

# Fast and Accurate Estimation of Shortest Paths in Large Graphs

@Authors: Andrey Gubichev, Srikantha Bedathur, Stephan Seufert, Gerhard Weikum

@Published in: ACM International Conference on Information and Knowledge Management, 2010 :499-508

@Presented by: Yina Lv , Time: Oct 26, 2018

@Action: October 22, 2018 11:08 AM

## 1.工作概述

1. 提出a scalable sketch-based index structure, 引入 `path-sketches`, 可以有效地用于小直径的large graphs (例如,大多数在线社交网络)。 `path-sketches` 保留每个节点与选定的一组 `landmark nodes` 之间的完整路径信息, 作为图形预处理步骤的一部分
2. 开发了一套轻量级但高效的技术, 这些技术使用 `path-sketches` 显著提高最短路径估计的质量
3. 随着对最短路径距离的估计, 我们展示了如何在没有计算开销的情况下按照其长度的递增顺序生成 `路径队列`, 这是许多应用程序在社交和生物网络上的重要需求
4. 最后, 我们在一个功能齐全的大规模图形处理引擎RDF-3X中实现所有提出的方法, 并对许多真实世界的大规模图进行评估,实现查询响应时间比传统路径计算快几个数量级, 同时将平均估计误差保持在0%和1%之间

## 2.背景

在现代数字世界中, 图的广泛应用, 比如在线社交网络中(LinkedIn,Facebook,MySpace), 大规模知识中的综合实体关系, 边缘存储库, 生物交互模型, 交通网络, 万维网文档之间巨大的超链接图, XML数据, 等等

例如:

1. 社交网络中, 一个人可以从他朋友那里得知他最喜欢的雇主的招聘经理。我们可以想到六度空间理论, 你和任何一个陌生人之间所间隔的人不会超过六个, 也就是说, 最多通过五个中间人你就能够认识任何一个陌生人
2. 生物网络模拟生物体内的复杂化学过程。生物学家可能对识别两种目标代谢物之间的生物转化途径感兴趣, 以帮助设计实验

在许多情况下, 图中节点数和边的数目达到了数百万个, 出于可伸缩性的原因, 为了让内存和处理器资源消耗最小化, 每一个最短路径的查询实例必须要尽可能快。

- 是什么让最短路径计算在大图上特别难？
  - Dijkstra算法是计算图中两个节点间最短路径的经典方法，时间复杂度为 $O(m+n\log(n))$ ，其中， $n$ 是节点数， $m$ 是边数
  - 对于节点数和边数非常多的情况计算时间很大。由于Dijkstra的算法必须构建和维护到图中所有节点的最短路径

## 3.算法描述

### 3.1 准备工作

$G = (V, E)$ 是一个有向图，节点集 $V$ ，边集 $E$ 。一条路径 $p$ 的长度为 $l$ ，且 $l \in \mathbb{N}$ ，那么节点数为 $l+1$ 。如果一条路径存在，那么可表示为如下形式：

$$p = (v_1, v_2, \dots, v_{l+1}), v_i \in V, 1 \leq i \leq l+1; (v_i, v_{i+1}) \in E, 1 \leq i \leq l$$

对于一个节点 $v \in V$ ，用 $S(v)$ 表示 $v$ 的后继节点集合，即具有 $(v, w) \in E$ 的顶点集 $w \in V$ 。其中 $(v, w)$ 表示节点 $v$ 指向节点 $w$ 。因此，我们也可得到 $v_{i+1} \in S(v_i)$

$|p|=l$  表示路径 $p$ 的长度。对于节点 $u, v \in V$ ， $\mathcal{P}(u, v)$  表示从 $u$ 到 $v$ 的所有路径集合。 $\text{dist}(u, v)$  表示从 $u$ 到 $v$ 的距离,是最短路径中的边数，如果不能从 $u$ 到达 $v$ ，则为无穷大。

$$\text{dist}(u, v) := \begin{cases} \arg \min_{p \in \mathcal{P}(u, v)} |p| & \text{if } \mathcal{P}(u, v) \neq \emptyset, \\ \infty & \text{else.} \end{cases}$$

### 路径估计

设两个顶点 $u, v \in V$ ，设 $p$ 表示从 $u$ 到 $v$ 的最短路径(注意可以有多条)，即从 $u$ 开始，以 $v$ 结尾的路径，长度为 $|p| = \text{dist}(u, v)$ 。此外， $q$ 作为从 $u$ 到 $v$ 的任意路径，通过将 $q$ 看作最短路径 $p$ 的逼近，我们可以定义这条路径的逼近误差为

$$\text{error}(q) := \frac{|q| - |p|}{|p|} = \frac{|q| - \text{dist}(u, v)}{\text{dist}(u, v)} \in [0, \infty].$$

### 路径级联

设 $p = (u_1, u_2, \dots, u_{L_1}, u_{L_1+1})$ 和 $q = (v_1, v_2, \dots, v_{L_2+1})$ 路径长分别是 $L_1$ 和 $L_2$ 。假设 $u_{L_1+1} = v_1$ ，即路径 $p$ 中的最后一个节点等于路径 $q$ 中的第一个节点。然后，通过将路径 $p$ 和 $q$ 连接起来，可以创建长度为 $L_1 + L_2$ 的新路径(以 $p \circ q$ 表示)：

$$p \circ q = (u_1, u_2, \dots, u_{l_1}, u_{l_1+1}) \circ (v_1, v_2, \dots, v_{l_2+1})$$

$$:= (u_1, \dots, u_{l_1}, v_1, \dots, v_{l_2+1}).$$

## 3.2 Sketch Algorithm

Sketch Algorithm 使用 landmark-based 的方法近似计算一般图中两个给定节点之间的最短路径长度

为了实时地回答一对节点的距离查询，该算法采用了两阶段的方法：

1. 预处理生成 Sketch (从所有顶点到landmark节点的距离)
2. 使用这些预先计算的数据，在查询时可快速得到节点近似距离。通过查询的源节点s和目标节点t与landmark节点l之间的距离  $\text{dist}(s, l)$  、  $\text{dist}(l, d)$  。那么  $s \rightarrow t$  近似距离为：

$$\tilde{d}(s, d) \leq \text{dist}(s, l) + \text{dist}(l, d)$$

在提出 Sketch Algorithm 的文章中，作者建议为每个节点v存储节点到某些landmark节点的距离  $\text{dist}(v, \cdot)$  和  $\text{dist}(\cdot, v)$  作为预处理的结果。这组 node-landmark 距离称为 sketch of a node

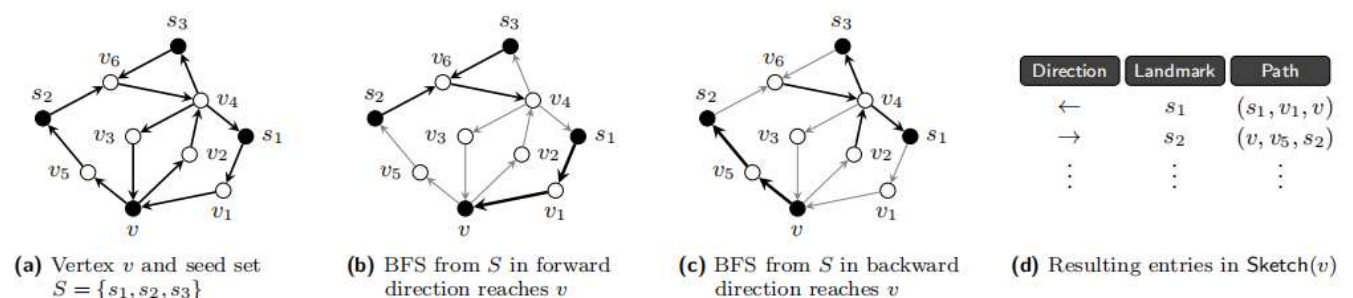
我们的做法（在原先算法基础上做了修改）：

- 在预处理阶段不仅计算距离还存储实际路径！
- 查询节点之间距离的估计，返回的不仅仅是查询节点之间的一条近似路径，即这些路径的队列（按路径长度按照顺序排序）。是因为有时候必须满足对某些节点/边的约束。

### 3.2.1 预处理

预计算步骤涉及对一些节点集进行采样，为图中的每个节点计算进出该集合中成员的最短路径，并将获得的路径集存储在外部的存储器上，这些路径将在稍后的近似步骤中使用。

Figure 1: Precomputation Example



1. 种子集抽样

$r := \lfloor \log(n) \rfloor, n=|V|$ 。我们抽取 $r+1$ 个节点集样本（叫做 seed sets ），并且大小依此为

$1, 2, 2^2, \dots, 2^r$ 。标号  $S_0, S_1, \dots, S_r$

## 2. 最短路径计算

对于每一个样本种子集和每一个节点  $v \in V$ ，我们计算一条从  $S_i$  到  $v$  的最短路径  $P_{S_i \rightarrow v}$ ，也就是计算出了任何一个节点到种子集的路径（可能有多条）。接下来，我们从  $S_i$  出发 BFS，构造一颗完全最短路径树。因此，对于每个节点  $v$ ，我们获得了最接近的种子节点，标记为  $l_1$ 。同理，我们可以计算  $P_{v \rightarrow S_i}$  得到一个  $l_2$ 。那么  $l_1, l_2$  就叫做节点  $v$  的关于  $S_i$  的 landmarks。

- 对于一个节点而言，在一个种子集最多可得到两个 landmarks。
- 每个节点的 landmarks 不一定相同
- $2rk$  个 landmarks 和路径组成了  $\text{Sketch}(v)$

### 3.2.2 近似最短路径

输入  $(s, d)$ ，输出  $p_{s \rightarrow d}$  良好的近似最短路径。

从磁盘导入  $\text{Sketch}(s)$  和  $\text{Sketch}(d)$ ，并且计算两者的公共 landmarks，如果不是公共 landmarks 去计算的话，要么路径不存在，要么误差相对大。

#### Algorithm 1: $\text{SKETCH}(s, d)$

Input:  $s, d \in V$

Result:  $Q$  – priority queue of paths from  $s$  to  $d$ , ordered by path length

```
1 begin
2   Load sketches  $\text{Sketch}(s), \text{Sketch}(d)$  from disk
3    $L \leftarrow$  common landmarks of  $\text{Sketch}(s)$  and  $\text{Sketch}(d)$ 
4   foreach  $l \in L$  do
5      $p \leftarrow$  path from  $s$  to  $d$  through  $l$ 
6     Add  $p$  to queue  $Q$ 
7   return  $Q$ 
```

在上述算法中，第五句  $p_{s \rightarrow d} := p_{s \rightarrow l} \circ p_{l \rightarrow d}$  存储在路径队列中(按路径长度升序排序)，返回回路径的优先级队列。

## 3.3 提高精度

- Cycle Elimination
- Shortcutting

### 3.3.1 Cycle Elimination

在Algorithm 1 提出的算法得到的路径可能是次优的，也就是比真实的最短路径更长,因为它们包含 环。

$$p_{s \rightarrow l} \circ p_{l \rightarrow d} = (s, v1, v2, l, v3, v1, d)$$

---

**Algorithm 2: SKETCHCE( $s, d$ )**

---

**Input:**  $s, d \in V$

**Result:**  $Q$  – priority queue of paths from  $s$  to  $d$ , ordered by path length

```
1 begin
2    $Q \leftarrow \text{SKETCH}(s, d)$ 
3   foreach  $p = (p_1, p_2, \dots, p_l) \in Q$  do
4     for  $i = 1$  to  $l - 1$  do
5       for  $j = 0$  to  $l - i - 1$  do
6         if  $p_i = p_{l-j}$  then
7            $Q \leftarrow Q \cup \{(p_1, \dots, p_i, p_{l-j+1}, \dots, p_l)\}$ 
8           break ▷ continue in line 3
9   return  $Q$ 
```

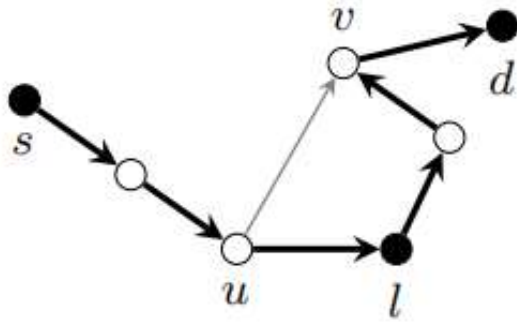
---

### 3.3.2 Shortcutting

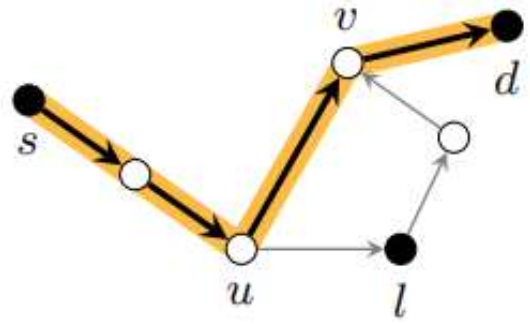
假设算法返回的队列 $Q$ 包含路径 $p_{s \rightarrow l \rightarrow d}$ ，从 $s$ 到 $d$ 经由landmark  $l$ 。路径中的两个节点 $u, v$ 实际上可能比 $p_{s \rightarrow l \rightarrow d}$ 的相应子路径中包含的节点更紧密。考虑下面描述的示例：当节点 $u$ 和 $v$ 通过长度为3的 $p_{s \rightarrow l \rightarrow d}$ 的子路径连接时，原始图形包含边缘  $(u, v)$ 。然后我们可以通过单边  $(u, v)$  轻松替换此子路径。

请注意，在此方式优化的所有情况下，landmark  $l$  一定位于从 $u$ 到 $v$ 的子路径中。





(a) Path from  $s$  to  $d$ . The original graph contains the arc  $(u, v)$



(b) Path from  $s$  to  $d$  after shortcutting

### Algorithm 3: SKETCHCESC( $s, d$ )

Input:  $s, d \in V$

Result:  $Q$  – priority queue of paths from  $s$  to  $d$ , ordered by path length

```

1 begin
2    $Q \leftarrow \text{SKETCHCE}(s, d)$ 
3   foreach  $p = (p_1, p_2, \dots, p_{i-1}, p_i, p_{i+1}, \dots, p_l) \in Q$  do
4     for  $j = 1$  to  $i - 1$  do
5        $\mathcal{S} \leftarrow$  set of successors of  $p_j$ 
6       for  $k = 0$  to  $l - i + 1$  do
7         if  $p_{l-k} \in \mathcal{S}$  then
8            $Q \leftarrow Q \cup \{(p_1, \dots, p_j, p_{l-k}, \dots, p_l)\}$ 
9           break
10  return  $Q$ 

```

## 3.4 Tree Algorithm

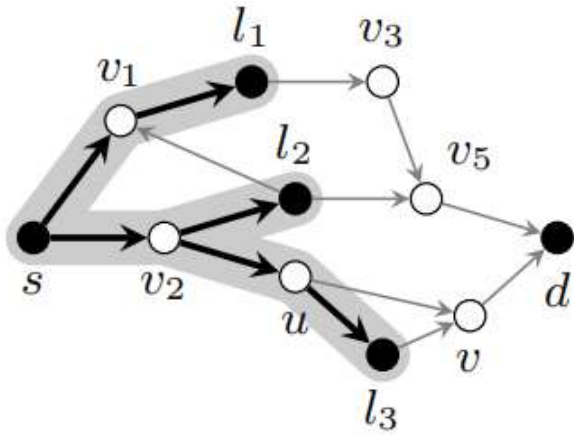
一种新的最短路径近似算法，它也利用了precomputed sketches。

节点 $v$ 的Sketch ( $v$ ) 包含两组路径：

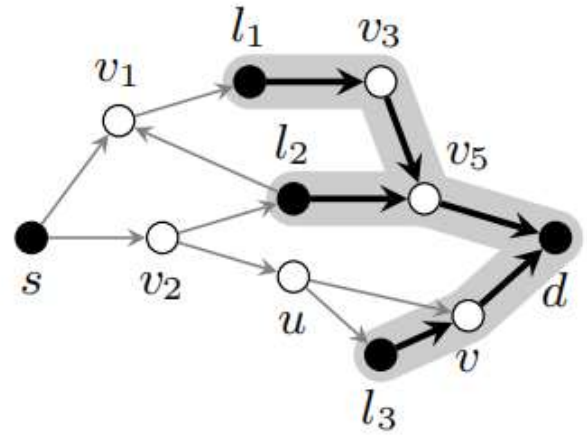
1. 将 $v$ 连接到landmarks的路径集 (称为 forward-directed paths , 前向路径)
2. 将landmarks连接到 $v$ 的路径集 (称为 backward-directed paths , 后向路径)

在无向设置中，两个集合将作为树的叶子和根节点 (landmarks作为叶子，节点 $v$ 作为根)。每棵

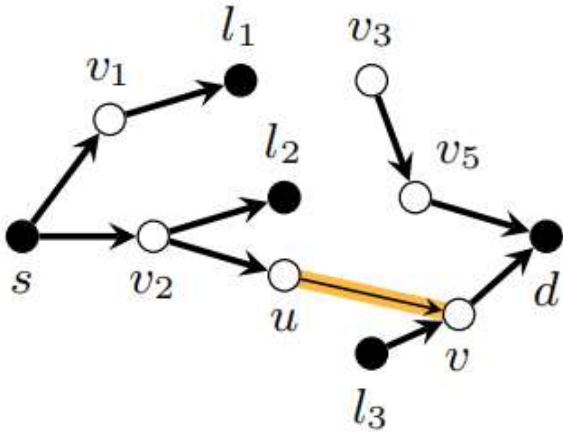
树的每个内部节点都对应于sketch中一条路径中的一个节点。在有向设置中，只有sketch的 forward-directed部分（从 $v$ 到landmarks）形成一棵树，而backward-directed path产生一个带有“reversed edges”的树,例如，下图a和b：



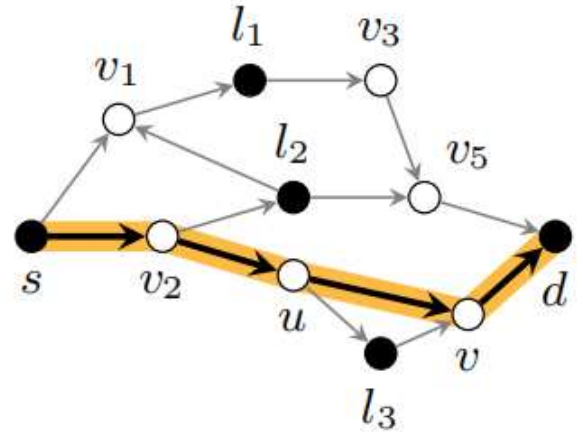
(a) Tree  $T_s$  of paths from  $s$  to landmarks



(b) “Reversed Tree”  $T_d$  of paths from landmarks to  $d$



(c) Arc discovered in line 13 of algorithm TREESKETCH



(d) Obtained path from  $s$  to  $d$

我们提出的新算法，称为 `TreeSketch`，将两个查询节点 $s, d$ 作为输入，从磁盘加载sketch

`Sketch(s)`，`Sketch(d)` 并构造以 $s$ 为根的树 $T_s$ ，其中在`Sketch(s)`包含存储的所有前向路径。

同样地，以 $d$ 为根的“reversed tree”  $T_d$ ，包含从landmarks到 $d$ 的所有后向路径。

然后，该算法同时在树上开始两个BFS：在 $T_s$ 中以 $s$ 为根节点的BFS( $T_s, s$ )和 $T_d$ 中以 $d$ 为根节点的RBFS( $T_d, d$ )。 $V_{BFS}$ 和 $V_{RBFS}$ 表示在各自BFS运行期间访问过的节点集。 $l_{shortest}$ 表示 $Q$ 中最短路径的长度。

对于在BFS( $T_s, s$ )过程中遇到的每个顶点 $u \in V_{BFS}$ ，该算法在原图中加载其继承者的列表 $S(u)$ 。然后，检查在RBFS期间发现的任何顶点( $T_d, d$ )存于 $S(u)$ 。如果存在这样一个顶点 $v \in S(u) \cap V_{RBFS}$ (见图c)，我们已经找到了从 $s$ 到 $d$ 的路径 $p$ ，如下：

$$p = p_{s \rightarrow u} \circ (u, v) \circ p_{v \rightarrow d}$$

其中 $p_{s \rightarrow u}$ 和 $p_{v \rightarrow d}$ 分别表示 $T_s$ 中从 $s$ 到 $u$ 的路径和 $T_d$ 中从 $v$ 到 $d$ 的路径(图d)。

代码中10-16行,  $v \rightarrow d$  路径确定的情况下, 找  $T_s$  中已访问过的节点中是否有某个节点  $x$  的后继节点集包括了  $v$ , 若包括, 那么从  $x$  到  $v$  存在一条路径, 即路径  $s \rightarrow x \rightarrow v \rightarrow d$  存在。放入路径优先队列  $Q$  中, 并与  $l_{\text{shortest}}$  进行比较

代码17-23行,  $s \rightarrow u$  确定的情况下, 在 RBFS 已经访问过的节点集  $V_{\text{RBFS}}$  中找。若存在某个节点  $x$  属于节点  $u$  的后继节点, 那么存在路径  $s \rightarrow u \rightarrow x \rightarrow d$ 。放入路径优先队列  $Q$ , 并与  $l_{\text{shortest}}$  进行比较

代码24, 如果没有找到比  $Q$  中当前最短路径更短的路径, 则算法终止。当两个 BFS 运行的深度总和超过  $l_{\text{shortest}}$  时就是这种情况。

#### Algorithm 4: TREESKETCH( $s, d$ )

**Input:**  $s, d \in V$

**Result:**  $Q$  – priority queue of paths from  $s$  to  $d$ , ordered by path length

```

1  begin
2  | Load Sketch( $s$ ), Sketch( $d$ ) from disk
3  |  $T_s \leftarrow$  tree of paths from  $s$  ▷ taken from Sketch( $s$ )
4  |  $T_d \leftarrow$  tree of paths to  $d$  ▷ taken from Sketch( $d$ )
5  |  $Q \leftarrow \emptyset$ 
6  |  $l_{\text{shortest}} \leftarrow \infty$ 
7  |  $V_{\text{BFS}} \leftarrow \emptyset$ 
8  |  $V_{\text{RBFS}} \leftarrow \emptyset$ 
9  | foreach  $u \in \text{BFS}(T_s, s)$  and  $v \in \text{RBFS}(T_d, d)$  do
10 | |  $V_{\text{BFS}} \leftarrow V_{\text{BFS}} \cup \{u\}$ 
11 | |  $p_{v \rightarrow d} \leftarrow$  path from  $v$  to  $d$  in  $T_d$ 
12 | | foreach  $x \in V_{\text{BFS}}$  do ▷ iteration in order of visits
13 | | | if  $v \in \mathcal{S}(x)$  then ▷  $\mathcal{S}(x)$  is set of successors of  $x$  in  $G$ 
14 | | | |  $p \leftarrow p_{s \rightarrow x} \circ (x, v) \circ p_{v \rightarrow d}$ 
15 | | | |  $Q \leftarrow Q \cup \{p\}$ 
16 | | | |  $l_{\text{shortest}} \leftarrow \min\{l_{\text{shortest}}, |p|\}$ 
17 | |  $V_{\text{RBFS}} \leftarrow V_{\text{RBFS}} \cup \{v\}$ 
18 | |  $p_{s \rightarrow u} \leftarrow$  path from  $s$  to  $u$  in  $T_s$ 
19 | | foreach  $x \in V_{\text{RBFS}}$  do ▷ iteration in order of visits
20 | | | if  $x \in \mathcal{S}(u)$  then ▷  $\mathcal{S}(u)$  is set of successors of  $u \in G$ 
21 | | | |  $p \leftarrow p_{s \rightarrow u} \circ (u, x) \circ p_{x \rightarrow d}$ 
22 | | | |  $Q \leftarrow Q \cup \{p\}$ 
23 | | | |  $l_{\text{shortest}} \leftarrow \min\{l_{\text{shortest}}, |p|\}$ 
24 | | if  $\text{dist}(s, u) + \text{dist}(v, d) \geq l_{\text{shortest}}$  then break
25 | return  $Q$ 

```

## 4.实现



我们实现了所有方法“Dijkstra”在线最短路径算法以及sketch-based的技术-在RDF-3X中，这是最近提出的用于存储和查询大型RDF图形存储库的高性能数据库系统。在介绍实现细节之前，我们简要介绍一下RDF-3X的背景，重点介绍它是否能够存储高度压缩的大图，同时易于查询。请注意，我们在RDF-3X中省略了查询处理和优化策略的详细信息，因为它们与我们的设置无关。

## 4.1 RDF-3X: RDF Graph Processor

RDF-3X是一个功能齐全的高性能存储引擎和查询处理器，专门用于存储。它基本上将整个图保持为一个巨大的三元组表，这与最近流行的属性表方法形成了对比。

边的三元组形式  $\langle u, e, v \rangle$

在RDF中， $\langle \text{Subject}(S), \text{Predicate}(P), \text{Object}(O) \rangle$  ,六个关于S、P、O的组合存于6棵B+树索引中。

每个(SPO, SOP, OSP, OPS, PSO, POS)都有一个索引。通过使用三元组的delta-coding可以显著地压缩每个索引，这是将倒排列表索引中使用的类似压缩的概括推广到id三元组。此外，RDF-3X还构建了6个索引，以便有效地支持分析查询。相反，整个数据库大小包含12个不同的集群B+树索引，上面提到的所有压缩都比三元组形式的原始数据要小。

在我们的实现中，用RDF-3X边存储图，每个边表示为一个三元组  $\langle s, e, t \rangle$  。我们不限RDF-3X自动构建所有12个索引。我们只利用SPO和OPS顺序索引。

**Dijkstra的算法：**在RDF-3X上实现Dijkstra算法基本上涉及在SPO索引上打开扫描，以针对在算法的执行期间访问的每个节点确定所有后继节点。请注意，我们在计算草图时需要反向Dijkstra算法，主要是为了“简单地”打开OPS索引上的扫描，并让算法运行。Dijkstra算法期间所需的优先级队列使用GNU-C++ STL中可用的实现在内存中维护。

## 4.2 Sketch 实现

为简单起见，我们将草图也作为RDF三元组存储在RDF-3X下的单独数据库中。由于sketch（距离和路径sketch）与有向路径相关联，因此我们将其格式设置如下：

$$\langle v_i \rangle \langle [to|from] \rangle \langle l_{ij} : d_{ij} \rangle$$

其中， $v_i$ 是源节点的id，往返是指示路径方向的字符串文字， $l_{ij}$ 是来自种子集 $S_j$ 的节点 $v_i$ 的界标， $d_{ij}$ 是到地标的相应最短距离。同样，路径草图也存储为以下形式的三元组：

$$\langle v_i \rangle \langle [to|from] \rangle \langle l_{ij} : p_{ij} \rangle$$

其中 $p_{ij}$ 指的是节点 $v_i$ 和地标节点 $l_{ij}$ 之间路径中节点ID的序列。

应该再次注意的是，RDF-3X在这个三重草图数据库上建立了12个索引，但我们只需要其中一个索引用于我们的最短路径估计算法，即在置换SPO上的B+树。因此，我们在后面的实验部分提供的磁盘空间消耗可以进一步大大减少，尽管草图的相对大小保持不变。在我们所有的实验中默认设置 $k=2$ 预计算回合。

# 5.实验评估

步骤:

- 概述所使用的数据集
- 描述在后续评估中使用的测试实例的生成
- 评估不同方法的逼近质量
- 执行查询运行时间测量
- 简要讨论了路径多样性问题
- 总结，量化了草图预计算的空间和时间要求

## 5.1 使用的数据集

Table 1: Used Datasets

Dataset	$ V $	$ E $	$\bar{\delta}$	$ S / V $	$d_{0.9}$
Slashdot	77,360	905,468	23.4	90.9 %	5.59
Google	875,713	5,105,039	11.7	49.6 %	9.02
YouTube	1,138,499	4,945,382	8.7	44.7 %	7.14
Flickr	1,715,255	22,613,981	26.4	69.5 %	7.32
WikiTalk	2,394,385	5,021,410	4.2	4.6 %	4.98
Twitter	2,408,534	48,776,888	40.5	57.5 %	5.52
Orkut	3,072,441	223,534,301	145.5	97.5 %	5.70

Datasets with no. of vertices  $|V|$ , no. of edges  $|E|$ , average degree  $\bar{\delta}$ , percentage of nodes in the largest strongly connected component  $S$  and effective diameter  $d_{0.9}$  [15].

## 5.2 测试实例生成

为了评估算法的逼近质量和运行次数，我们使用了一组表的测试三元组 $(x,y,dist(x,y))$ ,其中，包括了一对节点 $x, y$ 和 $x, y$ 之间的距离 $dist(x,y)$

我们使用Dijkstra的算法，通过均匀采样100个顶点并计算每个顶点的最短路径树(forward and backward direction)来生成这些三元组。作为输出，我们为每个采样顶点 $v$ 获得一个将顶点连接到每个其他可到达节点的树和一个“反转”树，将存在采样顶点的路径的每个节点连接到 $v$ 。

结果，我们获得了 $(x,y,dist(x,y))$ 中所示结构的一组三元组。然后，我们将这些三元组分组为与距离 $dist(x, y)$ 对应的类别。从每个这样的类别，我们最多采样50个三元组（测试）作为我们的测试集。每个网络的实际测试次数各不相同，因为组的数量以及所包含的三元组的数量可能不同。

## 5.3 结果

### 5.3.1 逼近质量

为了评估不同算法生成的路径的逼近质量，我们对每一个三元组(s,d,dist(s,d))用SKETCH, SKETCHCE, SKETCHCESC, TREESKETCH四种算法查询。对于每一个最短路径查询(s,d),我们将得到的路径长度 $l_{\text{shortest}}$ 与三元组中真实的最短距离进行比较，按照如下公式计算误差：

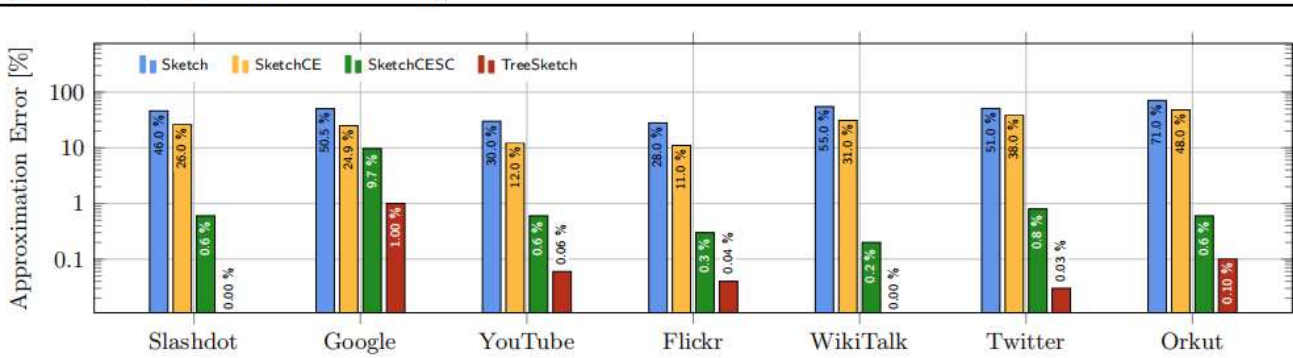
$$\text{error}(p_{s \rightarrow d}) = \frac{l_{\text{shortest}} - \text{dist}(s, d)}{\text{dist}(s, d)} \in [0, \infty].$$

对于每个算法，我们记录所有测试三元组的平均近似误差。所获得的误差值在表2中提供并绘制在图4中，使用对数刻度显示误差值

Table 2: Approximation Quality and Running Times

Dataset	Average Approximation Error					Average Running Times				
	Tests	Sketch	SketchCE	SketchCESC	TreeSketch	Sketch	SketchCE	SketchCESC	TreeSketch	Dijkstra (Queue)
Slashdot	910	46.0 %	26.0 %	0.6 %	0.00 %	140 ms	140 ms	198 ms	193 ms	4 s (46K)
Google	3,526	50.5 %	24.9 %	9.7 %	1.00 %	932 ms	932 ms	1,270 ms	1,339 ms	35 s (157K)
YouTube	758	30.0 %	12.0 %	0.6 %	0.06 %	872 ms	872 ms	1,282 ms	1,318 ms	48 s (380K)
Flickr	934	28.0 %	11.0 %	0.3 %	0.04 %	1,217 ms	1,217 ms	2,177 ms	1,951 ms	73 s (696K)
WikiTalk	780	55.0 %	31.0 %	0.2 %	0.00 %	703 ms	703 ms	1,400 ms	1,680 ms	101 s (2M)
Twitter	897	51.0 %	38.0 %	0.8 %	0.03 %	1,932 ms	1,932 ms	3,900 ms	4,000 ms	119 s (1.1M)
Orkut	385	71.0 %	48.0 %	0.6 %	0.10 %	1,090 ms	1,090 ms	2,595 ms	2,751 ms	503 s (2.5M)

Figure 4: Average Approximation Error error(.)



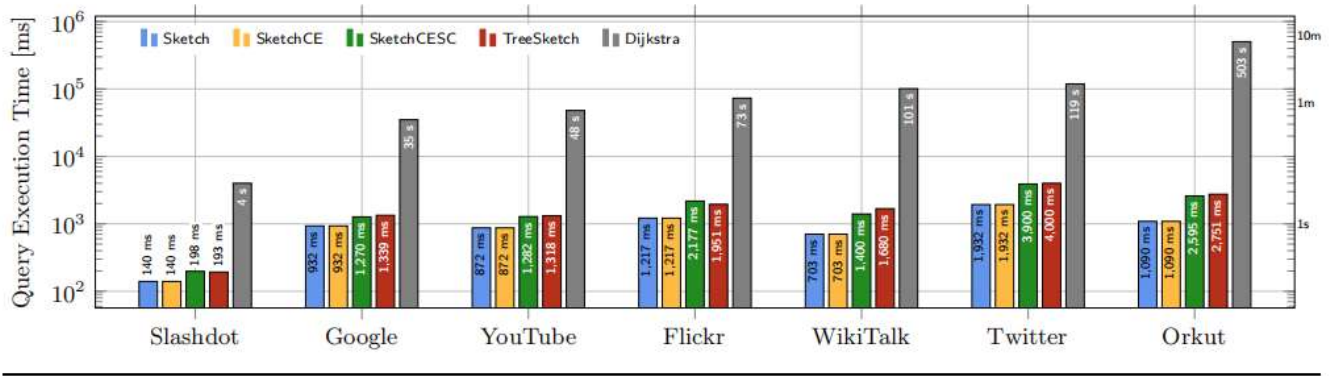
结果显示，我们提出的TREESKETCH算法最优，对于正在考虑的所有数据集，我们的算法能够返回查询节点的最短路径，且在最坏的情况下，平均估计误差为1%，同时为Slashdot和WikiTalk网络提供几乎所有情况下的精确解决方案。

### 5.3.2 查询执行时间

将我们的方法的结果与Dijkstra算法的平均运行时间进行了比较，如上表2

描述了查询执行时间的半对数图：

Figure 5: Query Execution Times



观察到，使用算法SketchCESC，我们能够在最小的数据集（Slashdot）中平均190毫秒内以极高的准确度回答最短路径查询，大型Twitter数据集中4秒。使用TreeSketch，我们可以提供更好的路径估计，提供比Dijkstra算法快一到两个数量级的结果。注意，对于Slashdot和WikitalK数据集，查询执行速度极快，同时实现0%的近似误差。

### 5.3.3 路径多样性

在许多应用程序设置中，不仅要快速生成最短路径的精确逼近，而且还必须尽可能多地返回候选路径。

我们的算法被设计成能够满足这个目标。返回一个有序的路径队列，例如，该队列可用于根据特定用户指定的约束筛选候选人。表4给出了生成的最短路径的平均数目。

Table 4: Number of paths obtained

	Slashdot	Google	YouTube	Flickr	WikiTalk	Twitter	Orkut
SketchCESC	15.0	2.4	19.3	33.3	18.6	45.5	9.5
TreeSketch	31.5	3.1	40.8	55.6	50.7	92.0	29.6

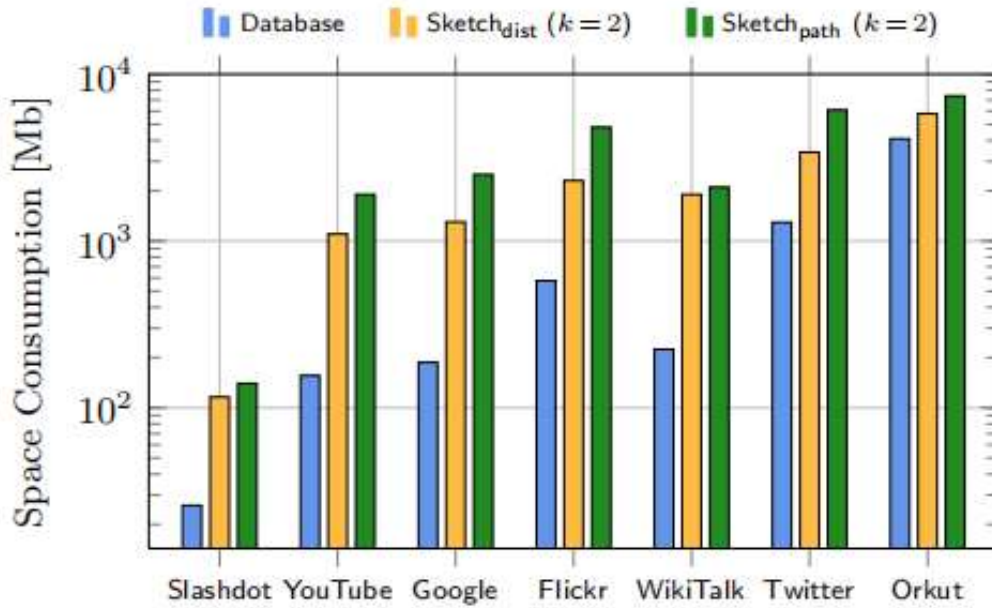
TreeSketch 创建的候选路径数总是大于其他变量生成的路径数。对于twitter数据集，我们平均可以生成92条路径，是我们SketchCESC算法的两倍。

### 5.3.4 评估预处理空间和时间要求

我们评估了预处理步骤(the sketch computation)的空间和时间要求。

空间要求：



**Figure 6:** Sketch Space Consumption for  $k = 2$ **Table 3:** Sketch Space and Time Consumption for  $k = 2$ 

Dataset	Database	Sketch <sub>dist</sub>	Sketch <sub>path</sub>	$t_{\text{precomp}}$ [s]
Slashdot	26 Mb	112 Mb	139 Mb	267.2
Google	0.15 Gb	1.1 Gb	1.9 Gb	4,156.6
YouTube	0.19 Gb	1.3 Gb	2.5 Gb	3,159.0
Flickr	0.57 Gb	2.3 Gb	4.4 Gb	2,223.4
WikiTalk	0.22 Gb	1.9 Gb	2.1 Gb	5,430.0
Twitter	1.30 Gb	3.4 Gb	6.1 Gb	12,806.1
Orkut	5.70 Gb	5.8 Gb	7.4 Gb	36,486.2

### 预处理时间:

我们测量所有数据集的预处理阶段的运行时间。表4提供了 $k = 2$ 预处理迭代所需时间的概述。请注意，所需时间以 $k$ 为线性增加。对于像Slashdot这样的小型数据集，可以在五分钟内获得sketch，而最大的数据集（Orkut）需要大约11个小时的预处理

## “I want to say”

这篇文章就是对原先的sketch算法进行了改进，比如消除环路的存在和shortcutting，以获得更短的路径。并且提出在相同的预先计算的数据上操作的一种全新的Tree sketch算法。此外，还可以生成多条候选路径，作为对一些真实世界数据集的实验评估的基础。

由于文中讨论的是有向图，那么对于上述改进之后，实验效果很可观，但是对于无向图的话，上述方法自然是不可行的。