
ESPEI Documentation

Release 0.1.3

Brandon Bocklund

Aug 17, 2017

1	Installation	3
2	Usage	5
2.1	Full run	5
2.2	Single-phase only	5
2.3	Multi-phase only	6
2.4	Customization	6
2.5	FAQ	6
2.5.1	Q: There is an error in my JSON files	6
2.5.2	Q: How do I analyze my results?	6
2.5.3	Q: Can I run ESPEI on a supercomputer supporting MPI?	7
3	Module Hierarchy	9
4	License	11
4.1	Writing input files	11
4.1.1	JSON Format	11
4.1.2	Phase Descriptions	11
4.1.3	Single-phase Data	13
4.1.4	Multi-phase Data	15
4.2	What's New	15
4.2.1	0.2.1 (2017-08-17)	15
4.2.2	0.2 (2017-08-15)	15
4.2.3	0.1.5 (2017-08-02)	16
4.2.4	0.1.4 (2017-07-24)	16
4.2.5	0.1.3 (2017-06-23)	16
4.2.6	0.1.2 (2017-06-23)	16
4.2.7	0.1.1 (2017-06-23)	16
4.2.8	0.1 (2017-06-23)	16
4.3	API Documentation	17
4.3.1	espei package	17
4.3.1.1	Subpackages	17
4.3.1.2	Submodules	17
4.3.1.3	espei.core_utils module	17
4.3.1.4	espei.datasets module	17
4.3.1.5	espei.paramselect module	17
4.3.1.6	espei.plot module	17

4.3.1.7	espei.run_espei module	17
4.3.1.8	espei.utils module	17
4.3.1.9	Module contents	17

5 Indices and tables **19**

ESPEI, or Extensible Self-optimizing Phase Equilibria Infrastructure, is a tool for automated thermodynamic database development within the CALPHAD method.

The ESPEI package is based on a fork of [pycalphad-fitting](#) and uses [pycalphad](#) for calculating Gibbs free energies of thermodynamic models. The implementation for ESPEI involves first fitting single-phase data by calculating parameters in thermodynamic models that are linearly described by the single-phase input data. Then Markov Chain Monte Carlo (MCMC) is used to optimize the candidate models from the single-phase fitting to multi-phase zero-phase fraction data. Single-phase and multi-phase fitting methods are described in Chapter 3 of [Richard Otis's thesis](#).

The benefit of this approach is the automated, simultaneous fitting for many parameters that yields uncertainty quantification, as shown in Otis and Liu [High-Throughput Thermodynamic Modeling and Uncertainty Quantification for ICME. Jom 69, \(2017\)](#).

The name and idea of ESPEI are originally based off of Shang, Wang, and Liu, [ESPEI: Extensible, Self-optimizing Phase Equilibrium Infrastructure for Magnesium Alloys Magnes. Technol. 2010 617-622 \(2010\)](#).

CHAPTER 1

Installation

Creating a virtual environment is highly recommended. You can install ESPEI from PyPI

```
pip install espei
```

or install in develop mode from source

```
git clone https://github.com/phasesresearchlab/espei.git
cd espei
pip install -e .
```


Run `espei -h` to see the options in the command utility.

ESPEI has two different fitting modes: single-phase and multi-phase fitting. You can run either of these modes or both of them sequentially.

To run either of the modes, you need to have a fit settings file that describes the phases in the system using the standard CALPHAD approach within the compound energy formalism. You also need to describe the data to fit. You will need single-phase and multi-phase data for a full run. Fit settings and all datasets are stored as JSON files and described in detail at the [Writing input files](#) page. All of your input datasets should be validated by running `espei --check-datasets my-input-datasets`, where `my-input-datasets` is a folder of all your JSON files.

The main output result is going to be a database (defaults to `out.tdb`) and an array of the steps in the MCMC chain (defaults to `chain.txt`).

Full run

A minimal run of ESPEI with single phase fitting and MCMC fitting would involve setting these two files

```
espei --datasets=my-dataset-folder --fit-settings=my-input.json
```

Single-phase only

If you have only heat capacity, entropy and enthalpy data and mixing data (e.g. from first-principles), you may want to see the starting point for your MCMC calculation. To do this, simply pass the `--no-mcmc` flag to ESPEI

```
espei --no-mcmc --datasets=my-dataset-folder --fit-settings=my-input.json
```

Multi-phase only

If you have a database already and just want to do a multi-phase fitting, you can specify a starting TDB file with

```
espei --datasets=my-dataset-folder --fit-settings=my-input.json --input-tdb=my-
↳starting-database.tdb
```

The TDB file you input must have all of the degrees of freedom you want as FUNCTIONS with names beginning with VV.

Customization

In all cases, ESPEI lets you control certain aspects of your calculations from the command line. Some useful options are

- `verbose` (or `-v`) controls the logging level. Default is Warning. Using `verbose once` gives more detail (Info) and twice even more (Debug)
- `tracefile` lets you set the output trace of the chain to any name you want. The default is `chain.txt`.
- `probfile` lets you set the output log-probability of each step in the chain to any name you want. The default is `lnprob.txt`.
- `output-tdb` sets the name of the TDB output at the end of the run. Default is `out.tdb`.
- `input-tdb` is for setting input TDBs. This will skip single phase fitting and fit all parameters defined as FUNCTIONS with names starting with VV.
- `no-mcmc` will do single-phase fitting only. Default is to perform MCMC fitting.
- `mcmc-steps` sets the number of MCMC steps. The default is 1000.
- `save-interval` controls the interval for saving the MCMC chain. The default is 100 steps.

Run `espei -h` to see all of the configurable options.

FAQ

Q: There is an error in my JSON files

A: Common mistakes are using single quotes instead of the double quotes required by JSON files. Another common source of errors is misaligned open/closing brackets.

To find the offending files, you can rename the datasets to anything not ending in `.json`, such as `my_datasets.json.disabled`. The renamed files will be ignored and it allows you to track down any problematic files.

Q: How do I analyze my results?

A: By default, ESPEI will create `chain.txt` and `lnprob.txt` for the MCMC chain at the end of your run and according to the save interval (defaults to every 100 iterations). These are created from arrays via `numpy.savetxt` and can thus be loaded with `numpy.loadtxt()`. Note that the arrays are preallocated with zeros. These filenames and settings (e.g. save interval) can be changed using the command line options, see `espei -h`. You can then use these chains and corresponding log-probabilities to make corner plots, calculate autocorrelations, find optimal parameters for databases, etc.. Finally, you can use `py:mod:espei.plot` functions such as `multiplot` to plot phase

diagrams with your input equilibria data and `plot_parameters` to compare single-phase data (e.g. formation and mixing data) with the properties calculated with your database.

Q: Can I run ESPEI on a supercomputer supporting MPI?

A: Yes! ESPEI has MPI support. Currently only single node processing is supported, but fixes are coming soon to support multiple nodes. To use ESPEI with MPI, you simply call ESPEI in the same way as above with `mpirun` or whichever MPI software you use. You also must indicate to ESPEI that it should create an MPI scheduler by passing the option `--scheduler='MPIPool'` to ESPEI.

CHAPTER 3

Module Hierarchy

- `fit.py` is the main entry point
- `paramselect.py` is where all of the fitting happens. This is the core.
- `core_utils.py` contains specialized utilities for ESPEI.
- `utils.py` are utilities with reuse potential outside of ESPEI.
- `plot.py` holds plotting functions

ESPEI is MIT licensed. See LICENSE.

Writing input files

JSON Format

ESPEI has a single input style in JSON format that is used for all data entry. Single-phase and multi-phase input files are almost identical, but detailed descriptions and key differences can be found in the following sections. For those unfamiliar with JSON, it is fairly similar to Python dictionaries with some rigid requirements

- All string quotes must be double quotes. Use "key" instead of 'key'.
- Numbers should not have leading zeros. 00.123 should be 0.123 and 012.34 must be 12.34.
- Lists and nested key-value pairs cannot have trailing commas. {"nums": [1, 2, 3,], } is invalid and should be {"nums": [1, 2, 3]}.

These errors can be challenging to track down, particularly if you are only reading the JSON error messages in Python. A visual editor is encouraged for debugging JSON files such as [JSONLint](#). A quick reference to the format can be found at [Learn JSON in Y minutes](#).

ESPEI has support for checking all of your input datasets for errors, which you should always use before you attempt to run ESPEI. This error checking will report all of the errors at once and all errors should be fixed. Errors in the datasets will prevent fitting. To check the datasets at path `my-input-data/` you can run `espei --check-datasets my-input-data`.

Phase Descriptions

The JSON file for describing CALPHAD phases is conceptually similar to a setup file in Thermo-Calc's PARROT module. At the top of the file there is the `refdata` key that describes which reference state you would like to choose. Currently the reference states are strings referring to dictionaries in `pycalphad.refdata` only "SGTE91" is implemented.

Each phase is described with the phase name as they key in the dictionary of phases. The details of that phase is a dictionary of values for that key. There are 4 possible entries to describe a phase: `sublattice_model`, `sublattice_site_ratios`, `equivalent_sublattices`, and `aliases`. `sublattice_model` is a list of lists, where each internal list contains all of the components in that sublattice. The BCC_B2 sublattice model is `[["AL", "NI", "VA"], ["AL", "NI", "VA"], ["VA"]]`, thus there are three sublattices where the first two have Al, Ni, and vacancies. `sublattice_site_ratios` should be of the same length as the sublattice model (e.g. 3 for BCC_B2). The sublattice site ratios can be fractional or integers and do not have to sum to unity.

The optional `equivalent_sublattices` key is a list of lists that describe which sublattices are symmetrically equivalent. Each sub-list in `equivalent_sublattices` describes the indices (zero-indexed) of sublattices that are equivalent. For BCC_B2 the equivalent sublattices are `[[0, 1]]`, meaning that the sublattice at index 0 and index 1 are equivalent. There can be multiple different sets (multiple sub-lists) of equivalent sublattices and there can be many equivalent sublattices within a sublattice (see FCC_L12). If no `equivalent_sublattice` key exists, it is assumed that there are none.

Finally, the `aliases` key is used to refer to other phases that this sublattice model can describe when symmetry is accounted for. Aliases are used here to describe the BCC_A2 and FCC_A1, which are the disordered phases of BCC_B2 and FCC_L12, respectively. Notice that the aliased phases are not otherwise described in the input file. Multiple phases can exist with aliases to the same phase, e.g. FCC_L12 and FCC_L10 can both have FCC_A1 as an alias.

```
{
  "refdata": "SGTE91",
  "components": ["AL", "NI", "VA"],
  "phases": {
    "LIQUID": {
      "sublattice_model": ["AL", "NI"],
      "sublattice_site_ratios": [1]
    },
    "BCC_B2": {
      "aliases": ["BCC_A2"],
      "sublattice_model": ["AL", "NI", "VA"], ["AL", "NI", "VA"], ["VA"],
      "sublattice_site_ratios": [0.5, 0.5, 1],
      "equivalent_sublattices": [[0, 1]]
    },
    "FCC_L12": {
      "aliases": ["FCC_A1"],
      "sublattice_model": ["AL", "NI"], ["AL", "NI"], ["AL", "NI"], ["AL", "NI"], [
↪ "VA"]],
      "sublattice_site_ratios": [0.25, 0.25, 0.25, 0.25, 1],
      "equivalent_sublattices": [[0, 1, 2, 3]]
    },
    "AL3NI1": {
      "sublattice_site_ratios": [0.75, 0.25],
      "sublattice_model": ["AL"], ["NI"]
    },
    "AL3NI2": {
      "sublattice_site_ratios": [3, 2, 1],
      "sublattice_model": ["AL"], ["AL", "NI"], ["NI", "VA"]
    },
    "AL3NI5": {
      "sublattice_site_ratios": [0.375, 0.625],
      "sublattice_model": ["AL"], ["NI"]
    }
  }
}
```


Single-phase Data

Two example of ESPEI input file for single-phase data follow. The first dataset has some data for the formation heat capacity for BCC_B2.

The `components` and `phases` keys simply describe those found in this entry. Use the `reference` key for book-keeping the source of the data. In `solver` the sublattice configuration and site ratios are described for the phase.

`sublattice_configurations` is a list of different configurations, that should correspond to the sublattices for the phase descriptions. Non-mixing sublattices are represented as a string, while mixing sublattices are represented as a lists. Thus an endmember for BCC_B2 (as in this example) is `["AL", "NI", "VA"]` and if there were mixing (as in the next example) it might be `["AL", ["AL", "NI"], "VA"]`. Mixing also means that the `sublattice_occupancies` key must be specified, but that is not the case in this example. Regardless of whether there is mixing or not, the length of this list should always equal the number of sublattices in the phase, though the sub-lists can have mixing up to the number of components in that sublattice. Note that the `sublattice_configurations` is a *list* of these lists. That is, there can be multiple sublattice configurations in a single dataset. See the second example in this section for such an example.

The `conditions` describe temperatures (T) and pressures (P) as either scalars or one-dimensional lists. Most important to describing data are the `output` and `values` keys. The type of quantity is expressed using the `output` key. This can in principle be any thermodynamic quantity, but currently only CPM*, SM*, and HM* (where * is either nothing, `_MIX` or `_FORM`) are supported. Support for changing reference states planned but not yet implemented, so all thermodynamic quantities must be formation quantities (e.g. `HM_FORM` or `HM_MIX`, etc.).

The `values` key is the most complicated and care must be taken to avoid mistakes. `values` is a 3-dimensional array where each value is the `output` for a specific condition of pressure, temperature, and sublattice configurations from outside to inside. Alternatively, the size of the array must be `(len(P), len(T), len(subl_config))`. In the example below, the shape of the `values` array is (1, 12, 1) as there is one pressure scalar, one sublattice configuration, and 12 temperatures. The formatting of this can be tricky, and it is suggested to use a NumPy array and reshape or add axes using `np.newaxis` indexing.

```
{
  "reference": "Yi Wang et al 2009",
  "components": ["AL", "NI", "VA"],
  "phases": ["BCC_B2"],
  "solver": {
    "sublattice_site_ratios": [0.5, 0.5, 1],
    "sublattice_configurations": [["AL", "NI", "VA"]],
    "comment": "NiAl sublattice configuration (2SL)"
  },
  "conditions": {
    "P": 101325,
    "T": [ 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 110]
  },
  "output": "CPM_FORM",
  "values": [[ [ 0 ],
               [-0.0173 ],
               [-0.01205],
               [ 0.12915],
               [ 0.24355],
               [ 0.13305],
               [-0.1617 ],
               [-0.51625],
               [-0.841  ],
               [-1.0975 ],
               [-1.28045],
               [-1.3997 ]]]
}
```

In the second example below, there is formation enthalpy data for multiple sublattice configurations. All of the keys and values are conceptually similar. Here, instead of describing how the output quantity changes with temperature or pressure, we are instead only comparing HM_FORM values for different sublattice configurations. The key differences from the previous example are that there are 9 different sublattice configurations described by `sublattice_configurations` and `sublattice_occupancies`. Note that the `sublattice_configurations` and `sublattice_occupancies` should have exactly the same shape. Sublattices without mixing should have single strings and occupancies of one. Sublattices that do have mixing should have a site ratio for each active component in that sublattice. If the sublattice of a phase is ["AL", "NI", "VA"], it should only have two occupancies if only ["AL", "NI"] are active in the sublattice configuration.

The last difference to note is the shape of the values array. Here there is one pressure, one temperature, and 9 sublattice configurations to give a shape of (1, 1, 9).

```
{
  "reference": "C. Jiang 2009 (constrained SQS)",
  "components": ["AL", "NI", "VA"],
  "phases": ["BCC_B2"],
  "solver": {
    "sublattice_occupancies": [
      [1, [0.5, 0.5], 1],
      [1, [0.75, 0.25], 1],
      [1, [0.75, 0.25], 1],
      [1, [0.5, 0.5], 1],
      [1, [0.5, 0.5], 1],
      [1, [0.25, 0.75], 1],
      [1, [0.75, 0.25], 1],
      [1, [0.5, 0.5], 1],
      [1, [0.5, 0.5], 1]
    ],
    "sublattice_site_ratios": [0.5, 0.5, 1],
    "sublattice_configurations": [
      ["AL", ["NI", "VA"], "VA"],
      ["AL", ["NI", "VA"], "VA"],
      ["NI", ["AL", "NI"], "VA"],
      ["NI", ["AL", "NI"], "VA"],
      ["AL", ["AL", "NI"], "VA"],
      ["AL", ["AL", "NI"], "VA"],
      ["NI", ["AL", "VA"], "VA"],
      ["NI", ["AL", "VA"], "VA"],
      ["VA", ["AL", "NI"], "VA"]
    ],
    "comment": "BCC_B2 sublattice configuration (2SL)"
  },
  "conditions": {
    "P": 101325,
    "T": 300
  },
  "output": "HM_FORM",
  "values": [[[-40316.61077, -56361.58554,
               -49636.39281, -32471.25149, -10890.09929,
               -35190.49282, -38147.99217, -2463.55684,
               -15183.13371]]]]
}
```

Multi-phase Data

The difference between single- and multi-phase is data is in the absence of the `solver` key, since we are no longer concerned with individual site configurations, and the `values` key where we need to represent phase equilibria rather than thermodynamic quantities. Notice that the type of data we are entering in the `output` key is ZPF (zero-phase fraction) rather than `CP_FORM` or `H_MIX`. Each entry in the ZPF list is a list of all phases in equilibrium, here `[["AL3NI2", ["NI"], [0.4083]], ["BCC_B2", ["NI"], [0.4340]]]` where each phase entry has the name of the phase, the composition element, and the composition of the tie line point. If there is no corresponding tie line point, such as on a liquidus line, then one of the compositions will be `null`: `[["LIQUID", ["NI"], [0.6992]], ["BCC_B2", ["NI"], [null]]]`. Three- or n-phase equilibria are described as expected: `[["LIQUID", ["NI"], [0.752]], ["BCC_B2", ["NI"], [0.71]], ["FCC_L12", ["NI"], [0.76]]]`.

Note that for higher-order systems the component names and compositions are lists and should be of length $c-1$, where c is the number of components.

```
{
  "components": ["AL", "NI"],
  "phases": ["AL3NI2", "BCC_B2"],
  "conditions": {
    "p": 101325,
    "T": [1348, 1176, 977]
  },
  "output": "ZPF",
  "values": [
    [
      ["AL3NI2", ["NI"], [0.4083]], ["BCC_B2", ["NI"], [0.4340]]],
      [
        ["AL3NI2", ["NI"], [0.4114]], ["BCC_B2", ["NI"], [0.4456]]],
        [
          ["AL3NI2", ["NI"], [0.4114]], ["BCC_B2", ["NI"], [0.4532]]]
      ]
    ],
  "reference": "37ALE"
}
```

What's New

0.2.1 (2017-08-17)

Fixes to the 0.2 release plotting interface

- `multiplot` is renamed from `multi_plot`, as in docs.
- Fixed an issue where phases in datasets, but not in equilibrium were not plotted by `dataplot` and raised an error.

0.2 (2017-08-15)

- New `multiplot` interface for convenient plotting of phase diagrams + data. `dataplot` function underlies key data plotting features and can be used with `eqplot`. See their API docs for examples. Will break existing code using `multiplot`.
- MPI support for local/HPC runs. Only single node runs are explicitly supported currently. Use `--scheduler='MPIPool'` command line option. Requires `mpi4py`.
- Default debug reporting of acceptance ratios
- Option (and default) to output the log probability array matching the trace. Use `--probfile` option to control.
- Optimal parameters are now chosen based on lowest error in chain.

- Bug fixes including
 - py2/3 compatibility
 - unicode datasets
 - handling of singular matrix errors from pycalphad's `equilibrium`
 - reporting of failed conditions

0.1.5 (2017-08-02)

- Significant error checking of JSON inputs.
- Add new `--check-datasets` option to check the datasets at path. It should be run before you run ESPEI fittings. All errors must be resolved before you run.
- Move the `espei` script module from `fit.py` to `run_espei.py`.
- Better docs building with mocking
- Google docstrings are now NumPy docstrings

0.1.4 (2017-07-24)

- Documentation improvements for usage and API docs
- Fail fast on JSON errors

0.1.3 (2017-06-23)

- Fix bad version pinning in `setup.py`
- Explicitly support Python 2.7

0.1.2 (2017-06-23)

- Fix dask incompatibility due to new API usage

0.1.1 (2017-06-23)

- Fix a bug that caused logging to raise if bokeh isn't installed

0.1 (2017-06-23)

ESPEI is now a package! New features include

- Fork <https://github.com/richardotis/pycalphad-fitting>
- Use `emcee` for MCMC fitting rather than `pymc`
- Support single-phase only fitting
- More control options for running ESPEI from the command line
- Better support for incremental saving of the chain

- Control over output with logging over printing
- Significant code cleanup
- Better usage documentation

API Documentation

espei package

Subpackages

espei.tests package

Submodules

espei.tests.test_datasets module

espei.tests.test_paramselect module

espei.tests.test_utils module

Module contents

Submodules

espei.core_utils module

espei.datasets module

espei.paramselect module

espei.plot module

espei.run_espei module

espei.utils module

Module contents

CHAPTER 5

Indices and tables

- `genindex`
- `modindex`
- `search`