

# CS 764: Topics in Database Management Systems Lecture 6: Buffer Management

Xiangyao Yu 9/23/2025

# Today's Paper: Buffer Management

#### Robust External Hash Aggregation in the Solid State Age

CWI. Amsterdam. Netherlands laurens.kuiper@cwi.nl

CWI. Amsterdam. Netherlands peter.boncz@cwi.nl

CWI, Amsterdam, Netherlands hannes.muehleisen@cwi.nl

Abstract—Analytical database systems offer high-performance in-memory aggregation. If there are many unique groups, temporary query intermediates may not fit RAM, requiring the use of external storage. However, switching from an in-memory to an external algorithm can degrade performance sharply.

We revisit external hash aggregation on modern hardware, aiming instead for robust performance that avoids a "performance cliff" when memory runs out.

To achieve this, we introduce two techniques for handling temporary query intermediates. First, we propose unifying the memory management of temporary and persistent data. Second, we propose using a page layout that can be spilled to disk despite being optimized for main memory performance. These two techniques allow operator implementations to process largerthan-memory query intermediates with only minor modifications.

We integrate these into DuckDB's parallel hash aggregation. Experimental results show that our implementation gracefully degrades performance as query intermediates exceed the available memory limit, while main memory performance is competitive with other analytical database systems.

Index Terms-relational databases, database query processing, aggregation

#### I. INTRODUCTION

operators. As a result, these systems could process workloads illustrated in Figure 1 that were larger than the small amount of available memory.

systems optimized for main memory [2], for both persistent on modern hardware [6], [8], OLAP systems should be able data and temporary query intermediates [3]. In these systems, to perform more robustly when intermediates exceed main main memory access became the bottleneck, and techniques memory. However, traditional buffer managers have fixed-size were devised to make better use of CPU caches [4]. DBMSes pages and a statically allocated pool. This inflexibility makes have now evolved into large monolithic database servers, often them undesirable for intermediates. Therefore, temporary data with large amounts of RAM at their disposal.

Efficient utilization of secondary storage, e.g., by caching, is help systems better utilize available memory [13]. without sacrificing in-memory performance [9].

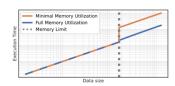


Fig. 1. Conceptual aggregation performance vs data size (log-log scale). When switching from an in-memory strategy to an external strategy that minimizes memory usage, performance degradation is harsh and sudden (a "performance cliff"). A unified strategy for in-memory and external aggregation that utilizes all available memory degrades more gracefully (performance

This body of research has focused on using storage for persistent data but, for the most part, ignored temporary query intermediates. Analytical (OLAP) systems, which frequently process large volumes of data and often have large query intermediates, became mainstream after DBMSes optimized for main memory. As systems became able to process queries Until late in the 20th century, main memory was expensive; on arbitrary-sized persistent tables, intermediate results can therefore, traditional database management systems (DBMS) - depending on the query - also grow to arbitrary sizes. In optimized for disk access, as this was their major bottleneck. these cases, many modern OLAP systems either abort queries "Spillable" data structures like B-trees [I] were used not only or switch to a traditional disk-based algorithm that is orders to speed up retrieval of persistent data but also inside query of magnitude slower, introducing a "performance cliff", as

Given the advancements of OLAP systems in the past two Around the 2000s, RAM prices decreased, and database decades [10]-[12], and the research into buffer management is allocated differently. Managing the entire memory pool, i.e., Pure in-memory systems are not economical, however. persistent and temporary data, in a cooperative manner may

key to providing good performance at a low cost [5]. In recent In this paper, we go beyond Cooperative Memory Manyears, there has been a renewed interest in buffer management, agement and take a unified approach to memory managespecifically for solid-state memory, that offers much higher ment for persistent and temporary data. We have developed bandwidth and lower latency than magnetic disk [6]-[8]. Data a specialized page layout specifically for temporary data to management systems are now reverting to being disk-based accommodate this. We have integrated this into the hash aggregation operator of DuckDB, our in-process analytical Algorithmica (1986) 1: 311-336



#### An Evaluation of Buffer Management Strategies for Relational Database Systems<sup>1</sup>

Hong-Tai Chou<sup>2,3</sup> and David J. DeWitt<sup>2</sup>

Abstract. In this paper we present a new algorithm, DBMIN, for managing the buffer pool of a relational database management system. DBMIN is based on a new model of relational query behavior, the query locality set model (QLSM). Like the hot set model, the QLSM has an advantage over the stochastic models due to its ability to predict future reference behavior. However, the QLSM avoids the potential problems of the hot set model by separating the modeling of reference behavior from any particular buffer management algorithm. After introducing the QLSM and describing the DBMIN algorithm, we present a performance evaluation methodology for evaluating buffer management algorithms in a multiuser environment. This methodology employed a hybrid model that combines features of both trace-driven and distribution-driven simulation models. Using this model, the performance of the DBMIN algorithm in a multiuser environment is compared with that of the hot set algorithm and four more traditional buffer replacement algorithms.

Key Words. Buffer management, Database systems, Page replacement strategies, Hybrid simulation, Performance evaluation

1. Introduction. In this paper we present a new algorithm, DBMIN, for managing the buffer pool of a relational database management system. DBMIN is based on a new model of relational query behavior, the query locality set model (OLSM.) Like the hot set model [Sacc 1], the QLSM has an advantage over stochastic models due to its ability to predict future reference behavior. However, the QLSM avoids the potential problems of the hot set model by separating the modeling of reference behavior from any particular buffer management algorithm. After introducing the QLSM and describing the DBMIN algorithm, the performance of the DBMIN algorithm in a multiuser environment is compared with that of the hot set algorithm and four more traditional buffer replacement algorithms.

A number of factors motivated this research. First, although Stonebraker [Ston 2] convincingly argued that conventional virtual memory page replacement algorithms (e.g., least recently used (LRU)) were generally not suitable for a

Received March 15, 1986; revised July 7, 1986. Communicated by Dale Skeen.

Parts of this article have been reprinted with permission by the "Very Large Data Base Endowment."

<sup>&</sup>lt;sup>1</sup> This research was partially supported by the Department of Energy under Contract No. DE-AC02-81ER10920 and the National Science Foundation under grant MCS82-01870.

<sup>&</sup>lt;sup>2</sup> Computer Sciences Department, University of Wisconsin, Madison, Wisconsin, USA.

<sup>&</sup>lt;sup>3</sup> Current Address: Microelectronics and Computer Technology Corporation, Austin, Texas, USA.

### Outline

#### Table data vs. temporary data

Buffer management for table data

Page-based buffer management

Buffer management for temporary data

Unified memory management

Future work

### Table Data vs. Temporary Data

#### **Table data**

- Data from input tables
- Page-based buffer management is widely used for table data
- Transparently handle data movement between disk and memory

#### **Temporary data**

- E.g., hash tables (for join or aggregation), sort buffers, etc.
- Traditionally malloc-based, not in page granularity
- Must explicitly move data between disk and memory

### Outline

Table data vs. temporary data

#### **Buffer management for table data**

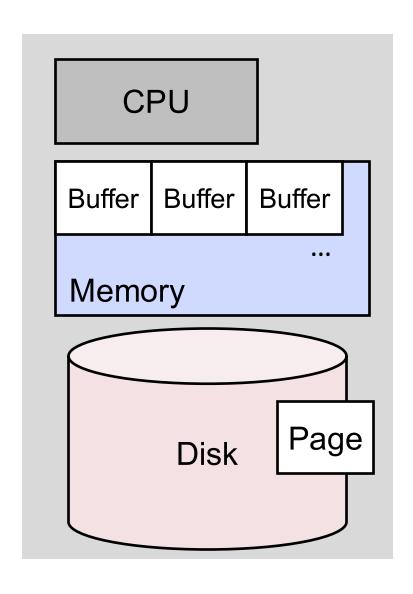
Page-based buffer management

Buffer management for temporary data

Unified memory management

Future work

### System Architecture

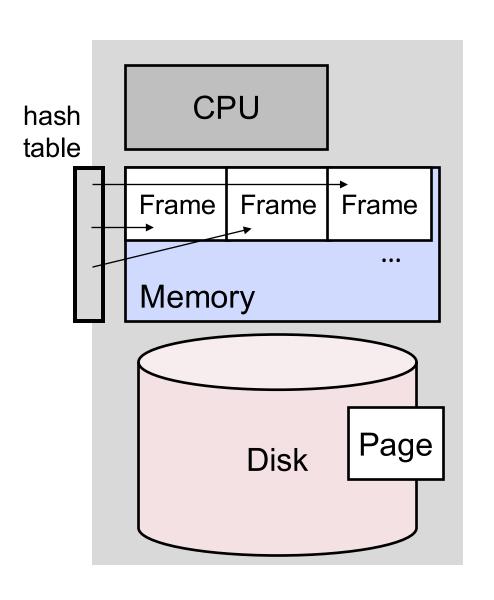


A database management system (DBMS) manipulate data in memory

 Data on disk must be loaded to memory before processed

The unit of data movement is a page

## Page-Based Buffer Management

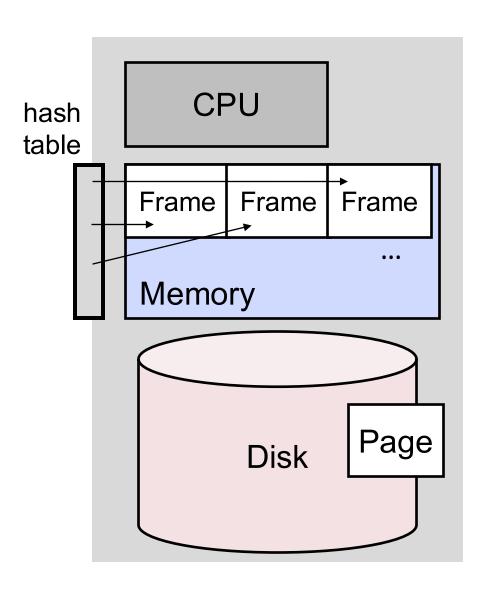


Page granularity: Data managed in page granularity

Indirection: Each page is identified with a page ID; a hash table stores whether the page is in memory or on disk.

Page placement is handled by the buffer manager and not controlled by operators

### Page-Based Buffer Management



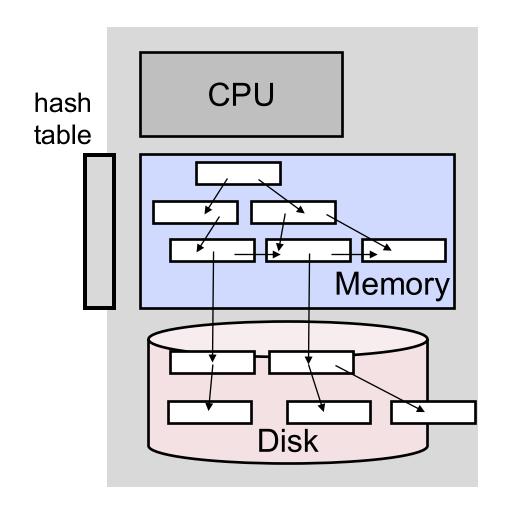


Table data can be organized in a B-tree

### Page Replacement Policy

Important question: what pages should stay in memory?

- LRU (Lease recently used)
- Clock
- MRU (Most recently used)
- FIFO, Random, ...

**Insight**: the optimal buffer replacement and allocation policies depend on the data access pattern, which is relatively easy to predict in a DBMS compared to hardware or OS

### LRU Replacement

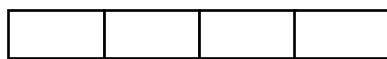
Replace the least-recently used (LRU) item in the buffer

Intuition: more recently used items will more likely to be used again in the future

### LRU Replacement Example

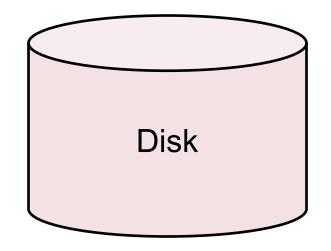
Example: memory contains 4 buffers. LRU replacement policy





Incoming requests

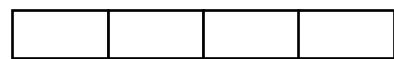
0, 1, 2, 3, 0, 1, 2, 4, 0, 1, 2, 5, ...



### A Different Access Pattern

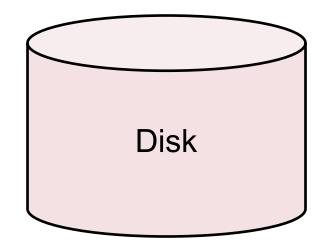
Example: memory contains 4 buffers. LRU replacement policy

Memory



Incoming requests

0, 1, 2, 3, 4, 0, 1, 2, 3, 4, ...



### MRU Replacement

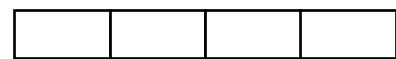
Replace the most-recently used (LRU) item in the buffer

Intuition: avoid the cache thrashing problem in the previous example

### MRU Replacement Example

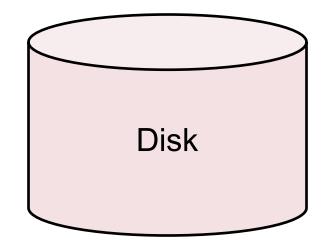
Example: memory contains 4 buffers. MRU replacement policy

Memory



Incoming requests

0, 1, 2, 3, 4, 0, 1, 2, 3, 4, ...



## Query Locality Set Model

#### **Observations**

- DBMS supports a limited set of operations
- Data reference patterns are regular and predictable
- Complex reference patterns can be decomposed into simple patterns

#### Buffer allocation decisions:

- 1. Locality set: The appropriate buffer pool size for a query
- 2. Replacement policy

### QLSM – Sequential References

Straight sequential (SS): each page in a file accessed only once

- E.g., select on an unordered relation
- Locality set: one page
- Replacement policy: any

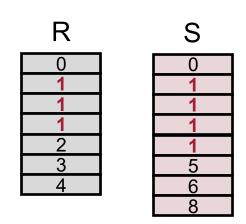
### QLSM – Sequential References

Straight sequential (SS): each page in a file accessed only once

- E.g., select on an unordered relation
- Locality set: one page
- Replacement policy: any

Clustered sequential (CS): repeatedly read a "chunk" sequentially

- E.g., sort-merge join with duplicate join keys
- Locality set: size of largest cluster
- Replacement policy: LRU or FIFO (buffer size ≥ cluster size), MRU (otherwise)



### QLSM – Sequential References

Straight sequential (SS): each page in a file accessed only once

- E.g., select on an unordered relation
- Locality set: one page
- Replacement policy: any

Clustered sequential (CS): repeatedly read a "chunk" sequentially

- E.g., sort-merge join with duplicate join keys
- Locality set: size of largest cluster
- Replacement policy: LRU or FIFO (buffer size ≥ cluster size), MRU (otherwise)

Looping Sequential (LS): repeatedly read something sequentially

- E.g. nested-loop join
- Locality set: size of the file being repeated scanned.
- Replacement policy: MRU

### QLSM – Random References

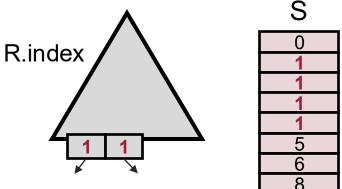
#### Independent random (IR): truly random accesses

- E.g., index scan through a non-clustered (e.g., secondary) index
- Locality set: one page or **b** pages (**b** unique pages are accessed in total)
- Replacement: any

### Clustered random (CR): random accesses with some locality

 E.g., join between non-clustered, non-unique index as inner relation and clustered, non-unique outer relation

- Locality set: size of the largest cluster
- Replacement policy :
   LRU or FIFO (buffer size ≥ cluster size)
   MRU (otherwise)



### Outline

Table data vs. temporary data

Buffer management for table data

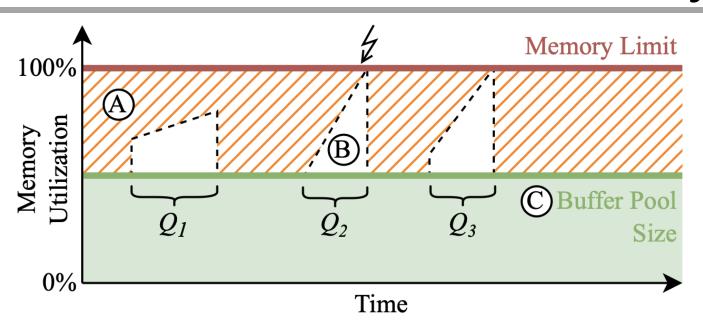
Page-based buffer management

#### **Buffer management for temporary data**

Unified memory management

Future work

### Traditional DB: Two Memory Pools



Sizes of pools are statically determined

### Page-based buffer pool for persistent data

Data spilling is implicitly handled by buffer manager

### Malloc-based buffer pool for temporary data

- Data spilling is not supported, or
- Data spilling is explicitly handled by operator

### Outline

Table data vs. temporary data

Buffer management for table data

Page-based buffer management

Buffer management for temporary data

**Unified memory management** 

Future work

### DuckDB Unified Memory Management

Insight: Try to use page-based data structures as much as possible

- Operators (e.g., group-by aggregation) are designed accordingly

#### **Benefits**:

- Disk spilling is implicitly managed through buffer manager
- Graceful performance degradation when data exceeds memory

#### Challenges:

- Allow arbitrary eviction without corrupting pointers
- Some data structures are non-trivial to store in pages (e.g., hash tables)

## **Allocation Types**

#### Persistent data

Fixed page size (256 KB)

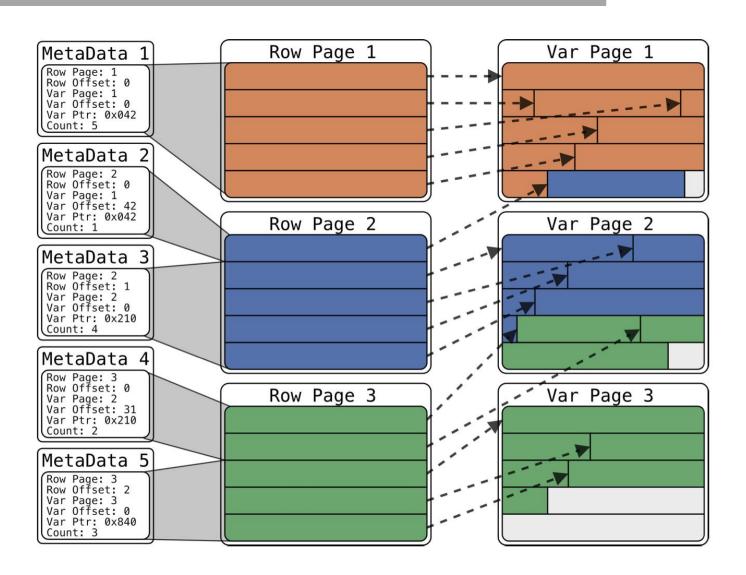
#### Temporary data

- Non-paged allocations (cannot be spilled)
- Paged fixed-size allocations => the most common use case
- Paged variable-size allocations

### Page Layout for Variable-Size Row

#### Design requirements:

- 1) Use a row-major data representation with fixedsize rows
- 2) Store variable-size data on separate pages
- 3) Use explicit addressing for variable-size data
- 4) Be spillable to storage without additional serialization.



# External Hash Aggregation

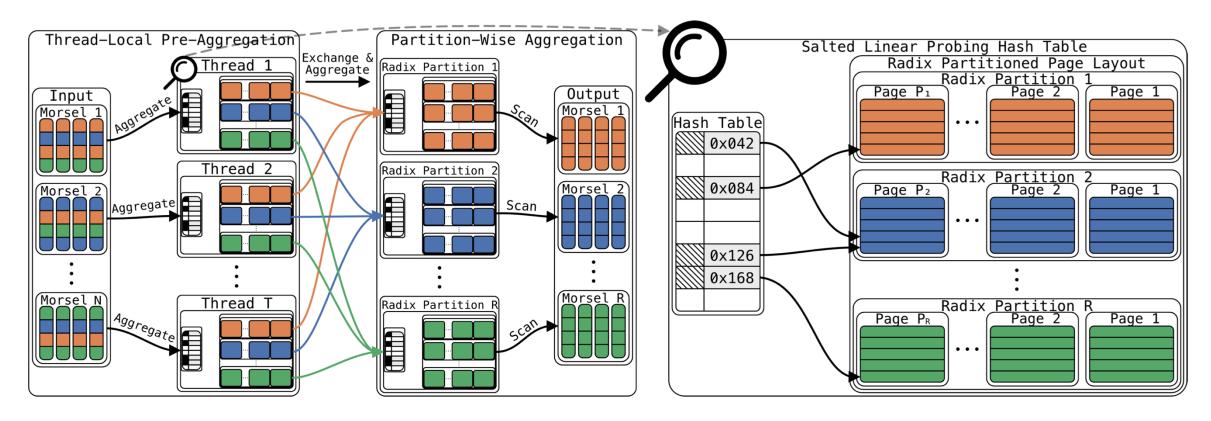


Fig. 3. DuckDB's hash aggregation. Morsels are assigned to threads until all input data has been read. During phase one, each thread pre-aggregates data in a small fixed-size linear probing hash table with one level of indirection, i.e., offsets obtained from hashes access an array of pointers pointing to tuples. Tuples are radix partitioned and stored using DuckDB's spillable page layout, enabling larger-than-memory aggregation. After pre-aggregation, partitions are exchanged and aggregated partition-wise in parallel. Fully aggregated partitions are immediately scanned, effectively becoming morsels in the next pipeline.

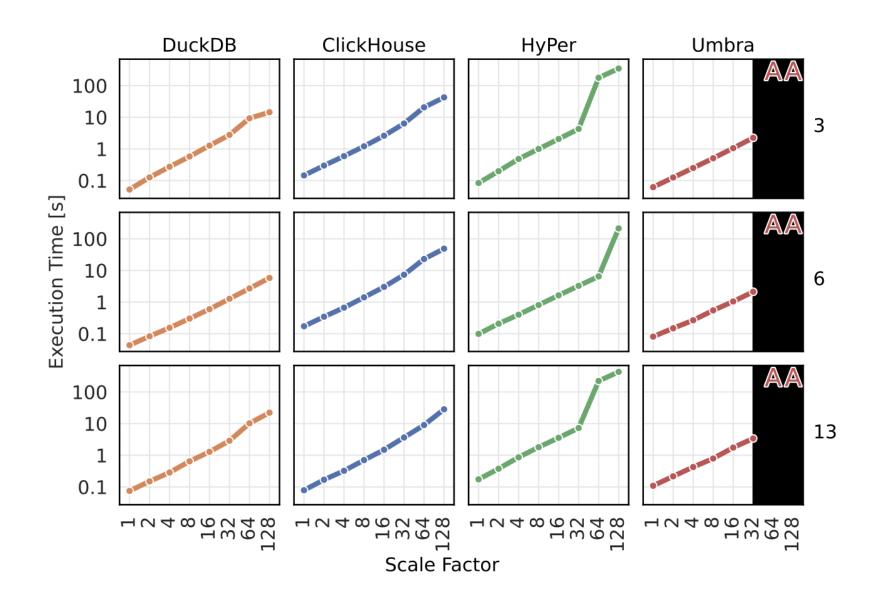
### Evaluation – Thin Variant

TABLE II

EXECUTION TIME IN SECONDS FOR THE THIN VARIANT OF ALL GROUPINGS AT SCALE FACTORS 2, 8, 32, AND 128. LOWER IS BETTER. THE LOWEST EXECUTION TIMES ARE HIGHLIGHTED IN BOLD. 'A' DENOTES THAT THE QUERY WAS ABORTED. WE SUMMARIZE BY NORMALIZING EXECUTION TIMES TO DUCKDB AND THEN TAKING THE GEOMETRIC MEAN.

SF	<b>SF</b>   2		8			}		32		ı		128			
System   D	ı Cl	Ну	Um	Du	C1	Hy	Um	Du	Cl	Ну	Um	Du	C1	Ну	Um
Grouping 1   0.0 2   0.0 3   0.1 4   0.1 5   0.0 6   0.0 7   0.1 8   0.2 9   0.1 10   0.2 11   0.1 12   0.1 13   0.1	8	0.13 0.20 0.16 0.07 0.21 0.28 0.27 0.30 0.41 0.37	0.02 <b>0.04</b> <b>0.13</b> <b>0.07</b> <b>0.05</b> 0.15 0.23 <b>0.20</b> 0.21 0.38 0.18 0.18 0.22	0.03 0.44 0.58 0.58 0.18 0.30 0.72 0.97 0.74 1.10 0.59 0.59 0.59	0.27 <b>0.16</b> 1.21 0.29 0.29 1.42 1.56 1.51 1.53 2.06 <b>0.51</b> <b>0.52</b> 0.70	0.05 0.66 1.00 0.83 0.36 0.81 1.36 1.39 1.58 1.71 1.73 1.80	0.04 0.19 <b>0.51</b> <b>0.28</b> <b>0.16</b> 0.55 0.87 <b>0.78</b> 0.78 1.50 0.65 0.65 0.79	0.14 2.03 2.78 2.86 0.74 1.27 3.42 5.25 3.59 5.78 2.72 2.65 2.87	1.10 0.75 6.35 1.63 1.57 7.32 8.32 8.11 8.01 11.34 2.59 2.59 3.64	0.14 2.86 4.28 3.38 1.66 3.27 5.61 5.72 6.42 169.63 6.94 6.80 7.24	0.16 0.68 2.24 1.22 0.54 2.12 4.32 3.44 3.47 14.98 2.86 2.84 3.39	0.54 10.80 14.49 22.06 3.17 5.80 24.62 65.97 32.77 89.27 21.38 21.89 22.20	4.06 <b>4.04</b> 42.47 <b>9.80</b> 9.41 48.79 51.57 <b>49.49</b> 46.51 <b>77.27</b> <b>18.43</b> <b>18.22</b> 28.31	0.54 12.83 348.10 231.86 6.86 213.41 457.52 412.86 444.50 576.68 413.54 411.58 432.33	A A A A A A A A A A A A A A A A A A A
Geometric Mean Normalized to DuckDB 1.0	<b>0</b> 1.69	1.80	1.13	1.00	1.53	1.98	0.98	1.00	1.69	2.17	0.94	1.00	1.48	8.74	A

### Evaluation – Thin Variant



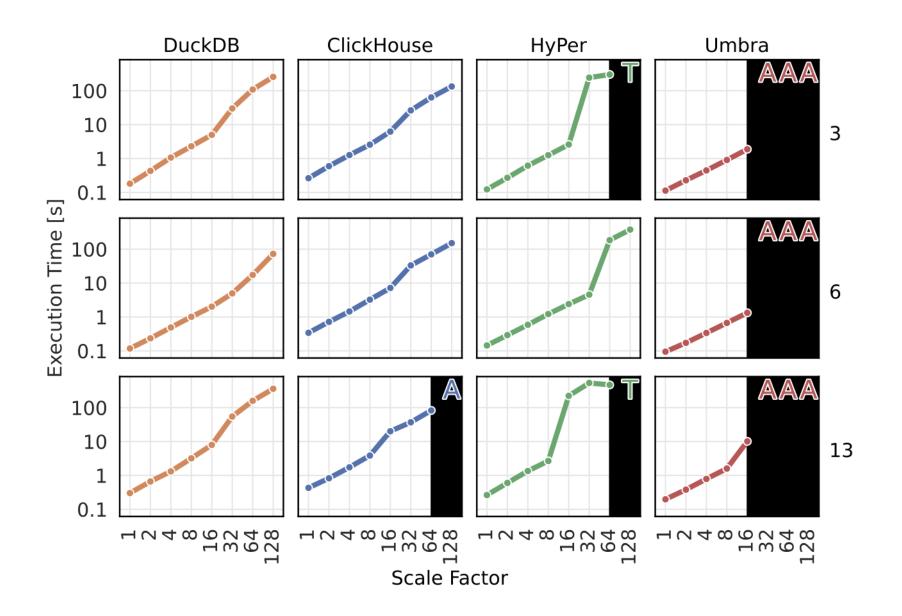
### Evaluation – Wide Variant

TABLE III

EXECUTION TIME IN SECONDS FOR THE WIDE VARIANT OF ALL GROUPINGS AT SCALE FACTORS 2, 8, 32, AND 128. LOWER IS BETTER. THE LOWEST EXECUTION TIMES ARE HIGHLIGHTED IN BOLD. 'A' DENOTES THAT THE QUERY WAS ABORTED. 'T' DENOTES THAT THE QUERY TIMED OUT AFTER 600 SECONDS. WE SUMMARIZE BY NORMALIZING EXECUTION TIMES TO DUCKDB AND THEN TAKING THE GEOMETRIC MEAN.

SF		2		[		8		1		32		1		128		
System	Du	<b>C</b> 1	Hy	Um   D	u C	C1	Ну	Um	Du	Cl	Hy	Um	Du	Cl	Hy	Um
Grouping 1 2 3 4 5 6 7 8 9 10 11 12 13	0.04 0.42 0.43 0.53 0.25 0.23 0.51 0.87 0.59 0.75 0.64 0.62	0.19 0.34 0.59 0.51 0.58 0.72 0.75 1.01 0.88 1.00 0.83 0.79 0.83	0.03 0.23 0.27 0.27 0.23 0.29 0.35 0.42 0.47 0.61 0.62 0.60	0.02     0.       0.22     2.       0.23     2.       0.23     2.       0.13     0.       0.17     1.       0.29     2.       0.36     3.       0.32     3.       0.42     3.       0.38     3.       0.38     3.       0.38     3.       0.38     3.	25 1.2 30 2.5 77 2.4 80 2.9 01 3.2 56 3.3 94 4.5 92 4.2 4.2 59 4.9 32 3.7	45 56 43 96 24 34 57 24 27 57	0.10 1.07 1.25 1.28 0.89 1.22 1.76 1.96 2.18 2.69 2.71 2.70 2.67	0.06 0.77 0.91 0.91 0.44 0.67 1.22 1.47 1.41 1.74 1.60 1.58 1.59	0.63 52.79 30.00 53.46 4.63 4.96 <b>30.61</b> 68.49 46.12 64.84 60.04 60.73 54.55	2.57 18.41 26.39 26.19 30.16 33.20 31.85 38.79 36.37 43.70 36.41 36.00 37.00	0.38 211.11 243.30 255.55 3.13 4.56 382.80 487.08 451.81 585.54 530.67 533.54 534.56	A A A A A A A A A A A A A A A A A A A	2.57 347.28 256.29 350.96 <b>67.56</b> <b>72.69</b> 260.12 407.26 331.08 399.46 396.53 382.09 359.95	9.91 111.66 133.50 122.97 A 150.75 A A A A A	1.52 499.04 T T 287.87 378.23 T T T T T	A A A A A A A A A A A A A A A A A A A
Geometric Mean Normalized to DuckDB	1.00	1.53	0.77	<b>0.54</b>   1.	00 1.4	45 (	0.70	0.45	1.00	1.06	4.54	A	1.00	A	Т	A

### Evaluation – Wide Variant



### Outline

Table data vs. temporary data

Buffer management for table data

Page-based buffer management

Buffer management for temporary data

Unified memory management

**Future work** 

### **Future Work**

- Apply the idea to OLTP setting
- Better eviction policy for unified memory management
- Adapt other blocking operators besides hash aggregation
  - E.g., join, sort, window operators
- Coordinate when multiple memory-intensive operators are active at the same time

### Questions – Buffer Management

- How common do real-world intermediates exceed memory?
- Challenges in extending to join, sorts, etc.?
- Multi-tenancy? (budget memory between simultaneous operators)
- Effective on skewed data?
- Row-major format for temporary tables? Good idea?
- How does the idea extend to distributed system?

### Discussion

Problem: How to manage spilling for temporary data

DuckDB approach: Page-based data structures as much as possible

- Benefit: Data spilling is implicitly handled by buffer manager
- Challenge: Nontrivial for certain data structures and operators

Traditional approach: Use malloc for temporary data

- Benefit: intuitive data structures and operators
- Challenge: Data spilling must be explicitly handled

Discussion question: Is there a framework that allows intuitive data structures & operators and yet supports efficient spilling?

### Before Next Lecture

Submit review for

David DeWitt, Jim Gray, <u>Parallel Database Systems: The Future of High Performance Database Processing</u>. Communications of the ACM, 1992