

High-Performance Computing in Theoretical Chemistry: Today and Tomorrow

Axel Kohlmeyer

Lehrstuhl für Theoretische Chemie

Ruhr-Universität Bochum

`<axel.kohlmeyer@theochem.rub.de>`

Blue Gene/L Workshop, Jülich, Sep 2004



Outline

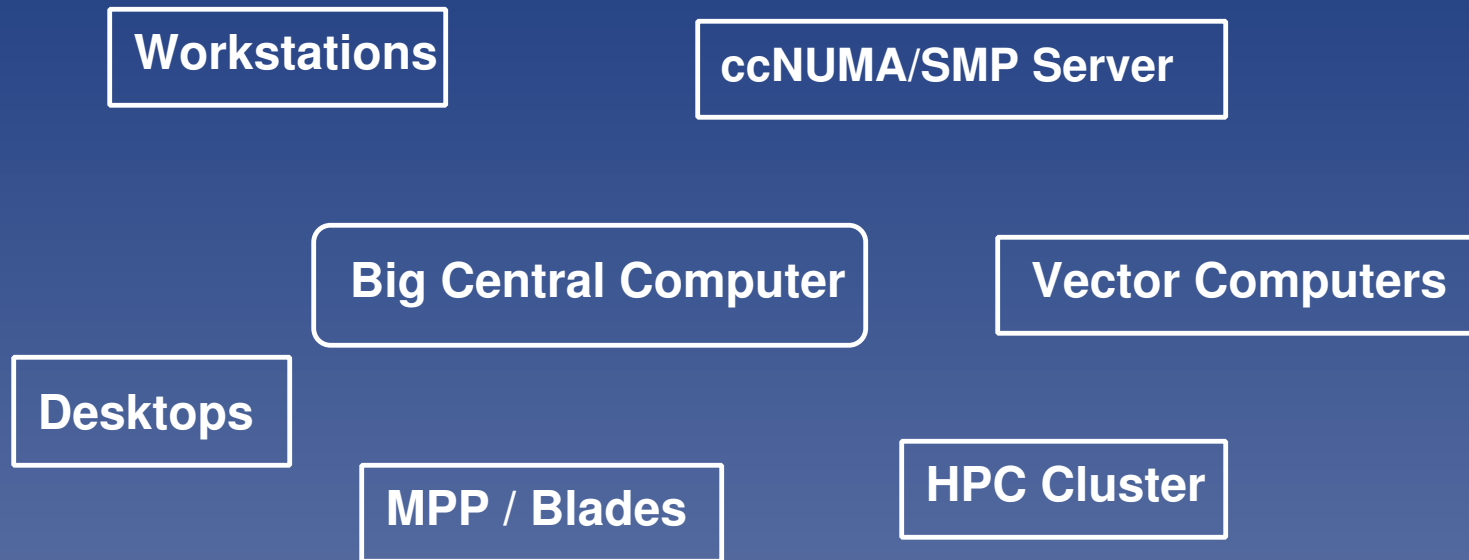
- Motivation
- Hardware
- Benchmarks
- Software
- Trends
- Examples
 - Solvated DNA Fiber
 - Surface Hopping
 - Catalytic Triade
- Summary



Motivation

- “Sysadmin-Perspective”:
 - ★ Planning and managing hardware acquisition
 - ★ Compiling, porting, and optimization on several platforms
 - ★ Different packages from all areas of theoretical chemistry
- Inverse approach to select topics and tools:
 - ★ Interest to use hardware efficiently
 - ★ Show how to use tools well with available hardware
 - ★ Select problems according to tools

Hardware Overview



Diversification trend:

from 'one size fits all' machine to diverse architectures

Unification trend:

Unix-like operating systems, MPI/OpenMP programming models



Hardware Differences

- not much difference in pure CPU speed
- main difference in memory and I/O bandwidth
- differences in reliability under heavy duty use
- differences in connectivity

for imbalanced machines (e.g. Linux PC): careful deployment planning and architecture specific optimizations most effective

CPMD Serial Runs (10 Ry)

Machine	Wall Time / s
AMD Athlon XP1600+, 1.4GHz, PC133	545
AMD Athlon MP1600+, 1.4GHz, PC266-ECC	443
Compaq Alpha EV6, 600MHz, XP1000	435
HP SuperDome 32000, HPPA 8700,750MHz	388
AMD Athlon XP2500+, 1.83GHz, PC333	361
Compaq Alpha EV67, 677MHz, ES40	284
AMD Opteron, 1.6GHz, 32-bit	287
AMD Opteron, 1.6GHz, 64-bit	254
Intel P4 Xeon, 2.4GHz, PC266	236
Compaq Alpha EV68AL, 833MHz, DS20	234
Intel Itanium2, 900MHz, HP zx6000	206
AMD Athlon64 3200+, 2.0GHz, PC333, 64-bit	173
IBM Power4+ 1.7 GHz, Regatta H+	171

CPMD Serial Runs (30/50 Ry)

Machine	Wall Time / s
AMD Athlon XP1600+, 1.4GHz, PC133	2878
HP SuperDome 32000, HPPA 8700, 750MHz	2672
Compaq Alpha EV6, 600MHz, XP1000	2624
AMD Opteron, 1.6GHz, PC266 memory, 64-bit	1292
Intel Pentium 4 Xeon, 2.4GHz,	1275
AMD Opteron, 1.6GHz, PC266 memory, 32-bit	1157
IBM Power4+ 1.7 GHz, Regatta H+	997
AMD Athlon XP1800+, 1.53GHz, PC266	5878
AMD Athlon XP2500+, 1.83GHz, PC333	5196
AMD Athlon XP2500+, 1.83GHz, dual-channel PC333	3848
Intel Itanium2, 900MHz, HP zx6000	3145
AMD Opteron, 1.6GHz, PC266 memory, 32-bit	3143
AMD Athlon64 3200+, 2.0GHz, PC333, 64-bit	3134
IBM Power4+ 1.7 GHz, Regatta H+,	2259

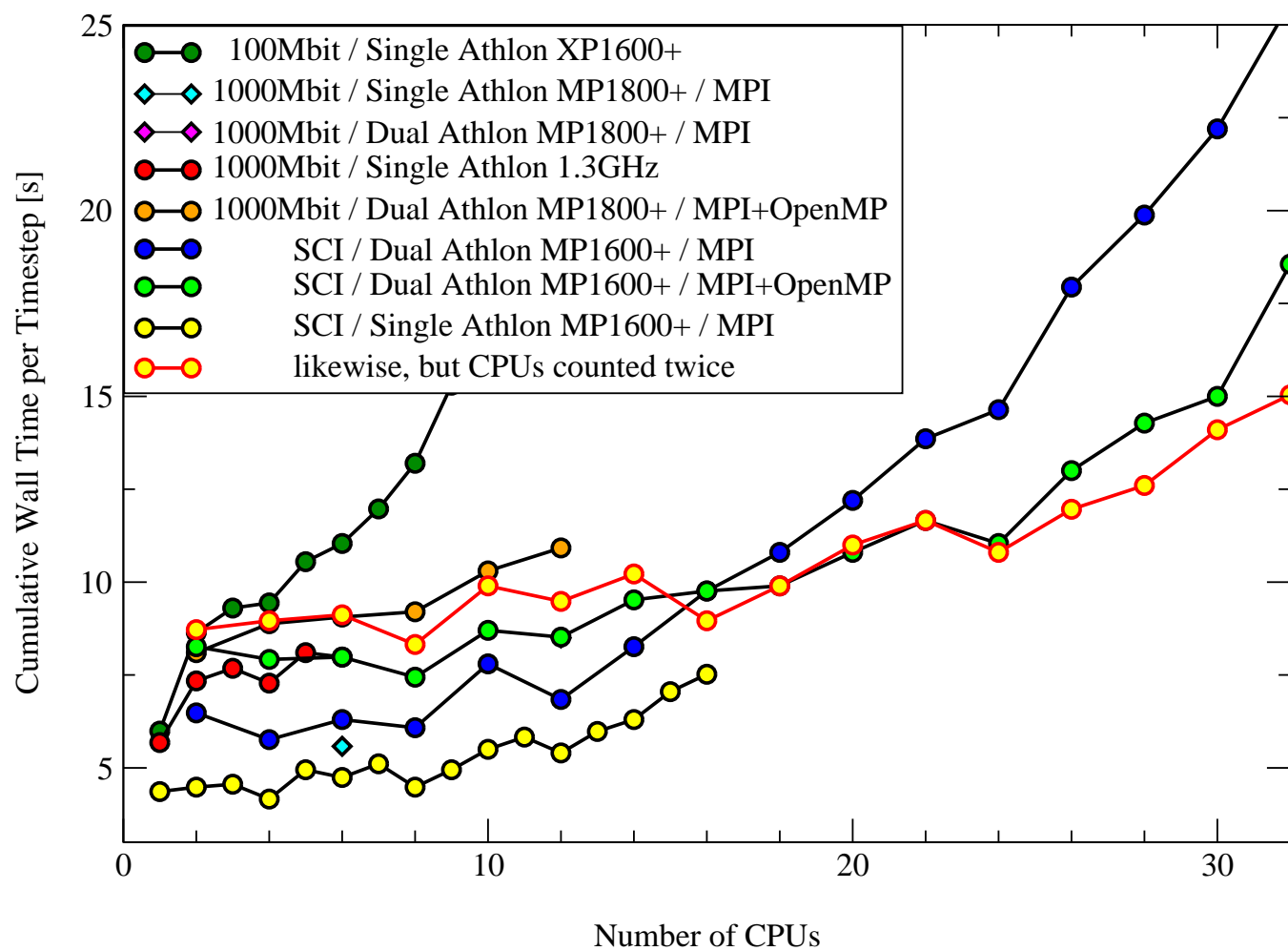


Library Optimizations

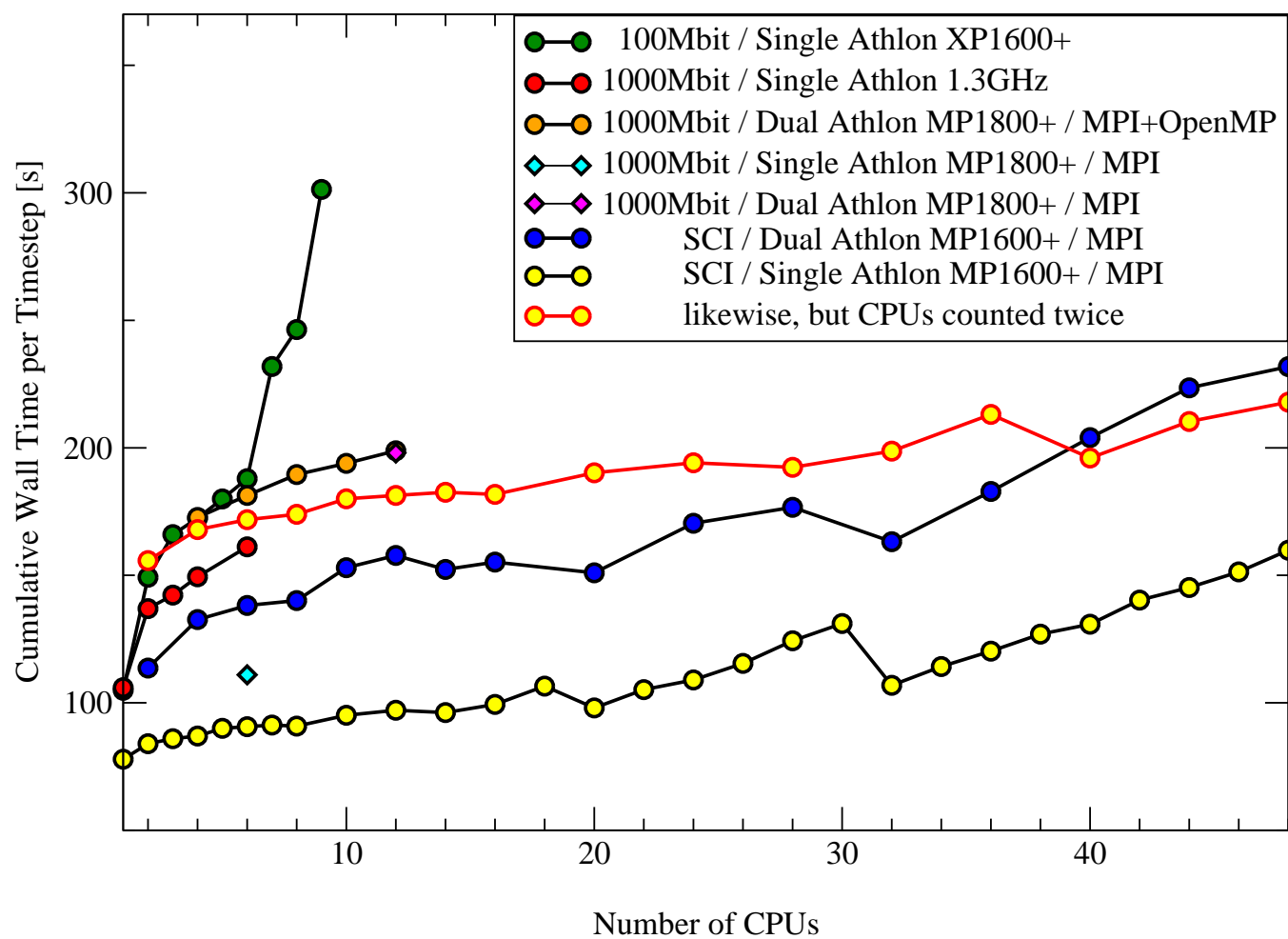
CPMD, 100 steps CP-MD: 63 Si-Atoms, 10Ry

Machine	BLAS	generic ATLAS	specific ATLAS
Athlon XP1800+ (1.53 GHz)	950 s 251%	428 s 113%	378 s 100%
Pentium IV 2GHz	765 s 173%	493 s 112%	441 s 100%
P4 Xeon 2.4GHz	471 s 171%	316 s 118%	276 s 100%
Pentium M 900MHz	716 s	430 s	-

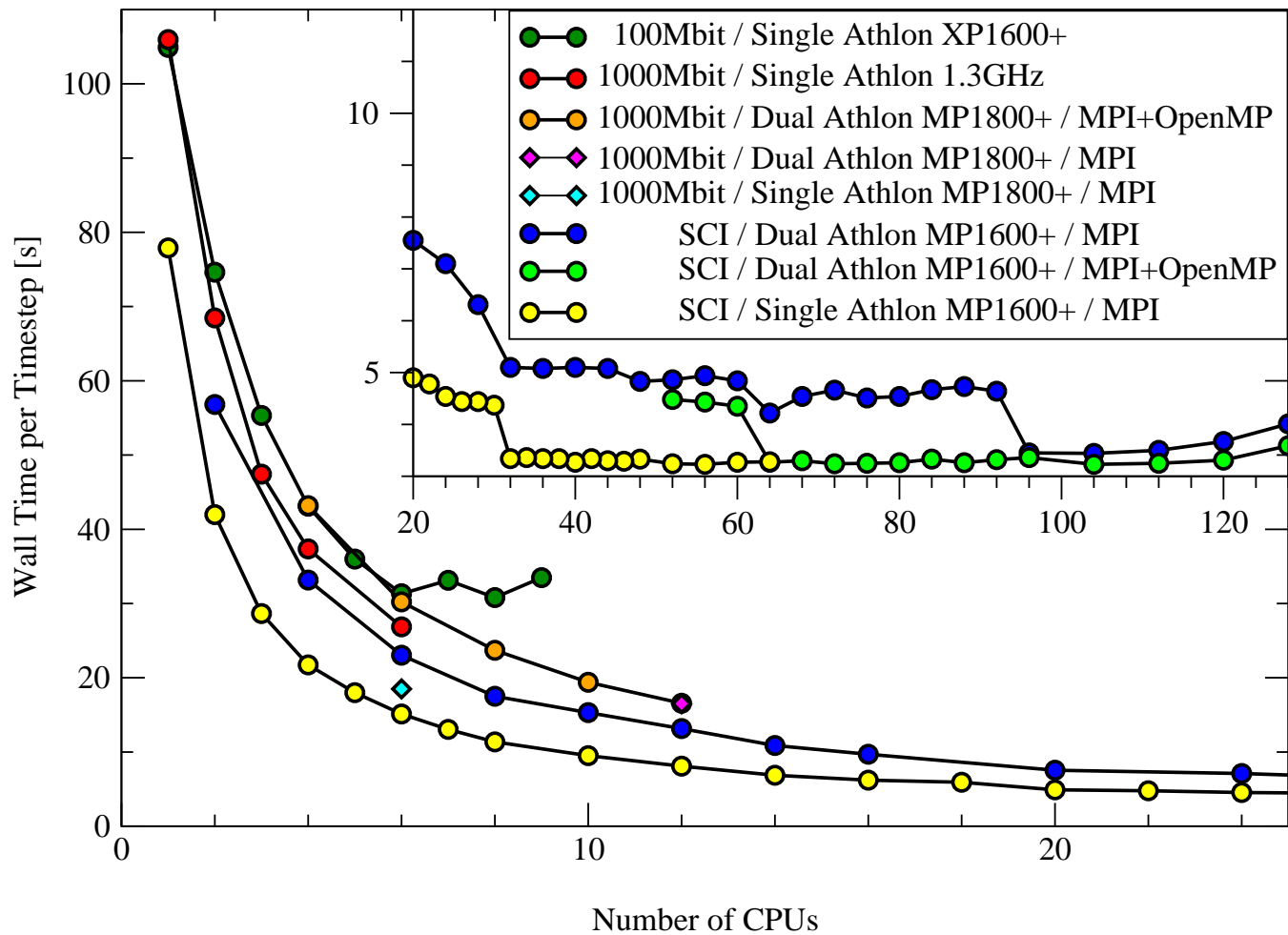
Si_63 bulk / PBC / 10 Ryd



Si_63 bulk / PBC / 70 Ryd



Si_63 bulk / PBC / 70 Ryd





Hardware Summary

- different types of hardware for different problems
 - better parallel scaling with larger problems
 - bandwidth matters (memory, I/O, network)
 - optimal throughput and optimal performance (“capacity” and “capability”) for different hardware
- applications in theoretical chemistry need them all

Application Overview

Computations in Theoretical Chemistry cover wide spectrum:

Method	Application	Scaling
Configuration Interaction	Energetics, Spectra	N^6
Hartree-Fock	Molecular Structure	N^4
DFT FP Molecular Dynamics	Structure Dynamics	N^3
classical Molecular Dynamics	Biomolecules, Liquids	N^2

High demand for more CPU power to:

- treat problems more accurately
- treat larger problems
- get better statistics

⇒ larger/faster computers *and* better algorithms



How To Improve The Software

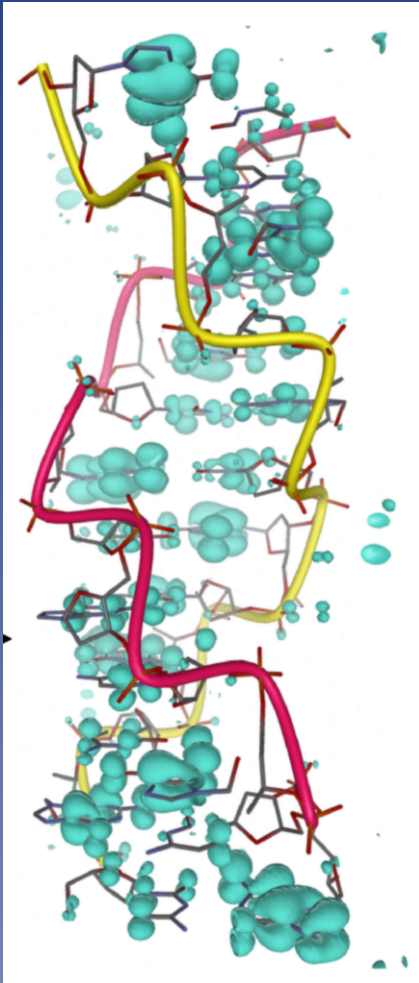
- new algorithms needed to approach 'linear scaling':
 - ignoring or approximating faraway contributions
 - only works for "large enough" problems
 - parallelization of existing software:
 - good parallel scaling difficult if not trivially parallelizable or written to be parallelizable
 - code development takes long time:
 - extensive testing needed to ensure correct results,
 - few good scientists are also good programmers
- ⇒ "new" software projects examples: NWChem, CP2k, NAMD



Where Use A Massively Parallel Machine?

- parallelizable² applications:
 - first principles molecular dynamics with path-integrals
 - replica-exchange (classical) molecular dynamics
- massive job farming through an external 'driver':
 - scanning of potential energy surfaces
 - 'combinatorial quantum chemistry'
- higher throughput for well scaling applications
 - for some applications (much) better statistics are needed
⇒ longer and/or multiple trajectories.

Example 1: DNA Fiber



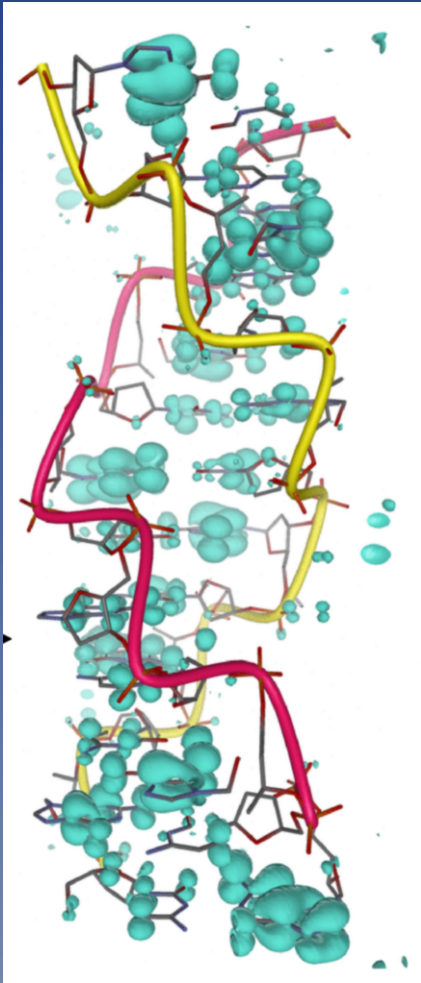
Current project on Earth Simulator
by Prof. Dr. Mauro Boero

Hydrated double strand poly G-C fiber
- 1194 atoms, 12227.7 \AA^3 supercell
- periodic boundary conditions

Car-Parrinello MD with Local Spin Density
approx. including BLYP gradient correction

Norm-conserving pseudopotentials
with 70 Ry plane wave cutoff.

Example 1: DNA Fiber Computational Details



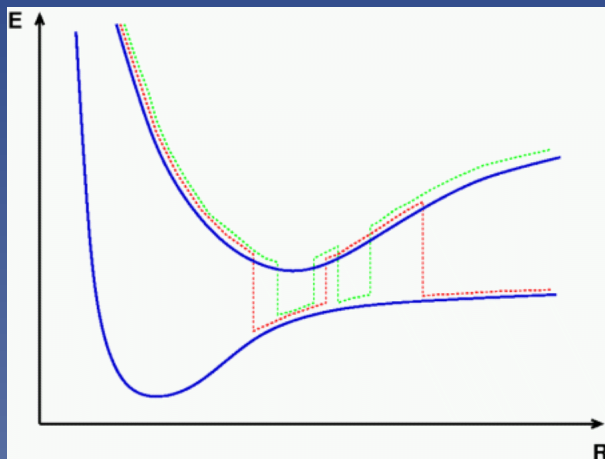
CPU time / MD step with 32 nodes: 110s

more than an order of magnitude more effort than previous “Benchmark System”:
64 water, 192 atoms, 1912.6 \AA^3 supercell

proton transport with $32 \text{ H}_2\text{O} + 1 \text{ H}_3\text{O}^+$
with ultrasoft (25 Ry) pseudopotentials
needs 1 single PC for a week

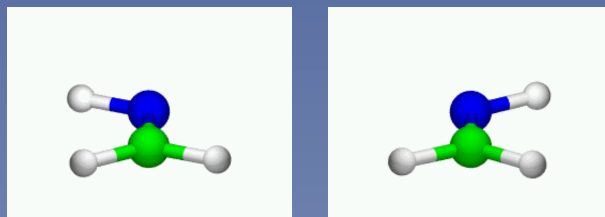
Example 2: CPMD with Surface Hopping

Current project on JUMP
by Dr. Nikos Doltsinis



BO-MD with two wavefunctions
simultaneously: ground state and first
excited singlet

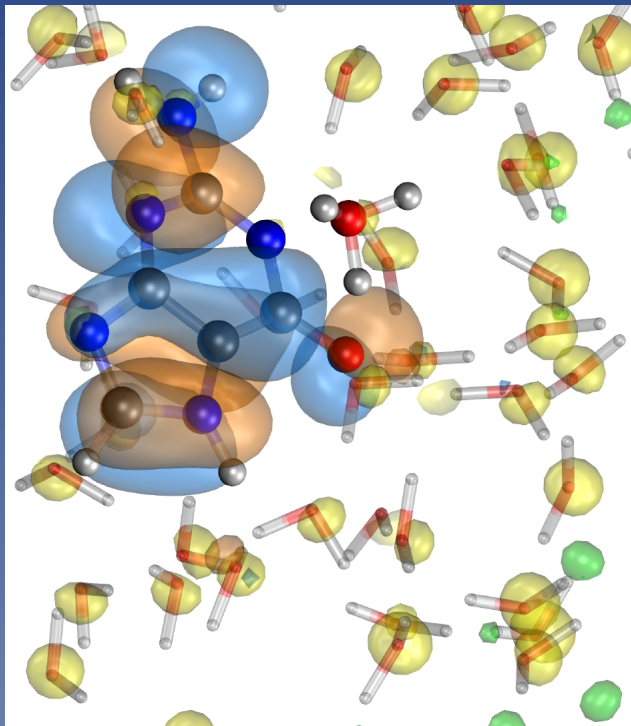
propagation needs (5-10x) smaller
timestep than conventional BO-MD



switch between potential surfaces
(hopping) during simulation

final result (e.g. quantum yield) based
on averaging over many trajectories

Example 2: CPMD with Surface Hopping Details



molecule in gas phase only benchmark,
alternatives are more precise or efficient

advantage of CPMD:
treatment of solvated molecules

singly occupied MOs separate

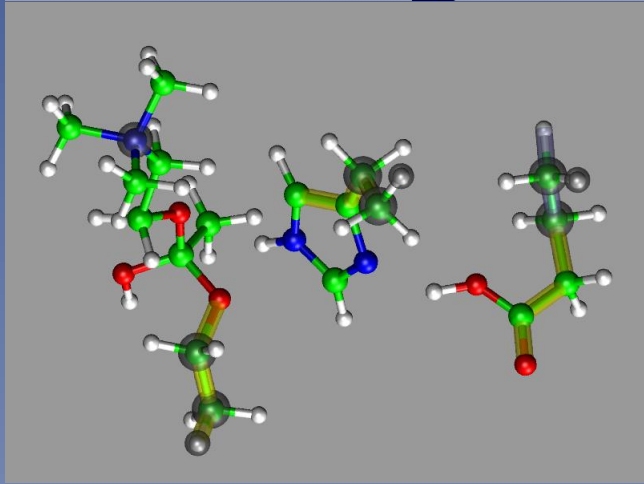
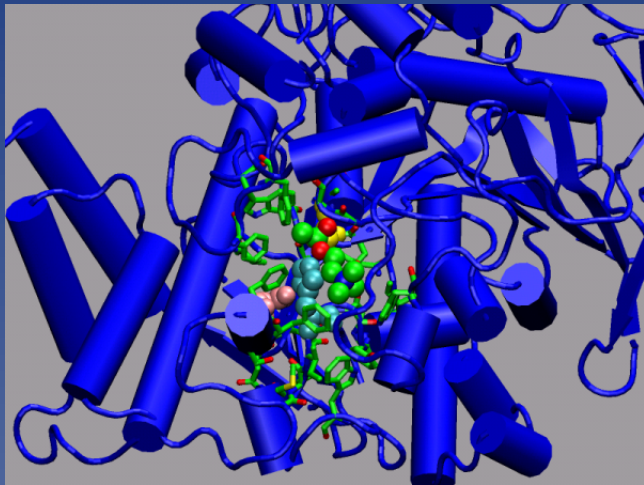
- one stays on guanine
- one delocalized in water

one MD step:

\approx 140 seconds on one 32-CPU frame

\approx one month per single trajectory

Example 3: Catalytic Triade

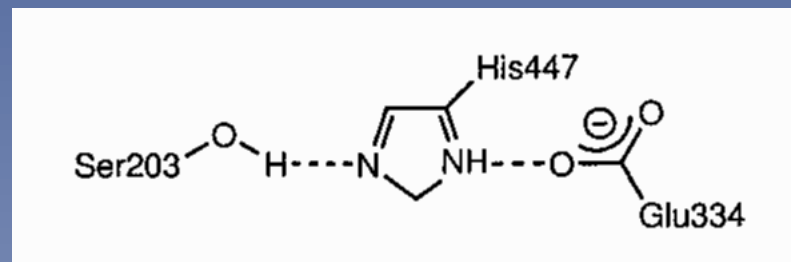


Deacylation of Acetylcholine in AChE,
undergraduate research project at RUB

Chemical neurotransmitter

Inhibition by nerve poisons: snakes
venom, chemical weapons (Sarin)

Prototype system for catalytic triade
mechanism





Example 3: Catalytic Triade Details

simple CPMD simulations of triade subsystem, show influence of whole protein structure on reactivity.

⇒ need at least QM/MM simulation on the timescale of fully classical system MD.

Protonation state of various groups unknown, yet electrostatics play important role.

⇒ full QM treatment of whole (≈ 530 residues) or large part of enzyme desirable.

Trends

- increasing size and complexity of studied topics:
from elementary steps to complex systems of interaction
 - ★ full systems instead of cut-down subsystems
 - ★ complex chemical reactions in condensed phase
 - ★ biochemically relevant processes
 - twofold increase in cpu time demand:
larger systems, longer trajectories for better statistics
- ⇒ *massive* increase in cpu time requirement



Summary

- theoretical chemistry uses and needs large variety of machines
 - software adapts to new hardware, but slowly
 - new algorithms strive for better scaling
 - some codes already scale extremely well
 - new codes are being developed
 - some trivially parallelizable approaches
only viable through massively parallel resources
 - old tools stay in use nevertheless.
- ⇒ diversification of platforms needed to satisfy all demands.



Thanks

- Prof. Dr. Dominik Marx, RUB
- Prof. Dr. Mauro Boero, Tsukuba University
- Dr. Nikos Doltsinis, Holger Langer, RUB
- too many people to name them individually who:
 - ★ provided access to their computer(s)
 - ★ discussed optimization and portability issues

(this page was intentionally left blank)