

Adapting Neural Networks to Fixed-Point and Integer Arithmetic

Hanane BENMAGHNA

hanane.benmagnia@etudiant.univ-perp.fr

Supervisors: Matthieu MARTEL and Yasmine SELADJI

Laboratory of Mathematics and Physics (LAMPS)

University of Perpignan Via Domitia

June 20, 2019



Summary

- 1 Introduction
- 2 Neural Networks
- 3 Fixed-Point and Integer Arithmetic
- 4 Contributions
 - Operations between operands with the same format
 - Operations between operands with different formats
 - Format of weighted sum
- 5 Conclusion & Perspectives

1 Introduction

2 Neural Networks

3 Fixed-Point and Integer Arithmetic

4 Contributions

- Operations between operands with the same format
- Operations between operands with different formats
- Format of weighted sum

5 Conclusion & Perspectives

Introduction



Difficulty to run tasks in real time



Embedded systems have often low memory



Difficulty maintaining efficiency with limited resources

New Problem:

Artificial Intelligence and especially Neural Networks are increasingly used in Embedded Systems !!!

Contribution

Input: A trained NN, working at some precision (Floating-Point Format)

Output: Code generation (the new NN) to "simulate" the original NN using integers and fixed-point arithmetic

Correctness: The new NN with smaller data types behaves **almost** like the original NN

- Function Approximation: results stay in a desired hull
- Classification : only $x\%$ decisions differ from original NN

"Simulate":

Same behavior as the original network with $x\%$ error

1 Introduction

2 Neural Networks

3 Fixed-Point and Integer Arithmetic

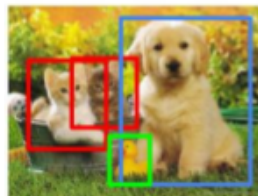
4 Contributions

- Operations between operands with the same format
- Operations between operands with different formats
- Format of weighted sum

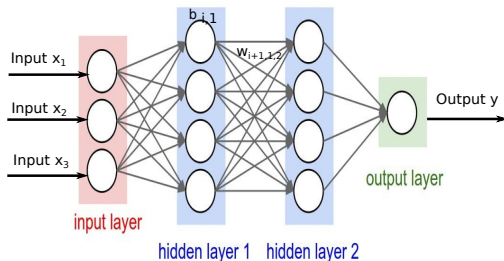
5 Conclusion & Perspectives

Neural Networks (NN)

- The conception is inspired by the biological functioning of neurons
- NN create fast classifications and generalizations
- NN are used on a variety of tasks (object recognition, speech and handwriting recognition, machine translation, medical diagnosis...)



CAT, DOG, DUCK



1 Introduction

2 Neural Networks

3 Fixed-Point and Integer Arithmetic

4 Contributions

- Operations between operands with the same format
- Operations between operands with different formats
- Format of weighted sum

5 Conclusion & Perspectives

Fixed-Point Format (FPF)

Float numbers are represented by **integers** in a format $\langle m, l \rangle$ or $\langle I, F \rangle$

Definition

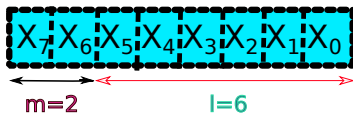
For $X = (x_{n-1}x_{n-2}\dots x_1x_0)_\beta$ and $x = X.\beta^{-l}$, we have: $p = m + l + 1$.

x : Fixed-Point number with implicit scale factor β^{-l} ,

X : Integer representation, p : Word **length**, m : The **most** significant bit,

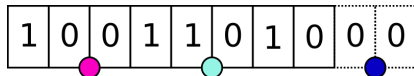
l : The **least** significant bit

Representation of X in a format $\langle m, l \rangle = \langle 2, 6 \rangle$



⚠ No exponent, the programmer must take care of the formats !

Different interpretations of the same integer representation



A Fixed-Point number in different format

Scaling Factor Symbol	Format	Value of X	Value of x
purple	$\langle 2, 6 \rangle$	$(10011010)_2 = (154)_{10}$	$(10.011010)_2 = (2.40625)_{10}$
cyan	$\langle 5, 3 \rangle$	$(10011010)_2 = (154)_{10}$	$(10011.010)_2 = (19.25)_{10}$
blue	$\langle 9, -1 \rangle$	$(10011010)_2 = (154)_{10}$	$(100110100)_2 = (308)_{10}$

Representation of 154 in different formats

1 Introduction

2 Neural Networks

3 Fixed-Point and Integer Arithmetic

4 Contributions

- Operations between operands with the same format
- Operations between operands with different formats
- Format of weighted sum

5 Conclusion & Perspectives

- 1 Introduction
- 2 Neural Networks
- 3 Fixed-Point and Integer Arithmetic
- 4 Contributions
 - Operations between operands with the same format
 - Operations between operands with different formats
 - Format of weighted sum
- 5 Conclusion & Perspectives

Addition

Let : x_1 and x_2 FP numbers with FPF $\langle i, f \rangle$:

$$x_1 + x_2 = x_3$$

$$\langle i, f \rangle + \langle i, f \rangle = \langle i, f \rangle \quad (1)$$

$$\begin{cases} 0 < i < p \\ 0 < i + f < p \\ f = p - i \end{cases} \quad (2)$$

Example :

$$\begin{array}{r} 000000.010 \quad \langle 6, 3 \rangle \\ + \\ 001000.000 \quad \langle 6, 3 \rangle \\ \hline 001000.010 \quad \langle 6, 3 \rangle \end{array}$$

Multiplication

Let : x_1 and x_2 FP numbers with FPF $\langle i, f \rangle$:

$$x_1 \times x_2 = x_3$$

$$\langle i, f \rangle \times \langle i, f \rangle = \langle i - f, 2f \rangle \quad (3)$$

$$\left\{ \begin{array}{l} 0 < i < p \\ 0 < i + f < p \\ f = p - i \end{array} \right. \quad (4)$$

$$x_1 \gg k = x_2$$

$$\langle i, f \rangle \gg k = \langle i+k, f-k \rangle \quad (5)$$

$$\left\{ \begin{array}{l} 0 < i < p \\ 0 < i+f < p \\ f = p-i \end{array} \right. \quad (6)$$

Shift Left

$$x_1 \ll k = x_2$$

$$\langle i, f \rangle \ll k = \langle i - k, f + k \rangle \quad (7)$$

$$\left\{ \begin{array}{l} 0 < i < p \\ 0 < i + f < p \\ f = p - i \end{array} \right. \quad (8)$$

AND, OR, XOR

$$x_1 \text{ AND } x_2 = x_3$$

$$\langle i, f \rangle \text{ AND } \langle i, f \rangle = \langle i, f \rangle \quad (9)$$

$$\begin{cases} 0 < i < p \\ 0 < i + f < p \\ f = p - i \end{cases} \quad (10)$$

Same format for OR and XOR!

- 1 Introduction
- 2 Neural Networks
- 3 Fixed-Point and Integer Arithmetic
- 4 Contributions
 - Operations between operands with the same format
 - Operations between operands with different formats
 - Format of weighted sum
- 5 Conclusion & Perspectives

Addition

$$x_1 + x_2 = x_3$$

$$\langle i_1, f_1 \rangle + \langle i_2, f_2 \rangle = \langle i_3, f_3 \rangle \quad (11)$$

$$\left\{ \begin{array}{l} 0 < i_3 < p \\ 0 < i_3 + f_3 < p \\ i_3 = \max(i_1, i_2) \\ f_3 = p - i_3 \end{array} \right. \quad (12)$$

Example 1: fixed format

00000.01	<5,2>	
+		
1000.0001	<4,4>	
<hr/>		
?	<6,3>	<6,3> 

Example 1

$$\begin{array}{r} 00000.01 \quad \langle 5, 2 \rangle \\ + \\ 1000.0001 \quad \langle 4, 4 \rangle \\ \hline ? \quad \langle 6, 3 \rangle \end{array} \quad \xrightarrow{\langle 6, 3 \rangle} \quad \begin{array}{r} 000000.010 \quad \langle 6, 3 \rangle \\ + \\ 001000.000 \quad \langle 6, 3 \rangle \\ \hline 001000.010 \quad \langle 6, 3 \rangle \end{array}$$

Example 2

Example 2: fixed precision

$$\begin{array}{r} 00000.01 \quad \langle 5, 2 \rangle \\ + \\ \hline 1000.0001 \quad \langle 4, 4 \rangle \\ \hline ? \quad \langle ?, ? \rangle \end{array} \quad \begin{array}{l} \text{---} \rightarrow \\ p=9 \end{array}$$

Multiplication

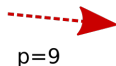
$$x_1 \times x_2 = x_3$$

$$\langle i_1, f_1 \rangle \times \langle i_2, f_2 \rangle = \langle i_3, f_3 \rangle \quad (13)$$

$$\left\{ \begin{array}{l} 0 < i_3 < p \\ 0 < i_3 + f_3 < p \\ i_3 = i_1 + i_2 - 1 \\ f_3 = p - i_3 \end{array} \right. \quad (14)$$

Example: fixed precision

00000.01	<5,2>
*	
1000.0001	<4,4>
<hr/>	
00000010.000001	<8,6>
= 2.0156	



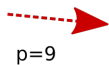
Example

$$\begin{array}{r} 00000.01 \\ * \\ 1000.0001 \\ \hline 00000010.000001 \\ = 2.0156 \end{array}$$

<5,2>

<4,4>

<8,6 >



p=9

$$\begin{array}{r} 00000.01 \\ * \\ 1000.0001 \\ \hline 00000010.0 \\ = 2.0 \end{array}$$

<5,2>

<4,4>

<8,1 >

error = 0.0156

AND, OR, XOR

$$x_1 \text{ AND } x_2 = x_3$$
$$\langle i_1, f_1 \rangle \text{ AND } \langle i_2, f_2 \rangle = \langle i_3, f_3 \rangle \quad (15)$$

$$\left\{ \begin{array}{l} 0 < i_3 < p \\ 0 < i_3 + f_3 < p \\ i_3 = \max(i_1, i_2) \\ f_3 = p - i_3 \end{array} \right. \quad (16)$$

Same format for OR and XOR!

- 1 Introduction
- 2 Neural Networks
- 3 Fixed-Point and Integer Arithmetic
- 4 Contributions
 - Operations between operands with the same format
 - Operations between operands with different formats
 - Format of weighted sum
- 5 Conclusion & Perspectives

Format of weighted sum

Let :

$$\begin{aligned}
 f(\hat{\mathbf{x}}) &= W\hat{\mathbf{x}} + \mathbf{b} \\
 f(\hat{\mathbf{x}}) &= \sum_{j=0}^n w_{ij} \hat{x}_j + b_i \\
 f(\hat{\mathbf{x}}) &= \underbrace{w_{i1} \hat{x}_1}_{\langle i_2, f_2 \rangle \times \langle i_1, f_1 \rangle} + \underbrace{w_{i2} \hat{x}_2}_{\langle i_2, f_2 \rangle \times \langle i_1, f_1 \rangle} + \dots + \underbrace{w_{in} \hat{x}_n}_{\langle i_2, f_2 \rangle \times \langle i_1, f_1 \rangle} + \underbrace{b_i}_{\langle i_3, f_3 \rangle} \\
 &= \underbrace{\langle i_1+i_2-1, f_1+f_2 \rangle}_{\langle i_4, f_4 \rangle = \langle i_1+i_2-1, f_1+f_2 \rangle} + \dots + \langle i_1+i_2-1, f_1+f_2 \rangle + \dots + \langle i_1+i_2-1, f_1+f_2 \rangle + \dots \\
 &= \langle i_5, f_5 \rangle = \langle \max(i_4, i_3), p - \max(i_4, i_3) \rangle = \langle \max((i_1+i_2-1), i_3), p - \max((i_1+i_2-1), i_3) \rangle
 \end{aligned}$$

$$\boxed{\langle \mathbf{i}_5, \mathbf{f}_5 \rangle = \langle \mathbf{max}((\mathbf{i}_1 + \mathbf{i}_2 - 1), \mathbf{i}_3), \mathbf{p} - \mathbf{max}((\mathbf{i}_1 + \mathbf{i}_2 - 1)) \rangle}$$

(17)

\mathbf{W} : matrix of weights, $\hat{\mathbf{x}}$: inputs and \mathbf{b} : bias

$\langle \mathbf{i}_1, \mathbf{f}_1 \rangle$: representation format of the input vector $\hat{\mathbf{x}}$

$\langle \mathbf{i}_2, \mathbf{f}_2 \rangle$: representation format of \mathbf{W}

$\langle \mathbf{i}_3, \mathbf{f}_3 \rangle$: representation format of bias \mathbf{b}

$\langle \mathbf{i}_5, \mathbf{f}_5 \rangle$: representation format of the result $f(\hat{\mathbf{x}})$

- 1 Introduction
- 2 Neural Networks
- 3 Fixed-Point and Integer Arithmetic
- 4 Contributions
 - Operations between operands with the same format
 - Operations between operands with different formats
 - Format of weighted sum
- 5 Conclusion & Perspectives

Conclusion & Perspectives

- Reducing the size and execution time of large NN with this approach
- Implementing a code transformation tool
- Experiments on large NN
- Handling classifiers and interpolators.

Related Work

- Dutta S. Jha S., Sankaranarayanan S. Tiwari A. (2018). “Output Range Analysis for Deep Feedforward Neural Networks”. In: *NFM*.
- Gehr T. Mirman M., Drachler-Cohen D. Tsankov P. Chaudhurri S. Vechev M.T. (2018). “AI²: Safety and Robustness Certification of Neural Networks with Abstract Interpretation”. In: *SP*.
- Martel, Matthieu (2017). “Floating-Point Format Inference in Mixed-Precision”. In: *NFM*.
- Singh G. Gehr T., Puschel M. Vechev M.T. (2019). “An Abstract Domain for Certifying Neural Networks”. In: *POPL*.

- Gehr T. 2018 Singh G. 2019

- Gehr T. 2018 Singh G. 2019
- Martel 2017
- Dutta S. 2018