

Doppelgänger: A Cache for Approximate Computing

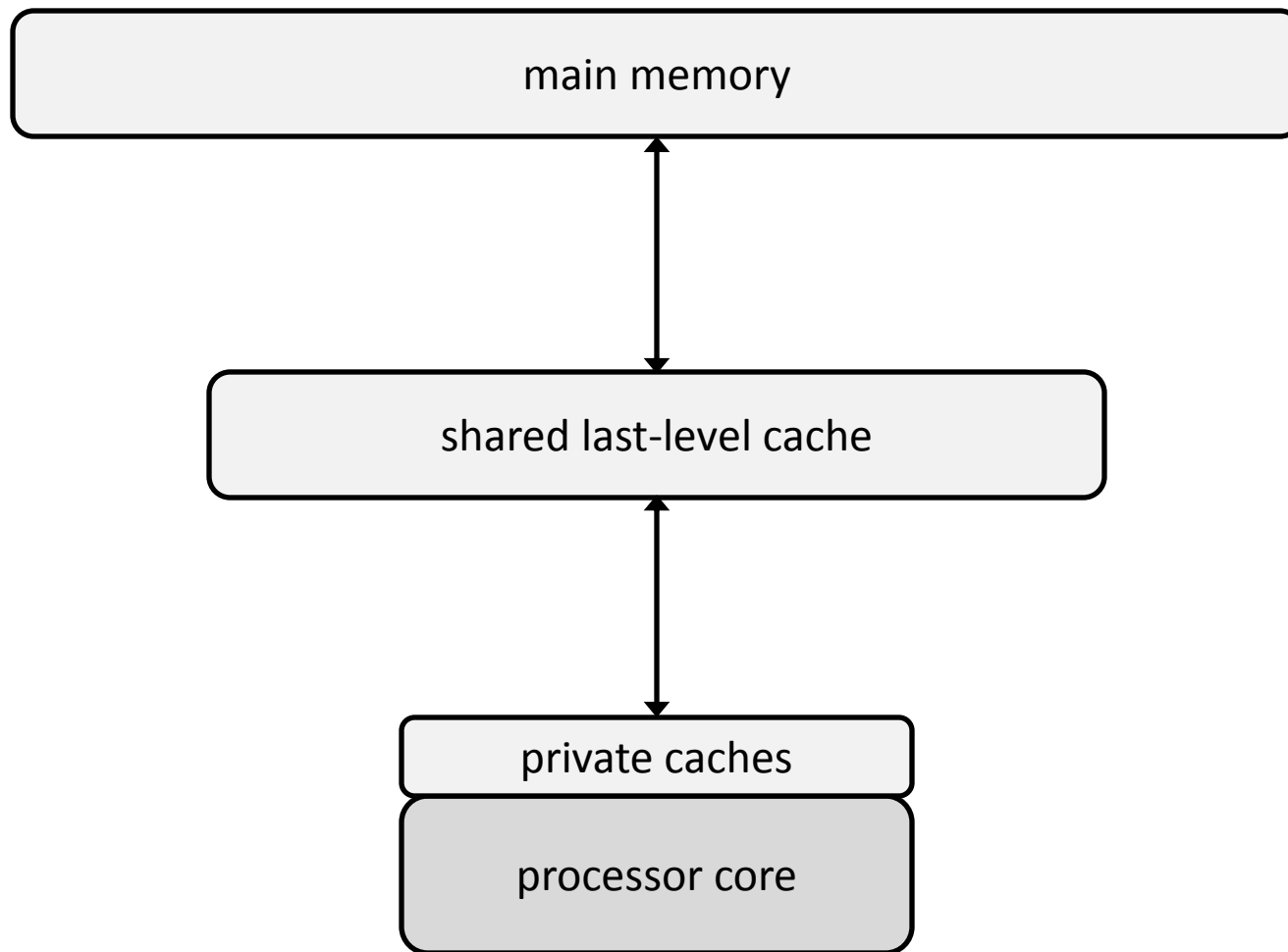
Joshua San Miguel

Jorge Albericio

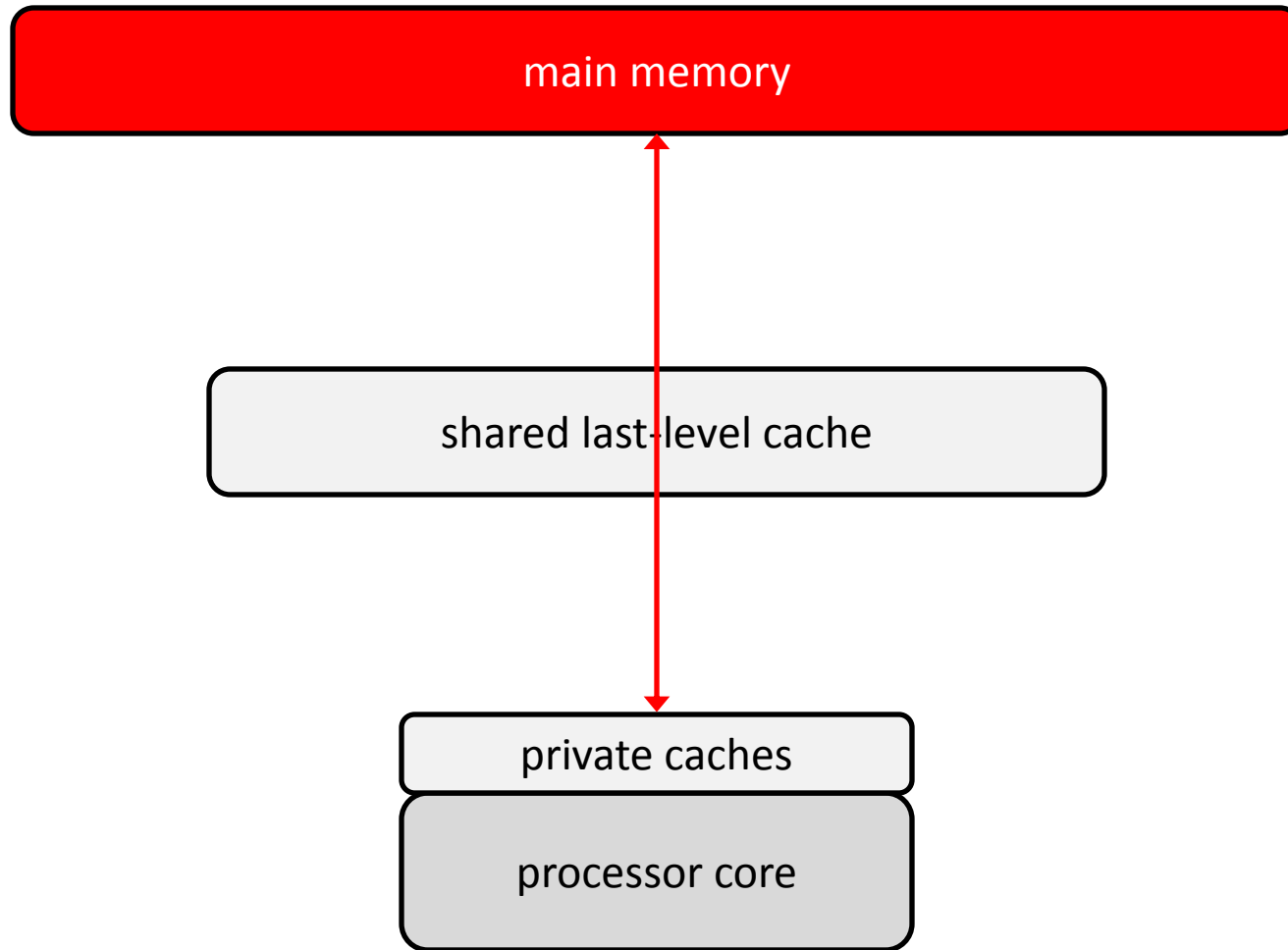
Andreas Moshovos

Natalie Enright Jerger

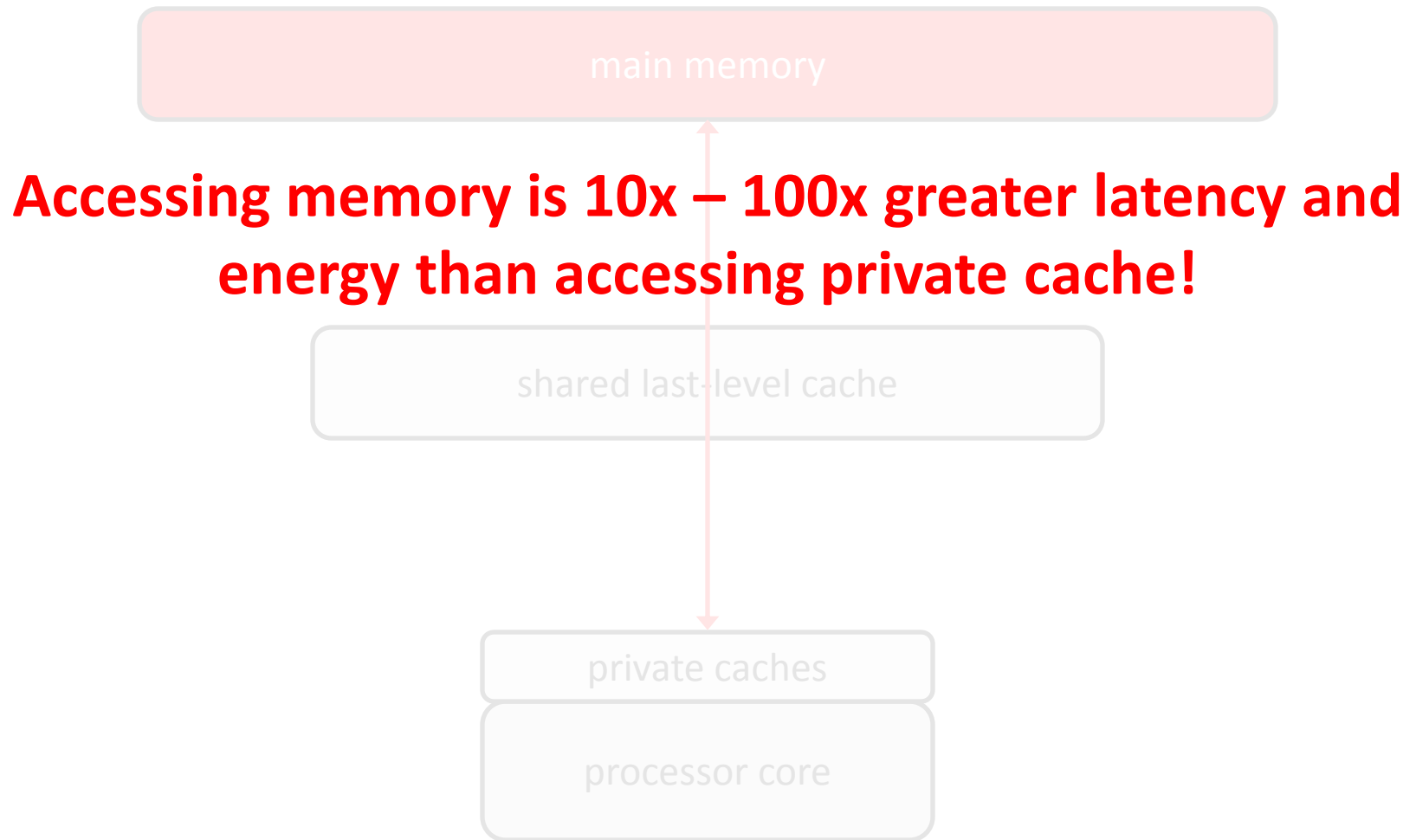
Cache Hierarchy



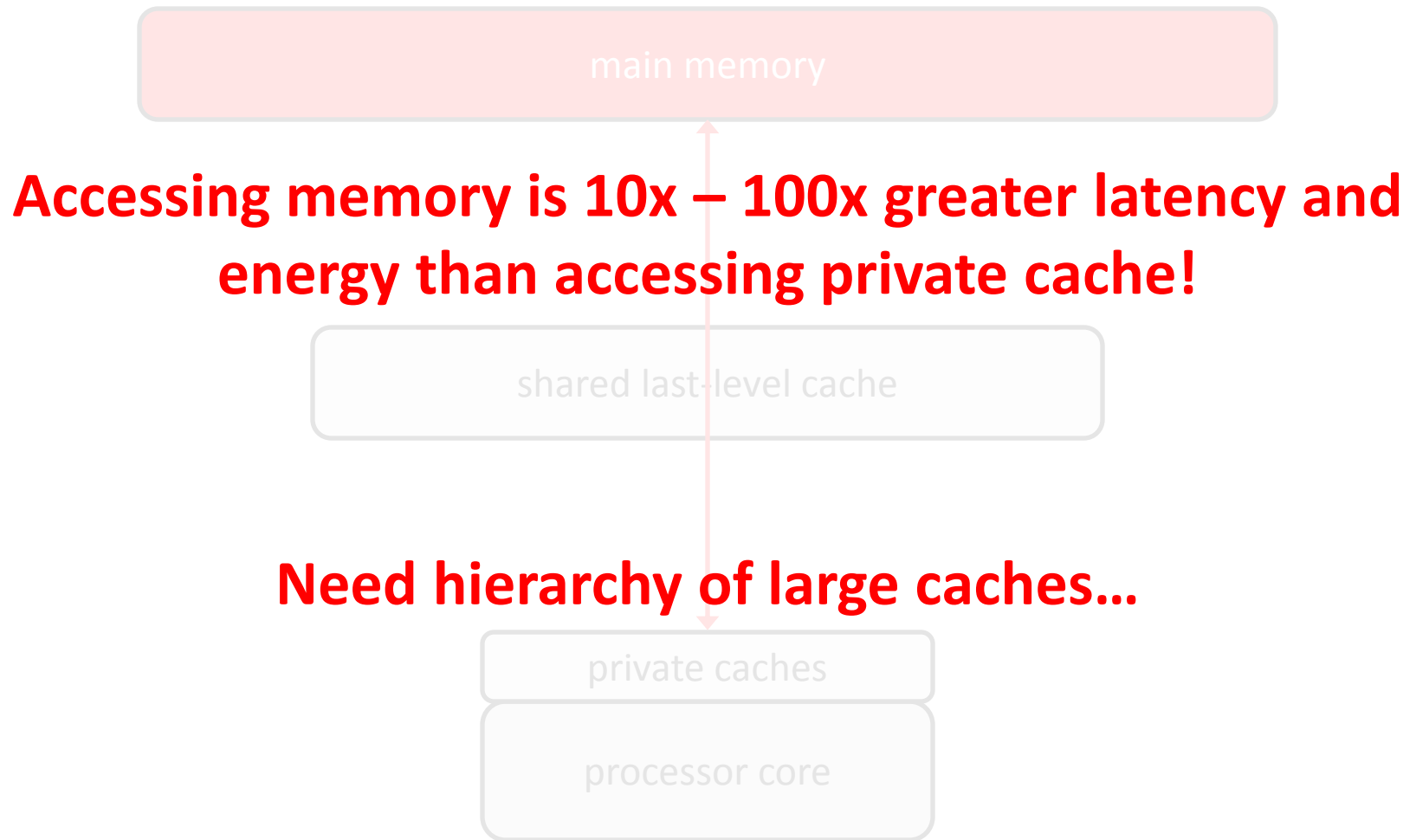
Cache Hierarchy



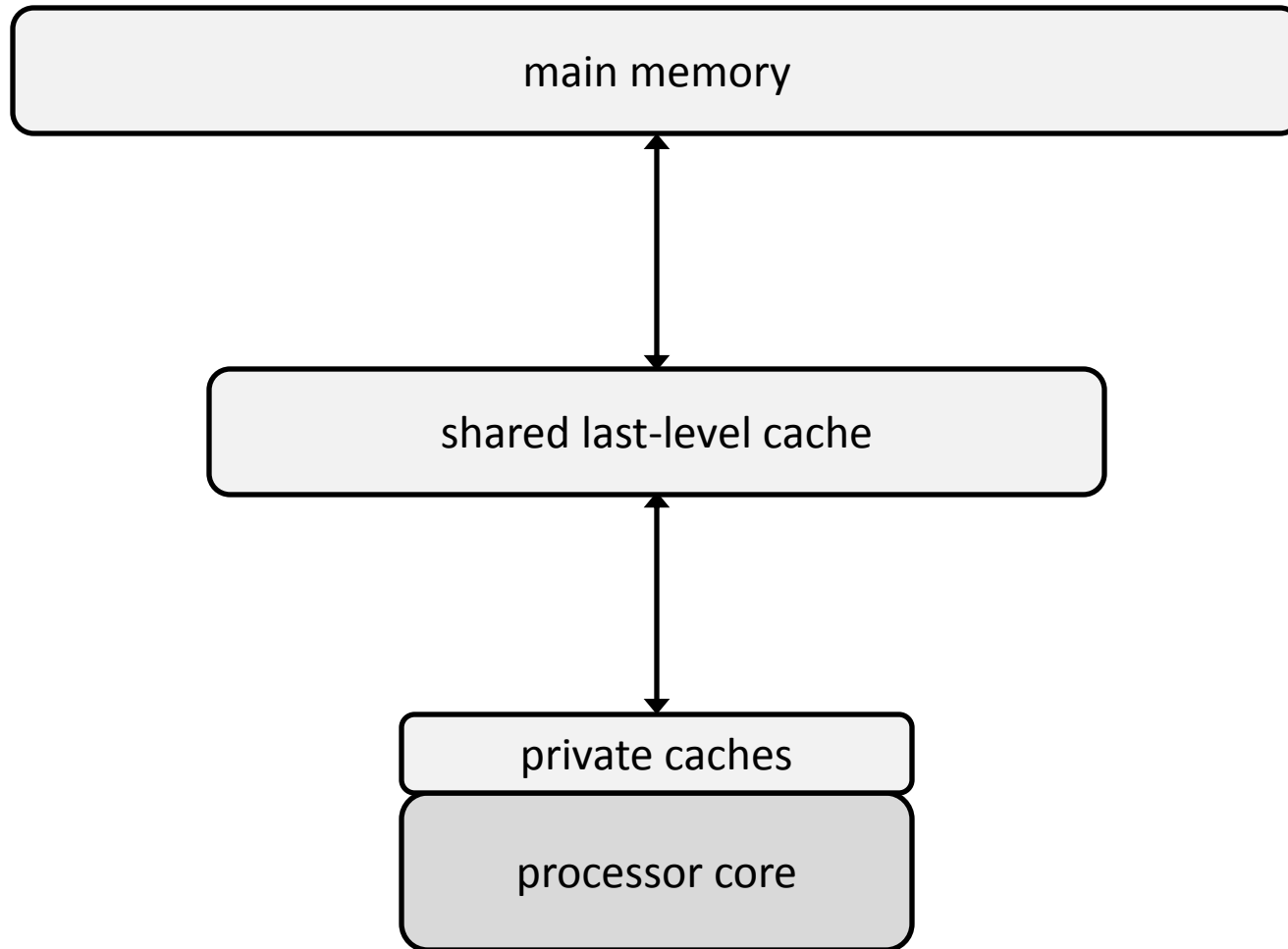
Cache Hierarchy



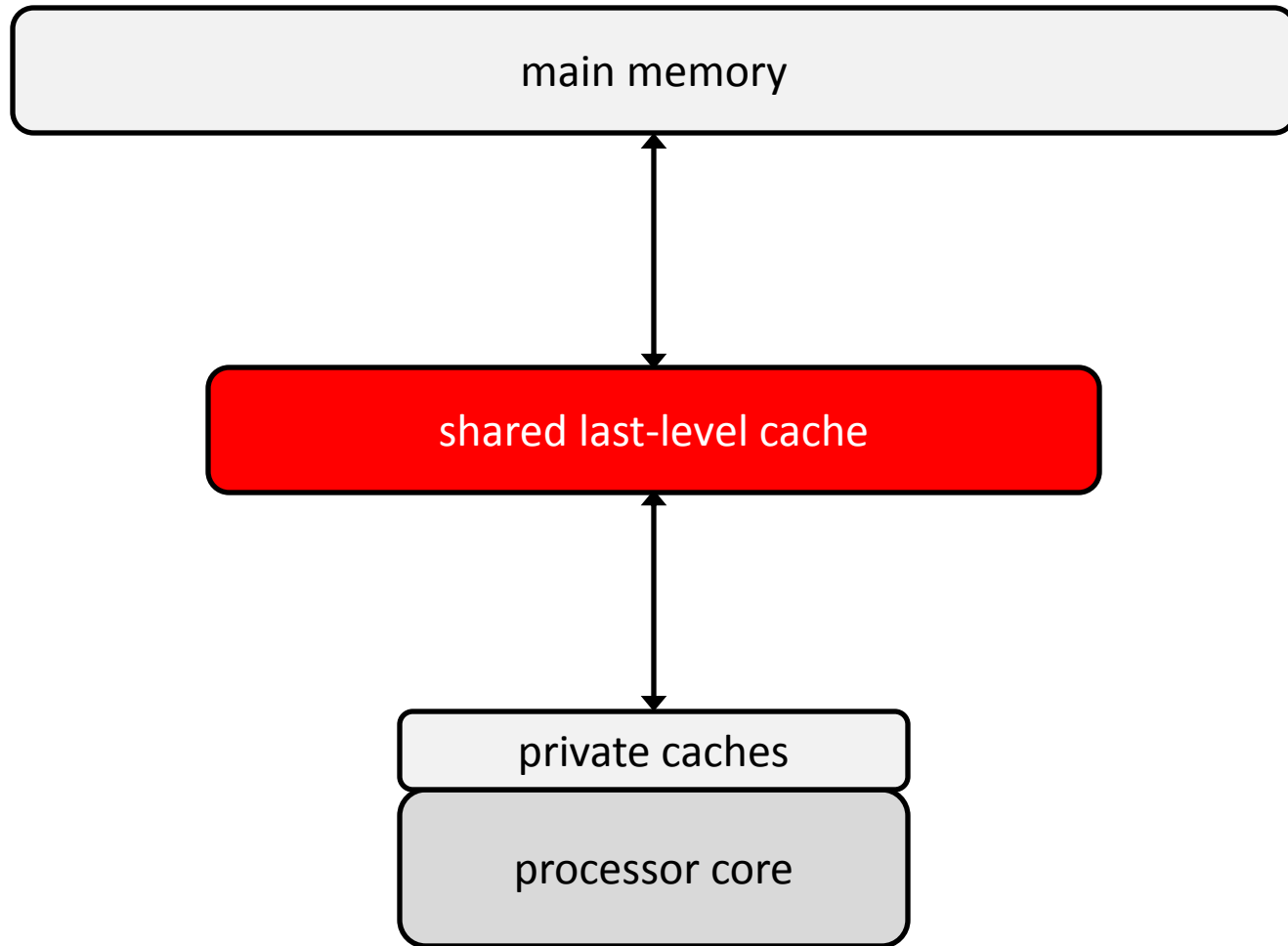
Cache Hierarchy



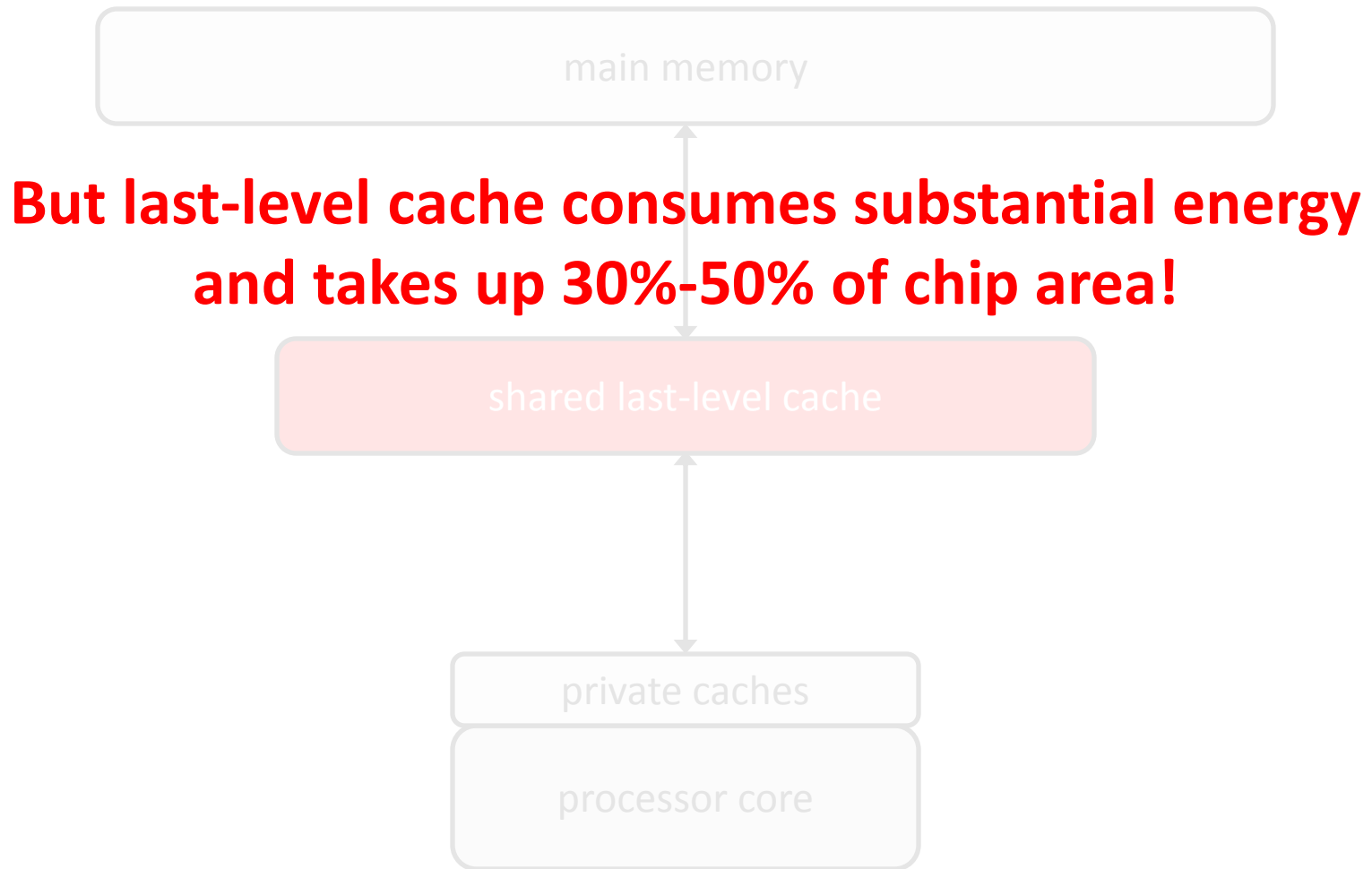
Cache Hierarchy



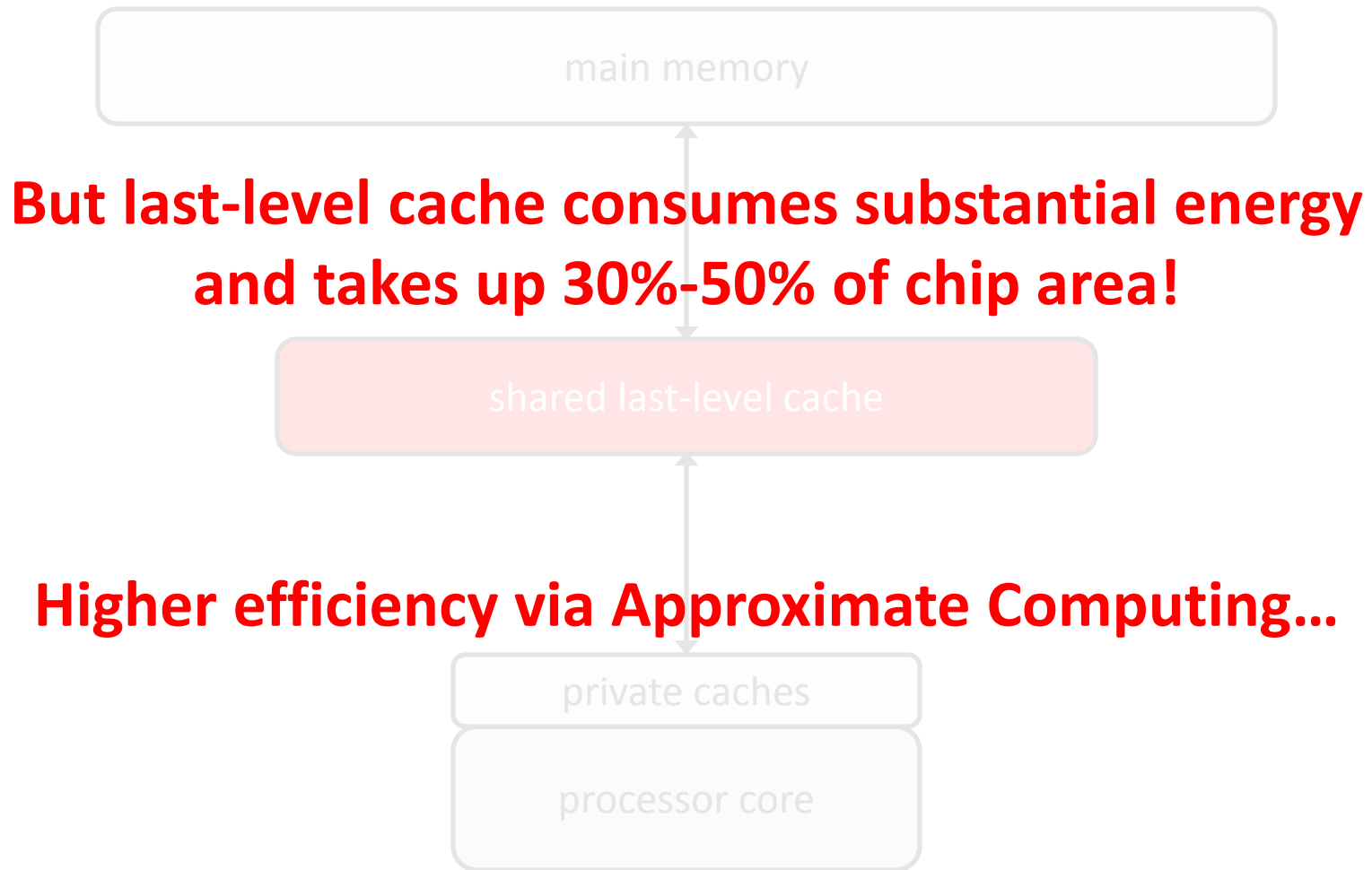
Cache Hierarchy



Cache Hierarchy



Cache Hierarchy



Summary

Doppelgänger Cache:

- Identifies **approximate similarity** in data block values.
 - **77%** cache storage savings of approximable data.

Summary

Doppelgänger Cache:

- Identifies **approximate similarity** in data block values.
 - **77%** cache storage savings of approximable data.
- Effectively compresses storage of approximately similar blocks.
 - **3x** better compression ratio than state-of-the-art techniques.

Summary

Doppelgänger Cache:

- Identifies **approximate similarity** in data block values.
 - **77%** cache storage savings of approximable data.
- Effectively compresses storage of approximately similar blocks.
 - **3x** better compression ratio than state-of-the-art techniques.
- Significantly reduces area and energy consumption.
 - Reduces total on-chip cache area by **1.36x**.

Outline

- Approximate Computing
 - Approximate Similarity
- Doppelgänger Cache
 - Cache Architecture
 - Similarity Mapping
- Evaluation

Approximate Computing

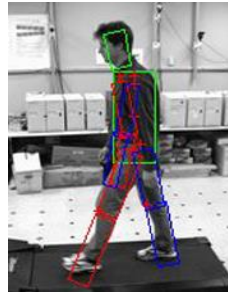
Not all data/computations need to be precise.

Data mining



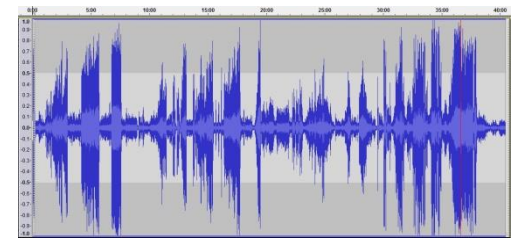
<http://www.zentut.com/>

Computer vision



<http://www.cc.gatech.edu/~cnieto6/>

Audio and video processing



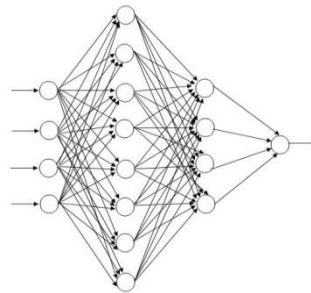
<http://themicparlour.blogspot.ca/>

Gaming



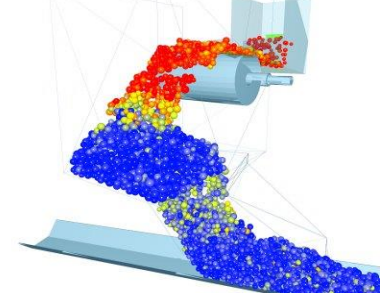
<http://www.businessweek.com/>

Machine learning



<http://www.analyticbridge.com/>

Dynamical simulation



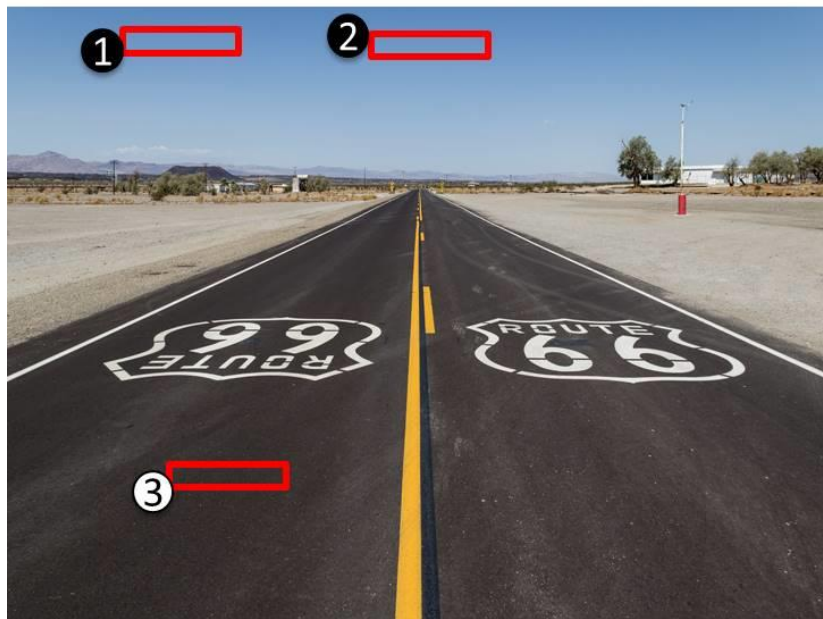
<http://www.scientific-computing.com/>

Approximate Similarity

Two data blocks are ***approximately similar*** (i.e., ***doppelgängers***) if replacing the values of one with the other still results in acceptable application output in the end.

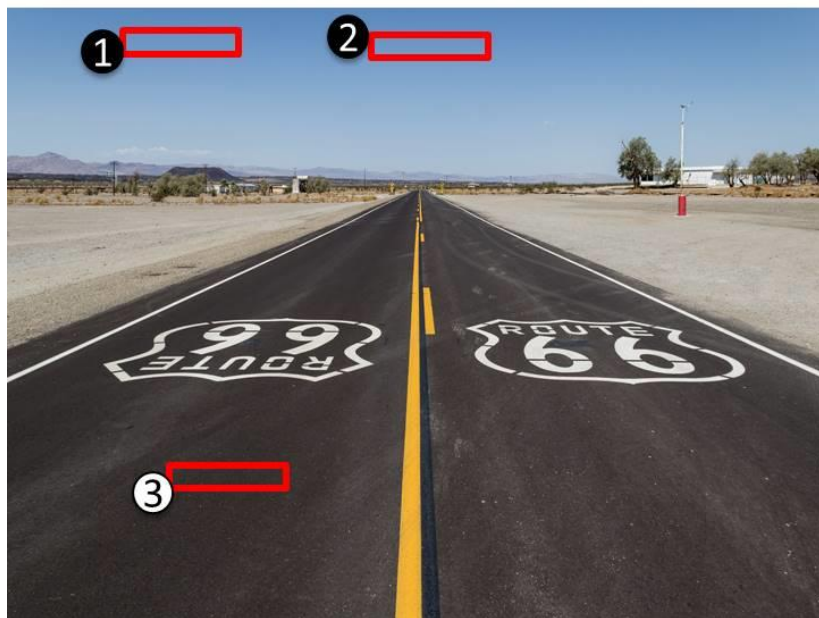
Approximate Similarity

Two data blocks are *approximately similar* (i.e., *doppelgängers*) if replacing the values of one with the other still results in acceptable application output in the end.



Approximate Similarity

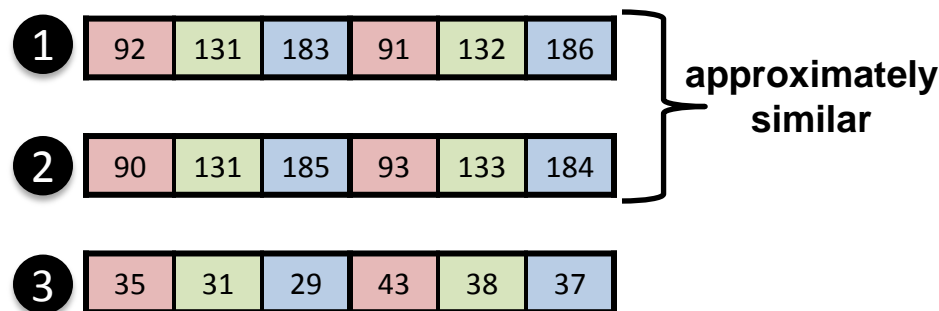
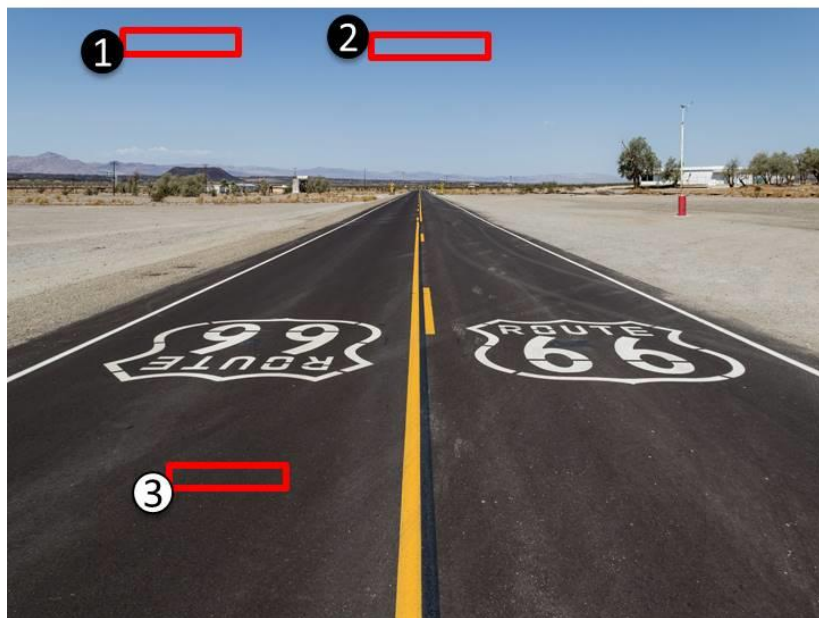
Two data blocks are *approximately similar* (i.e., *doppelgängers*) if replacing the values of one with the other still results in acceptable application output in the end.



1	92	131	183	91	132	186
2	90	131	185	93	133	184
3	35	31	29	43	38	37

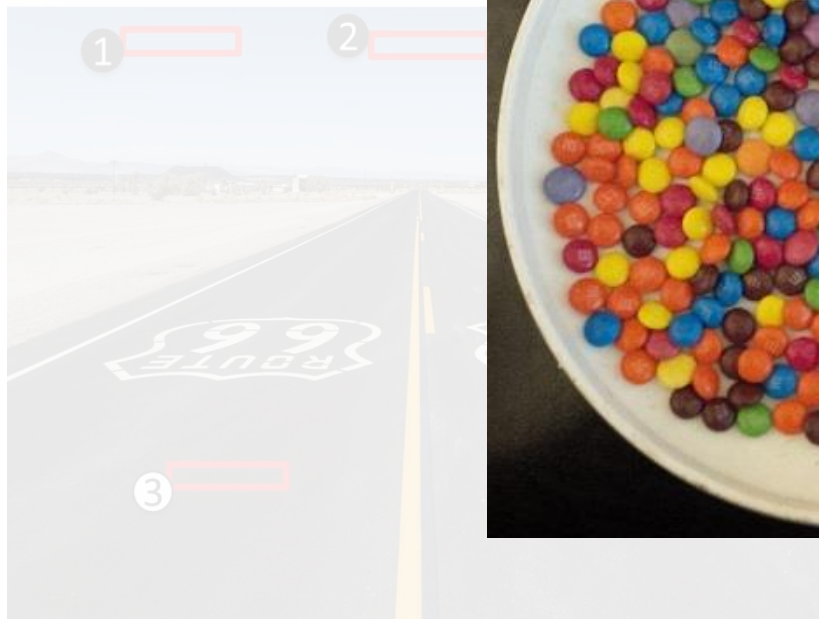
Approximate Similarity

Two data blocks are *approximately similar* (i.e., *doppelgängers*) if replacing the values of one with the other still results in acceptable application output in the end.



Approximate Similarity

Two data blocks are *approximately similar* (i.e., *doppelgänger*) if replacing the values in one block with the values in the other results in an acceptable application.

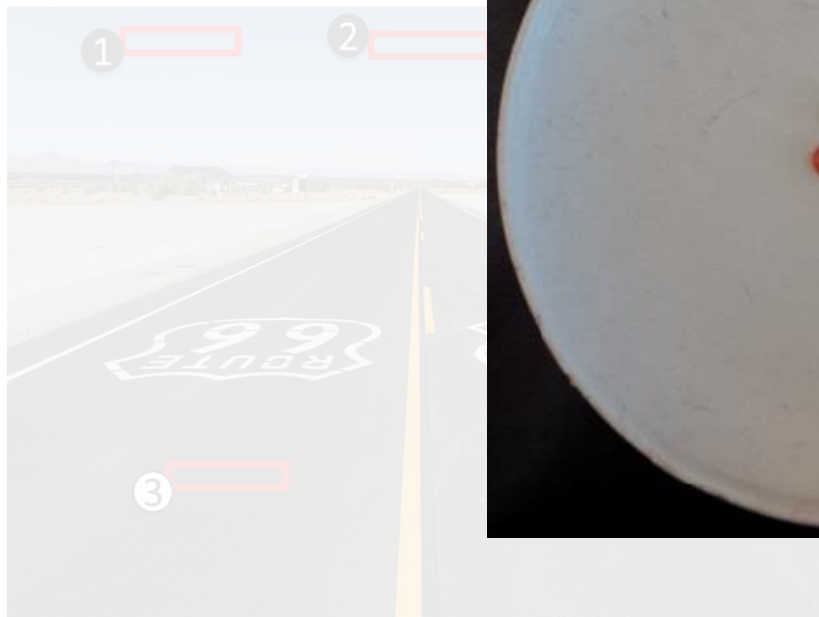


91	132	186
93	133	184
43	38	37

approximately similar

Approximate Similarity

Two data blocks are *approximately similar* (i.e., *doppelgänger*) if replacing the values in one block with the values in the other results in an acceptable application.



91	132	186
----	-----	-----

93	133	184
----	-----	-----

43	38	37
----	----	----

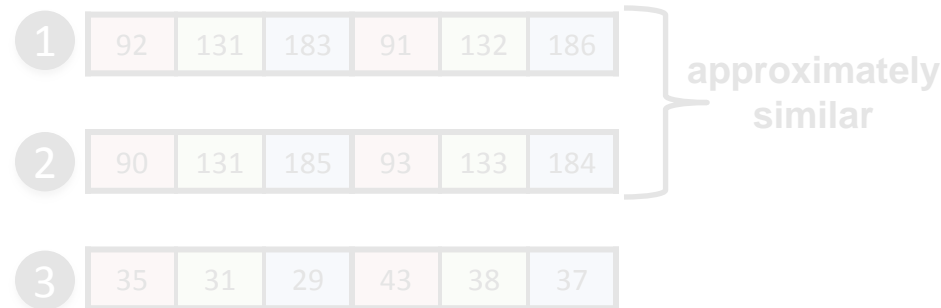
approximately similar

Approximate Similarity

Two data blocks are *approximately similar* (i.e., *doppelgängers*) if replacing the values of one with the other still results in acceptable application output in the end.



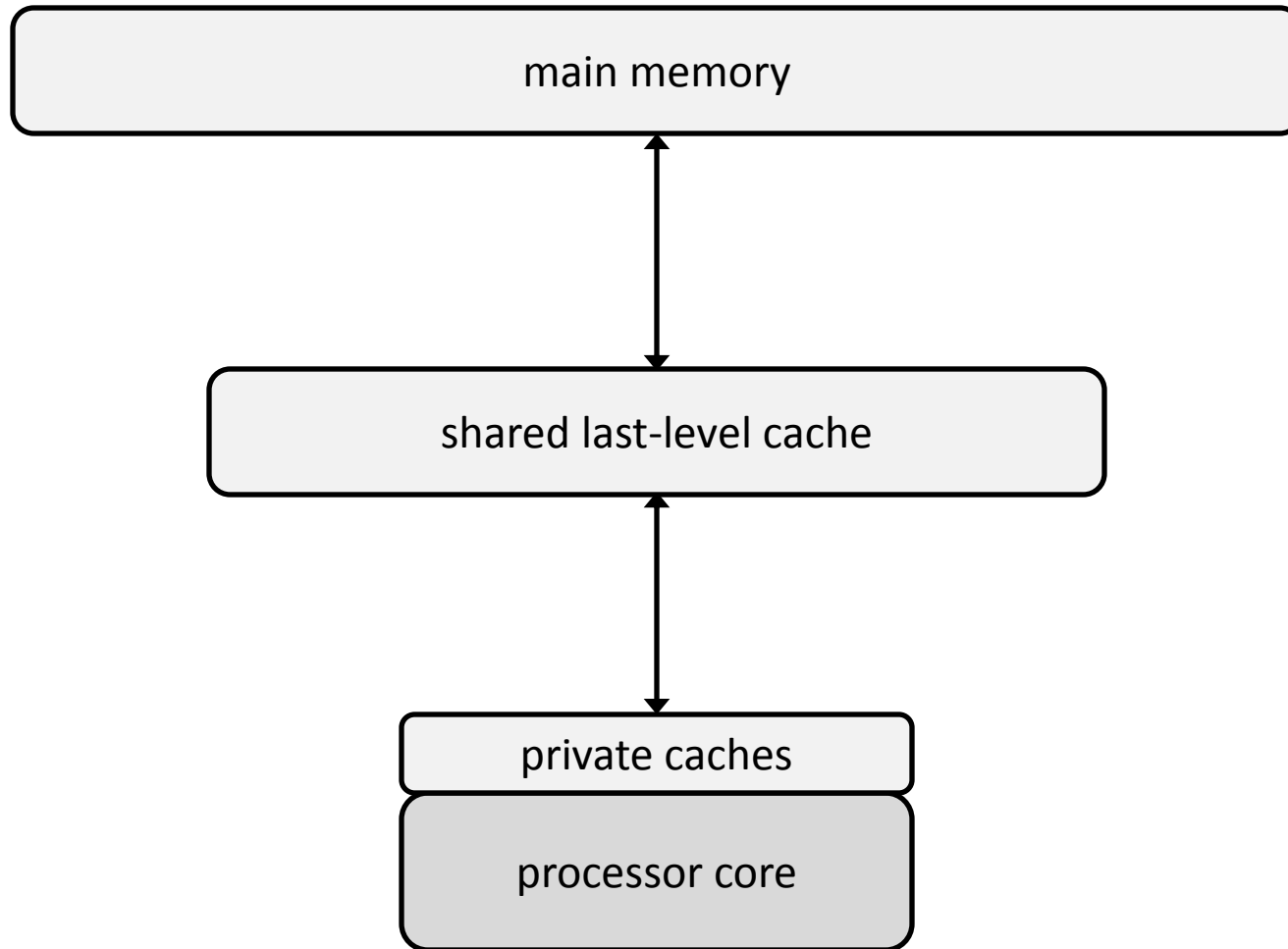
Allows for 77% cache storage savings of approximable data!



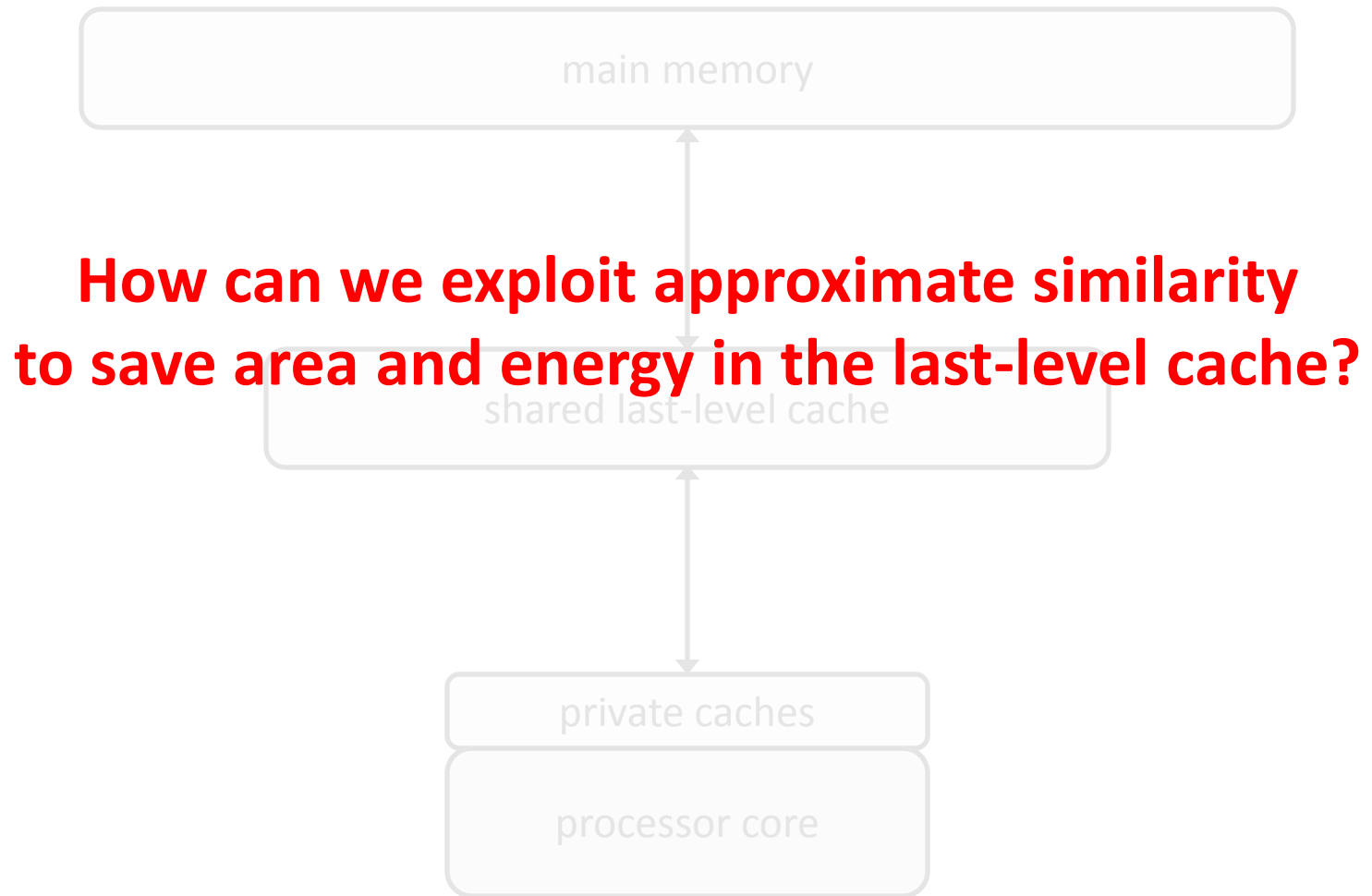
Outline

- Approximate Computing
 - Approximate Similarity
- **Doppelgänger Cache**
 - **Cache Architecture**
 - **Similarity Mapping**
- Evaluation

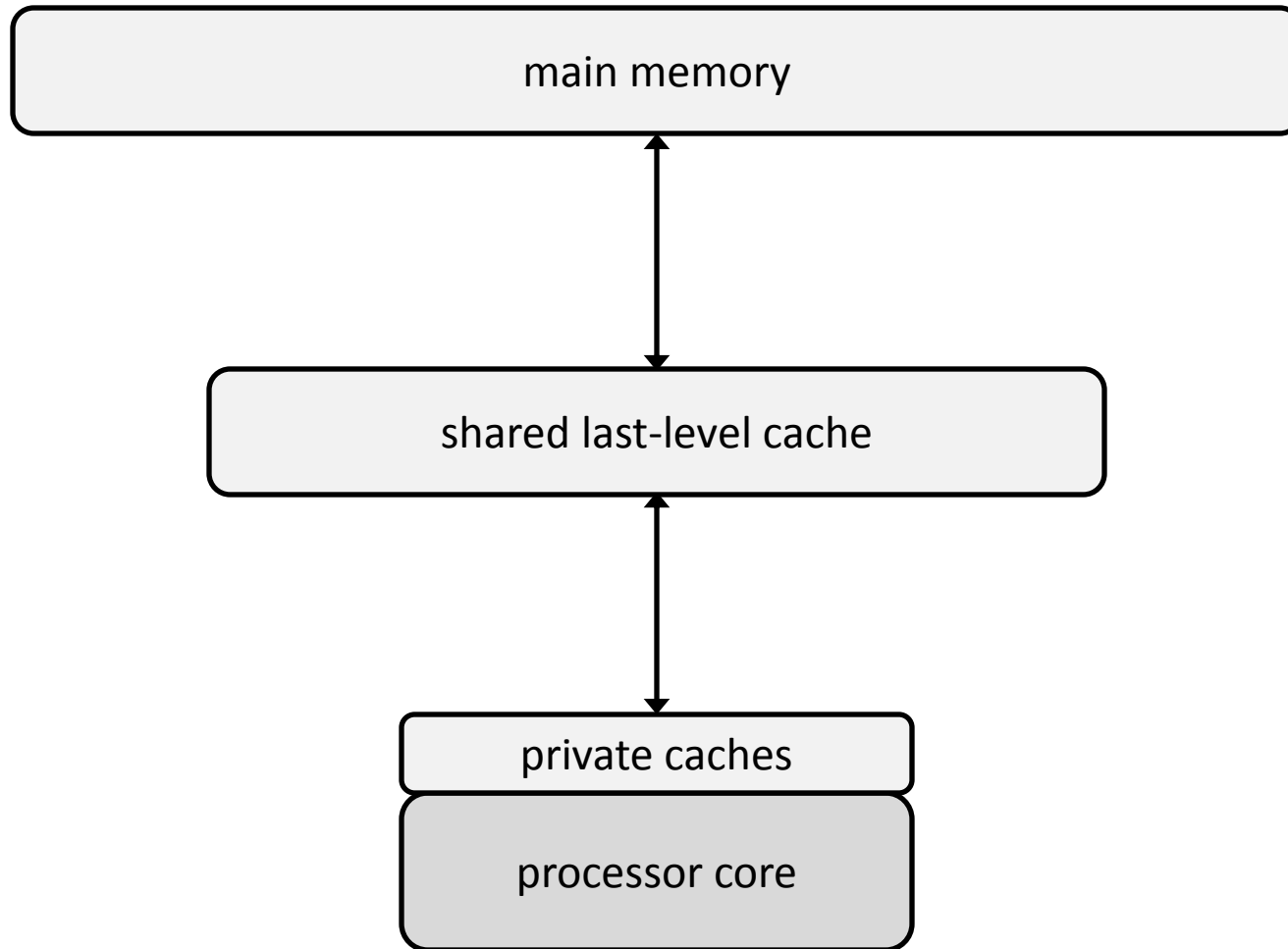
Doppelgänger Cache



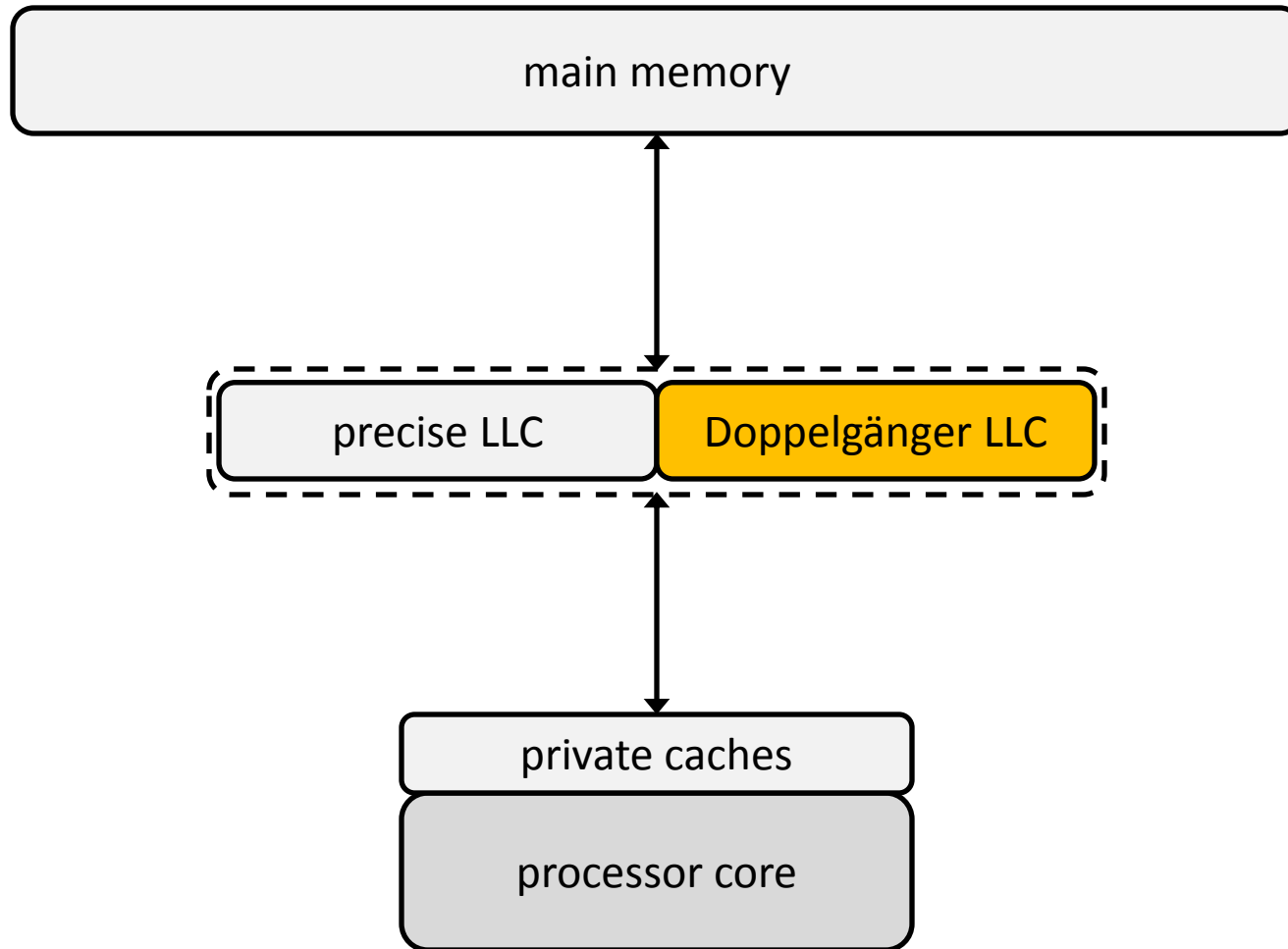
Doppelgänger Cache



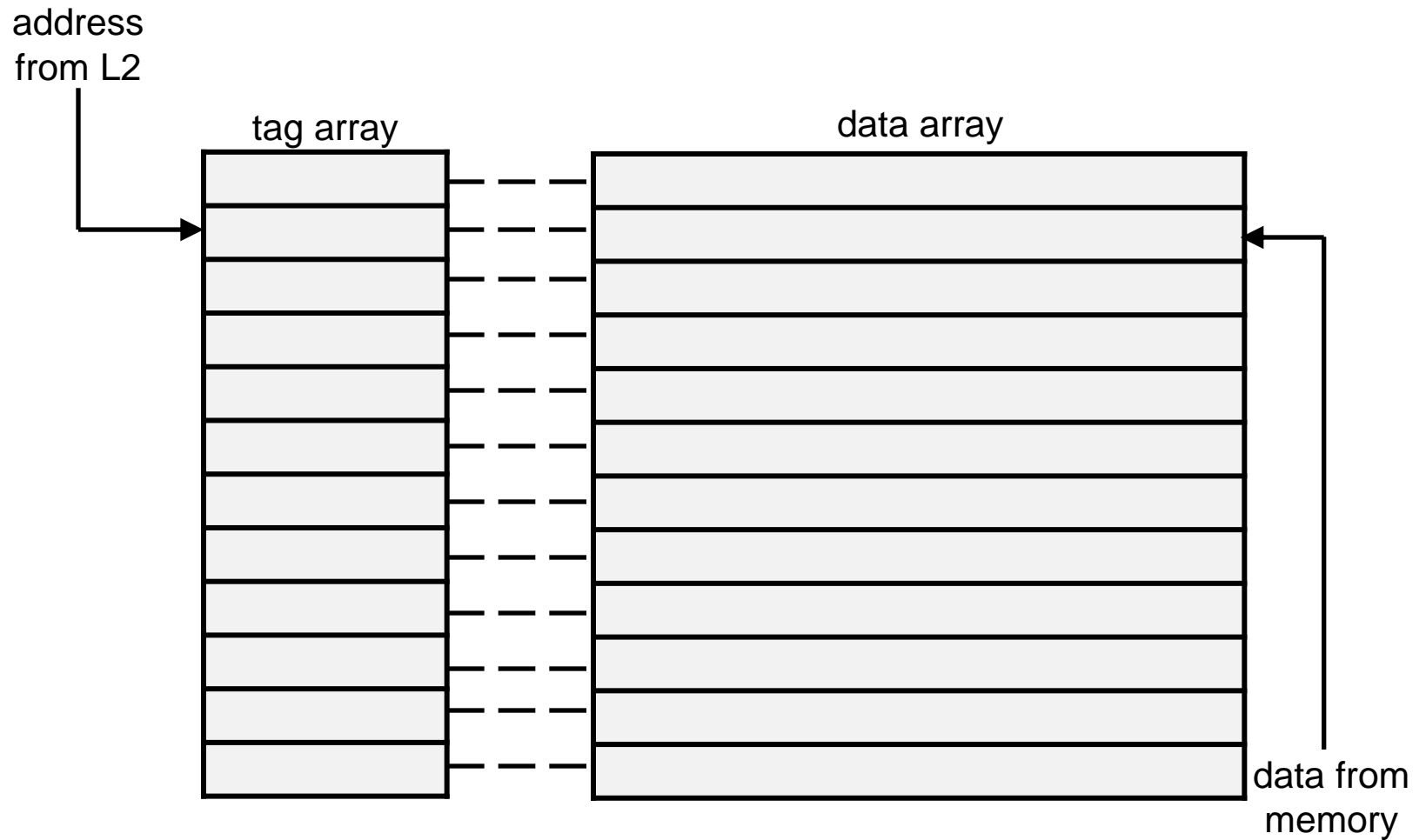
Doppelgänger Cache



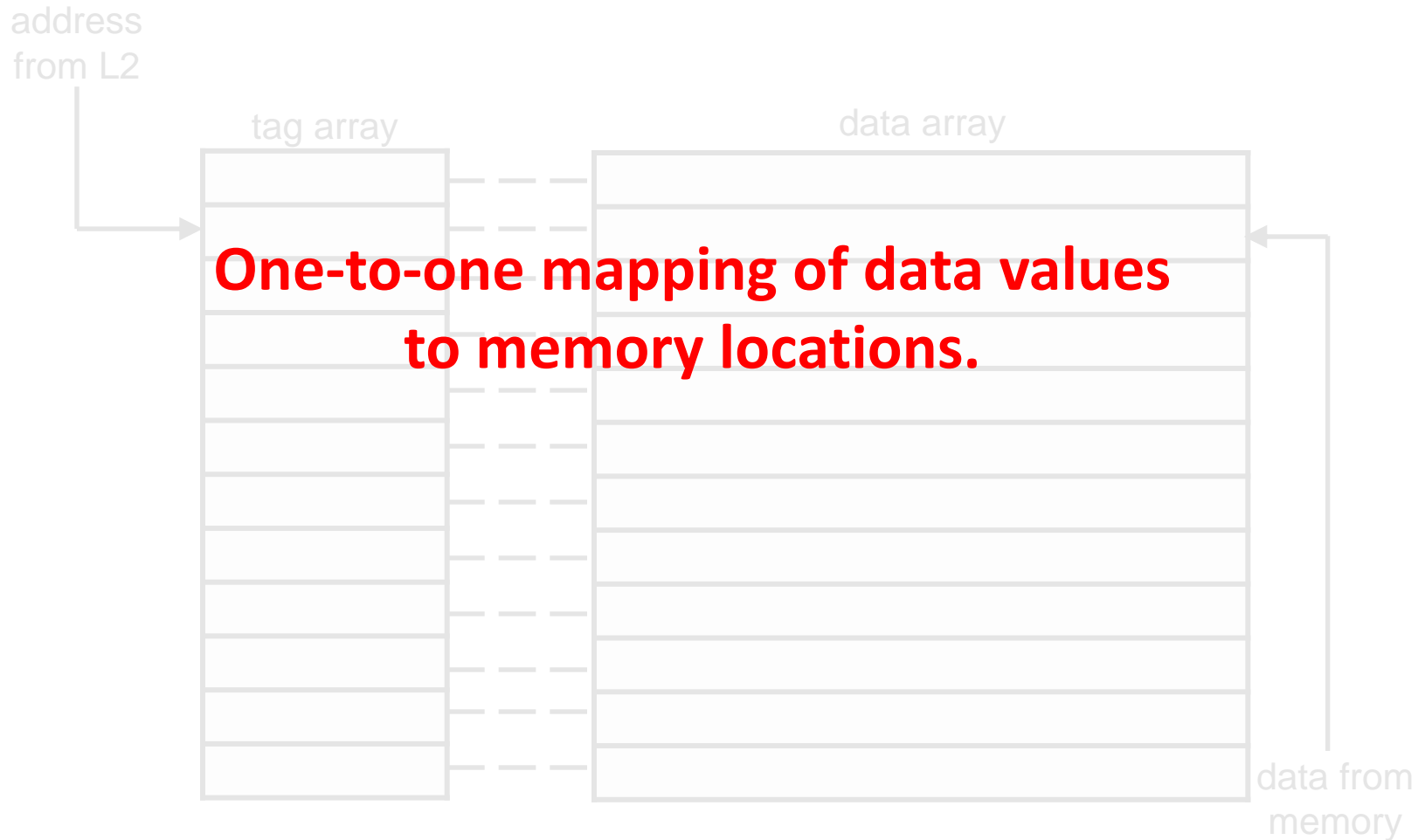
Doppelgänger Cache



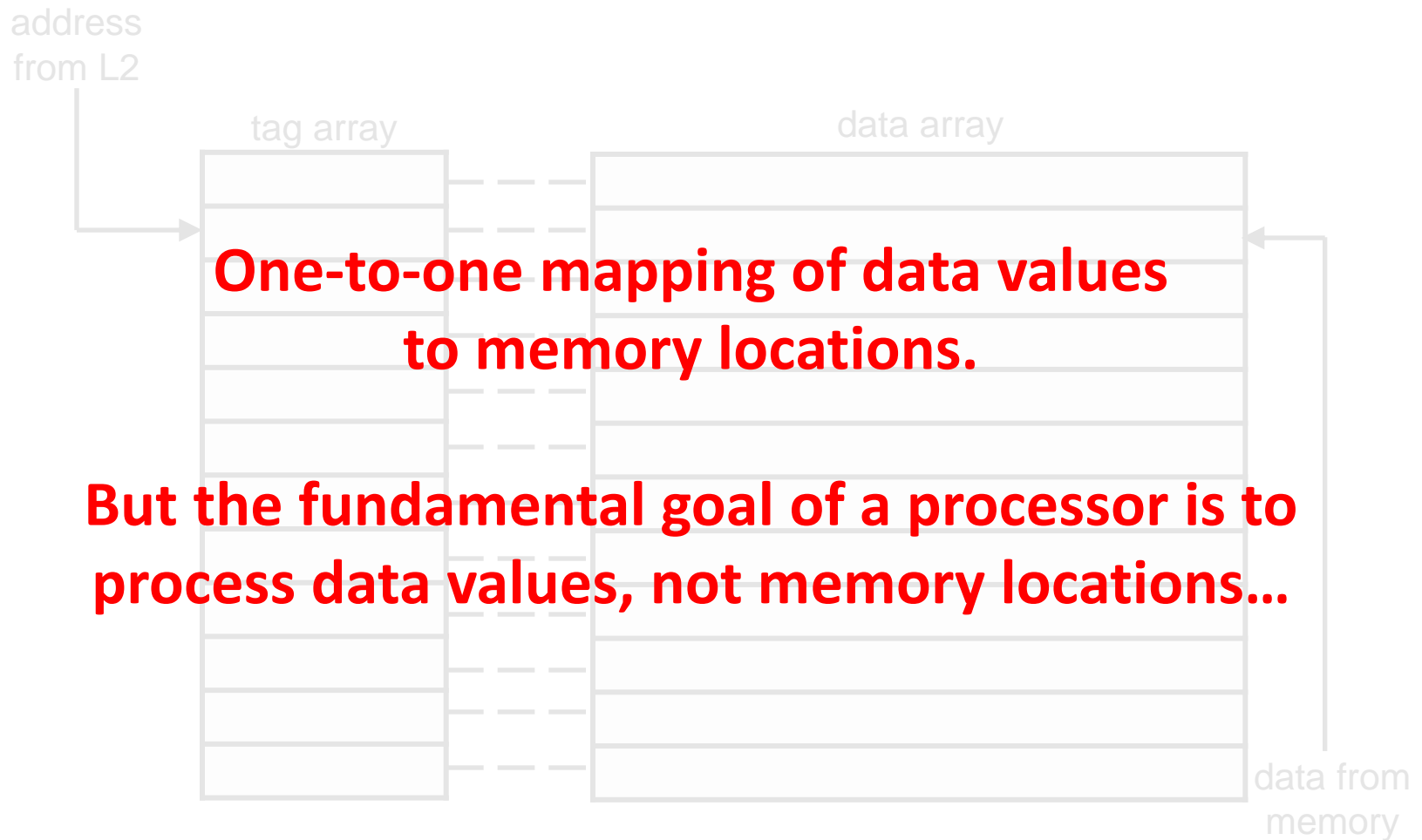
Conventional Cache



Conventional Cache



Conventional Cache



Conventional Cache

address
from L2

tag array



data from
memory

Conventional Cache

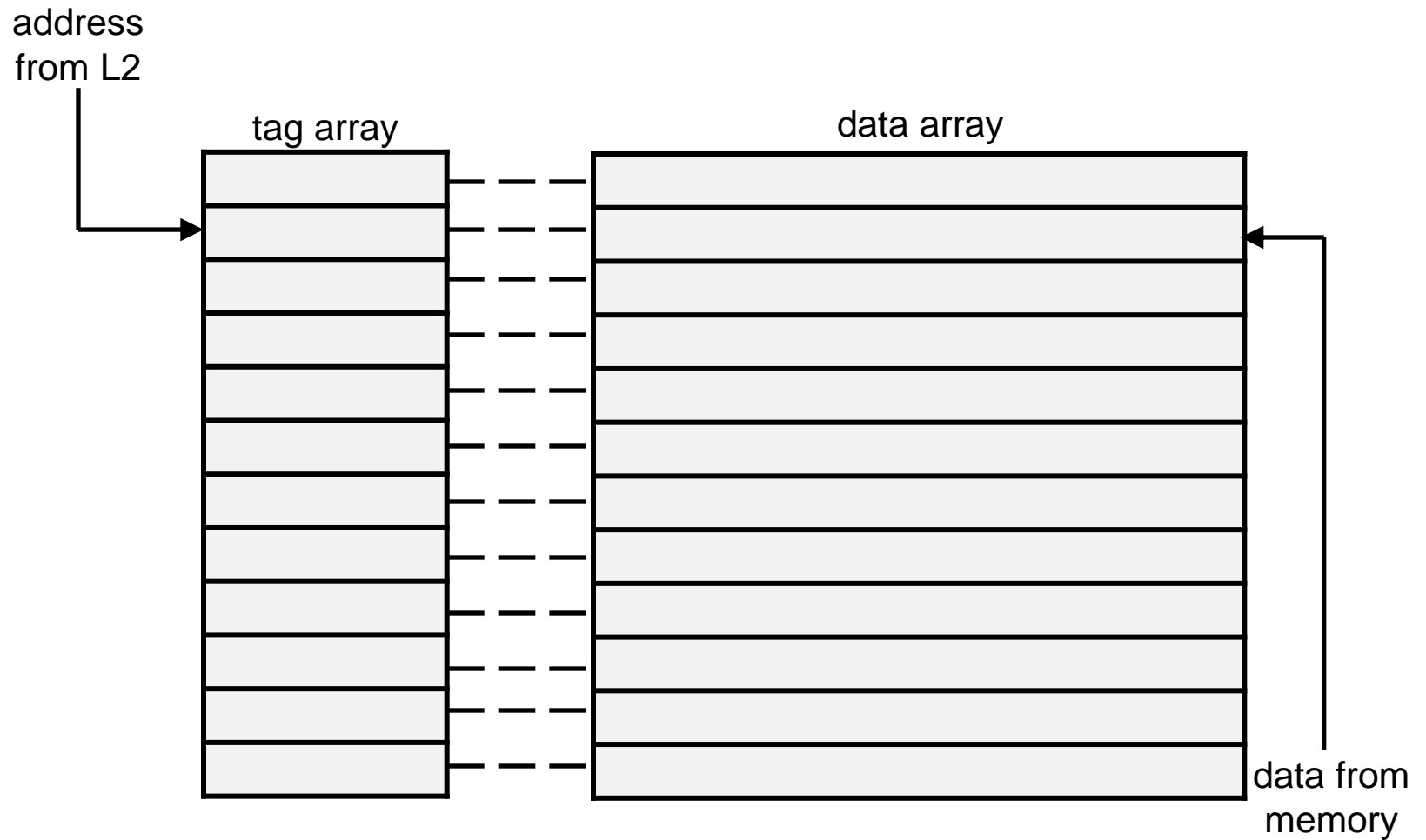
address
from L2

tag array

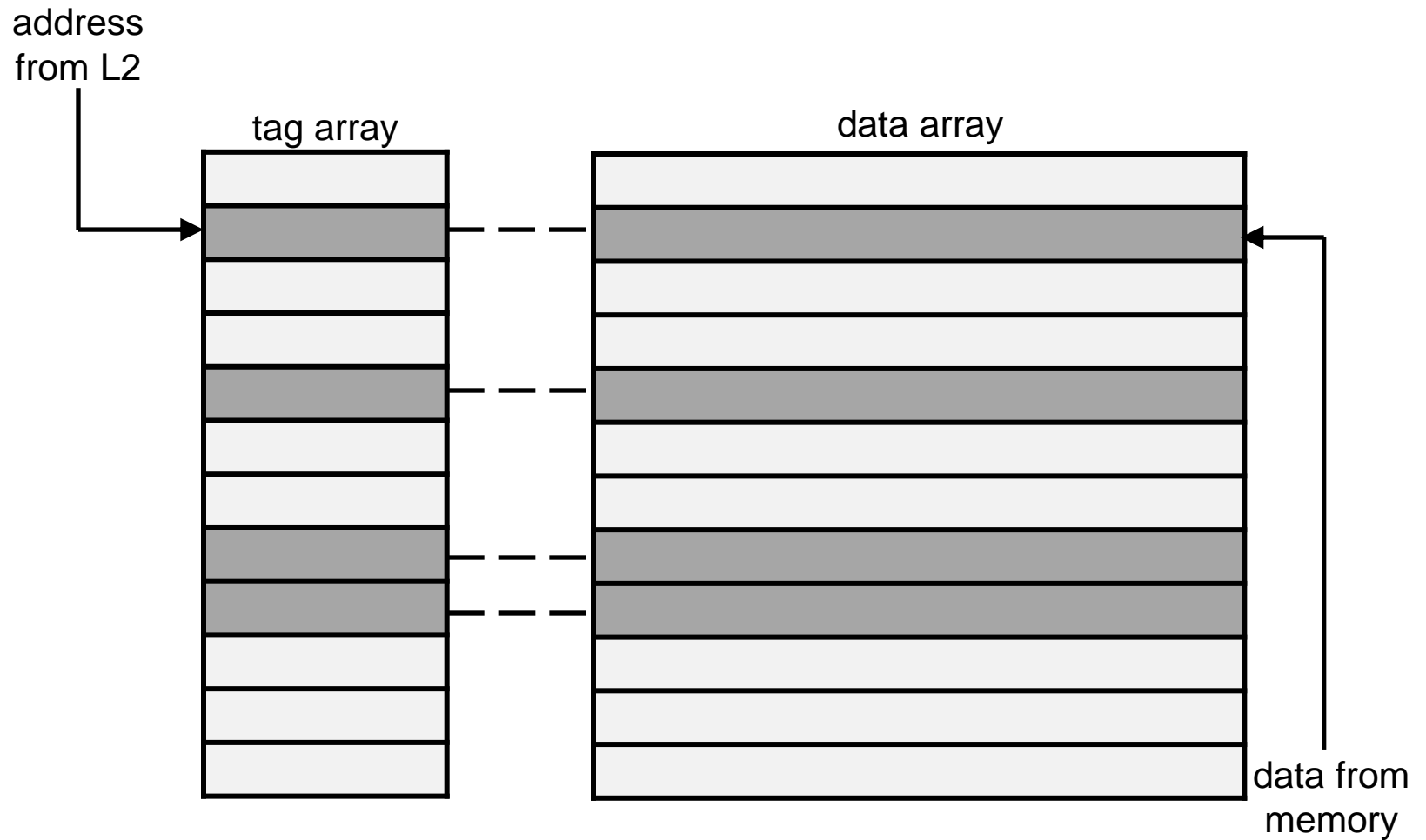


data from
memory

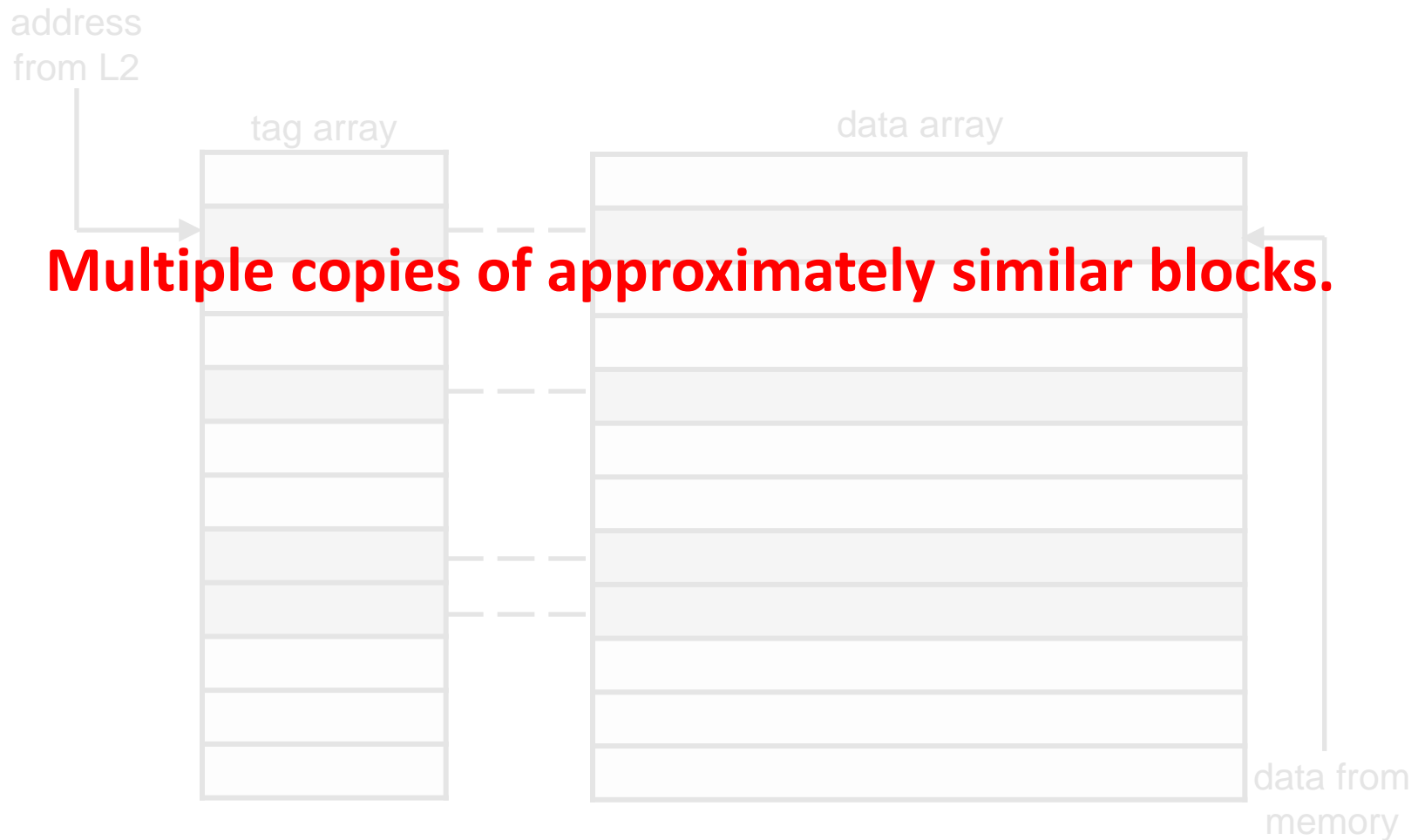
Conventional Cache



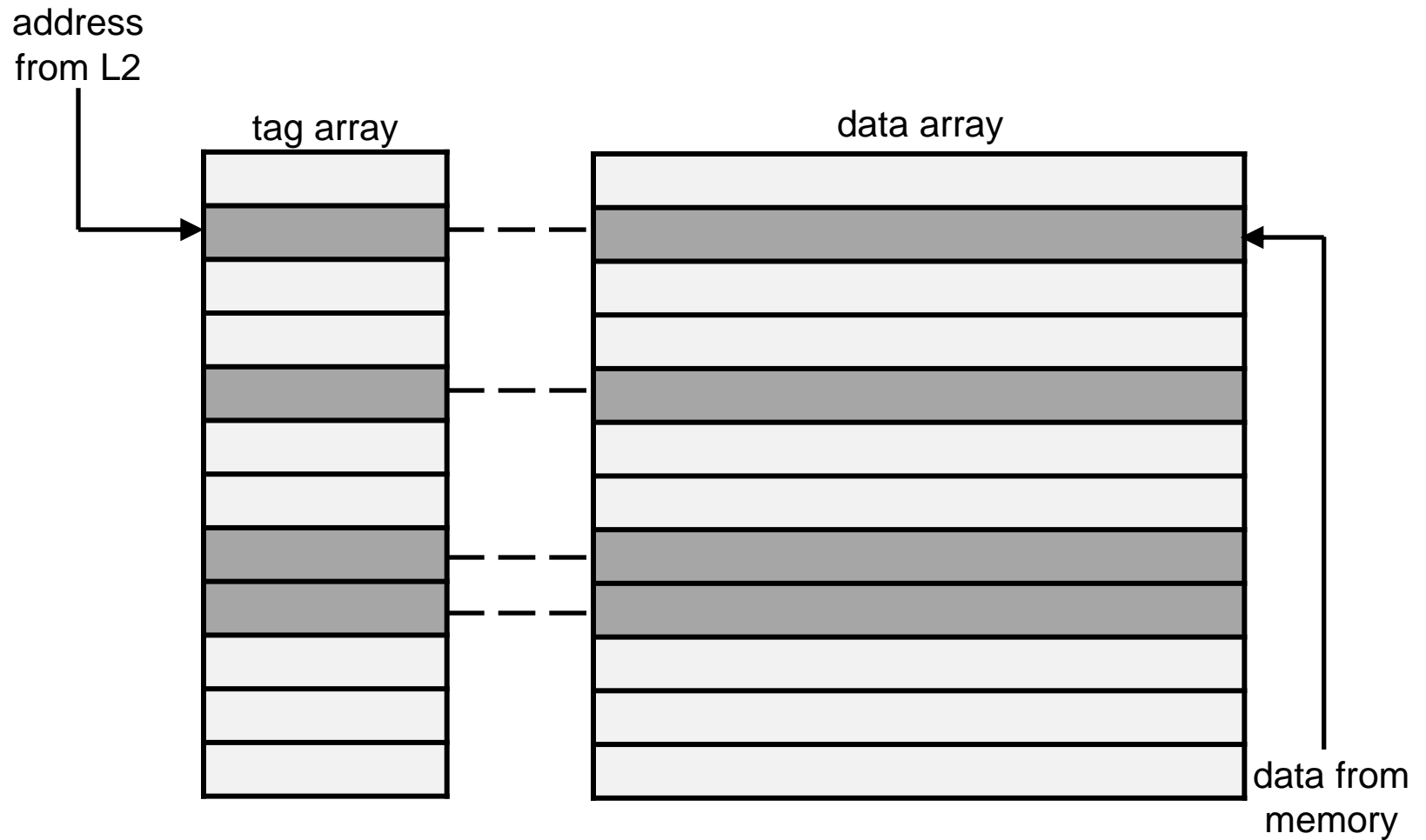
Conventional Cache



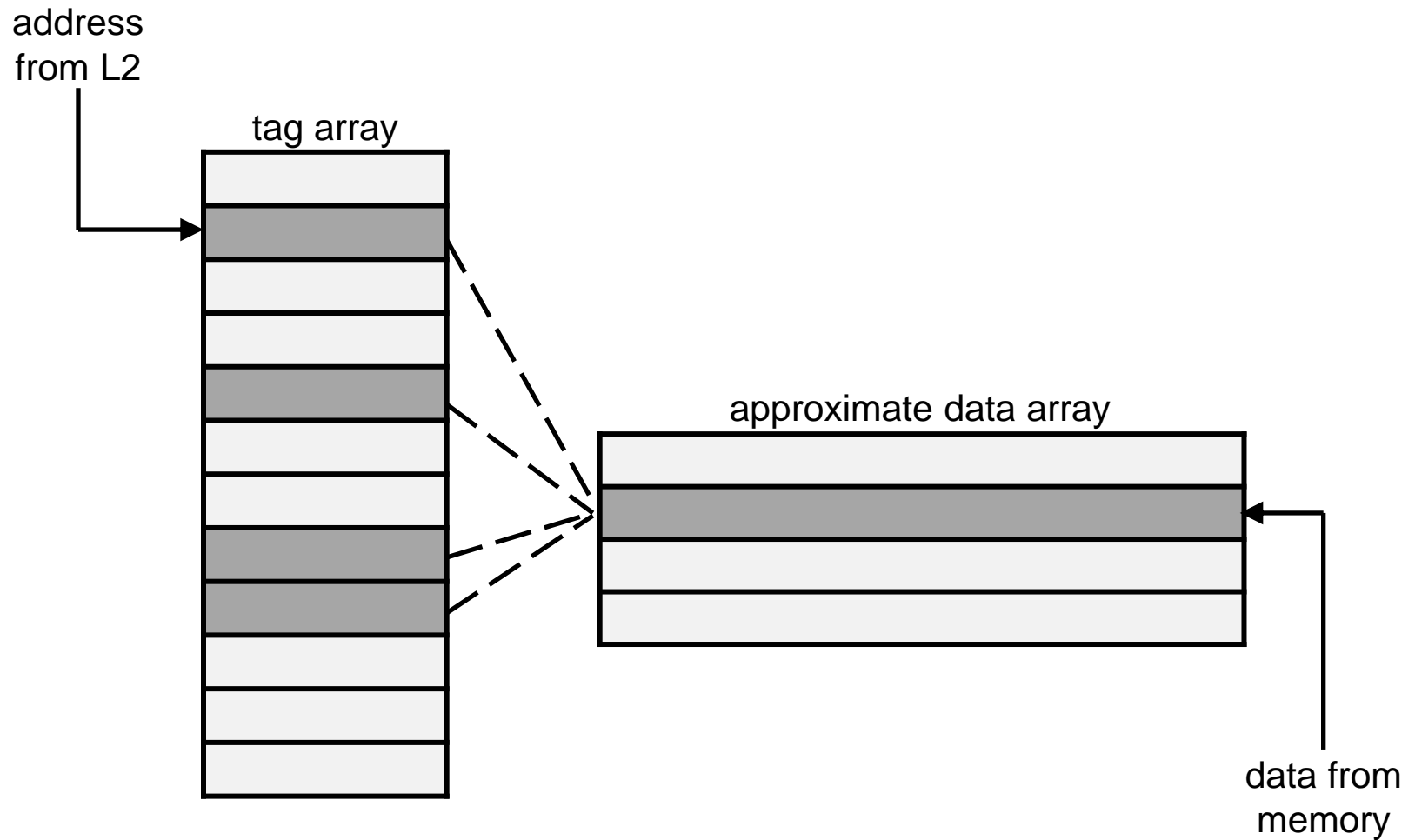
Conventional Cache



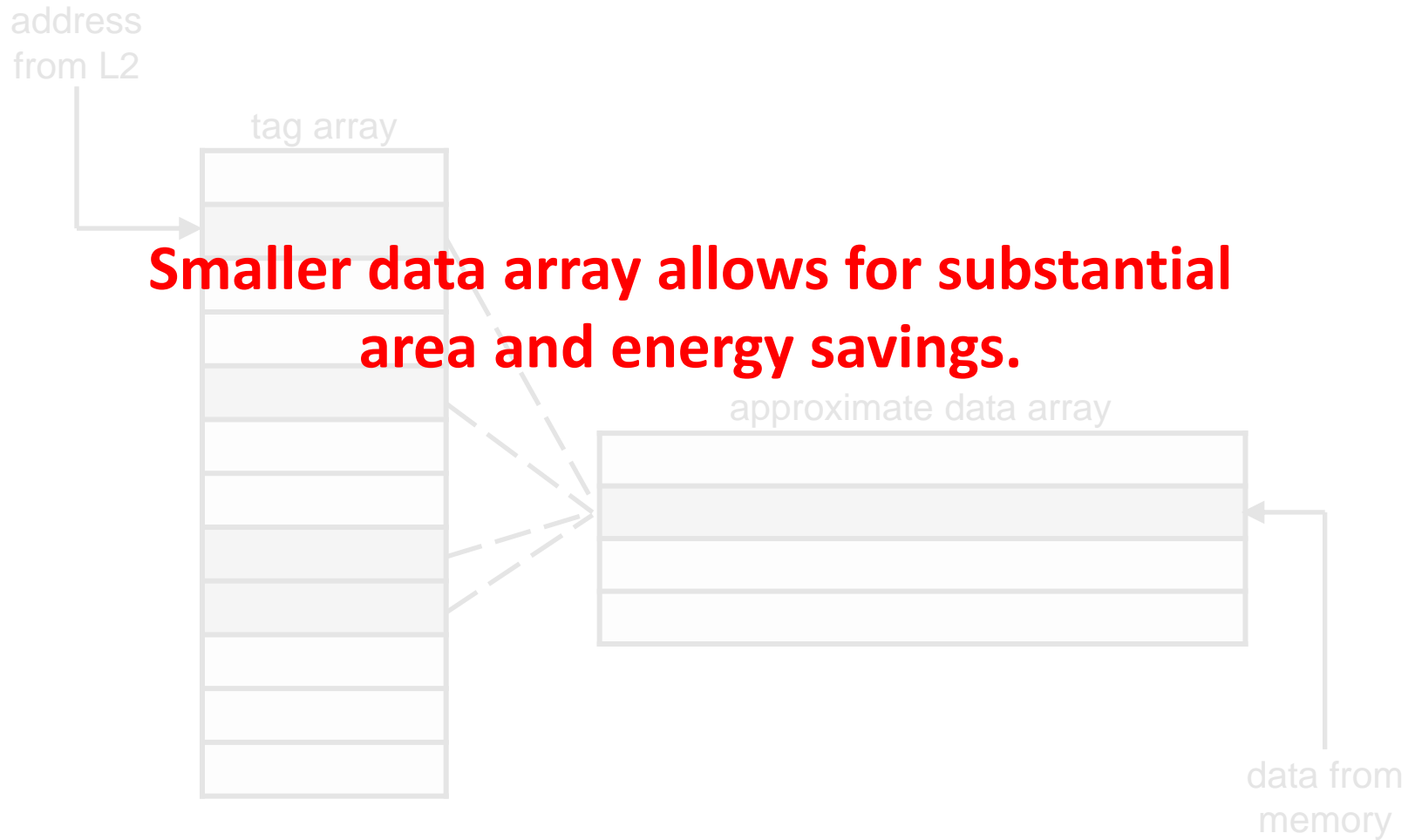
Conventional Cache



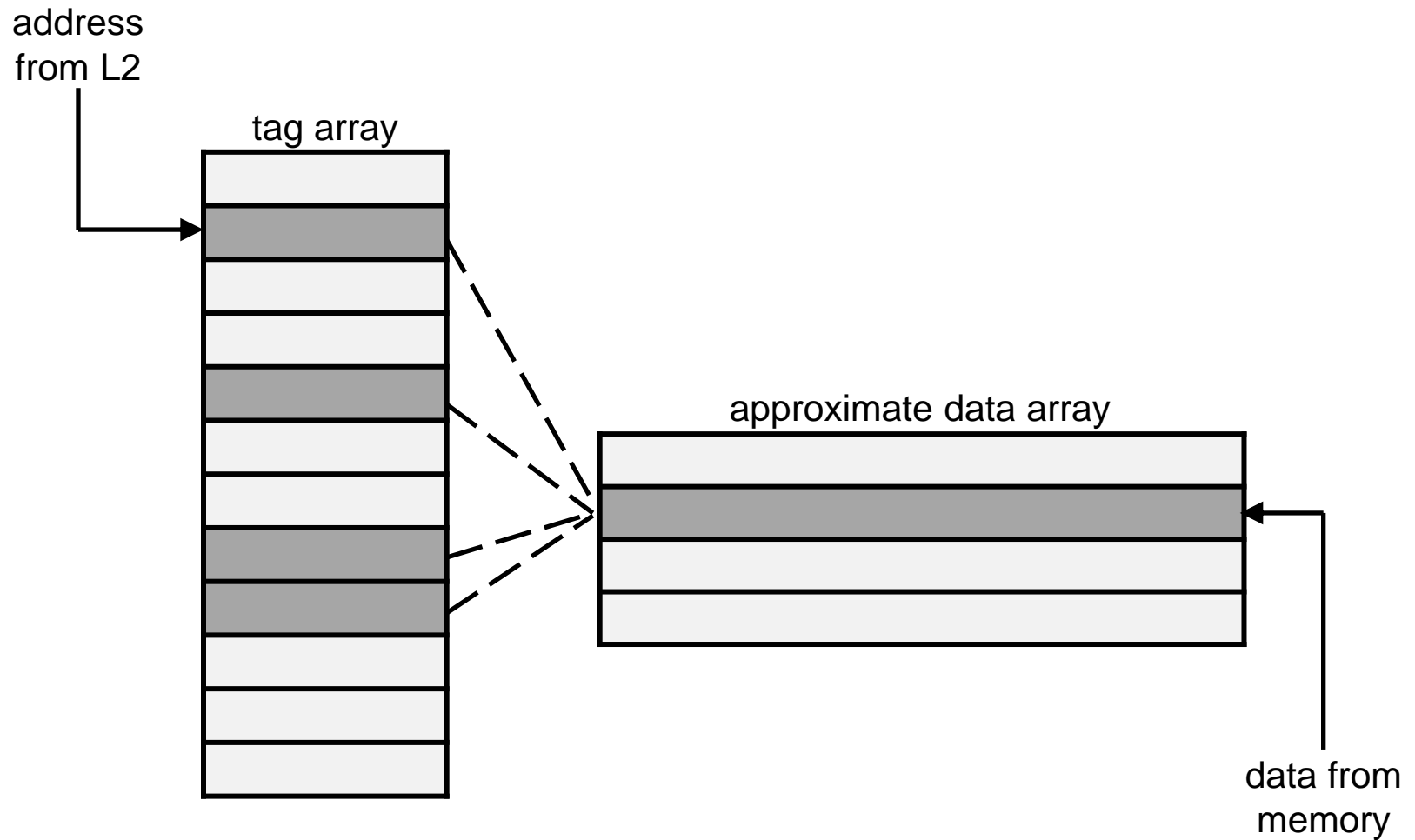
Doppelgänger Cache



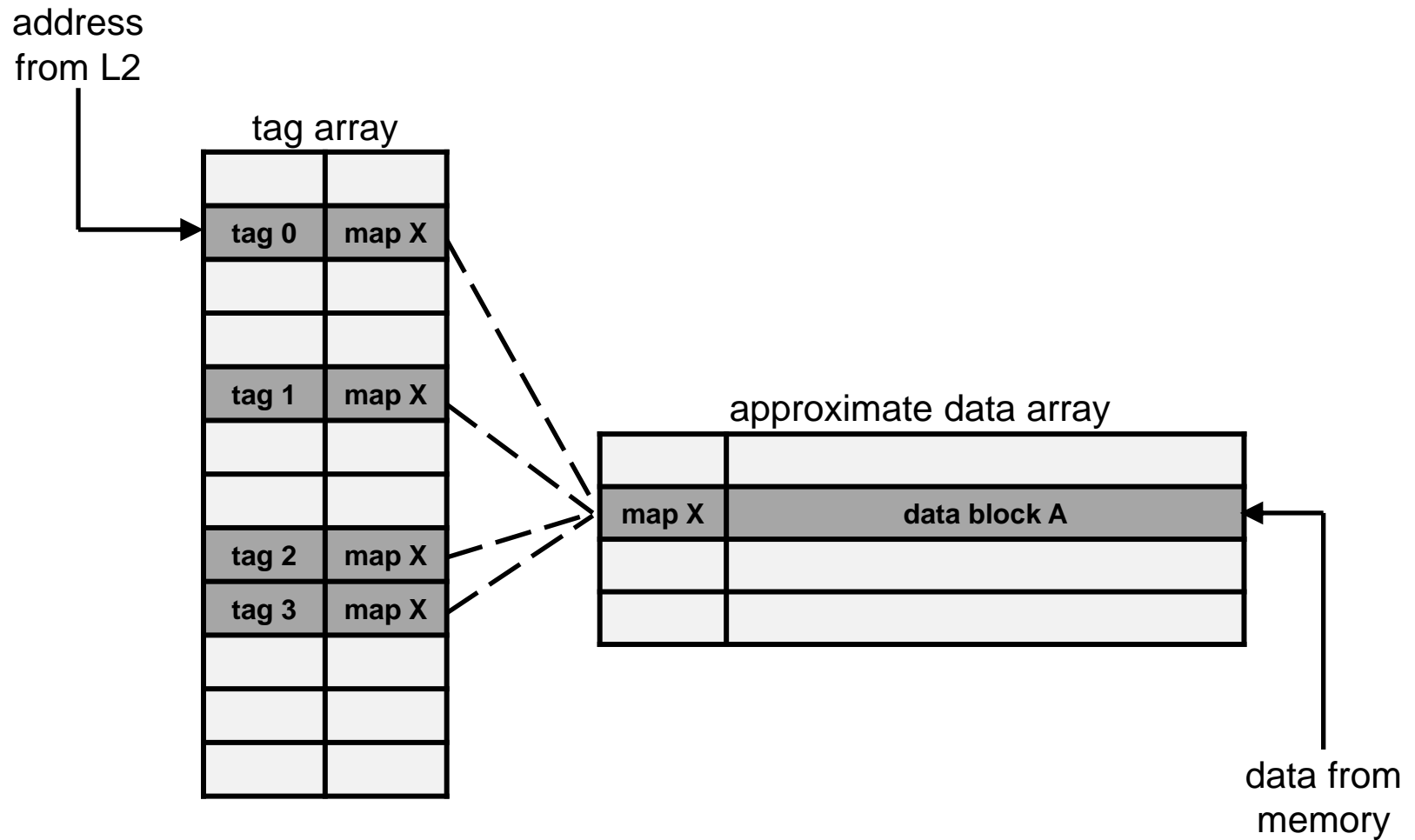
Doppelgänger Cache



Doppelgänger Cache



Doppelgänger Cache



Doppelgänger Cache - Lookups

tag array

tag 0	map X
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

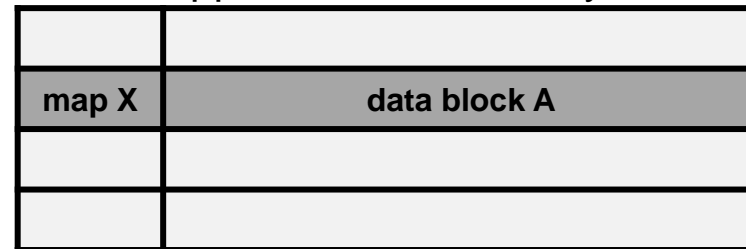
map X	data block A

Doppelgänger Cache - Lookups

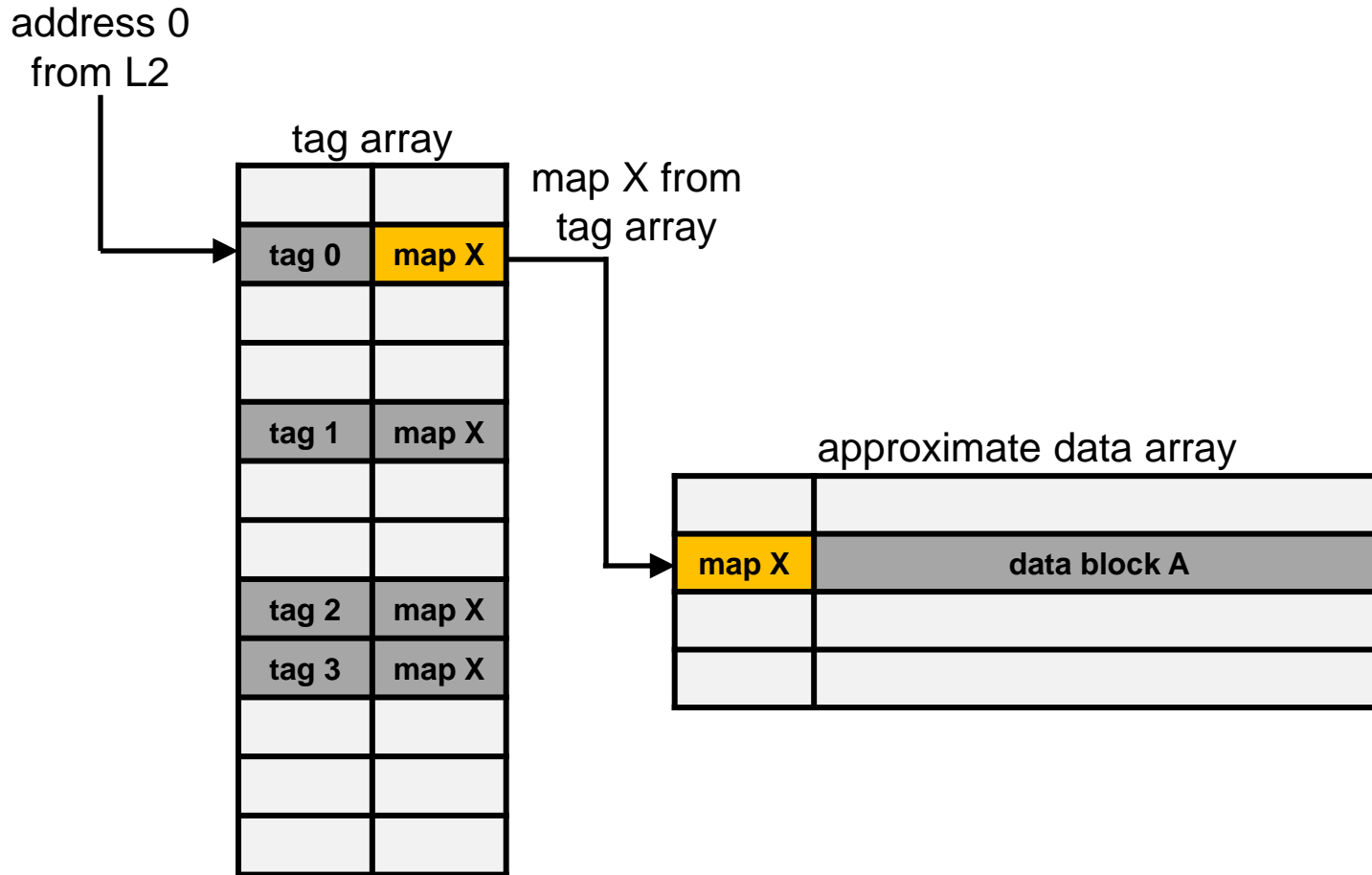
address 0
from L2



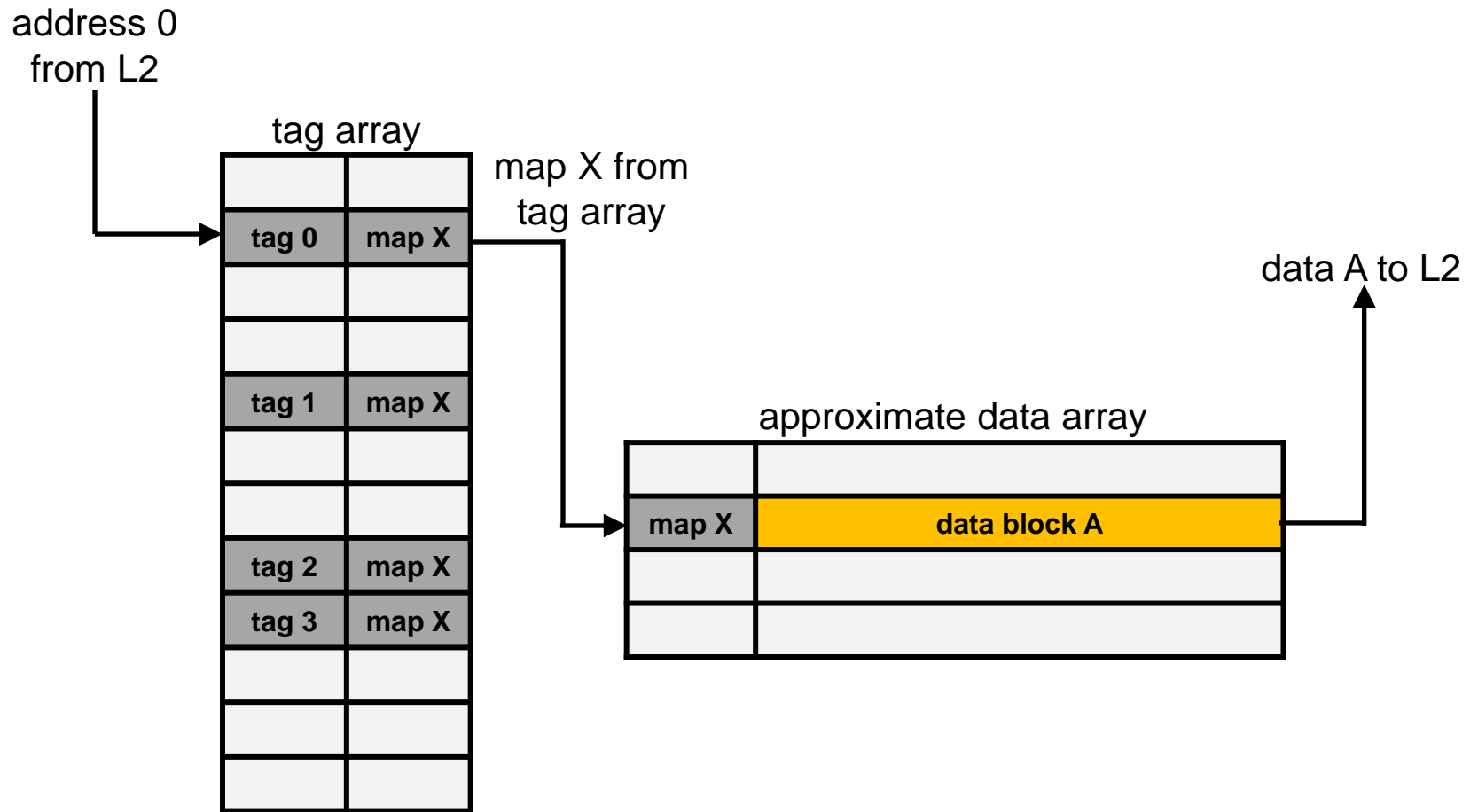
approximate data array



Doppelgänger Cache - Lookups



Doppelgänger Cache - Lookups



Doppelgänger Cache - Insertions

tag array

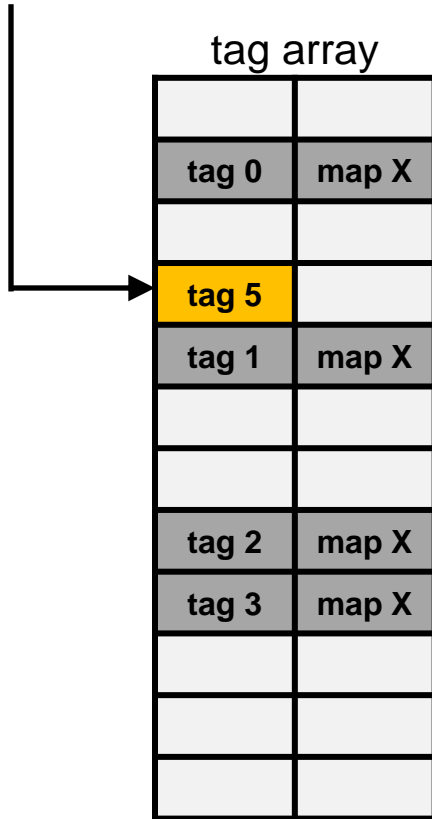
tag 0	map X
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

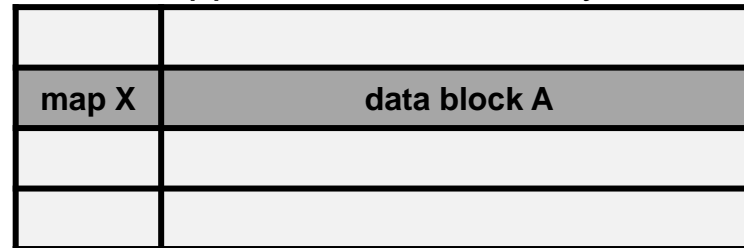
map X	data block A

Doppelgänger Cache - Insertions

address 5
from L2



approximate data array



Doppelgänger Cache - Insertions

address 5
from L2

tag array

tag 0	map X
tag 5	
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

map X	data block A

data B to L2

data B from
memory

Doppelgänger Cache - Insertions

address 5
from L2

tag array

tag 0	map X
tag 5	map Y
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

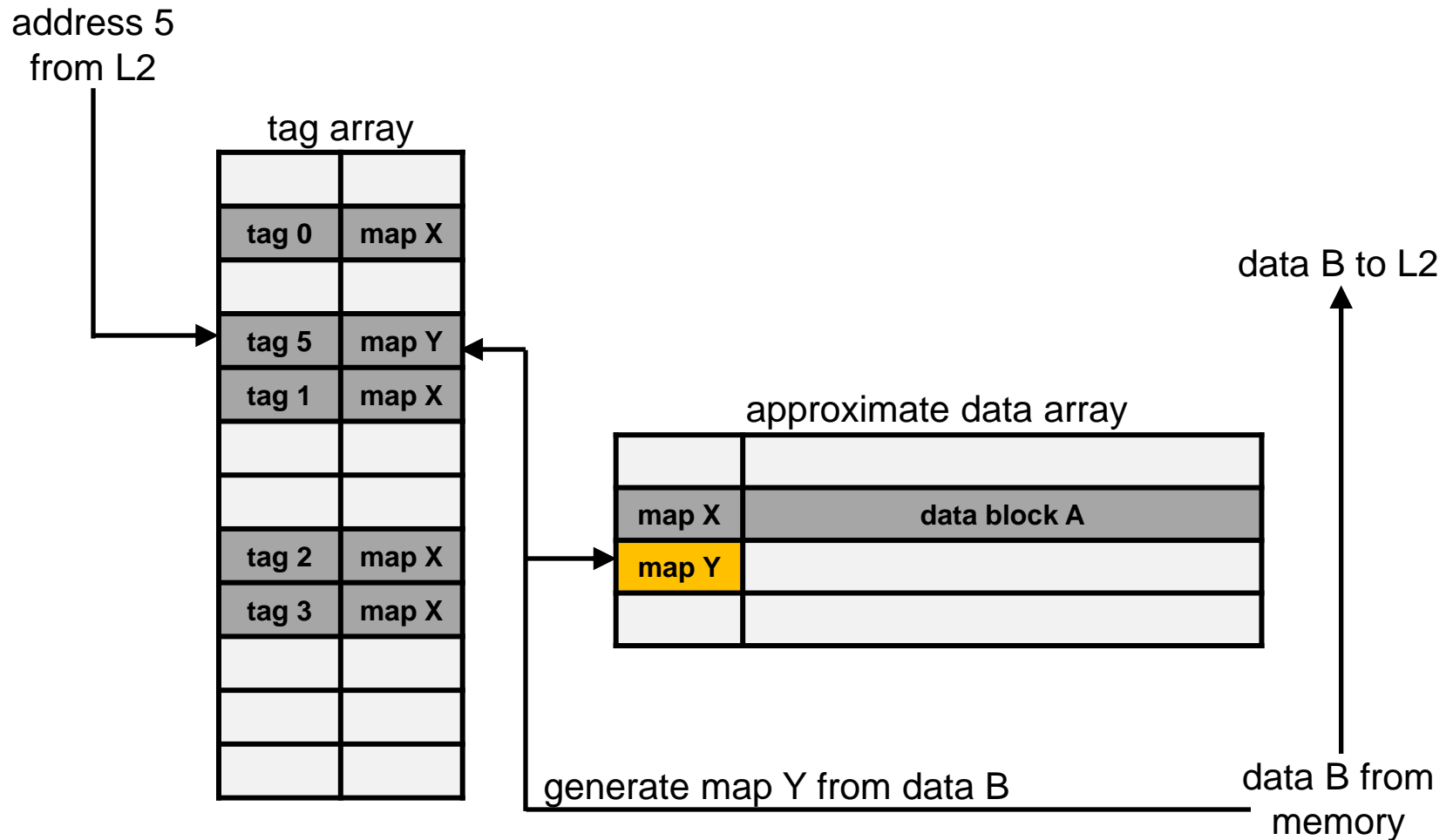
map X	data block A

data B to L2

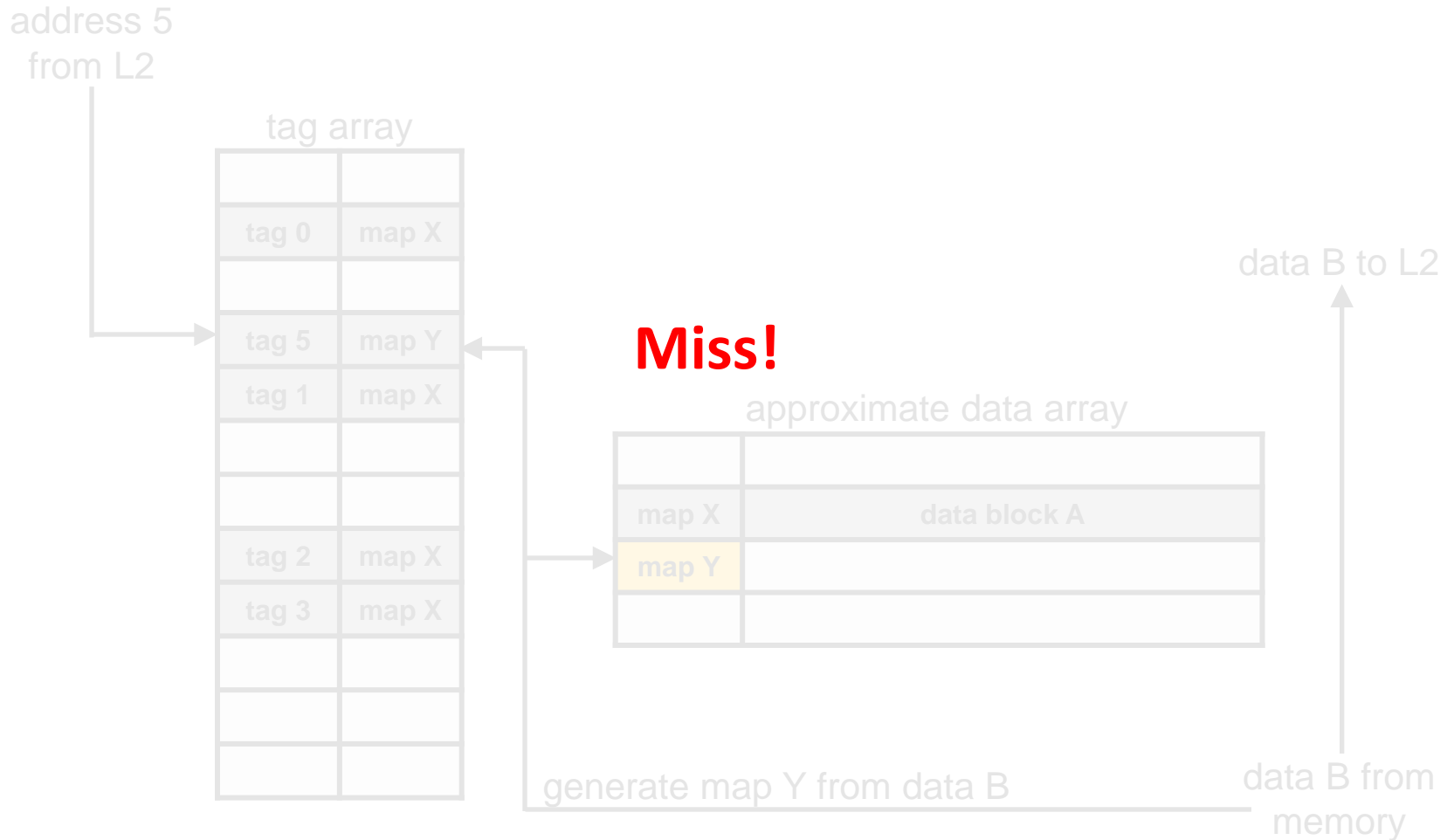
generate map Y from data B

data B from
memory

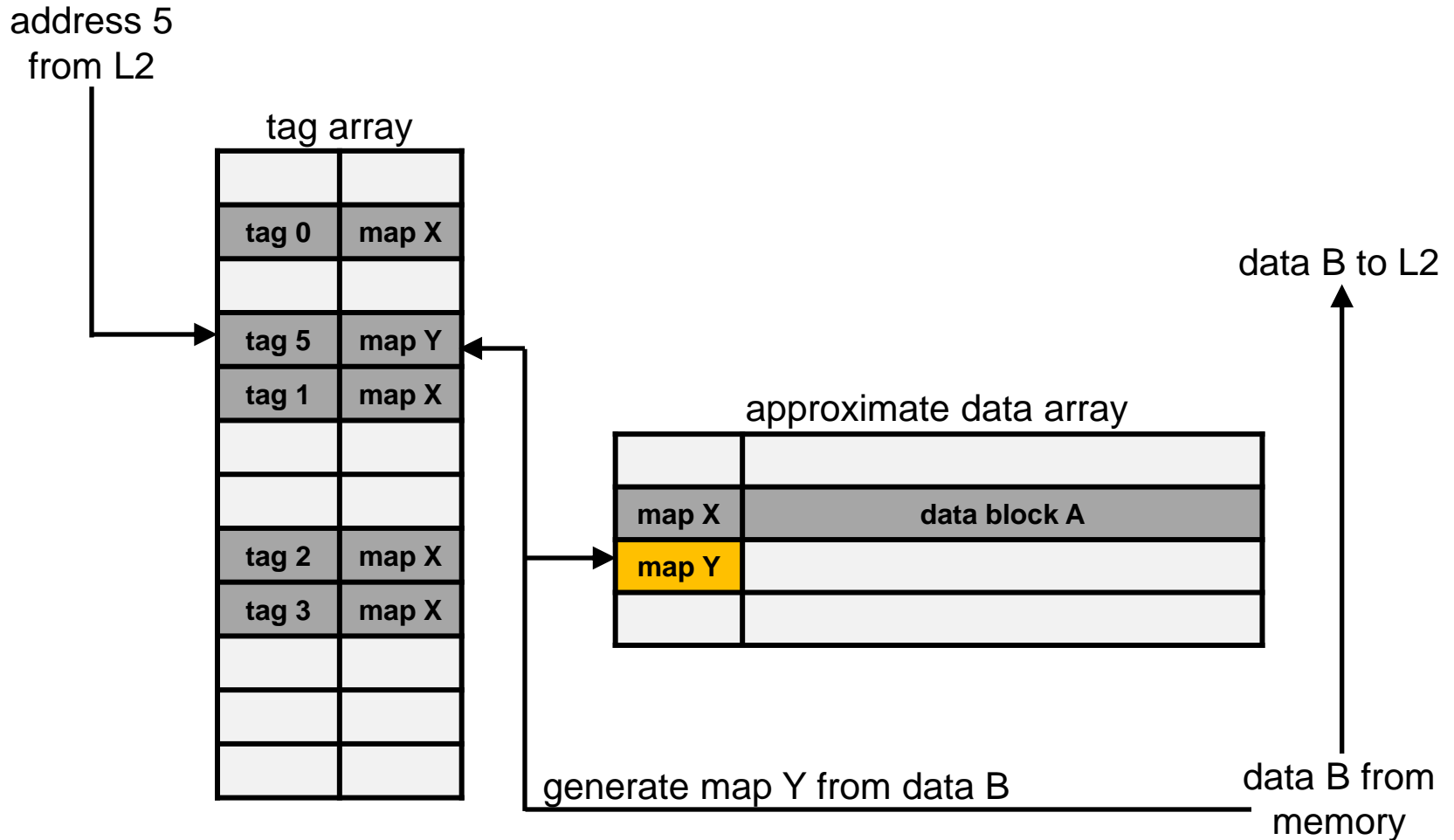
Doppelgänger Cache - Insertions



Doppelgänger Cache - Insertions



Doppelgänger Cache - Insertions (Miss)



Doppelgänger Cache - Insertions (Miss)

address 5
from L2

tag array

tag 0	map X
tag 5	map Y
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

map X	data block A
map Y	data block B

generate map Y from data B

data B to L2

data B from
memory

Doppelgänger Cache - Insertions

tag array

tag 0	map X
tag 5	map Y
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

map X	data block A
map Y	data block B

Doppelgänger Cache - Insertions

address 6
from L2

tag array

tag 0	map X
tag 6	
tag 5	map Y
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

map X	data block A
map Y	data block B

Doppelgänger Cache - Insertions

address 6
from L2

tag array

tag 0	map X
tag 6	
tag 5	map Y
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

map X	data block A
map Y	data block B

data C to L2

data C from
memory

Doppelgänger Cache - Insertions

address 6
from L2

tag array

tag 0	map X
tag 6	map X
tag 5	map Y
tag 1	map X
tag 2	map X
tag 3	map X

approximate data array

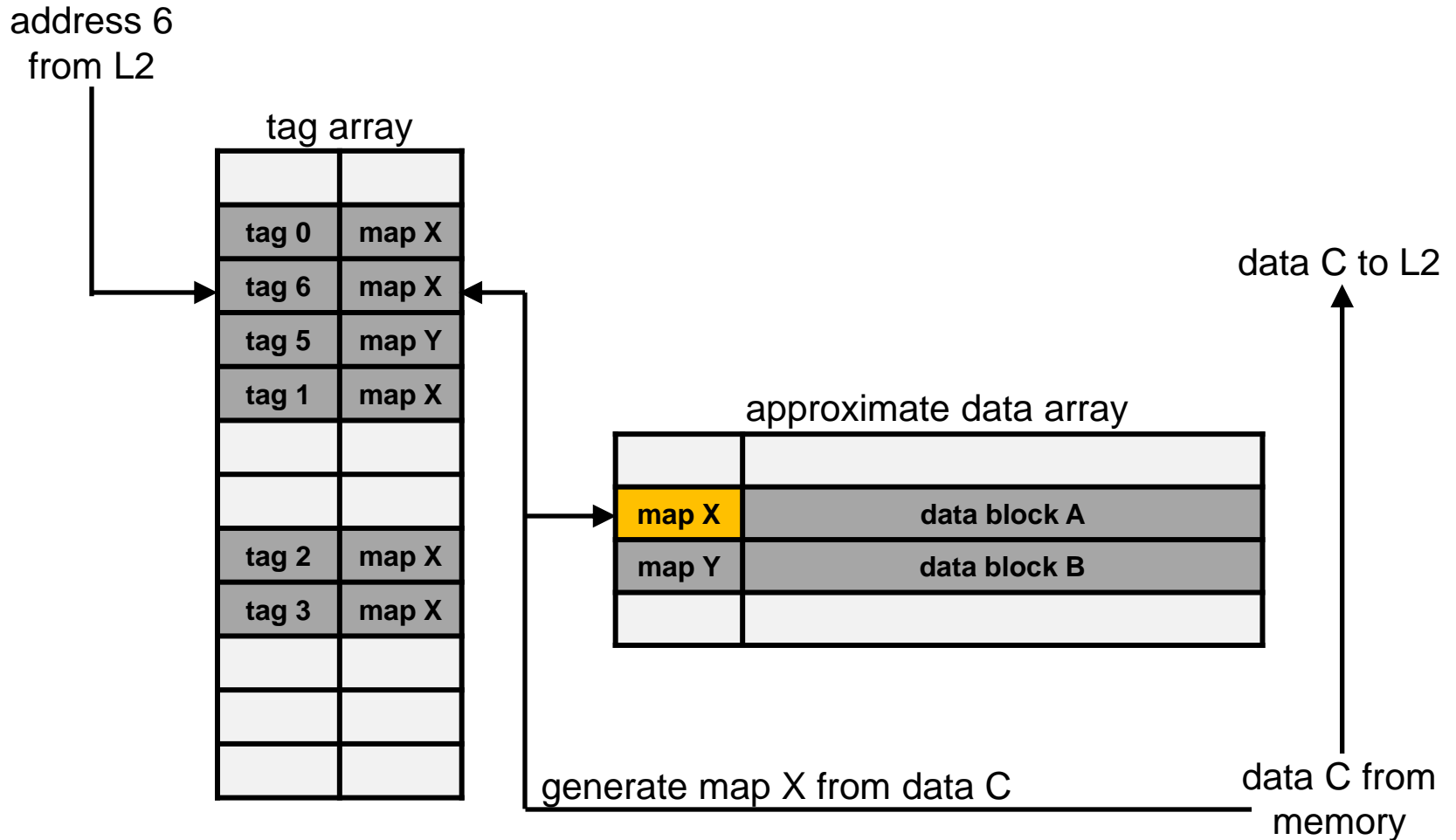
map X	data block A
map Y	data block B

generate map X from data C

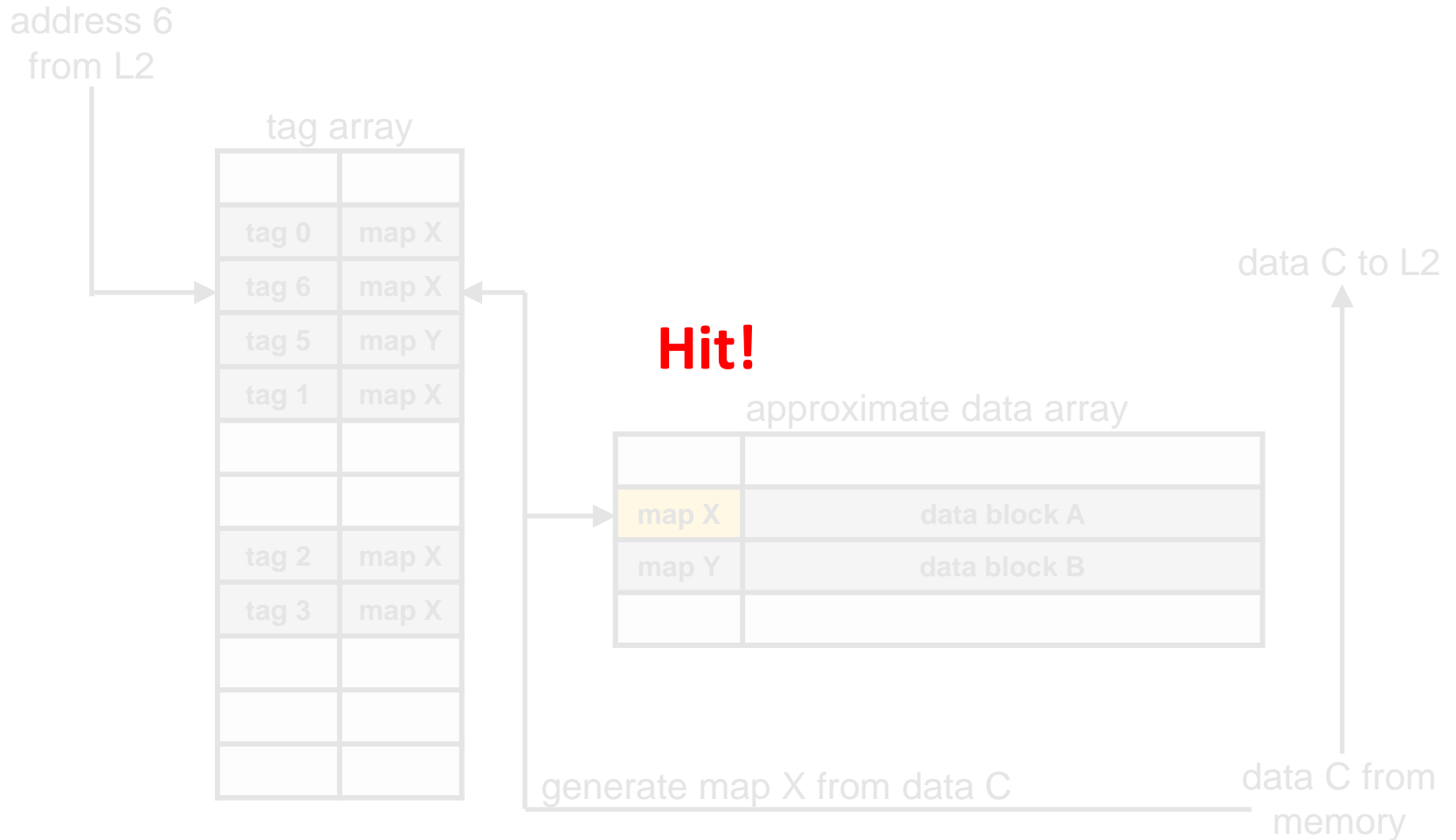
data C to L2

data C from
memory

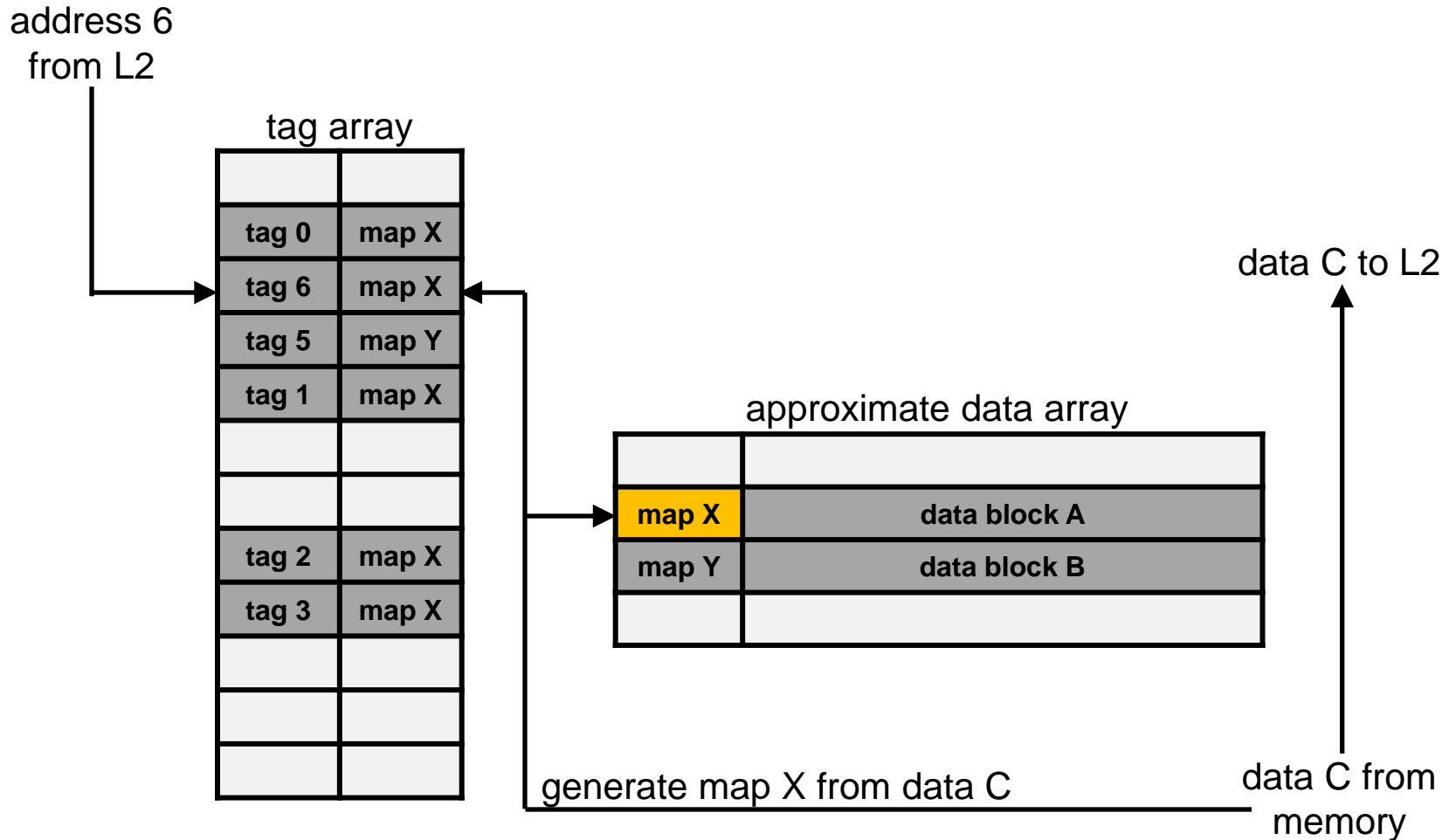
Doppelgänger Cache - Insertions



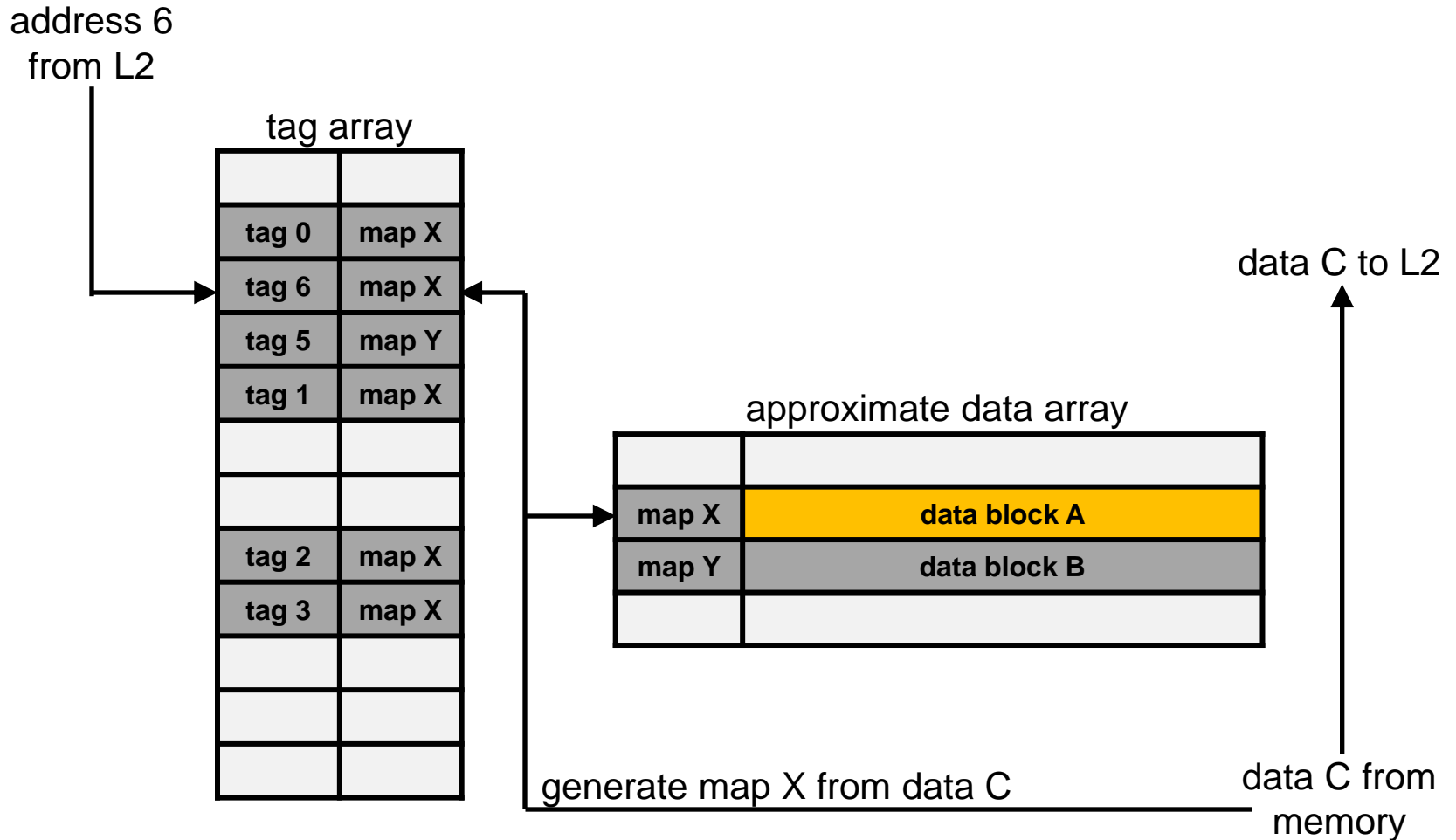
Doppelgänger Cache - Insertions



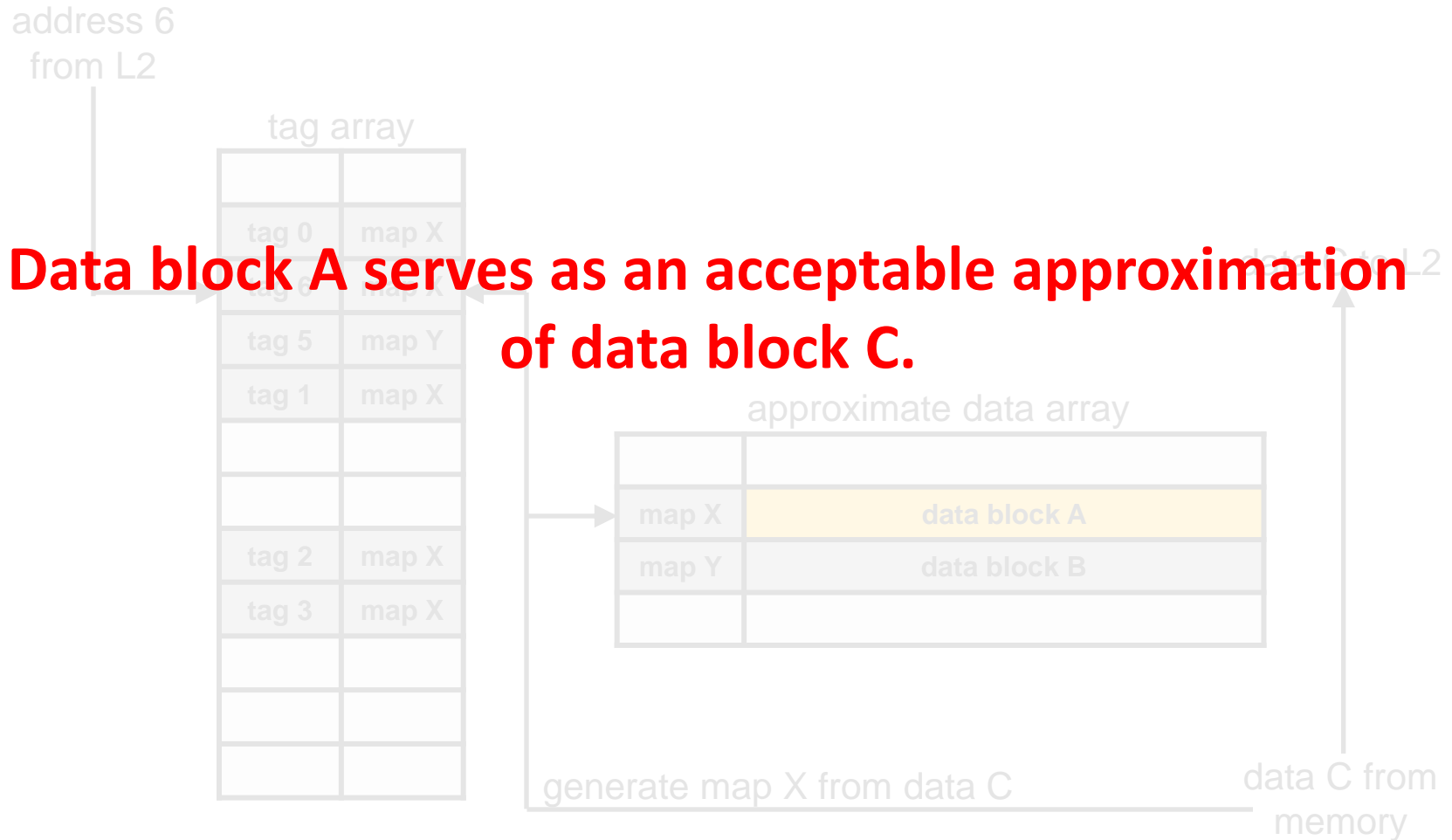
Doppelgänger Cache - Insertions (Hit)



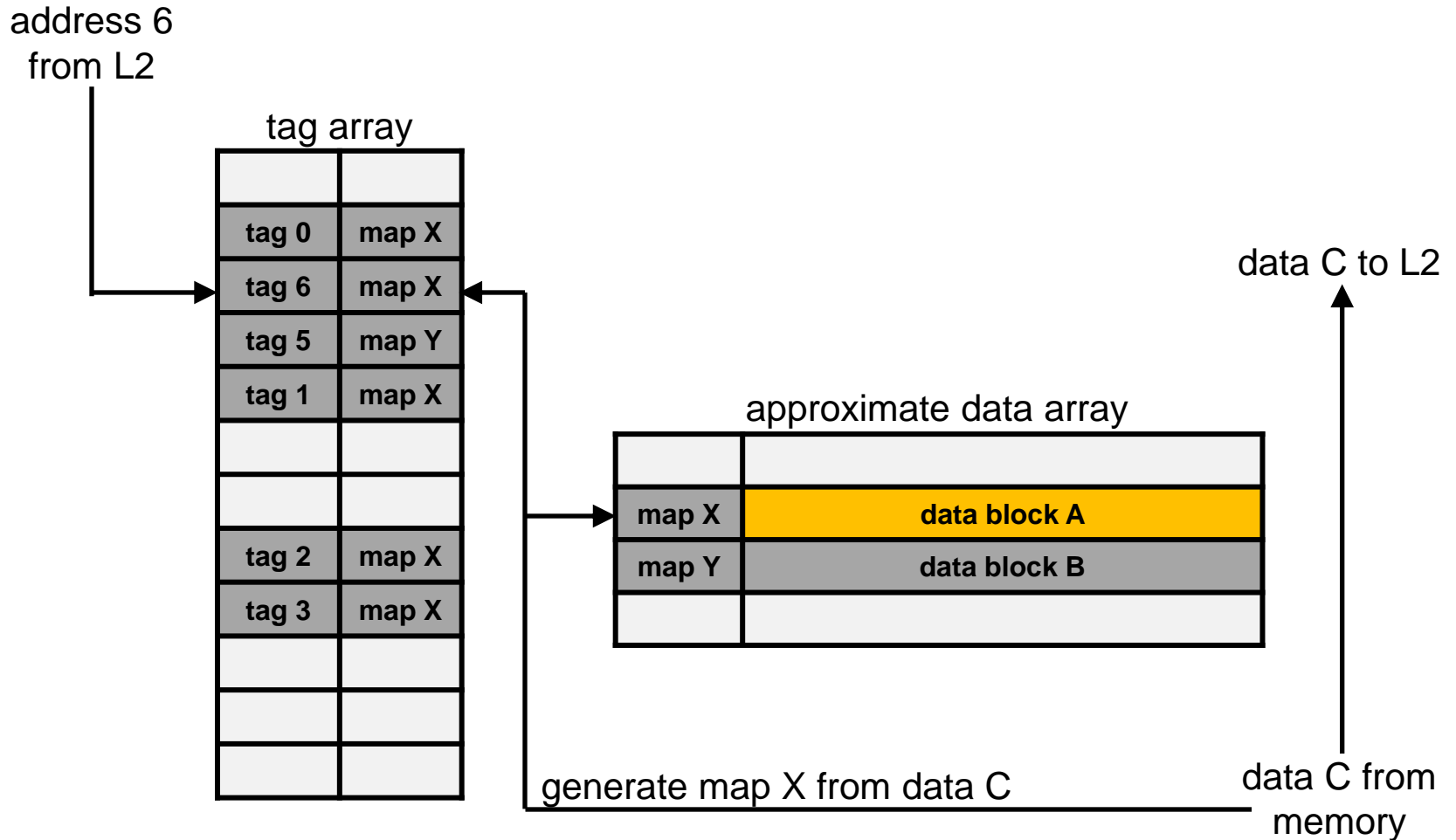
Doppelgänger Cache - Insertions (Hit)



Doppelgänger Cache - Insertions (Hit)



Doppelgänger Cache - Insertions (Hit)

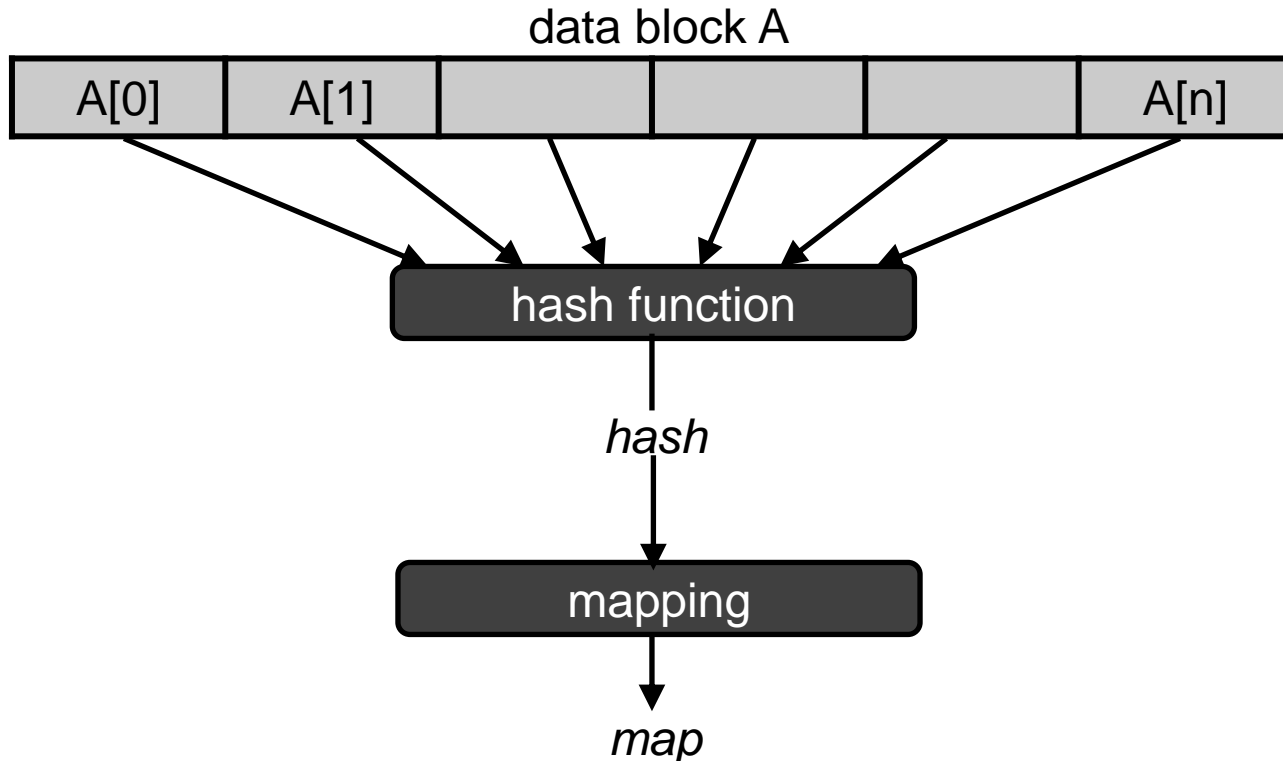


Doppelgänger Cache - Similarity Mapping

The **map** value represents the signature (or **likeness**) of a block. Blocks that generate the same map value are approximately similar.

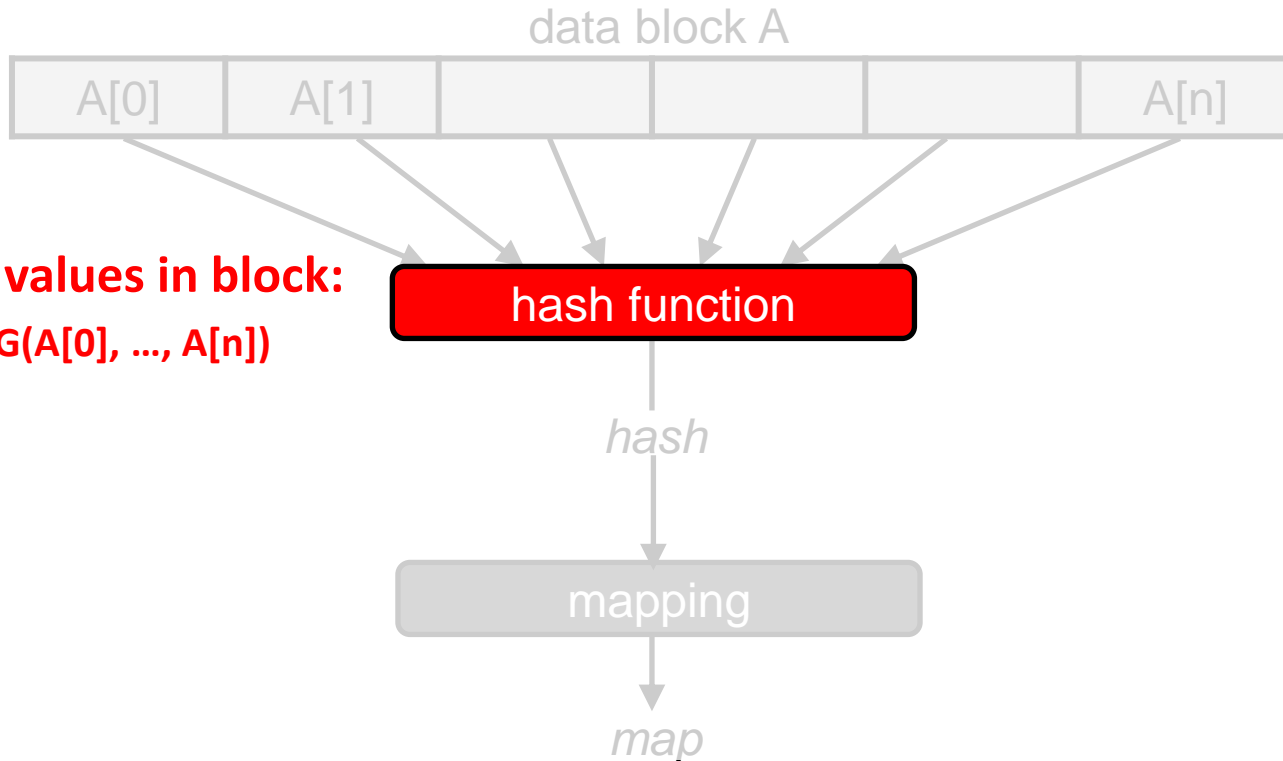
Doppelgänger Cache - Similarity Mapping

The **map** value represents the signature (or **likeness**) of a block. Blocks that generate the same map value are approximately similar.



Doppelgänger Cache - Similarity Mapping

The **map** value represents the signature (or **likeness**) of a block. Blocks that generate the same map value are approximately similar.

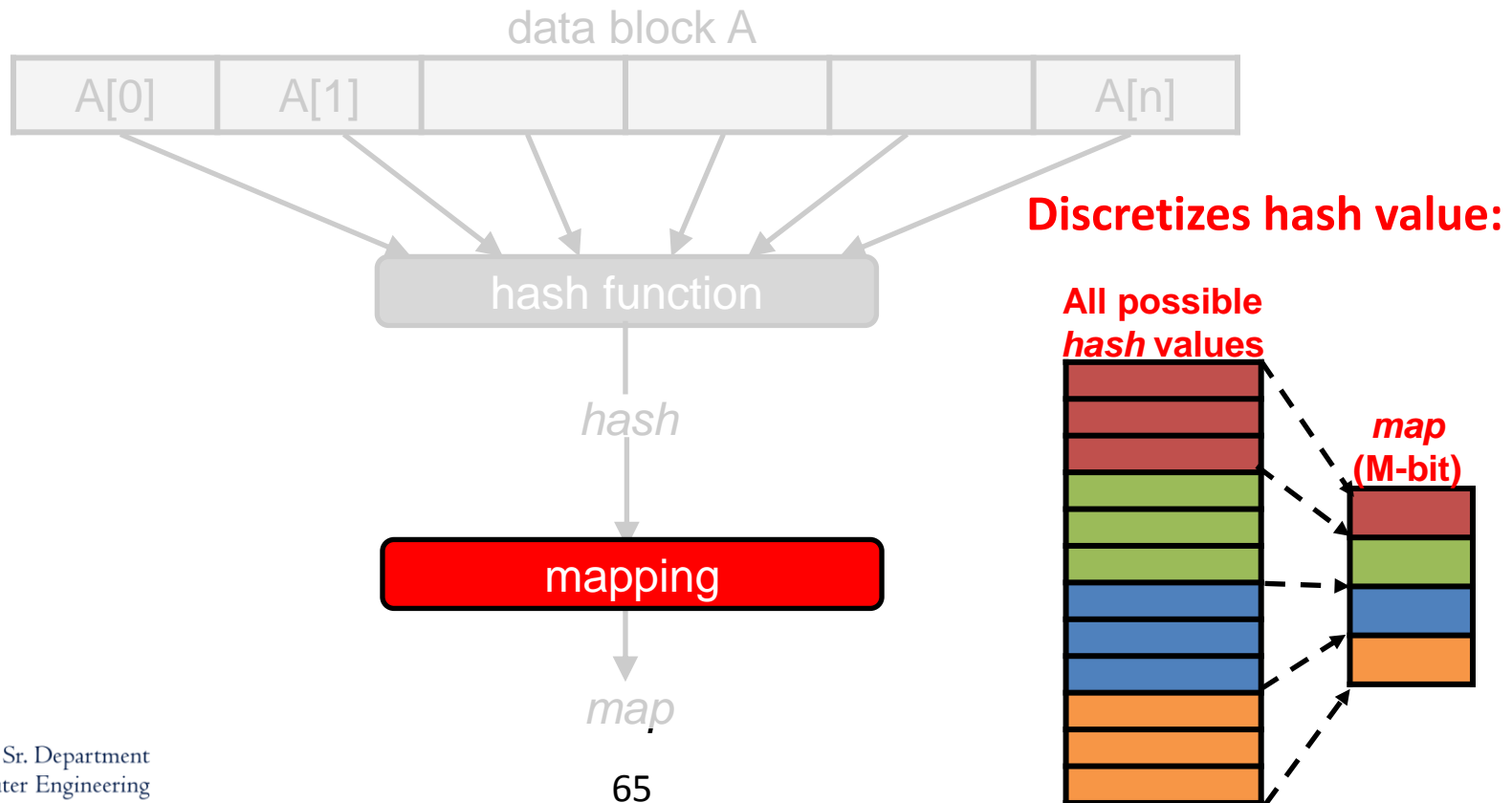


Aggregates values in block:

$$\text{hash} = \text{AVG}(A[0], \dots, A[n])$$

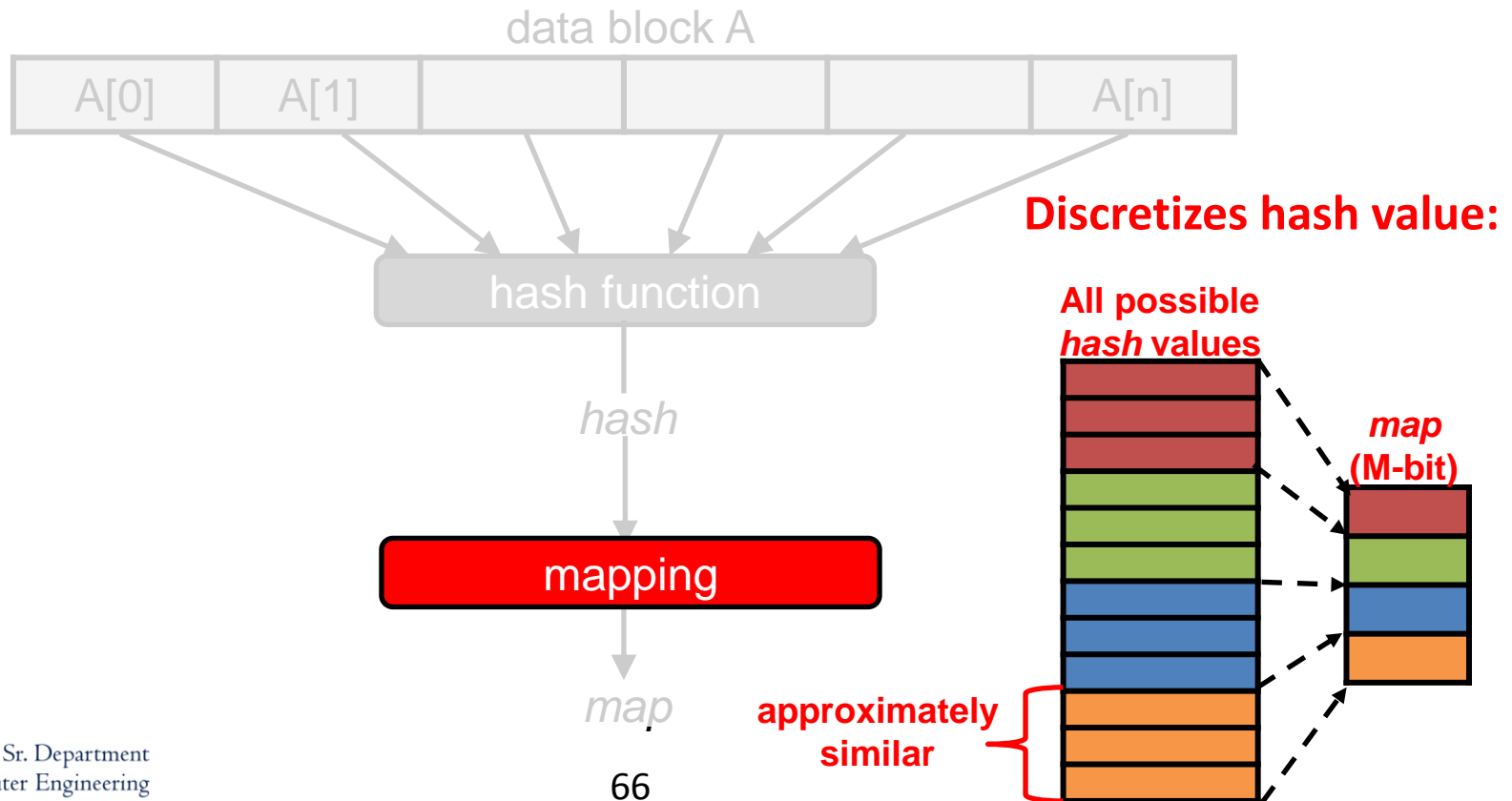
Doppelgänger Cache - Similarity Mapping

The **map** value represents the signature (or **likeness**) of a block. Blocks that generate the same map value are approximately similar.

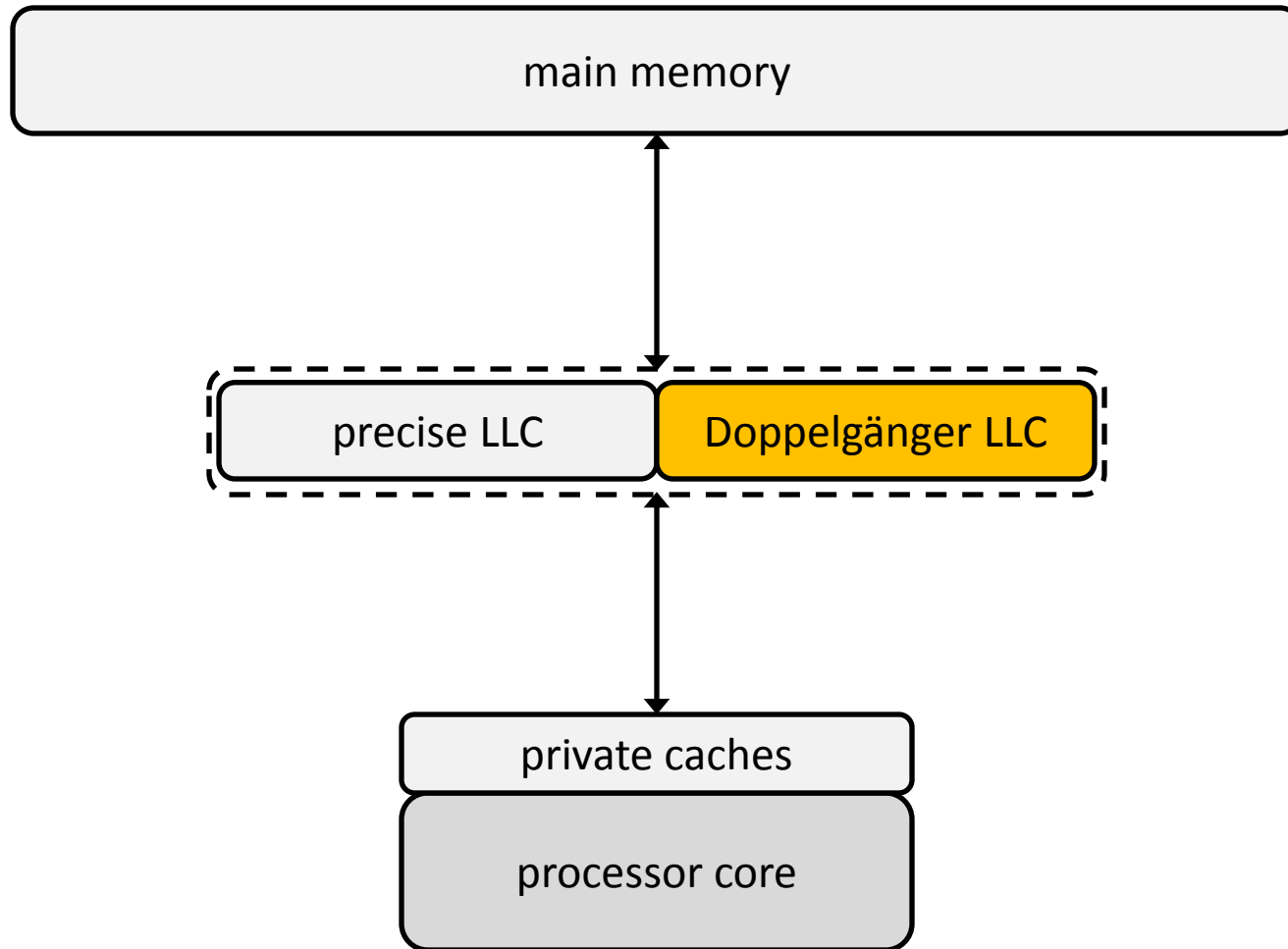


Doppelgänger Cache - Similarity Mapping

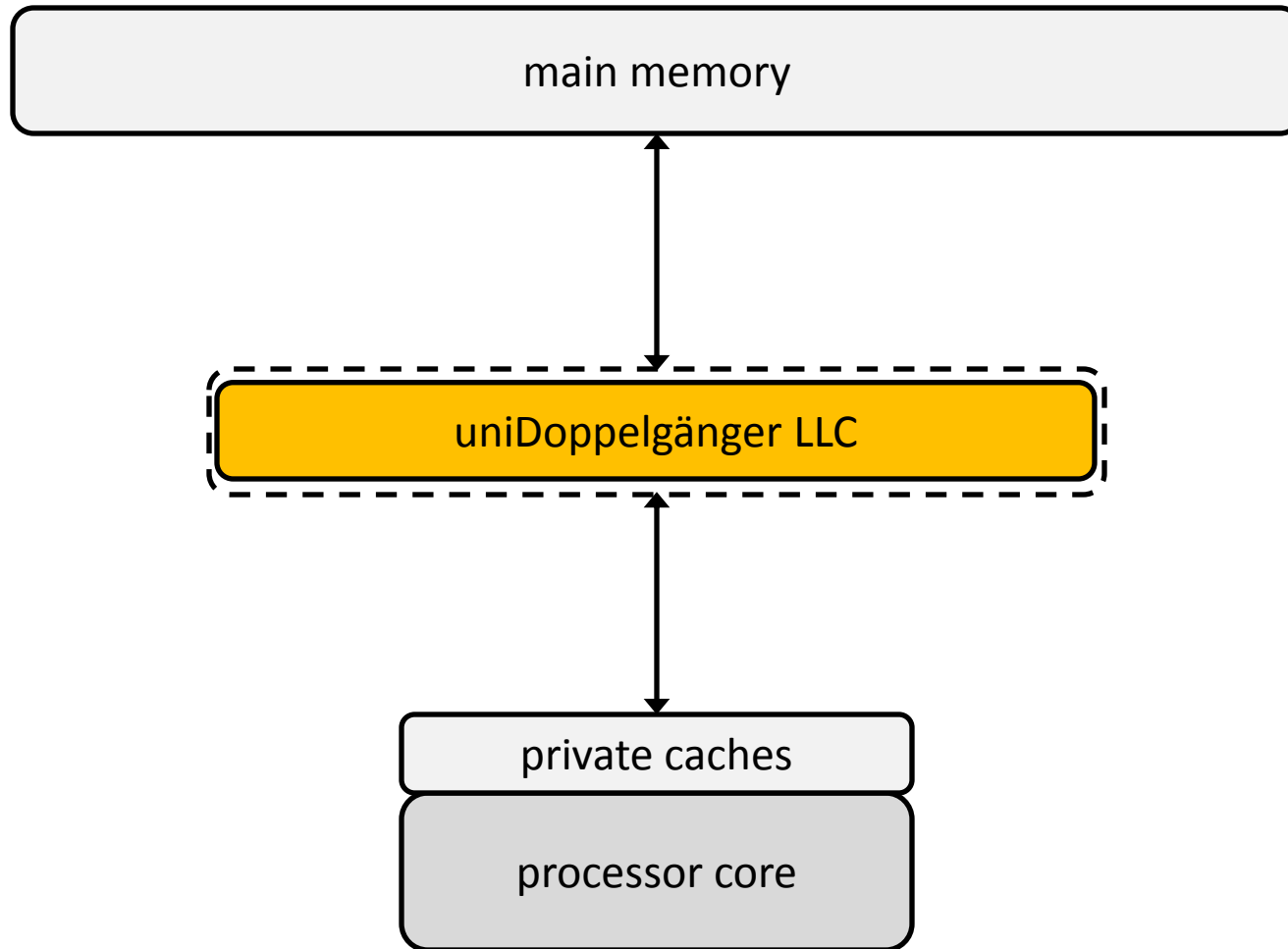
The **map** value represents the signature (or **likeness**) of a block. Blocks that generate the same map value are approximately similar.



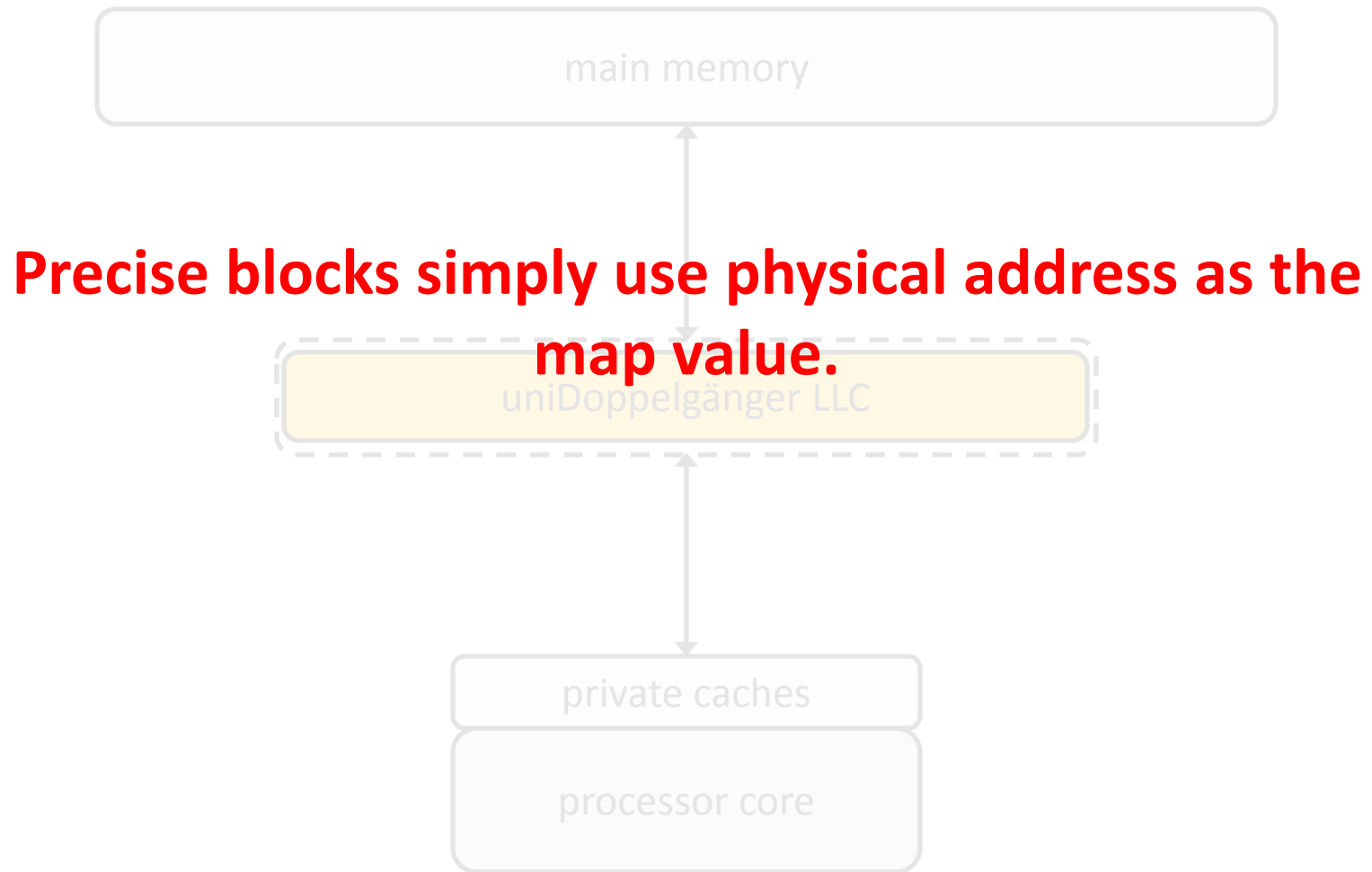
Doppelgänger Cache



uniDoppelgänger Cache



uniDoppelgänger Cache



Doppelgänger Cache

More details in paper:

- Cache writes, replacements and coherence.
- Details on hash functions and mapping.
- Sensitivity to size of map space and data array.
- Evaluation of uniDoppelgänger.

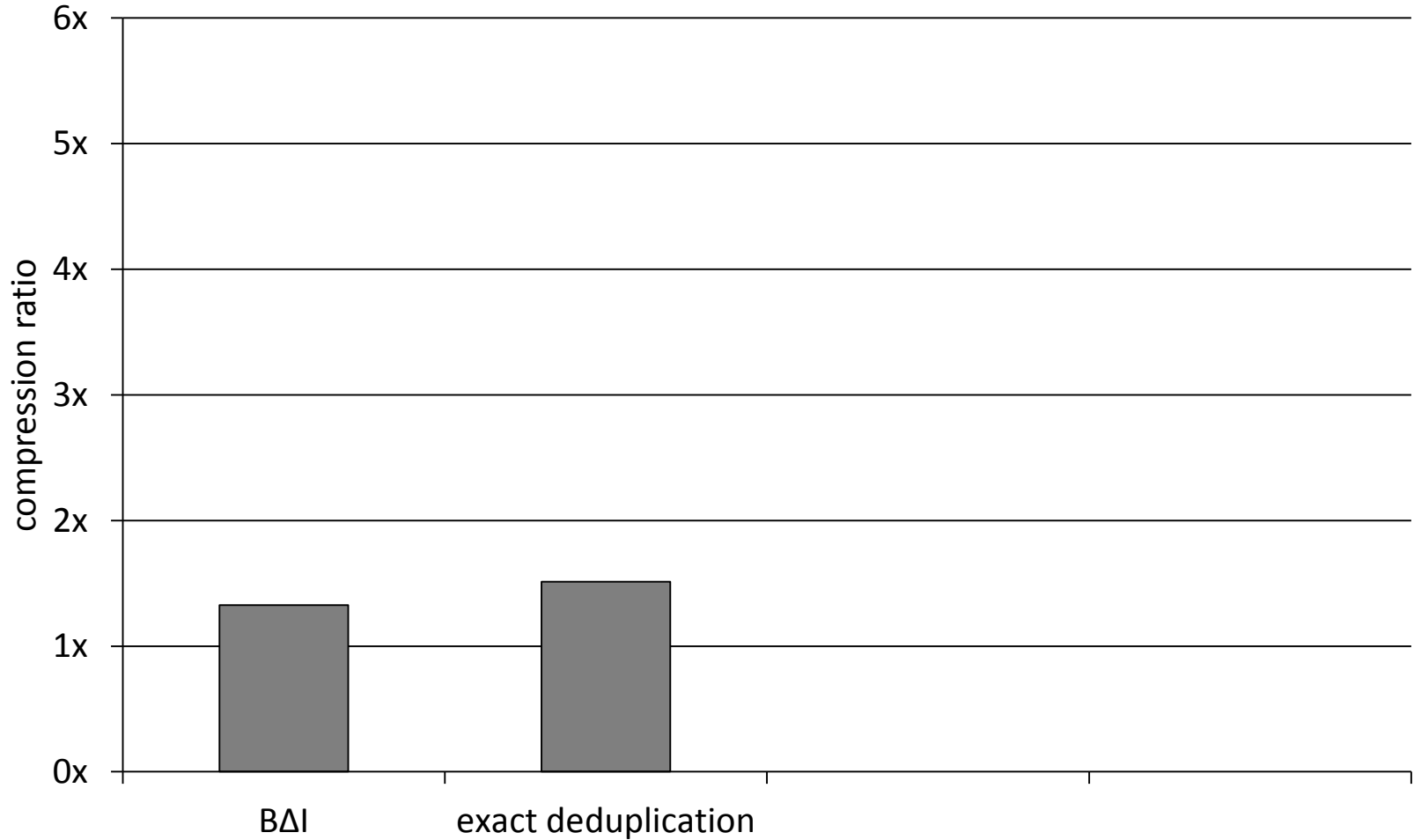
Outline

- Approximate Computing
 - Approximate Similarity
- Doppelgänger Cache
 - Cache Architecture
 - Similarity Mapping
- **Evaluation**

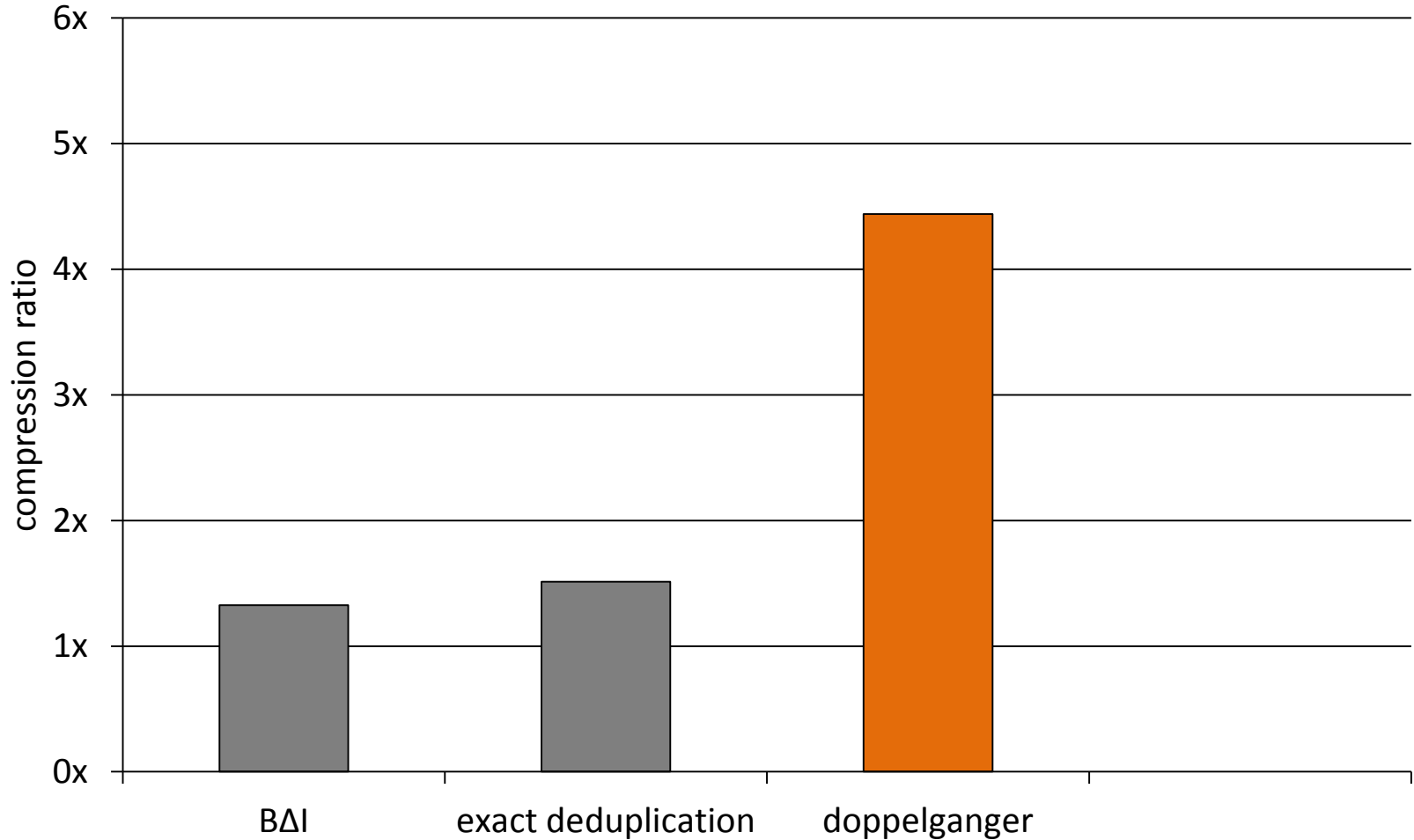
Evaluation

- **Applications:** PARSEC and AxBench
- **Performance:** Full-system cycle-level simulation
- **Error:** Pin simulation
- **Area and Energy:** CACTI
- **Configuration:**
 - 4 cores, private L1 and L2
 - 2MB shared LLC (1MB precise, 1MB Doppelgänger)
 - Doppelgänger: 14-bit similarity map, 1/4 data array

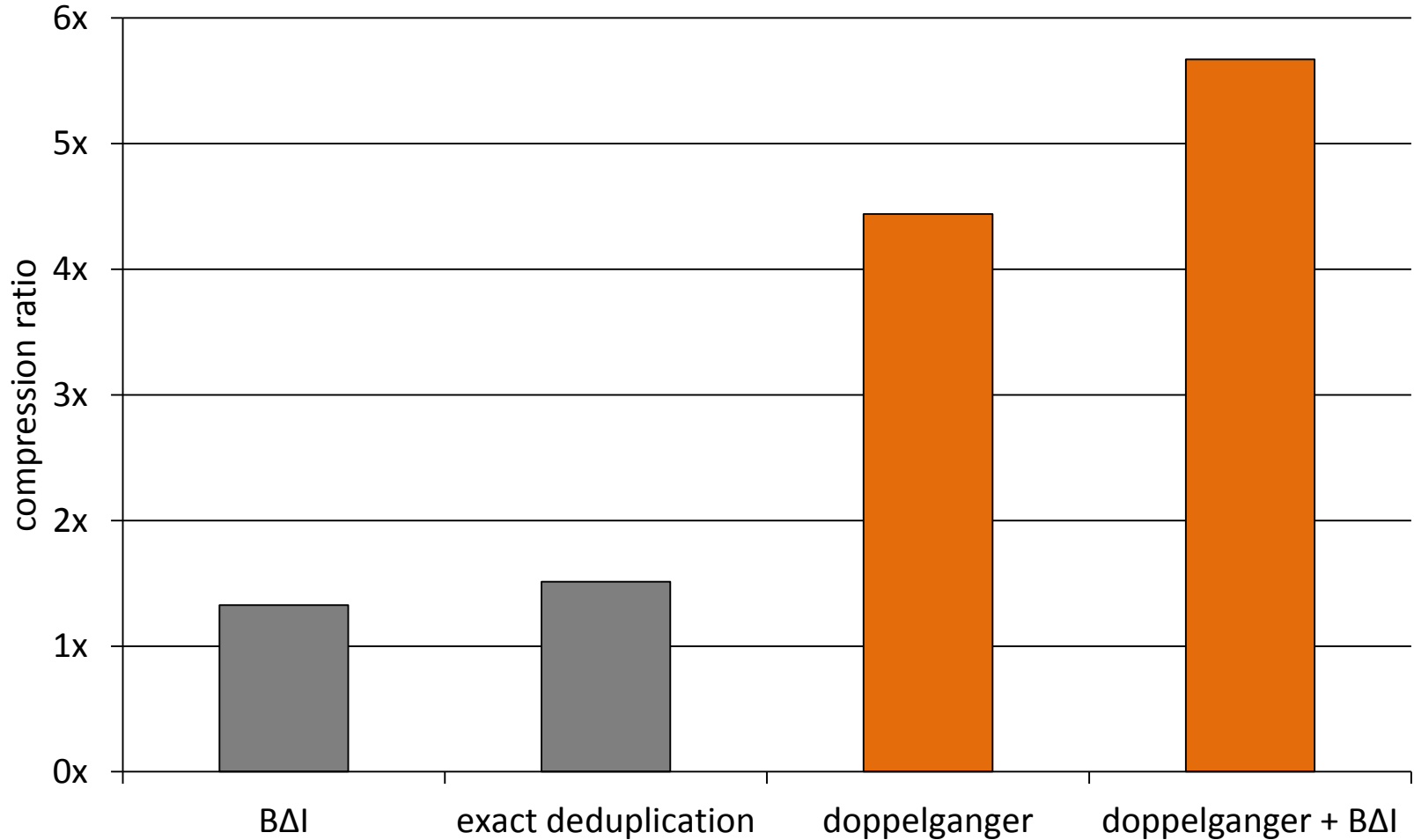
Evaluation - Compression Ratio



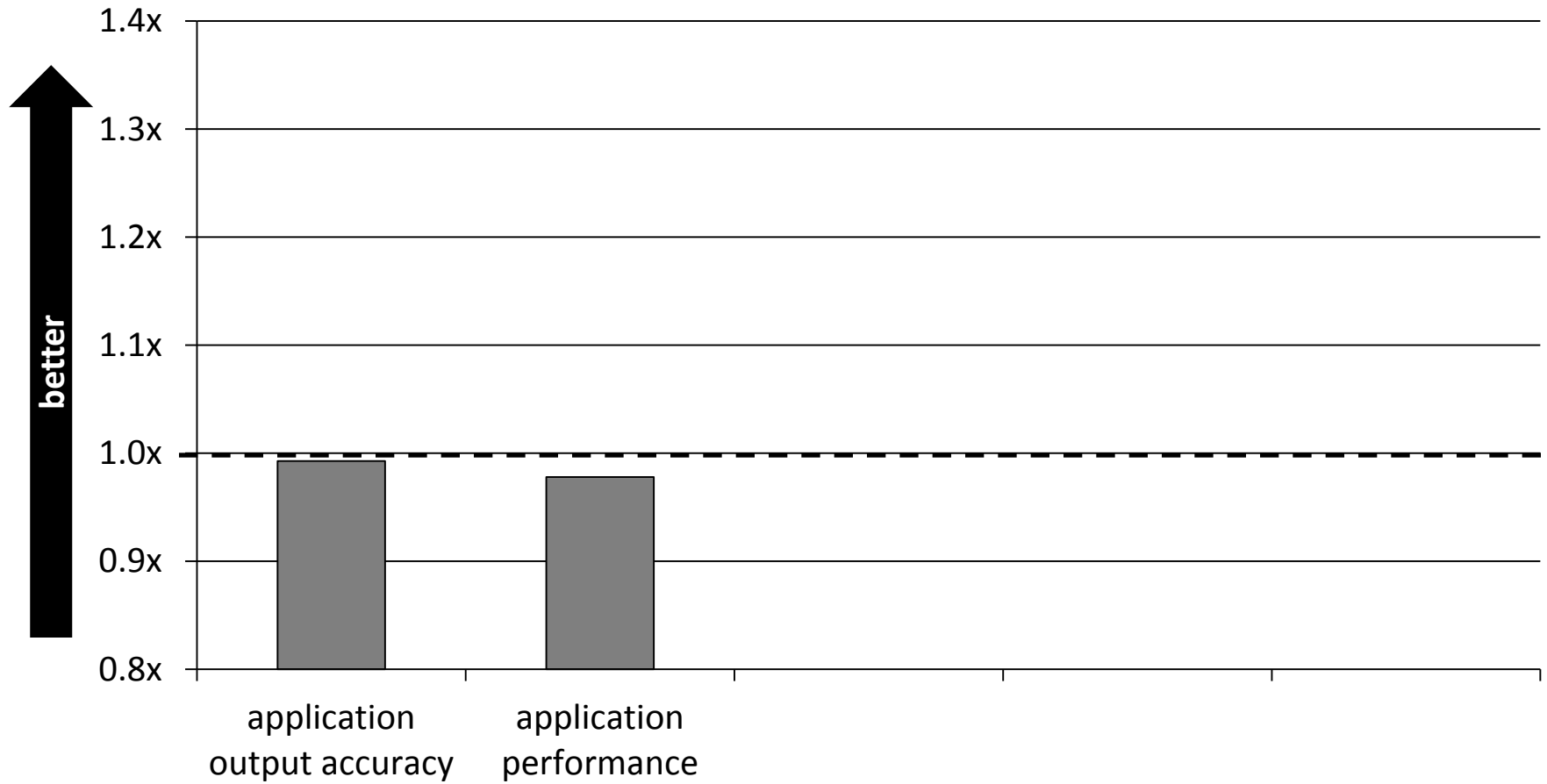
Evaluation - Compression Ratio



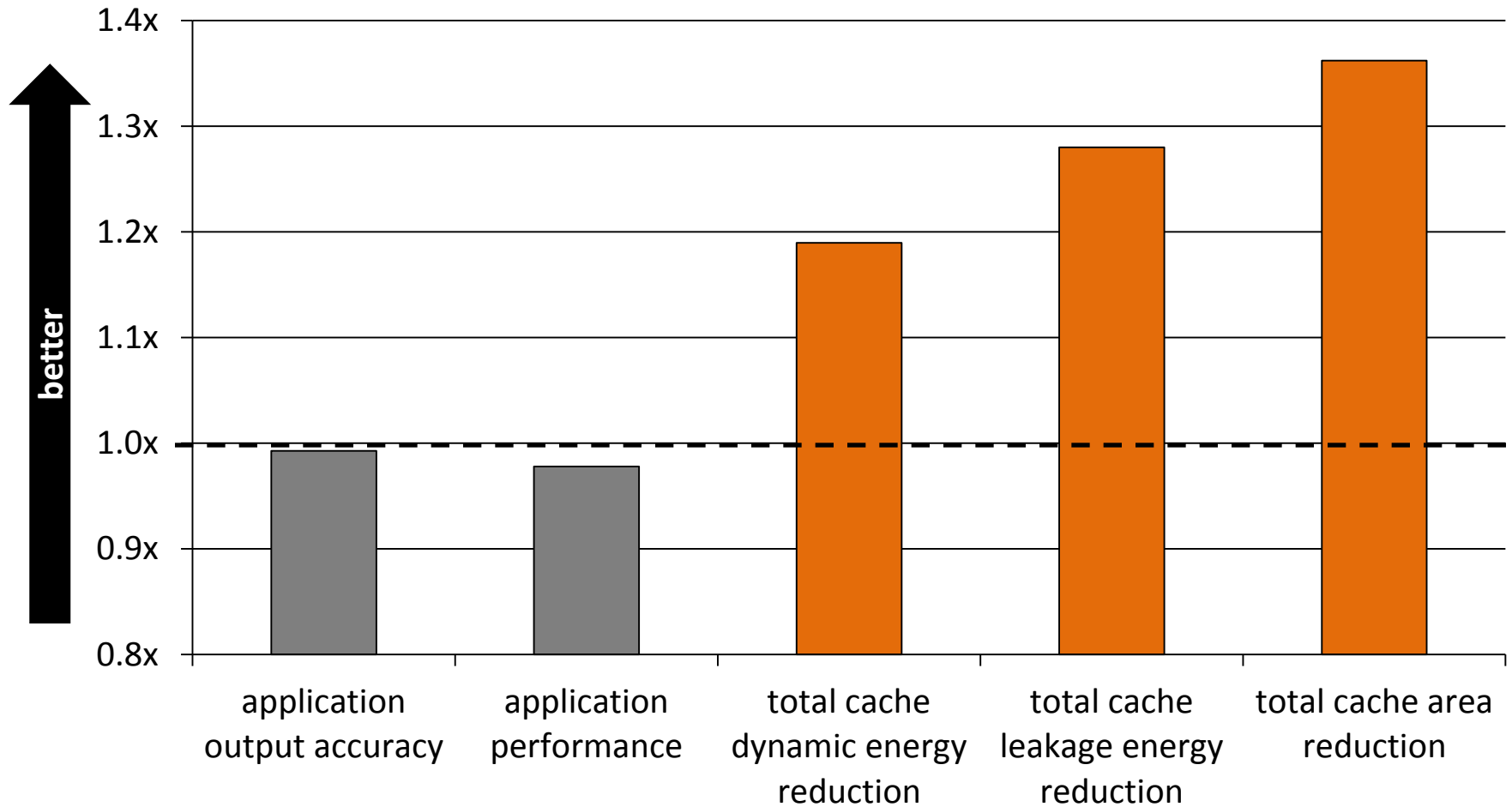
Evaluation - Compression Ratio



Evaluation



Evaluation



Conclusion

Doppelgänger Cache:

- Identifies **approximate similarity** in data block values.
 - **77%** cache storage savings of approximable data.
- Effectively compresses storage of approximately similar blocks.
 - **3x** better compression ratio than state-of-the-art techniques.
- Significantly reduces area and energy consumption.
 - Reduces total on-chip cache area by **1.36x**.



Thank you

Doppelgänger: A Cache for Approximate Computing

Joshua San Miguel

Jorge Albericio

Andreas Moshovos

Natalie Enright Jerger