



**ARESDB**

A GPU-Powered Real-time Analytics Engine

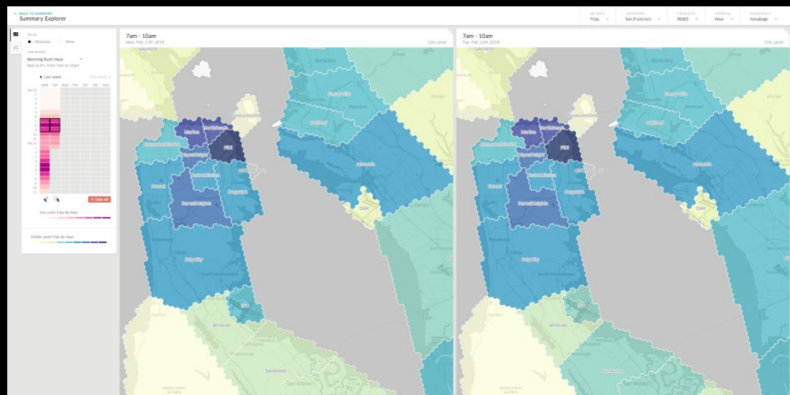
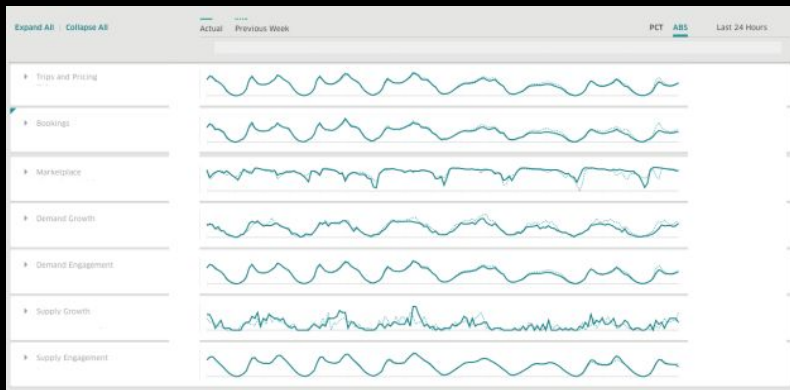
Uber

# Agenda

- Real-time analytics at Uber
- Leveraging GPU for real-time analytics
- AresDB Architecture and Features
- Learnings From GPU Programming
- Future Directions

# Real-time analytics at Uber

# Real-time Analytics Use Cases at Uber



Rider eyeballs



Open car information



Dynamic Pricing

# Real-time Analytics Use Cases Categorization

	Dashboards	Decision Systems	Ad hoc Queries
Dataset	Subset	Subset	All data
Ingestion Latency	Seconds to Minutes	Seconds to Minutes	Minutes
Query Pattern	Well known	Well known	Arbitrary
Query QPS	Medium	High	Low
Query Latency	Sub seconds	Sub seconds	Minutes
Target Users	City OPS, Executives	Engineers (application developers)	Data Scientists, Analytics, City OPS

# Mission of AresDB

- ❑ sub-sec level query latency
- ❑ second to min level ingestion latency

- ❑ High availability (4 9s)
- ❑ High data accuracy (3 9s)

- ❑ Uber scale and beyond

**Build a **fast**, **reliable** and **scalable** analytics platform solution to power Uber's Real-Time business intelligence**

# Leveraging GPU for Real-Time Analytics

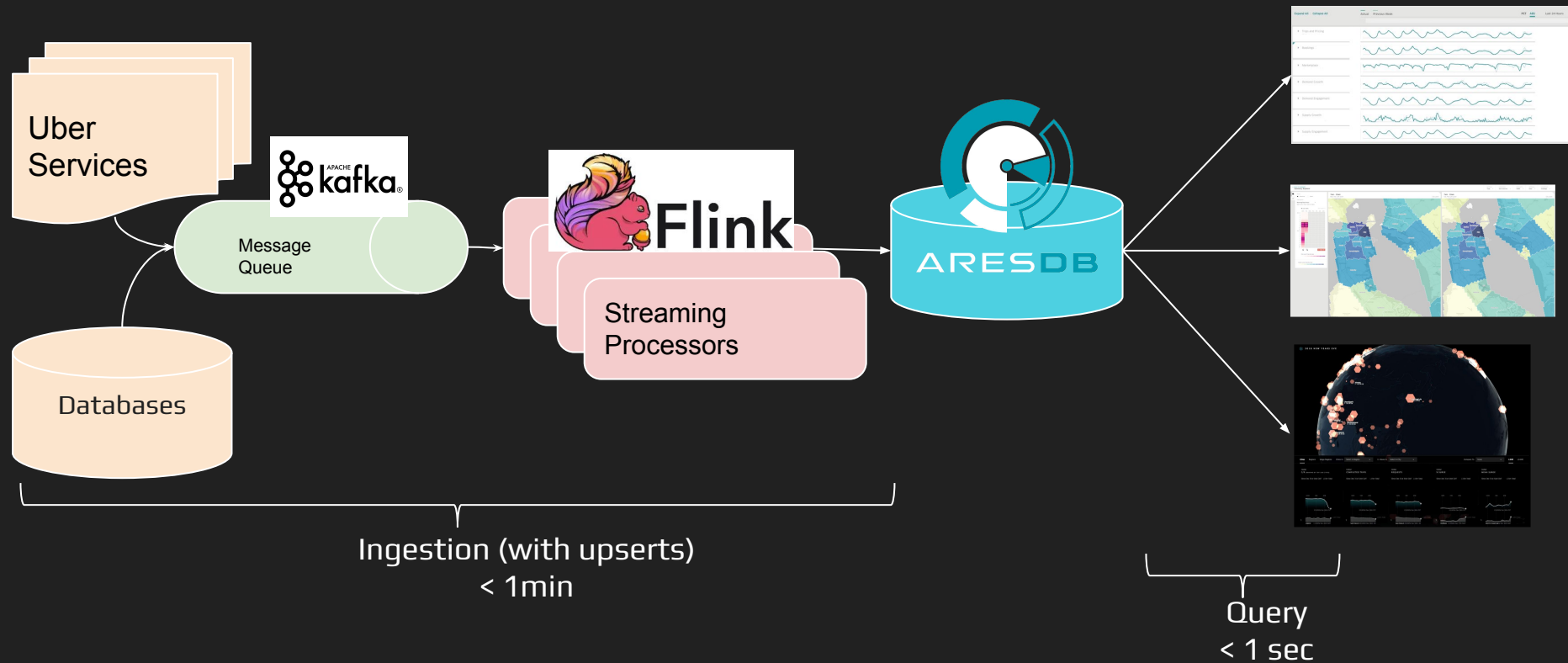
# AresDB: A GPU-Powered Real-time Analytics Engine



- High-efficiency storage
- Low-latency ingestion
- Sub-second query response time
- Feature set for real-time analytics



# How is AresDB used at Uber



# The Problem: Time-series Analytics

Computing **measures** by **dimensions** on time series data

request_at	city_id	fare
2017-04-13 10:25	1	15.3
2017-04-13 11:10	1	7.5
2017-04-14 10:35	1	20.1
2017-04-14 11:40	5	12.1
2017-04-14 15:45	5	5.6



day(request_at)	city_id	sum(fare)
2017-04-13	1	22.8
2017-04-14	1	20.1
2017-04-14	5	17.7

Why are GPUs well-suited ?

# GPU vs CPU

Intel® Xeon® Processor  
E5-2620 v3

Tesla P100

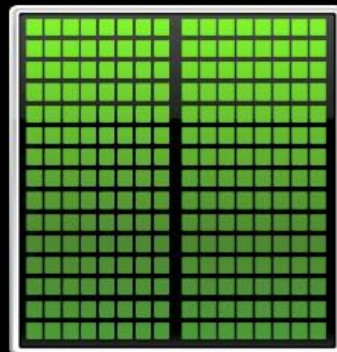
12 Threads

59GB/s

500 GFLOPS



CPU



GPU

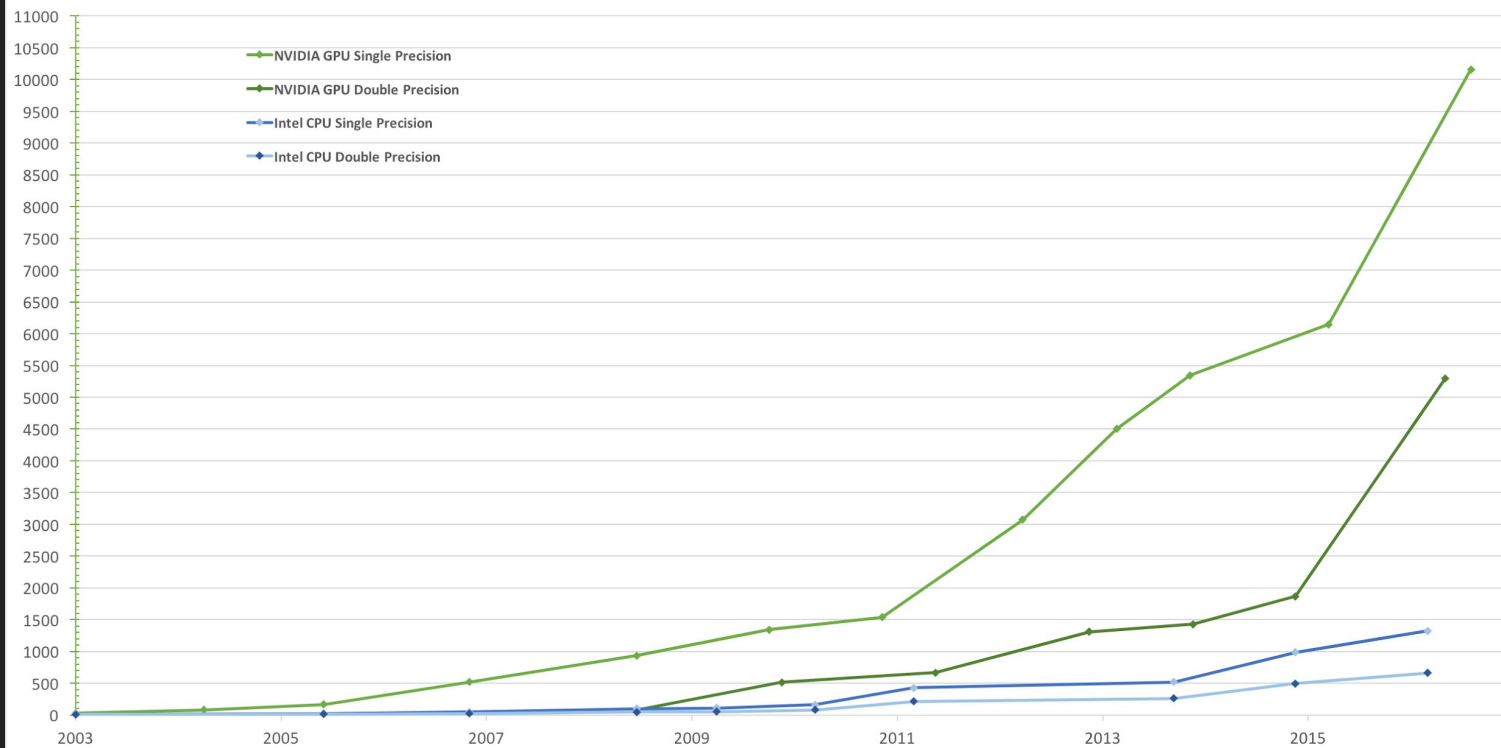
3584 CUDA  
cores

549G/s

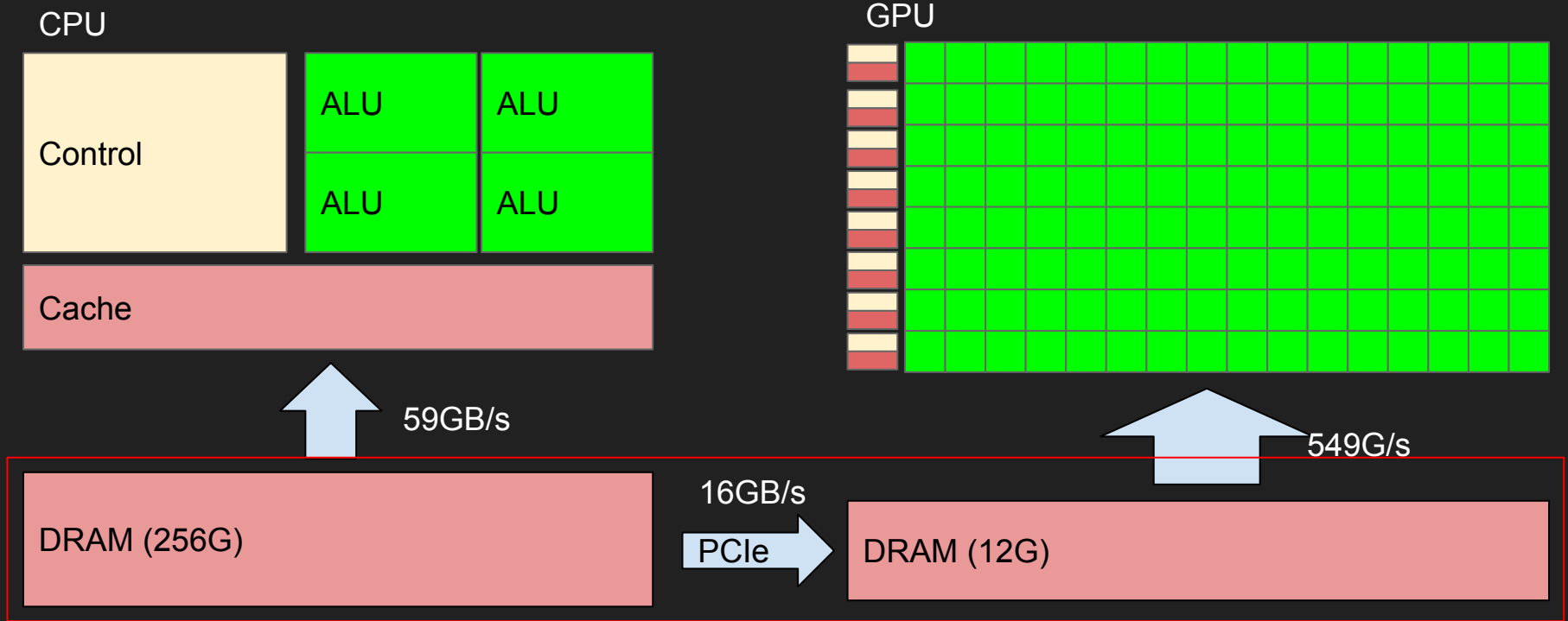
9.3 TFLOPS

# GPU vs CPU

Theoretical GFLOP/s at base clock



# Hardware Storage Choices



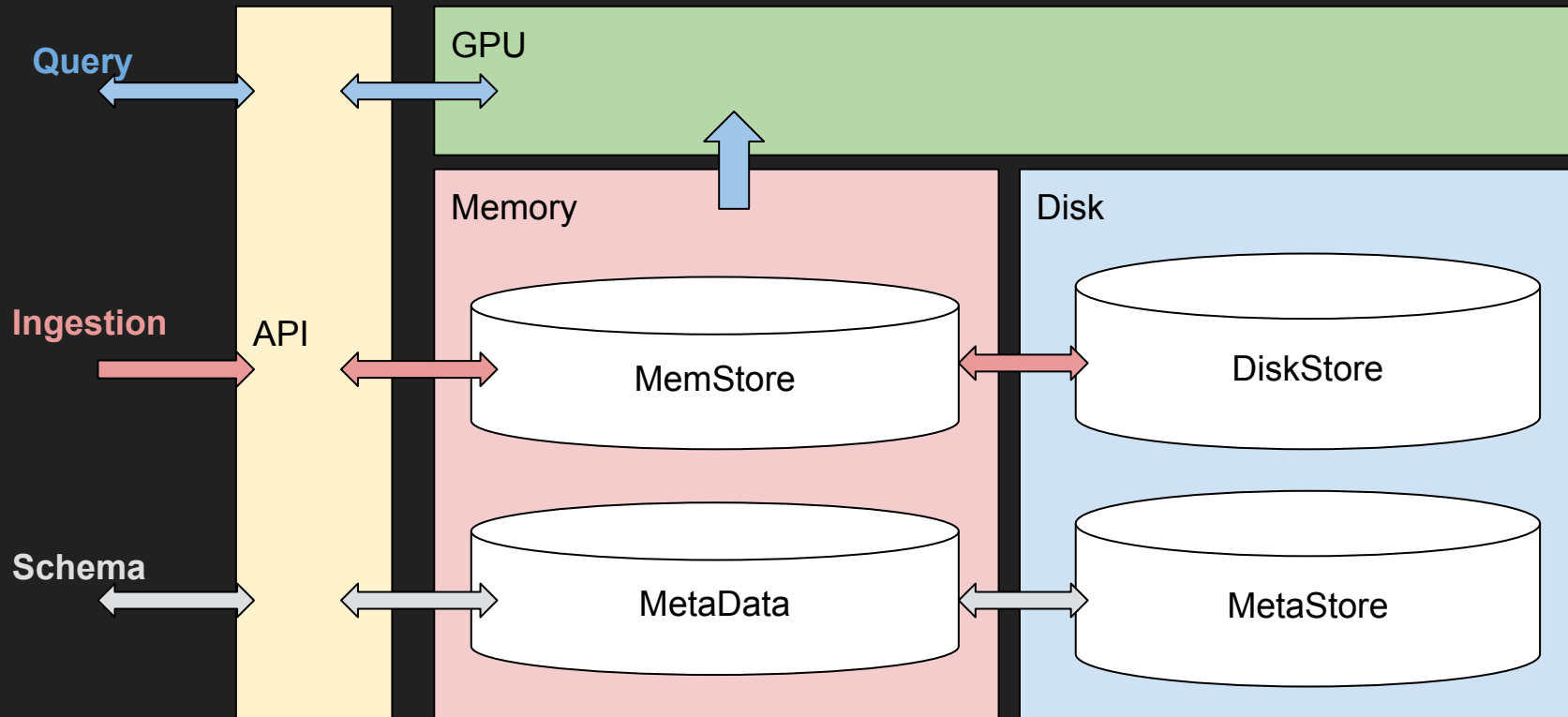
# Hardware Storage Choices

	Capacity	Bandwidth to GPUs	Ingestion Method	Comment
HDD	<b>8TB</b>	<b>100MB/s</b>	File writing similar to traditional databases	Unable to feed data fast enough to fully utilize GPUs
SSD	<b>4TB</b>	<b>600MB/s</b>		
NVMe	<b>2TB</b>	<b>3GB/s</b>		And also expensive
Host Memory	<b>256GB</b>	<b>15GB/s</b> per side; <b>30GB/s</b> two sides	Memory writing	Limited by PCIe bandwidth
GPU Memory	<b>12GBx8</b>	<b>500GB/s</b> on the same GPU; <b>15GB/s</b> across GPUs	Sharding across multiple GPUs; Complex memory writing	Tight coupling of storage and computation; ingestion is challenging

# AresDB Architecture and Features



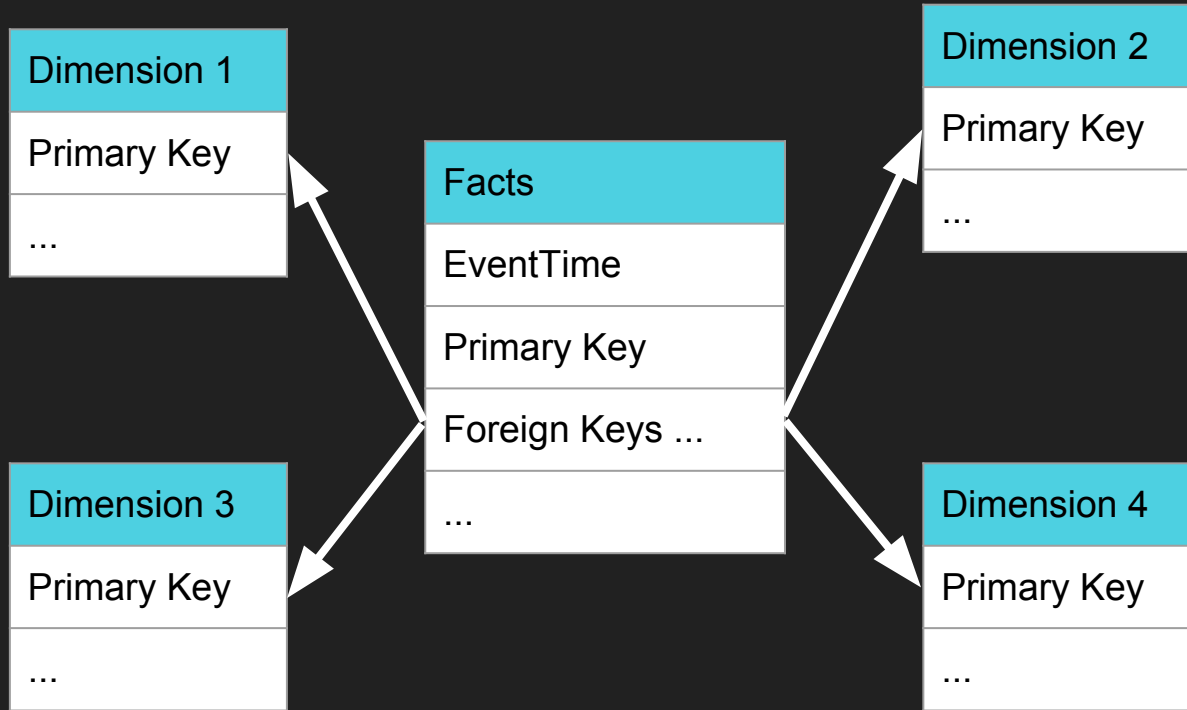
# AresDB Architecture: Single Instance



# Fact/Dimension Table

- Fact table
  - Facts about a business process
  - Each record associated with an event time (grows with time)
  - E.g. trips, orders, ...
- Dimension table
  - Descriptive attributes/dimensions
  - E.g. product catalogs, cities, ...

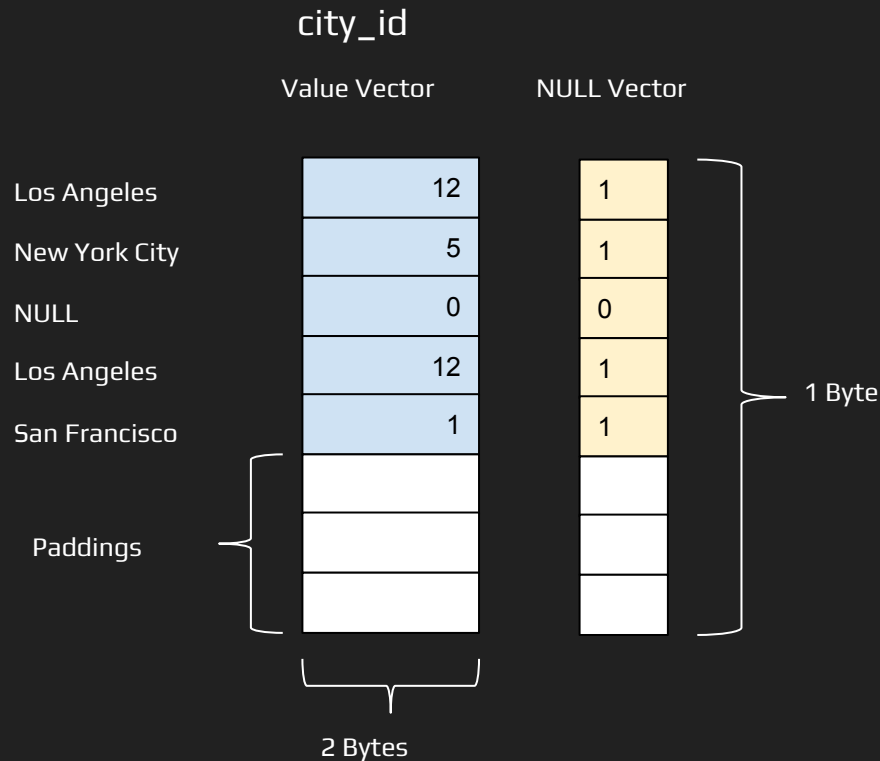
# Star Schema



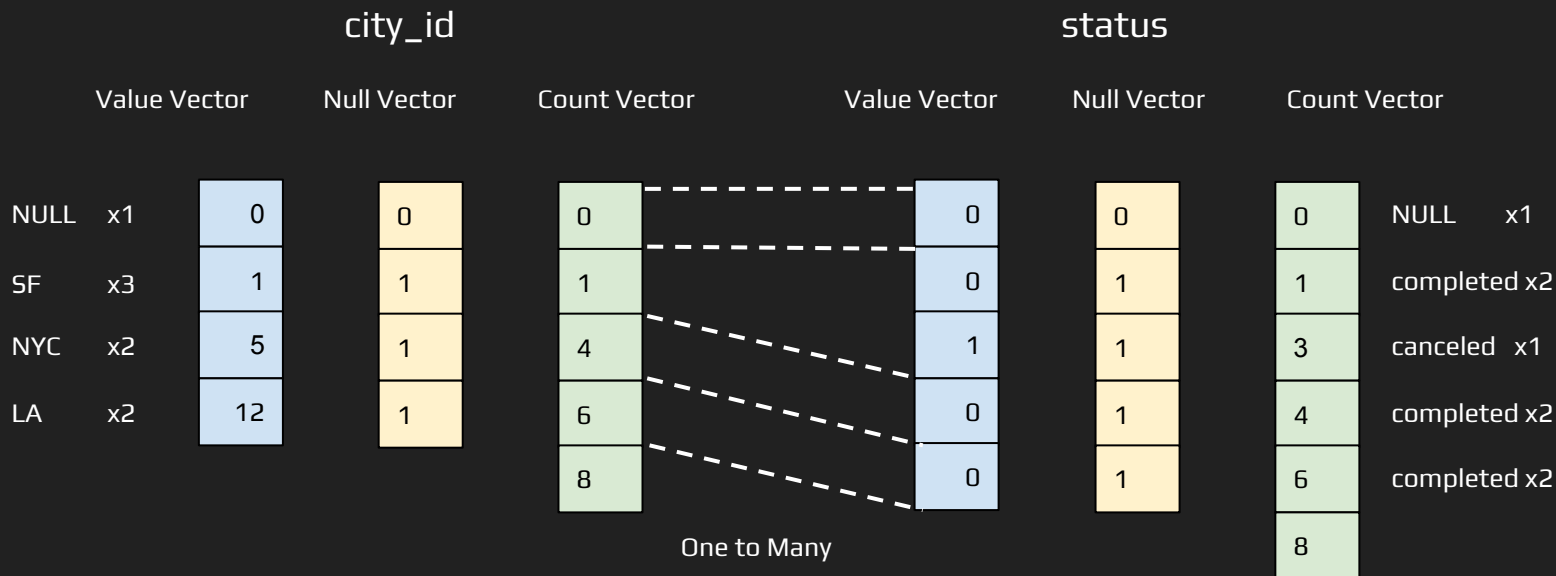
# Feature Highlights

- In-Memory Columnar Storage
- Real-time upserts
- GPU powered query engine
- Analytical Query Feature Set
  - Time zone, Time Filter, Time Bucketization
  - Geospatial analytics
  - Fact/dimension table joins
  - Hyperloglog

# Columnar Storage: uncompressed

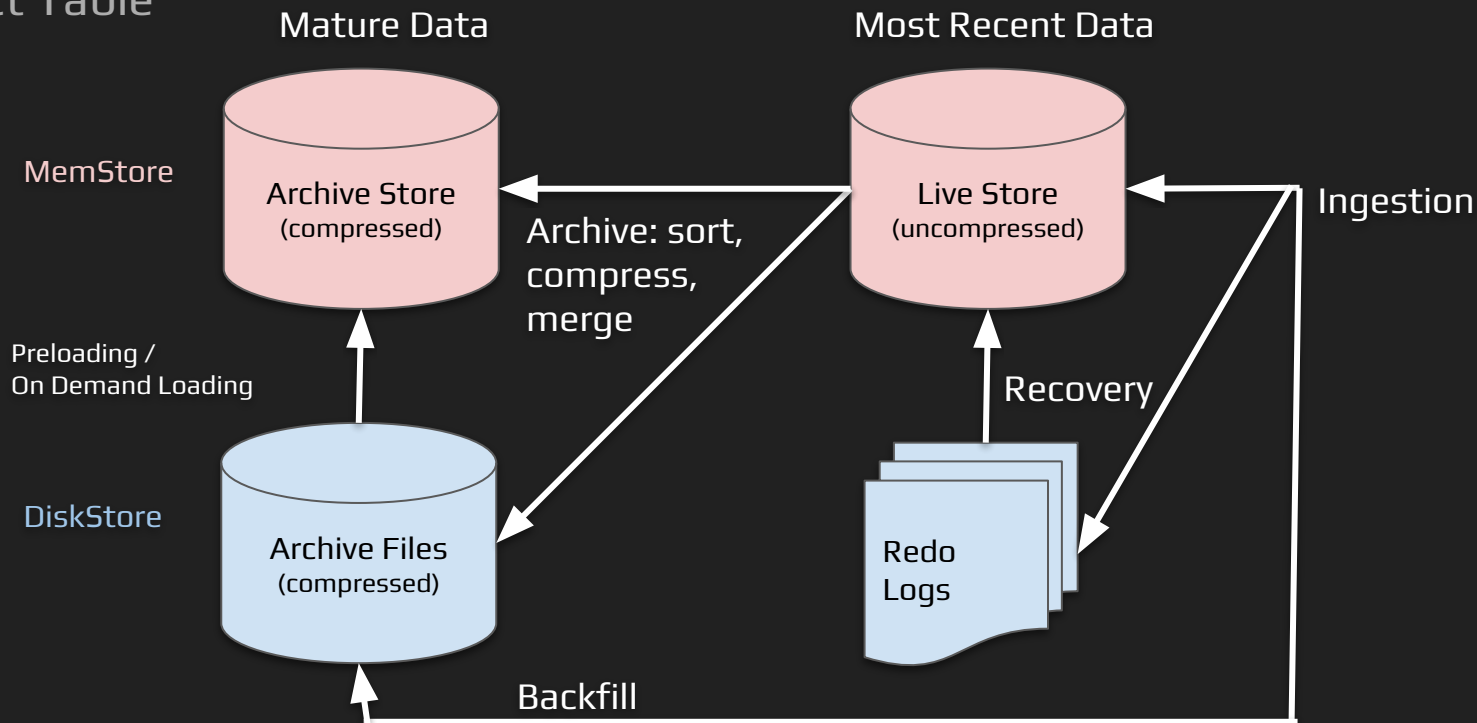


# Columnar Storage: compressed



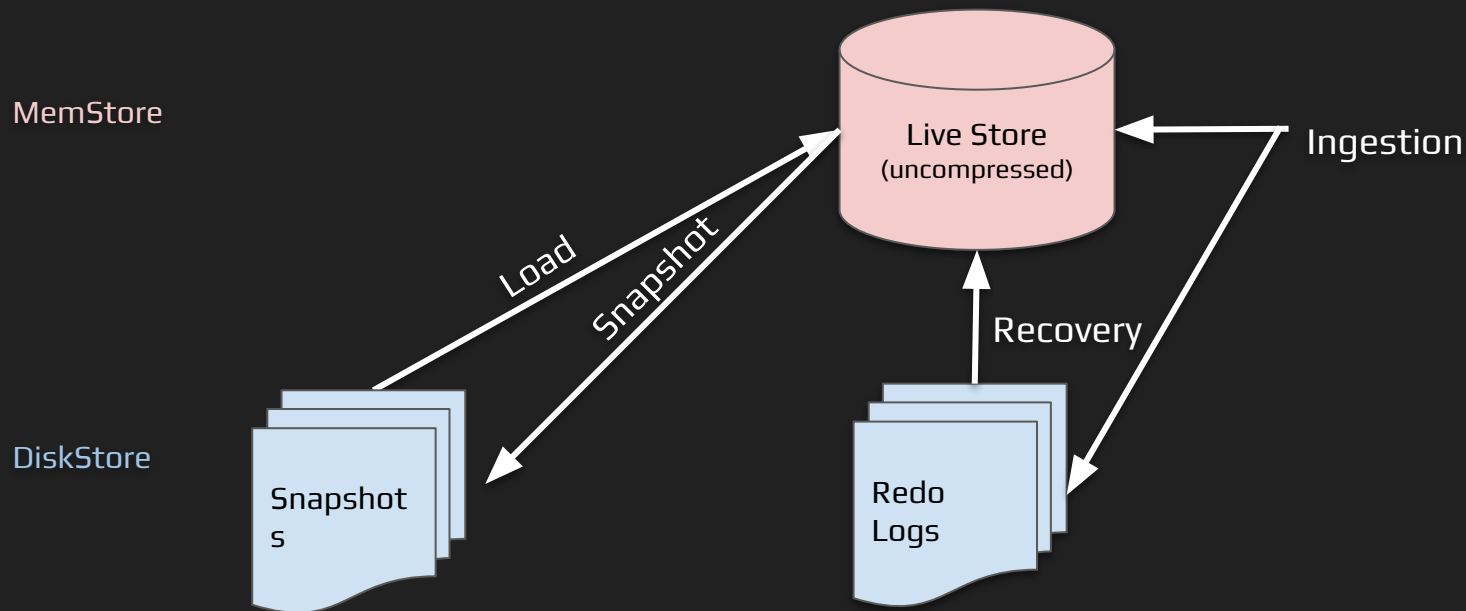
# Columnar Storage: fact table

## Fact Table



# Columnar Storage: dimension table

## Dimension Table

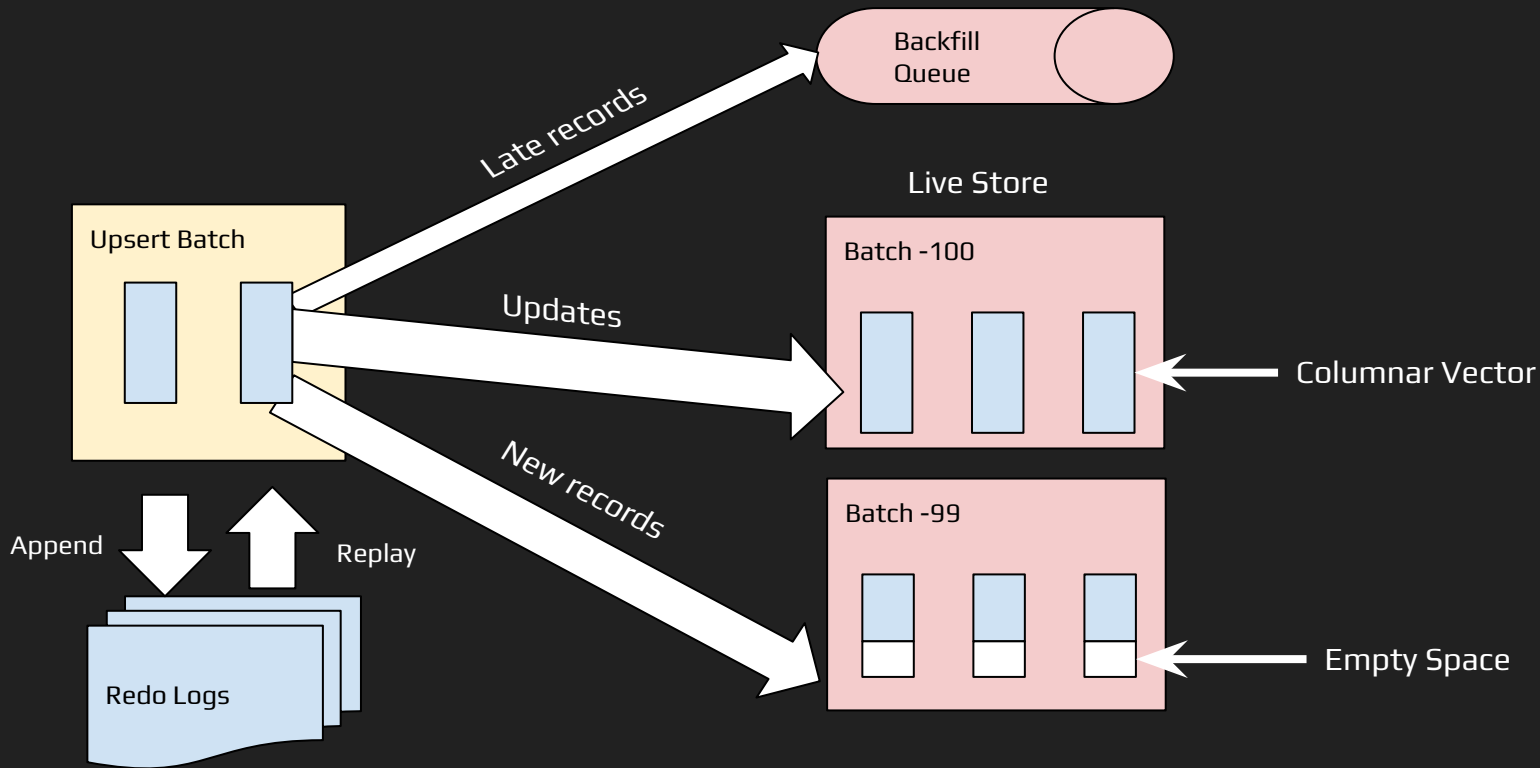




# Feature Highlights

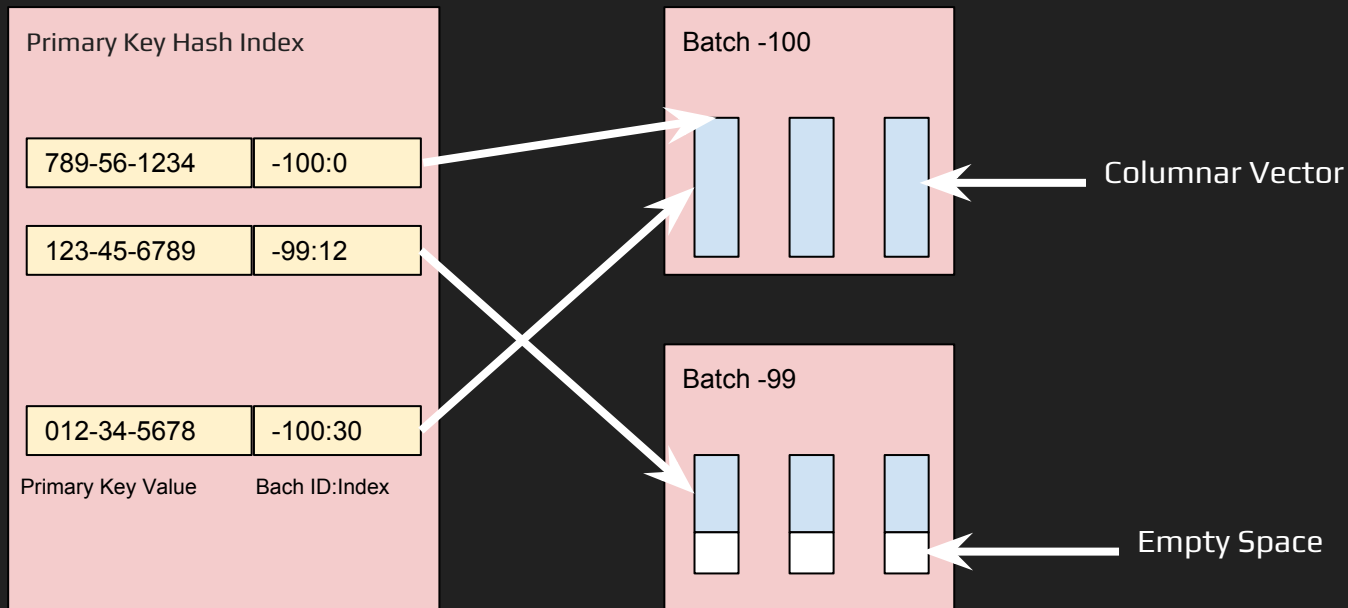
- In-Memory Columnar Storage
- Real-time upserts
- GPU powered query engine
- Analytical Query Feature Set
  - Time zone, Time Filter, Time Bucketization
  - Geospatial analytics
  - Fact/dimension table joins
  - Hyperloglog

# Real-time upserts: ingestion flow



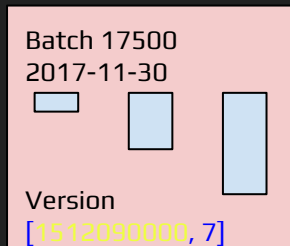
# Real-time upserts: deduplication

Live Store



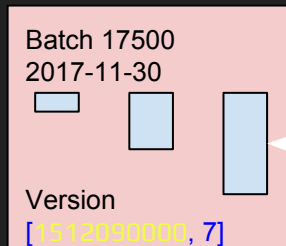
# Real-time upserts: archiving

Archive Store Version  
1512090000

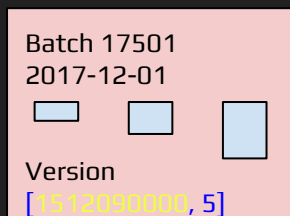


Reuse Unaffected  
Batches

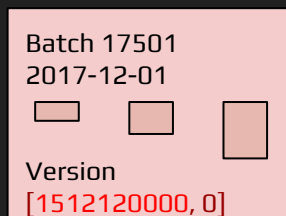
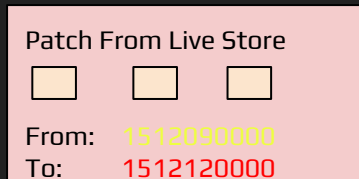
Archive Store Version  
1512120000



Sorted And Compressed  
Columnar Vector



Archive Merge



# Real-time upserts: event timeline

Ingestion



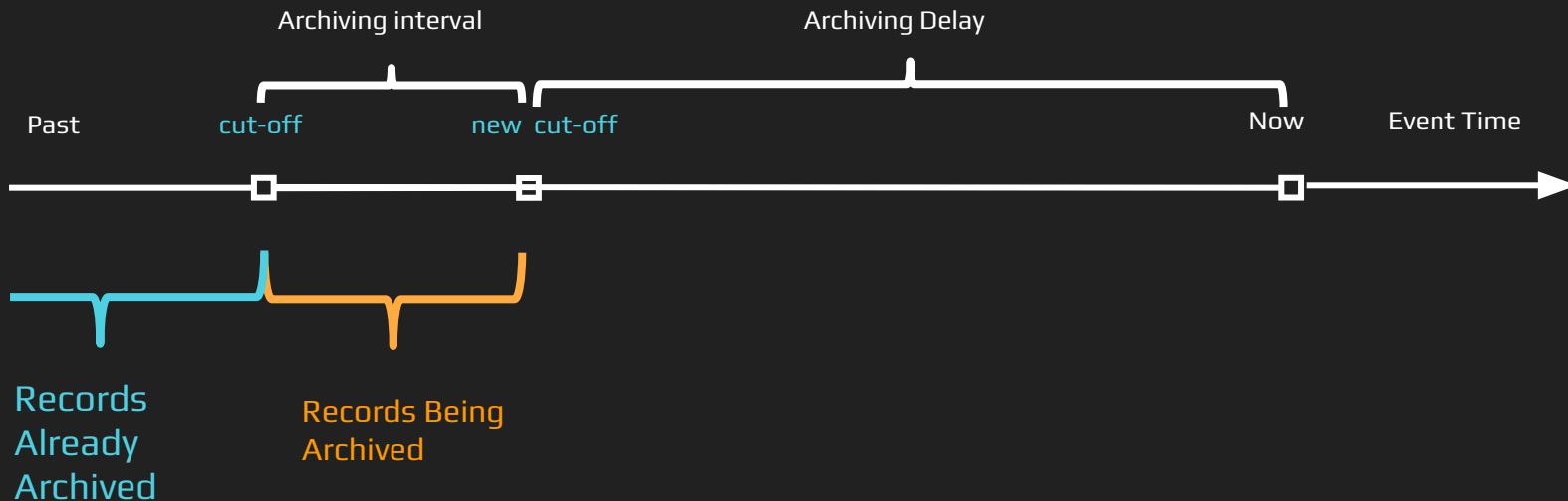
# Real-time upserts: event timeline

Query



# Real-time upserts: event timeline

## Archiving



# Real-time upserts: event timeline

Archiving



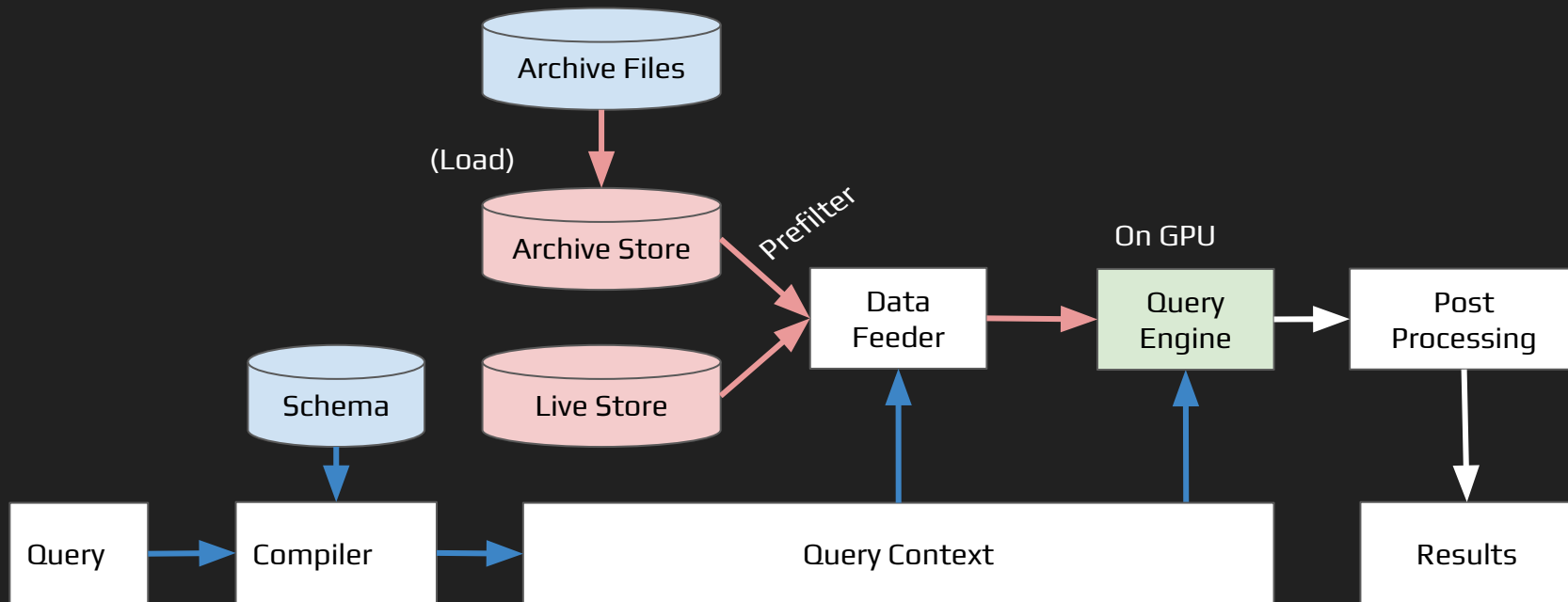


# Feature Highlights

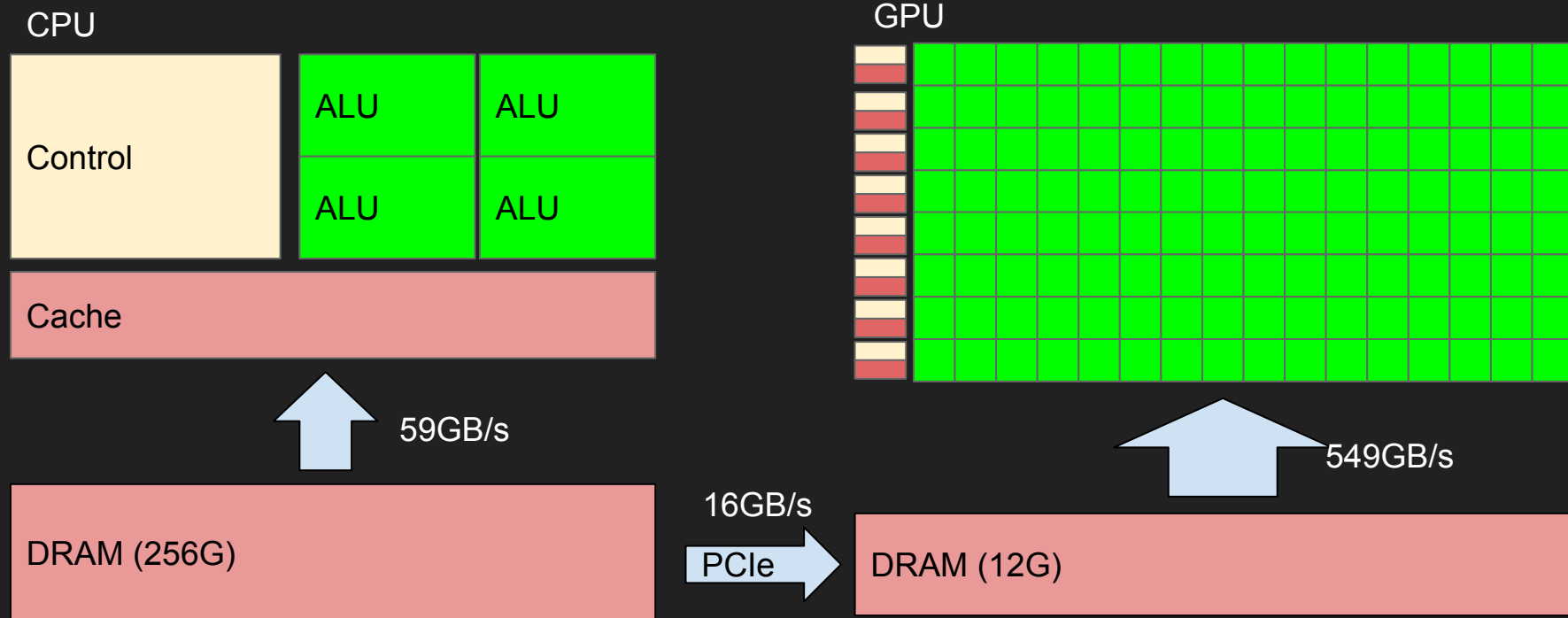
- In-Memory Columnar Storage
- Real-time upserts with deduplication
- GPU powered query engine
- Analytical Query Feature Set
  - Time zone, Time Filter, Time Bucketization
  - Geospatial analytics
  - Fact/dimension table joins
  - Hyperloglog

# Query Engine

High level architecture

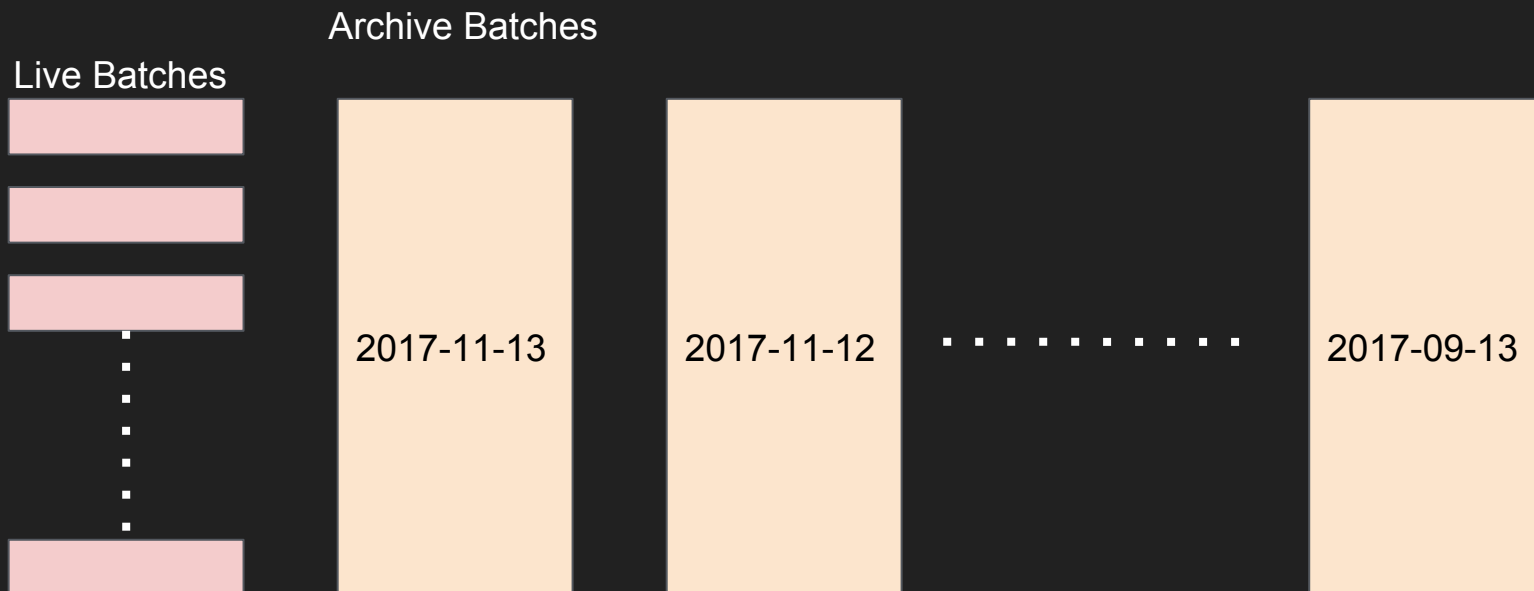


# Data Feeding



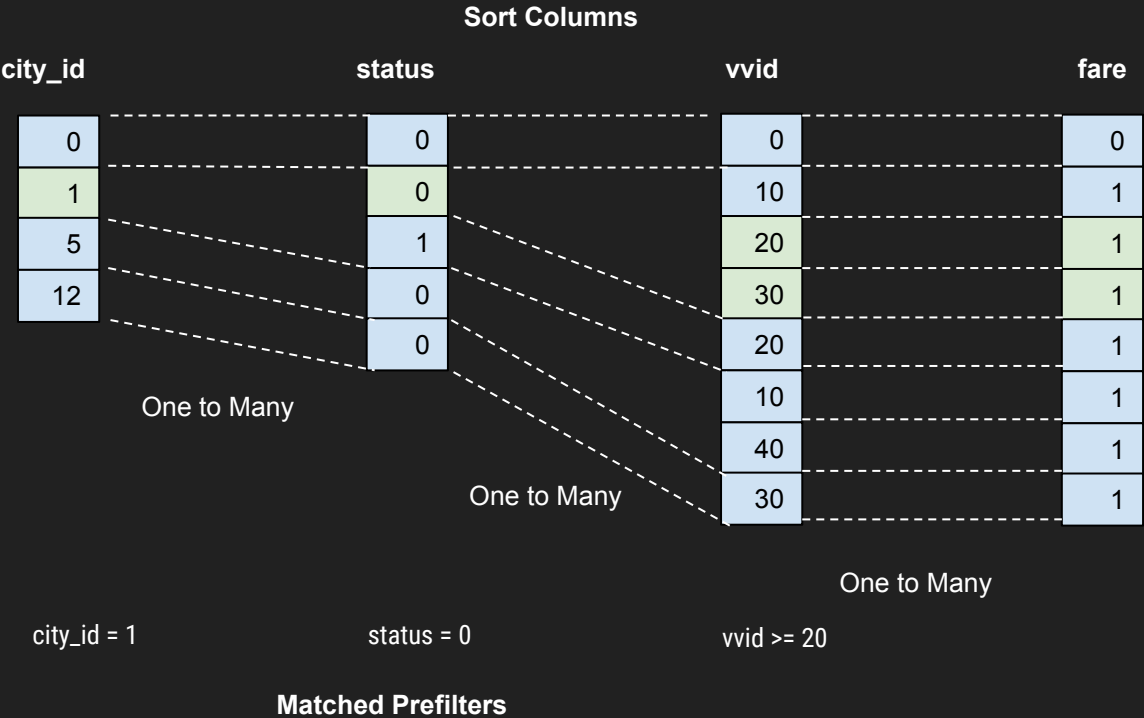
# Data Feeding

## Partitioned Data



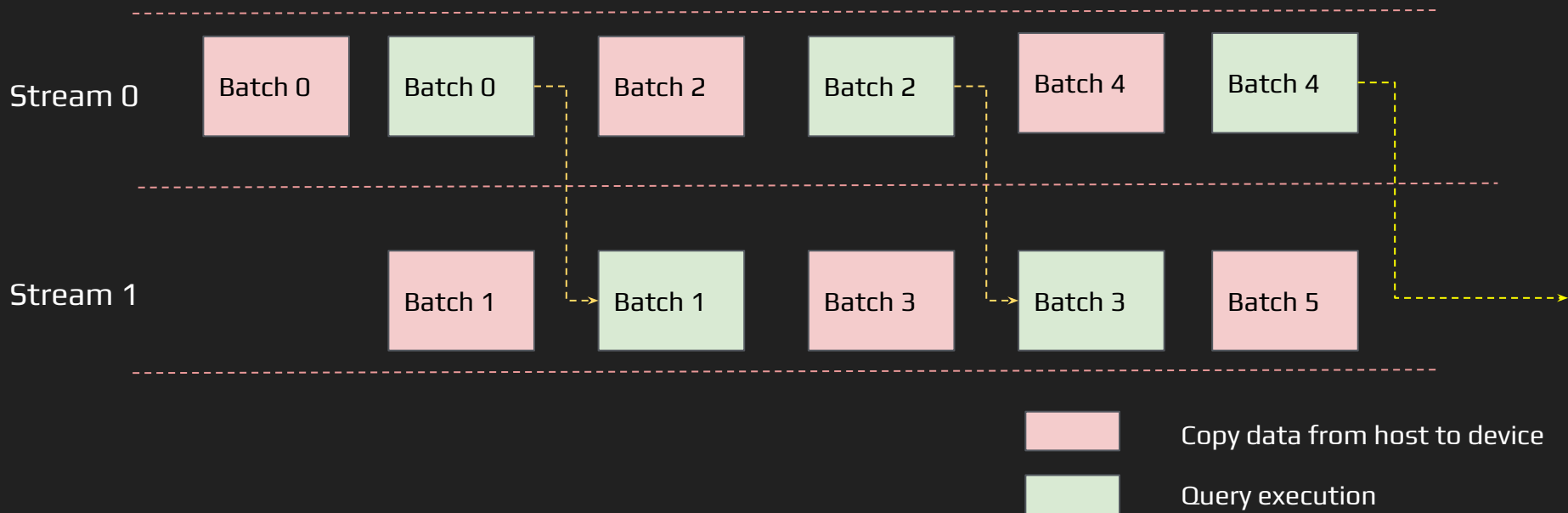
# Data Feeding

## Prefilter



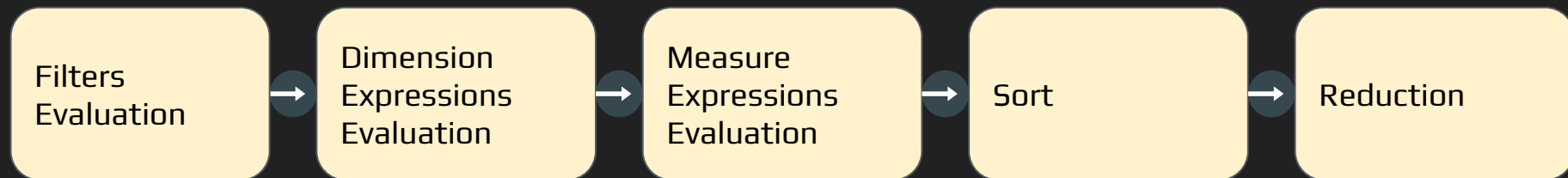
# Data Feeding

## Pipelining



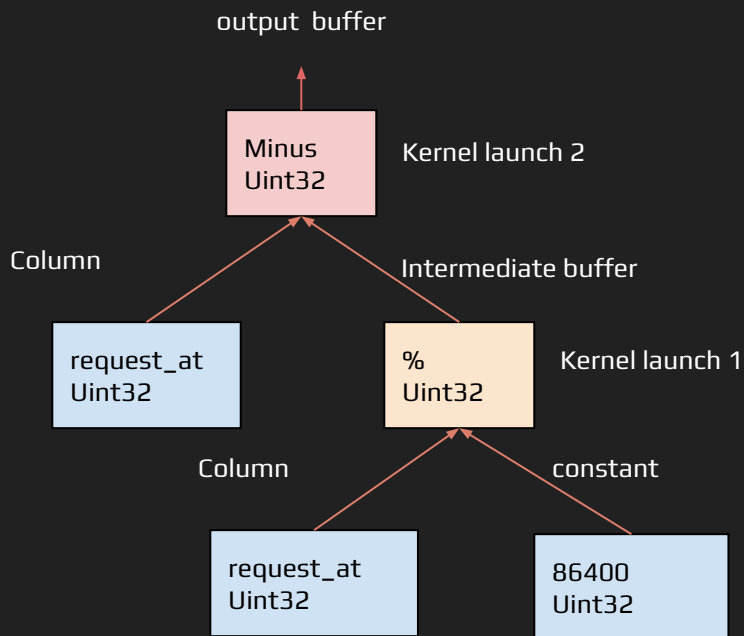
# Query Execution

## Execution Stages



# Query Execution

## Expression Evaluation

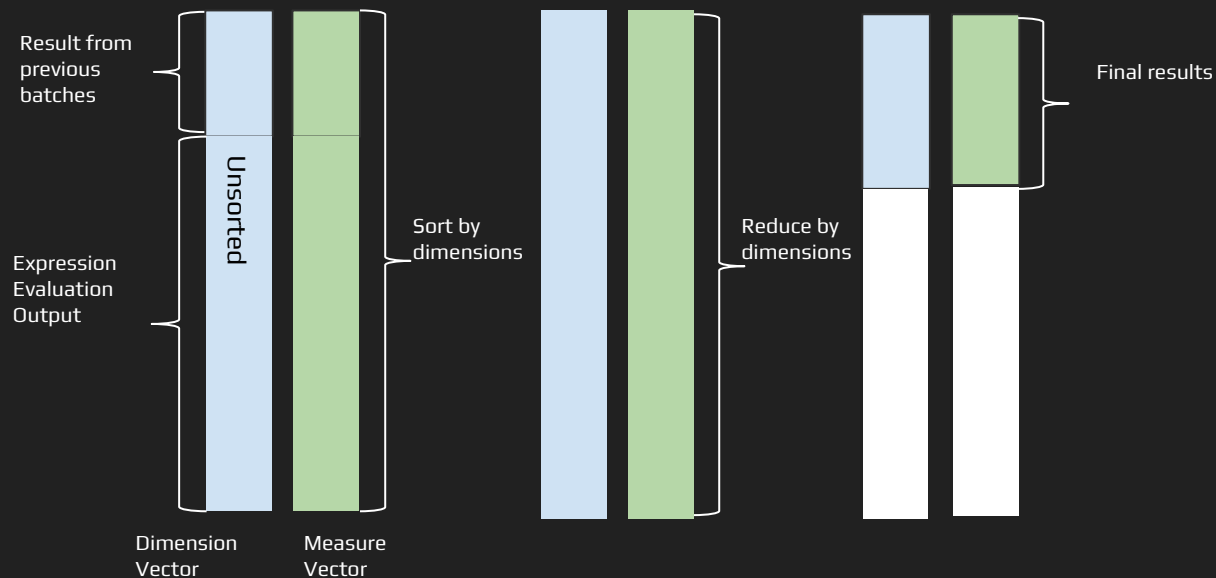


- One operator per kernel on non-leaf nodes
- Each leaf node is one of
  - column/constant
- Non-root, non-leaf node
  - kernel launch
  - output to intermediate buffer
- Root node
  - Kernel launch
  - Write to output buffer
- E.g., request\_at - request\_at % 86400

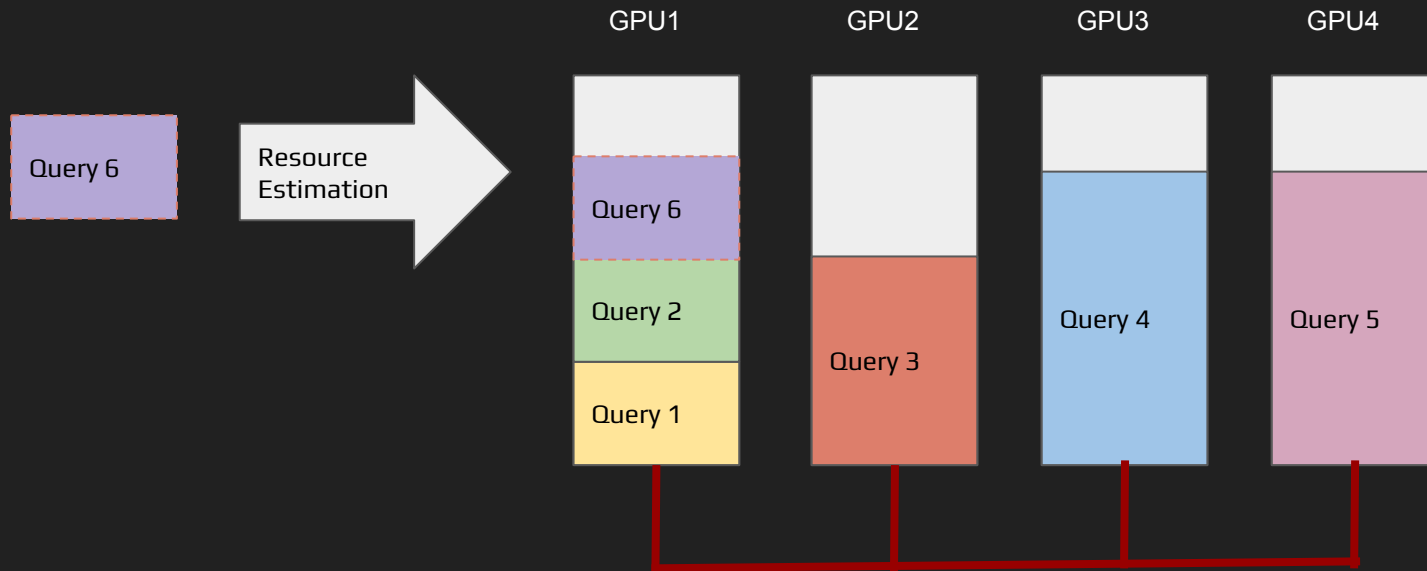


# Query Execution

## Aggregation (Sort and Reduction)



# Device Resource Management



# Feature Highlights

- In-Memory Columnar Storage
- Real-time upserts with deduplication
- GPU powered query engine
- Feature set for analytical queries
  - Time zone, Time Filter, Time Bucketization
  - Fact/dimension table joins
  - Geospatial analytics
  - Hyperloglog

# Timezone, Time Filter, Time Bucketization

```
SELECT    Count(*),
Unix_timestamp(Convert_tz(Concat(Date_format(Convert_tz(From_unixtime(((driver_info.first_active_at) -
(driver_info.first_active_at) % 900000) /
1000), 'GMT', 'America/Los_Angeles' ), '%Y-%m-%d
%H:'),
Lpad(15*Floor(Minute(Convert_tz(From_unixtime(((driver_info.first_active_at) - (driver_info.first_active_at)
%900000)/1000), 'GMT', 'America/Los_Angeles' ))/15),
2, '0')), 'America/Los_Angeles', 'UTC')) AS
time_dimension, driver_info.flow_type
FROM driver_info
WHERE driver_info.first_active_at >=
1534810500000
AND    driver_info.first_active_at <
1534813200000
GROUP BY 2,3
```

V.S.

```
{
  "table": "driver_info",
  "measures": [
    {
      "sqlExpression": "count(*)"
    }
  ],
  "dimensions": [
    {
      "alias": "ts",
      "sqlExpression": "first_active_at",
      "timeBucketizer": "day"
    },
    {
      "sqlExpression": "flow_type"
    }
  ],
  "timeFilter": {
    "column": "first_active_at",
    "from": "7 days ago"
  },
  "timezone": "America/Los_Angeles"
}
```

# Analytical Query Features

- Fact/Dimension Table Join
  - E.g. `trips.city_id = cities.id`
- Hyperloglog Cardinality Estimation
  - `countDistinctHLL(driver_id)`
  - Dedicated hll column
- Geospatial analytics
  - `GeoPoint`, `GeoShape`
  - `GeoIntersect(point, shape)`

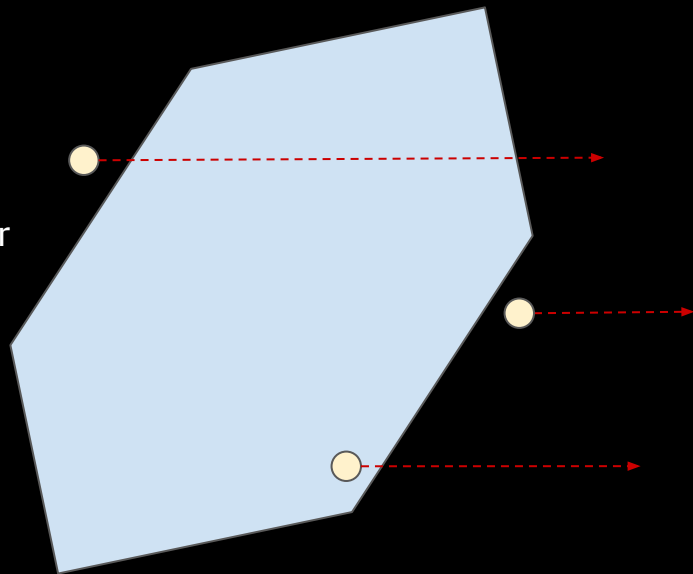
# Learnings from GPU Programming

# Learnings from GPU Programming

- Maximize parallelism
- Optimize memory access
- Maximize arithmetic intensity
- Reduce data transfer between GPU/CPU
- Profile, profile, profile

# Maximize Parallelism

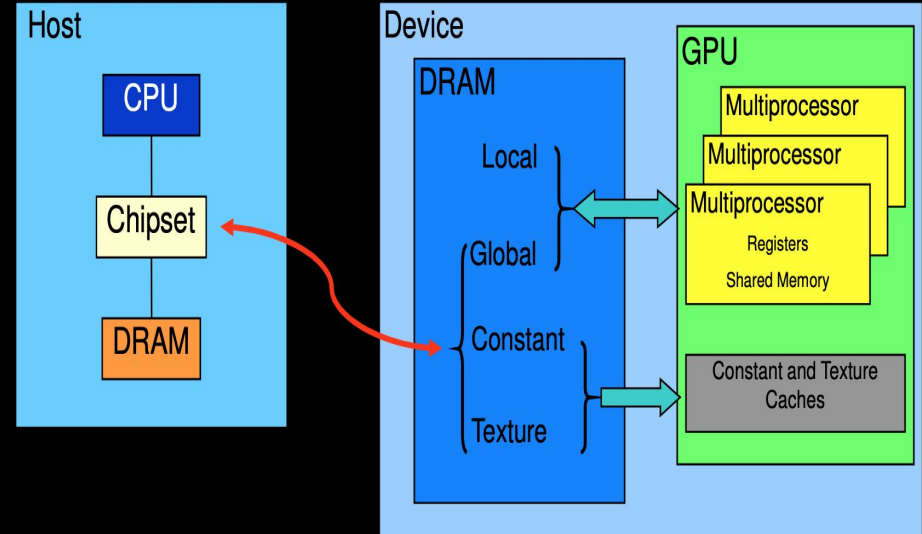
- Partition your computation to keep the GPU multiprocessors equally busy
  - Many threads, many thread blocks
  - E.g. Inclusion test for a point and polygons
    - One shape per thread vs one edge per thread
- Keep resource usage low enough to maximize occupancy
  - Register, shared memory
  - Careful design of data structure
    - Use less wide data type
      - Int64 -> uint32
    - Reuse memory space
      - Union
    - Passing offsets instead of pointers





# Optimize Memory Access

- Coalesced vs. non-coalesced = order of magnitude
- Global/Local device memory
- Shared memory
- Constant memory



[http://developer.download.nvidia.com/CUDA/training/NVIDIA\\_GPU\\_Computing\\_Webinars\\_CUDA\\_Optimization\\_April-2009.pdf](http://developer.download.nvidia.com/CUDA/training/NVIDIA_GPU_Computing_Webinars_CUDA_Optimization_April-2009.pdf)  
retrived date: 01/10/2019

## Maximize Arithmetic Intensity

- GPU spends its transistors on ALUs, not memory
- Sometimes it's better to recompute than to cache
- Do more computation on GPU instead of transferring back to CPU

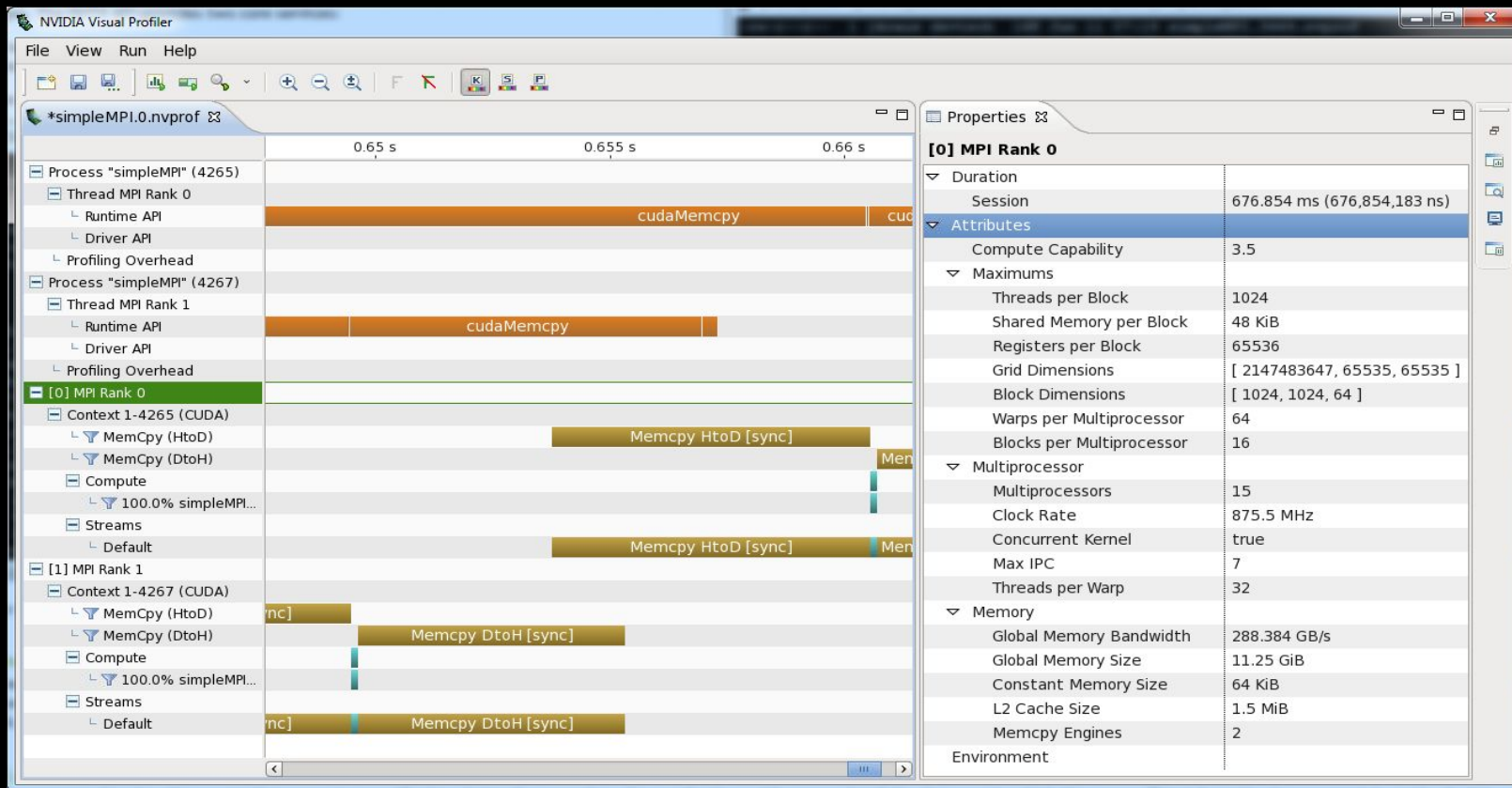
# Minimize CPU/GPU Transfers

- Group transfers
- Overlapping data transfers and computation
  - Async and stream api
  - Stream = sequence of operations that execute in order on GPU
  - Pipeline execution
- Pinned memory vs. pageable memory

# Profiling GPU Program NVVP

- Nvidia visual profiler
- Unified CPU/GPU timeline
- Automated performance analysis
- Guided application analysis

# NVVP cont'd



# Future Directions

- Beyond single instance
  - Sharding
  - Replication
- Ease of adoption
  - SQL interface
  - Native Kafka support
- More query features (eg. fact to fact table joins)
- Query engine optimizations (eg. GPU memory caching)
- Grow AresDB together with the community

# Questions ?

Tech blog: <https://eng.uber.com/aresdb/>

Git repo: <https://github.com/uber/aresdb>

Questions: email [uberopen@uber.com](mailto:uberopen@uber.com)

Follow our Facebook page:

[www.facebook.com/uberopensource](https://www.facebook.com/uberopensource)

# Thank you

Questions: email [uberopen@uber.com](mailto:uberopen@uber.com)

**Follow our Facebook page:**  
**[www.facebook.com/uberopensource](https://www.facebook.com/uberopensource)**

Proprietary © 2018 Uber Technologies, Inc. All rights reserved. No part of this document may be reproduced or utilized in any form or by any means, electronic or mechanical, including photocopying, recording, or by any information storage or retrieval systems, without permission in writing from Uber. This document is intended only for the use of the individual or entity to whom it is addressed. All recipients of this document are notified that the information contained herein includes proprietary information of Uber, and recipient may not make use of, disseminate, or in any way disclose this document or any of the enclosed information to any person other than employees of addressee to the extent necessary for consultations with authorized personnel of Uber.

