CMP103

# Edge Intelligence: How to Leverage Silicon Labs AI/ML to Improve Efficiency and Performance

Rich Lysaght

*Senior Product Marketing Manager*

SILICON LABS

# TinyML vs AI/ML



AI/ML - Broader machine learning models that require more computational power and are typically run on servers or high-performance edge devices.

- Complex tasks such as deep learning, large-scale data analysis, and sophisticated pattern recognition.
- Greater accuracy, advanced capabilities, and ability to handle complex data processing tasks.

TinyML - A subset of machine learning designed for deployment on microcontrollers and low-power devices.

- Ideal for simple tasks like anomaly detection and basic classification in resource-constrained environments.
- Low power consumption, real-time processing on-device, and cost-effective for embedded applications.

works with | SILICON LABS

# Why Machine Learning on Microcontrollers?

### Reduce Decision Latency

- Make more real time decisions closer to where the data is collected

### Lower data and device security risk

- Keeping data local to devices reduces risk of exposure during transmission

### Bandwidth Constraints

- Bandwidth limited IoT networks cannot transmit large amounts of data required for cloud centric architectures

### Offline Mode Operation

- Allows for nodes to operate autonomously and make decisions even when network is unavailable

### Lower Device and Service Cost

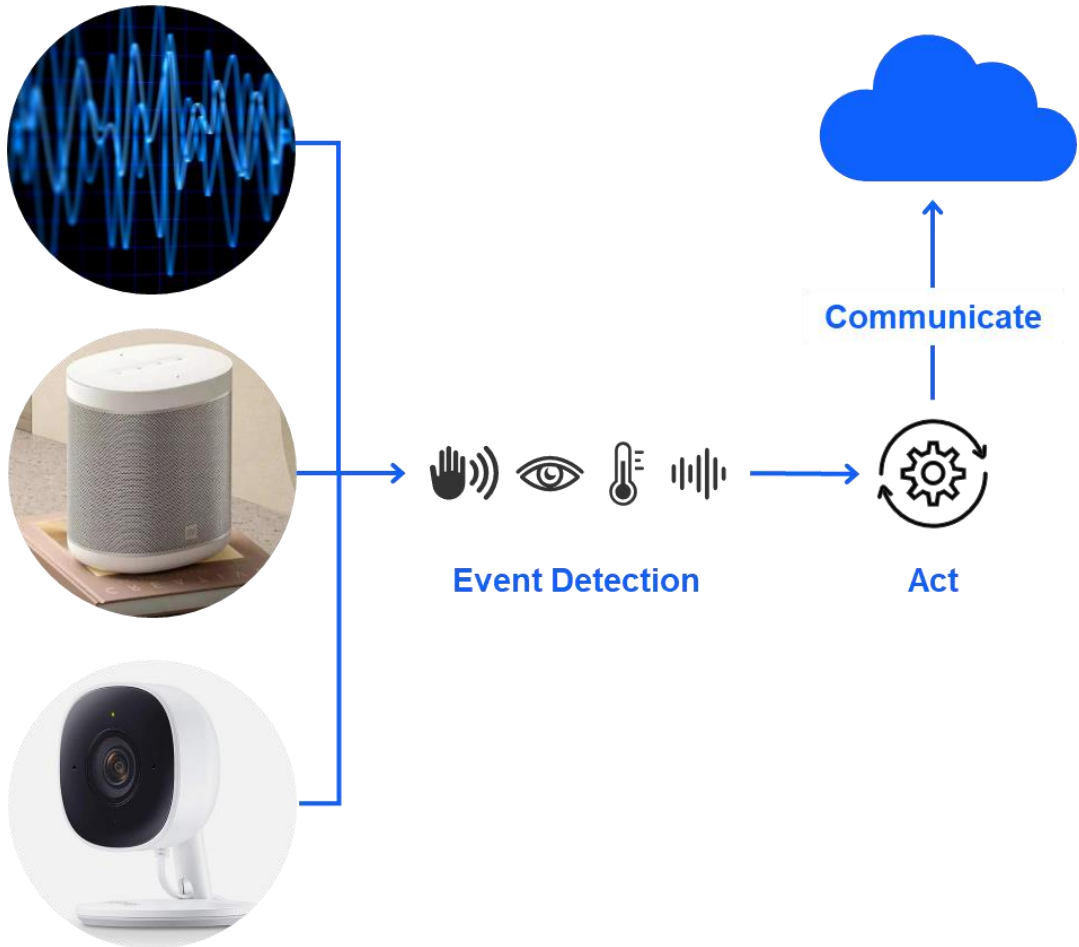- Lowers performance requirements for sensor devices and limits recurring costs

### Improved Low Power Operation

- Reduces number of network transmissions to improve overall battery life

**Data processing is more efficient with Machine Learning at the sensor level**

w/ works with | SILICON LABS

# Reducing Decision Latency

**Communicate**

**Event Detection**          **Act**

## More Real-Time Responses:

- Rapid decision-making allows systems to respond instantly to changes.

## Increased Accuracy:

- Quick decisions based on real-time data help in making timely adjustments, leading to more accurate outcomes
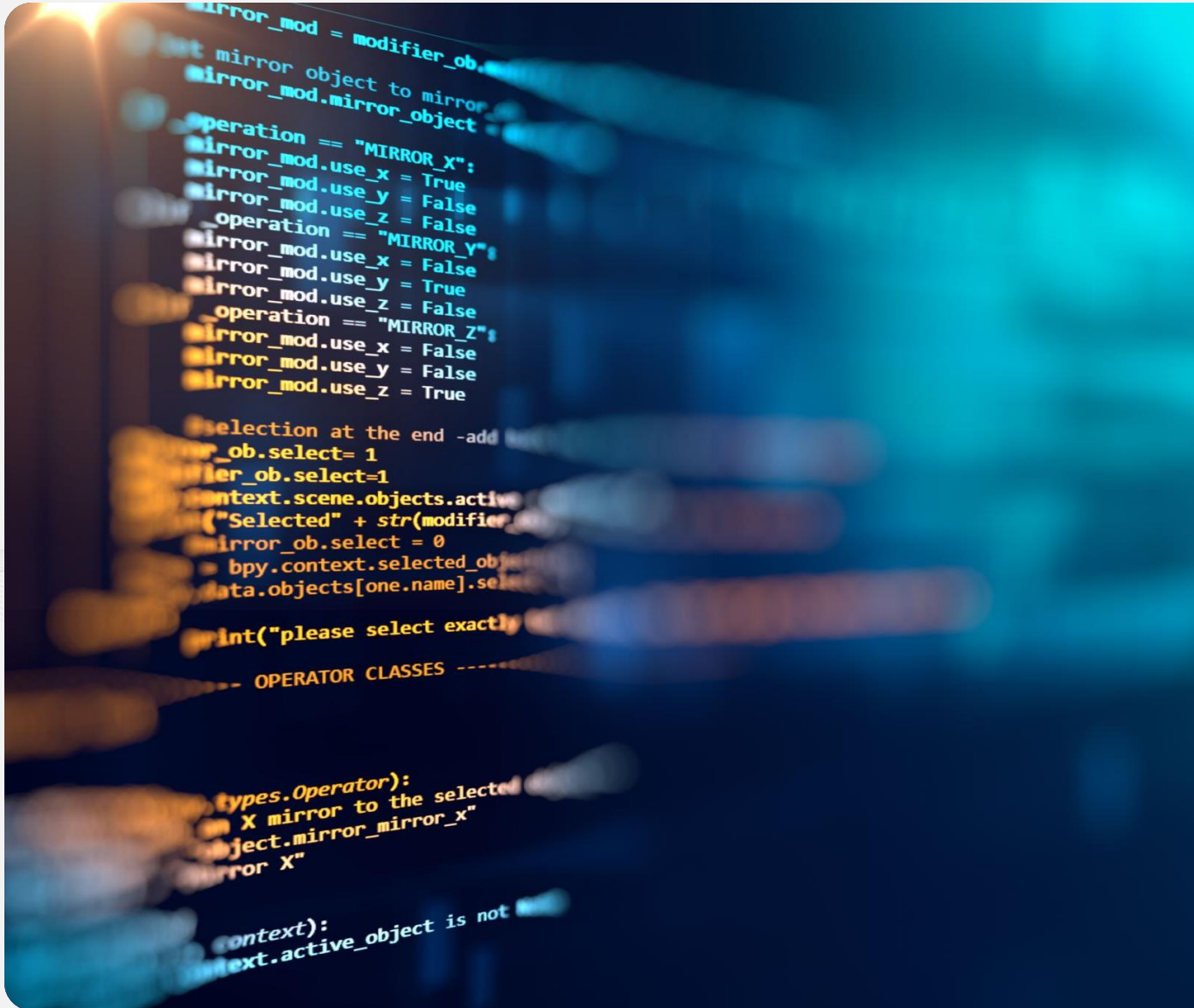
## Adaptive Learning locally:

- Systems can continuously learn and adapt to new data more effectively, improving model accuracy and reliability over time.

## Parallel Processing Capabilities:

- Addition of AIML accelerators enable simultaneous task execution, resulting in higher throughput and faster response times for real-time applications

works with | SILICON LABS

# Lowering Data Security Risks



**Local Data Processing:**
- Reduces the need to transmit sensitive information over networks, minimizing exposure to interception.
- Results are sent to the cloud rather than the data.

**Privacy Compliance:**
- Keeping data on-device aligns with privacy regulations, as less personal data is transmitted or stored in the cloud, thereby minimizing compliance risks.

**Data Encryption:**
- Secure Vault with PSA L3 certification ensures sensitive data to be encrypted before storage or transmission, adding an additional layer of security.

**Firmware Updates:**
- Secure over-the-air (OTA) updates for firmware can help ensure that devices are running the latest security protocols and patches, protecting against vulnerabilities.

**Real-Time Anomaly Detection:**
- TinyML can continuously monitor for unusual patterns locally, enabling immediate responses to potential threats.

works with | SILICON LABS

# Addressing Bandwidth Limitations of IoT Networks



## Reduce Data Volume:

- Only relevant data is transmitted instead of raw data. This is crucial in low-bandwidth settings like LPWAN, which typically supports data rates of 0.1 to 50 kbps.

## Optimize Use of Limited Bandwidth:

- Given LPWAN's constraints, minimizing the data sent not only conserves bandwidth but also enhances the reliability of communication.

## Selective Data Transmission:

- Edge devices send updates only when specific conditions are met (e.g., threshold breaches), which is crucial for LPWAN, where frequent transmissions can quickly exhaust bandwidth.

## Scalability:

- BLE mesh allows devices to communicate directly, reducing reliance on a central hub and enabling a scalable network where many devices can relay information.
- Wi-Fi networks effectively segment traffic, and by processing data at the edge, devices reduce network load and improve efficiency, allowing for better scaling without performance degradation.

works with | SILICON LABS

# Making Decisions Without Network Access

### Self-Sufficient Devices:

- Edge devices can operate independently, making them ideal for applications in remote or challenging locations where connectivity may be limited.

### Reduced Infrastructure Cost:

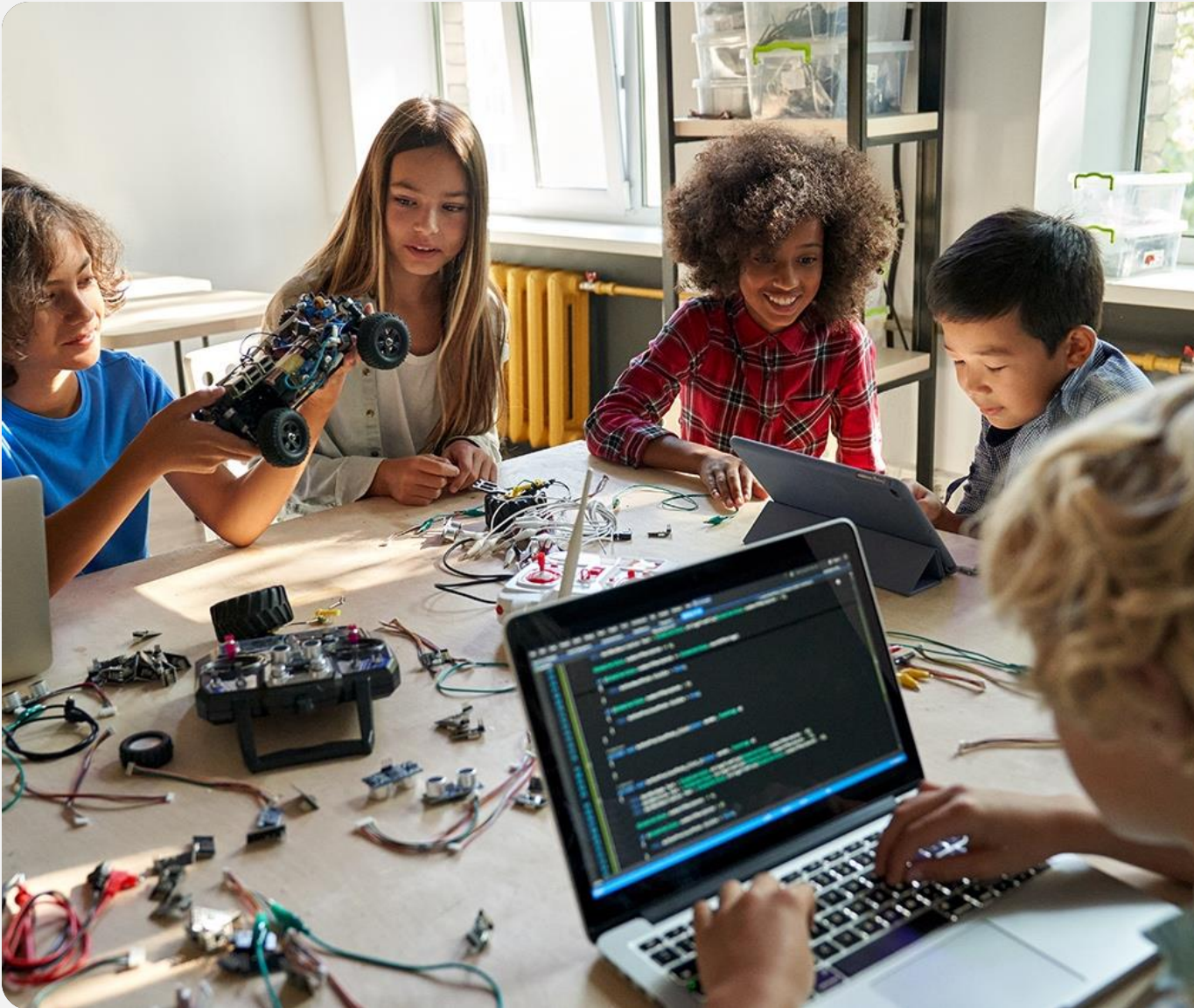- Fewer data transmissions lead to lower costs for cloud services and bandwidth, making IoT solutions more economical.

### Removes Dependence on Cloud Processing:

- Without offline capabilities, devices rely on stable internet connectivity, which can lead to downtime and reduced effectiveness in areas with poor connectivity.

### Lower Risk of Data Loss:

- Unstable connections can result in data being lost or corrupted during transmission, compromising the integrity of the information and potentially leading to poor decisions.

works with | SILICON LABS

# Lowering Costs to Make AI/ML More Accessible

**Affordable Hardware:**

- Advancements in technology have created cheaper, more efficient hardware like MCUs with accelerators for running AI/ML algorithms, reducing initial investments and broadening application possibilities.

**Wider Adoption:**

- More entities can integrate AI/ML into their operations, enhancing innovation and competition across sectors.

**Diverse Applications:**

- Smaller companies and startups can leverage AI/ML for various applications, from automation to data analysis, driving economic growth.

**Reduced Operating Expenses:**

- Streamlined algorithms and optimized hardware lead to lower energy consumption and maintenance costs, decreasing ongoing operational expenses.

**Increased ROI:**

- Lower recurring costs improve the return on investment for AI/ML projects, making them more attractive for businesses.

**Sustainability:**

- Reduced energy and operational costs contribute to more sustainable practices, appealing to environmentally conscious organizations.

works with | SILICON LABS

# Optimized Performance for Low Power Devices



## Compact Design:

- Lower power requirements enable smaller batteries or energy storage solutions, leading to more compact device designs, which is crucial in space-limited applications like wearables and small sensors.

## Efficient Algorithms:

- Edge AI/ML models are optimized for resource-constrained environments through techniques like quantization and pruning, reducing computational complexity and lowering energy consumption.

## Adaptive Sampling:

- Edge devices use adaptive sampling to collect and process data only when needed, minimizing unnecessary computations and data transfers, thereby conserving power.

## Extended Battery Life:

- Reduced power consumption results in longer battery life, decreasing the need for replacements or recharging, which lowers maintenance efforts and costs, enhancing user-friendliness.

## Network Efficiency:

- Low-power devices can operate effectively in diverse environments and network conditions, enabling scalable IoT solutions without overwhelming infrastructure.
- Enhanced Scalability:
- Efficient Resource Utilization, local processing reduces the need for centralized resources, allowing for more IoT devices to be deployed without overwhelming network infrastructure.

works with | SILICON LABS

# Adding Local Acceleration for AI/ML Inferencing

works with | SILICON LABS

# MVP Math library

Accelerate and do more efficient linear algebra operations with internal MVP subsystem

## Math **APIs (alternative to CMSIS_DSP)** available in GSDK

**VECTOR OPERATIONS**

- Vector Add
- Vector Absolute Value
- Vector Clip
- Vector Dot Product
- Vector Multiply
- Vector Negate
- Vector Offset
- Vector Scale
- Vector Sub
- Complex Vector Conjugate
- Complex Vector Dot Product
- Complex Vector Magnitude
- Complex Vector Magnitude Squared
- Complex Vector Multiply
- Complex Vector Multiply Real
- Vector Copy
- Vector Fill

**MATRIX OPERATIONS**

- Matrix Initialize
- Matrix Multiply
- Matrix Scale
- Matrix Sub
- Matrix Transpose
- Matrix Multiply Vector
- Matrix Add
- Complex Matrix Multiply
- Complex Matrix Transpose

✓ **Faster and more efficient** **execution of many algorithms with large data for example filtering algorithms**
✓ **Saving CPU cycles, saving power, resulting longer battery life**
✓ **Option to win sockets against faster CPUs**

**CortexM only**

| Matrix dims. | | CMSIS f32 cpu-cycles | CMSIS f16 cpu-cycles | MVP cpu-cycles | instr | stalls |
|---|---|---|---|---|---|---|
| 2x2 | 2x2 | 226 | 304 | 403 | 8 | 0 |
| 4x2 | 2x4 | 602 | 913 | 424 | 32 | 0 |
| 6x2 | 2x6 | 1210 | 1921 | 464 | 72 | 0 |
| 8x2 | 2x8 | 2050 | 3321 | 516 | 128 | 0 |
| 10x2 | 2x10 | 3122 | 5113 | 592 | 200 | 0 |
| 12x2 | 2x12 | 4426 | 7297 | 676 | 288 | 0 |
| 14x2 | 2x14 | 5962 | 9873 | 784 | 392 | 0 |
| 16x2 | 2x16 | 7730 | 12841 | 904 | 512 | 0 |
| 18x2 | 2x18 | 9730 | 16201 | 1036 | 648 | 0 |
| 20x2 | 2x20 | 11962 | 19953 | 1192 | 800 | 0 |
| 20x4 | 4x20 | 17962 | 27956 | 1593 | 1200 | 1 |
| 20x6 | 6x20 | 23742 | 39956 | 2193 | 1600 | 201 |
| 20x8 | 8x20 | 27562 | 47556 | 2793 | 2000 | 400 |
| 20x10 | 10x20 | 33162 | 59556 | 3393 | 2400 | 601 |
| 20x12 | 12x20 | 37162 | 67156 | 3993 | 2800 | 801 |
| 20x14 | 14x20 | 42762 | 79156 | 4593 | 3200 | 1000 |
| 20x16 | 16x20 | 46762 | 86756 | 5193 | 3600 | 1201 |
| 20x18 | 18x20 | 52362 | 98756 | 5793 | 4000 | 1401 |
| 20x20 | 20x20 | 56362 | 106356 | 6393 | 4400 | 1600 |

**~ 9x less cycles**

works with | SILICON LABS

# Benefits of the MVP ML Hardware Accelerator

Dedicated ML computing subsystem next to the CPU: Matrix Vector Processor (MVP)

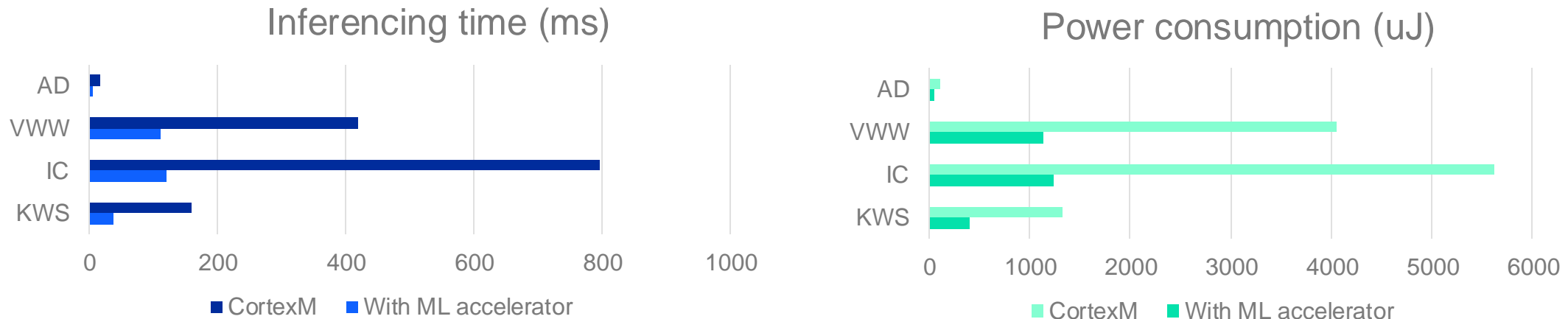Optimized MVP to accelerate ML inferencing with a lot of processing power **offloading the CPU**

**Up to 8x faster** inferencing over Cortex-M (see below perf. benchmark)

Up to **6x lower power** for inferencing (see below perf. benchmark)

Dedicated OPNs for MVP accelerated parts → EFR32MG24B[2]… or [3]

### Performance data with ML hardware accelerator vs. pure SW on CortexM*



Inferencing time (ms)

Power consumption (uJ)

*Standardized performance benchmark validated by independent benchmarking body **MLCommons.org**. Published in MLPerf Tiny v1.0. Results are for inferencing only (not for the complete application). You can refer to MLCommons as validated results-

works with | SILICON LABS

# MVP – Matrix Vector Processor Demo

# MVP – Matrix Vector Processor (AI/ML Accelerator)

AI/ML Hardware Accelerator Key Features

## Matrix Processor Accelerates ML Inferences

- Multi-dimensional array operations
- Handles real and complex data
- Offloads MCU

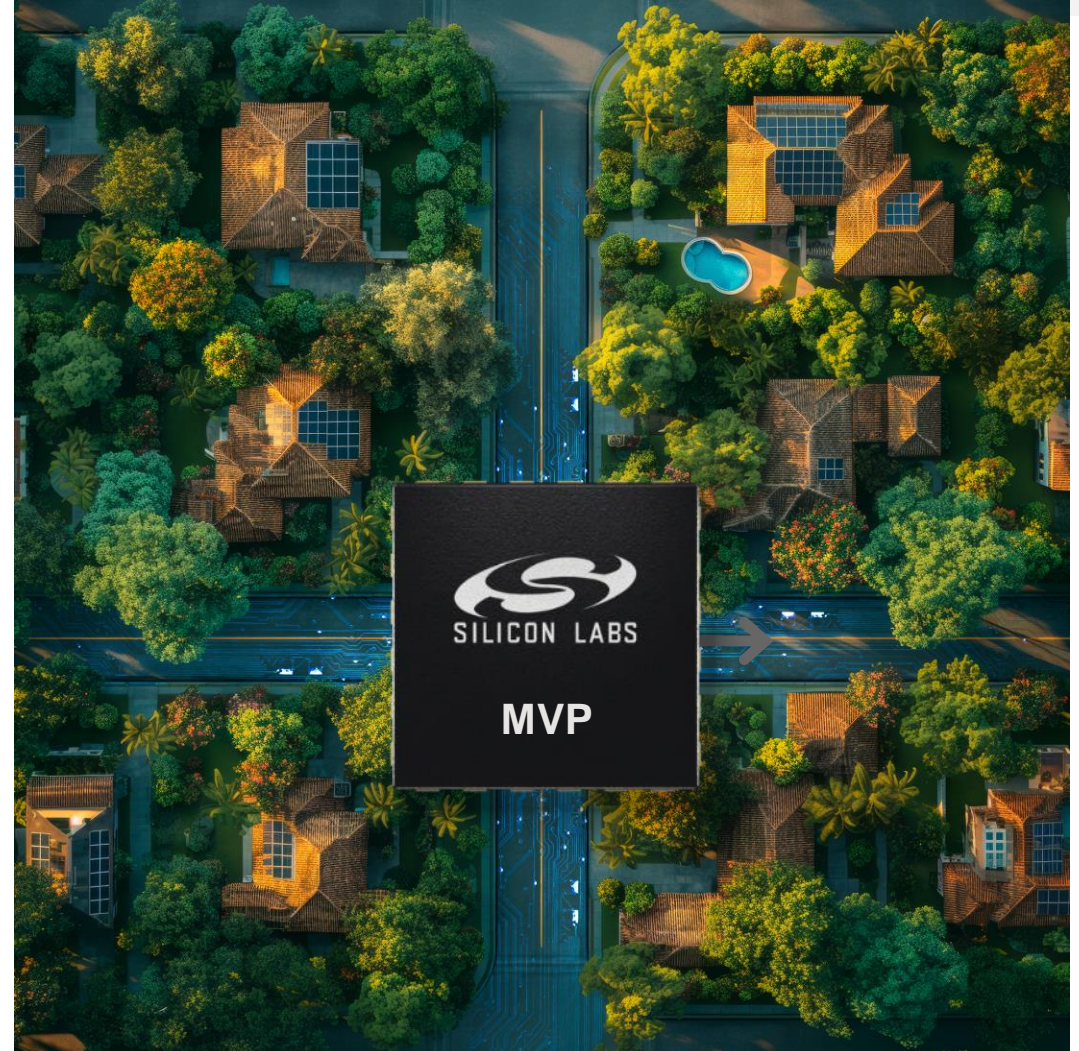## Up to 8x faster inference over Cortex-M

- Lower latency

## Up to 6x lower power for inferencing

- Longer battery life

## MVP Math Library

- Can be used for non-ML applications

- AI/ML Hardware Accelerator enables efficient Edge ML inferencing

works with | SILICON LABS

# MVP Math Library API



**https://docs.silabs.com/d/platform-compute-math/4.3/**

## MVP MATH LIBRARY API

**Vector Functions**
- sl_math_mvp_vector_clip_f16
- sl_math_mvp_complex_vector_dot_product_f16
- sl_math_mvp_vector_copy_f16
- sl_math_mvp_vector_sub_f6
- sl_math_mvp_vector_mult_f16
- sl_math_mvp_vector_abs_f16
- sl_math_mvp_vector_scale_f16
- sl_math_mvp_vector_add_f16
- sl_math_mvp_vector_add_i8
- sl_math_mvp_complex_vector_mult_real_f16
- sl_math_mvp_complex_vector_mult_f16
- sl_math_mvp_vector_negate_f16
- sl_math_mvp_complex_vector_conjugate_f16
- sl_math_mvp_vector_fill_f16
- sl_math_mvp_complex_magnitude_squared_f16
- sl_math_mvp_vector_dot_product_f16
- sl_math_mvp_clamp_i8
- sl_math_mvp_vector_offset_f16

## MVP MATH LIBRARY API

**Matrix Functions**
- sl_math_mvp_matrix_mult_f16
- sl_math_mvp_matrix_scale_f16
- sl_math_mvp_matrix_transpose_f16
- sl_math_mvp_complex_matrix_transpose_f16
- sl_math_mvp_matrix_add_f16
- sl_math_mvp_matrix_sub_f16
- sl_math_mvp_matrix_init_f16
- sl_math_mvp_matrix_vector_mult_f16
- sl_math_mvp_complex_matrix_mult_f16

**Utility Functions**
- sl_math_mvp_clear_errors
- sl_math_mvp_get_error

works with | SILICON LABS

# Demo

**MVP Math Library**

**xG24-RB4186**
- EFR32MG24B210F1536IM48
- +10dBm
- 1536kB Flash
- 256kB RAM
- MVP Equipped

**WSTK**
- SLWMB4002A

EFR32xG24 2.4 GHz 10 dBm Radio Board (BRD4186C Rev A00)

OVERVIEW    **EXAMPLE PROJECTS & DEMOS**    DOCUMENTATION    COMPATIBLE TOOLS

Run a pre-compiled demo or create a new project based on a software example.

Filter on keywords

mvp

1 resources found

Demos

Example Projects

Solution Examples

**Platform - Demonstrate the MVP math library**
This example project shows how to use the MVP math library.

View Project Documentation

CREATE

works with | SILICON LABS

# Demo

**MVP Math Library Demo Multiply**

```
288    printf("Fill matrix A with values:\n");
289    input_a[0] = 1.0;
290    input_a[1] = 2.0;
291    input_a[2] = 3.0;
292    input_a[3] = -4.0;
293    input_a[4] = 5.0;
294    input_a[5] = 6.0;
295    input_a[6] = 7.0;
296    input_a[7] = -8.0;
297    input_a[8] = 9.0;
298    input_a[9] = 10.0;
299    input_a[10] = 11.0;
300    input_a[11] = -12.0;
301    sl_math_matrix_init_f16(&matrix_a, 3, 4, input_a);
302    print_matrix(&matrix_a);
303
304    printf("Transpose matrix A into matrix B:\n");
305    sl_math_matrix_init_f16(&matrix_b, 4, 3, input_b);
306    sl_math_mvp_matrix_transpose_f16(&matrix_a, &matrix_b);
307    print_matrix(&matrix_b);
308
309    printf("Multiply matrix A with matrix B:\n");
310    sl_math_matrix_init_f16(&matrix_z, 3, 3, output);
311    sl_math_mvp_matrix_mult_f16(&matrix_a, &matrix_b, &matrix_z);
312    print_matrix(&matrix_z);
```

works with | SILICON LABS

# Demo

**MVP Math Library Demo Output**

# Demo

**Matrix Multiply Example**



$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 10 & 11 \\ 20 & 21 \\ 30 & 31 \end{bmatrix}$$

$$= \begin{bmatrix} 1\times10 + 2\times20 + 3\times30 & 1\times11 + 2\times21 + 3\times31 \\ 4\times10 + 5\times20 + 6\times30 & 4\times11 + 5\times21 + 6\times31 \end{bmatrix}$$

$$= \begin{bmatrix} 10+40+90 & 11+42+93 \\ 40+100+180 & 44+105+186 \end{bmatrix} = \begin{bmatrix} 140 & 146 \\ 320 & 335 \end{bmatrix}$$

works with | SILICON LABS

# Demo

**Matrix Multiply Example - Initialization**

```
55  //Variables for RGL matrix math
56  static float16_t    rgl_input_a[4][4];
57  static float16_t    rgl_input_b[4][4];
58  static float16_t    rgl_output[4][4];
```

```
77  void rgl_fill_matrix_a (void)
78  {
79      //ROW 0 Fill
80      rgl_input_a[0][0] = 1.00;
81      rgl_input_a[0][1] = 2.00;
82      rgl_input_a[0][2] = 3.00;
83      rgl_input_a[0][3] = -4.00;
84
85      //ROW 1 Fill
86      rgl_input_a[1][0] = 5.00;
87      rgl_input_a[1][1] = 6.00;
88      rgl_input_a[1][2] = 7.00;
89      rgl_input_a[1][3] = -8.00;
90
91      //ROW 2 Fill
92      rgl_input_a[2][0] = 9.00;
93      rgl_input_a[2][1] = 10.00;
94      rgl_input_a[2][2] = 11.00;
95      rgl_input_a[2][3] = -12.00;
96  }
```

Init

```
101 void rgl_fill_matrix_b (void)
102 {
103     //ROW 0 Fill
104     rgl_input_b[0][0] = 1.00;
105     rgl_input_b[0][1] = 5.00;
106     rgl_input_b[0][2] = 9.00;
107
108     //ROW 1 Fill
109     rgl_input_b[1][0] = 2.00;
110     rgl_input_b[1][1] = 6.00;
111     rgl_input_b[1][2] = 10.00;
112
113     //ROW 2 Fill
114     rgl_input_b[2][0] = 3.00;
115     rgl_input_b[2][1] = 7.00;
116     rgl_input_b[2][2] = 11.00;
117
118     //ROW 3 Fill
119     rgl_input_b[3][0] = -4.00;
120     rgl_input_b[3][1] = -8.00;
121     rgl_input_b[3][2] = -12.00;
122 }
```

works with | SILICON LABS

# Demo

**Matrix Multiply Example – Multiply & Print**

```
145  void rgl_multiply_matrix (int num_rows_a, int num_cols_a, int num_rows_b, int num_cols_b)
146  {
147
148      for (int i = 0; i < num_rows_a; i++)
149      {
150          for(int j = 0; j < num_cols_b; j++)
151          {
152              rgl_output[i][j] = 0;
153              for(int k = 0; k < num_rows_b; k++)
154              {
155                  rgl_output[i][j] += rgl_input_a[i][k] * rgl_input_b[k][j];
156              }
157          }
158      }
159  }
```

```
164  void rgl_print_output_matrix (int num_rows, int num_cols)
165  {
166    float16_t my_data;
167
168      printf("Output of Matrix Multiply.\n");
169      for (int r = 0; r < num_rows; r++)
170      {
171          for(int c = 0; c < num_cols; c++)
172          {
173              my_data = rgl_output[r][c];
174              printf("(%6.2f),  ", my_data);
175          }
176          printf("\n");
177      }
178  }
```
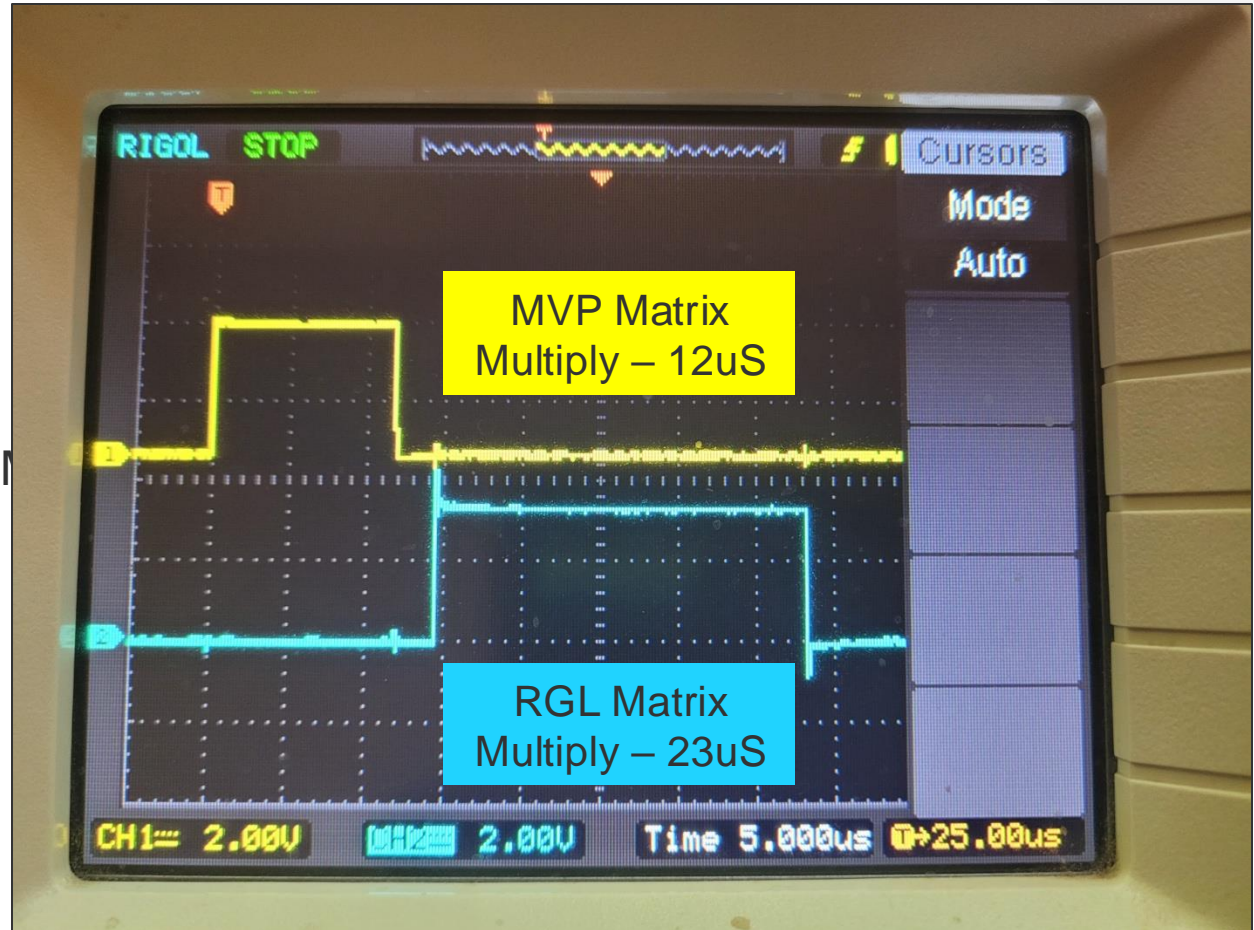
works with | SILICON LABS

# Demo

**Matrix Multiply Example – Compare Functions**

```c
462  void app_process_action(void)
463  {
464      //printf("\n");
465      //printf("Fill matrix A with values:\n");
466      input_a[0] = 1.0;
467      input_a[1] = 2.0;
468      input_a[2] = 3.0;
469      input_a[3] = -4.0;
470      input_a[4] = 5.0;
471      input_a[5] = 6.0;
472      input_a[6] = 7.0;
473      input_a[7] = -8.0;
474      input_a[8] = 9.0;
475      input_a[9] = 10.0;
476      input_a[10] = 11.0;
477      input_a[11] = -12.0;
478      sl_math_matrix_init_f16(&matrix_a, 3, 4, input_a);
479      //print_matrix(&matrix_a);
480
481      //printf("Transpose matrix A into matrix B:\n");
482      sl_math_matrix_init_f16(&matrix_b, 4, 3, input_b);
483      sl_math_mvp_matrix_transpose_f16(&matrix_a, &matrix_b);
484      //print_matrix(&matrix_b);
485
486      //printf("Multiply matrix A with matrix B:\n");
487      sl_math_matrix_init_f16(&matrix_z, 3, 3, output);
488      GPIO_PinOutSet (gpioPortB , LED0);
489      sl_math_mvp_matrix_mult_f16(&matrix_a, &matrix_b, &matrix_z);
490      GPIO_PinOutClear (gpioPortB , LED0);
491      //print_matrix(&matrix_z);
492
493      //Test RGL matrix multiply
494      rgl_fill_matrix_a();
495      rgl_fill_matrix_b();
496      GPIO_PinOutSet (gpioPortB , LED1);
497      rgl_multiply_matrix(3, 4, 4, 3);
498      GPIO_PinOutClear (gpioPortB , LED1);
499      rgl_print_output_matrix (3, 3);
500  }
```

```
Output of Matrix Multiply.
( 30.00),   ( 70.00),   (110.00),
( 70.00),   (174.00),   (278.00),
(110.00),   (278.00),   (446.00),
```

works with  |  SILICON LABS

# Demo

**Matrix Multiply – RGL vs MVP Performance**

Thank You

works with | SILICON LABS