



TECHNISCHE UNIVERSITÄT ILMENAU  
Faculty of Computer Science and Automation  
Institute for Practical Computer Science  
Department of Databases and Information Systems

Research Project

# **Definition of a Benchmark for Raster Data Processing**

Submitted By

Proteeti Prova Rawshan  
Matrikel Nr. 63752  
Winter Semester 2023/24  
Research in Computer and Systems Engineering

Supervised By:

Dr. Marcus Paradies

Ilmenau, March 1, 2024

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background and Related Work</b>	<b>2</b>
<b>3</b>	<b>Project Overview</b>	<b>4</b>
3.1	Problem Definition and Requirements . . . . .	4
3.2	ERA5 Climate Reanalysis Data . . . . .	4
3.3	Query Functions and Their Role in Benchmarking . . . . .	4
<b>4</b>	<b>Methodology</b>	<b>5</b>
4.1	Requirements . . . . .	5
4.2	Dataset . . . . .	5
4.3	Workload Generator . . . . .	5
4.4	Query formulation . . . . .	6
4.5	Benchmark execution and collected measurements . . . . .	7
<b>5</b>	<b>Implementation: Benchmark Design</b>	<b>8</b>
5.1	Benchmark Goals and Metrics . . . . .	8
5.2	Query Definition and Parameterization . . . . .	8
5.3	Workload Generation and Execution . . . . .	9
<b>6</b>	<b>Evaluation</b>	<b>11</b>
6.1	Benchmark Setup and Performance Metrics . . . . .	11
6.2	Result and Discussion . . . . .	11
<b>7</b>	<b>Future Work and Enhancements</b>	<b>13</b>
<b>8</b>	<b>Conclusion</b>	<b>14</b>
<b>Appendix</b>		
<b>A</b>	<b>API Request</b>	
<b>B</b>	<b>ERA5 Climate Reanalysis Data Exploration</b>	
<b>C</b>	<b>Query Functions Implementation</b>	
<b>D</b>	<b>Workload Generation &amp; Execution Scripts</b>	

# 1. Introduction

Benchmarking [1] plays a critical role in assessing the performance of various raster data processing systems that are widely used in applications such as environmental modeling and geospatial analysis. These systems, which include open-source tools such as xarray [12], have diverse functional scopes and performance characteristics, making it difficult to match the right tool with a particular task. The lack of a standardized benchmark for raster data [4] processing bars unbiased comparisons of these tools, thereby preventing users from making informed choices. A comprehensive benchmark that covers a range of representative queries, from basic data access to complex analytics, is urgently needed to address this gap. This benchmark will offer an objective means for comparing different systems, and enable users to select the most suitable tool for their needs [2, 13, 15].

## 2. Background and Related Work

According to numerous researchers, benchmarking is crucial in evaluating the performance of different data processing systems. Raster data processing systems are no exception, since they are used in a wide variety of applications, including, but not limited to, environmental modeling, and geospatial analysis. Currently, however, existing works focus on different data types as well as do not provide a comprehensive set of queries and detailed workload generation and execution mechanisms.

Dayarathna et al. [3] conduct an extensive survey of the literature of graph data management and processing benchmarks and identify 20 different benchmarks that have been proposed in last 15 years. It identifies several open issues in this area. A major problem is that benchmarks do not exist that exercise the system with a very high workload. In another study by [8], benchmarks spatially enabled RDF stores for geospatial data integration within Spatial Data Infrastructures (SDIs), identifying advancements in GeoSPARQL compliance and query performance. It highlights despite these advancements, if the challenge to ensure query correctness across different databases for interoperability is not resolved, there will remain a gap in achieving seamless integration and utilization of geospatial data across diverse data platforms. Meanwhile, Haynes et al. [6] introduces a benchmark used to compare how well big data platforms process raster geospatial data and to identify which systems might be best suited for different raster analysis tasks. It provides a means of identifying which platforms are particularly apt for certain prominence raster analysis tasks. It does not tackle the challenge of identifying which platforms are suitable for the large-scale processing of raster data. It does not deal with comparing these platforms in terms of their performance within "real-world" application scenarios.

Existing literature on database systems and benchmarks provides an exhaustive treatment of the area. However, none of these works addresses the unique needs of raster data processing in the context of climate data analysis. While the spatial database benchmark by Mauroux et al. [2] is enlightening about the nuances of spatial data handling, it does not specifically address the peculiarities of climate data and the challenges that raster data brings. Efficient storage and access mechanisms are pivotal in the context of array databases. The array database study by Merticariu et al. [9] emphasizes this need and highlights the gap in performance evaluation for complex data types, such as raster data. This is an important step in the array and raster database area. Optimization of distributed systems is also a topic of much research, as is the challenge of achieving high availability and consistency [15]. However, these databases typically do not address scientific data processing requirements. Lastly, [14], offers an alternative to current benchmarks that are fine tuned towards general-purpose big data systems but do not translate well to the computational and storage needs of scientific data processing.

The proposed benchmarking effort tackles the most significant gaps from prior studies, particularly in addressing handling of high workloads and efficient processing of large-scale raster data operations in real-world settings. By constructing a benchmark to evaluate how big data platforms perform during extensive raster analysis tasks, it gives a design for data management

and processing methodologies for geospatial information. Also, it tackles specific challenges of scalability and real-world applications which might be one of the critical void in current literature in geospatial data integration and analysis.

## **3. Project Overview**

### **3.1. Problem Definition and Requirements**

The project capitalizes on the preliminary development of a benchmark to undertake the performance assessment of data processing systems handling ERA5 climate reanalysis data. The justification for the need for a benchmark is the voluminous nature and complexity of climate data, which is too large and complex for tools to efficiently analyze, so analysis will be fast enough to support policy within the climate research time frame.

### **3.2. ERA5 Climate Reanalysis Data**

The European Centre for Medium-Range Weather Forecasts provides ERA5 climate reanalysis data [7], offering an updated record of comprehensive global atmospheric reanalysis from 1979. Nevertheless, its use in this research is bound by precision, resolution, and earlier-wide usage when it comes to climatic studies.

The ERA5 dataset [11] possesses complete atmospheric, land, and oceanic climate variables and, therefore, a well-articulated representation of climate patterns with time. With its high spatial and temporal resolutions, climate studies based on ERA5 will include diachronic weather forecasting to long-term climate change analysis.

### **3.3. Query Functions and Their Role in Benchmarking**

Query functions try to mimic the most basic simulation query structures common for data analysis on climate information and also other aspects. They offer simple benchmarks whereby performance evaluated by the systems should ensure effective service in executing these simple operations on big climate datasets. These include time and space aggregations, transformation of data, and derived variable calculations that are typical patterns in which climate data are analyzed.

## 4. Methodology

### 4.1. Requirements

The objective of this work is to critically evaluate the effectiveness and accuracy of different data processing systems in managing and analyzing large climate datasets. The benchmark has been crafted to stress test the systems' ability to survive and perform under heavy computational workloads common in climate research data analysis tasks, with the ultimate goal of identifying systems that can achieve both high throughput and low latency in processing and analyzing climate-related data.

### 4.2. Dataset

The ERA5-Land climate dataset was used as the foundation for this benchmarking study, owing to its extensive spatiotemporal coverage of climate variables. Specifically, this dataset includes land surface temperatures across Europe, spanning from January 2000 through the present day, detailing the climatic variations across the years. A request was formulated with which to retrieve this dataset, being sure to provide the Climate Data Store (CDS) API [10, 5] with all necessary parameters that they may return exactly the data necessary for this study's objectives. Successful acquisition of this dataset then allowed climate variables to be probed in detail, in which the xarray library was employed to facilitate the manipulation and analysis of the multi-dimensional data array.

The dataset was chosen for its comprehensive coverage of climate data across a vast temporal and spatial scale, mirroring the complexities encountered in real-world climate data analysis. This preliminary examination allowed to develop SQL queries, and also to establish a robust methodology for benchmarking, both with respect to the data's structure and granularity.

### 4.3. Workload Generator

The Workload Generator is conceived of as a tool to create a varied and representative set of analytical tasks meant to mimic the extensive and complex data analysis operations typical in climate research. This aspect of the process is essential to evaluate how well the data processing systems are able to manage diverse computational demands. The main objective of the Workload Generator is to produce tasks that simulate the key aspects of climate data analysis research, so that they provide a sound basis for evaluating the quality of the data processing systems' performance as the computational regime varies.

- **Configuration-Driven Approach:** The generator is implemented as a configuration-driven tool, which allows a set of tasks to be defined in a flexible manner, using a variety of established parameters so that the workload is able to include many different scenarios of data retrieval and transformation.
- **Representation of Tasks:** The tasks that are produced are chosen to reflect different aspects of climate data analysis, such as the retrieval of data, the performance of data aggregation, and carrying out various kinds of data transformation. Together, they are devised to put a full range of fascinating challenges in front of the systems, whose performance they are meant to probe.

## 4.4. Query formulation

A further exploration of the ERA5-Land dataset by formulating and executing a series of simple queries aimed at extracting specific climate data insights (with visualizations) for better understanding. This involved applying various data retrieval techniques to pinpoint trends, anomalies, and patterns in the temperature data across different time periods and geographical locations within Europe. Then, it was taken up to a higher level with more complicated stuff, such as figuring out long-term climate patterns and oddities. This step was pivotal in establishing a framework that tests both the efficiency and scalability of data management systems in handling raster data.

The queries detailed in Table 4.1 offer the means to tame the climate dataset in order to manipulate and explore it. Each query is intended to respond to an actual data operation commonly done in climate research.

Function	Query	Purpose
Average Temperature Over Time	TemporalRangeAggregation	Reflects the system's ability to handle time-series data.
Geographical Temperature Analysis	SpatialSubsetExtraction	Tests the system's spatial data handling efficiency
Adjusted Climate Data	TemporalDataModification	Evaluates the system's capacity for data manipulation
Derived Climate Metrics	AddDerivedVariable	Tests the system's computational functionality
Seasonal Weather Patterns	SeasonalAggregation	Assesses the system's ability to aggregate data.

Table 4.1.: Descriptive summary of the query functions implemented



## 4.5. Benchmark execution and collected measurements

In this phase, the aim is to convert the conceptual framework to an empirical method through which real-world climate data analysis scenarios can be executed and managed by each system.

The benchmark execution process itself is intended to measure how well each system runs through the workload. That is to say that the focus of the benchmark is on measuring execution time, CPU usage and memory consumption. These performance metrics are chosen because of their direct relevance in measuring the critical aspects of system performance such as computational speed, resource efficiency and the ability to scale under increased loads. By capturing and then analyzing these metrics, the study seeks to provide objective insight into the operational strengths and limitations of each system in the context of large-scale climate data set analysis.

# 5. Implementation: Benchmark Design

## 5.1. Benchmark Goals and Metrics

The benchmark was designed to assess the performance and efficiency of processing climate data. It evaluates the processing of ERA5-Land temperature data across Europe from 2000 to 2022.

The primary metrics for evaluation are execution time, CPU usage, and memory consumption, which will allow us to identify the most efficient data handling techniques for large-scale climate datasets. Such an assessment is critical to enhance climate data analysis workflows, which are essential to climate research and decision-making.

## 5.2. Query Definition and Parameterization

Query formulation is a critical step in translating scientific questions into code that is actually run. It involves defining a set of queries that collectively cover a wide spectrum of typical data analysis tasks found in climate datasets. In the *query\_fucntions.py*, these functions were used:

### 1. **temporal\_range\_aggregation**

- Purpose: This function aggregates (mean/sum) the 't2m' (2-meter temperature) variable over a given temporal range.
- Query Type: Analyzes how temperature changes over a given period, which is useful for understanding climatic trends/ variations within that timeframe.

### 2. **spatial\_subset\_extraction**

- Purpose: This function extracts the subset of 't2m' variable data within the given latitudinal and longitudinal ranges.
- Query Type: Essential to narrow down on a specific geographical area and hence observe climatic patterns or conditions localized to that area.

### 3. **temporal\_data\_modification**

- Purpose: This function modifies the 't2m' variable data over the specified month and year for all or a selection of latitude and/or longitude. It can increase/decrease the data by a specified multiple as well. It also checks for empty selections to ensure that the modifications are not applied where there is no data.

- Query Type: Could be used to see the impact of any condition/intervention on the temperature data, which will assist in understanding potential climate change scenarios/ conditions where the actual data values could be missing or any other anomalies in the dataset.

#### 4. `add_derived_variable`

- Purpose: This function adds a derived variable to the dataset. The new variable can be the existing variable 't2m' with any of the following operations performed with the given factor: multiply, add, subtract, divide.
- Query Type: Useful in generating new insights/ variables from the existing data and is critical for more complex climate analysis or model inputs.

#### 5. `seasonal_aggregation`

- Purpose: This function provides a helper for the user to aggregate the 't2m' data over the specified season (winter, spring, summer, autumn) for the given year. The helper internally uses a predefined month range for each season (e.g. winter is Dec, Jan, Feb).
- Query Type: Essential for looking at the conditions/ variations over a year across all the seasons and thus, the indication of each season's impact on the temperature.

### 5.3. Workload Generation and Execution

**Benchmark Configuration:** Defined in the *benchmark\_config.json*, is a configuration file that specifies parameters for benchmarks, including which functions to run, their parameters, and perhaps in what order. Here, the dataset was specified (GRIB file). Also, in an array of objects, the queries were defined, where each object has a type corresponding to a specific operation and a weight, which indicates the relative importance or frequency of each query type in the benchmark.

For example: `"type": "TemporalRangeAggregation", "weight": 30` - meaning this query will aggregate data over a specified time range, where it's 30% of the total queries performed.

**Workload Generation:** Workload generation is a key aspect of benchmarking, where the goal is to create a diverse and representative set of queries to evaluate the performance of data processing systems. This is done through the *workload\_generation.py* script, where a sequence of tasks are programmatically defined to mimic realistic data analysis scenarios on climate datasets such as the ERA5-Land temperature data, and the workload is generated in (*workload.csv*).

The script follows a configuration-driven approach where the types of queries, configuration parameters, and under which conditions they are run are specified in (*benchmark\_config.json*). Over-all, the workload consists of a balanced mix of data retrieval, aggregation, data transformation tasks which push on different aspect of the system's performance.

**Workload Execution:** The workload execution, outlined in *workload\_execution.py*, is responsible for running the generated workload against the data processing system. The script reads in the list of tasks created by the workload generator and executes each query sequentially. At the

conclusion of the query, key performance metrics (e.g., execution time, CPU, memory) are captured. These metrics are essential input for understanding the aforementioned system efficiency and scalability under differing computational loads.

To ensure accurate and reliable benchmark results, the execution mechanism is robust to errors and includes mechanisms to handle errors, logging and the ability to retry failed tasks. This functionality is particularly important for handling the complexity and potential unpredictability of processing large-scale climate datasets.

The exact order and relationship between these components would be like this: `benchmark_config` would be read by `workload_generation` to create a series of tasks, which are then executed by `workload_execution` using the functions defined in `query_functions.py`.

## 6. Evaluation

### 6.1. Benchmark Setup and Performance Metrics

The benchmark was designed to evaluate the performance and efficiency of various query functions, each aimed at testing different aspects of data handling and analysis capabilities.

- **Execution Time:** This metric denotes the total time for a query to complete its execution - from initiation to completion. This is important because understanding spiky or long latencies times for data processing can be important, since a rapid response helps build data analytics around issues and opportunities that are current, topical, and timely.
- **CPU Usage:** This is a measure of how intensively a query is exercising the computational elements of query execution. It expresses the percentage of the cycles the CPU is busy during the life cycle of the query.

But getting this exactly quantified is actually hard, as a modern operating system is likely to have thousands of tasks switching in and out, and in and out, and this will also make the CPU look as if it's executing other processes during the time in question, changing readings. So, CPU usage only gives a rough approximate idea of the load in terms of computation, but does not often wholly reflect the resource consumption of the query itself.

- **Memory Usage:** This is a measure of the extra memory (beyond the database memory) that is going to be utilized as this query executes. Memory usage is a critical measure in capacity planning, and a rich source of data for understanding the scalability of a new and growing data processing operation. Memory utilization that is excessive will lead to thrashing and will degrade all the performance of the system, not just the performance of the query execution.

### 6.2. Result and Discussion

Examination of 't2m' climate data operations reveals significant variability in computational demands. It is observed in [6.1](#) that, execution times varied from approximately  $1.46\text{E}+15$  to  $2.21\text{E}+16$  secs. Changes in CPU usage ranged up to 130%. Memory usage fluctuations varied from a decrease of 390625 MB to an increase of 3125 MB. Operations are associated with a 'query\_index' that categorizes different types. The operations types range from SpatialSubsetExtraction and AddDerivedVariable to various aggregations. Examining the data it appears that operations such as temporal modification and aggregation are particularly computationally demanding. The wide variability in resource usage clearly illustrates the complexity of performing analysis on climate data, with some types of operations consuming more resources (either due

## 6. Evaluation

to granularity of the data or the introduction of new variables).

APC result	123 execution_time_sec	123 cpu_usage_diff_percent	123 memory_usage_diff_mb	123 query_index	APC query_type
1 <xarray.DataArray 't2m' (time: 264, latitude: 0, longitude: 375)> [0 values with dtype=float32] Coordinates: number int32 ... * time (time) datetime64[ns] 2000-01-01 2000-02-01 ... 2021-12-01	2,211E+16	126	3125	0	SpatialSubsetExtraction
2 <xarray.DataArray 't2m' (time: 264, latitude: 0, longitude: 0)> [0 values with dtype=float32] Coordinates: number int32 ... * time (time) datetime64[ns] 2000-01-01 2000-02-01 ... 2021-12-01	5,575E+15	0	4296875	1	SpatialSubsetExtraction
3 <xarray.Dataset> Dimensions: (time: 264, latitude: 501, longitude: 701) Coordinates: number int32 ... * time (time) datetime64[ns] 2000-01-01 2000-02-01 ... 2021-12-01	1,66438E+16	47	42888671875	2	AddDerivedVariable
4 <xarray.DataArray 't2m' (latitude: 501, longitude: 701)> array([[259.49878, 259.58472, 259.66675, ..., nan, nan, nan], [259.83667, 259.8933, 259.948, ..., nan, nan, nan], [258.28662, 258.3706, 258.4546, ..., nan, nan, nan], [257.97607, 258.0874, 258.19482, ..., nan, nan, nan], [257.66667, 257.77889, 257.89111, ..., nan, nan, nan], [257.38889, 257.50111, 257.61333, ..., nan, nan, nan], [257.11111, 257.22333, 257.33556, ..., nan, nan, nan], [256.83333, 256.94556, 257.05778, ..., nan, nan, nan], [256.55556, 256.66778, 256.78, ..., nan, nan, nan], [256.27778, 256.39, 256.50222, ..., nan, nan, nan], [256.0, 256.11222, 256.22444, ..., nan, nan, nan], [255.72222, 255.83444, 255.94667, ..., nan, nan, nan], [255.44444, 255.55667, 255.66889, ..., nan, nan, nan], [255.16667, 255.27889, 255.39111, ..., nan, nan, nan], [254.88889, 254.99999, 255.11111, ..., nan, nan, nan], [254.61111, 254.72222, 254.83333, ..., nan, nan, nan], [254.33333, 254.44444, 254.55556, ..., nan, nan, nan], [254.05556, 254.16667, 254.27778, ..., nan, nan, nan], [253.77778, 253.88889, 254.0, ..., nan, nan, nan], [253.5, 253.61111, 253.72222, ..., nan, nan, nan], [253.22222, 253.33333, 253.44444, ..., nan, nan, nan], [252.94444, 253.05556, 253.16667, ..., nan, nan, nan], [252.66667, 252.77778, 252.88889, ..., nan, nan, nan], [252.38889, 252.5, 252.61111, ..., nan, nan, nan], [252.11111, 252.22222, 252.33333, ..., nan, nan, nan], [251.83333, 251.94444, 252.05556, ..., nan, nan, nan], [251.55556, 251.66667, 251.77778, ..., nan, nan, nan], [251.27778, 251.38889, 251.5, ..., nan, nan, nan], [251.0, 251.11111, 251.22222, ..., nan, nan, nan], [250.72222, 250.83333, 250.94444, ..., nan, nan, nan], [250.44444, 250.55556, 250.66667, ..., nan, nan, nan], [250.16667, 250.27778, 250.38889, ..., nan, nan, nan], [249.88889, 249.99999, 250.11111, ..., nan, nan, nan], [249.61111, 249.72222, 249.83333, ..., nan, nan, nan], [249.33333, 249.44444, 249.55556, ..., nan, nan, nan], [249.05556, 249.16667, 249.27778, ..., nan, nan, nan], [248.77778, 248.88889, 249.0, ..., nan, nan, nan], [248.5, 248.61111, 248.72222, ..., nan, nan, nan], [248.22222, 248.33333, 248.44444, ..., nan, nan, nan], [247.94444, 248.05556, 248.16667, ..., nan, nan, nan], [247.66667, 247.77778, 247.88889, ..., nan, nan, nan], [247.38889, 247.5, 247.61111, ..., nan, nan, nan], [247.11111, 247.22222, 247.33333, ..., nan, nan, nan], [246.83333, 246.94444, 247.05556, ..., nan, nan, nan], [246.55556, 246.66667, 246.77778, ..., nan, nan, nan], [246.27778, 246.38889, 246.5, ..., nan, nan, nan], [246.0, 246.11111, 246.22222, ..., nan, nan, nan], [245.72222, 245.83333, 245.94444, ..., nan, nan, nan], [245.44444, 245.55556, 245.66667, ..., nan, nan, nan], [245.16667, 245.27778, 245.38889, ..., nan, nan, nan], [244.88889, 244.99999, 245.11111, ..., nan, nan, nan], [244.61111, 244.72222, 244.83333, ..., nan, nan, nan], [244.33333, 244.44444, 244.55556, ..., nan, nan, nan], [244.05556, 244.16667, 244.27778, ..., nan, nan, nan], [243.77778, 243.88889, 244.0, ..., nan, nan, nan], [243.5, 243.61111, 243.72222, ..., nan, nan, nan], [243.22222, 243.33333, 243.44444, ..., nan, nan, nan], [242.94444, 243.05556, 243.16667, ..., nan, nan, nan], [242.66667, 242.77778, 242.88889, ..., nan, nan, nan], [242.38889, 242.5, 242.61111, ..., nan, nan, nan], [242.11111, 242.22222, 242.33333, ..., nan, nan, nan], [241.83333, 241.94444, 242.05556, ..., nan, nan, nan], [241.55556, 241.66667, 241.77778, ..., nan, nan, nan], [241.27778, 241.38889, 241.5, ..., nan, nan, nan], [241.0, 241.11111, 241.22222, ..., nan, nan, nan], [240.72222, 240.83333, 240.94444, ..., nan, nan, nan], [240.44444, 240.55556, 240.66667, ..., nan, nan, nan], [240.16667, 240.27778, 240.38889, ..., nan, nan, nan], [239.88889, 239.99999, 240.11111, ..., nan, nan, nan], [239.61111, 239.72222, 239.83333, ..., nan, nan, nan], [239.33333, 239.44444, 239.55556, ..., nan, nan, nan], [239.05556, 239.16667, 239.27778, ..., nan, nan, nan], [238.77778, 238.88889, 239.0, ..., nan, nan, nan], [238.5, 238.61111, 238.72222, ..., nan, nan, nan], [238.22222, 238.33333, 238.44444, ..., nan, nan, nan], [237.94444, 238.05556, 238.16667, ..., nan, nan, nan], [237.66667, 237.77778, 237.88889, ..., nan, nan, nan], [237.38889, 237.5, 237.61111, ..., nan, nan, nan], [237.11111, 237.22222, 237.33333, ..., nan, nan, nan], [236.83333, 236.94444, 237.05556, ..., nan, nan, nan], [236.55556, 236.66667, 236.77778, ..., nan, nan, nan], [236.27778, 236.38889, 236.5, ..., nan, nan, nan], [236.0, 236.11111, 236.22222, ..., nan, nan, nan], [235.72222, 235.83333, 235.94444, ..., nan, nan, nan], [235.44444, 235.55556, 235.66667, ..., nan, nan, nan], [235.16667, 235.27778, 235.38889, ..., nan, nan, nan], [234.88889, 234.99999, 235.11111, ..., nan, nan, nan], [234.61111, 234.72222, 234.83333, ..., nan, nan, nan], [234.33333, 234.44444, 234.55556, ..., nan, nan, nan], [234.05556, 234.16667, 234.27778, ..., nan, nan, nan], [233.77778, 233.88889, 234.0, ..., nan, nan, nan], [233.5, 233.61111, 233.72222, ..., nan, nan, nan], [233.22222, 233.33333, 233.44444, ..., nan, nan, nan], [232.94444, 233.05556, 233.16667, ..., nan, nan, nan], [232.66667, 232.77778, 232.88889, ..., nan, nan, nan], [232.38889, 232.5, 232.61111, ..., nan, nan, nan], [232.11111, 232.22222, 232.33333, ..., nan, nan, nan], [231.83333, 231.94444, 232.05556, ..., nan, nan, nan], [231.55556, 231.66667, 231.77778, ..., nan, nan, nan], [231.27778, 231.38889, 231.5, ..., nan, nan, nan], [231.0, 231.11111, 231.22222, ..., nan, nan, nan], [230.72222, 230.83333, 230.94444, ..., nan, nan, nan], [230.44444, 230.55556, 230.66667, ..., nan, nan, nan], [230.16667, 230.27778, 230.38889, ..., nan, nan, nan], [229.88889, 229.99999, 230.11111, ..., nan, nan, nan], [229.61111, 229.72222, 229.83333, ..., nan, nan, nan], [229.33333, 229.44444, 229.55556, ..., nan, nan, nan], [229.05556, 229.16667, 229.27778, ..., nan, nan, nan], [228.77778, 228.88889, 229.0, ..., nan, nan, nan], [228.5, 228.61111, 228.72222, ..., nan, nan, nan], [228.22222, 228.33333, 228.44444, ..., nan, nan, nan], [227.94444, 228.05556, 228.16667, ..., nan, nan, nan], [227.66667, 227.77778, 227.88889, ..., nan, nan, nan], [227.38889, 227.5, 227.61111, ..., nan, nan, nan], [227.11111, 227.22222, 227.33333, ..., nan, nan, nan], [226.83333, 226.94444, 227.05556, ..., nan, nan, nan], [226.55556, 226.66667, 226.77778, ..., nan, nan, nan], [226.27778, 226.38889, 226.5, ..., nan, nan, nan], [226.0, 226.11111, 226.22222, ..., nan, nan, nan], [225.72222, 225.83333, 225.94444, ..., nan, nan, nan], [225.44444, 225.55556, 225.66667, ..., nan, nan, nan], [225.16667, 225.27778, 225.38889, ..., nan, nan, nan], [224.88889, 224.99999, 225.11111, ..., nan, nan, nan], [224.61111, 224.72222, 224.83333, ..., nan, nan, nan], [224.33333, 224.44444, 224.55556, ..., nan, nan, nan], [224.05556, 224.16667, 224.27778, ..., nan, nan, nan], [223.77778, 223.88889, 224.0, ..., nan, nan, nan], [223.5, 223.61111, 223.72222, ..., nan, nan, nan], [223.22222, 223.33333, 223.44444, ..., nan, nan, nan], [222.94444, 223.05556, 223.16667, ..., nan, nan, nan], [222.66667, 222.77778, 222.88889, ..., nan, nan, nan], [222.38889, 222.5, 222.61111, ..., nan, nan, nan], [222.11111, 222.22222, 222.33333, ..., nan, nan, nan], [221.83333, 221.94444, 222.05556, ..., nan, nan, nan], [221.55556, 221.66667, 221.77778, ..., nan, nan, nan], [221.27778, 221.38889, 221.5, ..., nan, nan, nan], [221.0, 221.11111, 221.22222, ..., nan, nan, nan], [220.72222, 220.83333, 220.94444, ..., nan, nan, nan], [220.44444, 220.55556, 220.66667, ..., nan, nan, nan], [220.16667, 220.27778, 220.38889, ..., nan, nan, nan], [219.88889, 219.99999, 220.11111, ..., nan, nan, nan], [219.61111, 219.72222, 219.83333, ..., nan, nan, nan], [219.33333, 219.44444, 219.55556, ..., nan, nan, nan], [219.05556, 219.16667, 219.27778, ..., nan, nan, nan], [218.77778, 218.88889, 219.0, ..., nan, nan, nan], [218.5, 218.61111, 218.72222, ..., nan, nan, nan], [218.22222, 218.33333, 218.44444, ..., nan, nan, nan], [217.94444, 218.05556, 218.16667, ..., nan, nan, nan], [217.66667, 217.77778, 217.88889, ..., nan, nan, nan], [217.38889, 217.5, 217.61111, ..., nan, nan, nan], [217.11111, 217.22222, 217.33333, ..., nan, nan, nan], [216.83333, 216.94444, 217.05556, ..., nan, nan, nan], [216.55556, 216.66667, 216.77778, ..., nan, nan, nan], [216.27778, 216.38889, 216.5, ..., nan, nan, nan], [216.0, 216.11111, 216.22222, ..., nan, nan, nan], [215.72222, 215.83333, 215.94444, ..., nan, nan, nan], [215.44444, 215.55556, 215.66667, ..., nan, nan, nan], [215.16667, 215.27778, 215.38889, ..., nan, nan, nan], [214.88889, 214.99999, 215.11111, ..., nan, nan, nan], [214.61111, 214.72222, 214.83333, ..., nan, nan, nan], [214.33333, 214.44444, 214.55556, ..., nan, nan, nan], [214.05556, 214.16667, 214.27778, ..., nan, nan, nan], [213.77778, 213.88889, 214.0, ..., nan, nan, nan], [213.5, 213.61111, 213.72222, ..., nan, nan, nan], [213.22222, 213.33333, 213.44444, ..., nan, nan, nan], [212.94444, 213.05556, 213.16667, ..., nan, nan, nan], [212.66667, 212.77778, 212.88889, ..., nan, nan, nan], [212.38889, 212.5, 212.61111, ..., nan, nan, nan], [212.11111, 212.22222, 212.33333, ..., nan, nan, nan], [211.83333, 211.94444, 212.05556, ..., nan, nan, nan], [211.55556, 211.66667, 211.77778, ..., nan, nan, nan], [211.27778, 211.38889, 211.5, ..., nan, nan, nan], [211.0, 211.11111, 211.22222, ..., nan, nan, nan], [210.72222, 210.83333, 210.94444, ..., nan, nan, nan], [210.44444, 210.55556, 210.66667, ..., nan, nan, nan], [210.16667, 210.27778, 210.38889, ..., nan, nan, nan], [209.88889, 209.99999, 210.11111, ..., nan, nan, nan], [209.61111, 209.72222, 209.83333, ..., nan, nan, nan], [209.33333, 209.44444, 209.55556, ..., nan, nan, nan], [209.05556, 209.16667, 209.27778, ..., nan, nan, nan], [208.77778, 208.88889, 209.0, ..., nan, nan, nan], [208.5, 208.61111, 208.72222, ..., nan, nan, nan], [208.22222, 208.33333, 208.44444, ..., nan, nan, nan], [207.94444, 208.05556, 208.16667, ..., nan, nan, nan], [207.66667, 207.77778, 207.88889, ..., nan, nan, nan], [207.38889, 207.5, 207.61111, ..., nan, nan, nan], [207.11111, 207.22222, 207.33333, ..., nan, nan, nan], [206.83333, 206.94444, 207.05556, ..., nan, nan, nan], [206.55556, 206.66667, 206.77778, ..., nan, nan, nan], [206.27778, 206.38889, 206.5, ..., nan, nan, nan], [206.0, 206.11111, 206.22222, ..., nan, nan, nan], [205.72222, 205.83333, 205.94444, ..., nan, nan, nan], [205.44444, 205.55556, 205.66667, ..., nan, nan, nan], [205.16667, 205.27778, 205.38889, ..., nan, nan, nan], [204.88889, 204.99999, 205.11111, ..., nan, nan, nan], [204.61111, 204.72222, 204.83333, ..., nan, nan, nan], [204.33333, 204.44444, 204.55556, ..., nan, nan, nan], [204.05556, 204.16667, 204.27778, ..., nan, nan, nan], [203.77778, 203.88889, 204.0, ..., nan, nan, nan], [203.5, 203.61111, 203.72222, ..., nan, nan, nan], [203.22222, 203.33333, 203.44444, ..., nan, nan, nan], [202.94444, 203.05556, 203.16667, ..., nan, nan, nan], [202.66667, 202.77778, 202.88889, ..., nan, nan, nan], [202.38889, 202.5, 202.61111, ..., nan, nan, nan], [202.11111, 202.22222, 202.33333, ..., nan, nan, nan], [201.83333, 201.94444, 202.05556, ..., nan, nan, nan], [201.55556, 201.66667, 201.77778, ..., nan, nan, nan], [201.27778, 201.38889, 201.5, ..., nan, nan, nan], [201.0, 201.11111, 201.22222, ..., nan, nan, nan], [200.72222, 200.83333, 200.94444, ..., nan, nan, nan], [200.44444, 200.55556, 200.66667, ..., nan, nan, nan], [200.16667, 200.27778, 200.38889, ..., nan, nan, nan], [199.88889, 199.99999, 200.11111, ..., nan, nan, nan], [199.61111, 199.72222, 199.83333, ..., nan, nan, nan], [199.33333, 199.44444, 199.55556, ..., nan, nan, nan], [199.05556, 199.16667, 199.27778, ..., nan, nan, nan], [198.77778, 198.88889, 199.0, ..., nan, nan, nan], [198.5, 198.61111, 198.72222, ..., nan, nan, nan], [198.22222, 198.33333, 198.44444, ..., nan, nan, nan], [197.94444, 198.05556, 198.16667, ..., nan, nan, nan], [197.66667, 197.77778, 197.88889, ..., nan, nan, nan], [197.38889, 197.5, 197.61111, ..., nan, nan, nan], [197.11111, 197.22222, 197.33333, ..., nan, nan, nan], [196.83333, 196.94444, 197.05556, ..., nan, nan, nan], [196.55556, 196.66667, 196.77778, ..., nan, nan, nan], [196.27778, 196.38889, 196.5, ..., nan, nan, nan], [196.0, 196.11111, 196.22222, ..., nan, nan, nan], [195.72222, 195.83333, 195.94444, ..., nan, nan, nan], [195.44444, 195.55556, 195.66667, ..., nan, nan, nan], [195.16667, 195.27778, 195.38889, ..., nan, nan, nan], [194.88889, 194.99999, 195.11111, ..., nan, nan, nan], [194.61111, 194.72222, 194.83333, ..., nan, nan, nan], [194.33333, 194.44444, 194.55556, ..., nan, nan, nan], [194.05556, 194.16667, 194.27778, ..., nan, nan, nan], [193.77778, 193.88889, 194.0, ..., nan, nan, nan], [193.5, 193.61111, 193.72222, ..., nan, nan, nan], [193.22222, 193.33333, 193.44444, ..., nan, nan, nan], [192.94444, 193.05556, 193.16667, ..., nan, nan, nan], [192.66667, 192.77778, 192.88889, ..., nan, nan, nan], [192.38889, 192.5, 192.61111, ..., nan, nan, nan], [192.11111, 192.22222, 192.33333, ..., nan, nan, nan], [191.83333, 191.94444, 192.05556, ..., nan, nan, nan], [191.55556, 191.66667, 191.77778, ..., nan, nan, nan], [191.27778, 191.38889,					

## 7. Future Work and Enhancements

The exploration of ERA5 climate data with direct query execution and real-time visualization is a stepping stone to climate analytics. The possible improvements & future work might aim to build upon it in the following ways:

**Tool Evaluation:** Comparing the existing query execution and data analysis techniques with the likes of Dask or PySpark or similar tools for benchmarking and identifying improvements with respect to execution time, scalability and resource usage.

**Benchmarking Enhancements:** Developing a comprehensive benchmark setup with a variety of query types, larger datasets, and different computational environments to identify strengths and optimization opportunities.

**Dataset Expansion:** The scope of the work can be broadened to consider other variables from ERA5 as well as other climate datasets to enable the analysis of broader climate phenomena. To support these, new query functions will be needed that can extend the data volume processed and adapt the implemented data model to accommodate multi-dimensional datasets.

**Tool Comparison:** Systematic comparison with more climate data analysis tools to understand the salient features of the implemented tool in terms of speed, accuracy and scalability, for example.

**User Interface Development:** Development of a user-friendly interface that facilitates interactive queries, real-time visualization and easy data interpretation from complex queries and large datasets and hence, bridges the gap between the current tool and the climate science community.

## 8. Conclusion

In this report, a benchmarking framework was introduced for the analysis of ERA5 climate reanalysis data, highlighting the development, and particularly, the insights gathered from the exploration phase as well. The work has not only advanced the efficiency of climate data analysis, but has also laid the groundwork for furthering contributions to climate science research. The successes seen while working with these complex climate data should pave the way for a myriad of future work, ranging from comparisons against other tools to broadening the tool's applicability and user accessibility. It might also produce some make meaningful contributions to the understanding and addressing the challenges of the changing climate, with a continued focus on innovative thinking and broadening the analytical reach.



# A. API Request

In order to obtain the climate data, an API request was formulated and executed using the cdsapi library to the Climate Data Store (CDS), targeting the monthly averaged reanalysis data for the 2m\_temperature variable for a specific geographical bounding box over a particular time range (Europe in this case, 2000-2021).

```
1 import cdsapi
2
3 c = cdsapi.Client(key=cdsapi_key)
4
5 variables_list = ['2m_temperature']
6 for variable in variables_list:
7     c.retrieve(
8         'reanalysis-era5-land-monthly-means',
9         {
10             'format': 'grib',
11             'product_type': 'monthly_averaged_reanalysis',
12             'variable': variable,
13             'year': [str(year) for year in range(2000, 2022)],
14             'month': ['01', '02', '03', '04', '05', '06', '07', '08', '09',
15                     '10', '11', '12'],
16             'time': '00:00',
17             'area': [25, -25, 75, 45], # Europe bounding box
18         },
19         'ERA5-Land_2m_temperature_2000-2022_europe_2.grib'
```

The 'retrieve' method of the CDS API client accepts the dataset name, the request parameters, and the destination file path. Each parameter in the request serves a specific purpose:

**format** The format of the downloaded file.

**product\_type** The type of reanalysis product to download.

**variable** The specific climate variable of interest.

**year and month** The temporal range for the data retrieval.

**time** The time of day for which the data is required.

**area** The geographical bounding box specifying the spatial extent of the data.

## B. ERA5 Climate Reanalysis Data Exploration

In this study, the ERA5-Land dataset was interacted with using and explored with the xarray library. Xarray is an incredibly powerful Python package for using/manipulating multi-dimensional arrays especially when working with climate and meteorological data. The ERA5-Land dataset is made up of so many climate variables with high spatial and temporal resolution. The dataset comprises a lot of climate variables including '2m\_temperature' which is the temperature at 2 meters above the surface.

Key aspects of the ERA5-Land dataset include:

- Spatial resolution: Contains detailed latitude and longitude data. Consequently, allows for detailed geographical based analysis.
- Temporal resolution: Monthly averages of all variables from 2000, up until recently (up to the 22nd of February 2022).
- Variables: This study primarily focuses on the '2m\_temperature' variable.

The code below shows how the dataset was explored using xarray:

```
1 import xarray as xr
2
3 # Open the GRIB file using xarray with cfgrib engine for GRIB files
4 grib_file_path = "path_to_grib_file.grib"
5 ds_2 = xr.open_dataset(grib_file_path, engine='cfgrib')
6 print(ds_2)
```

There is also 'Attributes' section showing trustworthy metadata on the source of these data, the standard to which the data adhere, and the dataset's transformation history, making it a reliable, standards-compliant dataset that is suitable for studying climate and climate change.

```
1 <xarray.Dataset>
2 Dimensions:      (time: 264, latitude: 501, longitude: 701)
3 Coordinates:
4   * time          (time) datetime64[ns] 2000-01-01 2000-02-01 ... 2021-12-01
5   * latitude      (latitude) float64 75.0 74.9 74.8 74.7 ... 25.3 25.2 25.1
6   * longitude     (longitude) float64 -25.0 -24.9 -24.8 -24.7 ... 44.8 44.9
7 Data variables:
8   t2m            (time, latitude, longitude) float32 ...
9 Attributes:
10  GRIB_edition:      1
11  GRIB_centre:       ecmf
12  GRIB_centreDescription: European Centre for Medium-Range Weather
    Forecasts
```

## B. ERA5 Climate Reanalysis Data Exploration

---

```
13 Conventions: CF-1.7
14 institution: European Centre for Medium-Range Weather
    Forecasts
15 history: 2024-02-07T22:56 GRIB to CDM+CF via cfgrib
    -0.9.1...
```

As observed above, the dataset has 264 time points, covering 22 years including 501 latitudinal and 701 longitudinal points across Europe. The implementation focuses on variables like 2m temperature (t2m), which were vital for the analysis, giving detailed time and location information crucial for the comparisons.

## C. Query Functions Implementation

This section describes the implementation of core query functions within the project which are designed to analyze ERA5 climate reanalysis data. Each of these functions are tailored to address a specific analysis pattern which manipulates the queried data based on the xarray objects:

```
1     # Aggregate 't2m' over a specified temporal range.
2 def temporal_range_aggregation(ds, start_date, end_date, variable="t2m",
   operation="mean"):
3     # Selects the data slice based on the time range and performs the
   specified aggregation operation
4     # Returns: xarray.DataArray with the aggregation result
5
6 # Extract a spatial subset for 't2m'
7 def spatial_subset_extraction(ds, lat_range, lon_range, variable="t2m"):
8     # Selects a subset of the data based on the specified latitude and
   longitude range
9     # Returns: xarray.DataArray containing the spatial subset
10
11 # Modify 't2m' data for a specific month and year
12 def temporal_data_modification(ds, year, month, modification="increase",
   value=1.0, variable="t2m"):
13     # Applies a specified modification to 't2m' data for the given month
   and year, if data is available
14     # Returns: xarray.Dataset with the modified 't2m' data
15
16 # Add a derived variable based on 't2m'
17 def add_derived_variable(ds, new_variable_name, base_variable='t2m',
   operation='multiply', factor=1.0):
18     # Creates a new variable in the dataset based on the specified
   operation and factor
19     # Returns: xarray.Dataset with the new derived variable added
20
21 # Aggregate 't2m' over a specified season
22 def seasonal_aggregation(ds, year, season, variable="t2m", operation="mean"
   ):
23     # Performs an aggregation of 't2m' data over the specified season for
   the given year
24     # Returns: xarray.DataArray with the seasonal aggregation result
```

Each of these functions are critical for building a climate pattern analysis that is flexible and allows a user to build functionality that does several operations and allows to build up more and more complex climate analyses.

## D. Workload Generation & Execution Scripts

The workload generation script (`workload_generation.py`) creates a set of queries dynamically, with different types and parameters. It helps to simulate many different types of data processing tasks. Some of the key components of this script are:

- **Parameter Generation:** Functions to randomly generate different types of parameters for different types of queries to make sure that we test a variety of cases.
- **Workload File Creation:** The script creates a CSV file that contains the generated queries which can be run by the workload execution script.

```
1 # Example of parameter generation for a Temporal Range Aggregation query
2 def generate_parameters(query_type):
3     if query_type == "TemporalRangeAggregation":
4         # Generates start and end dates, variable, and operation
5
6 # Main function to generate the workload
7 def generate_workload(config_path="benchmark_config.json", workload_file="
    workload.csv"):
8     # Loads configuration, generates queries based on weights, and writes
    them to a CSV file
```

The workload execution script (`workload_execution.py`) reads the generated queries and executes them against the dataset, capturing and logging performance metrics such as execution time, CPU usage, and memory usage. It demonstrates integration with the `xarray` library for data manipulation and `psutil` for monitoring system resources.

- **Query Execution:** Implements logic to execute various types of queries, including temporal and spatial data manipulations.
- **System Usage Monitoring:** Captures system metrics before and after query execution to measure impact.
- **Results Aggregation:** Collects and stores the execution results for analysis in csv.

```
1 # Function to execute a single query and measure system usage
2 def execute_query(ds, query):
3     cpu_before, mem_before = get_system_usage()
4     # Executes the query based on its type (e.g., TemporalRangeAggregation)
5
6 # Main function to run the benchmark
7 def run_benchmark(workload_file, config):
8     # Opens the dataset, iterates through the workload, and executes
    queries
```

# Bibliography

- [1] Dina Bitton, David J DeWitt, and Carolyn Turbyfill. Benchmarking database systems-a systematic approach. Technical report, University of Wisconsin-Madison Department of Computer Sciences, 1983.
- [2] Philippe Cudre-Mauroux, Hideaki Kimura, Kian-Tat Lim, Jennie Rogers, Samuel Madden, Michael Stonebraker, Stanley B Zdonik, and Paul G Brown. Ss-db: A standard science dbms benchmark. *Under submission*, 114, 2010.
- [3] Miyuru Dayarathna and Toyotaro Suzumura. Benchmarking graph data management and processing systems: A survey. *arXiv preprint arXiv:2005.12873*, 2020.
- [4] QGIS 2.18 Documentation. Raster data), 2024. Accessed on 2024-02-27.
- [5] ECMWF. <https://cds.climate.copernicus.eu/#!/home>, 2024. Accessed on 2024-02-27.
- [6] David Haynes, Philip Mitchell, and Eric Shook. Developing the raster big data benchmark: A comparison of raster analysis on big data platforms. *ISPRS International Journal of Geo-Information*, 9(11):690, 2020.
- [7] Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András Horányi, Joaquín Muñoz-Sabater, Julien Nicolas, Carole Peubey, Raluca Radu, Dinand Schepers, et al. The era5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020.
- [8] Weiming Huang, Syed Amir Raza, Oleg Mirzov, and Lars Harrie. Assessment and benchmarking of spatially enabled rdf stores for the next generation of spatial data infrastructure. *ISPRS International Journal of Geo-Information*, 8(7):310, 2019.
- [9] George Merticariu, Dimitar Misev, and Peter Baumann. Towards a general array database benchmark: Measuring storage access. In *Big Data Benchmarking: 6th International Workshop, WBDB 2015, Toronto, ON, Canada, June 16-17, 2015 and 7th International Workshop, WBDB 2015, New Delhi, India, December 14-15, 2015, Revised Selected Papers 6*, pages 40–67. Springer, 2016.
- [10] Baudouin Raoult, Cedric Bergeron, A López Alós, Jean-Noël Thépaut, and Dick Dee. Climate service develops user-friendly data store. *ECMWF newsletter*, 151:22–27, 2017.
- [11] Copernicus Climate Change Service. Ecmwf reanalysis v5 (era5), 2023. Accessed on 2023-10-30.
- [12] Joe Hamman Stephan Hoyer. Xarray n-d labeled arrays and datasets in python, 2024. Accessed on 2024-02-28.
- [13] Gábor Szárnyas, Jack Waudby, Benjamin A Steer, Dávid Szakállas, Altan Birler, Mingxi Wu, Yuchen Zhang, and Peter Boncz. The ldbs social network benchmark: Business intelligence workload. *Proceedings of the VLDB Endowment*, 16(4):877–890, 2022.

- [14] Xinhui Tian, Shaopeng Dai, Zhihui Du, Wanling Gao, Rui Ren, Yaodong Cheng, Zhifei Zhang, Zhen Jia, Peijian Wang, and Jianfeng Zhan. Bigdatabench-s: an open-source scientific big data benchmark suite. In *2017 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pages 1068–1077. IEEE, 2017.
- [15] Ramon Antonio Rodrigues Zalipynis. Array dbms: past, present, and (near) future. *Proceedings of the VLDB Endowment*, 14(12):3186–3189, 2021.