

29-03-2021



D2.1 – Report on a first alpha release of the I/O library, ready for WP4

Version V1.0

GA no 952165

Dissemination Level

- PU: Public
- PP: Restricted to other programme participants (including the Commission)
- RE: Restricted to a group specified by the consortium (including the Commission)
- CO: Confidential, only for members of the consortium (including the Commission)

Document Information

Project Title	Targeting Real Chemical accuracy at the EXascale
Project Acronym	TREX
Grant Agreement No	952165
Instrument	Call: H2020-INFRAEDI-2019-1
Topic	INFRAEDI-05-2020 Centres of Excellence in EXascale computing
Start Date of Project	01-10-2020
Duration of Project	36 Months
Project Website	https://trex-coe.eu/
Deliverable Number	D2.1
Deliverable title	D2.1 – Report on a first alpha release of the I/O library, ready for WP4, as seen in GA
Due Date	M06 – 31-03-2021 (from GA)
Actual Submission Date	29-03-2021
Work Package	WP2 – Code modularization and interfacing
Lead Author (Org)	Sandro Sorella (Scuola Internazionale Superiore di Studi Avanzati di trieste (SISSA))
Contributing Author(s) (Org)	Anthony Scemama (Centre National de la Recherche Scientifique (CNRS))
Reviewers (Org)	François Coppens (Université de Versailles Saint-Quentin-en-yvelines (UVSQ)), Pablo Oliveira (UVSQ), Matúš Dubecký (Slovenská technická univerzita v Bratislave (STUBA)), Jan Beerens (Universiteit Twente (UT))
Version	V1.0
Dissemination level	PU
Nature	Report
Draft / final	Final
No. of pages including cover	29



Disclaimer



TREX: Targeting Real Chemical Accuracy at the Exascale project has received funding from the European Union Horizon 2020 research and innovation program under Grant Agreement No. 952165.

The content of this document does not represent the opinion of the European Union, and the European Union is not responsible for any use that might be made of such content.



Versioning

Version	Date	Authors	Notes
1.0	29-03-2021	Sandro Sorella (SISSA)	First Official Release



Abbreviations

AO	Atomic Orbital
API	Application Programming Interface
BSE	Basis Set Exchange
CI	Configuration Interaction
CNRS	Centre National de la Recherche Scientifique
CoE	Center of Excellence
ECP	Effective Core Potential
ERI	Electron Repulsion Integral
EZFIO	Easy Fortran Input/Output
GPFS	General Parallel File System
HDF5	Hierarchical Data Format
HPC	High Performance Computing
I/O	Input/Output
MO	Molecular Orbital
QMC	Quantum Monte Carlo
SISSA	Scuola Internazionale Superiore di Studi Avanzati di trieste
STUBA	Slovenská technická univerzita v Bratislave
TREX	Targeting REal chemical accuracy at the eXascale
TREXIO	TREX Input/Output
UT	Universiteit Twente
UVSQ	Université de Versailles Saint-Quentin-en-yvelines



Table of Contents

Document Information	i
Disclaimer	ii
Versioning.....	iii
Abbreviations	iv
Table of Contents	v
List of Figures	vi
List of Tables	vii
1 Introduction.....	1
2 Content of the files	2
Metadata	2
Nuclei	2
Electrons	2
Atomic Basis set	2
Atomic Orbitals	3
Effective core potentials	3
One-electron integrals in the AO basis set	3
Two-electron integrals in the AO basis set	4
One-electron integrals in the MO basis set	4
Two-electron integrals in the MO basis set	4
Summary	4
3 Design of the library	6
License	6
The front end	6
Conventions.....	9
Error handling.....	10
Safety	10
Prototype library.....	10
The back end	10
4 Illustrative example of usage of the library	11
References	I



List of Figures

- 1 Dependencies between the codes and the data. Codes are represented in gray, and data are represented in white. 1



List of Tables

1	Information that should be stored in the file.	4
---	---	---



1 Introduction

We build a library to help inter-operability between codes in the field of quantum chemistry, primarily focused on enabling the communication of data between the flagship codes of the Targeting REal chemical accuracy at the eXascale (TREX) Center of Excellence (CoE) (NECI, GammCor, Quantum Package, QMC=Chem, CHAMP, TurboRVB, QML). We expect this library to be also adopted by the community beyond the TREX CoE.

The data that needs to be stored is the electronic wave function, which is obtained from a post-Hartree-Fock calculation, or the one- and two-body density matrices, together with the one- and two-electron integrals that are necessary to compute the energy or other properties. As a wave function can be obtained by executing multiple codes in a complex workflow, the library should give the possibility to build the files incrementally using multiple codes.

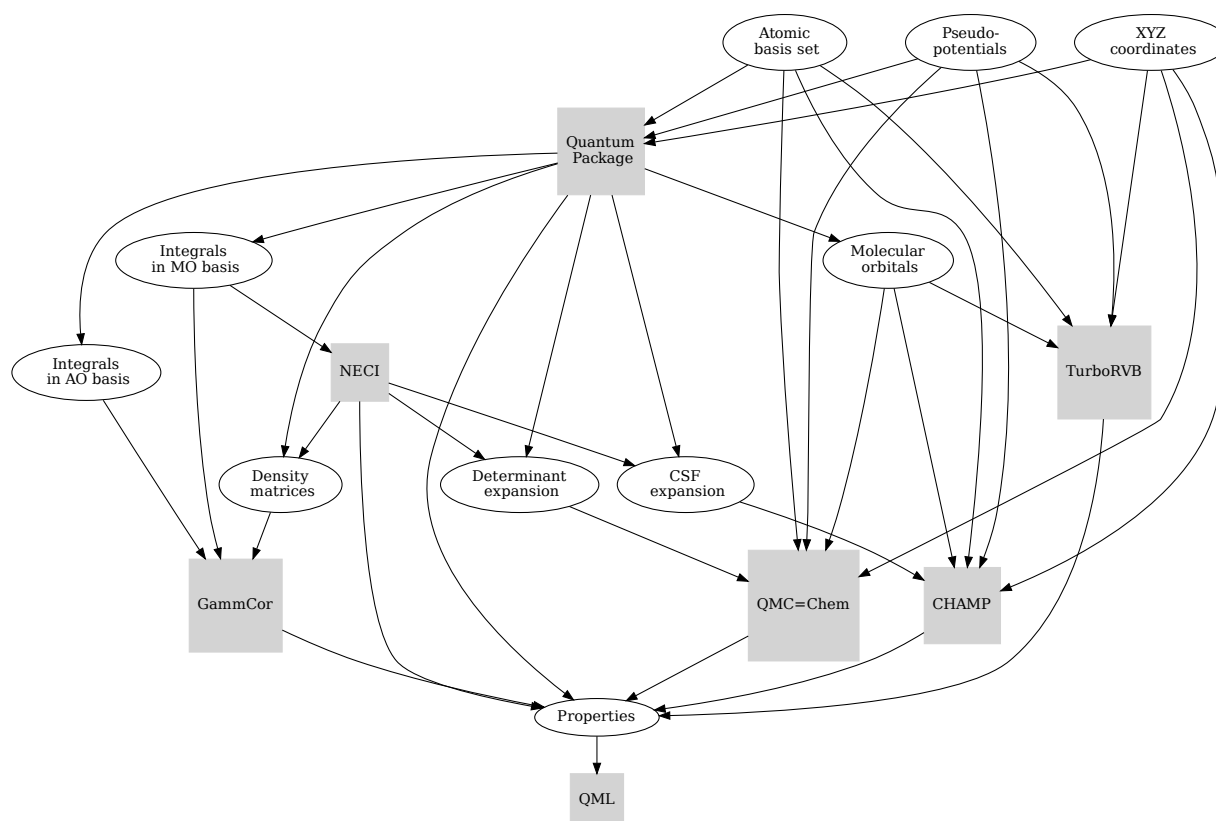


Figure 1: Dependencies between the codes and the data. Codes are represented in gray, and data are represented in white.

Fig. 1 shows which data are read and written by all of the codes of the TREX CoE. The objective of this library is to organize all the data in a file, and provide a common interface to make the data easily accessible to all the codes.

2 Content of the files

The files need to be self-contained: they should contain *all* the information needed to reconstruct the wave functions from an external program, without relying on any extra service to provide additional data. For instance, all the parameters of the atomic basis set should be explicitly stored instead of storing only the conventional name of the basis, which would require obtaining the parameters from another source. It is important for the files to be self-contained, as they are intended to be archived on an open-data repository: this reduces their dependencies on other services and therefore increases their re-usability.

Metadata

As we expect our files to be archived in open-data repositories, we need to give the possibility to the users to store some metadata inside the files. We propose to store the list of names of the codes which have participated to the creation of the file, a list of authors of the file, and a textual description.

Nuclei

We consider wave functions where the nuclei are considered as fixed point charges. The file should contain information describing the positions of the nuclei, their charges, their labels, and the point-group symmetry of the system.

Electrons

We consider wave functions expressed in the spin-free formalism, where the number of \uparrow and \downarrow electrons is fixed.

Atomic Basis set

We consider here basis functions centered on nuclei. Hence, we enable the possibility to define *dummy atoms* to place basis functions in random positions.

The atomic basis set is defined as a list of shells. Each shell s is centered on a center A , possesses a given angular momentum l and a radial function R_s . The radial function is a linear combination of N_{prim} *primitive* functions that can be of type Slater ($p = 1$) or Gaussian ($p = 2$), parameterized by exponents γ_{ks} and coefficients a_{ks} :

$$R_s(\mathbf{r}) = \mathcal{N}_s |\mathbf{r} - \mathbf{R}_A|^{n_s} \sum_{k=1}^{N_{\text{prim}}} a_{ks} \exp(-\gamma_{ks} |\mathbf{r} - \mathbf{R}_A|^p).$$

In the case of Gaussian functions, n_s is always zero. The normalization factor \mathcal{N}_s ensures that all the functions of the shell are normalized to unity. As this normalization requires the ability to compute overlap integrals, the normalization factors should be written in the file to ensure that the file is self-contained and does not require the client program to have the ability to compute such integrals.



Atomic Orbitals

Going from the atomic basis set to Atomic Orbitals (AOs) implies a systematic construction of all the angular functions of each shell. We consider two cases for the angular functions: the real-valued spherical harmonics, and the polynomials in Cartesian coordinates. In the case of spherical harmonics, the AOs are ordered in increasing magnetic quantum number ($-l \leq m \leq l$), and in the case of polynomials we choose the canonical ordering of the Libint[1] library, i.e

$$\begin{aligned} p &: p_x, p_y, p_z \\ d &: d_{xx}, d_{xy}, d_{xz}, d_{yy}, d_{yz}, d_{zz} \\ f &: f_{xxx}, f_{xxy}, f_{xxz}, f_{xyy}, f_{xyz}, f_{xzz}, f_{yyy}, f_{yyz}, f_{yzz}, f_{zzz} \end{aligned}$$

etc.

AOs are defined as

$$\chi_i(\mathbf{r}) = P_{\eta(i)}(\mathbf{r}) R_{\theta(i)}(\mathbf{r})$$

where i is the atomic orbital index, P encodes for either the polynomials or the spherical harmonics, $\theta(i)$ returns the shell on which the AO is expanded, and $\eta(i)$ denotes which angular function is chosen.

Effective core potentials

It is common to use Effective Core Potentials (ECPs) in Quantum Monte Carlo (QMC) calculations. An ECP V_A^{PP} replacing the core electrons of atom A is the sum of a local component V_A^{loc} and a non-local component $V_A^{\text{non-loc}}$. [2] The local component is given by

$$V_A^{\text{loc}}(r) = -\frac{Z_A^{\text{eff}}}{r} + \frac{Z_A^{\text{eff}}}{r} \exp(-\alpha_A r^2) + Z_{\text{eff}} \alpha_A r \exp(-\beta_A r^2) + \gamma_A \exp(-\delta_A r^2),$$

and the component obtained after localizing the non-local operator is

$$V_A^{\text{non-loc}}(r) = \zeta_A \exp(-\eta_A r^2) |0\rangle\langle 0| + \mu_A \exp(-\nu_A r^2) |1\rangle\langle 1|$$

where $r = |\mathbf{r} - \mathbf{R}_A|$ is the distance to the nucleus on which the potential is centered, Z_A^{eff} is the effective charge due to the removed electrons, $|0\rangle\langle 0|$ and $|1\rangle\langle 1|$ are projections over zero and one principal angular momenta, respectively (generalization to higher angular momenta is straightforward), and all the parameters labeled by Greek letters are parameters.

One-electron integrals in the AO basis set

The one-electron integrals are of the form

$$O_{ij} = \int \chi_i(\mathbf{r}) \hat{O} \chi_j(\mathbf{r}) d\mathbf{r},$$

where \hat{O} is a one-electron operator. The integrals needed to compute the energy are the overlap integrals \mathbf{S} with $\hat{O} = \mathbf{1}$, the kinetic energy integrals \mathbf{T} with $\hat{O} = \nabla^2$, the electron-nucleus potential integrals \mathbf{V} with $\hat{O} = -\sum_A -Z_A^{\text{eff}}/|\mathbf{r} - \mathbf{R}_A|$, and the effective core potential integrals \mathbf{V}^{PP} with $\hat{O} = V^{\text{PP}}$. It is also convenient to store the *core Hamiltonian* integrals, defined as the sum of all the previously mentioned one-electron integrals.



Two-electron integrals in the AO basis set

Electron Repulsion Integrals (ERIs) are given by

$$W_{pqrs} = \iint \chi_p^*(\mathbf{r}_1) \chi_q^*(\mathbf{r}_2) \frac{1}{|\mathbf{r}_1 - \mathbf{r}_2|} \chi_r(\mathbf{r}_1) \chi_s(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2$$

ERIs have (p, q, r, s) as indices, so their number grows formally as N^4 where N is the number of AOs. However, due to the locality of the AOs these integrals are negligible when the overlap between $\langle \chi_p | \chi_r \rangle$ or $\langle \chi_q | \chi_s \rangle$ is close to zero. Hence, this data structure is sparse, and it is preferable to store only the non-zero values. One can also remark that the ERIs are symmetric with respect to the exchange of electron coordinates \mathbf{r}_1 and \mathbf{r}_2 , and for real orbitals one can also exchange the indices p and r and/or q and s . This enables an even more compact storage if only unique values are stored.

One-electron integrals in the MO basis set

Post Hartree-Fock methods generally require integrals transformed from the AO basis set to the Molecular Orbital (MO) basis set, so it is convenient to be able to read them from a file. Some codes within TREX don't enforce the orthogonality between the MOs, so we provide the possibility to store the overlap matrix of the MOs in addition to all the other ones (kinetic, potential, ECP, core Hamiltonian).

Two-electron integrals in the MO basis set

The transformation of ERIs from the AO to the MO basis can be expensive, as it scales as N^5 . Therefore, storing these integrals is often necessary. The structure of these integrals is the same as the integrals in the AO basis set, with the same permutation symmetry in the indices:

$$W_{ijkl}^{\text{MO}} = \iint \phi_i^*(\mathbf{r}_1) \phi_j^*(\mathbf{r}_2) \frac{1}{|\mathbf{r}_1 - \mathbf{r}_2|} \phi_k(\mathbf{r}_1) \phi_l(\mathbf{r}_2) d\mathbf{r}_1 d\mathbf{r}_2$$

However, as the MOs are generally delocalized, these data structures are less sparse than their AO counterparts.

Summary

Table 1, summarizes the data that need to be stored in the files.

Table 1: Information that should be stored in the file.

Metadata	File description Code used to write the file Authors of the file
----------	--

Nuclei	Number of nuclei Atomic charges XYZ coordinates Atom labels Point-group Symmetry
Electrons	Number of \uparrow and \downarrow electrons
ECPs	Effective charge Exponents of the local component Coefficients of the local component Powers of r of the local component Exponents of the non-local component Coefficients of the non-local component Powers of r of the non-local component
Atomic basis set	Type : Gaussian or Slater Cartesian or spherical coordinates Normalization factors of the shells Nuclei on which the functions are centered Angular momenta Exponents of the primitives Coefficients of the primitives
One-electron integrals in the AO basis	Overlap integrals Kinetic energy Potential energy Local component of the ECP Non-Local component of the ECP Core Hamiltonian
Two-electron integrals in the AO basis	Array of indices of ERI Array of values of ERI

Molecular orbitals	Type : Hartree-Fock, Localized, Natural, . . . Coefficients Class : Core, Inactive, Active, Virtual, Deleted, . . . Symmetry Occupation number
One-electron integrals in the MO basis	Kinetic energy Potential energy Local component of the ECP Non-Local component of the ECP Core Hamiltonian
Two-electron integrals in the MO basis	Array of indices of ERI Array of values of ERI

3 Design of the library

The design of the library is split in two distinct sections: the front end, and the back end. The front end is the interface between the users and the library, and the back end is the interface between the library and the physical storage. The library is designed to decouple as much as possible the front end from the back end.

License

The library is licensed under the open-source 3-clause BSD license to facilitate its adoption in all quantum chemistry software, commercial or not.

The front end

From table 1, it appears clearly that the data can be organized in a tree structure, where the root of the tree is the file, elementary pieces of data are leaves of the tree, and the nodes between the root and the leaves constitute groups. Hence, the data is organized in *groups* (the left column of Table 1), and each group contains the associated data.

Most of the codes of the CoE are written in Fortran, with some scripts in Python. Therefore, the Application Programming Interface (API) should be such that the functions can be called easily in Fortran, and such that a Python interface to the library is easy to write. These constraints have lead us to the choice of implementing the library in C, with interfaces for Fortran and Python.



To maximize the portability of the library, the data exposed to the user in the API are reduced to the subset of elementary C types: 32-/64-bit integers and floats, scalars or arrays. The integer types defined in `<stdint.h>` (`int64_t` and `int32_t`) are used instead of the native C integer types. Boolean variables are stored as integers, 1 for `true` and 0 for `false`, and complex numbers are represented as an array of two floats, the real part at address 0 and the imaginary part at address 1.

The library allows the user to create, open or close a file, and check if the file exists. The file may be opened in read-only mode to protect it from accidental data corruption.

Once the file is open, any kind of data can be read by calling the corresponding function. The function names for data access obey the following structure:

```
trexio_<read|write|has>_<group>_<data>
```

where `<data>` is the data to access, `<group>` is the group in which the data belongs, and `<read|write|has>` can read, write, or test if the data exists in the file. For example, if a user wants to read the array of nuclear coordinates, the function `trexio_read_nucleus_coord` should be called.

As the different data types are rather limited and the rules to build the function names are simple, the best strategy is to generate automatically the C code for all these functions from a simple configuration file, `trex.json` as:

```

1 {
2   "metadata": {
3     "code_num"           : [ "int" , [                               ] ]
4     , "code"             : [ "char", [ "metadata.code_num" , "128" ] ]
5     , "author_num"       : [ "int" , [                               ] ]
6     , "author"           : [ "char", [ "metadata.author_num", "128" ] ]
7     , "description_length" : [ "int" , [                               ] ]
8     , "description"      : [ "char", [ "metadata.description_length" ] ]
9   },
10
11  "electron": {
12    "up_num"              : [ "int", [ ] ]
13    , "dn_num"            : [ "int", [ ] ]
14  },
15
16  "nucleus": {
17    "num"                  : [ "int" , [                               ] ]
18    , "charge"             : [ "float", [ "nucleus.num"           ] ]
19    , "coord"              : [ "float", [ "nucleus.num", "3"      ] ]
20    , "label"              : [ "char" , [ "nucleus.num", "32"     ] ]
21    , "point_group"       : [ "char" , [ "32"                     ] ]
22  },
23
24  "ecp": {
25    "lmax_plus_1"         : [ "int" , [ "nucleus.num"           ] ]
26    , "z_core"            : [ "int" , [ "nucleus.num"           ] ]
27    , "local_n"           : [ "int" , [ "nucleus.num"           ] ]
28    , "local_num_n_max"   : [ "int" , [                               ] ]
29    , "local_exponent"    : [ "float", [ "nucleus.num", "ecp.local_num_n_max" ] ]

```

```

30     , "local_coef"           : [ "float", [ "nucleus.num", "ecp.local_num_n_max" ] ]
31     , "local_power"        : [ "int" , [ "nucleus.num", "ecp.local_num_n_max" ] ]
32     , "non_local_n"       : [ "int" , [ "nucleus.num" ] ]
33     , "non_local_num_n_max": [ "int" , [ ] ]
34     , "non_local_exponent" : [ "float", [ "nucleus.num", "ecp.non_local_num_n_max" ] ]
35     , "non_local_coef"    : [ "float", [ "nucleus.num", "ecp.non_local_num_n_max" ] ]
36     , "non_local_power"   : [ "int" , [ "nucleus.num", "ecp.non_local_num_n_max" ] ]
37 },
38
39 "basis" : {
40     "type"                 : [ "char" , [ "32" ] ]
41     , "shell_num"         : [ "int" , [ ] ]
42     , "shell_factor"     : [ "float", [ "basis.shell_num" ] ]
43     , "shell_center"     : [ "int" , [ "basis.shell_num" ] ]
44     , "shell_ang_mom"    : [ "int" , [ "basis.shell_num" ] ]
45     , "shell_prim_num"   : [ "int" , [ "basis.shell_num" ] ]
46     , "prim_index"       : [ "int" , [ "basis.shell_num" ] ]
47     , "prim_num"        : [ "int" , [ ] ]
48     , "exponent"        : [ "float", [ "basis.prim_num" ] ]
49     , "coefficient"     : [ "float", [ "basis.prim_num" ] ]
50 },
51
52 "ao" : {
53     "num"                 : [ "int" , [ ] ]
54     , "cartesian"       : [ "int" , [ ] ]
55     , "shell"           : [ "int" , [ "ao.num" ] ]
56     , "normalization"   : [ "float", [ "ao.num" ] ]
57 },
58
59 "basis_1e_int" : {
60     "overlap"            : [ "float", [ "ao.num", "ao.num" ] ]
61     , "kinetic"         : [ "float", [ "ao.num", "ao.num" ] ]
62     , "potential"       : [ "float", [ "ao.num", "ao.num" ] ]
63     , "ecp_local"       : [ "float", [ "ao.num", "ao.num" ] ]
64     , "ecp_non_local"   : [ "float", [ "ao.num", "ao.num" ] ]
65     , "core_hamiltonian": [ "float", [ "ao.num", "ao.num" ] ]
66 },
67
68 "basis_2e_int" : {
69     "eri"                : [ "float sparse", [ "ao.num", "ao.num", "ao.num", "ao.num" ] ]
70 },
71
72 "mo" : {
73     "type"               : [ "char" , [ "32" ] ]
74     , "num"              : [ "int" , [ ] ]
75     , "coef"             : [ "float", [ "basis.shell_num", "mo.num" ] ]
76     , "class"           : [ "char" , [ "mo.num", "32" ] ]
77     , "symmetry"        : [ "char" , [ "mo.num", "32" ] ]
78     , "occupation"      : [ "float", [ "mo.num" ] ]
79 },
80
81 "mo_1e_int" : {

```



```

82     "kinetic"           : [ "float", [ "mo.num", "mo.num" ] ]
83     , "potential"      : [ "float", [ "mo.num", "mo.num" ] ]
84     , "ecp_local"      : [ "float", [ "mo.num", "mo.num" ] ]
85     , "ecp_non_local"  : [ "float", [ "mo.num", "mo.num" ] ]
86     , "core_hamiltonian" : [ "float", [ "mo.num", "mo.num" ] ]
87   },
88
89   "mo_2e_int" : {
90     "eri"           : [ "float sparse", [ "mo.num", "mo.num", "mo.num", "mo.num" ] ]
91   },
92
93   "rdm" : {
94     "one_e"         : [ "float", [ "mo.num", "mo.num" ] ]
95     , "two_e"       : [ "float sparse", [ "mo.num", "mo.num", "mo.num", "mo.num" ] ]
96   }
97 }

```

Such a simple file will be extremely valuable for the evolution of the library when some additional features will need to be implemented, such as the Configuration Interaction (CI) wave function, the Jastrow factor, data for periodic systems (cell tensor), range-separated ERIs, *etc.*

Conventions

To facilitate the understanding of the users, we have defined strong conventions. Indeed, using as few exceptions as possible makes it easier for users to guess the answers to their questions.

- All the data are stored in atomic units.
- Preprocessor constants are expressed in upper case.
- Pointers are always set to `NULL` when not attached to a memory block.
- The singular is always used for the names of the variables.
- The `num` suffix denotes counting. For example `apple_num` is the number of apples.
- All the functions will be provided with two versions. One with a 32-bit representation of integers, and one with a 64-bit. This will reduce the risk of integer overflows.
- The first argument of the functions is always the file handle.
- The functions always return an exit code.
- Reading functions pass the argument in which to read the data by address, similarly to `scanf`.
- Writing functions pass the argument to be written by value, similarly to `printf`.

Error handling

In case of error, the user of the library should be informed that the called function did not succeed in performing the requested work. The library should never make the calling program crash, nor decide to halt the execution. It should not even take the decision to print something on the terminal. The choice of how to handle the errors should be left completely to the code calling the function.

An effort is made in checking the validity of all the arguments of the functions to ensure that the preconditions are fulfilled. If some unexpected behavior happens, the function returns with an error code, which can be translated to a string using a function call.

Safety

Multiple precautions need to be taken to prevent users from accidentally corrupting files. For the alpha release, we have chosen the simple model in which when the file is opened by a process, the file is locked until it is closed. This ensures that the data previously read from the file has not been changed by another process, so it is easy to ensure that the data read from the file and cached in memory is consistent with the data stored in the file.

A second level of safety needs to be added for multi-threaded environments, to avoid situations where multiple threads write the same data at the same time, which could lead to inconsistent data being written. All the provided functions are thread-safe.

Prototype library

Before implementing the actual TREX Input/Output (TRESIO) library, a prototype was created using the design exposed in the previous section. This library is available in a repository under the TREX GitHub organization¹. This prototype uses the Easy Fortran Input/Output (EZPIO) library generator[3] as a back end to generate the stored files. This enabled us to concentrate on the front end, which is the direct interface to the users, while being able to use it to see if it can be conveniently used in practice.

The back end

We would like the data to be organized in the file, reproducing the hierarchy of the data. Using a binary format is desirable for the performance of large data sets, such as integrals or density matrices, but binary files are not necessarily compatible between different architectures because of the endianness of the binary representation. Hence, if we store data in binary format, the back end should make the files machine-independent by handling properly the endianness.

Finally, as the produced files are likely to be archived on open data repositories, it is desirable to have the possibility to compress them efficiently.

The EZPIO library, already used in QMC=Chem and Quantum Package, fulfills some of the requirements (automatically generated library and hierarchical storage). But this library is not well adapted to High Performance Computing (HPC) systems. EZPIO generates a large number of small

¹<https://github.com/trex-coe/trex-io-prototype>



files, and large supercomputers use distributed file systems such as Lustre[4], General Parallel File System (GPFS)[5] or BeeGFS[6] which all suffer from a huge performance impact when multiple files are used. In addition, EZFIO is intended to be easy to use, but Input/Output (I/O) performance is not the main objective of this library.

The Hierarchical Data Format (HDF5) file format and library[7] address properly all the required aspects. The first version of HDF5 was released in 1998, so this file format has been present in the HPC landscape for a long time. The data can be stored in a hierarchy similar to a file system, exactly in the way we described in table 1. The library also provides compression possibilities, and is reputed for its high performance read/write operations.

Although HDF5 fulfills all our needs, we need to be careful about some important side effects of using such a library. First, as data are written in a binary format, it is possible to corrupt a file if the program crashes during a write operation. Secondly, HDF5 is a complex piece of software, which might not be installed (or even difficult to install) on some systems. If our library only provides an HDF5 back end, the users unable to install HDF5 will not be able to use our library, and therefore will not be able to use any of the TREX codes. Hence, we need to protect our users from these situations.

We also provide a text-file back end. This back end is by far less efficient, but it has the advantage that it requires no dependencies. We also provide a tool to convert files from the HDF5 format to the text format, to ensure that if some large data have been prepared in the HDF5 format, it can be converted to the text file format for following calculations.

The separation of the front end and the back end makes it easy to implement new back ends in the future.

4 Illustrative example of usage of the library

We propose here to show as an example a Fortran program that reads x, y, z coordinates of a molecule and a basis set from a file obtained from the Basis Set Exchange (BSE) web site,[8] and stores it into a file using the TREXIO library.

```
1  subroutine fail_if_error(file,info)
2    use trexio
3    implicit none
4    integer*8, intent(in)      :: file
5    integer, intent(in)       :: info
6    character*(*), intent(out) :: message
7
8    if (info /= TREXIO_SUCCESS) then
9      call trexio_strerror(file, info, message)
10     print *, info, message
11     stop -1
12   end if
13 end subroutine check_success
14
15 !-----
16
```



```

17 subroutine read_xyz(trex_file, xyz_filename)
18   use trexio
19   implicit none
20   integer*8, intent(in)           :: trex_file
21   character*(128), intent(in)    :: xyz_filename
22   integer*8                       :: nucl_num           ! Number of nuclei
23   character*(256)                 :: title             ! Title of the file
24   character*(32), allocatable     :: nucl_label(:)     ! Atom labels
25   real*8, allocatable             :: nucl_charge(:)    ! Nuclear charges
26   real*8, allocatable             :: nucl_coord(:, :)  ! Nuclear coordinates
27   integer*8                       :: i
28   integer                         :: j
29   integer                         :: info
30   double precision, parameter     :: a0 = 0.52917721067d0
31
32   open(unit=10, file=xyz_filename)
33
34   read(10, *) nucl_num
35
36   allocate(nucl_label(nucl_num), &
37           nucl_charge(nucl_num), &
38           nucl_coord(3, nucl_num) )
39
40   read(10, '(A)') title
41
42   do i=1, nucl_num
43     read(10, *) nucl_label(i), nucl_coord(1:3, i)
44
45     info = trexio_element_number_of_symbol(trim(nucl_label(i)), j)
46     call check_success(info, 'Unable to convert symbol to number')
47
48     nucl_charge(i) = dble(j)
49   end do
50
51   close(10)
52
53   ! Convert into atomic units
54   nucl_coord = nucl_coord / a0
55
56   info = trexio_write_nucleus_num(trex_file, nucl_num)
57   call check_success(info, 'Unable to write number of nuclei')
58
59   info = trexio_write_nucleus_coord(trex_file, nucl_coord)
60   call check_success(info, 'Unable to write nuclear coordinates')
61
62   info = trexio_write_nucleus_charge(trex_file, nucl_charge)
63   call check_success(info, 'Unable to write nuclear charges')
64
65   info = trexio_write_nucleus_label(trex_file, nucl_label)
66   call check_success(info, 'Unable to write nuclear labels')
67
68   beta_num = int(sum(nucl_charge(:)))/2

```



```

69  alpha_num = int(sum(nucl_charge(:))) - beta_num
70
71  info = trexio_write_electron_up_num(trex_file,alpha_num)
72  call check_success(info, 'Unable to write up electrons')
73
74  info = trexio_write_electron_dn_num(trex_file,beta_num)
75  call check_success(info, 'Unable to write dn electrons')
76
77  end subroutine read_xyz
78
79  !-----
80
81  subroutine read_basis(trex_file, basis_filename)
82    use trexio
83    implicit none
84    integer*8, intent(in)          :: trex_file
85    character*(128), intent(in)    :: basis_filename
86
87    integer*8                      :: nucl_num          ! Number of nuclei
88    character*(32), allocatable    :: nucl_label(:)    ! Atom labels
89    integer*8                      :: shell_num, prim_num
90    integer*8, allocatable         :: shell_center(:)
91    integer , allocatable         :: shell_ang_mom(:)
92    integer*8, allocatable         :: shell_prim_num(:)
93    integer*8, allocatable         :: prim_index(:)
94    double precision, allocatable  :: shell_factor(:)
95    double precision, allocatable  :: exponent(:)
96    double precision, allocatable  :: coefficient(:)
97    character*(32)                 :: label
98    character*(80)                 :: buffer
99    integer                        :: i,j,k,n_shell,n_prim
100   integer                        :: info
101
102   open(unit=10, file=basis_filename)
103
104   info = trexio_read_nucleus_num(trex_file, nucl_num)
105   call check_success(info, 'Unable to read number of nuclei')
106
107   allocate(nucl_label(nucl_num))
108   info = trexio_read_nucleus_label(trex_file, nucl_label)
109   call check_success(info, 'Unable to read nuclear label')
110
111   shell_num = 0
112   prim_num  = 0
113   ! Find dimensioning variables
114   do i=1,nucl_num
115     info = trexio_element_name_of_symbol(nucl_label(i), label)
116     call check_success(info, 'Unable to get name of label')
117
118     ! Find element
119     rewind(10)
120   do

```



```

121     read(10,*, iostat=j) buffer
122     if (j < 0) exit
123     if (trim(buffer) == trim(label)) then
124         j=0
125         exit
126     end if
127 end do
128 if (j < 0) exit
129
130 ! Read shell
131 do
132     read(10,*,iostat=k) buffer, j
133     if (k /= 0) exit
134     shell_num = shell_num + 1
135     prim_num = prim_num + j
136     do k=1,j
137         read(10,*)
138     end do
139 end do
140
141 end do
142
143 buffer = 'Gaussian'
144 info = trexio_write_basis_type(trex_file, buffer)
145 call check_success(info, 'Unable to write basis type')
146
147 info = trexio_write_basis_shell_num(trex_file, shell_num)
148 call check_success(info, 'Unable to write basis shell_num')
149
150 info = trexio_write_basis_prim_num(trex_file, prim_num)
151 call check_success(info, 'Unable to write basis prim_num')
152
153 allocate(shell_center(shell_num), &
154          shell_ang_mom(shell_num), &
155          shell_prim_num(shell_num), &
156          prim_index(shell_num), &
157          shell_factor(shell_num), &
158          exponent(prim_num), &
159          coefficient(prim_num))
160
161 shell_num = 1
162 prim_num = 1
163 ! Read data
164 do i=1,nucl_num
165     info = trexio_element_name_of_symbol(nucl_label(i), label)
166     call check_success(info, 'Unable to get name of label')
167
168 ! Find element
169 rewind(10)
170 do
171     read(10,*, iostat=j) buffer
172     if (j < 0) exit

```



```

173     if (trim(buffer) == trim(label)) then
174         j=0
175         exit
176     end if
177 end do
178
179 ! Read shell
180 do
181     read(10,*,iostat=j) label, k
182     if (j /= 0) exit
183     select case (label(1:1))
184     case ('S')
185         shell_ang_mom(shell_num) = 0
186     case ('P')
187         shell_ang_mom(shell_num) = 1
188     case ('D')
189         shell_ang_mom(shell_num) = 2
190     case ('F')
191         shell_ang_mom(shell_num) = 3
192     case ('G')
193         shell_ang_mom(shell_num) = 4
194     case ('H')
195         shell_ang_mom(shell_num) = 5
196     case ('I')
197         shell_ang_mom(shell_num) = 6
198     case default
199         stop 'Too high angular momentum'
200     end select
201     shell_prim_num(shell_num) = k
202     shell_center(shell_num) = i
203     prim_index(shell_num) = prim_num
204     do j=1,shell_prim_num(shell_num)
205         read(10,*) buffer, exponent(prim_num), coefficient(prim_num)
206         prim_num = prim_num + 1
207     end do
208     shell_num = shell_num + 1
209 end do
210
211 end do
212
213 close(10)
214
215 info = trexio_write_basis_shell_center(trex_file, shell_center)
216 call check_success(info, 'Unable to write basis shell_center')
217
218 info = trexio_write_basis_shell_ang_mom(trex_file, shell_ang_mom)
219 call check_success(info, 'Unable to write basis shell_ang_mom')
220
221 info = trexio_write_basis_shell_prim_num(trex_file, shell_prim_num)
222 call check_success(info, 'Unable to write basis shell_prim_num')
223
224 info = trexio_write_basis_prim_index(trex_file, prim_index)

```



```

225  call check_success(info, 'Unable to write basis prim_index')
226
227  info = trexio_write_basis_exponent(trex_file, exponent)
228  call check_success(info, 'Unable to write basis exponent')
229
230  info = trexio_write_basis_coefficient(trex_file, coefficient)
231  call check_success(info, 'Unable to write basis coefficient')
232
233  return
234
235 10 continue
236  stop 'Unable to find element in basis write file'
237 end subroutine read_basis
238
239 !-----
240
241 program write_example
242  use trexio
243  implicit none
244
245  character*(128)           :: xyz_filename      ! Name of the xyz file
246  character*(128)           :: basis_filename   ! Name of the basis file
247  integer*8                 :: trex_file        ! Handle for the TREX file
248  integer                   :: i
249  integer                   :: info
250  character*(*), parameter  :: trex_filename = 'trex_file'
251
252  ! Get the xyz file name from the command line and user name
253  ! =====
254
255  i = command_argument_count()
256  if (i /= 2) then
257    print *, 'Expected:'
258    print *, ' - xyz file as 1st argument'
259    print *, ' - basis file as 2nd argument'
260    print *, './fortran_write cyanoformaldehyde.xyz cc-pvtz'
261    stop -1
262  end if
263
264  call get_command_argument(1,xyz_filename)
265  call get_command_argument(2,basis_filename)
266
267  ! Open the TREX file
268  ! -----
269
270  trex_file = trexio_open(trex_filename,'w', TREXIO_HDF5)
271
272  ! Write the data
273  ! -----
274
275  call read_xyz(trex_file, xyz_filename)
276  call read_basis(trex_file, basis_filename)

```




```

277
278  ! Close the file
279  ! -----
280
281  info = trexio_close(trex_file)
282  call check_success(info, 'Unable to close file')
283
284  print *, 'Wrote file '//trim(trex_filename)
285
286  end program write_example

```

To read back the data and print it, we can use the following program:

```

1  subroutine check_success(info,message)
2    use trexio
3    implicit none
4    integer, intent(in) :: info
5    character*(*)      :: message
6
7    if (info /= TREXIO_SUCCESS) then
8      print *, info, message
9      stop -1
10   end if
11  end subroutine check_success
12
13  !-----
14
15  subroutine read_electrons(trex_file, alpha_num, beta_num)
16    use trexio
17    implicit none
18    integer*8, intent(in)      :: trex_file
19    integer*8, intent(out)     :: alpha_num, beta_num
20    integer :: info
21
22    info = trexio_read_electron_up_num(trex_file,alpha_num)
23    call check_success(info, 'Unable to read up electrons')
24
25    info = trexio_read_electron_dn_num(trex_file,beta_num)
26    call check_success(info, 'Unable to read dn electrons')
27  end subroutine read_electrons
28
29  !-----
30
31  subroutine read_nuclei(trex_file, nucl_num, nucl_coord, nucl_charge, nucl_label)
32    use trexio
33    implicit none
34    integer*8, intent(in)      :: trex_file
35    integer*8, intent(in)      :: nucl_num
36    double precision, intent(out) :: nucl_coord(3,nucl_num)
37    double precision, intent(out) :: nucl_charge(nucl_num)
38    character*(32), intent(out) :: nucl_label(nucl_num)

```



```

39
40  integer :: info
41
42  info = trexio_read_nucleus_coord(trex_file,nucl_coord)
43  call check_success(info, 'Unable to read nuclear coordinates')
44
45  info = trexio_read_nucleus_charge(trex_file,nucl_charge)
46  call check_success(info, 'Unable to read nuclear charges')
47
48  info = trexio_read_nucleus_label(trex_file,nucl_label)
49  call check_success(info, 'Unable to read nuclear labels')
50
51  end subroutine read_nuclei
52
53  !-----
54
55  subroutine read_basis(trex_file, shell_num, prim_num, center, ang_mom, &
56     shell_prim_num, prim_index, expo, coef)
57  use trexio
58  implicit none
59  integer*8, intent(in)           :: trex_file
60  integer*8, intent(in)           :: shell_num, prim_num
61  integer*8, intent(out)          :: center(shell_num), shell_prim_num(shell_num)
62  integer , intent(out)           :: ang_mom(shell_num)
63  integer*8, intent(out)          :: prim_index(shell_num)
64  double precision, intent(out)   :: expo(prim_num)
65  double precision, intent(out)   :: coef(prim_num)
66
67  integer                         :: info
68
69  info = trexio_read_basis_shell_center(trex_file, center)
70  call check_success(info, 'Unable to read basis shell_center')
71
72  info = trexio_read_basis_shell_ang_mom(trex_file, ang_mom)
73  call check_success(info, 'Unable to read basis shell_ang_mom')
74
75  info = trexio_read_basis_shell_prim_num(trex_file, shell_prim_num)
76  call check_success(info, 'Unable to read basis shell_prim_num')
77
78  info = trexio_read_basis_prim_index(trex_file, prim_index)
79  call check_success(info, 'Unable to read basis prim_index')
80
81  info = trexio_read_basis_exponent(trex_file, expo)
82  call check_success(info, 'Unable to read basis exponent')
83
84  info = trexio_read_basis_coefficient(trex_file, coef)
85  call check_success(info, 'Unable to read basis coefficient')
86
87  end subroutine read_basis
88
89  !-----
90

```



```

91 program read_example
92   use trexio
93   implicit none
94   integer*8                :: nucl_num          ! Number of nuclei
95   character*(32), allocatable :: nucl_label(:)  ! Atom labels
96   real*8, allocatable      :: nucl_charge(:)    ! Nuclear charges
97   real*8, allocatable      :: nucl_coord(:, :)  ! Nuclear coordinates
98   integer*8                :: alpha_num        ! Number of alpha electrons
99   integer*8                :: beta_num         ! Number of beta electrons
100
101   integer*8                :: trex_file        ! Handle for the TREX file
102   integer                  :: i,j,k
103   integer                  :: info
104   character*(*), parameter :: trex_filename = 'trex_file'
105   double precision, parameter :: a0 = 0.52917721067d0
106
107   integer*8                :: shell_num, prim_num
108   integer*8, allocatable   :: shell_center(:)
109   integer , allocatable   :: shell_ang_mom(:)
110   integer*8, allocatable   :: shell_prim_num(:)
111   integer*8, allocatable   :: prim_index(:)
112   double precision, allocatable :: shell_factor(:)
113   double precision, allocatable :: exponent(:)
114   double precision, allocatable :: coefficient(:)
115   character*(32)           :: bastype
116   character*(32)           :: label
117   character, parameter     :: ang_mom(0:6) = (/ 'S', 'P', 'D', 'F', 'G', 'H', 'I' /)
118
119   ! Read the data from the TREX file
120   ! =====
121
122   ! Open the file
123   ! -----
124
125   trex_file = trexio_open(trex_filename,'r',TREXIO_HDF5)
126
127   ! Read the data
128   ! -----
129
130   call read_electrons(trex_file,alpha_num,beta_num)
131   print *, 'Electrons: ', alpha_num, ' up, ', beta_num, ' down'
132
133
134   info = trexio_read_nucleus_num(trex_file,nucl_num)
135   call check_success(info, 'Unable to read number of nuclei')
136
137   allocate(nucl_coord(3,nucl_num), &
138           nucl_charge(nucl_num), &
139           nucl_label(nucl_num) )
140
141   call read_nuclei(trex_file, nucl_num, nucl_coord, nucl_charge, nucl_label)
142

```



```

143 do i=1,nucl_num
144     print '(A4, 2X, F4.1,3(3X,F12.8))', nucl_label(i), nucl_charge(i), nucl_coord(1:3,i)
145 end do
146
147
148 bastype=''
149 info = trexio_read_basis_type(trex_file, bastype)
150 call check_success(info, 'Unable to read basis type')
151 print *, 'Basis type: ', trim(bastype)
152
153 info = trexio_read_basis_shell_num(trex_file, shell_num)
154 call check_success(info, 'Unable to read basis shell_num')
155
156 info = trexio_read_basis_prim_num(trex_file, prim_num)
157 call check_success(info, 'Unable to read basis prim_num')
158
159 allocate(shell_center(shell_num), &
160          shell_ang_mom(shell_num), &
161          shell_prim_num(shell_num), &
162          prim_index(shell_num), &
163          shell_factor(shell_num), &
164          exponent(prim_num), &
165          coefficient(prim_num))
166
167 call read_basis(trex_file, shell_num, prim_num, shell_center, &
168               shell_ang_mom, shell_prim_num, prim_index, exponent, coefficient)
169
170 k=0
171 do i=1,shell_num
172     if (shell_center(i) /= k) then
173         k = shell_center(i)
174         info = trexio_element_name_of_symbol(trim(nucl_label(k)),label)
175         call check_success(info, 'Unable to read name of element :')
176         print *, ''
177         print *, trim(label)
178     end if
179     print *, ang_mom(shell_ang_mom(i)), shell_prim_num(i)
180     do j=1,shell_prim_num(i)
181         print '(I3,X,E16.8,3X,E16.8)', j, &
182             exponent(prim_index(i)+j-1) , coefficient(prim_index(i)+j-1)
183     end do
184 end do
185
186
187 ! Close the file
188 ! -----
189
190 info = trexio_close(trex_file)
191 call check_success(info, 'Unable to close file')
192
193 end program read_example

```



References

- [1] E. F. Valeev, “Libint: A library for the evaluation of molecular integrals of many-body operators over gaussian functions,” <http://libint.valeev.net/>, 2020, version 2.7.0-beta.6.
- [2] M. Burkatzki, C. Filippi, and M. Dolg, “Energy-consistent small-core pseudopotentials for 3d-transition metals adapted to quantum Monte Carlo calculations,” *J. Chem. Phys.*, vol. 129, no. 16, p. 164115, Oct 2008.
- [3] “The Easy Fortran Input/Output (EZFIO) library generator,” Mar 2021, [Online; accessed 5. Mar. 2021]. [Online]. Available: <https://gitlab.com/scemama/EZFIO>
- [4] “Lustre file system,” Mar 2021, [Online; accessed 5. Mar. 2021]. [Online]. Available: <https://www.lustre.org>
- [5] “GPFS (General Parallel File System) - IBM,” Mar 2021, [Online; accessed 5. Mar. 2021]. [Online]. Available: https://researcher.watson.ibm.com/researcher/view_group.php?id=4840
- [6] “BeeGFS Documentation v7.2.1 — BeeGFS Documentation v7.2.1,” Feb 2021, [Online; accessed 5. Mar. 2021]. [Online]. Available: <https://doc.beegfs.io/latest/index.html>
- [7] The HDF Group. (2000-2010) Hierarchical data format version 5. [Online]. Available: <http://www.hdfgroup.org/HDF5>
- [8] “Basis Set Exchange (BSE),” Mar 2021, [Online; accessed 5. Mar. 2021]. [Online]. Available: <https://www.basissetexchange.org>

