# PipeRench Implementation of the Instruction Path Coprocessor

Yuan Chou, Pazhani Pillai, Herman Schmit, John Paul Shen
Department of Electrical and Computer Engineering
Carnegie Mellon University
Pittsburgh, PA 15213
{yuanchou,pillai,herman,shen}@ece.cmu.edu

## Abstract

*This paper demonstrates how an Instruction Path Coprocessor (I-COP) can be efficiently implemented using the PipeRench reconfigurable architecture. An I-COP is a programmable on-chip coprocessor that operates on the core processor's instructions to transform them into a new format that can be more efficiently executed. The I-COP can be used to implement many sophisticated hardware code modification techniques. We show how four specific techniques can be mapped to the PipeRench pipelined computation model. The experimental results show that a PipeRench I-COP used to perform trace construction and trace optimizations for a trace cache fill unit not only achieves good performance gains but can potentially be implemented in less than 10 mm$^2$ (assuming 0.18 micron technology) or approximately 3% of the die area of a current high-end microprocessor. We believe these results demonstrate the usefulness and feasibility of the I-COP concept.*

## 1 Introduction

### 1.1 Dynamic Code Modification

Spurred by relentless progress in VLSI design and fabrication, hardware design is evolving at a rapid pace and increasingly sophisticated microarchitectures are being implemented. On the other hand, software is changing much more slowly. One reason is the existence of a large installed base of legacy code that is too expensive to be replaced or recompiled. Another reason is that the deployment of new highly optimizing compilers usually lags behind the deployment of new microarchitectures. The end result is the increasing incompatibility between the compiler-produced object code and the most efficient implementations of fast execution cores that must execute these object code.

One recently proposed approach to solve this problem is to add hardware in the microarchitecture to dynamically modify the object code into an internal format that can be more efficiently processed by fast execution cores. We refer to this general approach as *hardware code modification*. For example, the Intel P6 [1] decoders translate the x86 instructions into an internal format called uops that are then executed by the execution core. Another example is the trace cache [2], which rearranges the ordering of instructions so that frequently executed sequences of instructions are stored in contiguous locations. The trace cache can reduce the complexity of instruction fetching and decoding. There are also proposals to optimize these traces [9][10] before loading them into the trace cache. Recently, there is a proposal to perform run-time program re-layout in hardware [25]. We believe that in the quest for ever higher performance, increasingly sophisticated hardware code modification techniques will be needed in the future.

An *Instruction Path Coprocessor* (I-COP), proposed in [3], is a programmable on-chip coprocessor that allows these hardware code modifications to be implemented in software much like microcode. An I-COP is analogous to a datapath coprocessor, except that it operates on the core processor's *instructions* themselves. The programmable nature of an I-COP affords several advantages. First, complex code modifications that are difficult to implement directly in hardwired logic may be more easily implemented in I-COP code. Second, it allows many code modification techniques to be implemented using the same engine, each of which can be selectively and adaptively invoked at run-time. Third, it allows specialization of microprocessors with the use of different I-COP code or even different I-COP implementations. Fourth, it makes it possible to modify and upgrade the machine simply by changing I-COP code without changing the hardware. We believe an I-COP can potentially be a valuable addition to the microarchitect's toolbox.

In evaluating the feasibility of the I-COP concept, [3]

showed that an I-COP prog rammed to i mplement tra ce construction and trace optimizations achieves good perfor- mance. T he l onger l atency (a s c ompared to hardwired logic) tha t the progra mmable I-COP t akes to pe rform the code m odifications ha d li ttle i mpact on pe rformance because the I-COP is located at the back-end of the core processor a nd be cause of t he fre quent re use of t he m odi- fied code. The prototype I-COP proposed consists of two VLIWs e ach wi th fou r ge neral func tion uni ts. Such a n I- COP i mplementation ca n r equire a s ignificant am ount o f chip area.

This paper proposes a novel and much more efficient I-COP implementation using a reconfigurable architecture called P ipeRench [4]. In s uch an im plementation, I-COP programs ar e actually *configuration bit s* that are down- loaded to t he reconfigurable fabric at run-time. Afte r c on- figuration, the fabric be comes a ha rdware design that implements th e de sired c omputation. Wha t distinguishes PipeRench from other reconfigurable fabrics is that it s up- ports v ery f ast re configuration a s we ll a s a vir tualization technique ca lled *pipeline reconfiguration*, which al lows a large logical design to be implemented on a small piece of hardware thro ugh ra pid c onfiguration of tha t ha rdware. This virtualization en ables s maller I -COP implementa- tions, and also allows complex I-COP programs to be writ- ten without the concern that they may not fit within the size of the reconfigurable fabric.

It was shown in [4] that the PipeRench reconfigurable fabric provides significant performance benefits for an application that exhibits one or more of the following fea- tures:

1. It operates on bit-widths that are different from a pro- cessor's basic word size.
2. Its data dependencies allow multiple function units to operate in parallel.
3. It is composed of a series of basic operations that can be combined into one specialized operation.
4. It can be pipelined.
5. Constant pro pagation can be p erformed, reducing the complexity of the operations.
6. The input values are reused many times within the com- putation.

The results in [3] suggest that the potential I-COP applications exhibit m any of th ese features. For e xample, the da ta bi t-widths in t he I-C OP applications a re odd a nd varied. There is al so ab undant p arallelism in these I-COP programs, thu s allowing multiple functional units t o ope r- ate in parallel. In addition, large portions of these programs are c omposed of basic operations t hat c an b e co mbined into specialized operations. In this paper, we show how hardware code modifications can be mapped to the PipeR- ench pipelined computation m odel and that the PipeRench

I-COP achieves good performance. Furthermore, we dem- onstrate tha t a P ipeRench I-COP c an be implemented a t very reasonable hardware cost, and in so doing, further val- idate the usefulness of the I-COP concept.

The rest of this paper is organized as follows. Section 2 familiarizes the r eader wi th the I -COP con cept an d the PipeRench reconfigurable architecture. Section 3 describes our PipeRench I-COP de sign and how I- COP applications are implemented in thi s design. S ection 4 presents the results of our exploration of the PipeRench design space as well a s d ie-area es timates o f s elected designs. S ection 5 concludes this paper.

## 2 Background

### 2.1 Instruction Path Coprocessors

An I-COP i s a p rogrammable coprocessor tha t oper- ates on the core processor's instructions to transform them into a new format that can be more efficiently processed by fast execution cores. These transformations can involve the ordering of instructions, the type of i nstructions (e.g. from the ori ginal i nstruction t o a s equence of s impler ins truc- tions) a nd even t he in struction s et (e.g. f rom the or iginal ISA to a new ISA tailored to the microarchitecture).

#### 2.1.1 Interface With Core Processor

The I-COP is located on the same chip as the core pro- cessor and runs c oncurrently wit h the core processor. In order no t t o negatively i mpact the c ore proc essor's cycle time, i t is situated at th e c ore p rocessor's b ack-end a nd interacts primarily with the core pr ocessor's comple- tion/retirement stage. The I-COP requires minimal explicit control by the core processor and rarely stalls the core pro- cessor. Figure 1 s hows t he i nterface be tween t he I -COP and the core processor.

An I-COP should be able to access non-architected entities of t he c ore proc essor, such as instruction and da ta caches, tr ace cache, branch and value predictor tables etc. Where s uch ac cesses a re allowed, careful c onsiderations are made to ensure that they do not affect the core proces- sor's critical timing paths.

In order for the I-COP to intelligently invoke the appropriate I-COP c ode ba sed on a pplication c haracteris- tics, t he c ore pro cessor has bui lt-in monitors to t rack it s currently executing application's behavior. The I-COP can either pol l t hese m onitors or the I-C OP can be in terrupt- driven. In the latter case, when the monitors exceed or dip below threshold levels, they interrupt the I-COP and cause it to vector to specific I-COP routines.
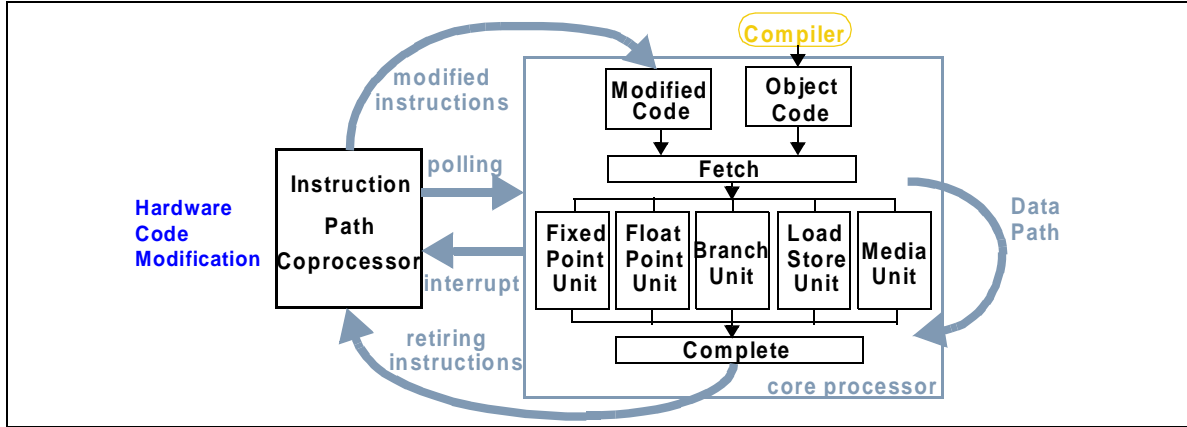
**Figure 1. Interface between I-COP and core processor.**

### 2.1.2 Initial Implementation

The ini tial I-COP im plementation [ 3] wa s based on conventional CPU de sign a nd c omprised of one or more VLIW engines (called slices) operating in parallel. For the I-COP a pplications s tudied, two VL IW s lices with four general functional units each represented a good c ost-performance trade-off. The VLIW organization was chosen to minimize hardware complexity, since I-COP programs are relatively small and can easily be statically scheduled. All the slices share a co mmon data memory. Since an I-COP replaces ha rdwired de signs wit h a pr ogrammable e ngine, slow-down can be e xpected. To e nsure a dequate performance, pa rallelisms i n I-C OP programs were e xploited; instruction-level parallelism was exploited within a VLIW slice and task-level parallelism was exploited across VLIW slices.

The instruction set for the I-COP VLIW slices consisted of 22 instructions. The core of the instruction set was a simple integer-based load/store architecture. In addition, ten specialized instructions were provided to facilitate writing efficient I-COP programs. The most important of these are powerful (a nd complex) pa ttern m atching instructions to enable re gular e xpression re cognition to be pe rformed quickly. Predication support as we ll as branch delay slots were also provided to eliminate the need for branch prediction. More details about this implementation can be found in [3].

The experimental results showed that this initial I-COP implementation achieved good performance for the I-COP applications studied. However, the drawback is that it requires a significant amount of hardware and can potentially consume sizable chip area.

### 2.2 PipeRench

PipeRench [4] is a reconfigurable fabric that supports the computational model shown in Figure 2. In this model, a computation on a data stream is expressed as a li nearly interconnected set of $S$ pipeline stages, where every stage is a function of the registered output of the previous stage and the registered output of the current stage. Many media and embedded c omputational ke rnels c an be m apped to this model with m any pi peline stages, whic h allows for high clock speeds and high throughputs. The small amount of feedback al lows f or ef ficient implementation. M any instruction tr ansformation techniques c an also be mapped to this model. In most instruction transformations, the particular t ransformations in itiated b y an y ins truction only affect subsequent instructions, which fits the limited feedback model.

Assuming that new inputs arrive every cycle, an implementation of this pipeline will require $S$ stages. In PipeRench, the technique of pipeline reconfiguration [5] is used t o s upport t he c ases whe n the inp ut s tream h as a n arrival ra te, or throughput $T$, which is less than one every cycle. In this c ase, $S$ physical stages cannot be kept busy. Alternatively the technique is also useful when the cost of $S$ stages is prohibitive. The technique is illustrated in Figure 3, where the number of stages in the application, $S$, is five and the number of physical pipeline stages $P$, is three. As the figure shows, the configuration of stages happens concurrently with the execution of other stages.

Using pipeline reconfiguration, the relationship between $S$, $P$ and $T$ is given by $T = max\left(\frac{P}{S}, 1\right)$. If $P \geq S$ and the input streams consist of a set of $N$ words, the entire
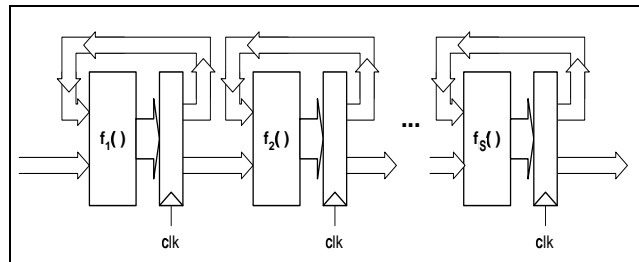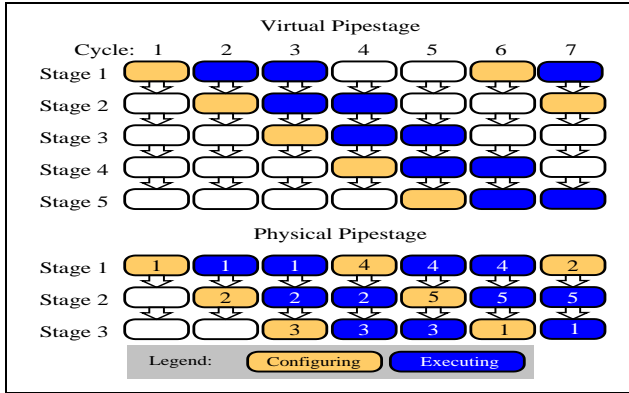


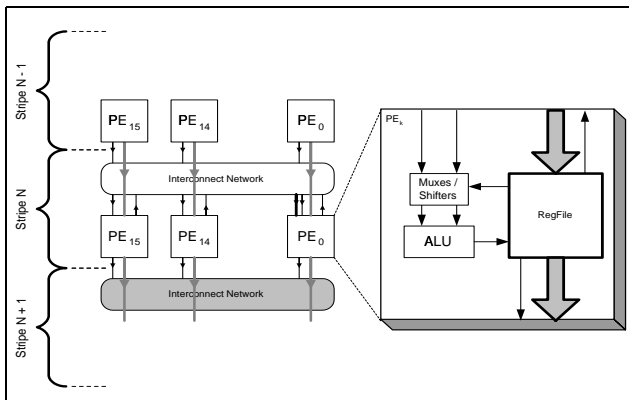**Figure 2. PipeRench computation model.**

**Figure 3. PipeRench pipeline reconfiguration.**

computation wi ll have a l atency of $N + S$ c ycles. If $P < S$ , then virtualization is necessary, and the computation will take $S \left\lceil \dfrac{N}{P} \right\rceil + P$ cycles to complete. In the rest of this pap er, we us e the te rm *virtual s tripes* to refer to the pipeline s tages re quired by t he a pplication a nd th e t erm *physical s tripes* to refer to the physical pipeline stages available. As any virtual stripe can be mapped to any physical stripe, all the physical stripes must have the same functionality and interconnect.

The current architecture of PipeRench is optimized by evaluating a set of media-centric applications and is illustrated in Figure 4. Each physical stripe consists of s ixteen ALUs (labelled PEs), which are each eight bits wide, connected w ith a byte-wise cr ossbar an d an el aborate s et o f shift r egisters. The A LUs are ca pable o f all p ossible b it-wise functions on two operands, as well as addition, subtraction and multiplexing. Each of the ALUs also contains an eight entry re gister f ile whi ch is pi pelined t o prov ide pipeline interconnect to do wnstream pipe line stages. State values (those feeding back in Figure 2) can only be stored in one s pecific register i n t he register fi le. An i nput a nd output bus moves operands on and off the execution fabric.

A physical de sign of t his architecture ha s be en com-



**Figure 4. PipeRench Architecture.**

pleted i n 0. 35 m icron and 0. 18 m icron process te chnologies. In 0. 18 micron, a single p hysical s tripe c onsumes 1.03 sq mm of silicon area, and operates at over 200MHz. Some small additional chip area is required for storage of configuration information a nd s tate t hat ne eds to b e h eld during vi rtualization. T his i s a ve ry c onservative de sign, with static CMOS circuits and fabricated in an ASIC process. We expect considerable headroom in improving both die area and clock speed.

PipeRench a pplications a re wr it ten in the *Dataflow Intermediate L anguage* (DIL), which is a single-assignment la nguage with C ope rators a nd a type s ystem t hat allows the bit-width of variables to be specified. The DIL compiler [6] c onverts the s ource in to a da taflow gra ph, decomposes t his gra ph i nto the na tive ope rators of the architecture and places and routes the operators on the PipeRench fabric. The output of t he compiler is a set of *configuration bi ts* (actually divided into a number of subsets, one subset per virtual stripe) that are used to configure the physical stripes at run-time.

## 2.3 PipeRench I-COP Advantages

In addition to being area-efficient, which we will demonstrate in Section 4.3, the PipeRench I-COP implementation a lso o ffers a num ber of oth er advantages. T he PipeRench architecture allows the designer to easily trade off the size of the reconfigurable fabric with other parts of the microarchitecture to optimize the overall design. Since the DIL code for the I-COP applications do not even need to be modified, changes to the number of physical stripes can be made very late in the design cycle. Moreover, when the same microarchitecture is implemented in the next process generation, the de signer ha s the o ption of inc reasing the number of physical stripes available to increase performance. Since physical stripes in the reconfigurable fabric are identical, this ca n be ac complished with m inimum redesign. The designer can also choose to upgrade the resident I-COP programs to further enhance performance. All in all, the PipeRench I-COP allows the designer to improve the performance of the c ore processor with minimal logic and circuit redesign.

The Pi peRench I-C OP a lso re tains t he ot her I -COP advantages li ke allowing c omplex ha rdware c ode modification techniques to be im plemented a s I-COP code a nd allowing m any ha rdware c ode modifications t o be imple-mented using the same engine, each of which can be selectively and adaptively invoked at run-time. In addition, the PipeRench I-COP makes it especially easy to specialize the core p rocessor b y va rying the s ize of the r econfigurable fabric to a chieve different performance goals and support different le vels of complexity in t he I -COP pr o-grams.
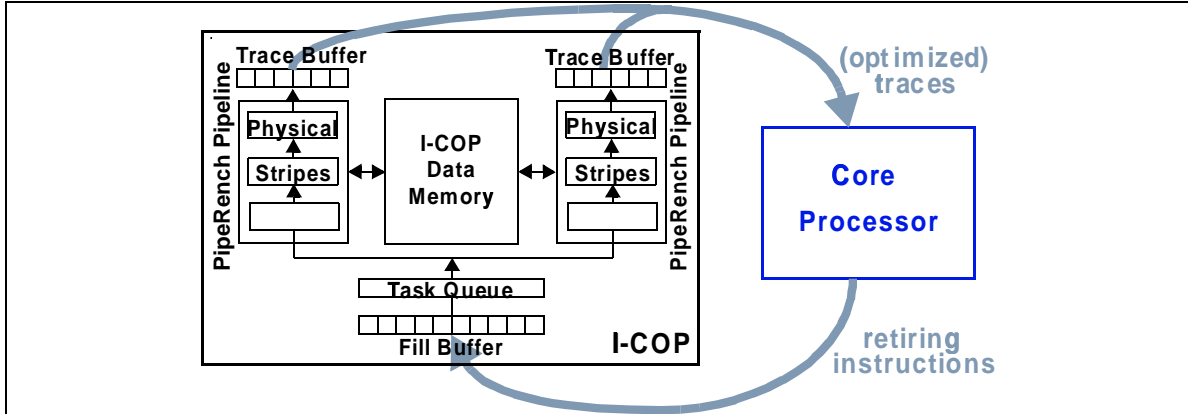
**Figure 5. PipeRench I-COP implementation.**

# 3    PipeRench I-COP Implementation

In this section, we describe the design of the Pipe-Rench I-COP. In order to assess its performance and die area re quirements, we s tudy how it  can im plement  four specific ha rdware c ode m odification t echniques: namely trace construction, register move trace optimization, stride data prefetch trace optimization, and linked data structure prefetch trace optimization. These are the same techniques implemented in the earlier study [3] and therefore allows us to compare the PipeRench and VLIW implementations in terms of performance and area efficiency. In Section 3.2, we describe how the se four c ode m odification te chniques are im plemented on  the  P ipeRench I-C OP a s DIL  pro -grams. W e  anticipate th at in   the fut ure,  many o ther advanced code modifications will be mapped to the PipeR-ench I-COP computation model.

## 3.1    PipeRench I-COP Design

The PipeRench I-COP implementation comprises of one or m ore P ipeRench p ipelines (e ach c onsisting o f one or  more phys ical s tripes) o perating i n pa rallel. A P ipe-Rench pipeline constructs and optimizes traces by treating the retiring instructions from the core processor as *streaming input data*. The outputs of each PipeRench pipeline are written  to it s lo cal  *trace buf fer*,  which ac ts a s t emporary storage to hol d a trace a s it i s be ing c onstructed. Whe n a trace is fully constructed, it is copied from the trace buffer to t he t race c ache. A P ipeRench I -COP i mplementation with two pipelines is shown in Figure 5. The fill buffer col-lects the retiring instructions from the core processor, and the task queue distributes them to the PipeRench pipelines. When  the fi ll buf fer i s ful l, i nstructions a re  dropped  at basic block boundaries. If a PipeRench pipeline has suffi-cient physical stripes to match the number of virtual stripes required  by  the  I-COP  applications,  it  accepts  one  fill buffer instruction per cycle as input and writes one instruc-tion to the trace buffer per cycle as output. Otherwise, the physical stripes are time multiplexed and the throughput of

trace processing will be less than one instruction per cycle. In Section 4, we evaluate the performance impact of vary-ing the number of Pipe-Rench pipelines and the number of physical stripes per pipeline.

## 3.2    Implementing Code Modifications Using Pip-eRench

To implement code modifications on the PipeRench I-COP, they are first mapped to the PipeRench computation model de scribed i n S ection 2.2. T hey a re t hen writ ten in the  DIL  la nguage a nd c ompiled  by  the  DIL   compiler to produce the configuration bits used to configure the physi-cal stripes of the PipeRench I-COP at run-time.

### 3.2.1    Trace Construction

The t race c ache [2][7 ][8] s tores  frequently  executed sequences of ins tructions in  physically c ontiguous storage locations, thus  a llowing hi gh ba ndwidth i nstruction fe tch without m ultiple c ache port s  nor ins truction a lignment logic. T his d ynamic re grouping  of i nstructions i s  per-formed by a hardware structure called the *fill unit* which is located at the back-end of the machine. A trace comprises not only of regrouped instructions but also the outcomes of the branches in the trace, the exit addresses of the trace (to facilitate pa rtial m atching [7]) a nd t he  type of t he  last instruction in the trace.

In  our  I-COP im plementation, logi c a ssociated with the fill buffer exa mines its first 16 e ntries and d etermines the  end of a  ne w  trace. It th en c opies  those in structions from the fill buffer to the I-COP memory and inserts a task into th e t ask que ue. Whe never a  P ipeRench pipeline i s free, it picks up a task from the front of the task queue and treating t he fi ll buf fer  instructions i n I-C OP m emory  as streaming input, processes one instruction in the trace at a time  and out puts t he pro cessed ins tructions t o the  tra ce buffer (see Figure 6). In the case of branch instructions, the PipeRench pipeline also outputs the branch outcome and
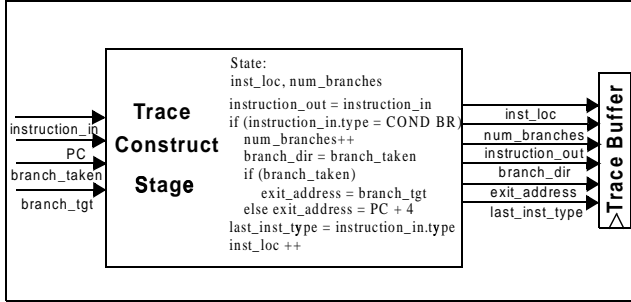
**Figure 6. Trace construction using PipeRench.**



**Figure 7. Register move optimization using PipeRench.**

exit address associated with that branch to the trace buffer. When the trace is fully constructed, the trace cache inside the core processor is read to check if there is an existing trace with the same starting PC. The new trace is written to the trace cache as long as it is not a subset of an existing trace. The PipeRench pipeline is then ready to pick up a new task. Based on the reconfigurable fabric's resource constraints (16 8-bit ALUs per stripe), the DIL compiler maps the trace construction logic to 11 virtual stripes. More details on the PipeRench implementation of trace construction and the other code modification techniques can be found in [26].

### 3.2.2   Register Move Trace Optimization

Beyond basic trace construction, the I-COP can perform optimizations on traces to achieve additional performance. Recently, there have been proposals for various trace optimizations [9][10]. The register move optimization [9] is one such example. In this optimization, instructions within a trace which move a value from one register to another register without modifying it are marked as *explicit move* instructions by the fill unit. Examples of such instructions are:

ADD Ra <- Rb + 0

SHIFT Ra <- Rb << 0

Instead of using execution resources to execute these instructions, their output registers are renamed to the same physical registers (or operand tags depending on the register renaming scheme used) as their input registers. Aside from saving execution resources, this also enables dependent instructions to execute earlier. The register renaming logic is modified to handle such explicit moves. A slight complication is that the input registers of dependent instructions within the same trace have to be substituted with the input register of the *explicit move* instruction.

In our PipeRench I-COP implementation shown in Figure 7, the register move optimization is performed after trace construction and before the trace is written from the trace buffer to the trace cache. Because this optimization is fairly expensive, it is not applied the first time a trace is
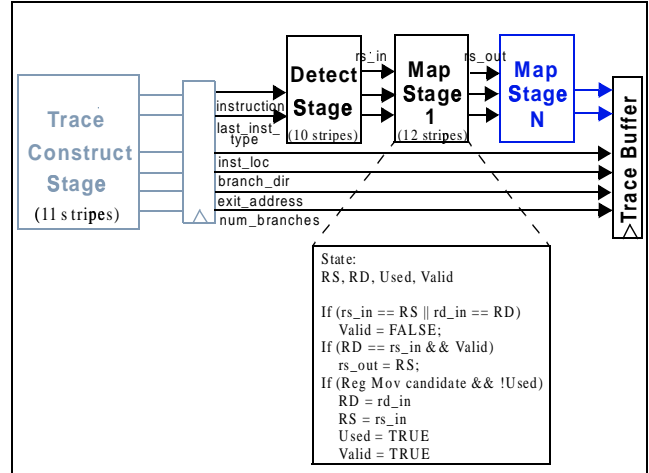
written into the trace cache. It is only applied to a trace that is found to be already in the trace cache and has been accessed *x* number of times. We found *x = 5* to be a good choice. Also, a trace is only optimized if it contains more than one conditional branch, since we assume the compiler already performs this optimization within a basic block. The input to the optimization is a stream of instructions from the result of trace construction and the output is a stream of optimized instructions that are written to the trace buffer. For every input instruction, two operations have to be performed. First, it must be determined if this instruction is a candidate for the optimization and if so, its type should be changed to that of an *explicit move*. Second, one or both of its source operand specifiers (i.e. register numbers) must be modified if that operand is dependent on an earlier register move candidate in the trace.

The first operation is essentially combinational logic and is performed by the Detect Stage shown in Figure 7. The DIL compiler produces a design of this stage that requires 10 virtual stripes.

The second operation is accomplished by keeping a set of mappings, labelled as the Mapping Stages in Figure 7. Each stage stores four values:

* a **valid** flag
* a 5-bit value **RD** which represents a register that is being mapped
* another 5-bit value **RS** which represents the register to which **RD** gets mapped
* a single bit, called the **used** flag, which is set if this stage has a valid mapping or ever had a valid mapping

At the beginning of each new trace, all of the stages are set to **invalid** and **unused**. When an instruction enters a stage, if the stage is **valid** and if a source register of the incoming instruction (rs_in) matches **RD,** then that source register will be renamed (rs_out) to **RS**. If the incoming

instruction is a register move candidate and the stage is **unused,** and if this instruction's mapping has not yet been stored, then the stage will be marked as **used** and **valid**. The source and destination of the instruction (rs_in and rd_in) will be stored in **RS** and **RD** respectively. A one bit flag will be sent to downstream stages indicating that the mapping for this instruction has already been stored.

If an instruction reaches a stage in which the destination of the instruction (rd_in) matches either **RD** or **RS**, the stage will be set to **invalid**. However, the stage will remain marked as **used,** since it previously had a valid mapping in it. This prevents future register move candidates from storing their mappings ahead (in stage order) of an already stored mapping and ensures that older mappings in the trace always appear earlier in stage order. This in turn ensures that when a new register move candidate stores its mapping, its source register will have already been correctly renamed. Each instruction only needs to pass through the pipeline just once, thus enabling a throughput of one instruction per cycle.

For a simple example of how this design works (it can also handle all the complex cases), consider the following example of a trace with just three instructions:

ADD r2 <- r1 + 0   *(1)*
ADD r4 <- r3 + r2   *(2)*
ADD r2 <- r10 + r11 *(3)*

Instruction *1* is eligible for the optimization and will create the mapping (**RD** = 2, **RS** = 1). Instruction *2* is not eligible for the optimization but one of its source operands matches the stored mapping (**RD** = 2, **RS** = 1) and so the instruction is transformed to ADD r4 <- r3 + r1. Instruction *3* is also not eligible for the optimization and since its destination matches the stored mapping (**RD** = 2, **RS** = 1), the mapping is invalidated.

The DIL compiler produces a design that requires 12 virtual stripes for each Map Stage. In our simulations, we found that having just one set of mapping (i.e. only one *explicit move* is allowed in a trace) achieves most of the performance gains of this trace optimization. This means this trace optimization takes a total of 10 + 12 = 22 virtual stripes in addition to the 11 virtual stripes for trace construction.

### 3.2.3  Stride Prefetch Trace Optimization

The stride prefetching scheme we implement in the I-COP is based closely on the hardware scheme proposed by Chen and Baer [11]. The basic idea is to record the effective addresses of loads as they are executed, compute the latest stride by comparing this address to the last effective address generated by the same static load, and update a 2-bit state machine. Depending on the resulting state of the

state machine, a prefetch request may be generated. All this information is recorded in a table called the Reference Prediction Table (RPT) stored in the I-COP data memory. Whenever a load is encountered during the construction of a trace, the 512-entry RPT is consulted to determine if it has a consistent stride. If so, a prefetch instruction is inserted in the trace before it is written to the trace cache. The prefetch instruction is only inserted if there is an empty slot in that particular trace cache line. This optimization is performed after trace construction and before the trace is written from the trace buffer to the trace cache. The input to the optimization is a stream of instructions from the result of trace construction and the output is a stream of the same instructions plus possibly one or more prefetch instructions. The DIL compiler produces a design that requires 14 virtual stripes.

### 3.2.4  LDS Prefetch Trace Optimization

Linked data structures (LDS) include linked lists, trees and graphs etc., where individual nodes are dynamically allocated from the heap and linked together through pointers to form the overall structure. The LDS prefetching we implement is based on that proposed by Roth et al. [12]. In this scheme, the goal is to correlate pairs of loads like the following, where the result of the first load is used as the base address for the second load:

LOAD r2 <- M[0(r1)]
LOAD r3 <- M[8(r2)]

After the correlation is established, whenever the first load is executed, a prefetch can be issued for the second load to hide the potential cache miss latency. Correlations are established by actual values rather than by symbolic means, with the help of two tables stored in the I-COP data memory: the 256-entry Potential Producer Window (PPW) and the 512-entry Correlation Table (CT). Whenever a load is encountered during the construction of a trace, it updates the PPW and CT. It also searches the CT and if it is found to be a producer, a prefetch instruction is inserted as part of the trace before the trace is written to the trace cache. The DIL compiler produces a design that requires 9 virtual stripes.

### 3.2.5  Comparison With VLIW-based I-COP

Table 1 compares the number of virtual stripes required by the PipeRench I-COP programs to the number of operations and cycles needed by the VLIW-based I-COP for the same programs. For the PipeRench I-COP, the number of cycles required to execute the program depends on the number of physical stripes available and is governed by the equations in Section 2.2. In Section 4.2, we study the performance impact of varying the number of physical stripes.

| I-COP Application | Virtual Stripes | VLIW ops | VLIW cycles |
|---|---|---|---|
| Trace construction | 11 | 50 | 18 |
| Register move trace optimization | 22 | 423 | 106 |
| Stride prefetch trace optimization | 14 | 130 | 33 |
| LDS prefetch trace optimization | 9 | 86 | 22 |

**Table 1: Comparison between PipeRench and VLIW-based I-COP implementations.**

# 4 Experimental Results

## 4.1 Simulation Methodology

Our performance simulator is built around Digital's ATOM tool [13] and uses the Alpha ISA [14]. Although it is trace-driven, it models the resource contention (but not cache effects) due to instructions on the mispredicted path.

The organization of the core processor is as follows. The trace cache contains 128KB of instructions (2048 lines of 16 instructions) and is 4-way set associative. Partial matching is implemented. The branch predictor is as described in [7]. It is an adaptation of the *gshare* predictor, and makes 3 predictions per cycle. We assume a perfect return address stack which is used to predict subroutine returns. The L1 instruction cache is 16KB and direct-mapped, with a 14 cycle miss latency. Because of the low instruction cache miss rates, an L2 instruction cache is not modeled.

| Functional Units | Units | Latency |
|---|---|---|
| Simple Integer | 8 | 1 |
| Complex Integer | 4 | 4 |
| Load/Store | 4 | 2/1 |
| Branch | 4 | 1 |
| Floating-Point Add/Multiply | 4 | 3 |
| Floating-Point Divide | 4 | 11(sp), 15(dp) |

**Table 2: Core processor execution resources.**

The front-end pipeline of the core processor, from fetch to dispatch, is four stages deep. Instructions are dispatched to a 512 entry centralized instruction window and are allowed to issue out-of-order. Perfect memory disambiguation is assumed. The functional unit mix and their execution latencies are shown in Table 2; all functional units are fully pipelined. The L1 data cache is 16KB and direct-mapped and the miss latency to the L2 data cache (assumed off chip) is 14 cycles. The L2 data cache is 256KB and 2-way set associative, with a miss latency to main memory of 75 cycles. In our data prefetching experiments, prefetched data are brought into a 64 entry fully-associative prefetch buffer. The PipeRench I-COP model is integrated with the core processor's simulator and is simu-

lated in detail at the machine cycle level.

Seven SPECint95 benchmarks [15] and three pointer-intensive Olden benchmarks [16] are used. Their input sets and dynamic instruction counts are shown in Table 3. The benchmarks are compiled using the default optimization flags of the SPEC distribution and are run to completion.

| Benchmark | Input Set | Inst Count |
|---|---|---|
| *compress* | 10000 e 2231 | 54M |
| *ijpeg* | tinyrose.ppm | 89M |
| *m88ksim* | dhry2tiny.lit | 99M |
| *go* | 5 9 | 78M |
| *gcc* | -O genoutput.i | 106M |
| *li* | queens 6 | 56M |
| *perl* | trainscrabbl | 47M |
| *health* | 5 levels, 500 iters | 176M |
| *perimeter* | 4K x 4K image | 43M |
| *treeadd* | 1024K nodes | 98M |

**Table 3: Benchmark characteristics**

## 4.2 Performance Data

In this section, we show the performance of the core processor under different PipeRench I-COP organizations. In particular, we vary the number of physical stripes per PipeRench pipeline as well as the number of PipeRench pipelines. Since the reconfigurable fabric may have to be clocked at a slower clock speed than the core processor, we show two sets of results. The first assumes the reconfigurable fabric is clocked at the same speed as the core processor while the second assumes it is clocked at half the speed.

While evaluating the different design points of the design space, it is helpful to bear in mind that each physical stripe in a 0.18 micron process occupies 1.03 sq mm of silicon area (approximately 1/300th the area of a 300 sq mm die used in current high-end microprocessors).

### 4.2.1 Trace Construction

Figure 8 shows the performance of the core processor with its I-COP implemented in different PipeRench organizations for trace construction. In particular, we vary the number of PipeRench pipelines as well as the number of physical stripes per pipeline. The upper graph assumes that the PipeRench I-COP runs at the same speed as the core processor while the lower graph assumes that it runs at half the speed. In both graphs, the *y* axis shows the harmonic mean of the IPCs of the seven SPECint95 benchmarks and the *x* axis represents the total number of physical stripes. The sets of data points on each graph represent varying the number of PipeRench pipelines. The number of physical stripes per pipeline can be derived by dividing the total number of physical stripes by the number of pipelines. For
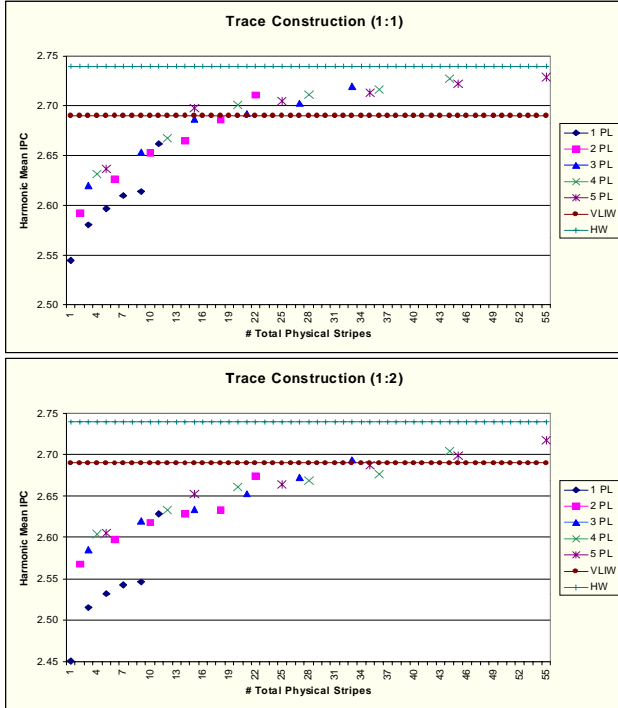
**Figure 8. Trace construction performance.**

example, the data point [2 pipelines (PL), 14 total physical stripes] implies there are seven physical stripes per PipeRench pipe line. For c omparison, the performance of the VLIW-based I-C OP (l abeled VL IW) a s we ll a s a ha rd-wired trace cache fill unit (labeled HW) are also shown.

The throughput at which the I-COP constructs traces is directly proportional to the total number of physical stripes available, w hile th e la tency o f tr ace construction i s inversely proportional to it. When the throughput of trace construction i s re duced, m ore i nstructions are dr opped from the fill buffer since the I-COP is not able to keep up with the rate at which in structions are retired by the core processor. However, because of the frequent reuse of previously constructed traces, these dropped instructions do not adversely a ffect o verall pe rformance. Mo reover, because the I-COP is located at the back end of the core processor, longer latencies in trace construction also do not seriously affect performance. Therefore, there is diminishing returns in performance as the total number of physical stripes is increased.

Given a fi xed t otal num ber of phys ical s tripes (a nd throughput), performance varies slightly depending on the exact PipeRench or ganization. This is due to several fac-tors. First, th e l atency of tr ace c onstruction ha s a c eiling function (see Section 2.2) that produces discontinuities. In particular, 11 physical stripes per pipeline results in partic-ularly good performance because the average trace length is approximately 11 instructions. Second, the number of

PipeRench pip elines a ffects trace s election because instructions are droppe d from the fi ll buf fer a t a different timing. When th ere ar e m ore p ipelines, th ere w ill b e a longer series o f c ontiguous t races followed by a larger number of dro pped ins tructions. Whe n t here a re f ewer pipelines, there will be a shorter series of contiguous traces followed by a smaller number of dropped instructions. The former situation is more desirable than the latter, so in gen-eral, for a given number of total physical stripes, it is better to have more pipelines and fewer physical stripes per pipe-line.

For a particular performance level (i.e. fixed value on *y* axis), the most desirable PipeRench I-COP organization is the one with the le ast tot al num ber of phys ical s tripes. For example, if we want to match the performance of the VLIW I -COP i mplementation, the design p oint [ 3 pi pe-lines, 15 total phy sical stripes] is t he best or ganization when the clock speed of the PipeRench I-COP matches the clock speed of the core processor. When the clock speed is half t hat of th e core p rocessor, t he de sign point [ 3 pi pe-lines, 33 total physical stripes] is the best organization.

### 4.2.2 Register Move Trace Optimization

Figure 9 s hows t he pe rformance of different PipeR-ench organizations when the register move optimization is applied in addition to basic trace construction.

The graphs in Figure 9 are organized in a similar fash-ion to those in Figure 8. Because the PipeRench I-COP is efficient i n i mplementing t his optimization ( 22 v irtual stripes in addition to 11 virtual stripes for basic trace con-struction; in c ontrast, the VL IW re quires 423 instructions in addition to 50 instructions for basic trace construction), fewer total physical stripes are needed to match the VLIW implementation. When the PipeRench I-COP runs at the same speed as the core processor, the [4 pipelines, 12 total physical stripes] organization m atches the performance of the VLIW I-COP. W hen it r uns at ha lf t he s peed, the [5 pipelines, 2 5 t otal physical stripes] organization a ccom-plishes the same goal.

From the results, we also observe that when compared to basic tra ce construction, give n the same PipeRench I-COP or ganization, a pplying thi s opt imization improves performance. For e xample, the harmonic mean IPC of the I-COP organization [3 pipelines, 15 total physical stripes] (assuming I-COP and core processor run at the same clock speed) increases from 2.69 to 2.72. Although these perfor-mance improvements are modest, no additional I-COP hardware was required; only the I-COP code, i.e. the Pip-eRench *configuration bits,* are changed.
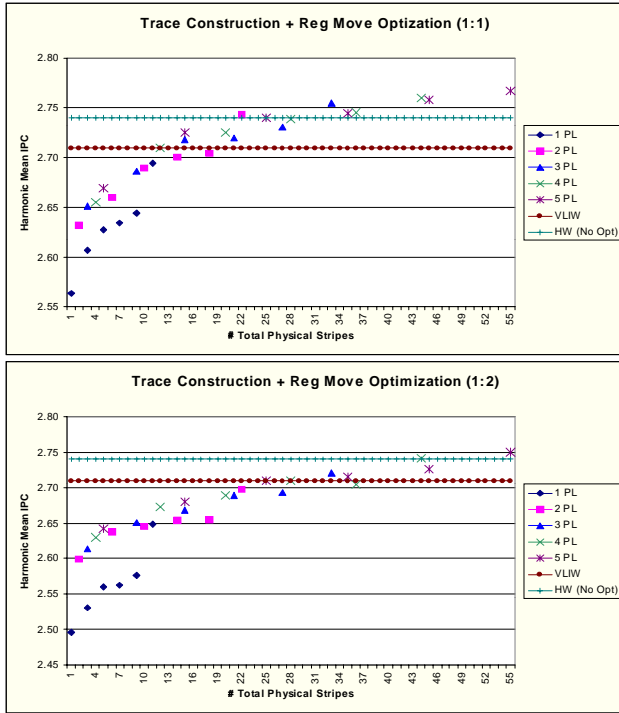
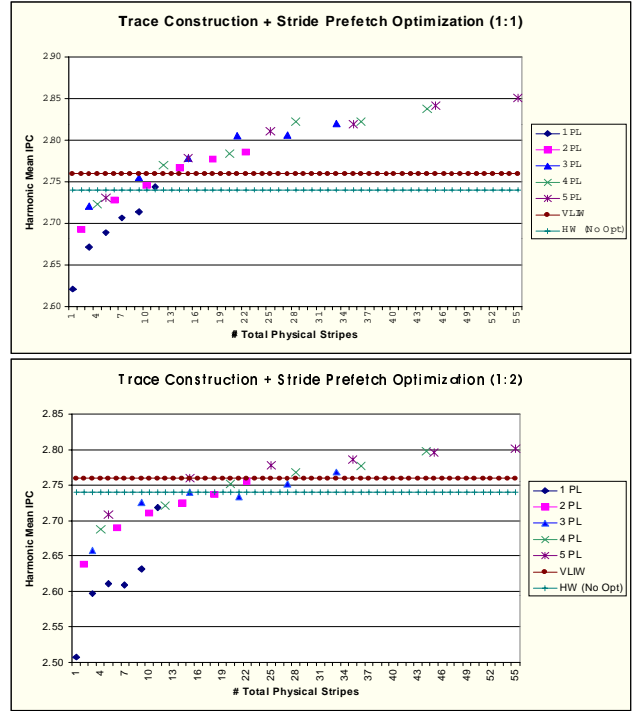**Figure 9. Register move optimization performance.**



**Figure 10. Stride prefetch performance.**

### 4.2.3  Stride Prefetch Trace Optimization

Figure 10 shows the performance of different PipeRench organizations when the stride data prefetch optimization is applied in addition to basic trace construction. This optimization is applied to all traces. The results are clearly superior to those in Figure 8, demonstrating the advantage of an I-COP in being able to improve core processor performance by modifying I-COP code and without changing the I-COP hardware. The PipeRench I-COP is efficient in implementing this optimization, requiring only a [3 pipelines, 9 total physical stripes] organization to match the VLIW I-COP when it is running at the same clock speed as the core processor. When it is running at half the speed, a [5 pipelines, 15 total physical stripes] organization is required. Note also that these I-COP organizations also handily exceed the performance of the hardwired trace cache fill unit performing trace construction with no trace optimization (labeled HW (No Opt) in Figure 10).

### 4.2.4  LDS Prefetch Trace Optimization

Figure 11 shows the performance of a PipeRench I-COP running at the same clock speed as the core processor when the LDS data prefetch optimization is applied. Because the IPC performance of the *health* benchmark is an order of magnitude lower than those of the other two Olden benchmarks, we avoid using the harmonic mean of their IPCs. Instead, the performance of each benchmark is shown separately. The I-COP organization shown is the

same one that matches the performance of the VLIW I-COP for the stride data prefetch optimization, i.e. 3 pipelines, 9 total physical stripes. We observe that a small PipeRench I-COP is able to match the performance of the VLIW I-COP. The performance of this PipeRench I-COP also exceeds that of the hardwired trace cache fill unit with no trace optimization by a considerable margin.
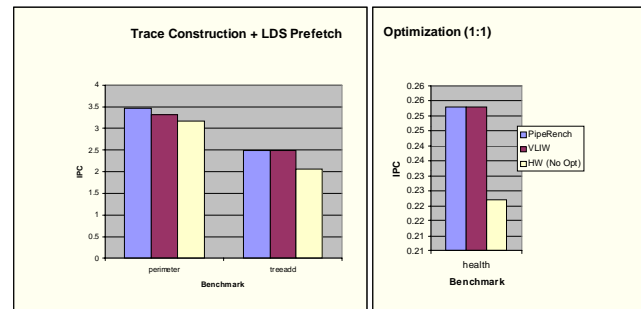


**Figure 11. LDS prefetch performance (1:1 clock speed).**

Figure 12 is similar to Figure 11 except that the results shown are for an I-COP that runs at half the clock speed of the core processor. The I-COP organization shown is the same one that matches the performance of the VLIW I-COP for the stride data prefetch optimization, i.e. 5 pipelines, 15 total physical stripes. Again, we observe that a small PipeRench I-COP is able to match the performance of the VLIW I-COP.
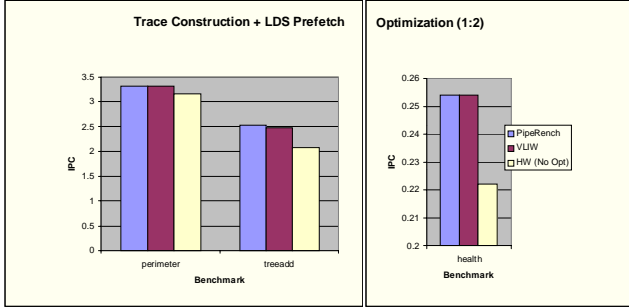
**Figure 12. LDS prefetch performance (1:2 clock speed).**

### 4.2.5 Estimated Area of PipeRench I-COP

To match or exceed the VLIW I-COP performance for trace construction and all three trace optimizations, and assuming that the PipeRench I-COP is only able to run at half the speed of the core processor, the estimated die area of the PipeRench I-COP (33 total physical stripes) fabricated in a 0.18 micron process is 33 x 1.03 = 34 sq mm. To put this in perspective, Table 4 shows the estimated die

| Component | % of die | Process (u) | Area (mm$^2$) | Area scaled for 0.18 u |
|---|---|---|---|---|
| IBM G6 [17] | FPU | 7.1 | 0.22 | 15.3 | 10.2 |
| Transmeta 3120 [18] | FPU | 12.3 | 0.22 | 9.5 | 6.3 |
| UltraSparc-2i [19] | FPU | 12.0 | 0.29 | 18.0 | 6.9 |
| AMD K6 [20] | FPU | 14.3 | 0.35 | 23.1 | 6.1 |
| NEC MP98 [21] | 64KB cache | 8.0 | 0.15 | 9.1 | 13.1 |
| NEC Cache SRAM [22] | 512K cache | 100.0 | 0.25 | 132.0 | 68.4 |

**Table 4: Die areas of microarchitecture structures.**

areas of other microarchitecture structures in the same process. The PipeRench I-COP is roughly equivalent in area to 256KB of fast SRAM, or about 11% of the die area of a current high-end microprocessor. If the I-COP is able to run at the same speed as the core processor, the die area required drops to approximately 15 x 1.03 = 15 sq mm, or roughly equivalent to 128KB of fast SRAM, or about 5% of the die area.

As noted in Section 2.2, these area estimates are likely to be conservative due to the conservative circuit design and fabrication process assumed. Also, the DIL compiler achieves relatively low utilization of PipeRench resources because of its fast and greedy approach to placement and routing, as illustrated in Table 5. The percentage resource utilization numbers are obtained by dividing the total number of native PipeRench operations in the application by

the total number of ALUs available in a design with the resultant number of physical stripes. Significant performance improvements can easily be obtained by optimizing the PipeRench architecture, circuit design and compiler for I-COP applications.

| I-COP Application | Resources Utilized Per Stripe |
|---|---|
| Trace construction | 54% |
| Register move trace optimization | 55% |
| Stride prefetch trace optimization | 70% |
| LDS prefetch trace optimization | 56% |

**Table 5: Utilization of PipeRench fabric resources.**

If one is willing to trade off a little performance (2.6% lower for trace construction, 1.8% for the register move trace optimization, 1.1% lower for stride prefetch trace optimization), and assuming the PipeRench I-COP is only able to run at half the speed of the core processor, one can implement the [3 pipelines, 9 total physical stripes] I-COP in 9 x 1.03 = 9.27 sq mm, which is roughly equivalent in area to 64KB of fast SRAM, or 3% of the die area of a current high-end microprocessor. In future fabrication processes (0.13 micron and beyond), the I-COP will occupy an even smaller fraction of the available die area.

## 5   Conclusions and Future Work

In this paper, we have described an efficient means of implementing an I-COP by using the PipeRench reconfigurable architecture. We also show how hardware code modifications can be mapped to the PipeRench pipelined computation model. In our experimental evaluation, we found that a PipeRench I-COP used to perform trace construction and trace optimizations for a trace cache fill unit not only achieves good performance but can be implemented in less than 11% of the area of a current high-end microprocessor. If one is willing to trade off only a little performance, this figure can be reduced to 3% or lower. We believe that this demonstrates that an I-COP can be implemented in a reasonable amount of chip area.

In addition to being area-efficient, the PipeRench I-COP implementation also allows the designer to easily trade off the size of the reconfigurable fabric with other parts of the microarchitecture to maximize overall performance. As the PipeRench configuration bits do not need to be modified, this trade-off can be changed very late in the design cycle. The PipeRench I-COP implementation is also highly scalable. As I-COP programs become more complex or more I-COP programs need to be run concurrently, the number of physical stripes in the reconfigurable fabric can be increased with minimal design effort.

In conclusion, we believe that we have demonstrated the I-COP concept to be useful and feasible. With the need for increasingly sophisticated hardware code modification techniques, we believe that an I-COP is a potentially powerful tool in the microarchitect's arsenal. We also believe that hardware code modification techniques enabled by the I-COP can be synergistically combined with software runtime code optimization techniques to further improve the performance of future high performance microprocessors.

Our current research focuses on studying other I-COP applications like using an I-COP to perform run-time trace scheduling [23] and completion-time branch prediction in the context of a trace cache [24]. We also plan to study the interface between the I-COP and the core processor in greater detail, and in particular how the core processor can selectively and adaptively invoke the appropriate I-COP programs based on application behavior. Finally, we hope that the demonstrated feasibility of the I-COP concept will serve to stimulate further research into advanced hardware code modification techniques.

### Acknowledgment

### References

[1] Linley Gwennap, "Intel's P6 Uses Decoupled Superscalar Design," in Microprocessor Report, Vol. 9, Issue 2, February 1995.

[2] E. Rotenberg, S. Bennett and J. Smith, "Trace Cache: a Low Latency Approach to High Bandwidth Instruction Fetching," in Proc. of 29th International Symposium on Microarchitecture, 1996.

[3] Y. Chou and J. Shen, "Instruction Path Coprocessors," in Proc. of 27th International Symposium on Computer Architecture, June 2000.

[4] S. Goldstein et al., "PipeRench: A Coprocessor for Streaming Multimedia Acceleration," in Proc. of 26th International Symposium on Computer Architecture, May 1999.

[5] H. Schmit, "Incremental Reconfiguration for Pipelined Applications," in Proc. of Workshop on FPGAs for Custom Computing Machines, April 1997.

[6] M. Budiu and S. Goldstein, "Fast Compilation for Pipelined Reconfigurable Fabrics," in Proc. of 7th International Symposium on Field Programmable Gate Arrays, February 1999.

[7] S. Patel, D. Friendly and Y. Patt, "Critical Issues Regarding the Trace Cache Fetch Mechanism," Technical Report CSE-TR-335-97, University of Michigan, May 1997.

[8] B. Black, B. Rychlik and J. Shen, "The Block-based Trace Cache," in Proc. of 26th International Symposium on Computer Architecture, May 1999.

[9] D. Friendly, S. Patel and Y. Patt, "Putting the Fill Unit to Work: Dynamic Optimizations for Trace Cache Microprocessors", in Proc. of 31st International Symposium on Microarchitecture, 1998.

[10] Q. Jacobson and J. Smith, "Instruction Pre-Processing in Trace Processors", in Proc. of 5th International Symposium on High Performance Computer Architecture, 1999.

[11] T. Chen and J. Baer, "Effective Hardware-Based Data Prefetching for High-Performance Processors", IEEE Transactions on Computers, Vol. 44, No. 5, 1995.

[12] A. Roth and G. Sohi, "Effective Jump-Pointer Prefetching for Linked Data Structures", in Proc. of 26th International Symposium on Computer Architecture, 1999.

[13] A. Srivastava and A. Eustace, "ATOM: A System for Building Customized Program Analysis Tools," in Proc. of SIGPLAN Conference on Programming Language Design and Implementation, 1994.

[14] Alpha Architecture Handbook, Digital Equipment Corporation, 1992.

[15] http://www.spec.org

[16] A. Rogers, M. Carlisle, J. Reppy and L. Hendren, "Supporting Dynamic Data Structures on Distributed Memory Machines", ACM Transactions on Programming Languages and Systems, 17(2), March 1995.

[17] K. Diefendorff, "Processors Penetrate Gigahertz Territory," Microprocessor Report, Vol. 14, Archive 2, February 2000.

[18] T. Halfhill, "Transmeta Breaks x86 Low-Power Barrier," Microprocessor Report, Vol. 14, Archive 2, February 2000.

[19] "Low-Cost UltraSPARC-2i Appears," Microprocessor Report, Vol. 12, No. 1, January 26, 1998.

[20] D. Draper et al., "Circuit Techniques in a 266 MHz MMX-enabled Processor," IEEE Journal of Solid State Circuits, Vol. 32, No. 11, November 1997.

[21] P. Glaskowsky, "NEC Decants Merlot," Microprocessor Report, Vol. 14, Archive 3, March 2000.

[22] H. Nambu et al., "1.8-ns Access, 550-MHz, 4.5-Mb CMOS SRAM," Vol. 33, No. 11, IEEE Journal of Solid State Circuits, Vol. 33, No. 11, November 1998.

[23] R. Nair and M. Hopkins, "Exploiting Instruction Level Parallelism in Processors by Caching Scheduled Groups," in Proc. of 24th International Symposium on Computer Architecture, June 1997.

[24] R. Rakvic, B. Black and J. Shen, "Completion Time Multiple Branch Prediction for Enhancing Trace Cache Performance," in Proc. of 27th International Symposium on Computer Architecture, June 2000.

[25] M. Merton et. al, "A Hardware Mechanism for Dynamic Extraction and Relayout of Program Hot Spots," in Proc. of 27th International Symposium on Computer Architecture, June 2000.

[26] P. Pillai, "The Instruction Path Coprocessor Implemented on the PipeRench Fabric," CMuART Tech. Report, Carnegie Mellon Univ., 2000.