

1
2
3
4
5
6
7
8

Supplementary Information:
**By how much can closed loop frameworks
accelerate computational materials discovery?**

Lance Kavalsky,^{1,*} Vinay I. Hegde,^{2,*} Eric Muckley,² Matthew S.
Johnson,³ Bryce Meredig,^{2,†} and Venkatasubramanian Viswanathan^{1,‡}

¹*Carnegie Mellon University, Pittsburgh, PA 15213*

²*Citrine Informatics, Redwood City, CA 94063*

³*Massachusetts Institute of Technology, Cambridge, MA 02139*

9 I. EFFECTIVENESS OF ACQUISITION FUNCTIONS

10 Here, we expand the set of acquisition functions and compare their effectiveness for the
 11 three sequential learning (SL) tasks considered in the main text: (a) finding candidates in
 12 the target window, (b) surfacing candidates of high quality, i.e., close to the target window,
 13 and (c) building minimal datasets for training ML models that are close in accuracy to those
 14 trained on the full dataset.

15 We define a new baseline for acquisition functions, a “space-filling” approach, that works as

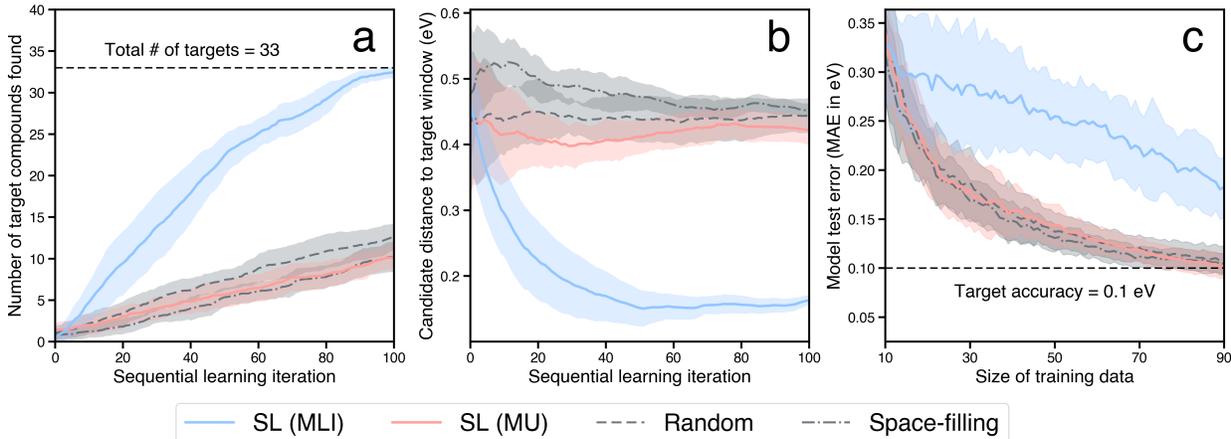


FIG. S1. A comparison of random search vs sequential learning (SL)-driven approach to find new bimetallic catalysts with a target property. (a) Overall, the SL-driven approach identifies all the 33 target candidates in the dataset within 100 iterations, $\sim 3\times$ faster than random search and a space-filling-based acquisition function. (b) Candidates surfaced via SL lie much closer to the target window on average, when compared to those chosen via random search or space-filling. (c) An SL-driven approach can help identify a much smaller number of examples that can be used to train ML surrogates to a desired accuracy, at a fraction of the overall dataset size. Here, the overall dataset has ~ 300 candidates, and an ML model trained on only $\sim 25\%$ of the candidates chosen via a SL-driven maximum uncertainty (MU)-based approach achieves the target accuracy. Notably, the other purely-exploratory acquisition functions considered here, random and space-filling, perform as well as MU for building datasets. In each case, the shaded region in the plots represent variation in the reported quantities estimated over 20 independent trials.

* These authors contributed equally to this work

† bryce@citrine.io

‡ venkvis@cmu.edu

16 follows: (1) Start with a randomly-chosen initial training set (e.g., ~ 10 , used previously),
17 (2) Build ML models, (3) From the design space of candidates, choose one that is *farthest*
18 from all the training examples (measured using Euclidian distance of each candidate, in the
19 145-dimensional Magpie feature vector-space, to the closest example in the current training
20 set) to augment the training set in each SL iteration, (4) Re-train ML models, (5) Perform
21 steps (1)–(5) for a predetermined number of iterations or until the termination condition
22 is met. This space-filling acquisition strategy performs similarly to the random (uniform)
23 acquisition strategy in finding candidates, candidate quality, as well as for building datasets
24 for ML surrogates (see FIG. S1).

25 Expectedly, MU, random, and space-filling perform worse than MLI for finding target
26 candidates and surfacing quality candidates (FIG. S1a,b), as both are SL tasks that benefit
27 from balancing exploration with exploitation. On the other hand, the space-filling and ran-
28 dom acquisition strategies perform as well as the MU-based one for building training datasets
29 (and all three exploratory acquisition strategies perform better than MLI), demonstrating
30 that the three approaches are nearly-equally effective for pure exploration (FIG. S1c). This
31 is consistent with previous reports comparing model accuracy as a function of SL iteration
32 using similar acquisition functions [1].

33 In terms of picking an acquisition strategy specifically for constructing training datasets,
34 there is additional nuance to consider related to cost. While the space-filling and random
35 acquisition functions do not require ML model building in each SL iteration (and are thus
36 faster and more economical compared to MU), the effectiveness of the three approaches is
37 seen to vary based on the data distribution in the particular design space of interest [1].
38 This can be the subject of future investigation.

39 II. HUMAN LAGTIME MODEL

40 When defining our human lag model, we invoke the following assumptions:

41 1. 3 windows of time:

42 (a) Researcher at work (9am-5pm): checks on job every couple of hours; average lag
43 of 1 hour.

44 (b) Researcher (partially) away from work (5pm-11pm): checks on job at the end of
45 the window; average lag of 3 hours.

46 (c) Researcher (completely) away from work (11pm-9am): checks on job at the end
47 of window; average lag of 5 hours

48 2. Uniform distribution of when jobs will finish through a week:

49 (a) Probability of job finishing during the week = $5/7$: researcher checks on job
50 according to 1, above.

51 (b) Probability of job finishing during the weekend = $2/7$: researcher checks on job
52 once during the weekend; average lag of 24 hours.

53 III. SURROGATE ACCURACY

54 The size of the unexplored design space shrinks as the simulated sequential learning (SL)
55 progresses, i.e., the number of candidates in the full dataset that the model has not “seen” yet
56 keeps continuously decreasing. So surrogate model accuracy estimates derived using model
57 predictions over the entire unexplored design space (the test set) in each SL iteration can be
58 affected by the continuously diminishing test set size. In order to mitigate this effect of test
59 set size on the surrogate model accuracy estimates, we employ a bootstrapping approach to
60 keep the test set size fixed in each SL iteration. At each SL iteration, 20 “bootstrap test
61 samples” are generated from the full unexplored design space. Each of the 20 bootstrap test
62 samples are generated by randomly sampling, with replacement, 100 candidates from the
63 unexplored design space. For each of the 20 bootstrap test samples (with 100 candidates
64 each) we calculate the mean absolute error (MAE). Finally, we run 20 independent trials
65 of the entire simulated SL pipeline. The accuracy of the surrogate model at a given SL
66 iteration is then defined by the mean and standard deviation of the MAEs of the bootstrap
67 test samples (generated at that SL iteration) aggregated over the 20 independent trials. We
68 use this final mean MAE of the surrogate model as the accuracy metric of interest with a
69 target value of 0.1 eV.

70 IV. SUMMARY OF ACCELERATION FACTORS

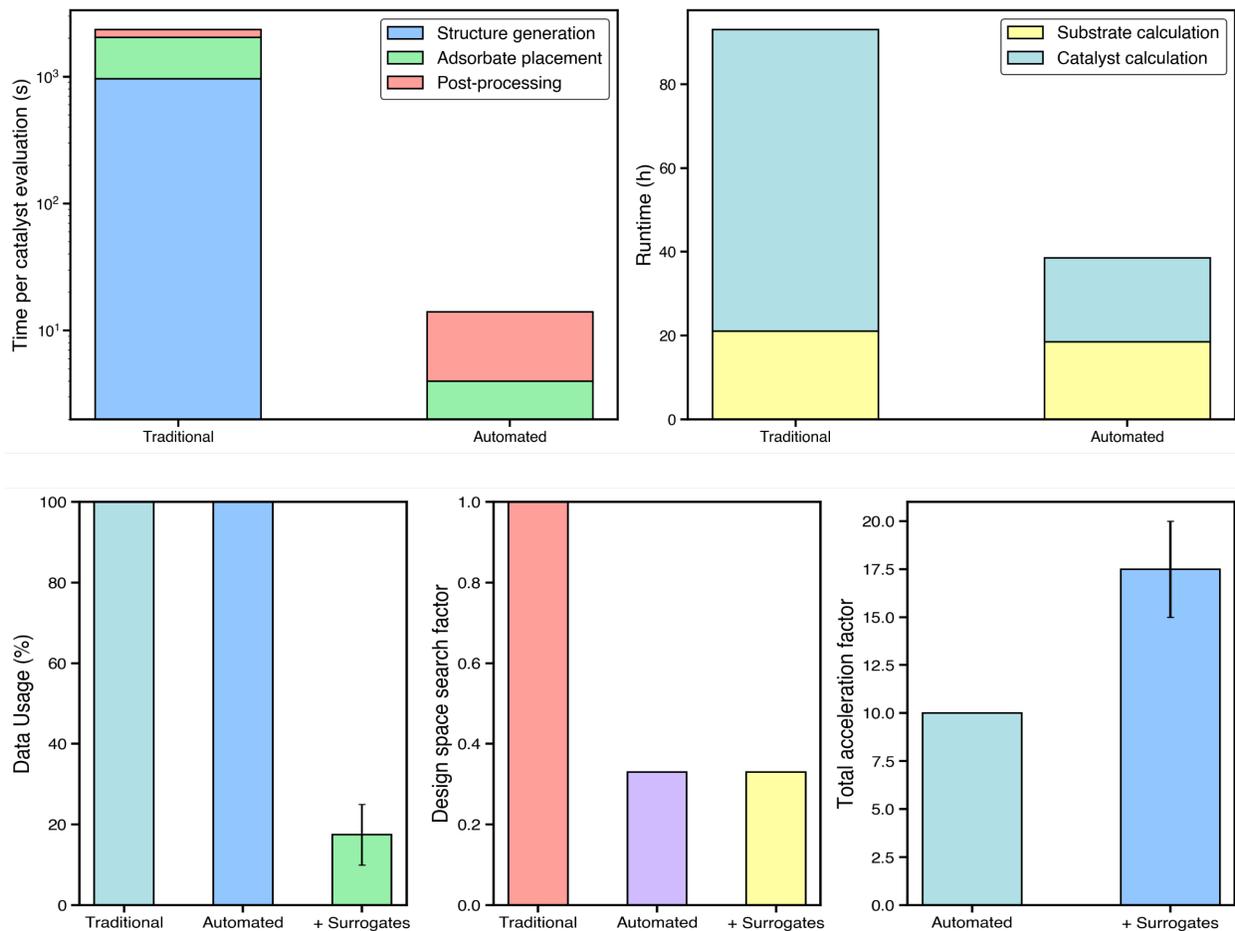


FIG. S2. Acceleration factors visualized as bar plots using the data from Table III of the main manuscript. The component-wise analysis highlighted here can be used to make informed decisions when designing closed-loop frameworks of varying topologies.

71 V. DATA AND SCRIPTS FOR REPRODUCIBILITY

72 All data and Python scripts required to perform the analysis presented in this work are
73 made available via the GitHub repository at [https://github.com/aced-differentiate/
74 closed-loop-acceleration-benchmarks](https://github.com/aced-differentiate/closed-loop-acceleration-benchmarks). The repository is organized as follows:

75 1. [data/](#)

76 • [benchmark_calculations_record.xlsx](#): Excel spreadsheet containing a record
77 of DFT calculations, associated raw timestamps, and a tabulation of the acceler-
78 ation estimates.

79 • [bimetallic_catalysts_dataset/](#)

80 – [ma_2015_bimetallics_raw.json.gz](#): Dataset of bimetallic alloys for CO2
81 reduction, in the [Physical Information File \(PIF\)](#) format, obtained from
82 [Dataset 153450](#) on Citrination.

83 Original data source: “Machine-Learning-Augmented Chemisorption Model
84 for CO2 Electroreduction Catalyst Screening”, Ma et al., *J. Phys. Chem.
85 Lett.* **6** 3528-3533 (2015). DOI: [10.1021/acs.jpcclett.5b01660](https://doi.org/10.1021/acs.jpcclett.5b01660)

86 – [transform.py](#): Python script for converting from the PIF format into tab-
87 ular data.

88 – [bimetallics_data.csv](#): Bimetallics catalysts dataset mentioned above in a
89 tabular format.

90 • [runtime_geometries/](#)

91 “Chemically-informed” and naive structures and settings in the form of `ase.traj`
92 files, corresponding to the discussion surrounding FIG. 3 in the main text. The
93 files can be read using ASE package (using `ase.io.read`).

94 2. [scripts/](#)

95 • [human_lagtime.py](#): Script for estimating human lagtime in job management,
96 calculated using a Monte Carlo sampling method.

97 • [sequential_learning.py](#): Script for running multiple independent trials of se-
98 quential learning (SL) and recording a history of training examples, model pre-
99 dictions and prediction uncertainties.

100 If run as-is, the script performs 20 independent trials of 100 SL iterations to op-
101 timize the `binding_energy_of_adsorbed` property in the bimetallic catalysts
102 dataset mentioned above, using four acquisition functions (results from each
103 recorded separately): random, maximum likelihood of improvement (MLI), max-
104 imum uncertainty (MU), and space-filling.

- 105 • `plot_acceleration_from_sequential_learning.py`: Script to aggregate re-
106 sults from the `sequential_learning.py` script, calculate and plot statistics re-
107 lated to acceleration from SL over a baseline.

108 If run as-is, the script reproduces the 3-paneled FIG. 5 in the main text.

- 109 • `plot_acceleration_from_sequential_learning__ALL_ACQ.py`: Similar to the
110 previous script; plots and compares statistics from all acquisition functions con-
111 sidered (MLI, MU, random, space-filling) for all SL tasks.

112 If run as-is, the script reproduces the 3-paneled FIG. S1 in the Supplementary
113 Information.

- 114 • `plot_levels_of_automation.py`: Script for plotting the cumulative time of ex-
115 ecuting the DFT pipeline at varying levels of automation.

116 If run as-is, the script reproduces the bottom panel from FIG. 2 in the main text.

117 A. Running the scripts

118 The required packages for executing the scripts are specified in `requirements.txt`, and
119 can be installed in a new environment (e.g. using [conda](#)) as follows:

```
120 $ conda create -n accel_benchmarking python=3.10  
121 $ conda activate accel_benchmarking  
122 $ pip install -r requirements.txt
```

123 The scripts are all in python, and can be run from the command line. For example:

```
124 $ cd scripts  
125 $ python sequential_learning.py
```

126 [1] C. K. Borg, E. S. Muckley, C. Nyby, J. E. Saal, L. Ward, A. Mehta, and B. Meredig, Quantifying
127 the performance of machine learning models in materials discovery, *Digital Discovery* **2**, 327
128 (2023).