

高频精细化气象格点数据实时处理系统设计与实现

李永生¹ 李高洁¹ 陈逸智¹ 张光宇²

(1 广东省气象探测数据中心,广州 510640; 2 广东省气象台,广州 510640)

摘要 以高频海量气象格点数据为研究对象,针对传统实时处理系统数据处理效率不高的问题,设计并实现了高频精细化气象格点数据实时处理系统的总体架构和业务流程,在分析海量高频次气象格点数据特点的基础上,设计和实现了符合气象业务需求的分布式存储模型;利用多通道动态感知技术实现了动态多通道的文件处理和文件到达的快速感知触发;利用实时数据快速处理技术实现基于精准位置寻址的快速数据块定位算法,进而实现数据块的精准定位;利用数据按需实时截取技术实现了在空间范围内按需进行裁剪的截取算法,进而实现数据按需抽取;实际业务应用表明,系统有效地提升了半结构气象数值预报产品数据的实时处理效率。

关键词 精准定位;实时处理;动态感知;按需截取

中图分类号: P409, TP311.1 DOI: 10.19517/j.1671-6345.20220018 文献标识码: A

引言

气象数值预报产品是典型的半结构化的格点类气象数据,具有 2 个显著的特点:①大数据量、大 I/O 访问,目前广东省气象局每天收集处理的数值预报产品的种类超过 30 种,数据频次涵盖 12 min/次到每天 2 次不等,产品的预报时效最长达到 64 天,数据量每天大约为 1 TB;②计算、存储等基础资源需求大,目前用于处理数值预报产品的服务器超过 50 台,存储资源近 600 TB,对于如此大量的资源需求,必须对产品进行集中统一处理后再供全省相关业务单位共享使用。特别是在全球气候持续异常的大背景下,各类极端天气气候事件频繁发生,气象灾害造成的损失和影响不断加重,这就对气象服务提出了更高的要求,为此气象服务业务更需要精细时间分辨率和空间分辨率的数值预报产品服务,这就使得气象数值预报产品呈现多频次高精细化的发展趋势,例如目前广东省气象局将逐步建立实时运行的广东重点区域精细模式,分辨率为 1 km,提供未来 12 h 精细预报产品,预报频次间隔时间为分钟级别,对产品的快速处理要求很高,但是广东省气象局

原有的数值预报业务处理系统是基于频次间隔数小时的数值预报产品业务进行设计开发的,面对频次达到分钟级别的数值预报产品处理就存在处理不及、业务展示慢等缺点,同时产品的种类呈现逐年增加的趋势,对现有的存储处理方式也提出了挑战,因此必须探索建立面向高频次、高精细化的气象格点数据实时处理系统,以满足实际业务需求。针对上述问题李永生等^[1]针对传统的关系型数据库在存储和管理数值预报产品数据方面存在效率低和存储能力不足的问题,初步实现了数据的分布式存储和处理;陈正旭等^[2]在综合气象网格化数据产品的结构及其应用场景上,设计了非结构化的列存储数据库;徐熙超等^[3]在 HBase 基础上提出了一个基于索引的气象结构化数据查询优化架构;本文在借鉴上述研究的基础上,聚焦非结构化气象数据的实时处理这一迫切的业务需求,通过对多通道动态感知、实时数据快速处理和数据按需实时截取等关键技术的探索,设计和实现了高频精细化气象格点数据实时处理系统,创新性地为高频精细化数值预报产品处理业务提供一体化的系统解决方案^[4-5]。

广东省科技计划项目(2018B020207012)资助

作者简介:李永生,男,1980 年生,硕士,高级工程师,研究领域为分布式数据处理、气象大数据分析应用技术,Email:879976993@qq.com

收稿日期:2022 年 1 月 14 日;定稿日期:2022 年 5 月 18 日

1 系统设计

1.1 系统架构设计

本系统功能上主要由3部分组成,分别是数据存储层、数据处理层和数据服务层^[6],如图1所示。其中,数据存储层采用分布式技术构建了分布式列数据库,基于分布式存储模型^[7]实现了不同格式的数值预报产品数据的持久化存储,并在完成存储系统管理流程设计的基础上实现数据的高效存储;数据处理层包括数据收集和数据解码2个主要功

能模块,其中数据收集功能主要完成Grib(Gridded Binary)格式数据、普通二进制格式数据和NetCDF(Network Common Data Form)格式数据的收集整理,然后按照系统时效要求完成数据源的准备,实现系统能够按照不同数据格式,分别采用不同的方式进行处理分析;数据解码功能通过多通道动态加载、数据动态感知触发、数据按需截取等技术完成数据快速解码的功能。数据服务层基于统一气象存储模型,开发RESTful(Representational State Transfer)查询接口支持数据的快速查询和访问控制。

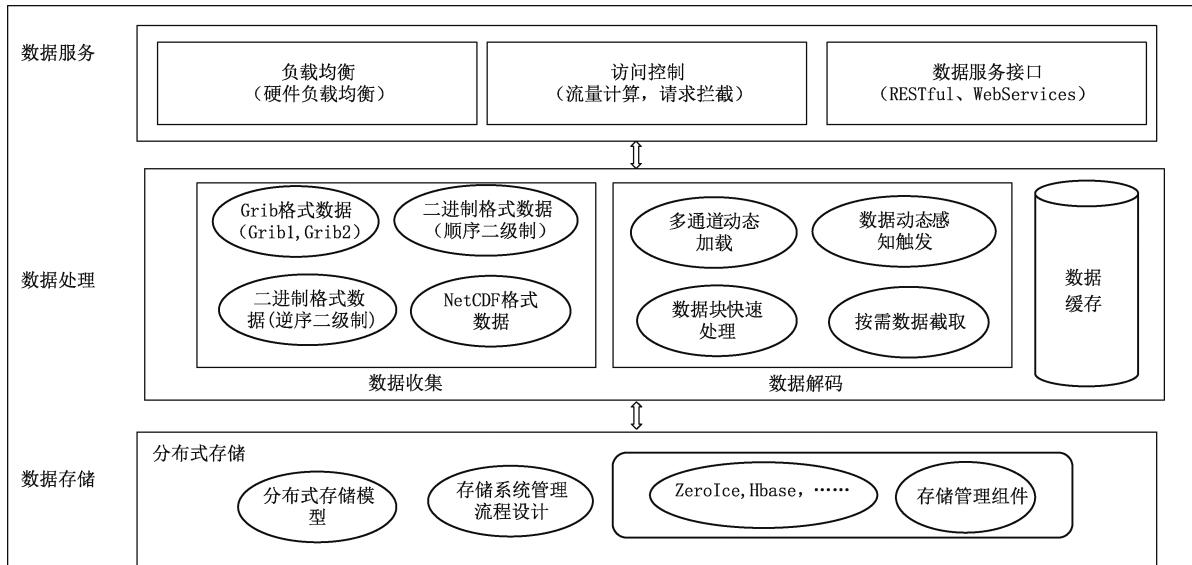


图1 系统功能结构

1.2 业务流程设计

系统的业务流程如图2所示,高频精细化的数值预报产品经收集处理模块收集整理完成后,按照不同的通道有序接入处理系统,数据处理程序动态感知不同通道数据接入后自动启动相应的处理程序,经数据解码器适配后启动快速解码功能模块完成数据的解码,形成单要素单层的平面数据块存入分布式数据库中^[8]。

1.3 分布式存储模型设计

高频精细化气象数值预报产品数据文件由若干个平面数据组成的Grib格式的文件,其中每个平面数据是一个非结构化的网格数据,网格点数由产品的空间分辨率决定,平面数据的示意如图3所示。

该平面数据Data[v]由产品的要素种类(VarName)、产品的生成时间(DataTime)、预报层次(LevelName)、预报时效(ForecastName)、成员变量

(MemberName)等维度信息唯一确定,表达为:

$$\text{Data}[v] = F(\text{VarName}, \text{DataTime}, \text{LevelName}, \\ \text{ForecastName}, \dots)$$

上述维度信息确定的平面数据Data[v]是网格点类的半结构化数据,其属性信息包括经度范围、纬度范围、空间分辨率以及数据维度上的南北走向,1代表维度上从南到北依次排列,0表示维度从北到南依次排列^[9]。

系统分布式存储采用整个平面压缩和平面分块压缩冗余存储相结合的策略^[10],整个平面压缩数据用于满足整个网格数据场检索需求,平面分块压缩数据用于满足定点定量预报数据检索需求,其中分块压缩每个列存储100×100个格点的压缩数据,在进行定点定量预报数据查询时,系统先根据定点经纬度计算出其在网格中的偏移位置,从而确定该经纬度对应的分布式存储数据表中的列,从该列中读

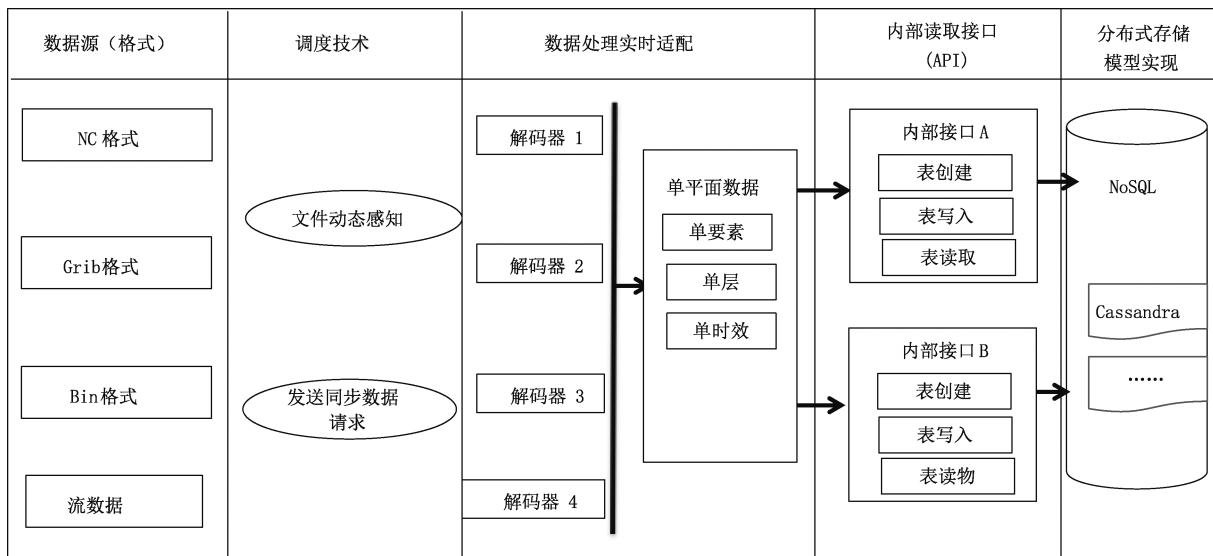


图 2 系统业务处理流程

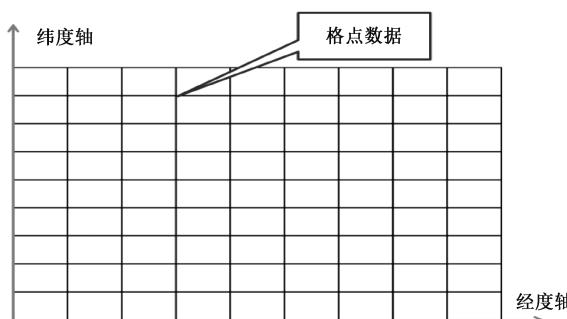


图 3 格点气象产品数据平面示意

取数据进行解压后, 获取相关格点的数据, 这样就大大减少数据量的传输, 提升数据处理性能。

1.4 分布式数据库结构设计

在分布式列存储数据库中, 数据采用列的方式进行存储^[11], 每一列对应着关系型数据库中的一个字段, 如何设计一个合理的主键是关键的问题, 下面分别从主键设计及列设计 2 个方面进行该系统的数据

库结构设计。

分布式列存储数据库通常仅支持主键进行数据的快速检索, 因此需要结合气象数据的检索查询特点, 以及不同格式、不同模式的气象数据特点进行相应的主键设计。以下以欧洲细网格模式数据为例阐述如何进行主键的设计^[12]。

气象数值预报模式产品数据一般包括产品要素名、产品制作中心名、产品制作方法名、资料时间、预报层次、预报时效 6 个维度的信息, 而且这些维度都可能会参与到气象数据的检索工作中去。因此, 对该模式数据进行如下的主键设计^[13]: ①主键由 6 个部分组成, 分别为产品要素名, 产品制作中心名, 产品制作方法名, 资料时间, 预报层次, 预报时效。②主键由 6 个部分组成, 并规定每部分的长度为 8 Byte, 对于不足 8 位的, 采取往前补 0, 形成该部分的主键, 整体的主键结构如图 4 所示。

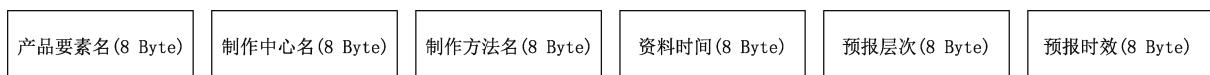


图 4 数据库主键结构

2 关键功能实现

高频精细化气象数值预报产品数据文件是由若干如图 3 所示的平面数据构成的组合体, 各个平面数据按特定顺序排列进而形成 Grib 格式的数据文

件。从技术原理分析可知, 将上述文件进行数据解码效率方面存在两方面的制约因素: ①要查找所需提取特定维度数据所在的位置, 即数据块的定位; ②提取后对内容的写操作。传统方法基本都是通过特定维度查找定位的方法, 产生了一定的时间损耗。

本系统重点从这两个方面着手提出了包括多通道动态感知触发技术、数据块快速处理技术、按需数据截取技术等多角度的创新设计,从而达到对数据快速解码处理的业务目的。

2.1 多通道动态感知触发

多通道动态感知触发在逻辑上由若干的通道单元组成,每个通道逻辑上包括通道源、触发逻辑单元、数据处理单元、Model(通道模型)配置单元4个部分,具有独立、解耦、容错、可插拔等特点。通道源是具体的文件接入路径,是程序的起始触发点,触发逻辑单元基于监听算法机制对触发事件进行实时监听,触发事件包括接入新的数值预报产品处理业务

和各通道对应的文件接入目录有新文件到达。具体处理机制是:新增产品需要处理则在系统 Model 配置文件中新增一个配置项,当系统动态检测到有新增配置项时,会通过 Model 配置逻辑单元读取 Model 对象属性配置文件,对通道配置文件中的所有新增通道配置信息进行解析,然后根据 Model 对象属性配置文件描述的通道属性信息建立相应的文件接入目录和加载处理算法,建立一个数据处理通道;当触发逻辑单元检测某通道源有新文件到达时,则触发数据处理单元对数据源进行解码处理,解析处理完毕后完成数据持久化处理^[14],具体实现流程如图 5 所示。

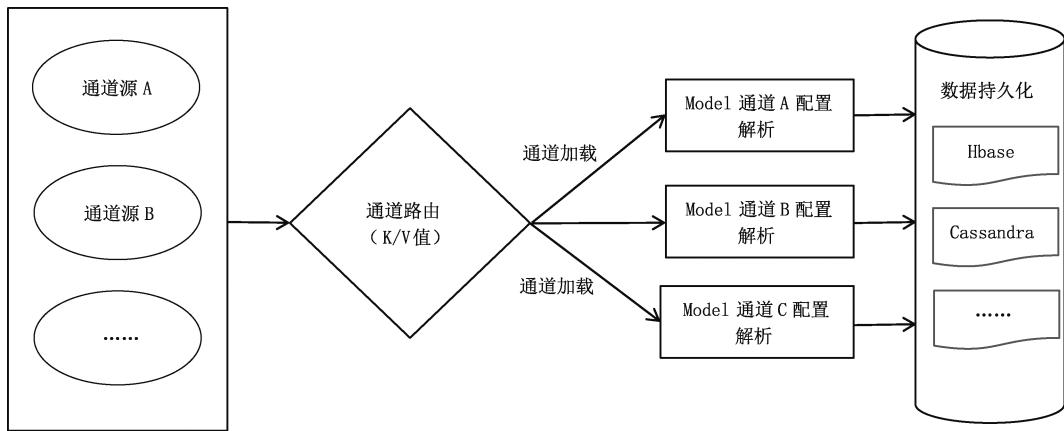


图 5 多通道动态感知技术流程

2.2 实时数据快速处理

每个 Grib 实体文件是若干平面数据 Data[v] 的组合^[15],每个平面数据对应 1 个数据块,文件结构示意如图 6 所示。步骤如下:①获取元信息体。获取整个 Grib 文件的所有数据块的元信息,包括每个“数据块”的起始位置,结束位置等,所获取的信息均

是轻量的元信息,形成元信息体。整个过程的耗时在毫秒级别。②获取实体信息体。过滤掉数据块的元数据信息,并一次性将所有 Grib 数据解码为二进制数据,包括数据的解压缩、解码,一次性将所有 Grib 数据解码为二进制数据,称为实体数据体。实体数据体存储的顺序和步骤①中元数据的顺序是一

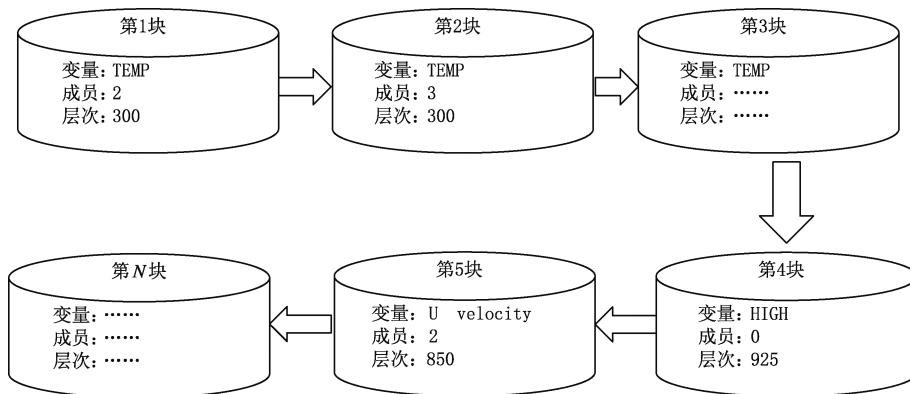


图 6 数据文件结构

一对应的。③目标数据块的位置。根据元数据体和实体数据体的对应关系,首先对元数据体进行整理,通过数据整理函数,将元数据信息进行分类,重点确定数据块的实际位置,位置确定后,根据属性信息计算数据的偏移量进而定位数据并对数据块进行循环读取。

2.3 数据按需实时截取

如果一次性读取整个特定空间范围(具体范围由产品空间属性决定)的全部平面数据,就会造成大量无用数据在业务系统上流转,从而浪费计算、存储资源,影响了处理时效^[16]。为此,设计实现了数据按需截取算法,其核心是能够在特定分辨率下的网格平面数据里面准确定位到需要获取的数据起点和数据的终点,进而从整体的大平面数据中截取出业务需要的数据块。

(1)每圈读取长度的确定。每圈数据的读取长度由经度范围 R_1 和整体的经度差 L_1 确定:

$$R_1 = (L_{\text{onstar}} - L_{\text{onend}}) + L_{\text{onInter}} \quad (1)$$

式中, L_{onstar} 为原始数据开始经度值, L_{onend} 为原始数据结束经度值, L_{onInter} 为原始经度间隔。

$$L_1 = (L_{\text{oriStar}} - L_{\text{cutStar}}) + |L_{\text{oriend}} - L_{\text{cutend}}| \quad (2)$$

式中, L_{oriStar} 为截取前经度起始值, L_{oriend} 为截取前经度起结束值, L_{cutStar} 为截取经度起始值, L_{cutend} 为截取经度结束值。

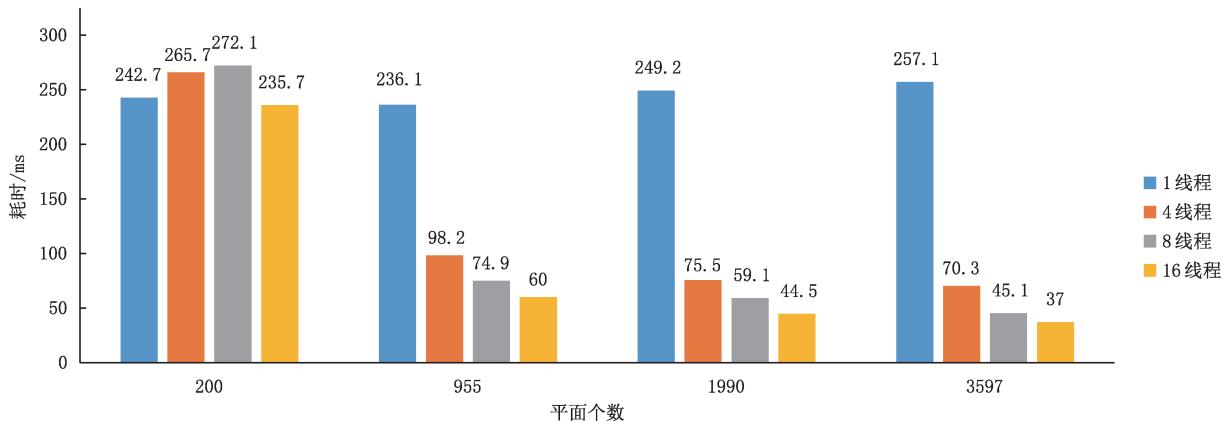


图 7 多线程多平面入库性能测试

从实际运行结果看,单个平面的解码处理时间在 200 ms 左右,采用通常的 8 进程运行解码程序,单个平面解码处理时间可以压缩到 50 ms 左右,以目前广泛应用的欧洲数值预报中心的全球确定性模式产品为例,每个预报时次生成的预报产品共约

计算单圈读取的数据长度 L_{en} :

$$L_{\text{en}} = 4(R_1 - L_1)/L_{\text{onInter}} \quad (3)$$

(2)确定数据截取读取的圈数。数据截取读取的圈数由整个纬度范围 R_2 和整体纬度差 L_2 确定:

$$R_2 = (L_{\text{atStar}} - L_{\text{atEnd}}) + L_{\text{atInter}} \quad (4)$$

式中, L_{atStar} 为原始数据纬度起始值, L_{atEnd} 为原始数据纬度结束值; L_{atInter} 为原始纬度间隔。

$$L_2 = (L_{\text{atiOriStar}} - L_{\text{atiCutStar}}) + |L_{\text{atiOriEnd}} - L_{\text{atiCutEnd}}| \quad (5)$$

式中, $L_{\text{atiOriStar}}$ 为截取前纬度起始值, $L_{\text{atiOriEnd}}$ 为截取前纬度结束值, $L_{\text{atiCutStar}}$ 为截取纬度起始值; $L_{\text{atiCutEnd}}$ 为截取纬度结束值。

计算读取圈数 N_{Cir} :

$$N_{\text{Cir}} = 4(R_2 - L_2)/L_{\text{atInter}} \quad (6)$$

3 性能分析和业务应用

3.1 系统数据处理性能测试

在多进程处理的情况下选取典型的气象数值预报产品,通过接入数量不同的文件对系统的解码入库性能进行测试^[17],即分别从线程数、文件个数,数据条数 3 个维度进行入库性能的测试,测试结果如图 7 所示,图中 4 组柱状图分别表示处理 4 组不同数量的平面数据组。

60 个文件,其中每个文件中约包含 200 个平面数据,实际运行结果显示:系统处理每个预报时次所产生的所有文件的耗时由原来的 30 min 左右压缩到 8 min,在数据处理环节的性能提升明显。

3.2 系统数据服务性能测试

业务系统通常以获取数值预报产品的平面数据为主,其中平面数据获取是指将整个平面数据返回到数据业务请求发起端,实际测试结果如图8所示,图中每3个性能耗时柱状图形成的组合自左到右分别表示数据到达远程业务发起端的处理总耗时、平均耗时和数据在本地服务器处理总耗时。测试平面的数据量在10 MB左右,实际测试结果表明,获取单个平面由现有业务的耗时由900 ms缩减至60~70 ms,批量取单平面平均耗时达30 ms左右^[18]。

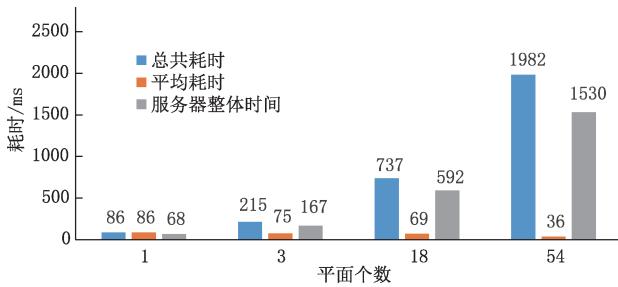


图8 不同平面数据查询效率

3.3 业务应用情况

基于本文设计实现的实时处理系统,实现了国家气象信息中心下发的近10类数值预报产品的快速数据处理,并通过数据服务接口为包括广东省气象行业的市县版格点预报服务系统和可视化图形制作等业务系统提供基础的数值预报产品数据服务。验证了系统的可靠性,具有较好的处理性能,其中数值预报可视化图形制作系统的可视化产品如图9所示。

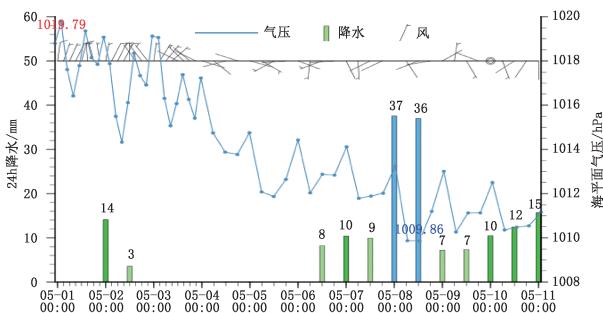


图9 数值预报可视化产品

4 结论与讨论

本文以高频次高精细化的数值预报产品实时快速处理业务应用为切入点^[19],兼顾现有业务现状和

未来的发展趋势,在分析高频精细化数值预报产品数据特点的基础上,提出一种高频海量格点气象数据实时快速处理方法,对高频次精细化数值预报产品的数据预处理、分布式存储和数据服务分布式架构处理进行探索研究,从而建立全新的基于分布式技术的处理架构和存储管理流程,目前系统已经业务化运行并面向全省用户提供数据服务,解决了如何快速高效地处理和应用海量高频精细化气象格点数据这一亟待解决的业务难题^[20],具有一定的业务推广应用价值,同时从实际应用情况反馈看,系统在动态适应原始产品数据格式更改等方面需要加强,未来探索以插件的方式对解码模块优化设计和开发。

参考文献

- [1] 李永生,曾沁,徐美红,等.基于Hadoop的数值预报产品服务平台设计与实现[J].应用气象学报,2015,26(1):122-128.
- [2] 陈正旭,李爽爽,孙晓燕.一种基于NoSQL的气象非结构化数据产品存储方法[J].气象科技,2017,45(3):430-434.
- [3] 徐熙超,杨铮,马廷淮.基于HBase的气象结构化数据查询优化[J].计算机工程与应用,2017,53(9):80-84.
- [4] 王彬,宗翔,魏敏.一个精细粒度实时计算资源管理系统[J].应用气象学报,2008,19(4):507-511.
- [5] 王彬,常飚,朱江,等.气象计算网格平台资源监视模块的设计与实现[J].应用气象学报,2009,20(5):642-648.
- [6] 李永生,曾沁,杨玉红,等.基于大数据技术的气象算法并行化研究[J].计算机技术与发展,2016,26(9):47-49.
- [7] 孔莉莎,吴换萍,刘秋锋,等.气候业务系统运行监控平台设计与实现[J].气象科技,2020,48(3):348-354.
- [8] 曾沁,李永生.基于分布式计算框架的风暴三维追踪方法[J].计算机应用,2017,37(4):941-944.
- [9] 刘媛媛,应显勋,赵芳.GRIB2介绍及解码初探[J].气象科技,2006,34(1):61-63.
- [10] 赵芳,薛蕾,刘媛媛.表格驱动码业务试验系统设计与实现[J].气象科技,2018,46(4):679-684.
- [11] 王甫棣,李湘,姚燕,等.北京GISC系统建模与实现[J].计算机技术与发展,2013,23(5):145-149.
- [12] 肖华东,孙婧,孙朝阳,等.中国气象局S2S数据归档中心设计及关键技术[J].应用气象学报,2017,28(5):633-639.
- [13] 孙周军,乔文文,侯灵,等.混合架构的可视化数据调度检索模型[J].计算机系统应用,2019,28(12):63-71.
- [14] 王兵,李杰.基于通用模型的GRIB格式数据读取技术[J].航空计算技术,2018,48(6):97-101.
- [15] 张瑞聪,任鹏程,房凯,等.Hadoop分布式物联网设备状态分析处理系统[J].计算机系统应用,2019,28(12):79-85.
- [16] 刘媛媛,何文春,王妍,等.气象大数据云平台归档系统设计及实现[J].气象科技,2021,49(5):697-706.
- [17] 王小兰,张雪芬,张婷,等.气象观测设备测试与试验系统设计

- 与实现[J]. 气象科技, 2021, 49(5): 707-715.
- [18] 姚莉, 林建, 李伟, 等. 全球逐时地面气温质量检测方法及应用[J]. 气象科技, 2022, 50(1): 1-8.
- [19] 田刚, 王继竹, 张华林, 等. 长江航运气象预报预警服务系统设计与应用[J]. 气象科技, 2020, 48(4): 503-510.
- [20] 胡洋. 基于深度学习的 SDN 虚拟蜜网路由优化[J]. 计算机系统应用, 2020, 29(10): 274-279.

Design and Implementation of a High-Frequency Fine Meteorological Grid Data Real-Time Processing System

LI Yongsheng¹ LI Gaojie¹ CHEN Yizhi¹ ZHANG Guangyu²

(1 Guangdong Meteorological Data Center, Guangzhou 510640; 2 Guangzhou Central Observatory, Guangzhou 510640)

Abstract: The overall architecture and operation flow of the high-frequency refined meteorological grid data real-time processing system are designed and implemented taking high-frequency massive meteorological grid data as the research object, focusing on the low data processing efficiency of the traditional real-time processing system. Based on the analysis of the characteristics of massive high-frequency meteorological grid data, a distributed storage model is designed and implemented to meet the needs of the meteorological service. Multi-channel dynamic sensing technology is used to implement dynamic multi-channel file processing and fast sensing triggers of file arrival. The fast data block positioning algorithm based on accurate location addressing is implemented using real-time data fast processing technology to realize the accurate positioning of data blocks. The data on-demand real-time interception technology is used to realize the interception algorithm, which can intercept the data on-demand, and then realize the data on-demand extraction. The practical application shows that the system can effectively improve the real-time processing efficiency of unstructured meteorological data.

Keywords: accurate positioning; real-time processing; dynamic perception; data interception