



北京大学

## 硕士研究生学位论文

题目： HybriG: 一种高效处理大量重  
边的属性图存储架构

姓 名： 黄权隆  
学 号： 1401214283  
院 系： 信息科学技术学院  
专 业： 计算机软件与理论  
研究方向： 分布式系统  
导 师： 崔斌教授

二〇一七年六月



## 版权声明

任何收存和保管本论文各种版本的单位和个人，未经本论文作者同意，不得将本论文转借他人，亦不得随意复制、抄录、拍照或以其他方式传播。否则一旦引起有碍作者著作权之问题，将可能承担法律责任。



## 摘要

在图中，起点和终点都相同的两条边称为重边。属性图是一种带标志和重边的有向图，图中的点和边可以拥有任意数目的属性值。属性图由于其丰富的表达能力而广泛应用于实际建模中。实际应用中一般用图数据库解决属性图的存储需求。相比于传统的关系型数据库，图数据库在做多跳邻域查询、路径查询等与图结构相关的查询时，具有更优异的性能。Titan是产业界日渐关注的一个开源的分布式图数据库，Titan的数据以邻接表的方式组织，每个点的邻接表存储了相邻的所有边，这使得与邻接点集相关的查询都需要遍历整个邻接表。当图中含有大量重边时，邻接表规模巨大，这种数据组织方式导致邻域查询性能严重受损。而邻域查询是大部分图查询的基础，如多跳邻域查询、路径查询、局部聚集系数查询（计算）等，这些查询往往由嵌套的邻域查询实现，随着邻域深度的增加，这种性能受损将被急剧放大。本文提出了一种基于Titan和列式存储数据库HBase的复合架构设计——HybriG，基于Titan和HBase建立存储层，用Titan来存储图的结构信息和点集的属性信息，HBase存储边集的所有属性信息。在HybriG中邻接表保持了项数和数据量上的精简，从而能克服上述图数据库的缺点。相比于传统图数据库Titan，HybriG在邻域点集相关查询以及边集数据批量导入上的性能提升一个量级以上。本文介绍了HybriG基于Titan和HBase的存储设计，并描述了在此存储设计基础上，如何高效地实现图查询以及图数据的插入操作。此外，本文还提出了图数据的高效导入方案，并保证导入过程中Titan与HBase存储数据的一致性。最后通过实验验证了HybriG在处理大量重边时的优异性能。

**关键词：**属性图，重边，图数据库，Titan，HBase，架构设计



# HybriG: A Distributed Architecture for Property Graph with Large Set of Multi-edges

Quanlong Huang (Computer Science)

Directed by Prof. Cui Bin

## ABSTRACT

The past decades have witnessed the massive growth of Internet. Vast amount of graph data was produced by the boom of social network, e-commerce and online education. To analyze and manage graph data, there are two branches of development. One focus on distributed graph computing. Solving problems like PageRank or other algorithms that fit into the Bulk Synchronous Parallel (BSP) model. Systems like Pregel and GraphLab are proposed for this branch. The other branch focus on management of graph data, providing support like OLTP and graph queries. In this branch, graph databases like Neo4j and Titan are developed for management of property graph.

The property graph is a directed, labeled graph with multi-edges, i.e., edges with the same source vertex and destination vertex. Vertices and edges can be associated with any number of properties. Since it can represent graph data in most scenarios, the property graph model has numerous applications in industry.

However, traditional graph databases encounter significant performance degradation when the graph contains large amount of multi-edges. We reveal the cause of this in Titan. It's an open source distributed graph database, which has attracted wide attention in industry. Furthermore, we propose HybriG, a better distributed architecture based on Titan and HBase for this scenario.

Titan stores graphs in adjacency list format, where a graph is stored as a collection of vertices with their adjacency lists. Each entry of the adjacency list stores an edge or a vertex property. When querying about adjacent vertices, Titan has to look through the entire adjacency list of the source vertex, which is a waste since we just need the connected vertices but not the multi-edges. This cost hurts the performance when the multi-edge set grows explosively. For high level queries who bases on adjacent vertices query, e.g. path queries, local cluster coefficient queries etc., the situation will be worse. HybriG implements property

graph model as well. It stores the vertex data and graph structure in Titan, the edge data in HBase, respectively. We chose HBase since it's one of the storage engines of Titan and is widely used in industry. Storing part of the graph data into HBase won't bring too much cost because it's also the data store of Titan. This separation helps HybriG to keep a concise adjacency list about the graph structure, which helps to gain an order of magnitude improvement in execution of adjacent vertices based queries and batch loading of edge set. The difficulty of this separation solution is that we should guarantee the consistency of edge data between Titan and HBase. In most of the scenarios with large multi-edges set, multi-edges are representing event-like data, e.g. phone call records between two people, which won't be modified after insertion. Thus we can relax the consistency constrain of edge data to just guarantee consistency for insertion. We leverage the transaction in Titan to achieve consistency with HBase in edge insertion, especially in batch loading of edge data. Finally, extensive experiments have been conducted to show the outstanding performance of HybriG.

**KEYWORDS:** property graph, multi-edge, graph database, Titan, HBase, architecture



# 目录

引言	1
第一章 相关工作	3
第二章 预备知识介绍	5
2.1 HBase简介 . . . . .	5
2.2 图数据库Titan . . . . .	5
2.3 图数据库Titan的局限性 . . . . .	6
第三章 HybriG的系统架构概述	7
第四章 Storage Layer	9
第五章 Query Adapter	11
5.1 查询设计 . . . . .	11
5.2 查询实现 . . . . .	11
5.2.1 仅依赖Titan数据的查询 . . . . .	12
5.2.2 关联Titan和HBase数据的查询 . . . . .	12
5.3 HybriG的查询优势 . . . . .	12
5.3.1 邻域点集相关查询的优势 . . . . .	12
5.3.2 边集相关查询的优势 . . . . .	13
第六章 Data Loader	15
6.1 边数据的高效导入 . . . . .	15
6.2 数据一致性 . . . . .	16
第七章 实验和分析	19
7.1 邻域点集相关查询 . . . . .	19
7.1.1 k-hop点集查询 . . . . .	19
7.1.2 局部聚集系数计算 . . . . .	20
7.2 边集相关查询 . . . . .	20
7.2.1 两点间边集查询 . . . . .	21
7.2.2 邻接边集查询 . . . . .	21
7.3 边数据的导入速度 . . . . .	22

结论	23
参考文献	25
附录 A 附件	27
致谢	29
北京大学学位论文原创性声明和使用授权说明	31

## 引言

图是一种表达能力丰富的数据结构，一般用 $G = (V, E)$ 表示一个图，其中 $V$ 是图中的点集， $E$ 是图中的边集。如果图中的边是有方向的，则称该图为有向图。在有向图中，如果两条边具有相同的起点和终点，则称这样的边为重边（multi-edge）。属性图[1]是一种带重边的有向图，图中的每个元素（点/边）可以拥有任意数量的属性值。特别地，每个元素都有一个label属性，标识该元素的类别。实际应用中，有些属性图会含有大量重边。比如在电信行业的通话关系图中，两人之间可以发生多次通话关系；又比如金融行业的交易关系图中，两人之间可以发生多次交易关系。更复杂的，在刑侦应用中，整个情报系统获得的信息融合成一个知识图谱[2]，其中的实体和关系类型丰富多样。比如实体类别有人、汽车、公司等，实体间的关系有亲属关系、共同出行关系、通话关系等。有的关系是可以多次重复的，频次可能成百上千，这使得两个实体间不仅有多种label（即多种关系）的边，还包含大量重边。图1是属性图在刑侦场景中的一个简单示例，显示了类型为人的结点以及通话关系和共同出行关系。其中共同出行关系是双向的，因此在图中以无向边表示。刑侦场景中的典型查询包括：从一个人出发查询与其相关的所有人（邻域点集查询）、查询两个人在几跳内是否有关系（路径查询）、查询两人之间的所有关系（两点间边集查询）等。传统解决方案将数据存储存储在关系型数据库中，如将所有通话数据存储存储在通话关系表中。关系型数据库能胜任图数据的存储，也能方便地检索元素的内容，但在图结构相关的查询上表现欠佳。比如两跳邻域查询，锁定一个嫌疑人后，要查看其两跳范围内的子图信息。为得到第二跳的邻域结点，需要对各种关系表做昂贵的JOIN操作。图数据库由于其以图的形式存储数据，在图的拓扑结构查询中具有更优异的性能。因此当应用场景中图的拓扑结构查询与元素的内容查询同等重要时，图数据库是最佳的选择。明略数据是一家大数据解决方案提供商，客户群涵盖刑侦、金融、电信等领域。其中SCOPA 是明略数据的重要产品，在海量刑侦数据上构建起一个数据分析挖掘平台，展现给领域专家的是一张属性图，关联了客户原有的所有数据。SCOPA平台即使用图数据库实现其存储核心。然而，在富含重边的属性图中，图数据库的查询性能不佳。在传统的图应用中，图中的边一般是稀疏的。如一些公开的图数据集中，边集的大小约为点集大小的十几倍。但在富含重边的应用场景中，边集大小可能为点集的上万倍。图数据库以邻接表的方式存储图，表中的每一项存储一条边。这使得查询邻接点集时，需要遍历整个邻接表中的所有边。当图中含有大量重边时，邻接表规模巨大，这种数据组织方式导致邻域查询性能严重受损。而邻域查询是大部分图查询的基础，如多跳邻域查询、路径

查询、局部聚集系数查询（计算）等，这些查询往往由嵌套的邻域查询实现，随着邻域深度的增加，这种性能受损将被指数级放大。一种简化重边的建模方式是，将所有同label的重边合并为一条边来表示，边上存储这些重边的所有信息，这样两个实体间的重边数能降低为边的label种类数。然而，将多条重边的数据合并一条边中存储，会使边的单位数据量骤增，图中的数据粒度过大，同时降低了单条重边的检索和插入效率，每次操作都要先读取相关的所有重边数据，再在其中做查询或插入操作。因此，处理大量重边是一个不可避免的需求。在许多富含重边的图应用场景中，尽管数据量较大，但数据的操作需求相对单一，而且没有很强的事务性要求。比如在知识图谱的构建中[2]，信息的来源是可靠的，即一旦信息进入知识图谱，则被认为是无需修改的，因此数据操作只需要插入和查询，更新操作很少被执行。在SCOPA的刑侦应用场景中，数据来源是可靠的，且图中数据都是一些客观的事实数据，不需要修改，数据操作主要是图查询和批量的图数据插入。数据插入操作包括原始的全量数据导入和定期的增量数据导入，没有并发的写冲突，也不需要很强的事务性要求。基于上述观察，在放宽了对强事务的支持后，我们可以设计一个相比传统图数据库更为高效的属性图处理系统，来应对含有大量重边的应用场景。本文提出了一种基于图数据库Titan和列式存储数据库HBase的复合存储架构——HybriG。HybriG能从容地处理富含重边的属性图，克服图数据库的局限性，具有更好的查询和插入性能。HybriG架构应用于SCOPA的底层核心设计，在实践中也验证了其性能的优越性。本文的贡献主要有：分析了图数据库处理大量重边时性能受损的原因，并基于此提出了解决方案HybriG及其存储、查询和高效数据导入设计，最后通过实验验证了HybriG的优异性能。本文组织如下：第2节介绍图数据库相关的预备知识，包括HBase及Titan的实现及其在处理大量重边时的局限性；第3-6节介绍HybriG的系统架构，包括查询模块的设计，数据高效导入方案的设计，以及数据一致性的设计；第7节展示实验结果。第8节介绍大规模图数据处理的相关工作。最后在第9节对本文进行总结。

## 第一章 相关工作

图数据分析与处理是大数据背景下的一大应用分支[10]，大数据背景下的图处理系统可以分为图计算框架和图数据库两大类。图计算框架旨在高效地进行全图量级的并行计算，如计算PageRank[11]、社区发现[12] (Community Detection)、子图匹配[13]等，提供的是对OLAP (Online Analytical Processing) 的支持。这类系统中，Pregel[14]、GraphLab[15]、PowerGraph[16]支持BSP[17]模型的计算，对图数据进行特定的分割，使得集群中的机器可以在内存中计算分区的数据，再通过交互完成全量数据的计算。上述系统中，Pregel没有开源，GraphLab和PowerGraph都需要单独部署，无法有机地融入已有的大数据生态系统中，需要导出数据以供计算。GraphX[18]解决了这个问题，其是基于分布式内存计算框架Spark[19]实现的图计算引擎，使得图计算可以融入整个数据流处理的框架。相对于分布式计算框架，GraphChi[20]则追求在单机完成大规模图数据的计算问题。在图计算领域还有许多研究旨在提高图计算框架在特定问题的处理性能，如[21-25]等。学术界对图计算的研究热情普遍高涨。图数据库专注于图数据的管理，提供高效的图遍历查询 (graph traversal)。图数据模型则是指图数据库如何以图的方式对数据进行抽象 [26]，应用较广泛的图数据模型有RDF 模型和属性图模型。RDF全称为Resource Description Framework，是由W3C制定的知识描述标准，使得按RDF表示的不同数据源可以进行数据交换或合并。RDF中的数据单元是由主语、谓语、宾语组成的三元组，可以直接对应为图上的一条边，因此将RDF模型归类为图数据模型。常见的RDF存储有Jena 、 AllegroGraph 等。当今的图数据库大都采用属性图模型设计[27]，如DEX[28] 、 GraphChi-DB[29]、Neo4j、Titan等，其中DEX和GraphChi-DB都只是单机数据库，Neo4j的开源版本也是单机的，只适合处理中等规模的图。Titan的底层实现了可插拔的存储接口，可部署在HBase、Cassandra或BerkeleyDB之上，因此可以很好地与Hadoop集群结合，组成统一的数据处理框架。然而，现有的图数据库都没有考虑含有大量重边的属性图应用场景。HybriG架构填补了这个空白，为含有大量重边的属性图应用场景提供了一种可行的解决方案。



## 第二章 预备知识介绍

### 2.1 HBase简介

HBase是Hadoop 生态系统中一个列式NoSQL数据库，基于BigTable[3]模型设计，搭建在HDFS（Hadoop分布式文件系统）之上。BigTable是Google在2008年提出的一个处理海量数据的NoSQL数据库。图 2是BigTable的模型示例。在BigTable的模型中，数据表以行为单位组织，每行由一个键值唯一标识，称为行键（rowkey）。一行中可以包含任意数目的单元格（cell），每个单元格由列名、版本号（时间戳）和值（value）组成。BigTable中的行是以行键排序的，每一行中的单元格又是以列名和版本号进行排序，这使得其与关系数据库一样能快速定位到目标单元格。区别于传统关系模型，BigTable中的列名不需要事先定义，因为其底层实际为一个Key-Value 存储，就如数据结构Map不需要事先定义所有key。每个列名从属于一个列族，只有列族需要在定义表结构时给出。HBase提供对某一行数据的Put、Get、Delete接口，能够具体操作该行中的给定列。另外，HBase提供给定行键范围的Scan接口，可以快速地获得连续几行的数据。Scan接口也可指定特定的列或附加Filter来过滤无关数据。HBase由于其优异的可扩展性和稳定性，在产业界得到了广泛应用。

### 2.2 图数据库Titan

图数据库是以图的形式来表示和管理数据的数据库[4]，与传统的关系型数据库相比，图数据库在图结构相关的查询上有更优异的性能，如多跳邻域查询、路径查询、局部聚集系数计算等。Titan是一个基于Blueprints 接口设计的开源图数据库，其实现了一个可插拔的存储接口，可以部署在BerkeleyDB、HBase或Cassandra[5]之上。相比于著名的图数据库Neo4j[6]，Titan是完全开源免费的，其受关注度正在与日俱增。而且由于Hadoop生态系统在产业界的广泛应用，在HBase上搭建Titan，即使用Titan on HBase应对图处理需求是较为常见的选择。在图数据库Titan中，基于BigTable模型，数据在HBase中以邻接表的形式存储，如图3所示。Titan为每个点分配了一个全局唯一的id，每个点占BigTable中的一行，行键就是点的id。每行存储了该点相关的属性和边，它们各占一个单元格。Titan在BigTable模型中设置最大版本数为1，从而使每行中的列名与单元格一一对应，根据行键和列名可以快速定位到目标单元格，实现对边和属性内容的快速检索。在Titan中，每个属性是一个key-value对。点的属性存储在该点所在的行，每个属性占一个单元格，并以属性名key作为列名，这使得每个点的

属性查询可以非常高效。每个点所在的行还存储了邻接的所有边数据，每条边占一个单元格。一条边的信息包含了邻接点、类别 (label)、方向、边的唯一id，以及边上的各属性。单元格的值用来存储边上的所有属性，列名则存储边上除属性外的其它信息。在BigTable模型中，同一行的单元格按列名排序。为了方便检索，每条边所在的单元格依次以label、方向、邻接点id以及边id拼接成为列名，即列名 = 边label + 方向 + adjVid + eid 这种列名设计使得一个点的所有边先按label排序，再按方向、邻接点id、边id等排序。其优点是当要检索该点指定label的所有边时，只需要扫描邻接表的一部分。图4是一个更详细的示例。

## 2.3 图数据库Titan的局限性

当图中的重边数量巨大时，Titan的邻接表列数也急剧变大，这会对邻域中点和边的检索带来严重影响。图4展示了一个富含重边的场景，为方便展示省略了边的方向。v1只有v2和v3两个邻接点，其中与v2有两种label的边，与v3有一种label的边。虽然邻接的点数不多，但v1与邻接点都有数量巨大的重边。由于邻接表存储的是每条边的信息，当要查询v1的所有邻接点集（即v1、v2）时，不得不遍历一次v1的所有边，即遍历整个邻接表。当重边数量巨大时，这是大量的无谓开销。为了规避上述情形，我们可以换一种列名设计来优化邻域点集查询，比如让邻接表先按邻接点id排序，令列名 = adjVid + 边label + 方向 + eid 这样当发现一个邻接点时，可以跳过相同邻接点的所有边，不用再遍历整个邻接表。然而，面对label相关的查询时又会面临同样的问题，比如查询该点总共有几种label的边，仍需要遍历该点的整个邻接表。因此，改变邻接表的列名设计并不能解决问题，本质原因是邻接表存储了所有的边集。另一方面，Titan为加快数据访问设计了缓存，存储最近访问的点及其邻接表内容。如果图中各点的重边都很多，则缓存空间被大量的重边占满。由于是重边，这些边中只有为数不多的邻接点。在做邻域点集相关查询时就会极大增加缓存的失效率。面对富含重边的属性图，把所有边集都存在邻接表中是Titan性能受损的原因所在。



## 第三章 HybriG的系统架构概述

为了更好地应对含有大量重边的属性图应用场景，规避传统图数据库的局限性，本文提出了一种复合存储架构HybriG。HybriG架构应用在明略数据的关联数据分析平台SCOPA中，提供属性图数据的查询与存储。HybriG将图的连接信息和点集的属性数据存储在Titan中，边集的所有属性数据则直接存储在HBase中。由于HybriG架构基于Titan和HBase实现存储，而Titan的数据本身也存储在HBase中，为了方便表述以及避免混淆，下文中的HBase指的都是不包含Titan数据表的部分。Titan在HBase中的数据表直接用Titan代称。图 5展示了HybriG的系统架构。Storage Layer是HybriG的存储层，Query Adapter和Data Loader分别为用户程序提供图查询和数据插入接口。其中，Query Adapter将图查询转换为对Titan和HBase数据的高效查询，再将结果汇总返回给用户程序。Data Loader实现了数据的高效导入，在面对大批量的数据导入时，实现了错误恢复、断点恢复机制，同时保证了Titan和HBase的数据一致性。下面将分别介绍HybriG各模块的实现。



## 第四章 Storage Layer

HybriG将图数据分开存储在Titan和HBase中。具体地，如果两点之间有相同label的重边，HybriG会在Titan中这两点间建立一条该label的边，将对应重边的数据都存储在HBase中的一个表，不妨称其为边表。每条边一行，行键是Titan图中分配的边id拼接上重边的主键，单元格的值则存放具体边的内容。重边的主键是这样定义的：在两点间同label的重边中，如果有一个属性能唯一确定一条边，则选该属性作为主键。比如交易关系的时间戳，两个人在同一个时间点只能发生一次交易，因此该时间戳可以唯一确定两人间的一次交易。如果某种label的边没有属性能唯一确定一条重边，则选该数据插入的时间戳作为主键。在HBase边表中，同label的重边数据的行键都有相同的前缀（即Titan中的边id），由于HBase中的数据是按行键排序的，它们在HBase中有相邻的位置。同label的重边经常会被同时检索，这种设计使得一次顺序扫描便可得到所有同label的重边。另外，HybriG利用Titan中边上的属性记录一些统计信息，如原图中实际有几条这样的重边，或者边上某个属性值的求和等，这些统计信息可以根据业务定制，用来加快业务相关的查询，不用再遍历相关的边。图6是数据存储的一个示例，左边是实际要存储的图数据， $e_1$ 、 $e_2$ 、 $\dots$ 、 $e_n$ 都是label为“交易”的重边，表示两人的所有交易记录。右边是数据在Titan和HBase中的存储情况。Titan中只存一条label为“交易”的边，该边有一个属性记录实际的边数。利用这条边的id作为前缀，拼接上重边的主键（交易的时间戳）作为HBase的行键，将每条边的数据都存入对应的单元格中。本文在引言中提到了一种避免产生重边的建模方式，即将所有同label的重边合并为一条边来表示，边上存储这些重边的所有数据。HybriG架构与其有相似之处，但本质区别是边集元素的数据粒度仍是每条重边，不会将多条重边的数据作为一个单元来处理。在HybriG中每条边的数据独占HBase边表中的一行，因此每条边的相关操作都会更加高效。得益于HBase的高可扩展性，这种存储方式的性能也不会受制于数据规模的增长。而前述建模方式将所有重边的数据合并在一边中存储，当重边量级巨大时，对每条边的操作势必更加低效。HybriG采用的是合并同label的重边，而不是所有的重边，这仍是基于数据粒度的考虑。图7是合并相同label重边后Titan的图示例（在HBase边表中的重边数据未显示）。在实际应用场景中，经常需要根据具体的关系种类（label）来进行邻域点集查找，对于具体边数据的查询也常按label进行。因此合并相同label的重边设计可以使很多查询在Titan中就得到满足，不需要再涉及HBase边表。



## 第五章 Query Adapter

Query Adapter模块将图查询转化为对应Titan和HBase的查询，再将结果汇总返回。下面分别叙述其查询设计、查询实现以及查询优势。

### 5.1 查询设计

我们在HybriG架构上实现了基本图查询接口，暴露给上层应用的仍是一张属性图。基本的图查询接口可分为如下几类：

- 以点为中心的查询：获取某点邻接的点、获取某点的属性、获取某点邻接的边等；
- 以边为中心的查询：获取某条边的端点、获取某条边的属性等；
- 从图中直接查询：查询符合某种条件的元素（点/边），其中的条件可以是label、限定返回集大小（limit）等。

更详细的属性图查询接口，可以查看Blueprints中的定义。

除了基本的属性图查询接口，HybriG还提供了重边的统计信息查询。比如查询两点间某种label的重边具体的边数、或者查询某个属性上的聚集信息（如MIN、MAX、AVG、SUM等），可以根据业务需要来定制具体统计的信息。比如通话关系所表示的边中，有记录通话开始时间戳的属性。可以设置HybriG统计该类重边在该属性的最小值，从而可以查询两人最早的通话发生时间。在HybriG架构中，这些统计信息存储在Titan中的边上，因而可以快速返回结果，无需再对具体的边数据（存在HBase边表中）进行统计。区别于传统Titan图数据库，我们没有实现事务性的接口，如newTransaction、commit和rollback等。因为业务场景可以规避对事务性的依赖，在富含重边的应用场景中，边集数据都是事件记录类型的数据，数据导入后就无需修改。而且数据导入可以分批定期执行，不会有多数据源写入导致的写-写冲突或读-写冲突。另外HybriG的数据分开存储在Titan和HBase两个数据库中，实现严格事务的代价比较大，会带来显著的性能牺牲。如Google的Percolator[7]在BigTable上实现了分布式事务，但写性能有75%的牺牲。

### 5.2 查询实现

查询的实现可分为两部分，一部分只依赖于Titan中的数据，另一部分需要联合Titan和HBase边表的数据来返回结果。下面分别叙述。

### 5.2.1 仅依赖Titan数据的查询

在HybriG架构中，点集和各点的邻接信息都存储在Titan中，因此只与点集相关的查询都可以直接转换为对Titan的查询，如点上属性的查询、查询给定点的邻域点集、在图中查询某种label的点等。对于重边的统计信息查询，其需要的统计结果都已在Titan的边中存储，故可以直接转换为对Titan边上的属性值查询。具体实现中需要维护一个映射关系作为元信息，以得知一个统计信息对应Titan边上的哪个属性。

### 5.2.2 关联Titan和HBase数据的查询

当查询具体的边数据时，就需要关联Titan和HBase来实现查询。在HybriG架构中， $\text{HybriG\_Edge} = \text{TitanEdgeID} + \text{HBaseRowData}$  上式表示每条边对象的两部分组成，所有边集数据存储在HBase的边表中，每条边的数据占据一行，存储其所有属性，该行的行键是Titan中对应的边id拼接上该边的主键值。因此查询某条重边的数据，需要先在Titan中找到对应边的id，再在HBase边表中找到对应行的数据。下面以两点间边集查询为例阐述HybriG的实现。两点间边集查询是指给定两个点，查询它们之间满足某种条件的边。比如已经从邻域查询得知v1与v2相邻，现在想得到它们之间某种label的所有边。查询的伪代码见算法1。算法1. 两点间给定label的边集数据查询伪代码。输入：点v1，点v2，eLabel 输出：两点间的所有符合条件的边集数据

```

1. e = v1.query().adjacent(v2).label(eLabel).limit(1).edges().next()
2. IF e == null THEN
3.   RETURN null
4. END IF
5. res = hbaseTable.scan(e.id, next(e.id))
6. RETURN wrapEdges(res)

```

算法1中，第1行查询两点间该label的边，并利用limit(1)修饰来加速，在HybriG架构中，Titan只在这两点间存储该label的一条边，因此我们可以利用边的唯一性来加速查询。第2-4行，若在Titan中查询到该label的边不存在，则不需要再在HBase中检索。第5行根据Titan中的边id，在HBase的边表中通过Scan查询得到所有边的具体数据。e.id是一个字符串，伪代码中的next(e.id)表示同等长度的下一个字符串，即将字符串e.id中最后一个字符的值加一。比如e.id为“abbb”，则next(e.id)为“abbc”。HBase的Scan函数返回连续的若干行数据，接收的两个参数是行键范围的起始和结束点，是一个左开右闭的区间。用e.id和next(e.id)作为该区间的左右端点，即可查询到所有行键以e.id作为前缀的数据。最后第6行wrapEdges函数将这些数据转换为具体的边对象，每行一条边。

## 5.3 HybriG的查询优势

### 5.3.1 邻域点集相关查询的优势

为了解释HybriG在图查询上的优势，图8展示了HybriG在HBase中的表结构，包

括Titan自身的HBase数据表和HybriG特殊设计的边表。当属性图的重边数量爆炸性增长时，Titan的数据表并不会被影响，因为两点间同label的边至多只会存在一条，每个点的邻接表列数仍为不同label邻接的点数。另外，Titan中的边只存放统计信息，因此每个单元格（即每条边）的数据量实际很小。面对含有大量重边的属性图，各点的邻接表仍保持了列数和数据量上的精简。即使面对邻域点集查询的全表扫描，也能提供优异的性能。另一方面，得益于邻接表的精简，Titan的缓存因此可以存放更多点集数据。对于邻域点集相关查询，如多跳邻域（k-hop）查询、路径查询、局部聚集系数计算等，具体的实现往往由多个基本的邻域点集查询组成，缓存的命中率就显得尤为重要。传统的解决方案中使用Titan存储所有的重边，使得缓存中只能存储少量点集的数据，大部分空间被边集数据占据，而具体的边集数据在查询中并不相关。在HybriG架构中，Titan的缓存能充分保留更多的点集数据，从而进一步提高邻域点集相关查询的缓存命中率。综上，HybriG对邻域点集相关的查询具有很好的表现，后续的实验结果将展示具体的数据。

### 5.3.2 边集相关查询的优势

在许多刑侦推演场景中，往往只需要查询两点间拥是否拥有某种label的边，而不需要查看具体各个重边的数据。比如得知两人之间有共同住宿酒店的边相连后，基本可以断定两人认识，领域专家可以在图中继续推演出其他相关的人，后续有需要再展开这条关系，查看具体的各条酒店住宿记录。HybriG架构为这种场景提供了便利，上述场景相当于在HybriG架构的Titan图中进行游走（Graph traversal），当有需要时再展开某条边，在HBase边表中读回相应的边数据。对于边集数据的读取，HybriG将所有边的数据存储在HBase边表中，而且每条边占一行，这使得对边的检索是行级别的检索，即在表中查找一行。传统的解决方案将重边数据都存储在Titan中，边集数据存储在各点的邻接表，而每个点的邻接表在HBase数据表中占据一行，因此对边的检索是选定行后的列级别检索。HBase中行级别的检索要略优于列级别的检索，因此HybriG会略占优势。然而，HybriG对边的检索需要跨Titan和HBase两个系统，会多一次交互的开销。实验表明，这两方面功过相抵，HybriG的边集数据查询性能与Titan相差不大。HybriG架构在边集统计信息的查询或计算上会有很大优势。在HybriG架构中，Titan中某种label的边存在，代表原属性图中两点间拥有该label的边，具体的重边数据存储在HBase中的边表，而Titan中的边上的则存储了该label声明时设定的统计信息，如具体的重边数目、边上某个属性的最大最小值等。对于这些统计信息，可以直接在Titan中查询得到结果。





## 第六章 Data Loader

Data Loader模块负责图数据的导入，包括点的导入和边的导入。在HybriG架构中，点集数据都存储在Titan中，故点集数据导入直接在Titan上完成。边集数据虽然存储在HBase中，但Titan中的对应边上存储了重边的统计信息，因此需要同时更新Titan和HBase，以保证二者数据的一致。下面详述边的数据导入，以及如何保证Titan和HBase的数据一致性。

### 6.1 边数据的高效导入

在HybriG架构中，边的插入既要更新Titan，又要更新HBase中的边表。对于给定label的一条边数据，首先要查看Titan中两点间是否已有一条边表示这样的关系存在。若这样的边不存在，则将其创建。然后更新Titan边上的统计信息，再利用其边id将具体边的数据存入HBase中。算法2展示了上述逻辑。第1行查询Titan中该label的边，由于至多只有一条，故可用limit(1)加速。第2-4行若这样的边不存在，则在Titan中添加一条边。第5行根据要插入的边数据来更新Titan中边上的统计信息。第6行提交对Titan的所有修改。第7行将边数据写入HBase边表中，行键为e.id拼接上边的主键。算法2. 插入一条边的伪代码输入：Titan图接口graph，HBase边表接口hbaseTable，点v1，点v2，edgeLabel，边数据realEdge 1. e = v1.query().adjacent(v2).label(edgeLabel).limit(1).edges().next() 2. IF e == null THEN 3. e = v1.addEdge(v2, edgeLabel) 4. END IF 5. updateStats(e, realEdge.properties) 6. graph.commit() 7. hbaseTable.put(e.id + realEdge.primaryKey, realEdge.properties) 传统的解决方案将所有重边存储在Titan中，边数据的导入逻辑只需要上述代码的第3、5、6行。HybriG架构中每条边的插入多加了一次查询操作（第1行），以及HBase边表的插入操作（第7行），带来了额外的时间开销。而且最耗时的也是这两步操作，分别要从底层存储读取数据以及持久化所有修改到底层存储中。如果有批量重边需要导入，这些开销可以让所有重边来平摊，即在数据导入时，对于两点间同label的重边作为一批来处理，这样就可以共享查询操作的结果，对Titan边上统计信息更新的持久化操作（commit）也只需进行一次。批量导入重边数据的伪代码如算法3所示。算法3. 重边数据的批量导入输入：Titan图接口graph，HBase边表接口hbaseTable，点v1，点v2，edgeLabel，重边数据集allRealEdges 1. e = v1.query().adjacent(v2).label(edgeLabel).limit(1).edges().next() 2. IF e == null THEN 3. e = v1.addEdge(v2, edgeLabel) 4. END IF 5. FOR realEdge IN allRealEdges DO 6. updateStats(e, realEdge.properties) 7. END

FOR 8. graph.commit() 9. FOR edge IN allRealEdges DO 10. hbaseTable.put(e.id + edge.pk, edge.properties) 11. END FOR 在算法3中，第1到8行是对Titan的更新，将两点间同label的重边作为一批处理，共享了对Titan的操作，从而平摊了额外的开销。另外第9到11行是对HBase边表的多次put操作，可以使用HBase的Bulk Loading 技术来进行加速。

## 6.2 数据一致性

在Titan的接口设计中，对图(graph)的初次操作将自动打开一个事务，执行graph.commit()时该事务提交，将事务里的修改持久化到底层存储中，Titan的事务保证了自身数据的一致性。HybriG架构由于把边数据的存储分开在Titan和HBase上，需要保证Titan和HBase上数据的一致性。比如Titan上关于重边的统计信息跟HBase边表的一致性，或者HBase边表里各行行键的边id部分跟Titan里的一致性。由于实现强一致性的代价过高，本架构保证的是Titan和HBase数据的最终一致性[8]，即系统保证在经过错误恢复后，数据最终会达到一致的状态。最终一致性足以满足应用场景的需求。由于只有边集数据是跨Titan和HBase存储的，下面讨论的都是边数据的导入。图 9演示了边数据导入的三个步骤，第①步对应前述算法3中的第1 7行，更新Titan数据；第②步对应第8行，提交修改；第③步对应第9 11行，利用边id将边数据插入HBase边表中。对数据的持久化修改是第②③步，数据会出现不一致的原因是②③并不是原子的。如果②成功但是③失败了，即成功持久化了对Titan数据的修改，但对HBase边表的修改却失败了，则两方的数据出现不一致。失败的原因是多种多样的，比如程序内存溢出（Out of Memory）、网络中断、硬件故障等。

图 9 边的数据导入分三步：①更新Titan数据；②提交Titan更改即graph.commit()；③将数据插入HBase边表中当数据导入程序在故障之后重启时，这种不一致性需要被修复。该批次的边数据将会被重新导入，为了得知②是否已经成功，我们需要知道Titan中的数据是否已经被修改了。这可以利用Titan的事务原子性来实现：在②提交之前，往Titan中写入一个成功标志（可以用一个点或属性来表示），然后再提交。根据原子性，该标志被成功写入当且仅当Titan的更新都被持久化。这样当数据导入程序重启时，我们先查看该标志是否存在，便可得知②是否已经成功了。若其已成功，则跳过这一步，直接进行第③步。该过程的伪代码如算法4所示。算法4. 保证一致性的边集数据批量导入输入：graph, batchID, HBase边表hbaseTable, 点v1, 点v2, edgeLabel, allRealEdges, 所有边的属性 1. e = v1.query().adjacent(v2).label(edgeLabel).limit(1).edges().next() 2. IF hasNotCompleted(graph, batchID) THEN 3. IF e == null THEN 4. e = v1.addEdge(v2, edgeLabel) 5. END IF 6. FOR realEdge IN allRealEdges DO 7. // 按需更新边上的统计信息

8. END FOR 9. markCompleted(graph, batchID) 10. graph.commit() 11. END IF 12. FOR edge IN allRealEdges DO 13. hbaseTable.put(e.id + edge.pk, properties) 14. END FOR

算法4中，第2行的函数hasNotCompleted判断给定的Titan图中是否存在给定批次的成功标志。若不存在则执行第3到10行更新Titan，其中第9行的markCompleted函数在Titan图中插入该批次的成功标志。值得一提的，我们在HBase中并没有放置成功标志。如果数据导入程序在第③步成功插入HBase数据后出现故障，则重启后还会将边集数据再次插入HBase中。但这是没有问题的，因为HBase边表中的数据不存在统计信息，因此不存在更新操作，所有操作都是新数据的插入操作。我们设置HBase的最大版本数为1，则多次插入同一内容到一个单元格，实际只保存一份，不会增加存储开销。



## 第七章 实验和分析

在现有的公开数据集（比如LDBC、SNAP、LAW）中，并没有富含重边的场景，这些图也不是属性图。富含重边的属性图普通存在于电信、金融、刑侦等行业中，数据是不公开的。由于实际的数据只能在相关部门内部查询，明略数据根据客户数据中统计出的特征构造测试用图，以支持SCOPA的开发与测试，验证HybriG架构的优秀性能。本组实验选用的图中有20万个顶点，48197700条边，平均度数为482，每个点平均与3.6个点相邻，两点间同label的重边的平均重数为134。实验对比的是直接将图存储在Titan中的传统方案以及将图存储在HybriG中的方案。实验环境为5台服务器组成的集群，每台机器安装Ubuntu14.04操作系统，物理配置均为一个Intel Xeon E3-1220 (3.10GHz)处理器、16GB内存、一个1Gbps网卡及一个4TB SATA接口硬盘。HBase部署在这5台机器上，每台机器的RegionServer设置最大堆内存为4GB。其中HBase版本为1.0.1.1，Titan版本为0.5.4。

### 7.1 邻域点集相关查询

邻域点集查询是指查询给定点的邻接点集。许多图查询基于邻域点集查询实现，如k-hop点集查询、局部聚集系数查询、广度优先搜索等。下面叙述两个基于邻域点集查询的图查询实验。

#### 7.1.1 k-hop点集查询

k-hop点集查询即查询给定点在k跳能到达的点集。实验在测试图中随机抽取100个点作为起点，查询它们的多跳邻域点集，测试它们的平均查询时间。同时又对小邻域（邻域点数小于5）的点集和大邻域（邻域点数大于20）的点集进行了同样的实验。表7.1、表7.2、表7.3 是Titan与HybriG的性能表现。

架构 hops	1	2	3	4
Titan	13.65	72.97	446.34	2755.86
HybriG	4.25	18.73	34.69	99.78

表 7.1 随机选100个点的k-hop查询平均耗时（毫秒）

可以看到，随着跳数的增加，两系统的查询时间都呈指数增长。HybriG不管在1跳的初始值还是在增长的倍数上都远优于Titan。正如2.3节分析的，对某个点的邻域点集查询需要遍历该点的邻接表。当图中的重边数目巨大时，Titan受累于其庞大的邻接

表，对邻域点集的查询速度大大降低。而HybriG由于极大简化了存储在Titan中的边数目，使得邻接表的数据量大大缩小，从而降低了邻接表的遍历耗时。

### 7.1.2 局部聚集系数计算

局部聚集系数 (Local Cluster Coefficient) 是图中每个点的一种度量，用于表示其邻域子图的紧密程度，即邻居节点之间有多大比例有边相连。简要的计算公式如下：

$$LCC(v) = \frac{|\{(u, w) : u, w \in N_v, e(u, w) \in E\}|}{|\{(u, w) : u, w \in N_v\}|}$$

其中 $N_v$ 表示点 $v$ 的邻域点集， $E$ 表示图中的边集。上式表示邻域的点对集合中，有边相连的点对比例。由于 $e(u, w) \in E$ 等价于 $w \in N_u$ ，故上式分子部分可转化为

$$|\{(u, w) : u, w \in N_v, w \in N_u\}| = \sum_{u \in N_v} |N_u \cap N_v|$$

从上式可知，局部聚集系数的计算需要对邻域点集中的点再做一次邻域点集查询，因此其复杂度相当于2跳邻域点集查询。实验对比的是HybriG与直接将图存储在Titan中的传统方案。测试点集的抽取方式同7.1节，即小邻域（邻域点数小于5）的点集、随机采样点集和大邻域（邻域点数大于20）的点集，每个点集拥有100个点。图10是实验结果。与k-hop查询一样，HybriG在局部聚集系数计算的性能上要远优于Titan，而且在大邻域点集上的优势更加明显。

## 7.2 边集相关查询

下面介绍边集相关查询的两个实验结果。7.2.1节是两点间边集查询，在刑侦场景中，图中的一类顶点代表人，人之间的边代表两人的关系，比如共同出行记录、共同住宿酒店记录、通话记录等。当研判专家锁定两个嫌疑人后，需要查询两人间的关

架构 hops	1	2	3	4
Titan	8.67	31.08	179.47	1140.0
HybriG	2.08	7.79	19.88	52.45

表 7.2 随机选100个小邻域点的k-hop查询平均耗时（毫秒）

架构 hops	1	2	3	4
Titan	17.79	110.34	721.25	4781.79
HybriG	3.52	14.49	39.98	163.15

表 7.3 随机选100个大邻域点的k-hop查询平均耗时（毫秒）

系数据，即为两点间边集查询。7.2.2节是邻接边集查询，在刑侦场景中，锁定嫌疑人后要查看其某种类型的关系数据，即为邻接边集查询。这两种查询的图查询语句是不同的，前者限定了边的两个邻接点，后者只给定了边一个端点，侧重于限定label。使用Blueprints图查询接口，两点间边集查询的示例语句如下：`v1.query().adjacent(v2).edges()` 邻接边集查询的示例语句如下：`v.getEdges(Direction.BOTH, eLabel)` 其中v1、v2、v是顶点，eLabel是边的一种label，Direction.BOTH是常量，代表边的方向可以是入边或出边。

### 7.2.1 两点间边集查询

实验测试的是给定两个点和一个label，查询两点间该label的所有边。在图中随机抽取100个有边相连的点对，实验测试查询这些边的耗时。图 11统计了两种系统的时间对比，两种方案的查询耗时非常接近。HybriG的存储层基于Titan和HBase实现，对边集的查询先要在Titan里查询边id，再在HBase的边表中查询具体的边，因此时间开销分两部分（详见5.1节）。图中展示了HybriG查询时间的两部分组成。两种方案的耗时相近主要有两方面的原因。一方面，HybriG检索边集数据需要在Titan和HBase这两个数据库中进行查询，且这两个查询不能并行。传统方案只需要在Titan中查询，HybriG多加了一轮查询的时间开销。然而，这部分的时间开销并不是很大。Titan的数据表本身也存储在HBase中，HybriG的查询实际上只是HBase中的跨表查询。跨表查询能共用HBase的一些缓存信息，如HMaster、RegionServer的位置信息、HBase中Root表和Meta表的信息等。不过尽管只是跨表查询，在这方面HybriG还是引入了时间开销。但另一方面，HybriG对边集的查询是转化为HBase边表中的行级别检索（详见5.1节）。而将图存储在Titan中的方案，对边集的查询实际是转化为HBase数据表中选定行后的列级别检索。当图中含有大量重边时，前者是在一张高窄（tall-narrow）表中查找连续的几行，后者是在一张扁宽（flat-wide）表中选定一行后查找连续的几列，前者的性能会略优一些。因此在这方面HybriG略优。

### 7.2.2 邻接边集查询

实验测试的是给定一个点，查询其某种label的所有边。点集的抽取方式同7.1节，即小邻域（邻域点数小于5）的点集、随机采样点集和大邻域（邻域点数大于20）的点集，每种点集抽取1000个点。图12展示的是两种系统的查询时间，以及查询耗时中具体的时间组成，其中HybriG的查询时间分为Titan中的查询耗时和HBase中的查询耗时两部分。当边集规模增大时，HybriG在Titan中的查询耗时并没有显著增加，而在HBase边表中的查询耗时增长的幅度也没有传统方案中的Titan大。

### 7.3 边数据的导入速度

关于边的数据导入速度，我们对不同规模的图进行了测试。对比直接将图存储在Titan以及将图存储在HybriG的两种方案。我们基于MapReduce[9]开发了分布式的导入程序。表7.4是边数据导入时间的统计：可以看到，在所有测试用图中，HybriG对

测试图号	点数	边数	重边的重数	Titan导入时间	HyBriG导入时间
1	50000	23786800	100	30min	14min
2	200000	48197700	100	45min	16min
3	200000	481977000	1000	9h 26min	1h 24min

表 7.4 边集数据的批量导入耗时

边集数据的导入拥有更高的速度。这主要有两方面的原因：一是HybriG对HBase边表的数据导入使用了HBase的Bulk Loading技术。二是Titan对新增的点和边都会分配一个全局唯一的id（这也是Titan没有实现HBase Bulk Loading的原因），HybriG大大减小了自身Titan中的边数目，从而节省了id分配的时间开销。



## 结论

本文提出了一种基于Titan和HBase的复合存储架构——HyBriG，可以高效处理含有大量重边的属性图。相比于传统图数据库Titan的实现，HyBriG在处理大量重边时拥有更优异的查询性能，在数据导入方面也有不错的表现。在获得高性能的同时，所带来的牺牲就是简化了对事务性的支持。然而，含有大量重边的许多应用场景并不需要很强的事务性支持，HybriG架构提供了最终一致性，足以满足这些场景的应用需求。在实践中，明略数据的关联数据分析平台——SCOPA基于HybriG架构开发了核心存储部件。证明了HybriG架构在处理大量重边时的优异性能。



## 参考文献

- [1] Author. “Title” [J]. *Journal*, 2014-04-01.
- [2] 作者。“标题” [J]。期刊，2014-04-01。



## 附录 A 附件

*pkuthss* 文档模版最常见问题:

`\cite`、`\parencite` 和 `\supercite` 三个命令分别产生未格式化的、带方括号的和上标且带方括号的引用标记: 1, [2]、<sup>[1,2]</sup>。

若要避免章末空白页, 请在调用 *pkuthss* 文档类时加入 `openany` 选项。

如果编译时不出参考文献, 请参考 `texdoc pkuthss` “问题及其解决”一章“其它可能存在的问题”一节中关于 `biber` 的说明。



## 致谢

*pkuthss* 文档模版最常见问题:

`\cite`、`\parencite` 和 `\supercite` 三个命令分别产生未格式化的、带方括号的和上标且带方括号的引用标记: 1, [2]、<sup>[1,2]</sup>。

若要避免章末空白页, 请在调用 `pkuthss` 文档类时加入 `openany` 选项。

如果编译时不出参考文献, 请参考 `texdoc pkuthss` “问题及其解决”一章“其它可能存在的问题”一节中关于 `biber` 的说明。





## 北京大学学位论文原创性声明和使用授权说明

### 原创性声明

本人郑重声明：所呈交的学位论文，是本人在导师的指导下，独立进行研究工作所取得的成果。除文中已经注明引用的内容外，本论文不含任何其他个人或集体已经发表或撰写过的作品或成果。对本文的研究做出重要贡献的个人和集体，均已在文中以明确方式标明。本声明的法律结果由本人承担。

论文作者签名：                    日期：      年      月      日

### 学位论文使用授权说明

(必须装订在提交学校图书馆的印刷本)

本人完全了解北京大学关于收集、保存、使用学位论文的规定，即：

- 按照学校要求提交学位论文的印刷本和电子版本；
- 学校有权保存学位论文的印刷本和电子版，并提供目录检索与阅览服务，在校园网上提供服务；
- 学校可以采用影印、缩印、数字化或其它复制手段保存论文；
- 因某种特殊原因需要延迟发布学位论文电子版，授权学校在□一年/□两年/□三年以后在校园网上全文发布。

(保密论文在解密后遵守此规定)

论文作者签名：                    导师签名：                    日期：      年      月      日