

the **abdus salam** international centre for theoretical physics

ICTP 40th Anniversary

SMR 1564 - 37

SPRING COLLEGE ON SCIENCE AT THE NANOSCALE (24 May - 11 June 2004)

**BIOMOLECULES; SIMULATION** 

Michele PARRINELLO

CSCS, ETH, Zürich, Switzerland

These are preliminary lecture notes, intended only for distribution to participants.





# Pushing back the frontiers of computer simulations

# Michele Parrinello

#### Department of Chemistry and Applied Biosciences ETH USI Campus, Lugano, Switzerland





# Molecular dynamics



Given a potential energy surface:

$$U(R_1, R_2, \ldots, R_N)$$

The dynamics can be determined from Newton's equation:

$$\boldsymbol{M}_{I}\boldsymbol{\ddot{R}}_{I} = -\nabla \boldsymbol{U}(\boldsymbol{R}_{1},\boldsymbol{R}_{2},\ldots,\boldsymbol{R}_{N})$$



# Empirical potentials

- Molecular mechanics: Intramolecular forces: bond stretch, bending, torsion
- Electrostatic interactions:
   Partial charges, dipoles, polarization
- Van der Waals interactions: Lennard-Jones, Buckingham potential
- Embedded-atom methods: Finnis-Sinclair, Glue model, Daw-Baskes







# Pros and cons

#### PROS

- Efficient
- Accurate in specific cases

#### CONS

- Not transferable
- No chemistry



### Dealing with the electrons





# Hartree-Fock

$$\Psi(x_{1},...,x_{N}) = \frac{1}{\sqrt{N!}} \begin{vmatrix} \varphi_{1}(x_{1}) & \cdots & \varphi_{1}(x_{N}) \\ \vdots & \ddots & \vdots \\ \varphi_{N}(x_{1}) & \cdots & \varphi_{N}(x_{N}) \end{vmatrix}$$
$$\left[ -\frac{1}{2} \nabla + V_{n}(x) + \sum_{j \neq i} \int \frac{1}{|x - x'|} |\varphi_{j}(x')|^{2} dx' - J \right] \varphi_{i}(x) = \varepsilon_{i} \varphi_{i}(x)$$

Exchange operator:

$$J\varphi_i(x) = \frac{1}{2} \sum_j \int \varphi_j(x) \frac{1}{|x-x'|} \varphi_j^*(x') \varphi_i(x') dx'$$



# **Beyond Hartree-Fock**

Adding correlations

Perturbation theory MP2, MP4, ... N<sup>4</sup>, N<sup>5</sup>
 Configuration interaction exp(N)
 Coupled clusters N<sup>6</sup>, N<sup>7</sup>

Small is beautiful!



# Hohenberg-Kohn

The energy of the ground state of a many-body system is a unique functional of the electron density:

 $E = E[\rho_e(r)]$ 

The functional is minimum for the ground state density:

 $E = E[\rho_e(r)] - \mu N$  $\frac{\delta E[\rho_e(r)]}{\delta \rho_e(r)} = \mu$ 



# Kohn-Sham

$$\rho_{e}(r) = 2 \sum_{n} \psi_{n}^{*}(r) \psi_{n}(r)$$

$$E = -\frac{1}{2} \sum_{n} \int dr \psi_{n}^{*}(r) \nabla \psi_{n}(r) + \int dr \rho_{e}(r) V_{ext}(r)$$

$$+ \frac{1}{2} \int dr dr '\rho_{e}(r) \frac{1}{|r-r'|} \rho_{e}(r') + E_{xc} [\rho_{e}(r)]$$

$$\left(-\frac{1}{2} \nabla + V_{ext}(r) + V_{H}(r) + V_{xc}(r)\right) \psi_{n}(r) = \varepsilon_{n} \psi_{n}(r)$$

$$V_{H}(r) = \int dr' \frac{1}{|r-r'|} \rho_{e}(r') \qquad V_{xc}(r) = \frac{\delta E_{xc} [\rho_{e}(r)]}{\delta \rho_{e}(r)}$$



# Born-Oppenheimer

The potential energy surface is defined by the instantaneous ground state electronic energy:

$$\phi(R_1,R_2,\ldots,R_N) = E_0(R_1,R_2,\ldots,R_N)$$

But

$$E_0(R_1,R_2,\ldots,R_N)$$

needs to be approximated. We shall choose a theory which has the right balance between accuracy and computational efficiency.



# Ab-initio MD

$$L = \frac{1}{2} \mu \sum_{n} \int dr \left| \dot{\psi}_{n}(r) \right|^{2} + \sum_{I} \frac{1}{2} M_{I} \dot{R}_{I}^{2} - E_{KS} \left[ \psi_{n}, R_{I} \right]$$
$$+ \sum_{n,m} \Lambda_{n,m} \left( \left\langle \psi_{n} \left| \psi_{m} \right\rangle - \delta_{n,m} \right) \right)$$
$$\omega_{e} \propto \sqrt{\frac{E_{g}}{\mu}} \ll \omega_{I}$$



### Car-Parrinello molecular dynamics







# Some problems





# Going to larger systems



**Field theoretical approach** 

S QM/MM

**DFT-based potentials** 



# Plane wave basis set





# Which basis set?

#### **Plane Waves**

- Orthogonal
- No chemical input
- Convergence easy to check
- Simple algebra
- Memory intensive
- Linear algebra and FFT
- No basis set superposition error

www.cpmd.org

#### Gaussians

- Non orthogonal
- Chemical input
- Convergence less easy to check
- More complex
- Reduced memory
- Quantum Chemistry know-how
- Basis set superposition error

cp2k.berlios.de



# Basis set expansion

Orbitals  $\Phi_i$  are expanded in a set of M basis functions  $\{\chi_{\alpha}\}$ 

$$\Phi_i(r) = \sum_{i=1}^M c_{\alpha i} \, \chi_\alpha(r)$$

**Basis functions** 

- Atomic orbital based (Gaussian, Slater, numerical)
- Plane waves
- Grid based, finite elements, wavelets



# Basis set expansion

$$\begin{split} S_{\alpha\beta} &= \int \chi_{\alpha}^{*}(r) \ \chi_{\beta}(r) \ dr & \text{overlap matrix} \\ H_{\alpha\beta} &= \int \chi_{\alpha}^{*}(r) \ \mathcal{H}(r) \ \chi_{\beta}(r) \ dr & \text{Hamiltonian matrix} \\ E_{ij} &= \delta_{ij} \ \epsilon_{i} & \text{Orbital energies} \\ \hline HC &= SCE \end{split}$$



# Orbitals





# Density matrix

$$P_{\alpha\beta} = 2\sum_{i=1}^{N_e/2} c_{\alpha i} c_{\beta i}^*$$

#### Properties of the density matrix

$$Tr(\mathbf{PS}) = \sum_{\alpha\beta}^{M} P_{\alpha\beta} S_{\alpha\beta} = N_e$$

normalisation

$$P = \frac{1}{2}PSP$$

idempotency



# Density matrix

- Unique
- Electron density

$$\rho(r) = \sum_{\alpha\beta} P_{\alpha\beta} \chi_{\alpha} \chi_{\beta}^{*}$$

• Expectation values

$$\langle \mathcal{O} \rangle = \operatorname{Tr}(P\mathcal{O})$$
  
 $\sum_{i} \epsilon_{i} = \operatorname{Tr}(PH)$ 



#### **Gaussian basis**

**Basis functions** 

$$\chi(r) = x^l y^m z^n \exp[-\alpha r^2]$$

Product of basis functions

$$\chi(r-A) \ \chi(r-B) = \tilde{\chi}(r-C)$$

Localization

$$\exp[-\alpha r^2] \to \mathsf{FFT} \to \exp[-\frac{G^2}{4\alpha}]$$



# DFT with gaussians





### The best of both worlds





# Linear scaling





V

# Standard approach

Solve by diagonalization

$$H_{\mu\nu}C^{\nu i} = S_{\mu\delta}C^{\delta i}\mathcal{E}_{i} \qquad i = 1...N$$
  
Construct P  
$$P^{\mu\nu} = \sum_{i} C^{\mu i}C^{\nu i}$$



# Orbital rotation

$$C(X) = C_0 \cos(\sqrt{X^T S X}) + X \frac{\sin(\sqrt{X^T S X})}{\sqrt{X^T S X}}$$
$$X^T S C_0 = 0 \qquad C(X)^T S C(X) = 1 \ \forall X$$

Minimize the energy with respect to X



# Overall cubic scaling





#### Accurate forces





# Checking the eigenstates





# DNA crystal





2388 atoms, 3960 orbitals DZV(d,p) 22596, TZV(2d,2p) 38688 675, 1100 sec / line search (SP4-32-1.3G) 2.5, 5 h /total Not yet fully cubic (45,43,8 % 3,2,1) Not yet sparse






#### One order of magnitude better



#### -----

## Why should one combine QM and MM?

Bio-systems are typically very large and catalyse

complicated reactions

- Proteins >1000 atoms
- Solvent >10000 atoms
- The active site ~ 100 atoms





## Mixed Quantum-Classical

- highly parallel QM/MM Car-Parrinello hybrid code
   Fully Hamiltonian
- MD driver: CPMD

QMI-Part: CPMID 3.3 (pbcs (2 boxes), plane waves, pseudo potentials, GGAs: BP86, BLYP, PW91, PBE...) n-1 nodes

MM-Part: GROMOS96 + P3M, AMBER) 1 node

#### **Interface Region**

Quantum Region (Car-Parrinello)

#### **Classical Region**

#### 77/

#### **QM/MM- Car-Parrinello Simulations**

- Development of improved QM/MM interfaces:
  - pseudo potentials for boundary atoms
  - efficient treatment of long-range electrostatics
  - electron spill out problem



#### The Bonded Part



•boundary atoms: monovalent pseudopotential

• distances - angles - torsions involving MM and QM atoms come from the force field



## Mixed Quantum-Classical Simulations











Explicit dependence of q<sup>ESP</sup> on the positions of all the MM and QM atoms and on the electronic density



#### QM/MM



A. Laio, J. VandeVondele and U. Röthlisberger, Chem. Phys. 116, 6941(2002)



#### Capturing the complexity





#### The way of the future?





## The quickstep team

J. Hutter, University of Zurich J. VandeVondele, University of Cambridge W. I-Feng Kuo, LLNL C. Mundy, LLNL Fawzi Mohammed, ETH Lugano M. Krack, ETH Lugano G. Lippert, BASF M. McGrath, LLNL I. Siepman, LLNL

## Why should one combine QM and MM?

Bio-systems are typically very large and catalyse complicated reactions



- Solvent >10000 atoms
- The active site ~ 100 atoms





# Mixed Quantum-Classical QM/MM- Car-Parrinello Simulations

- highly parallel QM/MM Car-Parrinello hybrid code
- Fully Hamiltonian
- MD driver: CPMD

QM-Part: CPMD 3.3 (pbcs (2 boxes), plane waves, pseudo potentials, GGAs: BP36, BLYP, PW91, PBE...) n-1 nodes

MM-Part: GROMOS95 + P3M, AMBER) 1 node **Interface Region** 

Quantum Region (Car-Parrinello)

#### **Classical Region**

#### 777

#### **QM/MM- Car-Parrinello Simulations**

- Development of improved QM/MM interfaces:
  - pseudo potentials for boundary atoms
  - efficient treatment of long-range electrostatics
  - electron spill out problem



#### The Bonded Part



boundary atoms: monovalent pseudopotential

• distances - angles - torsions involving MM and QM atoms come from the force field



## Mixed Quantum-Classical Simulations





ρ(**nr1i**,**nr2j**,**nr3k**)





 charges, located on the QM atoms, are fitted to the electrostatic fields on the MM atoms due to the electronic charge distribution



Explicit dependence of q<sup>ESP</sup> on the positions of all the MM and QM atoms and on the electronic density



#### QM/MM



A. Laio, J. VandeVondele and U. Röthlisberger, Chem. Phys. 116, 6941(2002)



## Capturing the complexity





#### The way of the future?





#### The time scale problem

Direct simulation allows only very short runs:  $\sim 10$  ps for ab-initio MD,  $\sim 10$  ns for classical MD

Many relevant phenomena take place on a larger time scale: chemical reactions, conformational changes, protein folding, etc.

Two-fold strategy:

a) Finding reactive pathsb) Exploring the free energy surface



#### Activated events





#### The quantum chemical approach

**×** Find the saddle point on the PES

**X** Use transition state theory

**×** Correct for zero point motion

#### 77;

## Many different solutions proposed

- Thermodynamic integration
- "Flattening" the surface (hyperdynamics, puddle-skimming, umbrella sampling, etc.)
- Trajectory-based schemes (reaction path sampling, Lagrangian action minimization, nudged elastic band, etc.)

## Finding the saddle points (eigenvalue following, dimer method, hessian-based methods, etc.)

- Temperature enhanced sampling (histogram reweighting, parallel tempering, etc.)
- Etc. etc.

## Driving the reaction

Suppose that the reaction coordinate q is known:

$$q(R_1, R_2, \ldots, R_N) = q$$

We force the reaction by adding a constraint term to the dynamics with a Lagrange multiplier:

$$\lambda (q(R_1, R_2, ..., R_N) - q)$$
  
$$\langle \lambda \rangle \approx \frac{\partial F}{\partial q}$$
  
The activation energy is: 
$$\Delta F = \int_{q_A}^{q^*} dq \langle \lambda \rangle_q$$



#### The right reaction path?



Activated events are often intrinsically multidimensional!!!



## Life is complicated





## Life is complicated





## How to explore a multidimensional free energy surface?

Need to be able to escape free energy minima

Our solution: Non-Markovian coarse-grained dynamics

A. Laio and M. Parrinello, PNAS



#### Collective variables

Choose a small set of slow collective variables:

i = 1, n

The s<sub>i</sub> :

- Discriminate between reactants and products
- Include all the relevant slow modes
- \* The reaction coordinate is a linear combination of the  $s_i$

#### 77/

#### Examples of collective variables

- Distances
- Angles: bending and torsional
- Coordination numbers: between individual atoms
- or between different species
- Local electric fields
- Number of n-fold rings
- Solvation energy
- Lattice vectors
- Energy
- Etc. etc.



### Probability distribution

$$P(S_1, S_2, \dots, S_n) = \frac{\int d\vec{R}_I \prod_{i}^n \delta\left(S_i - s_i(\vec{R}_I)\right) e^{-\beta V\left(\vec{R}_I\right)}}{\int d\vec{R}_I e^{-\beta V\left(\vec{R}_I\right)}}$$

We want to study the free energy as a function of these variables:

$$\beta F(S_1, S_2, ..., S_n) = -\ln P(S_1, S_2, ..., S_n)$$



#### The algorithm:

Wherever you go put a "small" Gaussian
Always move in the direction of the direction that minimizes the sum of V(s) and all the Gaussians



#### 77;





## NaCl in water

Two minima:	Contact ion pair (metastable) Dissociated
Collective coordinates:	Electric field on Na <sup>+</sup> Electric field on Cl <sup>-</sup> Distance Na <sup>+</sup> Cl <sup>-</sup>
Classical MD:	Amber force field




#### Free energy surface





#### Dialanine in water

1 dialanine in 287 TIP3P water

AMBER95 force field

Collective coordinates: backbone dihedral angles  $\Phi$  and  $\Psi$ 







## Stationary action principle

$$L = \frac{1}{2}\dot{q}^2 - V(q)$$

$$\ddot{a}S = \ddot{a}\int_{t_A}^{t_B} L(q(t), \dot{q}(t)) dt = 0$$

$$q(t_A) = q_A \qquad q(t_B) = q_B$$

Can we use this principle to find the trajectories?

Saddle point!



# Folding a small peptide





#### Motivation and aims

Oxidatitive damage to DNA is common and has fatal consequences

Guanine, having the lowest oxidation potential among the DNA bases, plays a fundamental role

Does the structure of DNA funnel the reactions towards a unique product?



#### A continuous metadynamics



M. Iannuzzi, M. Laio and M. Parrinello, PRL 2003



## From azulene to naphthalene







### Science fiction?



L.T. Scott, Acc. Chem. Res. 52, 15 (1982)



### G<sup>+</sup> localization how?

G<sup>+.</sup> can be formed: Directly 1. Hole 2. Via hole migration Migration

#### Chem. Rev. 98, 1109-1151 (1998)

77/

## Computational details

Model G:C decamer, Z conformation

X-ray structure available

Rich in G and the smallest cell Electronic structure treated by the Kohn-Sham method

BLYP and HCTH functionals, plane waves (70Ry cutoff)

Martin-Troullier pseudopotentials

Car-Parrinello molecular dynamics, CPMD code



777

#### Computational details

The hole localizes on G<sup>+</sup> and the proton moves to C





The model includes water and counter ions 1194 atoms Full quantum: 3,960 valence electrons, 408,238 PW!



### The HOMO of DNA

#### DNA laboratory specimen The full Monty

- 1194 atoms 1980 states
- 408.238 PW 13 Gb



F. Gervasio, P. Carloni and M.P., PRL

#### HOMO





### Fate of G<sup>+.</sup> in DNA

In duplex DNA the protonation is not clear

CH<sup>+</sup> has a pKa of 4.3 (G<sup>+</sup>: pKa 3.9)

\_ Sharing of proton expected



 $K_{exp} = 2.5$  is predicted from experiments in water Is this relevant for DNA?



## Bridging length scale

Most calculations were done in the QM/MM framework.

Laio A., VandeVondele J. and Roethlisberger U., J.Chem. Phys. 116, 6941(2002)

Due to high polarity of bonds H-capping was used. Added H were decoupled to MM.





Quantum sub-systems of increasing size were used.

The biggest quantum model was the full DNA.

## Protonation state of the bases: gas phase



## Protonation state of the bases: DNA



DNA $\Delta E(kcal/$	mol):
Gas phase:	$\sim 2./3.$ (inversion)
QM/MM	-4.5
No charge	-1.2
Quantum PO <sub>4</sub> <sup>3</sup>	-/full -3.7 /-3.9
Charge on PO <sub>2</sub>	$^{3-}$ and sugar $-5.3$
$\Delta\Delta E \sim 7-8$	

#### Sources

Backbone charges	~3
Error (QM/MM, exclusion)	<1
Geometry	3-4



#### Fate of G<sup>+.</sup> in DNA

In duplex DNA the oxidation product is 8-oxo-guanine (8oxoG)

30,000 8-oxoG in human genome! *Methods Enzymol.* **186**, *521 (1990)* 

In water a variety of products are formed



777

## Oxidation reaction 1



Meta coordinates used: N1 and water oxygen H-coordination number Carbon 8 O-coordination number Rate limited by water autoprotolysis Catalyzed by  $R_2PO_4^-$ N1 protonation state matters  $\Delta E$  N1 deprotonation of 8-OH-G:



## Oxidation reaction

