

# Low-Power Neural Network Accelerators: Advancements in Custom Floating-Point Techniques

Yarib Nevarez

Universität Bremen

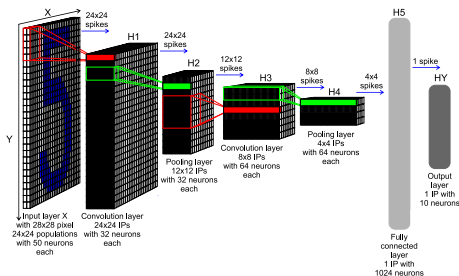
May 22, 2024



## 1 Hybrid 8-bit Floating-Point and 4-bit Logarithmic Computation

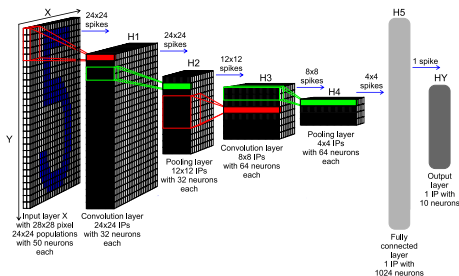
## 2 Hybrid 6-bit Floating-Point Computation

# Spike-by-Spike Neural Network

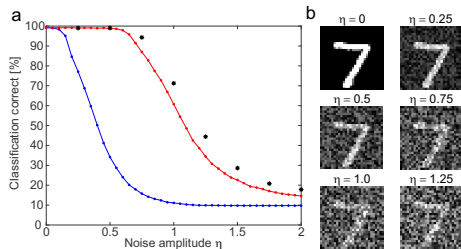


**Figure:** Spike-by-Spike (SbS) neural network architecture for handwritten digit classification task.

# Spike-by-Spike Neural Network

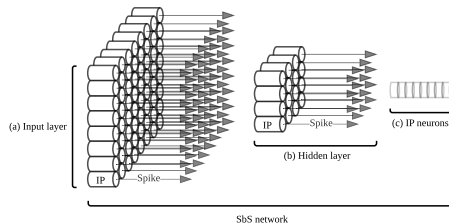


**Figure:** Spike-by-Spike (SbS) neural network architecture for handwritten digit classification task.



**Figure:** Performance classification of SbS NN versus equivalent CNN.

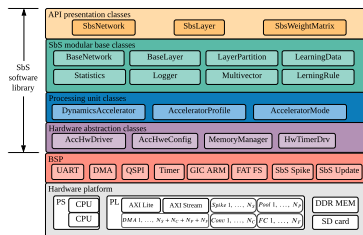
# Spike-by-Spike Layer Update



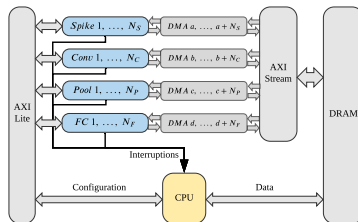
**Figure:** SbS inference population (IP) as independent computational entities.

$$h_{\mu}^{new}(i) = \frac{1}{1 + \epsilon} \left( h_{\mu}(i) + \epsilon \frac{h_{\mu}(i)W(s_t|i)}{\sum_j h_{\mu}(j)W(s_t|j)} \right)$$

# HW/SW Co-Design Framework



**Figure:** System-level overview of the embedded software architecture.



**Figure:** System-level hardware architecture with scalable number of heterogeneous processing units (PU).

# Processing Unit

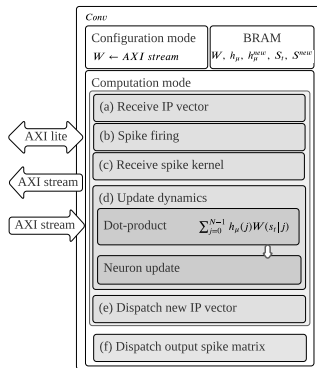


Figure: Conv processing unit.

# Processing Unit

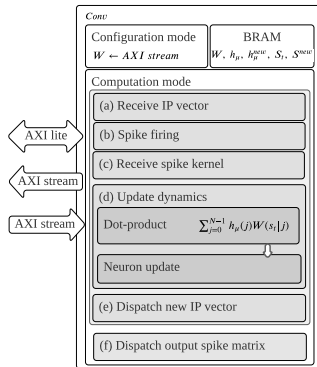


Figure: Conv processing unit.

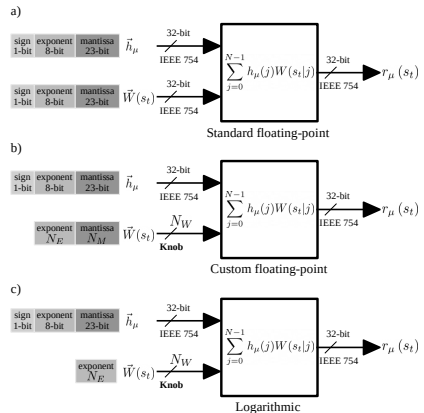


Figure: Dot-product hardware module.



# Hybrid Dot-Product Approximation

$$r_{\mu}(s_t) = \sum_{j=0}^{N-1} h_{\mu}(j) W(s_t|j) \quad (1)$$

$$E_{\min} = \log_2(\min_{\forall i}(W(i))) \quad (2)$$

$$N_E = \lceil \log_2(|E_{\min}|) \rceil \quad (3)$$

$$N_W = N_E + N_M \quad (4)$$

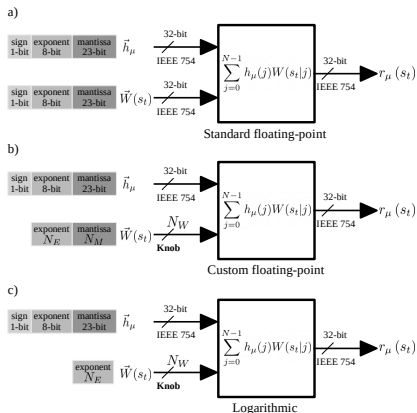


Figure: Dot-product hardware module.

# Dot-Product with Standard Floating-Point (IEEE 754)

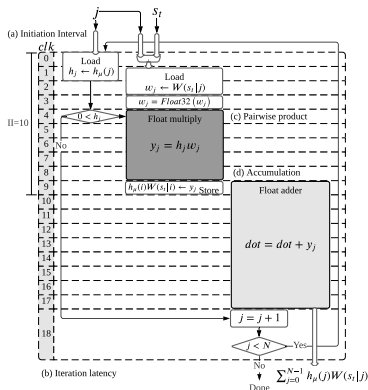


Figure: Dot-product hardware module with standard floating-point computation.

$$L_{f32} = 10N + 9$$

# Dot-Product with Hybrid Custom Floating-Point Approximation

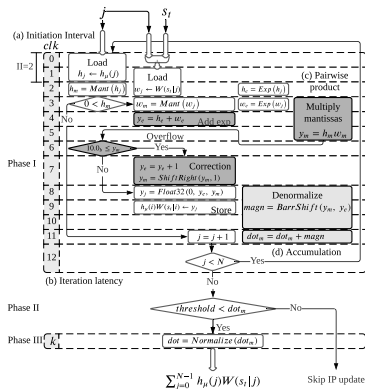


Figure: Dot-product hardware module with hybrid custom floating-point approximation.

$$L_{custom} = 2N + 11$$

# Dot-product with Hybrid Logarithmic Approximation

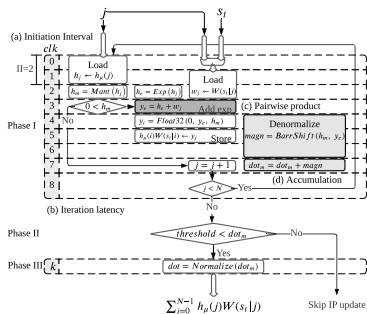
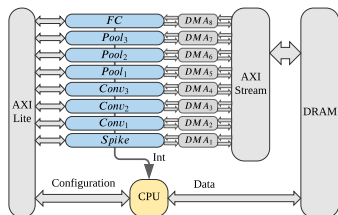


Figure: Dot-product hardware module with hybrid logarithmic approximation.

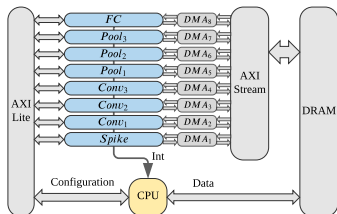
$$L_{custom} = 2N + 7$$

# Acceleration with Standard Floating-Point

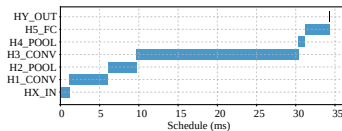


**Figure:** System overview of the top-level architecture with 8 processing units.

# Acceleration with Standard Floating-Point

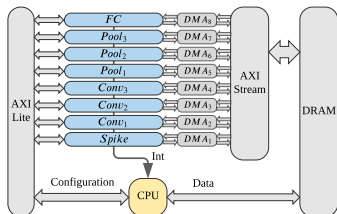


**Figure:** System overview of the top-level architecture with 8 processing units.

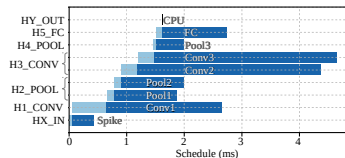


**Figure:** Computation on embedded CPU.

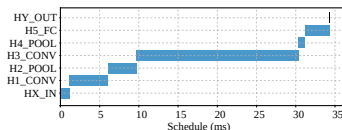
# Acceleration with Standard Floating-Point



**Figure:** System overview of the top-level architecture with 8 processing units.

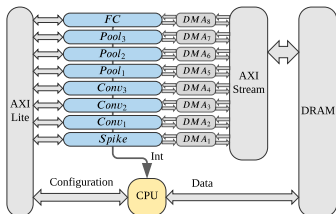


**Figure:** Performance of processing units with standard floating-point.

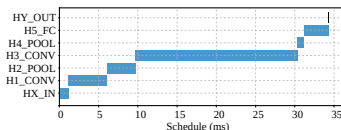


**Figure:** Computation on embedded CPU.

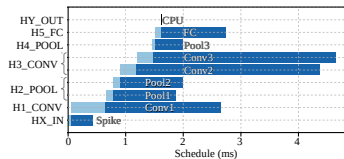
# Acceleration with Standard Floating-Point



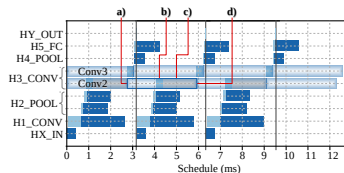
**Figure:** System overview of the top-level architecture with 8 processing units.



**Figure:** Computation on embedded CPU.



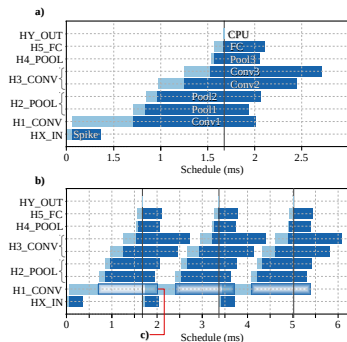
**Figure:** Performance of processing units with standard floating-point.



**Figure:** Performance bottleneck of cyclic computation on processing units with standard floating-point.

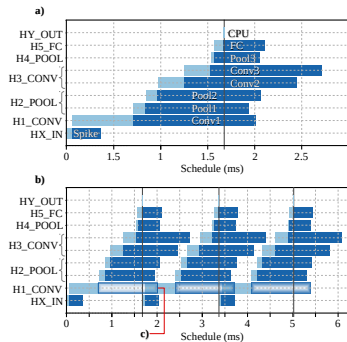


# Acceleration with Custom Floating-Point

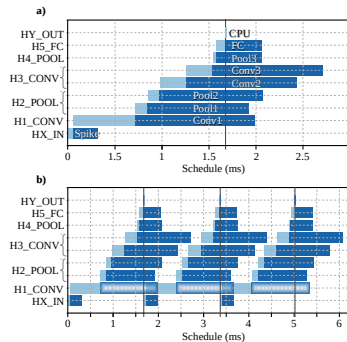


**Figure:** Performance on processing units with hybrid 8-bit floating-point.

# Acceleration with Custom Floating-Point

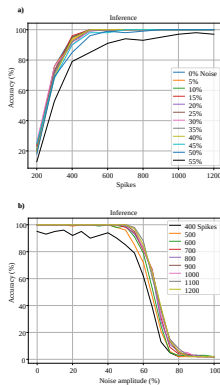


**Figure:** Performance on processing units with hybrid 8-bit floating-point.



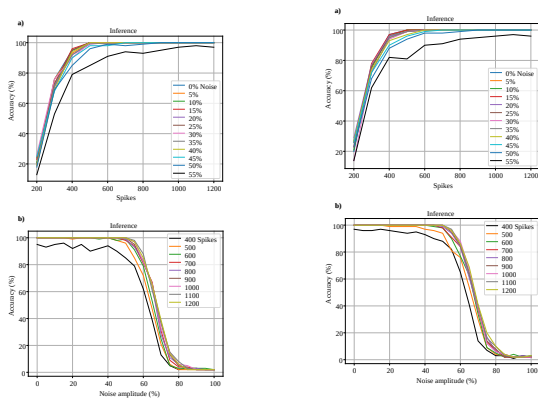
**Figure:** Performance of processing units with hybrid 4-bit logarithmic approximation.

# Noise tolerance



**Figure:** Noise tolerance with standard floating-point.

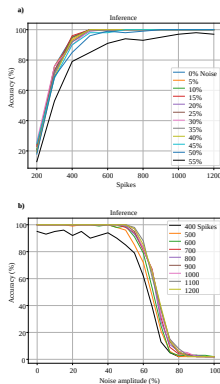
# Noise tolerance



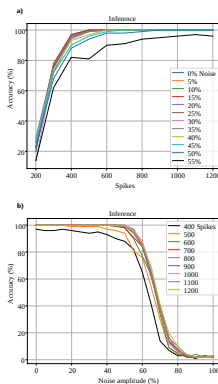
**Figure:** Noise tolerance with standard floating-point.

**Figure:** Noise tolerance with hybrid 8-bit floating-point approximation.

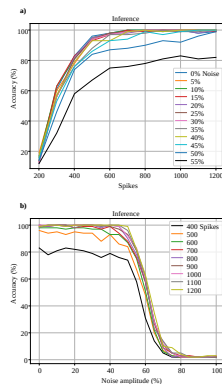
# Noise tolerance



**Figure:** Noise tolerance with standard floating-point.



**Figure:** Noise tolerance with hybrid 8-bit floating-point approximation.



**Figure:** Noise tolerance with hybrid 4-bit logarithmic approximation.

# Accelerator Implementations

Table: Accelerator implementations.

Platform implementation	Power (W)	Clk (MHz)	Latency (ms)	Acceleration	Accuracy (%)
Standard floating-point	2.420	200	3.18	10.7x	98.98
Hybrid floating-point 8-bit	2.369	200	1.67	20.5x	98.97
Hybrid Logarithmic 4-bit	2.324	200	1.67	20.5x	98.84

## 1 Hybrid 8-bit Floating-Point and 4-bit Logarithmic Computation

## 2 Hybrid 6-bit Floating-Point Computation

# Convolution Operation

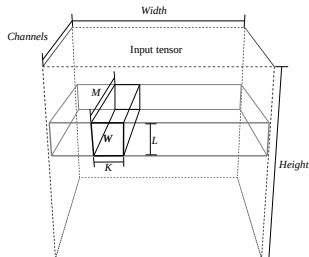
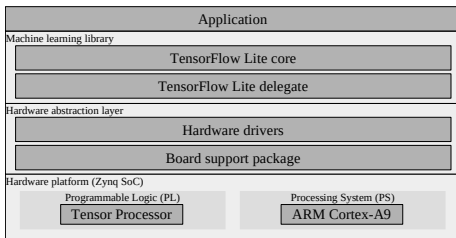


Figure: Two dimensional convolution operation.

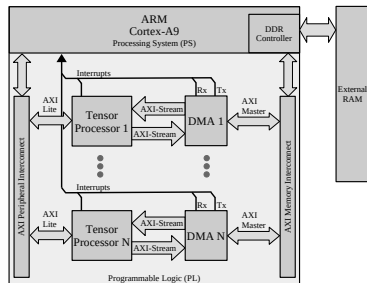
$$\text{Conv2D}(W, b, h)_{i,j,o} = \sum_{k,l,m}^{K,L,M} h_{(i+k,j+l,m)} W_{(o,k,l,m)} + b_o$$



# HW/SW Co-Design Framework



**Figure:** High level embedded software architecture.



**Figure:** Base embedded system architecture.

# Tensor Processor

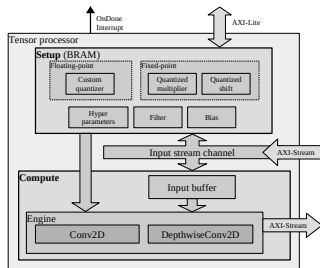


Figure: High level architecture.

# Tensor Processor

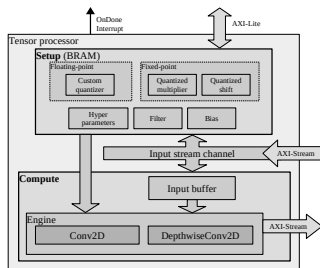


Figure: High level architecture.

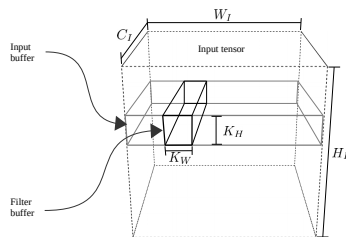


Figure: On-chip memory buffers.

# On-Chip Memory

$$TP_M = TP_B + V_M \quad (5)$$

$$TP_B = Input_M + Filter_M + Bias_M \quad (6)$$

$$Input_M = K_H W_I C_I BitSize_I \quad (7)$$

$$Filter_M = C_I K_W K_H C_O BitSize_F \quad (8)$$

$$Bias_M = C_O BitSize_B \quad (9)$$

$$C_O = \frac{TP_M - V_M - K_H W_I C_I BitSize_I}{C_I K_W K_H BitSize_F + BitSize_B} \quad (10)$$

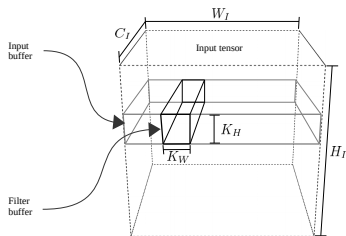


Figure: On-chip memory buffers.

# Dot-Product with Hybrid Floating-Point 6-bit

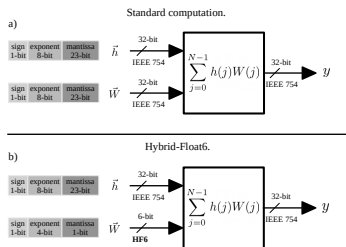


Figure: High level architecture.

# Dot-Product with Hybrid Floating-Point 6-bit

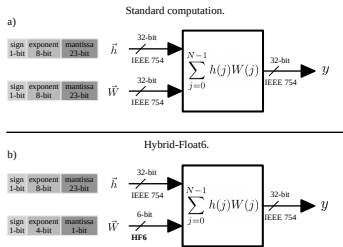


Figure: High level architecture.

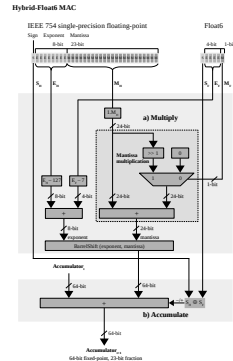


Figure: On-chip memory buffers.

$$L_{hf} = N + 7$$

# Custom Floating-Point Quantization

---

**Algorithm 1:** Custom floating-point quantization.

---

**Input:**  $MODEL$  as the CNN.  
**Input:**  $E_{size}$  as the target exponent bit size.  
**Input:**  $M_{size}$  as the target mantissa bits size.  
**Input:**  $STDM_{size}$  as the IEEE 754 mantissa bit size.  
**Output:**  $MODEL$  as the quantized CNN.

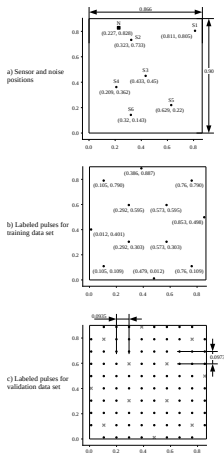
```

1 foreach  $layer$  in  $MODEL$  do
2   if  $layer$  is  $Conv2D$  or  $SeparableConv2D$  then
3      $filter, bias \leftarrow GetWeights(layer)$ 
4     foreach  $x$  in  $filter$  and  $bias$  do
5        $sign \leftarrow GetSign(x)$ 
6        $exp \leftarrow GetExponent(x)$ 
7        $fullexp \leftarrow 2^{E_{size}-1} - 1$  // Get full range value
8        $cman \leftarrow GetCustomMantissa(x, M_{size})$ 
9        $leftman \leftarrow GetLeftoverMantissa(x, M_{size})$ 
10      if  $exp < -fullexp$  then
11         $x \leftarrow 0$ 
12      else
13        if  $exp > fullexp$  then
14           $x \leftarrow (-1)^{sign} \cdot 2^{fullexp} \cdot (1 + (1 - 2^{-M_{size}}))$ 
15        else
16          if  $2^{STDM_{size}-M_{size}-1} - 1 < leftman$  then
17             $cman \leftarrow cman + 1$  // Above halfway
18            if  $2^{M_{size}-1} < cman$  then
19               $cman \leftarrow 0$  // Correct mantissa overflow
20               $exp \leftarrow exp + 1$ 
21           $x \leftarrow (-1)^{sign} \cdot 2^{exp} \cdot (1 + cman \cdot 2^{-M_{size}})$ 
22       $SetWeights(layer, filter, bias)$ 

```

---

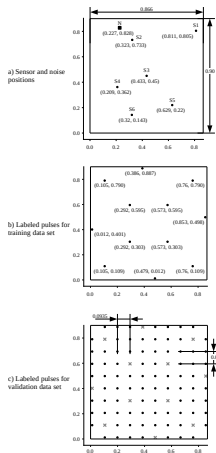
# Acceleration in Sensor Analytics (TinyML)



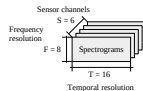
**Figure:** Structural health monitoring, all lengths are in metters (m).



# Acceleration in Sensor Analytics (TinyML)



a) Input tensor



CNN-regression model

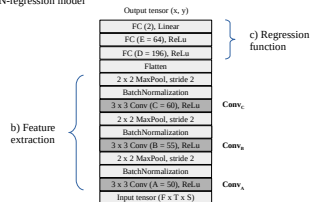


Figure: CNN-regression model for sensor analytics.

Figure: Structural health monitoring, all lengths are in meters (m).

# Training

---

**Algorithm 2:** Training with iterative early stop cycle.

---

**Input:**  $MODEL$  as the input model  
**Input:**  $D_{train}$  as the training data set  
**Input:**  $D_{val}$  as the validation data set  
**Input:**  $N_I$  as the stop patience for iterative training cycle  
**Input:**  $N_E$  as the early stop patience (epochs) for training  
**Input:**  $B_{size}$  as the mini-batch size  
**Output:**  $MODEL$  as the full-precision output model

```

// Initial training and evaluation
1  $Train(MODEL, D_{train}, D_{val}, N_E, B_{size})$ 
2  $mse_i \leftarrow Evaluate(MODEL, D_{val})$ 
3  $n_I \leftarrow 0$ 
4 while  $n_I < N_I$  do
5   // Iterative early stop cycle
6    $Train(MODEL, D_{train}, D_{val}, N_E, B_{size})$ 
7    $mse_v \leftarrow Evaluate(MODEL, D_{val})$ 
8   if  $mse_v < mse_i$  then
9      $Update(MODEL)$ 
10     $mse_i \leftarrow mse_v$ 
11  end
12  else
13     $MODEL \leftarrow LoadPreviousWeights()$ 
14     $n_I \leftarrow n_I + 1$ 
15  end
16 end
  
```

---

# Training

---

**Algorithm 2:** Training with iterative early stop cycle.

---

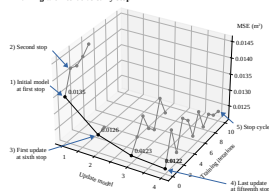
**Input:**  $MODEL$  as the input model  
**Input:**  $D_{train}$  as the training data set  
**Input:**  $D_{val}$  as the validation data set  
**Input:**  $N_I$  as the stop patience for iterative training cycle  
**Input:**  $N_E$  as the early stop patience (epochs) for training  
**Input:**  $B_{size}$  as the mini-batch size  
**Output:**  $MODEL$  as the full-precision output model

```

// Initial training and evaluation
1 Train(MODEL, D_train, D_val, N_E, B_size)
2 mse_i ← Evaluate(MODEL, D_val)
3 n_I ← 0
4 while n_I < N_I do
    // Iterative early stop cycle
    5 Train(MODEL, D_train, D_val, N_E, B_size)
    6 mse_v ← Evaluate(MODEL, D_val)
    7 if mse_v < mse_i then
    8   Update(MODEL)
    9   mse_i ← mse_v
10 end
11 else
12   MODEL ← LoadPreviousWeights()
13   n_I ← n_I + 1
14 end
15 end
  
```

---

Training with iterative early stop



# Training

**Algorithm 2:** Training with iterative early stop cycle.

---

```

Input:  $MODEL$  as the input model
Input:  $D_{train}$  as the training data set
Input:  $D_{val}$  as the validation data set
Input:  $N_I$  as the stop patience for iterative training cycle
Input:  $N_E$  as the early stop patience (epochs) for training
Input:  $B_{size}$  as the mini-batch size
Output:  $MODEL$  as the full-precision output model

// Initial training and evaluation
1  $Train(MODEL, D_{train}, D_{val}, N_E, B_{size})$ 
2  $mse_i \leftarrow Evaluate(MODEL, D_{val})$ 
3  $n_I \leftarrow 0$ 
4 while  $n_I < N_I$  do
    // Iterative early stop cycle
    5  $Train(MODEL, D_{train}, D_{val}, N_E, B_{size})$ 
    6  $mse_v \leftarrow Evaluate(MODEL, D_{val})$ 
    7 if  $mse_v < mse_i$  then
        8  $Update(MODEL)$ 
        9  $mse_i \leftarrow mse_v$ 
    10 end if
    11 else
        12  $MODEL \leftarrow LoadPreviousWeights()$ 
        13  $n_I \leftarrow n_I + 1$ 
    14 end if
15 end while

```

---

**Algorithm 2:** OnMiniBatchUpdate\_Callback.

---

```

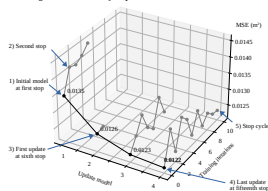
Input:  $MODEL$  as the full-precision input model
Input:  $E_{size}$  as the target exponent bits size
Input:  $M_{size}$  as the target mantissa bits size
Input:  $D_{train}$  as the training data set
Input:  $D_{val}$  as the validation data set
Input:  $N_{ep}$  as the number of epochs
Input:  $B_{size}$  as the mini-batch size
Output:  $MODEL$  as the quantized output model

// Quantize
1  $MODEL \leftarrow QuantizeTraining(MODEL, E_{size}, M_{size})$ 
2 if  $1 < epoch$  then
    // Update model after first epoch
    3  $mse_v \leftarrow Evaluate(MODEL, D_{val})$ 
    4 if  $mse_v < mse_i$  then
        5  $Update(MODEL)$ 
        6  $mse_i \leftarrow mse_v$ 
    7 end if
8 end if

```

---

**Training with iterative early stop**



# Training

---

**Algorithm 2:** Training with iterative early stop cycle.
 

---

```

Input:  $MODEL$  as the input model
Input:  $D_{train}$  as the training data set
Input:  $D_{val}$  as the validation data set
Input:  $N_I$  as the stop patience for iterative training cycle
Input:  $N_E$  as the early stop patience (epochs) for training
Input:  $B_{size}$  as the mini-batch size
Output:  $MODEL$  as the full-precision output model

// Initial training and evaluation
1  $Train(MODEL, D_{train}, D_{val}, N_E, B_{size})$ 
2  $mse_i \leftarrow Evaluate(MODEL, D_{val})$ 
3  $n_I \leftarrow 0$ 
4 while  $n_I < N_I$  do
    // Iterative early stop cycle
    5  $Train(MODEL, D_{train}, D_{val}, N_E, B_{size})$ 
    6  $mse_v \leftarrow Evaluate(MODEL, D_{val})$ 
    7 if  $mse_v < mse_i$  then
        8  $Update(MODEL)$ 
        9  $mse_i \leftarrow mse_v$ 
    10 end if
    11 else
        12  $MODEL \leftarrow LoadPreviousWeights()$ 
        13  $n_I \leftarrow n_I + 1$ 
    14 end if
15 end while
  
```

---



---

**Algorithm 2:** OnMiniBatchUpdate\_Callback.
 

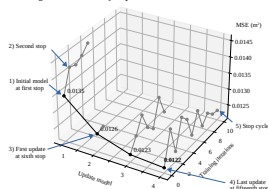
---

```

Input:  $MODEL$  as the full-precision input model
Input:  $E_{size}$  as the target exponent bits size
Input:  $M_{size}$  as the target mantissa bits size
Input:  $D_{train}$  as the training data set
Input:  $D_{val}$  as the validation data set
Input:  $N_{ep}$  as the number of epochs
Input:  $B_{size}$  as the mini-batch size
Output:  $MODEL$  as the quantized output model

// Quantize
1  $MODEL \leftarrow QuantizeTraining(MODEL, E_{size}, M_{size})$ 
2 if  $1 < epoch$  then
    // Update model after first epoch
    3  $mse_v \leftarrow Evaluate(MODEL, D_{val})$ 
    4 if  $mse_v < mse_i$  then
        5  $Update(MODEL)$ 
        6  $mse_i \leftarrow mse_v$ 
    7 end if
8 end if
  
```

---

**Training with iterative early stop**

**Quantization-aware training**
