Optimising Column Order for Better Performance and Maintainability

Discover how the arrangement of fixed-length and variable-length columns can significantly enhance database efficiency and maintainability.









A Little About Me

Amul Sul

Å

Database Developer EnterpriseDB



12 Years

Expertise in PostgreSQL internal development



Pune, India Office Location



Why It Matters







Faster Queries Improved user experience Reduced Costs Lower infrastructure needs Competitive Edge Outperform competitors





Agenda







Access Time from CPU to DISK





Latency Numbers Every Programmer Should Know

Operation	Latency
L1 cache reference	1 ns
L2 cache reference	4 ns
Main memory reference	100 ns
Read 1 MB sequentially from memory	3,000 ns
SSD random read	16,000 ns
Read 1 MB sequentially from SSD	49,000 ns (49 µs)
Disk seek	2,000,000 ns (2 ms)
Read 1 MB sequentially from disk	825,000 ns (825 µs)

LATENCY





Latency Numbers Every Programmer Should Know

Operation	Latency	1 ns = 1 sec
L1 cache reference	1 ns	1 sec
L2 cache reference	4 ns	4 sec
Main memory reference	100 ns	100 sec
Read 1 MB sequentially from memory	3,000 ns	3,000 sec (50 min)
SSD random read	16,000 ns	16,000 sec (4.4 hr)
Read 1 MB sequentially from SSD	49,000 ns (49 µs)	49,000 sec (13.6 hr)
Disk seek	2,000,000 ns (2 ms)	2,000,000 sec (23 days)
Read 1 MB sequentially from disk	825,000 ns (825 µs)	825,000 sec (9.5 days)





What is a Value?



- 4.0
- 'four'



Source: EDB internal talk by Robert Haas



What is a Value? (More Accurate Version)

- 4::pg_catalog.int4
- 4::pg_catalog.int8
- 4::pg_catalog.int2
- **4.0**::pg_catalog.numeric
- **4.0**::pg_catalog.float8
- 'four'::pg_catalog.text'
- 'four'::pg_catalog.varchar



Source: EDB internal talk by Robert Haas





Fixed Length vs. Variable Length

Fixed-Length

Predefined size: BOOL, CHAR, INT, FLOAT, etc

E.g. pg_catalog.int4 means a 4-byte integer, it should be stored using a fixed amount of storage, namely 4 bytes.

Variable-Length

Dynamic size: TEXT, BYTEA, VARCHAR, etc.

Can store as little as 0 bytes - that is, an empty string - and as much as 5 bytes less than 1GB, it should be stored using a variable amount of storage.



Understanding Type "Lengths"

pg_type.typlen

Indicates the length of a datatype.

— Positive values

Indicates that the data type is fixed-length.

- Negative value

Indicates that the data type is variable-length.

-1: Indicating that a value is stored as a varlena (most of).

-2: Indicating that a value is stored as a cstring (very few).

There are no other possibilities (currently).



2

З

4





Tuple Construction Strategy

- 23-byte header.
- Null bitmap, if the tuple contains any null columns, with 1 bit per column.
- Pad out to an 8 byte boundary.
- Column values, possibly separated by more alignment padding.
 - Column alignment is defined by **pg_type.typalign** 0 $\mathbf{c}' = 1, \mathbf{s}' = 2, \mathbf{i}' = 4, \mathbf{d}' = 8.$
 - A column must start at an offset which is a multiple of the 0 required alignment.



Tuple Construction Example

CREATE TABLE foo (a int2 not null, b int4 not null);

- Bytes 0-22: Tuple header.
- Byte 23: Padding byte, so that first column starts on a multiple of 8.
- Bytes 24-25: Column a.
- Bytes 26-27: Padding bytes, so that second column starts on a multiple of 4.
- Bytes 28-31: Column b.

```
tuple (24 + 8 => 32 bytes):
```

24	25	26	27	28	29	30	31
	SMAI	LINT			IN	IT	





CREATE TABLE example (

- cl SMALLINT,
- c2 BIGINT,
- c3 INT,
- c4 BIGINT);
- => (1, 1, 1, 1)

=> SELECT pg_column_size(t.*) - 24 AS row_size FROM example t; row_size 32 (1 row)

tuple (24 + 32 => 56 bytes):

24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49
SMALLINT										BIG	INT	-						IN	IT						





CREATE TABLE optimized_example (

- BIGINT, Cl
- c2 BIGINT,
- c3 INT,
- c4 SMALLINT);
- => (1, 1, 1, 1)

=> SELECT pg_column_size(t.*) - 24 AS row_size FROM optimized_example t; row_size 22 (1 row)

tuple	(24 +	22 => 46	bytes):
-------	-------	----------	---------

24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45
	BIGINT									BIG	INT					IN	IT		SMAI	LINT	



Variable Length – varlena

Simple varlena - 4 byte length word

- The simplest form of varlena is a 4-byte length word where the length value can range from 4 to 1 byte less than 1GB.
- The length word is followed by the payload bytes.
- This form of varlena requires 4-byte alignment.

— Short verlena - 1 byte length word

- The varlena length is 1-127 bytes, and so the payload is 0-126 bytes.
- If non-zero, it's the beginning of a varlena with a 1-byte header. If it's zero, skip to the next 4 byte boundary; a 4-byte varlena header begins there.
- If we are starting on a 4-byte boundary, read 1 byte initially. The value we read will tell us whether it's the first byte of a 4-byte header, or the only byte of a 1-byte header.





Variable Length – TOASTing

The Oversized-Attribute Storage Technique - TOAST

- PostgreSQL uses TOAST to store large TEXT values.
- By default, large variable length values are compressed and sometimes stored out-of-line.
- Why do we have TOAST?

2

3

- Limiting the size of a tuple to what can fit into a single 8kB page would be unacceptable.
- We need a way to take a tuple that might be quite large and turn it into one or more tuples each of which can fit into an 8kB page.
- How does TOAST work?
 - Replace larger varlenas with smaller ones.
 - Repeat until the tuple is as small as you need or want it to be, or until there's nothing else that can be done.
 - Fixed-size data types are not affected by TOAST.



CREATE TABLE foo (a int2, b text, c text); INSERT INTO foo VALUES (42, 'hello', repeat('test or something', 1000000));

SELECT pg_column_size(42::int2), pg_column_size('hello'::text), pg_column_size(repeat('test or something'::text, 1000000));

SELECT pg_column_size(a), pg_column_size(b), pg_column_size(c) from foo;

SELECT pg_column_size(a+0), pg_column_size(b || "), pg_column_size(c || ") from foo;



CREATE TABLE foo (a int2, b text, c text); INSERT INTO foo VALUES (42, 'hello', repeat('test or something', 1000000));

SELECT pg_column_size(42::int2), pg_column_size('hello'::text), pg_column_size(repeat('test or something'::text, 1000000)); => **2, 9, 1700004**

SELECT pg_column_size(a), pg_column_size(b), pg_column_size(c) from foo; => **2, 6, 194620**

SELECT pg_column_size(a+0), pg_column_size(b || "), pg_column_size(c || ") from foo; => **4**, **9**, **1700004**



CREATE TABLE foo (a int2, b text, c text); INSERT INTO foo VALUES (42, 'hello', repeat('test or something', 1000000));

SELECT pg_column_size(a), pg_column_size(b), pg_column_size(c) from foo; => 2, 6, 194620

tuple (24 + 26 => 50 bytes):







CREATE TABLE foo (a int2, b text, c text);



Variable Length - Optimization

STORAGE PLAIN for frequently accessed small values
If most of your TEXT values are small (<2 KB) and frequently accessed,
storing them with PLAIN storage avoids unnecessary TOAST
compression overhead.

STORAGE EXTERNAL for large values

If TEXT values are large (>2 KB) and infrequently accessed, storing them in TOAST reduces main table size, improving query performance.

STORAGE EXTENDED (default) if compression helps If compression significantly reduces storage.

Configuration

3

4

Some of this behavior is configurable using ALTER TABLE .. SET STORAGE or by setting the toast_tuple_target relation option.



Variable Length - Placement

Reduces Row Size and Improves Cache Efficiency

Placing variable-length columns at the end minimizes padding waste, reduces row size, and improves cache locality by allowing more rows per page.

Optimizes TOAST Usage & Minimizes Data Reads

Placing variable-length columns at the end allows PostgreSQL to access fixed-length data first, avoiding unnecessary TOAST retrieval, as large values (>2KB) are stored out-of-line.

Improves Index Efficiency

PostgreSQL indexes store row references (ctid) without large column values, so placing variable-length data at the end ensures compact row storage, optimizing indexed queries like WHERE id = 123.



2

3



Alignment and Padding

- Memory alignment and padding are fundamental concepts in low-level programming, particularly in C, C++, and system-level programming.
- Alignment is arranging data in memory at addresses that are multiples of the data type's size.
- PostgreSQL defaults to 8-byte and for the data type it uses pg_type.typalign alignment





0

Why Alignment and Padding

PostgreSQL is a performance-sensitive system that handles large amounts of structured data. Proper memory alignment ensures

Efficient CPU access to structured data in shared buffers and tuples.

Reduced cache misses and better CPU cache utilization.

Avoidance of unaligned memory access penalties, especially on architectures like ARM.

Optimization of struct layouts to minimize wasted memory while ensuring correctness.



Efficient CPU access

Modern CPUs read and write memory in fixed-size chunks known as cache lines or words, rather than byte-by-byte. These chunks are typically 4, 8, or 16 bytes, with 8 bytes (64-bit word) being the most common in modern architectures.

Efficiency in Memory Fetching:

CPUs fetch data from RAM in blocks to reduce the number of memory accesses. Instead of fetching a single byte at a time (which is slow), CPUs read multiple bytes in one go.

Alignment with CPU Registers:

The word size typically matches the CPU's register size (e.g., 4 bytes for 32-bit, 8 bytes for 64-bit), ensuring efficient aligned memory access and avoiding extra operations caused by misalignment.

Cache Optimization:

CPUs use L1/L2/L3 caches to store frequently accessed data, optimizing performance by fetching memory in 8-byte or larger blocks.



2

3

Example: CPUs Reads

- Suppose a CPU needs to read an BIGINT (8-byte integer) stored at memory address igodol0x1002 (unaligned).
- The CPU, working in 8-byte chunks, will need to read both 0x1000-0x1007 and 0x1008–0x100F, causing an extra fetch.
- If the data was aligned at 0x1000 or 0x1008, only one memory fetch would be needed—improving efficiency.

		CPL	J cycle	e 1				CPU c	ycle 2						
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	
0x1000	0x1001	0x1002	0x1003	0x1004	0x1005	0x1006	0x1007	0x1008	0x1009	0x100A	0x100B	0x100C	0x100D	0x100E	(
in	t2	int8													





Alignment In PostgreSQL

Tuple Storage Ensures proper field alignment for efficient access.

Shared Buffers

Memory pages must be aligned to avoid performance penalties.

Index Structures

Proper alignment improves lookup speed.

- WAL Logging

WAL records are aligned for efficient sequential writes.

Dynamic Memory Allocation

Ensures all allocated memory follows alignment constraints.



2

3

4

5





Advantages of Column Order Optimization

Pros

Optimised storage

Reduced I/O



Storage efficiency

- => CREATE TABLE example (c1 SMALLINT, c2 BIGINT, c3 INT, c4 BIGINT);
- => INSERT INTO example SELECT 1, 1, 1, 1 FROM generate_series(1, 10000000);
- => SELECT pg_size_pretty(pg_total_relation_size('example'));

pg_size_pretty

575 MB

=> CREATE TABLE optimized_example (c1 BIGINT, c2 BIGINT, c3 INT, c4 SMALLINT);

=> INSERT INTO optimized_example SELECT 1, 1, 1, 1 FROM generate_series(1,10000000); => SELECT pg_size_pretty(pg_total_relation_size('optimized_example')); pg_size_pretty

498 MB

←---- 77 MB less



ALLINT); 000000);

Reduced I/O

1 2 3

Faster Sequential Scans Improve data locality for faster retrieval.

Faster Maintenance Activity

Minimized total amount of read

Index Optimization

Optimize the index structure to enhance lookup speed. (E.g index only scans).







Limitations of Column Order Optimization

Cons

Increased query complexity.

Reduced maintainability.



Maintainability Considerations





Consider long-term costs of

When to Optimize Column Order



Frequently accessed columns Reduce I/O operations.

Fixed-width columns

Optimize storage usage.

High cardinality columns

For efficient filtering.



Modern Hardware Advancements

SSD Adoption

Increased adoption of solid-state drives (SSDs).

Memory Capacity

Larger memory capacities.





2

Key Takeaways

Strategic order Optimize performance.



Regular analysis

Adapt to changes.



