# DICE: Compressing DRAM Caches for Bandwidth and Capacity
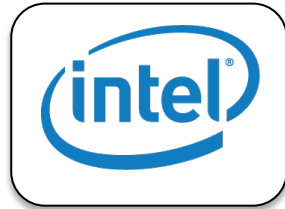
**Vinson Young**

*Prashant Nair*

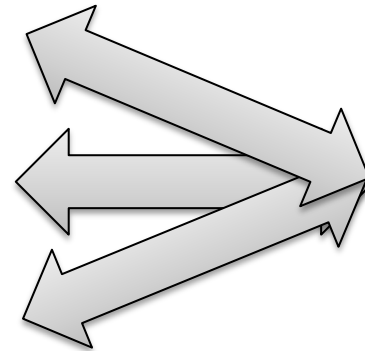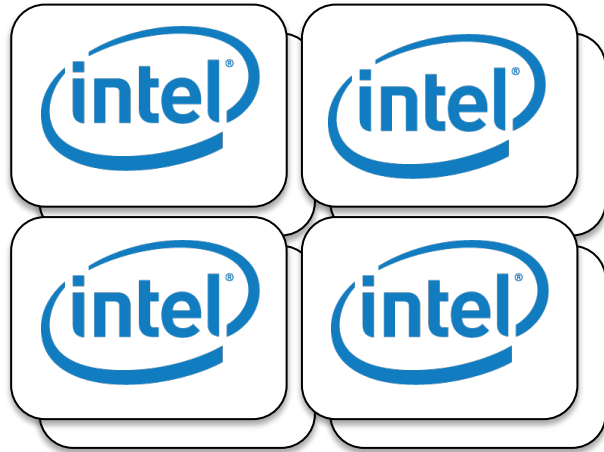*Moinuddin Qureshi*

Georgia Institute of Technology

**Moore's scaling encounters Bandwidth Wall**

## Moore's scaling encounters Bandwidth Wall

**Moore's scaling encounters Bandwidth Wall**

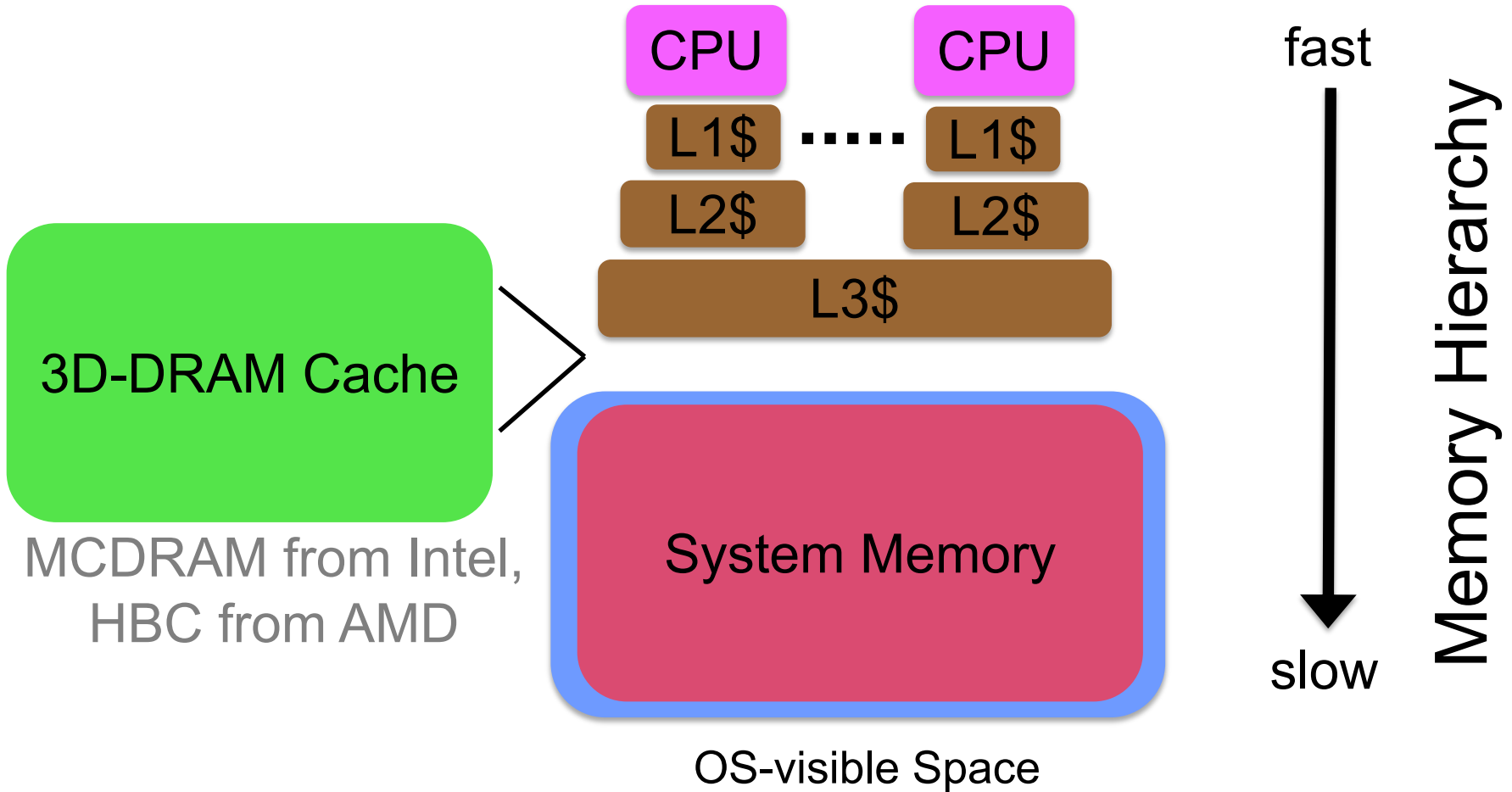# 3D-DRAM MITIGATES BANDWIDTH WALL



3D-DRAM

Hybrid Memory Cube (HMC) from Micron,
High Bandwidth Memory (HBM) from Samsung

# 3D-DRAM MITIGATES BANDWIDTH WALL



3D-DRAM

**3D-DRAM improves bandwidth, but does not have capacity to replace conventional DIMM memory**
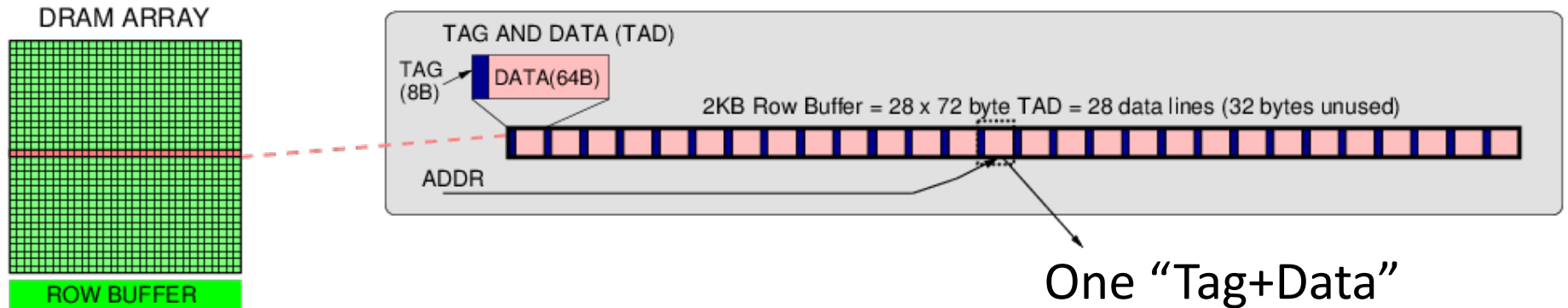
# 3D-DRAM AS A CACHE (3D-DRAM CACHE)



**Architecting 3D-DRAM as a cache can improve memory bandwidth (and avoid OS/software change)**

# PRACTICAL 3D-DRAM CACHE: ALLOY CACHE

Tags "part-of-line" ➔ Alloy **Tag+Data** ➔ Avoid **Tag Serialization**
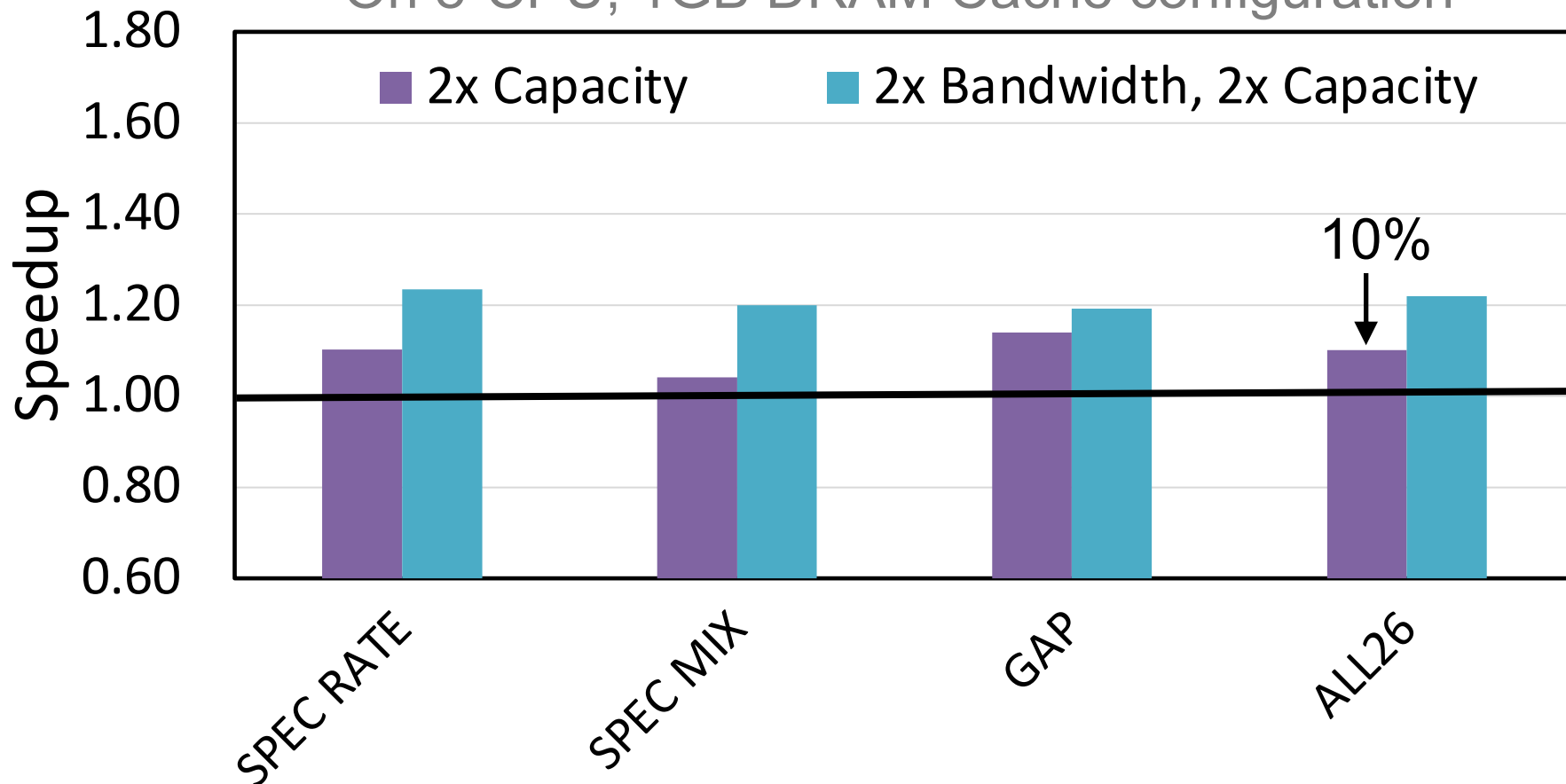


One "Tag+Data"

**Similar to DRAM Cache in KNL:** Direct-mapped, Tags in ECC

**Practical DRAM cache: low latency and bandwidth-efficient**

# 3D-DRAM CACHE BANDWIDTH IS IMPORTANT



On 8-CPU, 1GB DRAM Cache configuration

2x-capacity cache improves performance by 10%.
And, additional 2x bandwidth increases speedup to 22%.
Improving both bandwidth and capacity is valuable.
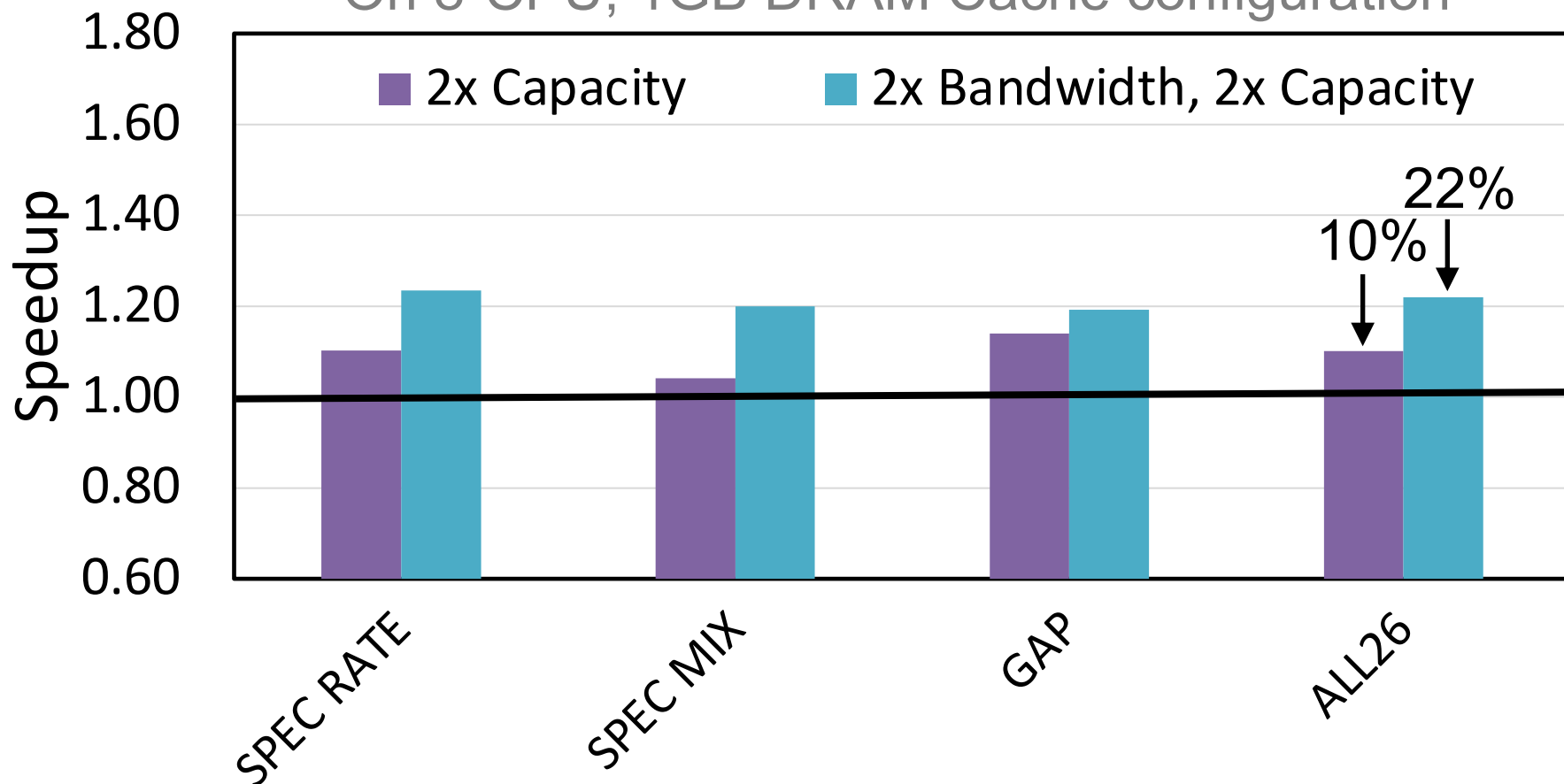
# 3D-DRAM CACHE BANDWIDTH IS IMPORTANT

On 8-CPU, 1GB DRAM Cache configuration
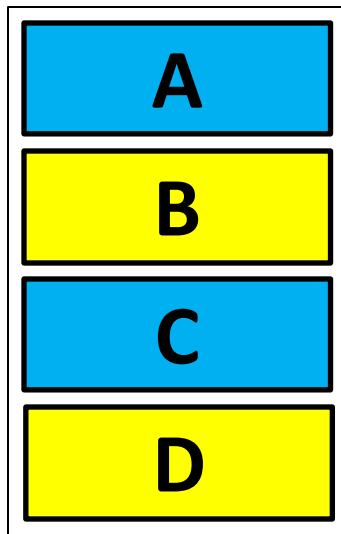


2x-capacity cache improves performance by 10%.
And, additional 2x bandwidth increases speedup to 22%.
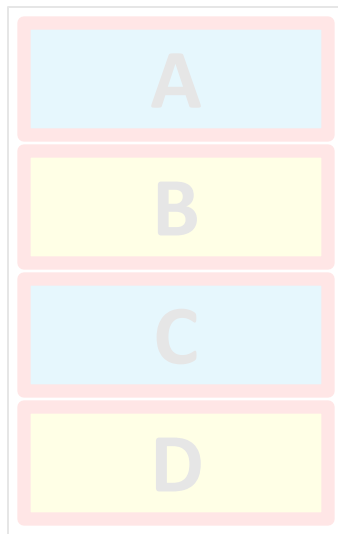Improving both bandwidth and capacity is valuable.

Baseline: Direct-Mapped, One Data Block in an access



A ➝ B ➝ C ➝ D

**Baseline**

Traditional
Compression
(**Incompressible**)

Spatial
Indexing
(**Compressible**)

Spatial
Indexing
(**Incompressible**)

Baseline: Direct-Mapped, One Data Block in an access



**Baseline**

Traditional
Compression
(**Incompressible**)

Spatial
Indexing
(**Compressible**)

Spatial
Indexing
(**Incompressible**)

Compression: Adds capacity



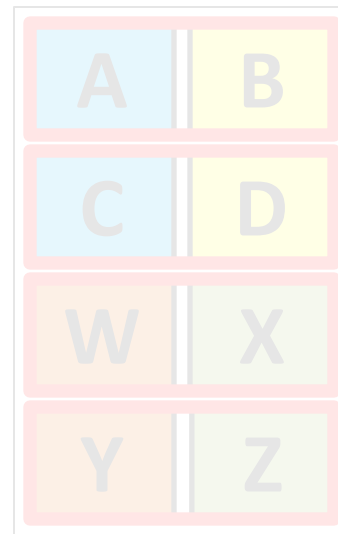Traditional Compression (**Compressible**)

Traditional Compression (**Incompressible**)

Spatial Indexing (**Compressible**)

Spatial Indexing (**Incompressible**)

Compression: Adds capacity



**Traditional Compression (Compressible)**

**Traditional Compression (Incompressible)**

**Spatial Indexing (Compressible)**

**Spatial Indexing (Incompressible)**
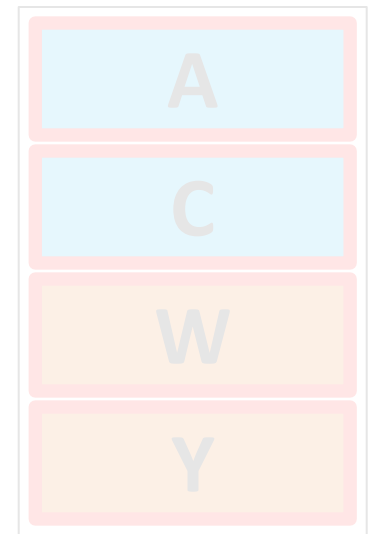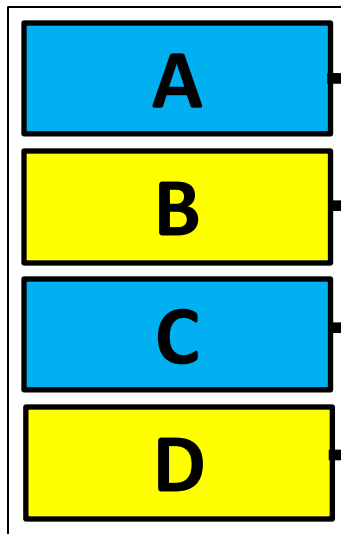
Compression: Adds capacity



Traditional Compression (Compressible)

**Traditional Compression (Incompressible)**

Spatial Indexing (Compressible)

Spatial Indexing (Incompressible)

Compression: Adds capacity

Compression: Adds capacity

A ➞ B ➞ C ➞ D

**4 accesses**
**@**
**1x-2x Capacity**
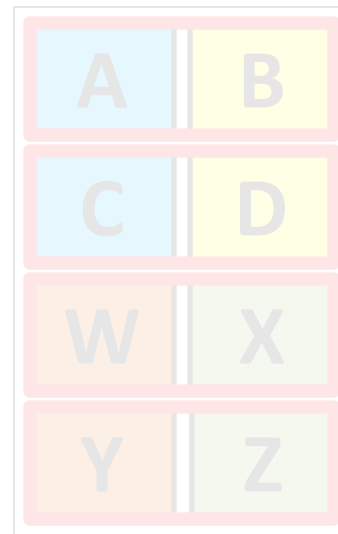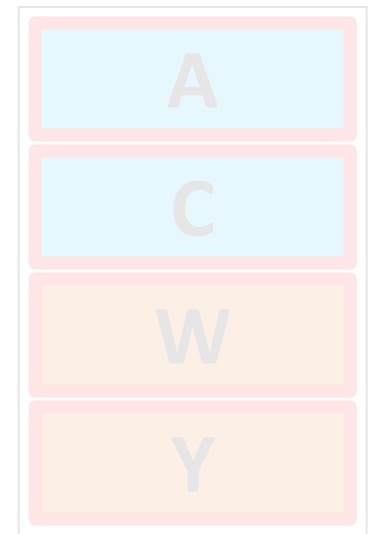
Traditional
Compression
(**Compressible**)

Traditional
Compression
(**Incompressible**)

Spatial
Indexing
(**Compressible**)

Spatial
Indexing
(**Incompressible**)

Compression: Adds capacity, improve bandwidth?



A ⟶ B ⟶ C ⟶ D

**Traditional Compression (Compressible)**

**Traditional Compression (Incompressible)**

**Spatial Indexing (Compressible)**

**Spatial Indexing (Incompressible)**

Compression: Adds capacity, improve bandwidth?

A ➝ B ➝ C ➝ D

| A | B |
| C | D |

*2x Bandwidth*

**Traditional Compression (Compressible)**

**Traditional Compression (Incompressible)**

**Spatial Indexing (Compressible)**

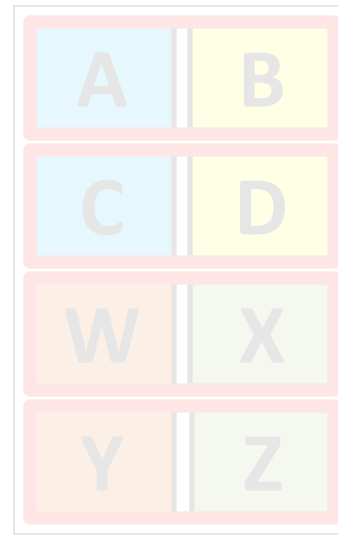**Spatial Indexing (Incompressible)**

11

Compression: Adds capacity, improve bandwidth?



Traditional Compression (Compressible)

Traditional Compression (Incompressible)

Spatial Indexing (Compressible)

**Spatial Indexing (Incompressible)**
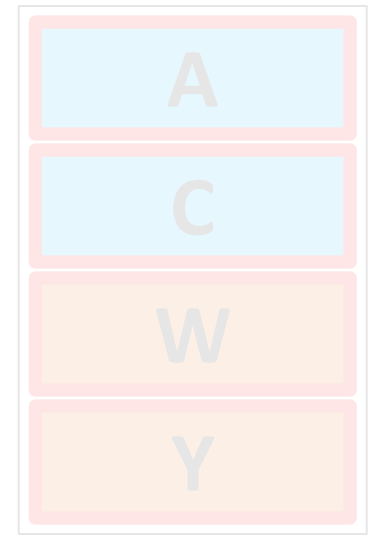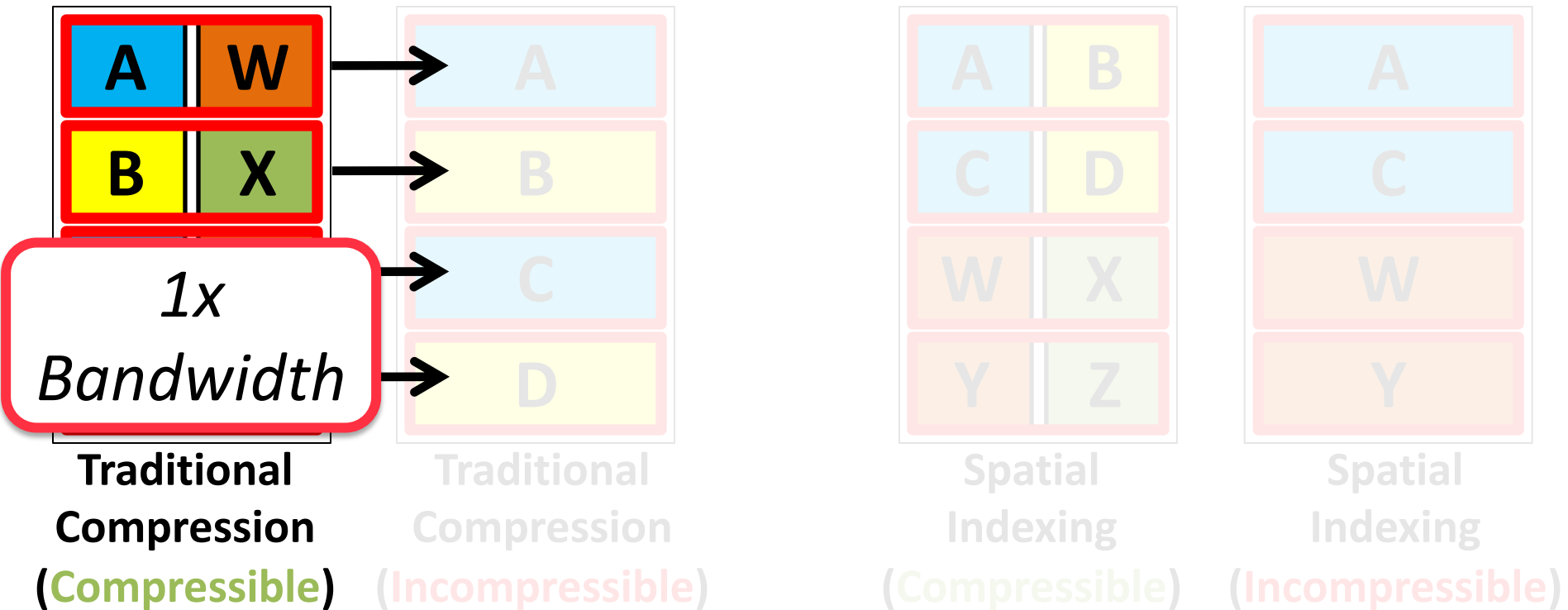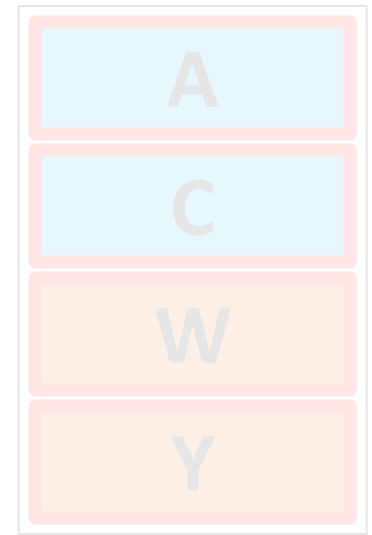
Compression: Adds capacity, improve bandwidth?



Traditional
Compression
(Compressible)

Traditional
Compression
(Incompressible)

Spatial
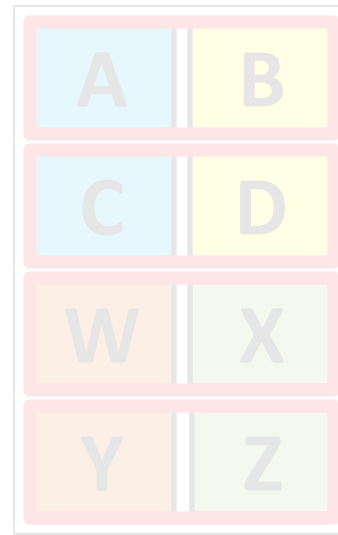Indexing
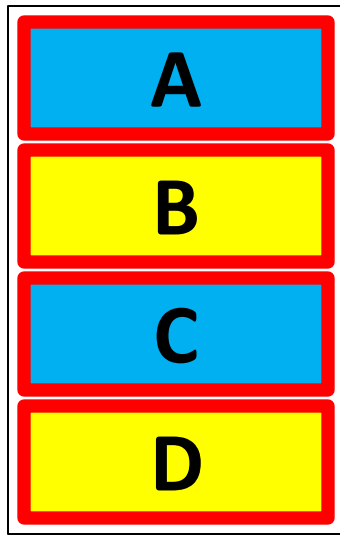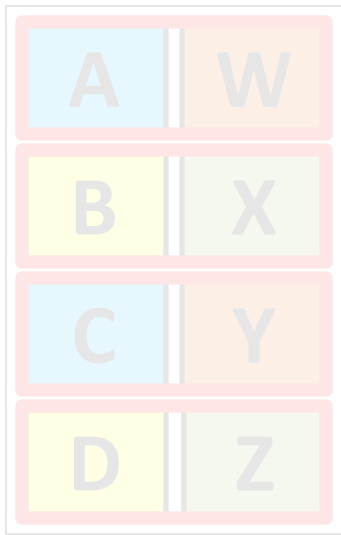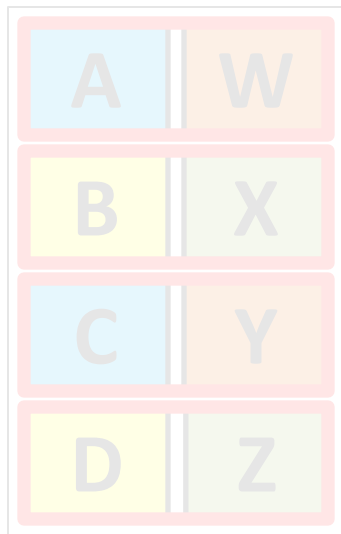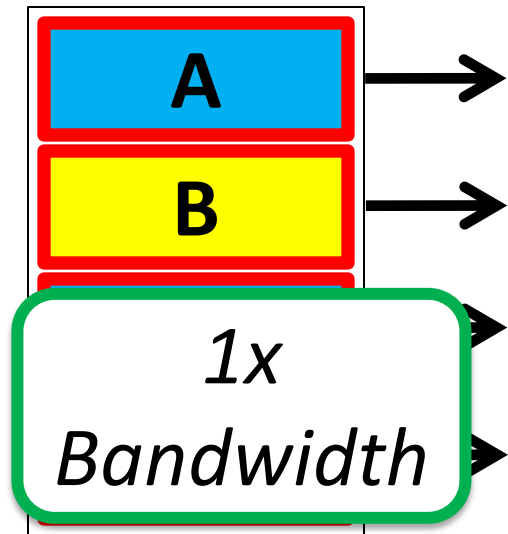(Compressible)

Spatial
Indexing
(Incompressible)

< 1x Bandwidth

B,D?

A

C

Compression: Adds capacity, improve bandwidth?

Traditional Compression

Spatial Indexing

| *1x Bandwidth* | *1x Bandwidth* | *2x Bandwidth* | *< 1x Bandwidth* |
|---|---|---|---|
| **Compressible** | **Incompressible** | **Compressible** | **Incompressible** |

Compression for capacity (TSI) sees little speedup (7%) due to diminishing returns on giga-scale caches

Compression for capacity (TSI) sees little speedup (7%) due to diminishing returns on giga-scale caches

Compression for capacity (TSI) sees little speedup (7%) due to diminishing returns on giga-scale caches

Spatial Indexing compression gets both benefits of bandwidth and capacity when lines are compressible. But, it hurts performance when lines are incompressible

Spatial Indexing compression gets both benefits of bandwidth and capacity when lines are compressible. But, it hurts performance when lines are incompressible

Spatial Indexing compression gets both benefits of bandwidth and capacity when lines are compressible. But, it hurts performance when lines are incompressible

Goal: Compression for Capacity **AND** Bandwidth



Traditional Compression

Spatial Indexing

*1x Bandwidth*

*2x Bandwidth*

*1x Bandwidth*

*< 1x Bandwidth*

**Compressible**

**Incompressible**

**Compressible**

**Incompressible**

DICE (Dynamic Index) → 19% Speedup + 36% ⬇ EDP

- Compressed DRAM Cache Organization ⬅

- Flexible Mapping for Quick Switching

- Dynamic Indexing ComprEssion (DICE)
  - Insertion Policy
  - Index Prediction

**Compression: Simple changes within the controller**

# PRACTICAL DRAM CACHE COMPRESSION



Compression: Simple changes within the controller

# DRAM CACHE TAG FORMAT

**Tag Boundary**

**Data**

**8 Bytes**

**64 Bytes**

| Tag A | Data A |

Cache controller receives 72B of tag+data. It can flexibly interpret bits as tag bits or data bits.

# PROPOSED FLEXIBLE TAG FORMAT

**Tag Boundary**

**Data**

Is Tag?

| A | | A |

We create Tag space as needed, for up to 28 lines.
Achieves 1.6x effective capacity.

# PROPOSED FLEXIBLE TAG FORMAT

Tag Boundary

Data

Is Tag?

| A | B | | B | A |
|---|---|---|---|---|

We create Tag space as needed, for up to 28 lines. Achieves 1.6x effective capacity.

# PROPOSED FLEXIBLE TAG FORMAT

**Tag Boundary**

**Data**

Is Tag?

| A | B | X | | X | B | A |

We create Tag space as needed, for up to 28 lines.
Achieves 1.6x effective capacity.

# PROPOSED FLEXIBLE TAG FORMAT

**Tag Boundary**

**Data**

Not Tag

| A | B | X | I | X | B | A |
|---|---|---|---|---|---|---|

We create Tag space as needed, for up to 28 lines.
Achieves 1.6x effective capacity.

# DICE OVERVIEW

- Compressed DRAM Cache Organization

- Flexible Mapping for Quick Switching ⬅

- Dynamic Indexing ComprEssion (DICE)
  - Insertion Policy
  - Index Prediction

# FLEXIBLE MAPPING (*TSI* OR *BAI*)



**Traditional Set Indexing (TSI)**

**Naïve Spatial Indexing**

**Bandwidth-Aware Indexing (BAI)**

Bandwidth-Aware Indexing (BAI) facilitates quick switching between two indices TSI and BAI.

22

# FLEXIBLE MAPPING (*TSI* OR *BAI*)



**Traditional Set Indexing (TSI)**

**Naïve Spatial Indexing**

**Bandwidth-Aware Indexing (BAI)**

Bandwidth-Aware Indexing (BAI) facilitates quick switching between two indices TSI and BAI.

# FLEXIBLE MAPPING (*TSI* OR *BAI*)



**Traditional Set Indexing (TSI)**

**Naïve Spatial Indexing**

**Bandwidth-Aware Indexing (BAI)**

Bandwidth-Aware Indexing (BAI) facilitates quick switching between two indices TSI and BAI.

23

# FLEXIBLE MAPPING (*TSI* OR *BAI*)



**Traditional Set Indexing (TSI)**

**Naïve Spatial Indexing**

**Bandwidth-Aware Indexing (BAI)**

Bandwidth-Aware Indexing (BAI) facilitates quick switching between two indices TSI and BAI.

24

# FLEXIBLE MAPPING (*TSI* OR *BAI*)



**Traditional Set Indexing (TSI)**

**Naïve Spatial Indexing**

**Bandwidth-Aware Indexing (BAI)**

Bandwidth-Aware Indexing (BAI) facilitates quick switching between two indices TSI and BAI.
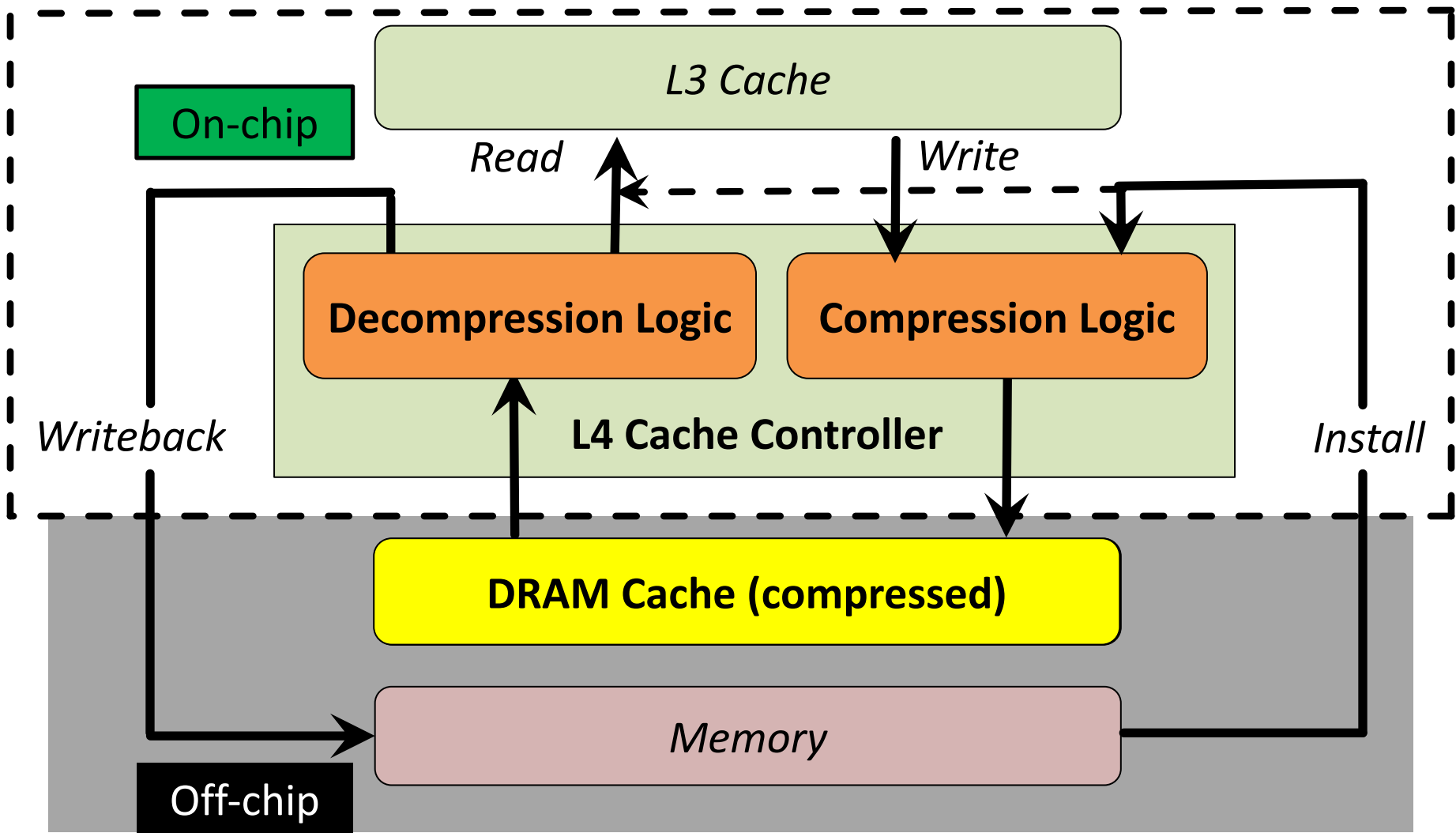
# DICE OVERVIEW

- Compressed DRAM Cache Organization

- Flexible Mapping for Quick Switching

- Dynamic Indexing ComprEssion (DICE)
  - Insertion Policy
  - Index Prediction
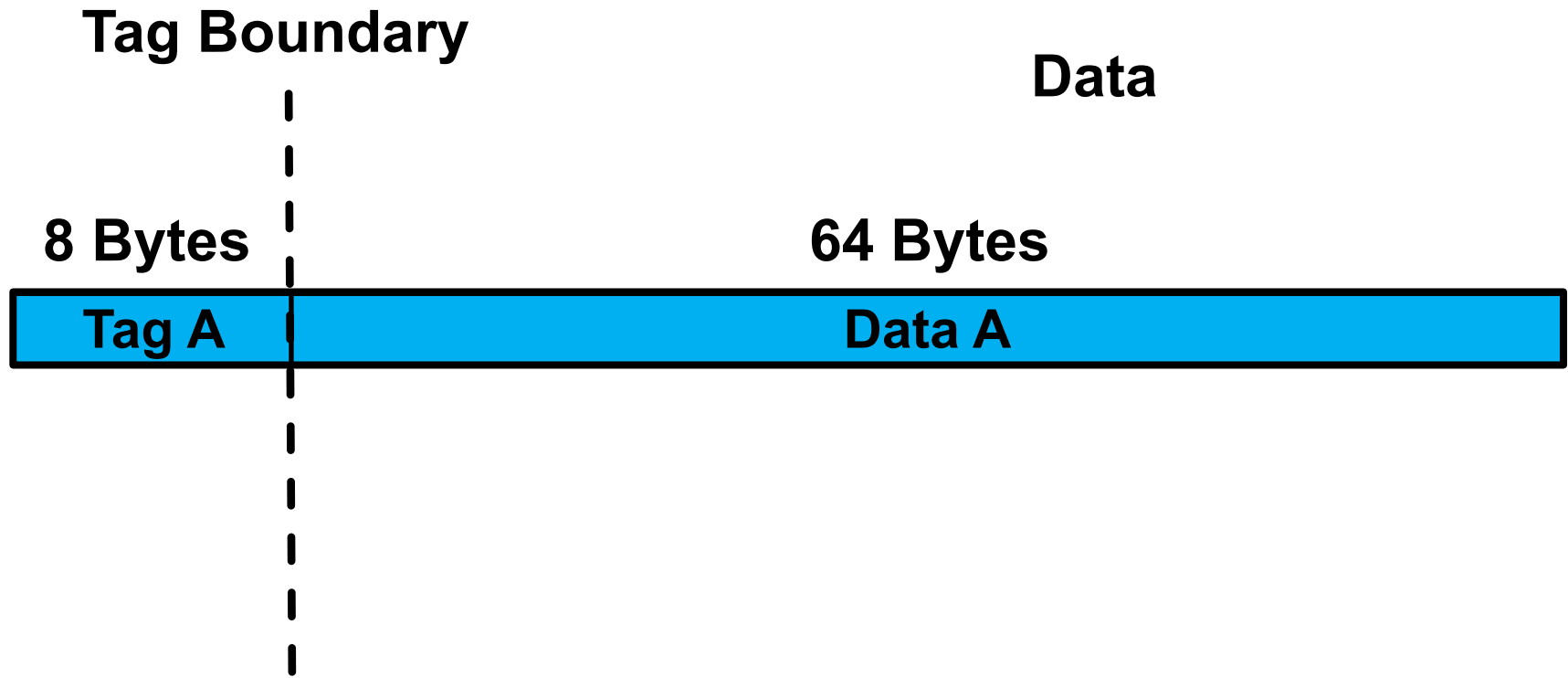
# DICE: DYNAMIC-INDEXED COMPRESSED CACHE

**DRAM Cache**



Compressibility Based Insertion

Cache Index Prediction

Install

Read

Traditional Set Index

Bandwidth-Aware Index

TSI = BAI

DICE: Dynamic-Indexing Cache comprEssion, decides index on install, and predicts index on read

# COMPRESSIBILITY-BASED INSERTION

**DRAM Cache**

Compressibility Based Insertion

Install

<= ½-size

| |
|---|
| Traditional Set Index |
| Bandwidth-Aware Index |

TSI = BAI

Compressibilty-based insertion uses Bandwidth-Aware Indexing when lines are compressible, and TSI otherwise

# COMPRESSIBILITY-BASED INSERTION

**DRAM Cache**

Compressibility Based Insertion

> ½-size

Install

<= ½-size

**Traditional Set Index**

**Bandwidth-Aware Index**

**TSI = BAI**

Compressibilty-based insertion uses Bandwidth-Aware Indexing when lines are compressible, and TSI otherwise

# COMPRESSIBILITY-BASED INSERTION

**DRAM Cache**

Compressibility Based Insertion

Install

> ½-size

<= ½-size

Traditional Set Index

Bandwidth-Aware Index

TSI = BAI

No explicit swaps. Eviction and install decides policy

Compressibilty-based insertion uses Bandwidth-Aware Indexing when lines are compressible, and TSI otherwise

# COMPRESSIBILITY-BASED INSERTION

**DRAM Cache**

Compressibility Based Insertion

**Install**

> ½-size

<= ½-size

**Traditional Set Index**

**Bandwidth-Aware Index**

**TSI = BAI**

But checking both wastes bandwidth

?

**Read**

No explicit swaps. Eviction and install decides policy

Compressibilty-based insertion uses Bandwidth-Aware Indexing when lines are compressible, and TSI otherwise

# SIMILAR INTRA-PAGE COMPRESSIBILITY

**Indices seen in a Compressible Page**



**Install**

<= ½-size

Lines within a page have similar compressibility

**Bandwidth-Aware Index**

**Read BAI**

DICE is likely to install lines of a page into similar index

# SIMILAR INTRA-PAGE COMPRESSIBILITY

**Indices seen in a Compressible Page**

Install

<= ½-size

Lines within a page have similar compressibility

Bandwidth-Aware Index

Bandwidth-Aware Index

Read BAI

DICE is likely to install lines of a page into similar index

# SIMILAR INTRA-PAGE COMPRESSIBILITY

**Indices seen in a Compressible Page**



DICE is likely to install lines of a page into similar index

**Indices seen in a Compressible Page**

Install

<= ½-size

Lines within a page have similar compressibility

Bandwidth-Aware Index

Bandwidth-Aware Index

Bandwidth-Aware Index

Read BAI

DICE is likely to install lines of a page into similar index

# SIMILAR INTRA-PAGE COMPRESSIBILITY

**Indices seen in an Incompressible Page**

**Traditional Set Index**

**Install**

> ½-size

Lines within a page have similar compressibility

**Read TSI**

Thus, page-based last-time prediction of index can be accurate (94%)

# SIMILAR INTRA-PAGE COMPRESSIBILITY

**Indices seen in an Incompressible Page**

**Install**

> ½-size

Lines within a page have similar compressibility

**Traditional Set Index**

**Traditional Set Index**

**Read TSI**

Thus, page-based last-time prediction of index can be accurate (94%)

# SIMILAR INTRA-PAGE COMPRESSIBILITY

**Indices seen in an Incompressible Page**

**Install**

> ½-size

Lines within a page have similar compressibility

| Traditional Set Index |
| Traditional Set Index |
| Traditional Set Index |

**Read TSI**

Thus, page-based last-time prediction of index can be accurate (94%)

# SIMILAR INTRA-PAGE COMPRESSIBILITY

**Indices seen in an Incompressible Page**



Install

> ½-size

Lines within a page have similar compressibility

Traditional Set Index

Traditional Set Index

Traditional Set Index

Traditional Set Index

Read TSI

Thus, page-based last-time prediction of index can be accurate (94%)

**Indices seen in an Incompressible Page**

Install

> ½-size

Lines within a page have similar compressibility

Traditional Set Index

Traditional Set Index

Traditional Set Index

Traditional Set Index

**Bandwidth-Aware Index**

Read TSI

Thus, page-based last-time prediction of index can be accurate (94%)

# SIMILAR INTRA-PAGE COMPRESSIBILITY

**Indices seen in an Incompressible Page**



Install

> ½-size

Lines within a page have similar compressibility

Traditional Set Index
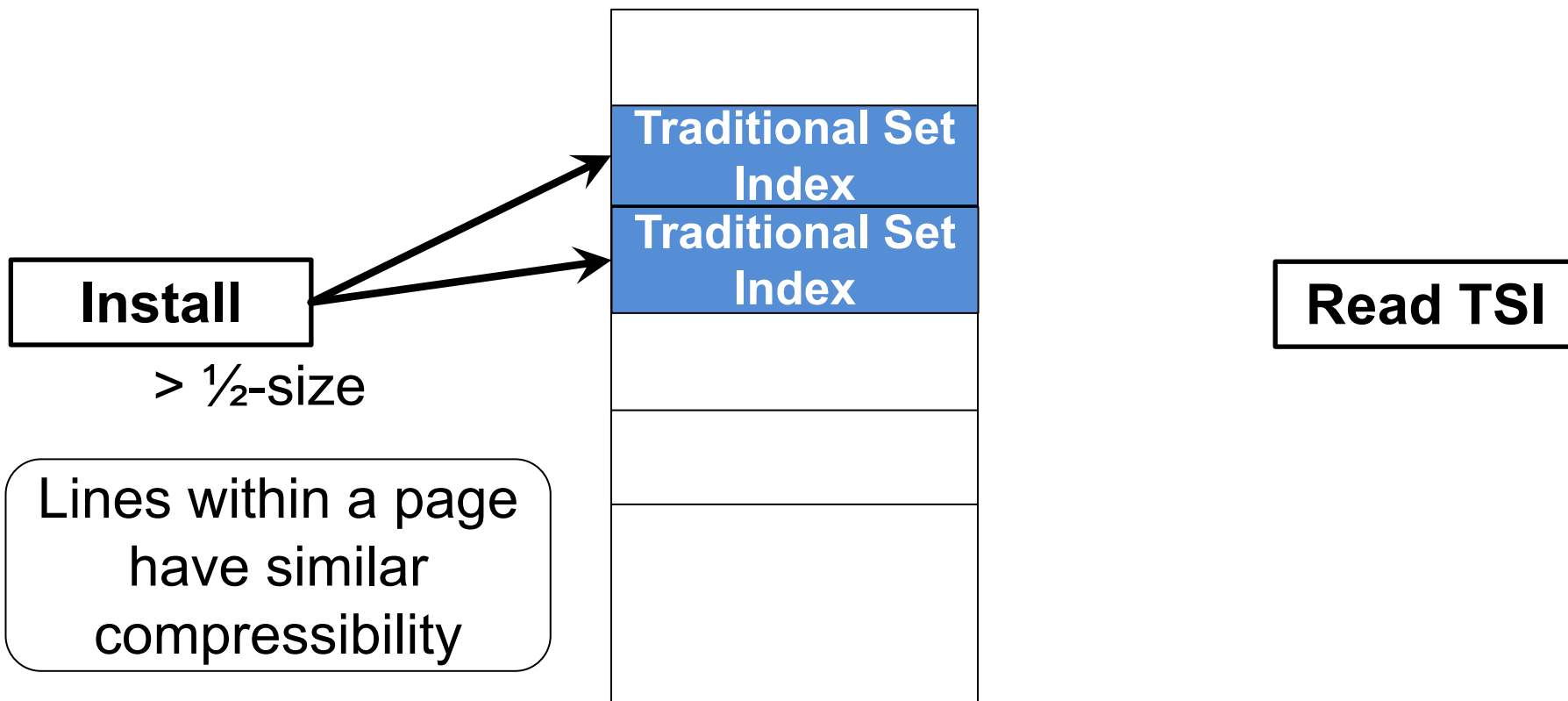
Traditional Set Index

Traditional Set Index

Traditional Set Index

Bandwidth-Aware Index

Read TSI

2nd access only on mispredict

Thus, page-based last-time prediction of index can be accurate (94%)

# PAGE-BASED CACHE INDEX PREDICTOR (CIP)



**Demand Access**

Page # → Hash

**Last-Time Table (LTT)**

| |
|---|
| 1 |
| 1 |
| 0 |
| 0 |
| 1 |
| 0 |
| 1 |

0 = **Traditional Set Index**
1 = **Bandwidth-Aware Index**

**Predict**
**Traditional Set Index**

Page-based last-time prediction exploits similar intra-page compressibility, to achieve high prediction accuracy (94%)

# DICE OVERVIEW

- Compressed DRAM Cache Organization

- Flexible Mapping for Quick Switching

- Dynamic Indexing (DICE)
  - Insertion Policy
  - Index Prediction

- Results ⬅

CPU

Stacked
DRAM

Commodity
DRAM

- Core Chip
  - 3.2GHz 4-wide out-of-order core
  - 8 cores, 8MB shared last-level cache
- Compression
  - FPC + BDI

# METHODOLOGY (1/8<sup>TH</sup> KNIGHTS LANDING)

Other sensitivities in paper

CPU

Stacked DRAM

Commodity DRAM

| | Stacked DRAM | Commodity DRAM |
|---|---|---|
| Capacity | 1GB | 32GB |
| Bus | DDR1.6GHz, 128-bit | DDR1.6GHz, 64-bit |
| Channels | 4 channels | 1 channel |
| Bandwidth | 100 GBps | 12.5 GBps |
| Latency | 35ns | 35ns |

# DICE RESULTS



DICE improves performance over both Spatial Indexing and Traditional Indexing with fine-grain decision (19%)

# DICE RESULTS



DICE improves performance over both Spatial Indexing and Traditional Indexing with fine-grain decision (19%)

# DICE RESULTS



DICE improves performance over both Spatial Indexing and Traditional Indexing with fine-grain decision (19%)

Goal: Compression for Capacity **AND** Bandwidth



**Traditional Compression**

**Spatial Indexing**

*1x Bandwidth*

*2x Bandwidth*

*1x Bandwidth*

*< 1x Bandwidth*

Compressible

**Incompressible**

**Compressible**

Incompressible

DICE (Dynamic Index) → 19% Speedup + 36% ⬇ EDP

# THANK YOU

# EXTRA SLIDES

- Extra Slides

# DIFFERENT CACHE SENSITIVITIES

**Table 8: Sensitivity of DICE on different caches**

|  | Base(1GB) | 2x Capacity | 2x BW | 50% Latency |
|---|---|---|---|---|
| SPEC RATE | +12.2% | +8.7% | +13.3% | +13.5% |
| SPEC MIX | +7.5% | +4.7% | +8.2% | +9.1% |
| GAP | +48.9% | +32.6% | +75.9% | +73.5% |
| GMEAN26 | +19.0% | +13.2% | +24.5% | +24.4% |

# COMPARISON TO PREFETCH

**Table 7: Comparison of DICE to Prefetch**

|  | 128B-PF | Nextline-PF | DICE | DICE+NL |
|---|---|---|---|---|
| SPEC RATE | +3.2% | +2.6% | +12.2% | +16.7% |
| SPEC MIX | +1.2% | +1.9% | +7.5% | +7.7% |
| GAP | -1.1% | -1.1% | +48.9% | +43.4% |
| GMEAN26 | +1.9% | +1.6% | +19.0% | +20.9% |

## Table 1: Comparison of different forms of compression

| Module to Compress | Improve Capacity Only? | Tag Overhead? | OS support Needed? |
|---|---|---|---|
| On-Chip Cache | Yes | Yes | No |
| Main Memory | No | No | Yes |
| **DRAM Cache** | No | No | No |

# DISTRIBUTION FOR INDEX DECISION

# DICE INSERTION THRESHOLD

## Table 4: Sensitivity to DICE threshold

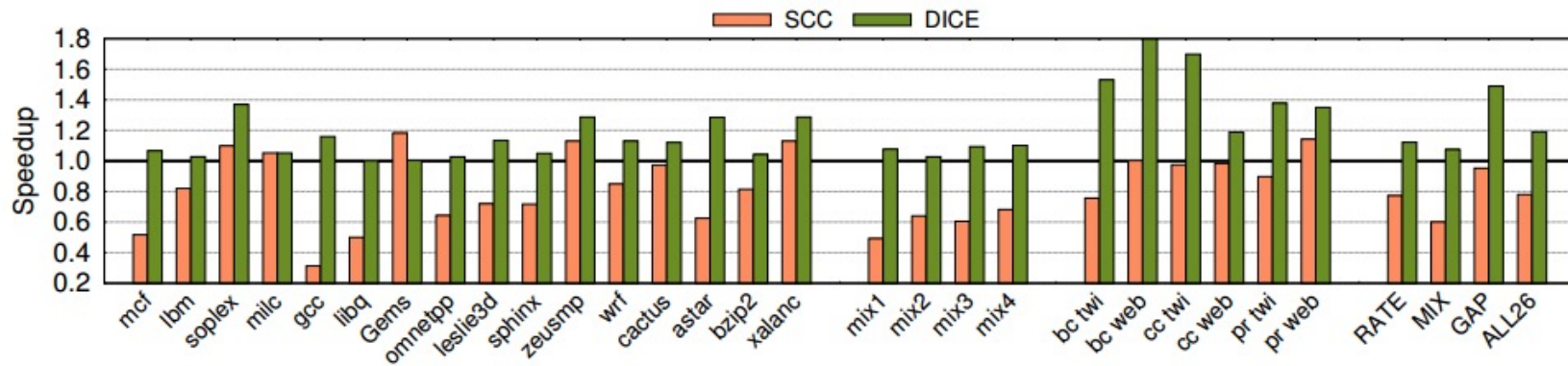|            | $\leq 32B$ | $\leq 36B$ | $\leq 40B$ |
|------------|-----------|-----------|-----------|
| SPEC RATE  | +10.6%    | +12.2%    | +11.1%    |
| SPEC MIX   | +6.4%     | +7.5%     | +7.4%     |
| GAP        | +47.6%    | +48.9%    | +49.1%    |
| GMEAN26    | +17.5%    | +19.0%    | +18.3%    |

# EFFECTIVE CAPACITY

**Table 5: Effective Capacity of TSI/BAI/DICE**

|  | TSI | BAI | DICE |
|---|---|---|---|
| SPEC RATE | 1.07x | 1.16x | 1.13x |
| SPEC MIX | 1.12x | 1.28x | 1.24x |
| GAP | 2.00x | 5.57x | 5.06x |
| GMEAN26 | 1.24x | 1.69x | 1.62x |

# L3 HIT RATE IMPROVEMENT

## Table 6: Effect of DICE on L3 hit rate

|  | BASE | DICE |
|---|---|---|
| SPEC RATE | 34.7% | 43.0% |
| SPEC MIX | 61.6% | 67.2% |
| GAP | 26.9% | 29.4% |
| AVG26 | 37.0% | 43.6% |

(a) TSI

| Set 0 | A0, A8 |
| Set 1 | A1, A9 |
| Set 2 | A2, A10 |
| Set 3 | A3, A11 |
| Set 4 | A4, A12 |
| Set 5 | A5, A13 |
| Set 6 | A6, A14 |
| Set 7 | A7, A15 |

(b) NSI

| Set 0 | **A0**, A1 |
| Set 1 | A2, A3 |
| Set 2 | A4, A5 |
| Set 3 | A6, A7 |
| Set 4 | A8, A9 |
| Set 5 | A10, A11 |
| Set 6 | A12, A13 |
| Set 7 | A14, **A15** |

(c) BAI

| Set 0 | **A0**, A1, |
| Set 1 | A8, **A9** |
| Set 2 | **A2**, A3 |
| Set 3 | A10, **A11** |
| Set 4 | **A4**, A5 |
| Set 5 | A12, **A13** |
| Set 6 | **A6**, A7 |
| Set 7 | A14, **A15** |