Large-scale electronic structure methods

- Introduction
- O(N) Krylov subspace method
- Applications
- Numerically exact low-order scaling method
- Outlook

Taisuke Ozaki (ISSP, Univ. of Tokyo)

The Winter School on DFT: Theories and Practical Aspects, Dec. 19-23, CAS.

Towards first-principle studies for industry



DFT calculations of thousands atoms is still a grand challenge. $O(N^3)$ Low-order



10^2 atom

Many applications done. There are many successes even for material design. DNA



System size

Battery



Materials properties

Materials properties of actual materials are determined by intrinsic properties and secondary properties arising from inhomogeneous structures such as grain size, grain boundary, impurity, and precipitation.
 In use of actual materials, the materials properties can be maximized by carefully designing the crystal structure and higher order of structures .



http://ev.nissan.co.jp/LEAF/P ERFORMANCE/



e.g., the coercivity of a permanent magnet of Nd-Fe-B is determined by crystal structure, grain size, and grain boundary.





神威·太湖之光: 125 Peta flops machine

Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway

NRCPC Cores: 10,649,600 Rmax: 93,014,593.9 (GFLOP/sec.) Pmax: 125,435,904 (GFLOPS/sec.)



According to Moore's law...



How large systems can be treated by Exa machines?



The applicability of the $O(N^3)$ DFT method is extended to only 5 times larger systems.

Linear scaling methods

Two routes towards O(N) DFT



- ψ: KS orbital
- ρ: density
- φ: Wannier function
- *n*: density matrix

Density functionals as a functional of ρ

Density functionals can be rewritten by the first order reduced density matrix: ρ

$$E_{\text{tot}}[n,\rho] = \text{Tr}(\rho H_{\text{kin}}) + \int d\mathbf{r} n(\mathbf{r}) v_{\text{ext}}(\mathbf{r}) + \int \int d\mathbf{r} d\mathbf{r} d\mathbf{r}' \frac{n(\mathbf{r})n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|} + E_{\text{xc}}[n]$$

where the electron density is given by ρ

$$n(\mathbf{r}) = \sum_{i,j} \rho_{ij} \chi_j(\mathbf{r}) \chi_i(\mathbf{r})$$

Locality of Wannier functions



J.Battacharjee and U.W.Waghmare, PRB 73, 121102 (2006)

Locality of density matrix



D.R.Bowler et al., Modell.Siml.Mater.Sci.Eng.5, 199 (1997)

Various linear scaling methods

×

Wannier functions (WF) Density matrix (DM) Variational (V) Perturbative (P)

At least four kinds of linear-scaling methods can be considered as follows:

DM+P WF+VWF+P $\mathbf{DM} + \mathbf{V}$ **Krylov subspace** Hoshi Orbital Density matrix by Li and Daw Mostofi Divide-conquer minimization by Galli, Parrinello, Recursion and Ordejon Fermi operator

O(N) DFT codes

OpenMX: (Krylov) Ozaki (U. of Tokyo) et al.

Conquest: (DM) Bowler(London), Gillan(London), Miyazaki (NIMS)

Siesta: (OM) Ordejon et al.(Spain)

ONETEP: (DM) Hayne et al.(Imperial)

FEMTECK: (OM) Tsuchida (AIST)

FreeON: (DM) Challacombe et al.(Minnesota)

Basic idea behind the O(N) method



Assumption

Local electronic structure of each atom is mainly determined by neighboring atomic arrangement producing chemical environment.

Convergence by the DC method

Just solve the truncated clusters → Divide-Conquer method



For metals, a large cluster size is required for the convergence.
→ Difficult for direct application of the DC method for metals

O(N) Krylov subspace method

Two step mapping of the whole Hilbert space into subspaces



Development of Krylov subspace vectors

The Krylov vector is generated by a multiplication of H by $|K\rangle$, and the development of the Krylov subspace vectors can be understood as hopping process of electron.



The information on *environment* can be included from near sites step by step, resulting in reduction of the dimension.

Generation of Krylov subspaces

The ingredients of generation of Krylov subspaces is to multiply $|W_n\rangle$ by S⁻¹H. The other things are made only for stabilization of the calculation.

$$|R_{n+1}\rangle = S^{-1}H|W_n\rangle$$

$$|W'_{n+1}\rangle = |R_{n+1}\rangle - \sum_{m=0}^n |W_m\rangle (W_m|\hat{S}|R_{n+1})$$

$$|W_{n+1}\rangle = S - \text{orthonormalized block vector of } |W'_{n+1}\rangle$$

Furthermore, in order to assure the S-orthonormality of the Krylov subspace vectors, an orthogonal transformation is performed by $U_{\rm K} = \mathbf{W}\mathbf{X}\lambda^{-1}$ $\lambda^2 = \mathbf{X}^{\dagger}\mathbf{W}^{\dagger}\hat{S}\mathbf{W}\mathbf{X}$

For numerical stability, it is crucial to generate the Krylov subspace at the first SCF step.

Embedded cluster problem

Taking the Krylov subspace representation, the cluster eigenvalue problem is transformed to a standard eigenvalue problem as:

$$Hc_{\mu} = \varepsilon_{\mu} Sc_{\mu} \longrightarrow H^{K} b_{\mu} = \varepsilon b_{\mu}$$

where H^K consists of the short and long range contributions.



- The embedded cluster is under the Coulomb interaction from the other parts.
- The charge flow from one embedded cluster to the others is allowed.

Relation between the Krylov subspace and Green's funtion

A Krylov subspace is defined by

$$\mathbf{U}_{\mathbf{K}} = \left\{ |W_0\rangle, (S^{-1}H)|W_0\rangle, (S^{-1}H)^2|W_0\rangle, \dots, (S^{-1}H)^q|W_0\rangle \right\}$$

A set of q-th Krylov vectors contains up to information of (2q+1)th moments.

$$\begin{aligned} \underline{\mathbf{H}}_{mn}^{K} &= (W_{0}|(A^{\dagger})^{m}HA^{n}|W_{0}) \\ &= (W_{0}|S(S^{-1}H)^{m+n+1}|W_{0}), \\ &= (W_{0}|S\mu^{(m+n+1)}S|W_{0}) \end{aligned} \qquad \begin{aligned} \mathbf{Definition of moments} \\ \mu^{(p)} &= c\varepsilon^{p}c^{\dagger}, \\ &= cc^{\dagger}Hcc^{\dagger}Hc\cdots c^{\dagger}Hcc^{\dagger}, \\ &= (S^{-1}H)^{p}S^{-1} \end{aligned}$$

The moment representation of G(Z) gives us the relation.

$$G_{ij}(Z) = \sum_{p=0}^{\infty} \frac{\mu_{ij}^{(p)}}{Z^{p+1}}$$

One-to-one correspondence between the dimension of Krylov subspace and the order of moments can be found from above consideration.

Convergence property

The accuracy and efficiency can be controlled by the size of truncated cluster and dimension of Krylov subspace.



In general, the convergence property is more complicated. See PRB 74, 245101 (2006).

Comparison of computational time

The computational time of calculation for each cluster does not depend on the system size. Thus, the computational time is O(N) in principle.



Parallelization

How one can partition atoms to minimize communication and memory usage?

Requirement:

- Locality
- Same computational cost
- Applicable to any systems
- Small computational overhead



T.V.T. Duy and T. Ozaki, CPC 185, 777 (2014).

Modified recursive bisection

If the number of MPI processes is 19, then the following binary tree structure is constructed.



In the conventional recursive bisection, the bisection is made so that a same number can be assigned to each region. However, the modified version bisects with weights as shown above.

Reordering of atoms by an inertia tensor

Atoms in an interested region are reordered by projecting them onto a principal axis calculated by an inertia tensor.



The principal axis is calculated by solving an eigenvalue problem with an inertia tensor:

Allocation of atoms to processes



Diamond 16384 atoms, 19 processes



Multiply connected CNT, 16 processes



Parallel efficiency on K



The parallel efficiency is 68 % using 131,072 cores.

Applications of the O(N) method

1. Interface structure between BCC Iron and carbides

H. Sawada et al., Modelling Simul. Mater. Sci. Eng. 21, 045012 (2013).

2. Desolvation of Li⁺

T. Ohwaki et al., J. Chem. Phys. 136, 134101 (2012). T. Ohwaki et al., J. Chem. Phys. 140, 244105 (2014).

3. Electronic transport of graphene nanoribbon

M. Ohfuchi et al., Appl. Phys. Express 7, 025101 (2014). H Jippo, T Ozaki, S Okada, M Ohfuchi, J. Appl. Phys. 120, 154301 (2016).

Precipitation in bcc-Fe

In collaboration wit Dr. Sawada (Nippon Steel)

Pure iron is too soft as structural material. Precipitation of carbide can be used to control the hardness of iron.



Precipitating materials: TiC, VC, NbC



Interface and strain energies



Diameter of precipitate

Resistance force and precipitate diameter

Y. Kobayashi, J. Takahashi and K. Kawakami, Scripta Mater. 67 (2012) 854



Diameter of precipitates R (nm)

Crossover from coherent to semi-coherent



Numerically exact low-order scaling method

Main difficulty: 'diagonalization'

O(N³) method - Numerically exact diagonalization Householder+QR method Conjugate gradient (CG) method Davidson method

Even if basis functions are localized in real space, Gram-Shmidt (GS) type method is needed to satisfy orthonormality among eigenstates, which results in $O(N^3)$ for the computational time.

O(N) method - can be achieved in exchange for accuracy. O(N) Krylov subspace method, DC, DM, OM methods, etc..

 $O(N^{2})$ method Is it possible to develop $O(N^{2})$ methods without introducing approximations? \rightarrow No more GS process.

Possible ways to avoid orthogonalization

Numerically exact low-order scaling method

- Numerically exact
- Applicable to insulators and metals
- Suitable for parallel computation
- ✓ Applicable to 1D, 2D, 3D systems
- Applicable to any local basis functions

TO, PRB 82, 075131 (2010)

Numerically exact low-order scaling method

1. Direct evaluation of the selected elements of ρ via a contour integration of the Green's function

$$\rho = M^{(0)} + \operatorname{Im}\left(-\frac{4i}{\beta}\sum_{p=1}^{\infty}G(\alpha_p)R_p\right)$$

2. Nested dissection of sparse matrix

TO, PRB 82, 075131 (2010)

Continued fraction rep. of Fermi function

$$\frac{1}{1 + \exp(x)} = \frac{1}{2} - \frac{x}{4} \left(\frac{1}{\left(\frac{x}{2}\right)^2} + \frac{\left(\frac{x}{2}\right)^2}{\left(\frac{x}{2}\right)^2} + \frac{\left(\frac{x}{2}\right)^2}{\left(\frac{x}{2}\right)^2} + \frac{\left(\frac{x}{2}\right)^2}{5 + \frac{\left(\frac{x}{2}\right)^2}{(2M - 1) + \cdots}} \right)$$

TO, PBR 75, 035123 (2007)

Contour integration

Cn Ζn 10000 (n-1) $\frac{Z(n-1)}{\beta}$ 5000 Imaginary **C**2 **Z**1 -5000 ß Z0 ß -10000 <u>-1</u>-1 0 Real -R+µ +R+µ μ 10⁵ Interval between Neighboring Poles 10⁴ **Continued fraction** 10³ Matsubara 10² 10¹ 10⁰ 20 80 100 40 60 Index of Poles

All the poles are located on the imaginary axis.

The form has a special pole structure, that is, the interval between neighboring poles increases in a faraway region from the real axis, which is very advantageous for the contour integration of Green's function.

Convergence of p w.r.t. poles

The calculation of ρ can be expressed by a contour integration:

$$\begin{aligned} \rho_{ij} &= \sum_{k} f(\frac{\varepsilon_{k} - \mu}{k_{\mathrm{B}}T}) \langle \chi_{i} | \phi_{k} \rangle \langle \phi_{k} | \chi_{j} \rangle, \\ &= -\frac{2}{\pi} \mathrm{Im} \int_{-\infty}^{\infty} dE f(\frac{E - \mu}{k_{\mathrm{B}}T}) G_{ij}(E + i0^{+}), \\ &= M_{ij}^{(0)} + \mathrm{Im} \left[-\frac{4i}{\beta} \sum_{p=1}^{\infty} G_{ij}(\alpha_{p}) R_{p} \right], \quad \stackrel{M_{ij}^{(0)} = \mathrm{Im} \left[-\frac{1}{\pi} \int_{-\infty}^{\infty} dE G_{ij}(E + i0^{+}) \right] \simeq iR G(iR)}{\alpha_{p} = \mu_{0} + i \frac{z_{p}}{\beta}} \\ &= M_{ij}^{(0)} + \mathrm{Im} \left[-\frac{4i}{\beta} \sum_{p=1}^{\infty} \sum_{k} \frac{\langle \chi_{i} | \phi_{k} \rangle \langle \phi_{k} | \chi_{j} \rangle}{\alpha_{p} - \varepsilon_{k}} R_{p} \right], \quad \text{Lehmann rep.} \\ &= M_{ij}^{(0)} + \sum_{k} \mathrm{Im} \left[-\frac{4i}{\beta} \sum_{p=1}^{\infty} \frac{\langle \chi_{i} | \phi_{k} \rangle \langle \phi_{k} | \chi_{j} \rangle}{\alpha_{p} - \varepsilon_{k}} R_{p} \right], \end{aligned}$$

The analysis shows that the number of poles for each eigenstate for a sufficient convergence does not depend on the size of system if the spectrum radius does not change. \rightarrow The scaling property is governed by the calculation of G.

Convergence property of the contour integration

Total energy of aluminum as a function of the number of poles by a recursion method at 600 K.

Nicholson et al., PRB **50**, 14686 (1994).

Poles	Proposed	$\frac{1}{1 + \left(1 + \frac{x}{n}\right)^n}$	Matsubara
10 20 40 60 80 150 200 250	-42.933903047211 -47.224346653790 -48.323790725570 -48.324441992259 -48.324441994952 -48.324441994952	$1+(1+\frac{x}{n})^n$ -33.734015919550 -33.623477214678 -33.346245616679 -33.143128624551 -32.870752577236 -33.837428496424 -33.418012271726 34.003411636691	-39.612354360046 -39.849746603905 -40.216055898502 -39.676965494522 -43.523770052176 -41.836938942518 -42.543354202255 43.024756221080
230 300 350 400 600 1000 2000 5000 10000	The energy completely converges using only 80 poles within double precision.	-34.003236479262 -48.324440028792 -48.324440274509 -48.324440847749 -48.324440847749 -48.324441306517 -48.324441650693 -48.324441857239 -48.324441926094	-43.466729654170 -43.834528739677 -44.185100655185 -45.233651519749 -46.331692884149 -47.202779497545 -47.921384128418 -48.122496320516

How can Green's function be evaluated?

• The Green's function is the inverse of a sparse matrix (*ZS-H*).

$$G(Z) = (ZS - H)^{-1}$$

• Selected elements of G(Z), which correspond to non-zero elements of the overlap matrix S, are needed to calculate physical properties.

- Our idea
 1. Nested dissection of (*ZS-H*)
 2. LDL^T decomposition for the structured matrix

\rightarrow a set of recurrence relations

TO, PRB **82**, 075131 (2010)

$$b$$

$$a b$$

$$b a b$$

The processes (i)-(v) are recursively applied to each domains with computational cost of $O(N(\log_2 N)^2)$ in total.

(i) Ordering:

The basis functions are ordered by coordinates along each direction.

(ii) Screening:

The basis functions with a long tail are assigned as part of the separator.

(iii) Finding of a starting nucleus:

Find a basis function having the smallest number of nonzero overlaps.

(iv) Growth of the nucleus:

Minimize $|N_0-N_1| + N_s$ by the growth of the nucleus.

 N_0 : # of bases in domain 0 N_1 : # of bases in domain 1 N_s : # of bases in separator

(v) Dissection:

Find a direction with the smallest $|N_0-N_1| + N_{s,}$, make the dissection along the direction.

Square lattice for the nested dissection

Inverse by LDL^T block factorization

A matrix X can be factorized using a Schur complement into a LDL^{T} form.

$$X = \begin{pmatrix} A & B^{T} \\ B & C \end{pmatrix} = \begin{pmatrix} I \\ L & I \end{pmatrix} \begin{pmatrix} A \\ & S \end{pmatrix} \begin{pmatrix} I & L^{T} \\ & I \end{pmatrix}$$
$$L = BA^{-1}$$
$$S = C - BA^{-1}B^{T}$$

Then, the inverse of *X* is given by

$$X^{-1} = \begin{pmatrix} A^{-1} + L^T S^{-1} L & -L^T S^{-1} \\ -S^{-1} L & S^{-1} \end{pmatrix}$$

Analysis of the computational cost

Timing result

SCF convergence

Parallel efficiency

Outlook

The locality of density matrix and basis function is a key to develop a wide variety of efficient electronic structure methods.

We have demonstrated three methods:

- O(N) Krylov subspace method
- Low-order scaling exact method
- O(N) *exact* exchange method

Plenty of developments of new efficient methods might be still possible.