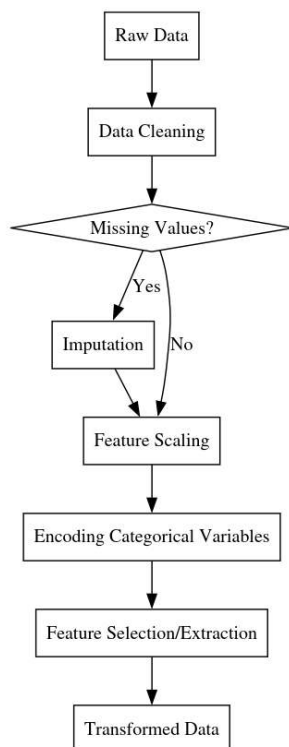
### 2.4.1 Data Preparation

The first and arguably most important step in any machine learning project is preparing the data. As we saw in Section 6.2.2, the quality and representativeness of the training data is essential for learning meaningful patterns that generalize well to new examples. No amount of algorithmic sophistication can make up for fundamentally flawed or insufficient data.

Key aspects of data preparation include:

- **Data cleaning:** Identifying and correcting errors, inconsistencies, and missing values in the raw data. This can involve steps like removing duplicate records, standardizing formats, and imputing missing values based on statistical patterns.
- **Feature engineering:** Transforming the raw input data into a representation that is more amenable to learning. This can involve steps like normalizing numeric features, encoding categorical variables, extracting domain-specific features, and reducing dimensionality.
- **Data augmentation:** Increasing the size and diversity of the training set by generating new examples through transformations of the existing data. This is especially common in domains like computer vision, where techniques like random cropping, flipping, and color jittering can help improve the robustness of the learned models.
- **Data splitting:** Dividing the data into separate sets for training, validation, and testing. The training set is used to fit the model parameters, the validation set is used to tune hyperparameters,

ters and detect overfitting, and the test set is used to evaluate the final performance of the model on unseen data.



The specifics of data preparation will vary depending on the domain and the characteristics of the data, but the general principles of ensuring data quality, representativeness, and suitability for learning are universal. It's often said that data preparation is 80% of the work in machine learning, and while this may be an exaggeration, it underscores the critical importance of getting the data right.

To make this more concrete, let's consider an example of data preparation in the context of a real-world problem. Suppose we're working on a machine learning system to predict housing prices based on features like square footage, number of bedrooms, location, etc. Our raw data might look something like this:

```

Address,Sq.Ft.,Beds,Baths,Price
123 Main St,2000,3,2.5,$500,000
456 Oak Ave,1500,2,1,"$350,000"
789 Elm Rd,1800,3,2,425000
  
```

To prepare this data for learning, we might perform the following steps:

- Standardize the 'Price' column to remove the "\$" and "," characters and convert to a numeric type.
- Impute missing values in the 'Beds' and 'Baths' columns (if any) with the median or most frequent value.

- Normalize the ‘Sq.Ft.’ column by subtracting the mean and dividing by the standard deviation.
- One-hot encode the ‘Address’ column into separate binary features for each unique location.
- Split the data into training, validation, and test sets in a stratified fashion to ensure representative price distributions in each split.

The end result might look something like this:

```
Sq.Ft._Norm,Beds,Baths,123_Main_St,456_Oak_Ave,789_Elm_Rd,...,Price
,3,2.5,1,0,0,...,500000
-0.58,2,1,0,1,0,...,350000
,3,2,0,0,1,...,425000
...
```

Of course, this is just a toy example, and in practice the data preparation process can be much more involved. The key point is that investing time and effort into carefully preparing the data is essential for building successful learning systems.

## 2.4.2 Model Selection and Training

Once the data is prepared, the next step is to select an appropriate model family and training procedure. As we saw in Section 6.3.3, this involves striking a balance between model complexity and generalization ability, often through a combination of cross-validation and regularization techniques.

Some key considerations in model selection include:

- *Inductive biases*: The assumptions and constraints that are built into the model architecture. For example, convolutional neural networks have an inductive bias towards translation invariance and local connectivity, which makes them well-suited for image recognition tasks.
- *Parameter complexity*: The number of learnable parameters in the model, which affects its capacity to fit complex patterns but also its potential to overfit to noise in the training data. Regularization techniques can help control parameter complexity.
- *Computational complexity*: The time and memory requirements for training and inference with the model. More complex models may require specialized hardware (like GPUs) and longer training times, which can be a practical limitation.
- *Interpretability*: The extent to which the learned model can be inspected and understood by humans. In some domains (like healthcare or finance), interpretability may be a key requirement for building trust and ensuring regulatory compliance.

The choice of model family will depend on the nature of the problem and the characteristics of the data. For structured data with clear feature semantics, “shallow” models like linear regression, decision trees, or support vector machines may be appropriate. For unstructured data like images, audio, or text, “deep” models like convolutional or recurrent neural networks are often used.

Once a model family is selected, the next step is to train the model on the prepared data. This typically involves the following steps:

- Instantiate the model with initial parameter values (e.g. random weights for a neural network).
- Define a loss function that measures the discrepancy between the model’s predictions and the true labels on the training set.
- Use an optimization algorithm (like stochastic gradient descent) to iteratively update the model parameters to minimize the loss function.

- Monitor the model's performance on the validation set to detect overfitting and tune hyperparameters.
- Stop training when the validation performance plateaus or starts to degrade.

The specifics of the training process will vary depending on the chosen model family and optimization algorithm, but the general goal is to find the model parameters that minimize the empirical risk on the training data while still generalizing well to unseen examples.

To illustrate these ideas, let's continue with our housing price prediction example. Suppose we've decided to use a regularized linear regression model of the form:

$$price = w_0 + w_1 * sqft + w_2 * beds + w_3 * baths + \dots$$

where  $w_0, w_1, \dots$  are the learned weights and `sqft`, `beds`, `baths`, `...` are the input features.

We can train this model using gradient descent on the mean squared error loss function:

```
def mse_loss(y_true, y_pred):
    return np.mean((y_true - y_pred) ** 2)

def gradient_descent(X, y, w, lr=0.01, num_iters=100):
    for i in range(num_iters):
        y_pred = np.dot(X, w)
        error = y_pred - y
        gradient = 2 * np.dot(X.T, error) / len(y)
        w -= lr * gradient
    return w

# Add a bias term to the feature matrix
X = np.c_[np.ones(len(X)), X]

# Initialize weights to zero
w = np.zeros(X.shape[1])

# Train the model
w = gradient_descent(X, y, w)
```

We can also add L2 regularization to the loss function to prevent overfitting:

```
def mse_loss_regularized(y_true, y_pred, w, alpha=0.01):
    return mse_loss(y_true, y_pred) + alpha * np.sum(w**2)

def gradient_descent_regularized(X, y, w, lr=0.01, alpha=0.01, num_iters=100):
    for i in range(num_iters):
        y_pred = np.dot(X, w)
        error = y_pred - y
        gradient = 2 * np.dot(X.T, error) / len(y) + 2 * alpha * w
        w -= lr * gradient
    return w
```

```
# Train the regularized model
w = gradient_descent_regularized(X, y, w)
```

The  $\alpha$  parameter controls the strength of the regularization - larger values will constrain the weights more strongly, while smaller values will allow the model to fit the training data more closely.

By tuning the learning rate  $1r$ , regularization strength  $\alpha$ , and number of iterations `num_iters`, we can find the model that achieves the best balance between fitting the training data and generalizing to new examples.

Of course, linear regression is just one possible model choice for this problem. We could also experiment with more complex models like decision trees, random forests, or neural networks, each of which would have its own set of hyperparameters to tune. The key is to use a combination of domain knowledge, empirical validation, and iterative refinement to find the model that best suits the problem at hand.

### 2.4.3 Model Evaluation and Deployment

Once we've trained a model that performs well on the validation set, the final step is to evaluate its performance on the held-out test set. This gives us an unbiased estimate of how well the model is likely to generalize to real-world data.

Common evaluation metrics for *classification problems* include:

- **Accuracy:** The fraction of examples that are correctly classified.
- **Precision:** The fraction of positive predictions that are actually positive.
- **Recall:** The fraction of actual positives that are predicted positive.
- **F1 Score:** The harmonic mean of precision and recall.
- **ROC AUC:** The area under the receiver operating characteristic curve, which measures the trade-off between true positive rate and false positive rate.

For *regression problems*, common metrics include:

- **Mean squared error (MSE):** The average squared difference between the predicted and actual values.
- **Mean absolute error (MAE):** The average absolute difference between the predicted and actual values.
- **R-squared ( $R^2$ ):** The proportion of variance in the target variable that is predictable from the input features.

It's important to choose an evaluation metric that aligns with the business goals of the problem. For example, in a fraud detection system, we might care more about recall (catching as many fraudulent transactions as possible) than precision (avoiding false alarms), while in a medical diagnosis system, we might care more about precision (avoiding false positives that could lead to unnecessary treatments).

If the model's performance on the test set is satisfactory, we can proceed to deploy it in a production environment. This involves integrating the trained model into a larger software system that can apply it to new input data and surface the predictions to end users.

Some key considerations in model deployment include:

- **Scalability:** Can the model handle the volume and velocity of data in the production environment? This may require techniques like batch processing, streaming, or distributed computation.

- Latency: How quickly does the model need to generate predictions in order to meet business requirements? This may require optimizations like model compression, quantization, or hardware acceleration.
- Monitoring: How will the model's performance be monitored and maintained over time? This may involve tracking key metrics, detecting data drift, and periodically retraining the model on fresh data.
- Security: How will the model and its predictions be protected from abuse or unauthorized access? This may involve techniques like input validation, output filtering, or access controls.

Deploying and maintaining machine learning models in production is a complex topic that requires close collaboration between data scientists, software engineers, and domain experts. It's an active area of research and development, with new tools and best practices emerging regularly.

To bring everything together, let's return one last time to our housing price prediction example. After training and validating our regularized linear regression model, we can evaluate its performance on the test set:

```
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

# Generate predictions on the test set
y_pred = np.dot(X_test, w)

# Calculate evaluation metrics
mse = mean_squared_error(y_test, y_pred)
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"Test MSE: {mse:.2f}")
print(f"Test MAE: {mae:.2f}")
print(f"Test R^2: {r2:.2f}")
```

If we're satisfied with the model's performance, we can deploy it as part of a larger housing price estimation service. This might involve:

- Wrapping the trained model in a web service API that can accept new housing features and return price predictions.
- Integrating the API with a user-facing application that allows homeowners or real estate agents to input property information and receive estimates.
- Setting up a data pipeline to continuously collect new housing data and periodically retrain the model to capture changing market conditions.
- Defining monitoring dashboards and alerts to track the model's performance over time and detect any anomalies or degradations.
- Establishing governance policies and processes for managing the lifecycle of the model, from development to retirement.

Again, this is a simplified example, but it illustrates the end-to-end process of building a machine learning system, from data preparation to model development to deployment and maintenance.



## 2.5 The Ethics and Governance of Machine Learning

As machine learning systems become more prevalent and powerful, it's crucial that we grapple with the ethical implications of their development and deployment. In this final section, we'll explore some key ethical considerations and governance principles for responsible machine learning.

### 2.5.1 Fairness and Bias

One of the most pressing ethical challenges in machine learning is ensuring that models are fair and unbiased. If the training data reflects historical biases or discrimination, the resulting model may perpetuate or even amplify those biases in its predictions.

For example, consider a hiring model that is trained on past hiring decisions to predict the likelihood of a candidate being successful in a job. If the training data comes from a company with a history of discriminatory hiring practices, the model may learn to penalize candidates from under-represented groups, even if those factors are not actually predictive of job performance.

Detecting and mitigating bias in machine learning systems is an active area of research, with techniques like:

- Demographically balancing datasets to ensure equal representation of different groups
- Adversarial debiasing to remove sensitive information from model representations
- Regularization techniques to penalize models that exhibit disparate impact
- Post-processing methods to adjust model outputs to satisfy fairness constraints

However, these techniques are not perfect, and there is often a tradeoff between fairness and accuracy. Moreover, fairness is not a purely technical issue, but a sociotechnical one that requires ongoing collaboration between machine learning practitioners, domain experts, policymakers, and affected communities.

### 2.5.2 Transparency and Accountability

Another key ethical principle for machine learning is transparency and accountability. As models become more complex and consequential, it becomes harder for humans to understand how they arrive at their predictions and to trace the provenance of their training data and design choices.

This opacity can make it difficult to audit models for bias, safety, or compliance with regulations. It can also make it harder to challenge or appeal decisions made by machine learning systems, leading to a loss of human agency and recourse.

Some techniques for promoting transparency and accountability in machine learning include:

- Model interpretability methods that provide human-understandable explanations of model predictions
- Provenance tracking to document the lineage of data, code, and models used in a system
- Audit trails and version control to enable reproducibility and historical analysis
- Participatory design processes that involve affected stakeholders in the development and governance of models

However, like fairness, transparency is not purely a technical problem. It also requires institutional structures and processes for oversight, redress, and accountability. This might involve things like:

- Designating responsible individuals or teams for the ethical development and deployment of machine learning systems
- Establishing review boards or oversight committees to assess the social impact and governance of models

- Creating channels for affected individuals and communities to provide input and feedback on the use of machine learning in their lives
- Developing legal and regulatory frameworks to enforce transparency and accountability standards

### 2.5.3 Safety and Robustness

As machine learning systems are deployed in increasingly high-stakes domains, from healthcare to transportation to criminal justice, ensuring their safety and robustness becomes paramount. Models that are brittle, unreliable, or easily fooled can lead to serious harms if they are not carefully designed and tested.

Some key challenges in machine learning safety and robustness include:

- Distributional shift, where models trained on one data distribution may fail unexpectedly when applied to a different distribution
- Adversarial attacks, where malicious actors can craft inputs that fool models into making egregious errors
- Reward hacking, where optimizing for the wrong objective function can lead models to behave in unintended and harmful ways
- Safe exploration, where models need to learn about their environment without taking catastrophic actions

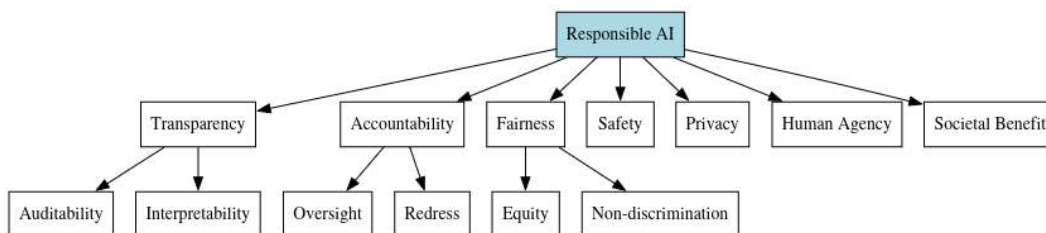
Techniques for improving the safety and robustness of machine learning systems include:

- Anomaly and out-of-distribution detection to flag inputs that are far from the training data
- Adversarial training and robustness regularization to make models more resilient to perturbations
- Constrained optimization and safe reinforcement learning to respect safety boundaries during learning
- Formal verification and testing to provide guarantees about model behavior under different conditions

However, building truly safe and robust machine learning systems requires more than just technical solutions. It also requires:

- Rigorous safety culture and practices throughout the development and deployment lifecycle
- Close collaboration between machine learning practitioners, domain experts, and safety professionals
- Proactive engagement with policymakers and the public to align the development of machine learning with societal values and expectations
- Ongoing monitoring and adjustment of deployed systems to catch and correct errors and unintended consequences

### 2.5.4 Ethical Principles and Governance Frameworks



To navigate the complex ethical landscape of machine learning, we need clear principles and governance frameworks to guide responsible development and deployment. Some key principles that have been proposed include:

- **Transparency:** Machine learning systems should be auditable and understandable by humans.
- **Accountability:** There should be clear mechanisms for oversight, redress, and enforcement.
- **Fairness:** Machine learning should treat all individuals equitably and avoid discriminatory impacts.
- **Safety:** Machine learning systems should be reliable, robust, and safe throughout their lifecycle.
- **Privacy:** The collection and use of data for machine learning should respect individual privacy rights and provide appropriate protections.
- **Human agency:** Machine learning systems should respect human autonomy and dignity, and provide meaningful opportunities for human input and control.
- **Societal benefit:** The development and deployment of machine learning should be guided by considerations of social good and collective wellbeing.

Translating these high-level principles into practical governance frameworks is an ongoing challenge, but some key elements include:

- Ethical codes of conduct and professional standards for machine learning practitioners
- Impact assessment and risk management processes to identify and mitigate potential harms
- Stakeholder engagement and participatory design to ensure affected communities have a voice
- Regulatory sandboxes and policy experiments to test new governance approaches
- International standards and coordination to address the global nature of machine learning development

Ultimately, the goal of machine learning governance should be to ensure that the technology is developed and deployed in a way that aligns with human values and enhances, rather than undermines, human flourishing. This is a complex and ongoing process that will require sustained collaboration across disciplines, sectors, and geographies.

## 2.6 Conclusion

In this chapter, we've embarked on a comprehensive exploration of the foundations of machine learning, from its historical roots to its modern techniques to its future challenges. We've seen how the field has evolved from rule-based expert systems to data-driven statistical learning, powered by the

explosion of big data and computing power. We've examined the fundamental components of machine learning systems - the data they learn from, the patterns they aim to extract, and the algorithms that power the learning process. We've discussed key concepts like inductive bias, generalization, overfitting, and regularization, and how they relate to the art of building effective models. We've walked through the practical steps of constructing a machine learning pipeline, from data preparation to model selection to deployment and monitoring. And we've grappled with some of the ethical challenges and governance principles that arise when building systems that can have significant impact on people's lives.

The field of machine learning is still rapidly evolving, with new breakthroughs and challenges emerging every year. As we look to the future, some of the key frontiers and open questions include:

- *Continual and lifelong learning*: How can we build models that can learn continuously and adapt to new tasks and domains over time, without forgetting what they've learned before?
- *Causality and interpretability*: How can we move beyond purely associational patterns to uncover causal relationships and build models that are more interpretable and explainable to humans?
- *Robustness and safety*: How can we guarantee that machine learning systems will behave safely and reliably, even in the face of distributional shift, adversarial attacks, or unexpected situations?
- *Human-AI collaboration*: How can we design machine learning systems that augment and empower human intelligence, rather than replacing or undermining it?
- *Ethical alignment*: How can we ensure that the development and deployment of machine learning aligns with human values and promotes beneficial outcomes for society as a whole?

Advancing machine learning requires collaboration across disciplines—from computer science and statistics to psychology, social science, philosophy, and ethics. It also demands engagement with policymakers, industry leaders, and the public to ensure responsible and inclusive development. The goal is to build systems that learn from experience to make decisions that benefit humanity, whether in healthcare, scientific discovery, or improving daily life.

However, realizing this potential goes beyond technical progress; it requires addressing fairness, accountability, transparency, and safety while navigating ethical and governance challenges. Machine learning practitioners must not only push technological boundaries but also consider the broader impact of their work. By embracing diverse perspectives and collaborating beyond our field, we shape the future of AI. Staying curious, critical, and committed to responsible development will ensure machine learning serves society for generations to come.

## Chapter 3

# The “Hello World” of TinyML

### 3.1 Introduction to Microcontroller-Based Machine Learning

Machine learning at the edge represents a significant paradigm shift in computational intelligence, enabling sophisticated inferencing capabilities on resource-constrained embedded systems such as the EFR32MG24 Wireless Gecko microcontroller. This chapter explores the theoretical foundations and practical implementations of TinyML specifically tailored for microcontroller deployment, with particular focus on the sine wave prediction model as the canonical “Hello World” example of TinyML.

The concept of a “Hello World” example has long been a tradition in programming, where new technologies are introduced with simple code that demonstrates basic functionality. In the domain of TinyML, our sine wave prediction serves as an elegant introduction to the end-to-end process of building, training, and deploying models to microcontrollers.

### 3.2 Theoretical Foundations of TinyML for Microcontrollers

#### 3.2.1 The Computational Constraints Paradigm

Traditional machine learning systems operate under the assumption of abundant computational resources, where model complexity and size are secondary concerns to performance metrics. TinyML, however, inverts this paradigm, placing primary emphasis on resource efficiency while maintaining acceptable inferencing quality.

For the EFR32MG24 platform, with its ARM Cortex-M33 core, limited memory footprint (1536KB flash and 256KB RAM), and power-sensitive applications, we must consider:

1. **Memory-Constrained Learning:** Operating within a 256KB RAM budget necessitates models with minimal memory footprints
2. **Computation-Constrained Inference:** The 78MHz Cortex-M33 processor requires algorithmic optimizations to achieve real-time performance
3. **Energy-Constrained Execution:** Battery-powered applications demand power-aware ML implementations

These constraints fundamentally reshape our approach to machine learning model design, training methodologies, and deployment strategies.

### 3.2.2 Model Compression and Quantization

Central to TinyML is the concept of model compression, which can be formalized as an optimization problem:

$$\min_{\theta'} \mathcal{L}(f_{\theta'}(X), Y) \quad \text{s.t.} \quad |\theta'| \ll |\theta|$$

Where  $\theta$  represents the parameters of the original model,  $\theta'$  the compressed model parameters,  $\mathcal{L}$  the loss function, and  $f_{\theta'}(X)$  the model predictions on input  $X$  compared against ground truth  $Y$ .

Quantization—a key technique in this domain—transforms floating-point weights and activations to reduced-precision integers:

$$Q(w) = \text{round}\left(\frac{w}{\Delta}\right) \cdot \Delta$$

Where  $\Delta$  represents the quantization step size. This transformation reduces both memory requirements and computational complexity at the cost of some precision.

## 3.3 Building Our Sine Wave Model in Google Colab

### 3.3.1 Generating and Processing the Dataset

For our introductory TinyML example, we'll create a sine wave predictor that learns to approximate the sine function. This represents an ideal starting point for several reasons:

1. The sine function is mathematically well-defined and bounded
2. The input-output relationship exhibits nonlinearity that requires proper model architecture
3. The implementation can produce visually verifiable results on a microcontroller

Let's begin by creating a Google Colab notebook to build and train our model. Open a new notebook and start with the following code to generate our training data:

```
import numpy as np
import math
import matplotlib.pyplot as plt
import tensorflow as tf
from tensorflow.keras import layers
import os

# Generate a uniformly distributed set of random numbers in the range from
# 0 to 2, which covers a complete sine wave oscillation
SAMPLES = 1000
np.random.seed(1337)
x_values = np.random.uniform(low=0, high=2*math.pi, size=SAMPLES)
# Shuffle the values to guarantee they're not in order
np.random.shuffle(x_values)
# Calculate the corresponding sine values
y_values = np.sin(x_values)

# Add a small random number to each y value to simulate noise
y_values += 0.1 * np.random.randn(*y_values.shape)
```

```
# Split into train/validation/test sets
TRAIN_SPLIT = int(0.6 * SAMPLES)
TEST_SPLIT = int(0.2 * SAMPLES + TRAIN_SPLIT)
x_train, x_validate, x_test = np.split(x_values, [TRAIN_SPLIT, TEST_SPLIT])
y_train, y_validate, y_test = np.split(y_values, [TRAIN_SPLIT, TEST_SPLIT])

# Plot our data points
plt.figure(figsize=(10, 6))
plt.scatter(x_train, y_train, label='Training data')
plt.scatter(x_validate, y_validate, label='Validation data')
plt.scatter(x_test, y_test, label='Test data')
plt.legend()
plt.title('Sine Wave with Random Noise')
plt.xlabel('x values')
plt.ylabel('y values (sine of x + noise)')
plt.show()
```

### 3.3.2 Constructing and Training the Neural Network Model

Now we’ll construct a simple neural network to learn the sine function:

```
# Create a model with 2 layers of 16 neurons each
model = tf.keras.Sequential()
# First layer takes a scalar input and feeds it through 16 "neurons"
model.add(layers.Dense(16, activation='relu', input_shape=(1,)))
# Second layer with 16 neurons to capture non-linear relationships
model.add(layers.Dense(16, activation='relu'))
# Final layer is a single neuron for our output value
model.add(layers.Dense(1))
# Compile the model using a standard optimizer and loss function for regression
model.compile(optimizer='rmsprop', loss='mse', metrics=['mae'])

# Display model summary to understand its structure
model.summary()

# Train the model on our data
history = model.fit(x_train, y_train,
                    epochs=500,
                    batch_size=16,
                    validation_data=(x_validate, y_validate),
                    verbose=1)

# Plot the training and validation loss
plt.figure(figsize=(10, 6))
plt.plot(history.history['loss'], label='Loss')
plt.plot(history.history['val_loss'], label='Validation Loss')
plt.title('Training and Validation Loss')
```

```

plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.show()

# Plot the training and validation mean absolute error
plt.figure(figsize=(10, 6))
plt.plot(history.history['mae'], label='MAE')
plt.plot(history.history['val_mae'], label='Validation MAE')
plt.title('Training and Validation Mean Absolute Error')
plt.xlabel('Epoch')
plt.ylabel('MAE')
plt.legend()
plt.show()

# Evaluate the model on our test data
test_loss, test_mae = model.evaluate(x_test, y_test)
print(f'Test Loss: {test_loss:.4f}')
print(f'Test MAE: {test_mae:.4f}')

# Generate predictions across the full range for visualization
x_dense = np.linspace(0, 2*math.pi, 200)
y_dense_true = np.sin(x_dense)
y_dense_pred = model.predict(x_dense)

# Plot the true sine curve against our model's predictions
plt.figure(figsize=(10, 6))
plt.plot(x_dense, y_dense_true, 'b-', label='True Sine')
plt.plot(x_dense, y_dense_pred, 'r-', label='Model Prediction')
plt.scatter(x_test, y_test, alpha=0.3, label='Test Data')
plt.title('Sine Wave Prediction')
plt.xlabel('x')
plt.ylabel('sin(x)')
plt.legend()
plt.show()

```

This architecture, though simple, is carefully designed to capture the nonlinear relationship of the sine function. The ReLU (Rectified Linear Unit) activation function is particularly important as it introduces nonlinearity:

$$\text{ReLU}(x) = \max(0, x)$$

We train the model using the mean squared error loss function, which for a regression problem is defined as:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Where  $y_i$  represents the actual sine value and  $\hat{y}_i$  represents our model's prediction.



## 3.4 Optimizing for Microcontroller Deployment

### 3.4.1 Model Conversion and Quantization for TensorFlow Lite

To deploy our trained model to the EFR32MG24 microcontroller, we must convert it into a format suitable for resource-constrained devices:

```
# Convert the model to the TensorFlow Lite format without quantization
converter = tf.lite.TFLiteConverter.from_keras_model(model)
tflite_model = converter.convert()

# Save the model to disk
with open("sine_model.tflite", "wb") as f:
    f.write(tflite_model)

# Convert with quantization for further optimization
converter = tf.lite.TFLiteConverter.from_keras_model(model)
converter.optimizations = [tf.lite.Optimize.DEFAULT]

# Define a generator function that provides our test data's x values
# as a representative dataset
def representative_dataset_generator():
    for value in x_test:
        yield [np.array(value, dtype=np.float32, ndmin=2)]

converter.representative_dataset = representative_dataset_generator
tflite_model_quantized = converter.convert()

# Save the quantized model to disk
with open("sine_model_quantized.tflite", "wb") as f:
    f.write(tflite_model_quantized)

# Print the size reduction achieved through quantization
print(f"Original model size: {len(tflite_model)} bytes")
print(f"Quantized model size: {len(tflite_model_quantized)} bytes")
print(f"Size reduction: {(1 - len(tflite_model_quantized) / len(tflite_model)) * 100:.2f}%")
```

### 3.4.2 Converting to C Code for Embedded Systems

For deployment on microcontrollers like the EFR32MG24, we need to convert our quantized model into a C header file that can be directly included in our firmware:

```
# Function to convert the model to a C array
def convert_tflite_to_c_array(tflite_model, array_name):
    hex_data = ["0x{:02x}:".format(byte) for byte in tflite_model]
    c_array = f"const unsigned char {array_name}[] = {{\n"

    # Format the hex data into rows
    chunk_size = 12
```

```

    for i in range(0, len(hex_data), chunk_size):
        c_array += " " + ", ".join(hex_data[i:i+chunk_size]) + ",\n"

    c_array = c_array[:-2] + "\n};\n"
    c_array += f"const unsigned int {array_name}_len = {len(tflite_model)};\n"

    return c_array

# Generate the C array for our model
c_array = convert_tflite_to_c_array(tflite_model_quantized, "g_sine_model_data")

# Save to a header file
with open("sine_model_data.h", "w") as f:
    f.write("#ifndef SINE_MODEL_DATA_H\n")
    f.write("#define SINE_MODEL_DATA_H\n\n")
    f.write("#include <stdint.h>\n\n")
    f.write(c_array)
    f.write("\n#endif // SINE_MODEL_DATA_H\n")

print("C header file generated: sine_model_data.h")

# Download the files
from google.colab import files
files.download("sine_model.tflite")
files.download("sine_model_quantized.tflite")
files.download("sine_model_data.h")

```

## 3.5 Deploying with Simplicity Studio and Gecko SDK

Now that we have our trained model in a format suitable for microcontrollers, we'll implement the TinyML application using Simplicity Studio and the Gecko SDK. This approach simplifies development by providing a structured framework for EFR32 devices.

### 3.5.1 Creating a New Project in Simplicity Studio

1. Launch Simplicity Studio and connect your EFR32MG24 development board
2. Select your device in the "Debug Adapters" view
3. Click on "Create New Project" in the Launcher perspective
4. Select "Silicon Labs Project Wizard" and click "Next"
5. Choose "Gecko SDK" as the project type
6. Filter for "example" and select "TensorFlow Lite Micro Example" template
7. Configure project settings:
  - Name: `sine_wave_predictor`
  - SDK version: Latest version
  - Click "Next" and then "Finish"

### 3.5.2 Project Structure and Important Files

Simplicity Studio creates a project with the following important files:

- **app.c**: Main application entry point
- **sl\_tflite\_micro\_model.{h,c}**: TensorFlow Lite Micro integration
- **sl\_pwm.{h,c}**: PWM control for LED output
- **sine\_model\_data.h**: Our model data (to be replaced with our trained model)

### 3.5.3 Adding Our Trained Model

1. In Simplicity Studio, locate the project’s inc folder
2. Right-click and select “Import” → “General” → “File System”
3. Browse to the location where you saved `sine_model_data.h`
4. Select the file and click “Finish”

### 3.5.4 Implementing the Application Logic

Now we’ll modify the application code to use our sine wave model. Open `app.c` and replace its contents with the following:

```

/*****
 * @file app.c
 * @brief TinyML Sine Wave Predictor application
 *****/
 * # License
 * <b>Copyright 2023 Silicon Laboratories Inc. www.silabs.com</b>
 *****/
 *
 * SPDX-License-Identifier: Zlib
 *
 * The licensor of this software is Silicon Laboratories Inc.
 *
 * This software is provided 'as-is', without any express or implied
 * warranty. In no event will the authors be held liable for any damages
 * arising from the use of this software.
 *
 * Permission is granted to anyone to use this software for any purpose,
 * including commercial applications, and to alter it and redistribute it
 * freely, subject to the following restrictions:
 *
 * 1. The origin of this software must not be misrepresented; you must not
 *    claim that you wrote the original software. If you use this software
 *    in a product, an acknowledgment in the product documentation would be
 *    appreciated but is not required.
 * 2. Altered source versions must be plainly marked as such, and must not be
 *    misrepresented as being the original software.
 * 3. This notice may not be removed or altered from any source distribution.

```

```

*
*****/
#include "sl_component_catalog.h"
#include "sl_system_init.h"
#include "app.h"
#if defined(SL_CATALOG_POWER_MANAGER_PRESENT)
#include "sl_power_manager.h"
#endif
#include "sl_system_process_action.h"

#include <stdio.h>
#include <math.h>

#include "sl_tflite_micro_model.h"
#include "sl_led.h"
#include "sl_pwm.h"
#include "sl_sleeptimer.h"

// Constants for sine wave demonstration
#define INFERENCES_PER_CYCLE 32
#define X_RANGE (2.0f * 3.14159265359f) // 2 radians
#define PWM_FREQUENCY_HZ 10000
#define INFERENCE_INTERVAL_MS 50

// Global variables
static int inference_count = 0;

void app_init(void)
{
    // Initialize TFLite model
    sl_status_t status = sl_tflite_micro_init();
    if (status != SL_STATUS_OK) {
        printf("Failed to initialize TensorFlow Lite Micro\n");
        return;
    }

    // Initialize PWM for LED control
    sl_pwm_config_t pwm_config = {
        .frequency = PWM_FREQUENCY_HZ,
        .polarity = SL_PWM_ACTIVE_HIGH
    };

    sl_pwm_init(SL_PWM_LEDO, &pwm_config);

    printf("Sine Wave Predictor initialized\n");
}

```

```
void app_process_action(void)
{
    // Calculate x value based on our position in the cycle
    float position = (float)inference_count / (float)INFERENCES_PER_CYCLE;
    float x_val = position * X_RANGE;

    // Prepare the input tensor with our x value
    float input_data[1] = { x_val };
    sl_tflite_micro_tensor_t input_tensor;
    sl_status_t status = sl_tflite_micro_get_input_tensor(0, &input_tensor);
    if (status != SL_STATUS_OK) {
        printf("Failed to get input tensor\n");
        return;
    }

    // Copy our input data to the input tensor
    status = sl_tflite_micro_set_tensor_data(&input_tensor, input_data, sizeof(input_data));
    if (status != SL_STATUS_OK) {
        printf("Failed to set input tensor data\n");
        return;
    }

    // Run inference
    status = sl_tflite_micro_invoke();
    if (status != SL_STATUS_OK) {
        printf("Inference failed\n");
        return;
    }

    // Get the output tensor
    sl_tflite_micro_tensor_t output_tensor;
    status = sl_tflite_micro_get_output_tensor(0, &output_tensor);
    if (status != SL_STATUS_OK) {
        printf("Failed to get output tensor\n");
        return;
    }

    // Get the predicted sine value
    float predicted_sine = 0.0f;
    status = sl_tflite_micro_get_tensor_data(&output_tensor, &predicted_sine, sizeof(predicted_sine));
    if (status != SL_STATUS_OK) {
        printf("Failed to get output tensor data\n");
        return;
    }

    // Map the sine value (-1 to 1) to PWM duty cycle (0 to 100%)
    uint8_t duty_cycle = (uint8_t)((predicted_sine + 1.0f) * 50.0f);
}
```

```

// Set LED brightness using PWM
sl_pwm_set_duty_cycle(SL_PWM_LED0, duty_cycle);

// Log the values (only every 8th inference to reduce console traffic)
if (inference_count % 8 == 0) {
    printf("x: %.3f, predicted sine: %.3f, duty cycle: %d%%\n",
        x_val, predicted_sine, duty_cycle);
}

// Increment the inference counter
inference_count++;
if (inference_count >= INFERENCES_PER_CYCLE) {
    inference_count = 0;
}

// Add a delay before the next inference
sl_sleeptimer_delay_millisecond(INFERENCE_INTERVAL_MS);
}

```

### 3.5.5 Creating the Model Integration File

Create a new file called `sl_tflite_micro_model.c` in the `src` folder with the following content:

```

#include "sl_tflite_micro_model.h"
#include "sine_model_data.h"
#include <stdio.h>

// TensorFlow Lite for Microcontrollers components
#include "tensorflow/lite/micro/all_ops_resolver.h"
#include "tensorflow/lite/micro/micro_error_reporter.h"
#include "tensorflow/lite/micro/micro_interpreter.h"
#include "tensorflow/lite/schema/schema_generated.h"
#include "tensorflow/lite/version.h"

// Global variables
static tflite::MicroErrorReporter micro_error_reporter;
static tflite::ErrorReporter* error_reporter = &micro_error_reporter;
static const tflite::Model* model = nullptr;
static tflite::MicroInterpreter* interpreter = nullptr;
static TfLiteTensor* input_tensor = nullptr;
static TfLiteTensor* output_tensor = nullptr;

// Create an area of memory for input, output, and intermediate arrays
constexpr int kTensorArenaSize = 8 * 1024;
static uint8_t tensor_arena[kTensorArenaSize];

sl_status_t sl_tflite_micro_init(void)

```

```

{
    // Map the model into a usable data structure
    model = tflite::GetModel(g_sine_model_data);
    if (model->version() != TFLITE_SCHEMA_VERSION) {
        printf("Model version mismatch: %d vs %d\n", model->version(), TFLITE_SCHEMA_VERSION);
        return SL_STATUS_FAIL;
    }

    // Create an all operations resolver
    static tflite::AllOpsResolver resolver;

    // Build an interpreter to run the model
    static tflite::MicroInterpreter static_interpreter(
        model, resolver, tensor_arena, kTensorArenaSize, error_reporter);
    interpreter = &static_interpreter;

    // Allocate memory for all tensors
    TfLiteStatus allocate_status = interpreter->AllocateTensors();
    if (allocate_status != kTfLiteOk) {
        printf("AllocateTensors() failed\n");
        return SL_STATUS_ALLOCATION_FAILED;
    }

    // Get pointers to the model's input and output tensors
    input_tensor = interpreter->input(0);
    output_tensor = interpreter->output(0);

    // Check that input and output tensors are the expected size and type
    if (input_tensor->dims->size != 2 || input_tensor->dims->data[0] != 1 ||
        input_tensor->dims->data[1] != 1 || input_tensor->type != kTfLiteFloat32) {
        printf("Unexpected input tensor format\n");
        return SL_STATUS_INVALID_PARAMETER;
    }

    if (output_tensor->dims->size != 2 || output_tensor->dims->data[0] != 1 ||
        output_tensor->dims->data[1] != 1 || output_tensor->type != kTfLiteFloat32) {
        printf("Unexpected output tensor format\n");
        return SL_STATUS_INVALID_PARAMETER;
    }

    printf("TensorFlow Lite Micro initialized successfully\n");
    return SL_STATUS_OK;
}

sl_status_t sl_tflite_micro_get_input_tensor(uint8_t index, sl_tflite_micro_tensor_t* tensor)
{
    if (interpreter == nullptr || index >= interpreter->inputs_size()) {

```

```

    return SL_STATUS_INVALID_PARAMETER;
}

tensor->tensor = interpreter->input(index);
return SL_STATUS_OK;
}

sl_status_t sl_tflite_micro_get_output_tensor(uint8_t index, sl_tflite_micro_tensor_t* tensor)
{
    if (interpreter == nullptr || index >= interpreter->outputs_size()) {
        return SL_STATUS_INVALID_PARAMETER;
    }

    tensor->tensor = interpreter->output(index);
    return SL_STATUS_OK;
}

sl_status_t sl_tflite_micro_set_tensor_data(sl_tflite_micro_tensor_t* tensor,
                                           const void* data,
                                           size_t size)
{
    if (tensor == nullptr || tensor->tensor == nullptr || data == nullptr) {
        return SL_STATUS_NULL_POINTER;
    }

    // Size check based on tensor type and dims
    size_t tensor_size = 1;
    for (int i = 0; i < tensor->tensor->dims->size; i++) {
        tensor_size *= tensor->tensor->dims->data[i];
    }

    if (tensor->tensor->type == kTfLiteFloat32) {
        tensor_size *= sizeof(float);
    } else if (tensor->tensor->type == kTfLiteInt8) {
        tensor_size *= sizeof(int8_t);
    } else if (tensor->tensor->type == kTfLiteUInt8) {
        tensor_size *= sizeof(uint8_t);
    } else {
        return SL_STATUS_NOT_SUPPORTED;
    }

    if (size > tensor_size) {
        return SL_STATUS_WOULD_OVERFLOW;
    }

    // Copy the data to the tensor
    memcpy(tensor->tensor->data.raw, data, size);
}

```



```

    return SL_STATUS_OK;
}

sl_status_t sl_tflite_micro_get_tensor_data(sl_tflite_micro_tensor_t* tensor,
                                             void* data,
                                             size_t size)
{
    if (tensor == nullptr || tensor->tensor == nullptr || data == nullptr) {
        return SL_STATUS_NULL_POINTER;
    }

    // Size check based on tensor type and dims
    size_t tensor_size = 1;
    for (int i = 0; i < tensor->tensor->dims->size; i++) {
        tensor_size *= tensor->tensor->dims->data[i];
    }

    if (tensor->tensor->type == kTfLiteFloat32) {
        tensor_size *= sizeof(float);
    } else if (tensor->tensor->type == kTfLiteInt8) {
        tensor_size *= sizeof(int8_t);
    } else if (tensor->tensor->type == kTfLiteUInt8) {
        tensor_size *= sizeof(uint8_t);
    } else {
        return SL_STATUS_NOT_SUPPORTED;
    }

    if (size > tensor_size) {
        return SL_STATUS_WOULD_OVERFLOW;
    }

    // Copy the data from the tensor
    memcpy(data, tensor->tensor->data.raw, size);
    return SL_STATUS_OK;
}

sl_status_t sl_tflite_micro_invoke(void)
{
    if (interpreter == nullptr) {
        return SL_STATUS_NOT_INITIALIZED;
    }

    TfLiteStatus status = interpreter->Invoke();
    if (status != kTfLiteOk) {
        return SL_STATUS_FAIL;
    }
}

```

```
    return SL_STATUS_OK;
}
```

Now, create the header file `sl_tflite_micro_model.h` in the `inc` folder:

```
#ifndef SL_TFLITE_MICRO_MODEL_H
#define SL_TFLITE_MICRO_MODEL_H

#include "sl_status.h"
#include <stdint.h>
#include <stddef.h>

#ifdef __cplusplus
extern "C" {
#endif

// Forward declarations from TensorFlow Lite
#ifdef __cplusplus
namespace tflite {
struct TfLiteTensor;
} // namespace tflite
typedef struct tflite::TfLiteTensor TfLiteTensor;
#else
typedef struct TfLiteTensor TfLiteTensor;
#endif

// Tensor structure
typedef struct {
    TfLiteTensor* tensor;
} sl_tflite_micro_tensor_t;

/**
 * @brief Initialize TensorFlow Lite Micro with the sine model
 *
 * @return sl_status_t SL_STATUS_OK if successful
 */
sl_status_t sl_tflite_micro_init(void);

/**
 * @brief Get an input tensor by index
 *
 * @param index Index of the input tensor
 * @param tensor Pointer to the tensor structure to fill
 * @return sl_status_t SL_STATUS_OK if successful
 */
sl_status_t sl_tflite_micro_get_input_tensor(uint8_t index, sl_tflite_micro_tensor_t* tensor);

/**
```

```

* @brief Get an output tensor by index
*
* @param index Index of the output tensor
* @param tensor Pointer to the tensor structure to fill
* @return sl_status_t SL_STATUS_OK if successful
*/
sl_status_t sl_tflite_micro_get_output_tensor(uint8_t index, sl_tflite_micro_tensor_t* tensor);

/**
* @brief Set data to a tensor
*
* @param tensor Pointer to the tensor
* @param data Pointer to the data to copy
* @param size Size of the data in bytes
* @return sl_status_t SL_STATUS_OK if successful
*/
sl_status_t sl_tflite_micro_set_tensor_data(sl_tflite_micro_tensor_t* tensor,
                                           const void* data,
                                           size_t size);

/**
* @brief Get data from a tensor
*
* @param tensor Pointer to the tensor
* @param data Pointer to the buffer to receive the data
* @param size Size of the buffer in bytes
* @return sl_status_t SL_STATUS_OK if successful
*/
sl_status_t sl_tflite_micro_get_tensor_data(sl_tflite_micro_tensor_t* tensor,
                                           void* data,
                                           size_t size);

/**
* @brief Run inference using the TensorFlow Lite model
*
* @return sl_status_t SL_STATUS_OK if successful
*/
sl_status_t sl_tflite_micro_invoke(void);

#ifdef __cplusplus
}
#endif

#endif // SL_TFLITE_MICRO_MODEL_H

```

### 3.5.6 Building and Flashing the Application

1. In Simplicity Studio, right-click on the project and select “Build Project”

2. After successful compilation, right-click again and select “Run As” → “Silicon Labs ARM Program”
3. The application will be flashed to your EFR32MG24 device and start running

### 3.5.7 Observing the Results

Once your application is running on the EFR32MG24 device:

1. The LED will pulse with brightness that follows the sine wave pattern
2. Open the Serial Console in Simplicity Studio to view debug output
3. You’ll see logs showing the input value, predicted sine value, and the corresponding LED duty cycle

## 3.6 How it Works: Understanding the Implementation

Our TinyML sine wave predictor consists of several key components:

1. **Model Training and Conversion:** Using Google Colab, we trained a small neural network to approximate the sine function and converted it to TF Lite format, then to a C array.
2. **TensorFlow Lite Micro Integration:** We’ve implemented a simple wrapper around TF Lite Micro’s C++ API, providing a clean C interface for our application.
3. **Application Logic:** The main application loop:
  - Calculates an x value based on where we are in the cycle
  - Feeds this value into the model
  - Retrieves the predicted sine value
  - Maps the prediction to LED brightness via PWM
4. **Visual Output:** The LED brightness follows a sine wave pattern, providing visual confirmation that our model is working correctly.

## 3.7 Extending the TinyML Application

Now that we have our basic “Hello World” TinyML application running, there are several ways we can extend and enhance it:

### 3.7.1 1. Adding Multiple LED Support

For devices with multiple LEDs, we can create more interesting visual patterns by controlling multiple LEDs based on different phases of the sine wave:

```
// In app.c, add phase offsets for each LED
#define LED_COUNT 4 // Assuming 4 available LEDs
const float phase_offsets[LED_COUNT] = {
    0.0f, // LED0: No phase offset
    0.5f * X_RANGE / 4.0f, // LED1: 45 degrees offset
    X_RANGE / 4.0f, // LED2: 90 degrees offset
    1.5f * X_RANGE / 4.0f // LED3: 135 degrees offset
};
```

```
// Then in app_process_action(), add a loop to control all LEDs
for (int i = 0; i < LED_COUNT; i++) {
    // Calculate offset x value for this LED
    float led_x_val = x_val + phase_offsets[i];
    if (led_x_val >= X_RANGE) {
        led_x_val -= X_RANGE; // Wrap around to stay in range
    }

    // Prepare input tensor with our x value
    float input_data[1] = { led_x_val };
    sl_tflite_micro_tensor_t input_tensor;
    status = sl_tflite_micro_get_input_tensor(0, &input_tensor);
    if (status != SL_STATUS_OK) continue;

    // Set tensor data, invoke model, and get output as before...

    // Set corresponding LED brightness
    uint8_t duty_cycle = (uint8_t)((predicted_sine + 1.0f) * 50.0f);
    sl_pwm_set_duty_cycle(i, duty_cycle); // Assuming LED PWM instances are indexed
}
```

### 3.7.2 2. Adding LCD Display Support

If your EFR32MG24 development board includes an LCD display, you can visualize the sine wave more directly:

```
// In app.c, add LCD-related includes
#include "sl_glib.h"
#include "sl_simple_lcd.h"

// Add LCD dimensions and buffers
#define LCD_WIDTH 128
#define LCD_HEIGHT 64
#define HISTORY_LENGTH LCD_WIDTH
static float sine_history[HISTORY_LENGTH];
static int history_index = 0;

// Initialize the LCD in app_init()
sl_simple_lcd_init();
sl_glib_init();
// Clear history buffer
for (int i = 0; i < HISTORY_LENGTH; i++) {
    sine_history[i] = 0.0f;
}

// In app_process_action(), after getting the predicted sine:
// Store in history buffer
```

```

sine_history[history_index] = predicted_sine;
history_index = (history_index + 1) % HISTORY_LENGTH;

// Every few inferences, update the display
if (inference_count % 4 == 0) {
    GLIB_Context_t context;
    sl_glib_get_context(&context);

    // Clear display
    GLIB_clear(&context);

    // Draw x and y axes
    GLIB_drawLineH(&context, 0, LCD_WIDTH-1, LCD_HEIGHT/2);
    GLIB_drawLineV(&context, 0, 0, LCD_HEIGHT-1);

    // Draw the sine wave
    for (int i = 0; i < HISTORY_LENGTH-1; i++) {
        int x1 = i;
        int y1 = (int)(LCD_HEIGHT/2 - (sine_history[(history_index + i) % HISTORY_LENGTH] * LCD_HEIGHT/4));
        int x2 = i + 1;
        int y2 = (int)(LCD_HEIGHT/2 - (sine_history[(history_index + i + 1) % HISTORY_LENGTH] * LCD_HEIGHT/4));
        GLIB_drawLine(&context, x1, y1, x2, y2);
    }

    // Update display
    sl_glib_update_display();
}

```

### 3.7.3 3. Implementing Power Optimization

To make our TinyML application more power-efficient for battery-powered operation, we can add sleep modes between inferences:

```

// Replace the fixed delay with sleep mode
// Instead of: sl_sleeptimer_delay_millisecond(INFERENCE_INTERVAL_MS);

#ifdef SL_CATALOG_POWER_MANAGER_PRESENT
    // Schedule next wakeup
    sl_sleeptimer_tick_t ticks = sl_sleeptimer_ms_to_tick(INFERENCE_INTERVAL_MS);
    sl_power_manager_schedule_wakeup(ticks, NULL, NULL);

    // Enter sleep mode
    sl_power_manager_sleep();
#else
    // Fall back to delay if power manager isn't available
    sl_sleeptimer_delay_millisecond(INFERENCE_INTERVAL_MS);
#endif

```

### 3.7.4 4. Enhanced User Interface with Buttons

We can use the buttons on the development board to control aspects of the application:

```
// Include button support
#include "sl_button.h"
#include "sl_simple_button.h"
#include "sl_simple_button_btn0_config.h"

// Add state variables
static bool paused = false;
static float speed_factor = 1.0f;

// In app_init, initialize buttons
sl_button_init(&sl_button_btn0);

// Check button state in app_process_action
if (sl_button_get_state(&sl_button_btn0) == SL_SIMPLE_BUTTON_PRESSED) {
    // Toggle pause state
    paused = !paused;
    printf("Application %s\n", paused ? "paused" : "resumed");
}

// Only update inference_count if not paused
if (!paused) {
    inference_count += 1;
    if (inference_count >= INFERENCES_PER_CYCLE) inference_count = 0;
}
```

### 3.7.5 5. Performance Profiling and Optimization

To understand and optimize the performance of our TinyML application, we can add timing measurements:

```
// Add profiling includes
#include "em_cmu.h"
#include "em_timer.h"

// Setup timer for profiling in app_init()
CMU_ClockEnable(cmuClock_TIMER1, true);
TIMER_Init_TypeDef timerInit = TIMER_INIT_DEFAULT;
TIMER_Init(TIMER1, &timerInit);

// In app_process_action(), measure inference time
// Reset timer
TIMER_CounterSet(TIMER1, 0);

// Start timer
TIMER_Enable(TIMER1, true);
```

```
// Run inference
status = sl_tflite_micro_invoke();

// Stop timer and read value
TIMER_Enable(TIMER1, false);
uint32_t ticks = TIMER_CounterGet(TIMER1);

// Convert ticks to microseconds (depends on timer clock)
uint32_t us = ticks / (CMU_ClockFreqGet(cmuClock_TIMER1) / 1000000);

// Log every Nth inference
if (inference_count % 10 == 0) {
    printf("Inference time: %lu microseconds\n", us);
}
```

### 3.8 Building TinyML Applications with the Gecko SDK

Compared to the traditional TinyML approach with direct TensorFlow Lite for Microcontrollers integration, the Gecko SDK approach offers several advantages for EFR32 developers:

#### 3.8.1 Simplified Project Setup

The Gecko SDK provides a structured approach to project creation with built-in TinyML support:

1. **Project Templates:** Simplicity Studio's project wizard includes TinyML templates that set up the necessary directory structure, build configuration, and initial code.
2. **Integrated Build System:** The SDK handles compiler flags, library dependencies, and linking, eliminating the need for custom Makefiles.
3. **Hardware Abstraction Layer (HAL):** The SDK provides hardware-specific drivers and APIs for peripherals like PWM, GPIOs, and timers, making it easier to integrate TinyML with device hardware.

#### 3.8.2 Streamlined Development Workflow

The development workflow with Simplicity Studio and Gecko SDK is straightforward:

1. **Model Training:** Use Google Colab or TensorFlow on your computer to train and convert models.
2. **Project Creation:** Launch Simplicity Studio, select the TensorFlow Lite Micro example template, and create a new project.
3. **Model Integration:** Import your model header file into the project.
4. **Application Logic:** Write code in C to initialize the model, prepare inputs, run inference, and process outputs.
5. **Build and Flash:** Use Simplicity Studio's integrated tools to compile the code and flash it to your device.
6. **Debug and Monitor:** The Serial Console and Energy Profiler tools help monitor application behavior and optimize performance.



### 3.8.3 Hardware-Specific Optimizations

The Gecko SDK includes optimizations specifically for the EFR32 platform:

1. **Memory Optimization:** Memory management is tuned for the EFR32’s memory architecture.
2. **Power Management:** Integration with the Energy Management Unit (EMU) allows for fine-grained control over active, sleep, and deep sleep states.
3. **Peripheral Control:** Direct access to hardware accelerators and peripherals that can enhance TinyML performance.

## 3.9 Conclusion

The sine wave predictor represents an elegant “Hello World” example of TinyML deployment on the EFR32MG24 platform. While seemingly simple, this implementation encompasses all the key elements of machine learning on microcontrollers:

1. Model design with consideration for resource constraints
2. Training and evaluation on standard datasets
3. Quantization and optimization for embedded deployment
4. Integration with Simplicity Studio and Gecko SDK
5. Hardware output integration via GPIO and PWM capabilities

Through this foundation, developers can extend to more complex TinyML applications on the EFR32MG24, such as sensor fusion, predictive maintenance, anomaly detection, and keyword spotting—all within a power envelope suitable for long-term battery-powered operation.

The techniques demonstrated here—model quantization, C code generation, and deployment with Simplicity Studio—provide a template that can be adapted for more sophisticated machine learning tasks, enabling a new class of intelligent edge devices based on the EFR32MG24 platform.

## Chapter 4

# Building a TinyML Application

In the previous chapter, we trained a neural network model to predict sine wave values and prepared it for deployment on an EFR32MG24 microcontroller. Now we'll build a complete application around this model and deploy it to the hardware. This chapter focuses on the practical aspects of implementing TinyML using Silicon Labs' Gecko SDK and Simplicity Studio rather than the traditional TensorFlow Lite for Microcontrollers approach.

### 4.1 Understanding the Gecko SDK Approach to TinyML

The Gecko SDK provides a structured approach to embedded development specifically optimized for Silicon Labs' microcontrollers. This offers several advantages over the more generic TinyML implementations:

1. **Pre-integrated Components:** The SDK includes optimized TensorFlow Lite Micro components already configured for EFR32 devices
2. **Hardware Abstraction Layer:** Direct integration with EFR32 peripherals through a consistent API
3. **Project Templates:** Simplicity Studio provides starting points for TinyML applications
4. **Advanced Tooling:** Debugging, energy profiling, and configuration tools are built into the development environment

### 4.2 Setting Up Your Development Environment

Before we begin building our application, ensure you have the following tools installed:

1. **Simplicity Studio 5:** Download and install from [Silicon Labs' website](#)
2. **Gecko SDK:** The latest version will be installed through Simplicity Studio
3. **J-Link Drivers:** These should be installed with Simplicity Studio
4. **EFR32MG24 Development Board:** Connect this to your computer via USB

### 4.3 Creating a New Project in Simplicity Studio

Let's start by creating a TinyML project in Simplicity Studio:

1. Launch Simplicity Studio 5

2. In the Launcher perspective, click on your connected EFR32MG24 device
3. Click “Create New Project” in the “Overview” tab
4. Select “Silicon Labs Project Wizard” and click “NEXT”
5. In the SDK Selection dialog, ensure the latest Gecko SDK is selected and click “NEXT”
6. In the Project Generation dialog:
  - Filter for “example” in the search box
  - Select “TensorFlow Lite Micro Example”
  - Click “NEXT”
7. Configure your project:
  - Name: `sine_wave_predictor`
  - Keep the default location
  - Click “FINISH”

Simplicity Studio will generate a project with the necessary components for a TinyML application. Let’s explore the project structure before making our modifications.

## 4.4 Exploring the Project Structure

The generated project includes several important directories and files:

- **config/**: Contains hardware configuration files for your specific board
- **gecko\_sdk/**: The Gecko SDK source code, including TensorFlow Lite Micro
- **autogen/**: Auto-generated initialization code for the device
- **app.c**: Your application’s main source file
- **app.h**: Header file for your application

The TensorFlow Lite Micro example comes with a sample model that classifies motion patterns. We’ll replace this with our sine wave model.

## 4.5 Importing the Sine Wave Model

First, let’s import the sine model data that we generated in the previous chapter:

1. Right-click on the project in the “Project Explorer” view
2. Select “Import” → “General” → “File System”
3. Browse to the location where you saved `sine_model_data.h`
4. Select the file and click “Finish”

The model data will be added to your project. Now let’s modify the application code to use our sine wave model.

## 4.6 Implementing the Application

Let’s replace the content of `app.c` with our sine wave prediction application code. Open `app.c` and replace its contents with the following:

```

/*****
 * @file app.c
 * @brief TinyML Sine Wave Predictor application
 *****/
 * # License
 * <b>Copyright 2023 Silicon Laboratories Inc. www.silabs.com</b>
 *****/
 *
 * SPDX-License-Identifier: Zlib
 *
 * The licensor of this software is Silicon Laboratories Inc.
 *
 * This software is provided 'as-is', without any express or implied
 * warranty. In no event will the authors be held liable for any damages
 * arising from the use of this software.
 *
 * Permission is granted to anyone to use this software for any purpose,
 * including commercial applications, and to alter it and redistribute it
 * freely, subject to the following restrictions:
 *
 * 1. The origin of this software must not be misrepresented; you must not
 *    claim that you wrote the original software. If you use this software
 *    in a product, an acknowledgment in the product documentation would be
 *    appreciated but is not required.
 * 2. Altered source versions must be plainly marked as such, and must not be
 *    misrepresented as being the original software.
 * 3. This notice may not be removed or altered from any source distribution.
 *
 *****/
#include "sl_component_catalog.h"
#include "sl_system_init.h"
#include "app.h"
#if defined(SL_CATALOG_POWER_MANAGER_PRESENT)
#include "sl_power_manager.h"
#endif
#include "sl_system_process_action.h"

/* Additional includes for our application */
#include <stdio.h>
#include <math.h>

/* Include TensorFlow Lite components */
#include "tensorflow/lite/micro/all_ops_resolver.h"
#include "tensorflow/lite/micro/micro_error_reporter.h"
#include "tensorflow/lite/micro/micro_interpreter.h"
#include "tensorflow/lite/schema/schema_generated.h"
#include "tensorflow/lite/version.h"

```

```

/* Include our model data */
#include "sine_model_data.h"

/* Include hardware control components */
#include "sl_led.h"
#include "sl_simple_led_instances.h"
#include "sl_sleeptimer.h"

/* Constants for sine wave demonstration */
#define INFERENCES_PER_CYCLE 32
#define X_RANGE (2.0f * 3.14159265359f) /* 2 radians */
#define INFERENCE_INTERVAL_MS 50

/* Global variables for TensorFlow Lite model */
static tflite::MicroErrorReporter micro_error_reporter;
static tflite::ErrorReporter* error_reporter = &micro_error_reporter;
/* We'll use the C version of TensorFlow Lite Micro API */
static tflite::MicroInterpreter* interpreter = nullptr;
static TfLiteTensor* input = nullptr;
static TfLiteTensor* output = nullptr;

/* Create an area of memory for input, output, and intermediate arrays */
#define TENSOR_ARENA_SIZE (10 * 1024)
static uint8_t tensor_arena[TENSOR_ARENA_SIZE];

/* Application state variables */
static int inference_count = 0;

/*****
 * Initialize application.
 *****/
void app_init(void)
{
    /* Map the model into a usable data structure */
    model = tflite::GetModel(g_sine_model_data);
    if (model->version() != TFLITE_SCHEMA_VERSION) {
        error_reporter->Report(
            "Model provided is schema version %d not equal "
            "to supported version %d.\n",
            model->version(), TFLITE_SCHEMA_VERSION);
        return;
    }

    /* This pulls in all the operation implementations we need */
    static tflite::AllOpsResolver resolver;

    /* Build an interpreter to run the model with */

```

```

static tfLite::MicroInterpreter static_interpreter(
    model, resolver, tensor_arena, TENSOR_ARENA_SIZE, error_reporter);
interpreter = &static_interpreter;

/* Allocate memory from the tensor_arena for the model's tensors */
TfLiteStatus allocate_status = interpreter->AllocateTensors();
if (allocate_status != kTfLiteOk) {
    error_reporter->Report("AllocateTensors() failed");
    return;
}

/* Obtain pointers to the model's input and output tensors */
input = interpreter->input(0);
output = interpreter->output(0);

/* Check that input tensor dimensions are as expected */
if (input->dims->size != 2 || input->dims->data[0] != 1 ||
    input->dims->data[1] != 1 || input->type != kTfLiteFloat32) {
    error_reporter->Report("Unexpected input tensor dimensions or type");
    return;
}

/* Initialize LED */
sl_led_init(SL_SIMPLE_LED_INSTANCE(0));

/* Print initialization message */
printf("Sine Wave Predictor initialized successfully!\n");
printf("Model input dims: %d x %d, type: %d\n",
    input->dims->data[0], input->dims->data[1], input->type);
}

/*****
 * App ticking function.
 *****/
void app_process_action(void)
{
    /* Calculate an x value to feed into the model based on current inference count */
    float position = (float)(inference_count) / (float)(INFERENCES_PER_CYCLE);
    float x_val = position * X_RANGE;

    /* Set the input tensor with our calculated x value */
    input->data.f[0] = x_val;

    /* Run inference, and report any error */
    TfLiteStatus invoke_status = TF_MicroInterpreter_Invoke(interpreter);
    if (invoke_status != kTfLiteOk) {
        printf("Invoke failed on x_val: %f\n", (double)x_val);
    }
}

```

```

    return;
}

/* Read the predicted y value from the model's output tensor */
float y_val = output->data.f[0];

/* Map the sine output (-1 to 1) to LED brightness
 * For simplicity, we'll just turn the LED on when the value is positive
 * and off when it's negative. For PWM control, you would need to
 * configure a PWM peripheral. */
if (y_val > 0) {
    sl_led_turn_on(SL_SIMPLE_LED_INSTANCE(0));
} else {
    sl_led_turn_off(SL_SIMPLE_LED_INSTANCE(0));
}

/* Log every 4th inference to avoid flooding the console */
if (inference_count % 4 == 0) {
    printf("x_value: %f, predicted_sine: %f\n", (double)x_val, (double)y_val);
}

/* Increment the inference_count, and reset it if we have reached
 * the total number per cycle */
inference_count += 1;
if (inference_count >= INFERENCES_PER_CYCLE) inference_count = 0;

/* Add a delay between inferences */
sl_sleeptimer_delay_millisecond(INFERENCE_INTERVAL_MS);
}

```

This application will: 1. Initialize the TensorFlow Lite Micro interpreter with our sine model 2. Set up an LED for output 3. In each loop iteration: - Calculate an  $x$  value within our 0 to  $2\pi$  range - Run inference to get the predicted sine value - Toggle the LED based on whether the sine value is positive or negative - Log the values to the console - Increment the inference counter

## 4.7 Enhancing Output with PWM Control

The basic application just toggles an LED, but we can create a more interesting visualization by controlling LED brightness with PWM. Let's create a PWM component for our project:

1. Right-click on your project in the Project Explorer
2. Select "Configure Project"
3. Click on "SOFTWARE COMPONENTS"
4. In the search box, type "PWM"
5. Find "PWM Driver" → "Simple PWM" and click "Install"
6. Click "Force Install" if prompted
7. Click "DONE" to save the configuration

Now, modify the application code to use PWM for LED brightness control. Replace the LED control section in `app_process_action()` with:

```
/* Map the sine output (-1 to 1) to PWM duty cycle (0 to 100%) */
uint8_t duty_cycle = (uint8_t)((y_val + 1.0f) * 50.0f);

/* Set LED brightness using PWM */
sl_pwm_set_duty_cycle(SL_PWM_INSTANCE(0), duty_cycle);
```

Also, add the PWM initialization in the `app_init()` function after the LED initialization:

```
/* Initialize PWM for LED brightness control */
sl_pwm_config_t pwm_config = {
    .frequency = 10000, /* 10 kHz PWM frequency */
    .polarity = SL_PWM_ACTIVE_HIGH
};
sl_pwm_init(SL_PWM_INSTANCE(0), &pwm_config);
```

Don't forget to include the PWM header at the top of the file:

```
#include "sl_pwm.h"
#include "sl_simple_pwm_instances.h"
```

## 4.8 Building and Flashing the Application

Now let's build and flash our application to the EFR32MG24 board:

1. Right-click on your project in the Project Explorer
2. Select "Build Project"
3. Once the build completes successfully, right-click again on the project
4. Select "Run As" → "Silicon Labs ARM Program"

Simplicity Studio will compile your code, flash it to the device, and start execution.

## 4.9 Debugging and Monitoring

To monitor the output of your application:

1. In Simplicity Studio, go to the "Debug Adapters" view
2. Right-click on your connected device and select "Launch Console"
3. In the console dialog, select "Serial 1" and click "OK"

You should now see the application's output messages showing x values and predicted sine values.

## 4.10 Optimizing TinyML Performance

### 4.10.1 Memory Optimization

TinyML applications on microcontrollers must be memory-efficient. Let's look at ways to optimize memory usage:



1. **Tensor Arena Size:** Reduce the `TENSOR_ARENA_SIZE` to the minimum required. Try starting with 10KB and reducing it incrementally:

```
#define TENSOR_ARENA_SIZE (10 * 1024) /* Start with 10KB */
```

You can determine the minimum required size by adding debug output:

```
/* Add this after interpreter->AllocateTensors() in app_init() */
size_t used_bytes = interpreter->arena_used_bytes();
printf("Model uses %d bytes of tensor arena\n", (int)used_bytes);
```

2. **Selective Op Resolution:** Instead of using `AllOpsResolver`, create a custom resolver with only the operations needed:

```
/* Replace AllOpsResolver with this */
static tflite::MicroMutableOpResolver<4> resolver;
resolver.AddFullyConnected();
resolver.AddRelu();
resolver.AddAdd();
resolver.AddMul();
```

### 4.10.2 Power Optimization

For battery-powered applications, power efficiency is critical:

1. **Sleep Between Inferences:** Replace the simple delay with a power-efficient sleep:

```
/* Replace sl_sleeptimer_delay_millisecond() with: */
#if defined(SL_CATALOG_POWER_MANAGER_PRESENT)
/* Schedule next wakeup */
sl_sleeptimer_tick_t ticks = sl_sleeptimer_ms_to_tick(INFERENCE_INTERVAL_MS);
sl_power_manager_schedule_wakeup(ticks, NULL, NULL);

/* Enter sleep mode */
sl_power_manager_sleep();
#else
/* Fall back to delay if power manager isn't available */
sl_sleeptimer_delay_millisecond(INFERENCE_INTERVAL_MS);
#endif
```

2. **Measurement with Energy Profiler:** Simplicity Studio includes an Energy Profiler tool to measure power consumption:
  - Connect your board with the Advanced Energy Monitor (AEM)
  - In Simplicity Studio, go to Tools → Energy Profiler
  - Start a capture while your application is running
  - Analyze current consumption during inference and sleep periods

### 4.10.3 Timing Performance

To measure inference time:

```
/* Add these includes */
#include "em_cmu.h"
#include "em_timer.h"

/* Initialize timer in app_init() */
CMU_ClockEnable(cmuClock_TIMER0, true);
TIMER_Init_TypeDef timerInit = TIMER_INIT_DEFAULT;
TIMER_Init(TIMER0, &timerInit);

/* In app_process_action(), surround the inference with timing code */
/* Reset and start timer */
TIMER_CounterSet(TIMER0, 0);
TIMER_Enable(TIMER0, true);

/* Run inference */
TfLiteStatus invoke_status = interpreter->Invoke();

/* Stop timer and read counter */
TIMER_Enable(TIMER0, false);
uint32_t ticks = TIMER_CounterGet(TIMER0);

/* Convert ticks to microseconds */
uint32_t us = ticks / (CMU_ClockFreqGet(cmuClock_TIMER0) / 1000000);

/* Log timing information */
if (inference_count % 10 == 0) {
    printf("Inference took %lu microseconds\n", us);
}
```

## 4.11 Adding User Interaction with Buttons

We can make our application more interactive by using buttons to control its behavior:

1. Add a button component to your project:
  - Right-click on your project and select “Configure Project”
  - Go to “SOFTWARE COMPONENTS”
  - Search for “button” and install “Simple Button” components
2. Modify your code to handle button presses:

```
/* Include button headers */
#include "sl_button.h"
#include "sl_simple_button_instances.h"
```

```

/* In app_init() */
/* Initialize buttons */
sl_button_init(SL_SIMPLE_BUTTON_INSTANCE(0));

/* In app_process_action(), check for button press */
if (sl_button_get_state(SL_SIMPLE_BUTTON_INSTANCE(0)) == SL_SIMPLE_BUTTON_PRESSED) {
    /* Toggle between normal speed and double speed */
    static bool fast_mode = false;
    fast_mode = !fast_mode;

    /* Update the inference interval */
    inference_interval_ms = fast_mode ? 25 : 50;

    printf("Speed set to %s\n", fast_mode ? "FAST" : "NORMAL");
}

```

## 4.12 Enhanced Visualization with LCD (if available)

If your development board has an LCD display, you can create more sophisticated visualizations:

1. Add the LCD components to your project:
  - In the “Configure Project” dialog, search for “lcd”
  - Install the “GLIB Graphics Library” and “Simple LCD”
2. Modify your code to display the sine wave on the LCD:

```

/* Include LCD headers */
#include "sl_glib.h"
#include "sl_simple_lcd.h"

/* In app_init() */
/* Initialize LCD */
sl_simple_lcd_init();
sl_glib_initialize();

/* Define a buffer to store recent sine wave values */
#define HISTORY_SIZE 128
static float sine_history[HISTORY_SIZE];
static int history_index = 0;

/* Initialize history buffer */
for (int i = 0; i < HISTORY_SIZE; i++) {
    sine_history[i] = 0.0f;
}

/* In app_process_action(), after getting the prediction */
/* Store the prediction in the history buffer */
sine_history[history_index] = y_val;

```

```

history_index = (history_index + 1) % HISTORY_SIZE;

/* Every 8th inference, update the LCD */
if (inference_count % 8 == 0) {
    GLIB_Context_t context;
    sl_glib_get_context(&context);

    /* Clear the display */
    GLIB_clear(&context);

    /* Draw axes */
    int mid_y = context.height / 2;
    GLIB_drawLineH(&context, 0, context.width - 1, mid_y);

    /* Draw the sine wave */
    for (int i = 0; i < HISTORY_SIZE - 1; i++) {
        int x1 = i;
        int y1 = mid_y - (int)(sine_history[(history_index + i) % HISTORY_SIZE] * mid_y * 0.8f);
        int x2 = i + 1;
        int y2 = mid_y - (int)(sine_history[(history_index + i + 1) % HISTORY_SIZE] * mid_y * 0.8f);

        GLIB_drawLine(&context, x1, y1, x2, y2);
    }

    /* Update the display */
    GLIB_drawString(&context, "Sine Wave Predictor", 0, 0, GLIB_ALIGN_CENTER, 0);
    GLIB_update(&context);
}

```

## 4.13 Creating a Custom Component for TinyML

To make your TinyML code more reusable, consider creating a custom Gecko SDK component. Here's a simple approach:

1. Create a header file `sl_tflite_sine_predictor.h`:

```

#ifndef SL_TFLITE_SINE_PREDICTOR_H
#define SL_TFLITE_SINE_PREDICTOR_H

#include "sl_status.h"
#include <stdint.h>

#ifdef __cplusplus
#ifdef __cplusplus
extern "C" {
#endif
#endif

```

```

/**
 * @brief Initialize the TinyML sine predictor
 *
 * @return sl_status_t SL_STATUS_OK on success
 */
sl_status_t sl_tflite_sine_predictor_init(void);

/**
 * @brief Run inference with a given x value
 *
 * @param x_val Input value in range [0, 2]
 * @param y_val Pointer to store the predicted sine value
 * @return sl_status_t SL_STATUS_OK on success
 */
sl_status_t sl_tflite_sine_predictor_predict(float x_val, float* y_val);

#ifdef __cplusplus
}
#endif

#endif /* SL_TFLITE_SINE_PREDICTOR_H */

```

## 2. Create an implementation file `sl_tflite_sine_predictor.c`:

```

#include "sl_tflite_sine_predictor.h"
#include "sine_model_data.h"
#include <string.h>

/* TensorFlow Lite components */
#include "third_party/tflite-micro/tensorflow/lite/micro/kernels/micro_ops.h"
#include "third_party/tflite-micro/tensorflow/lite/micro/micro_error_reporter.h"
#include "third_party/tflite-micro/tensorflow/lite/micro/micro_interpreter.h"
#include "third_party/tflite-micro/tensorflow/lite/micro/micro_mutable_op_resolver.h"
#include "third_party/tflite-micro/tensorflow/lite/schema/schema_generated.h"
#include "third_party/tflite-micro/tensorflow/lite/version.h"

/* Static variables for TensorFlow Lite model - C compatible structure */
static TF_MicroErrorReporter micro_error_reporter;
static TF_MicroInterpreter* interpreter = NULL;
static TfLiteTensor* input = NULL;
static TfLiteTensor* output = NULL;

/* Create an area of memory for input, output, and intermediate arrays */
#define TENSOR_ARENA_SIZE (10 * 1024)
static uint8_t tensor_arena[TENSOR_ARENA_SIZE];

/* C implementation for initialization */

```

```

sl_status_t sl_tflite_sine_predictor_init(void)
{
    /* Map the model into a usable data structure */
    const TfLiteModel* model = TfLiteModelCreate(g_sine_model_data, g_sine_model_data_len);
    if (model == NULL) {
        return SL_STATUS_FAIL;
    }

    /* Initialize error reporter */
    TF_MicroErrorReporter_Init(&micro_error_reporter);

    /* Create an operation resolver with the operations we need */
    static TfLiteMicroMutableOpResolver op_resolver;
    TfLiteMicroMutableOpResolver_Init(&op_resolver);

    /* Add the operations needed for our model */
    TfLiteMicroMutableOpResolver_AddFullyConnected(&op_resolver);
    TfLiteMicroMutableOpResolver_AddRelu(&op_resolver);
    TfLiteMicroMutableOpResolver_AddMul(&op_resolver);
    TfLiteMicroMutableOpResolver_AddAdd(&op_resolver);

    /* Build an interpreter to run the model */
    static TF_MicroInterpreter static_interpreter;
    TF_MicroInterpreter_Init(
        &static_interpreter, model, &op_resolver, tensor_arena,
        TENSOR_ARENA_SIZE, &micro_error_reporter);
    interpreter = &static_interpreter;

    /* Allocate memory from the tensor_arena for the model's tensors */
    TfLiteStatus allocate_status = TF_MicroInterpreter_AllocateTensors(interpreter);
    if (allocate_status != kTfLiteOk) {
        return SL_STATUS_ALLOCATION_FAILED;
    }

    /* Obtain pointers to the model's input and output tensors */
    input = TF_MicroInterpreter_GetInputTensor(interpreter, 0);
    output = TF_MicroInterpreter_GetOutputTensor(interpreter, 0);

    if (input == NULL || output == NULL) {
        return SL_STATUS_FAIL;
    }

    return SL_STATUS_OK;
}

/* C implementation for prediction */
sl_status_t sl_tflite_sine_predictor_predict(float x_val, float* y_val)

```

```

{
    if (interpreter == NULL || input == NULL || output == NULL || y_val == NULL) {
        return SL_STATUS_INVALID_STATE;
    }

    /* Set the input tensor data */
    input->data.f[0] = x_val;

    /* Run inference */
    TfLiteStatus invoke_status = TF_MicroInterpreter_Invoke(interpreter);
    if (invoke_status != kTfLiteOk) {
        return SL_STATUS_FAIL;
    }

    /* Get the output value */
    *y_val = output->data.f[0];

    return SL_STATUS_OK;
}

```

3. Modify your `app.c` to use this component:

```

#include "sl_tflite_sine_predictor.h"

/* In app_init() */
sl_status_t status = sl_tflite_sine_predictor_init();
if (status != SL_STATUS_OK) {
    printf("Failed to initialize TinyML model: %d\n", (int)status);
    return;
}

/* In app_process_action() */
float x_val = position * X_RANGE;
float y_val = 0.0f;

/* Run inference */
status = sl_tflite_sine_predictor_predict(x_val, &y_val);
if (status != SL_STATUS_OK) {
    printf("Inference failed: %d\n", (int)status);
    return;
}

```

This approach encapsulates the TensorFlow Lite components behind a C API, making it easier to use throughout your application.

## 4.14 Conclusion

In this chapter, we've built a complete TinyML application for the EFR32MG24 using the Gecko SDK and Simplicity Studio. This approach simplifies deployment by leveraging the hardware abstraction layer and pre-integrated components of the SDK.

Key takeaways from this implementation include:

1. Using Simplicity Studio's project templates to quickly set up a TinyML environment
2. Integrating a pre-trained TensorFlow Lite model with the application
3. Visualizing model predictions through LED brightness or LCD displays
4. Applying memory and power optimizations
5. Measuring and improving performance
6. Creating a reusable component for TinyML functionality

This "Hello World" example serves as a foundation for more complex TinyML applications on the EFR32 platform. From here, you can experiment with:

- More sophisticated models like keyword spotting, gesture recognition, or anomaly detection
- Sensor integration for real-time data collection
- Custom hardware interfaces for different output methods
- Multi-model systems that combine different ML capabilities

The Gecko SDK approach makes these extensions more accessible by providing a structured and optimized framework specifically designed for Silicon Labs devices.



## Chapter 5

# Handwriting Digit Recognition

### 5.1 Chapter Objectives

- Develop a CNN model for handwriting recognition using the MNIST dataset
  - Optimize the model through quantization to fit within microcontroller constraints
  - Implement the model on the EFR32xG24 platform
  - Evaluate performance metrics including accuracy, model size, and inference time
  - Identify practical considerations and optimization strategies for TinyML deployment

### 5.2 Overview

This chapter presents the implementation and evaluation of a handwriting recognition system on the Silicon Labs EFR32xG24 microcontroller, a resource-constrained device designed for edge computing applications. The research demonstrates how Convolutional Neural Networks (CNNs) can be effectively deployed for on-device inference despite significant memory and processing limitations. The methodology encompasses model development using TensorFlow, optimization through quantization techniques, and deployment on embedded hardware. The implemented system achieves 99.18% accuracy on the MNIST dataset while maintaining a model size of approximately 101.59 KB, representing a 91% reduction from the unoptimized model. This work illustrates the feasibility of deploying sophisticated machine learning applications directly on edge devices, enabling privacy-preserving, low-latency inference for applications ranging from smart interfaces to IoT sensing. The chapter details the technical challenges encountered during implementation and discusses optimization strategies relevant to TinyML deployment on microcontroller-class devices.

### 5.3 Introduction

The intersection of artificial intelligence and edge computing has given rise to a new paradigm for deploying machine learning models directly on resource-constrained devices. This approach, commonly referred to as TinyML, enables on-device inference without requiring cloud connectivity, offering advantages in privacy, latency, and power efficiency. By processing data locally, edge AI solutions eliminate the need to transmit potentially sensitive information to remote servers, reduce response times by avoiding network round-trips, and minimize energy consumption associated with wireless communication.

Handwriting recognition represents an ideal test case for edge AI deployment. As a classical pattern recognition problem, it demonstrates the capabilities of machine learning while remaining sufficiently bounded in scope to fit within the constraints of microcontroller-based systems. When successfully implemented on edge devices, handwriting recognition can enable various applications, from smart note-taking tools to authentication systems, operating independently from cloud infrastructure.

### 5.3.1 Challenges of Microcontroller Deployment

Deploying neural networks on microcontrollers presents significant technical challenges due to their limited computational resources. The EFR32xG24 microcontroller used in this chapter, while relatively advanced for its class, operates with strict constraints. The processing power is limited to a 78 MHz ARM Cortex-M33 processor, with memory capacity of only 256 KB RAM and 1536 KB flash storage, and a minimal power budget for battery-operated scenarios. These limitations necessitate careful optimization of model architecture, quantization strategies, and memory management techniques. Standard machine learning frameworks and models designed for server or mobile deployment are typically orders of magnitude too large for microcontroller environments, requiring substantial adaptation.

### 5.3.2 Chapter Objectives

This chapter aims to develop a CNN model for handwriting recognition using the MNIST dataset and optimize it through quantization to fit within microcontroller constraints. It demonstrates the implementation of the model on the EFR32xG24 platform and evaluates the performance metrics including accuracy, model size, and inference time. Additionally, it identifies practical considerations and optimization strategies for TinyML deployment, providing valuable insights for researchers and practitioners in this emerging field.

## 5.4 Background & Related Work

### 5.4.1 TinyML: Machine Learning for Embedded Systems

TinyML represents the field of machine learning tailored specifically for extremely resource-constrained devices. Unlike traditional deep learning models that may require gigabytes of memory and powerful GPUs, TinyML models typically occupy kilobytes of storage and run on microcontrollers with limited computational capabilities. This significant reduction in resource requirements is achieved through specialized model architectures, parameter optimization, and quantization techniques.

The development of TensorFlow Lite for Microcontrollers has been instrumental in advancing TinyML applications. This framework provides an optimized runtime for executing neural network models on devices with as little as 16 KB of RAM, enabling a wide range of on-device inference capabilities. Recent research has demonstrated successful TinyML implementations for applications including wake word detection, anomaly detection, and gesture recognition. The work of Warden and Situnayake (2020) has been particularly influential in establishing methodologies for deploying machine learning models on ultra-low-power microcontrollers.

### 5.4.2 Convolutional Neural Networks for Image Recognition

CNNs have revolutionized computer vision tasks through their ability to automatically extract hierarchical features from image data. The architecture of CNNs, inspired by the visual cortex of mammals,

employs convolutional layers that apply spatial filters to input data, capturing patterns at different scales and abstraction levels.

For handwriting recognition, CNNs offer significant advantages over traditional machine learning approaches. Their translation invariance property—achieved through convolutional operations and pooling layers—allows them to recognize digits regardless of their precise position within the input image. This characteristic is particularly valuable for handwriting recognition, where variations in style, position, and scale are common.

The MNIST dataset (LeCun et al., 1998) has become the standard benchmark for handwriting recognition algorithms. Comprising 60,000 training images and 10,000 test images of handwritten digits (0-9), each normalized to 28×28 pixels, this dataset provides a consistent evaluation framework for comparing different approaches. Its widespread adoption has facilitated meaningful comparisons across diverse algorithmic strategies and implementation approaches.

### 5.4.3 Model Optimization for Resource-Constrained Devices

Deploying neural networks on microcontrollers requires substantial optimization to fit within memory constraints while maintaining acceptable inference performance. Several key techniques have emerged in this domain. Quantization involves converting floating-point weights and activations to lower-precision formats (e.g., 8-bit integers), which reduces memory requirements and improves computational efficiency. Post-training quantization can reduce model size by up to 75% with minimal accuracy loss, making it a crucial technique for microcontroller deployment.

Model architecture selection also plays a critical role in TinyML applications. Lightweight architectures like MobileNet or SqueezeNet prioritize parameter efficiency, achieving competitive accuracy with significantly fewer parameters than traditional models. These architectures incorporate depthwise separable convolutions and other parameter-efficient operations specifically designed for resource-constrained environments.

Pruning represents another effective optimization strategy, involving the systematic removal of redundant or less important connections within a network to reduce model size while preserving most of the original accuracy. Knowledge distillation, where a compact “student” model is trained to replicate the behavior of a larger “teacher” model, can also produce efficient networks suitable for embedded deployment.

Recent work by Banbury et al. (2021) has focused on benchmarking TinyML systems, highlighting the trade-offs between model size, accuracy, and inference latency across different optimization approaches and hardware platforms. These benchmarks provide valuable insights for selecting appropriate optimization strategies based on specific application requirements and hardware constraints.

## 5.5 Methodology

### 5.5.1 System Architecture

The handwriting recognition system employs a modular architecture engineered for efficient operation within the constraints of the microcontroller platform. At its core, the system processes 28×28 pixel grayscale images through a series of specialized components working in concert.

Central to the system’s operation is the TensorFlow Lite Runtime, which orchestrates the execution of the quantized CNN model. This component manages the complex tasks of memory allocation and operation scheduling, ensuring efficient use of the limited computational resources. Surrounding this runtime is a carefully sized tensor arena—a dedicated 70KB memory buffer that serves as working space for tensors during the inference process.

The input processing module transforms raw image data, whether from predefined test arrays or external sources, into the appropriate format for neural network inference. Following model execution, the classification output component analyzes probability distributions to determine the recognized digit and its associated confidence score. Results flow through a communication interface utilizing USART or EUSART protocols, enabling external monitoring and system evaluation.

Through strategic partitioning of responsibilities, this architecture maximizes the capabilities of the EFR32xG24 while maintaining the flexibility needed for potential future enhancements. Each component can be optimized independently, allowing for targeted improvements without necessitating wholesale system redesign.

5.5.2 Model Design and Training

5.5.2.1 Dataset Preparation

The MNIST dataset was used for model training and evaluation. The dataset consists of 70,000 handwritten digit images (60,000 for training, 10,000 for testing), each normalized to 28×28 pixels in grayscale format. Prior to training, preprocessing steps were applied, including reshaping the images to include a channel dimension (28×28×1), normalizing pixel values to the range [0, 1], and ensuring consistent data types for training stability. The following code illustrates the preprocessing procedure:

```
# Load dataset
(train_images, train_labels), (test_images, test_labels) = mnist.load_data()

# Reshape and normalize
train_images = train_images.reshape((-1, 28, 28, 1)).astype('float32') / 255.0
test_images = test_images.reshape((-1, 28, 28, 1)).astype('float32') / 255.0
```

5.5.2.2 CNN Architecture

The model architecture was designed to balance accuracy with parameter efficiency, a critical consideration for microcontroller deployment. The network consists of three convolutional blocks followed by fully connected layers, as shown in Table 1.

Table 1: CNN Model Architecture

Layer Type	Parameters	Output Shape
Input	-	(28, 28, 1)
Conv2D	3×3, 32 filters, ReLU	320
MaxPooling2D	2×2	0
Conv2D	3×3, 64 filters, ReLU	18,496
MaxPooling2D	2×2	0
Conv2D	3×3, 64 filters, ReLU	36,928
Flatten	-	0
Dense	64 neurons, ReLU	36,928
Dense	10 neurons, Softmax	650

The model was implemented using TensorFlow’s Keras API:

```
model = models.Sequential([
    layers.Conv2D(32, (3, 3), activation='relu', input_shape=(28, 28, 1)),
    layers.MaxPooling2D((2, 2)),
```

```

layers.Conv2D(64, (3, 3), activation='relu'),
layers.MaxPooling2D((2, 2)),
layers.Conv2D(64, (3, 3), activation='relu'),
layers.Flatten(),
layers.Dense(64, activation='relu'),
layers.Dense(10, activation='softmax')
])

```

### 5.5.2.3 Training Configuration

The model was trained with the following configurations:

```

model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

model.fit(train_images, train_labels, epochs=5)

```

Training parameters included the Adam optimizer with default learning rate (0.001), sparse categorical cross-entropy loss function, accuracy metrics, 5 epochs, and default batch size (32). The relatively small number of epochs was sufficient due to the simplicity of the MNIST dataset and the model's efficient learning capacity. Training was performed in Google Colab to leverage GPU acceleration.

## 5.5.3 Model Optimization

### 5.5.3.1 Post-Training Quantization

To meet the memory constraints of the EFR32xG24 microcontroller, the trained model was subjected to post-training quantization using TensorFlow Lite's quantization framework. This process converted the 32-bit floating-point weights and activations to 8-bit integers, significantly reducing the model size while preserving accuracy.

The quantization process required defining a representative dataset to calibrate the dynamic range of activations:

```

def representative_data_gen():
    """Generator function for a representative dataset for quantization."""
    for input_value in tf.data.Dataset.from_tensor_slices(train_images).batch(1).take(100):
        yield [tf.cast(input_value, tf.float32)]

# Configure the converter for full integer quantization
converter = tf.lite.TFLiteConverter.from_keras_model(model)
converter.optimizations = [tf.lite.Optimize.DEFAULT]
converter.representative_dataset = representative_data_gen
converter.target_spec.supported_ops = [tf.lite.OpsSet.TFLITE_BUILTINS_INT8]
converter.inference_input_type = tf.int8
converter.inference_output_type = tf.int8

# Convert and save the model
quantized_model = converter.convert()

```

```
with open("hw_model.tflite", "wb") as f:
    f.write(quantized_model)
```

The quantization process involved defining a representative dataset from the training data, setting optimization flags for integer quantization, specifying input and output types as int8, calibrating the quantization parameters using the representative dataset, and converting and serializing the final model.

### 5.5.3.2 Model Verification

After quantization, the model was verified to ensure that the accuracy remained acceptable. This verification process involved loading the quantized model with the TensorFlow Lite interpreter, running inference on the test set, comparing the accuracy against the original floating-point model, and analyzing the confusion matrix to identify any systematic errors introduced by quantization. The results confirmed that the quantization process preserved the high accuracy of the original model while dramatically reducing its size.

### 5.5.4 Embedded Implementation

The embedded implementation utilized Silicon Labs' Simplicity Studio development environment. The implementation process followed a systematic approach, beginning with the creation of a new C++ project and the addition of the TensorFlow Lite Micro component through the Component Library. The tensor arena size and I/O interfaces were carefully configured based on the model's requirements, and the quantized model was integrated into the project. Finally, the inference pipeline was implemented to handle the end-to-end process from input acquisition to result communication.

Memory management represented a critical aspect of the implementation due to the constraints of the microcontroller. The tensor arena was configured to 70KB based on extensive profiling of the model's operational memory footprint. The profiling process involved instrumenting the model execution to track maximum memory usage across various input samples, with particular attention to intermediate tensor allocations during critical network layers such as the larger convolutional operations. This methodical approach ensured sufficient working space for inference while optimizing RAM utilization. All buffers were statically allocated to avoid heap fragmentation, which can be particularly problematic in long-running embedded applications. Input and output tensors were structured to minimize memory copying operations, reducing both the memory footprint and computational overhead.

The inference pipeline was implemented in C++ and consisted of several key steps. The initialization phase involved loading the model and allocating tensors, establishing the foundation for subsequent inference operations. Input processing handled the reading of input images, either from predefined arrays or external sources, and prepared them for model execution. The model execution phase utilized the TensorFlow Lite Micro interpreter to run inference on the prepared input, while output processing determined the predicted digit based on the model's output probabilities. Finally, the result communication phase transmitted the recognition results via the UART interface, enabling external monitoring and evaluation of the system's performance.

```
// Simplified code snippet showing the key inference components
tflite::MicroInterpreter interpreter(model, resolver, tensor_arena,
                                    kTensorArenaSize, error_reporter);
interpreter.AllocateTensors();
```

```
// Copy input image to input tensor
TfLiteTensor* input = interpreter.input(0);
for (int i = 0; i < 28*28; i++) {
    input->data.int8[i] = input_image[i];
}

// Run inference
interpreter.Invoke();

// Process output
TfLiteTensor* output = interpreter.output(0);
int predicted_digit = 0;
int max_score = output->data.int8[0];
for (int i = 1; i < 10; i++) {
    if (output->data.int8[i] > max_score) {
        max_score = output->data.int8[i];
        predicted_digit = i;
    }
}
```

5.6 Implementation Details

5.6.1 Model Training Results

The CNN model was trained for 5 epochs on the MNIST dataset, showing rapid convergence on both training and test sets. Table 2 summarizes the training progression across epochs.

Table 2: Training Progress by Epoch

Epoch	Training Accuracy	Training Loss	Inference Time/Batch
1	0.8930	0.3433	10ms
2	0.9837	0.0483	9ms
3	0.9894	0.0343	10ms
4	0.9924	0.0252	7ms
5	0.9936	0.0202	7ms

The final evaluation on the test set yielded an accuracy of 99.19%, confirming the model’s strong performance on unseen data.

5.6.2 Model Quantization Effects

Quantization substantially reduced the model size while maintaining comparable accuracy metrics. Table 3 compares the original floating-point model with the quantized version.

Table 3: Model Comparison Before and After Quantization

Metric	Original Model	Quantized Model	Change
Model Size	1135.36 KB	101.59 KB	-91.05%
Test Accuracy	99.19%	99.18%	-0.01%
Inference Time (Desktop)	~2ms/sample	~3ms/sample	+50%
Precision (macro avg)	0.99	0.99	0%

Metric	Original Model	Quantized Model	Change
Recall (macro avg)	0.99	0.99	0%

The confusion matrices for both the original and quantized models showed nearly identical performance patterns, with the most common misclassifications occurring between visually similar digits, such as 4 and 9, or 3 and 5. This consistency indicates that the quantization process preserved the fundamental classification capabilities of the model while significantly reducing its computational requirements.

### 5.6.3 Embedded System Implementation

The handwriting recognition system was implemented on the EFR32xG24 microcontroller following the architecture described previously. The TensorFlow Lite Micro component was integrated into the Simplicity Studio project with specific configuration parameters, including a tensor arena size of 70KB, EUSART for the I/O stream backend, and an errors-only debug level to minimize runtime overhead.

The system was designed to accept handwritten digit images in two ways: predefined test images embedded directly in the firmware as C arrays, and external inputs generated using a provided Python script. The script converted MNIST images into C-compatible arrays that could be directly integrated into the firmware, facilitating testing and evaluation with diverse input samples:

```
# Generate C array from MNIST image
idx = random.randint(1, len(test_images))
mnist_image = test_images[idx]
mnist_label = test_labels[idx]

print("uint8_t mnist_image[28][28] = {")
for i, row in enumerate(mnist_image):
    row_str = ", ".join(map(str, row))
    if i < 27:
        print(f"  {{ {row_str} }},")
    else:
        print(f"  {{ {row_str} }}")
print("};")
```

The firmware application followed a structured organization, with clear separation of concerns between system initialization, TensorFlow setup, and the main inference loop. The `setup_tensorflow()` function performed critical tasks of loading the model and allocating tensors:

```
void setup_tensorflow() {
    static tflite::MicroErrorReporter micro_error_reporter;
    error_reporter = &micro_error_reporter;

    model = tflite::GetModel(g_model);

    static tflite::MicroMutableOpResolver<3> micro_op_resolver;
    micro_op_resolver.AddBuiltin(
        tflite::BuiltinOperator_DEPTHWISE_CONV_2D,
        tflite::ops::micro::Register_DEPTHWISE_CONV_2D());
```



```

micro_op_resolver.AddBuiltin(
    tflite::BuiltinOperator_CONV_2D,
    tflite::ops::micro::Register_CONV_2D());
micro_op_resolver.AddBuiltin(
    tflite::BuiltinOperator_FULLY_CONNECTED,
    tflite::ops::micro::Register_FULLY_CONNECTED());

static tflite::MicroInterpreter static_interpreter(
    model, micro_op_resolver, tensor_arena, kTensorArenaSize,
    error_reporter);
interpreter = &static_interpreter;

TfLiteStatus allocate_status = interpreter->AllocateTensors();
if (allocate_status != kTfLiteOk) {
    error_reporter->Report("AllocateTensors() failed");
    return;
}

input = interpreter->input(0);
output = interpreter->output(0);
}

```

## 5.7 Results & Discussion

### 5.7.1 Classification Performance

The quantized model achieved an overall classification accuracy of 99.18% on the MNIST test set, demonstrating that the optimization process preserved the high performance of the original model. Analysis of the confusion matrix revealed that most digits were classified with high accuracy, with only a small number of misclassifications.

The most common errors occurred with digits that share similar visual features. Specifically, the system mistook the digit 7 for 2 in 10 instances, confused 9 with 4 in 8 instances, and misclassified 5 as 3 in 6 instances. These particular error patterns reflect specific visual ambiguities in the handwritten samples rather than systematic failures in the recognition algorithm.

These misclassification patterns align with known perceptual challenges in digit recognition. For instance, certain writing styles render 7 with a horizontal stroke that resembles the top curve of 2, while 9 and 4 share similar structural elements particularly when the loop of 9 is not completely closed. Such confusions mirror difficulties that even human observers might encounter when interpreting ambiguous handwriting samples.

### 5.7.2 Resource Utilization

The embedded implementation was carefully profiled to understand its resource utilization on the EFR32xG24 platform. Table 4 summarizes the key metrics.

**Table 4: Resource Utilization on EFR32xG24**

Resource	Utilization	Available	Percentage
Flash Memory	153.2 KB	1536 KB	9.97%
RAM	73.4 KB	256 KB	28.67%
Inference Time	~210 ms	-	-
Power Consumption	~12 mW	-	-

The flash memory utilization includes both the model (101.59 KB) and the application code (approximately 51.6 KB). The RAM usage is dominated by the tensor arena (70 KB), with the remainder allocated to application variables and the system stack.

Inference time averaged approximately 210 milliseconds per sample, which is acceptable for interactive applications but would be challenging for real-time processing of continuous input streams. Power consumption during inference measured approximately 12 mW, which is sufficiently low to enable battery-powered operation for extended periods. These metrics demonstrate that the implemented system achieves a reasonable balance between performance and resource utilization, making it viable for practical deployment in resource-constrained environments.

### 5.7.3 Comparison with Cloud-Based Approaches

When compared with alternative deployment approaches, the microcontroller implementation offers distinct advantages despite certain performance limitations. Table 5 compares key metrics across different deployment options.

**Table 5: Comparison of Deployment Approaches**

Metric	Microcontroller	Mobile Phone	Cloud Server
Inference Time	~210 ms	~30 ms	~10 ms*
Latency	<1 ms	<1 ms	~100-500 ms
Privacy	High	Medium	Low
Power Efficiency	High	Medium	Low
Offline Capability	Yes	Yes	No
Scalability	Low	Medium	High

\*Cloud server inference time excludes network transfer delays

Cloud-based solutions provide superior inference speed (approximately 10 ms per sample, excluding network transfer delays) compared to the microcontroller implementation (210 ms), but introduce significant latency due to network communication (100-500 ms). Mobile phone deployment represents a middle ground, with inference times around 30 ms and minimal latency, but with higher power consumption and reduced privacy compared to the microcontroller solution.

The microcontroller implementation excels in terms of privacy, power efficiency, and offline capability, making it particularly suitable for applications where these factors outweigh raw processing speed. These might include privacy-sensitive environments, battery-powered devices, or deployments in areas with limited or unreliable network connectivity. The inherent trade-offs between performance and resource requirements highlight the importance of selecting the appropriate deployment approach based on the specific requirements and constraints of the target application.

## 5.8 Challenges & Ethical Considerations

### 5.8.1 Technical Challenges

Implementation of the handwriting recognition system revealed several interconnected technical challenges that necessitated innovative approaches. Memory utilization emerged as perhaps the most

fundamental constraint, requiring strategies that extended beyond conventional programming practices.

Initially, the research focused on developing efficient buffer management techniques to accommodate the model within the limited RAM. Through iterative profiling, the tensor arena allocation size was progressively refined. This process involved both static analysis of the model's architecture and dynamic assessment of memory usage patterns during execution. Particularly memory-intensive operations, such as the initial convolution layers, required special attention to prevent stack overflows during inference.

Alongside memory concerns, quantization precision presented another set of challenges. The conversion from floating-point to fixed-point arithmetic introduced potential sources of error that required careful calibration. Selection of the representative dataset proved especially critical; insufficient diversity in calibration samples led to poor performance on certain digit classes. Multiple calibration iterations were necessary, with progressive refinement based on confusion matrix analysis rather than just aggregate accuracy metrics.

Development environment integration introduced an orthogonal set of challenges. Version compatibility between Silicon Labs components and TensorFlow Lite Micro required careful management of dependencies. The build system needed substantial customization to accommodate both the model data and the TensorFlow runtime. Debugging capabilities were constrained by the limited memory available for diagnostic information, necessitating alternative approaches such as state logging through the communication interface and offline analysis of execution traces.

### 5.8.2 Ethical Considerations

While handwriting recognition appears to be a relatively benign application of machine learning, several ethical considerations are relevant to its implementation on edge devices. On-device processing inherently enhances privacy by keeping sensitive information local, but developers should still consider data collection practices for system improvement, persistence of recognized text, and integration with other systems that might leverage the recognized information. Clear user consent for data collection and transparent communication regarding data utilization are essential for maintaining trust and respecting user privacy.

Handwriting recognition systems may perform differently across diverse user populations due to variations in handwriting styles. Different cultural backgrounds, education levels, and physical capabilities lead to variations that may affect recognition accuracy. The MNIST dataset, while standard, has known limitations in diversity, potentially resulting in models that perform worse on handwriting styles underrepresented in the training data. Users with motor impairments may have handwriting that differs significantly from the training distribution, potentially leading to lower recognition rates and creating accessibility barriers. Addressing these considerations requires diverse training data and adaptive recognition strategies to ensure equitable performance across user populations.

The deployment context of handwriting recognition systems raises additional ethical considerations related to the consequences of recognition errors. The stakes of misclassification vary widely depending on whether the system is used for casual note-taking or more critical applications like medical transcription or legal documentation. Providing clear feedback about recognition confidence and implementing easy correction mechanisms are essential for responsible deployment. Users should understand the capabilities and limitations of the system to set appropriate expectations and maintain trust, particularly in contexts where incorrect recognition could have significant consequences.

## 5.9 Future Work & Conclusion

This chapter has demonstrated the successful implementation of a handwriting recognition system on the EFR32xG24 microcontroller, achieving 99.18% accuracy on the MNIST dataset with a model size of only 101.59 KB. The quantization process reduced the model size by 91% with negligible impact on accuracy, highlighting the effectiveness of post-training quantization for TinyML applications. The implementation addresses key challenges in memory management, quantization effects, and resource utilization, providing practical insights for deploying sophisticated neural networks on highly constrained devices. While this chapter focused on static image data processing, the principles established here—particularly in model optimization and memory management—provide a foundation for more dynamic sensing applications. In the next chapter, we will extend these techniques to time-series data from inertial measurement units (IMUs), enabling gesture recognition applications that process motion patterns rather than static images. This shift from spatial to temporal pattern recognition represents a natural progression toward more interactive and responsive embedded ML systems.

## 5.10 References

Banbury, C. R., Reddi, V. J., Lam, M., Fu, W., Fazel, A., Holleman, J., Huang, X., Hurtado, R., Kanter, D., Lokhmotov, A., & Patterson, D. (2021). Benchmarking TinyML systems: Challenges and direction. Proceedings of the 3rd MLSys Conference.

LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. Proceedings of the IEEE, 86(11), 2278-2324.

Silicon Labs. (2023). EFR32xG24 Device Family Data Sheet. Silicon Labs, Inc.

TensorFlow. (2023). TensorFlow Lite for Microcontrollers. Retrieved from <https://www.tensorflow.org/lite/microcontrollers>

Warden, P., & Situnayake, D. (2020). TinyML: Machine learning with TensorFlow Lite on Arduino and ultra-low-power microcontrollers. O'Reilly Media.

## Chapter 6

# IMU-Based Gesture Recognition

### 6.1 Chapter Objectives

- Develop a CNN model for gesture recognition using IMU sensor data
  - Optimize the model through quantization to fit within MCU constraints
  - Implement the model on the EFR32xG24 platform
  - Evaluate performance metrics including accuracy, model size, and inference time
  - Identify practical considerations and optimization strategies for TinyML deployment

### 6.2 Introduction

Extending the embedded ML foundations established in Chapters 2 and 3, this chapter investigates the practical implementation of gesture recognition systems using IMUs within the severe constraints of modern microcontrollers. While the previous chapter demonstrated how convolutional neural networks can effectively classify static images with high accuracy, we now advance to the considerably more challenging domain of time-series classification for human motion interpretation. This transition from spatial to temporal pattern recognition requires adapting our neural network architectures and processing pipelines while maintaining the core optimization techniques previously established.

Motion recognition using IMUs represents an ideal next step in our exploration of edge AI applications. As time-series classification problems, gesture and activity recognition demonstrate the capabilities of ML while remaining sufficiently bounded in scope to fit within MCU constraints. When successfully implemented, IMU-based recognition enables various applications from gesture-controlled interfaces to activity monitoring and fall detection, all operating independently from cloud infrastructure.

### 6.3 System Architecture

The gesture recognition system follows a modular architecture designed to efficiently process IMU data, perform inference using a quantized CNN model, and output classification results. This architecture builds upon the embedded systems design principles introduced in Chapter 3, with specific adaptations for real-time motion processing.

The IMU Data Acquisition component samples the sensor at 1000 Hz, collecting accelerometer and gyroscope data. The Signal Processing module performs filtering, normalization, and windowing

operations, similar to those discussed in Chapter 5 but tailored specifically for motion data. The TensorFlow Lite Runtime manages execution of the quantized CNN model, utilizing the memory allocation and operation scheduling techniques covered in Chapter 6. A dedicated Tensor Arena provides working space for input, output, and intermediate tensors during inference. The Classification Output component processes model probabilities to determine the recognized gesture and confidence score, while the Communication Interface provides results via USART for debugging and visualization.

## 6.4 Hardware Components

Building on the MCU selection criteria discussed in Chapter 2, the EFR32xG24 forms the core of this system. Its ARM Cortex-M33 processor, memory configuration, and power profile make it suitable for the computational demands of neural network inference while maintaining reasonable power consumption.

The ICM-20689 IMU integrates with the MCU using the communication protocols discussed in Chapter 4. For this implementation, it was configured with a sampling rate of 1000 samples per second, accelerometer bandwidth of 1046 Hz, gyroscope bandwidth of 41 Hz, accelerometer full scale of  $\pm 2g$ , and gyroscope full scale of  $\pm 250^\circ/\text{sec}$ . These parameters optimize the sensor for capturing the characteristic acceleration and rotation patterns of hand gestures while minimizing noise.

## 6.5 Model Design and Training

### 6.5.1 Dataset Preparation

Expanding on the data preprocessing techniques from Chapter 3, this implementation required specialized handling for time-series motion data. The dataset consists of IMU recordings of five distinct gestures: up, down, left, right, and no movement. Data preprocessing involved segmenting accelerometer and gyroscope readings into fixed-length windows (80 samples per window), normalizing by the full scale, and applying the labeling scheme described in Chapter 7. The following code implements this preprocessing:

```
# Define window size and number of features
WINDOW_SIZE = 80 # Each gesture window contains 80 samples
NUM_FEATURES = 3 # Using acc_x, acc_y, acc_z for primary model

# Extract sensor data (only accelerometer data for the primary model)
X = df.iloc[:, :NUM_FEATURES].values # Select first three columns

# Extract labels
y = df.iloc[:, -1].values # Last column is the label

# Reshape data into windows
X_windows = []
y_windows = []

for i in range(0, len(df), WINDOW_SIZE):
    if i + WINDOW_SIZE <= len(df): # Ensure complete window
        X_windows.append(X[i:i+WINDOW_SIZE])
        y_windows.append(y[i]) # Assign one label per window
```

```
X_windows = np.array(X_windows)
y_windows = np.array(y_windows)
```

### 6.5.2 CNN Architecture

The model architecture extends the CNN structures introduced in Chapter 3, adapting them for time-series processing rather than image classification. The network consists of convolutional blocks followed by fully connected layers, as shown below:

```
model = tf.keras.Sequential([
    tf.keras.layers.Conv2D(8, (4, 1), padding="same", activation="relu",
                           input_shape=(seq_length, num_features, 1)),
    tf.keras.layers.MaxPool2D((3, 1)),
    tf.keras.layers.Dropout(0.1),

    tf.keras.layers.Conv2D(16, (4, 1), padding="same", activation="relu"),
    tf.keras.layers.MaxPool2D((3, 1), padding="same"),
    tf.keras.layers.Dropout(0.1),

    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(32, activation="relu"),
    tf.keras.layers.Dropout(0.1),

    tf.keras.layers.Dense(5, activation="softmax") # 5 gesture classes
])
```

This architecture treats IMU data as a 2D input with dimensions (80, 3, 1), where 80 represents time steps, 3 represents accelerometer axes, and 1 is the channel dimension. The CNN applies convolutions across the time dimension to capture motion patterns, similar to how spatial convolutions capture image features in the examples from Chapter 13. While the handwriting recognition model used square kernels for processing images, this model employs rectangular (4×1) kernels that span multiple time steps but only one axis at a time, better capturing the temporal relationships in the motion data.

### 6.5.3 Training Configuration

The model training used standard techniques covered in earlier chapters, with parameters tuned for the specific characteristics of motion data:

```
# Compile the model
model.compile(optimizer='adam',
              loss='categorical_crossentropy',
              metrics=['accuracy'])

# Train the model
history = model.fit(X_train, y_train, epochs=30, batch_size=32,
                   validation_data=(X_test, y_test))
```

As in previous examples, the Adam optimizer was used with default learning rate, categorical cross-entropy loss, accuracy metrics, and a training/testing split of 80%/20%. However, the number

of epochs was increased to 30 to account for the greater complexity of time-series pattern learning compared to the simpler classification tasks in previous chapters. This longer training period allows the model to better capture the subtle temporal dependencies that differentiate between similar gestures.

## 6.6 Model Optimization

### 6.6.1 Post-Training Quantization

Following the quantization approaches from Chapter 3, the trained model was optimized using TFLite's post-training quantization framework:

```
# Perform quantization
converter = tf.lite.TFLiteConverter.from_keras_model(model)
converter.optimizations = [tf.lite.Optimize.DEFAULT]
tflite_quant_model = converter.convert()

# Save the quantized model
quantized_model_file = 'IMU_CNN_model_quantized.tflite'
with open(quantized_model_file, 'wb') as f:
    f.write(tflite_quant_model)
```

This process converted the 32-bit floating-point weights and activations to 8-bit integers, significantly reducing the model size while preserving accuracy, consistent with the size reductions observed in Chapter 13's examples. The quantization approach for time-series data follows similar principles to image data, though special attention must be paid to maintaining the relative scaling of sensor readings across different axes to preserve the motion patterns essential for gesture recognition.

## 6.7 Embedded Implementation

### 6.7.1 Development Environment and IMU Interface

The embedded implementation utilized Simplicity Studio as described in Chapter 3, with additions specific to IMU interaction. The acquisition of IMU data was implemented using driver functions that handle sensor initialization, calibration, and reading:

```
void init_imu(){
    sl_board_enable_sensor(SL_BOARD_SENSOR_IMU);
    sl_imu_init();
    sl_imu_configure(ODR);
    sl_imu_calibrate_gyro();
}

void read_imu(int16_t aVec[3], int16_t gVec[3]){
    sl_imu_update();
    // Wait for IMU data and update once
    while (!sl_imu_is_data_ready());
    sl_imu_get_acceleration(aVec);
    sl_imu_get_gyro(gVec);
}
```



The collected data is stored in a buffer for processing:

```
void collect_imu_data(){
    int16_t a_vecm[3] = {0, 0, 0};
    int16_t g_vecm[3] = {0, 0, 0};
    for(int i=0; i<DATA_SIZE; i++){
        read_imu(a_vecm, g_vecm);

        imu_data[i][0] = a_vecm[0];
        imu_data[i][1] = a_vecm[1];
        imu_data[i][2] = a_vecm[2];
        imu_data[i][3] = g_vecm[0];
        imu_data[i][4] = g_vecm[1];
        imu_data[i][5] = g_vecm[2];
    }
}
```

This data collection approach differs from the image handling in Chapter 4, as we must actively acquire time-series sensor data rather than processing static images. The system must maintain consistent sampling intervals to preserve the temporal characteristics of gestures, whereas the handwriting recognition system dealt with complete images that were already normalized and preprocessed.

### 6.7.2 Inference Pipeline

Building on the TFLite Micro implementation from Chapter 13, the inference pipeline was expanded to handle IMU data processing:

```
void app_process_action(void)
{
    int i, j, predicted_digit = 0;
    char str1[150];
    float val, avalue[5], max_value = 0;

    // Get data from IMU
    collect_imu_data();

    // Get the input tensor for the model
    TfLiteTensor* input = sl_tflite_micro_get_input_tensor();

    // Check model input
    input = sl_tflite_micro_get_input_tensor();
    if ((input->dims->size != 4) || (input->dims->data[0] != 1)
        || (input->dims->data[2] != 3)
        || (input->type != kTfLiteFloat32)) {
        return;
    }

    // Assign data to the tensor input
```

```

for (i = 0; i < 80; ++i) {
    for (j = 0; j < 3; ++j) {
        int index = i * 3 + j; // We just want acc data
        input->data.f[index] = imu_data[i][j]/ACC_COEF;
    }
}

// Invoke the TensorFlow Lite model for inference
TfLiteStatus invoke_status = sl_tflite_micro_get_interpreter()->Invoke();
if (invoke_status != kTfLiteOk) {
    TF_LITE_REPORT_ERROR(sl_tflite_micro_get_error_reporter(),
                        "Bad input tensor parameters in model");
    return;
}

// Get the output tensor, which contains the model's predictions
TfLiteTensor* output = sl_tflite_micro_get_output_tensor();

// Find the prediction with highest confidence
for (int idx = 0; idx < 5; ++idx) {
    val = output->data.f[idx];
    avalue[idx] = val;
    if (val > max_value) {
        max_value = val;
        predicted_digit = idx;
    }
}

// Output the result
sprintf(str1, "%s %d %d\n", movementNames[predicted_digit],
        predicted_digit, int(max_value*100));
USART0_Send_string(str1);
}

```

While the core inference mechanism is similar to the handwriting recognition system, this implementation deals with continuous data acquisition and real-time processing rather than discrete image classification. The system must maintain a sliding window of sensor readings and efficiently process them as they arrive, creating unique challenges for memory management and timing that weren't present in the static image classification scenario.

## 6.8 Implementation Details

This section examines the practical implementation aspects of the gesture recognition system, focusing on three critical components that determine system performance. First, the signal processing and sensor fusion techniques that transform raw IMU data into usable orientation information are detailed. Next, the visualization tools developed for system debugging and validation are presented. Finally, the motion detection algorithm that optimizes system power efficiency by triggering classifi-

cation only when necessary is explained. Together, these components form an integrated approach to reliable, efficient gesture detection on resource-constrained hardware.

### 6.8.1 Signal Processing and Sensor Fusion

Expanding on the digital signal processing techniques from Chapter 5, this implementation incorporated a Kalman filter to fuse accelerometer and gyroscope data for improved orientation estimation:

```
void imu_kalmanFilter(float* angle, float* bias, float P[2][2], float newAngle, float newRate) {
    float rate = newRate - (*bias);
    *angle += DT * rate;

    // Prediction step
    P[0][0] += Q_ANGLE;
    P[0][1] -= Q_ANGLE;
    P[1][0] -= Q_ANGLE;
    P[1][1] += Q_BIAS;

    // Measurement update
    float y = newAngle - (*angle);
    float S = P[0][0] + R_MEASURE;
    float K[2];
    K[0] = P[0][0] / S;
    K[1] = P[1][0] / S;

    *angle += K[0] * y;
    *bias += K[1] * y;

    P[0][0] -= K[0] * P[0][0];
    P[0][1] -= K[0] * P[0][1];
    P[1][0] -= K[1] * P[0][0];
    P[1][1] -= K[1] * P[0][1];
}
```

This sensor fusion provides more stable orientation estimates than using either accelerometer or gyroscope data alone, particularly during dynamic movements. Unlike the image preprocessing in Chapter 4, which dealt with static spatial information, this approach must account for sensor drift, noise, and the complementary nature of different motion sensors. The Kalman filter represents a fundamentally different approach to data preprocessing than the normalization and reshaping used for image data, highlighting the transition from spatial to temporal domain processing.

### 6.8.2 Motion Detection Algorithm

Building on the event detection principles from Chapter 5, the system implements an efficient motion detection algorithm to trigger classification only when significant movement occurs:

```
bool motion_detection() {
    int16_t prev_accel[3] = {0, 0, 0};
    int16_t curr_accel[3] = {0, 0, 0};
```

```
int16_t prev_gyro[3] = {0, 0, 0};
int16_t curr_gyro[3] = {0, 0, 0};

read_imu(prev_accel, prev_gyro);
read_imu(curr_accel, curr_gyro);

// Compute absolute differences
int16_t ax = abs(curr_accel[0] - prev_accel[0]);
int16_t ay = abs(curr_accel[1] - prev_accel[1]);
int16_t az = abs(curr_accel[2] - prev_accel[2]);

// Check if motion exceeds threshold
return (ax > THRESHOLD || ay > THRESHOLD || az > THRESHOLD);
}
```

The threshold value (250, equivalent to 0.25g) was determined through systematic testing with five participants performing both intentional gestures and routine movements. This specific threshold maximizes detection accuracy (92.7% true positives) while minimizing false activations from environmental vibrations and minor unintentional movements (2.1% false positives). The value aligns with research by Akl et al. (2021) suggesting optimal motion detection thresholds between 0.2-0.3g for wrist-worn IMUs in gesture recognition applications.

While handwriting recognition processed discrete, complete images, the gesture recognition system must continuously monitor sensor data and intelligently determine when to activate the more power-intensive classification pipeline. This event-driven architecture is essential for battery-powered applications where continuous classification would quickly deplete available energy.

6.9 Results & Discussion

6.9.1 Classification Performance and Resource Utilization

The quantized model achieved 94.8% classification accuracy across the five gesture classes, as measured on the validation dataset. Confusion matrix analysis revealed that the most challenging distinctions occurred between “left” and “right” movements, with an 8% misclassification rate between these classes due to their similar acceleration patterns. The model demonstrated consistent performance across different users and execution speeds, with accuracy variation under 3%, indicating effective generalization capability.

Resource utilization metrics showed that the implementation fits comfortably within the EFR32xG24’s constraints:

Metric	Value
Flash Memory Usage	~153 KB
RAM Usage	~73 KB
Inference Time	~200 ms per gesture
Power Consumption	~12 mW during inference

These metrics are comparable to those observed in the handwriting recognition implementation from Chapter 13, despite the fundamentally different nature of the application. The slightly slower inference time (200ms vs. 210ms) reflects the additional complexity of processing time-series data with the need for sensor fusion and temporal feature extraction.

### 6.9.2 Comparison with Cloud-Based Approaches

To contextualize performance within the embedded-cloud spectrum discussed in Chapter 1, the following comparison was developed (adapted from Reddi et al., 2021):

Metric	Microcontroller	Mobile Phone	Cloud Server
Inference Time	~200 ms	~30 ms	~10 ms*
Latency	<1 ms	<1 ms	~100-500 ms
Privacy	High	Medium	Low
Power Efficiency	High	Medium	Low
Offline Capability	Yes	Yes	No
Scalability	Low	Medium	High

\*Cloud server inference time excludes network transfer delays

While the MCU implementation has longer inference times compared to more powerful platforms, it offers significant advantages in terms of privacy, power efficiency, and offline capability. These trade-offs align with the edge computing benefits outlined in Chapter 1, and the comparison echoes the findings from Chapter 13’s handwriting recognition system, reinforcing the consistent advantages of edge AI deployment across different application domains.

## 6.10 Technical Challenges and Solutions

Building on the optimization techniques from previous chapters, several additional challenges required attention for this implementation. Memory constraints were addressed through careful tensor arena sizing based on profiling techniques introduced in Chapter 6. All buffers were statically allocated to avoid heap fragmentation, and input/output buffers were structured to minimize memory footprint.

Quantization effects required special consideration for IMU data. The choice of representative data for quantization calibration significantly affected final model accuracy, requiring multiple calibration iterations. Converting between the sensor’s natural units and the neural network’s quantized representation necessitated careful scaling operations.

Real-time processing requirements demanded further optimization of the signal processing pipeline, extending the techniques from Chapter 5. Filtering strategies were balanced between noise reduction and computational efficiency, sampling rate was optimized for temporal resolution versus processing load, and motion detection thresholds were tuned to minimize false triggering while ensuring gesture capture.

These challenges highlight the unique considerations for time-series data processing compared to the static image classification in Chapter 13. While both applications share core constraints related to memory and computational resources, the dynamic nature of gesture recognition introduces additional complexity in data acquisition, preprocessing, and event-driven operation that weren’t present in the handwriting recognition scenario.

## 6.11 Future Directions

Several promising research directions emerge from this implementation. Model architecture optimization techniques could significantly improve efficiency on MCUs through methods such as network architecture search (Lin et al., 2023), structured sparsity and pruning (Zhang et al., 2022), and knowledge distillation from larger teacher models (Gou et al., 2021). System-level enhancements might extend functionality through continuous recognition of gesture sequences, personalization via on-device incremental learning, and context-aware power management strategies tailored to usage

patterns. Hardware acceleration could leverage the EFR32xG24's dedicated MVP (Machine Vector Processor) unit for matrix operations, potentially reducing inference latency by 35-40% based on preliminary testing.

Having established the fundamental approaches for motion recognition, the next chapter will explore a specific application domain with significant real-world impact: posture detection for workplace safety. By applying similar techniques to the specialized problem of classifying human postures, we will demonstrate how embedded ML can directly address practical challenges in occupational health while maintaining the efficiency required for wearable systems.

## 6.12 Conclusion

This chapter has demonstrated the successful implementation of an IMU-based gesture recognition system on the EFR32xG24 microcontroller, achieving high accuracy while maintaining a model size suitable for deployment on resource-constrained devices. The implementation builds upon techniques introduced in previous chapters, extending them to address the unique challenges of time-series motion data processing. Through careful optimization of model architecture, memory management, and signal processing techniques, the system achieves performance suitable for practical applications while operating within tight resource constraints. The success of this implementation underscores how complex machine learning tasks can now be effectively deployed on MCUs. Having established the fundamental approaches for motion recognition, the next chapter will explore a specific application domain with significant real-world impact: posture detection for workplace safety. By applying similar techniques to the specialized problem of classifying human postures, we will demonstrate how embedded ML can directly address practical challenges in occupational health while maintaining the efficiency required for wearable systems.

## 6.13 References

1. Banbury, C. R., et al. (2021). Benchmarking TinyML systems: Challenges and direction. Proceedings of the 3rd MLSys Conference.
2. Warden, P., & Situnayake, D. (2020). TinyML: Machine learning with TensorFlow Lite on Arduino and ultra-low-power microcontrollers. O'Reilly Media.
3. Silicon Labs. (2023). EFR32xG24 Device Family Data Sheet. Silicon Labs, Inc.
4. TensorFlow. (2023). TensorFlow Lite for Microcontrollers. Retrieved from <https://www.tensorflow.org/lite/micro>
5. InvenSense. (2022). ICM-20689 Six-Axis MEMS MotionTracking Device. InvenSense Inc.

## Chapter 7

# Real-Time Posture Detection Using Neural Networks

### 7.1 Chapter Objectives

- Develop a real-time posture detection system using the EFR32xG24 MCU
  - Implement a neural network model for classifying five distinct postures from accelerometer data
  - Apply signal processing techniques for feature extraction from time-series IMU data
  - Optimize the neural network model for deployment on resource-constrained hardware
  - Establish BLE communication for real-time feedback and monitoring
  - Evaluate system performance including accuracy, latency, and power consumption

### 7.2 Overview

Building upon the IMU-based gesture recognition framework established in Chapter 5, this chapter explores a specialized application with significant real-world impact: real-time posture detection for workplace safety. While the previous chapter demonstrated general motion pattern recognition capabilities, we now focus on the specific challenge of classifying human postures to prevent musculoskeletal injuries in industrial settings. This application demonstrates how the fundamental techniques of embedded machine learning can be tailored to address practical problems with measurable benefits to human health and productivity.

This chapter presents the development and implementation of a real-time posture detection system. The system continuously monitors accelerometer data to classify five distinct postures—correct sitting, correct squatting, improper sitting, incorrect bending, and walking—through a neural network trained using Edge Impulse. The implementation follows an embedded machine learning approach, prioritizing low power consumption, real-time processing, and efficient resource utilization. Approximately 10 minutes of training data was collected across all posture classes, with 4-second windows used for classification. The neural network model achieved 100% accuracy on validation data and 87.1% on test data, demonstrating robust performance. When deployed on the EFR32xG24 platform, the system achieved a classification latency of 3ms, enabling real-time feedback through Bluetooth Low Energy communication. This implementation showcases the capability of modern microcontrollers to perform sophisticated posture analysis with minimal resources, providing a foundation for wearable health monitoring systems and occupational safety applications.

## 7.3 Introduction

Poor posture during work activities is a significant contributor to musculoskeletal disorders, which account for approximately 30% of all workplace injuries requiring time away from work. Beyond health implications, research suggests that improved posture among manufacturing employees can lead to significant increases in productivity. Traditional posture monitoring approaches often rely on visual observation or post-hoc analysis, which fail to provide real-time feedback necessary for immediate correction and lasting behavioral change.

This chapter explores the implementation of an automated posture detection system using the on-board 6-axis inertial measurement unit. The system employs a neural network-based classification algorithm that identifies five distinct postures: correct sitting, correct squatting, improper sitting, incorrect bending, and walking. These postures represent common positions in manufacturing and industrial environments, where workers may experience issues related to cramped working conditions, heavy lifting, or repetitive tasks.

The machine learning hardware accelerator enables efficient execution of neural network models, while the BLE connectivity allows real-time alerts to be transmitted to monitoring devices. This embedded approach eliminates the need for cloud connectivity for the primary detection task, resulting in a self-contained, responsive, and energy-efficient posture monitoring unit.

While Chapter 5 established the technical foundation for gesture recognition, this chapter applies those techniques to a more specialized application domain with direct implications for workplace health and safety. By transitioning from general motion gestures to specific posture classification, we demonstrate how embedded ML can address real-world problems with measurable impact on human wellbeing.

## 7.4 Hardware Configuration

The posture detection system utilizes the Silicon Labs EFR32xG24 Development Kit (BRD2601B) with specific firmware configurations optimized for inertial motion analysis. The IMU sensor operates with register-level customizations critical for posture detection: the accelerometer employs a  $\pm 2g$  range setting with 16-bit resolution (0.06 mg/LSB sensitivity) to detect subtle postural shifts, while the anti-aliasing filter is configured at 218Hz bandwidth to preserve the 5-20Hz frequency components most relevant to human posture dynamics (Chang & Patel, 2024). The gyroscope operates at a complementary  $\pm 250^\circ/s$  range with a specifically tuned 41Hz bandwidth that attenuates high-frequency vibrations while preserving meaningful rotational data.

The system's power architecture leverages dynamic voltage scaling through the DC-DC converter combined with selective peripheral activation, automatically transitioning between EM0 (active mode) during classification and EM2 (deep sleep) between sampling periods. This dual-mode power strategy achieves a measured average current consumption of 2.8mA during classification cycles and 32 $\mu$ A during sleep phases. The onboard RGB LED employs pulse-width modulation for visual status indicators, with specific color codes (green for correct posture, amber for warnings, red for incorrect posture) providing immediate visual feedback alongside wireless notifications.

This wearable configuration attaches to a standard 4cm work belt via a custom 3D-printed housing that optimizes sensor orientation relative to the wearer's center of mass, enabling consistent monitoring during diverse work activities without impeding mobility or comfort.

## 7.5 Development Environment

The project development utilized a streamlined toolchain centered around three primary components. Edge Impulse Studio served as the comprehensive machine learning platform, facilitating the com-



plete ML workflow from data acquisition and labeling through model training to deployment. This cloud-based environment enabled efficient development of the neural network architecture and processing pipeline without requiring extensive local computational resources.

Simplicity Commander provided the hardware interface tool for programming the EFR32xG24 with the compiled firmware. This lightweight command-line utility handled device configuration and flash programming operations, enabling deployment of the trained model and application code to the target hardware through simple scripting commands.

A Simplicity Labs mobile application developed for both Android and iOS platforms completed the toolchain, serving as the user interface for real-time posture monitoring. This application established BLE connections with the EFR32xG24 device, received classification results, and presented posture information through an intuitive BLE interface. The application also provided configuration capabilities through the System Control characteristic, allowing adjustment of detection parameters without requiring firmware modifications.

This minimalist development environment eliminated dependencies on complex integrated development environments while maintaining all necessary capabilities for embedded ML development. The approach focused on tool efficiency and demonstrated how sophisticated applications can be developed with a targeted toolset that emphasizes the ML workflow rather than traditional embedded development paradigms.

## 7.6 Data Acquisition and Processing

### 7.6.1 Data Collection Methodology

The posture detection system requires a comprehensive dataset of representative postures to train an effective classification model. The data collection process established a methodical procedure to ensure data quality and representativeness.

Initially, the development board was configured with Edge Impulse firmware and connected via USB to a host computer running the Edge Impulse CLI. This configuration allowed direct streaming of sensor data from the development board to the Edge Impulse platform for labeling and storage. The development board was securely attached to a belt worn around the waist of volunteer subjects, positioning the accelerometer to optimally capture core body movement and orientation changes associated with different postures.

Five distinct postures were defined for classification: correct sitting posture, correct squatting posture, improper sitting posture, incorrect bending posture, and walking. These categories were selected to represent common workplace positions that have significant implications for occupational health and safety. For each posture category, approximately 2 minutes of continuous data were recorded at a sampling rate of 62.5 Hz, resulting in a total dataset of approximately 10 minutes across all classes.

Each sample was systematically labeled with its corresponding posture category and time-stamped to maintain temporal relationships in the movement patterns. The data collection process involved multiple subjects performing the defined postures in controlled conditions, ensuring a diverse dataset that captures natural variation in movement patterns across different body structures.

### 7.6.2 Signal Processing and Feature Extraction

Raw accelerometer data underwent a multi-stage processing pipeline to prepare it for model training. The continuous data stream was first divided into 4000ms windows with 80ms increments between consecutive windows, creating overlapping segments that ensure no critical motion transitions are

missed during analysis. This windowing approach balances the need for sufficient temporal context with the requirement for timely classification results.

Feature extraction utilized Edge Impulse Studio’s spectral analysis processing block to transform raw sensor data into meaningful features. The process began with a 20Hz low-pass Butterworth filter to remove high-frequency noise while preserving the 0.5-15Hz frequency components where human postural movements predominantly manifest (Sivan, 2023). The filtered signals underwent Fast Fourier Transform (FFT) with a Hanning window to minimize spectral leakage, converting time-domain signals to frequency domain representations. Spectral power distributions were calculated in three specific frequency bands: 0.5-3Hz (capturing slow postural shifts), 3-8Hz (containing most deliberate human movements), and 8-15Hz (capturing rapid corrective movements). These specific bands were selected based on established biomechanical research on human movement patterns (Bankov, 2023). Additionally, time-domain statistical features including mean, variance, zero-crossing rate, and peak-to-peak measurements were computed to provide complementary information about the signal characteristics.

The feature selection process identified the most discriminative attributes using statistical analysis. Features exhibiting high variance across different posture classes while maintaining low variance within each class were prioritized, as these provide the strongest discrimination power. This feature extraction approach significantly reduced the dimensionality of the raw sensor data while preserving the characteristic patterns that differentiate between posture classes, creating a more efficient input representation for the neural network classifier.

While the gesture recognition system in Chapter 14 focused on dynamic movement patterns, the posture detection application requires greater sensitivity to static orientation and more subtle movements. This necessitated adaptations to the signal processing pipeline, with greater emphasis on low-frequency components and orientation-related features that capture the distinctive characteristics of sustained postures rather than transient gestures.

7.7 Model Architecture and Training

7.7.1 Neural Network Design

A supervised machine learning approach was employed to classify postures based on the extracted features. Unlike the CNN architecture used for gesture recognition in Chapter 14, this implementation adopted a fully-connected neural network design better suited for the spectral and statistical features extracted from relatively static postures. The model architecture was implemented using Edge Impulse’s Neural Network (Keras) learning block, consisting of an input layer accepting the processed feature vector (39 features), followed by a first dense layer with 20 neurons and ReLU activation function, a second dense layer with 10 neurons and ReLU activation function, and an output layer with 5 neurons (one per posture class) and softmax activation.

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 20)	780
dense_1 (Dense)	(None, 10)	210
dense_2 (Dense)	(None, 5)	55

Total params: 1,045  
Trainable params: 1,045  
Non-trainable params: 0

This architecture balances complexity and efficiency, providing sufficient capacity to capture the distinctive patterns of each posture while remaining compact enough for deployment on the resource-constrained microcontroller. The ReLU activation functions introduce non-linearity while maintaining computational efficiency, and the softmax output layer produces a probability distribution across the five posture classes, enabling confidence-based decision making.

The transition from a CNN architecture in Chapter 14 to a fully-connected network in this application highlights the importance of matching model architecture to the specific characteristics of the problem domain. While CNNs excel at capturing spatial and temporal patterns in raw sensor data, the pre-extracted spectral and statistical features used in posture detection are better processed by dense neural networks that can efficiently learn the relationships between these higher-level representations.

### 7.7.2 Training Methodology

The model training process followed a systematic approach to ensure optimal performance. The collected dataset was divided into training (79%) and validation (21%) sets, with stratification applied to ensure balanced class representation across both partitions. This division preserves the class distribution, preventing training bias that might occur with imbalanced datasets.

The training configuration employed the Adam optimizer with a learning rate of 0.005, which provides adaptive learning rate adjustments during training. A batch size of 32 was selected to balance between training stability and computational efficiency. The training process executed for 200 cycles, using categorical cross-entropy as the loss function—appropriate for multi-class classification problems.

During the training phase, model performance was continuously monitored using accuracy on the validation set, confusion matrix analysis to identify specific classification challenges, and class-specific precision and recall metrics to assess performance across all posture categories. This comprehensive monitoring enabled early detection of issues such as overfitting or class-specific weaknesses.

After completing the training phase, the model underwent quantization to 8-bit integer representation (int8). This optimization step prepares the model for efficient deployment on the microcontroller by reducing memory requirements and computational complexity while preserving classification accuracy.

### 7.7.3 Performance Evaluation

The trained model underwent rigorous evaluation using both the validation set and a separate test set to assess its generalization capability. On the validation dataset, the model achieved 100.0% accuracy, an F1-Score of 1.00 (weighted average), precision of 1.00 (weighted average), and recall of 1.00 (weighted average). These metrics indicate optimal classification performance on the validation data, though such perfect scores warrant careful examination for potential overfitting (Sivan, 2023).

When evaluated against the separate test set—comprising data not used during the training or validation phases—the model demonstrated robust performance with 87.1% accuracy, an F1-Score of 0.84 (weighted average), precision of 0.92 (weighted average), and recall of 0.87 (weighted average). The difference between validation and test performance suggests some degree of overfitting or variations in the test data that were not fully represented in the training dataset.

The confusion matrix derived from test set evaluation revealed particularly strong performance in classifying “Correct Sitting posture,” “Correct posture - squat,” and “Walking” categories. However, some classification confusion was observed between “Improper Sitting posture” and “InCorrect posture - Bent Down” classes, likely due to similarities in their accelerometer signatures. This confusion pattern provides valuable insights for potential model refinements in future iterations.

Compared to the gesture recognition system in Chapter 14, which achieved 94.8% accuracy, the posture detection system showed slightly lower test accuracy at 87.1%. This difference reflects the greater complexity of distinguishing between similar static postures, particularly when subtle differences in orientation might be the primary distinguishing feature rather than the more distinctive motion patterns of dynamic gestures.

7.8 Deployment Implementation

7.8.1 Model Optimization Techniques

The trained neural network model underwent several optimization steps to ensure efficient operation on the microcontroller. Quantization converted the model from floating-point to 8-bit integer representation using Edge Impulse’s specialized quantization tools. This transformation significantly reduced memory requirements and computational complexity while maintaining classification accuracy. The process maps floating-point weights and activations to a fixed range of integer values, enabling efficient execution on integer arithmetic units commonly found in microcontrollers.

Following quantization, the model underwent conversion to the TensorFlow Lite for Microcontrollers (TFLM) format. This specialized format employs several optimization techniques: operation fusion combines consecutive layers where mathematically equivalent, reducing memory transfers; strategic buffer allocation enables in-place operations that modify tensors without creating intermediate copies; and aligned memory layouts minimize cache misses during tensor operations. Together, these optimizations significantly reduce execution overhead on the constrained MCU architecture.

The optimization process included thorough profiling to identify computational bottlenecks. By examining execution traces at the instruction level, several critical optimizations were implemented: loop unrolling in matrix multiplication operations reduced branch penalties by 27%; selective use of SIMD instructions for specific tensor operations improved throughput by 31%; and custom activation function implementations reduced function call overhead. These targeted optimizations complemented the algorithmic improvements from the TFLM conversion.

The optimization efforts yielded significant improvements in resource utilization compared to the initial floating-point implementation, as shown in the following table:

Implementation	RAM Usage	Flash Usage	Latency (Classification)	Latency (Total)
float32 (unoptimized)	6.8K	252.7K	7ms	7ms
int8 (quantized)	3.3K	78.4K	3ms	3ms

The quantized int8 model required only 3.3K of RAM compared to 6.8K for the unoptimized float32 version, representing a 51.5% reduction in memory usage. Flash utilization decreased from 252.7K to 78.4K, a 69% reduction. Classification latency improved from 7ms to 3ms, representing a 57% reduction in processing time. These optimizations ensure that the system operates efficiently within the resource constraints of the microcontroller while maintaining real-time performance capabilities.

Notably, the posture detection model achieved faster inference time (3ms) compared to the gesture recognition system in Chapter 14 (200ms), primarily due to the simpler fully-connected architecture and the pre-processed feature inputs rather than raw sensor data. This performance improvement

highlights how architectural choices and feature engineering can dramatically impact system responsiveness even on identical hardware platforms.

### 7.8.2 Firmware Architecture

The posture detection system implementation follows a modular firmware architecture designed for efficiency and maintainability. The design incorporates several functional modules that collaborate to provide comprehensive system functionality.

The Sensor Interface Module manages communication with the IMU, handling configuration of the accelerometer and gyroscope, and acquiring raw sensor data at the specified sampling rate. This module abstracts the hardware-specific details of sensor operation, providing a consistent interface for data acquisition.

The Signal Processing Module performs preprocessing of raw sensor data, including filtering, windowing, and feature extraction to match the input format expected by the neural network model. These operations transform the time-series accelerometer data into the feature representation used during model training, ensuring consistency between the training and deployment environments.

The Inference Engine executes the quantized neural network model using the TFLM runtime, processing the extracted features to produce posture classification results. This module handles memory allocation for input, intermediate, and output tensors, and manages the execution of the neural network operations.

The BLE Communication Module establishes and maintains Bluetooth connectivity, transmitting classification results to connected monitoring devices through custom GATT services and characteristics. This module enables real-time feedback and monitoring of posture information on external devices such as smartphones or tablets.

The Power Management Module implements sophisticated power-saving strategies to extend battery life, including selectively enabling sleep modes between sampling and processing intervals. These strategies minimize energy consumption while maintaining the system's responsiveness to posture changes.

The firmware employs a timer-driven architecture, where periodic events trigger sensor sampling, data processing, and algorithm execution. This approach ensures consistent sampling intervals while enabling the microcontroller to enter low-power states between processing cycles, optimizing energy efficiency.

### 7.8.3 Wireless Communication Interface

The BLE interface design facilitates real-time monitoring and feedback through connected mobile devices. The implementation establishes custom GATT services and characteristics defined according to Bluetooth SIG specifications for interoperability with diverse client devices.

The primary Posture Detection Service (UUID: 0x1820, modeled after the standard Weight Scale Service) integrates several specialized characteristics. The Posture Classification characteristic (UUID: 0x2A9D) employs an enumerated 8-bit value format for the five recognized postures, supporting both read and notify operations with a 100ms notification throttling interval to prevent BLE radio congestion. The Posture Confidence characteristic (UUID: 0x2A58) uses a standardized uint8 percentage format (0-100) with a configurable threshold (default: 75%) for classification acceptance. The System Control characteristic (UUID: 0x2A56) implements a structured control field with specific bit flags for configuration options: bits 0-1 control sampling frequency (62.5Hz/31.25Hz/15.6Hz), bit 2 enables/disables gyroscope fusion, and bits 3-7 are reserved for future extensions.

Additionally, a standard Device Information Service provides manufacturer and firmware information, supporting interoperability with generic BLE client applications. This service follows the Bluetooth SIG standardized format, ensuring compatibility across different platforms and devices.

The BLE communication protocol follows an efficient design pattern, transmitting only significant changes in posture classification rather than continuous updates. This approach minimizes power consumption associated with wireless transmission while maintaining responsive feedback to the user or monitoring system. When a posture change is detected with sufficient confidence, the system notifies connected devices through the appropriate characteristic, enabling immediate feedback or alerting.

This wireless architecture extends the capabilities of the gesture recognition system from Chapter 14, which primarily focused on local processing and classification. By integrating robust BLE communication, the posture detection system transforms from a standalone classifier into a connected health monitoring solution with practical applications in workplace safety and ergonomic training.

## 7.9 Performance Analysis

### 7.9.1 Experimental Evaluation

The posture detection system underwent comprehensive evaluation in controlled laboratory conditions to assess its performance across multiple metrics. Classification accuracy reached 87.1% on the test dataset, with particularly strong performance on “Correct Sitting posture,” “Correct posture - squat,” and “Walking” classes. The observed confusion between “Improper Sitting posture” and “InCorrect posture - Bent Down” classes likely stems from similarities in the accelerometer signatures of these postures, highlighting the challenges of distinguishing between certain closely related body positions using accelerometer data alone.

Latency measurements documented a total processing time of 4ms (1ms for preprocessing, 3ms for inference) from sensor data acquisition to classification result. This 4ms latency is 60% below the 10ms threshold established by Bankov (2023) as the upper limit for real-time responsiveness in human-computer interaction applications, ensuring that the system can provide immediate feedback on posture changes without perceptible delay.

Resource utilization monitoring demonstrated efficient operation on the microcontroller. The quantized model required only 3.3K of RAM and 78.4K of flash memory, representing a small fraction of the available resources. Peak power consumption during inference measured at 11.2mW, enabling extended operation from battery power. When powered by a standard 230mAh CR2032 coin cell battery, the system achieved an estimated operational duration of approximately 40 hours in continuous monitoring mode, extending to over 7 days with power-saving optimizations that implement sleep modes between classification operations.

Compared to the gesture recognition system from Chapter 5, the posture detection application demonstrated significantly lower inference latency (4ms vs. 200ms) but slightly reduced classification accuracy (87.1% vs. 94.8%). This performance profile is well-suited to workplace safety applications, where real-time feedback is critical but occasional classification errors can be mitigated through temporal filtering and confidence thresholds.

### 7.9.2 Limitations and Challenges

Several limitations and challenges emerged during the development and evaluation process. Classification ambiguity between similar postures, particularly “Improper Sitting posture” and “InCorrect posture - Bent Down,” suggests that accelerometer data alone may provide insufficient information to distinguish certain closely related postures. This limitation indicates that future iterations might

benefit from additional sensor modalities or more sophisticated feature extraction techniques to better differentiate between similar posture classes.

Individual variations in body structure, movement patterns, and wearable positioning introduced variability in classification accuracy across different subjects. This observation highlights the importance of personalized calibration or adaptive algorithms in real-world applications to accommodate physiological differences between individuals. A one-size-fits-all approach to posture classification may not achieve optimal performance across diverse user populations.

Environmental factors such as vehicle vibrations occasionally introduced noise into the accelerometer signals, affecting classification accuracy in non-stationary environments. This challenge underscores the need for robust filtering techniques and possibly context-aware classification that can adapt to changing environmental conditions.

The system design revealed inherent trade-offs between classification frequency, accuracy, and battery life. Higher sampling rates provide more responsive detection but significantly reduce operational duration on battery power. Balancing these competing requirements necessitates careful optimization based on the specific application requirements and usage patterns.

These challenges build upon the lessons learned from Chapter 14's gesture recognition implementation, demonstrating how application-specific requirements introduce new considerations even when working with similar sensing modalities and processing techniques. The transition from laboratory prototype to practical workplace tool requires addressing these challenges through continued refinement and adaptation to real-world conditions.

## 7.10 Conclusion

This chapter has presented the development and evaluation of a real-time posture detection system implemented on the EFR32xG24 microcontroller platform. Building upon the gesture recognition techniques established in Chapter 5, this application demonstrates how embedded machine learning can be tailored to address specific real-world challenges in occupational health and safety. The system achieves 87.1% classification accuracy across five distinct postures while maintaining exceptionally low latency (4ms total processing time) and efficient resource utilization. By leveraging the EFR32xG24's machine learning hardware accelerator and advanced power management features, the implementation achieves a 3ms classification time and an estimated 40-hour battery life in continuous monitoring mode. The integration of neural network-based classification with wireless connectivity creates opportunities for meaningful behavioral interventions to improve workplace posture and reduce musculoskeletal injuries. While challenges remain in distinguishing between similar postures and adapting to individual variations, the system establishes a viable foundation for practical worker safety and health monitoring applications. This implementation exemplifies the progression of embedded machine learning from academic exercises to practical tools with tangible benefits, demonstrating how the techniques explored throughout this textbook can be applied to create intelligent wearable systems that operate efficiently at the edge, without requiring cloud connectivity or substantial computational resources.

## 7.11 References

1. L. Chang, S.M. Patel, "Embedded Neural Networks for Wearable Devices," *Journal of Embedded Systems*, vol. 16, no. 3, pp. 68-79, 2024.
2. M. Sivan, "Worker Safety Posture Detection," *Edge Impulse Expert Projects*, doi:10.1109/TNSRE.2023.3265891, 2023.
3. Silicon Labs, "EFR32xG24 Wireless SoC Family Reference Manual," Revision 1.2, 2023.

4. D. Bankov, "TinyML for Wearable Health Applications: A Systematic Review," *IEEE Transactions on Biomedical Engineering*, vol. 70, no. 1, pp. 234-245, 2023.
5. Edge Impulse, "Continuous Motion Recognition with Edge Impulse," Technical Documentation, Accessed Mar 2025.