



Scalability of Weather & Climate Codes

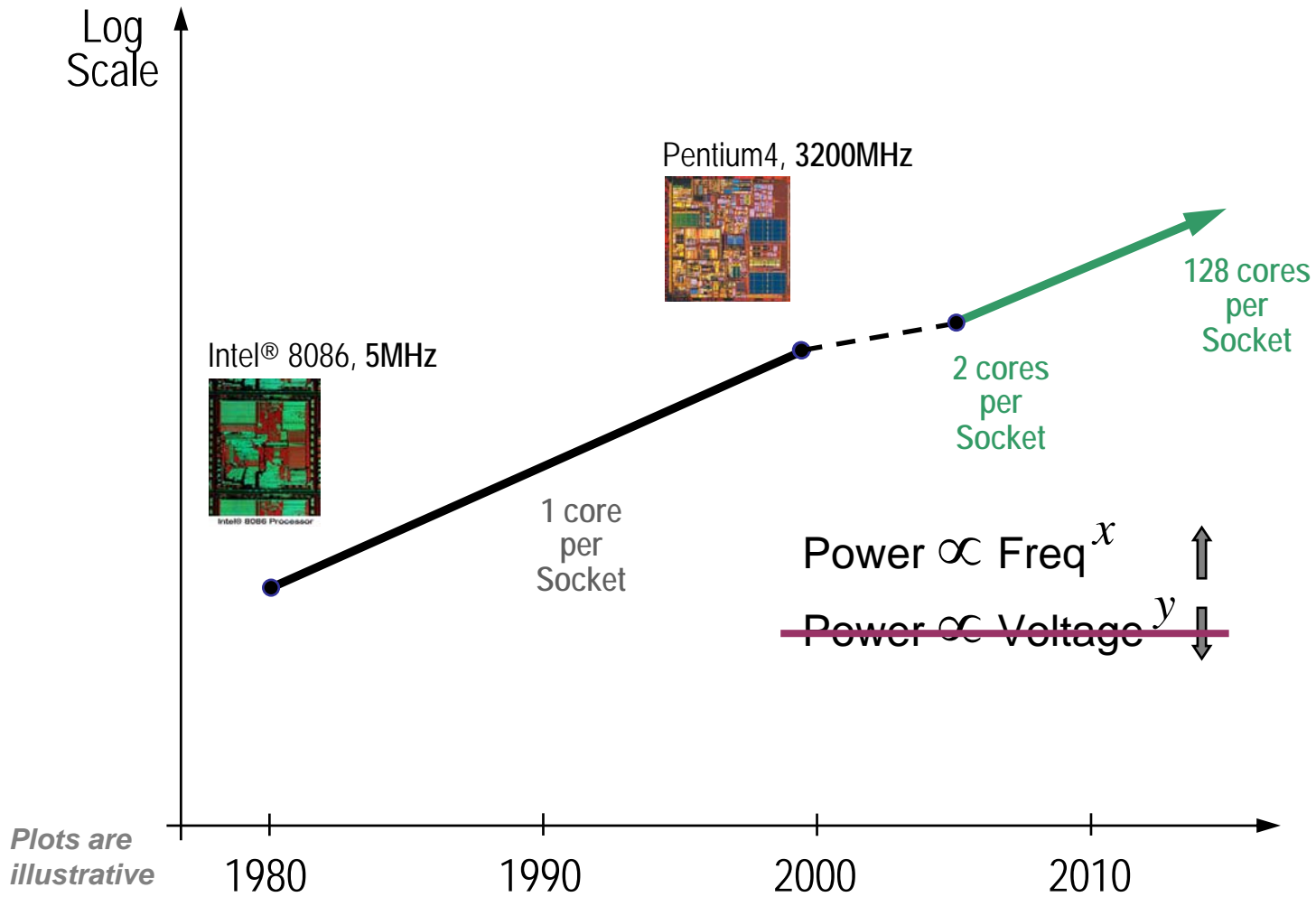
ECMWF Workshop Nov04'08

Back to Basics : Proposed New Architecture for Many-Core

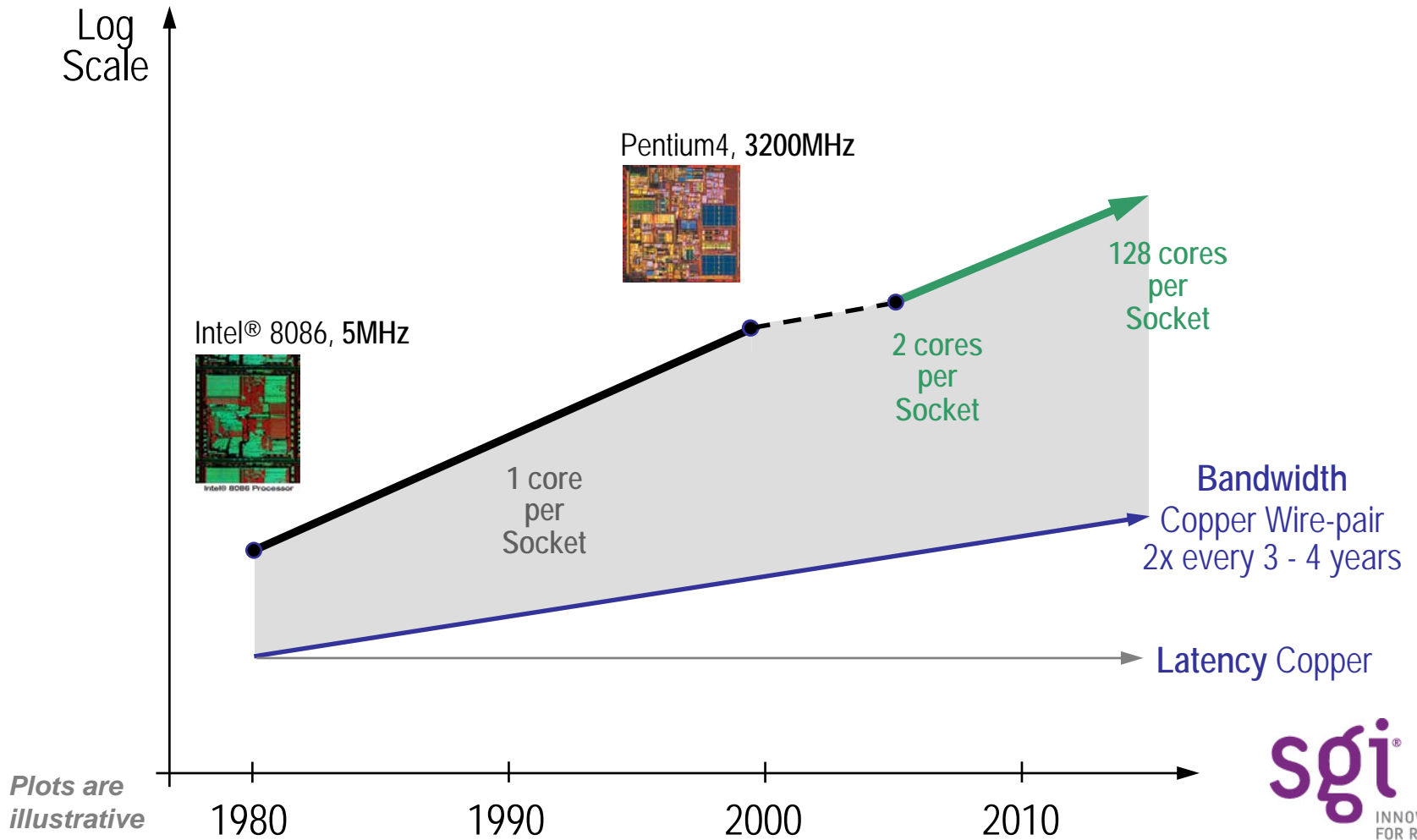
- Moore's Law Result Shift : Clock → Many-Core
- Problem Shift : Scaling Software Applications
- Load Balancing + Communications Costs
- New Architecture : Communications Focused

PGAS
+
**Programmable
Communications
Accelerators**

Trends

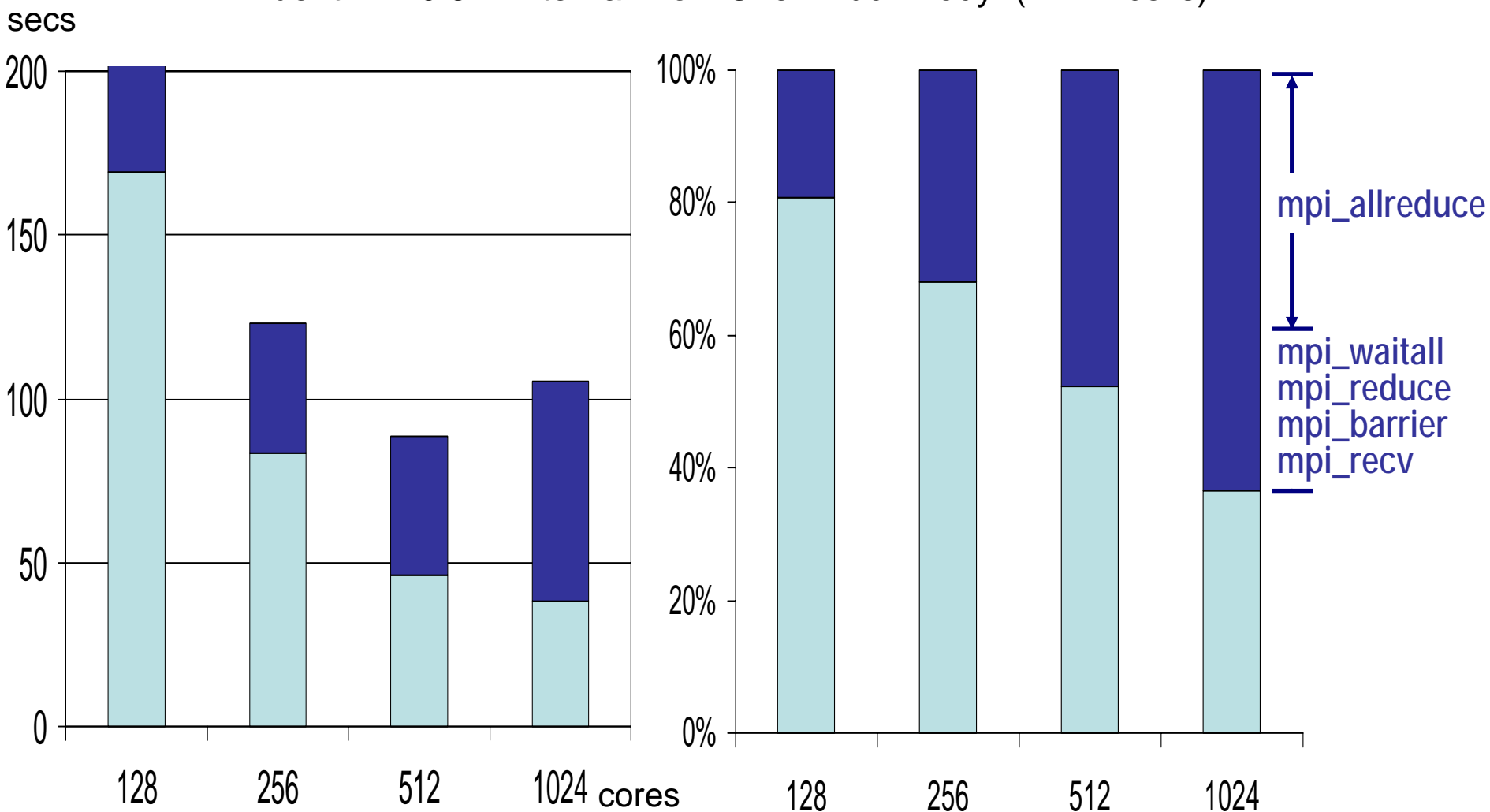


Trends

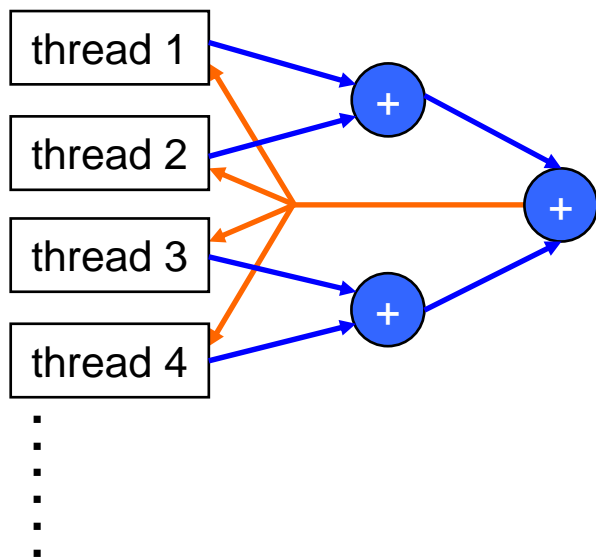


Total Run Time = Computation + Communications (bw & latency)

Fluent v12.0.5 : External Flow Over Truck Body (111M cells)



Latency : Pay a lot for relatively little gain (e.g. Global Collectives)



`mpi_allreduce ();`

Interconnect

- GigE ($64^3=256K$ mpi ranks)
- IB 4x DDR ($64^3=256K$ mpi ranks)

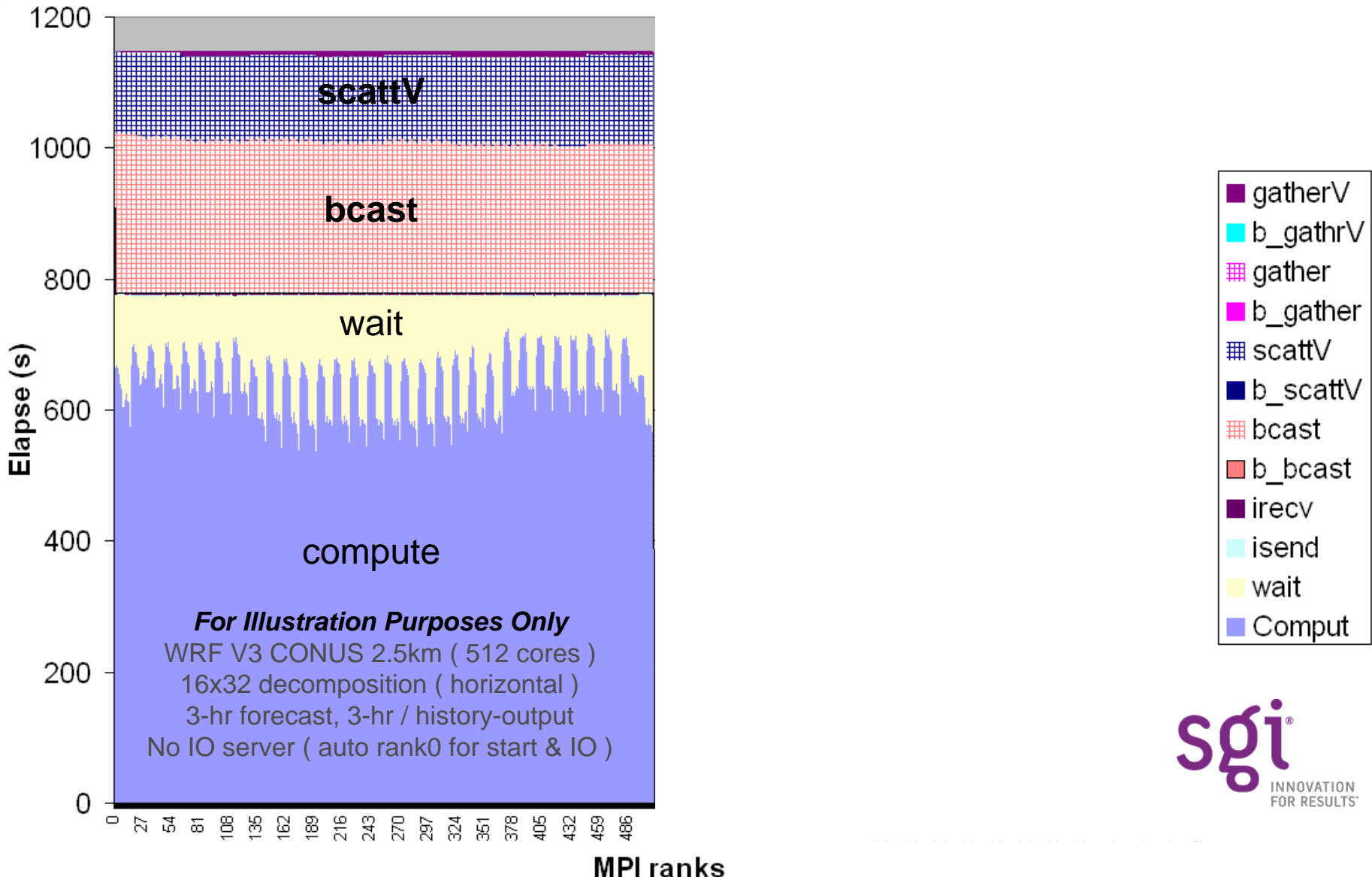
`mpi_allreduce ();`

100s us

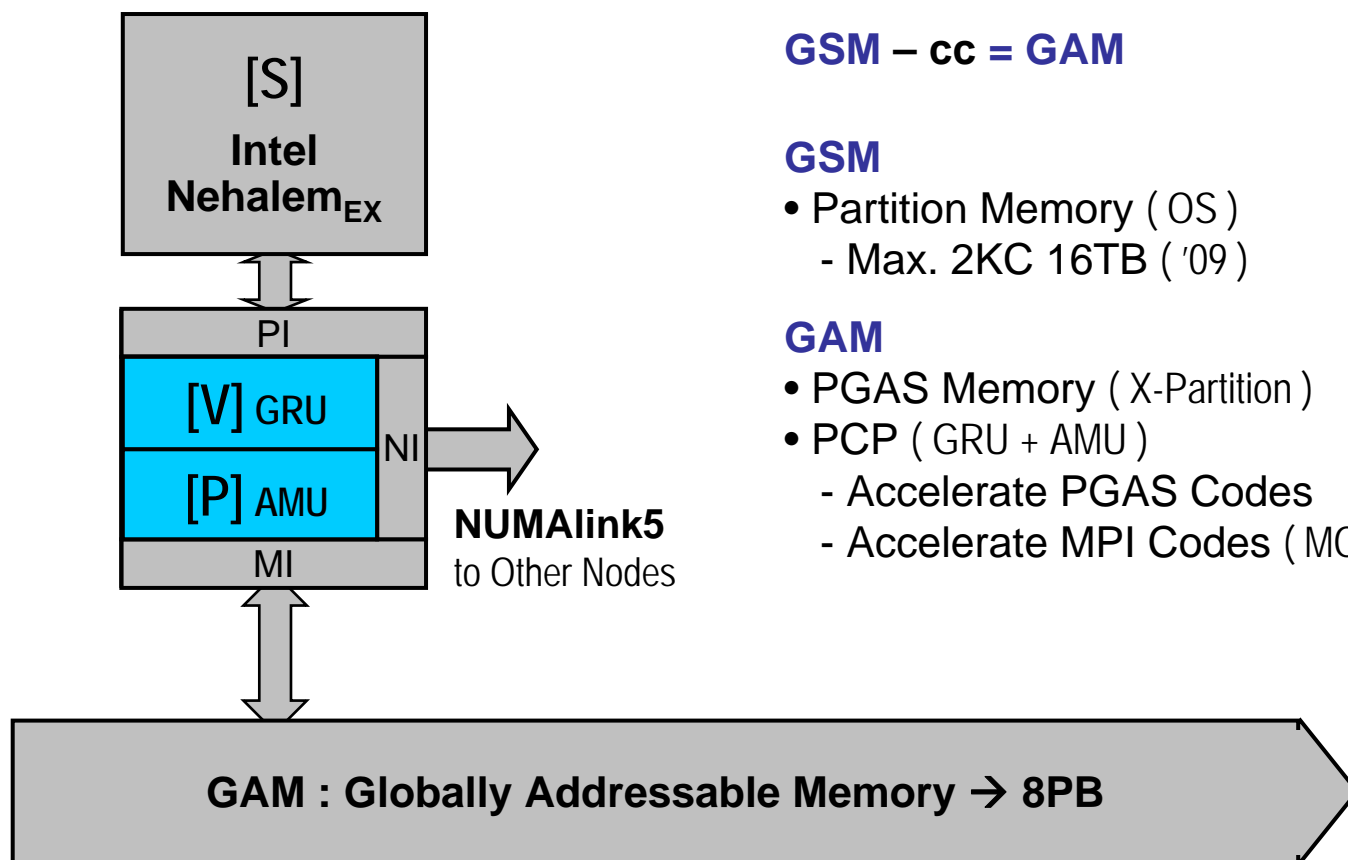
10s us

```
#pragma omp parallel \  
reduction(operator:list)  
upc_all_reduce
```

Total Run Time = Computation + Communications (bw & latency)



Ultraviolet : PGAS + PCP Architecture (accelerate communications)



GSM – cc = GAM

GSM

- Partition Memory (OS)
 - Max. 2KC 16TB ('09)

GAM

- PGAS Memory (X-Partition)
- PCP (GRU + AMU)
 - Accelerate PGAS Codes
 - Accelerate MPI Codes (MOE v.v. TOE)

Ultraviolet : PGAS + PCP Architecture for...

PGAS API

- UPC library support
- CAF F95 extensions (formally in Fortran-2008)
- Others welcome suggestions & support

Incumbent API

- MPI-1 library support
- MPI-2 library support
- SHMEM
- OpenMP

Back to Basics : Proposed New Architecture for Many-Core

- Moore's Law Result Shift : Clock → Many-Core
- Problem Shift : Scaling Software Applications
- Load Balancing + Communications Costs
- New Architecture : Communications Focused

PGAS
+
**Programmable
Communications
Accelerators**

**Scalable Commodity Cluster
but with all its memory PGAS**

