

# Detrimental task execution patterns in mainstream OpenMP<sup>®</sup> runtimes

Adam S. Tuft<sup>1</sup>[0009-0001-0251-7041],  
Tobias Weinzierl<sup>1</sup>[0000-0002-6208-1841], and  
Michael Klemm<sup>2,3</sup>[0000-0002-8634-4634]

<sup>1</sup> Department of Computer Science, Durham University, Durham, United Kingdom  
{adam.s.tuft,tobias.weinzierl}@durham.ac.uk

<sup>2</sup> Advanced Micro Devices GmbH, Dornach/Munich, Germany

<sup>3</sup> OpenMP Architecture Review Board, Beaverton, OR, USA  
michael.klemm@{amd.com,openmp.org}

**Abstract.** The OpenMP<sup>®</sup> API offers both task-based and data-parallel concepts to scientific computing. While it provides descriptive and prescriptive annotations, it is in many places deliberately unspecific how to implement its annotations. As the predominant OpenMP implementations share design rationales, they introduce “quasi-standards” how certain annotations behave. By means of a task-based astrophysical simulation code, we highlight situations where this “quasi-standard” reference behaviour introduces performance flaws. Therefore, we propose prescriptive clauses to constrain the OpenMP implementations. Simulated task traces uncover the clauses’ potential, while a discussion of their realization highlights that they would manifest in rather incremental changes to any OpenMP runtime supporting task priorities.

**Keywords:** OpenMP, Scheduling, Task-based programming

## 1 Introduction

With the advent of hundreds of cores on a contemporary computer chip in data centres, classic data parallelism reaches scalability limits. Even if we decompose algorithms into sequences of highly parallel steps, we will eventually fail to exploit the available parallelism of a machine. Task-based programming promises to ride to our rescue. From a programmer’s point of view, it imitates object-orientation’s success stories. Rather than reading an algorithm as a sequence of steps where each step exploits parallel capabilities over a large data set, we decompose an algorithm into many small “mini-algorithms” over well-defined data sets, i.e., all the data they actually read and write. The tasks can then spawn further child tasks or have inter-task dependencies.

A task logically encapsulates data plus operations on these data. While programming with tasks might reflect programming’s best practices, the HPC selling point behind tasks results from the fact that they help us to expose unprecedented scheduling freedom: task dependencies can often replace synchronization

in-between algorithmic steps. Tasks allow us to write code with a high theoretical concurrency level.

The OpenMP<sup>®</sup> API [7] has offered task directives since version 3.0, which was released in 2008. Since then, the task features have been refined and extended to provide a state-of-the-art task-parallel programming interface. The OpenMP specification does not require a specific execution mechanism or implementation strategy for OpenMP tasks. It is deliberately unspecific in several places, and therefore offers a certain degree of freedom to its implementations. A mainstream OpenMP implementation is provided with the Clang/LLVM compiler and all of our experimental data stem from this ecosystem. Though there are alternative popular implementations such as GNU’s OpenMP runtime, all share similar design rationale [6]. Whenever they yield similar execution patterns for a given task graph, our observations and concepts apply. They are generalizable.

In this paper, we study a numerical astrophysics code based upon adaptive Cartesian meshes [17] which heavily relies upon tasking [2,10]. We use it to highlight where the execution patterns from predominant OpenMP implementations are detrimental to the code’s performance. Through an artificial scheduling model, i.e., a task schedule simulator, we are able to quantify what better performance alternative schedules might be able to deliver.

Once we have introduced our demonstrator and the simulator (Sect. 2), we discuss, per runtime flaw, the degree to which it is a result of the OpenMP specification or arises from implementation decisions (Sect. 3). In Sect. 4, we propose extensions to the OpenMP tasking API which would allow an application to manipulate the task execution pattern and hence to run faster. Challenging the well-intended rationale of some OpenMP implementations, our work stands in the tradition of a transition from a descriptive to a prescriptive parallelization model. We conclude in Sect. 5 that programmers should, if they want, have a stronger say in how a task graph is actually mapped onto a task schedule.

## 2 A stationary black hole simulation analysed with Otter

We illustrate all OpenMP behaviour by means of a demonstrator from our astrophysical simulation suite ExaGRyPE [17]. It simulates a single, stationary black hole that is modelled via a first-order CCZ4 formulation [3]. Various numerical building blocks feed into this simulation, ranging from higher-order methods, adaptive mesh refinement, Sommerfeld boundary conditions, to tracer particles that allow us to evaluate global integrals over submanifolds.

This is an artificial yet numerically challenging setup [17]. To tackle the scenario, we need to simulate a large computational domain over a long time span exploiting all compute capabilities of the machine efficiently.

### 2.1 ExaHyPE’s code architecture

ExaGRyPE is a suite of solvers built on top of ExaHyPE [11] and its meshing framework Peano. Peano’s adaptive mesh refinement (AMR) is mandatory to

zoom into the area around the black hole. The arising adaptive mesh is static, i.e., it does not change over time. We use plain domain decomposition along the Peano space-filling curve to decompose the mesh for multiple processes using MPI. The same non-overlapping domain decomposition is then used once more to split up the rank-local domain and to distribute the arising chunks of the domain among the available threads. Per rank, this yields a classic fork-join parallelism. Within Peano, we map it onto an OpenMP `taskloop`. Each task traverses one subdomain on the rank, triggers all the simulation computations, and eventually synchronizes the subdomain-local data (for example, halos) with other tasks and ranks. The number of these traversal tasks is typically relatively small, as the data synchronization towards the end quickly eats up all efficiency gains if we make the subdomains too small.

---

**Algorithm 1** Pseudo code of the main traversal routine in ExaHyPE. Algorithmic steps marked with an asterisk host a `parallel for` loop.

---

```

1: function TRAVERSAL(...)
2:   #pragma omp taskloop nongroup untied
3:   for (int subdomain = 0; subdomain < K; subdomain++) do
4:     for (int cell_in_subdomain = 0; ...) do                                     ▷ Traverse local subdomain
5:       if ... then                                                               ▷ Only some cells define (produce) enclave tasks
6:         while database does not contain outcome yet do
7:           #pragma omp taskyield                                               ▷ Enclave outcomes from previous traversal
8:         end while
9:         ...
10:        #pragma omp task                                                         ▷ Spawn enclave task
11:        {
12:            ...
13:            ...                                                                     ▷ Actual work (*)
14:        }
15:        else
16:            ...                                                                     ▷ Process actual work immediately (*)
17:        end if
18:        ...
19:    end for
20:  end for
21:  #pragma omp taskwait                                                         ▷ Only wait for traversal tasks, not enclaves
22: end function

```

---

On top of the geometric data decomposition, we identify mesh cells or patches which are free of side effects [10], i.e., do not contribute towards global quantities, and are not urgent in the sense that they feed into MPI data exchange or AMR inter-resolution transfer operators. The remaining cells or patches can be spawned as separate tasks with no further in-dependencies or additional internal synchronization points. They form enclave tasks [2], which can, without a knock-on effect on MPI and the global simulation state, be executed at a later point, i.e., even after the actual traversal. The mesh traversal tasks therefore act as producers and consumers of tasks, as enclave tasks are spawned in one mesh traversal and contribute towards the solution in the subsequent mesh sweep (Alg. 1).

Several compute steps both within the traversal and within the enclave tasks exhibit further internal concurrency. This manifests as loop parallelism resulting

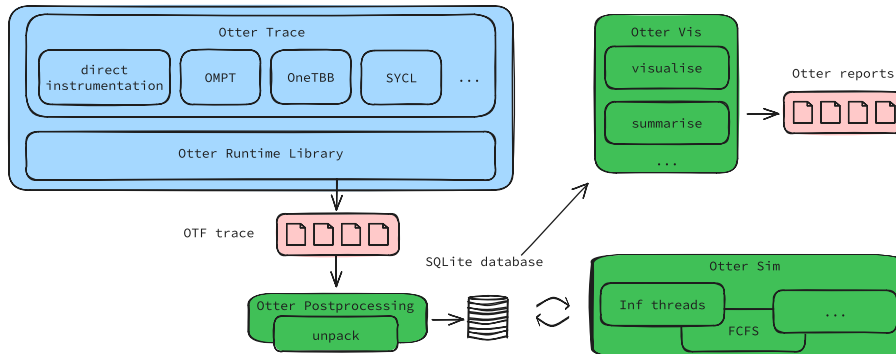


Fig. 1: The Otter tool suite and Otter’s trace-simulate-postprocess workflow.

from nested loops over different Finite Volumes or sample points of the solution that evaluate the differential equations, interpolate from one mesh resolution to another, or couple different numerical schemes.

Despite the static nature of the setup and the adaptive mesh, the compute cost per cell might change in each and every time step, as we employ non-linear equation solvers with a dynamic termination criterion per cell [3,11]. With the adaptive mesh refinement, the projections along AMR boundaries are expensive and make cells adjacent to AMR more costly than others. The cost per AMR transition depends on its orientation, i.e., whether the data layout favours a direction or not. Finally, some of our ExaGRyPE solvers switch to a subgrid model around the black hole. This yields a very high load for some mesh cells compared to others. W.l.g. we use the compute cell count as the cost metric and therefore renounce the construction of a bespoke geometric load balancing. We accept that the outermost fork-join parallelism due to traversal tasks is ill-balanced.

## 2.2 Otter tracing

To study task execution patterns, we rely on a tool suite called Otter (Fig. 1). Otter offers a macro API for annotating serial or partially parallelized code to highlight where tasks and loop parallelism could be introduced theoretically. It also can record existing OpenMP tasking through OMPT bindings. With both types of information, we run our simulation and let *Otter Trace* record the logical (hypothetical) task graph. This trace includes timing data, too. All trace data ends up in a modified OTF2 database [9].

Once postprocessed, scheduling simulators within the *Otter Sim* package allow us to re-play the recorded logical task graph within various idealized schedulers assuming infinite thread counts, infinite task queues serving all threads FCFS, different NUMA topologies, and so forth. *Otter Sim* can always (retrospectively) identify the critical path of a code as it has access to the whole execution trace. Therefore, we can make predictions of whether the execution

time would improve if manual task annotations were actually translated into OpenMP pragmas, or if alternative schedulers were available. Such statements are optimistic: For the present studies, we rely on FCFS scheduling with a global task queue. We ignore NUMA effects as well as the critical path analysis, and we also neglect further external factors such as bandwidth constraints or task activation latencies.

*Otter Vis* finally translates both the traced and hypothetical data into HTML reports containing runtime metrics, graphs and figures. It allows us to compare the recorded execution pattern of a code to different simulated, hypothetical traces. Such postprocessed data can guide the parallelization of a code, but also uncovers in hindsight unfortunate scheduling decisions from a real run.

### 3 Execution patterns

For our demonstrator runs, a standard 16-core AMD EPYC™ Processor model 7302 serves as testbed, although we intentionally limit the number of available OpenMP threads to four. This helps us to highlight execution patterns of interest. With a larger number of threads, the resulting data can become too multifaceted, obscuring details. In all of our experiments, we have validated that the tracing induces negligible runtime overhead of less than 2%. We can trust the tracing data.

#### 3.1 Task spawn guarantees

Task creation in the OpenMP API via a `task` directive introduces a task and also constitutes a Task Scheduling Point (TSP). At this point, it is at the OpenMP implementation’s discretion either to execute the task immediately “in situ” (undelayed in the same thread) or to actually spawn it as a deferred task. The latter sends the task to the task pool, from which it is picked up later; potentially by another thread. While programmers can push the behaviour towards undelayed execution via `if` and `final` clauses, they can not enforce a deferred task.

This freedom is intended by the OpenMP specification. It allows an implementation to easily deal with large numbers of tasks by switching between undelayed and deferred execution as needed, depending on runtime conditions such as number of available threads, current load of the system, etc. It allows for a throttling of the task creation [1,4].

As an implementation example, LLVM’s OpenMP runtime maintains a double-ended task queue per thread. Each thread enqueues created tasks in its own task queue, and always tries to pick a task for execution from the end. This leads to an effective last-in first-out execution behaviour of tasks. If there are no more ready tasks left in the thread-local queue, the thread attempts to steal tasks from queues of other threads. Stolen tasks are taken from the beginning of a queue to retrieve the longest-waiting tasks first.

In ExaHyPE, a traversal task spawns bursts of enclave tasks (Alg. 1). Switching to undelayed mode implies that these enclave task bursts might be artificially

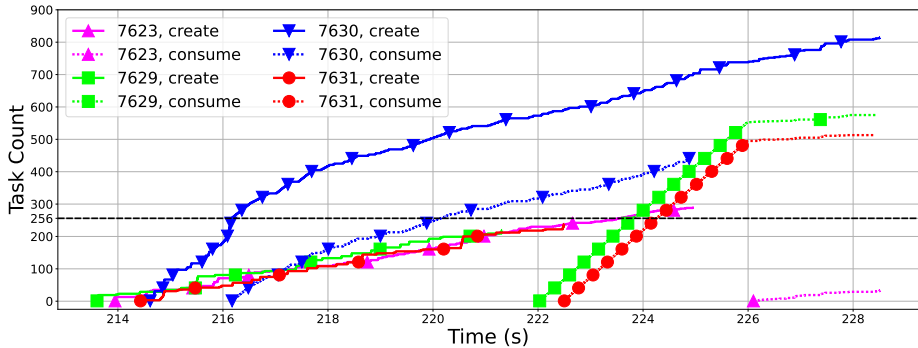


Fig. 2: Created vs. consumed enclave tasks on four threads (7623, 7629, 7630, 7631) over time for a single timestep. 7623, 7629 and 7631 produce tasks which are held in a task queue and then completed, i.e. consumed once the producing task has terminated. 7630 produces so many tasks that the producer task is suspended as further child tasks cannot be deferred. They are executed immediately. From 222s onwards, 7629 and 7631 start to process tasks, eventually steal from 7630 and hence allow 7630 eventually to stop interrupting the producer.

constrained. Traversal tasks terminating early due to geometric load imbalances consume the tasks they have spawned before they continue to steal enclaves from other threads’ queues. Traversal tasks running longer and producing many tasks see their tasks being stolen by otherwise idle threads. To enable this behaviour is the intention and motivation behind our enclave design [2]. Traversal tasks spawning a very high number of tasks might run into a situation where they exceed their task queue size, while no other threads are available to steal their tasks. They stop further task production, and instead immediately process child tasks [14]. Our testbed software stack seems to employ a queue threshold of 256 tasks. Once a thread has enqueued more than 256 tasks, the system switches into the undeferred mode (Figure 2).

If tasks producing the lion’s share of enclave tasks are also the critical tasks within their fork-join section, the switch to an immediate consumption introduces a bottleneck, as a failure to defer these tasks prolongs the critical path (Figure 3). With *Otter Sim*, we can simulate a world where the task queues have no upper threshold, i.e., tasks are always deferred. This would reduce the runtime by up to 4.7 %. There are enough threads available towards the end of each traversal to consume all deferred tasks spawned by the critical path.

There are two workarounds to enforce this behaviour: first, we can introduce helper queues on top of the actual OpenMP runtime [10,14]. Rather than spawning OpenMP tasks directly, we hold them back in a user-defined queue, releasing them only after the production task has terminated. While there are reasons for this approach besides the manual deferring—it allows us to fuse withheld tasks into one meta task to deploy them en bloc to a GPU [16] or to vectorize aggressively [10]—it replicates OpenMP core functionality. The alter-

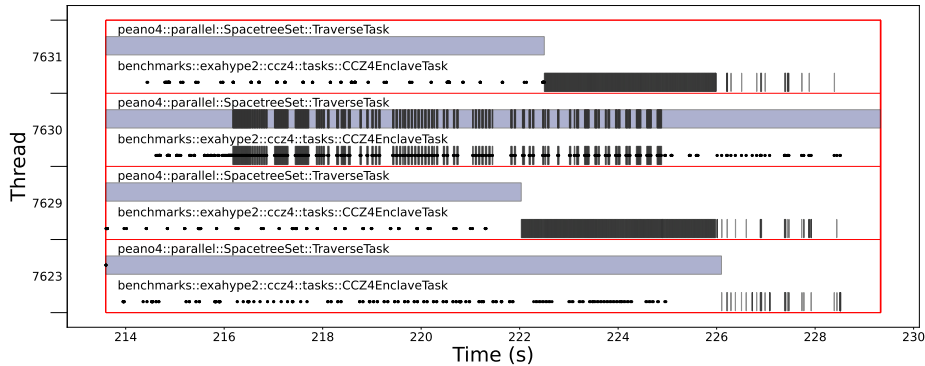


Fig. 3: Trace for task execution pattern from Figure 2 with grey bars illustrating the traversal tasks, dots showing the creation of enclave child tasks, and narrow black bars denoting the actual enclave task execution. Bars embedded in the traversal task on thread 7630 show the task being suspended when the thread immediately executes an enclave task. All data are recorded, i.e., not simulated. The trace illustrates that the traversal task of thread 7630 is on the critical path. Not deferring child tasks prolongs this path.

native second workaround instructs OpenMP to increase its task queues upon demand (`KMP_ENABLE_TASK_THROTTLING=0`). Yet, there are good reasons for limited queue sizes. If we increase them dramatically, we have to pay a runtime and memory overhead penalty. Our demonstrator selectively identifies tasks which benefit from flexible queue sizes without asking for globally dynamic queues.

### 3.2 Nested parallelism

High performance computing codes tend to combine task and data parallelism. Modern codes also tend to rely on hierarchical parallelism [12,15]. Such a code makes the concurrency fan out as the code descends along the call tree.

In the OpenMP API, the maximum number of nested parallelization levels is controlled through `OMP_MAX_ACTIVE_LEVELS` [7]. By default, OpenMP realizes strictly nested parallelism, i.e., our data-parallel region mapped onto a `parallel for` can only use a subset of the threads available to the enclosing section. A task that contains a `for` directive will therefore only utilize one thread to execute that region, as each task is tied to one thread.

In ExaHyPE, the subdomain traversals are realized as tasks. The domain decomposition geometrically makes some tasks responsible for the interpolation and restriction along AMR boundaries and/or the coupling of one physical model to the other. Both types of operations are very expensive. Logically, the arising projections between different solvers or mesh resolutions are embarrassingly parallel, i.e., could be mapped onto an embedded `parallel for`. Yet, we serialize any `parallel for` embedded into a task.

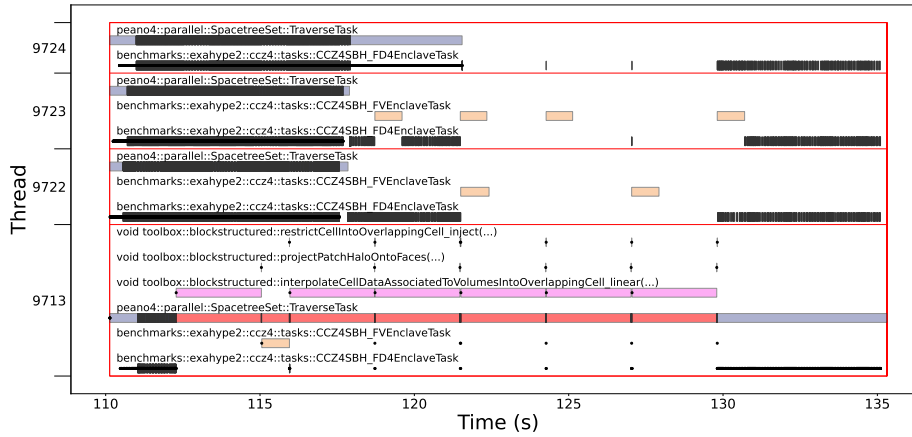


Fig. 4: Timeline of a single time step. The red bars highlight where the critical task runs into some embedded `parallel for` constructs. Recorded timings augmented by postprocessing data (critical path).

As the loops are expensive and as they often align along the critical path—in many ExaHyPE settings they are indeed responsible for making the owning task a member of the critical path—the lack of support for nested parallelism within tasks increases the makespan of the application (Fig. 4). For characteristic demonstrator setups, we could, in theory, gain up to a 43 % reduction in runtime if the loops were executed in parallel. This optimistic estimate assumes that sufficient idling threads are available whenever a critical task encounters a loop.

The OpenMP API provides no mechanism that allows a task to book multiple threads in a data-parallel fashion. This leaves us with two alternatives. On the one hand, we could use a `parallel for`, i.e., fire up a new team of threads. On the other hand, we could switch from a `for` loop to a `taskloop`. In an HPC context, a new team of threads is problematic. Typically, HPC codes spawn one thread per core right from the start to avoid oversubscription, i.e. hyperthreading resulting in overheads due to thread swapping. Mechanisms to reuse existing threads within nested regions are available (cmp. `KMP_HOT_TEAMS_MAX_LEVEL`), but do not address the present issue, if the threads are already booked out. Using a `taskloop` construct is hence more natural and, in combination with a `default(shared)` clause, a minimally invasive change to the code. Yet, creating many subtasks through a `taskloop` introduces overhead in itself, while the tasks might end up in the thread-local queue, i.e., continue to be serialized.

### 3.3 Fair yields

In ExaHyPE, there are typically few traversal tasks compared to the available core count. In our experiments, the traversal tasks therefore almost never wait for the enclave tasks. All enclave tasks are completed by threads not involved in the



traversal ahead of time. There are exceptions to the rule: NUMA-intensive systems such as the AMD EPYC™ processors make some developers deploy one MPI process per NUMA domain—effectively yielding a low core count per process—and dynamic load balancing might deploy more subpartitions to an MPI process than there are threads available. Our implementation therefore protects the access to enclave task outcomes with an atomic flag signalling whether the task is complete. If unset, the consumer thread yields and then polls again.

The result code starves in rare situations. Let  $T$  be the number of threads available and let  $C > T + 1$  be the number of consumer tasks that traverse the domain and consume one of the  $E$  enclave task outcomes which are pending in the system.  $E + C \gg T$  then is the total number of tasks. ExaHyPE can run into a situation where a consumer task yields, another untied consumer is swapped in, and the consumers all take turns checking the enclave task flags. The enclave tasks then starve.

In the OpenMP API specification, the `taskyield` construct is a hint to the implementation to introduce an additional TSP and, hence, to give other tasks a chance to be scheduled for execution. As a hint, an implementation can ignore the TSP or implement it in various ways with different performance implications [13]. For ExaHyPE, it is problematic that OpenMP’s `taskyield` does not provide a fairness guarantee and notably cannot be used to drain a pending task queue incrementally. One might argue that a robust realization of our producer-consumer pattern should employ task dependencies to avoid the polling for task readiness. Yet, we note that `taskwait` with a `depend` clause is not always straightforward to use from a programmer’s point of view, as all task dependency addresses have to reside within the user space. More severely, task dependencies are restricted to sibling tasks in OpenMP whereas we have parent–child relations here. ExaHyPE would benefit from a fair yielding mechanism to avoid occasional deadlocks on some hardware. At the moment, we have to manually work around such cases by adding user-defined task queues.

### 3.4 Taskwait semantics

The `taskwait` and `taskloop` directives used by our traversal tasks introduce synchronization points. We use `taskwait` with `taskloop nogroup` to synchronize the set of all child tasks with the same parent task, i.e., all traversal siblings, but not their spawned enclave tasks, while the synchronization set for the `taskloop` without `nogroup` includes all descendant tasks created from within the `taskloop` region, i.e., all enclave tasks produced, too. The `taskloop nogroup` can be replaced by a for loop spawning the traversal tasks individually, which is necessary for the NVIDIA software stack that lacks support for task groups. If supported, we find it to be more elegant and slightly faster than manual task spawning.

It is up to the OpenMP implementation how the synchronization points are realized. There are two basic options. First, the runtime can process further tasks while waiting until all tasks in the synchronization set have completed execution. Second, it can actively poll the synchronization construct, possibly deciding not to execute further tasks while it waits.

Let ExaHyPE spawn  $K$  (traversal) tasks at one point, and immediately after the end of that task group issue another one with  $K$  tasks. These are two time steps. If  $K - 1$  threads decide to process further (enclave) tasks at the first synchronization point, the final remaining thread might hit the end of the first task group while these threads still are busy. Only this one thread hence is available to immediately continue with the  $K$  tasks from the second task group.

While it makes sense for a thread to execute further tasks while waiting in a `taskwait` or at the end of a `taskgroup` region—this guarantees progress—it means that we add algorithmic latency to the second task group in the example. This latency is defined by the time it takes the  $K - 1$  threads to finish the currently active enclave task, and to join traversal tasks of the next task group, i.e., time step.

In ExaHyPE, the traversal tasks define the critical path. If one of the traversal tasks is not immediately kicked off at the start of a task group, we risk delaying the critical path. Furthermore, we observe that such latency can lead to low occupancy further down the road, as tasks have been processed at a scheduling point where it would have been better from a performance standpoint if the underlying thread had paused for a moment and then continued with the traversal task [14]. Yet, OpenMP provides no mechanism to stop a thread at the end of a task group from continuing with other tasks. We have no mechanism to flag to the system that we are aware that there are many ready tasks but that we also know that there will be a point reached soon with a low concurrency level where these tasks can all be handled without delaying any other time-critical task. In such a case, we might be willing to accept low occupancy temporarily, as long as we can trade this to an immediate continuation along the critical path—knowing that there will be enough resources later on to handle all the postponed tasks.

## 4 OpenMP extensions and their realization

---

**Algorithm 2** Domain traversal loop with modified OpenMP annotations.

---

```

function TRAVERSAL(...)
  #pragma omp ... nogroup latency
  for (int subdomain = 0; ...) do
    ...
    while ... do
      #pragma omp taskyield throughput
    end while
    ...
    #pragma omp task defer
    ...
  end for
end function

```

---

With a clear description of runtime flaws, we can propose some modifications to the OpenMP API that would help our demonstrator code. We distinguish between proposals for the specification API and suggestions how to implement an altered specification (Algorithm 2).

## 4.1 API modifications

ExaHyPE with its producer-consumer pattern would benefit from explicitly labelling tasks as “must be deferred”, This would complement the existing semantics of the OpenMP clauses `if` and `final`, i.e., allow developers to disable task throttling [1,4,5]. In ExaHyPE, our first and foremost goal is the reduction of the critical path requiring such a flag to be prescriptive. Yet, realization constraints and task pool overflows might require it to become a weakly prescriptive annotation (see below). A possible extension of the OpenMP API would be a new clause `defer` that extends the existing clauses of the `task` directive: `#pragma omp task defer(condition)`. If `condition` evaluates to `true`, the task shall be deferred; if it evaluates to `false`, the task maybe undeferred or deferred.

ExaHyPE would benefit from the introduction of tasks that roll out embedded loops over multiple threads. Providing such a feature with spreading guarantees is difficult [15]. A plain `taskloop` with the clause `priority(omp_get_max_task_priority())` would not facilitate the feature. It would assign the resulting loop chunks high priority compared to any other task in the system and load stealing would implicitly scatter the iteration range among the available OpenMP threads. However, there is no guarantee that they are stolen. We would require a `scatter` clause.

For a `taskyield` variant, we would envisage that programmers should be able to decide that a `taskyield` region should not be an no-op, but actually pick tasks from the task pool for execution. Furthermore, it would help to mark a TSP of a `taskyield` region as either high throughput or low latency: `#pragma omp taskyield latency|throughput`, with the default being the current implementation-defined behaviour.

A low latency TSP suspends the encountering task yet brings it back as soon as possible to minimize the probability that we lose all of its cache content, following the depth-first philosophy of OpenMP implementations. It reduces the algorithmic latency of polling realized through yield. In our case, we would rather use a throughput-oriented TSP which would cause a yielding task to go to the back of a queue. Such a yield could then also provide fairness guarantees.

To facilitate low latency `taskwait` or `taskgroup` regions, the above `latency` and `throughput` clauses would also be added to these directives. It instructs the implementation that threads hitting the corresponding synchronization point are to be kept free of other tasks, as they will be needed immediately afterwards for some high priority work (`latency`). While a suspended task is waiting for tasks to complete, the implementation shall not schedule other tasks to reduce wakeup latency. It is up to the user to guarantee that this clause does not induce a deadlock or starvation, for example by having all tasks encounter a latency optimized `taskwait latency` directive.

## 4.2 Realization

The opportunity to manually defer tasks to the task pool independent of the execution context means that we have to provide a dynamic task queue which

can grow without any constraints as a realization of an unbounded task pool. Otherwise, task pools might overflow. A weakened realization sticks to the existing implementation of task queues, but instead switches from task stealing to task distribution for “must be deferred” tasks: Whenever a thread spawns more tasks than its local task queue is able to accommodate, these tasks first are scattered over other threads’ queues (similar to how tasks with data affinity are distributed [8]). If and only if this task distribution fails as well, we fall back to undelayed execution for the created task. We hypothesize that this last variant provides a reasonable compromise between runtime efficiency, small changes to existing infrastructure, and improved runtime characteristics on the demonstrator side. To avoid that an active task distribution confuses the scheduling of victim threads, it is important that the created subtasks are assigned a very low priority, i.e., are not brought forward on the target thread. Otherwise, we would replicate the motivating problem once again.

To allow parallel loops that are embedded into a task to “invade” other threads, an implementation should go the other way regarding task priorities. The loop segments scattered over the queues have to have a higher priority than the highest priority task in any respective queue, and task stealing has to take priorities into account, too. This way, we can ensure that idle tasks steal the “right” tasks from the thread issuing the parallel for, and that the stealing does not delay the actual execution of the parallel loop that kicked it off.

A fair yield which can guarantee progress in our particular case could easily be mapped onto priorities, too: If a thread queue is aware of the lowest priority task, a throughput yield would label the suspended task with a priority that is by one smaller than the currently lowest priority.

The new “low latency” clause finally requires us to eliminate the task scheduling point at the end of the loop. This will let the used threads (logically) run idle. Once we ensure that a subsequent parallel loop schedules “its own” tasks or parallel segments first, we however obtain a taskwait which prioritizes low algorithmic latency over throughput. To ensure that the subsequent loop start does not issue any other task first, we can either hardcode the scheduling or make the scheduling biased. Given the spawned task higher priorities than all other pending ready tasks introduces the required bias.

### 4.3 Contextualization

There are many valid reasons and rationale why OpenMP implementations fall back to task throttling if too many tasks are created [5]. Our suggestion to introduce a dedicated new clause to avoid this accepts this fact, as it suggests localized modifications to few tasks. Such a clause should have no detrimental effect on existing codes.

Our nested parallelism within tasks aligns with existing developments within OpenMP and does not sacrifice threads [15]. It notably fits to OpenMP’s GPU kernel concept, where the `target` pragma lets the spawning code fan out into a new team with many threads. With a full support of nested parallelism on the host, massive tasks that should be offloaded to a GPU yet cannot be moved

there for one reason or the other benefit from the full concurrency of the host processor. The feature facilitates platform- and performance-portable code.

A fair yield is a natural cousin to the existing `detach` clause, which becomes useful if we cannot easily construct or realize a completion check and instead prefer (partial) draining of the task queue.

A “do not schedule” policy at an implicit synchronization point stands in the tradition of OpenMP to grant NUMA considerations high priority. The probability is high that people use it for subsequent parallel loops with the same granularity which run over related data.

## 5 Conclusion and outlook

OpenMP schedules are often not unique or enforced by the standard. We introduce four scenarios, where the LLVM implementation introduces runtime flaws: tasks are prematurely activated, tasks do not support embedded parallelism, tasks do not yield in a fair way, and waiting constructs always prioritize high throughput instead of low algorithmic latency.

These flaws result from common interpretations and rationale how to interpret and realise the standard efficiently for a magnitude of applications. Our work does not challenge the underlying implementation rationale of mainstream runtimes—indeed there are good reasons to implement things the way they are—but it suggests that users should be allowed to explicitly instruct OpenMP to realize things differently. While prescriptive OpenMP statements already enforce certain OpenMP behaviour, our work goes one step further and makes the prescriptive character cover certain realization decisions, too.

These modifications do not require major rewrites of the OpenMP runtime. Instead, the majority of them can be implemented combining task priorities with minor changes in the runtime’s logic. For all changes a mature implementations of priorities is a sine qua non which can induce further scalability challenges on massively parallel systems. In combination with few further tweaks, they however will provide multifaceted tuning opportunities to codes like ExaHyPE.

GNU and other runtimes share implementation rationale with LLVM. We may therefore expect that many of the documented flaws arise there, too, likely with quantitatively different characteristics. It is future work to assess these differences systematically. Our work uses one bespoke simulation code as demonstrator. We again expect other task-heavy codes to encounter similar flaws and, hence, to benefit from the proposed extensions. The scientific challenge for future work is to quantify these effects, but also to identify if the extensions and required modifications of the runtime could potentially harm the performance of other codes. They have the potential to make runtime implementations not backward compatible from a performance point of view.

## Acknowledgments

Tobias’ research has been supported by EPSRC’s Excalibur programme through its cross-cutting project EX20-9 *Exposing Parallelism: Task Parallelism* (Grant ESA 10 CDEL) and the DDWG projects *PAX-HPC* (Grant EP/W026775/1) as well as *An ExCALIBUR Multigrid Solver Toolbox for ExaHyPE* (EP/X019497/1). His group appreciates the support by Intel’s Academic Centre of Excellence at Durham University. The comparison of OpenMP vs. TBB and the assessment of early oneAPI OpenMP behaviour has led to some of the investigations reported here. This work has made use of the Hamilton HPC Service of Durham University.

AMD, the AMD Arrow logo, EPYC, and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other product names used in this publication are for identification purposes only and may be trademarks of their respective companies. ExaHyPE<sup>4</sup> and Otter<sup>5 6</sup> are open source.

## References

1. Spiros N. Agathos, Nikolaos D. Kallimanis, and Vassilios V. Dimakopoulos. Speeding Up OpenMP Tasking. In Christos Kaklamanis, Theodore Papatheodorou, and Paul G. Spirakis, editors, *Euro-Par 2012 Parallel Processing*, pages 650–661. Springer, 2012.
2. Dominic Etienne Charrier, Benjamin Hazelwood, and Tobias Weinzierl. Enclave Tasking for DG Methods on Dynamically Adaptive Meshes. *SIAM Journal on Scientific Computing*, 42(3):C69–C96, 2020.
3. Michael Dumbser, Federico Guercilena, Sven Köppel, Luciano Rezzolla, and Olindo Zanotti. Conformal and covariant Z4 formulation of the Einstein equations: Strongly hyperbolic first-order reduction and solution with discontinuous Galerkin schemes. *Phys. Rev. D*, 97:084053, 2018.
4. Alejandro Duran, Julita Corbalan, and Eduard Ayguade. An adaptive cut-off for task parallelism. In *SC ’08: Proceedings of the 2008 ACM/IEEE Conference on Supercomputing*, pages 1–11, 2008.
5. Thierry Gautier, Christian Pérez, and Jérôme Richard. On the Impact of OpenMP Task Granularity. In Bronis R. de Supinski, Pedro Valero-Lara, Xavier Martorell, Sergi Mateo Bellido, and Jesus Labarta, editors, *Evolving OpenMP for Evolving Architectures*, volume 11128 of *Lecture Notes in Computer Science*, pages 205–221, 2018.
6. Michael Klemm and Jim Cownie. *High Performance Parallel Runtimes: Design and Implementation*. De Gruyter, 2021.
7. Michael Klemm and Bronis R. de Supinski, editors. *OpenMP Application Programming Interface Specification Version 5.2*. OpenMP Architecture Review Board, 2021.
8. Jannis Klinkenberg, Philipp Samfass, Christian Terboven, Alejandro Duran, Michael Klemm, Xavier Teruel, Sergi Mateo, Stephen L. Olivier, and Matthias S.

---

<sup>4</sup> <https://gitlab.lrz.de/hpcsoftware/Peano/-/releases/2024OpenMPPaper>

<sup>5</sup> <https://github.com/Otter-Taskification/otter/releases/tag/2024-openmp-paper>

<sup>6</sup> <https://github.com/Otter-Taskification/pyotter/releases/tag/2024-openmp-paper>

- Müller. Assessing Task-to-Data Affinity in LLVM OpenMP. In Bronis R. de Supinski, Pedro Valero-Lara, Xavier Martorell, Sergi Mateo Bellido, and Jesus Labarta, editors, *Evolving OpenMP for Evolving Architectures*, volume 11128 of *Lecture Notes in Computer Science*, pages 236–251, 2018.
9. Andreas Knüpfer, Christian Feld, Dieter Mey, Scott Biersdorff, Kai Diethelm, Dominic Eschweiler, Markus Geimer, Michael Gerndt, Daniel Lorenz, Allen Malony, Wolfgang Nagel, Yury Oleynik, Peter Philippen, Pavel Saviankou, Dirk Schmidl, Sameer Shende, Ronny Tschüter, Michael Wagner, Bert Wesarg, and Felix Wolf. *Score-P: A Joint Performance Measurement Run-Time Infrastructure for Periscope, Scalasca, TAU, and Vampir*, pages 79–91. Springer, 2012.
  10. Baojiu Li, Holger Schulz, Tobias Weinzierl, and Han Zhang. Dynamic task fusion for a block-structured finite volume solver over a dynamically adaptive mesh with local time stepping. In *ISC High Performance 2022*, volume 13289 of *Lecture Notes in Computer Science*, pages 153–173, 2022.
  11. Anne Reinartz, Dominic E. Charrier, Michael Bader, Luke Bovard, Michael Dumbser, Kenneth Duru, Francesco Fambri, Alice-Agnes Gabriel, Jean-Matthieu Gallard, Sven Kppel, Lukas Krenz, Leonhard Rannabauer, Luciano Rezzolla, Philipp Samfass, Maurizio Tavelli, and Tobias Weinzierl. ExaHyPE: An engine for parallel dynamically adaptive simulations of wave problems. *Computer Physics Communications*, 254:107251, 2020.
  12. Sara Royuela, Maria A Serrano, Marta Garcia-Gasulla, Sergi Mateo Bellido, Jesús Labarta, and Eduardo Quiñones. The cooperative parallel: A discussion about runtime schedulers for nested parallelism. In Xing Fan, Bronis R. de Supinski, Oliver Sinnen, and Nasser Giacaman, editors, *OpenMP: Conquering the Full Hardware Spectrum*, volume 11718 of *Lecture Notes in Computer Science*, pages 171–185, 2019.
  13. Joseph Schuchart, Keisuke Tsugane, José Gracia, and Mitsuhsa Sato. The Impact of Taskyield on the Design of Tasks Communicating Through MPI. In Bronis R. de Supinski, Pedro Valero-Lara, Xavier Martorell, Sergi Mateo Bellido, and Jesus Labarta, editors, *Evolving OpenMP for Evolving Architectures*, volume 11128 of *Lecture Notes in Computer Science*, pages 3–17, 2018.
  14. Holger Schulz, Gonzalo Brito Gadeschi, Oleksandr Rudyy, and Tobias Weinzierl. Task inefficiency patterns for a wave equation solver. In Simon McIntosh-Smith, Bronis R. de Supinski, and Jannis Klinkenberg, editors, *OpenMP: Enabling Massive Node-Level Parallelism*, pages 111–124, Cham, 2021. Springer International Publishing.
  15. Jinghao Sun, Nan Guan, Feng Li, Huimin Gao, Chang Shi, and Wang Yi. Real-Time Scheduling and Analysis of OpenMP DAG Tasks Supporting Nested Parallelism. *IEEE Transactions on Computers*, 69(9):1335 – 1348, 2020.
  16. Mario Wille, Tobias Weinzierl, Gonzalo Brito Gadeschi, and Michael Bader. Efficient GPU Offloading with OpenMP for a Hyperbolic Finite Volume Solver on Dynamically Adaptive Meshes. In *ISC High Performance 2023*, volume 13948 of *Lecture Notes in Computer Science*, pages 65–85, 2023.
  17. Han Zhang, Christian Barrera-Hinojosa, Baojiu Li, and Tobias Weinzierl. ExaGRyPE: Numerical General Relativity Solvers Based upon the Hyperbolic PDEs Solver Engine ExaHyPE, 2024. (to be released on arXiv within the next few weeks).