

Optimizing NWP Model Physics on Next-Generation Processors

John Michalakes¹, Thomas Henderson², Michael J. Iacono³

¹NOAA National Centers for Environmental Prediction, College Park, Maryland

²NOAA Earth System Research Laboratory, Boulder, Colorado

³Atmospheric and Environmental Research (AER), Lexington, Massachusetts

Exponentially increasing supercomputing capability over the past half century has enabled a linear increase in forecast skill and resulting value to the public. Beginning around 2004, however, processor clock speed stopped scaling as rapidly as transistor density, which alone continues to double every several years. Manufacturers are now developing next-generation processors that are extremely floating-point capable – more than a Teraflop/second peak performance per chip – but that place new burdens on application developers to expose high degrees of parallelism in order to realize even a fraction of this potential. Examples include NVIDIA's General Purpose Graphics Processing Units (GPGPUs) and Intel's Many Integrated Core (MIC) architecture. One may also include new vector instruction sets that are finding their way onto new generations of conventional multi-core processors.

Many groups around the world, including ours and others in NOAA, are working to make efficient use of these new architectures by characterizing application behavior with respect to computational intensity, vectorization, concurrency, locality, and memory system performance; identifying and testing effective data organization and looping and other code restructuring strategies; and exploring new programming approaches that will leverage performance gains while minimizing impacts on development and maintenance of large NWP code bases. One area of our research has been to improve the performance of microphysics and radiative transfer physics, among the most expensive used in models that run operationally in NOAA.

The Rapid Radiative Transfer Model (RRTMG) developed at AER is in use within the Integrated Forecast System of the European Centre for Medium Range Weather Forecasts; the Community Earth System Model and the Weather Research and Forecast (WRF) model at NCAR; the Global Forecast System, the Climate Forecast System, and the NMM-B regional model at NOAA/National Centers for Environmental Prediction; and others. RRTMG is one of the most expensive physical processes in the NMM-B, costing upwards of eight percent of an overall forecast cost. RRTMG code and data were restructured to increase thread and vector parallelism necessary on the Knights Corner (KNC) version of the Intel MIC. Figure 1 shows the effect of optimizations on time spent in RRTMG relative to the baseline: adding an inner vector dimension to expensive shortwave radiation calculations (*chunk=8*); defining vector width at compile time (*+static*); and interleaving shortwave and longwave calculations on adjacent OpenMP threads (*+task interleave*) to reduce resource competition between threads. The restructured RRTMG code ran three times faster than the original code on the Intel MIC and thirty percent faster on the host Xeon processor. Similar improvements were seen for the GPU version of the RRTMG radiation developed by AER. We worked with NCAR to provide a combined GPU/MIC/Multi-core version of RRTMG that will be available to users in the WRF V3.7 release in Spring 2015.

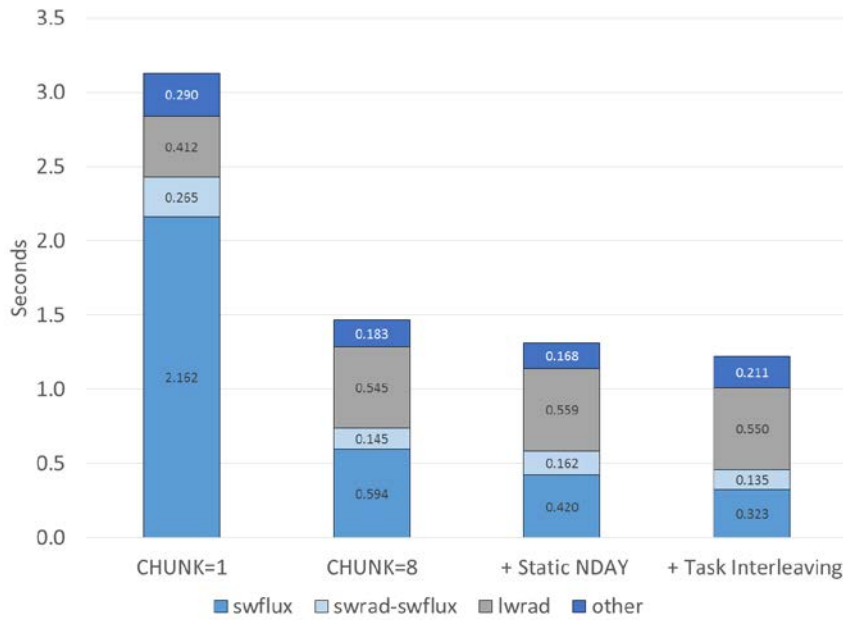


Figure 1. Effect of successive optimizations of the RRTMG kernel on the Intel MIC “Knights Corner” 60-core processor, relative to the baseline (CHUNK=1).

The WRF Single Moment 6-class Microphysics (WSM6) scheme is used in a variety of research and operational NWP models including NOAA’s Non-hydrostatic Icosahedral Model (NIM). WSM6 was optimized using a variety of techniques including threading, vectorization, array alignment, improving data locality, and use of compile-time constants for loop and array index bounds. We tested WSM6 on the KNC and on successively newer versions of Intel’s conventional multicore Xeon Processors: Sandybridge (SNB), Ivybridge (IVB) and Haswell (HSW). As with the RRTMG package, optimizations that improved WSM6 performance on the KNC coprocessor also provided benefit on SNB, IVB and HSW. WSM6 performance on the current KNC generation of MIC lagged behind its Xeon counterparts.

Device	Threads	Baseline (seconds)	Optimized (seconds)	Improvement
SNB	32	9.4	7.5	1.25
IVB	48	4.7	3.4	1.38
HSW	56	--	2.6	--
KNC	240	13.2	8.7	1.52

This work has also provided a testbed for investigating performance-portable programming models. The Intel MIC is programmed within the same overarching software environment as the Intel Xeon processor family. For GPU, directives-based approaches such as OpenACC and OpenMP extensions express fine-grained parallelism less invasively than device-specific languages such as CUDA and OpenCL, but with lower realized performance. Continued work is supported under a new Software Engineering for Novel Architectures (SENA) project that begins funding under the NOAA High Performance Computing Program in 2015.