

Mnemo: Boosting Memory Cost Efficiency in Hybrid Memory Systems

Thaleia Dimitra Doudali Ada Gavrilovska



Problem Statement

Georgia Tech

Capacity Sizing of the Hybrid Memory Components.



Facts:

- Future clouds will feature hybrid memory components.
- These components have different cost and access latencies.

Problem:

What is the ideal *capacity ratio* between the hybrid memory components?

Goals:

- Maximize system's cost efficiency.
- Keep performance guarantees.

Existing Solutions Data Tiering over *fixed* capacities.





how much DRAM vs NVM to use.

Solution Preview



Mnemo

Offline Profiling Tool for in-memory Key-Value Stores

- Data Tiering <u>Which</u> keys should be allocated in DRAM vs NVM?
- Capacity Sizing

<u>How many</u> keys should be allocated in DRAM, so that application performance remains high, but memory cost remains low?



Mnemo quickly generates an accurate trendline of application <u>performance</u> for incremental DRAM to NVM <u>capacity ratio</u>, thus memory cost.



How much capacity of each memory type?



Which parameters affect key-value store performance over hybrid memory systems?





Which parameters affect key-value store performance over hybrid memory systems?



Varying key access pattern

////// Mnemo @ HPBDC '19



Which parameters affect key-value store performance over hybrid memory systems?



Varying read:write request ratio



Which parameters affect key-value store performance over hybrid memory systems?



Varying key-value size

Motivation Takeaways



Which parameters affect key-value store performance over hybrid memory systems?

- Key access pattern
- Read:Write requests
- Key-Value sizes

These parameters determine the shape of the curve.

The height of the curve also depends on the latency difference in accessing DRAM vs NVM.



Takeaways:

In order to estimate performance we'll need to capture:

- The workload parameters.
- Performance baselines for DRAM vs NVM.





User Input



Provides just a workload description



Minimal User Effort No Application Modifications



User chooses the sweet spot that Output 🚺 Maximum Cost Efficiency brings the desired performance to User Desired performance levels under his cost budget.

Detailed Design IP of server with DRAM DRAM-NVM IP of server with NVM price difference **Sensitivity Engine** Workload Execution over DRAM-only and NVM-only. Performance 👉 baselines **Estimate Engine**

Georgia Tech



Evaluation Methodology





Estimate Accuracy ? Cost Efficiency ? Profiling Overhead ?





Evaluation Results



Mnemo successfully captures the trade-off between performance and cost



Evaluation Results Mnemo allows for significant cost reductions





For chosen performance sweet spot to be 10% slowdown from all-data-in-DRAM.

Evaluation Results







Profiling Overhead





Summary





https://github.com/Thaleia-DimitraDoudali/mnemo