Chemically-informed data-driven optimization (ChIDDO): Leveraging physical models and Bayesian learning to accelerate chemical research

(Electronic Supplementary Information)

Daniel Frey, Juhee Shin, Christopher Musco, and Miguel A. Modestino

Acquisition Function Calculations

Expected Improvement (EI):

$$EI(x) = \left(\mu(x) - f(x^+) - \xi\right)\psi\left(\frac{\mu(x) - f(x^+) - \xi}{\sigma(x)}\right) + \sigma(x)\phi\left(\frac{\mu(x) - f(x^+) - \xi}{\sigma(x)}\right)$$

where $\mu(x)$ is the mean of the regressor at x, $\sigma(x)$ is the variance of the regressor at x, f is the function to be maximized, x^+ is the location of the estimated maximum, ξ is the exploration/exploitation parameter, and $\psi(z)$ is the cumulative distribution function of a standard Guassian distribution, and $\phi(z)$ is the density function of a standard Gaussian distribution¹.

Probability of Improvement (PI):

$$PI(x) = \psi\left(\frac{\left(\mu(x) - f(x^+) - \xi\right)}{\sigma(x)}\right)$$

where $\mu(x)$ is the mean of the regressor at x, $\sigma(x)$ is the variance of the regressor at x, f is the function to be maximized, x^+ is the location of the estimated maximum, ξ is the exploration/exploitation parameter, and $\psi(z)$ is the cumulative distribution function of a standard Guassian distribution¹.

Upper Confidence Bound (UCB):

$$UCB(x) = \mu(x) + \beta\sigma(x)$$

where $\mu(x)$ is the mean of the regressor at x, $\sigma(x)$ is the variance of the regressor at x, and β is the exploration/exploitation parameter¹.



Objective Function Descriptions





BO, ChIDDO, and MISO algorithm descriptions

Algorithm 1: Bayesian Optimization with no physics model data **Procedure**: BO

Input: A set of N_{init} experimental points, X^{exp} , where $X^{exp} \in_R DS$ evaluated to give the objective function value, y^{exp} . P^{exp} = { X^{exp} , y^{exp} }

- 1. P^{exp} is passed to the GPR to make predictions, μ^{GPR} , and uncertainties, σ^{GPR} , at any point in DS.
- 2. P^{exp} , μ^{GPR} , and σ^{GPR} are passed to the acquisition function to select n_b new experiments, X^{next} , to evaluate.
- 3. Objective function is evaluated at X^{next} to give y^{next} (P^{next}) and then $P^{exp} = P^{exp} \cup P^{next}$.
- 4. Steps 1-3 is repeated until 50, N_{total}, experiments are evaluated.
- 5. Return: Pexp

Algorithm 2: Chemically-informed data driven optimization **Procedure**: ChIDDO

Input: A set of N_{init} experimental points, X^{exp} , where $X^{exp} \in_R DS$ evaluated to give the objective function value, y^{exp} . $P^{exp} = \{X^{exp}, y^{exp}\}$

Input: A model, M, of the system under study to approximate the objective function

- 1. $X^{phys} \in_R DS$ are evaluated by M to give y_{phys} ($P^{phys} = \{X^{phys}, y^{phys}\}$) and then $P^{tot} = P^{exp} \cup P^{phys}$
- 2. P^{tot} is passed to the GPR to make predictions, μ^{GPR} , and uncertainties, σ^{GPR} , at any point in DS.
- 3. P^{exp} , μ^{GPR} , and σ^{GPR} are passed to the acquisition function to select n_b new experiments, X^{next} , to evaluate.
- 4. Objective function is evaluated at X^{next} to give y^{next} (P^{next}) and then $P^{exp} = P^{exp} \cup P^{next}$.
- 5. P^{exp} is used to refine the parameters of M by using non-linear regression.
- 6. Steps 1-5 is repeated until 50, N_{total}, experiments are evaluated.
- 7. **Return**: P^{exp}

Algorithm 3: Multi-information Source Optimization as reported in Poloczek et al.²

Procedure: misoKG

Input: A set of N_{init} experimental points, X^{exp} , where $X^{exp} \in_R DS$ evaluated to give the objective function value, y^{exp} . $P^{exp} = \{X^{exp}, y^{exp}\}$

Input: A set of n_l information sources, l, that give different biases on the objective function.

- 1. $X^{exp} \in_{\mathbb{R}} DS$ are evaluated by each *l* to give $y^{exp}(l, X^{exp})$.
- 2. X^{exp} and y^{exp} are passed to the GPR to make the posterior prediction.
- 3. Until the budget for samples is consumed, determine the *l* and X^{next} that maximize the misoKG factor (Equation 1 in Reference 1)².
- 4. Objective function is evaluated at X^{next} and l and the posterior is updated with the new information.
- 5. Steps 2-4 are repeated until sample budget is exhausted.

Return: Pexp



Comparison between 2D simplified models and experimental objective function

the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.

10

20

Figure S1. dx versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search. (A) 2D Sphere, (B) 3D Sphere, (C) 4D Sphere. For each curve, 20 separate searches, S, were performed, and

Number of Experiments, N

30

[×] 1.00

0.75

0.50

0.25

0.00

20

Number of Experiments, N

30

40

50

10

40

50

BO and ChIDDO comparison to Edisonian search using MRB

30

40

50

20

Number of Experiments, N

10

€.0 `

0.6

0.4

0.2

0.0

0

1.0

0.8

0.4

0.2

0.0

٥.6 [×]



Figure S2. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 2D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S3. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 3D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S4. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 4D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S5. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 6D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S6. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 6D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S7. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 2D Camel function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S8. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 2D Mccormick function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S9. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 2D Branin function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S10. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 3D Hartmann function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S11. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 4D Hartmann function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S12. d_y and d_x versus number of experiments, N, comparing BO and ChIDDO with the Edisonian random and grid search on the 6D Hartmann function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For each of the BO/ChIDDO experiments, the MRB acquisition function was used.



Figure S13. d_x versus number of experiments, N, comparing different acquisition functions using BO on: (A) 3D Hartmann, (B) 4D Hartmann, (C) 6D Hartmann. For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S14. d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on: (A) 3D Hartmann, (B) 4D Hartmann, (C) 6D Hartmann. For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.

Acquisition function comparison for all objective functions



Figure S15. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 2D Branin function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S16. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 2D Branin function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S17. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 2D Camel function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S18. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 2D Camel function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S19. d_y and d_x versus number of experiments, *N*, comparing different acquisition functions using BO on the 2D Mccormic function. (A) d_y , (B) d_x . For each curve, 20 separate searches, *S*, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S20. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 2D Mccormick function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S21. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 2D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S22. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 2D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S23. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 3D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S24. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 3D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S25. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 4D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S26. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 4D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S27. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 6D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S28. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 6D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S29. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 2D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S30. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 2D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S31. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 3D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S32. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 3D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S33. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 4D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S34. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 4D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S35. d_y and d_x versus number of experiments, N, comparing different acquisition functions using BO on the 6D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.



Figure S36. d_y and d_x versus number of experiments, N, comparing different acquisition functions using ChIDDO on the 6D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation.

Noise level comparison for other objective functions using MRB

Figure S37. d_y and d_x versus number of experiments, N, comparing different noise levels using BO on the 2D Branin function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S38. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 2D Branin function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S39. d_y and d_x versus number of experiments, *N*, comparing different noise levels using BO on the 2D Mccormick function. (A) d_y , (B) d_x . For each curve, 20 separate searches, *S*, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S40. d_y and d_x versus number of experiments, *N*, comparing different noise levels using ChIDDO on the 2D Mccormick function. (A) d_y , (B) d_x . For each curve, 20 separate searches, *S*, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S41. d_y and d_x versus number of experiments, N, comparing different noise levels using BO on the 2D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S42. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 2D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S43. d_y and d_x versus number of experiments, N, comparing different noise levels using BO on the 3D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S44. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 3D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S45. d_y and d_x versus number of experiments, *N*, comparing different noise levels using BO on the 4D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, *S*, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S46. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 4D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S47. d_y and d_x versus number of experiments, *N*, comparing different noise levels using BO on the 6D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, *S*, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S48. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 6D Rosenbrock function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S49. d_y and d_x versus number of experiments, N, comparing different noise levels using BO on the 2D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S50. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 2D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S51. d_y and d_x versus number of experiments, N, comparing different noise levels using BO on the 3D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S52. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 3D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S53. d_y and d_x versus number of experiments, *N*, comparing different noise levels using BO on the 4D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, *S*, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S54. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 4D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S55. d_y and d_x versus number of experiments, N, comparing different noise levels using BO on the 6D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Figure S56. d_y and d_x versus number of experiments, N, comparing different noise levels using ChIDDO on the 6D Sphere function. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. The MRB acquisition function was used.

Physics model accuracy comparison for other objective function combinations using MRB acquisition function

Figure S57. d_y and d_x versus number of experiments, N, for different objective function mixing ratios. 2D Camel mixed with 2D Rosenbrock using Camel as physics model. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For all of these graphs, ChIDDO was used as the AL algorithm and MRB was used as the acquisition function.

Figure S58. d_y and d_x versus number of experiments, N, for different objective function mixing ratios. 2D Camel mixed with 2D Rosenbrock using Rosenbrock as physics model. (A) d_y , (B) d_x . For each curve, 20 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For all of these graphs, ChIDDO was used as the AL algorithm and MRB was used as the acquisition function.

Figure S59. d_y versus number of experiments, N, for different electrochemical physics model information, using a logarithmic exploration/exploitation rate. "Full" indicates the model is predicting the same information as the objective function. "No E", "No F", and "No EF" indicate the removal of Equations 8 and/or 9 from the physics model information, resulting in a less informative model. (A) 2D electrochemical model. (B) 3D electrochemical model. (C) 4D electrochemical model. For each curve, 25 separate searches, S, were performed, and the average of the results are the lines shown. The shadow around each of the lines represents the standard deviation. For all of these graphs, ChIDDO was used as the AL algorithm and MRB was used as the acquisition function.

Supporting Information References

1. Acquisition funcitons. <u>https://modal-</u>

python.readthedocs.io/en/latest/content/query_strategies/Acquisition-functions.html.

2. Poloczek, M.; Wang, J.; Frazier, P., Multi-information source optimization. *Advances in Neural Information Processing Systems* **2017**, 4288-4298.