

# 基于科学工作流的通用卫星数据处理调度系统

杨阳<sup>1,2</sup>, 张红梅<sup>1</sup>, 王爽<sup>1</sup>

(1. 中国科学院高能物理研究所, 北京 100049;

2. 中国科学院大学, 北京 100049)

**摘要:**天文卫星产生的海量观测数据下传至地面后,需要经过一系列的处理流程生成可供科学研究人员的数据产品。传统数据处理方式需要耗费大量的时间和人力成本,同时存在缺乏灵活性、任务管理困难及执行效率低下等缺陷。基于上述问题,为实现卫星数据产品生成的自动化和规范化,该文研究并设计了一种基于科学工作流的通用卫星数据处理调度系统。通过封装卫星数据处理算法,结合科学工作流、可视化编辑等相关技术,实现卫星数据处理的自动化和可配置,用户更新需求时无需修改代码,只需重新配置参数,利用可拔插组件设计机制实现系统的通用性,满足不断扩增的卫星的不同数据处理需求。根据已有卫星的数据处理需求进行应用测试,应用实例表明,该系统能够实现卫星数据的自动化处理,并且能够极大地提高处理流程的灵活性、可维护性和执行效率,有利于科学研究人员根据自己的需求获取卫星数据产品。

**关键词:**天文卫星数据处理;数据产品;科学工作流;任务调度;自动化

中图分类号:TP391

文献标识码:A

文章编号:1673-629X(2024)06-0164-07

doi:10.20165/j.cnki.ISSN1673-629X.2024.0075

## A Universal Satellite Data Processing and Scheduling System Based on Scientific Workflow

YANG Yang<sup>1,2</sup>, ZHANG Hong-mei<sup>1</sup>, WANG Shuang<sup>1</sup>

(1. Institute of High Energy Physics, Chinese Academy of Sciences, Beijing 100049, China;

2. University of Chinese Academy of Sciences, Beijing 100049, China)

**Abstract:** After the massive observation data generated by astronomical satellites is transmitted to the ground, it needs to go through a series of processing processes to generate data products that can be used by scientific researchers. Traditional data processing methods require a lot of time and labor costs, and there are shortcomings such as lack of flexibility, difficult task management, and inefficient execution. Based on the above problems, in order to realize the automation and standardization of satellite data product generation, we study and design a general satellite data processing and scheduling system based on scientific workflow technology. By encapsulating satellite data processing algorithms, combined with scientific workflow, visual editing and other related technologies, to realize the automation and configurability of satellite data processing, users do not need to modify the code when updating requirements, only need to reconfigure parameters. The pluggable component design mechanism is used to achieve a universal type of system to meet the different data processing needs of the expanding number of satellites. Application tests are carried out according to the data processing requirements of existing satellites, and the application examples show that the system can realize the automatic processing of satellite data, and can greatly improve the flexibility, maintainability and execution efficiency of the data processing process, which is conducive to scientific researchers to obtain satellite data products according to their own needs.

**Key words:** astronomical satellite data processing; data products; scientific workflow; task scheduling; automation

## 0 引言

在空间天文领域中,卫星数据处理是一个基础的研究方向,面临诸多问题和挑战。卫星通过观测设备

采集并产生了数百 TB 甚至 PB 级别的观测数据,这些数据被传输至地面接收站后,需要经过一系列处理流程转化为标准数据产品,以供科学研究人员使用并进

收稿日期:2023-09-04

修回日期:2024-01-05

基金项目:中国科学院战略性先导科技专项项目(XDA15020500)

作者简介:杨阳(1999-),女,硕士,研究方向为数据处理与数据管理;通讯作者:张红梅(1977-),女,硕士,正高级工程师,研究方向为数据库技术应用。

行后续的数据分析工作<sup>[1]</sup>。一方面,海量的观测数据会为后续的数据处理和分析工作带来巨大的挑战,采用传统数据处理方式需要耗费大量的时间和人力成本;另一方面,卫星的数据处理流程涉及多个领域的科学知识,不同的观测任务可能具有不同的处理流程和要求,传统的卫星数据处理方式缺乏灵活性,科研人员往往受制于系统默认的处理流程和任务调度策略,无法灵活地根据实际需求和优先级进行个性化设置。

随着卫星技术与应用的不断发展和卫星数量的增加,产生的数据量也会大幅增加,这将对数据处理和分析能力提出更高的要求。因此,设计一个具有卫星普适性、能够自定义数据处理流程、能够自动化调度执行的通用卫星数据处理调度系统成为了亟待解决的问题。为解决该问题,该文首先调研了工作流技术和该技术在科学研究领域的相关应用,其次基于科学工作流技术设计并实现了一种通用的卫星数据处理调度系统,最后根据已有卫星的数据处理需求进行应用测试。应用实例表明,该系统能够极大地提高卫星数据处理的灵活性、可维护性和执行效率。

## 1 相关技术与工作

### 1.1 工作流与科学工作流

工作流技术是一种用于组织、管理和自动化业务流程的方法和工具,是指一系列有序的任务、活动或步骤,按照特定的规则和逻辑顺序组织和执行,以完成特定的工作或业务流程<sup>[2]</sup>。通过工作流技术,可以将整个业务流程划分为不同的阶段或步骤,并确保每个步骤按照预定的顺序和规则执行<sup>[3]</sup>。随着工作流技术的飞速发展,工作流思想已经被广泛应用于科学实验中,它可以有效地描述和控制科学实验过程,从而帮助科学家更加轻松地处理和管理科学数据。这种以数据驱动、面向科学实验过程的工作流,被称为科学工作流<sup>[4-5]</sup>。通过科学工作流,科学数据的管理、分析、模拟和仿真等方面得到了有效的支持,从而为科学发现提供了更为便捷有效的环境。

### 1.2 科学工作流技术的应用

目前,工作流技术在科学研究领域的应用已经成为研究热点,主要集中在地球科学、生物信息学、天文学等领域,如 Kepler<sup>[6-7]</sup> 系统、Sunfall<sup>[8]</sup> 系统以及 Ewoks 系统等。

在卫星领域,Nguyen 基于科学工作流技术提出了 SDDS (Satellite Data Downloading System)<sup>[9]</sup>,该系统主要针对以下卫星数据处理阶段:连接数据源,下载原始数据,从原始数据提取二进制(0级)数据处理,处理0级数据并产生加载到数据库中的科学(1级)数据。

然而天文卫星的数据处理分析步骤并不是到此为止,针对不同卫星的科学目标,科学数据可能需要经过不同的处理方式生成不同的2级数据或其他数据,这些处理阶段具有不确定性,灵活性较高,因此需要更通用和具有定制功能的调度系统。

邓玉坤等人提出了一套轻量级科学工作流系统(C-SWF)的解决方案<sup>[10]</sup>,该系统具有任务定制和数据传输等基本功能,但仍存在局限性,不支持流程的可视化自定义以及程序跟踪调试。因此,科学工作流技术在空间天文领域的应用尚不成熟。

## 2 系统设计与实现

### 2.1 系统设计目标

基于科学工作流思想,该文将卫星的数据处理流程按步骤分解为一系列的任务,每个任务可以独立在计算机中调度执行。由于天文卫星数据处理流程具有数据量大、步骤可变化、高重复性、可视化展示需求的特点,该系统需要满足用户自定义数据处理流程、自动化调度执行和实时监控任务运行状态的需求。系统建设目标包括以下几点:

(1)实现可自定义配置的卫星数据处理流程,通过可视化的流程编辑面板,以编排任务执行顺序和填写表单的方式来定义个性化工作流。提高数据处理流程配置的灵活性。

(2)实现卫星数据处理流程的自动化管理和调度执行,根据流程定义和配置自动化分发并执行任务,支持分布式调度和定时调度。同时跟踪任务的执行进度并返回任务的执行状态,若任务失败,支持修改流程配置并重新调度执行。提高数据处理流程的执行效率。

(3)开发可视化的流程监控平台,实现监控卫星数据处理流程,提供查看任务状态、查看日志信息、定时调度、任务启停等功能,提高数据处理流程的可维护性<sup>[11]</sup>。

### 2.2 系统总体设计

根据系统的建设目标,将该系统从用户界面、子系统与功能模块以及数据存储等三个层面进行划分,如图1所示。

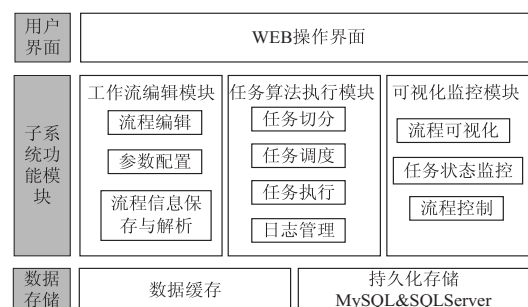


图1 系统架构

在子系统与功能模块层面根据职责划分为三个模块:工作流编辑模块、任务算法执行模块和可视化监控模块。三个模块之间的逻辑视图如图 2 所示。卫星数据处理算法由科研人员提供,并注册为前端页面中的算法组件和后端服务的执行器。每个算法组件和执行器形成一一映射的关系,算法组件以可视化的形式描述工作流中的任务节点,执行器则在计算机中实际执行任务。

注册完成后,用户在工作流编辑页面中编排算法组件,以绘制卫星数据处理流程并形成工作流定义,工

作流编辑模块负责保存和解析定义中包含的组件信息及前后依赖关系等元数据,并将其存储到数据库中。任务算法执行模块负责管理和调度执行工作流,调度器根据工作流定义切分任务节点,并自动分发任务给工作节点,调度对应的执行器执行任务,并将执行过程中的日志信息和结果返回给日志服务并存储到数据库中。可视化监控模块调用接口,实时获取整个工作流的运行状态和各个任务的执行结果,并以可视化方式展示给用户,提供更直观的管理模式。

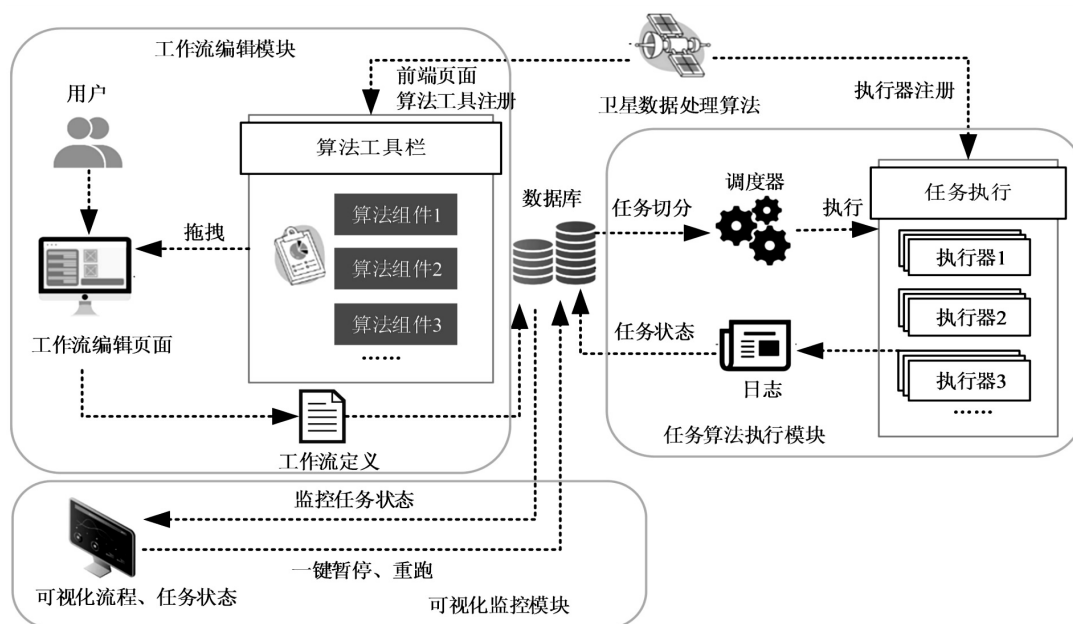


图 2 系统逻辑结构

### 2.2.1 工作流编辑模块

工作流编辑模块是卫星数据处理调度系统的重要模块,负责根据实际需求定制数据处理流程,实现卫星数据处理流程的可自定义配置。

天文卫星数据处理的目标是生成可供科研人员使用和分析的数据产品。通用的卫星数据处理流程如图 3 所示,从卫星传输到地面系统的科学数据经过解包等处理生成可发布的 1 级数据产品,然后通过计算和处理生成其他数据产品,如 2 级数据产品——光变曲

线(Continuous Light Curve, CLC)、能谱(Color Spectral Profile, CSP)等。大部分天文科学家的目标是获取可直接用于分析的科学数据,而不需要精通处理过程中涉及的算法程序,这些算法通常由专业的研究人员开发。然而目前流行的任务调度系统都需要编写代码来创建满足卫星数据处理需求的工作流。因此,该文提供了一种无需编写代码的方式来定义卫星数据处理流程,利用可视化编辑和拖放功能简化流程定义。

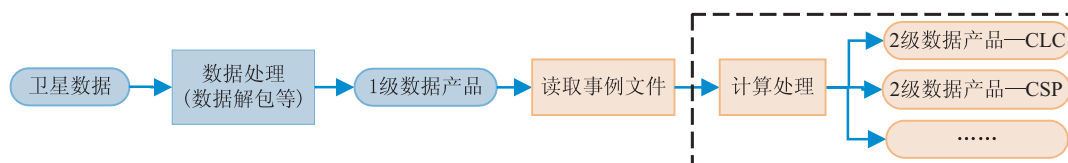


图 3 SVOM 卫星数据处理流程

工作流编辑模块提供可视化的编辑面板、组件库和参数表单。组件库包括算法组件和工具组件,文中算法组件以粗粒度的形式表示图 3 中虚线框部分的处理步骤并作为任务节点,可生成不同类型的 2 级数据产品;工具组件对应数据处理中更细粒度的算法。在

编辑面板中,用户通过编排组件,按执行顺序连接节点,形成有向无环图。每个节点具有唯一标识符并表示要执行的任务,有向边表示任务之间的依赖关系。用户可自定义配置任务参数,如输入输出文件路径和计算过程参数等,用户可通过更改参数进行特殊定制。

按照上述操作最终建立了由卫星 1 级数据产品生成各  
级数据产品的任务处理流程<sup>[12]</sup>。图 4 展示了该模块  
的界面设计原型。

在流程定义过程中,科研人员只需根据组的功能  
进行简单配置,无需关注具体的算法实现。当流程的  
逻辑需求发生变化时,通过重新编排和配置实现流程  
变更。

配置完成后,该模块会保存流程的元数据信息,包  
括工作流名称、描述、依赖关系、修改时间等。这些元  
数据方便用户管理和查找创建的工作流,同时为系统  
的其他模块提供相关信息支持。此外,该模块还会对  
流程定义进行解析,将工作流程分解成若干个任务节  
点,并将节点的组件信息和位置信息存储到数据库中,  
为任务算法执行模块的调度执行提供依据。

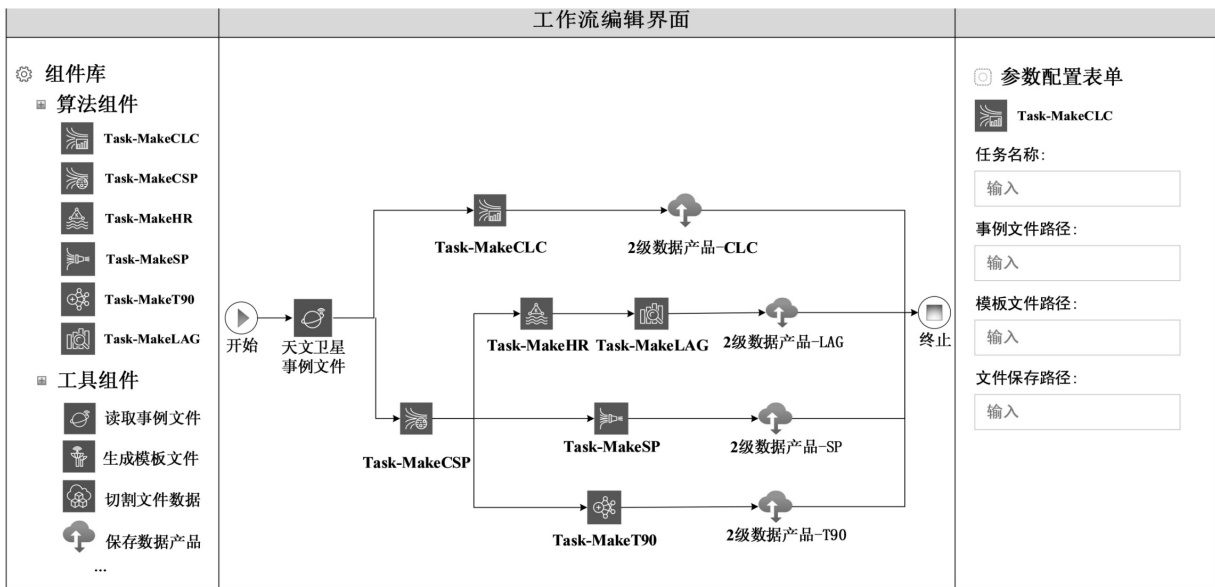


图 4 工作流编辑页面原型

### 2.2.2 任务算法执行模块

任务算法执行模块是卫星数据处理调度系统的核  
心模块,负责自动化管理和调度执行卫星数据处理  
流程。

模块由 API (Application Programming Interface) 接口、  
MasterServer、WorkerServer、LoggerServer 等组件构成。  
用户创建并启动流程后,API 接口处理前端请求并封  
装命令,MasterServer 根据工作流定义切分任务并分发  
任务给 WorkerServer。WorkerServer 是实际的计算节  
点,执行任务并返回执行进度和确认信息给  
MasterServer。任务执行结束后,LoggerServer 保存日  
志信息及执行输出结果。在整个过程中,任务的执行  
状态和日志信息实时反馈,确保数据处理的正确性以  
及发现问题的及时性。这些信息都存储在数据库中,  
并提供 API 接口与第三方系统集成。

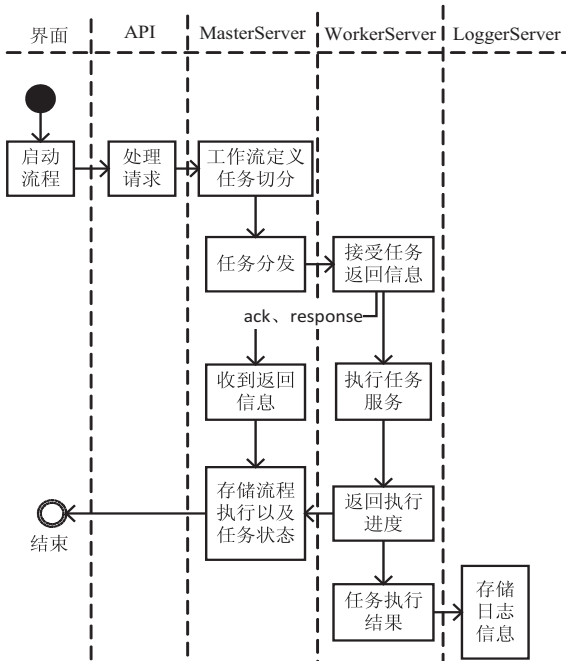


图 5 任务算法执行模块流程

图 5 展示了该模块的组件构成以及工作流程。该

随着卫星数据的数量和处理算法的不断增长,单  
节点计算环境无法满足所有卫星的需求,因此该系统  
需要具有可扩展性,能在集群环境中调度执行并管理  
计算资源。基于此需求,该模块的调度执行主要依托  
开源项目 DolphinScheduler 引擎<sup>[13]</sup>,该引擎不仅支持  
分布式环境,还采用去中心化结构,通过  
Zookeeper<sup>[14]</sup> 服务监控多个 MasterServer 和  
WorkerServer 的运行状态,实现负载均衡,即便单节点  
发生故障,调度系统仍能正常运转。

在具体计算环境中,调度系统将每个任务节点作  
为单独的执行器进行调度。因此,需要将算法组件对  
应的卫星数据处理算法进行服务封装,形成可在计算

机环境中调度执行的执行器,并与编辑模块中的组件库形成一一映射。采用可拔插组件设计机制<sup>[15]</sup>(Service Provider Interface, SPI),可以通过配置文件中的全限定名来加载并实现算法和工具的接口。服务封装的过程可概括以下几个步骤:

(1)创建 Maven 项目,添加 SPI 依赖以及服务的 API 接口依赖;

(2)创建服务的工厂类 TaskChannelFactory,帮助构建 TaskChannel 以及 TaskPlugin 参数,提供该服务的唯一标识,通过 TaskChannel 得到可执行的物理服务 Task,并为当前 Task 添加相应的算法实现;

(3)编译成 Jar 包上传至指定目录,形成可单独调度执行的任务组件。

基于可拔插组件设计机制,接口定义与业务代码实现分离,为该调度系统提供扩展功能,能够根据卫星的实际需求扩展或替换具体算法组件,同时降低系统的耦合度。

### 2.2.3 可视化监控模块

可视化监控模块是天文卫星数据处理流程的监控中心,负责调用 API 接口返回任务执行信息,实现对各个环节的实时监控和可视化展示<sup>[16]</sup>。

在 DolphinScheduler 引擎中,提供了监控中心功能,用于监测 MasterServer 和 WorkerServer 节点服务器的心跳、处理器使用量、内存使用量和平均负载量等信息。然而,该引擎未提供对具体任务组成和运行状态的可视化监控功能,无法直观地查看 workflow 发布状态、workflow 实例运行状态、任务执行状态以及任务的日志信息。因此,本研究设计并实现了可视化监控模块,该模块在 DolphinScheduler 监控中心的基础上进行集成和扩展,它不仅能实时掌握服务器的运行状态,还能对所有 workflow 的执行状态进行可视化监控。

可视化监控模块的页面展示结构分为项目信息饼

图、服务器信息图、workflow 定义列表、workflow 实例列表、任务状态流程图以及任务日志信息等六个区域。项目信息饼图展示当前用户权限下的所有项目;服务器信息图展示 MasterServer 节点和 WorkerServer 节点服务器的信息及状态;选择项目后,workflow 定义列表显示当前项目下的所有 workflow 定义,选择 workflow 定义后,workflow 实例列表区域展示该 workflow 定义下运行的所有实例(流程每运行一次,都会产生一个 workflow 实例),选择 workflow 实例后,任务状态流程图区域展示该 workflow 实例的任务组成、依赖关系和任务运行状态。同时提供重新配置、暂停和重跑操作控件,当 workflow 实例发生故障时,用户可以任务查看更详细的日志信息,并在重新配置后恢复整个 workflow。

通过本模块,用户可以查看卫星数据处理流程的任务执行状态和日志信息,实时了解流程的执行情况。有利于科学研究人员管理和维护卫星数据处理流程。

### 3 系统应用示例

为了提供良好的用户体验,并验证系统可用性,系统集成成了通用的卫星数据处理算法和工具,其中算法包括生成光变曲线(CLC)和计数谱(CSP)等,工具包括生成空数据模板 TOOL\_TEMPLATE 以及数据切割 SPLIT\_DATA。

天基多波段空间变源监视器(Space Variable Objects Monitor, SVOM)是中法两国合力研制的伽玛暴探测科学卫星,主要科学任务是监测和描述伽马射线暴的特征<sup>[17]</sup>。以该卫星的数据处理流程为例进行系统验证。用户登陆系统界面,在项目管理中创建项目和工作流,进入 workflow 编辑页面,按照 SVOM 卫星数据处理流程编排任务节点形成 workflow,测试效果如图 6 所示。

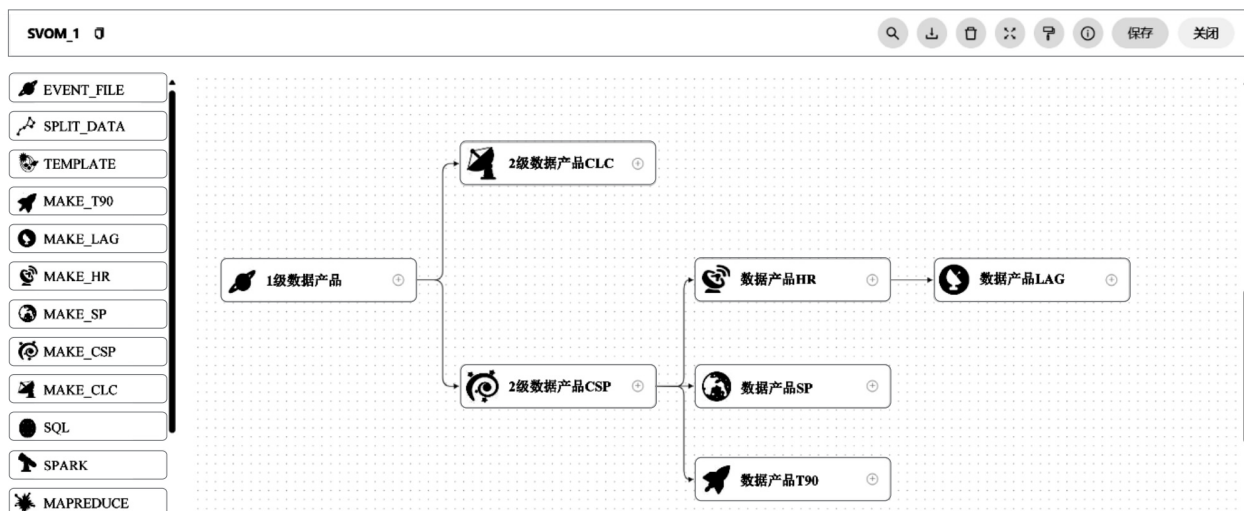


图 6 自定义数据处理流程

其中,CLC 和 CSP 数据产品能够直接由 1 级数据产品生成,HR、SP 及 T90 数据产品的生成依赖于 CSP 数据产品,LAG 数据产品的生成则依赖于 HR 数据产品。对于每个任务节点参数配置,以 2 级数据产品

CLC 节点为例,该节点的输入和输出参数为事例文件路径、CLC 模板文件路径和数据产品保存路径,表单如图 7 所示。

CLC事例文件路径\*

CLC模板文件路径\*

CLC数据产品保存路径\*

前置任务

取消 确定

图 7 参数配置

用户完成工作流编辑和节点配置后,保存工作流并运行。在可视化监控平台中,用户可以查看任务执行状态和日志信息,同时能够控制流程暂停和重跑,测

试效果如图 8 所示。其中任务节点显示绿色表示执行成功,红色表示执行失败,蓝色则表示未执行。



图 8 任务执行监控效果

#### 4 结束语

针对天文卫星传统数据处理方式存在的流程配置灵活性差、任务管理困难及执行效率低等问题,该文设计并实现了一种基于科学工作流技术的通用卫星数据处理流程调度系统。一方面,该系统在传统调度系统

的基础上封装了卫星数据处理算法,实现了流程的自动化调度执行,解决了卫星数据处理的任务管理问题。另一方面,该系统在传统调度系统的基础上,结合可视化技术,为科研人员提供了更友好的交互界面以及更灵活的流程配置方式,有利于科研人员根据实际需求配置卫星数据处理流程。最后,该系统还通

过可视化监控模块,实现了对卫星数据处理流程的实时监控。应用实例表明,该系统能够极大地提高卫星数据处理的灵活性、可维护性和执行效率。

在后续工作中,还需要继续扩展卫星数据处理算法库和工具库,在对未来新增卫星数据处理需求的支持、更“细粒度”的算法服务的支持等方面进行深入研究和改进。

#### 参考文献:

- [1] 郑世界,宋黎明,屈进禄,等. HXMT 卫星的观测规划与数据处理[J]. 现代物理知识,2016,28(4):38-45.
- [2] 罗海滨,范玉顺,吴澄. workflow 技术综述[J]. 软件学报,2000,11(7):899-907.
- [3] 肖飞,张为华,王东辉. 面向科学过程的工作流技术研究现状与趋势[J]. 计算机应用研究,2011,28(11):4013-4019.
- [4] QI L. Workflow management system based on WEB technology[J]. Cluster Computing,2017,20(2):941-947.
- [5] ZHAO Z, BELLOUM A, BUBAK M. Special section on workflow systems and applications in e-Science[J]. Future Generation Computer Systems,2009,25(5):525-527.
- [6] LI X, SONG J, HUANG R. A Kepler scientific workflow to facilitate and standardize marine monitoring sensor parsing and dynamic adaption[C]//International conference on software engineering and service science (ICSESS). Beijing: IEEE,2014:1023-1026.
- [7] ALTINTAS I, BERKLEY C, JAEGER E, et al. Kepler: an extensible system for design and execution of scientific workflows[C]//International conference on scientific and statistical database management. Santorini: IEEE,2004:423-424.
- [8] ARAGON C R, BAILEY S J, POON S, et al. Sunfall: a collaborative visual analytics system for astrophysics [C]//Symposium on visual analytics science and technology (VAST). Sacramento: IEEE,2007:219-220.
- [9] NGUYENM D. A scientific workflow system for satellite data processing with real-time monitoring [J]. arXiv:1812.02236v1,2018.
- [10] 邓玉坤,王锋,邓辉,等. 天文轻量级科学工作流系统的实现[J]. 天文研究与技术,2010,7(4):338-343.
- [11] WEN S T, JIN Z W, TAO F. Design of distributed timing task scheduling system for smart grid[J]. Journal of Physics: Conference Series,2021,2108:012049.
- [12] 胡伟峰,李炜. 一种拖拽组件的方法,装置,终端设备以及存储介质;202210702353[P]. 2023-08-07.
- [13] Apache Dolphin Scheduler. 开源工作流协调平台[CP/DK]. 2019. <https://dolphinscheduler.apache.org/zh-cn>.
- [14] 陈涛,索海燕. Apache ZooKeeper 设计理念和数据结构的研究[J]. 现代计算机,2022,28(21):63-68.
- [15] 唐开山. Java 类加载及 SPI 机制[J]. 电子制作,2020(24):55-57.
- [16] 邱秀连,康倩,王峥. 人物关系的可视化研究[J]. 计算机系统应用,2018,27(4):27-33.
- [17] 余舜京, GONZALEZ F, 魏建彦,等. 中法天文卫星(SVOM)伽玛暴联合探测任务[J]. 空间科学学报,2019,39(6):800-808.