

Recursively Embedded Atom Neural Network Package

-B. Jiang Group (2023)

Introduction

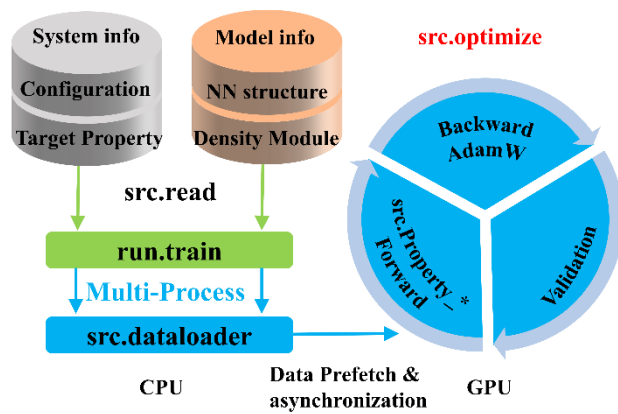
Recursively embedded atom neural network (REANN) is a PyTorch-based end-to-end multi-functional Deep Neural Network Package for Molecular, Reactive and Periodic Systems. Currently, REANN can be used to train interatomic potentials, dipole moments, transition dipole moments, and polarizabilities. Taking advantage of Distributed DataParallel features embedded in PyTorch, the training process is highly parallelized on both GPU and CPU. For the convenience of MD simulation, an interface to LAMMPS has been constructed by creating a new pair_style invoking this representation for highly efficient MD simulations.

Training Workflow

The training process can be divided into four parts: information loading, initialization, dataloader and optimization. First, the "src.read" will load the information about the systems and NN structures from the dataset and input files ("input_nn" and "input_density") respectively. Second, the "run.train" module utilizes the loaded information to initialize various classes, including property calculator, dataloader, and optimizer. For each process, an additional thread will be activated in the "src.dataloader" module to prefetch data from CPU to GPU in an asynchronous manner. Meanwhile, the optimization will be activated in the "src.optimize" module once the first set of data is transferred to the GPU. During optimization, a learning rate scheduler,

namely "ReduceLROnPlateau" provided by PyTorch, is used to decay the learning rate.

Training is stopped when the learning rate drops below 1×10^{-5} and the model that performs best on the validation set is saved for further investigation.



How to Use REANN Package

Users can employ geometries, energies, and atomic force vectors (or some other physical properties which are invariant under rigid translation, rotation, and permutation of identical atoms and their corresponding gradients) to construct a model to represent these physical properties via this package. There are three routines to use this package:

1. Prepare the environment
2. Prepare data
3. Construct a model

1. Prepare the environment

The REANN Package is built based on PyTorch and uses the "opt_einsum" package for optimizing einsum-like expressions frequently used in the calculation of the embedded density. In order to run the REANN package, users need to install PyTorch (version: 2.0.0) based on the instructions on the PyTorch official website (<https://pytorch.org/get-started/locally/>) and the package named opt_einsum (<https://optimized-einsum.readthedocs.io/en/stable/>).

2. Prepare data

There are two directories that users need to prepare, namely, "train" and "val", each of which includes a file "configuration" used to preserve the required information including lattice parameters, periodic boundary conditions, configurations, energy and

atomic forces (if needed). For example, users want to represent the methane system (CH₄) that has available atomic forces. The file "configuration" should be written in the following format.

```

point= 1
100e0 0e0 0e0
0e0 100e0 0e0
0e0 0e0 100e0
pbc 0 0 0
C 12.001 0.00000000 0.00000000 0.00000000 8.06859172166334 10.80335659070366 21.97603328930107
H 1.008 -1.17850750 0.26261528 -2.14262002 0.1165638588531 0.63359209279172 1.13009238957782
H 1.008 -0.32352777 -0.29187167 -0.67955432 -10.26965498566280 -11.38811672778598 -22.50791198138647
H 1.008 -2.26940807 0.81621901 0.06640190 1.34987449972433 -1.23095526897855 -0.28884762833028
H 1.008 -2.28911976 -1.50408332 0.67309787 0.73426882424044 1.18200349980833 -0.30898348884089
abprop: -1082.7845870920069
point= 2
100e0 0e0 0e0
0e0 100e0 0e0
0e0 0e0 100e0
pbc 0 0 0
C 12.001 0.00000000 0.00000000 0.00000000 -13.26713992103983 -3.24866948511857 0.53954932468248
H 1.008 -0.73133356 2.41878753 -0.94382677 0.34069960316167 -1.25051005991520 0.22970297323531
H 1.008 -1.75936836 -0.10797807 -0.07477593 4.89799124779970 0.61639911827382 0.18032182964539
H 1.008 0.91619771 0.20734120 -0.14431795 8.18578214187139 2.34304118663493 -1.45346429235577
H 1.008 0.25136836 -2.36005983 -1.13685075 -0.15740660537617 1.53965490790367 0.50385159822794
abprop: -1086.18078403197

```

The first line can be an arbitrary character other than a blank line. The next three lines are the lattice vectors defining the unit cell of the system. Users can copy the lattice vectors in POSCAR to this file. The fifth line is used to enable(1)/disable(0) the periodic boundary conditions in each direction. In this example, methane is not a periodic system, the fifth line should be “pbc 0 0 0”. For some gas-surface systems, only the x-y plane is periodic and the corresponding fifth line is “pbc 1 1 0”. Following N lines (N is the number of atoms in the system, here is 5): the columns from the left to right represent the atomic name, relative atomic mass, coordinates(x, y, z) of the geometry, atomic force vectors (if the force vector is not incorporated in the training, these three columns can be omitted). Next line: Start with "abprop:" and then follow by the target property (energy/dipole/polarizability).

3. Construct a model

In the former section “Prepare data”, file configuration has informed the package the information about the system. However, some hyperparameters about the embedded

density, training algorithms and architectures of neural networks are essential for obtaining an exact representation. Next, we will give an example to explain these hyperparameters in detail. There are two input files named "input_nn" and "input_density" saved in the "para" directory in the work directory.

Parameters in “input_nn” and “input_density” file

Although there are dozens of parameters in "input_density" and "input_nn", users can obtain a pretty good accuracy with the default set (the value shown in this manual) except for some parameters marked in red and “required” about the systems for example the "atomtype" that tells the elements involved in the system. The order of parameters is arbitrary and users can give parameters in any order they prefer. If users do not set the parameters "input_density" or "input_nn", code will automatically accept default values.

File “input_nn”

1. `start_table = 0` # required parameters type: integer 0/1/2/3/4;

(If the atomic force vectors are used in NN training. “1”: energies and force vectors will be used; “0”: only energies will be used; “2”: permanent dipole moment “3”: transition dipole moment; “4”: polarizability.)

2. `table_coor = 0` # required parameters type: integer 0/1;

(The form of the coordinates in “1” file, “0”: Cartesian coordinates. “1”: Fraction coordinates.)

3. `table_init = 0` # type: integer 0/1;

(Initialize parameters for a model, “0”: general method; “1”: restart from a previous training)

4. `nblock = 1` # type: integer

(Number of residual NN blocks)

5. `nl=[64,64]` # type: list [n_1, n_2, \dots, n_p] n_p : integer

(A list holds the number of neurons in each hidden layer of residual blocks. Each residual NN block will have the $n_1 \times n_2 \times \dots \times n_p \times n_1$. The last n_1 is inserted for identity mapping. The number of neuron should be a power of 2 for a better efficiency.)

6. `dropout_p=[0.0, 0.0]` # type: list [p_1, p_2, \dots] p_i : real number between 0~1

(A list holds the Bernoulli distribution probability of disabling the neuron in each hidden layer used in the dropout for regularization. The length of the list must be larger than or equal to the number of hidden layers.)

7. `table_norm=True` # type: bool True/False

(A switch enables (True)/disables (False) the layer normalization in hidden layers, which can help with regularization and acceleration of convergence.)

8. `re_coeff=0.0` # type: real number between 0~1

(L2 regularization coefficient. Default is 0; usually, a very small value or 0 are preferred)

Note: The dropout is a frequently-used technology in machine learning (ML) and has made a huge success. However, in many traditional ML models, only the outputs serve as the targets, which differ in the construction of potential energy surfaces. The minus gradients of the energy (atomic forces) are always used as an important term to improve the accuracy and reduce the required number of configurations. Meanwhile, the

introduction of atomic forces can act as regularization and reduce overfitting. So if users have a sufficient number of datasets including atomic forces and an appropriate size of NN, not that many or strong regularization needs to be included in models (including the L2 regularization coefficients).

9. Epoch=10000 # type: integer

(Maximum epochs in one fitting.)

10. patience_epoch=50 # type: integer

11. decay_factor=0.5 # type: real number between 0~1

(The learning rate will be decayed by a factor of “decay_factor” if the loss is not improved after “patience_epoch” epochs. $\text{new_lr} = \text{lr} * \text{decay_factor}$.)

12. print_epoch=1 # type: integer

(Calculate and print the training and validation rmse every “print_epoch” epochs.)

13. ratio = 0.9 # type: real between 0~1

(If an empty "validation" folder is provided, then the train will be divided into two parts (ratio:1-ratio) used for training and validation respectively.)

14. start_lr=0.01 # type: real number

(Initial learning rate. A comparatively high start_lr and small patience_epoch are preferred for better efficiency and accuracy.)

15. end_lr=0.00001 # type: real number

(Final learning rate)

(The learning rate will decrease from the “start_lr” to the “end_lr” by a factor of “decay_factor” according to the learning rate scheduler “ReduceLROnPlateau”

provided by PyTorch.)

16. `batchsize_train=64` # required parameters type: integer

17. `batchsize_val=128` # required parameters type: integer

(Number of configurations used in each batch for train (`batchsize_train`) and validation (`batchsize_val`). Note, this "`batchsize_train`" is a key parameter concerned with efficiency. Normally, a large enough value is given to achieve high usage of the GPU and lead to higher efficiency in training if you have sufficient data. However, for small training data, a large "`batchsize`" can lead to a decrease in accuracy, probably owing to the decrease in the number of gradient descents during each epoch. The decrease in accuracy may be compensated by more epochs (increase the "`patience_epoch`") or a larger learning rate. Some detailed testing is required here to achieve a balance of accuracy and efficiency in training. The value of "`batch_val`" has no effect on accuracy, and thus a larger value is preferred.)

18. `e_coeff=0.1` # type: real number

(Weight of energy)

19. `init_f=10` # type: real number

(Initial weight of atomic forces)

20. `final_f=0.5` # type: real number

(Final weight of atomic forces. The weight of the atomic forces decays with the decrease in the learning rate. Usually, a large initial weight of atomic forces is beneficial during the training.)

21. `activate = 'Relu_like'` # type: string options Relu_like/Tanh_like

(activation functions)

22. `queue_size=10` # type: integer

(The number of batch data prefetched to the GPU to ensure good accelerated performance.)

23. `DDP_backend="nccl"` # type: string options nccl/gloo

(The built-in backends included in PyTorch distributed. The nccl is preferred for distributed GPU training.)

24. `find_unused=False` # type: bool True/False

(False is the default value. If there is an error reporting unused parameters, change find_unused to True.)

25. `dtype='float32` # type: string "float32" / "float64"

26. `folder = "./"`

(The path save the directories "train" and "val")

The following parameters have a prefix "oc_" and play the same role as defined in the former that determines the NN structure. The difference is that the former represents the mapping from density to atomic energy, and the latter represents the mapping from density to orbital coefficients.

27. `oc_loop = 1` # type: integer

(Number of iterations used to represent the orbital coefficients.)

28. `oc_nl = [128,128]` # type: list [n₁, n₂, ..., n_p] n_p: integer

(Same as nl.)

29. `oc_nblock = 1` # type: integer

(Same as nblock)

30. oc_dropout_p=[0.0, 0.0] # type: list[p1,p2,...] pi: real number between 0~1

(Same as dropout_p)

31. oc_activate = 'Relu_like' # type: string options "Relu_like" / "Tanh_like"

(Same as activate)

32. oc_table_norm=False # # type: bool True/False

(Same as table_norm)

File "input_density"

1. cutoff = 4.5 # type: real number

(Cutoff distances)

2. neigh_atoms = 150 # type: integer

(Maximum number of neighbor atoms within the cutoff)

3. atomtype= ['O', 'H',] # required parameters type: list [elmenet1, element2,...]

(The involved elements in the systems)

4. nipsin= 2 # type: integer

(Maximal angular momenta determine the orbital type (s, p, d ..))

5. nwave=8 # type: integer

(Number of radial Gaussian functions. This number should be a power of 2 for better efficiency.)

With all needed files prepared, users can execute the training in a single node with multiple GPUs with the directive:

```
"python3 -m torch.distributed.run --nproc_per_node=$NPROC_PER_NODE --
```

`nnodes=1 --standalone $path "`, where "path" is the path that code locate in and the `NPROC_PER_NODE` is the number of GPUs in each node.

The multi-nodes multi-GPUs can be launched with a similar directive which need to specify the node and the port of the master node. Details can be found in <https://pytorch.org/docs/stable/elastic/run.html?highlight=distributed%20run#module-torch.distributed.run>.

In the training process, some iterative information will output to the file "nn.err".

- nn.err

First, some information about the model will be printed out.

```
REANN Package used for fitting energy and tensorial Property
2021-11-27-22_06_17
Epoch= 0 learning rate 1.000000e-03 train error: 17.27729 0.75929 test error: 6.40291 2.13923
Epoch= 1 learning rate 1.000000e-03 train error: 2.73391 0.16909 test error: 36.40526 2.13042
Epoch= 2 learning rate 1.000000e-03 train error: 2.05831 0.15300 test error: 93.88305 2.06302
Epoch= 3 learning rate 1.000000e-03 train error: 1.82869 0.14504 test error: 120.40647 1.83669
Epoch= 4 learning rate 1.000000e-03 train error: 1.88152 0.13968 test error: 97.24050 1.49304
Epoch= 5 learning rate 1.000000e-03 train error: 1.30521 0.13535 test error: 63.27221 1.17121
Epoch= 6 learning rate 1.000000e-03 train error: 1.54117 0.13272 test error: 37.29077 0.88717
Epoch= 7 learning rate 1.000000e-03 train error: 1.03487 0.12831 test error: 20.09610 0.64664
Epoch= 8 learning rate 1.000000e-03 train error: 1.48432 0.12623 test error: 9.96114 0.46306
Epoch= 9 learning rate 1.000000e-03 train error: 8.90729 0.15197 test error: 4.93163 0.33503
Epoch= 10 learning rate 1.000000e-03 train error: 0.87571 0.11939 test error: 0.59845 0.24723
Epoch= 11 learning rate 1.000000e-03 train error: 1.42815 0.12158 test error: 1.59155 0.19093
Epoch= 12 learning rate 1.000000e-03 train error: 1.25193 0.11799 test error: 2.49767 0.15582
Epoch= 13 learning rate 1.000000e-03 train error: 1.31798 0.11820 test error: 3.06294 0.13667
Epoch= 14 learning rate 1.000000e-03 train error: 1.42203 0.11633 test error: 3.41752 0.12621
Epoch= 15 learning rate 1.000000e-03 train error: 3.33254 0.11732 test error: 3.53356 0.12068
Epoch= 16 learning rate 1.000000e-03 train error: 6.30027 0.11942 test error: 3.17197 0.11669
Epoch= 17 learning rate 1.000000e-03 train error: 1.02547 0.11050 test error: 2.56956 0.11348
Epoch= 18 learning rate 1.000000e-03 train error: 1.68313 0.11190 test error: 2.12365 0.11125
Epoch= 19 learning rate 1.000000e-03 train error: 2.75737 0.11214 test error: 1.62302 0.10964
Epoch= 20 learning rate 1.000000e-03 train error: 2.22338 0.11100 test error: 1.15830 0.10831
Epoch= 21 learning rate 1.000000e-03 train error: 34.61486 0.24439 test error: 5.58866 0.14059
Epoch= 22 learning rate 1.000000e-03 train error: 5.48195 0.19695 test error: 30.95137 0.39146
Epoch= 23 learning rate 1.000000e-03 train error: 0.90179 0.12597 test error: 51.08136 0.38373
```

After that, the learning rate, RMSE of energy and atomic force in each direction of each epoch will be printed out. Note: The RMSE of training is estimated from the average of each batch rather than the realistic RMSE. The parameters of the model with the smallest RMSE in energy will be saved to "REANN.pth". Then they will be converted to serialized models ("PES.pt") via the torchscript, which can be loaded in python/C++ environments used for high-performance inference.

MD Simulations with LAMMPS

As mentioned before, the program will save “LAMMPS.pt”, which are two serializable models (float/double) designed for the interface as a new “pair_style” named “reann” in the LAMMPS (version: 10 Feb 2021). To run LAMMPS based on the “reann” pair_style, users need to build the LAMMPS from source with CMake. First, users need to download the LAMMPS and copy the file in each folder located in (“lammps-interface/” (for eann when oc_loop=0) or “lammps-REANN-interface/” (for reann when oc_loop >0)) to the corresponding directory in LAMMPS (create folder if it does not exist). Then users need to download the Torch C++ frontend (Libtorch) from the PyTorch website consistent with the environments in the targeted device (CPU or GPU with different versions of CUDA and cudnn). In addition, the downloaded Libtorch needs to be uncompressed to “libtorch” folder, which is located in the same directory with the src of LAMMPS. Some attention should be taken here that the pre-cxx11 ABI Libtorch (circled in the following picture) should be adopted and the version of CUDA should be exactly consistent with the device environments supported in Libtorch, which means only the “CUDA 10.2” and “CUDA 11.1” are available.

PyTorch Build	Stable (1.8.1)		Preview (Nightly)	
Your OS	Linux	Mac	Windows	
Package	Conda	Pip	LibTorch	Source
Language	Python		C++ / Java	
Compute Platform	CUDA 10.2	CUDA 11.1	ROCm 4.0 (beta)	CPU
Run this Command:	<p>Download here (Pre-cxx11 ABI):</p> <p>https://download.pytorch.org/libtorch/cu102/libtorch-shared-with-deps-1.8.1%2Bcu102.zip</p> <p>Download here (cxx11 ABI):</p> <p>https://download.pytorch.org/libtorch/cu102/libtorch-cxx11-abi-shared-with-deps-1.8.1%2Bcu102.zip</p>			

Finally, users can generate the Makefile using the directive “sh bulid.sh” in the “build” folder and then “make” to obtain the executable file “lmp_mpi”, more details about building the lammmps with cmake can be found on the website of Lammmps. Now, MD simulations in Lammmps can obtain a pretty good efficiency. There is an example in the folder “lammmps-interface/example”. The (sequence number) of elements should begin with 1 not 0 and the sequence of elements should be the same as the "atomtype" list in "input_density". Note: The multi-process parallelization algorithm for REANN is less straightforward as in classical force fields or conventional locally atomic descriptor-based ML models. Currently, the MD simulation based on the REANN model can only perform on a single process (GPU/CPU). However, if you set oc_loop=0, REANN will degrade to eann. The MD simulations can run on multiple CPUs or GPUs (the number of GPUs should be equal to the number of processes) and obtain a pretty good acceleration.

Inference Based on ASE

In the “ASE” folder, there is a python script "ase_reann.py" used as an example to

calculate the energy and atomic forces by invoking the model save in “PES.pt”. **Note:** The “atomtype” script in the inference should be same with that in the “input_density”. In this interface, we use a cell-linked algorithm to construct the neighbor list by a high efficient Fortran implementation. Thus, the Fortran code should be compiled by f2py, which generates a dynamic link library by running a "run" script, which can be called by python.