

Green Flash: Ultra-Efficient Climate Computing More Science Using Less Power

<http://www.lbl.gov/CS/html/greenflash.html>

Berkeley Lab researchers have proposed an innovative way to improve global climate change predictions by using a supercomputer with low-power embedded microprocessors, an approach that would overcome limitations posed by today's conventional supercomputers.

A paper published in the May 2008 issue of the *International Journal of High Performance Computing Applications* lays out the benefit of a new class of supercomputers for modeling climate conditions and understanding climate change. Using the embedded microprocessor technology found in cell phones, iPods, toaster ovens, and many other modern-day electronic conveniences, they propose designing a cost-effective machine for running these models and improving climate predictions.

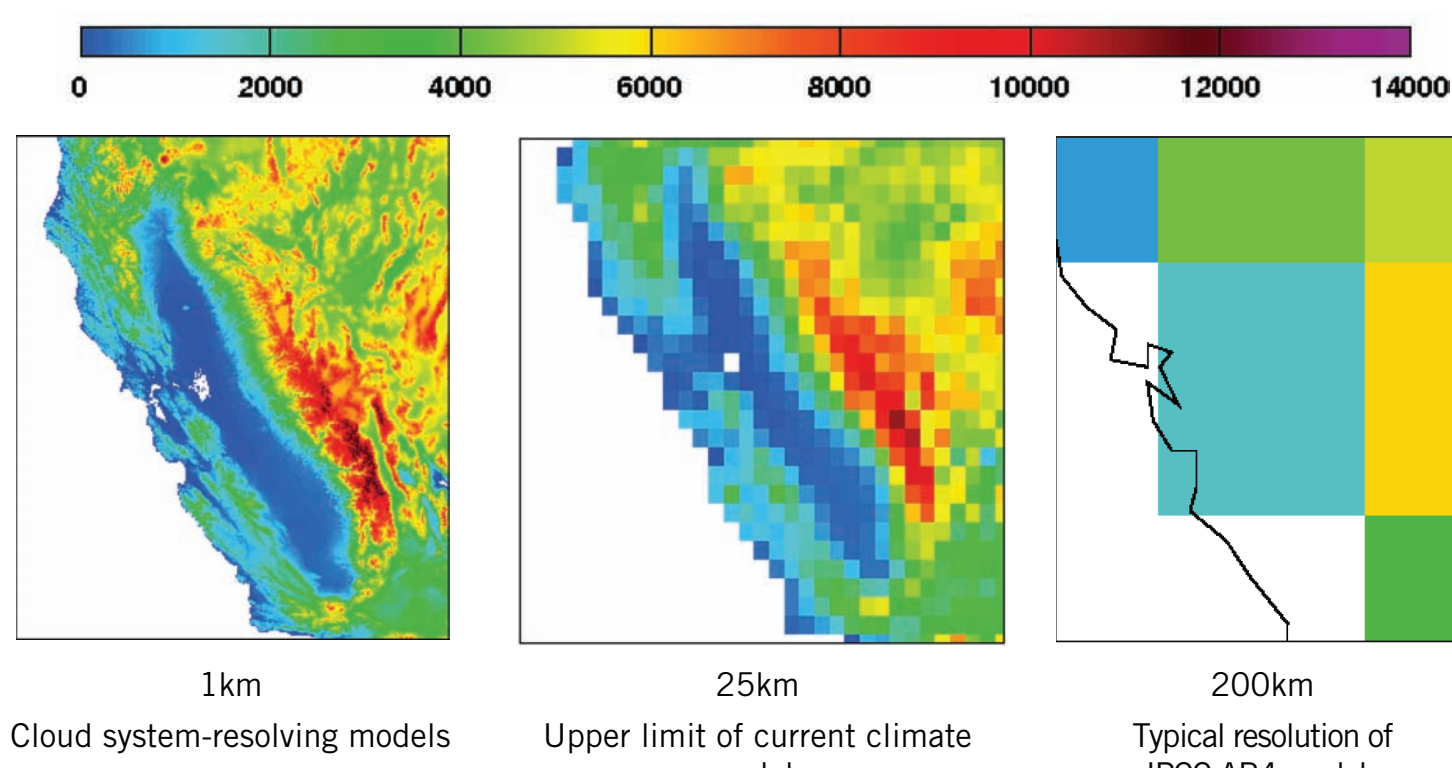
This research project has been named "Green Flash," after the *optical phenomenon* that sometimes appears on the horizon at sunset or sunrise.

Global Cloud-Resolving Model Require Exascale Computing

Currently, a major shortcoming to more accurate predictions is the quality of the cloud simulations. Namely, global models do not resolve individual clouds so subgrid scale statistical parameterizations are used to account for moist convective processes. Direct numerical simulation of cloud systems removes the need for cumulus parameterization but requires horizontal grid resolutions approaching one kilometer.

The challenge of constructing both a model and a machine to carry out century-scale global climate simulations at the kilometer scale requires transformational changes to computational technologies. Application-specific architecture design is the most efficient path to exascale computing for climate modeling as well as for many other high-end scientific disciplines.

Surface Elevation in Feet



Power Demand Is the Major Challenge to Exascale Computing

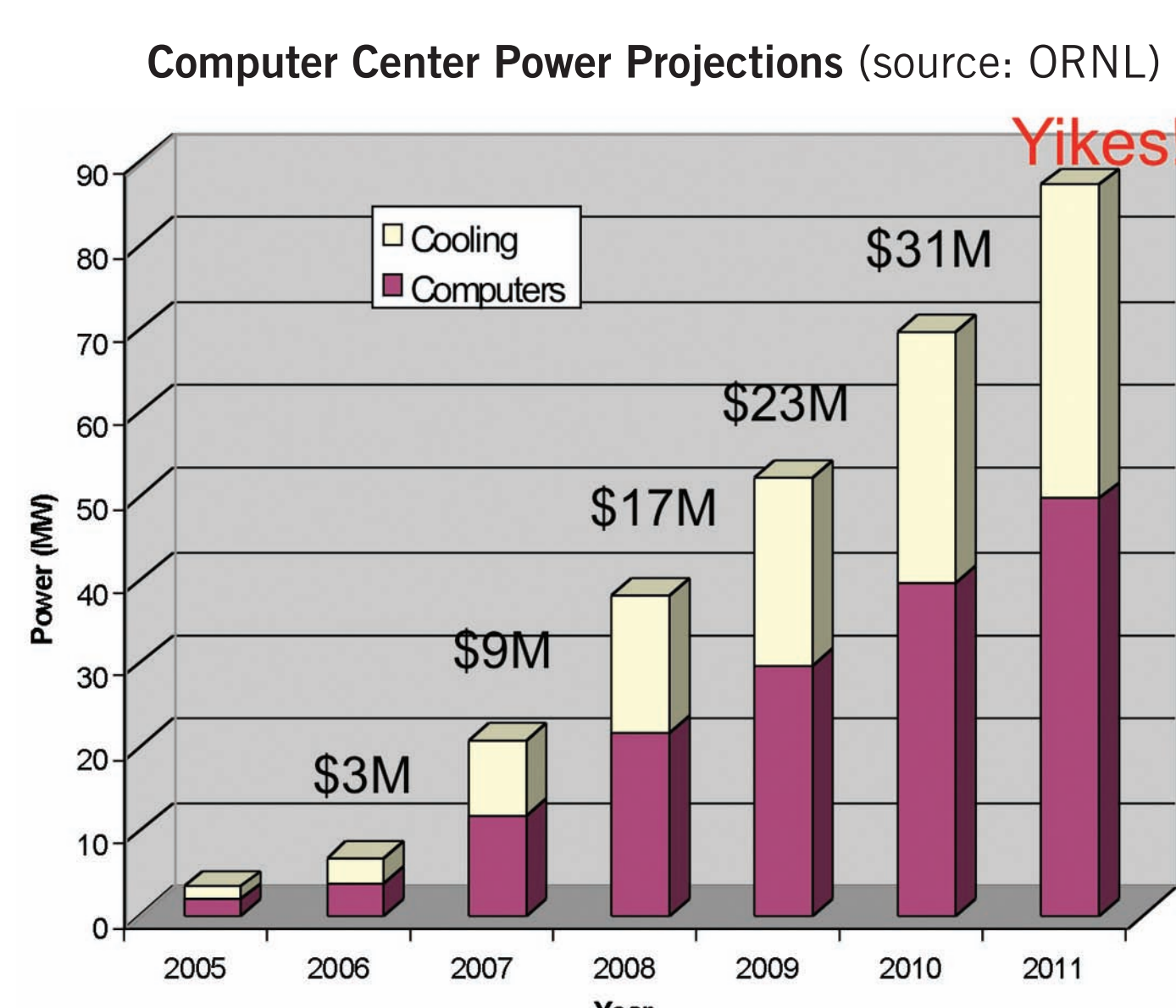
PROCESSOR	CLOCK	PEAK/CORE (Gflops)	CORES/SOCKET	SOCKETS	CORES	POWER
AMD Opteron	2.8GHz	5.6	2	890K	1.7M	179MW
IBM BG/P	850MHz	3.4	4	740K	3.0M	20MW
Climate Computer	650MHz	2.7	32	120K	4.0M	3MW

Even if each simple core is 1/4 as computationally efficient as a complex core, you can fit hundreds of them on a single chip and still be 100 times more power-efficient.

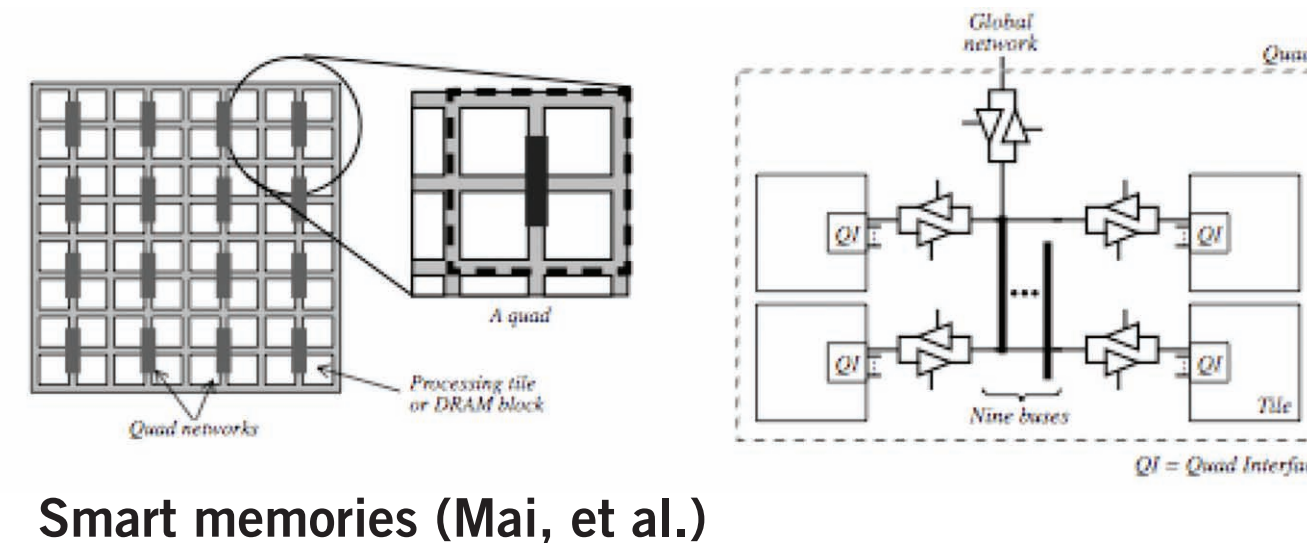
The cost of power required to run HPC systems is projected to match or exceed hardware costs — Berkeley Lab's extrapolation of Blue Gene and AMD design trends shows energy demands of 20MW and 179MW, respectively, when applied to climate science.

Application-Driven Architecture Design Improves Efficiency

Green Flash proposes tailoring the system architecture to climate modeling problems, replacing bloated serial processors with simple low-power cores. Borrowing ideas from embedded computing, the Green Flash team is co-designing hardware and software to gain greater efficiency. Using the Tensilica processor generator allows the fast creation of semi-customized cores, along with the compilers and other tools to use them.



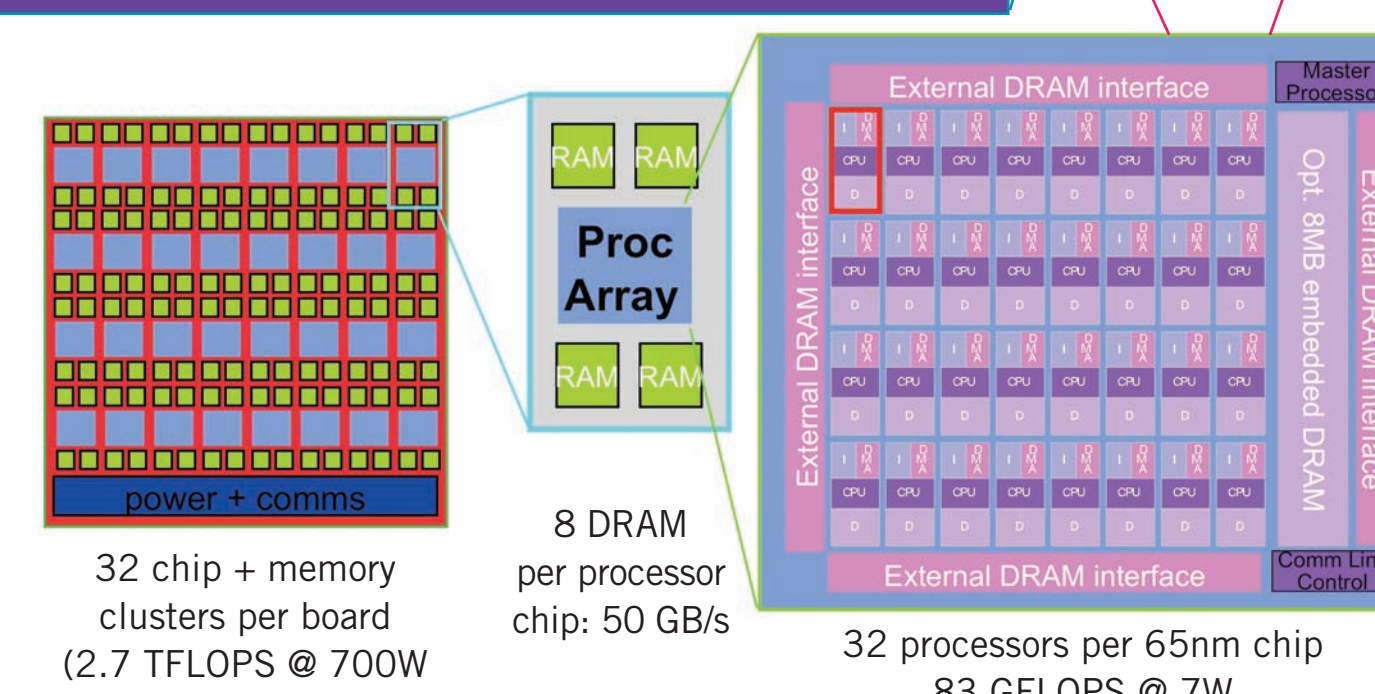
"The only way to lower power consumption is to reduce waste."
—Mark Horowitz



VLIW CPU:

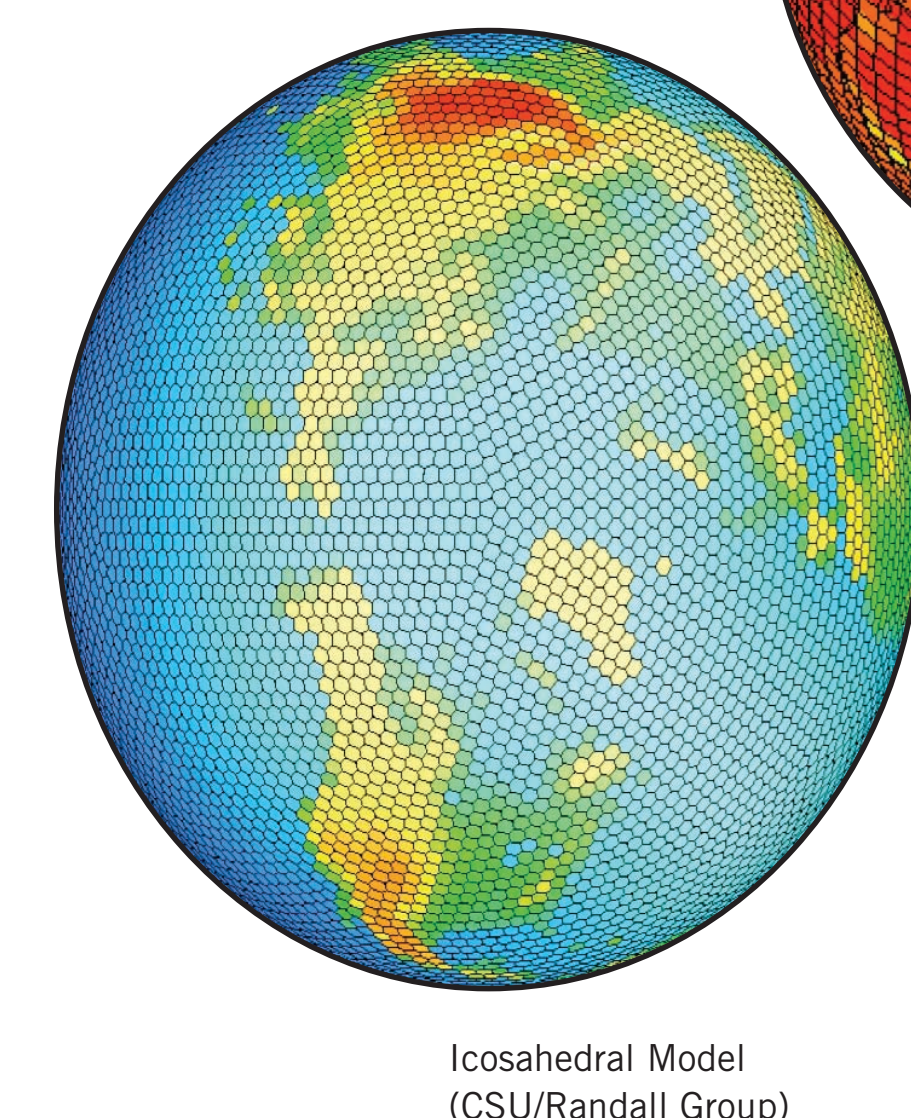
- 128b load-store + 2 DP MUL/ADD + integer op/ DMA per cycle
- Synthesizable at 650MHz in commodity 65nm
- 1mm2 core, 1.8-2.8mm2 with inst cache, data cache data RAM, DMA interface, 0.25mW/MHz
- Double precision SIMD FP : 4 ops/cycle (2.7GFLOPs)
- Vectorizing compiler, cycle-accurate simulator, debugger GUI (Existing part of Tensilica Tool Set)
- 8 channel DMA for streaming from on/off chip DRAM
- Nearest neighbor 2D communications grid

32K I
8 chan DMA
CPU
64-128K D
2x128b



Advanced Dynamics Algorithms Are Necessary for Global Cloud-System-Resolving Climate Models

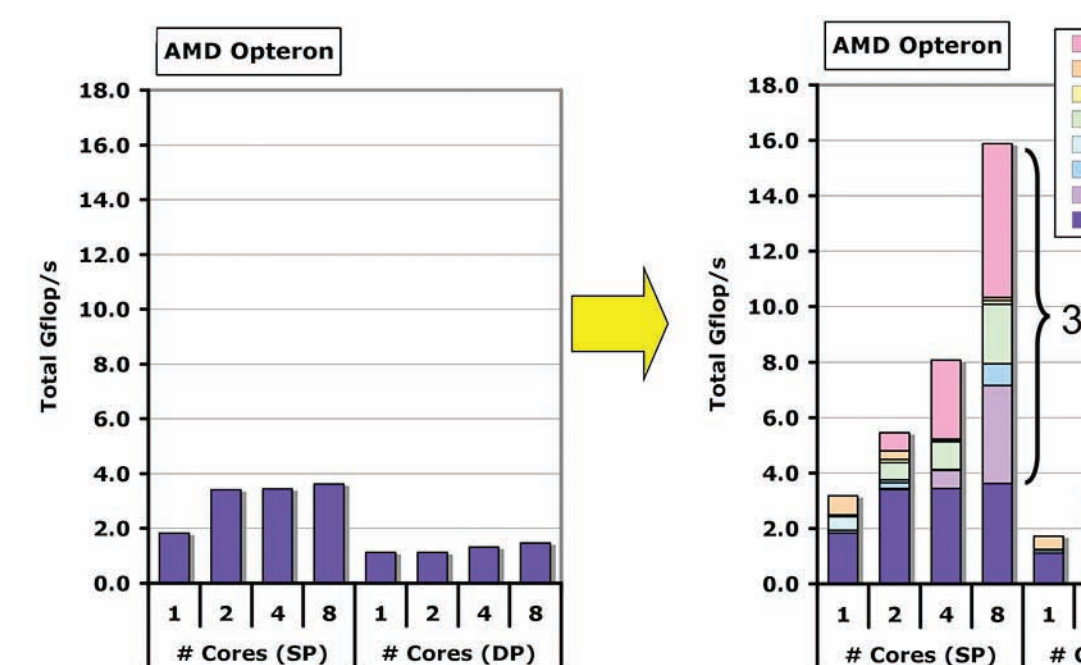
Climate models must simulate weather 1000 times faster than it actually occurs in order for scientists to use them effectively. At the high resolutions required to resolve cloud systems, advanced dynamics algorithms on novel grids such as icosahedral, cubed-sphere, or reduced grids are needed. We estimate that a one kilometer global model with 100 or more vertical levels can express 20 million way parallelism through a combination of horizontal and vertical domain decomposition. A machine with hundreds of thousands of large-scale multi-core chips would be ideally suited to this application.



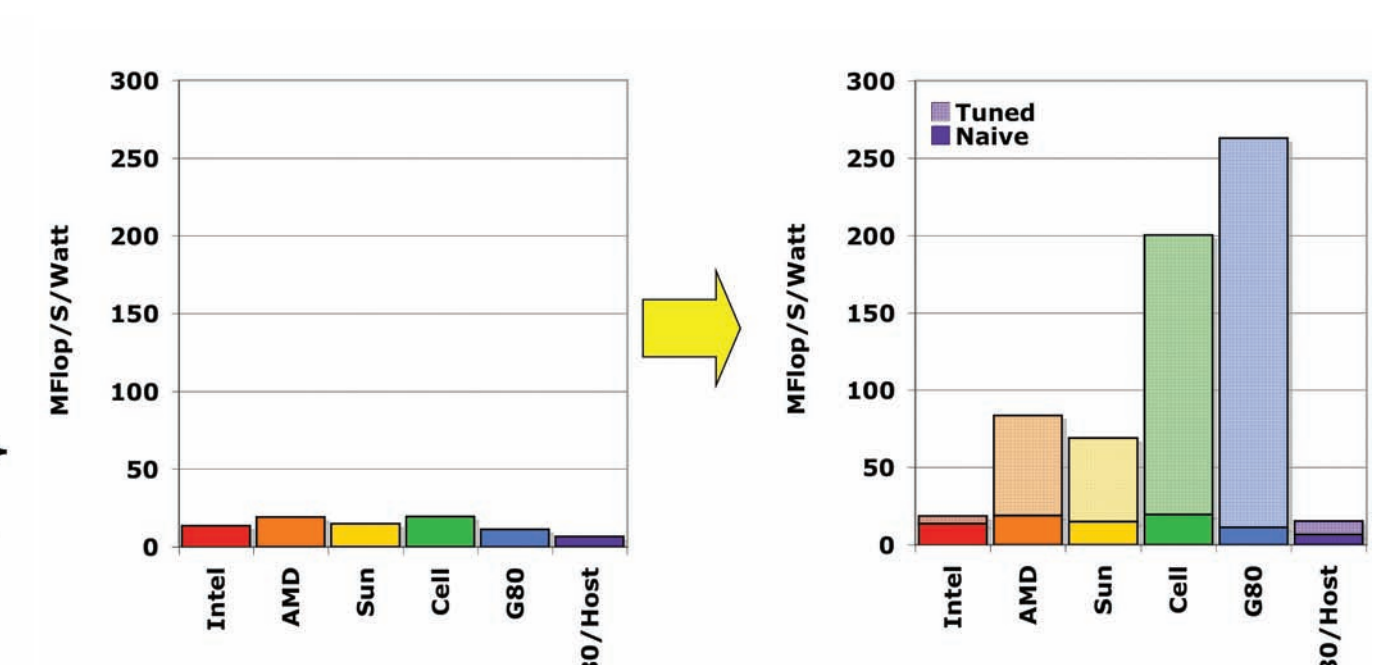
Auto-Tuning Framework

- Transitioning between different, complex architectures transparent to the user
- Allows performance portability across evolving architectural designs
- Domain-specific knowledge allows for better optimization than traditional compilers

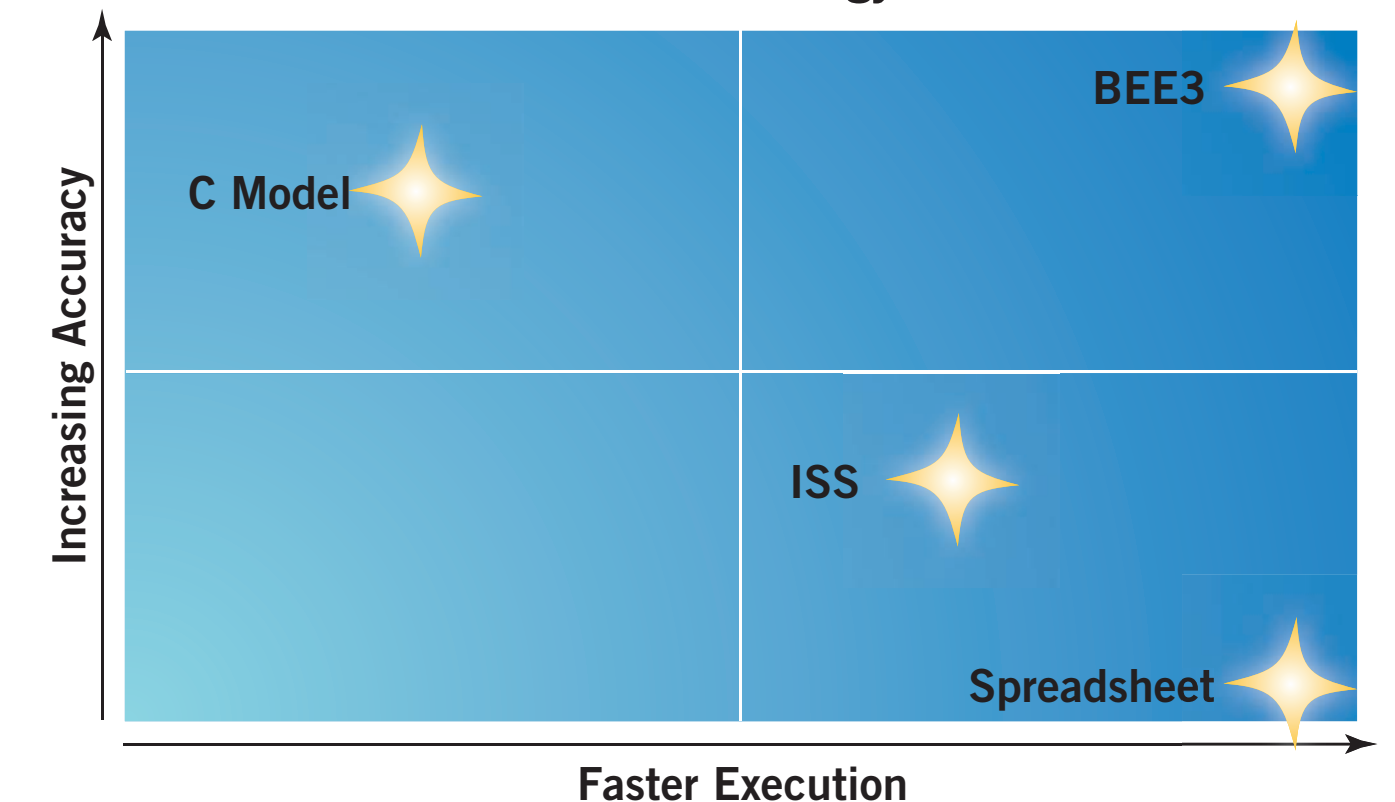
Auto-tuning increases performance, scalability . . .



and power-efficiency



Simulation Methodology Trade-Off



RAMP Provides Fast Hardware Emulation Using BEE3

The climate code is a long-running and highly complex workload. Current software performance simulators are slow and difficult to verify — and become intractable when high-flexibility, high-detail, and fast-execution are required! So, Green Flash is approaching the problem by using novel hardware emulation.



Green Flash leverages RAMP (the Research Accelerator for Multiple Processors), which is a multi-institutional effort to use reconfigurable logic for research into multi-processor architectures. RAMP has developed BEE3 (the Berkeley Emulation Engine, version 3), a 2U chassis with a tightly coupled 4 FPGA system for research computer architecture. By using BEE3 to perform direct hardware emulation, the Green Flash prototype is fast enough to allow the execution of a full climate model in a short amount of time.