

# Machine-learning accelerated geometry optimization in molecular simulation: Supporting information

Yilin Yang and John R. Kitchin\*

*Department of Chemical Engineering, Carnegie Mellon University*

Omar A. Jiménez-Negró

*Department of Chemical Engineering,*

*University of Puerto Rico-Mayagüez,*

*Mayagüez, PR 00681, Puerto Rico, USA*

(Dated: March 6, 2021)

## CONTENTS

Introduction	2
Example of NN_ensemble_relaxer utilization	2
Installation	2
Instantiate NN_ensemble_relaxer and run the relaxation	3
Figures for configurations used as the example in the manuscript	3
Code to access the datasets for reproducing the main figures in the manuscript	4
Acceleration of the geometry optimization for Acrolein/AgPd	4
Performance for various systems	6
Gaussian Process Regression	9
References	10

## INTRODUCTION

There are three parts in this supporting information:

- Example to use NN\_ensemble\_relaxer for molecule geometry optimization
- Figures for configurations used as the example in the manuscript
- Code to access the datasets for reproducing the main figures in the manuscript

## EXAMPLE OF NN\_ENSEMBLE\_RELAXER UTILIZATION

### Installation

We first download the NN\_ensemble\_relaxer repo from the github and compile the required file.

---

```
1 git clone https://github.com/yilinyang1/NN_ensemble_relaxer.git
2 cd utils && python libsymb_builder
```

---

## Instantiate NN\_ensemble\_relaxer and run the relaxation

Assume we have a ASE database called `db` that contains the configurations that need to be optimized. We can just feed it into the `Ensemble_relaxer` to conduct the relaxation.

---

```
1 from nn_optimize import Ensemble_Relaxer
2
3 # NN hyperparameters
4 nn_params = {'layer_nodes': [40, 40], 'activations': ['tanh', 'tanh'], 'lr': 1}
5
6 # confidence coefficients used to control to what extent we trust the NN model
7 alpha = 2.0
8
9 # feed ASE database db, set ground truth calculator,
10 # specify the folder name to store intermediate models and data
11 relaxer = Ensemble_Relaxer(db=db, calculator=EMT(), jobname='AuPd-nano-test',
12                             ensemble_size=10, alpha=alpha, nn_params=nn_params)
13
14 # relaxer.run() returns a ase-db containing relaxed configurations
15 relaxed_db = relaxer.run(fmax=0.05, steps=50)
```

---

## FIGURES FOR CONFIGURATIONS USED AS THE EXAMPLE IN THE MANUSCRIPT

In the second part of the Results section, we used 13 Acrolein/AgPd configurations as the example to show the advantage of multiple configurations during geometry optimization. Here, we provides the figures for these 13 configurations.

---

```
1 from ase.db import connect
2 from ase.io import write
3 from ase.visualize import view
4 import matplotlib.image as mpimg
5 from matplotlib import pyplot as plt
6 %matplotlib inline
7
8 path = 'Acrolein-AgPd-single-multiple-configs'
9 data_path = f'./{path}/initial-configs.db'
10 db = connect(data_path)
11 i = 1
12 for entry in db.select():
13     write(f'./{path}/images/image-{i}.png', entry.toatoms())
14     i += 1
15
```

---

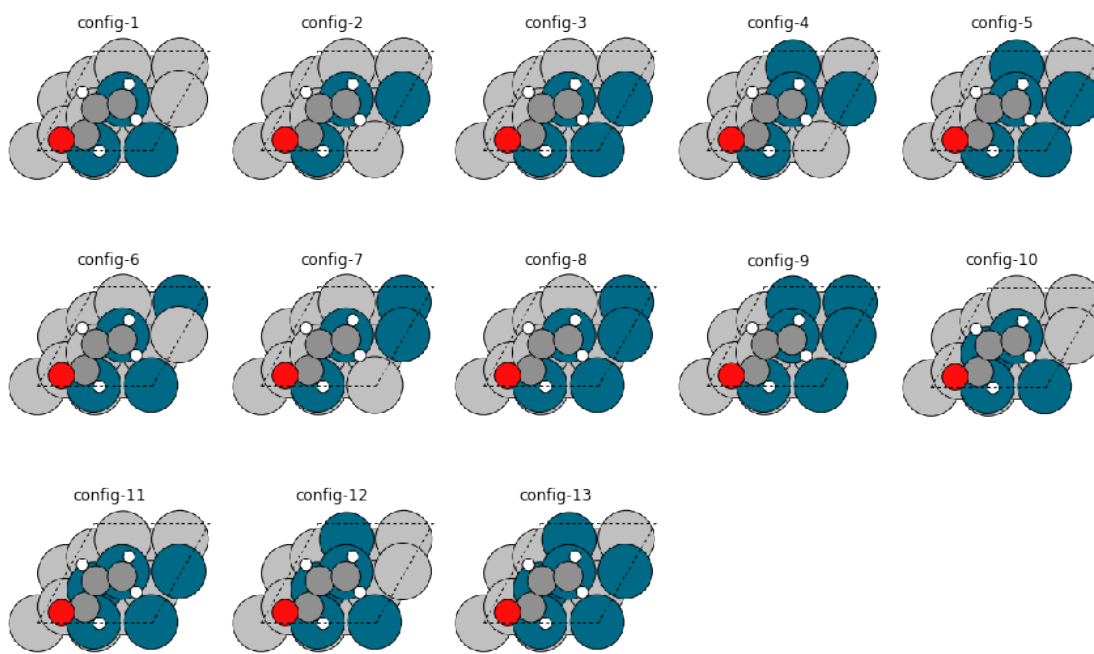
```

16 fig, axes = plt.subplots(3, 5, figsize=(12, 8))
17 for i in range(15):
18     row, col = i // 5, i % 5
19     if i < 13:
20         tmp_image = mpimg.imread(f'./{path}/images/image-{i+1}.png')
21         axes[row][col].imshow(tmp_image)
22         axes[row][col].set_title(f'config-{i+1}')
23     axes[row][col].axis('off')
24 fig.tight_layout()

```

---

<Figure size 864x576 with 15 Axes>



## CODE TO ACCESS THE DATASETS FOR REPRODUCING THE MAIN FIGURES IN THE MANUSCRIPT

This section shows how the figures in the manuscript were generated.

### Acceleration of the geometry optimization for Acrolein/AgPd

---

```

1 from ase.io.trajectory import Trajectory
2 import numpy as np
3 from matplotlib import pyplot as plt

```

```

4  %matplotlib inline
5
6
7  single_steps = []
8  multi_steps = []
9  warm_steps = []
10 gpr_steps = []
11
12 for i in range(13):
13     single_traj = Trajectory(f'./Acrolein-AgPd-single-multiple-configs/single-config-scratch-trajs/config-{i+1}.traj')
14     single_steps.append(len(single_traj))
15     multi_traj = Trajectory(f'./Acrolein-AgPd-single-multiple-configs/multi-config-scratch-trajs/config-{i+1}.traj')
16     multi_steps.append(len(multi_traj))
17     warm_traj = Trajectory(f'./Acrolein-AgPd-single-multiple-configs/multi-config-warmup-trajs/config-{i+1}.traj')
18     warm_steps.append(len(warm_traj))
19     gpr_traj = Trajectory(f'./Acrolein-AgPd-single-multiple-configs/gpr-trajs/config-{i+1}.traj')
20     gpr_steps.append(len(gpr_traj))
21
22 steps = [single_steps, multi_steps, warm_steps, gpr_steps]
23 labels = ['single', 'multiple', 'warm up', 'GPR']
24 mean_labels = ['single mean', 'multiple mean', 'warm up mean', 'GPR_mean']
25 xs = range(1, 14)
26
27 f, (ax, ax2) = plt.subplots(2, 1, sharex=True, figsize=(8, 5), )
28 ax.plot(xs, steps[3], '-o', label = labels[3], color=f'C{3}')
29 means = [round(np.mean(steps[3]), 1)] * len(xs)
30 ax.plot(xs, means, '--', color=f'C{3}', label = mean_labels[3])
31 ax.text(1, means[0] + 0.4, str(means[0]))
32 ax.set_xticks(range(1, 14))
33
34 for i in range(3):
35     ax2.plot(xs, steps[i], '-o', label = labels[i])
36     means = [round(np.mean(steps[i]), 1)] * len(xs)
37     ax2.plot(xs, means, '--', color=f'C{i}', label = mean_labels[i])
38     ax2.text(1, means[0] + 0.4, str(means[0]))
39
40 ax2.set_xticks(range(1, 14))
41 ax.spines['bottom'].set_visible(False)
42 ax2.spines['top'].set_visible(False)
43 ax.xaxis.tick_top()
44 ax.tick_params(labeltop=False)
45 ax2.xaxis.tick_bottom()
46 d = .015
47 kwargs = dict(transform=ax.transAxes, color='k', clip_on=False)
48 ax.plot((-d, +d), (-d, +d), **kwargs)
49 ax.plot((1 - d, 1 + d), (-d, +d), **kwargs)
50 kwargs.update(transform=ax2.transAxes)

```

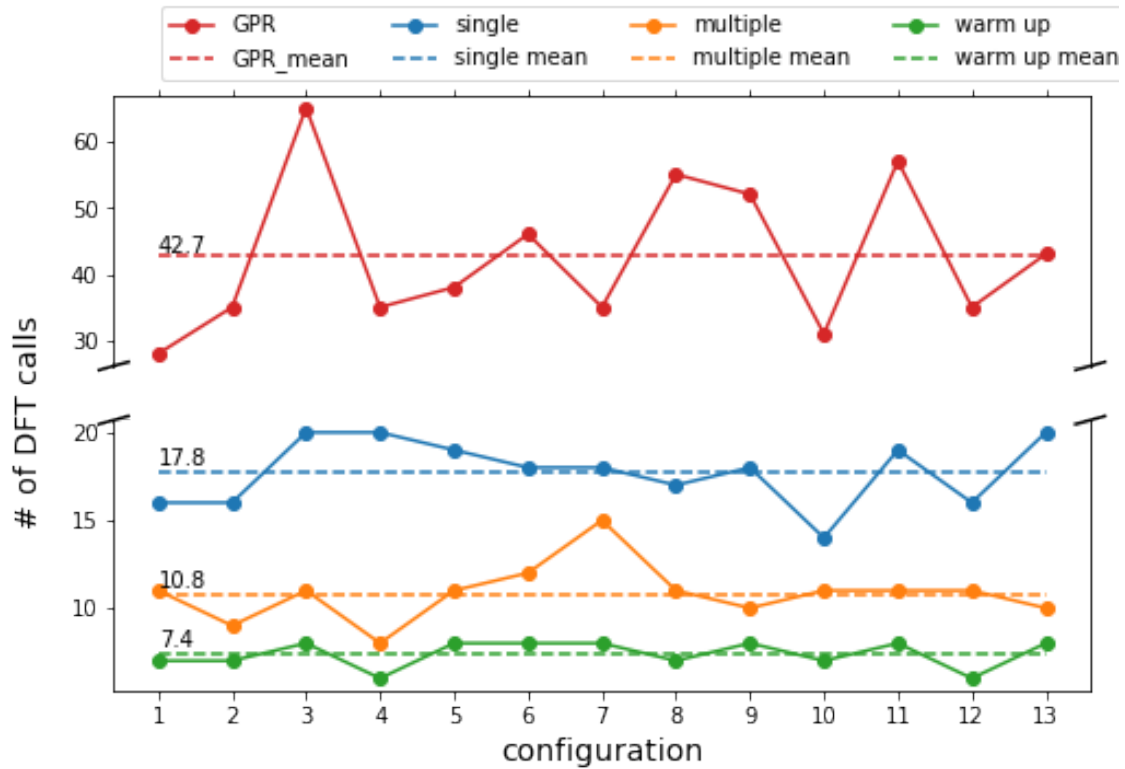
```

51 ax2.plot((-d, +d), (1 - d, 1 + d), **kwargs)
52 ax2.plot((1 - d, 1 + d), (1 - d, 1 + d), **kwargs)
53
54 ax2.set_xlabel('configuration', fontsize=14)
55 ax2.set_ylabel('# of DFT calls', fontsize=14)
56 ax2.yaxis.set_label_coords(-0.07, 1.0)
57 f.legend(bbox_to_anchor=(0.12, 1.0), ncol=4, loc='upper left')
58 f.savefig('./single-multi-warmup-gpr.png', dpi=300)

```

---

<Figure size 576x360 with 2 Axes>



## Performance for various systems

```

1 from matplotlib import pyplot as plt
2 import pickle
3 import numpy as np
4
5 path = 'more-geometry-optimization-data'
6
7 with open(f'./{path}/AuPd-slab-vasp-data-25.pkl', 'rb') as f:
8     slab_vasp_data = pickle.load(f)
9

```

```

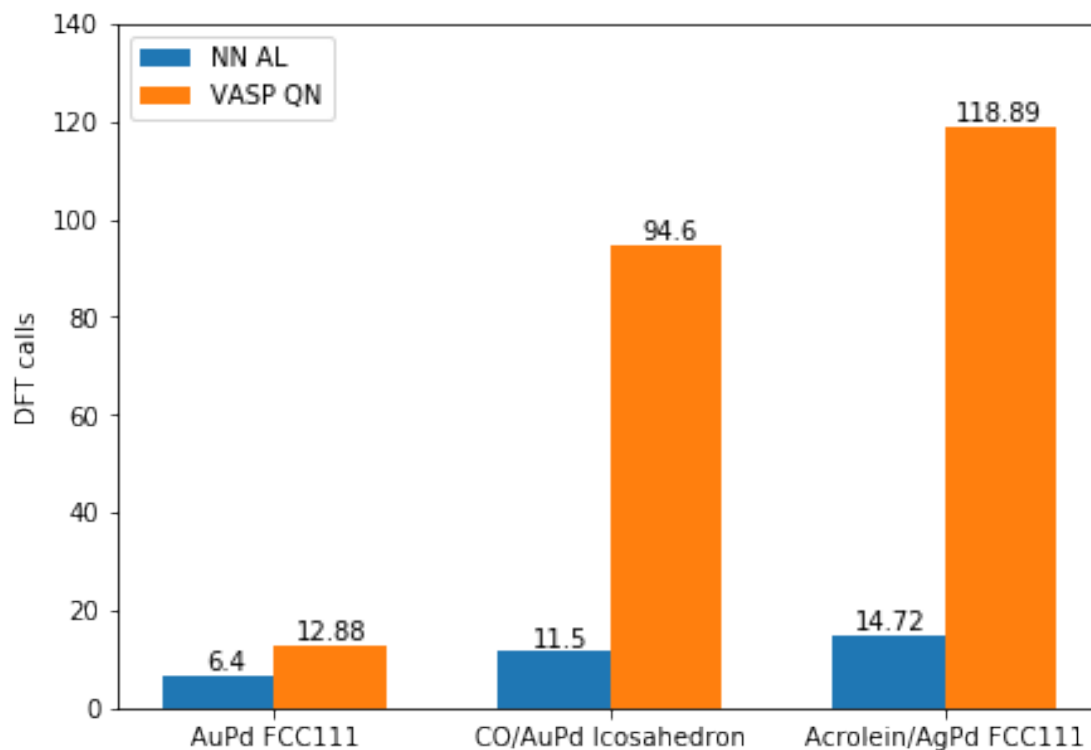
10 with open(f'./{path}/AuPd-slab-nn-data-25.pkl', 'rb') as f:
11     slab_nn_data = pickle.load(f)
12
13 with open(f'./{path}/Acrolein-AgPd-vasp-data-100.pkl', 'rb') as f:
14     ads_vasp_data = pickle.load(f)
15
16 with open(f'./{path}/Acrolein-AgPd-nn-data-scratch-100.pkl', 'rb') as f:
17     ads_nn_data = pickle.load(f)
18
19 with open(f'./{path}/CO-AuPd-Ico-vasp-data-10.pkl', 'rb') as f:
20     nano_vasp_data = pickle.load(f)
21
22 with open(f'./{path}/CO-AuPd-Ico-nn-data-10.pkl', 'rb') as f:
23     nano_nn_data = pickle.load(f)
24
25 ads_nn_step_mean = np.mean(ads_nn_data['steps'])
26 ads_vasp_step_mean = np.mean(ads_vasp_data['steps'])
27
28 slab_nn_step_mean = np.mean(slab_nn_data['steps'])
29 slab_vasp_step_mean = np.mean(slab_vasp_data['steps'])
30
31 nano_nn_step_mean = np.mean(nano_nn_data['steps'])
32 nano_vasp_step_mean = np.mean(nano_vasp_data['steps'])
33
34 step_mean = np.array([slab_nn_step_mean, slab_vasp_step_mean, nano_nn_step_mean,
35                       nano_vasp_step_mean, ads_nn_step_mean, ads_vasp_step_mean])
36
37 nn_steps = step_mean[[0, 2, 4]]
38 vasp_steps = step_mean[[1, 3, 5]]
39 nn_xs = [1, 2.5, 4]
40 vasp_xs = [1.5, 3, 4.5]
41
42 fig = plt.figure(figsize=(7, 5))
43 ax = fig.add_subplot(111)
44 ax.bar(nn_xs, nn_steps, width=0.5, label='NN AL')
45 ax.bar(vasp_xs, vasp_steps, width=0.5, label='VASP QN')
46 ax.set_ylabel('DFT calls')
47 ax.set_yticks(range(0, 160, 20))
48 ax.set_xticks([1.25, 2.75, 4.25])
49 ax.set_xticklabels(['AuPd FCC111', 'CO/AuPd Icosahedron', 'Acrolein/AgPd FCC111'])
50 ax.legend(loc='upper left')
51 ax.text(0.95, 7.5, str(nn_steps[0]))
52 ax.text(1.35, vasp_steps[0] + 1.1, str(vasp_steps[0]))
53 ax.text(2.4, nn_steps[1] + 1.1, str(nn_steps[1]))
54 ax.text(2.9, vasp_steps[1] + 1.1, str(vasp_steps[1]))
55 ax.text(3.85, nn_steps[2] + 1.1, str(nn_steps[2]))
56 ax.text(4.3, vasp_steps[2] + 1.1, str(vasp_steps[2]))

```

---

Text(4.3, 119.99, '118.89')

<Figure size 504x360 with 1 Axes>



---

```
1 from ase.db import connect
2 from matplotlib import pyplot as plt
3
4 init_db = connect('./Acetylene-hydrogenation-NEB/Acetylene-hydro-initial-configs.db')
5 vasp_db = connect('./Acetylene-hydrogenation-NEB/Acetylene-hydro-vasp-cnvg.db')
6 nn_db = connect('./Acetylene-hydrogenation-NEB/Acetylene-hydro-nn-cnvg.db')
7
8 init_nrgs = [entry.energy for entry in init_db.select()]
9 vasp_nrgs = [entry.energy for entry in vasp_db.select()]
10 nn_nrgs = [entry.energy for entry in nn_db.select()]
11 xs = range(len(vasp_nrgs))
12
13 fig = plt.figure()
14 ax = fig.add_subplot(111)
15 ax.plot(xs, nn_nrgs, '-o')
16 ax.plot(xs, vasp_nrgs, '-o')
17 ax.set_xlabel('Reaction coordinate')
18 ax.set_ylabel('energy (eV)')
```



```

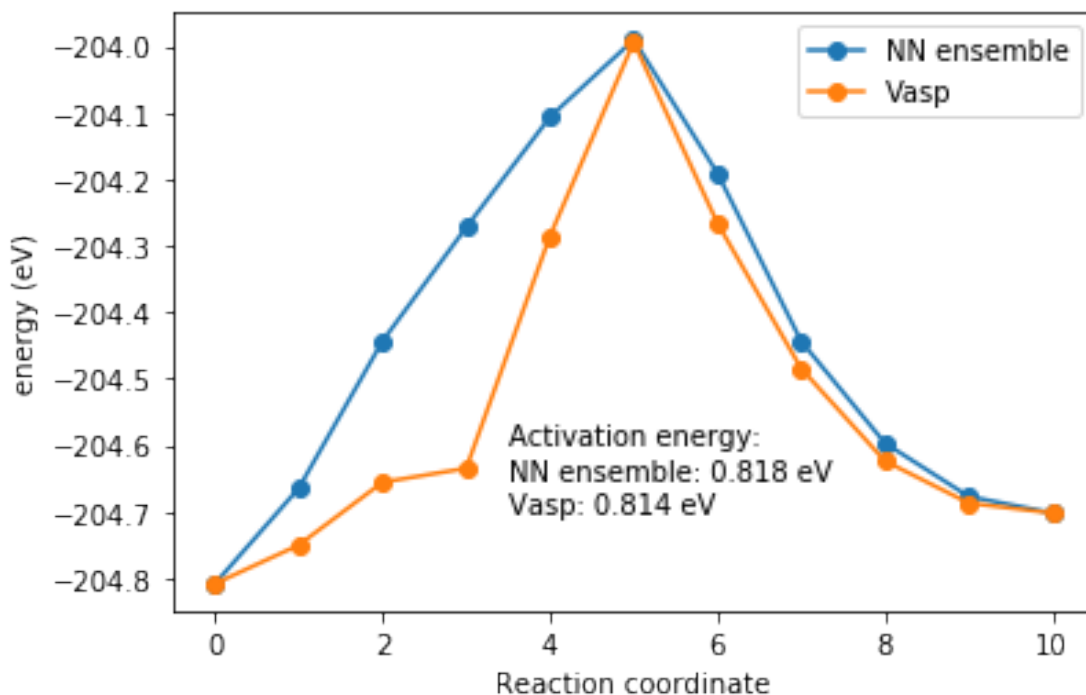
19 ax.legend(['NN ensemble', 'Vasp'])
20 vasp_act = vasp_nrgs[5] - vasp_nrgs[0]
21 nn_act = nn_nrgs[5] - nn_nrgs[0]
22 ax.text(3.5, -204.6, 'Activation energy:')
23 ax.text(3.5, -204.65, f'NN ensemble: {round(nn_act, 3)} eV')
24 ax.text(3.5, -204.7, f'Vasp: {round(vasp_act, 3)} eV')

```

---

Text(3.5, -204.7, 'Vasp: 0.814 eV')

<Figure size 432x288 with 1 Axes>



## GAUSSIAN PROCESS REGRESSION

We adapt the Gaussian Process Regression (GRP) method from previous literatures [1, 2] as one of the comparisons in this work. The GPR model uses the positions of the atoms as the feature  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ , and the model is trained on the corresponding energies (e) and the first order derivative which is the negative forces in this application. Thus,  $\mathbf{y} = [\mathbf{e}, -\mathbf{f}_1, \dots, -\mathbf{f}_N]$

Therefore, the prediction function could be sampled from the Gaussian Process defined by a prior mean and a kernel function:

$$f(\mathbf{x}) \sim GP(\boldsymbol{\mu}, k(\mathbf{x}, \mathbf{x}')) \quad (1)$$

where  $\boldsymbol{\mu}$  is the prior for the energies and forces and it is set as zero in our work. Given a training set  $D$ , the predicted mean and variance are

$$E[f(\mathbf{x}|D)] = \mathbf{k}(\mathbf{x}) [\mathbf{K}(\mathbf{x} + \sigma_n^2 \mathbf{I})]^{-1} \mathbf{y} \quad (2)$$

and

$$V[f(\mathbf{x}|D)] = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T [\mathbf{K}(\mathbf{x} + \sigma_n^2 \mathbf{I})]^{-1} \mathbf{k}(\mathbf{x}) \quad (3)$$

where  $\sigma_n$  is the noise of the data.

The kernel function could be partitioned into:

$$\mathbf{K}(\mathbf{x}) = \begin{pmatrix} \mathbf{K}_{ee}(\mathbf{x}, \mathbf{x}) & \mathbf{K}_{ef}(\mathbf{x}, \mathbf{x}) \\ \mathbf{K}_{ef}(\mathbf{x}, \mathbf{x}) & \mathbf{K}_{ff}(\mathbf{x}, \mathbf{x}) \end{pmatrix} \quad (4)$$

Squared exponential kernel is used in our implementation. Thus, the formula for these kernel function are:

$$k_{ee}(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left( -\frac{1}{2} \sum_{d=1}^D \frac{(x_d - x'_d)^2}{l_d^2} \right) \quad (5)$$

$$k_{fe}(\mathbf{x}, \mathbf{x}') = -\frac{\sigma_f^2 (x_d - x'_d)}{l_d^2} \exp \left( -\frac{1}{2} \sum_{j=1}^D \frac{(x_j - x'_j)^2}{l_j^2} \right) \quad (6)$$

$$k_{ff}(\mathbf{x}, \mathbf{x}') = \frac{-\sigma_f^2}{l_{d1}^2} \left( \delta_{d1d2} - \frac{(x_{d1} - x'_{d1})(x_{d2} - x'_{d2})}{l_{d1}^2} \right) \exp \left( -\frac{1}{2} \sum_{j=1}^D \frac{(x_j - x'_j)^2}{l_j^2} \right) \quad (7)$$

where  $\sigma_f$  is fixed as 1.0, and the bandwidth  $l$  is optimized isotropically.

---

\* [jkitchin@andrew.cmu.edu](mailto:jkitchin@andrew.cmu.edu)

- [1] J. A. G. Torres, P. C. Jennings, M. H. Hansen, J. R. Boes, and T. Bligaard, [Physical Review Letters](#) **122**, 156001 (2019).
- [2] O.-P. Koistinen, F. B. Dagbjartsdóttir, V. Ásgeirsson, A. Vehtari, and H. Jónsson, [The Journal of Chemical Physics](#) **147**, 152720 (2017).