

数据中心成本与利用率现状分析

包云岗 马久跃

中国科学院计算技术研究所

关键词：数据中心利用率 服务质量 PARD

2014年8月8日“百度最高奖”¹揭晓，“百度通用资源管理项目”（代号为Matrix），获得百万美元大奖。Matrix项目由百度数据中心管理团队开发，目标是克服当前数据中心管理所面临的世界级难题。这其实是该团队第二次获奖，2012年他们开发的“服务器潜能激发技术”也获得了“百度最高奖”。可见，百度对数据中心管理技术非常重视，因为正是百度数据中心的几十万台服务器支撑着百度公司每年数百亿的收入。

Matrix项目要解决的是什么样的世界级难题呢？宏观上来讲，就是实现成本与收益之间的平衡，用尽可能低的数据中心成本来实现尽可能大的企业收益。

规模与成本

数据中心已经成为我们生活的一部分，比如我们使用各种手机App上网，就是一个典型的“前

端移动计算+后端数据中心”的场景。我们触按手机屏幕后，手机就会将“触按”这个动作转变为一个网络请求，发送到后端运行云应用的数据中心进行处理，然后再通过网络传回来，在手机屏幕上显示（如图1所示）。这个过程大概需要2~4秒，其中一大半时间消耗在数据中心内。正是数据中心在后台支持着全世界数十亿网民的各种请求。

2013年，时任微软首席执行官的史蒂夫·巴尔默在全球合作

伙伴大会的主题演讲中提到，“我们（微软）的数据中心已超过100万台服务器，谷歌比我们更大一些，亚马逊比我们稍小一点。雅虎和脸谱等其他企业大概是10万台的量级。所以真正理解和管理如此大规模的数据中心，并提供公有云服务的企业是屈指可数的。”而国内的几家互联网巨头也都拥有超过30万台服务器的规模，并且正朝着100万台的规模快速发展。在很多企业看来，数据中心规模属于商业秘密，因



图1 数据中心已成为社会的基础设施

¹ “百度最高奖”是由百度公司首席执行官李彦宏于2010年7月提出的。这是百度公司最高级别的奖项，主要针对公司总监级别以下的对公司作出卓越贡献的基层员工，奖励对象为10个人以下的小团队。

表1 詹姆斯·皮恩根据占地面积、服务器尺寸估算的谷歌数据中心规模

年份	地点	服务器数量
2003	美国佐治亚州道格拉斯郡(Douglas County, Georgia, USA)	417600
2006	美国俄勒冈州达尔斯(The Dalles, Oregon, USA)	204160
2008	美国北卡罗来纳州勒努瓦(Lenoir, North Carolina, USA)	241280
2008	美国南卡罗来纳州蒙克斯科纳(Moncks Comer, South Carolina, USA)	139200
2008	比利时圣吉斯兰(St. Ghislain, Belgium)	250560
2009	美国爱荷华州康瑟尔布拉夫斯(Council Bluffs, Iowa, USA)	296960
2010	芬兰哈米纳港(Hamina, Finland)	116000
2011	美国俄克拉荷马州梅斯县(Mayes County, Oklahoma, USA)	125280
2012	爱尔兰都柏林(Profile Park, Dublin, Ireland)	46400
2013	新加坡裕廊西(Jurong West, Singapore)	200000
2013	香港九龙(Kowloon, Hong Kong)	200000

为很容易根据数据中心规模推测该企业的业务规模以及发展规划,所以很多企业都不愿意透露。于是有很多好奇者采用各种方式去推测超级互联网企业的数据中心规模,例如斯坦福大学教授乔纳森·库梅(Jonathan Koomey)曾根据数据中心的用电量比较准确地推断出谷歌在2010年大约拥有90万台服务器。而詹姆斯·皮恩(James Pearn)根据数据中心的占地面积、每个机架的尺寸、服务器的尺寸等数据,估算了谷歌在全球各个数据中心能容纳的服务器数量(如表1)²。根据皮恩的估算,谷歌的数据中心已经可以容纳230万台服务器(实际数量应该比估算值小)。

100万台服务器的成本是多

少呢?我们可以参考被称为文艺复兴式的黑客——詹姆斯·汉密尔顿在博客中的一些数据。汉密尔顿是负责亚马逊网络服务(Amazon Web Services, AWS)数据中心的杰出工程师,之前在微软负责必应(Bing)的数据中心。2013年《连线》网站上曾经有一篇文章《亚马逊为何聘用一名修车工管理云帝国?》专门介绍了这位开着游艇掌管亚马逊价值数十亿美元数据中心的极客³。

在巴尔默透露微软拥有100万台服务器的数字后,詹姆斯·汉密尔顿写了一篇博客“Counting Servers is Hard”,来估算100万台服务器的成本:如果每台服务器很便宜,只需2000美元的话,那么购买服务器需要20亿美元

(约125亿元人民币);此外还需要配备至少300MW的供电系统,还需要建15~30个机房。这些加起来大概需要22.5亿美元。所以100万台服务器的总建设成本大概是42.5亿美元(约260亿元人民币)。

这些只是建设成本,另外还需要服务器维护和更新成本,包括每年要消耗掉约26亿度电(约15亿元人民币),服务器一般3~5年就会淘汰,需要购置新服务器。所以,就算5年更新一轮,那么仅电费就需要75亿元,服务器更新又需要125亿元。相当于数据中心5年的运行维护成本是200亿元,平均每年40亿元。当然,这个估算是针对比较便宜的服务器,其实一般用于计算的服务器价格都会在2万元以上,所以仅购买服务器成本就会超过200亿元。

利用率现状

这么大规模的数据中心的利用率如何?图2(a)是2006年谷歌的5000台服务器的平均CPU利用率分布,这些服务器运行的是搜索、Gmail等在线应用。从图(a)可以看出,这些服务器的平均CPU利用率约为30%。

图2(b)是2013年1~3月份谷歌的2万台运行在线应用的服

² <https://plus.google.com/+JamesPearn/posts/VaQu9sNxJuY>.

³ 极客是美国俚语“geek”的音译。随着互联网文化的兴起,这个词含有智力超群和努力的语意,用于形容对计算机和网络技术有狂热兴趣并投入大量时间钻研的人。

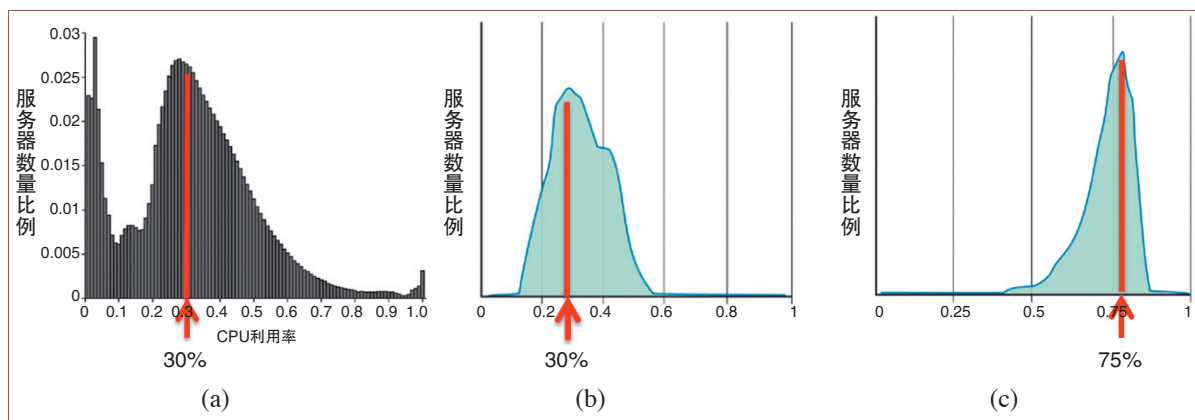


图2 谷歌数据中心CPU利用率：(a)2006年5000台在线应用服务器；(b)2013年2万台在线应用服务器；(c)2013年2万台批处理应用服务器⁴

服务器平均CPU利用率分布图。7年过去了，CPU利用率并没有显著提高，仍然只有30%左右。这意味着，假设100万台服务器中有50万台利用率只有30%，那么相当于5年100亿元人民币的运维成本中有70亿元被浪费掉了，只有30亿元真正产生了效益。

谷歌的数据中心技术是全世界领先的，2006年谷歌就已经拥有了45万台服务器。谷歌在2004年前后就已经开始研制数据中心管理系统Borg，将运行“在线应用”（加引号埋个伏笔）服务器CPU的利用率提高和维持在30%的水平。而更多的企业达不到这个水平，比如麦肯锡估计整个业界服务器的平均利用率大约是6%，而高德纳（Gartner）的估计稍乐观一些，认为大概是12%。这两个数字与笔者了解到的基本吻合，国内一些银行的数

据中心利用率大概是5%左右，而印度塔塔公司曾公布他们的服务器利用率大概是12%。

云计算不是可以通过虚拟化技术在一台服务器上运行多个虚拟机来提高CPU的利用率吗？其实谷歌的Borg系统早就采用了类似的容器技术，2013年的那2万台服务器上就是多个应用混合在一起运行后得到的结果。谷歌将数据中心分为了两类，一类是包含在线应用的数据中心，CPU利用率只有30%；而另一类是不包含在线应用，专门运行MapReduce等批处理的数据中心，CPU利用率平均为75%，如图2(c)所示。

数据中心面临的难题： 利用率vs.服务质量

为什么谷歌在线应用数据中

心的CPU利用率只有30%，而离线批处理数据中心的CPU利用率却可以达到75%？有没有可能把这两类数据中心统一起来，使整体利用率提高到75%？那就能省下一半的服务器，对于100万台服务器规模的数据中心，就可以累计节省上百亿元的采购与运维成本了。

对于这个问题，答案很可能是“基于现有的技术还做不到！”因为数据中心运维不仅要考虑成本，还要考虑收益。

当前流行的虚拟化技术可以让多个应用或虚拟机共享一台机器来提高服务器资源的利用率。但是这种共享会带来资源竞争，进而干扰应用程序的性能，影响在线应用的响应时间。快速的服务响应时间是衡量服务质量的关键指标，是让用户满意、留住用户的关键。

⁴ 图片来源：L. A. Barroso, J. Clidaras and U. Holzle, The datacenter as a computer: An introduction to the design of warehouse-scale machines, 2013。

现任雅虎首席执行官玛丽莎·梅耶尔 (Marissa Mayer) 曾经在谷歌做过一个实验, 把页面上的搜索结果从 10 个增加到 30 个, 希望能让用户一次性浏览到更多的信息。这样搜索结果的返回时间从 0.4s 增加到 0.9s。但是他们发现, 广告收入下降了 20%。梅耶尔对提升在线业务的用户体验总结为一条: 速度为王 (speed wins)。

微软、亚马逊也做过类似的实验。2009 年, 微软在必应搜索引擎上开展实验, 发现当服务响应时间增加到 2000ms 时, 每个用户带给企业的收益下降 4.3%。由于该实验对公司产生了负面影响, 最终不得不被终止。而亚马逊也发现其主页加载时间每增加 100ms 就会导致销售额下降 1%。这对于年营收达到数百亿美元的亚马逊而言, 1% 是很大的损失。

当前数据中心为了保障用户请求的服务质量, 不得不采用过量提供资源的方式, 哪怕是牺牲了资源利用率。资源浪费表现为两种形式, 一种是关键应用独占数据中心。国内大多数企业采用一个数据中心专门运行某个或某几个在线应用, 其他作业运行在其他数据中心上, 以减少对在线业务的干扰。另一种是夸大资源需求。谷歌为了提高服务器利用率, 采用先进的数据中心管理技术 Borg, 将多个应用混合运行在

一台服务器上。他们发现, 这样做, 对于那些离线作业来说是有效的, 因为用户提交查询后只要能给出结果就行, 哪怕慢一点也可以。但是对于那些在线应用开发的程序员 (比如 Gmail 开发人员) 来说, 如果知道自己开发的程序可能会和其他人开发的应用一起运行, 就会在最初夸大资源需求, 但其实际使用的要远少于申请的资源。这种现象在共享环境下很常见。图 3 是推特使用加州大学伯克利分校 AMPLab 开发的 Mesos 数据中心上运行了一个月的情况, 暗红色部分是申请的资源, 绿色为实际使用的资源, 不到 20%, 导致大量计算资源被浪费。

当前数据中心正面临资源利用率与应用服务质量之间的矛盾: 一方面, 在数据中心服务器上同时运行多个应用能有效提高资源利用率, 节省成本; 但另一

方面, 多个应用共享资源相互干扰, 影响应用的服务质量, 降低营收。目前数据中心为了保障营收, 而牺牲了资源利用率, 造成大量成本浪费。

现有解决方案与问题

国内企业之前还没有类似 Borg 这样的技术, 数据中心利用率很低, 有的甚至是个位数。百度的 Matrix 项目, 采用了与谷歌目前最先进的数据中心管理框架相似的设计理念, 朝这个方向迈出了一大步, 使在线数据中心利用率得到了有效提升。据百度内部消息, 随着 Matrix 在百度内部的部署, 2014 年百度节省了超过 5 亿元的成本。而随着百度数据中心的规模扩大到上百万台服务器, Matrix 对于百度的意义将更大。

既要保障在线应用的服务质量, 又要提高数据中心的资源利

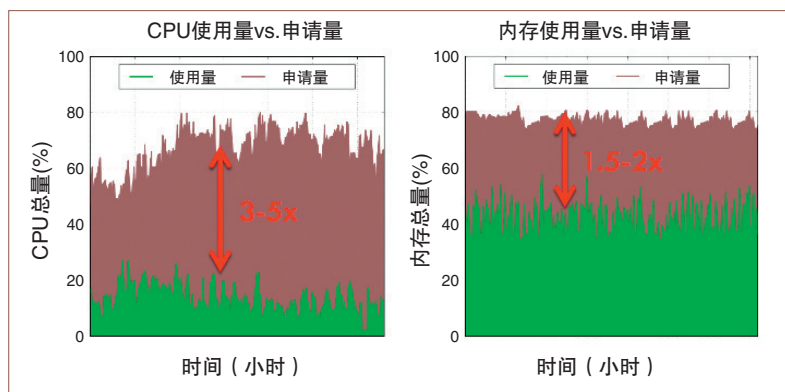


图3 推特使用数据中心管理系统Mesos后的资源使用状况⁵ (暗红色为程序员申请的资源, 绿色为实际使用的资源)

⁵ 图片来源: C. Delimitrou and C. Kozyrakis. Quasar: Resource-efficient and QoS-aware cluster management. ASPLOS, 2014。

用率，这是一个世界级难题，即使是数据中心技术最领先的谷歌也只能做到30%的利用率。学术界也在探索相关问题，加州大学伯克利分校开发 Mesos 就是一个例子。但这些技术都无法解决根本问题——底层硬件上的相互干扰。虽然现在服务器内部的 CPU 核心越来越多，理论上可以同时运行更多的应用，但是其他资源（比如 CPU 内部的高速缓存、内存带宽等）还处于一种“无管理的共享”状态。

以城市交通管理为例来理解当前服务器硬件层次“无管理的共享”状态。传统的数据中心服务器内部就像没有管理的城市道路交通，没有多个车道、没有红绿灯、没有交通规则。这必然会造成大量的冲突和混乱，也会导致一些关键车辆（如救护车、消防车等）的通行无法得到保障。当前数据中心保障在线应用服务质量的手段，其实是一种粗放的交通管制方式——因为管理部门无法区分车辆类别，所以只要得知某些道路上会出现关键车辆，就对其他车辆实施限行。这种方式显然极大地降低了道路的利用率。

由于底层硬件的局限，谷歌不得不耗费大量的人力财力，几乎对整个软件栈做了大手术，从操作系统内核到分布式系统架构进行了改造，率先提出了多项

多应用隔离技术，如 Container、CGroup 等。然而，最近谷歌在 EuroSys 2015 上发表的一篇介绍 Borg 系统的论文中明确指出，虽然使用了 CGroup 等机制来对混合后的应用进行性能隔离，但是对于底层共享缓存、访存带宽的竞争还是会对应用性能产生影响。2015年2月，美国工程院院士、谷歌的资深数据中心专家迪克·赛茨 (Dick Sites) 在威斯康辛大学麦迪逊分校的报告中介绍了数据中心对当前体系结构研究的四点挑战，呼吁亟须从硬件上支持性能隔离的体系结构创新。

对于这个难题，体系结构国际学术界事实上已有认识，2012年发表的学术共同体白皮书 *21st Century Computer Architecture*⁶ 中指出了问题的本质：“指令集 ISA、虚拟内存等传统抽象接口不能传递更多的信息到底层硬件（如应用安全级别、服务质量需求等），导致硬件无法区分不同安全级别或不同服务质量需求的应用，出现应用间相互干扰。因此，我们需要一种新的抽象接口来将程序员的需求和编译器的信息封装并传递给硬件，这会带来效率的极大提高，同时也会带来有价值的新功能。”然而，如何从体系结构层次上支持性能隔离？新的体系结构抽象接口该如何设计？这些仍然是未解的难题。

我们的研究工作—— PARAD

2012年夏天，笔者还在普林斯顿大学，思考过如何从计算机底层的体系结构入手支持资源管理，消除计算机硬件层次上的“无管理的共享”。当时普林斯顿计算机系有几位教授正在开展软件定义网络 (Software Defined Networking, SDN) 方面的研究，也邀请了很多顶级专家作报告，如软件定义网络主要发起人之一、加州大学伯克利分校教授斯科特·申克 (Scott Shenker) 等。平时和朋友也经常聊一些软件定义网络的技术问题。这些交流讨论让我了解到，早在20世纪90年代基于IP的因特网应用兴起时，网络界也曾面临过多应用共享与服务质量问题，因此提出了一系列网络服务质量技术（如区分服务 DiffServ 等）。而近年来崛起的软件定义网络则可以通过标识网络包、增加控制平面、增加可编程机制使网络管理变得灵活方便。

当时笔者产生了一个想法——“计算机内部各个部件之间的通信几乎都是基于包 (packet) 的方式，所以计算机内部其实就是一个小型网络，那是否可以将网络服务质量技术借用到计算机内部呢？（如图4所示）”2012年10月，笔者回到中国科学院计算技术研究所后组建了一个小

⁶ 中文翻译版《21世纪计算机体系结构——计算机体系结构共同体白皮书》，中国计算机学会通讯，第8卷，第12期，2012年12月。

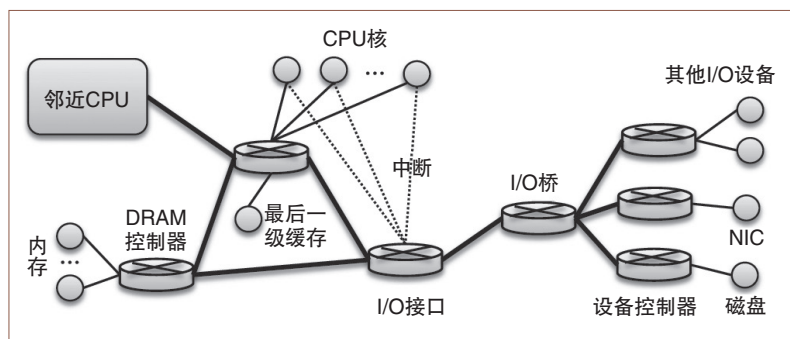


图4 计算机可以看做是一个网络,各部件之间通信已转为基于包的方式

团队,将这个想法付诸实践。我们经过大半年的调研与摸索,不断调整目标与技术路线,在2013年终于有了比较清晰的思路。我们将这个思路凝练为一种新的计算机体系结构——资源按需管理可编程体系结构(Programmable Architecture for Resourcing-on-Demand, PARD)。

这里用城市交通作为例子简单介绍 PARD 体系结构的核心设计理念:(1)将车辆根据不同的用途进行涂装并安装鸣笛,如救护车是白加红十字涂装等(对计算机内部流动的数据包贴上标签);(2)在一些关键路口设置红绿灯,在加油站、维修站等服务点设置管理装置(在计算机内部关键位置增加控制平面);(3)制定交通规则,红绿灯对救护车、消防车等关键车辆可以随时放行,而其他车辆则需要等待绿灯放行,服务点也是优先服务关键车辆(根据不同标签来区分处理数据包);

(4)交通规则可由管理部门根据需要进行调整,比如道路上出现一批武警巡逻车,临时为它们设立一些管理规则(管理员可以调整处理规则)。

事实上,我们日常的交通管理正是采用了上述理念。只要执行严格到位,这样的交通管理系统能在保障救护车等关键车辆顺利通行的前提下提高道路利用率。而 PARD 体系结构也正是希望通过相同的设计理念实现计算机内部高效灵活的资源共享与性能隔离,从而在多种应用混合的数据中心环境下,在保障关键应用的服务质量前提下提高资源利用率。

如果 PARD 技术被验证是可行的,将会有很多新的应用场景。比如对于云计算,可以做到更有效的分级服务管理,类似于航空公司的 VIP 服务,有的人愿意多交钱,享受更稳定的服务质量,甚至是特殊服务,比如硬件提供的加密或压缩服务。当然,这些

还需要更深入的研究。

PARD 在学术界与工业界得到了很好的认可。在学术界, PARD 阶段性研究成果已在体系结构领域顶级会议 ASPLOS 上发表,并得到评审的高度评价,认为 PARD 是在体系结构支持服务质量研究方向上走出了很好的第一步(a good first shot);因为国际学术界对 PARD 研究工作的肯定,笔者也收到邀请参加由加州大学伯克利分校大卫·帕特森(David Patterson)教授、瑞士洛桑联邦理工学院(EPFL)EcoCloud 中心主任巴巴克(Babak Falsafi)教授等组织的在德国举行的为期一周的关于数据中心架构(rack-scale computing)的论坛 Dagstuhl Seminar⁷;在工业界, PARD 项目得到华为公司的大力支持,中国科学院计算技术研究所与华为联合申请了包括多个高价值专利的专利群,并将进一步合作开展 PARD 原型服务器的研制。目前, PARD 全系统软件模拟器验证已完成,并已开源⁸,基于现场可编程门阵列(Field Programmable Gate Array, FPGA)的原型系统将于今年推出, PARD 软件栈也正在设计与实现中。

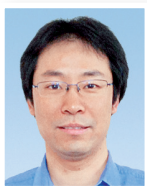
我们相信 PARD 为体系结构与系统软件研究提供了新的视角与平台,十分期待能与国内外同行开展更多的交流合作。■

⁷ <http://www.dagstuhl.de/de/programm/kalender/semhp/?semnr=15421>。

⁸ PARD 模拟器开源地址: <https://github.com/fsg-ict/PARD-gem5>。

致谢：

感谢 PARD 项目组和前沿系统组所有成员：隋秀峰、黄博文、余子濠、李文捷、李宇鹏、靳鑫、杨城、吕文彬、姚治成、任睿、徐天妮、曲宇鹏、展旭升、王聪、李品。感谢中科院计算所研究员孙凝晖与普林斯顿大学教授李凯的讨论与指导。感谢“计算所-华为联合实验室”对本项目的支持。感谢百度数据中心运维团队（特别是刘俊）。本文根据“静·沙龙”在线分享整理修改而成，感谢杨静的支持。



包云岗

CCF高级会员，2013年度CCF-Intel青年学者提升计划获得者，本刊编委、特邀专栏作家。中科院计算所副研究员。主要研究方向为计算机体系结构与计算机系统性能分析。
baoyg@ict.ac.cn



马久跃

中国科学院计算技术研究所博士生。主要研究方向为数据中心体系结构。
majuyue@ncic.ac.cn

CCF TC

2015 CCF大数据学术会议征文

第三届 CCF 大数据学术会议（会议编号：CCF-TC-15-11N）将于 10 月在合肥召开。CCF 大数据学术会议每年举行一次，是为促进大数据技术的研究与发展，推动大数据领域的学术研究与应用而由 CCF 主办的。2015 CCF 大数据学术会议将由中国科学技术大学和安徽大学联合承办。现面向全国征文。

范围要求：<http://dm.ustc.edu.cn/ccf-bigdata2015/submission.html>

提交截止：7 月 1 日

录用通知：8 月 20 日

联系：刘淇 1835 613 7539 qiliuql@ustc.edu.cn

全国高性能计算学术年会征文

全国高性能计算学术年会 (HPC China 2015)（会议编号：CCF-TC-15-32N）将于 11 月 6~8 日在无锡举行。全国高性能计算学术年会是中国高性能计算领域的盛会，是发布最前沿科研成果的平台。现面向全国征文。

范围要求：<http://hpcchina2015.csp.escience.cn/dct/page/1>

提交截止：7 月 15 日

录用通知：8 月 15 日

联系：李希代 010-6260 0662 xidai.niu@gmail.com

第六届全国软件测试学术会议
在南京举行

第六届全国软件测试学术会议（会议编号：CCF-TC-15-16N）4 月 18~20 日在南京召开，会议由 CCF 主办，CCF 容错计算专业委员会和解放军理工大学指挥信息系统学院承办。来自国内知名高校、军队以及地方工业部门的 100 多名代表参加了会议。

会议围绕军用软件及国产化软件的测评技术、软件测试工具的自主创新展开。北京邮电大学教授宫云战、解放军理工大学教授黄松等 8 位专家作了大会报告。与会者就各自的研究成果进行了交流。会议收到投稿论文 79 篇，录用 42 篇，并被推荐到核心期刊发表。