

DEPARTMENT OF INFORMATION TECHNOLOGY AND ELECTRICAL ENGINEERING

Spring 2024

Implementing an OpenMP runtime for MemPool in LLVM

Semester Project

Titlepage
Logo
Placeholder

Diego de los Santos Gausí dgausi@student.ethz.ch

1st of July 2024

Advisors: Samuel Riedel, sriedel@iis.ee.ethz.ch

Sergio Mazzola, smazzola@iis.ee.ethz.ch

Professor: Prof. Dr. L. Benini, lbenini@iis.ee.ethz.ch

Todo list

ref	iii
ref	iv
	14
	14
ref	14
ref	15
ref	15
ref	16
ref	18
ref	19
ref	19
	22
	25
	25
	26
	29
	30
	30
	31

ref		•	•		•	•					•		 	•			•	•	•	•	•	31
ref in footnote													 									32

Abstract

In order to effectively harness MemPool's parallelism, an initial GCC-compatible OpenMP runtime was developed for this architecture, supporting the most commonly used set of features. However, due to increasing reliance on the LLVM compiler infrastructure by the MemPool team, it is becoming crucial to have an OpenMP runtime compatible for it. This project aims to implement an OpenMP runtime that works with LLVM, supporting at least the same set of features as the previous runtime, while maintaining comparable performance and offering the possibility of future additions and improvements.

Add examples for improvements, perhaps cite the dynamic memory remapping thing

Acknowledgments

You rarely work in complete isolation, and this is the place to acknowledge contributions by others.

Declaration of Originality

Download the official declaration of originality¹, print it, and sign it. When printing your thesis, replace this sheet with the physically signed original paper. For the digital version, scan the filled declaration of originality (either before signing or remove your signature² from the image) and replace the text inside this file with \includepdf[scale=0.9]{declaration_of_originality}.

¹https://www.ethz.ch/content/dam/ethz/main/education/rechtliches-abschluesse/ leistungskontrollen/declaration-originality.pdf

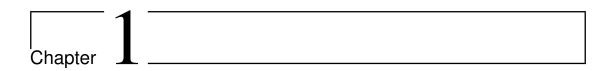
²By removing your signature from the digital version of your report, you enable us to share your report with collaborators without having to deal with privacy concerns.

Contents

1.	Intro	oductio	n	1
2.	Back	cgroun	d	3
			Pool	3
			Architecture Overview	3
	2.2.		MP	4
3.	Rela	ited Wo	ork	6
	3.1.	GCC I	mplementation	6
			Event Loop	6
		3.1.2.	Parallel Regions	8
		3.1.3.	Work Sharing Constructs	9
		3.1.4.	Barriers	11
	3.2.	HERO	LLVM Implementation	12
4.	Imp	lement	ation	13
	-		al Architecture	13
	4.2.	KMP I	Entrypoints	13
		4.2.1.	Parallel Construct	14
		4.2.2.	Work Sharing Constructs	15
		4.2.3.	Dynamic Loops	17
		4.2.4.	critical sections	19
		4.2.5.	Master and Single Constructs	20
		4.2.6.	Copyprivate Clause	21
		4.2.7.	Reduction Clause	22
		4.2.8.	Teams Construct	24
		4.2.9.	Barrier	26
		4.2.10.	Miscellaneous	26
	4.3.		lasses	27
			Task	27

Contents

		4.3.2.	Barrier	. 27									
		4.3.3.	Thread	. 29									
		4.3.4.	Team	. 33									
	4.4.	Suppo	orting Code	. 33									
			Runtime Namespace										
		4.4.2.	Heap Allocation	. 33									
		4.4.3.	Main Function Wrapper	. 33									
=	Resi	-14-0		2.4									
э.			a Campanton	34									
			ation setup										
			sections										
			arison to related work										
	5.4.	Currer	nt limitations	. 35									
6.	Con	clusion	and Future Work	36									
A.	Task	Descri	iption	37									
List of Acronyms													
List of Figures													
List of Tables													
Bil	oliog	raphy		47									



Introduction

MemPool [1] is an open-source manycore architecture developed at ETH. It is comprised of up to 1024 [2] 32-bit RISC-V Snitch [3] cores that share a common L1 cache. Unlike on GPUs, each core is individually programmable, therefore enabling the familiar shared memory programming model.

Although MemPool's parallelism can be harnessed by manually managing it using the primitives provided by MemPool's C runtime library, this can become a tedious and error prone task since the programmer has to be knowledgeable about the underlying architecture. Luckily, higher level abstraction layers such as OpenMP [4] allow the programmer to express this parallelism in a clear and concise way using the fork-join paradigm and without needing advanced architectural knowledge.

OpenMP is comprised of a set of compiler directives for both C/C++ and Fortran that enable this higher level abstraction. In order to use OpenMP, one has to use a compatible compiler that is able to roughly do the following:

- 1. The compiler must be able to parse the directives and compile the surrounding code according to their semantics.
- 2. The output of the compilation process must match the target platform in order to make use of its parallelism capabilities and primitives.

Both the GNU Compiler Collection (GCC) and LLVM provide an OpenMP runtime library interface that bridges the gap between both points: instead of directly generating a binary from OpenMP annotated code, the compiler generates calls to the OpenMP runtime library. This library has to be aware of how to parallelize the code given the platform, which makes it inherently non-portable. Naturally, popular compilers ship with their own OpenMP runtime library implementation, however, in order to get around the non-portability issue, they target a host running an operating system which is able to provide the necessary parallelism primitives. Since MemPool is designed to

1. Introduction

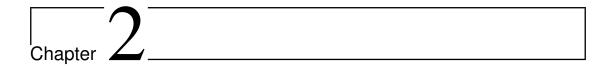
run bare-metal applications, this means that a custom OpenMP runtime library has to be developed.

Previously, the authors of MemPool developed a GCC compatible OpenMP runtime library with support for the most commonly used OpenMP constructs, such as *work sharing, critical sections* and *atomics*. However, as the MemPool team increasingly relies on the LLVM compiler infrastructure, it is essential to have a compatible OpenMP runtime library for LLVM as well.

The goal of this project is to implement an OpenMP runtime for LLVM that works on MemPool, supports at least the same set of features as the previous GCC runtime and achieves comparable performance.

The runtime will be implemented in C++ as detailed in Chapter 4. In Chapter 5, the new runtime will be evaluated in terms of performance and correctness against the GCC runtime. Finally, we will summarize the results and discuss possible shortcomings and future work in Chapter 6.

¹This implementation will be discussed in more detail in Section 3.1.



Background

2.1. MemPool

As demand for processing of highly parallel workloads grow over time, computing architectures with increasing core/processing element counts become more common. These can range all the way from custom accelerators with domain-specific architectures and programming models, to general-purpose multi-core processors. There is generally a trade-off between ease of programming and performance/efficiency.

One can achieve an interesting middle ground by grouping many relatively simple general purpose cores into a so called *compute cluster*. In practice, however, this approach does not scale well beyond a few tens of cores, which lead to the development of *multi-cluster* architectures comprised of many of such clusters. While this allows for greater core counts, it also introduces additional complexity and memory access overhead between each cluster.

MemPool [1] tries to address this issue by scaling a single cluster to the order of hundreds of cores that share a low-latency software managed L1 cache.

2.1.1. Architecture Overview

At the heart of MemPool is the *Snitch* core, which is a simple, in-order, 32-bit RISC-V core. This core was adapted for this architecture by adding the ability to retire out-of-order load instructions, which, in addition to its previously existing support for issuing multiple outstanding loads and stores through scoreboarding, allows MemPool to completely hide the L1 access latency.

MemPool's architecture is hierarchical and comprised of the following building blocks:

2. Background

Tile A tile is made up of a number of cores, some shared L1 banks, an L1 instruction cache shared amongst the tile's cores, as well as a set of remote ports to send/receive requests to/from other tiles trying to access the global L1 cache. Finally, it also has an Advanced eXtensible Interface (AXI) port connected to the system bus, which allows the tile to access the global L2 cache or system memory.

Group A group is made up of multiple tiles. Each tile's remote ports are connected with each other both within the same group as well as with other groups through a series of interconnects. Each tile's AXI port is connected to the rest of the system, also supporting Direct Memory Access (DMA).

Cluster Finally, a cluster is made up of multiple groups. Here, an AXI interface is used to connect MemPool to the rest of the system (e.g., L2 cache, system memory, etc.).

Another interesting aspect of MemPool's architecture is the *hybrid memory addressing scheme*. MemPool's shared L1 memory is interleaved at the word level across all banks. This means that accessing any given sequential region of memory would result in many strided accesses to remote tiles. This is not an issue for data shared across all cores, however, it is needlessly costly for local data (e.g., the stack). For this reason, MemPool is able to dedicate *sequential regions* of memory that are interleaved only within the same tile by re-interpreting the memory address, which is entirely done in hardware.

2.2. OpenMP

In her report [5], Qiu provides a comprehensive overview of OpenMP and its most commonly used constructs and clauses. In this work, we will implement all of the ones mentioned in her work, as well as the *teams* construct, described in the following:

The *teams* construct uses the following syntax:

```
Listing 2.1: teams construct
```

```
#pragma omp teams [clause[ [,] clause] ... ] new-line
structured-block
```

with the following clauses:

Listing 2.2: teams clauses

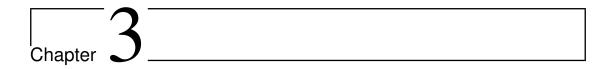
```
1 num_teams(scalar-integer-expression)
2
3 thread_limit(scalar-integer-expression)
4
```

2. Background

```
default(shared | firstprivate | private | none)
6
7
   private(list)
8
9
   firstprivate(list)
10
11
   shared(list)
12
13
    reduction([default ,] reduction-identifier : list)
14
15
   allocate([allocator :] list)
```

The *teams* construct is used to create a league of teams, where each team is a group of threads. The number of teams is less then or equal to the value defined by the *num_teams* clause, and the maximum number of threads in each team is less than or equal to the value defined by the *thread_limit* clause.

The *structured_block* is able to contain other OpenMP constructs such as *parallel*, meaning that this construct essentially allows the creation of independent groups of threads, each of which could work on different tasks. This is particularly useful in the context of MemPool, as it allows for the creation of multiple teams, each of which could be assigned to a different tile or group (or even to different clusters in a multi-cluster system).



Related Work

3.1. GCC Implementation

As previously mentioned, MemPool already had a working OpenMP runtime library for GCC. In the following, we will briefly discuss its architecture and implementation.

3.1.1. Event Loop

Because the entry point of the application is the same for all cores, it is necessary to distinguish between master and worker cores at runtime. Therefore, every application intended to use the GCC OpenMP runtime has to include the code shown on Listing 3.1 in its main file. Essentially, it distinguishes between two cases:

- 1. If the current core has ID 0 (i.e., it is the master core), it will continue executing standard OpenMP code (i.e., containing #pragma omp parallel directives, etc.).
- 2. Otherwise, it will enter an infinite loop, where the core first goes to sleep and waits for an interrupt, and then participates in the parallel execution of the program by running its assigned task. Waking up other cores and assigning tasks is done by the master core when it encounters OpenMP constructs that require it.

Listing 3.1: Main Event Loop

```
1  uint32_t core_id = mempool_get_core_id();
2
3  if (core_id == 0) {
    /* OpenMP Code */
5  } else {
    while (1) {
```

```
7  mempool_wfi();
8  run_task(core_id);
9  }
10 }
```

Listing 3.2: run_task Implementation

```
1
   typedef struct {
2
      void (*fn)(void *);
3
      void *data;
 4
      uint32_t nthreads;
5
      uint32_t barrier;
      uint8_t thread_pool[NUM_CORES];
6
7
   } event_t;
8
9
   void run_task(uint32_t core_id) {
      if (event.thread_pool[core_id]) {
10
11
        event.fn(event.data);__
        atomic_add_fetch(&event.barrier, -1, __ATOMIC_SEQ_CST);
12
13
      }
14 }
```

Listing 3.2 shows the implementation of run_task as well as event_t. Since this implementation assumes at most a single team (i.e., group of threads) working on a single task, it uses a global variable event of type event_t to store the task information.

The thread first checks if it is supposed to participate in the parallel execution by checking whether event.thread_pool[core_id] is set. If it is, it executes the function pointed to by event.fn and decrements the barrier, which makes sure that all threads have finished their work before starting the next task.

Listing 3.3: set_event Implementation

```
void set_event(void (*fn)(void *), void *data, uint32_t nthreads) {
 1
 2
      uint32_t num_cores = mempool_get_core_count();
 3
      event.fn = fn;
 4
      event.data = data;
 5
      if (nthreads == 0) {
 6
        event.nthreads = num_cores;
 7
        event.barrier = num_cores;
 8
      } else {
 9
        event.nthreads = nthreads;
10
        event.barrier = nthreads;
11
      }
12
13
      for (uint32_t i = 0; i < num_cores; i++) {</pre>
```

A task is set when the master thread calls set_event (Listing 3.3) after encountering an OpenMP construct that requires it. The function pointer and its arguments are passed through fn and data, respectively. The number of threads that should participate in the parallel execution is passed through nthreads. If nthreads is 0, all cores will execute it. This is controlled by the event.thread_pool array.

3.1.2. Parallel Regions

OpenMP parallel constructs will be transformed into calls to GOMP_parallel, which is responsible for setting the current task, waking up all cores, and waiting for them to finish. Note that this function is only called by the master thread.

Listing 3.4: GOMP_parallel Implementation

```
void GOMP_parallel_start(void (*fn)(void *), void *data,
1
2
                              unsigned int num_threads) {
      set_event(fn, data, num_threads);
3
4
     wake_up_all();
     mempool_wfi();
5
   }
6
7
   void GOMP_parallel_end(void) {
9
      uint32_t num_cores = mempool_get_core_count();
10
     while (event.barrier > 0) {
        mempool_wait(4 * num_cores);
11
12
     }
13
   }
14
15
   void GOMP_parallel(void (*fn)(void *), void *data,
16
                       unsigned int num_threads,
17
                       unsigned int flags) {
18
      uint32_t core_id = mempool_get_core_id();
19
      gomp_new_work_share();
20
      GOMP_parallel_start(fn, data, num_threads);
21
      run_task(core_id);
22
      GOMP_parallel_end();
23
```

3.1.3. Work Sharing Constructs

Work sharing constructs (e.g., for, sections, etc.) allow for threads in the current team to distribute the work to be done in a given parallel region. A global variable of type work_t is used for the bookkeeping related to these constructs and is initialized by calling gomp_new_work_share (Listing 3.4, Line 19).

Listing 3.5: work_t Struct

```
typedef struct {
 1
 2
      int end;
 3
      int next;
      int chunk_size;
 4
 5
      int incr;
 6
 7
      omp_lock_t lock;
 8
 9
      // for single construct
10
      uint32_t checkfirst;
11
      uint32_t completed;
12
      void *copyprivate;
13
14
      // for critical construct
15
      omp_lock_t critical_lock;
16
17
      // for atomic construct
18
      omp_lock_t atomic_lock;
19
   } work_t;
```

As an example for a work sharing construct, Listing 3.6 shows the implementation of the functions required for a dynamically scheduled for loop construct (i.e., #pragma omp parallel for schedule(dynamic)).

GOMP_parallel_loop_dynamic (Line 60) is very similar to GOMP_parallel, but it also initializes the loop bounds and chunk size with gomp_loop_init (Line 66).

Each thread calls GOMP_loop_dynamic_start (Line 1) once and GOMP_loop_dynamic_next (Line 36) continuously afterwards until the loop is finished. The former checks if another thread already initialized the loop bounds and chunk size, and initializes them if that is not the case. Then, it atomically increments the start of the next chunk for the next thread to use and sets the bounds of the current chunk for the current thread. This is done again with every call to GOMP_loop_dynamic_next.

3. Related Work

Listing 3.6: GOMP Dynamic For Loop Implementation

```
int GOMP_loop_dynamic_start(int start, int end, int incr,
                                 int chunk_size, int *istart,
 2
 3
                                 int *iend) {
 4
      int chunk, left;
 5
      int ret = 1;
 6
 7
      if (gomp_work_share_start()) { // work returns locked
 8
        gomp_loop_init(start, end, incr, chunk_size);
 9
10
      gomp_hal_unlock(&works.lock);
11
12
      chunk = chunk_size * incr;
13
14
      start = __atomic_fetch_add(&works.next, chunk, __ATOMIC_SEQ_CST);
15
16
      if (start >= works.end) {
17
        ret = 0;
18
19
20
      if (ret) {
21
       left = works.end - start;
22
        if (chunk > left) {
23
24
          end = works.end;
25
        } else {
26
          end = start + chunk;
27
        }
28
      }
29
30
     *istart = start;
31
      *iend = end;
32
33
      return ret;
34
35
36
   int GOMP_loop_dynamic_next(int *istart, int *iend) {
37
      int start, end, chunk, left;
38
39
      chunk = works.chunk_size * works.incr;
40
      start = __atomic_fetch_add(&works.next, chunk, __ATOMIC_SEQ_CST);
41
42
      if (start >= works.end) {
```

```
43
        return 0;
44
      }
45
46
      left = works.end - start;
47
48
      if (chunk > left) {
49
        end = works.end;
50
      } else {
        end = start + chunk;
51
52
53
54
      *istart = start;
55
      *iend = end;
56
57
      return 1;
58
   }
59
60
    void GOMP_parallel_loop_dynamic(void (*fn)(void *), void *data,
61
                                     unsigned num_threads, long start,
62
                                     long end, long incr, long chunk_size) {
63
      uint32_t core_id = mempool_get_core_id();
64
65
      gomp_new_work_share();
      gomp_loop_init(start, end, incr, chunk_size);
66
67
68
      GOMP_parallel_start(fn, data, num_threads);
69
      run_task(core_id);
70
      GOMP_parallel_end();
71
   }
72
73
   void GOMP_loop_end() {
74
      uint32_t core_id = mempool_get_core_id();
75
      mempool_barrier_gomp(core_id, event.nthreads);
76
   }
```

3.1.4. Barriers

To implement barriers, the runtime makes calls to mempool_barrier as implemented in MemPool's runtime library. It makes use of MemPool's interrupt capabilities to wake up all cores when the last core reaches the barrier. Note that cores that are not participating in the parallel execution will also be woken up, but they will immediately go back to sleep because of the check on Line 10 in Listing 3.2.

Listing 3.7: mempool_barrier Implementation

```
1
   void mempool_barrier(uint32_t num_cores) {
2
     // Increment the barrier counter
3
     if ((num_cores - 1) == __atomic_fetch_add(&barrier, 1,
 4
                                                 __ATOMIC_RELAXED)) {
        __atomic_store_n(&barrier, 0, __ATOMIC_RELAXED);
5
        __sync_synchronize(); // Full memory barrier
6
7
       wake_up_all();
8
     }
     // Some threads have not reached the barrier --> Let's wait
9
     // Clear the wake-up trigger for the last core reaching the barrier
10
11
     // as well
12
     mempool_wfi();
13
   }
```

3.2. HERO LLVM Implementation

In her work [5], Qiu also implements an OpenMP runtime library compatible with LLVM, however for the HERO [6] platform. The main difference between the GCC and LLVM implementations is that they use a different Application Programming Interface (API) (i.e., the compiler directives get translated into different function calls). The LLVM team provides a reference document [7] describing the internal runtime interface.



4.1. General Architecture

The OpenMP runtime library presented in this report is implemented in C++ and is structured into the following major components:

- KMP entrypoints.
- C++ classes for each relevant concept (e.g., threads, teams, barriers, etc.).
- Supporting code that bridges the gap between higher level concepts and the underlying MemPool architecture.

Each of them will be described in more detail in the following sections.

4.2. KMP Entrypoints

The KMP entrypoints are the functions that are called at runtime when running an OpenMP program. During the compilation process, the compiler converts OpenMP directives into calls to these functions as defined by the internal API between the LLVM compiler and the runtime library. An LLVM OpenMP runtime reference document is available at [7], however it is quite out of date and not very helpful in documenting the expected behavior of the runtime library. For this reason, using the source code of the default runtime library¹, as well as smaller third-party implementations² can be quite helpful.

¹https://github.com/llvm/llvm-project/tree/main/openmp

²https://github.com/parallel-runtimes/lomp

In the following, we will go trough each of the implemented entrypoints grouped by the OpenMP construct or clause that requires them.

4.2.1. Parallel Construct

__kmpc_fork_call

Description This function is called by the master thread of the current team (the default team contains all cores and the core with ID 0 runs the master thread). It is responsible for assigning the appropriate number of threads to the team and waking them up, as well as setting up the task that they will run.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- argc: Number of arguments passed to the microtask function.
- microtask: Function pointer to the task to be run by each thread in the team.
- ...: Variable number of arguments to be passed to the microtask function.

Implementation First, it creates a kmp::Task () object with the microtask, arguments casted to a void pointer array, and the number of such arguments. Then, it obtains the object representing the current thread using kmp::runtime::getCurrentThread() () and calls the forkCall () method on it with the task as an argument.

ref

Listing 4.1: __kmpc_fork_call

```
void __kmpc_fork_call(ident_t *loc, kmp_int32 argc,
1
2
                          kmpc_micro microtask, ...) {
3
    va_list args;
     va_start(args, microtask);
4
     kmp::Task kmpMicrotask(microtask, reinterpret_cast<void **>(args),
5
6
                             argc);
7
     kmp::runtime::getCurrentThread().forkCall(kmpMicrotask);
     va_end(args);
8
9
  };
```

__kmpc_push_num_threads

Description This function is called when using the num_threads clause in a parallel construct. It is used to request a specific number of threads to be used in the parallel section.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- global_tid: Global ID of the calling thread.
- num_threads: Requested number of threads.

Implementation This function essentially just calls the requestNumThreads method () ref on the object representing the current thread obtained by calling kmp::runtime::getThread () with global_tid.

```
Listing 4.2: __kmpc_push_num_threads
```

4.2.2. Work Sharing Constructs

The following entrypoints are used when running work sharing constructs such as static for loops or sections. Contrary to GOMP, LLVM uses the same entrypoints for both of them.

__kmpc_for_static_init_4

Description This function is called by every thread at the beginning of a static work sharing construct. It is responsible for setting the values of plastiter, plower, pupper, and pstride in order to assign a range of iterations to each thread. Because this assignment is static, the function only needs to be called once per thread. There is a variant of this function for handling the case where the loop iteration variable is unsigned. The semantics and implementation are the same except for the type of plower, pupper and plastiter, which becomes unsigned.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- gtid: Global ID of the calling thread.
- schedtype: Type of scheduling to be used. ³
- plastiter: Pointer to the *last iteration* flag. This is set to 1 if the calling thread is the one to execute the last iteration of the loop.
- plower: Pointer to the lower bound of the iteration range for the current thread.
- pupper: Pointer to the upper bound of the iteration range for the current thread.
- pstride: Pointer to the stride of the iteration range.
- incr: Increment amount of the loop iteration variable.
- chunk: Chunk size.

Implementation This function just calls the forStaticInit method () on the object representing the current team.

Listing 4.3: __kmpc_for_static_init_4

```
void __kmpc_for_static_init_4(ident_t *loc, kmp_int32 gtid,
1
2
                                  kmp_int32 schedtype,
3
                                  kmp_int32 *plastiter, kmp_int32 *plower,
4
                                  kmp_int32 *pupper, kmp_int32 *pstride,
                                  kmp_int32 incr, kmp_int32 chunk) {
5
6
     kmp::runtime::getThread(gtid).getCurrentTeam()->forStaticInit(
7
         loc, gtid, static_cast<kmp_sched_type>(schedtype), plastiter,
8
         plower, pupper, pstride, incr, chunk);
9
  };
```

__kmpc_for_static_fini

Description This function is called when the static work sharing construct is finished.

³https://github.com/llvm/llvm-project/blob/f28c006a5895fc0e329fe15fead81e37457cb1d1/ openmp/runtime/src/kmp.h#L357

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- global_tid: Global ID of the calling thread.

Implementation This function is not required for the correct implementation of the our runtime so it does nothing.⁴

4.2.3. Dynamic Loops

__kmpc_dispatch_init_4

Description This function is called once by each thread before running a dynamic for loop. It is responsible for storing the loop iteration variables and scheduling type so that they can be used by future calls to __kmpc_dispatch_next_4 (Section 4.2.3). Just like for the static loop, there is also a variant of this function that handles unsigned loop iteration variables.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- gtid: Global ID of the calling thread.
- schedtype: Type of scheduling to be used. ⁵
- lower: Lower bound of the iteration range for the current thread.
- upper: Upper bound of the iteration range for the current thread.
- incr: Increment amount of the loop iteration variable.
- chunk: Chunk size.

Implementation This function just calls the dispatchInit method () on the object refreeresenting the current team after removing any scheduling modifiers from the schedule variable using the SCHEDULE_WITHOUT_MODIFIERS macro since they are not supported by our runtime.

Listing 4.5: __kmpc_dispatch_init_4

```
1
   void __kmpc_dispatch_init_4(ident_t *loc, kmp_int32 gtid,
2
                                kmp_int32 schedtype, kmp_int32 lower,
3
                                kmp_int32 upper, kmp_int32 incr,
4
                                kmp_int32 chunk) {
5
     kmp::runtime::getThread(gtid).getCurrentTeam()->dispatchInit(
6
       loc, gtid,
7
       static_cast<kmp_sched_type>(SCHEDULE_WITHOUT_MODIFIERS(schedtype)),
8
       lower, upper, incr, chunk);
9
  }
```

__kmpc_dispatch_next_4

Description This function is called by each thread after calling dispatchInit and after completing the previously assigned chunk of iterations. It is responsible for setting the iteration variables for the current chunk of the loop that the calling thread should execute, as well as updating the loop information shared across the team so that future calling threads can know what work is left to do. There is also a variant of this function that handles unsigned loop iteration variables.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- gtid: Global ID of the calling thread.
- plastiter: Pointer to the *last iteration* flag. This is set to 1 if the calling thread is the one to execute the last iteration of the loop.
- plower: Pointer to the lower bound of the iteration range for the current thread.
- pupper: Pointer to the upper bound of the iteration range for the current thread.
- pstride: Pointer to the stride of the iteration range.

⁴Because the code generated from the OpenMP directives still calls this function, it is necessary to implement it even if it is left empty.

⁵https://github.com/llvm/llvm-project/blob/f28c006a5895fc0e329fe15fead81e37457cb1d1/ openmp/runtime/src/kmp.h#L357

Return Value 1 if there is still work to do for the calling thread, 0 otherwise.

Implementation This function just calls the dispatchNext method () on the object representing the current team.

Listing 4.6: __kmpc_dispatch_next_4

4.2.4. critical sections

```
__kmpc_critical and __kmpc_end_critical
```

Description These functions are used to implement critical sections. __kmpc_critical is called before entering a critical section and __kmpc_end_critical is called before leaving it.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- gtid: Global ID of the calling thread.
- crit: Pointer to a region of memory associated with the critical section. 8×32 bits are made available by the compiler at the location pointed to by this argument.

Implementation Both of these functions reinterpret the memory location pointed to by crit as a kmp::Mutex () object and use it to lock and unlock access to the critical section. We use a static assertion to make sure that the size of the kmp::Mutex object is smaller than the size of the memory region.

Listing 4.7: __kmpc_critical and __kmpc_end_critical

```
static_assert(sizeof(kmp::Mutex) <= sizeof(kmp_critical_name));</pre>
5
6
7
     kmp::Mutex *mutex = reinterpret_cast<kmp::Mutex *>(*crit);
8
     mutex->lock();
9
   };
10
11
   void __kmpc_end_critical(ident_t *loc, kmp_int32 gtid,
12
                              kmp_critical_name *crit) {
13
     Mutex *mutex = reinterpret_cast<kmp::Mutex *>(*crit);
14
     mutex->unlock();
15
   };
```

4.2.5. Master and Single Constructs

```
__kmpc_master and __kmpc_single
```

Description These functions are used to implement master and single constructs respectively. They are called by every thread before entering a master or single section. __kmpc_master needs to make sure that only the master thread of the team executes the section, while __kmpc_single needs to make sure that only a single thread executes it, however, it does not need to be the master thread.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- gtid: Global ID of the calling thread.

Return Value 1 if the calling thread should execute the section, 0 otherwise.

Implementation Both __kmpc_master and __kmpc_single are implemented identically for simplicity, since only letting the master thread execute single sections will always be correct, as there is only a single master thread per team. It essentially obtains the thread ID of the calling thread by calling getTid() on the thread object and compares it against 0, returning 1 if they are equal and 0 otherwise.

Listing 4.8: __kmpc_master and __kmpc_single

```
1 kmp_int32 __kmpc_master(ident_t *loc, kmp_int32 gtid) {
2   return static_cast<kmp_int32>(
3   kmp::runtime::getThread(gtid).getTid() == 0
```

```
4  );
5  };
6
7  kmp_int32 __kmpc_single(ident_t *loc, kmp_int32 gtid) {
8   return static_cast<kmp_int32>(
9   kmp::runtime::getThread(gtid).getTid() == 0
10  );
11 };
```

__kmpc_end_master and __kmpc_end_single

Description These functions are called after the execution of a master or single section respectively by the thread that executed it.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- gtid: Global ID of the calling thread.

Implementation These functions are not required for the correct implementation of our runtime so they do nothing.⁶

```
Listing 4.9: __kmpc_end_master and __kmpc_end_single
```

```
void __kmpc_end_master(ident_t *loc, kmp_int32 gtid){};

void __kmpc_end_single(ident_t *loc, kmp_int32 gtid){};
```

4.2.6. Copyprivate Clause

__kmpc_copyprivate

Description This function is used to implement the copyprivate clause associated with a single section. It is called by every thread that participates in the parallel section after the single section is executed and it is responsible for broadcasting the private data of the thread that executed it to the other threads.

⁶Because the code generated from the OpenMP directives still calls these functions, it is necessary to implement them even if they are left empty.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- gtid: Global ID of the calling thread.
- cpy_size: Size of the private data to be copied.
- cpy_data: Pointer to the private data to be copied.
- cpy_func: Pointer to a function that copies the private data.
- didit: 1 if the calling thread executed the single section, 0 otherwise.

Implementation This function just calls the copyPrivate method () on the object representing the current thread.

Listing 4.10: __kmpc_copyprivate

4.2.7. Reduction Clause

```
__kmpc_reduce and __kmpc_reduce_nowait
```

Description These functions are used to implement the reduction clause. The nowait version is called when the nowait clause is present in the reduction directive and indicates that no barrier should be executed after the reduction.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- global_tid: Global ID of the calling thread.
- num_vars: Number of variables to reduce.
- reduce_size: Size of the data to be reduced (in bytes).

- reduce_data: Pointer to the data to be reduced.
- reduce_func: Pointer to the function that performs the reduction on a pair of elements. The result is made available at the location pointed to by the first argument after calling it.
- lck: Pointer to a region of memory associated with the critical section used for the reduction. 8 × 32 bits are made available by the compiler at the location pointed to by this argument.

Return Value 1 for the master thread, 0 for others. 2 if the reduction should be performed using build in atomic operations.

Implementation Both functions are implemented identically since they only differ in the execution of a barrier, which is differentiated when calling <code>__kmpc_end_reduce_nowait</code> and <code>__kmpc_end_reduce</code> (Section 4.2.7). They always return 2 in order to use built-in atomic operations for the reduction.

Listing 4.11: __kmpc_reduce and __kmpc_reduce_nowait

```
1
   kmp_int32 __kmpc_reduce_nowait(ident_t *loc, kmp_int32 global_tid,
2
                                    kmp_int32 num_vars, size_t reduce_size,
3
                                    void *reduce_data,
                                    void (* reduce_func)(void *lhs_data,
 4
5
                                                          void *rhs_data),
6
                                    kmp_critical_name *lck) {
7
      return 2; // Atomic reduction
8
   }
9
10
   kmp_int32 __kmpc_reduce(ident_t *loc, kmp_int32 global_tid,
11
                             kmp_int32 num_vars, size_t reduce_size,
12
                             void *reduce_data,
13
                             void (*reduce_func)(void *lhs_data,
14
                                                 void *rhs_data),
15
                             kmp_critical_name *lck) {
      return __kmpc_reduce_nowait(loc, global_tid, num_vars, reduce_size,
16
                                   reduce_data, reduce_func, lck);
17
18
   }
```

__kmpc_end_reduce and __kmpc_end_reduce_nowait

Description These functions are called after the reduction is done. The first one performs a barrier, while the second one does not.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- global_tid: Global ID of the calling thread.
- lck: Pointer to a region of memory associated with the critical section used for the reduction. 8×32 bits are made available by the compiler at the location pointed to by this argument.

Implementation __kmpc_end_reduce_nowait does nothing⁷, while __kmpc_end_reduce performs a barrier using __kmpc_barrier (Section 4.2.9).

Listing 4.12: __kmpc_end_reduce and __kmpc_end_reduce_nowait

4.2.8. Teams Construct

__kmpc_fork_teams

Description This function is used to implement the teams construct. It is responsible for creating a league of teams as described in Section 2.2.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- argc: Number of arguments passed to the microtask function.
- microtask: Function pointer to the task to be run by each thread in the team.
- ...: Variable number of arguments to be passed to the microtask function.

⁷Because the code generated from the OpenMP directives still calls this function, it is necessary to implement it even if it is left empty.

Implementation This function works very similarly to __kmpc_fork_call (Section 4.2.1) except that it calls the forkTeams method () and that the microtask will only be executed by the master thread of each team.

Listing 4.13: __kmpc_fork_teams

```
void __kmpc_fork_teams(ident_t *loc, kmp_int32 argc,
1
2
                          kmpc_micro microtask, ...) {
3
    va_list args;
4
     va_start(args, microtask);
5
     kmp::Task kmpMicrotask(microtask, reinterpret_cast<void **>(args),
6
                             argc);
7
     kmp::runtime::getCurrentThread().forkTeams(kmpMicrotask);
8
     va_end(args);
9
  }
```

__kmpc_push_num_teams

Description This function is used to implement both the num_teams and thread_limit clauses of the teams construct. num_teams requests the number of teams to create and thread_limit limits the number of threads in each team.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- global_tid: Global ID of the calling thread.
- num_teams: Requested number of teams.
- num_threads: Requested maximum number of threads per team.

Implementation This function first checks if either num_teams or num_threads are greater than 0 and, if so, it sets the corresponding global variables () accordingly. The maximum number of teams is limited to half the number of cores.

Listing 4.14: __kmpc_push_num_teams

```
6  }
7
8  if (num_threads > 0) {
9   kmp::runtime::requestedThreadLimit = num_threads;
10  }
11 }
```

4.2.9. Barrier

__kmpc_barrier

Description This function is called to execute a barrier.

Arguments

- loc: Pointer to a struct containing information about the source code location of the call.
- global_tid: Global ID of the calling thread.

Implementation This function obtains the barrier object associated with the current team by calling the getBarrier method () and executes it by calling wait.

ref

Listing 4.15: __kmpc_barrier

4.2.10. Miscellaneous

__kmpc_global_thread_num

Description This function is used to obtain the global ID of the calling thread.

Arguments

• loc: Pointer to a struct containing information about the source code location of the call.

Return Value Global ID of the calling thread.

Implementation This function just returns the value returned by mempool_get_core_id, which is provided by the MemPool runtime library. It is used to obtain the current core ID, which is the same as the global thread ID since each thread runs on a single core with the same ID.

Listing 4.16: __kmpc_global_thread_num

```
1 kmp_int32 __kmpc_global_thread_num(ident_t *loc) {
2    return static_cast<kmp_int32>(mempool_get_core_id());
3 };
```

4.3. C++ Classes

4.3.1. Task

The Task class is essentially a wrapper around a KMP microtask. This microtask represents a section of code to be run by a thread (e.g., the body of a parallel section). Additionally, this class also stores the arguments that must be passed to the microtask. These are used by the compiler to access variables outside of the microtask's scope, which is why they are void pointers.

Attributes

- kmpc_micro func: Pointer to the microtask function.
- kmp_int32 argc: Number of arguments passed to the microtask function.
- void **args: Array containing the microtask arguments.

Relevant Methods

• void run(kmp_int32 gtid, kmp_int32 tid): Calls the microtask function with the arguments stored in args in addition to gtid and tid, which are always expected to be passed as the first and second argument respectively.

4.3.2. Barrier

This class is used to implement a thread barrier.

Attributes

- std::atomic<kmp_int32> barrier: Atomic counter used to keep track of the number of threads that have reached the barrier.
- std::atomic<kmp_int32> generation: Atomic counter used to distinguish between generations when using the *generation barrier*. This is not used in the *WFI barrier*.
- volatile int32_t numThreads: Number of threads participating in the barrier.

Relevant Methods

• void wait(): This function implements the actual barrier. It distinguishes between two cases: if there is currently only a single team running, then a *Wait For Interrupt* (*WFI*) *barrier* is used, otherwise a spinning *generation barrier* is used.

The first kind is essentially the same as in Section 3.1.4, with the only difference being that the barrier variable is not global, but a member of the Barrier class. This is important for the second kind of barrier, but it can also be used in this case for simplicity.

If there are multiple teams, then a *generation barrier* is used. This is because the WFI barrier does not work when multiple teams are running concurrently, since the call to wake_up_all would wake up all threads in all teams, which would interfere with other barriers. This barrier works similarly, but instead of waiting for an interrupt, the threads spin on the generation variable, which is incremented by the last arriving thread.

Listing 4.17: Barrier::wait

```
1
    inline void wait() {
 2
     if (runtime::numTeams == 1) {
 3
        // WFI barrier
 4
 5
        // Increment the barrier counter
 6
        if ((numThreads - 1) == barrier.fetch_add(1, std::memory_order\
           relaxed)) {
 7
          barrier.store(0, std::memory_order_relaxed);
 8
          std::atomic_thread_fence(std::memory_order_seq_cst);
 9
          wake_up_all();
10
        }
11
12
       // Some threads have not reached the barrier --> Let's wait
13
        // Clear the wake-up trigger for the last core reaching the \
           barrier as
```

```
14
        // well
15
        mempool_wfi();
16
17
      } else {
18
        // Spin generation barrier
19
        kmp_int32 gen = generation;
20
21
        // Increment the barrier counter
22
        if ((numThreads - 1) == barrier.fetch_add(1, std::memory_order\
           _relaxed)) {
23
          barrier.store(0, std::memory_order_relaxed);
24
          generation.fetch_add(1, std::memory_order_relaxed);
25
          std::atomic_thread_fence(std::memory_order_seq_cst);
        }
26
27
        while (gen == generation.load(std::memory_order_relaxed)) {
28
29
          // Spin
30
        }
31
      }
32
   };
```

4.3.3. Thread

This class represents a single OpenMP thread running on a core.

Attributes

- kmp_int32 gtid: Global ID of the thread.
- kmp_int32 tid: Team-local ID of the thread.
- std::optional<Task> teamsRegion: Optional task representing the body of the teams construct. This is only present in the master thread of the team.
- std::atomic<Team *> currentTeam: Pointer to the team the thread is part of.
- Mutex running: Mutex indicating whether the thread is running. It is locked if it is, unlocked otherwise. ____
- std::optional<kmp_int32> requestedNumThreads: Optional number of threads requested by a previous call to __kmpc_push_num_threads (Section 4.2.1).

Explain why this is necessary.

Relevant Methods

• void Thread::run(): This function implements the infinite loop that each thread except the thread with global ID 0 runs to wait for and execute work assigned to it.

First, the thread waits for an interrupt. After it is woken up, it acquires the running mutex and checks that it is part of a team by ensuring that currentTeam is not null, and whether it is the master thread of a team by checking if teamsRegion has a value.

If currentTeam is not null and teamsRegion is empty, meaning that the calling thread is a worker thread, the thread executes the implicit task of the team (), resets the currentTeam pointer, and executes the team's barrier associated with the end of the parallel region. The team pointer has to be reset before executing the barrier to avoid the case where a new team is assigned to the thread after the barrier but before resetting the pointer, which would cause the thread to miss the assignment to the new team.

If teamsRegion has a value, that means that the current thread is the master thread of a team and therefore has to execute the body of the teams construct. After doing that, it resets the teamsRegion pointer, deletes the team object, and executes the runtime's barrier associated with the end of the teams construct. ().

ret

Listing 4.18: Thread::run

```
1
   void Thread::run() {
 2
      while (true) {
 3
        mempool_wfi();
 4
        std::lock_guard<Mutex> lock(running);
 5
 6
        if (currentTeam != nullptr && !teamsRegion.has_value()) {
 7
 8
          (*currentTeam).getImplicitTask()->run(gtid, tid);
 9
10
          Team *prevTeam = currentTeam;
11
          currentTeam = nullptr;
12
13
          (*prevTeam).getBarrier().wait();
14
        } else if (teamsRegion.has_value()) {
15
16
          teamsRegion->run(gtid, tid);
17
18
          teamsRegion.reset();
19
20
          delete currentTeam;
21
          currentTeam = nullptr;
```

```
22
23          runtime::teamsBarrier.wait();
24      }
25     }
26     };
```

• void Thread::forkCall(Task parallelRegion): This function is called by the master thread of a team when it encouters a parallel construct.

It first checks if a specific number of threads has been requested and uses the total number of cores as the default if not. Then, it sets the number of threads for the team as well as the implicit task, which is then run on all threads in the team including the calling thread. Finally, it executes the team's barrier.

Listing 4.19: Thread::forkCall

```
void Thread::forkCall(Task parallelRegion) {
 1
 2
      kmp_int32 numThreads = this->requestedNumThreads
 3
                                      .value_or(NUM_CORES);
 4
      this->requestedNumThreads.reset();
 5
 6
      Team *team = currentTeam;
 7
 8
      // Setup
 9
      team->setNumThreads(numThreads);
10
      team->setImplicitTask(parallelRegion);
11
12
      // Run on all threads
13
      team->run();
14
      parallelRegion.run(gtid, tid);
15
16
      team->getBarrier().wait();
17
   };
```

• void Thread::forkTeams(Task teamsRegion): This function is called by the thread with global ID 0 when it encounters a teams construct.

First, it checks if a specific number of teams was requested and uses the number of groups (Section 2.1.1) as the default if not. This hints at the possibility of having different MemPool groups work on different parts of a given problem or on different problems altogether in the future. Then it sets both the global variable runtime::numTeams () and the number of threads participating in the global runtime::teamsBarrier () accordingly, and resets runtime::requestedNumTeams for the next call.

ref

ret

Next, it identifies the master threads of each team by dividing the total number of cores into runtime::numTeams equally sized chunks. Then, for each thread except for the calling thread⁸(), a new team is created with master thread coreId and team ID i. After setting the teams region and the requested thread limit, if any, the thread is woken up.

ref in footnote

The calling thread also sets its teams region and explicitly runs it after setting the requested thread limit if applicable. After executing the region, it is reset and the global teams barrier is executed, finally resetting the number of teams to 1.

Listing 4.20: void Thread::forkTeams

```
void Thread::forkTeams(Task teamsRegion) {
 2
      runtime::numTeams = runtime::requestedNumTeams
 3
                                     .value_or(NUM_GROUPS);
 4
      runtime::teamsBarrier.setNumThreads(runtime::numTeams);
 5
      runtime::requestedNumTeams.reset();
 6
 7
      kmp_int32 coresPerTeam = NUM_CORES / runtime::numTeams;
 8
 9
      for (kmp_int32 i = 1; i < runtime::numTeams; i++) {</pre>
10
        kmp_int32 coreId = i * coresPerTeam;
11
12
        Thread &thread = runtime::getThread(coreId);
13
14
        thread.setCurrentTeam(new Team(coreId, i));
15
        thread.setTeamsRegion(teamsRegion);
16
17
        if (runtime::requestedThreadLimit) {
18
          thread.requestNumThreads(
19
                  runtime::requestedThreadLimit.value());
        }
20
21
22
        thread.wakeUp();
23
24
25
      this->setTeamsRegion(teamsRegion);
26
      if (runtime::requestedThreadLimit) {
27
        this->requestNumThreads(
28
                runtime::requestedThreadLimit.value());
29
      }
30
31
      teamsRegion.run(gtid, tid);
```

⁸The thread with global ID 0 is always part of the default team, therefore it is not necessary to create a new team for it.

```
this->teamsRegion.reset();

runtime::teamsBarrier.wait();

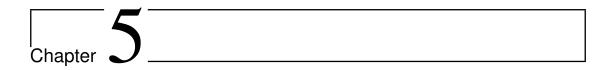
runtime::numTeams = 1;

};
```

4.3.4. Team

4.4. Supporting Code

- 4.4.1. Runtime Namespace
- 4.4.2. Heap Allocation
- 4.4.3. Main Function Wrapper



Results

- Performance against existing runtime
- Discussing correctness tests

In this chapter, you want to show that your implementation meets the objectives. Start by briefly describing the type of results you have collected and how these results relate to the objectives.

5.1. Evaluation setup

Precisely describe your measurement/evaluation setup in a top-down manner. For example:

- If you used a development platform, which one and how was it configured?
- Which tools did you use?
- What input data set / program / signals / ... did you use? If your data set is very large, you should fully define and describe it in an appendix chapter and refer to that definition here using simple labels.
- If your implementation is configurable, which configuration(s) did you evaluate?

The information given here (with references to external work) should be sufficient to reproduce the setup you used for your measurements.

5. Results

5.2. Other sections

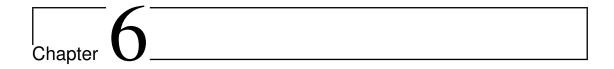
Structure the rest of this chapter into sections. For example, each section could discuss a different figure of merit or a different part of your implementation. As usual, discuss structure and content in advance with your advisors.

5.3. Comparison to related work

Compare your results to those achieved by others working on similar problems. This can be done in a separate section or directly in the sections above.

5.4. Current limitations

Demonstrating that your implementation meets the objectives is usually done by showing a lower bound on a given figure of merit. Conversely, you should also illuminate the limitations of your implementation by showing upper bounds on these (or other appropriate) figures of merit. Critically examining your ideas and their implementation trying to find their limits is also part of your work!



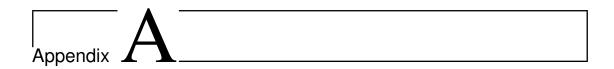
Conclusion and Future Work

- Interesting missing features (like nested parallelism and explicit tasks, as well as more clever ways to create teams) and how performance may be improved further.
- This can also be split into two chapters.

Draw your conclusions from the results and summarize your contributions. Point out aspects that need to be investigated further.

The conclusion can be structured inversely to the introduction: Summarize how the *evaluation* backs the *solution*, which solves the *problem*. Describe how your contributions improve the *situation* and what other (potentially newly discovered) problems have to be solved in the future.

Be concise: the conclusion normally fits on a single page and is rarely longer than two pages.



Task Description



TASK ASSIGNMENT FOR A SEMESTER THESIS AT THE DEPARTMENT OF INFORMATION TECHNOLOGY AND ELECTRICAL ENGINEERING

SPRING SEMESTER 2024

Diego de los Santos

Implementing an OpenMP runtime for MemPool in LLVM

March 18, 2024

Advisors: Samuel Riedel, ETZ J69.2, Tel. +41 44 632 65 69, sriedel@iis.ee.ethz.ch

Sergio Mazzola, ETZ J76.2, Tel. +41 44 632 81 49, smazzola@iis.ee.ethz.ch

Professor: Prof. Dr. L. Benini Handout: March 18, 2024 Due: July 01, 2024

The final report is to be submitted electronically. All copies remain property of the Integrated Systems Laboratory.

1 Introduction

Striving for high image quality, even on mobile devices, has led to an explosion in the pixel count of smartphone cameras over the last decade [1]. These image sensors, boasting tens of millions of pixels, create a massive amount of data to be processed on a tight power envelope as quickly as possible. Computational photography, computer vision, augmented reality, and machine learning are only a few of the possible applications. Highly specialized image signal processors (ISPs) harness the features of these workloads, which are highly parallelizable, to meet timing and power constraints. Google offers a prominent example: the Pixel Visual Core, employed for image processing in Google smartphones [2].

At ETH, we are developing our own manycore system that can be used as an ISP called MemPool [3]. It boasts 256 area-optimized 32-bit Snitch [4]. Snitch, developed at ETH as well, implements the RISC-V instruction set architecture (ISA), which is an open ISA targeting modularity and scalability [5]. Despite its size, MemPool gives all 256 cores low-latency access to the shared L1 memory, with a maximum latency of only five cycles when no contention occurs. This implements efficient communication among all cores, making MemPool suitable for various workload domains and easy to program.

A common framework for parallel programming is OpenMP [6]. It implements a task-based or fork-join approach to parallelization through annotations in the C code. The compiler will generate routines and runtime calls to distribute the parallel execution across the specific platform correctly. To this end, each system supporting OpenMP needs to implement a corresponding OpenMP Runtime. Specifically, there are two major runtimes, the GOMP and KMP runtime for GCC and LLVM, respectively. MemPool has minimal support for the GOMP runtime. However, since LLVM has become the standard compiler infrastructure for our research group, a KMP runtime to support LLVM is crucial.

Goal The goal of this thesis is to implement the KMP runtime for MemPool. To this end, the student will first identify the most important OpenMP functions for MemPool and then implement the runtime in conjunction with functional tests for verification. For benchmarking, the student will implement and evaluate digital signal processing (DSP) kernels and potentially full applications from the image processing or communications domain.

2 Milestones

The following are the milestones that we expect to achieve throughout the project:

• Become familiar with MemPool and the OpenMP runtime.

- Investigate the most useful features of the OpenMP standard.
- Implement the key features of the KMP runtime for MemPool and functionally verify them.
- Benchmark the newly implemented runtime on several DSP kernels.

2.1 Stretch Goals

Should the above milestones be reached earlier than expected and the student is motivated to do further work, we propose the following stretch goals to aim for:

- Advance the KMP implementation by implementing features beyond the basic set determined during the initial project phase.
- Build a full demonstrator application using OpenMP. The application can be chosen during the project, but examples include an HDR image processing pipeline, a ray tracing algorithm, or a 5G communication pipeline.
- Tune your benchmarks or demonstrator application for Heartstream, a 64-core chip implementing the latest MemPool architecture.

3 Project Realization

3.1 Time Schedule

The time schedule presented in Table 1 is merely a proposition; it is primarily intended as a reference and an estimation of the time required for each required step.

Project phase	Time estimate
Familiarization with MemPool and OpenMP to determine	2 weeks
key features	
Implementation and verification of key features	4 weeks
Evaluation of your work and implementation of DSP kernels	4 weeks
in OpenMP	
Write report	2 weeks
Prepare presentation	2 weeks

Table 1: Proposed time schedule and investment

3.2 Meetings

Weekly meetings will be held between the student and the assistants. The exact time and location of these meetings will be determined within the first week of the project in order to fit the student's and the assistants' schedule. These meetings will be used to evaluate the status and progress of the project. Beside these regular meetings, additional meetings can be organized to address urgent issues as well.

3.3 Weekly Reports

The student is advised, but not required, to write a weekly report at the end of each week and to send it to his advisors. The idea of the weekly report is to briefly summarize the work, progress and any findings made during the week, to plan the actions for the next week, and to bring up open questions and points. The weekly report is also an important means for the student to get a goal-oriented attitude to work.

3.4 Coding Guidelines

HDL Code Style Adapting a consistent code style is one of the most important steps in order to make your code easy to understand. If signals, processes, and modules are always named consistently, any inconsistency can be detected more easily. Moreover, if a design group shares the same naming and formatting conventions, all members immediately *feel at home* with each other's code. At IIS, we use lowRISC's style guide for SystemVerilog HDL: https://github.com/lowRISC/style-guides/.

Software Code Style We generally suggest that you use style guides or code formatters provided by the language's developers or community. For example, we recommend LLVM's or Google's code styles with clang-format for C/C++, PEP-8 and pylint for Python, and the official style guide with rustfmt for Rust.

Version Control Even in the context of a student project, keeping a precise history of changes is *essential* to a maintainable codebase. You may also need to collaborate with others, adopt their changes to existing code, or work on different versions of your code concurrently. For all of these purposes, we heavily use *Git* as a version control system at IIS. If you have no previous experience with Git, we *strongly* advise you to familiarize yourself with the basic Git workflow before you start your project.

3.5 Report

Documentation is an important and often overlooked aspect of engineering. A final report has to be completed within this project.

The common language of engineering is de facto English. Therefore, the final report of the work is preferred to be written in English.

Any form of word processing software is allowed for writing the reports, nevertheless the use of LATEX with Inkscape or any other vector drawing software (for block diagrams) is strongly encouraged by the IIS staff.

If you write the report in LaTeX, we offer an instructive, ready-to-use template, which can be forked from the Git repository at https://iis-git.ee.ethz.ch/akurth/iisreport.

Final Report The final report has to be presented at the end of the project and a digital copy needs to be handed in and remain property of the IIS. Note that this task description is part of your report and has to be attached to your final report.

3.6 Presentation

There will be a presentation (15 min presentation and 5 min Q&A) at the end of this project in order to present your results to a wider audience. The exact date will be determined towards the end of the work.

4 Deliverables

In order to complete the project successfully, the following deliverables have to be submitted at the end of the work:

- Final report incl. presentation slides
- Source code and documentation for all developed software and hardware
- Testsuites (software) and testbenches (hardware)
- Synthesis and implementation scripts, results, and reports

References

- [1] S. Skafisk, *This is How Smartphone Cameras Have Improved Over Time*, 2017 (accessed August 18, 2020). [Online]. Available: https://petapixel.com/2017/06/16/smartphone-cameras-improved-time/
- [2] J. Redgrave, A. Meixner, N. Goulding-Hotta, A. Vasilyev, and O. Shacham, "Pixel Visual Core: Google's fully programmable image, vision and AI processor for mobile devices," in 2018 IEEE Hot Chips 30 Symposium (HC30). Cupertino, US: IEEE Technical Committee on Microprocessors and Microcomputers, Aug. 2018.
- [3] S. Riedel, M. Cavalcante, R. Andri, and L. Benini, "MemPool: A scalable manycore architecture with a low-latency shared L1 memory," *IEEE Transactions on Computers*, vol. 72, no. 12, pp. 3561–3575, 2023.
- [4] F. Zaruba, F. Schuiki, T. Hoefler, and L. Benini, "Snitch: A tiny pseudo dual-issue processor for area and energy efficient execution of floating-point intensive workloads," *IEEE Trans. Comput.*, vol. 70, no. 11, pp. 1845–1860, Feb. 2021.
- [5] A. Waterman, Y. Lee, D. A. Patterson, K. Asanovic, A. Waterman, Y. Lee, and D. Patterson, "The RISC-V instruction set manual," 2014. [Online]. Available: https://github.com/riscv/riscv-isa-manual
- [6] O. A. R. Board, "OpenMP Application Programming Interface," 2021. [Online]. Available: https://www.openmp.org/wp-content/uploads/ OpenMP-API-Specification-5-2.pdf

Zurich, March 18, 2024

Prof. Dr. L. Benini

List of Acronyms

API Application Programming Interface

AXI Advanced eXtensible Interface

DMA Direct Memory Access

 $\mathsf{GCC}\,\ldots\ldots\mathsf{GNU}$ Compiler Collection

WFI Wait For Interrupt

List of Figures

List of Tables

Bibliography

- [1] S. Riedel, M. Cavalcante, R. Andri, and L. Benini, "Mempool: A scalable manycore architecture with a low-latency shared 11 memory," *IEEE Transactions on Computers*, vol. 72, no. 12, p. 3561–3575, Dec. 2023. [Online]. Available: http://dx.doi.org/10.1109/TC.2023.3307796
- [2] M. Bertuletti, S. Riedel, Y. Zhang, A. Vanelli-Coralli, and L. Benini, "Fast shared-memory barrier synchronization for a 1024-cores risc-v many-core cluster," 2023.
- [3] F. Zaruba, F. Schuiki, T. Hoefler, and L. Benini, "Snitch: A tiny pseudo dual-issue processor for area and energy efficient execution of floating-point intensive workloads," *IEEE Transactions on Computers*, vol. 70, no. 11, p. 1845–1860, Nov. 2021. [Online]. Available: http://dx.doi.org/10.1109/TC.2020.3027900
- [4] *OpenMP: The OpenMP API specification for parallel programming*. [Online]. Available: https://ww.openmp.org
- [5] W. Qiu, "Design and verification of llvm openmp runtime library on pulp/hero," 2019.
- [6] A. Kurth, A. Capotondi, P. Vogel, L. Benini, and A. Marongiu, "Hero: an open-source research platform for hw/sw exploration of heterogeneous manycore systems," in *Proceedings of the 2nd Workshop on AutotuniNg and ADaptivity AppRoaches for Energy Efficient HPC Systems*, ser. ANDARE '18. New York, NY, USA: Association for Computing Machinery, 2018. [Online]. Available: https://doi.org/10.1145/3295816.3295821
- [7] LLVM OpenMP Runtime Library Reference. [Online]. Available: https://raw.githubusercontent.com/llvm/llvm-project/main/openmp/runtime/doc/Reference.pdf