

# Analytic Modeling of Idle Waves in Parallel Programs: Communication, Cluster Topology, and Noise Impact

Supervisor: Jannek Squar

Phuong Thao Le Presentation date: 11.01.2022

Analytic Modeling of Idle Waves in Parallel Programs



- 1. Introduction
- 2. MPI Message Passing Library Interface Specification
- 3. Test bed and Experimental Methods
- 4. Related Works
- 5. Idle Waves Propagation Velocity for scalable code
- 6. Idle Waves interacting with MPI Collectives
- 7. Idle Waves Decay
- 8. Summary
- 9. References



# **1. Introduction**

Due to the need to compute enormous computer resources (time, memory).

- → Parallel Computing:
  - Solving a task by simultaneous use of multiple processors, which are components of a unified architecture.
  - There might be disturbances such as system and network noise, delays caused by one-off events, etc.





# **1. Introduction**

What are idle waves?

- Long one-off delays from one process cause periods of idleness, which later ripple through the system and affect other processes.
- Idleness: a state of delay, where the process needs to wait for information.

Consequences:

- Delay performance of the application
- Desynchronization among processes
- $\rightarrow$  Automatic communication overlap



# **1. Introduction**

#### **Goals of the study:**

- Analytic modelling of the propagation speed of Idle waves in scalable code with respect to:
  - Communication topology
  - Communication concurrency.
- Interaction of idle waves with MPI collectives.
- Initiation of idle wave decay and analyse of decay rate:
  - Topological Decay
  - Noise-induced Decay



#### **Interprocess communication:**

- Processes executing concurrently in the operating system may be independent or cooperating processes.
  - Independent: Processes not affecting each other
  - Cooperating: Processes share data with each other and thus have effects on each other.
- Advantages of cooperating processes: Information sharing, computation speedup, convenience, modularity.



- Cooperating processes require an interprocess communication (IPC) mechanism that allows exchanging data and information.
- Two fundamental models:
  - Shared memory
  - Message passing









# **Shared Memory**: Exchanged Data are written and read on the shared memory region.

#### Message Passing:

Communication takes place by means of messages exchanged between processes.



#### **MPI : Message-passing library interface specification**

- A standardized and portable message-passing standard designed to function on parallel computing architectures.
- Useful in distributed memory environment, where communication processes located on different computer nodes.
- Language bindings: C, C++, Fortran



## urg **2. MPI**

#### - Two main operations:

- Send(messages) and Receive(messages)

Point to Point communication



Collective functions





- Synchronous or asynchronous communication:
  - Synchronous (Blocking) : Send and receive must be completed before conducting another process
  - Asynchronous (Non-Blocking) : Computation and communication can be conducted at the same time (MPI\_Wait, MPI\_Test to see if the communication is finished)
- Rank is a unique identifier for each processor.



- Three clusters Emmy, SuperMUC-NG and Hawk
- Process-core affinity was enforced and multiple characteristics are adjusted in order to conduct the study.
- Process scalability: latency-bound communication and compute-bound workload
- One-off idle periods are deliberately generated by massively expanding one computational phase via doing extra work on one random MPI process, usually rank 5 in this study.



- Barrier free bulk synchronous parallel programs
- Distributed-memory parallel system using one MPI process per contention domain (typically ccNUMA)





## 4. Related works



**Figure 2:** Physical mechanism for the generation of idle waves in computation in 9 processes (2) **Figure 3:** The delay propagation mechanism in the most simple setting (3)



#### **5.1. Execution Characteristics**

- Traditional memory-bound algorithms such as stencil updates or SpMV with one MPI process per contention domain (ccNUMA node)
- In-core workload



### **5.2. Categorization of communication characteristics**

- Assuming P2P Homogeneous Situation

**Communication Topology:** A consequence of the physical problem underlying the numerical method and of the algorithm(discretization, geometry).

- Compact Topology:
  - Each process communicates with a dense, continuous array of neighbours with distances d (±1, ±2, ±3, .... ±12)
- Noncompact Topology:
  - Each process communicates with processes that are not continuous block d (±1,±12)





Figure 4: Communication topology with bidirectional open chain characteristics (1)



#### **Communication concurrency:**

- The number of P2P communications are grouped and subject to completion via MPI\_Waitall.

Cartesian Communicator example:

Dim [1] = 4

rank	$\begin{array}{c} 0 \\ (0,0) \\ 8/4 \end{array}$	1	2	3
(row, column)		(0,1)	(0,2)	(0,3)
source/des		9/5	10/6	11/7
	4	5	6	7
	(1,0)	(1,1)	(1,2)	(1,3)
	0/8	1/9	2/10	3/11
Dim[0] = 3	8	9	10	11
	(2,0)	(2,1)	(2,2)	(2,3)
	4/0	5/1	6/2	7/3

**Figure 5:** Example Cartesian communicator (6)



#### Assuming that all P2P communication is nonblocking



 $\stackrel{+}{\downarrow} P_i$  send to  $P_{i+dir \times d}$ ; §  $P_i$  receive from  $P_{i-dir \times d}$ 

Table 1: Selected algorithms for communication concurrency in the MPI Microbenchmarks (1)



#### **5.3. Analytical Model of Idle Wave Propagation**

- Propagation speed of an idle wave is the speed, in ranks per second, with which it ripples through the system.
- Restriction: Open boundary conditions across the MPI Processes
- $\rightarrow$  Affects the survival time and not the propagation speed of the wave



#### **Corner Cases:**

Minimum: Min Speed, Max Survival time

$$v_{\rm silent}^{\rm min} = 1 \left[ \frac{\rm ranks}{\rm iter} \right] \times \frac{1}{T_{\rm exec} + T_{\rm comm}} \left[ \frac{\rm iter}{\rm s} \right].$$

- Simple direct next-neighbour communication ( d = 1 )
- $T_{exec}$  and  $T_{comm}$ : Execution and communication times of one iteration of the bulk-synchronous program
- $\rightarrow$  The wave survives until system boundaries.



**Maximum:** Max Speed, The wave dying out quickly in a minimum of one time step.

$$v_{\text{silent}}^{\text{max}} = \alpha \left[ \frac{\text{ranks}}{\text{iter}} \right] \times \frac{1}{T_{\text{exec}} + T_{\text{comm}}} \left[ \frac{\text{iter}}{\text{s}} \right],$$
$$\alpha = \max \left( \text{size}_{\text{comm}} - r_{\text{inject}} - 1, r_{\text{inject}} - 1 \right).$$

 $r_{inject}$ : The rank where the idle wave originated  $size_{comm}$ : The total number of MPI processes



#### **Multi-neighbour Communication:**

$$v_{\text{silent}} = \kappa \cdot v_{\text{silent}}^{\min} \left[ \frac{\text{ranks}}{\text{s}} \right]$$

*k* is the distance (in processes) travelled by the wave in one time step (depending on communication topology and concurrency).

*j* is the longest-distance communication partner of a process.

$$\kappa = \begin{cases} \sum_{k=1}^{j} k = \frac{j(j+1)}{2} & \text{if compact MWSDim/MWSDir/blocking} \\ \sum_{k=1,j} k = j+1 & \text{if non-compact MWSDim/MWSDir/blocking} \\ j & \text{if MWMDim/SWMDim} \end{cases}$$



## **5.4. Experimental validation**

#### Microbenchmarks:

- One-off idle wave originated at rank 5, dark blue
- Execution time = 13ms (light blue) and a data volume of 1 KiB per message.



#### **Compact Communication:**



Figure 6: Validation for compact communication (1)



### **Compact Communication:**

- Propagation speed of idle waves is independent of the number of split-waits
- Higher speed:
  - Communication distance goes up with increasing number of communications of communication partners
  - Number of dimensions spanned within each MPI\_Waitall
- Higher  $k \rightarrow$  higher speed  $\rightarrow$  shorter survival time.
- Slower wave propagation  $\rightarrow$  more waiting time
  - $\rightarrow$  Resource utilization across processes.



#### **Noncompact Communication**



**Figure 7:** Validation for noncompact communication (1)



#### **Noncompact Communication**

- Difference between propagation speed of the "fast waves" and that of the "secondary waves" → Only fast waves remain.
- Higher speed of residual waves:
  - A larger number of split-waits
  - A smaller number of communication dimensions spanned by each MPI\_Waitall
  - A larger longest communication distance j
- The zigzag pattern dies out for MWSDim and MWMDim but remains for SWMDim



#### **Heterogeneous Communication**



**Figure 8:** Idle wave propagation with heterogeneous compact communication characteristics (a) Topology matrix; (b) Idle wave propagation for SWMDim concurrency. (1)



# 6. Idle Waves Interacting with MPI Collectives

- MPI supports Point-to-Point communication and collective functions for communication among multiple computer nodes
- Not all MPI collective routines eliminate a traveling idle wave, some may be permeable to it, depending on the implementation.
- Restriction: Intel MPI on Emmy



# 6. Idle Waves Interacting with MPI Collectives

- Globally synchronizing Primitives
- $\rightarrow$  Destroy idle waves completely
- Globally Non-synchronizing Primitives
  - MPI\_Reduce
  - MPI\_Gather
- Implementation Variants
  - I\_MPI\_ADJUST\_<opname>



# Figure 9: Interaction of idle waves with MPI Collectives (1)



Idle wave decay is the phenomenon where idle waves are damped in time.

#### 7.1. Topological Decay

- Three benchmarks Hawk, Emmy, SuperMUC-NG have different features in respect to system topology.
- Communication heterogeneities → Fine-grained noise → Idle wave decay.
- Decay rate: Emmy has the strongest effect then Hawk and lastly SuperMUC-NG





Figure 10: Topological idle wave decay on the benchmark systems. (1)



#### 7.2. Noise-Induced Decay

- Fine grained noise effects on idle wave decay with resourcescalable code
- Noise eliminates the trailing edge of the wave
- A small idle period of duration  $T_{noise}$  shortens the next by the exact amount of  $T_{noise}$
- →Cumulative process
- Only noise power and not the noise characteristics that has an impact on noise-induced decay rate.





**Figure 11:** Experiment comparing the average decay rate of an idle wave for 2 different noise characteristics (1)



#### 7.3. Experimental validation



Figure 12: Decay rate of an idle period in s/rank, comparing 3 different noise patterns (1)



- Analytic model of idle wave propagation speed based on communication topology and concurrency of resource-scalable MPI.
- MPI collectives can be permeable to idle waves depending on which collectives we use and how we implement or adjust them.
- System Topology → Fine-grained noise → Impact on idle wave decay rate.
- Only noise power and not noise characteristics has an impact on noise-induced decay rate.



- Afzal, A., et al: Analytic Modeling of Idle Waves in Parallel Programs: Communication, Cluster Topology, and Noise Impact (2021)
- (2) Markidis, S., et al: Idle Waves in high-performance computing. (2015)
- (3) Afzal, A., et al: Propagation and decay of injected one-off delays on cluster: a case study. (2019)
- (4) Hager, G., Wellein, G. : Introduction to High Performance Computing for Scientists and Engineers. (2010)
- (5) MPI: A Message-Passing Interface Standard Version 3.1 (2015)



(6) Prof.Dr-Ing Riedel, Morris: 2021 High Performance Computing
Lecture 2 Parallel Programming with MPI Part1. URL : <u>2021 High</u>
<u>Performance Computing Lecture 2 Parallel Programming with MPI Part1</u>
<u>— YouTube</u>

(7) Neso Academy: Interprocess Communication. URL : <u>Interprocess</u> <u>Communication – YouTube</u>

(8) NHR @FAU : Parallel Programming 2020: Lecture 7 - ccNUMA and wavefront parallelization with OpenMP. URL : <u>Parallel Programming</u> <u>2020: Lecture 7 - ccNUMA and wavefront parallelization with OpenMP -YouTube</u>



(9) Prof.Dr-Ing Riedel, Morris: 2021 High Performance Computing
Lecture 4 Advanced MPI Techniques Part1. URL: <u>2021 High Performance</u>
<u>Computing Lecture 4 Advanced MPI Techniques Part1</u> – <u>YouTube</u>

(10) Graphcore: Fundamentals of Bulk Synchronous Parallel Execution on the IPU

URL: <u>Fundamentals of Bulk Synchronous Parallel Execution on the IPU -</u> <u>YouTube</u>



# Thank you very much for listening !