

Deep Neural Networks on the Versat Reconfigurable Processor

João Pedro Costa Luís Cardoso
Electrical and Computer Engineering Department
Instituto Superior Técnico
Lisbon, Portugal
joao.pedro.cardoso@tecnico.ulisboa.pt

Abstract—This paper presents the problem of accelerating the execution of Deep Neural Networks (DNNs) using Coarse-Grained Reconfigurable Arrays (CGRAs), with special emphasis on compiling a DNN description into code that runs on a CPU/CGRA system. This topic is a vast one, so this paper focuses on simulating DeepVersat, a CPU/CGRA system suitable for DNN, and tools for running any Convolutional Neural Network on DeepVersat. The tools presented in this paper allow for architectural exploration and optimization, and a dramatic reduction of the development time.

Index Terms—CGRA, Versat, Darknet, Convolutional Neural Networks, Deep Neural Networks

I. INTRODUCTION

Neural Networks have been an object of study since the 1940s but until the beginning of this decade their applications were limited and did not play a major role in computer vision conferences. With its meteoric rise in research, several solutions to accelerate this algorithm have appeared, from Field Programmable Gate Arrays (FPGA) to Application Specific Integrated Circuits (ASIC) implementations.

Convolutional Neural Networks (CNNs) are a particular kind of DNN where the output values of the neurons in one layer are convolved with a kernel to produce the input values of the neurons of the next layer. This algorithm is compute bound, that is, its performance depends on how fast it can do certain calculations, and depend less on the memory access time. Namely, the convolutional layers take approximately 90% of the computation time.

The acceleration of these workloads is a matter of importance for today's applications such as image processing for object recognition or simply to enhance certain images. Other uses like instant translation and virtual assistants are applications of neural networks and their acceleration is of vital importance to bring them into the Internet of Things.

A suitable circuit to accelerate DNNs in hardware is the CGRA. A CGRA is a collection of Functional Units and memories with programmable interconnections to form computational datapaths. A CGRA can be implemented in both FPGAs and ASICs. CGRAs can be reconfigured much faster than FPGAs, as they have much fewer configuration bits. If reconfiguration is done at runtime, CGRAs add temporal scalability to the spacial scalability that characterizes FPGAs. Moreover, partial reconfiguration is much easier to do in

CGRAs compared to FPGAs which further speeds up reconfiguration time. Another advantage of CGRAs are the fact that they can be programmed entirely in software, contrasting with the large development time of customized Intellectual Property (IP) blocks. The Coarse Grain Reconfigurable Array (CGRA) is a midway acceleration solution between FPGAs, which are flexible but large, power-hungry, and difficult to reprogram, and ASICs, which are fast but generally not programmable.

However, mapping a specific DNN to a CGRA requires knowledge of its architecture, latencies, and register configurations, which may become a lengthy process, especially if the user wants to explore the design space for several DNN configurations. An automatic compiler that can map a standard DNN description into CPU/CGRA code would dramatically decrease the time to market of its users. Currently, there are equivalent tools for CPUs and GPUs and even for FPGAs.

The DeepVersat CGRA is the DNN accelerator to improve the performance of the DNNs in embedded hardware. Another objective is to increase the versatility of the Versat API and offer new functions to simplify the development of new software. One of these functions is a generic convolution for Versat which can, independently of the hardware configuration, configure the convolution to have the highest performance possible on the available functional units while being dynamic and to avoid developer work to adapt to new convolutions.

II. DEEP NEURAL NETWORKS

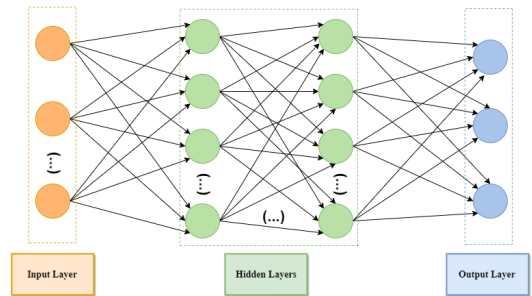


Fig. 1. Deep Neural Network Structure

A Neural Network (NN) is an interconnected group of nodes that follow a computational model that propagates data

forward while processing. The earliest NNs were proposed by McCulloch and Pitts [1], in which a neuron has a linear part, based on an aggregation of data and then a non-linear part called the activation function, which is applied to the aggregate sum. The issue with using only one neuron is that it is not able to be used in non-linear separable problems. By aggregating several neurons in layers and the input of each neuron as in figure 1 being based on the previous layers, that problem can be eliminated.

Each input to a neuron contributes differently to the output. The share is dependent on the weight value. These are obtained by training the network through various techniques, one of which is called Deep Supervised Learning [2]. For a certain input, there is an expected output and the real output of the NN. Then the loss function (the difference) is calculated and the weight values are iteratively modified for improving the outputs of the NN.

A Deep Neural Network (DNN) is a Neural Network that uses this approach for learning. It has multiple hidden layers and it can model complex non-linear relationships. If the activation function is non-polynomial, it satisfies the Universal approximation problem [3].

One of the limitations of traditional NNs is the complexity of layer interconnections. Using as an example the hand digit recognition problem and MNIST data set, composed of 28x28 grayscale images [4], in a traditional fully connected NN, a neuron from the second layer would have 28x28 weights. That is 3.136 kiloBytes per neuron of weight values while using 32-bit floating-point numbers (FP32). When building a more complex network for image recognition, the computational complexity grows quadratically with the number of neurons per layer.

A. Convolutional Neural Networks

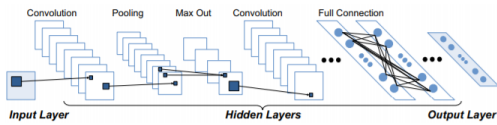


Fig. 2. CNN architecture example, taken from [5]

Convolutional Neural Networks (CNN) are a class of DNNs used in Image and Video recognition due to their shift invariance characteristic. They were first proposed in the 1980s but it was not until 2012 with AlexNet [6] that CNNs took off. Fundamentally, CNNs are a regularized version of Multilayer Perceptrons (MLP). These networks fix the complexity issue discussed as each neuron is only connected to a few neurons of the previous layer.

a) *Convolutional Layer*: In a typical CNN, not all layers are convolutional, but the convolutional layers are the most compute-intensive ones. CNNs take input images with 3 dimensions (width, height, and color space); for the following convolutional layers 3D arrays are used (width, height, and

number of channels). For the earlier example of the MNIST data set, the input would have dimensions 28x28x1 as it is a 2D image in grayscale.

To compute a neuron in the next layer we use the convolution equation 1 aided by Figure 3.

$$x_j^{l+1} = \delta \left(\sum_{i \in M_j} x_i^l * k_{ij}^{l+1} + b_j^{l+1} \right) \quad (1)$$

where x_j^{l+1} is the output, δ is the activation function, which depends on the architecture, x_i^l is the input of the convolution layer, k_{ij}^{l+1} is the kernel of the said layer which is obtained by training the network, and b_j^{l+1} is the bias.

Thus an output neuron depends only on a small region of the input which is called the local receptive field.

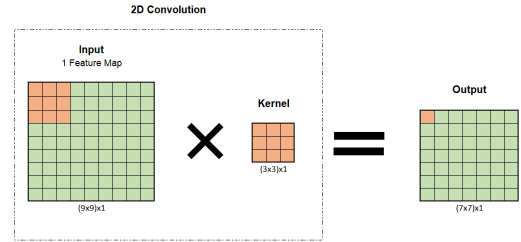


Fig. 3. 2D convolution with stride = 1 and without zero padding

The output's dimensions depend on several parameters of the convolution such as zero-padding and stride. The former means to add zeros around the edges of the input matrix. The latter means the step used for the convolution, if the value is e.g. 2, it will skip a pixel each iteration of the convolution. Equation ?? can be used to calculate the output size.

$$n^{l+1} = \frac{n^l - b^l + 2 \times p}{s} + 1 \quad (2)$$

where n is the width/height of the input of layer l , b is the width/height of the kernel, p is zero-padding while s is the stride.

The number of channels of the output is equal to the number of filters in the convolutional layer.

b) *Pooling Layer*: The MaxPool or AvgPool are layers used in Convolutional Neural Networks to downsample the feature maps to make the output maps less sensitive to the location of the features.

Maximum Pooling or MaxPool, like is suggested in its name groups $n \times n$ points and outputs the pixel with the highest value. The output will have its size lowered by n times. The Average Pooling or AvgPool, instead takes all of the input points and calculates the average. Downsampling can also be achieved by using convolutions with stride 2 and padding equal to 1. Upsample layers can be also used that turn each pixel into n^2 , where n is the number of times the output will be bigger than the input.

c) *Fully Connected Layer*: The fully connected layer is mostly used for classification in the final layers of the NN. It associates the feature map with the respective labels. It takes

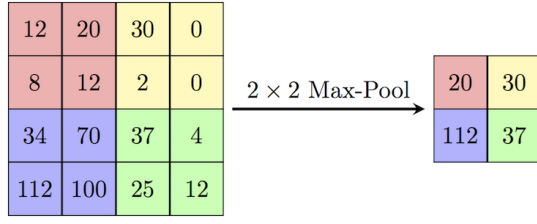


Fig. 4. Simple example of a max pool layer [7]

the 3D vector and outputs a single vector thus it is also known as flatten. Equation 3 describes the operation.

$$x_j^{l+1} = \delta\left(\sum_i (x_i^l \times w_{ji}^{l+1}) + b_j^{l+1}\right) \quad (3)$$

where w_{ji}^{l+1} are the weights associated with a specific input for each output.

d) *Route & Shortcut Layer*: The Shortcut layer or skip connection was first introduced in Resnet [8]. It allows connecting of the previous layer to another to allow the flow of information across layers. The Route layer, used in Yolov3 [9], concatenates 2 layers in depth (channel) or skips the layer forward. This is used after the detection layer in Yolov3 to extract other features.

e) *Dropout Layer*: This type of layer was conceived to avoid overfitting [10] by dropping the neurons with a probability below the threshold. In Figure 5, there is a graphical representation.

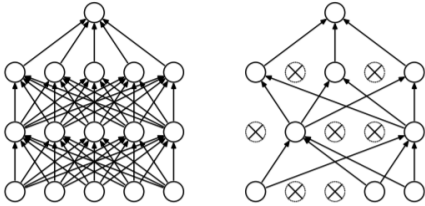


Fig. 5. Dropout if applied to all layers, adapted from [10]

f) *Activation Functions*: Activation Functions (AF) are functions used in each layer of a NN to compute the weighted sum of input and biases, which is used to give a value to a neuron. Non-linear AFs are used to transform linear inputs into non-linear outputs. While training Deep Neural Networks, vanishing and exploding gradients are common issues, in other words, after successive multiplications of the loss gradient, the values tend to 0 or infinity and thus the gradient disappears. AFs help mitigate this issue by keeping the gradient within specific limits. The most popular activation functions can be found in table I.

B. Frameworks for Neural Networks

To run a Neural Network model there are several popular frameworks like Tensorflow, PyTorch, Caffe, and Darknet. Their purpose is to offer abstraction to software developers

Activation Functions	Computation Equation
Sigmoid	$f(x) = \frac{1}{1 + e^{-x}}$
Tanh	$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
Softmax	$f(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$
ReLU	$f(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$
LReLU	$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x \leq 0 \end{cases}$
ELU	$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - 1 & \text{if } x \leq 0 \end{cases}$

TABLE I
POPULAR ACTIVATION FUNCTIONS

that want to run these networks. They also offer programming for different platforms like Nvidia GPUs by using the CUDA API.

1) *Darknet*: Darknet [11] is an open-source neural network framework written in C and CUDA. It is used as the backbone for Yolov3 [9] and supports several different network configurations such as AlexNet and Resnet. It utilizes a network configuration file (.cfg) and a weights file (.weights) as input for inference.

```
[convolutional]
batch_normalize=1
filters=32
size=3
stride=1
pad=1
activation=leaky
```

Listing 1. cfg code for a Convolutional Layer used in Yolov3 [9]

In Listing 1, there is a snippet of the file featuring a convolution layer with 32 kernels of size 3x3. It has stride 1 and zero padding of 1, meaning the output size equals the input size. The input size can be calculated by analyzing the previous layers and the network parameters. The network parameters in Listing 2 includes data to be used for training while only the first three parameters are needed for inference.

```
[net]
width=608
height=608
channels=3

learning_rate=0.001
burn_in=1000
max_batches=500200
policy=steps
steps=400000,450000
scales=.1,.1
```

Listing 2. cfg code for the network parameters

2) *Caffe*: Convolutional Architecture for Fast Feature Embedding (Caffe) [12] is also an open-source framework written

in C++ with a Python interface. Caffe exports a neural network by serializing it using the Google Protocol Buffers (ProtoBuf) serialization library. Each network has 2 prototxt files:

- `deploy.prototxt`- File that describes the structure of the network that can be deployed for inference.
- `train_val.prototxt`- File that includes structure for training. it includes the extra layers used to aid the training and validation process.

The Python interface helps generate these files. For inference only the deploy file matters. In Listing 3, there is a snippet of a deployed file.

```
name: "AlexNet"
layer {
  name: "data"
  type: "Input"
  top: "data"
  input_param { shape: { dim: 10 dim: 3 dim: 227 dim: 227 } }
}
layer {
  name: "conv1"
  type: "Convolution"
  bottom: "data"
  top: "conv1"
  param {
    lr_mult: 1
    decay_mult: 1
  }
  param {
    lr_mult: 2
    decay_mult: 0
  }
  convolution_param {
    num_output: 96
    kernel_size: 11
    stride: 4
  }
}
```

Listing 3. prototxt file for the input data and the first convolution layer of AlexNet [6]

III. DEEPVERSAT

Versat is a Coarse-Grained Reconfigurable Array (CGRA) Architecture. CGRAs are in-between Field Programmable Gate Arrays (FPGA) and general purpose processors (GPP). The former is fully reconfigurable and the highest performance for a workload can be achieved as the Architecture is tailored to the workload. GPPs on the other hand, are not reconfigurable and thus slower but are more generic and can process different workloads. While FPGAs have granularity at the gate level, CGRAs have granularity at the functional unit level. They are configurable at run-time and the datapath can be changed in-between runs.

A. Versat Architecture

The Versat Architecture [13]–[16] is depicted in Figure 6. It is composed of the following modules: DMA, Controller, Program Memory, Control File Registry, Data Engine, and Configuration module. The Controller accesses the modules through the control bus. The code made in assembly or C is loaded into the program Memory (RAM) where the user can write to the configuration module for the Versat runs. Between runs of the Data Engine, the Controller can start doing the next run configuration and calculations.

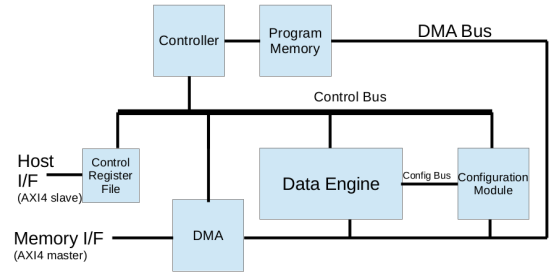


Fig. 6. Versat Topology [14]

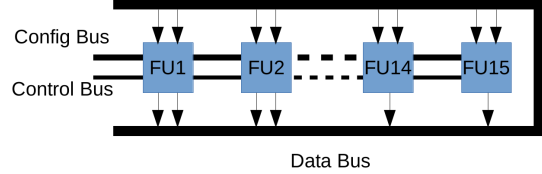


Fig. 7. Versat Data Engine Topology [15]

1) *Data Engine*: The Data Engine which is represented in Figure 7 carries out the computation needed on the data arrays. It is a 32-bit architecture with up to 11 Functional Units (FU): Arithmetic and Logic Unit (ALU), stripped down ALU (ALU-Lite), Multiplier and Accumulator (MAC) and Barrel Shifter. Depending on the project and calculations, a new type of FU or the existing ones can be altered to support the algorithm. The DE has a full mesh topology, which means that each FU can be the output to another, which leads to a decrease in operating frequency.

Each Input of a Functional Unit has a Mux with 19 entries, 8 of which are from the memories (2 from each Mem out of 4 total units) and the rest from the Functional Units (11).

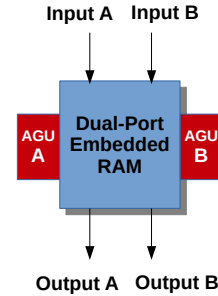


Fig. 8. Versat Memory Unit with one AGU per port [17]

The 4 Memories are dual port and for the input of both ports, there is an Address Generation Unit (AGU) that is able to reproduce two nested loops of memory indexes. The AGUs control which MEM data is the input of the FUs and where to store the results of the operation. Also, the AGUs support delayed start to line up timings due to latencies. The memory module is represented in Fig 8.

2) *Configuration Module*: Versat has several configuration spaces devised for each Functional Unit, with each space

having multiple fields to define the operation of the Functional unit (e.g. which op for the ALU). These are accessed before the run by the controller to define the datapath.

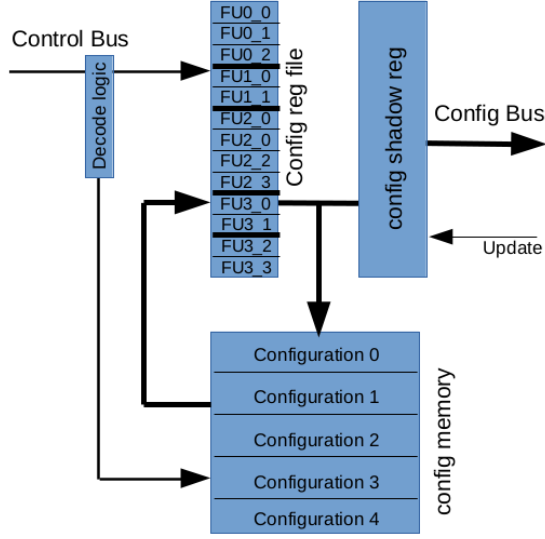


Fig. 9. Configuration Module [14]

The Configuration Module (CM), depicted in figure 9, has three components: configuration memory, variable length configuration register file and configuration shadow register. The latter holds the current configuration so the controller can change the values of the configuration file in-between runs. The decode logic finds which component to write or read, if it's the registers, it ignores read operations. Meanwhile, the configuration memory interprets both write and reads. When it receives a read, it writes into the register configuration data, when it's a write, it stores the data instead.

B. DeepVersat Architecture

The DeepVersat Architecture [18], in figure 10, decouples the Data Engine (DE) from all control and as such, it can be used with any CPU. It can be paired with hard cores in FPGA boards like the ZYNQ board with its A9 ARM dual-core CPUs or pair it with a soft core.

Its principle is to create the concept of a Versat Core: Configuration Module (CM) and its Functional Units (FU) connected with a control bus and a data bus. Instead of writing to memory, there is the option to write for the next Versat Core to create more complex and more complete Datapaths, to avoid having to reconfigure the cores.

The number of Layers and FUs are reconfigurable pre-silicon with the only limitation that each layer is identical. To program DeepVersat, an API is generated from the Verilog .vh files.

1) *DeepVersat System:* To make a complete system, a new controller is needed with a more robust toolchain. In a recent dissertation [18], the IOB-RV32 processor was used which uses the RISC-V Instruction Set (ISA) with 32-bit Integer base alongside Multiplication and Division extension

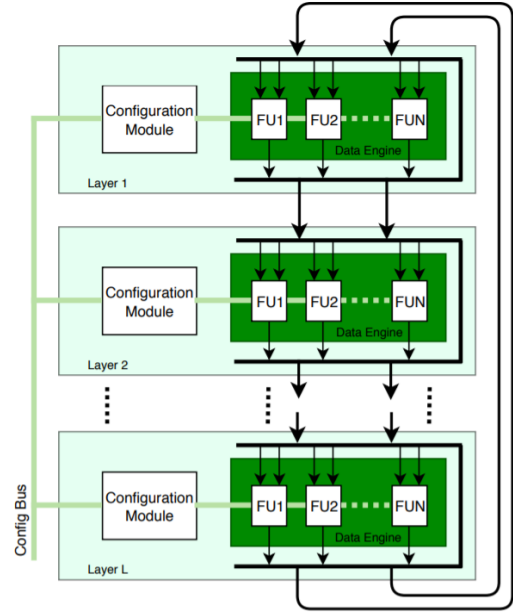


Fig. 10. DeepVersat Architecture [18]

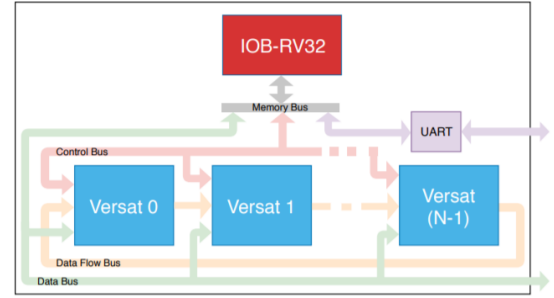


Fig. 11. DeepVersat System using a RISC-V RV32IMC soft processor [18]

and Compact Instruction extension. The core is derived from the open-source PicoRV32 CPU [19]. The IOB-RV32 uses its memory bus to access peripherals in which DeepVersat and the UART module are connected as such. The control bus is used to access the configuration modules of DeepVersat. The data bus is used to read and write a large amount of data into DeepVersat. The data flow bus is reserved for inter-Versat Core communication.

Peripheral	Memory address
UART module	12'h100xxxxx
DeepVersat control bus	8'h11xxxxxx
DeepVersat data bus	8'h12xxxxxx

TABLE II
DEEPPERSAT MEMORY MAP

The memory map to address the peripherals, including DeepVersat, is in table II. Each Versat has 15 bits of address while the CPU addresses the peripherals with 32 bits, with 8

of those occupied to choose the peripheral in question. That leaves 9 bits to address several Versat Cores which brings the theoretical maximum Versat cores to 512. The IOB-RV32 is compatible with the GNU toolchain to offer better portability of code and alongside the C++ Versat API the difficulty to code for the System diminishes.

IV. CNN COMPILING IN FPGAS

A. Toolflows for Mapping CNNs in FPGAs

Several software frameworks have been developed to accelerate development and execution of CNNs. The neural networks frameworks discussed in section II-B provides high-level APIs together with high performance execution on multi-core CPUs, GPUs, Digital Signal Processors (DSPs) and Neural Processing Units (NPU) [20]. FPGAs provide an alternative to these architectures as they provide high performance while also being low-power. FPGAs can meet several requirements including throughput and latency in the diversity of applications. Thus, several toolflows that map CNN descriptions into hardware to perform inference have been created. In table IV-A, a list of notable ones is presented.

Toolflow Name	Interface	Year
fpgaConvNet	Caffe & Torch	05/2016
DeepBurning	Caffe	06/2016
Angel-Eye	Caffe	07/2016
ALAMO	Caffe	08/2016
Haddoc2	Caffe	09/2016
DNNWeaver	Caffe	10/2016
Caffeine	Caffe	10/2016
AutoCodeGen	Proprietary Input Format	12/2016
Finn	Theano	02/2017
FP-DNN	Tensorflow	05/2017
Snowflake	Torch	05/2017
SysArrayAccel	C	06/2017
FFTCCodeGen	Proprietary Input Format	12/2017

TABLE III
CNN TO FPGA TOOLFLOWS, ADAPTED FROM [21]

1) *Supported Neural Network Models*: These toolflows support the most common layers in CNNs, which are discussed in section II-A. The acceleration target changes depending on the toolflow. For example, the fpgaConvNet [22] tool flow focuses more on feature extraction while offering nonaccelerated support for fully connected layers.

2) *Architecture & Portability*: As shown in figure 12, the fpgaConvNet architecture consists of a Front-End Parser that reads a (ConvNet) description of the network and a description of the target platform and produces, on the one hand, a Directed Acyclic Graph (DAG), which is then converted to a Synchronous Data Flow (SDF) hardware model, and on the other hand, a model of the target platform from which resource constraints are derived. The hardware model thus obtained goes into an Optimiser procedure, which produces a hardware mapping. Using hardware and software templates,

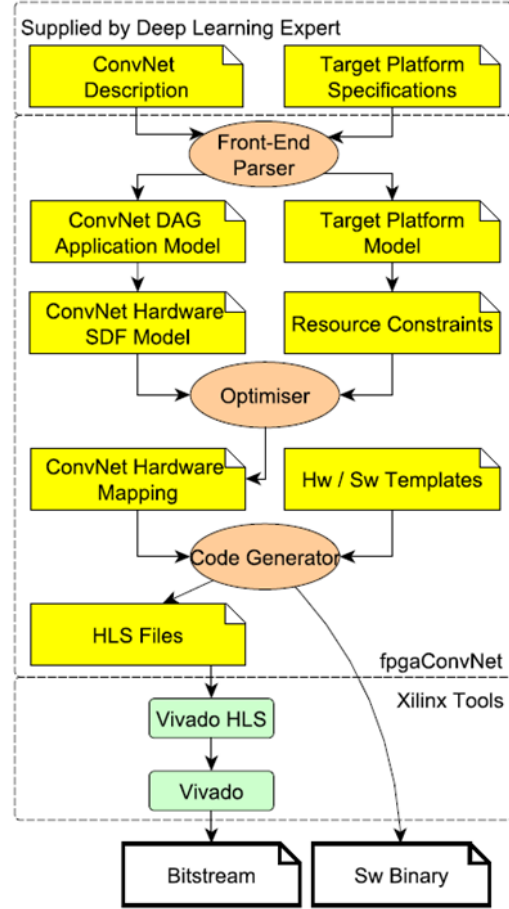


Fig. 12. fpgaConvNet Architecture. Taken from [22]

a Code Generator procedure, generates both the High Level Synthesis (HLS) input files and the software binaries that will run on the control CPU embedded in the FPGA. The HLS files go into the Xilinx (FPGA manufacturer) tools so that the configuration bitstream of the FPGA is produced.

V. DARKNET LITE

As mentioned in Section III, the DeepVersat system includes a RISC-V CPU to take out generic code and to write the configuration runs into Versat's memories. This means the first step into implementing software that can run any convolutional neural network on this system, the software must first run on the CPU then we offload Fixed Functions to Versat such as the convolutional layers, max pool, etc.

A. Porting Darknet to an embedded CPU

As mentioned in Section II-B is a framework for Neural Networks on C++ that uses dynamic memory and GPU acceleration option to get faster outputs. Also, the use of floats is prohibited in the embedded code as the RISC-V CPU only supports the extensions IM. I for Integer and M for multiplication. It also has a lot of features that are not needed in this work, such as training the CNN. By stripping

the features of Darknet we get a much simpler code framework appropriately named Darknet lite.

In the following figure, the data structure for a layer is shown. A CNN on Darknet lite is just an array of layers in which each has input, output, and layer parameters. Usually, the input is a past layer output or an image input.

```
struct layer{
    //Generic
    LAYER_TYPE type; //identifies layer's type
    ACTIVATION activation; // identifies layer's activation function
    void (*forward) (struct layer, struct network); // associated with
        forward method of each type of layer
    int groups;
    // Convolutional
    int batch_normalize; // indicates layer output must be normalized
        before applying activation function
    int batch; //always 1
    int inputs; // size of layer input
    int outputs; // size of layer output
    int h,w,c; // input dimensions
    int out_h, out_w, out_c; // output dimensions
    int n; //number of filters
    int size; // size of filter
    int stride; // indicates how many positions kernel moves
    int pad; // indicates size of padding surrounding image

    // Shortcut
    int index; //used in shortcut layer

    int classes; //used in yolo layer
    int *mask; //used in yolo layer
    int total; //used in yolo layer
    int * input_layers; //used in route layer
    int * input_sizes; //used in route layer
    fixed_t * biases; //used for convolutional and yolo layers
    fixed_t * scales; //used for convolutional layers with
        batch_normalize
    fixed_t * weights; // convolutional layer weights
    fixed_t * output; // layer output / result
    fixed_t * rolling_mean; //used for normalize_cpu
    fixed_t * rolling_variance; //used for normalize_cpu
    size_t workspace_size; // indicates max output size among all layers

    //Generic Var
    fixed_t f1; // float -> fixed 32 bit
};
```

Listing 4. Layer Struct Yolov3 [9]

By Parsing the .cfg file, a configuration file is written in C with the layer array and static position of the data for each layer. Each Layer has its definition in C to be run by the embedded CPU but for the sake of this project, several layers can be replaced by Functions that utilize Versat, the same way that the original Darknet framework had its functions written for CPU or GPU usage.

The following figure is an example of a CPU layer that computes the convolutional layer while using Fixed Point Logic.

```
void forward_convolutional_layer (layer l, network net) {
    int m = l.n; //number of filters
    int k = l.size*l.size*l.c; // filter dimensions * number of colours
    int n = l.out_w*l.out_h; //output dimension

    fixed_t *a = l.weights; //weight base address
    fixed_t *b = net.workspace; //max network's layer size
    fixed_t *c = l.output; // layer output
    fixed_t *im = net.input; // layer input
```

```
//Unroll image
if (l.size == 1) b = im;
else im2col_cpu(im, l.c, l.h, l.w, l.size, l.stride, l.pad, b);
//Perform convolution
gemm(0, 0, m, n, k, POINT, a, k, b, n, POINT, c, n);

//Normalize, scale and add bias
if (l.batch_normalize) forward_batchnorm_layer(l, net);
else add_bias(l.output, l.biases, l.batch, l.n, l.out_h*l.out_w);

//Apply activation method
activate_array (l.output, l.outputs*l.batch, l.activation);
// printf ("max=%f,min=%f\n",fixed_to_float(max),fixed_to_float(min));
}
```

Listing 5. Convolutional Layer using only CPU and fixed memory

B. Parsing CFG Files into the program

Caffe [12] is a deep learning framework as shown in section IV-A, using an open source tool [23], the output can be set to CFG. By using the network parser of Darknet, an array of layers is created with all its required parameters.

```
for (int i=0; i<net->n; i++)
{
    layer cur=net->layers[i];
    if (cur.workspace_size > workspace_size) workspace_size = cur.
        workspace_size;
    if (cur.outputs > outputs) outputs = cur.outputs;
    switch (cur.type)
    {
        case CONVOLUTIONAL:write_convolutional_IO(yoloc,i,cur,&base);
            break;
        case CONNECTED:write_connected_IO(yoloc,i,cur,&base);
            break;
        case MAXPOOL:write_maxpool_IO(yoloc,i,cur,&base);
            break;
        case DROPOUT:write_dropout_layer(yoloc,i,cur);
            break;
        case SOFTMAX:write_softmax_IO(yoloc,i,cur,&base);
            break;
        case AVGPOOL:write_avgpool_layer(yoloc,i,cur);
            break;
        case SHORTCUT:write_shortcut_IO(yoloc,i,cur,&base);
            break;
        case ROUTE:write_route_IO(yoloc,i,cur,&base);
            break;
        case RNN:write_rnn_layer(yoloc,i,cur);
            break;
        case YOLO:write_yolo_IO(yoloc,i,cur,&base);
            break;
        case UPSAMPLE:write_upsample_IO(yoloc,i,cur,&base);
            break;
        // Other layers needed to be addressed
        case GRU:
        case CROP:
        case REGION:
        case DETECTION:
        case COST:
        default : printf ("\nThis layer is not yet supported.Bye!\n");
            exit (0);
            break;
    }
}
```

Listing 6. For Loop for writing darknet layers

Afterward, by going through each layer, "yolo.c" will be written with all the data darknet lite will need. In listing 7,

the addresses of the data needed for the layer. In 8, the static parameters are defined as well.

```
/*Layer 2-CONVOLUTIONAL*/
#define FOUTPUT_2 BASE+3461616
#define FSCALES_2 BASE+4846064
#define FR_MEAN_2 BASE+4846096
#define FR_VARIANCE_2 BASE+4846128
#define FWEIGHTS_2 BASE+4846160
#define FBIASES_2 BASE+4850768
```

Listing 7. For Loop for writing darknet layers

```
/*GENERIC PARAMS-Layer 2*/
[2].type=0, [2].activation=7, [2].batch_normalize=1, [2].batch=1,
[2].inputs=692224, [2].outputs=1384448, [2].n=32,
[2].h=208, [2].w=208, [2].c=16,
[2].out_h=208, [2].out_w=208, [2].out_c=32,
[2].size=3, [2].stride=1, [2].pad=1,
[2].index=0, [2].classes=0, [2].total=0,
```

Listing 8. For Loop for writing darknet layers

VI. DEEPVERSAT SOFTWARE SIMULATOR

The need for a software simulator comes from the complexity of the configurations being written into Versat and the hardware simulation faults of taking too much time and being hard to debug.

The goal is to emulate what the hardware is doing much more efficiently than a simple Hardware simulation as the time of development for hardware is much higher than for simple software. The Simulator executes clock iteration per iteration getting the same results in each clock as the hardware. As Versat is a CGRA, different functional unit configurations are easy to accomplish in the simulator and the time to get results on performance for a specific program is a lot faster.

A. Architecture and Object Relation

The Simulator is made up of the Parent Class called Versat, which will be simulated itself, as each Versat instance is independent of the others, the simulations are also independent. The Versat is made up of 2 CStage Arrays, one is the "live" while the other is the shadow registers, where the configurations are held before the simulator is run. Each Stage is made up of its Functional Units, of which each is connected to the Databus. As it happens in the hardware, functional units can access the database which has the output of the current stage and the previous one.

1) *Functional Units*: The following table contains the functional units present in the simulator and is represented by "CFU" in figure 13. VI and VO represent CRead and CWrite classes respectively.

To add a new FU, it's as easy as creating a new class that will be used by CStage with a run(), update(), output(), and copy() method. Of course, if it has variables needed to be defined by the program, set param functions are also needed. Using the simulator, hardware development and program development can be parallelized to output a new program with more optimized performance.

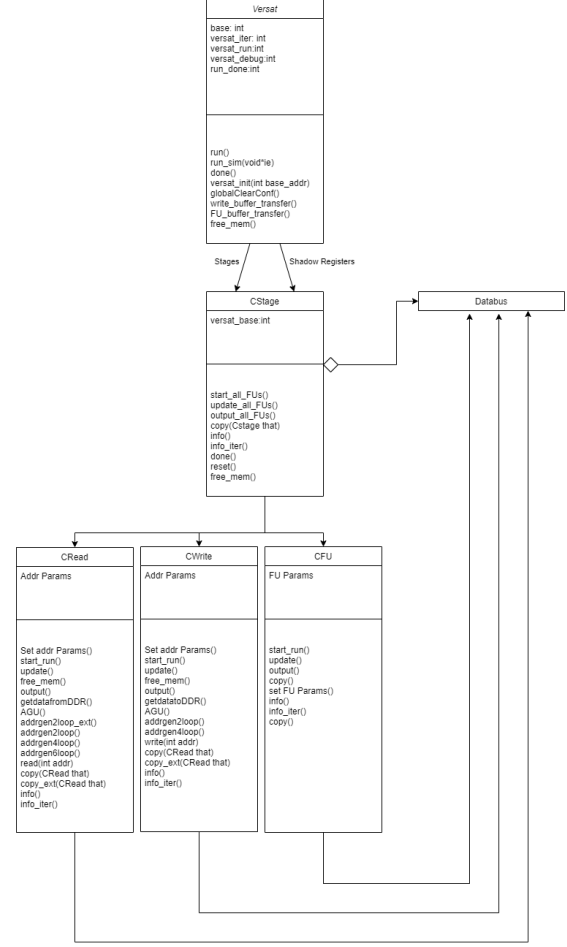


Fig. 13. Class Structure for the Versat Simulator

Functional Unit	Purpose
Read (VI) Mem Unit	Reads from DDR and sends Data to databus
Write (VO) Mem Unit	Reads from databus and sends Data to DDR
MulAdd (MAC)	Multiplication and Accumulate
Mul	Multiplication
Alu	Standard algorithmic and logic unit
AluLite	Stripped down algorithmic and logic unit
Barrel Shifter (BS)	Shifts to the right or to the left
Memory (Mem)	Sends/Receives data to/from the pipeline. Data is inserted through CPU communication

TABLE IV
VERSAT SIMULATOR FUNCTIONAL UNITS

In the next section, these methods will be explained in detail and their importance to the simulator.

B. Simulation

After the program that is running on the CPU finishes writing the configurations, it will call the run method of Versat. In figure 14, a sequence diagram is presented with the rundown of a typical program that uses Versat Simulator.

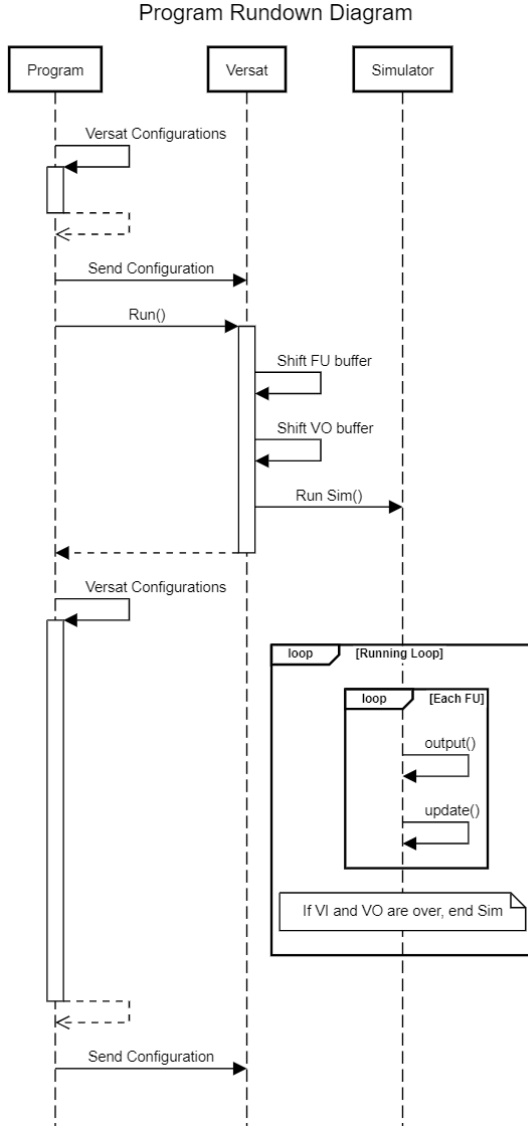


Fig. 14. Sequence Diagram of a Program using Versat Simulator

1) *Run() Function:* In the software API for Embedded Versat, the run function would write to a shadow register, which we can call "start" changing the value from 0 to 1. Similarly, another register would change the value to 0, which we can call "done". While this last register isn't turned to 1, Versat hasn't finished running with the previous configurations so all that can be done is to write configurations for future runs.

In the simulator, it works in a similar way to preserve compatibility as the goal is to have the same programs run on software simulators and the FPGA.

```

void CVersat::run()
{
    //MEMSET(base, (RUN_DONE), 1);
    run_done = 0;
    versat_iter = 0;

    //update shadow register with current configuration
    #if nVO > 0
        write_buffer_transfer();
    #endif
    #if nVI > 0
        FU_buffer_transfer();
    #else
        int i = 0;

        for (i = 0; i < nSTAGE; i++)
        {
            stage[i].reset();
            shadow_reg[i].copy(stage[i]);
        }
    #endif

    pthread_create(&t, NULL, run_simulator, (void*) this);
}

```

Listing 9. The Run function code

As we can see in the previous listing, we reset the state variables of the simulator, then shift the VO and FU shadow registers. This is done to simulate the pipeline delay in the FPGA. Because the data needs to come and go to the main memory (DDR), 1 run cycle is used just for fetching data and writing data. Using a small example: If a developer writes a configuration to do a 5x5 matrix multiplication, Versat will have to run 3 times. Once to fetch data from memory, the second for the actual use of Versat and the final one is to get data onto memory.

In the simulator, this is done using the same class instances and copying the configuration values. On the hardware, it's several flip-flop registers in a row. However, all these 3 stages can happen at once if you run multiple configurations in one program, e.g.: running a CNN through Versat, will have at least 1 run per layer. So, if it has 5 layers, Versat will have to run 5+2 times, the last 2 times are done to flush the Versat of any data.

After the shift, a new thread is created to run the simulator in parallel with the configurations, having the same behavior as the hardware.

2) *Start() Method:* At the beginning of the configuration run, the method "start run" of all FUs and memories are started. In this function, several functional units will have their state variables reset such as VI, VO and MAC FU.

3) *Databus:* The databus on Versat is a simple array that holds all the outputs of the functional units. The data type (versat_t) of the array depends on the width of Versat, which is part of the configuration file. Using higher width, e.g.: 64 bits, is useful for the same instruction, multiple data (SIMD) applications but requires the functional units to be adapted. For the purpose of this thesis, 16 bits and 32 bits are used depending on the neural network and how it is optimized.

When the Versat is instantiated in the program, the functional units constructor will point to the correct position of the databus as it's referenced in the following figure.

As mentioned in figure 10 from section III, each functional unit will be able to access the output from the functional units of the current stage and previous. Software-wise, each stage will be pointing to a part of the databus.

4) *Update() and Output() Method*: The update method's goal is to update the functional unit's value on the databus. Each functional unit has a pipeline delay to output or has a run delay configured, like the memories or MAC.

Meanwhile, the output method's goal is to, based on the inputs from the databus, calculate the result from the functional Unit.

For a compute functional unit such as the MAC or the ALU, this means reading from the databus for operands A and B and performing the selected operation. For the read memory (VI), it will output an address on the mem and performs a read operation. For the write memory, it will output an address and performs a write operation.

In the listing 11, the code of the Mul functional unit is used as an example.

```
void CMul::update()
{
    int i = 0;

    //update databus
    databus[sMUL[mul_base]] = output_buff[MUL_LAT - 1];
    // special case for stage 0
    if (versat_base == 0)
    {
        //2nd copy at the end of global databus
        global_databus[nSTAGE * (1 << (N_W - 1)) + sMUL[mul_base]] =
            output_buff[MUL_LAT - 1];
    }

    // trickle down all outputs in buffer
    for (i = 1; i < MUL_LAT; i++)
    {
        output_buff[i] = output_buff[i - 1];
    }
    // insert new output
    output_buff[0] = out;
}

versat_t CMul::output()
{
    // select inputs
    opa = databus[sela];
    opb = databus[selb];

    mul_t result_mult = opa * opb;
    if (fns == MUL_HI)
    {
        result_mult = result_mult << 1;
        out = (versat_t)(result_mult >> (sizeof(versat_t) * 8));
    }
    else if (fns == MUL_DIV2_HI)
    {
        out = (versat_t)(result_mult >> (sizeof(versat_t) * 8));
    }
    else // MUL_LO
    {
        out = (versat_t)result_mult;
    }

    return out;
}
```

Listing 10. Update and Output method of Mul

5) *Copy() and Info() Method*: Finally, the last 2 functions of the simulator, copy() the main purpose is to copy the configuration parameters from one instance to another, used mostly at the beginning of the run to simulate the shadow registers. Meanwhile, the Info method is a State printing function that outputs a string with the full data of the current iteration, this way, you can check iteration by iteration the progress of the simulation, just like in hardware. At this moment, if debugging is activated, each clock iteration output and the state of Versat will be in a file.

```
mul_add[0]
OpA=   -7
OpB=   -3
SelA=    1
SelB=    0
Addr=    3
Finished= 0
Out=    61
OUTPUT_BUFFER (LATENCY SIM)
Output[0]=61
Output[1]=40
Output[2]=60
Output[3]=45
```

Listing 11. Info output for the MAC functional unit

VII. VERSAT API 2.0

The Versat API, developed in a previous thesis [18], has the ability to conceal the calls to the hardware to avoid changing the program when the hardware changes.

In this chapter, the new functions that are part of the Versat API will be discussed. The goal is to make development for Versat just like writing normal code and to be easy to port code the same way CUDA has done the same to run SIMD code on Nvidia GPUs.

```
class CMemPort
{
public:
    int versat_base, mem_base, data_base;

    // Default constructor
    CMemPort()
    {
    }

    // Constructor with an associated base
    CMemPort(int versat_base, int i, int offset)
    {
        this->versat_base = versat_base;
        this->mem_base = CONF_BASE + CONF_MEM0A + (2 * i + offset) *
            MEMP_CONF_OFFSET;
        this->data_base = (i << MEM_ADDR_W);
    }

    // Methods to set config parameters
    void setIter (int iter)
    {
        MEMSET(versat_base, (this->mem_base + MEMP_CONF_ITER), iter);
    }
    void setPer (int per)
```

Listing 12. Sample Versat API implementation for the Hardware for Mem functional unit

A. API Architecture

In figure 15, a graphic representation of the new API is presented. It has 4 apparent layers (5 if you count the hardware):

- 1) Complex Mathematical API that is automatically optimized for the Versat Setup you chose. No dev work required
- 2) Read/Write using VI and VO for simpler setup of the data. Also includes easier FU functions to set up workloads.
- 3) Read/Write configurations for inside Versat Data (Int) or DDR to/from VI/VO (Ext).
- 4) Versat API 1.0 where each configuration variable needs to be set up individually
- 5) No API. Hardware registers where the values are used inside Versat.

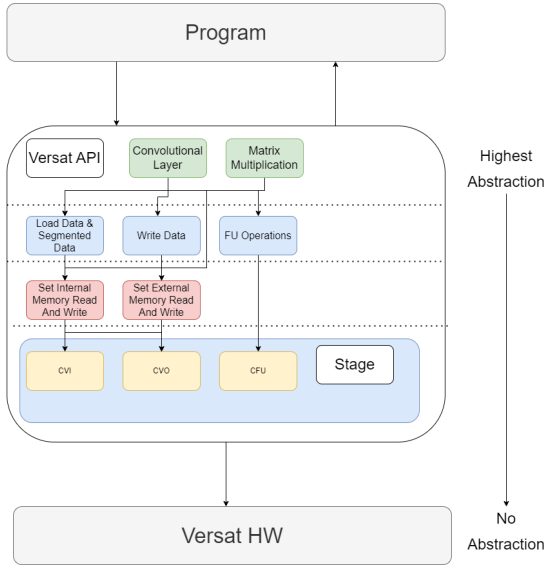


Fig. 15. Graphic representation of the new Versat API and its connections

B. Memory Operations API

When utilizing the VI instead of a MEM, the data transfer happens between the functional unit and direct memory access while on the mem, the CPU writes directly to Versat, wasting CPU cycles. For the API, this means going from a read method that is straightforward to more configuration methods to set up the read operation from DDR. The same happens to Write operations. To address this, 7 functions were created in 2 levels of abstraction: `load_data()`, `load_segmented_data()`, `write_data()` that use a lower level functions: `set_IntMem_Write()`, `set_ExtMem_Write()`, `set_IntMem_Read()` and `set_ExtMem_Read()`. The function of the higher abstraction memory functions is to abstract the parameters of the AGU. In the following listing, we have one of the implementations as an example.

```
int convolutional_layer_xyz (CVersat* Versat, int input_addr, int channels, int
height, int width, int kernel_size, int stride, int pad, int weights_addr,
int output_addr, int nkernels)
{
    Accumulator weights = Accumulator();
    Accumulator input[nOUTPUTS];
    Accumulator write_matrix[nOUTPUTS];

    int out_w=((width + 2*pad - kernel_size) / stride ) + 1;
    int out_h=((height + 2*pad - kernel_size) / stride ) + 1;
    int in_w=width+pad;
```

Listing 13. Load Segmented Data code

Although this means having to write code with the AGUs in mind and how they function. To avoid it, a new class was created, shown in listing ??listing:accumulator, to also abstract how the AGU counts loops and approximate the code to simple C++ code that runs on a CPU.

```
Accumulator()
{
    iter =per= shift =incr= iter2=per2= shift2=incr2= iter3=per3= shift3=incr3=nloops=
    delay=duty= start=extAddr=intAddr=0;
}
void add_loop(int per, int incr = 1 )
{
    switch(nloops)
    {
        case 5: this->iter3=this->per3;
                this->shift3=this->incr3;
        case 4: this->per3=iter2;
                this->incr3=shift2;
                if ( this->iter3==0)
                    this->iter3=1;
        case 3: this->iter2=this->per2;
                this->shift2=this->incr2;
        case 2: this->per2=iter;
                this->incr2=shift;
                if ( this->iter2==0)
                    this->iter2=1;
        case 1: this->iter=this->per;
                this->shift=this->incr;
        case 0: this->per=per;
                this->incr=incr;
                if ( this->iter==0)
                    this->iter=1;
        default : nloops++;
    }
}

void loop_settings (int start = 0, int duty = 0, int delay = 0, int extAddr = 0,
int intAddr = 0)
{
    this->duty=duty;
    this->start= start ;
    this->delay=delay;
    this->extAddr=extAddr;
    this->intAddr=intAddr;
}
};
```

Listing 14. Accumulator Class code

To transform from AGU parameters to for loop, it depends on the number of loops pretended to be done. VI AGU is 3 cascade Accumulators and as such, the increment on the second and third accumulators needs to be adjusted, just as shown below.

```
switch(loop.nloops)
{
```

```

case 6:
case 5: loop.incr2 += loop.shift * loop.iter + (loop.incr * loop.per) * loop.iter ; //
      4 + 2*2 = 8
      loop.incr3 += loop.shift2 * loop.iter2 + (loop.incr2 * loop.per2) * loop.iter2 ;
      break;
case 4:
case 3: loop.incr2 += loop.shift * loop.iter + (loop.incr * loop.per) * loop.iter ;
      default : break;
}

```

Listing 15. AGU parameters to Simple forloop parameters transform

C. Matrix Multiplication and Dot Product

As part of the new API, a matrix multiplication function was added. In listing 16 the code is presented. First, 2 Accumulator class variables are initialized. Afterward, using the 2 arrays address in DDR, the AGU configurations of the VIs to read from the main memory are set, then the AGU configurations of VI for the data handling inside the Data Engine. Finally, the function that will write the configuration of a MAC and the store AGU configurations. This last step is optional as the result of this matrix multiplication can be used in the same run to make other operations e.g.: adding a bias using one of the ALUs to the results.

```

int matrix_mult(CStage* Versat, int matrix_a, int matrix_b, int result_matrix, int
r_a, int c_a, int r_b, int c_b, bool store)
{
    Accumulator store_matrix_A = Accumulator();
    Accumulator store_matrix_B = Accumulator();

    // Send Data from DDR to Versat Memory
    store_matrix_A.add_loop(r_a * c_a);
    store_matrix_A.loop_settings(0, 0, 0, matrix_a, 0);
    store_matrix_B.add_loop(r_b * c_b);
    store_matrix_B.loop_settings(0, 0, 0, matrix_b, 0);

    set_ExtMem_Read(Versat, 0, store_matrix_A);
    set_ExtMem_Read(Versat, 1, store_matrix_B);

    Accumulator read_matrix_A = Accumulator();
    Accumulator read_matrix_B = Accumulator();

    // Read from Matrix A in Versat
    read_matrix_A.add_loop(r_a * c_a);
    read_matrix_A.add_loop(c_b - c_a);
    read_matrix_A.add_loop(c_a, 1);
    read_matrix_A.loop_settings(0);
    set_IntMem_Read(Versat, 0, read_matrix_A);

    // Read from Matrix B in Versat
    read_matrix_B.add_loop(r_a - c_b);
    read_matrix_B.add_loop(c_b - r_b * c_b + 1);
    read_matrix_B.add_loop(r_b, c_b);
    read_matrix_B.loop_settings(0);
    set_IntMem_Read(Versat, 1, read_matrix_B);

    // Do multiplication of the values and accumulate.
    muladd_operation(Versat, sVI[0], sVI[1], 0, MULADD_MACC, r_a * c_b, c_a,
MEMP_LAT, 0);

    // Store the results in Memory.
    if (store == true)
    {
        Accumulator write_matrix = Accumulator();
        write_matrix.add_loop(r_a * c_b, 1);
        write_matrix.loop_settings(0, 0, MEMP_LAT + MULADD_LAT, result_matrix, 0);
        set_ExtMem_Write(Versat, 0, write_matrix);
        write_matrix.add_loop(c_a, 0);
        set_IntMem_Write(Versat, 0, write_matrix, sMULADD[0]);
    }
}

```

```

return sMULADD[0];
}

```

Listing 16. Matrix Multiplication Configurations

The Dot product function is very similar, the configurations are identical for the data transfer from the main memory to the VIs. In the inside loops of the VIs, instead of 3 loops, we only need to use 1.

D. Generic Convolution

As explained in chapter ??chapter:Background, convolutional neural networks are a type of neural nets that are used mostly in image and object recognition by using convolutional layers. To run a convolutional layer on Versat with optimized performance, the configurations must be written with regard to several parameters:

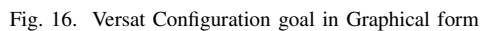
- 1) Memory Sizes used in VI and VO. The amount of data that can be stored at once. It determines the number of outputs done per run.
- 2) Functional Units used in the Data Engine. Here it's about the lowest common denominator, i.e. the bottleneck in the Data Engine determines the number of outputs done simultaneously.

This function has a total of 20 variables calculated at the start before the Versat configurations are written. The most important variables are the following:

- output height (h) and width (w) of the resulting matrix from the convolution.
- Number of outputs done simultaneously, also known as pipeline width (nOutputs). This value is pre-compiled as it depends on only Versat Configurations.
- Number of outputs that can be done per VI (y) in a single run and its variations. Outputs total (y₂), Output Lines per VI (y₃) Output Lines total (y₄). The value of y₄ and y₂ decide the different configuration scenarios.
- Resource Allocation Variables which are explained in subsection ??ConvolutionScenarios
- Address Variables
- AGU Configuration Variables

The hard part of the algorithm is to allocate the data in the most efficient way possible and to create the AGU configurations for the VIs and VO. For this algorithm, the CGRA will act like a GPU pipeline where several "threads" will exist that will output 1 point every k² cycles, where k is the kernel size used in the convolution.

1) *Loading Data:* Usually, when doing a convolution in CPU, the frameworks transform the convolution to a matrix multiplication by creating a new matrix that will multiply with a kernel vector. It's done this way as matrix multiply is a heavily optimized operation and can take advantage of a CPU's SIMD units or even call the GPU APIs and offset the workload there. On Versat, this is not needed as to calculate 1 output, we will need only enough space in mem to hold k²*ch where ch is the channels of the input. And as such, it means 9216 bytes per VI at least for YoloV3 CNN when using 16-bit operands.



On the code, this takes form in 1 single line, thus the importance of the previously written functions.

Listing 17. Load Input Matrix into VIs

$$size = w * (k + stride * (iter - 1))$$

2) *Convolution Scenarios:* When writing the configurations of the convolution runs, there are several cases the software needs to take into consideration. As explained in the previous

Figure 1 displays four 10x10 grids representing different output matrices for Case 1, Case 2, Case 3, and Case 4. The grids are labeled as follows:

- Output Matrix - Case 1:** A solid blue 10x10 grid.
- Output Matrix - Case 2:** A blue 10x10 grid with the bottom row (row 10) white.
- Output Matrix - Case 3:** A blue 10x10 grid with the top row (row 1) white.
- Output Matrix - Case 4:** A blue 10x10 grid with the top row (row 1) white.

Where y_2 is the blue squares

Fig. 17. Convolution Scenarios that Versat will have

In figure 18, the flowchart of each case is presented.

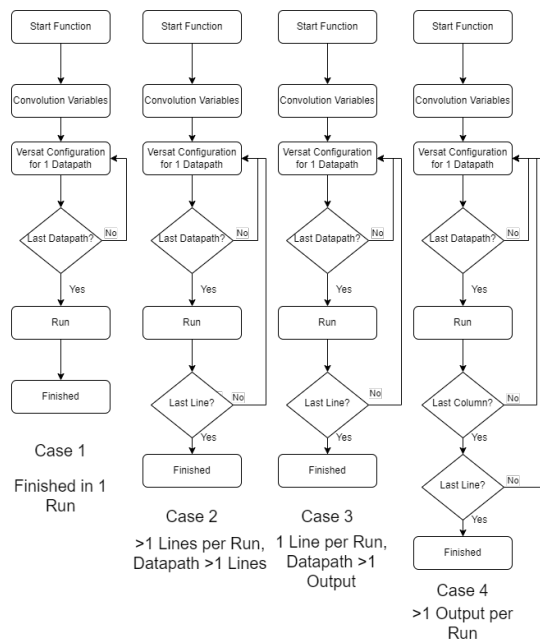


Fig. 18. Configuration Flowchart for the different scenarios

And on listing 18 the AGU configurations of the VIs that hold the input matrix, the MAC configuration, and finally the VO AGU configuration.

```
// h1 and h2
```



```

input[i]= Acumulator();
input[i].add_loop(nkernels,-in_w*((num_iter)* stride ));
input[i].add_loop(num_iter,(in_w*stride)- stride *out_w);
input[i].add_loop(out_w,-channels*size_per_channel+ stride );
input[i].add_loop(channels, size_per_channel-line_plus_one*
kernel_size+rewind_kernel);
input[i].add_loop(kernel_size , line_plus_one );
input[i].add_loop(kernel_size ,1);
input[i]. loop_settings (0);
set_IntMem_Read(stage,i+1,input[i]);

aux--;

muladd_operation( stage ,sVI[i+1],sVI[0],i,MULADD_MACC,(num_iter)*
out_w*nkernels,
kernel_size * kernel_size *channels,MEMP_LAT,0);

write_matrix[i] = Acumulator();
write_matrix[i].add_loop(nkernels ,out_w*(out_h+1)-(num_iter)*out_w);
write_matrix[i].add_loop((num_iter)*out_w,1);
write_matrix[i]. loop_settings (0,0,MEMP_LAT+MULADD_LAT,
output_addr_new,0);
set_ExtMem_Write(stage,i, write_matrix[i]);
write_matrix[i] = Acumulator();
write_matrix[i].add_loop(nkernels ,0);
write_matrix[i].add_loop((num_iter)*out_w,1);
write_matrix[i]. loop_settings (0,0,MEMP_LAT+MULADD_LAT,
output_addr_new,0);
write_matrix[i].add_loop(kernel_size * kernel_size *channels ,0);
set_IntMem_Write(stage,i, write_matrix[i],sMULADD[i]);

input_addr_new+=(in_w*(stride*(h1+aux_bool)))*(DATAPATH_W/8);
output_addr_new=output_addr_new+(h1+aux_bool)*out_w*(DATAPATH_W
/8);

```

Listing 18. Versat configurations for 1 datapath

VIII. RESULTS

In this chapter, the results from the simulator are presented and the new functions developed for Versat. In Section VIII-A, a setup used by a testbench used by DeepVersat is used to prove the simulator is working according to predictions. Afterward, in section VIII-B, a testbench is prepared to test the matrix multiplication and generic convolution with several hardware configurations available to test the performance of different scenarios in the simulator.

The tests were executed on a 64-bit machine, with an AMD Ryzen 7 5800H Processor and 16GB of RAM running Windows 11, version 22H2, WSL 2.0 with the image of Ubuntu 20.04. The compiler used is g++ version 9.4.0.

A. Simulator Testing

To test the simulator, a testbench was created that will create a random input matrix of 5x5 with a kernel size of 3. For each Stage defined in the headers file, a channel will be added and the result of the convolution will propagate through the stages.

To be more specific in the beginning, the configurations of the VIs are written to transfer the data from the program to Versat. The data uses the rand() function with seed using current time so the result is different every time. Both the input matrix and kernel map are randomized. The former value varies from -25 to 25 while the kernel varies from -5 to 5. Using the data, we calculate the result of the convolution in the CPU. Afterward, the configuration for the Bias mem is done and then stage by stage the configuration of the VI, MAC, and ALU is done. Finally, the configuration of the VO is written.

The estimated iterations needed are the following:

```

VERSAT TEST
Deep versat initialized in 974433 us
Data stored in versat mem in 2 us
Expected result of 3D convolution
88 -23 -12
-129 -272 -129
73 91 73
3D CONVOLUTION WITH 4-LOOP ADDRESS
Configurations (except start) made in 2 us
3D convolution with 4-LOOP ADDRESS
Configurations (except start) made in 8 us
Expected Versat Clock Cycles for this run 98
3D convolution done in 981 us
Simulation took 98 Versat Clock Cycles
Actual convolution result
88 -23 -12
-129 -272 -129
73 91 73

```

Fig. 19. Simulator test output in terminal

$$Est = Delay + Iter_2 * Per_2 * Iter_1 * Per_1$$

Where these are the AGU configurations of the VO where the results are written. The Delay is accumulated through the several stages by adding two due to the MACs and ALUs.

B. Testing the new API

In this section, the same method for the previous testbench is made. While the previous one relies on using API v1 for the configuration, these test benches run the new API.

1) *Testbench for Matrix Multiplication:* The Matrix Multiplication is a quite simple program. The only thing needed is an instance Versat, run versat_init(), create the matrixes, and then use the function matrix_multiplication() The data is also computed in the CPU the result to verify the output.

2) *Testbench for Generic Convolution:* Using the same method on the previous test benches, the following Convolution Layer was used with several Versat Configurations.

CNN Variable	Value
Kernel Size	2
Channels	2
Number of Kernels	2
Input Height	12
Input Width	12
Stride	1
Out Width	11
Out Height	11
Out Channels	2

TABLE V
CNN LAYER ON THE TESTBENCH

For this specific Versat hardware configuration the number of iterations needed is 711 using 3 Datapaths.

In Table VIII-B2, the different Datapath numbers and how it affects performance. A datapath is a combination of 1 VI, 1 MAC, and 1 VO. So the lower number in the Versat configuration file decides the number of valid datapaths, of course, VI needs +1 in numbers more than the functional units due to the Kernel memory.

The reason for these results is quite simple. In total, 11 output lines are divided by the datapaths. When the division is not a whole number, the remainder gets distributed by

Number of Datapaths	Iterations
1	1943
2	1063
3	711
4	535
6	359
8	359
11	183
16	183
22	183

TABLE VI

CNN LAYER ON THE TESTBENCH WITH SEVERAL VERSAT HARDWARE CONFIGURATIONS

available datapaths. The Consequence of this, when changing from 6 to 8 datapaths, the performance doesn't get any better. Datapath 0 will have to run twice to (2 lines) while Datapath 8 will run 1 line. To increase further the performance, the output channels would have to be divided through more datapaths.

IX. CONCLUSIONS

In this paper, several modules and tools were presented for neural networking development on Versat. The simulator is a significant improvement over normal hardware simulation for testing new software configurations and workloads. It also can predict the performance of the workloads and helps size how many functional units, stages, or how much the size of memories should be to achieve the highest performance. Furthermore, the new Versat API can bring new tools for development using the hardware and be able to write code for Versat akin to writing normal C++ code that runs on a CPU. Finally, the tools designed for Darknet give embedded development a boost by being able to run any CNN on embedded hardware, even if it's just a CPU. To test this, a suite of programs was planned and the results show the new software tools effectiveness.

A. Achievements

One achievement of this paper was the development of the simulator. The simulator is able to successfully emulate the output of the hardware, where a new program written for Versat can be tested in 5 seconds instead of several minutes to put the program into the FPGA. Another achievement is the generic convolution being able to run any type of convolution layer efficiently. By changing the Versat parameters, a new complete hardware configuration can be done and tested with the software to check new performance figures. Lastly, the darknet framework for embedded devices and the tools used to parse CFG are important for future work using the Versat CGRA.

REFERENCES

- [1] Gualtiero Piccinini. The first computational theory of mind and brain: A close look at mcculloch and pitts's "logical calculus of ideas immanent in nervous activity". *Synthese*, 141, 08 2004.
- [2] Yann LeCun, Y. Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521:436–44, 05 2015.

- [3] Moshe Leshno, Vladimir Ya. Lin, Allan Pinkus, and Shimon Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, 6(6):861 – 867, 1993.
- [4] mnist database of hand-written digits.
- [5] Masakazu Tanomoto, Shinya Takamaeda-Yamazaki, Jun Yao, and Yasuhiko Nakashima. A cgra-based approach for accelerating convolutional neural networks. pages 73–80, 09 2015.
- [6] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems* 25, pages 1097–1105. 2012.
- [7] Max-pooling / pooling.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- [9] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement, 2018.
- [10] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- [11] Joseph Redmon. Darknet: Open source neural networks in c. <http://pjreddie.com/Darknet/>, 2013–2016.
- [12] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross B. Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *CoRR*, abs/1408.5093, 2014.
- [13] Rui Santiago, João D. Lopes, and José T. de Sousa. Compiler for the versat reconfigurable architecture. REC 2017, 2017.
- [14] João D. Lopes, Rui Santiago, and José T. de Sousa. Versat, a runtime partially reconfigurable coarse-grain reconfigurable array using a programmable controller. *Jornadas Sarteeco*, 2016.
- [15] João D. Lopes and José T. de Sousa. Fast fourier transform on the versat cgra. *Jornadas Sarteeco*, 09 2017.
- [16] João D. Lopes and José T. de Sousa. Versat, a minimal coarse-grain reconfigurable array. In Dutra I., Camacho R., Barbosa J., and Marques O., editors, *High Performance Computing for Computational Science – VECPAR 2016*, pages 174–187. Springer, 2016. doi:10.1007/978-3-319-61982-8_17.
- [17] João Dias Lopes. Versat, a compile-friendly reconfigurable processor – architecture. Master's thesis, Instituto Superior Técnico, November 2017.
- [18] Valter Jorge Brás Mário. Deepversat: A deep coarse grain reconfigurable array. Master's thesis, Instituto Superior Técnico, November 2019.
- [19] Picorv32- a size-optimized risc-v cpu.
- [20] Andrey Ignatov, Radu Timofte, Przemyslaw Szczepaniak, William Chou, Ke Wang, Max Wu, Tim Hartley, and Luc Van Gool. Ai benchmark: Running deep neural networks on android smartphones, 10 2018.
- [21] Stylianos I. Venieris, Alexandros Kouris, and Christos-Savvas Bouganis. Toolflows for mapping convolutional neural networks on fpgas: A survey and future directions, 2018.
- [22] Stylianos I. Venieris and Christos-Savvas Bouganis. fpgaConvNet: A Framework for Mapping Convolutional Neural Networks on FPGAs. In *2016 IEEE 24th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pages 40–47. Institute of Electrical and Electronics Engineers (IEEE), May 2016.
- [23] Caffe2darknet python tool.