

MIR User Guide

February 10, 2015

1 Introduction

MIR is a task-based runtime system library written using C99 that supports detailed thread-based and task-based performance profiling. MIR scales well for medium-grained task-based programs. MIR supports subset of the OpenMP 3.0 tasks interface and a low level native interface for writing task-based programs. MIR allows the user to experiment with memory distribution policies and different scheduling policies. Example: Locality-aware scheduling and data distribution on NUMA systems.

2 Intended Audience

MIR is intended to be used by advanced task-based programmers. Knowledge of compilation and runtime system role in task-based programming is required to use and appreciate MIR.

3 Installation

3.1 Mandatory Requirements

- Machine with x86 architecture.
- Linux kernel later than January 2012.
- GCC.
- Binutils.
- Scons build system.
- R (for executing scripts)

- R packages
 - data.table (for data structure transformations)

3.2 Optional Requirements

Enabling extended features such as profiling, locality-aware scheduling and data distribution requires:

- libnuma and numactl (for data distribution and locality-aware scheduling on NUMA systems)
- GCC with OpenMP support (for linking task-based OpenMP programs)
- PAPI (for reading hardware performance counters during profiling)
- Paraver (for visualizing thread execution traces)
- Python 2.X and 3.X (for executing various scripts)
- Intel Pin sources (for profiling instructions executed by tasks)
- R packages:
 - optparse (for parsing data)
 - igraph (for task graph processing)
 - RColorBrewer (for colors)
 - gdata, plyr, dplyr, data.table (for data structure transformations)
- yEd (for task graph viewing, preferred)
- Graphviz (for task graph viewing)
- Cytoscape (for task graph viewing)

3.3 Source Structure

The source repository is structured intuitively. Files and directories have purpose-oriented names.

```
. : MIR_ROOT
|--docs : documentation
|--src : runtime system sources
|   |--scheduling : scheduling policies
|   |--arch : architecture specific sources
|--scripts
|   |--helpers : helpful scripts, dirty hacks
```

```

|--profiling : all things related to profiling
|   |--task
|   |--thread
|--programs : test programs, benchmarks
|   |--common : build scripts
|   |--native : native interface programs
|       |--fib
|           |--helpers : testing scripts
|--bots : BOTS port
|--omp : OpenMP interface programs
|   |--fib

```

3.4 Build

Follow below steps to build the basic runtime system library.

- Set MIR_ROOT environment variable.

```
$ export MIR_ROOT=<MIR source repository path>
```

Tip: Add the export statement to .bashrc to avoid repeated initialization.

- Build.

```
$ cd $MIR_ROOT/src
$ scons
```

Expert Tip: Ensure MIR_ROOT/src/SConstruct matches your build intention.

3.4.1 Enabling data distribution and locality-aware scheduling on NUMA systems

- Install libnuma and numactl.
- Create an empty file called HAVE_LIBNUMA.

```
$ touch $MIR_ROOT/src/HAVE_LIBNUMA
```

- Clean and rebuild MIR.

```
$ cd $MIR_ROOT/src
$ scons -c && scons
```

3.5 Testing

Try different runtime system configurations and program inputs on Fibonacci in `MIR_ROOT/programs/native/fib`. All programs under `MIR_ROOT/programs` can be used for testing.

```
$ cd $MIR_ROOT/programs/native/fib
$ scons -c
$ scons
$ ./fib-verbose
$ ./fib-debug
$ ./fib-opt
```

Note: A dedicated test suite will be added soon, so watch out for that!

4 Programming

4.1 OpenMP 3.0 Tasks Interface

A restricted subset of OpenMP 3.0 tasks — the `task` and `taskwait` constructs — is supported. Although minimal, the subset is sufficient for writing most task-based programs.

The `parallel` construct is deprecated. A team of threads is created when `mir_create` is called. The team is disbanded when `mir_destroy` is called.

Note: OpenMP tasks are supported by intercepting GCC translated calls to GNU libgomp. OpenMP 3.0 task interface support is therefore restricted to programs compiled using GCC.

4.1.1 Tips for writing MIR-supported OpenMP programs

- Initialize and release the runtime system explicitly by calling `mir_create` and `mir_destroy`.

- Do not think in terms of threads.
 - Do not use the `parallel` construct to share work.
 - Do not use barriers to synchronize threads.
- Think solely in terms of tasks.
 - Use the `task` construct to parallelize work.
 - Use clauses `shared`, `firstprivate` and `private` to indicate the data environment.
 - Use `taskwait` to synchronize tasks.
- Use `mir_lock` instead of the `critical` construct or use OS locks such as `pthread_lock`.
- Use GCC atomic builtins for flushing and atomic operations.
- Study example programs in `MIR_ROOT/programs/omp`.

A simple set of steps for producing MIR-supported OpenMP programs is given below:

1. When parallel execution is required, create a `parallel` block with `default(none)` followed immediately by a `single` block. The `default(none)` clause avoids incorrect execution due to assumed sharing rules.
2. Use the `task` construct within the `single` block to parallelize work.
3. Synchronize tasks using the `taskwait` construct explicitly. Do not rely on implicit barriers and taskwaits.
4. Parallelizing work inside a master task context is helpful while interpreting profiling results.
5. Compile and link with the native OpenMP implementation (preferably `libgomp`) and check if the program runs correctly.
6. Include `mir_public_int.h`. Call `mir_create` in the beginning of `main` and call `mir_destroy` at the end of `main`. Delete `parallel` and `single` blocks.
7. Compile and link with the appropriate MIR library (`opt/debug`). The program is now ready.

The native interface example rewritten using above steps is shown below.

```
int main(int argc, char *argv[])
{
    // Initialize the runtime system
    mir_create();
```

```

#pragma omp task
{
    // Now parallelize the work involved
    // Work in this case: create as many tasks
    // ... as there are threads
    int num_workers = mir_get_num_threads();
    for(int i=0; i<num_workers; i++)
    {
        #pragma omp task firstprivate(i)
        foo(i);
    }

    // Wait for tasks to finish
    #pragma omp taskwait
}
// Wait for master task to finish
#pragma omp taskwait
// Release runtime system resources
mir_destroy();

return 0;
}

```

4.2 Native Interface

Look at `mir_public_int.h` in `MIR_ROOT/src` for interface details and programs in `MIR_ROOT/programs/native` for interface usage examples. A simple program using the native interface is shown below.

```

#include "mir_public_int.h"
void foo(int id)
{
    printf(stderr, "Hello from task %d\n", id);
}

// Task outline function argument
struct foo_wrapper_arg_t
{
    int id;
};

// Task outline function
void foo_wrapper(void* arg)
{
    struct foo_wrapper_arg_t* farg = (struct foo_wrapper_arg_t*)(arg);
}

```

```

    foo(farg->id);
}

int main(int argc, char *argv[])
{
    // Initialize the runtime system
    mir_create();

    // Create as many tasks as there are threads
    int num_workers = mir_get_num_threads();
    for(int i=0; i<num_workers; i++)
    {
        struct foo_wrapper_arg_t arg;
        arg.id = i;
        mir_task_create((mir_tfunc_t) foo_wrapper,
                        &arg,
                        sizeof(struct foo_wrapper_arg_t),
                        0, NULL, NULL);
    }

    // Wait for tasks to finish
    mir_task_wait();

    // Release runtime system resources
    mir_destroy();

    return 0;
}

```

4.3 Compiling and Linking

Add `-lmir-opt` to `LDFLAGS`. Enable MIR to intercept function calls correctly by adding `-fno-inline-functions -fno-inline-functions-called-once -fno-optimize-sibling-calls -fno-omit-frame-pointer -g` to `CFLAGS` and/or `CXXFLAGS`.

4.4 Runtime Configuration

MIR has several runtime configurable options which can be set using the environment variable `MIR_CONF`. Set the `-h` flag to see available configuration options.

```

$ cd $MIR_ROOT/test/fib
$ scons
$ MIR_CONF="-h" ./fib-opt 3
MIR_INFO: Valid options in MIR_CONF environment variable ...

```

```

-h print this help message
-w=<int> number of workers
-s=<str> task scheduling policy. Choose among central, central-stack, ws, ws-de
    and numa.
-r enable worker recorder
-x=<int> task inlining limit based on num tasks per worker
-i collect worker statistics
-l=<int> worker stack size in MB
-q=<int> task queue capacity
-m=<str> memory allocation policy. Choose among coarse, fine and system.
-y=<csv> schedule policy specific parameters. Policy numa: data size in bytes
    below which task is dealt to private worker queue.
-g collect task statistics
-p enable handshake with Outline Function Profiler [Note: Supported only for a
    single worker}]

```

4.4.1 Binding workers to cores

MIR creates and binds one worker thread per core (including hardware threads) by default. Binding is based on worker identifiers — worker thread 0 is bound to core 0, worker thread 1 to core 1 and so on. The binding scheme can be changed to a specific mapping using the environment variable `MIR_WORKER_CORE_MAP`. Ensure `MIR_WORKER_EXPLICIT_BIND` is defined in `mir_defines.h` to enable explicit binding support. An example is shown below.

```

$ cd $MIR_ROOT/src
$ grep "EXPLICIT_BIND" mir_defines.h
#define MIR_WORKER_EXPLICIT_BIND
$ cat /proc/cpuinfo | grep -c Core
4
$ export MIR_WORKER_CORE_MAP="0,2,3,1"
$ cd $MIR_ROOT/programs/native/fib
$ scons
$ ./fib-debug 10 3
MIR_DBG: Starting initialization ...
MIR_DBG: Architecture set to firenze
MIR_DBG: Memory allocation policy set to system
MIR_DBG: Task scheduling policy set to central-stack
MIR_DBG: Reading worker to core map ...
MIR_DBG: Binding worker 0 to core 3
MIR_DBG: Binding worker 3 to core 0
MIR_DBG: Binding worker 2 to core 2
MIR_DBG: Worker 2 is initialized
MIR_DBG: Worker 3 is initialized
MIR_DBG: Binding worker 1 to core 1
...

```


5 Profiling

MIR supports extensive thread-based and task-based profiling.

5.1 Thread-based Profiling

Thread states and events are the main performance indicators in thread-based profiling.

Enable the `-i` flag to get basic load-balance information in a CSV file called `mir-worker-stats`.

```
$ MIR_CONF="-i" ./fib-opt
$ cat mir-worker-stats
```

TODO: Explain file contents.

MIR contains a `recorder` which produces execution traces. Use the `-r` flag to enable the recorder and get detailed state and event traces in a set of `mir-recorder-trace-*.rec` files. Each file represents a worker thread. The files can be inspected individually or combined and visualized using Paraver.

```
$ MIR_CONF="-r" ./fib-opt
$ $MIR_ROOT/scripts/profiling/thread/rec2paraver.py \
  mir-recorder-trace-config.rec
$ wxparaver mir-recorder-trace.prv
```

A set of `mir-recorder-state-time-*.rec` files are also created when `-r` is set. These files contain thread state duration information which can be accumulated for analysis without Paraver.

```
$ $MIR_ROOT/scripts/profiling/thread/get-states.sh \
  mir-recorder-state-time
$ cat accumulated-state-file.info
```

TODO: Explain file contents.

5.1.1 Enabling hardware performance counters

MIR can read hardware performance counters through PAPI during thread events.

- Install PAPI.

- Set the `PAPI_ROOT` environment variable

```
$ export PAPI_ROOT=<PAPI install path>
```

- Create a file called `HAVE_PAPI` in `MIR_ROOT/src`.

```
$ touch $MIR_ROOT/src/HAVE_PAPI
```

- Enable additional PAPI hardware performance counters by editing `MIR_ROOT/src/mir_recorder.c`.

```
$ grep -i "{PAPI_" $MIR_ROOT/src/mir_recorder.c
{"PAPI_TOT_INS", 0x0},
{"PAPI_TOT_CYC", 0x0},
/*{"PAPI_L2_DCM", 0x0},*/
/*{"PAPI_RES_STL", 0x0},*/
/*{"PAPI_L1_DCA", 0x0},*/
/*{"PAPI_L1_DCH", 0x0},*/
```

- Rebuild MIR.

```
$ scons -c && scons
```

Performance counter values will appear in the `mir-recorder-trace-*.rec` files produced by the recorder during thread-based profiling. The counter readings can either be viewed on Paraver or accumulated for analysis outside Paraver.

```
$ $MIR_ROOT/scripts/profiling/thread/get-events.sh mir-recorder-trace.prv
$ cat event-summary-*.txt
```

TODO: Explain file contents.

5.2 Task-based Profiling

Task are first-class citizens in task-based profiling.

Enable the `-g` flag to collect task statistics in a CSV file called `mir-task-stats`. Inspect the file manually or plot and visualize the fork-join task graph.

```
$ MIR_CONF="-g" ./fib-opt
$ Rscript ${MIR_ROOT}/scripts/profiling/task/plot-task-graph.R -d mir-task-
stats -c color
```

TODO: Explain file contents.

The `mir-task-stats` file can be further processed for additional information such as number of tasks and task lineage (UID for tasks).

```
$ Rscript ${MIR_ROOT}/scripts/profiling/task/process-task-stats.R -d mir-task
  -stats --lineage
$ cat mir-task-stats.info
$ head mir-task-stats.lineage
$ head mir-task-stats.processed
```

TODO: Explain file contents.

5.2.1 Instruction-level task profiling

MIR provides a Pin-based instruction profiler that traces instructions executed by tasks. Technically, the profiler traces instructions executed within outline functions of tasks in programs compiled using GCC. Follow below steps to build and use the profiler.

- Get Intel Pin sources and set environment variables.

```
$ export PIN_ROOT=<Pin source path>
$ export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:$PIN_ROOT:
  $PIN_ROOT/intel64/runtime
```

- Edit `PIN_ROOT/source/tools/Config/makefile.unix.config` and add `-fopenmp` to variables `TOOL_LDFLAGS_NOOPT` and `TOOL_CXXFLAGS_NOOPT`
- Build the profiler.

```
$ cd $MIR_ROOT/scripts/profiling/task
$ make PIN_ROOT=$PIN_ROOT
```

- View profiler options using `-h`.

```
$ $PIN_ROOT/intel64/bin/pinbin -t $MIR_ROOT/scripts/profiling/task/obj
  -intel64/mir_of_profiler.so -h -- /usr/bin/echo
...
-s specify outline functions (csv)
-c specify functions called (csv) from outline functions
-o [default mir-ofp] specify output file suffix
...
```

The profiler requires outline function names under the argument `-s`. The argument `-c` accepts names of functions which are called within

tasks. The argument `--` separates profiled program invocation from profiler arguments.

- The profiler requires handshaking with the runtime system. To enable handshaking, enable the `-p` flag in `MIR_CONF`.
- The profiler requires single-threaded execution of the profiled program. Provide `-w=1` in `MIR_CONF` while profiling.
- Information from the profiler becomes meaningful when correlated with task statistics information. Provide `-g` in `MIR_CONF` while profiling.
- Create a handy shell function for invoking the profiler and to enable task statistics collection.

```
function mir-inst-prof()
{
  "MIR_CONF='-w=1 -p -g' ${PIN_ROOT}/intel64/bin/pinbin -t ${
    MIR_ROOT}/scripts/profiling/task/obj-intel64/mir_of_profiler.so"
}
```

The profiler produces following outputs:

- Per-task instructions in a CSV file called `mir-ofp-instructions`. Example contents of the file are shown below. TODO: Add file contents. Each line shows instruction and code properties of a distinct task executed by the program. Properties are described below.
 - `task`: Identifier of the task.
 - `ins_count`: Total number of instructions executed by the task.
 - `stack_read`: Number of read accesses to the stack while executing instructions.
 - `stack_write`: Number of write accesses to the stack while executing instructions.
 - `ccr`: Computation to Communication Ratio. Indicates number of instructions executed per read or write access to memory.
 - `clr`: Computation to Load Ratio. Indicates number of instructions executed per read access to memory.
 - `mem_read`: Number of read accesses to memory (excluding stack) while executing instructions.
 - `mem_write`: Number of write accesses to memory (excluding stack) while executing instructions.

- `outl.func`: Name of the outline function of the task.
- Per-task events in a file called `mir-ofp-events`. Example contents of the file are shown below.

```
task,ins_count,[create],[wait]
14,446,[],[]
15,278,[],[]
10,60,[32,43],[47,]
```

Each line in the file shows events for a distinct task executed by the program. Event occurrence is indicated in terms of instruction count. Events currently supported are:

- **create**: Indicates when child tasks were created. Example: `[32,43]` indicates the task 10 created its first child at instruction 32 and second child at 43. Tasks 14 and 15 did not create children tasks.
- **wait**: Indicates when child tasks were synchronized. Example: `[47,]` indicates the task 10 synchronized with all children created prior at instruction 47.
- Program memory map in a file called `mir-ofp-mem-map`. This is a copy of the memory map file of the program from the `/proc` filesystem.

5.2.2 Visualization

MIR has a nice graph plotter which can transform task-based profiling data into task graphs. The generated graph can be visualized on tools such as Graphviz, yEd and Cytoscape. To plot the fork-join task graph using task statistics from the runtime system:

```
$ Rscript ${MIR_ROOT}/scripts/profiling/task/plot-task-graph.R -d mir-task-
stats.processed -p color
```

Tip: The graph plotter will plot in gray scale if `gray` is supplied instead of `color` as the palette (`-p`) argument. Critical path enumeration usually takes time. To speed up, skip critical path enumeration and calculate only its length using option `--cplengthonly`. Huge graphs with 50000+ tasks take a long time to plot. To save time, plot the task graph as a tree using option `--tree`.

The graph plotter can annotate task graph elements with performance information. Merge the instruction-level information produced by the instruction profiler with the task statistics produced by the runtime system, for the same run, into a single CSV file. Plot task graph using combined performance information.

```
$ Rscript ${MIR_ROOT}/scripts/profiling/task/process-task-stats.R -d mir-task
  -stats
$ Rscript ${MIR_ROOT}/scripts/profiling/task/merge-task-performance.R -l mir-
  task-stats.processed -r mir-ofp-instructions -k "task" -o mir-task-perf
$ Rscript ${MIR_ROOT}/scripts/profiling/task/plot-task-graph.R -d mir-task-
  perf -p color
```

5.3 Profiling Case Study: Fibonacci

The Fibonacci program is found in `MIR_ROOT/programs/native/fib`. The program takes two arguments — the number `n` and the depth cutoff for recursive task creation. Let us see how to profile the program for task-based performance information.

Compile the program for profiling — remove aggressive optimizations and disable inlining so that outline functions representing tasks are visible to the Pin-based instruction profiler. Running `scons` in the program directory builds the profiler-friendly executable called `fib-prof`.

```
$ cd $MIR_ROOT/programs/native/fib
$ scons
scons: Reading SConscript files ...
scons: done reading SConscript files.
scons: Building targets ...
scons: building associated VariantDir targets: debug-build opt-build prof-build
  verbose-build
...
gcc -o prof-build/fib.o -c -std=c99 -Wall -Werror -Wno-unused-function -
  Wno-unused-variable -Wno-unused-but-set-variable -Wno-maybe-
  uninitialized -fopenmp -DLINUX -I/home/ananya/mir-dev/src -I/home/
  ananya/mir-dev/programs/common -O2 -DNDEBUG -fno-inline-functions
  -fno-inline-functions-called-once -fno-optimize-sibling-calls -fno-omit
  -frame-pointer -g fib.c
...
gcc -o fib-prof prof-build/fib.o -L/home/ananya/mir-dev/src -lpthread -lm -
  lmir-opt
```

Tip: Look at the `SConstruct` file in `MIR_ROOT/test/fib` and build output to understand how the profiling-friendly build is done.

Identify outline functions and functions called within tasks of the `fib-prof` program using the script `of_finder.py`. The script searches for known outline function name patterns within the object files of `fib-prof`. The script lists outline functions as `OUTLINE_FUNCTIONS` and all function symbols within the object files as `CALLED_FUNCTIONS`.

```
$ cd $MIR_ROOT/programs/native/fib
$ $MIR_ROOT/scripts/profiling/task/of_finder.py prof-build/*.o
Using "..omp_fn.ol_" as outline function name pattern
Processing file: prof-build/fib.o
OUTLINE_FUNCTIONS=ol_fib_0,ol_fib_1,ol_fib_2
CALLED_FUNCTIONS=fib_seq,fib,get_usecs,main
```

Expert Tip: Ensure that `OUTLINE_FUNCTIONS` listed are those generated by GCC. Inspect the abstract syntax tree (use compilation option `-fdump-tree-optimized`) and source files.

The functions in the `CALLED_FUNCTIONS` list should be treated as functions potentially called within task contexts. Inspect program sources and exclude those which are not called within tasks. By looking at Fibonacci program sources, we can exclude `main` and `get_usecs` from `CALLED_FUNCTIONS`.

Tip: If in doubt or when sources are not available, use the entire `CALLED_FUNCTIONS` list.

Expert Tip: Identifying functions called by tasks is necessary because the instruction count of these functions are added to the calling task's instruction count.

Start the instruction profiler with appropriate arguments to profile `fib-prof`. Also collect task statistics at the same time.

```
$ mir-inst-prof \
  -s ol_fib_0,ol_fib_1,ol_fib_2 \
  -c fib,fib_seq \
  -- ./fib-prof 10 4
```

Inspect instruction profiler output.

```
$ head mir-ofp-instructions
$ head mir-ofp-events
```

Inspect task statistics.

```
$ head mir-task-stats
```

Summarize task statistics.

```
$ Rscript ${MIR_ROOT}/scripts/profiling/task/process-task-stats.R -d mir-task
  -stats
$ cat mir-task-stats.info
$ head mir-task-stats.processed
$ head mir-task-stats.lineage
```

Combine the instruction-level information produced by the instruction profiler with the task statistics produced by the runtime system into a single CSV file. Note that these files come from the same run.

```
$ Rscript ${MIR_ROOT}/scripts/profiling/task/merge-task-performance.R -l mir-
  task-stats.processed -r mir-ofp-instructions -k "task" -o mir-task-perf
```

Plot task graph using combined performance information and view on YEd.

```
$ Rscript ${MIR_ROOT}/scripts/profiling/task/plot-task-graph.R -d mir-task-
  perf -p color
$ yed task-graph.graphml
```