



Facultad
Informática
Universidad
Murcia



University of Murcia
Computer Science Faculty

MSc on New Technologies in Computer Science

**Towards a unified programming model
for diverse computing architectures:
Experiences using PHAST**

Author:

Eduardo José Gómez-Hernández
{eduardojose.gomez@um.es}

Advisor:

José Manuel García
{jmgarcia@um.es}

June, 2019

This page has been intentionally left blank

Declaration of Authorship

I, Eduardo José Gómez Hernández, declare that this master's thesis, titled "Towards a unified programming model for diverse computing architectures: Experiences using PHAST", and the work presented in it are mine.

I confirm that:

- **This work was done wholly or mainly while in candidature for a master's degree at this University.**
- **Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.**
- **Where I have consulted the published work of others, this is always clearly attributed.**
- **Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this dissertation is entirely my own work.**

Eduardo José Gómez Hernández
Murcia, June 2019

Contents

Abstract	iii
Resumen	iv
List of Figures	vi
List of Tables	vii
List of Listings	viii
1 Introduction	1
2 Background	2
2.1 CPU Revolution	2
2.2 Accelerator Revolution	2
2.2.1 GPU & GPGPU	2
2.2.2 FPGA	2
2.2.3 ASIC	2
2.2.4 The problem	2
2.3 ML & DNNs	2
2.3.1 Development and usage	2
2.3.2 Evolution with HPC	3
3 An overview of high performance parallel programming models	4
3.1 Performance Scaling Problems	4
3.2 High Performance Programming Models	4
3.3 Some High Performance Programming Models	4
3.3.1 OpenMP	4
3.3.2 OpenACC	5
3.3.3 OmpSs	5
3.3.4 OpenCL & SYCL	5
3.3.5 Kokkos	5
3.3.6 PHAST	5
3.4 Classification	6
3.4.1 Platforms	6
3.4.2 Type	6
3.4.3 Programming Paradigm	6
3.4.4 Memory Paradigm	7
3.4.5 Table	7
3.5 Portable Programming Models	9
3.5.1 Classification	9
4 PHAST applied to the Caffe Framework	10
4.1 Why PHAST	10
4.2 The Caffe Framework	10
4.3 Development	11
4.3.1 Blob	11
4.3.2 Convolution	12
4.3.3 Pooling	15

4.3.4	InnerProduct	16
4.3.5	ReLU	18
4.3.6	Softmax and Softmax Loss	19
4.3.7	Accuracy	21
4.3.8	Pending work	21
5	Results & Evaluation	22
5.1	Testbench	22
5.2	Experiences	22
5.3	Case Tests	22
5.4	Performance	23
6	Conclusions & Future Work	24
	Acknowledgments	24
	References	25
	Acronyms	27

Abstract

We are in an era that is reaching the limit of the multi- and many-core performance. Since the start of computing, the performance was always a limit to break, in the first place by scaling the frequency of each CPU generation, then Moore's Law forcing the increase of the transistors in an integrated circuit.

Nowadays, many tasks are delegated to specific devices capable to execute those tasks faster and more efficient than a traditional CPU, being the most known one the GPU. We have reached an accelerator boom, with the usage of FPGAs as accelerators instead of only prototyping and the creation of many ASICs for specific tasks, there are a vast of accelerators that can be used.

The problem of this accelerator boom is that to exploit the performance of each one, it is necessary to use its own programming language, a DSL. Therefore, to use multiple devices, programmers are forced to create multiple versions of the same source code for each device, increasing a lot the code complexity.

Machine learning is one of the fields where new accelerators are developed each day because being able to run more complex models in less time allows them to solve bigger problems.

To accomplish this big problem of multiple device programming, several programming models are developed. Therefore, we want to know which programming models are available, which is the main focus of each one, and their characteristics. There is not any complete classification of these high performance programming models, only an interesting survey from 2013, but many of the models classified are now deprecated.

We have developed a classification of 23 high performance models using 4 parameters, supported platforms, type of programming model, the programming paradigm and the memory paradigm. These models were selected from 50 reviewed models, but we discard discontinued models and the ones that have been merged into other models.

From this classification, we have extracted 8 models that can be fit into a single source programming model classification, it is only a sample of all models in the table that have this characteristic.

Then, we think that use one of these new models into a real application is interesting, therefore we selected PHAST, and as a joint collaboration with its creators, we have been porting the Caffe deep learning framework to PHAST.

In this work, we show the status of this ongoing project that we plan to finish and release freely. Using the Caffe CPU only code, we aim to port it using the performance portability approach, it means that one source code can be targeted to multiple devices. We have obtained to run two LeNet based networks for MNIST and CIFAR fully in PHAST, and passing almost all Caffe tests despite the unimplemented functionality.

We are not ready to give performance results, but our firsts tests give about 10x of performance loss, but with a few changes we got 2x improvement, and this can be greatly improved.

We continue working in this project, explore the remainder programming models is essential to complete our classification and finish the Caffe port extending it to other devices and getting better performance.

Resumen

En el inicio de la computación, aumentar el rendimiento ha sido la meta a conseguir. En el inicio, incrementar la frecuencia de la CPUs en cada generación era la clave. Adicionalmente incluyendo más transistores, creando la conocida Ley de Moore, que indicaba que cada 2 años se doblaba la cantidad de transistores en un circuito integrado. Así ha sido hasta principios del año 2000, donde los procesadores estaban alcanzando el límite de densidad energética. Con un cambio de paradigma, la creación de los procesadores multinúcleo, conseguimos seguir incrementando el rendimiento.

Durante este periodo de crecimiento de la capacidad de cómputo, distintos modelos de inteligencia artificial se iban creando, siendo las ANNs el algoritmo más conocido. Sin embargo, el uso de estos algoritmos conlleva una gran cantidad de tiempo, por lo que a mayor potencia de cómputo disponible, redes neuronales más complejas pueden ser desarrolladas, y de esta manera resolver problemas mayores.

A pesar de que aún seguimos incrementando el rendimiento de los procesadores multinúcleo, llevamos usando aceleradores hardware para ciertas tareas desde hace varios años. En un principio, los aceleradores más comunes, las GPUs, empezaron a ser usadas para cómputo de propósito general en lugar de ser exclusivamente para gráficos, pudiendo explotar así su paralelismo con muchos núcleos pero muy simples.

En esta revolución de aceleradores, empezamos a usar las FPGAs como hardware reconfigurable en tiempo de ejecución, en lugar de solo ser usadas como prototipado. A su vez, existen los ASICs, hardware específico aún más potente y eficiente que las FPGAs, siendo el ejemplo más conocido los minadores de bitcoins, ASICs desarrollados para hacer la función SHA256.

Con la gran cantidad de aceleradores disponibles, resulta extremadamente tedioso desarrollar aplicaciones capaz de usar varios aceleradores, ya que cada uno dispone de un lenguaje de programación propio denominado DSL. Para poder extraer el potencial disponible de un dispositivo hardware, resulta necesario el uso de su DSL. Por lo tanto, una aplicación que quiera extraer toda la potencia de cómputo que existe actualmente en los clústers heterogéneos, tiene que realizar distintas versiones del mismo código fuente, una para cada dispositivo que desee utilizar, incrementando radicalmente la complejidad del código.

Para facilitar el desarrollo de estas aplicaciones, existen los modelos de programación de alto rendimiento, los cuales facilitan la programación en entornos de alto rendimiento. Sin embargo, existen una gran cantidad de estos, y no hay ninguna clasificación actualizada de los modelos existentes. A pesar de esto, hay un pequeño estudio de 2013, sin embargo la mayoría de los modelos clasificados están en desuso.

A partir de esto, hemos desarrollado una clasificación de 23 modelos de programación de alto rendimiento de 50 que hemos revisado, siendo descartados la mayoría por estar descontinuados o haber sido mezclados con otros modelos ya existentes. Pero para saber las características principales que debe tener la clasificación, hemos empezado por coger varios modelos conocidos, y hemos examinado sus características para ver cuales son los puntos importantes a revisar. Entre estos modelos hemos elegido a OpenMP, OpenACC, OmpSs, OpenCL, SYCL, Kokkos, y PHAST. De esta selección, hemos concluido que una posible clasificación estaría basada en 4 características: las plataformas soportadas, el tipo de modelo, el paradigma de programación, y el paradigma de memoria.

De esa clasificación de 23 modelos, hemos extraído 8 que cumplen la característica de ser portables entre dispositivos, hay varios más, pero esto solo es una pequeña muestra. En estos 8 hemos observado los dispositivos que soportan, y los cambios que se deben realizar para usar el mismo código en otra plataforma.

Con toda esta información, hemos elegido PHAST para realizar este proyecto, en el cual hemos escogido el framework de deep learning Caffe, y hemos portado la funcionalidad del framework para ser portable entre plataformas. Este trabajo lo hemos realizado en colaboración los creadores de PHAST, pudiendo obtener así acceso temprano a la librería y teniendo su soporte durante el desarrollo.

En primer lugar Caffe tiene dos versiones del código fuente, una para CPU y otra para GPU, hemos eliminado el código de GPU para usar el CPU como base durante el desarrollo. Posteriormente hemos

observado el funcionamiento de Ca e como bloques, y hemos decidido portar la funcionalidad mínima necesaria para ejecutar dos redes basadas en LeNet disponibles como ejemplos de Ca e, una para la base de datos de MNIST y otra para CIFAR. Los bloques necesarios a portar son: Blob, Convolution, Pooling, InnerProduct, ReLU, SoftMax, SoftMax with Loss, y Accuracy.

Durante el trabajo, hemos usado la versión 1.0.1 de PHAST con GCC 6.30 y Cuda 9.0, usando como base Ca e obtenido del repositorio oficial en github. Con esta configuración, y un a máquina con dos Intel Xeon CPU E5-2603 v3 @ 1.60 GHz y una Geforce GTX 1080 8GB, hemos conseguido ejecutar satisfactoriamente ambas radas para CPU y GPU. Siendo el único cambio necesario el uso de un makefile u otro dependiendo de la plataforma deseada.

Adicionalmente, hemos decidido ejecutar los tests de Ca e para comprobar el estado actual de nuestro proyecto respecto a la implementación original, y a excepción de la funcionalidad no implementada, pasamos todos los tests. Y al ser un proyecto todavía en desarrollo, no tenemos resultados de rendimiento válidos, pero como primer vistazo, obtuvimos un 10x de perdida de rendimiento, que mejoró un 2x tras una pequeña revisión del código. Sin embargo, los casos de prueba son demasiado pequeños como para tenerlos en cuenta, y estamos convencidos de que pueden ser mejorados modernizando nuestro código modernizado.

Tras este trabajo, pensamos que nuestra clasificación puede ser usada como base para seguir clasificando el resto de modelos existentes, cosa que se queda como trabajo futuro. Y entre estos trabajos futuros también se encuentran pendientes terminar la adaptación de Ca e a PHAST, extender el trabajo a múltiples dispositivos, y mejorar el rendimiento obtenido.

List of Figures

1	Ca e framework as a block diagram	11
2	Blob block as a data container	12
3	Classic 2-D Convolution	12
4	2-D Convolution as a GEMM	13

List of Tables

1	High Performance Programming Models Classification	8
2	Single Source Programming Models	9
3	Case tests results	22

Listings

1	Ca e version: im2col	13
2	PHAST version: im2col	13
3	Ca e version: col2im	14
4	PHAST version: col2im	14
5	Ca e version: Max Pooling	15
6	PHAST version: Max Pooling	15
7	Ca e version: Back Max Pooling	16
8	PHAST version: Back Max Pooling	16
9	Ca e version: InnerProduct	16
10	PHAST version: InnerProduct	16
11	Ca e version: Back InnerProduct	17
12	PHAST version: Back InnerProduct	17
13	Ca e version: ReLU	18
14	PHAST version: ReLU	18
15	Ca e version: Back ReLU	18
16	PHAST version: Back ReLU	18
17	Ca e version: SoftMax	19
18	PHAST version: SoftMax	19
19	Ca e version: Back SoftMax	20
20	PHAST version: Back SoftMax	20
21	Ca e version: Accuracy	21
22	PHAST version: Accuracy	21

This page has been intentionally left blank

1. Introduction

Since the start of computing, the performance was always a limit to break. In the beginning, the way to break this limit started by increasing the frequency of each CPU generation. Additionally, CPUs started to include more and more transistors. This worked for a lot of time, creating the Moore's Law (every two years, the transistors in dense integrated circuit was doubled). But, near to the early 2000s, processors were very near to reach its limit in power density [1]. Changing to another paradigm was necessary to continue improving the performance, the creation of multi-cores.

After all the new unlocked performance for CPUs, there are several things that other hardware devices are able to do faster and more efficient. In instance, real-time graphics processing continue to be done by specific accelerators (GPUs). These devices are developed with different paradigms in mind, and normally to be very efficient in only one task. GPUs have many more cores than a CPU, but they are simpler.

Another well-known accelerator is the Field Programmable Gate Array (FPGA), created in the beginning as a hardware prototyping device. Nowadays, they are also used in many other tasks like a reprogrammable accelerator.

Low power high-performance devices are highly demanded, and specific architectures builtin inside an Application-Specific Integrated Circuit (ASIC) are getting a lot of attention. Many custom architectures are being developed for Neural Networks, Physics, Simulations, etc.

These accelerator devices have shown that CPUs are not always the best option to solve a problem. As the multi-core era is ending, Domain Specific Architectures (DSAs) are getting the testimony [2].

On the other hand, meanwhile CPUs started to get relevant performance, machine learning fired. One of the most powerful machine learning techniques is Deep Neural Networks (DNNs), but their algorithms required a lot of computing time. Until now, a vast machine learning algorithms have been developed and improved.

Once the technology was able to process DNN algorithms in reasonable times, they have become more and more ubiquitous in everybody's daily life. Image classification, language processing, economics, medicine, video games, among others, are some of the fields where DNNs are being used. They are trying to improve, but requiring more time, and accelerators are used to reduce the amount of time needed.

In this new era of accelerators, the creation of a programming model able to use all this kind of new accelerators is extremely important [3]. This new approach is called performance portability [4], and it is not yet completely solved.

In this work, we have focused on the following objectives.

- **Search for relevant programming models for multiple device programming:** There are a vast of programming models for multiple devices, therefore we want to find and classify them.
- **Learn to use a high performance programming model:** Like any other programming paradigm, these programming models have their own way to be used. Being able to use one of them could not be an easy task.
- **Apply the previously selected model to a real application:** A new programming model can be awesome, but if it is not used, it has meaningless.

The ultimate goal of this work is to have a programming language that allows that one source code can be run in CPU or GPU with only recompiling it with the right options, and in addition, that could be targeted to other devices with no changes (in the source code). When it is finished, we will release it to allow everyone to use it freely.

2. Background

2.1 CPU Revolution

The multi-core revolution, after reaching the power density limit in traditional processors, allowed the creation of a new paradigm where the increase of throughput is based in the parallel execution of multiple independent workloads. Most of the computer power in data centers was brought by mono- and multi-core CPUs.

From some time ago, some researchers, most of them in computer vision, noticed that the GPUs, that they were already using, were able to solve their problems faster. At the beginning, they started using programming languages for graphics, like OpenGL to implement some algorithms [5].

2.2 Accelerator Revolution

2.2.1 GPU & GPGPU

The most known accelerator is the GPU, able to generate graphics and compute data at high speed exploiting its Single Instruction Multiple Threads (SIMT) execution model, based in Single Instruction Multiple Data (SIMD) but with multithreading. Nowadays, GPUs are used like CPUs (GPGPU), we have specific languages to use them such as Cuda, or OpenCL. It is probably the most used accelerator.

2.2.2 FPGA

A FPGA is a configurable integrated circuit, built from an array of logic blocks and routing channels [6]. Each block is able to represent a little mathematical function, and connecting many of them, it is possible to build very complex devices.

In the beginning, it was mainly used for prototyping, but now it is also used as accelerators. There will no be as fast as a ASIC, but they can be reconfigured and can be used to implement any part of a program.

2.2.3 ASIC

To be faster than a FPGA and more power efficient, ASIC is developed and use in specific tasks. We have several ASICs currently in our devices. From the lowest computer expensive ones like the usb or keyboard controllers, to the bitcoin miners, which are ASICs developed to make the SHA256 function.

2.2.4 The problem

There are a vast of accelerators ready to be used and most of them are problem specific. Each one of these accelerators use its own programming language, a Domain Specific Language (DSL). Using the specific DSL for each device allows to benefit from the specific features of that device [2].

The problem with this accelerator boom is that each DSL is quite different. Therefore, making a portable software between accelerators is impossible without having multiple source files (at least one file per target device). Maintaining software with multiple source code files for the same functionality in different platforms is a very tedious task, greatly increasing the code complexity [7].

2.3 ML & DNNs

2.3.1 Development and usage

Theoretical and mathematical models of the artificial intelligence techniques were developed in the twentieth century, but the lack of computing power prevented it from progressing. Artificial Neural Networks (ANNs) are a type of brain-inspired learning algorithm, built from small units called neurons, being the perceptron the most common one. With the recent advances, deep neural networks have

been incorporated into numerous fields, such as image classification, language processing, economics, video games, and medicine. In some tasks, this new technology is being able to outperform human performance.

Medical image analysis has started to implement deep learning for screening and localization of malignant zones. Additionally, other medical areas are working with these kind of techniques as well, like the analysis of the genetic information inside DNA and RNA series [8]. The common objective is not replace physicians with deep learning techniques but supporting them to make better diagnoses.

2.3.2 Evolution with HPC

Machine learning techniques could be difficult to code and debug, therefore many frameworks have been developed to ease its use. Most of them are open source with software for most of the types of neural networks. The most known ones are Caffe, Caffe2, Tensorflow, Theano, PyTorch, Mxnet, and CNTK among others [9]. And there are frameworks like Keras, providing a more high-level experience, running on top of some of the aforementioned frameworks.

Specifically, the training phase is very time-consuming, since it is evaluating an optimization problem with hundreds, thousands, or even millions of parameters. Therefore, the reduction of the training phase execution time is a desirable feature for all frameworks. Thanks to this shorter training time, scientists using these frameworks can explore a wide solution space, and even develop more complex networks.

Therefore, most frameworks are able to use Nvidia's cuDNN, BLAS-like libraries, MKL-libraries, OpenCL, even specific libraries for custom integrated circuits, to speed up this computation. Also, many of them also allow other kind of parallelism for multiple nodes (using the MPI library).

3. An overview of high performance parallel programming models

3.1 Performance Scaling Problems

A High Performance Computing (HPC) cluster can run a vast of different jobs, each one specialized in different kinds of applications. Traditional clusters had many difficulties in scaling the applications they had to run. One of the mitigations taken was the use of heterogeneous clusters instead of the traditional ones. They are more flexible, allowing the use of the specific resources to fit the application better.

Now, we have several different architectures together to be used by the applications, but this is not as easy as it may appear. An application needs to have the code to be run in each one of this architectures to be able to use the cluster efficiently.

3.2 High Performance Programming Models

There are many things to keep in mind to be able to program an HPC cluster efficiently, such as dependencies, communications, locks, possible failures, without forgetting multiple versions of the same code for each available architecture.

To address this problem, there are many programming models to help developers writing efficient code. In instance, OpenMP tackles multi-core but locally in the same machine, and MPI tackles the communication between nodes. But, these are low-level approaches, there are others with a higher abstraction level.

We want to know which programming models are available, which is the main focus of each one, and their characteristics, such as the programming paradigm or the memory abstraction model. But only partial classification have appeared. There is an interesting survey from 2013 [10], but many of those models are now deprecated. Therefore, we start to classify some of them, but to know which features are important, we start by selecting a few and looking for their main details.

3.3 Some High Performance Programming Models

This first selection of programming models is done with the idea of having at least one well-known programming model with different approaches. But, also with two very similar approaches to be able to find their differences. We think that we are able to find the main features of the programming models to classify them.

3.3.1 OpenMP

OpenMP is a specification for a set of compiler directives, library routines, and environment variable that can be used to specify high-level parallelism in Fortran and C/C++ programs [11]. These directives are also known as pragmas, they give information to the compiler about what to do with a statement. And, as it is only an Application Programming Interface (API), compilers have the responsibility of implementing it.

Nowadays, OpenMP has improved a lot; from being able to parallelize loops to currently use devices as accelerators. The API has an offload option, allowing the usage of some accelerators compatible with the implementation. For instance, the Intel Compiler supported the Xeon Phi KNC Coprocessor as an offloading target. But now, a lot of compilers are making a lot of effort trying to accomplish the GPU support using offload.

Depending on the point of view, it is possible to refer to OpenMP as a kernel or a task programming model. However, the main feature of OpenMP is the use of directives to annotate the source code in an elegant and efficient way. It is focused on shared memory environments using multi- and many-core CPUs.

3.3.2 OpenACC

OpenACC is a specification for a set of compiler directives performance-portable parallel programming model designed for scientists interested in porting their codes to a wide variety of heterogeneous HPC hardware platforms and architectures [12]. Similar to OpenMP, it is only an API which needs to be supported by compilers.

Like OpenMP, it is a directive-based programming model but focused in heterogeneous hardware, which means that OpenACC is expecting to be run from a CPU to an accelerator device, like another kind of CPU or a GPU.

3.3.3 OmpSs

OmpSs is a programming model developed at Barcelona Supercomputing Center (BSC) with the objective of merging StarSs with OpenMP as a way of extending its directives to support asynchronous parallelism and heterogeneity [13]. Nowadays, it supports Intel CPUs, Nvidia GPUs, Intel Xeon Phi KNC (native and offload), IBM Power8, ARM CPUs, and OpenCL compatible MALI GPUs. Also, it is evolving into OmpSs-2, a new programming model with a different approach.

OmpSs is builtin in the Mercurium Compiler and executed by the Nanos++ runtime system. Also, similar to OpenMP and OpenACC, it is a directive-based programming model, but in this case, tasks are the main focus.

3.3.4 OpenCL & SYCL

OpenCL has a very different approach to OpenMP, it is a kernel based programming model, and it is nearer to a language instead of an API. It is an open standard originally developed by Apple, but currently maintained by Khronos [14]. OpenCL depends on a Installable Client Driver (ICD), which manages all of the OpenCL calls. First, a loader retrieves all the available ICDs. After preparing the context of the execution, kernels get compiled to that specific platform at runtime and executed through the ICD.

SYCL is a specification over OpenCL to enable single source C++ programming to OpenCL, and Khronos expects to be a step forward in standards convergence. SYCL is very similar to OpenCL, but with many features from new C++ standards (11, 14, and 17) and enabling the C++17 parallel Standard Template Library (STL).

3.3.5 Kokkos

With a totally different approach, Kokkos is one of the modernist approaches available. It is a new C++ library with a lot of ideas from STL containers and algorithms. The main objective of Kokkos is to maximize the amount of user code that can be compiled for diverse devices and obtain the same performance as a variant code specifically written for that device [15].

It defines a new template type called View, that is used to allocate the data that will be used later. In this View type, it is possible to build an n-dimensional structure to be used as a container to be modified by the functions.

Kokkos allows the creation of new functions to modify the data stored in Views, but it needs to be encapsulated into a class with some restrictions and with the main entry point made with the operator '()'. This is a functor, and it is a concept borrowed from the C++ STL. These functors are run from special functions with a similar nomenclature as the algorithms library in C++ STL.

3.3.6 PHAST

Continuing with the previous idea, Parallel Heterogeneous-Architecture STL-like Template (PHAST) library [16–19] is a modern C++ programming template library based on the classic STL containers with the performance portability philosophy. It currently supports multi-core CPUs and Nvidia GPUs,

allowing the users to write expressive and concise sequential-like code that can be automatically parallelized. Its main goal is to let the programmers code using high-level programming approaches without preventing them from applying low-level optimizations, keeping the main code at a higher level of abstraction.

Similar to STL containers, PHAST provides a vector container but extending it with matrix, cube, and grid, all of them working with a very similar interface. Also having vectors like primitives similar to an SSE or AVX vectors.

These containers can be modified using functors. Similar to Kokkos, the idea of a functor is borrowed from STL as a struct that inherits from a base functor struct and at least has the operator '()' defined. There are a lot of algorithms and factors predefined, but this allows the creation of functionality that is not defined by default in the library, without losing performance or portability.

In its roadmap, there is planned support for multiple devices in the same executable, lambda syntax, FPGA support, OpenCL backend, multi-GPU, and many other interesting features.

After writing the sequential code in the PHAST way, it is possible to change the device target changing a macro and the compiler. Therefore, having two different makefiles makes the trick. The most important thing is that the code has not changed, only the compiling process.

3.4 Classification

With this little overview of some programming models, we can proceed to classify them. But, first, it is necessary to specify which features we will be looking for, and a definition of it. There are some terms that could be misunderstandings or misinterpreted, and in this way, there will not be any kind of confusion.

3.4.1 Platforms

With the accelerators revolutions, there are a lot of devices that could be interesting, but in this case, we focus on CPU, GPU, and FPGA, without paying attention on any specific manufacturer.

3.4.2 Type

Depending on the requirements of modifying existing compiling architectures, we have defined 3 types of programming models:

- **Language:** A new set of primitives using a new compiler or a runtime system to run, such as Cuda or OpenCL.
- **Extension:** A modification to the compiler to add a set of directives or even primitives to be used in an existing language. These modifications may require a runtime system to run, but the original language is barely modified. An example of this could be OpenMP or AllScale.
- **Library:** Without any kind of modification to the existing compilers, it adds a set of primitives to be used in a program using the language specification. It could be a shared library or even a set of files to add to the compiling state. MPI, Phast, and TBB are good examples of a library.

3.4.3 Programming Paradigm

Each high-performance programming model has its own type of programming paradigm, i.e., how it is intended to be used. Most of the programming models can be classified into multiple programming paradigms. These are the programming paradigm we have taken into account:

- **Kernel:** Time-expensive parts are encapsulated into a special section that can be compiled apart, and it can be executed by other parts of the code. For instance, OpenCL creates kernels that are compiled in runtime and executed in the target platform.

- **Task:** The program is split into several tasks and they are run following the dependency graph. This can be very confusing because it is possible to use a kernel programming model as a task one, but it is not the same paradigm.
- **Directive:** Using a reduced set of primitives, it sets to the compiler what to do with a statement. It could be achieved with a set of functions or a set of pragmas.
- **C++ Template:** With the creation of C++11 and C++14 standard, the interest in the STL grew, new containers and the algorithm library extended its usage. In this paradigm, the data is stored in a set of containers, and only are modified using the functions available by an algorithm-like library, that usually allows to extend this set of functions. The functions executed by the algorithm-like library could seem like kernels, but the concept and data management is very different.
- **Skeleton:** The program has to be mapped to use a set of generic components with a specific pattern of computation and data dependence. For example, a dot product can be implemented as a MapReduce of multiplication and addition.
- **Wrapper:** In a wrapper programming model, it inherits the underlying programming paradigm, but encapsulating the functionality.
- **Threads:** This is not a programming paradigm by itself, as depending on the usage of the created threads it could be used as any paradigm.

3.4.4 Memory Paradigm

Memory paradigms are a bit tricky. Some of them are very similar, but the high-level concept is different, therefore we have split them into 6 categories.

- **Hierarchical:** Implicitly or explicitly copied, a hierarchical memory paradigm is shown when there are a host and at least one slave that is going to exchange information. It is normally associated with a CPU host and a device (like a GPU).
- **Shared:** In a shared memory environment, at any moment is possible to read or write at any part of the memory. Excluding conflicts between threads, devices, or compute units, all the memory is visible.
- **Implicit:** All memory transactions are handled by the programming model, without the involvement of the program. In some cases, the programming model allows to explicitly manage this memory, but it is not necessary.
- **Explicit:** This is exactly the opposite of implicit memory. The programmer must handle all memory transactions. In some cases, some of these transactions can be implicit if they are simple enough.
- **PGAS:** Partitioned Global Address Space (PGAS) assumes that there is a global memory address space, but it is logically partitioned to each thread, device, or computer unit.
- **Distributed:** Each thread, device, or compute unit has its own complete memory. Distributed memory is able to exchange data, but it is not necessary. In some cases, it is very similar to the explicit memory paradigm and hierarchical.

3.4.5 Table

With all this information, we have selected a total of 23 high-performance programming models from about 50 reviewed programming models to be classified in Table 1. Many of the discarded models are discontinued or have been merged into another one.

		Programming Paradigm							Memory Paradigm													
Type	Language	Kernel	Task	Directive	C++ Template	Skeleton	Wrapper	Threads	Hierarchical	Shared	Implicit	Explicit	PGAS	Distributed								
															OpenCL	OpenCL	Cilk	OpenMP	OpenAcc	Phast	Thrust	Boost.Compute
CPU	Extension		Cilk						OpenCL													
	Library	SyCL	Cilk	OpenMP	OpenAcc				OpenMP	Cilk	AIIScale	OpenAcc										
GPU	Extension								OpenCL													
	Library	SyCL		OpenMP	OpenAcc				Cuda	OpenMP		OpenAcc										
FPGA	Extension								OpenCL													
	Library	SyCL		OpenMP								OpenAcc										

Table 1: **High Performance Programming Models Classification Table:** HLS has been discarded because it depends heavily on the manufacturer. Libraries like: C++ AMP, ArrayFire, POCL, Halide, NOVA, Bolt, FastFlow, Muesli, Chapel X10, UPC, among many others, have been omitted because some of them have been discontinued and others have evolved into new models. **Some highlights:** - OpenMP is in two memory paradigms at the same time because the device target mode is hierarchical. - Cilk is in two model types because is an extension but has features of a new language. - PHAST is the only C++ template paradigm with implicit memory transference. - All FPGAs programming models are based in OpenCL.

3.5 Portable Programming Models

Single source programming targeting multiple devices is a development goal to be accomplished. Performance portability is the programming paradigm with the ability to resolve the problem. It recalls well-established highly-expressive techniques. For instance, programming is done using a set of primitives to develop complex programs, and most of the performance portability techniques use the same approach but with higher level primitives.

3.5.1 Classification

When trying to use multiple devices, some of these models could not be enough, but there are others with the characteristic of being single source programming models. We have extracted some of these models from our research (shown in Table 2) and find their principal device targeting and the changes needed to use the same source code in a different device.

Programming Model	Programming Paradigm	Memory Paradigm	Devices	Changes to other device
OpenMP	Directive	Hierarchical Shared	CPU, GPU Others	Few changes (only the directives)
OpenAcc	Directive	Hierarchical	CPU, GPU FPGAs, Others	Few changes (only the directives)
OmpSs	Directive	Hierarchical	CPU, GPU FPGAs, Others	Few changes (only the directives)
OpenCL	Kernel	Hierarchical	CPU, GPU FPGAs, Others	ICD, loader, few changes
SyCL	Kernel	Hierarchical	CPU, GPU FPGAs, Others	ICD and Loader
SkePU	Skeleton	Explicit	CPU, GPU	Recompile
Kokkos	Kernel	Explicit	CPU, GPU, Others	Recompile
PHAST	C++ Template	Implicit	CPU, GPU	Recompile

Table 2: **Single Source Programming Models:** Some of the models from Table 1 that are single source programming models, and the changes needed to use the same code to another device.

4. PHAST applied to the Caffe Framework

4.1 Why PHAST

PHAST code can be written once and targeted to different devices via a single macro. Its inner layers are implemented in Cuda C++ or `std::threads` so to allow targeting on Nvidia GPUs and multi-core CPUs. These layers are not part of the interface, so users can express their code in a platform-independent way. In fact, PHAST programmers can code their applications in terms of containers, iterators, and algorithms in an STL-like, thus using common sequential techniques [18].

As we have said in the previous section, PHAST use functors to modify the data stored in containers. This idea of a functor is borrowed from STL as a struct that inherits from a base functor struct and at least has the operator `()` defined. These functors are executed like kernels, therefore it is possible to see it like a loop calling the function in each iteration using each time the next element in the containers.

Despite PHAST promoting an architecture-agnostic programming style for productivity reasons, some small portions of code could benefit in performance from leveraging low-level architecture-specific features and optimizations. Well-known techniques like code replication on GPUs or the use of intrinsics on multi-cores are examples of this kind of low-level optimizations. Due to a possible performance gain, PHAST users could be willing to specialize some portions of their code according to the underlying device. PHAST does not prevent them from doing so, and low-level architecture-specific optimizations are still possible under the scope of `ifdefs` or similar constructs.

To test the PHAST library, it has been used to implement some applications such as AES [17], TRIAD [18] and DCT8x8 [18]. Also, an implementation of a histogram-stretching and an unsharp-mask filter is available as a set of tasks [19].

The PHAST library is available through request at its web page¹, but as part of a joint collaboration, we have early access to it.

4.2 The Caffe Framework

Caffe is the first DNN framework, developed by Berkeley AI Research. Nowadays, in the production environment, Caffe is replaced by Caffe2 and pyTorch, the Caffe successors, but in other cases, TensorFlow and MXNet are the selected ones. Caffe continues to be used in researching, due to the fact that it is very easy to modify, extend, or use, all of this without losing flexibility to run most of the state-of-the-art DNN models.

Caffe is a very good framework that runs on CPU or GPU only changing a flag, but this is done by having two implementations in two source codes files, one for CPU (.cpp) and other for GPU (.cu). Therefore, developers are forced to maintain two different versions of the same functionality. In this case, this is very well done and there are not a large number of differences between the files being the perfect target for a single source approach.

Before starting, an overview of the Caffe's internal structure, is given next. Caffe is built from multiple modules that work by themselves. It is possible to classify the blocks in two parts, containers and executors. Containers store data to be used by executors. Executors use the containers to exchange data and process it. For example, a layer gets a set of blobs, and with its own blobs, it computes the output ones.

A neural network has two different phases, inference and training. In the inference phase, data is passed through all the layers of the neural network in feed-forward mode to reach the last layer and get a result. But, in training mode, data is passed through all layers in feed-forward, like inference, but when it reaches the last layer of the network, that data is brought to a solver, it recalculates some values and starts the back-propagation through each layer but in reversed direction.

In this work, we are going to port all the blocks needed to be able to run two LeNet variations from example networks in Caffe for MNIST and CIFAR-10 databases. The blocks we need to port are

¹<https://www.phast-library.com/>

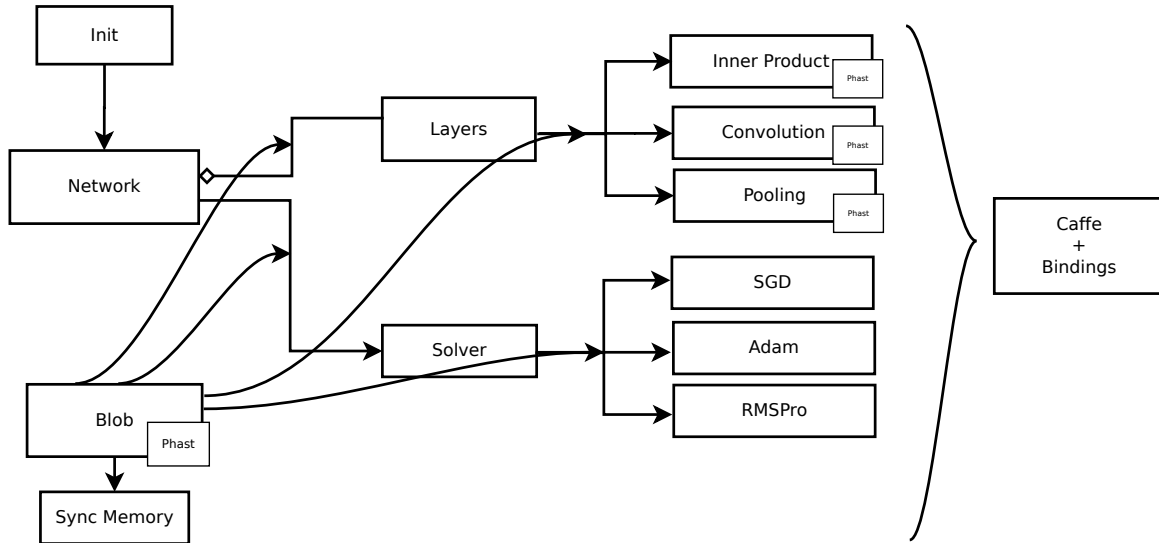


Figure 1: Part of the Caffe framework as a block diagram showing communications between blocks.

Blob, InnerProduct, Convolution, Pooling, ReLU, Accuracy, SoftMax, and SoftMax with Loss. We expect to run these neural networks through the Caffe binary in train and test mode using PHAST, therefore, we choose to modify only the blob and the lowest part of the main executors (Figure 1).

4.3 Development

We have removed the GPU code from Caffe to use the CPU version as our base for this project. In this way, we can focus on the algorithm itself and not in its GPU implementation.

From each block (or layer), we explain its most relevant part, showing its code compared to the original implementation. Sometimes, this new code is larger, but it is portable across the platforms that PHAST supports. We will show some functors related to its own block, but all of them are collected in one file, enabling the shared use of them.

4.3.1 Blob

The Blob block is a container for all the information in the neural network. It has to store two big arrays of data to be used at any required executor. The first array is named data and the second is diff. data is used as part of the feed-forward step. On the other hand diff is used to modify data at the back-propagation step.

Despite the fact that Caffe supports two data types, floating point in single and double precision, the Caffe binary only supports single precision, but using the python binding is possible to access to the double precision. Also, other data types are unsigned and signed integer that are used as extra information for the rest of the framework. One problem we found is that the current version of PHAST allows only single precision floating point numbers, therefore we have to override the Blob data type to use PHAST containers only when is possible.

The Blob block depends on the SyncMemory block to manage the memory, and depending on if the GPU is being used or not, it allocates the memory and manages the syncing process (Figure 2a). In our case, since we are only modifying a part of the framework, all the remainder part has to interact with our container, therefore we used the same trick as the GPU version of SyncMemory and have two containers, a host raw pointer and a PHAST vector (Figure 2b).

Although the Blob block is a container, it has some functionality inside. All of this functionality has been ported to PHAST, and we have added some extra functionality respecting to the PHAST vector container to convert it to another container (like a matrix or a cube) using a shallow copy.

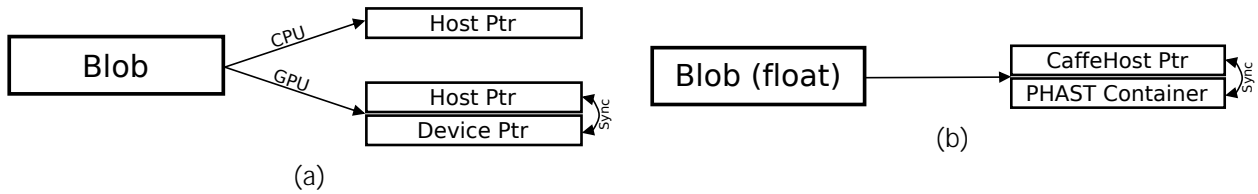


Figure 2: The Blob block from original Caffe (a) and the instanced version for single precision floating point numbers using PHAST (b).

4.3.2 Convolution

4.3.2.1 Feed-Forward

The Convolution block, or Convolution Layer, applies to the input a set of filters using a sliding window over the input. There are many ways to implement this sliding window, but applying the filter is doing a vector inner product for each sliding window. The most common variant of Convolution is Convolution 2-D (Figure 3), which is a simplification of a Convolution N-D.

As our sample network only use Convolution 2-D, we only port that specific variation. There is not much difference between a Convolution 2-D and a Convolution N-D, because we use the `im2col + gemm` implementation.

The `im2col + gemm` implementation is a way to map a convolution as a matrix multiplication (General Matrix Multiplication (GeMM)), but a data manipulation is needed to accomplish it. The `im2col` function maps the input matrix into columns to make the Convolution using a GeMM (Figure 4).

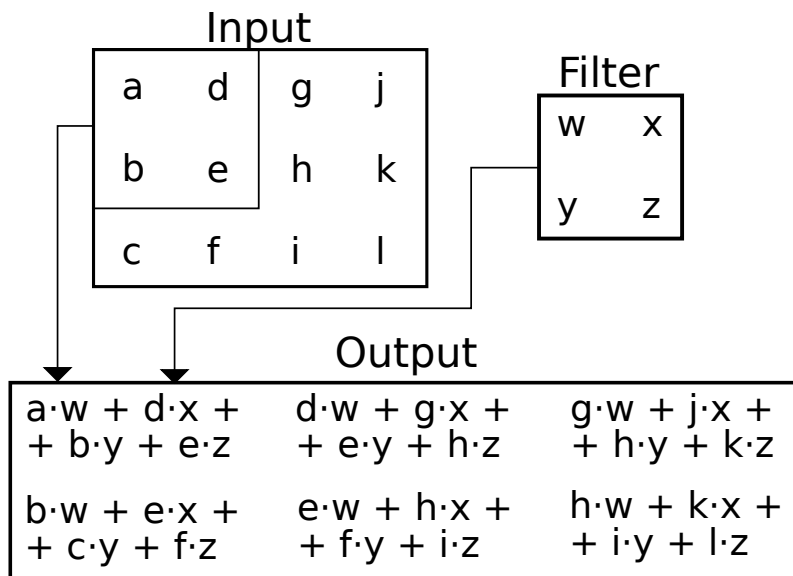


Figure 3: A convolution example using a 2x2 filter with stride 1 and padding 0 over a 4x3 input matrix.

The original Caffe's `im2col` function is a Penta-loop with dependencies in each iteration (Listing 1), therefore we decide to adapt it to be able to exploit a bit of parallelism. To create the PHAST version, we have merged all the loops and parametrized it with only one index (Listing 2). This change allows PHAST to use all the threads it sees appropriate.

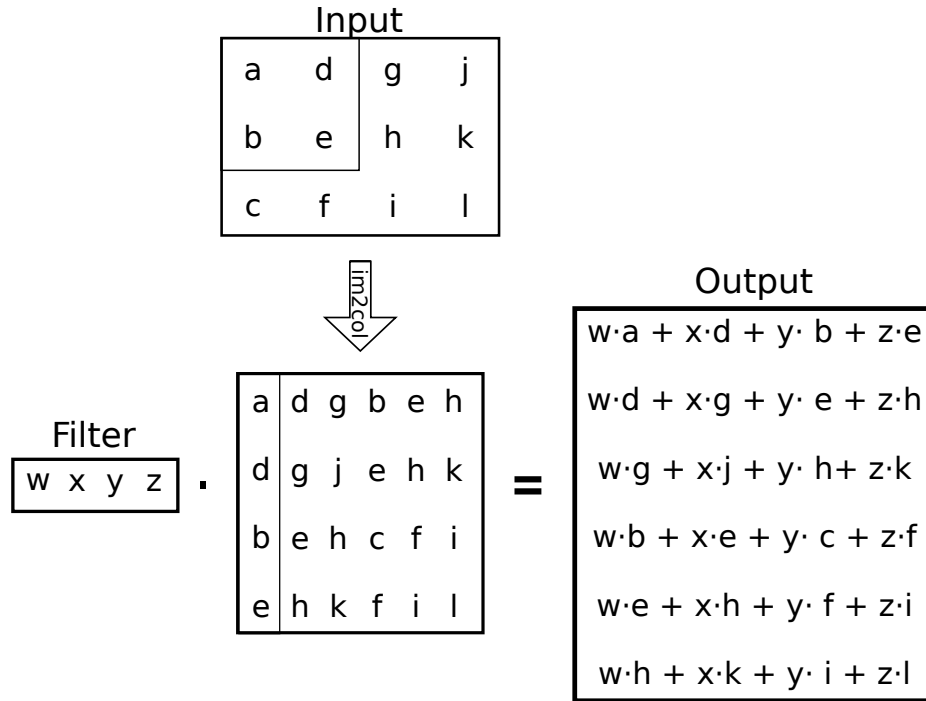


Figure 4: Convolution as a GeMM using the im2col function with 2x2 filter, stride 1, and padding 0.

```

1  template <typename Dtype>
2  void im2col_cpu(const Dtype data_im, const int channels,
3  const int height, const int width, const int kernel_h, const int
4  kernel_w,
5  const int pad_h, const int pad_w,
6  const int stride_h, const int stride_w,
7  const int dilation_h, const int dilation_w,
8  Dtype data_col) {
9  const int output_h = (height + 2 * pad_h -
10 (dilation_h * (kernel_h - 1) + 1)) / stride_h + 1;
11 const int output_w = (width + 2 * pad_w -
12 (dilation_w * (kernel_w - 1) + 1)) / stride_w + 1;
13 const int channel_size = height * width;
14 for (int channel = channels; channel-- > 0; data_im +=
15 channel_size) {
16 for (int kernel_row = 0; kernel_row < kernel_h; kernel_row++)
17 for (int kernel_col = 0; kernel_col < kernel_w; kernel_col
18 ++ ) {
19 int input_row = -pad_h + kernel_row * dilation_h;
20 for (int output_rows = output_h; output_rows > 0; output_rows
21 -- ) {
22 if (is_a_ge_zero_and_a_lt_b(input_row, height)) {
23 for (int output_cols = output_w; output_cols > 0; output_cols
24 -- ) {
25 (data_col++) = 0;
26 }
27 } else {
28 int input_col = -pad_w + kernel_col * dilation_w;
29 for (int output_col = output_w; output_col > 0; output_col
30 -- ) {
31 if (is_a_ge_zero_and_a_lt_b(input_col, width)) {
32 (data_col++) = data_im[input_row * width +
33 input_col];
34 } else {
35 (data_col++) = 0;
36 }
37 }
38 }
39 }
40 }
41 }
42 }
43 }

```

Listing 1: Caffe version: im2col

```

1  template <typename T, unsigned int policy = phast::
2  get_default_policy()>
3  struct im2Col : phast::functor::func_scal<T, policy> {
4  _PHAST_METHOD im2Col(int dilation_h, int dilation_w,
5  int stride_h, int stride_w,
6  int pad_h, int pad_w,
7  int kernel_h, int kernel_w,
8  int height, int width) {
9  dh = dilation_h; dw = dilation_w;
10 sh = stride_h; sw = stride_w;
11 ph = pad_h; pw = pad_w;
12 kh = kernel_h; kw = kernel_w;
13 h = height; w = width;
14 }
15 }
16 _PHAST_METHOD void operator()(phast::functor::scalar<T> &
17 col) {
18 int index = this->get_index();
19 const int oih = (h + 2 * ph - (dh * (kh - 1) + 1)) / sh + 1;
20 const int oiw = (w + 2 * pw - (dw * (kw - 1) + 1)) / sw +
21 1;
22 const int irow = -ph + (((index / (oih * oiw)) % (kh * kw)) /
23 kw) * dh + ((index % (oih * oiw)) / oiw) * sh;
24 const int icol = -pw + (((index / (oih * oiw)) % (kh * kw)) %
25 kw) * dw + ((index % (oih * oiw)) % oiw) * sw;
26 if (irow >= 0 && irow < h && icol >= 0 && icol < w)
27 col = in.at(((index / (oih * oiw)) / (kh * kw)), irow, w +
28 icol);
29 else col = 0;
30 }
31 }
32 int dh, dw;
33 int sh, sw;
34 int ph, pw;
35 int kh, kw;
36 int h, w;
37 }
38 phast::functor::matrix<T> in;
39 };

```

Listing 2: PHAST version: im2col

4.3.2.2 Back-Propagation

In the feed-forward stage, the `im2col` function duplicates some values to make a Convolution with a GeMM. In the back-propagation, we need to apply the reverse step to propagate the gradients to the previous layers. The most important part is the usage of `col2im` to map the gradients to the size of the input data.

Like in the feed-forward stage, the original implementation is also a Penta-loop (Listing 3), therefore we have followed the same approach as before and we have merged the loops and parametrized with only one index (Listing 4).

```

1  template <typename Dtype>
2  void col2im_cpu(const Dtype data_col, const int channels,
3                const int height, const int width, const int kernel_h, const int
4                kernel_w,
5                const int pad_h, const int pad_w,
6                const int stride_h, const int stride_w,
7                const int dilation_h, const int dilation_w,
8                Dtype data_im) {
9      ca_e.set(height width channels, Dtype(0), data_im);
10     const int output_h = (height + 2 * pad_h -
11                          (dilation_h * (kernel_h - 1) + 1)) / stride_h + 1;
12     const int output_w = (width + 2 * pad_w -
13                          (dilation_w * (kernel_w - 1) + 1)) / stride_w + 1;
14     const int channel_size = height * width;
15     for (int channel = channels; channel-- > 0; data_im +=
16         channel_size) {
17         for (int kernel_row = 0; kernel_row < kernel_h; kernel_row++)
18             for (int kernel_col = 0; kernel_col < kernel_w; kernel_col++)
19                 {
20                     int input_row = -pad_h + kernel_row * dilation_h;
21                     for (int output_rows = output_h; output_rows > 0; output_rows--)
22                         {
23                             if (!is_a_ge_zero_and_a_lt_b(input_row, height)) {
24                                 data_col += output_w;
25                             } else {
26                                 int input_col = -pad_w + kernel_col * dilation_w;
27                                 for (int output_col = output_w; output_col > 0; output_col--)
28                                     {
29                                         if (is_a_ge_zero_and_a_lt_b(input_col, width)) {
30                                             data_im[input_row * width + input_col] +=
31                                                 data_col;
32                                         }
33                                         data_col++;
34                                         input_col += stride_w;
35                                     }
36                                 }
37                             input_row += stride_h;
38                         }
39                 }
40     }
41 }

```

Listing 3: Caffe version: `col2im`

```

1  template <typename T, unsigned int policy = phast::
2  get_default_policy()>
3  struct col2im : phast::funcor::func_scal<T, policy> {
4      _PHAST_METHOD col2im(int dilation_h, int dilation_w,
5                          int stride_h, int stride_w,
6                          int pad_h, int pad_w,
7                          int kernel_h, int kernel_w,
8                          int height, int width) {
9          dh = dilation_h; dw = dilation_w;
10         sh = stride_h; sw = stride_w;
11         ph = pad_h; pw = pad_w;
12         kh = kernel_h; kw = kernel_w;
13         h = height; w = width;
14
15         oh = (h + 2 * ph - (dh * (kh - 1) + 1)) / sh + 1;
16         ow = (w + 2 * pw - (dw * (kw - 1) + 1)) / sw + 1;
17     }
18
19     _PHAST_METHOD void operator()(phast::funcor::scalar<T> &
20     col) {
21         int index = this->get_index();
22         int oc = index % ow;
23         int orp = (index / ow) % oh;
24         int c = index / (ow * oh * kh * kw);
25         int kr = ((index / (ow * oh)) % (kh * kw)) / kw;
26         int kc = ((index / (ow * oh)) % (kh * kw)) % kw;
27         int ir = -ph + kr * dh;
28         int ic = -pw + kc * dw;
29
30         if (ic >= 0 && ic < w && ir >= 0 && ir < h)
31             in.at(c, (ir + orp * sh) * w + (ic + oc * sw)) += col;
32     }
33
34     int dh, dw;
35     int sh, sw;
36     int ph, pw;
37     int kh, kw;
38     int h, w;
39     int oh, ow;
40
41     phast::funcor::matrix<T> in;
42 };

```

Listing 4: PHAST version: `col2im`

4.3.3 Pooling

4.3.3.1 Feed-Forward

The Pooling layer applies a mathematical function to a set of numbers to get only one value, such as maximum or minimum. Like the Convolution layer, it works using a sliding window over the input data and applying the function to each set.

As it can be seen in Listing 5, the structure is very similar to the Convolution block, but this time, we do not apply the same technique, because we have not enough time to verify that the merged version works properly. Therefore, we have only parallelized the outer loop (Listing 6).

```

1 // The main loop
2 for (int n = 0; n < bottom[0]->num(); ++n) {
3   for (int c = 0; c < channels_; ++c) {
4     for (int ph = 0; ph < pooled_height_; ++ph) {
5       for (int pw = 0; pw < pooled_width_; ++pw) {
6         int hstart = ph * stride_h_ - pad_h_;
7         int wstart = pw * stride_w_ - pad_w_;
8         int hend = min(hstart + kernel_h_, height_);
9         int wend = min(wstart + kernel_w_, width_);
10        hstart = max(hstart, 0);
11        wstart = max(wstart, 0);
12        const int pool_index = ph * pooled_width_ + pw;
13        for (int h = hstart; h < hend; ++h) {
14          for (int w = wstart; w < wend; ++w) {
15            const int index = h * width_ + w;
16            if (bottom_data[index] > top_data[pool_index]) {
17              top_data[pool_index] = bottom_data[index];
18              if (use_top_mask) {
19                top_mask[pool_index] = static_cast<Dtype>(index
20                );
21              } else {
22                mask[pool_index] = static_cast<Dtype>(index);
23              }
24            }
25          }
26        }
27      }
28      // compute o set
29      bottom_data += bottom[0]->o set(0, 1);
30      top_data += top[0]->o set(0, 1);
31      if (use_top_mask) {
32        top_mask += top[0]->o set(0, 1);
33      } else {
34        mask += top[0]->o set(0, 1);
35      }
36    }
37  }

```

Listing 5: Caffe version: Max Pooling

```

1 template <typename T, unsigned int policy = phast::
2   get_default_policy()>
3 struct poolingMax : phast::functor::func_mat_mat<T, policy> {
4
5   _PHAST_METHOD poolingMax(int stride_h, int stride_w,
6     int pad_h, int pad_w,
7     int kernel_h, int kernel_w,
8     int height, int width) {
9     sh = stride_h; sw = stride_w;
10    ph = pad_h; pw = pad_w;
11    kh = kernel_h; kw = kernel_w;
12    h = height; w = width;
13  }
14  _PHAST_METHOD void operator()(phast::functor::matrix<T>&&
15    in, phast::functor::matrix<T>&& out) {
16    auto outIT = out.begin_ij();
17    int index = this->get_index();
18    for (int i = 0; i < out.size_i(); ++i) {
19      for (int j = 0; j < out.size_j(); ++j, ++outIT) {
20
21        int hstart = i * sh - ph;
22        int hend = smin(hstart + kh, h);
23        hstart = smax(hstart, 0);
24
25        int wstart = j * sw - pw;
26        int wend = smin(wstart + kw, w);
27        wstart = smax(wstart, 0);
28
29        T num = outIT;
30        for (int y = hstart; y < hend; ++y) {
31          for (int x = wstart; x < wend; ++x) {
32            if (num < in[y][x]) {
33              num = in[y][x];
34              mask.at(index, i, j) = y * w + x;
35            }
36          }
37        }
38        outIT = num;
39      }
40    }
41  }
42
43  phast::functor::cube<T> mask;
44
45  int sh, sw;
46  int ph, pw;
47  int kh, kw;
48  int h, w;
49 };

```

Listing 6: PHAST version: Max Pooling

4.3.3.2 Back-Propagation

During the feed-forward stage, we have stored from where we have been taking each output value, therefore we have to map the values from the output to the input using that list of mapped values. Its objective is to apply all the gradients to its corresponding positions.

Like in the feed-forward stage, we have not merged all the loops (in Listing 7), because we did not verified that it will continue working properly. Therefore, we have only parallelized the first loop (Listing 8).

```

1 for (int n = 0; n < top[0]->num(); ++n) {
2   for (int c = 0; c < channels_; ++c) {
3     for (int ph = 0; ph < pooled_height_; ++ph) {
4       for (int pw = 0; pw < pooled_width_; ++pw) {
5         const int index = ph * pooled_width_ + pw;
6         const int bottom_index =
7         use_top_mask ? (int)top_mask[index] : (int)mask[index];
8         bottom_di [bottom_index] += top_di [index];
9       }
10    }
11    bottom_di += bottom[0]->o set(0, 1);
12    top_di += top[0]->o set(0, 1);
13    if (use_top_mask) top_mask += top[0]->o set(0, 1);
14    else mask += top[0]->o set(0, 1);
15  }
16 }

```

Listing 7: Caffe version: Back Max Pooling

```

1 template <typename T, unsigned int policy = phast::
2   < get_default_policy() >
3 struct poolingMaxBack : phast::functor::func_mat_mat<T, policy>
4   {
5   _PHAST_METHOD poolingMaxBack(int pooled_h, int pooled_w,
6     int width) {
7     ph = pooled_h; pw = pooled_w;
8     w = width;
9   }
10  _PHAST_METHOD void operator()(phast::functor::matrix<T>&
11    bottom, phast::functor::matrix<T>& top) {
12    int channel = this->get_index();
13    for (int i = 0; i < ph; ++i) {
14      for (int j = 0; j < pw; ++j) {
15        int index = mask[channel][i][j];
16        int y = index / w; int x = index % w;
17        bottom[y][x] += top[i][j];
18      }
19    }
20  }
21  phast::functor::cube<T> mask;
22  int ph, pw;
23  int w;
24 };

```

Listing 8: PHAST version: Back Max Pooling

4.3.4 InnerProduct

4.3.4.1 Feed-Forward

In neural networks, there is a layer usually known as Perceptron layer (or Dense layer), sometimes is also known as InnerProduct, because the output of the layer is the inner product of the input with the weights.

```

1 template <typename Dtype>
2 void InnerProductLayer<Dtype>::Forward_cpu(const vector<Blob
3   < Dtype > >& bottom,
4   const vector<Blob<Dtype> >& top) {
5   const Dtype bottom_data = bottom[0]->cpu_data();
6   Dtype top_data = top[0]->mutable_cpu_data();
7   const Dtype weight = this->blobs_[0]->cpu_data();
8
9   ca_e_cpu_gemm<Dtype>(CblasNoTrans, transpose_ ?
10     CblasNoTrans : CblasTrans,
11     M_, N_, K_, (Dtype)1.,
12     bottom_data, weight, (Dtype)0., top_data);
13
14   if (bias_term_) {
15     ca_e_cpu_gemm<Dtype>(CblasNoTrans, CblasNoTrans, M_,
16     N_, 1, (Dtype)1.,
17     bias_multiplier_cpu_data(),
18     this->blobs_[1]->cpu_data(), (Dtype)1., top_data);
19 }

```

Listing 9: Caffe version: InnerProduct

```

1 template <typename T, unsigned int policy = phast::
2   < get_default_policy() >
3 struct matrixPlusVectorRows : phast::functor::func_vec<T, policy
4   > {
5   _PHAST_METHOD matrixPlusVectorRows() {}
6
7   _PHAST_METHOD void operator()(phast::functor::vector<T>&
8     row) {
9     for (auto r = row.begin(), i = vec.begin(); r != row.end(); ++r,
10      ++i)
11       r += *i;
12 }
13
14 template <>
15 void InnerProductLayer<float>::Forward_cpu(const vector<Blob<
16   float > >& bottom,
17   const vector<Blob<float> >& top) {
18   phast::matrix<float> matA = bottom[0]->getDataAsMatrix(M_,
19     K_, false);
20   phast::matrix<float> matB = this->blobs_[0]->
21     getDataAsMatrix(K_, N_, ltranspose_);
22   phast::matrix<float> matC = top[0]->getDataAsMatrix(M_, N_,
23     false);
24   phast::dot_product(matA, matB, matC);
25
26   if (bias_term_) {
27     matrixPlusVectorRows<float> matrixPlusVectorRows;
28     matrixPlusVectorRows.vec.link(this->blobs_[1]->
29     getDataAsVector(N_));
30     phast::for_each(matC.begin_i(), matC.end_i(),
31     matrixPlusVectorRows);
32   }
33   if (!ltranspose_) matB.transpose();
34 }

```

Listing 10: PHAST version: InnerProduct

It is interesting that, in Listing 9 in line 13, under the condition of `bias_term_`, there is a matrix multiplication, but in Listing 10 under the same condition in line 21, there is a call to `matrixPlusVectorRows`. This is a trick made very often by Caffe, its creators have mapped all possible operations to matrix multiplications to be easier to port them to GPU, but now, we can make specific functors for these operations.

4.3.4.2 Back-Propagation

The InnerProduct back-propagation is a bit different from the other block we have seen. Instead of reversing the operation made in feed-forward to map to each value its own gradient and modify its weight, we add to the weights a scaled gradient based on the original data, then the same with the bias, and lastly, we propagate the changes to the previous layer. Despite of the trick to map all operations to matrix multiplications, this layer is very straight forward (Listings 11 and 12).

```

1  template <typename Dtype>
2  void InnerProductLayer<Dtype>::Backward_cpu(const vector<Blob<Dtype> >& top,
3         const vector<Blob<Dtype> >& bottom) {
4      const vector<bool>& propagate_down,
5      const vector<Blob<Dtype> >& bottom) {
6      if (this->param_propagate_down_[0]) {
7          const Dtype top_di = top[0]->cpu_di ();
8          const Dtype bottom_data = bottom[0]->cpu_data ();
9          // Gradient with respect to weight
10         if (transpose_) {
11             caffe_cpu_gemm<Dtype>(CblasTrans, CblasNoTrans,
12                                 K_, N_, M_,
13                                 (Dtype)1., bottom_data, top_di,
14                                 (Dtype)1., this->blobs_[0]->mutable_cpu_di ());
15         } else {
16             caffe_cpu_gemm<Dtype>(CblasTrans, CblasNoTrans,
17                                 N_, K_, M_,
18                                 (Dtype)1., top_di, bottom_data,
19                                 (Dtype)1., this->blobs_[0]->mutable_cpu_di ());
20         }
21     }
22     if (bias_term_ && this->param_propagate_down_[1]) {
23         const Dtype top_di = top[0]->cpu_di ();
24         // Gradient with respect to bias
25         caffe_cpu_gemv<Dtype>(CblasTrans, M_, N_, (Dtype)1.,
26                             top_di,
27                             bias_multiplier_cpu_data(), (Dtype)1.,
28                             this->blobs_[1]->mutable_cpu_di ());
29     }
30     if (propagate_down[0]) {
31         const Dtype top_di = top[0]->cpu_di ();
32         // Gradient with respect to bottom data
33         caffe_cpu_gemm<Dtype>(CblasNoTrans, transpose_ ?
34                             CblasTrans : CblasNoTrans,
35                             M_, K_, N_,
36                             (Dtype)1., top_di, this->blobs_[0]->
37                             cpu_data(),
38                             (Dtype)0., bottom[0]->mutable_cpu_di
39                             ());
40     }
41 }

```

Listing 11: Caffe version: Back InnerProduct

```

1  template <>
2  void InnerProductLayer<float>::Backward_cpu(const vector<Blob
3         <float> >& top,
4         const vector<bool>& propagate_down,
5         const vector<Blob<float> >& bottom) {
6      if (this->param_propagate_down_[0]) {
7          if (transpose_) {
8              phast::matrix<float> di = top[0]->getDi AsMatrix(M_,
9                  N_, false);
10             phast::matrix<float> data = bottom[0]->getDataAsMatrix
11                 (K_, M_, true);
12             phast::matrix<float> tmp(K_, N_, 0);
13             phast::matrix<float> wdi = this->blobs_[0]->
14                 getDi AsMatrix(K_, N_, false);
15             phast::dot_product(data, di, tmp);
16             phast::transform(tmp.begin_ij(), tmp.end_ij(), wdi.begin_ij(),
17                 wdi.begin_ij(),
18                 phast::plus<float>());
19             data.transpose();
20         }
21         else {
22             phast::matrix<float> di = top[0]->getDi AsMatrix(N_,
23                 M_, true);
24             phast::matrix<float> data = bottom[0]->getDataAsMatrix
25                 (M_, K_, false);
26             phast::matrix<float> tmp(N_, K_, 0);
27             phast::matrix<float> wdi = this->blobs_[0]->
28                 getDi AsMatrix(N_, K_, false);
29             phast::dot_product(di, data, tmp);
30             phast::transform(tmp.begin_ij(), tmp.end_ij(), wdi.begin_ij(),
31                 wdi.begin_ij(),
32                 phast::plus<float>());
33             di.transpose();
34         }
35     }
36     if (bias_term_ && this->param_propagate_down_[1]) {
37         phast::matrix<float> di = top[0]->getDi AsMatrix(N_, M_,
38             true);
39         phast::vector<float> weig = this->blobs_[1]->
40             getDi AsVector(N_);
41         phast::transform(di.begin_ij(), di.end_ij(), weig.begin(),
42             reduceMatrixVectors<float>());
43         di.transpose();
44     }
45     if (propagate_down[0]) {
46         phast::matrix<float> di = top[0]->getDi AsMatrix(M_, N_,
47             false);
48         phast::matrix<float> weig = this->blobs_[0]->
49             getDataAsMatrix(N_, K_, transpose_);
50         phast::matrix<float> bdi = bottom[0]->getDi AsMatrix(M_,
51             K_, false);
52         phast::dot_product(di, weig, bdi);
53         if (transpose_) weig.transpose();
54     }
55 }

```

Listing 12: PHAST version: Back InnerProduct

4.3.5 ReLU

4.3.5.1 Feed-Forward

ReLU or Rectified-Linear layer applies the ReLU function to each of its inputs. It is mainly used to remove negative numbers from a neural network, but Caffe has merged it with the Leaky-ReLU. The difference between a Leaky-ReLU and a ReLU is in negative numbers, ReLU outputs 0s, and Leaky-ReLU outputs the input value multiplied by a constant (Equation 1 and Listings 13 and 14).

$$\text{Leaky-ReLU}(x) = \begin{cases} x & \text{if } x > 0 \\ ax & \text{otherwise} \end{cases} \quad (1)$$

```

1 for (int i = 0; i < count; ++i) {
2   top_data[i] = std::max(bottom_data[i], Dtype(0))
3     + negative_slope * std::min(bottom_data[i], Dtype(0));
4 }
    
```

Listing 13: Caffe version: ReLU

```

1 template <typename T, unsigned int policy = phast::
2   < get_default_policy() >
3 struct reluFunc : phast::func::func_scal_scal<T, policy> {
4   _PHAST_METHOD reluFunc(T negative) {
5     slope = negative;
6   }
7
8   _PHAST_METHOD void operator()(phast::func::scalar<T>&
9     & in, phast::func::scalar<T>& out) {
10    out = smax(in, T(0)) + slope * smin(in, T(0));
11  }
12  T slope;
13 };
    
```

Listing 14: PHAST version: ReLU

4.3.5.2 Back-Propagation

Using the gradient as x , and the original data as y , we have a very similar equation to make the back-propagation algorithm (Equation 2). Like before, the code is very straight forward (Listings 15 and 16).

$$\text{Back-Leaky-ReLU}(x, y) = \begin{cases} x & \text{if } y > 0 \\ a & \text{otherwise} \end{cases} \quad (2)$$

```

1 for (int i = 0; i < count; ++i) {
2   bottom_di [i] = top_di [i] * ((bottom_data[i] > 0)
3     + negative_slope * (bottom_data[i] <= 0));
4 }
    
```

Listing 15: Caffe version: Back ReLU

```

1 template <typename T, unsigned int policy = phast::
2   < get_default_policy() >
3 struct reluBackFunc : phast::func::func_scal_scal_scal<T, policy
4   > {
5   _PHAST_METHOD reluBackFunc(T negative) {
6     slope = negative;
7   }
8   _PHAST_METHOD void operator()(phast::func::scalar<T>&
9     & in, phast::func::scalar<T>& di , phast::func::
10    scalar<T>& out) {
11    int m0 = in > 0 ? 1 : 0;
12    int m1 = in <= 0 ? 1 : 0;
13    out = di * (m0 + slope * m1);
14  }
15  T slope;
16 };
    
```

Listing 16: PHAST version: Back ReLU

4.3.6 Softmax and Softmax Loss

4.3.6.1 Feed-Forward

In classifiers, a probability distribution over a discrete variable with n possible values, is represented using the softmax function [20]. Indeed, we require not only that each element to be between 0 and 1, but also that the entire vector sums to 1 so that it represents a valid probability distribution [20].

We are going to apply the softmax function (Equation 3) to each value of the output. As we work with many outputs at the same time, we apply it over each vector of the data matrix (Listing 17). It is worth to notice that, most of the code in the PHAST version (Listing 18) is data preparation, which will be improved in future versions of PHAST.

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad (3)$$

```

1  template <typename Dtype>
2  void SoftmaxLayer<Dtype>::Forward_cpu(const vector<Blob<
   < Dtype> >& bottom,
3   < const vector<Blob<Dtype> >& top) {
4   const Dtype bottom_data = bottom[0]->cpu_data();
5   Dtype top_data = top[0]->mutable_cpu_data();
6   Dtype scale_data = scale->mutable_cpu_data();
7   int channels = bottom[0]->shape(softmax_axis_);
8   int dim = bottom[0]->count() / outer_num_;
9   ca e_copy(bottom[0]->count(), bottom_data, top_data);
10  // We need to subtract the max to avoid numerical issues,
   < compute the exp,
11  // and then normalize.
12  for (int i = 0; i < outer_num_; ++i) {
13   // initialize scale_data to the first plane
14   ca e_copy(inner_num_, bottom_data + i * dim, scale_data);
15   for (int j = 0; j < channels; j++) {
16   for (int k = 0; k < inner_num_; k++) {
17   scale_data[k] = std::max(scale_data[k],
18   bottom_data[i * dim + j * inner_num_ + k]);
19   }
20  }
21  // subtraction
22  ca e_cpu_gemm<Dtype>(CblasNoTrans, CblasNoTrans,
   < channels, inner_num_,
23   1, -1., sum_multiplier_cpu_data(), scale_data, 1., top_data
   < );
24  // exponentiation
25  ca e_exp<Dtype>(dim, top_data, top_data);
26  // sum after exp
27  ca e_cpu_gemv<Dtype>(CblasTrans, channels, inner_num_, 1.,
28   top_data, sum_multiplier_cpu_data(), 0., scale_data);
29  // division
30  for (int j = 0; j < channels; j++) {
31   ca e_div(inner_num_, top_data, scale_data, top_data);
32   top_data += inner_num_;
33  }
34  }
35  }

```

Listing 17: Caffe version: SoftMax

```

1  template <>
2  void SoftmaxLayer<float>::Forward_cpu(const vector<Blob<float
   < > >& bottom,
3   < const vector<Blob<float> >& top) {
4   int channels = bottom[0]->shape(softmax_axis_);
5   phast::cube<float> data = bottom[0]->getDataAsCube(
   < outer_num_, channels, inner_num_);
6   phast::cube<float> out = top[0]->getDataAsCube(outer_num_,
   < channels, inner_num_);
7   phast::vector<float> scale = scale->getDataAsVector(inner_num_
   < );
8
9   phast::copy(data.begin_ijk(), data.end_ijk(), out.begin_ijk());
10
11  auto dataIt = data.begin_i();
12  auto outIt = out.begin_i();
13
14  for (; dataIt != data.end_i() && outIt != out.end_i(); dataIt++,
   < outIt++) {
15   phast::matrix<float> dataM;
16   phast::matrix<float> outM;
17   dataM.set_dev(data.size_j(), data.size_k(), dataIt.get_dev() +
   < dataIt.get_abs_pos());
18   outM.set_dev(out.size_j(), out.size_k(), outIt.get_dev() + outIt
   < get_abs_pos());
19
20   auto dataMIt = dataM.begin_i();
21
22   phast::vector<float> tmpRow;
23   tmpRow.set_dev(dataM.size_j(), dataMIt.get_dev() + dataMIt
   < get_abs_pos());
24   phast::copy(tmpRow.begin(), tmpRow.end(), scale.begin());
25
26   for (; dataMIt != dataM.end_i(); dataMIt++) {
27   phast::vector<float> dataV;
28   dataV.set_dev(dataM.size_j(), dataMIt.get_dev() + dataMIt
   < get_abs_pos());
29   phast::transform(dataV.begin(), dataV.end(), scale.begin(),
   < maxFunc<float>());
30  }
31
32  matrixMinusVectorRows<float> matrixMinusVectorRows;
33  matrixMinusVectorRows.vec.link(scale);
34  phast::for_each(outM.begin_i(), outM.end_i(),
   < matrixMinusVectorRows);
35  phast::for_each(outM.begin_ij(), outM.end_ij(), expFunc<float
   < >());
36  outM.transpose();
37  phast::transform(scale.begin(), scale.end(), outM.begin_i(),
   < matrixReduceByRows<float>());
38  outM.transpose();
39  matrixDivVectorRows<float> matrixDivVectorRows;
40  matrixDivVectorRows.vec.link(scale);
41  phast::for_each(outM.begin_i(), outM.end_i(),
   < matrixDivVectorRows);
42  }
43  }

```

Listing 18: PHAST version: SoftMax

4.3.6.2 Back-Propagation

Like in the InnerProduct layer, the objective of the back-propagation stage is not to reverse the feed-forward stage, but to scale the data to the previous layers in the network (Listings 19 and 20). Same as before, most of the code is the PHAST version is data preparation.

```

1 template <typename Dtype>
2 void SoftmaxLayer<Dtype>::Backward_cpu(const vector<Blob<
3     Dtype> >& top,
4     const vector<bool>& propagate_down,
5     const vector<Blob<Dtype> >& bottom) {
6     const Dtype top_di = top[0]->cpu_di ();
7     const Dtype top_data = top[0]->cpu_data();
8     Dtype bottom_di = bottom[0]->mutable_cpu_di ();
9     Dtype scale_data = scale_.mutable_cpu_data();
10    int channels = top[0]->shape(softmax_axis.);
11    int dim = top[0]->count() / outer_num.;
12    ca e_copy(top[0]->count(), top_di , bottom_di );
13    for (int i = 0; i < outer_num.; ++i) {
14        // compute dot(top_di , top_data) and subtract them from the
15        // bottom di
16        for (int k = 0; k < inner_num.; ++k) {
17            scale_data[k] = ca e_cpu_strided_dot<Dtype>(channels,
18                bottom_di + i * dim + k, inner_num.,
19                top_data + i * dim + k, inner_num.);
20        }
21        // subtraction
22        ca e_cpu_gemm<Dtype>(CblasNoTrans, CblasNoTrans,
23            channels, inner_num., 1,
24            -1., sum_multiplier_.cpu_data(), scale_data, 1., bottom_di
25            + i * dim);
26    }
27    // elementwise multiplication
28    ca e_mul(top[0]->count(), bottom_di , top_data, bottom_di );
29 }

```

Listing 19: Ca e version: Back SoftMax

```

1 template <>
2 void SoftmaxLayer<float>::Backward_cpu(const vector<Blob<
3     float> >& top,
4     const vector<bool>& propagate_down,
5     const vector<Blob<float> >& bottom) {
6     int channels = top[0]->shape(softmax_axis.);
7
8     phast::cube<float> topDi = top[0]->getDi AsCube(
9         outer_num., channels, inner_num.);
10    phast::cube<float> topData = top[0]->getDataAsCube(
11        outer_num., channels, inner_num.);
12    phast::cube<float> botDi = bottom[0]->getDi AsCube(
13        outer_num., channels, inner_num.);
14    phast::vector<float> scaler = scale_.getDataAsVector(
15        inner_num.);
16    phast::copy(topDi .begin_ijk(), topDi .end_ijk(), botDi .
17        begin_ijk());
18
19    auto botIt = botDi .begin_i();
20    auto topDIt = topData.begin_i();
21
22    for (; botIt != botDi .end_i() && topDIt != topData.end_i();
23        botIt++, topDIt++) {
24        phast::matrix<float> botM;
25        phast::matrix<float> topDM;
26        botM.set_dev(botDi .size_j(), botDi .size_k(), botIt.get_dev()
27            + botIt.get_abs_pos());
28        topDM.set_dev(topData.size_j(), topData.size_k(), topDIt.
29            get_dev() + topDIt.get_abs_pos());
30
31        botM.transpose();
32        topDM.transpose();
33
34        reduceMatrixVectorByVectorDot<float>
35            reduceMatrixVectorByVectorDot;
36        reduceMatrixVectorByVectorDot.scal.link(scaler);
37        phast::transform(botM.begin_i(), botM.end_i(), topDM.begin_i
38            (), reduceMatrixVectorByVectorDot);
39
40        botM.transpose();
41        topDM.transpose();
42
43        matrixMinusVectorRows<float> matrixMinusVectorRows;
44        matrixMinusVectorRows.vec.link(scaler);
45        phast::for_each(botM.begin_i(), botM.end_i(),
46            matrixMinusVectorRows);
47    }
48
49    phast::transform(botDi .begin_ijk(), botDi .end_ijk(), topData.
50        begin_ijk(), botDi .begin_ijk(), phast::multiplies<float>
51        (>0));

```

Listing 20: PHAST version: Back SoftMax

4.3.6.3 With Loss

Regarding the SoftMax with Loss block, it is the same as SoftMax but adding similar code to the Accuracy Layer into the SoftMax with Loss to calculate the logarithmic error instead of accuracy. And, like the Accuracy layer, it has not back-propagation except the SoftMax itself.

4.3.7 Accuracy

4.3.7.1 Feed-Forward

The Accuracy layer is implicit in Caffe, it checks the numbers of good predictions and generates the corresponding percentage (Listings 21 and 22).

```

1 for (int j = 0; j < inner_num_; ++j) {
2   const int label_value =
3     static_cast<int>(bottom_label[i + inner_num_ + j]);
4   if (has_ignore_label_ && label_value == ignore_label_) {
5     continue;
6   }
7   DCHECK_GE(label_value, 0);
8   DCHECK_LT(label_value, num_labels);
9   if (top.size() > 1) ++nums_bu_er_.mutable_cpu_data()[
10    label_value];
11   const Dtype prob_of_true_class = bottom_data[i dim
12     + label_value inner_num_ + j];
13   int num_better_predictions = -1; // true_class also counts as "
14     better"
15   // Top-k accuracy
16   for (int k = 0; k < num_labels && num_better_predictions <
17     top_k_; ++k) {
18     num_better_predictions +=
19       (bottom_data[i dim + k inner_num_ + j] >=
20        prob_of_true_class);
21   }
22   // check if there are less than top_k_ predictions
23   if (num_better_predictions < top_k_) {
24     ++accuracy;
25     if (top.size() > 1) ++top[1->mutable_cpu_data()[label_value
26     ]];
27   }
28   ++count;
29 }

```

Listing 21: Caffe version: Accuracy

```

1 template <typename T, unsigned int policy = phast::
2   get_default_policy()>
3 struct doAccuracy : phast::functor::func_vec<T, policy> {
4   _PHAST_METHOD doAccuracy(int label_, int maxLabels, int
5     ignore, int ignoreValue, int isTop, int top_k) {
6     label = label_;
7     labels = maxLabels;
8     hasIgnore = ignore;
9     value = ignoreValue;
10    top = isTop;
11    topk = top_k;
12  }
13  _PHAST_METHOD T operator()(phast::functor::vector<T>&
14    row) {
15    T acc = 0;
16    if (hasIgnore && label == value) return 0;
17    if (label < 0 // label > labels) return 0;
18    if (top > 1) num[label]++;
19    T prob = row[label];
20    int predicts = -1;
21    for (auto it = row.begin(); it != row.end() && predicts < topk
22      ; ++it) {
23      if (( *it ) >= prob) predicts++;
24    }
25    if (predicts < topk) {
26      ++acc;
27      if (top > 1) out1[label]++;
28    }
29    return acc;
30  }
31 }
32 phast::functor::vector<T> out1;
33 phast::functor::vector<T> num;
34
35 int label;
36 int labels;
37 int hasIgnore;
38 int value;
39 int top;
40 int topk;
41 };

```

Listing 22: PHAST version: Accuracy

4.3.7.2 Back-Propagation

The Accuracy layer does not have any kind of back-propagation because it does not modify the data.

4.3.8 Pending work

There is some functionality that we have not implemented due to lack of time. In the Accuracy block test, we have not implemented yet multiple axes checking, it means that our version will not check against multi-labeled output. Also, due to some PHAST limitations that are going to be solved soon, we cannot implement the ignore label, which allows ignoring some parts of the batch in the accuracy measurement. On the other hand, the Convolution block has been only implemented the simplest 2D convolution using a 2D im2col. Also, we discarded the usage of independent filters using groups until we get more familiarized with the PHAST library.

5. Results & Evaluation

5.1 Testbench

In this study, we have used the PHAST library 1.0.1 compiled with GCC 6.3.0 and Cuda 9.0. Caffe was obtained from the official git repository², the commit we have used is 99bd99795dcdf0b1d3086a8d67ab1782a8a08383. Our compute machine is part of GACOP’s computer cluster [21], it is running CentOS Linux 7.5 1804 with Linux 3.10.0-862.14.4, powered by two Intel(R) Xeon(R) CPU E5-2603 v3 @ 1.60GHz with 64 GiB RAM memory, and a Geforce GTX 1080 8GB GDDR5X with Nvidia driver 410.48.

With the original Caffe code, we have trained two neural networks from Caffe examples, both are LeNet-based networks. The first one is able to classify the MNIST³ database, the network is built from 6 layers: 2 Convolutions, 2 Pooling, and 2 InnerProduct. The second one is designed for CIFAR10⁴ database, it has 8 layers: 3 Convolutions, 3 Pooling, and 2 InnerProduct. Additionally, both networks have a SoftMax layer with loss, an Accuracy layer, and at least 1 ReLU function.

We have successfully run both neural networks in CPU and GPU using PHAST. We obtain the same results as the original CPU Caffe’s implementation. To check that the results are the same, we have used the set of inputs that Caffe provides, then we have checked several parameters such as the output of the network, the accuracy, the loss, and some intermediate matrices. Despite we have noticed that the accuracy and the loss are enough to validate the results, we continue using the intermediate matrices and the outputs to be sure that all was working properly. Now, we are able to choose between the CPU or the GPU version depending on the makefile we use.

5.2 Experiences

During the development, we have confirmed that applying a high performance programming model is not much different from the programming methodologies we are used to. There are some exceptions, such as the Convolution block and the Accuracy block, where rethinking the algorithm is necessary, but despite that, most of the high performance models are based in well-known approaches, as we have mentioned before.

5.3 Caffe Tests

Once our testbench is running as intended, we decide to test how accurate is our implementation to the original implementation. Therefore, using Caffe tests files, we have checked all the blocks we have ported to PHAST that have a test (Table 3). We notice that all the functionality we have ported is working, and tests only fail on the unimplemented functionality.

Block	Passed	Not Passed	Total	%Passed
Convolution	3	12	15	20
Pooling	11	0	11	100
InnerProduct	9	0	9	100
SoftMax	4	0	4	100
SoftMax Loss	4	0	4	100
Accuracy	9	3	12	75

Table 3: Each available test for the block we have modify in single precision floating point numbers and the number of tests passed and not passed.

²<https://github.com/BVLC/caffe>

³More information: <http://yann.lecun.com/exdb/mnist/>

⁴More information: <https://www.cs.toronto.edu/~kriiz/cifar.html>

5.4 Performance

Since this is an ongoing project, we are not ready to give performance results because our current tests are very small. But, compared to Ca e's original CPU code using openblas, our code was running 10 times slower. After some little changes, we improved it by a factor of 2x. But these are not final results and they can be improved by finding and fixing the current hotspots. Also, we are trying to reduce the number of copies made at each operation, by adopting part of the code to use read-only data, or maintain some structures as long as possible. It is like modernizing our modernized code.

6. Conclusions & Future Work

In this work, we aim to find the most relevant high-performance portable programming models available and classify them into a table to be able to easily use the one that fits our needs. But, to actually understand how these models work, we selected one of them, PHAST, and applied it to a real case. Since deep learning and neural networks are very important today, we selected one of the firsts deep learning framework (Caffe) and adapted it to use PHAST.

From our high performance portable programming model classification, we found that there are many available models, each one with its own characteristics. Nowadays, all is merged and making a strict classification is more subjective that it should be, but we think that we have stablish a baseline that could be used to classify the majority of the available models.

In this work, we have seen that PHAST is suitable to be used in real applications, it is easy to use and remembers well-established highly-expressive techniques.

There are several works left behind due to lack of time, among them we have:

- Explore the remainder programming models and devices: There are many models we have left behind, exploring them can open new characteristics in the classification table.
- Complete the port of Caffe to PHAST: We plan to release Caffe freely when the port is finished. We think that the information gathers during this port is very useful for improving PHAST.
- Extend our work to another device: Increasing the portability is one of our main goals in this work, another compatible device goes to give us valuable information when developing using these programming models.
- Improve the performance: Our firsts results show that checking the code it can be greatly improved with no so much effort.

Acknowledgments

We are glad to thank Sandro Bartolini and Biagio Peccerillo, from University of Siena, for granting us the demo version of the PHAST library and all the support given during this period, allowing us to develop this project as a joint collaboration.

References

- [1] Antonio Gonzalez. Trends in processor architecture. *CoRR*, abs/1801.05215, 2018.
- [2] John L. Hennessy and David A. Patterson. A new golden age for computer architecture. *Commun. ACM*, 62(2):48–60, January 2019.
- [3] Intel. Intel’s ‘one api’ project delivers unified programming model across diverse architectures. <https://newsroom.intel.com/news/intel-s-one-api-project-delivers-unified-programming-model>.
- [4] JR Neely. Doe centers of excellence performance portability meeting. Technical Report LLNL-TR-700962, Lawrence Livermore National Lab.(LLNL), Livermore, CA (United States), 2016.
- [5] Christian-A. Bohn. Kohonen feature mapping through graphics hardware. In *In Proceedings of Int. Conf. on Compu. Intelligence and Neurosciences*, pages 64–67, 1998.
- [6] Vaughn Betz. FPGA Architecture for the Challenge. http://www.eecg.toronto.edu/~vaughn/challenge/fpga_arch.html.
- [7] V. Antinyan, A. B. Sandberg, and M. Staron. A pragmatic view on code complexity management. *Computer*, 52(2):14–22, Feb 2019.
- [8] Christof Angermueller, Tanel Pärnamaa, Leopold Parts, and Oliver Stegle. Deep learning for computational biology. *Molecular Systems Biology*, 12(7):878, 2016.
- [9] W. G. Hatcher and W. Yu. A survey of deep learning: Platforms, applications and emerging research trends. *IEEE Access*, 6:24411–24432, 2018.
- [10] Henry Kasim, Verdi March, Rita Zhang, and Simon See. Survey on parallel programming model. In *Proceedings of the IFIP International Conference on Network and Parallel Computing*, NPC ’08, pages 266–275, Berlin, Heidelberg, 2008. Springer-Verlag.
- [11] Tim Lewis. OpenMP FAQ. <https://www.openmp.org/about/openmp-faq/>.
- [12] OpenACC.org organization. Homepage | OpenACC. <https://www.openacc.org/>.
- [13] Barcelona Supercomputing Center. The OmpSs Programming Model | Programming Models @ BSC. <https://pm.bsc.es/ompss>.
- [14] The Khronos Group Inc. Opencl overview. <https://www.khronos.org/opencl/>.
- [15] H. Carter Edwards, Christian R. Trott, and Daniel Sunderland. Kokkos: Enabling manycore performance portability through polymorphic memory access patterns. *Journal of Parallel and Distributed Computing*, 74(12):3202 – 3216, 2014. Domain-Specific Languages and High-Level Frameworks for High-Performance Computing.
- [16] Biagio Peccerillo and Sandro Bartolini. PHAST library - enabling single-source and high performance code for gpus and multi-cores. In *2017 International Conference on High Performance Computing & Simulation, HPCS 2017, Genoa, Italy, July 17-21, 2017*, pages 715–718, 2017.
- [17] Biagio Peccerillo, Sandro Bartolini, and Çetin Kaya Koç. Parallel bitsliced AES through PHAST: a single-source high-performance library for multi-cores and gpus. *J. Cryptographic Engineering*, 9(2):159–171, 2019.
- [18] Biagio Peccerillo and Sandro Bartolini. PHAST - A portable high-level modern C++ programming library for gpus and multi-cores. *IEEE Trans. Parallel Distrib. Syst.*, 30(1):174–189, 2019.

REFERENCES

- [19] Biagio Peccerillo and Sandro Bartolini. Task-dag support in single-source PHAST library: Enabling flexible assignment of tasks to cpus and gpus in heterogeneous architectures. In *Proceedings of the 10th International Workshop on Programming Models and Applications for Multicores and Manycores, PMAM@PPoPP 2019, Washington, DC, USA, February 17, 2019*, pages 91–100, 2019.
- [20] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [21] Eduardo José Gómez-Hernandez and José Manuel García. Preparing and managing and hpc cluster: Lessons learned. In *Jornadas Sarteco 2019, Cáceres (Spain)*, September 2019.

Acronyms

ANN Artificial Neural Network. [2](#)

API Application Programming Interface. [4](#), [5](#)

ASIC Application-Specific Integrated Circuit. [1](#), [2](#)

BSC Barcelona Supercomputing Center. [5](#)

DNN Deep Neural Network. [1](#), [10](#)

DSA Domain Specific Architecture. [1](#)

DSL Domain Specific Language. [2](#)

FPGA Field Programmable Gate Array. [1](#), [2](#), [6](#)

GeMM General Matrix Multiplication. [12](#), [14](#)

HPC High Performance Computing. [4](#), [5](#)

ICD Installable Client Driver. [5](#)

PGAS Partitioned Global Address Space. [7](#)

PHAST Parallel Heterogeneous-Architecture STL-like Template. [5](#), [6](#), [10–12](#), [19–22](#), [24](#)

SIMD Single Instruction Multiple Data. [2](#)

SIMT Single Instruction Multiple Threads. [2](#)

STL Standar Template Library. [5–7](#), [10](#)

