

nimCSO: A Nim package for Compositional Space Optimization

Adam M. Krajewski¹, Arindam Debnath¹, Wesley F. Reinhart^{1,2},
Allison M. Beese¹, and Zi-Kui Liu¹

¹ Department of Materials Science and Engineering, The Pennsylvania State University, USA ² Institute for Computational and Data Sciences, The Pennsylvania State University, USA ¶ Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright
and release the work under a
Creative Commons Attribution 4.0
International License ([CC BY 4.0](#))

Summary

nimCSO is a high-performance tool implementing several methods for selecting components (data dimensions) in compositional datasets, which optimize the data availability and density for applications such as machine learning. Making said choice is a combinatorically hard problem for complex compositions existing in highly dimensional spaces due to the interdependency of components being present. Such spaces are encountered, for instance, in materials science, where datasets on Compositionally Complex Materials (CCMs) often span 20-45 chemical elements, 5-10 processing types, and several temperature regimes, for up to 60 total data dimensions.

At its core, nimCSO leverages the metaprogramming ability of the Nim language ([Rumpf, 2023](#)) to optimize itself at the compile time, both in terms of speed and memory handling, to the specific problem statement and dataset at hand based on a human-readable configuration file. As demonstrated in [Methods and Performance](#) section, nimCSO reaches the physical limits of the hardware (L1 cache latency) and can outperform an efficient native Python implementation over 400 times in terms of speed and 50 times in terms of memory usage (*not* counting interpreter), while also outperforming NumPy implementation 35 and 17 times, respectively, when checking a candidate solution.

nimCSO is designed to be both (1) a user-ready tool, implementing two efficient brute force approaches (for handling up to 25 dimensions), a custom search algorithm (for up to 40 dimensions), and a genetic algorithm (for any dimensionality), and (2) a scaffold for building even more elaborate methods in the future, including heuristics going beyond data availability. All configuration is done with a simple human-readable YAML config file and plain text data files, making it easy to modify the search method and its parameters with no knowledge of programming and only basic command line skills.

Statement of Need

nimCSO is an interdisciplinary tool applicable to any field where data is composed of a large number of independent components and their interaction is of interest in a modeling effort, ranging from social sciences like economics, through medicine where drug interactions can have a large impact on the treatment, to chemistry and materials science, where the composition and processing history are critical to resulting properties. The latter has been the root motivation for the development of nimCSO within the [ULTERA Project](#) ([ultera.org](#)) carried under the [US DOE ARPA-E ULTIMATE](#) program which aims to develop a new generation of ultra-high temperature materials for aerospace applications, through generative machine learning models ([Debnath et al., 2021](#)) driving thermodynamic modelling and experimentation ([Li et al., 2024](#)).

One of the most promising materials for such applications are the aforementioned CCMs, and their metal-focused subset of Refractory High Entropy Alloys (RHEAs) (Senkov et al., 2018), which are rapidly growing since first proposed by (Cantor et al., 2004) and (Yeh et al., 2004). Contrary to most of the traditional alloys, they contain a large number of chemical elements (typically 4-9) in similar proportions, in hope to thermodynamically stabilize the material by increasing its configurational entropy ($\Delta S_{conf} = \sum_i^N x_i \ln x_i$ for ideal mixing of N elements with fractions x_i), what encourages sampling from a large palette of chemical elements. At the time of writing, ULTERA Database is the largest collection of HEA data, containing over 6,300 points manually extracted from almost 550 publications. It covers 37 chemical elements resulting in extremely compositional spaces (Krajewski et al., 2024); thus, it becomes critical to answer questions like “Which combination of 15 elements will result in the largest dataset?” which has $\binom{37}{15}$ or 10 billion possible solutions.

Another significant example of intended use is to perform similar optimizations over large (many millions) datasets of quantum mechanics calculations spanning 93 chemical elements and accessible through OPTIMADE API (Evans et al., 2024).

Methods and Performance

Overview

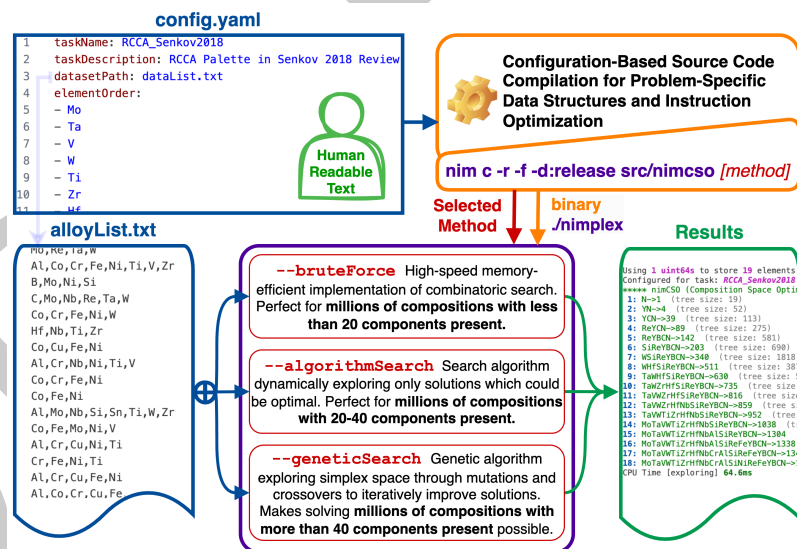


Figure 1: Schematic of core nimCSO data flow with a description of key methods. Metaprogramming is used to compile the software optimized to the human-readable data and configuration files at hand.

Under the hood, nimCSO is built around storing the data and solutions in one of two ways. The first is as bits encoded in an integer (uint64), which allows for highest speed and lowest memory consumption possible, but is limited to 64 dimensions and does not allow for easy extension to other use cases, thus as of publication it is used only in the special bruteForceInt routine. The second one, used in bruteForce, algorithmSearch, and geneticSearch is through a custom ElSolution type containing heuristic value (easily extensible) and BitArray payload, which is defined at compile time based on the configuration file to minimize necessary overheads. Both encodings significantly outperform both typical native Python and NumPy implementations, as shown in the Table ??.

Table 1: Benchmarks of (1) average time to evaluate how many datapoints would be lost if 5 selected components were removed from a dataset with 2,150 data points spanning 37 components, averaged over 10,000 runs, and (2) the size of the data structure representing the dataset. Values were obtained by running scripts in benchmarks directory on Apple M2 Max CPU.

Tool	Object	Time per Dataset	Time per Entry (Relative)	Database Size (Relative)
Python ^{3.11}	set	327.4 μ s	152.3 ns (<i>x1</i>)	871.5 kB (<i>x1</i>)
NumPy ^{1.26}	array	40.1 μ s	18.6 ns (<i>x8.3</i>)	79.7 kB (<i>x10.9</i>)
nimCSO ^{0.6}	BitArray	9.2 μ s	4.4 ns (<i>x34.6</i>)	50.4 kB (<i>x17.3</i>)
nimCSO ^{0.6}	uint64	0.79 μ s	0.37 ns (<i>x413</i>)	16.8 kB (<i>x52</i>)

Brute-Force Search

The brute force search is a naïve method of evaluating all possibilities; however, its near-zero overhead can make it the most efficient for small problems. In this implementation, all entries in the *power set* of considered elements are represented as a range of integers from 0 to $2^N - 1$, where N is the number of components, and used to initialize uint64/BitArrays on the fly. To minimize the memory footprint of solutions, the algorithm only keeps track of the best solution for a given number of elements present in the solution. Current implementations are limited to 64 elements, as it is not feasible for more than approximately 30 elements; however, the one based on BitArray could be easily extended if needed.

Algorithm-Based Search

For higher dimensional problems (20-50 components), the brute force search becomes suboptimal, prompting the need for a more efficient method. The algorithm implemented in the `algorithmSearch` routine iteratively expands and evaluates candidates from a priority queue (implemented through an efficient binary heap (Williams, 1964)) while leveraging the fact that *the number of data points lost when removing elements A and B from the dataset has to be at least as large as when removing either A or B alone* to delay exploration of candidates until they can contribute to the solution. Furthermore, to (1) avoid revisiting the same candidate without keeping track of visited states and (2) further inhibit the exploration of unlikely candidates, the algorithm *assumes* that while searching for a given order of solution, elements present in already expanded solutions will not improve those not yet expanded. This effectively prunes candidate branches requiring two or more levels of backtracking. In the authors' tests, this method has generated the same results as `bruteForce`, except for occasional differences in the last explored solution.

Genetic Search

Beyond 50 components, the *algorithm-based* method will likely run out of memory on most personal systems. The `geneticSearch` routine solves that through an evolution strategy to iteratively improve solutions based on custom mutate and crossover procedures. Both are of uniform type (Goldberg, 1989) with additional constraint of Hamming weight (Knuth, 2009) preservation in order to preserve number of considered elements in parents and offspring. In mutate this is achieved by using purely random bit swapping, rather than more common flipping, as demonstrated in the Figure 2.

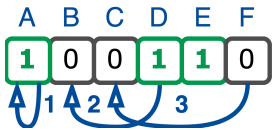


Figure 2: The schematic of mutate procedure where bits are swapping randomly, so that (1) bit can swap itself, (2) bits can swap causing a flip, or (3) bits can swap with no effect.

108 Meanwhile, in crossover, this constraint is satisfied by passing overlapping bits directly, while
109 non-overlapping bits are shuffled and distributed at positions present in one of the parents, as
110 shown in the Figure 3.

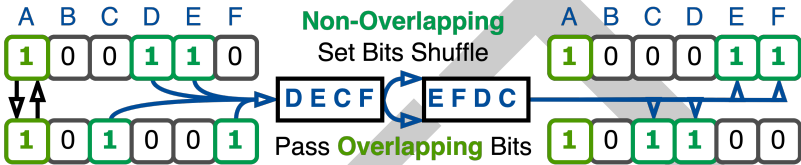


Figure 3: The schematic of uniform crossover procedure preserving Hamming weight implemented in nimCSO.

111 The above are applied iteratively, with best solutions carried to next generation, until the
112 solution converges or the maximum number of iterations is reached. Unlike the other methods,
113 this one is not limited by the number of components and lets user control both time and
114 memory requirements, either to make big problems feasible or to get a good-enough solution
115 quickly in small problems. However, it comes with no optimality guarantees.

116 **Use Examples**

117 Tracking up to 11.9 million solutions.

Method	Time (s)	Memory (MB)
nimCSO (-d:release -threads:on)	302s	488 MB
nimCSO (-d:danger -threads:off)	302s	488 MB
NumPy (Python 3.11)	302s	488 MB
Dict Python 3.11	302s	488 MB

118 **Contributions**

119 A.M.K. was responsible for conceptualization, methodology, software, testing and validation,
120 writing of manuscript, and visualization; A.D. was responsible for testing software and results
121 in training machine learning models; A.M.B., W.F.R., Z-K.L. were responsible for funding
122 acquisition, review, and editing. Z-K.L. was also supervising the work.

123 **Acknowledgements**

124 This work has been funded through grants: **DOE-ARPA-E DE-AR0001435, NSF-POSE**
125 **FAIN-2229690, and ONR N00014-23-2721.**

126 We would also like to acknowledge Dr. Jonathan Siegel at Texas A&M University for valuable
127 discussions and feedback on the project.

References

- Cantor, B., Chang, I. T. H., Knight, P., & Vincent, A. J. B. (2004). Microstructural development in equiatomic multicomponent alloys. *Materials Science and Engineering A*, 375–377, 213–218. <https://doi.org/10.1016/j.msea.2003.10.257>
- Debnath, A., Krajewski, A. M., Sun, H., Lin, S., Ahn, M., Li, W., Priya, S., Singh, J., Shang, S., Beese, A. M., Liu, Z.-K., & Reinhart, W. F. (2021). *Journal of Materials Informatics*, 1. <https://doi.org/10.20517/jmi.2021.05>
- Evans, M. L., Bergsma, J., Merkys, A., Andersen, C. W., Andersson, O. B., Beltrán, D., Blokhin, E., Boland, T. M., Balderas, R. C., Choudhary, K., Díaz, A. D., García, R. D., Eckert, H., Eimre, K., Montero, M. E. F., Krajewski, A. M., Mortensen, J. J., Duarte, J. M. N., Pietryga, J., ... Armiento, R. (2024). *Developments and applications of the OPTIMADE API for materials discovery, design, and data exchange*. <https://doi.org/10.48550/arXiv.2402.00572>
- Goldberg, D. E. (1989). *Genetic algorithms in search, optimization and machine learning* (1st ed.). Addison-Wesley Longman Publishing Co., Inc. ISBN: 0201157675
- Knuth, D. E. (2009). *The art of computer programming, volume 4, fascicle 1: Bitwise tricks & techniques; binary decision diagrams* (12th ed.). Addison-Wesley Professional. ISBN: 0321580508
- Krajewski, A. M., Beese, A. M., Reinhart, W. F., & Liu, Z.-K. (2024). *Efficient generation of grids and traversal graphs in compositional spaces towards exploration and path planning exemplified in materials*. <http://arxiv.org/abs/2402.03528>
- Li, W., Raman, L., Debnath, A., Ahn, M., Lin, S., Krajewski, A. M., Shang, S., Priya, S., Reinhart, W. F., Liu, Z.-K., & Beese, A. M. (2024). *Design and validation of refractory alloys using machine learning, CALPHAD, and experiments*. <https://doi.org/https://dx.doi.org/10.2139/ssrn.4689687>
- Rumpf, A. (2023). *Nim programming language v2.0.0*. <https://nim-lang.org/>
- Senkov, O. N., Miracle, D. B., Chaput, K. J., & Couzinie, J.-P. (2018). Development and exploration of refractory high entropy alloys—a review. *Journal of Materials Research*, 33, 3092–3128. <https://doi.org/10.1557/jmr.2018.153>
- Williams, J. W. J. (1964). Algorithm 232 - heapsort. *Communications of the ACM*, 7, 347–349. <https://doi.org/10.1145/512274.512284>
- Yeh, J. W., Chen, S. K., Lin, S. J., Gan, J. Y., Chin, T. S., Shun, T. T., Tsau, C. H., & Chang, S. Y. (2004). Nanostructured high-entropy alloys with multiple principal elements: Novel alloy design concepts and outcomes. *Advanced Engineering Materials*, 6, 299–303. <https://doi.org/10.1002/adem.200300567>