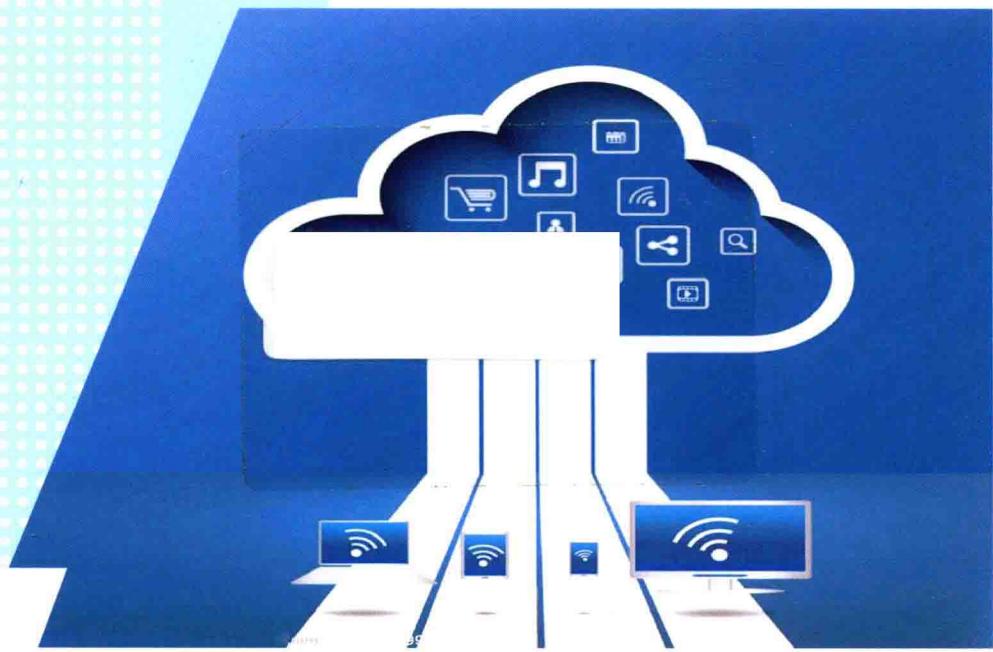


图书馆云的服务 等级协议

黎春兰 邓仲华 陆颖隽 著



 科学出版社

中 菁 卓 纯

图书馆云的服务等级协议

黎春兰 邓仲华 陆颖隽 著

科学出版社

北 京

内 容 简 介

本书分为三篇，包含 9 章内容。第一篇主要从定义、架构和服务过程等方面对云计算和图书馆云的相关研究进行介绍。第二篇先概述了服务质量的相关理论和管理方法，然后讨论了云计算的服务质量所面临的挑战，并寻求相关的质量保障措施。最后利用定量分析方法对图书馆云服务质量的影响因素做了分析。第三篇服务等级协议分别介绍了电信领域的服务等级协议基础理论和管理方法，在此基础上，重点讨论和比较不同云计算和图书馆云的服务等级协议的内容。

本书可供信息管理、图书馆学、云计算等领域教师、研究生阅读参考。

图书在版编目 (CIP) 数据

图书馆云的服务等级协议 /黎春兰, 邓仲华, 陆颖隽著. —北京: 科学出版社, 2016.10

ISBN 978-7-03-050254-4

I . ①图… II . ①黎… ②邓… ③陆… III. ①数字图书馆—图书馆服务—研究 IV. ①G250.76

中国版本图书馆 CIP 数据核字 (2016) 第 254374 号

责任编辑: 任 静 / 责任校对: 桂伟利

责任印制: 张 倩 / 封面设计: 迷底书装

科学出版社出版

北京东黄城根北街 16 号

邮政编码: 100717

<http://www.sciencep.com>

北京京华彩印刷有限公司 印刷

科学出版社发行 各地新华书店经销

*

2016 年 10 月第 一 版 开本: 720×1000 1/16

2017 年 1 月第二次印刷 印张: 15 1/2

字数: 300 000

定价: 85.00 元

(如有印装质量问题, 我社负责调换)

前　　言

国务院在 2015 年 7 月 4 日印发的《关于积极推进“互联网+”行动的指导意见》中指出：积极发挥我国互联网已经形成的比较优势，加快推进“互联网+”发展，有利于重塑创新体系、激发创新活力、培育新兴业态和创新公共服务模式，对打造大众创业、万众创新和增加公共产品、公共服务“双引擎”，主动适应和引领经济发展新常态、形成经济发展新动能、实现中国经济提质增效升级具有重要意义。实际上，“互联网+”是互联网的创新成果与经济社会各领域的深度融合，它催生了以互联网为基础设施和创新要素的经济社会发展新形态，如云计算、物联网、社会计算、大数据等新一代信息技术的新形态。的确，“云计算”是“互联网+计算资源服务化”的新形态，它先将互联网络中的资源（包括网络、服务器、存储、应用软件和服务等）虚拟载入一个可配置的计算资源池中，再以付费服务的方式将这些资源出租给企业或个人使用，是最大化资源利用率的一种商业服务模式，在各行各业中得到了广泛的应用。

图书馆向来注重将新技术应用于读者服务。在云计算的应用上也不例外。云计算向图书馆提供了“互联网+图书馆”的新的应用方式，图书馆通过租用云计算的基础设施服务、平台服务及软件服务，就能使用到云计算强大的计算实例、存储资源、平台及软件等。这种租用的云计算服务，实际上是把图书馆的服务、职能外包给云计算提供商，由云计算提供商代替图书馆承担保存图书馆资源、提供图书馆资源的责任。这样不仅可以使图书馆从其繁杂的基础设施管理和维护的活动中解放出来，还有利于提高工作效率和节约 IT 成本，图书馆从而可以把更多的精力专注于核心业务的创新工作。目前，这种“互联网+图书馆”的租用“云计算服务”的方式已逐步地推广应用，包括租用各种 IaaS、PaaS 和 SaaS。诸如 OhioLINK、DCPL、PITT 图书馆、EKU 图书馆等机构都采用了这种租用云计算服务的方式。“互联网+图书馆”的应用在实践中证明了它具有诸多的优势，如节约成本、提高工作效率等。但图书馆并不仅仅满足于此，它们想在“互联网+图书馆”的基础上，继续推进“云计算+图书馆”的应用，构建基于云计算的专门的图书馆服务管理平台即图书馆云。

“图书馆云”是“互联网+图书馆”的深入应用，体现了“互联网+”的思维，是“云计算+图书馆”的深度融合，是最大化图书馆的资源利用率和服务效率的共享服务平台。OCLC 于 2009 年构建了基于云计算的图书馆管理服务平台 WMS，标志着图书馆云的实现。目前 WMS 拥有 5 个数据中心，在欧洲、南美洲以及澳大利亚已经有超过 200 家图书馆正在使用 WMS 服务。除了 WMS 外，图书馆云服务平台还有 Sierra、Alma、Intota、Open Skies、Global Open Knowledge Base 等。

但是，包括 WMS 在内的图书馆云服务是一项复杂、综合的信息服务，底层 IT 基

基础设施（硬件、软件、系统等）、馆藏资源、各种集成服务及应用程序和工具，都是以组合服务的形式、由一群服务提供商共同提供。在这种复杂的组合服务环境里，一个用户面对多个提供商，服务关系交错复杂，相互间的职责难以理清，使包括资源安全、用户隐私和知识产权等方面的用户服务质量难有保障。为此，Robert Fox 提出使用类似购买软件许可证或签订服务等级协议（SLA）的形式，保证云计算环境图书馆的服务能以特定的价格在特定的时间交付特定的服务质量。SLA 是用于约定通信的服务质量指标和服务双方职责的正式协议，它具有服务质量目标及承诺、服务等级目标及保证的特征。SLA 被应用在图书馆领域后，图书馆云 SLA 也成为明确图书馆和用户双方职责、保证图书馆云服务质量的重要措施。如 WMS 保证每月正常运行时间是 99.8%，Alma 承诺每年正常运行时间至少为 99.5%。

本书尝试使用管理的方法（利用 SLA）来约束并解决图书馆云服务质量的问题。全书分为 3 篇，分别是云篇、服务质量篇和服务等级协议篇。

本书是国家自然科学基金项目（71173163）“云计算环境下图书馆信息服务等级协议研究”、教育部人文社会科学重点研究基地重大项目（11JJD630001）“信息资源云体系及服务模型研究”、教育部人文社会科学研究青年基金项目“图书馆云的服务质量模型”（13YJC870012）的研究成果。

本书在撰写过程中，参考了大量的文献，在此对这些前期的研究者表示感谢。本书的研究与出版得到多方的帮助与支持，特别是武汉大学信息管理学院、广西师范大学，在此深表感谢。同时，还要感谢张文萍、彭丽群、钱文静等的大力支持，他们前期研究的积累和在研究过程中的指导，使本书得以顺利完成。衷心感谢为本书付出辛勤劳动的各位老师和同学们！

由于作者水平有限，加上时间仓促，书中必定有诸多疏漏和错误，敬请广大读者与同行批评指正。

目 录

前言

第一篇 云

第1章 云计算基础	1
1.1 云计算的基本原理	2
1.1.1 云计算的特征分析	2
1.1.2 云计算的架构分析	3
1.1.3 云计算的实例分析	5
1.2 云计算的核心技术	9
1.2.1 虚拟化技术	9
1.2.2 分布式技术	14
1.2.3 浏览器技术	18
1.2.4 云计算与其他计算模式	20
1.3 云计算服务质量的挑战分析	22
1.4 基于SLA的云计算服务质量的保证措施	23
第2章 图书馆云	27
2.1 云计算给图书馆带来的机遇	27
2.1.1 图书馆的“云”思想	27
2.1.2 图书馆对云计算的需求	30
2.1.3 图书馆使用云计算的可能性	31
2.1.4 云计算在图书馆中的应用	32
2.2 图书馆云及图书馆云服务的定义	35
2.2.1 图书馆云的特征	35
2.2.2 图书馆云与图书馆租用云计算的解决方案	36
2.2.3 图书馆云与数字图书馆	36
2.2.4 图书馆云与万维网规模的图书馆	37
2.3 图书馆云的构建——以OCLC WMS为例	37
2.3.1 OCLC WMS简介	37
2.3.2 WMS的构建原则	38
2.3.3 WMS的架构	39

2.3.4 其他图书馆云的应用实例	40
2.4 图书馆云的服务过程	42
2.4.1 交互过程	42
2.4.2 业务关系	43
2.4.3 服务质量的挑战	44

第二篇 服 务 质 量

第 3 章 服务质量概述	47
3.1 QoS 的概念	47
3.2 QoS 的准则	51
3.2.1 制定 QoS 准则的原则	51
3.2.2 QoS 准则的制定	51
3.3 QoS 的参数	53
3.3.1 QoS 准则转换为 QoS 参数的规则	54
3.3.2 QoS 参数的测量	55
3.4 QoS 的管理	57
第 4 章 云计算的服务质量	60
4.1 云计算服务质量的定义	61
4.1.1 云计算 QoS 的四个视角	62
4.1.2 IaaS 服务质量驱动的云计算研究	62
4.2 云计算服务质量面临的挑战	63
4.2.1 障碍一 服务的可用性	64
4.2.2 障碍二 数据锁定	65
4.2.3 障碍三 数据安全与可审计性	65
4.2.4 障碍四 数据传输瓶颈	66
4.2.5 障碍五 性能的不可预测性	66
4.2.6 障碍六 存储的可扩展性	66
4.2.7 障碍七 大型分布式系统存在的缺陷	66
4.2.8 障碍八 快速伸缩	67
4.2.9 障碍九 信誉共享	67
4.2.10 障碍十 软件许可	67
4.3 云计算服务质量保障	68
4.3.1 SLA 保障 QoS 的基础	68
4.3.2 其他保障措施	70

第 5 章 图书馆云的服务质量	72
5.1 图书馆云服务质量的基础	73
5.1.1 信息交流的 SCR 模式	73
5.1.2 服务质量理论	74
5.1.3 SERVQUAL	75
5.1.4 LibQUAL+	76
5.1.5 DigiQUAL+	77
5.1.6 基于 Web 的图书馆服务质量模型	78
5.2 图书馆云服务质量的特征分析	79
5.2.1 焦点小组访谈	79
5.2.2 图书馆云服务质量特征的问卷设计与数据收集	82
5.3 图书馆云服务质量影响因素的因子分析	84
5.3.1 KMO 和 Bartlett 球形检验	84
5.3.2 因子和测度项选取的依据	85
5.3.3 因子分析	85
5.3.4 因子命名与信度检验	87
5.3.5 因子结构的优化	88
5.3.6 因子子维度的探索及其信度检验	90
5.3.7 探索性因子分析的结论	93
5.4 图书馆云服务质量影响因素的结构模型	94
5.4.1 结构模型的概念化	95
5.4.2 操作化定义	97
5.4.3 问卷设计与数据收集	98
5.4.4 模型适配度检验	101
5.4.5 结构模型的检验	103
5.4.6 整体模型的检验	106
5.4.7 本章小结	108

第三篇 服务等级协议

第 6 章 服务等级协议基础	111
6.1 SLA	111
6.1.1 SLA 的需求	112
6.1.2 SLA 的意义	112
6.1.3 SLA 的内容	113
6.2 SLA 的服务	126

6.2.1 服务的层次	126
6.2.2 服务的功能	127
6.3 SLA 的发展	128
第7章 服务等级协议的管理	130
7.1 SLA 管理的价值	131
7.1.1 对 SP 的价值	131
7.1.2 对客户的价值	132
7.1.3 对供应商的价值	133
7.2 SLA 的管理框架	133
7.2.1 eTOM 框架的概念视图	133
7.2.2 eTOM 商务过程框架的 CxO 级视图	135
7.2.3 eTOM 商务过程框架的 Level 2 和 Level 3 级视图	140
7.3 SLA 的生命周期管理	143
7.3.1 产品/服务开发阶段	143
7.3.2 谈判和销售阶段	145
7.3.3 实施阶段	147
7.3.4 执行阶段	148
7.3.5 评估阶段	153
7.3.6 关闭服务阶段	155
7.4 SLA 的参数管理	156
7.4.1 服务角度	157
7.4.2 技术特定参数	158
7.4.3 服务特定参数	158
7.4.4 服务/技术独立参数	158
7.4.5 SLA 参数框架与 SLA 的服务层次的关系	159
7.4.6 SLA 参数框架与 KQI/KPI 的关系	160
7.4.7 服务降级	161
7.5 SLA 的监测	161
7.5.1 QoS 和网络性能	162
7.5.2 网络性能数据的采集	164
7.5.3 数据采集的实现	165
7.5.4 QoS 参数与 NPM 的映射	166
7.6 SLA 的评价	168
7.6.1 SLA 监测的功能分析	168
7.6.2 SLA 报告的功能分析	169

第 8 章 云计算的服务等级协议	171
8.1 云计算 SLA 的内容	172
8.2 云计算 SLA 的需求	173
8.3 云计算 SLA 的管理	176
8.3.1 云计算 SLA 的类型	178
8.3.2 云计算的 SLA 链	179
8.3.3 云计算 SLA 的生命周期	179
8.3.4 SLA 的参数	180
8.4 云计算 SLA 的业务关系模型	181
8.5 云计算 SLA 的定价模型	183
8.5.1 云服务的成本设计	183
8.5.2 云计算服务的计费度量	183
8.5.3 云计算服务的定价模型	189
8.5.4 典型云计算服务价格策略的比较	190
8.6 典型云计算 SLA 案例	201
8.6.1 Google 云服务的 SLA	201
8.6.2 AWS 的 SLA	204
8.6.3 Microsoft Windows Azure SLA	206
8.6.4 典型云计算服务等级协议的比较	208
第 9 章 图书馆云的服务等级协议	211
9.1 图书馆云 SLA 的服务质量描述	211
9.1.1 通用 SLA 模型与服务质量参数	212
9.1.2 图书馆云 SLA 服务等级与服务质量水平	213
9.1.3 图书馆云 SLA 业务关系与服务质量责任	213
9.1.4 图书馆云 SLA 保证用户服务质量的价值	215
9.2 图书馆云 SLA 的组成要素	215
9.2.1 传统图书馆 SLA 的内容	216
9.2.2 图书馆云 SLA 的内容框架	219
9.2.3 图书馆云 SLA 的质量参数	221
9.2.4 图书馆云 SLA 的服务等级	222
9.2.5 图书馆云 SLA 的业务关系	223
9.3 图书馆云 SLA 的应用实例	225
9.3.1 OCLC WMS SLA	226
9.3.2 ExLibris Alma SLA	227
参考文献	229

第一篇 云

第1章 云计算基础

云计算是一个存储在某处的虚拟的、可扩展的资源池，通过网络可按需地向用户提供弹性的服务^[1]。正是这种灵活的、可扩展的、低成本的计算模式，激起了人们的兴奋和好奇。随着2006年Amazon推出EC2（elastic compute cloud）服务，让中小型企业按需（on-demand）购买Amazon数据中心的弹性计算能力，短短几年，云计算就成为IT行业的重要热题，包括各商业机构、科研机构和高校在内的团体组织纷纷投入到云计算的开发和应用研究热潮中来。但与此同时，人们在应用云计算服务的过程中，出现了多次云服务中断的事件，使云计算的服务质量面临巨大挑战。

Nelson^[2]在*Science*上预言，“未来5年内，世界上超过80%的计算和数据存储都发生在云端”。图书馆领域的计算和数据存储也不例外。云计算的应用和发展为图书馆提供了更好的信息服务平台，图书馆的基础设施由云计算提供商负责提供及管理，有利于图书馆把注意力集中在数据资源的自由共享上。云计算服务（包括IaaS、PaaS和SaaS等）在图书馆领域已有了初步的应用。图书馆不仅租用了包括数据存储、计算实例在内的云计算服务，还构建了多个基于云计算的图书馆服务平台，大大地缓解了图书馆的信息管理及成本压力，促进了图书馆的创新服务。

图书馆领域的学者意识到云计算的确为图书馆的信息服务提供了一个重要的环境，带来了重要的发展机遇，他们积极推荐云计算在图书馆信息服务中的应用研究。但同时也出现了许多新问题，服务质量的保证就是其中之一。特别是在云计算的环境下，图书馆对云计算的控制力以及其所在环境具有多方参与的特性，使用户的服务质量难有保障。美国圣母大学（University of Notre Dame）图书馆的高级系统管理员Fox^[3]提出使用类似购买软件许可证或签订服务等级协议的形式，保证云计算环境下图书馆的服务能以特定的价格在特定的时间交付特定的服务质量。

服务等级协议（service level agreement，SLA）是提供商和用户之间为保证服务质量而签署的一份关于服务内容、双方的责任与义务、质量水平与价格等服务细节的协议^[4]。用户与提供商签订SLA，就能获得SLA中规定的服务质量，或在没有获得规定的服务质量时，获得提供商给予的赔偿，使服务质量有保障。本书就基于云计算的图

图书馆服务的新环境对其服务质量的维度及结构进行探索，并根据服务质量维度和结构所反映的特征，用 SLA 来作为保证服务质量的手段，以保证用户使用图书馆云服务的满意度，使图书馆云服务具有竞争优势。

1.1 云计算的基本原理

云计算（cloud computing）是一种将计算资源通过网络交付使用的服务方式^[5]。它采用虚拟化技术将存储在大量分布式计算机上的资源构建成一个虚拟的数据中心，并将数据中心的资源以服务的方式通过网络交付给用户。用户只要连接网络，就可以按需地购买和使用网络中某处的计算资源，包括服务器、存储、应用程序和软件等资源。作者认为，该云计算的定义包含了三个基本概念。

(1) 云计算是互联网云^[6]（Internet cloud）的形象化，用于代表互联网（Internet）或一些大型的互联网环境（networked environment），不再局限在数据和应用程序这些软件资源的互联，更包含了基础设施方面的硬件资源。它描述的是这样的一种场景：客户端的数据和应用程序在某处被存储和获得^[7]。

(2) 云计算是一个数据中心^[8]，是自行维护和管理的虚拟化资源^[6]。用户的数据不再是存储在本地设备上，而是存储在通过互联网链接的远程虚拟数据中心里^[9]。用户的应用程序也不再是运行在个人的终端设备上，而是运行在通过互联网连接的远程大规模、分布式的服务器集群中。服务集群通过网络提供弹性的资源和服务^[10]，用户能随时随地方便地接入到该计算资源共享池并实现按需存取^[11]。

(3) 云计算是 XaaS（一切即服务）的服务模式，一切均以服务的方式交付，包括网络、服务器、存储、应用程序和数据等。它包括以服务的形式通过网络交付的应用程序和提供这些服务所需的数据中心的硬件和系统软件^[12]。

网络是云计算交付使用的工具，数据中心是云计算交付使用的内容，XaaS 是云计算交付使用的形式。提供商通过网络、以 XaaS 的方式将数据中心的资源提供给用户，实现云计算的交付使用，进而实现云计算的价值。用户通过网络、以 XaaS 的方式获取数据中心的资源或服务（称为云计算服务或云服务），实现云计算价值向自身竞争优势的转移。

1.1.1 云计算的特征分析

云计算的价值主要体现在其所具有的特征和优势上。美国标准化技术机构（National Institute of Standards and Technology, NIST）描述了云计算具有 5 大基本特征，即按需自服务（on-demand self-service）、广泛的网络接入（broad network access）、资源池化（resource pooling）、快速的弹性（rapid elasticity）和可计量的服务（measured service）^[11]。其中，按需自服务是指用户可以随时根据自己的需求，通过 GUI 或 API 自主地选择服务、设定服务的性能及服务时间等。广泛的网络接入是指不管客户所使

用的设备是什么，不管客户所处的位置在哪里，只要有连接的网络，用户就可以接入并使用云计算的计算能力。资源池化是指提供商的计算资源被联合起来、以多租户模型（multi-tenant model）的方式、根据客户的需求动态地分配和重新分配不同的物理资源和虚拟资源（包括存储、处理、内存、网络带宽和虚拟机等），同时服务许多客户，使巨大的数据中心具有规模经济^[13]。快速的弹性是指计算能力可以快速地、有弹性地提供，提供商利用多个冗余站点使提供具有更高的可靠性^[11]。对消费者来说，云计算可以提供无限的、可用的计算能力或计算资源，消费者可以根据需要快速地增加或减少资源，可以在任何时间以任何数量的形式购买这些计算能力。可计量的服务是指云计算使用各种计量表来自动控制和优化资源的使用，如存储能力的计量表、处理能力的计量表、带宽计量表、当前用户账户计量表等。通过这些计量表，可以监测、控制、报告资源的使用情况，使服务的使用具有透明性。

云计算的特征使它具备了一些独特的优势。

(1) 成本效益。云计算最突出的特征就是比部署在传统数据中心的方案要更具有成本效益，其中的成本包括硬件、软件、维护系统的人力资源等成本。云计算的即付即用或按需预定的价格模型，使用户只需支付使用的费用，就可以获得云端的各种丰富资源。而且资源存放在云端，由提供商的专门技术人员负责管理和维护，具有更高的安全性和可靠性。用户可以放心地把IT资源的基础设施交给云计算提供商，从而从繁重的IT投资、管理和维护等活动中释放出来，更能专注于自身核心竞争力和业务创新。

(2) 可扩展性。云计算是一个巨大的数据中心。用户只需增加花费就可以快速地扩展计算能力，也可以很容易地缩减计算能力，从而减少花费。实现这些快速扩展的能力，不需要用户放置大量的硬件和软件以等待投入生产，也不需要为了跟上所需求的资源情况而花费数周乃至数月来安装、部署或更新。

(3) 协作创新。存储在云端的资源具有随时随地的可获得性。只需一个简单的操作系统和完整功能的浏览器，用户就可以随时随地通过网络连接数据中心的资源和服务，也可以与他人进行资源的共享、协作，从而提高工作效率、实现业务创新。

1.1.2 云计算的架构分析

对用户来说，云计算就是通过网络获得的资源和服务，所谓的云计算与网络相近，只是这个云计算网络能提供包括基础设施在内的一切资源。就云计算内部的数据中心结构来说，它是由集群计算机使用分布式技术和虚拟化技术向网络中的各用户节点提供包括基础设施、平台和软件等资源的，如图1-1所示。

云计算内部的数据中心可以分为应用层、平台层和基础设施层，分别提供应用程序、平台和基础设施三种服务模式，即SaaS（cloud software as a service）、PaaS（cloud platform as a service）和IaaS（cloud infrastructure as a service），这些服务模式定义了云计算提供服务的职能和性质^[11]。

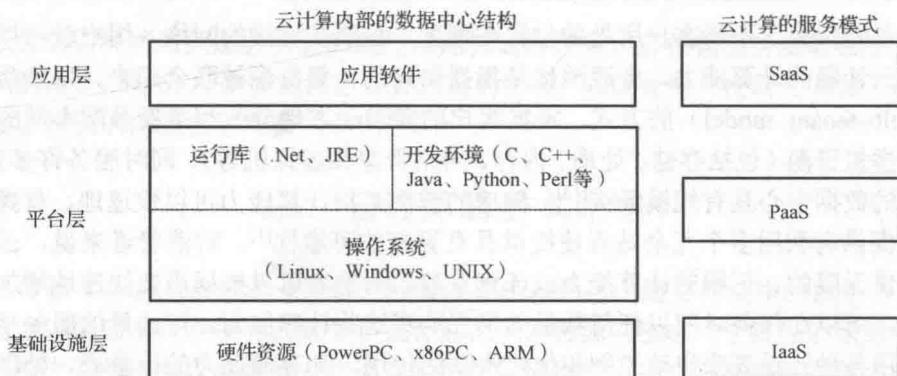


图 1-1 云计算的内部架构及服务模型

(1) 基础设施层是面向基础设施的提供者和管理者的。它们应用虚拟化技术、分布式技术、自动化部署技术等，通过互联网向用户提供 IaaS 的技术方案与服务模式。对它们来说，云计算是通过 IP 网络进行连接的、大规模的、分布式的数据中心基础设施，通过虚拟机按需地向用户提供处理、存储、网络和其他基本的计算资源。云计算的 IaaS 实现了资源动态、弹性地供应，实现了资源的整合与共享，提高了资源的利用率。Amazon EC2 和 S3 是 IaaS 的典型代表。

(2) 平台层是面向程序员和网络应用程序开发者的，是部署在基础设施上的应用程序，通过互联网向开发者提供运行平台的托管服务提供开发 SDK（软件开发工具包），包括应用程序开发、界面开发、数据库开发、存储、测试等。对开发者来说，云计算即 PaaS，是一个互联网规模的软件开发平台和运行环境，他们可在其中开发和部署自己的应用程序。如 Google App Engine (GAE)，Microsoft Azure Platform，Salesforce 的 Force.com 等均提供了 PaaS 服务。

(3) 应用层是面向最终用户的，它通过 Web 平台（如浏览器）将应用程序交付给最终用户。对最终客户来说，云计算就是 SaaS，它通过集中的数据中心向客户按需地提供可扩展的、弹性的应用服务。用户通过浏览器界面即可接入相应的应用程序，如 Google Apps，Salesforce CRM，Microsoft Online 和 IBM Lotus Live 等均属于 SaaS 服务。

目前市场上典型的云计算平台如表 1-1 所示。

表 1-1 典型的云计算平台

属性	Amazon EC2	Google App Engine	Microsoft Azure
专注	基础结构	平台	平台
服务类型	计算、存储 (S3)	Web 应用程序	Web 和非 Web 应用程序
虚拟化	操作系统层次运行于 Xen 管理程序上	应用程序容器	操作系统层次通过结构控制器
动态 QoS 参数	无	无	无
用户界面	Amazon EC2 命令行工具	基于 Web 的管理终端	Microsoft Windows Azure 入口

续表

属性	Amazon EC2	Google App Engine	Microsoft Azure
Web API	有	有	有
附加价值服务提供商	是	否	是
编程框架	基于 Linux 的个性化 Amazon 机器镜像 (AMI)	Python	.NET

根据数据中心的所有权及所面向用户，云计算可以划分为私有云（private cloud）、社区云（community cloud）、公有云（public cloud）和混合云（hybrid cloud）这4种部署模型^[11]。私有云是指云计算数据中心专供某组织内部运作，外部组织无法获取这些资源。它可以存在于组织的内部或外部，可由该组织管理，也可以委托第三方管理。社区云是指云计算数据中心由几个组织共享，支持一个具有共同利害关系（如使命、安全需求、政策和法规因素等）的特定社区，促进社区利用群体智慧进行协调和合作。公有云是指云计算数据中心是由某个组织拥有，并将该数据中心的资源按即付即用（pay-as-you-go）的形式销售给公众或大型行业团体使用。混合云是指云计算数据中心由以上两种或更多种云（私有云、社区云、公有云）构成，保留各种云的一些独特的本质，但是又被标准化或专有技术捆绑在一起，以促进数据和应用程序的移植性。

对用户来说，公有云与私有云的根本区别在于对资源的控制权上。公有云资源的控制权归提供商，私有云资源的控制权归用户或用户所在的单位。一般来说，大型企业或组织，通过使用新的工具和技术，就有能力把组织现有的基础设施改造为一个私有云或混合云。但对中小型企业个体客户来说，租用公有云或社区云提供的资源和服务，在经济和技术上都是比较符合实际的。

1.1.3 云计算实例分析

云计算通过虚拟化技术、分布式技术等将数据中心的资源按需部署，实现不同的服务模式（包括 IaaS, PaaS 和 SaaS）供用户选择。本书通过 IaaS 的计算实例、存储实例、网络传输实例及 PaaS, SaaS 分别分析云计算的应用实例。

1. IaaS

IaaS 利用虚拟化技术提供基本的计算资源，如计算实例、存储实例、网络传输实例等。

1) 计算实例

所谓计算实例就是逻辑上的计算机，也就是运行中的虚拟机。提供商利用虚拟化技术，将集群中的服务器虚拟为多个性能可配置的虚拟机（virtual machines, VM）^[14]，在这台 VM 上预先配置包括 CPU、内存、硬盘和 I/O 总线等的计算资源，然后根据池中资源使用的情况和用户请求资源的情况，灵活地分配和调度资源^[15]。用户租用一台 VM，就具有该 VM 资源的完整访问权限，包括针对此 VM 操作系统的管理员权限。例如，Amazon 的 EC2 为用户、开发人员提供了一个虚拟的集群环境，EC2 中的每一

个实例代表一个运行中的 VM。用户租用的实际上是虚拟的计算能力或性能。EC2 按计算能力划分了 7 种实例类型（包括标准实例、微型实例、高内存实例、高 CPU 实例、集群计算实例、集群 GPU 实例和高 I/O 实例）14 个等级，其中标准实例类型（包含 4 个等级）的计算能力情况如表 1-2 所示^[16]。

表 1-2 Amazon EC2 的标准计算实例

标准实例 (standard instances)	内存/GB	ECU (虚拟内核)	本地存储/GB	平台总线/bit
小型 (默认)	1.7	1(1)	160	32~64
中型	3.75	2(1)	410	32~64
大型	7.5	4(2)	850	64
超大型	15	8(4)	1690	64

其中，1 个 ECU 为 1.0~1.2 GHz 2007 Opteron or 2007 Xeon processor 的 CPU 性能。1 枚单核 CPU，相当于 1 个 ECU 的性能。

EC2 在后三种实例类型（即集群计算实例、集群 GPU 实例和高 I/O 实例）的配置中，还提供 10Gbit/s 的以太网连接。

目前为止，EC2 提供的实例等级最多，性能跨度大，可以满足不同层次的用户的需求。

2) 存储实例

存储实例利用分布式技术，将整个云计算的存储资源进行统一整合管理，为用户提供一个统一的存储空间。具体实现主要是异构传统的存储区域网络 (SAN)、网络附加存储 (NAS) 设备，将分散的存储资源按类型统一集中为一个大容量的存储资源，或将统一的存储资源通过分卷、分目录的权限和资源管理方法进行池化，再将虚拟存储资源分配给各个应用程序或最终用户使用^[15]。存储实例可提供集中存储、分布式扩展、虚拟本地硬盘、安全认证、数据加密、层级管理等功能。如 Amazon EBS (Amazon elastic block store) 专为 EC2 的计算实例提供额外的存储空间。用户可在 EBS 上创建 1GB~1TB 的存储卷到 EC2 计算实例的设备上，也可在同一实例上加载多个存储卷^[17]。Amazon S3 (simple storage service) 是一种面向 Internet 的存储服务^[18]。用户可通过它提供的 Web 服务界面，随时在 Web 上的任何位置存储和检索任意大小的数据。S3 将存储分为 6 个等级，并对每个等级根据每月的数据存储总量 (GB) 提供标准存储和去冗余存储 (reduce redundancy storage) 两种质量类型。前者数据存储的可靠性高达 99.99999999%，后者为 99.99%。Microsoft Windows Azure 也提供了 6 个等级两种冗余的质量类型^[19]。其中，地域冗余是将数据存储在同一区域内的另一个子区域中，以提供最高级别的持久性；本地冗余是在单个子区域中提供持久、可用的存储。当然，不同的等级、不同的质量，其价格是不同的。对相同的质量等级，存储的总量越大，每单位存储的价格就越低。

3) 网络传输实例

网络传输实例主要是配合虚拟机和虚拟存储空间为应用提供网络数据传输服务，具体实现是将一个物理的网络节点虚拟成若干个虚拟的网络设备，如交换机、负载均衡器等，同时进行资源管理。一般来说，为了实现数据的传输，用户在购买计算实例或存储服务时，也需要购买相应的网络传输服务。数据的传输包括数据的传入和传出。Amazon、Google 和 Microsoft 的云计算服务均提供了数据传输服务，并根据传入和传出的数据总量对数据传输收费，如表 1-3 所示。

表 1-3 Amazon (US Virginia 区)、Google (区域 1) 和 Azure (区域 1) 数据传出的价格

	Amazon(EC2/S3)				Google(GCE/GCS)			Microsoft Windows Azure			
	<10	<40	<100	<350	<1	<9	<90	<10	<40	<100	<350
等级/ (TB/月)	0.12	0.09	0.07	0.05	0.12	0.11	0.08	0.12	0.09	0.07	0.05
价格/ (\$ GB/月)											

大部分数据传出的流量总是远高于数据传入的流量的。几乎所有的提供商都免费提供数据传入服务，而对数据传出至不同区域的服务收取不同的费用。为鼓励用户多购买数据传输服务，提供商设计了数据传出的“阶梯价格”，每月传出的数据总量越大，每单位的费率就越低。

云计算的计算实例、存储实例和网络传输实例实现了动态地划分和部署资源，满足了客户的动态需求，降低了系统的复杂度，提高了资源的利用率，并且使用统一的资源池管理，使数据更安全，给用户带来了多方面的好处。用户可以获得应用所需的足够多的计算能力，而且无须对支持这一计算能力的 IT 基础设施付出相应的原始投资成本。用户在需要时可以像购买服务一样购买这种计算能力，按使用付费即可，不用担心计算设备与资源的日常维护开销和闲置成本。

2. PaaS

PaaS 提供的是应用程序的托管平台服务，通常是面向开发人员的。它通常是一个应用程序框架，让开发人员在基础设施上构建并部署 Web 应用程序。如 Google App Engine (GAE) 和 Amazon Elastic MapReduce (Amazon EMR) 是典型的云计算平台服务。Amazon EMR 运行在 Amazon EC2 和 S3 上，按照每机器实例 (machine-instance-hour) 收费。GAE 运行在其基础设施 Google Compute Engine (GCE) 上，按照每进程实例 (process-instance-hour) 收费^[20]。尽管两者同属平台服务，但两者的收费模式不同，这是非常关键的。因为 EMR 可以并行运行几十个进程处理大规模数据，而不需担心费用的问题。但在 GAE 中，即使进程在等待 I/O 传输的过程中，费用仍在增长。这就意味着，同时运行的进程越多，等待的时间就越长，花费也会越高。因此，对 GAE 来说，更少的 CPU 消耗时间，等于更少的花费。这对于需要多线程处理多个 Web 请求的 Python 开发者来说，GAE 每进程实例的收费要比 EMR 每机器实例的收费昂贵^[21]。