

# MAFw: A Modular Analysis Framework for Streamlining Data Analysis Workflows

## Statement of Need

In recent years, the complexity and volume of data generated in various scientific fields have increased exponentially. As a result, data analysis has become a critical component of scientific research, requiring efficient and flexible tools to manage and process large datasets. However, many existing data analysis frameworks are limited by their inflexibility, making it challenging to adapt to new research questions or integrate with other tools. MAFw addresses the need for a flexible and modular framework that enables data scientists to implement complex analytical tasks in a well-defined environment. Currently, data analysis workflows often require scientists to handle multiple tasks, such as data ingestion, processing, and visualization, which can be time-consuming and prone to errors. Moreover, the lack of standardization in data analysis pipelines can lead to difficulties in reproducing and sharing results. By providing a modular and customizable platform, MAFw aims to fill this gap and simplify the analysis workflow by providing a centralized location for storing and retrieving data.

## Comparison with Other Similar Libraries

MAFw is not the only framework available for data analysis, and several other libraries, such as Apache Beam, Luigi, and Nextflow, offer similar functionality. However, MAFw's conceptual design, which draws inspiration from the MARLIN framework, sets it apart from other libraries. MAFw's developers got inspired by MARLIN's modular environment, where particle physicists were developing their code in the form of shared libraries that could be loaded at run time in a plugin-like manner. By moving from C++ to Python and replacing the obsolete LCIO backend with a more flexible database-supported input/output, MAFw provides a unique combination of modularity, flexibility, and ease of use. Unlike other libraries, MAFw is designed to provide a comprehensive framework for data analysis, including data ingestion, processing, and visualization, making it an attractive choice for researchers with diverse analytical needs.

## Main Features and Functionalities

### Modularity via Subclassing and Plugin Mechanism

MAFw's core feature is its modularity, achieved through subclassing the generic `Processor` class and a plugin mechanism. This allows developers to create custom processors tailored to specific research questions, which can be easily integrated into the framework. The plugin mechanism enables seamless extension of the framework's functionality, making it an attractive choice for researchers with diverse analytical needs. By inheriting from the base `Processor` class, user-developed processors will come with some superpowers, like the ability to exchange data with the database back-end, displaying progress to the user, generating output graphs, and so on. The scientist's tasks will be limited to the implementation of the analysis code, making it easier to focus on the research question at hand. The modularity of MAFw also enables researchers to reuse and combine existing processors to create new workflows, reducing the time and effort required to develop new analytical tools.

## Cascading Processors and Steering Files

MAFw enables users to cascade different processors using a processor list, which can be executed via a simple steering file. This feature allows researchers to create complex data analysis workflows by combining multiple processors, making it an ideal tool for tasks that require multiple processing steps. The steering file provides a user-friendly interface for specifying the processor list, input parameters, and output files, making it easy to manage and reproduce complex data analysis workflows. The processor list can be easily modified and extended, allowing researchers to adapt their workflows to changing research questions or new data sources. Additionally, the steering file mechanism enables researchers to run their workflows in a consistent and reproducible manner, which is essential for scientific research.

## Database Integration and ORM

The database plays a crucial role in MAFw, providing a centralized location for storing and retrieving data. MAFw utilizes an Object-Relational Mapping (ORM) approach, specifically Peewee, to interact with the database. This allows for efficient data exchange between processors and enables seamless integration with various database management systems. The ORM also provides a high-level interface for database operations, making it easier to manage complex data relationships and perform queries. The database integration in MAFw simplifies the analysis workflow by providing a single location for storing and retrieving data, making it easier to manage and reproduce complex data analysis workflows. The use of an ORM also enables researchers to focus on their research questions, rather than the intricacies of database management, making it an attractive choice for researchers with limited database expertise.

## Real Case Application

MAFw has been successfully applied in a research study on autoradiography, published in the *Microchemical Journal* (1). The study utilized MAFw to analyze autoradiography images, demonstrating the framework's ability to streamline complex data analysis workflows. The researchers used MAFw to develop custom processors for image analysis, which were then integrated into a larger workflow using the steering file mechanism. The results of the study showcased the effectiveness of MAFw in facilitating reproducible and efficient data analysis.

## Conclusion

In conclusion, MAFw provides a flexible and customizable platform for data scientists to implement complex analytical tasks. Its modular design, plugin mechanism, and strong database integration make it an attractive choice for researchers with diverse analytical needs. By providing a comprehensive framework for data analysis, MAFw streamlines data analysis workflows, enabling researchers to focus on their research questions rather than the intricacies of data processing.

References:

(1) Krachler, M., et al. (2024). Potential of digital autoradiography for characterization of uranium materials. *Microchemical Journal*, 206, 111448.