

for GPU-Accelerated Queries on Out-of-Memory Datasets



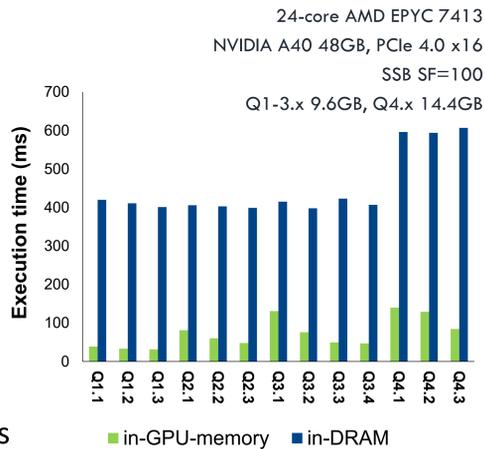
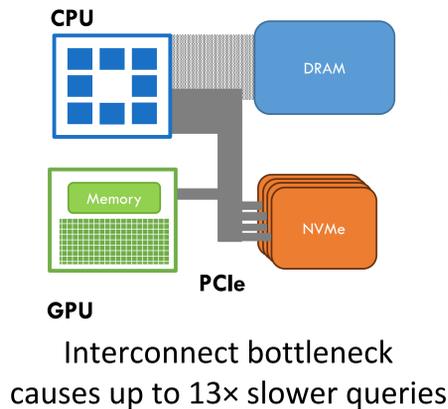
Hamish Nicholson, Konstantinos Chasialis, Antonio Boffa, Anastasia Ailamaki

hamish.nicholson@epfl.ch, konstantinos.chasialis@oracle.com, antonio.boffa@epfl.ch, anastasia.ailamaki@epfl.ch



Queries on out-of-memory datasets are limited by interconnect

Compression as a Solution to the Interconnect Bottleneck



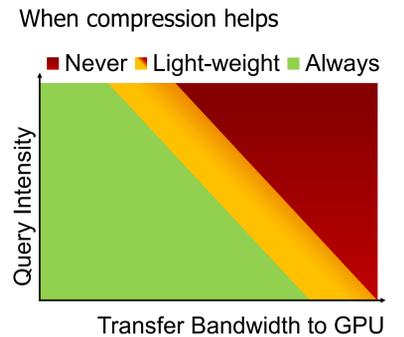
Compression reduces data movement but brings new performance trade-offs

Recent research:

1) in-memory engines introduce new lightweight encoders

2) out-of-memory systems rely on heavyweight compressors

Miss the full picture: whether compression helps depends jointly on bandwidth, query intensity, and selectivity

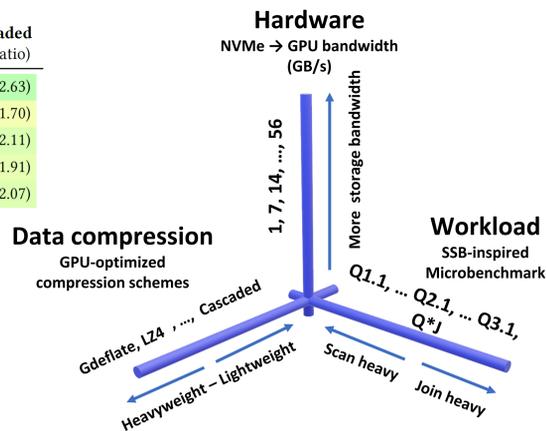


Prior research overlooks the conditional nature of compression benefits in GPU-accelerated systems

The effectiveness depends on more than just compression ratios

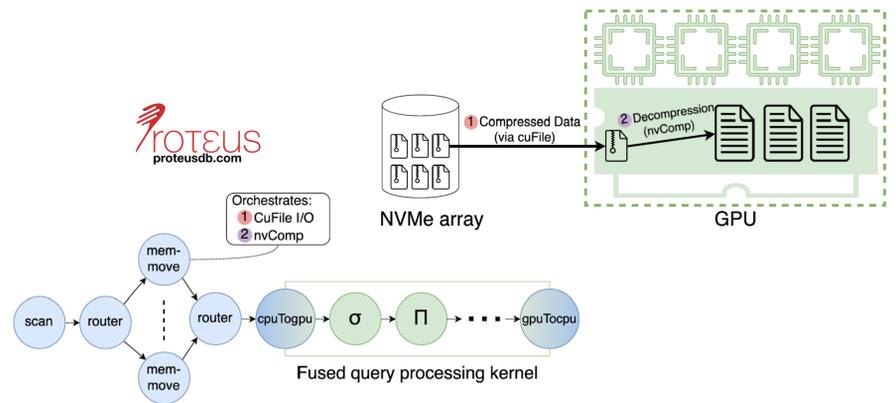
On-the-fly decompression inside mem-move operator

Query	Raw GB	LZ4 GB (ratio)	GDeflate GB (ratio)	Cascaded GB (ratio)
Q1.1	96	65.7 (1.46)	49.0 (1.96)	36.5 (2.63)
Q2.1	96	79.9 (1.20)	75.0 (1.28)	56.5 (1.70)
Q3.1	96	65.5 (1.47)	77.0 (1.25)	45.4 (2.11)
Q4.1	144	110.8 (1.30)	112.7 (1.28)	75.4 (1.91)
Q'J	144	107.3 (1.34)	103.6 (1.39)	69.4 (2.07)



Key factors that govern compression effectiveness?

Analysis of compression benefit on out-of-memory datasets along 3 axis



Extend GPU query engine to pipeline decompression, data movement, and execution

Results on Start Schema Benchmark

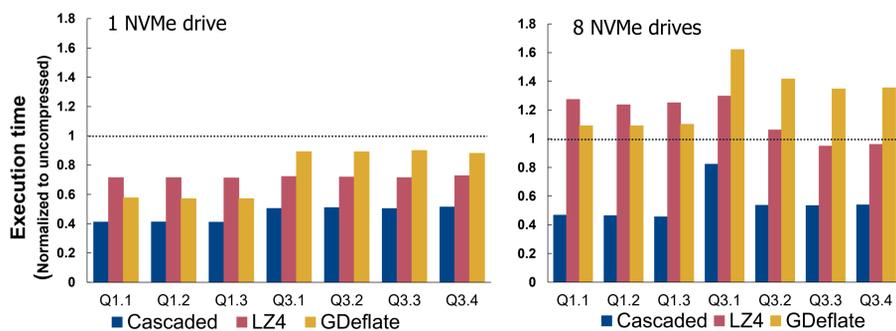
Roofline Analysis. Are we still interconnect bound?

Speed up ≈ Compression ratio

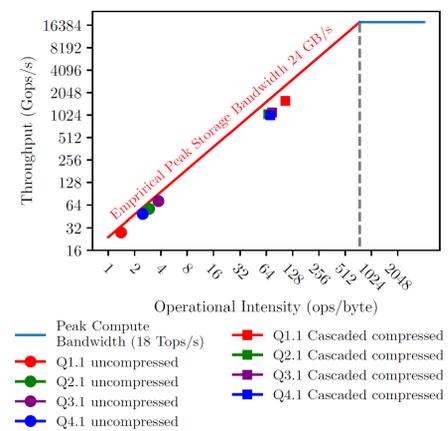
From 1-3 ops/byte To 67-104 ops/byte

28x more instructions executed

2x speed up



GDeflate and LZ4 degrade performance

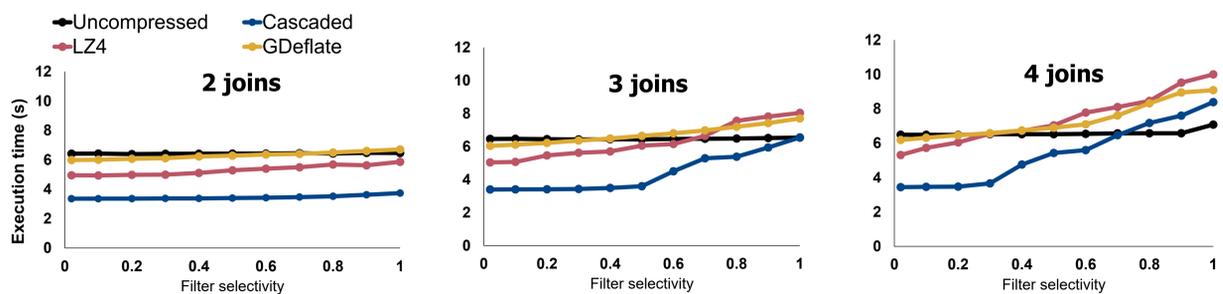


Even with decompression, analytical queries remain far from compute roof Leaving room for more (efficient) compression

When compression introduces a slow down

With high bandwidth, process intensive queries, high selectivity, decompression can stall the GPU

```
SELECT SUM(lo_revenue)
FROM date, customer, supplier, part, lineorder
WHERE lo_quantity < {FILTER_SELECTIVITY_PARAM}
AND lo_orderdate = d_datekey ← Q1J, one join
AND lo_suppkey = s_suppkey ← Q2J, two joins
AND lo_partkey = p_partkey ← Q3J, three joins
AND lo_custkey = c_custkey ← Q4J, four joins
```



Increasing query processing work per input tuple