# Cascade

## Burton Smith
## Cray Inc.

- DARPA's High Productivity Computer Systems program has three central objectives:
  - 0 Improve HPC system productivity
  - 0 Improve HPC programmer productivity
  - 0 Improve system robustness (reliability and security)
- Three phases are planned:
  - 0 Phase 1 (7/02–7/03): Define a system concept
    - ◆ Cray, HP, IBM, SGI, Sun
  - 0 Phase 2 (7/03–7/06): Prepare a development plan
    - ◆ Cray, IBM, Sun
  - 0 Phase 3 (7/06–7/10): Develop a system
    - ◆ Two awardees

- High system productivity
  - 0 Implement *very* high global bandwidth
  - 0 Use that bandwidth (the "wires") well
  - 0 Provide configurability to match user needs
- High human productivity and portability
  - 0 Support a mixed UMA/NUMA programming model
  - 0 Deliver strong compiler and runtime support
  - 0 Pursue higher level programming languages
- System robustness
  - 0 Virtualize all resources
  - 0 Make all resources dynamically reconfigurable
  - 0 Use introspection to detect bugs or intrusion

# Bandwidth is expensive

- High global system bandwidth is a "good thing"
  - It determines performance for many problems
  - It also makes improved programmability possible
- Sadly, connection costs badly trail Moore's law
  - Packages, circuit boards, wires, optical fibers…
  - Most of hardware cost is connection cost
- Cray builds systems with high global bandwidth
  - This strongly influences most of what we do
- Cray needs to stay competitive
  - We must make bandwidth cost less
  - We must use bandwidth wisely
- These ideas motivate much of Cascade's architecture

- Tune bandwidth ($\therefore$ cost) to match customer needs
- Use the cheapest link technology that fits each bill
  - $0$ Optical or electrical
- Make data rates as fast as the technology permits
  - $0$ Spending transistors in this cause is a bargain
- Design routers that use all the network links well
  - $0$ Randomized non-minimal routing, for example
- Use efficient network topologies
  - $0$ High degree routers

- If network link load is well balanced, node injection bandwidth B times average distance d (in hops) is bounded by link bandwidth $\beta$ times node degree $\Delta$
- Cost/node is proportional to $\beta$ times $\Delta$
  - $0$ Signaling rate and package pin count determine it
- Increasing the degree $\Delta$ lowers the average distance d
  - $0$ $\beta$ can be lowered to maintain (or even improve) cost
  - $0$ B will increase as d decreases
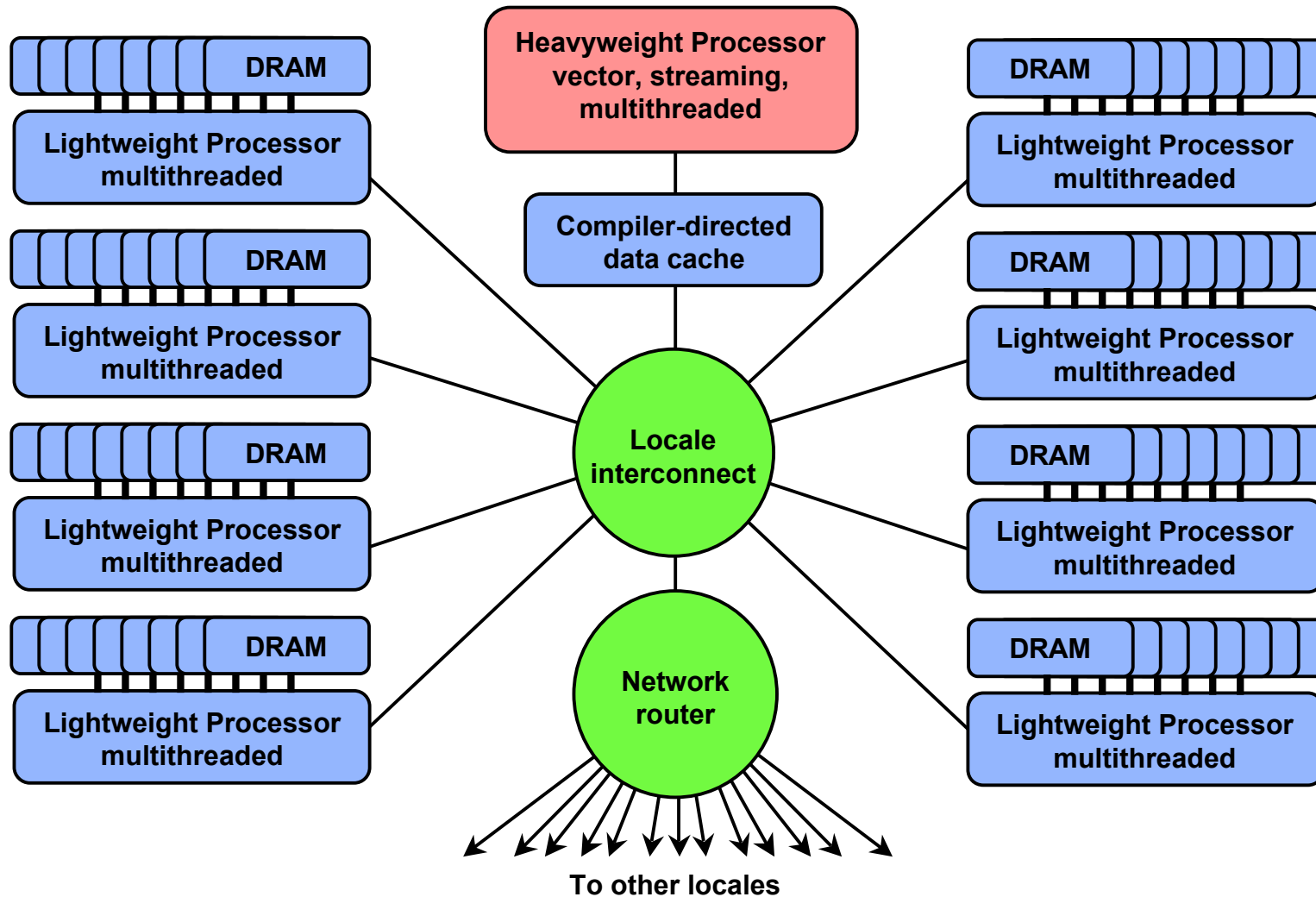- Conclusion: trading high node degree for link bandwidth can yield better injection bandwidth, latency, and cost

- Implement shared memory (UMA/NUMA hybrid)
  - 0 Eliminate overhead to enable small messages
- Tolerate memory latency
  - 0 CC-NUMA wastes bandwidth moving data around
  - 0 Use vectors and multithreading instead
- Exploit temporal locality in "heavyweight" processors
  - 0 Compiler-directed data cache
  - 0 Architectural support for streaming
- Exploit spatial locality in "lightweight" processors
  - 0 Threads migrate to follow the spatial locality
- Use other types of locality whenever possible
  - 0 *e.g.* atomic memory operations

# A locale building block

- The compiler often knows if data are safely cacheable, *i.e.* are temporarily private or temporarily constant
- It can tell the hardware what data to cache and when to flush or simply invalidate it
  - 0 Dead values as well as constants are invalidated
- Unnecessary coherence traffic is eliminated
- Latency and network bandwidth demand are reduced
- Threads within a processor can communicate and synchronize within cache to exploit streaming locality
- The cache becomes a much more general tool for exploiting many forms of temporal locality

- Lightweight threads in the memory can exploit spatial locality by migrating to memory they refer to
  - 0 Some remote references just block the thread
  - 0 Others cause migration to the remote memory
- Memory is block-hashed to provide a compromise between spatial locality and reference distribution
- Processor-in-memory (PIM) technology is an ideal implementation vehicle for this idea
- Threads are spawned by sending parcels to memory from either heavyweight or other lightweight threads
  - 0 Spawning and migration overheads are minimal
  - 0 In-memory operations are handled specially
- Generally, the compiler packages temporally local loops for heavy threads and the rest for light ones

# Sparse matrix-vector product

```
int nrows, rowp[];
val_idx_pair a[];
double b[], c[];
for (int i = 0; i<nrows; i++) {
  double sum = 0.0;
  for(int k = rowp[i]; k < rowp[i+1]; k++){
    sum += a[k].val*b[a[k].idx];
  }
  c[i]=sum;
}
```

**The data layout:**

$\qquad$ rowp[i] $\longrightarrow$ $\qquad$ rowp[i+1] $\longrightarrow$

a[k].val:  $\quad . \quad . \quad . \quad a_{i,15} \quad a_{i,42} \quad a_{i,53} \quad . \quad . \quad .$

a[k].idx:  $\quad . \quad . \quad . \quad 15 \quad\quad 42 \quad\quad 53 \quad . \quad . \quad .$

- There are three memory references for every two flops
  - 0 Two memory references are local and one isn't
  - 0 There are 2 flops per "global" memory reference
  - 0 (Dense) inner product is the same if one or both vectors is unit stride
- The thread context required for this loop comprises:
  - 0 The program counter (the "method")
  - 0 Exception flags, *etc.* (perhaps packed with the PC)
  - 0 The pointer to the vector of value-index pairs `a[]`
  - 0 A limit for loop control, set to `&rowp[i+1]`
  - 0 Two pointers to the vectors of doubles `b[]` and `c[]`
  - 0 The double accumulator `sum`
  - 0 A few temporaries, *e.g.* for `b[a[k].idx]`

CRAY

- Matlab lets scientists be productive programmers
  - 0 Execution performance is marginal at best
  - 0 Manual translation to Fortran is the typical fix
- We are developing *Chapel*, a programming language aimed at both programmability and performance
- Its key features:
  - 0 Interprocedural polymorphic type inference
  - 0 Locality abstraction via first-class domains
  - 0 Explicit parallel operations over domains
  - 0 Implicit parallelism packaging and optimization
  - 0 Automatic thread and memory management
  - 0 Open-source implementation
- We will also support mixed-legacy-language programs
  - 0 Fortran, C++, MPI, shmem, coarray languages

# NAS CG conj_grad() in Chapel

**Function return type elided (inferred from return statement)**

**Parameter types elided (inferred from callsite)**

**Local variable types elided (inferred from initializer, uses)**

**Built-in array reductions**

**Whole-array operations ⇒ data parallel implementation**

**Sequential iteration over an anonymous domain**

**Operate on sparse arrays as though dense, and independently of implementing data structures**

**Partial array reductions**

**Global view ⇒ processors not exposed in computation, array sizes**

**Separation of concerns ⇒ locale views, domain/array distributions & alignments, and sparse data structures are expressed elsewhere**

**Composable parallelism ⇒ this (parallel) function could be called from a parallel task (which in turn could be called from another…)**

**Promotion of scalar operators, values, and functions**

**Support for tuples**

```
function conj_grad(A, X): {
  const cgitmax = 25;

  var Z = 0.0;
  var R = X;
  var P = R;
  var rho = sum R**2;

  for cgit in (1..cgitmax) {
    var Q = sum(dim=2) (A*P);

    var alpha = rho / sum (P*Q);
    Z += alpha*P;
    R -= alpha*Q;

    var rho0 = rho;
    rho = sum R**2;
    var beta = rho / rho0;
    P = R + beta*P;
  }
  R = sum(dim=2) (A*Z);
  var rnorm = sqrt(sum (X-R)**2);

  return (Z, rnorm);

}
```

**Fortran+MPI  =  173-288 lines  (1265 tokens)**
**Chapel         =  20 lines          (150 tokens)**

- Operating system
  0 Scalability, robustness, utility
- System infrastructure
  0 RAS system, power, cooling
- Interconnect implementation
  0 Router design, network topology
- Productivity assessment
  0 Metrics, modeling, prediction, applications
- Implementation technology
  0 Interconnect, chip packaging, power, cooling
- Debugging
  0 Correctness, performance
- Marketing
  0 Costs to develop and manufacture, sales outlook

# Cascade collaborators

- Cray Inc.
  - 0 Burton Smith, David Callahan, Steve Scott, . . .
- Caltech/JPL
  - 0 Thomas Sterling, Hans Zima, Larry Bergman, . . .
- Notre Dame
  - 0 Peter Kogge, Jay Brockman, . . .
- Stanford
  - 0 Bill Dally, Christos Kozyrakis, . . .

The Cray team has experience in these technologies:
- Latency-tolerant vector NUMA systems
- Latency-tolerant multithreaded UMA systems
- Processor-in-memory technology
- High bandwidth interconnection networks
- High-productivity compiler technology
- Whole-program, incremental compilation
- Run-time systems for fine-grain synchronization
- Scalable, highly productive operating systems
- Supercomputer system integration

- HPCS matches Cray business objectives well
- We and our collaborators have expertise in the technological directions we intend to pursue
- We are confident of a successful outcome