

maze-dataset: Maze Generation with Algorithmic Variety and Representational Flexibility

Michael Igorevich Ivanitskiy¹, Aaron Sandoval⁴, Alex F. Spies²,
Tilman R  uker³, Brandon Knutson¹, Cecilia Diniz Behn¹, and Samy
Wu Fung¹

¹ Colorado School of Mines, Department of Applied Mathematics and Statistics ² Imperial College
London ³ UnSearch.org ⁴ Independent ¶ 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Published: unpublished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

Solving mazes is a classic problem in computer science and artificial intelligence, and humans have been constructing mazes for thousands of years. Although finding the shortest path through a maze is a solved problem, this very fact makes it an excellent testbed for studying how machine learning algorithms solve problems and represent spatial information. We introduce maze-dataset, a user-friendly Python library for generating, processing, and visualizing datasets of mazes. This library supports a variety of maze generation algorithms providing mazes with or without loops, mazes that are connected or not, and many other variations. These generation algorithms can be configured with various parameters, and the resulting mazes can be filtered to satisfy desired properties. Also provided are tools for converting mazes to and from various formats suitable for a variety of neural network architectures, such as rasterized images, tokenized text sequences, and various visualizations. As well as providing a simple interface for generating, storing, and loading these datasets, maze-dataset is extensively tested, type hinted, benchmarked, and documented.

```
cfg = MazeDatasetConfig(
    name = "test",
    grid_n = 5,
    n_mazes = 1,
    maze_ctor = gen_dfs,
    ... # many, many options
)
```

m:

```
ds = MazeDataset.from_config(cfg)
```

```
<ADJLIST_START> (2,4) <--> (3,4) ; (1,3) <--> (1,4) ;
(1,4) <--> (0,4) ; (2,0) <--> (3,0) ; (4,4) <--> (4,3)
(3,1) <--> (3,2) ; (3,0) <--> (4,0) ; (4,1) <--> (4,2)
(0,2) <--> (0,3) ; (2,2) <--> (1,2) ; (0,2) <--> (0,1)
(3,3) <--> (3,2) ; (1,0) <--> (0,0) ; (0,2) <--> (1,2)
(4,0) <--> (4,1) ; (2,1) <--> (2,0) ; (3,2) <--> (2,2)
(3,4) <--> (4,4) ; <ADJLIST_END>
```

```
<ORIGIN_START> (1,3) <ORIGIN_END> <TARGET_START> (2,3)
```

```
<PATH_START> (1,3) (0,3) (0,2) (1,2) (2,2) (2,1) (2,0)
(4,2) (4,3) (4,4) (3,4) (2,4) (2,3) <PATH_END>
```

Figure 1: Usage of maze-dataset. We create a MazeDataset from a MazeDatasetConfig. This contains SolvedMaze objects which can be converted to and from a variety of formats. Code in the image contains clickable links to [documentation](#). A variety of generated examples can be viewed [here](#).

Statement of Need

While maze generation itself is straightforward, the architectural challenge comes from building a system supporting many algorithms with configurable parameters, property filtering, and representation transformation. This library aims to greatly streamline the process of generating and working with datasets of mazes that can be described as subgraphs of an $n \times n$ lattice with boolean connections and, optionally, start and end points that are nodes in the graph. Furthermore, we place emphasis on a wide variety of possible text output formats aimed at evaluating the spatial reasoning capabilities of Large Language Models (LLMs) and other text-based transformer models.

For interpretability and behavioral research, algorithmic tasks offer benefits by allowing systematic data generation and task decomposition, as well as simplifying the process of circuit discovery (Räuker et al., 2023). Although mazes are well suited for these investigations, we found that existing maze generation packages (Cobbe et al., 2019; Ehsan, 2022; Harries et al., n.d.; Németh, 2019; Schwarzschild, Borgnia, Gupta, Bansal, et al., 2021) lack support for transforming between multiple representations and provide limited control over the maze generation process.

Related Works

A multitude of public and open-source software packages exist for generating mazes (Ehsan, 2022; Németh, 2019; Schwarzschild, Borgnia, Gupta, Bansal, et al., 2021). However, nearly all of these packages produce mazes represented as rasterized images or other visual formats rather than the underlying graph structure, and this makes it difficult to work with these datasets.

- Most prior works provide mazes in visual or raster formats, and we provide a variety of similar output formats:
 - `RasterizedMazeDataset`, utilizing `as_pixels()`, which can exactly mimic the outputs provided in `easy-to-hard-data` (Schwarzschild, Borgnia, Gupta, Bansal, et al., 2021) and can be configured to be similar to the outputs of Németh (2019)
 - `as_ascii()` provides a format similar to (Oppenheim, 2018; Singla, 2023)
 - `MazePlot` provides a feature-rich plotting utility with support for multiple paths, heatmaps over positions, and more. This is similar to the outputs of (Alance AB, 2019; Ehsan, 2022; Guo et al., 2011; Nag, 2020)
- The text format provided by `SolvedMaze(...).as_tokens()` is similar to that of (Liu & Wu, 2023) but with many more options, detailed in [section: Tokenized Output Formats](#).
- Preserving metadata about the generation algorithm with the dataset itself is essential for studying the effects of distributional shifts. Our package efficiently stores the dataset along with its metadata in a single human-readable file (M. Ivanitskiy, n.d.). As far as we are aware, no existing packages do this reliably.
- Storing mazes as images or adjacency matrices is not only difficult to work with, but also inefficient. We use a highly efficient method detailed in [section: Implementation](#).
- Our package is easily installable with source code freely available. It is extensively tested, type hinted, benchmarked, and documented. Many other maze generation packages lack this level of rigor and scope, and some (Ayaz et al., 2008) appear to simply no longer be accessible.

65

Features

66 We direct readers to our [examples](#), [docs](#), and [notebooks](#) for more information. Our package
67 can be installed from [PyPi](#) via `pip install maze-dataset`, or directly from the [git repository](#)
68 ([Michael I. Ivanitskiy et al., 2023a](#)).

69 Datasets of mazes are created from a `MazeDatasetConfig` configuration object, which allows
70 specifying the number of mazes, their size, the generation algorithm, and various parameters for
71 the generation algorithm. Datasets can also be filtered after generation to satisfy certain prop-
72 erties. Custom filters can be specified, and some filters are included in `MazeDatasetFilters`.

73

Visual Output Formats

74 Internally, mazes are `SolvedMaze` objects, which have path information and a tensor optimized
75 for storing sub-graphs of a lattice. These objects can be converted to and from several formats,
76 shown in [Figure 2](#), to maximize their utility in different contexts.

77 In previous work, maze tasks have been used with Recurrent Convolutional Neural Network
78 (RCNN) derived architectures ([Schwarzschild, Borgnia, Gupta, Huang, et al., 2021](#)). To
79 facilitate the use of our package in this context, we replicate the format of ([Schwarzschild,](#)
80 [Borgnia, Gupta, Bansal, et al., 2021](#)) and provide the `RasterizedMazeDataset` class which
81 returns rasterized pairs of (input, target) mazes as shown in [Figure 3](#).

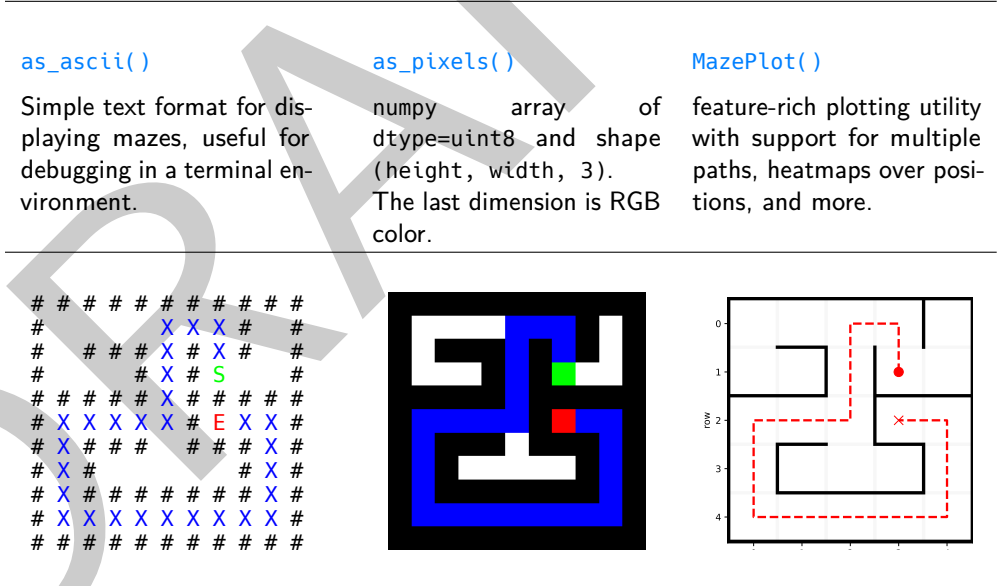


Figure 2: Various output formats. Top row (left to right): ASCII diagram, rasterized pixel grid, and advanced display tool.

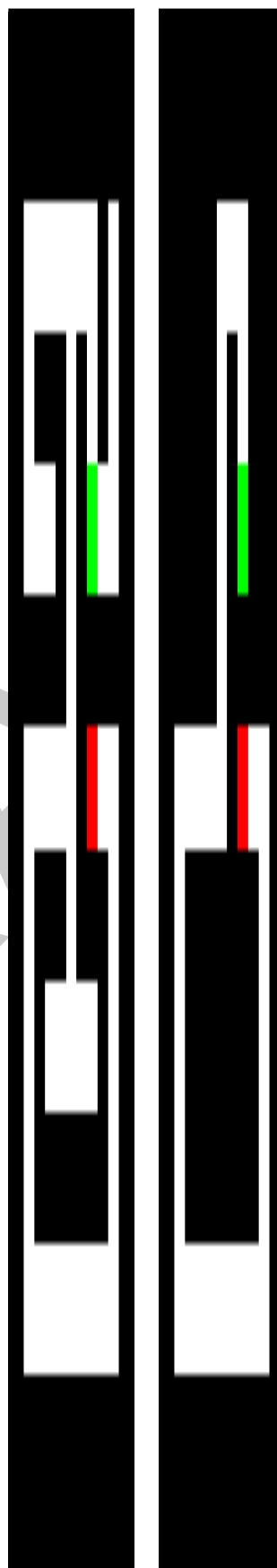


Figure 3: Input is the rasterized maze without the path marked (left), and provide as a target the maze

82 Tokenized Output Formats

83 Autoregressive transformer models can be quite sensitive to the exact format of input data,
84 and may even use delimiter tokens to perform reasoning steps (Pfau et al., 2024; Spies et
85 al., 2024). To facilitate systematic investigation of the effects of different representations of
86 data on text model performance, we provide a variety of text output formats, with an example
87 given in Figure 4. We utilize Finite State Transducers (Gallant, 2015) for efficiently storing
88 valid tokenizers.

```
<ADJLIST_START> (0,0) <--> (1,0) ; (2,0) <--> (3,0) ; (4,1) <--> (4,0) ; (2,0) <--> (2,1) ;
(1,0) <--> (1,1) ; (3,4) <--> (2,4) ; (4,2) <--> (4,3) ; (0,0) <--> (0,1) ; (0,3) <--> (0,2) ;
(4,4) <--> (3,4) ; (4,3) <--> (4,4) ; (4,1) <--> (4,2) ; (2,1) <--> (2,2) ; (1,4) <--> (0,4) ;
(1,2) <--> (0,2) ; (2,4) <--> (2,3) ; (4,0) <--> (3,0) ; (2,2) <--> (3,2) ; (1,2) <--> (2,2) ;
(1,3) <--> (0,3) ; (3,2) <--> (3,3) ; (0,2) <--> (0,1) ; (3,1) <--> (3,2) ; (1,3) <--> (1,4) ;
<ADJLIST_END> <ORIGIN_START> (1,3) <ORIGIN_END> <TARGET_START> (2,3) <TARGET_END>
<PATH_START> (1,3) (0,3) (0,2) (1,2) (2,2) (2,1) (2,0) (3,0) (4,0) (4,1) (4,2) (4,3) (4,4)
(3,4) (2,4) (2,3) <PATH_END>
```

Figure 4: Example text output format with token regions highlighted. **Adjacency list** : text representation of the graph, **Origin** : starting coordinate, **Target** : ending coordinate, **Path** : maze solution sequence. By passing an instance of `MazeTokenizerModular` to `as_tokens(...)` a maze can be converted to a text sequence. The `MazeTokenizerModular` class contains a rich set of options with 19 discrete parameters, resulting in over 5.8 million unique possible tokenizers.

89 Benchmarks

90 We benchmark for generation time across various configurations in Table 1 and Figure 5.
91 Experiments were performed on a [standard GitHub runner](#) without parallelism. Additionally,
92 maze generation under certain constraints may not always be successful, and for this we provide
93 a way to estimate the success rate of a given configuration, described in Figure 6.

maze_ctor	keyword args	all sizes	small $g \leq 10$	medium $g \in (10, 32]$	large $g > 32$
dfs		28.0	2.8	20.3	131.8
dfs	accessible_cells=20	2.3	2.2	2.4	2.2
dfs	do_forks=False	2.7	2.2	3.1	3.5
dfs	max_tree_depth=0.5	2.5	2.0	2.7	4.0
dfs_percolation	p=0.1	43.9	2.8	33.9	208.0
dfs_percolation	p=0.4	48.7	3.0	36.5	233.5
kruskal		12.8	1.9	10.3	55.8
percolation	p=1.0	50.2	2.6	37.2	242.5
recursive_div		10.2	1.7	8.9	42.1
wilson		676.5	7.8	188.6	3992.6
mean		559.9	13.0	223.5	3146.9
median		11.1	6.5	32.9	302.7

Table 1: Generation times in milliseconds for various algorithms and maze sizes. More information can be found on the [benchmarks page](#).

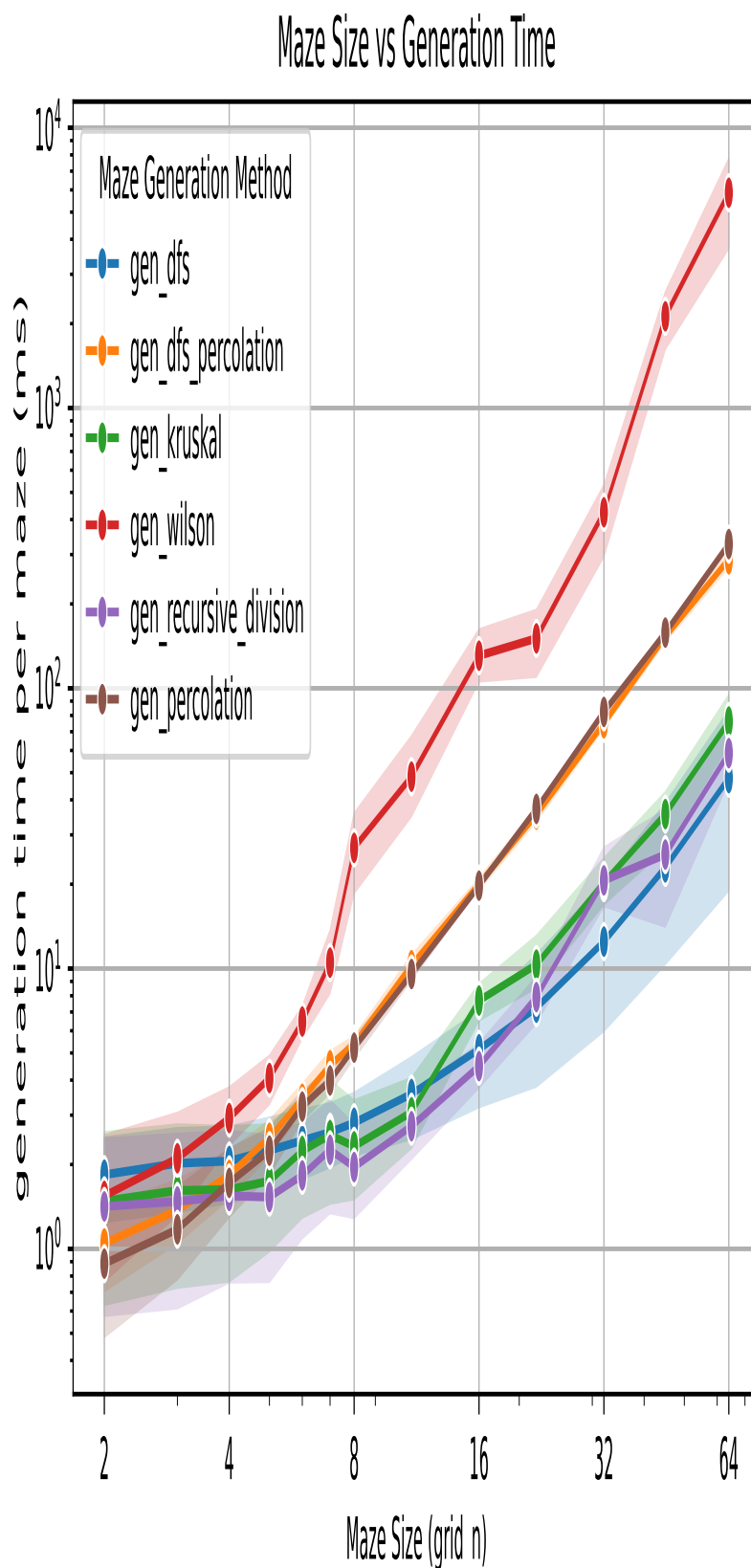


Figure 5: Plot of maze generation time. Generation time scales exponentially with maze size for all algorithms. Generation time per maze does not depend on the number of mazes being generated, and there is minimal overhead to initializing the generation process for a small dataset. Wilson's algorithm is notably less efficient than others and has high variance. Note that values are averaged across all parameter sets for that algorithm. More information can be found on the [benchmarks page](#).

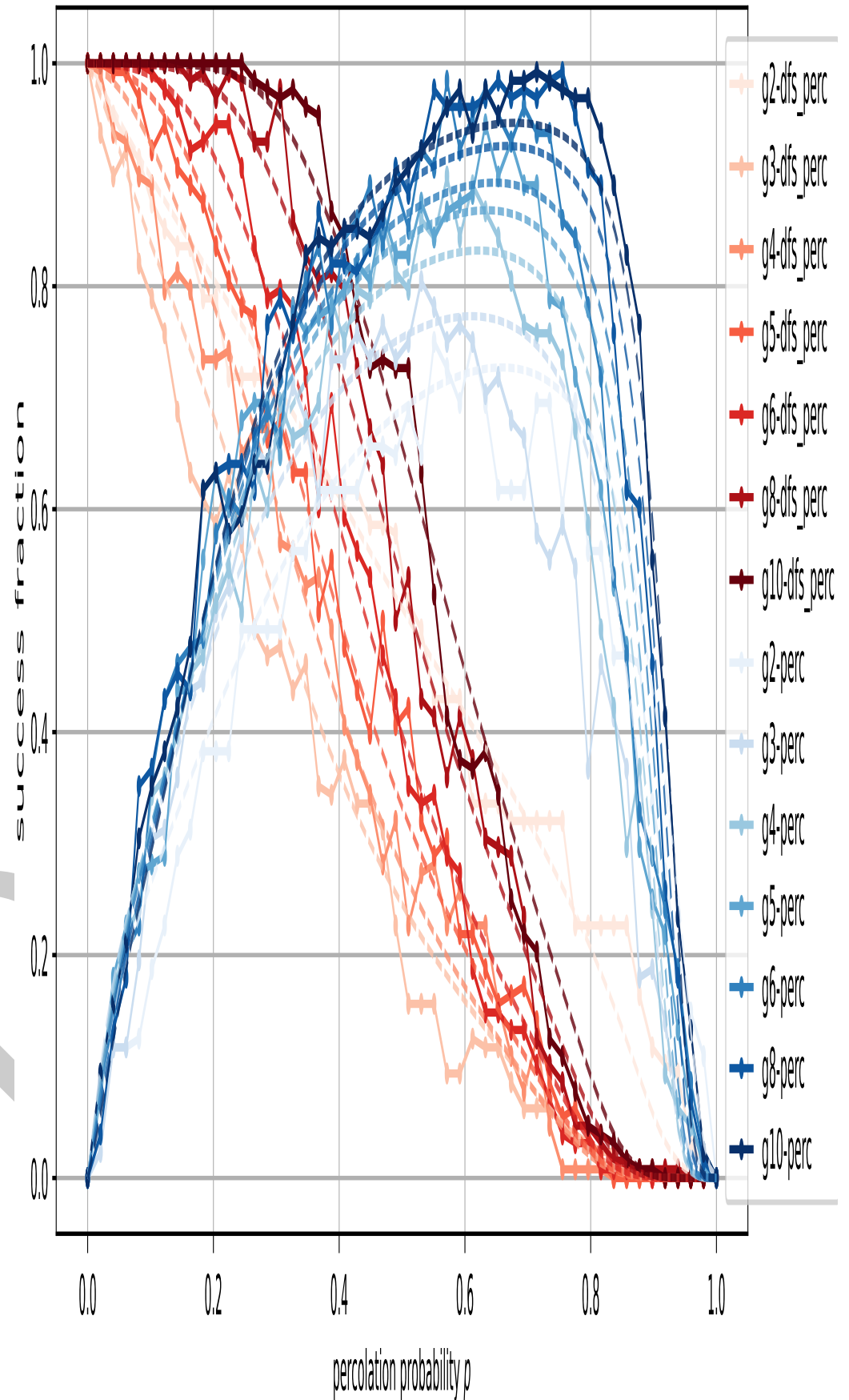


Figure 6: In order to replicate the exact dataset distribution of `?`, the parameter `start_endpoint` in `MazeGeneration` allows for additional constraints such as enforcing that the start and end points be in a “dead end” with only one accessible neighbor cell. However, combining these constraints with cyclic mazes can lead to an absence of valid start and end points. To deal with this, our package provides a way to estimate the success rate of a given configuration using a symbolic regression model trained with PySR `?`. An example of both empirical and predicted success rates as a function of the percolation probability p for various maze sizes, percolation with and without depth first search,

Implementation

Using an adjacency matrix for storing mazes would be memory inefficient by failing to exploit the highly sparse structure, while using an adjacency list could lead to a poor lookup time. This package utilizes a simple, efficient representation of mazes as subgraphs of a finite lattice, detailed in ??, which we call a [LatticeMaze](#).

Our package is implemented in Python([Rossum, 1995](#)), and makes use of the extensive scientific computing ecosystem, including NumPy ([Harris et al., 2020](#)) for array manipulation, plotting tools ([Hunter, 2007](#); [Waskom, 2021](#)), Jupyter notebooks ([Kluyver et al., 2016](#)), and PySR ([Cranmer, 2023](#)) for symbolic regression.

Usage in Research

This package was originally built for the needs of the ([Michael I. Ivanitskiy et al., 2023b](#)) project, which aims to investigate spatial planning and world models in autoregressive transformer models trained on mazes ([Michael Igorevich Ivanitskiy, Spies, et al., 2023](#); [Michael Igorevich Ivanitskiy, Shah, et al., 2023](#); [Spies et al., 2024](#)). It was extended for work on understanding the mechanisms by which recurrent convolutional and implicit networks ([Fung et al., 2022](#)) solve mazes given a rasterized view ([Knutson et al., 2024](#)), which required matching the pixel-padded and endpoint constrained output format of ([Schwarzschild, Borgnia, Gupta, Bansal, et al., 2021](#)). Ongoing work using maze-dataset aims to investigate the effects of varying the tokenization format on the performance of pretrained LLMs on spatial reasoning.

This package has also been utilized in work by other groups:

- By ([Nolte et al., 2024](#)) to compare the effectiveness of transformers trained with the MLM- \mathcal{U} ([Kitouni et al., 2024](#)) multistep prediction objective against standard autoregressive training for multi-step planning on our maze task.
- By ([Wang et al., 2024](#)) and ([Chen et al., 2024](#)) to study the effectiveness of imperative learning.
- By ([Zhang et al., 2025](#)) to introduce a novel framework for reasoning diffusion models.
- By ([Dao & Vu, 2025](#)) to improve spatial reasoning in LLMs with GRPO.

Acknowledgements

This work was partially funded by National Science Foundation awards DMS-2110745 and DMS-2309810. We are also grateful to LTFF and FAR Labs for hosting authors MII, AFS, and TR for a residency visit, and to various members of FAR's technical staff for their advice.

This work was partially supported by AI Safety Camp and AI Safety Support, which also brought many of the authors together. We would like to thank our former collaborators at AI Safety Camp and other users and contributors to the maze-dataset package: Benji Berczi, Guillaume Corlouer, William Edwards, Leon Eshuijs, Chris Mathwin, Lucia Quirke, Can Rager, Adrians Skapars, Rusheb Shah, Johannes Treutlein, and Dan Valentine.

We thank the Mines Optimization and Deep Learning group (MODL) for fruitful discussions. We also thank Michael Rosenberg for recommending the usage of Finite State Transducers for storing tokenizer validation information.

- <https://doi.org/10.48550/arXiv.2312.02566>
- Kitouni, O., Nolte, N. S., Williams, A., Rabbat, M., Bouchacourt, D., & Ibrahim, M. (2024). The factorization curse: Which tokens you predict underlie the reversal curse and more. *Advances in Neural Information Processing Systems*, 37, 112329–112355. <https://doi.org/10.48550/arXiv.2406.05183>
- Cluyver, T., Ragan-Kelley, B., Perez, F., Granger, B., Bussonnier, M., Frederic, J., Kelley, K., Hamrick, J., Grout, J., Corlay, S., Ivanov, P., Avila, D., Abdalla, S., & Willing, C. (2016). Jupyter notebooks - a publishing format for reproducible computational workflows. *Proceedings of the 20th International Conference on Electronic Publishing*, 87–90. <https://doi.org/10.3233/978-1-61499-649-1-87>
- Knutson, B., Rabeendran, A. C., Ivanitskiy, M., Pettyjohn, J., Diniz-Behn, C., Fung, S. W., & McKenzie, D. (2024). On logical extrapolation for mazes with recurrent and implicit networks. *arXiv Preprint arXiv:2410.03020*. <https://doi.org/10.48550/arXiv.2410.03020>
- Liu, C., & Wu, B. (2023). Evaluating large language models on graphs: Performance insights and comparative analysis. *arXiv Preprint arXiv:2308.11224*. <https://doi.org/10.48550/arXiv.2308.11224>
- Nag, A. (2020). MDL suite: A language, generator and compiler for describing mazes. *Journal of Open Source Software*, 5(46), 1815. <https://doi.org/10.21105/joss.01815>
- Németh, F. (2019). *Maze-generation-algorithms*. <https://github.com/ferenc-nemeth/maze-generation-algorithms>
- Nolte, N., Kitouni, O., Williams, A., Rabbat, M., & Ibrahim, M. (2024). Transformers can navigate mazes with multi-step prediction. *arXiv Preprint arXiv:2412.05117*. <https://doi.org/10.48550/arXiv.2412.05117>
- Oppenheim, J. (2018). *Maze-generator: Generate a random maze represented as a 2D array using depth-first search*. <https://github.com/oppenheimj/maze-generator/>; GitHub.
- Pfau, J., Merrill, W., & Bowman, S. R. (2024). Let's think dot by dot: Hidden computation in transformer language models. *arXiv Preprint arXiv:2404.15758*. <https://doi.org/10.48550/arXiv.2404.15758>
- Räuker, T., Ho, A., Casper, S., & Hadfield-Menell, D. (2023). Toward transparent ai: A survey on interpreting the inner structures of deep neural networks. *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, 464–483. <https://doi.org/10.1109/satml54575.2023.00039>
- Rossum, G. van. (1995). *Python reference manual* (CS-R9525). Centrum voor Wiskunde; Informatica (CWI). <https://ir.cwi.nl/pub/5008/05008D.pdf>
- Schwarzschild, A., Borgnia, E., Gupta, A., Bansal, A., Emam, Z., Huang, F., Goldblum, M., & Goldstein, T. (2021). *Datasets for Studying Generalization from Easy to Hard Examples* (No. arXiv:2108.06011). arXiv. <https://doi.org/10.48550/arXiv.2108.06011>
- Schwarzschild, A., Borgnia, E., Gupta, A., Huang, F., Vishkin, U., Goldblum, M., & Goldstein, T. (2021). Can you learn an algorithm? Generalizing from easy to hard problems with recurrent networks. *Advances in Neural Information Processing Systems*, 34, 6695–6706. <https://doi.org/10.48550/arXiv.2106.04537>
- Singla, A. (2023). Evaluating ChatGPT and GPT-4 for visual programming. *arXiv Preprint arXiv:2308.02522*. <https://doi.org/10.48550/arXiv.2308.02522>
- Spies, A. F., Edwards, W., Ivanitskiy, M. I., Skapars, A., Räuker, T., Inoue, K., Russo, A., & Shanahan, M. (2024). Transformers use causal world models in maze-solving tasks. *arXiv Preprint arXiv:2412.11867*. <https://doi.org/10.48550/arXiv.2412.11867>
- Wang, C., Ji, K., Geng, J., Ren, Z., Fu, T., Yang, F., Guo, Y., He, H., Chen, X., Zhan, Z., &

- 225 others. (2024). Imperative learning: A self-supervised neural-symbolic learning framework
226 for robot autonomy. *arXiv Preprint arXiv:2406.16087*. [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2406.16087)
227 [2406.16087](https://doi.org/10.48550/arXiv.2406.16087)
- 228 Waskom, M. L. (2021). seaborn: Statistical data visualization. *Journal of Open Source*
229 *Software*, 6(60), 3021. <https://doi.org/10.21105/joss.03021>
- 230 Zhang, T., Pan, J.-S., Feng, R., & Wu, T. (2025). T-SCEND: Test-time scalable MCTS-
231 enhanced diffusion model. *arXiv Preprint arXiv:2502.01989*. [https://doi.org/10.48550/](https://doi.org/10.48550/arXiv.2502.01989)
232 [arXiv.2502.01989](https://doi.org/10.48550/arXiv.2502.01989)

DRAFT