# DeepCAVE: A Visualization and Analysis Tool for Automated Machine Learning

**Sarah Segel**[1]                                                   S.SEGEL@AI.UNI-HANNOVER.DE

**Helena Graf**[1]                                                   H.GRAF@AI.UNI-HANNOVER.DE

**Edward Bergman**[2]                                         BERGMANE@CS.UNI-FREIBURG.DE

**Kristina Thieme**[1]                                 KRISTINA.THIEME@STUD.UNI-HANNOVER.DE

**Marcel Wever**[1,4]                                            M.WEVER@AI.UNI-HANNOVER.DE

**Alexander Tornede**[1]                                     A.TORNEDE@AI.UNI-HANNOVER.DE

**Frank Hutter**[2,3]                                                      FH@CS.UNI-FREIBURG.DE

**Marius Lindauer**[1,4]                                    M.LINDAUER@AI.UNI-HANNOVER.DE

[1]*Leibniz University Hannover,* [2]*University of Freiburg,* [3]*ELLIS Institute Tübingen,* [4]*L3S Research Center*

**Editor:** Zeyi Wen

## Abstract

Hyperparameter optimization (HPO), as a central paradigm of AutoML, is crucial for leveraging the full potential of machine learning (ML) models; yet its complexity poses challenges in understanding and debugging the optimization process. We present DeepCAVE, a tool for interactive visualization and analysis, providing insights into HPO. Through an interactive dashboard, researchers, data scientists, and ML engineers can explore various aspects of the HPO process and identify issues, untouched potentials, and new insights about the ML model being tuned. By empowering users with actionable insights, DeepCAVE contributes to the interpretability of HPO and ML on a design level and aims to foster the development of more robust and efficient methodologies in the future.

**Keywords:**   Automated machine learning, hyperparameter optimization, interpretable machine learning, human-centered machine learning, visualization

## 1 Introduction

Modern methods for hyperparameter optimization (HPO) encompass a diverse range of approaches beyond grid and random search (Turner et al., 2021; Bischl et al., 2023). However, their lack of insight and transparency can lead to a lack of trust by users (Drozdal et al., 2020) and make debugging such methods and their application difficult. We introduce Deep-CAVE, a visualization and analysis tool for HPO, aiding in understanding and debugging the application of HPO. At present, HPO is our main priority, wherein metadata in the format of configurations and performance values is generated in an iterative process, as many automated machine learning (AutoML) packages do, such as Auto-Sklearn (Feurer et al., 2022), TPOT (Olson and Moore, 2019), AutoPyTorch (Zimmer et al., 2021) or Auto-Keras (Jin et al., 2019). Similarly, leveraging DeepCAVE for neural architecture search (NAS) (Elsken
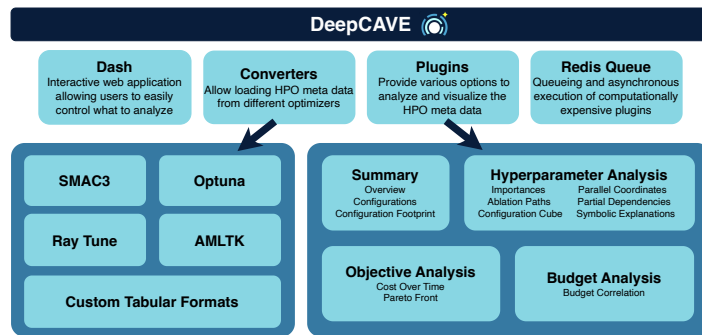
arXiv:2512.01810v1 [cs.LG] 1 Dec 2025

Figure 1: Overview of DeepCAVE's components and their functionalities.

et al., 2022) or combined algorithm selection and hyperparameter optimization (Thornton et al., 2013) is possible by encoding architectural decisions, algorithmic choices, or machine learning pipelines as part of the configuration. Advances in machine learning and HPO, including multi-objective or multi-fidelity methods, challenge interpretation and analysis tools. DeepCAVE addresses this by offering support for multiple objectives and fidelities, with specific visualizations tailored to the different scenarios. This tightly interacts with the vision of a more human-centered approach of AutoML (Lindauer et al., 2024), in which users can learn from HPO and ultimately even integrate their expertise into the system (Hvarfner et al., 2022; Mallik et al., 2023). Moreover, to facilitate a reproducible and transparent research process, users can export visualizations to incorporate into their research papers, thereby providing insight into their HPO process and optimized hyperparameters.

## 2 Features and Usage

DeepCAVE's architecture relies upon various components, as illustrated in Figure 1. The tool is exclusively written in Python, given its prevalence as the predominant programming language in machine learning (Raschka et al., 2020). Built on the Dash framework (Plotly Technologies Inc., 2015), it offers users a fully interactive environment within a web browser.

### 2.1 Converters

DeepCAVE utilizes run objects as a basic unit for data interpretation, where a run can be considered an HPO process consisting of a list of trials, each representing a hyperparameter configuration with associated objective value, budget, and seed. The interface allows for selection and grouping of runs, streamlining the analysis of many HPO runs at once. Converters are used to access the optimizer data stored in the file system and transfer them into a run object. They keep track of both finished and running optimization processes that consistently write new results to disk. At the time of writing, DeepCAVE supports loading runs created by the HPO packages SMAC3 (Lindauer et al., 2022), Ray Tune (Liaw et al., 2018), and Optuna (Akiba et al., 2019), as well as, the AutoML framework AMLTK (Bergman et al., 2024). Furthermore, to allow loading runs from arbitrary optimizers, DeepCAVE offers a generic tabular input format, as well as a programmatic interface.
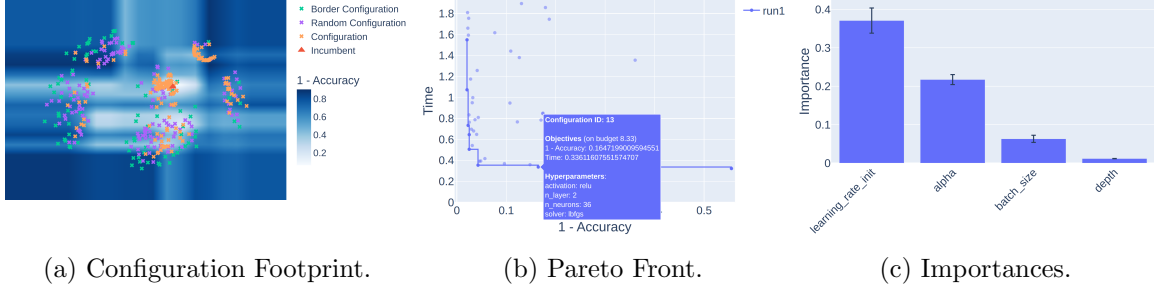
(a) Configuration Footprint.　　(b) Pareto Front.　　(c) Importances.

Figure 2: Examples of plots produced via DeepCAVE's plugins. A mouseover allows obtaining additional information, and clicking on single configurations in the plot opens the Configurations plugin, providing details regarding that configuration.

## 2.2 Analysis Plugins

DeepCAVE's strength lies in its modular plugin structure, providing diverse insights into the HPO run. Plugins allow analyzing overall runs, objectives, hyperparameters, and budgets through texts, tables, and Plotly (Plotly Technologies Inc., 2015) visualizations. Dynamic filtering enables the interactive selection of certain aspects, such as specific objectives or budgets. For computationally intensive plugins, Redis Queue (Redis Inc., 2024) is leveraged for efficient queueing and asynchronous execution, ensuring a responsive user experience.

**Overall Optimization Analysis** The *Overview plugin* offers a holistic view, showing optimizer choice, configuration space, and trial status. Detailed analysis of individual configurations can be performed via the *Configurations plugin*. Employing multi-dimensional scaling (Kruskal, 1964a,b) for dimensionality reduction, the *Configuration Footprint plugin* visually represents the optimizer's coverage of the configuration space (Biedenkapp et al., 2018), as shown in Figure 2a. This can help to identify limited exploration, suggesting the optimizer might be stuck in a local optimum.

**Objective Analysis** The *Cost Over Time plugin* visualizes optimization convergence, either over time or the number of trials evaluated, allowing to compare different optimization strategies such as the optimizer chosen, search space design, or other meta-level decisions. The *Pareto Front plugin* visually represents configurations that achieve different tradeoffs between conflicting objectives, such as runtime versus final loss, as shown in Figure 2b.

**Hyperparameter Analysis** As the relationship between hyperparameter configurations and objective values constitutes a crucial aspect of HPO, DeepCAVE provides a variety of plugins in this regard. *Parallel Coordinates* plots visualize hyperparameter configurations as lines, connecting their hyperparameter values and corresponding final scores (Golovin et al., 2017). Additionally, *Partial Dependence* plots (Friedman, 2001; Moosbauer et al., 2021) offer insights into the average marginal effect of single hyperparameters on the objective value. Moreover, *Symbolic Explanations* (Segel et al., 2023) provide explicit formulas capturing the relationship between hyperparameter and objective values. To understand the impact of individual hyperparameters on the objective, the *Importances plugin* can be used. With fANOVA (Hutter et al., 2014), see Figure 2c, and local (hyper)parameter importance (LPI) (Biedenkapp et al., 2018), users can explore both global and local importance.

3

*Ablation Paths* (Biedenkapp et al., 2017) provide insights into the impact of sequentially adjusting individual hyperparameter values towards the incumbent configuration. In practice, DeepCAVE has been used to analyze optimization runs with up to 39 hyperparameters, confirming its scalability for high-dimensional settings.

**Budget Analysis** Effective use of multiple fidelities crucially relies on a consistent rank or performance correlation to effectively allocate more resources to better configurations. The *Budget Correlation plugin* directly visualizes this correlation between fidelities, allowing users to assess the utility of their multiple fidelities.

Although DeepCAVE provides diverse analysis plugins and visualizations, their customization and composition are currently limited, outlining interesting future work.

## 3 Existing Tools for Explainable AutoML

With the rising interest in explainability and insight into optimization processes, a number of tools aiming to fill this gap have emerged. Among these, XAutoML (Zöller et al., 2023) focuses on AutoML and interpretability. Similar to DeepCAVE, it provides a number of different interactive visualizations. Notably, it allows users to extend its capabilities through custom "cards", similar to plugins in DeepCAVE, and supports a range of common (Auto)ML libraries. In contrast to our tool, it lacks support for (parallel) background computation and analysis of running HPO processes. Further, only one HPO run is considered at a time, making comparisons between runs cumbersome. Beyond XAutoML, platforms like RapidMiner (Altair Engineering Inc., 2024), Google Vizier (Golovin et al., 2017), and ATMseer (Wang et al., 2018) cater to specific facets of AutoML processes, such as algorithm selection, experiment tracking, and pipeline analysis, albeit often with limitations in terms of interpretability at the HPO level or support for run comparisons. Similarly, tools like Optuna (Akiba et al., 2019), IOHanalyzer (Wang et al., 2022), TensorBoard (Abadi et al., 2015), WandB (Biewald, 2020), MLFlow (Chen et al., 2020), MLJar (Płońska and Płoński, 2021), SigOpt (Clark and Hayes, 2019), IBM AutoAI (Wang et al., 2020), Pipeline Profiler (Ono et al., 2021), iraceplot (López-Ibáñez et al., 2025), and Boxer (Gleicher et al., 2020) contribute valuable insights through various forms of run analysis, visualization, and reporting. However, many of these tools depend on specific formats, lack extensibility, or fail to support comparisons across multiple HPO runs. DeepCAVE distinguishes itself in this field with its focus on HPO, interpretability, and interactivity. Its unique proposition lies in its ability to facilitate detailed comparisons across different HPO runs. Moreover, its extensive format compatibility, efficient computational parallelization, and browser-based accessibility mitigate major usability challenges found in comparable tools.

## 4 Conclusion

DeepCAVE is a novel tool for interactively analyzing and explaining HPO runs. It allows performing comprehensive analyses of HPO runs concerning their hyperparameters, objectives, and budgets via an interactive dashboard running in the web browser. DeepCAVE is accessible under Apache License Version 2.0[1] and can be installed via the instructions in the GitHub repository `https://github.com/automl/DeepCAVE`.

---

1. `https://www.apache.org/licenses`

**Acknowledgments and Disclosure of Funding**

**References**

M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL `https://www.tensorflow.org`.

T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proc. of KDD'19*, 2019.

Altair Engineering Inc. Altair Rapidminer data analytics and AI platform, 2024. URL `https://altair.com/altair-rapidminer`.

E. Bergman, M. Feurer, A. Bahram, A. R. Balef, L. Purucker, S. Segel, M. Lindauer, F. Hutter, and K. Eggensperger. AMLTK: A modular AutoML toolkit in Python. *Journal of Open Source Software*, 9(100):6367, 2024.

A. Biedenkapp, M. Lindauer, K. Eggensperger, C. Fawcett, H. Hoos, and F. Hutter. Efficient parameter importance analysis via ablation with surrogates. In S. Singh and S. Markovitch, editors, *Proceedings of the Thirty-First Conference on Artificial Intelligence (AAAI'17)*, pages 773–779. AAAI Press, 2017.

A. Biedenkapp, J. Marben, M. Lindauer, and F. Hutter. CAVE: Configuration assessment, visualization and evaluation. In R. Battiti, M. Brunato, I. Kotsireas, and P. Pardalos, editors, *Proceedings of the International Conference on Learning and Intelligent Optimization (LION)*, Lecture Notes in Computer Science. Springer, 2018.

L. Biewald. Experiment tracking with Weights and Biases, 2020. URL https://www.wandb.com/.

B. Bischl, M. Binder, M. Lang, T. Pielok, J. Richter, S. Coors, J. Thomas, T. Ullmann, M. Becker, A.-L. Boulesteix, D. Deng, and M. Lindauer. Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, page e1484, 2023.

A. Chen, A. Chow, A. Davidson, A. DCunha, A. Ghodsi, S. Hong, A. Konwinski, C. Mewald, S. Murching, T. Nykodym, P. Ogilvie, M. Parkhe, A. Singh, F. Xie, M. Zaharia, R. Zang, J. Zheng, and C. Zumar. Developments in MLflow: A system to accelerate the machine learning lifecycle. In *Proc. of International Workshop on Data Management for End-to-End Machine Learning*, 2020.

S. Clark and P. Hayes. Sigopt web page, 2019. URL https://sigopt.org.

J. Drozdal, J. Weisz, D. Wang, G. Dass, B. Yao, C. Zhao, M. J. Muller, L. Ju, and H. Su. Trust in AutoML: Exploring information needs for establishing trust in automated machine learning systems. In F. Paternò, N. Oliver, C. Conati, L. D. Spano, and N. Tintarev, editors, *Proceedings of the 25th International Conference on Intelligent User Interfaces (IUI'20)*, pages 297–307. ACM, 2020.

T. Elsken, A. Zela, J. Metzen, B. Staffler, T. Brox, A. Valada, and F. Hutter. Neural architecture search for dense prediction tasks in computer vision. *arXiv:2202.07242 [cs.CV]*, 2022.

M. Feurer, K. Eggensperger, S. Falkner, M. Lindauer, and F. Hutter. Auto-Sklearn 2.0: Hands-free automl via meta-learning. *Journal of Machine Learning Research*, 23(261): 1–61, 2022.

J. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, pages 1189–1232, 2001.

M. Gleicher, A. Barve, X. Yu, and F. Heimerl. Boxer: Interactive comparison of classifier results. *Computer Graphics Forum*, 39(3), 2020.

D. Golovin, B. Solnik, S. Moitra, G. Kochanski, J. Karro, and D. Sculley. Google Vizier: A service for black-box optimization. In S. Matwin, S. Yu, and F. Farooq, editors, *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'17)*, pages 1487–1495. ACM Press, 2017.

F. Hutter, H. Hoos, and K. Leyton-Brown. An efficient approach for assessing hyperparameter importance. In E. Xing and T. Jebara, editors, *Proceedings of the 31th International Conference on Machine Learning, (ICML'14)*, pages 754–762. Omnipress, 2014.

C. Hvarfner, D. Stoll, A. Souza, L. Nardi, M. Lindauer, and F. Hutter. $\pi$BO: Augmenting acquisition functions with user beliefs for bayesian optimization. In *Proc. of ICLR'22*, 2022.

H. Jin, Q. Song, and X. Hu. Auto-Keras: An efficient neural architecture search system. In A. Teredesai, V. Kumar, Y. Li, R. Rosales, E. Terzi, and G. Karypis, editors, *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD'19)*, pages 1946–1956. ACM Press, 2019.

J. Kruskal. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, 29:1–27, 1964a.

J. Kruskal. Nonmetric multidimensional scaling: A numerical method. *Psychometrika*, 29: 115–129, 1964b.

R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. Gonzalez, and I. Stoica. Tune: A research platform for distributed model selection and training. In R. Garnett, F. Hutter, J. Vanschoren, P. Brazdil, R. Caruana, C. Giraud-Carrier, I. Guyon, and B. Kégl, editors, *ICML workshop on Automated Machine Learning (AutoML workshop 2018)*, 2018.

M. Lindauer, K. Eggensperger, M. Feurer, A. Biedenkapp, D. Deng, C. Benjamins, T. Ruhkopf, R. Sass, and F. Hutter. SMAC3: A versatile bayesian optimization package for Hyperparameter Optimization. *Journal of Machine Learning Research*, 23(54):1–9, 2022.

M. Lindauer, F. Karl, A. Klier, J. Moosbauer, A. Tornede, A. Mueller, F. Hutter, M. Feurer, and B. Bischl. Position: A call to action for a human-centered AutoML paradigm. In *Proc. of ICML'24*, 2024.

M. López-Ibáñez, P. Oñate Marín, and L. Pérez Cáceres. *iraceplot: Plots for Visualizing the Data Produced by the 'irace' Package*, 2025. URL `https://auto-optimization.github.io/iraceplot/`.

N. Mallik, C. Hvarfner, E. Bergman, D. Stoll, M. Janowski, M. Lindauer, L. Nardi, and F. Hutter. PriorBand: Practical hyperparameter optimization in the age of deep learning. In *Proceedings of the 37th International Conference on Advances in Neural Information Processing Systems (NeurIPS'23)*. Curran Associates, 2023.

J. Moosbauer, J. Herbinger, G. Casalicchio, M. Lindauer, and B. Bischl. Explaining hyperparameter optimization via partial dependence plots. In M. Ranzato, A. Beygelzimer, K. Nguyen, P. Liang, J. Vaughan, and Y. Dauphin, editors, *Proceedings of the 35th International Conference on Advances in Neural Information Processing Systems (NeurIPS'21)*, pages 2280–2291. Curran Associates, 2021.

R. Olson and J. Moore. TPOT: A tree-based pipeline optimization tool for automating machine learning. In F. Hutter, L. Kotthoff, and J. Vanschoren, editors, *Automated Machine Learning: Methods, Systems, Challenges*, pages 151–160. Springer, 2019. Available for free at `http://automl.org/book`.

J. Ono, S. Castelo, R. Lopez, E. Bertini, J. Freire, and C. Silva. PipelineProfiler: A visual analytics tool for the exploration of AutoML pipelines. *IEEE Transactions on Visualization and Computer Graphics*, 27(2):390–400, 2021.

A. Płońska and P. Płoński. MLJAR: State-of-the-art automated machine learning framework for tabular data. version 0.10.3, 2021. URL `https://mljar.com`.

Plotly Technologies Inc. Collaborative data science, 2015. URL `https://plot.ly`.

S. Raschka, J. Patterson, and C. Nolet. Machine learning in Python: Main developments and technology trends in data science, machine learning, and artificial intelligence. *Information*, 11(4), 2020.

Redis Inc. Redis Queue, 2024. URL `https://github.com/rq/rq`.

S. Segel, H. Graf, A. Tornede, B. Bischl, and M. Lindauer. Symbolic explanations for hyperparameter optimization. In A. Faust, C. White, F. Hutter, R. Garnett, and J. Gardner, editors, *Proceedings of the Second International Conference on Automated Machine Learning*. Proceedings of Machine Learning Research, 2023.

C. Thornton, F. Hutter, H. Hoos, and K. Leyton-Brown. Auto-WEKA: combined selection and Hyperparameter Optimization of classification algorithms. In I. Dhillon, Y. Koren, R. Ghani, T. Senator, P. Bradley, R. Parekh, J. He, R. Grossman, and R. Uthurusamy, editors, *The 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'13)*, pages 847–855. ACM Press, 2013.

R. Turner, D. Eriksson, M. McCourt, J. Kiili, E. Laaksonen, Z. Xu, and I. Guyon. Bayesian optimization is superior to random search for machine learning hyperparameter tuning: Analysis of the Black-Box Optimization Challenge 2020. In H. Escalante and K. Hofmann, editors, *Proceedings of the Neural Information Processing Systems Track Competition and Demonstration*, pages 3–26. Curran Associates, 2021.

D. Wang, P. Ram, D. Weidele, S. Liu, M. Muller, J. Weisz, A. Valente, A. Chaudhary, D. Torres, H. Samulowitz, and L. Amini. Autoai: Automating the end-to-end ai lifecycle with humans-in-the-loop. In *Proceedings of the 25th International Conference on Intelligent User Interfaces (IUI'20)*, page 77–78, 2020.

H. Wang, D. Vermetten, F. Ye, C. Doerr, and T. Bäck. IOHanalyzer: Detailed performance analyses for iterative optimization heuristics. *ACM Transactions on Evolutionary Learning and Optimization*, 2(1), 2022.

Q. Wang, Y. Ming, Z. Jin, Q. Shen, D. Liu, M. J. Smith, K. Veeramachaneni, and H. Qu. ATMSeer: Increasing transparency and controllability in automated machine learning. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 2018.

L. Zimmer, M. Lindauer, and F. Hutter. Auto-Pytorch: Multi-fidelity metalearning for efficient and robust AutoDL. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43:3079–3090, 2021.

M. Zöller, W. Titov, T. Schlegel, and M. F. Huber. XAutoML: A visual analytics tool for understanding and validating automated machine learning. *ACM Transactions on Interactive Intelligent Systems*, 13(4), 2023.