

Hyperscaling of Plasma Turbulence Simulations in DEISA

H. Lederer, R. Hatzky, R. Tisma, A. Bottino*, F. Jenko*

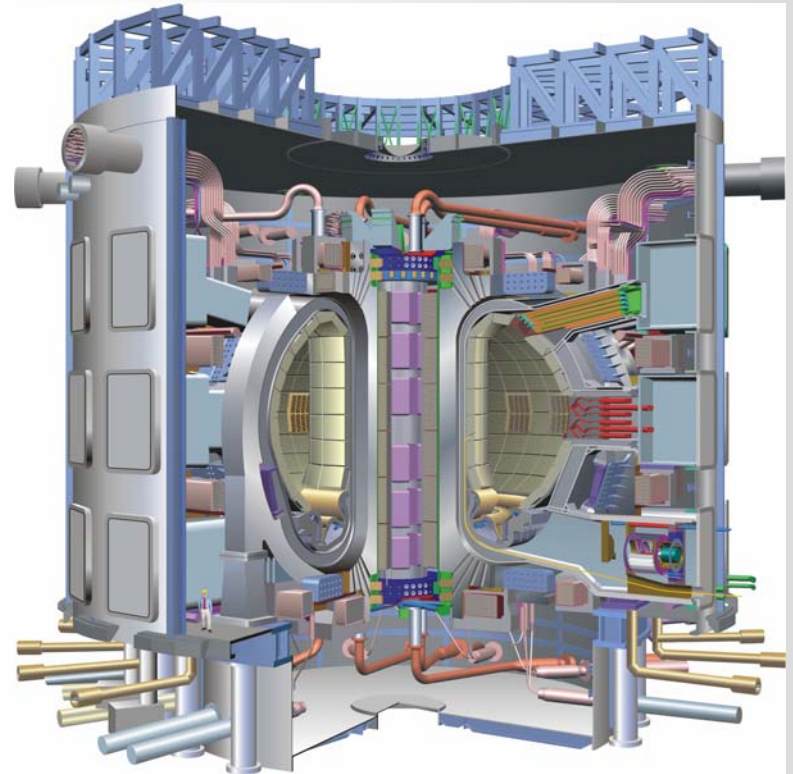
Garching Computing Center of the Max Planck Society

*Max Planck Institute for Plasma Physics

D-85748 Garching, Germany

ITER Experiment

www.iter.org



China



EU



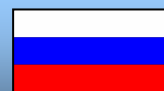
India



Japan



Korea



Russia



USA

incl. CH

CLADE 2007 Hermann Lederer
et al. (Lederer@rzg.mpg.de)

2



RZG IPP

Theory Support for ITER

- Large scale numerical simulations will be a necessity.
- Plasma turbulence simulations play a key role for the design, construction and optimization of the necessary fusion devices.
- The simulations will be compute and memory intensive
- Applications must be able to efficiently use tens of thousands of processors.
- Highly scalable applications are mandatory.

In Europe: Supercomputing support at continental level by DEISA (Distributed European Infrastructure for Supercomputing Applications)

DEISA

WWW.DEISA.ORG

11 sites in
7 countries

European Counterpart
of US TeraGrid

European Centre
for Medium-Range
Weather Forecasts

Forschungszentrum Jülich
in der Helmholtz-Gemeinschaft

CINECA
Consorzio Interuniversitario

epcc

lrz

iris

sara

RZG IPP

BSC
Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

CSC

Partners

Institut du Développement et des Ressources
en Informatique Scientifique, France
Forschungszentrum Jülich, Germany
Rechenzentrum Garching of the Max Planck
Society, Germany
Consorzio Interuniversitario, Italy
Edinburgh Parallel Computing Centre, UK
SARA Computing and Networking Services, The
Netherlands
Finnish Information Technology Center for
Sciences, Finland
European Centre for Medium-Range Weather
Forecasts, UK
High Performance Computing Center, Germany
Leibniz Computing Centre of the Bavarian
Academy of Sciences and Humanities, Germany
Barcelona Supercomputing Center (BSC), Spain

**Input from leading European
Plasma Physicists**

DEISA

**Application tuning,
hyperscaling,
porting to DEISA
architectures**

**DEISA
Joint Research
Activity in
Plasma Physics**

**DEISA Extreme
Computing Initiative**

DEISA & European Fusion Community

ITER support highly scalable applications required

DEISA: hyperscaling expertise
(unique service not offered by any
other European grid computing project)



**Address leading European
plasma turbulence codes**

Simulation Code ORB5

The ORB code family uses a particle-in-cell (PIC), time evolution approach, and takes advantage of all the recent techniques of noise reduction and control in PIC simulations.

ORB uses a statistical optimisation technique that increases the accuracy by orders of magnitude.

Initiated at CRPP (Centre de Recherches en Physique des Plasmas), Lausanne, ORB has been substantially upgraded at MPI for Plasma Physics (IPP), Garching.

Ongoing code development is made under a close collaborative effort between IPP and CRPP

Simulation Code GENE

GENE: so-called continuum (or Vlasov) code.

All differential operators in phase space are discretized via a combination of spectral and higher-order finite difference methods.

For maximum efficiency, GENE uses a coordinate system which is aligned to the equilibrium magnetic field and a reduced (flux-tube) simulation domain.

This reduces the computational effort by 2-3 orders of magnitude.

GENE can deal with arbitrary toroidal geometry (tokamaks or stellarators) and retains full ion/electron dynamics as well as magnetic field fluctuations.

At present, GENE is the only plasma turbulence code in Europe with such capabilities.

ORB5 Code: Single processor optimization

Two bottlenecks identified and improved.

- **implementation of the FFT**
- **cache sort of the Monte Carlo particles**

FFT: module written with different interfaces to specialized FFT libraries.
(original code included FFT source code with poor performance)

Interfaces to FFTs from IBM ESSL, Intel MKL, FFTW

-> no restrictions to vector lengths of powers of two

-> more flexibility to choose the grid resolution of the electrostatic potential

Cache sort:

- sorting the Monte Carlo particles relative to their position in the grid cells of the electrostatic potential, results in high cache reuse of the electrostatic field sampling
- overhead caused by introduction of the sort routine was minimized
- option to enlarge a work array for the sorting process (can speed up the isort routine by a factor of three)

ORB5 Code: Improving scalability

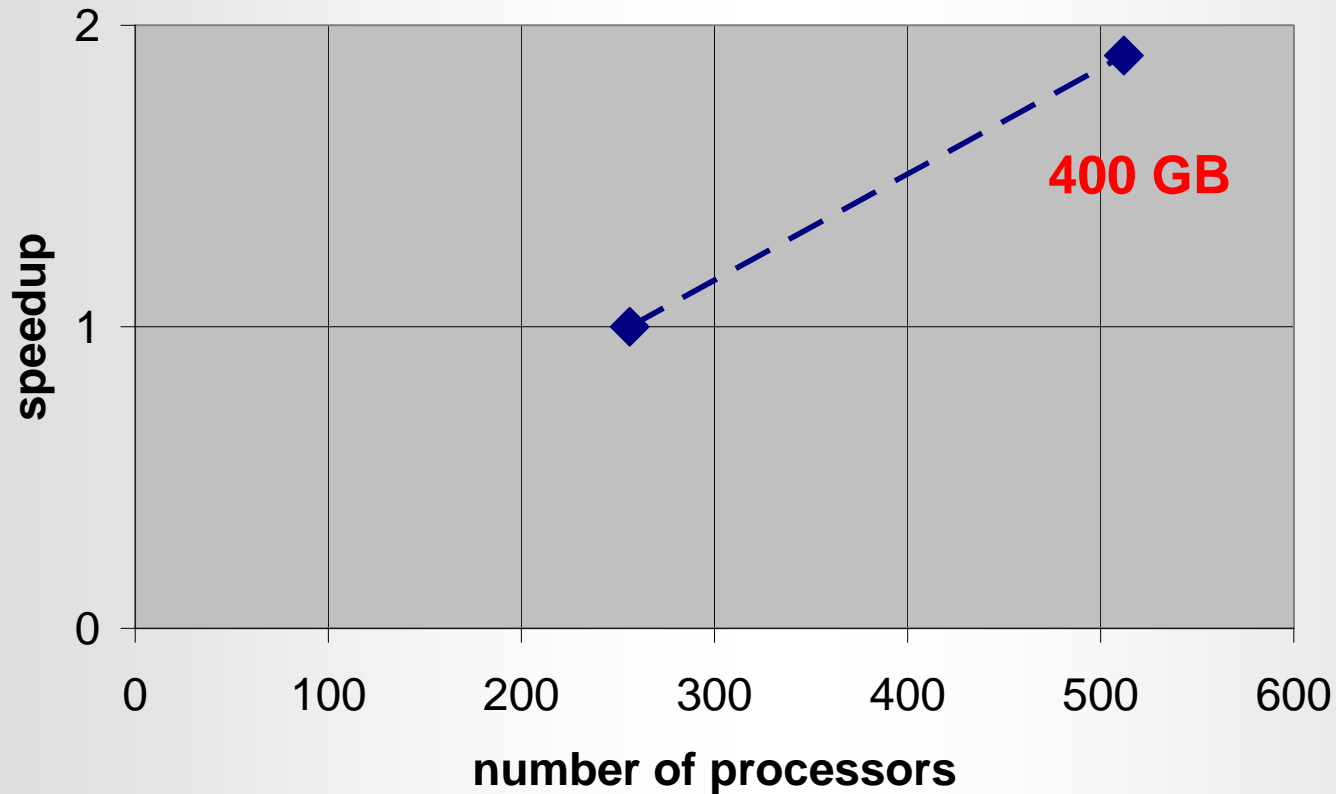
Implementation of the domain cloning concept:

- optimizes scaling
- decouple the selectable grid resolution from the number of processors used for the simulation.

Domain cloning builds on two decomposition techniques:

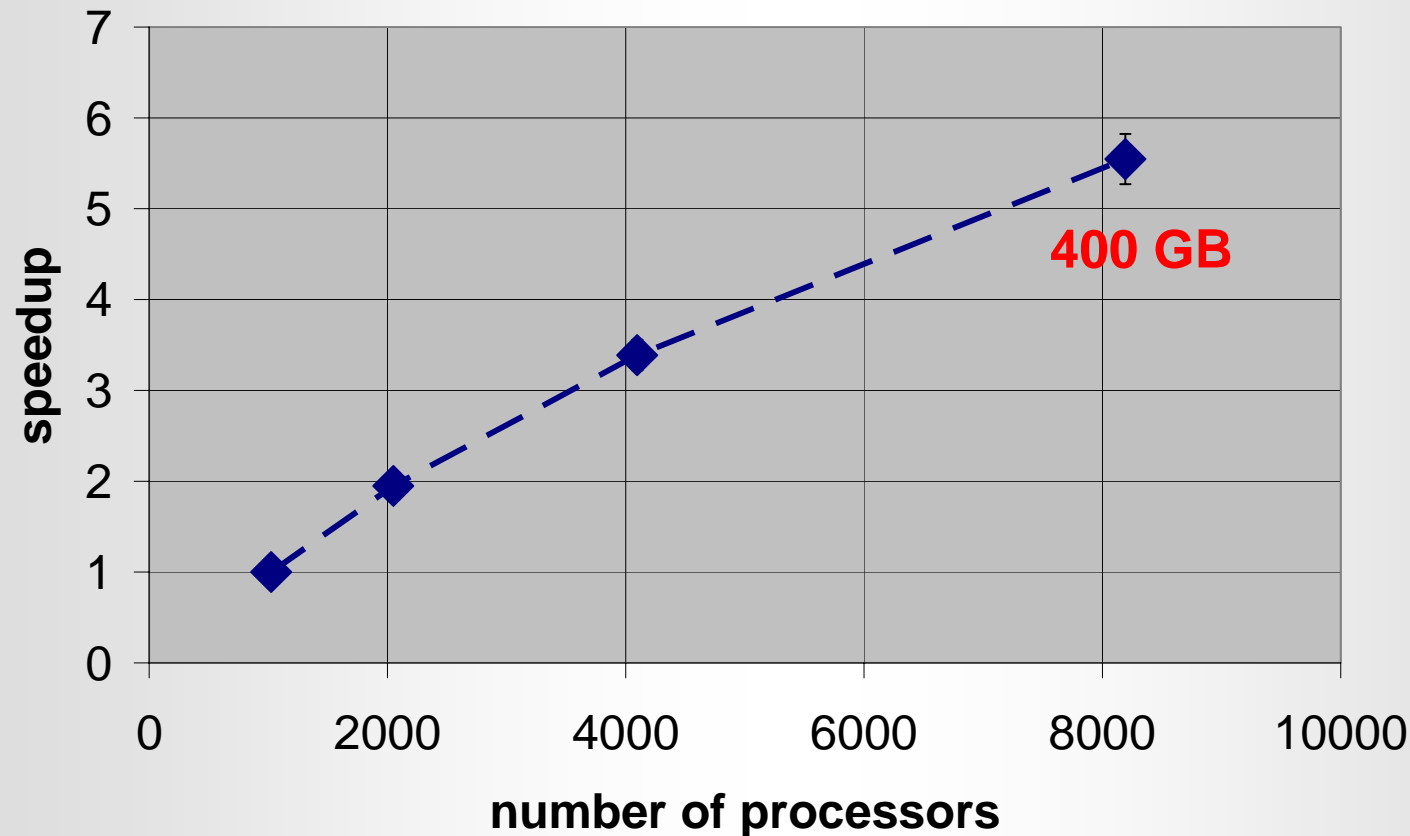
1. For the domain decomposition, different portions of the physical domain and of the corresponding electrostatic field grid are assigned to different processors together with the Monte Carlo particles that reside on them. As particles move from one region to another, they are communicated to the processor which is associated with the new region.
2. Particle decomposition: the whole spatial grid is assigned to every processor, but each processor takes care of only a subset of the particle population. Partial contributions to the ion density, which is required to update the electrostatic field, are communicated among processors and summed via global sum operations.

ORB5 Code: Test of scalability on IBM Power4



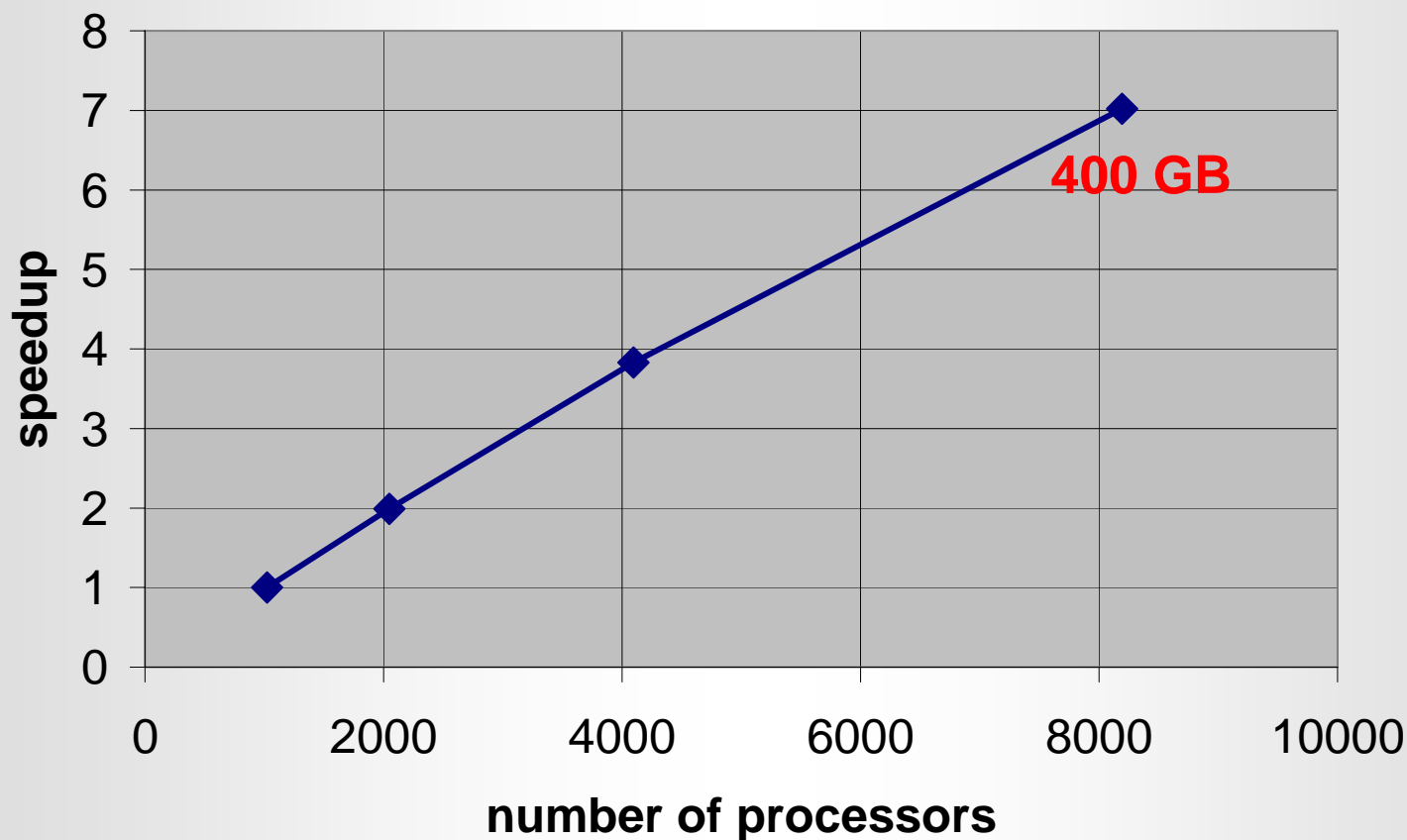
Strong scaling of ORB5 for an ITG simulation (~400 GB)
512 proc. result normalized on the 256 proc. result.
(Measurements on IBM [p690@1.3](#) GHz + HPS at RZG)

ORB5 Code: Scalability test on Cray XT3



Strong scaling of ORB5 for an **ITG** simulation (0.5 billion particles)
(results normalized on the result for 1024 processor-cores)
(Measurements on Cray XT3 at ORNL; courtesy of Cray Inc.)

ORB5 Code: Scalability test on BlueGene/L

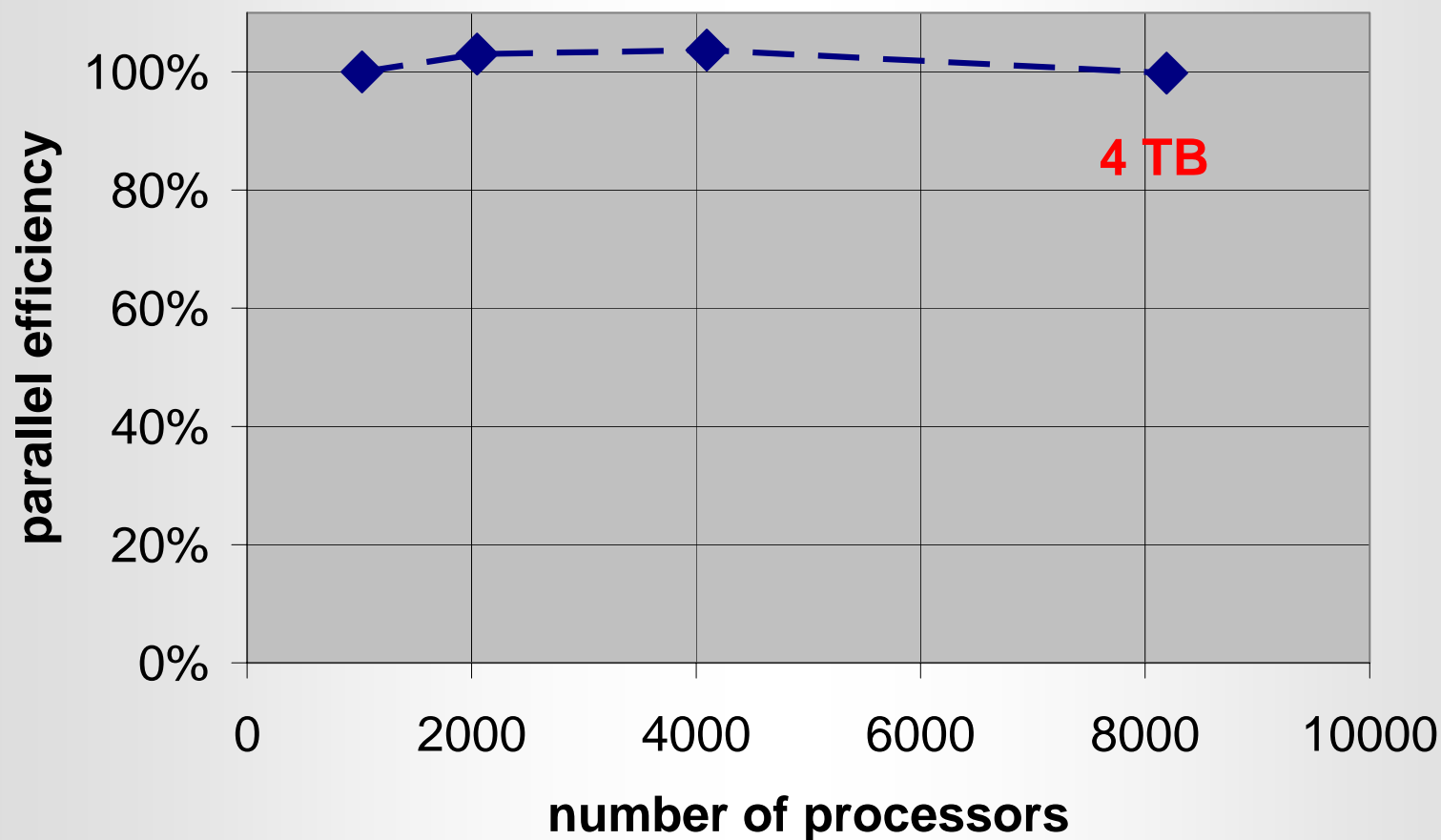


Strong scaling of ORB5 for an **ETG** simulation (0.8 billion particles)

Results normalized on the 1024 processor result

(Measurements on BlueGene/L at IBM Watson Research Center in co-processor mode)

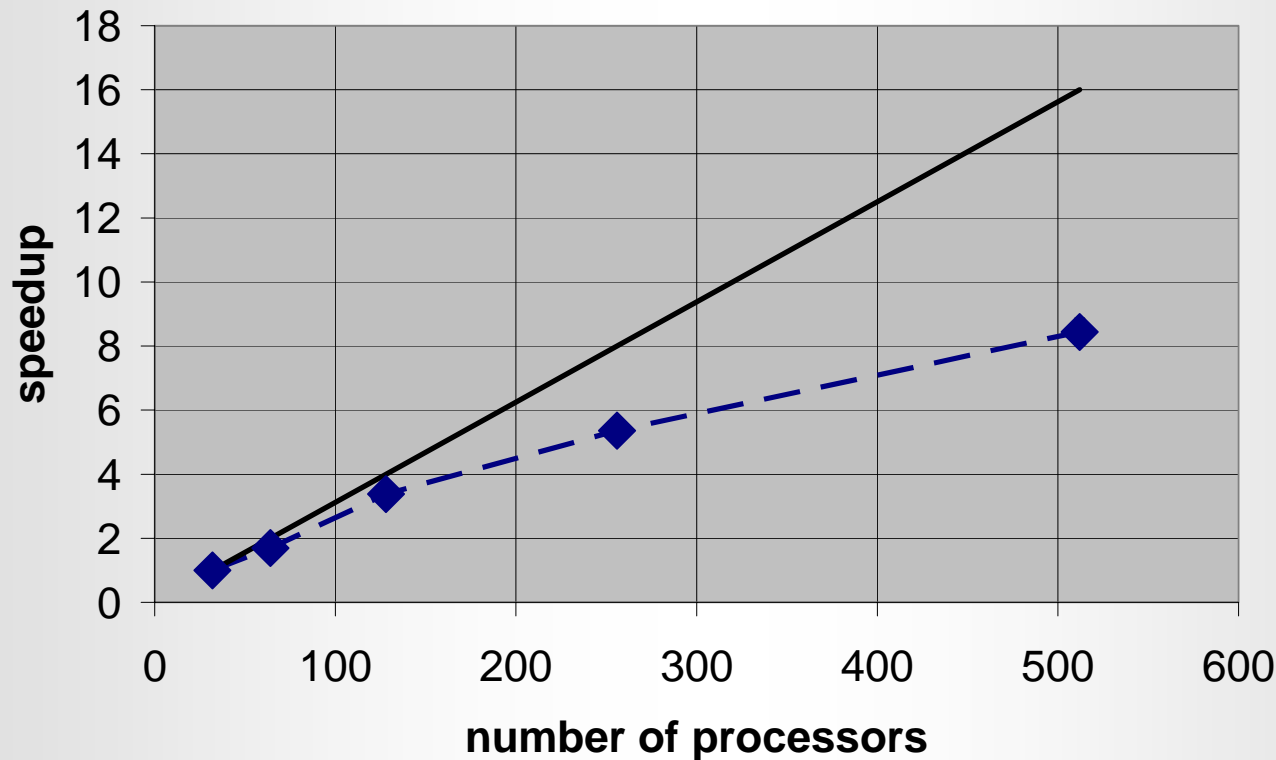
ORB5 Code: Scalability test on BlueGene/L



Weak scaling of ORB5 for **ETG** simulation (~ 0.8 M particles per processor)
Results normalized on the 1024 processor result
(Measurements on BlueGene/L at IBM Watson Research Center in co-processor mode)

Analysis of GENE v 9

- Mixed parallelization model MPI+OpenMP
- Upper limit of 64 MPI tasks



-> Limited scalability

Analysis of GENE v 9

6 dimensions available for parallelisation: species + 3 space coordinates. + 2 velocity coordinates

The spatial coordinates x and y , in GENE v 9 only treated serially, contain significant potential for domain decomposition.

A large number of 2-dimensional FFTs done on the xy planes. If the xy plane is distributed, it must be transposed in order to perform the FFT in the x and y directions.

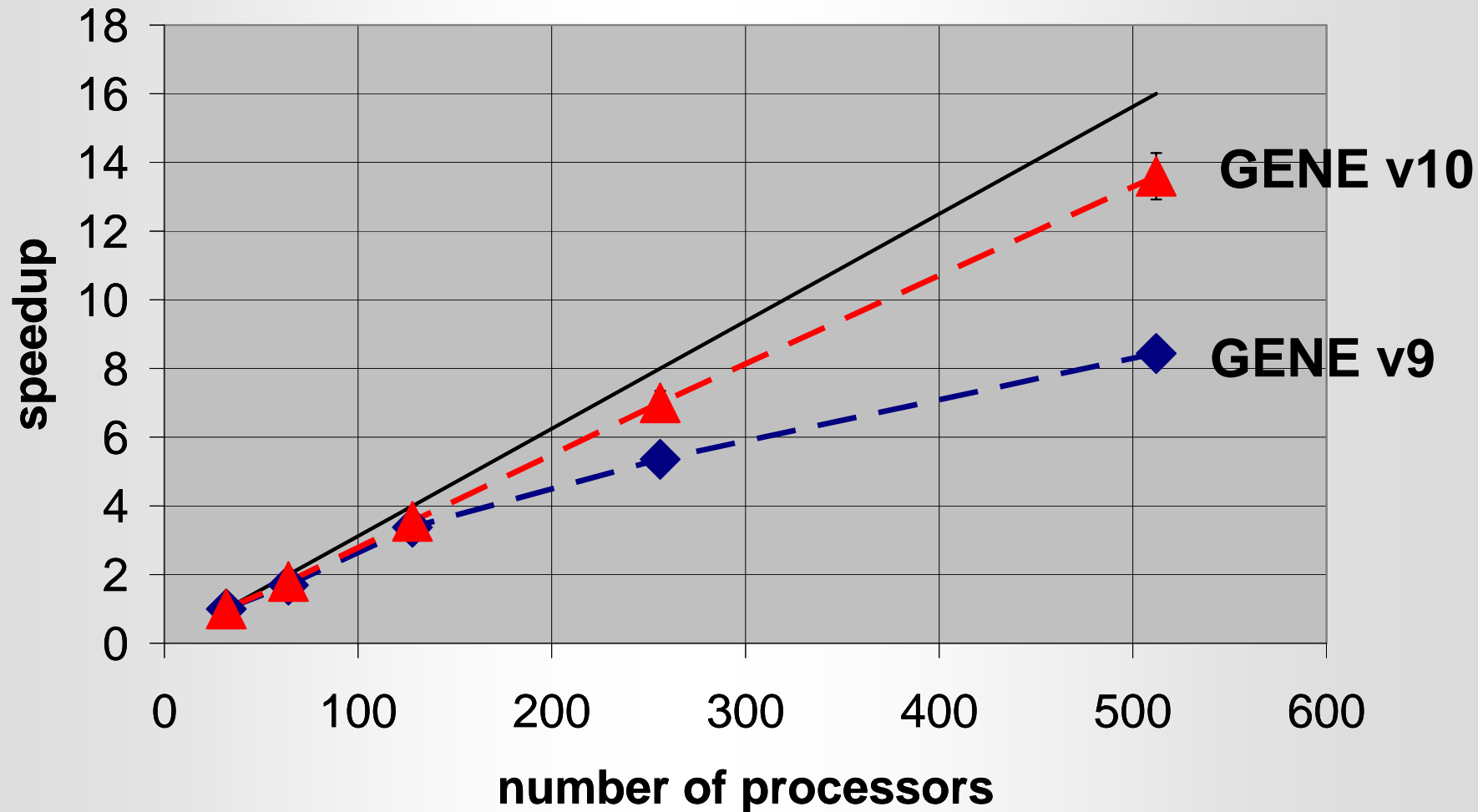
However, the transposition requires an all-to-all communication
⇒ high communication overhead.

If the y -coordinate could be left in k -space, the FFTs would have to be performed on the x -coordinate only, which is contiguous in memory. A transpose of the xy plane would **not** be necessary.

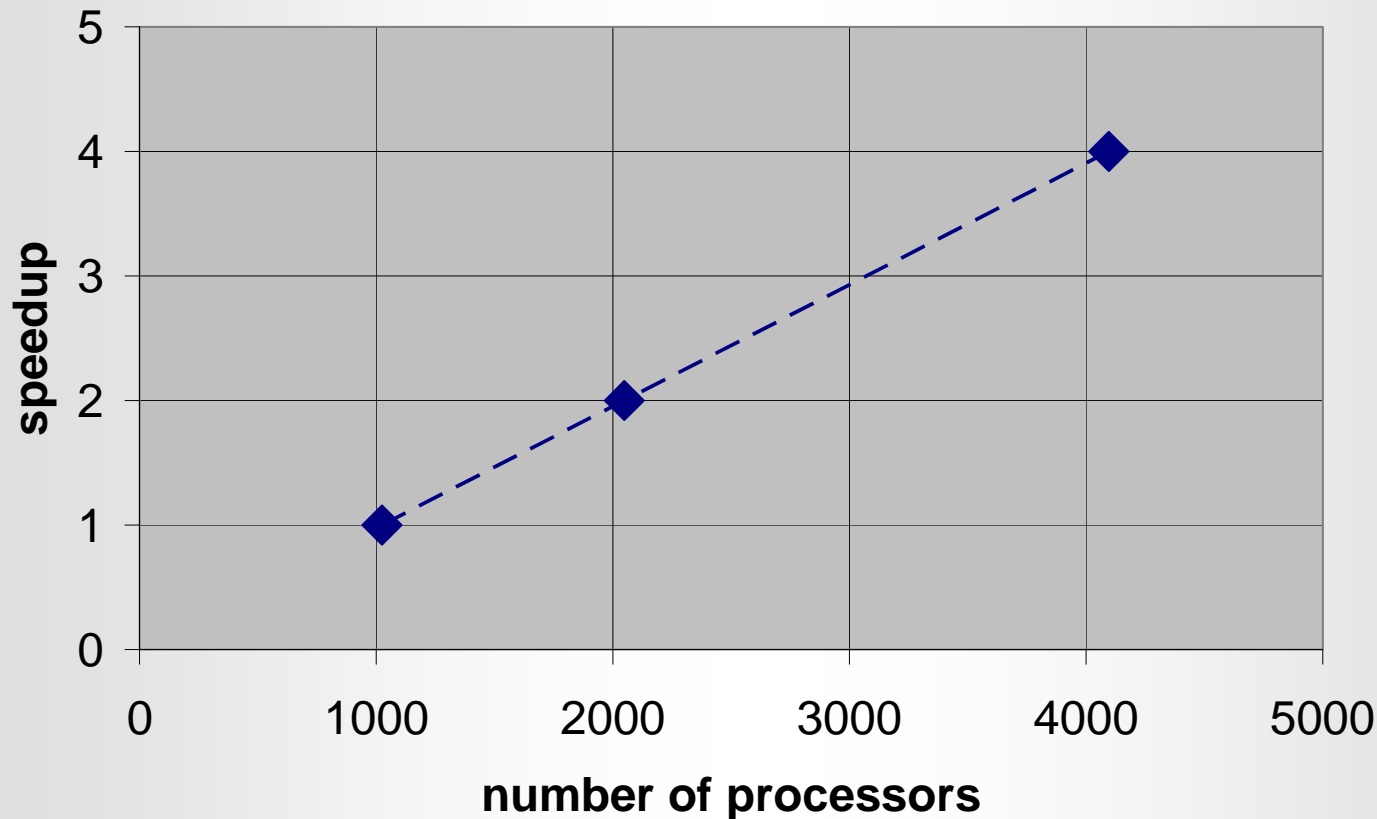
Improving the scalability of GENE

- Sections in the code where transformations to the configuration space are necessary.
- For these transformations, transpositions of the xy plane have to be performed. However, the number of these transformations is much smaller than the number of transformations to k -space.
- Implication of a major change to the overall data structure
-> consequences for many parts of the code.
- Prerequisite: code adaptation by the authors according to the new role of the y -coordinate.
- Main task: design and realization of the domain decomposition of the y -coordinate. Dynamical mapping of the number of points available in the y direction on the number of processors selected for treating the y direction (consequently applied to all loops in the y direction (~ 20 to 40)).

GENE scalability on IBM p690 at RZG



Verification of GENE v10 scalability on a larger number of processors on Cray XT3



Strong scaling of GENEv10 (problem size of ~300-500 GB)
normalized to 1024 processors (Measurements courtesy of Cray Inc.)

Enhancing portability of GENE

DEISA architectures: IBM Power4/5, PowerPC
SGI Altix (Intel Montecito), Cray XT4, NEC SX-8

FFT-routines: Interfaces to

- IBM proprietary ESSL-library (on IBM systems)
- Math Kernel Library (MKL) from INTEL (on SGI Altix)
- FFTW package

Bessel functions:

Routines originally used from the NAG library were replaced:
three new subroutines were written implementing these
Bessel functions

GENE v11

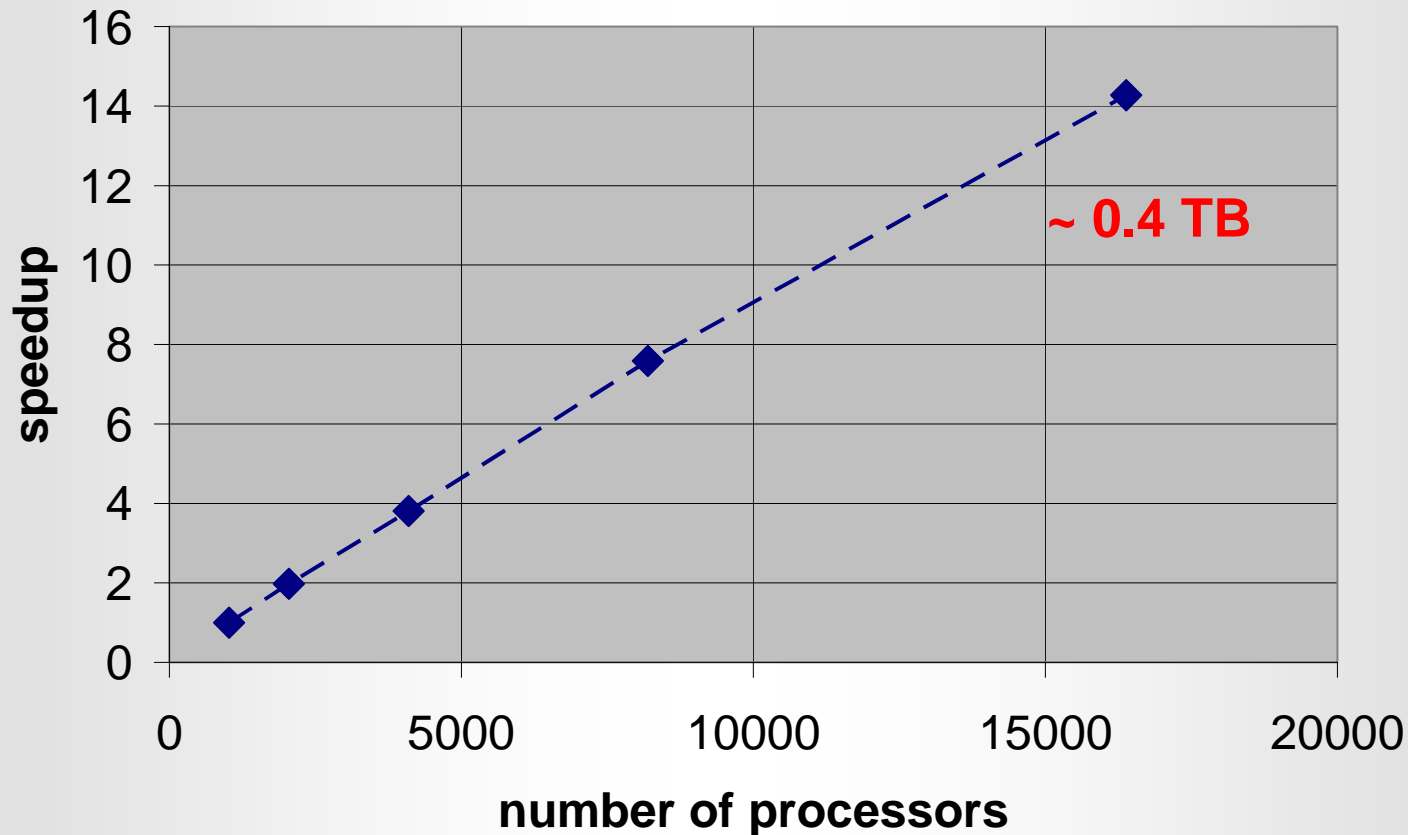
and tests on higher processor numbers

- GENE v11: Enhancements of algorithms and of functionality by authors
- Parallelization scheme of GENEv11 same as that of GENEv10
- Tests of GENE v11 scalability on higher processor numbers
- Porting of GENE v11 to IBM BlueGene/L
- Test of GENE v11 up to 8k processors on BG/L
- Results on BG/L for strong scaling measurements
 - on 4k processors: quasi-linear speedup
 - on 8k processors: degradation (parallel efficiency: 73%)

Further improvements of scalability?

- Still potential for parallelization: z velocity dimension
- Parallelization of the z dimension by domain decomposition
- New GENEv11+
- Test of GENEv11+ on BG/L
- Results on BG/L for strong scaling measurements:
 - on 4k processors: quasi-linear speedup
 - on 8k processors: quasi-linear speedup (parallel efficiency: 95%)

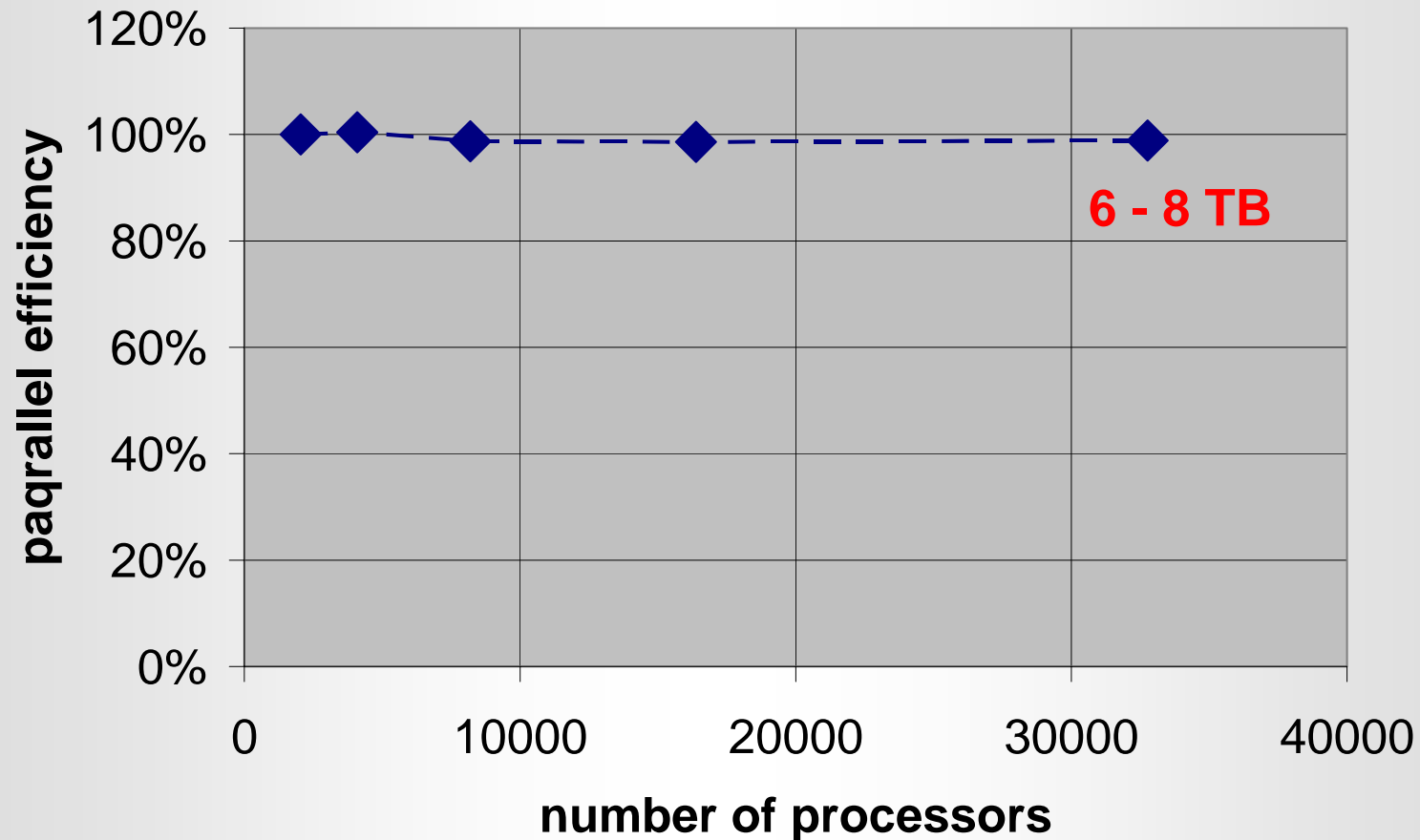
GENE v11+ on IBM BlueGene/L



Strong scaling of GENE v11+ normalized to 1k processors

(problem ~300-500 GB; measurements in co-processor mode at IBM Watson Research Center)

GENE v11+ on IBM BlueGene/L



Weak scaling of GENEv11+ normalized to 2k processors

(problem ~200 MB/proc; measurements in virtual node mode at IBM Watson Research Center)

Conclusions

Two important European plasma turbulence simulation codes, GENE and ORB5,

could be optimized and adapted to very high scalability requirements

as a major step towards efficient usage

of forth-coming new generations of petaflops supercomputers

required for realistic simulations of ITER.

Acknowledgements

We thank the European Commission for support through contracts FP6-508830 and FP6-031513.

We thank IBM for access to the BlueGene/L system at Watson Research Center, and J. Pichlmeier for support on using the system.

We thank Cray Inc. for scalability runs on the Cray XT3 system at ORNL.