# PLAID: Physics–Learning AI Datamodel

**Fabien Casenave** [1], **Xavier Roynard** [1], **and Alexandre Devaux–Rivière** [1,2]

**1** SafranTech, Safran Tech, Digital Sciences & Technologies, 78114 Magny-Les-Hameaux, France **2** EPITA, 14-16 Rue Voltaire, 94270 Le Kremlin-Bicêtre, France

## Summary

PLAID (Physics-Learning AI Datamodel) is a Python library and data format for representing, storing, and sharing physics simulation datasets for machine learning. Unlike domain-specific formats, PLAID accommodates time-dependent, multi-resolution simulations and heterogeneous meshes. The library provides a high-level API to easily load, inspect, and save data. Beyond basic I/O, PLAID includes utilities for machine-learning workflows. It provides converters to build PLAID datasets from generic tabular data, and a "Hugging Face bridge" to push/pull datasets via the Hugging Face hub. In short, PLAID couples a flexible on-disk standard with a software toolkit to manipulate physics data, addressing the needs of ML researchers in fluid dynamics, structural mechanics, and related fields in a generic fashion. Full documentation, examples and tutorials are available at plaid-lib.readthedocs.io.

## Statement of Need

Machine learning for physical systems often suffers from inconsistent data representations across different domains and simulators. Existing initiatives typically target narrow problems: e.g., separate formats for CFD or for finite-element data, and dedicated scripts to process each new dataset. This fragmentation hinders reproducibility and reuse of high-fidelity data.

PLAID addresses this gap by providing a generic, unified datamodel that can describe many physics simulation data. It leverages the CGNS standard (Poinot & Rumsey, 2018) to capture complex geometry and time evolution: for example, CGNS supports multi-block topologies and evolving meshes, with a data model that separates abstract topology (element families, etc.) from concrete mesh coordinates. On top of CGNS, PLAID layers a lightweight organizational structure.

By promoting a common standard, PLAID makes physics data interoperable across projects. It has already been used to package and publish multiple datasets covering structural mechanics and computational fluid dynamics. These PLAID-formatted datasets (hosted on Zenodo and Hugging Face) have supported ML benchmarks, democratizing access to simulation data.

## Functionality

- **Data Model and Formats:** A PLAID dataset is organized within a root folder (or archive), distinctly separating simulation data from machine learning task definitions, as illustrated in Figure 1. The dataset/ directory contains numbered sample sub-folders (sample_000...), each holding one or more .cgns files under meshes/ and a scalars.csv file. The dataset/infos.yaml file contains human-readable descriptions and metadata. The problem_definition/ folder provides machine learning context. It includes problem_infos.yaml (specifying the ML task inputs/outputs) and split.csv

39 (defining train/test splits). This design supports time evolution and multi-block/multi-
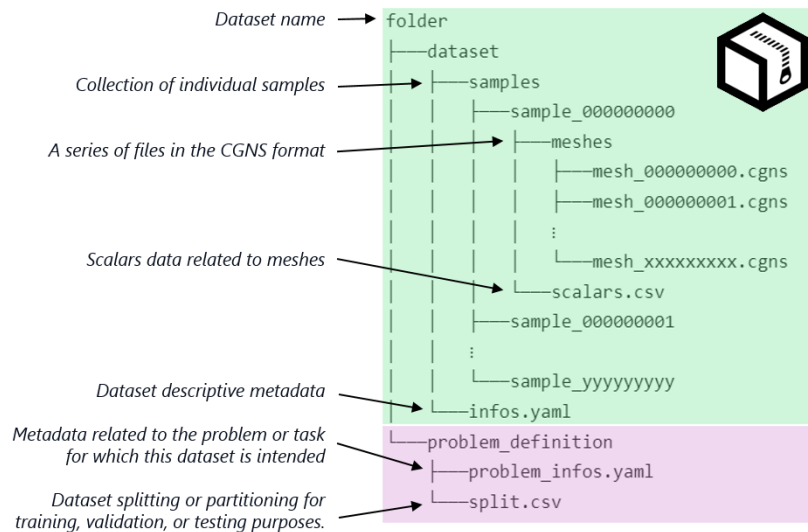40 geometry problems out of the box.



```
Dataset name ──────────→  folder
                          ├──dataset
Collection of individual samples ──→  │  ├──samples
                          │  │  ├──sample_000000000
A series of files in the CGNS format ──→  │  │  ├──meshes
                          │  │  │  │  ├──mesh_000000000.cgns
                          │  │  │  │  ├──mesh_000000001.cgns
                          │  │  │  │  ┊
                          │  │  │  │  └──mesh_xxxxxxxxx.cgns
Scalars data related to meshes ──→  │  │  └──scalars.csv
                          │  │  ├──sample_000000001
                          │  │  ┊
                          │  │  └──sample_yyyyyyyyy
Dataset descriptive metadata ──→  │  └──infos.yaml
                          └──problem_definition
Metadata related to the problem or task
for which this dataset is intended ──→  ├──problem_infos.yaml
Dataset splitting or partitioning for
training, validation, or testing purposes. ──→  └──split.csv
```

**Figure 1:** Overview of the PLAID dataset architecture.

41 • **Supported Data Types:** PLAID handles scalar, time-series and vector field data on
42 meshes, as well as sample-specific metadata. The `get_mesh(time)` method reconstructs
43 the full CGNS tree for a given timestep, with links resolved if requested (thereby returning
44 the complete mesh). Thus PLAID naturally supports mesh-based simulation outputs
45 with arbitrary element types and remeshing between time steps. Heterogeneity is allowed:
46 missing data is supported, and outputs on testing sets may be missing on purpose to
47 facilitate benchmark initiatives.

48 • **High-Level API:** The top-level `Dataset` class manages multiple `Sample` objects. Users can
49 create an empty `Dataset()` and add samples via `add_sample()`, or load an existing PLAID
50 data archive by calling `Dataset("path_to_plaid_dataset")`. The `Dataset` object sum-
51 marizes itself (e.g. printing "Dataset(3 samples, 2 scalars, 5 fields)") and provides access
52 to samples by ID. Batch operations are supported: one can `dataset.add_samples(...)`
53 to append many samples, or use the classmethods `Dataset.load_from_dir()` and
54 `load_from_file()` to load data from disk, with optional parallel workers. This high-level
55 interface abstracts away low-level I/O, letting users focus on ML pipelines.

56 • **Utilities:** PLAID includes helper modules for common tasks in data science work-
57 flows. The `plaid.utils.split` module provides a `split_dataset` function to partition
58 data into training/validation/testing subsets according to user-defined ratios. The
59 `plaid.utils.interpolation` module implements piecewise linear interpolation rou-
60 tines to resample time series fields or align datasets with different timesteps. The
61 `plaid.utils.stats` module offers an `OnlineStatistics` class to compute running
62 statistics (min, mean, variance, etc.) on arrays, which can be used to analyze dataset dis-
63 tributions. Moreover, a "Hugging Face bridge" (`plaid.bridges.huggingface_bridge`)
64 enables converting PLAID datasets to/from Hugging Face Dataset objects.

## Usage and Applications

66 PLAID is designed for AI/ML researchers and practitioners working with simulation data.
67 Various datasets, including 2D/3D fluid and structural simulations, are provided in PLAID
68 format in Hugging Face and Zenodo. Interactive benchmarks are hosted in a Hugging Face

69 community on these datasets, providing detailed instructions and PLAID commands for data
70 retrieval and manipulation, see (Casenave et al., 2025). These datasets are also used in recent
71 publications to illustrate the performance of the proposed scientific ML methods. In (Casenave
72 et al., 2024; Kabalan, Casenave, Bordeu, Ehrlacher, & Ern, 2025; Kabalan, Casenave, Bordeu,
73 & Ehrlacher, 2025), Gaussian-process regression methods with mesh morphing are applied
74 to these datasets. In (Carpintero Perez et al., 2024a, 2024b) the datasets are leveraged in
75 graph-kernel regression methods applied to fluid/solid mechanics.

76 In summary, PLAID provides a comprehensive framework for physics-based ML data. By
77 combining a unified data model, support for advanced mesh features, and helpful utilities, it
78 addresses the need for interoperable, high-fidelity simulation datasets. Future enhancements
79 involve developing general-purpose PyTorch dataloaders compatible with PLAID, along with
80 establishing standardized evaluation metrics and unified pipelines for training and inference
81 using the PLAID framework.

# References

83 Carpintero Perez, R., Da Veiga, S., Garnier, J., & Staber, B. (2024a). Gaussian process regres-
84     sion with Sliced Wasserstein Weisfeiler-Lehman graph kernels. *International Conference on*
85     *Artificial Intelligence and Statistics*, 1297–1305.

86 Carpintero Perez, R., Da Veiga, S., Garnier, J., & Staber, B. (2024b). Learning signals
87     defined on graphs with optimal transport and Gaussian process regression. *arXiv Preprint*
88     *arXiv:2410.15721*. https://doi.org/10.48550/arXiv.2410.15721

89 Casenave, F., Roynard, X., Staber, B., Akkari, N., Piat, W., Bucci, M. A., Kabalan, A., Nguyen,
90     X. M. V., Saverio, L., Perez, R. C., & others. (2025). Physics-learning AI datamodel
91     (PLAID) datasets: A collection of physics simulations for machine learning. *arXiv Preprint*
92     *arXiv:2505.02974*. https://doi.org/10.48550/arXiv.2505.02974

93 Casenave, F., Staber, B., & Roynard, X. (2024). MMGP: A Mesh Morphing Gaussian Process-
94     based machine learning method for regression of physical problems under nonparametrized
95     geometrical variability. *Advances in Neural Information Processing Systems*, *36*.

96 Kabalan, A., Casenave, F., Bordeu, F., & Ehrlacher, V. (2025). O-MMGP: Optimal Mesh
97     Morphing Gaussian Process regression for solving PDEs with non-parametric geometric
98     variations. *arXiv Preprint arXiv:2502.11632*. https://doi.org/10.48550/arXiv.2502.11632

99 Kabalan, A., Casenave, F., Bordeu, F., Ehrlacher, V., & Ern, A. (2025). Elasticity-based
100     morphing technique and application to reduced-order modeling. *Applied Mathematical*
101     *Modelling*, *141*, 115929. https://doi.org/10.1016/j.apm.2025.115929

102 Poinot, M., & Rumsey, C. L. (2018). Seven keys for practical understanding and use of CGNS.
103     *2018 AIAA Aerospace Sciences Meeting*, 1503. https://doi.org/10.2514/6.2018-1503