

# 从系统角度审视大图计算

吴城文,张广艳,郑纬民

清华大学计算机科学与技术系 北京 100084

## 摘要

大图计算已经成为学术界和工业界的一种基本计算模式,并且已经被应用到许多实际的大数据计算问题上,如社交网络分析、网页搜索以及商品推荐等。对于这些问题,大图的规模约有10亿级的点以及1 000亿级的边,这样的规模给大图的高效处理带来了诸多挑战。为此,介绍了大图计算的基本特征和挑战、典型的计算模型以及具有代表性的分布式、单机处理系统,同时对图处理系统中的关键技术进行总结,最后从系统的角度给出大图计算可能的一些研究方向。

## 关键词

大数据计算;大图计算;计算模型;计算系统

doi: 10.11959/j.issn.2096-0271.2015028

## *Reviewing Large Graph Computing from a System Perspective*

Wu Chengwen, Zhang Guangyan, Zheng Weimin

Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China

### *Abstract*

Large graph computing has been a fundamental computing pattern in both academic and industry field, and it was applied to a lot of practical big data applications, such as social network analysis, web page search, and goods recommendation. In general, most of large graphs scale to billions of vertices, and corresponding to hundreds billions of edges, which brings us challenges of efficient graph processing. Therefore, the basic feature and challenges of current large graph computing, typical computing models, and representative distributed, and single machine large graph processing systems were introduced. Then, some key technologies employed in large graph computing were summarized. Finally, some research directions in large graph computing from a system perspective were given.

### *Key words*

big data computing, large graph computing, computing model, computing system

## 1 引言

图可以用来表征不同实体间复杂的依赖关系。因而,在许多实际的应用当中,如社交网络分析、网页搜索、商品推荐等都可以使用图来进行问题的建模和分析。然而,在大数据时代,这类问题的规模通常十分庞大,以社交网络为例,Facebook在2014年7月的用户已经达到22亿户<sup>1</sup>,而用户之间的关系数量则更多,以数据的方式进行存储通常会占用几百GB甚至TB级的存储量。因此,大图计算不仅是计算密集型,同时也是存储密集型问题,如何在可以接受的时间内对大图进行计算,是需要解决的难题。

通常,为了快速地对大图进行处理,常常会使用分布式并行计算的思想,但是由于图计算本身特征使得在实现并行图计算时,不能使用传统科学计算领域的并行模式(计算偏微分方程)<sup>[1]</sup>;且以往在处理大数据问题上的map/reduce<sup>[2]</sup>模式,在处理图问题时效率极低;另外,并行图算法库Parallel BGL<sup>[3]</sup>或CGMgraph<sup>[4]</sup>没有容错机制。基于以上几点,需要一套符合大图计算特点的高效分布式并行计算框架。现在一些常见的分布式处理系统有Pregel<sup>[5]</sup>及其对应的开源实现Giraph<sup>2</sup>以及GraphLab<sup>[6]</sup>、PowerGraph<sup>[7]</sup>、GraphX<sup>[8]</sup>和Cyclops<sup>[9]</sup>。这些分布式系统大部分采用“think like a vertex”的思想,即以点为中心(vertex-centric)的计算模型,如图1<sup>[10]</sup>所示。在这种模型中,所有的点从其入边的邻点获取数据,执行用户自定义的函数对自己的状态进行更新,然后将自己的更新状态通过消息发给其出边的邻点。还有少数一些分布式系统采用了其他的计算模型,如PowerGraph的以边为中

心(edge-centric)的计算模型,如图2所示。在这种计算模型当中,首先依次遍历所有的边,将边的源点的更新值通过其出边传递给目的点,然后遍历所有的更新值,将更新值更新到目的点(在PowerGraph中将gather操作移到了scatter操作前面)。另外,还有以块<sup>[11]</sup>、路径<sup>[12]</sup>为中心的计算模型,在这类计算模型中,针对图结构来进行图划分,增加了计算的局部性,但是也存在图划分时间过长等问题。

分布式图处理系统随着问题规模的扩大具有很好的拓展性,但是在提高系统处理效率方面仍然面临许多挑战。比如图的划分,要提高系统性能需要在保证集群各节点负载均衡的情况下,使得集群内各节点的通信量最少,是一个NP难问题。此外,一个分布式系统需要解决集群内各节点协同工作、容错等一系列问题,而这类问题对系统的性能有重要的影响。另一方面,对于使用分布式系统的程序员来说,

1  
<http://tech.qq.com/a/20140725/000288.htm>

2  
<https://giraph.apache.org/>

```
vertex_scatter(vertex v)
  send updates over outgoing edges of v

vertex_gather(vertex v)
  apply updates from inbound edges of v

while not done
  for all vertices v that need to scatter updates
    vertex_scatter(v)
  for all vertices v that have updates
    vertex_gather(v)
```

图1 以点为中心的计算模型<sup>[10]</sup>

```
edge_scatter(edge e)
  send update over e

update_gather(update u)
  apply update u to u.destination

while not done
  for all edges e
    edge_scatter(e)
  for all updates u
    update_gather(u)
```

图2 以边为中心的计算模型<sup>[10]</sup>

环境的搭建、编写分布式程序比较复杂，而且程序的调试和优化又相对困难。基于此，最近一些大图计算的研究工作，在使用单台计算机进行大图计算处理上有了一些新的成果，如以点为中心的计算模型的GraphChi<sup>[13]</sup>和以边为中心的计算模型的X-Stream<sup>[10]</sup>，另外还有VENUS<sup>[14]</sup>、GridGraph<sup>[15]</sup>等。这些成果极大地降低了大图计算的成本开销，同时能够达到甚至好于一些分布式图计算系统处理时延。

本文将介绍当前大图计算的主要特征及挑战，从系统角度给出当前大图处理系统的主要特征及其研究成果，并对图处理系统中的关键技术进行总结，最后给出大图计算系统方面可能的研究方向。

## 2 大图计算的特征及挑战

大图计算是大数据计算中的一个子问题，除了满足大数据的基本特性之外，大图计算还有着自身的计算特性，相应地面临着新的挑战。

### (1) 局部性差

图表示着不同实体之间的关系，而在实际的问题当中，这些关系经常是不规则和无结构的，因此图的计算和访存模式都没有好的局部性，而在现有的计算机体系架构上，程序的性能获得往往需要利用好局部性。所以，如何对图数据进行布局 and 划分，并且提出相应的计算模型来提升数据的局部性，是提高图计算性能的重要方面，也是面临的关键挑战。

### (2) 数据及图结构驱动的计算

图计算基本上完全是由图中的数据所驱动的。当执行图算法时，算法是依据图中的点和边来进行指导，而不是直接通过程序中的代码展现出来。所以，不同的图结构在相同的算法实现上，将会有着不同

的计算性能。因此，如何使得不同图结构在同一个系统上都有较优的处理结果，也是一大难题。

### (3) 图数据的非结构化特性

图计算中图数据往往是非结构化和不规则的，在利用分布式框架进行图计算时，首先需要对图进行划分，将负载分配到各个节点上，而图的这种非结构化特性很难实现对图的有效划分，从而达到存储、通信和计算的负载均衡。一旦划分不合理，节点间不均衡的负载将会使系统的拓展性受到严重的限制，处理能力也将无法符合系统的计算规模。

### (4) 高访存/计算比

绝大部分的大图计算规模使得内存中无法存储下所有的数据，计算中磁盘的I/O必不可少，而且大部分图算法呈现出迭代的特征，即整个算法需要进行多次迭代，每次迭代需要遍历整个图结构，而且每次迭代时所进行的计算又相对较少。因此，呈现出高的访存/计算比。另外，图计算的局部性差，使得计算在等待I/O上花费了巨大的开销。

## 3 分布式大图计算系统

本节将介绍几个典型的大图处理的分布式系统，重点突出每个系统的特点。

### 3.1 Pregel

Pregel是由Google公司开发的分布式处理图系统，其主要的设计思想是基于BSP (bulk synchronous parallel)<sup>[16]</sup>。在此思想上，Pregel使用了以点为中心的计算模型，对整个图根据点进行划分，将不同的点以及相关的邻边存储到不同的计算机上。在Pregel中，用户可以自定义

点的compute()函数,每个点多次迭代执行这个函数,并最终得出整个图的计算结果。具体地,在每一次迭代(superstep)中,每个活跃的点(active vertex)会执行compute()函数,在这个函数中,该点读取在前一次迭代中其邻点发送的消息,通过这些消息计算自己新的状态,再将自己最新的状态通过出边发送给其邻点(邻点将会在下一次迭代中收到这些消息),然后该点会进入不活跃状态(inactive),如图3所示。当不活跃的点(inactive vertex)在下一轮收到消息时,就会重新处于活跃状态。当所有活跃的点执行完compute()函数之后,当前迭代结束,并且进入到下一次迭代。如果系统当中所有的点都处于不活跃状态,并且没有任何新的消息,算法结束。

Pregel使用了消息传递(message passing)的方式进行计算节点之间的通信,在一次迭代中每个点可以向其他点发送任意量的消息,而这些消息将会在下次迭代中被对应的点读取。在分布式的环境中,为了减少机器间的通信量,提升计算的性能,当点的compute()函数的操作符合交换律和结合律时,Pregel可以支持用户实现combiner()函数,把从机器 $M_i$ 到另一台机器 $M_j$ 上点 $v$ 的所有消息合并成一条消息。

## 3.2 Giraph

Giraph构建在Hadoop<sup>3</sup>之上,是对Google公司Pregel的开源实现。Facebook使用Giraph来进行社交关系图的分析。为了提升系统的性能,在原有Giraph基础上增加了一些优化的措施。Facebook在Giraph的加载图数据、写回图数据以及计算阶段引入了多进程,提升了系统的整体性能,尤其对计算密集型的应用,引入多线程

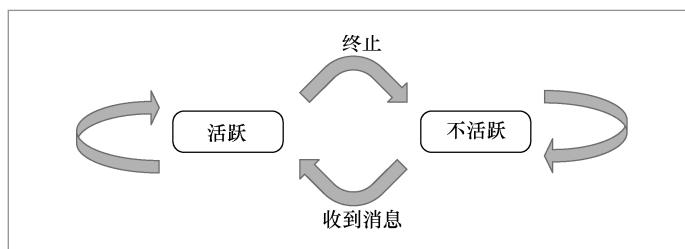


图3 Pregel点的状态机<sup>[5]</sup>

程可以使性能随着处理器的增加获得接近线性的加速比。

## 3.3 GraphLab和PowerGraph

与Pregel的同步数据推送的BSP模型不同,GraphLab使用异步的GAS(gather、apply、scatter)模型来实现大图分布式并行计算。GraphLab使用共享内存(shared memory)的方式来实现以点为中心的计算模式,在这种方式下,每个点可以直接读取和修改其邻点和邻边的值。在GraphLab上实现算法时,用户需要实现符合算法要求的GAS函数,在算法执行时,图的每个点都会执行该函数。

在gather阶段,每个执行GAS函数的活跃点从其邻点和邻边获取数据,然后使用这些值来计算自己的更新值,这里计算操作必须满足交换律和结合律。在apply阶段,活跃点将原来的旧值更新为计算得到的新值。在scatter阶段,活跃的点会通过邻边激活对应的邻点。如图4所示,在GraphLab中使用一个全局的调度器,各个工作节点通过从该调度器获取活跃的点来进行计算,这些正在被计算的点也可能将其邻点调入调度器中。最后当调度器中没有任何可调度的点时,算法终止。这种调度器的使用使得GraphLab同时支持算法的异步调度执行和同步调度执行。

在同步执行(synchronous execution)计算模式下,每个点或者边的更新不能

<sup>3</sup> <http://hadoop.apache.org/>

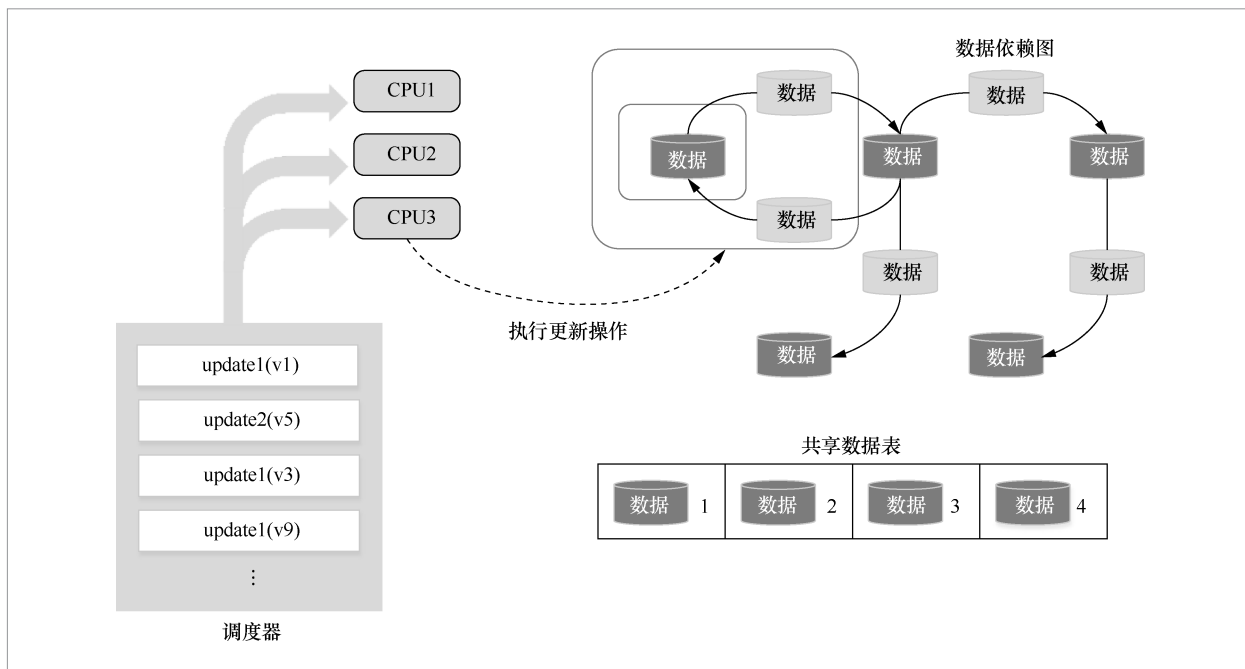


图 4 GraphLab 计算框架<sup>[7]</sup>

马上被当前迭代中接下来的计算感知到，直到当前迭代结束时，在下次迭代当中才能读取到更新的值。异步执行 (asynchronous execution) 与同步执行不同，点或者边的更新能够马上被接下来的计算所感知并使用到，这种计算模式可以使得如PageRank的一些算法收敛速度更快，但也同时会导致数据竞争，从而产生额外的计算开销。另外，在分布式系统中，

这种模式会产生随机的信息传递，因而也会产生较大的通信开销。一般来说，对于计算密集型的算法（如BP）来说，更适合使用异步计算的模式。

PowerGraph包含在GraphLab 2.2中，是在GraphLab的基础上对符合幂律分布 (power-law)<sup>[18]</sup>的自然图计算性能的改进，其主要改进是在图的划分上。如图5所示，PowerGraph使用了Vertex-cut的图

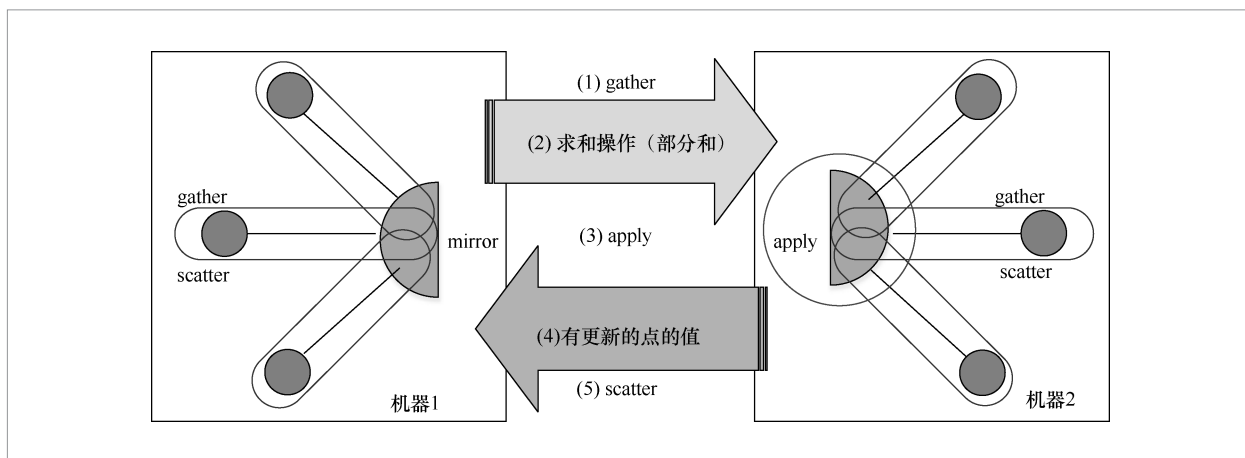


图 5 PowerGraph 切割点集划分及通信模式<sup>[7]</sup>

划分策略,将待处理的图以切割点集的方式进行划分,将那些度极大的点的边分割给不同的计算节点,同时,将对应的点也复制给这些计算节点作为镜像(mirror)点。具体计算时,每个主点及其对应镜像点在本地执行gather操作,随后镜像点将自己的计算结果发送给主点,收到全部计算结果后,主点执行apply操作,并且将更新值发送给所有镜像点,最后主点和镜像点进行scatter操作。

### 3.4 GraphX

如图6所示,GraphX是构建在分布数据流框架Spark<sup>4</sup>上的分布式图处理系统。GraphX支持Pregel和GraphLab的计算模型,并且拓展了Spark中的RDD(resilient distributed dataset,弹性分布数据集),引入了RDG(resilient distributed graph,弹性分布图),这种结构可以支持许多图操作,因此现有的大多数图算法都可以使用系统中提供的基本操作算子(如join、map和group-by)来实现,并且实现十分简单。为了利用Spark中这种算子操作,GraphX重构了新的vertex-cut图划分方法,将图划分成水平分区的顶点和边的集合。GraphX的性能比直接使用分布式数据流框架好一个数量级,稍差于GraphLab。另外,由于GraphX是构建在Spark之上的,所以GraphX能够得到低开销的容错和透明的错误恢复支持。

## 4 单机大图计算系统

随机单台计算机处理能力和存储能力的提升,再加上人们对于图计算模式研究的深入,一些在单机上处理大图计算的系统被提出,这些系统有着很好的图计算性

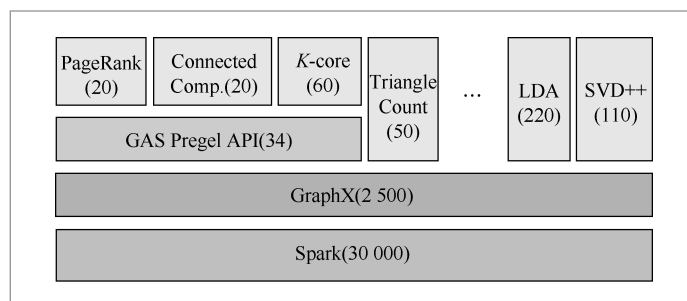


图6 GraphX的层次结构(括号中为代码行数)<sup>[8]</sup>

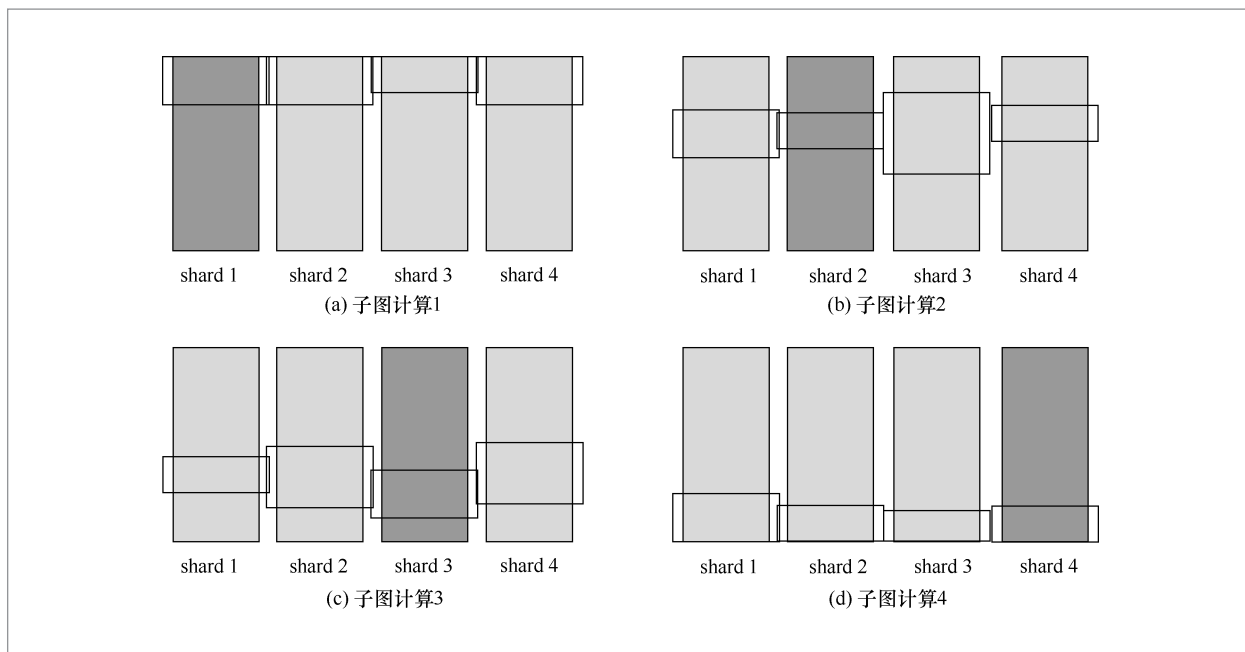
能,同时相比分布式系统,其低硬件成本和低功耗的优势明显。本节将介绍几个代表性的单机大图计算系统。

### 4.1 GraphChi

GraphChi是一个基于磁盘的单机大图处理系统。在大图计算中,计算的访存局部性非常差,严重影响到计算的性能。特别地,在单机情况下系统的计算能力十分有限,因此,为了提升计算性能,GraphChi使用了具有创新性的磁盘数据布局和对应的计算模型来减少磁盘的随机访问;使用选择性的调度来加速算法的收敛。

磁盘数据的布局和计算模型。GraphChi在计算前首先会对图数据进行预处理,将输入的图划分成多个shard,每个shard中存储对应点集的所有入边,并且将入边按照其源节点的ID进行排序,划分时需要保证每个shard中边的数量大致相同,每个shard都能够加载进内存。GraphChi使用以点为中心的计算模型,使用并行滑动窗口(parallel sliding window)来加载数据进行计算,如图7所示,每次(interval)计算一个子图,即一个shard所对应点集中所有点的值,需要顺序读取某个点集对应的入边(深灰色部分)以及该点集在其他shard中所对应的出边(黑色矩形框部分),这种数据布局和计

4  
<http://spark.apache.org/>

图7 并行滑动窗口计算模型<sup>[12]</sup>

算模型可以保证每次计算的I/O是顺序的。这样，一次迭代计算整个图中所有点的值，多次迭代，直到算法收敛。

选择性的调度。在GraphChi中可以使用选择调度性调度(selective scheduling)策略来加快图中某些点的收敛，尤其是对这些在两次相邻的迭代当中变化很显著的点。在点执行update()函数时，类似GraphLab中的apply()，可以将其

邻点加入调度器中，进行选择性的调度。

## 4.2 X-Stream

与GraphChi所使用的以点为中心的计算模型不同，X-Stream使用以边为中心的计算模型，并且所有的状态都保存在点中。X-Stream的计算过程主要分为3个阶段：scatter、shuffle和gather，如图8所示。在scatter阶段，X-Stream依次遍历每一条边，判断边的源节点是否产生更新，如果有更新产生，将边通过出边发送给目的节点。shuffle阶段是在对图进行划分之后，需要增加的一个不同划分块之间更新数据交换的阶段，主要是为了降低在scatter阶段的随机写开销。在gather阶段，X-Stream依次遍历在scatter阶段产生的所有更新，并更新对应点的状态值。X-Stream以边为中心的计算模型对边进行顺序访问，可以充分发挥磁盘的等二级存储介质的顺序访问高带宽加速图计算，但是在X-Stream中对点的访问还是随机的，为了对此进行优化，进一步提

```

scatter phase:
  for each streaming_partition p
    read in vertex set of p
    for each edge e in edge list of p
      edge_scatter(e):append update to Uout

shuffle phase:
  for each update u in Uout
    let p=partition containing target of u
    append u to Uin(p)
    destroy Uout

gather phase:
  for each streaming_partition p
    for each update u in Uin(p)
      edge_gather(u)
      destroy Uin(p)
  
```

图8 X-Stream 以边为中心的计算模型 (U<sub>in</sub>/U<sub>out</sub> 为输入 / 输出缓存)<sup>[13]</sup>

高计算性能, X-Stream对图的点集合均等划分成小的子点集合, 每个子点集合其每个点所有的出边也对应地组成一个边的划分集合。对点的划分主要满足每个子集合中的点都能够存储到内存中, 这样当计算每个划分块时, 对点的随机访问开销能够极大地降低, 为X-Stream进行划分后的计算模型。

在对图进行划分之后, 每个划分块在scatter阶段, 首先将所有的更新值写在本地的一个输出缓存中, 当所有的块都完成scatter之后, 进入一个shuffle阶段, 这个阶段的主要工作是将所有划分块的更新进行分配, 将更新分配到对应的划分块的输入缓存中, 作为gather阶段的输入, 对点的状态进行更新处理。相比于GraphChi, X-Stream对所有边进行顺序访问, 能够充分发挥磁盘等二级存储介质的顺序带宽的速度, 同时预处理阶段(简单的散列图划分操作)无须进行开销巨大的排序处理, 因此能够获得较好的图处理性能。

### 4.3 VENUS

尽管GraphChi在大图处理上能够取得

较好的计算效果, 但是也存在如下的缺陷: 预处理需要对边的源节点进行排序, 开销大; 图数据的加载和计算是分开的, 没有充分利用磁盘和I/O的并行来提高计算性能; 对shard内的边排序后, 每个点所对应的边不在相邻的位置, 缓存局部性不高。

基于以上的这几点观察, 笔者提出了如图9所示的以点为中心的流线型(vertex-centric streamlined)计算模型。在这种计算模型中, 笔者分别构建了g-shard和v-shard, 其中g-shard与GraphChi中shard的概念类似, 存储了一个子点集对应的所有入边, 但是不用对边进行排序, 而是将目的顶点相同的边存储在相邻的位置, v-shard存储对应一个g-shard中所有目的顶点和源顶点的值。另外, 使用了一个全局的点值表, v-shard从其中读取和写回对应的点值。系统计算点的更新值时, 无须像GraphChi将所有的入边和出边同时加载进内存, 只需将入边加载进内存, 同时节点更新后, 不用再将更新值写入出边, 这样可以极大地减少I/O。此外, 当加载完g-shard中一个点的所有入边时, 即可对该点的值进行计算, 重叠了I/O

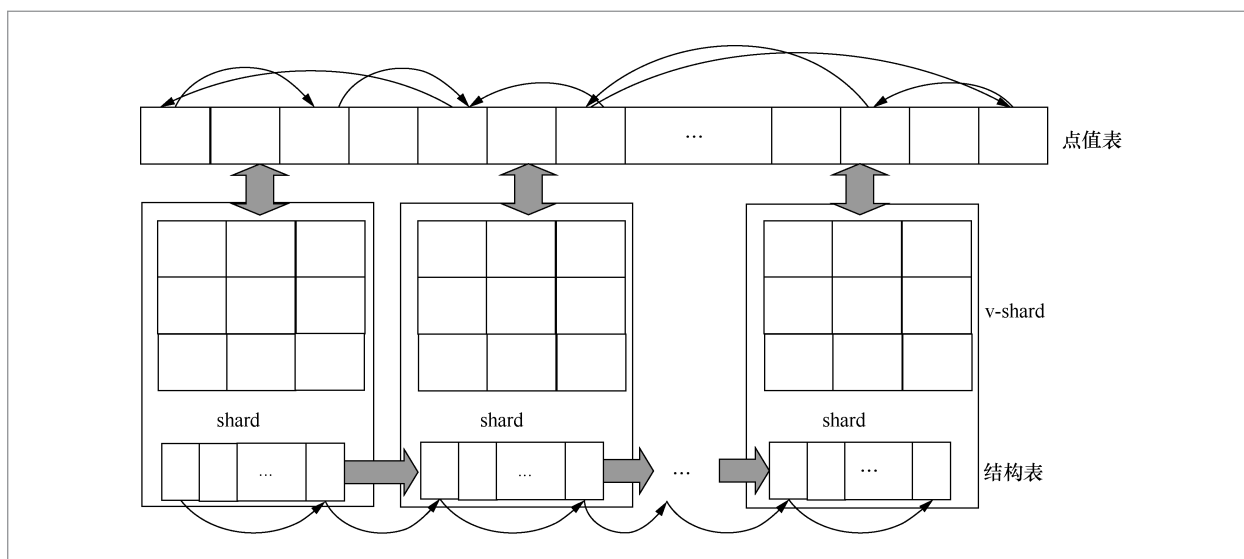


图9 以点为中心的流线型计算模型<sup>[14]</sup>

和CPU的时间开销，极大地提高了系统的性能。实验结果表明，VENUS的性能显著地好于GraphChi和X-Stream。

### 4.4 GridGraph

在X-Stream中，在scatter和gather阶

段之间，还需要一个shuffle阶段将每个划分在scatter阶段产生的更新值分配到对应划分的输入缓存中，供gather阶段进行计算。在scatter阶段，更新值会有 $O(|E|)$ 这样的规模，其中 $|E|$ 代表图中边的数量。所以，当内存不足时，需要将一部分缓存先写入磁盘，并且在gather阶段需要将写入磁盘的更新值重新读入内存，因此，在此过程中可能会触发较多的I/O，严重影响系统的性能。

为此，GridGraph提出了如图10所示的格子划分方式。首先，将整个点集划分成相同大小的 $P$ 份子点集，然后将边以行和列划分成格子，每一行对应在某个子点集内的点所对应的所有出边，每一列对应在某个子点集内的点所对应的所有入边。对应这种图的划分方法，笔者提出了双重滑动窗口的计算模型（如图11所示），是图10（a）中图结构的PageRank

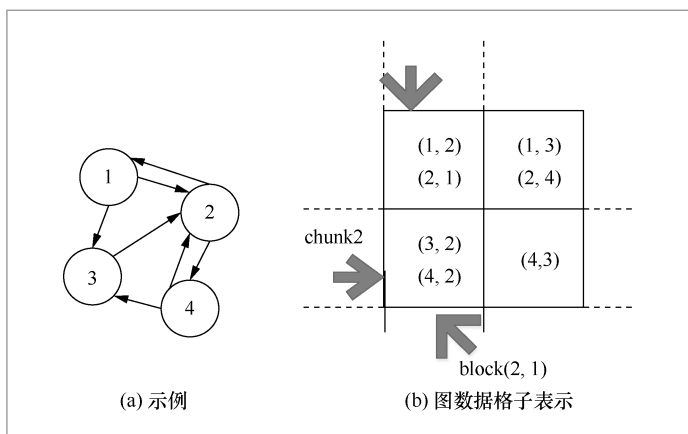


图10 GridGraph的图划分例子<sup>[5]</sup>

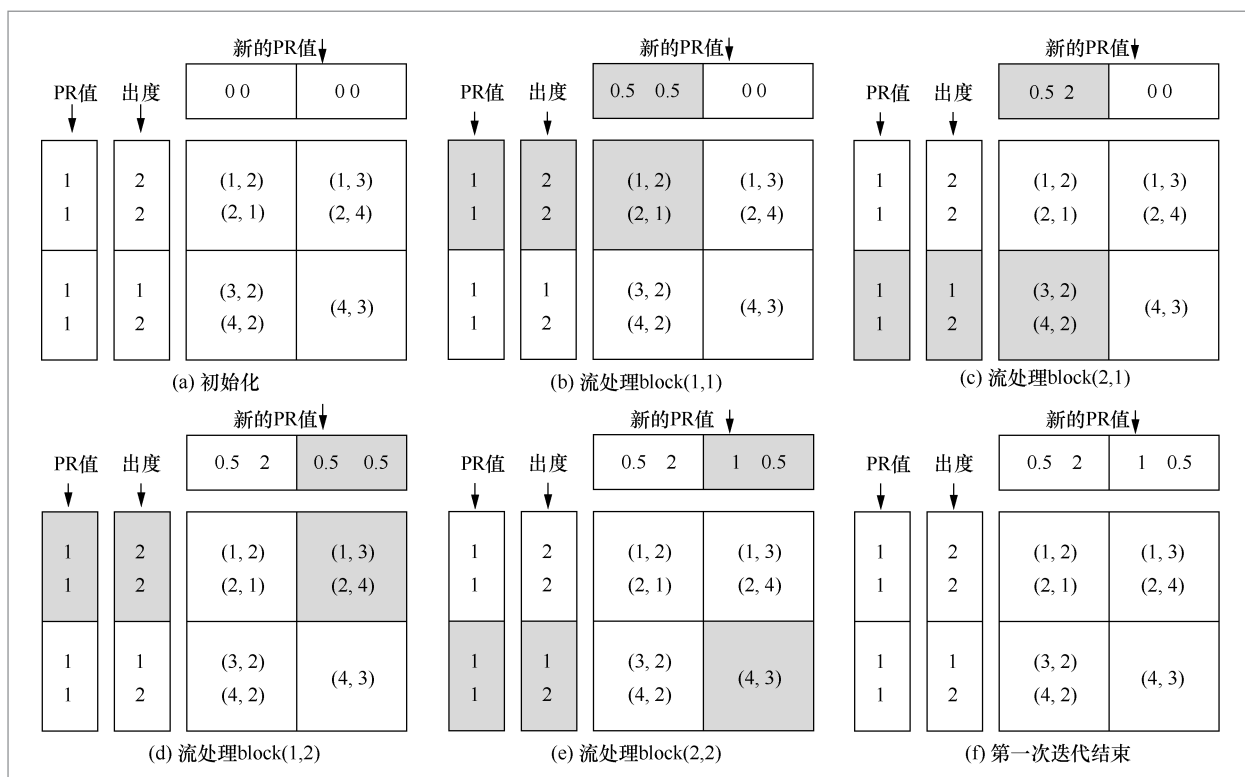


图11 双重滑动窗口计算模型示例<sup>[5]</sup>

第一次迭代过程, 计算点的更新值需要读取其入边源节点的值, 为此从上到下, 依次读取该列每个格子内的边进行计算, 然后当一列计算完毕后, 即完成一个子点集中点的值的计算, 窗口滑动到下一列, 继续进行计算, 直至所有的格子都遍历完毕。在这种计算模型中, 值的更新计算操作必须符合交换律, 另外, 这种方式点的更新是就地更新, 不会产生中间的更新结果, 极大地减少了I/O, 同时, 点的数据访问的局部性也有了提升。在进行图划分时, 使用二级的图划分策略, 即先将图划分成 $Q$ 份, 使得每个格子的边都能够存储进内存中, 然后再对刚才的每个格子进行划分, 使得每个小格子能够存储进最后一级cache (LLC) 当中。另外, GridGraph还支持选择性的调度, 在BFS和WCC这样的算法中, 可以极大地减少I/O, 提高计算性能。

## 5 大图计算中的关键技术

本节将介绍在分布式和单机图处理系统中常用的技术。

### 5.1 异构计算平台

在异构计算系统中, 存在着计算能力和计算特点不同的计算单元。比如, GPU具有比CPU更强的多线程并行计算能力, 因此在异构系统中, CPU会把一些或者全部的计算交给GPU来执行。在图计算领域, 相关的异构计算系统已经被开发出来。TOTEM<sup>[19]</sup>会将度高的点交给CPU计算执行, 而将度低的点交给GPU来执行。而另外一些系统, 如MapGraph<sup>[20]</sup>和CuSha<sup>[21]</sup>等, 会将整个图都交给GPU来执行。除了GPU和CPU的异构图计算平

台之外, 一些研究人员发现, solid-state drive (SSD) 有着与传统hard disk drive (HDD) 不同的访存特性。一些图计算系统 (如TurboGraph<sup>[22]</sup>和FlashGraph<sup>[23]</sup>) 针对SSD对计算系统进行了优化, 使得系统在SSD上有着很高的计算性能。目前使用异构计算的平台的图处理系统主要是单机图处理系统。

### 5.2 通信模型

在消息传递的通信模型中, 算法中点的状态保存在本地, 通过消息传递的方式更新在其他机器上点的状态。在Pregel和Giraph中, 使用了消息传递的通信模型, 为了确保所有更新的数据可用, 需要在前后两次迭代计算之间加入一个同步操作。

在共享内存的通信模型中, 各个处理单元允许并发访问和修改相同地址的数据。在一些分布式的计算系统 (如GraphLab和PowerGraph) 中, 使用了虚拟共享内存来实现各计算节点之间的透明的同步。在这些图处理系统中, 使用了假点 (ghost vertex) 的方式来实现虚拟共享内存。在假点的这种实现策略中, 图中的每个点有一个归属的工作节点, 另外有一些工作节点拥有该点的副本。因此, 在这种通信模型中, 当多个工作节点并发访问同一内存地址时, 需要考虑数据一致性的问题。

### 5.3 执行模型

#### (1) 同步执行

许多图算法由一系列迭代计算组成, 在前后两次迭代之间有一个全局的同步过程。这种执行模式将计算节点之间的通信控制在每次迭代的结束, 因此适合于那些

计算量小而通信量大的算法。

#### (2) 异步执行

在图中某个点的值有了更新值之后，立即将这个最新的更新值更新到该点上。在这种执行模式中，节点之间的通信是不规则的，因此这种模式对于计算量不均衡，并且节点之间通信量小的算法非常适用。

### 5.4 图的划分

图的划分是进行高效图计算的一个关键问题。通常，一个理想的图划分情况是各工作节点的任务量基本相同，同时各工作节点之间的通信量最小，但是这是一个NP难的问题。现在，常用的图划分算法分为3类。

第一类，首先对输入的图数据进行一次预处理，将初始的图数据转化为某个特定的存储格式，使得图计算的访存局部性更好或者使图数据的数据量占用更少。比如GraphChi使用shard以及shard内存源点的排序来增强磁盘访存的局部性。另外，X-Stream使用简单的流划分来降低预处理的开销。

第二类，在算法执行过程中使用动态的重划分，由于算法在执行之前行为是无法预测的，所以这种动态划分的策略可以根据现有算法的执行状态进行相应地划分，提高系统的性能。这种动态划分策略需要对图进行多次划分，引入了图划分开销。

第三类，使用edge-cut和vertex-cut划分。edge-cut将图中的点均匀地划分，并且保证跨不同划分块之间的边最少。vertex-cut将边均匀地划分，同时保证跨不同块之间的点最少。现实生活中的许多大图符合幂律分布<sup>[27]</sup>，因此，相比于edge-cut，使用vertex-cut有助于系统的负载均衡，但是图计算系

统需要使用以边为中心的计算模型，如PowerGraph。

### 5.5 负载均衡

负载均衡的算法分为静态负载均衡和动态负载均衡，静态负载均衡在算法执行之前进行任务的分配，但是由于算法在执行之前无法预测其具体的行为，因而在算法的执行过程中可能出现负载不均衡的情况。动态的负载均衡策略针对静态负载策略进行了改进，即在算法的运行过程中，系统中任务少的工作节点可以从任务量大的工作节点“偷取”任务来实现负载均衡，提高系统的整体性能。

### 5.6 容错

容错在分布式图处理系统中是需要解决的一个问题。在分布式处理系统中，每台机器都会有一定的概率出错失效，如果不加以处理，将对系统产生严重的影响。常见的分布式图处理系统使用主从节点的方式，在这种构建方式中，主节点负责整个系统的管理和调度，从节点负责具体的计算。主要的容错方式有多副本策略、日志重做策略等。在多副本策略中，当主工作节点执行其任务时，另外有一个工作节点作为副本工作节点会执行相同的任务；当主节点失效时，副本会接管主节点的工作任务，这种容错方式基本没有错误恢复时间，但是会消耗掉很多计算和内存资源。在日志重做的策略中，使用checkpoint或者log的方式记录工作节点的计算操作，当机器出现失效时，可以将记录的操作重做来进行恢复，这种恢复方式会消耗一定的恢复时间，但是对计算和内存资源的消耗相对较少。

## 6 结论及未来研究方向

本文介绍了几个典型的分布式大图处理系统和单机大图处理系统,这两种类型的系统有着各自的优点和缺点。对于分布式系统,其特点是计算能力强,能够应对不同的计算需求,但是编程模型和系统的构建(计算的协调和容错机制)比较复杂;对于单机系统,其特点是编程和计算模型简单,硬件开销很低,但是计算能力有限,无法满足某些计算需求。从计算模型来看,现在大图计算的计算模型主要分为两种:以点为中心的计算模型和以边为中心的计算模型。在分布式处理系统Pregel、GraphLab等以及单机系统GraphChi主要使用了以点为中心的计算模型,这种计算模型更易于编程和理解,以边为中心的计算模型主要用于单机的系统,如X-Stream。除了这两种主要的计算模型之外,还有一些系统从数据的局部性出发,提出一些新的计算模型来提升系统的性能,但从本质上来说,这些计算模型是基于以点为中心的计算模型,只是针对数据的布局,做出了相应的修改。

尽管现在有许多针对大图计算系统的研究工作被提出,但是从系统角度来看,在大图处理系统上还有许多值得深入研究的领域。在分布式图计算系统方面,设计一套高效、合理的图划分策略,不仅可以减少集群中各节点的通信开销,而且可以保证机器间的负载均衡,在这方面已经有一些相关的研究,但仍然值得更深入的研究。另外,容错也是分布式系统改善性能的一个重要方面,现在主要的容错方法有主副本备份容错、校验点容错等,目的是在减少容错开销的同时尽可能地提高错误恢复的速度。在单机图计算系统

方面,由于计算能力的限制,有效的图划分策略并且使用与划分策略相匹配的计算模型来增强计算的局部性是研究的热点。另一方面,应该充分发挥机器的多核特点,使得I/O和计算并行,并且提高计算时的并行度,这两点也是值得深入研究的方向。

## 参考文献

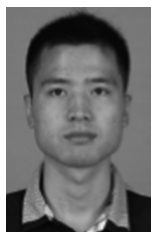
- [1] Lumsdaine A, Gregor D, Hendrickson B, *et al.* Challenges in parallel graph processing. *Parallel Processing Letters*, 2007, 17(1): 5~20
- [2] Dean J, Ghemawat S. MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 2008, 51(1): 107~113
- [3] Gregor D, Lumsdaine A. The parallel BGL: a generic library for distributed graph computations. *Proceedings of Parallel Object-Oriented Scientific Computing (POOSC)*, Glasgow, UK, 2005
- [4] Chan A, Dehne F, Taylor R. CGMGRAPH/CGMLIB: implementing and testing CGM graph algorithms on PC clusters and shared memory machines. *International Journal of High Performance Computing Applications*, 2005, 19(1): 81~97
- [5] Malewicz G, Austern M, Bik A J C, *et al.* Pregel: a system for large-scale graph processing. *Proceedings of ACM Special Interest Group on Management of Data*, Indianapolis, IN, USA, 2010: 135~146
- [6] Low Y C, Bickson D, Gonzalez J, *et al.* Distributed GraphLab: a framework for machine learning in the cloud. *Proceedings of the VLDB Endowment (PVLDB)*, 2012, 5(8): 716~727
- [7] Gonzalez J E, Low Y C, Gu H J, *et al.* Power graph: distributed graph-parallel computation on natural graphs.

- Proceedings of the 10th USENIX Symposium on Operating Systems Design and Implementation, Hollywood, CA, USA, 2012: 17~30
- [8] Gonzalez J E, Xin R S, Dave A, *et al.* Graphx: graph processing in a distributed dataflow framework. Proceedings of the 11th USENIX Symposium on Operating Systems Design and Implementation, Broomfield, CO, USA, 2014: 599~613
- [9] Chen R, Ding X, Wang P, *et al.* Computation and communication efficient graph processing with distributed immutable view. Proceedings of High-Performance Parallel and Distributed Computing, New York, USA, 2014: 215~226
- [10] Yan D, Cheng J, Lu Y, *et al.* Blogel: a block-centric framework for distributed computation on real-world graphs. Proceedings of the VLDB Endowment (PVLDB), 2014, 7(14): 1981~1992
- [11] Yuan P P, Zhang W Y, Xie C F, *et al.* Fast iterative graph computation: a path centric approach. Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, Piscataway, NJ, USA, 2014: 401~412
- [12] Kyrola A, Blelloch G, Guestrin C, *et al.* GraphChi: large-scale graph computation on just a PC. Proceedings of the 10th USENIX Symposium on Operating Systems Design and Implementation, Hollywood, CA, USA, 2012: 31~46
- [13] Roy A, Mihailovic I, Zwaenepoel W. X-stream: edge-centric graph processing using streaming partitions. Proceedings of ACM Symposium on Operating Systems Principles, Farmington, PA, USA, 2013: 472~488
- [14] Cheng J F, Liu Q, Li Z G, *et al.* VENUS: vertex-centric streamlined graph computation on a single PC. Proceedings of the 31st IEEE International Conference on Data Engineering, Seoul, Korea, 2015: 1131~1142
- [15] Zhu X W, Han W T, Chen W G. Grid graph: large-scale graph processing on a single machine using 2-level hierarchical partitioning. Proceedings of the 2015 USENIX Conference on Usenix Annual Technical Conference, Santa Clara, CA, USA, 2015: 375~386
- [16] Valiant Leslie G. A bridging model for parallel computation. Communications of the ACM, 1990, 33(8): 103~111
- [17] Low Y C, Gonzalez J, Kyrola A, *et al.* GraphLab: a new framework for parallel machine learning. Proceedings of Conference on Uncertainty in Artificial Intelligence, Catalina Island, California, USA, 2010
- [18] Barabasi A L, Albert R. Emergence of scaling in random networks. Science, 1999, 286(5439): 509~512
- [19] Gharaibeh A, Costa L B, Santos-Neto E, *et al.* On graphs, GPUs, and blind dating: a work load to processor matchmaking quest. Proceedings of IEEE the 27th International Symposium on Parallel and Distributed Processing, Washington DC, USA, 2013: 851~862
- [20] Fu Z S, Personick M, Thompson B. MapGraph: a high level API for fast development of high performance graph analytics on GPUs. Proceedings of Graph Data-management Experiences & Systems, Utah, USA, 2014: 1~6
- [21] Khorasani F, Vora K, Gupta R, *et al.* CuSha: vertex-centric graph processing on GPUs. Proceedings of the International ACM Symposium on High-Performance Parallel and Distributed Computing, Vancouver, Canada, 2014: 239~252
- [22] Han W S, Lee S, Park K, *et al.* TurboGraph: a fast parallel graph engine handling billion-scale graphs in a single PC. Proceedings of the 19th ACM SIGKDD Conference on Knowledge Discovery and

Data Mining, Chicago, USA, 2013: 77~85  
 [23] Zheng D, Mhembere D, Burns R, *et al.*  
 FlashGraph: processing billion-node  
 graphs on an array of commodity

SSDs. Proceedings of the 13th USENIX  
 Conference on File and Storage  
 Technologies, Santa Clara, CA, USA,  
 2015: 45~58

### 作者简介



吴城文, 男, 清华大学计算机科学与技术系硕士生, 主要研究领域为大数据图计算。



张广艳, 男, 博士, 清华大学计算机科学与技术系副教授, 中国计算机学会会员, 主要研究领域为大数据计算、网络存储、分布式计算。



郑纬民, 男, 清华大学教授、博士生导师, 中国计算机学会理事长, 目前主要从事并行与分布式计算、存储系统的研究工作, 主持和参与多项国家“973”计划、“863”计划、国家自然科学基金项目。近年来在IEEE TC/IEEE TPDS/ACM TOS/FAST等本领域顶级期刊与国际会议发表论文40余篇。

收稿日期: 2015-08-19

基金项目: 国家重点基础研究发展计划(“973”计划)基金资助项目(No.2014CB340402), 国家自然科学基金资助项目(No.61170008, No.61272055)

**Foundation Items:** The National Basic Research Program of China(973 Program)(No.2014CB340402), The National Natural Science Foundation of China(No.61170008, No.61272055)

论文引用格式: 吴城文, 张广艳, 郑纬民. 从系统角度审视大图计算. 大数据, 2015028

Wu C W, Zhang G Y, Zheng W M. Reviewing large graph computing from a system perspective. Big Data Research, 2015028