



# Introducing **LSM-tree** into PostgreSQL, making it as a data gobbler



By Shichao Jin

# About Me



- . 8 years experience in DBMS R&D
- . Redshift AWS and Facebook RocksDB
- . Ex-PhD student at UWaterloo
- . Fan of PostgreSQL
- . Founder of VidarDB





# Outline

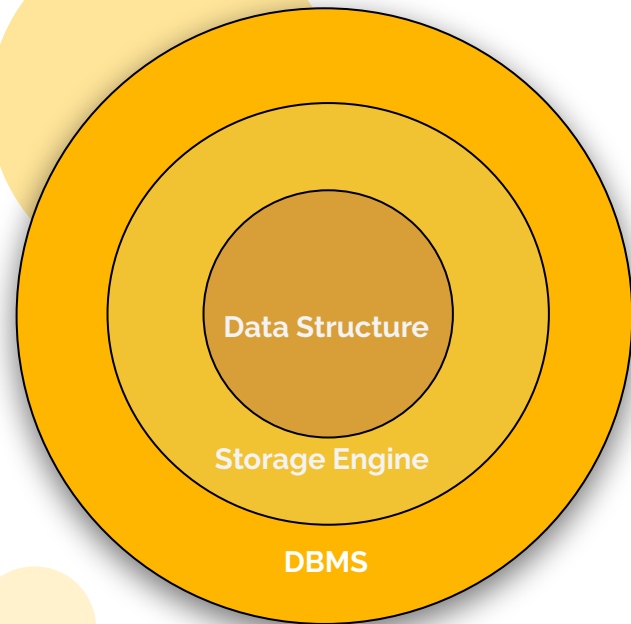
- **Experiment & Analysis**
- **LSM-tree Introduction**
- **Our Implementation**





# Background

- Storage engine accounts for a big chunk of the performance
- What is “storage engine”?  
incarnation of a data structure





# Background

- B+tree/B-tree: InnoDB, BerkeleyDB, WiredTiger
- LSM-tree: LevelDB, **RocksDB**





# Experiment Setting

## Hardware:

CPU: Xeon E3-1265L, 2.5GHz, 4 cores

Disk: 7200 RPM

RAM: 16GB DRR3

## Software:

TPC-H 10GB, **shuffled!**

Ubuntu 18.04, PostgreSQL11.6, RocksDB 6.2, default setting  
libpqxx



# Experiment Results

TPC-H Lineitem 1GB and 10GB insertion throughput (tuples per second)



PGCon 2020



VidardB  
Simplify Your Backend

7

# Why huge difference?



**Data structure!**

PostgreSQL uses **B+tree!**

So what's wrong with B+tree?

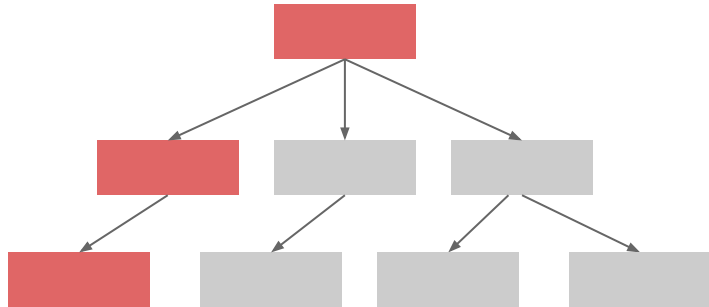




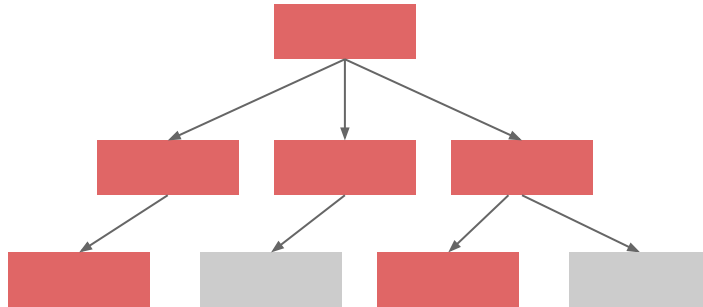


# B+tree Reading Pattern

Rare case: data come in order



Common case: data come in **random** order



Red rectangle means disk page cached!



PGCon 2020



VidarbDB  
Simplify Your Backend



# Problem with B+tree

I only talk about the insertion case here, but the root cause comes from the reading pattern of B+tree. Every insertion of B+tree needs to trigger **random reading** to locate the exact position of B+tree.

Note: Reading uses **cache**!





# Problem with B+tree (Cont.)

Data is large, and usually come not in the order of primary key, the previous cached disk blocks are useless!

**B+tree insertion causes too many disk access!**





# Disk vs. SSD vs. RAM in my laptop

## Sequential vs. Random

Seq

Ran

Magnetic Disk			SSD			RAM		
Seq			Seq			Seq		
Read (MB/s)			Read (MB/s)			Read (MB/s)		
Write (MB/s)			Write (MB/s)			Write (MB/s)		
SEQ1M Q8T1	99.61	102.33	SEQ1M Q8T1	529.80	302.44	SEQ1M Q8T1	5240.53	7353.85
SEQ1M Q1T1	100.05	102.34	SEQ1M Q1T1	500.25	403.33	SEQ1M Q1T1	6381.45	7376.68
RND4K Q32T16	0.67	0.99	RND4K Q32T16	147.04	238.80	RND4K Q32T16	1161.90	789.33
RND4K Q1T1	0.34	0.68	RND4K Q1T1	19.42	39.89	RND4K Q1T1	659.70	424.08

Magnetic Disk

SSD

RAM



PGCon 2020



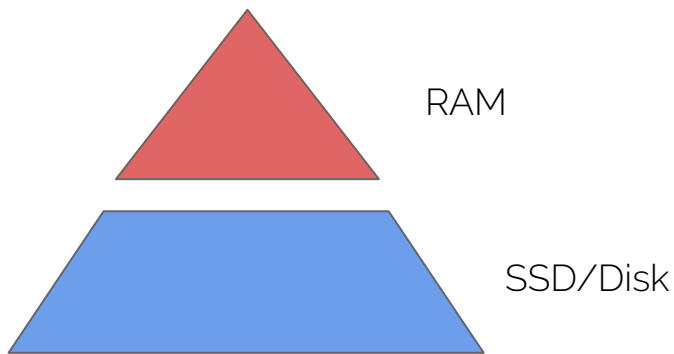
Vidardb 12  
Simplify Your Backend



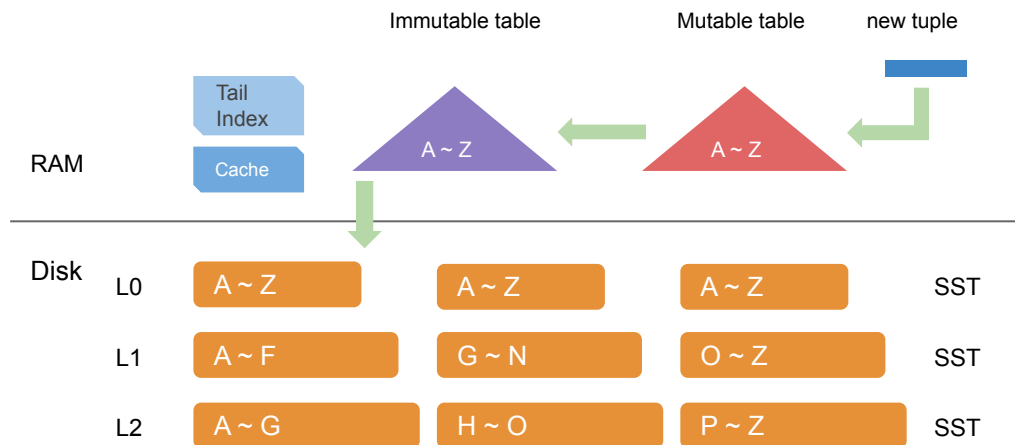
# Solution

No touch of disk (excluding WAL) when insertion!  
Tiered design.

**LSM-tree (Log-structured merge-tree)** comes to rescue!



# LSM-tree



SST: Sorted String Table

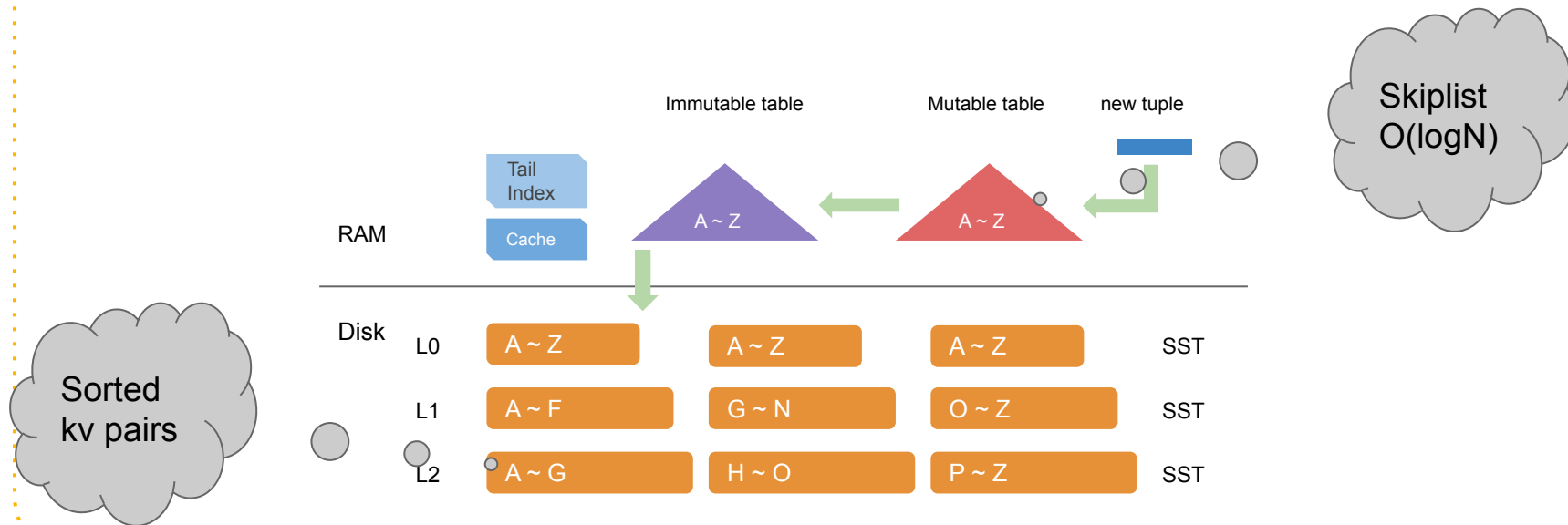


PGCon 2020



**Vidarb** 14  
Simplify Your Backend

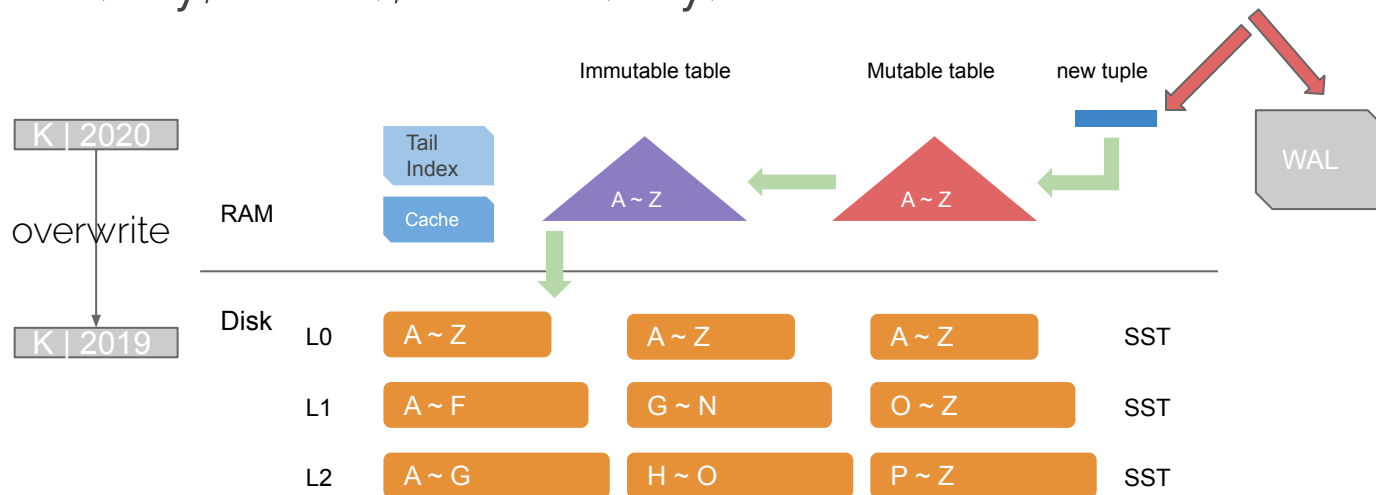
# LSM-tree



# Basic Operations



Put(key, value), Delete(key)



PGCon 2020



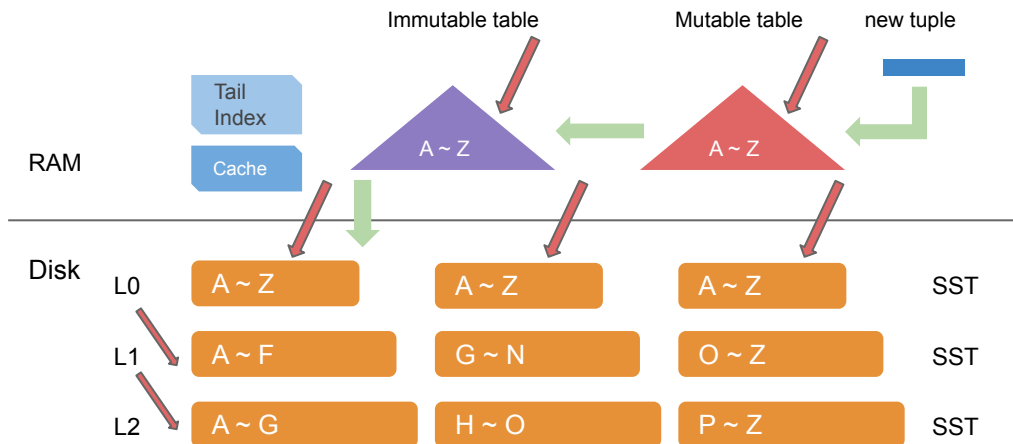
**VidardB 16**  
Simplify Your Backend



# Basic Operations



Get(Key), return value



PGCon 2020

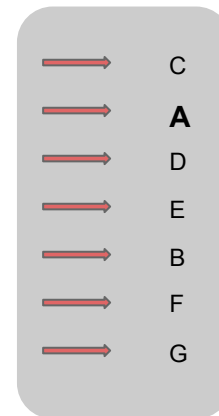
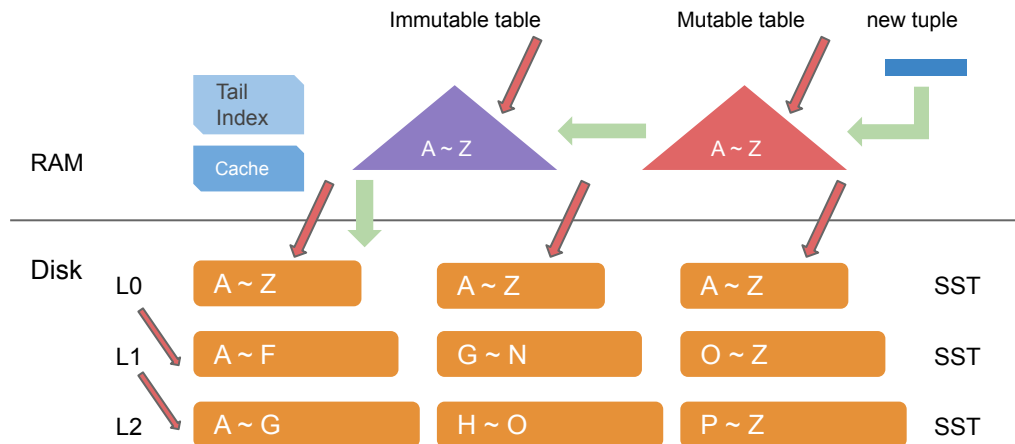


**VidardB** 17  
Simplify Your Backend

# Basic Operations



Iterator(): seek & next



PGCon 2020

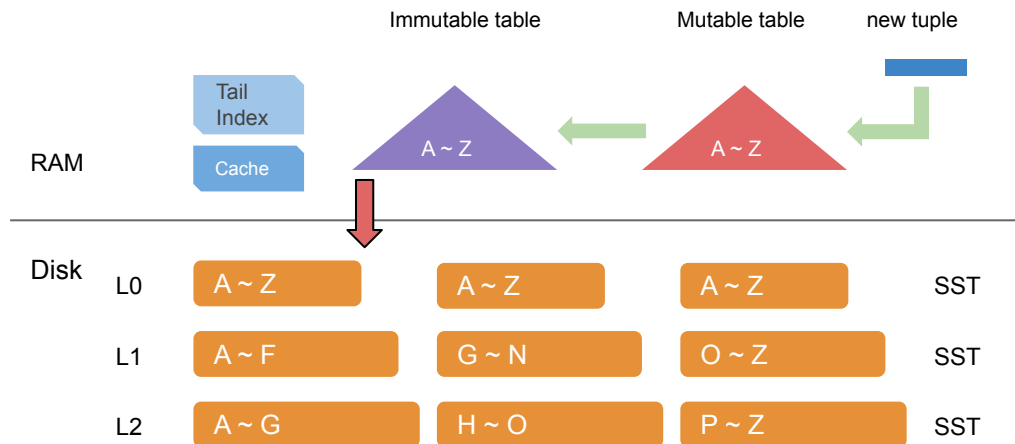


**VidarDB 18**  
Simplify Your Backend

# Advanced Operations



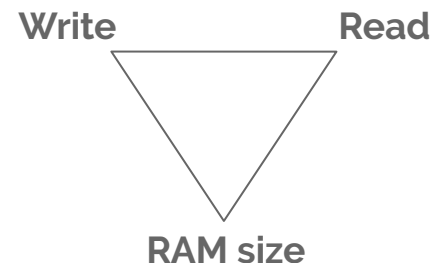
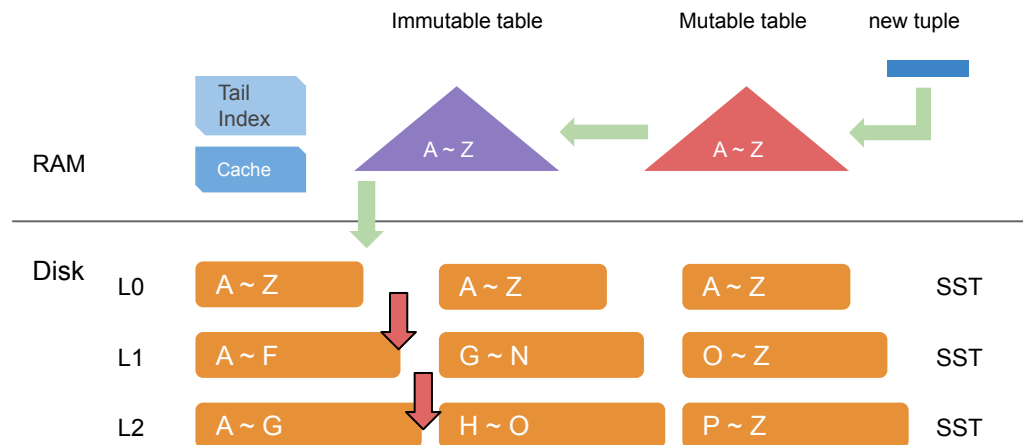
Flush(): RAM → L0, usually by the background thread





# Advanced Operations

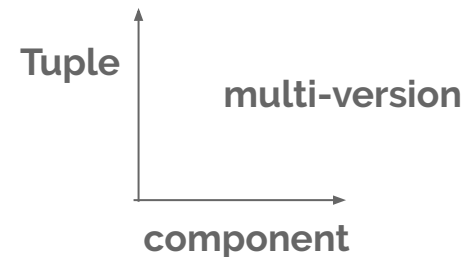
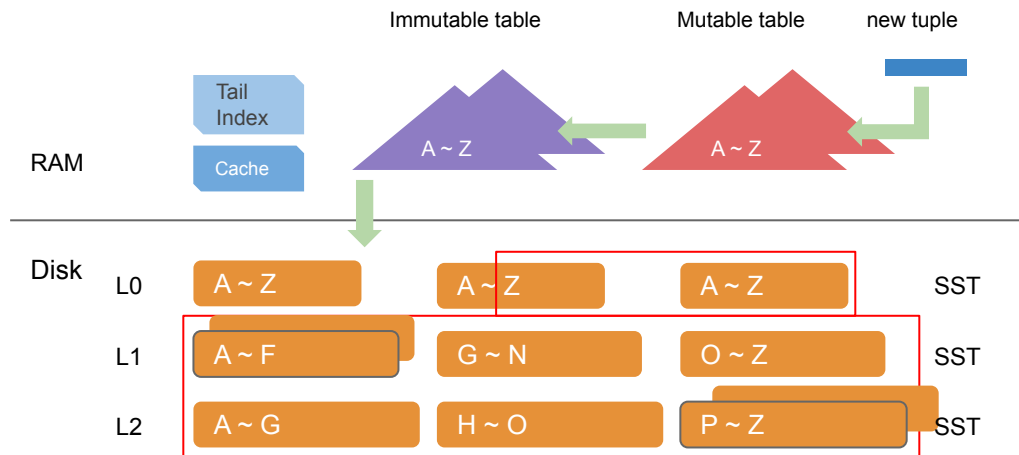
Compact(k1, k2):  $L0 \rightarrow L1$ ,  $L1 \rightarrow L2$ , usually called by the background thread





# Advanced Operations

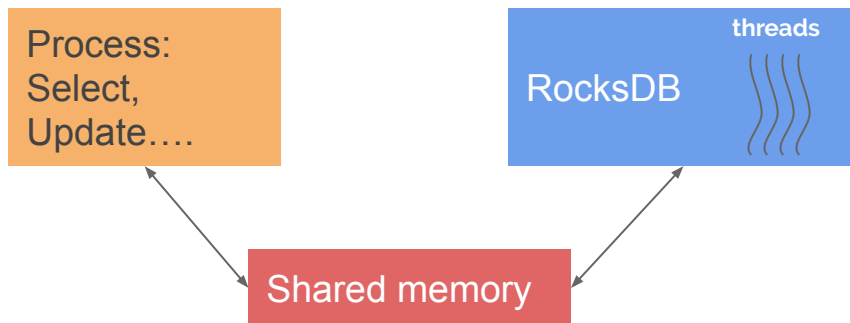
Snapshot(seq): sequence number, which is always increasing



# Implementation



- Foreign data wrapper (FDW)
- RocksDB is based on the multi-thread model
- Shared memory



# More To Do



- Engine recognizes the data types of PostgreSQL
- Transaction?
- Migrate to pluggable storage engine, table AM





# Motivation of designing a new data structure

1. random R/W of RAM is **2000X** faster than magnetic disk, but **120X** more expensive
2. Ratio of RAM/disk size is **40~200X** larger than 1970s, when B+tree is designed
3. New storage choice, NVRAM
4. Requires expertise to choose different engines, data systems





# Thanks!

[Shichao@vidardb.com](mailto:Shichao@vidardb.com)

Twitter: @jsc0218

[github.com/vidardb/vidardb](https://github.com/vidardb/vidardb)



[github.com/vidardb/PostgresForeignDataWrapper](https://github.com/vidardb/PostgresForeignDataWrapper)



PGCon 2020



**VidarDB** 25  
Simplify Your Backend