

FPL 2024

Fantastic arithmetic beasts and where to find them

Florent de Dinechin

Bogdan Pasca

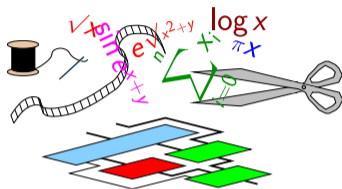


Where to find them? In your application...

Where to find them? In your application...

... but then you will need some help to bring them to light.

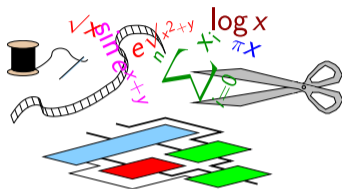
www.flopoco.org



Where to find them? In your application...

... but then you will need some help to bring them to light.

www.flopoco.org



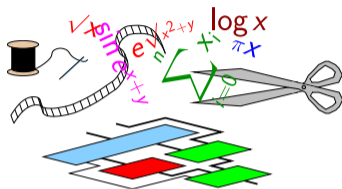
A generator of **application-specific** hardware arithmetic operators

- written in C++, outputting VHDL
- open and extensible

Where to find them? In your application...

... but then you will need some help to bring them to light.

www.flopoco.org



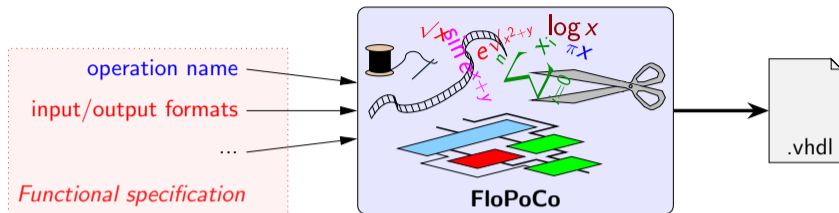
A generator of **application-specific** hardware arithmetic operators

- written in C++, outputting VHDL
- open and extensible

A philosophy of **computing just right**

- Interface: You ask for 17 bits, you get 17 *correct* bits.
- Inside: (try to) never compute bits that are not useful to the final result

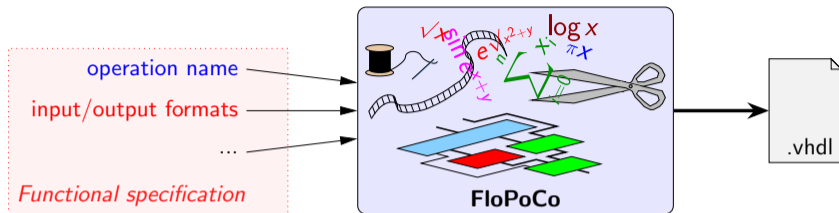
From the boring domestic animals ...



Everybody likes a single precision floating-point adder (here combinatorial)

```
./flopoco IEEEFPAdd wE=8 wF=23
```

From the boring domestic animals ...



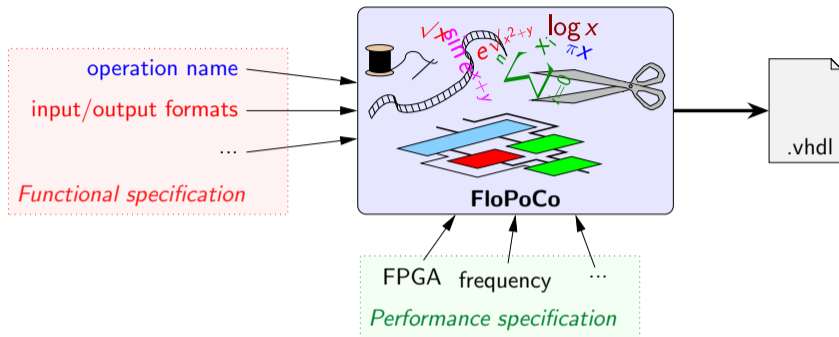
Everybody likes a single precision floating-point adder (here combinatorial)

```
./flopoco IEEEFPAdd wE=8 wF=23
```

... but it has many interesting little brothers...

```
./flopoco FPAdd wE=6 wF=31
```

From the boring domestic animals ...



Everybody likes a single precision floating-point adder (here combinatorial)

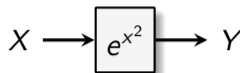
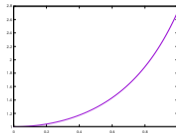
```
./flopoco IEEEFPAdd wE=8 wF=23
```

... but it has many interesting little brothers...

```
./flopoco FPAdd wE=6 wF=31 frequency=300 dualpath=true
```

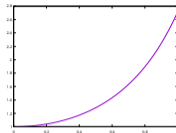

... to fantastic arithmetic beasts

Suppose you need to evaluate some function,
say $e^{(x^2)}$ on $[0, 1)$...



... to fantastic arithmetic beasts

Suppose you need to evaluate some function,
say $e^{(x^2)}$ on $[0, 1)$...
... with inputs and outputs on 24 bits.



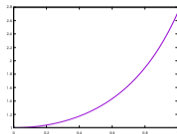
... to fantastic arithmetic beasts

Suppose you need to evaluate some function¹,
say $e^{(x^2)}$ on $[0, 1)$...

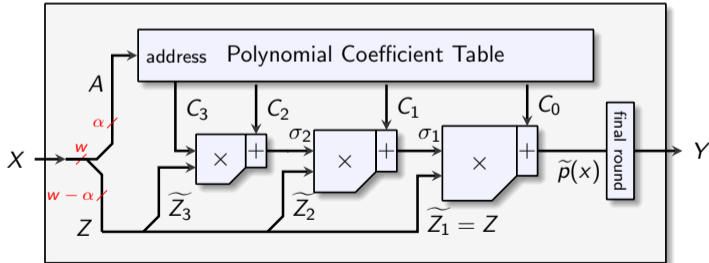
... with inputs and outputs on 24 bits.

There are several ways of doing this in FloPoCo.

Here is one of them.



```
./flopoco FixFunctionByPiecewisePoly f="exp(x*x)" lsbIn=-24 lsbOut=-24 d=3
```



¹ It works on the set of functions on which it works (TM)

Another one from one of yesterday's talk

Computing $X \bmod 3329$ for X a 24-bit integer:

```
./flopoco IntConstDiv wIn=24 d=3329 computeQuotient=false  
computeRemainder=true arch=3
```

(Not as good as yesterday's paper, though.)



Danila Gorodecky and Leonel Sousa

Scalable architecture of constant division on FPGA.

ARITH, 2023.

Florent is busy until retirement

We'll see more fantastic beast in this talk.

Not all: there is already an infinite number of them in FloPoCo...

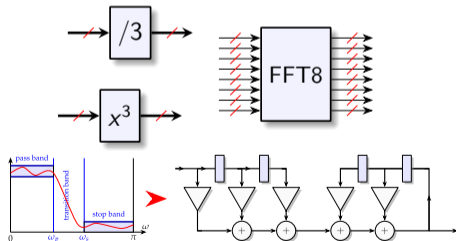
and a larger infinity still to be discovered.

Scope of FloPoCo

Hardware finite-precision implementations

of any computing kernel with a **clear mathematical definition**.

- We have only scratched the surface of function approximation
- We'll see many variants of classical operations
- Coarser kernels such as Fast Fourier Transforms
- From a frequency response to an IIR
- ...



Agenda

Careless PhD students and their pets gone wrong

Fantastic but not evil: circuits computing just right

Fantastic arithmetic beasts escaped to vendor tools

Bit heaps: the mutant biology of arithmetic beasts

Why fantastic arithmetic beasts didn't take over the world (and how to address it)

Backup slides

Careless PhD students and their pets gone wrong

Careless PhD students and their pets gone wrong

Fantastic but not evil: circuits computing just right

Fantastic arithmetic beasts escaped to vendor tools

Bit heaps: the mutant biology of arithmetic beasts

Why fantastic arithmetic beasts didn't take over the world (and how to address it)

Backup slides

All my life, I have been afflicted with very good students

Very good students tend to write kilolines of code...

All my life, I have been afflicted with very good students

Very good students tend to write kilolines of code...

Jérémie Detrey's PhD, 2004-2007

- FPLibrary: open-source VHDL for floating-point $+$, $-$, \times , $/$, $\sqrt{\quad}$

All my life, I have been afflicted with very good students

Very good students tend to write kilolines of code...

Jérémy Detrey's PhD, 2004-2007

- FPLibrary: open-source VHDL for floating-point $+$, $-$, \times , $/$, $\sqrt{\quad}$
- then parametric floating-point `sin`, `cos`, `exp`, `log`, ...

All my life, I have been afflicted with very good students

Very good students tend to write kilolines of code...

Jérémie Detrey's PhD, 2004-2007

- FPLibrary: open-source VHDL for floating-point $+$, $-$, \times , $/$, $\sqrt{\quad}$
- then parametric floating-point \sin , \cos , \exp , \log , ...
- using two novel generic techniques for hardware function approximation
 - multipartite tables
 - HOTBM (higher-order table method)

All my life, I have been afflicted with very good students

Very good students tend to write kilolines of code...

Jérémie Detrey's PhD, 2004-2007

- FPLibrary: open-source VHDL for floating-point $+$, $-$, \times , $/$, $\sqrt{\quad}$
- then parametric floating-point `sin`, `cos`, `exp`, `log`, ...
- using two novel generic techniques for hardware function approximation
 - multipartite tables
 - HOTBM (higher-order table method)
- then LNS (Logarithm Number System) operators for good measure

All my life, I have been afflicted with very good students

Very good students tend to write kilolines of code...

Jérémie Detrey's PhD, 2004-2007

- FPLibrary: open-source VHDL for floating-point $+$, $-$, \times , $/$, $\sqrt{\quad}$
- then parametric floating-point \sin , \cos , \exp , \log , ...
- using two novel generic techniques for hardware function approximation
 - multipartite tables
 - HOTBM (higher-order table method)
- then LNS (Logarithm Number System) operators for good measure

16 papers, thanks to a solid and well-tested agile development methodology

one paper, one random heap of quick-and-dirty code

All my life, I have been afflicted with very good students

Very good students tend to write kilolines of code...

Jérémie Detrey's PhD, 2004-2007

- FPLibrary: open-source VHDL for floating-point $+$, $-$, \times , $/$, $\sqrt{\quad}$
- then parametric floating-point \sin , \cos , \exp , \log , ...
- using two novel generic techniques for hardware function approximation
 - multipartite tables
 - HOTBM (higher-order table method)
- then LNS (Logarithm Number System) operators for good measure

16 papers, thanks to a solid and well-tested agile development methodology

one paper, one random heap of quick-and-dirty code

- All sorts of bits of Java/Python/C++ to generate some of the VHDL

All my life, I have been afflicted with very good students

Very good students tend to write kilolines of code...

Jérémie Detrey's PhD, 2004-2007

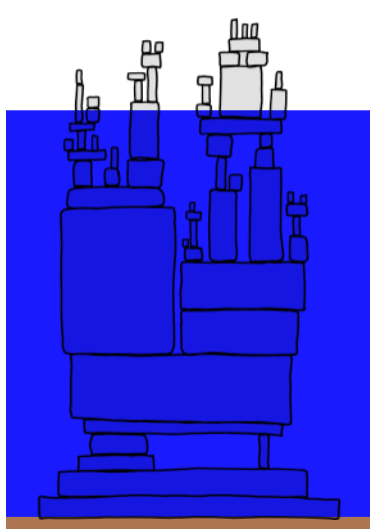
- FPLibrary: open-source VHDL for floating-point $+$, $-$, \times , $/$, $\sqrt{\quad}$
- then parametric floating-point \sin , \cos , \exp , \log , ...
- using two novel generic techniques for hardware function approximation
 - multipartite tables
 - HOTBM (higher-order table method)
- then LNS (Logarithm Number System) operators for good measure

16 papers, thanks to a solid and well-tested agile development methodology

one paper, one random heap of quick-and-dirty code

- All sorts of bits of Java/Python/C++ to generate some of the VHDL
- Design-space exploration scripts, test-benches, etc

Our scientific artifacts after Jérémie's PhD

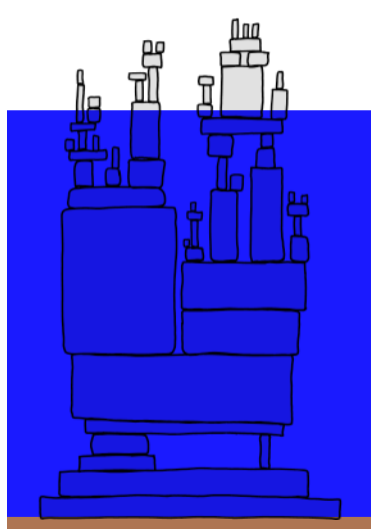


FPLibrary (VHDL available online)

stuff described in Jérémie's PhD

drawing from <https://xkcd.com/2347/>

Our scientific artifacts after Jérémie's PhD

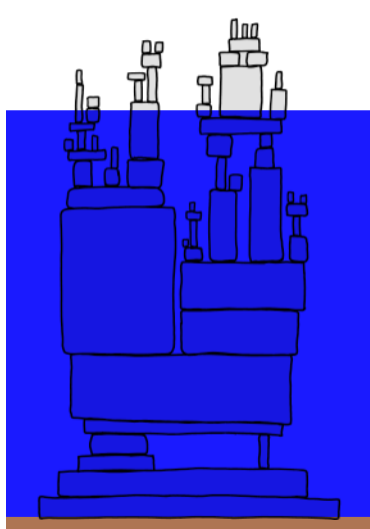


FPLibrary (VHDL available online)

stuff described in Jérémie's PhD
(in French)

drawing from <https://xkcd.com/2347/>

Our scientific artifacts after Jérémie's PhD



FPLibrary (VHDL available online)

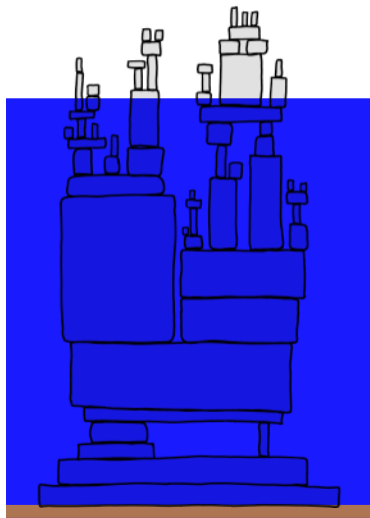
stuff described in Jérémie's PhD
(in French)

in other words:

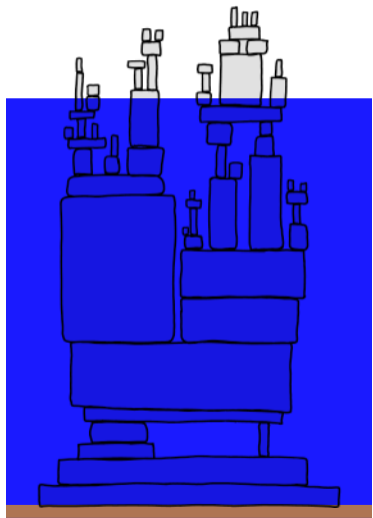
work doomed to oblivion when the student leaves
(after his PhD, Jérémie defected to finite-field arithmetic)

drawing from <https://xkcd.com/2347/>

Hence an Engineering Grand Plan



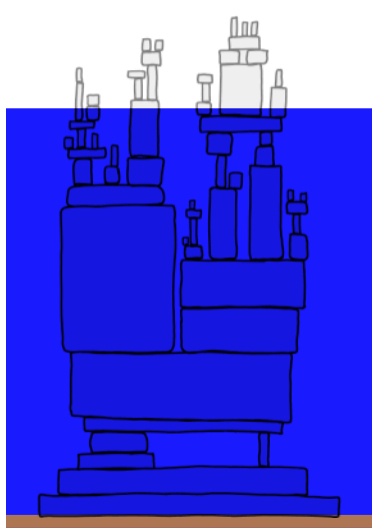
Hence an Engineering Grand Plan



Rewrite this from scratch,
and distribute it

and it shall be called FloPoCo:
Floating-**P**oint **C**ores (but not only)

Hence an Engineering Grand Plan

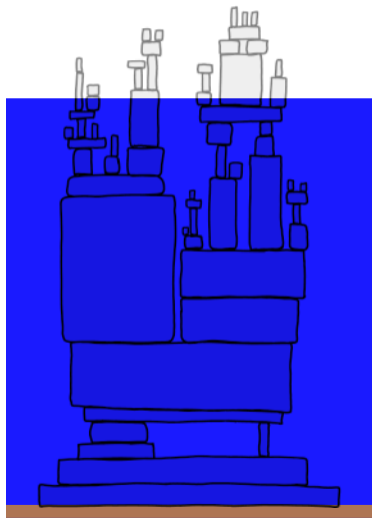


VHDL is generated,
no need to distribute it

Rewrite this from scratch,
and distribute it

and it shall be called FloPoCo:
Floating-**P**oint **C**ores (but not only)

Hence an Engineering Grand Plan



VHDL is generated,
no need to distribute it

Rewrite this from scratch,
and distribute it

and it shall be called FloPoCo:
Floating-**P**oint **C**ores (but not only)

OK, it doesn't really look like a winning move...
but wait a bit.
(and I need to hire a Really Good PhD student)

Historical excuses for all the bad technical choices

- It had to be **C++** because Jérémie had written HOTBM in C++
(and the thing is, at the time, I didn't understand half of it)

Historical excuses for all the bad technical choices

- It had to be **C++** because Jérémie had written HOTBM in C++
(and the thing is, at the time, I didn't understand half of it)
- **Generating VHDL** because FPLibrary was written in VHDL
(and to be frank, quite a lot of it was magical to me)

First version of FloPoCo was a superset of FPLibrary... by printing out FPLibrary code

Historical excuses for all the bad technical choices

- It had to be **C++** because Jérémie had written HOTBM in C++
(and the thing is, at the time, I didn't understand half of it)
- **Generating VHDL** because FPLibrary was written in VHDL
(and to be frank, quite a lot of it was magical to me)

First version of FloPoCo was a superset of FPLibrary... by printing out FPLibrary code

A ~~stupid primitive~~ modest approach to hardware generation, but immediate benefits

- better scaling, easier debugging than parametric VHDL
when you have many parameters
 - instead of a VHDL **generate if**, a C++ **if** \implies only the **true** branch in the VHDL
 - same for **generate for** loops
 - (compared to Jérémie's parametric recursive VHDL for tree-like structures)

Historical excuses for all the bad technical choices

- It had to be **C++** because Jérémie had written HOTBM in C++
(and the thing is, at the time, I didn't understand half of it)
- **Generating VHDL** because FPLibrary was written in VHDL
(and to be frank, quite a lot of it was magical to me)
First version of FloPoCo was a superset of FPLibrary... by printing out FPLibrary code

A ~~stupid primitive~~ modest approach to hardware generation, but immediate benefits

- better scaling, easier debugging than parametric VHDL
when you have many parameters
 - instead of a VHDL **generate if**, a C++ **if** \implies only the **true** branch in the VHDL
 - same for **generate for** loops
 - (compared to Jérémie's parametric recursive VHDL for tree-like structures)
- and very soon: **automatic pipelining** – because each submitted paper stated:
“the design will be pipelined in the final version”
and this is a perfect waste of good student's time

The engineering foundations to a Scientific Grand Plan

First written in this paper

When FPGAs are better at floating-point than microprocessors

First written in this paper

When FPGAs are better at floating-point than microprocessors

- When? As soon as the processor lacks hardware support:

Models	Instruction Distribution					
	Add	Mult.	Div.	Sqrt.	Exp.	Log
bjt	22	30	17	0	2	0
diode	7	5	4	0	1	2
hbt	112	57	51	0	23	18
jfet	13	31	2	0	2	0
mos1	24	36	7	1	0	0
vbic	36	43	18	1	10	4

SPICE Model-Evaluation,
cut from Kapre and DeHon (FPL 2009)

Dura Amdahl lex, sed lex.

- but also fused operations such as $\sqrt{x^2 + y^2}$, and more...
- In my humble opinion, this was a visionary paper: submitted to ISFPGA 2008

Lack of results, prove it

As we all know, a reviewer is always right.
Therefore, we stubbornly wrote

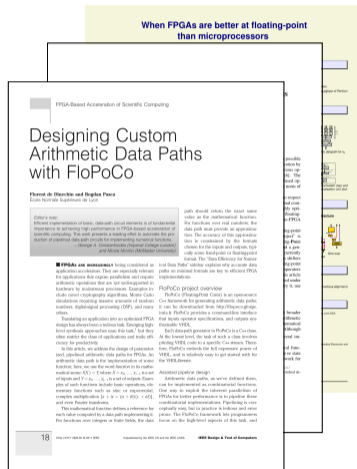
- An FPL 2009 paper:
Generating high-performance custom floating-point pipelines.



Lack of results, prove it

As we all know, a reviewer is always right.
Therefore, we stubbornly wrote

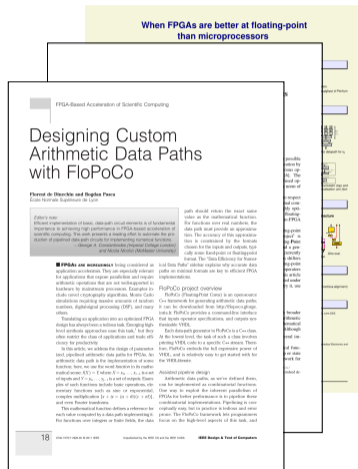
- An FPL 2009 paper:
Generating high-performance custom floating-point pipelines.
- Its journal version *Designing custom arithmetic data paths with FloPoCo.*
IEEE Design & Test of Computers, 2011.



Lack of results, prove it

As we all know, a reviewer is always right.
Therefore, we stubbornly wrote

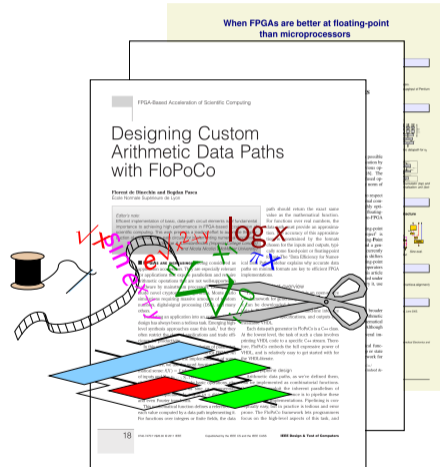
- An FPL 2009 paper:
Generating high-performance custom floating-point pipelines.
- Its journal version *Designing custom arithmetic data paths with FloPoCo.*
IEEE Design & Test of Computers, 2011.
- 400+ GScholar citations,
TCFPGA Hall of Fame this year



Lack of results, prove it

As we all know, a reviewer is always right.
Therefore, we stubbornly wrote

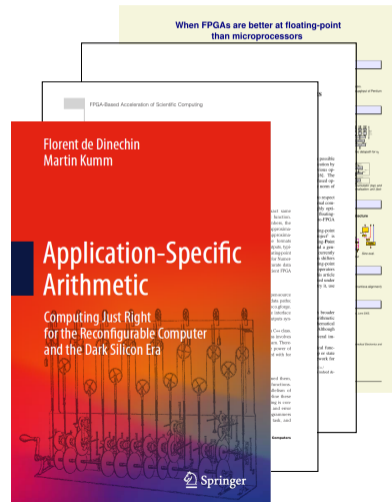
- An FPL 2009 paper:
Generating high-performance custom floating-point pipelines.
- Its journal version *Designing custom arithmetic data paths with FloPoCo.*
IEEE Design & Test of Computers, 2011.
 - 400+ GScholar citations,
TCFPGA Hall of Fame this year
 - How many people actually read it?
It is the “how to cite” paper for FloPoCo



Lack of results, prove it

As we all know, a reviewer is always right.
Therefore, we stubbornly wrote

- An FPL 2009 paper:
Generating high-performance custom floating-point pipelines.
- Its journal version *Designing custom arithmetic data paths with FloPoCo.*
IEEE Design & Test of Computers, 2011.
 - 400+ GScholar citations,
TCFPGA Hall of Fame this year
 - How many people actually read it?
It is the “how to cite” paper for FloPoCo
- And finally, an **800-page** book:
Application-Specific Arithmetic. Springer, 2024.



Refinement of the Grand Plan

If the final title of your PhD is the same as it was when you started,
your research is probably boring.

When FPGAs are better at floating-point than microprocessors

Refinement of the Grand Plan

If the final title of your PhD is the same as it was when you started,
your research is probably boring.

When FPGAs are better at floating-point than microprocessors

Not your neighbour's FPU

Refinement of the Grand Plan

If the final title of your PhD is the same as it was when you started,
your research is probably boring.

When FPGAs are better at floating-point than microprocessors

Not your neighbour's FPU

FPGA-specific arithmetic

(floating-point, but not only)

Refinement of the Grand Plan

If the final title of your PhD is the same as it was when you started,
your research is probably boring.

When FPGAs are better at floating-point than microprocessors

Not your neighbour's FPU

FPGA-specific arithmetic *(floating-point, but not only)*

All the operators you will never see in a processor (and how to build them)

Refinement of the Grand Plan

If the final title of your PhD is the same as it was when you started,
your research is probably boring.

When FPGAs are better at floating-point than microprocessors

Not your neighbour's FPU

FPGA-specific arithmetic *(floating-point, but not only)*

All the operators you will never see in a processor (and how to build them)

Save routing! Save power! Don't move around useless bits!

Refinement of the Grand Plan

If the final title of your PhD is the same as it was when you started,
your research is probably boring.

When FPGAs are better at floating-point than microprocessors

Not your neighbour's FPU

FPGA-specific arithmetic *(floating-point, but not only)*

All the operators you will never see in a processor (and how to build them)

Save routing! Save power! Don't move around useless bits!

Application-specific arithmetic *(FPGA, but not only)*

Refinement of the Grand Plan

If the final title of your PhD is the same as it was when you started

When FPGAs are better at floating-point than microprocessors

Not your neighbour's FPU

FPGA-specific arithmetic

All the operators you will never see in a processor (and

Save routing! Save power! Don't move around useless

Application-specific arithmetic

Circuits computing just right

Florent de Dinechin
Martin Kumm

Application-Specific Arithmetic

Computing Just Right
for the Reconfigurable Computer
and the Dark Silicon Era



Message to my younger self

If you believe in an idea, stick to it, whatever the reviews say.

(bad reviews just mean your good idea was badly explained...)

Fantastic but not evil: circuits computing just right

Careless PhD students and their pets gone wrong

Fantastic but not evil: circuits computing just right

Fantastic arithmetic beasts escaped to vendor tools

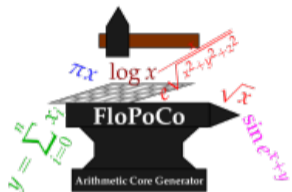
Bit heaps: the mutant biology of arithmetic beasts

Why fantastic arithmetic beasts didn't take over the world (and how to address it)

Backup slides

Computing just right?

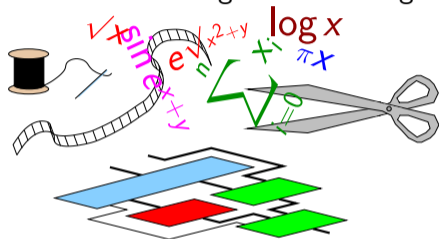
Bogdan said “great project, we need a logo” and designed this:



It was initially OK, but... soon we were using a *delicate chisel* more than a **hammer**.

Computing just right?

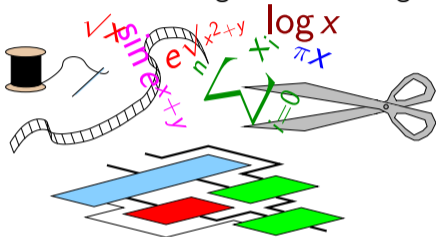
So as soon as Bogdan left I designed this:



(the proper term is probably *allogory*)

Computing just right?

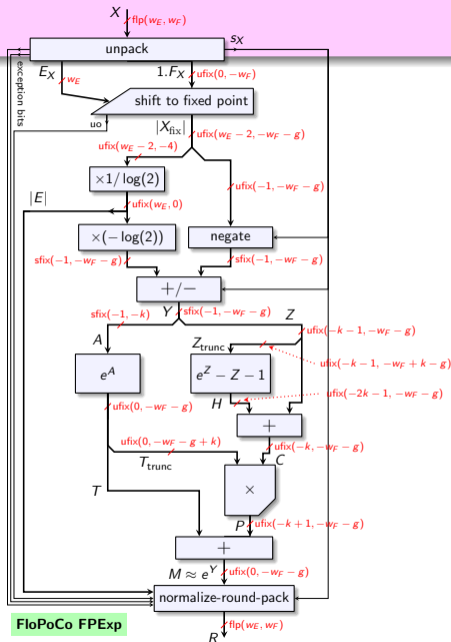
So as soon as Bogdan left I designed this:



(the proper term is probably *allogory*)

This is the kind of thing FloPoCo does →
It is a **floating-point exponential** operator
where each wire, each component is
tailored to its context with love and care.

(not a very good logo either)



Save resources! Save power! Don't move useless bits around!

In software, as soon as your result is correct, it is probably wasteful

Gustafson: Does *Angry birds* really need single precision (8 decimal digits of accuracy) considering that the trajectory was input using your fat fingers?

Save resources! Save power! Don't move useless bits around!

In software, as soon as your result is correct, it is probably wasteful

Gustafson: Does *Angry birds* really need single precision (8 decimal digits of accuracy) considering that the trajectory was input using your fat fingers?

Plain common sense

- If the lower bits carry useless noise, you don't want to compute them...
- ... and you want even less to store them, transmit them, compute on them.

In FPGAs, we have this freedom.

Save resources! Save power! Don't move useless bits around!

In software, as soon as your result is correct, it is probably wasteful

Gustafson: Does *Angry birds* really need single precision (8 decimal digits of accuracy) considering that the trajectory was input using your fat fingers?

Plain common sense

- If the lower bits carry useless noise, you don't want to compute them...
- ... and you want even less to store them, transmit them, compute on them.

In FPGAs, we have this freedom.

With great freedom come great opportunities

- In a circuit, we may choose, for each variable, how many bits are computed/stored/transmitted! → **the opportunities**
- Overwhelming freedom! Help! → **the challenges**

Opportunity #1: Computing Just Right makes interfaces simpler

Output format (number of output bits) specifies operator accuracy

The output format defines a quantum of precision (u or ulp for “unit in the last place”)



Opportunity #1: Computing Just Right makes interfaces simpler

Output format (number of output bits) specifies operator accuracy

The output format defines a quantum of precision (u or ulp for “unit in the last place”)

- No need to compute more accurately than u : we couldn't output it
- No sense in computing less accurately than u : we don't want to output garbage bits



Opportunity #1: Computing Just Right makes interfaces simpler

Output format (number of output bits) specifies operator accuracy

The output format defines a quantum of precision (u or ulp for “unit in the last place”)

- No need to compute more accurately than u : we couldn't output it
- No sense in computing less accurately than u : we don't want to output garbage bits



Inspired by IEEE-754

- define **quantization**(x) for $x \in \mathbb{R}$.
- define **operator**(X) = **quantization**(**operation**(X))

Example: round to the nearest.

Opportunity #1: Computing Just Right makes interfaces simpler

Output format (number of output bits) specifies operator accuracy

The output format defines a quantum of precision (u or ulp for “unit in the last place”)

- No need to compute more accurately than u : we couldn't output it
- No sense in computing less accurately than u : we don't want to output garbage bits



Inspired by IEEE-754

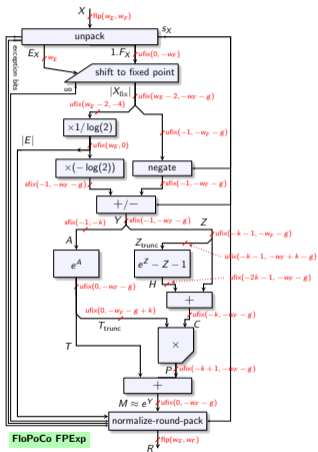
- define **quantization**(x) for $x \in \mathbb{R}$.
- define **operator**(X) = **quantization**(**operation**(X))

Example: round to the nearest.

If you add one bit to the output, you divide u by 2, hence double the accuracy.

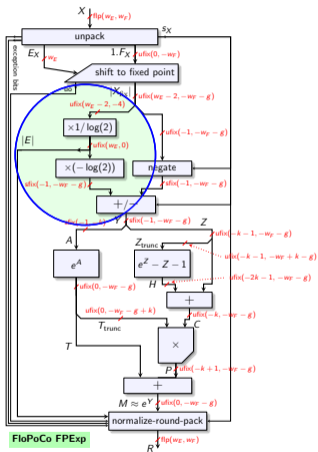
Computing just right \iff over-parameterization

Example:



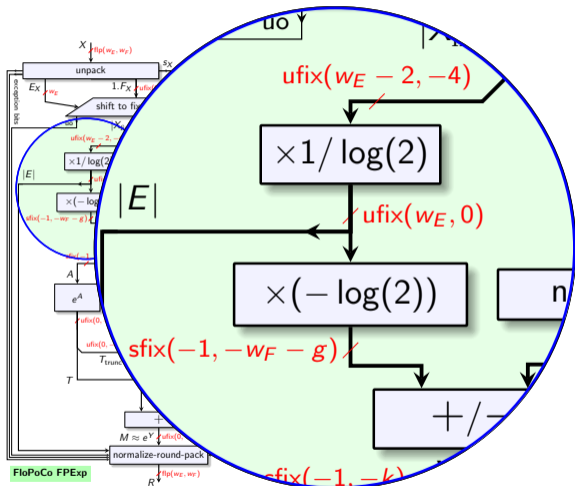
Computing just right \iff over-parameterization

Example:

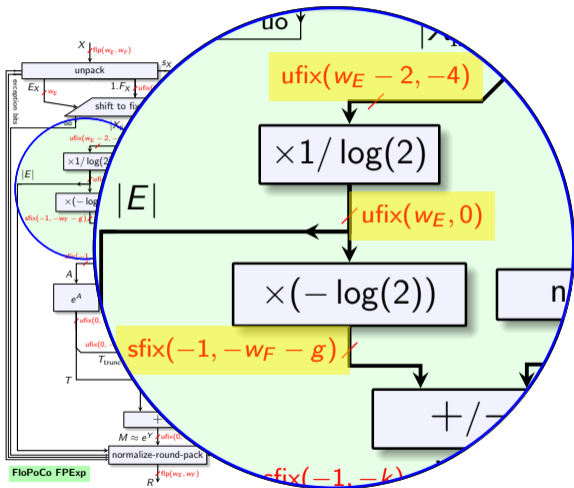


Computing just right \iff over-parameterization

Example:



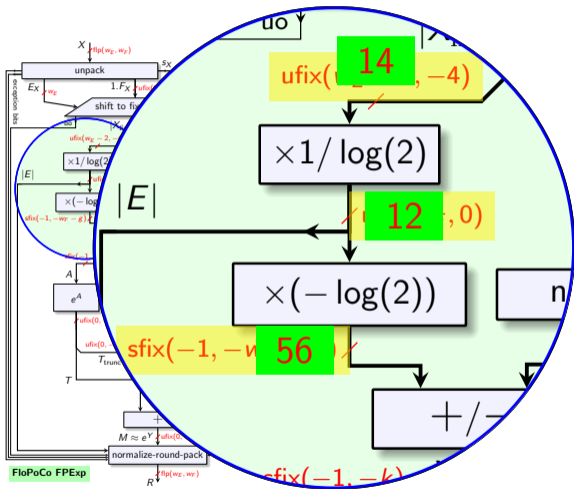
Computing just right \iff over-parameterization



Example:

Multipliers of all shapes and sizes

Computing just right \iff over-parameterization



Example:

Multipliers of all shapes and sizes

In a double-precision exponential,

- $w_E = 11$, $w_F = 52$,
- first multiplier 14-bits in, 12 bits out
- second multiplier 12-bits in, 56 bits out
... and truncated left and right

Not your neighbour's multiplier.
Very strange arithmetic beasts.

Over-parameterization sounds a bit like over-engineering, doesn't it?

- ⊖ OK, there is a bit more work involved in designing a parametric operator
 - To start with, it must be a hardware-generating program:
 - There is an infinite number of multipliers-by-a-constant.
- You cannot chain them all in a library.

Over-parameterization sounds a bit like over-engineering, doesn't it?

- ⊖ OK, there is a bit more work involved in designing a parametric operator
 - To start with, it must be a hardware-generating program:
 - There is an infinite number of multipliers-by-a-constant.
 - You cannot chain them all in a library.

- ⊕ Direct benefit to end-users: freedom of choice, application-specific, etc.

Over-parameterization sounds a bit like over-engineering, doesn't it?

- ⊖ OK, there is a bit more work involved in designing a parametric operator
 - To start with, it must be a hardware-generating program:
 - There is an infinite number of multipliers-by-a-constant.
 - You cannot chain them all in a library.

- ⊕ Direct benefit to end-users: freedom of choice, application-specific, etc.

- ⊕ More future-proof when the target hardware changes

Over-parameterization sounds a bit like over-engineering, doesn't it?

- ⊖ OK, there is a bit more work involved in designing a parametric operator
 - To start with, it must be a hardware-generating program:
 - There is an infinite number of multipliers-by-a-constant.
 - You cannot chain them all in a library.
- ⊕ Direct benefit to end-users: freedom of choice, application-specific, etc.
- ⊕ More future-proof when the target hardware changes
- ⊕ **It actually simplifies the design of composite operators** (e.g. the exponential)!
 - You don't know how many bits on this wire make sense? Keep it open as a parameter.
 - Then experiment: estimate cost and accuracy as a function of the parameters
 - Then program the choice of the best parameter values,
 - e.g. using ILP or common sense (whichever gives the best results)

Opportunity #2: Operator specialization

Not really fantastic arithmetic beasts... just the usual ones with special disabilities.

- Multiplication by a constant

- multiplication by integers: $17X = (X \ll 4) + X$;
 $8721X = ((17X) \ll 9) + 17X$
- but also by reals such as $\log(2)$ or $\sin(42\pi/256)$
- Two main techniques, tens of papers
- *An FFT mostly consists of constant multiplications*

Opportunity #2: Operator specialization

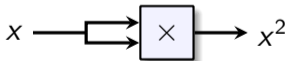
Not really fantastic arithmetic beasts... just the usual ones with special disabilities.

- Multiplication by a constant
 - multiplication by integers: $17X = (X \ll 4) + X$;
 $8721X = ((17X) \ll 9) + 17X$
 - but also by reals such as $\log(2)$ or $\sin(42\pi/256)$
 - Two main techniques, tens of papers
 - *An FFT mostly consists of constant multiplications*
- Division by 3 (for various values of 3)
 - in floating point for Jacobi and other stencils
 - integer (quotient and remainder) for addressing in 3 memory banks

Opportunity #2: Operator specialization

Not really fantastic arithmetic beasts... just the usual ones with special disabilities.

- Multiplication by a constant
 - multiplication by integers: $17X = (X \ll 4) + X$;
 $8721X = ((17X) \ll 9) + 17X$
 - but also by reals such as $\log(2)$ or $\sin(42\pi/256)$
 - Two main techniques, tens of papers
 - *An FFT mostly consists of constant multiplications*
- Division by 3 (for various values of 3)
 - in floating point for Jacobi and other stencils
 - integer (quotient and remainder) for addressing in 3 memory banks
- A squarer is a multiplier specialization



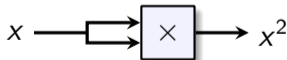
$$\begin{array}{r} 321 \\ \times 321 \\ \hline 321 \\ 642 \\ 963 \\ \hline 103041 \end{array}$$

Opportunity #2: Operator specialization

Not really fantastic arithmetic beasts... just the usual ones with special disabilities.

- Multiplication by a constant
 - multiplication by integers: $17X = (X \ll 4) + X$;
 $8721X = ((17X) \ll 9) + 17X$
 - but also by reals such as $\log(2)$ or $\sin(42\pi/256)$
 - Two main techniques, tens of papers
 - *An FFT mostly consists of constant multiplications*
- Division by 3 (for various values of 3)
 - in floating point for Jacobi and other stencils
 - integer (quotient and remainder) for addressing in 3 memory banks

- A squarer is a multiplier specialization



- Specialization of elementary functions to specific domains
- ...

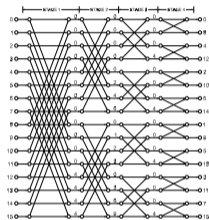
$$\begin{array}{r} 321 \\ \times 321 \\ \hline 321 \\ 642 \\ 963 \\ \hline 103041 \end{array}$$

A FloPoCo-enabled success story

 Mario Garrido, Konrad Möller, and Martin Kumm.

World's Fastest FFT Architectures: Breaking the Barrier of 100 GS/s.
IEEE Transactions on Circuits and Systems I, 66(4):1507–1516, 2019.

- Fully unrolled FFT (up to 256 points)
 - i.e. inputting 256 complex values per cycle, at 500 MHz
 - well above 10 TOp/s if you count all additions and multiplications
- 16-bit in/out, wider datapath inside
- Look, Ma: no multiplier !
 - each multiplier expanded as an adder graph (and optimally so)
- about 1/5th of LUT + registers of the target device (Virtex UltraScale 190)
 - ... leaving the 1800 DSP blocks free for more interesting things.



A good start, in FPGA design, is not to imitate the processor solution.

Division by 3 is simpler than exponential

TL;DR: multiplying X by 3 is computing $2X + X$;

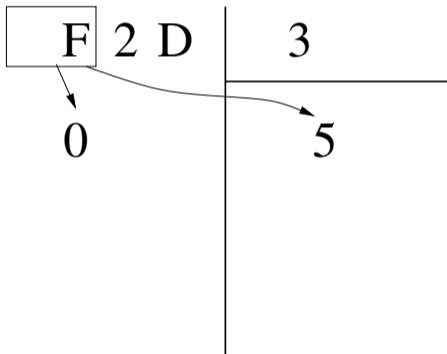
Dividing by 3 should not be much more complex.

Division by 3 is simpler than exponential

TL;DR: multiplying X by 3 is computing $2X + X$;

Dividing by 3 should not be much more complex.

Dividing an **hexadecimal** number by 3

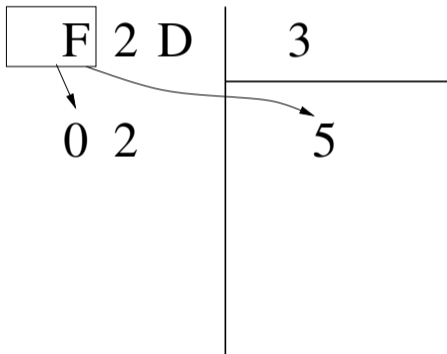


Division by 3 is simpler than exponential

TL;DR: multiplying X by 3 is computing $2X + X$;

Dividing by 3 should not be much more complex.

Dividing an **hexadecimal** number by 3

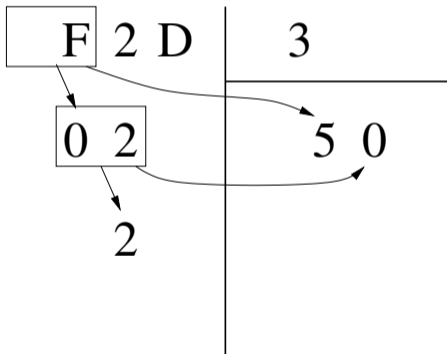


Division by 3 is simpler than exponential

TL;DR: multiplying X by 3 is computing $2X + X$;

Dividing by 3 should not be much more complex.

Dividing an **hexadecimal** number by 3

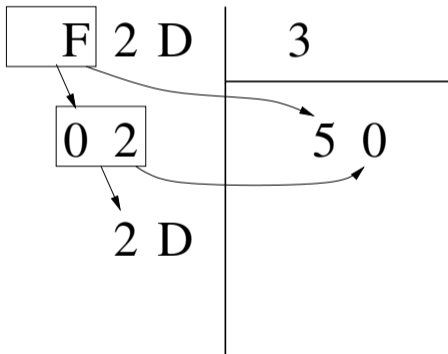


Division by 3 is simpler than exponential

TL;DR: multiplying X by 3 is computing $2X + X$;

Dividing by 3 should not be much more complex.

Dividing an **hexadecimal** number by 3

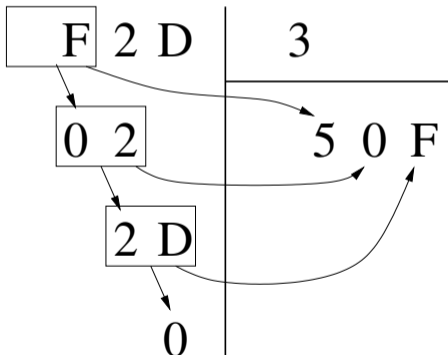


Division by 3 is simpler than exponential

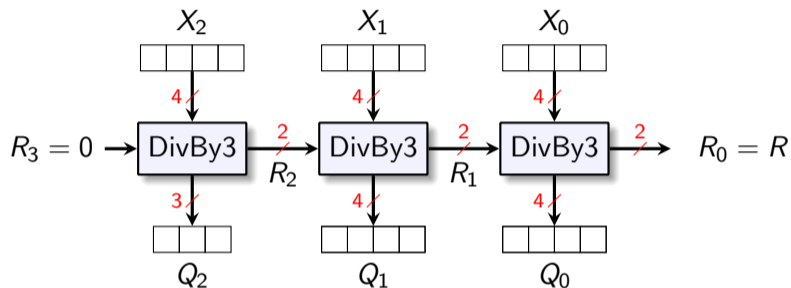
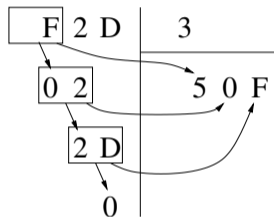
TL;DR: multiplying X by 3 is computing $2X + X$;

Dividing by 3 should not be much more complex.

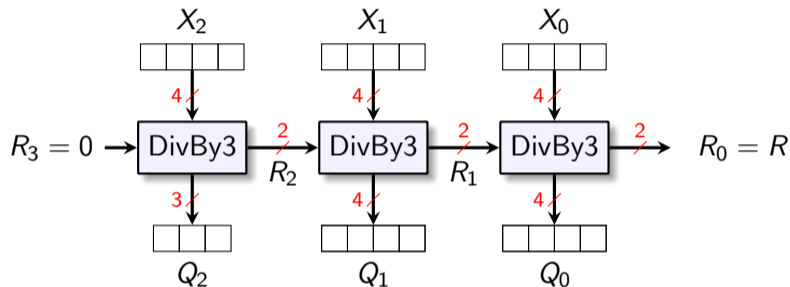
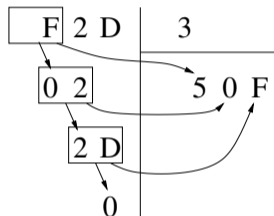
Dividing an **hexadecimal** number by 3



Getting inspiration from the vexations of childhood

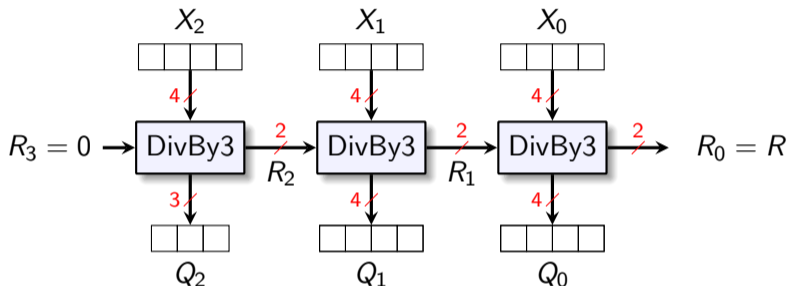
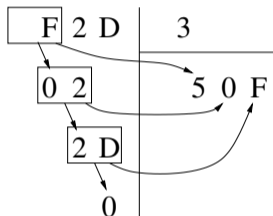


Getting inspiration from the vexations of childhood



OK, this looks like an architecture, but we still need to build this (smaller) DivBy3 box.

Getting inspiration from the vexations of childhood

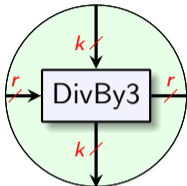


OK, this looks like an architecture, but we still need to build this (smaller) DivBy3 box.

If you're too lazy to compute, then tabulate

... here a table of 2^6 entries of 6 bits each.

Opportunity #3: target-specific optimizations

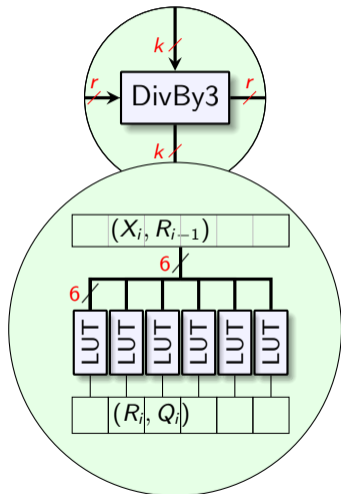


Generalizing hexadecimal to **radix 2^k**

... or, how **over-parameterization** allows for adaptation

- to various values of 3, like $D = 5$, or 7, or 9

Opportunity #3: target-specific optimizations



Generalizing hexadecimal to **radix 2^k**

... or, how **over-parameterization** allows for adaptation

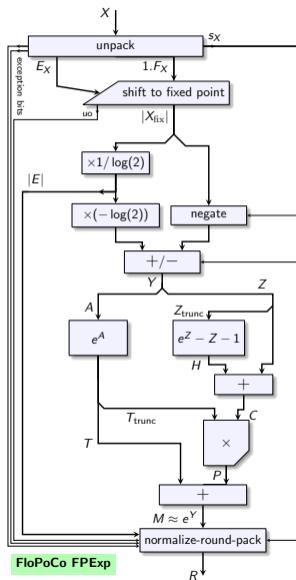
- to various values of 3, like $D = 5$, or 7, or 9
- to a given FPGA

Perfect match to modern FPGAs

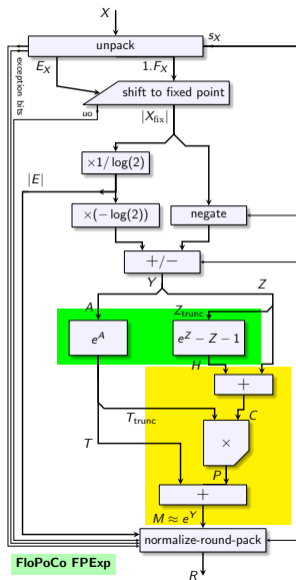
Unit of area: the LUT, with α input bits (here $\alpha = 6$)

Opportunity #3: target-specific optimizations

Modern FPGAs also have



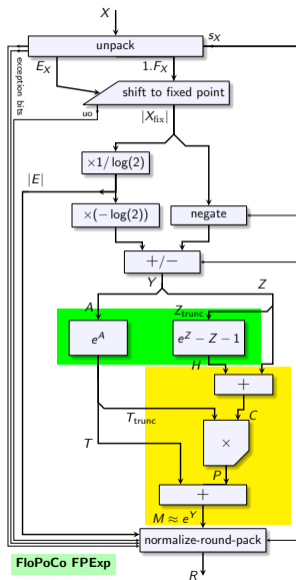
Opportunity #3: target-specific optimizations



Modern FPGAs also have

- small multipliers with pre-adders and post-adders
- ... and dual-ported small memories

Opportunity #3: target-specific optimizations



Modern FPGAs also have

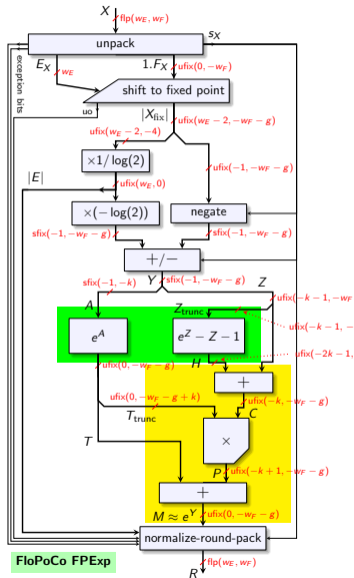
- small multipliers with pre-adders and post-adders
- ... and dual-ported small memories

Single-precision accurate exponential on Xilinx

- one block RAM (0.1% of the chip)
- one DSP block (0.1%)
- < 400 LUTs (0.1%, \approx one FP adder)

to compute one exponential per cycle at 500MHz
(\sim one AVX512 core trashing on its 16 FP32 lanes)

Opportunity #3: target-specific optimizations



Modern FPGAs also have

- small multipliers with pre-adders and post-adders
- ... and dual-ported small memories

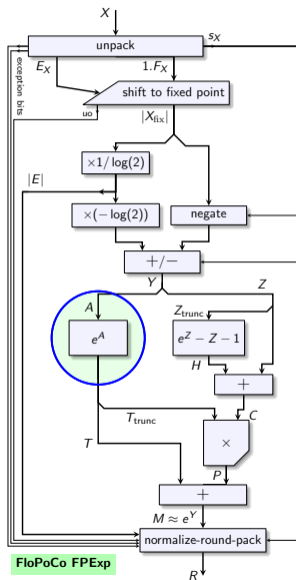
Single-precision accurate exponential on Xilinx

- one block RAM (0.1% of the chip)
- one DSP block (0.1%)
- < 400 LUTs (0.1%, \approx one FP adder)

to compute one exponential per cycle at 500MHz
(\sim one AVX512 core trashing on its 16 FP32 lanes)

*For one specific value only of the architectural parameter $k!$
(over-parameterization is cool)*

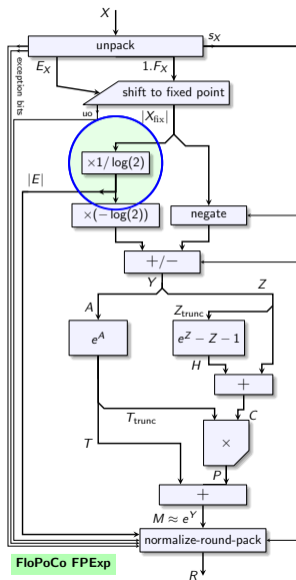
Opportunity #4: Tabulation



*Being unable to trust my reasoning, I learnt by heart
the results of all the possible multiplications
(E. Ionesco)*

- ... and all the possible exponentials

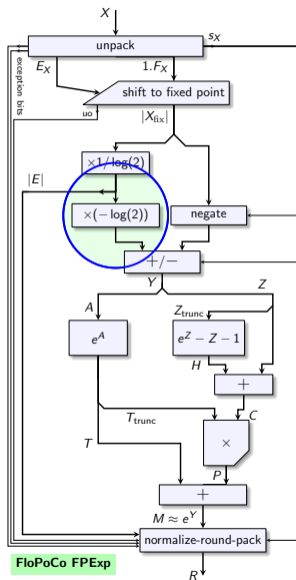
Opportunity #4: Tabulation



*Being unable to trust my reasoning, I learnt by heart
the results of all the possible multiplications
(E. Ionesco)*

- ... and all the possible exponentials
- ... and all the possible values of $e^Z - Z - 1$
- ... and indeed, all the possible multiplications

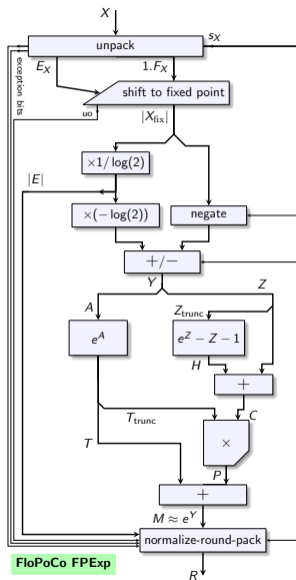
Opportunity #4: Tabulation



*Being unable to trust my reasoning, I learnt by heart
the results of all the possible multiplications
(E. Ionesco)*

- ... and all the possible exponentials
- ... and all the possible values of $e^Z - Z - 1$
- ... and indeed, all the possible multiplications

Opportunity #4: Tabulation

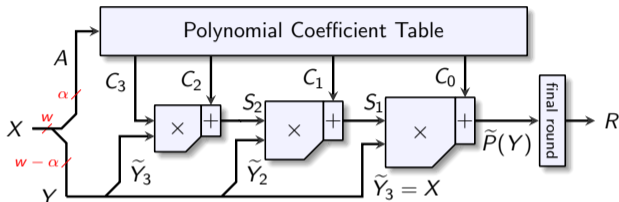
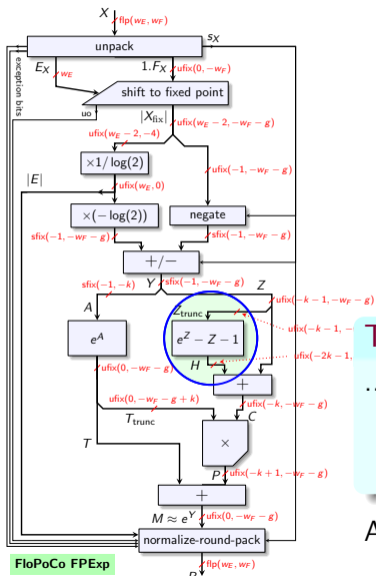


*Being unable to trust my reasoning, I learnt by heart
the results of all the possible multiplications*
(E. Ionesco)

- ... and all the possible exponentials
- ... and all the possible values of $e^Z - Z - 1$
- ... and indeed, all the possible multiplications

Reading a tabulated value is very efficient
when the table is close to the consumer.

Opportunity #5: Generic approximators (when tabulation won't scale)



The FloPoCo FixFunctionByPiecewisePoly operator

... has taken us so much time it is well worth a full part.

- state-of-the-art polynomial approximation
- each multiplier tailored with love and care

Also multipartite tables, filter approximators, and more to come.

Opportunity #6: merged arithmetic in bit heaps

... has taken us so much time it is well worth a full part.

Fantastic arithmetic beasts escaped to vendor tools

Careless PhD students and their pets gone wrong

Fantastic but not evil: circuits computing just right

Fantastic arithmetic beasts escaped to vendor tools

Bit heaps: the mutant biology of arithmetic beasts

Why fantastic arithmetic beasts didn't take over the world (and how to address it)

Backup slides



FloPoCo, PhD and Altera

- 2011 joined Altera European Technology Center
- brought the FloPoCo spirit along
- grafted into the DSP Builder team



FloPoCo, PhD and Altera

- 2011 joined Altera European Technology Center
- brought the FloPoCo spirit along
- grafted into the DSP Builder team
 - model-based design (Matlab Simulink frontend)
 - powerful mapping backend (using WYSIWYG)
 - floating-point support in its infancy

<https://agwaycapecod.com/the-art-of-grafting-plants/>

PhD life v.s. Industry

- PhD: highly efficient exponential implementation
- few months: approach analysis, design, implementation, test

- PhD: highly efficient exponential implementation
- few months: approach analysis, design, implementation, test

Floating-point exponential functions for DSP-enabled FPGAs

Florent de Dinechin, Bogdan Pasca

LIP (ENSIL-CMIS-Aristo-CICRA), Ecole Normale Supérieure de Lyon
46 allée d'Italie, F-69622 Lyon, France
{florent.de_dinechin, bogdan.pasca}@enscm.fr

Abstract—This article presents a generator of floating-point exponential operators targeting recent FPGAs with embedded memories and DSP blocks. A single-precision operator consumes just one DSP block, 18Kbits of dual-port memory, and 90 slices on Virtex-4. For larger precisions, a generic approach based on polynomial approximation is used and proves more resource-efficient than the iterates. For instance a double-precision operator consumes 4 BlockRAM and 12 DSP48 blocks on Virtex-4, or 10 36Kb and 22 36Kb multipliers on Stratix-III. This approach is flexible and is demonstrated to scale up to quadruple-precision, while enabling frequencies close to the FPGA's maximal frequency. All the proposed architectures are fast and accurate for all the floating-point range. They are suitable in the open-source FPGaCts framework.

1. INTRODUCTION

The exponential function is, after the basic arithmetic operations, one of the most most useful building block for floating-point applications. On FPGAs, it has been used for scientific or financial Monte-Carlo simulations [1], for SPICE simulation [2], in phylogenetic tree reconstruction, in quantum chemistry simulations, and in the implementation of the power function [3] among others.

A. Previous works

Several publications have described exponential implementations. We list them here, and will discuss more details the choices they made and their performance impact in Section IV. Earlier works targeted single precision, first by adapting to FPGAs a software algorithm based on floating-point operations [4], then by using a more efficient fixed-point architecture [5]. This architecture was later improved [1], however the table-based method used there doesn't scale up to double-precision, as the size of the tables grows exponentially with the mantissa size.

As FPGAs are increasingly being used for double-precision floating-point, iterative architectures that scale better [6], [7], [8] were adapted for FPGAs [9]. The architecture in [9] was designed with 5-input LUTs in mind, but is poorly suited to DSP-enabled FPGAs, as IV-8 will show. It was parametrized in precision, but to our knowledge was never pipelined. Another pipelined, but double-precision only implementation was proposed in [10], [11].

In [12], a COREIC-based approach using several parallel COREIC cores was proposed. It has a complex control including input and output FIFOs. Being racks-2 COREIC,

it computes one digit per iteration and thus has a very long latency. Moreover, it is based on a floating-point adder, whereas COREIC is inherently a fixed-point computation, so there is probably room for improvement there.

From a user point of view, the current state of the art is probably the floating-point exponential function `EXP7_232` provided with Altera MegaWizard since 2008 [13]. This implementation exploits the DSP blocks, is parametrized in exponent and mantissa size, and is fully pipelined. Being included in the standard Quartus releases, it is widely available, although only for Altera targets.

Many other publications have addressed the computation of exponential function in ASIC, e.g. [6], [7], [14], [8]. However, it is difficult to evaluate the relevance of such works on FPGAs.

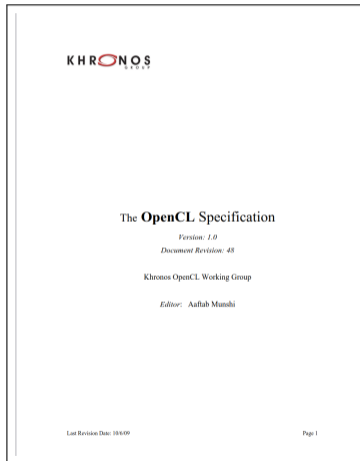
B. Contributions

In the present article, we propose yet another architecture for the floating-point evaluation of the exponential function, and its implementation in the open-source FPGaCts project¹. Its main specifications are the following:

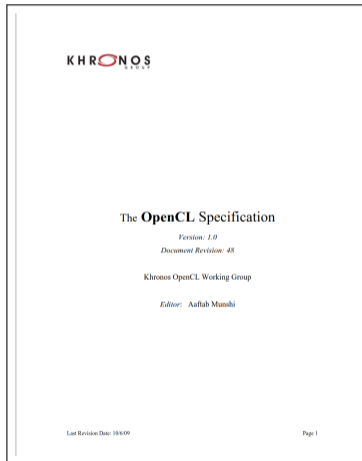
- The algorithm, based on the usual multiplicative stage selection followed by a polynomial approximation, was chosen with DSP blocks and embedded memories in mind, so it makes efficient use of these resources. For instance, the single-precision version now involves just one 17x17-bit multiplier and 18Kbits of dual-port memory, and runs at 375MHz on a Virtex-4, which is a large improvement in all respects over the state of the art [1].
- As we believe that floating point on FPGAs should exploit the flexibility of the target and therefore not be limited to IEEE single and double precision, the algorithm and implementation proposed here are fully parametrized in exponent and mantissa size. They scale to double-precision and beyond.
- The implementation is pipelined to a user-specified frequency. It is fast but accurate for all supported mantissa sizes.
- The architectures are generated as synthesizable VHDL portable to any FPGA target. In addition, many target-specific optimizations are performed by the FPGaCts framework [15], [16], [17].

¹<http://www.enscm.fr/~lip/Genetic/Work/FPGaCts/>

- OpenCL - first real driver for math.h coverage



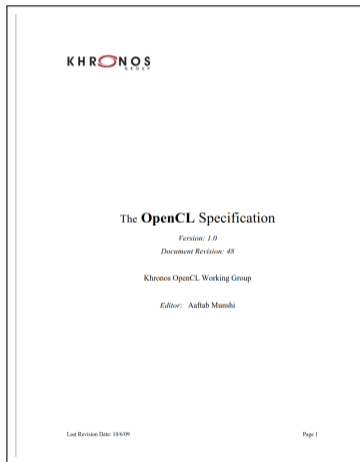
- OpenCL - first real driver for math.h coverage



The image shows the table of contents for the OpenCL Specification document. It lists various sections and their corresponding page numbers. The sections include:

- 6.5 Address Space Qualifiers.....150
 - 6.5.1 __global (or global).....150
 - 6.5.2 __local (or local).....151
 - 6.5.3 __constant (or constant).....151
 - 6.5.4 __private (or private).....152
- 6.6 Image Access Qualifiers.....153
- 6.7 Function Qualifiers.....154
 - 6.7.1 __kernel (or kernel).....154
 - 6.7.2 Optional Attribute Qualifiers.....154
- 6.8 Restrictions.....157
- 6.9 Preprocessor Directives and Macros.....160
- 6.10 Attribute Qualifiers.....162
 - 6.10.1 Specifying Attributes of Types.....163
 - 6.10.2 Specifying Attributes of Functions.....165
 - 6.10.3 Specifying Attributes of Variables.....165
 - 6.10.4 Specifying Attributes of Blocks and Control-Flow-Statements.....167
 - 6.10.5 Extending Attribute Qualifiers.....167
- 6.11 Built-in Functions.....168
 - 6.11.1 Work-Item Functions.....168
 - 6.11.2 Math Functions.....170
 - 6.11.2.1 Floating-point macros and programs for math.....172
 - 6.11.3 Integer Functions.....177
 - 6.11.4 Common Functions.....180
 - 6.11.5 Geometric Functions.....182
 - 6.11.6 Relational Functions.....184
 - 6.11.7 Vector Data Load and Store Functions.....186
 - 6.11.8 Image Read and Write Functions.....190
 - 6.11.8.1 Samples.....192
 - 6.11.8.2 Built-in Image Functions.....192
 - 6.11.9 Synchronization Functions.....199
 - 6.11.10 Explicit Memory Fence Functions.....200
 - 6.11.11 Access Copies from Global to Local Memory, Local to Global Memory, and Prefetch.....201
- 7. OPENCL NUMERICAL COMPLIANCE.....203
 - 7.1 Rounding Modes.....203
 - 7.2 INF, NaN and Denormalized Numbers.....203
 - 7.3 Floating-Point Exceptions.....204
 - 7.4 Relative Error as ULPs.....204
 - 7.5 Edge Case Behavior.....207
 - 7.5.1 Additional Requirements Beyond C99 TC2.....207
 - 7.5.2 Changes to C99 TC2 Behavior.....210
 - 7.5.3 Edge Case Behavior in Flush-To-Zero Mode.....211

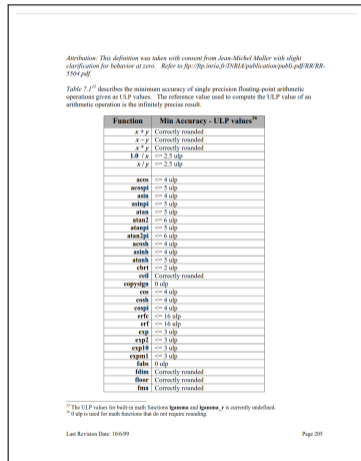
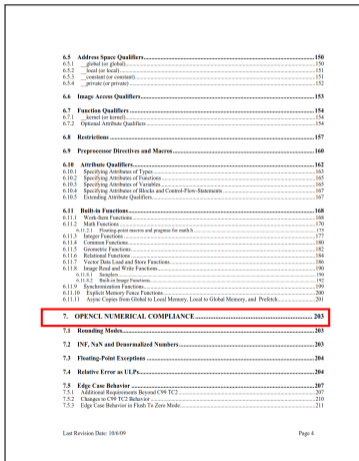
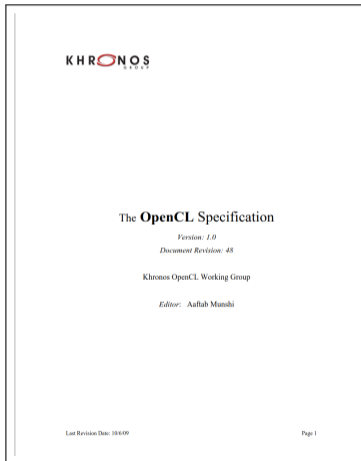
- OpenCL - first real driver for math.h coverage



The image shows the Table of Contents for the OpenCL Specification. The entries are as follows:

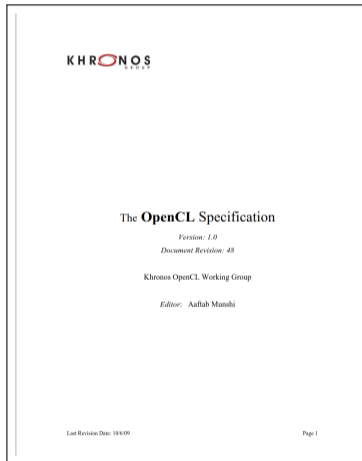
6.5 Address Space Qualifiers.....	150
6.5.1 _global (or global).....	150
6.5.2 _local (or local).....	151
6.5.3 _constant (or constant).....	151
6.5.4 _private (or private).....	152
6.6 Image Access Qualifiers.....	153
6.7 Function Qualifiers.....	154
6.7.1 _kernel (or kernel).....	154
6.7.2 Optional Attribute Qualifiers.....	154
6.8 Restrictions.....	157
6.9 Preprocessor Directives and Macros.....	160
6.10 Attribute Qualifiers.....	162
6.10.1 Specifying Attributes of Types.....	163
6.10.2 Specifying Attributes of Functions.....	165
6.10.3 Specifying Attributes of Variables.....	165
6.10.4 Specifying Attributes of Blocks and Control-Flow-Statements.....	167
6.10.5 Extending Attribute Qualifiers.....	167
6.11 Built-in Functions.....	168
6.11.1 Work-Item Functions.....	168
6.11.2 Math Functions.....	170
6.11.2.1 Floating-point macros and programs for math.....	172
6.11.3 Integer Functions.....	177
6.11.4 Common Functions.....	180
6.11.5 Geometric Functions.....	182
6.11.6 Relational Functions.....	184
6.11.7 Vector Data Load and Store Functions.....	186
6.11.8 Image Read and Write Functions.....	190
6.11.8.1 Samples.....	192
6.11.8.2 Built-in Image Functions.....	192
6.11.9 Synchronization Functions.....	199
6.11.10 Explicit Memory Fence Functions.....	200
6.11.11 Access Copies from Global to Local Memory, Local to Global Memory, and Prefetch.....	201
7. OPENCIL NUMERICAL COMPLIANCE.....	203
7.1 Rounding Modes.....	203
7.2 INF, NaN and Denormalized Numbers.....	203
7.3 Floating-Point Exceptions.....	204
7.4 Relative Error as ULPs.....	204
7.5 Edge Case Behavior.....	207
7.5.1 Additional Requirements Beyond C99 TC2.....	207
7.5.2 Changes to C99 TC2 Behavior.....	210
7.5.3 Edge Case Behavior in Flush-To-Zero Mode.....	211

- OpenCL - first real driver for math.h coverage



PhD life v.s. Industry

- OpenCL - first real driver for math.h coverage



The image shows the table of contents for the OpenCL Specification. The entries are listed on the left with their corresponding page numbers on the right. The entry '7. OPENC L NUMERICAL COMPLIANCE' is highlighted with a red box. At the bottom right, it says 'Page 4'.

6.5 Address Space Qualifiers.....	150
6.5.1 _global (or global)	150
6.5.2 _local (or local)	151
6.5.3 _constant (or constant)	151
6.5.4 _private (or private)	152
6.6 Image Access Qualifiers.....	153
6.7 Function Qualifiers.....	154
6.7.1 _kernel (or kernel)	154
6.7.2 Optional Attribute Qualifiers.....	154
6.8 Restrictions.....	157
6.9 Preprocessor Directives and Macros.....	160
6.10 Attribute Qualifiers.....	162
6.10.1 Specifying Attributes of Types.....	163
6.10.2 Specifying Attributes of Functions.....	165
6.10.3 Specifying Attributes of Variables.....	165
6.10.4 Specifying Attributes of Blocks and Control-Flow-Statements.....	167
6.10.5 Extending Attribute Qualifiers.....	167
6.11 Built-in Functions.....	168
6.11.1 Work-Item Functions.....	168
6.11.2 Math Functions.....	170
6.11.2.1 Floating-point macros and programs for math.....	172
6.11.3 Integer Functions.....	177
6.11.4 Common Functions.....	180
6.11.5 Geometric Functions.....	182
6.11.6 Relational Functions.....	184
6.11.7 Vector Data Load and Store Functions.....	186
6.11.8 Image Read and Write Functions.....	190
6.11.8.1 Samples.....	192
6.11.8.2 Built-in Image Functions.....	192
6.11.9 Synchronization Functions.....	199
6.11.10 Explicit Memory Fence Functions.....	200
6.11.11 Access Capabilities from Global to Local Memory, Local to Global Memory, and Prefetch.....	201
7. OPENC L NUMERICAL COMPLIANCE.....	203
7.1 Rounding Modes.....	203
7.2 INF, NaN and Denormalized Numbers.....	203
7.3 Floating-Point Exceptions.....	204
7.4 Relative Error as ULPs.....	204
7.5 Edge Case Behavior.....	207
7.5.1 Additional Requirements Beyond C99 TC2.....	207
7.5.2 Changes to C99 TC2 Behavior.....	210
7.5.3 Edge Case Behavior in Flush-To-Zero Mode.....	211



Hardware Testing
Exhaustive where possible

Pivotal tool for arithmetic function design

Plenty of previous works using hardware polynomial approximation

- most of it hand-tuned for a given function (not generic)
- not accessible (papers, not code)
- heuristics used do not scale to precisions larger than 32 bits

Pivotal tool for arithmetic function design

Plenty of previous works using hardware polynomial approximation

- most of it hand-tuned for a given function (not generic)
- not accessible (papers, not code)
- heuristics used do not scale to precisions larger than 32 bits

Highlights

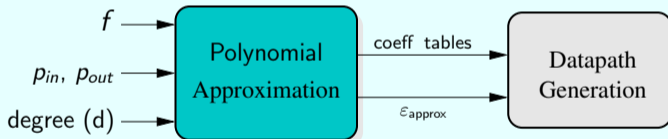
- scales up to **64bits and more**
- state-of-the art polynomial approximations thanks to Sollya
- finer **datapath optimization**
- **pipelined** to a user-specified frequency
- **fully automated** and integrated in open-source **FloPoCo**

How do we do it?

Consider the function $f(x)$, with $x \in [0, 1)$ and $IM(f(x)) \in [0, 1)$
Approximate it with the polynomial p of degree d (**given**) such that:

$$\epsilon_{\text{total}} = \max|f - p| \leq 2^{-p_{\text{out}}} = 1\text{ulp}$$

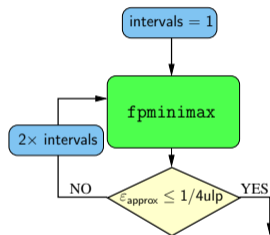
Two fold process



$$\epsilon_{\text{total}} = \epsilon_{\text{approx}} + \epsilon_{\text{eval}} + 1/2\text{ulp} \leq 1\text{ulp}$$

Polynomial approximation

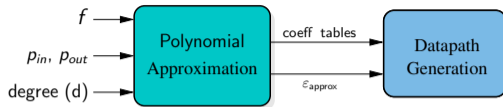
- use modified Remez algorithm from Sollya (fpminimax)
N. Brisebarre and S. Chevillard, *Efficient polynomial L^∞ - approximations*
- *input:*
 - function f , degree d
 - interval I
 - list of coefficient size constraints
- *output:*
 - polynomial with **precision-constrained** coefficients (no need to round them)
 - $\max(|f_i(y) - p_i(y)|) \leq \varepsilon_{\text{approx}} \forall i \in \{0..2^k - 1\}$



Advantages

- usually best polynomials given the input specifications
- might reduce by 1 polynomial degree for some intervals

Polynomial Evaluation



Use Horner (trade latency for size)

$$p(y) = a_0 + y \times (a_1 + y \times (a_2 + y \times (a_3 + y \times \underbrace{a_4}_{\sigma_0})))$$

$\underbrace{\hspace{10em}}_{\pi_1}$

$\underbrace{\hspace{10em}}_{\sigma_1}$

$\underbrace{\hspace{10em}}_{\pi_2}$

$\underbrace{\hspace{10em}}_{\sigma_2}$

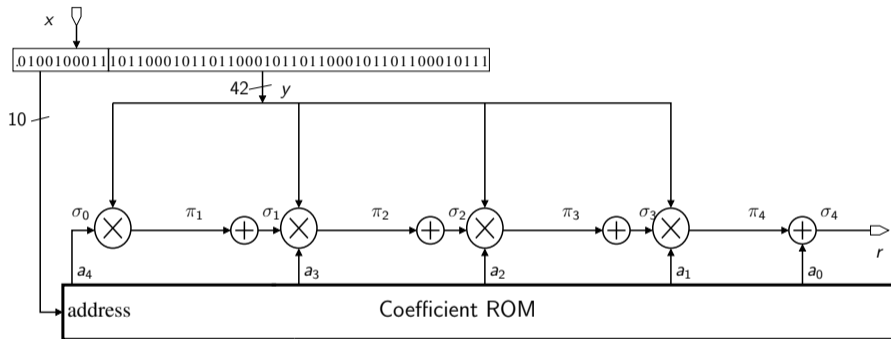
$\underbrace{\hspace{10em}}_{\pi_3}$

$\underbrace{\hspace{10em}}_{\sigma_3}$

$\underbrace{\hspace{10em}}_{\pi_4}$

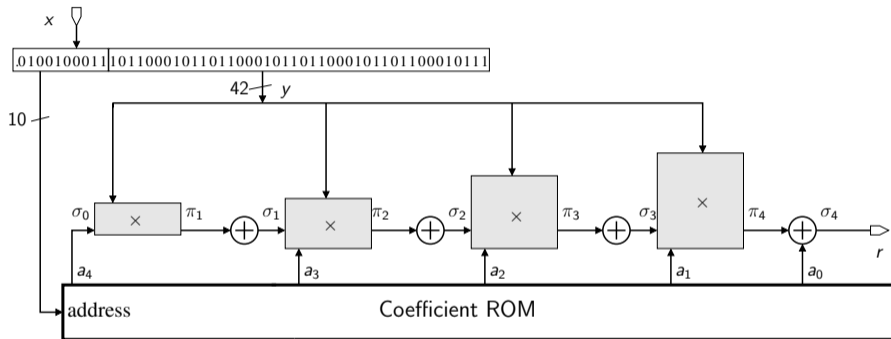
$\underbrace{\hspace{10em}}_{\sigma_4}$

The architecture for $\log_2(1+x)$, DP, $d=4$



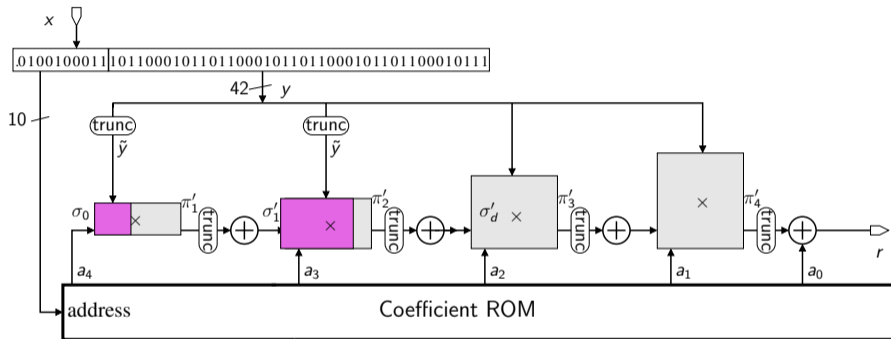
truncation on y and on π_j

The architecture for $\log_2(1+x)$, $DP, d=4$



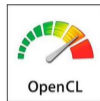
truncation on y and on π_j

The architecture for $\log_2(1+x)$, DP, $d=4$



truncation on y and on π_j

Where to find them? In Megawizard, DSP Builder, oneAPI, OpenCL



Same function, different VHDL:

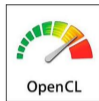
- pipelined to 300 MHz

```
./flopoco frequency=300 FPAdd wE=6 wF=31
```

- A larger but shorter-latency architectural variant:

```
./flopoco FPAdd wE=8 wF=23 dualpath=true
```

Where to find them? In Megawizard, DSP Builder, oneAPI, OpenCL



Same function, different VHDL:

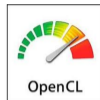
- pipelined to 300 MHz

```
./flopoco frequency=300 FPAdd wE=6 wF=31
```

- A larger but shorter-latency architectural variant:

```
./flopoco FPAdd wE=8 wF=23 dualpath=true
```

Where to find them? In Megawizard, DSP Builder, oneAPI, OpenCL



Same function, different VHDL:

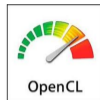
- pipelined to 300 MHz

```
./cmdPolyEval -frequency 300 FPAdd 6 31
```

- A larger but shorter-latency architectural variant:

```
./cmdPolyEval FPAddExpert 8 23 1 1 0
```


Where to find them? In Megawizard, DSP Builder, oneAPI, OpenCL



Same function, different VHDL:

- pipelined to 300 MHz

```
./cmdPolyEval -frequency 300 FPAdd 6 31
```

- A larger but shorter-latency architectural variant:

```
./cmdPolyEval FPAddExpert 8 23 1 1 0
```

- different function, another language (Verilog):

```
./cmdPolyEval -lang VERILOG FPArctan2 10 44
```

Bit heaps: the mutant biology of arithmetic beasts

Careless PhD students and their pets gone wrong

Fantastic but not evil: circuits computing just right

Fantastic arithmetic beasts escaped to vendor tools

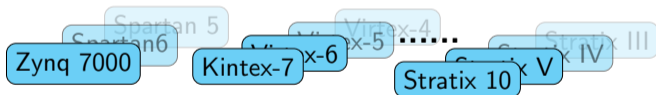
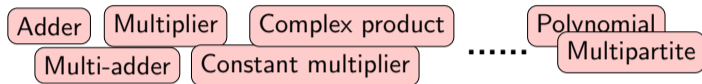
Bit heaps: the mutant biology of arithmetic beasts

Why fantastic arithmetic beasts didn't take over the world (and how to address it)

Backup slides

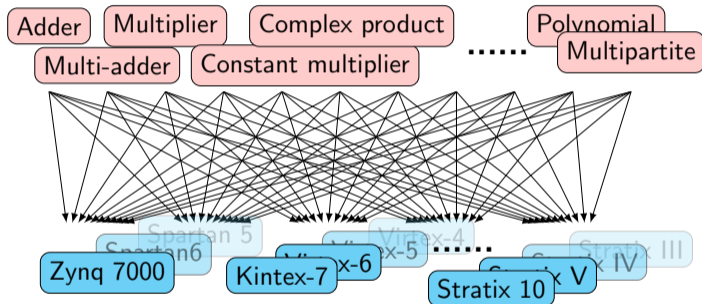
So much VHDL to write, so few ~~slaves~~ students to write it !

In theory, I know how to optimize by hand each operator for each target...



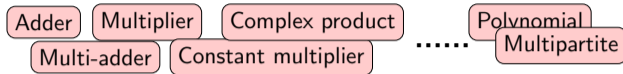
So much VHDL to write, so few ~~slaves~~ students to write it !

In theory, I know how to optimize by hand each operator for each target...

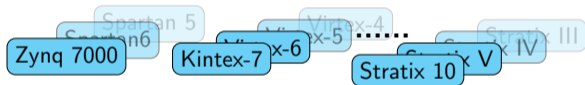


But I don't have the resources.

One data-structure to rule them all...

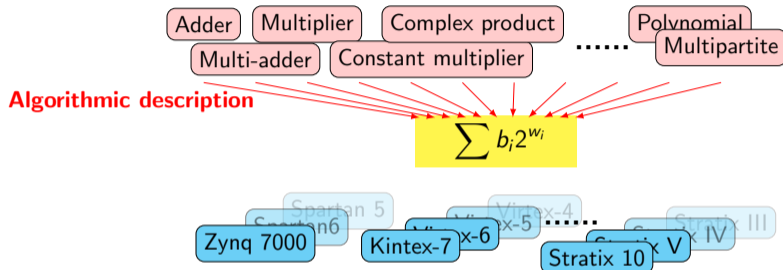


$$\sum b_i 2^{w_i}$$



The **sum of weighted bits** as a first-class arithmetic object

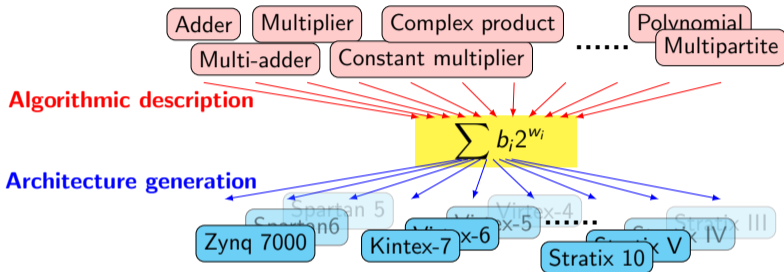
One data-structure to rule them all...



The **sum of weighted bits** as a first-class arithmetic object

- Captures the true binary math of an operation
(with all sorts of **bit-level optimization** opportunities)

One data-structure to rule them all... and in the hardware to bind them



The **sum of weighted bits** as a first-class arithmetic object

- Captures the true binary math of an operation
(with all sorts of **bit-level optimization** opportunities)
- The corresponding **compressor trees** can be optimized for each target
... and optimally so for practical sizes, thanks to M. Kumm

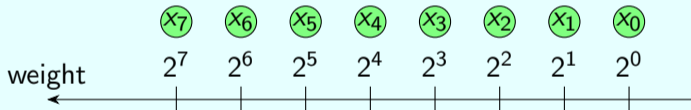
Sums of weighted bits?

- Integers or real numbers represented in **binary fixed-point**

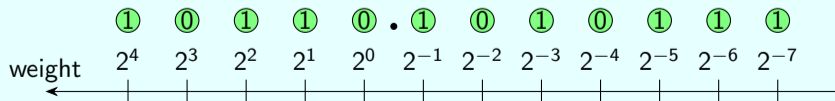
$$X = \sum_{i=i_{\min}}^{i_{\max}} 2^i x_i$$

- 2^i : “weight” \implies “ X is a sum of weighted bits”

Representation as a **dot diagrams**



Example: 17.42 written in binary



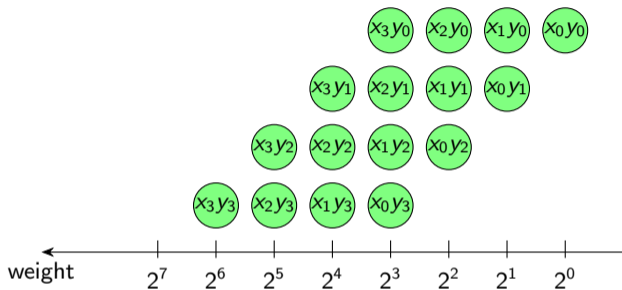
The historical bit heap was invented for building multipliers

$$\begin{aligned}XY &= \left(\sum_{i=0}^3 2^i x_i\right) \times \left(\sum_{j=0}^3 2^j y_j\right) \\ &= \sum_{i,j} 2^{i+j} x_i y_j\end{aligned}$$

A multiplier is an architecture that computes this sum.

The historical bit heap was invented for building multipliers

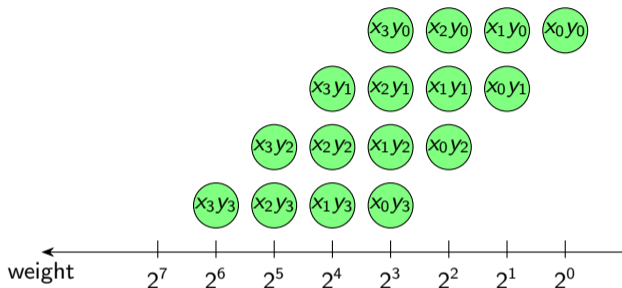
$$\begin{aligned} XY &= \left(\sum_{i=0}^3 2^i x_i \right) \times \left(\sum_{j=0}^3 2^j y_j \right) \\ &= \sum_{i,j} 2^{i+j} x_i y_j \end{aligned}$$



A multiplier is an architecture that computes this sum.

The historical bit heap was invented for building multipliers

$$XY = \left(\sum_{i=0}^3 2^i x_i \right) \times \left(\sum_{j=0}^3 2^j y_j \right)$$
$$= \sum_{i,j} 2^{i+j} x_i y_j$$



A multiplier is an architecture that computes this sum.

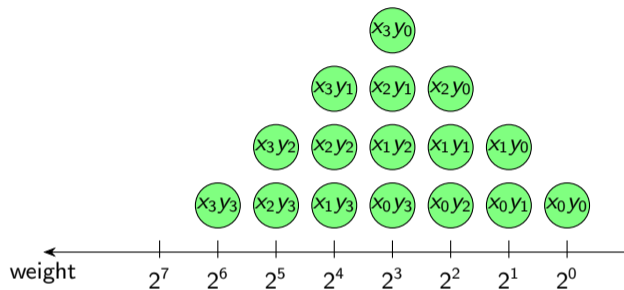
Historical motivation for bit heaps:

$\sum_{i,j} 2^{i+j} x_i y_j$ expresses the bit-level parallelism of the problem

... exposing design freedom thanks to **associativity** and **commutativity** of the \sum
(and a few other boolean tricks)

The historical bit heap was invented for building multipliers

$$\begin{aligned}XY &= \left(\sum_{i=0}^3 2^i x_i\right) \times \left(\sum_{j=0}^3 2^j y_j\right) \\ &= \sum_{i,j} 2^{i+j} x_i y_j\end{aligned}$$



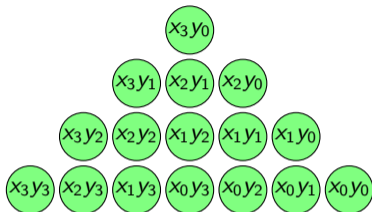
A multiplier is an architecture that computes this sum.

Historical motivation for bit heaps:

$\sum_{i,j} 2^{i+j} x_i y_j$ expresses the bit-level parallelism of the problem

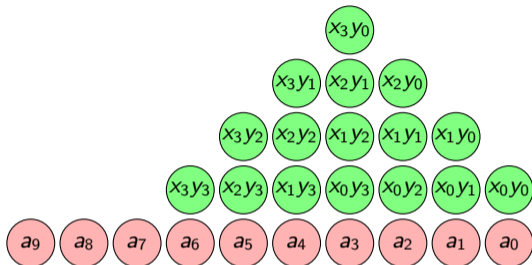
... exposing design freedom thanks to **associativity** and **commutativity** of the \sum
(and a few other boolean tricks)

$$XY = \sum_{i,j} 2^{i+j} x_i y_j$$



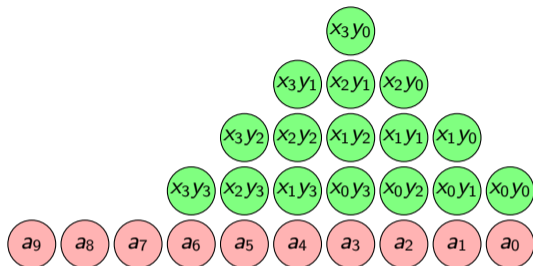
Beyond product

$$A + XY = \sum_i 2^i a_i + \sum_{i,j} 2^{i+j} x_i y_j$$

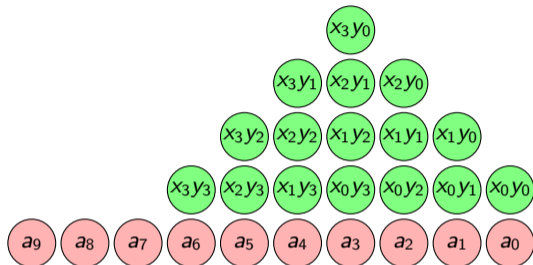


Beyond product

$$A + XY = \sum_{w,h} 2^w b_{w,h}$$



$$A + XY = \sum_{w,h} 2^w b_{w,h}$$



When generating an architecture

consider **only one big sum of weighted bits**

- get rid of artificial sequentiality (inside operators, and between operators)
- focus on true timing information (e.g. critical path delay of each weighted bit)
- A global optimization instead of several local ones (and solved by ILP)

Well beyond product

A bit heap is anything that can be developed as $\sum_{w,h} 2^w b_{w,h}$

- the sum of two bit heaps is obviously a bit heap
- the product of two bit heaps is also a bit heap

Well beyond product

A bit heap is anything that can be developed as $\sum_{w,h} 2^w b_{w,h}$

- the sum of two bit heaps is obviously a bit heap
- the product of two bit heaps is also a bit heap

Any polynomial of multiple variables is a bit heap

... where each $b_{w,h}$ is the AND of a few input bits.

This includes sums of squares, FIR filters, etc

Well beyond product

A bit heap is anything that can be developed as $\sum_{w,h} 2^w b_{w,h}$

- the sum of two bit heaps is obviously a bit heap
- the product of two bit heaps is also a bit heap

Any polynomial of multiple variables is a bit heap

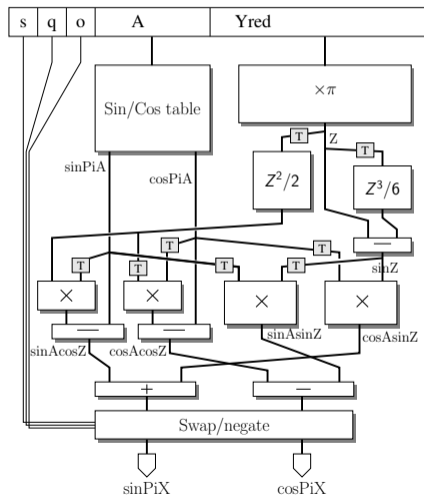
... where each $b_{w,h}$ is the AND of a few input bits.
This includes sums of squares, FIR filters, etc

And then more

- A huge class of function may be *approximated* by polynomials
- The $b_{w,h}$ may be read from arbitrary look-up tables
- An operator may include several bit heaps

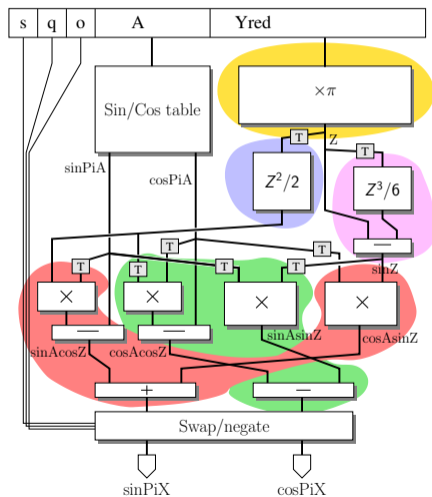
A good hammer transforms every problem into a nail

A sine/cosine architecture (HEART 2013)



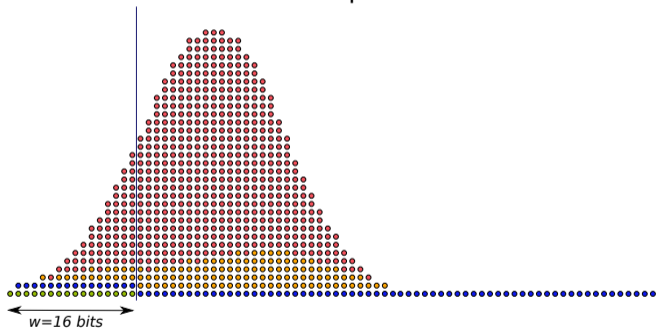
A good hammer transforms every problem into a nail

A sine/cosine architecture (HEART 2013) mostly consists of **5 bit heaps**

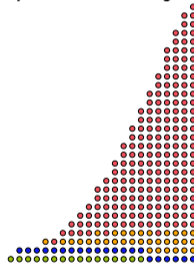


A bit heap for $Z - Z^3/6$ in the previous architecture

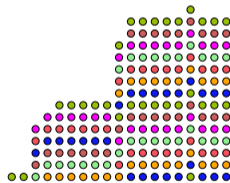
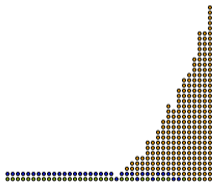
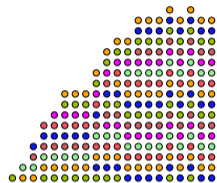
Full bit heap



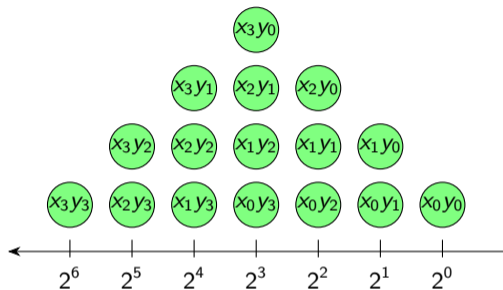
Bit heap truncated just right



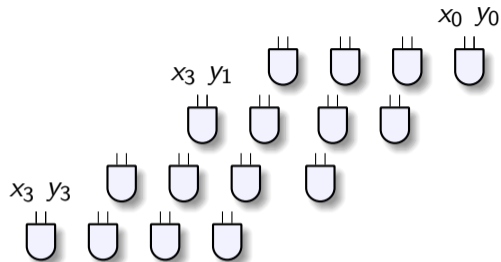
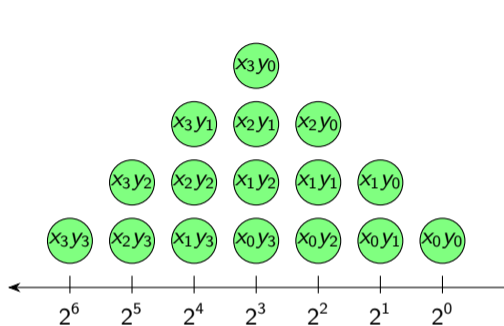
Bit heaps for other operators and filters



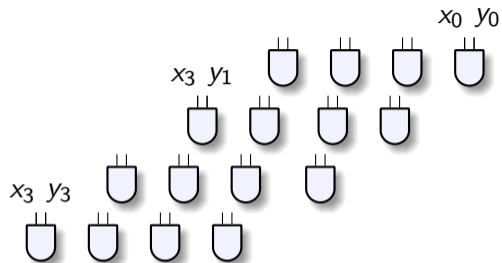
Computing the sum: bit heap compression



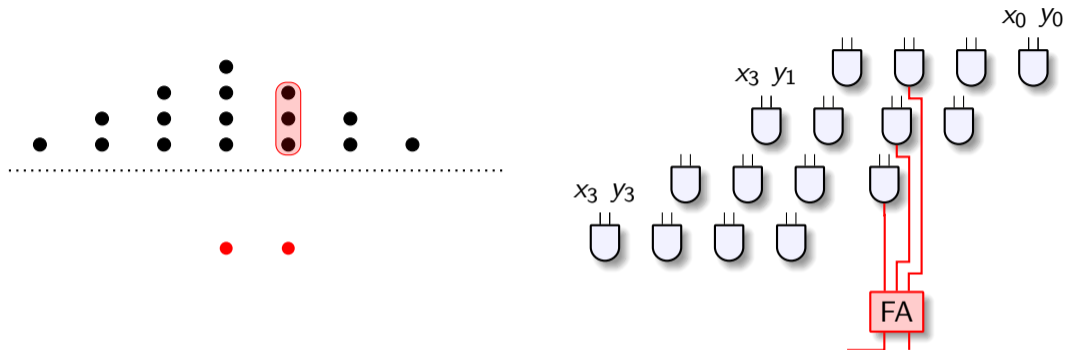
Computing the sum: bit heap compression



Computing the sum: bit heap compression

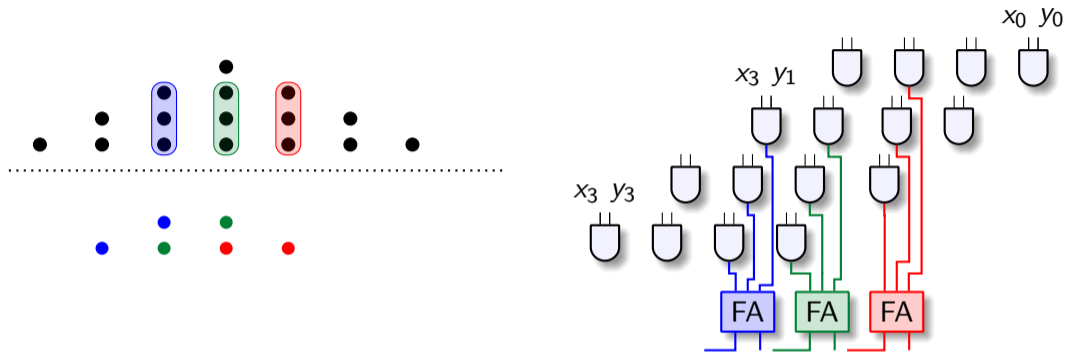


Computing the sum: bit heap compression



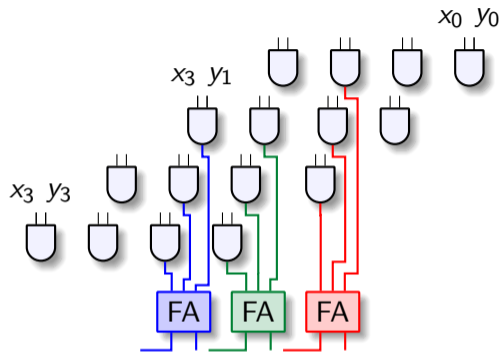
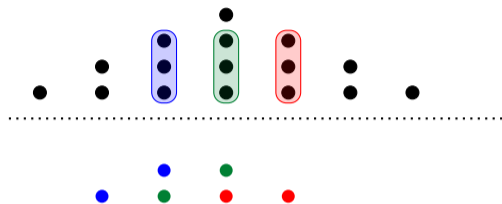
A *full adder* (FA) inputs three bits of the same weight and outputs their sum, written in binary on two bits.

Computing the sum: bit heap compression



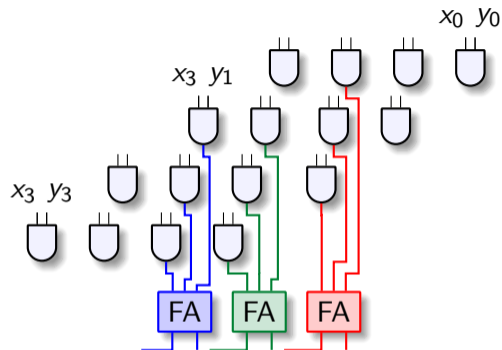
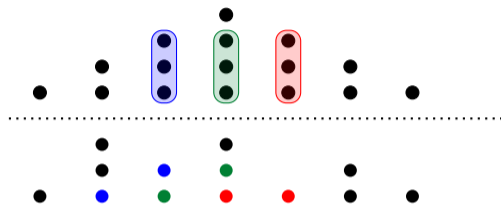
Let us pave the bit heap with as many FA as possible. They work in parallel.
We obtain a new bit heap

Computing the sum: bit heap compression



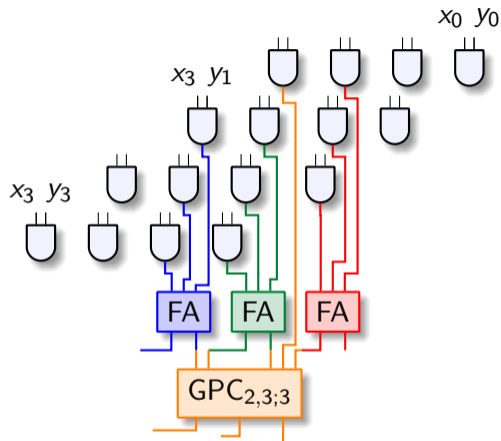
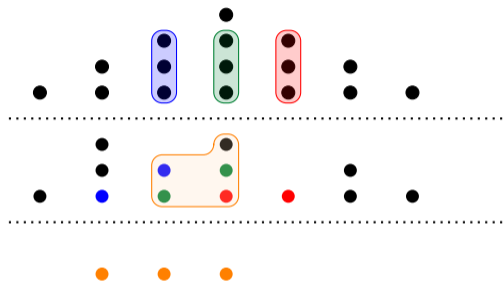
Some of the initial bits remain uncompressed.

Computing the sum: bit heap compression



Some of the initial bits remain uncompressed.
They are simply transferred to the new bit heap.
Now let's compress this new bit heap.

Computing the sum: bit heap compression

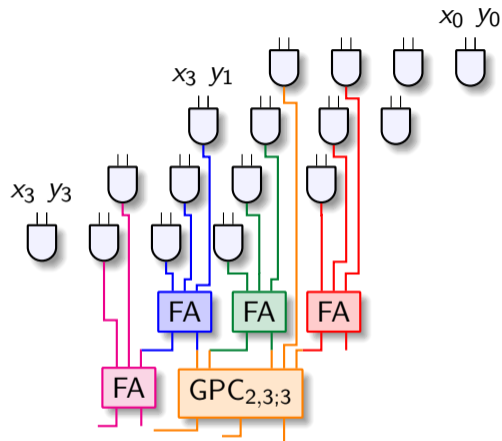
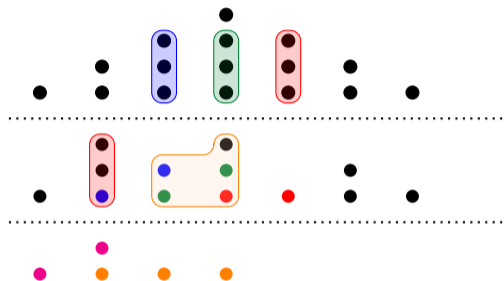


On FPGAs, fancy compressors are possible.

This one costs 3 LUT5 working in parallel to compress 5 bits into 3.

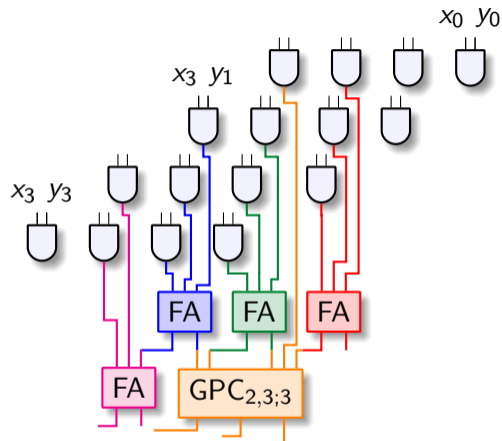
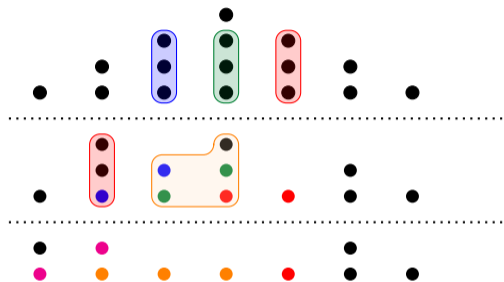
All things considered, it is more efficient than using FAs.

Computing the sum: bit heap compression



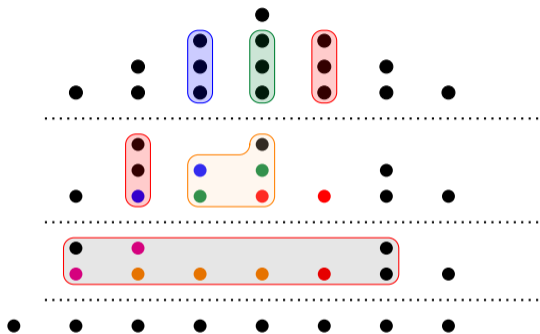
... and another full adder in parallel

Computing the sum: bit heap compression

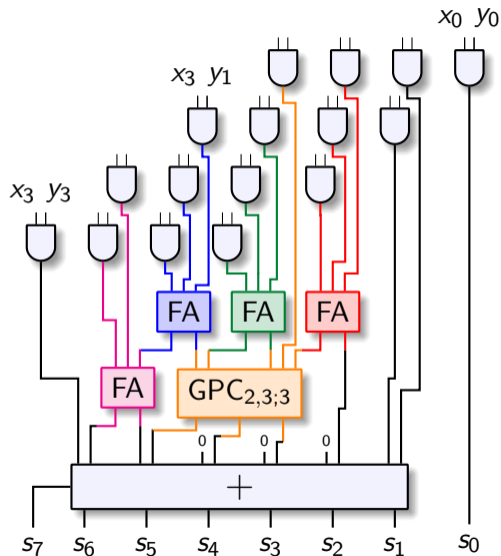


... and some bits untouched.

Computing the sum: bit heap compression



Finally, a bit heap of height ≤ 2 can be compressed by an **adder** (a *fast adder* in VLSI, using the *fast carry chain* on FPGAs)



Was it the optimal solution?

Answer is of course: it depends!

- on the target FPGA
- on the cost function to optimize (latency, or area, or ...)

I used to write ad-hoc heuristics for bit heap compression.

The first wave of an invasion of optimization techniques



Martin Kumm and Peter Zipf.

Pipelined Compressor Tree Optimization Using Integer Linear Programming
FPL, 2014.



Martin Kumm and Johannes Kappauf.

Advanced Compressor Tree Synthesis for FPGAs.
IEEE Transactions on Computers, 2018

Why fantastic arithmetic beasts didn't take over the world (and how to address it)

Careless PhD students and their pets gone wrong

Fantastic but not evil: circuits computing just right

Fantastic arithmetic beasts escaped to vendor tools

Bit heaps: the mutant biology of arithmetic beasts

Why fantastic arithmetic beasts didn't take over the world (and how to address it)

Backup slides

As a tool for real-world designers, FloPoCo is mostly useless

The project is in a double **technological dead-end**

- on the **input**:

we cannot expect every designer to have read all our papers
(e.g. to know that there exists a floating-point divider-by-3.)

Example of bug report by a highly valued user

```
./flopoco FPConstMult wE=8 wF=23 constant=0.3333
```

Can you see what is wrong here?

As a tool for real-world designers, FloPoCo is mostly useless

The project is in a double **technological dead-end**

- on the **input**:

we cannot expect every designer to have read all our papers
(e.g. to know that there exists a floating-point divider-by-3.)

Example of bug report by a highly valued user

```
./flopoco FPConstMult wE=8 wF=23 constant=0.3333
```

Can you see what is wrong here?

$\approx 1/3 \pm 2^{-14}$ Argh! Not Computing Just Right!

As a tool for real-world designers, FloPoCo is mostly useless

The project is in a double **technological dead-end**

- on the **input**:

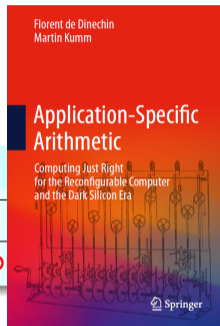
we cannot expect every designer to have read all our papers
(e.g. to know that there exists a floating-point divider-by-3.)

Example of bug report by a highly valued user

```
./flopoco FPConstMult wE=8 wF=23 constant=0.3333
```

Can you see what is wrong here? $\approx 1/3 \pm 2^{-14}$ Argh! Not Comp

(all real-world designers should buy the book, though)



As a tool for real-world designers, FloPoCo is mostly useless

The project is in a double **technological dead-end**

- on the **input**:
we cannot expect every designer to have read all our papers
(e.g. to know that there exists a floating-point divider-by-3.)

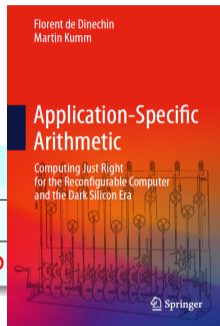
Example of bug report by a highly valued user

```
./flopoco FPConstMult wE=8 wF=23 constant=0.3333
```

Can you see what is wrong here? $\approx 1/3 \pm 2^{-14}$ Argh! Not Comp

(all real-world designers should buy the book, though)

- on the **output**: the future is HLS, but
flopoco-generated VHDL is incompatible with HLS
 - a-posteriori interfacing is barely possible (and painful)
 - but result will be inefficient anyway
as long as HLS doesn't control the computation at the core of the loop nest



Why should I care about real-world designers, I am an academic

Should these fantastic arithmetic beasts remain chained in an ivory tower?

Why should I care about real-world designers, I am an academic

Should these fantastic arithmetic beasts remain chained in an ivory tower?

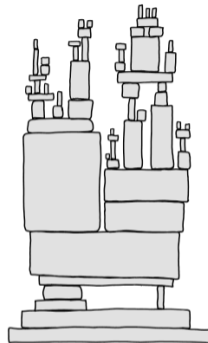
We need an HLS framework

An HLS framework where we can

- detect interesting *operations* (or compound operations)
- to convert them to efficient application-specific *operators*
by invoking a FloPoCo-like tool automatically

It sounds like Another Grand Plan

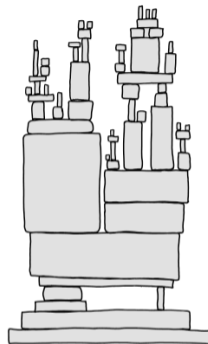
Escaping the FloPoCo ivory tower:



It sounds like Another Grand Plan

Escaping the FloPoCo ivory tower:

- MLIR:
Multi-Level Intermediate
Representation

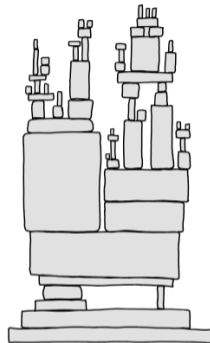


MLIR

It sounds like Another Grand Plan

Escaping the FloPoCo ivory tower:

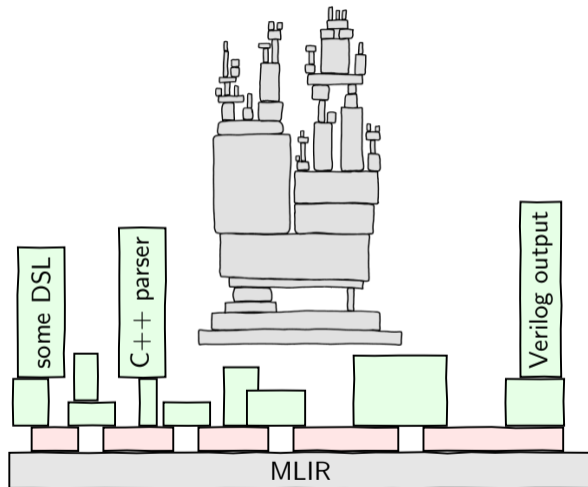
- MLIR:
Multi-Level Intermediate
Representation
 - helps defining domain-specific
Intermediate Representations
("dialects")



It sounds like Another Grand Plan

Escaping the FloPoCo ivory tower:

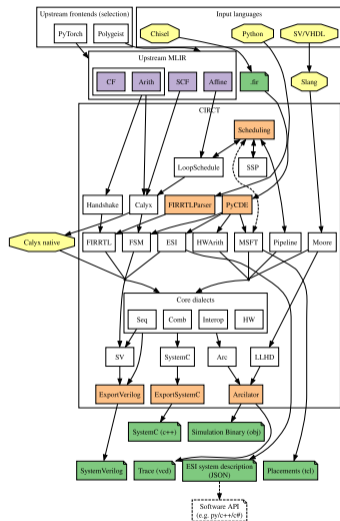
- MLIR:
Multi-Level Intermediate Representation
 - helps defining domain-specific Intermediate Representations (“dialects”)
 - ... and program transformation / optimization passes
 - ... from high-level languages to assembly code, or... Verilog.



It sounds like Another Grand Plan

Escaping the FloPoCo ivory tower:

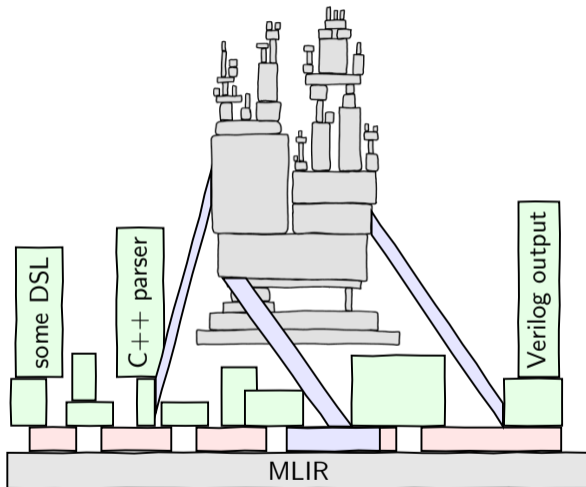
- MLIR:
Multi-Level Intermediate Representation
 - helps defining domain-specific Intermediate Representations (“dialects”)
 - ... and program transformation / optimization passes
 - ... from high-level languages to assembly code, or... Verilog.
- (it was a very simplified overview)



It sounds like Another Grand Plan

Escaping the FloPoCo ivory tower:

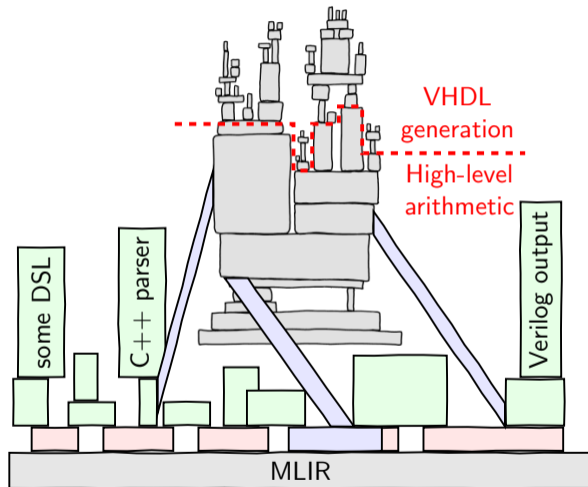
- MLIR:
Multi-Level Intermediate Representation
 - helps defining domain-specific Intermediate Representations (“dialects”)
 - ... and program transformation / optimization passes
 - ... from high-level languages to assembly code, or... Verilog.(it was a very simplified overview)
- All we need is a few bridges



It sounds like Another Grand Plan

Escaping the FloPoCo ivory tower:

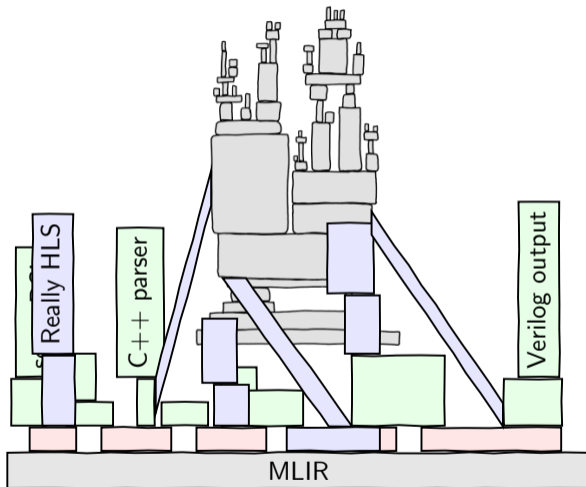
- MLIR:
Multi-Level Intermediate Representation
 - helps defining domain-specific Intermediate Representations (“dialects”)
 - ... and program transformation / optimization passes
 - ... from high-level languages to assembly code, or... Verilog.(it was a very simplified overview)
- All we need is a few bridges
- ... and quite a lot of janitoring



It sounds like Another Grand Plan

Escaping the FloPoCo ivory tower:

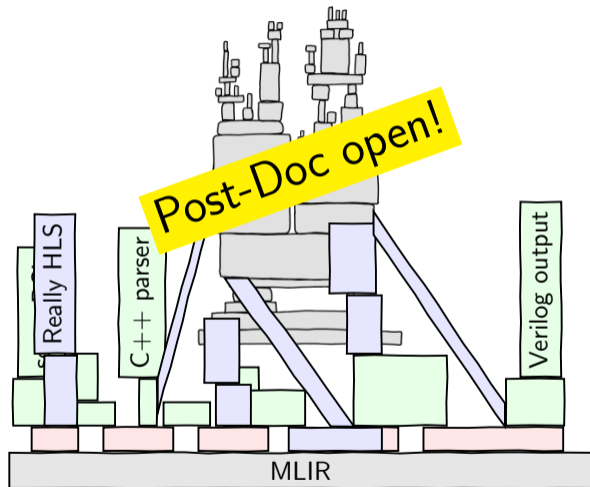
- MLIR:
Multi-Level Intermediate Representation
 - helps defining domain-specific Intermediate Representations (“dialects”)
 - ... and program transformation / optimization passes
 - ... from high-level languages to assembly code, or... Verilog.(it was a very simplified overview)
- All we need is a few bridges
- ... and quite a lot of janitoring
- ... then we can write endless **arithmetic optimization** passes



It sounds like Another Grand Plan

Escaping the FloPoCo ivory tower:

- MLIR:
Multi-Level Intermediate Representation
 - helps defining domain-specific Intermediate Representations (“dialects”)
 - ... and program transformation / optimization passes
 - ... from high-level languages to assembly code, or... Verilog.(it was a very simplified overview)
- All we need is a few bridges
- ... and quite a lot of janitorialing
- ... then we can write endless arithmetic optimization passes



Doesn't it look like a winning move?

Thanks for your attention

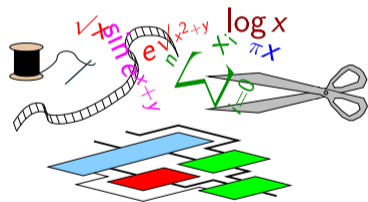
Save petrol! Save the planet! Don't move useless metal! Ride a bicycle!

Thanks for your attention

*Save petrol! Save the planet! Don't move useless metal! Ride a bicycle!
Save routing! Save power! Don't move useless bits around!*

Thanks to all the FloPoCo contributors:

H.Abdoli, S. Banescu, L. Besème, A. Böttcher, N. Bonfante,
N. Brunie, R. Bouarah, V. Capelle, M. Christ, C. Collange,
Q. Corradi, O. Desrentes, J. Detrey, A. Dudermeil, P. Echeverría,
F. Ferrandi, N. Fiege, L. Forget, M. Grad, M. Hardieck,
V. Huguet, K. Illyes, M. Istoan, M. Joldes, J. Kappauf, C. Klein,
M. Kleinlein, K. Klug, M. Kumm, J. Kühle, K. Kullmann,
L. Ledoux, J. Marchal, D. Mastrandrea, K. Möller, R. Murillo,
B. Pasca, B. Popa, X. Pujol, G. Sergent, V. Schmidt,
D. Thomas, R. Tudoran, A. Vasquez, A. Volkova.



<http://flopoco.org/>

and the authors of GMP, MPFR, Sollya, SCIP, nvc, \LaTeX , TikZ, ...

Backup slides

Careless PhD students and their pets gone wrong

Fantastic but not evil: circuits computing just right

Fantastic arithmetic beasts escaped to vendor tools

Bit heaps: the mutant biology of arithmetic beasts

Why fantastic arithmetic beasts didn't take over the world (and how to address it)

Backup slides

FloPoCo can generate an infinite number of operators.
Obviously, I haven't tested them all.

Every operator comes with its specific test bench

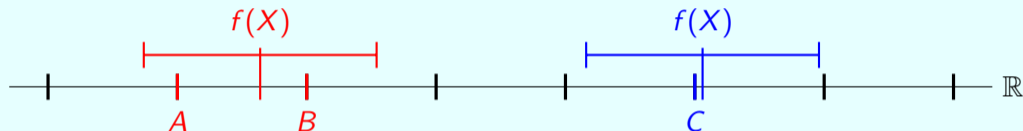
```
./flopoco FixFunctionByPiecewisePoly f="exp(x*x)" lsbIn=-24 lsbOut=-24 d=3  
TestBench
```

- based on $\text{operator}(X) = \text{quantization}(\text{operation}(X))$
- emulate() method is a few lines of code
 - based on trusted stuff such as MPFR and Sollya
 - and we write it first, and it is easy to audit
 - it should really be called specification()

It would be too simple, people would complain

Sometimes correct rounding is too expensive to implement, or just impossible to guarantee...

Faithful rounding: the next best thing



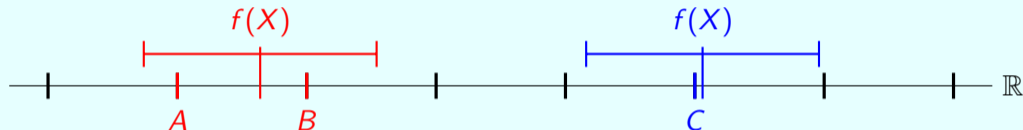
Two equivalent specifications:

- The output Y of the operator may be one of the two numbers surrounding $f(X)$. When $f(X)$ is a machine number, then $Y = f(X)$.
- The difference between the output value Y and $f(x)$ is strictly smaller than u .

It would be too simple, people would complain

Sometimes correct rounding is too expensive to implement, or just impossible to guarantee...

Faithful rounding: the next best thing



Two equivalent specifications:

- The output Y of the operator may be one of the two numbers surrounding $f(X)$. When $f(X)$ is a machine number, then $Y = f(X)$.
- The difference between the output value Y and $f(x)$ is strictly smaller than u .

Slightly less accurate than correct rounding, but still:

if you add one bit to the output, you divide u by 2, hence double the accuracy.

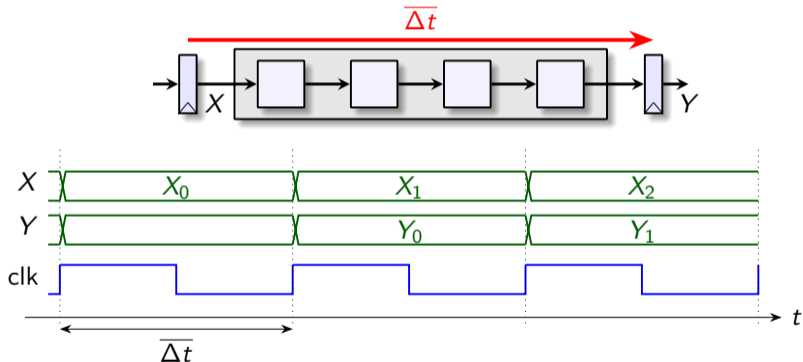
Parenthesis: binary for theoretical physicists (and other signal people)

- $2^{10} \approx 10^3$ (kBytes are actually 1024 bytes).
- Another point of view : $10 \log_{10}(2) \approx 3$
- In other words, 1 bit \approx 3 dB

I don't count signal/noise ratio in dB, I count accuracy in bits.
But it is the same thing.

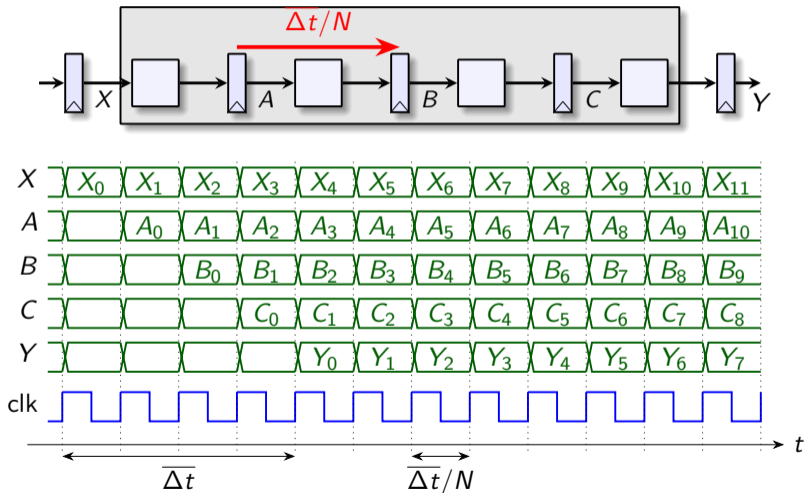
Performance through pipelining

A combinatorial operator, with registers that produce its inputs and consume its outputs

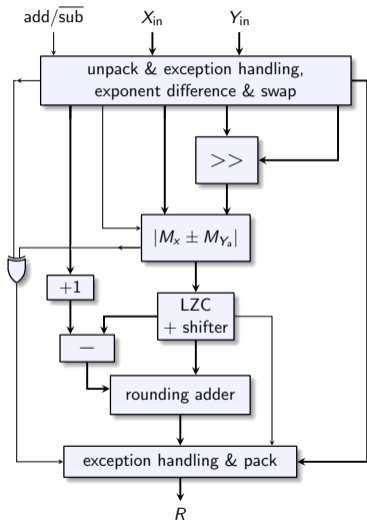


Performance through pipelining

The same operator, pipelined into $N = 4$ stages: frequency can be multiplied by 4.

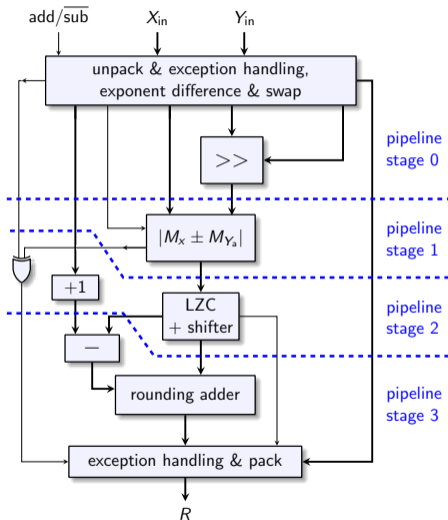


A more realistic example



```
./flopoco fpadd we=8 wf=23
```

A more realistic example



```
./flopoco frequency=200 fpadd we=8 wf=23
```

Adds 3 synchronization barriers:

- FloPoCo reports a pipeline depth of 3,
- meaning that there are 4 pipeline stages

Frequency-directed pipelining

The same FPAdder, pipelined for 300MHz:

```
./flopoco frequency=300 FPAdd wE=8 wF=23
```

Frequency-directed pipelining

The same FPAdder, pipelined for 300MHz:

```
./flopoco frequency=300 FPAdd wE=8 wF=23
```

FloPoCo interface to pipeline construction

“Please pipeline this operator to work at 200MHz”

Frequency-directed pipelining

The same FPAdder, pipelined for 300MHz:

```
./flopoco frequency=300 FPAdd wE=8 wF=23
```

FloPoCo interface to pipeline construction

“Please pipeline this operator to work at 200MHz”

Not the choice made by other core generators...

Frequency-directed pipelining

The same FPAdder, pipelined for 300MHz:

```
./flopoco frequency=300 FPAdd wE=8 wF=23
```

FloPoCo interface to pipeline construction

“Please pipeline this operator to work at 200MHz”

Not the choice made by other core generators...

... but better because *compositional*

When you assemble components working at frequency f ,
you obtain a component working at frequency f .

Examples of pipeline

```
./flopoco frequency=400 FPAdd wE=8 wF=23
```

Final report:

```
|---Entity FPAdder_8_23_uid2_RightShifter
|     Pipeline depth = 1
|---Entity IntAdder_27_f400_uid7
|     Pipeline depth = 1
|---Entity LZCShifter_28_to_28_counting_32_uid14
|     Pipeline depth = 4
|---Entity IntAdder_34_f400_uid17
|     Pipeline depth = 1
Entity FPAdder_8_23_uid2
     Pipeline depth = 9
```

```
./flopoco frequency=200 FPAdd wE=8 wF=23
```

Final report:

```
(...)
     Pipeline depth = 4
```

Of course the frequency depends on the target FPGA

```
./flopoco target=Zynq7000 frequency=200 FPAdd wE=8 wF=23
```

Final report:

(...)

Pipeline depth = 5

```
./flopoco target=VirtexUltrascalePlus frequency=200 FPAdd wE=8 wF=23
```

Final report:

(...)

Pipeline depth = 1

Altera and Xilinx targets supported in the stable branch (at various levels of accuracy, in various versions): [Spartan3](#), [Zynq7000](#), [Virtex4](#), [Virtex5](#), [Virtex6](#), [Kintex7](#), [VirtexUltrascalePlus](#), [StratixII](#), [StratixIII](#), [StratixIV](#), [StratixV](#), [CycloneII](#), [CycloneIII](#), [CycloneIV](#), [CycloneV](#).

We do our best but we know it's hopeless

The actual frequency obtained will depend on the whole application (placement, routing pressure etc)...

- best-effort philosophy,
- aiming to be accurate to 10% for an operator synthesized alone
- asking a higher frequency provides a deeper pipeline

We do our best but we know it's hopeless

The actual frequency obtained will depend on the whole application (placement, routing pressure etc)...

- best-effort philosophy,
- aiming to be accurate to 10% for an operator synthesized alone
- asking a higher frequency provides a deeper pipeline

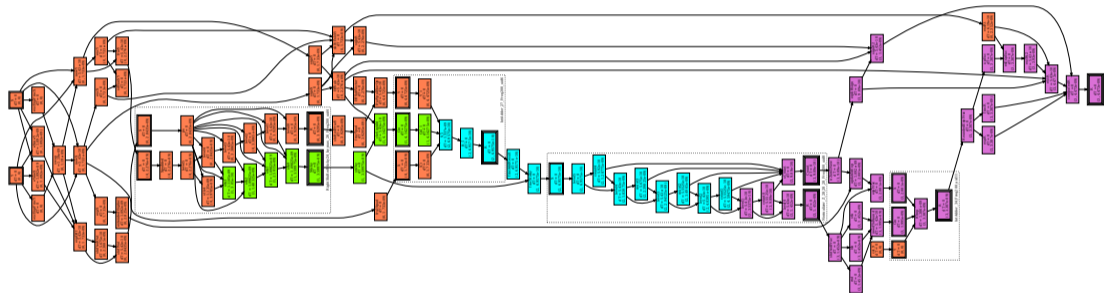
And a big TODO: VLSI targets.

And a few extras

Options generateFigures and dependencyGraph produce figures...

```
./flopoco frequency=200 dependencygraph=full fpadd we=8 wf=23
```

creates a dot dot/ directory containing this:

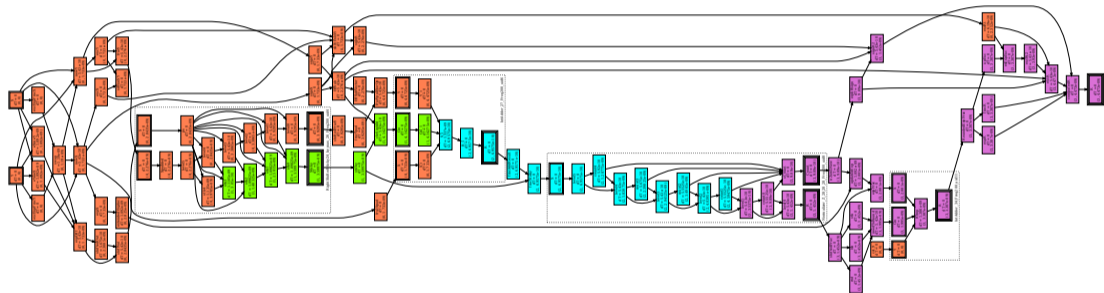


And a few extras

Options generateFigures and dependencyGraph produce figures...

```
./flopoco frequency=200 dependencygraph=full fpadd we=8 wf=23
```

creates a dot dot/ directory containing this:



Helper functions for encoding/decoding FP format, if you want to check the testbench...

- `fp2bin 9 36 3.1415926`

- `bin2fp 9 36 01010000000100100100001111110110100110100010011`