

POP methodology: Identifying Performance Bottlenecks through Efficiency Metrics

Judit Giménez (BSC)

HORIZON-EUROHPC-JU-2023-COE



1 January 2024–31 December 2026

Grant Agreement No 101143931



...with a perfect scaling?



- Strong scaling
- Unbalanced simulation

What about the reference case?



Measuring the efficiency of a run



- Model the application run on a given platform (app + SW stack + HW)
- Run allocation: using a set of resources during a time interval



Parallel run **efficiency** = %useful useful = time spent with the computations (outside parallel runtime)

- Reasons to be inside MPI parallel runtime (potential bottlenecks)
 - Data transfer
 - But also
 - Work unbalance
 - Temporal unbalances and serializations



Measuring the scaling efficiency



- Scaling of the parallel efficiency (and its components)
 - Is increasing/reducing the % useful time?
 - If it is reduced, which factor increases its weight?
- Scaling of the **computations**
 - There is code replication (increase of instructions)?
 - Does the scale affect the IPC?
 - Strong scaling vs. weak scaling

All factors can be expressed as a percentage or a positive number

• always the higher the better



Why is my code scaling?



CG-POP mpi2s1D - 180x120









POP's efficiency metrics



• The POP methodology defines hierarchies of metrics (<u>pop-coe.eu/sites/</u> <u>default/files/pop_files/metrics.pdf</u>)

- Global Efficiency (GE)
 - Parallel Efficiency (PE)
 - Load Balance Efficiency (LB)
 - Communication Efficiency (CommE)
 - Serialization Efficiency (SerE)
 - Transfer Efficiency (TE)
 - Computation Efficiency (CompE)
 - IPC Scaling
 - Instruction Scaling
 - Frequency Scaling







- Looking at the reported efficiencies:
 - Above 0.8 may be considered good
 - Execution vs. Scaling
 - These numbers will drive the analysis



	- 100	192	96	48	24	
Lower value for the	- 100	34.83	44.96	54.66	68.29	-
smallest core count and	- 80	42.52	46.06	57.82	68.29	-
main factor that limits	- 60 - 60 -	44.04	46.62	58.09	68.45	-
scalability	ntag	96.55	98.79	99.53	99.75	-
	- 40 eco	79.56	97.61	94.54	100.00	-
	- 20	95.02	95.57	98.63	100.00	-
	0	89.44	94.83	97.69	100.00	-
		7				

Global efficiency

-- Parallel efficiency

- -- Load balance
- -- Communication efficiency

-- Computation scalability

-- IPC scalability

-- Instruction scalability





Efficiency model analysis – examples



Other factors limiting scalability

- Data transfer
- Code replication



100.00 101.57 102.04 103.68 -- Computation scalability 100.00 103.63 110.81 115.60 -- IPC scalability

-- Instruction scalability

-- Communication efficiency -

-- Serialization efficiency

-- Transfer efficiency

Global efficiency

-- Parallel efficiency

-- Load balance

-- Frequency scalability

Efficiency model analysis – examples

384

94.39

91.03

97.04

93.81

97.72

96.00

89.60

100.10

768

88.23

88.42

93.02

95.06

99.20

95.83

99.79

119.74

83.47

99.84

1536

84.71

85.27

92.83

91.86

97.34

94.37

99.35

124.53

79.95

99.78

3072

78.86

81.20

92.84

87.47

93.97

93.08

97.12

127.48

76.59

99.47

- 60

- 40

- 20

· 0

96

95.95

94.46

96.50

97.89

99.06

98.82

98.02

99.99

48

97.58

97.58

98.30

99.26

99.78

99.49

100.00

100.00

192

95.50

93.59

98.33

95.18

97.92

97.20

92.11

99.97



Smooth

Main factor that will limit scalability (IPC cannot grow forever)

Very good efficiency with the smallest core count

Improvement on IPC compensates increase of instructions



Selecting FoA and understanding why

Parallel efficiency 93.28 Communication eff. 93.84

Parallel efficiency 77.93 Communication eff. 79.79

Parallel efficiency 28.84 Communication eff. 30.42



Understanding why – LB efficiency





Ex 2: MPI

2DH useful duration correlated with @ sem46P4_small_N02_FWI_all_waves.chop1.prv		2dh useful instructions @	sem	46P4	₽_sm	all_M02_F	lI_a	u_	waves.chop1	.prv
				_	-	_	_	_	-	
					-				-	
			_	·				-	-	·
			<u> </u>	_	_	<u> </u>	_	-		
			_	_	_			-		
			_	_		_		•••		
			· · ·	<u> </u>	_		_	-		
			-	_	_	-	_	-		
			_	_	_			÷.,		
			<u> </u>	<u> </u>	_	_	_	_		
			-	_	_			-		
			_	_	_	_	_	<u> </u>		
			· · ·	<u> </u>	<u> </u>	_		-		
				_	_	_		-		
				_	_			۰.		
					-	_		_		
					_	_		-		
					_		_			
	F							-		-
				_	_					
				_	_	_		1		
			- C			_				
				-			<u> </u>	_		
							<u> </u>			
	-									
							-			
	-							-		
				_	-		_			
							т			
	-		_				-	т		
	-									
	-				_			-		
	-			-						
	-									
	-				-					
			_	_						
	-				_			-		
	_									
				_	-					
						-		_		
						-		-		
	-		-	-	•		-	-		
				_	_			-		

duration

instructions

Regions not parallelized with OpenMP

Structured instructions unbalance



Understanding why – transfer eff.

Ex 1: Load Balance 76% Serialization 99%, Transfer 84%



With instantaneous communications

- Most of the point to point calls (red) almost disappear → data transferred in point to point is limited by network resources
- The all reduce (pink) increases a bit \rightarrow absorbs previous load imbalance

Understanding why – serialization eff. 2000

Ex 1: Audited version (left) vs. improved version (right)



 94,230 us
 142,744 us

MPI calls



Serialization detected as stepped computations \rightarrow running these phases in parallel obtain a reduction of 22% in the region.



Understanding why – IPC scaling



Ex 1: OpenMP



IPC reduction is on almost all the phases

Ex 2: MPI



Good one: light green region

3 regions suffer IPC degradation (arrows to the left)

[4 regions have code replication]





Efficiency model extensions



15

Applying the original model to hybrid codes 2000



- Efficiency computed as percentage of time outside the two parallel runtimes
- Useful as a first step to distinguish between Load Balance and Communication. But how to dig down?
 - Efficiencies are mixing inefficiencies from MPI and OpenMP → Need to distribute the blame between the two programming models



What do we want?





- Global load balance and communication concepts can be mapped to any parallel programming paradigm
- The efficiencies at hybrid level collapse the contributions from the two programming models

17

A very simple case





• From the MPI point of view, OpenMP runtime is as useful as computation



Whatever cannot be blamed on MPI is caused by OpenMP



OpenMP efficiencies [%] Parallel efficiency 49.05 Load Balance 52.69 Communication 93.11

18

Same approach applied to MPI+CUDA



- Bigger impact from CUDA component with a similar contribution in all the scales
- MPI contribution increases with the scale, but still lower than CUDA



A different approach for MPI+CUDA



 Some codes use efficiently the GPUs but intentionally do not use the CPUs → lower metrics focusing on the GPUs only

Global Parallel efficiency - 47.43 47.19 46.80 46.27 45.73 44.48 43.61 43.30 42.34 42.04	
	100
Global Load balance - 50.00 49.91 49.97 49.91 49.92 49.83 49.59 49.64 49.59 49.40	80
GPUs Global efficiency - 94.85 94.08 92.94 91.03 89.33 85.99 83.83 82.77 80.17 79.00	• 60 (%)
GPUs Parallel efficiency - 94.85 94.33 93.53 92.39 91.23 88.67 86.87 86.18 84.19 83.52	ntag
GPUs Load Balance - 100.00 99.79 99.86 99.65 99.59 99.33 98.77 98.80 98.62 98.62 98.15	- 40 au
GPUs Communication efficiency - 94.85 94.53 93.66 92.71 91.61 89.27 87.95 87.23 85.37 85.09	⊷ 20
GPUs Computation scalability - 100.00 99.73 99.37 98.53 97.92 96.98 96.50 96.05 95.22 94.59	0

- Considering all the resources allocated, the efficiencies are around 40-50% because the work is concentrated in the GPUs (reported as unbalance). The small degradation in the parallel efficiency is also observed in the GPUs metrics.
- The GPUs global efficiency reports a degradation with the scale that is related mainly with the communication efficiency but also with the computation scalability. The load balance between GPUs is ok for all the scales.

How to calculate the metrics?



- From Extrae/Paraver data
 - BSC's Basic Analysis module metrics (tools.bsc.es/downloads)
- From Score-P traces
 - Since CUBE version 4.5 (<u>www.scalasca.org/software/cube-4.x/download.html</u>)
- This training will show you how to use them!



Performance Optimisation and Productivity A Centre of Excellence in HPC





This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 676553 and 824080.