# map-D

# map-D
## data refined

**map-D** A GPU Database for Real-Time Big Data Analytics and Interactive Visualization

SC13 Denver
#mapDsc13

Tom Graham
Todd Mostak

**map-D?** super-fast database
built into GPU memory

**Do?** world's fastest
real-time big data analytics
interactive visualization

**Demo?** twitter analytics platform
1billion+ tweets
milliseconds

# #mapDsc13
# #NVIDIA
# #SC13

# Core Innovation

Map-D's database architecture is integrated into the memory on GPUs

Takes advantage of the memory bandwidth and massive parallelism on multiple GPUs and clusters

Runs 70-1000x faster than other in-memory databases and analytics platforms

Any kind of data

# #HAIYAN

**1billion+ tweets on 8 NVIDIA Tesla K40s**

2,880 x 8 = 23,040 cores            2.3 TB/sec memory bandwidth
12 x 8 = 96GB memory            >30 teraflops compute power

**Nothing is pre-computed!**
**Streaming live tweets**
**Interactive and real-time analytics**

# map-D overview

- SQL-enabled database (not a GPU accelerator)
- **Real-time search of any size dataset in milliseconds**
- Interactive visualizations generated on the fly
- **Compatible with any type of data**
- Scales to any size of dataset
- **Live data streams onto the system**
- Powered by inexpensive, off-the-shelf hardware
- **1000+ analytic/visualization queries per second**
- Optimized for GPUs but also runs on CPUs, Phi, AMD and mobile chips

**#mapDsc1**

# 1billion+ Tweetmap

**500 million tweets a day   =   7-10 million 'geocoded'**

**Tweet = more than just 140 characters:**
- **geo coordinates**
- **timestamp**
- **user and follower information**
- **reply information**
- **#hashtags**
- **host platform**

**Tweet volume and velocity is a massive challenge**

**Need new tools to interactively visualize data**

**#mapDsc1**

# 1billion+ Tweetmap

**Correlate with external and internal data sets**
- Brand preference vs census district income
- Tweet density by region (chloropleth)

**Deep analysis of content**
- What product, show, or person is discussed over time
- What opinion is being expressed 'sentiment analysis'

# "Shared Nothing" Processing
## Multiple GPUs, with data partitioned between them

**Filter**
text ILIKE 'rain'

**Filter**
text ILIKE 'rain'

**Filter**
text ILIKE 'rain'

**Node 1**

**Node 2**

**Node 3**

# Tweet Indexing on GPU

## Encode tweets using a "dictionary"

| Filter | Filter |
|---|---|
| text ILIKE 'rain' | SELECT tweetid FROM words WHERE id = 57663 |

| Word | Encoding |
|---|---|
| ... | ... |
| Rain | 57663 |
| Rainbow | 57664 |
| Rainman | 57665 |
| Rainy | 57666 |
| ... | ... |

# Filtering in Parallel

- Column-oriented execution
  - Avoids wasting memory bandwidth

**Filter:**
SELECT tweet id FROM words WHEREid = 57663

- Filter:
  - Produce bitmap of tweets to read
  - Read tweets, increment output bins in bitmap

| TweetId | WordId | TweetId | Lat | Lon |
|---------|--------|---------|-------|-------|
| ... | ... | ... | ... | |
| 1 | 57663 | 1 | -41.5 | 23.1 |
| 2 | 57664 | 2 | -41.7 | 77.4 |
| 2 | 27 | 3 | -37.4 | 48.2 |
| 3 | 8841 | 4 | 28.4 | -44.0 |
| ... | ... | ... | ... | |

**Data Tables Reside in GPU Memory**

# Filtering in Parallel

- **1000+ GPU threads**
- **Running in "warps"**
- **Threads in same warp run the exact same instructions**
  - **Need same amount of data to be efficient**

| TweetId | WordId |
|---------|--------|
| ... | ... |
| 1 | 57663 |
| 2 | 57664 |
| 2 | 27 |
| 3 | 8841 |
| ... | ... |

Warp 1

Warp 2

Warp 3

**Bitmap**

| |
|---|
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |

Tweet 1

Tweet n

# Filtering in Parallel

- **1000+ GPU threads**
- **Running in "warps"**
- **Threads in same warp run the exact same instructions**
  - **Need same amount of data to be efficient**

| TweetId | WordId |
|---------|--------|
| ... | ... |
| 1 | 57663 |
| 2 | 57664 |
| 2 | 27 |
| 3 | 8841 |
| ... | ... |

Warp 1
Warp 2
Warp 3

**Bitmap**

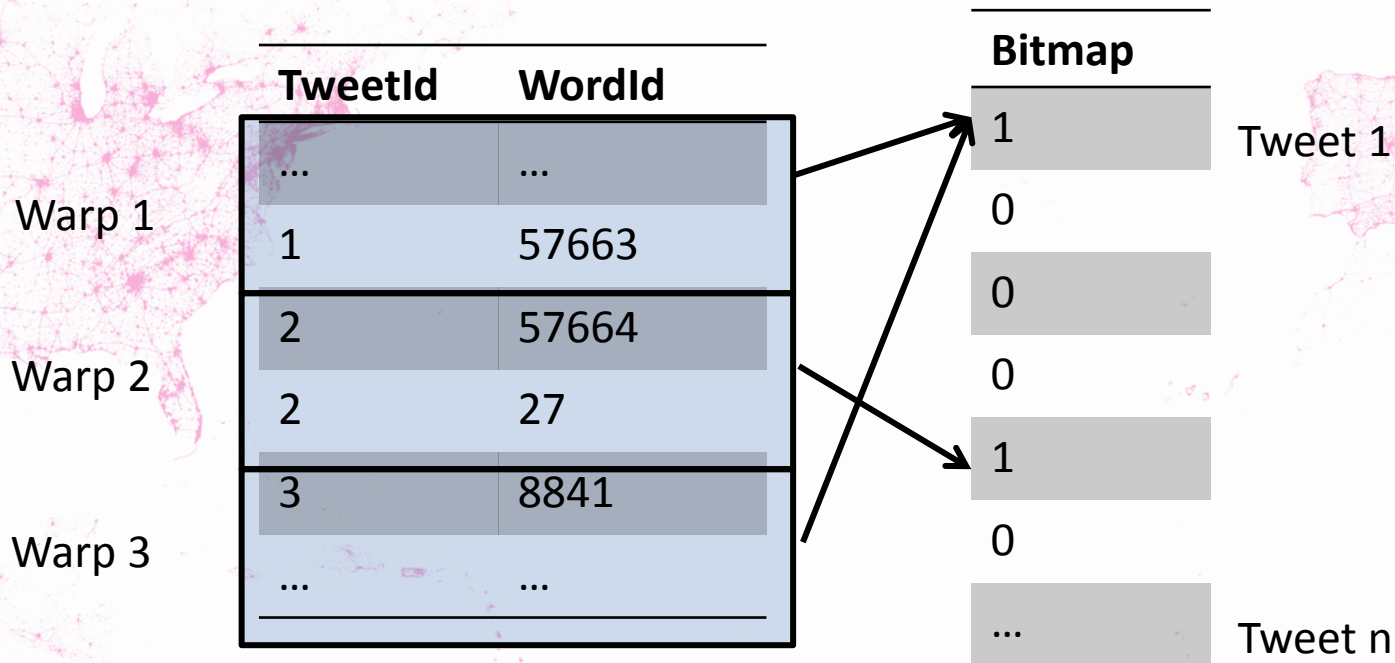| |
|---|
| 1 |
| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| ... |

Tweet 1

Tweet n

# Filtering in Parallel

- **1000+ GPU threads**
- **Running in "warps"**
- **Threads in same warp run the exact same instructions**
  - **Need same amount of data to be efficient**

| | TweetId | WordId |
|---|---|---|
| Warp 1 | ... | ... |
| | 1 | 57663 |
| Warp 2 | 2 | 57664 |
| | 2 | 27 |
| Warp 3 | 3 | 8841 |
| | ... | ... |

**Bitmap**

| Bitmap | |
|---|---|
| 1 | Tweet 1 |
| 0 | |
| 0 | |
| 0 | |
| 1 | |
| 0 | |
| ... | Tweet n |

# Filtering in Parallel

- **1000+ GPU threads**
- **Running in "warps"**
- **Threads in same warp run the exact same instructions**
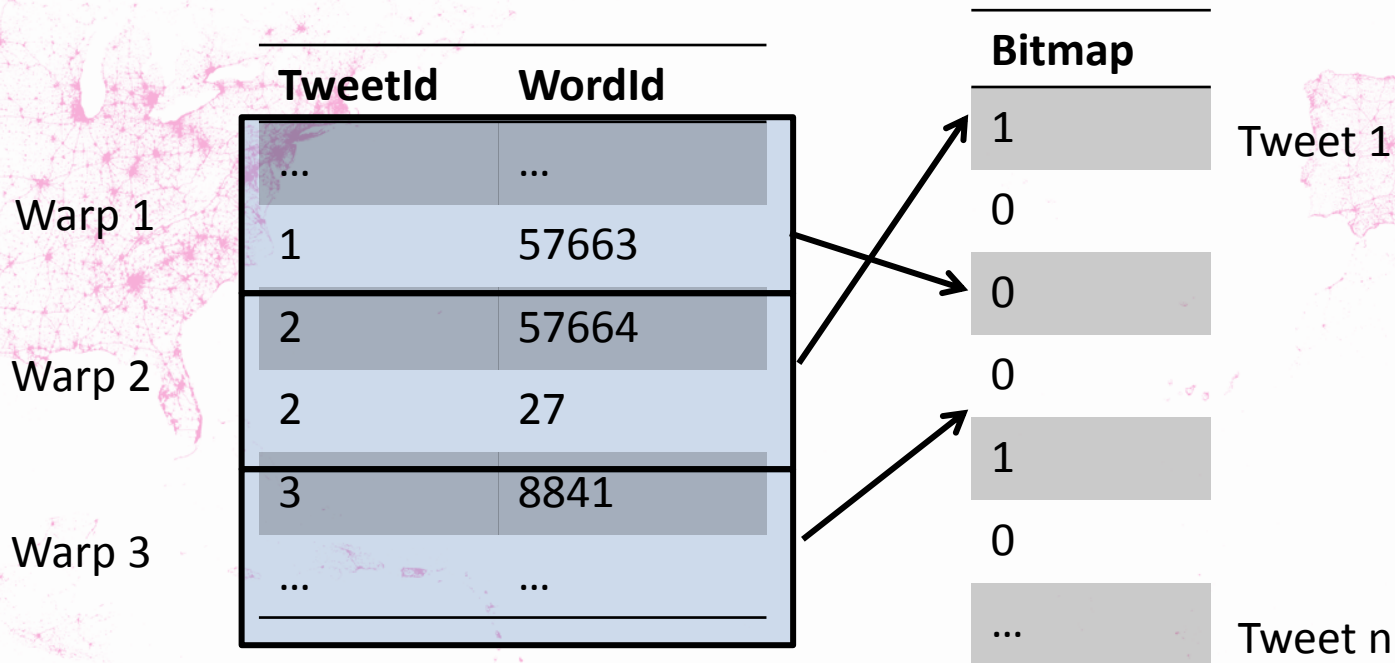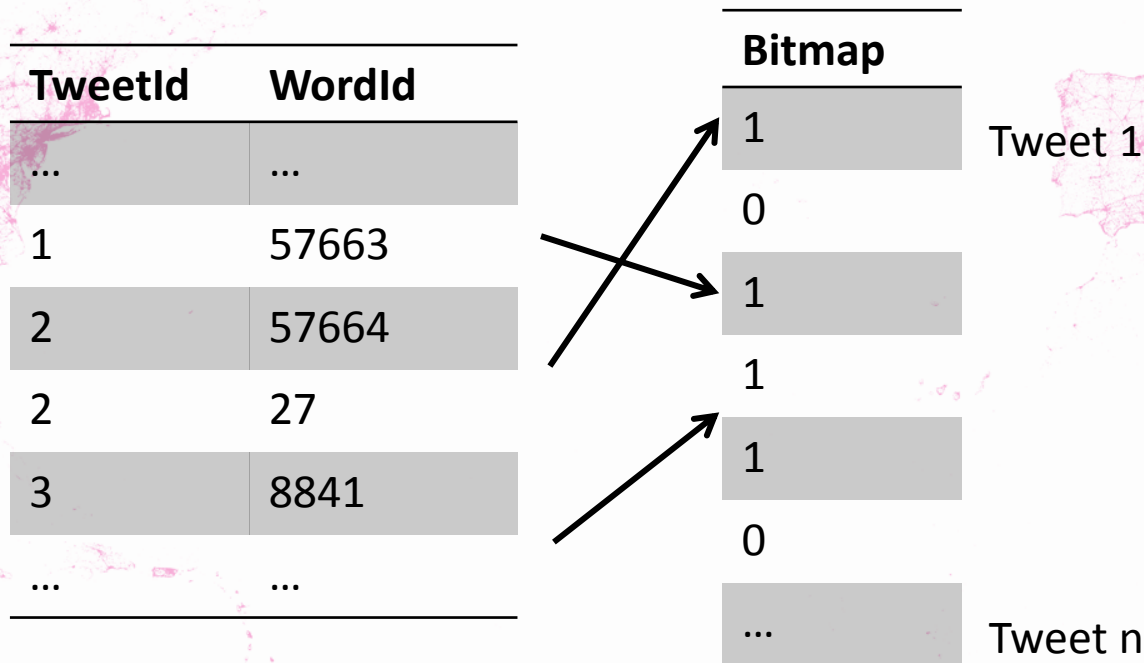  - **Need same amount of data to be efficient**

| TweetId | WordId |
|---------|--------|
| ... | ... |
| 1 | 57663 |
| 2 | 57664 |
| 2 | 27 |
| 3 | 8841 |
| ... | ... |

**Bitmap**

| |
|---|
| 1 |
| 0 |
| 1 |
| 1 |
| 1 |
| 0 |
| ... |

Tweet 1

Tweet n

# Filtering in Parallel

- **1000+ GPU threads**

- **Running in "warps"**

- **Threads in same warp run the exact same instructions**
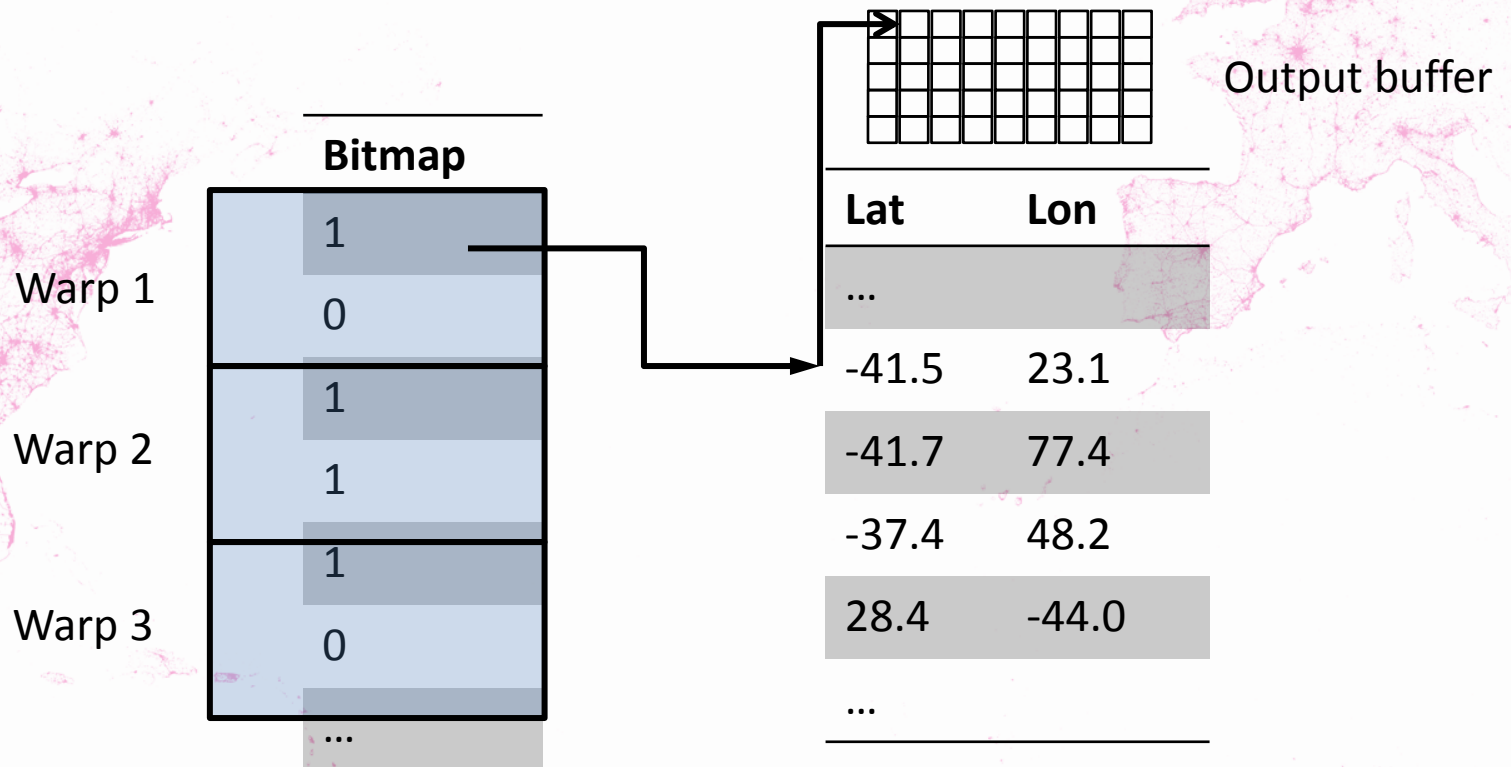  - **Need same amount of data to be efficient**

| Bitmap | |
|--------|--|
| 1 | Tweet 1 |
| 0 | |
| 1 | |
| 1 | |
| 1 | |
| 0 | |
| ... | Tweet n |

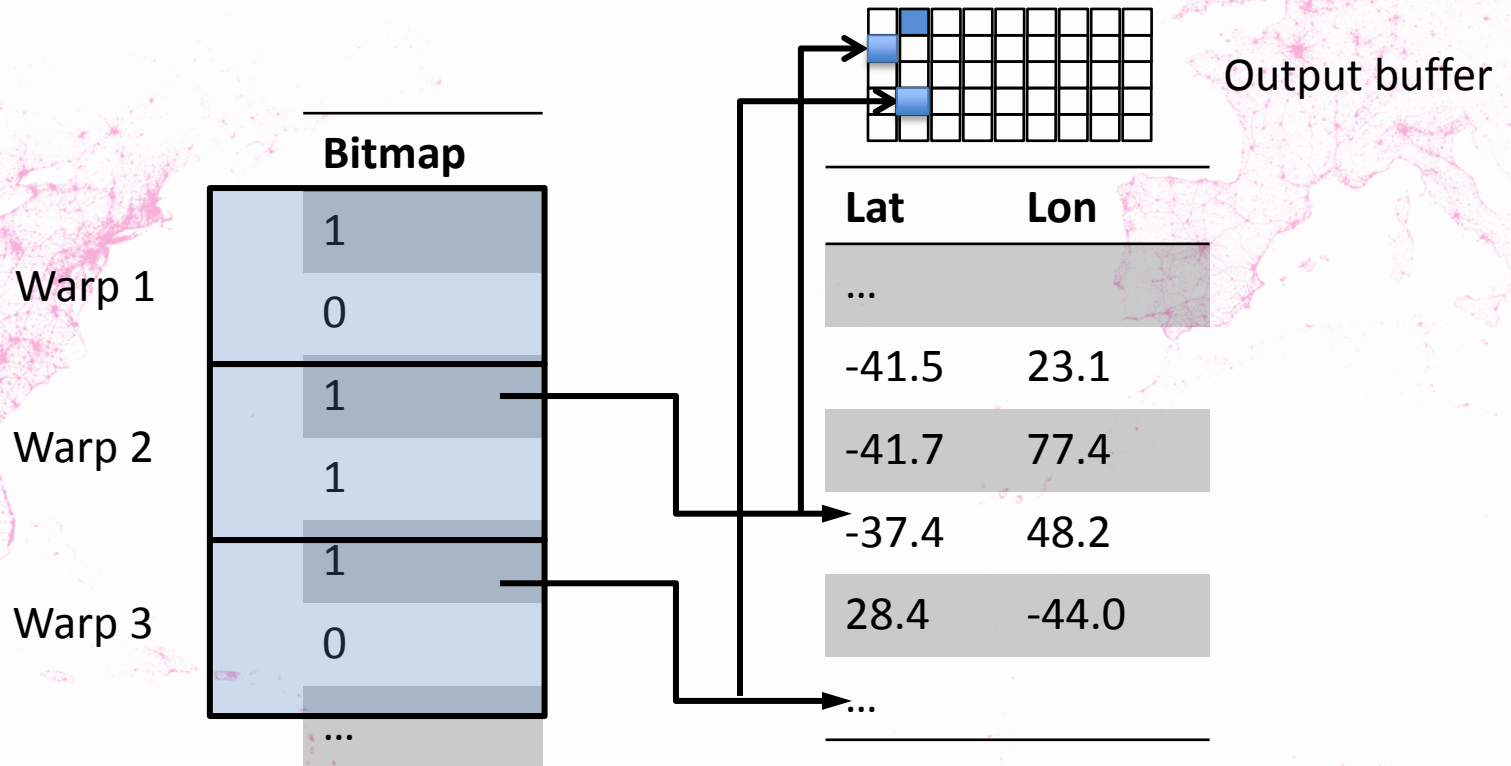| Lat | Lon |
|-----|-----|
| ... | |
| -41.5 | 23.1 |
| -41.7 | 77.4 |
| -37.4 | 48.2 |
| 28.4 | -44.0 |
| ... | |

# Filtering in Parallel

- **1000+ GPU threads**
- **Running in "warps"**
- **Threads in same warp run the exact same instructions**
  - **Need same amount of data to be efficient**

Output buffer

**Bitmap**

| | 1 |
| Warp 1 | |
| | 0 |
| | 1 |
| Warp 2 | |
| | 1 |
| | 1 |
| Warp 3 | |
| | 0 |
| | ... |

| Lat | Lon |
| --- | --- |
| ... | |
| -41.5 | 23.1 |
| -41.7 | 77.4 |
| -37.4 | 48.2 |
| 28.4 | -44.0 |
| ... | |

# Filtering in Parallel

- **1000+ GPU threads**
- **Running in "warps"**
- **Threads in same warp run the exact same instructions**
  - **Need same amount of data to be efficient**



Output buffer

| Bitmap | | Lat | Lon |
|--------|---|-----|-----|
| 1 | | ... | |
| 0 | | -41.5 | 23.1 |
| 1 | | -41.7 | 77.4 |
| 1 | | -37.4 | 48.2 |
| 1 | | 28.4 | -44.0 |
| 0 | | ... | |
| ... | | | |

Warp 1
Warp 2
Warp 3

# Effective big data tools

Democratization of big data analytics

**Interaction with live data streams**

Socialization of data driven insight

**Map-D is open source**

# Map-D is a startup

## Supported enterprise-grade database
- Appliance or in the cloud

## Platform integration
- Cloudera **|** NVIDIA **|** Software AG

## Tailored database and analytics solutions
- Twitter **|** Major League Baseball
  Sunlight Foundation **|** Leidos

## Free, public big data tools powered by Map-D
- Harvard's Worldmap **|** National Geographic
  Smithsonian Center for Astrophysics **|** MIT CSAIL

#mapDsc1

# Play with our live demo

# mapd.csail.mit.edu

# Who has been tweeting at SC13?

# #mapDsc13

#mapDsc1

# Special thanks



Prof Sam Madden, MIT CSAIL

# map-D

1billion+ Demo in NVIDIA booth

@datarefined

info@map-d.com

map-d.com