# Data Management in Heterogeneous Memory Systems

Thaleia Dimitra Doudali, Ada Gavrilovska
thdoudali@gatech.edu, ada@cc.gatech.edu

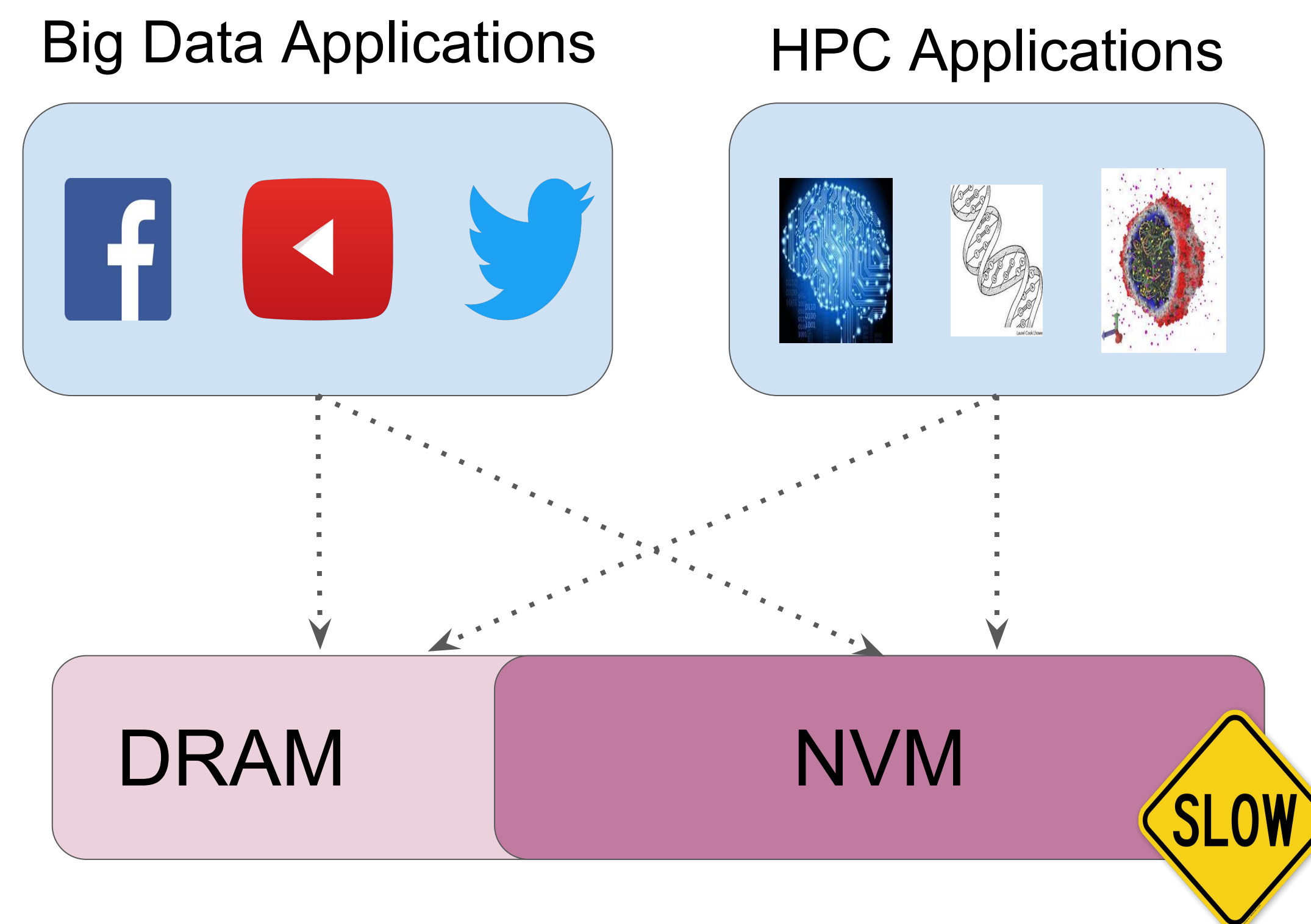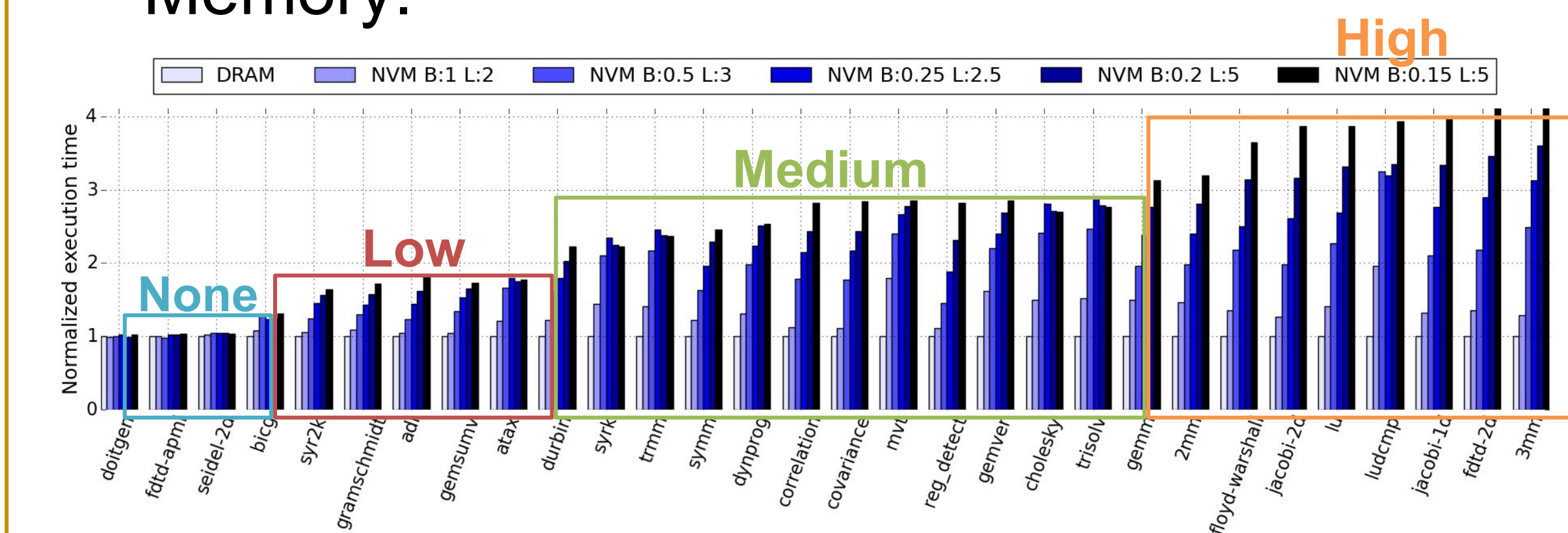## 1. Motivation

- ❖ High Performance Computing and Big Data applications have dataset sizes that often exceed the available DRAM capacities.

- ❖ Emerging memory technologies that are much cheaper, such as Non Volatile Memory, are used to extend the memory space creating a **heterogeneous memory subsystem**.

- ❖ Data in Non Volatile Memory will incur higher access latencies, affecting the application performance, slowing it down compared to an ideal case when all data could fit in DRAM.

- ❖ Existing solutions **reduce the performance slowdown** by prioritizing allocations of the most frequently accessed objects in DRAM. However, they assume fixed hardware capacities.
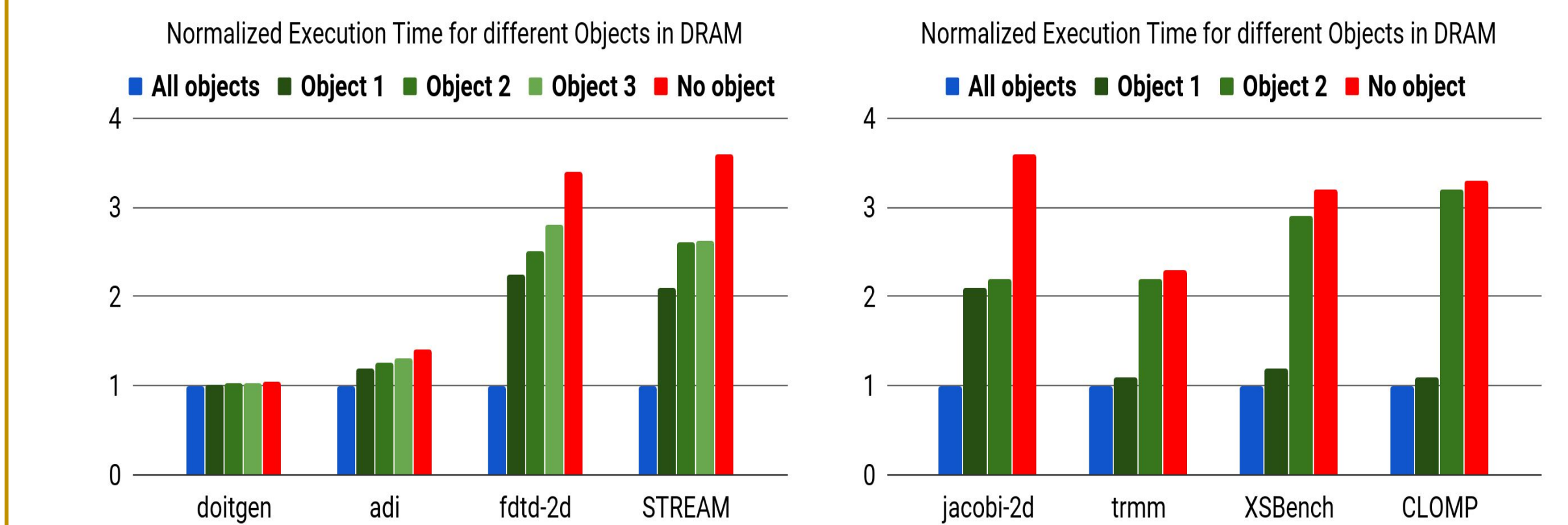
## 2. Problem Statement

Big Data Applications        HPC Applications



DRAM          NVM          SLOW

**Problem:** How do we _size_ the memory components and _manage_ data over the heterogeneous memory system, so as to _minimize_ the application performance slowdown?

## 3. Observations

- ❖ Not all applications are slowed down in the same degree when accessing Non Volatile Memory.



- ❖ Not all data objects of an application help reduce the performance slowdown when allocated in DRAM.
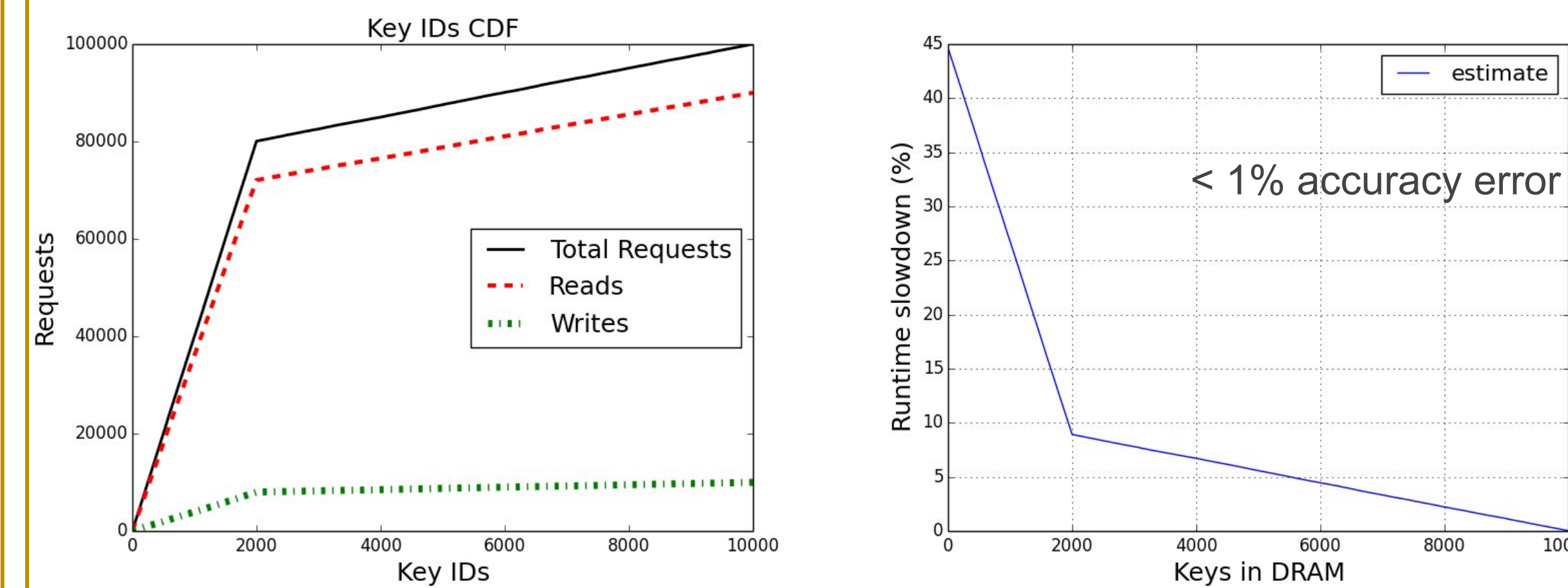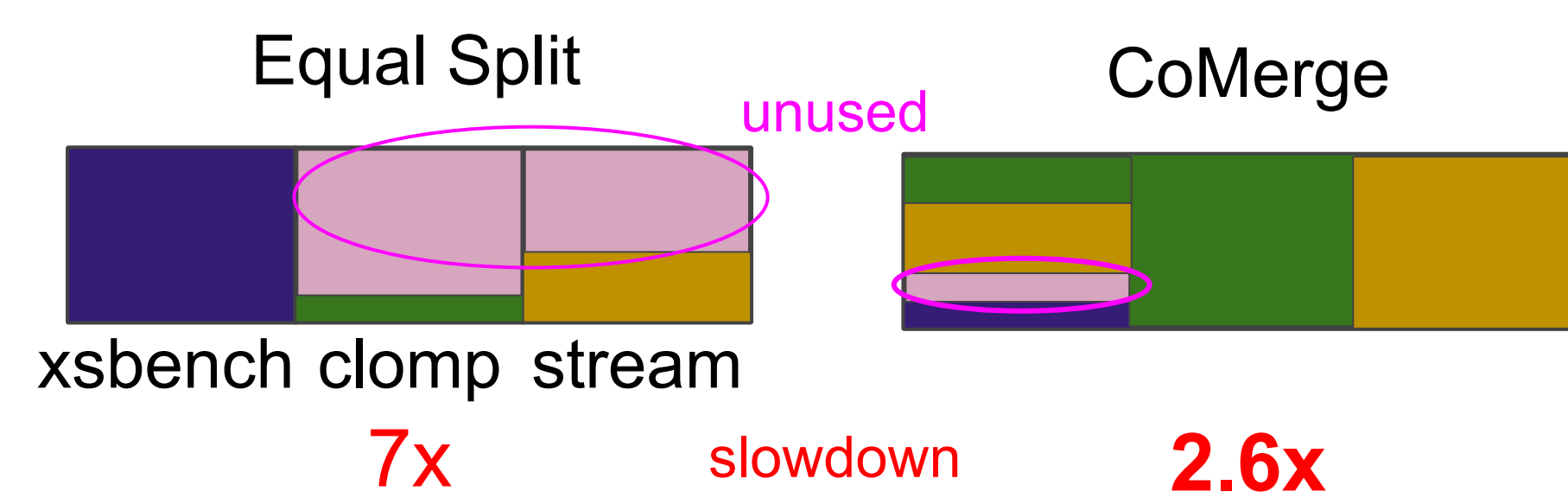


Normalized Execution Time for different Objects in DRAM



## 4. CoMerge Solution

**CoMerge:** Memory sharing policy that prioritizes DRAM allocations for critical data objects. Achieves:

- Lower runtime across all collocated applications.
- Higher DRAM utilization.



Equal Split — unused — CoMerge

xsbench  clomp  stream

7x    slowdown    2.6x

*Thaleia Dimitra Doudali and Ada Gavrilovska. 2017. CoMerge: Toward Efficient Data Placement in Shared Heterogeneous Memory Systems. In Proceedings of MEMSYS 2017, Alexandria, VA, USA, October 2–5, 2017, 11 pages.*

## 5. Mnemo Solution



Key IDs CDF

< 1% accuracy error

**Mnemo:** Profiling tool that estimates the application performance slowdown for incremental DRAM capacity on a heterogeneous memory system.

_under submission_

## 6. Future Directions



| How? | Use Cases | Applications |
|---|---|---|
| • Offline Profiling  Mnemo | • HW capacity sizing | • Key-Value stores |
| • Manual Exploration  CoMerge | • Efficient resource distribution | • HPC kernels |
| • OS dynamic solution | • Scheduling | • HPC apps |
|  |  | • Databases |
|  |  | • Graphs |
|  |  | • ML/DNN |

**Goal**: connect the rest of the dots.