

# OMPRTF

Brady Bangasser

Modern HPC Systems and Programs alike are increasingly reliant on aggressive compiler optimizations to speed up programs, often merely to obtain a realistically duration runtime. Most often, these aggressive techniques rely on static analysis, which, when given simple, small, and prolific programs, yields highly optimized machine code. In other cases, where programs are large or contain highly complex or runtime-dependent logic, static analysis becomes infeasible or logically impossible. The solution to this is dynamic analysis, which involves executing a program under varied conditions and parameters, and using the resulting timings and data to apply extremely aggressive optimization. OpenMP RunTime Fabric (OMPRTF) uses this technique to apply conditional optimizations to programs that utilize OpenMP's GPU Offloading functionality.

OMPRTF is broken into multiple stages for simplicity, modularity, and ease of concurrent and pipelined execution. The first of these stages is compiling the user's software into LLVM Intermediate Representation (IR). LLVM IR was chosen as the input representation for OMPRTF primarily because of the flexibility it offers with input languages; the user has to compile a program to LLVM IR, which is trivial to do with a vast set of existing languages. LLVM IR also benefits from an existing API due to its use within the LLVM Compiler Backend, which significantly simplifies the implementation of the latter stages of OMPRTF, given their reliance on program in-memory representation and modification. The existing API also eliminates the need for string manipulation. The IRs' use within the LLVM Backend also yields the benefit of being faster to compile than more complex languages, such as C or Fortran.

The second stage of OMPRTF, called cmeta, is responsible for the initial augmentation of the LLVM intermediate representation (IR). Its primary objective is to provide sufficient metadata to the profiler such that identifying problematic device offloads becomes as simple as indexing an instruction within a function. To achieve this, cmeta first comprehensively strips away all existing metadata, ensuring that the resulting representation is entirely detached from the source code. This process produces a more abstract, source-language-independent view of the program, highlighting one of LLVM's core strengths as a language-agnostic compilation framework.

After removing the original metadata, cmeta iterates through each function in the IR, assigning a unique identifier to every basic block and recording both an offset and a parent function for each instruction. The newly generated information is then formatted as DWARF debug metadata and written back into the IR, effectively overwriting the previous debug information. Although this procedure is conceptually elementary, it provides an elegant and robust solution to the problem of preserving precise IR-level location information within the final executable.

The third and fourth stages of OMPRTF involve compiling the augmented IR and executing the resulting program. Since the source code has already been lowered to Intermediate Representation, the compilation process is streamlined. The IR is passed directly to the LLVM backend to generate an executable ready for analysis. The profiling phase is then handled by a modified version of Luke Marzen's OMPDataPerf, which aggregates raw OpenMP target data operation logs to identify performance bottlenecks.

This profiling component detects various erroneous or wasteful operations, including duplicate data transfers, round-trip transfers, repeated allocations, and unused memory resources. By grouping operations by hashes, device identifiers, and memory regions, it computes both cumulative execution time and potential resource waste. While the original profiler was designed as a standalone interactive tool, the logic has been significantly refactored to function as a library. This adaptation enables programmatic access to structured analysis results, allowing OMPRTF to consume findings directly rather than parsing textual output.

A critical aspect of this refactoring was addressing the disconnect between runtime profiling and static analysis. The original tool relied on runtime instruction addresses, which are unusable in later offline optimization stages. To bridge this gap, the system now converts runtime addresses into relative, offline-stable representations using dynamic symbol resolution and base address subtraction. These normalized addresses are saved alongside specific error types in shared memory, allowing the pipeline to correlate runtime inefficiencies with IR-level metadata without requiring direct access to source-level debug information.

The pipeline culminates in the execution of the optimizer pass, a sophisticated LLVM transformation stage that operationalizes the intelligence gathered during profiling to enact surgical performance improvements. This pass does not merely apply generic heuristics; instead, it ingests the specific fault report generated by the profiler, identifying exact instructions associated with wasteful behaviors, such as redundant data movement or unused allocations. Upon locating the corresponding OpenMP offload directives within the static IR, the optimizer initiates a comprehensive expansion process. Standard OpenMP compilation typically lowers target regions into opaque, monolithic runtime library calls (such as `__tgt_target_kernel`) which encapsulate the entire lifecycle of an offload, argument mapping, host-to-device transfer, kernel execution, and device-to-host retrieval, into a single "black box" instruction. While efficient for code generation, this abstraction is detrimental to optimization because it hides the distinct operational steps from the compiler, treating the offload as an indivisible atomic unit.

To overcome this, the OMPRTF optimizer mechanically deconstructs these high-level calls, effectively "inline-expanding" the runtime logic directly into the user's control flow graph. This transformation breaks the abstraction barrier, exposing the constituent micro-operations that make up the offload event: the individual memory allocation requests, the specific pointer associations, the distinct data copy commands, and the synchronization barriers. By converting the opaque library call into a transparent sequence of discrete IR instructions, the pass gains fine-grained control over the execution flow that was previously impossible to achieve. It can now isolate the specific instruction responsible for a "duplicate transfer" or an "unused allocation" and remove it from the instruction stream without destabilizing the rest of the kernel launch. This localized expansion transforms the optimization problem from a binary choice of running or skipping an entire kernel into a precise editing process, where the compiler can retain the necessary computational logic while selectively pruning the specific data movement operations that the dynamic profiling stage identified as redundant.

The proposed system is subject to several inherent limitations, primarily stemming from the stochastic nature of profile-guided optimization. Because it is computationally infeasible to capture an execution trace that encompasses every possible input permutation and edge case, the generated profiles act only as approximations of typical runtime behavior. Consequently, optimizations tailored to specific input datasets may prove detrimental or functionally incorrect when subjected to divergent workloads. To mitigate this risk without compromising the semantic integrity of the program logic, the system must inject runtime guard clauses to ensure the program's integrity. These mandatory conditional checks, while necessary for correctness, introduce unavoidable control-flow overhead that can diminish net performance gains.

Furthermore, the architecture faces challenges regarding hardware heterogeneity and resource discovery. The current implementation does not robustly address device selection or load balancing in multi-GPU environments via OpenMP; conversely, in environments entirely lacking GPU resources, the profiler is prone to silent failure modes rather than graceful degradation. Finally, the system is engineered exclusively for the Linux ecosystem. Portability to Windows or macOS was explicitly deprioritized, a decision aligned with the High-Performance Computing landscape. As the Top500 list currently features no supercomputers running Windows or macOS, and given that current Apple silicon lacks the necessary driver support for OpenMP GPU offloading, cross-platform compatibility was deemed theoretically redundant.

Future research directions involve extending the system's capabilities beyond mere data collection into the realm of heuristic interpretation and analysis. Rather than simply reporting raw metrics, subsequent iterations of the profiler could perform logical analysis on the results, identifying patterns in execution flow that suggest deeper structural inefficiencies. Furthermore, the scope of the profiler should be expanded to include device-side kernel optimization. By analyzing the specific computational characteristics of the offloaded code, the system could suggest improvements to register usage or thread occupancy. Finally, moving from a reactive model, which targets only "troublesome" bottlenecks, to a holistic optimization strategy is a critical next step. By applying the optimization pass to all OpenMP offload regions comprehensively, the system could aggregate memory management operations. The comprehensive optimization of all OpenMP offloads would significantly amortize the cost of host-to-device transfers, leading to a drastic reduction in redundant alloc and delete calls and ensuring better data residency on the GPU.

A promising avenue for future research is the integration of adaptive data compression to mitigate the PCIe bandwidth bottleneck. While current optimizations focus on reducing the frequency of transfers, implementing a compression-aware runtime could significantly decrease the volume of data per transaction. We propose investigating lossless compression algorithms (such as LZ4 or Zstandard) for pointer-heavy data structures and lossy, error-bounded compression (like SZ or ZFP) for large floating-point arrays standard in scientific simulations. By offloading the decompression task to the GPU, the system could effectively "overclock" the available interconnect bandwidth. This work would require a sophisticated cost-benefit analyzer within the profiler to ensure that the computational overhead of compression does not exceed the time saved during transmission, particularly on high-speed interconnects like NVLink.