



# Great Experiences with PiP (Proecss-in-Process)

Atsushi Hori

March 3, 2023

Scattered wits take a long time in picking up.

Take nothing on its looks; take everything on evidence.  
Theres no better rule.

I must be taken as I have been made. The success is not mine, the failure is not mine, but the two together make me.

Charles Dickens, *Great Expectations*

# Acknowledgment

PVAS (Partitioned Virtual Address Space), created by Akio Shimada, is the predecessor to PiP. He came up with the concept of employing PIE. When Akio and I first met with Pavan Balaji, he already had a keen understanding of the potential of PVAS. He also provided me with some excellent inspiration for PiP. He suggested PiP as the name. Min Si spent a lot of time assisting me with my PiP papers. Excellent articles were written by Kaiming Ouyang to use PiP to enhance MPICH performance. I was able to spend the majority of my time building PiP thanks to Yutaka Ishikawa. Balazs Gerofi also provided me with some insightful feedback on PiP. Noriyuki Soda created the Github actions for testing PiP and really aided me in the development of PiP.

Finally, I want to thank Fusa Hori, my wife, for letting me dedicate my time to writing this book.

# Preface

## Motivation

I stopped working in April 2022. Since then, I've written this PiP tutorial while also maintaining the PiP library (which includes this document; see <https://github.com/procinproc>).

PiP offers a distinctive and novel execution model. I've been having a lot of problems developing the PiP library because Linux and the current Glibc don't support the new execution mechanism. I intend for the novel execution model that PiP offers to be significant in computer science in the future, and I will do my best to keep the PiP library up to date.

I am not confident in my ability to adapt PiP to the impending Glibc and Linux because of the fact that Glibc and Linux change much more frequently than PiP does due to manpower constraints. So I made the decision to compose this document to share my thoughts.

This is the reason why this can serve as both an internal document and a tutorial for utilizing PiP.

## Expected Readers

In this document, I made every effort to include as many sample programs as I could. They are primarily written in C. A Linux version of the latest PiP library has been tested (CentOS and Redhat, version 7 and 8). Therefore, readers must be conversant with Linux and C programming.

## PiP Versions

**PiP Version 1** This is the very first release of PiP but it is obsolete now.

**PiP Version 2** This is the stable version of PiP.

**PiP Version 3** This is a beta version of PiP that uses User-Level Process (ULP) and Bi-Level Thread (BLT). There will be no description of BLT and ULP in this paper because they are not stable at the time of writing.

## Other Documents

Not all of the PiP library's features are covered in this document. Consult the PiP reference manual (PDF file at <https://github.com/procinproc/PiP/blob/pip-2/doc/latex-inuse/libpip-manpages.pdf>), an HTML document, or man pages (man pages and HTML documents will be installed with PiP library) for this purpose.

## Sample Programs

Running the sample programs on a Docker environment running on Mac OSX allowed for thorough testing of each program's functionality and the production of all output examples.

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>PiP Basics</b>                           | <b>3</b>  |
| 1.1      | PiP Tasks                                   | 3         |
| 1.1.1    | pipcc and pip-exec Commands                 | 3         |
| 1.1.2    | Comparing MPI, OpenMP and PiP               | 4         |
| 1.1.3    | Export and Import                           | 6         |
| 1.2      | Spawning PiP Tasks and Waiting Terminations | 8         |
| 1.2.1    | Spawning PiP tasks                          | 8         |
| 1.2.2    | Waiting for Terminations of PiP tasks       | 13        |
| 1.2.3    | Terminating PiP tasks                       | 15        |
| 1.3      | Timing Synchronization among PiP Tasks      | 16        |
| 1.3.1    | Barrier Synchronization                     | 16        |
| 1.3.2    | Using PThread Synchronization               | 18        |
| 1.3.3    | pthread_barrier                             | 18        |
| 1.3.4    | pthread_mutex                               | 19        |
| 1.4      | PiP Commands                                | 20        |
| 1.4.1    | pip-man                                     | 20        |
| 1.4.2    | pipcc and pipfc                             | 20        |
| 1.4.3    | pip-exec                                    | 21        |
| 1.4.4    | pips  | 22        |
| 1.4.5    | pip-gdb                                     | 23        |
| 1.4.6    | pip-mode and printpipmode                   | 23        |
| 1.4.7    | libpip.so                                   | 23        |
| 1.5      | Summary                                     | 23        |
| 1.6      | Myths on PiP                                | 24        |
| <b>2</b> | <b>PiP Advanced</b>                         | <b>27</b> |
| 2.1      | Rationale                                   | 27        |
| 2.2      | Execution Mode                              | 28        |
| 2.2.1    | Differences Between Two Modes               | 28        |
| 2.2.2    | How to Specify Execution Mode               | 29        |
| 2.3      | Spawning Tasks - Advanced                   | 30        |
| 2.3.1    | Start Function                              | 30        |
| 2.3.2    | Stack Size                                  | 30        |

|          |   |           |
|----------|---|-----------|
| 2.3.3    | CPU Core Binding                                | 31        |
| 2.3.4    | File Descriptors and Spawn Hooks                | 31        |
| 2.4      | Execution Context                               | 33        |
| 2.5      | Debugging Support                               | 35        |
| 2.5.1    | PIP_STOP_ON_START                               | 35        |
| 2.5.2    | PIP_GDB_SIGNALS                                 | 36        |
| 2.5.3    | PIP_SHOW_MAPS                                   | 37        |
| 2.5.4    | PIP_SHOW_PIPS                                   | 37        |
| 2.5.5    | PIP_GDB_PATH and PIP_GDB_COMMAND                | 37        |
| 2.6      | Malloc routines                                 | 37        |
| 2.7      | XPMEM   | 39        |
| <b>3</b> | <b>PiP Internals</b>                            | <b>40</b> |
| 3.1      | PiP Implementation                              | 40        |
| 3.1.1    | Spawning Tasks                                  | 40        |
| 3.1.2    | Calling <code>clone()</code> System Call        | 41        |
| 3.1.3    | Execution Mode in Details                       | 42        |
| 3.1.4    | Name of PiP Tasks                               | 42        |
| 3.2      | Issues related to Linux Kernel, Glibc and Tools | 43        |
| 3.2.1    | Loading a Program                               | 43        |
| 3.2.2    | Glibc   | 44        |
| 3.2.3    | Glibc RPATH Setting                             | 48        |
| 3.2.4    | Linux   | 48        |
| 3.2.5    | Tools   | 51        |
| 3.3      | Remaining Issues                                | 51        |
| 3.3.1    | Retrieving Memory                               | 51        |
| <b>4</b> | <b>PiP Installation</b>                         | <b>52</b> |
| 4.1      | Building from Source Code                       | 52        |
| 4.2      | <b>pip-pip command</b>                          | 53        |
| 4.3      | Using Spack                                     | 54        |

# List of Figures

|     |   |    |
|-----|---|----|
| 1.1 | Differences of OpenMP, MPI and PiP . . . . .  | 6  |
| 2.1 | Cross-Malloc-Free Issue . . . . .             | 38 |
| 2.2 | Cross-Malloc-Free with Freeing List . . . . . | 38 |



# List of Tables

|     |  |    |
|-----|--|----|
| 2.1 | Differences between two modes . . . . .                    | 28 |
| 2.2 | Mode-Agnostic Functions . . . . .                          | 29 |
| 2.3 | Execution Mode Predicates . . . . .                        | 29 |
| 2.4 | Possible Signal Names for <b>PIP_GDB_SIGNALS</b> . . . . . | 37 |
| 3.1 | Command Name Setting (1st char.) . . . . .                 | 42 |
| 3.2 | Command Name Setting (Second char.) . . . . .              | 43 |
| 3.3 | Glibc functions wrapped by PiP library . . . . .           | 45 |
| 3.4 | Compatibility of Tools . . . . .                           | 51 |

# Listings

|      |   |    |
|------|---|----|
| 1.1  | Hello World ( <b>hello.c</b> ) . . . . .                              | 3  |
| 1.2  | Hello World - Compile and Execute . . . . .                           | 3  |
| 1.3  | Hello World having a static variable ( <b>hello-var.c</b> ) . . . . . | 4  |
| 1.4  | Hello World with a static variable - Compile and Execute . . . . .    | 4  |
| 1.5  | Hello World in OpenMP ( <b>hello-var-omp.c</b> ) . . . . .            | 5  |
| 1.6  | Hello World in OpenMP, PiP and MPI - Compile and Execute . . . . .    | 5  |
| 1.7  | Export and Import ( <b>export-import</b> ) . . . . .                  | 7  |
| 1.8  | Execution of Export and Import . . . . .                              | 7  |
| 1.9  | Spawn ( <b>spawn-root</b> ) . . . . .                                 | 9  |
| 1.10 | Spawn ( <b>spawn-task</b> ) . . . . .                                 | 9  |
| 1.11 | Spawn - Execution . . . . .   | 10 |
| 1.12 | Spawn Myself ( <b>spawn-myself</b> ) . . . . .                        | 10 |
| 1.13 | Spawn Myself - Execution . . . . .                                    | 11 |
| 1.14 | Starting from user-defined function ( <b>userfunc</b> ) . . . . .     | 11 |
| 1.15 | Starting from user-defined function - Execution . . . . .             | 12 |
| 1.16 | Starting from main function ( <b>mainfunc</b> ) . . . . .             | 12 |
| 1.17 | Starting from main function - Execution . . . . .                     | 13 |
| 1.18 | Waiting for specified PiP task terminations ( <b>wait</b> ) . . . . . | 13 |
| 1.19 | Waiting for specified PiP task terminations - Execution . . . . .     | 14 |
| 1.20 | Waiting for any PiP task terminations ( <b>waitany</b> ) . . . . .    | 14 |
| 1.21 | Waiting for any PiP task terminations - Execution . . . . .           | 15 |
| 1.22 | PiP Task Termination function ( <b>exit</b> ) . . . . .               | 15 |
| 1.23 | PiP Task Termination - Execution . . . . .                            | 16 |
| 1.24 | Barrier Synchronization ( <b>barrier</b> ) . . . . .                  | 16 |
| 1.25 | Barrier Synchronization - Execution . . . . .                         | 17 |
| 1.26 | Pthread Barrier ( <b>pthread-barrier</b> ) . . . . .                  | 18 |
| 1.27 | Pthread Barrier - Execution . . . . .                                 | 19 |
| 1.28 | Pthread Mutex ( <b>pthread-mutex</b> ) . . . . .                      | 19 |
| 1.29 | Pthread Mutex - Execution . . . . .                                   | 20 |
| 1.30 | <b>pip-check</b> - Execution Example . . . . .                        | 21 |
| 1.31 | <b>pip-exec</b> - Execution Example . . . . .                         | 22 |
| 1.32 | <b>libpip.so</b> - Execution Example . . . . .                        | 23 |
| 2.1  | Before and After Hooks . . . . .                                      | 32 |
| 2.2  | Before and After Hooks - Execution . . . . .                          | 33 |

|     |   |    |
|-----|---|----|
| 2.3 | Calling a Function of Another Task . . . . .        | 33 |
| 2.4 | Function Call of Another Task - Execution . . . . . | 34 |
| 2.5 | Stop-on-start Script Example . . . . .              | 35 |
| 2.6 | Stop-on-start Script Example - Execution . . . . .  | 36 |
| 3.1 | .interp Section of the ps command . . . . .         | 43 |
| 3.2 | Constructors and Destructors . . . . .              | 46 |
| 3.3 | Constructors and Destructors - Execution . . . . .  | 47 |
| 3.4 | A Memory Map Example . . . . .                      | 48 |
| 4.1 | Building from Source Code . . . . .                 | 53 |
| 4.2 | PiP-pip installation example . . . . .              | 53 |
| 4.3 | Spack installation example . . . . .                | 54 |

# Introduction

This document explains the PiP (short for Process-in-Process) library. This library with a little odd name offers a relatively new execution paradigm that combines the finest aspects of multi-process and multi-threaded execution models.

Multiple CPU cores in a CPU die or socket are becoming very common, and a parallel execution environment is essential for maximizing the capability of the many-core architecture. Despite accessing the same physical memory device, a process in the multi-process model, where numerous processes run concurrently on a node, cannot directly access data held by the other processes. Processes typically communicate in order to exchange information. In my opinion, all forms of data copying, whether carried out by hardware or software, are a part of communication. Data copying is energy-, memory-, and time-intensive and should be avoided if possible. What if processes could access other people's data without communicating with them first? In the multi-thread paradigm, static variables are shared across threads, and if threads attempt to update their values, race problems must be prevented. What if static variables were private for each thread?

This is what drives me to develop PiP. As the name *Process-in-Process* might imply, a process can spawn additional processes inside of its own address space. This sounds like the multi-thread execution model, but the name of the *process* in PiP indicates that, unlike the multi-thread approach, each new process has its own static variable set. As a result, while keeping the privatized static variables, the created processes that share the same address space can access data that belongs to the others. On this basis, data copying and communication can be prevented.

In essence, a multi-process model shares nothing, a multi-thread model shares everything, and PiP's execution paradigm allows for everything to be shareable.

Regarding some of my forebears, different implementations exist that offer this particular execution approach. PiP, however, stands out since it is implemented entirely at the user-level and doesn't require a new OS kernel, a patched OS kernel, or new language processing systems.

High performance computing (HPC) has been my area of focus, and I know very little about the other disciplines. HPC applications are all I

can think of. But I think that PiP's usability can be extended to other industries.

# Chapter 1

## PiP Basics

For those who are unfamiliar, let me first go through the fundamentals of PiP: 1) how to execute a PiP program, 2) how to write a PiP program, and 3) how to use PiP commands. This chapter's explanations don't get into specifics. Consult the other publications (man pages and PDF) or Chapter 2 for further information.

### 1.1 PiP Tasks

This section will explain how PiP tasks are created simply and how they operate differently from processes (made using MPI) and threads (created using OpenMP).

#### 1.1.1 pipcc and pip-exec Commands

The first example is the well-known C program “hello world” listed below;

Listing 1.1: Hello World (`hello.c`)

```
#include <stdio.h>
int main() {
    printf( "Hello World\n" );
    return 0;
}
```

As you can see, this program is a perfect match for a standard C program. If the `pipcc` command was used to compile the program, it can be run as a standard C program or as a PiP task by using the `pip-exec` command.

Listing 1.2: Hello World - Compile and Execute

```
$ pipcc --silent hello.c -o hello
$ ./hello
Hello World
$ pip-exec ./hello
```

```
Hello World
$
```

A true C compiler can be called with the proper options, such as `-I`, `-L`, and others, using the `pipcc` command, which is written as a shell script. You will see the options available when the `pipcc` script calls the backend C/C++ compiler if the `silent` option is omitted.

In this case, `pip-exec` is being used to run an executable file as PiP tasks rather than a standard Linux process. The hello program does not operate differently in the process and PiP task in this case. In the following section, we'll talk about this issue.

### 1.1.2 Comparing MPI, OpenMP and PiP

We slightly alter the “hello world” software as follows to clarify the distinction between the Linux process and PiP task;

Listing 1.3: Hello World having a static variable (`hello-var.c`)

```
#include <stdio.h>
int x;
int main() {
    printf( "Hello World (&x:%p)\n", &x );
    return 0;
}
```

Now, the address of the static variable `x` in the “Hello World” program is printed out along with the message “Hello World.” The number of PiP jobs to be created and run concurrently can be set using the `pip-exec` command option. In the execution example that follows, the number three (3) is supplied. The same `a.out` execution using MPI's output is also included. It should be noted that the “Hello World” program runs in parallel with `pip-exec` and `mpiexec`.

Listing 1.4: Hello World with a static variable - Compile and Execute

```
$ pipcc --silent hello-var.c -o hello-var
$ ./hello-var
Hello World (&x:0x555555601030)
$ pip-exec -n 3 ./hello-var
Hello World (&x:0x7ffff67d9030)
Hello World (&x:0x7ffff48db030)
Hello World (&x:0x7ffffe92c030)
$ mpiexec -n 3 ./hello-var
Hello World (&x:0x555555601030)
Hello World (&x:0x555555601030)
Hello World (&x:0x555555601030)
$
```

The variable `x` is found at the address `0x5555556010301` according to the first execution of `a.out`. With MPI execution, this circumstance is same<sup>1</sup>. The variable `x` is however executed at various places for each PiP activity. This is due to MPI jobs not sharing the same address space as PiP tasks.

Readers who are curious in the distinction between PiP and OpenMP may observe that threads also share the same address space. The “Hello World” program with a static variable is shown in the example below written in OpenMP.

Listing 1.5: Hello World in OpenMP (`hello-var-omp.c`)

---

```
#include <stdio.h>
int x;
int main() {
    #pragma omp parallel
    printf( "Hello World (&x:%p)\n", &x );
    return 0;
}
```

---

The output of program 1.5’s execution is displayed below. The addresses of variable `x` in this case are the same for MPI and OpenMP executions. The variable’s addresses with PiP execution, however, are different pairings.

Listing 1.6: Hello World in OpenMP, PiP and MPI - Compile and Execute

```
$ pipcc --silent -fopenmp hello-var-omp.c -o hello-var-omp
$ export OMP_NUM_THREADS=2
$ ./hello-var-omp
Hello World (&x:0x555555601038)
Hello World (&x:0x555555601038)
$ pip-exec -n 2 ./hello-var-omp
Hello World (&x:0x7ffff67d9038)
Hello World (&x:0x7ffff67d9038)
Hello World (&x:0x7ffffefde3038)
Hello World (&x:0x7ffffefde3038)
$ mpiexec -n 2 ./hello-var-omp
Hello World (&x:0x555555601038)
Hello World (&x:0x555555601038)
Hello World (&x:0x555555601038)
Hello World (&x:0x555555601038)
$
```

These variations are explained in Figure 1.1. The variable `x` is shared by all of the OpenMP threads, and they all use the same address space. Each MPI process in an MPI environment has its own address space, and two (2) threads can execute in each address space while sharing a variable in an MPI process. However, each PiP task has its own variables, thus threads

---

<sup>1</sup>For simplicity, we disabled ASLR (Address Space Layout Randomization) in this example.



0 and 1 only share variables within the same PiP task; they do not share variables inside any other PiP tasks. In PiP, all PiP tasks share the same address space.

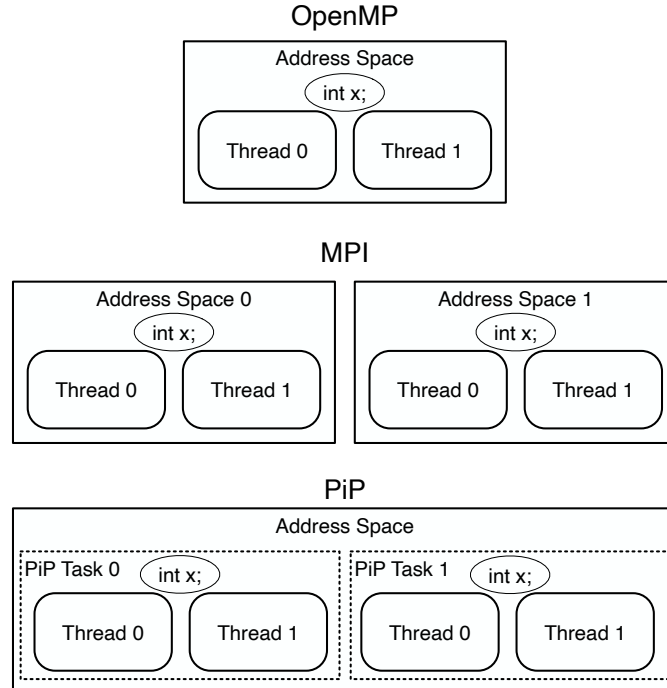


Figure 1.1: Differences of OpenMP, MPI and PiP

Static variables are associated to an address space in the traditional process model and thread model. Because of this, each process has its own static variables, which are shared by all threads using the same address space. Each PiP task is guaranteed to have its own static variable set under the PiP execution model, decoupling from the address space while maintaining address space sharing. **variable privatization** is what this is.

It is simple to share information among PiP tasks while maintaining the independence of each PiP task's execution thanks to the nature of PiP, which includes privatized variables and sharing an address space. The "Hello World" program has so far been proved to be able to execute as PiP tasks concurrently, although this program is quite basic and there is no information exchange between PiP tasks. We'll demonstrate how information can be shared between PiP processes in the section after this one.

### 1.1.3 Export and Import

If the address of the information to be exchanged is known, sharing an address space allows data owned by a PiP task to be accessed. The address

of the shared data can be broadcast by a PiP task, and the other PiP task(s) can obtain the published address.

To start, each PiP task has a **PIPID** that helps it stand out from the rest. By giving the **PIPID** of the exporting PiP task, other PiP tasks that share the same address space can import the exported address.

Listing 1.7: Export and Import (**export-import**)

---

```
#include <pip/pip.h>
#include <stdlib.h>
int x;
int main( int argc, char **argv ) {
    int pipid, *xp;
    pip_get_pipid( &pipid );
    if( pipid == 0 ) {
        x = strtol( argv[1], NULL, 10 );
        pip_named_export( &x, "export" );
    } else {
        pip_named_import( 0, (void**) &xp, "export" );
        printf( "%d: %d\n", pipid, *xp );
    }
    return 0;
}
```

---

In this program, after adjusting the value of `argv[1]`, a PiP task with a PIPID of zero (zero) exports the address of the variable `x` by using **pip\_named\_export()**. By using the **pip\_named\_import()** function, the remaining PiP tasks import the address that PiP task 0 exported. An outcome of this program's execution is shown below. The other PiP tasks can view the value that PiP task 0 exported, as demonstrated.

Listing 1.8: Execution of Export and Import

```
$ pip-exec -n 4 ./export-import 1234
1: 1234
2: 1234
3: 1234
$ pip-exec -n 4 ./export-import 18526
1: 18526
2: 18526
3: 18526
$
```

The address with the specified name is published by the **pip\_named\_export()** function. The **pip\_named\_import()** function blocking-waits for the specified PiP job by **PIPID** to reach the defined address. To avoid a race condition, it is not permitted to export an address with the same name more than once for the purpose of updating the address.

The PiP library's functions almost always return an integer value as an error code. A return code of zero (zero) denotes success. This error code

is identical to those that Linux defines. Due to simplicity and clarity, the returned code is not tested in the examples presented thus far and moving forward.

It is forbidden in MPI to access the data that is held by other processes running on the same node<sup>2</sup>. In MPI, communication is the sole permitted method. Communication fundamentally entails copying data in some way (done by software or hardware). Data copying consumes memory, power, and time.

## 1.2 Spawning PiP Tasks and Waiting Terminations

The `pip-exec` command starts PiP tasks. (PiP) root process is the name of the process that creates PiP tasks. The `pip-exec` process is a root process. The root process's address space is used to map and execute PiP tasks that it has spawned. This chapter will describe how to spawn PiP tasks.

### 1.2.1 Spawning PiP tasks

#### Spawning a program as PiP tasks

Listing 1.9 is an example of a PiP root program. It spawns  $N$  PiP tasks, where  $N$  is specified by the first parameter of the program. The `pip_init()` function must be called to initialize the PiP library before calling any other PiP functions, although there are some exceptions to this. The first argument is output returning `PIPID` of the calling task. In this case, `PIP_PIPID_ROOT` is returned, since the function is called by the root. The second input argument is to specify the maximum number of spawning PiP tasks. The other arguments will be explained in Chapter 2.

The `pip_spawn()` function is called after then. The first and second arguments are the same with the Linux's `execve()` function; the first is to specify the executable file to be executed and the second argument is to specify the parameters executing the program. The third is to specify environment variables. When it is `NULL`, then value of the Glibc global variable `environ` is taken. The fourth argument is to specify the CPU core number to bind the spawned PiP task and which CPU core. In this example, the value of `PIP_CPUCORE_ASIS` means that the (CPU) core-bind should be the same with the one when calling `pip_spawn()`. The fifth is an input and output argument, and you can specify a `PIPID` or set to the value

---

<sup>2</sup>Strictly speaking, some MPI implementations based on the thread model may allow this. Major MPI implementation, such as MPICH, Open MPI, and many other MPI implementations provided by vendors are based on the process model, and there is no way to access data owned by the other MPI process.

of **PIP\_PIPID\_ANY** so that PiP library can choose any. After calling **pip\_spawn()**, the argument returns the actual **PIPID**.

The **pip\_wait()** is to wait the termination of the spawned task. Its first argument is to specify the **PIPID** of the terminating task. The second parameter, although **NULL** is set in this example, is the same with the Linux's **wait()** function, returning the terminating status of a task.

The **pip\_fin()** function works as the opposite of **pip\_init()**, finalizing PiP library and freeing allocated resources. After calling **pip\_fin()**, most PiP library functions return an error code (**EPERM**).

Listing 1.9: Spawn (spawn-root)

---

```
#include <pip/pip.h>
#include <stdlib.h>
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, NULL, 0 );
    for( i=0; i<ntasks; i++ ) {
        pipid_task = PIP_PIPID_ANY;
        pip_spawn( argv[2], &argv[2], NULL,
                  PIP_CPUCORE_ASIS, &pipid_task,
                  NULL, NULL, NULL );
        pip_wait( pipid_task, NULL );
        printf( "PiP task (PIPID:%d) done\n",
               pipid_task );
    }
    pip_fin();
    return 0;
}
```

---

Listing 1.10 is very similar to the “Hello World” program in the previous section. The major difference here is calling the **pip\_init()** function. The **pip\_init()** may look strange because this function behaves differently depending on if it is called from a PiP root or PiP task. Unlike root, this function call is optional in the PiP task program. By calling this, you can get **PIPID** and the number of maximum PiP tasks which are specified by the root. Listing 1.11 shows an example of the execution of Listing 1.9 and 1.10.

Listing 1.10: Spawn (spawn-task)

---

```
#include <pip/pip.h>
int main( int argc, char **argv ) {
    int pipid, ntasks;
    pip_init( &pipid, &ntasks, NULL, 0 );
    printf( "\"%s\" from PIPID:%d/%d\n",
           argv[1], pipid, ntasks );
    pip_fin();
}
```

---

```

    return 0;
}

```

---

Listing 1.11: Spawn - Execution

```

$ ./spawn-root 4 ./spawn-task "What's up?"
"What's up?" from PIPID:0/4
"What's up?" from PIPID:1/4
"What's up?" from PIPID:2/4
"What's up?" from PIPID:3/4
PiP task (PIPID:0) done
PiP task (PIPID:1) done
PiP task (PIPID:2) done
PiP task (PIPID:3) done
$

```

## Spawning myself

A program can be either the PiP task or the PiP root. Combining the programs from Listings 1.12 and 1.10 is demonstrated in Listing 1.12 as an example. We hope you can comprehend `pip_init()`'s peculiar behavior. The PiP root process functions similarly to a PiP task. It has a unique **PIPID** called the **PIP\_PIPID\_ROOT**. Listing 1.13 provides an illustration of this execution.

Listing 1.12: Spawn Myself (spawn-myself)

```

#include <pip/pip.h>
#include <stdlib.h>
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, NULL, 0 );
    if( pipid == PIP_PIPID_ROOT ) {
        /* PiP root */
        for( i=0; i<ntasks; i++ ) {
            pipid_task = i;
            pip_spawn( argv[0], argv, NULL,
                      PIP_CPUCORE_ASIS, &pipid_task,
                      NULL, NULL, NULL );
            pip_wait( pipid_task, NULL );
            printf( "PiP task (PIPID:%d) done\n",
                   pipid_task );
        }
    } else {
        /* PiP task */
        printf( "\"%s\" from PIPID:%d/%d\n",
               argv[2], pipid, ntasks );
    }
}

```

```

    pip_fin();
    return 0;
}

```

---

Listing 1.13: Spawn Myself - Execution

```

$ ./spawn-myself 4 "Learning PiP."
"Learning PiP." from PIPID:0/4
"Learning PiP." from PIPID:1/4
"Learning PiP." from PIPID:2/4
"Learning PiP." from PIPID:3/4
PiP task (PIPID:0) done
PiP task (PIPID:1) done
PiP task (PIPID:2) done
PiP task (PIPID:3) done
$

```

### Starting from other than main

In the preceding instances, PiP tasks start from the `main()` function. PiP enables tasks to launch user-defined functions other than the `main()`. Use the `pip_task_spawn()` function in this situation rather than executing `pip_spawn()`.

Listing 1.14: Starting from user-defined function (`userfunc`)

```

#include <pip/pip.h>
#include <stdlib.h>
int user_func( void *arg ) {
    char *msg = (char*) arg;
    int pipid, ntasks;
    pip_get_pipid( &pipid );
    pip_get_ntasks( &ntasks );
    printf( "USER-FUNC: \"%s\" from PIPID:%d/%d\n",
           msg, pipid, ntasks );
    return 0;
}

int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, NULL, 0 );
    if( pipid == PIP_PIPID_ROOT ) {
        pip_spawn_program_t prog;
        pip_spawn_from_func( &prog,
                             argv[0],      /* exec file */
                             "user_func",  /* func name */
                             (void*) argv[2], /* arg */
                             NULL,         /* environ */

```

```

                                NULL );/* explained later */
for( i=0; i<ntasks; i++ ) {
    pipid_task = i;
    pip_task_spawn( &prog, PIP_CPUCORE_ASIS, 0,
                    &pipid_task, NULL );
    pip_wait( pipid_task, NULL );
}
} else {
    /* NEVER REACH HERE */
    printf( "MAIN: \"%s\" from PIPID:%d/%d\n",
            argv[2], pipid, ntasks );
}
pip_fin();
return 0;
}

```

The software for this example is listed in Listing 1.14. The `pip_spawn_program_t` structure is defined to minimize the number of arguments needed to spawn a PiP process. All necessary data for starting a program, such as the function name and path to the executable file, are stored in this structure. The `pip_spawn_from_func()` function is also defined to set this information in order to mask the structure's specifics. The user-defined function must take a single argument (`void*`) and return an integer value that is identical to the `main()` function's return value.

Listing 1.15: Starting from user-defined function - Execution

```

$ ./userfunc 4 "Calling user_func"
USER-FUNC: "Calling user_func" from PIPID:0/4
USER-FUNC: "Calling user_func" from PIPID:1/4
USER-FUNC: "Calling user_func" from PIPID:2/4
USER-FUNC: "Calling user_func" from PIPID:3/4
$

```

The `pip_spawn()` was initially released (from version 1). After that, I found users could launch PiP tasks other than the `main()`, and the `pip_task_spawn()` function had been introduced (from version 2 or later). When beginning from the `main()` function, the `pip_spawn_program_t` structure must be set by calling the `pip_spawn_from_main()` function. Listing 1.16 is a program that updates Listing 1.12 to use the `pip_task_spawn()` and `pip_spawn_from_main()` functions.

Listing 1.16: Starting from main function (`mainfunc`)

```

#include <pip/pip.h>
#include <stdlib.h>
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, NULL, 0 );

```

```

if( pipid == PIP_PIPID_ROOT ) {
    pip_spawn_program_t prog;
    pip_spawn_from_main( &prog,
                        argv[0],      /* exec file */
                        argv,         /* argv */
                        NULL,         /* environ */
                        NULL ); /* explained later */
    for( i=0; i<ntasks; i++ ) {
        pipid_task = i;
        pip_task_spawn( &prog, PIP_CPUCORE_ASIS, 0,
                        &pipid_task, NULL );
        pip_wait( pipid_task, NULL );
    }
} else {
    printf( "MAIN: \"%s\" from PIPID:%d/%d\n",
           argv[2], pipid, ntasks );
}
pip_fin();
return 0;
}

```

Listing 1.17: Starting from main function - Execution

```

$ ./mainfunc 4 "Calling main"
MAIN: "Calling main" from PIPID:0/4
MAIN: "Calling main" from PIPID:1/4
MAIN: "Calling main" from PIPID:2/4
MAIN: "Calling main" from PIPID:3/4
$

```

## 1.2.2 Waiting for Terminations of PiP tasks

As readers may have already noted, the `pip_wait()` function is used to watch for PiP tasks that have been started to finish. The `pip_wait()` function functions similarly to the `wait()` function in Linux. Linux's `wait()` function frequently cooperates with PiP jobs, however there is one instance where it does not. Therefore, it is advised that users use the `pip_wait()` function. Similar to the `wait()` function in Linux, the argument of `pip_wait()` is a pointer to an integer variable. The Linux `WIFEXITED`, `WIFSIGNALED`, `WEXITSTATUS`, `WIFSIGNALED`, and `WTERMSIG` macros can be used to inspect the returned integer.

Listing 1.18: Waiting for specified PiP task terminations (`wait`)

```

#include <pip/pip.h>
#include <stdlib.h>
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i, exitval = 0;

```



```

ntasks = strtol( argv[1], NULL, 10 );
pip_init( &pipid, &ntasks, NULL, 0 );
if( pipid == PIP_PIPID_ROOT ) {
    for( i=0; i<ntasks; i++ ) {
        int status;
        pipid_task = i;
        pip_spawn( argv[0], argv, NULL,
                    PIP_CPUCORE_ASIS, &pipid_task,
                    NULL, NULL, NULL );
        pip_wait( pipid_task, &status );
        printf( "PiP task (PIPID:%d) done: %d\n",
                pipid_task, WEXITSTATUS(status) );
    }
} else {
    exitval = pipid;
}
pip_fin();
return exitval;
}

```

---

Listing 1.19: Waiting for specified PiP task terminations - Execution

```

$ ./wait 4
PiP task (PIPID:0) done: 0
PiP task (PIPID:1) done: 1
PiP task (PIPID:2) done: 2
PiP task (PIPID:3) done: 3
$

```

**pip\_wait()** waits for the PiP task termination specified by **PIPID**. **pip\_wait\_any()** function can wait for any PiP tasks and **PIPID** and exit status are returned when terminated (See Listing 1.20 and 1.21). **pip\_trywait()** and **pip\_trywait\_any()** are the non-blocking versions of **pip\_wait()** and **pip\_wait\_any()**, respectively.

Listing 1.20: Waiting for any PiP task terminations (**waitany**)

```

#include <pip/pip.h>
#include <stdlib.h>
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i, exitval = 0;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, NULL, 0 );
    if( pipid == PIP_PIPID_ROOT ) {
        int status;
        for( i=0; i<ntasks; i++ ) {
            pipid_task = i;
            pip_spawn( argv[0], argv, NULL,
                        PIP_CPUCORE_ASIS, &pipid_task,
                        NULL, NULL, NULL );

```

```

    }
    for( i=0; i<ntasks; i++ ) {
        pip_wait_any( &pipid_task, &status );
        printf( "PiP task (PIPID:%d) done: %d\n",
                pipid_task, WEXITSTATUS(status) );
    }
} else {
    exitval = pipid;
}
pip_fin();
return exitval;
}

```

Listing 1.21: Waiting for any PiP task terminations - Execution

```

$ ./waitany 4
PiP task (PIPID:0) done: 0
PiP task (PIPID:1) done: 1
PiP task (PIPID:2) done: 2
PiP task (PIPID:3) done: 3
$

```

### 1.2.3 Terminating PiP tasks

By using the **pip\_exit()** function, PiP tasks and root can stop running. Comparable to Linux's **exit()** function is this function. As stated above, it is advised to use **pip\_exit()** rather than **exit()** since, while the Linux **exit()** function typically works, there are some instances where it does not. Listing 1.22 and 1.23 provide a working example of **pip\_exit()**.

Listing 1.22: PiP Task Termination function (**exit**)

```

#include <pip/pip.h>
#include <stdlib.h>
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i, exitval = 0;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, NULL, 0 );
    if( pipid == PIP_PIPID_ROOT ) {
        for( i=0; i<ntasks; i++ ) {
            int status;
            pipid_task = i;
            pip_spawn( argv[0], argv, NULL,
                        PIP_CPUCORE_ASIS, &pipid_task,
                        NULL, NULL, NULL );
            pip_wait( pipid_task, &status );
            printf( "PiP task (PIPID:%d) done: %d\n",
                    pipid_task, WEXITSTATUS(status) );
        }
    }
}

```

```

    pip_exit( 100 );
    /* NEVER REACH HERE */
} else {
    exitval = pipid * 10;
    pip_exit( exitval );
    /* NEVER REACH HERE */
}
}

```

---

Listing 1.23: PiP Task Termination - Execution

```

$ ./exit 4; echo $?
PiP task (PIPID:0) done: 0
PiP task (PIPID:1) done: 10
PiP task (PIPID:2) done: 20
PiP task (PIPID:3) done: 30
100
$

```

## 1.3 Timing Synchronization among PiP Tasks

This section will go over how PiP tasks' timing synchronization works.

### 1.3.1 Barrier Synchronization

The PiP library only supports barrier synchronization as a synchronization mechanism at this time. Barrier synchronization in PiP uses an API that was taken from the PThread library. PiP has three functions that are equivalent to `pthread_barrier_init()`, `pthread_barrier_wait()`, and `pthread_barrier_destroy()`, respectively: **pip\_barrier\_init()**, **pip\_barrier\_wait()**, and **pip\_barrier\_fin()**.

Listing 1.24: Barrier Synchronization (barrier)

```

#include <pip/pip.h>
#include <stdlib.h>
pip_barrier_t barr;
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i;
    pip_barrier_t *barrp = &barr;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, (void**) &barrp, 0 );
    if( pipid == PIP_PIPID_ROOT ) {
        pip_barrier_init( barrp, ntasks );
        for( i=0; i<ntasks; i++ ) {
            pipid_task = i;
            pip_spawn( argv[0], argv, NULL,
                       PIP_CPUCORE_ASIS, &pipid_task,

```

```

        NULL, NULL, NULL );
    }
    for( i=0; i<ntasks; i++ ) {
        pip_wait_any( &pipid_task, NULL );
        printf( "PiP task (PIPID:%d) done\n", pipid_task );
    }
    pip_barrier_fin( barrp );
} else {
    if( argv[2] == NULL ) {
        pip_barrier_wait( barrp );
    }
    printf( "PIPID:%d %f [S]\n", pipid, pip_gettime() );
}
return 0;
}

```

Listing 1.24 gives the third argument to the **pip\_init()** function a non-NULL value this time. This is an additional method of exporting a pointer to spawned PiP jobs from the root. Children in this example are given the address of the **pip\_barrier\_t** static variable so they can synchronize by using **pip\_barrier\_wait()**.

To make the impact of the barrier synchronization more clear, the synchronization only occurs when the second parameter of the program execution is left blank; in this case, PiP tasks display the **pip\_gettime()** return values. The **pip\_gettime()** function returns the current value of the **gettimeofday()** function in double format with a second unit.

Listing 1.25 displays an example of this application being run. The barrier synchronization is absent in the initial run, and the values returned by **gettimeofday()** exhibit significant variation. A lesser difference can be detected in the second run, where the barrier synchronization occurs.

Listing 1.25: Barrier Synchronization - Execution

```

$ ./barrier 4 NOBARRIER
PIPID:0 1661152530.820478 [S]
PIPID:1 1661152530.847696 [S]
PIPID:2 1661152530.878638 [S]
PIPID:3 1661152530.902075 [S]
PiP task (PIPID:0) done
PiP task (PIPID:1) done
PiP task (PIPID:2) done
PiP task (PIPID:3) done
$ ./barrier 4
PIPID:3 1661152531.178360 [S]
PIPID:0 1661152531.178426 [S]
PIPID:1 1661152531.178457 [S]
PIPID:2 1661152531.179018 [S]
PiP task (PIPID:3) done
PiP task (PIPID:0) done

```

```
PiP task (PIPID:1) done
PiP task (PIPID:2) done
$
```

### 1.3.2 Using PThread Synchronization

The PThread library's synchronization features for PiP tasks are likewise available to users. This is due to the fact that PiP tasks and threads both use the same address space.

### 1.3.3 pthread\_barrier

The PThread's barrier functions can also be used to provide the same barrier synchronization. Listing 1.26 shows the program's simple substitution of the PThread's barrier functions for the PiP's barrier functions.

Listing 1.26: Pthread Barrier (pthread-barrier)

```
#include <pip/pip.h>
#include <stdlib.h>
#include <pthread.h>
pthread_barrier_t barr;
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i;
    pthread_barrier_t *barrp = &barr;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, (void**) &barrp, 0 );
    if( pipid == PIP_PIPID_ROOT ) {
        pthread_barrier_init( barrp, NULL, ntasks );
        for( i=0; i<ntasks; i++ ) {
            pipid_task = i;
            pip_spawn( argv[0], argv, NULL,
                      PIP_CPUCORE_ASIS, &pipid_task,
                      NULL, NULL, NULL );
        }
        for( i=0; i<ntasks; i++ ) {
            pip_wait_any( &pipid_task, NULL );
            printf( "PiP task (PIPID:%d) done\n", pipid_task );
        }
        pthread_barrier_destroy( barrp );
    } else {
        if( argv[2] == NULL ) {
            pthread_barrier_wait( barrp );
        }
        printf( "PIPID:%d %f [S]\n", pipid, pip_gettime() );
    }
    return 0;
}
```

Listing 1.27: Pthread Barrier - Execution

```
$ ./pthread-barrier 4 NOBARRIER
PIPID:0 1661152533.381032 [S]
PIPID:1 1661152533.410988 [S]
PIPID:2 1661152533.451032 [S]
PIPID:3 1661152533.566188 [S]
PiP task (PIPID:0) done
PiP task (PIPID:1) done
PiP task (PIPID:2) done
PiP task (PIPID:3) done
$ ./pthread-barrier 4
PIPID:3 1661152534.127914 [S]
PIPID:1 1661152534.127852 [S]
PIPID:0 1661152534.127925 [S]
PIPID:2 1661152534.128041 [S]
PiP task (PIPID:1) done
PiP task (PIPID:3) done
PiP task (PIPID:0) done
PiP task (PIPID:2) done
$
```

### 1.3.4 pthread\_mutex

Similarly, `pthread_mutex()` also works with PiP.

Listing 1.28: Pthread Mutex (`pthread-mutex`)

```
#include <pip/pip.h>
#include <pthread.h>
#include <stdlib.h>
#include <unistd.h>
#define NITERS (1000)
typedef struct sync_tasks {
    pthread_barrier_t    barr;
    pthread_mutex_t      mutex;
    int                  count;
} sync_t;
sync_t sync_tasks;
int lock;
void increment( sync_t *syncp ) {
    int tmp;
    if( lock ) pthread_mutex_lock( &syncp->mutex );
    tmp = syncp->count;
    usleep( 10 );
    syncp->count = tmp + 1;
    if( lock ) pthread_mutex_unlock( &syncp->mutex );
}
int main( int argc, char **argv ) {
    int pipid, ntasks, i;
```

```

sync_t *syncp;
pip_get_pipid( &pipid );
pip_get_ntasks( &ntasks );
lock = ( argc == 1 );
if( pipid == 0 ) {
    syncp = &sync_tasks;
    pthread_barrier_init( &syncp->barr, NULL, ntasks );
    pthread_mutex_init( &syncp->mutex, NULL );
    syncp->count = 0;
    pip_named_export( syncp, "sync" );
} else {
    pip_named_import( 0, (void**) &syncp, "sync" );
}
pthread_barrier_wait( &syncp->barr );
for( i=0; i<NITERS; i++ ) increment( syncp );
pthread_barrier_wait( &syncp->barr );
if( pipid == 0 ) {
    printf( "count=%d (%d*%d)\n", syncp->count,
           ntasks, NITERS );
}
return 0;
}

```

---

Listing 1.29: Pthread Mutex - Execution

```

$ pip-exec -n 10 ./pthread-mutex
count=10000 (10*1000)
$ pip-exec -n 10 ./pthread-mutex NOLOCK
count=992 (10*1000)
$

```

## 1.4 PiP Commands

The PiP package's PiP commands are covered in this section. Some of these have already been demonstrated but just briefly discussed. The specifics of the PiP commands will be described in this section.

### 1.4.1 pip-man

The PiP man pages are displayed by this command. Although using this command does not require users to be concerned with the `MAN_PATH`, it does execute the Linux `man` command with the `MAN_PATH` option to the PiP man pages (if installed successfully).

### 1.4.2 pipcc and pipfc

**pipcc** is the compiler script for compiling PiP applications for C and C++, and **pipfc** is for Fortran, as was already mentioned in Section 1.1.1. The **--which** option will display the real back-end compiler's pass. For **pipcc** or **pipfc**, users can also choose the back-end compiler by changing the environment variable **CC** or **FC**.

By default, **pipcc** and **pipfc** compile programs to create code that can execute as either a PiP task or the root process. Users can specify the **--piproot** or **--piptask** options to run an application that only supports PiP tasks or the PiP root. Any PiP software that has been compiled as PiP tasks can, in fact, run as a PiP root. The **--piptask** option is therefore the same as the **--pipboth** option (to be both root and task).

The **--cflags** option displays the actual compilation options that should be provided to the back-end compiler, and the **--lflags** option displays the linker options. The actual compiling and/or linking processes are disabled by the **--cflags** or **--lflags** options.

Listing 1.30: **pip-check** - Execution Example

```
$ pip-check /usr/bin/ps
/usr/bin/ps : not a PiP program
$ pip-check /usr/bin/ls
/usr/bin/ls : not an ELF file
$ cat /usr/bin/ls
#!/usr/bin/coreutils --coreutils-prog-shebang=ls
$ pipcc --silent pip.c -o pip
$ pip-check ./pip
./pip : Root&Task
$ pipcc --silent --piptask pip.c -o pip-task
$ pip-check ./pip-task
./pip-task : Root&Task
$ pipcc --silent --piproot pip.c -o pip-root
$ pip-check ./pip-root
./pip-root : Root
$
```

The **pip-check** application can be used to determine whether a file can be executed as a PiP program. A program may not actually run as a PiP program even if it explicitly states that it will. Linux commands cannot be run as PiP tasks, even if they are created as PIE programs. Any shell script (*shebang*, the file with the extension “#!”) follows the same rules. Listing 1.30 shows an example of a shell script that uses the **ls** command.

### 1.4.3 pip-exec

The **pip-exec** command is used in the examples thus far to execute PiP tasks that are derivations of a single program. Though all PiP tasks created from those programs share the same address space, **pip-exec** has the ability



to invoke several programs. Programs are divided in this way by colons (:) (Listing 1.31).

Listing 1.31: `pip-exec` - Execution Example

```
$ cat prog.c
#include <pip/pip.h>
int main( int argc, char **argv ) {
    int pipid;
    pip_get_pipid( &pipid );
    printf( "This is %s [%d]\n", argv[0], pipid );
    return 0;
}
$ pipcc --silent prog.c -o a.out
$ cp a.out b.out
$ cp a.out c.out
$ pip-exec -n 2 ./a.out : -n 3 ./b.out : -n 1 ./c.out
This is ./a.out [0]
This is ./a.out [1]
This is ./b.out [2]
This is ./b.out [3]
This is ./b.out [4]
This is ./c.out [5]
$
```

#### 1.4.4 pips

Similar to how the Linux `ps` command works, the `pips` command outputs a list of all presently active PiP tasks and PiP roots. Here is an example where three (3) `pip-exec` are launched, with each one executing one of the PiP tasks (a, b, or c).

```
$ pips
PID    TID    TT      TIME      PIP  COMMAND
18741  18741  pts/0   00:00:00  RT   pip-exec
18742  18742  pts/0   00:00:00  RG   pip-exec
18743  18743  pts/0   00:00:00  RL   pip-exec
18741  18744  pts/0   00:00:00  OT   a
18745  18745  pts/0   00:00:00  OG   b
18746  18746  pts/0   00:00:00  OL   c
18747  18747  pts/0   00:00:00  1L   c
18741  18748  pts/0   00:00:00  1T   a
18749  18749  pts/0   00:00:00  1G   b
18741  18750  pts/0   00:00:00  2T   a
18751  18751  pts/0   00:00:00  2G   b
18741  18752  pts/0   00:00:00  3T   a
```

As you can see, this output resembles the one of the `ps` command quite a bit. If this is a PiP root or PiP task, it is shown in the unfamiliar PIP

column (first character). The other numerical digit, “0-9,” stands for the PiP task, while “R” stands for root. In Section 2.2, PiP execution mode is represented by the second character.

The **pips** command has many options. Refer PiP man page (Section 1.4.1) for more details.

### 1.4.5 pip-gdb

PiP-aware **gdb** (GNU debugger) is known as **pip-gdb**. PiP tasks are introduced as inferiors of **GDB**. This is an illustration of a PiP-gdb debugging session. Note that the **pip-gdb** can only debug PiP tasks in the process mode.

```
(pip-gdb) info inferiors
  Num  Description                      Executable
*  4    process 1904 (pip 2)           /somewhere/pip-task-2
   3    process 1903 (pip 1)           /somewhere/pip-task-1
   2    process 1902 (pip 0)           /somewhere/pip-task-0
   1    process 1897 (pip root)        /somewhere/pip-root
```

### 1.4.6 pip-mode and printpipmode

The **pip-mode** command is to set PiP execution mode and the **printpip-mode** outputs the current execution mode (refer to Section 2.2 and 3.1.4).

### 1.4.7 libpip.so

The PiP library **libpip.so** can also run as a program, showing the information how the library was build and installed.

Listing 1.32: libpip.so - Execution Example

```
$ ${PIPLIBDIR}/libpip.so
Package:      Process-in-Processs
Version:      2.4.1
License:      the 2-clause simplified BSD License
Build OS:     Linux 5.10.104-linuxkit #1 SMP Thu Mar 17 17:08:06 UTC 2022
Build CC:     gcc (GCC) 8.5.0 20210514 (Red Hat 8.5.0-4)
Prefix dir:   /home/ahori/git/pip-2/install
PiP-glibc:    /home/ahori/pip-glibc/install/lib
ld-linux:     /home/ahori/pip-glibc/install/lib/ld-2.28.so
Commit Hash:  2485d3f923302ef03432bc52a5ddc3c4b0398fca
Debug build:  no
URL:          https://github.com/procinproc/PiP/
mailto:       procinproc-info@googlegroups.com
$
```

## 1.5 Summary

### PiP root and PiP task

- PiP applications must be compiled using the **pipfc** or **pipecc** commands, respectively, for C, C++, and Fortran.

- The `pip-exec` command allows PiP programs to be executed as PiP tasks.
- PiP programs can be invoked as regular programs to execute as non-PiP tasks.
- Static variables in a PiP application are privatized, meaning that each PiP task has its own set of the static variables, in contrast to the traditional multi-thread paradigm (i.e. OpenMP).
- In contrast to the traditional multi-process approach (i.e. MPI), PiP tasks have access to each other's data and share the same address space.

## PiP API

- The majority of PiP functions return a Linux-defined error code.
- Per address space, each PiP task has a unique `PIPID`.
- PiP root must call `pip_init()` to initialize PiP library (). The initialization function may or may not be called by child PiP tasks.
- By using the `pip_spawn()` or `pip_task_spawn()` functions, the PiP root can start new PiP tasks.
- Use the `pip_named_export()` and `pip_named_import()` methods to find the address for gaining access to the data of the other PiP tasks.
- `pip_barrier_wait()` is used for tasks to synchronize. Synchronization can also be accomplished using PThreads' synchronization functions.
- The invoking PiP process and PiP root are terminated by the `pip_exit()` method.
- By using a member of the `pip_wait` function family, the PiP root can wait for a created PiP task to finish.

## 1.6 Myths on PiP

### I fail to distinguish what PiP performs from shared memory.

The shared memory model enables access to the data of the other process. The XPMEM<sup>3</sup> shared memory mode facilitates the sharing of any existing memory region with the other process, in contrast to the POSIX shared

---

<sup>3</sup><https://github.com/hpc/xpmem>

memory paradigm, which only permits the sharing of newly allocated memory regions. I use the term "shared address space model" instead of "shared memory model" to describe PiP's technique for mapping memory areas. The shared address space and shared memory appear to be the same in order to access the data owned by the other.

However, the mechanisms behind the two memory theories differ greatly. To request shared memory, one must issue a system call to the OS kernel. The system calls that alter the memory mapping are rather expensive. When the addresses of the data are known, a PiP job can access any data owned by the others without making an expensive system call because PiP tasks, on the other hand, are mapped to a single memory address.

If the shared data are allocated dynamically or are distributed over the address space and difficult to pack into a memory area, the shared address space solution is advantageous.

### **A serious security risk might arise when many programs share the same address**

A serious security risk might arise when many programs share the same address. Programs that share an address space can run thanks to PiP. Facilitating efficient information sharing between programs is the main goal here. If there is no interaction between them, running them in the PiP environment is useless.

Programs for communication fundamentally have the same outcome. Even in the most straightforward case, if two programs are connected by a Linux/Unix pipe and one of the programs crashes, the other programs also crash after receiving the `SIGPIPE` signal. Communicating processes agree on when and how to communicate with others. The PiP scenario is not exceptional.

### **Sharing address space makes debugging difficult**

It is true that if one operation in a PiP environment destroys data that belongs to another, the consequences could be disastrous. If done maliciously, this cannot be prevented (see also above). It may be more difficult to debug the destruction if a software flaw caused it than the multi-process paradigm. There are two things to note in this situation: 1) the increased risk of actually destroying data rather than accessing an invalid memory location (`SIGSEGV`), and 2) the existence of numerous execution entities.

For the earlier point, the ASLR can be helpful. The bug's phenomena can change over time if ASLR is enabled. With the multi-thread example, the last point's situation is essentially the same.

Anyway, I haven't had any problems with issues based on this scenario yet.

**My program does not have any static variables and I do not need PiP.**

Programs can be written without using static variables. However, Glibc's routines use a lot of static variables. Some of the Glibc functions may be used by your runtime system. Therefore, it is generally quite difficult to write programs without using any static variables.

### **PiP may consume more memory than the other execution models**

A single address space is used to map the memory segments of active PiP activities, as shown in Listing 3.4. For instance, Listing 3.4 only displays the memory segments connected to Glibc (`libc-2.28.so`), which is loaded three (3) times in this case to operate one PiP root and two PiP tasks, each of which depends on Glibc. A memory map of `libc-2.28.so` is contained in each segment set. The same file is used to map all three segment sets, and only one set uses the same amount of RAM.

Each address space of a process in the multi-process paradigm only has one `libc-2.28.so` segment set, although another process also uses the same memory mapping as Glibc. Therefore, practically speaking, the memory needed to run PiP tasks is comparable to the memory needed to operate multiple processes. However, regardless of the number of threads, the multi-thread model only has one variable segment that is shared by all of the threads. Additionally, performing PiP jobs requires more memory than running multiple threads does.

### **There must be some hidden overhead for running PiP programs**

As of now, one overhead that is greater than the multi-process model is known. System calls for address space change, such `mmap()` and `brk()`. This is because any change to an address space inside the OS kernel must be locked, and lock contention causes more overhead. The scenario is the same when using the multi-thread model, and while `mmap()` has a higher overhead than the multi-process model, it is nearly identical when using the multi-thread model. Up to now, no additional overhead in PiP has been identified.

## Chapter 2

# PiP Advanced

So far, the basic of PiP is described, In this chapter, more detailed functionalities provided by PiP will be explained.

### 2.1 Rationale

The steps to spawn a PiP task are as follows (a more complete technique is provided in Section 3.1.1);

The steps to spawn a PiP task are: 1. construct a new name space by using the Linux `dlmopen()` function, 2. create a PiP task process (or thread) by using the Linux `clone()` system call, and 3. enter the starting function of a user application.

1. build a new name space by using the Linux `dlmopen()` function,
2. create a PiP task process (or thread) by using the Linux `clone()` system call, and
3. enter the starting function of a user application.

In contrast to `dlopen()`, the `dlmopen()` function can construct a new name space. Here, the global symbol names (functions and global variables) that must be resolved upon program loading make up the name space. Functions and variables may be privatized from the other PiP tasks by establishing a distinct name space.

It's crucial to call `dlmopen()` and `clone()` in the correct order. As it looked pretty reasonable, at first I tried to call them in the following order: `clone()`, followed by `dlmopen()`, but this did not work at all. This is the rationale behind why only PiP root can start PiP tasks and watch for their terminations.

The loaded address of a program is fixed by default in some circumstances (or in the majority of circumstances prior to CentOS/Redhat 8).

If so, PiP is unable to load more than one program in the same address space. The PiP executables must be built as PIE(Position Independent Executable) files in order for the programs to be able to be loaded at any arbitrary address. All PiP tasks must be built as PIE, which means they must be built using the `--piptask` option using `pipcc` or `pipfc` (or nothing to use the default). Keep in mind that the PiP root program could not be PIE.

PiP task can be formed by running the loaded program in a new name space with another thread. Unfortunately, nothing is ever that easy. Many problems are caused by Glibc. These problems are discussed in the following section.

## 2.2 Execution Mode

PiP library is made to function on Linux. According to Section 2.1, it is very dependent on the `dlopen()` and `clone()` functions. Particularly, the `clone()` function is called with a unique set of `CLONE` flags. Many Linux variants exist, and some of them do not support the specified `CLONE` flag combination (for example, McKernel<sup>1</sup>). There are two ways to execute PiP in such an environment: one involves executing `clone()` with a particular flag combination, and the other involves calling `pthread_create()` to start a PiP task while using the standard `CLONE` flag combination. **Process mode** refers to the former and **pthread mode** to the latter. Both modes keep the fundamental characteristics of the PiP, including address space sharing and flexible privatization.

### 2.2.1 Differences Between Two Modes

`CLONE` flag combinations ultimately determine the difference in the PiP execution mode. The `CLONE` flag values `CLONE_FS`, `CLONE_FILES`, `CLONE_SIGHAND`, and `CLONE_THREAD` are reset, but `CLONE_VM` and `CLONE_SYSVSEM` are set, in contrast to the pthread mode.

Table 2.1: Differences between two modes

|                        | Process Mode | Pthread Mode |
|------------------------|--------------|--------------|
| Address Space Sharing  | yes          | yes          |
| Variable Privatization | yes          | yes          |
| File Descriptors (FDs) | not shared   | shared       |

The primary distinctions between the two modes are shown in Table 2.1. There are a lot more variations, however the PiP library offers mode-agnostic

<sup>1</sup><https://github.com/ihkmckernel/mckernel>

functions so that users can develop PiP programs without worrying about the variations in mode.

Table 2.2: Mode-Agnostic Functions

| Mode-Agnostic                  | Process Mode               | Pthread Mode                   | note             |
|--------------------------------|----------------------------|--------------------------------|------------------|
| <code>pip_exit()</code>        | <code>exit()</code>        | <code>pthread_exit()</code>    | termination      |
| <code>pip_wait()</code>        | <code>wait()</code>        | <code>pthread_join()</code>    | wait termination |
| <code>pip_kill()</code>        | <code>kill()</code>        | <code>pthread_kill()</code>    | send signal      |
| <code>pip_sigmask()</code>     | <code>sigprocmask()</code> | <code>pthread_sigmask()</code> | signal mask      |
| <code>pip_signal_wait()</code> | <code>sigwait()</code>     | <code>sigwait()</code>         | wait signal      |
| <code>pip_yield()</code>       | <code>sched_yield()</code> | <code>pthread_yield()</code>   | yield            |

There are also predicate functions for users to know the current mode listed in Table 2.3.

Table 2.3: Execution Mode Predicates

| Function name                   | note              |
|---------------------------------|-------------------|
| <code>pip_is_threaded()</code>  | if pthread mode   |
| <code>pip_is_shared_fd()</code> | if FDs are shared |

In the current implementation, `pip_is_threaded()` and `pip_is_shared_fd()` have the same meaning. There may be situations in which those two terms have different meanings, so having those functions is necessary.

## 2.2.2 How to Specify Execution Mode

When calling `pip_init()` or setting the `PIP_MODE` environment variable at runtime, the execution mode can be selected. The function prototype for the `pip_init()` shown below. Up until this point, the first three arguments have been discussed.

```
int pip_init( int *pipidp,      [IN/OUT]
              int *ntasksp,     [IN/OUT]
              void **root_expp, [IN/OUT]
              int opts );       [IN]
```

The final `opts` parameter can take one of the following values: `PIP_MODE_PROCESS`, `PIP_MODE_THREAD`, or both, or zero. Zero has the same value as `PIP_MODE_PROCESS|PIP_MODE_PTHREAD`. The `PIP_MODE` environment can be a string that contains the words "process" or "pthread." The value of the `PIP_MODE` environment variable is checked when the `opts` value is zero or the value of either oring the two. If no environment is specified, the PiP library picks a suitable one. The environmental value and the ethical worth cannot be at odds with one another.



## 2.3 Spawning Tasks - Advanced

Further features of `pip_spawn()` and `pip_task_spawn()` that have not yet been explained will be covered in this section. The function prototypes for these functions are included below for your convenience;

---

```
int pip_task_spawn( pip_spawn_program_t *progp, [IN]
                  uint32_t coreno,             [IN]
                  uint32_t opts,               [IN]
                  int *pipidp,                 [IN/OUT]
                  pip_spawn_hook_t *hookp );   [IN]

int pip_spawn( char *filename,                 [IN]
               char **argv,                   [IN]
               char **envv,                   [IN]
               int coreno,                    [IN]
               int *pipidp,                   [IN/OUT]
               pip_spawnhook_t before,        [IN]
               pip_spawnhook_t after,         [IN]
               void *hookarg );               [IN]
```

---

In Section 1.2.1, the first argument of the `pip_task_spawn()` function is already covered in Section 1.2.1. To compact the first three arguments of `pip_spawn()`, use the `pip_spawn_program_t` structure.

### 2.3.1 Start Function

These PiP spawn routines eventually jump into the start function (`main()` or a user-specified one), as already demonstrated in Listing 1.2.1. PiP requires knowledge of the start function's address in order to enable this. Here, one of the two prerequisites must be satisfied:

- the starting function is defined as a global symbol, or
- if the starting function is defined as a local symbol then the executable file must not be stripped.

The `-rdynamic` option is used during program compilation with the `pipcc` and `pipfc` to make the symbol global for the `main()` function. When it comes to user-defined local symbols, the PiP library reads the executable file to spawn and attempts to locate the starting function using the ELF metadata. Unfortunately, if the local symbol information is stripped, it is lost, and PiP is unable to locate the starting function.

### 2.3.2 Stack Size

The `PIP_STACKSIZE` environment variable can be used to control the stack size of created PiP tasks. Its value can be prefixed with "T," "G,"

“M,” “K,” or “B” which, like the `OMP_STACKSIZE` environment provided by OpenMP, stands for the *TiB*, *GiB*, *MiB*, *KiB*, or *Byte* unit, respectively. *KiB* is presumed in the absence of a suffix. The environment variables `KMP_STACKSIZE`, `GOMP_STACKSIZE`, or `OMP_STACKSIZE` are similarly effective with the priority in this order unless `PIP_STACKSIZE` is provided. While `PIP_STACKSIZE` only impacts the stack size of PiP processes, `KMP_STACKSIZE`, `GOMP_STACKSIZE`, and `OMP_STACKSIZE` also affect the size of OpenMP threads.

### 2.3.3 CPU Core Binding

The spawned PiP task is bound to the designated CPU core by the `coreno` argument. This is the *N*th core number by default. The absolute core number should be 0Red with the `PIP_CPUCORE_ABS` flag if users want to specify the absolute core number. Specifying this flag might not have an impact on the core number specification if the core numbers are continuous. Only certain CPU architectures with non-contiguous core numbers exhibit this disparity (e.g., Fujitsu A64FX). To connect to the same CPU cores as the root process executing the spawn function, the `coreno` argument might be `PIP_CPUCORE_ASIS`.

### 2.3.4 File Descriptors and Spawn Hooks

Similar to how `fork()` works, file descriptors from the root process are duplicated and provided to the newly spawned child in the **process mode**. File descriptors are simply shared by PiP root and PiP tasks in the **pthread mode**.

The file descriptor is closed after calling the **before hook** mentioned below (if any), and then the start function is called if the close-on-exec flag of a file descriptor owned by the root process is set in process mode.

The structure containing the final three arguments of `pip_spawn()` is the final argument of `pip_task_spawn()`. Calling the `pip_spawn_hook()` function will set the `pip_spawn_hook_t` structure. Here is the prototype;

---

```
void pip_spawn_hook( pip_spawn_hook_t *hook,    [OUT]
                    pip_spawnhook_t before,    [IN]
                    pip_spawnhook_t after,     [IN]
                    void *hookarg ) {           [IN]
typedef int (*pip_spawnhook_t)( void* );
```

---

Before invoking the start method (such as `main()`), this structure’s **before** function is called when a PiP task is created. And when the PiP task is ready to finish, the **after** function is invoked. Both functions are invoked with the `hookarg`-specified argument, which can pass any kind of data.

In Linux/Unix, calling `fork()` and `execve()` generally results in the creation of a new process. Here, the parent process's file descriptors are transferred to the newly created child process. In many instances, such file descriptors are closed or duplicated, and other settings are changed in between calls to `fork()` and `execve()`. The task is only created by one function in PiP, hence it is not possible to use the same settings as when `fork&exec` is used. These hook functions are offered in order to accomplish this. Here is an illustration of one of these hook functions;

Listing 2.1: Before and After Hooks

---

```

#define _GNU_SOURCE
#include <pip/pip.h>
#include <stdlib.h>
#include <unistd.h>
#define FD_TASK          (10)
int before_hook( void *argp ) {
    int *fdp = (int*) argp;
    printf( "PID:%d Before Hook: fd=%d\n", getpid(), *fdp );
    fflush( stdout );
    dup2( 1, *fdp );
    return 0;
}
int after_hook( void *argp ) {
    int *fdp = (int*) argp;
    printf( "PID:%d After Hook:  fd=%d\n", getpid(), *fdp );
    fflush( stdout );
    close( *fdp );
    return 0;
}
int main( int argc, char **argv ) {
    int pipid, ntasks, pipid_task, i;
    int arg = FD_TASK;
    ntasks = strtol( argv[1], NULL, 10 );
    pip_init( &pipid, &ntasks, NULL, 0 );
    if( pipid == PIP_PIPID_ROOT ) {
        pip_spawn_program_t prog;
        pip_spawn_hook_t    hooks;
        pip_spawn_from_main( &prog, argv[0], argv,
                             NULL, NULL );
        pip_spawn_hook( &hooks,
                        before_hook,
                        after_hook,
                        &arg );
        printf( "PID:%d MAIN\n", getpid() );
        fflush( stdout );
        for( i=0; i<ntasks; i++ ) {
            pipid_task = i;
            pip_task_spawn( &prog, PIP_CPUCORE_ASIS,

```

```

        0, &pipid_task, &hooks );
    pip_wait( pipid_task, NULL );
}
} else {
    char *msg;
    asprintf( &msg, "Hello from PIPID:%d\n", pipid );
    write( FD_TASK, msg, strlen( msg ) );
    free( msg );
}
pip_fin();
return 0;
}

```

In this illustration, the **before hook** copies the root process' FD1 to FD10. Next, the job that was started uses FD10 to write a message. The **after hook** closes the FD10 finally. You should be aware that the root calls those hook functions to run them.

Listing 2.2: Before and After Hooks - Execution

```

$ ./hook 2
PID:37231 MAIN
PID:37232 Before Hook: fd=10
Hello from PIPID:0
PID:37232 After Hook:  fd=10
PID:37233 Before Hook: fd=10
Hello from PIPID:1
PID:37233 After Hook:  fd=10
$

```

## 2.4 Execution Context

Readers should be familiar with the execution context under PiP before I discuss the remaining arguments. The state of the CPU, or the contents of hardware registers, can be thought of as the general definition of the execution context. This definition might not be sufficient on PiP. Let's look at an example. Assume the same application contains a function called `foo` and is running as two PiP tasks. One PiP job can call the function of the other PiP task by supplying the function pointer and utilizing the **pip\_named\_export()** and **pip\_named\_import()**. The function `foo` also makes use of a static variable, let's say `var`. When task A calls task B's function `foo`, task B's variable, not task A's, is accessed by the called function (Listing 2.3 and 2.4).

Listing 2.3: Calling a Function of Another Task

```

#include <pip/pip.h>
int var;

```

```

int foo( void ) { return var; }
int main( int argc, char **argv ) {
    int pipid, ntasks, prev;
    int(*funcp)(void);
    pip_get_pipid( &pipid );
    pip_get_ntasks( &ntasks );
    prev = ( pipid == 0 ) ? ntasks - 1 : pipid - 1;
    var = pipid * 100;
    pip_named_export( foo, "foo%d", pipid );
    pip_named_import( prev, (void**) &funcp, "foo%d", prev );
    printf( "PIPID:%d foo(%d)=%d\n", pipid, prev, funcp() );
    return 0;
}

```

Sorry, this may detour off the main road, but in this example program, `pip_named_export()` and `pip_named_import()` are called in a different way than they were previously. Similar to Linux's `printf()`, the second and third arguments to the `pip_named_export()` and `pip_named_import()` functions, respectively, are really format strings followed by the argument(s) required by the format.

Listing 2.4: Function Call of Another Task - Execution

```

$ pip-exec -n 4 ./context
PIPID:1 foo(0)=0
PIPID:2 foo(1)=100
PIPID:3 foo(2)=200
PIPID:0 foo(3)=300
$

```

The execution context in a PiP environment may differ from that in a common sense environment, as shown in Listing 2.4, and occasionally its behavior changes in a very subtle way. This occurs pretty frequently in the PiP library and makes debugging quite challenging.

Furthermore, the CPU architecture and PIE implementation have a significant impact on the relationship between static variables and function addresses. For `x86_64` and `AArch64`, the aforementioned description is accurate, but not on `x86_32`. As a result, it is not advised to do this.

## Rationale

Some readers may be curious as to why this occurs. Let me elaborate. The address map contains this ruse. Listing 2.4 displays a section of the address map with three jobs executing, with a focus on the Glibc (`/lib64/libc-2.28.so`) segments.

```

...
7ffff53f8000-7ffff55a4000 r-xp 00000000 fe:01 3049686 /lib64/libc-2.28.so
7ffff55a4000-7ffff57a4000 ---p 001ac000 fe:01 3049686 /lib64/libc-2.28.so
7ffff57a4000-7ffff57a8000 r--p 001ac000 fe:01 3049686 /lib64/libc-2.28.so
7ffff57a8000-7ffff57aa000 rw-p 001b0000 fe:01 3049686 /lib64/libc-2.28.so

```

```

...
7fff69f6000-7fff6ba2000 r-xp 00000000 fe:01 3049686 /lib64/libc-2.28.so
7fff6ba2000-7fff6da2000 ---p 001ac000 fe:01 3049686 /lib64/libc-2.28.so
7fff6da2000-7fff6da6000 r--p 001ac000 fe:01 3049686 /lib64/libc-2.28.so
7fff6da6000-7fff6da8000 rw-p 001b0000 fe:01 3049686 /lib64/libc-2.28.so
...
7fff73d5000-7fff7581000 r-xp 00000000 fe:01 3049686 /lib64/libc-2.28.so
7fff7581000-7fff7781000 ---p 001ac000 fe:01 3049686 /lib64/libc-2.28.so
7fff7781000-7fff7785000 r--p 001ac000 fe:01 3049686 /lib64/libc-2.28.so
7fff7785000-7fff7787000 rw-p 001b0000 fe:01 3049686 /lib64/libc-2.28.so
...

```

Glibc segments come in three sets. The final (readable and write) segment of each set is where the static variables are located. By using the offset from the instruction (program counter relative addressing mode) to the variable, an instruction can access a static variable. In order to keep all offsets equal, the gap size between the code segment (at the top of the set) and the variable segment (at the bottom of the set) is crucial, and all gap sizes must be the same. Variables and instructions are connected in PIE<sup>2</sup> in this fashion, allowing PIE applications and shared objects that were created with the PIC option to be loaded in any place.

Unfortunately, not all CPU architectures implement this addressing mode<sup>3</sup>. For example, x86\_32 doesn't. This architecture sacrifices one general purpose register in order to point the variable segment, causing performance to suffer as a result. Moreover, Listing 2.3 program displays a distinct behavior.

## 2.5 Debugging Support

Some environment variable settings may help debugging PiP programs.

### 2.5.1 PIP\_STOP\_ON\_START

This environment variable is to stop (by sending SIGSTOP to spawned PiP task just before calling the **before hook** (Section 2.3.4), or jumping to the starting function if before hook is not given. The value of this environment must meet the following format;

---

```
PIP_STOP_START=[<script-file>]@<PIPID>
```

---

The optional `<script>` is a shell script to be executed on the suspension, and `<PIPID>` is the **PIPID** to be suspended. If `<PIPID>` is -1, then all spawned PiP tasks will be stopped. The `<script>` is invoked with three parameters; PID and **PIPID** of the stopped PiP task, followed by a path to the program of the task. Do not forget to set the *executable* bit on this `<script-file>` file.

---

Listing 2.5: Stop-on-start Script Example

---

<sup>2</sup>This is not the case if not compiled as PIE.

```
#!/bin/sh
PID=$1
PIPID=$2
PROG='basename $3'
echo "###" $0 "###" "${PROG} ${PID} ${PIPID}"
pips -f ${PID} # strace, ltrace, pip-gdb, ...
kill -CONT ${PID}
```

Listing 2.5 shows an example of the script for the **PIP\_STOP\_ON\_START**. Here, **pips** command is invoked instead of some debugging command<sup>3</sup>. Note that the target task is already stopped by delivering the SIGSTOP signal. Somehow you have to explicitly deliver the SIGCONT signal to the task if you want to resume the task. Listing 2.6 shows the result of **PIP\_STOP\_ON\_START** execution with this script file.

Listing 2.6: Stop-on-start Script Example - Execution

```
$ pipcc --silent hello.c -o hello
$ echo $PIP_STOP_ON_START
onstart.cmd@2
$ pip-exec -n 4 ./hello
Hello World
Hello World
PiP-INFO[36954(R):R] PiP task[2] (PID=36957) is SIGSTOPed and executing 'onstart.cmd' script
### onstart.cmd ### hello 36957 2
Hello World
File "/home/ahori/git/pip-2/install/bin/pips", line 228
    from __future__ import print_function
    ^
SyntaxError: from __future__ imports must occur at the beginning of the file
Hello World
$
```

### 2.5.2 PIP\_GDB\_SIGNALS

This environment variable **PIP\_GDB\_SIGNALS** is to set the signals to trigger some actions by specifying the **PIP\_SHOW\_MAPS** and **PIP\_SHOW\_PIPS**, followed by the PiP-gdb invocation. The value of this environment is as follows;

---

```
PIP_GDB_SIGNALS=[ <SIGNAME> ] { "+"|"-" <SIGNAME> }
```

---

The possible **<SIGNAME>** vale are listed below;

Here, ‘‘ALL’’ means all signals list in this table. Each signal name in this table can be concatenated by using the plus (+) and/or minus (-) symbols. For example, “ALL-SIGUSR1” indicates the all signals excluding SIGUSR1. “SIGUSR1+SIGUSR2+SIGINT” represents SIGUSR1, SIGUSR2 and SIGINT.

---

<sup>3</sup>All examples are executed on a Docker environment but **ptrace** (and other commands using **ptrace**) was unable to run in this example even withthe **--cap-add=SYS\_PTRACE** Docker option (I confirmed **gdb** worked). So **pips** was used instead in this example.

Table 2.4: Possible Signal Names for **PIP\_GDB\_SIGNALS**

SIGHUP  
 SIGINT  
 SIGQUIT  
 SIGILL  
 SIGABRT  
 SIGFPE  
 SIGINT  
 SIGSEGV  
 SIGPIPE  
 SIGUSR1  
 SIGUSR2  
 ALL

### 2.5.3 PIP\_SHOW\_MAPS

If **PIP\_SHOW\_MAPS** environment is set to “on” and the a signal specified by the **PIP\_GDB\_SIGNALS** is delivered, then the address map (Listing 3.4, for example) will be shown.

### 2.5.4 PIP\_SHOW\_PIPS

If **PIP\_SHOW\_PIPS** environment is set to “on” and the a signal specified by the **PIP\_GDB\_SIGNALS** is delivered, then **pips** command (Section 1.4.4) is invoked to show the status of the other PiP tasks in the same address space.

### 2.5.5 PIP\_GDB\_PATH and PIP\_GDB\_COMMAND

When **PIP\_GDB\_PATH** is set to the path to **pip-gdb** a signal specified by the **PIP\_GDB\_SIGNALS** is delivered, then PiP gdb (Section 1.4.5) will be invoked. If the value of the **PIP\_GDB\_COMMAND** environment is set to a valid filename and if the filename contains some GDB commands, then PiP-gdb will be invoked to work with this command file.

## 2.6 Malloc routines

Suppose that we are making a producer-consumer style program using PiP, a PiP task is a producer and another PiP task is a consumer. Unlike the conventional process model, there is no need of calling IPC (Inter Process Communication) system call in PiP. All we have to do is just passing pointers pointing data to be passed from the producer to the consumer.



Here, an issue arises. If the passing data is allocated by a `malloc()` routine, then the passed data is `free()`ed by the consumer. As described so far, each PiP task has its own `malloc()` and `free()` routines associated with static variables holding and maintaining a memory pool. The consumer receives the data allocated from the memory pool of the producer and tries to `free()` it when it becomes unnecessary. However, the `free()` routine on the consumer has no knowledge about the producer-allocated memory region and fails (Figure 2.1). I named this situation *cross-malloc-free*.

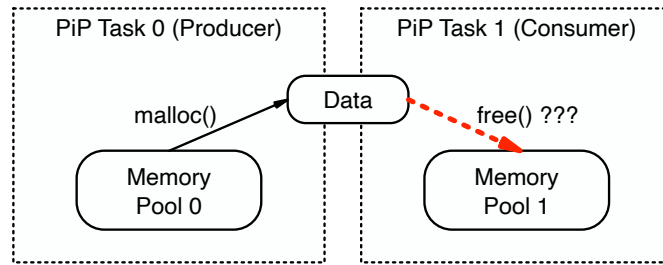


Figure 2.1: Cross-Malloc-Free Issue

I tried this by using the `malloc()` routines provided Glibc and I found that this works in most cases, not always. I do not know why this works (again, **in most cases**) with the Glibc `malloc()` routines, but I believe this situation must be avoided.

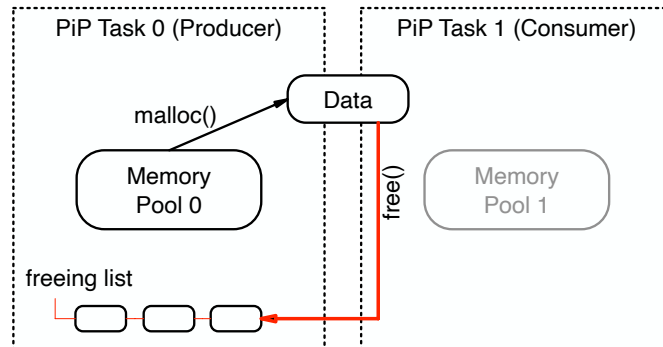


Figure 2.2: Cross-Malloc-Free with Freeing List

To deal with this, PiP library wraps `malloc()` routines as shown in Table 3.3. When a memory region is allocated, the `malloc()` wrapper function embeds the information who allocates that region. When this region is to be `free()`ed, the `free()` wrapper function connects the region to the freeing list of the task allocating the region. The regions in the freeing list are eventually `free()`ed when one of the `malloc()` wrapper functions (Figure 2.2).

## 2.7 XPMEM

As described in Section 1.6, XPMEM is known to provide a shared memory model which is more convenient than the POSIX shared memory. Again, the shared address space which PiP provides includes the shared memory model which XPMEM and POSIX shared memory provide. Thus, the same functionalities of XPMEM can also be implemented by using PiP.

The PiP library provides the same functions which are provided by XPMEM. Those who have programs using XPMEM can easily switch to using PiP. By using PiP, there is no need of installing XPMEM kernel module. Most importantly, XPMEM functions provided by PiP work much faster than those of XPMEM. This is because no system call is involved to map memory segment of the other process(es) in PiP. Indeed, most XPMEM functions almost do nothing since the other processes are already mapped from the beginning in PiP.

## Chapter 3

# PiP Internals

### 3.1 PiP Implementation

#### 3.1.1 Spawning Tasks

Before PiP version 2.4, PiP tasks were created with the procedure as follows;

1. The spawned program is loaded by calling `dlopen()`,
2. Glibc is initialized in the execution context of the loaded program,
3. Call `clone()` or `pthread_create()` (chosen by the **PIP\_MODE** environment setting) to spawn the PiP task,
4. The before hook is called if any, and finally
5. Jump into the starting function.

From PiP version 2.4, the wrapper functions listed in Table 3.3 were introduced. When implementing the wrapper functions, I noticed that wrapping the `dlsym()` is almost impossible.

A function wrapper is usually implemented as; 1) obtain the wrapping function address by calling the `dlsym()` with the `RTLD_NEXT` argument, 2) do the wrapping job before and/or after calling the original function. The most of the Glibc `malloc()` routines has the other weak symbols (`malloc()` and `__libc_malloc()`, for example) and users can call the Glibc `malloc()` routines without calling `dlsym()`. If there is no such weak symbol, we cannot create a wrapper function for `dlsym()`. How can I wrap a Glibc function without calling `dlsym()`?

To solve this issue, I implemented another program, so called **ldpip.so** to load the PiP library and user program. here is the details of new spawning process;

1. Load **ldpip.so** in the PiP library package by calling `dlopen()` and jump into a function defined inside of it,

2. The starting function of **ldpip.so** initializes Glibc,
3. Obtain Glibc function addresses to wrap them later by **libpip.so**,
4. Load **libpip.so** by calling `dlopen()`,
5. Load a user program by calling `dlopen()`,
6. Call `clone()` or `pthread_create()` (chosen by the **PIP\_MODE** environment setting) to spawn the PiP task,
7. Jump into a function inside of PiP library and initialize the PiP library,
8. The before hook is called if any, and finally
9. Jump into the starting function in the user program.

At the time of loading **ldpip.so**, no wrapper functions are defined in this program and obtaining the Glibc function addresses is easy, just referencing them. After loading the **libpip.so** and jumping into a function defined in **libpip.so** where the wrapping functions are defined, the Glibc functions to be wrapped are now wrapped by using the function table created by **ldpip.so**<sup>1</sup>.

The Glibc initialization<sup>2</sup> must be done with the execution context (Section 2.4) of the spawned PiP task. In the older version of PiP library, this was done by; 1) calling `dlsym()` to the loaded handle, returned by `dl[m]open()`, to obtain the initialization function and then 2) call the function. In the new implementation, the initialization was done by simply calling the initialization function from the **ldpip.so** where the execution context is the same with that of PiP task.

Thus, by introducing PiP loader program (**ldpip.so**), things can go in a simpler way.

### 3.1.2 Calling `clone()` System Call

As described in Section `refsec:spawn-details`, PiP library calls the `clone()` system call with a special flag combination. The `clone()` system call has many arguments and some of them are hard to implement, I decided to wrap the `clone()` system call to modify only the flag setting.

One issue to wrap the `clone()` system call is that the `clone()` is called not only the PiP library, but also some other libraries (e.g., PThread library). A simple function wrapping cannot handle both situations where

---

<sup>1</sup>If actual dynamic linking would be done in the order of `dl[m]open()`, then the wrapping functions in **libpip.so** would not work as described here. As long as I checked, the Glibc (**libc.so**) is at the last of the search order of **ld-linux.so**, and this works.

<sup>2</sup>Calling `__ctype_init()`

the flags needs to be changed when it is called by the PiP library but it should not change the flag when called by some other libraries.

To solve this issue and protect the `clone()` system call from being called from PiP library and from another simultaneously, a special locking mechanism was implemented by using the `test-and-set` atomic instruction. When the PiP library is about to call `pthread_create()` call, which eventually calls the `clone()` system call, it locks by using the `test-and-set` instruction with the value of current thread ID (TID). The wrapper function of the `clone` firstly tries to lock, but it fails with value of the current TID, then it is the case of calling from the PiP library. If the lock succeeds, then it is called by some other library. In the former case, the original `clone()` is called with the modified flags. In the latter case, the `clone()` is called with the same argument with the wrapper function call. Needless to say, the lock is unlocked after returning from the original `clone()` system call.

### 3.1.3 Execution Mode in Details

There are two sub-modes in the PiP's `process mode`, `process:preload` and `process:piplone`<sup>3</sup>. The `process:preload` mode is implemented by wrapping the `clone()` function described above and the `process:piplone` is implemented to have another `pthread_create()` like function implemented in the patched-glibc (Section 3.2.2). If PiP library is configured to use the patched-glibc, then `process:piplone` is taken, otherwise `process:preload` is taken.

### 3.1.4 Name of PiP Tasks

Some readers may wonder how the `pips` command (Section 1.4.4) can distinguish the PiP tasks and the other normal processes and/or threads. In Linux, each process and thread can have a name, which can be seen by the `top` command in the `COMMAND` column. The PiP library sets the command name by calling the `prctl()` system call (in `process mode` or `pthread_setname_np()` call (in `pthread mode`. The PiP library uses the first two characters of the name (Table 3.1 and 3.2)).

Table 3.1: Command Name Setting (1st char.)

| First char. | Distinction | Note  |
|-------------|-------------|---|
| R           | PiP Root    |   |
| 0..9        | PiP Task    | the least significant digit of <code>PIPID</code> |

<sup>3</sup>In PiP implementation earlier than version 2.4, there was another mode `process:got`. But this becomes obsolete in the newer versions.

Table 3.2: Command Name Setting (Second char.)

| 2nd char. | Execution Mode         | Abbreviation | Note     |
|-----------|------------------------|--------------|----------|
| :         | <b>process:preload</b> | L            | obsolete |
| ;         | <b>process:piplone</b> | C            |          |
| .         | <b>process:got</b>     | G            |          |
|           | <b>pthread</b>         | T            |          |

The name may have up to 16 characters and PiP occupies the first two characters. The remaining 14 characters are used for representing original command name. The abbreviation column in Table 3.2 shows the characters used by the **pip-mode** command (Section 1.4.6) to specify the PiP execution mode. The **pips** command can now distinguish the normal processes or threads from the PiP roots and PiP tasks by using those first 2 characters of the command.

## 3.2 Issues related to Linux Kernel, Glibc and Tools

This section will explain about the issues when implementing PiP.

### 3.2.1 Loading a Program

Before going into the details about the Glibc issues when implementing PiP, readers should understand how a program is loaded into memory. This subsection describes only about the program loading procedure of Linux, apart from PiP implementation.

When the Linux's `execve()` system call to run a program, the Linux kernel opens and reads the executable file, searching the ELF section named `“.interp.”`

Listing 3.1: `“.interp”` Section of the `ps` command

```
$ readelf -a /usr/bin/ps | grep interpreter
[Requesting program interpreter: /lib64/ld-linux-x86-64.so.2]
$
```

Listing 3.1 shows the value of the `“.interp”` section, `/lib64/ld-linux-x86-64.so.2`. The Linux kernel invokes the loader specified by the `“.interp”` section and asks the loader to load a program specified by the `execve()` parameter. Then the loader loads the program and additionally loads and links the shared libraries required to run the program. Once everything is loaded, the loader jumps into the starting function defined in Glibc to initialize Glibc and finally user-defined `main()` function is called.

The program loader, often simply called `ld-linux.so`, is loaded once per address space and kept in memory until the end of the process (see also Listing 3.4). This is responsible for any loading process by resolving the

external symbol references. The Glibc functions defined in the `libdl.so` (`-ldl`), such as `dlopen()`, `dlopen()`, `dlsym()` and so on, are just API and their functional bodies exist in this program loader.

### 3.2.2 Glibc

PiP provides a new execution model which cannot be categorized into neither the process model nor the thread model. In this new model, although its name is not yet given, tasks share the same address space like the thread model, but maintaining the variable privatization like the process model. This execution model is novel and not yet recognized by most of the tool chains provided by Linux and others. Indeed, the most of the time to develop PiP was devoted to find niches in Glibc.

#### PiP Task is Unable to Spawn PiP Task

As described in Section 2.1, the order of calling `dlopen()` and `clone()` is important. This restriction also means that a PiP task cannot spawn a PiP task as a child of spawning PiP task because this breaks the restriction. Thus, the current PiP implementation inhibits for a PiP task to call `pip_spawn()`.

#### Recycling PiP Tasks

As far as I tested, the resources; name space, loaded PIE program, and shared libraries required by the PIE program, are not released by calling the `dlclose()`. I believe this issue can be fixed by patching the Glibc, however, I decided not to do so. Thus, once a PiP task is created, then the **PIPID** of the task will not be recycled even the PiP task terminates. The reason of my decision will also be discussed in Section 3.3.1.

#### Number of name spaces

The number of name spaces which the `dlopen()` can create is hard-coded as 16. Considering PiP tasks run in parallel and the number of CPU cores nowadays, this number of 16 is apparently too small. The PiP package provides PiP-glibc where the number of name spaces is increased, up to 300 PiP tasks<sup>4</sup>.

The name space table resides in the `ld-linux.so` and this means that the `.interp` ELF section of PiP programs must be changed so that the program is loaded by the new `ld-linux.so`. This can be done by specifying `--dynamic-linker` option of the GNU linker and the `pipcc` and `pipfc` do this.

---

<sup>4</sup>Once I asked Glibc development members to increase the size, but they did not accept my opinion. Refer [https://sourceware.org/bugzilla/show\\_bug.cgi?id=23978](https://sourceware.org/bugzilla/show_bug.cgi?id=23978)

The name space table resides at the top of a structure in `ld-linux.so`. Some Glibc functions refer to the members in this structure directly. This causes another problem. Once the size of the name space table is changed, the addresses of the other members in the same structure are also changes. As described, only one `ld-linux.so` can be loaded in an address space. As a result, all PiP programs sharing the same address must be linked with the same Glibc.

## PiP-gdb

The `ld-linux.so` embeds some information for debugging into the loaded program. Unfortunately, I found that this code fragment resides on the pass calling the `ld-linux.so` from the top (by the kernel), not on the pass called from `dlopen()` and `dlopen()`<sup>5</sup>. The patched PiP-glibc fixed this issue. Thus, the `pip-gdb` command (Section 1.4.5) can only work with the PiP programs linked with the patched PiP-glibc.

## Global lock

Most programs are linked with Glibc and PiP programs are no exception. PiP allows to run multiple PiP programs in the same address space. This means that each PiP task has its own Glibc. And the simultaneous calls of some Glibc functions may not work because of a race condition.

To avoid this condition, PiP library provides the functions, `pip_glibc_lock()` and `pip_glibc_unlock()`, to serialize the Glibc function calls. The following Glibc functions are wrapped by PiP library to introduce the lock and users do not have to care the race.

Table 3.3: Glibc functions wrapped by PiP library

|                           |                           |                             |
|---------------------------|---------------------------|-----------------------------|
| <code>dlsym</code>        | <code>dlopen</code>       | <code>dlopen</code>         |
| <code>dlinfo</code>       | <code>dlclose</code>      | <code>dlerror</code>        |
| <code>dladdr</code>       | <code>dlvsym</code>       | <code>getaddrinfo</code>    |
| <code>freeaddrinfo</code> | <code>gai_strerror</code> | <code>pthread_create</code> |
| <code>pthread_exit</code> |                           |                             |
| <code>malloc</code>       | <code>free</code>         | <code>calloc</code>         |
| <code>realloc</code>      | <code>memalign</code>     | <code>posix_memalign</code> |

The functions `pthread_exit()` and below in this table have another reason to have function wrappers. The wrapping reason of `pthread_exit()` will be explained in the Section 2.2 and the reason of wrapping `malloc` routines will be explained in Section 2.6.

<sup>5</sup>I guess the loaded code by using `dlopen()` or `dlopen()` cannot be debugged.



These listed functions may not be complete. There can be a case where some other Glibc functions may suffer from the race condition. This problem can be avoided by introducing the above locking functions. This lock can be used recursively and users can avoid deadlock situation easily.

## Constructors and Destructors

The constructors and destructors are used in C++ programs. Constructors and destructors are list of functions. Generally, constructor functions are called just before the program begins, and destructors functions are called when the program is about to exit.

In PiP, the behavior of the constructors and destructors is somewhat different. To explain this, I should start explaining how the constructors and destructors are implemented in general. The constructor functions are listed in the `.init_array` section of an ELF file. The destructor functions are listed in the `.fini_array` section. Constructors are called when `ld-linux.so` finishes loading and linking objects. Destructors are called when `dlclose()` is called.

Now back to PiP. Again, constructors are called inside of the call of `dlmopen()` when spawning a PiP task. The `dlmopen()` is called by the PiP root process. Thus, the constructors of a program are called by the root. Here is the example;

Listing 3.2: Constructors and Destructors

---

```
#include <pip/pip.h>
#include <iostream>
#include <unistd.h>
#include <sys/types.h>
char *pipidstr( void ) {
    static char idstr[32];
    int pipid;
    if( pip_get_pipid( &pipid ) != 0 ) {
        sprintf( idstr, "[R] PID:" );
    } else {
        sprintf( idstr, "[%d] PID:", pipid );
    }
    return idstr;
}
static int x = 0;
class Hello {
public:
    Hello(void ) {
        std::cout << pipidstr() << getpid() << " Hello" <<
            std::endl;
    }
    ~Hello(void ) {
        std::cout << pipidstr() << getpid() << " Bye" <<
```

```

        std::endl;
    }
};
Hello hello;
int main() {
    std::cout << pipidstr() << getpid() << " MAIN " <<
        std::endl;
    return 0;
}

```

Listing 3.2 is a C++ program having a constructor and destructor. When the constructor of this program is called, the PiP library is not yet initialized, and the `pip_get_pipid()` return an error (EPERM). So, the function `pipidstr()` takes care of this situation. Listing 3.3 shows the execution example of this program. As shown, the PIDs output by the constructors are not the same with the ones of the PiP tasks.

Listing 3.3: Constructors and Destructors - Execution

```

$ ./hello
[R] PID:37019 Hello
[R] PID:37019 MAIN
[R] PID:37019 Bye
$ pip-exec -n 2 ./hello
[R] PID:37020 Hello
[0] PID:37021 MAIN
[0] PID:37021 Bye
[R] PID:37020 Hello
[1] PID:37022 MAIN
[1] PID:37022 Bye
$

```

## LD\_PRELOAD

LD\_PRELOAD only works with PiP root, not PiP tasks. This is because `dlopen()` simply ignores the LD\_PRELOAD environment setting.

## Shared Objects

Some shared objects, such as GCC related runtime libraries, must be located in the same directory where the `ld-linux.so` does. The `pipnlibs` shell script found in the PiP-glibc package makes symbolic links of the shared objects in the `/lib64` directories to meet with the restriction.

## Loading Program by dlmopen()

The Glibc in CentOS/RedHat 8 (and possibly newer ones) does not allow to load a program by the `dlmopen()` function<sup>6</sup>. The `pip-unpie` program is to cheat this Glibc restriction. This program is automatically executed by the `pipcc` or `pipfc` when creating a PiP executable, and not to be invoked by users directly.

### 3.2.3 Glibc RPATH Setting

When using Spack<sup>7</sup>, it automatically adds RPATHs for every program and PiP is no exception. A problem arises when to install the PiP-glibc by using Spack. When the PiP-glibc is built by using Spack, Spack adds the RPATH setting to the compiled Glibc, but it is not allowed in CentOS/Redhat 8 to load Glibc with the RPATH setting. To avoid this, PiP-glibc has a program (`annul_rpath`) to unset the RPATH setting of the compiled Glibc.

### 3.2.4 Linux

#### Heap Segment

There is another issue which comes from Linux kernel, not from Glibc. Before explaining this issue, let us start from the how an address space is composed.

Listing 3.4: A Memory Map Example

```
00400000-00402000 r-xp 00000000 00:71 77130812 /PiP/bin/pip-exec
00601000-00602000 r--p 00001000 00:71 77130812 /PiP/bin/pip-exec
00602000-00603000 rw-p 00002000 00:71 77130812 /PiP/bin/pip-exec
00603000-00624000 rw-p 00000000 00:00 0 [heap]
7ffe8000000-7ffe8021000 rw-p 00000000 00:00 0
7ffe8021000-7ffefc000000 ---p 00000000 00:00 0
7ffefc000000-7ffefc000000 ---p 00000000 00:00 0
7ffefc000000-7ffff0000000 rwxp 00000000 00:00 0
7ffff0000000-7ffff0021000 rw-p 00000000 00:00 0
7ffff0021000-7ffff4000000 ---p 00000000 00:00 0
7ffff4000000-7ffff46db000 r-xp 00000000 00:71 79448679 /PiP/example/a.out
7ffff46db000-7ffff48da000 ---p 00001000 00:71 79448679 /PiP/example/a.out
7ffff48da000-7ffff48db000 r--p 00000000 00:71 79448679 /PiP/example/a.out
7ffff48db000-7ffff48dc000 rw-p 00001000 00:71 79448679 /PiP/example/a.out
7ffff48dc000-7ffff48f6000 r-xp 00000000 00:71 77130790 /PiP/lib/libpip.so.0
7ffff48f6000-7ffff4af6000 ---p 0001a000 00:71 77130790 /PiP/lib/libpip.so.0
7ffff4af6000-7ffff4af7000 r--p 0001a000 00:71 77130790 /PiP/lib/libpip.so.0
7ffff4af7000-7ffff4af8000 rw-p 0001b000 00:71 77130790 /PiP/lib/libpip.so.0
7ffff4af8000-7ffff4bf8000 rw-p 00000000 00:00 0
7ffff4bf8000-7ffff4da4000 r-xp 00000000 fe:01 3445390 /lib64/libc-2.28.so
7ffff4da4000-7ffff4fa4000 ---p 001ac000 fe:01 3445390 /lib64/libc-2.28.so
7ffff4fa4000-7ffff4fa8000 r--p 001ac000 fe:01 3445390 /lib64/libc-2.28.so
7ffff4fa8000-7ffff4faa000 rw-p 001b0000 fe:01 3445390 /lib64/libc-2.28.so
7ffff4faa000-7ffff4fae000 rw-p 00000000 00:00 0
7ffff4fae000-7ffff4fc5000 r-xp 00000000 fe:01 3446072 /lib64/libpthread-2.28.so
7ffff4fc5000-7ffff51c4000 ---p 00017000 fe:01 3446072 /lib64/libpthread-2.28.so
7ffff51c4000-7ffff51c5000 r--p 00016000 fe:01 3446072 /lib64/libpthread-2.28.so
7ffff51c5000-7ffff51c6000 rw-p 00017000 fe:01 3446072 /lib64/libpthread-2.28.so
7ffff51c6000-7ffff51ca000 rw-p 00000000 00:00 0
7ffff51ca000-7ffff51cc000 r-xp 00000000 fe:01 3445775 /lib64/libdl-2.28.so
7ffff51cc000-7ffff53cc000 ---p 00002000 fe:01 3445775 /lib64/libdl-2.28.so
7ffff53cc000-7ffff53cd000 r--p 00002000 fe:01 3445775 /lib64/libdl-2.28.so
```

<sup>6</sup>Refer [https://sourceware.org/bugzilla/show\\_bug.cgi?id=11754#c15](https://sourceware.org/bugzilla/show_bug.cgi?id=11754#c15). I tested this situation but I cannot find this problem with PiP.

<sup>7</sup>`PiPKW{pip-unpie}program`

```

7fff53cd000-7fff53ce000 rw-p 00003000 fe:01 3445775 /lib64/libdl-2.28.so
7fff53ce000-7fff53d6000 r-xp 00000000 00:71 77130791 /PiP/lib/ldpip.so.0
7fff53d6000-7fff55d5000 ---p 00008000 00:71 77130791 /PiP/lib/ldpip.so.0
7fff55d5000-7fff55d6000 r--p 00007000 00:71 77130791 /PiP/lib/ldpip.so.0
7fff55d6000-7fff55d7000 rw-p 00008000 00:71 77130791 /PiP/lib/ldpip.so.0
7fff55d7000-7fff55d8000 ---p 00000000 00:00 0
7fff55d8000-7fff55d8000 rwxp 00000000 00:00 0
7fff55d8000-7fff55d9000 r-xp 00000000 00:71 79448679 /PiP/example/a.out
7fff55d9000-7fff567d8000 ---p 00001000 00:71 79448679 /PiP/example/a.out
7fff567d8000-7fff567d9000 r--p 00000000 00:71 79448679 /PiP/example/a.out
7fff567d9000-7fff567da000 rw-p 00001000 00:71 79448679 /PiP/example/a.out
7fff567da000-7fff567f4000 r-xp 00000000 00:71 77130790 /PiP/lib/libpip.so.0
7fff567f4000-7fff569f4000 ---p 0001a000 00:71 77130790 /PiP/lib/libpip.so.0
7fff569f4000-7fff569f5000 r--p 0001a000 00:71 77130790 /PiP/lib/libpip.so.0
7fff569f5000-7fff569f6000 rw-p 0001b000 00:71 77130790 /PiP/lib/libpip.so.0
7fff569f6000-7fff56ba2000 r-xp 00000000 fe:01 3445390 /lib64/libc-2.28.so
7fff56ba2000-7fff56da2000 ---p 001ac000 fe:01 3445390 /lib64/libc-2.28.so
7fff56da2000-7fff56da6000 r--p 001ac000 fe:01 3445390 /lib64/libc-2.28.so
7fff56da6000-7fff56da8000 rw-p 001b0000 fe:01 3445390 /lib64/libc-2.28.so
7fff56da8000-7fff56dac000 rw-p 00000000 00:00 0
7fff56dac000-7fff56dc3000 r-xp 00000000 fe:01 3446072 /lib64/libpthread-2.28.so
7fff56dc3000-7fff56fc2000 ---p 00017000 fe:01 3446072 /lib64/libpthread-2.28.so
7fff56fc2000-7fff56fc3000 r--p 00016000 fe:01 3446072 /lib64/libpthread-2.28.so
7fff56fc3000-7fff56fc4000 rw-p 00017000 fe:01 3446072 /lib64/libpthread-2.28.so
7fff56fc4000-7fff56fc8000 rw-p 00000000 00:00 0
7fff56fc8000-7fff56fca000 r-xp 00000000 fe:01 3445775 /lib64/libdl-2.28.so
7fff56fca000-7fff571ca000 ---p 00002000 fe:01 3445775 /lib64/libdl-2.28.so
7fff571ca000-7fff571cb000 r--p 00002000 fe:01 3445775 /lib64/libdl-2.28.so
7fff571cb000-7fff571cc000 rw-p 00003000 fe:01 3445775 /lib64/libdl-2.28.so
7fff571cc000-7fff571d4000 r-xp 00000000 00:71 77130791 /PiP/lib/ldpip.so.0
7fff571d4000-7fff573d3000 ---p 00008000 00:71 77130791 /PiP/lib/ldpip.so.0
7fff573d3000-7fff573d4000 r--p 00007000 00:71 77130791 /PiP/lib/ldpip.so.0
7fff573d4000-7fff573d5000 rw-p 00008000 00:71 77130791 /PiP/lib/ldpip.so.0
7fff573d5000-7fff57581000 r-xp 00000000 fe:01 3445390 /lib64/libc-2.28.so
7fff57581000-7fff57781000 ---p 001ac000 fe:01 3445390 /lib64/libc-2.28.so
7fff57781000-7fff57785000 r--p 001ac000 fe:01 3445390 /lib64/libc-2.28.so
7fff57785000-7fff57787000 rw-p 001b0000 fe:01 3445390 /lib64/libc-2.28.so
7fff57787000-7fff5778b000 rw-p 00000000 00:00 0
7fff5778b000-7fff577a2000 r-xp 00000000 fe:01 3446072 /lib64/libpthread-2.28.so
7fff577a2000-7fff5779a1000 ---p 00017000 fe:01 3446072 /lib64/libpthread-2.28.so
7fff5779a1000-7fff5779a2000 r--p 00016000 fe:01 3446072 /lib64/libpthread-2.28.so
7fff5779a2000-7fff579a3000 rw-p 00017000 fe:01 3446072 /lib64/libpthread-2.28.so
7fff579a3000-7fff579a7000 rw-p 00000000 00:00 0
7fff579a7000-7fff579a9000 r-xp 00000000 fe:01 3445775 /lib64/libdl-2.28.so
7fff579a9000-7fff579a9000 ---p 00002000 fe:01 3445775 /lib64/libdl-2.28.so
7fff579a9000-7fff57baa000 r--p 00002000 fe:01 3445775 /lib64/libdl-2.28.so
7fff57baa000-7fff57bab000 rw-p 00003000 fe:01 3445775 /lib64/libdl-2.28.so
7fff57bab000-7fff57bc5000 r-xp 00000000 00:71 77130790 /PiP/lib/libpip.so.0
7fff57bc5000-7fff57dc5000 ---p 0001a000 00:71 77130790 /PiP/lib/libpip.so.0
7fff57dc5000-7fff57dc6000 r--p 0001a000 00:71 77130790 /PiP/lib/libpip.so.0
7fff57dc6000-7fff57dc7000 rw-p 0001b000 00:71 77130790 /PiP/lib/libpip.so.0
7fff57dc7000-7fff57dea000 r-xp 00000000 fe:01 3446237 /lib64/ld-2.28.so
7fff57dea000-7fff57fe4000 rw-p 00000000 00:00 0
7fff57fe4000-7fff57fe8000 r--p 00000000 00:00 0 [vvar]
7fff57fe8000-7fff57fea000 r-xp 00000000 00:00 0 [vdso]
7fff57fea000-7fff57feb000 r--p 00023000 fe:01 3446237 /lib64/ld-2.28.so
7fff57feb000-7fff57fff000 rw-p 00024000 fe:01 3446237 /lib64/ld-2.28.so
7fffffde000-7fffffde000 rwxp 00000000 00:00 0 [stack]
ffffffffff600000-ffffffffff601000 r-xp 00000000 00:00 0 [vsyscall]

```

Listing 3.4 shows an example of the output of doing “cat /proc/⟨PID⟩/maps.” Here, “pip-exec -n 2 ./a.out” was executed, resulting one **pip-exec** process and two ./a.out tasks. The file, /proc/⟨PID⟩/maps, lists all memory segments in an address space of the process ⟨PID⟩ usually. A loaded shared object has consecutive three or four segments; executable, gap (not accessible, if any), constants and data. The rightmost column of a line indicates the `mmap()`ed filename, the second from the left column indicates the permission of the memory segment. ‘r’ is readable, ‘w’ is writable, ‘x’ is executable and ‘p’ is private (copy-on-write). There are also some special segments whose filename is in a pair of square brackets; [stack], [heap], and so on. These are created by the Linux kernel for special purposes as their names suggest. The segments having no filename are created by the `mmap()` system call.

Remember, this is the address space of running one PiP root and two PiP tasks, resulting to have all the segments of the three tasks. Note that the all three tasks have exactly the same `/proc/<PID>/maps` content, and there are three sets of a shared library and only one set of `ld-linux.so` (`ld-2.28.so`) can be seen in Listing 3.4.

Usually, the heap segment, mainly used by `malloc()`, exists only one per address space. As shown in this example, there is only one heap segment, meaning the heap segment is shared by three tasks (one for root and two for PiP tasks).

The size of the heap segment can be increased or decreased by calling the `brk()` system call. Most cases, there are two calls of `brk()` to allocate or deallocate heap memory, one for obtaining the current heap end address and another for setting the new heap end address. This is exactly what the Glibc's `sbrk()` does. Apparently, this API is not thread-safe at all, and thus, the shared heap memory cannot be used safely by PiP tasks.

Fortunately, the `malloc()` routines in Glibc is designed to check if there are two or more name spaces and if so they do not use the `brk()` system call, use `mmap()` instead. So, the Glibc `malloc()` routines can work with PiP without any problem. However, if some other routines use the `brk()` system call (or `sbrk()` Glibc function), for example, replacing the Glibc `malloc()` routines with some other `malloc()` implementation, then this shared heap may result in a problem.

## Core File

Suppose that we have a catastrophic situation and all PiP tasks and their root process dump core files of their own. On the current Linux, a core file is associated with a process (including threads inside of it). Thus, each PiP task and the root may produce core, resulting to have many core files. Here, the address space of them are shared and the created core files and all of them are almost the same excepting the CPU state.

Let me explain this with an example. Suppose that we have PiP task *A* and *B* running on the same address space of the root *P*, and an error happens resulting all *P*, *A*, and *B* produce core files. There can be a small time difference when to produce each core file. When the first core file, of *A* for instance, is being created, the other *P* and *B* are still running and the memory of the shared address space can be altered by those running tasks. `gdb`, however, assumes that a core file is a consistent snapshot of memory and CPU state. The above PiP situation breaks this assumption. If *B* produces another core file, may or may not be caused by the error on *A*, the same situation can happen. Thus, the `pip-gdb` command (Section 1.4.5) does not support for debugging from a core file. To solve this issue, PiP-aware OS kernel to have the consistent core files is needed.

### 3.2.5 Tools

As described, PiP sets a special combination of the `clone()` flags. As a result of this, some tools do not work. Here is the list of tools which are known to work or not at the time of this writing<sup>8</sup>.

Table 3.4: Compatibility of Tools

| Compatible                                 | Incompatible          |
|--|-----------------------|
| <code>strace</code><br><code>ltrace</code> | <code>valgrind</code> |

## 3.3 Remaining Issues

### 3.3.1 Retrieving Memory

Let us suppose a case where PiP task *A* pass a pointer to PiP task *B* (Listing 1.8, for example). After then, task *A* terminates for some reason. What if task *B* tries to dereference the pointer to access data which task *A* had? This situation can also happen if the string obtained by calling `getenv()` is passed to the other task. The consequence of this may introduce difficult situation hard to debug. This situation must be detected by compilers and/or tools which are aware of PiP-style execution model.

So, I decided not to reclaim any memory resources when a task terminates, not calling `dlclose()` nor `free()`. In the current PiP implementation, **PIPID** can be allocated only once. And not releasing memory resource will not cause further problem.

---

<sup>8</sup>`ltrace` depends on its version

## Chapter 4

# PiP Installation

There are several ways to install PiP listed below;

- Building from source code
- **pip-pip** command
- Using *Spack*

It was scheduled to use RPM (yum) and Docker, but they are not available at the time of this writing.

### 4.1 Building from Source Code

Usually, building full PiP package consists of the following steps;

1. Building PiP-glibc (optional)
2. Building PiP library
3. Building PiP-gdb (optional)

The **1** and **3** steps are optional, but PiP-gdb requires PiP-glibc. So the possible combinations are;

- PiP library only
- PiP-glibc and PiP library, and
- PiP-glibc, PiP library and PiP-gdb.

Listing **4.1** shows a typical case of building full set of PiP software package.

Listing 4.1: Building from Source Code

```
$ PIPTOP=$PWD
$ git clone https://github.com/procinproc/PiP-glibc.git
...
$ mkdir glibc-build
$ pushd glibc-build
$ ../PiP-glibc/build.sh ${PIPTOP}/install
...
$ popd
$ git clone https://github.com/procinproc/PiP.git
...
$ pushd PiP
$ ./configure --prefix=${PIPTOP}/install --with-glibc-libdir=${PIPTOP}/install/lib
...
$ make install
...
$ popd
$ git clone https://github.com/procinproc/PiP-Testsuite.git
...
$ pushd PiP-Testsuite
$ ./configure --with-pip=${PIPTOP}/install
...
$ make test
...
$ popd
$ git clone https://github.com/procinproc/PiP-gdb.git
...
$ pushd PiP-gdb
$ ./build.sh --prefix=${PIPTOP}/install --with-pip=${PIPTOP}/install
...
$ ./test.sh
...
$ popd
$
```

## 4.2 pip-pip command

The procedure to install full set of PiP package might be cumbersome, but installing PiP package by using the **pip-pip** (<https://github.com/procinproc/PiP-pip>) command is much easier.

Listing 4.2: PiP-pip installation example

```
$ git clone https://github.com/procinproc/PiP-pip.git
$ cd PiP-pip
$ ./pip-pip --yes

RedHat/CentOS:      8
CPU Architecture:   x86_64
```



```

List of installations
GITHUB PiP-v2
  Prefix dir:    ${PWD}/install/x86_64_centos-8_github_pip-2
  Work dir:     ${PWD}/work/x86_64_centos-8_github_pip-2

  .....

Summary
OK      git https://github.com/procinproc/PiP.git@pip-2 ${PWD}/
        install/x86_64_centos-8_github_pip-2
$

```

## 4.3 Using Spack

Spack<sup>1</sup> is another installation tool designed for the HPC software packages and PiP can also be installed by using Spack. Listing 4.3 shows the example of installing PiP (including PiP-glibc)<sup>2</sup>.

Listing 4.3: Spack installation example

```

$ git clone https://github.com/spack/spack.git
$ cd bin
$ ./spack install process-in-process
...
$

```

---

<sup>1</sup><https://spack.io>

<sup>2</sup>Unfortunately, the current version does not install PiP-gdb for some reason.

# Bibliography

- [1-Hori18] Atsushi Hori, Min Si, Balazs Gerofi, Masamichi Takagi, Jai Dayal, Pavan Balaji, and Yutaka Ishikawa. Process-in-process: Techniques for practical address-space sharing. In *Proceedings of the 27th International Symposium on High-Performance Parallel and Distributed Computing, HPDC '18*, page 131143, New York, NY, USA, 2018. Association for Computing Machinery. (**Note: This is the first paper proposing Process-in-Process.**).

**URL** <https://doi.org/10.1145/3208040.3208045>

**Abstract:** The two most common parallel execution models for many-core CPUs today are multiprocess (e.g., MPI) and multithread (e.g., OpenMP). The multiprocess model allows each process to own a private address space, although processes can explicitly allocate shared-memory regions. The multithreaded model shares all address space by default, although threads can explicitly move data to thread-private storage. In this paper, we present a third model called process-in-process (PiP), where multiple processes are mapped into a single virtual address space. Thus, each process still owns its process-private storage (like the multiprocess model) but can directly access the private storage of other processes in the same virtual address space (like the multithread model). The idea of address-space sharing between multiple processes itself is not new. What makes PiP unique, however, is that its design is completely in user space, making it a portable and practical approach for large supercomputing systems where porting existing OS-based techniques might be hard. The PiP library is compact and is designed for integrating with other runtime systems such as MPI and OpenMP as a portable low-level support for boosting communication performance in HPC applications.

We showcase the uniqueness of the PiP environment through both a variety of parallel runtime optimizations and direct use in a data analysis application. We evaluate PiP on several platforms including two high-ranking supercomputers, and we measure and analyze the performance of PiP by using a variety of micro- and macro-kernels, a proxy application as well as a data analysis application.

- [2-Hori20] A. Hori, B. Geroft, and Y. Ishikawa. An implementation of user-level processes using address space sharing. In *2020 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pages 976–984, 2020. **(Note: This is the second paper on PiP experimental version (v3) implementing Bi-Level Thread and User-Level Process.)**

**Abstract:** There is a wide range of implementation approaches to multi-threading. User-level threads are efficient because threads can be scheduled by a user-defined scheduling policy that suits the needs of the specific application. However, user-level threads are unable to handle blocking system-calls efficiently. To the contrary, kernel-level threads incur large overhead during context switching. Kernel-level threads are scheduled by the scheduling policy provided by the OS kernel which is hard to customize to application needs. We propose a novel thread execution model, *bi-level thread*, that combines the best aspects of the two conventional thread implementations. A bi-level thread can be either a kernel-level thread or a user-level thread at runtime. Consequently, the context switching overhead of a bi-level thread is as low as that of user-level threads, but thread scheduling can be defined by user policies. Blocking system-calls, on the other hand, can be called as a kernel-level thread without blocking the execution of other user-level threads. Furthermore, the proposed bi-level thread is combined with an address space sharing technique which allows processes to share the same virtual address space. Processes sharing the same address space can be scheduled with the same technique as user-level threads, thus we call this implementation a *user-level process*. However, the main difference between threads and processes is that threads share most of the kernel state of the underlying process, such as

process ID and file descriptors, whereas different processes do not. A user-level process must guarantee that the system-calls always access the appropriate kernel information that belongs to the particular process. We call this *system-call consistency*. In this paper, we show that the proposed bi-level threads, implemented in an address space sharing library, can resolve the blocking system-call issue of user-level threads, while at the same time it retains system-call consistency for the user-level process. A prototype implementation, ULP-PiP, proves these concepts and the basic performance of the prototype is evaluated. Evaluation results using asynchronous I/O indicate that the overlap ratio of our implementation outperforms that in Linux.

- [3-Ouyang20] Kaiming Ouyang, Min Si, Atsushi Hori, Zizhong Chen, and Pavan Balaji. Cab-mpi: Exploring interprocess work-stealing towards balanced mpi communication. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC '20*. IEEE Press, 2020. (**Note: This is the first MPI optimization paper using PiP.**)

**Abstract:** Load balance is essential for high-performance applications. Unbalanced communication can cause severe performance degradation, even in computation-balanced BSP applications. Designing communication-balanced applications is challenging, however, because of the diverse communication implementations at the underlying runtime system. In this paper, we address this challenge through an interprocess work-stealing scheme based on process-memory-sharing techniques. We present CAB-MPI, an MPI implementation that can identify idle processes inside MPI and use these idle resources to dynamically balance communication workload on the node. We design throughput-optimized strategies to ensure efficient stealing of the data movement tasks. We demonstrate the benefit of work stealing through several internal processes in MPI, including intranode data transfer, pack/unpack for noncontiguous communication, and computation in one-sided accumulates. The implementation is evaluated through a set of microbenchmarks and proxy applications on Intel Xeon and Xeon Phi platforms.

- [4-Ouyang21] Kaiming Ouyang, Min Si, Astushi Hori, Zizhong Chen, and Pavan Balaji. Daps: A dynamic asynchronous progress stealing model for mpi communication. In *2021 IEEE International Conference on Cluster Computing (CLUSTER)*, pages 516–527, 2021. **(Note: This is the another MPI optimization paper using PiP.).**

**Abstract:** MPI provides nonblocking point-to-point and one-sided communication models to help applications achieve communication and computation overlap. These models provide the opportunity for MPI to offload data transfer to low level network hardware while the user process is computing. In practice, however, MPI implementations have to often handle complex data transfer in software due to limited capability of network hardware. Therefore, additional asynchronous progress is necessary to ensure prompt progress of these software-handled communication. Traditional mechanisms either spawn an additional background thread on each MPI process or launch a fixed number of helper processes on each node. Both mechanisms may degrade performance in user computation due to statically occupied CPU resources. The user has to fine-tune the progress resource deployment to gain overall performance. For complex multiphase applications, unfortunately, severe performance degradation may occur due to dynamically changing communication characteristics and thus changed progress requirement. This paper proposes a novel Dynamic Asynchronous Progress Stealing model, called Daps, to completely address the asynchronous progress complication. Daps is implemented inside the MPI runtime. It dynamically leverages idle MPI processes to steal communication progress tasks from other busy computing processes located on the same node. The basic concept of Daps is straightforward; however, various implementation challenges have to be resolved due to the unique requirements of inter-process data and code sharing. We present our design that ensures high performance while maintaining strict program correctness. We compare Daps with state-of-the-art asynchronous progress approaches by utilizing both microbenchmarks and HPC proxy applications.

- [5-Hori22] Atsushi Hori, Kaiming Ouyang, Balazs Gerofi, and Yutaka Ishikawa. On the difference between shared memory and shared address space in hpc communication. In Dhabaleswar K. Panda and Michael Sullivan, editors, *Supercomputing Frontiers*, pages 59–78, Cham, 2022. Springer International Publishing. (**Note: This is the third paper on PiP comparing shared memory model and shared address space model which PiP provides.**).

**Abstract:** Shared memory mechanisms, e.g., POSIX shmem or XPMEM, are widely used to implement efficient intra-node communication among processes running on the same node. While POSIX shmem allows other processes to access only newly allocated memory, XPMEM allows accessing any existing data and thus enables more efficient communication because the send buffer content can directly be copied to the receive buffer. Recently, the shared address space model has been proposed, where processes on the same node are mapped into the same address space at the time of process creation, allowing processes to access any data in the shared address space. Process-in-Process (PiP) is an implementation of such mechanism. The functionalities of shared memory mechanisms and the shared address space model look very similar – both allow accessing the data of other processes –, however, the shared address space model includes the shared memory model. Their internal mechanisms are also notably different. This paper clarifies the differences between the shared memory and the shared address space models, both qualitatively and quantitatively. This paper is not to showcase applications of the shared address space model, but through minimal modifications to an existing MPI implementation it highlights the basic differences between the two models. The following four MPI configurations are evaluated and compared; 1) POSIX Shmem, 2) XPMEM, 3) PiP-Shmem, where intra-node communication is implemented to utilize POSIX shmem but MPI processes share the same address space, and 4) PiP-XPMEM, where XPMEM functions are implemented by the PiP library (without the need for linking to XPMEM library). Evaluation is done using the Intel MPI benchmark suite and six HPC benchmarks (HPCCG,

miniGhost, LULESH2.0, miniMD, miniAMR and mpiGraph). Most notably, mpiGraph performance of PiP-XPmem outperforms the XPmem implementation by almost 1.5x. The performance numbers of HPCCG, miniGhost, miniMD, LULESH2.0 running with PiP-Shmem and PiP-XPmem are comparable with those of POSIX Shmem and XPmem. PiP is not only a practical implementation of the shared address space model, but it also provides opportunities for developing new optimization techniques, which the paper further elaborates on.

[6-Ouyang22] Kaiming Ouyang. PhD thesis, University of California Riverside, 2022. (**Dr. Ouyang’s Ph.D. Thesis using PiP.**)

**Abstract:** In exascale computing era, applications are executed at larger scale than ever before, which results in higher requirement of scalability for communication library design. Message Passing Interface (MPI) is widely adopted by the parallel application nowadays for interprocess communication, and the performance of the communication can significantly impact the overall performance of applications especially at large scale. There are many aspects of MPI communication that need to be explored for the maximal message rate and network throughput. Considering load balance, communication load balance is essential for high-performance applications. Unbalanced communication can cause severe performance degradation, even in computation-balanced Bulk Synchronous Parallel (BSP) applications. MPI communication imbalance issue is not well investigated like computation load balance. Since the communication is not fully controlled by application developers, designing communication-balanced applications is challenging because of the diverse communication implementations at the underlying runtime system. In addition, MPI provides non-blocking point-to-point and one-sided communication models where asynchronous progress is required to guarantee the completion of MPI communications and achieve better communication and computation overlap. Traditional mechanisms either spawn an additional background thread on each MPI process or launch a fixed number of helper processes on each node. For

complex multiphase applications, unfortunately, severe performance degradation may occur due to dynamically changing communication characteristics. On the other hand, as the number of CPU cores and nodes adopted by the applications greatly increases, even the small message size MPI collectives can result in the huge communication overhead at large scale if they are not carefully designed. There are MPI collective algorithms that have been hierarchically designed to saturate inter-node network bandwidth for the maximal communication performance. Meanwhile, advanced shared memory techniques such as XPMEM, KNEM and CMA are adopted to accelerate intra-node MPI collective communication. Unfortunately, these studies mainly focus on large-message collective optimization which leaves small- and medium-message MPI collectives suboptimal. In addition, they are not able to achieve the optimal performance due to the limitations of the shared memory techniques. To solve these issues, we first present CAB-MPI, an MPI implementation that can identify idle processes inside MPI and use these idle resources to dynamically balance communication workload on the node. We design throughput-optimized strategies to ensure efficient stealing of the data movement tasks. The experimental results show the benefits of CAB-MPI through several internal processes in MPI, including intranode data transfer, pack/unpack for noncontiguous communication, and computation in one-sided accumulates through a set of microbenchmarks and proxy applications on Intel Xeon and Xeon Phi platforms. Then, we propose a novel Dynamic Asynchronous Progress Stealing model (Daps) to completely address the asynchronous progress complication; Daps is implemented inside the MPI runtime, and it dynamically leverages idle MPI processes to steal communication progress tasks from other busy computing processes located on the same node. We compare Daps with state-of-the-art asynchronous progress approaches by utilizing both microbenchmarks and HPC proxy applications, and the results show the Daps can outperform the baselines and achieve less idleness during asynchronous communication. Finally, to further improve MPI collectives performance, we propose Process-in-Process



based Multiobject Interprocess MPI Collective (PiP-MColl) design to maximize small and medium-message MPI collective performance at a large scale. Different from previous studies, PiP-MColl is designed with efficient multiple senders and receivers collective algorithms and adopts Process-in-Process shared memory technique to avoid unnecessary system call and page fault overhead to achieve the best intra- and inter-node message rate and throughput. We focus on three widely used MPI collectives MPI Scatter, MPI Allgather and MPI Allreduce and apply PiP-MColl to them. Our microbenchmark and real-world HPC application experimental results show PiP-MColl can significantly improve the collective performance at a large scale compared with baseline PiP-MPICH and other widely used MPI libraries such as OpenMPI, MVAPICH2 and Intel MPI.

# Index

- ASLR, 25
- close-on-exec, 31
- Linux Command
  - gdb, 50
  - ltrace, 51
  - top, 42
- Linux Define Symbol
  - CLONE\_FILES, 28
  - CLONE\_FS, 28
  - CLONE\_SIGHAND, 28
  - CLONE\_SYSVSEM, 28
  - CLONE\_THREAD, 28
  - CLONE\_VM, 28
  - EPERM, 9, 47
  - RTLD\_NEXT, 40
  - SIGPIPE, 25
  - SIGSEGV, 25
  - WEXITSTATUS, 13
  - WIFEXITED, 13
  - WIFSIGNALED, 13
  - WTERMSIG, 13
- Linux Environment Variable
  - GOMP\_STACKSIZE, 31
  - KMP\_STACKSIZE, 31
  - LD\_PRELOAD, 47
  - OMP\_STACKSIZE, 31
- Linux Function
  - \_\_ctype\_init(), 41
  - \_\_libc\_malloc(), 40
  - brk, 26
  - brk(), 50
  - clone, 27, 28
  - clone(), 40–42, 51
  - dl[m]open(), 41
  - dlclose(), 44, 46
  - dlopen, 27, 28
  - dlopen(), 40, 44–48
  - dlopen, 27
  - dlopen(), 41, 44, 45
  - dlsym(), 40, 41, 44
  - execve, 8, 32
  - execve(), 43
  - exit, 15
  - fork, 31, 32
  - free(), 38
  - gettimeofday, 17
  - malloc, 38, 40, 50
  - malloc(), 38, 40, 50
  - mmap, 26
  - mmap(), 49
  - prctl(), 42
  - printf, 34
  - pthread\_exit(), 45
  - pthread\_barrier\_destroy, 16
  - pthread\_barrier\_init, 16
  - pthread\_barrier\_wait, 16
  - pthread\_create, 28
  - pthread\_create(), 40–42
  - pthread\_mutex, 19
  - pthread\_setname\_np(), 42
  - sbrk(), 50
  - wait, 9, 13
- Linux Variable
  - environ, 8
  - main, 11, 12, 30, 31
  - MPI, 3, 5, 8
    - mpiexec, 4
  - MPI process, 8
  - OpenMP, 3, 5

PIC, 35  
 PIE, 28, 35  
 PiP Command  
   annul\_rpath, 48  
   pip-check, 21  
   pip-exec, 3, 4, 8, 21, 22, 24, 49  
   pip-gdb, 23, 37, 45, 50  
   pip-mode, 23, 43  
   pip-pip, 52, 53  
   pip-unpie, 48  
   pipcc, 3, 4, 21, 28, 30, 44, 48  
     --cflags, 21  
     --lflags, 21  
     --pipboth, 21  
     --piproot, 21  
     --piptask, 21, 28  
     --which, 21  
   pipfc, 21, 28, 30, 44, 48  
     --cflags, 21  
     --lflags, 21  
     --pipboth, 21  
     --piproot, 21  
     --piptask, 21, 28  
     --which, 21  
   pipnlibs, 47  
   pips, 22, 23, 36, 37, 42, 43  
   printpipmode, 23  
 PiP Define Symbol  
   PIP\_CPUCORE\_ABS, 31  
   PIP\_CPUCORE\_ASIS, 8, 31  
   PIP\_MODE\_PROCESS, 29  
   PIP\_MODE\_PTHREAD, 29  
   PIP\_MODE\_THREAD, 29  
   PIP\_PIPID\_ANY, 9  
   PIP\_PIPID\_ROOT, 8, 10  
   PIP\_STACKSIZE, 30  
 PiP Environment Variable  
   PIP\_GDB\_COMMAND, 37  
   PIP\_GDB\_PATH, 37  
   PIP\_GDB\_SIGNALS, 36, 37  
   PIP\_MODE, 29, 40, 41  
   PIP\_SHOW\_MAPS, 36, 37  
   PIP\_SHOW\_PIPS, 36, 37  
   PIP\_STACKSIZE, 31  
   PIP\_STOP\_ON\_START, 36  
   pipcc  
     CC, 21  
   pipfc  
     FC, 21  
 PiP Function  
   pip\_spawn(), 24  
   pip\_barrier\_fin, 16  
   pip\_barrier\_init, 16  
   pip\_barrier\_wait, 16, 17  
   pip\_exit, 15  
   pip\_exit(), 24, 29  
   pip\_fin, 9  
   pip\_get\_pipid(), 47  
   pip\_gettime, 17  
   pip\_glibc\_lock(), 45  
   pip\_glibc\_unlock(), 45  
   pip\_init, 8–10, 17, 29  
   pip\_init(), 24  
   pip\_is\_shared\_fd, 29  
   pip\_is\_shared\_fd(), 29  
   pip\_is\_threaded, 29  
   pip\_is\_threaded(), 29  
   pip\_kill(), 29  
   pip\_named\_export, 7, 33, 34  
   pip\_named\_export(), 24  
   pip\_named\_import, 7, 33, 34  
   pip\_named\_import(), 24  
   pip\_sigmask(), 29  
   pip\_signal\_wait(), 29  
   pip\_spawn, 8, 9, 11, 12, 30, 31  
   pip\_spawn(), 44  
   pip\_spawn\_from\_func, 12  
   pip\_spawn\_from\_main, 12  
   pip\_spawn\_hook, 31  
   pip\_task\_spawn, 11, 12, 30, 31  
   pip\_task\_spawn(), 24  
   pip\_trywait, 14  
   pip\_trywait\_any, 14  
   pip\_wait, 9, 13, 14  
   pip\_wait(), 29  
   pip\_wait\_any, 14  
   pip\_yield(), 29  
 PiP Term

- after hook, 33
- before hook, 31, 33, 35
- ldpip.so, 40, 41
- libpip.so, 23, 41
- pip\_barrier\_wait(), 24
- PIPID, 7–10, 14, 24, 35, 42, 44, 51
- process:got, 42, 43
- process:pipeclone, 42, 43
- process:preload, 42, 43
- Process mode, 28
- process mode, 31, 42
- pthread, 43
- pthread mode, 28, 31, 42
- PiP Type
  - pip\_barrier\_t, 17
  - pip\_spawn\_hook\_t, 31
  - pip\_spawn\_program\_t, 12, 30
- test-and-set, 42
- variable privatization, 6