

# Scaling an Augmented RISC-V Processor Design with High-Level Synthesis

Fourth International workshop on RISC-V for HPC

Johannes Schoder, Martin Bücker

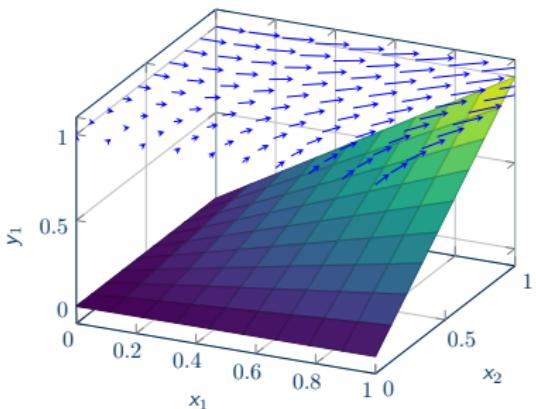
Friedrich Schiller University Jena, Department of Computer Science

Hamburg, 16th May 2024

# Motivation

# What we want in the long run:

Accelerate and Embed Automatic Differentiation into a Domain Specific RISC-V Processor.  
Work in progress..



$$y = f(x_1, x_2) = x_1 \cdot x_2$$
$$\nabla f(x_1, x_2) = \nabla y = \begin{pmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \end{pmatrix} = \begin{pmatrix} x_2 \\ x_1 \end{pmatrix}$$

We want: The gradient of a differentiable function without having to work for it.

- RISC-V - open source, expandable ISA
- HLS - High-Level Synthesis: fast prototyping
- FPGA - scalable and reconfigurable
- AD - Automatic Differentiation - program transformation in hardware

# Why do we want that?

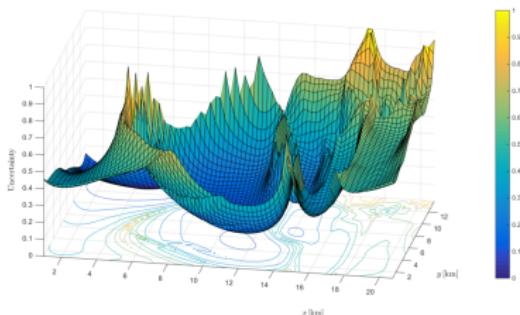


Figure: MeProRisk II: Optimization strategies and risk analysis for deep geothermal reservoirs

<https://www.ac.uni-jena.de/research/project/meprorisk>

We need the information about the gradient for, e.g.,

- Simulations
  - solve ODEs
  - solve nonlinear systems (newton method)
- Optimization
  - parameter estimation
  - minimize cost functions
  - train neural networks
- Examples:
  - physics simulation (shape optimization)

# Automatic differentiation

Example  $f: \mathbb{R}^3 \rightarrow \mathbb{R}^2$ :

- Typically done with software tools
- Two distinct modes:
  - "reverse mode automatic differentiation (AD)"
  - forward mode AD
- This study focuses only on forward mode AD

$$y_1(x_1, x_2, x_3) = \frac{\sin(x_1 \cdot x_2)}{\exp(x_2 - x_3)}$$
$$y_2(x_2, x_3) = \exp(x_2 - x_3) - 2 \cdot x_3$$

- Forward mode:  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  preferred if  $n \ll m$
- Reverse mode:  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  preferred if  $n \gg m$

# Forward Mode AD

# C Code

$$y_1 = f(x_1, x_2) = (a - x_1)^2 + b(x_2 - x_1^2)^2$$

```
1 int Rosenbrock (int x_1, int x_2){  
2     int a = 1;  
3     int b = 10;  
4     int result=((a-x_1)*(a-x_1))+(b*(x_2-x_1*x_1)*(x_2-x_1*x_1));  
5     return result;  
6 }
```

Listing 1: C code example for the Rosenbrock function

# ASM Code

```
1 Rosenbrock: #a0 = x_1; a1 = x_2 (procedure call standard)
2 li      a2, 1      #a2 = a
3 sub    a2, a2, a0  #a2 = a-x_1
4 mul    a2, a2, a2  #a2 = (a-x_1)^2
5 mul    a0, a0, a0  #a0 = x_1^2
6 sub    a0, a1, a0  #a0 = (x_2-x_1^2)
7 mul    a0, a0, a0  #a0 = (x_2-x_1^2)^2
8 li      a1, 10     #a1 = b
9 mul    a0, a0, a1  #a0 = b*(x_2-x_1^2)^2
10 add   a0, a0, a2  #a0 = (a-x_1)^2+b*(x_2-x_1^2)^2
11 ret          #return a0 (procedure call standard)
```

Listing 2: RV32IM assembly code example for the Rosenbrock function

Step 1: Decomposing complex mathematical functions in "elementary functions"  
- (SUB, ADD, MUL)

# Extended Evaluation Procedure

$$\begin{array}{l}
 v_{-1} = x_1 \\
 v_0 = x_2 \\
 \hline
 v_1 = a - v_{-1} \\
 v_2 = v_1 \cdot v_1 = v_1^2 \\
 v_3 = v_{-1} \cdot v_{-1} = v_{-1}^2 \\
 v_4 = v_0 - v_3 \\
 v_5 = v_4 \cdot v_4 = v_4^2 \\
 v_6 = b \cdot v_5 \\
 v_7 = v_2 + v_6 \\
 \hline
 y_1 = v_7
 \end{array}
 \quad
 \left| \begin{array}{l}
 \overrightarrow{\Delta v_{-1}} = \overrightarrow{\Delta x_1} \\
 \overrightarrow{\Delta v_0} = \overrightarrow{\Delta x_2} \\
 \hline
 \overrightarrow{\Delta v_1} = \overrightarrow{\Delta v_{-1}} \cdot \frac{\partial v_1}{\partial v_{-1}} = -\overrightarrow{\Delta v_{-1}} \\
 \overrightarrow{\Delta v_2} = \overrightarrow{\Delta v_1} \cdot \frac{\partial v_2}{\partial v_1} = 2 \cdot \overrightarrow{\Delta v_1} \cdot v_1 \\
 \overrightarrow{\Delta v_3} = \overrightarrow{\Delta v_{-1}} \cdot \frac{\partial v_3}{\partial v_{-1}} = 2 \cdot \overrightarrow{\Delta v_{-1}} \cdot v_{-1} \\
 \overrightarrow{\Delta v_4} = \overrightarrow{\Delta v_0} \cdot \frac{\partial v_4}{\partial v_0} + \overrightarrow{\Delta v_3} \cdot \frac{\partial v_4}{\partial v_3} = \overrightarrow{\Delta v_0} - \overrightarrow{\Delta v_3} \\
 \overrightarrow{\Delta v_5} = \overrightarrow{\Delta v_4} \cdot \frac{\partial v_5}{\partial v_4} = 2 \cdot \overrightarrow{\Delta v_4} \cdot v_4 \\
 \overrightarrow{\Delta v_6} = \overrightarrow{\Delta v_5} \cdot \frac{\partial v_6}{\partial v_5} = b \cdot \overrightarrow{\Delta v_5} \\
 \overrightarrow{\Delta v_7} = \overrightarrow{\Delta v_2} \cdot \frac{\partial v_7}{\partial v_2} + \overrightarrow{\Delta v_6} \cdot \frac{\partial v_7}{\partial v_6} = \overrightarrow{\Delta v_2} + \overrightarrow{\Delta v_6} \\
 \hline
 \overrightarrow{\Delta y_1} = \overrightarrow{\Delta v_7}
 \end{array} \right. \quad (1)$$

# Extended Evaluation Procedure

$$\begin{aligned} v_{i-n} &= x_i & i = 1 \dots n \\ \overrightarrow{\Delta v_{i-n}} &= \overrightarrow{\Delta x_i} & i = 1 \dots n \end{aligned}$$

---

$$\begin{aligned} v_i &= \phi_i(v_j)_{j \prec i} & i = 1 \dots l \\ \overrightarrow{\Delta v_i} &= \sum_{j \prec i} \frac{\partial}{\partial v_j} \phi_i(v_j) \cdot \overrightarrow{\Delta v_j} & i = 1 \dots l \end{aligned} \tag{2}$$

---

$$\begin{aligned} y_{m-i} &= v_{l-i} & i = m - 1 \dots 0 \\ \overrightarrow{\Delta y_{m-i}} &= \overrightarrow{\Delta v_{l-i}} & i = m - 1 \dots 0 \end{aligned}$$

Automatic Differentiation [1]:

---

[1] Andreas Griewank and Andrea Walther: Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation, <https://doi.org/10.1137/1.9780898717761>

# Extended Evaluation Procedure

$$v_{-1} = x_1 = 0$$

$$v_0 = x_2 = -0.5$$

...

$$y_1 = 3.5$$

$$\overrightarrow{\Delta v_{-1}} = \overrightarrow{\Delta x_1} = 1$$

$$\overrightarrow{\Delta v_0} = \overrightarrow{\Delta x_2} = 0$$

...

$$\overrightarrow{\Delta y_1} = -2$$

with  $\overrightarrow{\Delta x_1}, \overrightarrow{\Delta x_2} \in \mathbb{R}^1$

results in a single scalar entry of the gradient:  $\overrightarrow{\Delta y_1} = \frac{\partial y_1}{\partial x_1}$

$\overrightarrow{\Delta y_1} = \frac{\partial y_1}{\partial x_1} = -2$  partial derivative of  $y_1$  wrt.  $x_1$

- because  $x_1, x_2, \overrightarrow{\Delta x_1}$ , and  $\overrightarrow{\Delta x_2}$  are independent of each other, and these values aren't overwritten in the extended evaluation procedure, we can also say, that  $\overrightarrow{\Delta y_1}$  contains the total derivative of  $y_1$  wrt.  $x_1$ ;  $\frac{dy_1}{dx_1}$

# Seed matrix $S$

$$\begin{array}{l|l}
 \begin{array}{l}
 v_{-1} = x_1 = 0 \\
 v_0 = x_2 = -0.5 \\
 \hline
 \cdots \\
 \hline
 y_1 = 3.5
 \end{array}
 &
 \begin{array}{l}
 \overrightarrow{\Delta v_{-1}} = \overrightarrow{\Delta x_1} = (\textcolor{green}{1} \quad \textcolor{orange}{0}) \\
 \overrightarrow{\Delta v_0} = \overrightarrow{\Delta x_2} = (\textcolor{green}{0} \quad \textcolor{orange}{1}) \\
 \hline
 \cdots \\
 \hline
 \overrightarrow{\Delta y_1} = (-2 \quad -10)
 \end{array}
 \end{array} \tag{4}$$

results in  $\overrightarrow{\Delta y_1} = \left( \frac{\partial y_1}{\partial x_1} \quad \frac{\partial y_1}{\partial x_2} \right) \Big|_{0, -0.5}$

$p = 2$

# Seeding - Gradient

Seeding  $S \in \mathbb{R}^{(n \times p)}$  with  $p = n$ , and  $S = I_n$  (Identity)

$$S = \begin{pmatrix} \overrightarrow{\Delta x_1(1)} & \overrightarrow{\Delta x_1(2)} & \dots & \overrightarrow{\Delta x_1(p)} \\ \overrightarrow{\Delta x_2(1)} & \overrightarrow{\Delta x_2(2)} & \dots & \overrightarrow{\Delta x_2(p)} \\ \vdots & & \ddots & \vdots \\ \overrightarrow{\Delta x_n(1)} & \dots & \dots & \overrightarrow{\Delta x_n(p)} \end{pmatrix} \quad \left| \quad S = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \dots & \dots & 1 \end{pmatrix} \in \mathbb{R}^{(n \times n)} \right.$$
(5)

will result in the gradient at a specific position:

$$\overrightarrow{\Delta y_1} \Big|_{(x_1, x_2, \dots, x_n)} = \nabla f(x_1, x_2, \dots, x_n) = \left( \frac{\partial y_1}{\partial x_1} \quad \frac{\partial y_1}{\partial x_2} \quad \dots \quad \frac{\partial y_1}{\partial x_n} \right) \Big|_{(x_1, x_2, \dots, x_n)}$$

# HLS Processor Design

# Starting Point

## Reminder: what we want?

We want information on the gradient for free by performing the program transformation in hardware

## How can we get it?

Using an augmented RISC-V processor design based on existing RISC-V processor designs written for HLS.

Thanks to the DALI-LIRMM team and Prof. Goossens<sup>[2]</sup>

---

<sup>[2]</sup>Bernard Goossens - Guide to computer processor architecture : a RISC-V approach, with high-level synthesis, <https://doi.org/10.1007/978-3-031-18023-1>

# AD operations RV32I ISA - Changes to ALU

RISC-V Inst.	RISC-V ISA	Arithmetic Operation $v_i = \phi_i(v_j)_{j < i}$	Simultaneous Transformed Operation $\vec{\Delta v}_i = \sum_{j < i} \frac{\partial}{\partial v_j} \phi_i(v_j) \cdot \vec{\Delta v}_j$
ADD , ADDI	I	$v_i = v_{j_1} + v_{j_2}$	$\vec{\Delta v}_i = \vec{\Delta v}_{j_1} + \vec{\Delta v}_{j_2}$
SUB , SUBI	I	$v_i = v_{j_1} - v_{j_2}$	$\vec{\Delta v}_i = \vec{\Delta v}_{j_1} - \vec{\Delta v}_{j_2}$
SLL , SLLI	I	$v_i = v_{j_1} \cdot 2^{v_{j_2}}$	$\vec{\Delta v}_i = \vec{\Delta v}_{j_1} \cdot 2^{v_{j_2}} + \vec{\Delta v}_{j_2} \cdot v_{j_1} \cdot 2^{v_{j_2}} \cdot \log_e 2$
SRA , SRAI , (SRL , SRLI)	I	$v_i = \frac{v_{j_1}}{2^{v_{j_2}}}$	$\vec{\Delta v}_i = \vec{\Delta v}_{j_1} \cdot \frac{1}{2^{v_{j_2}}} - \vec{\Delta v}_{j_2} \cdot v_{j_1} \cdot \frac{1}{(2^{v_{j_2}})} \cdot \log_e 2$

# Augmented Processor Design

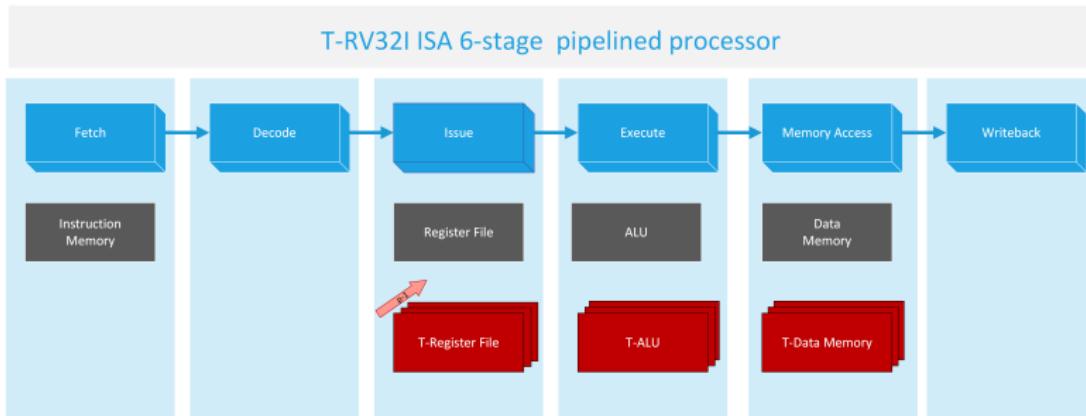
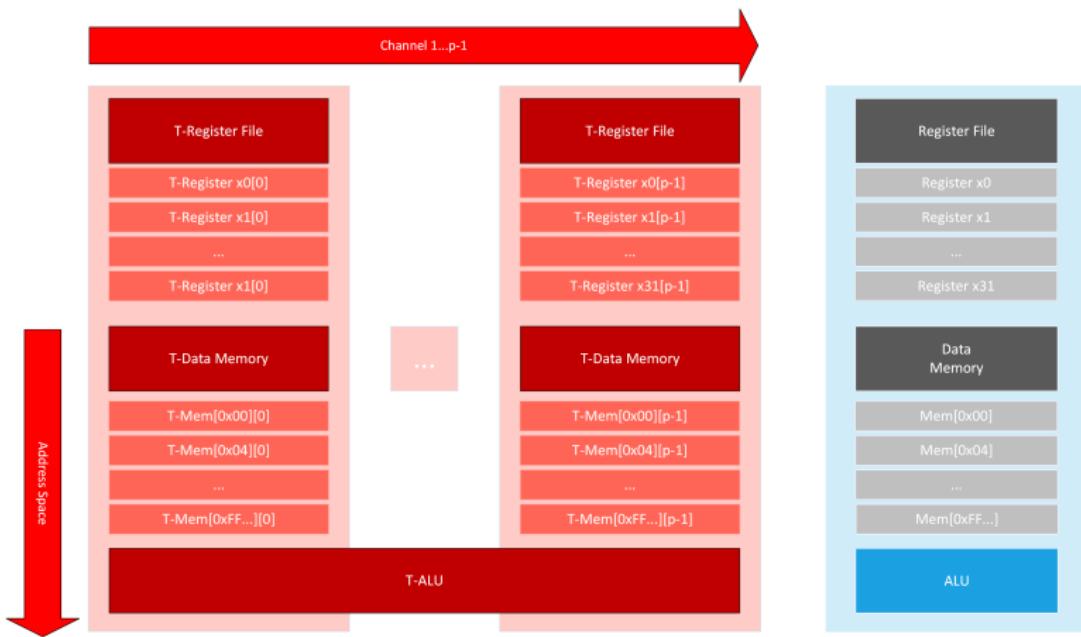


Figure: Augmented Processor Design

# Memory Layout - "Channels"



# Custom instructions

Added custom RV32-Instructions for:

1. Setting seed in Memory TSW
2. Loading gradient information from AD-Memory to Register TLW
3. Copying gradient information from AD-Register to Register TMV

# Alveo U50DD



Figure: Alveo U50 <sup>a</sup>

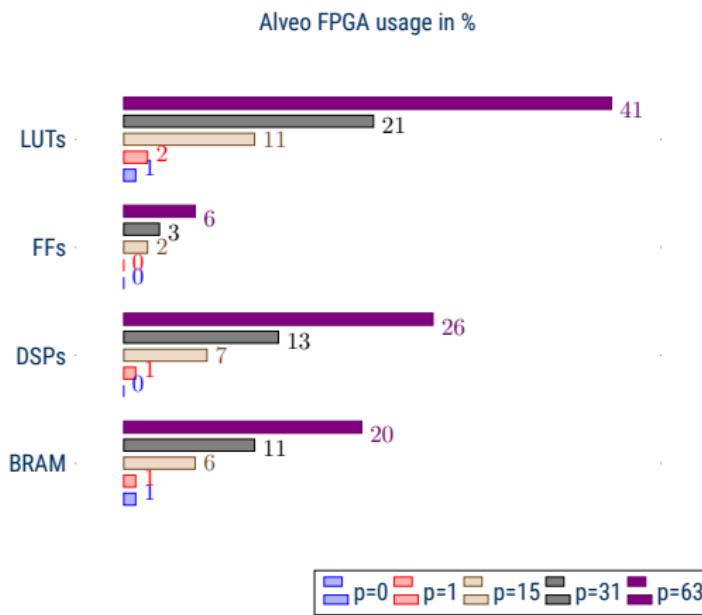
## UltraScale+™ FPGA

- 8 GB of high-bandwidth memory (HBM2)
- LUTs: 872K
- FFs (Registers): 1,743K
- DSPs: 5,952
- BRAM (36Kb): 1344 (47.3 Mb)
- UltraRAM (288 Kb): 640 (180.0 Mb)

---

<sup>a</sup><https://www.xilinx.com/products/boards-and-kits/alveo/u50.html>

# HLS (T-RV32I) for Alveo U50DD



## Memory size:

- I-RAM:  $2^{13}$  Bytes = 8KB
- D-RAM:  $2^{14}$  Bytes = 16KB (+16KB\*p)

## Estimated Clock:

- 0: 2.97ns
- 1: 8.677ns
- 15: 11.626ns
- 31: 14.481ns
- 63: 14.590ns

# Code size reduction:

## unaugmented RISC-V processor

```
1 #ad_v2_0 = ad_v1_0 + ad_v0_0;
2 lw a4,-44($0)
3 lw a5,-36($0)
4 add a5,a4,a5
5 sw a5,-52($0)
6 #ad_v2_1 = ad_v1_1 + ad_v0_1;
7 lw a4,-48($0)
8 lw a5,-40($0)
9 add a5,a4,a5
10 sw a5,-56($0)
11 #v2 = v1 + v0;
12 lw a4,-28($0)
13 lw a5,-24($0)
14 add a5,a4,a5
15 sw a5,-32($0)
```

## augmented RISC-V processor

```
1 #v2 = v1 + v0;
2 lw a4,-32($0)
3 lw a5,-28($0)
4 add a5,a4,a5
5 sw a5,-36($0)
```

# Speedup for one arithmetic operation

$$T_c^{\text{orig}}(p) = (c \cdot p + 1) \cdot \text{CPI} \cdot T_c^{\text{orig}}$$

$$T_c^{\text{aug}}(p) = 1 \cdot \text{CPI} \cdot T_c^{\text{aug}}(p)$$

$$S(p) := \frac{T_c^{\text{orig}}(p)}{T_c^{\text{aug}}(p)} = (c \cdot p + 1) \cdot \frac{T_c^{\text{orig}}}{T_c^{\text{aug}}(p)}$$

$$S(63) \approx (c \cdot 63 + 1) \cdot \frac{2.97\text{ns}}{14.590\text{ns}} \approx (c \cdot 63 + 1) \cdot 0.2 \geq 12.8$$

$p$  = channels,  $c$  = software "overhead" of non-augmented design,

$T_c^{\text{orig}}$  = clock cycle time of the original processor design,

$T_c^{\text{aug}}(p)$  = clock cycle time of the augmented processor design with  $p$  channels,

CPI = equals the clock cycles per instruction

# AD operations RV32I ISA - Changes to ALU

RISC-V Inst.	RISC-V ISA	Arithmetic Operation $v_i = \phi_i(v_j)_{j < i}$	Simultaneous AD-Operation $\overrightarrow{\Delta v_i} = \sum_{j < i} \frac{\partial}{\partial v_j} \phi_i(v_j) \cdot \overrightarrow{\Delta v_j}$
ADD , ADDI	I	$v_i = v_{j_1} + v_{j_2}$	$\overrightarrow{\Delta v_i} = \overrightarrow{\Delta v_{j_1}} + \overrightarrow{\Delta v_{j_2}}$
SUB , SUBI	I	$v_i = v_{j_1} - v_{j_2}$	$\overrightarrow{\Delta v_i} = \overrightarrow{\Delta v_{j_1}} - \overrightarrow{\Delta v_{j_2}}$
SLL , SLLI	I	$v_i = v_{j_1} \cdot 2^{v_{j_2}}$	$\overrightarrow{\Delta v_i} = \overrightarrow{\Delta v_{j_1}} \cdot 2^{v_{j_2}} + \overrightarrow{\Delta v_{j_2}} \cdot v_{j_1} \cdot 2^{v_{j_2}} \cdot \log_e 2$
SRA , SRAI , (SRL , SRLI)	I	$v_i = \frac{v_{j_1}}{2^{v_{j_2}}}$	$\overrightarrow{\Delta v_i} = \overrightarrow{\Delta v_{j_1}} \cdot \frac{1}{2^{v_{j_2}}} - \overrightarrow{\Delta v_{j_2}} \cdot v_{j_1} \cdot \frac{1}{(2^{v_{j_2}})} \cdot \log_e 2$

# Speedup

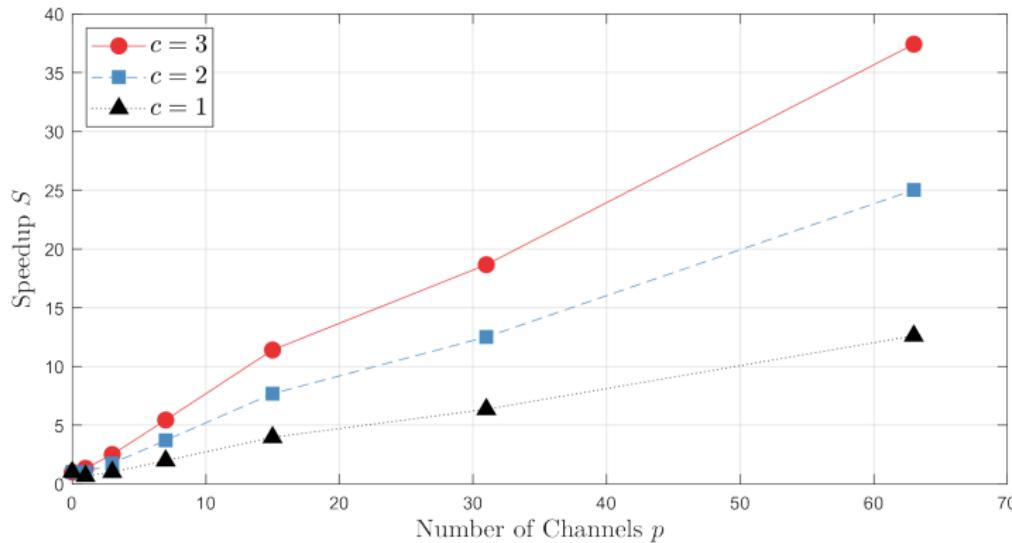


Figure: Exptected Speedup

# What works quite well

## Advantages

1. "Easy" to use, user doesn't have to think much about AD
2. AD adaptions don't alter the "normal" RISC-V instructions/ operations
3. Performance looks promising on paper
4. Scalable/ configurable design (FPGA)

# Next steps:

## Todo

1. Extend to T-RV32IMF (Mostly done)
2. Memory Hierarchy; so far BRAM instead of DRAM
  - Quantify effects on cache misses?
  - Applicability to larger workloads?
3. Ideally: Extend to T-RV32IMFV

# Limitations

## Todo

1. Clock frequency reduced
2. Processor size increased (FPGA usage)
3. Needs  $p$ -times more memory by design
4. Limited amount of arithmetic operations of the ISA
5. Useful larger applications?
6. Extensive testing needed
7. Usability of the prototype
8. Bit-Tricks

# Thanks!