FinalThesis

# TheCountingAlgorithmforsimulationof million-gatedesigns

by

# KlasArvidsson

LITH-IDA/DS-EX--04/046--SE

2004-05-12

Linköpingsuniversitet
Institutionenfördatavetenskap

virtutech®

FinalThesis

# TheCountingAlgorithmforsimulationof million-gatedesigns

**by**

# KlasArvidsson

LiTH-IDA/DS-EX--04/046--SE

2004-05-12

Supervisor:  BengtWerner
Examiner:    ErikLarsson

## Abstract

Akeypartinthedevelopmentandverificationofdigitalsystemsis simulation.Buthardwaresimulatorsareexpensive,andsoftware simulationisnotfastenoughfordesignswithalargenumberofgates.As today'sdigitaldesignsconstantlygrowinsize(numberofgates),andthat trendshowsnosignstoend,fastersimulatorshandlingmillionsofgatesare needed.

Weinvestigatehowtocreateasoftwaregate-levelsimulatorableto simulateahighnumberofgatesfast.Thisinvolvesatrade-offbetween memoryrequirementandspeed.Acompactnetlistrepresentationcan utilizecachememoriesmoreefficientbutrequiresmoreworktointerpret, whilehighmemoryrequirementscanlimittheperformancetothespeedof mainmemory.

WehaveselectedtheCountingAlgorithmtoimplementtheexperimental simulatorMICA.Themainreasonsforthischoiceisthecompactwayin whichgatescanbestored,butstillbeevaluatedinasimpleandstandard way.

Thereportdescribestheissuesandsolutionsencounteredandevaluatethe resultingsimulator.MICAsimulatesaSPARCarchitectureprocessor calledLeon.Largernetlistsareachievedbysimulatingseveralinstancesof thisprocessor.Simulationof128instancesisdoneataspeedof9million gatespersecondusingonly3.5MBmemory.InMICAthisdesign correspondto2.5milliongates.

## Content

# 1 Introduction

## 1.1 Motivation

Today's digital systems are growing larger and larger in terms of complexity and the number of transistors that can fit on one chip. With the concept of system-on-chip entire systems can fit on a single chip. As entire platforms of processor with cache, ASIP, ASIC, memories etc. now fit on one chip, simulation grows important as crucial part in several important design steps. Many subcomponents are found as IP-packages, and early simulation of the system can help detecting problems and making better design decisions. Effects of exchanging a part of a system with a new or updated component can be studied in a simulator, without building a prototype.

Design verification and software development is made easier by letting the engineer see what actually happens inside the chip. Simulation also speed up system development time by enabling co-design. Software development can be started already with an early description of the system, for example an instruction set, which can be simulated before anything is actually built.

This leads to a need for fast simulators capable of handling a high number of gates. One might argue that there are lots of simulators already, capable of simulating almost any abstraction level with high observability, such as ModelSim for VHDL simulation, and we have SPICE if we really want details. But are they fast, and do we need all that observability for all purposes? The answer to both questions is no. With that kind of observability they are not fast, and no, in most cases we do not need to see all glitches and timing issues, timing analysis are better and faster for this [Jennings91].

But software or hardware LCC (Levelized Compiled Code) simulators then? Yes, but software LCC simulation does not handle very large designs, or is slow for large designs, and hardware simulators are typically very expensive.

In the past decade there has been a tremendous increase in computer performance and complexity. Approximately ten years ago, a workstation would be a 66MHz 486 with 8MB memory. Today, the equal amount of money would give several GHz and almost hundred times the memory, a memory equal to old times hard drive sizes. Still, not much seems to have happened when it comes to software simulator development. In fact most work found in the area is from the late 1980 or early 1990 decade.

## 1.2 Purpose

### 1.2.1 Goal

As briefly mentioned above simulators are practical tools during software and hardware development. This work is intended to investigate how to construct a simulator that can simulate large netlists, and how fast it can be made on a regular workstation. We will implement and evaluate MICA, a Multi-queue Interpretive Count Algorithm simulator. In a previous work a simulator named EMIL was developed. That simulator serves as a reference to verify and compare the result of this work against, and some parts of EMIL will be reused in this work. To get a clearer picture of EMIL, I would like to refer the reader to [EMIL2002]. Only a brief description is given in section 2.1.

### 1.2.2 Limitations

The original plan was to implement MICHEL (Multi-queue Interpretive and Compiled Hierarchical Event-driven Levelized simulator), where selected parts of the netlist could be compiled for higher speed, while other parts was to use a compact interpreted representation. This plan also involved optimizing compiled code for the x86 architecture. But lack of time canceled the compiled part, and instead we got MICA, still assuming the x86 architecture.

## 1.3 Problem Definition

When building a simulator for small to average sized designs we are faced with several problems, such as [1]:

- o Which abstraction level should be simulated?

- o What signal and timing models should be used?

- o Is event-driven or oblivious simulation the best choice?

The answers depend on what we want from our simulator; how much details of the simulated circuit we want to monitor. Secondly, they depend on how long we are prepared to wait; how fast simulation we need. This is a tradeoff between observability and performance. Most often, what we would like is all information in no time. Thus we have a fourth question:

- o How do we make it faster?

---

[1] The reader not familiar with the different concepts such as event-driven, oblivious, zero-delay etc. that is mentioned in this section is referred to section 2.2 for further explanations.

4

EMIL tries to answer these questions, and as a continuation of EMIL many answers remain the same in this work. We will, like EMIL, use event-driven gate-level simulation with zero-delay timing and only two signal levels. For many purposes, such as evaluating a component together with a full system, or as part of design verification, the observability given by this suffices. Instead, the fourth question, simulation speed is what really matters.

But with a large design, this grows more complex. First we have the obvious expectation that more gates will lead to longer execution time. If $10^3$ gates took $10^3$ time units we would naturally expect $10^6$ gates to take $10^6$ time units. Secondly, more gates will naturally need more memory. Modern computers have huge amounts of memory, so by size this is not a problem, but by performance it is. The larger the memory the slower it is, and this breaks that first and obvious expectation. Having a data structure too large to fit in L1 cache will lead to slower execution, and having an algorithm with bad locality for memory accesses will make it even worse. Thus, the fourth question can be reformulated:

- How to simulate large designs fast?

And this is the question this work will investigate. Let us take EMIL as an example and try to picture it with a large design.

## 1.4 ImagineEMILwithanXmilliongatedesign

Letusfirstexaminesomedatafrom[EMIL2002]in Table1andTable2. TheEMILbenchmarkswasrunona500MHzPentiumII Iprocessor [1].

|  | Total | Eachcycle | Eachsecond |
|---|---|---|---|
| Executiontime | 30s | 0.6ms | - |
| Simulatedcycles | 50000 | - | 1667 |
| Gateevaluations | 64M | 1278 | 2M |
| Changedsignals | 22M | 445 | 739K |
| Realinstructions | 5.3G | 106K | 176M |
| Memoryused | 348KB | | |
| Gateactivity | 14.8% | | |
| Totalgates | 8637 | | |

**Table1EMILbenchmarkstatistics**

|  | Size | Latency(ns) | Latency(cycles) |
|---|---|---|---|
| Level1data | 64KB | 6 | 3.0 |
| Level2cache | 512KB | 45 | 22.5 |
| Memory | 256MB | 147 | 73.5 |

**Table2Cachestatistics**

EMILcachemisseswasreportedtobelow,almostno L2missesandonly about5%L1misses,synchronousgatesexcluded(all searchedallcycles). Thissuggeststhatthesimulationnetliststructure fitintheL2cache,and themostfrequentlyusedgatesfitintheL1cache. Thefirstsuggestionis clearlytrue,sincethememoryusedforthesimulat ionnetlistisonly384 KBandtheL2cacheis512KB.Makingthesimplifie dassumptionthatit isthesame15%(thegateactivity)ofthegatesth atevaluateseachcyclewe calculatethecriticalnetlistsizetoabout58KB $(384 \times 15\%)$.Withasmall marginthisfitsnicelyintheL1datacache,soou rassumptionseemstrueto atleastsomeextent.Ifitweremainlydifferent1 5%ofthegatesthatwas evaluatedfromcycletocycle,moreL1misseswould occur.

---

[1]ThisinformationwasnotmentionedintheEMILreport,butisfromtheauthorof EMILandconfirmedbythecachedatainTable2;3cyclesdividedby6nsequal 500MHz.

Now,letusimagineEMILwithan X milliongatesdesign.Thismeana
simulationnetlistsizeof $XM \times 384K/8637 \approx 40X$ millionbytes.Atthe *same*
speedandgateactivitythiswouldtake $XM \times 50000 \times 15\%/2M \approx X$ hoursper
milliongates.But 40X MBwillnotfitinanycache,thecriticalsizeal one
wouldbeabout 6X MB,byfarlargerthananycachelevel,andsince we
nowwillneed73.5cyclesforeverymemoryaccessi nsteadof3cycleswe
willexecuteupto24.5timesslower.Thus,wecan expectmanyhoursof
simulation.

## 1.5  Gate-levelsimulation

Digitalsystemscanbedescribedinseveraldiffere ntdomainsondifferent
abstractionlevels.Onepopularwayofshowingthem areGajski'sY-chart
inFigure1.



**Figure1Gajski'sY-chart**

Therearethreedomains,behavioral,structuraland physical,eachwithan
abstractionhierarchycontainingsystem,registert ransfer,logicandcircuit
level.Gate-levelsimulationtakesastructuraldes criptionofasystematthe
logiclevelandreturnsthebehavioralresponseto anychoseninput
sequence.

Themostimmediateviewofagateisaunitimpleme ntingoneofthebasic
booleanfunctions(AND,OR,NOT),butwithincrease dabstractionto
handlecomplexitytheconceptofagatehasbecome somewhatblurred.For

7

example, a 32-bit adder can also be viewed as a gate, or super-gate, although it consists of 32 full adders, in turn consisting of several basic gates. So let us say that gate-level means a circuit description built with logic complexity from basic boolean functions up to units implementing more complex but still conceptually simple functions, such as flip-flops, multiplexers and full adders. Now we know what a gate-level simulator should take as input, but what should it be able to do?

We can try to define a gate-level simulator as a simulator *that takes a circuit design at gate-level and evaluates all gates in the circuit so that correct output is produced*.

However, this is not a good definition, since the important thing with a simulator is to be able to correctly view the internal state of a circuit at any point in simulation. Thus, a better definition would be:

> *A simulator that enable the user to correctly, at gate-level and any point in simulated time, view any part of the circuit's internal state.*

This is a better definition, since it defines only how we would like to use the simulator, and not how it should work internally. With this definition the simulator has the opportunity to cheat, and only evaluate gates that change, or only simulate the monitored parts of the circuit at gate-level, the rest of the circuit could be simulated faster on a higher level.

## 1.6 Computer architecture

Simulation speed depends on the implementation, the algorithms used and the computer architecture. If we want maximum speed they have to be coordinated. We have to use algorithms and coding that is adapted to the target architecture. This mean the properties of the target computer architecture must be known to make an efficient implementation. A brief overview of a typical x86 superscalar architecture is given in Figure 2. We assume the execution units are pipelined. To keep the instruction pool full instructions are prefetched from memory and decoded. When a control transfer (jump, call, ret) instruction appears the prefetcher tries to determine the new address at where to fetch instructions by prediction, often based on branch history statistics kept in a branch prediction table. If the prediction is wrong, all instructions following the jump instruction must be canceled and the correct instructions fetched. This typically costs nearly as many cycles as there are pipeline steps.

Thus,firstconclusionisthatafastimplementatio nshouldavoidjump instructionsthatarehardtopredictcorrectly.



**Figure2Atypicalcomputerarchitecture**

Now,ifeitherinstructionordataisnotintheL1 cachetheexecutionmust waitupto73cycles,accordingtoTable2insecti on1.4.Thecache contentsareupdatedbasedonusagehistory.So,se condly,afast implementationeitherfitcompletelyinL1cache,o rhavehighlocality. Andlast,inordertouseseveraloftheavailable executionunitsineach cycleweneedlowinterdependencebetweeninstructi ons.

## 1.7 Leon

LeonisafreelyavailableimplementationofaSPAR Cv.8processorfrom GaislerResearch[Gaisler].SPARCisaRISCprocess orarchitecture formulatedbySunMicrosystemsin1985:"SPARCwas designedasa targetforoptimizingcompilersandeasilypipeline dhardware implementations.SPARCimplementationsprovideexce ptionallyhigh executionratesandshorttime-to-marketdevelopmen tschedules." [SPARC]

In[EMIL2002],Leonwaschosenastestinput"beca useitisagood exampleofhardwarethatonecouldwanttosimulate atgate-levelin cooperationwithafullsystemsimulator".Inthis workwesticktothis choicealsoinordertocomparetheresultsagainst EMIL.Weuseexactly thesameLeon(intheformofanEDIFnetlist[EDIF ])totestand benchmarkthesimulator.Forfurtherinformationon howLeonwas configuredandsynthesizedthereaderisreferredt o[EMIL2002].

## 2  RelatedWork

### 2.1  EMIL

EMILisanEvent-drivenMulti-queueInterpretiveLevelizedgate-level simulator[EMIL2002].ThesimulatorengineinEMILconsistsofa scheduler,aqueuesystemandadispatcher.APERLscriptcreatesa specificevaluationfunctionforeachgatetypefromthedescriptioninthe logiclibrary(asynch_OR3inFigure3).Thenetlistislevelized,andthe levelisstoredwitheachgatetogetherwithpointerstofan-in,fan-outand evaluationfunctioninagatestructure,seeFigure3.Thesizeofthis structurevariesdependingonfan-inandfan-out.Duringsimulationthe dispatcherpickgatesfromthequeuesystemandcallsitsevaluation function(withthegatestructureasargument).Theevaluationfunctionin turncallstheschedulerifthegateoutputchanged.



**Figure3EMILgatestructureforanOR3gate**

Thisgivesfastgateevaluationsinceeachgatehasaspecificcompiled function.ItalsogivesasmallcodesizecomparedtoatraditionalLCC simulatorsinceonlyoneinstanceofeachgatefunctioniscreated,andthe numberofgatetypesinthelogiclibrarylimitstheenumberoffunctions. However,itcanbearguedifEMILreallyisinterpretive,sinceitactually generateC-codeforeachgatetypeandcompileittogetherwiththe simulatorengine.

InEMIL,theproblemsoccurringwithsynchronousgateswhenitcomes toevaluationorderandasynchronousinputs(seesection2.2.3,Gate orderingbelow)aresolvedbyaspecialorderedsynchronousqueue,anda specialasynchronouspartscheduledandevaluatedwiththeasynchronous gates.

10

## 2.2 Simulator concepts

Several simulator concepts have emerged as an effect of previous work in the area, mostly aiming at higher simulation speed but also for more precise simulation (exacter models) or reduced preprocessing time. These can be divided in Time and signal modeling, Algorithm choices, Gate ordering, Implementation strategies, and Netlist representation.

### 2.2.1 Time and signal modeling

In reality, a digital signal is nothing but an analog voltage pulse. It needs time to rise, to fall, and to propagate from one place to another. How this is described in a simulator is first question of what resolution one need, which abstraction level to model. At gate level the voltage level is typically described by a discrete set of values, in VHDL {U, X,0,1,Z,W,L,H,-}. Modeling of more signal levels or more timing accuracy requires more calculations, and is as a result slower.

For timing there are three methods, nonunit-delay, unit-delay, or zero-delay. Nonunit-delay means that it takes a certain number of time units for a signal to propagate through each gate. How long is determined for each gate. In unit delay this time is fixed to one time unit for all gates, and in zero-delay no delay is used. The delayed models allow us to see glitches and timing issues of signals, but we need to simulate with a higher time resolution than clock cycles; we must simulate each time step any input changed. With zero-delay we can cancel all glitches and simulate a gate only when its inputs are stable for the cycle.

### 2.2.2 Algorithm choices

Traditionally two types of simulation algorithms exist, *event-driven* and *oblivious*. In event-driven simulation each changed signal generates an event, triggering evaluation of affected gates whose output signal possibly changes, in turn generating new events. Thus only gates involved in signal changes are ever evaluated. Oblivious simulation on the other hand evaluates all gates each cycle. It does not care nor remember (explaining the name oblivious) that nothing changed for a particular gate since last evaluation, but simply evaluate it again. The advantage with this method over event-driven simulation is that no event system is needed, no code determining new gates to schedule after each gate evaluation, and no code managing the scheduled gates. Each gate evaluation can be made small and fast, but we have to evaluate more of those.

The normal way to describe a gate is by giving the boolean function, where to find each input, and where to save the output. But this can be modeled differently. Maurer uses a method of describing gates originally discovered by D. Schuler to create the Inversion Algorithm [Maurer97].

Themethod,namedcountingalgorithm,usesthefact thatmostnormal gatescanbedescribedbyacountandadominantva lue(D.Schulerfrom [Maurer94]).Whennoinputhavethedominantvalue ,theoutputis0(or 1)elseitis1(or0).Forexample,whennoinput ofanANDis0 (dominant),theoutputis1,else0.Orwhennoinp utofaNORis1 (dominant),theoutputis1,else0.Sothealgorit hmsimplycountsthe numberofdominantinputs.Thecountingalgorithmt henevaluatesgates eachtimeaninputchanges:

```
if changed input is dominant then
   increment gate count
   if gate count is 1 then
     output = not output
   end if
else
   decrement gate count
   if gate count is 0 then        (no input dominant)
     output = not output
   end if
end if
```

Thismethodhassomeadvantages,andsomedrawbacks :

+ Allinputsarerepresentedbyonecount.

+ Asinglesimpleevaluationfunction.

– Requirescorrectinitializationofgateoutputand count.

– Doesnotsupportalltypesofgates,onlynormalsi mplegateswhere allinputshavethesamemeaningtothegate.

Anotherwaytogoisbybranchingprograms.Basedo ntheboolean functionaBDD(BinaryDecisionDiagram)iscreated foreachoutputofa gatenetwork,andabranchingprogramisderivedfr omtheBDD.The methodhastheadvantageofhavingaworst-caseeva luationcomplexity proportionaltothenumberofinputsandoutputs,b uttheBDDsizegrows exponentiallyintheworstcase(althoughovercome inthegeneralcase) [Ashar95].Otherdisadvantagesarethatonlyinput andoutputvaluesare known;intermediatenetshavenoobservability,sin cenogatesare simulated,andlocalityispoor.[Jiang2003]prese ntsageneralizedmethod ofthisintheformofaGeneralizedCofactoringDi agram(GCD).

### 2.2.3 Gateordering

ConsiderthegatesinFigure4.Clearlywehaveto evaluatethegatesin correctorder,inthiscaseA,B,C,Disonesuch order.IfweevaluateA,B, D,Cwemightgetwrongresult,orwillhavetoeva luateCagain.The commonwayforasimulatortodetermineacorrecto rderisbylevelization.

12

Eachgateisplacedinaleveldependingonitsdep thinthecircuit.Allgates
inonelevelaremutuallyindependent,andcanbee valuatedinanyorder.
Levelizationisdonebyfirstassigninglevel0to gateswhoseinputsare
known(connectedtocircuitin-portsorregisters), inthiscaseAandB
receiveslevel0.Thengateswhosefan-ingatesall havealevelaregiven
thenexthigherlevel[SSIM87].ThusDcannotbeg ivenalevelsinceits
fan-ingateCdoesn'thavealevelyet,butCisgi venlevel1.LastDcanbe
givenlevel2.



**Figure4Levelizationexample**

Aproblemwithlevelizationoccurswhenacircuith asfeedbackpaths
(loopsorsometimesalsocalled *stronglyconnected* gates).Afeedbackpath
ischaracterizedbythefactthattheoutputofeac hgateinthepathdepends
ontheoutputofallothergatesinthepath.Consi deriftheoutputfromD
wasconnectedtooneoftheinputsofB,creatinga nasynchronousloopB,
C,D,B.Bcannotbegivenalevelsinceitsfan-in gateDhavenoleveland
thusCandDcannotbegivenaleveleither.

Inthework[LECSIM90],suchpathsarehandledby detectingandgroup
thegatesinasinglegate.Ingate-levelmostasyn chronousloopsarealready
grouped,formingdifferentkindsofflip-flops.Loo psincludinga
synchronousgate(synchronousloops)arecommon,bu tnotreallyaloop
sincethesynchronousgatebreakstheloopuntilne xtcycle.

Asecondproblemwhenitcomestogateorderingis synchronousgates.
Sinceallsynchronousgatesevaluatesimultaneously ,theymustbe
evaluatedincorrectorder.Startingevaluationwit hthefirstoftwo
consecutivesynchronousgateswillproduceanewou tputfortheinputto
thesecond,butitshouldusetheoldvalue;inrea litythenewvaluehaveno
timetopropagate,sincebothgatesreceivesthecl ocksignal
simultaneously.Startingwiththesecond(lastoft heconsecutive)gate,or
insertingabuffergatebetween,solvesthis[EMIL 2002].Synchronous

gatescanalsohaveanasynchronousinputthathas tobelevelizedand evaluatedcorrectly.

Anotherorderingpossibility,atleastforobliviou ssimulation,wouldbeto searchthenetlistdepth-firstfromeachout-port, somethinglike:

```
mark all design inport nets as storage points
for all design out-ports
   recurse driving gate
end for
```

Gaterecursionroutine:

```
for each gate fan-in
   if net is marked as storage point then
      output load operation
   else if net have > 1 fan-out then
      recurse driving gate
      output gate evaluation operation
      output store operation
      mark net as storage point
   else
      recurse driving gate
      output gate operation
   end if
end for
```

Theoutputfromthealgorithmwouldbeaninstructi onsequencewhere eachgatecanreaditsinputfromthestack,andsa veoutputtostack, hopefullyresultinginhighlocality.Detectingfee dbackloopswouldbe straightforward,handlingthemsomemoretrouble.

## 2.2.4 Implementationstrategies

Talkingaboutaprogramminglanguage,suchasC,th ereareaclear differencebetweenthetwoconceptsofinterpreted versuscompiled.An interpretedlanguagesuchasLISPorPERLisexecut eddirectlybya program(interpreter),whileacompiledlanguagear eputthrougha compilerandexecutedbyhardware.Thosetwoconcep tsoccuralsointhe areaofsimulators,anetlistcaneitherbereadby thesimulatorand interpreted,orcompiledwithasimulationenginet oexecutedirectly.A compiledsimulatortriestoremoveasmuchaspossi bleoftheruntime translationfromnetlisttoexecutedinstructions.

Usingoptimizationtechniquessuchasloopunrollin g,directaddressing insteadofindirect,andthreadedcode(arrangecod esegmentsinexecution ordertoremovejumps,calls,andreturns)canmake acompiledsimulator muchfaster.Wecouldalsoimagineanevaluationro utineevaluatingtwo (ormore)independentgatessimultaneously,toincr easeinstruction

14

parallelism,andtherebyhopefullyuseasuperscala rarchitecturemore effectively.

Attemptsatcondition-freesimulationhavealsobee nmade.TheInversion Algorithm[Maurer94]doesthisbytogglingthegat eprocessingroutine betweentwoversionsaftereachcall,andlaterwor ks[Maurer2000] similarlybytakingtheaddressoflabels.Theimme diatethoughtofthisis nice,afunctioncallisafixedjump,andtheproc essorcancorrectlypredict fixjumps. *Butitdoesnotwork* .Whatwegetisaregisterindirectjump,the processingroutineaddressisloadedtoaregister andtheprocessoristold tojumptotheaddressinthatregister.Andthisi saconditionaljump,it dependontheaddressintheregister.Further,sin cetheaddressistoggled aftereachcall,thebranchpredictiontablewilla lwayscontainthewrong address,andwegetacostlymiss-prediction *each*call.

Now,doesitworktofixthisbytakingtheaddress ofalabelandthenuse gotowiththatlabeladdress?Again,thethoughtis nice,butwehavethe sameproblem,andfurther,itisnotsupportedinC (GNUextensionstogcc supportsithowever).Onlyiftheactualgotoinstr uctionbitsarereplaced withotherit *might*work,justsomeissuestoomany;itrequiresassem bly,it isnotportable,itisunreadable,doestheprocess orreallypredictfixjump frominstructiononly?Andareweallowedtomodify thecodesectionby hardwareandOS?

Otherwaysofoptimizingbycodewouldbetopackm anysignalstogether anddoonelogicoperationonallsimultaneously.H owever,theoverheadin packing,unpackingandrepackingwouldbesubstanti alatgate-level.At registertransferlevelitwouldbemorenatural.

Last,implementinganevent-drivenlevelizedsimula torcanuseoneor severalqueuestostorescheduledgates.Onequeue isfastindetermining thequeue,butneedtoorderscheduledgatesbylev eltoavoid reevaluations.Severalqueues(multi-queue)needmo reoperationsin determiningthequeueanddeterminewhenallqueues areempty,butneed noorderamonggatesinthequeues.Ifwedistingui shfull-queue(asmany queuesaslevels),multi-queue(severalqueues,but somelevelsmightshare queue)andsingle-queue(onlyonequeue),MICAwoul dbeafull-queue simulator.

### 2.2.5 Netlistrepresentation

Anetlistisoftendescribedhierarchically;acomp onentisdescribedin detailinoneplaceandthenusedinseveral.Aphy sicalimplementationof thenetlistmustflattenthenetlist,bycopyingth edescriptionofeach componenttoeachplacethecomponentisused.

Asimulatorhowever,canchoosetoexploitthehier archytosavespaceby creatingonedescriptionofhowtosimulatethecom ponentandusethis withdifferentdataeachtimethecomponentoccurs inthenetlist,muchlike aprogrammercancreateonefunctionforacommont askandthenuseitin severalplaceswithdifferentdata.Comparedtoaf latversionthiscreates someoverheadincallingthehierarchicalcomponent s,butwesavespace. [Lewis91]

[Maurer99]distinguishahierarchicalcomponentus edonlyonceasa partition.Suchpartitionscanbeusedtoreducega teevaluations,andalso becreatedinordertoreducetheschedulingoverhe adinanevent-driven simulator.Thelatterisdoneby[Blaauw93]witha clusteringalgorithm. [DeVane97]observesthatgatesthatprecedeandfo llowsynchronousgates needonlybeevaluatedwiththese,and[Maurer99] thatexistingpartitions canbeswitchedoffsomecycles,forexampleaCPU withbothALUand FPUneedonlysimulatetheunitactuallyusedeach cycle.(Inmyopinion thecircuitdesignershouldthinkofdisablingsuch unitswhennotused, sincethiswillalsosavepower,whichisanincrea singissueduetotheheat generatedandbatterylifetime,butthisisasidet rack.)Partitioningisalso usedtoovercomelargeBDDsizesinthecreationof branchingprograms [Ashar95].

### 2.2.6 Combinations

Thedifferentconceptscanbecombinedinseveralw ays.Some combinationsarestraightforward.Historicallythe compiledandoblivious conceptsdidgohandinhand,asdidtheinterpreti veandevent-driven.The LCC(LevelizedCompiledCode)istraditionallycons ideredthefastesttype ofsoftwaresimulator,suchas[SSIM87].Asthene edforfastersimulators increasednewcombinationswheresought.Thefastc ompiledcode concept,andtheevent-drivenideareducingthenum berofgateevaluations requiredwascombinedinseveralsimulators[SLS88 ],[COSMOS87]. [LECSIM90]alsoaddedthelevelizationtechniquet othiscombinationto minimizegatereevaluations,while[Lewis91]hase xploredtheeffectsof utilizingnetlisthierarchyaswellastheeffects ofcaches.

16

# 3 Approach

## 3.1 Solutionstrategies

As we have seen there are a lot of different approaches to simulation, and how to make it fast. To return to the question in section 1.3, how do we simulate a high number of gates fast? Obviously we want to evaluate as few gates as possible with as few instructions as possible, using a minimal datastructure to store the netlist. We have some options:

- o Use alternate evaluation methods to reduce or simplify gate evaluations.

- o Store only the information absolutely needed in simulation, as tight as possible.

- o Exploit hierarchy; use the same subcomponent description in several places, but with different data.

- o Compile the most frequently used gates to increase speed for those.

- o Use algorithms that increase data locality.

Some of these options contradict, i.e. storing only absolutely needed data in a tight way will most likely mean some extra instructions packing and unpacking bits, and compiled gates typically takes more space. However, looking back at Table 2 (section 1.4) we note that we can spend up to 19 instructions avoiding each L1 cache miss, and up to 70 instructions avoiding L2 misses, meaning some extra work involved in reducing data size might be worthwhile.

This can be formally expressed. Assume we evaluate a total of $x$ gates. Each evaluation need $g$ cycles, and $w$ cycles are wasted each time a gate is not in cache. Let $m$ be the miss rate, the total percentage of gates not in cache. The total cycles needed for evaluation is now $y = xg + xmw$. Suppose we can achieve the new miss rate $m_{new} = m(1-d)$ where $d$ is the percent miss decrease, by adding $i$ cycles per evaluation. Then we get $y_{new} = x(g+i) + xm_{new}w$. Of course we want:

$$y_{new} \leq y$$

$$x(g+i) + xm_{new}w \leq xg + xmw$$

$$i + m_{new}w \leq mw$$

$$i \leq mw - m(1-d)w$$

$$i \leq mwd$$

Thatis,adding2instructionspergatehaving20% missrateandwasting
20cyclespermissweneedtodecreasethemissesb yatleast50%,which
wouldbehard.Butsaywecouldmovefrommemoryto theL2cache(50
cyclessave)by30%memorysizedecrease.Thismean sthatadding3
instructionspergatewillbetoourbenefitifthe originalmissrateismore
than20%.Andifwecandoitwithoutaddinginstru ctionswecan'save'3
instructionspergate.

## 3.2 Netlistrepresentation

Thisworkfocusesonexploringthethreetopoption sinsection3.1by
usingideasfromthecountingalgorithmdescribede arlier.Maurerusesthis
methodfocusingonreducinguselesssimulations,un propagatedchanges,
andforunconditionalsimulation[Maurer97].Howev erwearenot
interestedinthis,themethodofreducingunnecess aryevaluationsin
practicemeanweevaluatethegateforeachinputc hange,andthe
unconditionalmethodusedactuallymeanwegetlot ofregisterindirect
branches.

Thefeatureofthecountingalgorithmwewouldlike touseisthefactthat
agatedoesnothavetoknowwhateachandeveryin putvalueis,orwhere
tofindthem.Onlytheamountofhighinputsandga tetypeisenoughto
determinetheoutput.Thismeanwegetasimpleeva luationfunction,anda
compactnetlistdescription.Asweshallselater, over70%oftheswitching
gatesarerepresentedbyonlyonestandardtype.Th isisachievedby
combiningthecountideawiththetraditionalevalu ationmethodusedin
EMIL.Thisalsoremovestheneedforproperinitial izationofcountvalues
andoutputs.

Toenableuseofhierarchyandreducetheamountof informationstoredin
thesimulatednetlist,wedistinguishthreekindso fdataandseparatethem
indifferentareas, *management*, *structure*and *data*.Themanagementarea
storesinformationusefultofindandaccessindivi dualgatesandnets,but
notneededforsimulation,forexampleinstancenam es.Thestructurearea
storesdatathatisreadonlyduringsimulation,e. g.howgatesandnetsare
connected,butnotvaluesthatarechangeddurings imulation.Duetothe
natureofthecountwayofstoringinputsthisisa one-waystructure,gates
knowthelocationoftheirfan-out,butnotofthei rfan-in.Finally,thedata
areastorealldynamicdatachangedduringsimulati on.

## 3.3  Mainproblems

Itishardlyasurprisethatthissolutioncreates someproblems.Thereare mainlytwoofthem.

1. Asmentionedinsection2.2.2,Algorithmchoices ,onlysimplegates withindependentinputssuitthecountingalgorithm .Othergates havetobeidentifiedandhandleddifferently.

2. Separatinggatesinstructureanddataareamean weneedtwo differentpointerstoeachgate,ormusthaveequal sizedgatessothe samepointercanaccessbothareas.Butgateshave differentoutput netsizes,fromjustoneorafewfan-outsuptoth esizesoftheclock orresetnetwork.

### 3.3.1  Problem1,incompatiblegates

Wecancreateaspecialorhierarchicalversionof thegate,eithercompiled insamemannerasEMILdoes,ordescribedwithsimp lercountcompatible gatesandusingourstandardcountevaluation.But thisinvolvesextrawork storingandpassinginputstoandfromthespecial orhierarchicalgate,and morechoiceswhendecidingwhatevaluationfunction touse.

Ifsuchgatesinsteadaredecomposedtoseveralsim plergatessupporting thecountingalgorithm,anddirectlyreplacedinth enetlistwemightstill benefit,despitemoregates.Sincewegetseverals implergates,insteadof onecomplex,weexpectonlysomepartsofthecompl exgatetobe evaluatedwhenaninputchange,thusreducinggate activity.



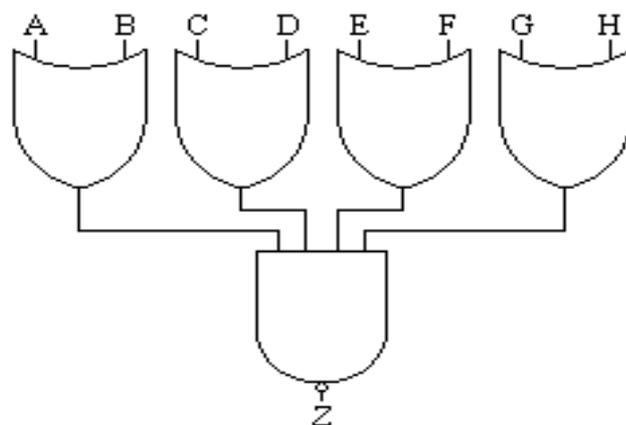**Figure5TheAO12Pgatefromthelogiclibraryused**
**{ Z=((A+B)(C+D)(E+F)(G+H))'}**

TheAO12Pgateforexample,viewedwithsimplergat esinFigure5,has eightinputsandalargebooleanfunction.Withas pecialorhierarchical versionallinputsmustbefoundandtheentirefun ctionevaluatedassoon asoneinputchange.Butifwedecomposethisgate tofivecount compatiblegates(oneforeachORandoneforthef inalNANDasthe

figure show), we only need to evaluate one OR and possibly the NAND as a result of one input changing. So decompose the gate seems more promising.

Still, this solution result in some really bad cases, and for some synchronous gates this is still not enough. In short, synchronous gates are represented with one standard synchronous gate and a front function to get special behavior, but this is described more in section 4.3.2, Synchronous gates.

Here we will instead mention one of the gates really bad to decompose: the two-to-one multiplexer. There are four special properties of this gate making it a bad case. First, it is the third most common gate in Leon, and because it is part of both data and control path it is likely to change often. Second, it splits to as many as four gates. Third, a change on the select signal activates all of these splits. And fourth, the function is really simple, a special routine could evaluate the gate with one if-statement. Although there was not time to implement and investigate the effect of a special type this is really needed, as we will see later.

### 3.3.2 Problem 2, nets are not equal sized

Using two pointers for one gate is both space consuming and inconvenient, using fixed size gates is much more appealing. The problem is then to represent the fan-out pointers. One way to go would be to have one list with the first fan-out of each gate, one with the second etc, and sort the gates by number of fan-outs. The fan-outs can then be accessed in turn using the same offset into the fan-out lists as the gate offset in the structure and data area. If it points outside the list no more fan-outs exist for that gate. This would not waste any memory, but would be bad from a cache point of view; fan-outs would be scattered in memory. It also imposes a fix order, but we might want to order gates to reduce the distance in memory to its fan-outs (enabling smaller fan-out pointers), or orders of frequently used gates are clustered in memory (to increase locality).

Each gate can also store a pointer to a fan-out list. But according to statistics most gates have only one fan-out, and storing this directly would be more efficient. This lead us to the solution of storing a fix number of fan-outs in each gate, and use a special gate type for gates with many fan-outs, storing a pointer to a fan-out list.

The separation of structure and data in fixed size gates, as well as using the count evaluation method, also mean a restriction of only one output per gate. To represent gates with multiple outputs we have identical choices as in Problem 1, and we solve it by creating one gate per original output.

20

## 4 Implementation

## 4.1 Overview

Figure6givesanoverviewofhowMICAworks.Theo valsdescribe
actionstakenbyMICA,andtherectangularboxesde pictdatausedasinput
toandoutputfromthoseactions.Thearrowsindica tewhichdataisusedas
inputandwhatwegetasresult.Thedarkeractions (ovals)indicateparts
thatwererereusedfromtheEMILimplementation.



**Figure6Implementationoverview**

Thefollowingsectionswilldescribemostofthese actionsandsomeofthe
dataformatsusedinmoredetail.Theinputtothe simulatorisanEDIF
netlistandaLogicLibrarydescribingthegatesus edbythenetlist.Thegate
functionsareextractedfromtheLogicLibraryand classifiedintocount
types.Thisclassificationisaddedtotheparsedn etlist,possiblysplitting
somegates,yieldingacountcompatiblenetlist.Qu eueandpriorityisthen
addedforeachgatebylevelization.Finallythene tlistisoptimizedtothe
MICAinternalsimulationformat,andreadytosimul ate.

## 4.2 FunctionExtraction

TheLogicLibrarycontainsmoreinformationthanne ededforourpurpose.
InEMILaPERLscriptwasusedtoextracttheboole anfunctionofeach

21

gateandcreateitsC-codeevaluationfunction.Thisscriptwasstrippedand slightlymodifiedtoinsteadcreateatextfiledescribingonlytheimportant propertiesofeachgate.Thesimplegrammardescribingeachgateisgiven inAppendixB,Grammar.

Synchronousgatescannotyetbeextractedcorrectly,andhavetobe manuallydescribedinthetextfile.Othergatescaneasilybehandtunedif desired,forexampletodecreasethenumberofgatescreatedwhen splitting,butthiswasnotdone,sincetheautomaticsplittingproducebest possibleresultinmostcases(exceptionislargemultiplexers).

## 4.3 CountClassification

### 4.3.1 Asynchronousgates

Thecountingalgorithmusesthefactthatmostgateshaveeitherhighor lowoutputiff(ifandonlyif)exactlynoneoralloftheinputsarehigh. ANDisforexamplehighiffexactlyallinputsarehigh,NORiffexactly noneoftheinputsarehigh.Weneedtodeterminehoweachgatetypecan bedescribedinthisway.

Assumetheinputsarerepresentedbyacountvaluetellinghowmany inputsarehigh.Whenaninputrises(lowtohigh)weincrementthiscount, andwhenaninputchangesfromhightolowwedecrement.Dependingon gateweletheexactsituationwhenoutputisknownberepresentedwith thecountvaluezero,andcalltheoutputinthissituationOaZ(Outputat Zero),i.e.forathreeinputANDwelletthecountvaluereachzeroiffall inputsarehigh,andsetOaZtohigh.Thismeanthecountvalueoffour ANDmuststartat–3whennoinputishigh.AtwoinputNORwillalso haveOaZhigh,butstartwithcountequalto0fornoinputhigh(andreach count2whenallinputsarehigh).Wecallthegatesthatcanbedescribed likethis *AtZero*gates(theyswitchatzerocount).

Similarlywecanderive *AtEven*and *AtPos*gates(outputisOaEwhenthe countvalueiseven,orOaPwhencountispositive).Thosetwotypes correspond,respectively,tothesumandcarryfunctionofafulladdergate. Inall,wecandistinguishthreepropertiesdescribingagate,the *count type*, the *output at*value,andthecount *start*.

ClassificationofgatestooneoftheAtZero,AtPos orAtEventypesisof coursedoneautomatically.Thebooleanfunctionofeachgateoutputisread fromthetext-filethatwasextractedfromthelogiclibrary,andevaluated foreachpossibleinputcombination(sincebasicgatesaresmallthisisnota problem,butitcouldbecomeforgateswithmanyinputs).

Theresultofeachinputcombinationisstoredinanarrayintheslot determinedbythenumberofhighinputs.Ifdifferentoutputvalueoccupies

thesameslotweknowthatthegateinputscannotberepresentedbya count.Next,theoutputarrayisinvestigated.Ifoutputhigh(orlow)existin onlyoneoftheslots(withallotherslotstheoppositevalue)wehavean AtZerogate.Ifonlyoneswitchfromhightolow(orviceverse)isfound whenwesearchtheoutputarray,anAtPosgateisfound,andifitswitches betweeneveryvaluewehaveanAtEvengate.

Ifagatecannotberepresentedbyacountitsbooleanfunctionissplitto partsthatcan.Thisisdonebyadepthfirsttraversalofthebinary expressiontreebuiltforevaluationofthegate.Eachnodeinthetreehas thefunctionofOR,AND,orNOT,witheachleaveasoneinputvalue.The treeisthensplitateachpointwherethechildrenofORandANDnodesis notofsametypeasitsparentnode.Someoptimizationisdonetotransform (A'+B')sub-treesto(AB)'andA'B'sub-treesto(A+B)',asthiscreates oneinsteadofthreesplits.Noweachgateisdescribedbyoneormore outputfunctions,eachwitha *counttype* , *outputat* valueand *start*value.

Letusinvestigateacarryfunctionasanexample, viewedinFigure7.The functionisevaluatedforeachinputcombination.Threedifferent combinationswithonlyonehighinputexist,andthreecombinationswith twohighinputs,butwithonehighinputtheoutput ( *Co*)isalwayslow,and withtwoalwayshigh,sonoconflictoccurswhenjoiningtherowswiththe same *#Ones*count.Afterjoiningtheserowswehaveoutputhigintwo casesandlowintwocases,butonlyoneswitchpointexist.Thuswehave an *AtPos*gate.Nowwebiasthe *#Ones*tostartatzeroattheswitchpointto getthe *start*countvalue(–2),andthe *outputat* value(high).Ifweinstead imagineathree-inputANDgate(Co=CiAB)theoutputwouldofcourse behighonlywhen *#Ones*equal3,givingusan *AtZero*gate,andthe *Count* wouldstartat–3.

| A B Ci | #Ones | Co | Co=CiAB'+CiA'B+Ci'AB+CiAB |
|---|---|---|---|
| 0 0 0 | 0 | 0 | |
| 0 0 1 | 1 | 0 | |
| 0 1 0 | 1 | 0 | |
| 0 1 1 | 2 | 1 | |
| 1 0 0 | 1 | 0 | |
| 1 0 1 | 2 | 1 | |
| 1 1 0 | 2 | 1 | |
| 1 1 1 | 3 | 1 | |

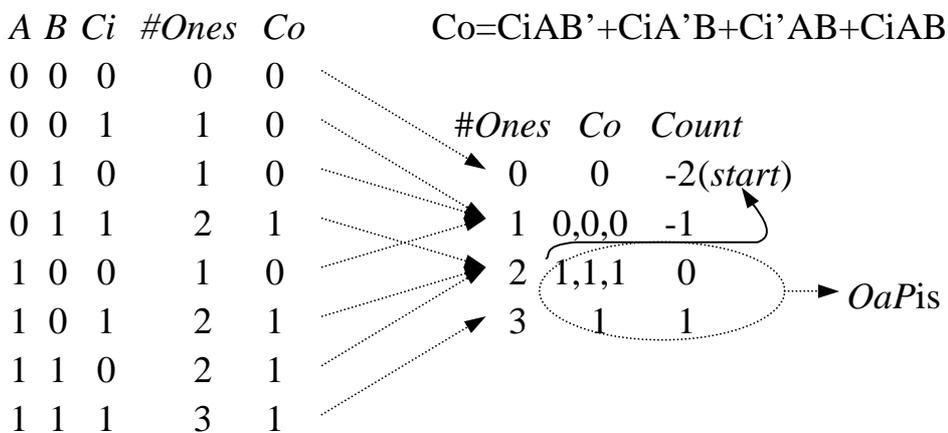| #Ones | Co | Count |
|---|---|---|
| 0 | 0 | -2(*start*) |
| 1 | 0,0,0 | -1 |
| 2 | 1,1,1 | 0 |
| 3 | 1 | 1 |

*OaP*is

**Figure7Classificationexample**

### 4.3.2 Synchronousgates

Synchronousgatesarespecialinseveralways.Firsttheyshouldonly
evaluateonrisingorfallingclockedge,secondtheyshouldevaluateallat
once(ifconnectedtomatchingclock),thirdsomehavesynchronousor
asynchronousreset,presetorbothresetandpreset,andatlasttheyoften
haveanextrainvertedoutput.Clearlynotallofthisiscompatiblewiththe
countingalgorithmschemeusedforasynchronousgates.Mostimportant,
thecountingalgorithmdoesnotknowwhichinputiswhich,andthus
cannotdistinguishbetweendata,resetorpresetports.

Withonlyadatainput(asimpleD-FF)wehavenothingbutanAtZero
gate(withOaZlow)updatedonlyateitherrisingorfallingclockedge.The
extrainvertedoutputcouldbecreatedbysimplyaddingasiblingflip-flop
usingOaZhigh.However,thatwouldalsodoubletheschedulesand
evaluationsofsynchronousgates,andstillnothandleresetand/orpreset
inputs.

Instead,weletthesynchronousgatebeadouble-gate,twoconsecutive
gateswherethefirst(master)holdthemaingatetype,thedatainputand
thenon-invertedoutput,andthesecond(slave)theinvertedoutputand
presetorresetinput.Theslavegateisgivenaspecialtypetoinstructthe
simulatortoeitherevaluatethemastergatedirectly(asynchronous
reset/preset)oronlymakesurethemasterisscheduled(synchronous
reset/preset),whilethemastergatetakescareoftheevaluationandfan-out
updateofboth.(SlavegatematchtheasynchronouspartinEMIL.)

Nowallsynchronousgatesarereformulatedtousethis'generic'double-
gate.Forexample,aJKflip-flopisbuiltwiththeJKfunctionasan
asynchronouspartinfront:

JKfrontfunction:             D=J'K'QN'+JK'+JKQN

genericdoublegateflip-flop:   QN=(Q=D)':onCP

Whenitcomestothetimingissues,MICAcanhandlethosegatesinsame
fashionasEMILwithaspecialorderedsynchronousgatequeue,orby
insertingbuffersbetweenconsecutivesynchronousgates.

## 4.4  GateSplitting

Here,thecountclassificationisjoinedwiththenetlist(splitofboolean
functionisalreadydoneintheclassificationstep).Foreachgateinthe
netlisttheclassificationofthatgatetypeisexamined.Iftheclassification
consistsofseveraloutputfunctions,newgatesarecreatedwithinputsand
outputaccordingtothebooleanfunctionofthatoutput.Thenewgatesthen
replacetheoldgateinthenetlist.

24

In this step we also add buffer gates as necessary to handle large fan-outs, and between synchronous gates if we want the faster synchronous evaluation switch.

## 4.5 Levelization

This is a technique to group gates independent of each other together, see section 2.2.3, Gate ordering. MICA reuses the levelization from EMIL. Inputs and synchronous gates are given a start level. Then gates whose fan-in gates all have a level are given the level next higher than the level of the highest fan-in gate, until all gates have a level assigned.

## 4.6 Netlist Optimization

This step creates the data structures actually used during simulation, removing all information not needed for simulation. The levelized netlist is kept with pointers from each gate to corresponding gate in the simulation netlist. Thus, we can still easily access the netlist by gate and net names in order to get observability.

The first step in this process is to determine the position of each gate in the simulation netlist, and the amount of memory needed. Currently the gates are simply ordered as they appear in the netlist, in the future however, other orders might be considered, see section 7, Future Optimizations. Once the position of each gate is decided the structure area is created, for each gate adding count type, output at, level and position of fan-out gates. Last, the data area is created, initiating and scheduling every gate. The initialization is done for each gate by setting the old output to low (also for complementary outputs), and storing the start value for the gate type as gate count (since all old outputs are set low, there are no high inputs, and the start value applies). As a final step an initialization half-cycle is run (using standard simulation functions), making all outputs consistent (making complementary outputs complementary, and outputs from logic_1 [1] high). Strictly, this is nothing but the first simulation cycle, but it is needed to get in synch with EMIL.

## 4.7 Simulation Netlist

The netlist used during simulation consist of two parts, one read only structure area and one data area containing all values that change during simulation. A third management part is used to be able to access the

---

[1] Logic_1 and logic_0 are special gates used to store the constant values 1 and 0 in the netlist. They do not have any inputs, nor do they change during simulation. They have their own FixAt count type.

simulationnetlistbygateornetnames,butthisw illnotbedescribed.A gateisrepresentedbyanoffsetusedtoindexinto botharesasviewedin Figure8below.

This separationofthenetlistintostructure,desc ribingfixedproperties suchashowgatesareconnected,anddataholdingd ynamicthingssuchas outputvaluesallowustoexploit hierarchy.Severalindependentdata areascansharethesamestructure area.Forexample,anetlistdescribing aprocessormightcontainseveral integerunits.The implementationofthis canthenbe describedwith one

Structure : : : Data : : :

Gate at offset i

structurearea,but eachinstancegetits owndata.Thisway memoryissaved,but thedrawbackisthat structureanddatais separatedinmemory, reducingthedatalocality.

Structure
```
count type
output at
queue
fan-out 1
fan-out 2
fan-out 3
```

Data
```
count value
old output
if scheduled
```

**Figure8Simulationnetlist**

Figure8showthepropertiesstoredinrespectivea reaforeachnormalgate andthenetitdrives [1].Foreachgateinthestructureareathe *counttype* tells whichevaluationmethodtouse(AtZeroetc.), *outputat* containtheoutput whentheconditiondeterminedbycounttypeistrue , *queue*holdthequeue thisgateshouldbescheduledinandlastwehaves lotstostoretheoffsetto fan-outs.Thedataareacontainhowmanyinputsto thegatearehighinthe *countvalue* (butwithzerobiasdependingonoriginalgatetyp e),theold output,andifthegateisscheduled.Inthefigure weseetheadvantageof

---

[1]Toavoidconfusionsomeclarificationisinorder.InMICAagateisrestrictedto haveonlyoneoutput.Thenetdrivenbythisoutputisstoredtogetherwiththegateas fan-outsandoldoutput.Thusthestoragespaceallocatedtoagatecanequallycorrect bereferredtoasanet.Althoughthisdocumenttriestorefertoitasagate,theter mnet ismorenaturalinsomecases.

26

thecountingalgorithm,eachgatecanbestoredwithoutknowledgeofitsfan-in,andwithoutknowingwhichinputagivenfan-outisconnectedto,onlyknowingtheoffsettoeachfan-outgateisenough.ThiscanbecomparedtoEMIL,whohadtostorebothfan-inandfan-outpointers.Anotherdifferenceimpliedbythecountwayofstoringgatesisthatwheneachgateisevaluateditupdateseachofitsfan-outgatecountswhileithavetoaccessthemforschedulinganyway.EMILdidtheotherwayaround,eachgatehadtobothfetchtheinputvaluesfromfan-ingatesbeforeevaluationandaccessthefan-outgatesforscheduling.

Thecurrentimplementationuseseightbytesstructurepergateandonebytedata,totalninebytespergate.Thegoalwastoreachfourbytesstructure,buttherewasnotenoughtimetoinvestigatethis.Sixbytescanbeeasilyachieved,butisnottested.Usingtwo,four,oreightbyteshastheadvantagethatthegateoffsetcanbeeasilyscaledinhardwareaddresscalculation,anditdividesonecachelineevenly.

## 4.8 Simulation

### 4.8.1 Gateevaluation

Thesimulatorusesthecountingmethodasdescribedearliertocalculateeachgateoutput.Thismeanwewillhaveonlyafewgatetypesofwhichonlyone,AtZero,standforthemajorityoftheevaluations.Todecidewhichevaluationfunctiontouseforeachgateweneedamainsimulatorswitch.Thismightseemreallybad,withapotentialbranchmiss-predictionforeachgate,butsince73%oftheevaluationsarefromtheAtZerogatetypetheprocessorwillbeabletodoagoodbranchprediction.Fortheasynchronousgatesandqueuestheevaluationlooptoreachastablestatelooksasfollow:

```
loop forever
  pop gate from queue
  case AtZero
    evaluate AtZero gate
  case AtEven
    evaluate AtEven gate
  case AtPos
    evaluate AtEven gate
  case UserDefined
    call supplied function
  case QueueSwitch
    if scheduled asynchrounous gates left
      switch queue
    else
      break
  end loop
```

Forreadabilityonlythemaintypesareincludedin theswitch.The
correspondingswitchforthesynchronousgatesexis tintwoversions.One
examinesallsynchronousgatesinapredefinedorde rjustlikeEMIL,and
evaluatesthechangedones.Theother,whichisfas ter,assumesno
consecutivesynchronousgatesexistandhasthesam estructureasthe
asynchronousversionabove.Theonlydifferenceis thequeue-switch(we
shouldevaluateonlythesynchronousqueue),andth echoicesintheswitch.

EvaluationofthedominatingAtZerotype(thiscan beimplemented
mostlyunconditional)isdescribedbythispseudoc ode:

```
if count value equals zero then
   new output = output at
else
   new output = not output at
end if
change = new output – old output
old output = new output
if change not equals zero then
   for all fan-outs
      fan-out count += change
      schedule fan-out
   end for
end if
```

Firstwedeterminethenewoutputdependingonthe *countvalue* and
*outputat* value(seeFigure8above).Thenwedetermineift heoutput
changed.Ifitrosethefan-outsshouldbeincremen ted,andifitfellwemust
decrement.Lastwestorethenewoutputandapplyt hechange(ifany)to
eachfan-out.Onlythefirstconditiondiffersint heothercounttypes,AtPos
andAtEven.

### 4.8.2 Mainsimulationloop

PerhapsrecognizedfromEMIL,wefinallygivethem ainsimulationloop:

```
for each half-cycle
   clock = not clock
   interface Leon ports
   evaluate asynchrounous gates
   if clock is high then
      evaluate synchrounous rising edge gates
   else
      (no falling edge gates)
   end if
   evaluate asynchrounous gates
end for
```

28

First, assume we are in a stable state, meaning all values have propagated to a register input or circuit port. Nothing change, everything is waiting for circuit in-ports to change. That tells it is a good time to look at the out ports, and update the in-ports.

In the case of Leon we update the data bus according to Leons wishes expressed on all out-ports (address bus etc.), or we can reset Leon by changing the reset port. As the last port interaction we let the clock tick a half cycle and update Leons clock port. After this some asynchronous gates depending on the updated ports need to update (they possibly affect inputs to some synchronous gates), so we must evaluate them. Now all synchronous gates have correct inputs and depending on whether the clock rose or fell those are updated (Leon has only rising edge flip-flops).

The synchronous gates will in turn change the data on the fan-out inputs, and careful here, so we call the evaluation of asynchronous gates a second time. Again we have reached a stable state, all circuit in-port changes has been processed, and we are waiting for new changes (most likely the clock).

Now, what was that about careful? As said before the synchronous gates have to be evaluated in correct order. The synchronous gates update all at once, so output changes from one does not have time to affect the inputs of a consecutive synchronous gate, immediately following the first. But if we update the second after the first precisely that will happen, that the inputs might have changed. This does not occur when an asynchronous gate is placed between the synchronous gates, since the evaluation of this one occurs after all synchronous gates.

A potential problem with the synchronous scheme used here and in EMIL is that at the same once as synchronous gates update, also latches (handled with the asynchronous gates) that are just enabled update, potentially causing same order problem for latches sooner or later affecting a synchronous (or latch) input. Luckily Leon latches is enable low, and synchronous gates rising edge, meaning latches always are disabled when synchronous gates change, and we can ignore this potential problem for now.

### 4.8.3 Queue system

Gate queues are handled in the simplest manner. Each queue receives a memory area fitting all gates in that queue, implemented as a stack with the bottom at the beginning of the memory area. A special queue-switching gate is placed at stack bottom (much like queue trailer in [Maurer97]), automatically switching to next queue when the current is empty. A scheduled gate count is used to keep track of scheduled gates. This has the benefit of not having to check all queues for gates every cycle, and it

handlesgateloops(gatesscheduledinqueuealreadyevaluated)by
simplywraparoundtofirstqueueiftherearegateslefttoevaluatewhen
lastqueueisfinished.Thedrawbackisextraincrementanddecrement
instructionsfor *each*evaluatedgate(thecheaperwaywouldbetoidentify
allloopgatesandcreateaspecialexceptiontypeforthose).MICAhas61
queues,evaluatedtwiceeveryhalfcycle,meaningchecking12million
queuesforgatesifallqueueshavetobechecked,withthegatecountonly
4.5millionisneeded.

## 4.9 Testsetup

Thegoalofthisworkwastobeabletosimulatelargenetlistsreasonably
fast.Toknowhowthesimulatorperformswithalargenetlistwehaveto
loadsuchanetlistinthesimulatorandactuallytest.Onlyreasoningasof
whattoexpect,asdonewithEMILabove,mightgiveinaccurateresult.

Herewefacesomeproblems,thefirstisthatwehavenolargefreely
availablenetlist,nextisthatloadingsuchlargenetlistwouldtaketoolong
time(readingthenetlistfastwasnotthegoalof neitherthisworknor
EMIL,thustheslow(doublelinkedlist)EMILcodewasreused).Third,it
wouldrequiremuchextraworktosynthesizeanewnetlist,adaptspecial
components,andwritetestprogramsetc.

Allofthoseproblems,aswellaslettinguscomparetoMICAagainst
EMIL,weresolvedbyinstantiatingtheLeonnetlistupto256times,and
simulateinstancebyinstancecyclebycycle.Since MICAwas
implementedwithhierarchicalnetlistsinmindthis showedtobeeasy.
(Hierarchicalsupportisnotentirelycompletedor testedsinceLeonhasno
sharedhierarchicalinstances,andtheoldnetlist loadingcodeflattensthe
netlist.)

Thesolutionresultedinalargenetlistconsistingofupto256Leondata
instancessharingthesamestructure.Notethatin thecaseofmorethana
coupleofLeoninstancesitisprobably,incomparisontoarealmillion-gate
netlist,notrealisticthatthedegreeofsharedstructureisthishigh.When
examiningtheresultwemustkeepthisinmind.EachLeonwasconnected
toitsownmemoryandprom,thelatterloadedwith oneoutofthreeslightly
differentmatrixadditionprograms.ThuseachinstanceofLeonwasrun
independentoftheothers.ForthepurposeofsimulatingseveralLeon
instancesthemainsimulatorloopwasslightlymodified:

```
for each halfcycle
   clock = not clock
   for each instance
      interface Leon ports
      evaluate aschrounous gates
      if clock is high then
         evaluate synchrounous rising edge gates
      else
         (no falling edge gates)
      end if
      evaluate aschrounous gates
   end for
end for
```

As shown above we simulate as if it was one large netlist possibly interchanging data between different hierarchical instances each cycle. Simulate the first instance all cycles, then switching to the next instance and simulate all cycles again etc. (i.e. exchange the two outer loops) would have been possible in our case since no data exchange took place, but would not have been of practical interest since it would give unrealistic good cache behavior (it would be equal to multiplying the one instance execution time).

To conclude, we simulated several instances of a gate-level netlist implementing the SPARCv8 processor Leon, each Leon independently running a matrix addition program.

## 4.10 Verification

Of course we have to be sure the simulator is correct, i.e. that the loaded netlist has exactly the same state (viewed from gate level) as the corresponding physical implementation in each stable state. Otherwise the simulator would be useless. This comparison was done against the previous simulator EMIL. However, due to the old implementation of the netlist in EMIL, which was reused in this work, finding the logic value of a net is slow. So this check was not done in each stable state, but only at sparse points and last simulated cycle, based on the assumption that most errors would remain and multiply in the next cycle (which do tend to be true from debug experiences). Also, an exhaustive check would only be useful if we are sure the reference simulator is correct.

The result of the program running on the simulated netlist Leon was also checked to be correct (though this does *not* guarantee the correctness of the simulator).

31

# 5 Result

MICAwastestedwithdifferentnumberofLeoninsta  ncesand50000 cycleswithresultslistedinTable3andFigure9.      Thetablelistsinturn numberofLeoninstancessimulated,totalsimulatio  ntimeforallinstances, totalmemory,totalmemorytimesgateactivity,exe  cutiontimeforeach Leoninstance,memoryusedperLeoninstance,andf  inallyhowmany byteseachgateusesonaverage.

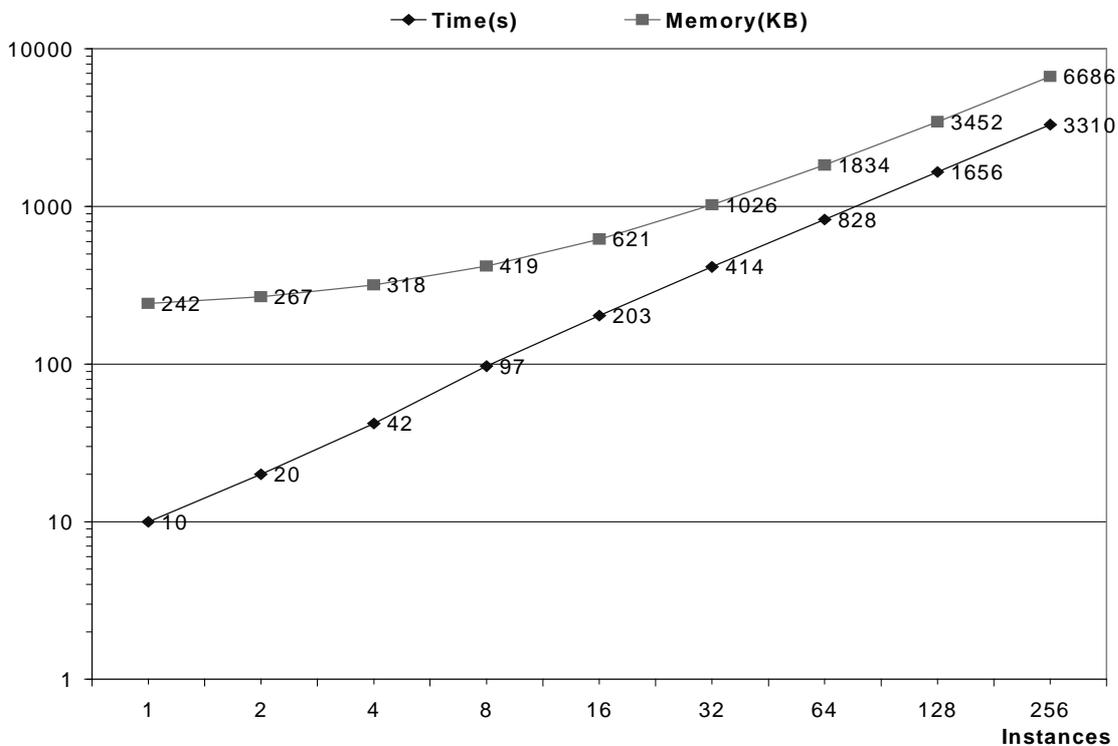| Leon Instances | Time (mm:ss) | Memory (KB) | Critical Mem(KB) | Time/Inst. (s) | Mem/Inst. (KB) | Mem/Gate (Bytes) |
|---|---|---|---|---|---|---|
| 1 | 00:10 | 242 | 28.4 | 10.0 | 242.1 | 11.1 |
| 2 | 00:20 | 267 | 31.4 | 10.0 | 133.7 | 6.1 |
| 4 | 00:42 | 318 | 37.3 | 10.5 | 79.5 | 3.6 |
| 8 | 01:37 | 419 | 49.2 | 12.1 | 52.4 | 2.4 |
| 16 | 03:23 | 621 | 72.9 | 12.7 | 38.8 | 1.8 |
| 32 | 06:54 | 1026 | 120.4 | 12.9 | 32.0 | 1.5 |
| 64 | 13:48 | 1834 | 215.3 | 12.9 | 28.7 | 1.3 |
| 128 | 27:36 | 3452 | 405.2 | 12.9 | 27.0 | 1.2 |
| 256 | 55:10 | 6686 | 785.0 | 12.9 | 26.1 | 1.2 |

**Table3Benchmarks**



**Figure9Executiontime**

32

The benchmarks was run on an Athlon Thunderbird 1466MHz (64+64KB exclusive L1 cache, 256KB L2 cache, and 512MB RAM), compiled with gcc 3.2.2 –O3 –march=athlon. Old EMIL was also reru non an identical configuration, 50000 cycles took 7.25s (from Table 4 on page 35).

For the purpose of comparison, and since a buffered solution can be implemented in EMIL as well, MICA was run with ordered synchronous gates (the order problem was explained in section 2.2.3, Gate ordering). Using the buffered solution gave an 11% performance increase (8.9s). 26 buffer gates had to be added after gate split, causing a total of only 30 extra evaluations (in EMIL we must expect much more extra buffers (~800) and evaluations).

The graph in Figure 9 shows the execution time and memory requirements from Table 3. We use a logarithmic scale on both axes to get a unified distance between the plots (to avoid all plots ending up in lower left corner). The memory requirements behave as we can expect from the implementation. With few instances structure memory dominates, whereas data dominates with many instances. Looking at the execution time it seems really good, increasing only linearly with the number of Leon instances. However, looking closer, things are a little worse.
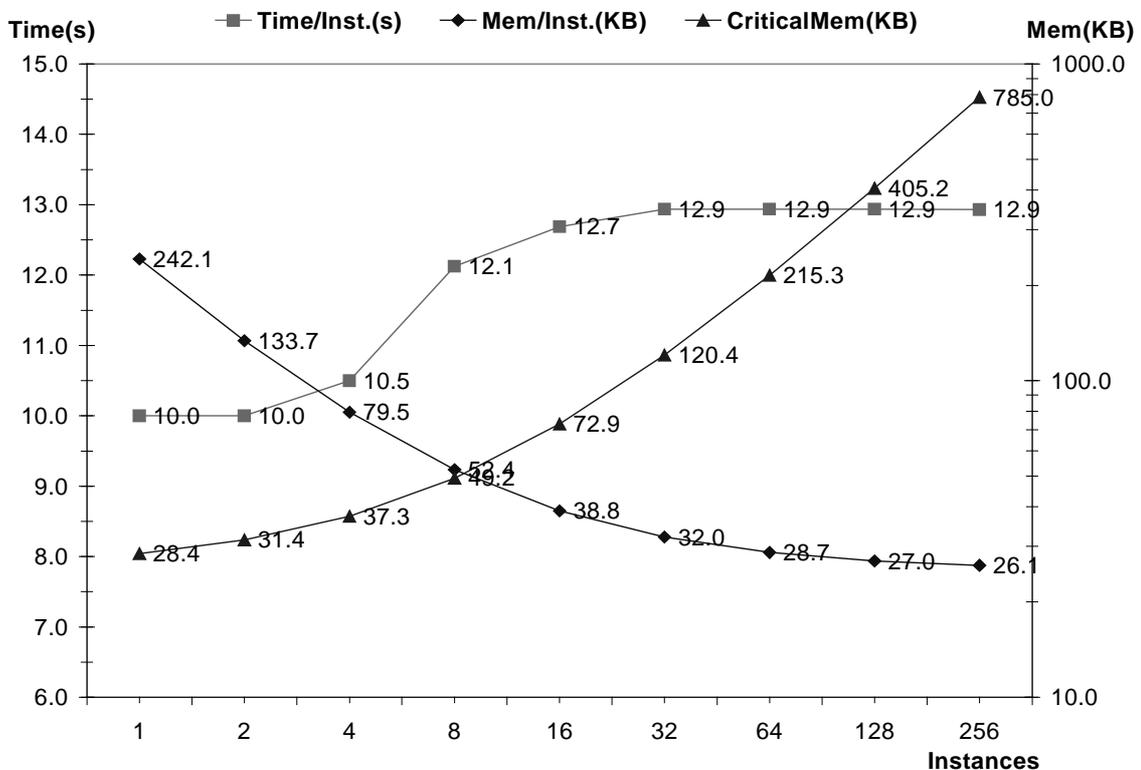


**Figure 10 Per instance graph**

InFigure10thegraphovermemoryandexecutiontimeperinstanceshowsusweactuallyhavea30%performancedecreasegoingfromapproximately4to16instances.WhatweseecanbeinterpretedasthemovefromutilizingmostlyL1cachetoutilizingtheL2cache.Withuptofourinstancesmostfrequentlyusedgatesfitinthe64KBL1datacache,butwith16to32instancesthisisnotenough.Thisinterpretationwouldexpectasimilarperformancedecreasewhenreachingabove128instances,whennoteventheL2cacheisenough.Amazinglythisdoesnothappen,andtheexplanationisprobablythatwithseveralLeoninstances,thestructureareaconstantlyoccupiestheL1cache,whilethedataareasarethrashedbackandforthbetweenmainmemoryandtheL2cache.

ThisisaneffectofthefactthatthememorydemandwithmanyLeoninstancesisstronglyidealized,sharingonestructureforallinstancesisgrowingunrealisticthemoreinstanceswehave.Thus,arealmillion-gatedesignmostlikelyneedsmorememoryforstructurearea,increasingmemoryrequirementanddecreasingperformancebymorethrashingincachesandmemory.

Themeasureofmostfrequentlyusedmemory(criticalmemory)canbediscussed.Itisnotlikelythatitisthesamesubsetofthegatesthatswitcheachcycle.Butitisasimplewaytogetaroughhint.Fromthegraphwecansethattheperformancedegradationoccursintheintervalfrom30to90KB,justaroundtheL1datacachesize.

Usingthestatisticstodosomecalculationswithamorereasonablestructuresharingofthreetoone,wefindthat3194KBmemorywouldberequiredalreadyat32instances.Thatmeanacriticalsizeof374Kandmorecannotbeexpectedtogivegoodcachebehavior.Thus,wecannotexpectthespeedof117M/12.9=9Mgatespersecondwithmorethan~275000gates(~635000countgates).

TheobservantreadermightalreadyhavenoticedthatMICAactuallyis38%slowerthanEMILsimulatingoneLeon50000cycles.Isthisallabigfailurethen?Well,notreally.Themainpurposewastosimulatelargedesignsfast,andfromthediscussionofwhattoexpectfromEMILwithmanygates,MICAwouldbefasteratthistask.(OfcourseitwouldbeveryinterestingtoactuallytestEMIL'sperformanceonalargedesign,oronseveralLeoninstances,butduetothelackoflargefreelyavailabledesignsandtheimplementationofEMILthiswouldbetoocumbersome,andthediscussionhavetosuffice.)

WhatwecandoistolookatthestatisticsfromEMILandMICAinTable4andtrytoexplainwhyMICAisslower,andhowtoovercomethis.

|  | EMIL | | | MICA | | |
|---|---|---|---|---|---|---|
|  | Total | Each cycle | Each second | Total | Each cycle | Each second |
| Executiontime | 7.25s | 145ns | - | 10.00s | 200ns | - |
| Simulatedcycles | 50000 | - | 6897 | 50000 | - | 5000 |
| Gateevaluations | 64M | 1278 | 8.8M | 117M | 2331 | 11.7M |
| Synchronous | 4.0M | 79 | 548K | 3.4M | 69 | 343K |
| Queueswitches | - | - | - | 4.5M | - | - |
| Schedules[1] | 22M | 447 | 3.1M | 131M | 2614 | 13.1M |
| Dummyschedules | - | - | - | 13M | 266 | 1.3M |
| Memoryused [2] | 348KB | | | 242KB | | |
| Gateactivity | 14.8% | | | 11.7% | | |
| Criticalmemory [3] | 52KB | | | 28KB | | |
| Totalgates | 8637 | | | 19856 | | |
| Totalnets | 10858 | | | 21896 | | |

[1] *CalculateddifferentlyinMICA*

[2] *Notincluding11328bytesforLeoninternalregist ersandcache*

[3] *Memorytimesgateactivityasasimpleapproximati onoffrequentlyusedmemory*

**Table4Comparison**

Thefirstwenoteisthehugedifferenceinschedul esandevaluations.For schedulesthisispartlybecausetheyarecalculate ddifferently.EMIL countscallstoascheduleroutinethatschedulesa llfan-outforagate(one callperchangedoutsignal).MICAcountseachgate actuallyputin evaluationqueue.DummyschedulesoccurinMICAwhe nthegateto scheduleisalreadyscheduled(equaltothenumber ofgateswithmultiple changedinputs).Theotherreason,validforboths chedulesand evaluations,isthatMICAhasmoregatestoevaluat eduetosplitgates. SomethinginterestingisthatMICAevaluates11.7m illiongatespersecond wileEMILonlyevaluate8.8million,thusMICAis3 3%fasterpergatebut still38%slowerintotal.Theexplanationisalso heretobefoundinthe totalnumberofgates.MICArequirescountcompatib legateswithonlyone outputpergate,meaningalotofgateshavetobe splittosimplerversions.

DespitethislargenumberofgatesMICAismuchmor ememoryefficient. MICArequires30%lessmemoryforoneinstanceofL eon.Thus,the countingmethodofstoringandevaluatinggatesare verymemoryefficient. Themainreasonforthisisthatallfan-inpointer scanberemoved,saving about50%memorypergate.

ToseewhereallextragatescomefromwehavetolookinAppendixA, Statistics.Sinceitisthenewnetsthatcauseevaluationsnotdonebefore[1], welookatwhichnetsthatareaddedratherthangates.Becausegatesand netsprimarilyhaveaone-to-onerelationshipaddednetscloselymatch addedgatesanyway.(Onlypads,synchronousgatesandsomeaddershave severaloutputnets.)Originallythereare10858netsinLeon.10584are addedbecauseofsplitcountincompatiblegates.454areaddedtobuffer gateswithhugenets(sinceweneedAtZerogatestoholdhugenets).The statisticsincludeonlythegatetypesoutofthe84usedthatneedspecial treatment.TogglestatisticsfromEMILisalsoincludedtogetherwitha worst-caseestimationofhowmanyextraevaluationsthesplitgatewould causeforasingleinputchange.(Schedulestatisticsshowthatonly10%of thegateshadmorethanoneinputchange,butwemustrememberthisis collectedinMICA,aftersomegateswithmanyinputsaresplit.)Adding thetotalsumofestimatedextraevaluationswiththetotalevaluationsin EMILgiveusatotalsumof $64M + 67M = 131M$ evaluations,slightlymore thanthoseofMICA.Itisthusreasonable(anddisappointing)tosaythat theworstcaseisgenerallywhatweget.Earlyestimationswasto optimistic,andtheassumptionthatachangeinoneinputtoacomplexgate wouldnotleadtotheevaluationofallpartsinthesplitversion,leadingto reducedgateactivity,hadonlylittleeffect.

ButlookingatthestatisticswealsoseethatthemainproblemisFDS2L andthesmallmultiplexers(discussedinsection3.3.1,Problem1). Togethertheystandforoverhalfoftheextraevaluations.Creatingspecial typesforthoseisexpectedtopayoff.Somelateexperimentsindicated 11MlessevaluationsiftheFDS2Lgateusesaspecialtype.

---

[1]Itisthenetsthataregiventheresultsofgateevaluations.Thecalculationofnet valuesmusthavebeendoneinsomewayforallnetsthatexistedinEMILsimulation. NewgatesstoringtheresulttooldnetsjustmovethiscalculationtothenewgateThus itismainlythenewnetsthatcreateworkinMICAthatwasnotdoneinEMIL.

# 6 Conclusion

Simulation is an important part of digital circuit development and verification. But hardware simulators are expensive and software simulation is slow on today's large designs.

We have implemented MICA to simulate large designs. One major part of this work has been to reduce the memory required for such design, in order to use the faster cache memory in a higher degree. The Count Algorithm used has shown to be very space efficient, by not storing fan-in information. We use 30% less memory than the reference simulator EMIL, and according to new statistics from MICA this can be further reduced by as much as 40%, see Reduce structure size further in section 7.

This reduction is despite the major drawback of the algorithm, that incompatible gates have to be split or specially handled. The MICA simulation netlist contain twice as many nets as the original netlist. This is also the reason that the reference simulator is in total 28% faster than MICA, at least for small netlists. Seen from evaluation speed MICA is 33% faster, evaluating 12M gates per second compared 9M per second for EMIL. We note that the overall speed of MICA depends on the Logic Library. When using a library containing only simple gates MICA would be faster also in total.

To emulate a netlist with a high number of gates we simulated up to 256 instances of our testing netlist. This corresponds to 2.2M gates (5.1M gates in MICA) and was simulated at 9M gates per second, which is the same speed as EMIL despite a much larger netlist. However, due to unrealistic hierarchy sharing, we can't expect this speed with more than 250000-300000 gates (0.6M in MICA). A larger cache would raise this limit.

From simulator perspective this is reasonable fast, on average 163 processor cycles are used per evaluated gate. From a user perspective however, the speed attained is very slow. A real circuit can be clocked in MHz, while we in simulation of 300000 gates can measure this in Hz only. With this slowdown simulation of one second circuit time would take 12 days. If we need to simulate the circuit for more than some hundred thousands cycles this is not feasible. It seems software gate-level simulation of huge netlists have a long way to go, and much work remain.

To conclude, MICA shows interesting features when it come to memory requirement and gate evaluation speed. Gates can be stored without fan-in information, and evaluated fast with a standard function. The major drawback is that not all gates are directly compatible, requiring special handling. In section 7, Future Optimizations below we briefly describe some ideas to reduce the penalty of this and other ideas.

# 7 FutureOptimizations

TherearemanywaystomodifyMICAinordertoincreasespeedorreducememoryneeds.Someofthesewasintendedtobepartofthiswork,butwascanceledduetolackoftime.Thoseimprovementsareinsteadpresentedinthissection.

## 7.1 RemoveFixAtgates

TheLeonnetlisthas662logic_1gatesand320logic_0gates.InMICAthosearekeptasis.Oneoptimizationwouldbetoremovethosegates,sincetheyarenotneeded.Inthecaseoflogic_0gatestheycanjustberemoved,andlogic_1shouldbepossibletoremoveafterfirstincrementingeachfan-outgatescountstartvalue.Thiswillnotaffectsimulationspeed,butreducethestructureanddatasizebyalmosteightandonekilobyterespectively.Fromobservabilitypointofviewthismeanthosegatesandnetswillnotexistinthesimulator,whichmightbeundesiredorconfusing.

## 7.2 Removeunnecessarygates

SomegatesinLeon,mostnotablyflip-flopsandadders,havetwooutputs,butinsomecasesoneoutputisunused.Whenthosegatesaresplittoonecountgateforeachoutputthecountgatewiththeunusedoutputcanberemoved.Atotalof428suchunusedoutputsexistinLeon,butmanyofthosearefromflip-flops,andsincetheyaredouble-gatesthisdonotapplytothem.ThusnonoticeableimprovementisexpectedforLeon.

Second,allinvertersbeingtheonlyfan-outofagatecanberemovedbysimplyinvertingtheoutputatvalueforthefan-ingate,butweloseobservabilityforthenetconnectingtheinverteranditsfan-ingate,sincethenetisremoved.Mostsuchinvertersareremovedinthesynthesis[EMIL2002].

## 7.3 Reducestructuresizefurther

Theoriginalplanwastouseasmartgateorderingtominimizethedistancetofan-outgates,andthenstorethefan-outoffsetsineachgaterelativetoitsownposition.Thiscanenableuseoflessstorage(smallerpointers)perfan-out,andisstilltobeinvestigated.

Oneinterestingthingwhenitcomestosizereductionisthefan-outstatisticsaftergatesaresplit.Itshowsthat77%ofallgateshaveonlyonefan-out.Sincethosegatesusesonlyhalfthespaceallocatedforagateinstructureareaalmost40%ofthestructureareaisunused.Moreover,forthemostcommongateAtZero,thisfactistruefor89%(65%oftotalevaluatedgates).Ahighlyoptimizedonefan-outAtZerotypemightbeworthwhile.

Wecouldforexampleshrinkthestandardgatestora gesizetoincludeonly
onefan-out,andstorelargerfan-outspecial.

Whattakesspaceandtimeistostoreandloadpoin terstoinputsand
outputs.Sowecouldpartitionthenetlistinsmall asynchronousnetworks,
forexamplenetworksbetweensynchronousgates,and usesmallerpointers
withineachnetwork.Thiswouldbebestdonetored uceinputandoutputs.
Groupingxgateswithyinputstoonewithxyinput sdonotgainvery
much.

## 7.4 SpecialMUX21andFDS2Ltypes

Thisissomethingreallyneeded.Thoseareverycom mongatesinLeon,
andprobablyanynetlist.Despitetheirverysimple functiontheyaresplitto
4-5gates.Estimationsshowthat26%ofallevaluat ionscouldberemoved
ifeachofthemcouldberepresentedwithasingle gate.Thedrawbackof
addingextragatetypesismorechoicesinthemain switch.

## 7.5 Circularqueuesystem

AsimpleversionofthiswasimplementedinEMIL,b utitwasnot
successfulduetothewaythelevelizationisdone. Ittriedtoexploitthat
gatesmostoftenscheduletheirfan-outinnextlev eltoreducethenumber
ofqueues.Hereisavagueideaofpreparingforth isalreadywhen
levelizingthenetlist.Ifweareabletorestructu rethelevelssonogatehave
theirfan-outfurtherawaythansay8levels,theq ueuejustevaluatedcould
bereusedasthequeue8awayfromthecurrent.As is,levelizationplaces
allsynchronousgatesinlevelzero,causinggates withlowleveltohave
manyfan-outsinalllevels.Insteadallowingsever alsynchronouslevels
couldincreasethelevelofmanygates,decreasing thegaptofan-outs.

## 7.6 Takecareofclocktiming

Apotentialproblemregardingevaluationorderofs ynchronousgatesand
latcheswereobservedinsection4.8.2,Mainsimula tionloop.Thisshould
besolvedsomehow,probablyinvolvingsomeclassifi cationofhowgate
netsdependondifferentclockevents.Twophasecl ockingisshortly
mentionedin[Jennings91].

## 7.7 Assemblemainevaluationloop

Byassemblingthemainsimulationswitchloopandt hemostcommon
typeswecouldprobablyutilizetheCPUregistersm oreefficient,and
removesomeinstructions.Performancewouldincreas eatcostof
portability.

## 7.8 Partitions and compiled super-gates

In previous works gates have been partitioned in groups to reduce scheduling overhead, for example [Blaauw93]. If we manage to partition the netlist into many small asynchronous nets, compilation could be done selectively for partitions with high activity.

Creating a dumb oblivious compiled version of a structure area seems fairly easy. We simply replace the count type and output at value with code for corresponding condition, each fan-outs with an add instruction, and order all code fractions by queue. The data area would not need to change. Interesting is that for such compiled version all inverters can be removed, simply by issuing a subinstruction instead of add when a fan-out points to an inverter. All we need is a flag identifying inverters since the count type do not tell.

We can also think of partitioning all gates whose output later on only affect one specific input to a multiplexer into one cluster. Such cluster would only need evaluation when that specific multiplexer input is selected.

## 7.9 Order gates by evaluation frequency

Currently gates have no specific order in the simulation netlist, and the locality of event-driven simulation is poor. MICA uses only one byte per gate in the data area, and 8 bytes per gate in the structure area, while a cacheline of an Athlon CPU is 64 bytes (32 bytes for a Pentium III CPU) [AMD]. This means seldom evaluated gates will waste much of each cache line. By placing gates in order of evaluation frequency the locality could be increased to utilize the cache more effective.

## 7.10 Hardware support

Calculating the new output of a gate is a rather simple operation, and so is the update and schedule operation of a single fan-out. Nevertheless those actions need quite a few assembler instructions each. Not that it will ever happen, but just for the interesting thought we could imagine if we had hardware support to do these operations in one clock cycle or instruction each.

# 8 Explanations

| | |
|---|---|
| ASIC | ApplicationSpecificIntegratedCircuit. |
| ASIP | ApplicationSpecificInstructionsetProcessor. |
| co-design | Simultaneousdesignofhardwareandsoftwareto speeduptimetomarket. |
| criticalsize | Dumbestimationofmostfrequentlyusedmemory bytheproductofgateactivityandtotalmemory. |
| EDIF | ElectronicDesignInterchangeFormat. |
| fan-in | Thegatesconnectedtoagatesinputs. |
| fan-out | Thegatesconnectedtoagatesoutput. |
| gateactivity | Theaveragepercentageofgateschangedduringone cycle. |
| gcc | GNUCompilerCollection. |
| half-cycle | Thetimeduringwitchtheclockishigh orlow,that is,thetimefromrisingtofallingclockedge,or from fallingtorising. |
| IP-package | IntelligentPropertypackage,circuitscanbe deliveredasnothingbutasoftdescriptioninsome languageorfileformat. |
| LCC | LevelizedCompiledCode. |
| level | Maximumnumberofgatestofurthestpreceding circuitportorsynchronousgate. |
| levelization | Methodtodetermineeachgateslevel. |
| logiclibrary | Acollectionofgatetypesandafunctional descriptionforeachtype. |
| net | Thewireconnectingagateoutputwiththeinputsof itsfan-outgates. |
| netlist | Designdescriptioninformofgatenetworks. |
| PCB | PrintedCircuitBoard. |
| PERL | PracticalExtractionandReportLanguage,atext manipulatingscriptlanguage. |
| queue | Storagespaceforoneorasetoflevels. |
| SPARCv.8 | ScalableProcessorArchitectureversion8. |

| | |
|---|---|
| stable-state | Time-point when all signals have propagated to a latch or register. |
| system-on-chip | Everything integrated in one chip, no more PCBs crammed with chips. |
| VHDL | VHSIC Hardware Description Language. |
| VHSIC | Very High Speed Integrated Circuit. |

## 9 Acknowledgements

## *10 References*

[COSMOS87] Randal E. Bryant, Derek Beatty, Karl Brace, Kyeongsoon Cho, Thomas Sheffler
*COSMOS: A Compiled Simulator for MOS Circuits*
24[th] ACM/IEEE Design Automation Conference, pp.9-16, 1987

[EMIL2002] Tomas Holmberg
*EMIL: An Event-Driven Multi-Queue Interpretive Levelized Gate-Level Simulator*
Thesis work, Microelectronics and Information Technology, Royal Institute of Technology, Sweden, 2002

[LECSIM90] Zhicheng Wang, Peter M. Maurer
*LECSIM: A Levelized Event Driven Compiled Logic Simulator*
27[th] ACM/IEEE Design Automation Conf., pp.491-496, 1990

[SLS88] Zeev Barzilai, Daniel K. Beece, Leendert M. Husmain, Vijay S. Iyengar, Gabriel M. Silberman
*SLS–A Fast Switch-Level Simulator*
IEEE Transactions on Computer-Aided Design vol.7, no.8, pp.838-849, August 1988

[SSIM87] Lang-Terng Wang, Nathan E. Hoover, Edwin H. Porter, John J. Zasio
*SSIM: A Software Levelized Compiled-Code Simulator*
24[th] ACM/IEEE Design Automation Conference, pp.2-8, 1987

[Ashar95] Pranav Ashar, Sharad Malik
*Fast Functional Simulation Using Branching Programs*
Proceedings of the 1995 IEEE/ACM International Conference on Computer-aided Design, pp.408-412, December 1995

[Bertacco99] Valeria Bertacco, Maurizio Damiani, Stefano Quer
*Cycle-based Symbolic Simulation of Gate-level Synchronous Circuits*
36[th] ACM/IEEE Design Automation Conf., pp.391-396, 1999

[Blaauw93] David T. Blaauw, Larry G. Jones
*Reducing the Scheduling Cost in Event-Driven Simulation Through Component Clustering*
4[th] European Conference on Design Automation, pp.18-22, 1993

[DeVane97] Charles J. DeVane
*Efficient Circuit Partitioning to Extend Cycle Simulation Beyond Synchronous Circuits*
Proceedings of the 1997 IEEE/ACM International Conference on Computer-aided Design, pp.154-161, November 1997

[French95]     RobertS.French,MonicaS.Lam,JeremyR.Levitt,Kunle Olukotun
*AGeneralMethodforCompilingEvent-DrivenSimulations*
32<sup>nd</sup>ACM/IEEEDesignAutomationConference,1995

[Gateways94]  PeterM.Maurer,YunSikLee
*Gateways:ATechniqueforAddingEvent-DrivenBehaviorto CompiledSimulators*
IEEETransactionsonComputer-AidedDesignofIntegrated CircuitsandSystems,vol.13,no.3,pp.338-352,March1994

[Jennings91]  GlennJennings
*ACaseagainstEvent-DrivenSimulationforDigitalSystem Design*
24<sup>th</sup>annualSymposiumonSimulation,pp.170-176,April1991

[Jiang2003]  YunjianJiang,SlobodanMatic,RobertK.Brayton
*GeneralizedCofactoringforLogicFunctionEvaluation*
40<sup>th</sup>ConferenceonDesignAutomation,pp.155-158,June2003

[Jones94]    LarryG.Jones,DavidT.Blaauw
*ACache-BasedMethodforAcceleratingSwitch-Level Simulation*
IEEETransactionsonComputer-AidedDesignofIntegrated CircuitsandSystems,vol.13,no.2,pp.211-218,February1994

[Lewis91]    DavidM.Lewis
*AHierarchicalCompiledCodeEvent-DrivenLogicSimulator*
IEEETransactionsonComputer-AidedDesign,vol.10,no.6, pp.726-737,June1991

[Maurer94]   PeterM.Maurer
*TheInversionAlgorithmforDigitalSimulation*
Proceedingsofthe1994IEEE/ACMInternationalConferenceon Computer-aidedDesign,pp.258-261,1994

[Maurer96]   PeterM.Maurer
*IsCompiledSimulationReallyFasterthanInterpreted Simulation?*
9<sup>th</sup>InternationalConferenceonVLSIDesign,pp.303-306,1996

[Maurer97]   PeterM.Maurer
*TheInversionAlgorithmforDigitalSimulation*
IEEETransactionsonComputer-AidedDesignofIntegrated CircuitsandSystems,vol.16,no.7,pp.762-769,July1997

[Maurer99]   PeterM.Maurer
*EfficientSimulationforHierarchicalandPartitionedCircuits*
12<sup>th</sup>InternationalConferenceonVLSIDesign,pp.236-241, 1999

[Maurer2000]  Peter M. Maurer
*Event Driven Simulation Without Loops or Conditionals*
Proceedings of the 2000 IEEE/ACM International Conference on
Computer-aided Design, pp. 23-26, 2000

[Olukotun98]  Kunle Olukotun, Mark Heinrich, David Ofelt
*Digital System Simulation: Methodologies and Examples*
35[th] Conference on Design Automation, pp. 658-663, June 1998

[Shadow93]  Peter M. Maurer
*The Shadow Algorithm: A Scheduling Technique for both
Compiled and Interpreted Simulation*
IEEE Transactions on Computer-Aided Design of Integrated
Circuits and Systems, vol. 12, no. 9, pp. 1411-1413, September
1993

[YeeAu91]  Wing Yee Au, Daniel Weise, Scott Seligman
*Automatic Generation of Compiled Simulations through
Program Specialization*
28[th] ACM/IEEE Design Automation Conference, pp. 205-210,
1991

[AMD]  Advanced Micro Devices Inc., www.amd.com, 2004

[EDIF]  www.edif.org, 2004

[Gaisler]  Gaisler Research, www.gaisler.com, 2004

[SPARC]  SPARC International Inc.
*The SPARC Architecture Manual version 8* , 1992
www.sparc.org, 2004

# 11 Appendix A, Statistics

## 11.1 Split statistics

| Gatetype | InLeon | Splits | Extra nets* | EMIL evaluations | EMIL Activity | WCET** | Extra evaluations |
|---|---|---|---|---|---|---|---|
| MUX21L | 730 | 4 | 2,190 | 6,992,736 | 19% | 4 | 20,978,208 |
| FDS2L1 | 676 | 5 | 2,950 | 3,808,294 | 11% | 5 | 15,233,176 |
| MUX31L | 34 | 7 | 204 | 1,176,269 | 69% | 5 | 4,705,076 |
| AO12P | 237 | 5 | 948 | 4,062,416 | 34% | 2 | 4,062,416 |
| AO7 | 536 | 2 | 536 | 3,698,926 | 14% | 2 | 3,698,926 |
| MUX21H | 78 | 4 | 234 | 965,651 | 25% | 4 | 2,896,953 |
| AO2 | 383 | 3 | 766 | 2,426,294 | 13% | 2 | 2,426,294 |
| AO1P | 284 | 2 | 284 | 2,032,912 | 14% | 2 | 2,032,912 |
| AO3 | 163 | 2 | 163 | 1,524,714 | 19% | 2 | 1,524,714 |
| AO4 | 266 | 3 | 532 | 1,520,091 | 11% | 2 | 1,520,091 |
| AO11 | 152 | 4 | 456 | 1,519,467 | 20% | 2 | 1,519,467 |
| FJK1S | 26 | 8 | 195 | 151,252 | 12% | 7 | 907,512 |
| EO1 | 76 | 3 | 152 | 855,586 | 23% | 2 | 855,586 |
| FA1A | 61 | 2 | 0 | 845,849 | 28% | 2 | 845,849 |
| AO6 | 109 | 2 | 109 | 775,939 | 14% | 2 | 775,939 |
| MUX81P | 8 | 12 | 88 | 67,956 | 17% | 12 | 747,516 |
| EON1 | 122 | 3 | 244 | 705,495 | 12% | 2 | 705,495 |
| IVDA | 52 | 2 | 0 | 401,923 | 15% | 2 | 401,923 |
| FD1S | 98 | 5 | 392 | 67,976 | 1% | 5 | 271,904 |
| FDS2LP | 3 | 5 | 12 | 58,395 | 39% | 5 | 233,580 |
| AO7P | 7 | 2 | 7 | 229,651 | 66% | 2 | 229,651 |
| B2I | 19 | 2 | 0 | 176,862 | 19% | 2 | 176,862 |
| HA1 | 40 | 2 | 0 | 145,418 | 7% | 2 | 145,418 |
| AO3P | 4 | 2 | 4 | 72,300 | 36% | 2 | 72,300 |
| FDS2 | 102 | 2 | 102 | 63,853 | 1% | 2 | 63,853 |
| FD4S | 1 | 5 | 4 | 5,539 | 11% | 5 | 22,156 |
| B3IP | 4 | 2 | 0 | 5,879 | 3% | 2 | 5,879 |
| IVDAP | 5 | 2 | 0 | 1,990 | 1% | 2 | 1,990 |
| FJK1 | 1 | 6 | 6 | 342 | 1% | 6 | 1,710 |
| MUX41 | 1 | 7 | 6 | 1 | 0% | 7 | 6 |
| FD1 | 575 | 1 | 0 | 624,193 | 2% | 1 | 0 |
| LD2 | 128 | 1 | 0 | 13,031,101 | 204% | 1 | 0 |
| FD2 | 27 | 1 | 0 | 20,848 | 2% | 1 | 0 |
| FD4 | 14 | 1 | 0 | 5,992 | 1% | 1 | 0 |
| FD1P | 1 | 1 | 0 | 67 | 0% | 1 | 0 |
| Total | 5,023 | | 10,584 | 48,042,177 | 19% | | 67,063,362 |

*FDS2L, FJK1S and FJK1 add nets for some earlier u    nused inverted outputs, thus adding slightly
more than expected (splits-1)*in_leon

**WCET (Worst Case Evaluations Triggered by single i    nput change)

I

## 11.2 Schedule and evaluation statistics

| Type | Schedules | Percent | Explanation |
|---|---|---|---|
| Real | 117,401,953 | 89.83% | Not scheduled before |
| Dummy | 13,287,387 | 10.17% | Scheduled before |
| Total | 130,689,340 | | |

| Type | Evaluations | Percent | Explanation |
|---|---|---|---|
| AtZero | 88,674,199 | 73.27% | Standard type |
| Eval1Up | 12,800,128 | 10.58% | Change on latch enable input |
| HugeNet | 9,000,661 | 7.44% | Standard type with large fanout |
| QueueSw | 4,473,446 | 3.70% | Queue switches |
| SynFF | 3,431,113 | 2.83% | Change on synchronous flip flop data input |
| AtEven | 871,013 | 0.72% | Full adder sum function |
| AtPos | 845,849 | 0.70% | Full adder carry function |
| Compiled | 672,523 | 0.56% | User-defined (Pads, registers and dcache) |
| Latch | 231,229 | 0.19% | Change on latch data input |
| FlipFlop | 26,762 | 0.02% | Change on asynchronous flip flop data input |
| FixAt | 982 | 0.00% | Logic 0/1 gates |
| Eval1UpIf | 42 | 0.00% | Change on flip flop clear or preset input |
| Total | 121,027,947 | | |

## 11.3 Fan-out statistics

*Fanout statistics for all gates*

| Fanout | Orignets | Percent | Aftersplit | Percent | Afterbuf | Percent |
|---|---|---|---|---|---|---|
| 0 | 306 | 2.82% | 306 | 1.43% | 306 | 1.40% |
| 1 | 6,483 | 59.71% | 16,480 | 76.86% | 16,934 | 77.34% |
| 2 | 1,948 | 17.94% | 2,254 | 10.51% | 2,254 | 10.29% |
| 3 | 833 | 7.67% | 940 | 4.38% | 940 | 4.29% |
| 4 | 381 | 3.51% | 448 | 2.09% | 448 | 2.05% |
| >5 | 907 | 8.35% | 1,014 | 4.73% | 1,014 | 4.63% |
| Added | 0 | | 10,584 | | 454 | |
| Total | 10,858 | | 21,442 | | 21,896 | |

*Fanout statistics for AtZero type*

| Fanout | Amount | Percent | Evaluations | Percent | Of total evaluations |
|---|---|---|---|---|---|
| 1 | 14,662 | 88.43% | 78,965,438 | 89.05% | 65.25% |
| Other | 1,919 | 11.57% | 9,708,761 | 10.95% | 8.02% |
| Total | 16,581 | | 88,674,199 | | 73.27% |

## 12AppendixB,Grammar

```
gate:       NAME { stmtlist }

stmtlist:   stmtlist stmt
|           stmt

stmt:       set ';'

assign:     OUTPUT '=' or

or:         or '+' and
|           and

and:        and not
|           not

not:        expr '\''
|           expr

expr:       '(' assign ')'
|           '(' or ')'
|           INPUT
|           LOGIC

clock:      assign ON INPUT
|           assign WHEN INPUT
|           assign

set:        clock RESET INPUT
|           clock PRESET INPUT
|           clock
```

| | |
|---|---|
| **Titel**<br>Title | Simulering av miljoner grindar med Count Algoritmen<br><br>The Counting Algorithm for simulation of million-gate designs |
| **Författare**<br> Author | Klas Arvidsson |

**Sammanfattning**
Abstract
A key part in the development and verification of digital systems is simulation. But hardware simulators are expensive, and software simulation is not fast enough for designs with a large number of gates. As today's digital designs constantly grow in size (number of gates), and that trend shows no signs to end, faster simulators handling millions of gates are needed. We investigate how to create a software gate-level simulator able to simulate a high number of gates fast. This involves a trade-off between memory requirement and speed. A compact netlist representation can utilize cache memories more efficient but requires more work to interpret, while high memory requirements can limit the performance to the speed of main memory. We have selected the Counting Algorithm to implement the experimental simulator MICA. The main reasons for this choice is the compact way in which gates can be stored, but still be evaluated in a simple and standard way. The report describes the issues and solutions encountered and evaluate the resulting simulator. MICA simulates a SPARC architecture processor called Leon. Larger netlists are achieved by simulating several instances of this processor. Simulation of 128 instances is done at a speed of 9 million gates per second using only 3.5MB memory. In MICA this design correspond to 2.5 million gates.

**Nyckelord**
Keyword
simulation, count-algorithm, gate-level, interpretive, event-driven, hierarchical, multi-queue, computer-aided design, logic evaluation, zero-delay