# Towards Memory Oriented Scalable Computer Architecture and High Efficiency Petaflops Computing

Thomas Sterling

Center for Advanced Computing Research
California Institute of Technology

**Abstract**. The separation of processor logic and main memory is an artifact of the disparities of the original technologies from which each was fabricated more than fifty years ago as captured by the "von Neumann architecture". Appropriately, this separation is designated as "the von Neumann bottleneck". In recent years, the underlying technology constraint for the isolation of main memory from processing logic has been eliminated with the implementation of semiconductor fabrication foundries that permit the merger of both DRAM bit cells and CMOS logic on the same silicon dies. New classes of computer architecture are enabled by this opportunity including: 1) *system on a chip* where a conventional processor core with its layers of cache are connected to a block of DRAM on the same chip, 2) *SMP on a chip* where multiple conventional processor cores are combined on the same chip through a coherent cache structure, usually sharing the L3 cache implemented in DRAM, and 3) *processor in memory* where custom processing logic is positioned directly at the memory row buffer in a tightly integrated structure to exploit the short access latency and wide row of bits (typically 2K) for high memory bandwidth. This last, PIM, can take on remarkable physical structures and logical constructs and is the focus of the NASA Gilgamesh project to define and prototype a new class of PIM-based computer architecture that will enable a new scalable model of execution. The MIND processor architecture is the core of the Gilgamesh system that incorporates a distributed shared memory management scheme including in-memory virtual to physical address translation, a lightweight *parcels* message-driven mechanism for invoking remote transaction processing, multithreaded single cycle instruction issue for local resource management, graceful degradation for fault tolerance, and pinned threads for real time response. The MIND architecture for Gilgamesh is being developed in support of "sea of PIMs" systems for both ground based Petaflops scale computers and scalable space borne computing for long term autonomous missions. One of its specific applications is in the domain of symbolic computing for knowledge management, learning, reasoning, and planning in a goal directed programming environment. This presentation will describe the MIND architecture being developed through the Gilgamesh project and its relation to the Cray Cascade Petaflops computer being developed for 2010 deployment under DARPA sponsorship.