

# 数据库发展的第二次浪潮

文 ■ SAP全球数据库精进中心技术总监 卢东明

数据库并不是一个新的话题，其产生于上个世纪70年代。到80年代时，数据库的发展开始进入一个高峰期，尤其在80年代后期时，Oracle公司所开发的用于商业应用的关系型数据库开始进入到大型企业，甚至到中小型企业里，产生了非常大的经济利益。

上个世纪80年代后期，Sybase作为一家新公司，在数据库市场进入相对稳定的阶段后，又挑起了新的战争。Sybase采用当时非常新的Client/Server的结构，第一次分割了硬件使用的模式。

上个世纪90年代初期，数据库进入了乱战期。当时在数据库市场，不仅有IBM这样综合性IT厂商，也有Oracle这样专业化的数据库公司，还有各种各样的数据库公司在各个方向上发展。事实上，每当一个产业里有多家公司在百花齐放时，也一定是技术创新最活跃的阶段。因此，我认为90年代初期是数据库发展的第一次浪潮，因为那时有各种各样的理念，有各种各样的技术都在不断的竞争，谁也没有垄断这个市场。

从上个世纪90年代后期，到2000年这10年里，数据库进入一个相对稳定、相对整合的时代，直接的表现是很多有特点的数据库公司都纷纷被大企业收购，这种变化使得到2000年之后，企业选择数据库的

时候已经不再问技术了，更多的是在问品牌：比如你的公司在市场上占有率多高，谁的市场占有率高就选谁的产品。

## 新技术推动第二次浪潮出现

不过，现在这个时间点非常有趣，在数据库市场第一次浪潮出现20年后，我个人认为数据库又进入了第二次浪潮。为什么这样讲呢？我们看一下过去10年，甚至更多的时间里统治数据库市场的技术，几乎都属于行式数据库。行式数据库的优点是易于进行数据的存储、删除、查询。但是现在已经进入了大数据的时代，是分析的时代，决定商战最终胜负的因素，是看谁能最快地挖掘出数据的价值。这也意味着，我们必须有新的技术出现，来满足大数据时代用户的新需求。

那么，在各种各样大数据的场景下，现在有哪些新的技术产生呢？一个是列式数据库，二是数据流处理，三是嵌入式数据库，四是内存数据库。可以看到，现在数据库市场又出现了各种各样丰富的技术，就好像上个世纪90年代初期一样，它又会是一个百花齐放、百家争鸣的时代，各种不同的技术在影响着我们的应用。所以，我认为数据库市场现在又到了一个新浪潮的起落点。



## 全球最好的“ROI”

下面，我们先来看看列式数据库。传统的行式数据库按行来存储数据，在数据存储、查询时很方便，但如果要对数据进行分析，尤其是我们只分析其中的某些字段的时候，就会产生一些问题。

而应运而生的列式数据库则可以非常好的解决这个问题。第一，它只在你所需要的字段上发生I/O，I/O效率可以提高10倍。而且，因为列式数据库是按列来存储，所以很多信息是可以压缩的。比如在一个会场里有几百个人，如果我想知道这些人中，拿苹果手机的人有多少人。行式的做法是每个人把他包里所有的东西都给我，我一个一个来翻，显而易见，这种做法非常慢。但是如果我知道有可能做什么样的分析，在这些人一进门时就放几个箱子，让大家把不同种类的物品放在不

▲SAP全球数据库精进中心技术总监卢东明：“在数据库市场第一次浪潮出现20年后，我个人认为数据库又进入了第二次浪潮。”



同的箱子里,那么我只需要到手机的箱子里去找,就可以得到准确的数据。这就是列式数据库的基本原理,非常方便,非常迅速,可以成量级的提高分析的速度,同时降低分析所存储数据的规模。

目前,列式数据库在中国已经得到大量的使用,比如电信行业里用于信令的存储和分析。在国外也有类似的案例,比如美国的税务局,在分析全美国纳税人七年的报税记录,总共十几亿条信息时,也是非常好地使用了列式数据库。

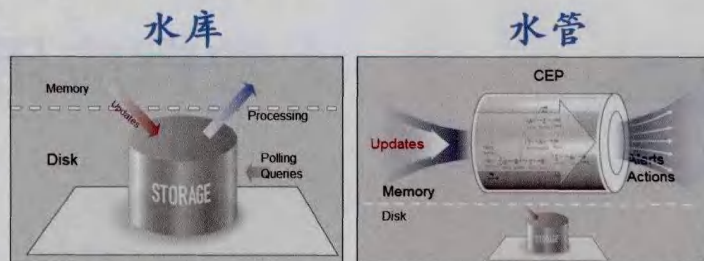
美国税务局由于采用列式数据库进行存储、分析,使得效率大大提升,其2006年的审计数量为2000年的3倍,收缴的税款达到592亿美元,上升了75%。这可能是全球最好的一个ROI(投资回报率),想象一下,如果我是税务局的CIO,我要给领导打一个报告,说想购买一种新型的数据库,可能要花一千万美元。领导可能会问:购买这一数据库大概会带回多少回报?而如果我的回答是:大概500亿美元左右,那么领导肯定会毫不犹豫地批准。

### 用数据流技术“捕鱼”

接下来我们看看数据流技术, CEP(Complex Event Processing, 复杂事件处理)。什么叫复杂事件处理?

举个例子,一个汽车里会有各种各样的传感器,监控着汽车的车速、车压、油压、水温、刹车等等,每一个传感器会不断产生新的信号,而每一个信号都是一个简单事件,但是这些简单事件组合在一起就可能成为复杂事件。比如在过去的两秒钟里,一辆汽车的速度从每小时80公里,忽然间降到了15公里。汽车的胎压从2.5,急速的降到了1.5。座椅重量从160斤降成25斤。把这几个事件都综合在一起,

图 数据库和数据流技术的区别



1. 基于硬盘实现
2. 先存储, 然后查询、处理
3. 先查询 → 再获得数据

1. 基于内存计算
2. 事件定义→过滤、识别
3. 数据始终在流动→同时支持查询

汽车里的电脑会实时地分析这一个复杂事件,做出一个实时判断:汽车可能撞车了,司机可能已经离开座位。可以看到,这时我们关注的不再是数据的存储,而是数据的使用,它要触发下一个事件,在这个例子中,所要触发的事件就是打开安全气囊来保护司机的安全。需要注意的是,数据库和数据流之间有很大的区别:数据库就如同水库,首先是要把数据存储下来,然后以后再使用,更多是用于事后的分析。但是对数据流而言,就如同一个水管,它希望做到的是当有鱼经过这个水管的时候,能够实时把鱼捕获。因此,这是一个非常重要的方法论的变化,事实上,目前在很多行业都在使用这一技术实现对数据的实时分析。

### 嵌入式数据库助力移动应用

接下来再看看嵌入式数据库。如今,每个人的手机从能力而言都超过了10年前的一台笔记本电脑。事实上,通过手机获取大量的数据正成为大数据的来源之一。

那么,在手机或者一些嵌入式设备里如何使用数据库呢?这本身就是一个非常大的挑战,也是一个非常大的商机。嵌入式的环境或者

移动办公的环境和企业中的机房是完全不一样的,环境可能十分复杂,也没有专业的人员进行管理,但是它又确实起着非常关键的作用,比如在物流、零售、餐饮等行业,会有很多移动应用的场景,在这些应用中,都可以使用类似的嵌入式数据库来帮助系统进行数据的分析。

### 内存数据库带来效率大幅提升

最后,让我们来看看内存数据库。内存数据库并不是一个全新的概念,但是随着HANA技术的推广和宣传,大家对内存技术有了新的认识。内存数据库对我们来说意味着什么?

事实上,相对于传统的磁盘数据库管理方式,内存数据库是将数据放在内存中直接进行操作。相对于磁盘,内存的数据读写速度要高出几个数量级。因此,将数据保存在内存中相比从磁盘上访问能够极大地提高应用的性能。这也是为什么我们看到,使用HANA技术的全内存计算,与行式数据库相比,在效率方面可以提高几个量级。

可以看到,上述这些新的数据库技术的出现,将会推动数据库市场第二次浪潮的出现。S