

TORSTEN HOFELER, ROBERT ROSS, TIMOTHY ROSCOE

Distributing the Data Plane for Remote Storage Access

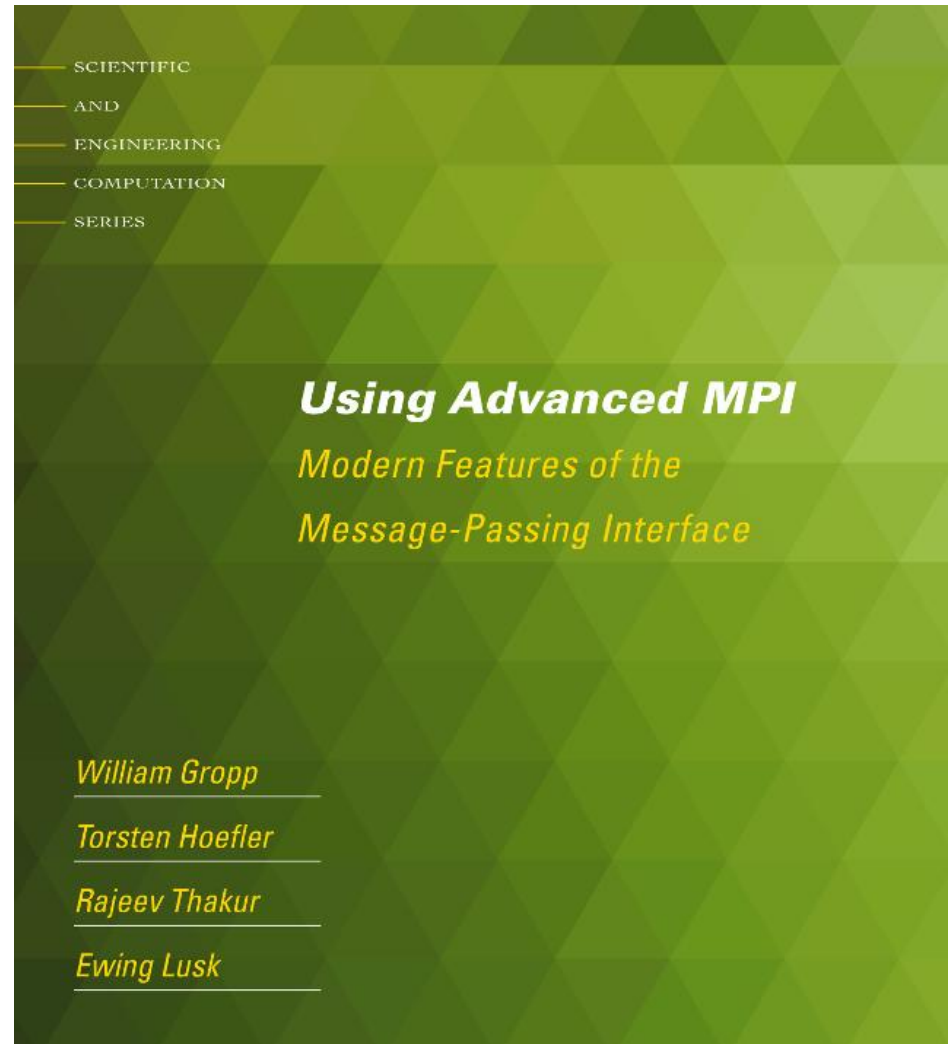


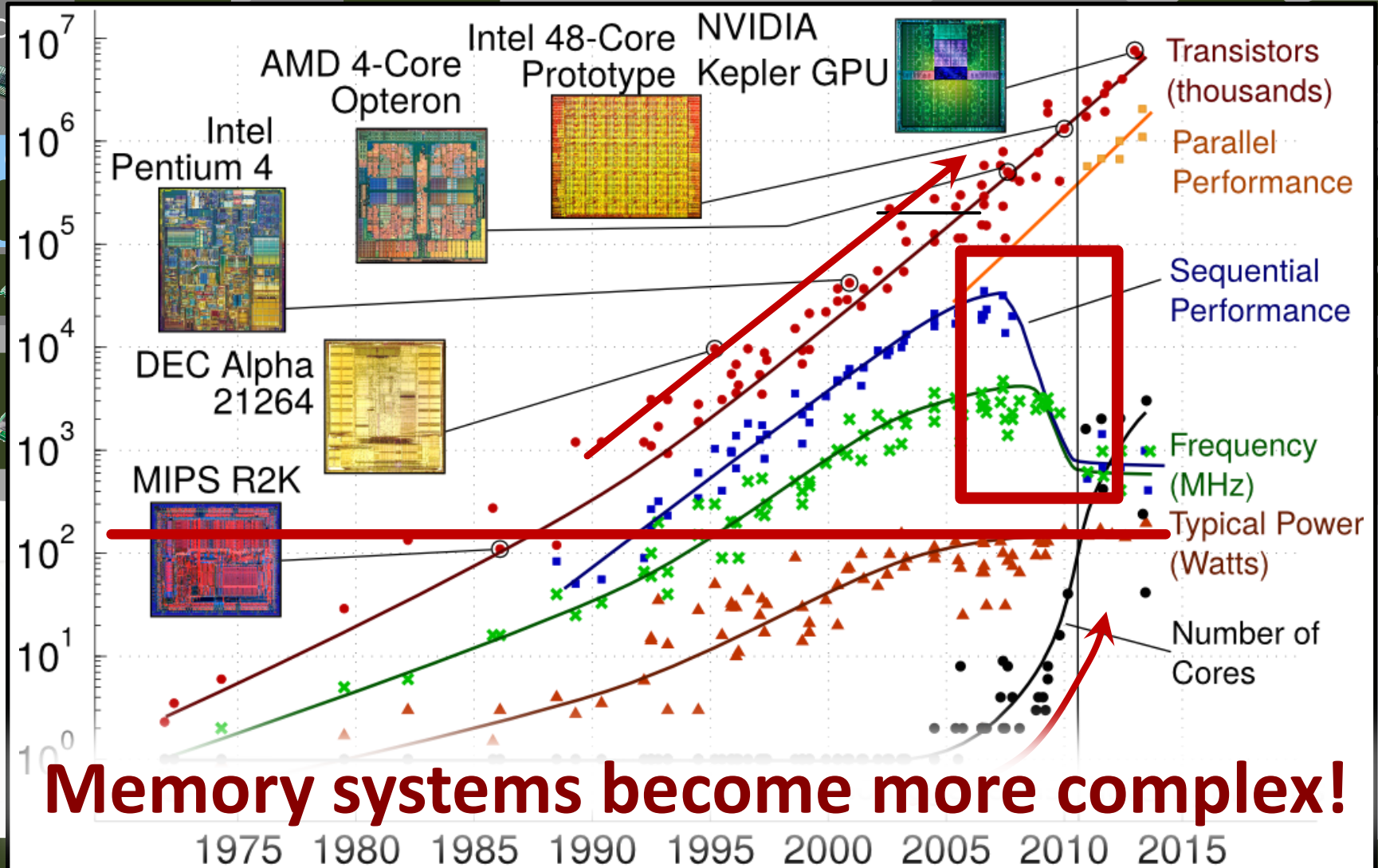
My outing

- I'm an HPC (systems) guy



- Programming Models
- Performance Models
- Network (Models)

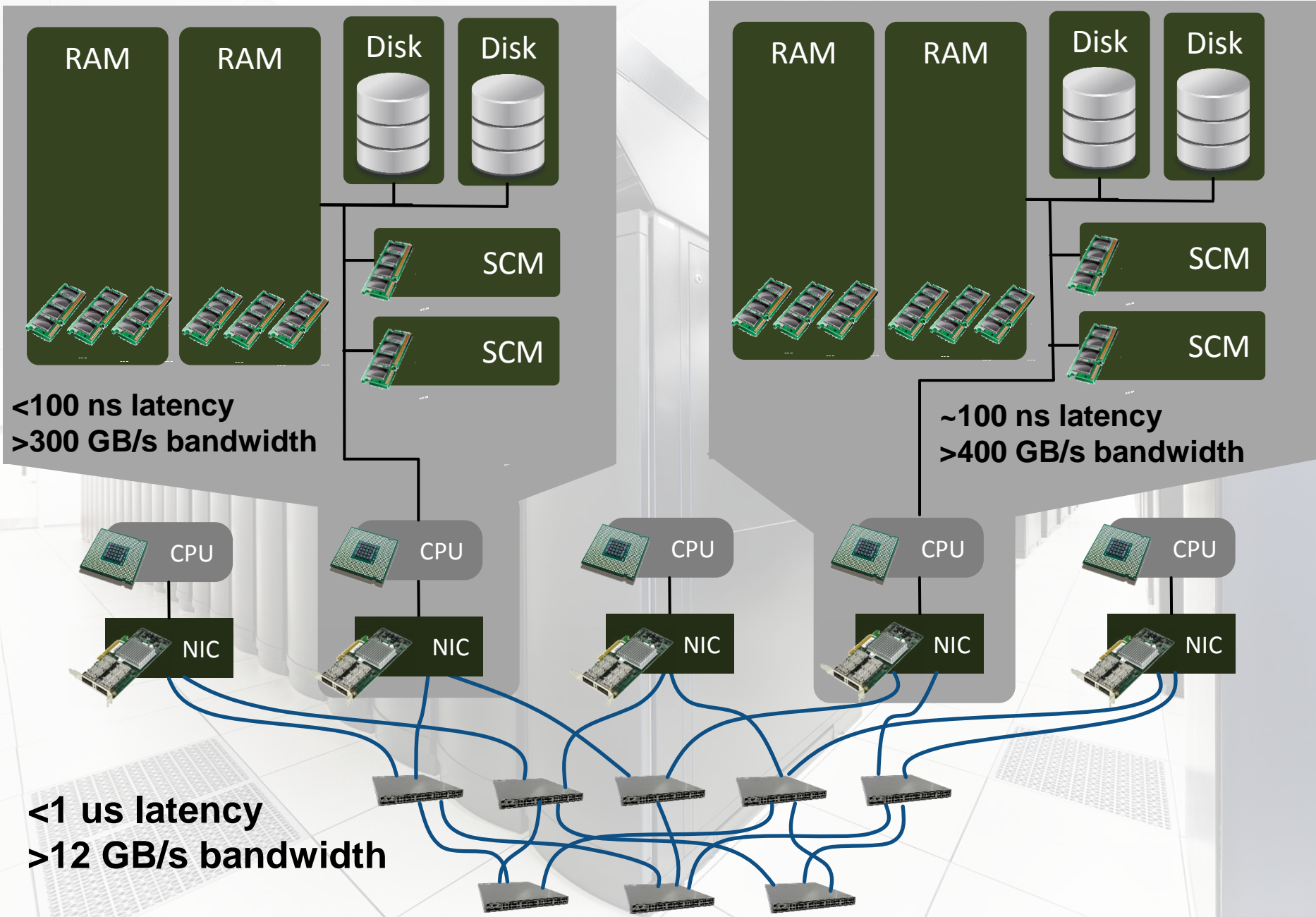




Memory systems become more complex!

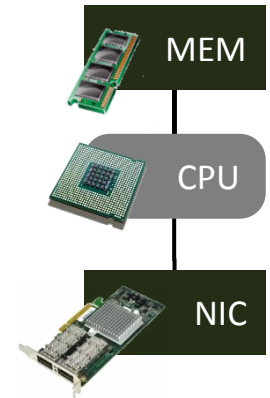
Data partially collected by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond

Data movement is the new challenge!



The diagram illustrates the decomposition of a query into three parallel subqueries (A, B, and C) and their corresponding data partitions. The data is distributed across five cylinders representing different partitions.

- Subquery A:** Involves tables AB1, AB2, A1, and A2. The data is distributed across the first two cylinders.
- Subquery B:** Involves tables B1, BC1, B2, and AB4. The data is distributed across the third cylinder.
- Subquery C:** Involves tables C1, C2, AC1, and C3. The data is distributed across the fourth and fifth cylinders.



- [illegible]

[3]: Murray, D. G., McSherry, F., Isaacs, R., Isard, M., Barham, P., and Abadi, M. Naiad: A Timely Dataflow System. In SOSP'13

A (30-year-old) networking idea revived

- **Separating data and control plane**

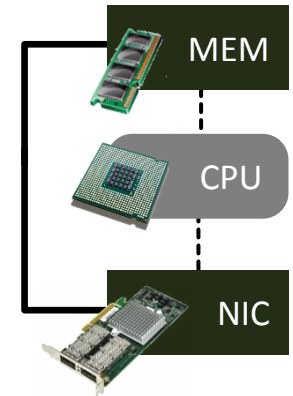
- Get the CPU/software out of the way!
- Software-defined IO (SDIO, cf. NASD [1], Aerie [2])
- *We set up a network route right into the device!*

- **Key point is (direct) access to storage**

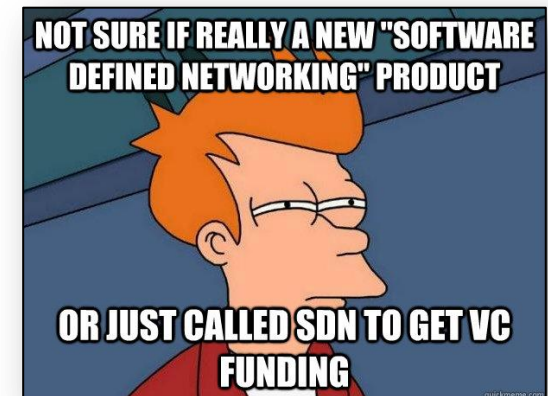
- Allocation
- Read/Write
- Protection
- Caching
- Consistency, coherence, durability

- **Our central research question:**

- *How can we design a fast software/hardware data plane for safe, secure, direct remote access to persistent storage devices?*



<1 us latency

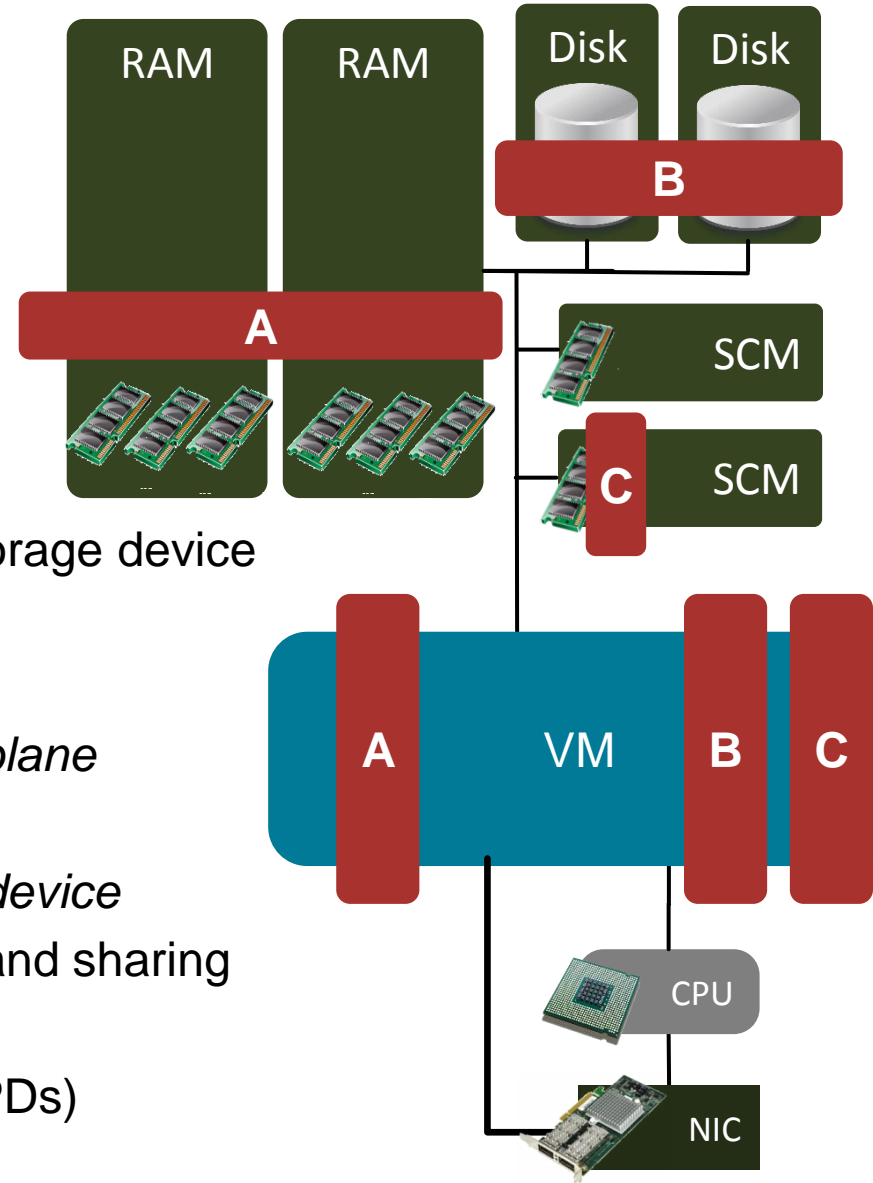


[1]: Gibson, G., et al. A Cost-effective, High-bandwidth Storage Architecture. SIGPLAN Not. 33, 11 (Oct. 1998)

[2]: Volos, H., et al.. Aerie: Flex-ible File-system Interfaces to Storage-class Memory, EuroSys'14

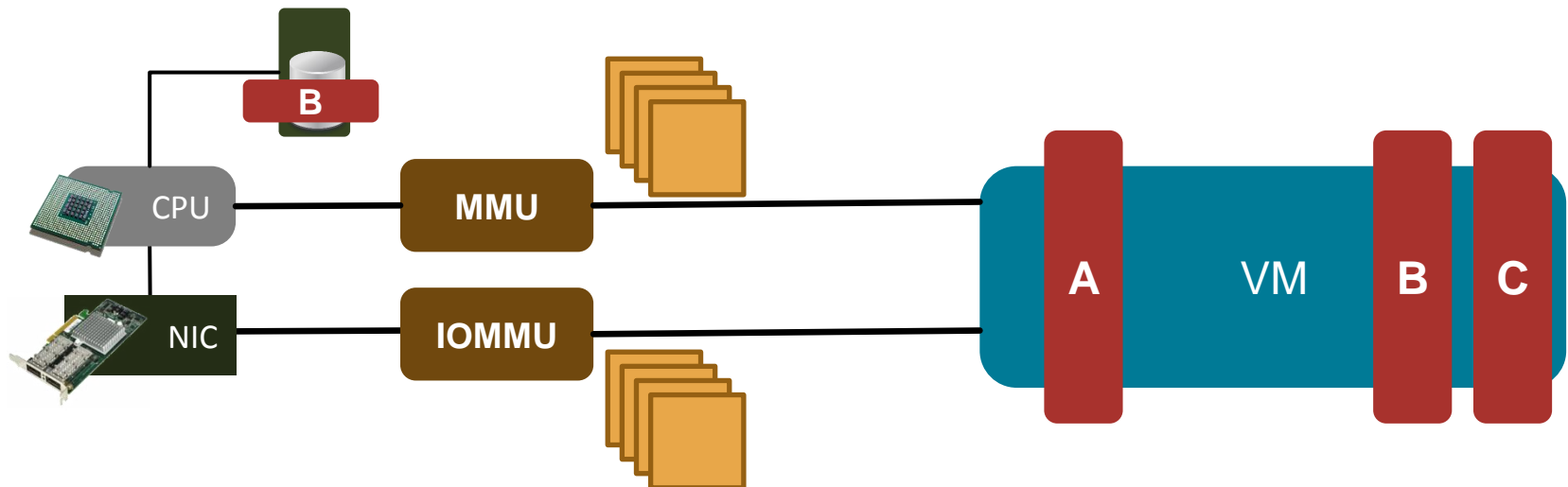
(1) Device allocations

- **Exokernel-like filesystem library**
 - Data path 100% in user-level
Manages metadata, data, and coordination
 - Data is stored in allocations
- **An allocation ...**
 - Is an area of main memory or on a storage device
Placed explicitly
 - Can be created, opened, or closed
Using a central or distributed control plane
 - Has a contiguous address space
Block translation implemented in the device
 - Is the smallest unit of access control and sharing
Named in a global namespace
 - Access through capabilities (e.g., IB PDs)



(2) Read/Write and protection

- **Allocations are accessed ...**
 - Locally via MMU-mappings (cf. Aerie)
 - Remotely via IOMMUs or other address translations
- **Not all devices are part of the physical address space**
 - Mainly legacy two options:
 1. Software fallback (monitor RDMA regions and keep consistency)
 2. Use IOMMU logging schemes (cf. Active Access [1])



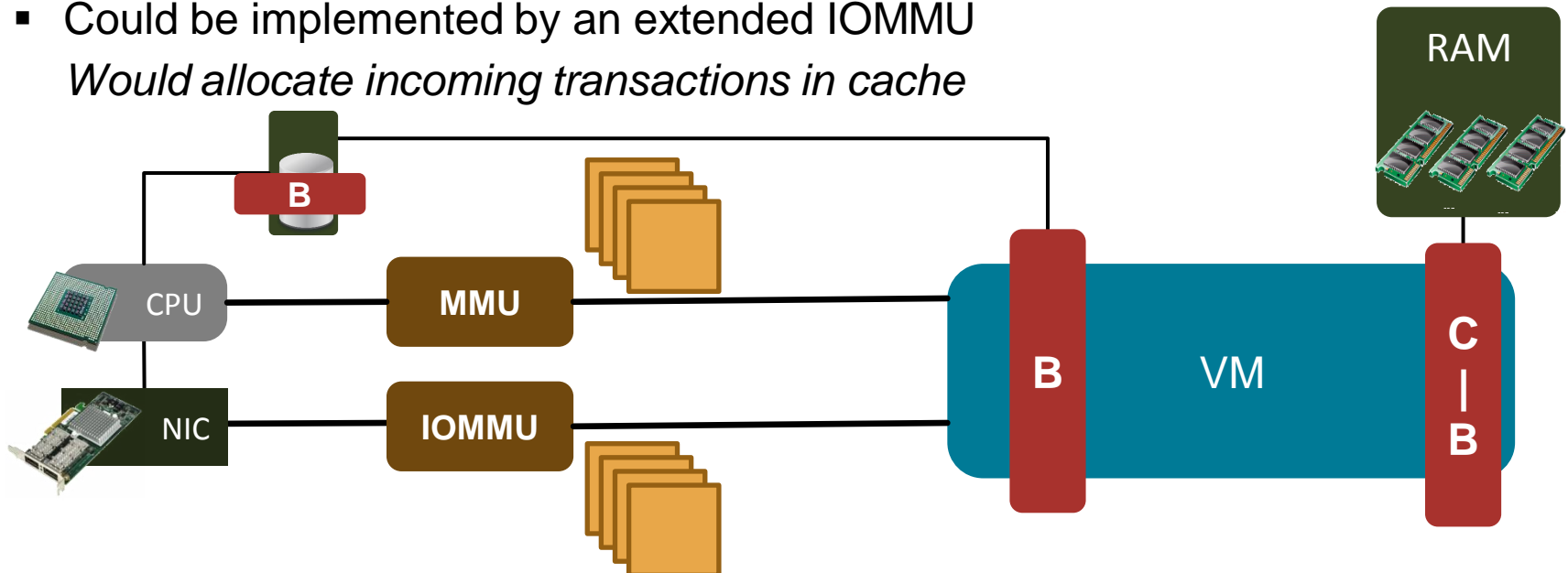
(3) Caching

■ Software caching

- Explicit caching, similar to RMA programming
Allocate local cache and access it
- Can be application-specific or caching library (standard techniques)

■ Hardware caching

- Set up local memory as cache for remote allocations
Use standard (e.g., LRU) replacement policies
- Could be implemented by an extended IOMMU
Would allocate incoming transactions in cache



(4) Consistency and coherence

- **Current RDMA does not support consistent atomic access**

- At least not large enough

- **We propose a weaker consistency model**

- All read/write accesses are nonblocking!

- Arrange accesses into epochs separated by fence operations

Modified data is only valid at the end of an epoch

- The type determines the isolation level

Shared: only consistency after epoch ends

Exclusive: consistency + atomicity

Persistent: consistency + durability

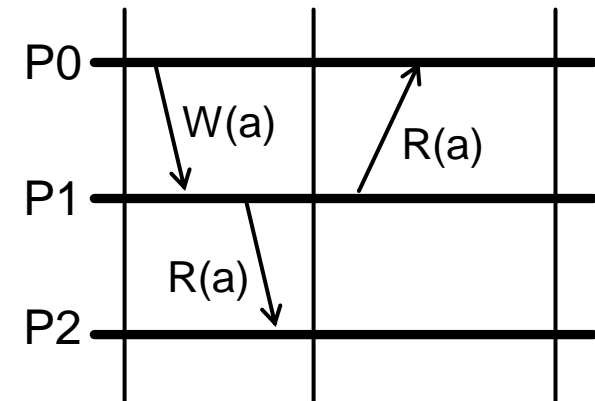
Optimistic: consistency + atomicity but can fail

- Types can be combined

e.g., persistent + exclusive

- **Implemented similar to other RMA programming models [1]**

- May require remote flushes (in the worst case RPCs)



Other filesystem requirements

■ Crash recovery

- Use transactional (optimistic, exclusive, persistent) epochs for metadata
- Must ensure that locks time out if processes disappear

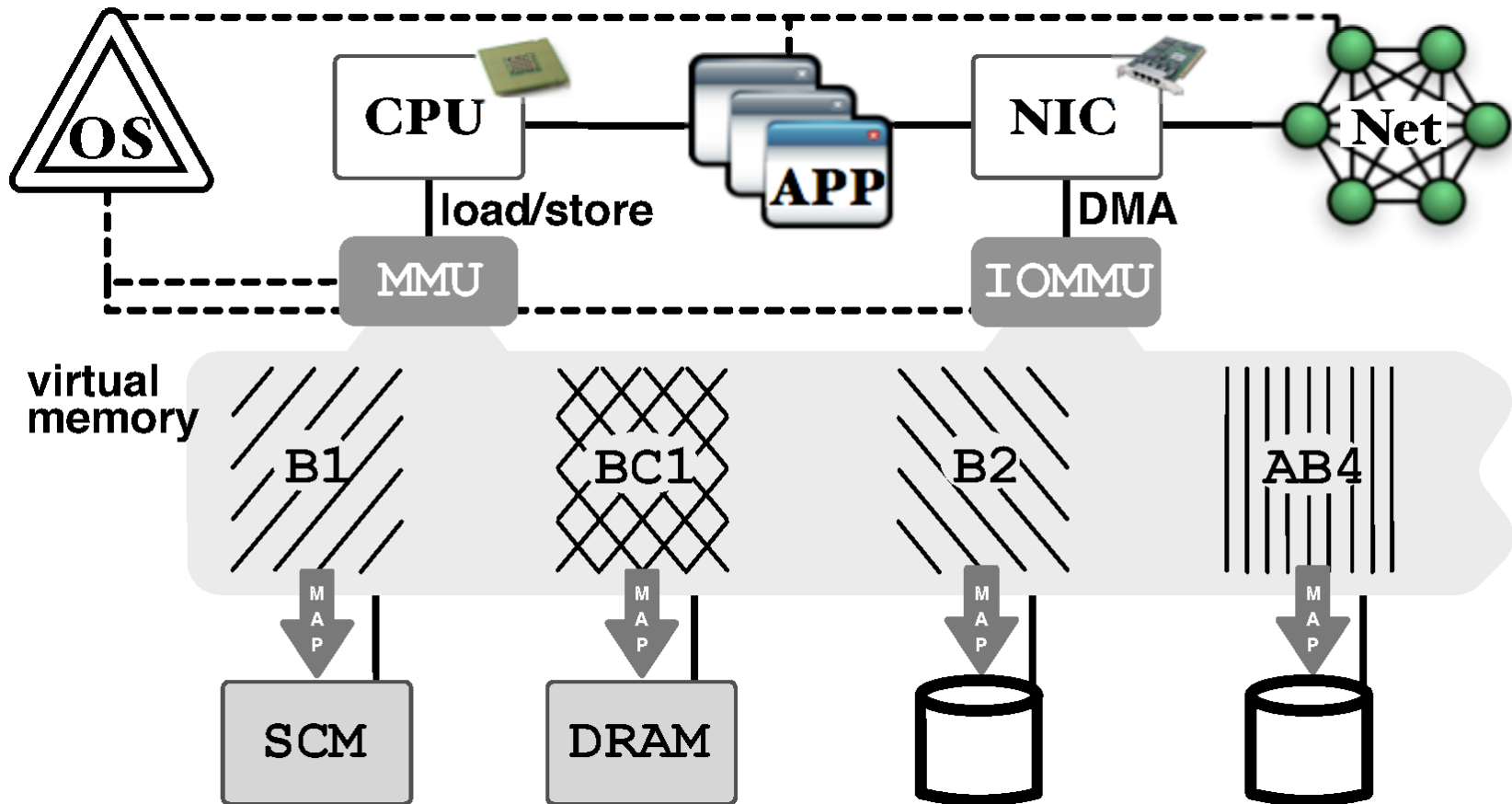
■ Scalability

- Scoping limits context of coherency/epochs (e.g., a shared file)
- Integration with programming model (e.g., MPI-3 RMA)

■ Compatibility

- Provide standard library of user-level file systems
- POSIX consistency with single-operation exclusive, persistent epochs
- Magic byte in allocation allows automated “mounting” like files

Discussion



Merge the network transparently into the file system.