Scalable Processing of Spatial-Keyword Queries

Walid G. Aref Ahmed Mahmood Purdue University and Google LLC. Alexandria University-Egypt



Outline

- Introduction and Background
- Querying Spatial-Keyword Data
- Scalable Spatial-Keyword Systems
- Research Directions



Introduction and Background

- The Scale of Spatial-Keyword Data
- Spatial-Keyword Applications



Big Spatial-Keyword Data

- Spatial-keyword data is composed of:
 Geo-location
 - Textual content, i.e., set of keywords
 - Other attributes, e.g., timestamp, user identifier
- Large amounts of spatial-keyword data are being generated
 - Location-aware devices, e.g., smart-phones
 - Social networks
 - Micro-blogging applications



Spatial-Keyword Data is Everywhere



UNIVERSITY

The Scale of Spatial-Keyword Data



http://www.internetlivestats.com/ 4/15/21 **PURDUE**

Observations on Spatial-Keyword Data

 The nonuniform and dynamic nature of spatialkeyword data



General Purpose Big Data Systems

 These systems are not equipped with spatial-keyword-specific indexing and query processing algorithms



Introduction and Background

- The Scale of Spatial-Keyword Data
- Spatial-Keyword Applications



Ad Targeting

- Online location-based advertising campaigns
- Location-based advertising, e.g., mobile ecoupons
- Also known as location-aware publish/subscribe systems



Identification of Nearby Attractions



Trip Planning

 A tourist wants to identify groups of points that are close to each other





Traffic and Navigation Systems

• Waze







Micro-Blogs Analysis



https://www.trendsmap.com/



Outline

- Introduction and Background
- Querying Spatial-Keyword Data
- Scalable Spatial-Keyword Systems
- Research Directions



Querying Spatial-Keyword Data

- Spatial-Keyword Queries
- Spatial-Keyword Query Languages



The Spatial-Keyword Filter Query

 Identifies objects located inside the spatial range of the query and matches the textual criteria of the query



The Spatial-Keyword Top-K Query

- Retrieve the top-k spatial-keyword objects ranked based on a function of:
 - Distance between the object and the query location
 - Similarity between the keywords of the object and the query
- Example: Find the top three tuples that have any of the keywords: restaurant, seafood



The Boolean-kNN Query

- A special version of the spatial-keyword top-K query
 - Identify all objects that contain the query keywords
 - Rank these objects based on the distance between objects and the location of the query



The m-Closest Keywords (mCK) Query

 Identify the group of spatial-keyword objects that collectively cover a set of keywords and have a minimum diameter



Zhang et al., Keyword search in spatial databases: Towards searching by document, ICDE 2009 Zhang et al., Locating mapped resources in web 2.0, ICDE 2010 Guo et al., Efficient algorithms for answering the m-closest keywords query, SIGMOD 2015



The Collective Spatial-Keyword Group Query (CoSKQ)

- Find a group of objects that collectively cover the keywords of the query and minimize a function of:
 - The distance between the query and the group
 - The distances among the objects in the group



The Top-k Groups Query

- Identifies multiple groups of spatial-keyword objects that collectively cover the keywords of the query
- Groups are ranked based on a function of
 - Distance of the group to the location of the query
 - The distances among the objects in the group
 - The multiplicity of objects covering the same keywords
- A group may contain multiple objects covering the same keyword



Variations of Spatial-Keyword Queries

- Spatial-Keyword Similarity Join
- Spatial-Keyword Skyline Query
- Spatial-Keyword queries on road-networks
- Continuous queries over spatial-keyword data streams
- Direction-aware spatial-keyword filter query



Querying Spatial-Keyword Data

- Spatial-Keyword Queries
- Spatial-Keyword Query Languages



Querying Spatial-Keyword Data

- Spatial-Keyword Queries
- Spatial-Keyword Query Languages

 MQL
 - Atlas



MQL Micro-Blogs Query Language

 A micro-blogs management system along with a micro-blogs query language



Magdy et.al., Towards a microblogs data management system, MDM 2015



MQL Micro-Blogs Query Language

```
* SELECT [CONTINUOUS] attr_list
FROM stream_name1 [,stream_name2,...]
[WHERE condition]
ORDER BY F(arg_list)
LIMIT k
ON {LAST T {MINUTES | DAYS} | (T_start,T_end)}
```

```
* SELECT [CONTINUOUS] grouping_attr_list,
COUNT(attr_list)
FROM stream_name1 [,stream_name2,...]
[WHERE condition]
GROUP BY grouping_attr_list
LIMIT k
ON {LAST T {MINUTES | DAYS} | (T_start,T_end)}
```

UNIVERSITY

MQL Micro-Blogs Query Language

 Example: Retrieve the most frequent 10 keywords from tweets in Ukraine since February 18, 2014:

```
SELECT CONTINUOUS keyword, COUNT(*)
FROM twitter_stream
WHERE location WITHIN (52,44.7,39.91,21.8)
GROUP BY keyword
LIMIT 10 ON ("18 Feb 2014",∞)
```



Querying Spatial-Keyword Data

- Spatial-Keyword Queries
- Spatial-Keyword Query Languages
 - MQL
 - Atlas



The Atlas Query Language Spatial-Keyword Building Blocks

- Design building-block operators that can be composed to realize complex spatialkeyword queries
 - Similar to the relational SELECT, PROJECT, JOIN building-block operators

Mahmood et al., Atlas: on the expression of spatial-keyword group queries using extended relational constructs, SIGSPATIAL, 2016



The Specification of Atlas

```
SELECT {*|attr1 [AS alias][,attr2,..]}
FROM source_name1 [,source_
name2, ...]
[WHERE condition]
[ORDER BY F(arg_list)]
[LIMIT {k|condition}]
```

```
SELECT grp_attr1[AS alias][,grp_attr2,..}, AGGR_F [AS alias](attr_list)
FROM source_name1 [,source_ name2, ...]
[WHERE condition]
[PARTITION BY grp_att_list AS grp_alias]
[ORDER BY F(grp_arg_list)]
[LIMIT k]
[HAVING grp_condition]
```



Spatial-keyword Grouping Operators

- CONDITIONAL LIMIT
 - Retrieves a list of items until a condition is satisfied
- PARTITION BY
 - Similar to the traditional group-by, yet returns the group tuples (not aggregates over the groups)
- WITHIN_DIST(D)
 - Identifies groups of tuples such that the distance between every pair of tuples within a group is upperbounded by D



Example 1: Spatial-Keyword TOP-K Query

- Retrieve the top-k spatial-keyword objects ranked by a function of:
 - Distance between the object and the query location
 - Similarity between the keywords of the object and the query
- Example: Find the top three tuples that have any of the keywords: pizza, seafood, pasta
- Q.keywords and Q.loc



Example 2: Spatial-Keyword Group Query

Find groups of objects that collectively contain the keywords "cinema, restaurant, cafe" such that objects in a group are within 3 miles of each other. The groups are to be ranked by a function of the groups' inter-object distances and the distance from each of the groups to a specific location



Outline

- Introduction and Background
- Querying Spatial-Keyword Data
- Scalable Spatial-Keyword Systems
- Research Directions



Outline

- Introduction and Background
- Querying Spatial-Keyword Data
- Scalable Spatial-Keyword Systems
- Research Directions



Scalable Spatial-Keyword Systems

- Query-specific algorithms
- Extensions to general-purpose big data systems
- Spatial-keyword-only systems



Boolean kNN Processing over MapReduce



Li et al., Evaluating spatial keyword queries under the mapreduce framework, DASFAA 2012



Big Spatial-Keyword Processing Systems

- Query-Specific Algorithms
- Extensions to general-purpose big data systems
- Spatial-keyword-only systems



Tornado

- Adaptive and Distributed System for Spatial-Keyword Stream Processing
 - Extend the general-purpose Storm streaming system
 - Distribute the input streams to multiple worker processes

Mahmood et al., Adaptive processing of spatial-keyword data over a distributed streaming cluster, SIGSPATIAL 2018



The Storm Streaming System



http://www.business-software.com/wp-content/uploads/2014/09/Storm.png



Tornado: Adaptive and Distributed System for Spatial-Keyword Stream Processing





The Routing Layer

- Partition the space into non-overlapping rectangles
- Identify evaluators to which data objects and queries belong to



Tornado: Adaptive and Distributed System for Spatial-Keyword Stream Processing

 Evaluators use FAST^{*} to internally index continuous spatial-keyword queries



ERSITY

Limitations of Existing Spatial-Keyword Indexes

- Integrate a spatial index with a textual index
 - Do not account for the popularity of keywords across different regions
- Spatial Indexes
 - R-tree
 - High update overhead
 - Grid
 - High memory requirements due to replication of textual structures in multiple grid cells



Limitations of Existing Spatial-Keyword Indexes

- Textual Indexes
 - Ranked inverted list
 - Poor search performance for long posting lists

- Ordered keyword trie
 - Similar to the traditional trie
 - High memory requirements
 - Requires having a total order on the keywords
 - Requires knowing the entire vocabulary





- A textual index that accounts for the popularity of keywords
- Starts as a ranked inverted list (RIL)
- The frequent-keyword threshold θ is used to prevent having long posting lists
- A query is attached to its least frequent keyword
- Frequencies of keywords are not known apriori
 - progressively maintained



 When the number of queries attached to a keyword exceeds the frequent-keyword threshold θ, mark it as frequent







 Use more keywords of queries to improve textual discrimination



Frequencies Map

K1	3
K4	2



What happens when many queries have the exact same keywords?





FAST



Replication that increases the memory overhead



Memory Optimization in FAST



Object Matching: Searching FAST



Object Matching: Searching FAST



Object Matching: Searching FAST



UN

IVERSITY

PS²Stream

- Distributed publish/subscribe system
- Extends the storm streaming system





Scalable Spatial-Keyword Systems

- Query-Specific Algorithms
- Extensions to general-purpose big-data systems

Spatial-keyword-only systems



TwitterStand

Detection of news from Twitter



Sankaranarayanan et al., Twitterstand: news in tweets, SIGSPATIAL 2009



TwitterStand

Diplomats deliver ultimatum on Honduras coup Less that 1 hour ago - geardlankews 30 Read ANDER 875/14 Diplomats deliver ultimatum on Honduras coup-1011 tweets - Similar Stories - Original Source - Locations The news is all over Michael Jackson while North Korea is threating us and testing missiles North America with no news coverage Less than 1 hour age - Twitter-Streaming The news is all over Michael Jackson while North Korep is threating us and testing missiles with no news coverage EARDEN 1833 tweets - Similar Stories - Original Source - Locations 144.0 thinks jacksons dad Joe should not be anywhere near those kids he abused Michael and will do the same to Michaels kids Give Debbie rowe a go Lease these 1 hour ago - Twitten Streaming thinks acksons dad Joe should not be anowhere near those kids he abused Michael and will do the same to Nichaels kids Give Debbie rowe a go 574 basets - Similar Stories - Original Source - Locations Lough Marines establish positions in Afghan assault feed-fod Ocean Less than 1 hour ago - timesnews Eastrolle Marines establish positions in Alghan assault. Atlantic Ocean 209 tweets - Similar Stories - Original Source - Locations Ocean 62 IranElection Tehran Mousavi Iran neda Neda How to send an anonymous email 1 hours ago - Twitten-Streaming Southern Goran 62 IranElection Tehran Mous av Iran neda Neda Haw to send an anotymous email. Muccial! 5422 tweets - Similar Stories - Original Source - Locations See Anteretica Sanford to reveal schedule details Less than 1 hour ago - thestate Sanford to reveal actiedule details 7000 mil Virtual Earth" 1175 tweets - Similar Stories - Original Source - Locations



TAGHREED

- Answers spatial-keyword queries over micro-blogs including:
 - Spatial-keyword filter
 - Top-k frequent keywords
 - Top-k active users
 - Top-k famous users
 - Aggregates over micro-blogs
 - Top-k languages used while tweeting

Magdy et al., Taghreed: A System for Querying, Analyzing, and Visualizing Geotagged Microblogs, SIGSPATIAL 2014



TAGHREED



Outline

- Introduction and Background
- Querying Spatial-Keyword Data
- Scalable Spatial-Keyword Systems
- Research Directions



Research Directions

- Benchmarking
- Pipelined Evaluation
- Big Spatio-Temporal-keyword Processing



Benchmarking Big Spatial-Keyword Systems

- Spatial-only benchmarks
 - Jackpine
 - BerlinMOD
- Relational benchmarks – TPC-H



http://csng.cs.toronto.edu/projects/20 http://dna.fernuni-hagen.de/secondo/BerlinMOD/BerlinMOD.html



Benchmarking Big Spatial-Keyword Systems

- Existing spatial-keyword benchmarks
 - Limited set of queries
 - Specific use cases, e.g., social-network analysis

Chen et al., Spatial keyword query processing: An experimental evaluation, PVLDB 2013 Panda et al., Performance Evaluation of Social Network Using Data Mining Techniques, Computational Social Networks 2012



Benchmarking Big Spatial-Keyword Systems

- Large spatial-keyword datasets
- Support of a wide range of spatial-keyword queries
- Support of several evaluation scenarios
 - Batching
 - Streaming



Research Directions

- Benchmarking
- Pipelined Evaluation
- Big Spatio-Temporal-keyword Processing



Pipelined Evaluation of Big Spatial-keyword Queries

- Design building-block operators that can be composed to realize complex spatialkeyword queries
 - Similar to the relational SELECT, PROJECT, JOIN building-block operators





Research Directions

- Benchmarking
- Pipelined Evaluation
- Big Spatio-Temporal-keyword Processing



Big Spatio-Temporal-keyword Processing

- Spatial-keyword data are often associated with temporal attributes

 Time-stamp of tweets
- We need to investigate
 - Indexing
 - Processing
 - Query processing and Query languages



Thank you!





The authors acknowledge the support of the National Science Foundation under Grant Numbers III-1815796 and IIS-1910216.

