# Declarative Memory Services

Jeronimo Castrillon

Jana Giceva

Yu Hua
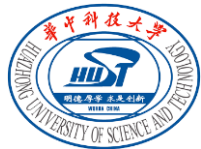
Kimberly Keeton

Akhil Shekar

Kevin Skadron

**Tianzheng Wang**
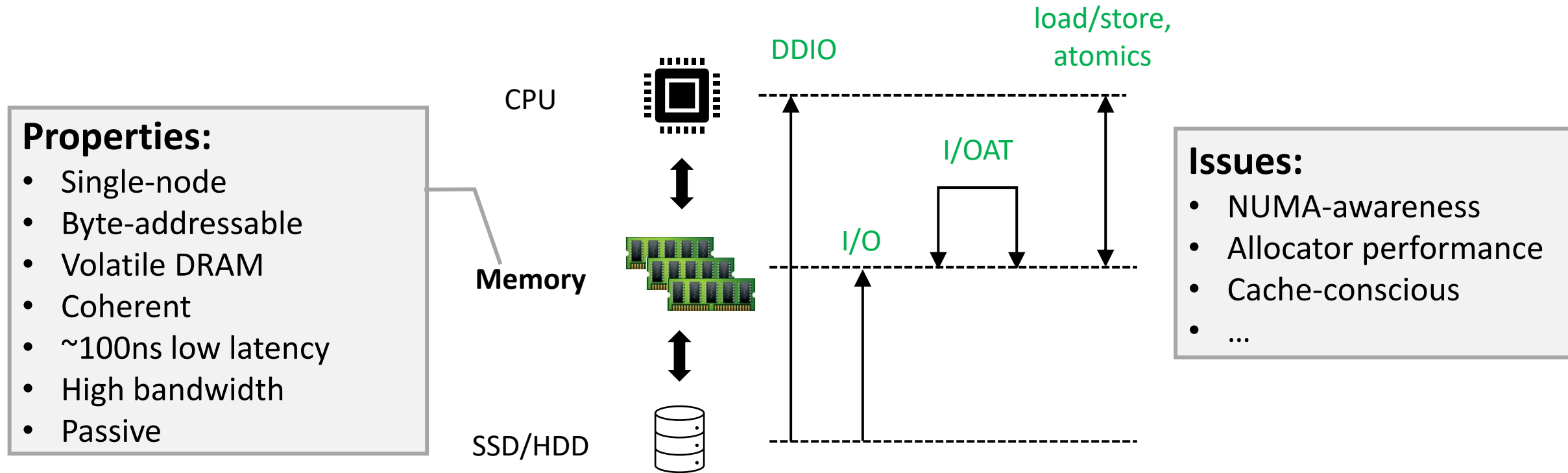
Huanchen Zhang

# "Memory" traditionally…
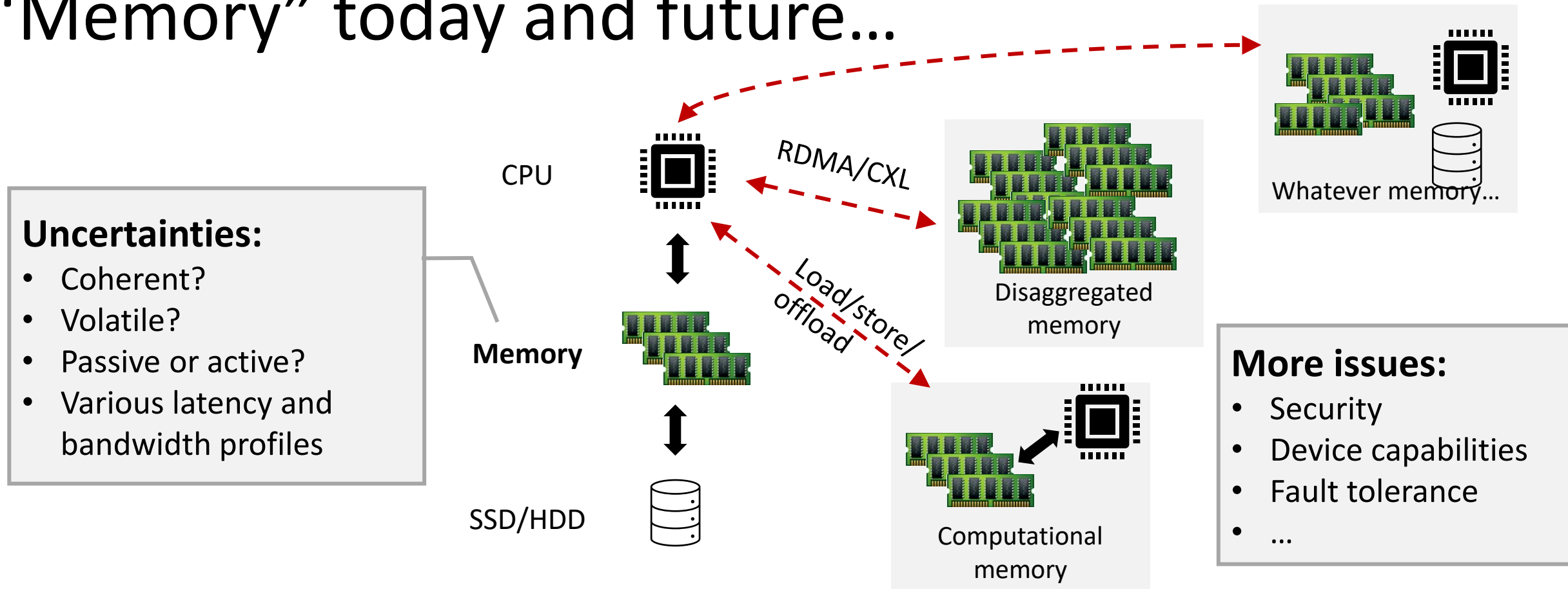
**Properties:**
- Single-node
- Byte-addressable
- Volatile DRAM
- Coherent
- ~100ns low latency
- High bandwidth
- Passive

CPU

Memory

SSD/HDD

DDIO

load/store, atomics

I/OAT

I/O

**Issues:**
- NUMA-awareness
- Allocator performance
- Cache-conscious
- …

Relative tractable primitives and tools + imperative programming

Life was ok.

# "Memory" today and future…



**Uncertainties:**
- Coherent?
- Volatile?
- Passive or active?
- Various latency and bandwidth profiles

CPU

RDMA/CXL

Load/store/ offload

Memory

SSD/HDD

Disaggregated memory

Whatever memory…

Computational memory

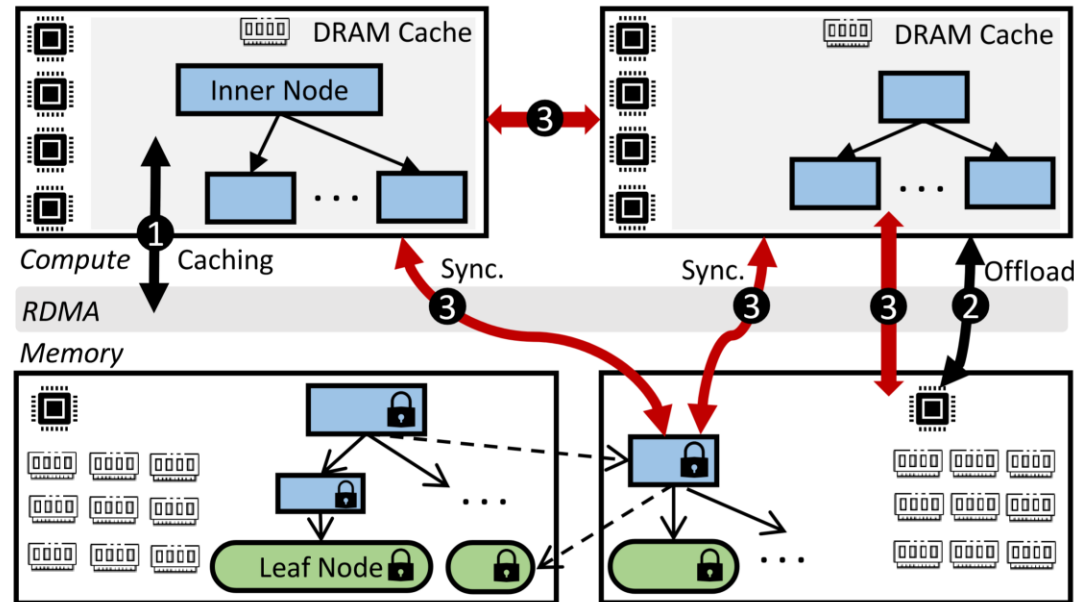**More issues:**
- Security
- Device capabilities
- Fault tolerance
- …

Intractable primitives → highly complex, imperative programming

Life is hard.

# Case Study:
# Adapting a B+-Tree for disaggregated memory

**(1) Longer latency, should cache:**

- Which B+-tree nodes to cache?
- Is there coherence between compute servers?



**(2) Memory has CPU, should offload:**

- How much CPU do I have?
- What operations to offload?

**(3) Data placement + replication:**
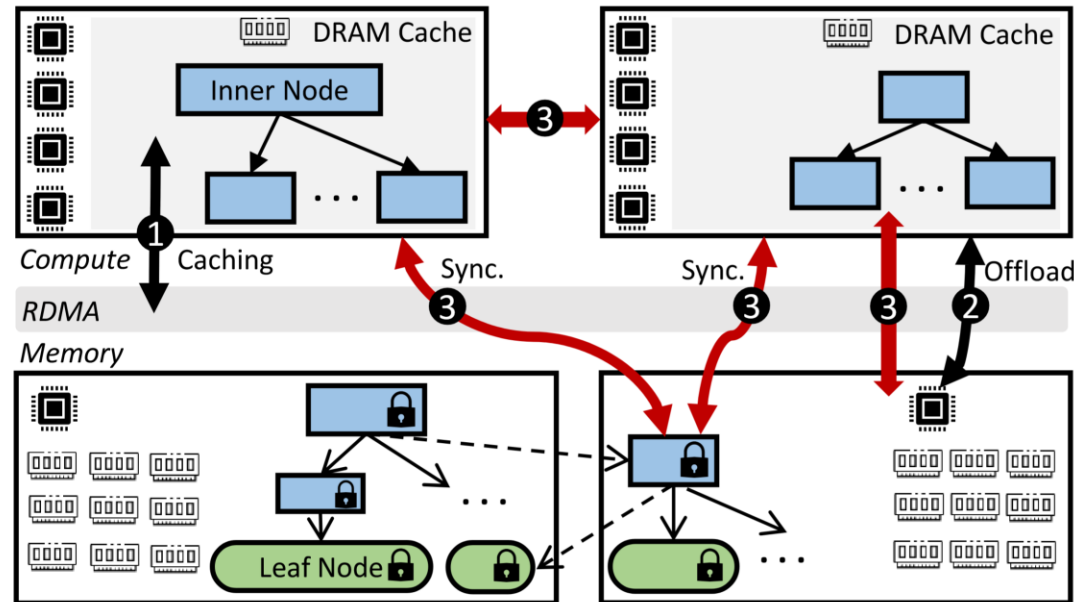
- Who can access which data?
- How to partition?

\* DEX: Scalable Range Indexing on Disaggregated Memory, VLDB 2024

# Case Study:
# Adapting a B+-Tree for disaggregated memory

**(1) Longer latency, should cache:**
- Which B+-tree nodes to cache?
- Is there coherence between compute servers?



**(2) Memory has CPU, should offload:**
- How much CPU do I have?
- What operations to offload?

**(3) Data placement + replication:**
- Who can access which data?
- How to partition?

Hand-coded decisions
Unsustainable (*more cases in paper*).

Case Study:
# Adapting a B+-Tree for disaggregated memory



**(1) Longer latency, should cache:**
- Which B+-tree nodes to cache?
- Is there coherence between compute servers?

**(2) Memory has CPU, should offload:**
- How much CPU do I have?
- What operations to offload?

**(3) Data placement + replication:**
- Who can access which data?
- How to partition?

Hand-coded decisions
Unsustainable (*more cases in paper*).

# Would be nice to be more _declarative_

- Decouple device-specific logic from high-level design
  - "I want this function to be offloaded, if possible"
  - "Latency to access this memory block should not exceed 5ms"

- Simplify programming for today and future, unknown architectures
  - Same DBMS design, any hardware

- Better cross-device optimizations

How to get there?

# Vision: Declarative Memory Services

**Three-layer design:**

- Abstraction Layer
  - Developers work with "logical memory regions" and data flows
  - Annotate with desired properties

- Calibration Layer
  - Discover and index device capabilities
  - Expose device primitives and APIs

- Memory Services Layer
  - A set of generic "memory services" that well use memory devices
  - Jointly optimize for the application based on annotations

Caveat: yet to implement, this is pure vision!

# Declarative Abstraction layer



**B+-tree node definition:**

```
struct InternalNode {
    KV kv_pairs[MAX_KV];
    int key_count;
    …
};
```

Declare desired properties

**Previously:**

```
InternalNode *n = allocate(…)

// hand-made decision to cache it
cache.insert(n);
```

**Now with DMS:**

```
[cacheable, coherent, latency < 10μs]
InternalNode *n = allocate(…)

// placed in coherent, compute-side memory, by DMS
cache.insert(n);
```

Data flows work similarly:
- Properties attached to tasks, enforced by DMS runtime

Physical design and logical functionality decoupled

# Calibration Layer

- Discovers and track device capabilities, provide APIs

- Key component: device catalogue
  - A table that evolves with hardware changs

| Device | Capabilities | APIs | Characteristics |
|--------|-------------|------|-----------------|
| Local DRAM | Coherence Byte-addressable | dram-load, dram-store, dram-dsa, atomics… | … x GBps within socket, under y load… |
| CXL DRAM | Partial coherence Byte-addressable | cxl-load, cxl-store… | … 300ns best - 1us worst latency… |
| Membrane (computational memory) | Compute Byte-addressable | pim-load, pim-store, pim-offload… | … x ns latency with host… |

Implemented and maintained by DMS developers

Challenging

Declarative Memory Services

# Memory Services Layer

- Use device catalogue APIs to build services



- DEX example:
  - Services needed: data placement and caching
  - Upon allocation: place data based on annotated desired properties
  - Runtime: lightweight metadata tracking for caching

- Customized policies possible
  - "Please don't evict parent node before child node"
  - "Please use this encoding scheme for such and such data"

# Research Challenges and Agenda

- Device Characterization
  - Beyond simple stats: e.g., latency behaviour under varying load levels
  - Self-evolving the device catalogue with new hardware

- Properties → Services: When to pick which implementation?

- SLA Guarantees
  - Memory services monitor metrics, and migrate between services to meet SLO
  - How to deal with conflicting SLAs?
    - E.g., tenants prioritizing throughput vs. latency

- DMS Deployment
  - DMS requires non-trivial information (global and local server) to work

- Correctness and Debugging
  - DMS-based programs are declarative
  - How to verify their correctness and debug them? Tools for exploring why an SLO was missed?

# Summary

- Memory is heterogeneous: complexity arises with more features
  - Current approach to leveraging memory devices is unsustainable
  - Hand-crafted with low-level primitives
  - Getting worse as hardware evolves

- **Declarative Memory Services**
  - <u>Developers</u> specify logical functionality
  - <u>Calibration layer</u> discovers and characterises devices
  - <u>Memory services</u> provide physical implementations and optimizations

*Thank you!*