

国家级气象高性能计算机管理与应用 网络平台设计*

宗翔 王彬

(国家气象信息中心,北京 100081)

摘要

该文首先介绍了国家气象信息中心在高性能计算能力建设方面的成就及存在的一些问题,阐述了计算网络的主要思想和技术特点,并具体介绍了国家级气象高性能计算机管理与应用网络平台的架构和采用的主要技术路线和方法。

关键词: 网格;元调度;资源信用审计;业务监控

引言

经过多年的持续投资和大力建设,中国气象局国家气象信息中心目前已建成为计算能力国内第一、气象领域内国际先进的超级计算中心,拥有国内外多个品牌的高性能计算机系统。实现了计算能力跨越式发展后,系统管理如何从单个同构集群层面提升到跨多个异构管理域的多系统全局管理就成了面临的主要问题。

秉承计算网络的无缝资源接入、虚拟性、协同和共享的元计算思路,基于成熟的网格工具软件和开放软件技术,经过半年多的自主研发,已经初步建成了基于网格技术的国家级气象高性能计算管理与应用计算平台,提供了全局的用户管理、资源记账、分配管理、状态监视、作业管理、统计分析、安全认证、共享服务等功能。实践证明,采用成熟的网格、Web等技术,因地制宜自主开发,完全可以实现高性能计算资源集成、可控的高效管理。

1 国家气象信息中心高性能计算资源现状和存在的问题

1.1 能力建设取得的辉煌成就

国家气象信息中心是国内首屈一指的超级计算中心,拥有国内一流、国际先进的高性能计算资源,总

体峰值能力达 23 T FLOPS 以上。其中,2005 年引进并投入到业务运行的 IBM 高性能计算机系统,计算能力达 21.76 T FLOPS,在最新的世界 TOP500 中排第 35 位,是目前国内计算能力最强的计算机系统。

国家气象信息中心还是国内唯一同时以持续不间断业务方式运行向量巨型机、大规模并行机和高性能机群的国家级计算中心,国内外主要高性能计算机厂商的产品均在中心落户安装,并建立了一支有较丰富经验的高性能计算机系统维护、应用开发和场地运行的人才队伍。目前国家气象信息中心正在运行的计算机及运算能力如表 1 所示。

国家气象信息中心从 1999 年开始面向国内(各研究所、大学及大企业等单位)开放共享其计算资源,积累了丰富的资源共享经验。

1.2 系统管理服务的现状和问题

与国外顶尖的超级计算中心相比,国家气象信息中心在系统、用户、资源管理水平、服务质量上都存在着较大的差距。

系统管理急需从单个同构集群层面提升到跨多个异构管理域的多系统整体管理上。目前,国家气象信息中心购置了多家厂商、异构平台的高性能计算机系统,由于各厂商提供的计算机管理工具的局限性,无法了解这些计算机系统整体上的使用情况,也无法进行全局一致的用户管理,存在一定的管理死角;同时也无法科学地对多个计算机系统作为一个整体进行长期使用规划和需求分析。

* 2006-05-12 收到,2006-08-14 收到再改稿。

表 1 国家气象信息中心的主要计算机系统

计算机系统	安装时间	峰值速度	CPU 数	内存/磁盘
CRAY C92	1995-01	2 G FLOPS	2	1 GB/127 GB
IBM SP2	1995-04	8.4 G FLOPS	32	2 GB/156 GB
曙光 1000 A	1998-10	3.2 G FLOPS	9	2 GB/37 GB
神威 I	1999-08	384 G FLOPS	384 (96 节点)	48 GB/1200 GB
IBM SP	1999-11	70 G FLOPS	80 (10 节点)	26 GB/324 GB
YH III	1999-12	18 G FLOPS	17	9 GB/100 GB
神威 新世纪 32P	2003-12	153 G FLOPS	32 (16 节点)	128 GB/4 TB
神威 新世纪 32I	2004-05	166 G FLOPS	32 (16 节点)	128 GB/7 TB
IBM P690/ P655	2004-11	707 G FLOPS	104 (12 节点)	208 GB/15 TB
IBM CLUSTER 1600	2004-12	21.76 T FLOPS	3152 (376 节点)	9 TB/30 TB

资源管理手段粗糙,无法进行精细粒度的资源管理。欠缺系统管理和统计方面的业务管理工具,无法精确地跟踪记录用户占用资源和使用状况的科学统计数据。这些在客观上也造成了目前系统的吞吐率不高和资源调度机制不完善的问题。为此,国家气象信息中心急需研发资源管理的业务软件,加大对计算和磁盘存储资源使用的监视、管理和分配,以清楚地知道用户在高性能计算机系统上使用了多少资源,了解资源是否满足需要,并能够合理地分配资源,实现资源公平、均匀、可控制的使用。

在服务水平上,缺乏与用户的有效沟通和交流,不能充分和及时地了解用户的实际需要和遇到的问题,无法为用户提供深层次的应用支持,IT 管理系统还没有完全建立。

2 计算网络技术的兴起和发展

网络计算技术出现于 20 世纪 90 年代,主要源于某些领域的科学计算对高性能计算能力无止境的需求。计算机适用于解决计算性很强的复杂问题,但计算机在技术和建造上的发展永远赶不上一些实际问题对计算机能力的需求。计算环境不能满足要求常常导致计算机无法解决复杂的实际问题。科学和工程计算的迅猛发展,如生物基因组测算、气象预报和气候预测、核爆炸模拟、高能物理计算等,不断对高性能计算能力和体系结构提出了更高的要求和挑战。尽管更高峰值的超级计算机被不断制造出来, TOP500 列表不断被刷新,但还是无法跟上时代发展的需求。

于是计算机和科学计算领域的研究者想到,既然单个高性能计算机系统提供的计算能力终究是有限的,如果能有一种技术可以把多个(可能地理上分布的、分属不同管理域、具有不同体系架构的)计算

机联合起来,把所有的计算能力聚合到一起,作为一个整体协同工作,为用户提供服务,那么它的计算能力就会超过任何一种单独的超级计算机或者高性能计算机的能力。这个模型就是元计算(meta-computing)模型,而实现元计算的技术就是网格技术^[1]。

以计算资源的共享和协同工作为目标构建的网格计算环境通常称为计算网格。计算资源的类型包括能够直接提供计算能力的高性能计算机、科学计算程序、应用、服务、数据,支持计算的存储器、网络、科学仪器和高端显示设备等。计算资源协同工作和共享是网格计算研究的最初动因,也是网格计算发展的持续推动力^[2]。

3 国家级气象高性能计算机管理与应用网络平台设计

3.1 架构设计

考虑到国家气象信息中心的主要高性能计算资源均位于中国气象局园区网内部,网络连接可靠,属于典型的本地域网格(Local Area Grid)环境。于是平台采用了传统的集中型网格架构,分为网格管理层和 HPC 本地两个层次。把主要的功能集中在网格管理层次上,而由 HPC 本地管理层提供基层支撑。

如图 1 所示,由网格中央节点汇总管理各个加入的高性能计算机系统,并与之交互,负责资源、作业、用户和信息的管理。

如图 2 所示,从整体上看,自上而下,国家级气象高性能计算机管理与应用网络平台可分为用户接口层、网格管理层、HPC(高性能计算机)本地管理层和 HPC 资源层。

平台面向的用户可分为 3 大类:高性能计算资

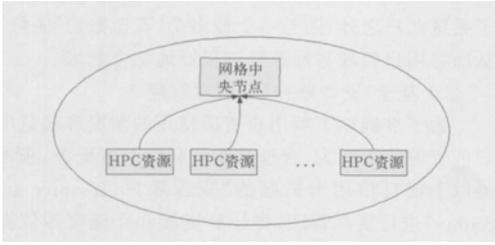


图1 国家级气象高性能计算机管理与应用网络平台的网络架构

源的最终用户、管理人员和决策者。平台的用户接口层位于最上层,直接面向用户,提供两种接口:Web 服务系统和命令行界面(Command Line Interface)。平台的网络管理层与 HPC 本地管理层构成了网格的功能管理层,二者之间除了直接通讯之外,很多的信息交换通过资源信息数据库实现。平台的资源层是国家气象信息中心拥有的各种异构的超级计算机系统,包括 IBM 高性能计算机系统、神威新世纪集群系统、IBM 大 SP 系统等等。

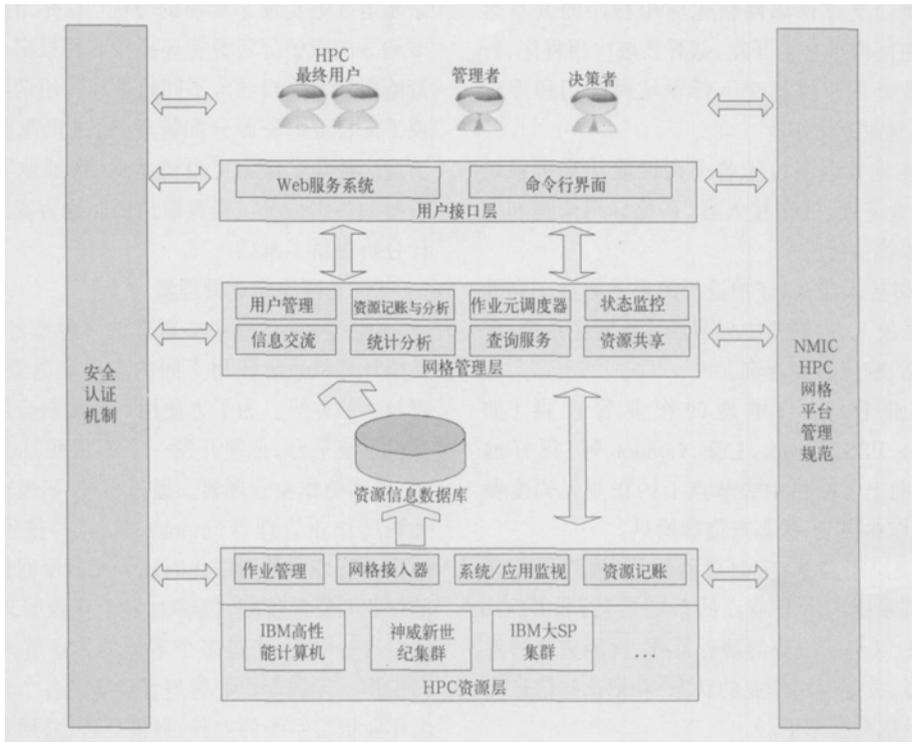


图2 国家级气象高性能计算机管理与应用网络平台的方案设计

3.2 功能模块

用户接口层可分为 Web 和命令行等多种形式。平台的绝大部分功能通过 Web 门户系统提供,支持 J2EE 和网络服务^[3]标准。同时还提供用户和管理人员习惯的 Telnet/SSH 命令行环境,提供基于 Shell, PERL 的命令接口,完成提交作业和系统管理、状态察看等功能。此外,还提供了安全的数据访问服务。

网络管理层位于网络中央节点上,是整个平台的核心,从网格多站点的角度,对国家气象信息中心管理的所有计算资源进行全局的管理和监视。具体可分为用户管理、资源记账与分配、作业元调度、状态监控、信息交流、统计分析、查询服务、资源共享等

功能模块。

① 用户管理:将国家气象信息中心管理的多个高性能计算机系统上的用户管理起来,实现网格环境下的跨集群系统的用户信息和用户一致性管理。

② 资源使用记账和分配管理:精确、实时地进行跨集群的资源使用记账,基于精细化的量化统计,按照一定的分配策略进行公平的资源分配。

③ 作业元调度:用户通过使用统一、友好的网格作业提交工具,就可以向任何一个集群提交作业,而无需与某个集群系统上面的作业提交系统交互。

④ 状态监控:监视网络上各种资源的状态,对各个计算机系统和核心业务/关键进程的状态进行实时的监视。

⑤统计分析:基于对各种资源使用情况的细粒度的记账和监视,提供对不同用户、不同单位、不同项目的周报/月报/季度报告/年报,以友好的图表形式提供。同时,进行分析,对决策者提供决策支持服务。

⑥查询服务:提供用户自定制的接口,可查询各种用户资源使用、分配余额等各种信息。

⑦资源共享:在解决异构系统互联、安全资源分配使用、协同和作业监视等问题的基础上,逐步开放共享资源。支撑保障网格化应用程序的共享运行。支持在网格平台上开发、部署领域应用程序,行业内和行业外用户可提交运行领域内应用程序作业,并可定制业务化作业。

HPC本地管理层位于各个高性能计算机系统上,包括作业管理、网格接入器、系统应用监视和资源使用记账等功能模块:

①网格接入器:以守护进程的形式运行于高性能计算机系统上,向网格中央节点注册所在系统,定时汇报状态,响应各种查询。

②作业管理:与本地的作业管理器(如LoadLeveler, PBS/Torque, LSF, Condor等)很好地集成起来,向上与网格中央节点上的作业元调度器通信,接受作业提交、状态查询和操纵。

③系统应用监视:定时启动运行,检测所在高性能计算机系统内所有节点的系统信息,如CPU、I/O及网络、文件系统等的动态负载,检测重要的系统进程状态,关键应用作业的状态,并把这些信息插入到资源信息数据库里。

④资源使用记账:记录本地高性能计算机系统上的计算资源(主要是CPU时间)和存储资源(主要是磁盘文件系统)的使用情况,动态地插入到资源信息数据库里。

3.3 关键点和解决方案

①全局一致的集中式用户管理

用户管理是整个平台的核心和基础。采用了全局一致的集中式用户管理方案,以适应跨多站点的网格环境。建立了用户全局一致性数据库和关系表,对国家气象信息中心管理的各个高性能计算机系统上的所有系统账户进行整理清理,统一入库。实行用户信息的全局一致性管理,包括统一的用户系统账户/UID、账户名和密码的同步广播、丰富了用户个人信息内容、级联嵌套的用户层次管理。此外,与传统意义上用户管理不同的是,每一个用户除

了系统账户之外,还与一个或多个“资源账户”关联,从而将用户管理与资源管理很好地融合起来。

②基于“资源账户”的资源管理

为了准确地了解用户资源使用的情况和规范用户的资源使用行为,合理分配和科学调度资源,采用类似于银行信用卡机制的“资源账户(resource account)”进行资源管理,将计算资源和存储资源依据它们不同的性能换算成虚拟的资金单元。通过大家熟悉的信用卡操作方法,如存款、取款、查询余额、转账和退款等实现了资源的分配、消费、消费查询、调节和分配撤销等资源管理操作。而根据一定的分配策略向资源账户注入不同数量和使用期限的资金实现了细粒度的资源分配管理,通过资源账户资金的消费次数和金额就可以清楚地、精细地了解了用户资源的使用情况。资源账户机制也为资源使用的统计分析提供了基础。

③元调度作业管理器

国家气象信息中心管理的多种型号、品牌的高性能计算机系统使用不同的本地资源管理器,存在着较大的差异。为了方便用户使用整合后的异构计算机系统平台,必须开发一个不依赖某个计算机系统、透明的作业管理器。通过引进、研发成熟的网格元调度作业管理器(meta-scheduler)技术来解决本地作业管理器的网格化问题。元调度器提供了集中的作业和资源管理,能够对多个作业管理器之间的通信进行协调,使得多个不同的作业管理器一起协同工作。元调度作业管理器屏蔽了各个计算机系统上作业提交系统的差异,提供一致、易友好友好的作业提交工具,允许用户在不了解每种类型的作业管理器的情况下提交作业,使用户感觉不到平台管理机制的差异,实现了整个网格内各个计算机系统的作业提交和管理。能够基于一定的调度策略(例如负载均衡)把提交的作业均匀地分配到具体的高性能计算机系统上。

④实时、友好的综合业务监控

借鉴国外先进气象业务中心运行监视的成功经验,引进可视化的综合业务监控开发工具,开发满足国家气象中心实际需要的综合业务监控系统。综合业务监控系统首先实现了对各计算机系统运行状态实时监控,包括节点运行是否正常、CPU资源使用情况、文件系统利用率监视(通过阈值限制来报警)等。另一方面,综合业务监控系统抓住数值天气预报系统等的业务流程特点和共性,实时监控关键系

统作业(如作业管理系统)和业务作业的运行。

4 实现方案和进展

4.1 技术路线

引进与研发相结合。网格计算技术的飞速发展和国外成熟的高性能计算资源管理技术是国家气象信息中心研发的基础,为解决问题提供了可能。从国家气象信息中心的实际需要和研发力量出发,秉承“引进、消化吸收和再创新”的技术路线,在目前已有的成熟商业化网格平台软件技术和国外相关开源自由软件的基础上自主研发构建一个国家级气象信息高性能计算网格平台。

力求在多个方面做出新意,取得突破,在国内高性能计算中心中领先,具有气象信息处理和业务应用特色。经过业务化试验,最终可切实用于业务生产,作为国家气象信息中心高性能计算机系统的主要管理平台。

4.2 实现情况

通过利用已有的工作基础,引进国外较为成熟的网格平台软件和自由软件,国家气象信息中心已经初步搭建了一个国家级气象信息高性能计算机管理与应用网络平台。

通过高性能计算机系统上部署安装网格平台软件,把计算机系统聚集起来。目前,已在主要的计算机系统上安装了 MOAB 网格平台软件^[4],包括 IBM Thunder 集群、IBM Typhoon 集群、神威新世纪 32I 集群、神威新世纪 32P 集群、神威新世纪 48I 集群、IBM SP 系统等。

在自由软件 Gold^[5]的基础上,开发资源使用记账和分配管理业务软件,目前已经实现了对 IBM Thunder 集群、Tempest 集群、Typhoon 集群的 CPU 时间全面实时不间断计费,正在实验磁盘空间资源的计费。

利用 ECMWF 研发的 SMS^[6]软件,实现对计算机系统和业务系统状态的监视控制和调度。目前已经实现了对 IBM Thunder、Typhoon、Tempest 集群、NSF 服务器和通信系统计算机的系统状态监视,实现了对 T213 等大部分数值模式的作业状态监视和控制。

资源共享方面,在网格平台上已成功开发完成了 MM5 业务模式共享系统。MM5 业务模式共享系统通过提供气象用户定制的 MM5 模式预报产品

(单站、区域图形和 MICAPS 格式),从而共享使用国家气象信息中心的计算资源,提供基于 Web 技术的门户系统和 FTP 数据分发服务的用户接口。目前已经实现了每天两次定时对武汉区域中心气象台、青海省气象台的业务化运行。

5 结 语

对一个国家级的超级计算中心来讲,利用网格技术和思想来构造国家级气象高性能计算机管理与应用平台是一个新的发展方向和挑战,是提升管理和服务水平等软实力的一个重要标志。但由于目前在该领域还没有一些成熟的商业软件和成功的解决方案可以提供借鉴和引用,所以,国家级气象高性能计算机管理与应用网络平台的设计与实现是一个非常好的尝试,特别是将业务化运行的理念和目标贯穿到了架构的总体设计、开源软件的二次开发等。比如,对引进的 MOAB、Gold、SMS 等软件源码都进行了深入的了解和学习,同时,根据国家气象信息中心的实际现状和具体的用户需求,对这些软件都进行了大量的修改和新增功能的开发。实践证明,对国家气象信息中心这样一个业务化程度高的超级计算中心来讲,利用一些较成熟的开源软件,经过消化、吸收和再创新,建立一个拥有部分自主知识产权的国家级气象高性能计算机管理和应用平台系统是比较成功的。

目前该平台还有一些关键的技术急需进一步的开发和实现,比如,基于不同策略的作业元调度系统、服务质量(QoS)管理系统和数据网格等。特别是基于元数据管理的数据网格的实现,才能使目前的国家级气象高性能计算机管理和应用平台发挥更大的作用和效益。

致 谢:感谢所有参与“国家气象网络应用节点建设项目”领导和科技人员的大力支持与辛勤劳动。

参 考 文 献

- [1] Foster I, Kesselman C. The Grid: Blueprint for a New Computing Infrastructure. Morgan Kaufman, 1999.
- [2] Foster I, Kesselman C, Tuecke S. The anatomy of the grid: enabling scalable virtual organizations. *International J Supercomputer Applications*, 2001, 15(3).
- [3] Czajkowski K, Ferguson D F, Foster I, et al. The WS-Resource Framework. <http://www.globus.org/wsrf/specs/ws-wsrf.pdf>, March 5, 2004.

- [4] MOAB Grid Suite . <http://www.clusterresources.com/pages/products/moab-grid-suite.php> .
- [5] GOLD home page . <http://www.emsl.pnl.gov/docs/mscf/gold/> .

- [6] SMS home page . <http://www.ecmwf.int/products/data/software/sms.html> .

Design and Practice of National Meteorological HPC Management and Application Network Platform

Zong Xiang Wang Bin

(High Performance Computing Division , National Meteorological Information Center ,
China Meteorological Administration , Beijing 100081)

Abstract

National Meteorological Information Center (NMIC) of CMA operates the fastest high performance computer (HPC) system in China and the total computing capacity of NMIC also ranks best in China and keeps a leading position among meteorological information centers all over the world . With the great-leap-forward development of capability construction , the “ soft ability ” , characterized by system and resource management , user support and quality of service , is left behind . So efforts must be made on development enhancement and construction in this area , bringing into full play the HPC resources in NMIC . Since its birth , grid technology has seen a rapid growth and been an influential direction in information technologies . Grid brings distributed and heterogeneous computer systems together , works cooperatively as a whole , and provides nontrivial quality of service , which enables the management and sharing of HPC resources . Based on the computational grid concept , a national meteorological HPC management and application network platform is put forward . The platform adopts centralized grid architecture , and consists of four levels , namely , user interface , grid management , HPC local management and HPC resources . The platform finds ideal solutions to four key aspects : globally consistent and centralized user management , resource management based on “ resource accounts ” , meta-scheduler and comprehensive operation monitor . Utilizing existing work , with the introduction of mature grid software and open source software , the platform is preliminarily implemented . In the future , research and development efforts will continue in job scheduling policy , quality of service management and data grid , so as to build and perfect the national meteorological HPC management and application network platform and to put it into actual operation finally .

Key words : grid ; meta-scheduling ; resource credits accounting ; operation monitor