

---

# **Genarris**

***Release 2.0***

**Rithwik Tom, Timothy Rose, Imanuel Bier**

**Sep 20, 2019**



---

## Contents

---

<b>1</b>	<b>Installation</b>	<b>1</b>
<b>2</b>	<b>Introduction to Running Genarris</b>	<b>5</b>
2.1	Configuration File . . . . .	5
2.2	Option Category . . . . .	6
2.3	Output Formats . . . . .	6
<b>3</b>	<b>Genarris 2.0 Procedures for Robust Workflow</b>	<b>7</b>
3.1	Description . . . . .	7
3.2	Genarris Procedures . . . . .	7



# CHAPTER 1

---

## Installation

---

1) Setup MPI and MKL If already installed and modules exist, load them after unloading all conflicting modules. Note, in this installation tutorial we will use intel including intel's parallel studio package, but other program environments such as gnu will also work. e.g.:

```
module unload gnu
module unload openmpi
module load intel
module load impi
```

If MKL and MPI are already installed but modules do not exist, include the MPI and MKL directories in your environment variables. e.g.:

```
#Change to your parallel studio path
export $intel=/opt/ohpc/pub/intel/intel18/compilers_and_libraries_2018.3.222/linux

export PATH="$intel/mpi/intel64/bin_ohpc:\
$intel/mpi/intel64/bin:$intel/bin/intel64:$PATH"

export LD_LIBRARY_PATH="$intel/mpi/intel64/lib:$intel/mpi/mic/lib:\
$intel/compiler/lib/intel64:$intel/compiler/lib/intel64_lin:\
$intel/ipp/lib/intel64:$intel/mkl/lib/intel64_lin:\
$intel/tbb/lib/intel64/gcc4.1:\
/opt/ohpc/pub/intel/intel18/debugger_2018/iga/lib:\
/opt/ohpc/pub/intel/intel18/debugger_2018/libipt/intel64/lib:\
$intel/daal/lib/intel64_lin:$intel/tbb/lib/intel64_lin/gcc4.4"

export LD_LIBRARY_PATH="$intel/mpi/intel64/lib:$intel/mpi/mic/lib:$intel/compiler/lib/\
↪intel64:$intel/compiler/lib/intel64_lin:\
$intel/ipp/lib/intel64:$intel/mkl/lib/intel64_lin:\
$intel/tbb/lib/intel64/gcc4.1:\
/opt/ohpc/pub/intel/intel18/debugger_2018/iga/lib:\
/opt/ohpc/pub/intel/intel18/debugger_2018/libipt/intel64/lib:\
$intel/daal/lib/intel64_lin:$intel/tbb/lib/intel64_lin/gcc4.4"
```

Also export LD\_PRELOAD to load the parallel studio MKL and Scalapack so importing FHI-aims and numpy does not cause conflict. e.g.:

```
export LD_PRELOAD="$intel/mkl/lib/intel64_lin/libmkl_intel_lp64.so:\
$intel/mkl/lib/intel64_lin/libmkl_sequential.so:\
$intel/mkl/lib/intel64_lin/libmkl_core.so:\
$intel/mkl/lib/intel64_lin/libmkl_blacs_intelmpi_lp64.so:\
$intel/mkl/lib/intel64_lin/libmkl_scalapack_lp64.so:\
$intel/mpi/intel64/lib/libmpi.so.12"
```

2) create a python 3.5+ virtual environment e.g.:

```
#Change this to your desired anaconda install path
export $anaconda=${HOME}/anaconda
mkdir $anaconda
cd $anaconda
```

download and install anaconda e.g.:

```
wget https://repo.anaconda.com/archive/Anaconda3-2019.07-Linux-x86_64.sh
chmod +x Anaconda3-2019.07-Linux-x86_64.sh
./Anaconda3-2019.07-Linux-x86_64.sh
```

Include anaconda's binary in PATH e.g.:

```
export PATH=$anaconda/anaconda3/bin:$PATH
```

Make a python environment called e.g. genarris\_env by installing intelpython3\_core. e.g.:

```
conda config --add channels intel
conda create -n genarris_env intelpython3_core python=3
```

3) direct your path variables to include the new env e.g.:

```
export PYTHONPATH="$anaconda/anaconda3/envs/genarris_env/lib/python3.6:\
$anaconda/anaconda3/envs/genarris_env/lib/python3.6/site-packages:\
$PYTHONPATH"

export PATH="$intel/mpi/intel64/bin_ohpc:$intel/mpi/intel64/bin:\
$intel/bin/intel64:$anaconda/anaconda3/envs/intelpython3_full/bin:\
$anaconda/anaconda3/bin:$PATH"
```

4) Extract Genarris\_v2.tar.gz into a desired directory and enter it e.g.:

```
export $genarris=${HOME}/genarris
mkdir $genarris
cp Genarris_v2.tar.gz $genarris
cd $genarris
tar -xzf Genarris_v2.tar.gz
```

5) Install Genarris. Note, one reason we recommend to create a python virtual env earlier is that running this installation script will remove the ase installation (if any) in the currently active python environment. e.g.:

```
cd $genarris/Genarris
python setup.py install
```

Genarris is now installed. We will first test that Genarris imports and MPI is working correctly with the following test and then the next step will be to compile FHI-aims as a python-importable library if you desire to use FHI-aims.

6) Test that Genarris imports and MPI is working correctly. Modify the submission script for your backend (here, we used slurm).:

```
cd $genarris/documentation/mpi_and_genarris_test
sbatch mpi_and_genarris_test.sh
```

The desired output is that each rank reports a unique number.

#### 7) Compile libaims into a python-importable library

Set ulimit to avoid any possible memory problems:

```
ulimit -s unlimited
ulimit -v unlimited

# Set OMP_NUM_THREADS to 1
export OMP_NUM_THREADS=1
```

Obtain FHI-aims from <https://aims-git.rz-berlin.mpg.de/aims/FHIaims> If you don't have permissions, ask Volker Blum at [volker.blum@duke.edu](mailto:volker.blum@duke.edu):

```
export $aims=${HOME}/aims #Change to your desired location for FHI-aims
```

In its src directory (\$aims/src), make sure the Makefile has all compilation flags (user defined settings) commented out. Copy the make.sys file in the documentation directory of Genarris into FHI-aims' src directory. The make.sys is pasted here for reference.:

```
cp $genarris/documentation/make.sys $aims/src
```

Note, this make.sys assumes you are using intel's parallel studio and that your cluster's backend is intel. If this isn't the case, you'll need to set the flags accordingly.:

```
# make.sys
#####
# Basic Flags #
#####
FC = mpiifort
FFLAGS = -O3 -ip -fp-model precise -fPIC
F90FLAGS = $(FFLAGS)
ARCHITECTURE = Generic
LAPACKBLAS = -L${MKLROOT}/lib/intel64 \
             -lmkl_intel_lp64 \
             -lmkl_sequential \
             -lmkl_core \
             -lmkl_blacs_intelmpi_lp64 \
             -lmkl_scalapack_lp64
F90MINIFLAGS = -O0 -fp-model precise -fPIC

#####
# Parallelization Flags #
#####
USE_MPI = yes
MPIFC = ${FC}
SCALAPACK = ${LAPACKBLAS}

#####
# C,C++ Flags #
#####
```

(continues on next page)

(continued from previous page)

```
CC = icc
CFLAGS = -O3 -ip -fp-model precise -fPIC
```

Compile FHI-aims as a shared library object:

```
cd $aims/src
make -j 20 libaims.scalapack.mpi
```

where the "20" is however many cores you'd like to use for compilation.

Make a directory for compiling FHI-aims as a python library e.g.:

```
mkdir $aims/aims_as_python_lib
cd $aims/aims_as_python_lib
```

# Copy the Makefile and aims\_w.f90 in the Genarris documentation directory to this directory. A copy of it has been pasted here for reference. Note that you will need to change the libaims version (currently shown as 190522). Again, you'll need to change the f90exec and/or fcompiler flags if your backend is not intel. aims\_w.f90 is a wrapper script to interface with FHI-aims. e.g.:

```
cp $genarris/Genarris/documentation/Makefile $aims/aims_as_python_lib
cp $genarris/Genarris/documentation/aims_w.f90 $aims/aims_as_python_lib
```

Create the Makefile with the following contents:

```
LIBAIMS=${aims}/lib/libaims.190522.scalapack.mpi.so
include_dir=${anaconda}/anaconda3/envs/genarris_env/include

aims_w.so: aims_w.f90
    f2py --f90exec=mpiifort --fcompiler=intelem -m aims_w \
        -c aims_w.f90 ${LIBAIMS} -I${include_dir}

clean:
    rm aims_w.*.so
```

Compile FHI-aims as an importable python library!:

```
make
```

8) Test that FHI-aims can run a job Modify the submission script in the \$genarris/documentation/aims\_test directory to run on your backend (here we used slurm).:

```
export PYTHONPATH=$PYTHONPATH:$aims/aims_as_python_lib
cd $genarris/documentation/aims_test
sbatch aims_test.sh
```



---

## Introduction to Running Genarris

---

### 2.1 Configuration File

Genarris is a random crystal structure generation code that can be adapted to perform *ab initio* crystal structure prediction. The modularity of Genarris is achieved through the sequential execution of procedures. The execution of Genarris is controlled by a **configuration** file. Below is a small example of a configuration file for Genarris.:

```
[Genarris_master]
procedures = ["Pygenarris_Structure_Generation"]

[pygenarris_structure_generation]
# Path to the single molecule file to used for crystal structure generation
molecule_path = relaxed_molecule.in
# Number of cores (MPI ranks) to run this section with
num_cores = 56
# Number of OpenMP Threads
omp_num_threads = 2
num_structures = 5000
Z = 4
sr = 0.85
tol = 0.00001
max_attempts_per_spg_per_rank = 1000000000
geometry_out_filename = glycine_4mpc.out
output_format = json
output_dir = glycine_4mpc_raw_jsons
```

**Sections** of the configuration file are denoted by square brackets, [...]. All parameters that are specified below a section are called **options**. The workflow of Genarris can be precisely controlled by the user by specifying the order of the desired procedures in [Genarris\_master]. The user must also include the corresponding section for each procedure listed in [Genarris\_master]. Each section may have many options which are required, optional, or inferred.

This document details the options for procedures that are executed in the Genarris 2.0 *Robust* workflow. In order these are:

```
[ "Relax_Single_Molecule",  
  "Estimate_Unit_Cell_Volume",  
  "Pygenarris_Structure_Generation",  
  "Run_Rdf_Calc",  
  "Affinity_Propagation_Fixed_Clusters",  
  "FHI_Aims_Energy_Evaluation",  
  "Affinity_Propagation_Fixed_Clusters",  
  "Run_FHI_Aims_Batch"]
```

There are many options that can be specified and modified for each section. All of these options are specified in this document under the **Configuration File Options** section of each procedure.

## 2.2 Option Category

There are three *categories* of **Configuration File Options**. These are *required*, *optional*, and *inferred*. In the **Configuration File Options**, these categories are specified after the *type* of the option, such as *int*, *float*, or *bool*.

1. *Required* options have no category placed after the type in the documentation. These options are required to be in the configuration file for execution of Genarris.
2. *Optional* arguments are specified after the option *type*. These arguments have default settings built into the code perform well in general. The user may specify these *optional* arguments in the configuration file to have more control over the program executing.
3. *Inferred* options are specified after the option *type*. These options may be present in multiple different procedures. For example, the option `aims_lib_dir` is needed in the `Relax_Single_Molecule`, `FHI_Aims_Energy_Evaluation`, and `Run_FHI_Aims_Batch`. But, because it is an inferred parameter, it only needs to be specified once in the earliest procedure in which occurs and then it will be inferred by all further procedures. Options which are inferred are thus optional in all proceeding sections.

## 2.3 Output Formats

There are three output formats supported within the Genarris source code. These are *json*, *geo*, or *both*.

- The *json* file format is the native structure file format for Genarris. This file format supports storing the structure ID, the geometry, and property information.
- The *geo* file format is the file format support by FHI-aims. Additionally, this file format is support by [Jmol](#) , a 3D chemical structure visualizer, and by [ASE](#), the atomic simulation environment tools written for Python.
- The user may also specify *both*, in which case both the *json* files and *geo* file for every structure will be produced.

---

## Genarris 2.0 Procedures for Robust Workflow

---

### 3.1 Description

This section details all arguments and configuration file options for the procedures executed by the Robust Genarris 2.0 workflow. Each procedure is a class function of the of the `Genarris` master class. The documentation follows a standard format for each procedure. The name of the procedure is given first followed by a short description of the function the function it performs. Below the description is the the configuration file options subsection. This section gives the name, the data type, the *Option Category*, and a description of each option which is accepted by the procedure. By referencing this documentation, the user can obtain precise control over the execution of Genarris procedures.

### 3.2 Genarris Procedures

**class** `Genarris.genarris_master.Genarris` (*inst\_path*)

Master class of Genarris. It controls all aspects of the Genarris workflow which can be executed individually or sequentially. Begins by reading and interpreting the configuration file. Calls the defined procedures with the options specified in the configuration file. Some options may be inferred from previous sections if they are not present in every section.

**Affinity\_Propagation\_Fixed\_Clusters** (*comm*)

AP that explores the setting of preference in order to generate desired number of clusters.

#### Configuration File Options

**output\_dir** [str] Path to the directory where the chosen structures will be stored.

**preference\_range** [list] List of two values as the [min, max] of the range of allowable preference values.

**structure\_dir** [str, inferred] Path to the directory of files to be used for the calculation. Default is to infer this value from the previous section.

**dist\_mat\_input\_file** [str, inferred] Path to the distance matrix output from the descriptor calculation. Default is to infer this value from the previous sections.

**output\_format** [str, optional] Format the structure files should be saved as. Default is both.

**cluster\_on\_energy** [bool, optional] Uses energy values to determine exemplars. Structures with the lowest energy values from each cluster are selected. Default is False.

**plot\_histograms** [bool, optional] If histogram plots should be created of the volume and space groups. Default is False.

**num\_of\_clusters** [int or float, optional] Float, must be less than 1. Selects a fraction of the structures. Int, selects specific number of structures equal to int. Default is 0.1.

**num\_of\_clusters\_tolerance** [int, optional] Algorithm will stop if it has generated the number of clusters within the number of desired clusters and this tolerance. Default is 0.

**max\_sampled\_preferences** [int, optional] Maximum number of preference values to try.

**output\_without\_success** [bool, optional] Whether to perform output procedures if the algorithm has reached the maximum number of sampled preferences without finding the correct number of clusters. Default is False.

**affinity\_type** [list, optional] List of [type of affinity, value] argument Scikit-Learn AP algorithm.

**affinity\_matrix\_path** [str, optional] Path to the affinity matrix to use for the AP algorithm. Default is `affinity_matrix.dat`.

**damping** [float, optional] damping argument for Scikit-Learn AP algorithm. Default is 0.5.

**convergence\_iter** [int, optional] convergence\_iter argument for Scikit-Learn AP algorithm. Default is 15.

**max\_iter** [int, optional] max\_iter argument for Scikit-Learn AP algorithm. Default is 1000.

**preference** [int, optional] preference argument for Scikit-Learn AP algorithm. Default is None.

**verbose\_output** [bool, optional] verbose argument for Scikit-Learn AP algorithm. Default is False.

**property\_key** [str, optional] Key which the AP cluster will be stored in the properties of each structure object. Default is `AP_cluster`.

**output\_file** [str, optional] Path where info about the AP algorithm execution will be stored. Default is `./AP_cluster.info`.

**exemplars\_output\_dir** [str, optional] If provided, will output the exemplars of each cluster to this folder. Default is None.

**exemplars\_output\_format** [str, optional] File format of structures to be output. Default is both.

**structure\_suffix** [str, optional] Suffix to apply to structure files which are written. Default is `.json`.

**output\_dir\_2: str, inferred** Code automatically looks for the option `output_dir_2` if the output directory already exists. This is how the code currently identifies that AP is running for a second time. Default behavior is to not use this option if `output_dir` does not already exist.

**num\_of\_clusters\_2: int or float, optional** num\_of\_clusters for second clustering step. Default value is 0.1.

**output\_file\_2** [str, inferred] Use if running AP algorithm twice, such as in the Robust workflow. Default is to use `output_file`.

**exemplars\_output\_dir\_2** [str, inferred] Exemplars output directory if second clustering step is used. Default is to use `exemplars_output_dir`.

**cluster\_on\_energy\_2** [str, inferred] How to choose exemplars for the second clustering step. Default is to use `cluster_on_energy` value.

**energy\_name\_2** [str, inferred] Energy name to use for second clustering step. Default is to use `energy_name`.

#### **Estimate\_Unit\_Cell\_Volume** (*comm*)

Performs volume estimation using a machine learned model train on the CSD and based on Monte Carlo volume integration and topological molecular fragments. See Genarris 2.0 paper for full description.

#### Configuration File Options

**volume\_mean** [float, optional] If provided, uses this value as the volume generation mean without using the ML model to estimate the volume.

**volume\_std** [float, optional] If provided, uses this value for structure generation, otherwise a default value of 0.075 multiplied by the prediction volume per unit cell is provided.

#### **FHI\_Aims\_Energy\_Evaluation** (*comm*, *world\_comm*, *MPI\_ANY\_SOURCE*, *num\_replicas*)

Runs Self-Consistent Field calculation on a pool of structures.

#### Configuration File Options

See `Run_FHI_Aims_Batch()`

#### **Pygenarris\_Structure\_Generation** (*comm*)

Uses the Genarris module written in C to perform structure generation. This module enables generation on special positions.

#### Configuration File Options

**molecule\_path** [str] Path to the relaxed molecule geometry.

**output\_format** [str, optional, default="json"] Determines the type of file which will be output for each structure. Can be one of: json, geo, both.

**output\_dir** [str] Path to the directory which will contain all generated structures which pass the intermolecular distance checks.

**num\_structures** [int] Target number of structures to generate.

**Z** [int] Number of molecules per cell to generate.

**volume\_mean** [float, optional] See `Estimate_Unit_Cell_Volume()`

**volume\_std** [float, optional] See `Estimate_Unit_Cell_Volume()`

**sr** [float, optional] Defines the minimum intermolecular distance that is considered physical by multiplying the sum of the van der Waals radii of the interacting atoms by sr. Default value is 0.85.

**tol** [float, optional] Tolerance to be used to identify space groups compatible with the input molecule.

**max\_attempts\_per\_spg\_per\_rank** [int] Defines the maximum number of attempts the structure generator makes before moving on to the next space group.

**num\_structures\_per\_allowed\_SG\_per\_rank** [int] Number of structures per space group per rank which will be generated by Pygenarris.

**geometry\_out\_filename** [str] Filename where all structures generated by Pygenarris will be found.

**omp\_num\_threads** [int] Number of OpenMP threads to pass into Pygenarris

**truncate\_to\_num\_structures** [bool] If true, will reduce pool to exactly the number defined by num\_structures.

**Run\_Rdf\_Calc** (*comm*)

Runs RDF calculation for the pool of generated structures. RDF descriptor is similar to that described in Behler and Parrinello 2007. Then calculates the structure difference matrix.

### Configuration File Options

**structure\_dir** [str, inferred] Path to the directory of structures to evaluate.

**dist\_mat\_fpath** [str] Path to file to write distance matrix to.

**output\_dir** [str] Path of directory to write structures to (will create if it DNE). If 'no\_new\_output\_dir' then input structures will be overwritten.

**normalize\_rdf\_vectors: bool, optional** Whether to normalize the rdf vectors over the columns of the feature matrix before using them to compute the distance matrix. Default is False.

**standardize\_distance\_matrix: bool** If True, standardizes the distance matrix. The method is to divide all elements by the max value in the distance matrix. Because it is a distance matrix and thus all elements are positive, the standardized elements will be in the range [0, 1]. Default is False.

**save\_envs: bool, optional** Whether to save the environment vectors calculated by the RDF method in the output structure files. Default is False.

**cutoff** [float, optional] Cutoff radius to apply to the atom centered symmetry function. Default is 12.

**n\_D\_inter** [int, optional] Number of dimensions to use for each type of pair-wise interatomic interaction found in the structure. Default is 12.

**init\_scheme** [str, optional] Can be centered or shifted, as described in Gastegger et al. 2018. Default is shifted.

**eta\_range** [list, optional] List of two floats which define the range for eta parameter in Gastegger et al. 2018. Default is [0.05, 0.5].

**Rs\_range** [list, optional] List of two floats which define the range for Rs parameter in Gastegger et al. 2018. Default is [[0.1, 12].

**pdist\_distance\_type** [str, optional] Input parameter for the pdist function. Default is Euclidean.

**Relax\_Single\_Molecule** (*comm, world\_comm, MPI\_ANY\_SOURCE, num\_replicas*)

Calls run\_fhi\_aims\_batch using the provided single molecule path.

### Configuration File Options

See `Run_FHI_Aims_Batch()`

**Run\_FHI\_Aims\_Batch** (*comm, world\_comm, MPI\_ANY\_SOURCE, num\_replicas*)

Runs FHI-aims calculations on a pool of structures using num\_replicas.

### Configuration File Options

**verbose** [bool] Controls verbosity of output.

**energy\_name** [str] Property name which the calculated energy will be stored with in the Structure file.

**output\_dir** [str] Path to the directory where the output structure file will be saved.

**aims\_output\_dir** [str] Path where the aims calculation will take place.

**aims\_lib\_dir** [str, inferred] Path to the location of the directory containing the FHI-aims library file.

**molecule\_path** [str] Path to the geometry.in file of the molecule to be calculated if called using harris\_single\_molecule\_prep or relax\_single\_molecule.

**structure\_dir** [str, inferred] Path to the directory of structures to be calculated if calculation was called not using harris\_single\_molecule\_prep or relax\_single\_molecule.

**Z** [int, inferred] Number of molecules per cell.