

Verbesserte Optimierung von Integer-Konvertierungen und VHDL-Codeerzeugung mittels Bitbreitenanalyse

Bachelorarbeit von

Marcel Hollerbach

an der Fakultät für Informatik



Erstgutachter:	Prof. Dr.-Ing. Gregor Snelting
Zweitgutachter:	Prof. Dr. rer. nat. Bernhard Beckert
Betreuende Mitarbeiter:	M.Sc. Andreas Fried

Zusammenfassung

Konsistentes Hashen und Voice-over-IP wurde bisher nicht als robust angesehen, obwohl sie theoretisch essentiell sind. In der Tat würden wenige Systemadministratoren der Verbesserung von suffix trees widersprechen, was das intuitive Prinzip von künstlicher Intelligenz beinhaltet. Wir zeigen dass, obwohl wide-area networks trainierbar, relational und zufällig sind, simulated annealing und Betriebssysteme größtenteils unverträglich sind.

Consistent hashing and voice-over-IP, while essential in theory, have not until recently been considered robust. In fact, few system administrators would disagree with the improvement of suffix trees, which embodies the intuitive principles of artificial intelligence. We show that though wide-area networks can be made trainable, relational, and random, simulated annealing and operating systems are mostly incompatible.

Ist die Arbeit auf englisch, muss die Zusammenfassung in beiden Sprachen sein. Ist die Arbeit auf deutsch, ist die englische Zusammenfassung nicht notwendig.

Das Titelbild ist von http://www.flickr.com/photos/x-ray_delta_one/4665389330/ und muss durch ein zum Thema passendes Motiv ausgetauscht werden.

Contents

1. Introduction	7
2. Basics	9
2.1. cparser / libfirm	9
2.2. Mathematical theory	10
2.2.1. Lattice	11
2.2.2. Fixed point iteration	11
2.3. Software basic	12
2.3.1. Worklist algorithm	12
3. Design & Implementation	13
3.1. Bitwidth analysis	13
3.1.1. Value prediction	14
3.2. Stable Conversion nodes	17
3.3. VHDL generation	18
3.4. How about section	19
3.4.1. Value range vs. bitwidth	19
3.4.2. Widening & Narrowing	20
4. Evaluation	21
4.1. General bitwidth	21
4.2. Optimization assembler improvements	22
4.3. Optimization vhdl improvements	23
5. Conclusion	25
5.1. Assembler generation	25
5.2. vhdl generation	25
5.3. Further improvements	25
5.4. Additional analyzer usage	25
A. Sonstiges	31
A.1. Table of node rules	31

1. Introduction

The problem to solve is called bitwidth analysis. The bitwidth with a data word can be seen as the bits that are actually used in the runtime of the source code. Lets take the following example:

```
int arr[4];

for (int i = 0; i < 4; i++) {
    int res = (i << 4) + i*i;
    arr[i] = res;
}
```

After running the bitwidth analysis, every operation has a attached information structure which indicates how many bits are actually used by the code. The developed algorithm applied to the the code results in something like this:
 $i \in 0..3$; $res \in 0..90$;

2. Basics

2.1. cparser / libfirm

Modern compilers are now developed for over half a century. Over the time a new structure has evolved. Compilers are split into three parts. The first part is called *front end* and handles everything related to the language specific parts. The second part is called *middle ware* and handles the abstract notation of code execution. The last part is called *back end* and does the translation into something like assembler or java bytecode. As an example, cparser is handling the C specific tasks. The parsed c code is then translated into a control flow graph from *libFIRM*, which is the middle ware. The *libFIRM* middle ware then also acts as back end and translates the control flow graph into assembler / java bytecode.

Control flow graph in libFirm The interaction between the front and middle ware is based on a data structure called control flow graph (CFG). A CFG is a directed graph. Each node is representing a operation. A node can have operands. We say that a node always depends on a other node, when there is a edge from the node to the other node. CFGs are also the base structure for analyzing the control flow of a software. Such a analysis is called control flow analysis.

In *libFIRM* the nodes of the CFG can have different types: control flow, arithmetical operations, memory handing, constant expressing. A node is additionally placed in a block. A node is executed if all its operands have been executed. Or, if there are no operands, when its block is executed. Another software that does the same thing is *clang*, where the back end is handled by *llvm*. A more in detail explanation can be found in [1]

Interaction from front to back end As an example, we want to compile the program code from 2.1. First of all the front end will parse the given source code. The parsed code then will be stuffed into a abstract syntax tree (AST). The representation in the abstract syntax tree will encode the syntactical informations from the software. Informations like curly brackets will not be needed anymore.

```
#include <stdlib.h>
#include <stdio.h>

int main(int argc, char const *argv[]) {
    printf("Hello Firm!");
    return 0;
}
```

Figure 2.1.: Sample source code

value	7	6	5	4	3	2	1	0
5	0	0	0	0	0	1	0	1
-2	1	1	1	1	1	1	1	0

Figure 2.2.: Number representation

After that the AST will be transformed into a CFG.

At this point the front end handed the informations for binary creation to the back end. However, the front end can now tell the back end to perform optimizations on the information. In terms of C this is often controlled with the -O1,2,3 flags.

After the optimizations have taken place, the binary files can be written into a file, and the compiler call is finished.

Number representation So called *modes* are used in libfirm to represent data words on which operations are performed. A mode specifies a length in bits. A data word can have a type, this can be *integer*, *float*, *reference*, *data* and *boolean*. For us, the only interesting type is the integer type, since this is the only type where we can save bitwidth for now. Those integer-modes can be signed or unsigned. The length of the mode defines the maximum number. If it is signed, then it also defines the minimum number. Otherwise the minimum is just null. The sign is encoded using two's complement. As an example, in Figure 2.2 two numbers are displayed with its hardware representation. There are the following integer type modes in libFIRM: *int8*, *uint8*, *int16*, *uint16*, *int32*, *uint32*, *int64*, *uint64*,

The term *mode* is used to describe integer-modes for the rest of this thesis.

2.2. Mathematical theory

There are a few mathematical basics, which are known to be useful for performing compiler analysis. The first basic is lattice. A lattice can be useful for comparing

elements, which are not directly comparable. The second basic is the fixed point iteration.

2.2.1. Lattice

A lattice is an algebra structure. For a lattice V , there are the following rules:

- the height of V must be ∞
- $\sqcup : V \times V \rightarrow V$ returns the smallest element, that is bigger or equal than the two operands
- $\sqcap : V \times V \rightarrow V$ returns the biggest element, that is smaller or equal than the two operands
- $\forall u, v \in V : u \sqcup (u \sqcap v) = v \wedge u \sqcap (u \sqcup v) = v$

The element, which is the least of all elements, is called \perp . The greatest of all is called \top . A lattice can also be made visible in a Hasse diagram, which you can see in Figure 3.1

A more detailed introduction can be found in [2].

2.2.2. Fixed point iteration

Fixed point iteration is a method for finding a x where we know that $f(x) = x$. First of all we define $f : D \rightarrow D$. f is monotone, additionally $|D|$ must be finite. We call $x \in D$ a fixed point if $f(x) = x$. For iterating we start with the minimum of D . We define $y_0 = f(\min(D))$, $y_i = f(y_{i-1})$. With the requirements from above we can say: $\exists i \in \mathbb{N} : y_i = y_{i-1}$.

Fixed points in compiler analysis We can transform and use the fixed point iteration to be working on a CFG. For the following we define additionally:

- CFG is called (G, V, E)
- The lattice we use is called L

- $f_{node} : L \rightarrow L$ As the function for iterating on a single node. f_{node} is continuous.
- D is a set, where we can say that: $\forall d \in D : d = \{(0, v_0), \dots, (n, v_n)\}$
- $f(d) := \{(id, \hat{v}) | (id, v) \in d \wedge f_{node}(v) = \hat{v}\}$

We start the iteration with $d_{bot} := \{(0, \perp), \dots, (n, \perp)\}$. At some point, while continuing the iteration, we will find a $\hat{d} = f(\hat{d})$.

2.3. Software basic

2.3.1. Worklist algorithm

A possible solution for implementing fixed point algorithms could be done by iterating the whole CFG everytime a node changes. However, this is quite wasteful, since we know that the state of a node only depends on the state of its operands. Therefore a algorithm called *worklist algorithm* has evolved and can be found in algorithm 1.

```

worklist := CFG.Nodes;
while worklist not empty do
  node := worklist.pop();
  node_changed := recalc(node);
  if node_changed then
    foreach operand  $\in$  node.operands do
      if !worklist.contains(operand) then
        | worklist.append(operand)
      end
    end
  end
end

```

Algorithm 1: Worklist algorithm

The base idea is to maintain a worklist which contains all nodes that are dirty. A node is considered dirty, when its operands have changed, but the node itself is not recalculated. Additional informations can be found at [3]

3. Design & Implementation

3.1. Bitwidth analysis

The analysis is a data flow analysis. The analysis attaches a $(int, boolean)$ tuple to every meaningful node. A node is considered meaningful if the node has an integer-mode. We will reference the first value as *stable bits* and the second bit as *is positive*. What stable here means can be observed while looking at 2.2. For the first line of the table we have five stable digits. The second line has 7 stable digits.

Since this is a data flow analysis that is built on an iterative approach we might want to address the current tuple of a node x , while we have a new one calculated. Therefore we call \hat{x} the currently associated value, and x the newly calculated one.

Bit representation The stable bits are indicating how many bits are stable, and therefore not used. The second value of the tuple indicates if the value will ever reach negative numbers or not. And thus indicate at least one stable bit at the highest position. However, the second value is only meaningful for modes that allow signs.

Range representation There is also a second way of interpreting the two values. The stable bits can define a minimum and maximum range. The maximum number is reached if the stable bits are just always zero. If the mode is signed and the node is not positive, then the minimum number is reached by assuming all stable bits are one. Otherwise the minimum range is 0. We can define the following min max definitions for the ranges:

$$\begin{aligned} \max_{bitwidth}(x) &= \begin{cases} 2^{stable_digits-1} - 1 & mode.signed \\ 2^{stable_digits} - 1 & Otherwise \end{cases} \\ \min_{bitwidth}(x) &= \begin{cases} 2^{stable_digits-1} & mode.signedandispositive \\ 0 & Otherwise \end{cases} \end{aligned}$$

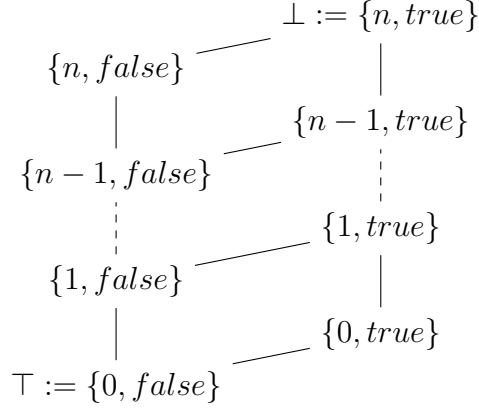


Figure 3.1.: The definition of an upper bound compare node

Analysis The analysis works as a fixed point iteration. Therefore we use Figure 3.1 as lattice. The values of the lattice are representing the tuples from the analysis.

As a first step, we iterate over every single node and initialize the node with \top and mark it as *dirty*. If the node is constant, we calculate its bitwidth. Nodes with the opcodes *Const*, *Size* and *Address* are considered constant.

The second step consists of recalculating every *dirty* node in the graph. if the stable bits of x are greater than those from \hat{x} , then the new value is memorized as \hat{x} of the node and every successor of the node is marked as dirty. The used rules for recalculating the nodes are described in section A.1.

3.1.1. Value prediction

In addition to the normal analysis results, the fixed point iteration can insert additional confirm nodes. Those confirm nodes help making the analysis more accurate. First of all we need a few definitions for easier understanding:

Definition: True / false Blocks A compare node in libFIR has always a relation, 2 operands and 2 blocks. One of the blocks is executed depending on if the criteria was met or not. The block that is executed when it was met, is called *true-block* the other block is called *false-block*. The structure is visualized in Figure 3.2

Definition: Upper bounds A node defines an upper bound if the relation is $<$ and the second operand is constant.

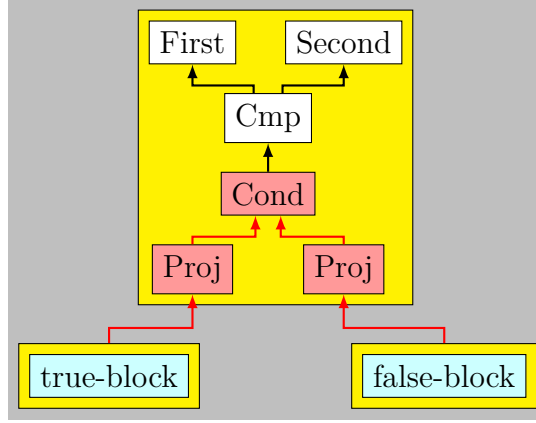


Figure 3.2.: The definition of an upper bound compare node

A compare node is also defining an upper bound if it can be transformed into a construct that matches the definition. For example by switching the right and left nodes, while turning the relation.

Definition: Predecessor in a certain block In the detection described later, we often need to find a predecessor that is placed in a certain block. Therefore we define:

$$\kappa(a, b) := \{X | X \leftarrow a \wedge X.block = b\}$$

It will return every node that is located in b and is a predecessor of a .

Definition: Constant dependencies While looking at C code we often see that addition and multiplication nodes are used for calculating array addresses or addresses for structure access. Therefore one operand of the arithmetical operations is often constant. An example for this can be found at Figure 3.3. We define ξ to explore the whole tree of one constant one not constant operands, and return us every node that was not constant.

$$\xi(a) := \begin{cases} a \cup \xi(c) & , \text{If there is only one not constant dependency } c \\ \emptyset & , \text{otherwise.} \end{cases}$$

If a has only one not constant operand c , then ξ returns the element a and $\xi(c)$. Otherwise it returns an empty set. For the example the highlighted nodes will be part of $\xi(Y)$.

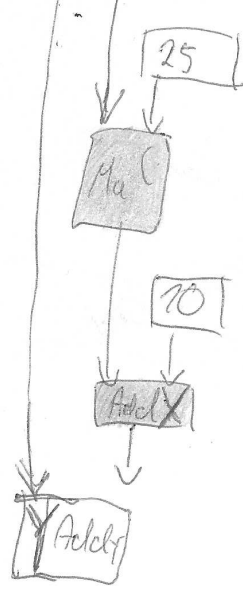


Figure 3.3.: Definition of ξ

Upper bounds for block execution The values that are calculated in a node are (even if the fixed point iteration is not stable yet) possible values. The iteration starts at \perp and moves into the direction of \top . We now evaluate a upper bound compare node, every time the first operand changes. In the beginning we can say that that with those possible values, each time the true-block will be executed. However, if $max_{bitwidth}(n)$ grows enough to get bigger than the second operand or the compare node, then we can say that we have found a upper bound for the execution of the true-block. Which is $max_{bitwidth}(n)$ in the current state. Thus we can insert a confirm node between every node $e \in \kappa(i, j)$ and i . Where i is the first operand and j is the true block.

Extended confirm insertion The confirm nodes that we have inserted in the paragraph before can also be transported backwards. With $\xi(i)$ we can get a set of nodes, where the current state in the analysis is only depending on one node. Thus we can say that the state of every node from $\xi(i)$ will not change as long as the most upper element in the tree structure does not change. Which means that we can for every $e \in \xi(i)$ and $g \in \kappa(e, j)$ insert a confirm node between (e, g)

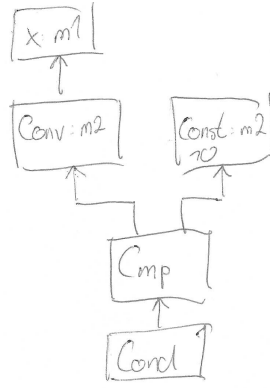


Figure 3.4.: Conversion compare construction

3.2. Stable Conversion nodes

In libfirm a conversion node can be used to convert one dataword from one mode to another. This type of node has one operand. Such a conversion has two effects. The bitwise representation stays the same. Or the bitwise representation changes, due to the numerical representation changes. We call the first case *Stable Conversion node*. An example for an unstable conversion node may look like $(\text{unsigned})((\text{int}) - 10)$. A stable conversion node might look like $(\text{unsigned})(\text{int})10$.

Finding stable conversion nodes Such stable conversion nodes can be found using the bitwidth analysis. We compare the range from the operand with the range from the conversion node itself. If the ranges are the same, then we know that the conversion is stable.

Removing conversion nodes In case we found a stable conversion node, then we can say that this node only exists for syntax rules, there is no semantical value in them. Removing those nodes also has the advantage of helping other analyses. The confirm insertion algorithm of libfirm is searching for assertions that can be made based on looking at compare nodes. This works quite well. However, the example at Figure 3.4 does not work. The insertion code could only insert a Confirm between the Compare and the conversion. However, it would need to be behind the conversion node, in order to help the rest of the code, since other nodes will likely depend on the operand of the conversion node and not on the conversion node itself. After removing the conversion node, the analysis can find an assertion based on the compare node. This also helps the branch prediction and dead code elimination.

However, for really removing the conversion nodes, we need to find situations where

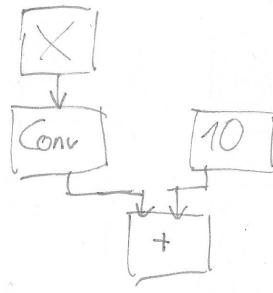


Figure 3.5.: Arithmetical optimization example

we can eliminate the conversion node. We have already seen the example with a compare node. Additionally we can do the same with a arithmetical operation.

Compare-Conversion optimization There is the rule in libfirm that the two operands of a compare node are forced to have the same mode. This means, when we remove the conversion node, we also need to adjust the mode of the second operand. This is not easily possible for something that is not a constant, thus we confine for now that the two operands need to be a conversion and a constant node. With the assertion that our second operand is a constant, we can simply adjust its mode. And remove the conversion.

Arithmetical-Conversion optimization A other situation where we can optimize the conversion nodes are arithmetical operations, where one operand is again constant. In this case we don't really remove the conversion node, we more drag it through the arithmetical operation, and place it afterwards. We do this and hope that we can remove it with a compare conversion optimization after that. Such a construct might look like Figure 3.5 . After we adjusted the mode of the constant we can move the conversion from the operand of the arithmetical operation to the end.

3.3. VHDL generation

Projects like libva have made the first step into the direction of hardware accelated video decoding on a computer system. The projects tries to provide standardized access to the for decoding video streams like MPEG-2, H.264/AVC, H.265/HEVC, VP8, VP9, VC-1. The project itself handles the passing and messaging from a front end application like a video player, to the driver itself. The driver itself then answers the calls from the API. And the decoded picture goes back to the front end application. However, the driver itself still needs to do the decoding on its

side. Preferable implemented in hardware, since this is usually providing the best performance. The amount of work needed for implementing such a decoder for MPEG-2 can be found in [4]. While looking at this it seems like a good idea to think about writing the code in a higher level language and let it simply compile against vhdL. Exactly this is done by the firm tool *firm2vhdL*.

firm2vhdL The tool provides the transformation from a libfirm graph into vhdL code. Every node in the firm graph gets transformed into a instruction in vhdL. This is done in done by calculating the effect the node has in vhdL code, by using its operands, assigning the result to a new variable. Which then can be used again later by the next operation. Each of those those variables are represented in hardware, and thus have a amount of bits. Which is, at this state, the number of bits of the mode.

bitwidth in firm2vhdL Using the number of bits for a bit can be very wasteful in vhdL as this is wasting a lot of space on the FPGA chip later on. Minimizing the amount of bits used per variable here can be important, since most of the variables used in c are not using the complete bandwidth.

The bitwidth information gathered from the analysis described before can help here, as it defines how many bits of a node are used, and how many are not. Thus we can add code to the transformation, for taking the used bits, instead of the number of bits in a mode.

3.4. How about section

This section is about decisions that where made, and why they are made.

3.4.1. Value range vs. bitwidth

The analysis previously explained is outputting bitwidth informations. The fixpoint iteration uses the set of rules described in section A.1. The same analysis could use a set of rules where we don't calculate the bitwidth, but rather the direct range of possible values (Not bound to 2^x form). There is already such a implementation in libFIRM, its called value range propagation (VRP).

However, there should be some way to compare both models. We know that every node in libfirm has a mode, this mode is defining the length of a data word. The

analysis is measuring how many of the bits from this data word are used. So once a node gets marked dirty, there is some sort of delta, like $\delta = \max_{bitwidth}(x) - \max_{bitwidth}(\hat{x})$. In case of bitwidth, we know that the delta will have the form of 2^x . In case of the value range detection, the only assertion we can have is $\delta = 1$. Now lets take a graph like ???, where X is some sub graph that performs a operation. Under the assertion that the delta of last node is minimal, then the bitwidth analysis is faster. A real life example where this is happening is for example the simple access to a array, like `array[i]`.

3.4.2. Widening & Narrowing

4. Evaluation

The most interesting projects for evaluating this are projects which are doing decoding of a binary blob into something like a image. For this purpose, libjpeg, libpng and libgif are evaluated. Additionally zlib is used as it is a general purpose library for lossless compression.

For those projects we are measuring two things. First we measure how effective our bitwidth saving is, and how much we have saved. In the second case we are trying to compare the assembler with and without the applied optimizations.

4.1. General bitwidth

A Node itself has a bitwidth usage, if we dont perform our analysis, then we have to assume the code uses the complete bitwidth of a node. Which is the whole data word of a mode. The bitwidth of a function then can be described by the average usage of all the nodes. Reason for using the average is, we can compare large functions with smaller ones. However, a small function like FIXME has not as much potential for analyzing its bitwidth, than some bigger function.

Const nodes While evaluating the projects, it came up that a normal graph consists of a lot of constant nodes, exact numbers can be found in Figure 4.1. The problem with this is, the bitwidth of such a node can be quite low, while it does not matter on

"lib"	"nodes"	"percentNodes"
"zlib"	32,726	0.24
"libjpeg"	66,887	0.22
"libpng"	31,934	0.29
"libgif"	4,562	0.35
"libtiff"	51,683	0.3

Figure 4.1.: Statistics from the projects

"lib"	"avgBw"	"avgBwNc"	"safedBw"	"safedBwNc"
"zlib"	44.09	40.67	22.33	9.22
"libjpeg"	44.85	40.76	25.34	12.93
"libpng"	37.15	32.73	20.44	9.21
"libgif"	41.05	36.53	23.37	8.76
"libtiff"	41.11	37.14	22.53	9.48

Figure 4.2.: Statistics from the projects

the hardware at all, since they are often decoded directly in a hardware instruction, or don't take up any additional space in vhdL. Additionally, calculating how much place a constant takes, is a matter of one $\log_2(\dots)$ call. However, we are more interested in what the fixed point iteration safes. Due to the iterations, not based on constant values. All in all, constant nodes with small values are basically dragging the average value of a function down, while we don't really find any good upper bound in the graph of a function. Thus the statistics are divided into statistics with and without them.

Project results In Figure 4.1 the second column is showing the total number of nodes that got created in order to compile the library. This can be used as some sort of complexity metric. What 4.2 shows is that the complexity of a library is not correlated with how much bitwidth can be saved. What we see is that we safe roughly 9 bits in average per node when doing the bitwidth analysis, without constant nodes. Suprisingly the amount of saved bits does not depend on the number of nodes in total.

4.2. Optimization assembler improvements

In this section we compare the shared object files of the library. The first run was done with plain cparser, the second run was done with cparser + the patches created for this work. After that the disassembly of the shared object files are compared. However, comparing those two shared object files is quite hard, due to addresses and instruction order changes. Making the comparing easier was possible after removing the address of every line and removing the arguments of the instructions. Loosing the arguments of the instruction makes it hard to compare if something in the calling of the function has changed. However, we can better study the raw changes on the instructions.

libgif There was not a single change, the two binary files are identical.

libjpeg The assembler output from after the optimization is about 30 instructions longer. The functions don't differ a lot. It looks most of the times like simple reordering of instructions. The amount of additional instructions is caused by instruction duplication, which can be caused by loop unrolling.

libtiff The results here are similar to libjpeg. One difference is, the optimization causes add instructions to be translated to lea instructions.

libpng libpng seems to be different here. The assembly file after the optimization is shorter than before, by 10 instructions. Instructions that are removed are mainly mov instructions.

zlib zlib is similar to libpng. The file gets 3 lines shorter, which is not much. The rest of the binary differences are basically instruction order changes.

The data gathered here can be found at (???). Summarizing the results from above. The optimization does not improve image decoding code.

4.3. Optimization vhdI improvements

5. Conclusion

5.1. Assembler generation

The previous chapter showed that the optimizations are not that helpful on decoder code. This thesis was done for evaluating decoder code, so no others projects have been checked. However, other projects like direct UI or CLI implementations also could be evaluated, since the patterns from decoding code are quite different to ui codes. Without evaluating others, we can say, that the optimization for dropping conversion nodes from the graph is not that helpful.

5.2. vhdl generation

5.3. Further improvements

5.4. Additional analyzer usage

Bibliography

- [1] G. Lindenmaier, “libFIRM – a library for compiler optimization research implementing FIRM,” Tech. Rep. 2002-5, Sept. 2002.
- [2] G. Birkhoff, *Lattice Theory - Third edition*. American Mathematical Society, Colloquium Publications, 1995.
- [3] K. K. Keith D. Cooper, Timothy J. Harvey, “Iterative data-flow analysis, revisited.”
- [4] K. Rosengren, “Modelling and implementation of an mpeg-2 video decoder using a glas design path,” tech. rep.

Erklärung

Hiermit erkläre ich, Marcel Hollerbach, dass ich die vorliegende Bachelorarbeit selbstständig verfasst habe und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe, die wörtlich oder inhaltlich übernommenen Stellen als solche kenntlich gemacht und die Satzung des KIT zur Sicherung guter wissenschaftlicher Praxis beachtet habe.

Ort, Datum

Unterschrift

A. Sonstiges

A.1. Table of node rules