

Self-consistent field methods for large molecular systems

Trygve Helgaker, University of Oslo, Norway

Filip Pawłowski, Simen Reine, Erik Tellgren

Stinne Høst, Branislav Jansík, Poul Jørgensen, Pekka Manninen, Jeppe Olsen,
University of Aarhus, Denmark

Sonia Coriani, University of Trieste, Italy

Paweł Sałek, Royal Institute of Technology, Sweden

Central European Symposium on Theoretical Chemistry 2006

September 24–27, 2006

Hotel Skalny, Zakopane, Poland

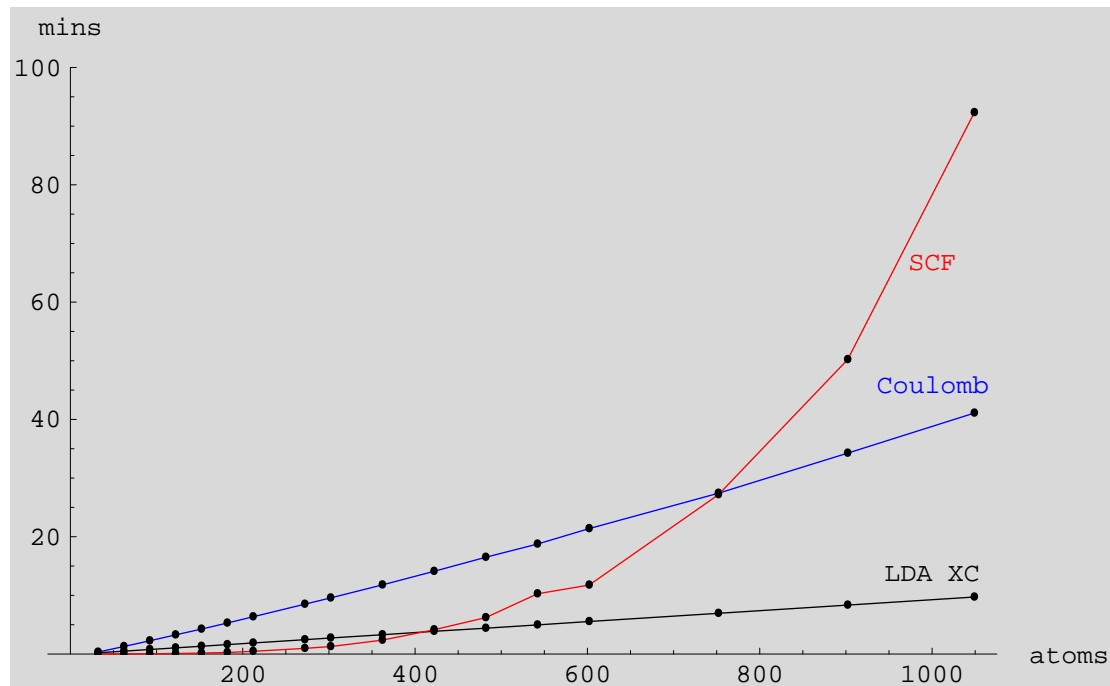
Self-consistent field (SCF) theory

- Consider the optimization of the SCF energy (here LDA) of molecular systems:

- small systems dominated by KS-matrix evaluation, with **linear scaling**
- large systems dominated by SCF diagonalization, with **cubic scaling**

$$\mathbf{FC} = \mathbf{SC}\epsilon$$

$$\mathbf{D}_{\text{new}} = \mathbf{C}_{\text{occ}} \mathbf{C}_{\text{occ}}^{\text{T}}$$



- To achieve linear scaling, we must avoid diagonalization and MOs!
- We shall here consider an **alternative to diagonalization**:
 - it optimizes the density matrix directly, avoiding MOs
 - involves only additions and multiplications of (sparse) one-electron matrices
 - for large (sparse systems), the calculations scale linearly with system size

Direct optimization of the density matrix

- Consider the direct optimization of the density matrix:

$$E(\mathbf{D}) = \text{Tr } \mathbf{D}\mathbf{h} + \text{2-el. part}$$

- there are constraints on the density matrix:

$$\underbrace{\mathbf{D} = \mathbf{D}^T}_{\text{symmetry}}, \quad \underbrace{\text{Tr } \mathbf{D} = N}_{\text{trace}}, \quad \underbrace{\mathbf{D}^2 = \mathbf{D}}_{\text{idempotency}} \quad (\text{orthonormal basis})$$

- any optimization must obey these constraints

- Many strategies are based on purification of the density matrix

$$\tilde{\mathbf{D}} = 3\mathbf{D}^2 - 2\mathbf{D}^3 \quad (\text{McWeeny purification, 1960})$$

- Li, Nunes and Vanderbilt (1993)

$$\tilde{E} = \text{Tr } \tilde{\mathbf{D}}\mathbf{h} + \mu(\text{Tr } \mathbf{D} - N) + \text{2-el. part}$$

- Millam and Scuseria (1997), Challacombe (1999)
- Palser and Manolopoulos (1998), Niklasson (2002)

- We shall pursue a different approach, based on an explicitly parameterization of \mathbf{D}

Exponential parameterization of the density matrix

- In a real, nonorthogonal AO basis, with $\mathbf{S} \neq \mathbf{I}$, let \mathbf{D} be a valid HF/KS matrix:

$$\underbrace{\mathbf{D} = \mathbf{D}^T}_{\text{symmetry}}, \quad \underbrace{\text{Tr} \mathbf{D} \mathbf{S} = N}_{\text{trace}}, \quad \underbrace{\mathbf{D} \mathbf{S} \mathbf{D} = \mathbf{D}}_{\text{idempotency}}$$

- Any other valid density matrix $\mathbf{D}(\mathbf{X})$ can then be generated from this matrix:

$$\underbrace{\mathbf{D}(\mathbf{X}) = \exp(-\mathbf{X} \mathbf{S}) \mathbf{D} \exp(\mathbf{S} \mathbf{X})}_{\text{exponential parameterization}}, \quad \underbrace{\mathbf{X}^T = -\mathbf{X}}_{\text{antisymmetric}}$$

- Helgaker, Jørgensen and Olsen: *Molecular Electronic-Structure Theory* (Wiley, 2000)
- Head-Gordon and coworkers, MolPhys **101**, 37 (2003), JCP **118**, 6144 (2003)
- We can obtain any valid density matrix, in the AO basis, without recourse to MOs!
 - in particular, we may optimize the energy by freely varying $X_{\mu\nu}$ with $\mu > \nu$:

$$E_{\min}(\mathbf{X}) = \min_{\mathbf{X}} [\text{Tr} \mathbf{D}(\mathbf{X}) \mathbf{h} + \text{2-el. part}]$$

- Is the use of $\mathbf{D}(\mathbf{X})$ a practical proposition?
 - we shall in this talk demonstrate that it is indeed so
 - we shall consider energy optimizations and property calculations

Two questions about $\mathbf{D}(\mathbf{X}) = \exp(-\mathbf{X}\mathbf{S})\mathbf{D}\exp(\mathbf{S}\mathbf{X})$

- Can it be evaluated efficiently?

- we use a generalized Baker–Campbell–Hausdorff (BCH) expansion:

$$\mathbf{D}(\mathbf{X}) = \mathbf{D} + [\mathbf{D}, \mathbf{X}]_S + \frac{1}{2} [[\mathbf{D}, \mathbf{X}]_S, \mathbf{X}]_S + \dots$$

- we have here introduced the S commutator

$$[\mathbf{D}, \mathbf{X}]_S = \mathbf{D}\mathbf{S}\mathbf{X} - \mathbf{X}\mathbf{S}\mathbf{D}$$

- converges rapidly (purification may be necessary), in about 10 matrix multiplications

- Are redundancies a problem?

- the AO space consists of two parts: the occupied space and the virtual space

$$\mathbf{P} = \mathbf{D}\mathbf{S} \text{ (onto occupied space), } \quad \mathbf{Q} = \mathbf{I} - \mathbf{D}\mathbf{S} \text{ (onto virtual space)}$$

- only rotations between the occupied and virtual spaces are nonredundant:

$$\mathbf{X} = \underbrace{\mathbf{P}\mathbf{X}\mathbf{P}^T + \mathbf{Q}\mathbf{X}\mathbf{Q}^T}_{\text{redundant}} + \underbrace{\mathbf{P}\mathbf{X}\mathbf{Q}^T + \mathbf{Q}\mathbf{X}\mathbf{P}^T}_{\mathbf{X}_{\text{ov}}}$$

- to avoid problems with redundancies, we use the projected parameterization

$$\mathbf{D}(\mathbf{X}) = \exp(-\mathbf{X}_{\text{ov}}\mathbf{S})\mathbf{D}\exp(\mathbf{S}\mathbf{X}_{\text{ov}}), \quad \mathbf{X}^T = -\mathbf{X}$$

Diagonalization-free Roothaan–Hall SCF optimization

- The SCF (Fock or Kohn–Sham) energy may, in principle, be optimized directly:

$$E_{\min} = \min_{\mathbf{X}} E(\mathbf{X}) \quad \Leftrightarrow \quad \underbrace{\mathbf{F}(\mathbf{D})\mathbf{D}\mathbf{S} = \mathbf{S}\mathbf{D}\mathbf{F}(\mathbf{D})}_{\text{stationary condition}}$$

– a difficult global minimization problem!

- In MO theory, the Roothaan–Hall SCF scheme works well, especially with DIIS:

$$\mathbf{F} = \mathbf{h} + \mathbf{g}(\mathbf{D}) \quad \begin{array}{c} \xrightarrow{F} \\ \xleftrightarrow{D} \\ \xleftarrow{D} \end{array} \quad \mathbf{F}\mathbf{C} = \mathbf{S}\mathbf{C}\epsilon; \quad \mathbf{D}_{\text{new}} = \mathbf{C}_{\text{occ}}\mathbf{C}_{\text{occ}}^T$$

– each diagonalization is equivalent to minimizing the sum of the (occ.) orbital energies

$$\epsilon(\mathbf{X}) = \sum_I \epsilon_I = \text{Tr } \mathbf{D}(\mathbf{X})\mathbf{F}$$

- By analogy with MO theory, we set up the following Roothaan–Hall SCF scheme:

$$\mathbf{F} = \mathbf{h} + \mathbf{g}(\mathbf{D}) \quad \begin{array}{c} \xrightarrow{F} \\ \xleftrightarrow{D} \\ \xleftarrow{D} \end{array} \quad \epsilon_{\min} = \min_{\mathbf{X}} \text{Tr } \mathbf{D}(\mathbf{X})\mathbf{F}; \quad \mathbf{D}_{\text{new}} = \mathbf{D}(\mathbf{X}_*)$$

– at each SCF iteration, we minimize $\text{Tr } \mathbf{D}(\mathbf{X})\mathbf{F}$ with respect to \mathbf{X}

– the new density is then obtained by expansion of $\mathbf{D}(\mathbf{X})$ with the minimizer \mathbf{X}_*

- We thus avoid MOs and diagonalization but retain the SCF iterations

Newton minimization of the Roothaan–Hall energy function

- At each SCF iteration, our task is to minimize the Roothaan–Hall energy function

$$\varepsilon(\mathbf{X}) = \text{Tr } \mathbf{D}(\mathbf{X})\mathbf{F} = \text{Tr } \mathbf{D}\mathbf{F} + \text{Tr } [\mathbf{D}, \mathbf{X}]_S \mathbf{F} + \frac{1}{2} \text{Tr } [[\mathbf{D}, \mathbf{X}]_S, \mathbf{X}]_S \mathbf{F} + \dots$$

- Truncating at second order and setting the gradient to zero, we obtain the Newton step:

$$\mathbf{HXS} + \mathbf{SXH} = \mathbf{G} \leftarrow \text{the Roothaan–Hall (RH) Newton equation}$$

- where the (negative) gradient and Hessian matrices are given by

$$\mathbf{G} = \mathbf{F}^{\text{vo}} - \mathbf{F}^{\text{ov}}$$

$$\mathbf{H} = \mathbf{F}^{\text{vv}} - \mathbf{F}^{\text{oo}}$$

$$\mathbf{F} = \mathbf{F}^{\text{oo}} + \mathbf{F}^{\text{ov}} + \mathbf{F}^{\text{vo}} + \mathbf{F}^{\text{vv}}$$

- A RH diagonalization corresponds to an exact minimization (many Newton steps)
 - however, a **partial minimization** will do
 - in fact, one RH Newton step is usually sufficient
- Because of their large dimensions, the Newton equations cannot be solved directly
 - we use an iterative scheme: **the conjugate-gradient method**
 - key step: repeated evaluation of the residual $\mathbf{R} = \mathbf{G} - \mathbf{HXS} - \mathbf{SXH}$
 - all operations are elementary (sparse) matrix manipulations

Solution of the Roothaan–Hall Newton equations

- At each SCF iteration, we solve the Roothaan–Hall Newton equations

$$\mathbf{HXS} + \mathbf{SXH} = \mathbf{G}$$

- a naïve application of the CG method converges slowly
- the equations are ill-conditioned since $\kappa(\mathbf{H})\kappa(\mathbf{S}) \gg 1$ (κ is the condition number)

- The equations may be made well-conditioned by a **Löwdin orthonormalization**

$$\tilde{\mathbf{H}}\mathbf{Z} + \mathbf{Z}\tilde{\mathbf{H}} = \tilde{\mathbf{G}}, \quad \tilde{\mathbf{A}} = \mathbf{S}^{-1/2}\mathbf{A}\mathbf{S}^{-1/2}$$

- convergence is greatly improved since $\kappa(\tilde{\mathbf{H}}) = \kappa(\mathbf{S}^{-1/2}\mathbf{H}\mathbf{S}^{-1/2}) \ll \kappa(\mathbf{H})\kappa(\mathbf{S})$
- we obtain $\mathbf{S}^{-1/2}$ by an iterative scheme (B. Jansík)
- orthogonalization is also possible by **Cholesky decomposition** $\mathbf{S} = \mathbf{U}^T\mathbf{U}$
- Millam and Scuseria (1996), Challacombe (1998), Head-Gordon *et al.* (2003)

- Further **diagonal preconditioning** cuts the number of iterations by one half

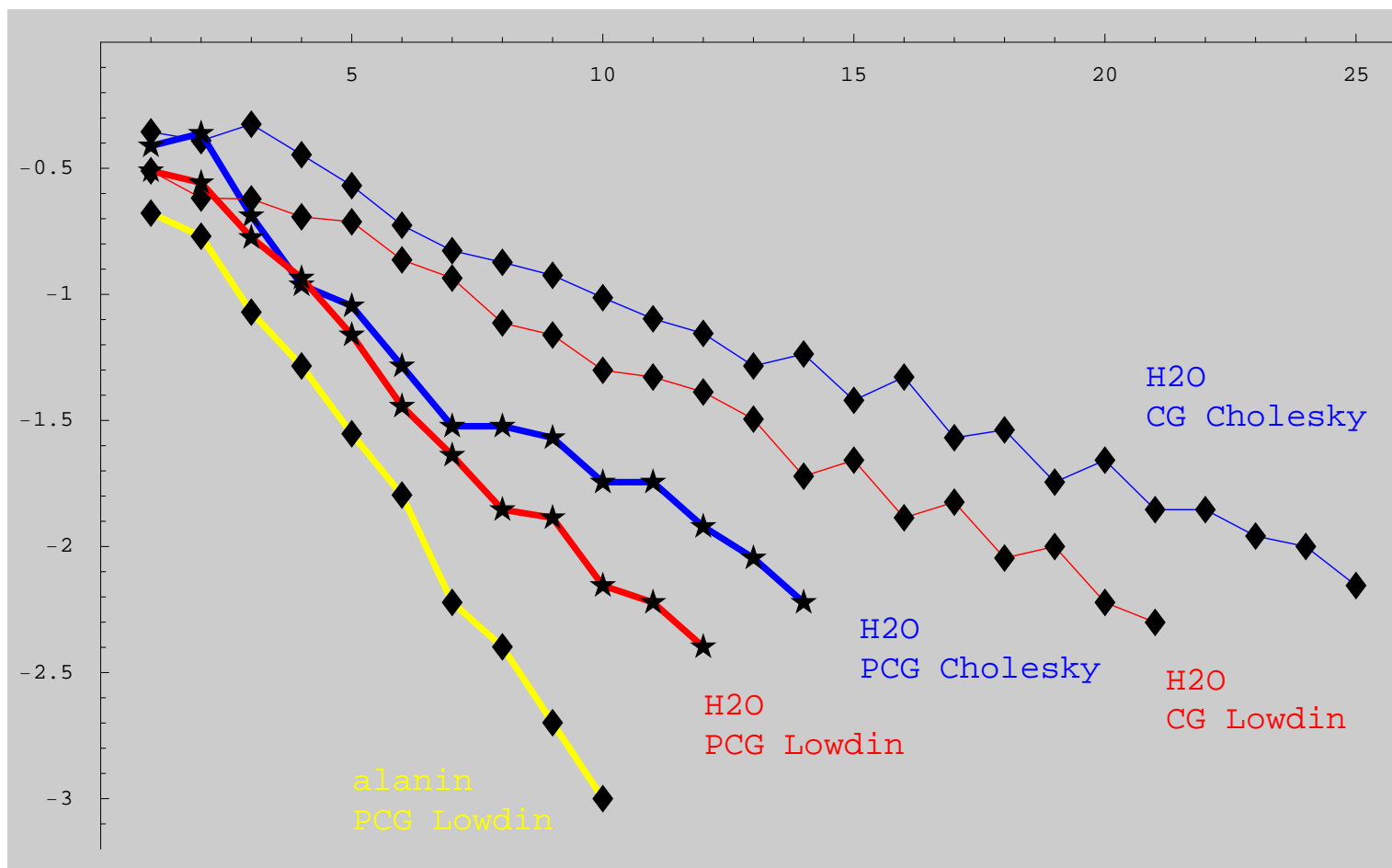
$$\tilde{\tilde{\mathbf{H}}} = \tilde{\mathbf{H}}_{\text{diag}}^{-1} \tilde{\mathbf{H}}$$

- 10 iterations typically reduce the residual by two orders of magnitude

- About 70 multiplications needed for one Newton iteration with density-matrix generation

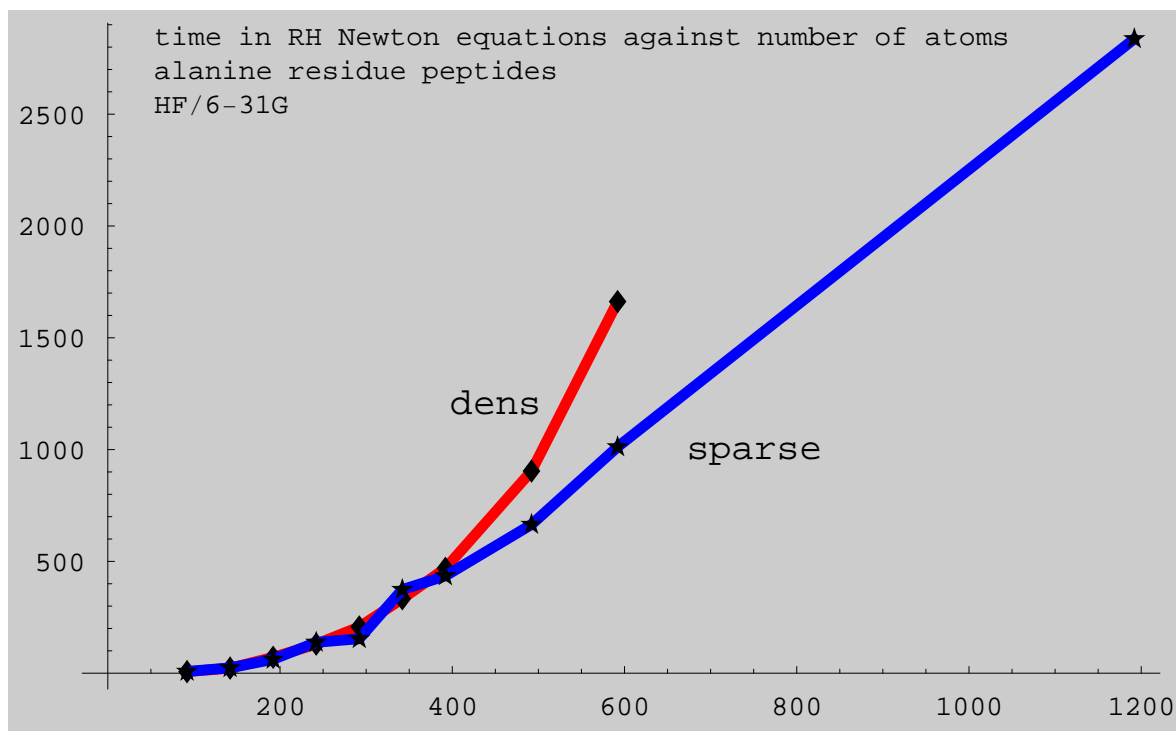
Iterative solution of Roothaan–Hall Newton equations

- Logarithmic plots of the residual against the number of iterations
 - H₂O, LDA/t-aug-cc-pVTZ
 - 99 alanine residue peptides, LDA/6-31G (5449 AOs)



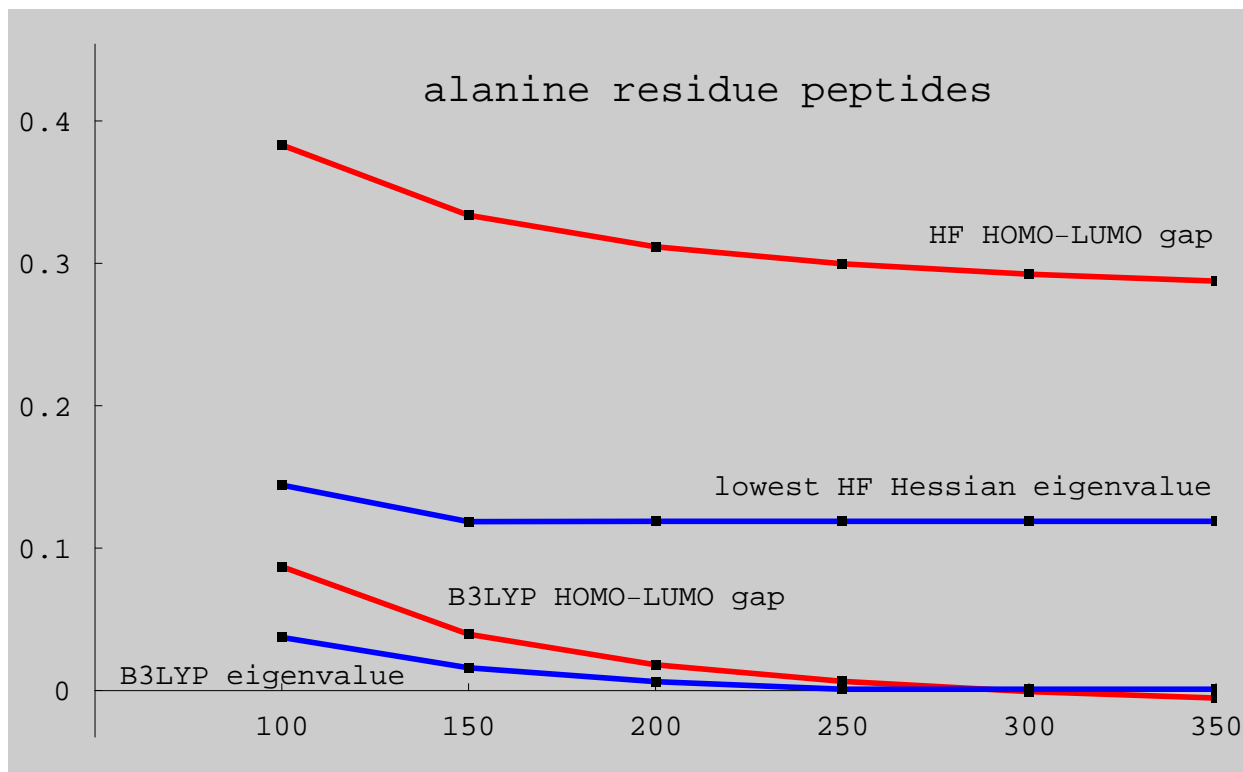
CPU time spent in Roothaan–Hall Newton equations

- We have successfully avoided Fock/Kohn–Sham diagonalization
 - minimization rather than the solution of a generalized eigenvalue problem
 - rapidly convergent: 50–100 sparse matrix multiplications needed
- Linear scaling is obtained by employing sparse-matrix algebra
 - compressed sparse-row (CSR) representation of few-atom blocks



SCF optimizations in small and large molecules

- Diagonalization can be avoided by solving Newton equations
- However, SCF convergence is typically more difficult in larger systems
 - small (or negative) HOMO–LUMO gaps and small Hessian eigenvalues in DFT
 - 6-31G HOMO–LUMO gap and lowest Hessian eigenvalue in alanine residue peptides



- We have modified the standard SCF scheme, to make it more robust

The trust-region self-consistent field (TRSCF) method

- SCF optimizations have two ingredients
 1. Roothaan–Hall minimization (diagonalization): $\min_{\mathbf{X}} \text{Tr } \mathbf{D}(\mathbf{X})\mathbf{F}$
 2. DIIS-type averaging of density matrices: $\bar{\mathbf{D}} = \sum_{i=0}^n c_i \mathbf{D}_i$, $\sum_i c_i = 1$
- In the Roothaan–Hall step, we minimize subject to a constraint on the step size:
 - this amounts to a simple level shifting of the Fock/Kohn–Sham matrix

$$\mathbf{F} \rightarrow \mathbf{F}(\mu) = \mathbf{F} - \mu \mathbf{S}^{\text{oo}} \quad (\text{only occupied–occupied part shifted})$$

- for $\mu > 0$, the HOMO–LUMO gap increases, making large steps unfavourable
 - μ is adjusted until step is of desired length, during the iterative solution
- In the averaging step, we construct a second-order model of the SCF energy

$$E_{\text{DSM}}(c_i) \approx E_{\text{SCF}}(3\bar{\mathbf{D}}\mathbf{S}\bar{\mathbf{D}} - 2\bar{\mathbf{D}}\mathbf{S}\bar{\mathbf{D}}\mathbf{S}\bar{\mathbf{D}}), \quad \bar{\mathbf{D}} = \sum_i c_i \mathbf{D}_i$$

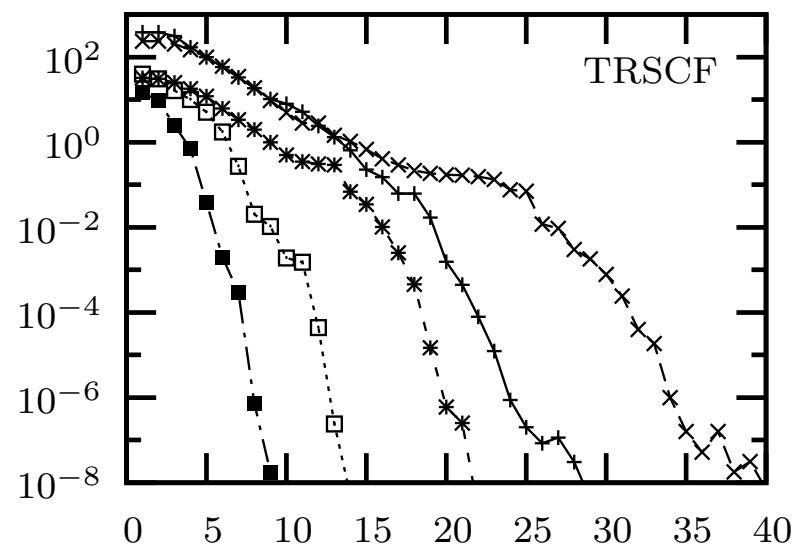
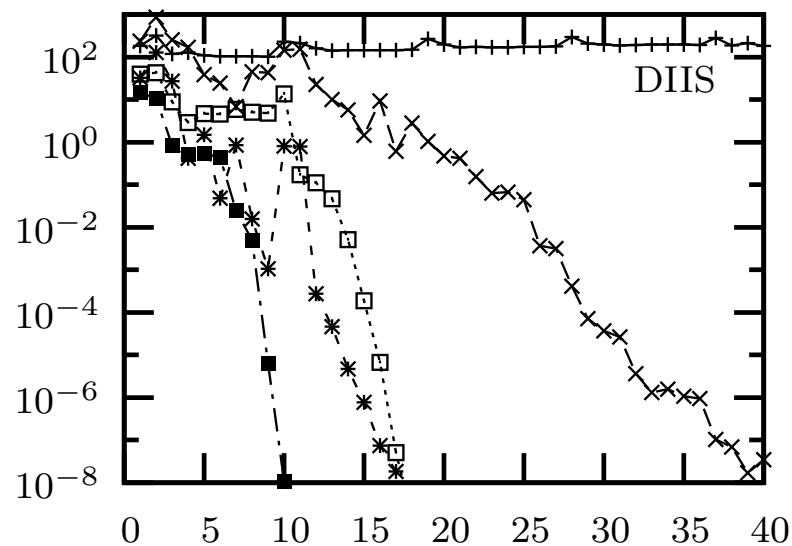
- it has the correct gradient but an approximate Hessian
 - the best density matrix is obtained by minimization, subject to a step-size constraint

$$\min_{c_i} E_{\text{DSM}}(c_i) \leftarrow \text{density-subspace minimization (DSM)}$$

- JCP **121**, 15 (2004); JCP **123**, 074103 (2005)

The TRSCF method (continued)

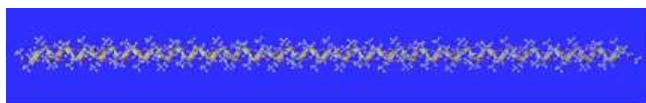
- The TRSCF method gives a stable and uniform convergence towards the SCF minimum
- Convergence of LDA calculations for a variety of molecules
 - zinc complex(+), rhodium complex(\times), cadmium complex(*), CH_3CHO (\square) and H_2O (\blacksquare)



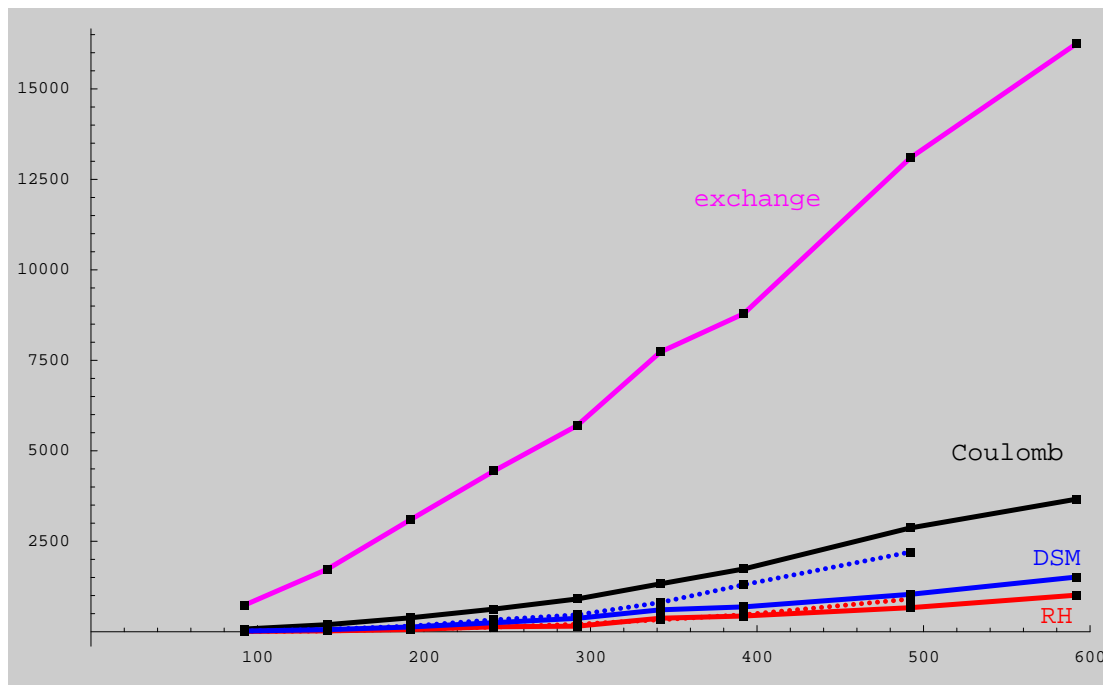
- Problems arise with small Hessian eigenvalues and small (negative) HOMO–LUMO gaps
 - Note: a Newton step of $\text{Tr } \mathbf{FD}$ is a quasi-Newton step of $E_{\text{SCF}}(\mathbf{D})$
 - we are doing a quasi-Newton optimization of the energy
 - the quasi-Newton step is poor for small or negative HOMO–LUMO gaps
 - revert to full Newton if necessary: $\min \text{Tr } \mathbf{FD} \rightarrow \min E_{\text{SCF}}(\mathbf{D})$

Illustration: alanine residue peptides

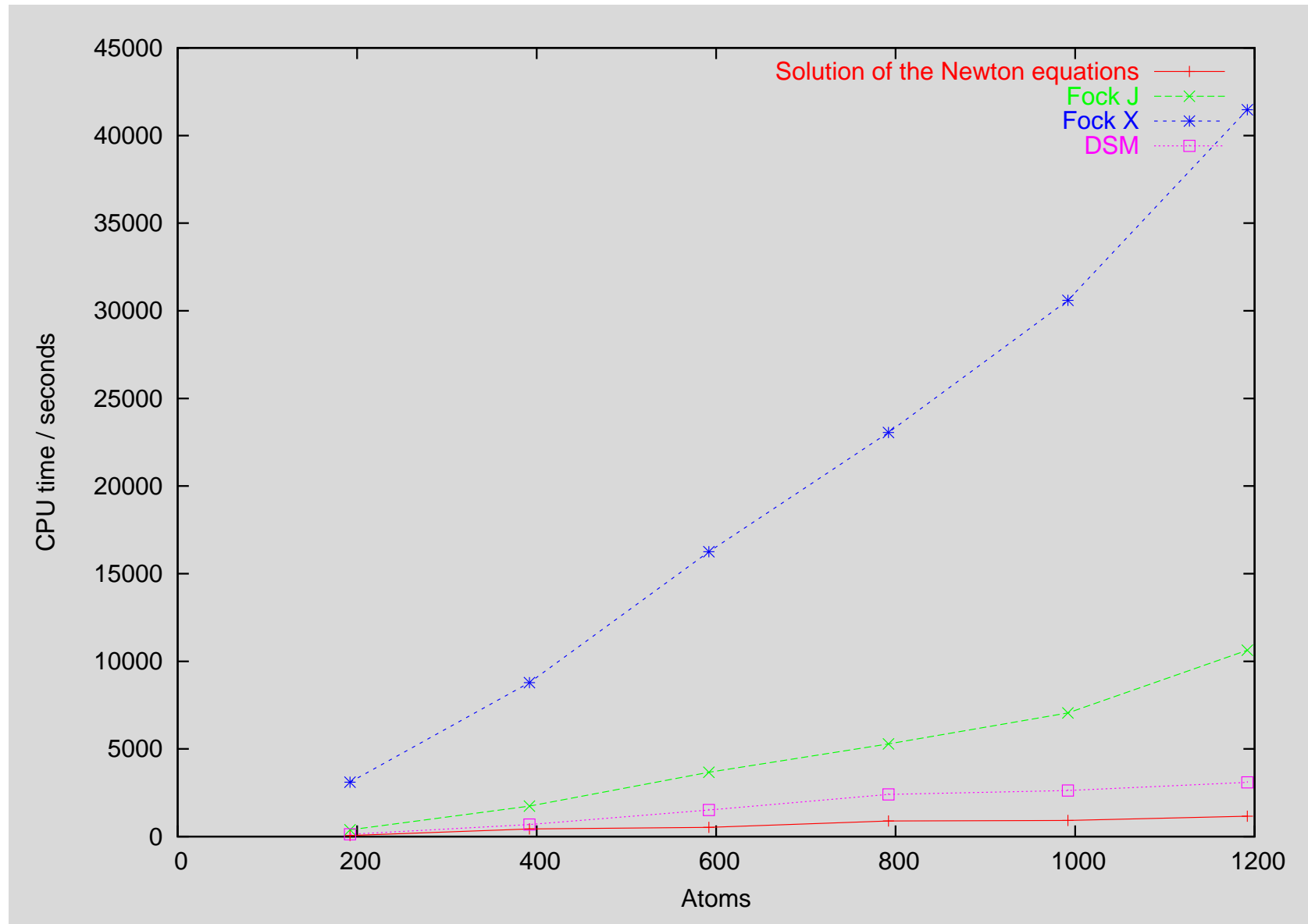
- Features of the code
 - diagonalization-free trust-region Roothaan–Hall (TRRH) energy minimization
 - trust-region density-subspace minimization (TRDSM) for density averaging
 - boxed density-fitting with FMM for Coulomb evaluation (Simen Reine)
 - LinK for exact exchange, linear-scaling exchange-correlation evaluation
 - compressed sparse-row (CSR) representation of few-atom blocks



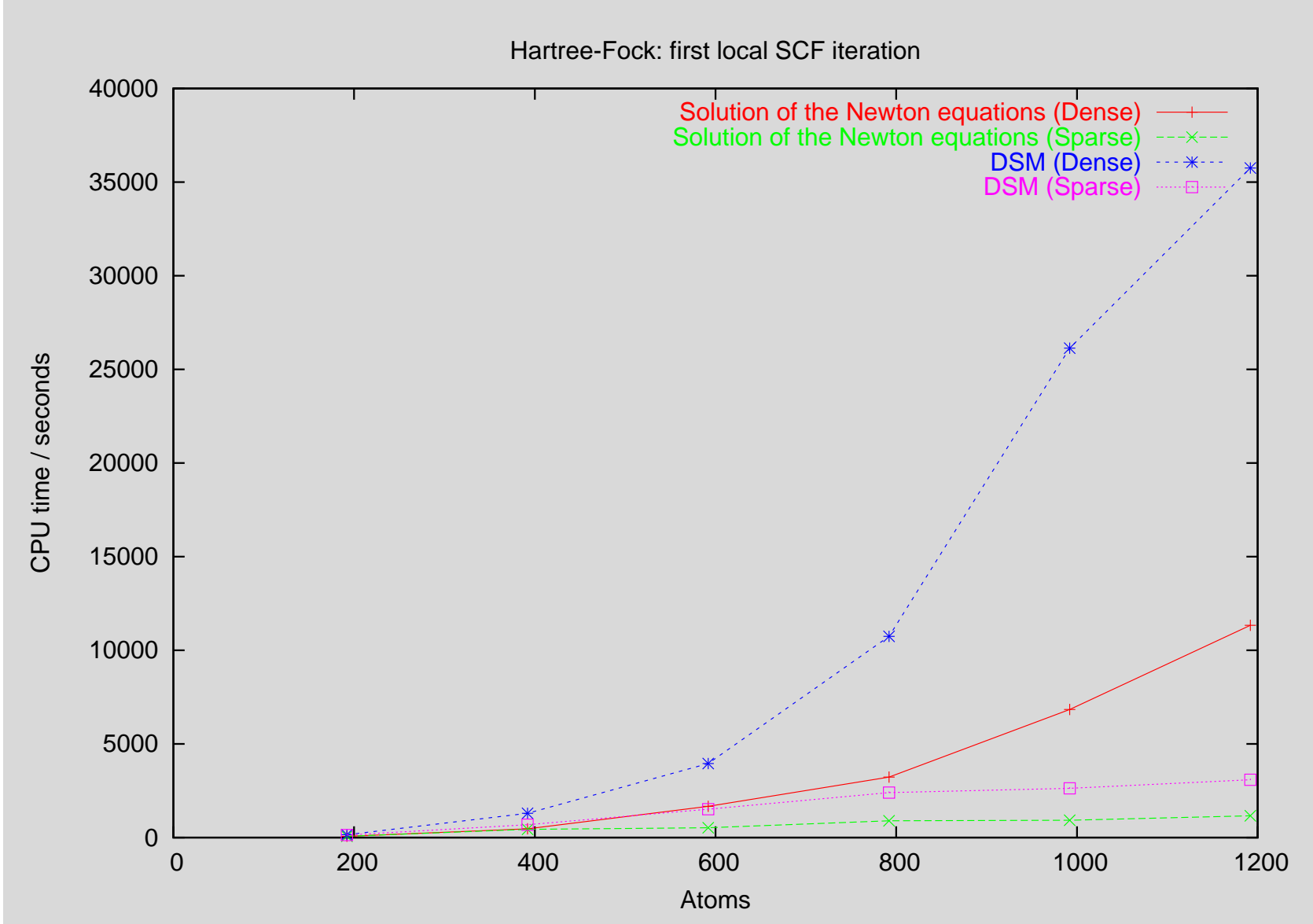
- alanine residue peptides
 - CPU time against atoms
 - HF/6-31G
 - 5th SCF iteration
 - dominated by exchange
 - RH step least expensive
 - full lines: sparse algebra
 - dashed lines: dens algebra



Alanine residue peptides: timings in first local SCF iteration (HF/6-31G)



Alanine residue peptides: dense vs. sparse matrix algebra (HF/6-31G)



Response theory

- The expectation value of \hat{A} in the presence of a perturbation \hat{V}_ω of frequency ω :

$$\langle t | \hat{A} | t \rangle = \langle 0 | \hat{A} | 0 \rangle + \int \langle\langle \hat{A}; \hat{V}_\omega \rangle\rangle_\omega \exp(-i\omega t) d\omega + \dots$$

- the linear-response function $\langle\langle \hat{A}; \hat{V}_\omega \rangle\rangle_\omega$ carries information about the first-order change in the expectation value

- The linear-response function may be represented compactly as:

$$\langle\langle \hat{A}; \hat{V}^\omega \rangle\rangle_\omega = -\mathbf{A}^{[1]T} \underbrace{(\mathbf{E}^{[2]} - \omega \mathbf{S}^{[2]})^{-1} \mathbf{V}_\omega^{[1]}}_{\text{linear equations}} \leftarrow \begin{cases} \mathbf{E}^{[2]} & \text{electronic Hessian} \\ \mathbf{S}^{[2]} & \text{metric matrix} \\ \mathbf{A}^{[1]} = \text{vec}(\mathbf{ADS} - \mathbf{SDA}) \end{cases}$$

- In practice, the response functions are evaluated by solving a set of linear equations

$$\begin{aligned} (\mathbf{E}^{[2]} - \omega \mathbf{S}^{[2]}) \mathbf{N}^{[1]} &= -\mathbf{V}_\omega^{[1]} \\ \langle\langle \hat{A}; \hat{V}^\omega \rangle\rangle_\omega &= \mathbf{A}^{[1]T} \mathbf{N}^{[1]} \end{aligned}$$

- can this be accomplished efficiently in the AO basis?

Solution of the response equations

- The response equations are solved in the same manner as the RH Newton equations:

$$(\mathbf{E}^{[2]} - \omega \mathbf{S}^{[2]}) \mathbf{x} = \mathbf{V}^{[1]}$$

- transformation to orthogonal basis (Cholesky or Löwdin)
- generation of an iterative subspace until the residual is sufficiently small

$$\mathbf{R} = (\mathbf{E}^{[2]} - \omega \mathbf{S}^{[2]}) \mathbf{x} - \mathbf{V}^{[1]}$$

- Key step: multiplication of Hessian and metric matrices with trial vectors

$$\mathbf{E}^{[2]}(\mathbf{X}) = \mathbf{HXS} + \mathbf{SXH} + \mathbf{g}^{\text{vo}}([\mathbf{D}, \mathbf{X}]_s) - \mathbf{g}^{\text{ov}}([\mathbf{D}, \mathbf{X}]_s)$$

$$\mathbf{S}^{[2]}(\mathbf{X}) = \mathbf{S}^{\text{oo}} \mathbf{XS}^{\text{vv}} - \mathbf{S}^{\text{vv}} \mathbf{XS}^{\text{oo}}$$

- requires recalculation of Fock/Kohn–Sham matrix with modified AO density matrix

- For rapid convergence, the residual vector is preconditioned

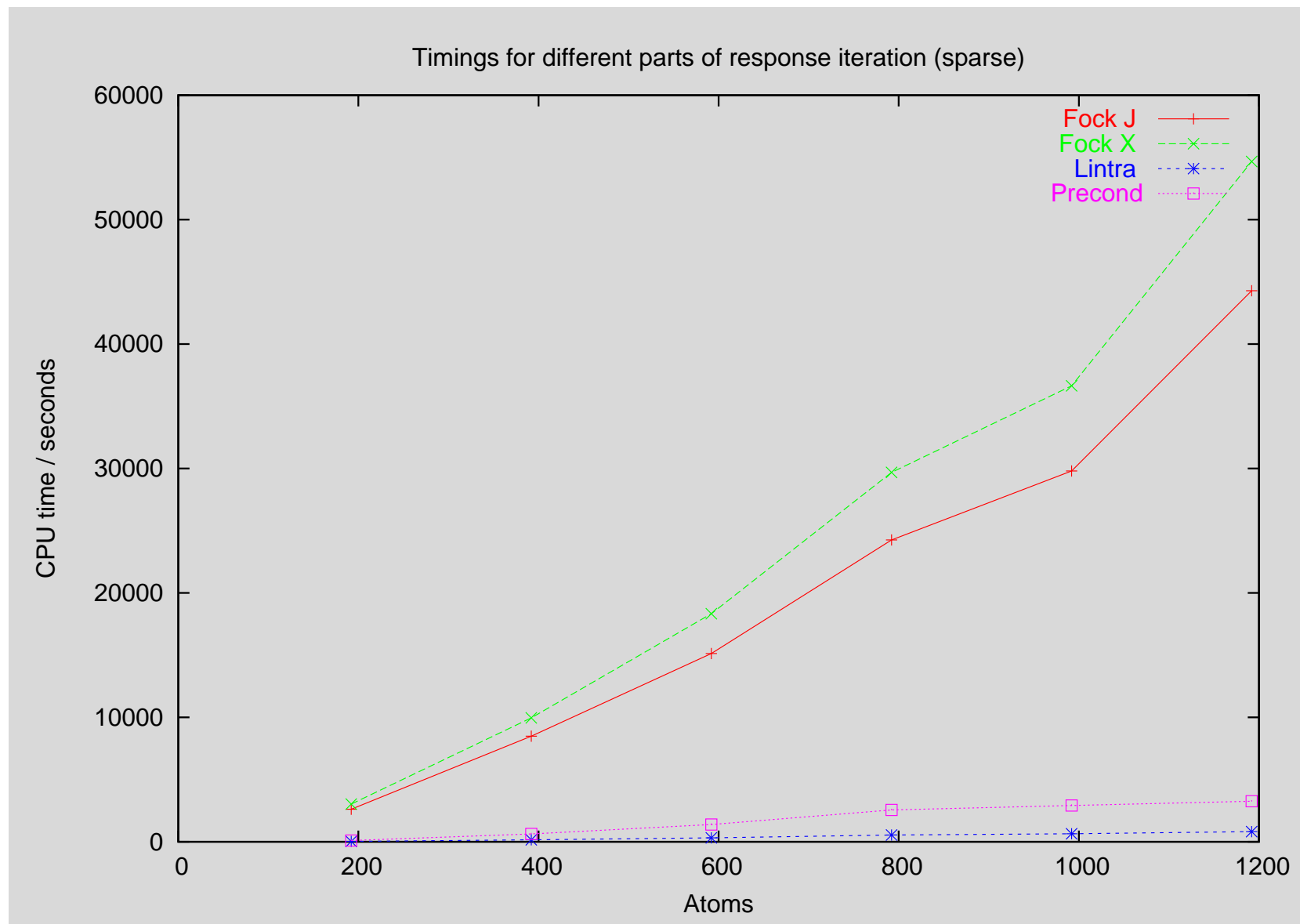
$$\tilde{\mathbf{R}} = \mathbf{M}^{-1} \mathbf{R}, \quad \mathbf{M} = \mathbf{E}^{[2]} (\text{without red part}) - \omega \mathbf{S}^{[2]}$$

- nondiagonal preconditioning requires about 5 conjugate-gradient steps

- With this preconditioner, the response equations converge in about 4 iterations

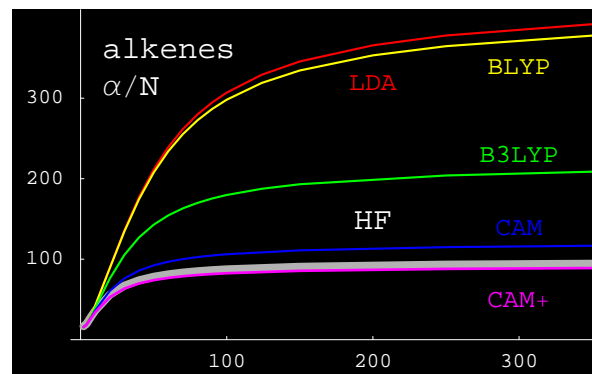
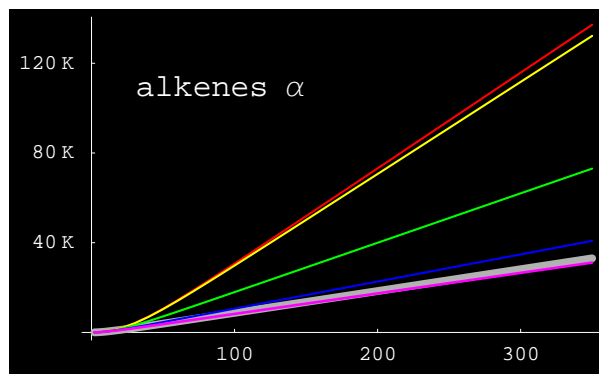
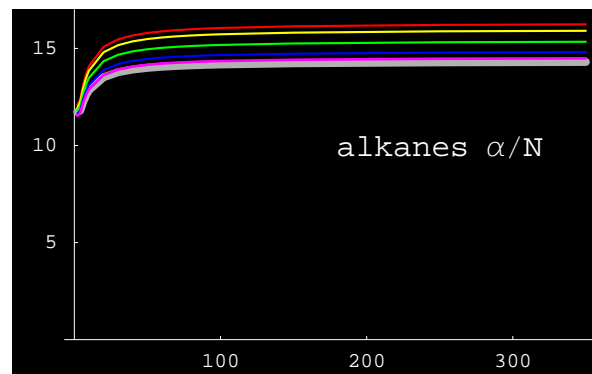
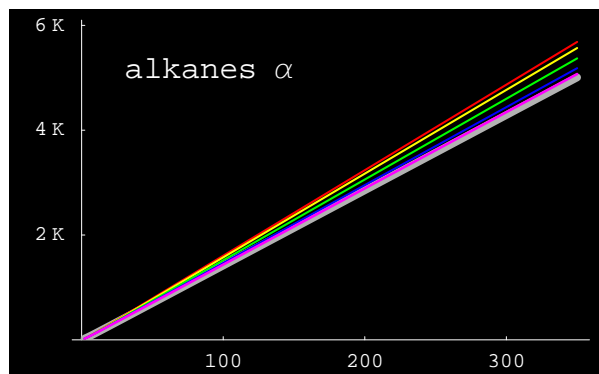
- indeed, this is the same convergence as in the canonical MO basis
- total cost: 4 Fock/Kohn–Sham evaluations, 100 matrix multiplications

Alanine residue peptides: timings in one conjugate-gradient iteration (HF/6-31G)



Polarizabilities of linear alkanes and alkenes

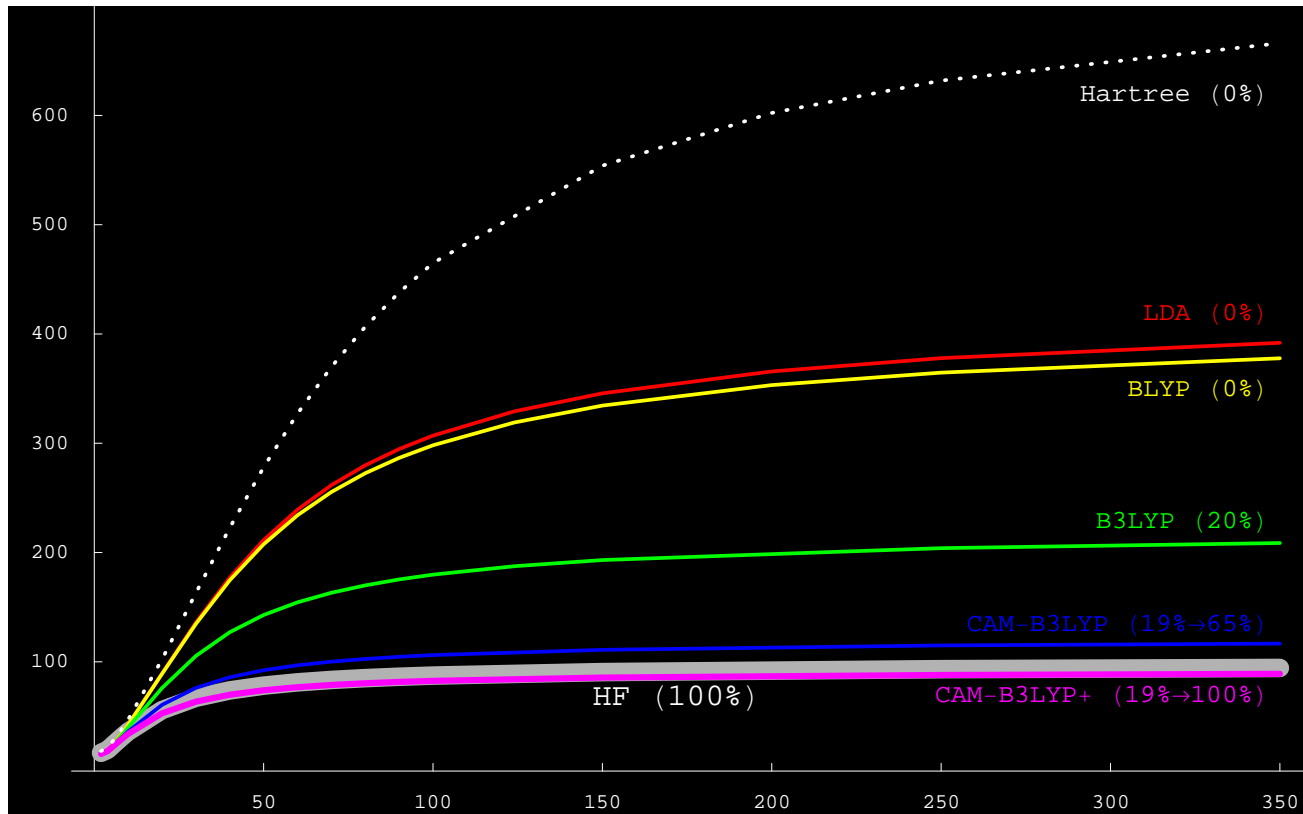
- To illustrate, we have calculated longitudinal polarizabilities in linear polymeric chains
 - HF and DFT α and α/N in 6-31G basis, plotted against the number of carbons N



- The alkenes are about an order of magnitude more polarizable than the alkanes
 - all models agree on alkanes (α/N -limit: HF 14.4; LDA 16.3)
 - widely different results for alkenes (α/N -limit: HF 97; LDA 427)

The importance of exact exchange for longitudinal polarizabilities

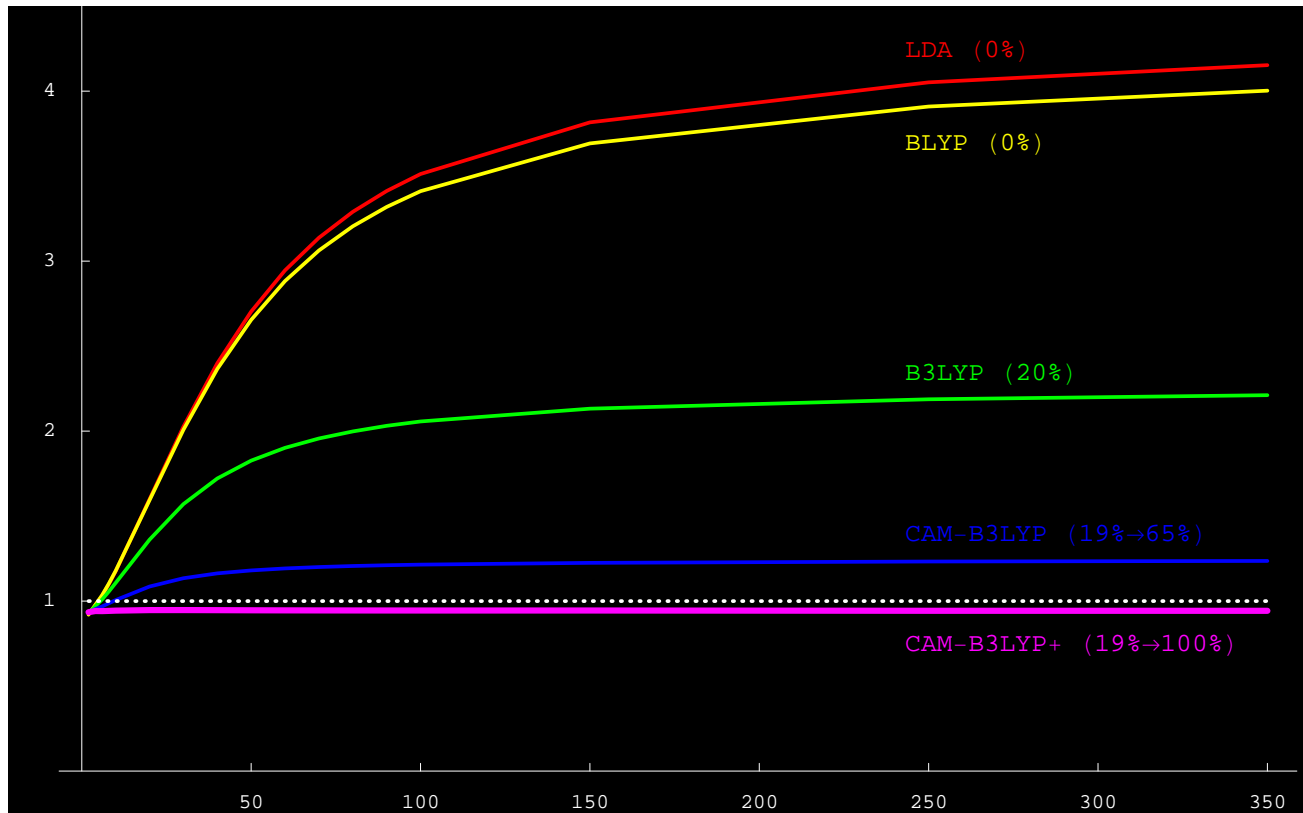
- Without a good description of long-range exchange, the systems become too polarizable



- the Hartree model neglects all exchange and overestimates by a factor of eight
- pure DFT has a poor long-range exchange and overestimates by a factor of four
- hybrid functionals improve the situation, introducing some exact exchange
- compromise solution: standard DFT at short range, full exchange at long range

The importance of exact exchange for longitudinal polarizabilities

- Without a good description of long-range exchange, the systems become too polarizable



- the Hartree model neglects all exchange and overestimates by a factor of eight
- pure DFT has a poor long-range exchange and overestimate by a factor of four
- hybrid functionals improve the situation, introducing some exact exchange
- compromise solution: standard DFT at short range, full exchange at long range

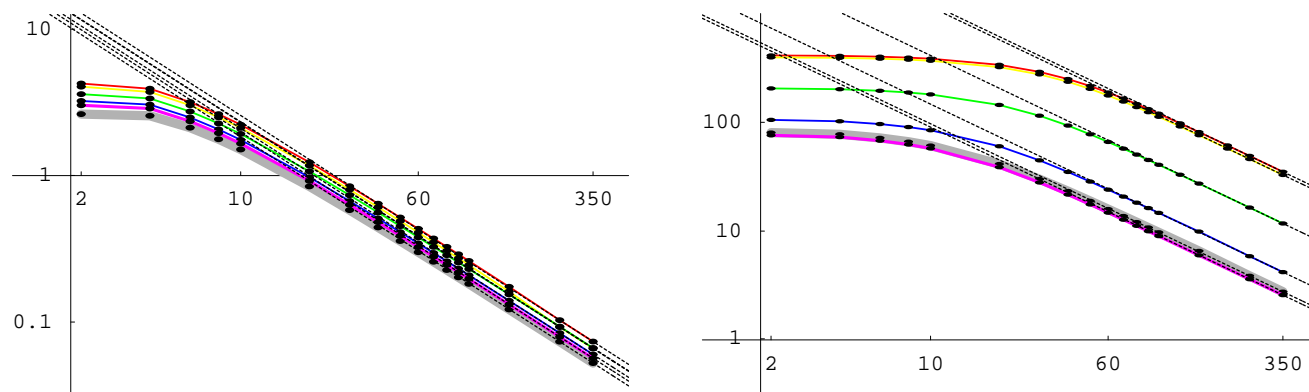
Asymptotic behaviour of group polarizabilities

- How does the group polarizability converge towards the infinite limit?

$$\bar{\alpha}_{\infty} - \bar{\alpha}_N = eN^{-1} + \mathcal{O}(N^{-2}) \quad \text{Kudin *et al.*, JCP **122**, 134907 (2005)}$$

– this behaviour is universal, holding at all levels of theory

- Log–log plots of $\bar{\alpha}_{\infty} - \bar{\alpha}_N$ for alkanes and alkenes:



– limit obtained by extrapolation $\bar{\alpha}_{\infty} = (\alpha_N - \alpha_M)/(N - M)$

– straight lines of slope -1 superimposed through the points at $N = 350$

- The asymptotic region is reached with $\text{C}_{30}\text{H}_{62}$ (alkanes) and $\text{C}_{60}\text{H}_{62}$ (alkenes)

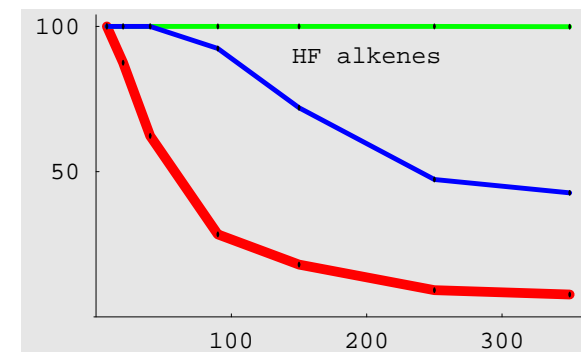
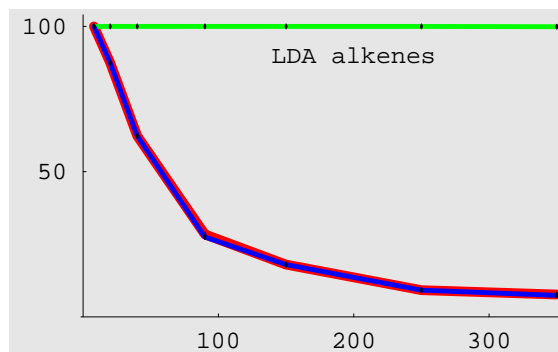
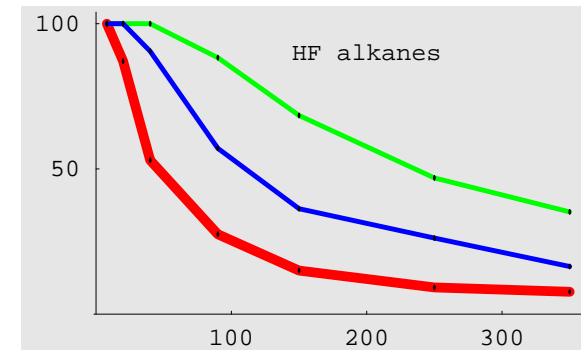
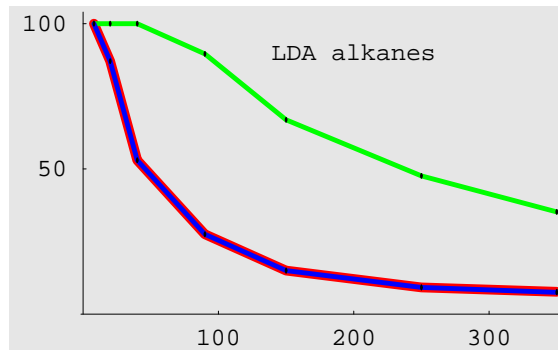
– alkane $\bar{\alpha}_{\infty}$ predicted to within 1% from $\text{C}_{30}\text{H}_{62}$

– alkene $\bar{\alpha}_{\infty}$ predicted to within 1% from $\text{C}_{60}\text{H}_{62}$ for HF and from $\text{C}_{150}\text{H}_{152}$ for LDA

Sparsity of linear alkanes and alkenes

- Each energy optimization was converged in 6–14 SCF iterations
 - about 70 matrix multiplications for each TRRH step (diagonalization)
 - about 50 matrix multiplications for each TRDMS step (DIIS)
- Each polarizability component required 3–4 response iterations
 - about 20 matrix multiplications in each iteration
- Percentage of matrix elements greater than 10^{-6} in alkane and alkene chains

- **overlap matrix:**
 - sparse
- **density matrix:**
 - nonsparse for alkenes
 - sparse for alkanes
- **Fock/KS matrices:**
 - KS matrix like overlap
 - Fock matrix intermediate between overlap and density matrices



Conclusions

- We have discussed the optimization of SCF energies without MOs
 - in each SCF iteration, we replace diagonalization by minimization
 - minimization by Newton's method, one step is usually enough
 - minimization stable and fast, highly competitive with diagonalization
 - 50–100 sparse matrix multiplications required
- Large molecules represent a more difficult minimization problem
 - small Hessian eigenvalues for pure DFT
 - trust-region SCF: careful step-size control
 - revert to second-order if necessary
- Linear-response is straightforward in the AO basis
 - one Fock/Kohn–Sham matrix build and 20 matrix multiplications pr. iteration
 - stable convergence in 3–5 iterations