# [544] Caching and PyArrow

Tyler Caraza-Harter

# Learning Objectives

- write cache-friendly code with PyTorch and PyArrow

- use memory mappings via PyArrow to access data that is larger than physical memory

- enable swapping to alleviate memory pressure

- configure Docker memory limits on physical memory used

# Outline

CPU: L1-L3

Demos: PyTorch+PyArrow...

OS (Operating System): Page Cache

Demos: PyArrow+Docker

# Granularity

If a process reads 1 byte and misses, *how much data should the CPU bring into the cache?*

- **too little:** we'll have many more misses if we read nearby bytes soon
- **too much:** wasteful to load data to cache that might never be accessed

L1-L3 cache data in units called **cache lines**

- modern CPUs typically 64 bytes (for example, 8 int64 numbers)
- M1/M2 uses 128

# Cache Lines and Misses

cacheline

**int64**
**int64**
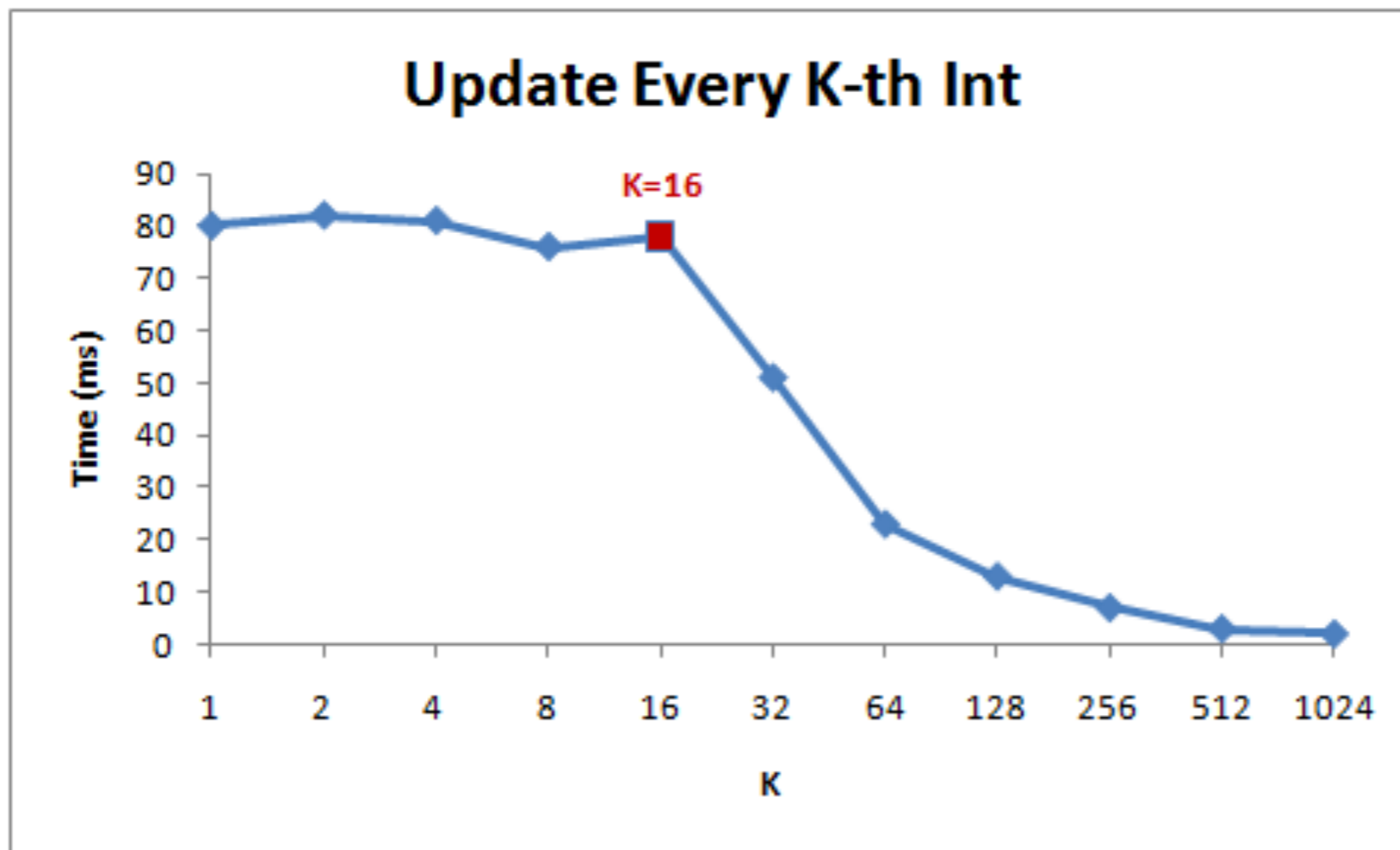**int64**
**int64**
int64
int64
int64
int64

cacheline

int64
int64
int64
int64
int64
int64
int64
int64

how many
misses?

cacheline

**int64**
int64
int64
int64
int64
int64
int64
int64

cacheline

**int64**
int64
int64
int64
int64
int64
int64
int64

how many
misses?

cacheline

int64
int64
int64
int64
**int64**
**int64**
**int64**
**int64**

cacheline

**int64**
**int64**
**int64**
**int64**
**int64**
**int64**
**int64**
**int64**

how many
misses?

# Example 1: Step and Multiply

```
for (int i = 0; i < arr.Length; i += K) arr[i] *= 3;
```

**Update Every K-th Int**



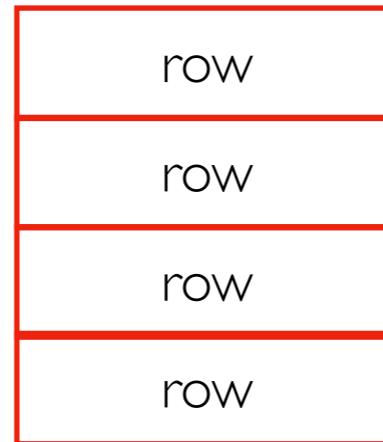[Gallery of Processor Cache Effects](http://igoro.com/archive/gallery-of-processor-cache-effects/)
http://igoro.com/archive/gallery-of-processor-cache-effects/

# Example 2: Matrices
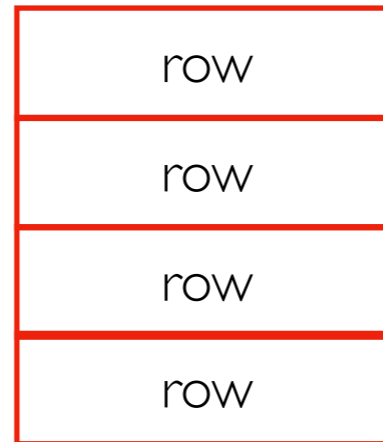
matrix of numbers
**logically**, 2-dimensional

| row |
|-----|
| row |
| row |
| row |

**physically**, those rows are arranged along
1-dimension in the virtual address space

| | code | | row | row | row | row | ... | stack | |
|-|------|-|-----|-----|-----|-----|-----|-------|-|

0                                                                N

virtual address
spaces

# Example 2: Matrices

matrix of numbers
**logically**, 2-dimensional

| row |
| --- |
| row |
| row |
| row |

summing over row:
data consolidated over few cache lines

| | code | | row | row | row | row | ... | stack | |
|---|---|---|---|---|---|---|---|---|---|

0 ............................................................................................................................ N
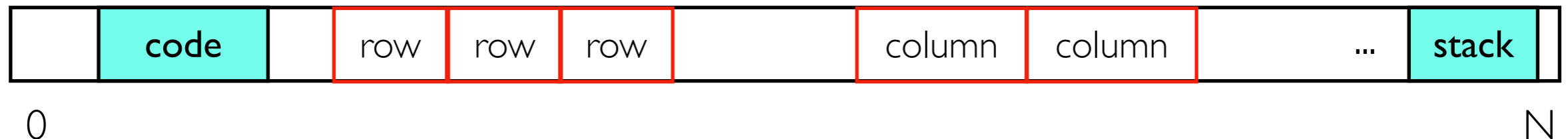
summing over column: each number is in its own cache line and triggers a cache miss

# PyTorch: Controlling Layout with Transpose

for efficiency, transpose doesn't actually move/copy data,
meaning we can get fast column sum by (a) putting
column data in rows and (b) transposing

```
torch.tensor([[1,2],
              [3,4],
              [5,6]])
```
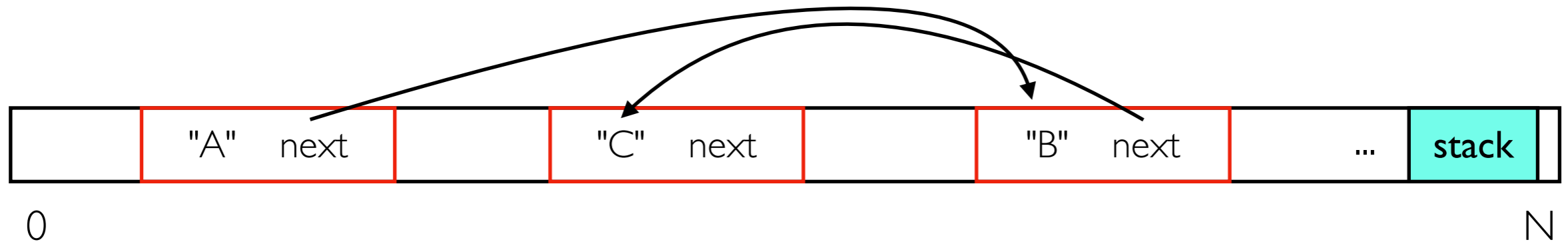
```
torch.tensor([[1,3,5],
              [2,4,6]]).T
```

| | code | | row | row | row | | column | column | ... | stack | |
|---|---|---|---|---|---|---|---|---|---|---|---|

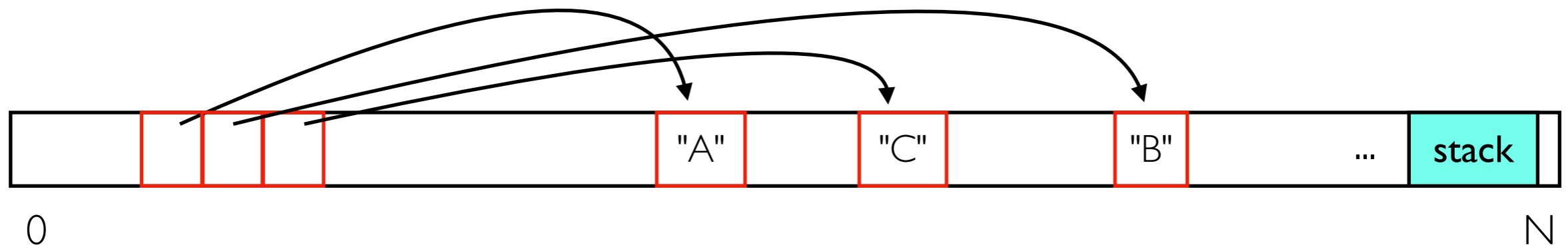0                                                                              N

any calculations on the two tensors will produce the same results,
but they'll each be faster for different access patterns!
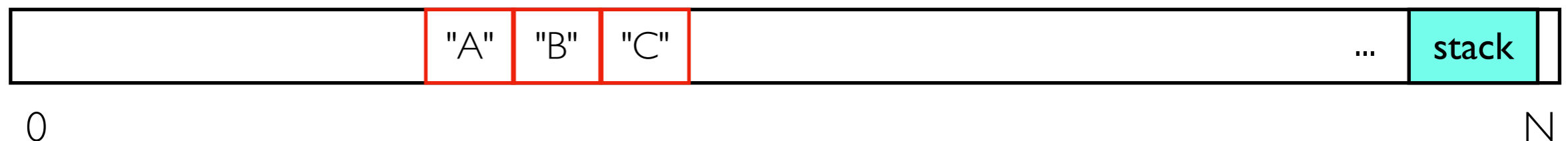
# Example 3: Ordered Collections of Strings
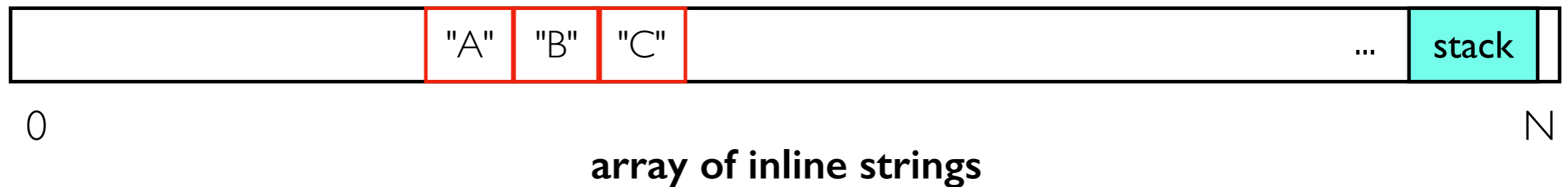
which layout is most cache friendly?



**linked list**
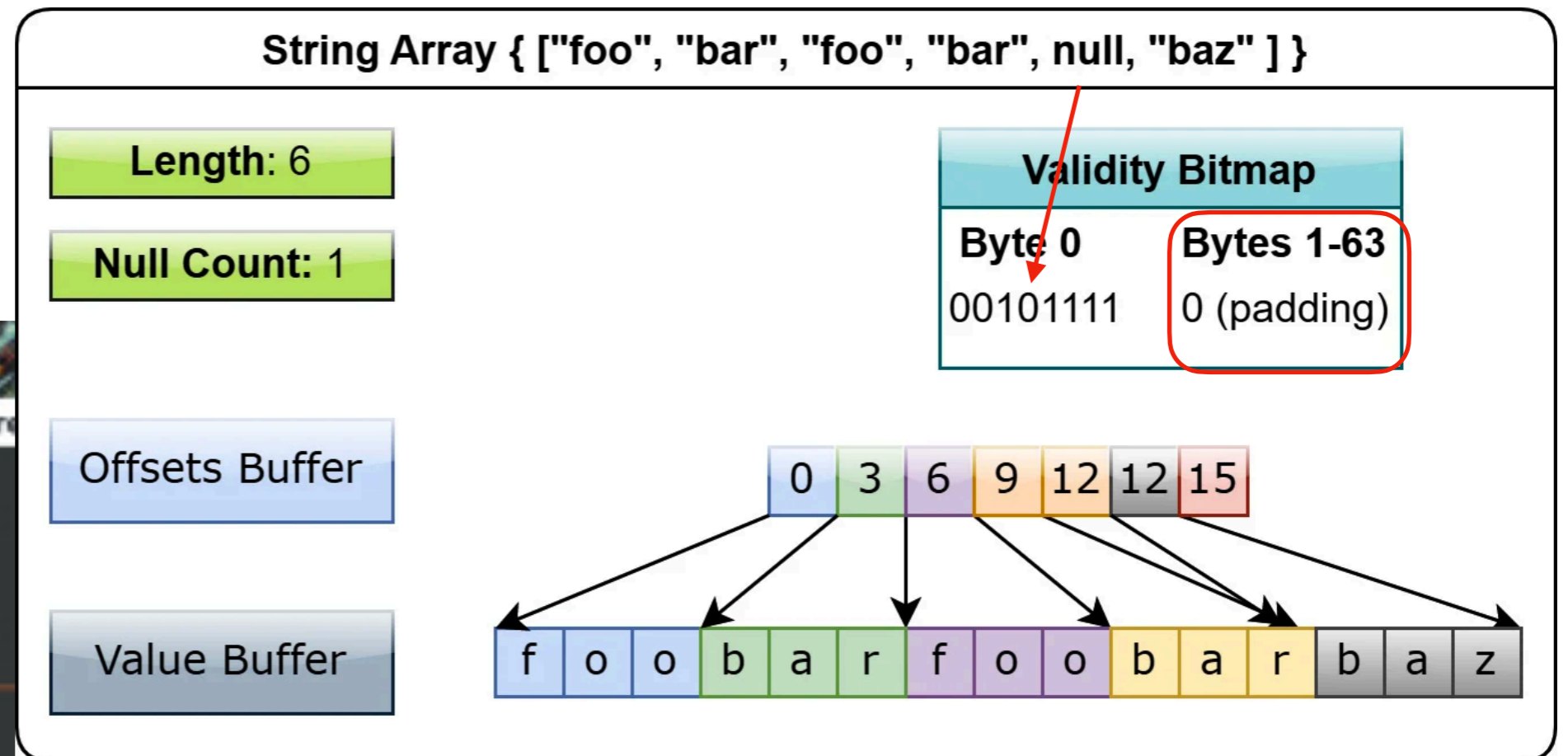
**array of references to strings**

**array of inline strings**

# Example 3: Ordered Collections of Strings

how to tell the end of one string from the start of the next?
how to jump immediately to string at index i?
how support null/None?

| | | "A" | "B" | "C" | | ... | stack | |
|---|---|---|---|---|---|---|---|---|

0                                             N
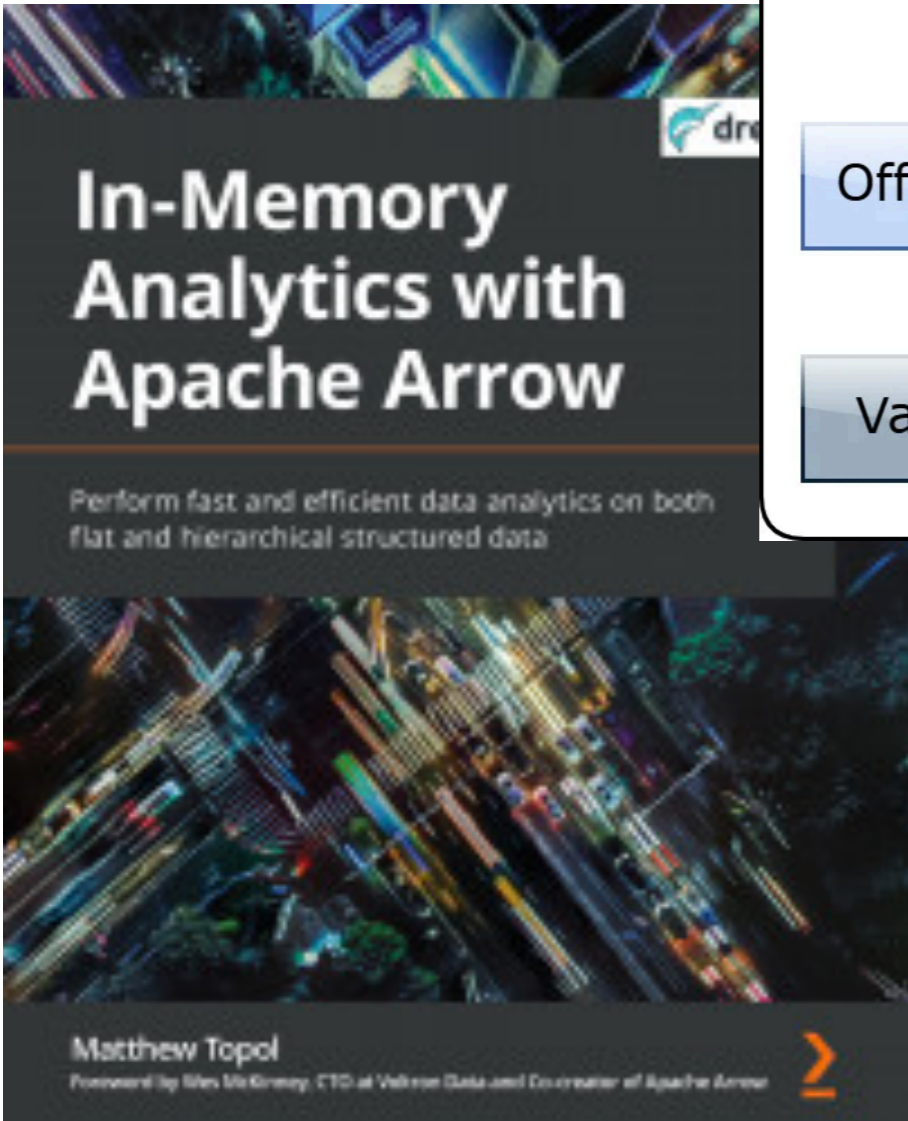
**array of inline strings**

# PyArrow String Array Data Structure



data is packed into fewest possible cache lines

- collection of named arrays is a Table
- arrays for different types, each cache friendly
- null support for types like int (not forced into floats)

https://www.packtpub.com/product/in-memory-analytics-with-apache-arrow/9781801071031

# Outline

CPU: L1-L3

Demos: PyTorch+PyArrow...

OS (Operating System): Page Cache

Demos: PyArrow+Docker
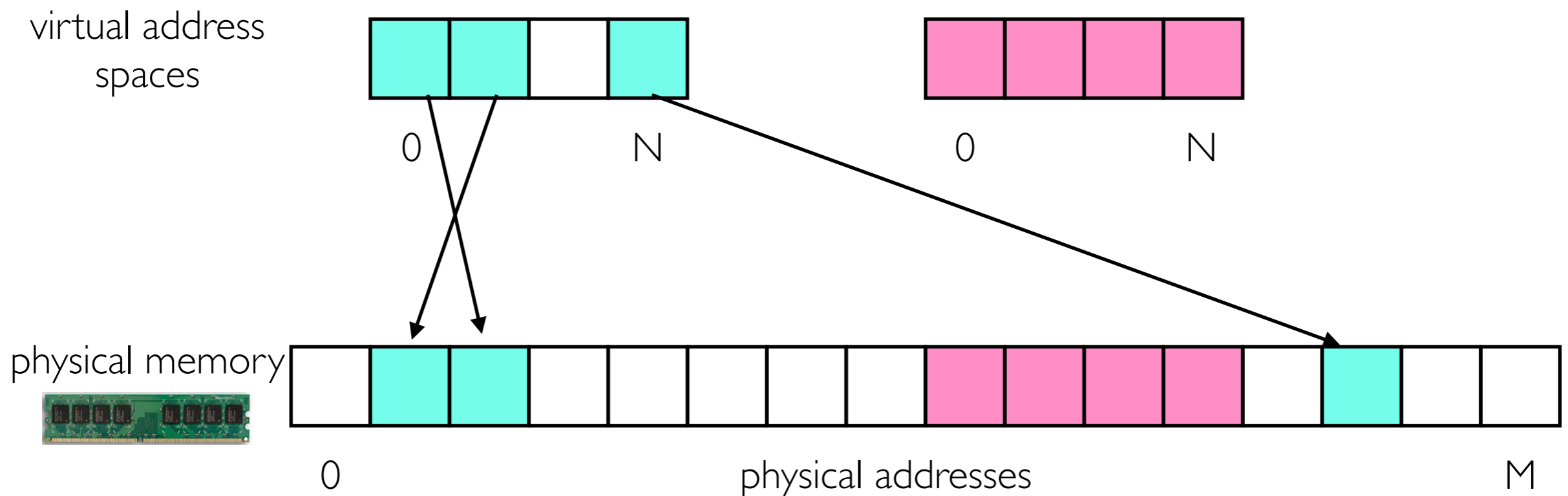
# Outline

CPU: L1-L3

Demos: PyTorch+PyArrow...

OS (Operating System): Page Cache

Demos: PyArrow+Docker

# Review Processes and Address Spaces

Address spaces

- Each process has it's own virtual address space
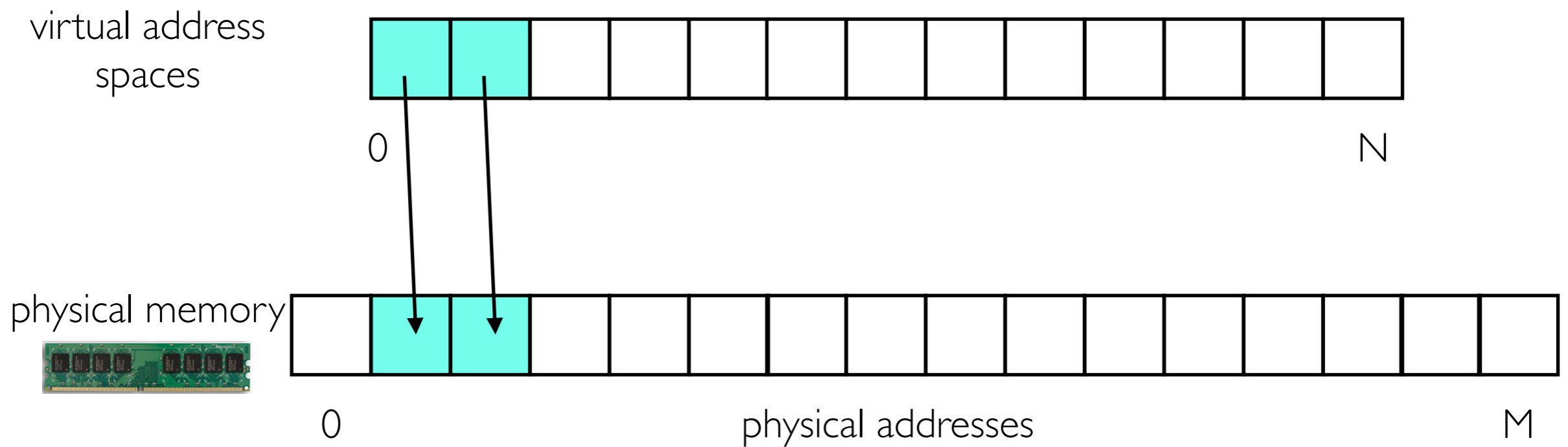- pages (usually 4 KB) of memory are mapped to physical memory

# mmap (Memory Map)

An mmap call can add new regions to a virtual address space. Two varities:
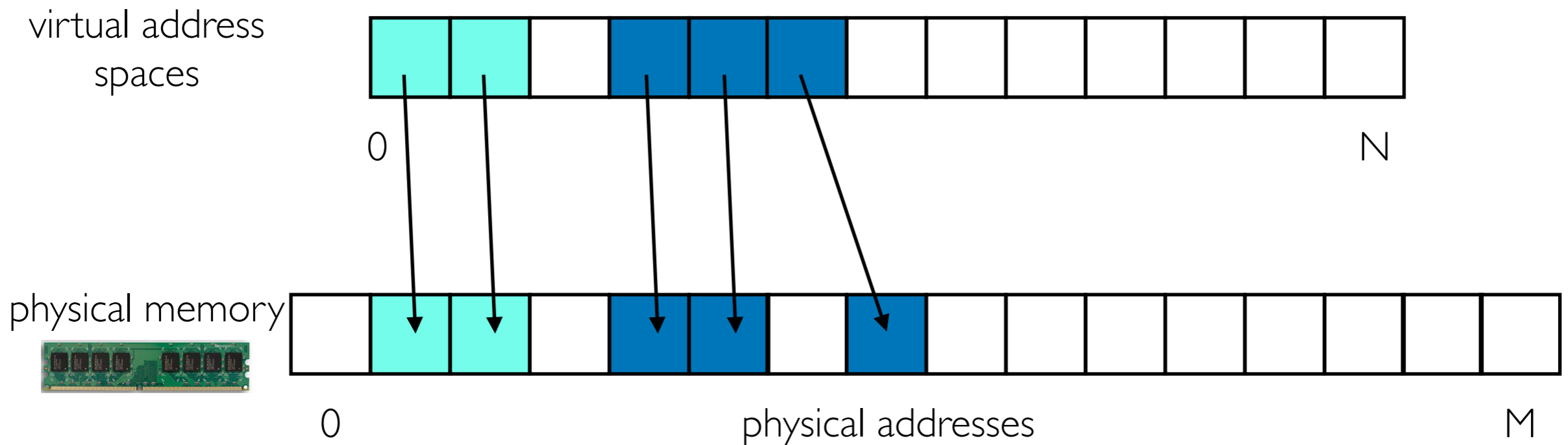- anonymous
- backed by a file

virtual address spaces

0                                                                    N

physical memory

0                         physical addresses                         M

# Anonymous mmap

An mmap call can add new regions to a virtual address space. Two varities:

- anonymous
- backed by a file

anonymous      3 pages

```
import mmap
mm = mmap.mmap(-1, 4096*3)
```



virtual address spaces

0                                    N

physical memory

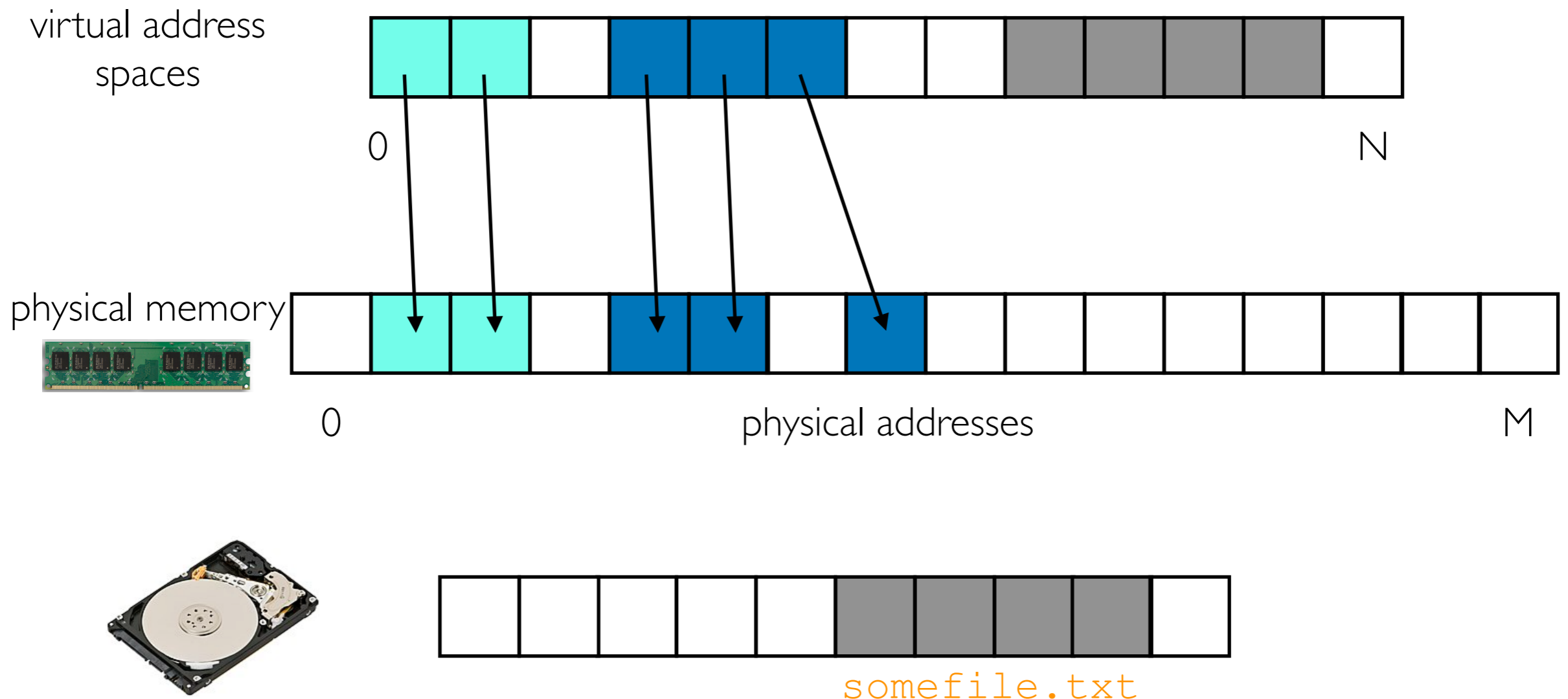0            physical addresses            M

- Python (and other language runtimes) will mmap some anonymous memory when they need more heap space
- this will be used for Python objects (ints, lists, dicts, DataFrames, etc.)

# File-Backed mmap

An mmap call can add new regions to a virtual address space. Two varities:
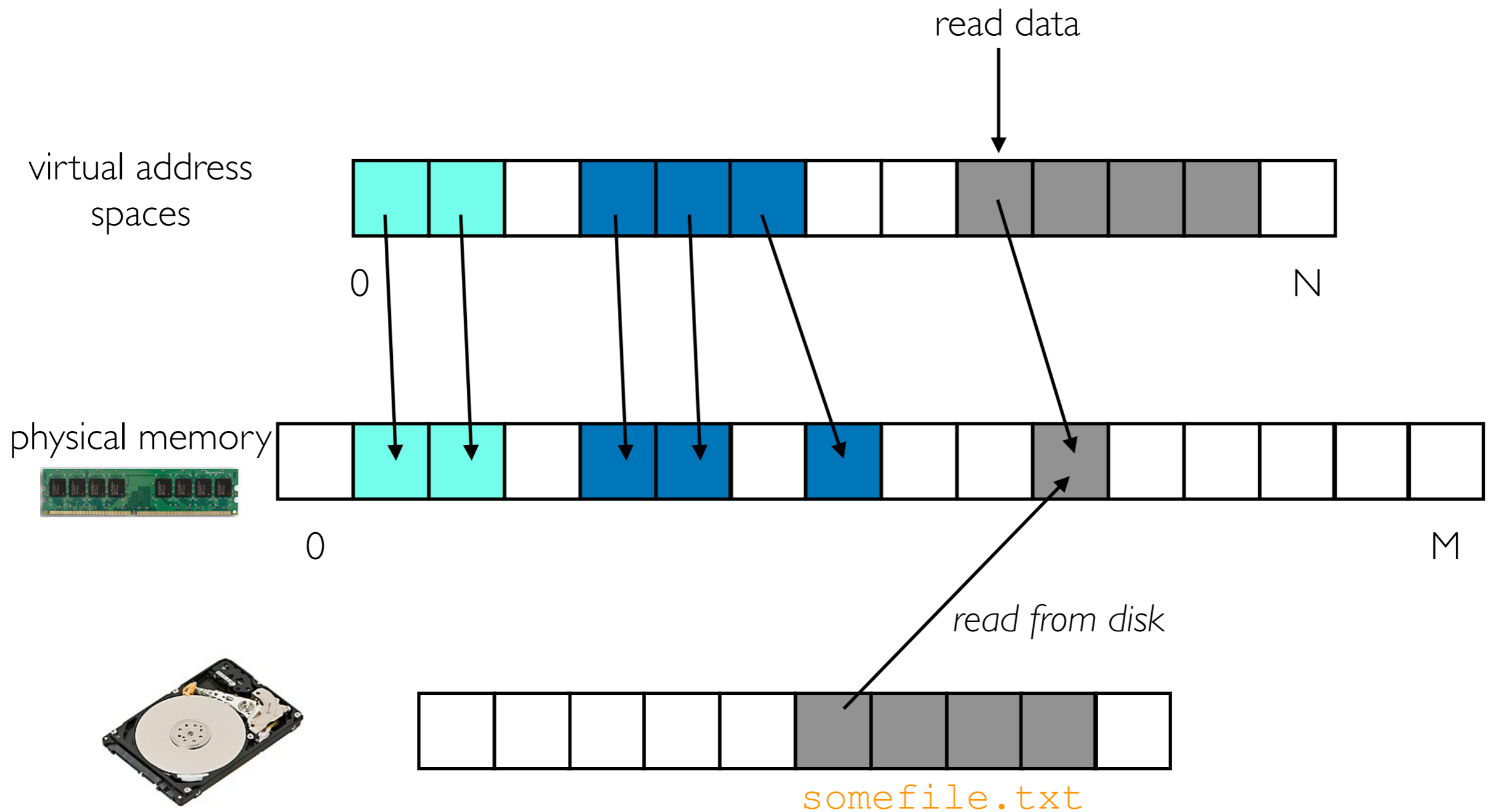
- anonymous
- backed by a file

```
import mmap
f = open("somefile.txt", mode="rb")
mm = mmap.mmap(f.fileno(), 0,    # 0 means all
                access=mmap.ACCESS_READ)
```

virtual address
spaces

0                                                                N

physical memory

0                    physical addresses                          M

somefile.txt

# File-Backed mmap

An mmap call can add new regions to a virtual address space. Two varities:
- anonymous
- backed by a file

read data

virtual address spaces

0                                    N

physical memory

0                                    M

read from disk

somefile.txt

# File-Backed mmap

An mmap call can add new regions to a virtual address space. Two varities:
- anonymous
- backed by a file

read data

virtual address spaces

0                                    N

physical memory

0                                                M

read from disk

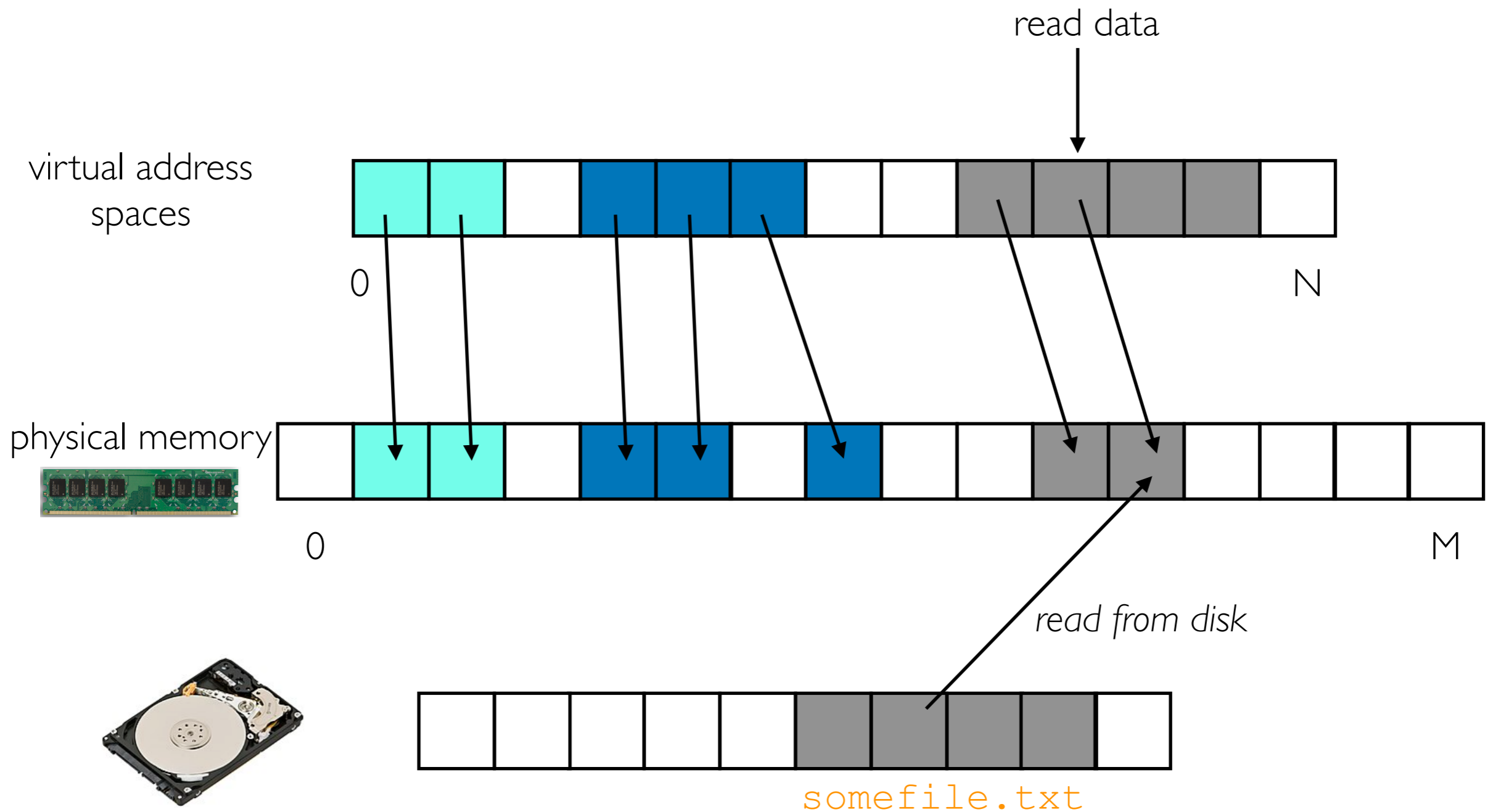somefile.txt

# File-Backed mmap

An mmap call can add new regions to a virtual address space. Two varities:

- anonymous
- backed by a file

- **virtual** memory used: 9*pagesize = 36 KB
- **physical** memory used: 7*pagesize = 28 KB

virtual address spaces

0                                                                     N

physical memory
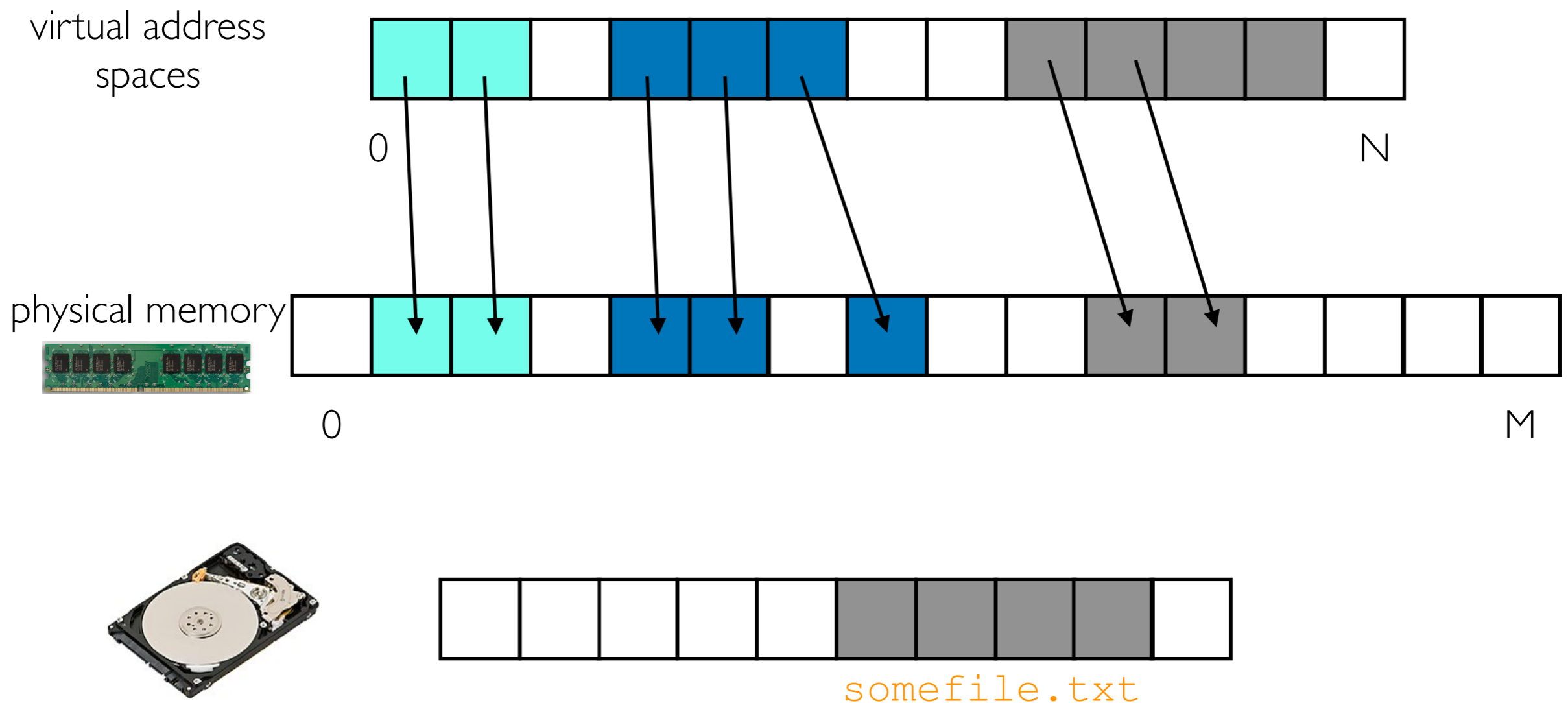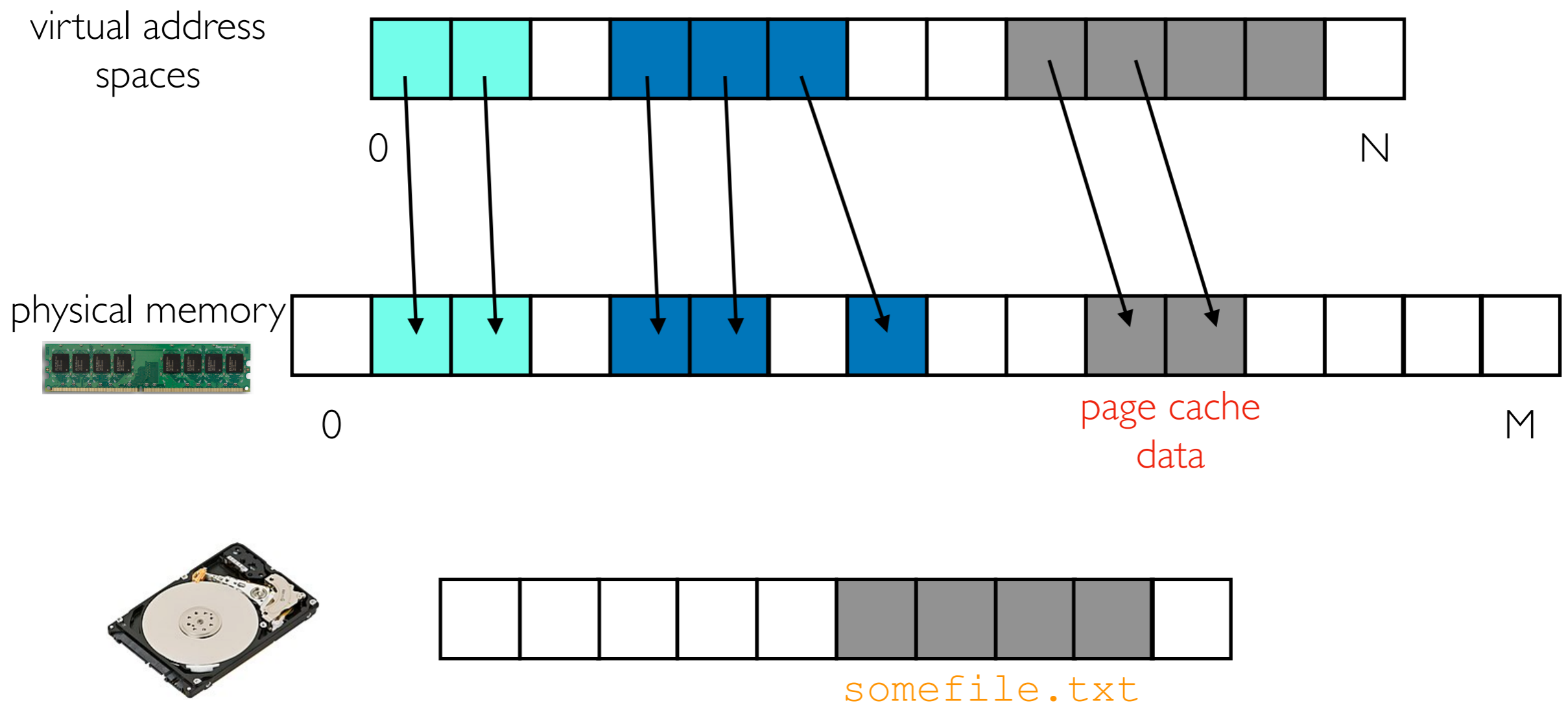
0                                                                     M

somefile.txt

# File-Backed mmap

An mmap call can add new regions to a virtual address space. Two varities:

- anonymous
- backed by a file

- data loaded for accesses to file-backed mmap regions are part of the "page cache"

virtual address spaces

0                                                                          N

physical memory

0                                                    page cache
                                                     data                 M

somefile.txt

# File-Backed mmap

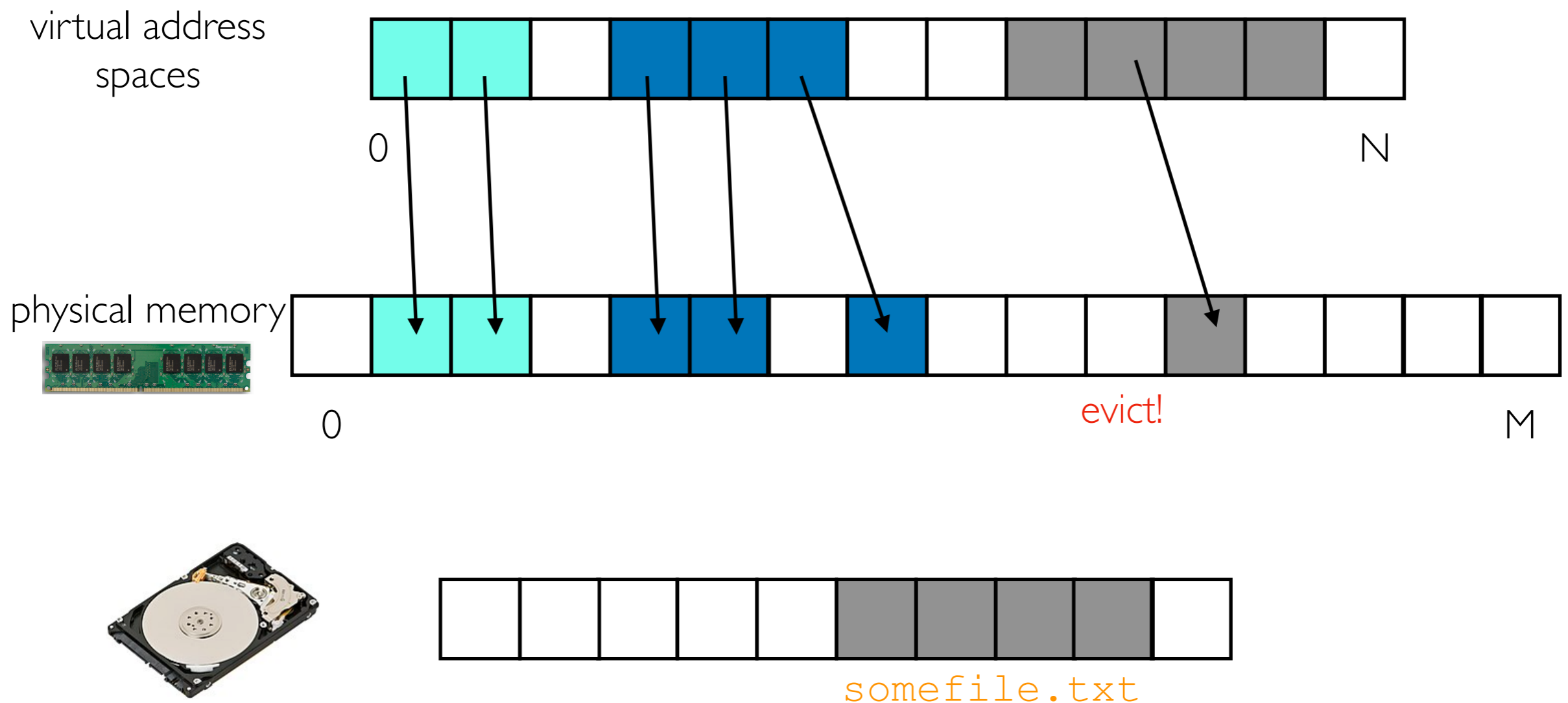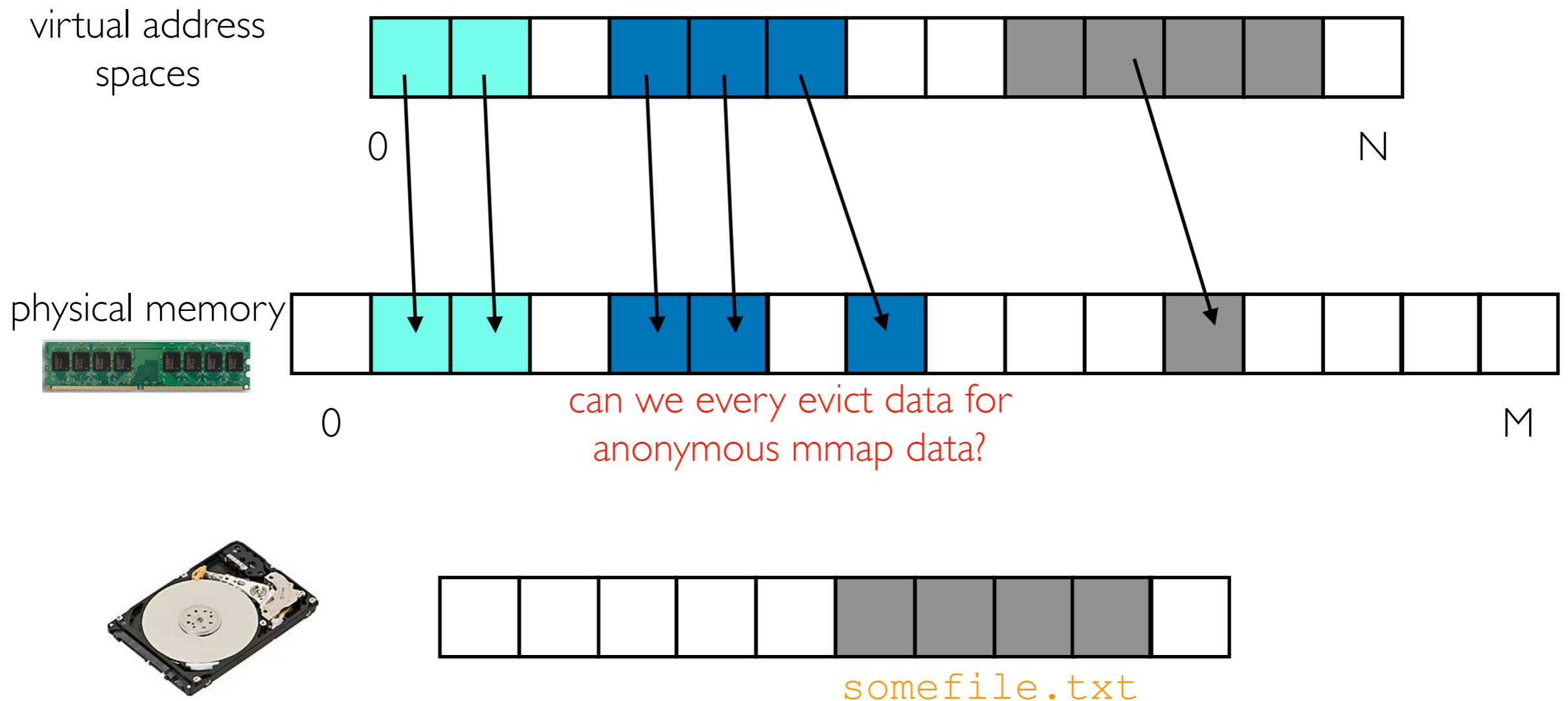An mmap call can add new regions to a virtual address space. Two varities:

- anonymous
- backed by a file

- data loaded for accesses to file-backed mmap regions are part of the "page cache"
- it works like a cache because there is another copy on disk, so we can evict under memory pressure

virtual address spaces

0                                                              N

physical memory

0                                   evict!                      M

somefile.txt

# Swap Space

An mmap call can add new regions to a virtual address space. Two varities:
- anonymous
- backed by a file

virtual address spaces
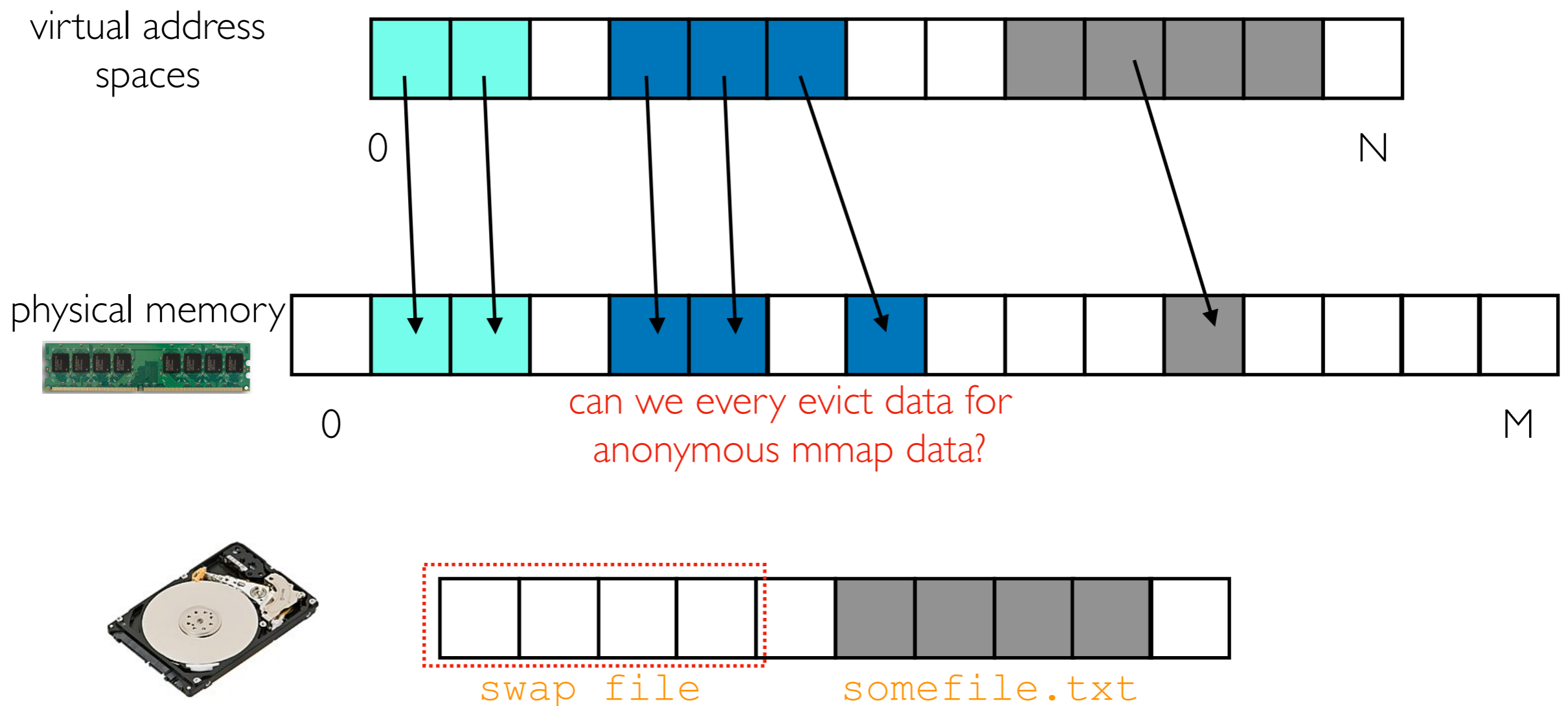
0

N

physical memory

0

can we every evict data for anonymous mmap data?

M

somefile.txt

# Swap Space

An mmap call can add new regions to a virtual address space. Two varities:

- anonymous
- backed by a file
- we can create same space (a swap file) to which the OS can evict data from anonymous mappings

virtual address spaces

0                                                                   N

physical memory

0                                                                   M

can we every evict data for
anonymous mmap data?

swap file          somefile.txt

# Swap Space
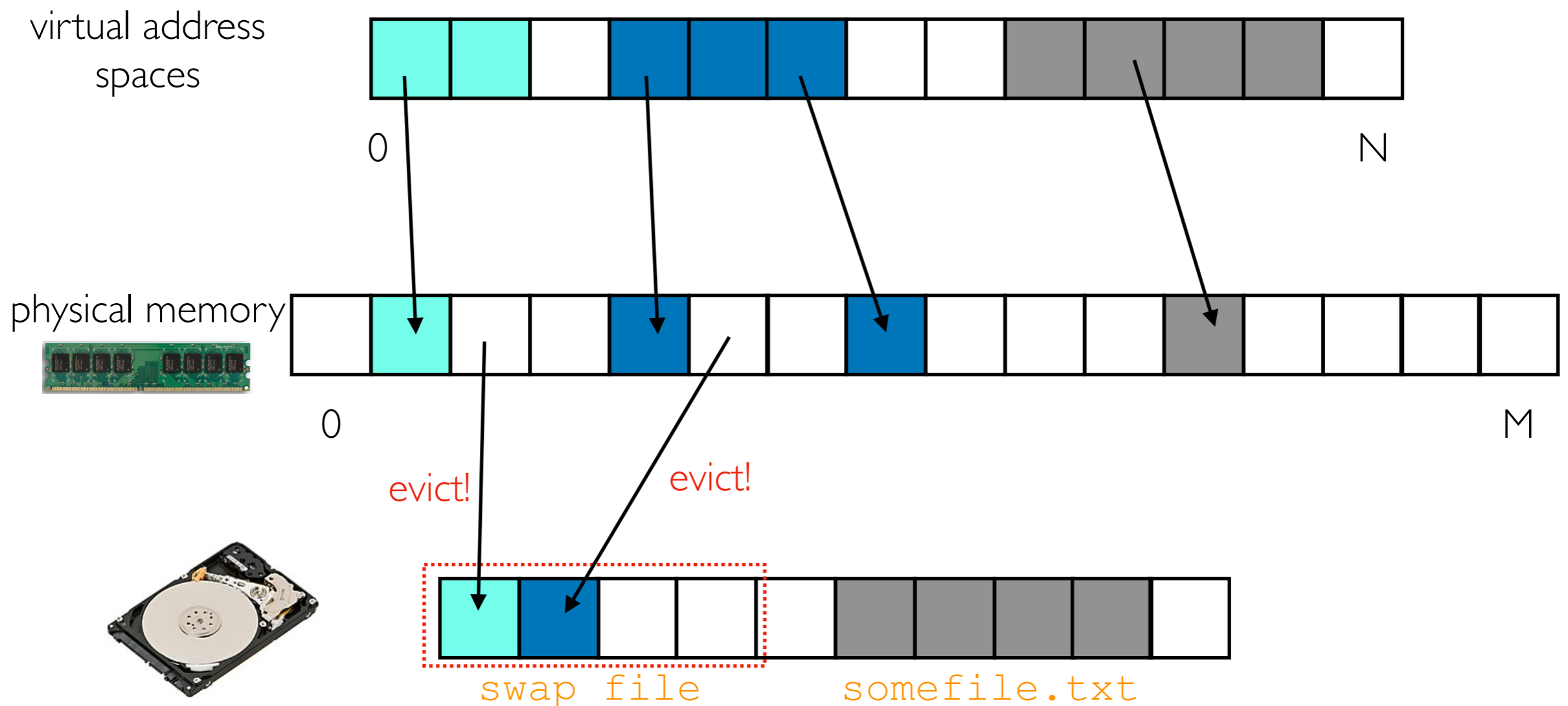
An mmap call can add new regions to a virtual address space. Two varities:

- anonymous
- backed by a file

- we can create same space (a swap file) to which the OS can evict data from anonymous mappings
- of course, if we access these virtual addresses again, it will be slow to bring the data back

virtual address spaces

0                                                                    N

physical memory

0                                                                    M

evict!          evict!

swap file        somefile.txt

# Outline

CPU: L1-L3

Demos: PyTorch+PyArrow...

OS (Operating System): Page Cache

Demos: PyArrow+Docker