European Summer School of Quantum Chemistry 2024 Torre Normanna Sicily

Lectures 1 & 2: Algorithm Design

Frank Neese



MAX-PLANCK-GESELLSCHAFT MPI für Kohlenforschung Kaiser-Wilhelm Platz 1 45470 Mülheim an der Ruhr

You have an Equation, what now?

Let us assume that you have derived an equation, e.g. using Wick's theorem



Or diagrams ...



i,j,k,l=occupied a,b,c,d=virtual

$$\sigma_{ab}^{ij} \leftarrow P(ij)P(ab) \sum_{kl} \sum_{cd} \langle kl || cd \rangle t_{ac}^{ik} t_{db}^{lj}$$

You are charged (or simply want) to implement that. What do you do?

Goals of Quantum Chemical Method Development

 \checkmark To be able to calculate ,things' (energies, properties) that could not be calculated

before or on systems that were not accessible before

✓ To develop a better (more accurate, more elegant, more compact, more transparent,

...) theory for a known property.

Develop new approximations to known equations

√ ...

✓ To obtain the same number faster than before

To obtain an approximate number faster (and in ,improved scaling') than before

Prelude

A little information on Computers

Computer Architecture and Algorithm Design



Main Memory: Stack vs Heap

Stack: Fixed Size Reserved Main Memory region of static variables

Heap: Variable Size Main Memory region of <u>dynamic</u> variables



Main memory

Heap

NOTE

- memory allocated on the Stack will be released automatically when the variable goes out of scope.
- Memory allocated on the Heap you need to de-allocate yourself or you are created a memory leak (very common mistake)

The Stack can Overflow

✓ The Stack is usually small (e.g. default on a Mac is 8MB)

- If you are putting a lot of data on the stack if may overflow
- Core dump / Segmentation fault

```
void SignOfLife()
{
    double x[1024];
    printf("Sign of Life - still hanging in there!\n");
    fflush(stdout);
};
frankneese@MacBook-Pro-von-Frank 2024 - ESQC-24 % ./stack
Sign of Life - still hanging in there!
void SignOfLife()
{
    double x[1024*1024];
    ...
```

frankneese@MacBook-Pro-von-Frank 2024 - ESQC-24 % ./stack
zsh: segmentation fault ./stack

✓ The Heap can overflow too of course. In this case there might be chance to recover using "exception handlers" (e.g. C++ try ... catch)

Side Remark: Call by Reference vs Call by Value

An argument **passed by value:**

- creates the copy of the variable on the stack
- unchanged after the function call is done
- Potentially creates overhead from copying or stack overflows

void foo(int i) { i=i+1; }; int i=1; foo(i); printf("i=%d"); The printout will be i=1

An argument **passed by reference:**

- What is passed is a pointer to the storage location of the variable (Heap or Stack)
- Potentially changed after the function call is done
- Creates no overhead

```
void foo(int &i) { i=i+1; };
int i=1; foo(i); printf("i=%d"); The printout will be i=2
```

Summary: Hardware components

Hard-disk:	Potentially large (assumed infinite in the "von Neumann machine" or the abstract "Turing machine" Very slow compared to the rest. Avoid as much as you can!
Main memory:	Small Stack (static), large Heap (dynamic) Relatively slow. It takes time to dig data out from there.
Bus:	Transfer from main memory to CPU/Cache Often rate limiting in actual calculations
Cache:	Fast memory directly attached to the CPU Relatively small. You want all the data the is being processed there
Register:	Central part of the CPU that holds the data or instructions that are being processed next
CPU:	Runs instructions in sequence (e.g. GHz means nsec for individual instructions) - make sure it does not run idle!

PART 1

Before we talk about how to compute things efficiently:

"The greatest performance gains are coming from the calculations that you <u>don't</u> do at all"

3 Ways to Avoid Unproductive Computation

1. Use of symmetry Integrals have selection rules, e.g. in $(\mu\nu|\kappa\tau)$

The direct product $\Gamma(\mu) \otimes \Gamma(\nu) \otimes \Gamma(\kappa) \otimes \Gamma(\tau)$

Must contain the **totally symmetric irrep**, provided μ, ν, κ, τ

Are adapted to the irreps of the point group

2. Use of permutation symmetry

Integrals have permutation symmetry that *usually* should be used $(\mu\nu|\kappa\tau) = (\mu\nu|\tau\kappa) = (\nu\mu|\kappa\tau) = (\nu\mu|\tau\kappa)$ $= (\kappa\tau|\mu\nu) = (\tau\kappa|\mu\nu) = (\kappa\tau|\nu\mu) = (\tau\kappa|\nu\mu)$

3. Avoid terms that are (near) zero or factors that are (near) one

... here is where the art & science of thresholding starts!

Avoid small numbers but make sure errors don't add up!

Self Consistent Field

$$\mathbf{F}(\mathbf{c})\mathbf{c}_i = \varepsilon_i \, \mathbf{S} \mathbf{c}_i$$

$$S_{\mu\nu} = \left\langle \mu \,|\, \nu \right\rangle$$





"Mean Field" Hartree-Fock

$$F_{\mu\nu} = h_{\mu\nu} + \sum_{\kappa\tau} P_{\kappa\tau}(\mu\nu \mid \kappa\tau) - \frac{c_x}{2} \sum_{\kappa\tau} P_{\kappa\tau}(\mu\kappa \mid \nu\tau) + \int \mu(\mathbf{r}) \frac{\delta E_{XC}[\rho]}{\underbrace{\delta\rho(\mathbf{r})}_{V_{XC}(\mathbf{r})}} \nu(\mathbf{r}) d\mathbf{r}$$

One Electron

Coulomb

HF Exchange

XC Potential

Almlöf's Revolutionary Proposals

For decades progress in quantum chemistry was prevented by the O(N⁴) of twoelectron integrals.



Even if the integrals can be stored for a 1000 basis function calculation, the I/O penalty is huge and the CPU remains largely idle while waiting for data to arrive from the hard drive

The integral bottleneck was finally overcome by Almlöfs revolutionary proposals

- 1. Do <u>NOT</u> store integrals. Recalculate when needed (direct SCF)
- 2. Split the calculation of the Coulomb and exchange terms and use the most efficient approximation for each rather than use the same integrals for both.

Let's take a look at Electron Repulsion Integrals

Look at an ERI:
$$(\mu\nu|\kappa\tau) = \int \int \frac{\mu(\mathbf{r}_1)\nu(\mathbf{r}_1)\kappa(\mathbf{r}_2)\tau(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2$$

This can be viewed as the electrostatic interaction of two smeared out charge distributions: $a_{1}(r_{1}) = \mu(r_{1})\nu(r_{2})$

$$\rho_{\mu\nu}(\boldsymbol{r}_1) = \mu(\boldsymbol{r}_1)\nu(\boldsymbol{r}_1)$$
$$\rho_{\kappa\tau}(\boldsymbol{r}_2) = \kappa(\boldsymbol{r}_2)\tau(\boldsymbol{r}_2)$$

And it is advantageous to take the basis functions themselves as Gaussians:

$$\mu_A(\mathbf{r}) = S_{l_\mu m_\mu}(\mathbf{r} - \mathbf{R}_A) N_\mu \sum_k d_{k\mu} \exp\left(-\alpha_k r_A^2\right)$$

Negligible Integrals: Gaussian Product Theorem



- In a large system there are only O(N) ,significant' Gaussian products.
- → The should be precomputed and stored as a list (e.g. cut-off K_{AB} >= T_{cut})
- The significant bra- and ket-products interact via the 1/r operator (never small!).
- There are O(N²) non-negligible integrals

The principle of "Direct SCF"

$$\begin{split} F_{\mu\nu} &= h_{\mu\nu} + G_{\mu\nu} \quad G_{\mu\nu} = \sum_{\kappa\tau} P_{\kappa\tau} \left[\left(\mu\nu \mid \kappa\tau \right) - \left(\mu\kappa \mid \nu\tau \right) \right] \\ \\ \begin{array}{c} \texttt{G=0} \\ \texttt{loop } \mu \\ \texttt{loop } \kappa \\ \texttt{loop } \mathsf{t} \leq \texttt{K} \ (\mu\nu \leq \mathsf{K}\mathsf{T}) \\ \texttt{test} = \boxed{\texttt{IntegralEstimate} \left(\mu, \nu, \mathsf{K}, \mathsf{T} \right) \\ \texttt{test} = \underbrace{\texttt{IntegralEstimate} \left(\mu, \nu, \mathsf{K}, \mathsf{T} \right) \\ \texttt{if } (\texttt{test} \leq \underbrace{\texttt{Thresh}}) \ \texttt{skip} \\ \texttt{else} \\ \texttt{Calculate} \ (\mu\nu \mid \mathsf{K}\mathsf{T}) \end{split}$$

end (else)

add $G(\mu,\nu) += P(\kappa,\tau) (\mu\nu | \kappa\tau)$ (Coulomb)

(and permutations of indices)

add $G(\mu,\kappa) \rightarrow = P(\nu,\tau) (\mu\nu | \kappa\tau)$ (Exchange)

end loops $\mu,\nu,\kappa,$

Only contributions >= Thresh go into the Fock matrix

Better than testing for small integrals alone since P can be large

Kohn's Conjecture and the Density Matrix



➡ The decay is exponential, but slow. 10⁻¹⁰ is only reached at 20-25 Angström!

Nevertheless, in insulators, there are only O(N) significant density matrix elements

Intrinsic Scaling of Coulomb and Exchange

Assuming exponential decay of the density, Almlöf realized that the intrinsic scaling of the Coulomb and exchange terms is different:

Coulomb:





Conclusion: Use the most efficient way to calculate or approximation each term separately!

Integral Estimates (I): Almlöf's estimate

In order to not decide that we do not calculate an integral, we need an estimate for it

Look at an ERI:
$$(\mu\nu|\kappa\tau) = \int \int \frac{\mu(r_1)\nu(r_1)\kappa(r_2)\tau(r_2)}{|r_1 - r_2|} dr_1 dr_2$$

Let us preted for a moment that r_{12}^{-1} is not there. Then:

$$\int \int \mu(\boldsymbol{r}_1) \nu(\boldsymbol{r}_1) \kappa(\boldsymbol{r}_2) \tau(\boldsymbol{r}_2) d\boldsymbol{r}_1 d\boldsymbol{r}_2 = \int \mu(\boldsymbol{r}_1) \nu(\boldsymbol{r}_1) d\boldsymbol{r}_1 \int \kappa(\boldsymbol{r}_2) \tau(\boldsymbol{r}_2) d\boldsymbol{r}_2 = S_{\mu\nu} S_{\kappa\tau}$$

Now assume that the bra- and ket distributions are centered at

$$\boldsymbol{R}_{\mu\nu} = \langle \mu | \boldsymbol{r} | \nu \rangle \quad \boldsymbol{R}_{\kappa\tau} = \langle \kappa | \boldsymbol{r} | \tau \rangle$$

Now set
$$r_{12}^{-1} \approx |\mathbf{R}_{\mu\nu} - \mathbf{R}_{\kappa\tau}|^{-1} \equiv R^{-1}$$

And arrive at

$$(\mu\nu|\kappa\tau) \approx \frac{S_{\mu\nu}S_{\kappa\tau}}{R}$$

Estimate
$$(\mu\nu|\kappa\tau)_{est-I} = max\left(\left|\frac{S_{\mu\nu}S_{\kappa\tau}}{R}\right|\right)$$

(Over the members of the 4 shells)

- ➡ NOT a rigorous upper bound
- Reasonably cheap to compute
- Does take the R dependence into account to some extent

Integral Estimates (II): Ahlrich's estimate

Häser and Ahlrichs used the Schwartz inequality to show:

$$(\mu\nu|\kappa\tau) \leq \sqrt{(\mu\nu|\mu\nu)} \sqrt{(\kappa\tau|\kappa\tau)} = Q_{\mu\nu}Q_{\kappa\tau}$$

Estimate
$$(\mu\nu|\kappa\tau)_{est-II} = max(Q_{\mu\nu}Q_{\kappa\tau})$$

(Over the members of the 4 shells)

- Is a rigorous upper bound
- Is cheap to compute
- Does NOT depend on R and hence will strongly overestimate integrals with well separated bra and ket distributions

Integral Estimates (III): Multipole Estimate

Let us take two expansion points $R_{\mu\nu} = \langle \mu | r | \nu \rangle$ $R_{\kappa\tau} = \langle \kappa | r | \tau \rangle$

And express the two charge distributions in terms of their (real, spherical) multipoles:

$$M_{LM}^{\mu\nu} = \sqrt{\frac{4\pi}{2L+1}} \langle \mu | r^L S_{LM} | \nu \rangle$$

For one-center charge distributions $L = |l_{\mu} - l_{\nu}| \dots l_{\mu} + l_{\nu}$

Assuming the two local coordinate systems are aligned and the charge distributions are not overlapping, the **binolar expansion** vields:

1.

$$(\mu\nu|\kappa\tau)_{multipole} = \sum_{L_1} \sum_{L_2} \sum_{m=-l_{<}}^{l_{<}} f_{L_1L_2m} \frac{M_{L_1m}^{\mu\nu} M_{L_2m}^{\kappa\tau}}{R^{L_1+L_2+1}}$$
$$f_{L_1L_2m} = \frac{(-1)^{L_2+|m|}}{\sqrt{(L_1+L_2)! (L_1+m)! (L_2+m)! (L_1-|m|)! (L_2-|m|)!}}$$

Bühler, RJ, Hirschfelder, JO Bipolar Expansion of Coulombic Potentials, Phys. Rev., 1951, 83, 628-633

The multipole formula **becomes fully accurate** (at least 16 digits) once the charge distributions don't overlap.

- Evaluating the multipole formula exactly is too costly the estimate may become as expensive or more expensive than the actual integral calculation
- ✓ For the purpose of pre-screening, one should only be interested in the **lowest multipole** interaction, because it is the one that covers the longest distances:

Estimate
$$(\mu\nu|\kappa\tau)_{est-III} = max \left(\left| f_{L_1^{min}L_2^{min}m} \right|_{max} \frac{\left| M_{L_1^{min}m}^{\mu\nu} \right| \left| M_{L_2^{min}m}^{\kappa\tau} \right|}{R^{L_1^{min}+L_2^{min}+1}} \right)$$

- Will break down for overlapping charge distributions overlap
- Not cheap to compute
- Misses higher order multipole contributions.
- NOT an upper bound, i.e. Will perhaps dramatically underestimate the integral @medium R

Integral Estimates (IV): The "QQR" and "CSAM"

- Lambrecht and Ochsenfeld J. Chem. Phys., 2005, 123, 184102 derived rigorous upper bounds on the basis of the multipole expansion (too expensive in practice)
- Maurer, Ochsenfeld et al. J. Chem. Phys., 2012, <u>136</u>, 144107 realized that higher multipoles can be simulated by the Schwartz integral and proposed the "QQR" estimate:

Estimate
$$(\mu\nu|\kappa\tau)_{est-IV} = max\left(\frac{Q_{\mu\nu}Q_{\kappa\tau}}{R - ext_{\mu\nu} - ext_{\kappa\tau}}\right)$$

The **extent** of a charge distribution is defined by:

$$ext_{\mu\nu} = \sqrt{\frac{2}{\alpha + \beta}} erf^{-1}(1 - \tau)$$

 $\alpha,\beta=exponents,\tau\approx 10^{-4}-10^{-6}, erf^{-1}=inverse$ error function

Thompson and Ochsenfeld et al. J. Chem. Phys., 2017, 144101 further tweaked the QQR by realizing that the distance dependence can be simplified

$$(\mu\nu|\kappa\tau) \leq \sqrt{(\mu\mu|\kappa\kappa)}\sqrt{(\nu\nu|\tau\tau)} = T_{\mu\kappa}T_{\nu\tau}$$
$$(\mu\nu|\kappa\tau) \leq \sqrt{(\mu\mu|\tau\tau)}\sqrt{(\nu\nu|\kappa\kappa)} = T_{\mu\tau}T_{\nu\kappa}$$

Which features the distance dependence of the interacting bra/ket distributions.

Defining:

$$\tilde{T}_{\mu\kappa} = \frac{T_{\mu\kappa}}{\sqrt{Q_{\mu\mu}Q_{\kappa\kappa}}}$$

Gives the final (CSAM) estimate:

Estimate
$$(\mu\nu|\kappa\tau)_{est-V} = Q_{\mu\nu}Q_{\kappa\tau}max(\tilde{T}_{\mu\kappa}\tilde{T}_{\nu\tau},\tilde{T}_{\mu\tau}\tilde{T}_{\nu\kappa})$$

Comparison of Estimates



Performance in practice: (Gly)₁₅/def2-SVP



	#(Cycles)	Energy (Eh)	#(Fock	time/sec)
Schwartz	11 -31	75.70618041870	09 1	822
QQR	11 -31	75.7061 <u>791467</u>	<u>75</u> 1	659 (9%)
CSAM	11 -31	75.7061 <u>774913</u>	<u>13</u> 1	561 (15%)
Almloef	Wild div	vergence		

... 0.3 Microhartree loss of accuracy for 15% performance gain (Will be more for larger systems)

Pre-screening: Wrapping up

- ✓ The best way to speed up a computation is to not do it :-)
- ✓ Identifiying near zero's is and art & science that is not done even after 30+ years
- ✓ In skipping small contributions:
 - It is good but not strictly necessary to have rigorous upper bounds
 - Numerical stability must never be sacrificed

Always remember:

- Computing a bad number fast is useless because it is still a bad number
- First the approximation has to meet a specified accuracy goal, <u>then</u> it *can* be fast

PART 2

How to compute things you cannot avoid efficiently

Chapter 1:

Scaling Laws and Their Impact on Algorithms

Scaling Laws

A quantum chemical algorithm can be characterized by it's scaling behavior:

Scaling with respect to system size (#(Atoms), #(Basis functions),...) Scaling with respect to basis set (Size, Angular momentum,...)

A scaling law can be written as:

$$T = aN^b$$

- T Time taken by algorithm
- *a* ,Prefactor'
- *b* Scaling Exponent

Optimizing an algorithm:

Holy grail:

Bring down the prefactor Bring down the scaling *Linear scaling* with a small prefactor

Figuring out the Scaling Law

General:

Dimensionality of target quantity x Scaling of loops required to obtain it

Example:

$$\psi_p(\mathbf{r}) = \sum_{\mu} c_{\mu p} \varphi_p(\mathbf{r})$$

The number of occupied and virtual MOs is proportional to system size



 $(\mu \nu \mid \kappa \tau)$ Number of AOs integrals proportional to N⁴ (O(N⁴))



Prefactor vs Scaling



For many applications nonlinear scaling with a small prefecture is the preferred choice

In developing reduced scaling algorithms one shoots for *early crossover*

Golden Law of Development

- In general, the workflow of a quantum chemical algorithm contains many steps (e.g. localization, integral transformation, equation solution, perturbative correction, ...),
 Each step will have its own scaling law.
- Each step will have its own scaling law



Profile your Program!

Total execution time	•••	153019.575 sec
Localization of occupied MO's		7516.449 sec (4.9%)
Fock Matrix Formation		11392.614 sec (7.4%)
First Half Transformation		37824.285 sec (24.7%) // I his is worth your
RI-PNO integral transformation	• • •	17832.376 sec (11.7%) <i>while!</i>
Initial Guess	• • •	5376.961 sec (3.5%)
DIIS Solver		8855.850 sec (5.8%)
State Vector Update	• • •	1.744 sec (0.0%)
Sigma-vector construction	• • •	8177.969 sec (5.3%)
<0 H D>	• • •	0.072 sec (0.0% of sigma)
<0 H S>	•••	0.003 sec (0.0% of sigma)
<d h d>(0-ext)</d h d>	•••	575.591 sec (7.0% of sigma)
<d h d>(2-ext Fock)</d h d>	•••	1.921 sec (0.0% of sigma)
<d h d>(2-ext)</d h d>	•••	1512.608 sec (18.5% of sigma)
<d h d>(4-ext)</d h d>	• • •	684.157 sec (8.4% of sigma)
<d h d>(4-ext-corr)</d h d>	•••	2880.920 sec (35.2% of sigma) Gain from optimizing
CCSD doubles correction	•••	33.534 sec (0.4% of sigma) these steps?
<s h="" s="" =""></s>	•••	78.695 sec (1.0% of sigma)
<s h d>(1-ext)</s h d>	•••	79.135 sec (1.0% of sigma)
<d h s>(1-ext)</d h s>	•••	5.117 sec (0.1% of sigma) 🦯 🛛
<s h d>(3-ext)</s h d>	•••	28.949 sec (0.4% of sigma)
CCSD singles correction	•••	0.108 sec (0.0% of sigma) $/$
Fock-dressing	•••	1541.152 sec (18.8% of sigma) /
Singles amplitudes	•••	15.255 sec (0.2% of sigma) /
(ik jl)-dressing	•••	441.823 sec (5.4% of sigma)/
(ij ab),(ia jb)-dressing	•••	213.171 sec (2.6% of sigma)
Pair energies	• • •	1.235 sec (0.0% of sigma)
Total Time for the density		632.934 sec (0.4% of ALL) [*]
Total Time for computing (T)	• • •	32529.433 sec (21.3% of ALL)

Chapter 2

Writing Efficient Programs
The Do's and Don't's of Programming: Overview

Some rules for scientific programming that are relevant for obtaining high performance:

- Avoid short, nested Loops
- Avoid Multidimensional Arrays
- Access arrays in "Unit Stride"
- Avoid indirect addressing
- Make use of matrix multiplications and BLAS
- Make use of LAPACK
- Move redundant work out of the inner loops
- Minimize disk I/O, do it in larger chunks and do it as far ,outside' as possible
- Watch out of Load Balancing in parallel programming

Instruction Pipelines and Logic

Ideal: The CPU has preloaded a ,pipeline' of instructions and the data required to perform the next operations is in the CACHE



A logical instruction whose outcome can not be predicted at compile time brings the CPU and CACHE out of the ,groove'



Chapter 2.1

Unit stride and avoiding short loops

Unit Stride Access

The CACHE has a finite size that is rather small. If one loads an array into the CACHE that is larger than the CACHE size, one should avoid ,jumping' around in the array but rather only access consecutive positions in the array (**unit stride access**)

Example: Say, the CACHE holds 1024 array elements and we want to add up the elements of an array y that contains 2048 elements.

Good: x=0 for (i=0;i<2048;i++) x=x+y[i]</pre>

The compiler can optimize well: load the first 1024 elements of y and the next 1024 elements.
 Performs optimally without any ,CACHE misses'

Bad:

```
x=0
```

for (i=0;i<2048;i++) x=x+y[yorder[i]]</pre>

or for (i=0;i<2048;i++) x=x+y[i]-y[N-i-1]

Two problems:

- yorder[i] may be anything in the range 0..2047 for any i and hence we may have to reload y into the CACHE multiple times
- We use ,indirect addressing'. There is no way for the compiler to know the value of yorder[i] and hence after each addition we have to look again which element of y we need next.

Example: Loop Unrolling

Time critical routines should not contain logic and should not contain nested loops. The process of eliminating short loops in favor of hand optimized, explicit code is called ,Loop unrolling'

Example: Calculation of integrals using the McMurchie/Davidson method

In the MD method, molecular integrals can be very elegantly calculated using an expansion of the Gaussian product in a Gaussian Hermite basis

Cartesian Gaussian on center A: $G^A_{abc;\alpha} = (x - X_A)^a (y - Y_A)^b (z - Z_A)^c \exp(-\alpha r_A^2)$

Repulsion integral in MD:

const

$$\left(G_{abc;\alpha}^{A}G_{a'b'c';\beta}^{B} \mid G_{def;\gamma}^{C}G_{d'e'f';\delta}^{D}\right) = f_{\alpha\beta\gamma\delta}\sum_{t=0}^{a+a'}\sum_{u=0}^{b+b'}\sum_{v=0}^{c+c'}E_{t}^{AB}E_{u}^{AB}E_{v}^{AB}\sum_{t'=0}^{d+d'}\sum_{u'=0}^{e+e'}\sum_{v'=0}^{f+f'}(-1)^{t'+u'+v'}E_{t'}^{CD}E_{u'}^{CD}E_{v'}^{CD}R_{t+t',u+u',v+v'}^{CD}E_{v'}^{CD}E_{v'}^{CD}E_{v'}^{CD}E_{v'}^{CD}R_{t+t',u+u',v+v'}^{CD}E_{v'}^{$$

Expansion of G^AG^B in Hermite basis Expansion of G^cG^p in Hermite basis

Integrals in Hermite basis

Example: Short Loops and Multidimensional Arrays

Pseudocode for a general MD integral routine

```
Calculate Array EAB
Calculate Array ECD
                          recursive formulas. Nested loops of length \sim I_A + I_B (or I_C + I_D)
Calculate Array R
loop ixyz over Cartesian components of A
  loop jxyz over Cartesian components of B
     loop kxyz over Cartesian components of C
       loop lxyz over Cartesian components of D
        \mathbf{x}=\mathbf{0}
        loop t =0..a+a'
         loop u =0..b+b'
                                               10 nested loops!
           loop v = 0..c+c'
                                               For s and p functions these run basically from 0 to 1
            loop v' = 0..f+f'
             loop t' =0..d+d'
              loop u' =0..e+e'
                 x=x+ EAB[x][a][a'][t]*EAB[y][b][b'][u]*EAB[z][c][c'][v]
                     *ECD[x][d][d'][t']*ECD[y][e][e'][u']*ECD[z][f][f'][v']*(-1)<sup>t'+u'+v'</sup>
                     *R[t+t'][u+u'][v+v']
              end loops t',u',v'
         end loops t,u,v
          ELREP[ixyz][jxyz][kxyz][lxyz]=x
end loops i, j, k, lxyz
```

Example: Short Loops and Multidimensional Arrays

Alternative: For low angular momenta create hand optimized routines and store integrals in linearized arrays

	()	Calc_sssp	p()
CalC_SSSS		ab	= a+b
db		cd	= c+d
	$= C + \alpha$	abcd	= ab+cd;
abcd	= ab+cd;	pprim	<pre>= 4.0*ab*cd*sqrt(abcd);</pre>
pprim	= $4.0*ab*cd*sqrt(abcd);$	SR	= Kab*Kcd/pprim;
SR	= Kab*Kcd/pprim;	PQX	= (PX-QX);
PQX	= (PX-QX);	POY	= (PY-OY);
PQY	= (PY-QY);	PÕZ	= (PZ - OZ);
PQZ	= (PZ-QZ);	RPO2	= POX*POX+POY*POY+POZ*POZ;
RPQ2	= PQX * PQX + PQY * PQY + PQZ * PQZ;	~ W	= $ab^*cd/abcd;$
W	= ab*cd/abcd;	RT	= W*RPO2;
RT	= W*RPQ2;	Calc F	Function(F)
Calc_F_Function(F)		t1 – –	= W/cd*F[1];
ELREP[0] = F[0] * SR;	ELREP[(0] = (ODZ * F[0] + POZ * t1) * SR;
		ELREP[1	1] = (ODX*F[0]+POX*t1)*SR;
		ELREP[2	2]= (QDY*F[0]+PQY*t1)*SR;

NO logic, **NO** short loops > The compiler can optimize this code most efficiently

Efficient modern integral libraries (e.g. libint) make use of machine generated, highly unrolled code

Numerical Example

	unoptimized code	unrolled code	libint	speedup
(SS SS) (10 ⁷ times)	1.8	1.2	0.7	(3x)
(pp pp) (10 ⁶ times)	8.3	2.6	0.4	(21x)
(dd dd) (104 times)	4.1	0.4	0.1	(41x)
(ff ff) (10 ³ times)	9.1	0.5	0.2	(45x)

"to a large extend the efficiency of a computer code is a result of the care taken during the implementation stage and not due to the particular method selected for implementation." - Roland Lindh

A Real Life Example

In the theory of multipole interactions, There occurs the multipole interaction tensor:

$$T_{LM,L'M'} \left(\mathbf{R}_{P} - \mathbf{R}_{Q} \right) = (-1)^{L'} I_{L+L'}^{M+M'*} \left(\mathbf{R}_{P} - \mathbf{R}_{Q} \right)$$
$$= (-1)^{L'} \sqrt{\frac{4\pi (L+M)! (L-M)!}{(2L+1)}} R_{PQ}^{-L-1} Y_{L+L'}^{M+M'*} \left(\theta_{PQ}, \phi_{PQ} \right)$$

(\mathbf{R}_{P} , \mathbf{R}_{Q} are the multipole expansion centers, we need all I=0..L, I'=0..L')

Let us write a subroutine for that - easy, we only have to evaluate some spherical harmonics, right?

... Better version

void MultipoleInteractionTensor_ILMreal_Fame(double *PQ, double &RPQ, int LmaxBra, int LmaxKet, TRMatrix &T){

```
//Load all required ILM at once
Load_AllIlm_stddef(LmaxBra+LmaxKet,X,Y,Z,R,ILM.p); no redundant computation
// Start looping over angular momenta
int64 LM1=0;
double *TLM1 = \&(T(0,0));
for (int64 L1=0;L1<=LmaxBra;L1++){</pre>
                                      Correct loop order for unit stride access (Bra indices)
  for (int64 M1=-L1;M1<0;M1++){</pre>
    int64 LM2 = 0;
    for (int64 L2=0;L2<=LmaxKet; L2++){</pre>
                                          Correct loop order for unit stride access (Ket indices)
     // Pointer to correct ILM vector
                                      ));
      double *ILre = &(ILM(L12,0
      double phase1 = (L2%2==0)? 1.0 : -1.0; No silly array access for phase factors
        // Now we have to distinguish four cases to form Trygve's real interaction tensor
      for (int64 M2=-L2;M2<0;M2++){</pre>
                                              No Logic in the innermost loop
        // Common prefactor
        double PreFac
                         = phase1 * 2.0;
        double ILM P re = ILre[M12Paddr]; // Real part
                                                              of I( L1+L2, M1+M2)
        double ILM M re = ILre[M12Maddr]; // Real part
                                                              of I( L1+L2, M1-M2)
                         = PreFac* (-ILM P re + phase2 * ILM M re);
        TLM1[LM2]
        LM2++:
        }// M2
```

... three more loops like this for the three other cases

MultipoleInteractionTensor Shame vs Fame



The version of the same function that has memory optimized unit stride access performs a FACTOR of 2-4 faster for non-trivial angular momenta

Chapter 2.2

Making the most out of using External Libraries

Libraries: The only ones you really need

Relying on third party software that may or may not be maintained in long term or may or may not be portable between platforms can be dangerous! There are three you likely cannot avoid:

1. BLAS (Basic Linear Algebra System)

- a) Level 1: Vector/Vector operations
- b) Level 2: Matrix/Vector operations
- c) Level 3: Matrix/Matrix operations

2. LAPACK (Linear Algebra Package)

Linear algebra routines (Diagonalization, Linear equation systems, Cholesky decomposition, singular value decomposition, ...)

3. MPI (Message Passing Interface)

Low level routines for parallelization using a distributed memory paradigm

These are highly efficient, standardized and portable libraries.

(In ORCA, we nevertheless have put one software layer above them in order to have no direct calls to third party software whatsoever)

Example: The power of BLAS

Let us look at two ,innocent' matrix multiplications:

$$\mathbf{C} = \mathbf{A}\mathbf{B} \quad C_{ij} = \sum_{k} A_{ik} B_{kj}$$

$$\mathbf{C} = \mathbf{A}\mathbf{B}^T \quad C_{ij} = \sum_k A_{ik} B_{jk}$$

Which we can program as follows:

```
loop i = 1 ... N
loop j = 1 ... N
x=0.0;
loop k = 1 ... N
x=x+A(i,k)*B(k,j); or x=x+A(i,k)*B(j,k)
end loop k
C(i,j)=x;
end loop j
end loop i
```

Example: The power of BLAS (II)

For two densely filled essentially random, square matrices A and B with N=2750

		directly programmed	BLAS (dgemm)
$\mathbf{C} = \mathbf{A}\mathbf{B}$:	99	1.7
$\mathbf{C} = \mathbf{A}\mathbf{B}^T$:	11	1.7
$\mathbf{C} = \mathbf{A}^T \mathbf{B}$:	104	1.7
		Why that?	USE BLAS LEVEL 3 (DGEMM) WHENEVER YOU CAN!)
$\mathbf{A} =$	• • •	 The matrices are arr places. Hence A(i,k) A(k,i) is not! Huge (factor 10!) per Even worse would be the main memory (error) 	ange row-wise in contiguous memory is accessing the matrix in unit stride while rformance penalty! e to have rows scattered somewhere in .g. Numerical Recipes matrix routines in C)
			Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz

Example: The power of LAPACK

Example: 3000x3000 matrix

	Hand written	Intel-MKL		
Diagonalization	28.1 sec	dsyevr 5.3 sec ~5x		
Cholesky decomposition	2.4 sec	dpotrf 0.2 sec ~12x		
Singular value decomposition	315.0 sec	dgesvd 21.7 sec ~25x		

Application: Integrals over Gaussians

The form of a Gaussian

$$\varphi_{\mu}^{A}(\mathbf{r}) = N_{abc}(x - X_{A})^{a_{\mu}}(y - Y_{A})^{b_{\mu}}(z - Z_{A})^{c_{\mu}}\exp(-\alpha_{\mu}r_{A}^{2})$$

Multiplying two Gaussians





The McMurchie Davidson Method (I)

If we have a nasty polynomial, we can expand it in terms of nice polynomials

Orthonormal harmonic oscillator eigenfunctions (Hermite polynomials)

$$\begin{split} H_n(x) &= (-1)^n \exp(x^2) \frac{d^n}{dx^n} \exp(x^2) & H_0(x) = 1 \\ H_1(x) &= 2x \\ H_2(x) &= 4x^2 - 2 \\ H_3(x) &= 8x^3 - 12x \end{split}$$

In one dimension:
$$(x - X_A)^{a_{\mu}} (x - X_B)^{a_{\nu}} = \sum_{t=0}^{a_{\mu} + a_{\nu}} E_t^{a_{\mu} a_{\nu}} H_t (x - X_P)$$

Recursion relation $E_t^{i+1,j} = \frac{1}{2p} E_{t-1}^{ij} + (X_P - X_A) E_t^{ij} + (t+1) E_{t+1}^{ij}$



Potential and (Nuclear Wavefunction) 2

(a) McMurchie, L.E.; Davidson, E. (1978) J. Comp. Phys. <u>26</u>,218

 $E_0^{00} = 1$

(b) Helgaker, T.; Raylor, P.R. (1995) in: Yarkony (Ed.) Modern Electronic Structure Theory, World Scientific, 725ff

McMurchie-Davidson in Matrix Form



Obviously, not all tuv combinations occur for each member of the shell pair

#(tuv) combinations as a function of angular momenta

L_{μ}/L_{ν}	0	1	2	3	4
0	1	4	10	20	35
1		10	20	35	56
2			35	56	84
3				84	120
4					165

Advantage of Factorisation



Never worse than O(L⁸) which is better than O(L¹⁰) in the original MD

FN The SHARK integral generation and digestion system, J. Comp. Chem., 2022, 1-16 (DOI: 10.1002/jcc.26942)

SHARK vs Libint: "in vitro"

I∖D	1	2	3	4	5	6
0	1.371	1.724	2.073	2.268	2.347	2.430
1	1.100	1.003	1.052	1.169	1.157	1.174
2	1.199	0.765	0.778	0.793	0.802	0.791
3	5.471	3.628	3.610	3.355	3.787	3.406
4	8.758	5.284	5.210	5.229	5.264	5.227
5	9.792	5.589	4.409	4.360	4.655	4.788
6	18.567	5.912	4.852	4.704	4.751	4.604
7	31.497	6.749	5.360	5.227	5.412	5.031

Contraction Depth

<1 Libint is faster >1 SHARK is faster

Notes: The numbers in the table give tm(Libint)/tm(SHARK), where tm(X) is the time taken by algorithm X to calculate the indicated number of integral batches.

FN The SHARK integral generation and digestion system, J. Comp. Chem., 2022, 1-16 (DOI: 10.1002/jcc.26942)

Angular Momentum

Chapter 2.3

Finding Algorithms with Minimal <u>FLoating Point Operations</u>

(... and whether this is the ultimate goal)

Design of an algorithm: FLOP count

In the early days of algorithm design, developers were carefully minimizing the number of **floating point operations (FLOPs)** required to accomplish a given task

Example: Partial integral transformation $(\mu \nu \mid \kappa \tau) \rightarrow (ia \mid jb)$

i,j= occupied MOs (#=O), a,b, unoccupied MOs (#=V), μ , ν , κ , τ =basis functions (#=B)

$$\psi_p(\mathbf{r}) = \sum_{\mu} c_{\mu p} \varphi_p(\mathbf{r})$$

Naive:
$$(ia \mid jb) = \sum_{\mu} \sum_{\nu} \sum_{\kappa} \sum_{\tau} c_{\mu i} c_{\nu a} c_{\kappa j} c_{\tau b} (\mu \nu \mid \kappa \tau)$$
 $FLOPS = B^4 O^2 V^2$

O(N⁸) scaling

Must be possible to do better than that

FLOP Count: Partial Integral transformation

Algorithm A: occupied indices first

Algorithm B: virtual indices first

$$(i\nu \mid \kappa\tau) = \sum_{\mu} c_{\mu i}(\mu\nu \mid \kappa\tau) \qquad (B^{4}O) \quad 3125 \qquad (\mu a \mid \kappa\tau) = \sum_{\nu} c_{\nu a}(\mu\nu \mid \kappa\tau) \qquad (B^{4}V) \quad 28215$$

$$(i\nu \mid j\tau) = \sum_{\kappa} c_{\kappa j}(i\nu \mid \kappa\tau) \qquad (O^{2}B^{3}) \quad 312 \qquad (\mu a \mid \nu b) = \sum_{\tau} c_{\tau b}(\mu a \mid \kappa\tau) \qquad (V^{2}B^{3}) \quad 25312$$

$$(ia \mid j\tau) = \sum_{\nu} c_{\nu a}(i\nu \mid j\tau) \qquad (O^{2}VB^{2}) \quad 281 \qquad (ia \mid \nu b) = \sum_{\mu} c_{\mu i}(\mu a \mid \nu b) \qquad (OV^{2}B^{2}) \quad 2531$$

$$(ia \mid jb) = \sum_{\tau} c_{\tau b}(ia \mid j\tau) \qquad (O^{2}V^{2}B) \quad 253 \qquad (ia \mid jb) = \sum_{\nu} c_{\nu j}(ia \mid \nu b) \qquad (O^{2}V^{2}B) \quad 253$$
Four O(N⁵) steps
ratio of FLOP counts:
$$\frac{\#(FLOPS)_{A}}{\#(FLOPS)_{B}} = \frac{O}{V} \frac{(2B^{3} - V^{3})}{(B^{2} + 3B^{2}V - 3BV^{2} + V^{3})} < 1 \qquad 0.07$$

Always transform the index first that offers the largest data reduction!

Example: GFLOPS for B=500, O=50, V=450

FLOP count versus Performance

In order to capitalize on the efficiency of the BLAS routines, it is sometimes advantageous to sacrifice optimal FLOP count.

Example: Integral direct partial integral transformation for MP2

$$E_{MP2} = -\frac{1}{4} \sum_{i,j,a,b} \frac{\left[(ia \mid jb) - (ib \mid ja) \right]^2}{\varepsilon_a + \varepsilon_b - \varepsilon_i - \varepsilon_j}$$

Key step: integral transformation

$$(ia \mid jb) = \sum_{\mu} \sum_{\nu} \sum_{\kappa} \sum_{\tau} c_{\mu i} c_{\nu a} c_{\kappa j} c_{\tau b} (\mu \nu \mid \kappa \tau)$$

FLOP count optimized algorithm

```
Full eightfold permutation symmetry used
loop ibatch over batches of occupied MOs
   loop p=1..NBas
                                                              have to be able to store N_{Bas}^3 integrals for each
     loop q=1...p
                                                              occupied MO. Hence need batches of occupied
        loop r=1...p
                                                                                MOs
          loop s=1..r|q
              Calculate (pq|rs)
              loop i=1..Nocc (in ibatch)
              ITMP[p,q,r,i]+= C<sub>occ</sub>[s,i]*(pq|rs) and non-redundant permutations of indices
          end i in ibatch
    end loops p,q,r,s
                                                                 Transformation of 2<sup>nd</sup> index
    loop p=1..NBas
         loop r=1..NBas
           loop i=1,...Nocc (in ibatch)
               loop j=1..i
                loop q=1..NBas
                    JTMP[p,j,r,i] += C_{occ}[q,j] * ITMP[p,q,r,i]
                 end loop q
    end loops j,i,r,p
                                                                   Transformation of 3<sup>rd</sup> index
    loop i=1..Nocc (in ibatch)
       loop j=1..i
         loop p over AO's
           loop b=1..NVirt
              loop r over AO's
                                                                Transformation of 4<sup>th</sup> index
                ATMP(p,b) +=C[r,b] *JTMP[p,j,r,i]
         end loops r,b,p
         loop a=1..Nvirt
           loop b=1..Nvirt
              loop p over AO's
                KIJ[a,b] += C[p,a] *ATMP[p,b]
         end loops p,a,b
         Evaluate MP2 amplitudes and pair energy
    end loops i,j
        end loop i
end loop ibatch
```

BLAS optimized algorithm



Performance Test



Diclophenac

def2-TZVP (667 basis functions)

4 GB main memory used

FLOP optimized algorithm

(25 batches necessary)

>100,000 sec

1732 sec

BLAS optimized algorithm

TOTAL TIME for half transformation:	1697.0	sec
AO-integral generation :	1078.9	sec
Half transformation :	354.0	sec
K-integral sorting :	60.4	sec

Chapter 2.4

Using Factorizations and Finding the ones with the best FLOP count

Example: Factorization in Coupled Cluster

The scaling of an algorithm can sometimes be reduced through **factorization.** This happens if intermediates can be defined that only depend on a subset of the summation indices. In this case the summations can be carried out in two steps:

Look at one nonlinear term in the CCSD equations:

$$\sigma_{ab}^{ij} \leftarrow \sum_{kl} \sum_{cd} \langle kl \mid \mid cd \rangle t_{cd}^{ij} t_{ab}^{kl}$$
 O(N⁸) scaling

- 4 target indices
- 4 summation indices
- ... But any quantity depends on only 2 target indices at a time
- Must be able to re-arrange loops more cleverly

Two possibilities:

$$\sigma_{ab}^{ij} \leftarrow \sum_{cd} t_{cd}^{ij} \underbrace{\sum_{kl} t_{ab}^{kl} \left\langle kl \mid \mid cd \right\rangle}_{Y_{cd}^{ab}}$$

Example: Factorization in Coupled Cluster

$$\sigma_{ab}^{ij} \leftarrow \sum_{kl} t_{ab}^{kl} \sum_{\underline{cd}} \langle kl \, || \, cd \rangle t_{cd}^{ij} \qquad : \quad X_{kl}^{ij} = \sum_{cd} \langle kl \, || \, cd \rangle t_{cd}^{ij} \qquad N_{\text{occ}}^4 \text{ Storage}_{N_{\text{occ}}^4 \text{N}_{\text{virt}^2} \text{ FLOPS}}$$

$$\sigma_{ab}^{ij} \leftarrow \sum_{kl} t_{ab}^{kl} X_{kl}^{ij} \qquad N_{\text{occ}}^4 \text{N}_{\text{virt}^2} \text{ FLOPS}_{N_{\text{occ}}^4 \text{N}_{\text{virt}^2} \text{ FLOPS}}$$

$$\sigma_{ab}^{ij} \leftarrow \sum_{cd} t_{cd}^{ij} \sum_{\underline{kl}} t_{ab}^{kl} \langle kl \, || \, cd \rangle \qquad : \quad Y_{cd}^{ab} = \sum_{kl} t_{ab}^{kl} \langle kl \, || \, cd \rangle \qquad N_{\text{virt}^4} \text{ Storage}_{N_{\text{occ}}^2 \text{N}_{\text{virt}^2} \text{ FLOPS}_{N_{\text{occ}}^2 \text{N}_{\text{virt}^4} \text{ FLOPS}}$$

$$\sigma_{ab}^{ij} \leftarrow \sum_{cd} t_{cd}^{ij} \sum_{\underline{kl}} t_{ab}^{kl} \langle kl \, || \, cd \rangle \qquad : \quad Y_{cd}^{ab} = \sum_{kl} t_{ab}^{kl} \langle kl \, || \, cd \rangle \qquad N_{\text{virt}^4} \text{ Storage}_{N_{\text{occ}}^2 \text{N}_{\text{virt}^4} \text{ FLOPS}_{N_{\text{occ}}^2 \text{N}_{\text{virt}^4} \text{ FLOPS}}$$

$$\sigma_{ab}^{ij} \leftarrow \sum_{cd} t_{cd}^{ij} Y_{cd}^{ab} \qquad N_{\text{occ}^2 \text{N}_{\text{virt}^4} \text{ FLOPS}_{N_{\text{occ}^2} \text{N}_{\text{virt}^4} \text{ FLOPS}_{N_{\text{occ}^2} \text{N}_{\text{virt}^4} \text{ FLOPS}_{N_{\text{virt}^4} \text{ Storage}_{N_{\text{virt}^4} \text{ Storage}_{N_{\text{virt}$$

Let us return to our initial question!

We had:

i,j,k,l=occupied a,b,c,d=virtual



STEP 1: What is your data and how do you store what?

How many occupied (n_o) and virtual (n_v) orbitals do I have?

- Say $n_0=50$, $n_v=400$, then $(n_0*n_v)^{2\sim 3}$ GB storage, $(n_v)^4=190$ GB
- Probably need to store that on disk and retrieve in portions

STEP 2: Rewrite the equations in matrix form

Integrals: $\langle kl || cd \rangle = (kc|ld) - (kd|lc) \equiv K^{kl}(c,d) - K^{kl}(d,c)$ Amplitudes: $t_{ab}^{ij} = T^{ij}(a,b)$

Series of matrices orderd by internal label pairs

Always look for matrices and vectors!

$$\widetilde{K}^{kl}(c,d) \equiv K^{kl}(c,d) - K^{kl}(d,c) \qquad \left(\boldsymbol{\sigma}^{ij} \leftarrow P(ij)P(ab) \sum_{kl} \boldsymbol{T}^{ik} \widetilde{\boldsymbol{K}}^{kl} \boldsymbol{T}^{lj} \right)$$

How it is REALLY NOT done

```
Get array SIGMA(i,j,a,b)
Get array T(i,j,a,b)
Get array KS(i,j,a,b)
Loop over i \ge j O(N^2)
  Loop over a,b
                        O(N^2)
                           O(N^2)
    Loop over k,l
     Loop over c,d
                             O(N^2)
        SIGMA(i,ja,b,)+= T(i,k,a,c)
                         *T(j,l,d,b)
                         *KS(k,l,c,d)
      End c,d
    End k,l
  End a, b
```

Overall O(N⁸) With heavy Memory demands and no BLAS

End i,j

How it is ALSO NOT done

```
Loop over pairs i>=j
                            O(N^2)
  Get matrix SIGMA(i,j)
  Loop over k,l
                            O(N^2)
    Get matrix T(i,k)
    Get matrix T(j,1)
                                                          Overall O(N<sup>7</sup>)
    Get matrix KS(k, 1) (KS=K-squiggle)
                                                          With heavy
                                                              I/O
     Form intermediate X= KS(k, 1) *T(j, 1)
                                               O(N^3)
     Add to SIGMA(i,j) += X(k,j) *T(j,l)
                                                O(N^3)
  End k,l
  Store matrix SIGMA(i,j)
End i,j
```

How it IS done



Chapter 2.5

Precompute what you can afford to avoid redundant re-computation
Precomputed quantities

Example: shell pair data

```
Loop ish>= jsh Highly redundant since independent of ish, jsh, A or B!

Calculate KIJ=di*dj*exp(-ai*aj/(ai+aj)*R_{AB}^2)

Calculate P = 1/(ai+aj)*(ai*R_A+aj*R_B)

Loop ksh, lsh(ish==ksh?jsh:ksh)

Calculate KKL=dk*dl*exp(-ak*al/(ak+al)*R_{CD}^2)

Calculate Q = 1/(ak+al)*(ak*R_C+al*R_D)

Calculate (IJ|KL) {P,Q,KAB,KCD,...}
```

•••

Better: Precompute shell pair data AND screen for negligible shell pairs

```
Loop ish>= jsh
Calculate KIJ=di*dj*exp(-ai*aj/(ai+aj)*R<sub>AB</sub><sup>2</sup>)
if |KIJ|<TCut then reject shell pair
Calculate P = 1/(ai+aj)*(ai*R<sub>A</sub>+aj*R<sub>B</sub>)
Store KIJ, P in memory or on disk
```

Move Work out of the Inner Loops: Split-J

Choosing intermediates wisely such that redundant work is move out of the inner loops helps performance

Example: Integrate integral evaluation as early as possible into the target quantities. For the Coulomb matrix, (Ahmadi & Almlöf): When we calculate the integrals one by

$$\begin{split} J_{\mu\nu} &= \sum_{\kappa\tau} P_{\kappa\tau} \left(\mu\nu \mid \kappa\tau \right) & \text{one, we repeatedly re-calculate this quantity N2 times although it is independent of μ,ν .Likewise:

$$&= \sum_{\kappa\tau} P_{\kappa\tau} \sum_{twv} E_{twv}^{\mu\nu} \sum_{t'u'v'} (-1)^{t'+u'+v'} E_{t'u'v'}^{\kappa\tau} R_{t+t',u+u',v+v'}^{t+t',u+u',v+v'}$$

$$&= \sum_{twv} E_{twv}^{\mu\nu} \sum_{t'u'v'} R_{t+t',u+u',v+v'} \sum_{\substack{\kappa\tau \\ = P_{t'u'v'} \text{ independent of } \mu\nu,tuv}} (-1)^{t'+u'+v'} P_{\kappa\tau} E_{t'u'v'}^{\kappa\tau}$$

$$&= \sum_{twv} E_{twv}^{\mu\nu} \sum_{t'u'v'} R_{t+t',u+u',v+v'} \sum_{\substack{\kappa\tau \\ = P_{t'u'v'} \text{ independent of } \mu\nu,tuv}} (-1)^{t'+u'+v'} P_{\kappa\tau} E_{t'u'v'}^{\kappa\tau}$$$$

FN J. Comp. Chem. 2003, 24, 1740-1747; FN J. Comp. Chem., 2022, 1-16

Performance example



Coulomb term (sec) (20-builds)

5796 sec

Traditional treatment

Split-J algorithm

2834 sec

=Ahmadi-Almlöf =Head-Gordon J-engine



Identical numerical result, same scaling, but significant speedup realized through thoughtful structuring of the entire computational process

Chapter 2.6

Be careful with Input/Output

Example: I/O Heavy Algorithms

The I/O system is the slowest part of your computer!

- ➤ Use it as little as possible
- > Move its usage as far outside in the loop structure as reasonably possible
- Avoid reading small chunks of data

Example: Integral symmetrization in EOM-CCSD

```
6641 sec
                                                                           31 sec
Loop i=1...Nocc
                                                          Loop i=1..N<sub>occ</sub>
                                                              Initialize buffer K<sup>ib</sup> for all b
  loop a=1..N<sub>vir</sub>
   Write NULL matrix Kia into buffer IABC
                                                              loop a=1..Nvir
  end loop a
                                                                read matrix Kia(b,c) from IABC
  loop a=1..N<sub>vir</sub>
                                                                loop b=1..Nvir
    Read matrix K^{ia}(b,c) = (ib|ac) from IABC
                                                                  loop c=1..Nvir
                                                                     K<sup>ib</sup>(a,c)+=Kia(b,c);
    loop b=1..N<sub>vir</sub>
       Read matrix K<sup>ib</sup>(c,d) = (ic|bd) from IABC
                                                                   end loop c
       loop c=1..N<sub>vir</sub>
                                                                end loop b
         K^{ib}(a, c) = +K^{ib}(a, c) + K^{ia}(b, c);
                                                               end loop a
                                                               Write entire buffer Kib into IABC
       end loop c
       Store matrix Kib in IABC
                                                          ∎end loop i
    end loop b
                                                                    SAME operation count!
  end loop a
end loop i
                                                            Factor 200 performance difference!!
```

Chapter 2.7

Parallelization in a nutshell

Single CPU Clockspeed / Single Thread Performance



Growth in clock rate of microprocessors. Between 1978 and 1986, the clock rate improved less than 15% per year while performance improved by 25% per year. During the "renaissance period" of 52% performance improvement per year between 1986 and 2003, clock rates shot up almost 40% per year. Since then, the clock rate has been nearly flat, growing at less than 1% per year, while single processor performance improved at less than 22% per year.





Relative Performance (%)

Performance: Moore's Law

Microprocessor Transistor Counts 1971-2011 & Moore's Law



Consequence's of Moore's Law

Paradigm Change:

Requires explicit parallelization by the programmer!

"From this historical perspective,

it's startling that the whole IT industry has bet its future that programmers will finally successfully switch to explicitly parallel programming"

(Patterson, Hennessy: The Hardware/Software Interface, 2009)

Amdahl's Law of Diminishing returns



Speedup:

$$S(N) = \frac{1}{(1-P) + \frac{P}{N}}.$$

P: Parallel portion of code N: Number of Processors

Parallelization in a Nutshell

Principle idea: let a number of processors, say *n*, work on parts of the computational problem in parallel and combine sub results into the final result.

Ideal Scenario: The problem breaks down perfectly and the time required to solve the problem is 1/*n*.

Shared Memory Models:

- Open MP, POSIX threads
- Efficient use of resources, no memory replication
- Difficult to debug large programs
- Can only be used on one machine with common memory

Message Passing Models

- Communication via messages between processes
- Choice between replicated and distributed memory
- Distributed memory difficult to implement efficiently
- Can be used between machines

Hybrid Models:

- Threads + MPI
- Combines shared memory on one machine with message passing between machines
- Adaptation into official standards is slow





Parallelization

Parallelization is of vital importance in modern high-performance computing, yet a LOT can go wrong here! We can only scratch the surface of this complex subject.

A few rules:

- Each process should have roughly the same amount of work to do (Load Balancing).
- 2. Do the parallelization as far ,**outside**⁴ as possible (e.g. distribute the *outer*most loop).
- 3. Excessive **communication** (e.g. sending large chunks of data) between processes should be avoided as much as possible.
- 4. Synchronization should not happen inside time critical loops and there should be as little of it as possible.
- 5. I/O in parallel applications is difficult if several processes access the same file.

Load Balancing Example

Integral calculation: The time required to calculate a given integral batch is a complicated function of angular momenta, contraction depth, orbital exponents and prescreening efficiency

- Load balancing difficult to guarantee
- Random distribution of batches among processors.
- Uneveness will average out in the limit of many batches





Parallelization and Ahmdal's Law

Group Parallelization

Idea:

Divide the processes into groups that scale well (e.g. up to 8) and then parallelize independently over the groups

! PAL(8x8)



Everolimus 151 atoms, def2-tzvp (2606 basis functions)



Chapter 3

Some Useful Programming Techniques for Writing "good" Software

Prelude: Who are you writing code for?

✓ For yourself because you want to check out some ideas

- ✓ Just for a paper, but not to be used later
- ✓ For your boss because you want to get a Ph.D.

✓ For a program package that is supposed to be long lived

- ... it needs to be well documented (in english)
- ... don't try to be funny!
- ... Write the FM (so that users can avoid reading the FM)
- ... Make sure it compiles on any platform
- ... Minimize the dependence on elements that are outside your control
- ... put effort into making it as efficient as possible

Everything is ok!

Mostly anything is ok!

... depends on your boss

Good Programming Habits

- ✓ Show some **respect** for the work that went into the package:
 - Communicate with the team and dont' go off on lonely tangents
 - Be respectful of the code organization and try to fit in
 - Do not rewrite code of others without prior conversation
- ✓ Write lots of **comments** (in english!) including references to papers.
- ✓ Use the **existing infrastructure** of the package
- ✓ **Debug** carefully before checking into the main branch. Delete debug code
- ✓ Provide plenty of **test jobs** with reference results.
- Profile and optimize the code carefully
- ✓ Be feature complete (e.g. not just closed shell)
- ✓ Program as simple as possible and only as fancy as necessary. Do not "show off"
- ✓ Do not write 50000 line functions break it down, be **modular**, make it **reusable**
- ✓ Use logical, recognizable **file names**

... bad habits are the opposite of everything on the list

Chapter 3.1

Using Recursion

Recursion

Recursion can be a very effective way to arrive at compact, elegant code.

- 1. Define the starting point of the recursion
- 2. Define the recursive conditions strategy
- 3. Define the termination conditions

Easy example: the factorial number:

$$n! = \prod_{i=1}^{n} i \qquad \qquad 0! = 1$$

Straightforward linear programming:

```
int fac_n = 1;
for (int i=1;i<=n;i++) fac n = fac n*i;</pre>
```

Recursive programming: int factorial(int n) { termination if (n>0) return n*factorial(n-1); condition else return 1; };

Recursion: Finding the shortest path from A to B



```
// FN 08/18
// Find a pathway between A and B (if it exists)
//
// ON INPUT A

    the first atom one is looking for

             В
                          - the target atom for the pathway
//
                          - the list of connectivities XAB(A,k) - k'th atom connected to A; k<NAB(A)
11
            XAB
            NAB
                          - the number of connectivities for each atom NAB(A)=number of bonds at A
11
            ActualLength – the recursion depth
11
                          - the maximal length of the path that is allowed
            MaxLength
//
            abortAtLength - abort the search if a path of this length or shorter has been found
11
// ON OUTPUT Length
                          - the length of the pathway found
void GEO_FindPathway(int A, int B, TIMatrix &XAB, TIVector &NAB, int ActualLength, int &Length, int &MaxLength, int abortAtLength)
{
 // Abort if a path with abortAtLength is found
 if (Length <= abortAtLength && Length > 0) return;
 // Abort, if no path is found or we have exceeded the allowed maximal length
 if (ActualLength> MaxLength) return;
 // The number of connections made by atom A
 int NA= NAB(A);
 for (int k=0;k<NA;k++){</pre>
    // C is the actual k'th atom connected to A
   int C= XAB(A,k);
    // If we have found our target atom, we stop the recursion
   if (C==B){
     // first pathway found
     if (Length<0){</pre>
                                                            Start condition
        Length= ActualLength;
     }
     // if we had one before, check whether this one
     // is shorter. We take the shortest
                                                            Check whether actual walk is shorter
     else{
        if (ActualLength<Length) Length=ActualLength;</pre>
                                                           Termination
     }
     // There can be no path shorter than one bond.
                                                                                        Recursion. New
     if (Length==1) return;
     // In any other case, we keep recursing
     GE0_FindPathway(C,B,XAB,NAB,ActualLength+1,Length,MaxLength, abortAtLength);
                                                                                        origin of search is C
     // .... but stop here, since no other atom in that list
     // can match B
      return;
    // and otherwise continue the recursion with the actual atom along the pathway
                                                                                        Recursion. New
    else{
      GE0_FindPathway(C,B,XAB,NAB,ActualLength+1,Length,MaxLength, abortAtLength);
                                                                                        origin of search is C
   }
 }
}
```

Recursion: Setting up ORMAS CI Spaces

ORMAS (Occupation Restricted Multiple Active Spaces)

GAS (Generalized Active Space)

- ✓ The orbital space is divides into N subspaces
- ✓ Subspace K can have Kmin ... Kmax electrons
- ✓ Subspace K has NORBK orbitals
- ✓ Inside each subspace a CAS(NELK,NORBK) is formed
- ✓ All combinations of subspaces that give the correct NEL are wanted

```
* FN 08/2024
*
   Recursive function to figure out the combination of subspace lists that lead
   to the correct number of electrons in the active space
*
                            - the information about the subspaces
   ON INPUT ORMAS
              NELTARGET

    target number of electrons

             ACTNEL- actual number of electron during recursionACTLEVEL- actual level of recursion (e.g. LEVEL=subspace level)ACTCOMB0- actual combination
 *
 *
                            - only count combinations or also store them
              JustCount
   ON OUTPUT NCOMBINATIONS - number of combinations that lead to the correct NEL

    if we store, the combinations themselves

              COMBO
*
*/
void ORMAS FindSubSpaceCombos(TIVector &ORMAS,
                                                                                      // -----
// NOT bottom level: recurse to the next level
                                         NELTARGET,
                              int
                              int
                                         ACTNEL,
                              int
                                         ACTLEVEL,
                                                                                      else{
                              TIVector &ACTCOMBO,
                                                                                        ORMAS_FindSubSpaceCombos(ORMAS,
                                        &NCOMBOS,
                              int
                                                                                                                  NELTARGET,
                              TIMatrix &COMB0){
                                                                                                                  ACTNEL,
 int64 NLEVELS= ORMAS(0);
                                                                                                                  ACTLEVEL+1,
 int64 norb = ORMAS(1+3*ACTLEVEL+0);
                                                                                                                  TCOMBO,
 int64 nmin = ORMAS(1+3*ACTLEVEL+1);
                                                                                                                  NCOMBOS,
 int64 nmax = ORMAS(1+3*ACTLEVEL+2);
                                                                                                                   COMBO,
 TIVector TCOMBO:
                                                                                                                   JustCount);
 TCOMB0.CopyVec(ACTCOMB0);
                                                                                      };
 // Loop over all electron numbers of this level
                                                                                    };
 int ACTNEL0= ACTNEL;
                                                                                  };
 for (int64 n=nmin; n<=nmax; n++){</pre>
   // Remember where we are at this level
                                                                                                           Recursion to next
   TCOMBO(ACTLEVEL)=n:
                                                      Reached last
    // current number of electrons
                                                                                                          next subspace
   ACTNEL= ACTNEL0 + n;
                                                      subspace
   // if that is already too large we can skip
   if (ACTNEL>NELTARGET) break;
   // We arrived at the bottom level: see what we have
   if (ACTLEVEL == NLEVELS -1){
      // That is the correct number of electrons
      if (ACTNEL==NELTARGET){
                                                           Correct number of
        for (int64 i=0;i<NLEVELS;i++){</pre>
          COMBO(NCOMBOS, i) = TCOMBO(i);
                                                           electrons found
       };
       NCOMBOS++;
     } else{
        // nothing to do: number of electrons is incorrect
      }
   }
```

ORMAS(14: 6 10 12 / 2 0 4 / 50 0 2)

no of electrons (total) Subspace 1 Subspace 2 Subspace 3 (domos) (active) (virtuals) 6 MOs, 0-2 holes any 50 MOs occupation 0-2 electrons

Subspace CAS lists

Subspace	1:	10 electrons in 6 orbitals => 21 configu	irations
Subspace	1:	11 electrons in 6 orbitals => 6 configur	ations
Subspace	1:	12 electrons in 6 orbitals => 1 configur	ations
Subspace	2:	0 electrons in 2 orbitals => 1 configur	ations
Subspace	2:	1 electrons in 2 orbitals => 2 configur	ations
Subspace	2:	2 electrons in 2 orbitals => 3 configur	ations
Subspace	2:	3 electrons in 2 orbitals => 2 configur	ations
Subspace	2:	4 electrons in 2 orbitals => 1 configur	ations
Subspace	3:	0 electrons in 50 orbitals => 1 configu	ations
Subspace	3:	1 electrons in 50 orbitals => 50 configu	irations
Subspace	3:	2 electrons in 50 orbitals => 1275 confi	gurations

Subspace combinations

Number of valid subspace combinations = 9

COMBO	1:	10	2	2	=>	80325	CFGs
COMBO	2:	10	3	1	=>	2100	CFGs
COMBO	3:	10	4	0	=>	21	CFGs
COMBO	4:	11	1	2	=>	15300	CFGs
COMBO	5:	11	2	1	=>	900	CFGs
COMBO	6:	11	3	0	=>	12	CFGs
COMBO	7:	12	0	2	=>	1275	CFGs
COMBO	8:	12	1	1	=>	100	CFGs
COMBO	9:	12	2	0	=>	3	CFGs

Chapter 3.1

Object Oriented Programming

Object Oriented Programming

In object oriented programming you celebrate the unity of code and data by creating **classes**, that contain

- Data that are private to the object
- Functions that work on these data
- MUCH safer than having global data that is passed around the program
- MUCH easier to build up and administrate complicated data structures

```
class TGaussianShell{
 private:
                 // number of primitives
  int nprim;
                 // angular momentum
  int 1;
                                                                       data
  int ofs; // position in the basis function list
  double *a, *d; // exponents, contraction coefficients
 public:
  // Constructor and Destructor
  TGaussianShell() {
     nprim =0; l=0; ofs=0; a=nullptr; d=nullptr;
                                                                           Constructor and Destructor
  };
  ~TGaussianShell() {
     if (a!=nullptr) { delete[] a; a=nullptr;
     if (d!=nullptr) { delete[] d; d=nullptr;
   };
  // Setters and Getters
                                                                         Access to data
  int GetNPrim() { return nprim; };
  double *GetA() { return a; };
                                                                         Productive functions
  . . .
  // Productive functions
  void Copy(int xnprim, int xl, int xofs, double *xa, double *xd) {
          nprim= xnprim; .... };
 void Copy(TGaussianShell &SH) { Copy(SH.GetNPrim(),SH.GetL(),SH.GetOfs(),SH.GetA(), SH.GetD());};
  void Normalize();
 void Store(FILE *f);
 void Read(FILE *f);
 void Print();
};
```

```
class TGaussianAtom{
 private:
  int lmaxA; // Highest L for this atom
  int *NShells; // Number of shells in each angular momentum
  TGaussianShell **Shells;
 public:
  TGaussianAtom() { ... }; // initialize
  ~TGaussianAtom() { ....}; // delete data
  // Getters and Setters
  TGaussianAtom *GetShell(int 1, int ish) { return & (Shells[1][ish]); };
  // Productive functions
  void GetMemory(int xlmaxA, int *xNShells;) {
    lmaxA= xlmaxA; NShells= new int[lmaxA+1];
    for (int l=0;l<=lmaxA;i++) NShells[l]= xNShells[l];</pre>
    Shells= new TGaussianShell *[lmaxA+1];
    for (int l=0; l<=lmaxA; l++) Shells[l]=new TGaussianShell[NShells[l];</pre>
 };
  void Normalize() {
    for (int l=0; l<=lmaxA; l++)</pre>
      for (int ish=0; ish<NShells[1];ish++) Shells[1][ish].Normalize();</pre>
 };
 void Copy(TGaussianAtom &GA) {
    GetMemory( GA.GetLmax(), GA.GetNShells() );
    for (int l=0; l<=lmaxA; l++)</pre>
      for (int ish=0; ish<NShells[1];ish++)</pre>
        Shells[1][ish].Copy(*GA.GetShell(1, ish));
  };
  void Store(FILE *f);
 void Read(FILE *f);
 void Print();
};
```

```
class TSegmentedBasisSet{
```

private:

```
// Number of atoms
  int NAtoms;
  TGaussianAtom *Gauss; // the actual Gaussians for each atom
 public:
  TSegmentedBasisSet() { ... }; // initialize
  ~TSegmentedBasisSet() { ....}; // delete data
  // Productive functions
 void GetMemory(int xNAtoms) {
    NAtoms= xNAtoms;
    Gauss= new TGaussianAtom[NAtoms];
  };
  void SetAtomBasis(int A,TGaussianAtom &GA) {
    if (xAtom>=0 and xAtom<NAtoms) Gauss[A].Copy(GA);
    else <throw exception>;
 };
 void Normalize() {
    for (int A=0;A<NAtoms; A++) Gauss[A].Normalize();</pre>
  };
 void CalcOverlap(TSymmetricMatrix &S);
 void Read(FILE *f);
 void Print();
};
```

Inheritance and Virtual Functions

- ✓ Very often, you have a bunch of tasks to do that have something in common and they are embedded in larger tasks, where they perform one specific action.
- ✓ In order to help you streamlining such situations, C++ let's you design "virtual functions"

✓ Example:

Making Virtual Functions Concrete

```
/* _____
* Now we define one concrete task.
*
* The concrete task is an "heir" of Abstract task and
* overloads the virtual function "DoSomething" which is
* will do something concrete.
* NOTE: the virtual function here is declared with "=0"
 * _____
*/
class WindowCleaner : public AbstractTask{
public:
 virtual void DoSomething(const char *house);
};
void WindowCleaner::DoSomething(const char *house)
{
 printf("Cleaning Windows in house %8s\n",house);
}
/* _____
* Another concrete task
*
  _____
*/
class FloorCleaner : public AbstractTask{
public:
 virtual void DoSomething(const char *house);
};
void FloorCleaner::DoSomething(const char *house)
{
 printf("Cleaning floors in house %8s\n",house);
}
```

Incorporating virtual functions

```
/* _____
* Here is how to use it:
* We have a function that is in need of performaing a
* certain task, and pass the concrete task onto it
   _____
 */
void CleaningForce(AbstractTask &MyTask)
{
 int NHouses = 5;
 const char *House[5] = {"Miller", "Smith", "Jones", "Mayer", "Trump"};
 for (int i=0;i<NHouses; i++){</pre>
   printf("Now working on house %d, the home of %s\n",i+1,House[i]);
   printf(" task is -> ");
   MyTask.DoSomething(House[i]);
   printf("\n");
 };
};
```

How it looks in practice

```
int main(int argc, char **argv)
{
 // These are our "workers" - the ones who do something concrete
 WindowCleaner W;
 FloorCleaner F:
 JewelThief J;
 // Now we can drive the performance of concrete tasks easily
 printf("First cleaning task\n");
 printf("-----\n");
 CleaningForce(W);
 printf("Second cleaning task\n");
 printf("-----\n");
 CleaningForce(F);
 printf("Third cleaning task\n");
 printf("-----\n");
 CleaningForce(J);
 // But this will work too
 AbstractTask *X;
 printf("Dynamic cleaning task\n");
 printf("-----\n");
 X= new JewelThief;
 CleaningForce(*X);
}
```
Program output

First cleaning task

Now working on house 1, the home of Miller task is -> Cleaning Windows in house Miller Now working on house 2, the home of Smith task is -> Cleaning Windows in house Smith Now working on house 3, the home of Jones task is -> Cleaning Windows in house Jones Now working on house 4, the home of Mayer task is -> Cleaning Windows in house Mayer Now working on house 5, the home of Trump task is -> Cleaning Windows in house Trump Second cleaning task Now working on house 1, the home of Miller task is -> Cleaning floors in house Miller Now working on house 2, the home of Smith task is -> Cleaning floors in house Smith Now working on house 3, the home of Jones task is -> Cleaning floors in house Jones Now working on house 4, the home of Mayer task is -> Cleaning floors in house Mayer Now working on house 5, the home of Trump task is -> Cleaning floors in house Trump

Third cleaning task

Now working on house 1, the home of Miller task is -> Stealing jewelery in house Miller while the other idiots are cleaning Now working on house 2, the home of Smith task is -> Stealing jewelery in house Smith while the other idiots are cleaning Now working on house 3, the home of Jones task is -> Stealing jewelery in house Jones while the other idiots are cleaning Now working on house 4, the home of Mayer task is -> Stealing jewelery in house Mayer while the other idiots are cleaning Now working on house 5, the home of Trump task is -> Stealing jewelery in house Trump while the other idiots are cleaning ... he he, at least he deserves it

Example: Calculating One-Electron Integrals

```
calculates
SUBROUTINE OneElectronLoop(TIntegralKernel
                                                   &K, -----
                                                           integrals for one
                                                           primitive pair
                               TIntegralConsumer
                                                   &C)
                                                            Does something
                                                            with one batch of
BEGIN
                                                           integrals
  Loop ish over shells
    Loop jsh<=ish over shells
       Compute or get la, lb, P, PA, PB
       TCART=0
      For iprim=0..nprim-1
         CALL K.PrimitiveIntegrals(IPRIM)
         Sum ICART += da*db*IPRIM
       end
       Tranform ICART to Spherical harm. ILM
         CALL C.IntegralConsumer(ILM)
    End jsh
  End ish
END
```

Only one such loop (and one for general contraction) in the entire program package! Covers ALL one-electron integrals

Example of an Integral Kernel

```
// Calculation of the overlap integral
void TOverlapIntegral::CalcPrimitiveIntegrals(int64 la, int64 lb, double a, double b,
                                               double R2, double *P, double *PA, double *PB,
                                               Tensor<5, double> &ET, Tensor<5, double> &RT,
                                               Tensor<3, double> &INTS)
{
 auto E = ET.SubTensor<4>({0}, {0});
 // Compute the constant SAB
 double ab = a+b;
 double ABI = 0.5/ab;
  double SAB = pow(2.0 * M PI * ABI, 1.5) * exp(-2.0 * ABI * a * b * R2);
  // Compute the auxiliary array E
 E Function(la, lb, ABI, SAB, PA, PB, E);
  // Loop over Cartesian components and make S
  int64 dimi= CDIM(la);
  int64 dimj= CDIM(lb);
  for (int64 i = 0; i < dimi; i++) {</pre>
    int64 l1 = GTO xyz[la][i][0];
    int64 m1 = GTO xyz[la][i][1];
    int64 n1 = GTO xyz[la][i][2];
    for (int64 j = 0; j < dimj; j++) {</pre>
      int64 l2 = GTO_xyz[lb][j][0];
      int64 m2 = GTO_xyz[lb][j][1];
      int64 n2 = GTO xyz[lb][j][2];
      INTS(0, i, j) = E(0,11,12,0) * E(1,m1,m2,0) * E(2,n1,n2,0);
    }; // j
  }; // i
};
```

Example of an Integral Consumer

```
/* ______
 * FN 03/2021
 *
 * For most integrals the right action is to just
 * simply store it in a matrix.
 * _____
 * /
class TSymOneElectronIntegralStorer
                   : public TOneElectronConsumer{
 int64 KernelLength;
 TSharkBasis *BAS;
 TRMatrixSym *IOUT;
public:
 TSymOneElectronIntegralStorer() {
   KernelLength=0;
   BAS=0;
   IOUT=0;
 };
 void SetKernelLength(int64 x){
   KernelLength=x;
 void SetBasis(TSharkBasis *x){
   BAS=x;
  }
 void SetOutput(TRMatrixSym *x){
   IOUT=x;
  }
 virtual void DigestIntegrals(int64 ish, int64 jsh,Tensor<3> &INTS);
};
```

```
void TSymOneElectronIntegralStorer::DigestIntegrals
                                  (int64 ish, int64 jsh,Tensor<3> &INTS)
{
  int64 li = BAS->BG[ish].l;
  int64 lj = BAS->BG[jsh].l;
  int64 ofsi= BAS->BG[ish].ofs;
  int64 ofsj= BAS->BG[jsh].ofs;
  int64 dimi= SHARK LDIM(li);
  int64 dimj= SHARK LDIM(lj);
  for (int64 i=0;i<dimi; i++){</pre>
    int64 addri= ofsi+i;
    int64 jend= dimj;
    if (Diagonal) jend=i+1;
    for (int64 j=0;j<jend; j++){</pre>
      int64 addrj= ofsj+j;
      for (int64 nk=0; nk<KernelLength; nk++){</pre>
        IOUT[nk](addri, addrj) = INTS(nk, i, j);
      }
    }
  }
}
```

PART 4

Automatic Code Generation

Problems with Method Development



Conclusions:

- The technicalities of development occupy most of our time
- Humans make mistakes, Debugging takes a lot of time
- The human brain can only deal with so much complexity. Beyond it is hopeless
- We need programming tools that take us directly from the Ansatz (our idea) to efficient, production level code
- Automatic Code Generation

Code Generation Tools

- ✓ Janssen & Schaefer, ROCCSD, pioneering work 1991
- ✓ Tensor contraction engine in NWCHEM, various CC (Hirata, Auer & Co)
- ✓ Diagram based arbitrary order CC/MRCC (Kallay)
- ✓ Gecco Internally contracted MRCC (Köhn)
- ✓ Genetic algorithm based code generator, MRCC (Hanrath)
- ✓ Automatic code generator, FIC-MRCI (Knizia, Werner)
- ✓ MREOM-CC (Huntington, Nooijen)
- ✓ General active space EOM CC (Kong, Demel, Shamsundar, Nooijen)
- ✓ Bagel/Smith CASPT2 gradient, (Shiozaki)
- ✓ Yanai, Saitow, DMRG-CASPT2, various contracted variants
- ✓ ACES III programming ,super-language' (Deumens, Bartlett & Co)
- ✓ Cyclops (Solomonik)
- ✓ Tiled Arrays (Valeev)
- ✓ many others

Simple & Straightforward Equation Generation

Any Ansatz (single- or multi-reference) that can be formulated in terms of 2nd quantization, quickly leads to expectation values of the form

$$\Big\langle \Psi_{_0} \, | \, E^n_{_m} E^q_{_p} ... E^s_{_r} \, | \, \Psi_{_0} \Big
angle, \qquad \quad E^q_{_p} = a^+_{_q eta} a_{_p eta} + a^+_{_q lpha} a_{_p lpha}.$$

Or, in terms of elementary spin-orbital operators:

$$\left\langle \Psi_{_{0}}\mid a_{_{m}}^{^{n}}a_{_{p}}^{^{q}}...a_{_{r}}^{^{s}}\mid \Psi_{_{0}}
ight
angle ,$$

If the orbital space is divided in internal (i,j,k,l), active (t,u,v,w) and virtual (a,b,c,d), the important commutation relations apply:

$$\begin{split} \left[E_{p}^{q}, E_{r}^{s}\right] &= E_{p}^{s} \delta_{qr} - E_{r}^{q} \delta_{ps}, \\ E_{i}^{p} \left|\Psi_{0}\right\rangle &= 2\delta_{ip} \left|\Psi_{0}\right\rangle, \quad \left\langle\Psi_{0}\right| E_{p}^{i} = 2\delta_{ip} \left\langle\Psi_{0}\right|, \\ E_{a}^{p} \left|\Psi_{0}\right\rangle &= 0, \quad \left\langle\Psi_{0}\right| E_{p}^{a} = 0, \end{split}$$

Thus:

Equation Generation

Strategy:

- ✓ Use the commutation relation to change the order of operators
- ✓ Move lower internal labels to the right
- ✓ Move upper internal labels to the left
- ✓ Move lower external labels to the right
- ✓ Move upper external labels to the left
- Creates 0's, Kronecker deltas and ,pre-densities' (MR case)

$$\gamma_{tv...x}^{uw...y} = \left\langle \Psi_0 \mid E_t^u E_v^w \dots E_x^y \mid \Psi_0 \right\rangle.$$

Issues: \checkmark redundant terms are generated

- ✓ terms that cancel each other are generated
- Equivalent terms may have inequivalent labels

√ ...

Postprocessing required

Awkward by hand, easy for a computer

Code Generation Chain

1. Equation Generator:

- ✓ Takes the Ansatz and generates equations
- ✓ Identifies identical, redundant and cancelling terms
- ✓ brings all labels into a ,canonical form'

2. Factorizer

- ✓ Identifies possible intermediates
- ✓ Finds the best possible intermediates and contraction order
- ✓ Finds common intermediates in different terms
- Ensures that all terms have their correct formal scaling

3. Code generator

- ✓ Writes code for a specific electronic structure package
- Recognizes patterns/contractions for which highly optimized code exists
- Ensures that all terms have their correct formal scaling
- Ensures minimal I/O and maximal use of BLAS
- ✓ Generates parallel code, code for specific machines,

Realization of a Code generation chain (AGE)



Cost model

In order to find the best possible intermediates and factorization, we need to have a prediction how long each contraction should take.



Generated vs Hand Written Code



Where does the hand written Code win?

Hand code:

$$\sigma^{ij} \leftarrow \sum_{k} \left[\left(2\mathbf{C}^{ik} - \left(\mathbf{C}^{ik}\right)^{\dagger} \right) \left(\mathbf{K}^{kj} - \frac{1}{2} \mathbf{J}^{kj} \right) - \frac{1}{2} \left(\left(\mathbf{C}^{ik} \right)^{\dagger} \mathbf{J}^{kj} \right) - \left(\left(\mathbf{C}^{ik} \right)^{\dagger} \mathbf{J}^{kj} \right)^{\dagger} + \left(\mathbf{K}^{ik} - \frac{1}{2} \mathbf{J}^{ik} \right) \left(2\mathbf{C}^{kj} - \left(\mathbf{C}^{kj} \right)^{\dagger} \right) - \frac{1}{2} \left(\mathbf{J}^{ik} \left(\mathbf{C}^{kj} \right)^{\dagger} \right) - \left(\mathbf{J}^{ik} \left(\mathbf{C}^{kj} \right)^{\dagger} \right) + \right] \right] \cdot \mathbf{4} \operatorname{dgemm/k}$$

Generated code: $\sigma^{ij} \leftarrow \sum_{k} -\mathbf{J}^{ik} \mathbf{C}^{kj} - \mathbf{C}^{kj} \mathbf{J}^{ik} - \mathbf{J}^{kj} \mathbf{C}^{ik} - \mathbf{C}^{ik} \mathbf{J}^{kj} - \mathbf{C}^{ki} \mathbf{K}^{kj} - \mathbf{K}^{ik} \mathbf{C}^{jk} + 2\mathbf{K}^{ik} \mathbf{C}^{jk} + 2\mathbf{C}^{ik} \mathbf{K}^{kj} - \mathbf{K}^{kj} \mathbf{K}^{kj} \mathbf{K}^{kj} - \mathbf{K}^{kj} \mathbf{K}^{kj} \mathbf{K}^{kj} - \mathbf{K}^{kj} \mathbf{$

Coupled Cluster Gradients



- Canonical Coupled Cluster gradients with perhaps >500 basis functions possible
- Parallel Scaling is good
- More than 10x faster than numerical gradients

Reduced Scaling FIC-MRCC Implementation



- Accessible molecular size roughly the same as single reference CCSD
- Without reduced scaling limit about 8 active orbitals
- ➡ With reduced scaling limit about 12 active orbitals





Complexity: Example

Fully internal contracted MRCI (or MRCC, also CASPT2/NEVPT2) works with contracted functions in the first-order interacting space (FOIS)

$$\left|\Phi_{ij}^{ta}\right\rangle = E_{ij}^{ta}\left|\Psi_{0}\right\rangle = \sum_{I} C_{I}^{(CASSCF)} E_{ij}^{ta}\left|\Phi_{I}^{(CAS)}\right\rangle$$

- ✓ 10 Excitation classes -> 100 Blocks of matrix elements
- Not orthogonal
- ✓ Not linearly independent
- Extremely complicated matrix elements
- ➡ 1945 equations including up to four body density
- Factorized into 3674 equations
- Removed 1222 redundant intermediates
- Nearly hopeless to program by hand. Readily done with code generator as a matter of hours (perhaps days)

MQM 2016

Influence of the choice of projection manifolds in the CASPT2 implementation

Takeshi Yanai S, Yuki Kurashige, Masaaki Saitow, Jakub Chalupský, Roland Lindh & Per-Åke Malmqvist Pages 2077-2085 | Received 12 Oct 2016, Accepted 01 Dec 2016, Published online: 27 Dec 2016

... found a (small) bug in the hand coded version of the CASPT2 method

THE JOURNAL OF CHEMICAL PHYSICS 142, 051103 (2015)



Communication: Automatic code generation enables nuclear gradient computations for fully internally contracted multireference theory

Matthew K. MacLeod and Toru Shiozaki Department of Chemistry, Northwestern University, 2145 Sheridan Rd., Evanston, Illinois 60208, USA

(Received 11 January 2015; accepted 27 January 2015; published online 5 February 2015)

... Fully automated, large scale nuclear gradient for CASPT2. Optimizations of metalloporphyrins

Analytical Gradient Theory for Strongly Contracted (SC) and Partially Contracted (PC) N-Electron Valence State Perturbation Theory (NEVPT2)

Jae Woo Park

J. Chem. Theory Comput., Just Accepted Manuscript • DOI: 10.1021/acs.jctc.9b00762 • Publication Date (Web): 10 Sep 2019

Code generation: Summary

- ✓ Code generation enables the implementation of ,impossibly complicated' methods
- Code generation reduces development times from years to hours/days
- Code generation can produce code for specific hardware, thus ensuring optimal performance
- Code generation can ensure that all parts of the code have consistent quality
- ✓ Once the code generation chain produces correct results, it is extremely reliable (e.g. a small bug was identified in the original CASPT2 code in 2015, CASPT2 is from 1990!)
- Code generation will play an important part in future quantum chemistry
- Generated code can be made almost as efficient as the best hand optimized code
- In the future we keep just a wavefunction Ansatz in the source code repository and generate the code during compile time. Any improvement in the code generation chain is the immediately applied to all parts of the program.