

OPTIMIZE YOUR PMDK APPLICATION'S PERFORMANCE WITH THE HELP OF INTEL[®] VTUNE[™] AMPLIFIER PROFILER

Dmitry Ryabtsev

Sergey Vinogradov

LEGAL DISCLAIMER & OPTIMIZATION NOTICE

INFORMATION IN THIS DOCUMENT IS PROVIDED "AS IS". NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. INTEL ASSUMES NO LIABILITY WHATSOEVER AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO THIS INFORMATION INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT.

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit www.intel.com/benchmarks.

Copyright © 2018, Intel Corporation. All rights reserved. Intel, Pentium, Xeon, Xeon Phi, Core, VTune, Cilk, and the Intel logo are trademarks of Intel Corporation in the U.S. and other countries.

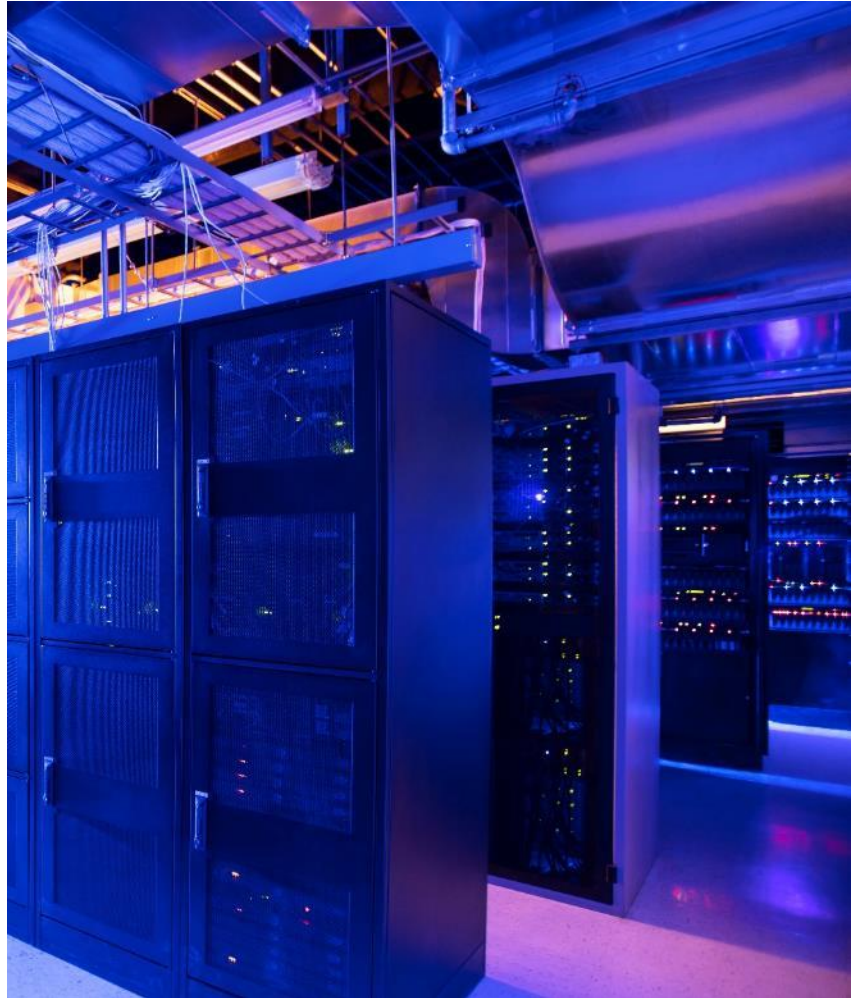
Optimization Notice

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Notice revision #20110804

AGENDA

- Overview of VTune capabilities useful for PMEM profiling
 - Characterization
 - Data Profiling
 - Bandwidth monitoring
- Case Study: PMEMKV Optimization



MAIN VTUNE CAPABILITIES FOR PMEM

3 key ingredients:

- **Application characterization from uArch perspective**

Is accessing memory really a bottleneck? How big it is?

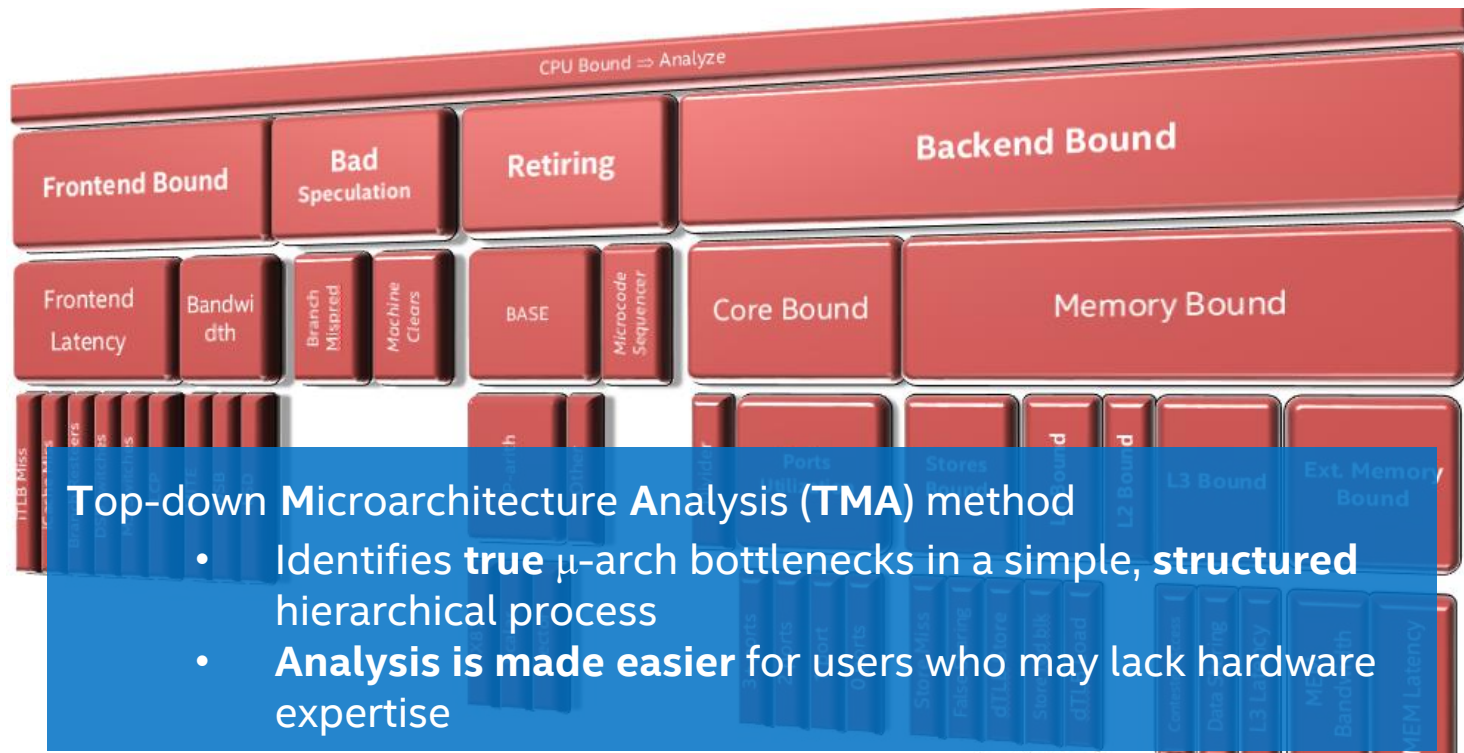
- **Bandwidth analysis**

Is my workload approaching memory bandwidth limits? What code is responsible?

- **Data profiling**

Accessing what data structures is the most problematic?

APPLICATION CHARACTERIZATION FROM UARCH PERSPECTIVE



APPLICATION CHARACTERIZATION FROM UARCH PERSPECTIVE

New metric in TMA: **“Persistent Memory Bound”**

Shows percentage of cycles CPU stalled while waiting from PMEM

Complements existing Memory Bound, DRAM Bound, etc. metrics

⌵ Memory Bound [?] :	75.0%	🚩 of Pipeline Slots
➤ L1 Bound [?] :	1.4%	of Clockticks
➤ L2 Bound [?] :	0.1%	of Clockticks
➤ L3 Bound [?] :	3.1%	of Clockticks
➤ DRAM Bound [?] :	0.0%	of Clockticks
➤ Persistent Memory Bound [?] :	57.9%	🚩 of Clockticks
➤ Store Bound [?] :	7.7%	of Clockticks

DATA PROFILING

Attributes performance metrics not only to the code but to the **data structures** as well

Works by:

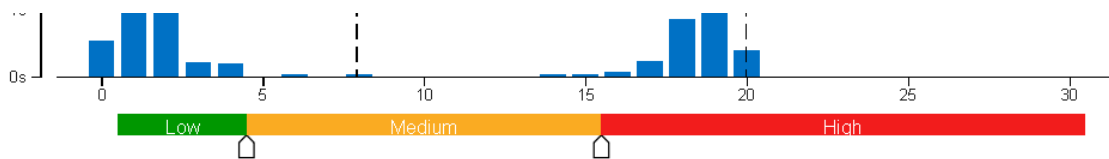
- Instrumenting memory API
- Capturing data address for PEBS facility on each memory event sample
- Correlating above two things

Latest VTune versions support PMDK memory allocation API (pmemobj_tx_alloc, pmemobj_tx_realloc, etc.)

BANDWIDTH ANALYSIS

VTune is able to measure both regular DRAM and PMEM bandwidth. Based on uncore events.

Automatically correlates with the code running at the same time – easy to see who is responsible for high bandwidth consumption



Top Memory Objects with High Bandwidth Utilization

This section shows top memory objects, sorted by LLC Misses, that were accessed when bandwidth utilization was high for the domain selected in the histogram area.

Memory Object	LLC Miss Count
array.c:95 (366 MB)	37.4%
array.c:132 (366 MB)	36.0%
array.c:134 (366 MB)	15.8%

CASE STUDY

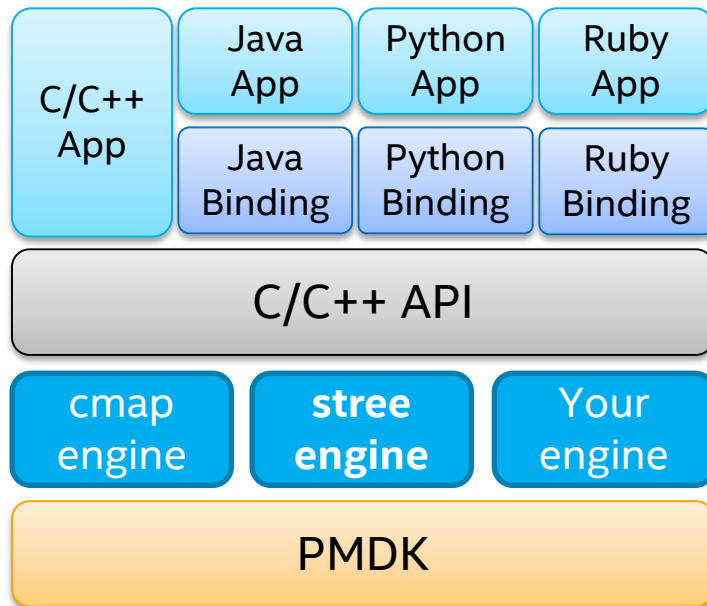
PMEMKV: FILLRANDOM BENCHMARK ON STREE ENGINE

PMEMKV

- Embedded Key/Value storage optimized for persistent memory.
- **pmemkv_bench** – db_bench ported from RocksDB.

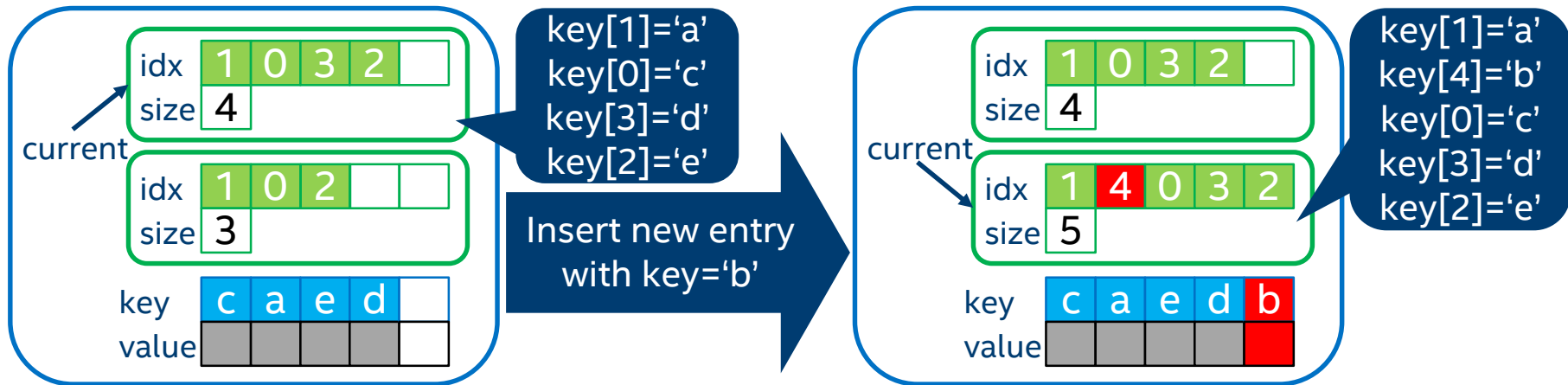
Benchmark

- Fillrandom benchmark using stree engine.
- 4 000 000 entries inserted.
- Key size -16 bytes, Value size – 200 bytes.



```
pmemkv_bench --engine=stree --num=4000000 --value_size=200 --benchmarks=fillrandom
```

INSERT INTO B+TREE LEAF NODE



- Insert key-value pair to the end.
- Merge current idx array with index of new element to the second idx array.
- Switch current pointer to merged idx array.

```
class leaf_node_t {
private:
    value_type entries[slots];
    leaf_entries_t v[2];
    uint32_t current;
};
```

VTUNE UARCH CHARACTERIZATION

Clockticks:	86,458,898,535	
Instructions Retired:	35,338,687,890	
CPI Rate [?] :	2.447 ▾	
MUX Reliability [?] :	0.899	
Retiring [?] :	11.3%	of Pipeline Slots
Front-End Bound [?] :	12.2%	of Pipeline Slots
Bad Speculation [?] :	3.1%	of Pipeline Slots
Back-End Bound [?] :	73.4% ▾	of Pipeline Slots
Memory Bound [?] :	69.8% ▾	of Pipeline Slots
L1 Bound [?] :	0.0%	of Clockticks
L2 Bound [?] :	0.6%	of Clockticks
L3 Bound [?] :	18.8% ▾	of Clockticks
Contested Accesses [?] :	0.0%	of Clockticks
Data Sharing [?] :	0.0%	of Clockticks
L3 Latency [?] :	5.0% ▾	of Clockticks
SQ Full [?] :	0.0%	of Clockticks
DRAM Bound [?] :	0.0%	of Clockticks
Persistent Memory Bound [?] :	49.1% ▾	of Clockticks
Store Bound [?] :	28.2% ▾	of Clockticks
Store Latency [?] :	24.8% ▾	of Clockticks
False Sharing [?] :	0.0%	of Clockticks
Split Stores [?] :	0.0%	of Clockticks
DTLB Store Overhead [?] :	0.3%	of Clockticks
Core Bound [?] :	3.6%	of Pipeline Slots

Our workload is memory bound

All data stored in persistent memory
We are persistent memory bound

WHAT IS WRONG WITH OUR DATA STRUCTURES?

VTune Memory Access analysis can aggregate performance data per memory objects

leaf_node_t
Memory
objects

Current field:
Offset inside
leaf_node_t

Functions that access
the field

There are a lot of LLC
misses

Grouping: Memory Object / Offset / Function / Allocation Stack

Memory Object / Offset / Function / Allocation Stack	Hardware Event Count by Hardware Event Type	
	MEM_LOAD_RETIRED.L3_MISS...	MEM_TRANS_RETIRED.LOAD_LATENCY_GT_4
db_bench.cc:629 (16 KB)	10,257,175	1,802,391
▶ 16656	1,090,763	87,261
▶ persistent::internal::leaf_node_t<pstring<(unsigned long)20>, pstring<(unsigned long)200>, pstring<(unsigned long)63>::consistent	1,050,785	7,021
▶ persistent::internal::leaf_node_t<pstring<(unsigned long)20>, pstring<(unsigned long)200>, pstring<(unsigned long)63>::consistent	40,028	0
▶ 16136	560,392	21,063
▶ 4000	240,168	1,003
▶ 3752	170,119	7,021
▶ 8464	160,112	2,006
▶ 5984	160,112	3,009
▶ 1768	140,098	3,009
▶ 6480	140,098	1,003

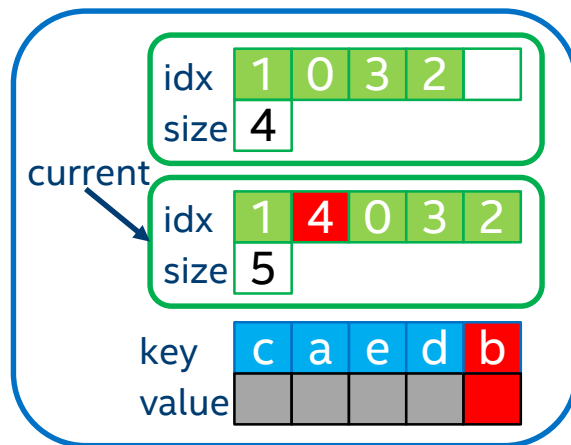
READ-AFTER-FLUSH PATTERN IN LEAF NODE

Problem:

- Reading **current** value causes a lot of LLC misses.
- Read-Modify-Flush pattern when updating **current** value.

Solution:

- Introduce shadow copy **p_current** to avoid flushes.
- Application reads and modifies **current** value. But never flush it.
- Changes of **current** value propagated to **p_current**.
- After restart **current** value restored from **p_current**.



```
void insert(...) {  
    ...  
    update(current);  
    p_current = current;  
    flush(p_current);  
    ...  
}
```

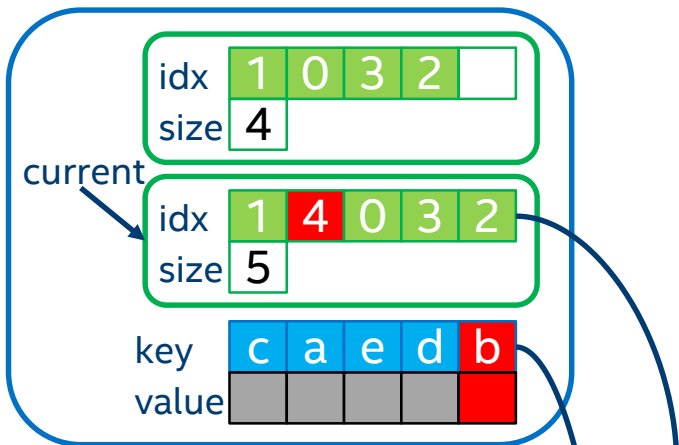
FALSE SHARING IN LEAF NODE

Problem:

- Reading **current** value still cause a lot of LLC misses.
- **current** shares the same cache line with **leaf_entries_t v[2]**.

Solution:

- Change layout of the **leaf_node_t**.
- Keep **leaf_entries_t v[2]** and **current** in different cache lines.
- `alignas(64)` for **v[2]** array and **entries** array.



```
class leaf_node_t {  
private:  
    uint32_t p_current;  
    value type entries[slots];  
    leaf_entries_t v[2];  
    uint32_t current;  
};
```

PERFORMANCE IMPROVEMENTS IN PERSISTENT B+ TREE

Version	Ops/Sec	Descriptions
0	168066	Initial version
1	170586	Introduce shadow copy to avoid read-after flush pattern
2	178589	Change layout to avoid false sharing

Throughput of stree engine on fillrandom benchmark was increased by ~10K Ops/Sec.

```
class leaf_node_t {  
private:  
    uint32_t current;  
    alignas(64) value_type entries[slots];  
    leaf_entries_t v[2];  
    alignas(64) uint32_t p_current;  
};
```


SUMMARY

- Optimizing persistent memory bound applications is a challenging task.
- Hotspot might be blurred across multiple functions accessing the same memory.
- Intel VTune helps to identify inefficiency in data structures layout.
 - Memory Analysis allows to aggregate performance data per memory objects.

Contacts:

- If you need any help to enable or optimize your applications for persistent memory feel free to contact us:
 - Dmitry Ryabtsev Dmitry.Ryabtsev@intel.com
 - Sergey Vinogradov sergey.vinogradov@intel.com

SOFTWARE ANALYSIS TOOLS FOR INTEL® OPTANE™ DC PERSISTENT MEMORY DOWNLOADS & TECHNICAL ARTICLES


software.intel.com/persistent-memory


MEMORY AND STORAGE CONVERGE

Develop innovative solutions that maximize memory capacity, data resiliency, and performance using Intel® Optane™ DC persistent memory.


Library


Training


Tools


Get Help

Click 

THANK YOU

Questions

