

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC BÁCH KHOA
KHOA KHOA HỌC VÀ KỸ THUẬT MÁY TÍNH



BÁO CÁO
ĐỒ ÁN CHUYÊN NGÀNH

XÂY DỰNG CƠ SỞ DỮ LIỆU TÍCH HỢP
CHO CÁC HỆ THỐNG DỰ BÁO NGẮN HẠN
TRONG KHÍ TƯỢNG THUỶ VĂN

Ngành: Khoa học máy tính

HỘI ĐỒNG: Hội đồng 5
GVHD1: Lê Hồng Trang
GVHD2: Trương Quỳnh Chi
TKHD: Lê Thị Bảo Thu
—o0o—
SVTH1: Trần Hà Tuấn Kiệt (2011493)
SVTH2: Nguyễn Đức Thụy (2012158)



Thành phố Hồ Chí Minh, Tháng 12/2023

Lời cam đoan

Chúng tôi cam đoan rằng đồ án này là kết quả của công việc nghiên cứu của chúng tôi, được thực hiện dưới sự hướng dẫn tận tâm của PGS. TS. Lê Hồng Trang tại Trường Đại Học Bách Khoa - Đại học Quốc gia TP.HCM . Chúng tôi xác nhận rằng tất cả các thành viên trong nhóm đã đóng góp công sức và kiến thức của mình để hoàn thành công trình này.

Chúng tôi cam đoan rằng đồ án này không phải là bản sao từ bất kỳ công trình nào khác và tất cả các nguồn thông tin đã được chú thích đầy đủ theo quy tắc trích dẫn khoa học. Mọi ý kiến, thông tin hay dữ liệu từ nguồn ngoại lai đều được công bố và thực hiện theo quy tắc đạo đức và quy định của Trường Đại Học Bách Khoa - Đại học Quốc gia TP.HCM .

Chúng tôi xác nhận rằng tất cả các thành viên trong nhóm đã thực hiện nghiên cứu này với lòng tận tụy và tính chân thành cao, tuân thủ mọi nguyên tắc đạo đức nghiên cứu khoa học. Chúng tôi cũng đảm bảo rằng không có hành vi gian lận hoặc vi phạm nào đã xảy ra trong quá trình thực hiện đồ án.

Chúng tôi chấp nhận mọi trách nhiệm và hậu quả pháp lý liên quan đến nội dung của đồ án và đồng thời biết ơn sự hướng dẫn, góp ý quý báu của giáo viên PGS. TS. Lê Hồng Trang đã giúp chúng tôi hoàn thiện công trình này.

Thành phố Hồ Chí Minh, Ngày 29 tháng 12 năm 2023

NGƯỜI THỰC HIỆN

Trần Hà Tuấn Kiệt

Nguyễn Đức Thụy

Lời cảm ơn / Lời ngỏ

Chúng tôi, nhóm nghiên cứu gồm hai thành viên, xin gửi lời cảm ơn chân thành đến tất cả mọi người đã đóng góp vào thành công của đồ án này.

Đầu tiên và quan trọng nhất, chúng tôi muốn bày tỏ lòng biết ơn đến PGS. TS. Lê Hồng Trang . Sự tận tâm và kiến thức sâu rộng của thầy không chỉ hướng dẫn chúng tôi qua những thách thức của đồ án mà còn giúp chúng tôi phát triển kỹ năng nghiên cứu và phê bình.

Chúng tôi cũng muốn bày tỏ lòng biết ơn đặc biệt đến tất cả những người bạn, đồng nghiệp, và gia đình đã hỗ trợ chúng tôi trong suốt quá trình nghiên cứu. Sự góp ý và ý kiến của mọi người đã làm giàu thêm nội dung và chất lượng của đồ án.

Cuối cùng, nhưng không kém phần quan trọng, chúng tôi cảm ơn nhau - đối tác nghiên cứu đồng hành trong mỗi bước của đồ án. Sự hợp tác và đóng góp chung của chúng tôi đã tạo nên một sản phẩm mà chúng tôi tự hào.

Chúng tôi tin rằng đồ án này là một bước tiến quan trọng trong sự phát triển của chúng tôi và không thể đạt được mà không có sự hỗ trợ và đóng góp của tất cả mọi người.

Tóm tắt nội dung

Trong tài liệu này, chúng tôi giới thiệu một đề xuất toàn diện để tích hợp, điều phối và giám sát dữ liệu khí tượng thủy văn tại Việt Nam. Hệ thống được đề xuất có thể làm nền tảng cho việc xây dựng một kho dữ liệu quốc gia chuyên sâu về thông tin khí tượng.

Để làm rõ đề xuất của chúng tôi, chúng tôi đã xây dựng một bản thử nghiệm toàn diện, trong đó chúng tôi đã thực hiện mô phỏng về quy trình thu thập và truyền phát dữ liệu từ trạm thời tiết tại Nhà Bè. Đồng thời, chúng tôi đã tập trung vào quá trình chuyển đổi và tổ chức lưu trữ dữ liệu để đảm bảo sự minh bạch và hiệu suất tối ưu trong quá trình xử lý thông tin.

Bài nghiên cứu này là biểu hiện của những nỗ lực đáng kể trong việc tối ưu hóa quản lý dữ liệu khí tượng tại Việt Nam. Hệ thống đề xuất này không chỉ giải quyết các thách thức về tích hợp và điều phối mà còn đặt nền tảng cho một kho dữ liệu quốc gia vững mạnh, phục vụ cho những nhu cầu đa dạng của nghiên cứu và ứng dụng trong lĩnh vực khí tượng.

Bảng phân công công việc

STT	Họ và tên	MSSV	Phân chia công việc
1	Trần Hà Tuấn Kiệt	2011493	<ul style="list-style-type: none">- Phân tích dữ liệu- Xây dựng kiến trúc hệ thống- Nghiên cứu, tìm hiểu công nghệ và áp dụng- Xây dựng API
2	Nguyễn Đức Thụy	2012158	<ul style="list-style-type: none">- Phân tích dữ liệu- Xây dựng kiến trúc hệ thống- Nghiên cứu, tìm hiểu công nghệ và áp dụng- Xây dựng Data Flow

Mục lục

1	Giới thiệu	9
1.1	Phát biểu bài toán	9
1.2	Động lực thực hiện	10
1.3	Mục tiêu dự án	11
1.4	Phạm vi dự án	12
2	Cơ sở lý thuyết	13
2.1	Hệ Cơ sở dữ liệu	13
2.2	Hệ thống dữ liệu	15
2.3	Lý thuyết về khí tượng	16
2.3.1	Các Khái Niệm Cơ Bản	17
2.3.1.1	Radar thời tiết	17
2.3.1.2	Phương trình Radar và độ phải hồi vô tuyến	20
2.3.1.3	Vận tốc xuyên tâm	22
2.4	Định dạng dữ liệu trong phân tích khí tượng	22
2.4.1	Định dạng SIGMET - raw format (Vaisala)	22
2.4.2	Định dạng NETCDF - Network Common Data Form	23
2.5	Công nghệ sử dụng	24
3	Phân tích và thiết kế hệ thống	25
3.1	Tổng quan	25
3.2	Thăm dò dữ liệu	25
3.3	Mô hình cơ sở dữ liệu	25
3.4	Kiến trúc hệ thống	25



4	Hiện thực	26
4.1	Luồng dữ liệu	26
5	Kiểm thử	28
5.1	Unit testing	28
5.2	Integrated testing	28
6	Hướng phát triển	29

Danh sách hình vẽ

2.1	Hệ thống radar thời tiết - [10]	16
2.3	So sánh kết quả thu được từ phương pháp PPI và RHI - [3]	19
2.4	Minh hoạ hệ số phản xạ từ dữ liệu radar Nhà Bè	21
2.5	Minh hoạ các tình huống vận tốc mà radar Doppler có thể quan sát. (a) Phương của gió trùng tại M trùng với đường kính đường tròn có tâm tại radar, radar có thể xác định được vận tốc tại đây. (b) Phương của gió trùng với tiếp tuyến của đường tròn, radar không thể xác định được vận tốc. (c) Phân tích hướng gió tại M thành 2 vận tốc vuông góc nhau, radar chỉ xác định được vector vận tốc theo M_r	22
2.6	Thông tin radar ở định dạng NETCDF. Số chiều của bộ dữ liệu tổng cộng là 2975 chiều, được phân nhóm cho 4 nhãn khác nhau.	24

Danh sách bảng

2.1	Tương quan hệ số phản xạ của radar và giáng thủy - Stull [10]	21
-----	---	----

Chương 1

Giới thiệu

1.1 Phát biểu bài toán

The National Center for Hydrometeorological Forecasting (NCHMF), abbreviated as "Trung tâm Dự báo khí tượng thủy văn quốc gia" in Vietnamese, is an organizational unit under the General Department of Meteorology and Hydrology, Ministry of Natural Resources and Environment[6]. The National Hydro-Meteorological Forecasting Center has several crucial missions, including the establishment and presentation of standards and technical regulations for meteorological and hydrological forecasting, the operation of the national forecasting and warning system, monitoring and reporting on weather conditions and climate change, issuing and disseminating forecast bulletins and warnings, and participating in international meteorological agreements. Additionally, the center is responsible for conducting research, application, and technology transfer related to forecasting and warning, and implementing administrative reform and anti-corruption measures. These key missions contribute significantly to the center's role in ensuring public safety and providing essential meteorological and hydrological information.

There is a deliberate focus on those aspects of climate data management that are of interest to NMHSs wishing to make the transition to a modern climate database management system and, just as important, on what skills, systems and processes need to be in place to ensure that operations are sustained. In the context of the ever-growing complexity of climate change, the task of creating an integrated database for short-term forecasting systems in the fields of meteorology and hydrology poses a considerable challenge. The question at hand is how we can

optimize the management of weather information from multiple sources and store it efficiently in a database. This optimization is crucial to ensure the provision of synchronized and high-quality information to support forecasting systems.

1.2 Động lực thực hiện

Information about the weather has been recorded in manuscript form for many centuries. The early records included notes on extreme and, sometimes, catastrophic events and also on phenomena such as the freezing and thawing dates of rivers, lakes and seas, which have taken on a higher profile with recent concerns about climate change. Specific journals for the collection and retention of climatological information have been used over the last two or three centuries (WMO 2005). The development of instrumentation to quantify meteorological phenomena and the dedication of observers to maintaining methodical, reliable and well-documented records paved the way for the organized management of climate data. Since the 1940s, standardized forms and procedures gradually became more prevalent and, once computer systems were being used by NMHSs, these forms greatly assisted the computerized data entry process and consequently the development of computer data archives. The latter part of the twentieth century saw the routine exchange of weather data in digital form and many meteorological and related data centers took the opportunity to directly capture and store these in their databases. Much was learned about automatic methods of collecting and processing meteorological data in the late 1950s, a period that included the International Geophysical Year and the establishment of the World Weather Watch. The WMO's development of international guidelines and standards for climate data management and data exchange assisted NMHSs in organizing their data management activities and, less directly, also furthered the development of regional and global databases. Today, the management of climate records requires a systematic approach that encompasses paper records, microfilm/microfiche records and digital records, where the latter include image files as well as the traditional alphanumeric representation.

Before electronic computers, mechanical devices played an important part in the development of data management. Calculations were made using comptometers, for example, with the results being recorded on paper. A major advance occurred with the introduction of the Hollerith system of punch cards, sorters and tabulators. These cards, with a series of punched holes

recording the values of the meteorological variables, were passed through the sorting and tabulating machines enabling more efficient calculation of statistics. The 1960s and 1970s saw several NMHSs implementing electronic computers and gradually the information from many millions of punched cards was transferred to magnetic tape. These computers were replaced with increasingly powerful mainframe systems and data were made available online through developments in disk technology.

Aside from advances in database technologies, more efficient data capture was made possible through the mid-to-late 1990s with an increase in automatic weather stations (AWSs), electronic field books (i.e. on-station notebook computers used to enter, quality control and transmit observations), the Internet and other advances in technology. Not surprisingly, there are a number of trends already underway that suggest there are many further benefits for NMHSs in managing data and servicing their clients. The Internet is already delivering greatly improved data access capabilities and, providing security issues are managed, we can expect major opportunities for data managers in the next five to ten years. In addition, Open Source⁷ relational database systems may also remove the cost barriers to relational databases for many NMHSs over this period.

1.3 Mục tiêu dự án

It is essential that both the development of climate databases and the implementation of data management practices take into account the needs of the existing, and to the extent that it is predictable, future data users. While at first sight this may seem intuitive, it is not difficult to envisage situations where, for example, data structures have been developed that omit data important for a useful application or where a data centre commits too little of its resources to checking the quality of data for which users demand high quality.

In all new developments, data managers should either attempt to have at least one key data user as part of their project team or undertake some regular consultative process with a group of user stakeholders. Data providers or data users within the organization may also have consultative processes with end users of climate data (or information) and data managers should endeavour to keep abreast of both changes in needs and any issues that user communities have. Put simply, data management requires awareness of the needs of the end users.

At present, the key demand factors for data managers are coming from climate prediction, climate change, agriculture and other primary industries, health, disaster/emergency management, energy, natural resource management (including water), sustainability, urban planning and design, finance and insurance. Data managers must remain cognizant that the existence of the data management operation is contingent on the centre delivering social, economic and environmental benefit to the user communities it serves. It is important, therefore, for the data manager to encourage and, to the extent possible, collaborate in projects which demonstrate the value of its data resource. Even an awareness of studies that show, for example, the economic benefits from climate predictions or the social benefits from having climate data used in a health warning system, can be useful in reminding senior NMHS managers or convincing funding agencies that data are worth investing in. Increasingly, value is being delivered through integrating data with application models (e.g. crop simulation models, economic models) and so integration issues should be considered in the design of new data structures.

1.4 Phạm vi dự án

The project will focus on Ho Chi Minh City, a large urban area with a unique climate and significant impact on the daily lives of the community.

We will conduct research and collect data from multiple meteorological and hydrological monitoring stations in the city to ensure diversity and representation of local weather conditions. Once we have gathered sufficient data, we will proceed with preprocessing and standardizing the data to ensure accuracy and consistency of the dataset. Finally, we will build an integrated database to store and manage information from various sources, creating a unified and reliable data source. The integrated information source will be used to provide multidimensional and detailed data on current weather conditions and short-term forecasts. The goal is to help the community and relevant entities better prepare for unpredictable weather fluctuations.

Having an integrated database is expected to improve the effectiveness of weather forecasting models, enabling individuals and businesses to intelligently and safely cope with challenging weather conditions.

Chương 2

Cơ sở lý thuyết

2.1 Hệ Cơ sở dữ liệu

Database systems perform vital functions for all sorts of organizations because of the growing importance of using and managing data efficiently. A database system consists of a software, a database management system (DBMS) and one or several databases. DBMS is a set of programs that enables users to store, manage and access data. In other words database is processed by DBMS, which runs in the main memory and is controlled by the respective operating system

A database is a logically coherent collection of data with some inherent meaning and represents some aspects of the real world. A random assortment of data cannot be referred to as a database. Databases draw a sharp distinction between data and information. Data are known facts that can be recorded and that have implicit meaning. Information is data that have been organized and prepared in a form that is suitable for decision-making. Shortly information is the analysis and synthesis of data. The most fundamental terms used in database approach are "entity", "attribute" and "relationship". An entity is something that can be identified in the users' work environment, something that the users want to track. It may be an object with a physical or conceptual existence. An attribute is a property of an entity. A particular entity will have a value for each of its attributes. The attribute values that describe each entity become a major part of data stored in the database

Database Management System is a general-purpose software system designed to manage large bodies of information facilitating the process of defining, constructing and manipulating databases for various applications. Specifying data types, structures and constraints for the data

to be stored in the database is called defining a database. Constructing the database is the process of storing data itself on some storage medium that is controlled by the DBMS. Querying to retrieve specific data, updating the database to reflect changes and generating reports from the data are the main concepts of manipulating a database. The DBMS functions as an interface between the users and the database ensuring that the data is stored persistently over long periods of time, independent of the programs that access it [7]. DBMS can be divided into three subsystems; the design tools subsystem, the run time subsystem and the DBMS engine.

The design tools subsystem has a set of tools to facilitate the design and creation of the database and its applications. Tools for creating tables, forms, queries and reports are components of this system. DBMS products also provide programming languages and interfaces to programming languages. The run time subsystem processes the application components that are developed using the design tools. The last component of DBMS is the DBMS engine which receives requests from the other two components and translates those requests into commands to the operating system to read and write data on physical media [5].

Database approach has several advantages over traditional file processing in which each user has to create and define files needed for a specific application. In these systems' duplication of data is generally inevitable causing wasted storage space and redundant efforts to maintain common data up-to date. In database approach data is maintained in a single storage medium and accessed by various users. The self-describing nature of database systems provides information not only about database itself but also about the database structure such as the type and format of the data. A complete definition and description of database structure and constraints, called meta-data, is stored in the system catalog. Data abstraction is a consequence of this self-describing nature of database systems allowing program data independence. DBMS access programs do not require changes when the structure of the data files are changed hence the description of data is not embedded in the access programs. This property is called program-data independence. Support of multiple views of data is another important feature of database systems, which enables different users to view different perspective of database dependent on their requirements. In a multi-user database environment users probably have access to the same data at the same time as well as they can access different portions of database for modification. Concurrency control is crucial for a DBMS so that the results of the updates are correct. The DBMS software is to ensure that concurrent transactions operate correctly when several users

are trying to update the same data

Using a DBMS also eliminates unnecessary data redundancy. In database approach each primary fact is generally recorded in only one place in the database [6]. Sometimes it is desirable to include some limited redundancy to improve the performance of queries when it is more efficient to retrieve data from a single file instead of searching and collecting data from several files, but this data duplication is controlled by DBMS so as to prohibit inconsistencies among files. By eliminating data redundancy inconsistencies among data are also reduced [5]. Reducing redundancy improves the consistency of data while reducing the waste in storage space. DBMS gives the opportunity of data sharing to the users. Sharing data often permits new data processing applications to be developed without having to create new data files. In general, less redundancy and greater sharing lead to less confusion between organizational units and less time spent resolving errors and inconsistencies in reports. The database approach also permits security restrictions. In a DBMS different types of authorizations are accepted in order to regulate which parts of the database various users can access or update.

2.2 Hệ thống dữ liệu

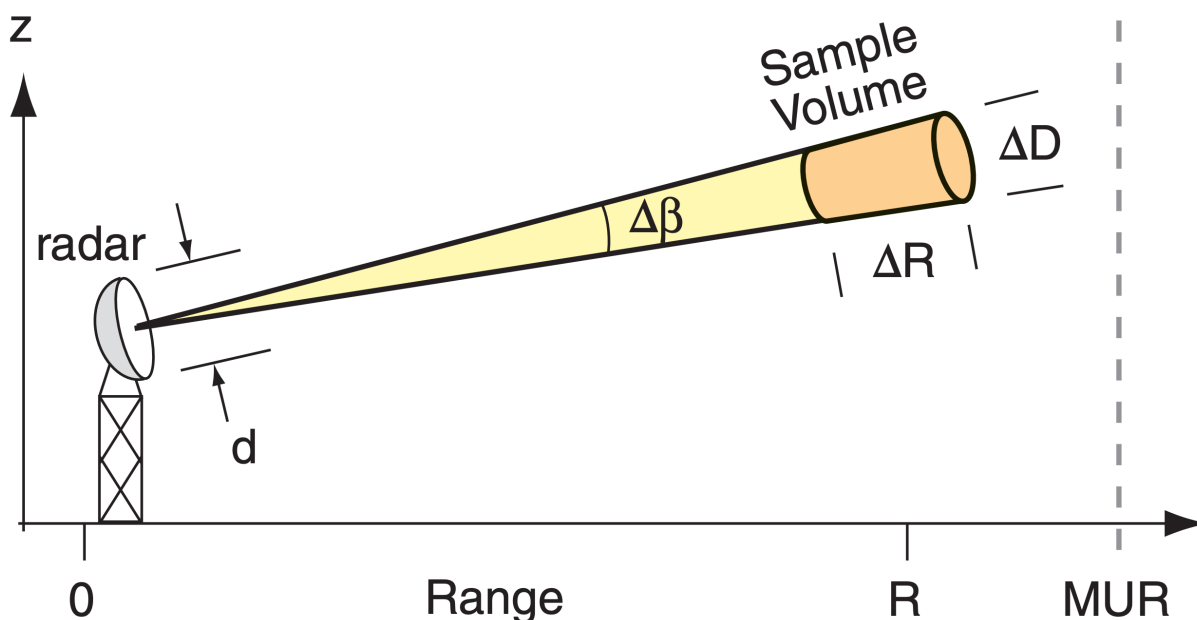
In contemporary computing environments, the landscape is increasingly dominated by data-intensive applications, a departure from the traditional emphasis on compute-intensive tasks. The limiting factor for these applications seldom resides in the sheer computational power of the CPU; rather, the primary challenges typically revolve around the magnitude of the data, its intricate structures, and the rapidity with which it undergoes change. Unlike compute-intensive operations that heavily rely on processing speed, data-intensive applications, dealing with extensive datasets, intricate data structures, or swiftly evolving information, necessitate adept strategies for storage, retrieval, and manipulation. Consequently, effectively addressing the multifaceted dynamics of data becomes paramount, highlighting the imperative for sophisticated data management and processing techniques to optimize performance in the face of these intricate challenges.

Why should we amalgamate these diverse elements within the overarching label of data systems? Recent years have witnessed the emergence of a plethora of novel tools for data storage and processing, each meticulously optimized for an array of distinct use cases, rendering them

incompatible with conventional categorizations [9]. Consider, for instance, datastores that concurrently function as message queues (e.g., Redis) or message queues equipped with database-like durability assurances (such as Apache Kafka). The demarcation lines between these categories are progressively fading, reflecting a landscape where boundaries are increasingly ambiguous.

Moreover, a growing number of applications now present challenges of such magnitude or diversity that a solitary tool is no longer sufficient to fulfill all its data processing and storage requisites. Instead, the workload is deconstructed into tasks amenable to efficient execution by individual tools. These disparate tools are then intricately interwoven using application code, offering a nuanced and adaptable approach to the multifaceted demands of contemporary data management and processing.

2.3 Lý thuyết về khí tượng



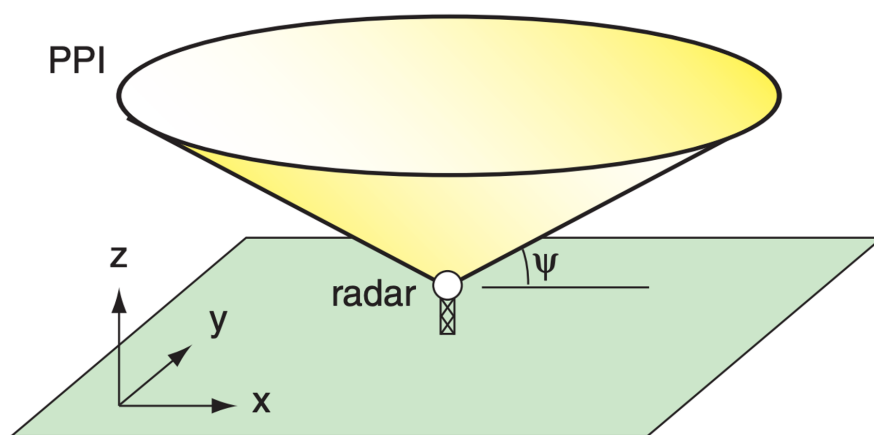
Hình 2.1: Hệ thống radar thời tiết - [10]

2.3.1 Các Khái Niệm Cơ Bản

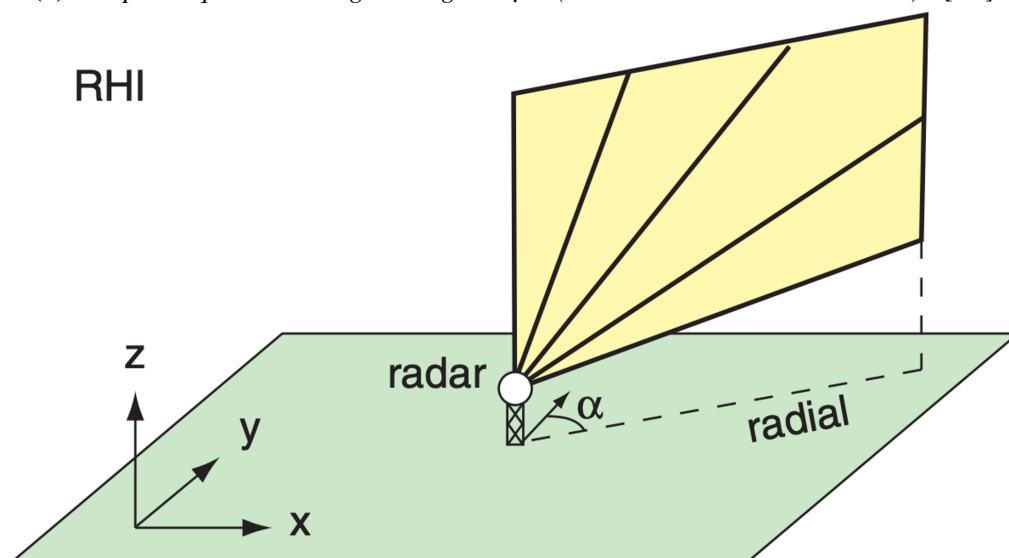
2.3.1.1 Radar thời tiết

¹ Radar thời tiết là một loại cảm biến có khả năng phát sóng vô tuyến (bước sóng trong phạm vi từ 250 - 1000 kW) [10]. Để gia tăng cường độ sóng, một chảo anten (antenna dish) hình parabol được sử dụng nhằm hội tụ bước sóng. Radar có thể nâng và hạ (tùy theo yêu cầu) để thu nhập thông tin tại các vị trí chỉ định trong không gian 3 chiều.

Thông thường, các radar được lập trình để quét theo góc hướng (azimuth) 360° , mỗi vòng sẽ quét ở một góc nâng khác nhau. Như vậy, radar sẽ mất khoảng từ 4 đến 10 phút để hoàn thành một lần quét.



(a) Sản phẩm quét tròn với góc nâng cố định (Plan-Position Indicator - PPI) - [10]

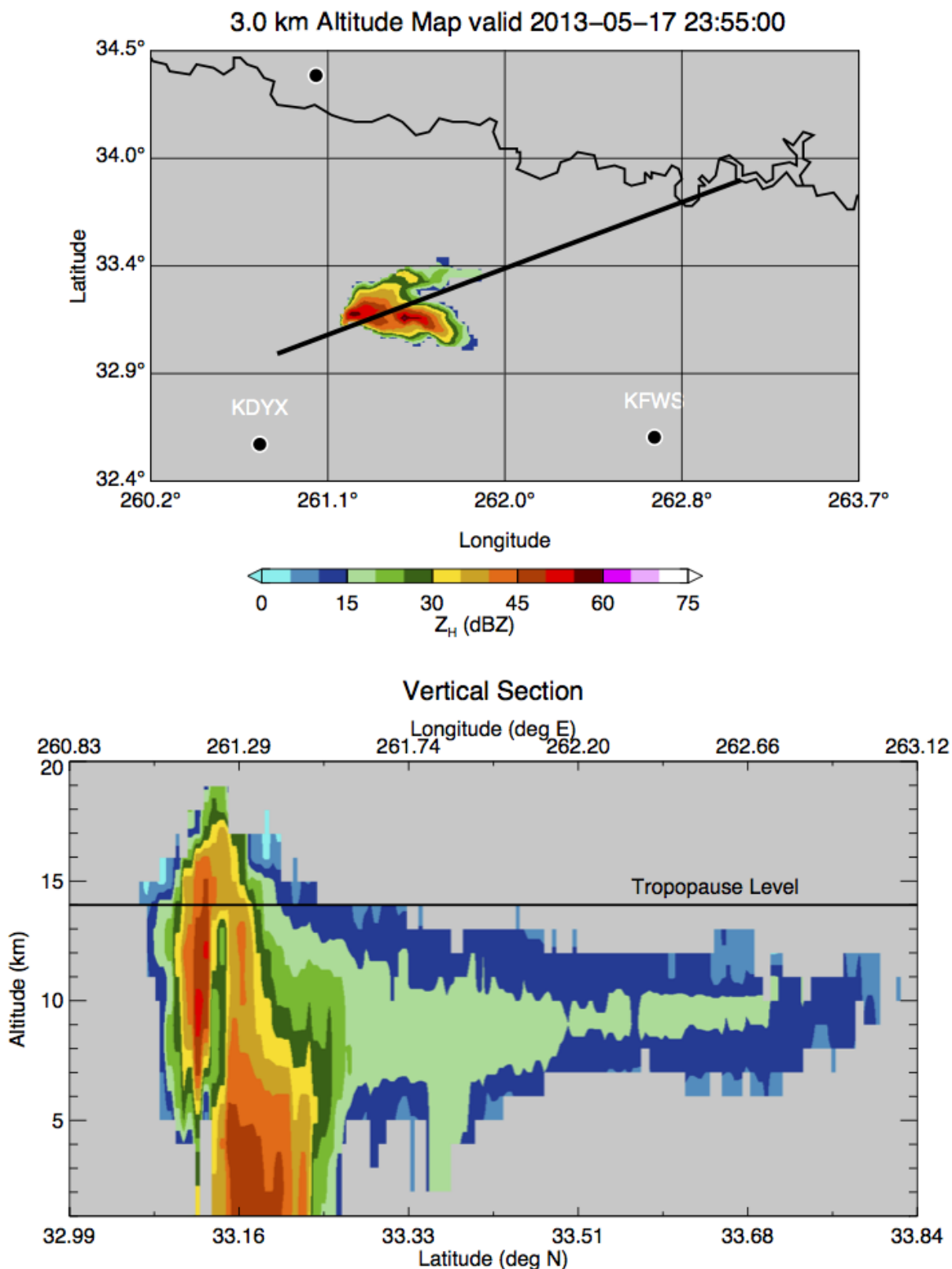


(b) Sản phẩm quét thẳng đứng ở một góc phương vị nhất định (Range Height Indicator) - [10]

¹Tên tiếng Việt của các thuật ngữ sẽ được căn cứ dựa trên TCVN 12636-12 : 2021 [4]



Đối với biểu diễn PPI, radar sẽ quét toàn bộ góc hướng, nhưng chỉ ở một góc nâng nhất định. Kết quả thu được tương tự một bản đồ trên mặt phẳng. Với RHI, radar sẽ giữ nguyên góc hướng nhưng thay đổi về góc nâng. Kết quả thu được giúp người xem có cái nhìn rõ nét hơn về chiều cao, kích thước của hiện tượng khí tượng.



Hình 2.3: So sánh kết quả thu được từ phương pháp PPI và RHI - [3]

2.3.1.2 Phương trình Radar và độ phản hồi vô tuyến

Tại một thời điểm, radar sẽ phát ra một luồng sóng trong khoảng thời gian ngắn ($\Delta t = 0.5 - 10\mu s$). Lúc này, tùy thuộc mật độ các phân tử tự do trong không khí (hơi nước, khói bụi, ...), năng lượng của bước sóng này sẽ bị hấp thụ một phần. Cường độ bước sóng mà radar nhận được sẽ nhỏ hơn cường độ sóng ban đầu. Tỷ lệ này được thể hiện thông qua **Phương trình radar** [10]:

$$\left[\frac{P_R}{P_T} \right] = [b] \cdot \left[\frac{|K|}{L_a} \right]^2 \cdot \left[\frac{R_1}{R} \right]^2 \cdot \left[\frac{Z}{Z_1} \right]$$

Trong đó, các biến của phương trình gồm có:

- $|K|$ không có đơn vị:
 - $|K|^2 \approx 0.93$ cho các hạt nước lỏng
 - $|K|^2 \approx 0.208$ cho tinh thể băng
- $R(\text{km})$: khoảng cách
- $R_1 = \sqrt{Z_1 \cdot c \cdot \Delta t / \lambda^2}$: hệ số khoảng cách
- Z : Hệ số phản hồi vô tuyến của Radar
- $Z_1 = 1 \text{ mm}^6 \text{ m}^{-3}$: hệ số đơn vị phản hồi vô tuyến

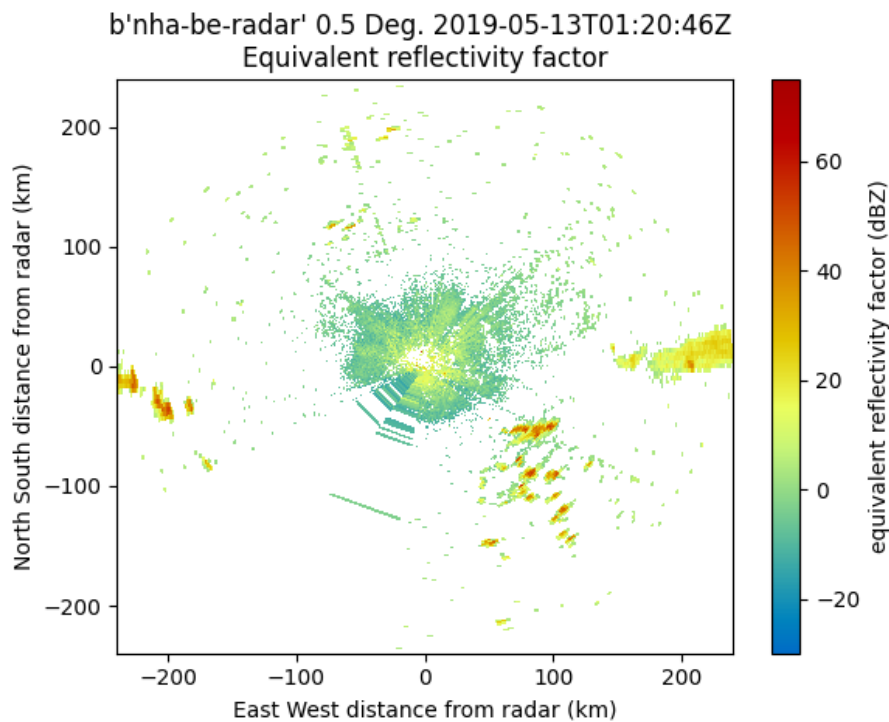
Từ phương trình Radar, ta suy ra được công thức tính độ phản hồi vô tuyến:

$$\text{dBZ} = 10 \left[\log \left(\frac{P_R}{P_T} \right) + 2 \log \left(\frac{R}{R_1} \right) - 2 \log \left| \frac{K}{L_a} \right| - \log(b) \right]$$

Các nhà khí tượng thủy văn học thường quan tâm đến con số này vì nó tỉ lệ thuận với mức độ giáng thủy (precipitation).

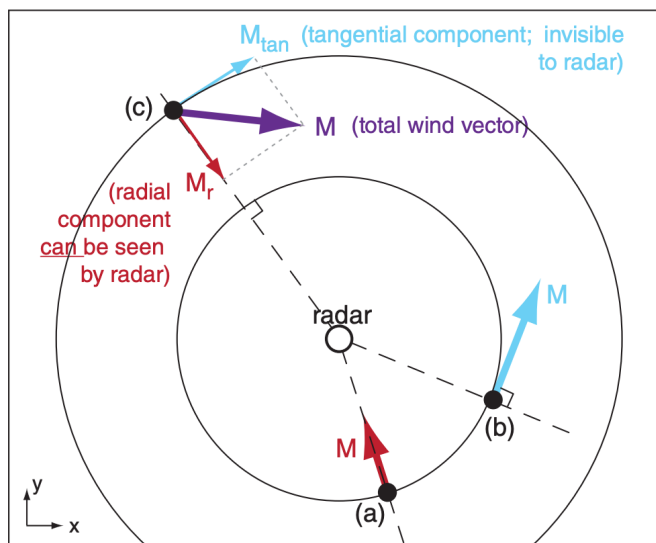
Giá trị (dBZ)	Thời tiết
-28	Sương mù
-12	Không khí trong lành
25 - 30	Tuyết khô / mưa nhẹ
40 - 50	Mưa lớn
75	Mưa đá khổng lồ

Bảng 2.1: Tương quan hệ số phản xạ của radar và giáng thủy - Stull [10]



Hình 2.4: Minh họa hệ số phản xạ từ dữ liệu radar Nhà Bè

2.3.1.3 Vận tốc xuyên tâm



Hình 2.5: Minh họa các tình huống vận tốc mà radar Doppler có thể quan sát. (a) Phương của gió trùng tại M trùng với đường kính đường tròn có tâm tại radar, radar có thể xác định được vận tốc tại đây. (b) Phương của gió trùng với tiếp tuyến của đường tròn, radar không thể xác định được vận tốc. (c) Phân tích hướng gió tại M thành 2 vận tốc vuông góc nhau, radar chỉ xác định được vector vận tốc theo M_r .

Khi các sóng vô tuyến từ các radar Doppler này truyền đến các phân tử trong không khí, sự chuyển dịch vị trí của các hạt này làm lệch pha giữa tín hiệu truyền đi và tín hiệu nhận lại được. Các radar sẽ căn cứ vào thông tin này để tính toán vận tốc gió tại các điểm trong không gian.

2.4 Định dạng dữ liệu trong phân tích khí tượng

2.4.1 Định dạng SIGMET - raw format (Vaisala)

Vaisala là một công ty Phần Lan, chuyên về các lĩnh vực thuộc môi trường và khí tượng thủy văn. Định dạng RAW (trong một số tài liệu còn gọi là SIGMET [1]) là một trong những định dạng lưu trữ mà công ty này đã phát triển ra nhằm thực hiện tổ chức dữ liệu xuất ra từ các thiết bị radar của họ.

Một số điểm nổi bật về định dạng này có thể kể đến như:

- Nội dung của file được phân thành một **block** lớn. Mỗi block có kích thước đúng 6144 bytes. Kích thước này vừa đúng bằng kích thước lưu trữ chính trên các thiết bị băng từ cũ.
- File thường là tổng hợp của tất cả những lần radar tiến hành quét dữ liệu.

- Các phần dữ liệu (record) sẽ được sắp xếp trong phạm vi 1 block (6144 bytes). Trong trường hợp phần block còn dư, dữ liệu sẽ được đệm thêm các số 0.

Với những đặc điểm kể trên, có thể nhận thấy các ưu điểm chính từ việc lưu trữ định dạng RAW bao gồm: [11]

- Thân thiện với các loại băng từ. Đây thường là các thiết bị phổ biến trước đây, hiện nay vẫn được sử dụng rộng rãi nhờ vào hiệu quả từ mức dung lượng / chi phí.
- Nhờ sử dụng cơ chế block, SIGMET giúp các hệ thống lưu trữ thực hiện các biện pháp hồi phục (error recovery) trên mức block.

Nhược điểm chính mà nhóm quan ngại là khả năng mapping (liên kết) giữa cấu trúc khi lưu trữ trong ổ cứng và trên băng từ.

2.4.2 Định dạng NETCDF - Network Common Data Form

NetCDF (Network Common Data Form) là một định dạng tệp tin linh hoạt được thiết kế chặt chẽ để lưu trữ dữ liệu khoa học đa chiều. Trong hệ thống thư viện netCDF, có nhiều định dạng nhị phân được hỗ trợ, mỗi định dạng đóng góp vào tính linh hoạt và khả năng mở rộng của quản lý dữ liệu [8]. Đáng chú ý, các định dạng này bao gồm:

1. Định dạng Classic: Ban đầu được sử dụng trong phiên bản đầu tiên của netCDF và vẫn là lựa chọn mặc định cho việc tạo tệp tin.
2. Định dạng 64-bit Offset: Giới thiệu từ phiên bản 3.6.0, định dạng này hỗ trợ kích thước biên và tệp tin lớn hơn.
3. Định dạng netCDF-4/HDF5: Xuất hiện từ phiên bản 4.0, sử dụng định dạng dữ liệu HDF5 với một số hạn chế.
4. Định dạng HDF4 SD: Hỗ trợ chủ yếu cho việc đọc dữ liệu.
5. Định dạng CDF5: Hỗ trợ được đồng bộ với dự án parallel-netcdf.

Tất cả các định dạng này đều thể hiện tính tự mô tả, với một phần tiêu đề chi tiết mô tả cấu trúc của tệp tin, bao gồm các mảng dữ liệu và siêu dữ liệu tệp tin dưới dạng thuộc tính tên/giá

trị. Thiết kế này đảm bảo tính độc lập với nền tảng, với các vấn đề như endianness được giải quyết một cách linh hoạt thông qua các thư viện phần mềm.

Hãy xem xét ví dụ cụ thể về việc lưu trữ các thông số khí tượng quan trọng như nhiệt độ, độ ẩm, áp suất, tốc độ và hướng gió trong các tệp tin netCDF. Điều này minh họa khả năng của định dạng này trong xử lý các bộ dữ liệu khoa học đa dạng, cung cấp một phương tiện mạnh mẽ và linh hoạt để quản lý thông tin đa chiều.

```
● → titan2023 ncdump -h radar.nc
netcdf radar {
dimensions:
    time = UNLIMITED ; // (1748 currently)
    range = 1198 ;
    sweep = 5 ;
    string_length = 32 ;
```

Hình 2.6: Thông tin radar ở định dạng NETCDF. Số chiều của bộ dữ liệu tổng cộng là 2975 chiều, được phân nhóm cho 4 nhãn khác nhau.

Bắt đầu từ phiên bản 4.0, API netCDF giới thiệu khả năng sử dụng định dạng dữ liệu HDF5. Sự tích hợp quan trọng này cho phép người dùng netCDF tạo tệp tin HDF5, mở khóa những lợi ích như kích thước tệp tin lớn hơn đáng kể và hỗ trợ cho nhiều chiều không giới hạn. Bước tiến này đánh dấu một bước quan trọng hướng tới việc tận dụng những ưu điểm mở rộng của định dạng HDF5.

NetCDF Classic và Định dạng 64-bit Offset là tiêu chuẩn quốc tế của Open Geospatial Consortium[2], thể hiện sự chắc chắn và độ tin cậy trong việc đảm bảo khả năng tương thích và mở rộng của định dạng netCDF trên toàn cầu.

2.5 Công nghệ sử dụng

Chương 3

Phân tích và thiết kế hệ thống

3.1 Tổng quan

3.2 Thăm dò dữ liệu

3.3 Mô hình cơ sở dữ liệu

3.4 Kiến trúc hệ thống

Chương 4

Hiện thực

4.1 Luồng dữ liệu

Phần hiện thực của nhóm sẽ nằm trong năm bước còn lại trong mô tả tại hình ??.

Tại bước 3, nhóm sẽ setup (cài đặt) một server SFTP đơn giản. SFTP là một giao thức đơn giản và phổ biến. Hiện nay, có rất nhiều những thư viện và công cụ để hỗ trợ giao tiếp dựa trên giao thức này. Ngoài ra, so với FTP, giao thức kể trên còn đảm bảo tính bảo mật trong suốt quá trình chuyển dịch dữ liệu. Tùy thuộc vào mức độ cho phép, nhóm có thể hỗ trợ phía trạm quan trắc xây dựng các scripts (kịch bản) để tự động forward (chuyển tiếp) các file sau khi đã được xử lý tại đây. Hoặc ngược lại, phía trạm quan sát có thể gửi file đến server trên một cách thủ công.

Khi file đã được upload đến server SFTP, nhóm sử dụng Airflow để điều hành tất cả các luồng ETL hiện có trong hệ thống chung. Ở thời điểm hiện tại, nhóm chỉ dừng lại với một DAG duy nhất, để xử lý dữ liệu đến từ trạm quan trắc Nhà Bè. Airflow sẽ tiến hành quan sát những file được thêm mới vào server SFTP của chúng ta và khởi chạy ETL. Việc lựa chọn Apache Spark ở đây dựa trên khối lượng dữ liệu và độ phức tạp được đặt ra. Nếu lượng dữ liệu là không quá nhiều cho mỗi file SIGMET mới, và bản thân Python có xử lý được, không cần thiết phải sử dụng Spark ở bước này.

Dữ liệu về khí tượng khi được đưa đến cơ sở hạ tầng của nhóm sẽ được phân ra hai luồng chính: Những metadata (thông tin mở rộng) của dữ liệu gốc như ngày tạo ra, kích thước, thời điểm ghi nhận, ... sẽ được lưu trong một RDBMS (hệ quản trị cơ sở dữ liệu quan hệ) truyền thống. Cụ thể ở đây, nhóm lựa chọn PostgreSQL nhờ vào độ phổ biến và mức độ am hiểu của

nhóm. Các metadata lưu trữ ở đây giúp cơ sở dữ liệu của nhóm nhanh chóng phản hồi các query (truy vấn) mà chưa cần trực tiếp phải sử dụng đến dữ liệu gốc. Một số query phổ biến có thể kể đến như:

- Các mốc thời gian đang được ghi nhận bao gồm những gì? (Ví dụ: từ ngày 21/11/2023 cho đến ngày 17/12/2023)
- Tại thời điểm x , radar có tọa độ địa lý là bao nhiêu?
- Các trường dữ liệu đang được lưu trữ là gì?

Bên cạnh đó, DB (cơ sở dữ liệu) trên còn đóng vai trò như mục index (chỉ mục) giúp hệ thống nhanh chóng xác định vị trí lưu trữ của dữ liệu gốc.

Với các dữ liệu về khí tượng thủy văn cụ thể, nhóm nhận thấy rằng sẽ không thật sự hiệu quả khi lưu trữ chúng trực tiếp trên các DBMS trên. Đồng thời, nhóm nhận thấy việc lưu trữ dữ liệu trên file vẫn đem đến một kích thước tổng quan hợp lý. Vì vậy, nhóm quyết định sẽ tách phần dữ liệu thô ra và lưu trữ trực tiếp trên các files. Đồng thời kết hợp với các index (đã đề cập ở trên) để tăng tốc quá trình truy xuất.

Để tạo cửa ngõ cho việc truy vấn dữ liệu, phục vụ cho các bên về model, machine learning và AI, ... nhóm sẽ phát triển một server Backend đơn giản, sử dụng FastAPI của Python để giúp tăng tốc độ phát triển giải pháp. Tại bước 5, backend sẽ nhận dữ liệu truy vấn dưới định dạng REST API (tại bước 6), truy vấn dữ liệu trong DB của metadata và trong các file dữ liệu và trả về kết quả đạt được. Ở những lần train khác nhau, các bên của Machine Learning có thể kết nối đến server này để lấy dữ liệu.

Cần nói thêm, toàn bộ hệ thống sẽ được phát triển và vận hành theo hướng containerize (đóng gói) và sẽ được deploy (triển khai) trên nền tảng Kubernetes. Việc này thể hiện khả năng của hệ thống trong việc duy trì tính sẵn sàng cao (High-Availability) cũng như dễ dàng trong việc duy trì giải pháp. Trong phạm vi phần minh họa này, nhóm sẽ chỉ dừng lại với việc triển khai trên một cụm máy tính nhúng Raspberry Pi.

Sau cùng, tại bước 7, nhóm muốn đề xuất thêm một vấn đề. Nếu phù hợp, nhóm có thể xây dựng thêm một DataLoader (bộ nạp dữ liệu) để phục vụ nhanh chóng đến các nhóm làm model khác. Một trong những thư viện phổ biến hiện nay của các bên AI là Pytorch, nên nhóm sẽ tiếp cận với nền tảng này trước.

Chương 5

Kiểm thử

5.1 Unit testing

5.2 Integrated testing

Trong nghiên cứu này, chúng tôi đã thành công trong việc xây dựng một Proof-of-Concept (chứng minh khái niệm) mang tính ứng dụng cao, nhằm mục đích giảm thiểu các công đoạn trong quy trình làm việc thông thường. Đây là một bước quan trọng để tối ưu hóa và cải thiện hiệu suất làm việc trong các ngữ cảnh nghiên cứu và thực tế.

Chúng tôi đã đặt ra mục tiêu tạo ra một hệ thống linh hoạt có khả năng thích ứng cao, giúp giảm bớt những bước phức tạp trong quy trình công việc. Bằng cách này, chúng tôi không chỉ giúp tăng cường hiệu suất mà còn giảm áp lực công việc đối với nhân sự, tạo điều kiện thuận lợi cho sự sáng tạo và tập trung vào các nhiệm vụ chính.

Chúng tôi không chỉ dừng lại ở việc phát triển hệ thống mà còn đề xuất các chiến lược triển khai linh hoạt, nhấn mạnh sự tích hợp dễ dàng vào môi trường làm việc hiện tại của những người đang thực hiện công việc thu thập thông tin và dự báo thời tiết.

Chương 6

Hướng phát triển

Phát triển thành Hệ thống nền tảng dữ liệu thời tiết

Hệ thống nền tảng dữ liệu thời tiết (Weather Data Platform - WDP) được phát triển với mục tiêu trở thành một giải pháp toàn diện cho việc khai thác sức mạnh của dữ liệu thời tiết. Hệ thống được thiết kế để đáp ứng những yêu cầu cụ thể của các nhà khí tượng, học giả, nghiên cứu học thuật và các nhà phát triển từ nhiều lĩnh vực khác nhau, bao gồm cả freelancers, doanh nghiệp và tổ chức phi chính phủ (Non-governmental Organizations - NGOs). WDP đóng vai trò như một trung tâm tập trung cho việc tích hợp dữ liệu thời tiết, phân tích, và nhiều tính năng khác.

Với mong muốn phát triển thành Hệ thống nền tảng dữ liệu thời tiết, chúng tôi hướng đến sự hoàn thiện và đa chiều hoá thông tin thời tiết. Không chỉ là một bảng số liệu, mà là một trải nghiệm toàn diện. Trong tương lai, bên cạnh việc tiếp tục xây dựng cơ sở dữ liệu tích hợp theo hướng đã đề xuất, chúng tôi hứa hẹn sẽ tiếp tục nghiên cứu để mở rộng và phát triển cơ sở dữ liệu tích hợp này thành hệ thống nền tảng dữ liệu thời tiết với những hướng phát triển như sau:

1. **Dữ liệu phi tuyến:** Mở rộng từ việc tích hợp dữ liệu cơ bản, chúng tôi sẽ chú trọng vào việc cung cấp dữ liệu phi tuyến, chi tiết và đa nguồn, giúp người dùng khám phá thêm về môi trường xung quanh.
2. **Trí tuệ nhân tạo thấu hiểu:** Sử dụng trí tuệ nhân tạo để thấu hiểu ngôn ngữ của thời tiết, từ những biến đổi nhỏ đến những sự kiện lớn, tạo nên một nguồn thông tin thời tiết sâu sắc và thông minh.

3. **Giao diện người dùng tương tác:** Không chỉ là việc truy cập thông tin, mà còn là việc tương tác với dự báo thời tiết. Giao diện người dùng sẽ là nơi người dùng thể hiện sự tò mò và tương tác trực tiếp với dữ liệu.
4. **Kết Nối Thông Tin Địa Lý:** Tận dụng hệ thống thông tin địa lý để mang đến cái nhìn thực tế hóa, địa bàn hóa cho dự báo thời tiết. Điều này giúp người dùng hiểu rõ hơn về tác động thời tiết đối với môi trường xung quanh họ.
5. **Tối ưu hoá hiệu suất:** đảm bảo khả năng đáp ứng nhanh chóng và đồng đều trong mọi điều kiện.
6. **Bảo mật dữ liệu:** Tăng cường an toàn dữ liệu để đảm bảo tính bảo mật và toàn vẹn của thông tin thời tiết.
7. **Hệ thống dự báo nâng cao:** Nghiên cứu và tích hợp trí tuệ nhân tạo để cải thiện khả năng dự báo và đưa ra thông tin dự báo cáo chính xác.
8. **Kiểm thử và tối ưu hoá:** Tiến hành kiểm thử hệ thống để đảm bảo tính ổn định và xử lý mọi vấn đề tiềm ẩn. Tối ưu hóa hiệu suất nếu cần.
9. **Triển khai và duy trì:** Triển khai hệ thống và duy trì một chu kỳ cập nhật đều đặn để đảm bảo rằng nó luôn cung cấp thông tin thời tiết chính xác và đáng tin cậy.

Tài liệu tham khảo

- [1] Radxconvert - Irose wiki. URL <http://wiki.lrose.net/index.php/RadxConvert>.
- [2] Ogc standard netcdf classic and 64-bit offset. <https://www.opengeospatial.org/standards/netcdf>, Accessed: 2017-12-05. Archived from the original on 2017-11-30. Retrieved 2017-12-05.
- [3] Casey. Using range height indicator scan of radar. URL <https://earthscience.stackexchange.com/questions/7222/using-range-height-indicator-scan-of-radar>. Truy cập lần cuối ngày 17/12/2023.
- [4] Tổng cục Khí tượng Thủy văn. Quan trắc khí tượng thủy văn - phần 12: Quan trắc ra đa thời tiết. *TCVN 12636-12 : 2021*, 2021.
- [5] Ramez A. Elmasri and Shamkant Navathe. *Fundamentals of Database Systems*. Addison Wesley, third edition, 1998.
- [6] TRUNG TÂM DỰ BÁO KHÍ TƯỢNG THỦY VĂN QUỐC GIA, Jun 2203.
- [7] G. Latisen and G. Vossen. *Models and Languages of Object Oriented Databases*. Addison Wesley, 1998.
- [8] Russ Rew, Glenn Davis, Steve Emmerson, Cathy Cormack, John Caron, Robert Pincus, Ed Hartnett, Dennis Heimburger, Lynton Appel, and Ward Fisher. Unidata netcdf, 1989. URL <http://www.unidata.ucar.edu/software/netcdf/>.
- [9] Michael Stonebraker and Uğur Çetintemel. 'one size fits all': An idea whose time has come and gone. In *Proceedings of the 21st International Conference on Data Engineering (ICDE)*, April 2005.

- [10] Roland Stull. Weather Radars, 12 2022. [Truy cập lần cuối ngày 17/12/2023].
- [11] *RAW Product Format - IRIS Programming Guide - IRIS Radar*.
Vaisala. URL https://ftp.sigmet.vaisala.com/files/html_docs/IRIS-Programming-Guide-Webhelp/raw_product_format.html.