The Din eration Committee for Narth & eyas San barding an certifies that this is the approved version of the following dissertation:

Polymorphism: A Unifying Approach for Managing r Complexity, Enabling Technology Scale and Extracting Concurrency

Committee:

Steel or W. Kedder, Superstor

Table of Contents

Acknow is aggreent

Abstract

List of Tables List of Figures

Chapter 2 Introduction
3.3 Principle of Polymorph im
2.5 System design genellen
2.1 Granularity of Personer
2.2 Granularity of Personer
2.2 Granularity of Personer
2.3 Thurst Architectur
3.3 Thurst Architectur
3.5 Thurst Architectur
3.5 Thurst Architectur
3.7 Thurst Architectur
3.8 Thurst Architectur
3.9 Thurst Architectur
4.9 Thurst Architectur
5.9 Thurst Architectur
5.0 Thurst Architectur
5.0 Thurst Architectur
5.0 Thurst Architectur
5.0 Thurst Architec

Chapter 2. Related Work

I.3. Scalable architectures I.4. Microarch bed are techniques for ILP

7.1 Kem el central behavion. S.
7.3 Microard bed archived diagram. F.
7.3 Microard bed archived diagram. F.
7.3 Microard bed archived diagram metassium. Seftware managed codes, finitely and class sets and control form. F

7.4 Execution core and constrainment anisms, a) in struction, agreemed rectification and EF data storage, b) Local PC and EF to struction, storage is a model of MIMIO Secution.

7.5 Species only different mechanisms, relative to haveless archi-

To the contact clause of any brailings can be represented by different

Instruction-level Parallellem | E.P.: The professional type of published

to among the bridge marking or writing a red as memory leads, there

nith metic operations. The operations are since h RISC of the operation

and the system is handed a single program written with a sequentia

The differences had were application domains in dodes several other for

· Memory access sattems which include stream hardide roughly or man

. Tree of aditionally engration, namely fixed rotal or floating role is

Thread-level Parallellem |TLP|: Parallellem between multiple threads of

centag, etc. on a all he houses to a three types of parallels on:

the name or the far operation applied to them.

irmeniar a creum in nigel of receptor data or sectors

promore is mist [26].

tice at explaining threat-level parallellom and data-level parallellom. When execution of the mode on a charle a process on, sign Michaelle his her levels of a more or estimation are two, IPCs are in the range of \$2 to \$25 for an application mix manisting of SEMBC and SPEC CPU2000 workloads. When exeming programs with DEP, the pulyment has a mechanisms we propose provide has operage of L.IX across a set of DLP workloads, compared to a there mechanisms provide competitive performance using a single execution

List of Figures

1.1 Grandwitty of parallel processing elements on a chip. Number of cores that can the same typical fram chip. 4.1 Execution the management on the control of the 6.1 TLP-made performance | stitution| - SPEC CPUINS only. | IN-6.2 TLP-made (parkey compared to orbifold execution - SPEC CPUINS only. | ---

CPUINT sale.

TEP-ma de extention efficiency - SPEC CPUINT sale.

TEP-ma de performance | ordination | - EEMBC sale.

TEP-ma de symday compared to retalized execution - EEMBC de execution efficience - E EMB C onite.

5. Some a uniter of altra-wide interpresented the the Grid Processor | 1 III |-

a TRUPS of it like configuration we promote in this distortiation The fire contained the contemporal when applied to contain the

are played by typick makes had be a creshold a realiting from aggregating and high of these saids together. Course grain and declares using conventional wide-for the crain condition.

come unit come cas he subdicided when the unit combined with . On two key lanights are: If the the datathor graph as a hardeleted of abstract to express our correscy to the hundware to elastrate the hundware's seed for red bear wing concurrency, and reduce the hardware or otherdy of instruction Polymorphism: A Unifying Approach for Managing essor Complexity, Enabling Technology Scalabilit and Extracting Concurrency

Karthikeyan Sankaralingam, B.Tech., M.S.

DESERTATION Presented to the Faculty of the Graduate School of The University of Texas at Aprils is Parial Palifilment for the Degree of

DOCTOR OF PHILOSOPHY

THE UNIVERSITY OF TEXAS AT AUSTIN

Chapter 1. TRIPS Architecture

Chapter 4. Polymorphism in the TRIPS Architecture 4.1 Principles of Polymorphism

Mechanism
42.1 Execution are management
42.2 Control flow management
42.3 Data o longe management
42.4 Summary
Instruction-Level Parallelon

Acknowledgments

Me tittet om med light general med å rekladled by a syn k from another person. Kall af ar har some tit titte eith deep gratitate of titte eks hare lightet tit flores eith i

Many people have contributed to my dimeriation research and my exnorthway at U.T. I would like to that though may advise, Stone Keekler, for his advice, on ideace, and impining that helped me net to this exist in my ora excurregement and technical experies that made this dissertation possible State has been an excellent ments rand throb addisorand than and thank him

Date Berner to the reclinion of the CART cross and a dediction afe her fire my dio estation research has also played an important yest in m professional decelopment. I am thankful for the nomenos supportanties have had to interact with him, and for his many in sight follows ments on tenics maging from microarchitestore pipelines to herwing heer.

I would like to a danged edge the moral and technical copport that have received from the cividents in my research group, CART, I would like to thank them for their feedback on my research, aften ding many of my practice talks, prooffeeding my papers, and indulging in an memor technical and a sa-

4.4.1 Execution are management 4.4.2 Control flew management 4.4.3 Data of engagement 5 Threat-Level Parallelium 4.5.1 Execution are management 4.5.2 Control flew management

Chapter 5. Performance Evaluation: ILF

Chapter 5. Performance Evaluation: TLP

Rendi 62.1 SPEC CPUINN beschmark

makes of peneral perpose percesses has grown from an doubly over the pa

two decides. This improvement has come from deeper simplines and factor

and tere. Decke integration has played a large role in improving process

performance or well, enabling large nach hi multi-megabyte cocker, multiple

the ring point units on chip, and makes architecture structures to improve per-

females. Due to technology imitation of wire delays MI newer F21, and

transition is likely to clow down. Decke integration has already recited a

tion to extract more performance. As a result, performance growth in th

fiture must ome from extracting more concurrency from applications. As

difference and it extract concurrency at all levels, including thread level and

competents data level condiction, and not selven out the crain in time

tion level parallelian. But conventional architectures are poor at extracting so chi different granularities of parallelian and furthermore sely primarily on

large controlled structures like resister files, we are tables, and are deter-

to extract concerns on. Due to the afterment learnite technology limitations

technologies to in Burtille. There is instead a desire for scalable and modular

Breadly, the two treads that processor architects flor are: If emerging

architects res.

point where convenies and architectures are usually to a tilize more ex-chip tran

52 Beschmark 53 Resth

technical discussions that kept me entertained, informed, and always proceded my thinking. I would expectably like to though Romadon Nagaraja a with whom I had a

joyed working ingelber with him and confl not have asked for a l collaborator. Starting from a printible work on writing a PowerPC part of the Since left calculate since he see that a single beginning of the section of a faithful ideas for the TRIPS agreement and design and replication of the agreement chip, I have always support Manudas's insights and hearing his p The instruction on described in this dissertation and many aspects of the microarchited are were jointly designed by as-

I would like to thank the staff in the Computer Sciences department for The facilities staff in the department were noted anding, and the rare time they expected their BOFB (dect me were completely facilities) I would like to thank Gen Nathar for help in a with graduate orbital microflags, several travel in sec. as 4 ker immenter attence and tolerance.

Last, but not the least, I would like thank my family. My mather and for a core and in grows when I was going to graduated Thanks Range for parting ay with me as a commute for two years and to bending all my excention to

as a hid burdler for 25 years. Thanks Dad for your words of whiten

for design met had alogie that man achieve economies of scale, provide support for heteroceness annion tions, and combat the appreciar complexity arising na these technology fron 6. In this discernition, we introduce polynocybise computation. The key likes behind polymorphism is to our figure course grain

microard lied are blocks to provide an adaptive and diexible amount only

realish is an doma do his microarch it we are it in the . plication demands is to hold a betem promos chip, which contains multiple processing corps, and designed to many distinct class of workleads effect to do. The Taxasials processer is one example of integrated belong under \$3\$ The wa na ajor dawa i idao ta tibib ay pranch and a muni ad hardwa ne canay locky ciaco there is little design reason between the trace of processing and more recognic stilization when the application mix contains a balance different than that

ha magenessa process on , that militaring the afterm or thosed complexity prob-lem. The polymorphose nature of the process recens allows the hardware to be madigated to provide special purpose behavior on an application-byrelication havis, that adapting to a wide cange of application chares. Since

ributed architecture. The mechanisms to implement polymorphism are do-

religied in the mintext of this architecture. We chose this architecture as our haveling a real which to devote the medianisms for columns him because feature of the architecture are:

L. Dut affect day ender on are so coded in the LSA to enable direct instruction harmonia communication and reduce the everheads of detection and managing day enden der ihnt en aventional unterfereter ym een om mant pay. This new class of IS As ended EDGE Explicit Data Gray's Executing committee being darafter to the ISA, without having to charge programming models. Unlike VLIW architectures, the execution order determining the dynamic execution order.

2. The program is partitioned into well-defined blocks to limit the scape of the dependences on that the name or dependence are does not exceed the instruction space. Dependences inside on this block are exceded diently in the instructions, while dependences nones thinks are expressed through architectural registers or store-load pain. This execution madel fit does executed and committee for black of in the client attentically to relace the overheads of instruction management like reakter renam ing, dependence checking, and branch prediction. These merheads are

cation make and hardware capability does not arise since the hardware can be

In this discertation, we define a reliterized polymorphism and describe

a care of afterhick to which we halfd a san to decide mechanism to imple-

ment polymorphism. We describe the TRIPS architecture which is a technology combile and partitioned design. The TRIPS IS A to use in times of a new

charaf ISAs called Explicit Data Graph Execution | E DGE | which we propose

the hard ware. The polymorphonomedications are described in the context of the TRIPS making one. In the remainder of this chapter we provide a short

aversion of a signar skips, the TRUPS architecture, and conclude with a thesi

. We define architectural polymorphism as $h \in \operatorname{shifty} \operatorname{demonstrate}_{h} h$

An disnakh a' caras srain mi mar à i lethra blada ai rantin s. h. di mains

control ligic for leaving datapath and a large class onto largely normal/feel, in facility regression while predictors that our to specialized on an application

by application leads. The main principle of polymorphism are the following

which are developed in detail through the remainder of this discertation:

. Mississis army different emerlative of contletion

adapted at con-time to any application mix.

1.1 Principles of Polymorphism

I would also like to asknowledge the institutions that helped support

my research in am dance others! National Science Penndation for the initial

ted me and DARPA for the TRIPS for dia-

Appendix A. tsim-proc and GPA simulator comparison
A.1 Description
A.2 Resets

219

Appendix B. IPC reduction from apeculation depth

1. To makage dwigs or apleadly and address wire delay scaling, the compatation core is completely distributed as including the derive control and

1.4 Implementation of Polymorphism

ware at ma-time to perform different fractions. Unlike a remarks made arcompanents in some diefersolbesteinig finactions firm primitive logic blocks at

processor design and distinctions the solementation are much from at her This taxes one or mides a 4-1 as le that may be used to charify architectures accept a parties of this space. In chapter 2 which discusses related work, we classify other architectures a coording to this taxonomy. If show we briefly

Polymorphism: A Unifying Approach for Managing Processor Complexity, Enabline Technology Scalability and Extracting Concurrency

Karib kepas Saskanilaga m, Ph.D. The University of Texas at Astria, 2000

Supercior: Stephen W. Kenkler

ing applications with historogeneous computation needs, and technology limit tations of power, wire delay, and process rapidities. Designing marking and pleading, a self-wave programmability problems, and reduces economic of scale There is a previous a well for come device whilecombs that can appeld a comn or for heterocea was any likelings, combat a recessor complexity, and achie polymorphism to build each endable procession that provide copport for each het graceness come at ation by the series in a different granularities of carallelle adaptive and Bookle processor substrate. Technology scalability is additived

by a dwigning an architecture using scalable and made for microarchitecture

We not be dutaffere graph as the satisfying abstract is a layer across the s three types of parallelism. All programs are expressed in terms of dataffic required by the granularity of parallelion. We introduce EDGE IS As a chiof ISAs, as an architectural column for efficiently extremine camblelous for ha ili ing techa ah ay scalah le arch il ect ares.

We developed, the TRIPS architecture as an implementation of EDGE the TRIPS microarchitecture are its heartly partitioned and machalar design pared of two presenter core and a distributed RMB associations |NUCA sand it memory crotem.

Our performance receive class that the TRIPS makes architecture ca or take good in the clien-levely amiliation. On a set of has day limited benef IPCs in the maps of 4 to 6 are seen, and on a set of highly data pumbel range of \$.5 to \$2. With more approxime compiler optimization we expect

A.1 Companies of GPA constances THEPS constance in the DEP

List of Tables

Different parenters des similared.
 Bretchard may in 1-1 hand on figurities. SPEC CPCRIS
 Different may in 1-1 hand on figurities. SPEC CPCRIS
 Different may in 1-1 hand of the constraint of the constraint.
 Different may in 1-1 hand on figurities. SPEC CPCRIS
 Different may in 1-1 hand on figurities. SPEC CPCRIS
 Different may in 1-1 hand on figurities. SPEC CPCRIS
 Different may in 1-1 hand on figurities. SPEC CPCRIS
 Different may in 1-1 hand on figurities.
 Different may in 1-1 hand on figurities.
 Different may in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 hand on figurities.

Description of the man in 1-1 han

6.4 Benchmark in Kin 3-Th end could not list. EEMBC color. First oft may in the workload mix number and the second colors list the heachmarks exceed concerns the aspens of the multipregnation of workload.

• Encomy of medicalism to that different microarchitecture structures

. Reconfigura course conta blacks to a proble different fractionality in

different fractionality. Polymorphism refers to our figuring microarchi

lecture blocks to provide different functionality and in different from

Before applying this abstract definition of architectural polymorphic

to a processor architectures to decides the resources and medications for its

dressed: the granularity of processor cores, granularities of parallelism, and

The grand being of procession space the following spectrum chows in

1. Office the grained FP GAs which countries a first array of paint or con-

Egorable lookup tables interconnected through a configurable of work

1.2 System design questions

1.2.1 Gramiarity of Processors

technology scalability.

Figure 14.

are need differently at different times, rather than application-specific

B2 IPC companies with 6-dwg and 2-deep speculation - EEMB

Figure 12: Granularly of parallel processing elements on a chip. Number a care that can fit on a typical Glam chip.

1. Late of hade proceeding core like in Pipersock [8] or PACT-XPP [8]

The primitive processor elements are more powerful than gates and

le skar table: like is an FPGA. Because they are a marammed at a higher

they are programmed at a higher level of a lot motion than PP GA: they

1. Many simple in-order process on the in the BAW architecture [87, 149]

or Son Nagara chip \$85. Each processing core is a fall Bedged processor

that may application completed own to the ISA of the processor, BAW

also has the ability to me conhictionted compiler techniques to man a

4. Many a sweetly in the Purpler processors like in the Powers of its 10 th. The

Chapter 1

Introduction

agly diview application domains, producing distinct markets for deaktop, network, ceres, orient fire, graphics, and digital rignal processors. While clearly a period is a satellication - secific a effect size into proceeding, there are onbecause they are toxed to exploit specific types and granularities of paralleken, and to come extent due to instruction of co-cialization. Emerging an eligation, with hetero one can computational problems to such as impotables cycleme that may capper outs between computation. Prince systems can be betergeness at the hardware level and can be built using ropositik processors in copport this application had ampen They in Mer from 1900 problems: relaced economics of scale compared to a single general purpose design and design-time freezing of the processor make and composition. These two problems motivate themsed for a flexible or such marphana pencennar deniga i bar can adapt i a different application demando

ally spirate in the milli-watt regime, whereas server work hads a perat-

difference between different upp kentien demains, it skontd ake he noted that there choses of our correspy are not materially exclusive, in fact, 2 is common to extract come amount of LEP in conditional multichroaded work hade like detailing weaklends. As example of simultaneously using TEP and DEP : the IBM Cell processor, where madrith reading is extensively used to partition work among eight Synergictic Proceeding Englans which are SIMD execution to see of a amiliation. While ILP and TLP are well as denotes 4, the differences het were neverance with DEP intercepting and serviced. In charter 7 we as decision the behar to raf these programs

1.2.1 Technology Scalability

a destinative simpline.

creature, like region Mes, branch prediction table, and rename table, to extract concurrency. Al. increasing wire debry and the finite on nineline death from a performance and naver personality points; the scalability of these architectures \$1,74,72,77, \$42]. Consequently, restaulting distribution efficiency in processor architectures.

 Scalability and Modularity: The basic ideas habited polymorphism lead to the contraction of embble and reconfigurable modular blocks to repriet match tempolitation demains.

. Complexity: The engage of medical modular is extend to and her

. Energy efficiency: By ming a small set of mechanisms and adapting the processor is an application's needs, polymorph an architectures can be energy efficient for wide class of domain compared to general purpose programmable processors. Mowever, it is not clear low close polymorthen contemp can get to the enemy efficient of openialized agreement.

1.3 TRIPS Architecture

In this discontainer, we develop a technology scalable and ited are called TRIPS which not a new dataflew exceding IS A to express concurrency more efficiently to the hardware. The hardware is implemented using a distribute

data and works with only a named originate links and for communication

Architechral polymorphism provides the capability to configure hardmartin s

morphism from other approaches to reconfiguration like FPGAs which

a programmable architecture to application useds.

· Processor type: While polymery him does not require or imply a chip made of them opened to a more our organita this dissertation we gest girt our rek to discussing and realizating pulposaryhism for homogene The Smart Memorie chip is another example of a homogeneous polymarginary and heaters.

· Processor granularity: Artifictural pulposerphism can be imple mu ted as the grain corn the simple in-order processors or course grain over the the TROPS architecture. Designing polymorph on machasisms Braccerating the crain care to execute a large or crain order) (dif-Bres : challenger from our illesia ou cour e crais core for core or lan fin mis coomman.

· Configuration granularity: And testeral polymorphism is defined as configuration of course grain and examplification blocks and is (Bereat Romery it whites (Bereat Caucillas) from the grain yield the compa-nents like dates at littles, like and FPGA, or primitive processing ele-

In this diversaries, we discuss reference him in the context of the TRUPS are cover to server different contribution of condiction. The main

Table 44: A taxonomy of architecture

was on what calculater 5 to , the black some moins but a said the emichia memory

morphism below. Using polymorphism the reservation stations can be recon-Spaced in the following way to day! the processor to different granularities of rallelina: I configure the reservation stations like an instruction winds and denote all entries to one thread to oppose R.P. 2; those the expension cratical amount matriple threads for TEP, and I provide instruction to manufacture error at every ALU title to content the emis DLP that is been

1.5 Thesis Statement

This discention introduces the concept of architectural polymery him -the monthline to confine come unit migraphic state blocks to an-

presents the design and implementation of a scalable processor that can be configured to one out different on aniarities of carallelism prince of contrast on and walk also pulymorphs to medical mother topp or ing different granularities of parallelon on the TRIPS procesor.

1.6 Dissertation Contributions

Architectural Polymorphism: We him do so the sun cap of archited ara a alternatively man of develop the main a placial or and a vet of mechanisms of time by there yets signer that configure course grain microarchitesture blocks to appoint different grant lattice of parallella m. Compared to reconfigurable as chitectures which attempt to provide capper for diverse workloads using a seatherly as arough of he filled different front land blacks from a resilience inciple behind polymorphism is to adapt marce grain blocks to

TRUPS Architecture: Wedgeschie the TRUPS among a constitution, in IS A just include of an EDGE IS Aj, and interpretable one'. EDGE IS

*The principle behind EDGE ISAs and the implementation of the TRPS ISA and it is the two contribution of the contribution of t

the wire-delay and complexity to see that a hope realists of the and efficient micourts lived are. As a result, large amounts of DLP will have to partitioned into threads and distributed across a set of narrow-issue cares. TRIPS and other hand provides a scalable very wide is see design that can be tailored to application needs using pulpmarphism.

Finally, But hat et al. introduced a nervatibly metric to quantity the ability of an architecture to effectively execute a broad set of applications [22]. that captured directs behavior. This rematikiy metric is simply a quantitati metric for comparing different types of architectures and does not describe as characterize the architecture itself. Then formally define yen stills as: 'the geometric mean of the speed up of each of the applications in the Versallean on the relative to the architecture which provides the heat execution time for

Extensions to conventional designs: Is addition to reconfiguration for performance, adaptivity has been need to increase energy efficiency. Although et al. El introduce admitre proposine where sa-chie structures are disagnically recited to provide power efficient execution. This can be thought of any adjunction within the LLP domain that a section time up plication behavior. ier is improve margy efficiency. Other examples of specific microarchitectors mechanism to emelde adaptability include the following: ad in the cacke size

age which is configured as a data cache. The Piranka architecture explore illed architectures targeted at server workloads and took an extremal position for the time \$71. It interruted eight very simple core along with a complace much blemucky, memory controllers, coherence bundware, and network controller, all on a single skip built using ASIC LISper technology. Association of the state of the tiled architecture that ones ham apeneous tiles in Smart Memories [Mil]. The Stackrescalar [132] and A-AP [353] architectures are other example of he at DSP upplications. Emerging fine-grain of CMP architectures, onth at Sun' Nic para [65, 52] or ISM's Cell [63], on also be eleved as tiled architectur. care [M2], Picchip [4], Clean year [63], and Silica Mire [64] many of which are recieved here [63].

ner tile. In ceneral, these tiled architectures are interconnected at the memare interfered, although RAW allows register-based inter-processor commucation. TRIPS differ in two ways: [4] different types of tile are compared t create an approximation of [2] TRIPS over distributed control network protects to implement for client that would at herwise be controlled in a convention a

and consumed within the block. Which temporaries can be forwarded direc from producers to concesses, without mer being written back to any central comps. The dataflew graph is someted in the block through instructions in traction communication of these black tomporaries. Black outputs, however, must be written to a central stampe like a register file whom the black committe. The black enter to a few chiledy and the black involved it to receive create the dataflewages for the entire program. They are at a final of transfer instructions which specify the address of the succeeding block are also t as block outputs. Medification to memory are maintained in a temporary compensed the black is committed.

1.2.1 Block Execution

The compiler calkally acign each increasing in a block to one of th named ALU days. Each ALU can have multiple in the clien days as occupied with it. Special read instruction, used to read block hypot, are assigned to the region file. Expension of an instruction block proceeds as follows: λ black is the checked and mapped onto the ALCs in the execution orbitals at sace. Each is direction in the black is started in the instruction slot at the ALU | imilar to a recommend of station | to which it was stationary as igned. The read instructions in sed at the regions file, and block inputs and tripper the duration execution by injecting the values to appropriate ALUs.

When all of an increasion's operate have a mixed at an ALU, the

or educity express concurrency to the hardware by someting programs as senament of atomic blocks of executing with blocks on ordine a dataflew entry processor care procides a MIN-eatry instruction window and can been up t M instructions every cycle. We have also halfs a prototype skip in Minn ASIC technology companied of two TRIPS are enter companied a distrib IMB consistent [NUCA] social memory system.

Data-lewi Program Attributes. We arried a detailed characterization of the fundamental behavior of data-parallel yragman based on their memory accompathers, program control behavior, and available concurrency.

Experimental Evaluation: Our performance reaches how that the TROPS microarch de are un cucials good later classical parallelium. On a set of has deplimized benedic IPCs in the range of 4 in 6 are seen, and as a set of highly data parallel bundmarks with compiler puremed orde IPCs is the many of 1 to 6 are two. On the EEMBC and SPEC CPUBLISH exchangly w we expect these sumbers to improve.

The polymorphus mechanism proposed in this discriminance effec-ce at exploiting thread level parallelism and data-level parallelism. When executing 4 threads on a single processor, significantly higher levels of proces or eliterial awares 1900 awas the reason of \$7 to \$4 for an applicable

tia wayo Fil, tizing a consensations (64), adjusting the too the load interest or earliest her the 1996, ad not be to see with about with the from the part of the first of t

At a marrier granularity, single-15A beterogeneous processor attempt is provide copport for different granularities of parallellom by integrating ma tiples types of corns which all ose the came ISA [64]. In a cimilar cela, Kuma et al. discour the architectoral traductival abutang varying degrees of hard war her were processors and threads in a SMT/CMP by held design to explore the

Course-grained Reconfigurable architectures: Fisher et al. propose Cartons-II processes where processes core are spatherized at design time memory sizes and hierarchy, number and types of fined hand suits on a all h actionized. There et a harden et lafterare exclusive a recession also income optimized for that applies lies is presented. The final processor is fully penera common and margin all other are likely in allow, all of any an efficiently as emplete too khala flow for 17 ath editing yet on our and an 15 A, harel on a set of applications [256].

2.4 Microarchitecture techniques for ILP

We conlead this literal are retieve by discussing work related t ing in tirection level carallelism. The dataflew execution model and tradable techniques for extra clair ILP are the month dought related areas.

Databor: The execution model and ISA design for the TRUPS appropriate heavily impired by prior dutaflow computers. Dennis and Missian proposed italic data flow architecture in their cominal paper on dataflow computing [88] ance data tokens could not be you dured by an instruction until the produced by it during a precious dynamic instance were consumed. As a result the levels of concurrency that can achieved by an erhapsian matricle items in of a loop is Smith. Dynamic dataflew addresses this problem by dynamic cally labeling dataflew are und managing there in a back table of datafle tak mar [13]. Can tinning this work on dynamic data flow Arris d and Nikhily re named a Tangeri-Taken Dataffew architecture with sometridate-driven in time tion inheduling the programs expressed in a data three language [2]. Call et al., later proposed a by held dataflew execution model where program car partitioned into code blocks made up of instruction sequences, called threads with detailment award in human threads [16]. The rich his are of detailments directors to review by Arrivé and Coller [15]. The TROPS approach diff. execution for a limit of window of instructions, and rely on compiler instruction

restitional taraffew machine [2,38]. When the intraction compl must be forwarded to the ALCs holding constraint interactions, and/or to the register He if the result is a black on type.

Operands are delivered directly from a reducer to consumers in sint to might in the AEC are well maker than being his already to a CAEC. As a conventional and become, which require complex bypa tor were ALU, a simple point to equippe a count will reflect the EDGE and the tures. Since all operands are forwarded to the location where instructions or ha Great, as in traction does not encode the correctionalists arregister name of its inputs, only its aniputs. The physical declinations of the instancion result are easeded explicitly into an instruction.

When all of the interestings in a black have completed, the black is committed. Black outputs are written back to the register the and applace to memory are carried only. Subsequently, the black is removed from the ALCs, and the next block is mapped outside execution colorests. In the erest of as exception being mated by any instruction is a block, the entire black is re-executed after the the exception is certical. Similar to pipuline the onlinequent block with the execution of the current block. With this typ of everlay, multiple blocks can be in Hight simultaneously and the ALCs in mix consisting of REMBC and SPEC CPU2001 workloads. When executing starrams with DIP, the risk man have mechanisms we arrange a maide har sicmens operation of 2.1X across a set of DLP workloads, comention may delusf extension goodly ILP. Compared to apecialized archi there mechanisms provide competitive performance using a single execution

1.7 Dissertation Organization

Therest of this dissertation is organized as follows. Chapter 2 discusses minimi work and along this discentializate the material or print work. Chapter I describes the TRIPS architecture and the emission of RIPS data. We e the main features of EDGE ISAs, describe the microardide the THOPS thip and briefly describe the logic design, resilication, synthesis and physical design of the prototype TRIPS chip.

Chapter 4 describes architectural polymorphism. We describe the three principle heliad out mornion and a classification scheme for processor reto arrest in to fixed, one childred, and not reported our resources. We then describe vi and recorrected to implement polymorphism to copport HEP, TEP, and DEP.

Chapter Specials a performance evaluation of the THPS processor ficusted up in struction-level parallelism. The performance evaluation is based is an erest defree radicated processes charlator. Chapter 6 presents a perform ance evaluation of extracting throat-level variableion of

medianisms in the TRIPS processes.

Chapter? precents a detailed application characterization of data paral-

let program chared on their findamental behavior. Haved on this characterist

proposed. This chapter also includes a performance evaluation of these mes

anicus on a high lend processor simulator that us dels the TRIPS processo

Finally, that for 5 concludes and a plate to come from directions in the coff-

processors, circum processing and other by held architectures. The key differ

sace between many of these architectures and polymorphism is the ability to

to so out different cross brities of so miletims as dithe crossfurity of reconfic

Vector processors: Early data parallel architectures were chinic testor

processors which were held union security \$10.500 for his broad management

and large vector register thes \$5, 305, 30 \$1. These markins were designed t

programs with regular control and data behavior, but much indicate some di-

one of inecular has simplered memory accesses using scatter and cathe

es malica e. Processos with frequent impendar memory references or account

to history taking performed pointy. A number of architectures have been

proposed or built to overcome the limitation of the right vector execution

medeland to allow for denamic instruction scheduling and conditional execu

tectures widely applicable as they provided support only for a cohect of data

different nector have and specification of the nector have readers sequentia

and a sa-rectarizable code nery in efficient on this architecture. Short ned or

proceeding has found to very him commercial interpresentation to the firm a interaction extends as the fact of MMX, 55 E2, All free and VIS. These architects

parallel programs. The Verter BRAM architecture is another restor process

to aptimize other technology constraints like power and new

et i of polymorphosis systems and the application of polymorphic

printing which does detailible computing in the empty [8]. Vertication techniques are used to praemite configuration states for this arms for large blacks of most hire colds. One of the drawbacks in the architecture is the lacof copport the executing requestial programs efficiently and lack of access t mademation memory. The Michelan \$7 processor belongs to a new class of thips called Field Programmatic Object Army FPOA), in which, instead of configuration of gates like an FPGA, designers work with a massively parallel army of process figured fraction scales like 16-bit ALUs, muchiply-accounts to sailt, and register the which can communicate through an interconnect fab

ric. Marientela has written a literature carego af other reconfigurable grade et architecture carpetel at a single application domain [68,70]. is the ASH and the are, the analogs in madel and dataflow concerare limiter to the TRIPS approach \$7}. The main difference being that, ASH targets applicable as pecific hardware for small programs, as appliced to compliant large on smart late a resolution of configuration may sed to a long grammable in hitmiss. The Gary architecture and the BRASS project as of an PPGA has alseconfiguration approach to affined compute in motive regions of as application to an exchip FPGA [2].

2.2 Data parallel architectures

Several and here have proposed architectures and marchanisms for data illel ambitectume. In this cention we dismost be work most desidy related to some group of and or restory more one, syctolic arrays, \$1 MD/MIMD

this dissertation. Nagarajan presents a more detailed somey of approaches to tar is an expension burt.

constraints that have led architects to build scalable and modular designs. Promore parallelism, and openifically instruction level parallelism. Extracting IL. creater three requirements for processor architectures: It a large window of niefnig regram in it met is no. 2) a realable s execute a large somber in the client concurrently, and 2 | a high hand width and lewistency memory system. Rangan ath an and Frank in described an empirical study of decembras

ILP: Technology limits of power, design complexity, and wire delays are the

ired ILP execution models [124]. Subject at proposed Multicodor procesone, in which a charle prompt is broken an into a collection of colorability tails [64]. A different approach to counting a distributed with the number races for the execution partitions [258], in that work, Yajapeyana and Mira proposed renaming temporary registers within a trace to reduce the Small propored the HLEP architecture where a distributed microarch using PIPO has all lastered has loose queen execute lastered has which have has broken into strands of descendent instructions [87]

Other correction with effort than ethic (LP) are focused on large-window samildies he mean efekerkeels like and over lating \$5,043, hebrid datafles

1.2.2 Key Advantages

The block-promit model will be effect to fithe number of incinitions is the block is large except to civil large design desce chains that may be well from the ALU chaining in the grid. The experimental results in Chapter 5 there that compiler-yearment black sizes are significant, when prelication is ared to eliminate on tradition hazards.

When we started this research we performed several empirical studies is that the dy, we seed the Trimoma compiler infrastructure [20] to measure the amounties of blocks that are important for EDGE IS As; of the size of blocks, b) number of block in mis, of number of block on in mis, d) number of block (emponente), as de financi of block temponetes. Our failful evaluation indicated that existing programs were well saited for this architecture. Typical black sizes maked from 27 to 125 decamination executed in the cities, which are refficiently large to amortize the deling or exheads. The number of layet and nation to pair of first large fraction of the blocks was lost than M in more of the beach marks, is dicating that the amount of register file communimaps in dicates that a substantial amount of communication to the controlled as block temporaries, the grid. Finally, the average number of our comes of produced radical cally 33, which chows that the autwork within the execution is hitter to do so not require harpe has dwill in first in making our manifestion

This execution model addresses several of the dialogue for micro processor performance scaling. In particular, an implement register recursing table and there are fower register the read and writes. Deup the like lack of these structures, instructions can execute in an order deterhaved no sa true data den sa dences, with sat expensive have r checking or a broad carring by parring and forwarding network. Palachada e al. demonstrated that broadcast bypass networks scale possity and typically of saidily networks, and propose Routed Inter-ALC networks. RIANs in a scalable communication and work for factors on occurs (1931).

The explicit concurrency expressed in the ISA, and static mapping of in traction to records an endy allow for a colable and made by microsedifference in olemon tarion. Purchemone, if the obtained instruction have in men will take place along the et, point to equint wires. In tractions off of the ALU. The physical layout of ALU is expected to the instruction schedule

imize the critical path. Other publications extendinely characterize and and section of stating peaties \$2, 28, 28,

Related Work

to the form of this divine ation. The related work is proposed around the four

main them et of this discentialist palymorphism, data parallel and scalable and design, and microarchitecture techniques for ILP.

carryrine work that has focused an uppy of for different type of a on a single solutente oring recordinate or other means.

2.1 Polymorphism

This shapper discusses and differentiates prior work must decide related

Below we discuss the previous work related to pulymorphism. We dis-

Multithreading. While multithread is not directly related to an exercise

different trace of applications, submornation like hebation has been used

rapport malithreading is modern processes. We briefly trace the bis

of multithreading before describing these systems. Multithreading has been

Malifi heraded p is elia ing was used in the Peripheral and Control Process

of the Control Data 6600 computer architecture of the early 1860 to provide

for their medical country of all a phicometric forms of the IS A extension a contract

 $\log \pi \pi \sin \pi \sin \theta$ -word of a MMX/SSE2 regions whose using a confer regions as

processing elements are interest sected in a pipelinel manner, with such el-

ement performing the same operation is negations; and passing along the

processed data to its a sinkhors FSL Prior to this formal definition and to se-

ilkatina of cyclolic armys, the British Colores computer employed an arch i

tecture similar to contails arrays for code breaking [33]. In general, syctolic

armen have estimatible been used to build special extreme annilousing specific

hardware [29]. The Warp machine and a spiritic army to construct a pro-grammable data parallel architecture to copport calculation computing and eig-

and proceeding applications [30]. The iWarp architecture excended the design

of the Whey muchine, by designing on 1986ay block that could be replicated

of a procenting core and a communication upon twich archeolrated the com-

maximatica between different (When chies, The (When architecture was also

targeted at scientific and image processing applications, executing parallel programs on a large (Warp sprima consisting of many (Warp shicks, using a

Chapter 3

TRIPS Architecture

MAKE and SS E2, have noncompact for scalar-vector operations, unit or

wilds and tackers compute monarce between multiple program threads [67].

3.3 The TRIPS IS A

The TRIPS ISA is an example of an EDGE architecture, which aggregates up to US instructions have a single block that whey the block-number execution model, meaning that ablack is logically fricked, executed, and committed as a classe exists. While details of the TRIPS ISA maybe found in

1 1 1 THIPS Blocks

Early THIRS black manifes of 195 hors land one for early of the annual le-25 in irractions. The complete constructs blocks and using as each instruction to a location. Each block is compared of the west two and five 125-byte chanks by the microarch here are. As there in Picture L.L. every hind, include a header chank which on codes no to 32 read and no to 32 price instruction (that access the MS architectural registers. The read instructions pull values and afthe register and send them to compute instruction in the block, whereas the write lastered load return only all from the block to the specified and ited and instructions are distributed across the fluoresignier banks, as described in the sem rection.

second circuit peripheral processors [201]. More recently, the MEP makiprocessor states had builted notes and the number of processes dynamically in order to take advantage of carrying amount of parallelism is a problem [23]." Other recent options that prorided mad threading rapport on a ringle YLSI chip in chide Mydra [65], th MIT M-Mada is \$11, MIT Alerife machine \$1, and the Pires to making reces-

Fire grain much librarding to chare processor recourses be has been explored uping different techniques. The Term comparies or stem had special for the coming multitheresting interference LIW instructions from 4 F ferral (herad) every cycle [6]. Reckler and Dally proposed an arch that incorporated both compile time and ron time information to interfere matriple VLIW in traction can indicate at fractional and c [64]. Both of these Tallies et al. described their approach of supporting multiple thread context mati à breut lag [25]. They you need as lan waster moths dof replicating cor-tain architectural strange elements in the your count, but charlen most at he measures to consider the execution of multiple throads simultaneously in the processor simpling. Yam amoto and Nemiros sky a manused an arth itertures imifor to SMT had with a grand a instruction quarter for each thread [264]. Unperson

hybrid multithreading and systellic processing model.

SIMB/MIMD processors: The SIMD and MIMD term was mined by Fire in the town over of manual or problems and \$10. The early fire-over SIMD machine like the CM-2 [04] and MacPar MP-1 [04] provided high AL per he memore accesses. Modern a meramma ble cras bles a messon consis of a negrowide SLMD execution on size to perform fractional and neglectors ing \$1. AND ADD architectures have typically have used to build hope scale parallel architectures. Other examples in the degraphics pipelines [8] and rides processing [84]. The Universitä architecture is a the grain AUAND architecture that one register changes to communicate between independent processing units, and he makes there changed citible to the compiler allows that he tween the index on deat of reason [60]. The one of register this area in this arch is Proceeding architecture \$7). The most precalest are of the grain MIMID preen sing is in mardem graph in processors which contain vertex shadons that are MIMO and Declare [82, 82].

Stream processors: Stream processing, which has similarlike to rector procentry and SIMD computation, is being explored in several architectures turgeted at multimed in proceeding. The stream proceeding paradigm is based on

The THIP Samb Heriare defines a new instruction set architecture calle Explicit Data Graph Execution | EDGE| which allows dependence to be explicitly eare ded in the instructions. The microarchitecture is heartly partilieu ed and non-well defined communication and works to build a larger compe gradued processors [aks known as Grid Processors] to achieve high perfe mance on ringle-threaded applications with high IEP. These cares are anymental with sales as has a factored that each letter or mailer or may interest. of parties the core for explicitly consumers to applications at different companied that are difficult to reak, the TRIPS architecture is howelfy partitle and to avoid these large contralized structures and long wire runs. These

point communication channels that are exposed to coffware schedulers for

The TEOPS architecture is constructed of modelar blocks and hence specifies a conditional to disc for exclude a neb magnificat. The key distleane in defining refermentant from the TRIPS is to belong their arms

DLP can maximize their new of the available recovered, and at the came time are id would be come lead to and a pay calculate structures. The TRIPS statem and instruction storage, to minimize both software and hardware complex and configuration acceptants. The remainder of this chapter first describes EDGE SAC and the execution madel for such SAs. We then describe the TRIPS intraction or which a new intraction of EIRIE ISAs, the TRIPS processor microarchitecture, an overview of polymorphism, and a description of the prototype TROPS chip.

Explicit Data Graph Execution |EDGE| architectures allow compilerprocessed distribute graphs to be mapped to an execution cabilitate. The two defining features of an EDGE ISA are:

1. Block-stem is execution.

* Emvious de taffer-like execution enabled by direct in traction to destruct in found assessed dayer of advisaction to express one correctly to the hardware

Support for Polymorphism: We are this architect and in so of all after exceding in the ISA to exploit different grant lattice of parallelon efficiently

management affect the could billy of hardware and limit achievable perfe manys. The hinds attended to marries there are bendered acress many last malings and expresses concurrence efficiently to the hardware ing over a large number of instructions and reduce the number of branch

et al. provide a detailed corner of mobilities of ing liberators [187].

arable army that reled as a lot a Behip communication.

Novel architectures: Brown et al. developed the Texas Becomfigurable

proceeding on a single-rationals [47, 68]. The TRAC project was free out to

ha lifting intercensection networks and a stimulator or manages in the first on the

approach works well for thread level and data level parallelism, single threaded

execution coffee on this architecture. The main concepts of difference between

Sm am Memories and TROP'S is that TROP'S has a well defined set of sylectalized resources that can be need to support specific application

needs. For example, TROPS has a traditional 2-way set an achilice instruction

tion freck. In function dues not change with application behavior. A occupi example is the exit predictor need in TWPs, which is need to predict contro

flow for respectful programs, in Smart Memories on the other hand, there are

as such fixed measures like the last motion cache or concluded measures like

the exit predictor, factored the architecture straptly provides an array of tiles

with each tile containing multiple SRAM hashs, as interconnection are work

and a cimple processor care, 5 parketizing efficient instruction suche behavior

which concerns and a produce of grame of data, while concentrate through these

bened fractions. There kennel fractions are in turn assisted to such demon

is the circum. Imagine dabbed a circum processor, is a \$1MD/rector bybeid

using a \$1 MD control unit complet with a memory upstem recembling a vector

machine [127]. Other saids MIMD and hed are such as Mercinac and BAW

yek and programming language techniques \$7,57}. The Brook programmin

larguage provides copport for stream computation on graphics hardware [26].

Mybrid architectures: Recent amounts have concerted combining vector

computation on its with modern and aborder processors. The Toronta hards i

tecture uner a betem process; computation up prouds and integrates a 22 will)

rector core and a high northman constructor EVS core to turn of data loss

rector mandel of execution with global synchmolimities between the differen

nector laws with partitioned nector registers and optimized accounts to the

regular L2 code for vector loads. The designers went to great lengths to provide the high bandwidth required out of the L2 code with an innumitie

conflict-free a differe peneration estem et a maximize the number of concurren

accesses to different cache hanks for many traces of studed accesses [226]. Pa-

juelo et a L proposed operatories dynamics exterimitien to which sectorimble and exegutes to are detected in request in London, and are operatories presente

The dutaffer exceding is efficient at expressing ILP, TLP, and DLP, Per ex-

tracting ILP, the dataflew exceding express or the limit of parallelism in this is:

on all regions of a program, directly to the hardware. The hardware not con-

dependence of sching. Perest fact ing T.E.P., the dat affew exceding expresses the

limited namifelism in such thread, and the hardware can interience multin

ing multiple instructions from different thread contexts. The dataflew graph

abstraction amortizes the overheads of instruction management across overs

instruction in a full black of instructions. For extracting DLP, the datables graph abstraction directly expenses then bus data parallelism to the hardware

-typically the graphs are very large when programs have data level paral

lekem. In conventional RISC and CISC ISAs which require the hardware to

amildion and instruction level parallelism \$5]. Tamatch pr

The Stanford Smart Memories employs any led with marriags much

the a large core from a modelar obstrate [10.0]. While this

1.3.2 Direct Instruction-instruction communication

register like het west the producer and consumer.

Direct is dractical communication, in which induced in a block and

At those in Figure 3.2, the TRIPS ISA organic direct instruction

their or emails directly to consumer instructions within the same block in a

dutaffier facking, permits distributed execution by eliminating the need for

any interesting chared, contradered constitutes such as an increwed downer of

the consumer relies and forward a produced operand directly to its targe

instruction in . The placehit turner fields of I and T is shown in the encoding

each rescribe the secretary tree Seft, gight, predicated with two bits and the

larget hatroction with the remaining seems. A microarchitecture apporting

this ISA will determine where each of a block's US instructions is mapped

thereby determining the distributed flow of an erand values the data flow crash

A record aspect of the instruction or ording to phosment. While, the

while each block. As instruction's number is implicitly determined by i position in the church shown in Figure 3.3.

behavior out of the memory (for a almost impossible, White more home posees and perhaps simpler than the TEUPS design, the lask of any operations

coupled processor cares \$15. This architecture provides an innovative way a hadding a realishly, tightly integrated MLMD agray for data integrity street, ing. Clearly this architecture was excellent vector order and the grain MIMD parallelion. This architecture lacks many mechanisms that are required for extractice (LP, for example, it is do not only ordering medications for leady torm-ordering. As a mould it is no dear how well this architecture will a scheme

So contract of an annual and professions called \$1.0 to consum the TEP, and DEP for media applications [332]. They immediate a DEP such signs called SIMD section and creams [SVector/SStreams], which is into crated within a conventional constructor based CMP/SMF architecture with ink-word SIMD parallelian. The technique is able to explain the simple in plements that of oil-word SIMD already common in many markins and it provides the headin of full Hedged vector processing. The primary focus of AIP is to connect multiple trace of numberior an operational architecture or any changes to the 15 A and microarchitecture. The main draw tack of ALP to that it augment on conventional processor core and as a resulit does not reade to large into e width s. The techniques proposed in ALP extend

as a dedicated exclusivation and [136]. This architecture is also beterapes ma date it a period two dedicated datas athis.

fatefacily is an embedded percent or that in dodes a high performs cular MI P532 cominingment with a narmy haved pure fellows or much unit [13] The exclor math said could need as array of ALU; examened in each other oring a high handwidth in ter-ALC are work field by a high handwidth L2 code The L2 cache can cortain a handwidth of 64 Chyles/sec, when remains at Title. The interesting control is the army in order SIMD with med. HEL eth er example of a hybrid and it ed are that in do dec an on to Sorder you con a and an in eight SIMD amontum, do thing committee processor engine, with and were managed memory over \$1.50. The out-of-order an entra manager memory for the SPEs and is used to program DMA engines that enterin is DRAM is eastly memory transfer.

2.3 Scalable architectures

With transister course approaching one billion, thed architecture arraping as an approach to manage design complexity. The RAW architecture pleasured research into many of the house facing tiled architecture such as ween the tiles, in the clien scheduling notes tille, and efficient memory notes action (186) \$8, 507 and , 500]. In the RAW and Declare, all lifes are it write. and includes preserver sers, a mater, memory ordering logic, and data stor-

mort bemade juck as first or commit j, prociding latency tolerance to make distributed execution a marking. The case then et al. or entity the branch are

3.2 Execution Model The execution model for EDGE (SA) (res) and look of instructions as as

stamicant for firsting, executing, and committing. The execution substimia collection of ALCs, each of which is architecturally visible and named ficity, we assume that all ALUs are homogeneous and can execute any instruction.

Block-atomic execution: In the block-atomic execution model, have no inar are placed him blocks by the compiler. Blocks may high depredicate increasions for horse as internal transfer of control; taken branches (in 4 th) last in the clien is a black) transfer control to a secreting black. A black confi that he about block, aprecious of hyperblock [80], or a real-line trace [80]

Dataflow graph abstraction: The ISA allow the dataflow graph of exemilian to be directly exceeded in the blocks. The data and and concerned by a block are of three types: |1| black inputs, which are calculated by [2] block on the by which are rather created with in the block and need by on b-

March States Company OPCION CONT TO C

(1 to 1) =

justed all, describe the other appeals of this planement problem and introduce ology of classifying architectures based on when jointic or dynamic in sed [188]. Marger et al. [18] clandit other architectures according to this

Other assemblished elements of the SA lichde the 'PR' field which receipes whether such instruction is predicated an an incoming ten or false predicate, and the load/stare identifier [LSID] field, which specifies requestion order in which loads and stares must execute. The TRIP If A manual contains a complete description of the instruction set und bec-

Instruction Charlet Instruction Charle 1 Instruction Charakil Ingenerion Charles Cinaturations Figure 1.2: TRIPS Block Femal The header doubt also holds there trace of coarms state for the block:

a 32-bit force mark that indicate which of the provide 32 memory in the tion care core, block execution Hage that indicate the execution made of the block, and the number of instruction "budy" chanks in the block. The store mark is used for distributed detection of block completion.

A thick may contain up to four help cheats—and contolling of 12 induction—he a maximum of 25 induction, at most 22 of which can be lead and curve. In addition, all provide execution of a given block must always out the same number outputs (times, register writes, and onbranchi mearifest of the amilitated such takes through the block. This conchains in accounts to detect black completion on the distributed of The compiler is may as tible for perenting blocks that conform to these con-

1-12 targets simply create the linkages, the anderlying processor microarchi tecture is exposed to the compiler to it can generate efficient placement, with the real of minimizer communication distance among instructions. Natura

3.4 TRIPS Microarchitecture Principles

The goal of the TROP Smikroarchitecture is to achieve high concurrency whether ILP, T.LP, or DLP, on a technology-regulation distributed care. Our definition of a colable and share in the iran processor that has no global wires in halfs from small set of record components sitting on record networks, and can be extended to a wider-kinne implementation without recompling ourse

The three spacegable principles habited that tiple of microsophisectors Modularity: The microarchitecture is constructed with a small set of tiles

replicated and connected to gether at necessary. Tiled enture: The microscolitecture is physically purished and tiled in sature. The logical organization of the tiles has a physically tiled of gasiration or well. The tiled nature allowed a hierarchical design flow at all stages of the design a design specification through RTL coding conflication, and physical during. Made being refer simply to the legiall construction of the architecture through a small set of and a tilent refer to a physically regular placement and interesonection among the

Interconnection networks: The tiles [modules] communicate through well

predicated branches. The block predictor as a branch instruction's three-b exit field to construct exit his tories in their of the tax takes /ant-taken hits. The predictor has two major early; as easy predictor and a target predictor. The playing a termament heat/graine predictor similar to the Alpha 2 224 [63] with SK, SK, and EK hit in the local, clobal, and to an amont self-resident

When the extraormer is analytical it is combined with the analytical the first address to access the target predictor to predict the next-block address.

The target predictor contains four major structures: a branck target business. a his ack tracered in or 120% bits. The BTB agolich tamen for branches, th CTB for calls and the BAS for returns. The branch type greater or predict

into banks, with one bank in such RT. Like the other tile, register banks are and on an the DPN is leading the case other to a law instructions that read an

Since many defense rain of instruction are converted to intro-block register has dwidth requirements are reduced by approximately 7.0%, as arenage, compared to a MISC or CISC process or. The floor distributed banks on

Controlled register the contepower and delay problems is large, and

Picage 3.3: T 80 P5 Providence Chin Schematic

MA SOC ESC N I G R R R

N M M N I D E E E E E N M M N I D E E E E

N M M N I D E E E E

SDRAM 1 C2C (x4)

that provide to Mickel register handwidth with a small number of parts; in the TRUPS agreement such RT hash has two read and and one write and Six or the THEP 5 IS A transition HES architectural recition, each of the four HT contains one 32-register hank for each of the four SAIT shreads that the our appears for a total of MS regulers per RT.

RT contains two other major structures: a read given and a write given as shown in Figure 2.6. These queens contain the eight read and eight werd : ared to forward register writer dynamics by to collection to block reading from there existen. The real and write comes perform an emphasis fraction to resister manuals of the a physical resister file is a supercuber resonant, but les complex is implement due to the ISA support for read and webintenction.

1.5.4 Execution Tile |ET|

As shown in Figure 23, each of the SEETs manifes of a fairly standard ringle-irrae pipeline, a hank of 64 more at innotations, an integer unit, and a fleating-ratio and . All pairs are fully simplified except for the interest in the na it, which takes 2 depoins. The 64 mornal has stations hold eight last medians for each of the eight in Hight THEP 5 blocks. Each mornal has station has field for two \$4-bit spens to data open ato and a sale bit predicate.

are the motion of the Complete many of the traditional motion ording techniques, and as the central foliable implements a microarchitecture Anches Sal [-G-R-R-R protocols and [225]. Taylor et al. describe another taxon only for classifying h in a a a a the network from course to decidation [845]. 100000 1-0-0-0-0 Each of the processer cape is hardemented using five up has tilled up.

Global control natwork (GCN)

0-0-0-0

global coal mid like |GT| . We execution time |ET| , four regions time |RT| , four data (dec |DT|), and the instruction (the |IT|). The major processor core micreate work to the operand a stwork [or OPN], thowa in Figure 2.4. It connects all the tile except for the IT; is a two-dimensional, wormhole-mated, but mesh topology. The OPN has reported control and data channels, and can deliner use \$4-bit data up mand p m link per cyclic a control header packet is

Each processor core contains six other micron elevate as described in Table 3.4 Links in each of these networks connect only nearest neighbors networks in Figure 3.4 and discuss their stage later in this section.

The particular arms general of the that we implemented in the protetype produces a core with. Newske and of order into, \$480 of L limitacion

cark e, JIKB of L1 is in cacke, and 4 SMT threads. The makes architecture cap

tire retwirk ESN

parts up to eight TRIPS blocks in High trianglian mostly, seven of them specstatics. En place thread is manian, or two blocks are thread if four threads \$424 instructions.

The two processors on the chip have independent microsofworks. To commutate, they must go through the secondary memory system, in which the On-Chip Network |OCN| is embedded. The OCN is a 4xH, we mission mated mesh network, with M-hyte data links and four circuit channels. The network to splitssland. Our cack older cland transfers (and housing acket th Bowel by their Midyle data packets), although at her request class are capped of the so grations like leads and other to an entirely maper. The OCN acts as the transport fabric for all interests on on an L2 codes DRAM, U.O. and DMA inte-

In the cort of this section, we describe the contents of each processor

a given Aff . By a lighting the OCN with the DT s, each IT/DT pair has its own

spirate and late the occupier memory colons, connecting high handwidth

into the care for streaming arcting inns. The Network Tiles [NTs] carroand

ing the mamney cyclem and an translation agents for determining where to

mate memory system requests. Each of them contains a programmable mat ing table that determines the destination of such memory system request. By in the grade mapping functions within the TLBs and the network intenfa

of ways including as a single IMB chared level-2 cache, as two independent

directation is on the processor case, we refer the moder to [66] form on de-

table on the cache organization, and [88] for details on the TRIPS On-Chin

Network. The other six (fee on a skir's OCN are I/O clients, namely two

SDRAM cut miles, two DMA cut miles, are Chip-to-Chip on to Ser, and

one external has one to Ber that one interface to a PowerPC 448GP chip, which

LUND level-1 cycles in a parameter of as a 1948 and in about a manu-



Figure 18: TRIPS Tile-level Discount: Resister Tile - RI core tile, and then in Section 2.6, those how clobal energines among the tileis it as think and committees implemented by distributed microard iterior

1.5.1 Global Control Tile [GT]

ore 13, it holds the black PC and handles all TRIPS black management prediction. Dock . disputch. completion detection . Dock is a miss rediction; and interrupte) and community it who holds the control registers that configure the

processor into different speculation, execution, and threading mades. Thus the GT interacts with all of the material networks, as well as the GPN for madian and writing the black PC. The major circulars in the GT are th increasing cache ing arrays, the increasing TLB, and the next-block prediter [225].

on e of the black that are free, the GT access on the black predictor, which takes three cycles and emits the predicted turnet address of the next black. Each

| DECTION AND SEPTEMBER | DECTION A SERVICE OF SERVICE AND ADDRESS OF THE S type of the branch currently being predicted $|m\,I|/m\, ara/b\, mach/sequential-$

branch). The type predictor is necessary because of the architecture's distributed for the most scoke the arrelled or sever sees the acts of broady instruction tings they are two tiles to from the IT; to the ET;, to the branch type may

earted properties and promes deadlock or oldance and scalability prop-

The TRIPS prototype chip implements an EDGE ISA called the TRIPS

IS A. In the following paragraphs we describe the microarchitecture of this

proletype. The three major companions the skip are two present on and the occupany memory system. The processor cores accopy the top- and bottom

right quadrants of the ship, and the saich is memory system accupies the lef

half of the diff. Each operator care is a Movide inner TRUPS care that

cas have up to M24 instructions in Might. The secondary memory system includes a set of the that are configured to firm a NCC A cashe, two integrated

SDRAM controller, a DMA controller, two disposed by [C2C] on to ben that

The tile is the promour core and the tiles in the na-skip network or

EBC | that is used to interface to a PowerP C chip.

projected this. Figure 3.3 charm the tile-level black diagram of the TRIPS

As a rest in the above principle, this microscoline is one possible, permitting different numbers and topologics of title in new implementation with only materials disagree to the title lagin, and no change in the collection.

Lill-make. Since each TRIPS black concerns as many as \$40 bytes' worth of instructions, the migrouph become breaks blocks into the MS-instruction hanks, cashing such chunk in one respective IT. Each MKO IT hank co

2-way, SKB hash of the 22KB Lithau cache, as shown in Figure 28 . Virtual address or are interleased arms; the Delile at the granularity of the Delile 64B code back to addition to the Libraria back, each DT contains a copy of the land/store given $\beta 5\,Q$, a dependence predictor, a surveying back-side coalering write buffer, a data TLB, and a MSHR that can repport up to M

Figure 1 4: TRIPS Margar elevate (GRD, DEN, and ESN and above

Because the DTs are distributed in the network, we implemented a mann man si deden su deuce a prolicia prolicie propied with each data cache hank Loads ince from the ETs, and a dependence prediction accors in parallel with the cashe access only when the hard arrive at the DT. The dependence predictor is such DT to ma 2024 on by 10 vector. When an apprenishly in sed where predictor entry contains a set hit is stalled and all prior stores has completed. Since there is no was to dear individual his sector entries in this

The hardest challenge in designing a distributed data cashe was the memory disambly sailes hardware. The TESPS ISA moticin each black in 32 maximum in med leads and storm. Since eight blacks can be in Hight at sace, up to 256 memory operation may be in High t. However, the mapping

of memory operations to DTs to enhance and their effective addresses are compared. The two group and oppositions are laid eleminated here to distribute the LSO among the DTs, and the determinant when all earlier cores has completed section all DT (see that a held-hard load can be se-

We taked the LSO digital stills ambles density by brate force. Con complexity, as leads would have to be musted to two places and then spack mnice on the appropriate artists. Partitioning the USD amount the DTs was on earfale partitions, factoral, we replicated from copies of a 256-eater LSQ, salest such DT. This relation is no might be and was reful bin certic maximum occurance of all LS Or is 25 St. but was the least complex all grant is often the protetype. The LSQ manaccept use load or class per cycle, forwarding data from surface stores as necessary. If there is a partial in-Hight match, [e.g. mattink store have instructions feed in a six ale, later land word in struction is the

114 Secondary Memory System

The TRUPS and stone or sensor a RAM coats NUCA Billiams, and nized into M Memory Tites [MTs], such one of which holds a 4-way, \$4800 hank. Each MT who includes up no chip nework [OCN] noter and a singleeatry MSHIR. Each hash may be configured as an L2 cache hash or as a stratches of memory, by sending a configuration command a grow the OCN to

3.6 Microarchitecture Execution Model

As defined by the ISA, blinds execution is atomic, and the leage in the copyright his logical risew of atomic block execution with operability

2. Execution: the actual execution of the individual instruction in the

der of the certion block execution refer to all three citys, while the balking phrase execution of block instructions refer to this second city

Since the processor care to physically distributed, different parts of the black are fricked from 4 flowers tile, and execution has near in a distributed this is a series the different tiles, and the problem to reach a cell is street across different tiles. Table 3.2 summarises the timeline of block execution and the how the different microsoft interact to create the logical view of atomic block

Below we distract with a detailed example, the example of the Alexander instructions alone. A detailed description of the his discreme and the in elementation of the microarchitecture pipeline can be final in [234].

Figure 3. 13 shows an example of how a code sensence is executed an instruction to reservation station in an ET, AI of the operator described

Figure 1.8: Recoding of a single instruction and mapping instruction in orderation stations.

are delicered over the OFN. The code starts when the conditations of [1] is in sed to RTI. It reads the rate without from and destoral register R4 or from the write space of a prior in-flight block that write to R4. That cake have to the left enemand of two interestings, the sec. N M and the midd. N M

When the test in the claim receives the register value and the immedia To note from the most intention, it field and yet does a yetlicate which is mated in the product of motion X[N]. Since X[N] is producted as fixing, if the motion of product data is not as a material payment has a rate of I, the modis will fixe if the predicate's nature is 1 N N I will not to be 10 to be an N N marketing the switcher left as word in for grand and the confit to the address field of the Lo find word j. Note that if N | does not five due to a mix match of predicate, the dependent load wi so t fire, as it will sever receive its left up eras d.

with the cubes of the load and motive it to N[12]. The DT non-the load/core

rand R4 N[1,L] N[2,L] mov1 80 N[1] feq N[2,p] N[3,p] f. mail 84 N[22,L] f. mail N[44,L] N[24,R] lw 88 N[23,L] mov N[24,L] N[24,R] mv 80 callo \$foo GT RT0 RT1 RT2 RT3 ET0 ET1 ET2 DT1 ET4 ET5 ET6 ET7 DT3 ET12 ET13 ET14 ET15 ID: | the the last and I for the time, in this example; is ensure that they execute in the properly regram a clerift key share the same address. The result of the lead to find and on the the men in the state to the address and data fields that although two instructions are targeting such operand of the store, only as end there is respective ratif for the to the products. When the rare is real the black has preduced all of its suspens and is ready to commit. Note that if the stage is a sliffled, it does not affect memory, but simply signals the DT that the core has horsel. Notified register writer and stores are used to easure

3.7 TRIPS Prototype Chip

They by conditions as dissylementation of the TRIPS chip were driven e principles of partit is ning and replacation. The physical design and floor plan directly represents the logical blerarchy of TRIPS (decreased of only by solution ship, a surest-neighbor networks. The migrouphilests propriet deleof modularity, liling, and communication through well defined not works, are directly reflected in the physical design and simplified the physical design pro-

The sale exemiles to our search solely or communication religious

extensit has contratter [EBC], and the "processor half" command from th ERC to the GTs. All of these signals are latency to breat, however, and all Distanchical desire has been common studies for pairs come time

mality recense designs, in which a processor can be my brated many time on the chip, both one hierarchical and distributed design. THOP'S differs from SOC: and CMP: in that the individual like are designed to have divened for elimination but to cooperate register to implement a more powerful and designrealish in an intercernant. In the following two colored land, we first employed detailed specification of the TRIPS thip and then briefly discuss the physic

The THOP's chip is in please red in the HIM CU-DLASIC process, which has a drawn feature the of Minn and Tipper of metal. The chip it of hich demore than \$70 million (marks) on in a chip area of \$8.5mm by \$8.57mm, which is aboved in a 47 Sm m concern ballsorid arms, nucleone. The TRIPS obtained in is pieces in a communique estategra amy pecuage, the trainer cap compa-term included faculty, craft, and graduate students at UT-Acotta and an IRM Microelectronics ASIC design from lacated in Acotta, T.X. UT-Acotta was more applied for all and decision, basic decision, profitation, and timing, 10 M.

for the physical design tasks including test influstracture inventor, the final skyrical Bromian, slack coast matter of all cells, and the turesat process

The final clock period at won't care process parameter is 4.5m; which accounts for a colonistic clock skew as 4 wirds a parasitics from the fit of larger To this order, this corresponds to approximately 32 factors of 4 (where I FO4 is the latency for a single inventor to drive four capies of instity. By companies, leading edge colons microprocessors are in the mage of \$3-28 EDG M. A major dealer state country with a more proprietared during term THIPS and the time into that is meregime. All fing a more apprecise process and his conservative rates than a classifierd ASIC process was blue also a TRIPS clesh mis competitive with that of a high- and commercial microprocess

Figure 3. II shows an annual sted Bourston discrement the TRUPS chin takes directly from the design database as well as a course area breakdown by finction. The diagram shows the boundaries of the TRIPS (the , as well as the placement of regions and SRAM arrays within each tile. We did not label the national titles (NTs) that correspond the DCN state that are as small. Also for more of viewing, we have a milled the individual logic cells from this plat Table 3.3 lim the area breakdown of the major companies to of the chin.

Controllers: In addition to the care tiles, the TRIPS this ake includes

المستقالة المستحد Moco 88#.88**m**eni

connect to an indicional SGB SDBAM DEAM. The chips a ship controller [CSC] extends the naship network to a four-part nach contexts in glories by

connects to other TRIPS ships. There links nominally man at an exalf the surprocessor clock and an in 266MHA. Each TRIPS arministra heard in dude of to transfer data to and from any two regions of the physical address space including a director manned to other THIP'S procedure: the clabal aboutes

server, OCN, the two D1/188MIR DDR SDRAM counties (SDC) with

Finally the external has controller ERCL to the interfere to an an

efectoring twiling and STAG 1/O houndary cons. In addition, we and a IRM partners added a sone controller to enable the sone chains to be now for cities debug is fractional mode by allowing can accord to most of th started trace. The THIPS chie also include two absorbated force IPLLs click, and two clocks for the DDR SDRAM controller. We as a methal these

processor. The EBC allows the PowerPC to read and write all TRIPS chip

architects mi state imemory, pointers, etc., and relate intermed ten nests from

10s and Test: The TRIPS day include nearly \$50 algorithm, including

US for each SDRAM interface, 3 12 for the chips webly controller [35 play per channel × fore direction (× input/ontput per direction), and 68 play for the

ERC. Not shown in Figure 3. II are the individual 1/10 cells, which are place

near the negleters of the chip. Some of ET t. MTt. and DTt are larger than

Finally, the ASIC methodology maniput LSSD continuous for man-

clock, CIC clock and SDRAM clock burndaries. The CIC interface C2C surkets are exactrosized into the local demain before being need.

a mada lar ASIC design flow, As a part of their ASIC services, IBM yes to me and memory, but also for branch prediction tables, instruction open or and programing clatical. Through a nairemity licence Sympactic apprished their Deriga Ware raite which in Saded syntherizable integer units, fi point units, queens, and P.O.C. The design-time advantages of the ASIC flow are officed by greater area and clower clock rates relative to a contemple ign Manager the advantages of tile heat portionals small upon directly to

Table 3.4 (see add lead death as the deags of each TRIPS (de. The Cell minime column three the number of phendle initiates is end-tlie, which precion a minime minime of the complexity of the the. Areas White indicates the total number of hits found in dense register and SRAM arrays as a per tile basis, while Six echaws the area of a representative of such type of tills. It led so to see a how the total analyer of capies of that tile arrays the entire chip, and N Chip Armindicates the fraction of the total skip arm

At shown in Table 1.8, the DT is sensinly the most employed the tiles, day in home part to the demands of an extender memory review mather than the distributed nature of the TRUPS agrees on its cell count and arm to skewed to member by the CAM arrays for the maximum cited load/store

т.,	Francisco	Call instances	Dit.	S law	Inclusive	Acres
GT	Process contint	2,554	53 K	3.2	- 2	- 2
BT	Register \$16	26.2.54	1400	1.3	5	2.2
IT.	Instruction cache	5,445	1.35 K	1.3		2.91
DT	Li Data saske	1192 86	55 %	5.5	5	3.3
ET	Indication reporting	57.557	1330	2.3	32	25.7
MIT	L 2 Data rasks	80215	5 C K	6.5	1.6	31.2
XI	OCX NW horrace and	23,467		1.3	24	7.3
	continue.					
SDC	DDE SDEAM ITERE	94,440	630	5.5		12
AMO	DMA continue	38,3 65	480	1.3	2	1.
EDC	External has continued	29.3.47		1.3	1	1.2
CIC	Chiefarchie communic	6214		2.2	1	1.2
	ratio continue					
	Total for exist skin!	7.01	11.50	338	. 116	

dealer army ofractions was available. We saw the same phonomenous in OPN and OCN matter. The large cell matter in the ET are due largely to the the standard cell library rather than implemented using a custom datapath.

Verification. The partition of payors of the TRUP's this facilitated a highly blemed ind verticalise strategy. Each of the 12 tile design teams created a op his tice ted self-th whing twith each for their the that employed hat hall excess and madem tests to exercise as many of the compression as a sociale. The ran To some oversige, we as greened each tile design with event considers, an

emal state machines hit all of the pertinent states. The the design up proach also provided apportunity for concurrent development and verification of the tiles before putting the tile together and certification of the processor We also used for remarkable or reformance region in a Color

a raite of microb exchanges, with some randomly generated programs, we reduced the average error between the low-level performance simulator and the BTL simely for from MS or average to 25. This effort as exceed sixt wa performance bugs, sen of which turned out to be worth the effort to fix. The three most significant on a worse fixing the innepriority in the ET, reducing the Back penalty by an expele, and reordering predictor operation to eliminate

In this chapter we described the EDGE ISAs, the TRIPS ISA which i as els crasce of as EDGE architecture, it makes architecture design, and outlight the implementation of the TRUPS amontone chip. The dataflow graph abstraction in the ISA and the scalable, partitioned, made he nature of microradities are provide anteral support for polymorphism. The prototype chip provides landed polymorphism support, namely, explicit thread level pur allelies by onb-dividing the hot method window, providings the of memory

sions in the context of the TRIPS processor architecture.

Chapter 4

Polymorphism in the TRIPS Architecture

Emerging applications with heterogeneous computation needs and for and combat processor complexity. Architectural polymorphism achieves this by altering the helanter of course grained components to support differen also requires an anderlying architecture that can scale with technology and is held to the model or make surplified are blocks. In the arm has distinct we described the TRUPS architecture which provides such a smith leased most star processing substrate and in this shapter we are TROPS as the baseline and the one for developing the mechanisms for pulpmorphism.

The seed for architectural mechanisms for distinct application domains has been exident for many years and has in fact been emiliable for almost a Pewer C Alices [14], SPARC VIS [82], PA-RISC MAX2 [84], MIPS MIDAX [74] and Alpha MVI [2] provide general purpose architectures with a means to exwhile much reals data berel variablelium. All of the laster clies not extension

of polymorphism. The final-end of the processor is configured slightly diffemally to read from a reparate physical register file, whereas the execution way. Typically memory distmbiguation hardware and caching appears differently. Simultaneous matrick reading |SMT| is a second form of polymorphic which is convenient a prevalence in the despectation and chies and then onore 1851, in an SMT are cover, the register files, in traction frich busic, an core of the microarchitectors uporates the came whether executing one thread or multiple threads.

While this limited polymorphism has been sufficient that far, for one apelication from demoist to a crown big the inherent bid empressity of an elications

• Multimedia databases: The amount of multimedia data is proving modify and different types of computation are required on these fatabann \$25

and constation \$15 is games all have different computation needs, with

phones and handhold game decises are expected to perform multiple

and integrates up to circumscensum, each being dedicated to a reparat starturit. (Daniel I Demantics ararmatas, and a print gral I (Dan stratio Second hand held mane find are expect a medited of proceeding to be on a single decker whed Ethernell, whether (WSF), and colour [16] om manication, et empe mana com ont, à lem et de técnit Mention, records and distribute the management. 3D cound field, and 3D cides represented to come a few [86].

Designing markiple operities rate than in tends on a processor complexirablem. Architectumi palymanyhim cairm this application betempensity problem and addresses technology constraints in a complexity-effective man per. We defined a strangerships in chapter I as "the ability to modify the Da cileadire of course grain microards bed are blocks, by disaging course legic but leaving data path and clorage elements largely named Med, to build a programmable architecture that can be openialized on an application-by elimities having. We are complexity-effective in the came reason at Moure'

A complexity effective duties it a duties that: Hembrace a rela-tively matter afavorable principle and a recited mech-ation, and Hard bear mitted in collapsing susceously com-justily like there may find annual and depart mediation.

phicm, and explain why there mechanisms are final amental building blocks for sek men tion.

The TRUPS architecture is used as one specific architecture and macoursided are to implement and realisate these medications. Chaosing a eddic ISA and micronribitecture is necessary for quantitative evaluation bis ISA and micronribitecture are also inherestly suited to support polymorphism. The dutaffew graph abstraction in the TRIPS IS A directly headcascorrancy at different grass briller. The distributed and model areas are at the microarchitestore already provides the coarse grain building blocks that

are received the architectural and empression. ers architecture. The specific implementation of the mechanisms are sted is THEPS processor microarchitecture, but the haste mechanisms muld b

4.1 Principles of Polymorphism

Adaptivity across granularities of parallelisms: Polymorphism is intended to specify between easy commentation canability and adopt to dram ing an elication behavior and demands. As described in Chapter 1, we identify the differences in granulations of parallels may the find amental architecture

difference between applications. Based on granularities of purcleions, pregrams can be broken down into three catemories; instruction-level carallellan tage must be able to adapt to these three grass betites of parallelism.

Economy of mechanisms. To be or malexiat effective, the automatables mechanisms man the few in number and they should provide a set of p reconfigurable fluctionally to microarchitecture blocks that can be used or exhibite an architecture on an architecture by sometication back, instead of a heigh a net of fixed function extensions. As a short case state, consider a limites that has simple data beed parallelum and operates on two lon army). One fixed fraction extension is to build a vector core and interface i to a convenies alpes on or as 4 compile programs into vector in tractions. On the other has 4, the spirit of polymorph has it to create mechanisms for a con-centional processor to modify the horescling fields, solest and execution legito provide incinacion efficiency and modify the memory cyclem to provid is so of for recular memory accesses. The design challenge is to determine

Our are much to determining the emechanisms was to identify the haic properties of programs as d how they affect the microarchitecture. Nased so this analysis, we determine a find amental set of mechanisms that configure th microard ited are differently to coppure different grass brities of parallelans

Parallelism.	Resources	Politica
	Execution	ore management
ILP	Reierrat les i	May making to distract ow graphs
	ctation c	
TLP	Reservation	May malify be dat after gray to from
	ctation (different thread :
TLP	Is or each talk a	Printing between threads
	Sidect legic	
DLP	Reservation	May begressrated dutation graphs
	(1818/8)	
	Datastor	ge management
ILP	Regater film	Register resuming acress his day
TLP	Regater film	Storage for and hed are state from
		many threads
DLP	Regater tile	Bigh register tile has dwidth
DLP	Memory 1710	that has dwidth and outbears on-
	1 ms	in Red memory management
	Control fi	ow management
ILP	factors of the s	Costrol operatation
	fetch	
TLP	Is or each talk a	Control operatories and fetch mat-
	frich	tiple thread:
DLP	factors of the s	Optimize regular coatrol flow -
	frich	properly deal instructions

The said is memory is configured as a new saiderin cache access. NUC Albanks. The banks have miss-handling legic, and ofting arrays, and status his to behave the a cache. The saidily a dweet also provides a high-handwidt em is alogy introduced by Kim et al., the THIPS ship implements as

To community, the fixed recorner, namely the data codes and inornation makes, the specialized recourses, namely, the next-black predictor and MSHRD, LS On, and other member looks and the reference has been seen

4.5 Thread-Level Parallelism

cover a likesting can be achieved by mapping multiple throats of control or to a place expression. Taking et al., introduced the terminature of simplic

haire many important. The SIMD approxima paradies is one efficient a of instructions and relacing during complexity, for exactly these type of pro grams. Polymerylant medicalism made and to taller as articles are interesting material legio. Executing the time dutaffirm graph in a long wit many iterations, can be viewed as Single Instruction Multiple Data SIMD execution, where the dataflow crash can be viewed as one single SIMD in metics executed acress multiple ALC sites. The electronic of repetitive instruction frick and an accounty speculation must be removed, to reach th efficiencies that a true SIMO model ma provide. We develop a medianism called insets diam residuates that augments the lasteration selection logic at each individual ET to reasonapped instructions and augment the frich

logic to fitch interaction in a loop just on or. Also, with come type of DLP programs, a fine gain motivate safe model that executes a Makin let normel by Multiple Data, MIMD), execution model is preferred. The ILP and TLP execution model of requesting a program conster that fricker and may concernive dataffer graph | jumetime through materily emilities; is not very efficient mappined to this approach because mandels. By adding instruction storage copy or and sequencing th ALC: independently the execution core can be tailed to look like a MIMIT

We classify the type of resorce in polymorphus architecture into there categories haved on their fraction. In the next section we identify these

Fixed resources. Some resources in the amount open to in the same way monthly of the application execution on the property. For examp dractions as possible and provides how-latency access to the program's in clear than cleans. If it all recourses are found mental to the hards age

Polymorphous resources: The configurable resource in the processor perfirm different types of operations or change their operation policies, deichedolog policy to fetch from multiple instruction streams if the proover a maticaged to execute multiple threads simultaneously.

Specialized resources: Some resources in the processor are specialized for plications never needing on the functionality. The replicated register the durage is an SAM procedural an example of cach a recourse. In an SMT processor which support up to 4 simultaneous threads, there are

4.3 Mechanisms The THIPS IS A expresser constrainty to the hardware by herables programs into his data and exceeding instruction depondences within these blocks by much in gride during the grouph explication the USA. This durinflow group high rise tion to need no the nu Sping them energical different granularities of parallelic and the mechanisms are built around this dataflow execution model. Helow we describe the automorphism mechanisms with respect to the three main

four capies of the architectural replater Me. When only one thread is see-

The sense absolves we can up it so has maken a recover continue ab ma

phase architectures the capability of a dayling to application used a. Manu-

persons and historopersons systems can be analyzed in terms of this resource

chinification. Meterogenessis; quiens have sulpfixed recorded and specialize

minument of the example the vector register file in the Turnstala and beet a

is a specialized resources, whereas the execution core is a fixed resource. The

IBM Cell a more of SPE, may be maddered uncolded more my date than

ready here brought into neighboring memory banks [665]. Today's madion

thin and the Minchel Micha be niewed as humanesses a systems with only

on time on the processor, three of the register files are completely unuse da la chicago -both memo er and rechten. To be efficient, these upplication operitie reconstructional be minimize

4.1.1 Execution core management

The THIPS IS A breaky programs into block and so notes duration graph in thee block. The execution care precises a construction the time exist which there during up to make dynamically may yet. These concention cuttime who referred to as block that place blocks are may set to them) form one a elementation or ourse and are managed differently haved on

processor companies the execution care, instruction frick and control, and

Acres different granulardies of parallelland, the surface of these datables graph) can vary, and the types of communication between these dataflew graph) can change as well. With reposition of es, where ILP is the damin an type of parallelium, the circ of the graph circumstite small wof the order of 28 to 48 instructions. When even the much into recently dataflow one to from 416. when there is ample data-level parallelons, these graphs can be very large. To extract SLP efficiently, the reservation stations are used to man a numher of conceptationin fricked durather granks, there there could use insignife mall and many outh graphs are needed to \$\$ the recernition station space To extract T.D., the recentation clutters are partitioned acress programs and dataflee graphs from multiple programs are mapped to the reservable a sta-

Depending on the type of parallelism, the control behavior of applica-

statically properties at committee times 41.2 Control flow management

tions may quite dramatically. These control they mechanisms mytore all of fetch across threads for TLP, and 3| Optimized instruction fetch to exploi most live coarse flow for DLP. With proofing with most is in the close-level sarallelien, it is gracial to have highly accomic control flows the control flow is very keepstar and is hard to determine statically at compile time. With thread-level parallelism, to up timize the performance across threads, the instruction flow management between threads is an important quarties to address and introduces policy decisions in hadding the instruc-tion feets madels. With programs dominated by data-level parallelism, the central flow behavior is very repetitive and early predictable. Using control they operation technique can unnecessarily place instruction frick on the extitual path to execution, fastend, we design an optimized instruction frick medianism that reason instructions that have been fetched once for a those

the progression stations are used to hold one single large graph that can be

There mained flow techniques are not matually exchairs. By noing limited amount of control speculation within each thread white extracting TLP, the process refliciency can be further increased. Some programs with

DLP are less apported by a fine grain MIMD on transport the control flow mechanisms to configure the emorator like a MMD maddine are similar to he instruction flow management for TLP.

Based on the liveness, data values in opporants may be classified as is a program is within a few last of code, and in the TRIPS compiler such data are live only within a block or dataflow graph. Long-term data is data where liveness is trackable within a function, and in TRUPS such data are liv data is written to memory. In a RISC and before chart-term and hasp-term values are stored in registers, and persistent data in memory. Polymorphism haved on application needs.

Short-term date: Dataflew graphs are directly mapped to reservation stalinas and short-term data are data approach placed between under in the datative graph. There are mapped to reservation stations and the ISA explic tily arrigan there calses to specific reservation stations.

Long-term data: Long-term data are called passed between dataflew graphs

that the compiler has placed in different blocks. These are mapped to the ar

our design of configuration of \$2 storage as send dayad to

difference register clarage and depending on granularity of parallelism, the

resider once can be managed differently. When executing only one throat

and while executing multiple threads, the physical register space is partitioned

Persistent data $|Memory|\colon$ Conventions tyrogram ming models and in C

C++, and Java have a simple view of memory used for storing persistent data

with the hardware and the operation contempressible for eaching policies

and marine. This otratees work well for irrestlar resonant where denomin

However, when the program behavior is regular and well of mit and, then

is benefit to explicitly managing memory through software. In the TRIPS

drip, the na-drip memory to construct of noing at the of interconnect of memor

hanks. These memory hanks are expected to software and can be can be con

Spared to behave at NUCA type L2 code banks [66], contributed memory, or

has dwidth interface that enhances access to persistent storage. The IBM Cel

ement. The Streaming Register File architecture of Imagine [235] in sig-

processor and impacts are other process on that is chall excited memory man-

among mality is thread to

plementing these mechanisms. In the fillnesing sections we describe these mechanisms for each tree of carallelist

Parallelism	Resources	Policies
	Execution	ore management
H.P	Receptables:	May mali ble dutaffew graph:
	(1413/8)	
TLP	Reservation	May mality be dat after gray to from
	(141303)	different thread)
TLP	Is or each to a	Prioritize between threads
	Select legic	
DLP	Reservation	May begressrated dutation graphs
	(141303)	
	Datastors	ge management
ILP	Regater file	Register resuming acres s his day
TLP	Regater file	Storage for and hed are state from
		many threads
DLP	Regater film	Bigh register the handwidth
DLP	Memory 150	that has fulfill and others on
	1 (0)	in led memory management
	Control fi	w management
ILP	Sant rand in a	Control operatation
	frich	
TLP	factors of to a	Control operatories and fetch mal
	frick	tiple threads
DLP	factors of to a	Optimize regular coatrol flow
	fetch	peace felicked last mid-last

Granularity of configuration: Apolymorphus cardifecture after behan

or of competency microardited are modules, by charging the matrol loca

ties by configuring Amegrain blocks may be a challenge. Remarigamble as

dillectures perform the grain reconfiguration to synthetize blocks with diffe

wars. They have all mainly provided application specific hardware and not

programmable hardware. As recieved in shapters I and I, example includ-

derigns werk well fire a must domain of peoblems where the upplication can be easily mapped to the bardware, typically 'regular' upplication but perform

and a complementary of the property of the pro

ter at adapting to different types of programs. This shapter describes the

the management, and a configurable memory cyclem. In the chapter, we qualitatively justify this approach in terms of decign complexity. In the next

three charters we discuss the apparationic performance results that each an

provide a broader discussion comparing pulymorphism to other approaches

targets hop tend paralleloss only.

over a inella e. The Gam architecture performs fine on it on fire m

EPGA: Twoller Park-TPP Modeling Planmach and MSH All of their

4.4 Instruction-Level Parallelism

the clost broaded and an efficiently by see letting in struction design and letting. Proriors publications have referred to come of these techniques by referring to them to the Demorph made of the processor [223].

tion wholey and receives to exploit commency in the intention stream. To exploit the in the TROP's processes, the resembles station in the care are configured as a large, distributed, instruction into window. The direct in eget encoding in the TMPS IS A combine out-of-order execution while arold ing the non-civil ire in or window looks prof convent is not marchines. To one the instruction to them effectively as a large window, the processor must provide high-has dwidth last mad in a ferchiag, an appropriate control as a data consolation and a high-handwidth, low-kitency mamony cyclem that presented expension quest rection revel excelle the implementation of the mechanisms for exploit

4.4.1 Execution cars management

The polymorphus creasures in the execution care are the recognition stations that practic instruction and operand starge space. To extract ILP, these reservation stations are configured to behave the an instruction window. Such a configuration now the reservation stations at each Execution Tile to for dament ally a three-dimension all whede this gregion. The x- and y-dimension corresponde to multiple instruction data at such ET, as shown in Figure 4.1.

into a 2-D region, arrigaing such instruction to one under in the 2-D space Several policies can be implemented to man the instructions in the ISA to these hardware data a proided by the micros malifecture. In the TRIPS and strate we an asset fixed rise blocks, and break the instruction window in to groups of 225, with such tack group being uniqued one blick of interestions. Recall that with 64 progration stations at such tile and a total of 36 percentian tiles the

may duraflow graphs directly to the ETs. This physically distributed has wis down stread across the ETs allows orders of many its delia greater in wis down

the compared to conventional constraint anterior decision of the TRIPS

mplementation we achieve an order of magnitude increase. Since there ar

multiple reservation stations at each ET and multiple ETs, this window i

Process & Death and in the principal process of the Park I was not a first the Park rights 4, in 1000 1 to Statements into purpose the first terminal mapping of in-tention to recognition stations in a group. All communication within the black is determined by the compiler which assign instructions to recognition cations and or emple are denomically musting directly from E.T. to ET. Con-

The number of him that can be excelled in the target field, implicitly limits the size of the dataflow graphs that the compiler can construct, and hence the circ of the block of the number of hits in the carget field also directly omegands to the amount of state the microarchitecture needs to support predict branches to be predicated, thus hid in ground there in side these graphs. the coffware darlings in halfling large dataflow graphs where the number of a smed instructions at maxime is small. Purple TRUPS projets a chiefer charge a 7-bit to goet field since any experimental prophy showed block sign were allow as to pack the compiler to its limits and explore the design opine.

4.42 Control flow management

they logic employs two mechanisms; control operatation to build large instructhe window and high has dwidth in order to Brid.

Control speculation: The compiler is able to prove at a blocks comprised of dataflee graph that are between 28 and 60 in the client on average. However, to extract LIP, a much larger wholey of interesting must be examined and

this is achieved by operatoring on control flow between blocks. The basis ha ill a next-block predictor that can predict the next-block to be fetched and and may madilyle blocks to the instruction window and execute instruction action there has do character words. The approblem a predictor is a specialized erce and the receivables clubbes from a polymorphose receiver, both a

Next-block predictor. The next-block prediction is made using a smilet-up Communications of explicit PAL which english a bit are rather to distribute the branch that is a redicted to be the exit of the block-scall each block

The value presented by the exit amplicant is used to index into a set of ST St to ab talk the next predicted block address. The branch to also predicted by the exit predictor, and is need to select an address from the multiple ST St. Manganuthan et al. describe the predictor in for her detail cheeds are [225]. This predictor organization explains the the predictor to be decoupled from the instruction first engine. The He servation stations Specifically is the TROPS processer, the total laster

tion window size provided by the handware is \$124, with \$4 stats available at each of the M ET: [M + 64 - M24], There 64 date at each ET, are brokes him groups of 8. The TNIPS processor allows only fixed size blocks, with each block containing US instructions [consol instruction are exceded as NOPs by the compiler). The hardware manages the remap one black of MS inciracions | S + M - MS|. The regroups are used to map speculative blocks. These groups are managed the a circular to ther with the non-speculative black successively being married to control to \$ 2, and or one.

Migh-bandwidth instruction fetch: Tellill be larged in that ed in tration what we the processor includes high-handwidth in creation frich mechasions through the one of a set of partitioned instruction caches. There has he which are in the last modion Tile $|\Pi|$ are a fixed measure, meaning that the but arter & the came is dependent of the type of parallelism. These suche has be are interferent or do that much hank holds 32 of the 325 instructions in a block and the 32 hotmodism in each bank correspond to hotmotism that have been migned to ET (hother same row to that IT. The processor area a program outer that points to the clart of blocks where a block header is exceded When there are free reservation stations to map in structions, the control log new, with a single access and streams if the this hank's respective new

4.4.1 Data storage management

Short-term data: To extract high ILP, the short-term data apena de armapped to the reservable stations. The management of these short-term at a special of forms another fixed resource in the processor. Short-term dat were level, this communication maps to operand spaced between recentled

through registers and their life time in the program space and hyle dutation cracks. Read for magning in marcalingal processing create light between adeat inclinacione in the inclination window. Similarly, when ex LP by speculatively executing duration graphs is an EDGE architecture, must create links between dataflow graphs dynamically, in that the court of execution of a dataflow crash dues not have wall small its predection to data aperands efficiently, the microarchitecture must implement black level register resonate at a silver rapid planting of ration is dressed data flow graphs

Persistent data: To convert high ILP, the processor memory system many eraride a biol-band width, low-lateure data codes, and mort majora is respec physically distributed data storage in the processor core, comprised of Data Tiles [DT], a configured to behave like a first level data cacks, and the on ctuading cache makes, LSQc which detect load/ctore dependences and enforce the correct ordering of leads and stars in the program, and store merging

In a principle, a programs, with DEP, can be executed an the TRUPS are or eredying an emitral flow speculation and having the hardware extract IL is chapter 7 we present a detailed characterization of DLP programs and a der the following three categories: the execution core, contradition, and data

4.4.1 Execution core management

For programs with ILP and TLP, the dataflow graphs were typicall mail and control-flow speculation or explicit multithreading to a secondly to with large iteration country. As a result, the hardware overheads of special ties and reflying every such a femalicity reading on a be rightly reduced or mpletely removed. Instead, the most efficient way of managing the see large saralled dataflow graphs to the reservation station, without relying on

4.5.2 Control flow management

Control flow or on lating is refut to be a loss of any for DLP program. with power efficiency is in ordering frick and high handwidth in ordering from

gether, provide a highly effective distributed processor substrate for extraction

implementations of SMT have focused on extension and modification to

present a set of polymorphies mechanisms that can be used for fine grain interfereign of hot medical in a processer's pixeling. He harden sharing data ys than distance elements, our implementation of SMT eliminates some of the mylicated simultanes of precious implementations the multiple monder hallow

The basic principle for copporing throughtend parallelism is to cylin monetar class per en agree the tweet modify to the mode, and an gravest the central legic to dynamically share datapath components, like the fluctional eterant that he call. We break the arrange or charge manages into all or with such the chair gating sed to a different thread of material. The controllagion is an granted to implement a fairn or policy to allow each thread of control to access the datasets. And finally, the architecturally visible storage, name the register film, are replicated. With he such thread, the process or still extracts
(E.P., but as each sike is nameworthous when maning a single program, the ILP extracted per thread is lower. In the following only rections, we discuss the mechanicme that implement this tech

led his ordered between the chronic that have a result in order has be execute.

4.5.1 Execution core management

fair and of holding a navigoculative and opeculative blocks for a ringle

. When the mark line: Static partitioning is straight-forward and easy to in plen ent, but can bere processer resurres paurty nillinet when dif-front threads have different nore notice of principle. While, dynamic partitioning can be aware of onth application needs, it increases both the hardware and coffware complexity. Expressing a ser priorities an policies to the hardware introduces software complexity and dynam partition of processor resource introduces hardware complexity. Mantware ambling havel assumation on implement denamic partitionin without any charges to coffware.

• New to purities: The reservation stations form a 2-D instruction space which can be dired in different ways to may multiple threads. Figure 42 shows a spectrum of partitioning strategies. The main differences be-tween the partitioning schemes are implementation on appeals; showed distance from the register thes across threads, chewed distance from the data tiles across the thready shewed in struction for dub and width and hitency. The partitioning strategies shown in [a], [i], and [c] in figure 4.5 add complexity to the instruction fields logic as the natural align 32 instructions per hank many by changed, or the instruction fields agrewith most be assumed at the party intending agree news. Figure 4.24 suchus ped and requires modifications only to the instruction selection

legic in the core. Since the TRIPS ISA has fixed \$25-instruction blocks any hind of partitioning courtegy must provide at least 125 class for each through and any additional class can be used for conceptation within a To keep the microarchitecture's execution as close to the ILP model

as you his, and to reduce implementation complexity, in the THOPS prote-type chip we implemented a simple sharing scheme denoted in Figure 4.24. Each thread gets 1/4th of the resources and up to 4 threads can be executing

45.2 Control flow management

Control flow management mechanisms to support throughout puralletters to not very different from the mechanisms used for ILP, namely control they operatation and high has dwidth instruction frich. The added requirements to that both must be done formultiple programs.

Control flow speculation: To special TLP, costed flow speculation in required for each thread, which can be achieved by hadding multiple next black prediction, one for each thread, or chapty chaning on epitedictor between mali is to bread a. Since the exit his tary is constal to high according prediction we replicate the global history shift registers and maintain one copy for each particular thread is used to make a prediction using the chared exit predictor

[18 him per thread], the mostling replicated stampe is quite small. Mish-bandwidth instruction-fitch: The manufactor of the instruction

cack et an dit be net work to tires miliotro clination to the processor it again lifes lical to what it required for copporing LEP. The only difference being that fercise of blocks are initiated from different throads every cycle, which is dependent on the rate of which threads complete. Tollien et al., investiga erem I pulicie that can implemented for instruction first between multipl contending threads [B4], in the TRIPS protetype we implemented a simple mand-make takene which give easel artesis to all executing threads and

4.5.1 Data store management

Short-term date: The management of short-term data is identical to what is done to extract LEP, since within such thread the process restract LEP to to a lower extent. The microscolic rate of a using concentration of a person is rank that there short-term data calcer parred between an decia the datafles graphs can never be sent to raises from one thread to unother thread

Long-term data: To opport multiple through executing on the same procenter core, en ough replicated register storage must be provided to maintain the architecture state of each executing thread. One copy of the arch

the storage is expensive and not in the spirit of polymorphism if that storage case of he need for an ething size. The replier renaming hardware many h my lication of temp energy through or dutup with in required to create this recon tion rat to reck ter tile.

Persistent-data: The memory system operate much the same as whe extracting R.P. Similar to modification to the register renaming logic, the control horiz in the data tile is madified to ensure that lead/core checking in erformed only within a thread and not acres o threads.

4.6 Data-Level Parallelism

Data-level purificion is most commonly found in streaming media an rejeablific annileations and in characterized by the fillerwise main attribute predictable hop-based coated flow with large devailed coate, large data sets engalar access patterns, poor locality hat talenace to memory latency, and high computation intentity [M4]. The dutation graph abstraction already hands that the efficients connection this bind of sorollation, since the con-RISC or CISC ISAs. We halfd polymorphon conclusions to further optimifor the months control and dataflew behavior exhibited by these applications

Chapter 5

Performance Evaluation: ILP

One of the primary goals of the TROPS architecture and the ISA is to which we have califored to be within MS of the hardware.

We are a ret of heartmark rates with different levels of complexit and different types of behavior to quantitatively evaluate the TEEPS design and demonstrate in effectiveness. We start with a set of hand-written and contenchment became which we bearily hand up timized and toned haved on proffing the hemelt and understanding the internations between the code and the microarchitecture. This microbenchmark analysis demonstrates the potential of the architecture. We then employ a set of data parallel kernels and the EE ABC maked of heartmark on he to explore the performance of pr the DLP benefit and the EEMBC programs to quite regular and the memory

to Section 5.2 we describe the methodology of this ILP study and tools Section 52 discusses the performance rounds.

5.1 Methodology

illing, we declayed a detailed cyclodered simulator, called him year, which madeb the hardware at a much more detailed level than high solved simulaters the SimpleScalar [28]. Our performance in Educion effort chowed that mademity generated test programs. We use a critical path analysis to all bires mid cal (188) to attribute percentages of the critical path of the groups to different microarchitectoral activities using the technique first proposed b Fields et al. [45]. These results provide holght him the effectives of and seen head of different components of the microard declars. To place the TRIPS processor is the context of a one on tional microarch heritare, Table 5.2 lists

Our baseline or mearly a point is a 467MHz Aleka 21284 processor with all programs compiled using the native Gens compiler with the 504 -arch set? they set. We chare the Alpha because it has an approximate the care

Memory accesses in DEP agreement are deminated by recolar action typica by and or fixed stride. Mewerer, sign fica at number of other types of data accessor are also present, including lengths accesses to small lookup tables, account to a large number of reasting countage, in efficients of an FIR Eler &r example). The combination of structured and nastroacces i partiems i megalmes a dat automa penyi tematik at can pimelide high han desid th regular data and low lates of irregular account.

Short, been dated. The management of characters date is ideal fed to what DLP does not make any difference to the way must of these operands are managed. The fixed cities regular memory access shown by these emerans present an apportunity for optimizing some short-term data accesses. Whe that implement this strated across in the dataflew graph, show resolution as well. Vector intention sets include some form of a lead instruction that read and higherwords of data from memory and write to a next or register the Similarly, a multi-word had instruction can be used to find multiple words from memory and send the operands to reservation stations in the ETs. This instructions, the communication or ordered of rooting multiple address to the

ack m, and the memory accors overheads of reading such word from the cache Long-term date: According register value can become a buildeneck, if one

resister raise has a lot of facest. For enormy, with among DLP this is a commande abouted the omenon. Purthermore, the programming made dataffew graph introduces in efficiency when the register values du not change acmort such dynamic in the scribble dataflew graph executed. Pergangram on a DLP this type of read-only behavior can be determined by the compl whereas it can be more challenging for all programs . We propose a mechanism called agreem directals at on whereby agreements that do not change during one b tiple iterations of a dataflew graph are read once and rested multiple times registerrend and rename. This mechanism is not restricted to DLP, and can b at titred while extracting ILP or TLP if the compiler can statically determine

Persistent date: To opport DLP, a offware managed cicks memory to I using the saichly memory tile is better than hardware managed or avention a carbine. Other deploys the Smort Memories Imprine and the ISM Cellans the reconfiguration of the memory tiles includes turning off tag checks to allaw direct data amon access as disconnection the cache line prolacement stat

clude black transfer between the tile and remote starage (main memory or other tile), strikel access to remote storage [pather/scatter], and indirect is host of the the's storage. Instead of using the processor is not better these transfers, a now level DMA controller is integrated on skip to perform these for client.

In this continue we have decepted the original or of polymership mand a care not of final amountal moderations to copport instruction local paradolism thread level parallellam, and data level parallellam. Granularity of parallellar erace between application types and how it affects the microarchit

The dataflow graph is used as a naifying abstraction to express our conare efficiently used to held upon latine last mutical, with a next-block predictor is a recipiled resource) and in perform on to I flow a rediction. For TLP. which is competent to compare or agree and little through, the processor reconnect are divided up between the thready and polymorphism control logic is the processor care exacts all threads get to one the processor datapath months in a fair facility. For DEP, which is characterized by our

the specification of polymery has a mechanisms for DLP. Polymorphism series as a natural way to address process or complexit

different amondamine of camillelying The simplicity is implied a table of the medications and economy of these mechanisms corrects solven an horse and it architectures to copport fators application needs, in the following shapters we evaluate the performance that can be attained using these polymorphism

extract large amount of concurrings, in this chapter we fixed an instruction level pumiled on and demonstrate that the TRIPS processor has the potential is explain greater community than the heat-of-breed ILP processors. On eralization in haned on the prototype design to bog a syste accorate character

fortural of many of the heartmarks is small. Finally, we evaluate the new

Processor parame-	Cenfiguration
ter	
Lit faitt acthia Cuche	Der 16KB banks, Seny ter atterfate, I part piet
	hank
L Data Cark e	Pear SKill bankt, 2-ay ter nit neinte, i part per
	hank
Meghten	4 regate hanks, 22 regitter pie hanks, 1 part
	per hank
Is the dies Perd	16 in traction per cycle
tante dina to se	If in traction per cycle
In title clien Commit	le in trecilent per cycle
Lund und Store parit	deffert be had and core para
Cantral Pin- Predic-	Predictor using east histories to predict the uses.
ties	block, employing a ten resement herel/gob are pre-
	dictor challe to the Alpha 21264 with SK, MK,
1	and IR has in the book global, and tear nament
1	mit prelictors, te penirely
L2 Carke	100 B 42 carbe, with 5 parts

Table 5, & TRIPS processor parameters

that will research has BOM clark entirely as IS 8 that heads itself to all the cots. We not him Alpha, a cimulator wild ated against the Alpha hard ware to take the hardine measurements or that we could normalize the level-2 carb as a memory and small allow better our earlies of the process and primar

5.2 Benchmarks

Since a key goal in this discertation is to explore technique	
th become to different types of workloads, we choosprograms	from different

described helps.

ndomork		TRIPS TO			lebi	Speedup
	190	Cycles [18 88s]			C yeles G III s	
manite/ažtimeli manite/aifitili	1.32	7516	35	131	9793 97	1.79
manite/aif:ff() manite/aif()) manite/harefy()	1.0	783 7834	6	1.6	5227 715	1.16
manice/bitmap#i manice/cackeb#i	1.34	932 746	32 22	13	1855	1.13
manite/castdr#1 manite/ideita#1 manite/iidi#1	1.07	1453 521 683	26 23 21	12 13 12	515 511	1,54 1,17 1,84
manice/matrice)	1.00	7752	41	14	4578	1.53
manice/ps - access (0.93	2262 785	21 22	1.0	1199 535	1.53
manates/th brokel (manates/trapek H)	1.6	1073	26	12	185 663 61496	1.22 1.62 1.24
nmer/djreg nmer/djreg orkins/oref	1.38	75 37	24 26	12	61496 65276 7167	1.57
- 11king/phid =-	1. 8	11155	24	1.4	6215	1.62

tion cacks mits more and combined applications with similar and discimina

In previously published work, we adopted this approach to evaluate

is to start a SPE C C PU2000 beaching its by creating on the workload mixer [23]

We do tilled programs into two cologodes namely, here we way in tensive and high memory intensive charef on the L2 code mits not an and concention than

of all I mixed high/low, low/low, and high/high. Other features of program

that can't affect awardian afficience in mall three ded made in dade the arm

In this dissertation, we undertake a more thorough analysis of ma

breaded execution. We have a large application space which includes II

EEMBC programs, 11 SP EC CPU2888 programs, and 22 DEP betacle, it is

hand to determine a-priori what application characteristics are important as

infate (he phase behælte of these application . Por this study, we decided a the approach of using a large on wher of rand on program mixes and generalis

ears at mixer to create different trace of overlassing a marray behavior. If

carering a significantly larger parties of the program behavior, this appears
practice a more comprehensive evaluation of multithreading efficiency. This

evaluation of makegy to similar to the methodology model by Todous et al. and

Speedon: Figure 6.7 shows speed to achieved be expected to TEP and ompared to certained execution of the multi-programmed workloads in IL

made. The workland misse are corol in the came order to for Pipere 6. L. In

dewiews, and 22 slow a speciar, op to 221%, and so arenge 42%. This speciary or clowdown exhibited by a program mixto primarily a function o

the available variables in the errors my. When there is a lot of carolleles in

the thready, the TLP made dues not follow this processor because, on

I simultaneous blocks from a single thread can be executing at a time | |h-

effective instruction window size is 250), while in the SEP made the effective

instruction window size in 1824, Hence, a short was in the TEP must show on

likely to occur for programs with ample concurrency. For each of the program mixes, we examined IPC in the U.P-sonde and over that the occupy IPC of th

programs in the makes that section a showless is 3.24, while that of the makes

that exhibit a speed up to 2.4. A more cophisticated partitioning of reservation

lized, with 256 entries assigned to each thread, only one program mix does

were in the TLP made compared to certal execution. The primary reason

health the research achieved by the TLP made, in that the effects of home ch operatation depth per thread. In this, examining the characters continue we

this spendage.

We exclude the four mirror such marks from this cody, as they are pri-

other authorities on SMT [354].

characteristics to study the seasificity of the architecture to the work had-

is to cover different grant latities of parallelland, 19 per of in 11 m et is a mixer, an

havie program behavior. We are four reparate ratios of bourdaments: If a re-

of kand caned keer by optimized microb exchanges, 21 and of data yard be been marks (DLP), becade we developed, 21 the EERM Control \$4\$ and 4

the SPEC CPU2000 table [844]. Table 5.2 fert the benchmarks which are

without being hampered by compiler technology, we not four operate ma-

contractments that are very specific in their behavior, who is a backing algo-

si him and it a very request is by regions with finited amounts of concernor

I of the 11 SPEC CP C200 the exhausts perform before as TRIPS, One of the main reasons for the lower performance is that the average block sizes that the complex is able to construct is much imaler for these benchmarks. I addition, the control mapped in increse b higher is the SPEC benchmarks is there have more irregular control flow than the simple DLP has dimarks and

is consent the same community much mars in the cell to the level of or shirtication in the compiler, at they are built from large or de-bases and refer

a meaningful heachmark in its for its dying multithreading efficiency. Pari he more come of the extinumation incolemental in those heartmarks as one is the thirteend of leasest time as only with all 1924 properties of at logs artifield in the program. All programs are multi-completion and when a program fin is a while other are still executing, it is restorted. When every programs has

completed execution succe, we stop the simulation and collect simulation data

Since the EEMING only SPEC CRITISH only and the DLP benefit for

tery different beharior and ran-times, we chose program mixes such that al

the programs was as a multi-programmed workload were from the same saids

The three performance metric that we are for evaluation are:

1. Processor Utilization: The fanctional resources in the processor that are kent have. We meature the number of instruction retired nor cycle [PC] to measure processor of Maniles. We compare the processor at a finalist between the T.P.-ande and H.P.-ande of the processor, in the ILP-made we assume the programs in the workload mix are execute certaily, and the IPC reported for the IEP-mode for that application mix is the total agrader of instructions executed a gross all the applica-

1. Processor Speedup: The oresion, ormanel is executing the mix of application is a recipied made, executing one offer another exploiting

most cases, it is also spending fewer cycles in wasted upon latine work.

Efficiency: We make the efficiency by comparing performance and in the sensity. Figure 6.2 shows speedup achieved during executing in TLP made compared to concurrent execution of the multi-programmed workloads in IEP unde sa makipleym on an. Mon kikat, while the ideal configuration bith linky effemance parchite, the war configuration is maximum performanc that can realistically be achieved given the physical resource constraints of the TLP made. While executing 2 threads, no average an efficiency of \$450 is achieved compared to the most configuration, and no average of 4550 is achieved compared to the ideal configuration. While executing 4 threads the efficiency to minimum in groups. The manual being that, the event and of contention between the thready for managers is appropriate by the reduce DLP formets. We developed the data parallel hearkmarks to and estand

from cimulation and existed path analysis.

for data level carallelism. For the cake of our limits we or wat the rationale we analyze DEP behavior in chapter 7 and we include a brief commany bere The DLP kernels cover a large, if not entire, space of data parallel applications as during map of into four broad categories with a total of 12 kers ob .

de ellell is an field optimized discrete estine imaginess estapa ted that not only

integer math. marks is a straight-forward matrix multiplication agreem

radd does rector addition of two 2008-element rectors. All of these bereich

are quite much and are provided to has desprimized and an feel back of taised

ERMING and SPEC CPU2888: We need all 18 of the EEMING heartmarks which are split into five sategories salled: automotive, concerner, networking, office, and orderon. They are all how it loss havel with small working set sizes and instruction flore data. We adjusted the iteration manner of the EEMBC beachmarks to reduce their execution time and hence cimulation time. We need a ratio of of SPEC CPU2000 benchmarks for which the ends and input set circum ade simulation (ractable,

All the election arts were compiled using the TRIPS compiler toolchain which takes C or BORT NA NZZ code and aredo so complete T-BUPS binaries

to build large his orbitely. Sociald, their designic behavior is former of mon-

are a scenier, coal ratios casted in the enemand are work, loaded are desended

conflicts, and control speculation all may significantly and can conserper

mance loom. In 19th of these drawbacks, our results show and write amounts

of concernacy being explained by the case. Since the code quality from our

We conclude from this analysis that the TRIPS microards but are can

we conclude from the analysis that the section is necessarily and the distributed areas beauty, given between with self-seal concerning and appreciate bandending.

Whether the care will be able to exclude LP on full bonds marks, or whether

the compeller will be able to reperate to Michaelb and haded only promits over

qualitativitat are subjects of auguing work in the TRIPS project. Even in,

microb exchange. On complex programs the the SPEC CP COSS beach-

time of a compiler for a new processor is not chart, but we anticipate signifi-

1. Processor Efficiency: Efficiency of the TLP made in accroming re-

come and contention conflicts. We come are the execution of multiple

threads an one single processor in TLP made, to executing each thread halopendently on its own delicated TROPS processor. We measure of

Sciency by comparing performance against two configurations, called

cated processers mem. The first configuration, idealis the defa

ILP made of the processor is which up to eight speculative blocks can

execute charles would still size of of the 1824 recent to station in

the processor. The second configuration, was, stiffice only a quarter of the procession stations in the processor with at most one speculative

A placer the prompte conflicts from the materalian conflicts by opening

threads. While the ideal configuration is the limit performance possible

achieved given the physical resource constraints of the TLP mode. Note

that, compared to the TLP-words, both the ideal and max configuration

are 2 full processors for executing 2 (head) and 4 full processors when

that will read up the hardware. Although the TRIPS compiler and letter comare many TROP to profit up limitations that are correctly height eveloped and incorporated is to the compiler. Prior to completion of those optimizations, the TRIPS compiler performs are will be in adequate because many of the TRIPS

5.3 Results 5.1.1 Microbenchmarks

Table 5.2 skews the perfermance of the TRIPS processor compare

over the Alaka. We computed to exten by computing the number of codes arried to can each a marram on the two simulations. The third column shows the speedup of the hand-generated TRIPS code over that of Alpha, Colons

There are cereal accel for are in this ISA, execution model, and macountries on a Realization there are not in detail a become the common fibble work, and Nacamia a provider a detailed analysis covering many of the eth to a his discension [97]. Need butters is the S.A. has are studied includfaces t optimizations and prelimites optimizations. The different microset protectly and their everheads are the two main features of the microard-line , or that can affect performance and a detail critical path analysis of different microarchitecture exemit (howe the hot tlenech (in the detign).

In this chapter we have focused on demonstrating the potential for the architecture and making the case for this days of ISAs and partitioned microard deduce from a performance standards. The emobile show that the architecture can perform well on a broad chart of programs and can exce hand optimized programs. It seems as our charles point for evaluating polymorphism to the law TROPS made multiplied using polymorphism to matic re-relative amornion activity broad class of any ficulties;

compare execution time to a configuration where such program to ma repumight an it componer with all among a resource dended in extraction and ILP from that single program. In the remainder of this chapter we refer

6.2 Results

We discuss the performance roughs for each of the three saites, as and SPEC CPU2000, EEMBC, as d DLP hemels, individually. Our workland con-tion of random mixes of gragams, all picked from the same coits.

Figure 6.1 through 6.2 show reads for the SPEC CPU2001 onto Figure 6.4 through 6.8 show months for the EEMBC color, and Figures 6.2 through \$3 three received and relative parallel has dissurbly. Tables \$.2 through \$.7 there the experimentary that were executed.

6.2.1 SPEC CPU2111 bruchmarks

Otilization: Figure 6.1 shows the IPC for the 3-Thread and 4-Thread con-Francisco with the workland mayor contact by the difference between IDC in the TLP made and IPC is ILP made. For each program mix, the IPC when executing in TLP made to shown along with the normal IPC when the pro-

The this propose of considerary is welling that directables now that its minings of XP - $x \neq b$. Declarate publications have observed it such a configuration as the D-act $p \neq b$ and $p \neq b$.

concurrency such muchine is exploiting. The dispurity between the compiled and hand-up im hed THIPS code indicates the correct inefficiencies in the The mode of an that for the hand estimized emerges, the TRIPS

(4) there the instruction throughput, financialists per clock or 10°C; of the

three configurations. The ratio of these IPCs do not correlate directly to

Cord of discount desired a labera with two couldests, raiging from 3.3 to 6.3. The opendays over the Alpha core mage from 1.3 to 3.36, who cen a devices on TRIPS became it as almost entirely original burdenard. What little on correscy there is, is mixed out by the Alpha core. The wider reprovides an add biosoft enems, and in stead on the THOP 5 provides performs slightly were because of the black overheads, such as inter-black register forwarding, would have providing discrete two horsesses the TREPS core has exactly double the L2 memory hardwidth that the Alpha don Harr peri-ar appared to (ws), con 2 log is an approximately prefers aftiwe. These conditi-demonstrate the parential of the TMPS core and chow that it is parential to halff a altra-wide is one distributed a resource to efficients union concernance

The compiler reason ted version of these microbes shought do not seeform so well so the hand a plimited vertice. For more is and each the compiler and others to the memory system becomes a significant buttleneck. For the

Chapter 6

Performance Evaluation: TLP

In this chapter we evaluate the performance of polymorphism mechanism implemented in the TROPS are because for TLD. We briefly a utilize the method share a sed for shipping theoretically and then discout the newform and realts. The salemarsh are marken into the entert thread-level samileds min

- . Execution core: Partitioning of the reservation stations in the execuns care between multiple throads. The TMPS prototype chip implements a static partitioning approach in which each threads can utilize up to 256 of the available 1824 recent to a clusters, State each blick rependunt programs can exemite concurrently on the processing
- . Control flow: Polymerph as mechanisms are implemented in the black Bick digk and next block prelicion. The block fleck digk is negmented to cycle between the different program (bread) or they commit their blocks and first state become empty. Next blick prediction to provided Be such thread with a repursive 20-bit philal bittery register for such

For the 2-Thread configuration, on accuracy the IPCs are the same beres the TLP-made and ILP-made. We can clearly see 4 distinct types of helanies. Recall that the main difference to a program's execution entiresment in the TAP 2-Thread configuration compared to the IAP-mode in: II recentree like data tiles, sperand network, and register files.

- L E.P-mode much better than TLP-mode: is it of 40 mice up to mix II, the ILP made of execution provides better processor utiliza-tion than the TLP-made. With only two threads executing, only half the process of the structure stations are used because of the simple partitioning strates s. For her each throughout to execute one so explaine block and one can specialize block only. Specifically 4 programs in this calls, fp/II.kowim, fp/II.koppla, fp/II.kop.uks, and fp/III.coprid dow as almost 2X drop in performance when they execute in such a on Boundies with an effective instruction window of only 1924. The 13
- 1. M.P-mode slightly better than TLP-mode: $\lambda_{\rm BS} \ker \lambda_{\rm BS}$ into , from N through 26 perform slightly better in the SLP-mode. These are mixed the programs have small amounts of SLP and not very good contra

thip. Mowever, for the purpose of this evaluation, they are written in C assum ing a requestial programming model and compiled using the TRIPS toolchain to produce block atomic TROP 5 binaries. No hand optimization or architecture

nect to obtain, then do nice come incicht into the quantilities of TROPS. The

Table 5.4 cheep the performance obtained on the data parallel heach-

mark raite. These applications have ample DEP and are typically coded in

operational ISAs. For example, the graphics become will be coded in the assem-

his happage of the seriex shader or fraction; shader processor is a graphic

thread. The other corage constant in the next-black predictor which in da de the hour de turnet haffen, call tamet ha flen, and the return a differe tiack are thated between all threads. Data storage: The register tiles have support for perferming register remaining only between blocks that belong to one thread. The data till

Control Register [TCR] and Processor Control Register [PCR] that can he and in configure the processor. The PCR register can be set to configure the processoriate a multithreaded mode and the TCR register can be used to set the number of threads that must execute. In this dissertation, we refer to this multithread of made as the TLP made of

the processor, while other publications have used the term. Towarph to

6.1 Methodology

The cycle-action is a limitator, but server described in the apprior, thanter also maded the salemarsh as medications for TLP. We seed this simple to for the recell precented in this chapter. The compilation strategy used and the binaries are identical to what we need for our LEP study described in the

precious displer. All the heartmarks used were compiled using the TRIPS

- sificantly reduce performance. 1. TLP-mode slightly better than SLP-mode: Mixe 11 (from) 10 perform (lightly better in TLP-mads. These are mixed where one ap-plication's performance is secretly limited by the reduced instruction window, where conditions and
- 1. The P-mode much better than IL P-mode: Finally mixed 11th much If perform much being in TLP-mode than on ILP-mode, or arem to execute together. There are mixed where the IPC of both applications is unite low to clast with, and then have not control operators as sifficulty. Instead, the presence of 2 threads, and hence two course of creful and open brice work every or de improves the averall process

The results show her diverse behavior is the 4-Though configuration with the TLP model sing wone for only on spragm maix. The average IPC is see making telt. I and makes from \$27 to \$46. Purplements compared to the of charing recourses between multiple threads quite effectively. With five

show low IPCs. The network processing headmarks perform a significan amount of computation for every network packet, each of which typically conthat of SSE have of data. The company has been dealers have for an evenconcurrence on the an empression markets in marklet, or empression in deep states

compiler toolchain which takes C or FORTRANC? code and produces complete

TRUPS binaries. We adjusted the forming counts of the EEMBC benchmarks

to reduce their execution time and honor simulation time. We need a subject

of SPEC CPU2000 has demarks for which the reduced input of time made

is all configuration: \$466 of next-block predictors storage tables are pre-

sided to each program with resurate 12-bits of clobal bits or desired to each

regram. The SThead configuration and the 2-Thead configuration leave

1/40 h and half of the processor torage monages and tilized, respectively. This

is an artifact of the static resource partitioning decision that was made for

the projective implementation and does not imply the noteman horometha-

We exert a different mixes of programs in both the 2-Th-end configuration and 4-Th-end configuration. A key methodological provides in solution

is what the enforcementation is there for each a study. Previous resurd on

have charified a marama pains different criteria as charamement hebatiar char-

acterized by L2 cache min m to, on im by emisting behavior characterized by

sceful historetry cycle. The workload mix is which the TEP made does were

consider such to Figure in 12 Shake to by Figure is, and the Statement

m control speculation. In the TLP-mode, the contention effects overcome

stilled the proposer when executing 4 threads. When executing 2 threads

recall the TLP-scale has better allianting than the ILP-scale in only half of the program mixes. These recalls copped a more ophistical of partitioning

an em arti cas tido de amene e tilization et il forther when each a could accorde

the benefits of having more notful announcementative work,

distribution (metal) is

are available.

opecific to sing of the course code was performed for these experiments. This

her demark in the last maps are his instead before than the net of migration de

of SEP + the SPC competered \$25 to \$452. One of the resonal for the high

surforms are in that the complier mantly peachter programs with large his day

Low ILP: The three network appropriate hearing are not been up then

marks according to common behavior.

organic of markets in marallel. In the season tight and of version of the org gram (he compiler or the hardware to mashle to much the parallelism that to available a cross such distant regions in the program and the only concurrence that can be mixed to ILP in the dynamic interesting window. In chapter 7 we

mine more on correscy in each contactor,

sarallel architecture to a structure on DLP mechanisms 1.1.1 REMBC and SPEC CPUBLIC henchmarks

we describe on perceptions to that commany the performance of specialized data Tables 5.5 and 5.6 show the performance obtained on the EEEMD and SPEC CPU2000 benchmarks, Mart of the EE MBC beachmarks are ver

Memory intensity: The two minutes proceeding benefit, files of AV, and similar in that they make how you earlike memory system. Although the blick sizes that the compiler can penerate are quite large \$8 and 22], the final PC

during program execution to quite low - around A. Both / B and A.U. have a large

as maker of farem are accessed. Coffering ately, here one the scheduler is an awar

there in tending in rack a way that their our tention for the TEOPS open of

serwork links in low. The reald microbes should chow conduct charter with

High REP: Most of the youghness have high REP with IPCs to high to 63 4

Using data flow graphs and building a large dynamic instruction sequence

through control flow open lating to effective at expering data level parallelism to the hardware. Alternate approaches affectorization or SIMO computation

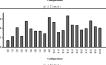
that are meant for DLP computation are likely to perform better. In chapter 7

The EENING handwards mustic share a tread new site for in the reis its for the SPEC CPUSHIB has demarks. We briefly to managine the result and our abordation below. Figure 6.4 and show the IPC comparison to twee ILP and TLP made while maxing 2 threads and 4 threads. Fewer of the benchmarks have lots of parallelling, com one program mixes shown benefit from molitik model execution |mixes I through II in I-Theord on Apont ha and lithrough 28 in 4-Thouse model.

is remarked speedup and efficiency also the excellence similar to the SPEC CPUBLICATION Figure 6.5 and 6.6 show the speedup and efficiency data for the 2-Th-end and 4-Th-end multiprocious. The secure speedup is 27%, and 88% for the 2-Thornal and 4-Thornal multiparations, respectively But hits 2-Thread and 4-Thread configuration care carpitals ply efficient, comtimes, overgreeforming the ideal and max configurations. The reasons they are shie to exceed the performance of stand-alone execution is that, because of rediscorded. The efficiency is the 2-Th wead scoring scatter compared to the idea configuration is NOS, and NOS while compared to the configuration.

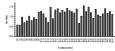
Figure 6.6: TEP-mode performance [collimited] - EEMBC collin

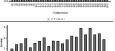




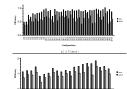


\$14 (kess)













mix is 2.Therefore figuration dmix comber and the record o



Overall the TIP made is quite effective at utilizing the processor of

control to execute multiprogrammed workloads. The polymorphous mecha-

alone a partide an execution window with reduced conculation death for such

or is able to effectively generate north work and is often significantly before

that have a neg control to explain a helanting, counded with contlib in the size out of

and in extract parallelism from different threads.

heart marks in the 4-Thread configuration.

In thirt, the reduced opinion has been thinked processed and the programs

The primary hindrance is performance that we expected was resource

contention for the shared resources between the threads and we find of that

while the recorder content ion did grow classificantly, only in the case of pro-

grams with large an easts of parallelism did it affect performance. We mos-

cacke parts, operand network, and register the. As a result of contention at

all there is not are in the amount of the execution a help eval evaluation of

makke. Table 65 the the percentage of cycles that the execution the are utilized to the measure conflict. We can see that while measure controlled chardy increase, a rightficul difference home only firster SPEC CPU2888

dateli, recensi ikkaz Inizebili, alber satut

Overall the data parallel headmark benefit heads very little from making in TEP made. Figure 62 class the IPC comparison for the data parallel heach marks. Overall only 4 of 40 yra gram mixes y ordems better while musting in the DTh end configuration and by only 6% better on average, as If of 20 program mixes perform better white reading in the 4-Th-end TLP

6.2.1 Data parallel benchmarks

structure parallelies in them, and executing them in TLP mode is to do so a hit of coatestica between the enumers for charel resource like the data cacke, opening network, and projeter file. For the eprograms the efficiency of the TLP mode is also less compared to the pretion (we called . A chosen in Figure 6.3, in the 3-T band configuration, the area preficiency compared to the ideal on figuration to a h 4.26,

and in 74% compared to the max configuration. The results for the 4-Th-end configuration are similar, the areas or efficiency command to the Lifest on the comparison are minimized to the presenting compared to the max on figuration. While one would expect the efficiency of 4-T head configuration to be worse than the 2-Th-end configuration because of more contention, that really is not the care. Between the 2-Th-end and 4-Th-end configurations, the amount of rewhich is why the efficiency remains relatively suchanged. The 2-Th-said TLP made area cally by \$1 be any on our \$12 of the \$124 instanction window cloth in



tire at creating an illusing of a full empression for each program. In terms of implementation complexity the change grouped are units small, on tril lost change is the instruction select logic, register renaming logic, and modifies tions to come table look up in the branch predictor. Going against the spirit of your mary blood, adding TEP copport requires addition of extra architectural

It will be interesting to emiliate in detail the could be of the TLP node, in this cody, we did not realisate two, and little moded workloads with

made. Due to cimulation constmines and constraints of the design, we embsaled a maximum of 4 threads execution. Studying how desily this can be realed to an interesting appearing to explore. Also in this study we did not may

one the power contemption aspects of the TLP mode. While the implication for power caring techniques like clock-guing are not drastically different from the ILP, the heads tick may need to be changed a little compared to the ILP for the execution is another interesting for my direction to explore.

Chapter 7 Data level Parallelism

tion copy lied to a large number of data records. Mictorically, cyclemaclargeted at DLP have been receive architectures like verture processon, contails arrays and SIMD army, on timized for six ole control and explaining the recoloriis the intraction treum and data stream. Such architectures had surner application domains, but more recently by held SIMD-VLIW architectures like the Imagine architecture and multimed h. ISA extensions have here to meted

The main forms of this chapter is a systematic analysis of DLP in the by characteristic (their find amental bakarter in terms of memory behavior, central hebarier, and competition. We then exactly sixtle has but the heldenecks in conventional microarchitectures for DLP processing. Bused on th polymorphon conscious me to oppose data-level parallelon.

we matical either sed for a detailed analysis of DLP workloads and community

Section 7.2 we are tile a detailed characterization of the fundamental helanter of DEP workloads and it Section 72 we evaluate DEP workloads using conventional execution model to determine the but fleneth; that his der DEP execution, in Section 7.4 we use the application characterization to develop a set of Beside microarchitecture mechanisms. Finally, in Section 7.5 we present et it 1820ere menannannen.

7.1. DLP Overview and History

unit-jarant y ngrami an guwa g is ing in ate, itensi ig is inte-illy, as é émis adia g incusseé performis or from baréwire. S peclafized hi ré-ware it commo place in the roo filme graph ke, ilga al process is g, a elwerk pro contact and high-performance orientific companies domains. Madem crash to processors have markly evolved from 21 GFloor in: 431 MHz; in 2112 237 or 260 GP(sp.) |st 650 AM2; in the latest AT | Radeon R360, in late 2006, Hatest on these levels of performs we can conclude that they have at least firsty 32-bit Beating point and a. Software radio for 1G wireless has shaul receivers are later contains there eight nector physica and delines you's performance of en in A Cities. The 1855 Cell property in the Plantation Landon has a Nicolate contemplating estimated and performance of 1 Triang Mt. White there domains of data surable as pling tions have many common characteristic they typically these differences in the types of memory accesses, computation requirements, and controlly charter.

|a| 4 Thread)

More data-parallel and disclared target a colorer of data-parallel pri grams, and have poor upper for application on take of that orbits. We ter architecture, area like efficient areas tien for areasym, with mostle results les effective as programs that require computation across multiple vector ele ments or a creat memory in an analysis stared or irrection for him. SIMD and by tects and a particle on a set of the communication had seem a second to a said. It has beopacker sized as d kes or provide poor capport for application with one dition a execution and data describes branches. MIMD architectures have to piculo hose constructed of course-grain processors and operate an larger cheads of data using the single-grapman, multiple data | 5 PMD| execution madel, with poor opport for Sue-grain synchronization. Emerging applications, such as makilma emakira ankik it rantral kabanian that maning fine emin MIMI

Many data-parallel and lications which consist of common or on that exmits. For example, most real-time graphic processing systems are specialize

ment rely as complete sensiting to execute such leaps.

. Annime less bands: Figure 7. It shows the most constitution for the line

raited to fine-grain MIMD mark iner, since each processing element can

Manufacture of the state of the contract of the state of ment, can make any of these long control templates more simplex. Data-parallel architectures have traditionally implemented conditionals by using

eredication, conditional circums \$21, or vector marks [344]. Piper exercises

int of control, each as a model by a flavormin MMD architecture on a reduce

or eliminate these overheads that conditionals have in highly synchronized

The TMSSISC 64 M DSP chip integrates two specialized units targeted at conrelation exceeds and forward error correction errors into Walk many of there specialized accelerators have been defined to a chigle so more fraction and become one new emerging that consist of multiple programmable data parallel processor that are specialized in different ways. The Stury Equation En sine in the deal two or exhibited median and to was traved for respective process cincolation \$6). The Sony Mandheld Engine in tegrators DSP core, a 2D graph

ice core and an ARM RESC core on a ringle chip, each targeted at a distinct

Figure 6.5: TLP-mode execution efficiency + DLP on the

Design Convergence: Integrating many such specialized DLP care leads to increased design cost and arm, since different to not of processing many to desirated and interrated together. While data lend camillelian is one for to mental area wit that affect the agreement area signified. DEP workings to carled earnigh that a detailed analysis of these workloads is required to under count their behavior.

In this discription, we identify and characterize the application demande of different data partied program cheese. While these cheese have come common affections, namely high competitional intensity and high nom-ory has dwidth, we show that they also have important differences in their memory acono behavior, including control behavior and including clashardware capabilities carring from simple or has remeats, the efficient look up tables, to different execution models, or disp \$1MD or MIMD.

oregalized processes in each application demain.

Dataflow graph abstraction: The TRIPS are countil and coine for data parallel execution with its high fractional and density, efficient ALUALU communication, high memory has dwall be, and technology scalability. Th data flore state 15.3 day in a serviday serveral extensión con abilities, includin esta in By broad execution, that enable a circle bi-forward implementation of th ment of these are leading to Remaining time to the coirt of a democration

Based on the program of this is identified, we propose a set of poly mareken microarchitetami mechanima for an amenting the execution cor-DLP behavier and can be applied to diverse architectures ranging from vector comparison of the performance of these mechanisms to correct hest-of-bree

helianter without adding more dutapath or corage demonst

imati mixis 4-Theat matigamilia - DLP is t

7.2 Application Behavior

Data-earabel workleads can be classified into domain, based on the main and across the different demains. The applications cary from simple computations on image data converting one rator space to another [comprising Mr of harmerican, to complex energetical restinct on network such part of this spectrum: digital signal processing, scientific, network/security applications sategorized by three parts of the architecture they affects same ory, includeling control, and execution core. We then describe our cuite of datas amiliel a moramy and a moral their attributes.

7.2.1 Program Attributes

At an abstract level, data-parallel programs countrief a loop body executing on different party of the input data. In a data parallel architecture this long hady is typically executed on different execution units, a penaling on different parts of memory in parallel. We refer to this long hady as a kernel Typically the iterations of a loop are independent of each other and can execut

execution to the computation of a 2D discrete codes transfers. [DCT] and first blocks of an image, In this case, carallelists can be explained by a more into the The processing of each instance of the hersel is identical and can be performe is a globally syndrouser manner across different computation as des. Amore complex data-parallel computation is a technique or lied skinning which is need for animal in in graphics proceeding. A dynamically carping number of matrix vector multiplies are performed at each polygon vertex in a ID model. The different vertices in the model can be operated upon in parallel, completely

Memory behavior: The memory behavior of data-parallel applications can be described into four different types: [3] regular memory accesses, [2] inregular mem ery accesses, [3] a amed constant scalar operands, and [4] is desc constant operands, in characteristic DLP programs, we are interested in th frequency of a correspond of each of the four types of access to a hersel. The

. A sailer on one end Data-sarable benefit to siculate and from more or in a need to refer to a group of elements on which a single iteration of a of I elements, corresponding to I primary rolar components. Because of the resplacion of these accesses, microarchitecture, on a six elas excesses nethe address sales his and other cresheads as occlased w according memory, by in sing one intending to first one or more full

- ery is a random access fulles similar to conventional sequential pro-grams. One example of such behavior is texture accesses in graphics programs. Unlike regular memory accesses, the averbands of these asorner cannot be amortised by aggregating them. Typical texture data
- . Scalar constants: Many operations in data parallel bernels are real inc mantants that are namedified through the full execution of the hersel each as the mastrant and is consolution filter applied to an image.

 The number of coefficients is often imad and can typically be stored in much be registers rather than memory.
- . And good complete: Many DEP an elications grow to small look on tables with the index determined at realisms. Energytics bemock one cock Stakey table with between 256 and 3024 S-bit entries to calculate and byte for another byte during computation. These accesses on he fre-

corage space, but themseld on each chard width.

Figure 7.3 (how) the three different types of our to the decion possible.

- ciractions with an internal control flow. A decrease to care in a challe rector opension, but the 2D DCT made transformed into this made iteration of these berneh executes in the exact came finding, to these because any well-ration for restor or \$1MD control. Figure 7. in these this type of control has not or with example MGB to YiQ or kernel prendered s.
- · Simples take lays: A dightly more complex type of control behavior of can when the kem of contains loops with static loop bounds. Pipere 7. In shows this type of control behavior with an example exception bemel the hernel is the same and on the executed in a nector or SIMD stelle ach kemely can be narolled at compile time increasing the code size pend thirteely large instruction storage requirements. Architectures that







Figure 7.2: Kernel control behavio

architects res.

Table 7.1 describes a cuite of DEP berneb relacted from four major application domains. Tables 7.2 and 7.3 characterize these between according to the computation, memory and control orders arrowded are lands. The divided by the dataflew crash height; when the loss hound was rapidle, the

Benchmark	Description
	Mal im els a premits
DI 1 + 811	BIG II is 'NIQ con-enter.
ter	A 20 DCT of an SaS knage bleck.
bigbyan 62 m	A 20 high pass filter.
	ncering, receib (3588 by tep ackeb)
210.2	MDS dedren.
Mijsésé	Rifud mel AES purker exceppion.
Blo- 5-1	Ble-S. b parker except that.
	Samble cale
PPT	1924-print complex DVT.
LU Demmysthis	LU decomposition of a dense 1824-1824 m
	reis.
	megraphic premits See(48).
en en en krapte	Back once lighting - 24 ambiest, differ
	specular and employed lighting.
fragmente imple	Barie fingment lighting - bb umblent, differ
	specular and employed lighting.
on intendential	Veries ibuder for a reflective in iface.
fragmente eller ins	Programs skuder readering a reflective sarfa
	tring only maps.
* 10 1 10 4 5 5 5 5 5	A series thader sied for azimation - 21 mm
	the transfermation matrices.
	A fingment chader implementing unicotrop

the monet [in 64-bit words] that such best of reads and writes, the occupacolumn gives the number of irregular memory accesses, and the third and for the sum are valence, describe the supplication reading to the his the bare and the time of the leakup table for indexed contracts, if one is unded. The control column is dicates the number of key item has within the kernel || | and and whether the love house is are radiable arm to keep of instances, in which care the kerneli exhibit data desendent control and operer a fine crain MIME of instructions executed runter from about 188 to 1888 for each instructs. In sector or SIMD architectures, which hick capport for fine grain branching each in tance would execute all MOS instructions, using predication or other exhibit wide custofies in each of the attelliation, demonstrating diversity in the the description of the basis of DEP analysis to the soul this soulies in the fe-

7.3 Microarchitecture analysis

is the arction period we described the back attribute of DLP on hattles echi for data level parallelius, la the sext certies we may these procenter but the seck t bank to program behavior and derive a net of polymorph on

gram behavior and processor but thencek analysis provides wider application coverage and more thesis dily to the resulting architecture than simply creat increased animal to one floors the arrange or like at her architect area—COMD area

We compile the applications coded using a requestful programming model and compiled using the TRIPS compiler to create TRIPS binaries. We simulate these binaries on TRIPS simulator and use bins or it only which can quantify different microarch best are even to that contribute to a grag run's critimily sit, to identify buttlesecht. We madeled a perfect L2 cache to minimi can also determine the maximum speedup possible gives the processor recontrol and compiler, by ten or his all every end microarchitest are events from the critical path and recompating the critical path. We track three groups of microarchited are executed by the critical for the three classes of mechanisms. samely, processor control, execution core, and data corage.

. Fetch: All the black sequencing/prediction, field, and deallecate event age crossed together ander this heading. For DLP workloads, since large repetitive execution is common, uptimized block requencing logic co-rigationally reduce the averband introduced by many of their execution

- . Married and married All accounts to recition and included in this organization reading, writing, regiter renaming, delays to roote operands from the register film to a consumer, and the delays to mate black on typic to the probler file. We as alree recivier accesses as a separate category because DEP y myramı often uccen ibe regiter film repeatedly to read real ime muntants. Six ce this is a read only access, it provides an appartually for optimization, since the register tiles are designed for the common case of the came register being read and written across blocks.
- . Moreover received 111 the microscophistories are not that contribute to the data tiles, delays to restend dresses and raises to the data tiles, and debric to make calmy back to company for leads, in this manifolding five types requires applicated compiler analysis that can determine on time constants and data structure analysis, in addition this was

7.3.2 Applysis

Table 7.4 shows the percentage of the critical path that it upon in such of the three main organic of events. The remain third, and fough columns have the match side to the critical such from fetch, resister the accesses and memory accesses, and the last column shows maximum speedup passible

Benchmark	Per	5 predtg		
	Fret	10000	Memery	1
DSP/cm+en	35.5	4.7	37.0 [35.21]	14.5
DSP/4ct	41.5	4.2	325 (1535)	11.5
DSP/high panistres	15.4	5.7	31.3 (23.54)	5.6
graphic/fragmented ection	12.0	37.4	12.1 [11.55]	2.5
graph kr/fragments knylelight	21.1	31.4	26.4 [19.21]	4.4
graphics/retrestedion	12.1	D.5	32.4 (28.2)	5.4
graphics fremesmaplelight	17.4	0.5	226 (16.52)	4.3
graphics/nerverblasing	25.5	1.7	\$2.1 \$5.35	7.6
2 H = 21 k / h 24 m 2 x h	2.1	33.4	15.5 (21.44)	3.5
101-012/0165	17.1	7.5	12 9.35	31.3
101-112/11/11/11	35.2	1.2	1.5 (41.26)	21.3
releatise/fb	75.7	8.4	11.5 (43.15)	5.1
U.D. SEITEMST	5.5	L1	55.5 (53.5)	34.7
Avenue	25.2	5.5	25.7	11.2

The number within parenthe is in the fourth column, shows the percentage of so what derive the critical cools to ear in the time on what is not address or from

Fetch: Column two shows that an average, the instruction frick related erm is account for door to \$65 of the program cycles. For programs like tipsduel, where the compiler is able to produce only small blocks. It is structed to an average , more than \$550 of the program cycles are decoded to managing inotraction Both. By examining the program concessed and analysing program had actor we determined that rijudual provides an apport and y for concurrence

can be exploited by providing a very fine-grained MIME execution in bitrate

Register processes. The average coalchaling of reciter account to the senas shown in the third column. As expected, programs with few operation na registers, see little of their critical path denoted to register access m. Pag example #1 and LV are duminated by memory access and their register ac-cess contribution are less than L Register accesses became a buttleneck for application that are a large number of matime continuts, which are register

Memory access & Second programs are deminated by the numbers feycle spent in moreovery. This delay includes the constention delays at the matern and he hanks to much the data (de cuche hanks, and mater contention del while musting replies hash turke our samen, intrinsic suche access delays, TLI lanks pro and land-core can flict detection delays.

We can see a correlation between the number of memory accesses t instruction ratio presented in Table 7.2 and the fraction of critical cycles can tributed to by memory account. Now Set, richdest, or buildinging. Ft and AU are all dominated by a large as maker of memory accessor. Recall that the

memory according well. Correspondingly the memory according to be to the critical path nation from 4.6% to over 75%. Forthermore for programs with ry cyclem, as shown by the numbers within parenthesis in the fourth column So sed includi the reacce on the providence Econi performance inconvenients

Speedup: The last colons in Table 7.4 shows the speedup that can b) he physical resources are still the cause \$250-wide instruction window, \$6 wide into a and 125 registers). We are a broad definition of microscal inchr instruction, and the delays incorred as a most of these means is somehead The creeder depiced from this definition of creekend does not account for any natestial change to the coffware model or a normal name and of

The speedup values range from 2.5X to almost 35X, in discoing there are chalikaat microarchitecture overheads while executing DIP yragmus and that the potential improvement from microarchitecture mechanisms tor pried at these combands is quite large. These large parential speedups an out-of-order oppercular person or the the Alpha 2 1264.

7.3.3 Summery

The quantitative analysis and the detailed program characterization those that DEP programs those a set of common at trabules. The quantitative as aly in above that heliding makes architecture mechanisms targeted at the specific at titlentes can provide algalificant improvements. For example, if w reduced all of the feets overhead ; for FFT, a 4X improvement in perform and based on the critical path of microarchitecture events, it is likely that if s erforms are in a revenues; from mad internets and more till be additive. Fig. all by calcing disaging the apparamentar and execution model, it is not this to achiere spiedays beyond what is passible by simply medating micron whi seerbooks. For example, some programs with the grain concurrency, can be dramatically speeded by using decoupled execution between "threads" that

7.4 Data-Parallel Microarchitectural Mechanisms

The program analysis presented in Section 7.2 provided to with in sigh into program behavior and the critical path analysis in the precises section quantified the buttlenecks in the execution cone, but not be non-trad, and mon-ory system. In this section we do only the microarchitecture mechanisms w developed based on these insights. Figure 7.2 shows a block diagram of an abthat microschiletors. We exclude the referenchest mechanisms is ten

er cotto cestali

Table 7.5: Data-yamilely regram at rith one and the set of naiven al microsc chirecteral mechanisms. Mechanisms is narrathesic indicate TRIPS seems

components as shown in the third column. The hot column has the beachis the memory (y) em: [2] a collection. For measured memory (x) by (1) m and in cup part high bandwidth regular memory accorder, and [2] a bardwa managed cacked memory or harden is need to consider efficient investiblemen. ore a crease. The execution care is subspaced with additional in cal overage Maily report sometiments opened according to 400.00 of ware managed local data storage for according indexed named constant Finally examining control behavior, in territor stamps at each ALU is the execution core is added for copy ording than simple begs, and a lead program construct such ALC is added to provide data dependent branching behavior

the resting overhead of MIMD execution degrades performance slight

- . SIMD + arefor engrand arrays (S.O.: The performance of man out that is not be eight to the example). These perform belong the On addition configuration as chosen by the set of 7 a morano in Picture 7.5.
- SEMD + scalar operand + lookup table access |S-O-D| ; S to job proper than the M-D marking.
- MIMD MI: The bureline MIMD configuration degrades performance concept at relative to \$40.40 for all an alteration, even to other alternations This degradation arises because in the MEMD model is a food instruction.

 From such ALC must be consed through the network to reach the memor hiterface. In the armino three SIMD configuration, market aims a black has adulted, a multi-west lead instruction could be placed near and expension is depositently in the MEMER and el, in this includes a si-
- · MEMD + lookup table access [M-D]; The MMD mining with

THIPS: We decelored the THPS and iterate to an implementation of

Figure 7.2: Microard Hed are black diagram

of there abstract measures and specifically in the context of the TRIPS processor. The mechanisms are not implemented in the TRIPS protetype chip and

this study where we explore what mechanism could pass highering beneated

has dwidth receipt memory access and low lates or irreceipt memory accesses

potentially produce speedups up to SX for the DLP programs. We propose

off ware managed cacke and a kard ware managed cacked memory system for

Software managed cache: Figure 7.1 (how) the configuration of the mem-

ory system that provides a high-handwidth access for expolar access pattern.

Portions of the semi-dary-level code should us be reconfigured as a fully sef-

were managed carby SMC i, in this confirmation, the hardware reducement

While we described these mechanisms using the TRIPS processor as

the hardine, they are extremal and applicable to other architectures. The

SMC, come halfer and the L White tractions can be added in a straight forms of

tectures by adding direct channels from the L2-cucket to the finational unit

and an gravesting the pipeline to wakeup instructions dependent on the hado

milariach appart for two den from the LI memory to the rectarregate

the, so high and ware techniques to generate conflict free addresses to different

hanks in memory, in contrast to our approach of pucking all the regular as

accesses in this architecture, the L Scatche memory must be addressable usin

er edal om tier) rather instructions. Most om restional on en miar processes

to recording stations is expension and between and both the intention and opens deer ballon in a mechanism can be applied to provide intention

and upone director. Many DSP processors have implemented genesic related

control locks mad first to freek from a despita struction story halfer. Conven-

he as greated with a local PC and charage to flen to provide a MEMO made

The programme clarks on TRIPS have a construct or correct sades or

To provide MEMO copport local PCs are added and the local ALU

erter-kinning. With heal looping materi, these programs require for

tilling many sum liellem. Dem nie nie berektiming nier data denenden

can thing, the overheads of predicated execution [or candition abve

• Flexibility: The last single har labeled Flexible in Figure 2.5 shows

the harmonic mean of 19 cedays achieved by a Meable and lived are when a subset of mechanisms are combined according to application needs

in a in a fit and it as \$ consent through parter simple behing \$40.

and the rest on M-D). Averaged arran the different application, this Beable dynamic toping presides 55% befor performance over a fine

S configuration, 20% better than fixed S-O and 1% better than a fixe

Table 7.8 shows the results of a mugh companion between the per-

formulate of the configurable TEOPS architecture to published performance

mistir on specialized hardware. Columns I and I show performance, colum

4 shows the performance metrics [which rary], and only madel certific the spirit

chaired hardware. For each of the are lication, we nicked the hest combination

down and violate any microarch field are accomplished, clace the TEMPS process

7.5.5. Communicate americalized architectures

M-D markins

vi en (heir operand) a mire from (he S MC. T he Turns) als arch

statife and install for Librarie memories.

branches in different words to common tight house

7.4.1 Memory system mechanisms

Polymorphism: This discretation introduces and ited and polymorphism the capability to madigate the hardware of roa-time to perform different for tion). Unlike reconfigurable archiver are that synthetize complex logic from primitive fractions, the polymorphism principle to to build course grain re configurable microarchitectural blocks whose finaction can be changed at ma time. We need the TRIPS architecture to the baseline for developing an implementing these polymorphers mechanisms. The TRIPS architecture is med star dwigs with well defined microsoph iteriors blocks and it is technolog realist to device. Therefor persian as a conditionalize starting mobile for teaching mention and marghing. We are used and evaluated medication (largered as three processor on surses: the execution case, control flow unit, and memory

Beselfa: Our performance receive show that the TRUES retreated test or contains IPCs in the cause of 4 to 5, and on a set of block data nameds

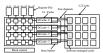


Figure 73: Memory (ye) on mechanisms, Sufferire managed on dis, fact dissipated and describe the firm

acces patterns while providing high handwidth. The DAM engine are not seem that provided here has been a functionally provided. here has been a functional provided to the contract of the con

cheme and tag shocks in these suche banks are disabled. The SMC banks each contain a DMA on give that is explicitly programmed by software. These hanks are excessed to and are fully managed by the programmer or compiler Only the recolar memory agrees with talkally identifiable by the compiler of the SMC, and they also bypass the L bracke since temporal locality is poor Using the data tiles which from the Librarche is also passed in between emanaging calered sy at that ferel becomes a shallenge. The programming abstraction to manager his SAC . Providing on the offware managed caches [referred to as a circum register the or SRP; is a natural configuration to exploit the regular

of execution . While adding each local storage goes against the spirit of poly-

merchian and could demonstrally increase the design ormalexity of vector

as if \$1MD much large, there much identities increase the domain one or they can

This section grown to the compilation strategy, simple tion methodology,

ating and managing the following: [4] performance improvement provided by

each mechanism, [2] benefit from different mechanisms for each application

For the ILP and TLP evaluation (12 dy we need the him-grac sycle no ate simulators, For the evaluation of the DLP mechanisms we need different

infrastructure, primarily because medifying him-procts medelall the mecha-

ners more abouted charletor, which has been described by Desibus [35] as

the GPA simple ter, that models the TROPS processor. This simple term on

28 aprilyria Tanadan)

binaries processed by the LMPACT compilers ad translates in traction in

7.5 Results

operatized architectures.

144

TRIPS-like instruction set, and now a scheduler that has similar hearistics to the TRUPS (cheduler, The different mechanism) were interrated into this simplayer for the performance experiments. Appendix A describer more details

high bandwidth transfer from main memory is to the SRP.

Wide lands: Overhead and latency to access the SMC can be reduce

by using a LW [final multiple word] instruction for reads, As LW instruction

ALC: or metaple recognition stations in the same ALC in a single new in the

the army. To reduce the write part provides, a stage hafter malesces stage

High-bandwidth streaming channels. To deliver these specials at a

fact rate to the execution corp. dedicated channels are provided from the SM

Cade d L1-memory: Irregular memory according to a be efficiently handle

by aring the level-1 cache and that a bunks in the level-2 art can figured at SMC

hanks, in applications such as graph to readering, such a cacking medianism

The branching behavior of data-parallel bernels dictate in struct is a fetc

and making improviments which are: [1] maket of first include in and it are five

ael lastra d'ion la morralles (la lèsa, molling la lastración cade percur and d'yn maic cade acces (power, and | || MIMD praces in grapper for bernel

that exhibit fine grain data dependent branching. To avoid repeatedly fetch-

ing instruction of a large the ALC are enhanced to reason formed in the

for the irregular texture is shapt can provide low latency access \$25.

or and by one ALU fit their modify it was right to realize one disease.

from different ander together before writing them back to the SAMC.

hanks to a corresponding row of ALCs. The army haved design po

anteral partitioning of the cache banks to rows of ALUs.

7.4.2 Instruction Fetch and Control Mechanisms

on this simulation in fine tracture and compares this simulator to bim-proc. All the arrange were band-to ded in a TRUP'S the instruction of these to Michael infrastructure and datasets for a realistic since belon of unio etropic-Allering, we exclude it from all our performance tables and figures. All the thip. The sub-difference being that this THOP's like ISA and the THOP'S IS A is that the file form at for the binaries were different. Hence some instruction cache helantar would be different. Where movible we station in parallel the $kem\,sto\,(a\,m))$ by the instruction stars persons the ALUs. We mean recent its speedups in terms of execution cycles between the buseline and the different marking configurations. The simulations are med that all data was resident in the software managed cash o SMC | or L2 (torage for all application). Exc for h, the datasets of all applications the entirely in the SMC.

7.5.2 Baseline TRIPS performance

Our hardine configuration models the TMP 5 princippeckly with the GPA simulator. We assume each data cache bank à connected to a S4KU SMC bank. The functional sait and cache access latencies are mulipared to mat store Albhu 2 1284. Each under it the processor country of an integer ALC

It sty high or than the typical high PON designs of these specialized process on On the signal processing order, the TRIPS core in the S.O madigam tion, it up to 5 time. Buter than the MPC 7447, with the improvement coming from the 4X higher in ne-width 18 to . 181. The TRUPS are on take much half the number of functional and car the imagine are

For the scientific or derive compare performance to the Tarantols and b tecture. The TRIPS S configuration is store has dwidth limited and about a factor of two worse than the Tamatola architecture. The TROPS yeak memory has dwidth from the process or to the memory system for store to 4 words/cycle for an execution array with 18 execution units, whereas Tamasala allows 32 words/cycle as an execution array with 12 execution units.

For the network proceeding programs we compare perfermance to Cry temastic, a programm able specialized network processor. By exploiting the extensive data level purplishes in network flows, the TRIPS S-0 and S-0-0 configuration y extern an order of magnitude better (has 19 existing back-wars, where the packets are processed sensity [maller nambers in the table for these programs indicates better performance). Cryptomaniac could also patentially exploit concurrency across packet flows, and in fact many network processors do exactly that by providing multiple simple cores on obly and

Table 7.6 (how) the performance of the banding measured in terms of as maker of stellal common tation, operations to taked over cools, and lack tile rethead doubtect has the address compute and load and store instruction Only the DSP programs contain a very high computation throughput, averaging a boar 2 op s/cycle, while all other applications contain low throughputs,

interest multiplier, and an EP Clark had a marking and diskle constitute

is consider iterations reading from a local charage. To efficiently apport data

Instruction revitalization: is the TRUPS on over, the ALU; almade

we arguest the ALC with opposit for re-using instruction mappings for

security detailed of a hop. This mechanism, which we call instruct on

residabation, words as follows: before the start of a bernel, a setup. Neck executes a repeal instruction specifying the maximal map have distribe bernel

which is raced to any exist bandware count register CTR. Then the instruction

of the kernel are mapped to the execution core and execute their first iteration

When the trending completes | listermined by the block control logic), the

pendent bene thing, each ALU is augmented with a local pengram counter

Six on the handle of RIPS per on our is up timized for ILP, converting the data level parallellom in these applications to LEP exists in in efficiencies for DEP programs. For example, loops cannot be to ficiently samiled to provide large enough blocks to efficiently utilize the array of ALUs, and every scalar spenial or memory reference must approve through shaped structure such as the L1 make as 4 the common register Me. Share many DLP yragrams have large demands on these prompted, the finited handwidth presents the and her are from achieving its potential performance.

We program med the graphics kernels for the NVIDLA Quadra FX chi and measured performance on a 2-4 GHz Perform4 hased system. In the we proposed a circumst coffware managed cache memory along with a bardnumber extra signal and high and heating. TREPS and not have dedicated hardware managed lend-1 cades. For the execution care and instruction control wa ure primarily because of the much higher issue width and functional : const. On fragment want pleas the other hand the specialized hardware on performs TRIPS by roughly SX, Although the exact details on the number of stome can be combined in different ways based on application demand and for critical on its (fixed upday + floating upday on the Oracles FX are not are a sweetly in such to a model with a \$1MD and MIMD execution model on to the largers and er office clienal and s. The other graphics proceeding bernels modified in 55 455 befor performance than a fixed yet could be architecarem ore complex | to hig more holds of high, more constants, and data depon-

7.6 Summary

In this display we agree tell a comprehensive treatment of a marrana ceering a large spectrum of the DLP application space, including signal processing, scientific, network/secrity, and restrine graphics application While there may be DLP application on hide these domains, the flor studled in this divienation provide community coverage over the application op ace. We ideal Bed the key memory, control, and computation demands of DLP application and characterized the helperior of the DLP application in its

at hearthing in one case, thus the two we heachmarked, and will perform

at here as well as the other better), and likely poorer.

We then a many of a set of complementary and reput independent decision mechanisms targeted at the memory system, is struct is a control, and execution

There distributed arminous have maked as in construct a Newsde. 3024-bit models with day, a stee Familier processor, which works quite well as a that code can be compled efficiently for this architecture, or that the onecerror will be commedified over with blok-spaller or de on real application Ouce systems are up and maning in the Pall of 2006, a detailed emitration of the capabilities of the TROPS design will help and entand the circuiths and

is terms of VLSI design complexity, the homogeneous approach has effaite administrates, to this discertation, we introduced a principled approach of noing polymorphism to achieve design convergence and have focused on specialize district fluctionality spins an economy of medicalizate drives by block control logic broadcasts a global restraine signal to all the nodes in the execution army s which resets the status hits of the instructions in the recertation stations, priming them for executing another iteration. When the COR register reaches nero, the next hernel's execution commences. To amortize the cost of the global resisulate broadcast delay, blocks are

as miled as march as result in an determined by the sampler of the reservables. stations, is as to reduce the number of resistinctions. Figure 7.4a shows the datapath and control path medification added by this mechanism. The this dedication as past to the presentation (tables) indicate the state this presided for rectalization. In the TRUPS processes, saint instruction probabilization

Local program counters: To uncost the count falls described branching the execution core is configured as a MIMID processing army by adding local PC car the ALC. To simplify the duting ath weathe add a separate All instructions along a from which instructions are finished and executed sequentially. [A slightly more complex, but area off deat implementation is to re-use the local instruction storage already present in the ALCs and not the PC to read this clarage.) Prior to executing homels in a MEMES mode, their instructions are founded into this store by executing a setup block, which capter historicians from memory into this storage and resets the heat PC to seem at every ALU Once this wat philock term is not, the army of ALUs begin executing in MIMID faching. Each and e independently on season itself by fetching from its loca

2.5.1. Configuration of Marhanisms

the application set we examined.

The medianisms described in Section 7.4 can be combined in different

ways according to application requirements to produce as many as 20 d β ferent manifest machine configurations of a single flexible architecture. The

frequency of each type of memory access, the control behavior of the kernels

and the instruction size of hemsels, measured in Table 7.2 and 7.2 determine

the ideal combination of mechanisms on the TRIPS processor, in this disce

tation we focus on three public configurations, shown in Table 72 , that over

and an another managed cache. The SMC has been dether the hafter and

the high speed channels to common instead with the execution core. We describe the free configuration in detail below:

core, that man appart such type of DLP behavior. For the memory cycles

y nord local aperand clarage, local lactus clien cliens carage, a coffware un anap

mir. We found the approach of contemizing the

tage. The approach is this discertation of customizing the architecture to

hat the contamination we propose making different execution models on the

came in hitrate and can be performed after fabrication. When compared to

are likely and reciffic processors in each of the domains, the architecture hads

using the medicalism in this discertation achieves performance in a similar mage, whose assuming for check rate and ALU count. White each applica-

the provide a property and property will be the own do make, and a bare size ifficant

processor described in this dissertation. For example the helpful of STMD an

to me. Potential term that must execute multiple dayon of DLP analog to

the e-grain MIMD execution may deb to a reasonable goal for other DLP unit bec-

The mechanisms that we propose are not strictly limited to the TRIPS

flexit fire to perform well on DLP as elications and side its domain

local corage, and head program constent at each ALU cite. These mecha-

is all the configurations, one memory bank per row is configured to be

COR register is decremented. If the counter has not yet reached need, the instruction stars. The opening stars gold flors are used as read/write registers providing a chapte in-order for th/register-read/execute pipelass. Figure 7.45 shows a schematic of the modified ALU datapath to support such a MIMO model. While this MIMD model has a one time started delay, instruction mentalization in care a rectalization delay between every iteration

Multiple under out he apprepaied impelher to execute oue iteration of a bemelia the MMD makel ereciding a locked wide-tone marking for such itentitus of the hem of using the inter-ALC activate for the grain ALCALC appearance in the continue of the co into multiple desaminable to sed core. As other mode of correction is to excare different bemode on the ALUs, and increasing orders as high element them through dering pipeline can be implemented by partitioning the ALU: among vertex proceeding, carretination, and fragment processing kernels. Since the ALC: on earlike limitations of correct graphics pipelines in which the cortex, mater-

7.43 Execution core mechanisms

Efficient realize up reand as d leid seed realize up rean d access must be rep-parted for data-parallel execution. For large, statically narolled loops, reading culses from the registers for each iteration of the loop to expensive in term of power, register file handwidth, and other overheads of register file acce-

· SIMD machine: Combining offware managed memory system with

as instruction rectalization mechanism greater a hazeline model that is

similar i e S IMD and rector mark inco. In traction recitalization adde th

up port for instruction and control efficiency that make SIMD and recto

machines efficient at DLP. The reservation stations distributed accor-

the tile can be thought of a forming a diviriated vector register the

and the latteration mapped according different tiles from one larg

 \bullet NEMD + scalar operand access. This haveless matrice |S| can be

accurated with operand on taling in the create the S.O markine. This

on Ego ration optimizes the injection of radios into the execution array.

• SEMED + scalar operand + lookup table access: The S-O-D ma

discussion local Ed data corrupt to each ALC of the \$40 marking, This

na figa mation departs the mass from the spirit of polymorphism as it ad

. MEMO: Combining the memory update with local PCs creater a has

MEMID + lookup table access: Addition of heal Ed data change

of DLP heliation needs to be supported, the Book dity can be sacrified for

simplicity by implementing a subset of the machanisms on a fixed architectur.

marter to one is no configuration counter the M-D machine

to a MIMID marking [M], in addition the control legic at the ALU; in

augusent ed to requence in it en el icur (no tend of execution in y nee dat affice

Using the memory cyclem for indexed scalar operands in one cache access over heads and on some cache handwidth. Two medianisms implemented at the

Operand revitalization: This mechanism reason register raises ance the have howe received at an ALU, providing permittent register-the like starage at each reservation station. Successive iterations of the loop rease the values from ment appeared residualities we add claims bits to the reservation coloring, as drew is Figure 7 de.

All data storage: A college managed Life data (brane at each ALU provide rapport for indexed contar constants is ne example to the lookey to bles and written to the morrothic stations. The index to mad the LR data store is provided by the ALCs and the results are written but into the legal register.

Table 7.5 is minurized the program of the stee that we identified in our program districted and and map there to the mechanisms we deectived above. The first column of Table 7.5 little there at this atom. The record column & to the grap and machanism targeted at different microard-lists

Figure 7.5: Speedup using different mechanisms, relative to has disc archive-ture. Pengrama group of by hort machine configuration.

7.5.4 Performance Evaluation

Figure 7.3 shows the application speedups obtained by these differ ent mank his one fire notions relative to the hardine. The following sum one ha preferred the S , times preferred the S-O and for cyreferred M-D configuration

. SIMD execution |S|: f| and h are rector-oriented but damain and require high momenty handwidth and high instruction forth man. Comy ared to the har else a fine offild sympley to achieved because of the higher ALU orithanton and higher memory hand width of the S configuration

Chapter 8

Conclusions

ing applications with betternyment computation and card technology lim tations of power, wire-feloy, and process cartains. Designing multiple applicalles results arounder or rescalated applications introduces decide comn of for helemorason and licating a combat a recent of complexity, and achiev namies of sale. In this discretation, we introduce undistribution in white paration by expecting different granulation of purallelians. The hand idea is polymorphical in to configure course grain microarchitecture blocks to pr eide an ndaptive and Mexible processor scholmite. Technolo achieved with scalable and modelar microarchitecture blocks.

in this discentiation, we identified the cranslative of similarity as the fundamental difference between application classes and use it categorise app

and data level parallelam. To provide and best and copport action all these types of parallelism, we propose architectural polymorphism drives by three of mechanisms, and microarchitectural reconfiguration at a course granularity

We not be data flow on the antifered about the form of the breeze action then three types of parallelism. We have done EDGE 15As, a class of 15 As, as an and became include for efficiently expressing parallelism for healting tech using presisting architectures. All programs are expressed in terms of dataflew graphs and directly mapped to the hardware which is partitioned depending on the grant britt of camillelons.

EDGE ISA e: EDGE IS As excelle description directly in the opening his sary and employ a block atomic execution model. The explicit dependence exceding efficiently express on the dataffew graph [as d beace on corrects], sh ciating the used for complex hardware to rediscover parallelism. The block at amic execution model, color the grandardy of execution and state man agences in the hardware and eliminates instruction level evertends, for tend of tracking and hed and change at an instruction level which leads to a lot of EDGE with a hearily paritioned and distributed makes architecture implementation to achieve technology realability. The two most sign Meant Deats to of the THIPS microard had are are it there is no mit based and mudaler design

of its distributed protocols. On a set of band-optimized formula, the process

On the EEMBC and SPEC CPU2000 benchmarks, with compiler penerate code we see IPCs is the range of \$3 to \$3, with an average IPC of \$3 for the PERMIC on the and the fire the SPEC CREATER only the head on the lead microb exchanges, the TMPS process or is up to 4 times better than an Alphu 13244. With compiler generated or de for large outsticked to be demarks like the E EARL Can dispersion Consists heartmarks, the T RUPS processor performs were that the Alpha 2 256 is most care.

Hand up timized version (of the EEMS C beachmarks perform up to 5 times better than the Alpha 2004 and many beachmarks share several of the meny temperature of the meny timing times. Some of these hand up timing times, which is choice better instruction merging, land /store dependence elimination through better regio There are correctly has duptinalmities and not yet in the compiler for two rea man. If the heart tier are led for these sector but has care from heart mark to beachmark and are at times haved an examining microarchitecture orbits yath meats, and $2\|$ our cycle accurate simulation are too slow and we expect to and entered more of the hardware's behavior on complex to debutes when we have offices. As the complier materies and we develop a better and ential diag of the headstin, we exceed maps of these antimization to be

or compiler and the compiler generated code performance to improve The pulymery has a mechanisms proposed in this discertains are effec-

ting 4 threads on a single processor, sign Wennetty higher tenets of process or stillenties are over, IPCs are in the range of \$2 to 3.8 for an application mix consists of FEMBC and SPEC CRITISE workloads. Compared to a arongo (PC of Line 4.14, there application mixed have much higher (PC \rightarrow 2.2 when remaining with 2 application concerns thy, and 3.3 when making with 4. ann lications.

> When executing programs with data level parallelism, compared to an execution model of extracting only LEP in the TREPS processor, the DEP meckasions year its average speedups of 5.6 across a set of DEP workloads The speedup procised by the individual mechanisms range from 11 o 152. The pulpmany has smeak anisms can block of MRPS and doctors to match the perfecman or of course level arrangement consisted at all flowers consecuted DEP worklands to match the performance of hert-of-treel DSP chips, graphics ships, and red or dale on workloads to evidlent for each.

and at the time of this discentation, we expect systems back at the end of full 2006, in 2002 we carried with premating receip has all as high level

weaknesses of the system and the technology and answer these questions.

as a homogeneous computing substrate to satisfy the computation needs of fiture upplication (that are likely to have between eason computation needs. We believe this upproach is reperior to building a betweeneessering close comproof of multiple operation) processes. For designers who with to half s alemaraha as se stema, the three main shallenger are VLSI design complexity of ware complexity, and resim sluggeous emints of performance, power, area and reliability-all of which imposite into market constraints.

example, we demonstrated a dear instruction control but he set on adentifi comparing beauty like /ft and AU decomposition by program analysis. On critical path analysis showed that more than half the program sycles are spent is fielding the processor are with instruction. This material calculus baseometric that caudial fielded instruction in her consid in the processor. care without introducing any new corner conclusion. The number of mechaalone to cover L.P. T.L.P. and DLP are few in an other and into hones; in a there would be compler than building multiple cores on chip.

As an illustrative case study, we compare the Tamatala processor is a betergeness design, to TMP5. The Tarantaln processor comprior a 22 wife rector core and a high performance out of order EVS core interested an a cincle skip. Mal returns the solemen has a THIP'S do in a laof polymorphism in the TEOPS are in during reason in the processor core, the memory or dem, as 4 the register files.

• In the THOPS approach there is significant series and rece in datapath design since one core is replicated instead of having to design two

design cost. Tama tala provides global synchronization between the 4.5ferent elector have with partitioned centur registers and optimized a c

make. In THEP 5, we simplified the memory system and instead provide the L2 cache hashs as sensith put memories. While the Tamats has prouck to allow rector access to the L2 cache in do des a complex could the address repenting others to maximize has dwidth [236], to create contributed memories at each TRIPS memory tile, the tast checks are mply dashled. The IBM Cell processor uses a similar approach to management sy.

. Unlike To match which contains sector register files which a set in his med and written for every instruction, we showed first did not imple ment in the prototype obly | polymorphon mechanisms that can use the mornation stations closely integrated with each ALU to create sector project the the helanier with operator by carrier on ability.

number of constants and other image for data structures prefern posts; is the TMP 5 approach, as application can chose to continue using the Literation for each irregular account, while using the collegence managed memory for high-handwidth regular memory accesses.

and then arthere restlication storage. The mechanisms are by definition on related and ma he used reparately or to pelier. For example, the fire DLI number of specialized designs to known and the certification methodology fo them is well defined. The certification complexity of such a heterogramma deica compared to a polymerol our device is an interesting operator to a direcwhile deciding an which relation to pick. While this discontains leaves th quarties upon, we do not nice it so no introduction or hard challenge. The TRIPS projetype this implements a limited amount of such polymorphos processor, and the configuration of the memory tiles as scratchy of memorie We satisfied these markes is much as demanded of the a recessor through randomized through mademization. The level of coverage achieved in this process had to be believe that the configurates is not much more difficult than configuramali is in het grangen enn en en.

5.2.2 Software complexity

Designing , decodoping, and compiling applications with heterogeness committee and amount dallered for the entire coffeens that Whe

211 215 the target to a heterogeneous processor with multiple specialized processor one many decide which are Scatten to here to ited for which represents. When must decide as the configuration of the different microarchitectural blocks to compiling for each homogeneous systems more complex than compiling for

Some offering to be in to on any manual to both colonic, sample the termining application behavior, determining the groundrity of the parallelism and mapping of processor mapubility to the application. On the other hand to me toffware deck into are different because the two contents are to radically different. Examples in this title following: If white come this and designing for beterogeneous systems knowing the application mix is important, 2 grating applications from an expectalized core to another can pase a challenge since each case is toned to a specific type of application, and 2) application homogeneous systems y non-different shallenger: 1) determining the mapping

in this directation, we did not address this software complexity shall leage. We said the well that a mong a set of partible configuration, there was a natural and preferred configuration for some applications. We did not address how the compiler or reactime cyclem can determine these properties or the

are predicated, each that during program execution exactly one branch in other than't predicate to each left.

- C. Operand network: The THPS simulator models the exact severe network protocol by modelling the control-packet and data-packet pro-tocols of the network. The GPA simulator simply has a communication
- 5. Fetch, commit, and flush networks: The GPA simulator does not model the fetch, commit, and this networks and instead near thord dela-te model their behavior.
- 5. Memory system: The GPA simulator simulator the distributed data to and the LSQ high by madeling 4 parts in a centralized cache which are all equidicisal from the left edge of the processor wee. As a resul only the horizontal moting delays are accounted for, in the case of althe leads in a program going to one single data tile, the GPA (marinto lating a data tile with 4 poets and 4 special distribute links

To summarize, the GPA simulator models to me microarch decisive blocks at a high level of abstraction which could result in over estimating the performance. Secondly the risker IS A med by the IMPACT compiler allows it to

- bl GPGPU: www.com.com.
- and Donald Years. The MIT Alemie marking architecture and per femance. In 25 Ct 155: Proceedings of first find some clinterest on a gregorium en Compa le orditation quigni 2-43, 28 55.
- [6] Viko Agarwal, M. S. Britkikesk, Stephen W. Keckler, and Dong Burger Clock mixers, ip c: The end of the mod for conventional mixing more on la Pracocinea e/ 5 e 175 dom a l/o tem al ma l/o messiam en Cen estr
- [6] Kuri Abdey, Reality engine graphics, in Proc. of Sci. 28 Ann. Conf. m Consider Gradia, rape 38-18, 381.
- 81 David St. Alberte i States States has eater. Steen G. Dreede, San Grigorio Maghio, Michael L. Scott, Greg Semerara, Pradip Boos, Alper Bayektorangh, Peter W. Cook, and Study E. Schatter. Dynami

derignen skanne in halfding heterogeneens systems av homogeneens systems years, several application specific compilers have been proposed to deal with enwise a necour our stexts. Application specific constitution that it awar of groups, appearing our enterform course commission, FFTW i recent example include, PLAME [69] and ATLAS [60] targeted at times sigebra, SPIRAL [23] which not a dynamic programming approach to ap imire the compilation of OSP routher, and the Armstway compiler mean for scientific computing librarie \$ \$, Programming hazang effort in du de Streamit [57] targeted at streaming and multimed is programs, Cg targeted a recensive \$100 Shapers to recens of a comment are received \$10 or generate optimized parallel code [8].

The common characteristics of all these effects are the following: a | as on demina ding of application behavior at an algorithmic level, by important proportion of the microarchitecture are exposed to coffesive layers, of coarci many and other program properties are expensed through the dia goag elec-ter the complex or hardware to a strongly bardened;

While not related to these domain specific compilation and happage

bemain and in the DEP study in shapter?. They were compiled using the TRIPS compiler for the TRIPS (implained at the Trimorus IMPACT compiler for the GPA simulator. The code and instruction many for each simple to and the last two columns show the ratio of cycle and t in the class of the TRLPS classifier to the GPA classic or. The notation T/Gdeparter ratio of FRIPS to GPA.

perents more compact code than the TEOPS compiler which contributes to

The GPA charleter acceptalmates performance by anywhere between L4X to L5X, and an average over-solution performance by 1X compare to the TMIPS simply or Same of this performance difference is a result of contri- le TRIPS time la reparente su aven pe L4X mare intra citata. Ti remain der of the needermance difference is a result of the about dis a error in

To fear a set the contributions from the compiler and contributions the comparious of the two simulation on these bernet. For the GPU simula ters these kernels were compiled using the the Trimmum SMPACT compiler

approaches, the compilation of ming pile the IBM CELL processer shows some of the exhanact epit is an d has successfully employed technique. The companies scalar codes on SIMD units, an ismatic peneration of SIMD codes, and data and orde partition is parrow the multiple processor core to generate high qualtre on de 1821. With respect to better consumer and firstless models and the increasing capability of processors, we believe the lessons of such compiler and imaging of effects will grow in importance and most be used to address the software com-

5.2.1 Technology constraints

This dissertation has feed only as realisating the performance of pulymorphism and the TRIPS architecture. Other technology constraint include arm, power, and increasingly reductibly. We have not quantitatively addressed computions to other design with respect to those constraints. Clearly, a specialized processor will be more area and power efficient, but how much better liable than programmable processors. Studying polymorphism from a power arm, and reliability per sective is an exciting arms of research consist with the others complexity in no.

8.3 A Brave New World of Polymorphous Architectures

morphism is a sate mil design convergence colonies for follore and is tectures that must provide massive compostation owner and support for her matica aceds. A partitioned design leads inclinationally halfilling a reals ble and modular microarchited are with one correct, expressed

This divineration opens by two lessel arms of fitters work:

- 1. Complies for solvmention. Expense a microard her an execute solvmay him techniques to the compiler introduces several challenges; if which microscolifecture medications to expose to the software layer, 2 how to explore these mechanisms, 2) how to determine and choidy proomm behavior, and 40 how to as temptimily man emerge behavior to
- polymorphism to improve performance. These principles of polymor which may be used for other objective like II achieving different level and I | improving reliability. In a more general cents, a comprehensive a nativity of polymers bloss with respect to all technology constraints will

In this discension we developed and evaluated the idea of polymor shipm and a many of a set of mechanisms to meted at connection all on aniartier of parallels to . ILP, TLP, and DLP. One application of th directables to to see the principle here to determine what mechanism are required when the need is to coppure only a specific set of applications. The plications between easity challengs, from tomes tallian Antions that plague the realing of conventional makes architectures, and the technology limitation of power, wire-delay, and process radiation present significant challenges to the rreary explicity exceded in them, and the principles of tiled design with well defined migrouphly extend privately area and in this divisitation appeals a punick a relation. We furnes serond of these elements is micros most on

down not determine the conclusion of the DLP condy which not the GPA

Appendix A

tsim-proc and GPA simulator comparison

In this dissertation we used two simulation for vary offices ance evaluating. One is him errorwhich is a detailed coole-level simulater; has madels th IPS per cere er ar a much muse detailed level than higher-level rimulation SimpleScalar [25]. Our performance calidation effect thoused that perfor makes much from him-proc were as average within 18% of three shinks from the RT L-level simulator, across a large number of craft of and maxism i generated test programs. Because this models the hardware at such a d lend it is not very extensible and we need a more abstract classiator called the GPA simulator for our DLP study in shapper 7. This climitation is with under generated by the IMPACT compiler and translates instruction into a TMPS-like instruction format and note a ched alor that has similar benefit to to the TRIPS tcheduler, in this section we compare these two simulation and describe the difference between the two

The quantitative conductor of this study is that the GPA constance in the worst case over-estimates performance by IX compared to the california THEPS classifier and it is a average within 2X of this validated classifier The poor code quality from the TEOPS compiler and the abstraction error

contribute moghly in equal measure to this over estimation.

The main differences between the two simulation in do de-

- 1. ISA: The GPA simulator and the IMPACT compiler whose hotes: tions are different from the TRIPS ISA. Specifically the implementation of predication in IMPACT which includes general to a of complementary militates and associated associates \$550 is much different from the en TRIPS is typically higher.
- 2. Compiler quality: The IMPACT compiler is a replicate of and hear piler are cometime a factor of two less than the TRIPS compiler.
- 1. Control flow: The control flow in cleans to be a the GPA sinch made in certal order to the taken branch and the architecture change affected by inclinations beyond it are conceded ont. Since this is a high level cimulator we do not model the exact mechanisms by which tion are processed for mandling out such execution and all branches

2.14

Appendix B IPC reduction from speculation depti

the ST bend configuration, where a chiple program to ros in the TLP-and of the amoretes. As a result, the speculation death of the around it reduced



whose the least that	12	1.1	9.2
extensite sakelit	17	11	- 11
extensive multiple	13	1.7	12
estructura (1904)	13	10	1.0
estending outstander	13	1.1	131
office of the	12	- 11	17.
remark the Mil	- 13	- 17	17.
extensitive removals	13	127	17.
extensive effect	- 13	- 67	15
antomics and	- 13	- 10	18
and a married of the party is	13	- 17	23.
cannon ton	13	3	21
extractive left)	12	1,5	23,
anticontine operation	1.3	6.7	26.
actionates (block)	1.5	ž.	37.
All or health	- 12	1,1	27.
artmarking myl	1.0	8.7	25.
estrenative (415-91)	1.5	1.4	31
admidter herelet.	9.5	8.4	33.
office (Disell)	1.5	1,3	11.
	1.6	1.4	33.
antenutive alimen	9.5	6.7	33,
antomotive hitmayli	1.3	8,5	16.
artenting pitter	1.2	8.7	16.
1001103 2011011	9.5	6.7	16.
automatics Materia	1.0	8,5	33.
office baint2	1.2	8.7	41.
reserves followith	1.5	1,1	45.
textiti (111-ed)	1.5	8,5	54
(66) D D H	2.5	1.1	51.

Brachwark				
	ILP-met.	I-THEEM		
oleman III	8.7	1.3	+13.3	
antennik rijaduri	1,1	1.3	5,8	
antwork blooks	1.2	1.2	35.5	
nimon fi	1.4	12	2.7	
graphic longerentteretten	1.5	1.3	2.5	
graybbr emtearhaptelight	2.4	13	94,3	
results) did	- 0	1,5	55.	
graphics longerestifftylelight	2.4	1.0	56.6	
graphic contended the	1.7	- 13	8.3	
graphs certes blacks	41	1.4	55.5	
eradi bigkyaniller	5,5	3.1	3.3	
arteack and b	8,5	9.2	74.2	
result i succest	5.1			

Bibliography

- Id Aleks Architecture Bandhook, Venius 3, October 1996.
- [6] Jeff Andrews and Nick Baker. Marx 160 system architecture. JEEE Marc 16 (1911-27 1116)

16 12 1:45-45, 2111.

[8] Micro Anneston, Emmanuel A. Arnold, Thomas Greek, B. T. Konj Mester S. Lam, One: Menticingle, and Jon A. Webb. The Wary Comsales. Architecture, implementation, and Performance. 26 (4): 521-

Fl D.B. Allis and, Selective cache ways: On-ten and cache more resulting

[8] Rater Alvence, David Callabau, David Commisse, Scho Robbers

ilias. In Proceedings of the Unit International Symposium on Microso dillochro, pages 248 –255, December 1855.

Allas Percentels, and Burnes Smith. The Tera companies contemp. In

- 331. December 387.
- [2] Arrivé and R. S. Nakit. Executing a program on the MIT Tagged Takes Dataflew Ard Loci and HEEF Transactions on Community, 19 11:110-
- [B] A. Amin i and K.P. Gordeler. The Obvion mer, 1995

[36] IL its Bakar and Stracks Manne. Power and energy relaction to pipeline to bo city. In Proceedings of the 18th International Symposium m Computer Architecture races 100-415, 2011.

whereas for the TRIPS simple for these binaries were hear in hand on timined

upo the TRIPS is tractica count to 8.5 times the Trimana is tractica count,

whereas on compiler renembed orded two L4X to first only 2 heroeb have

which likely matches the code quality generated by the Trimana compiler

for the GPA timelator, patentially created a chaption in which the difference

hetween the two simulation can know on it primarily the microarchitector

no dellag. New, or average, the GPA simulator over-estimates performance

The month from the elem controlled experiments, lend as to conclude

The land optimization referes instruction countrigationally—a area

ments, grip I and many J. Using each optic

carting from compiler penerated code.

- m Computer Architecture, pages 202-4-12, 2016.
- December 2005.
- [2] L.A. Barron, K. Gharacharles, R. McNamara, A. New 12th, S. Qui or 176 Anternational Superstant on Computer Artificture pages 252-
- [B] V. Burngarts, P. May, A. Nfelst, M. Verback, and M. Weishardt. PACT XPP -A Sel-Recorderable Data Processing Archite m deligration, Jane 2012.
- [20] G. Barm carrier, A. Arer, D.E. Bern boldt, A. Billiterra, V. Charrello dayappan, and A. Sibiryakov. Symbolic of High-Performance Parallel

ing: Data-flew hased speculative parallelisation of methods in sequential programs. In Proceedings of the 18rd Anna at international opinional contents to the contents of the

- [36] Max Barna. MP Come for Bandheld Appr. Mor spreamer Aspert,
- 201, June 2011.
- Is the least and Confirmed in Engineering of Assembly arable by time

D. Co diero, X. Gao, R.J. Hamico, S. Himio, S. Rrida amountly, S. Rrida an, C. Lam, Q. Lu, M. Norijea, R.M. Pilver, J. Ramasa jan, P. So-

Programs for a Class of Al Initia Quantum Chemistry Michels. Procontinues of the CENES, \$1 β (27% \rightarrow \$2, 2005).

- [8] Salaman and Ralabriahana and Garried and Sala Program deposit in lawon Swine 1881 save N=24
 - $[12] \ S. \ Herker, \, H. \ Cake, \, G. \, Cax, \, S. \, Gleaner, \, T. \, Green, \, H. \, T. \, Keng, \, M. \, Lam,$ J. Cring dd. and J. Webb. Bram: An interrated edition to high-rane

 - Mill has Book. These sorrelled rames starting a companion is and some. In the order
 - Mil Ing Buck, Tim Poles, Daniel Burn, Jersey Stroman, Karney Par

[10] Tom Black, The Margar MP - Lambington, in Proc. of IEEE Comp-

- [t,t] June than Blow. Game Decemperate Barder Than You Think. ACM
- B. Moore, C. Peteron, J. Pierer, L. Ruskin, P. S. Torne, J. Satton.
- 115, November 1866. [23] Shokhar Y. Bucker. Designing reliable systems from na reliable comp nears: The challenges of transists regulability and describation. ARRE An era, 25 (6): 08-16, 2 115.
- [8] V. Bere and J. Wallington. Chemps: A reconfigurable data-flowing for tides proceeding. IEEE Transactions on Circuits and Systems for Date Zedanden, 1 (1:38 8-98, 1981).
- halian, Mike Henrica, and Par Haurakan. Brook for GPU: cream 226

21 (1):777-786 (2014)

- \$77 Mile at Burdin, Gigleb Ven in turnmant, Till enter Chelens, and Seth Conen Galdirela. Spatial computation. In International Conference on Ar-ditectural Support for Programming Languages and Operating Systems, pager 18-35, October 2014.
- 281 Dans Barner and Told M. Aprila. The simple calculated out region 2.6. Technical Report 1342, Computer Science Department, University of Wiccomia, June 2017.
- [39] Dong Burger et al. Scaling to the end of clicks with EDGE architecis no. JESE Computer, 17 (7):44-45 , July 2004.
- [10] Michael K. Chen, Kine Peng Li, Ruigi Lina, Jacon M. Liu, Linia Lin compiled a closely are lication; while making one of group mains. In PLDI '15: Procedure of he 1115 ACM SIGPLAN on france in Pr Pero, 2115.
- hal The Connection Machine CAS2 Technical Summary, April 1987.
- 32 K. Cana, X. Ches. S. Karburda, K. S. McKinler, and D. Rowse. A

Programm in Language and Operating Systems (ASPLOS), 2006. 321 S.A. Constant. Colorest to exists and existences. ASSA dense of

- [14] NVIDIA Corp. NVIDIA GPU programming pole, vL1.3, Nov.
- [33] Alirina Cristal, Oliverio J. Santana, Mater Valera, and June F. Martines Toward bile-in the clien a measurem. A CM. Transport may an Ar & Hechrie
- md Cede Optimization, 1|4|:355-417, December 210 [66] Dier if E. Culler, Annen g Sub, Maus Erik Schunser, Thursten von Kicken and John Wavergook. Pine-grain parallelism with minimal hardware manuel: A complete controlled threaded abound machine. In Proceed-
- hal William I Diale Pairiet Barrier Marian Erec Timorie I Kaleki Pennis Labout, Jung Ho Alia, Nawa Jayarena, Upul J. Kayasi, Ali histel Day, Japan't Germanije, and Inc Book. Merrimor: Super companies with Streams, in The Proceeding of the 2003 Antomation of n/orance for High Parlarm on as Computing, Note orking, Storage, on Anapois SC 13, November 2012

[88] J. Donaio and D. Manano. A preliminary architecture for a h they processes. In Proceedings of the Ind Amnual Symposium on Commier Architecture pages 12 6-02, January 1975.

- of Computer Sciences, December 2005. Hell Railance skin Decition, Danie Harrier, and Stephen W. Kenkley, Measure-
- af fieldfi Anmelfotenet met Spropoisor in Corps te Architectere pages 266-427, July 2013. N.S. Krait District of Protect K. Date: Non-Harborn or and House Altiere extration to powerpe accelerates media processing
- IEEE M e s, 21 p ≤ 1-41, 2111. N2 Practic Daties. Recognition, Mining and Spathers Mores Companies to
- the Em of Tem. Andel Jechnoles Massains Polyman 2005. [11] A. Eichenberger, K. O'Brien, K. O'Brien, P. Wu, T. Chen, P. Oden D. Prener, J. Shepherd, H. So., Z. Sara, A Wang, T. Zhong, P. Zhoo, and M. Godwins. Onlinizing Compiler for the Cell Property. In Property
- ings of S. e. U.S. Anterna S. mal. Conference on Parallel Architectures and Completion Techniques, pages 65—76, September 2005.

[44] Embelded Meraphers of Burdmark Convenient, EEMBG 2011. 225

[66] Raper Espain, Maior Valera, and J. E. Smith. On inforder vector and a tectures. In Press, of the 18th Ann. In Playings, on Microsophischer pager 201-27, December 207. [47] M.C. Seja sweki et al. Oceaniew of the Texas Becombigurable Army Computer, in AFIPS Conference Proceedings, pages 621-642, 1866. Mil Bandon Demonts and Mark I. Klauer. The Co. Tatarial. Addition-

We by Pat Ithing Company, 2013.

[65] Ruger Ropaus, Peterio Arlanaus, Just Rmer, Stephen Petin, Julio Gag

Super Grammer, from Hermanies, Toni Jane, Good Lawney, Markey

Martin a, and Andre Senter. Tamatistic A certar extension to the alpha

architecture. In Proceedings of The 15th Internal and Symposium of Computer Architecture, pages 25 1-10 2, May 2012.

- [49] Brins Fields, Shai Rabin, and Rasthian Budik. Preming percessor esticie via chi kal-sath aredictica. In Proceedings of the 18th John o Anism at mail Summarian on Committee Architecture sugar 74-65. Jul.
- As free Chara, Years and Greening, and What S. Lee, The M-Machine multicomputer. In MICSO 25: Procedup of he25th mem alinto k multiproposium on Microsoft lichten, pages 166—265, 1855.

fix 1.8 ft. In term at an all Symposium on After our distactor 4 pages 224-225 December 2556

- 521 January A. Pinter, Paula Parat socks, and Gineers e Dendi, Carriers-92 processors: Letting applications define architectures. In Pres. of the 15h Ann. Int'l Symp. on Microsoft ledars, page 124-225, December
- Team, Company Vol. 0-23, pp. 848, 3872. [64] Daniele Pelegrani and Antonio Germiler. Energy-effective to be legi
- dileter, pape 231-235, 2013. [5] Matter Prips. A fact Poster (mattern compiler. in PLIN 1999) Proceedings of the ACM SIGPLAN 1889 conformer on Proceedings house of dates and involve misting, same \$3-30. ACMP on . 1991
- Matt Mes, and Reed Taylor. Piperends: A reconfigurable and iterators and compiler ORDE Committee 1140-21422 April 1141 57 Michael Gentes, William Thies, Michael Sarrom arck, Javer Liu, All S
- Mei, Christopher Leger, Andrew A. Lamb, Jersey Wong, Henry Hoff

- not Divid T. Mars and Samue Americans & Street Compiler to t mal Conform or in Architech ral Support for Programming Longs up a mi Operating Systems, pages 25 1440, October 2012. [8] Paul Grat, Charglys Kim, Bater McDonald, Stealer W Newton
- and Dong Storger, Implementation and Evaluation of On-Chip Network Architectures. In Proceedings of the Lift International Conference on Comparing Resign name Tri-TT Detail or 1985 [31] John A. Grando, Fred G. Garraroux, Greg M. Henry, and Bahler A. can de Geijn. FLAME: Fermal Linear Algebra Methods Environment.
- ACM 71 may act may m. M athem at m.15 of to m.; 27 (4):422-433, 2013. [61] Raje Gapia, A flar-grained MIMO architecture based anno resolver channels, in Proceedings of the Live Amount Workshop and Symposium on Manageogramming and Manageditecture, page 25–27, 2008.
- Not have Garrer and Calvin Lin. Househour: a committee for exclusion to remaining of reference his ratios. Proceedings of the IEEE 81 p | :142-417, 2105 .
- [62] Right S. Make en and Annuay Chaptu. The Design and Analysis of a Cuch Architecture for Texture Mapping, in Proc. of the Light Arth Sympo on Computer Architecture, page 185-185, June 1857.

- hel Tan R. Halbill. PicaChip makes a Big MAC. Microgrammer Aspe F | B |: F = B, October 2011.
- [63] Tom R. Halffall. Clearly and Mile Design Targets. Milespe score Appr. 18 (1:38-2) June 2014. [6] Tom B. Ballidi. Buy her at Silcon Mars. Mix spraces or Asper-
- B | 6 | 17 = 21 , June 2115 . [67] Ten R. Halffl. Marking Challenger PPGAs. Microscope & Series
- 85 Laure Hamm out, Baren A. Narfel, and Kanle Oblates, A Statle Chip Maltip morning /ASSE Computer, 38 p p79-45, 3997.
- 1891 Beiner W. Bartentein. Course crain recording table architecture less testes committe a ASP-DAC, once 384-471, 2013. Fell Beiner W. Harristonia. A decade of reconficing the companies: a cities

are noting on the 1st DATE, pages 44240, 2014.

Fill Alles Harrieta and Thomas R. Pamb. The section made size death for siona province. In Proceedings of The 19th Anternational Supressin Fil Alia Burriela and Thomas B. Parak. Online in a president masse

- Fil. Into B. Hanner and Into Winnerson. Gare: J.MIPS Propriate with Beconfigurable Copraces on In Proceedings of the IEEE Symposium of Preli-Programm also Caston. Comparing Machines, pages 18-06, April
- F41 Jan Helping, MIPS BISC Archive up. Volume I: Introduction to the 18A [0.4 st.]. Diction Number 117-15-11-11-11-117-11-1, Pri 5, [5] John Henney and David Patterinn. Computer Architecture: A do m-
- e Approach. Marga a Kaufta ann Publishen, Inc., 1886. [6] M. Peter Mattee. Power efficient processor and iterate and the celprocesses, in Proceedings of the 13th Antonian and Conference on Migh-Pajarn m m Canyo te Árchitechre, pages 258-462, Palenary 200
- [7] M. S. Heisk ike k, Dong Hurger, Stephen W. Kenkler, Premk k here Skie stomer, Neman P. Jospyi, and Keith I. Parker. The optimal legic death per pineline clare in 6 to 5 fed inverter delays. In Proceeding of The 18th International Symposium on Computer Architecture, pag N=24, June 2012.
- FS| M.T. Kang and C.E. Leisen in. Spirite army: [for VLS1], in Span-Motely Proceedings 1979 Fel Onias Jacobium, Steve Bounett, Nikhii Sharma, and James E. Smith

- hrs, pages 238-229, Petroary 3897. [0] B. M. Janesia and J. C. Brewn, A contribution of the a contrib
- stable attac companies. In Proceedings of the 5th Ames at International sainn an Caspete Árditadire, juga 8 148, 1882. [61] J. A. Kalle, M. N. Day, H. P. Hoft et, C. R. John, T. R. Masses and D. Shippy. In trade that to the Cell multiprocessor. IBM James
- of A so on th. on 4 D or elegen on 1, 43 (4/3) , 5 eptemb or 2005. 32 Cent J. Karasi, William J. Dally, Scott Blower, Peter B. Matting John D. Owen, and Branck Khadany. Efficient Conditional Open tion for Data-parallel Architectures. In Pres. of the 13rd Ann. Int'l Symp. on Mile words technology pages 201-271, 2001.
- [12] Upul J. Kapat, Peter Mattins, William J. Dally, John D. Owen, , as Briss Towler. Stream Scheduling. In Proceedings of the 2rd Workshop m Notes and Stromon Property rates 31 1-115. Dromb et 211 h. 841 Stockers W. Kerkler and William J. Dally, Processor contiling interesting compile time and restime the dalling for parallelism. In Proceeding of the 18th Annual Followish must Symposium on Computer Architecture,
- THE RESERVE ACMINING THE

- REL Characters Non-Danie Borner, and Stanton W. Kartler. An admitted and an illerial course to refer the delay described and the course is Proceedings of the 18th Symposium on Architectural Support for Pro cramming Lamp as as at Complete States, page 111-111. October
- [67] He-Seep Kim and James E. Smith. As increasing out and in dillecture for instruction level distributed processing. In ISCA, page 7140, 200, [68] Percarka Kongelin, Kart irgana r Ainga ma, and Konte Olskeins. No
- agam : Å 12-way maliiki maded S parcyricenii n. /ARR M i n. p. 23 [2]:23 Ξ [88] C. Koryrakis, J. Gala, D. Marin, S. Williams, I. Marmidis, S. Pope D. Jones, D. Patternen, and K. Yelick. Vector BRAM: A Media-enterted Very set Property or with E-market of DR ANL, in 12 to Well China Conference,
- According [6] Christon Kommitte and David Patternes. Oversming the limitation of contentional vector processors. In Proc. of the 18th Int'l Sprage on Computer Architecture, pages 199-419, June 2011.
- [84] Banny Kenshinsky, Christopher Butten, Mark Mampion, Store Geeling Briss Phagra, Jared Canner, and Krote Acasonic. The Vector-Thread size in Computer Architecture, pages 52-43, 2014.

- 82 Kerk Krewit. Sun't Nagara paun in the come. Managenessa
- [0] Kee is Krewell, Startey Argis Accelerate Beatly. Microprocessor Re. per (April 2015). Bell Radon & Kom or, Roll L. Parkas, Norman P. Joseph Partin paratic Ranparathus, and Dean M. Tullom. Single-by belonger may mad

A apor 1, 15 p : 11=13, 5 eptember 2014 .

- ings of 5 c 16 5 International Symposium on Microarditecture, pages 35] Sales Kemar, Norman P. Jonesi, and Dear M. Tellora. Combined
- con day malays manage, in M/CMO, pages 305 = 106, 2014. hal A. Kraimann et. at. Weine thin Architecture For Employ Stations ARES M(e s, 21 2):41 €7, March 21111.
- [67] L.E. Shar and E.S. Davidous. A Multiminiprocessor System imple seated Through Pipelining, Computer, Peb. 1874, pp. 42-61 88 Bate B. Lee, Subward carallelements may 1, AREE More, 18 91 S In
- [65] Walter Lee, Rajeet Barns, Matthew Frank, Dembhaktani Szikriskan Januthan Babb, West Surtur, and Samus Americingles. Space-time schede in cofficient of a clered complete monorary machine. In ASPLOS VVV : Proceedings of the eighth international conference on Architectural 237

New Yor, NY, USA, 1888. ACMP 1881. [80] S.A. Mallin, D.C. Liu, W.Y. Chin, R.E. Hall, and R.A. Bengman

- Effective compiler to pport for predicated execution using the by park lock is Proceedings of Section American and Symposium on Microsc [10.1] Ken Mai, Tim Parche, Nawan Jaranesa, Ban He, William J. Dalle, an Mark Microsite. Smart memories: a modular reconfigurable archite
- tom. In Proceedings of the 27th Amount International Symposium or Comparier Architecture, page 20 1407 1, June 2011. [107] William B. Mark and Donald Bouell. Benjalma we derive contents in The University of Texas of Annia , Annia, T.X. May 2005.
- [83] William R. Mark, R. Sterm Glassife, Kort Abeley, and Mark J. Kil-
- [B4] Baber McDerald, Bamadao Nagarajan, Karibiteyan Santamilagam Dang Burger, and Stephen W. Keckler. THEPS Increasing Ser Ardirectors [55] Marcal, Tedralcal Report TR-65-9, Decarding to

22.1

siam on Morror & Hochr & pages 127-225, December 2003. [63] Tem R. Halfalli, IDF Deliver Extreme Surprise, October 2003. 222

Control flow upon latter in multicolary more ent. In Proceedings of \hbar s 224

16, March/April 1888.

235

and Co. I States for Programming Great to Marianas is a Collin

- [83] B. Moors, A. Patigo, R. Smith, and W. Buchele. Concepts of th Sysm/17 t Ambitectum. In Proc. of the 11th Int'l Sums. on Commite
- [B6] Charles B. Moore. Managing the Transition from Complexity to Elemany Desira Conservors, IEEE Mary 11 (129-4), 2014.
- [97] Bamadan Naparajaa. Design and Anakais of John sign Scalable Arof Computer Sciences, December 2006.
- [38] Bamaian Namania, Xia Ches, Robert G. McDonald, Donn Burger. and Stephen W. Keckler. Critical path analysis of the trips architecture. In Proceedings of the IEEE International Symposium on Parlows once Anabas (ISPASS), race 27-47, March 2006.
- [305] Bamadan Nagam jan, Sandwy K. Kothwalin, Dong Berger, Kathryn S McKinley, Calcin Liu, and Stephen W. Keckler. Static Placement Oynamic from [SPDI] Schooling for EDGE Architectures. In 13th Je tem at mal Conference in Parallel Architecture in 4 Compilation Ted-
- [130] Bamadan Naman ina Karib kerna Suntarah man, Sies ken W. Keckler and Dong Burger. A Dongs. Space Evaluation of Grid Processor Archiv tectums. In Proceedings of the 14th Arms of Internal en al Symposis
 - 233
- [B1] J. M. Testler, J. S. Doton, Jr. J. S. Pieto, H. Le, and H. Sinkany POWER i colon m kranchitettare. IBM Jarrael of Score dom i Derdgmet, 26 [1:5-86, Junior 2013.
- is AFIP 5 Conference Proceedings , 1864 Full John Computer Conference, Stanton Banks for Work haven D.C. (1965)
- [32] Mary Trembler, J. Mickard O'Conner, Venin test Narayanan, and Line He. Vis specia new media processing . IEEE Marro, 18 $|\mathbf{i}| : \mathbf{0} - \mathbf{0}$, 1886.
- [83] Trimurus : As in fact mit are for received in in circle in a select parallel in
- [D4] Data M. Tellora, Street J. Hopen, Jud S. Huer, Black M. Leev, Jack L. Lo, and Behoven L. Stamm. Explicing choice: increasing freek and ione on an implementable charles on a multithreading processor. In IS Ci. 1961: Proceedings of the 10th sensel in terms is an allow receive
- [B3] Dan M. Tellen, Spon J. Erren, and Heart M. Leen. Simplicaon Makinkredier: Maximizias On-Chie Pamilekon. In 15 Ot 1931 vånga af å stård anns at internationatiopsyssism on Conyeter erditechre, page 20 1401, 2011.

- [113] Param Jr. S. Oberel and Gardadar S. Sakis, Para Selica in the first tend. n /8C4, rate 111-141, 1111.
- [122] John Oliver, Marichankar Man, Paul Saliana, Jel Rink Crandall, Erik Cremitowski, Leide Josev IV, Diana Prantin, Venkarek Akelia, and Prei etic T. Che se. Se achre color: A Ma his le Clock Domais. Perse Aware, Tile-Based Embedded Processer. in Fracesdings of the 18th Annual International Symposium on Architecture, pages 188–18 %, June
- [13] Tim Olon, Advanced Processing Techniques Using the Intrincity Fact AATH Promoter, in Embedded Process or Forum, May 2002.
- [18] Alex Pajuete, Antonio Gunnier, and Mates Valence. Specialize Dy namic Verterin ties. In Proceedings of the 19th Anternational Surger sion in Computer Architecture, pages 27 1-481, May 2012.
- [18] S. Paladada, N. P. Jacppi, and J. E. Smith. Complexity-effective expension pressure. In Proc. of Sciil Science (ed.), Smith.
- Computer Architecture, page 285-418, June 2017. [136] Alex Peter and Uni Weber. Many rechaptors extending to the intell
- [12] D. Pann, S. Acass, M. Heiliger, M. N. Day, H. P. Hefter, C. Jakon, J. Kakle, A. Kameyama, J. Kenty, Y. Macch schi, M. Biley, D. Shippy D. Stolisk, M. Stmiski, M. Wang, J. Warnesk, S. Weisel, D. Weisel,
 - 241
- [BS] Jim Tarley. Tencilca CPU Bend: to Designers's Will. Mcc spraces or 1871 Then Charges Berns Relation 289; and June & #232 Mrs. A same 11 [1:29-43, 2013.

Agent 11 (4), March 1999.

- [25] Seram Value com and Trika Mira. Improving on greater instruction dig atch and it me by explaining dynamic code requesces. In Proceedings of the Alfa Association of metallicities on the Computer Architecture. pages 140, June 2017.
- [BS] V.Kathall, M.Schlander, and B.R.Ran. Hylogd architecture openiica Laboratorie, Petrony 2001.
- [30] Eller Wainpeld , Michael Taylor, Devah haki an i Srik risk an, Vivek Surkar Walter Lee, Victor Lee, Jane Kim, Matthew Prank, Peter Finck, Baller Barra, Janutian Bath, Samar Amendingke, and Amer Agarest. Bar ing it All in Saftware: RAW Modeline. Computer, 11 p | 16-41, 1887.
- [30] S. Clin Waster and Judy J. Den comp. As terms the fer repet threat alcebra cofigure. In Supercomparing 1911: Proceedings of the 1991 $A\,CM/I\,EEE$ conjected in Superconjustry page (40). HEEE Conjected patter Suckey, 1986.

- T. Yamazaki, and K. Yamwa. The dwigs and implementation of a firstpromition cell processes. In IEEE International Solid-State Circuits
- [18] Mait Phare and Greg Humphreys. Design and implementation of a Physical Princet Reports System, Drug edition, April 2011.
- [139] Dunity Penemany, Garban Kacab, and Kanad Ghous, Reducing news mality is datapath measures. In Proceedings of the SUS International Symposium on Microsophileches, pages 58-68 3, 2002
- [20] Marko: Podes, Jo Morra, Jeremy Johnson, David Pades, Massels res Versaeska, Kana Chen, Balteri W. Johnson, 1 and Nick Birrole SPIRAL: Cuts Georgies for DSP Tranforms, \$1:222-675, 2015.
- [Bill Steen E. Rauck, Nathan L. Binker, and Steven K. Reinkantt. Accalable instruction queue design using dependence chains. In Proceedings of the 19th International Symposium on Computer Architecture pages 19-217 1111
- [22] Batric Ratta L, Inc Bratt, Kni e Acasoric, and Assar Aparent. Venstilling and Verrallench: A New Metric and a floud mark Suite for Flexitie Artitled tra . Test sical Report MIT-LC 5-T M-6461, Manuali met it
- [82] O. Welf and J. Blom. Star Com Lanacher Pirt. Architecture. Microgram morr Americ 2 (14): 2-41, October 2015.
- And Opening on Computer Architecture, pages 100-418, June 200 L.
- [384] W. Yamam ete and M. Nemimerky, Increasing our excellent effects and through multistruming. In Proceedings of Sec. 5 (normal engl Con pages 49-45, June 2015.
- [85] Zhiyi Yu, Michael Memanea, Ryan Apparina, Omar Saltari, Michael Lai, Jermy Web, Erk Werk, Timorh Mehenin, Mandeep Singh, and Berne M. Bern. An Arynchronen Array of Simple Procession for DSP Applications, in Proceedings of the IEEE International Sold State Cir. min Conference, \$55 CC 183, pages 428-428, 2.888.
- [93] Liu We, Cara Woose, and Told Applia. Cross-Maniac: A Part

om pa kep, 7 | 1:3-41, 1991.

Karibibeyan Santaralis yani was born in Chenna Lindia on Sib Petroary 1978, the con of Paritheil Sunkaralla part and Alyacamini Sunkara man the Indian Institute of Technology, Midden in 1888. He entered the gradual erann m is Camputer Schoon of the University of Texas at Aprila is April 2000. The received a Manner of Science degree in Assault 200

Vita

[22] S. Rajagopal, S. Rixner, and J. Comillaro. A programmable baseband

[04] N Ranganathan and M Franklin. As empirical only of discontained

[33] Niya Banganathan, Ramadan Nagarajan, Dong Borger, and Stophen W. Keekler. Combining by port-looks and odd prediction to increase front-

m d has dwidth and performance. Technical Report TR-62-6 % Depart

[26] B. Ramaksiskas Max and Joseph A. Phier. Instruction-level parallel

[27] Scott Rissor, William J. Dully, Upul J. Kapati, Bracel Khadany, Abelards

Laper-Lague to, Peter B. Mattine, and John D. Owens. A handwidth

m al international Symposium on Microschiteches, page 1-12, 1888.

[BS] Scott Riemer, William J. Dady, Upu I J. Kapaci, Bracek Khadany, Abelarda

Layer-Layer to, Peter B. Martins, and July D. Owens, Advantages to

efficient architecture for media processing. In Proceedings on the 1-bit

proceeding: Mintary, attention, and purspective. The Jaronal of Super-

ment of Computer Sciences, The University of Texas at Annin, Annin,

statement desire for coffware defined to disc. In Press of the IEEE Intil

ILP execution models. In 6th International Conference on Architectural

Support for Programming Languages and Operating Systems, pages 272—

Perman and address: M 12 William Creek Dr. April 200 Annin , Texas 797 41

This disconsists was typical with $D^{\alpha}(\mathbf{x}^{\beta})$ by the arthur. TOTAX is a deciminal proporation system developed by Leafle Languages or a special continual Bandid Kantidy TyX Programs

Aniem al mai Symposium im Microst del technic, pages 1-13, December mar, Stephen W. Keckler, and Dong Burger. Distributed Microarch -[23] Raddy Ungaban and Will Mean and Andrew McCabe. Systate Arrays.

- In this are of Physics Publishing, 1887. [33] Karthibe as Santaralianan, Vincent Alay Sinch, Stephen W. Kerbler, [20] E. Rateshen, O. Jacobson, Y. Samiles, and J. Smith. Trace on one rs. In Proceedings of the 18th Arms of distinct mad Supersism of
- m Compa to Design, pages 178-27, October 2001. [BB] Richard M. Ramell. The CRAN-I Computer System. Communications
- # 5 c.l (W. 11 0.8 (-1), June 1975.
- [B2] R.Sanaka, ALP: Energy Egic on Support for All Locals of Paralleless for Complete Medic Applications, PhD thesis, Chicoccity of Hinch at Urban a-Champuigs, Department of Computer Sciences, July 2005.
- [22] Karthikera Sankanilaran, Banadao Nasaraja, Bainia dia, Chanten Kim, Jurdynk Hab, Stephen W. Reckler, Dang Hanger, and Charles R. Moore. Exploiting LP, TLP and DLP with the Polymorphone TRIPS Architecture. In Proceedings of the 19th Annual Journal on al Sunsasion in Computer Architecture, pages 422-413, June 2013.
- [B4] Karibikeyan Sankaralingana, Ramadan Nagarajan, Makeri McDanaki, Rajagop sha Desikas, Saam bh Dreila, M.S. Garini as, Paul Grais, Dieya Galati, Beather Banco, Changleya Kim, Bhiming Liu, Nitya Rungan athan, Simba Sethamad haran, Sudia Sharif, Premk b here Shiraka-

- tectural Protocole in the TRUPS Prototy on Processes. In Proceedings of 41-43, Dromber 2016.
- and Dang C. Burger. Basted inter-ALU Networks for ILP Scalability and Performance. In Proceedings of the 2 Jul International Conference
- [23] A. Senare and R. Espain, Conflict-free account a classical vectors on a d cards. IEEE Transactions on Computers, 54 F | S 12-4 M, 2015
- min, Bill Yoles, Dans Barner, and Karlarya S. McKinley, Complian for EDGE and ledges. In Fact Internal maj &CN / IEEE Summann m Cute Commutes and Optimization (CCO), pages 251-281
- [23] Aaron Smith, Romoton Nagarajan, Karib Reyan Sankara Ingam, Robert McDanid, Deng Berger, Seplen W. Keetler, and Kaitera S. McKin-ley. Efficient Dataflew Produktion. In Proceedings of he 35h Annu al Antom at mat Ograposium on Microsochitecher, pages 41–43, December
- $[0.9]\ 0.3.5\ \mathrm{min}$. Acts become and applications of the BEP multiprocessor 244

- computer (priem. In SPIE Real Time Signal Proceeding 177, pages 24 le-[80] James E. Smith, Grey Passer, and Rathin A. Suprimer. Vector intens-tion set oupport for continuous operation. In Proc. of the 47th 2019. Same, on Computer Architecture page 26 1-461, June 2011.
 - [101] Garmeier S. Sehl, Scott E. Breach, and T. N. Wincksman, Malticular Computer Architecture, pages 436-435, Jane 2885.
 - [162] Eric Semante and Dear Campean, Increasing processor perfermance by implementing dampar pipelines. In Precedings of Tee 28th Internation of Symposium on Computer Architecture, pages 25-46, June 2002.
 - [M3] S. Srin burns, H. Hajenr, H. Akkary, A. Ghandi, and M. Opina. Continnal flow pipelines. In do term at small Conference on Architectural Suppor or Programming Languages and Operating Systems, yages 207–228, Oc
 - 1861 Standard Performance Employing Compression, SPEC CP 9210 L 2001.
 - [MS] Comm. Takah ash i, Scott Cottler, Sang M. Dhang, Brian Flacks, and Just Silberman. Pawes Cancion: Design of the Cell Processor's Synegical Processor Element. IEEE Micro, 25 p.: 38-35, 2005.

243

- [96] Deeps Talls, Lim John, and Dong Burger. But describe in multimedia processing with SIMD sink extensions and architectural enhancements
- [M7] Michael Beiffeel Taylor, Javon Kim, Javon Miler, David West rint, Fac Gladent, Ben Gresswitt, Henry Hoffman, Part Johnson, Walter Lee Jas-Wook Lee, Albert Ma, Archal Sund, Mark Senecki, Nichan Shaid Approxi. The flaw microprocess in λ computational fatric for sufware circular and renoral-someon enough no. ARRE More, 11 (121-15)
- Agarwal. Scalars penal activities the chip intercaused for ity is part in tioned and iteriores, in Proceedings of h of h Internal multigrapoise
- [33] Michael Beifferd Taylor, Walter Lee, Janua Miller, David West Half, Inc. Bratt, Bes Greswald, Heavy Hoffmans, Paul Johnson, Janua Kim James Pest a, Arché Saraf, Nathan Shaid man, Volber Stramp m, Matthew Prank, Saman P. Amarasingke, and Anast Agarwat. Evaluation of the Now Memoryon of the Execute Win-Date Architecture for U.C. and Stream. In Proceedings of the 13st Annual International Sympo-on Computer Architecture, pages 2-10, 2014.

246