

DEBUGGING AND ANALYSIS OF LARGE-SCALE PARALLEL SOFTWARE

by
Saeed Taheri

A dissertation submitted to the faculty of
The University of Utah
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy

School of Computing
The University of Utah
July 2021

Copyright © Saeed Taheri 2021
All Rights Reserved

The University of Utah Graduate School

STATEMENT OF DISSERTATION APPROVAL

The dissertation of Saeed Taheri
has been approved by the following supervisory committee members:

<u>Ganesh Gopalakrishnan</u> ,	Chair(s)	<u>?? July 2021</u> Date Approved
<u>Zvonimir Racamaric</u> ,	Member	<u>?? July 2021</u> Date Approved
<u>Hari Sundar</u> ,	Member	<u>?? July 2021</u> Date Approved
<u>Alexander Lex</u> ,	Member	<u>?? July 2021</u> Date Approved
<u>Martin Burtscher</u> ,	Member	<u>?? July 2021</u> Date Approved

by Mary Hall , Chair/Dean of
the Department/College/School of Computer Science
and by David B. Keida , Dean of The Graduate School.

ABSTRACT

I will add an abstract

For my family

CONTENTS

ABSTRACT	iii
LIST OF FIGURES	vii
LIST OF TABLES	viii
NOTATION AND SYMBOLS	ix
CHAPTERS	
1. THE FIRST	1
1.1 The first section	1
1.1.1 The first subsection	2
1.1.2 The second subsection	2
1.1.3 The third subsection	2
1.1.3.1 The first subsubsection	2
1.1.3.2 The second subsubsection	2
1.1.3.2.1 The first numbered paragraph	2
1.1.3.2.2 The second numbered paragraph	3
1.2 The second section	3
1.3 The third section	5
1.4 Free software packages	6
1.5 Resizing figures	9
1.6 Summary and conclusions	13
2. THE SECOND	15
2.1 Introduction	15
2.2 Background and Related Work	18
2.2.1 Binary Instrumentation	18
2.2.2 Efficient Tracing	19
2.3 Design of PARLOT	20
2.3.1 Tracing Operation	21
2.3.2 Incremental Compression	22
2.3.3 Compression Algorithm	23
2.3.4 PIN and Call-Stack Correction	24
2.4 Evaluation Methodology	25
2.4.1 Benchmarks and System	25
2.4.2 Metrics	25
2.4.3 Tracing Tools	26
2.5 Results	28
2.5.1 Tracing Overhead	28
2.5.2 Required Bandwidth	31

2.5.3	Compression Ratio	31
2.5.4	Overheads	35
2.5.5	Compression Impact	37
2.6	Discussion and Conclusion	37
APPENDICES		
A.	THE FIRST	39
B.	THE SECOND	40
C.	THE THIRD	41
REFERENCES		44

LIST OF FIGURES

1.1	The first figure.	2
1.2	The second figure.	3
1.3	The third figure.	8
1.4	The fourth figure (at 50% scale).	10
1.5	The fifth figure (at 75% scale).	10
1.6	The sixth figure (at native size).	10
1.7	The seventh figure (at 125% scale).	10
1.8	The eighth figure (at 175% scale).	10
1.9	The ninth figure (at 50% scale)	12
1.10	The tenth figure (at 75% scale)	12
1.11	Using \LaTeX picture mode	13
2.1	Overview of PARLOT	21
2.2	Average tracing overhead on the NPB applications - Input B	28
2.3	Average tracing overhead on the NPB applications - Input C	28
2.4	Average required bandwidth per core (kB/s) on the NPB applications - Input B	28
2.5	Average required bandwidth per core (kB/s) on the NPB applications - Input C	30
2.6	Average compression ratio of PARLOT on the NPB applications - Input B	33
2.7	Average compression ratio of PARLOT on the NPB applications - Input C	33
2.8	Tracing overhead breakdown - Input B	33
2.9	Tracing overhead breakdown - Input C	34
2.10	Variability of PARLOT(M) overhead on 16 nodes - Input B	34
2.11	PARLOT-NC tracing overhead breakdown - Input B	35
2.12	PARLOT-NC tracing overhead breakdown - Input C	35

LIST OF TABLES

1.1	Lowercase Greek letters.	4
1.2	Uppercase Greek letters.	6
2.1	Overhead added by each tool.	26
2.2	Required bandwidth per core (kB/s)	29
2.3	Compression ratio	32
2.4	Tracing overhead of versions of PARLOT(M)- Input B	36
2.5	Tracing overhead of versions of PARLOT(A)- Input B	36

NOTATION AND SYMBOLS

α	fine-structure (dimensionless) constant, approximately $1/137$
α	radiation of doubly-ionized helium ions, He^{++}
β	radiation of electrons
γ	radiation of very high frequency, beyond that of X rays
γ	Euler's constant, approximately $0.577\,215 \dots$
δ	stepsize in numerical integration
$\delta(x)$	Dirac's famous function
ϵ	a tiny number, usually in the context of a limit to zero
$\zeta(x)$	the famous Riemann zeta function
\dots	\dots
$\psi(x)$	logarithmic derivative of the gamma function
ω	frequency

CHAPTER 1

THE FIRST

This is a chapter. Remember that there should *always* be at least of few lines of prose after each sectional heading: failure to do so is a disservice to your readers, and also produces incorrect vertical spacing.

1.1 The first section

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

In **Figure 1.1** on the following page, we have a picture, and the \LaTeX markup to include it looks like this:

```
\begin{figure}[t]
  \centerline{\includegraphics{fig1}}
  \caption{The first figure.}%
  \figlabel{fig1}
\end{figure}
```

We intentionally omitted an extension on the filename, so that this document can be processed with `latex` to get an output `.dvi` file, or with `pdflatex` to get an output `.pdf` file. The first case uses the file `fig1.eps`, and the second uses `fig1.pdf`. The `distill` or `ps2pdf` commands can be used to convert from *Encapulated PostScript* files to *Portable Document Format* files.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.



Figure 1.1. The first figure.

1.1.1 The first subsection

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

1.1.2 The second subsection

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

1.1.3 The third subsection

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

1.1.3.1 The first subsubsection

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

1.1.3.2 The second subsubsection

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

1.1.3.2.1 The first numbered paragraph Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah.

blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah.

1.1.3.2.2 The second numbered paragraph Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah.

1.2 The second section

Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah.

In **Figure 1.2**, we have another picture.

Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah.



Figure 1.2. The second figure.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

In **Table 1.1**, we show the 24-character lowercase Greek alphabet.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah.

Table 1.1. Lowercase Greek letters.

α	alpha
β	beta
γ	gamma
δ	delta
ϵ, ε	epsilon
ζ	zeta
η	eta
θ, ϑ	theta
ι	iota
κ	kappa
λ	lambda
μ	mu
ν	nu
ξ	xi
\omicron	omicron
π	pi
ρ	rho
σ, ς	sigma
τ	tau
υ	upsilon
ϕ, φ	phi
χ	chi
ψ	psi
ω	omega

Table 1.2. Uppercase Greek letters. Notice that several have the same letter shapes as Latin letters, and for those, \TeX does not define macro names. For convenience, we supply our own definitions of these macros: `\Alpha`, `\Beta`, `\Epsilon`, `\Zeta`, `\Eta`, `\Iota`, `\Kappa`, `\Mu`, `\Nu`, `\Omicron`, `\Rho`, `\Tau`, and `\Chi`.

A	Alpha
B	Beta
Γ	Gamma
Δ	Delta
E	Epsilon
Z	Zeta
H	Eta
Θ	Theta
I	Iota
K	Kappa
Λ	Lambda
M	Mu
N	Nu
Ξ	Xi
O	Omicron
Π	Pi
P	Rho
Σ	Sigma
T	Tau
Y	Upsilon
Φ	Phi
X	Chi
Ψ	Psi
Ω	Omega

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah.
 Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

1.4 Free software packages

The Free Software Foundation offers almost 300 software packages, most easily portable to many different operating systems and CPU platforms. They include at least these:

a2ps, acct, acm, adns, alive, anubis, apl, archimedes, aris, aspell, auctex,
 autoconf-archive, autoconf, autogen, automake, avl, ballandpaddle, barcode, bash,
 bayonne, bc, binutils, bison, bool, bperl2owfn, c-graph, ccaudio, ccd2cue, ccrtcp,

ccscript, cfengine, cflow, cgicc, chess, cim, classpath, classpathx, clisp, combine, commoncpp, complexity, config, coreutils, cpio, cppi, cssc, cursynth, dap, datamash, ddd, ddrescue, dejagnum, denemo, dico, diction, diffutils, dionysus, direvent, dismal, dominion, easejs, ed, edma, electric, emacs, emms, enscript, fdisk, ferret, findutils, fisicalab, flex, fontutils, freedink, freefont, freeipmi, gama, garpd, gawk, gcal, gcc, gcide, gcl, gcompris, gdb, gdbm, gengen, gengetopt, gettext, gforth, ggradebook, ghostscript, gift, gleem, glibc, global, glpk, gmp, gnash, gnats, gnatsweb, gnu-c-manual, gnu-crypto, gnu-pw-mgr, gnubatch, gnubik, gnucap, gnucobol, gnudos, gnue, gnugo, gnuit, gnupop, gnuplot, gnupod, gnuprologjava, gnuradio, gnurobots, gnuschool, gnushogi, gnusound, gnuspeech, gnuspool, gnustep, gnutls, gnutrusion, gnuzilla, goptical, gperf, gprolog, greg, grep, groff, grub, gsasl, gsegrafx, gsl, gslip, gsrc, gss, gtypist, guile-gnome, guile-gtk, guile-ncurses, guile-opengl, guile-rpc, guile-sdl, guile, gv, gvpe, gxmessage, gzip, halifax, health, hello, help2man, hp2xx, httptunnel, hurd, hyperbole, idutils, ignuit, indent, inetutils, intlfonts, jacal, jel, jwhois, kawa, less, libcdio, libextractor, libffcall, libiconv, libidn, libmatheval, libmicrohttpd, librejs, libsigsegv, libtasn1, libtool, libunistring, libxmi, lightning, lilypond, liquidwar6, lsh, m4, macchanger, mailman, mailutils, make, marst, maverik, mc, mcron, mcsim, mdk, metahtml, mifluz, mig, miscfiles, mit-scheme, moe, motti, mpc, mpfr, mpria, mtools, myserver, nano, ncurses, nettle, non-gnu, ocrad, octave, oleo, orgadoc, osip, paperclips, parallel, parted, patch, pem, pexec, phantom, pies, plotutils, proxyknife, pspp, psychosynth, pth, pyconfigure, radius, rcs, readline, recutils, reftex, remotecontrol, rottlog, rpge, rush, sather, sauce, savannah, scm, screen, sed, serveez, sharutils, shishi, shmm, shtool, sipwitch, slib, smalltalk, solfege, spacechart, spell, sqlltutor, src-highlight, stow, superopt, swbis, tar, termcap, termutils, tesc, teximpatient, texinfo, thales, time, tramp, trueprint, unifont, units, unrar, userv, uucp, vc-dwim, vcdimager, vera, wb, wdiff, websocket4j, wget, which, windows, xaos, xboard, xhippo, xlogmaster, xnee, xorriso, and zile.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.
 Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
 blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

1.5 Resizing figures

In **Figure 1.4** through **Figure 1.8** on the following page, we show how graphics files can be rescaled to convenient sizes, with input like this:

```
\begin{figure}[p]
  \centerline{\includegraphics[scale = 0.5]{fig1}}
  \caption{The fourth figure (at 50\% scale).}%
  \figlabel{fig4}
\end{figure}

\begin{figure}[p]
  \centerline{\includegraphics[scale = 0.75]{fig1}}
  \caption{The fifth figure (at 75\% scale).}%
  \figlabel{fig5}
\end{figure}
```



This is Figure 1

Figure 1.4. The fourth figure (at 50% scale).



This is Figure 1

Figure 1.5. The fifth figure (at 75% scale).



This is Figure 1

Figure 1.6. The sixth figure (at native size).



This is Figure 1

Figure 1.7. The seventh figure (at 125% scale).



This is Figure 1

Figure 1.8. The eighth figure (at 175% scale).

```

\begin{figure}[p]
  \centerline{\includegraphics{fig1}}
  \caption{The sixth figure (at native size).}%
  \figlabel{fig6}
\end{figure}

\begin{figure}[p]
  \centerline{\includegraphics[scale = 1.25]{fig1}}
  \caption{The seventh figure (at 125\% scale).}%
  \figlabel{fig7}
\end{figure}

\begin{figure}[p]
  \centerline{\includegraphics[scale = 1.75]{fig1}}
  \caption{The eighth figure (at 175\% scale).}%
  \figlabel{fig8}
\end{figure}

```

You can include multiple images, each with its own caption inside a single *unbreakable* figure environment, like this example shown in **Figure 1.9** and **Figure 1.10** on the next page, although you might want to adjust interfigure vertical space with a `\vspace{}` command:

```

\begin{figure}[t]
  \centerline{\includegraphics[scale = 0.5]{fig1}}
  \caption{The fourth figure (at 50\% scale).}%
  \figlabel{fig9}
  \vspace{3ex}
  \centerline{\includegraphics[scale = 0.75]{fig1}}
  \caption{The fifth figure (at 75\% scale).}%
  \figlabel{fig10}
\end{figure}

```

Blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah.



This is Figure 1

Figure 1.9. The ninth figure (at 50% scale), boxed with the tenth figure.



This is Figure 1

Figure 1.10. The tenth figure (at 75% scale), boxed with the ninth figure.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
blah blah blah blah blah blah blah blah blah blah blah blah blah blah.

As a final example in this chapter, **Figure 1.11** on the following page shows how you can use \LaTeX picture mode for annotating and positioning graphics images prepared outside \LaTeX . The input that produced that figure looks like this:

```
\begin{figure}[t]
  %% The original image is 216bp wide by 72bp high, but we
  %% rescale it to 150 picture units divided by \unitlength:
  %% 150 / 0.75 = 112.5 mm
  \newcommand {\myfig} {\includegraphics[width = 112.5mm]{fig1}}

  \begin{center}
    %% The \unitlength is chosen to make the complete picture fit
    %% within the page margins

    \setlength{\unitlength}{0.75mm}

    %%%      insert (width,height)(lower-left-x,lower-left-y)
    \begin{picture}(170,70)(10,10)
      %% Place the included image FIRST!
      \put(10,10) {\myfig}

      %% Everything that follows OVERLAYS the original image!

      \graphpaper[10](0,0)(170,70)

      %% Mark the image center and corners by centered bullets
      \newcommand {\thetodot} {\makebox (0,0) {$\bullet$}}
      \put( 85, 35) {\thetodot}
      \put( 10, 10) {\thetodot}
```


blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah
blah blah blah blah blah blah blah blah blah blah blah.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
blah blah blah blah blah blah blah blah blah blah blah.

CHAPTER 2

THE SECOND

2.1 Introduction

Understanding and debugging HPC programs is time-consuming for the user and computationally inefficient. This is especially true when one has to track control flow in terms of function calls and returns that may span library and system codes. Traditional software engineering quality assurance methods are often inapplicable to HPC where concurrency combined with large problem scales and sophisticated domain-specific math can make programming very challenging. For example, it took months for scientists to debug an MPI laser-plasma interaction code [?].

HPC bugs may be a combination of both flawed program logic and unspecified or illegal interactions between various concurrency models (e.g., PThreads, MPI, OpenMP, etc.) that coexist in most large applications [?]. Moreover, HPC software tends to consume vast amounts of CPU time and hardware resources. Reproducing bugs by rerunning the application is therefore expensive and undesirable. A natural and field-proven approach for debugging is to capture detailed execution traces and compare the traces against corresponding traces from previous (stable) runs [?,?]. A *key requirement* is to do this collection *as efficiently as possible* and in *as general and comprehensive a manner* as possible.

Existing tools in this space do not meet our criteria for efficiency and generality. The highly acclaimed STAT [?] tool has helped isolate bugs based on building equivalence classes of MPI processes and spotting outliers. We would like to go beyond the capabilities offered by STAT and support the collection of *whole-program* traces that can then be employed by a gamut of back-end tools. Also, STAT is usually brought into the picture when a failure (e.g., a deadlock or hang) is encountered. We would like to move toward an “always on” collection regime, as we cannot anticipate when a failure will occur – or, more

importantly, *whether the failure will be reproducible*. There are no reported debugging studies on using STAT in continuous collection (“always on”) mode. In CSTG [?], the collection is orchestrated by the user around chosen collection points and employs heavy-weight unix backtrace calls. These again are different from PARLOT, where collection points would not be a priori chosen.

The thrust of the work in this paper is to avoid many of the drawbacks of existing tracing-based tools. We are interested in avoiding source-code modifications and recompilation — thus making binary instrumentation-based tools the only practical and widely deployable option. We also believe in the value of creating tools that are *portable across a wide variety of platforms*.

Our goal is to use *compression* for trace aggregation and to offer a generic and low-overhead tracing method that (1) collects dynamic call information during execution (all function calls and returns) for debugging, performance evaluation, phase detection [?], etc., (2) has low overhead, (3) and requires little tracing bandwidth. *Providing all these features in a single tool that operates based on binary instrumentation is an unsolved problem*. In this paper, we describe a new tracing approach that fulfills these requirements, which we implemented in our proof-of-concept PARLOT tool.

With PARLOT, users can easily build a host of post-processors to examine executions from many vantage points. For instance, they can write post-processors to detect unexpected (or “outlier”) executions. If needed, they can drill down and detect abnormal behaviors *even in the runtime and support library stack* such as MPI-level activities. In HPC, it is well-known (especially on newer machines) that bugs are often due to broken libraries (MPI, OpenMP), a broken runtime, or OS-level activities. Having a single low-overhead tool that can “X-ray” an application to this depth is a goal met by PARLOT— a unique feature in today’s tool eco-system.

To further motivate the need for whole-program function call traces, consider the expression $f() + g()$. In C, there is no sequence point associated with the $+$ operator [?]. If these function calls have inadvertent side-effects causing failure, it is important to know in which order $f()$ and $g()$ were invoked—something that is easy to discern using PARLOT’s traces. One may be concerned that such a tool introduces excessive execution slowdown. PARLOT goes to great lengths to minimize these overheads to a level that we believe most

users will find acceptable. The mindset is to “*pay a little upfront to dramatically reduce the number of overall debug iterations*”.

As proof of concept, we gathered preliminary results from using the PARLOT tracing mechanism to compare different runs. We injected various bugs into the MPI-related functions of ILCS [?], a parallelization framework for iterative local searches. We ran PARLOT on top of executions of buggy and bug-free versions of ILCS and collected traces. Since PARLOT’s traces maintain the order of the function calls, we were able to split the traces at multiple points of interest and to feed different chunks of traces to a Concept Lattice data structure [?] [?]. Having the totally ordered sequence of function calls of the whole program for each active process/thread enabled us to quickly narrow down the search space to locate the cause of the abnormal behavior in the buggy version of ILCS.

This paper does not pursue debugging per se but rather a thorough benchmarking of PARLOT. It makes the following main contributions:

- It introduces a new tracing approach that makes it possible to capture the whole-program call-return, call-stack, call-graph, and call-frequency information, including all library calls, for every thread and process of HPC applications at low overhead in both space and time.
- It describes a new incremental data compression algorithm to drastically reduce the required tracing bandwidth, thus enabling the collection of whole-program traces, which would be infeasible without on-the-fly compression.
- It presents PARLOT, a proof-of-concept tool that implements our compression-based low-overhead tracing approach. PARLOT is capable of instrumenting x86 applications at the binary level (regardless of the source language used) to collect whole-program call traces.

The remainder of this paper is organized as follows. Section ?? introduces the basic ideas and infrastructure behind PARLOT and other tracing tools. Section ?? describes the design of PARLOT in detail. Sections ?? and ?? present our evaluation of different aspects of PARLOT and compare it with a similar tool. Section ?? concludes the paper with a summary and future work.

2.2 Background and Related Work

Recording a log of events during the execution of an application is essential for better understanding the program behavior and, in case of a failure, to locate the problem. Recording this type of information requires instrumentation of the program either at the source-code or the binary-code level. Instrumenting the source code by adding extra libraries and statements to collect the desired information is easy for developers. However, doing so modifies the code and requires recompilation, often involving multiple different tools and complex hierarchies of makefiles and libraries, which can make this approach cumbersome and frustrating for users. Instrumenting an executable at the binary level using a tool is typically easier, faster, and less error prone for most users. Moreover, binary instrumentation is language independent, portable to any system that has the appropriate instrumentation tool installed, and provides machine-level insight into the behavior of the application.

2.2.1 Binary Instrumentation

Executables can be instrumented *statically*, where the additional code is inserted into the binary before execution, which results in a persistent modified executable, or *dynamically*, where the modification of the executable is not permanent. In dynamic binary instrumentation, code behavior can be monitored at runtime, making it possible to handle dynamically-generated and self-modifying code. Furthermore, it may be feasible to attach the instrumentation to a running process, which is particularly useful for long-running applications and infinite loops.

Many different tools for investigating application behavior have been designed on top of such Dynamic Binary Instrumentation (DBI) frameworks. For instance, Dyninst [?] provides a dynamic instrumentation API that gives developers the ability to measure various performance aspects. It is used in tools like Open-SpeedShop [?] and TAU [?] as well as correctness debuggers like STAT [?]. Moreover, VampirTrace [?] uses it to provide a library for collecting program execution logs.

Valgrind [?] is a shadow-value DBI framework that keeps a copy of every register and memory location. It provides developers with the ability to instrument system calls and instructions. Error detectors such as Memcheck [?] and call-graph generators like

CALLGRIND [?] are built upon Valgrind.¹

We implemented PARLOT on top of PIN [?], a DBI framework for the IA-32, x86-64, and MIC instruction-set architectures for creating dynamic program analysis tools. There is also version of PIN available for the ARM architecture [?]. PARLOT mutates PIN to trace the entry (call) and exit (return) of every executed function. Note that our tracing and compression approaches can equally be implemented on top of other instrumentation tools. For example, PMaC [?] is a DBI tool for the PowerPC/AIX architecture upon which PARLOT could also be based.

2.2.2 Efficient Tracing

When dealing with large-scale parallel programs, any attempt to capture reasonably frequent events will result in a vast amount of data. Moreover, transferring and storing the data will incur significant overhead. For example, collecting just one byte of information per executed instruction yields on the order of a gigabyte of data per second on a single high-end core. Storing the resulting multi-gigabyte traces from many cores can be a challenge, even on today’s large hard disks.

Hence, to be able to capture whole-program call traces, we need a way to decrease the space and runtime overhead. *Compression* can encode the generated data using a smaller number of bits, help reduce the amount of data movement across the memory hierarchy, and lower storage and network demands. Although the encoded data will later have to be decoded for analysis, compressing them during tracing enables the collection of *whole-program* traces.

The use of compression by itself is not new. Various performance evaluation tools [?, ?, ?] already employ compression during the collection of performance analysis data. Tools such as ScalaTrace [?] also exploit the repetitive nature of time-step simulations [?]. Aguilar et al. [?] proposed a lossy compression mechanism using the Nami library [?] for online MPI tracing. Mohror and Karavanic [?] investigated similarity-based trace reduction techniques to store and analyze traces at scale.

¹Given the absence of tools similar to PARLOT, we employ CALLGRIND as a “close-enough” tool in our comparisons elaborated in §2.4.3. In this capacity, CALLGRIND is similar to PARLOT(M), a variant of PARLOT that only collects traces from the main image. We perform such comparison to have an idea of how we fare with respect to one other tool. In §??, we also present a self-assessment of PARLOT separately.

Many performance and debugging tools for HPC applications [?,?] rely on mechanisms such as MRNet [?] to accelerate the collection and aggregation of traces based on an overlay network to overcome the challenge of massive data movement and analysis. However, our experiments show that, due to the high compression ratio of PARLOT traces, such mechanisms for data movement and aggregation may be unnecessary.

The novelty offered by PARLOT lies in the combination of compression speed, efficacy, and low timing jitter made possible by its *incremental* lossless compression algorithm, which is described in §??. It immediately compresses all traced information while the application is running, that is, PARLOT does not record the uncompressed trace in memory. As a result, just a few kilobytes of data need to be written out per thread and per second, thus requiring only a small fraction of the disk or network bandwidth. The traces are decompressed later when they are read for offline analysis. From the decompressed full function-call trace, the complete call-graph, call-frequency, and caller-callee information can be extracted. This can be done at the granularity of a thread, a group of threads, or the whole application. We now elaborate on the design of PARLOT that makes these innovations possible.

2.3 Design of PARLOT

Our experimental results in §?? highlight why *compression* is essential to make our approach work. We used PARLOT to record a unique 16-bit identifier for every function call and return. Tracing just this small amount of information without compression when running the Mantevo miniapps [?] on Stampede 1 resulted in about 2 MB/s of data per core on average. Extrapolating this value to all 102,400 cores of Stampede 1 (not counting the accelerators) yields 205 GB/s of trace data, which exceeds the Lustre filesystem’s parallel write performance of 150 GB/s. Enabling PARLOT’s compression algorithm reduced the emitted trace data by a factor of 100 on average, a ratio that is quite stable w.r.t scaling, making it possible to trace full-scale programs while leaving over 98% of the I/O bandwidth to the application. Therefore, PARLOT should also work for codes with higher bandwidth requirements than the ones we tested.

Figure 2.1 provides a general overview of PARLOT’s workflow. Basic blocks within program executables are *dynamically* instrumented before being executed. The collected

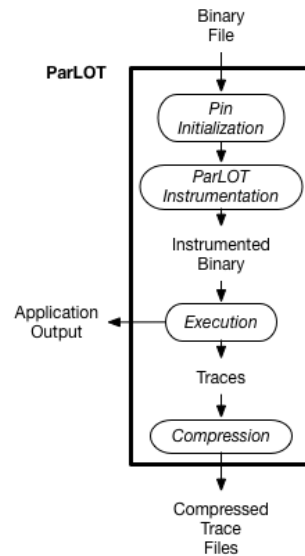


Figure 2.1. Overview of PARLOT

data are compressed on-the-fly at runtime.

2.3.1 Tracing Operation

PARLOT uses the PIN API as its instrumentation mechanism to gather traces. In particular, it instructs PIN to instrument every thread launch and termination in the application as well as every function entry and exit. The thread-launch instrumentation code initializes the per-thread tracing variables and opens a file into which the trace data from that thread will be written. The thread-termination code finalizes any ongoing compression, flushes out any remaining entries, and closes the trace file. PARLOT assigns every static function in each image (main program and all libraries) a unique unsigned 16-bit ID, which it records in a separate file together with the image and function name. This file allows the trace reader to map IDs back to function-name/image pairs.

For every function *entry*, PARLOT executes extra code that has access to the thread ID, function ID, and current stack-pointer (SP) value. Based on the SP value, it performs call-stack correction if necessary (see §2.3.4), adds the new function to a data structure it maintains that holds the call stack (which is separate from the application’s runtime stack), and emits the function ID into the trace file via an incremental compression algorithm (see §2.3.2). All of this is done independently for each thread. Similarly, for every function *exit*,

PARLOT also executes extra code that has access to the thread ID, function ID, and current SP value. Based on the SP value, it performs call-stack correction if necessary, removes the function from its call-stack data structure, and emits the reserved function ID of zero into the trace file to indicate an exit. As before, this is done via an incremental compression algorithm. We use zero for all exits rather than emitting the function ID and a bit to specify whether it is an entry or exit because using zeros results in more compressible output. This way, half of the values in the trace will be zero.

2.3.2 Incremental Compression

PARLOT immediately compresses the traced information even before it is written to memory. It does, however, keep a sliding window (circular buffer) of the most recent uncompressed trace events, which is needed by the compressor. It compresses each function ID before the next function ID is known. The conventional approach would be to first record uncompressed function IDs in a buffer and later compress the whole buffer once it fills up. However, this makes the processing time very non-uniform. Whereas almost all function IDs can be recorded very quickly since they just have to be written to the buffer, processing a function ID that happens to fill the buffer takes a long time as it triggers the compression of the entire buffer. This results in sporadic blocking of threads during which time they make no progress towards executing the application code. Initial experiments revealed that such behavior can be detrimental when one thread is polling data from another thread that is currently blocked due to compression. For example, we observed a several order of magnitude increase in entry/exit events of an internal MPI library function when using block-based compression.

To remedy this situation, the compressor must operate incrementally, i.e., each piece of trace data must be compressed when it is generated, without buffering it first, to ensure that there is never a long-latency compression delay. Few existing compression algorithms have been implemented in such a manner because it is more difficult to code up and probably a little slower. Nevertheless, we were able to implement our algorithm (discussed next) in this way so that each trace event is compressed with similar latency.

2.3.3 Compression Algorithm

We used the CRUSHER framework [?, ?, ?, ?] to automatically synthesize an effective and fast lossless compression algorithm for our traces. CRUSHER is based on a library of data transformations extracted from various compression algorithms. It combines these transformations in all possible ways to generate algorithm candidates, which it then evaluates on a set of training data. We gathered uncompressed traces from some of the Manteco miniapps [?] for this purpose. This evaluation revealed that a particular word-level Lempel-Ziv (LZ) transformation followed by a byte-level Zero-Elimination (ZE) transformation works well. In other words, PARLOT's trace entries, which are two-byte words, are first transformed using LZ. The output is interpreted as a sequence of bytes, which is transformed using ZE for further compression. The output of ZE is written to secondary storage.

LZ implements a variant of the LZ77 algorithm [?]. It uses a 4096-entry hash table to identify the most recent prior occurrence of the current value in the trace. Then it checks whether the three values immediately before that location match the three trace entries just before the current location. If they do not, the current trace entry is emitted and LZ advances to the next entry. If the three values match, LZ counts how many values following the current value match the values following that location. The length of the matching substring is emitted and LZ advances by that many values. Note that all of this is done incrementally. The history of previous trace entries available to LZ for finding matches is maintained in a 64k-entry circular buffer.

ZE emits a bitmap in which each bit represents one input byte. The bits indicate whether the corresponding bytes are zero or not. Following each eight-bit bitmap, ZE emits the non-zero bytes.

As mentioned above, we had to implement the two transformations incrementally to minimize the maximum latency. This required breaking them up into multiple pieces. Depending on the state the compressor is in when the next trace entry needs to be processed, the appropriate piece of code is executed and the state updated. If the LZ code produces an output, which it only does some of the time, then the appropriate piece of the ZE code is executed in a similar manner.

2.3.4 PIN and Call-Stack Correction

To be able to decode the trace, i.e., to correctly associate each exit with the function entry it belongs to, our trace reader maintains an identical call-stack data structure. Unfortunately, and as pointed out in the PIN documentation [?], it is not always possible to identify all function exits. For example, in optimized code, a function's instructions may be inlined and interleaved with the caller's instructions, making it sometimes infeasible for PIN to identify the exit. As a consequence, we have to ensure that PARLOT works correctly even when PIN misses an exit. This is why the SP values are needed.

During tracing, PARLOT not only records the function IDs in its call stack but also the associated SP values. This enables it to detect missing exits and to correct the call stack accordingly. Whenever a function is entered, it checks if there is at least one entry in the call stack and, if so, whether its SP value is higher than that of the current SP. If it is lower, we must have missed at least one exit since the runtime stack grows downwards (the SP value decreases with every function entry and increases with every exit). If a missing exit is detected in this manner, PARLOT pops the top element from its call stack and emits a zero to indicate a function exit. It repeats this procedure until the stack is empty or its top entry has a sufficiently high SP value. The same call-stack correction technique is applied for every function exit whose SP value is inconsistent. Note that the SP values are only used for this purpose and are not included in the compressed trace.

The result is an internally consistent trace of function entry and exit events, meaning that parsing the trace will yield a correct call stack. This is essential so that the trace can be decoded properly. Moreover, it means that the trace includes exits that truly happened in the application but that were missed by PIN. Note, however, that our call-stack correction is a best-effort approach and may, in rare cases, temporarily not reflect what the application actually did. For example, this can happen for functions that do not create a frame on the runtime stack. When implementing PARLOT on top of another DBI framework, call-stack correction may not be needed, resulting in even lower PARLOT overhead.

2.4 Evaluation Methodology

2.4.1 Benchmarks and System

We performed our evaluations on the MPI-based NAS Parallel Benchmarks (NPB) [?]. NPB includes four inputs sizes. To keep the runtimes reasonable, we show results for the class *B* (small-medium) and class *C* (medium-large) inputs.

We compiled the NPB codes with the mpicc and mpif77 wrappers of MVAPICH 2.2.1, which are based on icc/ifort 14.0.2 using the prescribed -g and -O1 optimization flags. Quick tests showed that higher optimization levels do not significantly improve the performance.

We ran all experiments on Comet at the San Diego Supercomputer Center [?], whose filesystem is NFS and Lustre. Comet has 1944 compute nodes, each of which has dual-socket Intel Xeon E5-2680 v3 processors with a total of 28 cores (14 per socket) and 128 GB of main memory. Note that we only used 16 cores per node as many of the NPB programs require a core count that is a power of two. To study the scaling behavior, we ran experiments on 1, 4, 16 and 64 compute nodes, i.e., on up to 1024 cores.

2.4.2 Metrics

We use the following metrics to quantify and compare the performance of the tracing tools. Unless otherwise noted, all results are based on the median of three identical experiments.

- The **tracing overhead** is the runtime of the target application when it is being traced divided by the runtime of the same application without tracing. This lower-is-better ratio measures by how much the tracing (and the compression when enabled) slows down the target application.
- The **tracing bandwidth** is the size of the trace information divided by the application runtime. To make the results easier to compare, we generally list the tracing bandwidth per core, i.e., the tracing bandwidth divided by the number of cores used. This lower-is-better metric is expressed in kilobytes per second (kB/s) per core. It specifies the average needed bandwidth to record the trace data.
- The **compression ratio** is the size of the uncompressed trace divided by the size of

the generated (compressed) trace. This higher-is-better ratio measures the factor by which the compression reduces the trace size. In other words, without compression, the tracing bandwidth would be higher by this factor.

2.4.3 Tracing Tools

We compare our PARLOT tool, implemented on top of PIN 3.5, with CALLGRIND 3.13. PARLOT was compiled with gcc 4.9.2 using PIN’s make system and CALLGRIND with Valgrind’s make system. We created the following versions of PARLOT to evaluate different aspects of its design.

- **PARLOT(M)** is the normal PARLOT tool configured to only collect function-call traces

Table 2.1. Overhead added by each tool

Input	Tool	# Nodes	bt	cg	ep	ft	is	lu	mg	sp	GM
B	PARLOT(M)	1	1.6	1.8	2.6	2.1	2.5	1.3	2.5	1.3	1.9
		4	1.8	1.9	1.9	1.7	1.8	1.8	1.5	1.7	1.8
		16	2.2	2.6	2.0	1.9	1.8	2.7	2.4	2.2	2.2
		64	2.1	2.2	2.4	2.0	4.3	4.4	2.0	2.1	2.5
		AVG	1.9	2.1	2.2	1.9	2.6	2.6	2.1	1.8	2.1
	PARLOT(A)	1	1.8	2.7	4.2	2.8	4.2	1.7	4.8	1.7	2.8
		4	2.6	3.1	3.4	2.8	3.0	2.8	2.8	2.7	2.9
		16	3.5	4.2	3.4	2.9	2.8	4.3	4.5	3.7	3.6
		64	3.1	3.3	3.8	3.0	5.4	4.7	3.2	3.3	3.7
		AVG	2.8	3.3	3.7	2.9	3.9	3.4	3.8	2.8	3.2
	CALLGRIND	1	8.6	6.0	8.9	10.1	2.5	7.5	3.3	6.6	6.1
		4	6.0	3.6	2.9	3.5	1.5	5.2	1.2	5.8	3.2
		16	4.3	3.3	2.2	2.2	1.7	4.6	1.8	4.3	2.8
		64	2.3	2.0	1.7	2.1	4.1	4.0	1.5	2.5	2.3
		AVG	5.3	3.7	3.9	4.5	2.4	5.3	2.0	4.8	3.6
C	PARLOT(M)	1	1.4	1.3	2.5	1.9	2.3	1.1	1.7	1.1	1.6
		4	1.6	1.7	1.8	1.6	1.7	1.3	1.8	1.4	1.6
		16	1.8	2.4	2.5	1.5	1.8	2.2	2.4	1.8	2.0
		64	2.2	2.7	2.4	1.6	4.5	3.4	2.4	2.2	2.6
		AVG	1.8	2.0	2.3	1.7	2.6	2.0	2.1	1.6	1.9
	PARLOT(A)	1	1.5	1.6	3.2	2.0	2.8	1.2	2.5	1.2	1.9
		4	1.9	2.4	2.6	2.1	2.6	1.7	3.1	1.7	2.2
		16	2.7	3.5	4.1	2.1	2.8	3.2	4.0	2.5	3.0
		64	3.6	4.1	4.2	2.2	5.5	4.4	4.2	3.0	3.8
		AVG	2.4	2.9	3.5	2.1	3.4	2.6	3.5	2.1	2.7
	CALLGRIND	1	8.5	4.4	13.2	13.1	3.3	7.9	5.9	5.1	6.9
		4	8.7	4.5	4.8	6.4	1.7	6.4	2.8	6.3	4.6
		16	6.9	3.9	3.1	2.8	1.8	6.4	2.1	6.1	3.7
		64	4.4	3.5	2.1	2.5	4.2	5.2	2.1	3.8	3.3
		AVG	7.1	4.1	5.8	6.2	2.8	6.5	3.2	5.3	4.6

from the main image of the application.

- **PARLOT(A)** is the normal PARLOT tool configured to collect function-call traces from all images of the application, including library function calls.
- **PIN-INIT** is a crippled version of PARLOT from which the tracing code has been removed. The purpose of PIN-INIT is to see how much of the overhead is due to PIN.
- **PARLOT-NC** is the normal PARLOT tool but with compression disabled. It writes out the captured data in uncompressed form. The purpose of PARLOT-NC is to show the performance impact of the compression.

It proved surprisingly difficult to find a tool that is similar to PARLOT because there appear to be no other tools that generate whole program call traces. In the end, we settled on CALLGRIND as the most similar tool we could find and used it for our comparisons. CALLGRIND is based on the Valgrind DBI tool. It collects function-call graphs combined with performance data to show the user what portion of the execution time has been spent in each function.

Each CALLGRIND trace file contains a sequence of function names (or their code) plus numerical data for each function on its caller-callee relationship with other functions. Moreover, it contains cost information for each function in terms of how many machine instructions it read. This information is collected using hardware performance counters. The format of the file is plain ASCII text. Interestingly, all numerical values are expressed relative to previous values, i.e., they are delta (or difference) encoded. This simple form of compression is enabled by default in CALLGRIND.

We believe the information traced by CALLGRIND is reasonably similar to the information traced by PARLOT(M). Whereas CALLGRIND's traces include performance data that PARLOT does not capture, PARLOT records the whole-program call trace, which CALLGRIND does not capture. The full function-call trace is a strict superset of the call-graph information that CALLGRIND records because the call graph can be extracted from the function-call trace but not vice versa. In particular, CALLGRIND cannot recreate the order of the function calls a thread made whereas PARLOT can.

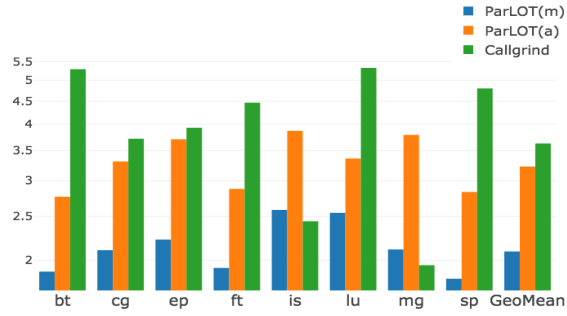


Figure 2.2. Average tracing overhead on the NPB applications - Input B

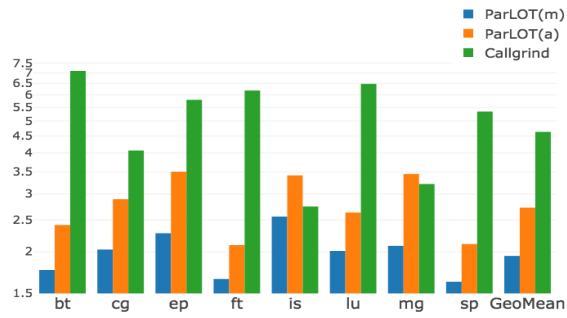


Figure 2.3. Average tracing overhead on the NPB applications - Input C

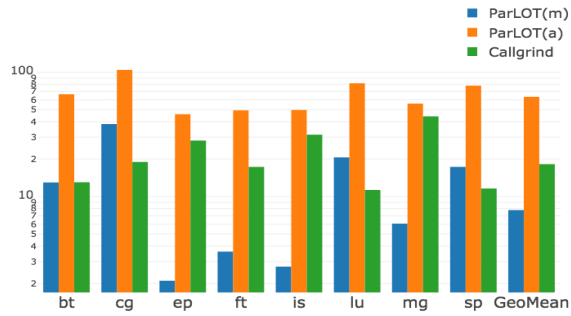


Figure 2.4. Average required bandwidth per core (kB/s) on the NPB applications - Input B

2.5 Results

2.5.1 Tracing Overhead

Table 2.1 shows the tracing overhead of PARLOT(M), PARLOT(A), and CALLGRIND on each application of the NPB benchmark suite for different node counts. The last column

Table 2.2. Required bandwidth per core (kB/s)

Input	Tool	# Nodes	bt	cg	ep	ft	is	lu	mg	sp	GM
B	PARLOT(M)	1	4.7	21.9	3.8	1.5	0.8	2.4	5.6	5.4	3.7
		4	14.3	41.1	1.9	3.5	2.2	21.5	6.5	15.9	8.1
		16	14.3	46.6	1.5	4.9	3.4	31.8	6.5	18.6	9.4
		64	18.6	43.6	1.3	4.6	4.5	27.1	5.6	29.6	9.9
		AVG	13.0	38.3	2.1	3.6	2.7	20.7	6.1	17.4	7.8
	PARLOT(A)	1	48.7	89.4	47.2	45.6	60.0	53.6	60.8	54.3	56.2
		4	61.8	101.2	45.2	55.1	53.2	71.1	54.9	73.6	62.7
		16	74.0	116.9	47.4	48.9	47.8	100.9	55.8	84.6	68.0
		64	81.8	110.2	44.2	48.0	37.8	100.3	52.7	99.9	66.5
		AVG	66.6	104.4	46.0	49.4	49.7	81.5	56.0	78.1	63.3
	CALLGRIND	1	1.6	7.7	7.4	4.6	39.5	2.6	34.4	2.7	6.7
		4	6.5	16.0	22.1	15.7	45.5	8.6	45.5	7.8	16.3
		16	17.2	24.6	37.4	23.8	29.9	16.2	51.5	15.8	24.9
		64	26.8	27.7	45.9	25.1	11.0	17.8	45.3	20.2	25.0
		AVG	13.0	19.0	28.2	17.3	31.5	11.3	44.2	11.6	18.2
C	PARLOT(M)	1	1.8	17.0	5.2	1.2	0.7	0.8	3.6	1.4	2.2
		4	7.5	44.9	3.0	2.5	2.1	20.1	7.1	13.7	7.6
		16	16.3	55.0	1.8	6.1	3.4	34.1	7.2	20.7	10.7
		64	17.5	61.4	1.3	5.9	4.4	38.3	5.6	26.1	10.9
		AVG	10.8	44.6	2.8	3.9	2.7	23.3	5.9	15.5	7.8
	PARLOT(A)	1	17.8	53.4	26.3	20.9	48.3	25.3	52.6	19.5	30.0
		4	51.8	95.8	36.8	43.8	51.4	58.4	54.2	65.8	55.2
		16	75.4	121.0	44.3	61.4	46.9	101.1	56.5	101.3	71.4
		64	80.6	135.2	43.5	46.3	37.1	117.9	54.1	99.0	69.0
		AVG	56.4	101.4	37.7	43.1	45.9	75.7	54.3	71.4	56.4
	CALLGRIND	1	0.4	3.1	2.0	1.1	14.6	0.7	7.0	0.8	1.9
		4	1.8	8.9	7.7	4.5	31.7	2.8	21.0	2.8	6.4
		16	6.0	15.8	22.9	10.8	26.5	7.5	39.1	7.0	13.7
		64	14.3	19.6	35.8	12.2	11.1	11.9	40.7	12.8	17.4
		AVG	5.6	11.8	17.1	7.1	21.0	5.7	26.9	5.8	9.8

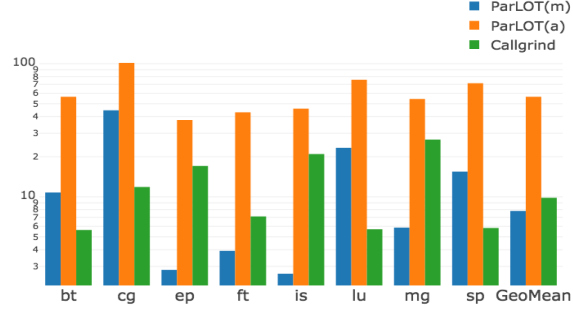


Figure 2.5. Average required bandwidth per core (kB/s) on the NPB applications - Input C

of the table lists the geometric mean over all eight programs. The AVG rows show the average over the four node counts.

On average, both PARLOT(M) and PARLOT(A) outperform CALLGRIND. The bolded numbers in Table 2.1 for input C show that the average overhead is 1.94 for PARLOT(M), 2.73 for PARLOT(A), and 4.63 for CALLGRIND. Figures 2.2 and 2.3 show these results in visual form.

The key takeaway point is that the overhead of PARLOT is roughly a factor of two to three, which we believe users may be willing to accept, for example, if it helps them debug their applications. This is promising especially when considering how detailed the collected trace information is and that most of the overhead is due to PIN (see §2.5.4). Note that PARLOT’s overhead is typically lower than that of CALLGRIND, which collects less information.

The overhead of PARLOT increases as we scale the applications to more compute nodes. However, the increase is quite small. Going from 16 to 1024 cores, a 64-fold increase in parallelism, only increases the average overhead by between 1.3- and 2.1-fold. In contrast, CALLGRIND’s overhead decreases with increasing node count, making it more scalable. Having said that, CALLGRIND’s overhead is larger for the C inputs whereas PARLOT’s overhead is larger for the smaller B inputs. In other words, PARLOT scales better to larger inputs than CALLGRIND.

PARLOT’s scaling behavior can be explained by correlating it with the expected function-call frequency. When distributing a fixed problem size over more cores, each core receives

less work. As a consequence, less time is spent in the functions that process the work, resulting in more function calls per time unit, which causes more work for PARLOT. In contrast, when distributing a larger problem size over the same number of cores, each core receives more work. Hence, more time is spent in the functions that process the work, resulting in fewer function calls per time unit, which causes less work for PARLOT and therefore less tracing overhead. Hence, we believe PARLOT's overhead to be even lower on long-running inputs, which is where our tracing technique is needed the most.

In summary, PARLOT's overhead is in the single digits for all evaluated applications and configurations, including for 1024-core runs. It appears to scale reasonably to larger node counts and well to larger problem sizes.

2.5.2 Required Bandwidth

Table 2.2, Fig. 2.4 and Fig. 2.5 show how much trace bandwidth each tool requires during the application execution. On average, PARLOT(M) requires less bandwidth than CALLGRIND, especially for smaller inputs. PARLOT(A)'s bandwidth is much higher as it collects call information from all images and not just the main image like PARLOT(M) does.

We see that the required bandwidth for different input sizes of the NPB applications are almost equal in PARLOT. According to the NPB documentation, the number of iterations for inputs B and C are the same for all applications. They only differ in the number of elements or the grid size. It is clear that the required bandwidth of PARLOT is independent of the problem size, unlike CALLGRIND, where the input size has a linear impact on the results.

2.5.3 Compression Ratio

Table 2.3 shows the compression ratios for all configurations and inputs. On average, PARLOT stores between half a kilobyte and a kilobyte of trace information in a single byte. We observe that the average compression ratio for PARLOT(A) on input C is 644.3, and its corresponding required bandwidth from Table 2.2 is 56.4 kB/s. This means PARLOT can collect **more than 36 MB** worth of data per core per second while only needing 56 kB/s of the system bandwidth, *leaving the rest of the available bandwidth to the application*. In comparison, CALLGRIND collects **less than 100 kB** of data but still adds more overhead

Table 2.3. Compression ratio

Input	Tool	# Nodes	bt	cg	ep	ft	is	lu	mg
B	PARLOT(M)	1	3 035.9	94.4	12 456.2	12 173.5	9 718.4	167.7	99.1
		4	586.6	82.5	10 368.4	1 737.1	909.2	140.3	255.0
		16	346.7	113.3	8 563.9	1 077.4	1 200.6	179.0	387.6
		64	252.2	147.8	7 611.0	1 122.6	1 908.0	366.8	437.3
		AVG	1 055.4	109.5	9 749.9	4 027.6	3 434.0	213.5	294.7
	PARLOT(A)	1	514.5	137.4	3 335.8	1 226.7	543.2	314.6	260.9
		4	315.7	137.2	1 266.9	436.2	316.2	287.3	329.6
		16	226.9	181.6	1 246.7	1 026.5	927.1	299.3	469.3
		64	329.2	247.3	1 394.1	1 043.9	1 984.6	410.3	548.5
		AVG	346.6	175.9	1 810.9	933.3	942.8	327.9	402.1
C	PARLOT(M)	1	8 619.0	111.2	13 068.0	21 335.6	21 856.5	350.0	247.4
		4	1 910.6	110.5	12 418.7	6 520.3	2 256.6	112.8	268.0
		16	580.8	133.2	11 017.4	1 239.3	1 347.9	164.5	396.9
		64	322.8	131.9	9 155.0	1 065.1	1 896.3	223.7	465.7
		AVG	2 858.3	121.7	11 414.7	7 540.1	6 839.3	212.7	344.5
	PARLOT(A)	1	2 579.4	181.8	7 377.0	5 143.1	1 520.4	408.2	314.8
		4	448.6	161.3	3 194.6	1 062.9	527.3	274.7	319.4
		16	285.1	185.7	1 765.5	588.9	1 106.3	273.6	467.4
		64	290.0	214.7	1 512.9	1 237.3	2 038.7	329.0	496.2
		AVG	900.8	185.9	3 462.5	2 008.1	1 298.2	321.4	399.4

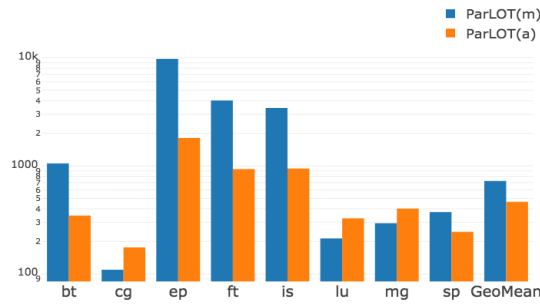


Figure 2.6. Average compression ratio of PARLOT on the NPB applications - Input B

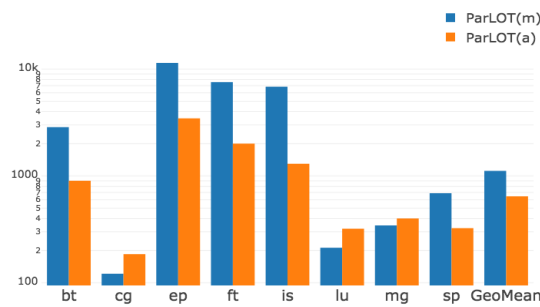


Figure 2.7. Average compression ratio of PARLOT on the NPB applications - Input C

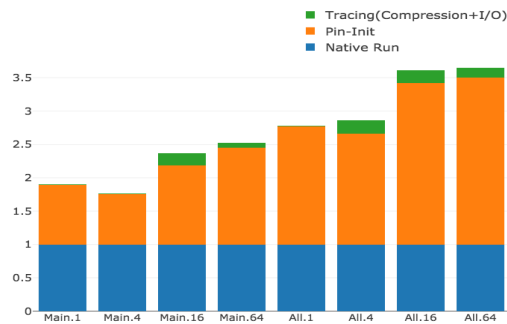


Figure 2.8. Tracing overhead breakdown - Input B

compared to either PARLOT(A) or PARLOT(M). The average amount of trace data that can be collected by PARLOT(A) is **360x** (85x for PARLOT(M)) larger than that for CALLGRIND. In the best observed case, the compression ratio of PARLOT exceeds 21000. This is particularly impressive because it was achieved with relatively low overhead and incremental

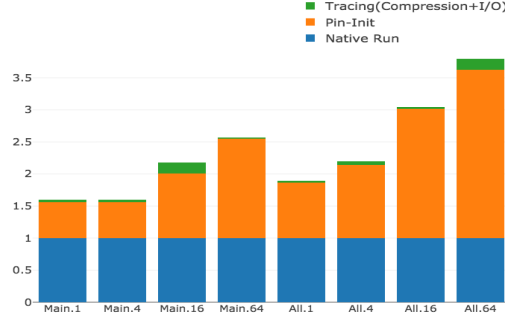


Figure 2.9. Tracing overhead breakdown - Input C

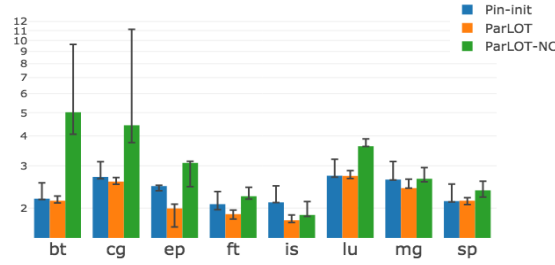


Figure 2.10. Variability of PARLOT(M) overhead on 16 nodes - Input B

on-the-fly compression. Generally, the compression ratios of PARLOT(M) are higher than those of PARLOT(A) because the variety of distinct function calls on the main image is smaller than when tracing all images, thus compression performs better on PARLOT(M). Also by looking at Fig. 2.4, Fig. 2.5, Fig. 2.6 and Fig. 2.7, we find EP to have the highest compression ratio of the NPB applications. At the same time, it has the minimum required bandwidth. The opposite is true for CG, which exhibits the lowest compression ratio and the highest required bandwidth. CG is a conjugate gradient method with irregular memory accesses and communications whereas EP is an embarrassingly parallel random number generator. CG's whole-program trace contains a larger number of distinct calls and more complex patterns than that of EP, thus resulting in a higher bandwidth and lower compression ratio.

PARLOT's compression mechanism works better on larger input sizes because larger inputs tend to result in longer streams of similar function calls (e.g., a call that is made for

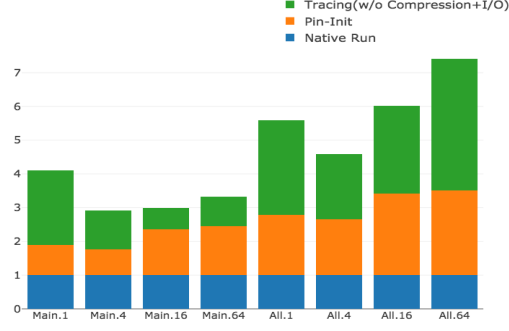


Figure 2.11. PARLOT-NC tracing overhead breakdown - Input B

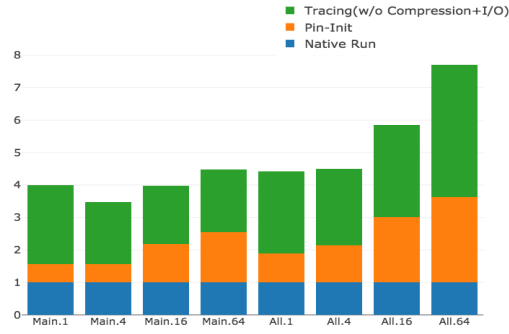


Figure 2.12. PARLOT-NC tracing overhead breakdown - Input C

every processed element).

2.5.4 Overheads

Tables 2.4 and 2.5 present the average overhead added to each application for different versions of PARLOT. The last row of these two tables presents the geometric mean. This information captures how much each phase of PARLOT slows down the native execution.

In general, one expects the following inequality to hold: the overhead of PIN-INIT should be less than that of PARLOT, which should be less than that of PARLOT-NC. This is not always the case because of the non-deterministic runtimes of the applications. In fact, the variability across three runs of each experiment is shown in Fig. 2.10 where we present the minimum, maximum and median overheads. These overheads are for input size B and 16 nodes. This variability explains the seeming inconsistencies in Tables 2.4 and 2.5.

On average, PIN-INIT adds an overhead of 3.28 and PARLOT(A) adds an overhead of

Table 2.4. Tracing overhead of versions of PARLOT(M)- Input B

Input: B	Nodes :	1			4			16		
	Detail Tools:	PIN-INIT	PARLOT	PARLOT-NC	PIN-INIT	PARLOT	PARLOT-NC	PIN-INIT	PARLOT	PARLOT-NC
Main	bt	1.5	1.5	5.6	1.7	1.7	5.0	2.1	2.1	5.0
	cg	1.7	1.8	2.3	1.8	1.8	2.6	2.7	2.5	2.6
	ep	2.9	2.6	20.4	1.9	1.8	5.3	2.4	1.9	5.3
	ft	1.8	2.1	6.1	1.7	1.7	2.7	2.0	1.8	2.7
	is	2.4	2.4	4.8	1.7	1.7	2.0	2.1	1.7	2.0
	lu	1.3	1.3	1.4	1.7	1.7	2.2	2.7	2.7	2.2
	mg	2.5	2.5	2.7	1.5	1.5	1.5	2.6	2.4	1.5
	sp	1.3	1.3	2.4	1.7	1.7	3.5	2.1	2.1	3.5
	GM	1.8	1.9	4.1	1.7	1.7	2.9	2.3	2.1	2.9

Table 2.5. Tracing overhead of versions of PARLOT(A)- Input B

Input: B	Nodes :	1			4			16		
	Detail Tools:	PIN-INIT	PARLOT	PARLOT-NC	PIN-INIT	PARLOT	PARLOT-NC	PIN-INIT	PARLOT	PARLOT-NC
All	bt	1.7	1.8	6.1	2.3	2.5	6.1	3.2	3.5	6.1
	cg	2.6	2.7	3.8	2.8	3.0	4.4	4.0	4.2	4.4
	ep	4.3	4.1	22.2	3.1	3.4	7.1	3.1	3.3	7.1
	ft	2.8	2.7	6.8	2.6	2.7	3.8	2.8	2.9	3.8
	is	4.4	4.2	7.0	2.8	2.9	3.4	2.9	2.8	3.4
	lu	1.7	1.7	2.3	2.5	2.7	4.8	3.9	4.3	4.8
	mg	4.8	4.7	5.3	2.5	2.7	3.0	4.3	4.4	3.0
	sp	1.7	1.7	3.0	2.4	2.6	5.0	3.2	3.6	5.0
	GM	2.7	2.7	5.5	2.6	2.8	4.5	3.4	3.6	4.5

3.42. This means that **almost 96% of PARLOT(A)’s overhead is due to PIN**. The results of PARLOT(M) and other inputs follow the same pattern as shown in Fig. 2.8 and 2.9. The overhead that PARLOT (excluding the overhead of PIN-INIT) *adds* to the applications is very small. If we were to switch to a different instrumentation tool that is not as general as PIN but more lightweight, the overhead would potentially reduce drastically.

2.5.5 Compression Impact

Fig. 2.11 and Fig. 2.12 show the overhead breakdown of PARLOT-NC, which illustrate the impact of compression. They also highlight the importance of incorporating compression directly in the tracing tool. On average, PARLOT-NC slows down the application execution almost **2x** more than PARLOT(A). The average overhead across Table 2.5 for PARLOT(A) is **3.4**. The corresponding factor for PARLOT-NC is **6.6**. The numbers of PARLOT(M) and input C follow the same pattern. For example, PARLOT-NC slows down the application execution almost **1.66x** more than PARLOT(M).

Clearly, compression not only lowers the storage requirement but also the overhead. This is important as it shows that the extra computation to perform the compression is more than amortized by the reduction in the amount of data that need to be written out.

This result validates our approach and highlights that incremental, on-the-fly compression is likely essential to make whole-program tracing possible at low overhead.

2.6 Discussion and Conclusion

In this paper, we present PARLOT, a portable low overhead dynamic binary instrumentation-based whole-program tracing approach that can support a variety of dynamic program analyses, including debugging. Key properties of PARLOT include its on-the-fly trace collection and compression that reduces timing jitter, I/O bandwidth, and storage requirements to such a degree that whole-program call/return traces can be collected efficiently even at scale.

We evaluate various versions of PARLOT created by disabling/enabling compression, not collecting any traces, etc. In order to provide an intuitive comparison against a well known tool, we also compare PARLOT to CALLGRIND. Our metrics include the tracing

overhead, required bandwidth, achieved compression ratio, initialization overhead, and the overall impact of compression. Detailed evaluations on the NAS parallel benchmarks running on up to 1024 cores establish the merit of our tool and our design decisions. PARLOT can collect more than 36 MB worth of data per core per second while only needing 56 kB/s of bandwidth and slowing down the application by 2.7x on average. These results are highly promising in terms of supporting whole program tracing and debugging, in particular when considering that most of the overhead is due to the DBI tool and not PARLOT.

The traces collected by PARLOT cut through the entire stack of heterogeneous (MPI, OpenMP, PThreads) calls. This permits a designer to project these traces onto specific APIs of interest during program analysis, visualization, and debugging.

A number of improvements to PARLOT remain to be made. These include allowing users to selectively trace at specific interfaces: doing so can further increase compression efficiency by reducing the variety of function calls to be handled by the compressor. We also discuss the need to bring down initialization overheads, i.e., by switching to a less general-purpose DBI tool.

Acknowledgment

This research was supported by the NSF. We thank our colleague Dr. Hari Sundar from the University of Utah who provided insight and expertise that greatly assisted the research. We also thank the Texas Advanced Computing Center (TACC) and the San Diego Supercomputer Center (SDSC) for the infrastructure they provided for running our experiments.

APPENDIX A

THE FIRST

This is an appendix. Notice that the `\LaTeX` markup for an appendix is, surprisingly, `\chapter`. The `\appendix` command does not produce a heading; instead, it just changes the numbering style from numeric to alphabetic, and it changes the heading prefix from **CHAPTER** to **APPENDIX**.

Blah blah blah blah blah blah blah blah blah blah blah blah blah blah.
Blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah
blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah
blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah
blah blah blah blah blah blah blah blah blah. Blah blah blah blah blah
blah blah blah blah blah blah blah blah.

APPENDIX B

THE SECOND

This is an appendix.

Blah blah blah blah blah blah blah blah blah blah blah blah blah.
Blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah
blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah. Blah blah blah blah
blah blah blah blah blah blah blah blah blah blah blah blah blah.

APPENDIX C

THE THIRD

This is an appendix.

There are several books [12, 19–21, 23–25, 27–30] listed in our bibliography.

We also reference several journal articles [1, 2, 4, 8–10, 13–18, 22, 31, 32] and three famous doctoral theses of later winners [3, 6, 7] of the Nobel Prize in Physics (1922, 1933, and 1921):

Notice that, even though those citations appeared in `\LaTeX \cite{...}` commands with their `BIBTeX` citation labels in reverse alphabetical order, thanks to the `citesort` package, their reference-list numbers have been sorted in numerically ascending order, and then range-reduced.

Mention should also be made of a famous Dutch computer scientist's first publication [5].

Font metrics are an important, albeit low-level, aspect of typesetting. See the *Adobe Systems* manual about that company's procedures [26].

The bibliography at the end of this thesis contains several examples of documents with non-English titles, and their `BIBTeX` entries provide title translations following the practice recommended by the American Mathematical Society and SIAM. Here is a sample entry that shows how to do so:

```
@PhdThesis{Einstein:1905:NBM,  
  author =      "Albert Einstein",  
  title =      "{Eine Neue Bestimmung der Molek{\u}ldimensionen}.  
                ({German}) [{A} new determination of molecular  
                dimensions]",  
  type =      "Inaugural dissertation",  
  school =     "Bern Wyss.",  
  address =    "Bern, Switzerland",  
  year =      "1905",  
  bibdate =    "Fri Dec 17 10:46:57 2004",  
  bibsource =  "http://www.math.utah.edu/pub/tex/bib/einstein.bib",  
  note =      "Published in \cite{Einstein:1906:NBM}."}
```

```

acknowledgement = ack-nhfb,
language =       "German",
advisor =        "Alfred Kleiner (24 April 1849--3 July 1916)",
URL =            "http://en.wikipedia.org/wiki/Alfred_Kleiner",
remark =         "Received August 19, 1905 and published February 8,
                  1906.",
Schilpp-number = "6",
}

```

The `note` field in that entry refers to another bibliography entry that need not have been directly cited in the document text. Such cross-references are common in `BIBTEX` files, especially for journal articles where there may be later comments and corrigenda that should be mentioned. Embedded `\cite{}` commands ensure that those possibly-important other entries are always included in the reference list when the entry is cited. The last bibliography entry [32] in this thesis has a long `note` field that tells more about what some may view as the most important paper in mathematics in the last century.

When entries cite other entries that cite other entries that cite other entries that ..., multiple passes of `LATEX` and `BIBTEX` are needed to ensure consistency. That is another reason why document compilation should be guided by a `Makefile` or a batch script, rather than expecting the user to remember just how many passes are needed.

`BIBTEX` entries are *extensible*, in that arbitrary key/value pairs may be present that are not necessarily recognized by any bibliography style files. The `advisor`, `acknowledgement`, `bibdate`, `bibsource`, `language`, `remark`, and `Schilpp-number` fields are examples, and may be used by other software that processes `BIBTEX` entries, or by humans who read the entries. `DOI` and `URL` fields are currently recognized by only a few styles, but that situation will likely change as publishers demand that such important information be included in reference lists.

In `BIBTEX` `title` fields, braces protect words, such as proper nouns and acronyms, that cannot be downcased if the selected bibliography style would otherwise do so. In German, all nouns are capitalized, and the simple way to ensure their protection is to brace the entire German text in the title, as we did in the entry above.

The world's first significant computer program may have been that written in 1842 by Lady Augusta Ada Lovelace (1815–1852) for the computation of Bernoulli numbers [16, 18]. She was the assistant to Charles Babbage (1791–1871), and they are the world's

REFERENCES

- [1] H. P. BABBAGE, *Babbage: Babbage's analytical engine*, Monthly Notices of the Royal Astronomical Society, 70 (1910), pp. 517–526, 645. Reprinted in [27, §2.3].
- [2] N. H. F. BEEBE AND R. P. C. RODGERS, *<PLOT79>: a comprehensive portable Fortran scientific line graphics system, as applied to biomedical research*, Computers in Biology and Medicine, 19 (1989), pp. 385–402.
- [3] N. H. D. BOHR, *Studier over Metallernes Elektronteori. (Danish) [Studies on the electron theory of metals]*, doktor disputats, Københavns Universitet, København, Danmark, 1911. Afhandling for den filosofiske doktorgrad. [Thesis for the Doctor of Philosophy].
- [4] W. J. CODY, JR., *Analysis of proposals for the floating-point standard*, Computer, 14 (1981), pp. 63–69.
- [5] E. W. DIJKSTRA, *Functionele beschrijving van de ARRA. (Dutch) [Functional description of the ARRA]*, Tech. Rep. 12, Mathematisch Centrum, Amsterdam, The Netherlands, 1953.
- [6] P. A. M. DIRAC, *Quantum Mechanics*, Ph.D. thesis, Cambridge University, Cambridge, UK, June 1926. According to [12, p. 101], this is the first thesis to be submitted anywhere on the subject of quantum mechanics.
- [7] A. EINSTEIN, *Eine Neue Bestimmung der Moleküldimensionen. (German) [A new determination of molecular dimensions]*, inaugural dissertation, Bern Wyss., Bern, Switzerland, 1905. Published in [8].
- [8] ———, *Eine neue Bestimmung der Moleküldimensionen. (German) [A new determination of molecular dimensions]*, Annalen der Physik (1900) (series 4), 324 (1906), pp. 289–306. See corrections [9, 10]. This is a slightly revised version of Einstein's doctoral dissertation [7].
- [9] A. EINSTEIN, *Bemerkung zu meiner Arbeit: Eine Beziehung zwischen dem elastischen Verhalten. (German) [Remark on my paper: "A relationship between the elastic behavior ..."]*, Annalen der Physik (1900) (series 4), 339 (1911), pp. 590–590. See [11].
- [10] A. EINSTEIN, *Berichtigung zu meiner Arbeit: Eine neue Bestimmung der Moleküldimensionen. (German) [Corrections to my work: a new determination of molecular dimensions]*, Annalen der Physik (1900) (series 4), 339 (1911), pp. 591–592. See [8].
- [11] ———, *Eine Beziehung zwischen dem elastischen Verhalten und der spezifischen Wärme bei festen Körpern mit einatomigem Molekül. (German) [A relationship between the elastic behavior and the specific heat of solid bodies with monatomic molecules]*, Annalen der Physik (1900) (series 4), 339 (1911), pp. 170–174, 590. See remarks [9, 10].

- [12] G. FARMELO, *The Strangest Man: The Hidden Life of Paul Dirac, Mystic of the Atom*, Basic Books, New York, NY, USA, 2009.
- [13] H. H. GOLDSTINE AND A. GOLDSTINE, *The Electronic Numerical Integrator and Computer (ENIAC)*, Mathematical Tables and Other Aids to Computation, 2 (1946), pp. 97–110. Reprinted in [27, §7.7].
- [14] P. HALL AND P. PATIL, *Properties of nonparametric estimators of autocovariance for stationary random fields*, Probability Theory and Related Fields, 99 (1994), pp. 399–424.
- [15] J. L. HEILBRON AND T. S. KUHN, *The genesis of the Bohr atom*, Historical Studies in the Physical Sciences, 1 (1969), pp. vi, 211–290.
- [16] V. R. HUSKEY AND H. D. HUSKEY, *Lady Lovelace and Charles Babbage*, Annals of the History of Computing, 2 (1980), pp. 299–329.
- [17] S. C. JOHNSON AND M. E. LESK, *Language development tools*, The Bell System Technical Journal, 57 (1978), pp. 2155–2176.
- [18] E. E. KIM AND B. A. TOOLE, *Ada and the first computer: The collaboration between ada, countess of lovelace, and computer pioneer Charles Babbage resulted in a landmark publication that described how to program the world’s first computer*, Scientific American, 280 (1999), pp. 76–81.
- [19] D. E. KNUTH, *The T_EXbook*, vol. A of Computers and Typesetting, Addison-Wesley, Reading, MA, USA, 1986.
- [20] ———, *The METAFontbook*, vol. C of Computers and Typesetting, Addison-Wesley, Reading, MA, USA, 1986.
- [21] ———, *Digital Typography*, CSLI Publications, Stanford, CA, USA, 1999.
- [22] S. N. LAHIRI, Y. LEE, AND N. CRESSIE, *On asymptotic distribution and asymptotic efficiency of least squares estimators of spatial variogram parameters*, Journal of Statistical Planning and Inference, 103 (2002), pp. 65–85.
- [23] L. LAMPORT, *L^AT_EX—A Document Preparation System—User’s Guide and Reference Manual*, Addison-Wesley, Reading, MA, USA, 1985.
- [24] F. MITTELBACH, M. GOOSSENS, J. BRAAMS, D. CARLISLE, C. ROWLEY, C. DETIG, AND J. SCHROD, *The L^AT_EX Companion*, Tools and Techniques for Computer Typesetting, Addison-Wesley, Reading, MA, USA, second ed., 2004.
- [25] F. W. J. OLVER, D. W. LOZIER, R. F. BOISVERT, AND C. W. CLARK, eds., *NIST Handbook of Mathematical Functions*, Cambridge University Press, Cambridge, UK, 2010.
- [26] POSTSCRIPT DEVELOPER TOOLS & STRATEGIES GROUP, ADOBE SYSTEMS INC., *Adobe font metric files specification — Version 3.0*, Mountain View, CA, USA, Mar. 1990.
- [27] B. RANDELL, ed., *The Origins of Digital Computers: Selected Papers*, Texts and monographs in computer science, Springer-Verlag, Berlin, Germany / Heidelberg, Germany / London, UK / etc., third ed., 1982.

- [28] A. ROBBINS AND N. H. F. BEEBE, *Classic Shell Scripting*, O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472, USA, 2005.
- [29] D. SALOMON, *The Advanced T_EXbook*, Springer-Verlag, Berlin, Germany / Heidelberg, Germany / London, UK / etc., 1995.
- [30] S. SINGH, *Fermat's Enigma: The Epic Quest to Solve the World's Greatest Mathematical Problem*, Walker and Company, 435 Hudson Street, New York, NY 10014, USA, 1997.
- [31] R. TAYLOR AND A. WILES, *Ring-theoretic properties of certain Hecke algebras*, *Annals of Mathematics*, 142 (1995), pp. 553–572. This paper is a companion to [32], providing the remedy for the flaw in Wiles' 1993 proof of Fermat's Last Theorem. See also [30].
- [32] A. WILES, *Modular elliptic curves and Fermat's Last Theorem*, *Annals of Mathematics*, 142 (1995), pp. 443–551. This paper contains the bulk of the author's proof of the Taniyama–Shimura conjecture and Fermat's Last Theorem, carried out at Princeton University. The companion paper [31] contains the solution to the flaw discovered in the proof that Wiles announced on June 23, 1993, in Cambridge, England. See also [30]. In March 2014, now Royal Society Research Professor Sir Andrew John Wiles of Oxford University was awarded the prestigious Abel Prize in Mathematics for this proof — an award that also carries a cash prize of six million Norwegian crowns, or about US\$722,000.