



Adding Machine Learning to the Management of Heterogeneous Resources

Thaleia Dimitra Doudali

The Era of Data

“More than **59 ZB** of data will be created, captured, copied, and consumed in the world this year.”

Source: International Data Corporation, May 2020.

Exploded
Data Sizes



Data Analytics Pipeline

ZBs of data

Capture

Process

Store

Analyze

Use

Need for speed and massive storage capacities!

The Era of Heterogeneous Hardware

Emerging technologies across layers and vendors.

Compute Acceleration

Nvidia GPUs

Titan X
VS
GTX 980
VS
Tesla M40
VS
Tesla K80

ARE DPUS THE NEXT DATACENTER REVOLUTION?

NVIDIA

Cloud TPU v2
180 teraflops
64 GB High Bandwidth Memory (HBM)



Data Storage Acceleration

intel OPTANE™
PERSISTENT MEMORY

AMD

HIGH BANDWIDTH MEMORY

V-NAND SSD 980 PRO
PCIe 4.0 NVMe M.2
2TB
SAMSUNG

Network Acceleration

Mellanox Innova™-2 Flex
Open Programmable SmartNIC



Interconnection Standards



Gen-Z Consortium

Industry Leaders developing a memory-semantic interconnect

AMD ARM BROADCOM CAVIUM CRAY

DELL EMC Hewlett Packard Enterprise HUAWEI IBM IDT

Lenovo Mellanox Micron Microsemi redhat

SAMSUNG SEAGATE SK hynix WDC Western Digital XILINX

The Era of Heterogeneous Hardware

Across computing platforms.

Supercomputers

Datacenters



- Home
- Technologies
- Sectors



Number of Nodes	4,608
Node performance	42 TF
	512 GB DDR4 + 96 GB HBM2
	1600 GB
	>10 PB DDR4 + HBM2 + Non-volatile
	2 IBM POWER9™ 9,216 CPUs 6 NVIDIA Volta™ 27,648 GPUs
	250 PB, 2.5 TB/s, GPFS™
	13 MW
	Mellanox EDR 100G InfiniBand
	Red Hat Enterprise Linux (RHEL) version 7.4



Available first on Google Cloud: Intel Optane DC Persistent Memory

A2 VMs now GA—the largest GPU cloud instances with NVIDIA A100 GPUs



Personal Devices

70% faster ML accelerators

New image signal processor

80% faster Neural Engine

50% faster CPU Faster than any other smartphone chip

First 5 nm chip A14 in a smartphone

16-core Neural Engine

Apple A14

Machine learning controller Best machine learning platform in a smartphone

6-core CPU

50% faster GPU Faster than any other smartphone chip

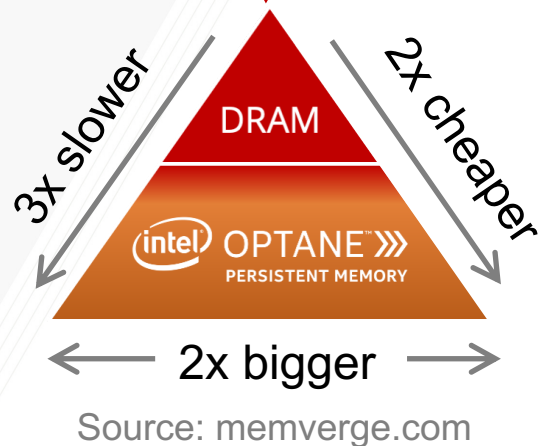
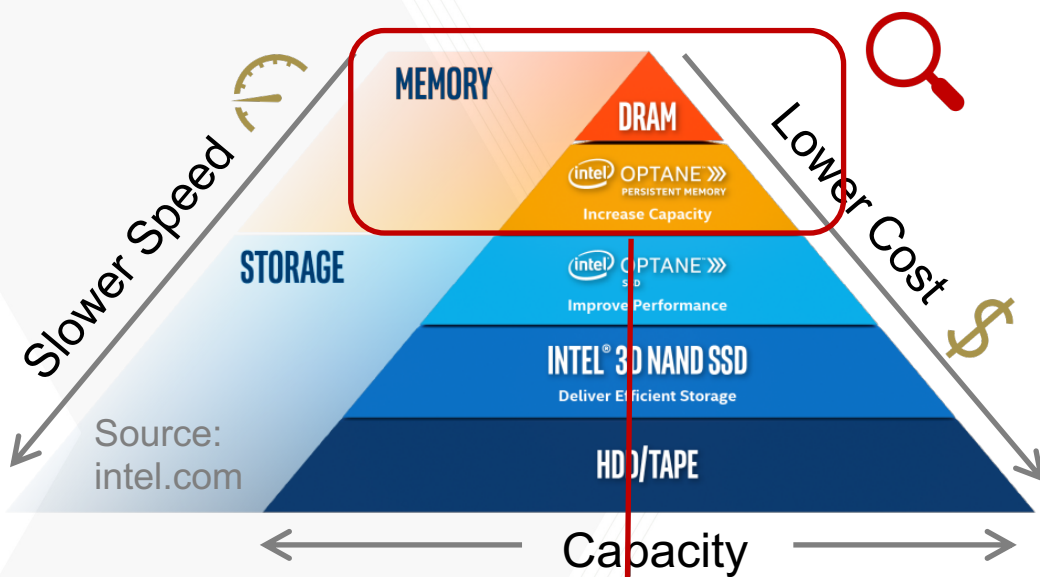
11 trillion operations per second on the Neural Engine

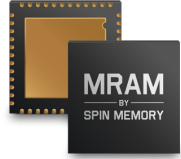



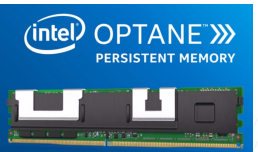
11.8 billion transistors

Improved memory compression

Secure Enclave

Heterogeneity Trade-offs



Heterogeneity Characteristic	Emerging Technology	Hardware Vendors
Low Latency	MRAM	 
High Bandwidth	HBM	 
Persistence	PMEM	

Real Impact on Real Applications

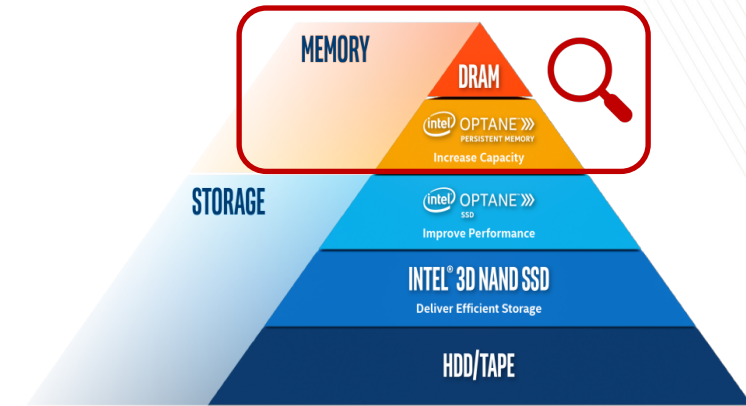
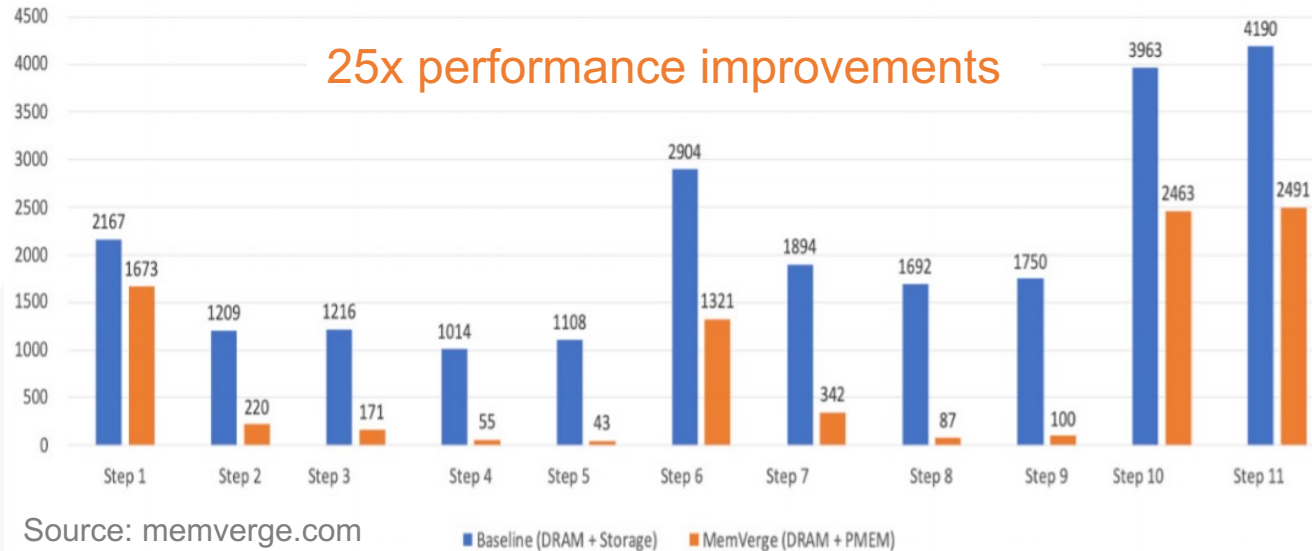
When using heterogeneous (hybrid) memories.

SOLUTION BRIEF

Big Memory Accelerates Single-Cell RNA Sequencing



Execution Time (s) of Each Analysis Stage (Compute + I/O or Snapshot)



How to boost performance?
Dynamic Data Allocations
across the hardware layers.

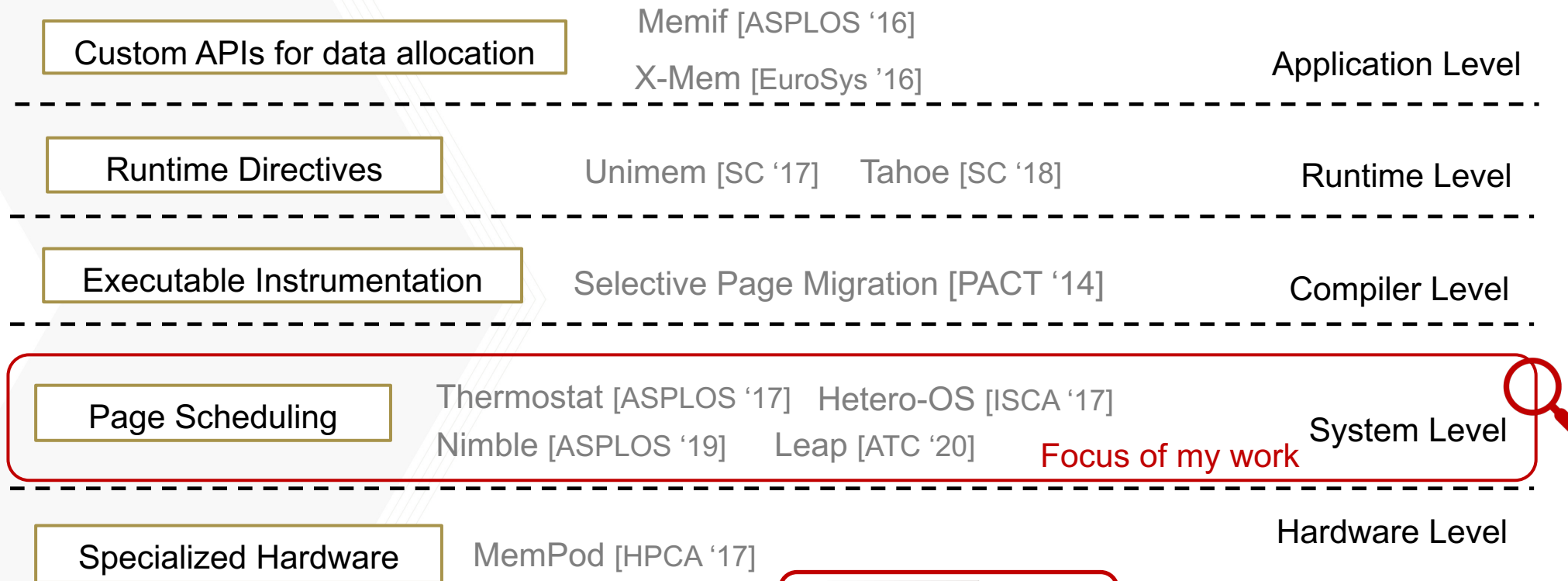
Complex decision mix:

- Which / How much / Where / When to move data?
- Capacity sizing / sharing?

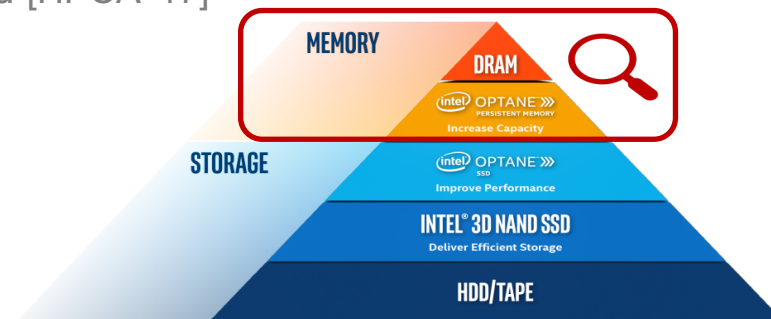


Solutions across the Software Stack

Selective Publications.



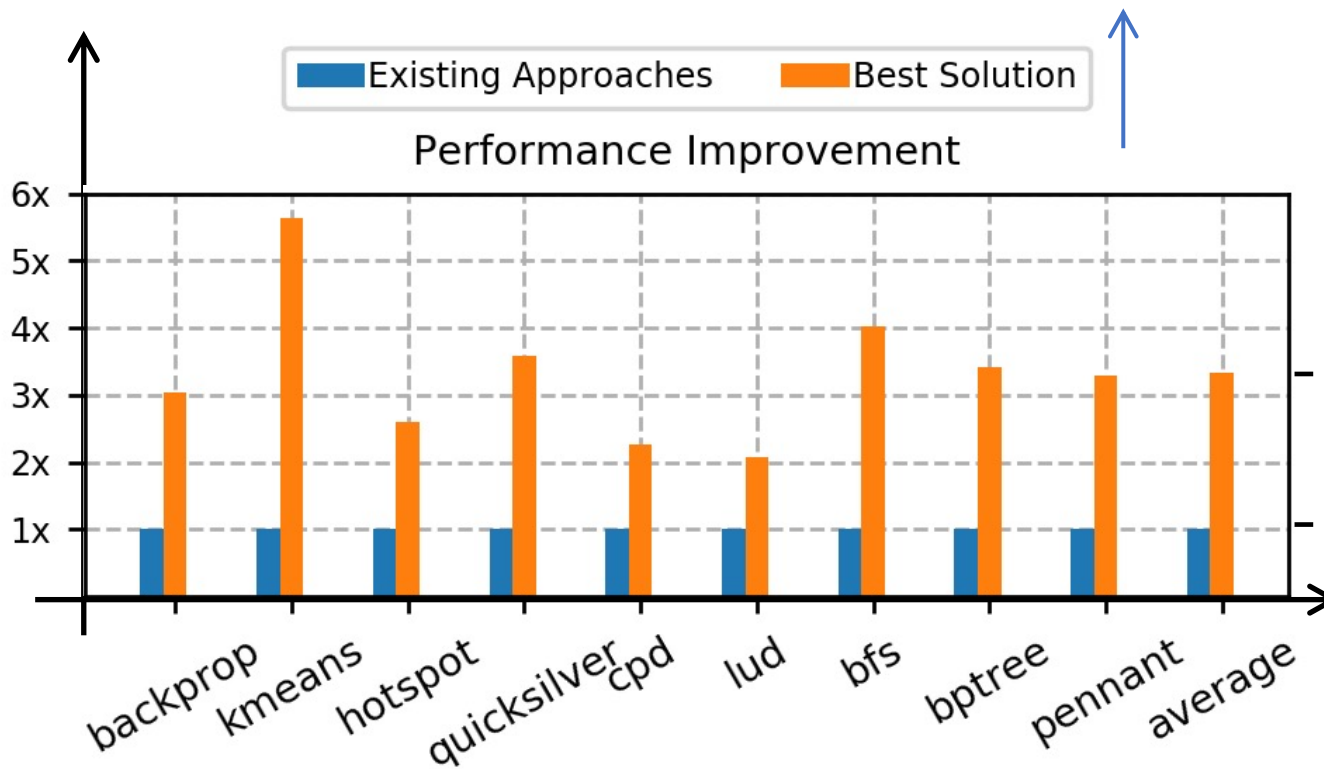
Focus of my work



Room for Performance Improvement

Left by existing approaches.

The higher, the better



- A-priori knowledge
 - Data Access Patterns
- Fine-tuned operation
 - Extensive experimentation

Best Solution

> 3x attainable improvement 

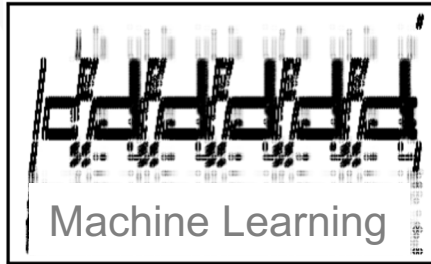
Existing Approaches

- Heuristics
 - Technology-specific
 - Application-specific
- Fixed configuration knobs
 - Empirically tuned

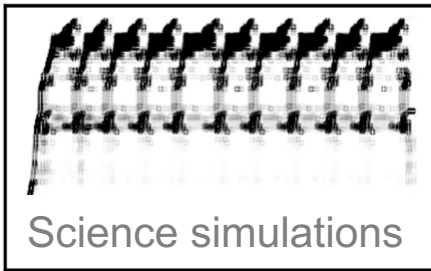
Research Contributions



Video Analytics



Machine Learning



Science simulations

Data access patterns

Applications



Resource Sharing

CoMerge - MEMSYS '17



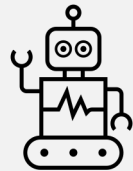
Cost Efficiency

Mnemo - HPBDC '19



Operational Frequency
Tuning

Cori – MEMSYS '20,
IPDPS '21



Practical Machine Learning (ML) Integration

Design Foundations

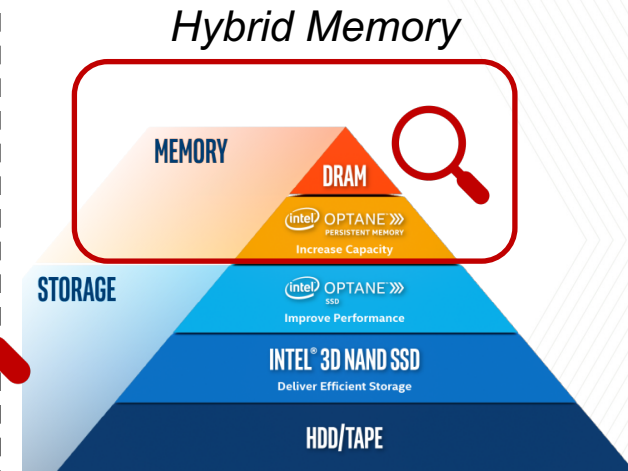
Kleio – HPDC '19

Reducing ML Overheads

Under Submission

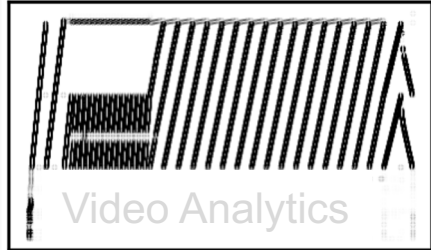
...to be continued

System-level Resource Manager

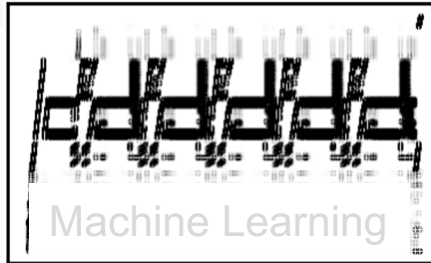


Heterogeneous Hardware

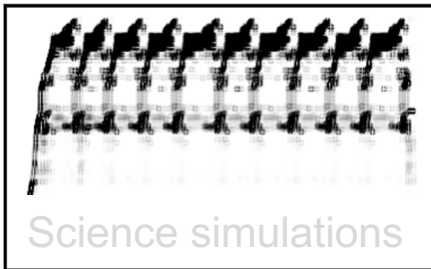
Research Highlight



Video Analytics



Machine Learning



Science simulations

Data access patterns

Applications



Resource Sharing

CoMerge - MEMSYS '17



Cost Efficiency

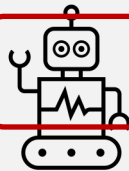
Mnemo - HPBDC '19



Operational Frequency
Tuning

Cori – MEMSYS '20,
IPDPS '21

Practical Machine Learning (ML) Integration



Design Foundations

Kleio – HPDC '19

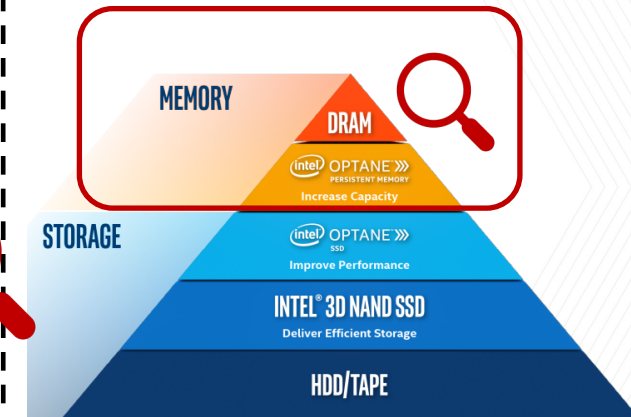
Reducing ML Overheads

Under Submission

...to be continued

System-level Resource Manager

Hybrid Memory

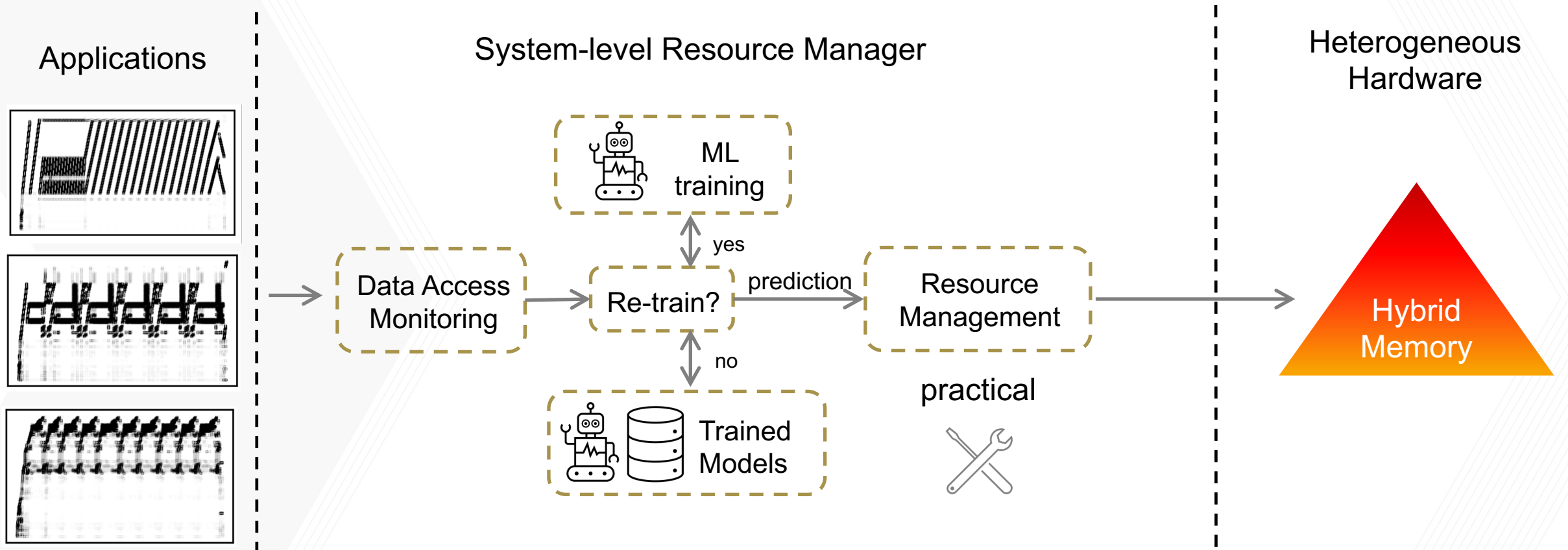


Heterogeneous Hardware

Research Contributions (Kleio)

The Goal

ML-augmented heterogeneous resource manager.

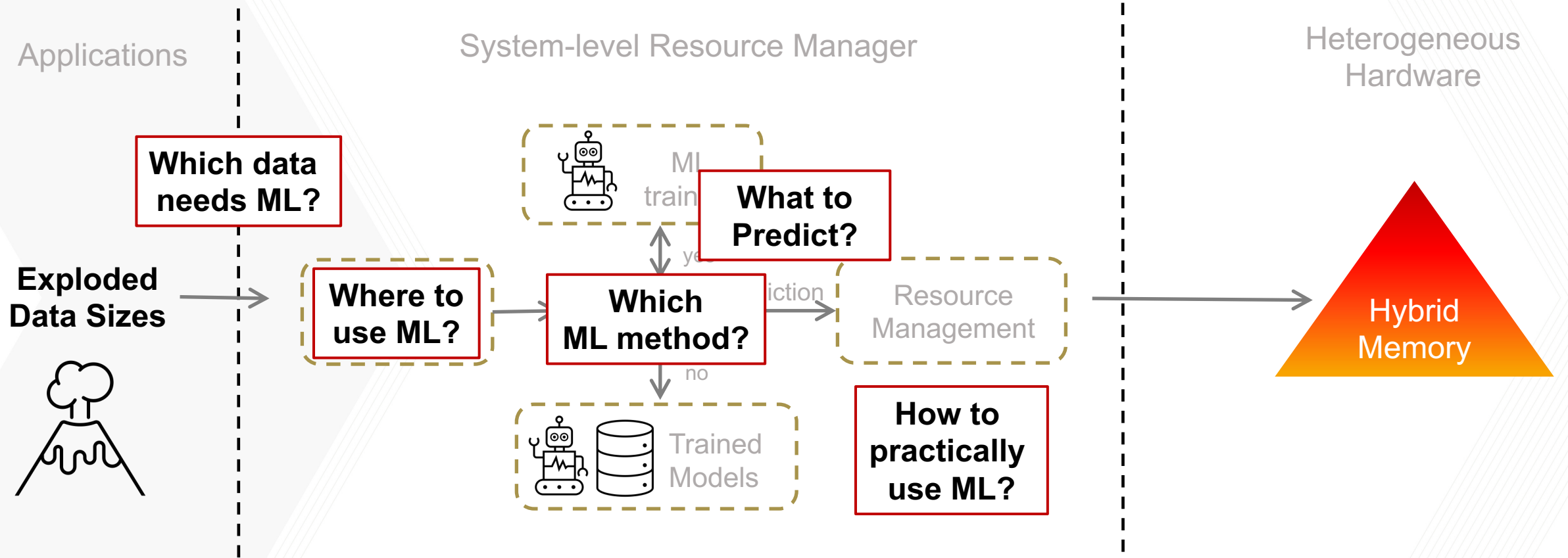


Who cares for such a solution?

Hardware Vendors (e.g., AMD).
Datacenters, Supercomputers.
System for ML.

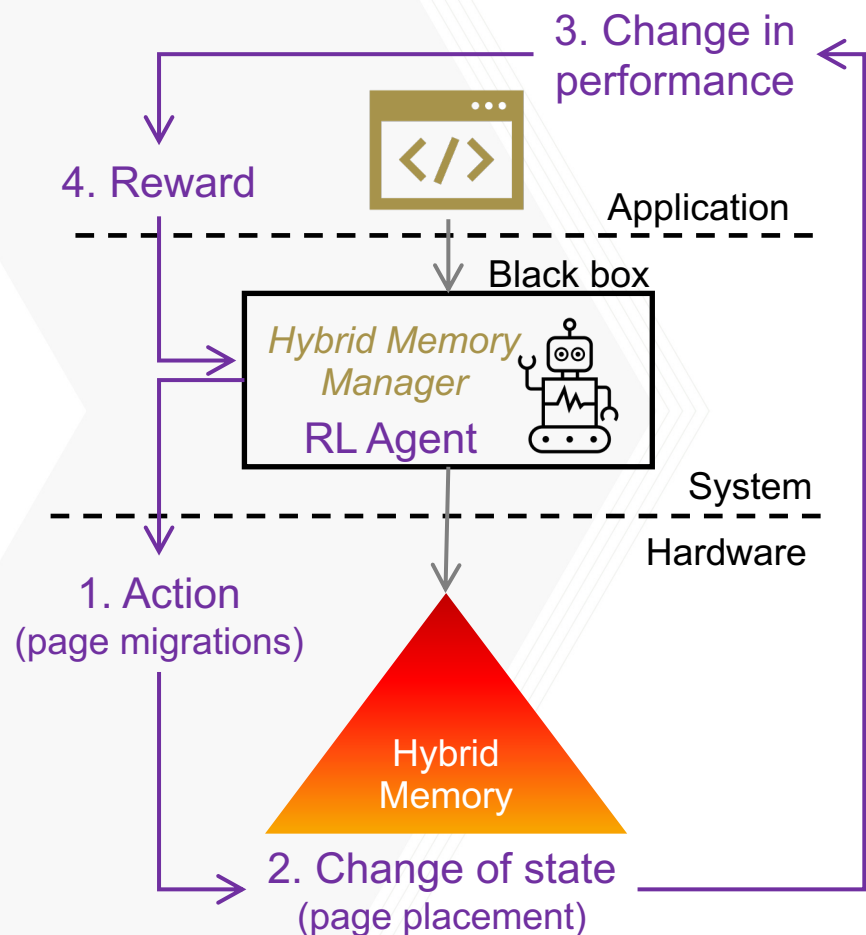
Contributions Towards the Goal

Laying the grounds for the *practical* integration of ML.



Where to use ML?

Learn which pages to move. *Replace* the memory manager with ML.



Learn the Action: Learn from moving pages across hybrid memory using **Reinforcement Learning**. Learn from mistakes (e.g., cold pages in DRAM).

Why it is not a good fit:

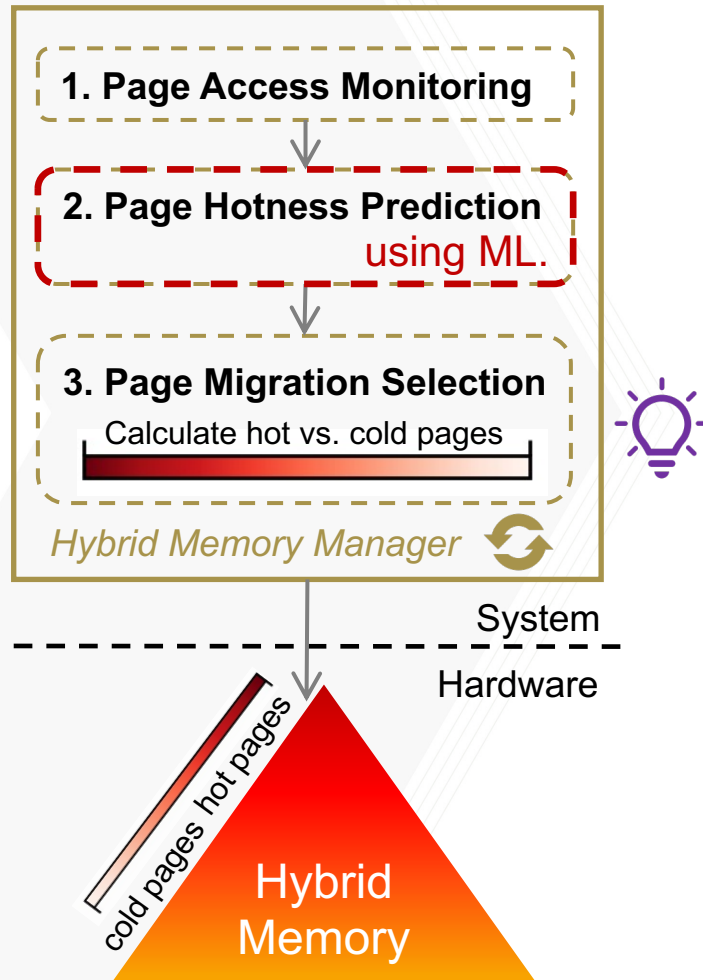
- Exponential Action Space = 2^N
- Need to re-train if configuration of hybrid memory changes.
 - Number of memory units.
 - Difference in access speeds / capacities.

Not practical / scalable.



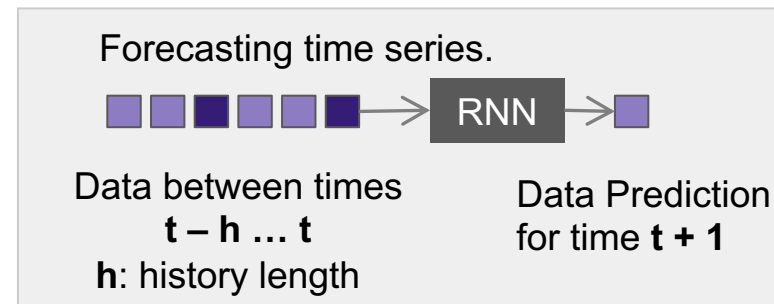
Where to use ML?

Learn which pages will be accessed in the future. *Augment* the memory manager with ML.



Learn the Behavior: Learn which pages will be accessed in the future. The manager will then move hot and cold pages appropriately.

Recurrent Neural Networks



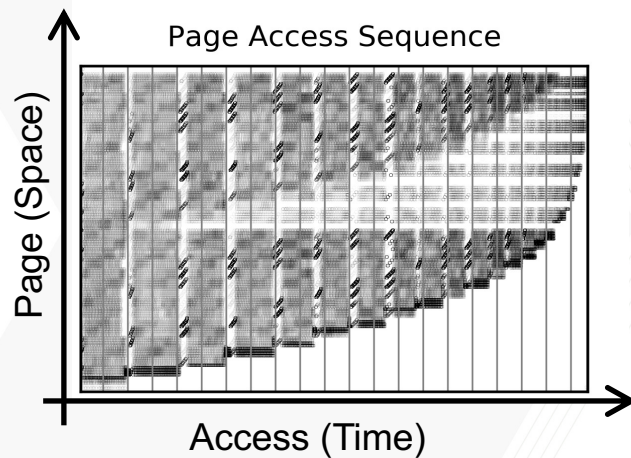
What to Predict with RNNs?

Next page accessed vs. page hotness.

Exploded
Data Sizes






Memory access trace = Time series of memory access.



Learn which page will be accessed next.

Page Access Sequence
e.g., 100, 101, 102.. → **RNN** → Next Page to be Accessed
e.g., 103

No. models	Overheads	Accuracy
1 per app 	Days to train. Months to fine-tune. 	Low. Top-k predictions not useful. 

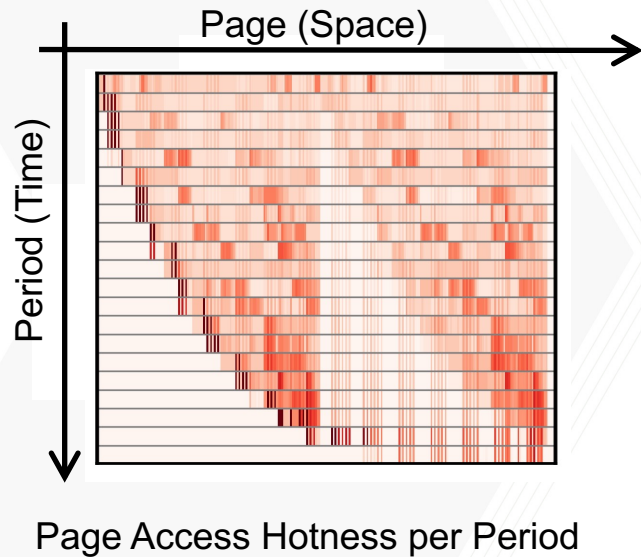
What to Predict with RNNs?

Next page accessed vs. page hotness.

Exploded
Data Sizes






Flip the view of the problem!



Learn how hot a page will be in the future.

Page Access Hotness
Across previous Periods
e.g., 100, 0, 0, 100.. → **RNN** → Page Access Hotness
in the next Period. e.g., 0

No. models	Overheads	Accuracy
Many per app 	Parallel training. Smaller models. 	High accuracy. 

How to practically use ML?

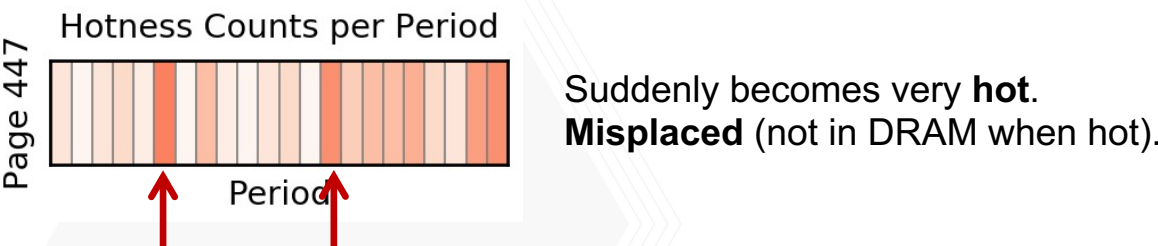
Augment existing approaches with ML.

Greek Trivia: According to the ancient Greek mythology, Kleio was the muse of history, daughter of Mnemosyne, goddess of memory.



[HPDC '19] **Kleio**: a Hybrid Memory Page Scheduler with Machine Intelligence.

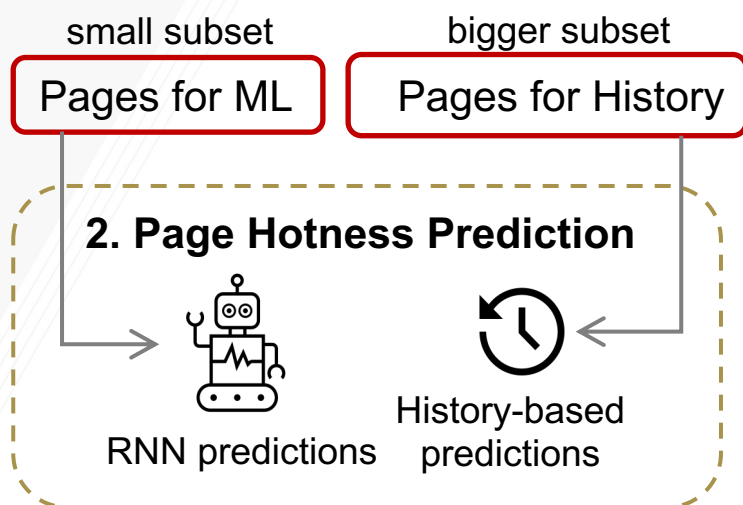
Do all pages need ML? No! Only the ones that current history-based solutions manage inefficiently.



} Small subset of such pages.

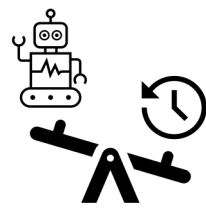


Key Idea

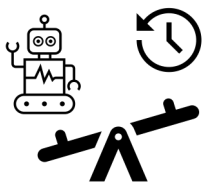


Hybrid Memory Manager Component augmented with ML.

Prediction Accuracy

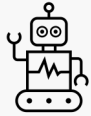


Practical Lightweight



Evaluation

Kleio delivers on average 80% of the attainable performance improvements.



100% Accurate ML
for Selected Pages



History-based
Rest of Pages

“Best” Solution



RNNs for
Selected Pages



History-based
Rest of Pages

Kleio

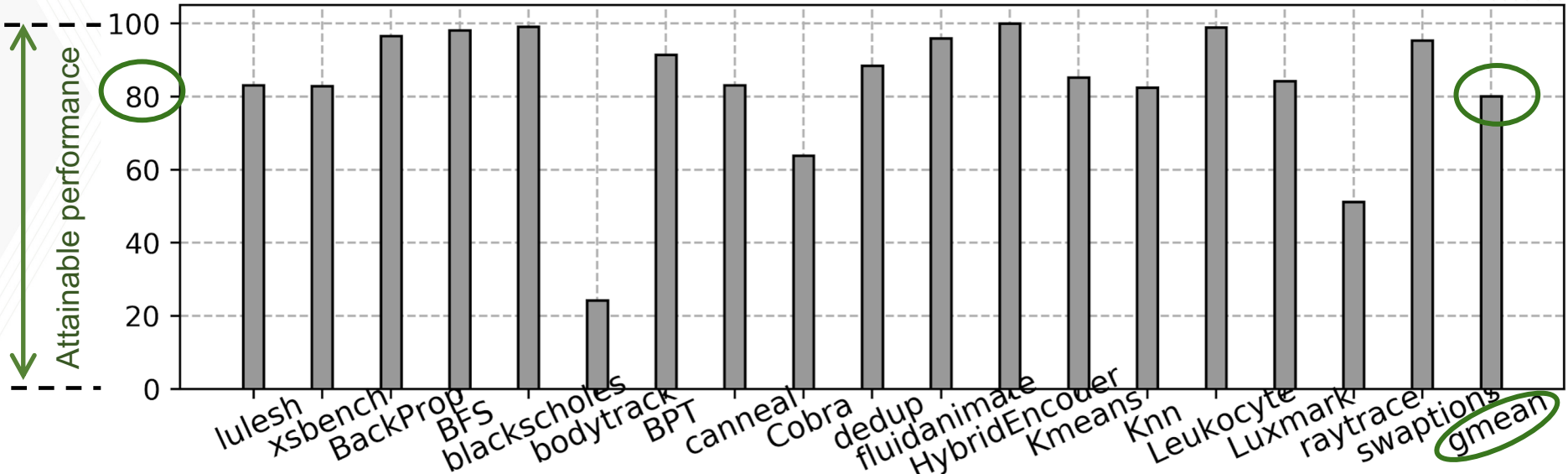
Baseline Solution



History-based
Predictions for all Pages

The higher
The better ↑

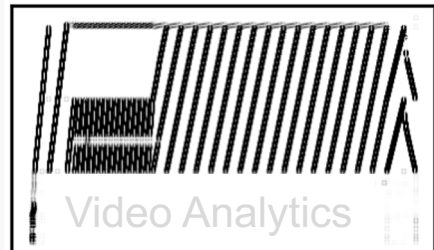
More than **95%** for **half** of the applications!



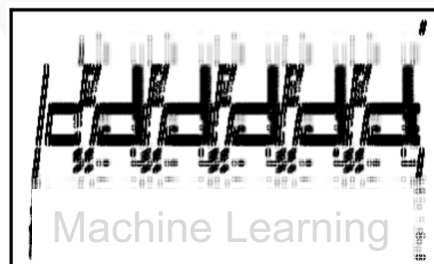
For 100 selected pages.

Research Contributions (Other)

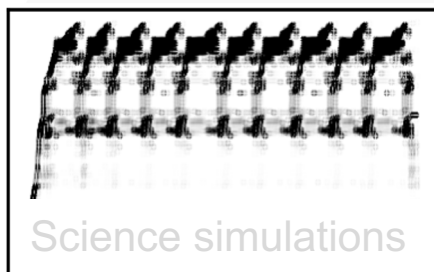
Can we do better?



Video Analytics



Machine Learning



Science simulations

Data access patterns

Applications



Resource Sharing

CoMerge - MEMSYS '17



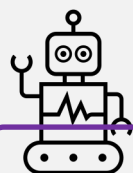
Cost Efficiency

Mnemo - HPBDC '19



Operational Frequency
Tuning

Cori – MEMSYS '20,
IPDPS '21



Practical Machine Learning (ML) Integration

Design Foundations

Kleio – HPDC '19

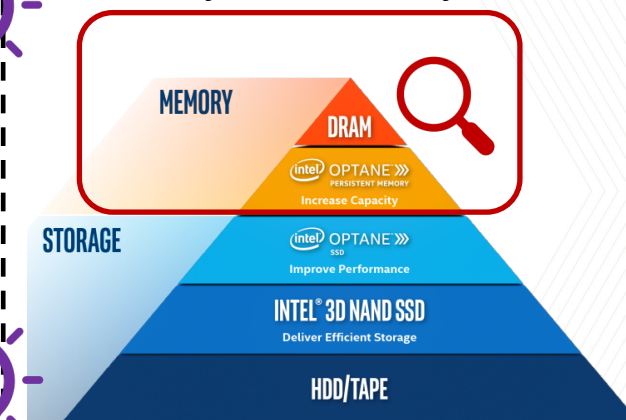
Reducing ML Overheads

Under Submission

...to be continued

System-level Resource Manager

Hybrid Memory



Heterogeneous Hardware

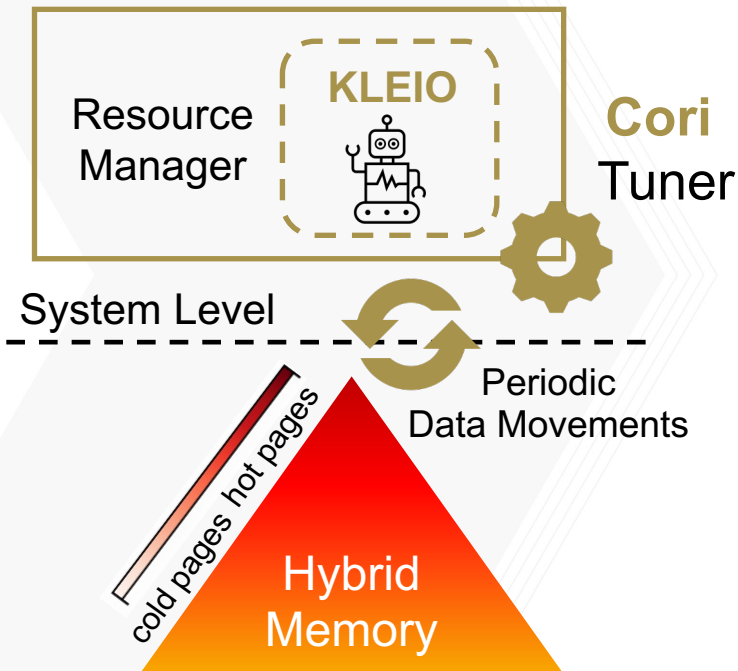
Boosting the Effects of Machine Learning

Using observation-driven insights.

Greek Trivia: According to the ancient Greek mythology, Cori (short for Terpsichore) was the muse of dance, sister of Kleio, daughter of Mnemosyne, goddess of memory.



[IPDPS '21] Cori: Dancing to the Right Beat of Periodic Data Movements over Hybrid Memory Systems.



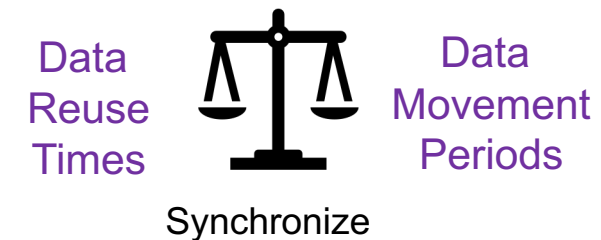
Cori tunes the frequency of data movements.



Performance boost.



Key Idea



It Is All About The Right Granularity

[Under Submission] Clustering Patterns for Practical Machine Intelligent Hybrid Memory Management.

Tuning the resource manager's
operational frequency (period) ...
... tunes the patterns for ML!

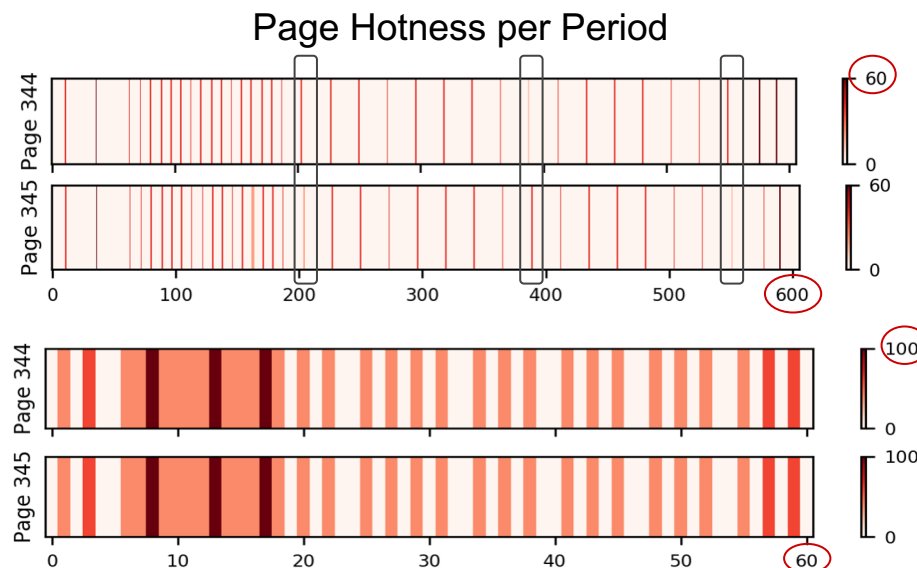



Key Idea


Group pages with *identical* patterns
under a *single* ML model.



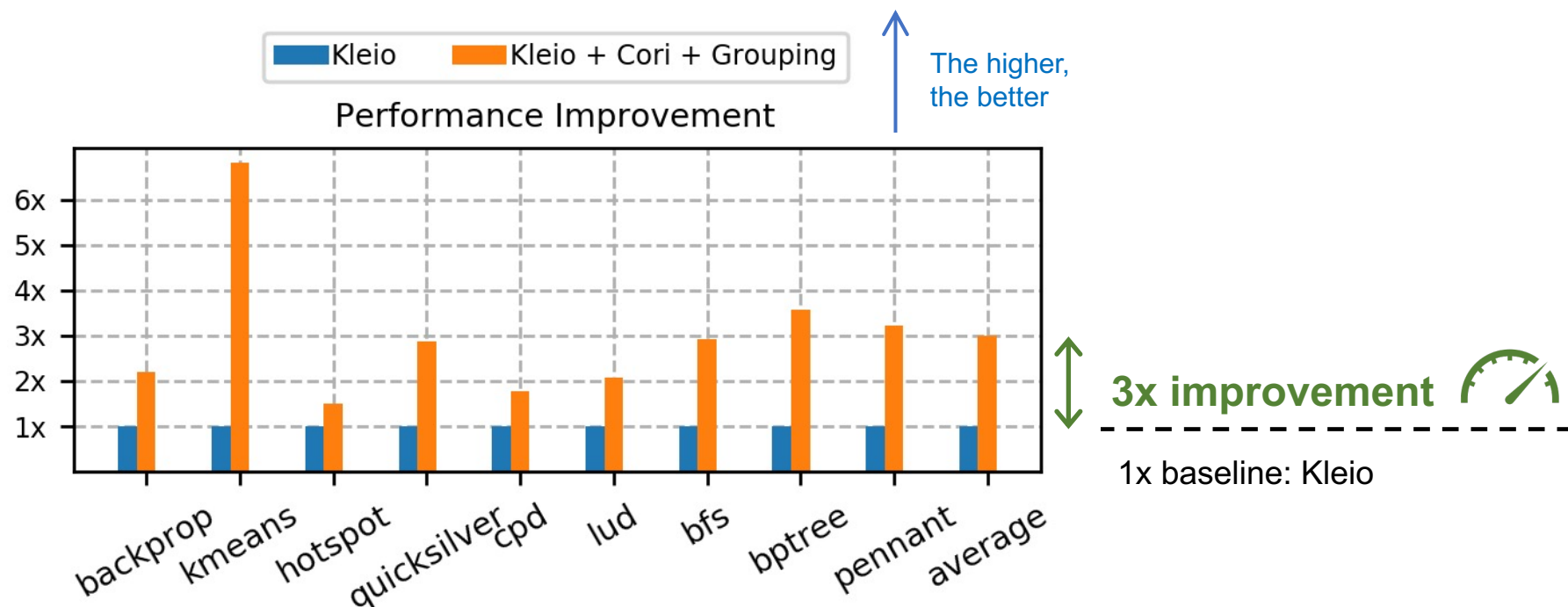
ML overheads



Kleio 
Sequences are
slightly different.

Cori 
Sequences
are **identical**.

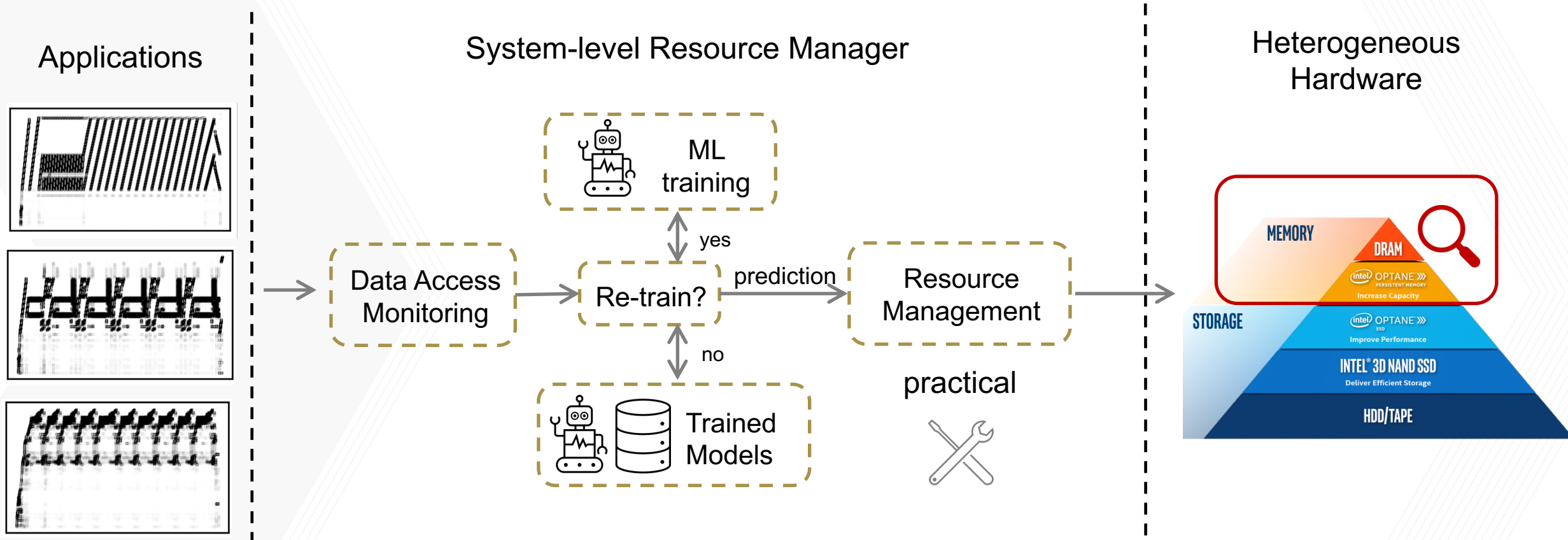
Boosting Application Performance



Fine-tuned operation further boosts the effects of ML in resource management.

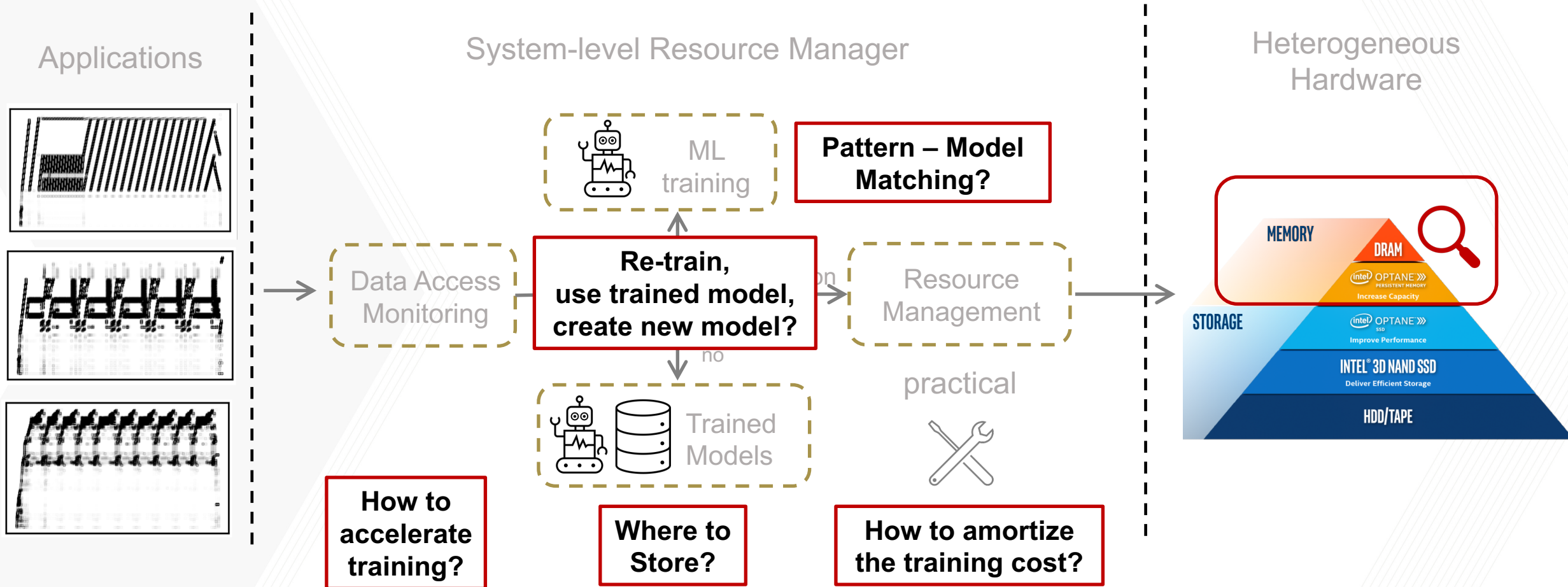
Future Research Directions

ML-augmented Heterogeneous Resource Manager



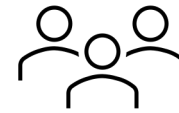
Immediate Future Contributions

Fully integrated adaptive resource manager.



Intelligent Management of Extreme Heterogeneity

Hardware configuration?
Data / Resource Management
across layers / nodes?



Users

Multi-tenancy?
Isolation?

Performance?
Cost / Energy /
Resource Efficiency?



High-Speed Interconnects

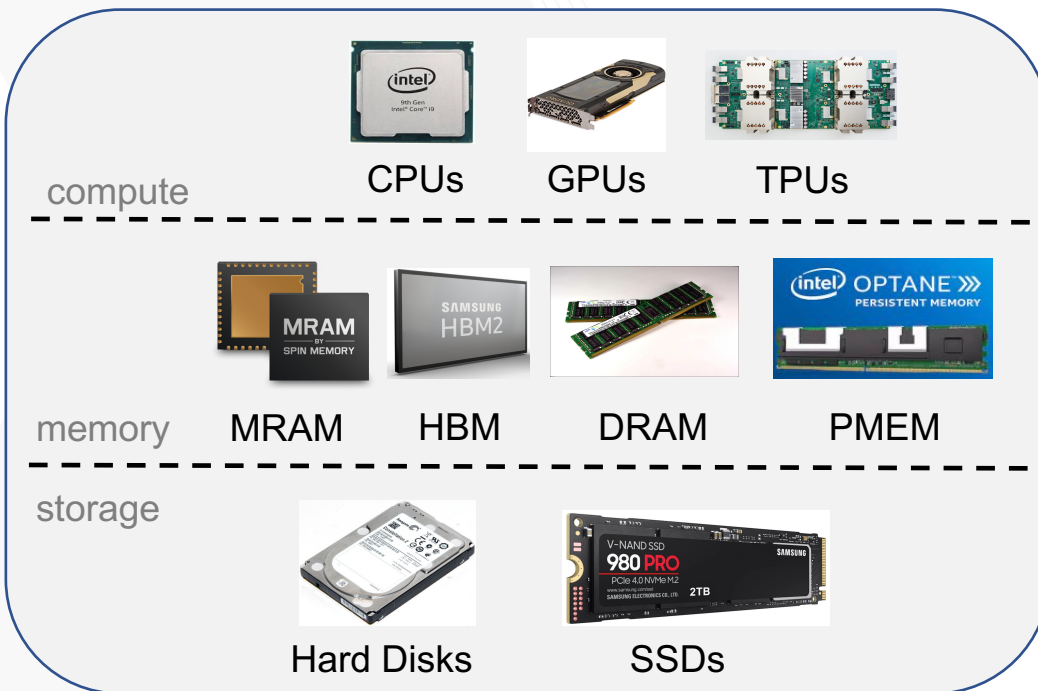


Datacenter



Supercomputer

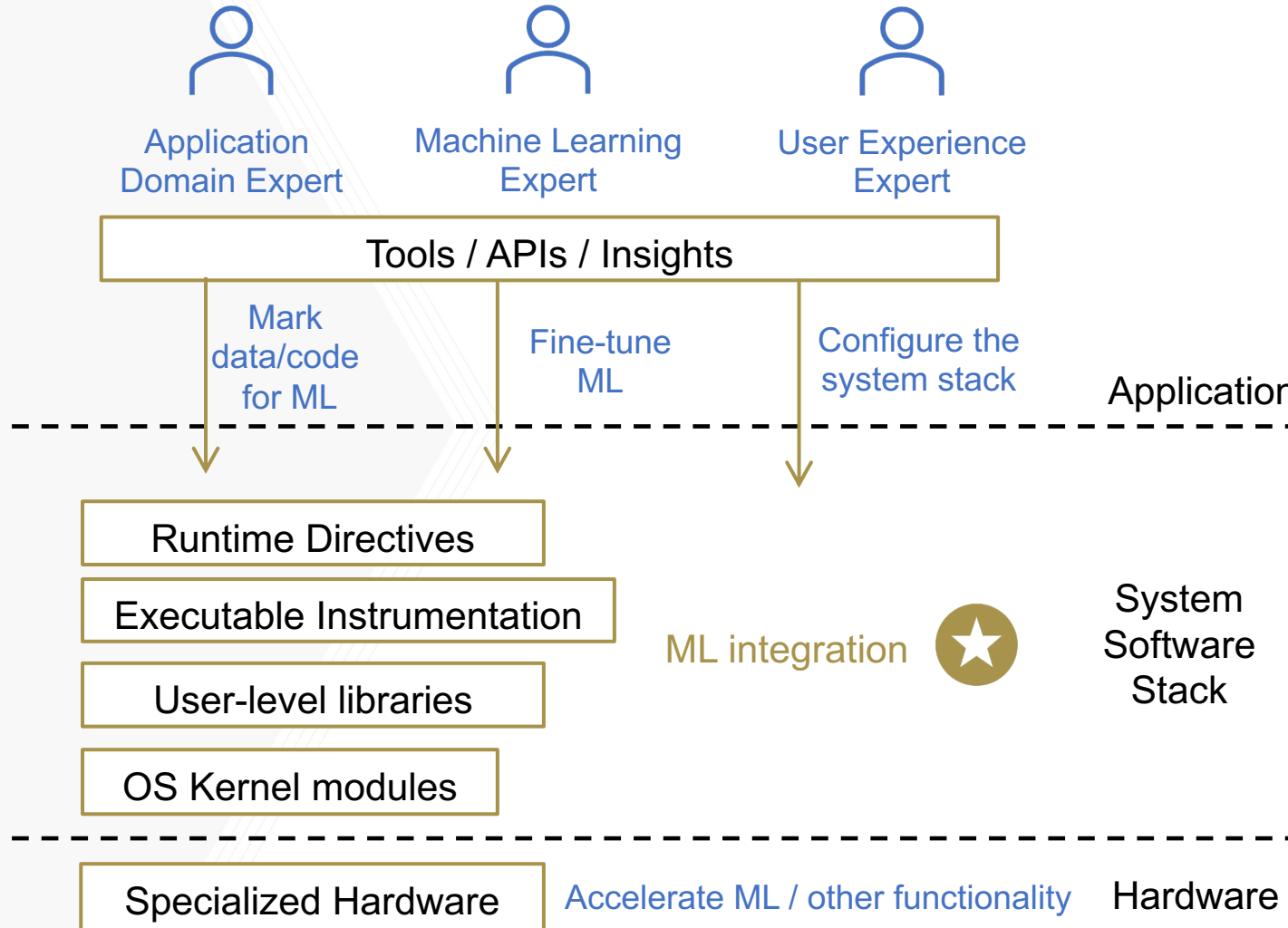
Massive Node Clusters
Disaggregated Resources



Local Node

ML integration Aspects:
Necessity Effectiveness Practicality Interpretability

Cross-Stack Synergies for ML integration



Summary

Greek Trivia: According to the ancient Greek mythology, **Thaleia** was the muse of comedy, daughter of Mnemosyne, goddess of memory.



THALIA.



CLIO.



TERPSICHORE.



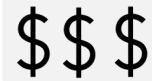
Mnemosyne

First-author Publications



Resource Sharing

CoMerge - MEMSYS '17



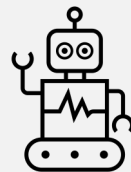
Cost Efficiency

Mnemo - HPBDC '19



Operational Frequency Tuning

Cori – MEMSYS '20,
IPDPS '21



Practical Machine Learning (ML) Integration

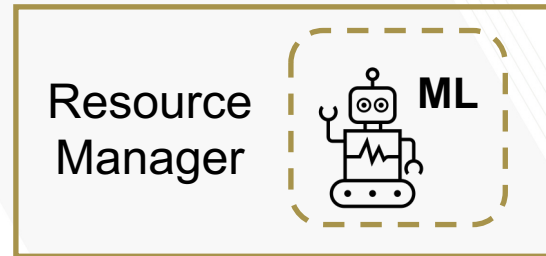
Design Foundations

Kleio – HPDC '19

Reducing ML Overheads

Under Submission

...to be continued



System Level

Heterogeneous Hardware

