

# MPI+OpenMP Paralelization of Density Functional Theory Method in GAMESS

Vladimir Mironov<sup>1</sup>, Yuri Alexeev<sup>2</sup>, Dmitri G. Fedorov<sup>3</sup>

<sup>1</sup>Department of Chemistry, Lomonosov Moscow State University; <sup>2</sup>Computational Science Division, Argonne National Laboratory; <sup>3</sup>National Institute of Advanced Industrial Science and Technology

## Abstract

In this work, Density Functional Theory (DFT) method is parallelized with MPI-OpenMP in quantum chemistry package GAMESS. It has been implemented in both regular and Fragment Molecular Orbital (FMO) based DFT codes. The scalability of the FMO-DFT code was demonstrated on Cray XC 40 Theta supercomputer. We demonstrated excellent scalability of the code up to 2,048 Intel Xeon Phi nodes (131,072 cores). Moreover, the developed DFT code is about twice faster than the original code because of our new grid integration algorithm.

## GAMESS (General Atomic and Molecular Electronic Structure System)

- A general purpose comprehensive *ab initio* quantum chemistry package
- Maintained by the research group of Prof. Mark Gordon at Iowa State University (<http://www.msg.ameslab.gov/games>).
- ✓ Enables most major quantum mechanical methods (Hartree-Fock, Møller-Plesset perturbation theory, coupled cluster, multiconfiguration self consistent field, configuration interaction, density functional theory).
- ✓ Ported to many computer architectures.
- ✓ Free and widely used on everything from laptops to supercomputers.
- ✓ Highly scalable, including many distributed data algorithms.
- ✗ More than a million lines of legacy Fortran 77 code, with an associated parallelization library comprising 15,000 lines.
- ✗ Parallelization library does not support multithreading

## Fragment Molecular Orbital (FMO) Method

□ In FMO, a molecular system is divided into fragments and the total properties, such as the energy or its gradient, is calculated from those fragments and (in FMO2) their pairs, computed in the embedding potential. FMO treats the electrostatics for the whole system while the exchange-repulsion and charge transfer are only accounted for at the local level of fragment pairs.

□ The total energy can be written as:

$$E = \sum_I^{N_{frag}} E_I + \sum_{I < J}^{N_{frag}} \sum_{J < I}^{N_{frag}} (E_{IJ} - E_I - E_J)$$

where the monomer (I) [FMO1], dimer (IJ) [FMO2] energies can be obtained using any standard QM method

- Current applications:
  - geometry optimization (~2,000 atoms)
  - *ab initio* level single point energies (~20,000 atoms)
  - pair interaction analysis (~20,000 atoms)
  - drug design, ligand docking, chemical reactions in solution etc
- Massively parallel FMO utilizes the Generalized Distributed Data Interface (GDDI):
  - With the molecule divided into fragments, each fragment or pair is computed on a distinct subset, or group, of the available processors
  - Each fragment/pair is then run in parallel in each group
  - This provides two levels of parallelization, greatly increasing the parallel efficiency

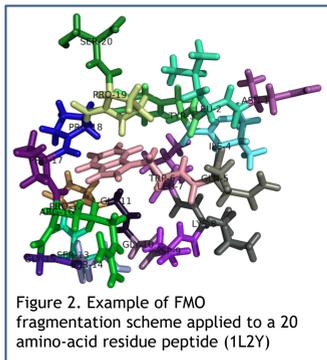


Figure 2. Example of FMO fragmentation scheme applied to a 20 amino-acid residue peptide (1L2Y)

## Density functional theory (DFT) method

- One of the most popular quantum chemistry methods
- High precision results for relatively low computational cost
- Very similar to the Hartree-Fock method

### Legacy DFT implementation in GAMESS:

1. Inefficient algorithms of numerical grid calculation, which are slow and result in too large grid
2. In-house algorithms for Fock-matrix update
3. For the nuclear gradients the unnecessary grid weight contribution is always computed
4. C-style access to multidimensional array with grid data

### DFT algorithm modernization strategy:

- Introduce modern DFT algorithms
- Use optimized BLAS subroutines if possible
- Optimize access to the data

```

1: Set up numerical grid (r_i, ω_i)
2: repeat                                     ▶ Begin SCF loop
3:   Compute Coulomb and optionally exact exchange terms
4:   for each (r_i, ω_i) do
5:     Compute electron density ρ(r_i)
6:     (Optional, GGA) compute density gradient ∇ρ(r_i)
7:     (Optional, mGGA) compute kinetic energy density τ(r_i)
8:     Compute f_i^XC(ρ, ∇ρ, τ)
9:     Update the Fock matrix using
       F_mn+ = ω_i f_i^XC φ_m(r_i) φ_n(r_i) ▶ rank-2 update
10:   end for
11:   Diagonalize the Fock matrix, update ρ and ε
12: until |D - D^old| < conv or |ε - ε^old| < conv
    
```

Figure 1. Generic algorithm of DFT method

## Goals of the study

- Develop multithreaded DFT code in GAMESS to reduce following drawbacks of the MPI implementation:
  - Excessive memory utilization on massively parallel CPU architectures for legacy codes like GAMESS.
  - Only half of GAMESS processors are doing computation, other are used as data-servers.
- Bring recent advances of the DFT method to GAMESS users
- Create a framework for DFT method development in GAMESS

## FMO algorithm in GAMESS

### Basic steps of FMO (FMO2) calculation:

1. Compute one-electron (core) Hamiltonian of each fragment
2. Construct initial guess wavefunction of each I fragment
3. Monomer self consistent charge (SCC) field calculations. Compute the wavefunction of fragments in a field of other fragments. The latter is described by the embedding electrostatic potential (ESP). This step loops until all wavefunctions converge
4. Dimer calculations. They are done once at the end
5. Aggregate the results and compute properties

### Problem of FMO-DFT calculation:

- Computing XC contribution to the Fock matrix is a major bottleneck for small chemical systems, like FMO fragments.

## Hardware description

	Theta (ALCF)	Intel Xeon single node
Processor	Intel Xeon Phi 7230, 64 cores, 1.3 GHz	2x Intel Xeon Gold 6254
Memory per node	16 GB MCDRAM + 192 GB DDR4	1.5 TB DDR4
Number of nodes	3,624	1
Interconnect type	Aries interconnect with Dragonfly topology	-
Compiler	Intel Parallel Studio XE 2018.2	Intel Parallel Studio XE 2018.2

```

1: Calculate one-electron Hamiltonian of each fragment
2: Guess initial density matrix for each fragment D_i
   ▶ Monomer calculation:
3: repeat
4:   for i = 1, NumFrag do
5:     GDDI load balancing
6:     Compute embedding ESP quasi-Fock matrix
7:     Compute wavefunction of the fragment i, update
       energy and density matrix of fragment i (ε_i, D_i)
8:     Compute fragment properties
9:     Broadcast results (energy, density, charges)
10:   end for
11: until max(|D_i^old - D_i|) < conv
   ▶ Dimer calculation:
12: for i = 1, NumFrag do
13:   for j = 1, i - 1 do
14:     GDDI load balancing
15:     Compute embedding ESP of the dimer
16:     Compute wavefunction, get energy (ε_ij) and
       density matrix(D_ij) of the fragment pair ij
17:   end for
18: end for
19: Gather results of fragment pairs calculation
20: Calculate whole system properties, print results
    
```

Figure 3. Algorithm of FMO (FMO2) method

## New DFT implementation

- Store grid as a list of clusters of points:
  - Discard clusters of points whose contribution to the Fock matrix is less than a threshold
  - Clustering allows to use BLAS level 3 operations when updating Fock matrix
- Whenever possible skip computing and storing Becke's factors to quadrature weights
- MPI+OpenMP parallelization both on grid construction and XC integration steps
- Static/dynamic load balancing for MPI processes and dynamic load balancing for OpenMP threads

```

1: Zero XC contribution to Fock matrix: F(:, :) = 0
2: !$omp parallel do schedule(dynamic) reduction(F)
3: for m = 1, N_clusters, MPI_ChunkSize do
4:   MPI static/dynamic load balancing
5:   Get data of all points in Q_m
6:   for each point p in Q_m do
7:     Compute AO densities ρ_i^AO
8:     (Opt.) Compute ∇ρ_i^AO
9:     (Opt.) Compute ∇(∇ρ_i^AO)
10:    if all ρ_i^AO < p_thr then
11:      Mark point p for removal
12:    end if
13:  end for
14:  if cluster Q_m is empty then
15:    Mark cluster m for removal
16:  end if
17:  Compute ρ_i^MO for remaining points using BLAS
18:  (Opt.) Compute ∇ρ_i^MO
19:  Compute XC values on a grid
20:  Update Fock matrix using BLAS
21: end for
22: if grid was changed then
23:   Synchronize grid among MPI processes
24: end if
    
```

Figure 4. Integration of XC functional

## Benchmark description

- DFT setup:**
- Phenylalanine aminoacid:
    - Moderate-size aminoacid
  - B3LYP/6-31G\* energy + gradient calculation
    - Precise DFT grid (96 radial × 1202 angular grid points per atom)
  - Run on single node to compare the performance of both codes

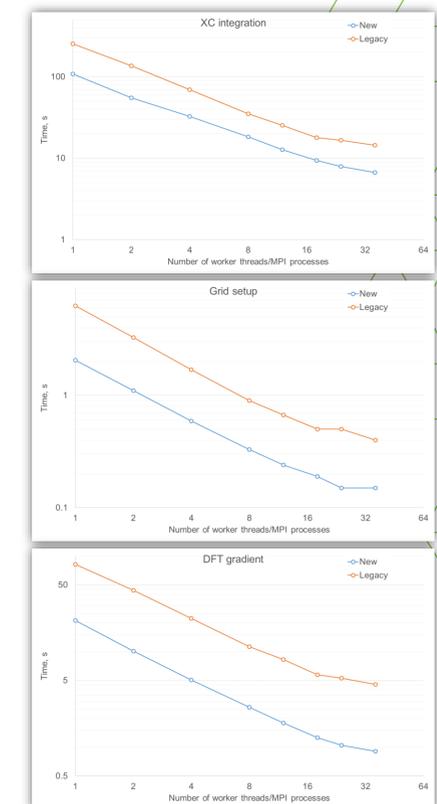
- FMO-DFT setup:**
- Spherical water droplet with 30 Å radius
  - FMO2-B3LYP/6-31G\* energy calculation
  - Large-scale benchmarking on Theta supercomputer

## Conclusions

In this work, Density Functional Theory (DFT) method has been parallelized with MPI-OpenMP in quantum chemistry package GAMESS. The parallel performance has been measured on a single Intel Xeon Phi node and supercomputer Theta (Cray XC40 based on Intel Xeon Phi). The developed hybrid MPI+OpenMP code scales well up to 2,048 Intel Xeon Phi nodes. Moreover, we also rewrote DFT grid integration code, resulting in about 2x performance improvement compared to the original code.

## Benchmark results

### Single node benchmarks



```

1: Set up partitioning function σ(μ)
2: Set up grid point clusters:
3: for each atom do
4:   Compute spherical grid S for an atom
5:   Compute radial grid R for an atom, split on intervals R_i
   ▶ Atomic "sphere of influence":
6:   Compute r_min: ∀ r < r_min, σ(μ(r)) = 1
7:   ▶ Construct clusters of points for an atom:
8:   for each radial interval R_i do
9:     Divide points of spherical grid to spherical sectors S_j
10:    Store clusters of grid points Q_ij = R_i × S_j in memory
11:  end for
12: end for
13: ▶ Process grid point clusters:
14: !$omp parallel do schedule(dynamic)
15: for m = 1, N_clusters, MPI_ChunkSize do
16:   MPI static load balancing
17:   if r_max(Q_m) < r_min then
18:     Mark cluster as non-modified
19:   else
20:     for each spherical direction sph in cluster do
21:       for each radii value rad in cluster do
22:         ▶ Process point p(rad, sph):
23:         Compute w_p^b = w_rad w_sph
24:         if p in sphere of influence of another atom then
25:           w_p = 0
26:         exit radii loop
27:       end if
28:       Compute Becke's factors f(σ) for point p
29:       and all other atoms
30:       Compute w_p from w_0 and f(σ)
31:       if w_p = 0 then
32:         exit radii loop
33:       end if
34:     end for
35:   end for
36:   if all w_p = 0 then
37:     Drop cluster Q_m
38:   end if
39: end for
40: Synchronize grid among MPI processes
    
```

Figure 5. Molecular grid construction

### FMO-DFT benchmarks on Theta

