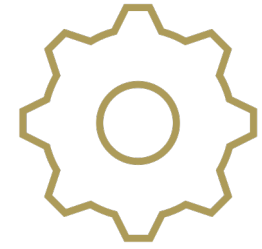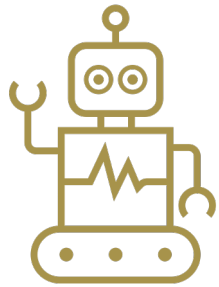Smart

Fast

Systems

Machine Learning

Computer Vision

# Building Smart and Fast Systems using Machine Learning and Computer Vision.

Thaleia Dimitra Doudali

Assistant Research Professor @IMDEA Software Institute

# About Me

2015

2021

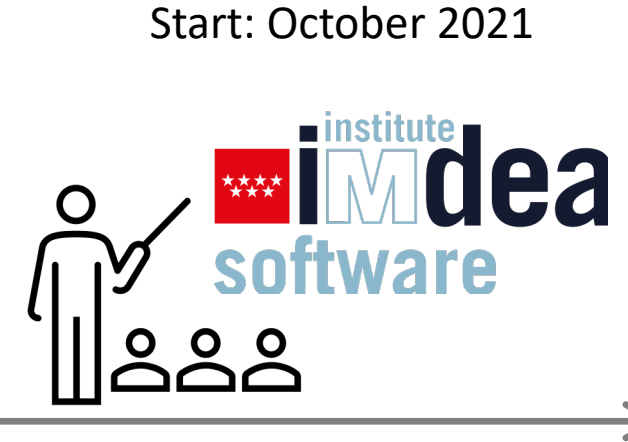Start: October 2021

Born and raised
in Greece.

Undergrad in ECE at
NTUA, Athens, Greece.
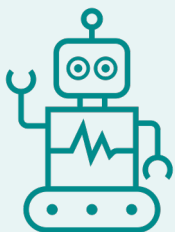
PhD in CS at
Georgia Tech, Atlanta, USA.

Advised by Ada Gavrilovska.

Assistant Professor at
IMDEA, Madrid, Spain.

# About My Research

My research lies at the intersection of Machine Learning and Systems.



Machine Learning (**ML**)

**ML _for_ Systems**

e.g., RNNs for system-level pattern prediction.

**Systems _for_ ML**

e.g., design new systems to optimize ML workloads.

Operating Systems (**OS**) Software



Computer Vision (**CV**)

**ML + CV _for_ Systems**

e.g., image-based ML for pattern recognition and prediction.
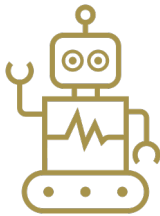
# Talk Outline

**Why do we need Smarter and Faster Systems?**
The evolution of the hardware technologies, calls for software improvements.

**Building *Smart* Systems**
Using machine and human intelligence to build practical ML-based systems.

**Building *Fast* Systems**
Reducing ML-based management overheads with visualization.
Building image-based system pipelines.

**Future Research Directions**

# Talk Outline

**Why do we need Smarter and Faster Systems?**
The evolution of the hardware technologies, calls for software improvements.

**Building *Smart* Systems**
Using machine and human intelligence to build practical ML-based systems.

**Building *Fast* Systems**
Reducing ML-based management overheads with visualization.
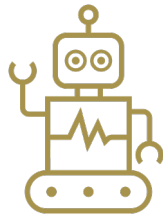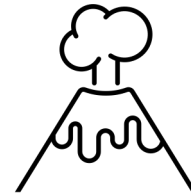Building image-based system pipelines.

**Future Research Directions**

# The Era of Data

"More than **65 ZB** of data will be created, captured, copied, and consumed in the world this year."

Source: International Data Corporation, March 2021.

**Exploded Data Sizes**

Scientific Simulations

Big Data

Artificial Intelligence

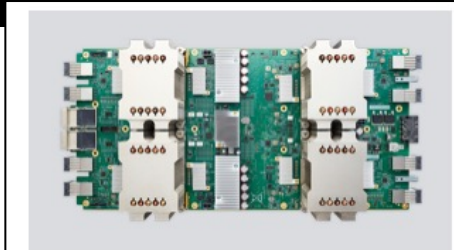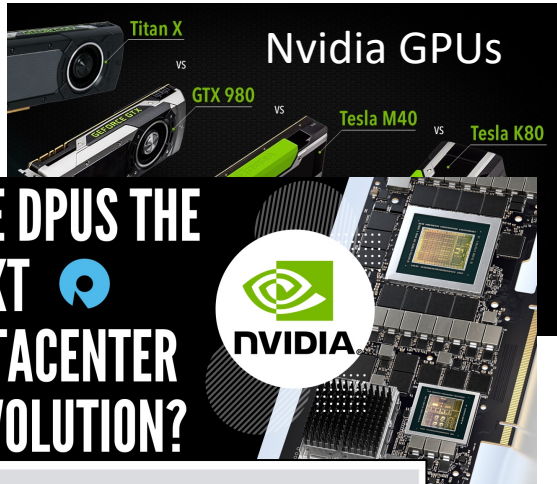Video Analytics

*Data Analytics Pipeline*

**ZBs of data**  →  Capture → Process → Store → Analyze → Use

Need for speed and massive storage capacities!

# The Era of Heterogeneous Hardware

## Compute Acceleration

Nvidia GPUs

Titan X
GTX 980
Tesla M40
Tesla K80

**ARE DPUS THE NEXT DATACENTER REVOLUTION?**

NVIDIA

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

Google

## Data Storage Acceleration

intel OPTANE »»»
PERSISTENT MEMORY

AMD

HIGH BANDWIDTH MEMORY

AMD

V-NAND SSD
980 PRO
PCIe 4.0 NVMe M.2
www.samsung.com/ssd
SAMSUNG ELECTRONICS CO., LTD.
SAMSUNG
2TB

## Network Acceleration

**Mellanox Innova™-2 Flex Open Programmable SmartNIC**

†

Mellanox
Innova™ 2 FLEX

## Interconnection Standards

CXL Compute Express Link ™

Gen-Z Consortium                    GEN Z

• Industry Leaders developing a memory-semantic interconnect

AMD    ARM    BROADCOM    CAVIUM    CRAY

DELL EMC    Hewlett Packard Enterprise    HUAWEI    IBM    IDT

Lenovo    Mellanox TECHNOLOGIES    Micron    Microsemi    redhat.

SAMSUNG    SEAGATE    SK hynix    WD Western Digital    XILINX.

10/11/2016                    © Gen-Z Consortium 2016                    2

# Heterogeneity Across Computing Platforms

**Supercomputers**

Exascale Era



**Datacenters**

Available first on Google Cloud: Intel Optane DC Persistent Memory

A2 VMs now GA—the largest GPU cloud instances with NVIDIA A100 GPUs
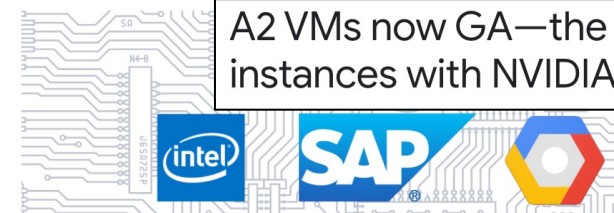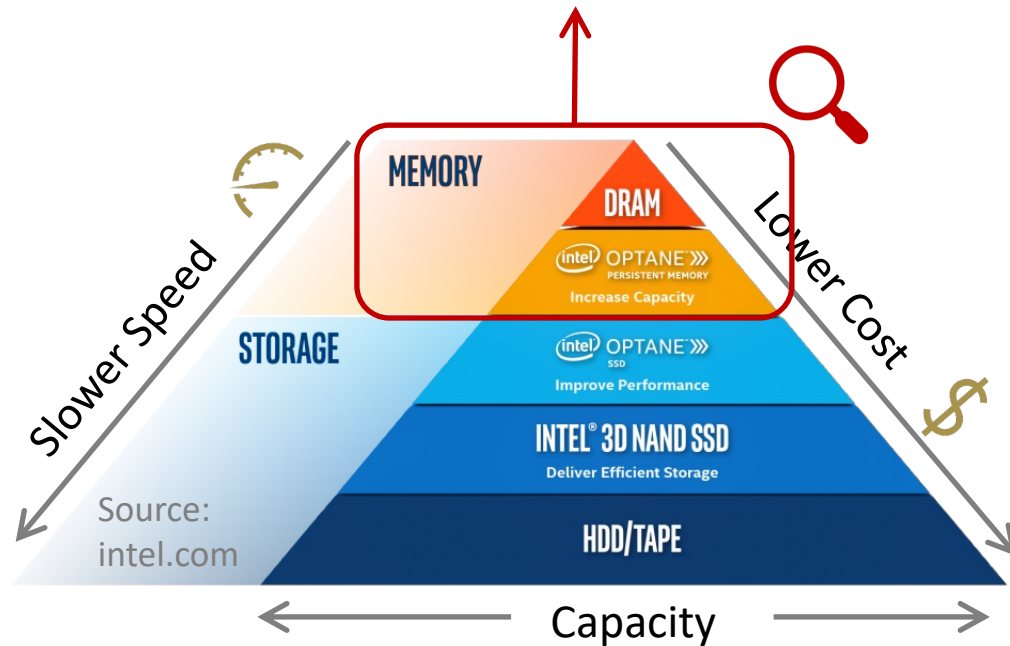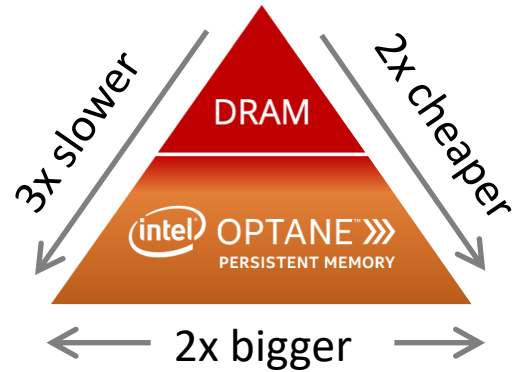
| Application Performance | 200 PF |
|---|---|
| Number of Nodes | 4,608 |
| Node performance | 42 TF |
| Memory per Node | 512 GB DDR4 + 96 GB HBM2 |
| NV memory per Node | 1600 GB |
| Total System Memory | >10 PB DDR4 + HBM2 + Non-volatile |
| Processors | 2 IBM POWER9™ 9,216 CPUs 6 NVIDIA Volta™ 27,648 GPUs |
| File System | 250 PB, 2.5 TB/s, GPFS™ |
| Power Consumption | 13 MW |
| Interconnect | Mellanox EDR 100G InfiniBand |
| Operating System | Red Hat Enterprise Linux (RHEL) version 7.4 |

**Personal Devices**

3x slower

2x cheaper

DRAM

(intel) OPTANE »»
PERSISTENT MEMORY

2x bigger

Source: memverge.com

Slower Speed

Lower Cost $

MEMORY

DRAM

(intel) OPTANE »»
PERSISTENT MEMORY
Increase Capacity

STORAGE

(intel) OPTANE »»
SSD
Improve Performance

INTEL® 3D NAND SSD
Deliver Efficient Storage

HDD/TAPE

Source: intel.com

Capacity

| Characteristic | Technology | Vendors |
|---|---|---|
| Low Latency | MRAM | |
| Uniform Latency | DRAM | |
| High Bandwidth | HBM | |
| Persistent / Non Volatile | PMEM / NVM | |

We are in the era of **Hybrid Memory** Systems.
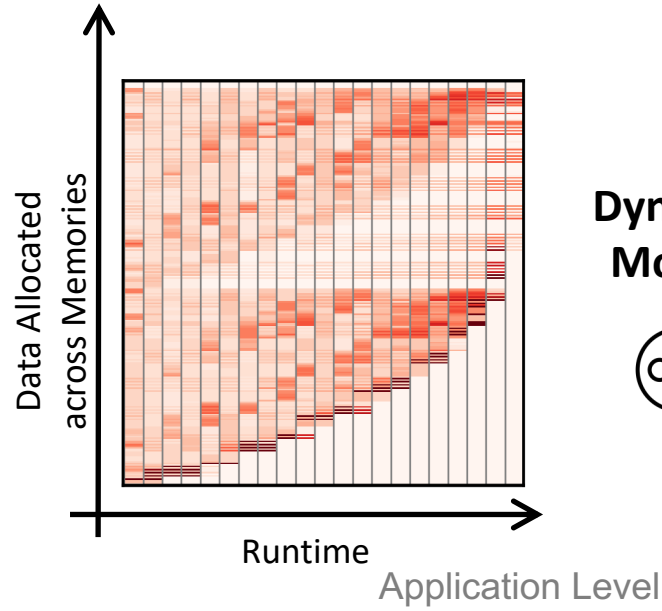A mix of different technologies at different speeds / capacities / costs.

The OS should move pages dynamically across hybrid memory to maximize the efficiency.

Time Period

Page

Runtime

*hot* page

= # Accesses of a Page during a Time Period

*cold* page

Periodically moves pages

**Resource Manager**

Operating System (OS)

Memory Hardware (HW)
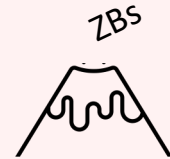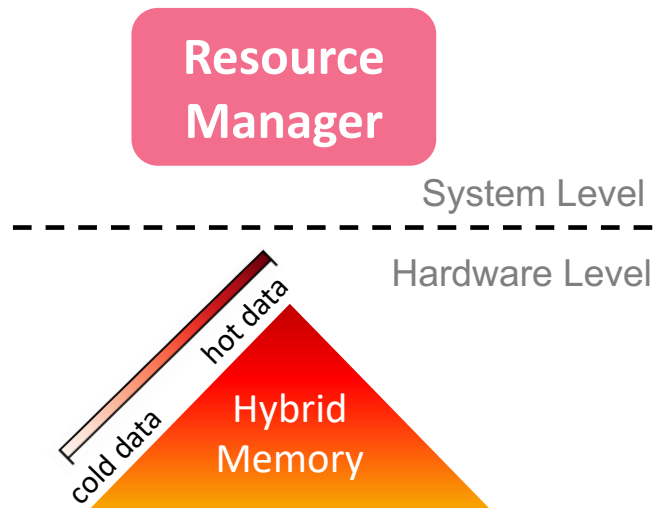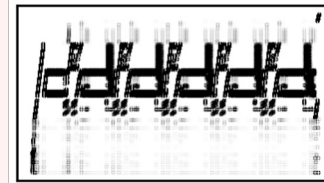
hot pages

cold pages

Hybrid Memory

**Dynamic Data Movements!**

It is a **complex decision mix** to manage the data allocated across memories.

E.g., Which / How much / Where / When to move data?

Runtime

Data Allocated across Memories

Application Level

**Resource Manager**

System Level

Hardware Level

hot data

cold data

Hybrid Memory

Why do we need smarter and faster systems?

ZBs

Application data sizes

Complex data access patterns

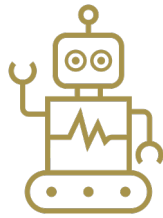Exploded system parameter space

performance / practicality

Hard to balance

# Talk Outline

**Why do we need Smarter and Faster Systems?**
The evolution of the hardware technologies, calls for software improvements.

**Building *Smart* Systems**
Using machine and human intelligence to build practical ML-based systems.
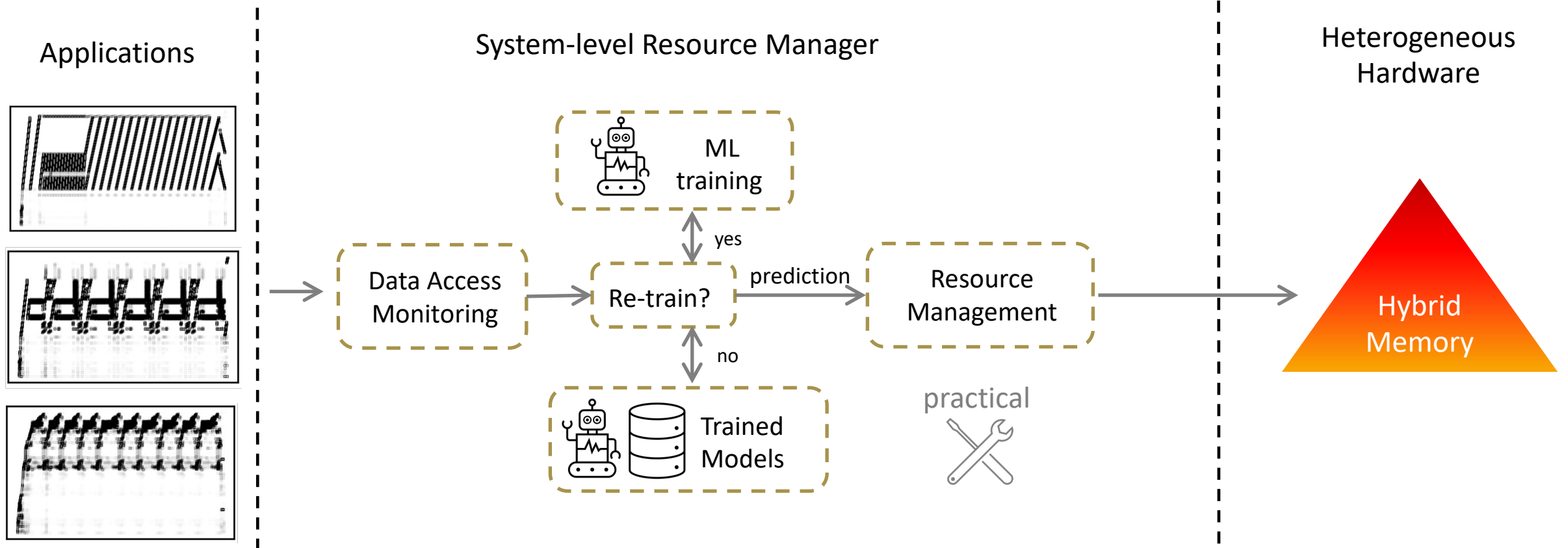
**Building *Fast* Systems**
Reducing ML-based management overheads with visualization.
Building image-based system pipelines.
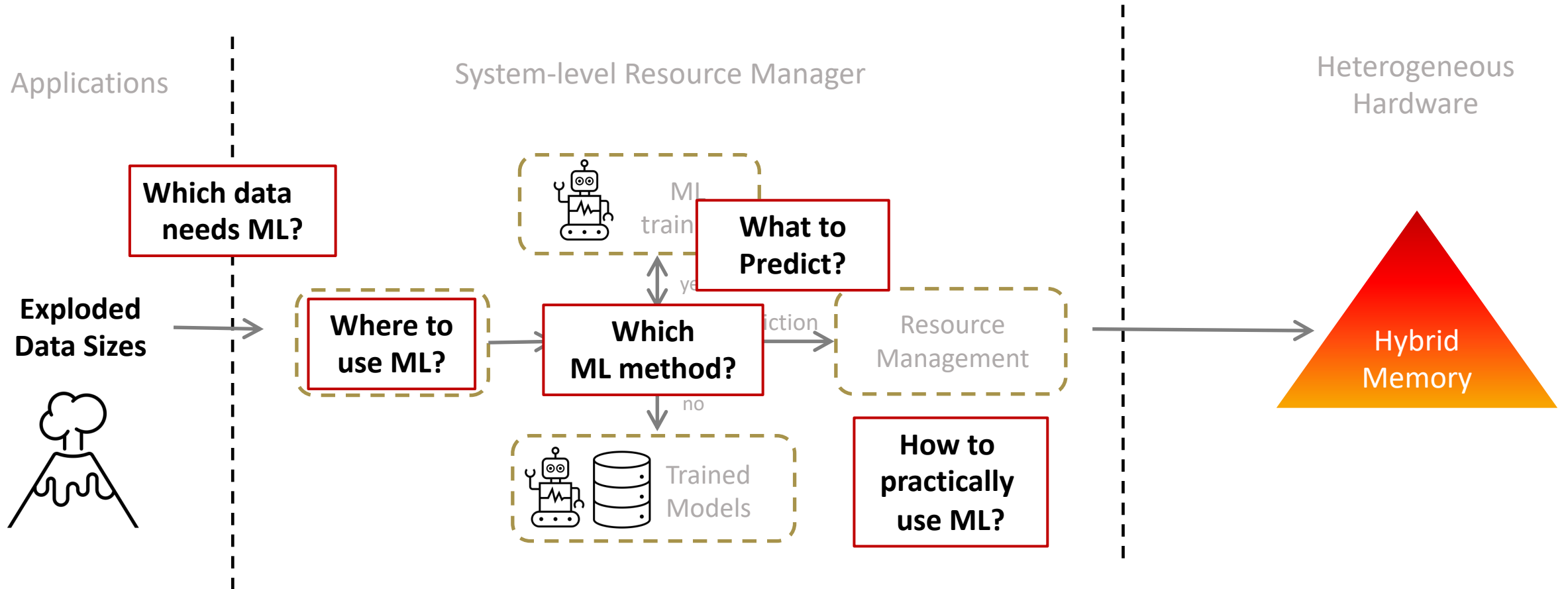
**Future Research Directions**

Applications

System-level Resource Manager

Heterogeneous Hardware

ML training

Data Access Monitoring

Re-train?

yes

no

prediction

Resource Management

practical

Trained Models

Hybrid Memory

Applications

System-level Resource Manager

Heterogeneous Hardware

**Which data needs ML?**

**Exploded Data Sizes**

**Where to use ML?**

ML training

**What to Predict?**

yes

**Which ML method?**

prediction

Resource Management

Hybrid Memory

no

Trained Models

**How to practically use ML?**

# System design of Kleio

**Kleio:** a hybrid memory page scheduler with machine intelligence. [HPDC 2019]

Applications

System-level Resource Manager

Heterogeneous Hardware

Page Selector

1. Page Access Monitoring

Page Hotness

small subset

Pages for ML

bigger subset

Pages for History

**2. Page Hotness Prediction**

ML-based predictions
(**Per page RNN models**)

History-based predictions

Page Hotness Prediction

Page Migrations

Hybrid Memory

cold pages  hot pages

3. Page Migration Selection

Calculate hot vs. cold pages

*Not all pages "need" ML.*

**Result:** Kleio bridges 80% of the performance gap between existing and oracular solutions.

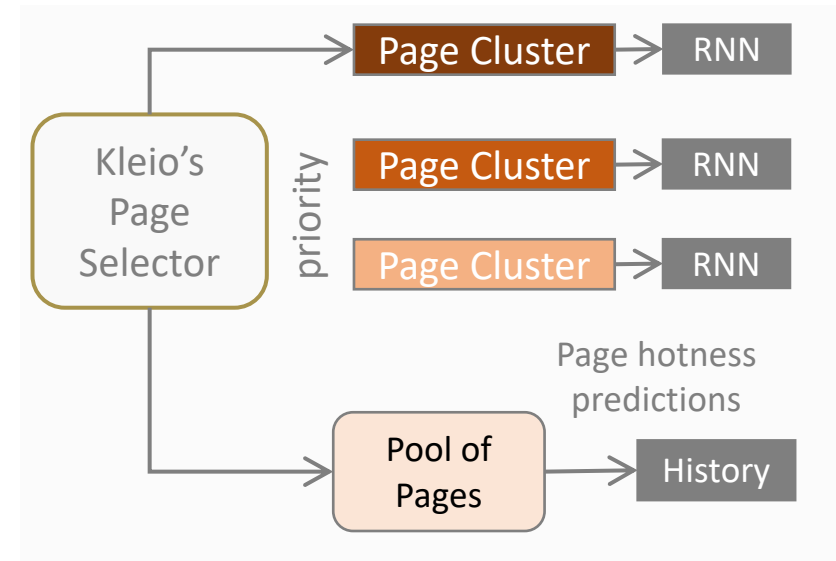Apply ML **when** and **where** necessary.



Apply ML on a small page subset.

↳ Foundations for practical use of ML.

Carefully select pages for ML.
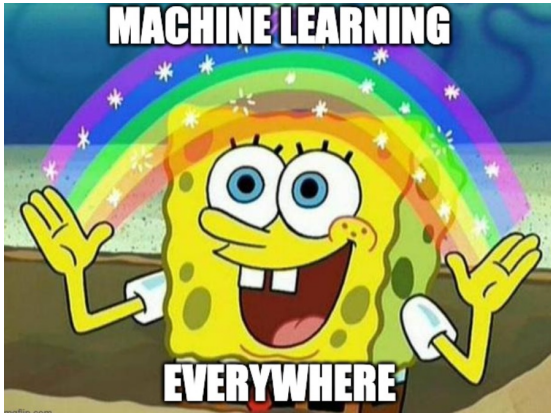
↳ Application performance boost.

Small can still mean thousands of pages, because of the massive memory footprints of modern workloads.

**Can we reduce the number of pages via clustering?**

**Coeus:** Clustering (A)like Patterns for Practical Machine Intelligent Hybrid Memory Management . [CCGrid 2022]

Page Hotness per Period



**Clustering?** Let's use ML!

For example, K-means.
- How many clusters?
- Clustered input to ML?

Not trivial to configure.

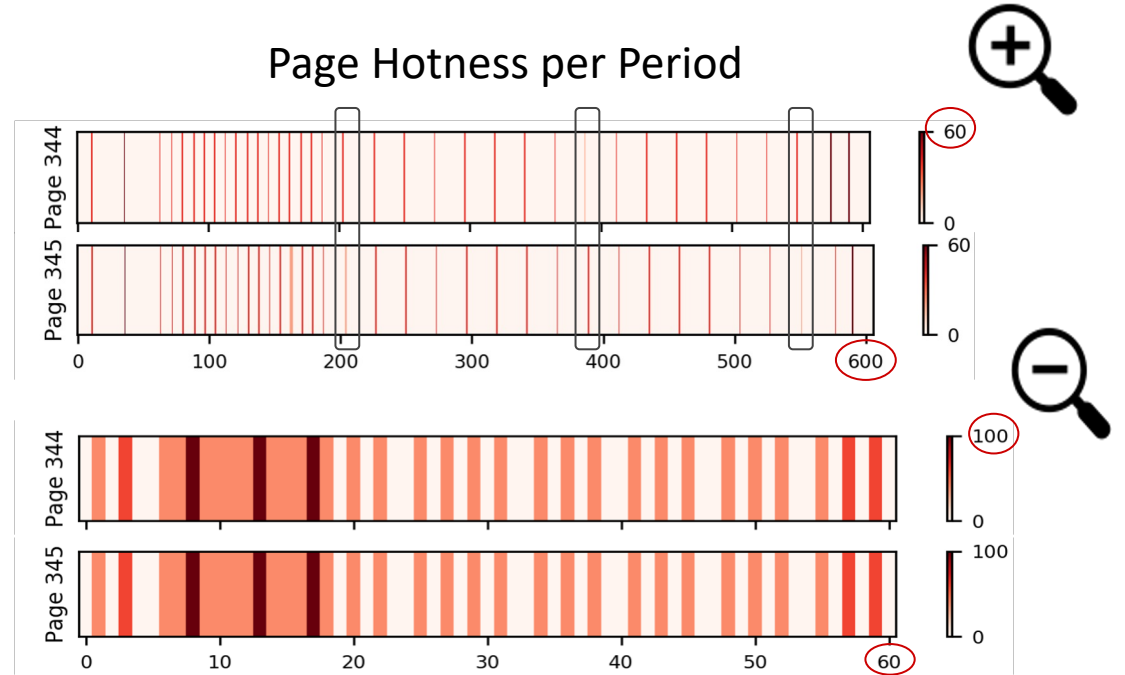Let's use our human intelligence..

.. Kleio learns the patterns of page hotness across time periods.

So what if I increase the duration of the period?

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

Key Idea

3x less RNNs

3x more performance

# Talk Outline

Apply ML **when** and **where** necessary.



Apply ML on a small page subset.

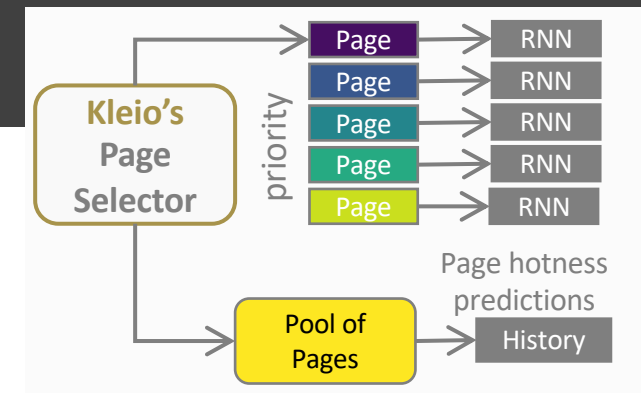↳ Foundations for practical use of ML.

Carefully select pages for ML.
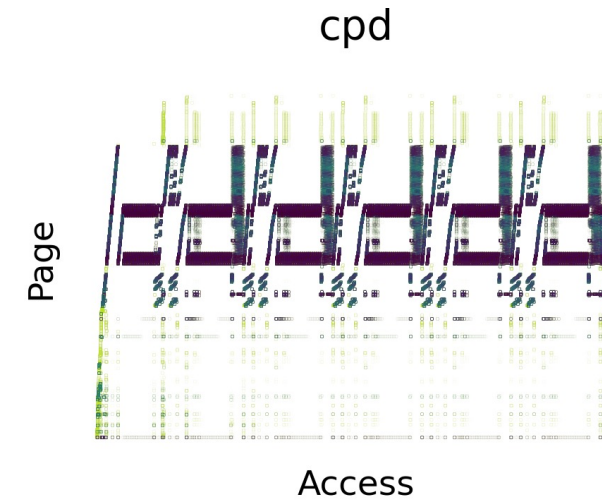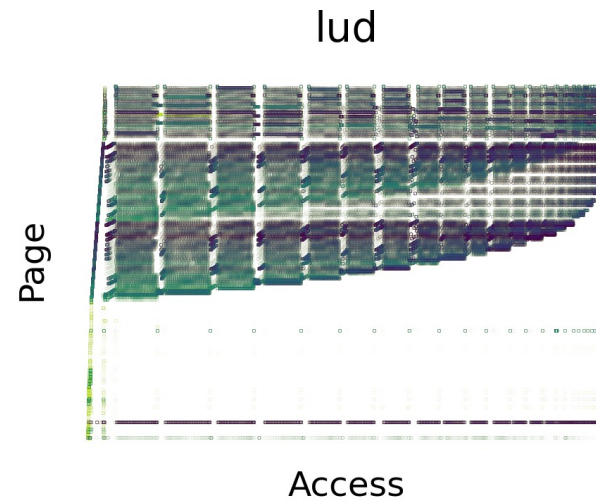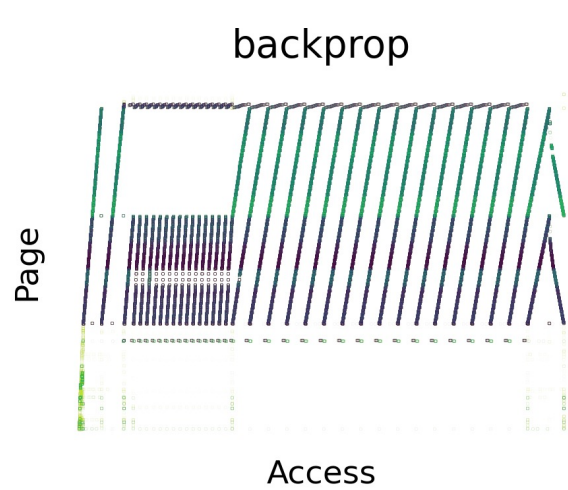
↳ Application performance boost.

**The page selection is not a lightweight process.**
Performance modeling and estimations are used to maximize the effects of ML on application performance.

**Can we accelerate the page selection process?**

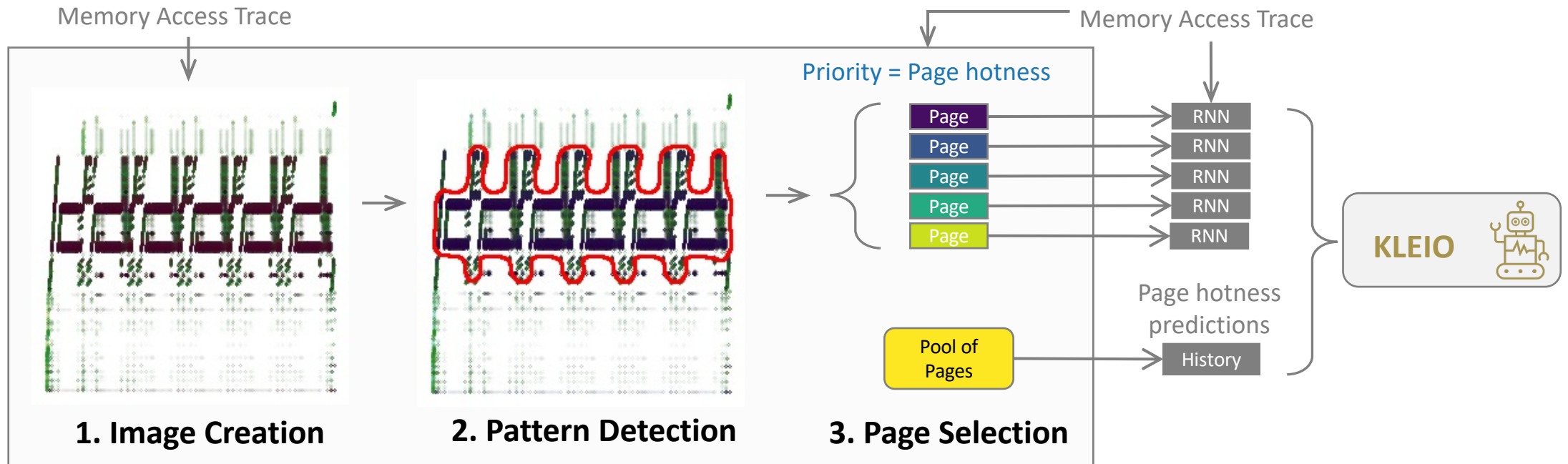High Priority — Low Priority

backprop

lud

cpd

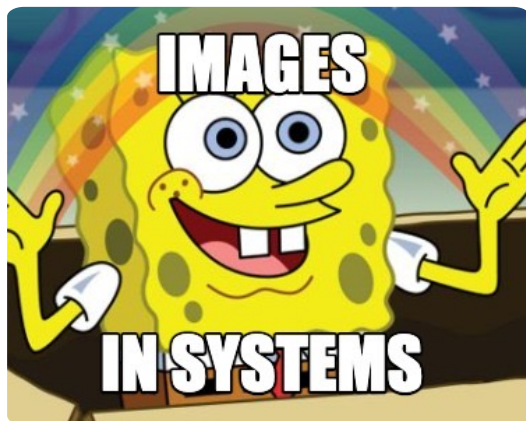*Neighboring* pages that are part of distinct access patterns across *time* receive similar priority for ML.

# Towards Image-based Page Selection

**Cronus:** Computer Vision-based Machine Intelligent Hybrid Memory Management. [MEMSYS 2022]



**1. Image Creation**   **2. Pattern Detection**   **3. Page Selection**

Cronus reduces by **400x** the page selection times, from minutes down to seconds.

**Creating images helps:**

- Another way to represent data, reducing their dimensionality to a 2D / 3D space.

- Captures spatial and temporal correlations.

- Leverage computer vision and image-based algorithms.

**Feature Extraction**
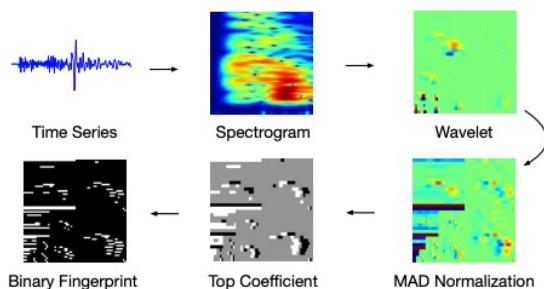
**Image-based ML Classifiers**



**Figure 3:** The fingerprinting algorithm encodes time-frequency features of the original time series into binary vectors.

Source: Kexin Rong et al. at VLDB '18.

**Earthquake Detection:**
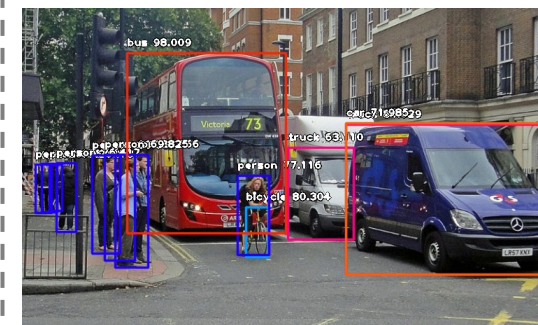Extract Frequencies of Seismic Waves.



Figure 1: Typical workstation of a professional trader. Credit: Photoagriculture / Shutterstock.com.

Source:
J.P. Morgan AI labs.

Figure 4: Various visual representations of the same time-series data.

**Stock Market Forecasting:**
Trading by learning time series data as images.



**Autonomous Driving:**
Object Detection & Recognition

What can an image-based system pipeline look like?

**E.g., predicting future resource utilization.**



Input Image

e.g., memory utilization

**Pattern Recognition**

Class = "sinusoid"

*Choose based on pattern.*

Pre-trained ML models

**Pattern Prediction**

Prediction

**Resource Management System**

*Not Accurate Prediction? Retrain.*

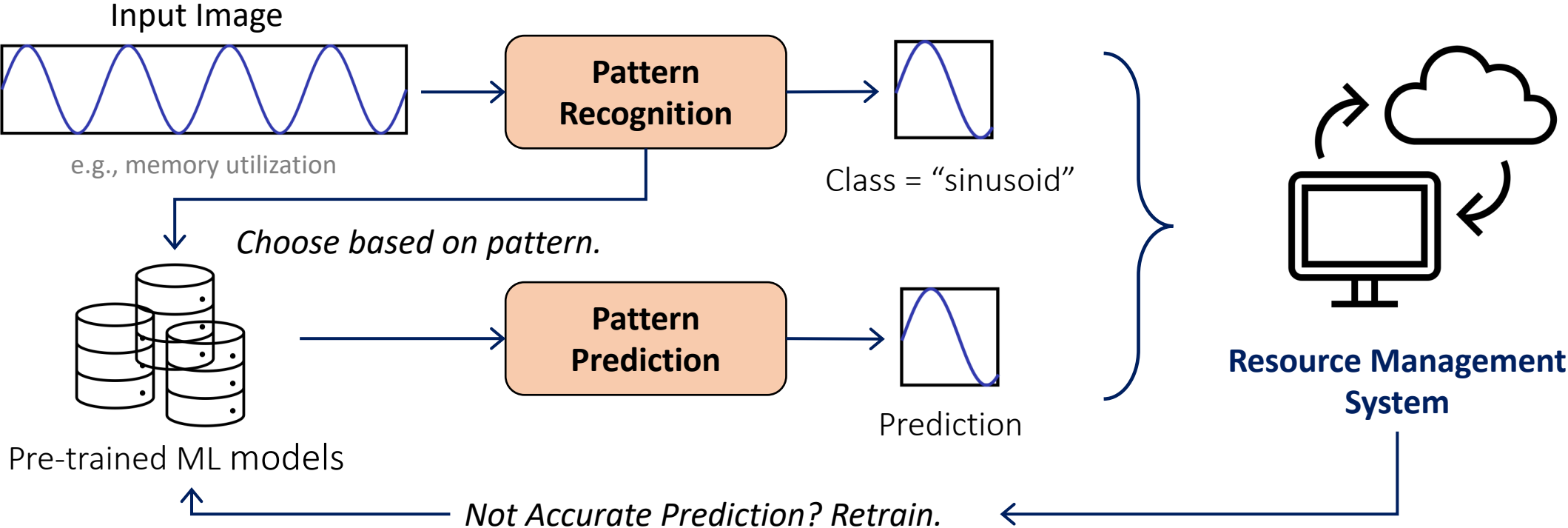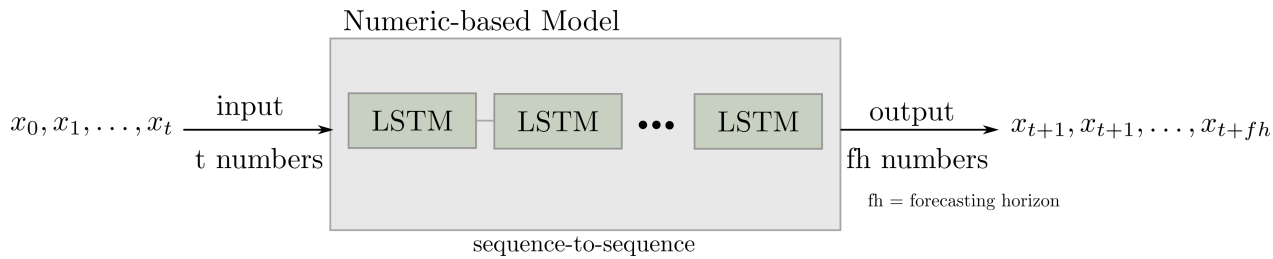# Image-based vs. Number-based Machine Learning

Research paper under submission.

Numeric-based Model

$x_0, x_1, \ldots, x_t$ $\xrightarrow[\text{t numbers}]{\text{input}}$ [LSTM] — [LSTM] $\bullet\bullet\bullet$ [LSTM] $\xrightarrow[\text{fh numbers}]{\text{output}}$ $x_{t+1}, x_{t+1}, \ldots, x_{t+fh}$

fh = forecasting horizon

sequence-to-sequence

**Number-based LSTM model**

Image-based model

$\xrightarrow[\text{1 image}]{\text{input}}$ [ConvLSTM] [ConvLSTM] [ConvLSTM] [ConvLSTM] [ConvLSTM] $\xrightarrow[\text{1 image}]{\text{output}}$
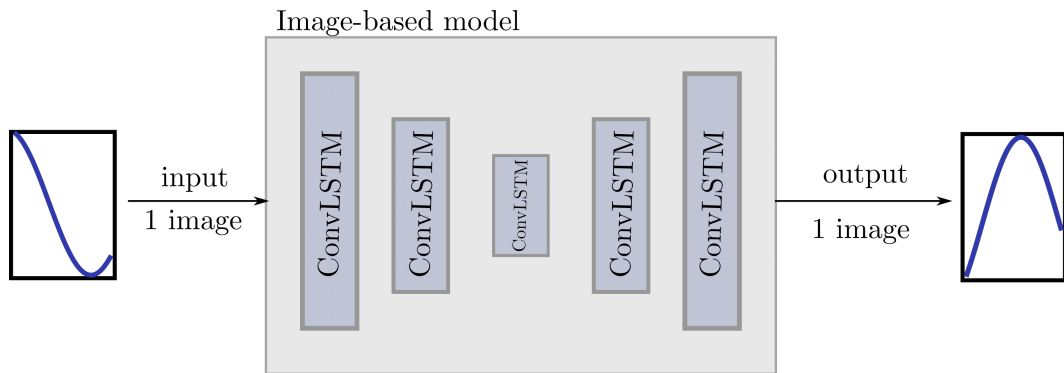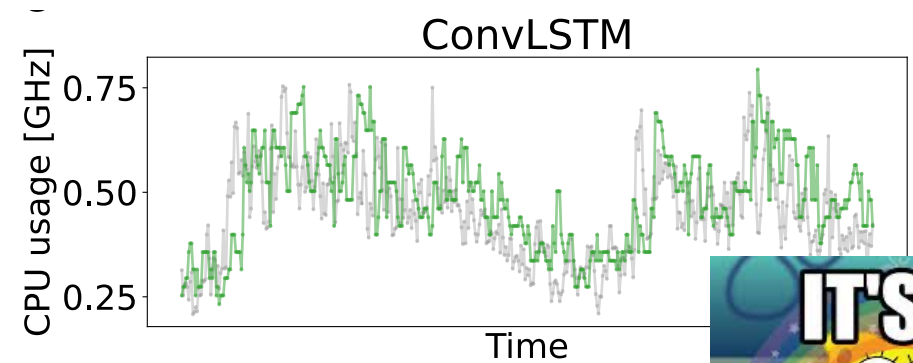
**Image-based ConvLSTM model**

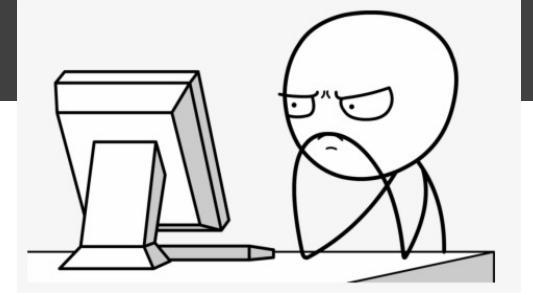The image-based ConvLSTM makes more accurate predictions.

What can an image-based system pipeline look like?

**E.g., learning memory access patterns.**



image + metadata

**Pattern Recognition**

*Choose based on pattern.*

Pre-trained ML models

**Pattern Prediction**

stride

stride

Class = "stride"

Prediction

*Not Accurate Prediction? Retrain.*

Memory Manager

System Level

Hardware Level

hot data

cold data

Hybrid Memory

WORK IN PROGRESS
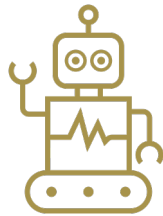
Ground Truth



Prediction



More challenging, since the data access patterns are more complex.

# Talk Outline

**Why do we need Smarter and Faster Systems?**
The evolution of the hardware technologies, calls for software improvements.

**Building *Smart* Systems**
Using machine and human intelligence to build practical ML-based systems.
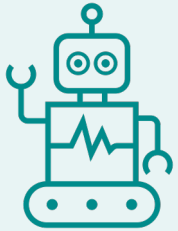
**Building *Fast* Systems**
Reducing ML-based management overheads with visualization.
Building image-based system pipelines.

**Future Research Directions**

# Future Research Directions

My research lies at the intersection of Machine Learning and Systems.

**Machine Learning (ML)**

**Computer Vision (CV)**

**Operating Systems (OS) Software**

ML *for* Systems

E.g., Online practical training, ML for different systems problems.

Systems *for* ML

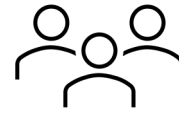E.g., Optimize memory management for RNNs / ML workloads.

ML + CV *for* Systems

E.g., Image-based pattern recognition and prediction of resource usage.

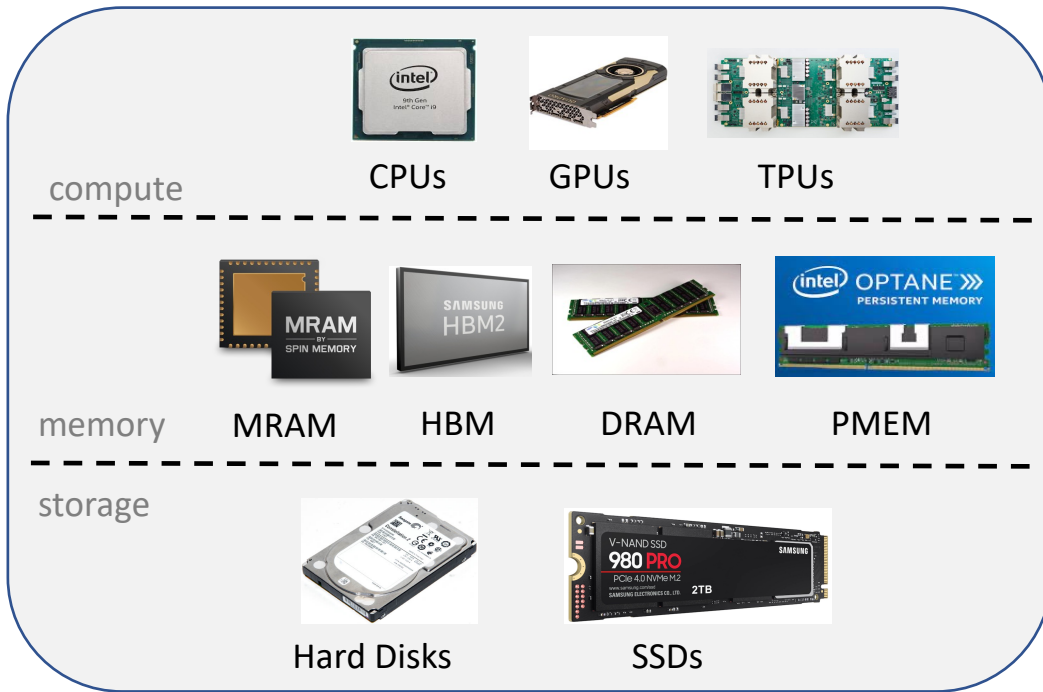# Intelligent Management of Extreme Heterogeneity

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

Multi-tenancy?
Isolation?

Users

Performance?
Cost / Energy /
Resource Efficiency?

compute

CPUs    GPUs    TPUs

memory

MRAM    HBM    DRAM    PMEM

storage

Hard Disks    SSDs

*Local Node*

High-Speed Interconnects

Datacenter

Supercomputer

Massive Node Clusters
Disaggregated Resources

System vs. HW / SW co-design?

**ML integration Aspects:**
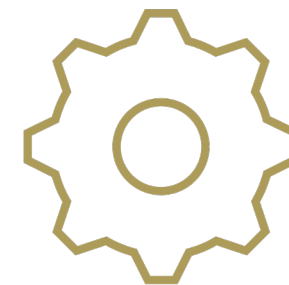Necessity    Effectiveness    Practicality    Interpretability

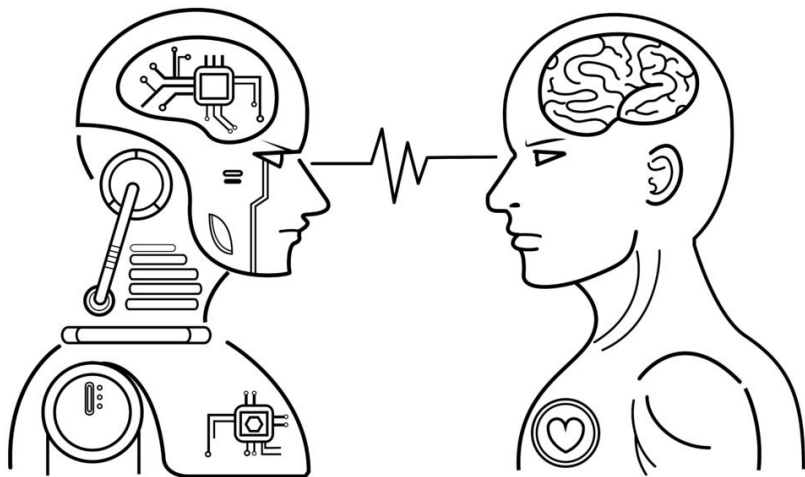Scan this to find more about my work.
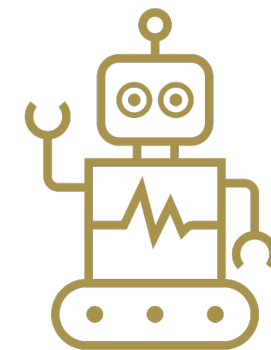
Smart

Fast

Systems

Artificial vs HUMAN Intelligence

How can we use our human intelligence to build **practical** systems that leverage machine learning and computer vision?

Machine Learning

Computer Vision