

# X-Device Query Processing by Bitwise Distribution

Holger Pirk

Stefan Manegold

Thibault Sellam

Martin Kersten

CWI, Amsterdam, Netherlands

# The Problem

# The Data



# The Data

- A moving objects database of vehicles in NW Europe
  - 240M Tuples (lat int, long int, trip bigint, timestamp bigint)
- Applications like traffic monitoring, forecasting or planning

# The Queries



# The Queries

- Lot's of spatial range queries:

select \* where x between  $x_{low}$  and  $x_{high}$  and

y between  $y_{low}$  and  $y_{high}$

- Focus on Query Throughput

# Let's kill it with Hardware

- GPUs have high bandwidth and are massively parallel
  - Many queries can be evaluated concurrently
- GPU memory is never large enough

# What can you do?

- Host-to-Device Streaming suffers from PCI bottleneck
  - Horizontal Distribution leaves most data/costs on the CPU
  - Sampling only yields approximate results
  - Decomposed Storage is not enough
- **The existing techniques just don't suffice**

# Approach

Remember gifs through a modem connection?



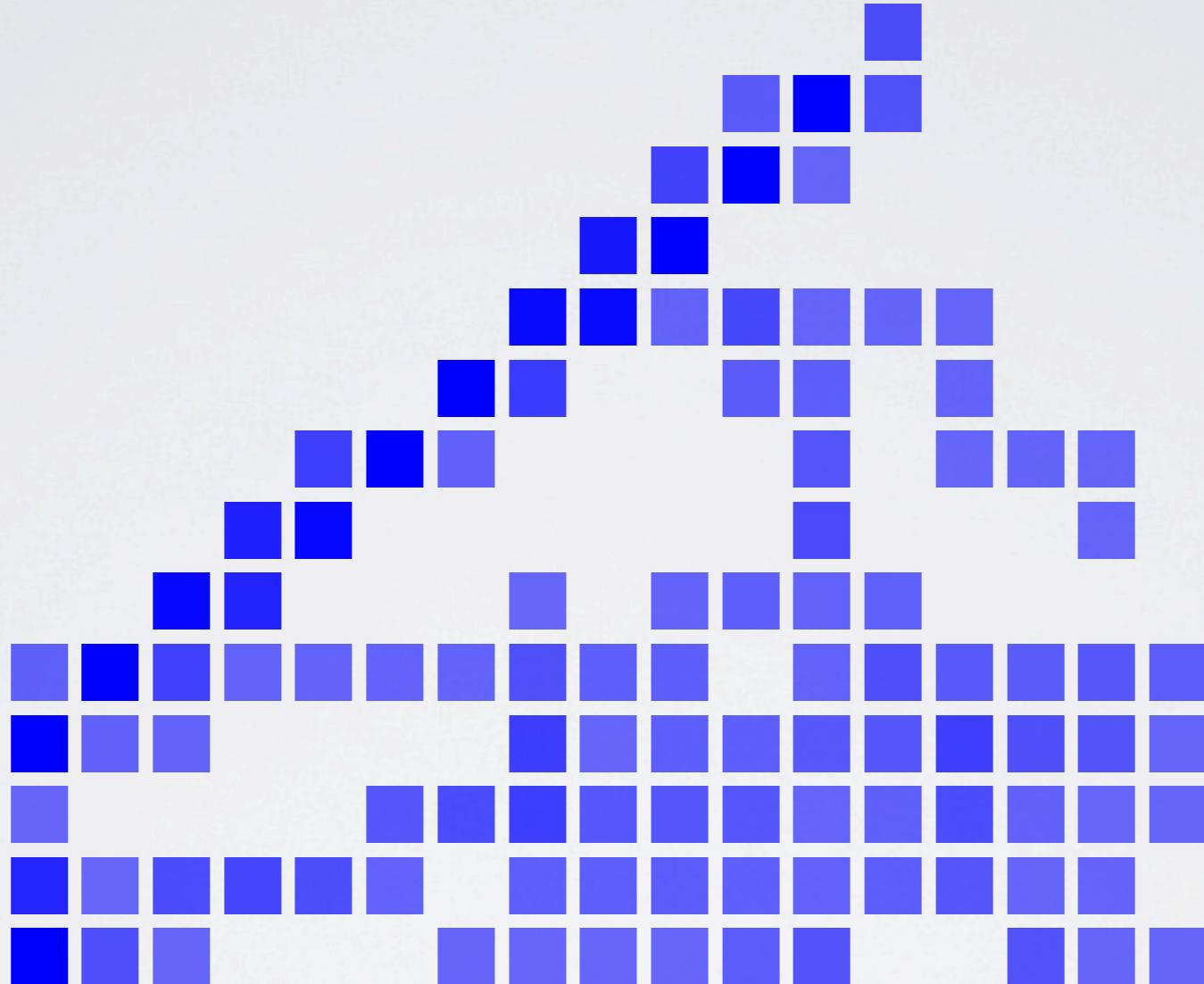
# The Approach



# The Approach



# Approximate

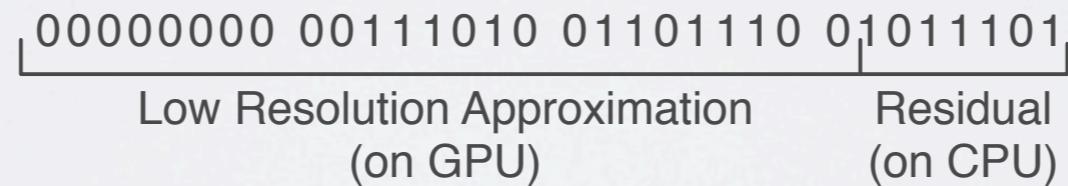


# Refine



# Implementation

- Vertical Partitioning and Distribution
  - But with finer slices: Bitwise Decomposition & Distribution

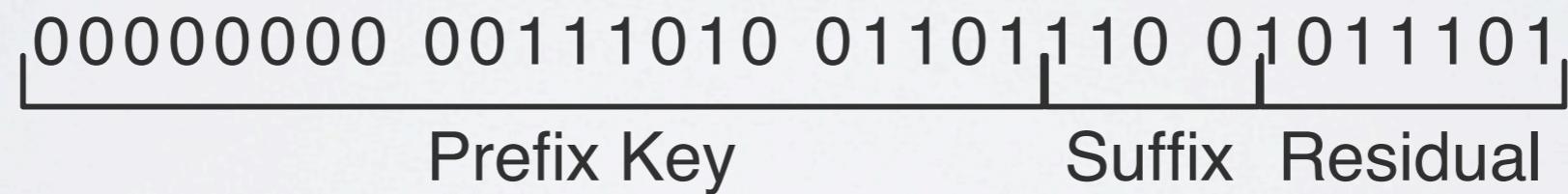


- Precision on GPU is selected appropriately (see paper)
- Low Resolution Data is well compressible

# Compression

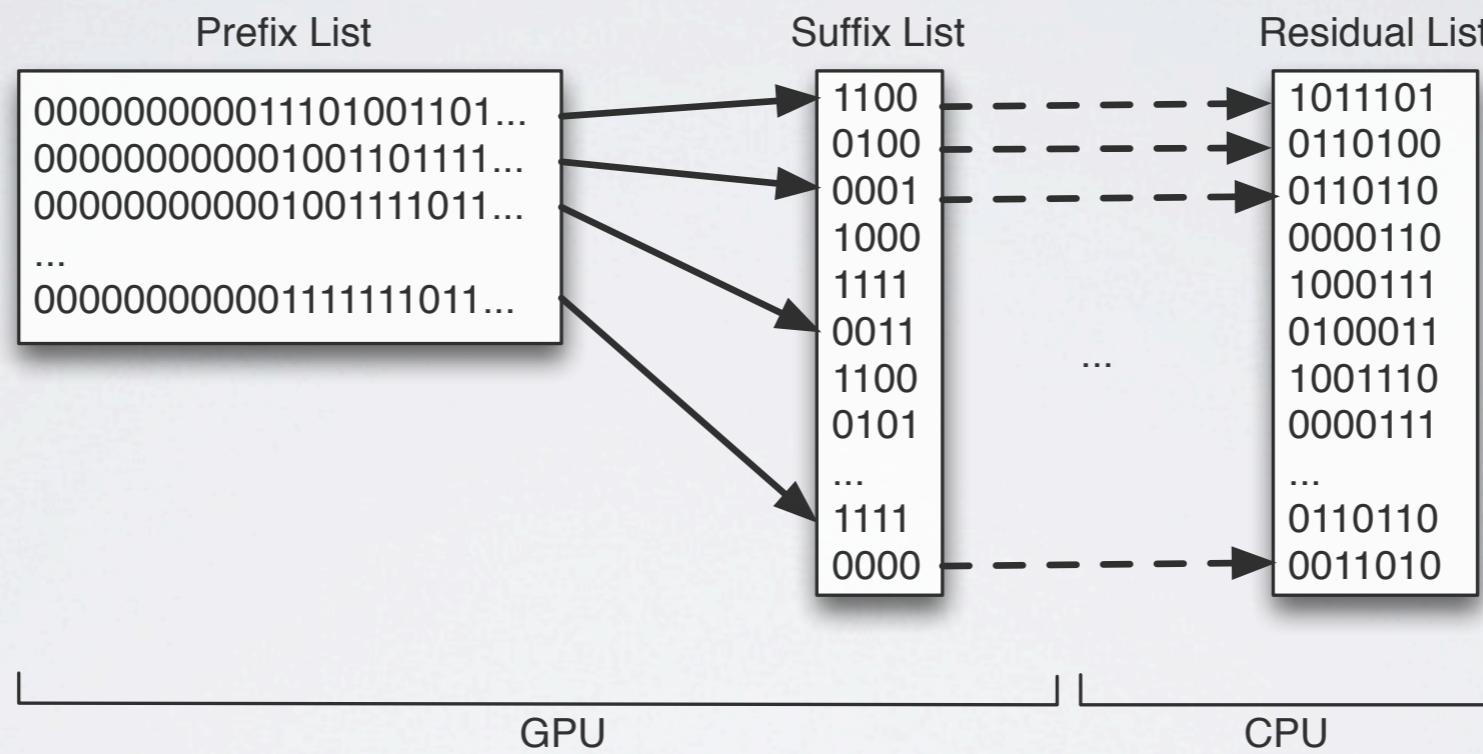
- Has to be Lightweight and parallel decompression friendly:

- Prefix Compression
- Values are split into:



- Spatial Locality helps compression

# Compression



- Radix-Clustering clustering improves compression further

# Query Evaluation

- Queries are evaluated in phases, tuples reconstructed lazily
- Every device does it's best with the available data

Prefix 00000000 00111010 01101000 00000000

+ Suffix 110 0

= Partial Tuple 00000000 00111010 01101110 00000000 GPU

# GPU-Phase

- Geared towards throughput
- Parallelized over Queries and Prefixes
- Prefix-Clusters are skipped
  - This may cause thread divergence
  - Output-Regions are overallocated

# Query Evaluation

- Queries are evaluated in phases, tuples reconstructed lazily
- Every device does it's best with the available data

Prefix 00000000 00111010 01101000 00000000

+ Suffix 110 0

= Partial Tuple 00000000 00111010 01101110 00000000 GPU

# Query Evaluation

- Queries are evaluated in phases, tuples reconstructed lazily
  - Every device does its best with the available data

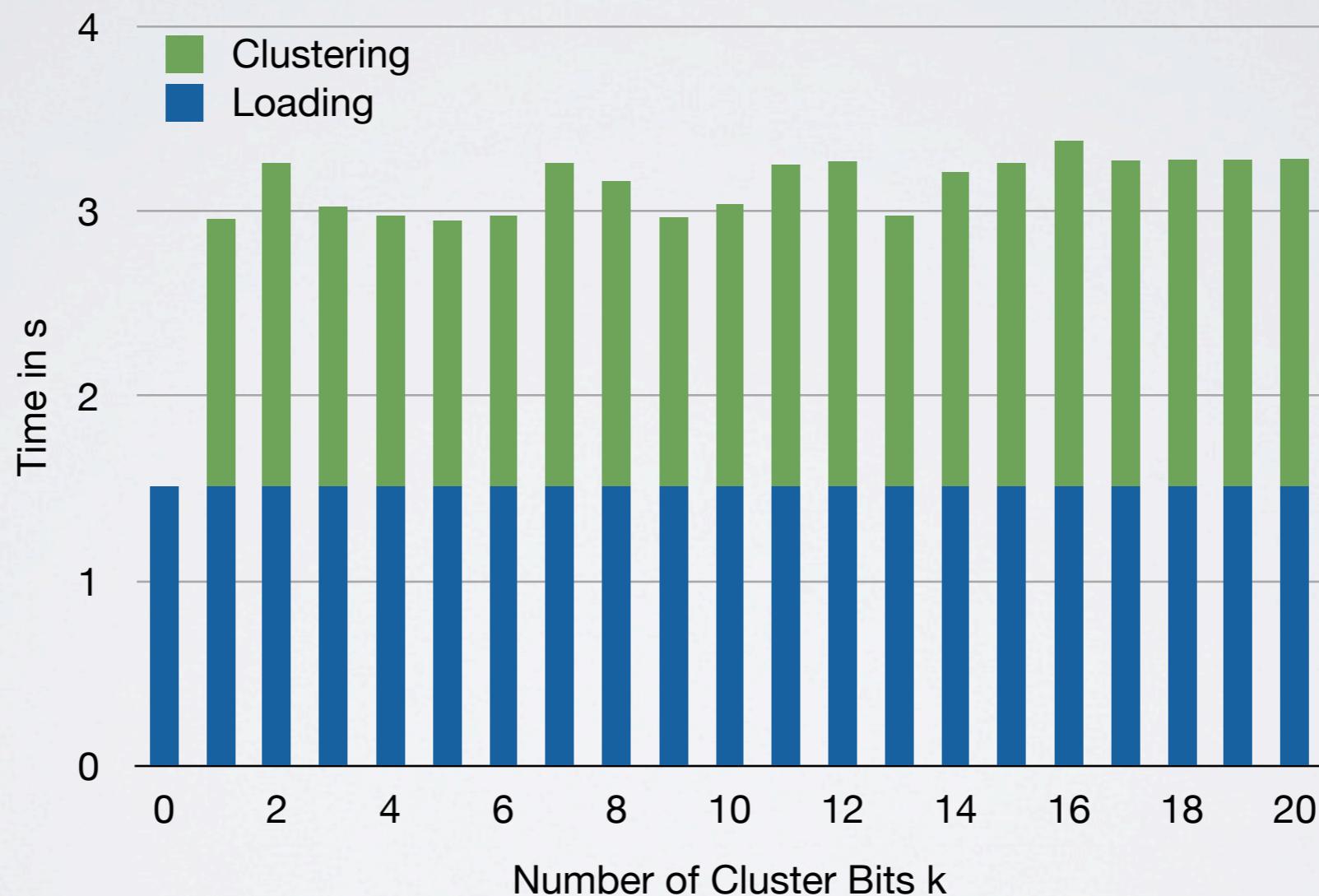
Prefix	<span style="border: 1px solid black; padding: 2px;">00000000 00111010 01101</span>	000 0000000
+ Suffix		<span style="border: 1px solid black; padding: 2px;">110 0</span>
= Partial Tuple	00000000 00111010 01101110 00000000	GPU
	-----	-----
Partial Tuple	00000000 00111010 01101110 00000000	CPU
+ Residual		<span style="border: 1px solid black; padding: 2px;">1011101</span>
= Full Tuple	00000000 00111010 01101110 01011101	

# CPU-Phase

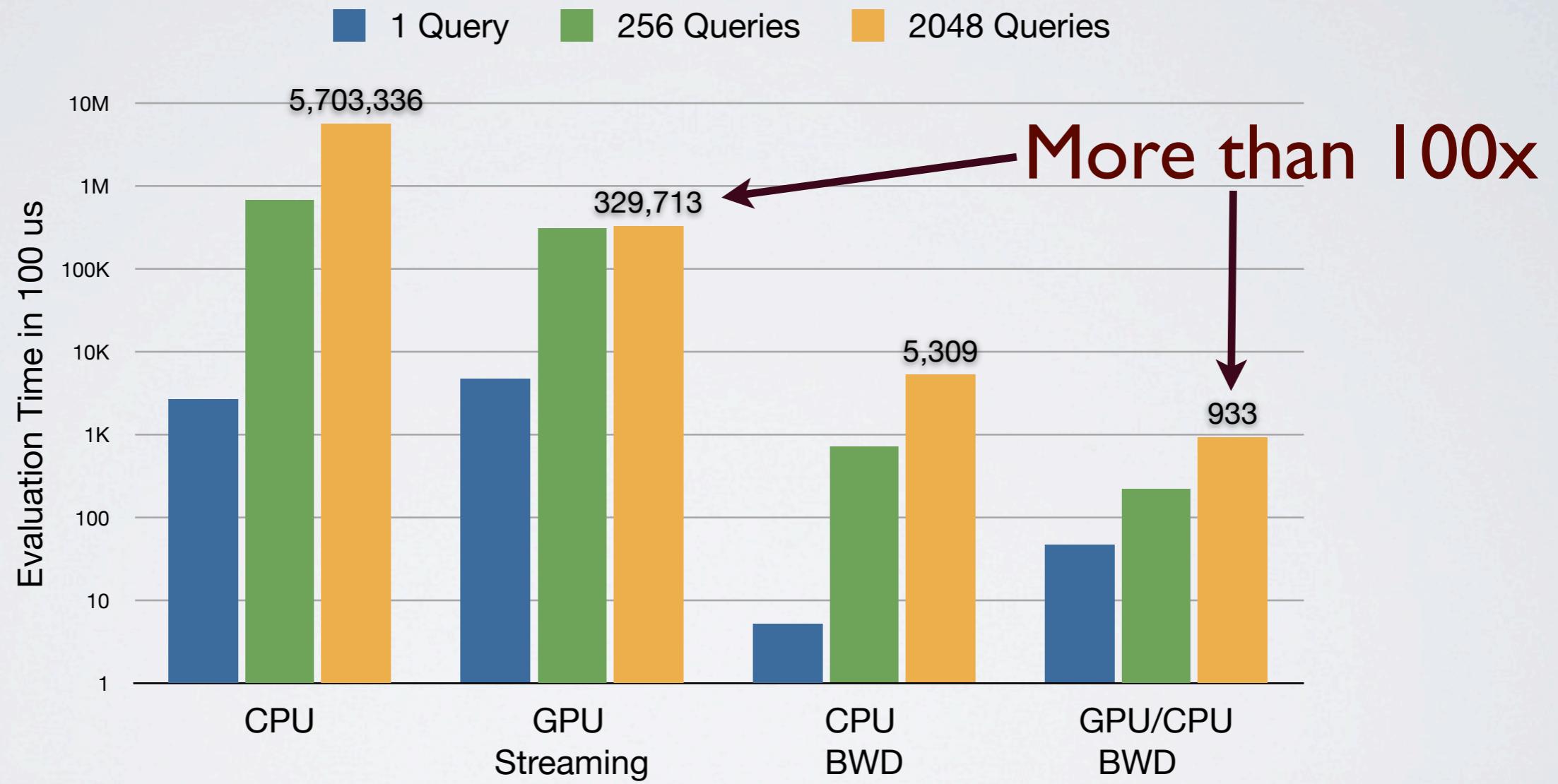
- Approximation may include false positives
- Post-filtering is done on the CPU
- **But:** they have to be transferred to the CPU
- We found around 25% false positives in the first phase

# Evaluation

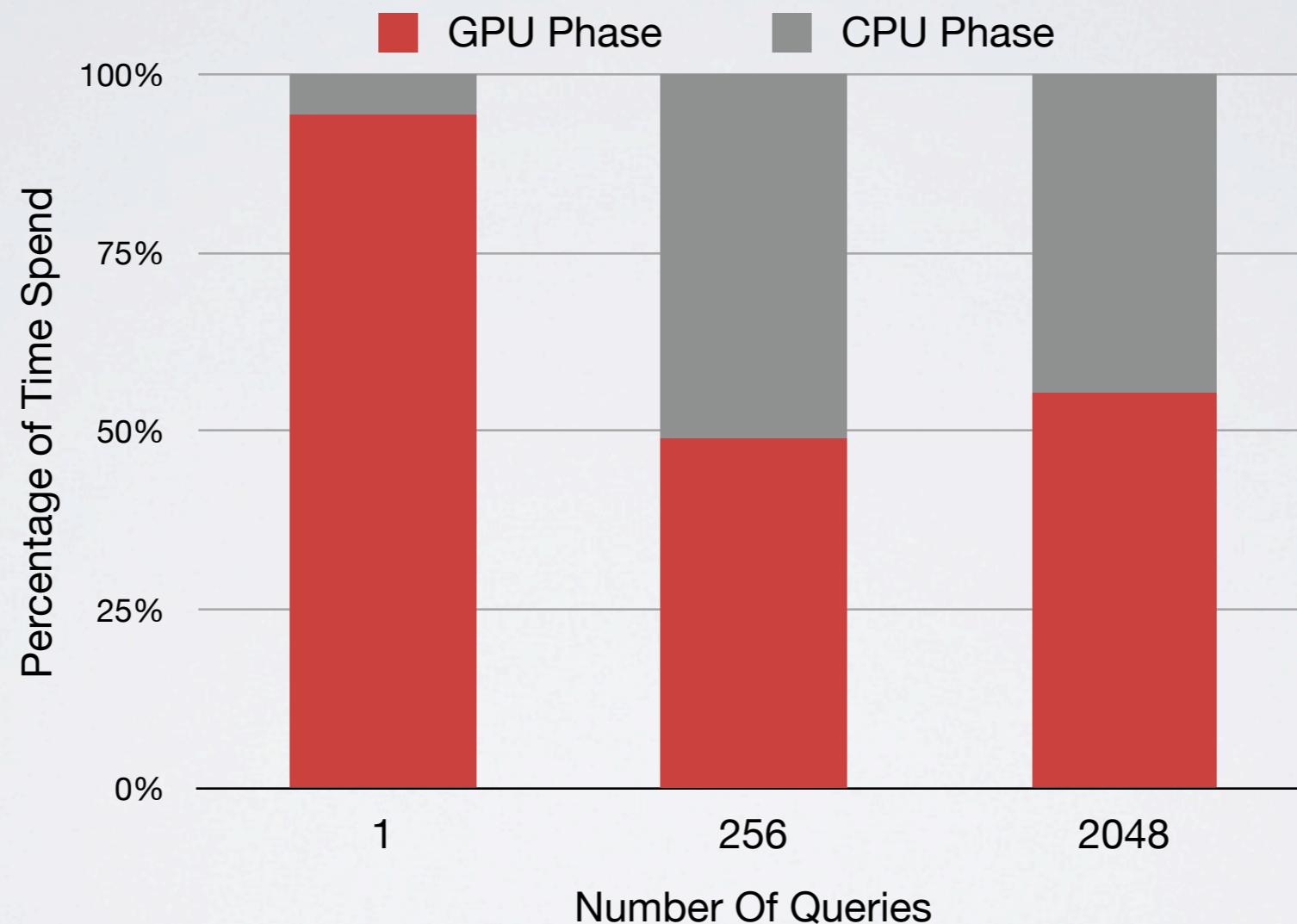
# Loading Performance



# Query Performance



# Load Balancing



# Conclusions and Outlook

# Conclusion

- **Bitwise Decomposition** offers high query throughput on devices with different characteristics like GPUs and CPUs
- Largely targeted at scans
  - Doesn't really compete with indices in latency

# Outlook

- Other X-Device Setups:
  - SSD-HDD (*remember the Keynote?*)
  - Mobile - Server
  - GPU - CPU - SSD - HDD - Server
- Other Applications:
  - Relational Queries (Joins, Aggregations, ...)
- Other Data: Categorical, Float, ...

# Thank you !

Visit us at the poster session on Wednesday 2pm

# Questions?

Visit us at the poster session on Wednesday 2pm