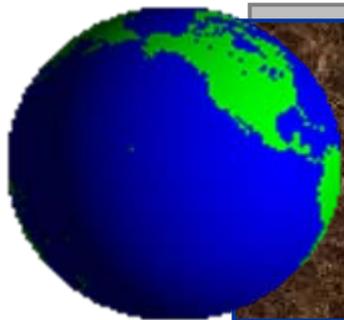


HPC和云计算

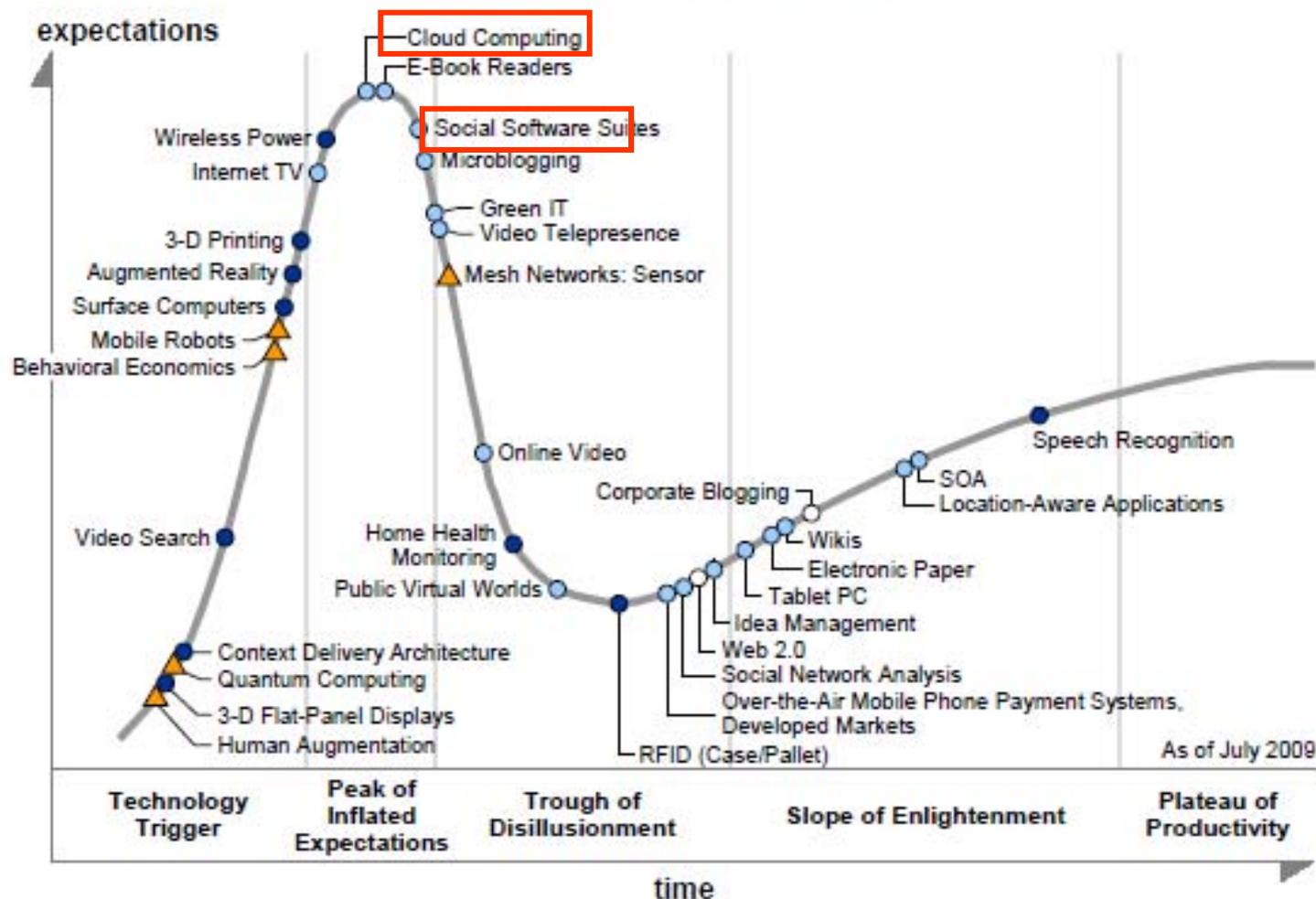
—兼谈应重视计算机系统研究



李国杰，中科院计算所

2010.10.29

Emerging Technologies Hype Cycle 2009



Years to mainstream adoption:

- less than 2 years
- 2 to 5 years
- 5 to 10 years
- ▲ more than 10 years
- ⊗ obsolete before plateau

云计算的市场很大

- 赛迪顾问最新研究显示，国内2010年数据中心与云计算应用市场规模将分别达到839亿元和119亿元，同比增长率分别达到17.8%和29.5%。预计到2012年，中国数据中心和云计算应用市场规模将分别达到**1170亿元和213亿元**。
- IDC预计：未来五年内**亚太地区**在IT云服务上的支出将增长4倍，2014年达到**46亿美元**。整个亚太区2010-2014年的年均复合增长率将达到**40%以上**。
- 2010年10月20日，Gartner公布了**2011年十大战略技术**，**云计算位居首位**。
- 目前全球有**1000多万台**数据中心的服务器还不是以云的方式工作，还有**3000多万台**服务器不在托管中心，发展空间很大。

国内数据中心市场



云计算的优势

- Patterson: 云计算技术和商业模式的优势在于:
 - 1) 给人以可按需立即获得无限计算资源的想象
 - 2) 消除云用户的预先承诺, 有利于企业从小开始
 - 3) 可根据短期的计算资源需要付费, 不需要时就可释放资源.
- 过去几年utility computing 没有成功是由于不能全部满足以上三个条件。例如Intel公司2000-2001年的计算服务需要签订长期服务合同而不是按小时付费。
- 长距离网络带宽成本的下降比其他条件慢很多。云服务提供商可能采取新的办法降低大规模数据传送的成本, 比如可以采用邮寄快递到数据中心（**邮政系统是目前带宽最大的网络系统**）

Berkeley 白皮书对云计算的积极评价

- 云计算很可能像制造商对硬件工业的带来的影响一样影响软件业。TSMC这样的foundry 可以使nVidia这类Fab-less公司在芯片业取得成功，而不需要拥有现代半导体生产线。同样，云计算可以摊销许多“**无数据中心**” (Datacenter-less) 的公司取得成功。
- 如1000台服务器工作1小时的成本与1台服务器1000小时相当。资源的这种弹性使得用户不必为扩展花费多余的成本，这在IT历史上是史无前例的变化。
- 构建并操作一个极大规模、商用低成本的数据中心是云计算的关键，因为在非常大经济规模的情况下，可使电力、网络带宽、操作、软件及硬件等成本**降低5~7倍**。

对云计算的乐观预期

- 云计算将扩大HPC 服务范围。随着虚拟化技术的提高,通信延迟降低, 紧耦合的计算将在更大范围内具有吸引力。
- 多数人认为计算是达到某一目的的工具而本身不是目的。云计算可以使我们熟悉的桌面工具和人机界面潜移默化地改变与扩张, 可使**计算和分析扩展到科研工作者每天工作的环境**。
- Dan Reed: 云计算使得计算和海量数据特别便宜, **云平台最终将取代传统的HPC基础设施**。成功的技术是不可见的技术。
- 如果主机 (mainframe) 是**跳棋**, PC an和 Internet 是**象棋**, 云计算则是要眼光全局的**围棋**。

对云计算的另一种声音

- Oracle CEO Ellison：“云计算有趣的地方在于我们已经重新定义了云计算以将所有我们已经做的东西都包含进来...我不明白我们在云计算下会做什么不一样的事情，除了改变我们的一些广告用词以外。”
- Free软件创始人Richard Stallman：“这样是愚蠢的，比愚蠢还糟糕：这是一个市场炒作活动。有人说这是不可避免的——无论何时你听到有人这么说，很有可能是商业活动使它成真的。”
- Patterson：我不采用“X as a service” (XaaS)”这种术语，X可以是硬基础设施、硬件和平台。我不清楚它们之间的严格的差别在何处。

今天的数据中心与未来的HPC

- 云计算的易用性会影响传统的HPC计算模式。传统的排队批处理方式很难实现按需即时响应的科学计算， On-demand 的云计算给HPC提供了**更易交互的计算模式**。
- 如同是几年前用大众化的PC服务器搭建**集群**以及最近用**GPU**加速科学计算一样，云计算对于HPC也是一次**模式转变**（game changers）。
- 构建百万节点数量级的数据中心与今天构建petascale及今后构建 exascale的系统有许多相同的困难。Dr. Reed认为他们是一对**“双胞胎”**。
- 共同的挑战包括高速互连、存储分层(包括flash, PCM等) 异构多核处理器、系统可靠性和恢复能力、机柜、冷却、能耗效率、和编程等等
- 今天mega-datacenter 的经验将可用于未来的exascale 超级计算机设计。

云计算是超级计算中的新发展

- 对高性能计算（HPC）而言，云计算并不是一个新的概念。事实上，已经发展近30年的超级计算中心也是一种早期的云计算模式：昂贵的计算资源集中部署，多个领域的用户通过互联网远程使用计算服务并依据使用量支付费用。但这种HPC服务和当前所谈论的云计算又有着一些明显的区别，如**没用充分采用虚拟化技术、没有良好的用户界面**等。
- 位于高端计算和桌面计算之间的众多对高性能计算有潜在需求的用户。调研表明，阻碍这些潜在用户使用高性能计算的主要障碍包括：缺乏HPC人才、建设和运维的成本以及使用HPC应用的复杂度。而云计算正是应对这些挑战的最佳途径。
- 云计算将扩大HPC服务的范围。随着虚拟化即时的提高，通信延迟降低，紧耦合的计算将在更大范围内具有吸引力。

传统HPC平台与“HPC云”的区别

	传统架构	云架构
资源管理	作业管理系统：为作业找资源，只管理处理器、应用软件	为用户、作业进行 动态地资源创建和回收 ，管理处理器、内存、存储、网络和应用软件
虚拟化	不支持	服务器虚拟化、存储虚拟化、网络虚拟化
用户管理	独立的用户管理系统，用户无法独享资源	统一用户管理，用户可以独享资源
平台支持	无法修改已安装平台，无法动态修改	可以同时支持多种平台，可以动态修改
数据存储	没有备份机制，不支持异构存储	完善的备份、恢复机制，支持异构存储平台
用户使用	无资源审批流程，无法自定义资源配置	审批、拒绝、预留机制，可以自定义资源平台、软件等

云计算还不适合做尖端的超级计算

- Dan Reed: 云计算绝对不是为特定目的构造的性能顶尖计算机的替代品。如果一种petascale计算需要极低的任务间通信延迟，今天的云计算肯定不适合。但是对于大多数使用较小规模设备的研究者，云计算是有吸引力的替代品。
- 目前的云模型并不支持顶尖的超级计算。动员 grand challenge 应用的人做云计算就如同要说服驾驶第一方程式赛车的深受去乘公共汽车。
- HPC主要执行计算密集型的任务，CPU的利用率已经很高虚拟化技术对提高HPC的CPU利用率作用不大。

目前的云计算做HPC效率较低

- 基于云计算理念来构建超级计算中心，除了满足传统的或现有的HPC用户需求外，更重要的是创造并吸引众多新领域的用户。
- 美国德州先进计算中心（TACC）的 Edward Walker 对 Amazon EC2 上HPC应用的性能表现进行了研究，应用选择常用的基准测试程序NPB，测试结果表明：几乎相同的硬件条件下，对OpenMP版本的8个测试程序EC2性能**下降7%至21%**不等，MPI版本性能则下降**40%至1000%**不等。
- 虚拟化对计算密集型（如果数据能全部放进内存）应用的影响很小，而I/O密集型应用的性能则会有一定下降

在Amazon EC2 上运行MPI性能不高

- Performance is below the level seen at dedicated, supercomputer centers, however, **performance is comparable with low-cost cluster systems.**
- Significant performance deficiency arises from messaging performance where **latencies and bandwidths are between one and two orders of magnitude inferior to big computer center facilities.**

System	latency	uni-bw	bi-bw
LAM	81.20 μ s	57.85MB/s	81.98MB/s
GridMPI	83.46 μ s	54.60MB/s	77.07MB/s
MPICH2 nem	300 μ s	15.72MB/s	26.08MB/s
MPICH2 sock	85.87 μ s	58.49MB/s	83.42MB/s
OpenMPI	300 μ s	16.44MB/s	17.99MB/s

LAM/ACES	35.83 μ s	117.64MB/s	198.59MB/s

---MIT Constantinos Evangelinos and Chris N. Hill CCA-08 paper

核心问题是突破计算机系统技术

目前推广云计算的重点在转变商务模式

- 用户感觉到的云计算的好处主要是减少购买硬件软件的信息化开支，更好地满足动态变化的需求，降低用户端软件升级的维护成本和管理成本。这种好处主要来自**商务模式的转变**，其核心技术是**虚拟化技术**。
- 目前数据中心转向云计算平台的动力主要是**服务器的统计复用**，可以降低数据中心的运行成本，提高服务器和海量存储的**利用率**，其关键技术也是虚拟化技术。
- 虚拟化技术是一种相对门槛较低的技术，因此各大公司和各地政府都可以在较短时间内建立“云计算平台”。
- 实际上真正支持云计算的是**计算机系统技术**，这些技术用户看不见，媒体也很少宣传。
 - 与李开复的会面

用户感觉不到的云计算技术

- 云计算系统的本质可以看成是：
资源虚拟化 + 并行计算
- 云计算不等于虚拟化。虚拟服务器并不能组成一朵云，云计算的能力远远超出一般的虚拟化解决方案。
- 并行技术是藏在云计算背后的核心技术，也是Google等云计算公司具有竞争力的关键技术。
- 各个层次的互连网络在数据中心起到一个非常核心的作用，其作用可能超过服务器本身。
- 虚拟化技术已经开始改变企业对服务器、操作系统以及计算资源的重新部署，甚至导致全新的IT管理模式，这一变化无疑将对操作系统产生重大影响。

云计算的挑战与机遇

	问题	机会
1	服务的可用性	选用多个云计算提供商；利用弹性来防范DDOS攻击
2	数据丢失	标准化的API；使用兼容的软硬件以进行波动计算
3	数据安全性和可审计性	采用加密技术，VLANs和防火墙；跨地域的数据存储
4	数据传输瓶颈	快递硬盘；数据备份/获取；更加低的广域网路由开销；更高带宽的LAN交换机
5	性能不可预知性	改进虚拟机支持；闪存；支持HPC应用的虚拟集群
6	可伸缩的存储	发明可伸缩的存储
7	大规模分布式系统中的错误	发明基于分布式虚拟机的调试工具
8	快速伸缩	基于机器学习的计算自动伸缩；使用快照以节约资源
9	声誉和法律危机	采用特定的服务进行保护
10	软件许可	使用即用即付许可；批量销售

—— 引自Berkeley白皮书

计算机系统技术面临转折性挑战

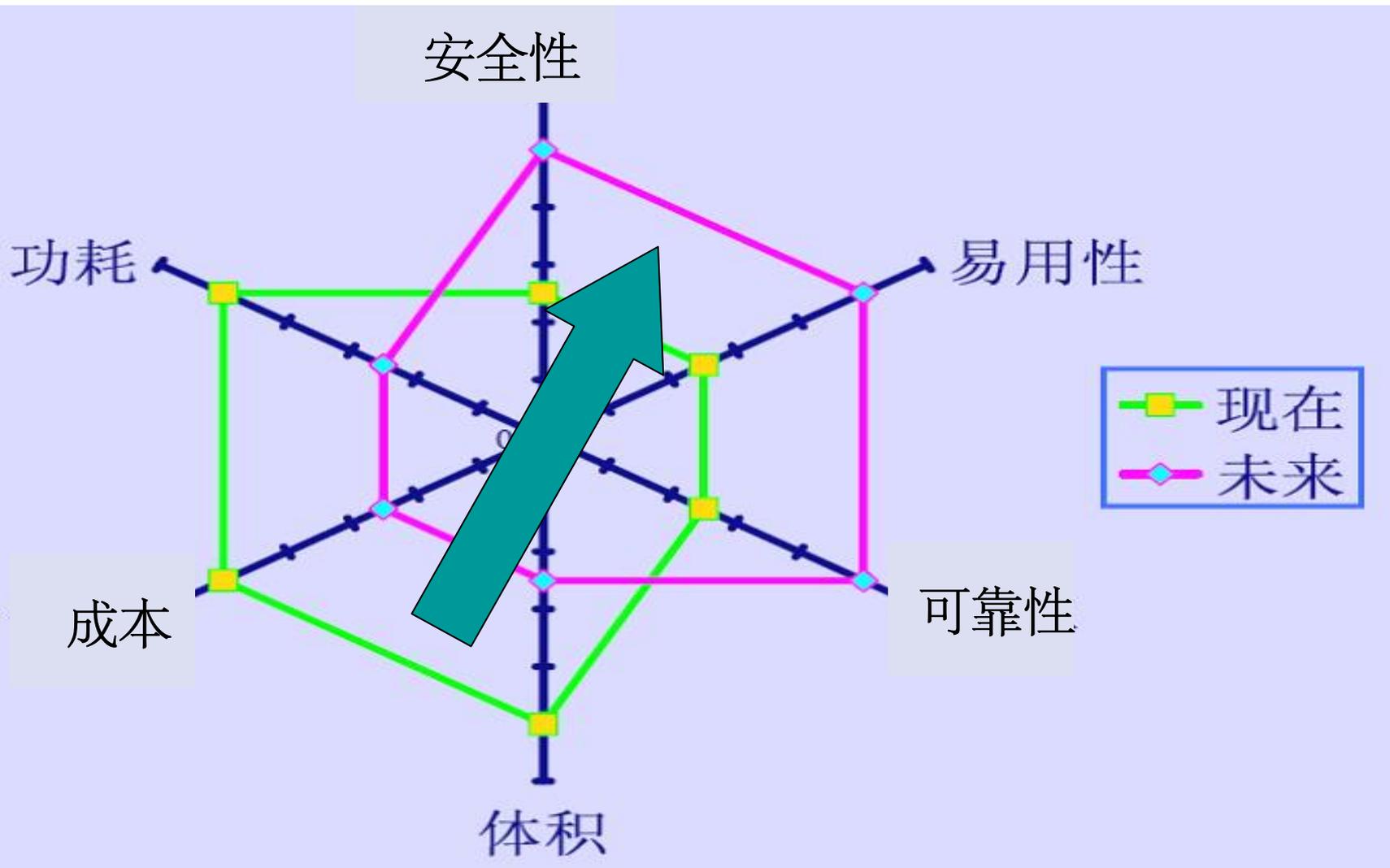
- 计算机系统结构（Architecture）的概念是从上世纪60年代研制IBM360时提出来的，重要的贡献是**区分了硬件与软件，定点与浮点，提出了系列计算机**的概念，这些概念一直沿用至今。
- 现在和未来的云计算的workloads的特征与过去 HPC、事务处理等应用有很大区别，从新的应用中（大量网络服务）应**归纳出新的基本指令集**。
- 今后决定一个云计算平台是否能存活不是光看虚拟化技术，而是看它的**资源利用率，成本和可靠安全**等系统因素。
- 认为计算机系统技术已经差不多了，国家可以不支持了，是一种非常缺乏远见的判断。把对云计算的支持偏重于中间件也是一种片面的政策。**高效可靠的后台设备是云计算的关键**。

计算机系统研究的基本问题

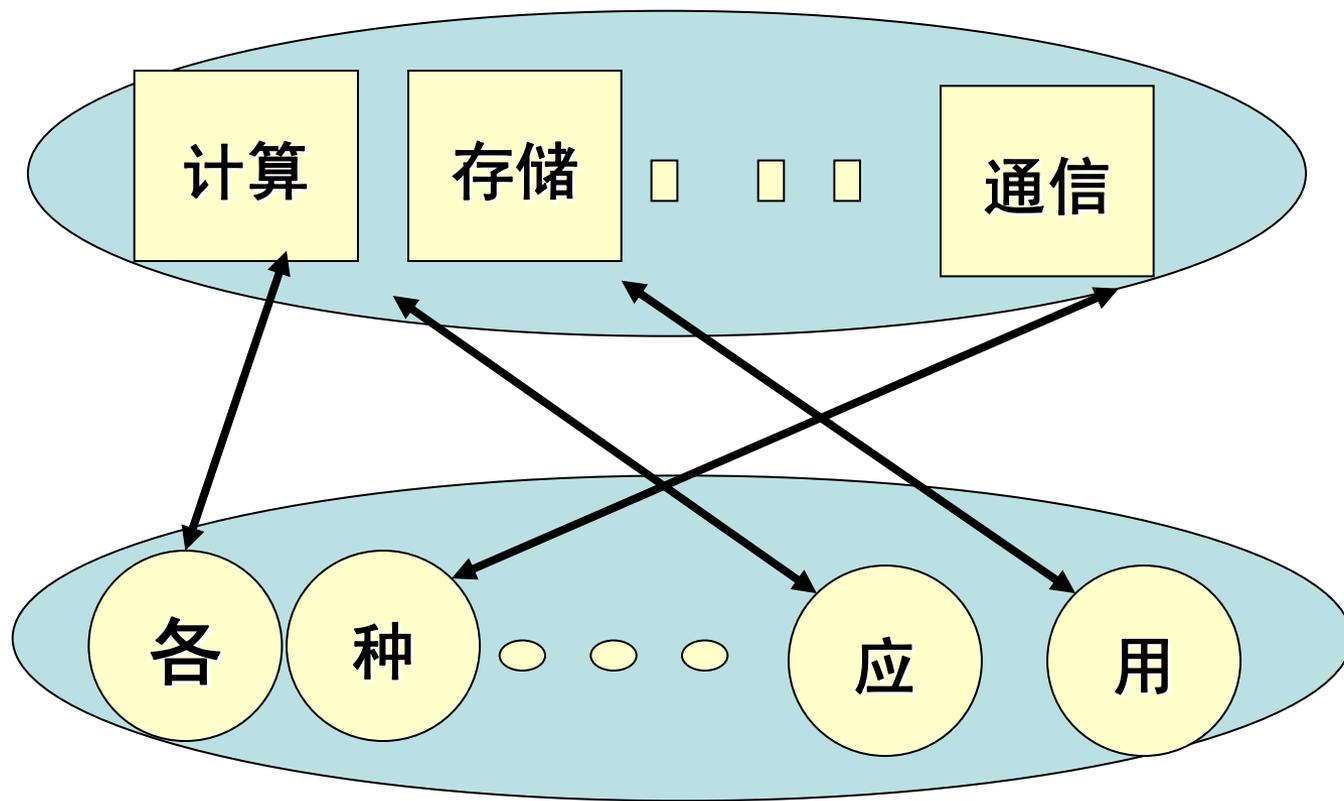
- 如果你问Patterson等国外学者：“你是研究什么的？”他们一般会简洁地回答：“System”。我国计算机学者一般很少讲自己是研究系统的。研究计算机系统的学者关心什么？
- 计算机系统研究关注的主要问题包括：
 - 计算机指令系统
 - 适应各种应用的系统结构（通用与专用）
 - 资源利用的效率（包括虚拟化技术）
 - 计算机使用的方便性和灵活性
 - 编程的效率(尤其是并行变成)
 - 计算机的性能和可扩展性
 - 系统的可靠性与安全性
 - 降低计算机生命周期的总成本
 - 减少计算机的能耗

目前的云计算只涉及其中少数问题，许多基本问题有待解决。

改变计算机系统技术的研究方向



计算机系统研究最基本的问题—— 满足应用需求的计算资源配置



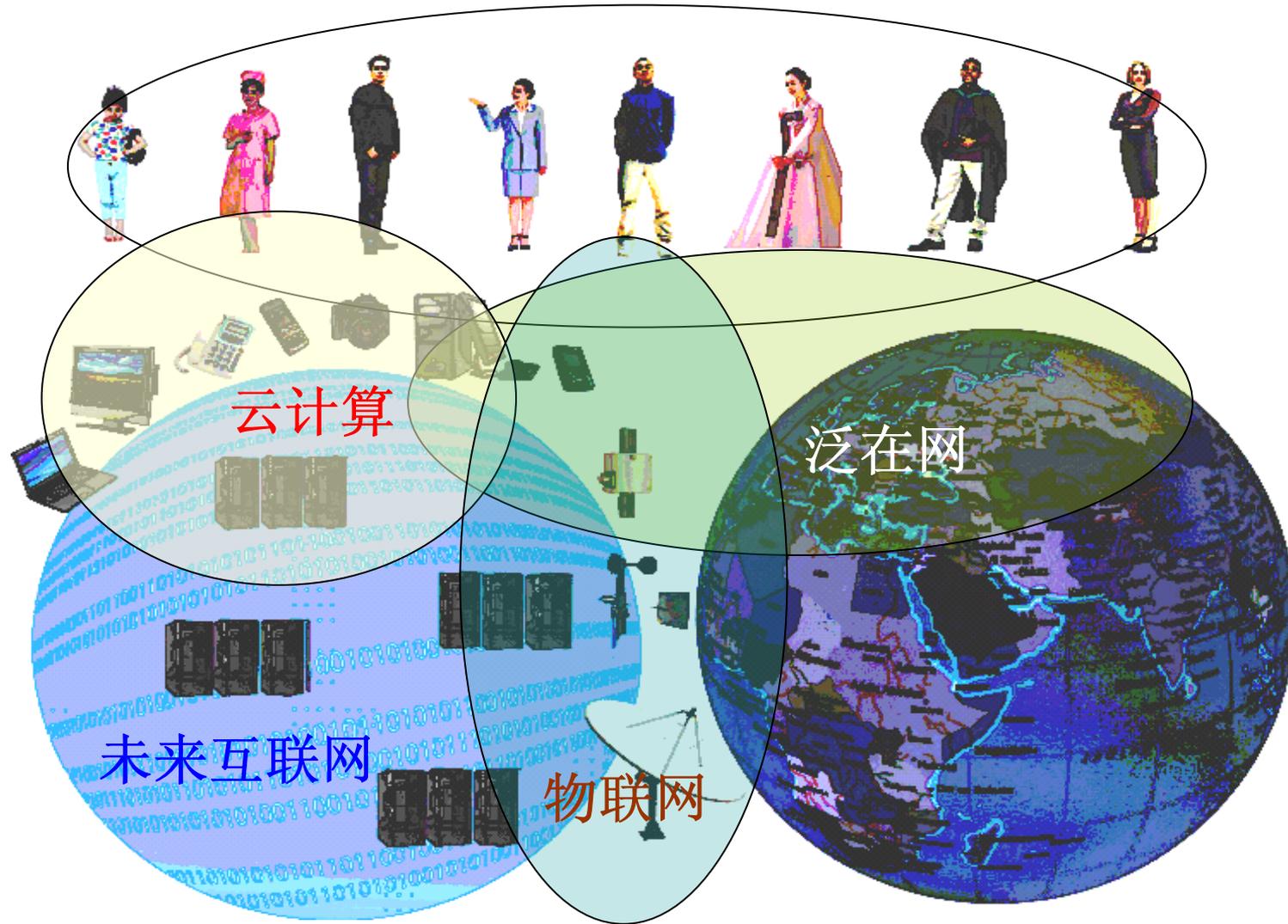
计算机应用的“昆虫悖论”

- 用一句话概括21世纪信息技术的发展趋势就是“为大众计算(Computing for the Masses”)。
- 数十亿用户和各行各业的应用需求一定千差万别。日本东京大学的坂村健教授曾把PC机、手机和物联网的应用种类分别比喻成**哺乳类动物(2万种)**、**鱼类(3万种)**和**昆虫类(100万种)**。计算机系统如何满足如此多的应用。
- 为每一种应用设计一种专用芯片和系统对供应商不经济，采用同一种通用计算机对用户而言效率不高。通用和专用是计算机系统发展中永恒的矛盾，也是最大的挑战。
- 可重构芯片可计算机，可复用的软硬件模块是计算机系统研究的追求目标。虚拟化技术也是解决此矛盾的途径之一

云计算是MRMT系统

- 我们可以仿照Flynn的计算机分类，将计算机系统按**资源自治域**（Resource）—**任务**（Task）分成4类（资源自治域的概念是借用互联网自治域的提法，需要认真定义和研究）
 - SRST（单资源单任务系统）1对1，如低端手机
 - SRMT（单资源多任务系统）1对多，如mainframe
 - MRST（多资源单任务系统）多对1，如某些HPC
 - MRMT（多资源多任务系统）多对多，如云计算
- 人们常说网格是**多对1**的系统，云计算是**1对多**的系统。实际上云计算的后台有许许多多资源，很难在一个操作系统控制之下，是一个典型的多对多的系统。
- 数据中心的庞大的计算、存储资源如何高效的调配是计算机系统研究的大问题。关键是资源是不是在统一的调度系统管理之下。

未来互联网、物联网、泛在网和云计算



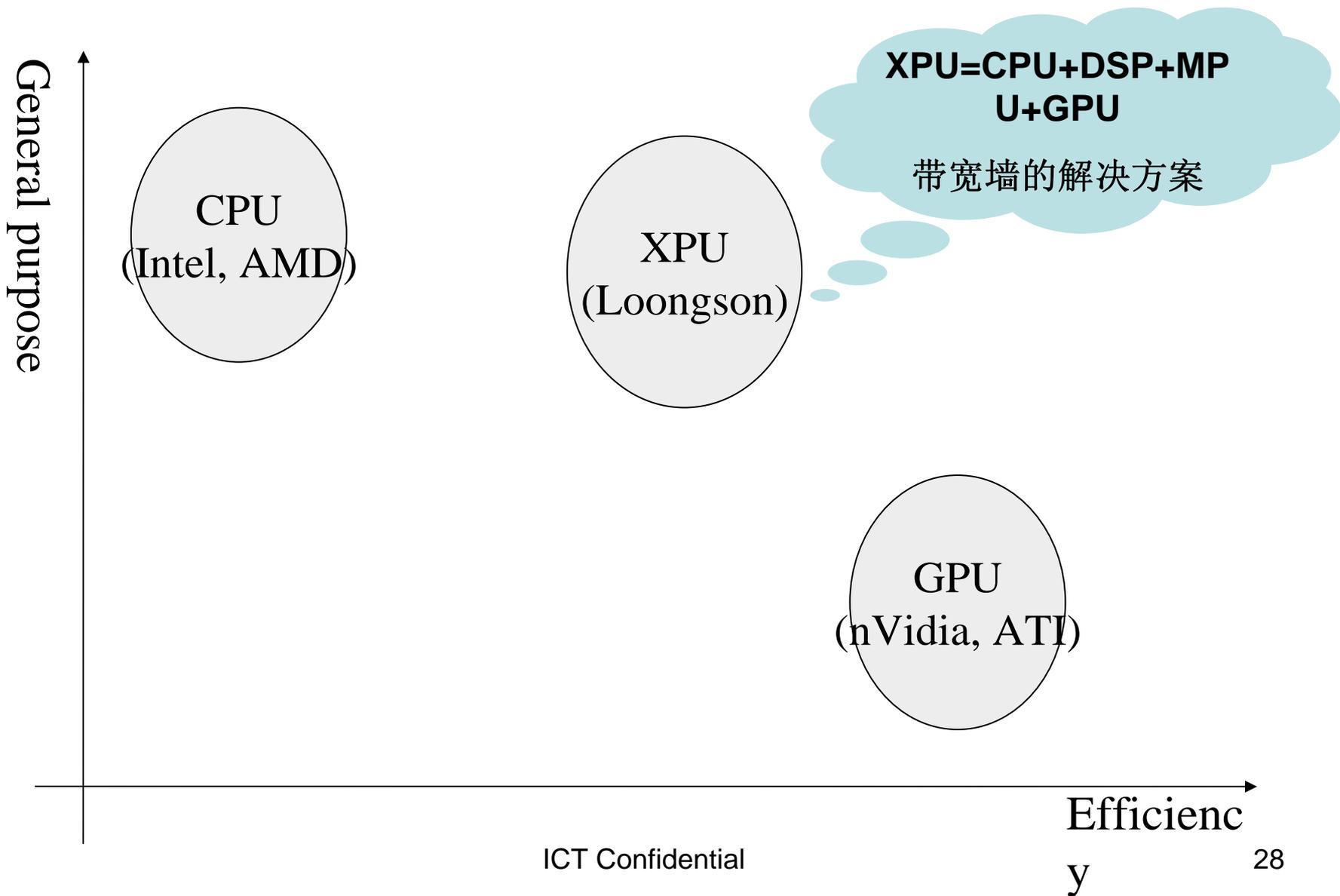
网络问题逐渐变成计算机系统问题

- 电信业正在进入“后电信时代”。通信技术与业务正在趋向计算技术与应用；计算技术与应用正在趋向网络与服务提供，CT、IT正在真正走向融合。联通研究院将这种融合模式称为“**公众计算通信网（PCCN）**”
- 在原有公众通信网的接入、交换、路由、传输要素的基础上，公众**计算**通信网还将实现计算处理能力、虚拟分配、调度管理以及业务开发等主要技术。
- 华为、中国移动等公司正在大量吸收懂系统结构的高端计算机人才。既懂计算机系统又懂通信协议的人才是目前最稀缺的人才。我国通信和计算机教育的分离不利于人才培养
- 通信领域两院院士陈俊亮教授到计算机学会担任服务计算专业委员会主任可看成一个标志性的事件。

计算机系统的难点在并行处理

- 并行处理已研究了几十年，论文多如牛毛，但进展不大。
- 云计算号称用1000台服务器工作1小时的成本与用一台服务器工作1000小时相当。问题是效率怎样，如果只完成了单服务器的1/10工作，仍然不合算。
- 并行计算最关心的“**如何提高计算机的性能和效率**”，这个问题从来没有改变，但答案在不断变化。
- 影响并行效率的障碍一是**编程（人工效率）**，二是通信延迟与带宽。**“带宽墙”可能比“存储墙”和“功耗墙”更高。**
- 历史上计算机设计的匹配规律是完成一次浮点运算需要保证一个字节的供数能力。目前主流CPU的运算速度与供数带宽之比是**1: 0.3-0.5**，即**100Gflops的芯片需要50GBps左右的内存带宽（4~8个DDR3）**。GPU芯片一个字节供数要完成十次以上浮点运算，典型的“茶壶煮饺子”。

XPU处理器核体系结构



“互连为王”——数据中心现状



互连带宽和延迟是必须啃下的硬骨头

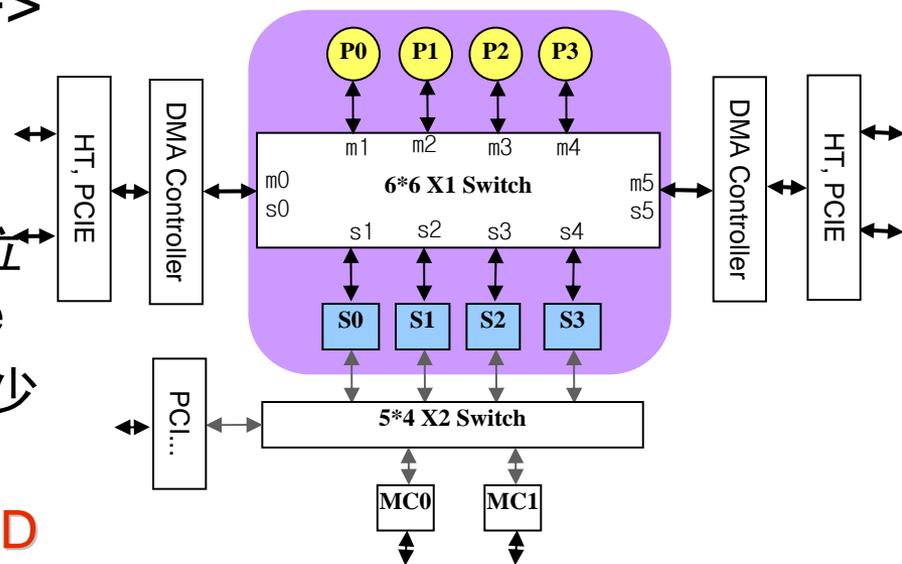
- 未来的数据中心需要与目前主流服务器不同的计算机系统来满足云计算的要求。
- 数据中心需要具有大量接口的网络交换器，其价格远远高于市面上流行的交换器，比普通交换器对分带宽(bisection bandwidth) 高**10倍**的交换器的价格要高出**100倍**。
- 目前的以太网技术无法让数据的传输速率超过每秒100G，主要因为没有这么多的能量来给提供这种数据传输速度的网络系统提供电力和进行冷却。
- 3D芯片和使用硅基光电子学来制造低成本、集成、传输速率达TB级的互连是解决互连问题的希望

Multicore Is Bad News For Supercomputers

- Throughput = Concurrency/Latency
 - Exploiting parallelism
 - Exploiting locality
- Multi-core Cannot Deliver Expected Performance as It Scales
- “The Troubles with Multicores”, David Patterson, IEEE Spectrum, July, 2010
- “**Multicore Is Bad News For Supercomputers**”, Samuel K. Moore, IEEE Spectrum, Nov, 2008

I/O为中心的处理器体系结构

- 以互联网络应用为代表的高吞吐量计算成为主流
- 处理器设计阶段：计算为中心->存储为中心->I/O为中心
- 面向I/O为中心的结构
 - 提升I/O系统在处理器存储层次中的位置，使I/O系统高于处理器二级cache
 - 优点：能够利用片上二级cache来减少对内存访问，提高性能
 - 应用于龙芯3号设计，测试表明：**SSD磁盘访问性能提高了40%，相对Intel平台有15%的性能提高**
- 相关文章发表于IEEE Micro 2009、HPCA'2010



图：四核龙芯3号结构

挖掘并行性是计算机系统的巨大挑战

时间	2020年	2030年	2050年
器件	CMOS	纳米量子器件	量子、生物分子
计算速度	Exaflops (10^{18})	Zettaflops (10^{21})	>Yottaflops (10^{24})
并行度	10^{8-9}	$10^{10} - 10^{12}$	$10^{13} - 10^{15}$
内存容量	25PB	EB (10^{18} B)	ZB (10^{21} B)
功耗	40MW	MW	MW
用途	核聚变模拟 蛋白质折叠等	地球模拟 生命科学等	MEMS优化 脑科学模拟等

2010

2020

2030

2050



这一场“并行革命”可能失败

- ▣ 人生三件很不愿做又不得不做的事：纳税、死亡，并行处理！
- ▣ 2007年1月，Stanford大学校长，计算机体系结构领域的权威学者 John Hennessy在ACM杂志上指出：“当我们谈论并行性和轻松地使用真正的并行计算机时，我们是在谈论一个计算机科学家面对的最困难的问题，如果我在计算机企业，我将感到恐慌。”
- ▣ IT产业从一个高成长的产业变成一个等待替代产品的产业，我们怎么办？如果软件不能有效地利用几十甚至上千个片内CPU核，计算机就不可能更新换代了，这是一个巨大的危机。

被扔进历史垃圾桶的并行计算机

Corporations Vanishing



Myrias
1991



BBN
1997



Multiflow
1990



ESCD
1990



Convex
1994



Kendall
Square
Research
1996



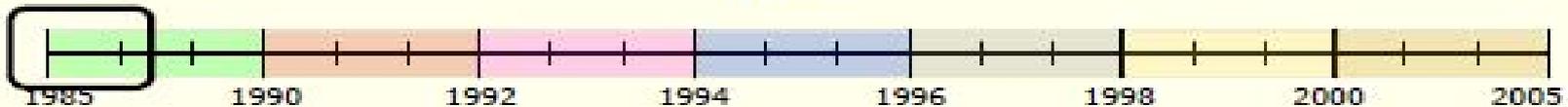
MasPar
1996



Cray Research
1996



nCube
2005



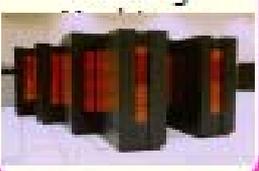
1989
ETA



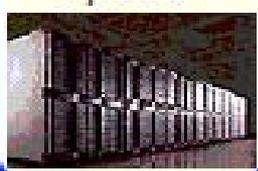
1992
Meiko Scientific



1994
Thinking



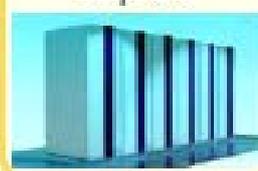
1995
Pyramid



1998
DEC



1999
Sequent



信息为什么这么“重”？

—解决计算机功耗问题的联想

- 假设要传送200TB的数据从北京到西安（1200公里），按200GB的硬盘一公斤计算，大约一吨重，按目前货运市价（每吨公里）0.1元左右计算，运费大约120元，加上其他开销不会超过**500元**。
- 若是用租用1Gps的专线，年租费70—180万元，按天算2000—5000元，每天满打满算可传送约10TB，需要20天，租金需要**2—10万元**。
- 这就相当于火车运送硬盘的“原子”只有1吨重，但需要送过去的“BIT”按运费算超过100吨，**BIT比原子“重”100倍**。
- 如果把全国的网络带宽增加**100倍**，即从10Gb到1Tb，可能全国发的电**一半**要用在网络上（现在占5%左右）

降低系统功耗的多种途径

—途径非常多意味着没有找到真正的途径

如何降低运算所消耗的功耗？

$$P = KC_{out}V_{dd}^2f_{clk}/2 + I_{sc}V_{dd} + I_{leak}V_{dd}$$

Layer	Technique	Reductions	Factors	References
System	System shutdown	41%~99%	V, k	[CSB94,CIC94]
	Dynamic voltage & frequency scaling	10%~73%	V, f	[SEO05]
	Algorithm selection	33%	k	[OY94]
	Compiler optimization	13%~20%	k	[Lee00]
Architecture	Data representation	13%~32%	k	[yu02][STD94]
	Parallel processing with low voltage	51%~80%	V, C, f	[CSB92]
	Cache design	20~80%	V, C, f	[Yang02] [BAF94,PR95]
	Bus encoding	15%~48%	k	[Lyu02]
	Operand isolation	30%~40%	k	[Banerjee06][Munch00]
Logic	Logic synthesis	<70%	C, k	[Hsu02][IP94][TMA95]
	Clock gating	20%~75%	k	[Li02][Monica03]
	Technology mapping	<47%	C, k	[LM93][TAM93]
	Path balancing	9%~41%	C, k	[Kim01][BCH94]
Circuit	Low swing clock	30%~63%	V	[ELE00][HK98]
	MTCMOS,VTCMOS,DTCMOS	20%~80%	I_{leak}	[Li99][Far97][Tad96]
	Power gating	<100%	V, I_{leak}	[David02]
	Device stacking	1%~56%	I_{leak}	[Halter97][Rahul03]

龙芯3B CPU在性能功耗比上 达到了目前世界先进水平

芯片型号	频率 (GHz)	工艺 (nm)	核数	Die面积 (mm ²)	功耗 (W)	双精度浮点 峰值 (GFLOPS)	性能功耗比 (GFLOPS/W)
Intel Core i7 980 XE	3.2	32	6	240	130	107.55	0.827
Intel Sandy Bridge	3 ~ 4	28	8	370	130	256	1.96
AMD Opteron X12	2.4	45	12	346	130	152	1.16
IBM Power7	3~4.1	45	8	567	100	264.96	2.64
IBM PowerXCELL	3.2	45	9	221	80	100	1.25
Fujitsu SPARC fxVIII	2.2	42	8	513	50	128	2.56
龙芯3B	1.0	65	8	300	40	128	3.2

可靠性设计面对的科学问题

自测试自诊断自修复—3S原理

如何在由数十亿个元件组成的芯片上构造稳定可靠的系统？

缺陷容忍

成品率 可靠性

故障容忍

“Killer”缺陷

高能粒子、电源噪声、温度波动

“Latent”缺陷

制造周期

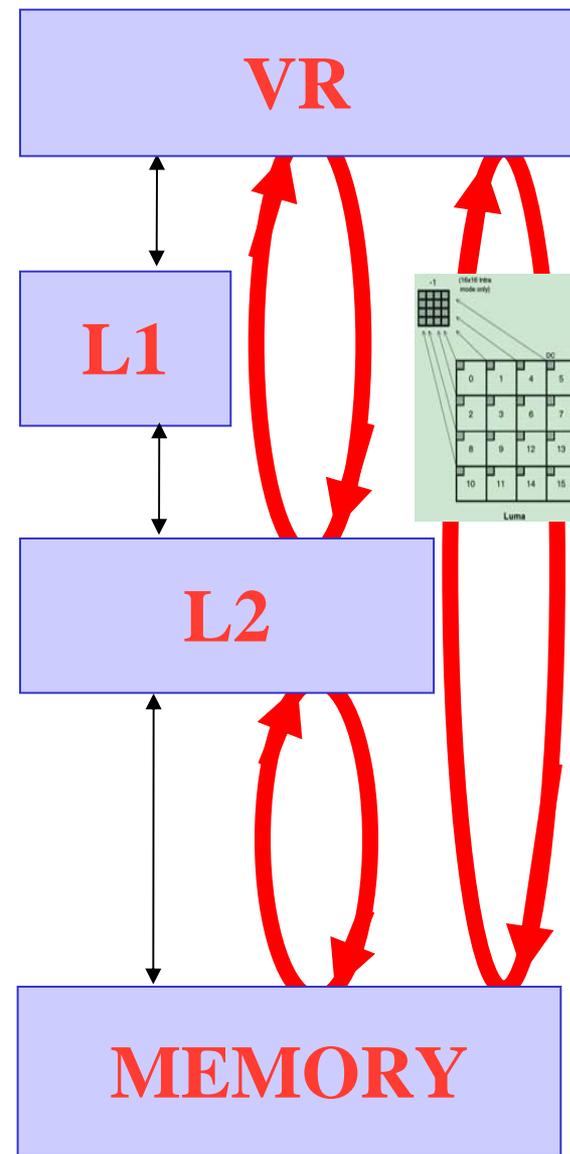
生命周期

计算机系统必须解决存储层次问题

- L1 cache reference: 0.5 ns
- Branch mis-predict: 5 ns
- L2 cache reference: 7 ns
- Mutexlock/unlock: 25 ns
- Main memory reference 100 ns
- Compress 1K Bytes with Zippy 3000 ns
- Send 2K Bytes over 1 GBPS network 20000 ns
- Read 1 MB sequentially from memory 250000 ns
- Round trip within data center 500000 ns
- Disk seek 1000000 ns
- Read 1MB sequentially from disk 2000000 ns
- Send one packet from CA to Europe 15000000 ns

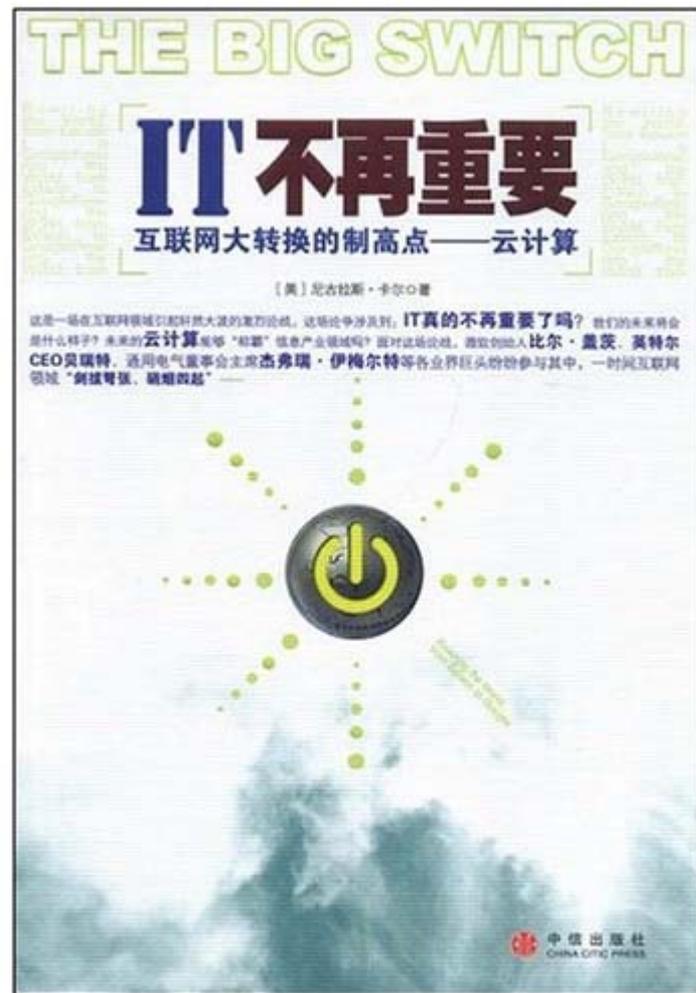
龙芯XPU 体系结构-GS464V

- 功能强大的向量运算部件
每个核两个256位向量部件
新增300多条SIMD 指令 (Linpack, FFT, filter, media.....)
- 专门的数据通道
VR、L2、MEM之间的专门链路，类似DMA
运算和数据引擎并行
数据传送过程中的重构：矩阵转置、FFT位反、媒体熵解码....
- 融合通用性与运算效率
Linpack效率>90%，
FFT效率>85%，
单核1080p的H.264解码>100帧/秒
在HotChips'2010上发表

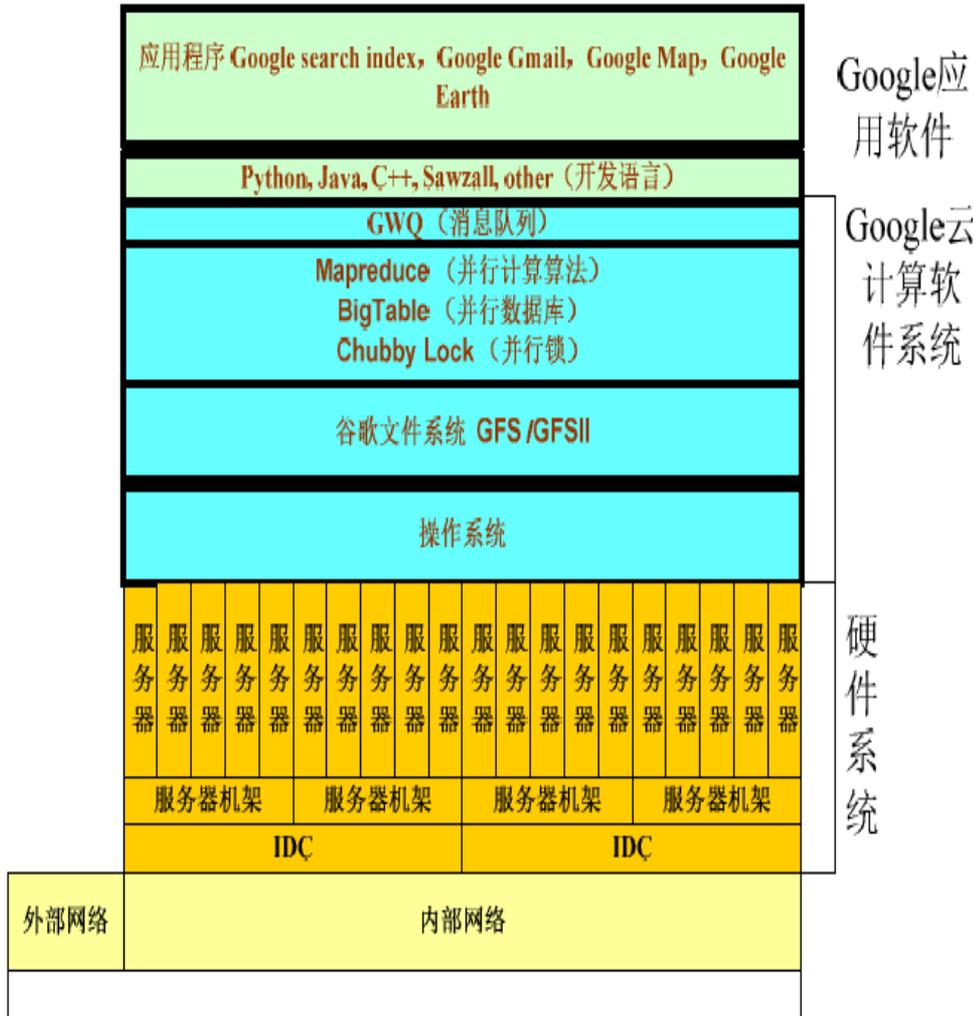


互联网不同于电网

- 对云计算的技术转换意义讲述最明白的书是Nicholas Carr写的“The Big Switch”,国内翻译成“IT不再重要”。
- 如同小型（直流）电站必然转变为大型（交流）电站一样，个人电脑必将让位于公共运算时代。
- 但云计算与电力网的不同，**电力网只送能量不传送应用**，所有的应用都是客户的责任；而**信息服务应用可以通过网络传送**。
- 电脑运算比发电**更具模块性**，数据存储、处理、传送可分拆成不同的服务。由不同公司提供，减少供应方的垄断。
- 云计算使万维网确实变成了**万维电脑**，成为我们感官和心灵的延伸。



Google的系统架构



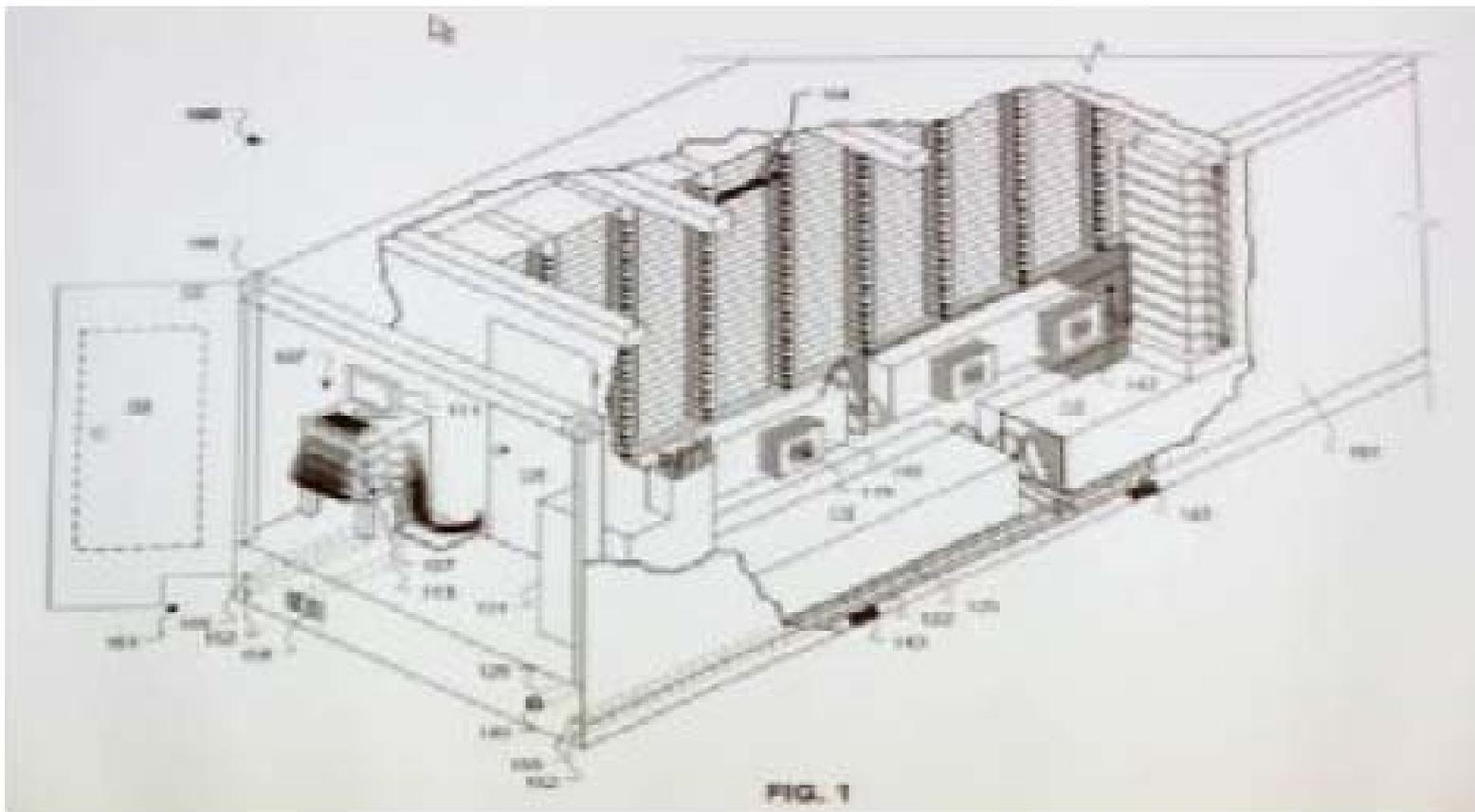
- Google 最大的优势在于它建造了既性价高(并非廉价)又能承受极高负载的高性能系统，具有较高的成本优势。
- 其IT 系统运营成本约为其他互联网公司的60%左右。同时 Google 程序员的效率比其他Web 公司同行们高出50%~100%，原因是Google 已经开发出了一整套专用于支持大规模并行系统编程的定制软件库。

Google 达拉斯数据中心



- 占用了附近一个**180万千瓦**（长江三峡发电站的1/10）水力发电站的大部分电力输出，利用河水冷却服务器。

Google模块化标准集装箱



- 每个机箱装有**30个机架**，**1160台服务器**，每个机箱（？）的功耗**250KW**，便于快速部署，降低对机房基建的要求。

网上公开的Google技术

- Google的每个服务器机架内部连接每台服务器之间网络是**100M以太网**，在服务器机架之间连接的网络是**1000M以太网**。
- 目前 Google 已经在全球运行了**38 个大型的IDC中心**，超过**300 多个GFSII 服务器集群**，超过**80万台计算机**。
- Google 所拥有的八十万台服务器都是自己设计打造的，Google 认为这是公司的核心技术之一。
- 每个服务器刀片**自带12V 的电池**来保证在短期没有外部电源的时候运行

改变不触动“核心技术”的科研模式



获图灵奖的与系统有关的 计算机科学家

- 2009 Thacker, Charles P **Alto personal computer, Ethernet .**
- 2008 Liskov, Barbara **fault tolerance, and distributed computing.**
- 2007 Clarke, Edmund M, Emerson, E Allen, Sifakis, Joseph
effective verification technology
- 2006 Allen, Frances **optimizing compiler techniques**
- 2005 Naur, Peter **ALGOL 60 compiler**
- 2004 Cerf, Vinton, Kahn, Robert E **internetworking, TCP/IP**
- 2002 Adleman, Leonard M. Rivest, Ronald L, Shamir, Adi
public-key cryptography
- 1997 Douglas Engelbart **mouse GUI**
- 1992 Bulter Lumpson **PC environment**
- 1990 Fernando J. Corbato **资源共享的计算机系统开发**
- 1987 John Cocke **开发RISC计算机**
- 1983 Ken Thompson、Dennis M. Ritchie **UNIX操作系统。**
- 1967 Maurice V. Wilkes **存储程序的计算机**
- 1966 A.J. Perlis **先进编程技术和编译架构**



请批评指正!