



HAL
open science

Formal description of ML models for unambiguous implementation

Adrien Gauffriau, Iryna De Albuquerque Silva, Claire Pagetti

► **To cite this version:**

Adrien Gauffriau, Iryna De Albuquerque Silva, Claire Pagetti. Formal description of ML models for unambiguous implementation. 12th European Congress on Embedded Real Time Software and Systems (ERTS 2024), Jun 2024, Toulouse, France. hal-04167435v1

HAL Id: hal-04167435

<https://hal.science/hal-04167435v1>

Submitted on 20 Jul 2023 (v1), last revised 17 Jun 2024 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Formal description of ML models for unambiguous implementation

Adrien Gauffriau and Claire Pagetti
Airbus, France – ONERA, France

July 20, 2023

Abstract

Implementing deep neural networks in safety critical systems, in particular in the aeronautical domain, will require to offer adequate specification paradigms to preserve the semantics of the trained model on the final hardware platform. We propose to extend the NNEF language in order to allow traceable distribution and parallelisation optimizations of a trained model. We show how such a specification can be implemented in CUDA on a XAVIER platform.

1 Introduction

Machine learning (ML) applications have been gaining considerable attention in the field of transportation. However, their use in real-life operational safety-critical products, in particular in the aeronautical domain subject to stringent certification, raises several issues regarding functional correctness, compliance with requirements, formal verification, safety or implementation. In order to tackle those issues, new guidelines – named AS6983 standard [EUR21] – are currently drafted by the EUROCAE WG-114/SAE G-34 joint working group that cover the whole spectrum of the system development including the data sets composition, the ML model design and its implementation. In this paper, we focus only on the implementation of the ML model.

1.1 Context

In the ML current practices, a *training framework* is used to design an ML model and the resulting ML model is then deployed on the target with a *deployment framework*. It is up to the training framework or a designer to export the trained model description in an exchange format and up to the deployment framework to import the ML model description. The left part of figure 1 shows those practices. The deployment framework is most of the time an ML model interpreter, that can accommodate any type of ML model architecture, and that allocates at runtime the execution on the different available accelerators (e.g. GPU, FPGA) of the target. These ML frameworks have been designed 1) to ease as much as possible the deployment of models for the users and 2) to optimize as much as possible the execution performance (usually expressed in trillion operations per second – TeraOp/s). As a result, they are very impressive and allow for complex deployments and optimizations. Those are hard, if not impossible, to reproduce for a programmer without using ML libraries (e.g. CUDNN on NVIDIA).

The counterpart of such an approach is 1) the absence or limited information of internal computation and allocation; 2) small modifications and adaptations of the ML model (e.g. when exporting the ML model description or when quantizing on the fly some matrix multiplication to execute on deep learning accelerator – DLA). If this grey/black box approach is acceptable and suitable for mainstream applications, it is a blocking point for highly safety-critical applications. As a result, the main objective of the AS6983 standard with respect to the implementation process is the *semantics preservation* of the off-line trained model on the final hardware platform. This means that the execution of the ML model in the training framework should be exactly reproduced on the embedded target during execution. To reach that objective, an alternative development process is proposed as illustrated in the right hand part of

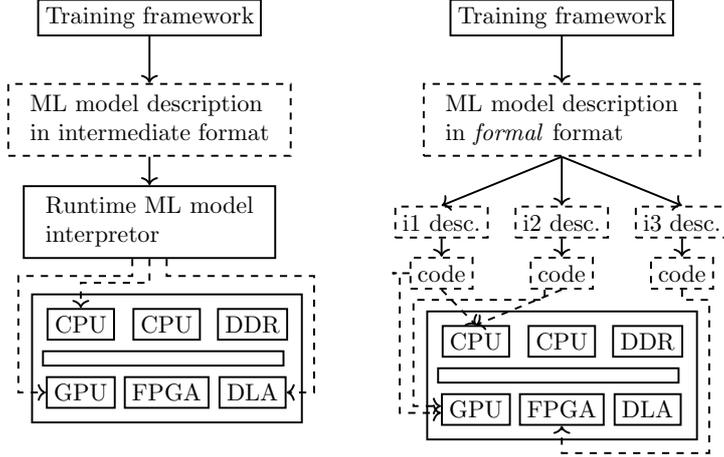


Figure 1: Left part: current ML deployment practice and right part: proposed aeronautical practice

figure 1. The cost is the reduced performance and deployment capacities compared to ML deployment approaches.

1.2 Assurance development process to ensure semantic preservation

The principle of this trustable development process (figure 1) is to ensure the semantic preservation at each step of the development cycle (e.g. between training framework and description, between description and code, between code and executable). This is guaranteed thanks to a series of requirements provided in the AS6983.

Requirement 1. First, the trained model must be formally and unambiguously defined in what we call an adequate *format*. Such a *format* must come with a formal syntax and semantic, and should be agnostic of any (training and/or deployment) framework.

Requirement 2. Second, the implementation process should allow several deployments on hardware platforms and it must be known beforehand how the ML model will be mapped on the accelerators and when. This entails in particular that the *format* should allow several types of deployment such as distribution, parallelization or pipelining. Thus, there is a step that consists in splitting the ML model description (in the chosen format) as a series of *item* descriptions. Indeed, in the avionics domain, a target processor is decomposed as a set of software (SW) and/or hardware (HW) *items*. Let us consider for instance an ULTRASCALE+ (ZCU102) platform [Xil19]: it is composed of several ARM cores, a GPU and an FPGA. If the ML model is spread over the different components, in particular on one ARM and the FPGA, there will necessarily be several items. Indeed, the ARM associated code will be considered as one SW item and will go through the ED-12C/DO-178C [RTC11] development process. Whereas the hardware design of the FPGA will be considered as another item (HW) and will go through the ED-80/DO-254 [RTC00] development process. This is the reason why we need to be able to decompose a global ML model description into several pieces equivalent to the global model description.

Requirement 3. Third, the implementation has to follow the usual aeronautical development standards (e.g. ED-12C/DO-178C [RTC11] for software).

1.3 Contributions

Our general objective is to define an approach compatible with the AS6983 requirements presented in the prior section. We focus on a representative subset of deep neural networks (DNNs) that is feed-forward neural networks trained off-line (see section 2.1). Our main goal is the definition of an adequate format, with a formal syntax and semantics, able to describe both 1) a global ML model, and 2) any parallelized allocation on several items, the behaviour of which is equivalent to the global model.

For requirement 1. There are several initiatives to propose an intermediate format between trained models and their implementation such as ONNX [BLZ⁺19] (Open Neural Network Exchange). After a thorough evaluation of different existing formats, we identified NNEF (Neural Network Exchange Format) [The22] as the most suited for our purpose: syntax and semantics are public and moreover the authors

made a strong effort to provide a formal specification.

Since NNEF provided a potential candidate, we decided to construct our format on top of it. As it is now, NNEF describes formally the global ML model. Indeed, the semantics of NNEF is almost fully defined (see section 2). What is however missing in NNEF is the clear formalization of the *execution model*, that is the formal behaviour behind a series of NNEF *instructions*. To fix this missing element, we rely on Petri nets, a usual representation of program behaviours [Pet77]. This choice is also consistent with the need to express distributed behaviour on items, as Petri nets also allow to model all combinations of execution: sequence, pipeline, recursion or parallel [Old86].

For requirement 2. We illustrate in section 3 why decomposing ML model into items is of interest. The first contribution (see section 4) is the definition of an extension of NNEF to allow the distribution of an ML model onto different items. To that end, we extend the syntax and semantics of NNEF. We rely in particular on *logical data exchange* among items to distribute the computation and express the semantics with *coloured Petri nets*. We formally show that the synchronization of the Petri nets behaves as the *original* non distributed Petri net.

For requirement 3. We do not propose a complete DO-178C compatible approach but instead show a reasonable implementation approach on the XAVIER platform [NVI19a] that we believe could be with some effort compliant with the ED-12C/DO-178C [RTC11]. To that end, we define an automatic CUDA-based code generation from an ML description in the extended NNEF format and assess in our evaluation its code traceability, the deterministic behaviour on the platform and the semantics preservation from the trained framework. All of that is detailed in section 5.

We will discuss the related work in section 6 before concluding.

2 Format of description – NNEF

KHRONOS standardization group¹ has defined the NNEF (Neural Network Exchange format) format with a specification that provides a syntax and a semantic. After outlining the definition of deep neural network, we will present NNEF. We focus on the NNEF syntax and semantics elements needed for our purpose. The reader can refer to [The22] for a complete description of NNEF. We end the section by defining the execution model semantics with a translation into Petri net.

2.1 Brief Reminder on Neural Network

There are multiple ways to define DNNs: directed graphs, computational graphs or simply the mathematical functions transforming the input into the output. NNEF relies on a computation graph approach. In any case, the input of a neural network can be seen as a multi-dimensional vector also called *tensor*.

Definition 1 (Tensor) A 3D-tensor T is represented by its size (n_c, n_h, n_w) where n_c the number of channels (or feature maps), n_h is the height and n_w the width. We denote by T_{x_1, x_2, x_3} the value of T for the indices x_1, x_2, x_3 .

A feed-forward deep neural network consists of a series of layers executed without recursion.

Definition 2 (Feed-forward Deep Neural Network) A feed-forward neural network $N = \langle Q, V \rangle$ is a directed graph where:

- $Q = \{l_1, \dots, l_n\}$ the nodes are the layers l_i . A layer can be of many types such as an activation, a convolution or a pooling. A layer comes with a set of parameters (e.g. weights or stride);
- $Q_i \subseteq Q$ is the set of input layers and $Q_f \subseteq Q$ is the set of final layers, in particular $Q_i \cap Q_f = \emptyset$;
- $V \subseteq Q \times Q$ the arrows represent the data flow within the neural network. $\text{Pre}(i) = \{j | (l_j, l_i) \in V\}$ is the set of layers l_j such that output(l_j) is an input of l_i . An input layer $j \in Q_i$ only consumes input tensors, i.e $\text{Pre}(j) = \emptyset$;

We can compute the ancestors of a layer, $\text{Pre}^*(j) = \cup \text{Pre}^n(j)$ where $\text{Pre}^n(j) = \cup_{j_i \in \text{Pre}^{n-1}(j)} \text{Pre}(j_i)$ and $\text{Pre}^1(j) = \text{Pre}(j)$

Because the network is feed-forward, $j \notin \text{Pre}^*(j)$.

¹<https://www.khronos.org/>

Example 1 The LE_{NET}-5 [LBD⁺ 89] is a usual CNN developed for hand written digit images recognition. Its structure (see figure 2) consists of 8 consecutive layers: convolution, pooling, flat and dense. $Q_i = \{\text{conv1}\}$, $Q_f = \{\text{softmax}\}$, and for all layers $i \geq 2$, $\text{Pre}(i) = \{i - 1\}$.

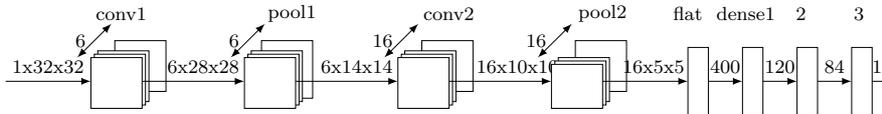


Figure 2: LE_{NET}-5 neural network

The size of the input / output tensors are shown on the figure. The first 2D-convolution conv1 which takes inputs of size $1 \times 32 \times 32$, is composed of 6 kernels of size $1 \times 5 \times 5$ and a stride of $(1, 1)$. Because the sizes of the input and the kernel are well balanced taking into account the stride, no padding is needed. The activation function ReLu is applied to the outputs. The first pooling layer pool1 is a max pooling with stride $(2, 2)$ and window $(2, 2)$. The second 2D-convolution conv2 is composed of 16 kernels of size $6 \times 5 \times 5$ and of a stride $(1, 1)$. The activation function ReLu is applied to the outputs. The second pooling layer pool2 is a max pooling with stride $(2, 2)$ and window $(2, 2)$. The 3D-tensor of size $16 \times 5 \times 5$ is flattened in a 1D-tensor of size 400. Three dense layers are then applied, each followed by a ReLu. A dense layer consists in applying a linear transformation $W \cdot I + B$ where W are the weights and B the bias. Finally, the post-processing is a softmax.

Remark 1 In this article we are only focusing on the inference topic. Thus, whereas NNEF provides the possibility to express the batch size (number of inputs computing in one inference), we only consider a batch size of 1. The possibility to increase the batch is motivated by the learning phase to improve the optimization convergence (stochastic gradient).

Definition 3 (Function associated to a DNN) The function f_N computed by a DNN $N = \langle Q, V \rangle$ is the composition of the functions computed by each layer. In case of a directed DNN - i.e. where $Q_i = \{l_1\}$, $Q_f = \{l_n\}$ and $\forall i \geq 2$, $\text{Pre}(i) = \{i - 1\}$ - the function is simply given by $f_N = f_{l_n} \circ \dots \circ f_{l_1}$ where \circ is the composition operator. Otherwise, $f_N = (f_{l_{j_1}}, \dots, f_{l_{j_p}})$ where $Q_f = \{l_{j_1}, \dots, l_{j_p}\}$ and the functions are defined as follows: $\forall i$ such that $l_i \notin Q_i$, $f_{l_i}(f_{l_{j_1}}, \dots, f_{l_{j_n}})$ where $\{j_1, \dots, j_n\} = \text{Pre}(i)$. The semantics of each elementary function (e.g. a 2D-convolution) is given for instance in NNEF by providing its mathematical function.

2.2 Neural Network Description in NNEF

Our first goal concerns the definition of a format, with a formal syntax and semantics, able to describe any ML model such as the LE_{NET}-5 of example 1. Let us explain why NNEF fulfils this goal. An NNEF description is composed of two parts: 1) A computation graph described in a human readable text file; and 2) the parameters provided in multiple raw data file. The fact that the description is textual is important for the traceability between the specification (output of the training framework) and the embedded code.

The computation graph file describes all parameters needed and operations to be done. More precisely, each line of the computation graph is an elementary instruction (NNEF *compound fragment*) that may be split into several atomic operations (NNEF *primitives*). The result of each instruction is stored in a named variable, that represents a tensor, and which can be used as input for other instruction(s). To compute an operation all its inputs variables shall be available.

Remark 2 NNEF description follows a SSA (static single assignment) form [CFR⁺ 89] which helps and simplifies the implementation process. It is usual to translate a program in its associated SSA form before compilation or optimization passes [BBD⁺ 17].

Example 2 Let us illustrate how to specify a neural network with an example. The NNEF textual specification of the LE_{NET}-5 detailed in example 1 is given in the listing 1. First, all parameters are declared and stored as variables. e1 is the input tensor of size $1 \times 32 \times 32$, v1 contains the 6 kernels of size $1 \times 5 \times 5$ of the first convolution and v2 is the bias applied at the end of the first convolution. The parameters needed by the layers should all be defined as variables in the description file.

```

graph torch_jit_export(e1) -> (out)
{
  e1 = external<scalar>(shape = [1, 1, 32, 32]);
  v1 = variable<scalar>(shape = [6, 1, 5, 5], label = 'v1');
  v2 = variable<scalar>(shape = [1, 6], label = 'v2');
  v3 = variable<scalar>(shape = [16, 6, 5, 5], label = 'v3');
  v4 = ...; v5 = ...; v6 = ...; v7 = ...;
  v8 = ...; v9 = ...; v10 = ...;

  o1 = conv(e1, v1, v2, stride = [1, 1], dilation = [1, 1], padding = [(0, 0),
    (0, 0)], groups = 1);
  o2 = relu(o1);
  o3 = max_pool(o2, size = [1, 1, 2, 2], stride = [1, 1, 2, 2], dilation = [1,
    1, 1, 1], padding = [(0, 0), (0, 0), (0, 0), (0, 0)], border = 'ignore
    ');
  o4 = conv(o3, v3, v4, ...; o5 = relu(o4);
  o6 = max_pool(o5, ...; o7 = reshape(o6, shape = [0, -1]);
  o8 = linear(o7, v5, v6); o9 = relu(o8);
  o10 = linear(o9, v7, v8); o11 = relu(o10);
  o12 = linear(o11, v9, v10);
  out = softmax(o12, axes = [1]);
}

```

Listing 1: LUNET-5 with NNEF syntax

After the variables declaration, comes the computation graph itself. The output of the first convolution is stored in the variable `o1`. When calling the function / compound fragment `conv`, the user must instantiate the full set of parameters for this type of layer: input tensor, kernels, bias, stride, dilation, padding and groups. Every parameter appears explicitly in the definition and there is no ambiguity. For instance, the way to declare the padding explicitly states how the padding applies on top / bottom / right / left. After the convolution, the activation function has to be explicitly applied to `o1` and thus is not hidden in the convolution, ensuring again an unambiguous description. The first pooling layer results in variable `o3`. Reading the description, we recognize the LUNET-5 detailed before. The flat layer is encoded with a more expressive function `reshape` that allows several reshaping. The dense layer is called `linear`. The post-processing `softmax` is also explicit.

The NNEF specification also provides the link between instructions (e.g. `conv`) appearing in the file and their associated mathematical functions. We illustrate this with the max pooling layer only.

2.3 Max Pooling Layer Semantics

In order to illustrate issues that we may without having a formal description we propose to use the Max Pooling example. Let us first remind the functional semantics of a max pooling layer where a padding is applied and no dilatation is applied. Thus, the function is defined by $\mathcal{P}ool_{k,s} \circ \mathcal{P}_{p,v}$ where each function is defined below.

Definition 4 (Padding associated function – $\mathcal{P}_{p,v}$) Let $p = (p_t, p_b, p_l, p_r)$ be a 4-tuple of integers representing the padding to be applied on each border of a 3D-tensor and v the float value to be used for the padding. The padding function $\mathcal{P}_{p,v}$ applied on a 3D-tensor I of size (n_h, n_w, n_c) outputs a 3D-tensor $O = \mathcal{P}_{p,v}(I)$ of size (o_h, o_w, o_c) with $o_h = n_h + p_t + p_b$, $o_w = n_w + p_l + p_r$ and $o_c = n_c$ such that

$$O_{x,y,z} = \begin{cases} v & \text{if } (x \leq p_t) \text{ or } (x > n_h + p_t) \text{ or} \\ & (y \leq p_l) \text{ or } (y > n_w + p_l) \\ I_{x-p_t, y-p_l, z} & \text{otherwise} \end{cases}$$

Definition 5 (Pooling layer associated function – $\mathcal{P}ool_{k,s}$) Let $s = (s_h, s_w)$ be the stride parameters and $k = (k_h, k_w)$ be the height and width of the window. The pooling applied on a 3D-tensor I of size (n_h, n_w, n_c) outputs a 3D-tensor $O = \mathcal{P}ool_{k,s}(I)$ of size (o_h, o_w, o_c) with $o_h = \left\lfloor \frac{n_h - k_h}{s_h} + 1 \right\rfloor$, $o_w = \left\lfloor \frac{n_w - k_w}{s_w} + 1 \right\rfloor$ and $o_c = n_c$ with $O_{x,y,z} = \max(I[s_h \cdot (x-1) + 1 : s_h \cdot (x-1) + k_h + 1][s_w \cdot (y-1) + 1 :$

$s_w \cdot (y - 1) + k_w + 1][z]$. Here, $I[s_{11} : s_{21}, \dots, s_{1k} : s_{2k}]$ represents the slice of I of all the values $I_{s_{11}+x_1, \dots, s_{1k}+x_k}$ with $i \in [1, k]$ and $x_i \in [1, s_{2i} - s_{1i}]$.

The NNEF syntax of a max pooling layer describes the padding values with an enumerate string. Ignoring the border for a max pooling layer is equivalent to pad with the minimum float value (MIN_F). This corresponds to the neutral operator of the max function. Looking now at the *max_pool* elementary instruction according to NNEF documentation [The22], it is defined in the pseudo-code by 2 atomic operations. 1. *argmax_pool* that computes an array of index (corresponding to the max in each pool); 2. *sample* that returns for an array of index, an array with corresponding values. This pseudo-code indeed encodes the expected function $\mathcal{P}ool_{k,s} \circ \mathcal{P}_{p,v}$.

Of the importance of unambiguous description. We propose to highlight the importance of making unambiguous textual descriptions of ML models by comparing two state-of-the-art training frameworks, namely PYTORCH and KERAS (with TENSORFLOW). More specifically, we compare their way of encoding max pooling. A first important remark is that it is not possible to pad the *borders* (top / bottom / left / right) with any combination.

- In KERAS, a padding inside a max pooling layer can only be declared by a string $\in \{$ "valid", "same" $\}$. "valid" means no elements to add, while "same" means that padding is added on the right and on the bottom borders only to fit the size of the pool.

```
# KERAS SYNTAX for pool1
MaxPool = tf.keras.layers.MaxPooling2D(pool_size=(2,2), strides=(2,2),
padding = 'same')

# NNEF for KERAS SYNTAX
max_pool_keras = max_pool(input, size = [1, 1, 2, 2], stride = [1, 1,
2, 2], dilation = [1, 1, 1, 1], padding = [(0, 0), (0, 0), (0, 1)
, (0, 1)], border = 'ignore');
```

This corresponds to $\mathcal{P}ool_{[2,2],[2,2]} \circ \mathcal{P}_{[0,1,0,1],MIN_F}$.

- In PYTORCH, a padding inside a max pooling layer can only be declared by one or two integers. In case of one integer, this defines the number of elements to add to each border (top, bottom, left and right). In presence of 2 integers, the first gives the number of elements to add at the top and bottom, and the second at the left and right.

```
#PYTORCH SYNTAX for pool1
MaxPool = NN.MaxPool2d(2,          2,          1)
# Kernel Size, Stride, Padding

#NNEF for PYTORCH SYNTAX
max_pool_torch = max_pool(input, size = [1, 1, 2, 2], stride = [1, 1,
2, 2], dilation = [1, 1, 1, 1], padding = [(0, 0), (0, 0), (1, 1),
(1, 1)], border = 'ignore');
```

This corresponds to $\mathcal{P}ool_{[2,2],[2,2]} \circ \mathcal{P}_{[1,1,1,1],MIN_F}$

The semantics of KERAS and PYTORCH are not equivalent, and there is no possibility to convert one into another at once. The only way to make a valid conversion is to explicitly add a padding layer before the max pooling. This difference of semantic between the two frameworks leads to issues when we want to convert a model from one framework to the other.

2.4 NNEF Execution Model Semantics

The semantics of the *execution model*, that is the formal behaviour behind a series of NNEF *instructions*, is not explicitly given by the standard. It assumes that one instruction can be executed when all its inputs are computed and available. Thus, executing the instructions in sequence following the order of the NNEF guarantees a correct execution. There are two types of instructions: those reading parameters from binary files or input tensors and layer-associated instructions based on one or several atomic operations. An instruction is always of the form

$$var = operation(v_1, \dots, v_k, cst_1, \dots, cst_j);$$

where var is a variable computed by the *operation* (any NNEF *fragment*), v_i are either variables computed beforehand, the input tensor or fixed parameters (e.g. kernels), and cst_i are constant (e.g. stride). The NNEF execution model can be formally expressed using the Petri net formalism [Pet77].

Translation 1 A NNEF description, composed of n instructions, generates a Petri net (P, T, V) where:

- the set of places P corresponds to all variables appearing in the NNEF description (i.e. n places for a description of n instructions);
 - a token in a place means that the associated variable is available for computation;
 - initially there are as many tokens in each parameter-associated place as the parameter is needed in the instructions and as many tokens in the input tensor place as the input tensor is used by the instructions;
 - there is a unique final place corresponding to the last variable computed in the NNEF file;
- the set of transition names T corresponds to all instructions appearing in the NNEF description;
- $V \subseteq 2^P \times T \times \mathbb{N} \times P$ defines the set of transitions. A transition can be fired if there is a token in all input places. When it fires, the transition removes a token from each of these places and generates as many tokens as defined on the edge in the unique output place.
 - each instruction $var = op(v_1, \dots, v_k, cst_1, \dots, cst_j)$ generates the transition given in figure 3 where p is the number of time var will be consumed by other instructions;
 - when only one token is produced by a transition, we omit the integer.

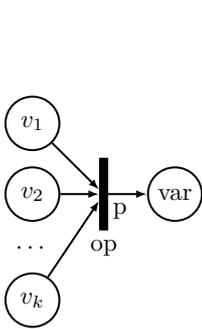


Figure 3: Translation of one instruction

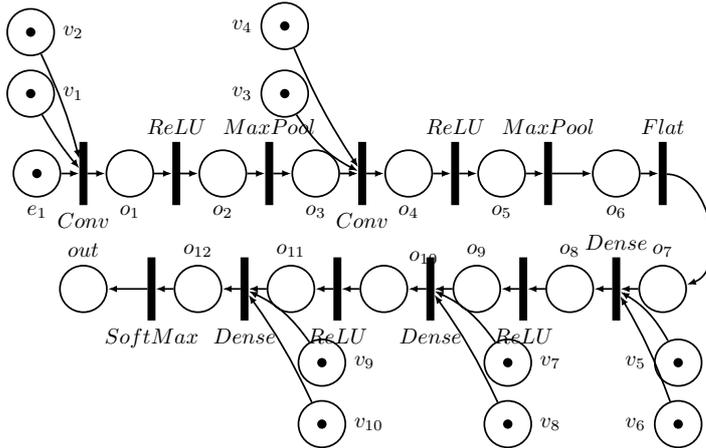


Figure 4: NNEF semantics of the LEnET-5 express with Petri net. Initial marking

The semantic of the Petri net clarifies the execution order that is unclear in the NNEF formalism. Especially, the order of the textual file should not impose a unique order when several valid ones may exist.

Example 3 The LEnET-5 model of figure 2 with its associated NNEF description in listing 1 has the associated Petri net given in figure 4. We recognize the instructions sequence that describes the computation of the neural network graph. There is a unique possible schedule for this NNEF description, but we will see later other NNEF models that accept several schedules (see section 3).

Remark 3 The Petri net of figure 4 only defines the semantics of a single inference pass. It is usual to repeat the inference pass in order to process new inputs (e.g. in a periodic manner). This can also be represented using the Petri formalism by sending back tokens to $e1$ and v_k places. For demonstration and clarity, subsequently in the paper, we always consider the semantics of a single inference.

A way to define the semantics of Petri nets is to compute the set of reachable *markings*, where a marking defines the number of tokens in each place at a given instant.

Definition 6 (Marking) Considering a Petri net with n places, a marking is defined as a vector $v \in \mathbb{N}^n$ giving the number of tokens $v[i]$ in the i -th place. This initial location of tokens is called the initial marking, this represents the starting state of the system. A final marking is a marking such that there is one token only in each final place and from which no transition can be fired any more.

Example 4 In the Petri net of figure 4, there are 24 states. There are 11 tokens in the initial making and there is a single token in the unique final marking (in the place out).

Property 1 (Initial marking and unique final marking) The unique initial marking is defined by the translation and consists of token(s) in the input tensor variable place and parameters-associated places. Because we consider feed-forward neural networks, there is unique final marking.

Definition 7 (Paths and semantics) A valid path starts from the initial marking m_i , lists a series of fireable transitions and ends in a final marking m_f , i.e. $m_i \xrightarrow{t_1} m_1 \xrightarrow{t_2} \dots \xrightarrow{t_i} m_f$. The semantics of a Petri net is the set of valid paths.

Property 2 (Possible executions of an NNEF description) Because we consider feed forward neural networks, the number of valid paths is finite and the valid paths correspond to all possible execution orders respecting the semantics of the ML model.

Semantics preserving code generation could lead to any implementation the path of which is valid. Full sequential code following the order of instructions of the NNEF file is one of them.

3 Is distribution needed?

Because we consider highly distributed platforms, a designer may choose to split the ML model into several parts in order to accelerate the execution and reduce the execution time. In particular, it could lead to developing parts independently and on different *items* following the aeronautics standards [SAE10]. In such a case, there should be a formal description for each item that will be the input specification for HW/SW item implementation.

Remark 4 Note that if the designer considers its platform as a unique item, the NNEF description will be the specification.

Looking at the LEnET-5 model or other classical models, it is unclear why distributing across into items would be beneficial.

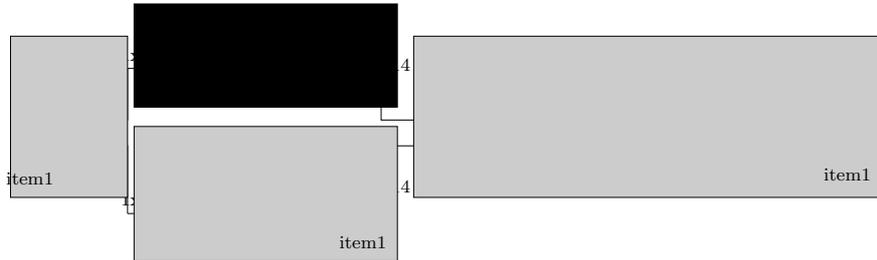
Example 5 (Off-loading convolution) Let us consider for instance the implementation of the LEnET-5 on an ULTRASCALE+ (ZCU102) platform [Xil19] composed of several ARM cores, a GPU and an FPGA. Let us assume that we choose to execute the convolution on the FPGA and all other layers on the ARM, there will necessarily be several items. The idea will be to offload the input tensor of each convolution on the FPGA and retrieve the feature maps from the FPGA (see figure 5).



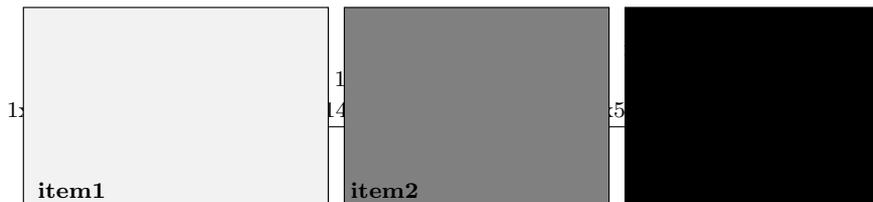
Figure 5: LEnET-5 with off-loading

A second type of distribution could be to refine the layers and exhibit more parallelism by distributing the computation of a layer across several items. It will be up to the designer to show the semantic preservation at this refinement level.

Example 6 (Parallelizing the layers) Looking again at the L_{ENET}-5, we can split the computation of the first two layers on two different items. Each item will do the convolution+maxpool on a part of the input image. An idea is to split the input image along the height and execute the two sub-parts on two different items (see figure below). In order to keep the semantics for the convolution, some overlap exists between the two sub-images.



Example 7 (Pipelining the computation) A third type of distribution could be the pipelining of the DNN execution. In this case, each item is in charge of computing a specific layer (or a group of layers). The first item computes the first layer(s) on the first tensor input and sends its output to the second item that will compute the second group of layer(s) while the first item starts processing the second input tensor. This classical mechanism enables to reduce (e.g. for video processing) the frame rate while degrading the latency. The depth of the pipeline is the number of inputs that can be handled at the same time among the pipeline. Taking into account the pipelining should be done in the Petri net by restoring the tokens in the initial places at the end of item 1 execution.



4 NNEF extension for multiple items

The purpose of this section is to propose a manner to separate the specification of each item so that any execution of the items respecting the specification will properly encode the global ML model. To that extent, we propose first to extend the syntax of NNEF to allow explicit parallelization and then to express the associated semantics with colored Petri nets.

4.1 Extension for item splitting

We first need to specify the item on which the description will be implemented. To that end, we enrich the graph definition with the keyword *graphitem* to provide the logical id of the HW or SW item.

Syntax1 GraphItem

```
<graph-definition> ::= <graph-declaration> <graph-declaration-item> <body>
```

```
<graph-declaration-item> ::= "graphitem" <identifier> <identifier>
    "(" <identifier-list> ")" "-$>$" "(" <identifier-list> ")"
```

Semantics 1 The first <identifier> refers to the item id, the second <identifier> is the name of the local node and elements of the <identifier-list> will be input or output of the graphitem. The semantic of the <graph-declaration-item> is such that all instructions within the <body> are executed by the item.

We also need to *exchange data* between several items and ensure that those exchanged data are available before computation. To that end, we introduce a new type of variable, namely *variablesync*. This references a variable that could be read from or write to another *graphitem*. We use a *fragment* to define this new type.

Syntax2 VariableSync

```
fragment variablesync<? = scalar>(shape: integer[]) -> ( output: tensor<?> )
```

Semantics 2 Each *variablesync* is a shared variable with a unique writer and possibly multiple readers. Writer is in charge to transmit the variable via the instruction *send_var* and each reader can access this data via *get_var* instruction.

We then define new NNEF instructions to send or get data between several *graphitem*.

Syntax3 get_var

```
fragment get_var<? = scalar>( source : graphitem, data : variablesync)
-> ( output: tensor<?> )
```

Semantics 3 *get_var* appears in each reader description. The output of this instruction is a local variable that gets the content of the shared variable and which is available for the caller item. Source is the item that writes and provides the shared variable the name of which is given by data.

Similarly, the writer must define the instruction to send a shared variable to other items.

Syntax4 send_var

```
fragment send_var<? = scalar>( dest : graphitem[], data : scalar)
-> ( output: variablesync )
```

Semantics 4 *send_var* appears in the writer description only. It takes as input the list of reader items and the name of the tensor that shall be sent. The output tensor is a global *variablesync* that will support the synchronization.

The rest of the NNEF syntax remains unchanged.

4.2 Splitting an NNEF description into multi-item descriptions

Initially, the DNN is described in a unique NNEF description as shown in section 2. Such a description contains 3 types of instructions:

- definition of input tensor(s);
- definition of fixed parameters;
- variables computed by each layer.

A splitting consists in partitioning the last type of instructions among the items, add the adequate definition(s) of tensors / fixed parameters and add the adequate *send_var* / *get_var*. The composition of item descriptions shall respect the semantics of the full NNEF description.

Example 8 Let us consider the following DNN with its associated Petri net (figure 6).

Let us assume that the DNN is allocated on 3 items such that o_1, o_6, o_7 and *out* are computed on item 1; o_2, o_3 are computed on item 2 and o_4, o_5 are computed on item 3. Thus, the description on the items are given as follows:

```

graph DNN(e1) -> (out)
{
  e1 = external...
  v1 = variable... 'variable1'); v2
    = ...; v3 = ...; v4 = ...;
  v5 = ...; v6 = ...; v7 = ...; v8
    = ...; v9 = ...; v10 = ...;
  o1 = conv(e1,v1,v2,... ;
  o2 = conv(o1,v3,v4,...; o3 = conv
    (o2,v5,v6,...;
  o4 = conv(o1,v7,v8,...; o5 = conv
    (o4,v5,v6,...;
  o6 = concat (o3, o5); o7 =
    flatten(o6,...;
  out = gemm(o7, v9, v10...;
}

```

Listing 2: Complete DNN NNEF

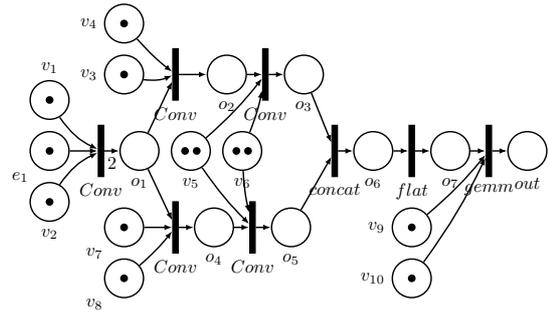


Figure 6: Petri associated to DNN of listing 2. Initial marking

```

graphitem ITEM1 DNN1([e1, vsync2, vsync3]) -> ([vsync1, out])
{
  e1 = external...;
  v1 = variable... 'variable1'); v2 = ...; v9 = ...; v10 = ...;
  vsync1 = variablesync<scalar> (shape ...;
  vsync2 = variablesync ...; vsync3 = variablesync ...;

  o1 = conv(e1,v1,v2,...;
  vsync1 = send_var ([ITEM2, ITEM3], o1);
  o3 = get_var(ITEM2, vsync2);
  o4 = get_var(ITEM3, vsync3);
  o6 = concat (o3, o4); o7 = flatten(o4,...;
  out = gemm(o7, v9, v10...;
}

graphitem ITEM2 DNN2([vsync1]) -> ([vsync2])
{
  v3 = variable... 'variable3'); v4 = ...; v5 = ...; v6 = ...;
  vsync1 = variablesync ...; vsync2 = variablesync ...;

  o1 = get_var(ITEM1, vsync1)
  o2 = conv(e1,v3,v4,...; o3 = conv(o2,v5,v6,...
  vsync2 = send_var ([ITEM1], o3);
}

graphitem ITEM3 DNN3([vsync1]) -> ([vsync3])
{
  v5 = ...; v6 = ...; v7 = ...; v8 = ...
  vsync1 = variablesync ...; vsync3 = variablesync...;

  o1 = get_var(ITEM1, vsync1)
  o4 = conv(o1,v7,v8,...; o5 = conv(o1,v5,v6,...
  vsync3 = send_var ([ITEM1], o5);
}

```

Listing 3: NNEF for all items

The union of the instructions of each item corresponds to the complete NNEF description with the additional variablesync and the communication instructions. Locally in the item, the pointers to the input tensor and fixed parameters must also be declared.

4.3 Petri-based semantic

We define the execution model semantics of multi-item descriptions using coloured Petri nets [JK09]. We associate a colour to each item where the colour is set to the tokens and edges (on which the coloured tokens transit).

Translation 2 *Let assume there are N items. We first apply the translation 1 for each item leading to N independent Petri nets (P_i, T_i, V_i) , each with a unique and distinct colour. For the new instruction, the translation is extended as follows:*

- a varsync does not generate any place;
- a get_var does not generate any transition;
- a send_var produces a transition sync with an incoming edge from the variable the content of which is transmitted.

The set of NNEF item descriptions generates a Petri net (P, T, V) which is roughly speaking the union of the N Petri nets (P_i, T_i, V_i) . More precisely,

- any input tensor or fixed parameter that is duplicated in the NNEF files appears in the Petri net of the item. Those duplicated places are merged. Because we use the same naming convention, $P = \cup P_i$;
- the initial tokens in the places are also merged leading to places with possibly multiple tokens and multiple colours;
- $T = \cup T_i \cup T'$ where T' are the transitions connecting places of one item to other items thanks to the sync transition. More precisely,
 - for each writer, there are k edges back from the sync label where k is the number of readers. The colour of the each arrow is the one of the reader and the number of tokens sent back corresponds to the number of time the shared variable appears in the reader item instructions;
 - for each reader, there is an edge from the writer place before sync to each transition requiring the shared data;
- As before, when a coloured token is present in a place, it means that the associated variable is available for the item identified by the colour and can be used by transition.

Example 9 *The Petri net associated to the example 8 is given in figure 7. We present the initial marking with colored tokens. Each color represents the state of one item. Compare to the figure 4, we expressend here the multi items semantics with synchronizations.*

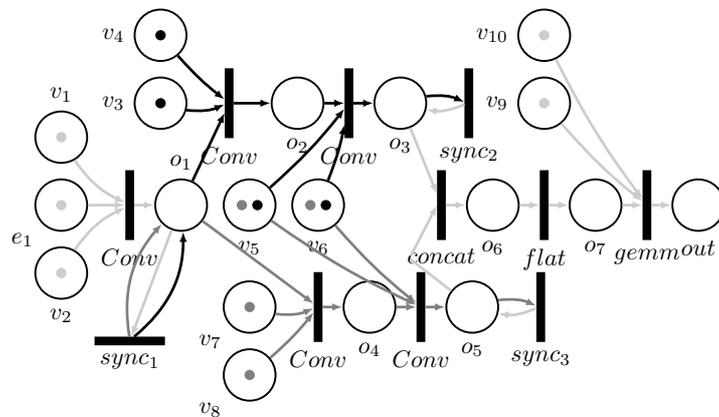


Figure 7: Semantic of the items synchronization

Property 3 (Equivalence between Petri nets) *The semantics of the multi-items behaviour is equivalent to the complete original ML model.*

Remark 5 *It is important to explicitly describe the `send_var` and `get_var` either in the NNEF files but also in the Petri net because items are supposed to be independent and segregated. In particular, an item X is not allowed to access the memory space of an item Y and interfere with its execution. This is classical in aeronautics, see Arinc 653 specification. The XTRATUM hypervisor [CRM⁺09] is an example of time and space partitioning hypervisor that provides communication with sampling and queuing ports (close to Arinc 653 requirements).*

5 Implementation and experiments

The previous sections showed how to fulfill the requirements 1 and 2 of the introduction. The purpose of this section is to give some hint of how a DO-178C compatible implementation process, as required per requirement 3, could be defined taking as input an extended NNEF specification. The considered target, a Jetson XAVIER TX system-on-chip, is composed of 6 Carmel ARM cores, a GPU, 2 deep learning accelerators (NVDLA) and other dedicated circuits. The use of a NVIDIA platform is mainly motivated by its availability and the ease to quickly deploy neural networks application. We will not discuss the adequation of GPU and CUDA implementation with DO-178C objectives because it is an open problem.

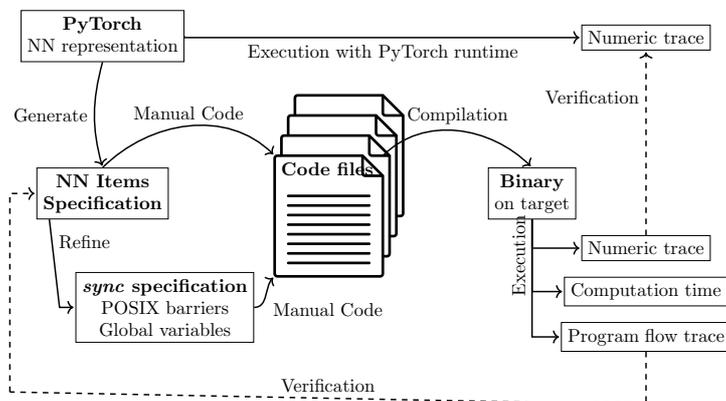


Figure 8: Implementation workflow

The proposed implementation process workflow is summarized in figure 8. The *sync* implementation relies on barrier mechanism and each item NNEF description is manually coded. In order to validate the semantic preservation, we made some instrumentation to show that: (i) the execution trace is included in all possible execution traces defined by the Petri net; (ii) the numerical precision is kept; (iii) the measured execution time does not vary. We will use the multi-items example presented in example 8 as the specification. For our experiments, each item is allocated to one CPU and all the CUDA cores of the GPU (grouped in a CUDA stream). As a consequence, we do not guarantee a segregation between items (as they share the GPU) we instead focus on a way to implement parallel operations of neural networks with a static code scheduling while preserving the semantics. As for the Petri semantics, we only developed a code for a single inference. Nevertheless, it is easy to slightly modify the code by adding loops to handle several input tensors.

5.1 Get/Send specification

We chose to implement `get_var` and `send_var` with 1) global variables stored in the SRAM of the XAVIER and 2) the POSIX barrier mechanism of the *pthread* library. A barrier b , shared among several processes, will block them as long as not all of them reach b . Such a behavior is strictly included in the semantics of the *sync* transition within the Petri net. However, it is not the most efficient as it prevents the sending item to proceed until all the receiving items reach b , whereas the semantics of *sync* transition only requires a receiver to wait for the sender (not the sender to wait for all receivers). Nonetheless, the barrier mechanism is optimal for our example because no sender has to process any instruction before a further `get_var` or stop execution.

5.2 Manual code generation

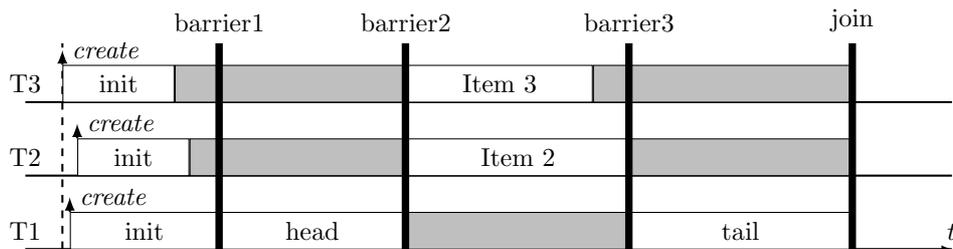
There are C and PYTHON interpreters of the NNEF format [NNE18] but only for traditional CPU target. Consequently, no existing tool supports our syntax extension nor state-of-the-art GPU. Thus we developed the code for each item using C++ and CUDA using the CUDNN library. Basically the C++ code is executed by the ARM processor whereas CUDA allows the definition of kernels that are executed synchronously by all CUDA cores. The CUDNN library is built on top of CUDA for executing common neural networks layers.

5.2.1 Software architecture

Practically, each type of layer is implemented using a dedicated C++ class that inherits from the abstract *Layer* class that defines common attributes and methods to be implemented (*init()* and *forward()*) by child classes. In effect, *init* statically allocates tensors and CUDNN descriptors while *forward* launches the layer computation based on CUDNN for Convolution and max pooling layers. Each item contains one object implementing a static scheduler. More precisely, during the *init* phase, each item creates an object for each layer which are stored in ordered lists. Thus, items 2 and 3 need one single ordered list whereas item 1 needs two ordered lists (one for the first part and one for the second part). During *forward*, the scheduler calls in order the *forward* of objects stored in ordered lists.

5.2.2 Scheduling

We define one separate thread for each item allocated to one CPU + CUDA cores. More precisely, synchronizations between threads use `pthread_barrier_t` and associated APIs (`barrier_init` and `barrier_wait`). Barriers synchronize accesses to shared variables.



The execution sequence starts with the 3 threads creation on the CPUs and then reaches the first synchronization barrier. Then Item1 thread calls the *forward* method of layers of the head (until *send_var*) while Items 2 and 3 threads wait for the second synchronization barrier. After, the second synchronization barrier, Items 2 and 3 threads call *forward* method of their layers while Item1 thread waits for the third synchronization barrier. After the third synchronization barrier, Item 1 thread calls *forward* method of layers of the tail. At last, threads join and exit.

5.3 Semantic preservation of the Petri net

The first analysis aims at verifying that all observed scheduling of layers on the XAVIER respects the Petri net semantics. Because we use a static scheduler, all schedules should behave as shown in section 5.3 which is included in the semantics of coloured Petri net of figure 7. For that, we logged each start/end of branches and layers and we stressed the robustness of the implementation by adding some temporal noises (sleep in the code).

All observed traces respected the schedule of section 5.3 with some timing variations. When observing the implementation with no noise, execution traces of Item 2 and 3 are interleaved on the GPU. When adding a wait of 1s at the beginning of Item 2 (just after barrier1), all layers of Item 3 were executed before those of Item 2.

5.4 Semantics preservation of the function

The second instrumentation mechanism aims at checking that the functional semantics of the DNN is preserved. We achieve this by re-implementing the NNEF specification in PYTORCH. Then we define 100 random vectors that we run both on the PYTORCH implementation and on the C++ implementation on

the NVIDIA target. Finally, we compute the overall average error mean between both executions for the 100 runs.

We were not able to find the exact convolution algorithm of PYTORCH. We think that it exists a non documented heuristic that calls the best algorithm (considering execution time) depending of the convolution parameters and available hardware. According to the CUDNN documentation [NVI19b], it is possible to select the convolution algorithm among a list, but details of the implementation are not given. Thus, there may be a discrepancy between convolutions that we cannot fix. The average error mean is extremely small 1.10^{-7} for FLOAT32 using 3 CUDNN algorithms (namely gemm, Winograd and direct). Nevertheless numeric precision results for this experiment are in an acceptable range that is very close to the available numeric precision of floating point representation and this also is observed by other frameworks [SCGP22].

5.5 Measured Execution Time (MET)

One objective of the DO-178C that we did not mention until now is the capacity to estimate the Worst Case Execution Time (WCET). Due to the complexity of NVIDIA target, a formal demonstration using static analysis may be difficult. But at least, a good property is a low variation of the measured execution time among several executions. In our case, the generated code does not contain any IF-THEN-ELSE patterns or dynamic loop conditions. Thus, the variability is only linked to the hardware behavior. We measured the MET of the complete DNN and of the first convolution of Item 1 over 10 runs. We rely on the *nsys* tool from NVIDIA to get timing measurements.

	Mean(MET)	MIN(MET)	(MAX(MET))	STD(MET)
First Conv	324 2976 ns	322 688 ns	326 656 ns	1.45 ns
DNN	24 257 μ s	16 285 μ s	53 950 μ s	13 526 μ s

The MET of the first convolution is very stable with a very low jitter. The MET distribution of the DNN is large and to understand why, we need to investigate the low level behaviour. NVIDIA GPUs are black-boxes processors on which we cannot guarantee worst-case execution time [AA21, ACR22].

6 Related work

There are plenty of different formats but no consensus within the community. State-of-the-art frameworks like PYTORCH, or TENSORFLOW rely on custom black-box formats built on top of *protocol buffers* [Goo01] developed by Google. A *protocol buffer* is a structured binary format that is not human readable and requires conversion tools and template files to be interpreted. Thus, the syntax is not formally defined and specification of layers are only available through documentation website. For example, TENSORFLOW proposes the *.h5* [HDF01] format and KERAS format [Ten15] builds on top of protocol buffer. Moreover the way to save a neural network is not unified among frameworks and may evolve every updated version with poor backward compatibility. When moving from caffe to caffe2 (known today as PYTORCH), the caffe [JSD⁺14] format was not supported anymore.

All previous formats were developed specifically for training frameworks (open source or proprietary) without any objective for sharing models. Their main purpose was to allow saving and reloading previous trained models without too much consideration on syntax and semantics. ONNX [BLZ⁺19] and NNEF [The22] were developed with the objective to be independent from frameworks. ONNX is still based on binary protocol buffers [Goo01] (so without a textual syntax) but is proposing a functional semantics through its github site. On the contrary NNEF is proposing a textual format with a syntax and a semantics that is formally defined in a specification. Unfortunately, the NNEF format suffers from a small community and tools supporting the format.

NNet [The22] format is an example of ad-hoc format. RELUPLEX neural networks examples [KBD⁺17] are in NNet. It is based on textual files but without definition of syntax and semantics. Import and export tools are provided, but its utilization for sharing neural networks between teams remains supremely painful.

Some other formats like [CE-17, Apa18, LAB⁺21a] tackle the need for intermediate representation between a neural network description and an implementation on a target. Especially this supports different optimizations passes like layers folding or low level tensor manipulation description like in LLVM [Lat02]. We consider that we are closest to programming language than from a neural network description format. Most of the time ONNX or NNEF are used as input like in [LAB⁺21b, PBCPB20, JBL⁺20].

Because DNN are data-flow, it is natural to translate them into *synchronous languages*. There are some works such as [LFG20] that proposed to encode them as Synchronous Dataflow Graphs or SCADE tool suite [CPPP18] which is currently developing a DNN libraries. Once the translation is done, it is then possible to reuse all the qualified code generation tools. This is complementary to our work as we could use the NNEF description to generate the SCADE program for instance. To the best of our knowledge, none of actual available neural network description formats propose solution for describing multi items implementation with concerns on sharing variables among them.

7 Conclusion

We have proposed a formal extension of NNEF that takes into account the *execution model* of a description and allows for the modification of a description of a trained model to define traceable distribution and parallelisation optimizations that preserves the semantics while improving the execution time compared to a pure sequential approach. We have also proposed a code generation strategy based on barriers for exchanging data between items.

As a future work, we would like to extend the NNEF tools to parse our extensions and define an automatic code generation in CUDA. We will also improve the GPU utilization to ensure more predictable behaviour. We will in particular reuse works from [AA21, ATV⁺21].

8 Acknowledgement

This work has benefited from the AI Interdisciplinary Institute ANITI, which is funded by the French “Investing for the Future – PIA3” program under the Grant agreement ANR-19-P3IA-0004. The authors gratefully acknowledge the support of the PHYLOG 2 project.

References

- [AA21] Tanya Amert and James H. Anderson. Cupid^{rt}: Detecting improper GPU usage in real-time applications. In *24th IEEE International Symposium on Real-Time Distributed Computing, ISORC 2021, Daegu, South Korea, June 1-3, 2021*, pages 86–95. IEEE, 2021.
- [ACR22] Michaël Adalbert, Thomas Carle, and Christine Rochange. PasTiS: building an NVIDIA Pascal GPU simulator for embedded AI applications. In *11th European Congress on Embedded Real-Time Systems (ERTS 2022)*, 2022.
- [Apa18] Apache. TVM, 2018. <https://tvm.apache.org/>.
- [ATV⁺21] Tanya Amert, Zelin Tong, Sergey Voronov, Joshua Bakita, F. Donelson Smith, and James H. Anderson. Timewall: Enabling time partitioning for real-time multicore+accelerator platforms. In *42nd IEEE Real-Time Systems Symposium, RTSS 2021, Dortmund, Germany, December 7-10, 2021*, pages 455–468. IEEE, 2021.
- [BBD⁺17] Timothy Bourke, Lélío Brun, Pierre-Évariste Dagand, Xavier Leroy, Marc Pouzet, and Lionel Rieg. A formally verified compiler for lustre. In Albert Cohen and Martin T. Vechev, editors, *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017, Barcelona, Spain, June 18-23, 2017*, pages 586–601. ACM, 2017.
- [BLZ⁺19] Junjie Bai, Fang Lu, Ke Zhang, et al. Onnx: Open neural network exchange. <https://onnx.ai/>, 2019.
- [CE-17] CE-List. N2D2 frameworks, 2017. <https://github.com/CEA-LIST/N2D2>.
- [CFR⁺89] Ron Cytron, Jeanne Ferrante, Barry K Rosen, Mark N Wegman, and F Kenneth Zadeck. An efficient method of computing static single assignment form. In *Proceedings of the 16th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 25–35, 1989.

- [CPPP18] Jean-Louis Colaço, Bruno Pagano, Cédric Pasteur, and Marc Pouzet. Scade 6: From a kahn semantics to a kahn implementation for multicore. In Hiren Patel, Tom J. Kazmierski, and Sebastian Steinhorst, editors, *2018 Forum on Specification & Design Languages, FDL 2018, Garching, Germany, September 10-12, 2018*, pages 5–16. IEEE, 2018.
- [CRM⁺09] Alfons Crespo, I. Ripoll, M. Masmano, P. Arberet, and Metge Jean-Jacques. Xtratum: An open source hypervisor for tsp embedded systems in aerospace. pages 31–, 05 2009.
- [EUR21] EUROCAE WG-114/SAE joint group. Certification/approval of aeronautical systems based on AI, 2021. on going standardization.
- [Goo01] Google. Protocol Buffers, 2001. <https://portal.hdfgroup.org/display/HDF5/Introduction+to+HDF5>.
- [HDF01] HDF group. Introduction to HDF5, 2001. <https://www.hdfgroup.org/solutions/hdf5/>.
- [JBL⁺20] Tian Jin, Gheorghe-Teodor Bercea, Tung D. Le, Tong Chen, Gong Su, Haruki Imai, Yasushi Negishi, Anh Leu, Kevin O’Brien, Kiyokuni Kawachiya, and Alexandre E. Eichenberger. Compiling onnx neural network models using mlir, 2020.
- [JK09] Kurt Jensen and Lars M. Kristensen. *Coloured Petri Nets: Modelling and Validation of Concurrent Systems*. Springer Publishing Company, Incorporated, 1st edition, 2009.
- [JSD⁺14] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [KBD⁺17] Guy Katz, Clark W. Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. ReLuplex: An efficient SMT solver for verifying deep neural networks. *CoRR*, abs/1702.01135, 2017.
- [LAB⁺21a] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. MLIR: Scaling compiler infrastructure for domain specific computation. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, pages 2–14, 2021.
- [LAB⁺21b] Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques A. Pienaar, et al. MLIR: scaling compiler infrastructure for domain specific computation. In Jae W. Lee, Mary Lou Soffa, and Ayal Zaks, editors, *International Symposium on Code Generation and Optimization (CGO)*, pages 2–14, 2021.
- [Lat02] Chris Lattner. LLVM: An Infrastructure for Multi-Stage Optimization. Master’s thesis, Computer Science Dept., University of Illinois at Urbana-Champaign, Urbana, IL, Dec 2002. See <http://llvm.cs.uiuc.edu>.
- [LBD⁺89] Yann LeCun, Bernhard E. Boser, John S. Denker, Donnie Henderson, Richard E. Howard, Wayne E. Hubbard, and Lawrence D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural Comput.*, 1(4):541–551, 1989.
- [LFG20] Daniel Luenemann, Maher Fakh, and Kim Gruettner. Capturing neural-networks as synchronous dataflow graphs. In *MBMV 2020 - Methods and Description Languages for Modelling and Verification of Circuits and Systems; GMM/ITG/GI-Workshop*, pages 1–10, 2020.
- [NNE18] NNEF project. NNEF GitHub tools repository, 2018. <https://github.com/ChronosGroup/NNEF-Tools/>.
- [NVI19a] NVIDIA - Dustin Franklin. Introducing the Jetson Xavier Nx, 2019. <https://developer.nvidia.com/blog/jetson-xavier-nx-the-worlds-smallest-ai-supercomputer/>.
- [NVI19b] NVIDIA developer team. CudNN library documentation, 2019. <https://docs.nvidia.com/deeplearning/cudnn/api/index.html>.
- [Old86] Ernst-Rüdiger Olderog. Operational petri net semantics for ccsp. In *European Workshop on Applications and Theory in Petri Nets*, pages 196–223. Springer, 1986.

- [PBCPB20] Hugo Pompougnac, Ulysse Beaugnon, Albert Cohen, and Dumitru Potop-Butucaru. From SSA to Synchronous Concurrency and Back. Research Report RR-9380, INRIA Sophia Antipolis - Méditerranée (France), December 2020.
- [Pet77] James L Peterson. Petri nets. *ACM Computing Surveys (CSUR)*, 9(3):223–252, 1977.
- [RTC00] RTCA/EUROCAE. DO-254/ED-80 - Design Assurance Guidance For Airborne Electronic Hardware, 2000.
- [RTC11] RTCA/EUROCAE. DO-178C/ED-12C - Software Considerations in Airborne Systems and Equipment Certification, 2011.
- [SAE10] SAE/EUROCAE. Aerospace Recommended Practices ARP4754a/ed-79a- development of civil aircraft and systems, 2010.
- [SCGP22] Iryna De Albuquerque Silva, Thomas Carle, Adrien Gauffriau, and Claire Pagetti. ACE-TONE: predictable programming framework for ML applications in safety-critical systems. In *34th Euromicro Conference on Real-Time Systems ECRTS 2022*, pages 3:1–3:19, 2022.
- [Ten15] TensorFlow. Keras Model, 2015. https://www.tensorflow.org/api_docs/python/tf/keras/Model.
- [The22] The Khronos NNEF Working Group. Neural Network Exchange Format – Version 1.0.5, 2022. <https://registry.khronos.org/NNEF/specs/1.0/nnef-1.0.5.pdf>.
- [Xil19] Xilinx. ZCU102 Evaluation Board – User Guide, 2019. <https://docs.xilinx.com/v/u/en-US/ug1182-zcu102-eval-bd>.