# Benefits of NVMe NVRAM vs NVDIMM, a database application example

## Jerome Gaysse – IP-Maker

# Goal

- OLTP database performance comparison with different storage options

  - Full flash SSD

  - NVMe NVRAM

  - NVDIMM

# Methodology

- Part 1 - Existing hardware
  - Flash SSD, NVMe NVRAM, NVDIMM
  - MS SQL server 2016, HammerDB

- Part 2 - Estimation with new product design for higher capacity

# OLTP performance

- A question of latency

- Many small read/write accesses to the DB file
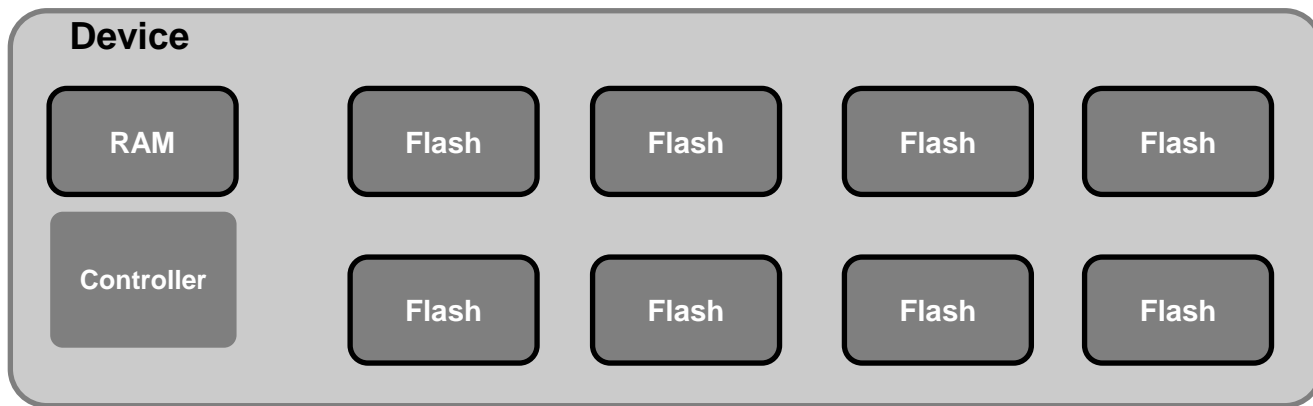
- Write accesses to the LOG file

# Part 1

- Flash SSD
  - Read 100μs latency
  - Write 500μs latency
- NVMe NVRAM
  - Read/write 12μs latency
- NVDIMM
  - Read/write 3μs latency

# Flash SSD

- Latency : 100µs/500µs
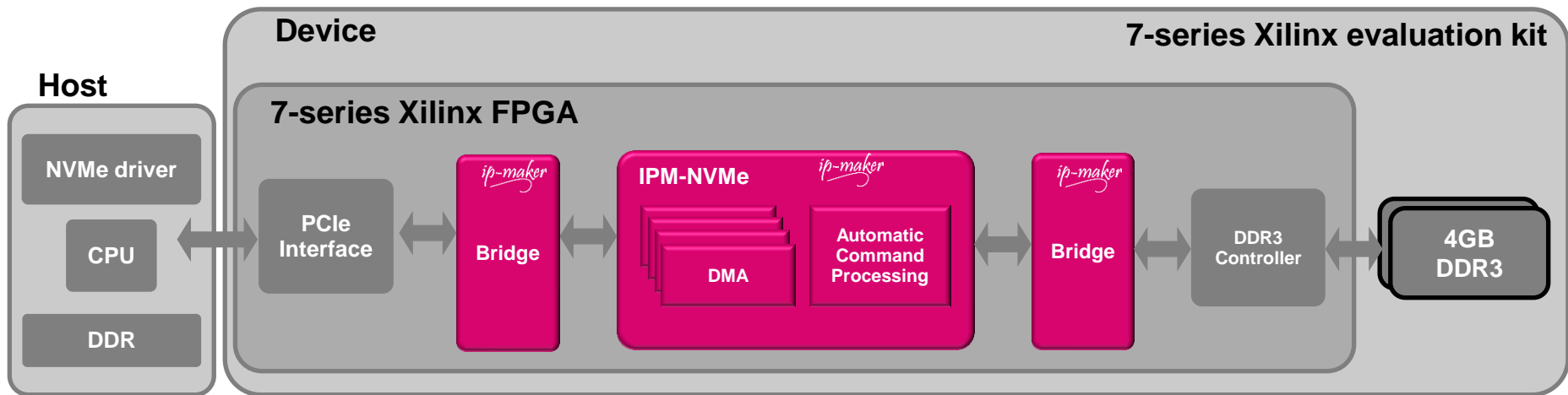
# NVMe RAM

- 12µs latency

# Latency details

12μs

**Gen3x4, FPGA demo, OS IRQ**

**NVMe IP (Command fetch + data management): 0,6μs (clock 125MHz)**

File system+ NVMe driver: 2,85

Doorbell, Command read: 1.5μs

Data transfer: 1.1μs (Gen3 speed, 4 lanes)

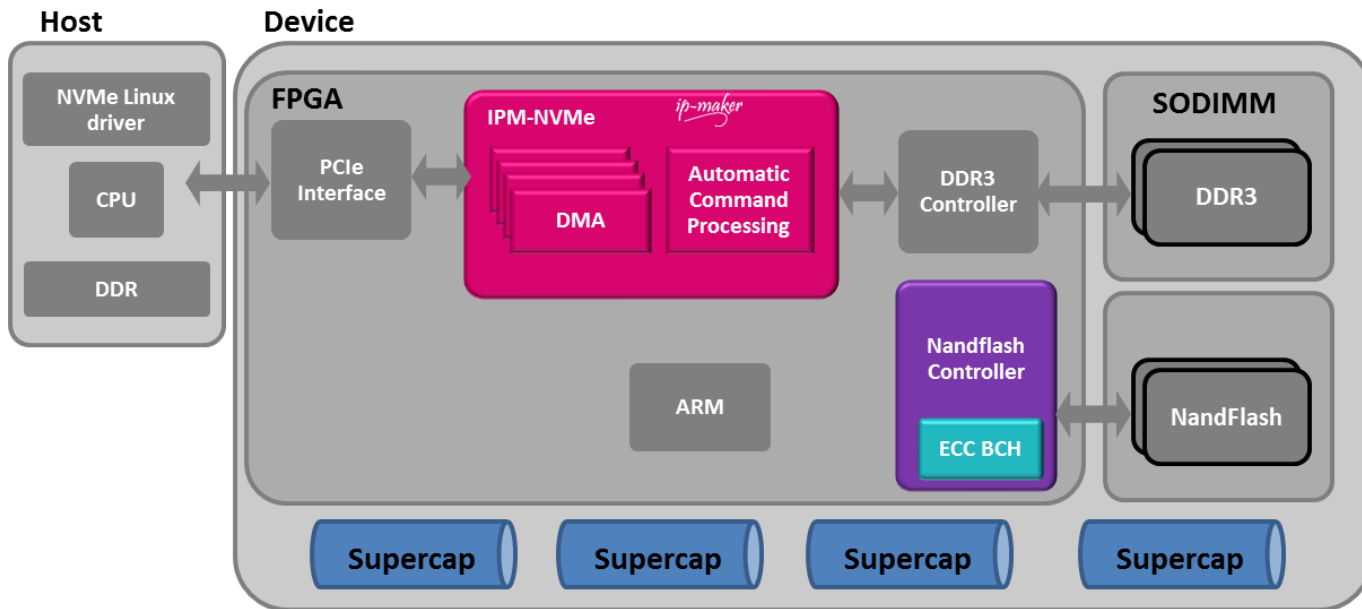Host IRQ management + PCIe latency: 7.7

# Options to reduce latency

- PCIe gen 4
- Command Memory Buffer (CMB)
- Command Memory Buffer (CMB) with persistent memory
- Polling mode

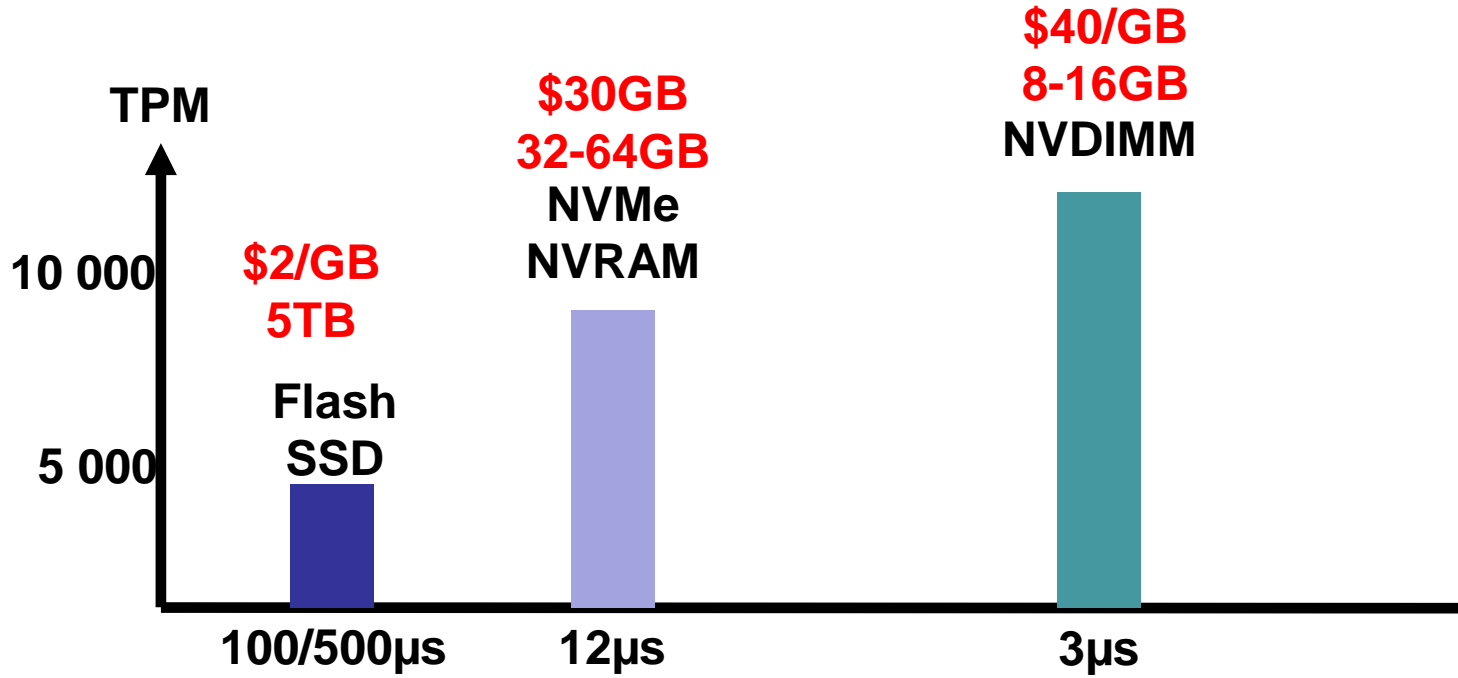- => from 12 to 5µs

# NVMe NVRAM Implementation (NVDIMM-N like)

# NVDIMM

- NVDIMM simulation using:
  - 4GB LRDIMM
  - RAM disk software
- Latency measured with FIO: 3µs

# Performance results

# The price for performance

- Flash: $2/GB, 5TB
    - 4K TPM
- NVMe NVRAM: $30/GB, 32GB
    - 10K TPM
- NVDIMM: $40/GB, 8GB
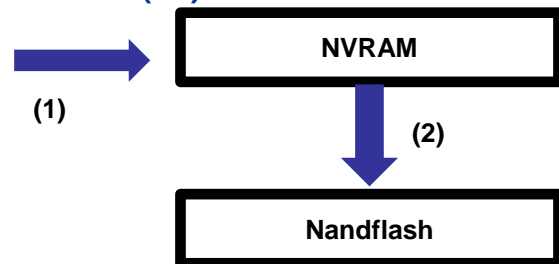    - 14K TPM

- What about TCO for TB database?

# Part 2

- NVMe NVRAM
  - High storage capacity ?

- NVDIMM
  - High storage capacity ?

# NVMe NVRAM Product design

- Achieving high capacity and low write latency
  - Non-volatile buffer for low latency
  - Nandflash storage for high capacity
  - Highly parallel implementation for high throughput
- Based on pairs of NVRAM and nandflash memories.
  - The data is first coming from the controller (1).
  - Then it is copied in the nandflash (2).

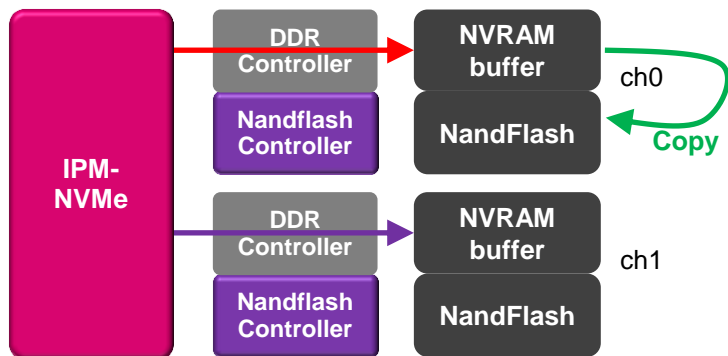**(1)** → | NVRAM |

**(2)** ↓

| Nandflash |

# Theory of operation (1/2)

- The first 4 IOs are sent to the NVRAM buffer 0.

- The second 4 IOs are sent to the NVRAM buffer 1.

- During this time, the data is read from the NVRAM buffer 0 and written into the nandflash channel 0.
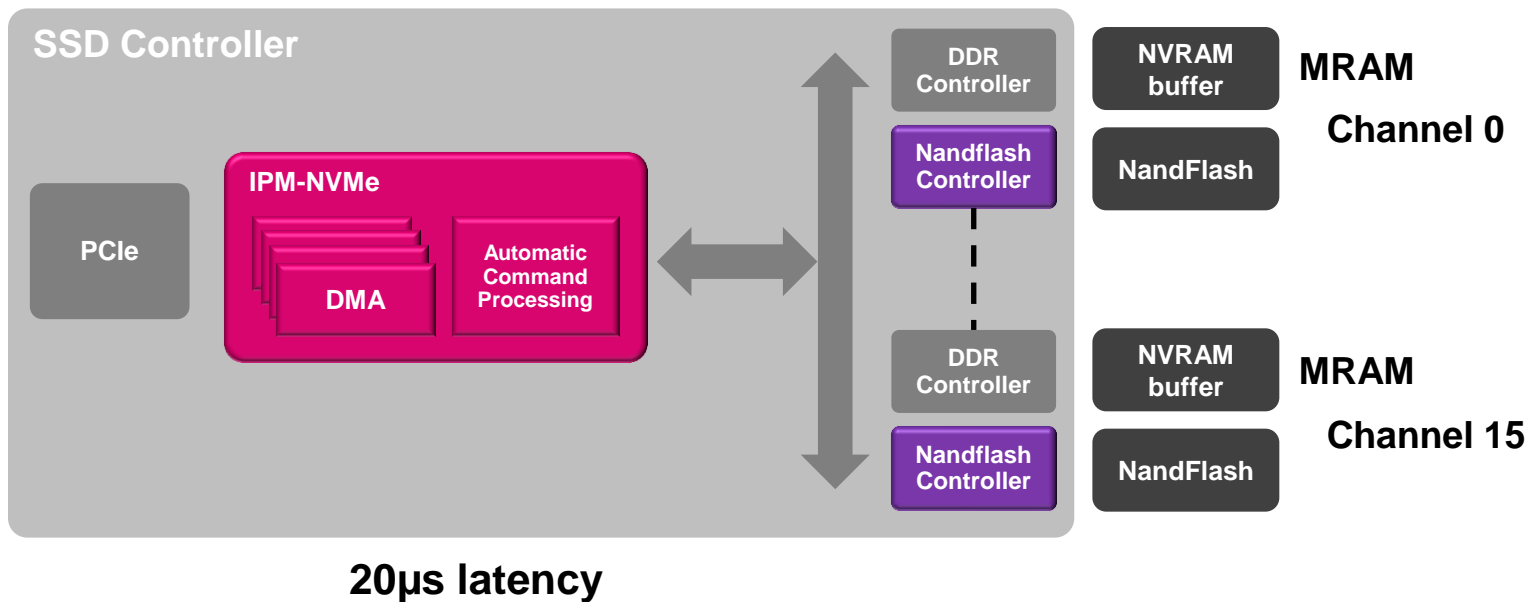
# Theory of operation (2/2)

**Sequence:**
**1=write to NVRAM CH0**
**2=write to NVRAM CH1**
**3=copy from NVRAM to Flash**
**4=prog flash**
**…**

IPM-NVMe

DDR Controller → NVRAM buffer    ch0
Nandflash Controller → NandFlash    **Copy**

DDR Controller → NVRAM buffer    ch1
Nandflash Controller → NandFlash

4IOs, IO=4kB, QD=1,
total = 4x 20µs**=80µs**

16kB@800MB/s**=20µs**

NVRAM ch0 — 4 IOs

NVRAM ch1 — 4 IOs

Prog time = **1200µs**

UNFC ch0 — Prog

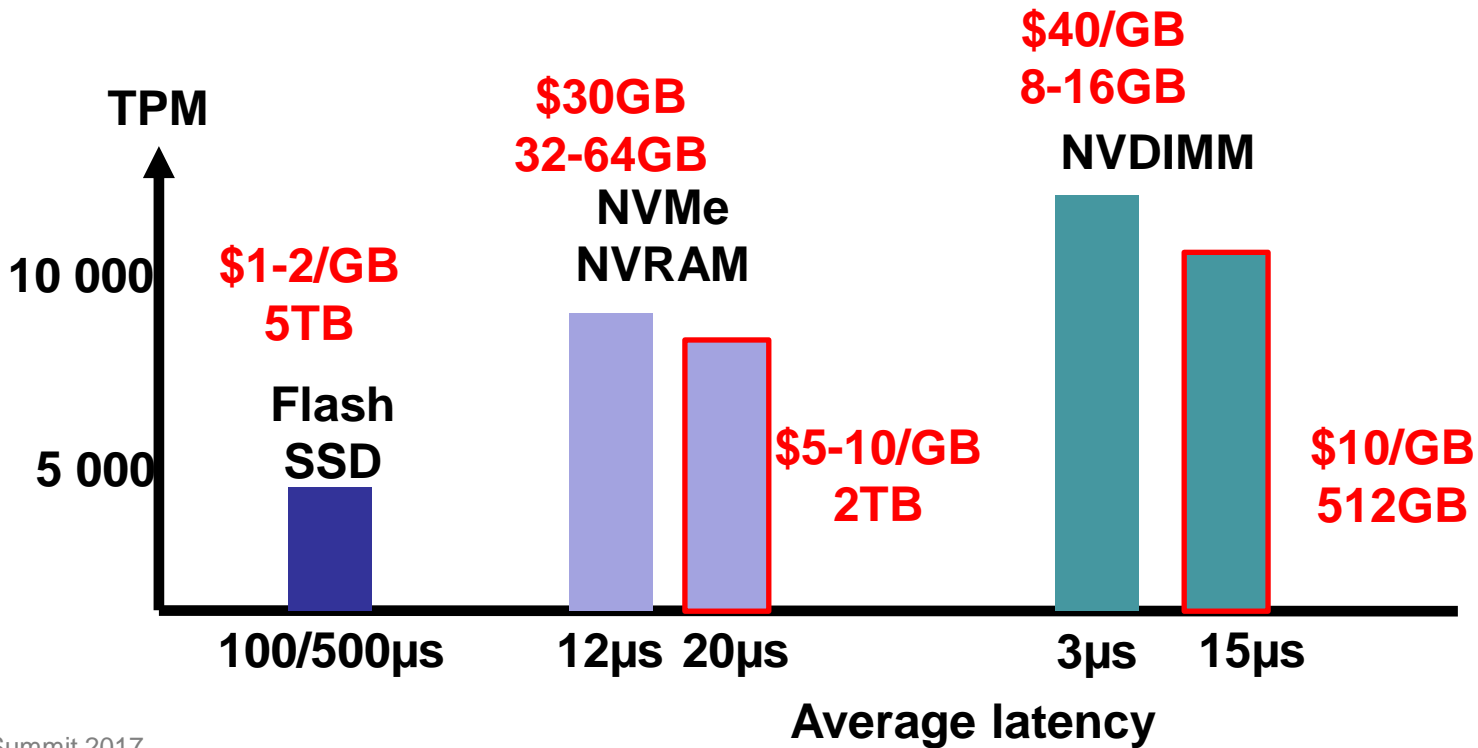# NVMe NVRAM Implementation (with MRAM)



**20µs latency**

# NVDIMM

- Higher storage capacity?
  - Yes, few hundreds of GB of Flash can be added
- Highly parallel design?
  - No, limited by PCB area
  - Average latency to increase
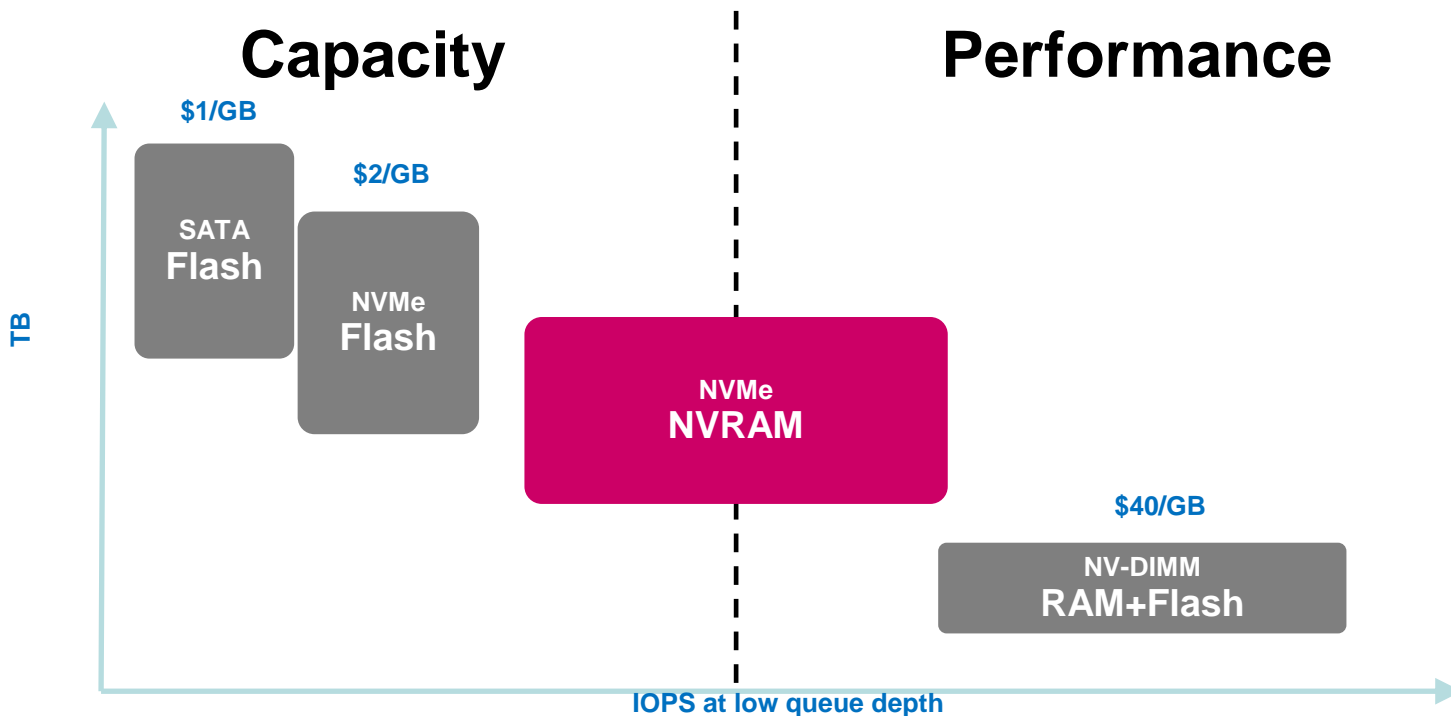
# Performance estimation



TPM

$40/GB
8-16GB
**NVDIMM**

$30GB
32-64GB
**NVMe
NVRAM**

$1-2/GB
5TB

10 000

**Flash
SSD**

$5-10/GB
2TB

5 000

$10/GB
512GB

100/500µs        12µs  20µs        3µs        15µs

**Average latency**

# The price for performance

- Flash: $2/GB, 5TB
  - 4K TPM
- NVMe NVRAM: $5/GB, 2TB
  - 9K TPM
- NVDIMM: $10/GB, 512GB
  - 13K TPM

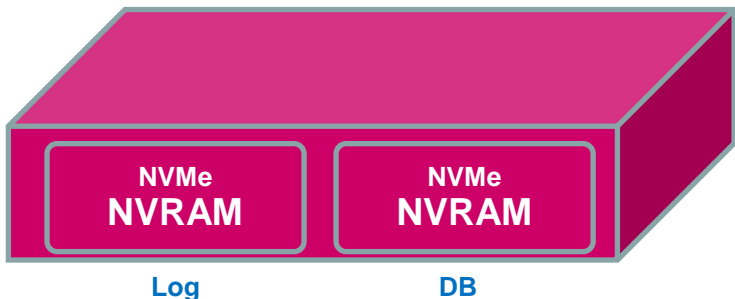# Synthesis

# NVMe NVRAM vs NVDIMM

**OLTP application**



Log            DB

**NVMe NVRAM:
for both Logs and DB files**

Jerome Gaysse

jerome.gaysse@ip-maker.com

www.ip-maker.com