

High-Performance Analytics on Large-Scale GPS Taxi Trip Records in NYC

Jianting Zhang

Department of Computer Science The City College of New York

Outline

•Background and Motivation

•Parallel Taxi data management on GPUs

•Efficient shortest path computation and applications



Taxicabs

- •13,000 Medallion taxi cabs
- •License priced at \$600, 000 in 2007

•Car services and taxi services are separate

Taxi trip records

•~170 million trips (300 million passengers) in 2009

•1/5 of that of subway riders and 1/3 of that of bus riders in NYC





Over all distributions of trip distance, time, speed and fare (2009)



Other types of Origin-Destination (OD) Data



Cellular phone calls

Social network activities



- How to manage OD data?
 - Geographical Information System (GIS)
 - Spatial Databases (SDB)
 - Moving Object Databases (MOD)
- How good are they?
 - Pretty good for small amount of data $\textcircled{\odot}$
 - But, rather poor for large-scale data \otimes

- Example 1:
 - Loading 170 million taxi pickup locations into PostgreSQL
 - UPDATE t SET PUGeo = ST_SetSRID(ST_Point("PULong","PuLat"),4326);
 - 105.8 hours!
- Example 2:
 - Finding the nearest tax blocks for 170 million taxi pickup locations using open source libspatiaindex+GDAL
 - 30.5 hours!

I do not have time to wait... Can we do better?





Cloud computing+MapReduce+Hadoop



Multicore CPUs



GPGPU Computing: From Fermi to Kepler







Nvidia GTX Titan GPU sold at Amazon

EVGA EVGA GeForce GTX TITAN SuperClocked 6GB GDDR5 384bit, Dual-Link DVI-I, DVI-D, HDMI,DP, SLI Ready Graphics Card Graphics Cards 06G-P4-2791-KR by EVGA

Price: \$1,128.90

Only 1 left in stock. Ships from and sold by <u>DVR Sales</u>.

- Base Clock: 876 MHz
- Boost Clock: 928 MHz
- Memory Clock: 6008 MHz Effective
- CUDA Cores: 2688
- Memory: 6144MB GDDR5 384bit
- Show more
- 4 new from \$1,121.99

4.5 Teraflops of single precision and 1.3 Teraflops of double precision



ASCI Red: 1997 First 1 Teraflops (sustained) system with 9298 Intel Pentium II Xeon processors (in 72 Cabinets) \$\$\$?

Space?

Power?



•Feb. 2013

- •7.1 billion transistors (551mm²)
- •2,688 processors
- •Max bandwidth 288.4 GB/s
- •PCI-E peripheral device
- •250 W (17.98 GFLOPS/W -SP)
- Suggested retail price: \$999

What can we do today using a device that is more powerful than ASC I read 16 years ago?

- The goal is to design a data management system to efficiently manage large-scale OD/trajectory data on massively data parallel GPUs
- With the help of new data models, data structures and algorithms
- To cut the runtimes from hours to seconds on a single commodity GPU device
- And support interactive queries and visual explorations

Outline

•Background and Motivation

•Parallel Taxi data management on GPUs

•Efficient shortest path computation and applications



(Zhang, Gong, Kamga and Gruenwald, 2012)





⁽Zhang, Gong, Kamga and Gruenwald, 2012)



(Zhang, Gong, Kamga and Gruenwald, 2012)



Single-Level Grid-File based Spatial Filtering

- Data
 - Taxi trip records: 300 million in two years (2008-2010), ~170 million in 2009 (~150 million in Manhattan)
 - NYC DCPLION street network data: 147,011 street segments
 - NYC Census 2000 blocks: 38,794
 - NYC MapPluto Tax blocks: 735,488 in four boroughs (excluding SI) and 43,252 in Manhattan
- Hardware
 - Dell T5400 Dual Quadcore CPUs with 16 GB memory
 - Nvidia Quadro 6000 with 448 cores and 6 GB memory

Top: grid size =256*256 resolution=128 feet Right: grid size =**8192*8192** resolution=4 feet

Spatial Aggregation

9,424/326=30X (8192*8192)

Temporal Aggregation

1709/198=8.6X (minute)

1598/165 = 9.7X (hour)

(Zhang, You and Gruenwald, 2012)

	P2N-D	P2P-T	P2P-D
CPU time	_	15.2 h	30.5 h
GPU Time	10.9 s	11.2 s	33.1 s
Speedup	_	4,900X	3,200X

Algorithmic improvement: 3.7X Using main-memory: 37.4X Parallel Acceleration:

24.3X

(Zhang, You and Gruenwald, 2012), (Zhang and You 2012a) Zhang and You 2012b)

Outline

•Background and Motivation

•Parallel Taxi data management on GPUs

•Efficient shortest path computation and applications

•Mapping Betweenness Centralities

•Outlier Detection

-overview

- Shortest path computation
 - Dijkstra and A*
 - New generation algorithms
 - Contraction Hierarchy (CH) based

• Network Centrality (Brandes, 2008)

- •Can be easily derived after shortest paths are computed
- •Mapping node/edge between centrality can reveal the connection strengths among different parts of cities

Mapping Betweenness Centralities

- •166 million trips, 25 million unique
 - •Shortest path computation completes in less than 2 hours (5,952 seconds) on a single CPU core (2.0 GHZ)
 - •4200 pairs computation per second
 - •3 orders of magnitude faster than ArcGIS NA

Mapping of Computed Shortest Paths Overlaid with NYC Community Districts Map

Mapping Betweenness Centralities (All hours)

Mapping Betweenness Centralities (bi-hourly)

00H

02H

06H

18H

12H

- 10001 - 100000

22H

Legend:

14H

The data is not as clean as we had thought...outlier detection

•Existing approaches for outlier detection for urban computing

•Thresholding: e.g. 200m < dist < 30km

•Locating in unusual ranges of distributions

•Spatial analysis: within a region or a land use type

•Matching trajectory with road segments – treat unmatched ones as outliers

•Some techniques require complete GPS traces while we only have O-D locations

•Large-scale shortest path computing has not been used for outlier detection

The data is not as clean as we had thought...outlier detection

- In addition:
 - Some of the data fields are empty
 - Pickup and drop-off locations can be in Hudson River
 - The recorded trip distance/duration can be unreasonable

Outlier detections for data cleaning are needed

-Discussion on outlier detection

- The approach is approximate in nature
 - Taxi drivers do not always follow shortest path
 - Especially for short trips and heavily congested areas
 - But we only care about aggregated centrality measurements and the errors have a chance to be cancelled out by each other
- Increasing D₀ will reduce # of type I outliners, but the locations might be mismatched with segments
- Reducing D₁ and/or W will increase # of type II outliers but may generate false positives.

-Results of outlier detection

Ongoing Research @ GeoTECI

•CudaGIS - a general purposed GIS on GPUs

•More GIS modules in addition to indexing/query processing

•Trajectory data management on GPUs

- •Online moving point location updating
- •Segmentation/simplification/compression
- •Matching with road networks
- •Aggregation and warehousing
- •Indexing and query processing →similarity join
- •Data mining (moving cluster, convoy, swarm...)

http://www-cs.ccny.cuny.edu/~jzhang/ jzhang@cs.ccny.cuny.edu