

# A L<sup>A</sup>T<sub>E</sub>X template for CPC Computer Physics Descriptions

First Author<sup>a,\*</sup>, Second Author<sup>a,b</sup>, Third Author<sup>b</sup>

<sup>a</sup>*First Address*

<sup>b</sup>*Second Address*

---

## Abstract

A submitted program is expected to be of benefit to other physicists or physical chemists, or be an exemplar of good programming practice, or illustrate new or novel programming techniques which are of importance to some branch of computational physics or physical chemistry.

Acceptable program descriptions can take different forms. The following Long Write-Up structure is a suggested structure but it is not obligatory. Actual structure will depend on the length of the program, the extent to which the algorithms or software have already been described in literature, and the detail provided in the user manual.

Your manuscript and figure sources should be submitted through the Elsevier Editorial System (EES) by using the online submission tool at <http://www.ees.elsevier.com/cpc>.

In addition to the manuscript you must supply: the program source code; job control scripts, where applicable; a README file giving the names and a brief description of all the files that make up the package and clear instructions on the installation and execution of the program; sample input and output data for at least one comprehensive test run; and, where appropriate, a user manual. These should be sent, via email as a compressed archive file, to the CPC Program Librarian at [cpc@qub.ac.uk](mailto:cpc@qub.ac.uk).

*Keywords:* keyword1; keyword2; keyword3; etc.

---

## PROGRAM SUMMARY/NEW VERSION PROGRAM SUMMARY

---

\*Corresponding author.

*E-mail address:* firstAuthor@somewhere.edu

*Program Title:*

*Licensing provisions(please choose one): CC0 1.0/CC By 4.0/MIT/Apache-2.0/BSD 3-clause/BSD 2-clause/GPLv3/CC BY NC 3.0*

*Programming language:*

*Supplementary material:*

*Journal reference of previous version:*

*Does the new version supersede the previous version?:*

*Reasons for the new version:*

*Summary of revisions:\**

*Nature of problem(approx. 50-250 words):*

*Solution method(approx. 50-250 words):*

*Additional comments including Restrictions and Unusual features (approx. 50-250 words):*

[1] Reference 1

[2] Reference 2

[3] Reference 3

\* Items marked with an asterisk are only required for new versions of programs previously published in the CPC Program Library.

## 1. GPU Architecture

While general purpose computing on graphics processing units (GPGPU) (GPUs) has been adopted by the high performance computing (HPC) community for quite some time, it can seem quite complex to the uninitiated. While most quantum chemists who decide to dip their toes into computational waters can get away with having little to no understanding of what is actually going on under the hood when programming for a central processing unit (CPU), the same can not at all be said for GPUs. Therefore, in this section I will explain what makes a GPU tick in order to help understand the terms and techniques used in the following chapter. As CUDAProphet was optimized for (and is currently hardcoded for :( ) the Tesla C2050, I will be referring to its specs for examples when needed. I will begin with the most

granularity possible, and zoom out so that by the end of this section, the reader should come away prepared for the rest of this chapter.

### 1.1. *Threads, Blocks, and Grids*

The most granular element of computation on a GPU is the thread. When a CUDA function or subroutine (hereafter called a kernel) is called, it is a thread that actually executes the code. What makes GPGPU so powerful is that while a thread in a kernel is always executing the same code, the data a thread works with can be completely different between threads. This method of parallelism is called single instruction, multiple data (SIMD).

32 threads are organized into a structure called a warp. Within a warp, all threads execute code in lock-step. If a condition arises where some threads in warp must execute some code, and other threads in the same warp execute some other code (for instance in a **if** **elses** statement), this leads to serialization of code execution, a process called warp divergence. In a worst case scenario, this could cause all 32 threads to execute serially which could lead to a massive hit to performance. It can also cause warps to fall behind other warps and has the potential to cause race conditions to appear. If one warp need data generated by another warp that is several steps behind, this can cause all kinds of confusing errors to appear. Therefore, it is generally advisable to avoid warp divergence whenever possible. But, for programs of any considerable usefulness, warp divergence will be inevitable. Thankfully, the good people of Nvidia have included the **syncthreads** command. This provides a barrier that all threads within a block must meet before any can continue.

The next highest structure is the thread block, or block for short. Simply put, a block is a collection of 1 or more threads. When a kernel is called, the block scheduler assigns each block to a streaming multiprocessor (SM). How many blocks a SM can run at once depends on the resources each block uses (more on this in section ??), but no more than 8 block per SM can ever be run at once on our C2050. Blocks also add an element of dimensionality to threads through the thread index (*threadIdx*) variable. This variable has three parts, *threadIdx.x*, *threadIdx.y*, and *threadIdx.z*. These variables make dealing with matrixes and arrays much simpler and can help with paralyzing large nested loops. The size of dimensions of the grid can be accessed through the *blockDim* variable. *blockDim* also has *x*, *y*, and *z* parts. These variables are all integer type, and the *threadIdx* variables range from 1 in Fortran (or 0 in C/++) to the relevant part of *blockDim* in

Fortran (or *blockDim* -1 in C/++). The *threadIdx* set of variables are all reset at the boundary of each block.

Finally there is the grid. The grid is the entire collection of blocks that are launched for a kernel. They also give a dimensionality to blocks through the *blockIdx* variable which has similar parts to the *threadIdx* variable. It also adds in the *gridDim* variable which is analogous to *blockDim*.

## 2. Program flow

In this section, the control flow of a typical calculation is given.

*Step 1.* The input file is read by the *intin* subroutine. If needed, the basis set and open-shell configurations are calculated by *formbs* and *find\_bin\_configurations* respectively. The options set in the input file are rewritten to stdout.

*Step 2.* The calculation of small arrays and other constants is performed by *calc\_parameters* and *bsnorm*. *calc\_parameters* calls *bcoef* and *setvc*.

*Step 3.* The mapping of threads to the unrolled one and two-electron integral matrixes are calculated by *lmpqrsa* on the GPU. All other data from step 1 and 2 is then uploaded to the GPU.

*Step 4.* The one and two electron integrals are calculated on the GPU by *eint1gpu* and *eint2gpu* respectively.

*Step 5.* The initial guess of the density matrix is calculated by *guess* on the GPU. This is done by the diagonalization of the one-electron Hamiltonian matrix.

*Step 6.* *scfiter* then performs the SCF until convergence of the density matrix has been reached, or the maximum number of iterations has been reached. SCF is performed with cuSOLVER functions, as well as a few custom helper functions. The converged eigenvectors, values, and energies are downloaded from the GPU and then written to stdout.

*Step 7 (optional).* If *jobtyp*='bsopt', then an optional basis set optimization is then carried out. This starts by assigning pointers to variables on the CPU and GPU through *hookup\_cpu* and *hookup\_gpu*. Then, the four wtbs parameters are optimized by *newuoa*. *newuoa* calls *calc\_energy* which essentially reconstructs the basis set from new parameters from *newuoa*, then repeats steps 2 - 6 and feeds the energy back into *newuoa*. This repeats until optimal wtbs parameters have been found.

### 3. Alterations for CUDA

While almost all of the code from the original DFRATOM was modified in some way, the most extreme changes were the integral evaluation and the formation of the P and Q matrixes. Therefore, we will discuss these changes in detail in this section.

The original DFRATOM calculates the two electron integrals with a long series of nested for loops. Working from the outside in, the indices of the loops are  $L$ ,  $P$ ,  $Q$ ,  $M$ ,  $R$ , and  $S$  where  $L$  and  $M$  are spinor symmetries,  $P$  and  $Q$  are basis functions of symmetry  $L$ , and  $R$  and  $S$  are basis functions of symmetry  $M$ . The pseudocode for these is shown in Algorithm 1. Where  $n_{sym}$  is the total number of symmetries, and  $nbs(i)$  is the number of basis functions for symmetry species  $i$ . Thus, each set of J and K integrals are uniquely defined by their  $L$ ,  $M$ ,  $P$ ,  $Q$ ,  $R$ , and  $S$  values. In order to have an effective CUDA implementation of this algorithm, we to find a way to map these six numbers to CUDA threads. The simplest approach would be to map the values of  $L$ ,  $P$ , and  $Q$  onto the  $threadIdx.x$ ,  $threadIdx.y$ , and  $threadIdx.z$  variables for each thread, launch the needed number of blocks, and then have each thread loop over the remaining indices. The pseudocode for this can be seen in Algorithm 2. A similar method is employed by many other programs, and for larger systems it works perfectly well. But because this program is for only single atoms, problems begin to appear.

The first is the problem of warp divergence. Because the maximum values of  $M$ ,  $R$ , and  $S$  depend on  $L$ ,  $P$ , and  $Q$ , different threads will have a different number of loops to complete than others. Because a streaming multiprocessor (SM) must finish the block it is currently working on before it can grab another, there is the possibility that most of the threads in a block are idling while waiting for others in the same block to finish. The second problem is that with this method, there will always be several blocks that have threads that remain idle throughout the block's runtime, no matter what. This can be seen more clearly in [Figure whatever](#). The third problem with this is that with the maximum values for  $L$ ,  $P$ , and  $Q$  available for a single atom, there might not even be enough combinations to completely fill the GPU. While all of these issues begin to disappear once the number of integrals to evaluate becomes large enough, we are still very much in the range where they are in play. Therefore a smarter algorithm had to be used.

The second problem can be solved by opting for a one dimensional solution instead of the three dimensional one of Algorithm 2. By restricting ourselves

---

**Algorithm 1** The original

---

```
for  $L = 1$  to  $nsym$  do
  for  $P = 1$  to  $nbs(L)$  do
    for  $Q = 1$  to  $P$  do
      for  $M = 1$  to  $L$  do
        if  $L = M$  then
           $maxr = P$ 
        else
           $maxr = nbs(M)$ 
        end if
        for  $R = 1$  to  $maxr$  do
          if  $(L = M)$  and  $(P = R)$  then
             $maxs = Q$ 
          else
             $maxs = R$ 
          end if
          for  $S = 1$  to  $maxs$  do
            compute the J and K integrals of  $L, M, P, Q, R,$  and  $S$ 
          end for
        end for
      end for
    end for
  end for
end for
```

---

---

**Algorithm 2** Easy Code

---

$L = threadIdx.x + (blockIdx.x - 1) * blockDim.x$

$P = threadIdx.y + (blockIdx.y - 1) * blockDim.y$

$Q = threadIdx.z + (blockIdx.z - 1) * blockDim.z$

**if** ( $L \leq nsym$ ) **and** ( $P \leq nbs(L)$ ) **and** ( $Q \leq P$ ) **then**

**for**  $M = 1$  to  $L$  **do**

**if**  $L = M$  **then**

$maxr = P$

**else**

$maxr = nbs(M)$

**end if**

**for**  $R = 1$  to  $maxr$  **do**

**if** ( $L = M$ ) **and** ( $P = R$ ) **then**

$maxs = Q$

**else**

$maxs = R$

**end if**

**for**  $S = 1$  to  $maxs$  **do**

    compute the J and K integrals of  $L$ ,  $M$ ,  $P$ ,  $Q$ ,  $R$ , and  $S$

**end for**

**end for**

**end for**

**end if**

---

to only using `threadIdx.x`, we ensure that only the last block to run will have the possibility of having threads remaining idle throughout the block's runtime. The third problem can be solved by having each thread calculate one and only one set of J and K integrals. If we start with a valid combination of  $L$ ,  $M$ ,  $P$ ,  $Q$ ,  $R$ , and  $S$  values, we can very easily figure out which thread will calculate that set of integrals by using the following equations.

$$n'(j) = \frac{j^2 + j}{2} \quad (1)$$

$$y = n'(nbs(L - 1)) + n'(P - 1) + Q \quad (2)$$

$$x = n'(nbs(M - 1)) + n'(R - 1) + S \quad (3)$$

$$i = n'(x) + y \quad (4)$$

$$i_{max} = n'(\sum_{j=1}^{nsym} n'(nbs(j))) \quad (5)$$

$nbs(0) = 0$  and  $i$  would be equal to  $threadIdx.x + (blockIdx.x - 1) * blockDim.x$ . Starting with a value of  $i$  and working our way back though is a much more challenging task. It becomes easier if we reframe it in the following way.

Consider [Figure whatever](#). It shows the top half of the symmetric two-electron integral matrix for a problem with  $nsym = 2$  and 3 basis function for symmetry one, and two for symmetry two. In each element of this matrix, there is a set of seven numbers. The top two are  $L$  and  $M$ , then  $P$  and  $Q$ , then  $R$  and  $S$ , and the last number is the value of  $i$  for the thread calculating that integral. With this, it can be seen that each element in the same row have the same  $M$ ,  $R$ , and  $S$  values and the elements in the same column have the same  $L$ ,  $P$ , and  $Q$  values. Therefore, finding out which column the element belongs to gives us our value for  $y$ , and finding the row give us the value of  $x$ . This can be done with the following binary search algorithm.

Where variables with the  $s_$  prefix refer to those in the shared memory, and  $lownum$  and  $highnum$  refer to the minimum and values the  $i$  could be for the current guess ( $mid$ ) of  $y$ . From here,  $L$  can be found with Algorithm 4, and  $P$  and  $Q$  can be found with Algorithm 5. The same set of algorithms



---

**Algorithm 3** Binary Search for  $x$  and  $y$ 

---

```
if  $theadIdx.x \leq nsym$  then
     $s\_nsym = nsym$ 
     $s\_nbs(theadIdx.x) = nbs(theadIdx.x)$ 
     $s\_nprime(theadIdx.x) = n'(s\_nbs(theadIdx.x))$ 
end if
call syncthreads
 $i = threadIdx.x + (blockIdx.x - 1) * blockDim.x$ 
if  $i \leq i_{max}$  then
     $low = 1$ 
     $high = \text{sum}(s\_nprime(1 : s\_nsym))$ 
    while  $low \leq high$  do
         $mid = \frac{(low+high)}{2}$ 
         $lownum = \frac{(mid-1)(mid-2)}{2} + mid$ 
         $highnum = lownum - 1 + mid$ 
        if  $(i \leq highnum)$  and  $(i \geq lownum)$  then
             $y = mid$ 
            exit
        else if  $i > highnum$  then
             $low = mid + 1$ 
        else if  $i < lownum$  then
             $high = mid - 1$ 
        end if
    end while
     $x = i - lownum + 1$ 
end if
```

---

---

**Algorithm 4** Binary Search for  $L$ 

---

$i = threadIdx.x + (blockIdx.x - 1) * blockDim.x$

**if**  $i \leq i_{max}$  **then**

$low = 1$

$high = s_{nsym}$

**while**  $low \leq high$  **do**

$mid = \frac{(low+high)}{2}$

$lownum = 1 + \text{sum}(s_{nprime}(1 : mid - 1))$

$highnum = lownum - 1 + s_{nprime}(mid)$

**if**  $(i \leq highnum)$  **and**  $(i \geq lownum)$  **then**

$L = mid$

**exit**

**else if**  $y > highnum$  **then**

$low = mid + 1$

**else if**  $y < lownum$  **then**

$high = mid - 1$

**end if**

**end while**

**end if**

---

---

**Algorithm 5** Binary Search for  $P$  and  $Q$ 

---

```
 $i = \text{threadIdx}.x + (\text{blockIdx}.x - 1) * \text{blockDim}.x$   
if  $i \leq i_{\max}$  then  
   $low = 1$   
   $high = s\_nbs(L)$   
  while  $low \leq high$  do  
     $mid = \frac{(low+high)}{2}$   
     $lownum = \frac{(mid-1)(mid-2)}{2} + mid + \text{sum}(s\_nprime(1 : L - 1))$   
     $highnum = lownum + mid - 1$   
    if  $(i \leq highnum)$  and  $(i \geq lownum)$  then  
       $P = mid$   
      exit  
    else if  $y > highnum$  then  
       $low = mid + 1$   
    else if  $y < lownum$  then  
       $high = mid - 1$   
    end if  
  end while  
   $Q = y - lownum + 1$   
end if
```

---

can then be used to get  $M$ ,  $R$ , and  $S$  by substituting the relevant variables. If there is sufficient global memory available, these values can be stored and referred to later as needed. Otherwise, they could be calculated on the fly as needed. Because binary search scales as  $\mathcal{O}(n \log n)$ , this should ensure that this remains a fast method of mapping threads to integrals for large problems as well. **(CHECK WITH MARIUSZ TO BE SURE THE NEXT IS TRUE!!!!!!)** With some alterations, this method could also apply to molecular symmetries other than a single atom. For instance in C1, all possible combinations of four basis functions must be used (ignoring those that appear on the bottom triangle of the two electron integral matrix of course). We could simply remove the search for  $L$  and  $M$ , have the initial value of *high* in Algorithm 5 be the total number of basis functions, and remove the `sum(s_nprime(1 : L - 1))` term from *lownum*.

From here, the code for actually evaluating the integrals remains largely the same as the original code, except for some minor changes to allow for more efficient global or shared memory access. We also use a process referred to as "grid-stride looping" where all these binary search algorithms have their `if  $i \leq i_{max}$  then` removed, and then are placed within the following loop: `for  $i = threadIdx.x + (blockIdx.x - 1) * blockDim.x$  to  $i_{max}$ ,  $i += blockDim.x * gridDim.x$  do`. If we know the occupancy of the algorithm on the GPU beforehand, we can launch exactly the number of blocks that will fill the GPU. This reduces the overhead of block swapping and lets us further eke out some performance.

#### 4. Input Description

The program can be excited on Unix-like systems in the following way

```
$ path_to_executable input_file > output_file
```

The input file must end in ".inp" or an error will be given. Redirection of stdout to an output file is optional, but is recommended to save the results of a calculation. The input file is read using the namelist functionality of Fortran. A description of what must appear on each line of the input file is given below. Sample input files are also given at the end of this document.

1. A title of no more than 200 characters.
2. \$contrl

jobtype		The type of calculation to be performed.
	= 'energy'	Will do a single point energy calculation.
	= 'bsopt'	Will optimize the basis set. Can only be used if bastype equals 'wtbs'.
c		The speed of light. If not given, the default is set to 137.03599976 au.

### 3. \$nuc

znuc		The charge of the nucleus.
nucmdl		The nuclear model to use.
	= 1	Point nucleus (default).
	= 2	Finite sphere (not yet supported).
	= 3	Gaussian.
rnuc		The radius of the finite sphere nucleus.
alpha		The exponent for the gaussian nucleus.
		Defaults are given in litdata.f90.

### 4. \$bas

nsym		The number of symmetries to be used.
bastype		The type of basis set given.
	= 'wtbs'	Use a wtbs.
	= 'rdin'	Read in the basis set from the input file.
ngroup		The number of different groups to use in the wtbs scheme (default 1).

The next line will depend on what bastype was set to. If bastype equals 'rdin' then the following lines must be the number of functions for the S+ symmetry, followed by the exponents to use, each on a new line. The pattern repeats for each new symmetry. See the sample input files for further clarification. Otherwise, if bastype equals 'wtbs' the \$wtbs group is read next.

### 5. \$wtbs has to be given if bastype='wtbs'.

wtbspara	The $\alpha$ , $\beta$ , $\delta$ , and $\gamma$ wtbs parameters. If there is more than one group, the order would be $\alpha_1, \beta_1, \delta_1, \gamma_1, \alpha_2, \beta_2, \delta_2, \gamma_2$ and so on.
nbs	The number of functions used in each symmetry.
start	Where in the $\zeta$ pool each symmetry starts taking exponents from (default=1,1,1,1,1,1).
groups	What group each symmetry belongs to (default=1,1,1,1,1,1).

The next line depends on what the jobtype was set to. If jobtype equals 'energy' \$newuoa is skipped and \$seconfi will be read next. If jobtype equals 'bsopt', then the \$newuoa group will be needed.

6. \$newuoa has to be given if jobtype='bsopt'. Refer to the newuoa documentation for more information if needed.

rhobeg	The initial value of the trust region used by newuoa (default=0.1).
rhoend	The final value of the trust region used by newuoa. Must be smaller than rhobeg (default= $1.0 \times 10^{-4}$ ).
iprint	The print level for newuoa.
= 0	No printing from newuoa (default).
= 1	Print only when newuoa has finished.
= 2	Print only when the trust region has decreased by an order of magnitude.
= 3	Print every iteration of newuoa.
maxfun	The maximum number of calls to calfun newuoa will make before terminating (default=500).

## 7. \$seconfi

nclose	The number of closed spinors for each symmetry.
nopen	The number of open spinors for each symmetry. There is a limit to one open orbital per symmetry.
freeel	The number of electrons available in the open spinors.
autogen	Automatically generates all possible combinations of spinor occupancies (default=.false.).
nconf	The number of configurations to be read in (needed if autogen is false).

If autogen is false, then the next nconf lines will be the spinor occupancies. They will be given as real numbers with one configuration per line.

8. \$scf

maxitr	The maximum number of SCF iterations (default=50).
ixtrp	The method of extrapolation.
= 0	No extrapolation (default).
= 1	Extrapolate the Fock matrix.
dfctr	Damping factor for Fock maxtrix (default=0.3).
thdll	Convergence limit for the large-large components of the density matrix (default= $1.0 \times 10^{-5}$ )
thdsl	Convergence limit for the small-large components of the density matrix (default= $1.0 \times 10^{-7}$ )
thdss	Convergence limit for the small-small components of the density matrix (default= $1.0 \times 10^{-9}$ )