### GPU Acceleration of Automatic Differentiation in C++ with Clad

Ioana Ifrim, Princeton University

### Content

- Motivation 1.
- 2. Automatic Differentiation and applications ML Case Study
- 3. Clad.AD Plugin for Clang
- 4. C++ Compilation Pipeline
- 5. Clang Compilation Pipeline. Clad
- 6. GPU Accelerated AD
- 7. Clad & CUDA as a Service
- 8. Summary
- 9. Future Steps



### Motivation

In mathematics and computer algebra, automatic differentiation (AD) is defined as a set of techniques used for numerically evaluating the derivative of a function specified by a computer program. Automatic differentiation is an alternative technique to Symbolic differentiation and Numerical differentiation (the method of finite differences) and has applications ranging form the Machine Learning areal of domains to High Energy Physics.

The aim of Clad is to provide automatic differentiation for C/C++ which works without code modification (including legacy code)

The range of automatic differentiation (AD) application problems are defined by their high computational requirements and thus can greatly benefit from parallel implementations on graphics processing units (GPUs).







### Case Study : ML Application

In machine learning, we use gradient descent to update the parameters of our chosen model. A set of training inputs x are fed forward into the model generating corresponding activations yi. We define an error E as the difference computed between the data target output t and the model output y<sub>3</sub>. The error adjoint is propagated backward, resulting in the gradient with respect to the weights:

$$\nabla_{w_i} E = \left(\frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_6}\right)$$

This procedure is the central player of the gradient-descent optimisation algorithm. The required gradient is obtained by the backward propagation of the susceptibility of the objective value at the output, using the chain rule to compute partial derivatives of the objective wrt each of the weights. In this way, the resulting algorithm can be interpreted as transforming the network evaluation function composed with the objective function under reverse mode AD (generalisation of the back-propagation procedure)





## Case Study : ML Application

#### Manual Differentiation

It was historically the case that Machine Learning researchers would dedicate considerable amounts of time to the process of manual derivation of analytical derivatives which in turn were used in gradient descent procedures.



**Differentiation Methods** 

Ioana Ifrim - GPU Acceleration of Automatic Differentiation in C++ with Clad - CaaS Monthly Meeting - 06 May 2021

## Case Study: ML Application

#### Numerical Differentiation

Numerical differentiation is the finite difference approximation of derivatives using values of the original function evaluated at some sample points with  $e_i$  being i-th unit vector and h > 0 is a small step size

The introduction of round off errors bring forth issues of consistency, convergence, and stability of the numerical solution.

Moreover, the O(n) complexity of numerical differentiation for a gradient in n dimensions is the main obstacle to its usefulness in machine learning, where n can be as large as millions or billions in state-of-the-art deep learning models

 $\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}$ 

## Case Study: ML Application

#### Symbolic Differentiation

Symbolic differentiation is the automatic manipulation of expressions for obtaining derivative expressions and it is carried out by applying transformations representing rules of differentiation such as

$$\frac{d}{dx}(f(x) + g(x)) \rightsquigarrow \frac{d}{dx}f(x) + \frac{d}{dx}g(x)$$
$$\frac{d}{dx}(f(x)g(x)) \rightsquigarrow \left(\frac{d}{dx}f(x)\right)g(x) + f(x)\left(\frac{d}{dx}g(x)\right)$$

Symbolic differentiation can easily produce exponentially large symbolic expressions which take correspondingly long times to evaluate. This problem is known as expression swell. Moreover, it may require transcribing result back into code.



### Case Study : ML Application

#### Automatic Differentiation

Automatic generation of a C++ program able to compute the derivative of a given function

The procedure involves applying the chain rule of differential calculus throughout the semantics of the original program

Example Usage: Fitting a logistic regression model by minimising the binary cross-entropy loss of the logistic output

 $\mathcal{L} \circ \sigma(X, y, \beta) = p \log \sigma(X\beta) + (1-p) \log(1 - \sigma(X\beta))$ 

$$\nabla_{\beta} \mathcal{L} = p \nabla_{\beta} \log a(h) + (1 - p) \nabla_{\beta} \log(1 - a(h))$$

$$= p \operatorname{diag} \left( \frac{1}{a(h_i)} \right) \nabla_{\beta} a(h) + (1 - p) \operatorname{diag} \left( \frac{1}{1 - a(h_i)} \right) \nabla_{\beta} (1 - a(h))$$

$$= p \operatorname{diag} \left( \frac{1}{a(h_i)} \right) \nabla_{h} a(h) \nabla_{\beta} h(X, \beta)$$

$$- (1 - p) \operatorname{diag} \left( \frac{1}{1 - a(h_i)} \right) \nabla_{h} (1 - a(h)) \nabla_{\beta} h(X, \beta)$$



## Case Study: ML Application

#### Automatic Differentiation - Forward Mode

In forward mode auto differentiation, we start from the left-most node and move forward along to the right-most node in the computational graph – a forward pass

We calculate elementary derivatives using the expressions and leveraging the chain rule to obtain the intermediate derivatives at each step, obtaining the desired derivative with respect to the first variable. A forward pass is needed for each desired derivative - – one derivative with respect to each of the n input parameters.

#### $h(X,\beta), \nabla_{\beta}h(X,\beta) \to a(h), \nabla_{h}a(h), \to \log a(h), \nabla_{a}\log a(h)$

Derivative function created by the forward-mode AD is guaranteed to have at most a constant factor (around 2-3) more arithmetical operations compared to the original function. clad::differentiate(f, ARGS) takes 2 arguments:

- f is a pointer to a function or a method to be differentiated
- 2. ARGS is either:

a single numerical literal indicating an index of independent variable (e.g. 0 for x, 1 for y) a string literal with the name of independent variable (as stated in the *definition* of f, e.g. "x" or "y") Generated derivative function has the same signature as the original function f, however its return value is the value of the derivative.

**Forward Pass** 



**Deep Multi Layer Neural Network Forward Pass** 



### Case Study : ML Application

#### Automatic Differentiation - Reverse Mode

In reverse mode autodiff, we start from the outer-most node

Suppose that we are interested in calculating the log derivative with respect to a particular activation, we employ the chain rule

 $h(X, \beta) \rightarrow a(h) \rightarrow \log a(h)$ 

 $\nabla \log a(h) \rightarrow \nabla a(h) \rightarrow \nabla h(X, \beta)$ 

Machine learning tasks involve a large number of feature space parameters who are to be tuned, thus reverse mode AD fits perfectly the task of calculating the derivatives of the cost function wrt model parameters

```
auto f_grad = clad::gradient(f);
     double result1[2] = {};
     f_grad.execute(x, y, result1);
     std::cout << ''dx: '' << result1[0] << '' ' << '' dy: '' << result1[1] << std::endl;</pre>
 5
     auto f_dx_dy = clad::gradient(f, 'x, y'); // same effect as before
 6
 7
     auto f_dy_dx = clad::gradient(f, ''y, x'');
     double result 2[2] = \{\};
      f_dy_dx.execute(x, y, result2);
10
     //·note·that·the·derivatives·are·mapped·to·the·"result"·indices·in·the·same·
     order.as.they.were.specified.in.the.argument:
     std::cout << ''dy: '' << result2[0] << ' ' ' << ''dx: '' << result2[1] << std::endl;</pre>
12
```

Ioana Ifrim - GPU Acceleration of Automatic Differentiation in C++ with Clad - CaaS Monthly Meeting - 06 May 2021

Backpropagation





# Clad.AD Plugin for Clang

- Supports derivatives (partial and higher order), gradients, hessians and jacobians.
- Provides low-level derivative access primitives
- Allows embedding in frameworks

Clad is a compiler plugin extending Clang able to produce derivatives in both forward and reverse mode:

11

### Typical C++ Compilation Pipeline





V. Vassilev, L. Moneta

Automatic Differentiation in C++ with clad. Integration in ROOT

Ioana Ifrim - GPU Acceleration of Automatic Differentiation in C++ with Clad - CaaS Monthly Meeting - 06 May 2021

94th ROOT PPP Meeting



### Clang Compilation Pipeline. Clad



Clad — Automatic Differentiation for C++ Using Clang V. Vassilev

Mini-Workshop: Differentiable Programming for High-Performance, Data-Intensive Computations



## **GPU Accelerated AD**

Considering our Machine Learning case study, the goal is to compute the gradient of the cost function with respect to a transformation parameter vector x.

From an AD perspective, this can be done either by invoking the forward mode derivative once for every dimension in the parameters space or by a single pass of the reverse mode derivative.

These passes are bounded by access to and computations performed on the transformation parameters, hence this process is an excellent candidate for acceleration through GPU support implementation.

Tasks featuring heavy computations increase their time consumption proportional with the data sets magnitude. These applications can thus profit from the usage of threads and in this sense GPU acceleration brings a new layer of optimisation and a proportional speed up.



Cuda Thread

#### GPU Accelerated AD **Original Function**

The CUDA support for Clad includes extensions that allow one to execute functions on the GPU using many threads in parallel.

Function attributes cloning has been introduced for <u>device</u> \_\_host\_\_ to be carried forward in the Clad gradient definition

\_\_host\_\_ declaration as well as previous dependencies on the standard library functionalities not supported by CUDA, were reimplemented. (falling back on Thrust, the template library for CUDA based on the Standard Template Library (STL) has proved to be an issue in the context of Clang compilation)

Clad uses Tape Records for the execution that is replayed such that the gradient is produced in one pass - this also required extensions for the CUDA context and removal of dependencies on the standard library.

- \_\_\_device\_\_\_\_host\_\_\_\_double exponential\_pdf(double x, double lambda, double x0){
- $\cdots$  if ((x-x0) < 0)
- ·····return·0.0:
- ....return lambda \* std::exp (-lambda \* (x-x0));

#### Clad Gradient

1 void **exponential\_pdf\_grad**(double x, double lambda, double x0, double \*\_result) \_\_attribute\_\_((device)) · \_\_attribute\_\_((host)) · { ....bool.\_cond0; ....double.\_t0; ....double.\_t1; ....double.\_t2; ....double.\_t3; double t4;  $\cdots$ \_cond0 = (x - x0) < 0; $\cdots$  if  $(\_cond0) \cdot {$ ....double.exponential\_pdf\_return = 0.; .....goto.\_label0; 11 12 · · · · } 13 ....\_t1.=.lambda;  $\cdots t3 = -lambda;$ 15  $\cdot \cdot \cdot \cdot _t 2 \cdot = \cdot (x \cdot - \cdot x 0);$ 16  $\cdots t4 = t3 + t2;$ 17  $\cdots$ \_t0 = std::exp(\_t4); 18 ....double exponential\_pdf\_return = .\_t1 \* .\_t0; 19 ....goto \_label1; 20 -\_label1: 21 • • • • { 22  $\cdots$  double r0 = 1 \* t0;23  $\cdots$  result [1UL]  $\cdot$  +=  $\cdot$  r0; 24  $\cdots$  double r1 = 1 + 1; 25 ....double.\_r2.=.\_r1.\*.custom\_derivatives::exp\_darg0(\_t4); 26  $\cdots$  double r3 = r2 \* t2;27  $\cdots$  result [1UL] += -\_r3;  $\cdots$  double r4 = t3 \* r2;28  $\cdots$  result [0UL]  $\cdot$  +=  $\cdot$  r4; 29 30  $\cdots$  result [2UL]  $+=-_r4;$ 31  $\cdots$ 32 ····if·(\_cond0) ····· label0: 33 34 . . . . . . . . ; 35



## GPU Accelerated AD



- Benchmark showcases how using CUDA can influence the overall AD performance in computation of a gauss gradient with different dimensions
- GPU : Tesla P100-PCIE-16GB
- CPU : Intel(R) Xeon(R) Gold 6148 CPU @ 2.40GHz

Ioana Ifrim - GPU Acceleration of Automatic Differentiation in C++ with Clad - CaaS Monthly Meeting - 06 May 2021

#### **Original Function**

\_\_device\_\_\_host\_\_.double.sum(double\*.p,.int.dim).{  $\cdots$  double r = 0.0; $\cdot \cdot \cdot \cdot \cdot \mathbf{for} \cdot (\mathbf{int} \cdot \mathbf{i} = \cdot \mathbf{0}; \cdot \mathbf{i} < \cdot \mathbf{dim}; \cdot \mathbf{i} + +)$  $\cdots$   $r \cdot += \cdot p[i];$ return r;

#### **Clad Gradient**

1	<pre>void sum_grad(double * p, int dim, double * result) _ attribute_((device))</pre>
	<pre>attribute((host)) {</pre>
2	$\cdots$ double $d_r = 0;$
3	<pre>unsigned.longt0;</pre>
4	$\cdots$ int $d_i = 0;$
5	<pre>clad::tape<int>t1.=.{};</int></pre>
6	$\cdots$ double $r = 0.;$
7	$\cdot \cdot \cdot \cdot \_t0 \cdot = \cdot 0;$
8	$\cdots$ for (int $i = 0; i < dim; i++)$ {
9	•••••_t0++;
10	$\cdots \cdot r + = \cdot p[clad::push(_t1, \cdot i)];$
11	$\cdots$
12	$\cdots$ double sum_return = r;
13	····goto·_label0;
14	label0:
15	$\cdots d_r + = 1;$
16	<pre>for (;t0;t0) {</pre>
17	$\cdots$ double $r_d0 = d_r;$
18	$\cdots d_r + = c_r d0;$
19	<pre>result[clad::pop(_t1)] += _ r_d0;</pre>
20	$\cdots d_r -= c_r d0;$
21	$\cdots$

#### More benchmarks coming soon \*

22 }



## GPU Accelerated AD

```
#include."clad/Differentiator/Differentiator.h"
     #define • N • 10485760
     typedef void(*func) (double *x, double *p, double sigma, int dim, double
     *_result);
      ___device____host___double gaus(double* x, double* p, double sigma, int dim) {
     \cdot \cdot \cdot double \cdot t \cdot = \cdot 0;
      \cdots for ( int i = 0; i < dim; i++)
10
      \cdots + (x[i] - p[i]) + (x[i] - p[i]);
11
12
      ...t.=.-t./.(2*sigma*sigma);
      ...return.std::pow(2*M_PI, -dim/2.0) * std::pow(sigma, -0.5) * std::exp(t);
13
14
15
      __device___host___void gaus_grad(double *x, double *p, double sigma, int dim,
16
                                                                                             body generated by Clad
     double *_result);
17
18
     auto gaus_g = clad::gradient(gaus);
19
20
      ___device___func p_pow2 = sum_grad;
21
                                                                                             device function pointer
```



17

### Clad & CUDA as a Service



The demo shows cling usage of clad as a plugin to produce a derivative on the fly and send it to a CUDA kernel for execution

Ioana Ifrim - GPU Acceleration of Automatic Differentiation in C++ with Clad - CaaS Monthly Meeting - 06 May 2021



## Summary

- The generated Clad derivatives are now supported for computations on CUDA kernels thus allowing for further optimisation
- Scheduling still requires a certain degree of user input which we aim to further automize
- on GPUs.
- Challenges in terms of:
  - CUDA based on the Standard Template Library (STL)) due to compatibility issues with Clang) (fixed)
  - Passing the gradient function by pointer when compiling with Clang (fixed) / Cling (wip)

Clad can now handle a hybrid GPU/CPU setup, where the generation is currently done on the CPU, while the execution can be parallelised

- CUDA version Clang & Cling compatibility can be observed in implementation choices (e.g. not using Thrust (C++ template library for





# Future Steps

- Full support of arrays
- optimisation
- optimisation and global memory constraints in mind

Enable AD support for second order derivatives for HEP analysis through ROOT (data analysis software package) via Clad : GSoC Project

In the interactive case - provide an api which synthesises the CUDA kernel automatically and alleviates the user's need to deal with memory

- Currently the scheduling procedure requires a certain degree of user input to make it suitable for a hybrid CPU/GPU setup. A future aim is to fully automate this last step for complete CUDA integration, where the full toolchain process needs to be formalised with both scheduling





#### Thank you!

Ioana Ifrim - GPU Acceleration of Automatic Differentiation in C++ with Clad - CaaS Monthly Meeting - 06 May 2021

