

Friedrich-Alexander-Universität Erlangen-Nürnberg

The Role of Idle Waves in Modeling and Optimization of Parallel Programs

AYESHA AFZAL^{1,2}, GEORG HAGER¹, GERHARD WELLEIN^{1,2}

 ¹ Erlangen National High Performance Computing Center
 ² Department of Computer Science, University of Erlangen-Nürnberg Germany

NHR PerfLab Seminar April 26, 2022







Vision: white-box first-principle performance modelling





It's intricated (bottlenecks interact, systems are noisy, etc.)

Ayesha Afzal <ayesha.afzal@fau.de> | NHR PerfLab Seminar 2022

Motivation



Pr	rocess 0	T _{exec}	T _{non-exec}	T _{exec}	T _{non-exec}	•••
Process 1		T _{exec}	T _{non-exec}	T _{exec}	T _{non-exec}	
Memory-bound MPI-parallel programs: timeline view		Ν	"Lock-step" behavi To network contention o	ior at start r load imbalance		

Motivation











Wall-clock time [s]

Computational wavefront

302³ lattice cell, 8 GB data set, non-blocking, 1D domain decomposition, distance-1 communication, 10 Emmy@RRZE sockets @2.2 GHz





Lessons learned (1): Impact of idle wave on overlap



Ayesha Afzal <ayesha.afzal@fau.de> | NHR PerfLab Seminar 2022

7

FÂU

NHR

$$v_{min} = \kappa \times \sigma \times 1 \left[\frac{r \tan \kappa}{iter} \right] \times \frac{1}{T_{comp} + T_{comm}} \left[\frac{r er}{s} \right]$$

$$\kappa = \frac{j(j+1)}{2} \text{ or } (j+[i]) \text{ or } j$$

j /*i*: longest /shorter-distance partner

resource-scalable MPI programs

Trank

8

FÅU

NHR

[itor]

Communication pattern and concurrency: compact



1: while $d \le dims$ do2: while $dir \le bi$ do	1: while $dir \le bi$ do 2: while $d \le dims$ do
3: MPI_Isend #	3: MPI_Isend‡
4: MPI_Irecv §	4: MPI_Irecv §
5: end while	5: end while
6: MPI_Waitall	6: MPI_Waitall
7: end while	7: end while



distance in processes travelled in one time step by the idle wave



Communication pattern and concurrency: compact







distance in processes travelled in one time step by the idle wave



Communication pattern and concurrency: heterogeneous





Idle wave propagation: 3D jacobi smoother





Idle wave propagation: 3D jacobi smoother







Receiver rank

Idle wave propagation: sparse MVM with HPCG matrix







Receiver rank

Idle wave propagation: sparse MVM with HPCG matrix







Idle wave propagation: sparse MVM with HPCG matrix





15: MPI_Allreduce; 16: rNorm = sqrt (rNorm);

ement - miniAMR





Lessons learned: Impact of idle wave on overlap









Ayesha Afzal <ayesha.afzal@fau.de> | NHR PerfLab Seminar 2022

April 26, 2022

Lessons learned: Impact of idle wave on overlap







Non-linear interactions of idle wave





Non-linear interactions of idle wave with system topology





Lessons learned: Impact of idle wave on overlap





Lessons learned: Impact of idle wave on overlap





Scalable versus contented processes



Ayesha Afzal <ayesha.afzal@fau.de> | NHR PerfLab Seminar 2022

FAU

NHR

Scalable versus contented processes



SuperMUC-NG @2.3 GHz, non-temporal stores, bi-dir, 1024 B, close chain, distance-1 communication

NHR



Ayesha Afzal <ayesha.afzal@fau.de> | NHR PerfLab Seminar 2022

FAU

Lessons learned: Impact of idle wave on overlap









Algorithmic dependency: some collectives can be permeable to idle waves



4: tmp = 0.0; 5: for $idx = rowPtrA[rcc6: tmp += valA[idx7: end for Listi1: dx2: in3: fo4: tmp = 0.0;5: for idx = rowPtrA[row] += tmp;2: in3: fo4: tmp = 0.0;5: for idx = rowPtrA[row] += tmp;5: for idx = rowPtrA[row] : rowPtrA[row+1] - 1 do6: tmp += valA[idx] * x[colldxA [idx]];7: end for8: b[row] += tmp;9: end for$	<pre>bw] : rowPtrA[row + 1] - 1 do] * x[colIdxA [idx]] ; 4: local_spMVM (A, x, b); 5: MPI_Wait ; * 6: remote_spMVM (A, x, b); 7: MPI_Barrier; 8: swap (b, x); 9: end while</pre>	4: local_spMVM (A 5: MPI_Wait; * 6: remote_spMVM 7: MPI_Barrier List 8: swap (b, x); 1: 9: end while 4: MPI_Wait; * 5: local_spMVM (A, x, b); 6: remote_spMVM (A, x, b); 7: MPI_Barrier; 8: swap (b, x); 9: end while	A, x, b); (A, x); (A, x);
* Two MPI_Wait routines wait for both MPI receiv	e and send requests to complete.	t to top!	1279 0 500 12790 500 1279
Matrix-orderBandwidth n_{elec} HHQ-large- pe highHHQ-large- ep lowHHQ-small- pe highHHQ-small- ep low	$\begin{array}{c} {}_{\rm trons} - n_{\rm sites} - n_{\rm phonons}{}^{\rm s} n_{\rm r} = n_{\rm c}{}^{\rm s} \\ \hline 3 - 8 - 10 & 60988928 \\ 3 - 8 - 10 & 60988928 \\ 6 - 6 - 15 & 6201600 \\ 6 - 6 - 15 & 6201600 \end{array}$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	(a) pe order (b) ep order 200
Phase vs. Rank-order 96-pe 144-pe 240-pe	480-pe 720-pe 960-pe 1296-pe 96-ep 144	-ер 240-ер 480-ер 720-ер 960-ер 1296-ер	
Exec median [ms] 53.5 35.74 18.51 Comm median [ms] 29.48 24.56 16.17	9.05 6.53 5.06 3.85 48.47 30 14.75 11.28 9.78 7.99 16.17 15	.87 17.57 8.23 5.7 4.54 3.03 .01 14.41 12.17 8.96 6.62 6.13	$1279 \begin{bmatrix} & & & \\ & & & & \\ & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & &$
CER median 0.55 0.69 0.87	1.63 1.73 1.93 2.08 0.33 0.	49 0.82 1.48 1.57 1.46 2.02	Receiver rankReceiver rank(a) pe order(b) ep order





NHR

FAU



Local (global) domain in each dimension

Large amount of time spent in the MPI library is not not set in the multiple of the set of the set



Pseudo implementation of parallel LBM

1: while iter \leq nIters do

- 2: stream_collide_update (lattice, u_lid, omega);
- 3: set_boundary_condition (u_lid);
- 4: MPI_Isend;*
- 5: MPI_Irecv;*
- 6: MPI_Wait;
- 7: ghost_cells_update ();
- 8: **if** ((iter % collective_step) == 0) **then**
- 9: MPI_Allreduce;
- 10: **end if**
- 11: swap (local_src_lattice, local_dst_lattice) ;
- 12: end while



Speedup with higher automatic overlap

NHR FAU

Performance: MPI parallel LULESH proxy application



Performance: Hybrid parallel Chebychev Filter Diagonalization





(a) TOPI-EHN

Ayesha Afzal <ayesha.afzal@fau.de> | NHR PerfLab Seminar 2022

Performance: Hybrid parallel Chebychev Filter Diagonalization



B M1-NON-SPLITŮ Ů M1-SPLITŮ Ů M1-PIPELINEŮ Ů M2-NON-SPLITŮ Ů M2-SPLITŮ Ů M2-PIPELINEŮ ® M3-NON-SPLITŮ Ů M3-SPLITŮ Ů M3-PIPELINE ■ M4-NON-SPLITŮ Ů M4-SPLITŮ Ů M4-PIPELINEŮ Ů M5-NON-SPLITŮ Ů M5-SPLITŮ Ů M5-PIPELINEŮ Ů M6-NON-SPLITŮ Ů M6-SPLIT ■ M7-NON-SPLITŮ Ů M7-SPLITŮ Ů M7-PIPELINEŮ Ŵ M8-NON-SPLITŮ Ů M8-SPLITŨ Ů M8-PIPELINEŮ № M9-NON-SPLITŮ Ů M9-SPLIT



Better overlap for decomposition of slow ide wave and large communication overhead

WITH BARRIER

Pipeline mode better (non-split suffer more for more saturating case)

WITHOUT BARRIER

Non-split mode on-par with pipeline mode or even better for more saturating case

Overlapping via explicit programming techniques may not be necessary for strongly bandwidth-saturating code with large (but not dominant) communication overhead due to the presence of natural overlap by desynchronization







DisCostiC: A DSL-based Parallel Simulation Framework

Using First-Principles Analytic Performance Models





TitleThe Role of Idle Waves in Modeling and
Optimization of Parallel Programs

Contact

Ayesha Afzal ayesha.afzal@fau.de

Acknowledgement

Georg Hager, Gerhard Wellein

