



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Predicting applications performance far far away

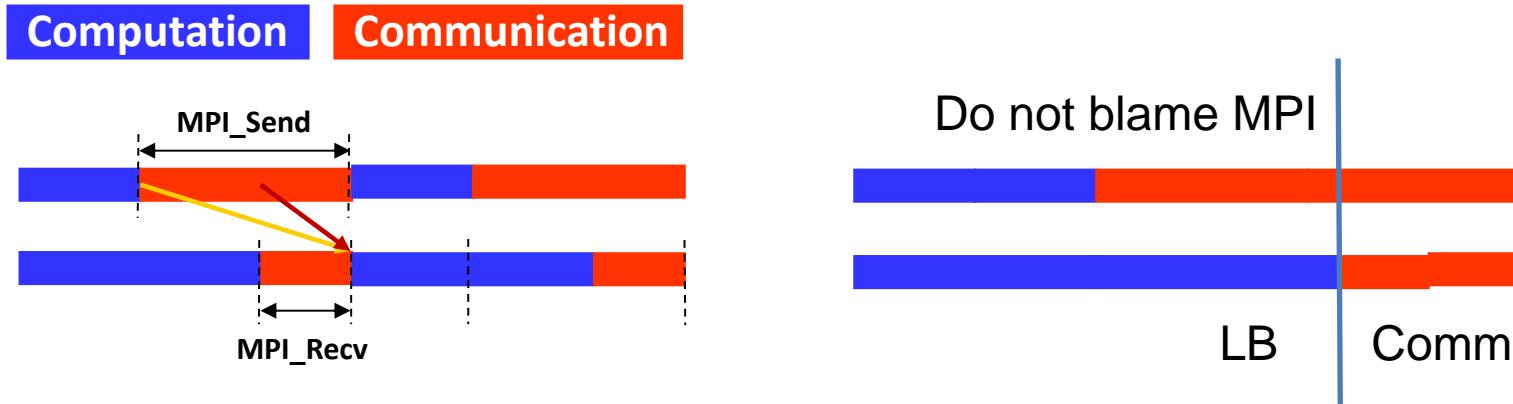
Judit Gimenez (judit@bsc.es)

Scalable Tools Workshop (Tahoe)
August 2015

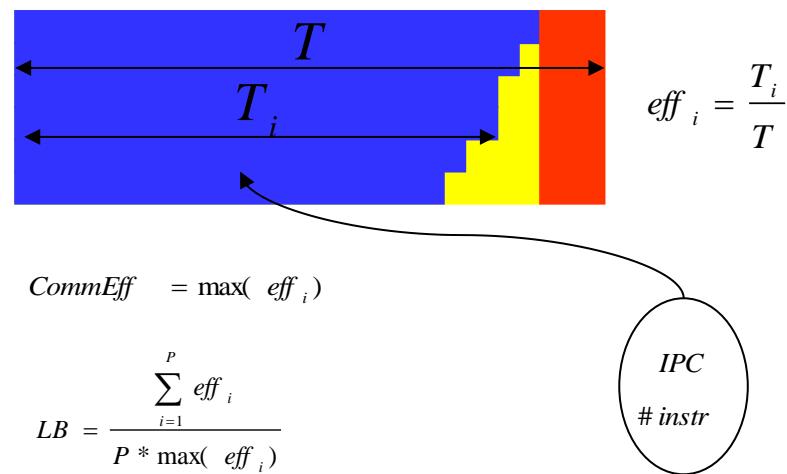
Our believe

Traces for small core counts express the symptoms that would produce performance degradation at larger scale even they have no impact

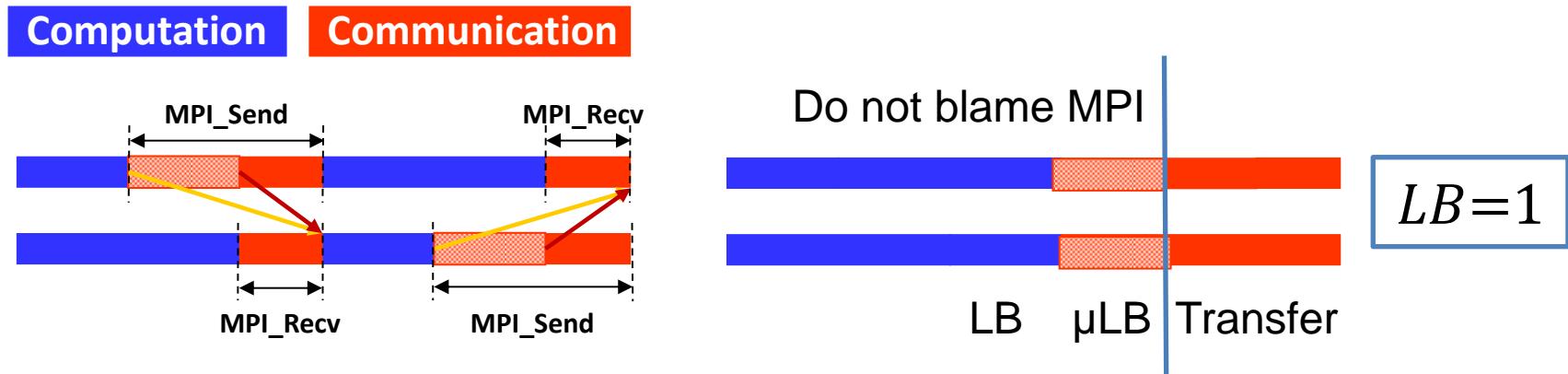
Parallel efficiency model



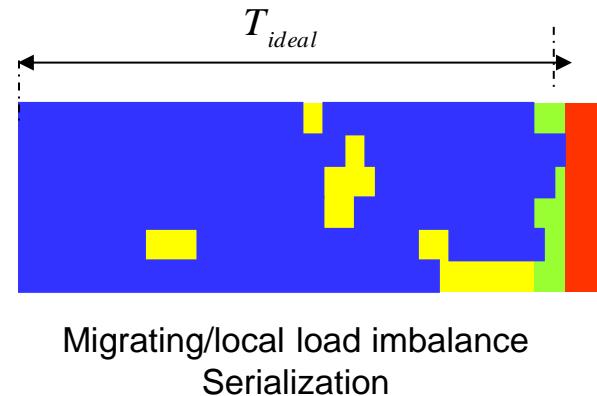
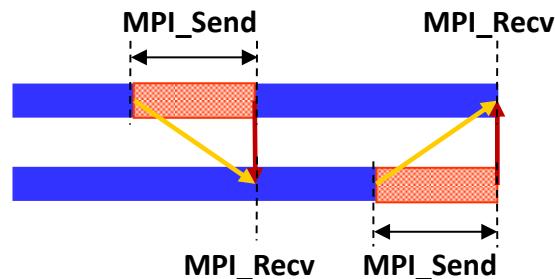
Parallel efficiency = LB eff * Comm eff



Parallel efficiency refinement: LB * μ LB * Transfer



- « Serializations / dependences (μ LB)
- « Dimemas ideal network → Transfer efficiency = 1



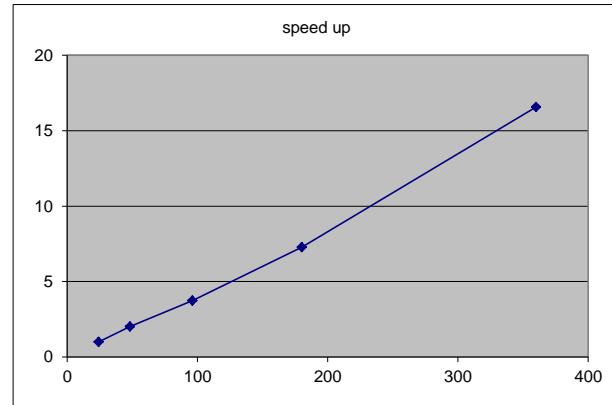
$$\mu LB = \frac{\max(T_i)}{T_{ideal}} \quad \text{Transfer} = \frac{T_{ideal}}{T}$$

Why scaling?

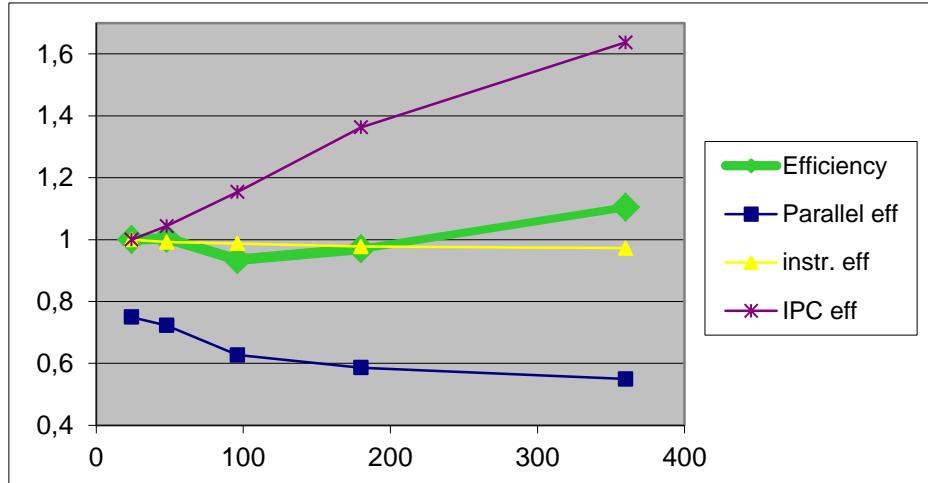
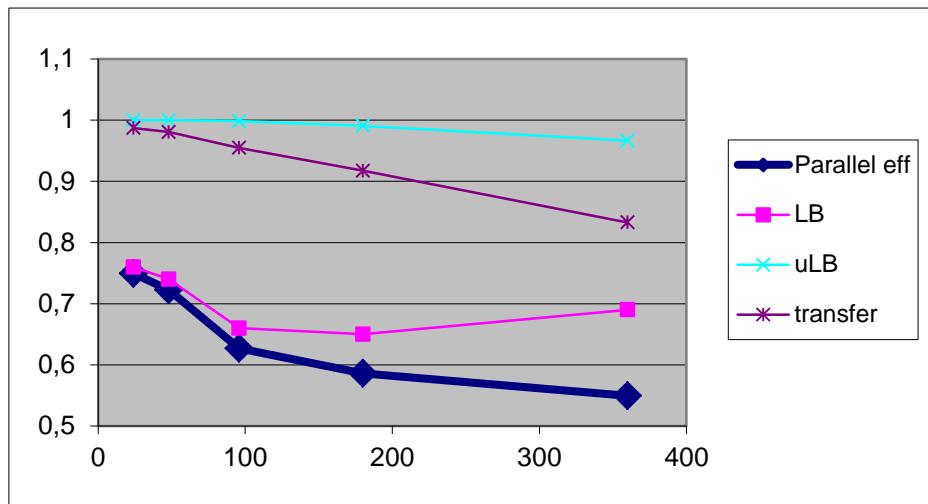
$$\eta_{\parallel} = LB * Ser * Trf$$

CG-POP mpi2s1D - 180x120

Good scalability !!
Should we be happy?



$$\eta = \eta_{\parallel} * \eta_{instr} * \eta_{IPC}$$



Scalability prediction

Efficiency extrapolation

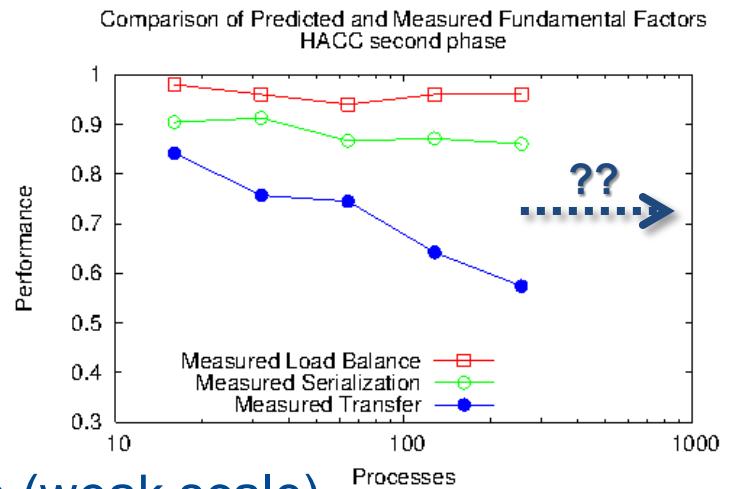
- From few executions @ low core counts

Fit based on

- Reasonable fundamental behavioral models (e.g. Amdahl)
- Guided by observed internal application behavior (tools)

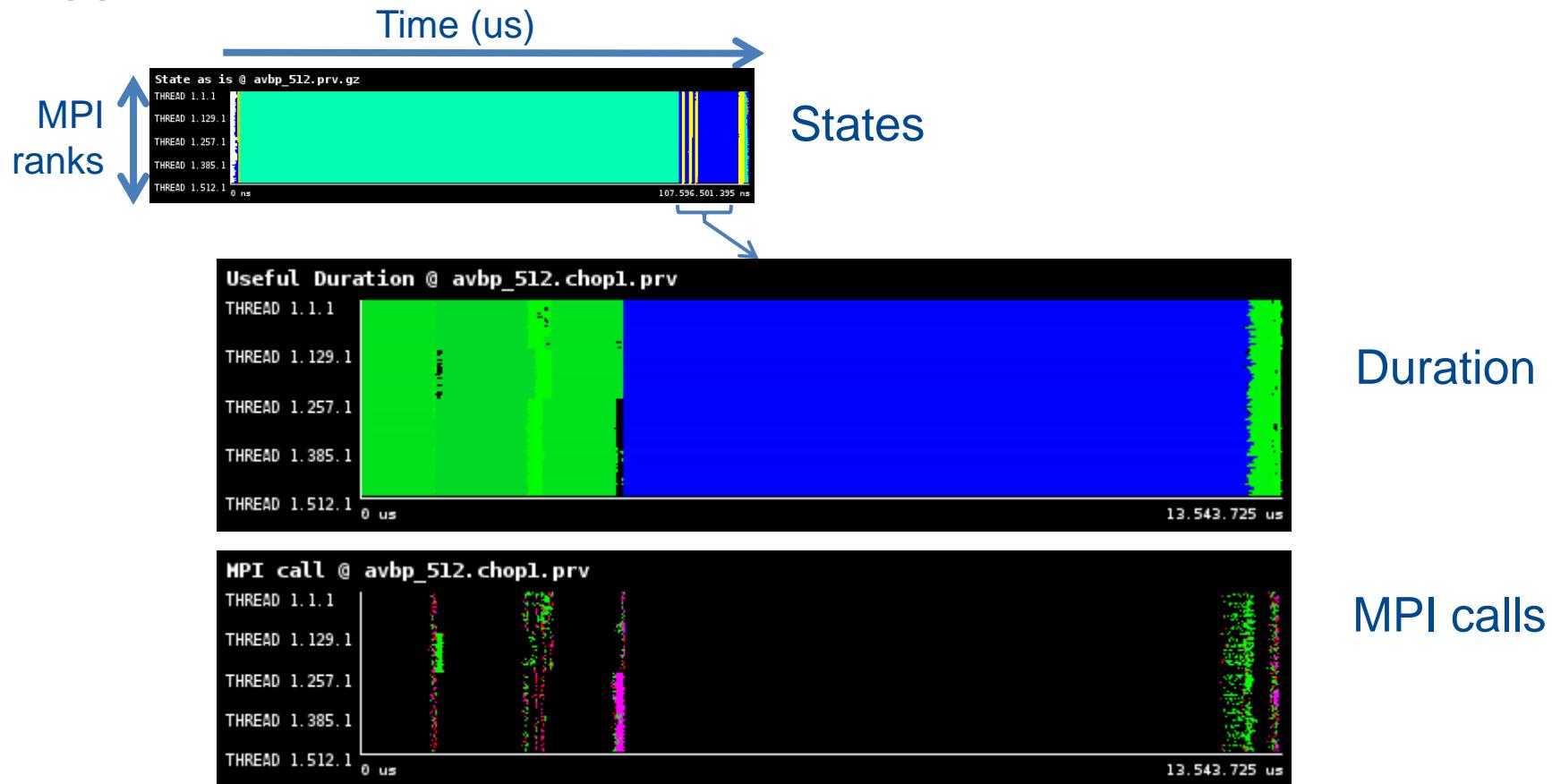
Examples from DEEP project

- AVBP – CFD code (strong scale)
- TURBORVB – Montecarlo simulation (weak scale)



AVBP structure

Application structure



AVBP scalability

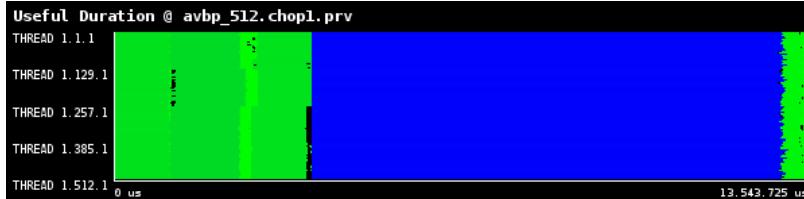
Input

- BG/Q runs 512, 1024, 2048, 4096
- Pure MPI executions, strong scale
- CFD simulations

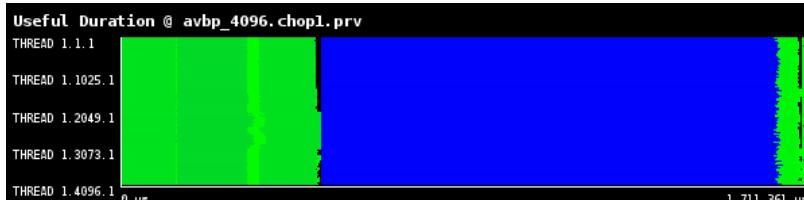
Analysis

- Load Balance: aprox. constant

512



4096

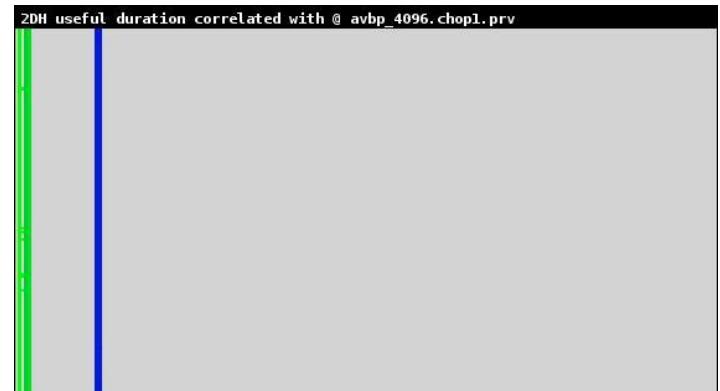


Computations histogram

512



4096

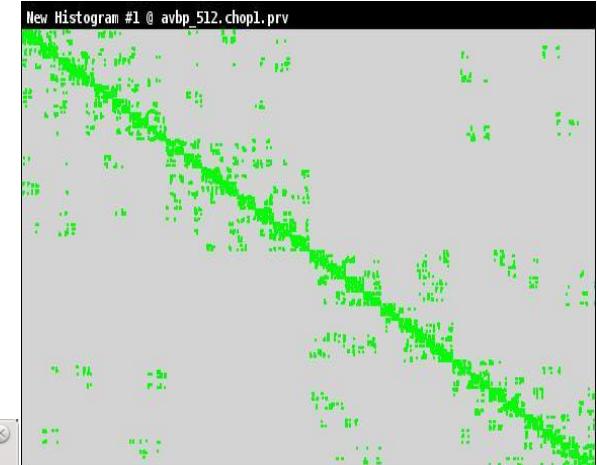


AVBP scalability

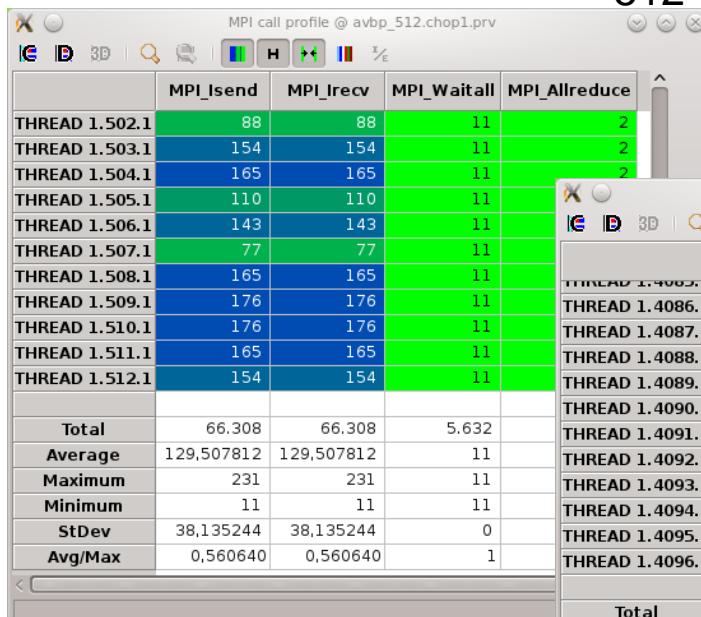
Analysis (cont.)

- Communications – similar pattern
↑ avg #point 2 point calls

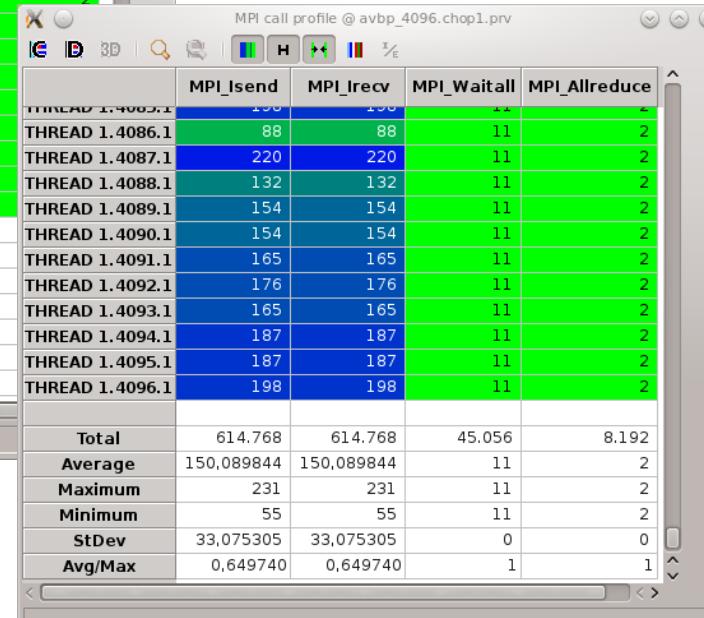
512



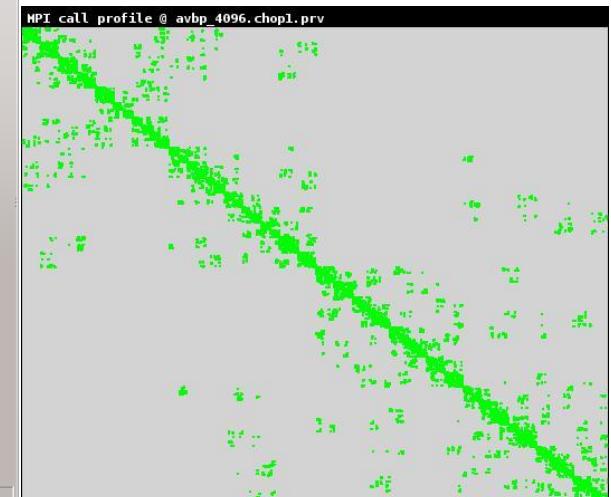
512



4096

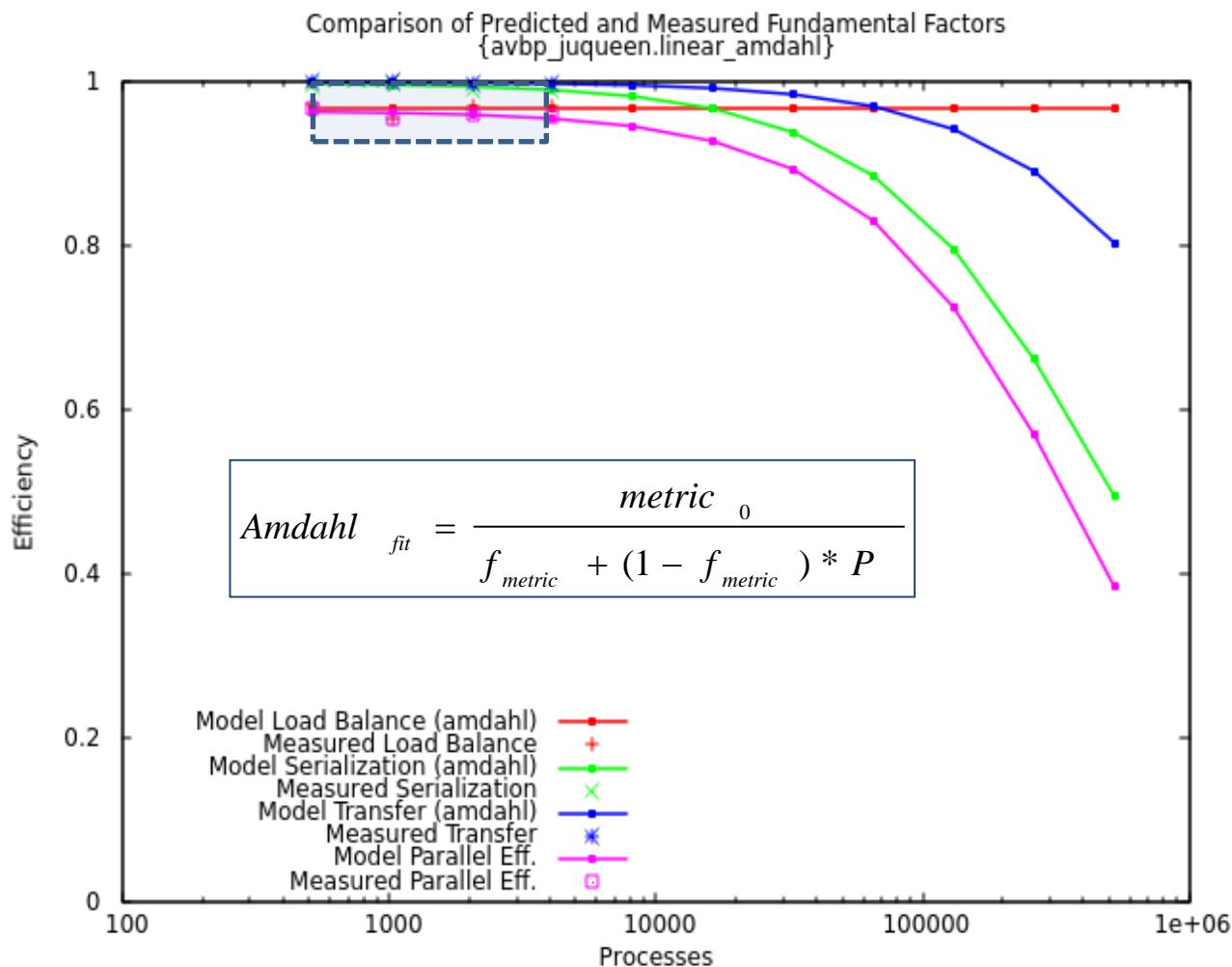


4096



AVBP extrapolation

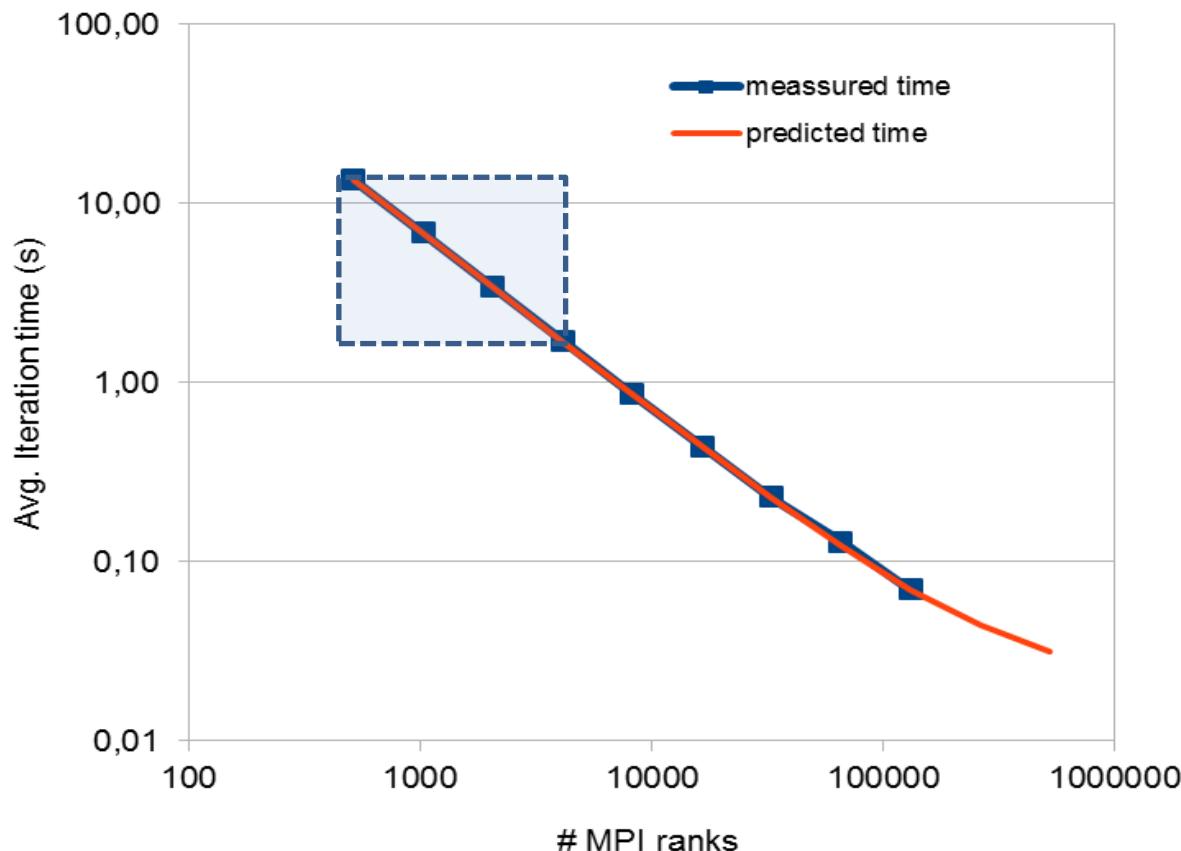
« Input: runs 512 – 4096



Strong scale →
small
computations
between MPI
calls become
more and more
important

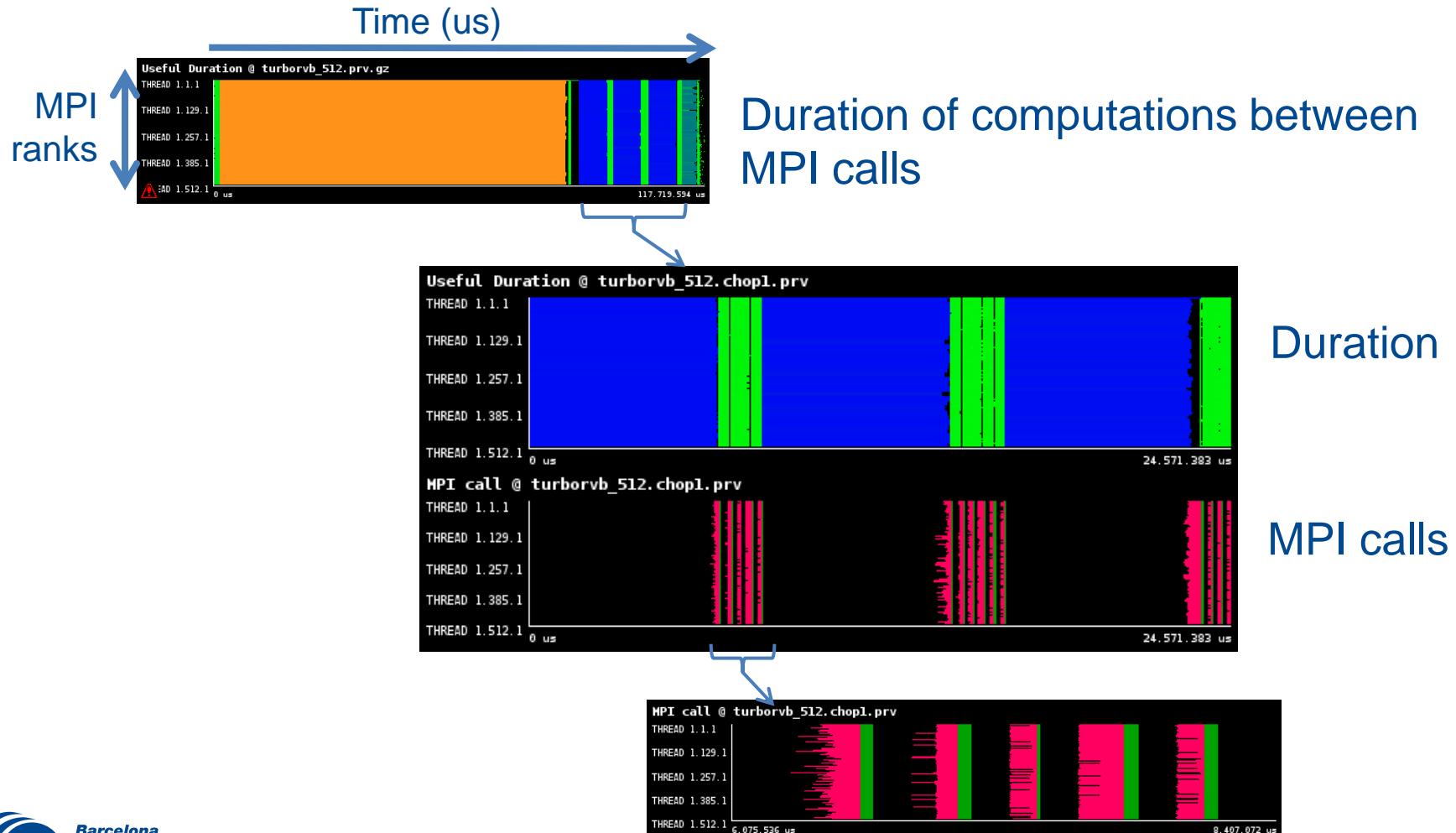
« Average iteration time (seconds)

- From Parallel efficiency (no instr., no IPC)



TURBORVB structure

Application structure



TURBORVB scalability

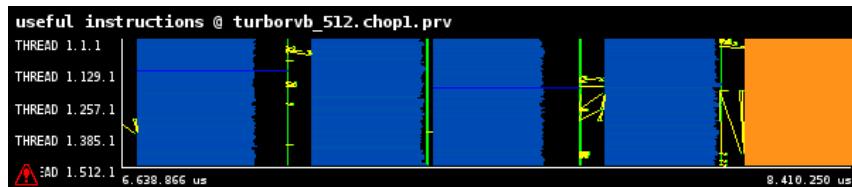
Input

- runs 512, 1024, 2048, 4096
- Pure MPI executions
- Montecarlo simulations

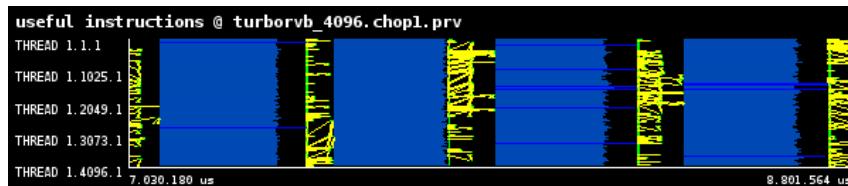
Modeling

- Load Balance: Small random unbalance

512



4096



Computations histogram

512



4096



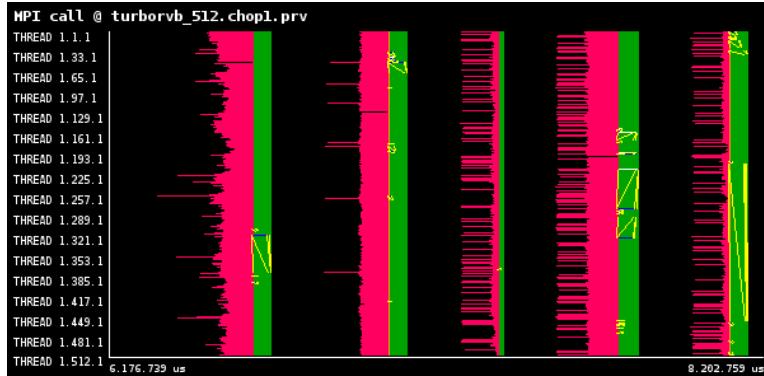
TURBORVB scalability

Modeling (cont.)

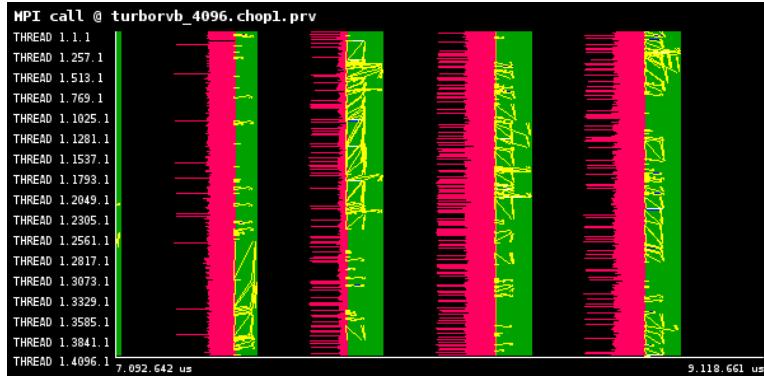
- Communications (Serialization, Transfer)
 - Collectives: Constant (size, number)
 - Point to point
 - Few processes do chains of 7 communications
 - Average calls per iteration per process constant (0.25)
 - Probability to generate contention on node network devices smoothly increases with the number of cores

$$Amdahl_{fit} = \frac{metric_0}{f_{metric} + (1 - f_{metric}) * P^{1/3}}$$

512

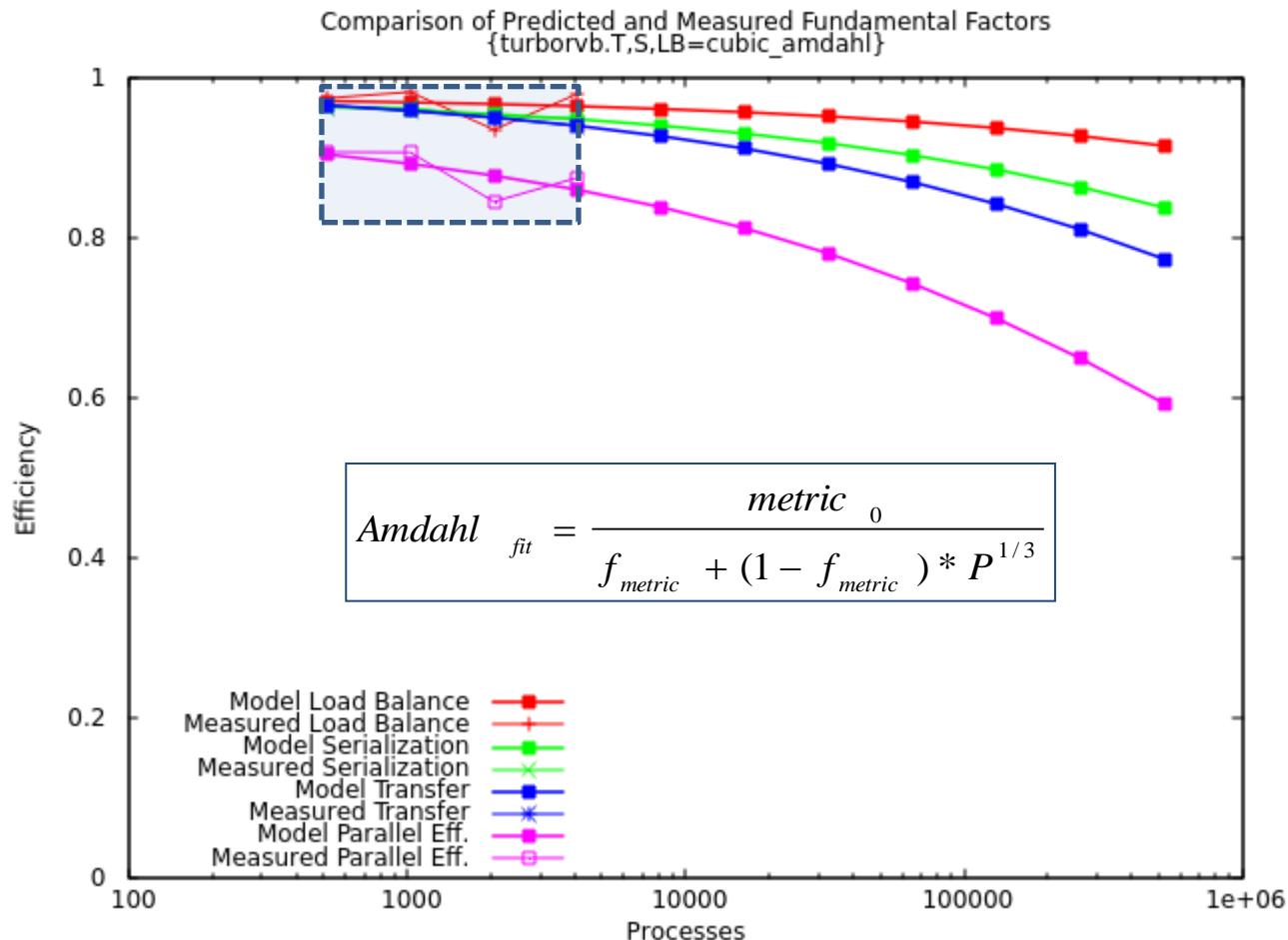


4096



TURBORVB extrapolation

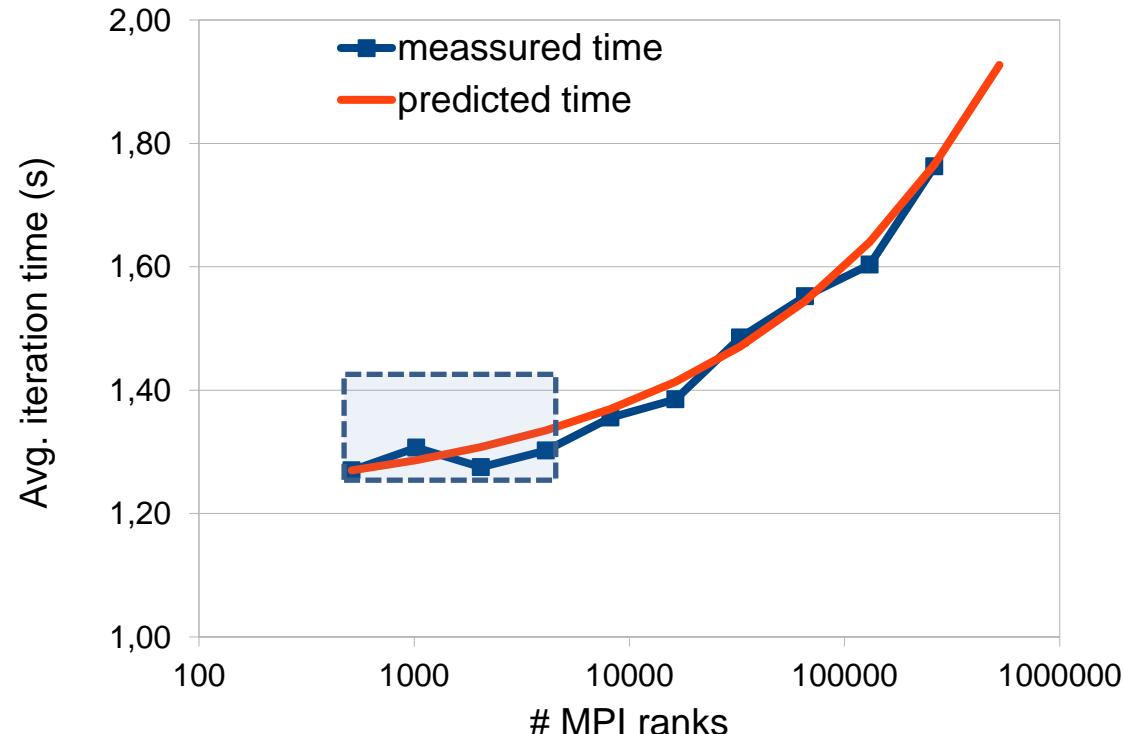
« Input: runs 512 – 4096



Network contention can be reduced balancing / limiting the random selection within a node

« Average iteration time (seconds)

- From Parallel efficiency (no instr., no IPC)



Conclusions

- « 4 measurements can predict performance at 100x
- « A lot of insight is obtained digging into traces for small core counts
- « Important side effect: identify relevance of different efficiency factors (where to improve)