

Problem: Analytic Performance Model for Distributed-Memory Applications

Distributed Runtime Model

Lock-step Synchronicity Assumption Failure

Execution Models

Fastest, Fast, Less Fast, Slow

Roofline

Communication Models

LogP

Hockney

Demote load imbalances and network contention

Application: Computation and Communication Characteristics (Nomenclature)

Influence Factors: Hierarchical Topology (affinity matters): Node-level: ccNUMA SMP nodes with multi-socket, multi-core CPUs of Hyperthreads and SIMD cores. Network-level: multi-routes of different characteristics (like latency and asymptotic bandwidth), multi-ways of handling interconnect protocols for intra-chip, inter-chip/node.

Delay Flow Graphs: Basic Flavors and Flow Speed on Silent Systems

Small Message 16 KB (Eager Protocol)

Large Message 132 KB (Rendezvous Protocol)

Communication Topologies

Environmental Setup to Demote:

- Execution Delays (Compute-bound vdivpd, fixed cy/inst., no data transfer boundness)
- System Noise (Single, fully nonblocking; IB/OP leaf switch)
- Communication Delays (PPN=1 → no intranode comm. full non-blocking bandwidth with small Inter-node comm. overhead → no network contention)

Propagation Speed

$$v_{silent} = \frac{\sigma \cdot d}{T_{exec} + T_{comm}} \left[\frac{ranks}{s} \right]$$

σ : { 2, if Bidirectional; R, B/SNB mode; 1, else }

T_{comm} : Runtime of comm. happening in delay flow

d : distance between communicating processes

Future Work

Network Topology

- Complex Applications with Advanced P2P and Collective Communication Patterns
- Multidimensional Structured and Unstructured Grids

Analytical Modeling

Software Development

Analytical-Model Based Simulator for Chip-level Bottleneck and System Topology

References

- A. Afzal et al. Propagation and Decay of Injected One-Off Delays on Clusters: A Case Study. Accepted for IEEE Cluster 2019, Albuquerque, NM, September 23-26, 2019. Preprint: arXiv:1905.10603.
- T. Hoefler et al. LogGOPSim - Simulating Large-Scale Applications in the LogGOPS Model. DOI: 10.1145/1851476.1851564.
- S. Markidis et al. Idle waves in high-performance computing. DOI: 10.1103/PhysRevE.91.013306.
- IB. Peng et al. Idle period propagation in message-passing applications. DOI: 10.1109/HPCC-SmartCity-DSS.2016.0134.
- J.D. McCalpin et al. Memory bandwidth and machine balance in current high performance computers. IEEE computer society technical committee on computer architecture (TCCA) newsletter, 2(19-25), 1995.

Global Non-synchronism Phenomenology: Noise-assisted Long-distance Correlations and Structure Formation

Communication Variations

Network Bandwidth Bottleneck:

Execution Variations

Memory Bandwidth Bottleneck:

STREAM Triad MPI Process Parallelism per Core

CR, d=±1, E Mode, NB, CF

$$v_{silent} = \frac{\sigma \cdot d}{T_{exec} + T_{comm}} \left[\frac{ranks}{s} \right]$$

Variations: Random/Periodic with System Topology/ Application Time step

System Noise Characteristics

PPS=2, PPS=4 (=k), PPS=6, PPS=8, PPS=10 (=n)

Synchronism

$T_{exec} < T_{idle}$ ↔ **Partial Non-synchronism** ↔ $T_{exec} > T_{idle}$ ↔ **Perfect Non-synchronism**

Motivation: Model the Mystery of Non-synchronism

Distributed Runtime Model

Execution Runtime Models

Roofline ◯ ECM

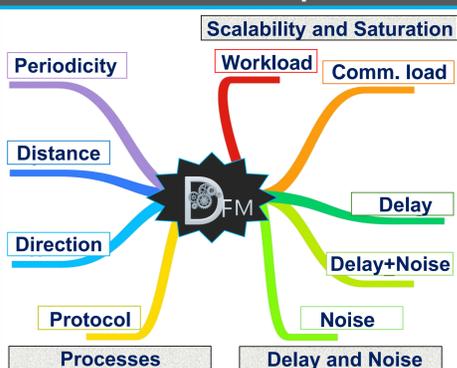
Communication Runtime Models

Hockney ◯ LogP

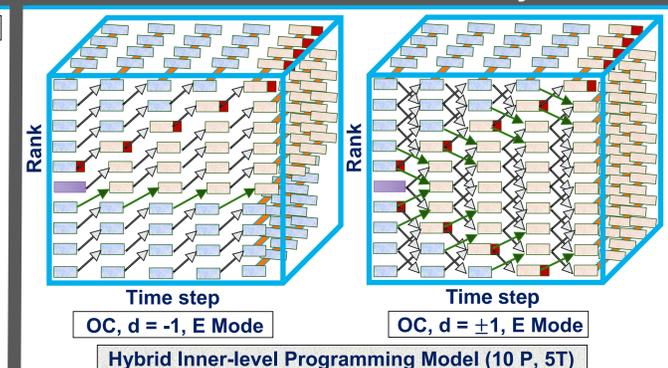
Non-synchronicity Mechanism

Accelerator ◯ Slowdown

Parameter Space



Communication-assisted Delay Flow



Interaction & Noise-assisted Damping + Elimination of Delays

	Lock-step Synchronism	Non-synchronism
Execution-comm. Overlap	Perfect Non-overlap	Partial or Perfect Overlap
Memory Bus/Network Utilization	Intermittent	More Continuous
Processes in Saturated Resource(s) Regime	P = n	P < n
Resource(s) Bandwidth Per Process	Less (contented)	More/Full (less/no contented)
Execution time, $T_{exec}^{BW(PPS)}$	Slowest	Faster
Idle time, $T_{idle}^{BW(PPS)}$	Minimum	Lengthen
Idlewave Decay	Weak	Strong

Wall-clock Time [s]

OC, d = ±1, R Mode, NB **CR, d = ±1, R Mode, NB**

Triad (2 MPI PPS * 5 OpenMP T, V_{mem} = 4.8 GB, Msg = 29 MB, b_{IB} ≈ 3 GB/s)

KONWIHR OMI4papps

Workload	Periodicity for Boundary Conditions	Direction	Distance of Communicating Processes	Communication Protocol	Communication Routine	Communication Load
◻ : Compute Bound	OC : Open Chain ◻ ◻ ◻ ◻	+ or - : Uni directional	d = 1 : Direct Neighbour Processes	E : Eager	B : Blocking (sendrecv, sendrecv,...)	CF : Contention Free
◻ : Memory Bound	CR : Close Ring ◻ ◻ ◻ ◻	± : Bi directional	d = 2, 3, ..., n : Indirect Neighbour Processes	R : Rendezvous	NB : Non-blocking (isend, irecv, isend, irecv, ..., waitall)	CF : Contention in Network
					SNB : NB + split wait (isend, irecv, waitall, isend, irecv, waitall, ...)	