# Research Topics in Database Management Bigger, Faster, and Stronger Systems Niv Dayan

# Who am I?

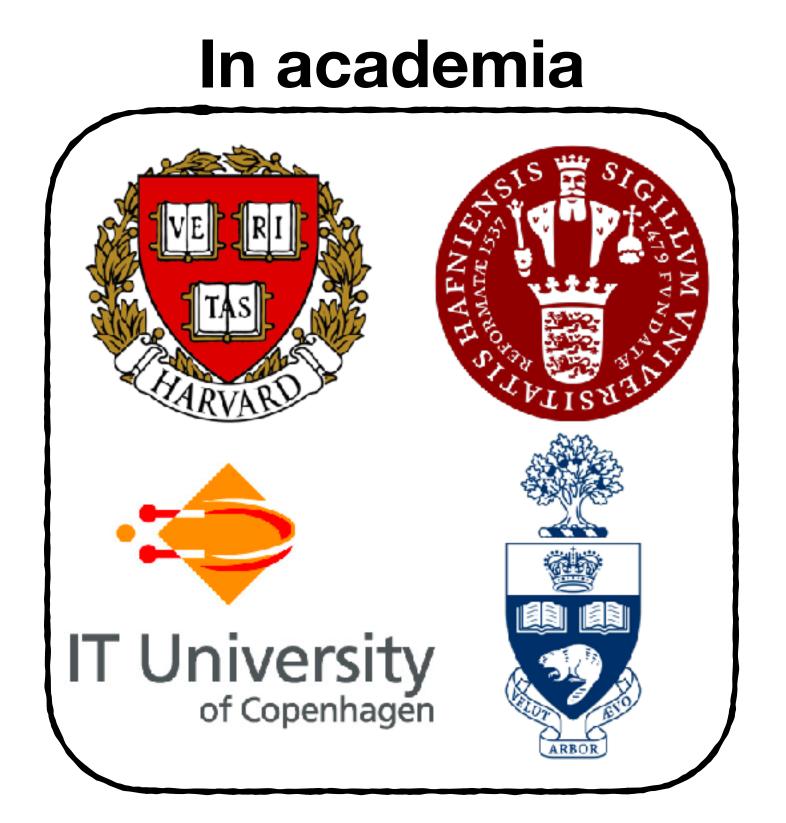
> 12 years of research experience in data structures & algorithms for databases



https://www.nivdayan.net/

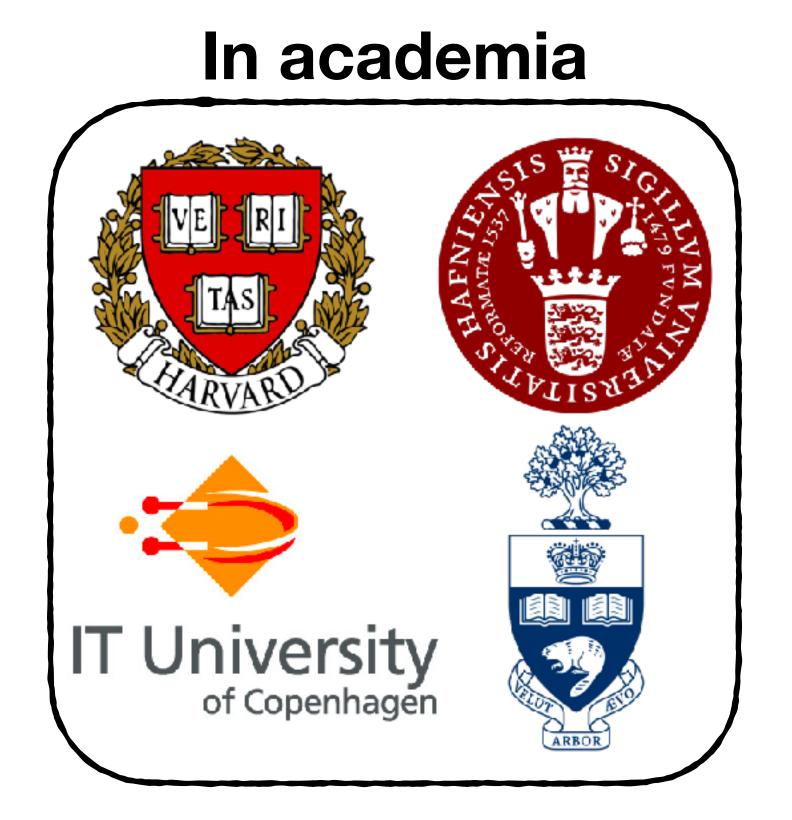
### Who am I?

> 10 years of research experience in data structures & algorithms for databases



### Who am I?

> 10 years of research experience in data structures & algorithms for databases



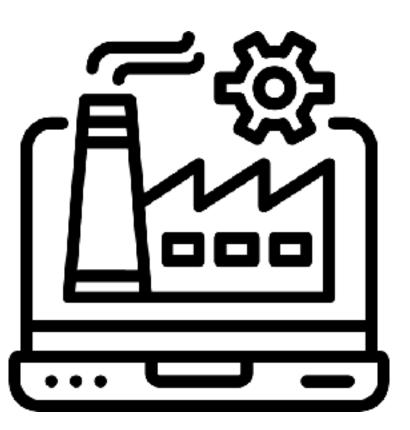


### This course combines both

**Theory** 



**Practice** 



For data structures & algorithms for databases

# Who are you?



# Who are you?

Grad Undergrad

### **Introduction to Databases**

Relational algebra, data modeling, SQL

### **Introduction to Databases**

Relational algebra, data modeling, SQL

### **Operating Systems**

Concurrency & synchronization File systems, virtual memory

### **Introduction to Databases**

Relational algebra, data modeling, SQL

### **Operating Systems**

Concurrency & synchronization File systems, virtual memory

### **Design and Analysis of Data Structures**

Binary trees, sorting, hash tables, priority queues, Big-O analysis

### **Introduction to Databases**

Relational algebra, data modeling, SQL

### **Operating Systems**

Concurrency & synchronization File systems, virtual memory

### **Design and Analysis of Data Structures**

Binary trees, sorting, hash tables, priority queues, Big-O analysis

### Database Internals e.g., (CSC443)

Storage, buffer pools, B-trees, transactions, write-ahead logging, query processing, etc.

### **Introduction to Databases**

Relational algebra, data modeling, SQL

### **Operating Systems**

Concurrency & synchronization File systems, virtual memory

### Design and Analysis of Data Structures

Binary trees, sorting, hash tables, priority queues, Big-O analysis

### Database Internals e.g., (CSC443)

Storage, buffer pools, B-trees, transactions, write-ahead logging, query processing, etc.

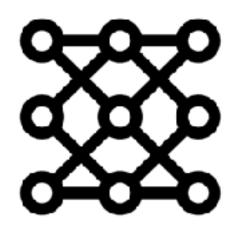
Solid programming skills in C, C++, Java, or at least Python

# **Background Knowledge**

CSC443 is important background

All lectures recorded

Will let you know what to catch up on





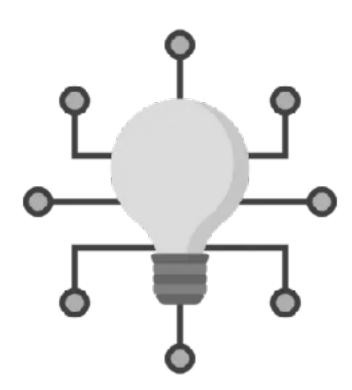


https://www.nivdayan.net/database-system-technology-csc443

# Seminar



Reading ≈20-30 Papers

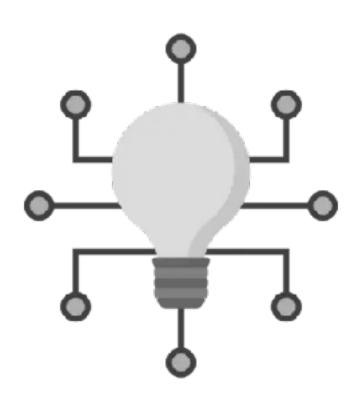


Small Research Project

# Seminar



Reading ≈20-30 Papers



Small Research Project

# Papers

topical

Classics

"Best papers"

influential (citations)









# Papers

topical

Classics

"Best papers"

influential (citations)









New skill & insight

# Why read papers?

Reading papers is a skill

Get research ideas

Employ the state of the art

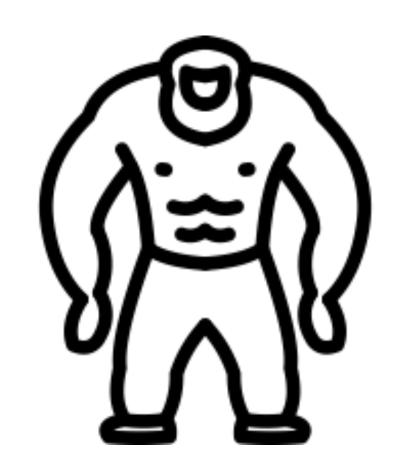






### Focus

DB is a huge field



Data structures for modern hardware

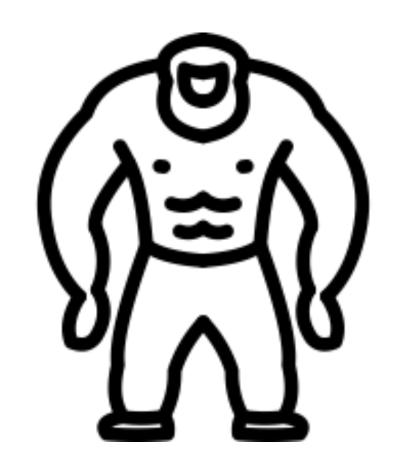


For a specialized course, I may as well teach my specialty



# Focus

DB is a huge field



Data structures for modern hardware

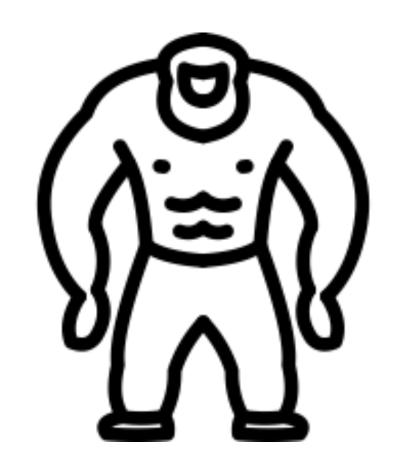


For a specialized course, I may as well teach my specialty



# Focus

DB is a huge field



Data structures for modern hardware



For a specialized course, I may as well teach my specialty



# Website

### CSC2525

Research Topics in Database

Management

**Instructor:** Niv Dayan

**Lectures:** Wednesday 13:00-15:00 (BA1200)

Office Hours: after each class



https://www.nivdayan.net/research-topics-in-database-management-csc2525



# First 5 ish are lectures by me



First 5 ish are lectures by me



Subsequent 7 ish

1-1.5 hour lecture



Student Presentation



First 5 ish are lectures by me



Subsequent 7 ish

1-1.5 hour lecture

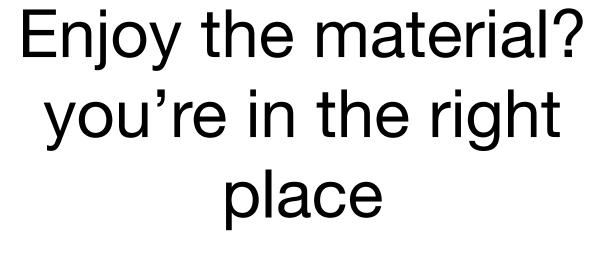
Student Presentation

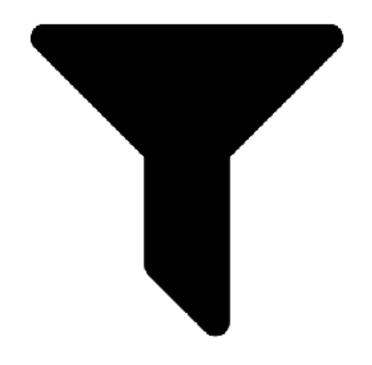


More on presentations later

# Use the first two lectures wisely

Dynamic Arrays & Filter Data
Structures



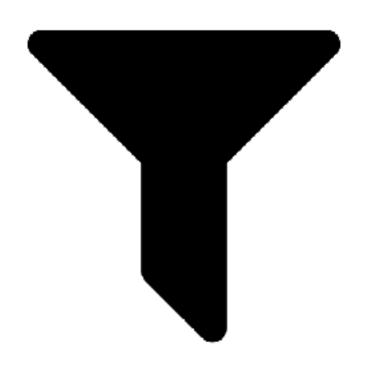




# Use the first two lectures wisely

Dynamic Arrays & Filter Data
Structures

Enjoy the material? you're in the right place





# Participation

You are required to attend each class

Read papers in advance

Participate in class discussions



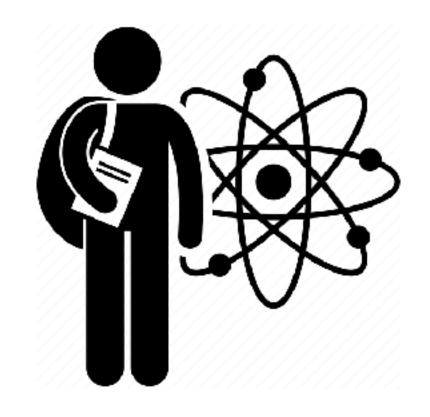




# **Project**

Implement & evaluate

Proposals due by mid-Feb

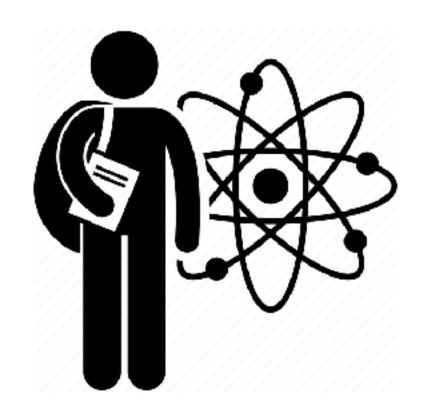




# Project

Implement & evaluate

Proposals due by mid-Feb





More on this later

# **Oral Exam**

20-30 min per student

Papers & Material





# **Grade Components**

- (1) Project Report & Code
- (2) Paper presentation
- (3) Oral exam

# **Grade Components**

- (1) Project Report & Code
- (2) Paper presentation
- (3) Oral exam

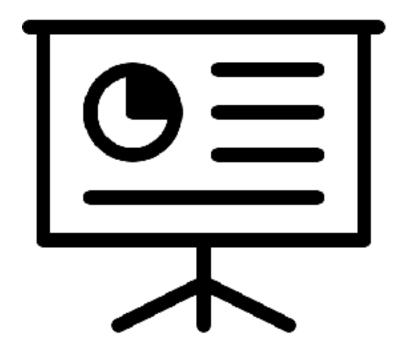
Precise breakdown to be announced later

# Office Hours

Right after class

You're encouraged to reach out before your presentation







Post questions for everyone's benefit!



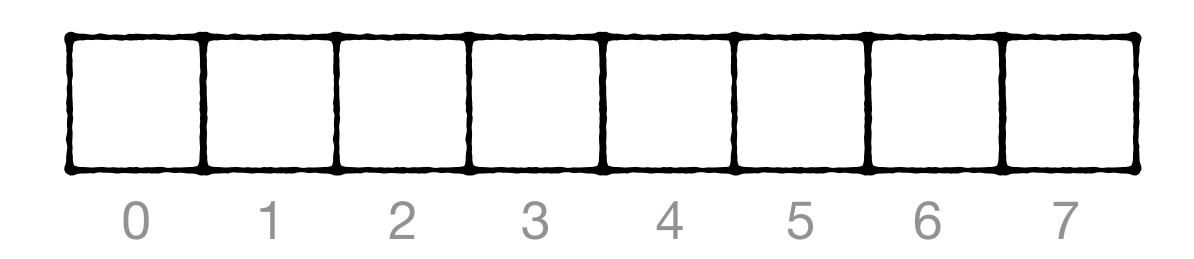
We'll record classes, but you must still attend.

### And now to our first lecture

# Dynamic Arrays CSC2525 Research Topics in Database Management

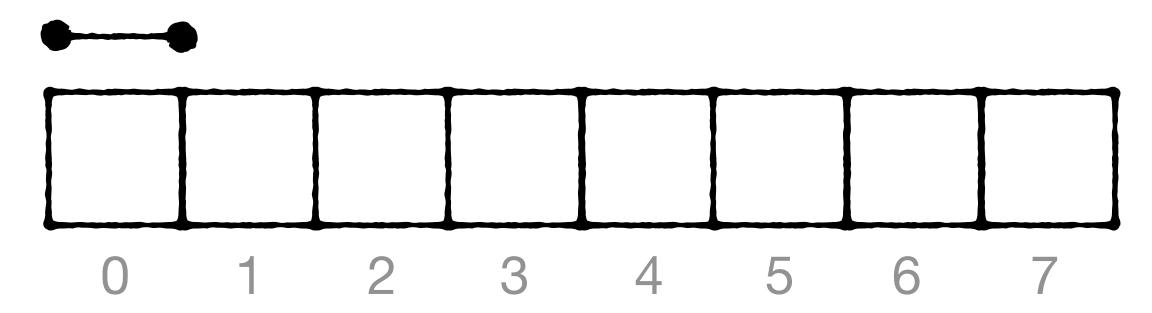
Niv Dayan

# Arrays



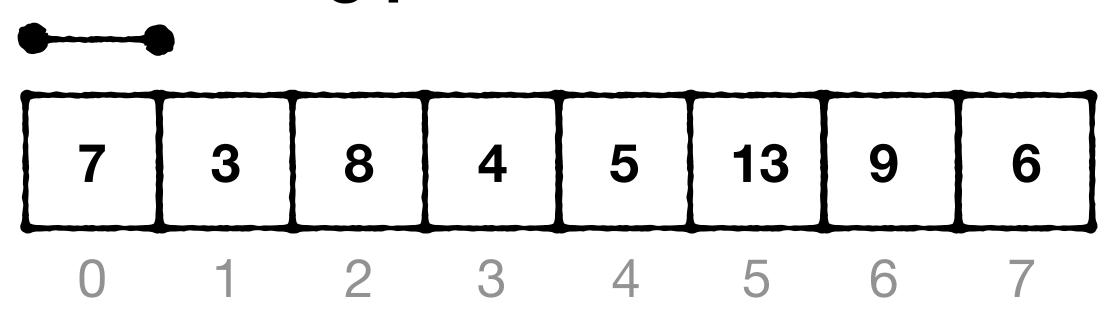
# Arrays

#### Fixed width slots



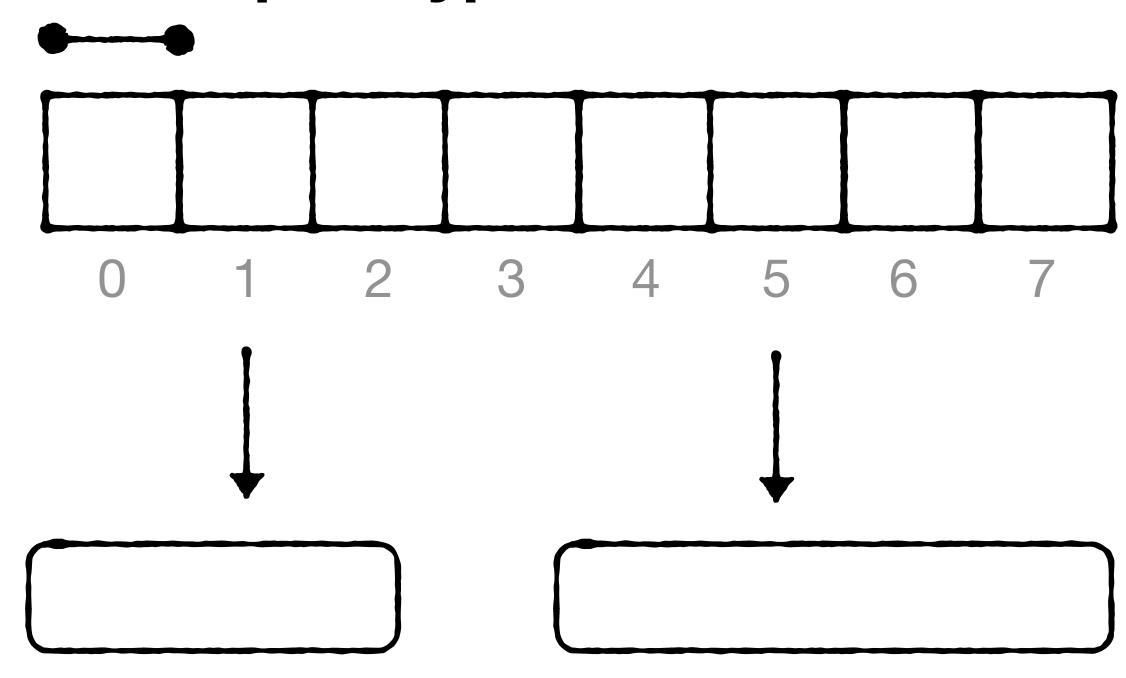
Fixed width slots

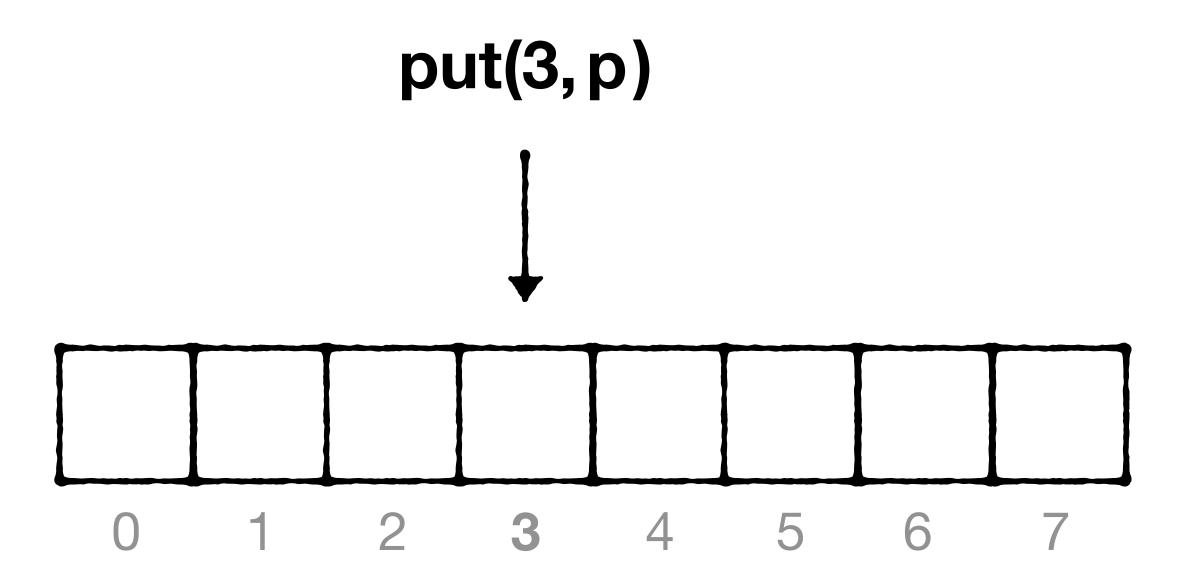
#### e.g., integers or floating points



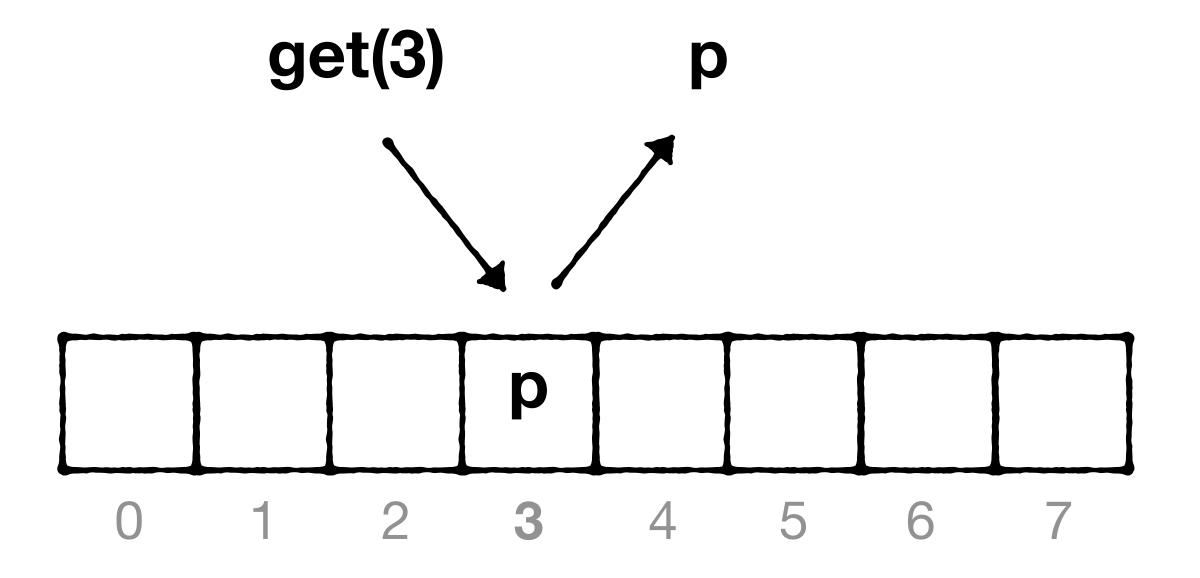
Fixed width slots

### Or pointers to complex types

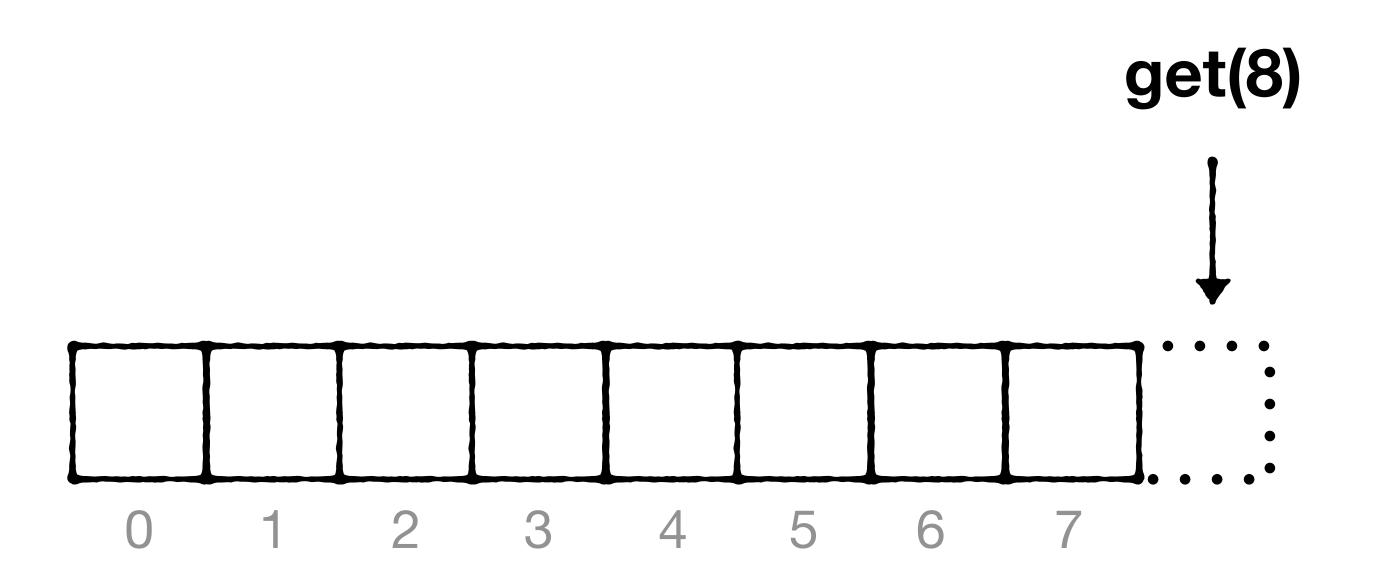




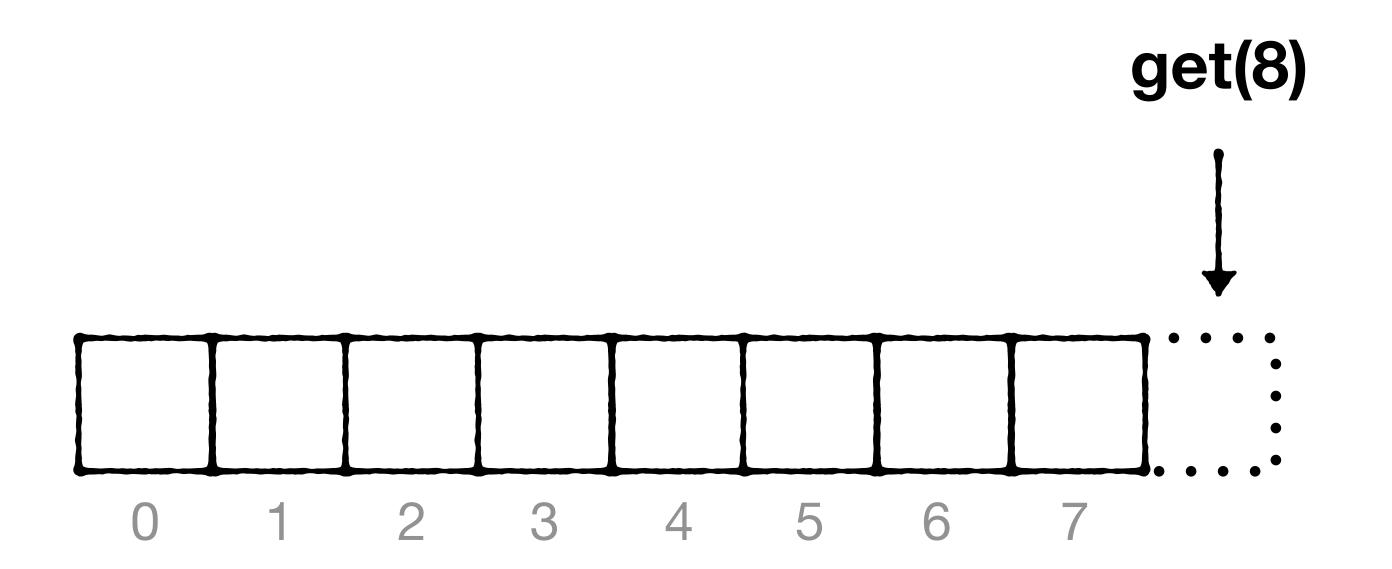
Supports random access



Supports random access



Overflow error (e.g., java)



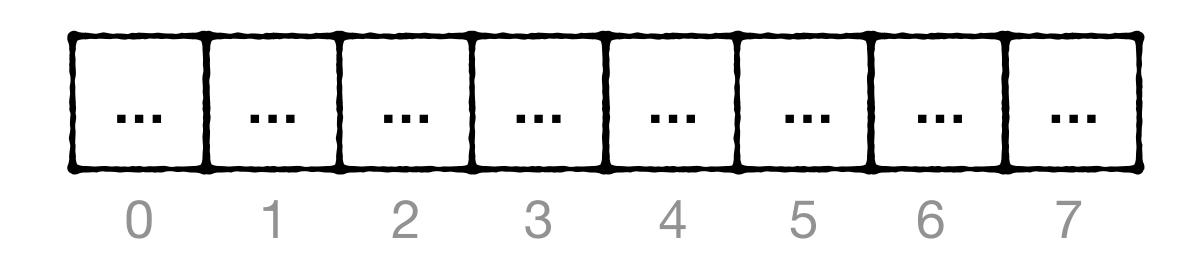
Overflow error (e.g., java)

Undefined behavior (e.g., C++)

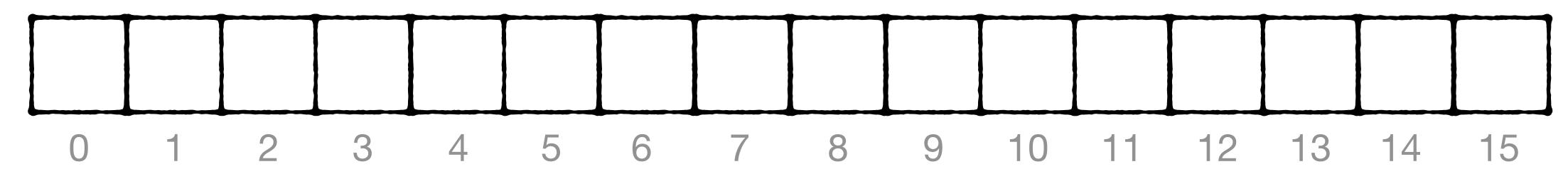
# How to keep inserting when out of space?

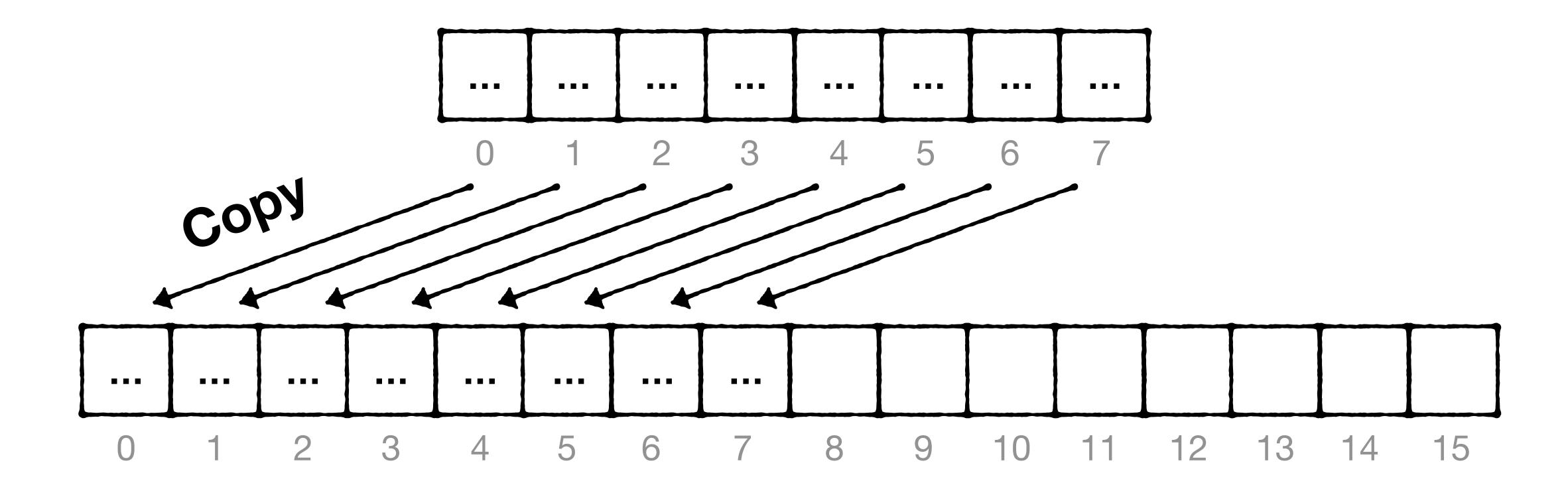


# How to keep inserting when out of space?

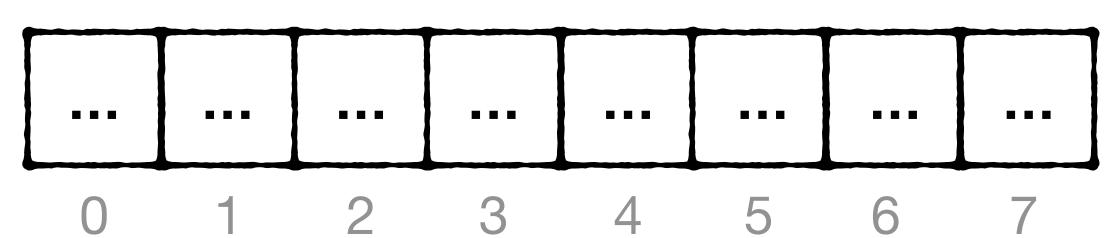




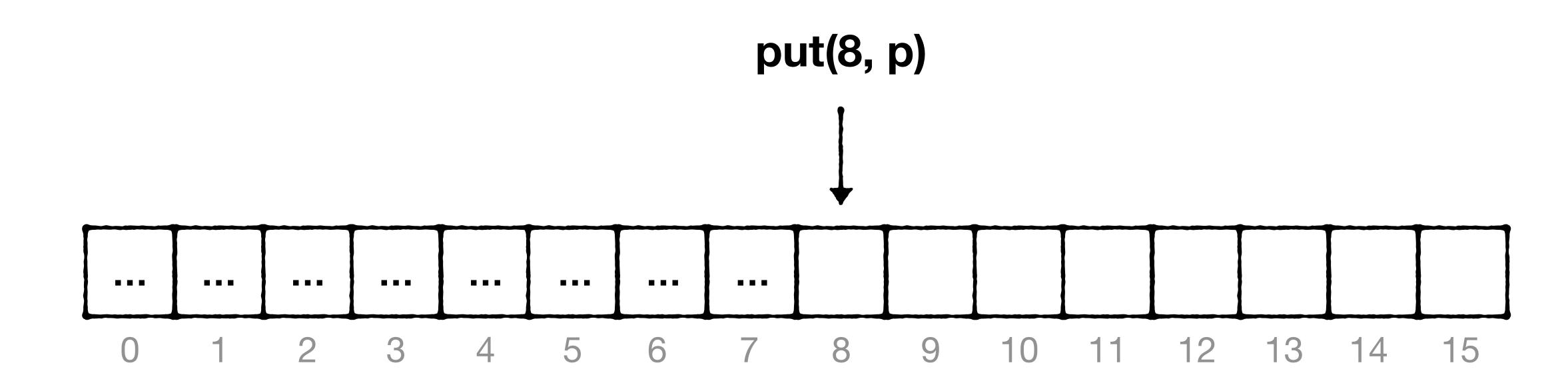




#### Deallocate







Static int nums  $[4] = \{0, 0, 0, 0\};$ 

Dynamic std::vector<int> vec;

```
std::vector<int> vec;
for (int i = 0; i < 10; ++i) {
   vec.push_back(i);
   std::cout << "Added " << i << ", size: " << vec.size()
        << ", capacity: " << vec.capacity() << std:end;
}</pre>
```

Element	Size	Capacity
0	1	1
1	2	2
2	3	4
3	4	4
4	5	8
5	6	8
6	7	8
7	8	8
8	9	16
9	10	16

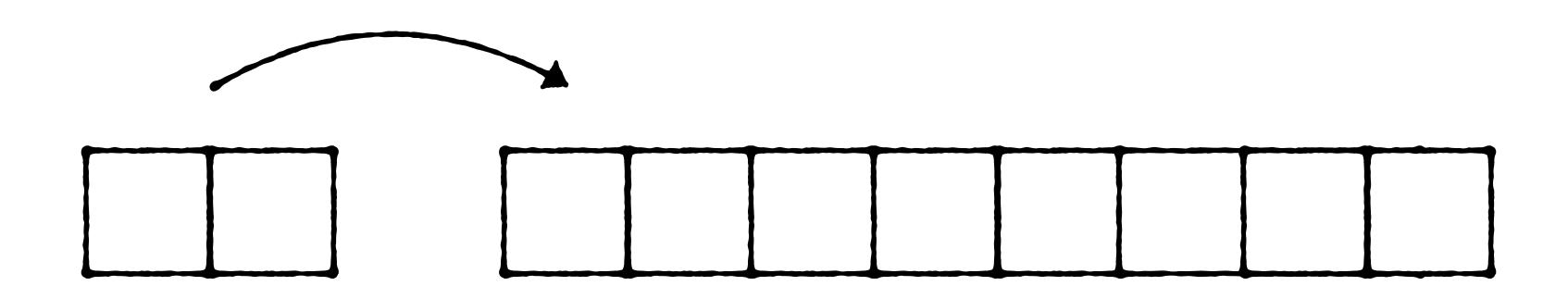
Element	Size	Capacity
0	1	1
1	2	2
2	3	4
3	4	4
4	5	8
5	6	8
6	7	8
7	8	8
8	9	16
9	10	16

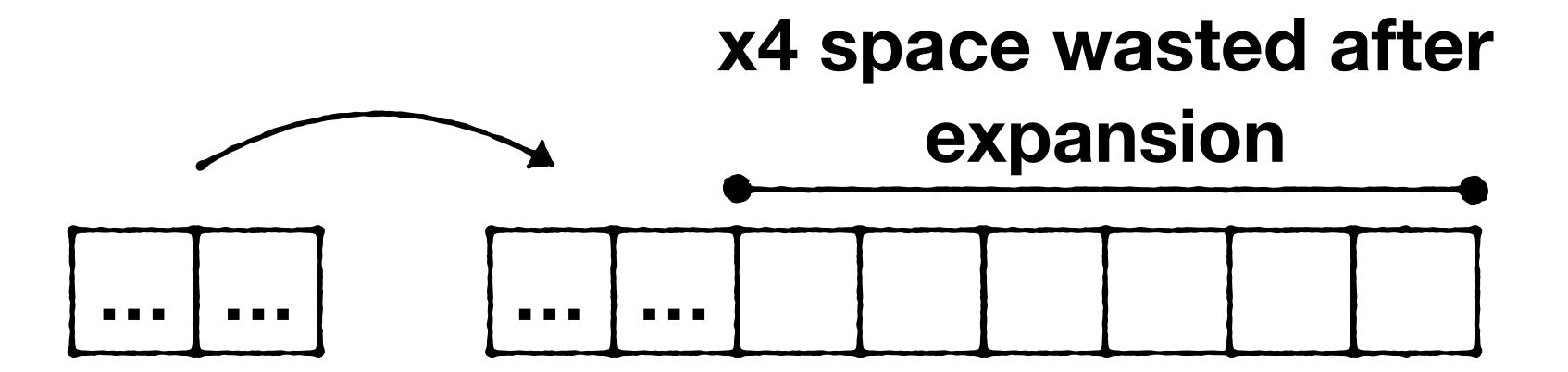
#### **Uses Growth Factor 2**

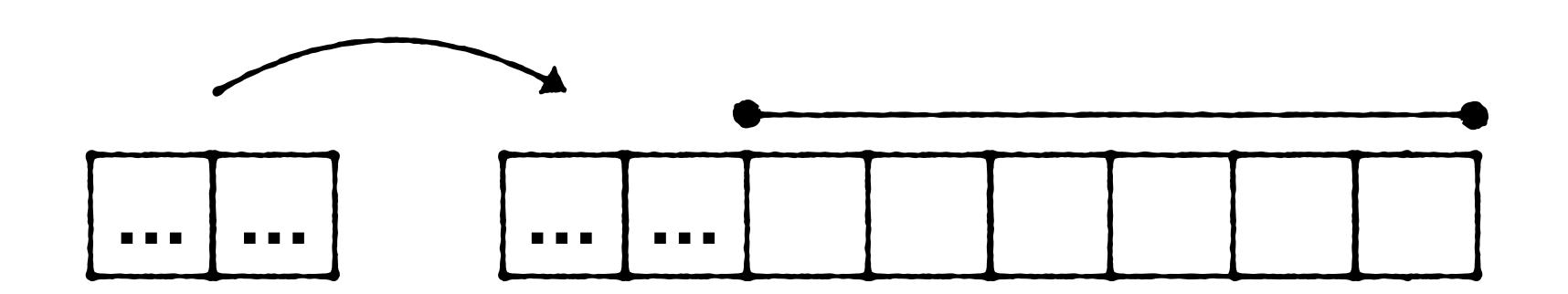
#### What if Growth Factor G is

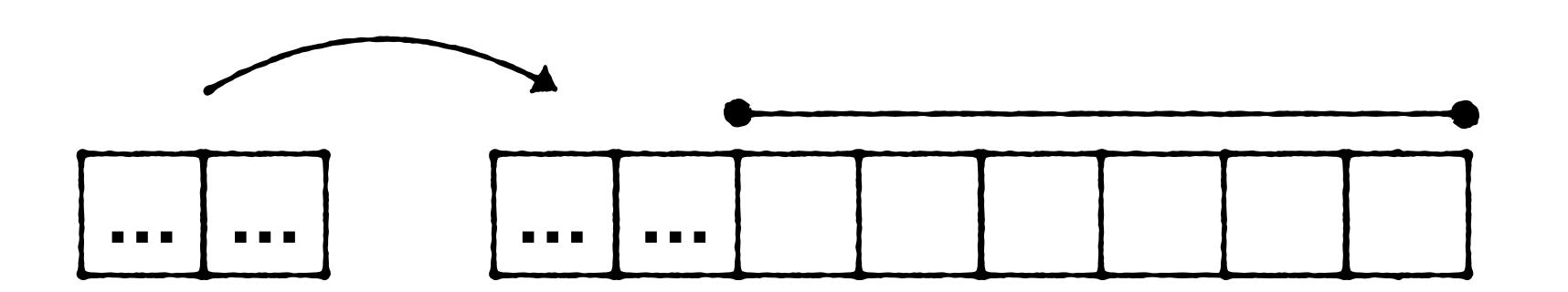
too high?

too low?



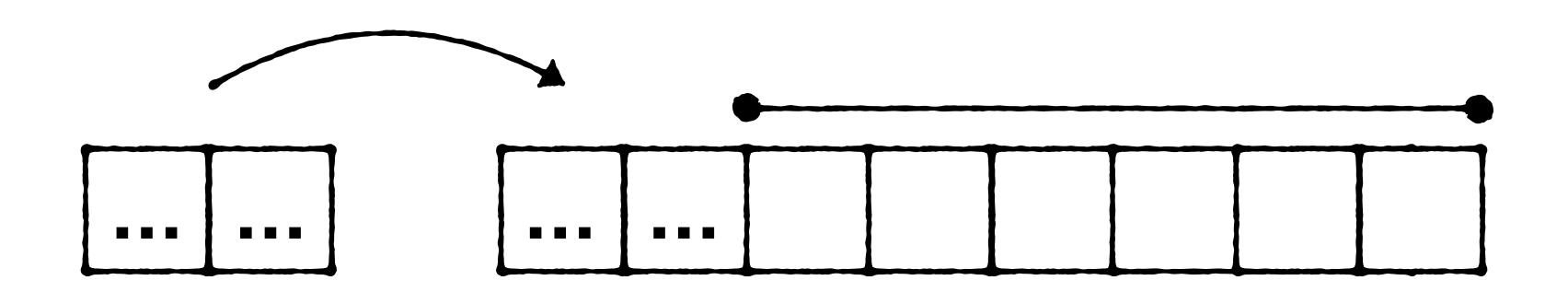






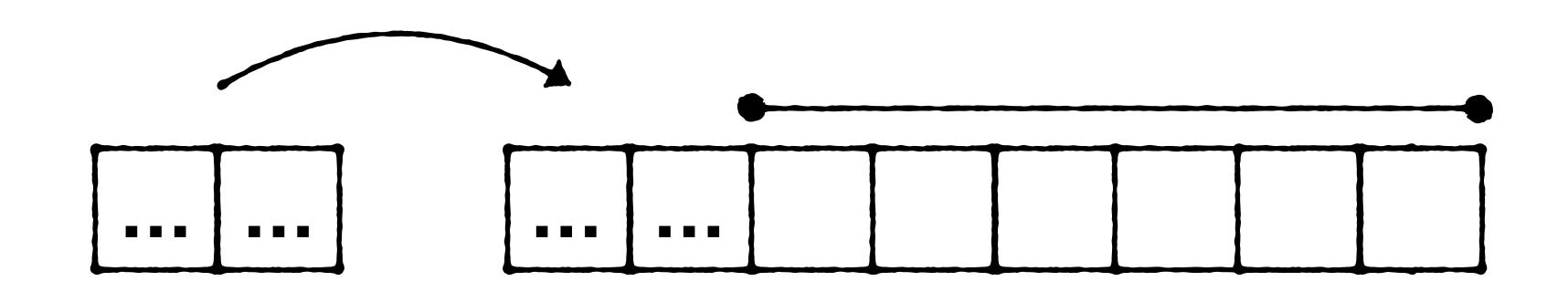
Space-amplification = G

Right after expansion



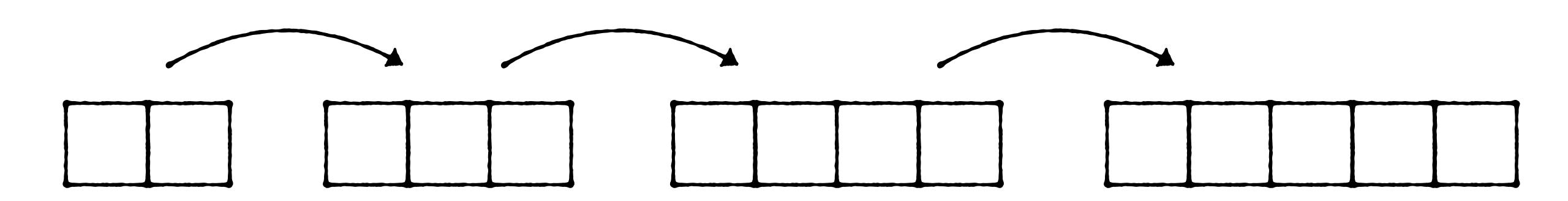
$$= G + 1$$

During expansion



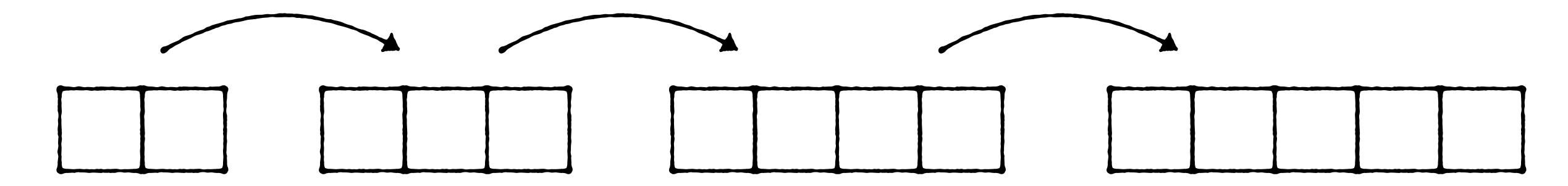
What if Growth Factor G is **too low e.g., 1.2** 

# What if Growth Factor G is too low e.g., 1.2

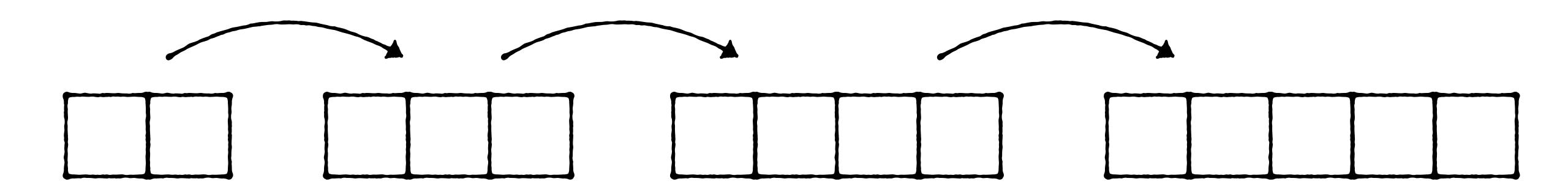


# What if Growth Factor G is too low e.g., 1.2

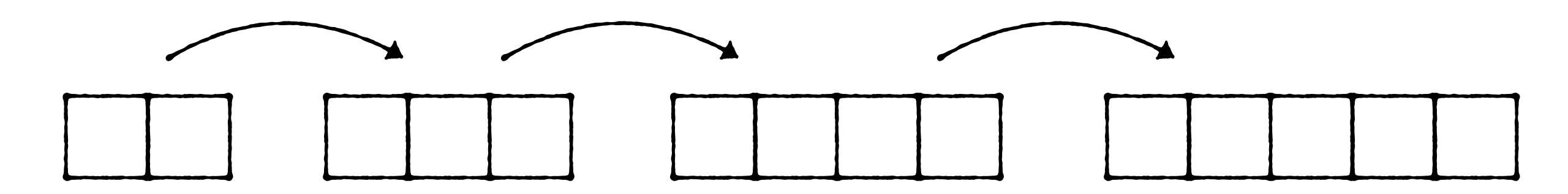
#### Insertion overheads increase



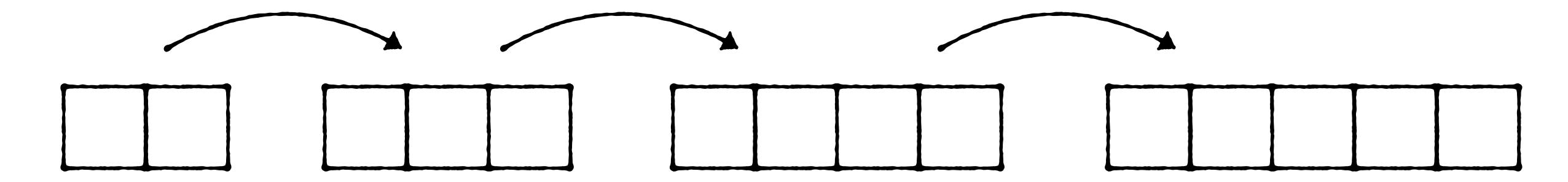
#### What if Growth Factor G is too low



#### What if Growth Factor G is too low



#### What if Growth Factor G is too low



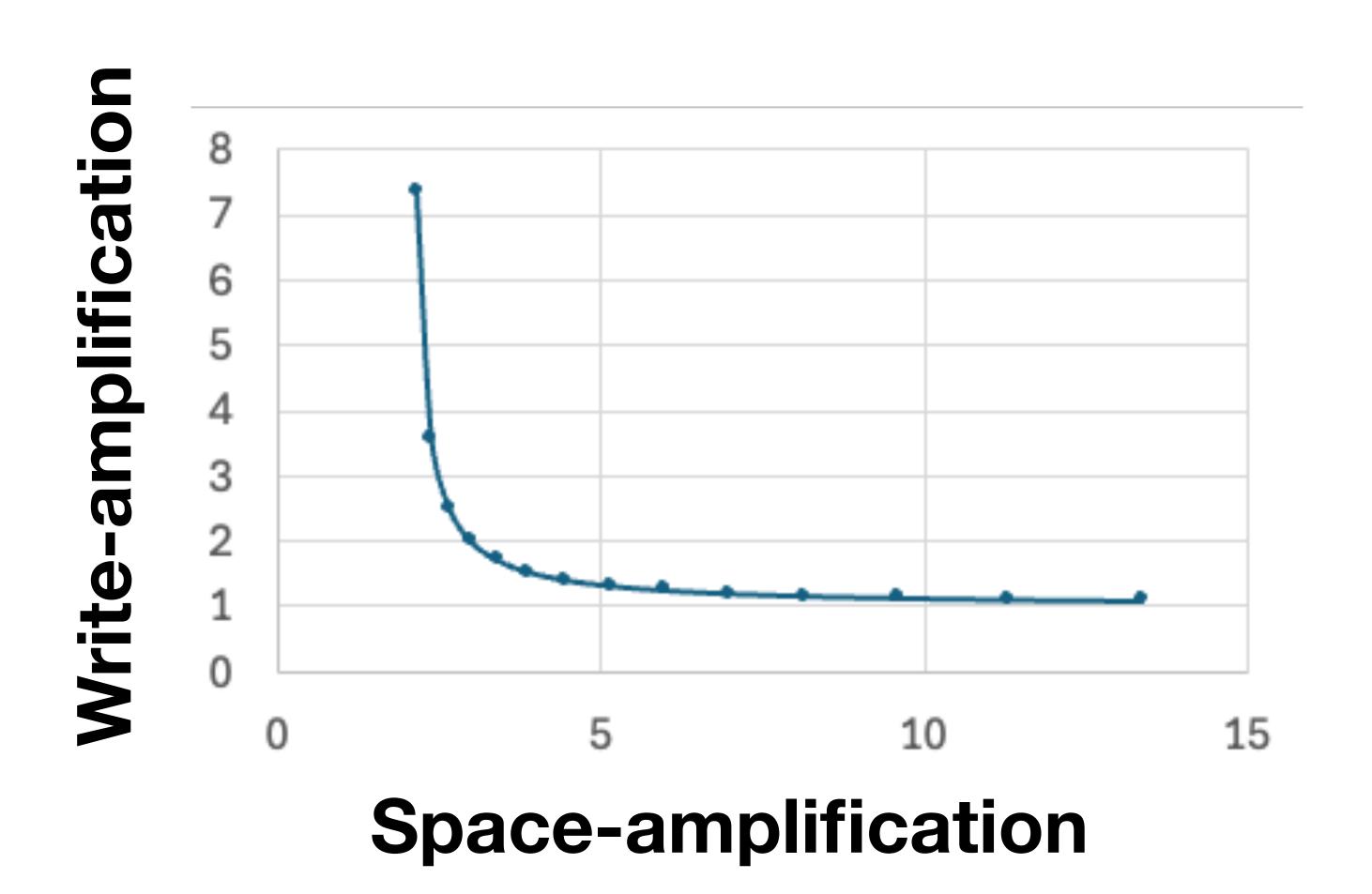
# Growth factor G impact

Space-amplification

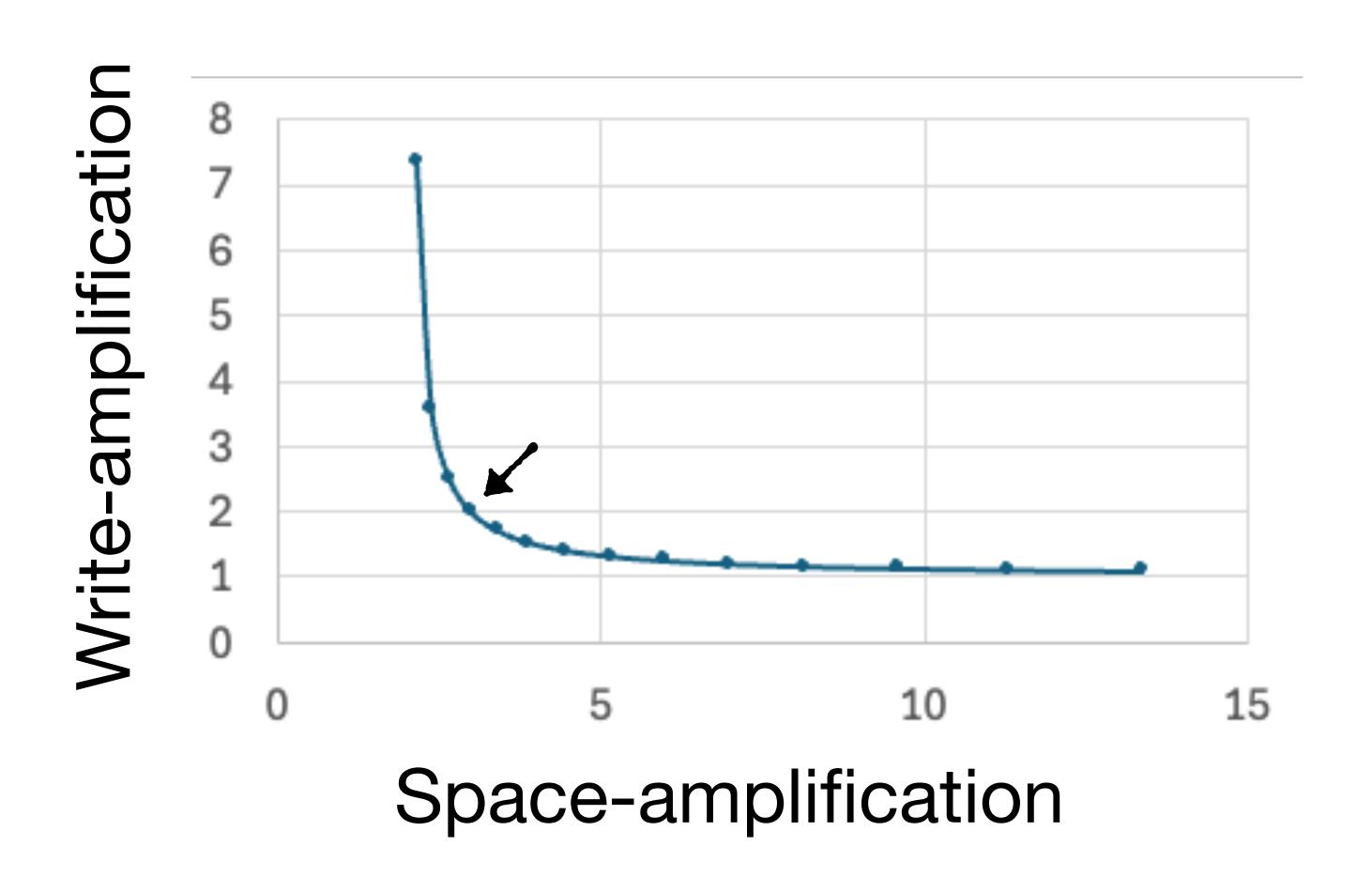
G+1

Write-amplification

# Growth factor G impact

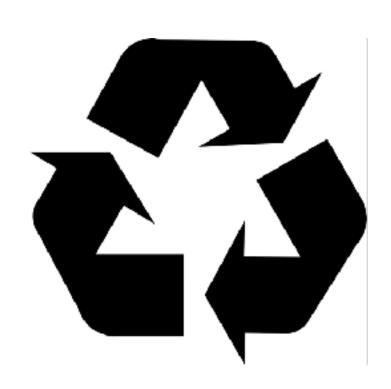


#### Growth factor 2 achieves a good balance

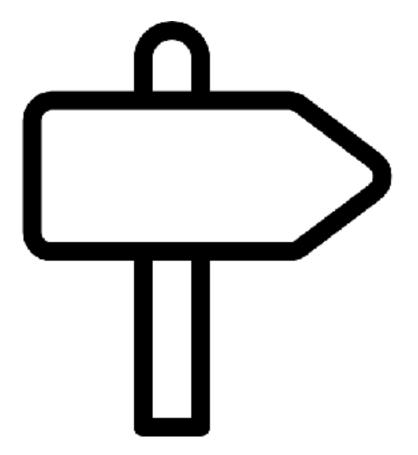


#### And now to new stuff

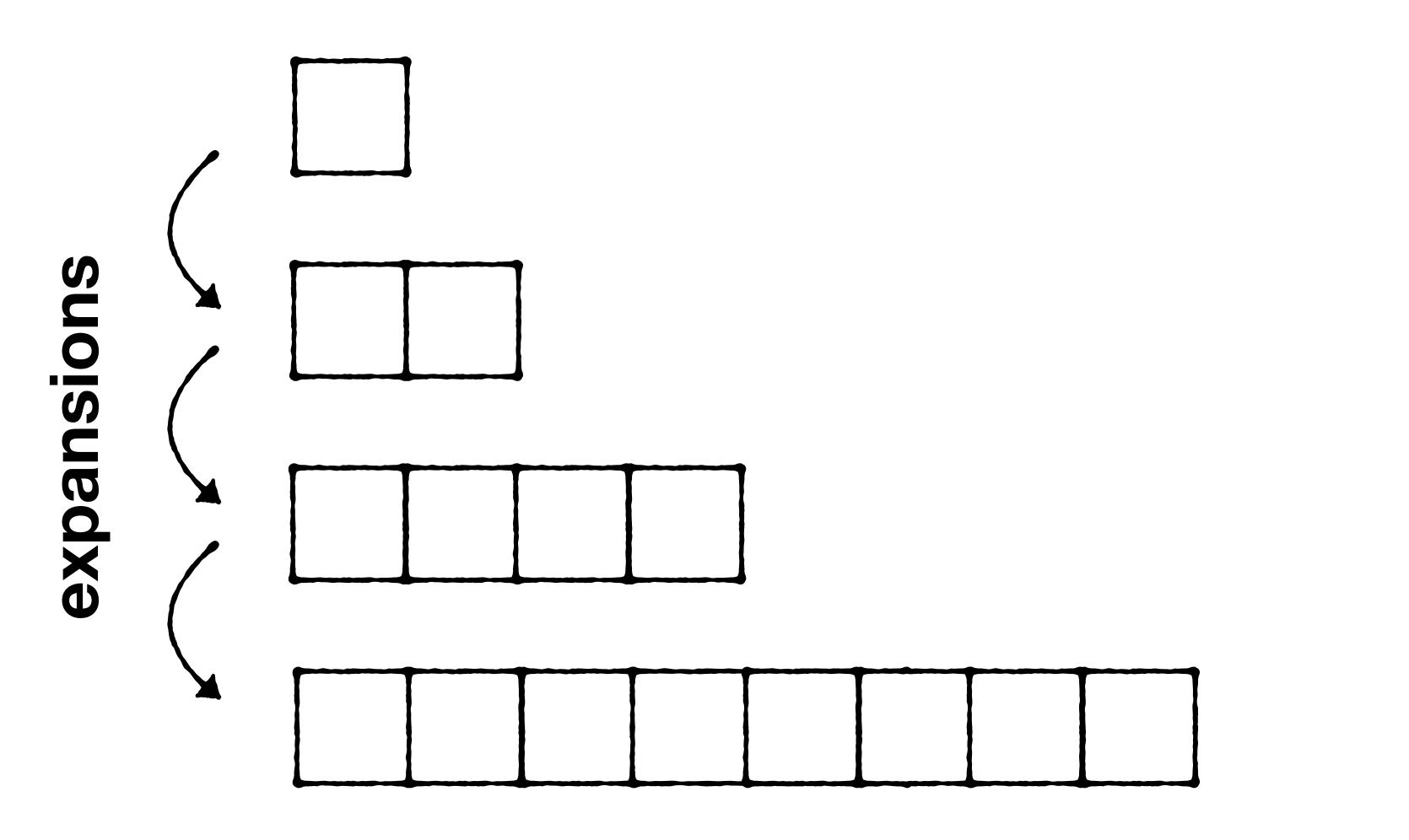
# Reusing Deallocated Space



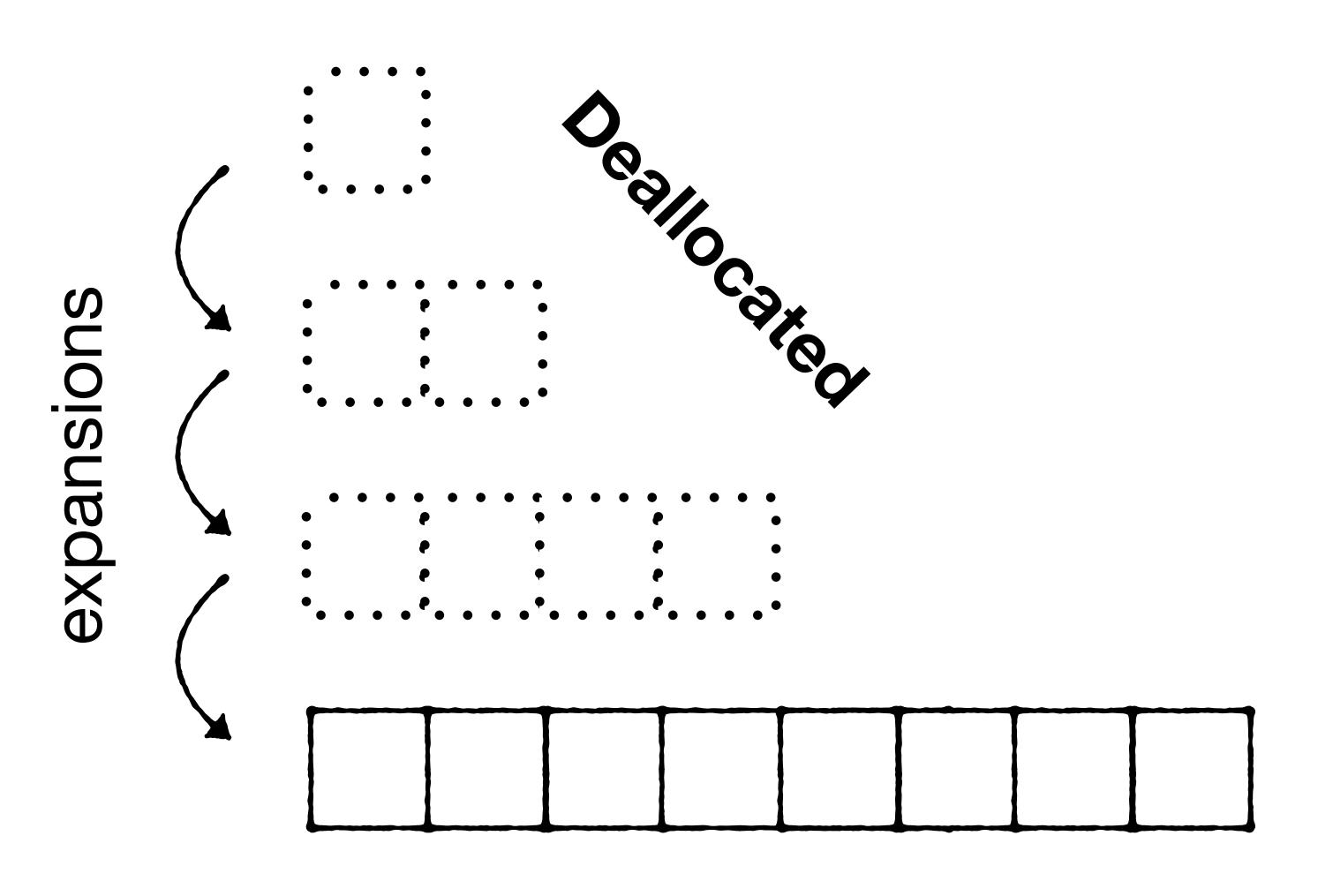
# Alleviating tradeoff via indirection



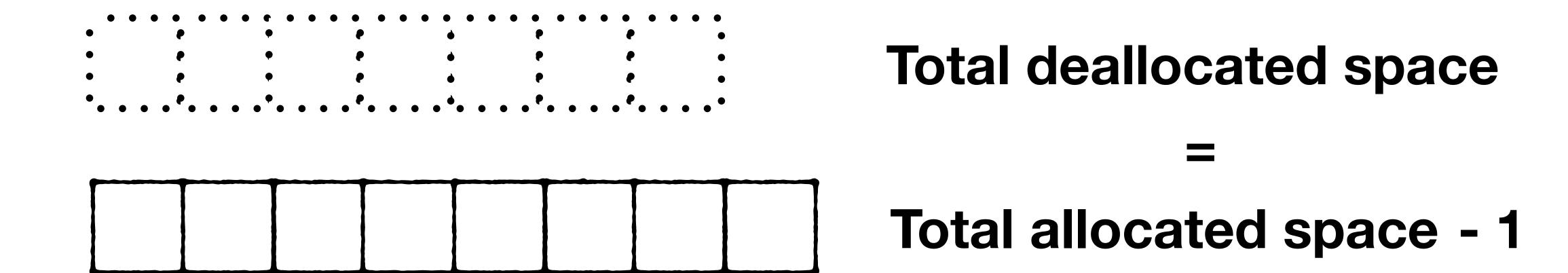
# Reusing Deallocated Space Assume G ≥ 2



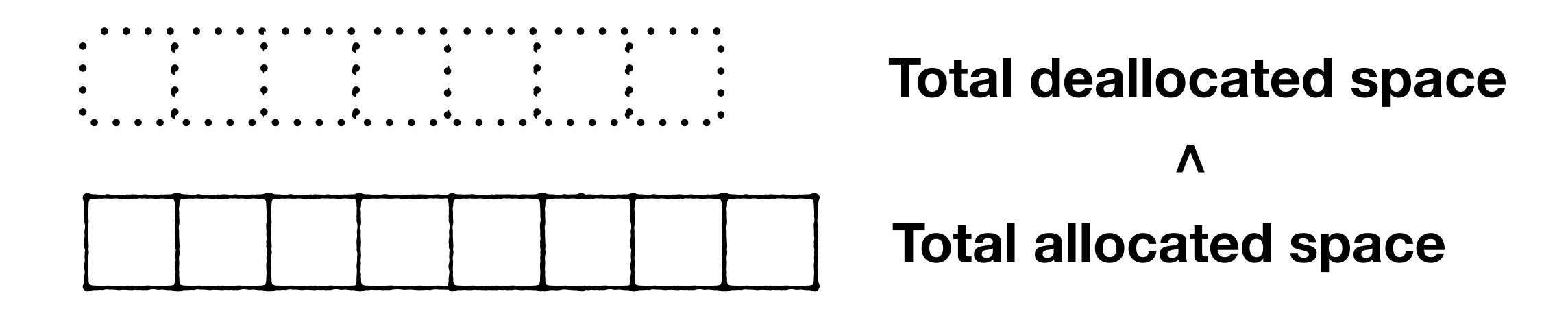
### Reusing Deallocated Space (G ≥ 2)



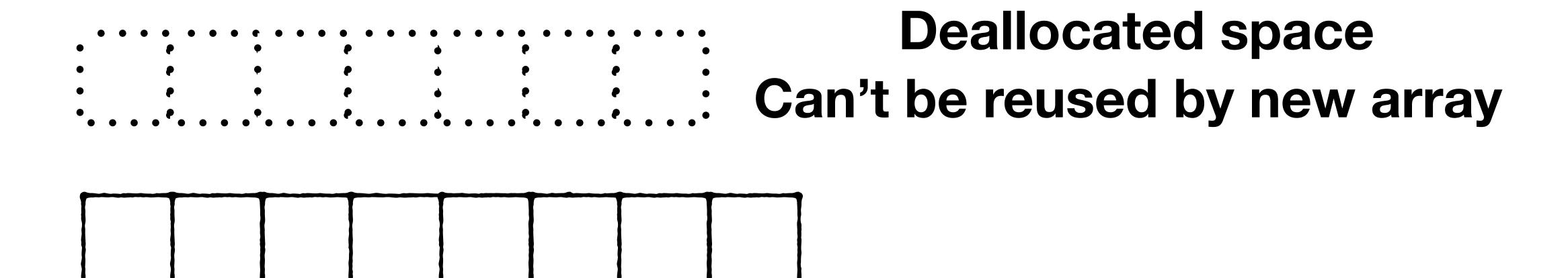
#### Reusing Deallocated Space



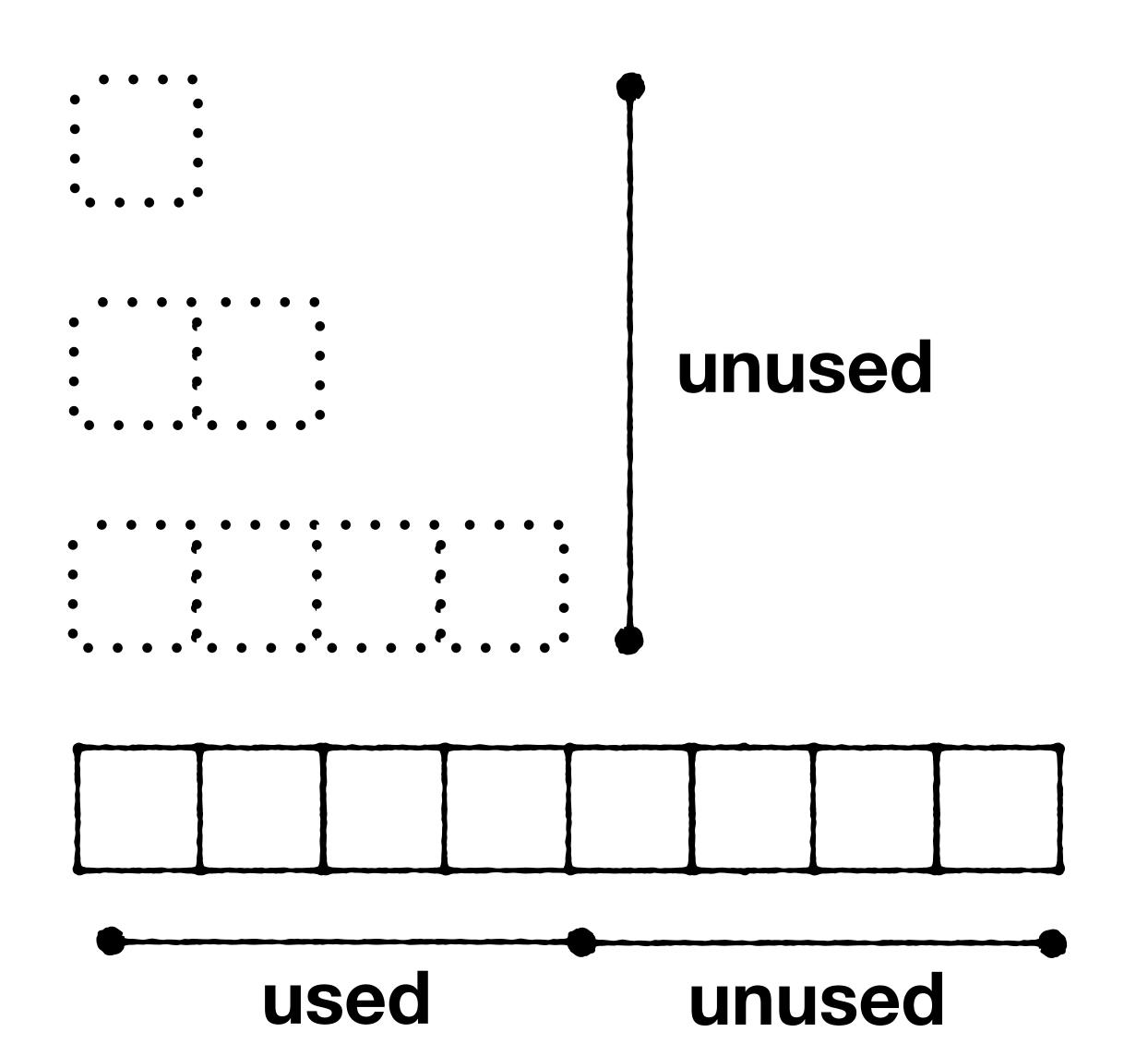
#### Reusing Deallocated Space



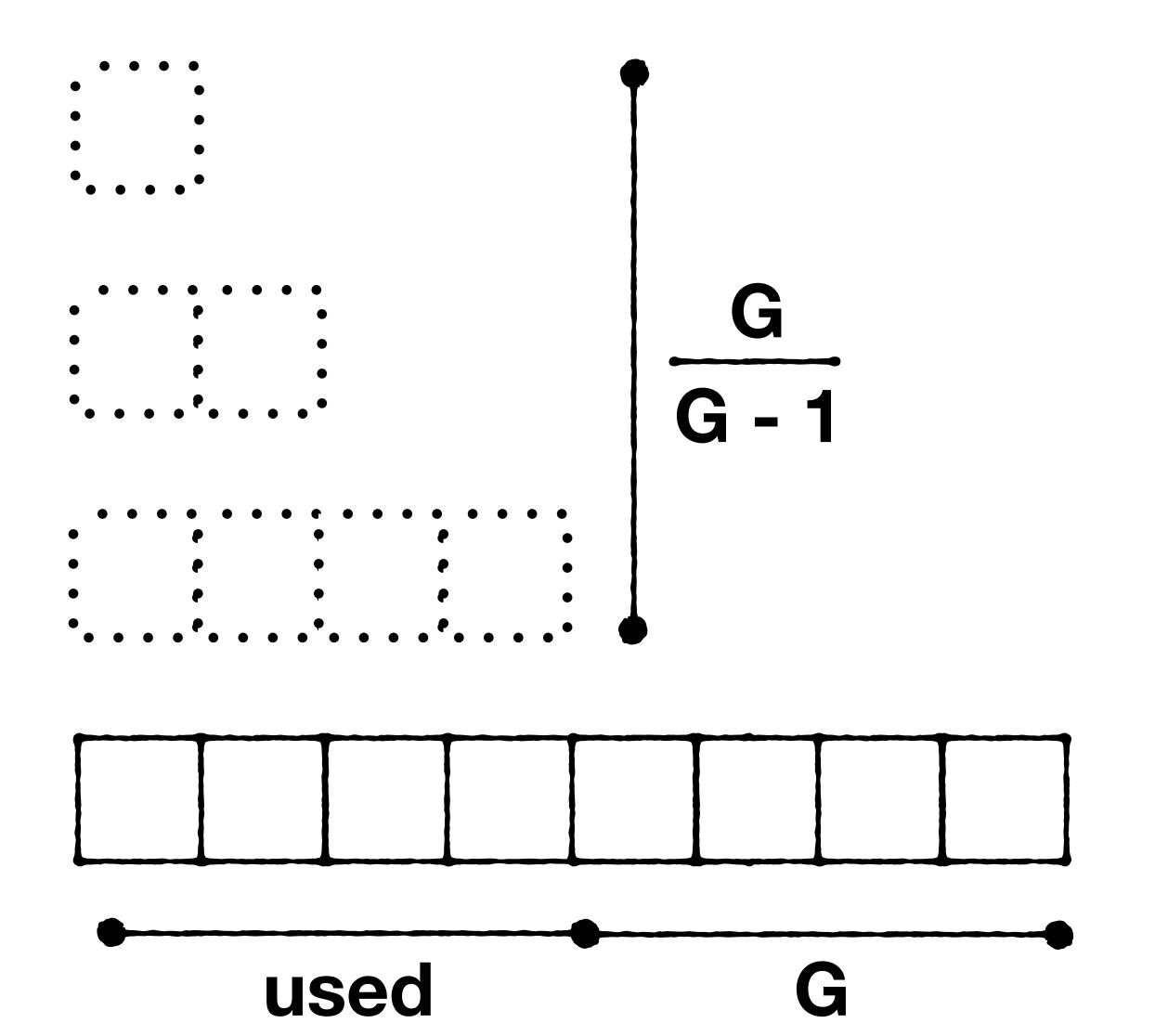
#### Reusing Deallocated Space



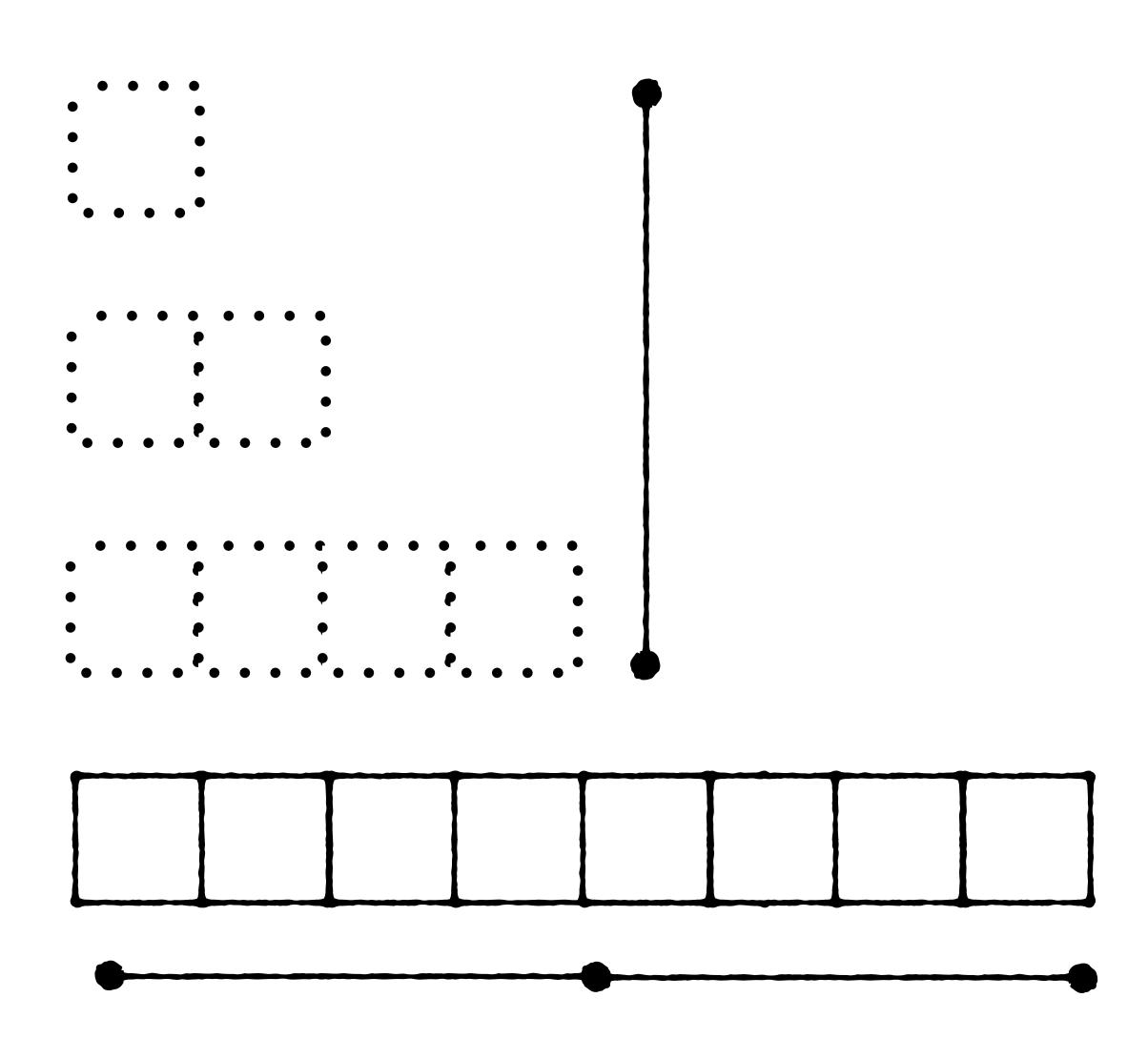
#### Deriving Max Space-Amp



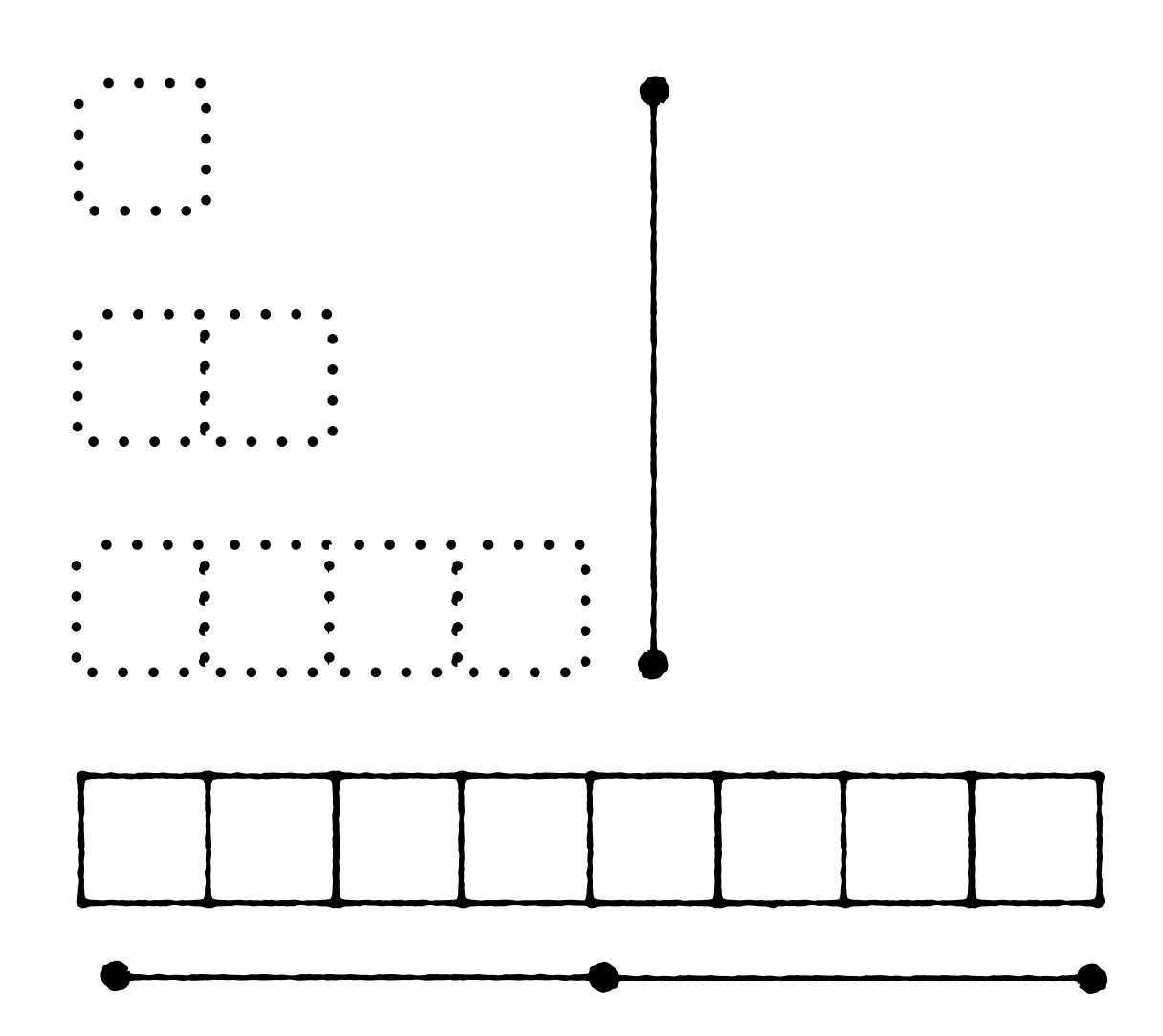
## Deriving Max Space-Amp



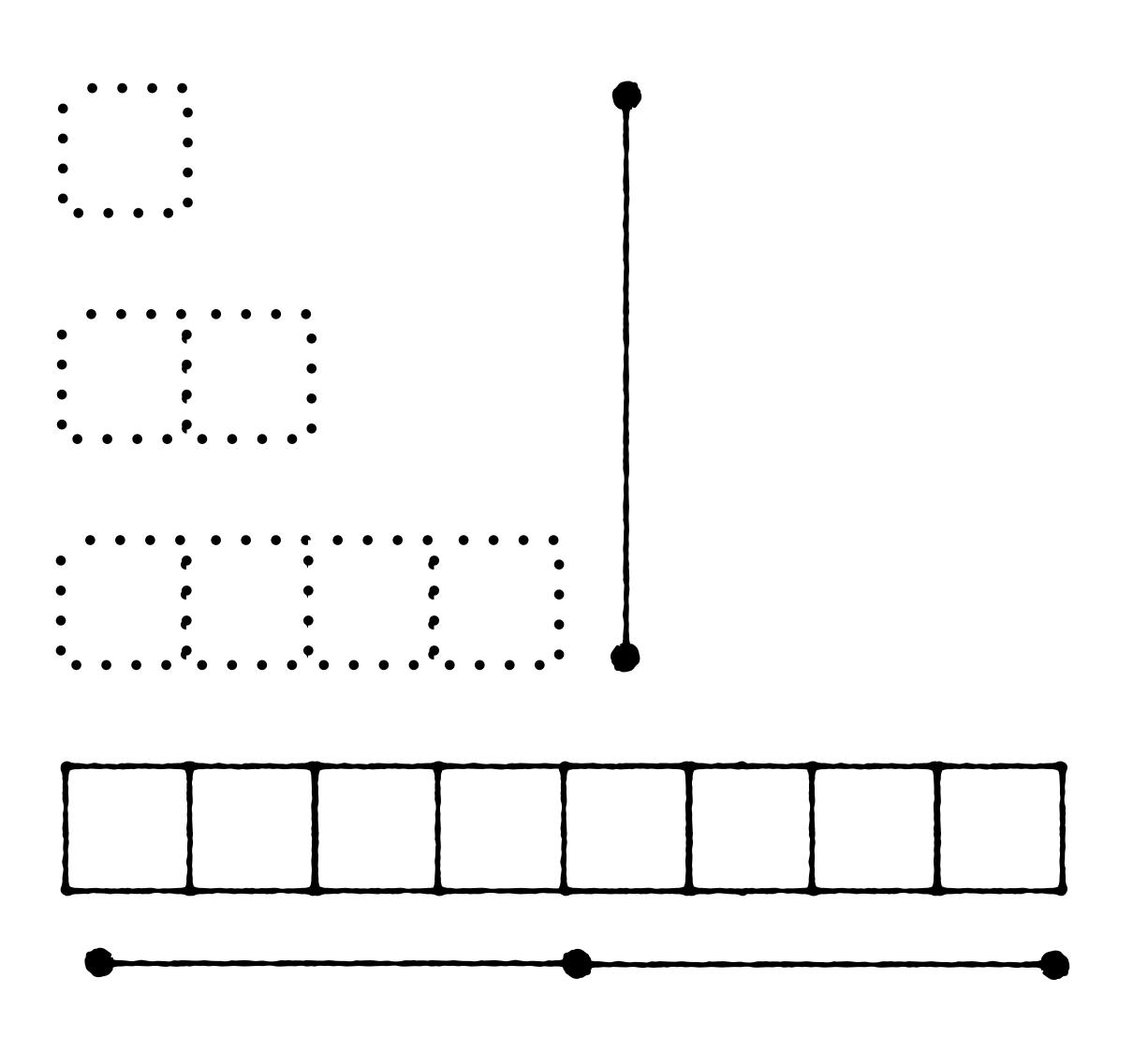
Max Space-Amp = 
$$\frac{\text{used} + \text{unused}}{\text{used}} = \frac{G}{G-1} + G$$



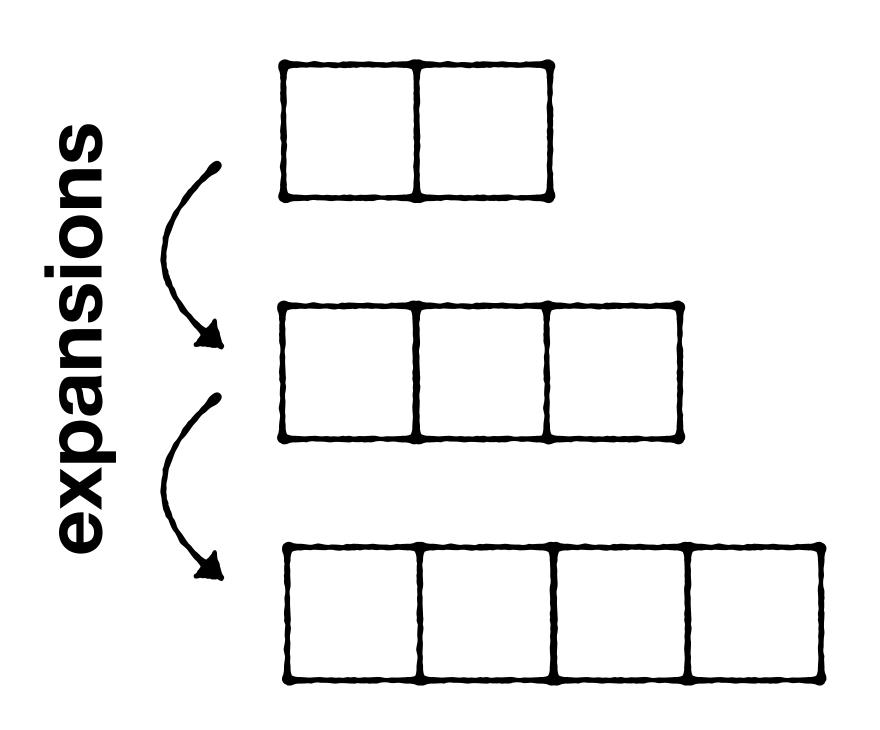
Max Space-Amp = 
$$\frac{\text{used} + \text{unused}}{\text{used}}$$
 =  $\frac{G^{-}}{G - 1}$ 

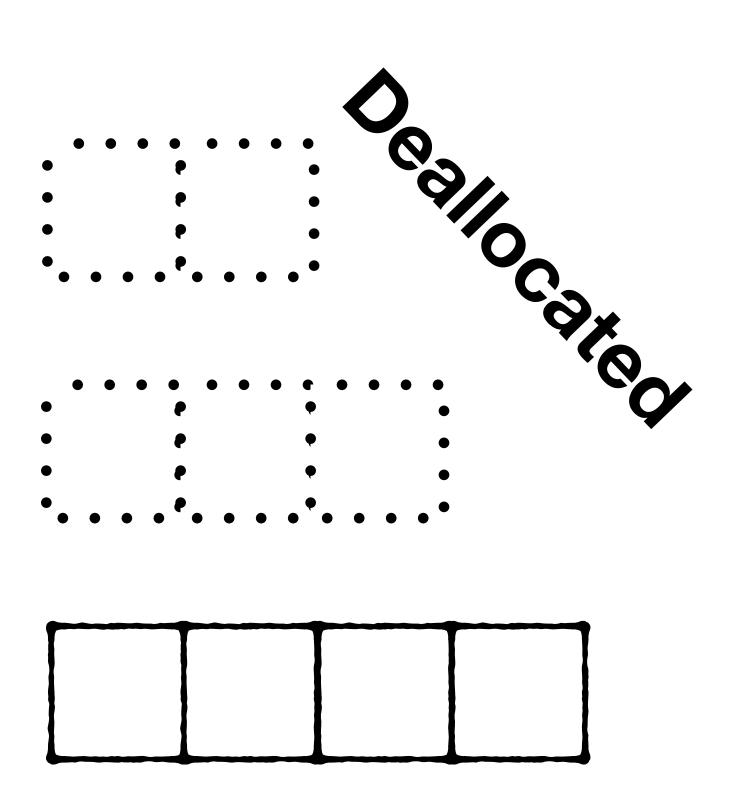


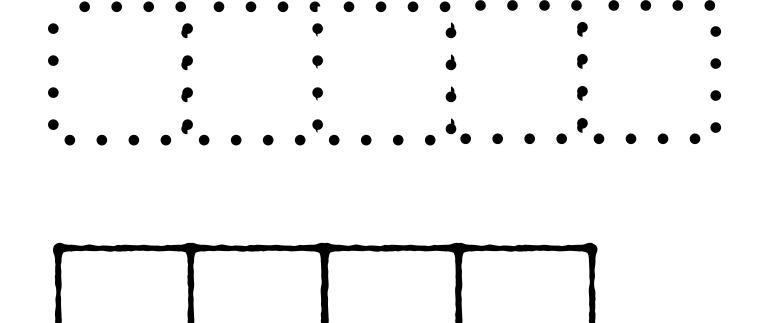
Max Space-Amp = 
$$\frac{\text{used} + \text{unused}}{\text{used}} = \frac{G^2}{G - 1} = 4$$
 for  $G = 2$ 







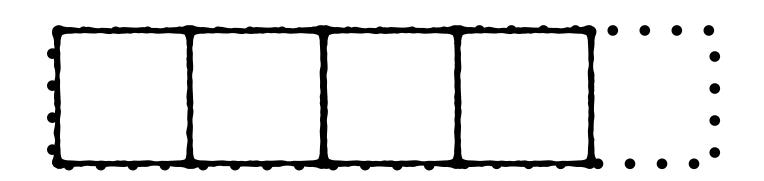




Total deallocated space

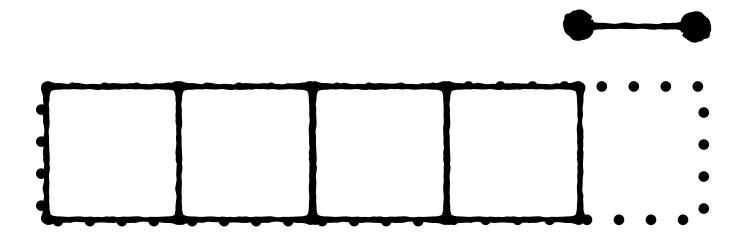
V

Total allocated space

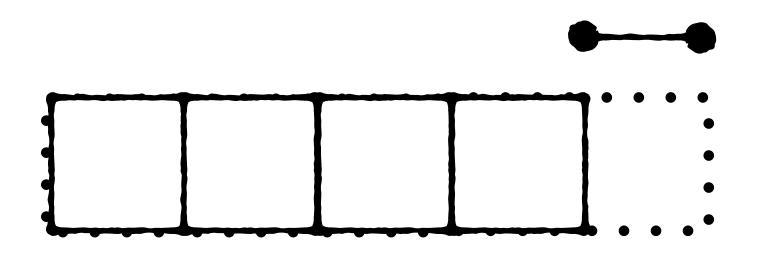


Reuse is possible

# Wasted space

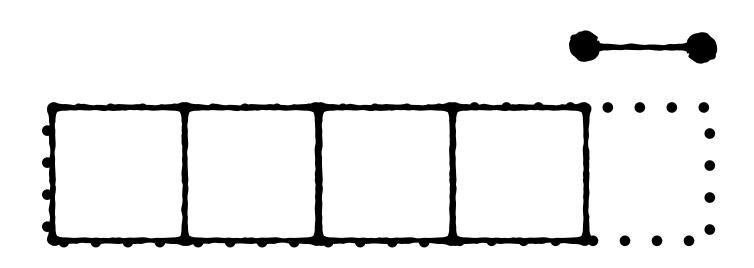


Wasted space



Could have expanded by larger factor

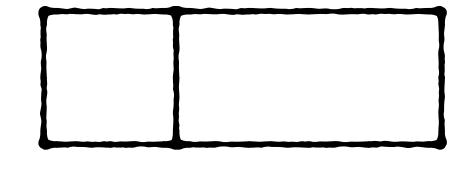
### For which growth factor, do we perfectly reuse the space?

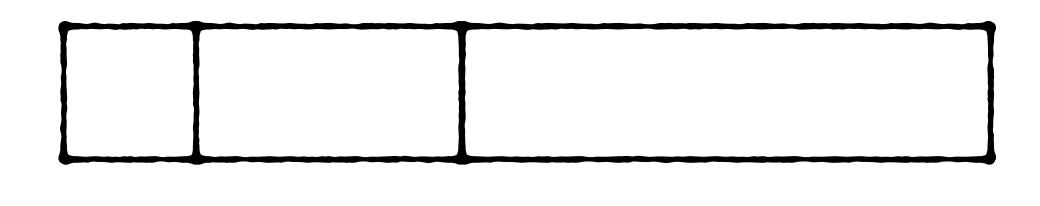


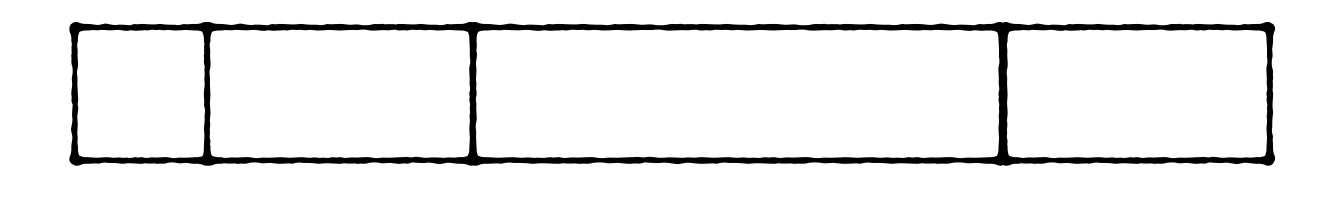


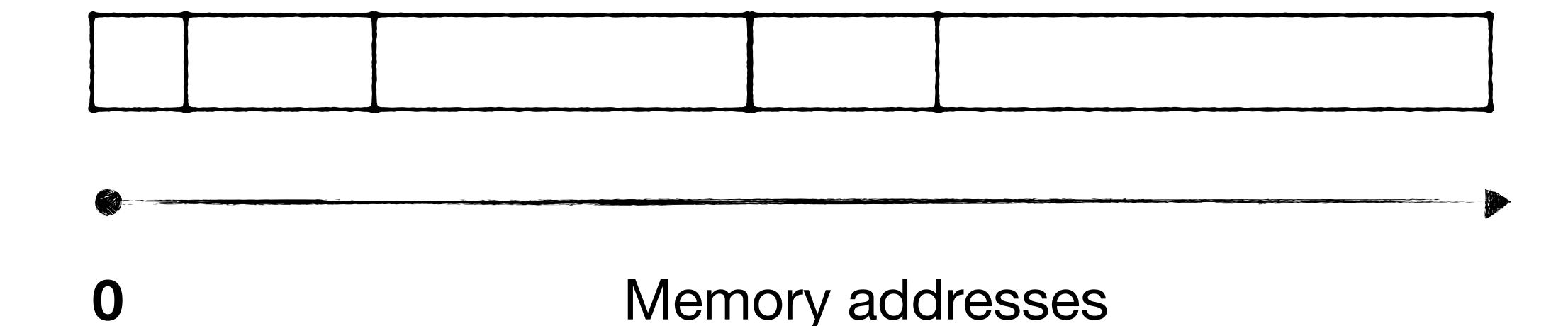
#### Assumptions on memory allocator



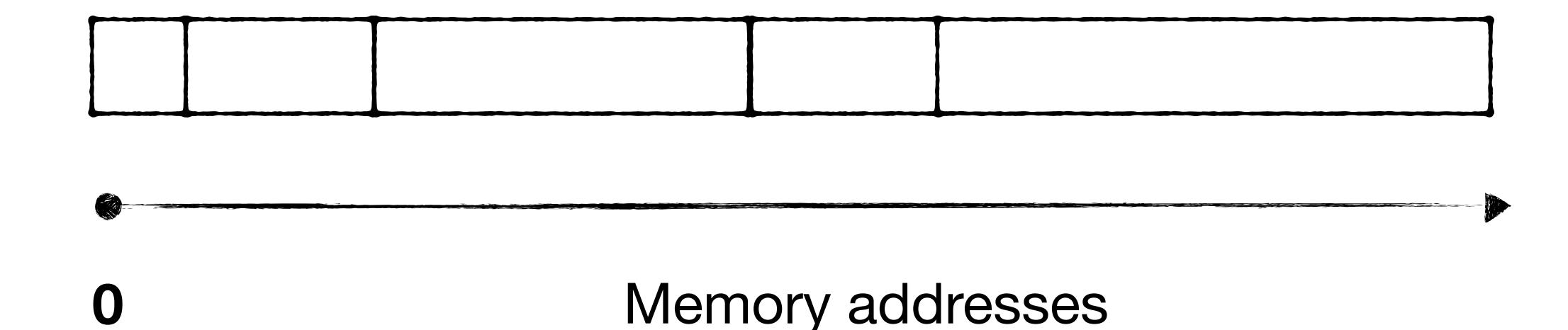








Reuse when possible

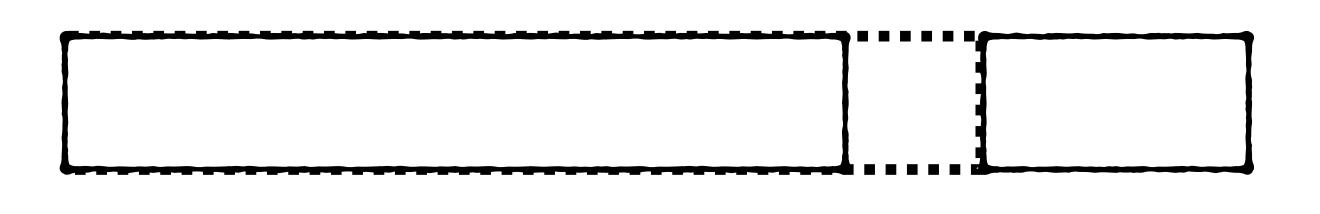


Reuse when possible

Deallocated

Memory addresses

Reuse when possible



Reuse when possible

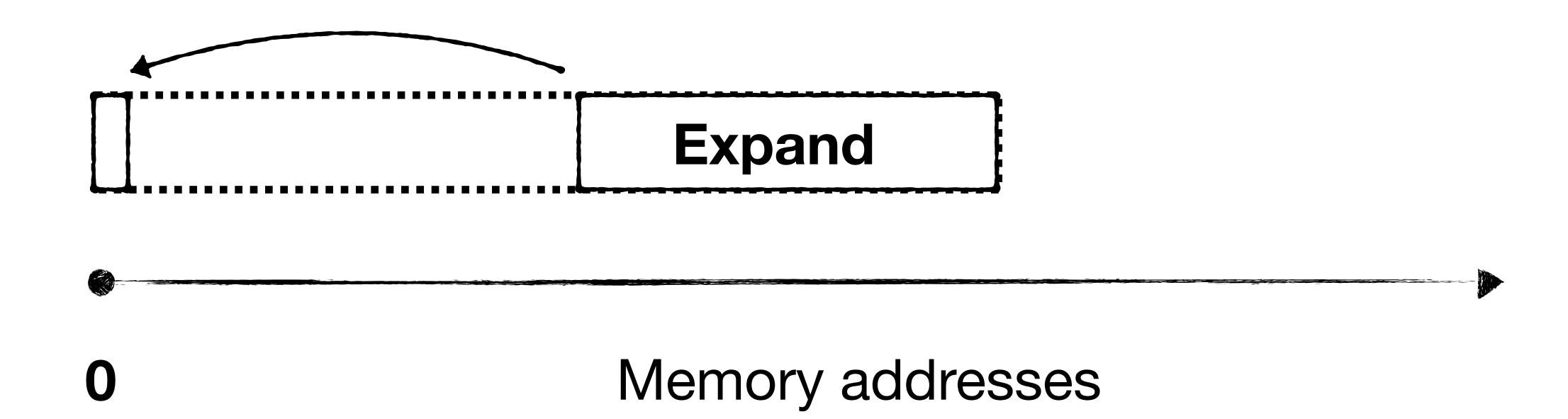
Deallocate as we copy (simplistic)

Expand

Memory addresses

Reuse when possible

Deallocate as we copy (simplistic)



Reuse when possible

Deallocate as we copy (simplistic)

copied →

Memory addresses

Reuse when possible

Deallocate as we copy (simplistic)

Reuse when possible

Deallocate as we copy

Can't expand in-place

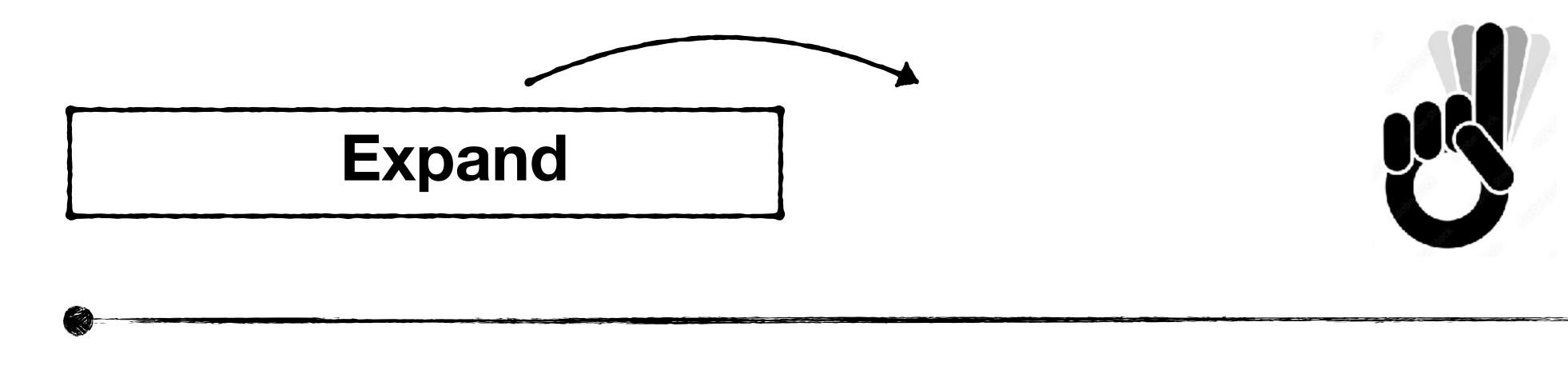
Assumptions:

Contiguous allocation

Reuse when possible

Deallocate as we copy

Can't expand in-place



Assumptions:

Contiguous allocation

Reuse when possible

Deallocate as we copy

Can't expand in-place

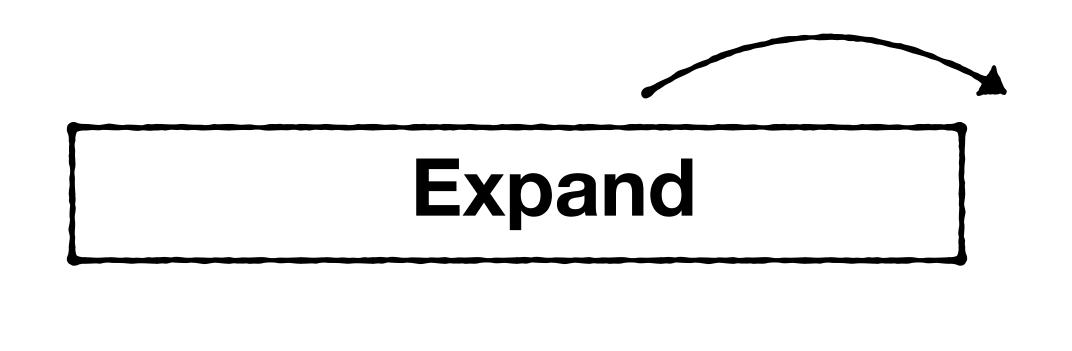


Assumptions: Contiguous allocation

Reuse when possible

Deallocate as we copy

Can't expand in-place



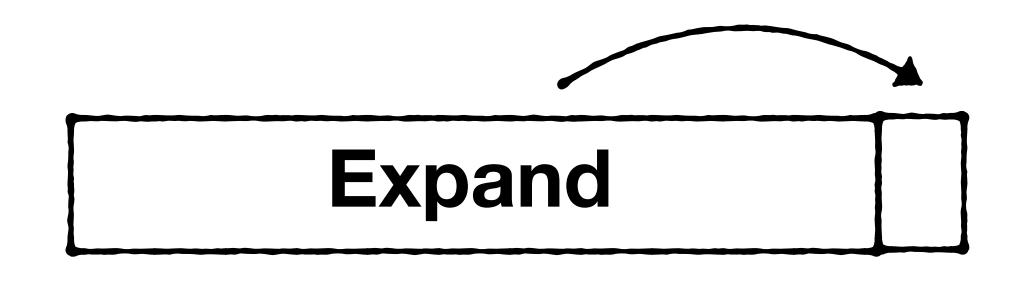
Memory addresses

Assumptions: Contiguous allocation

Reuse when possible

Deallocate as we copy

Can't expand in-place



Memory addresses

Assumptions: Contiguous allocation

Reuse when possible

Deallocate as we copy

Can't expand in-place

Memory addresses

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

Subject to: 
$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

Subject to: 
$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Clever ideas?



Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

1 1 2 3 5 8 13 21



1170 - 1250 Italy

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Satisfies this: 
$$\rightarrow$$
 Size  $_{i-2}$  + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

1 1 2 3 5 8 13 21 34 55 **89 144** 

**1.618** 

Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

$$\frac{\text{Size }_{i-1}}{\text{Size }_{i-2}} = \frac{\text{Size }_{i}}{\text{Size }_{i-1}} = G$$

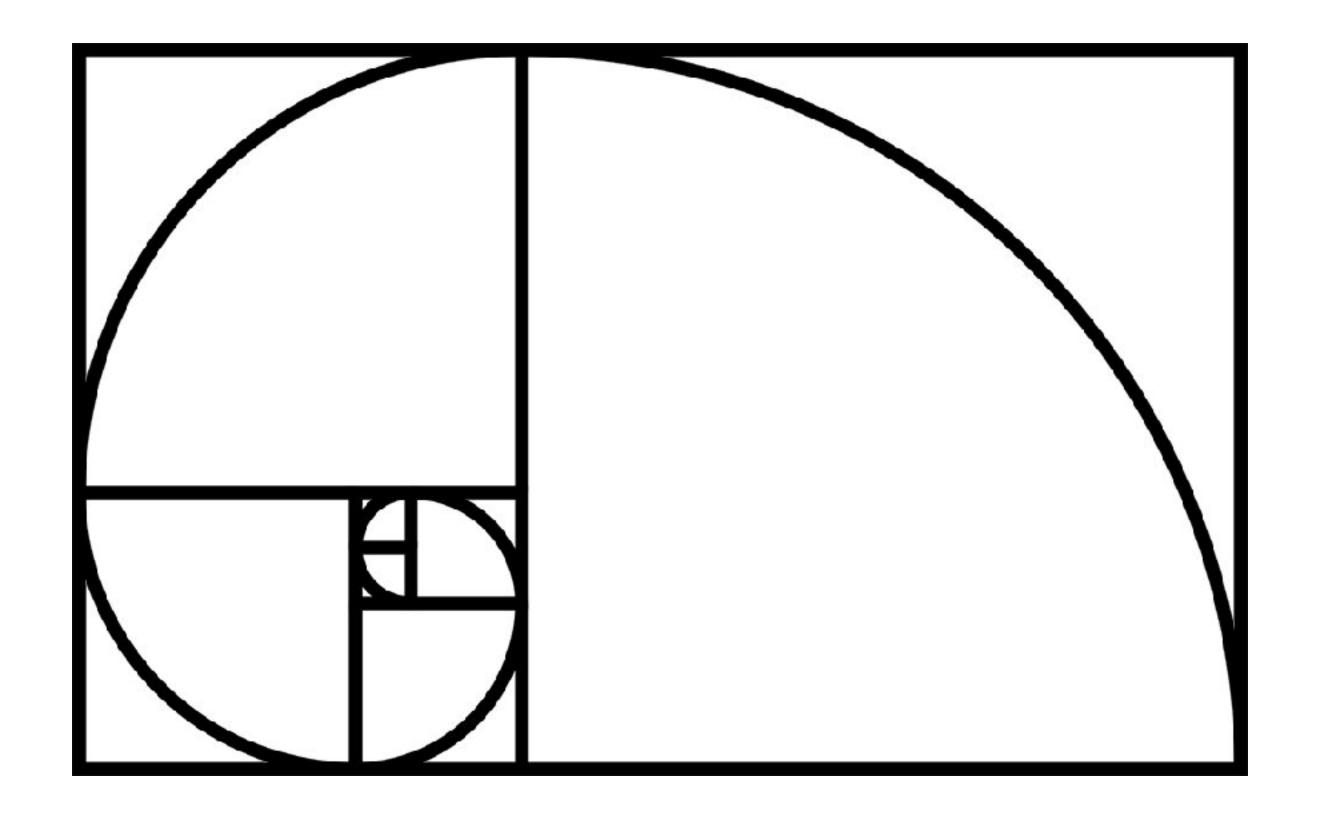


Ratio converges to the "Golden Ratio" ф

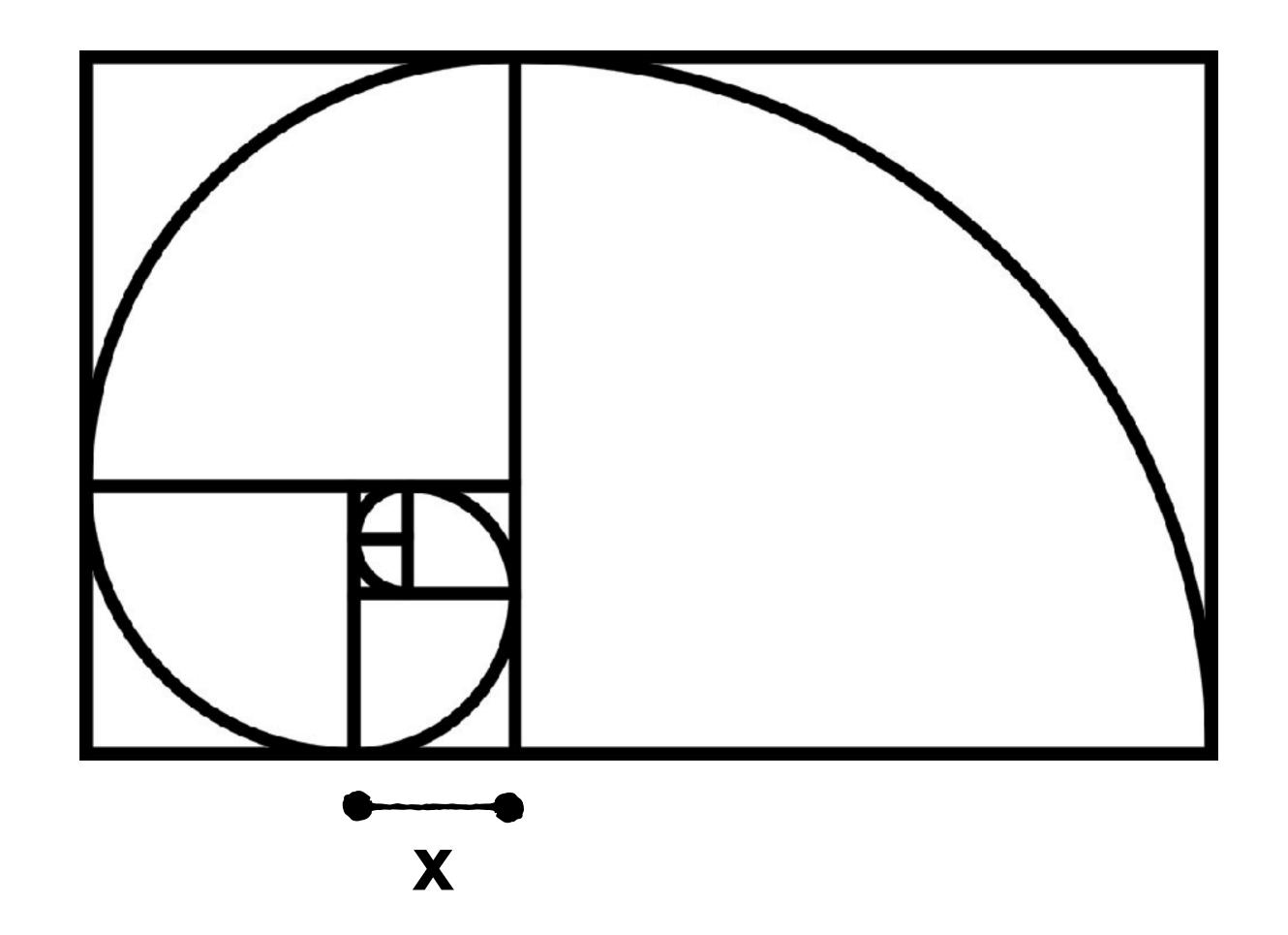
1 1 2 3 5 8 13 21 34 55 89 144

Ratio converges to the "Golden Ratio"  $\phi = 1.618033988749...$ 

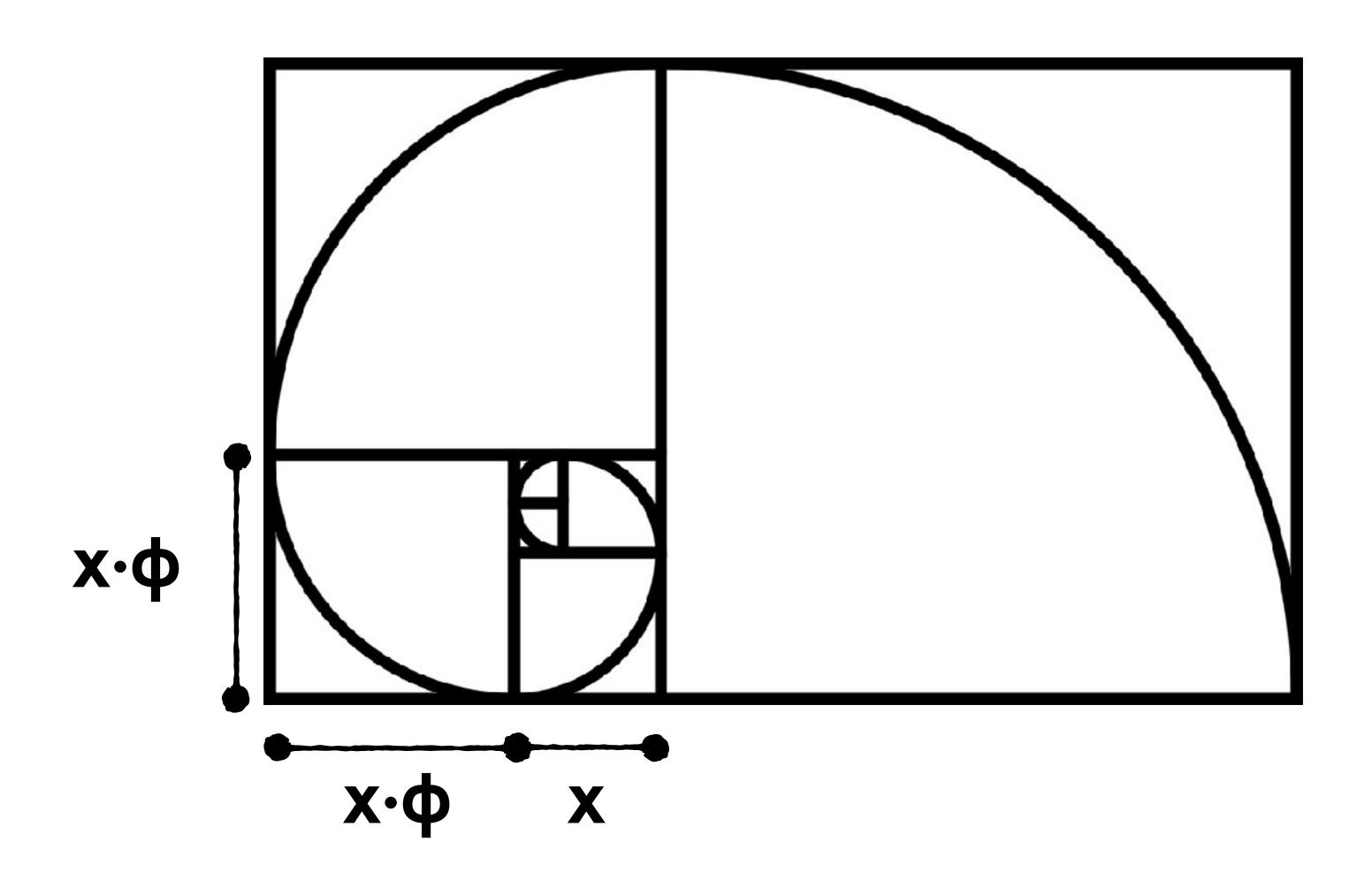
Ratio converges to the "Golden Ratio"  $\phi = \frac{1 + \sqrt{5}}{2}$ 

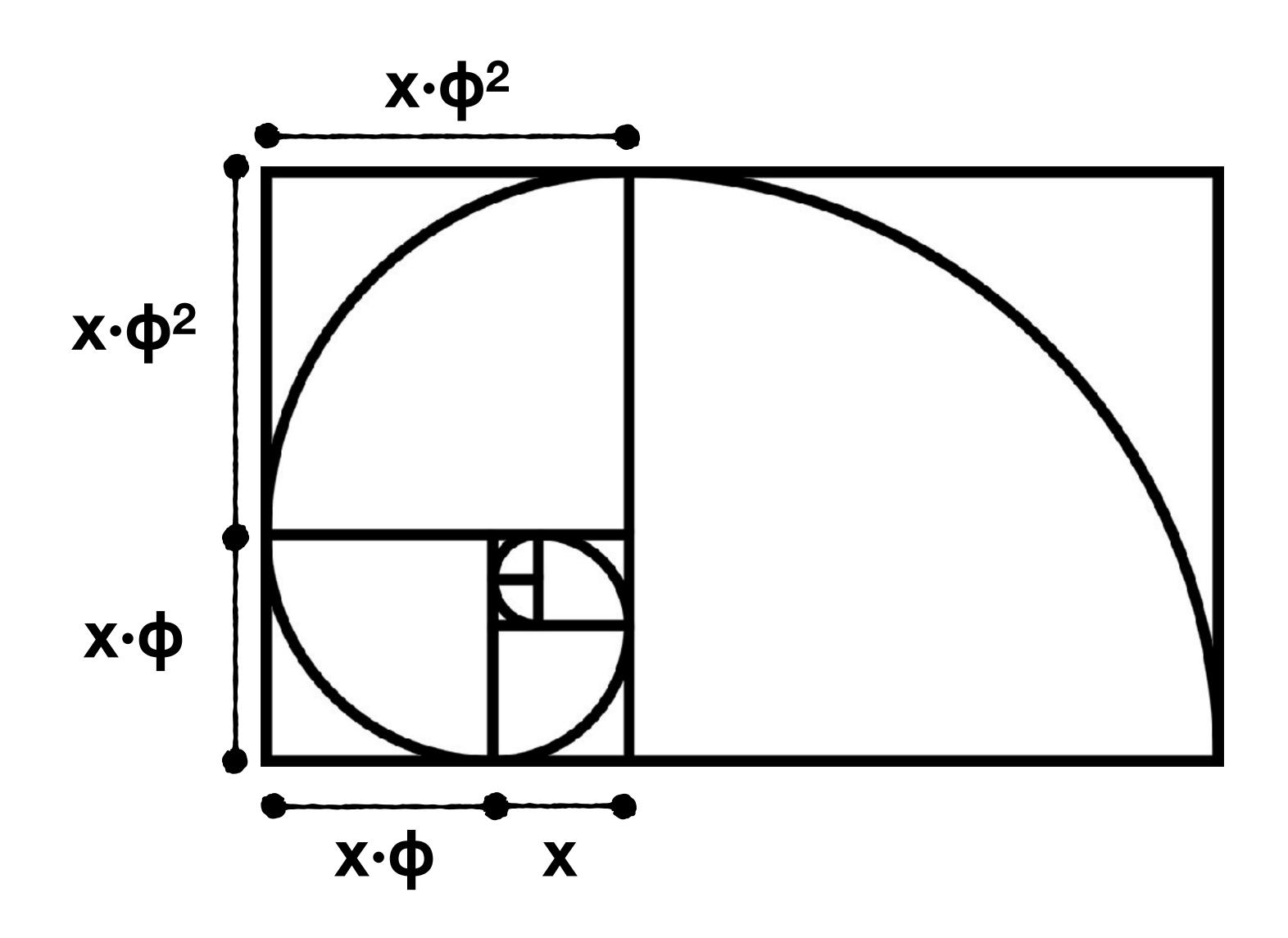


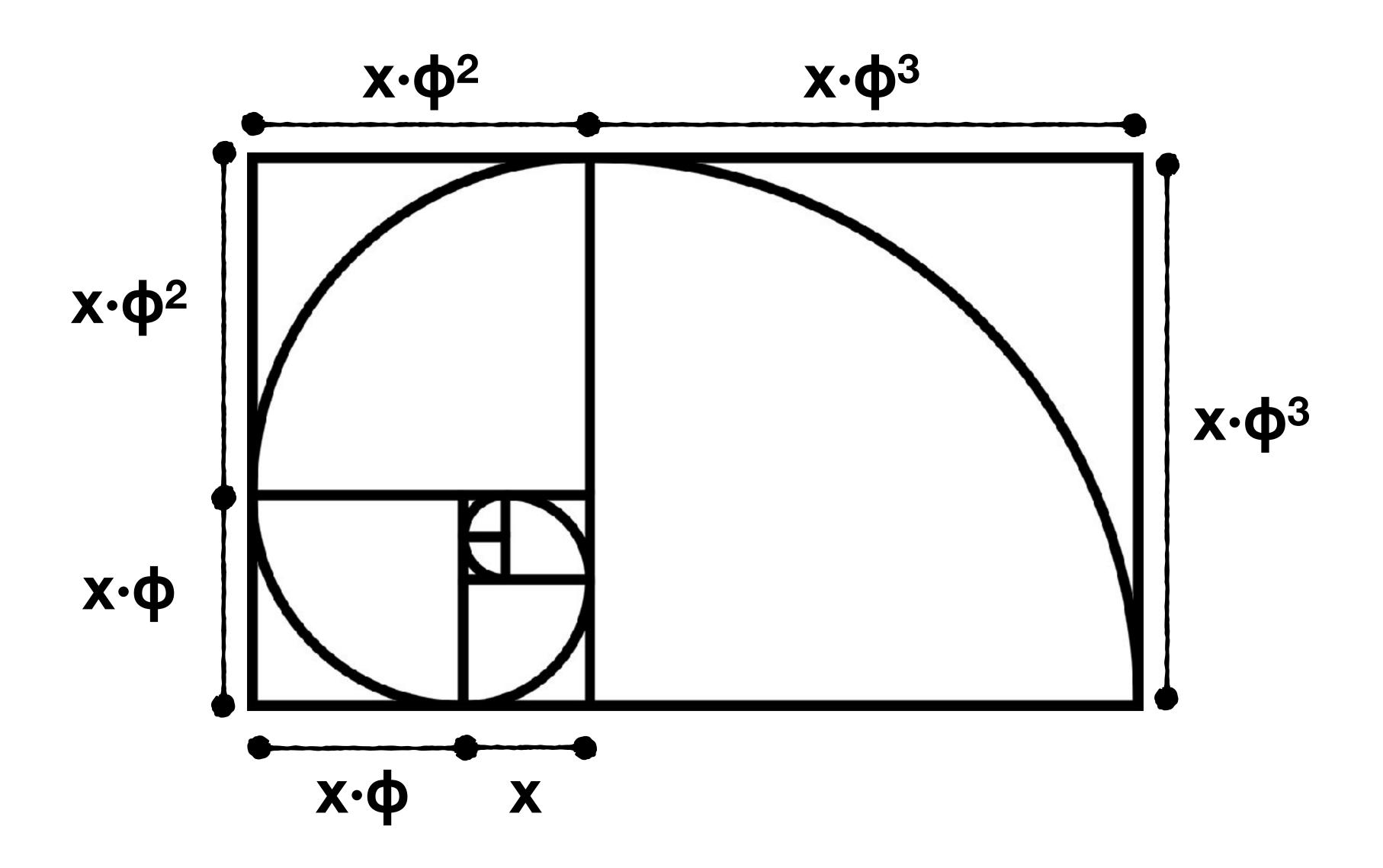
Ratio converges to the "Golden Ratio"  $\phi = \frac{1 + \sqrt{5}}{2}$ 

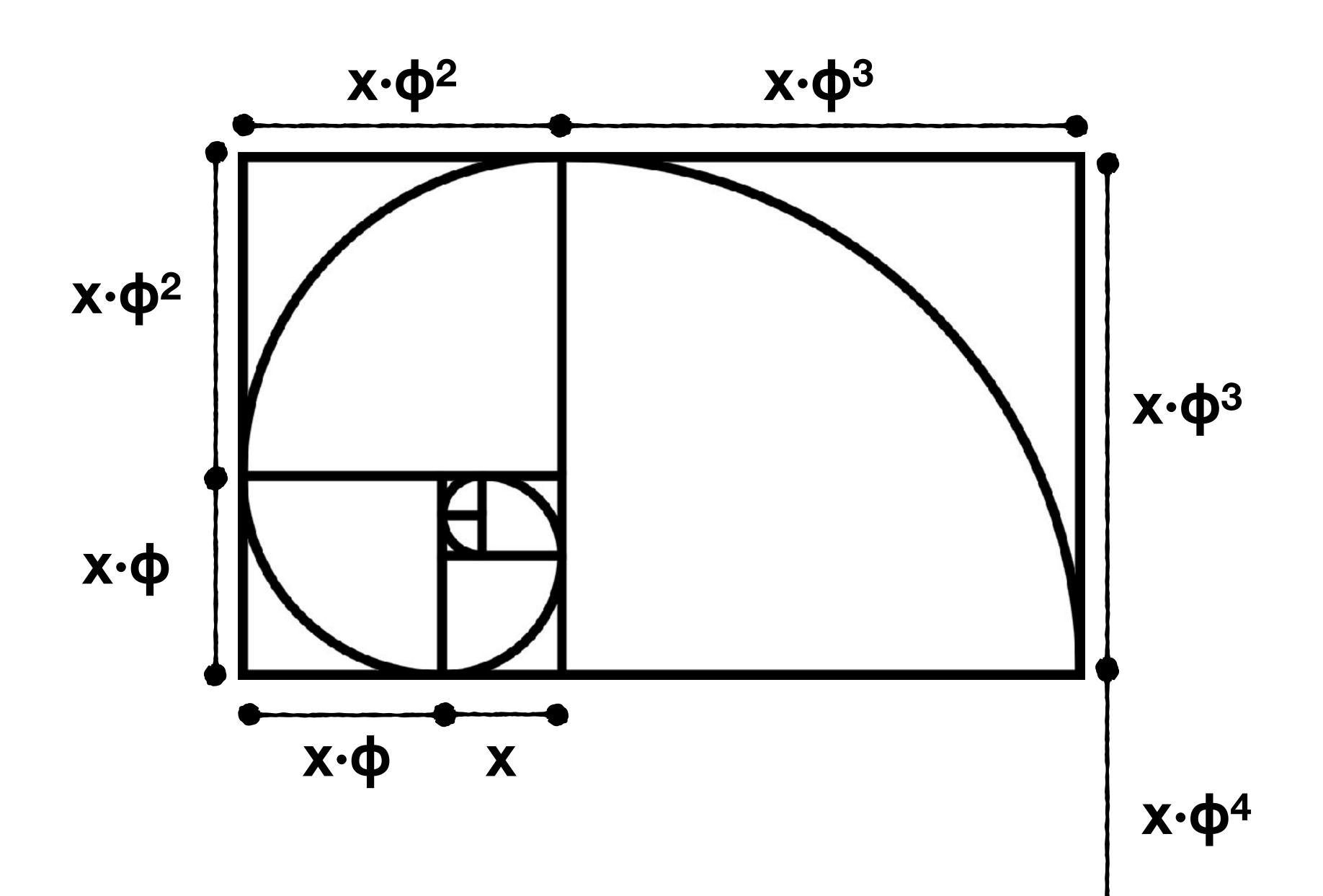


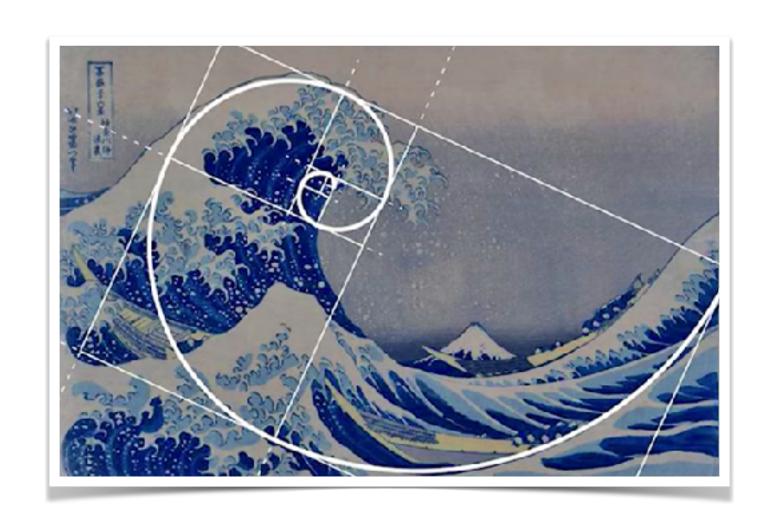
Ratio converges to the "Golden Ratio"  $\phi = \frac{1 + \sqrt{5}}{2}$ 





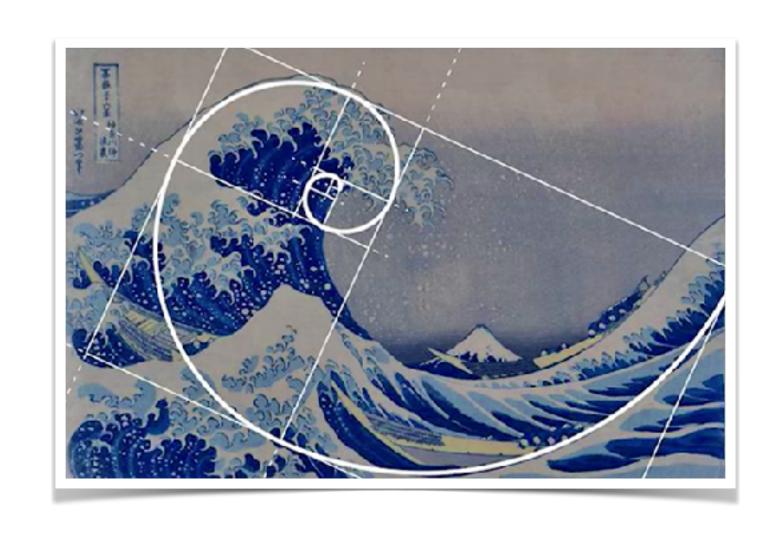


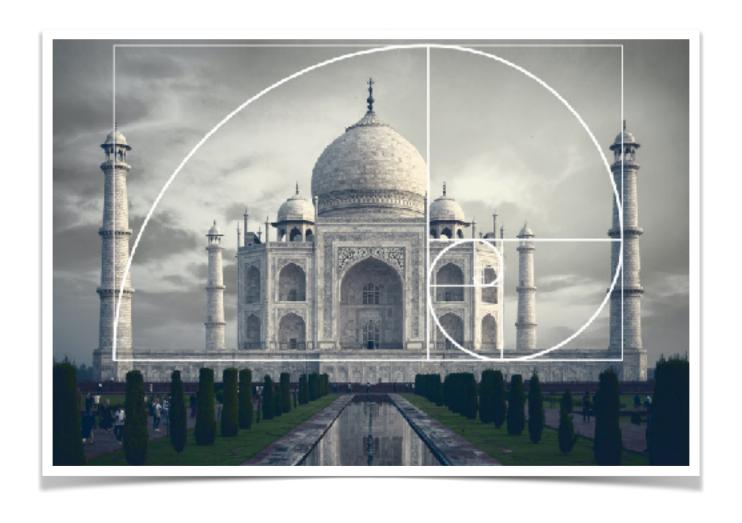




### Art

The Great Wave off Kanagawa





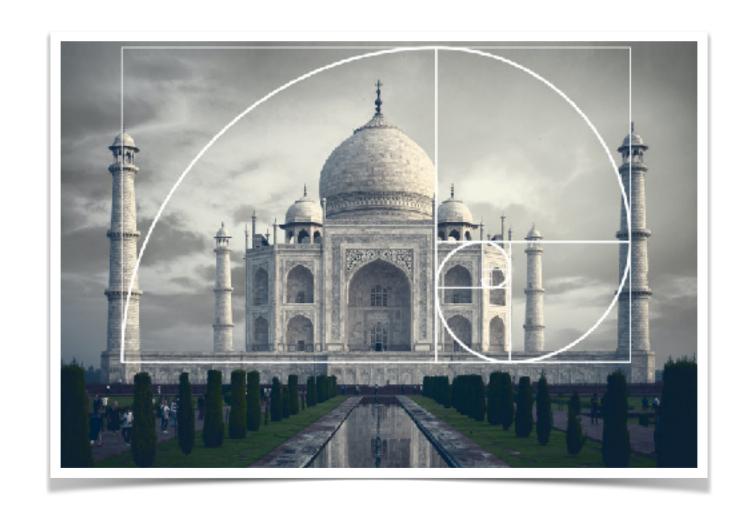
Art

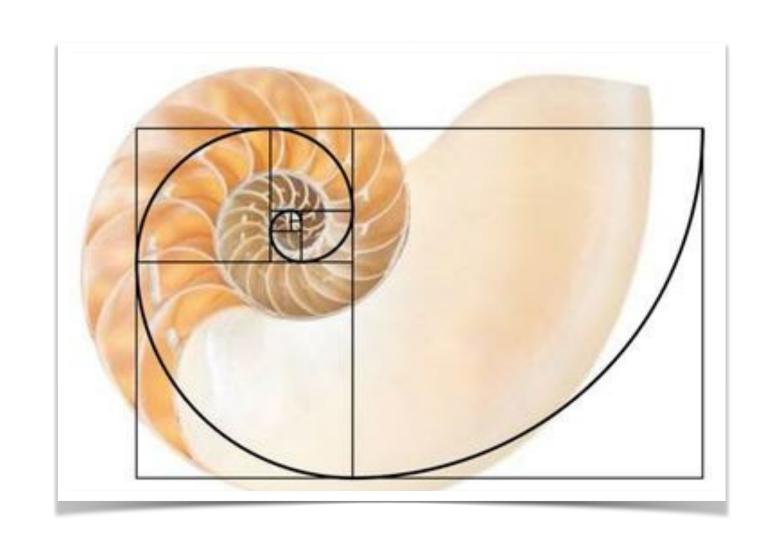
The Great Wave off Kanagawa

Architecture

Taj Mahal







Art

The Great Wave off Kanagawa

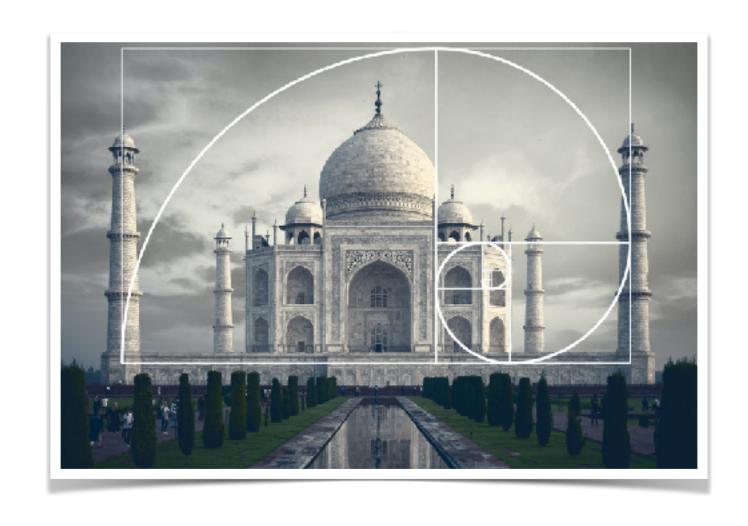
Architecture

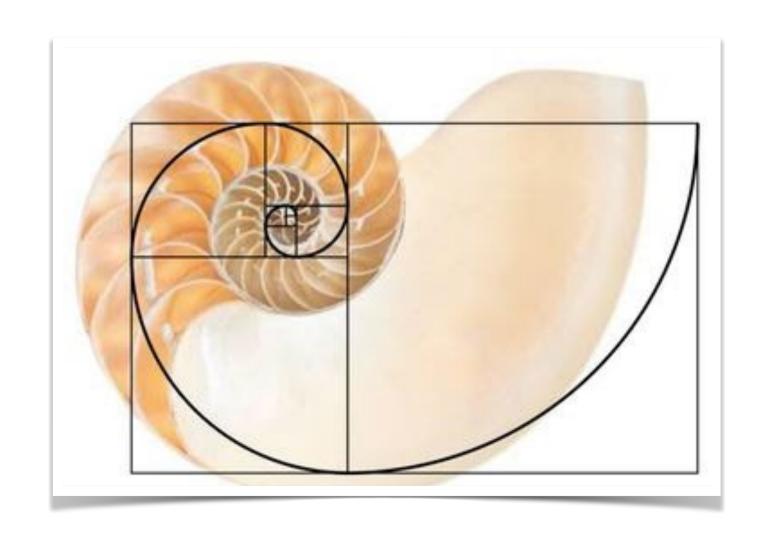
Taj Mahal

Nature

**Nautilus Shell** 







Art Architecture

Nature

And now also in computer science:)

### **Weird Properties**

$$\Phi = \Phi^2 - 1 = \frac{1}{\Phi - 1}$$

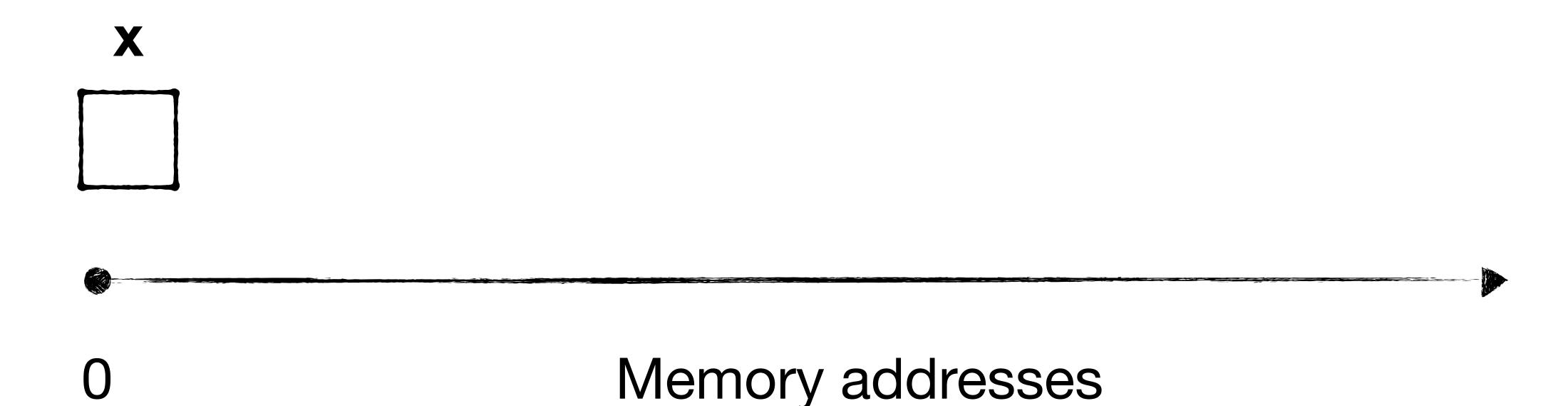
Size 
$$_{i-2}$$
 + Size  $_{i-1}$  = Size  $_i$ 

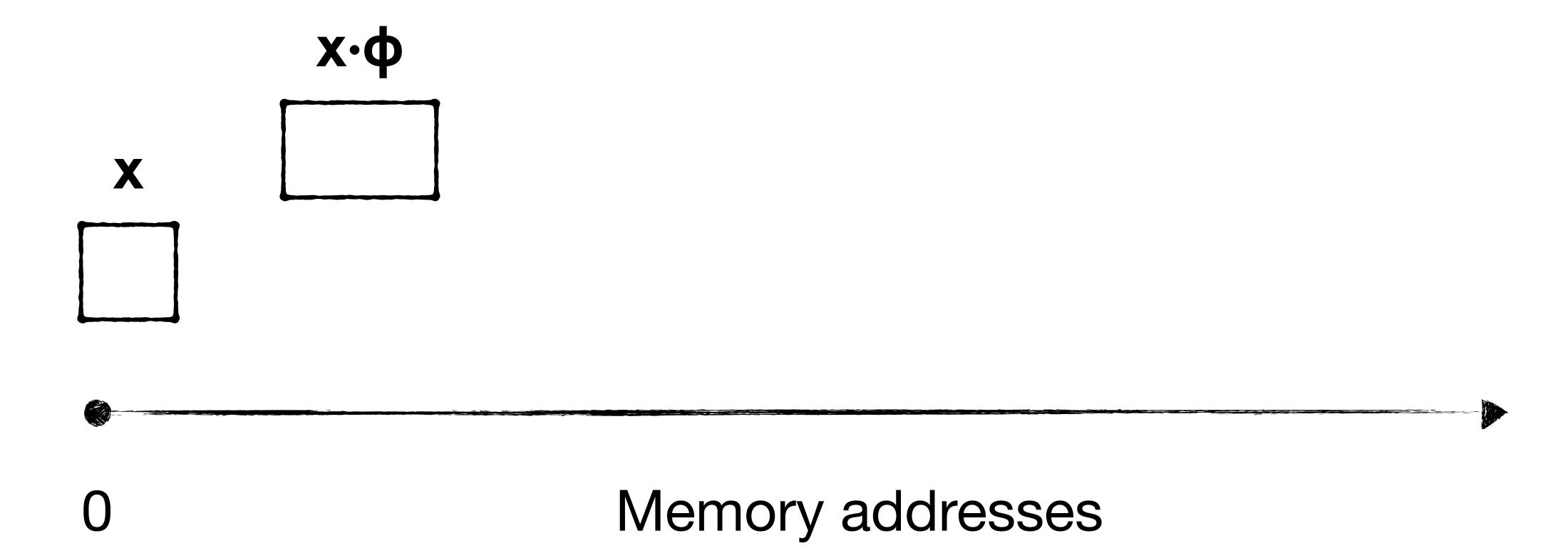
$$Size_{i-2} + Size_{i-1} = Size_i$$

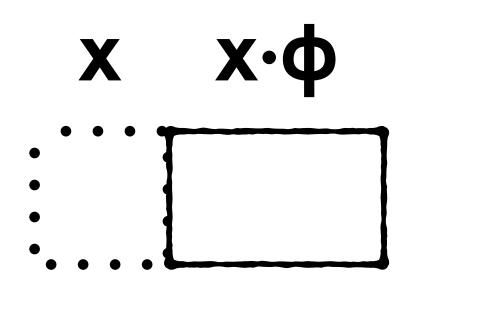
$$Size_{i-1} = Size_i$$

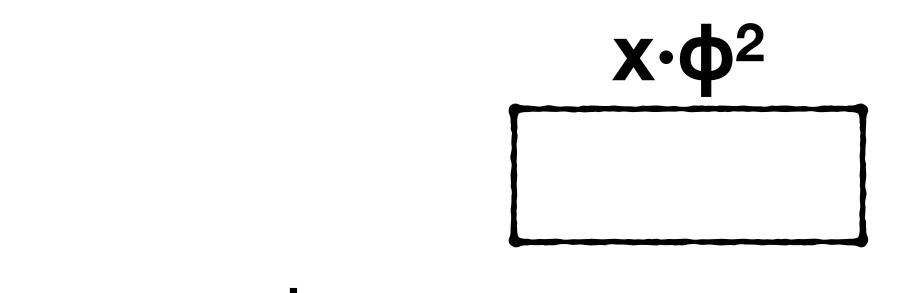
$$Size_{i-2} = Size_i$$

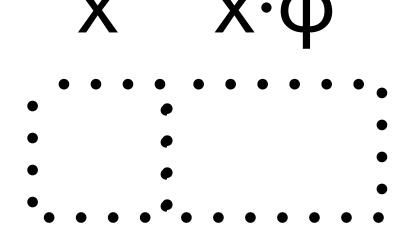
$$Size_{i-1} = G$$

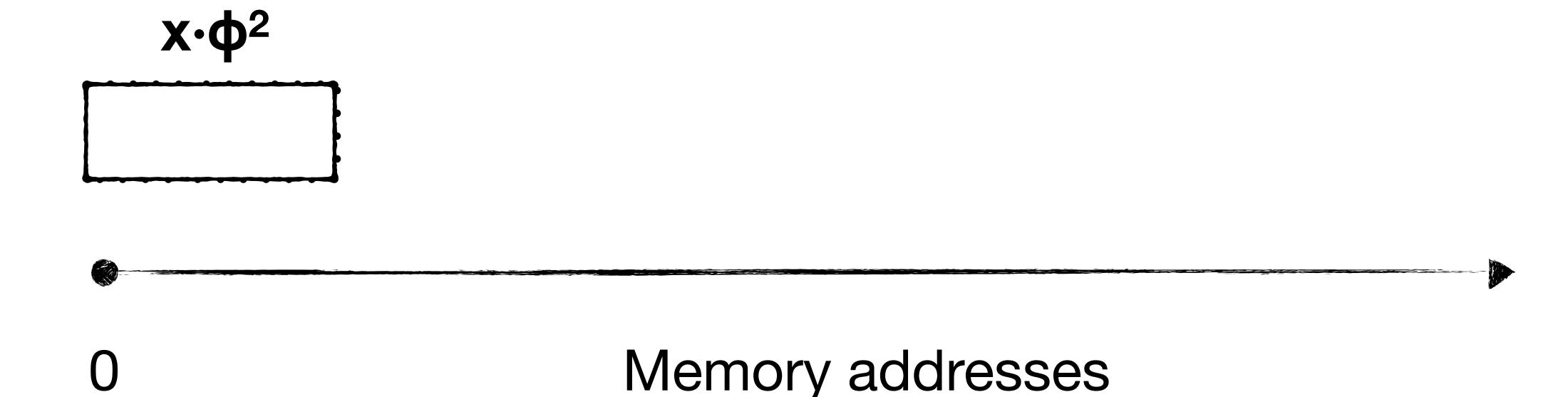


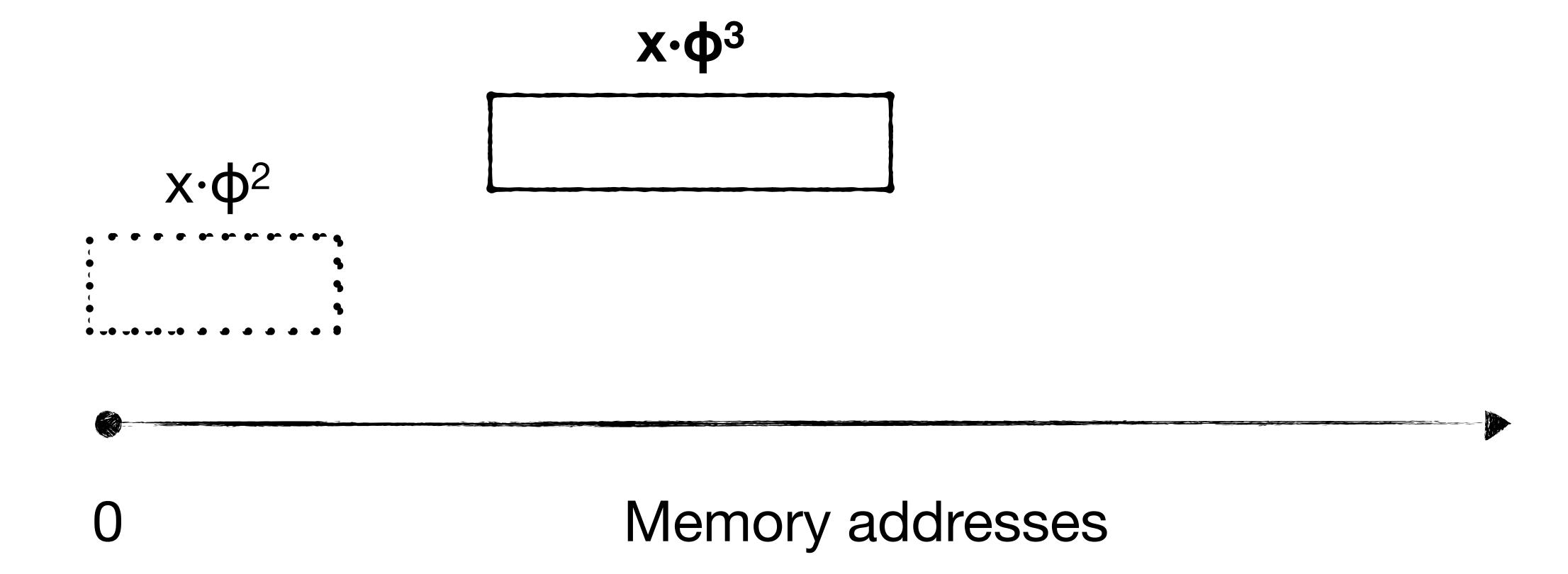


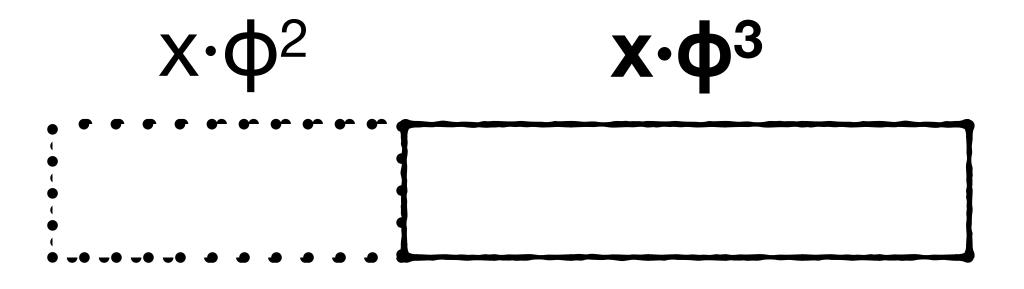


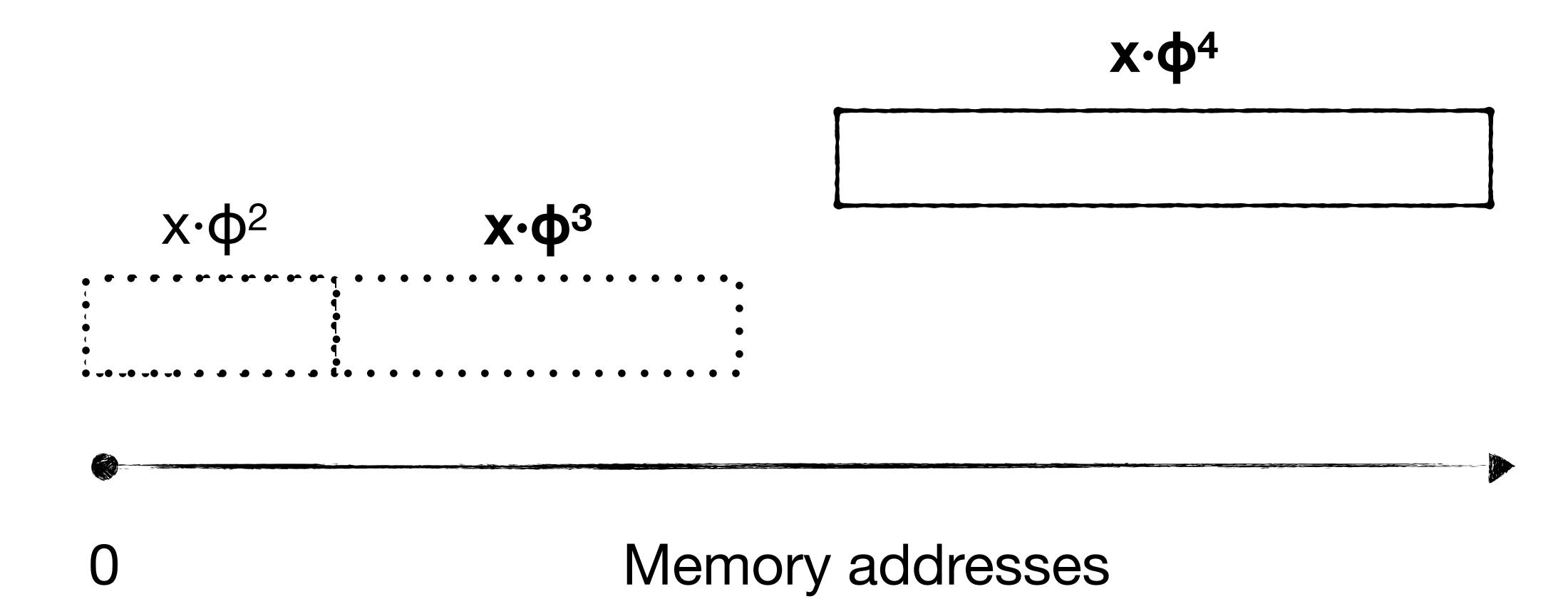


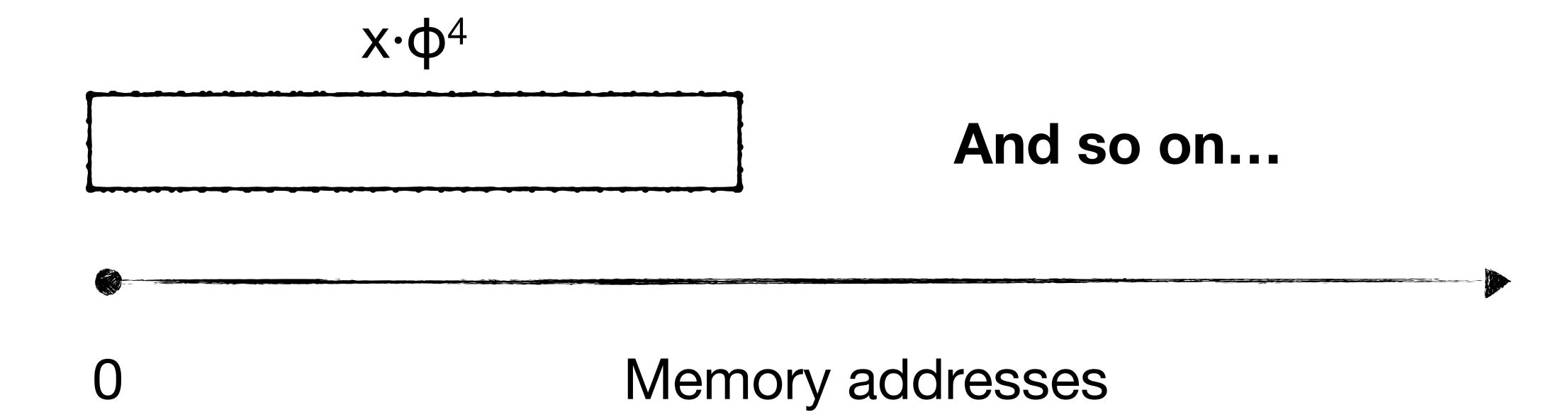




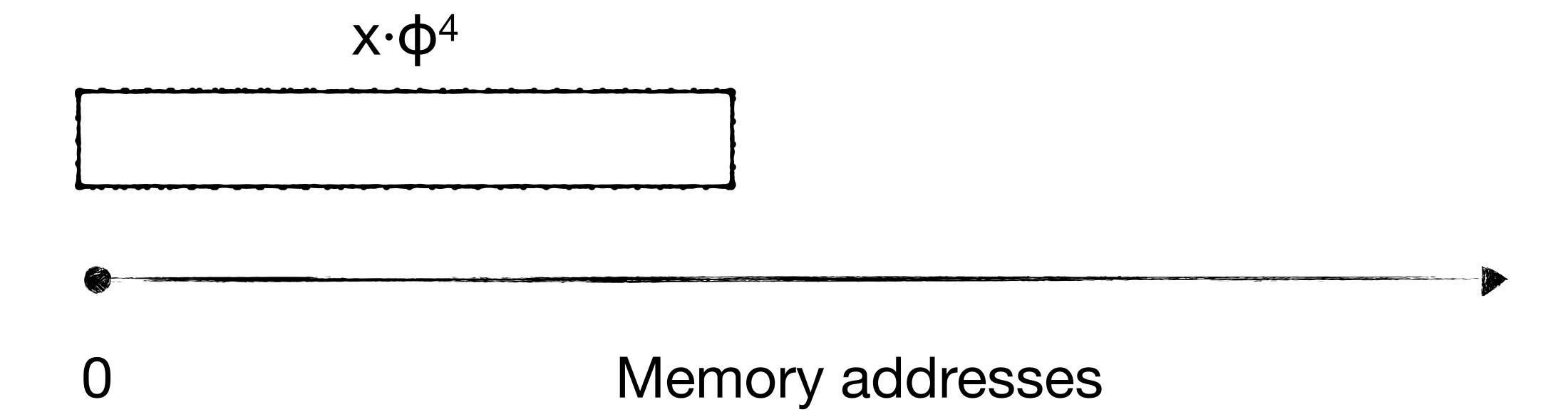








#### Write-Amplification

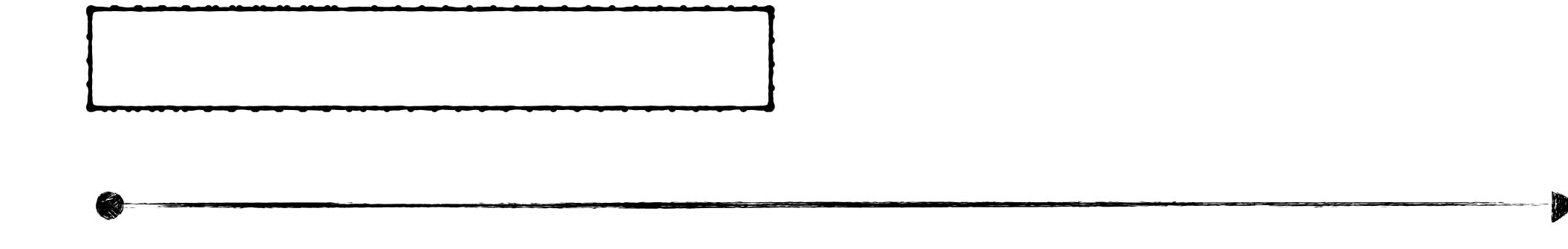


Write-Amplification = 
$$\frac{G}{G-1}$$

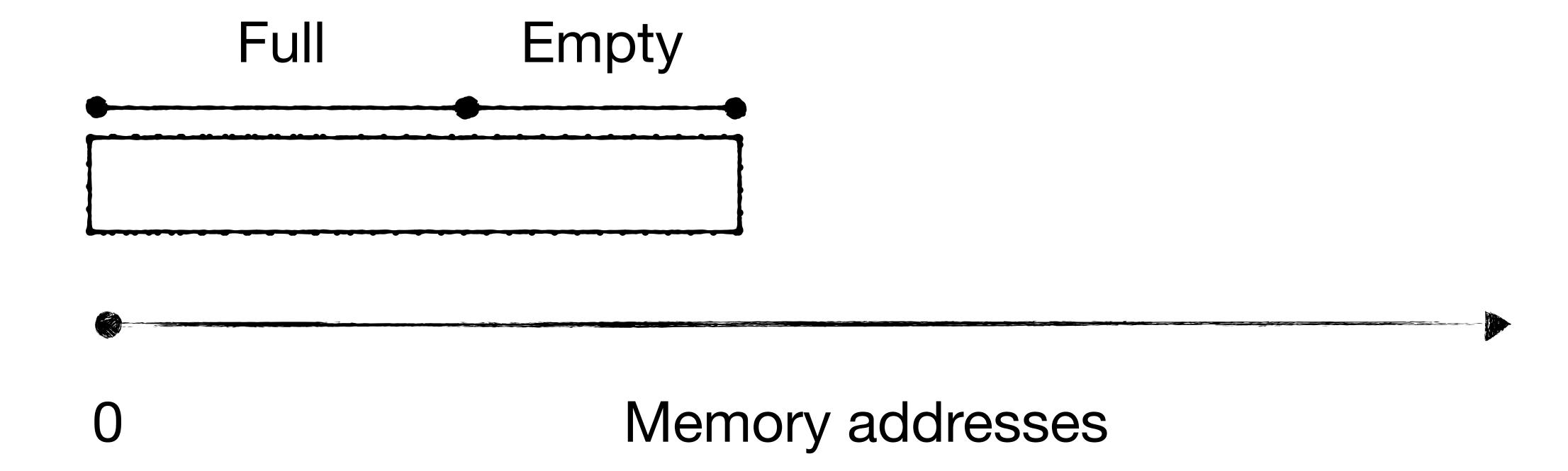
Write-Amplification = 
$$\frac{G}{G-1} = \frac{\Phi}{\Phi-1}$$

Write-Amplification = 
$$\frac{G}{G-1} = \frac{\varphi}{\varphi-1} = \varphi+1$$

#### Space-Amplification?

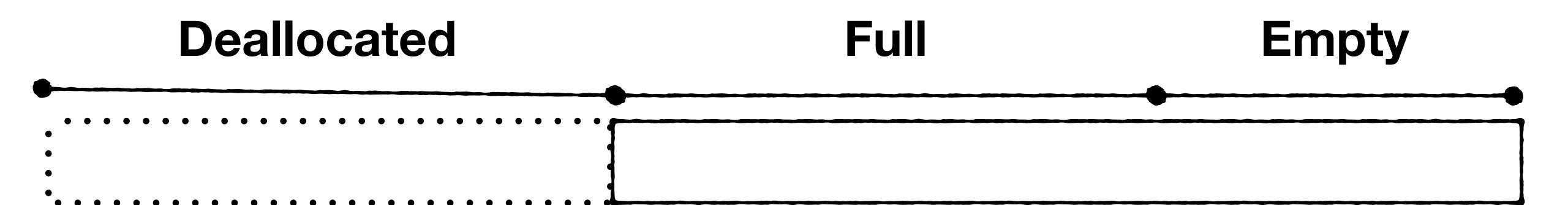


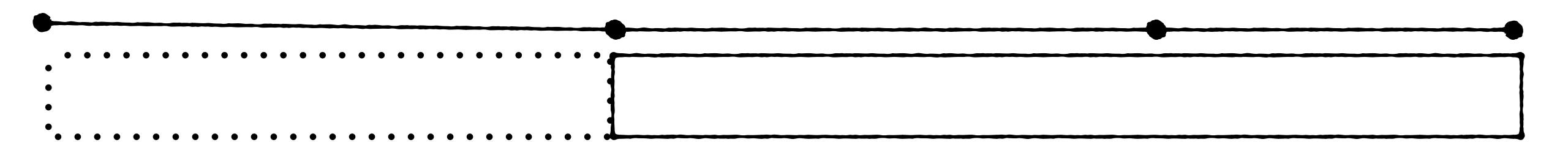
#### **Space-Amplification?**



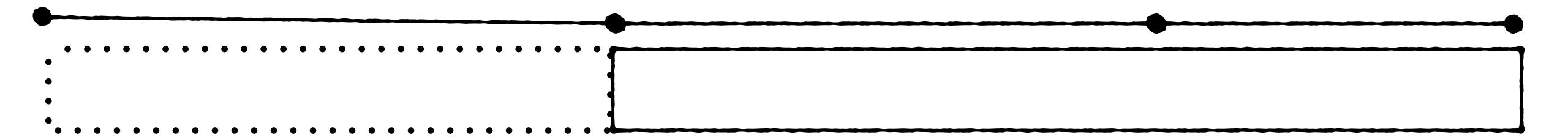
Space-Amplification = 
$$\frac{\text{Full} + \text{Empty}}{\text{Full}} = G = \Phi$$

#### Max Space-Amp?





Max Space-Amp =  $1 + \phi$  = 2.61



#### Write-amp

Space-amp

$$\frac{G}{G-1} + G$$

# Alternates G to G+1

#### In the wild

Implementation	Growth factor
Java ArrayList	1.5
Python PyListObject	~1.125
Microsoft Visual C++ 2013	1.5
G++ 5.2.0	2
Clang 3.6	2
Facebook folly/FBVector	1.5
Rust Vec	2
Go slices	between 1.25 and 2
Nim sequences	2
SBCL (Common Lisp) vectors	2
C# (.NET 8) List	2

Source: https://en.wikipedia.org/wiki/Dynamic\_array#Growth\_factor

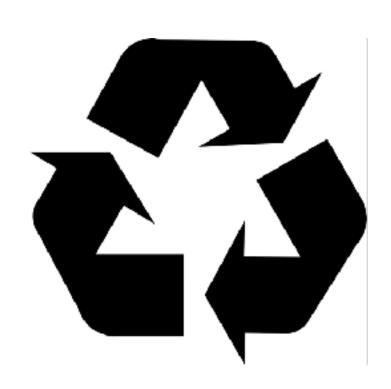
#### Facebook folly/FBVector

https://github.com/facebook/folly/blob/main/folly/docs/FBVector.md

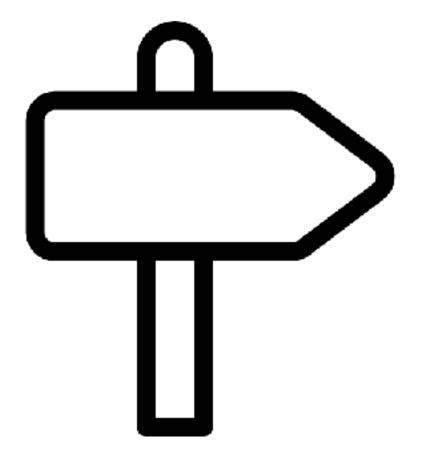
Real-world discussion of these issues

#### And now to new stuff

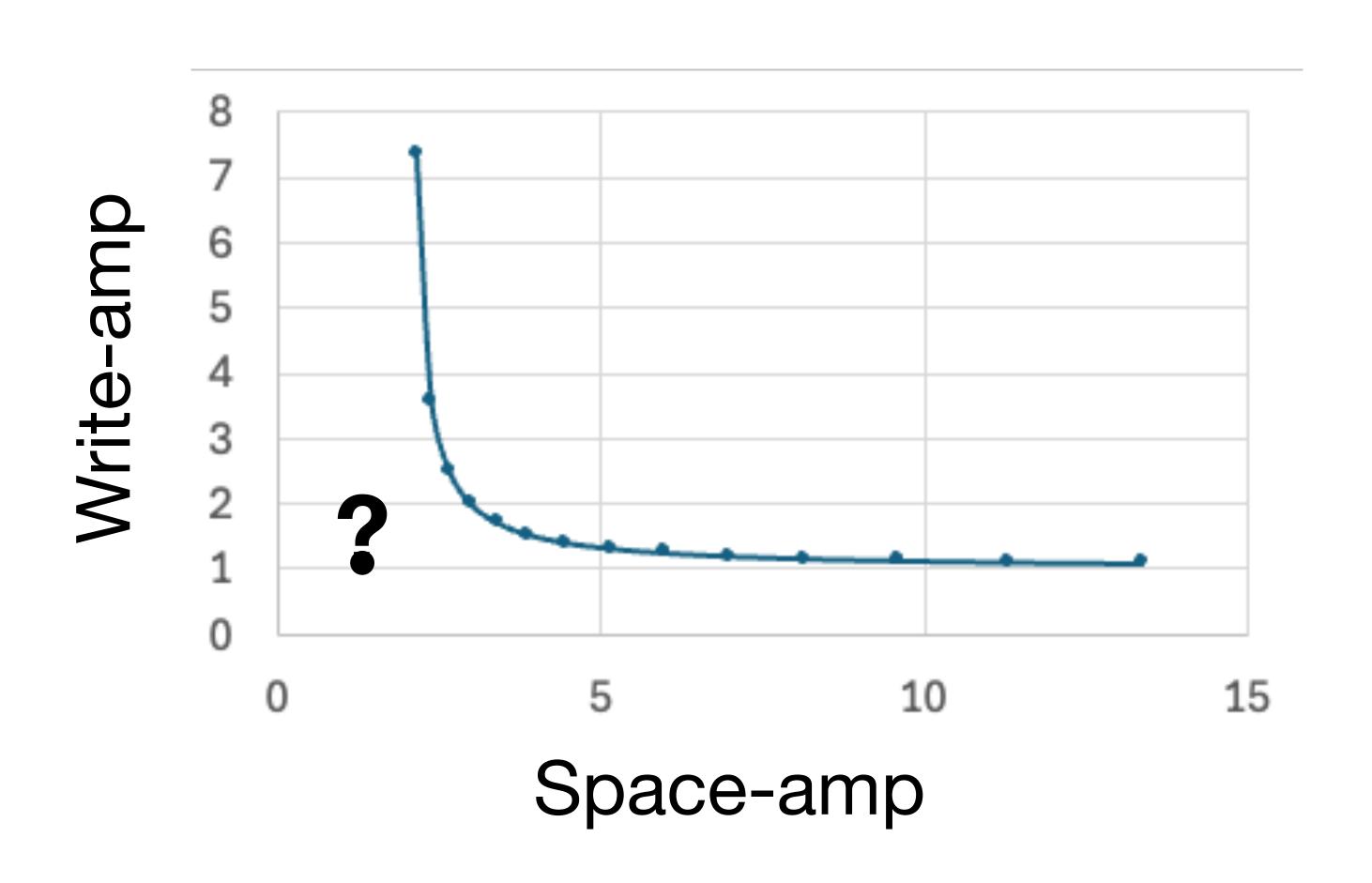
# Reusing Deallocated Space



### Alleviating tradeoff via indirection



#### Can we completely overcome this trade-off?



Array

Array Extra

Array Extra

Promise: write-amp of ???

Array Extra

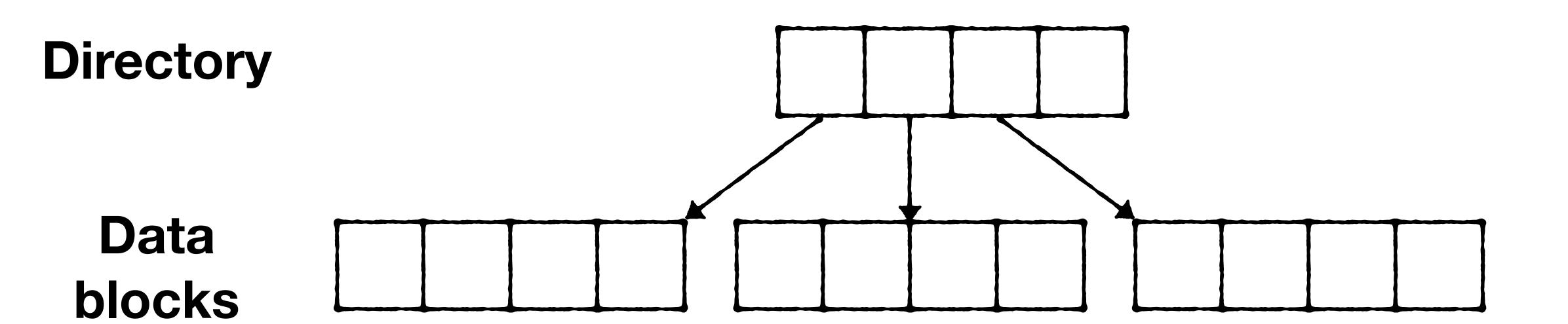
Promise: write-amp of 1

Array Extra

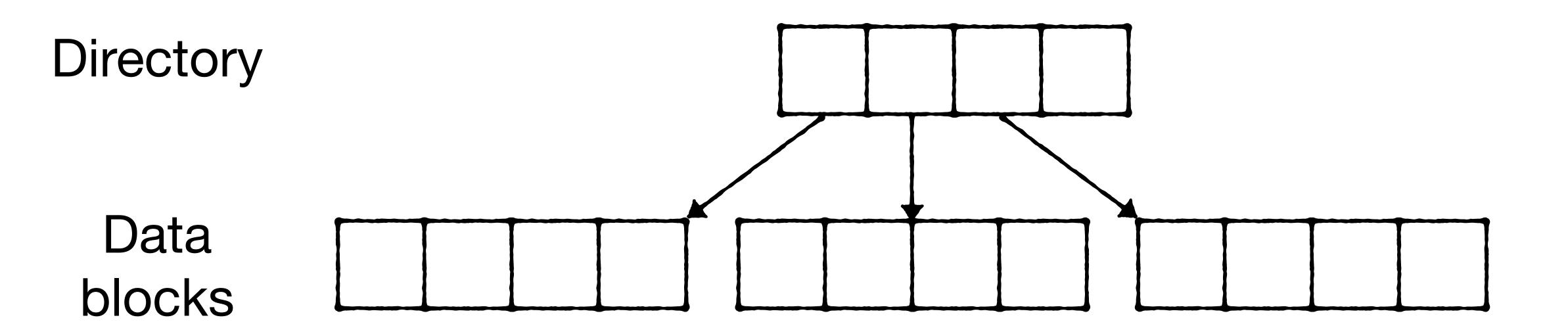
Promise: write-amp of 1

space-amp? We shall see :)

#### Add a layer of indirection

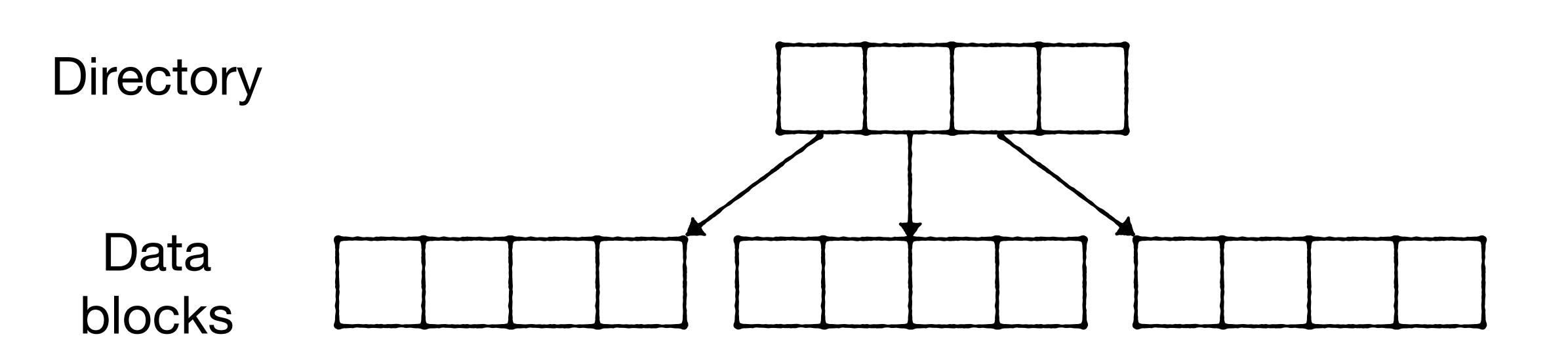


## get(i)



get(i)

#### Data block = [i / data block size]



get(i)

Data block = [i / data block size]

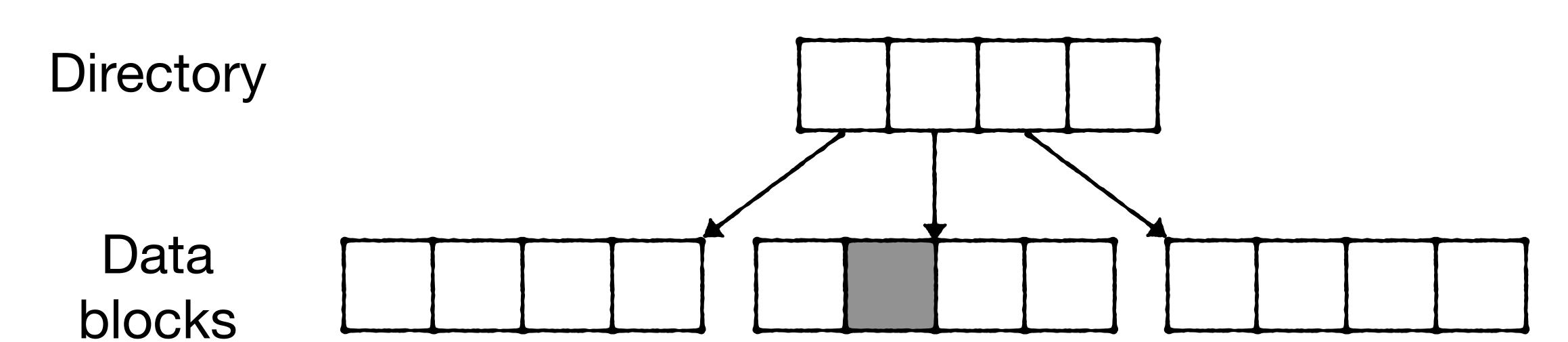
offset within = i % data block size

Directory

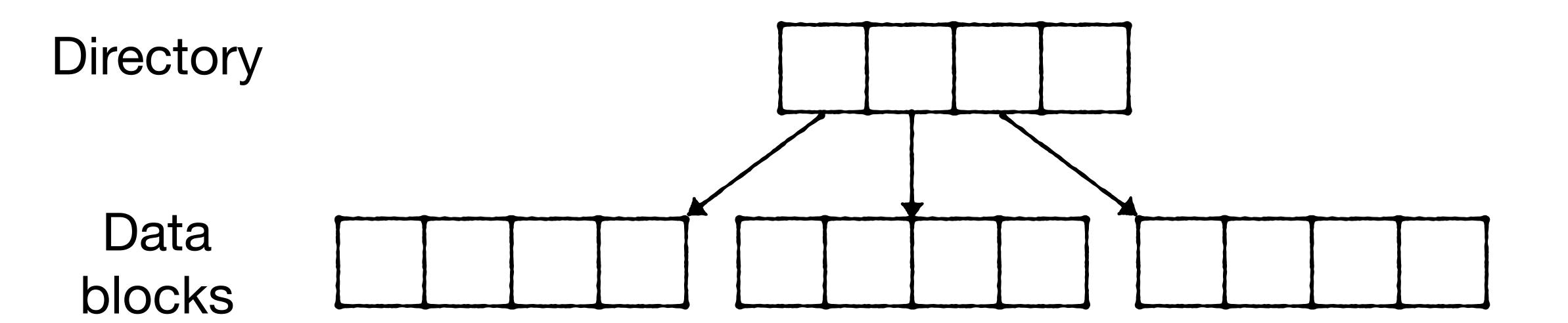
Data
blocks

**get(5)** 

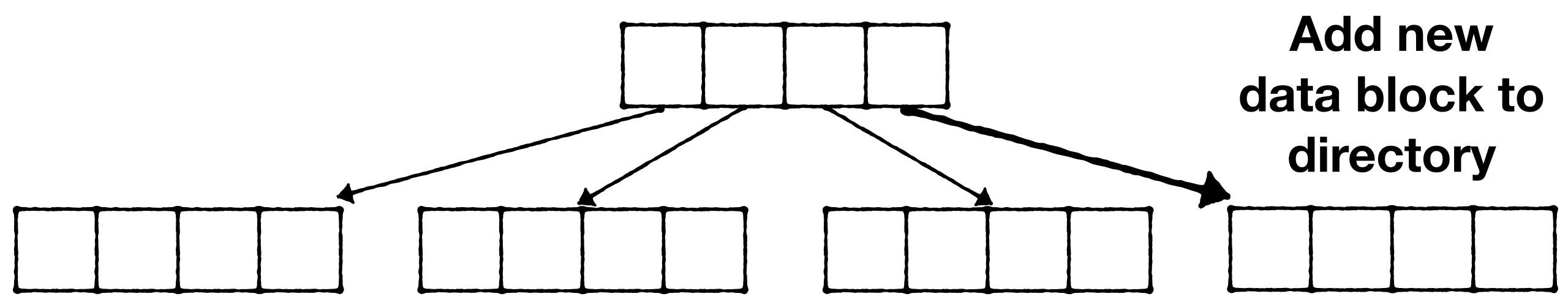
Data block =  $\lfloor 5/4 \rfloor = 1$ offset within = 5 % 4 = 1



## Expand?



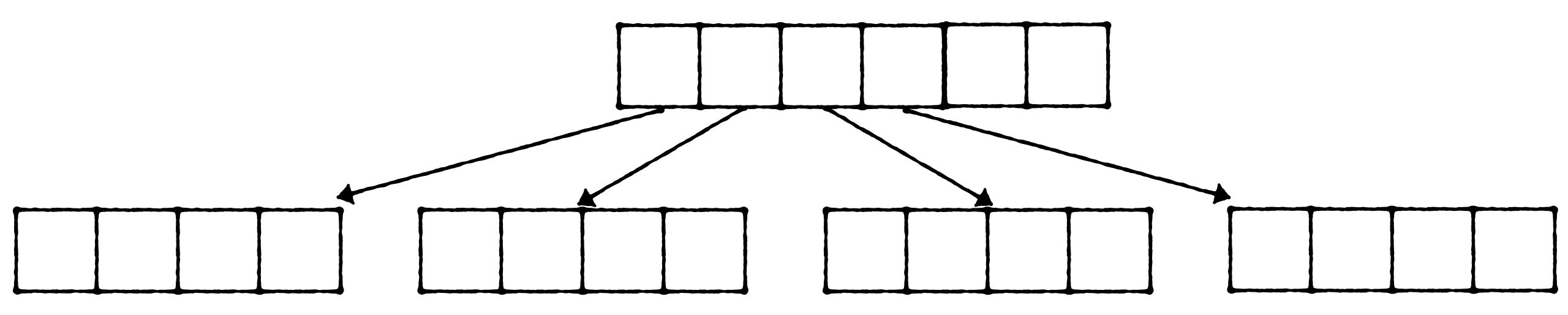
## Expand?



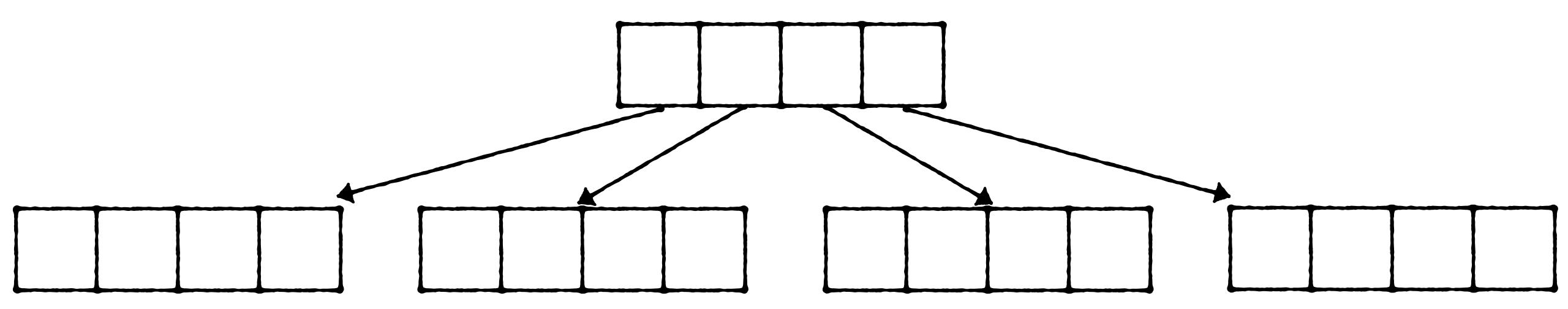
Expand? Expand directory if we need more space

Expand?

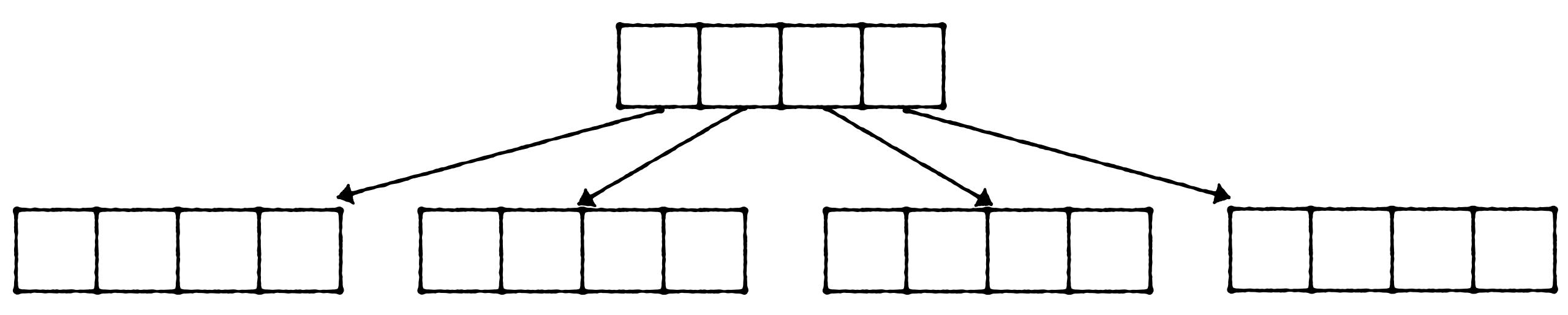
# Expand directory if we need more space



## Downside?

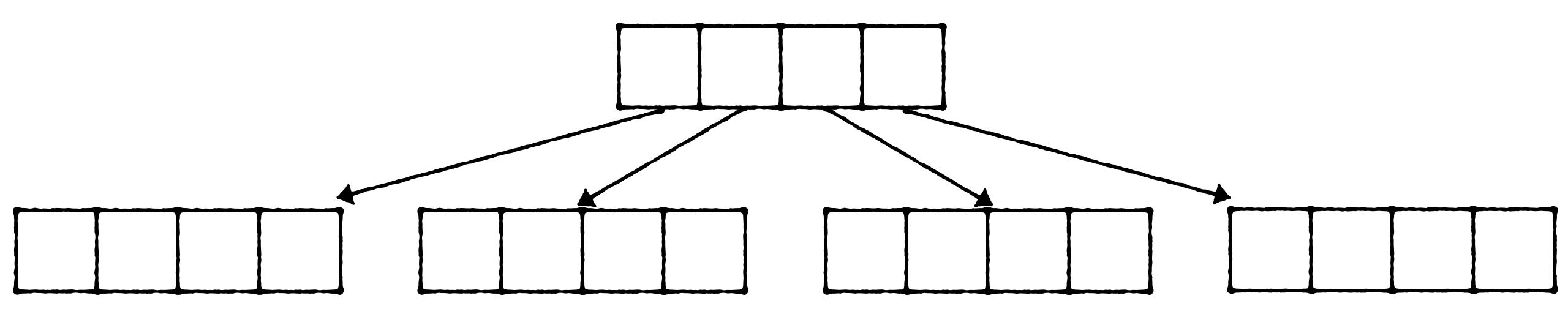


## Downside: 2 memory hops per access



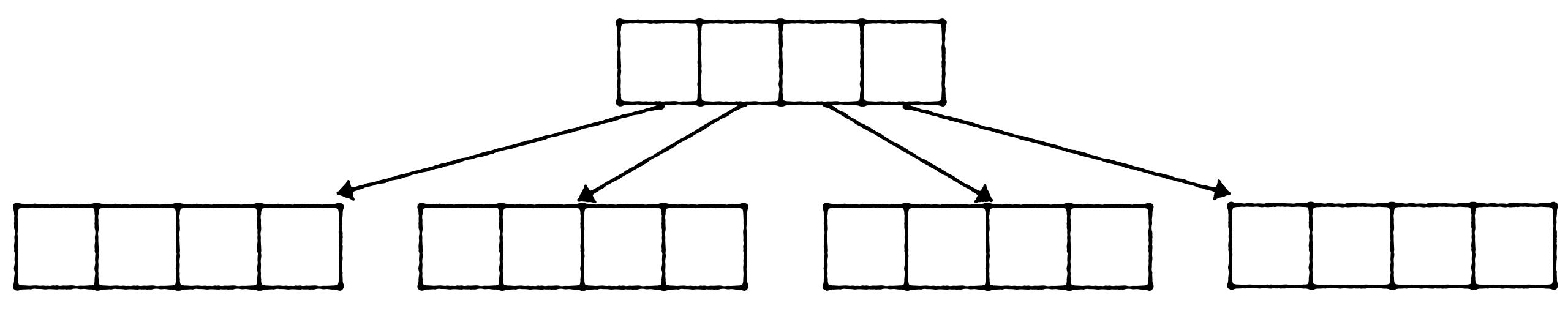
## Downside: 2 memory hops per access

## Mitigation?

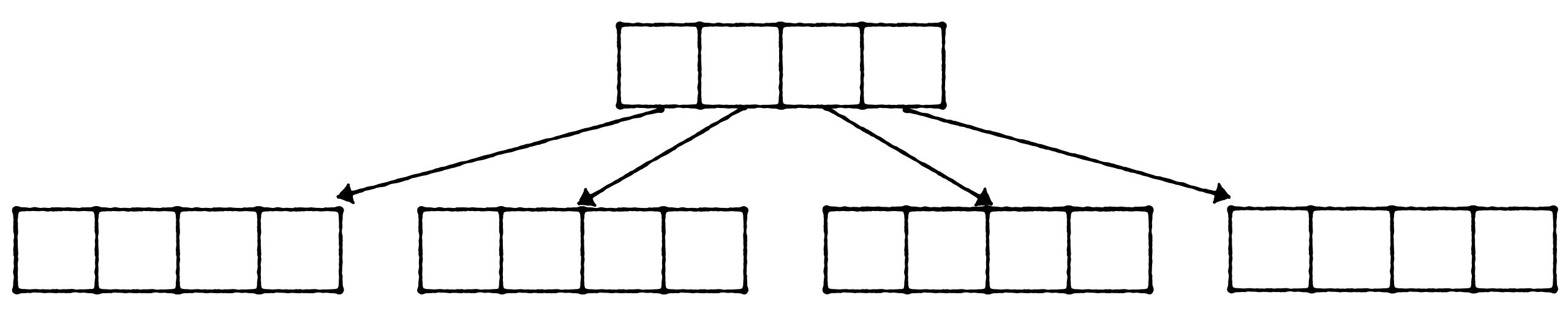


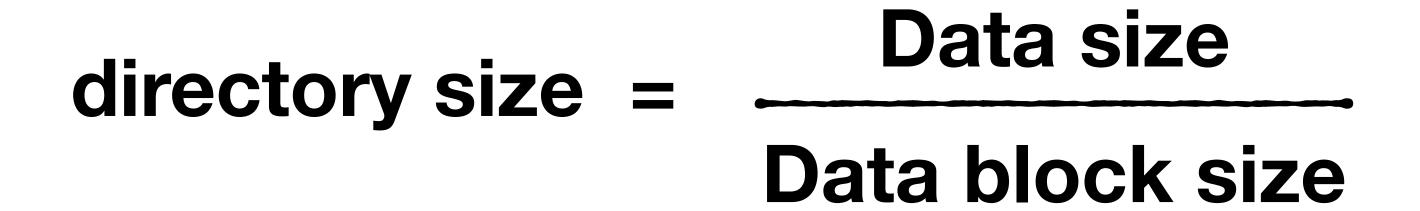
Downside: 2 memory hops per access

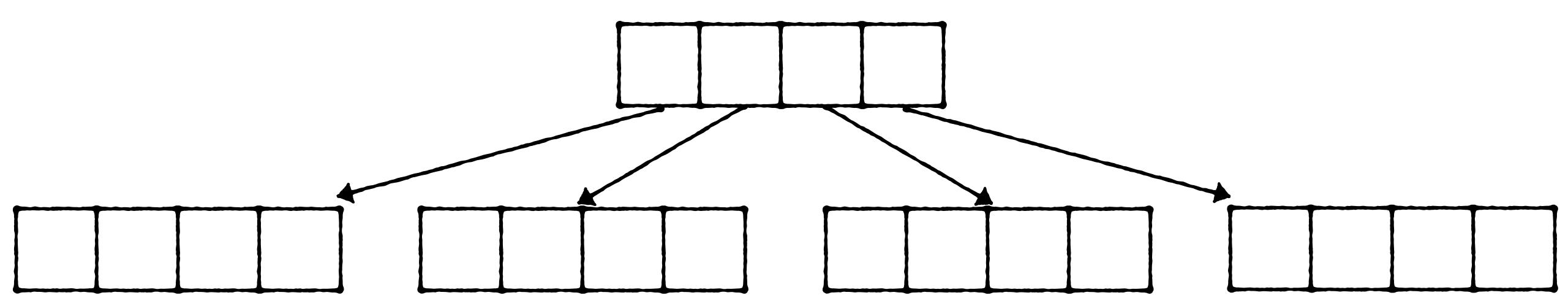
## Mitigation: directory must fit in L1 cache



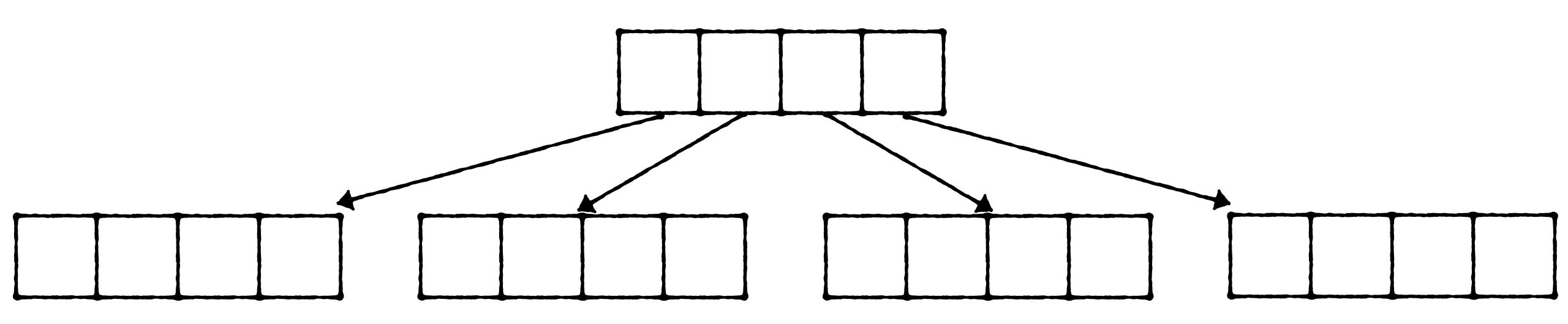
# directory size?





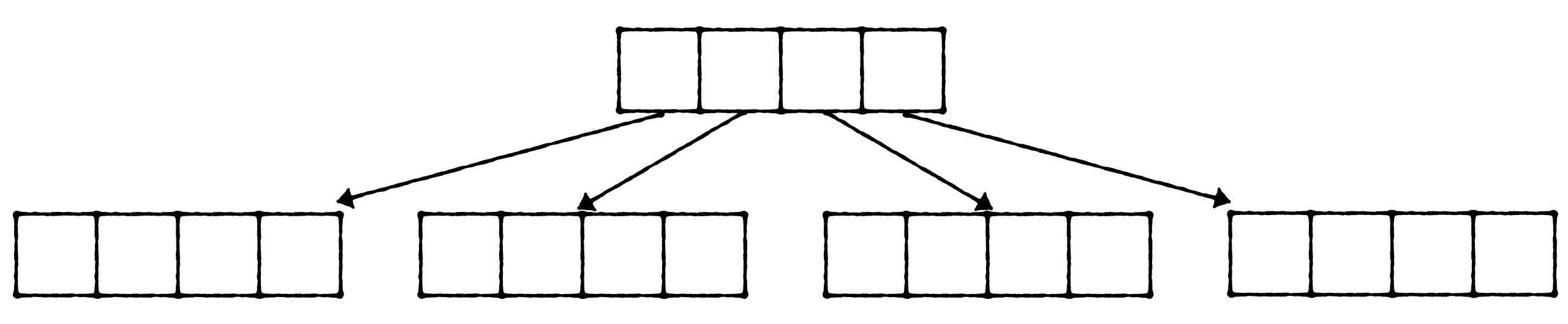






Risk: data blocks are initialized too small

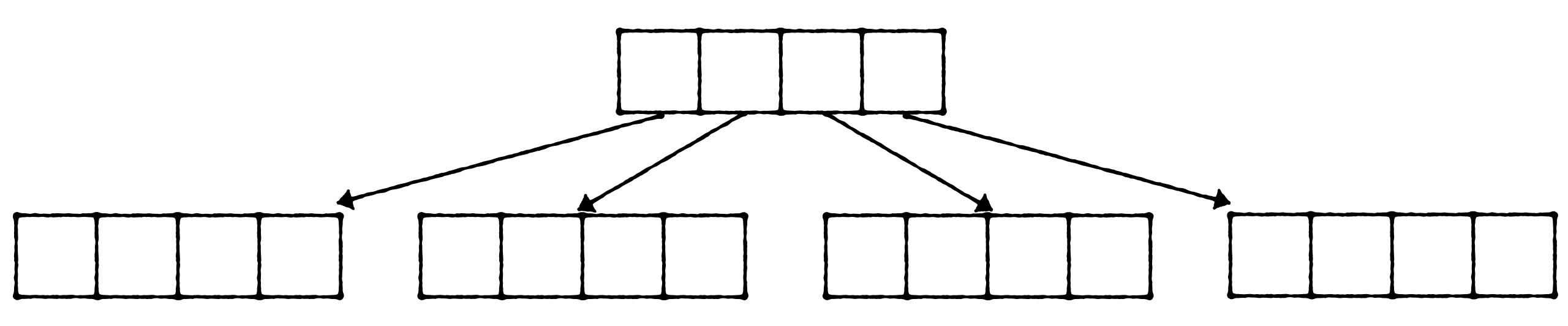




Risk: data blocks are initialized too small

Directory may outgrow the L1 cache

directory size = 
$$\frac{\text{Data size}}{\text{Data block size}} = O(N)$$

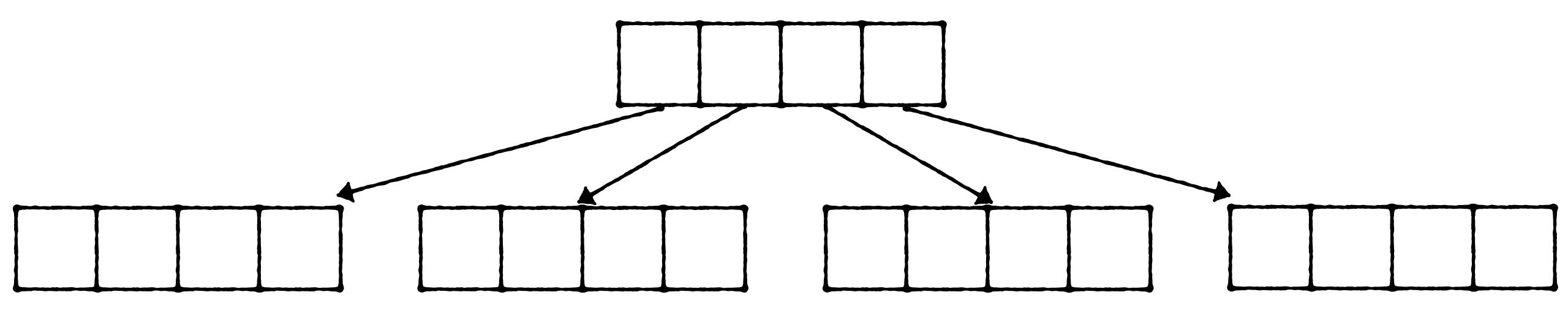


Risk: data blocks are initialized too small Directory may outgrow the L1 cache

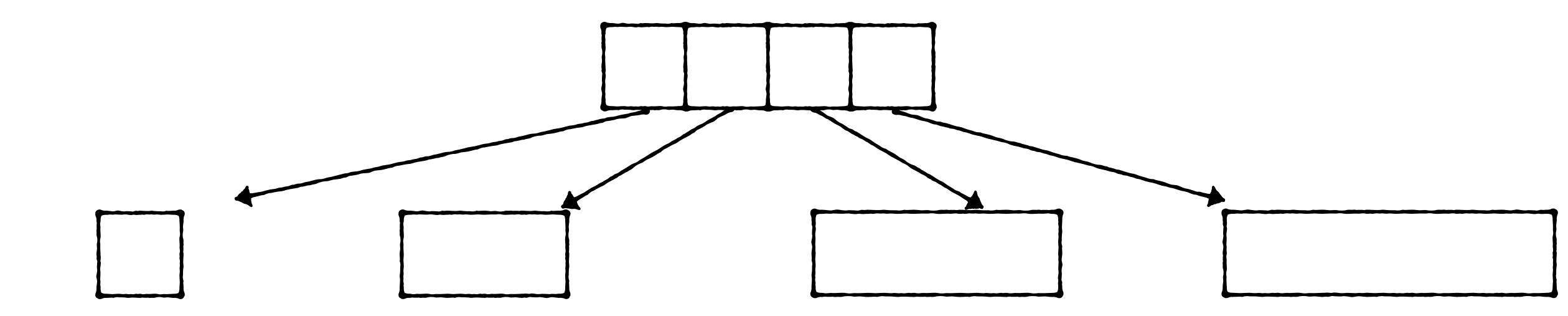
Solution?

## Resizable Arrays in Optimal Time and Space

Algorithms and Data Structures Symposium, 1999

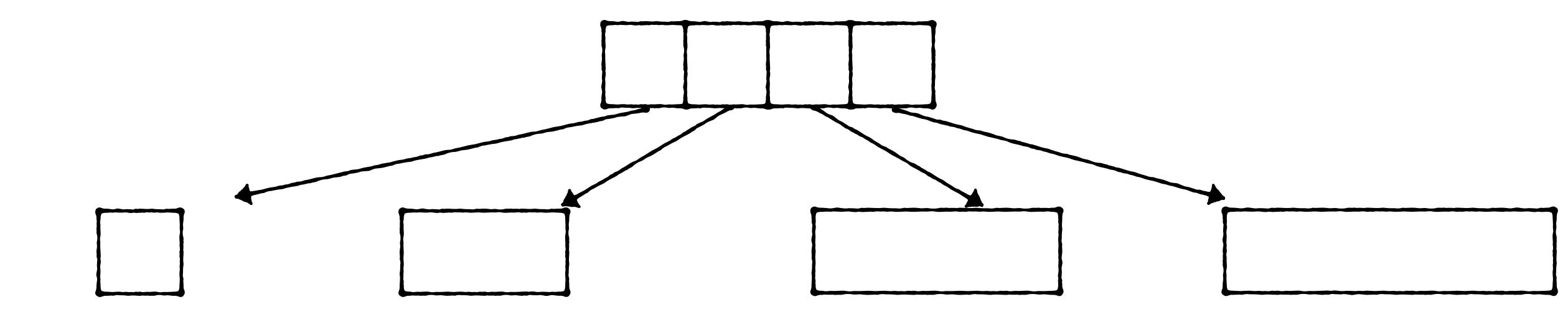


## Resizable Arrays in Optimal Time and Space



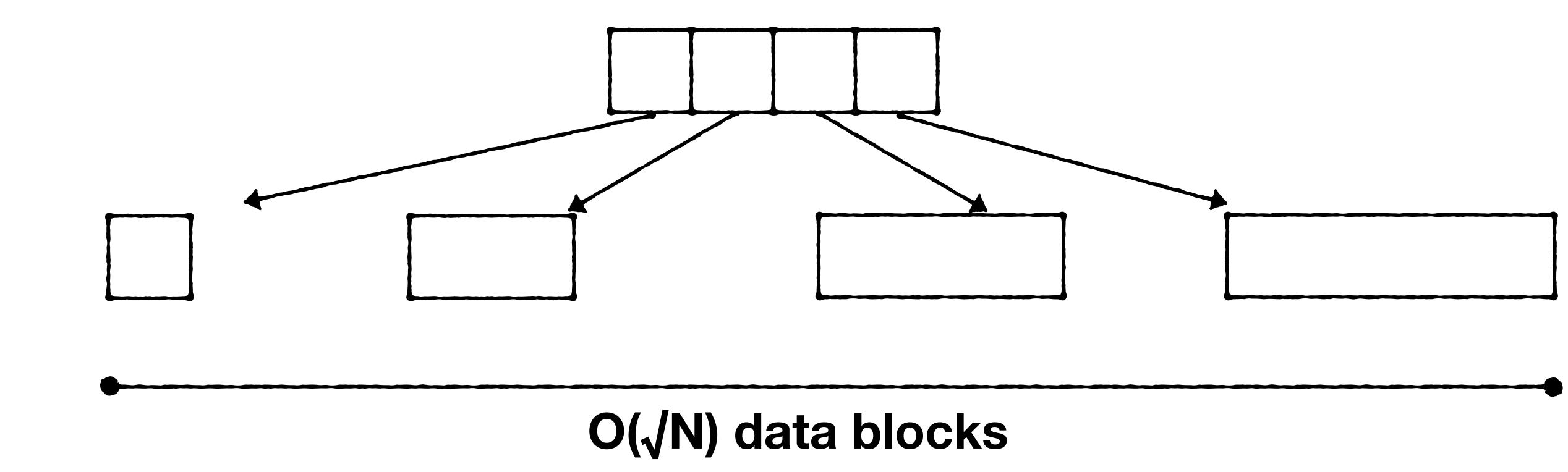
# Data blocks should grow in size

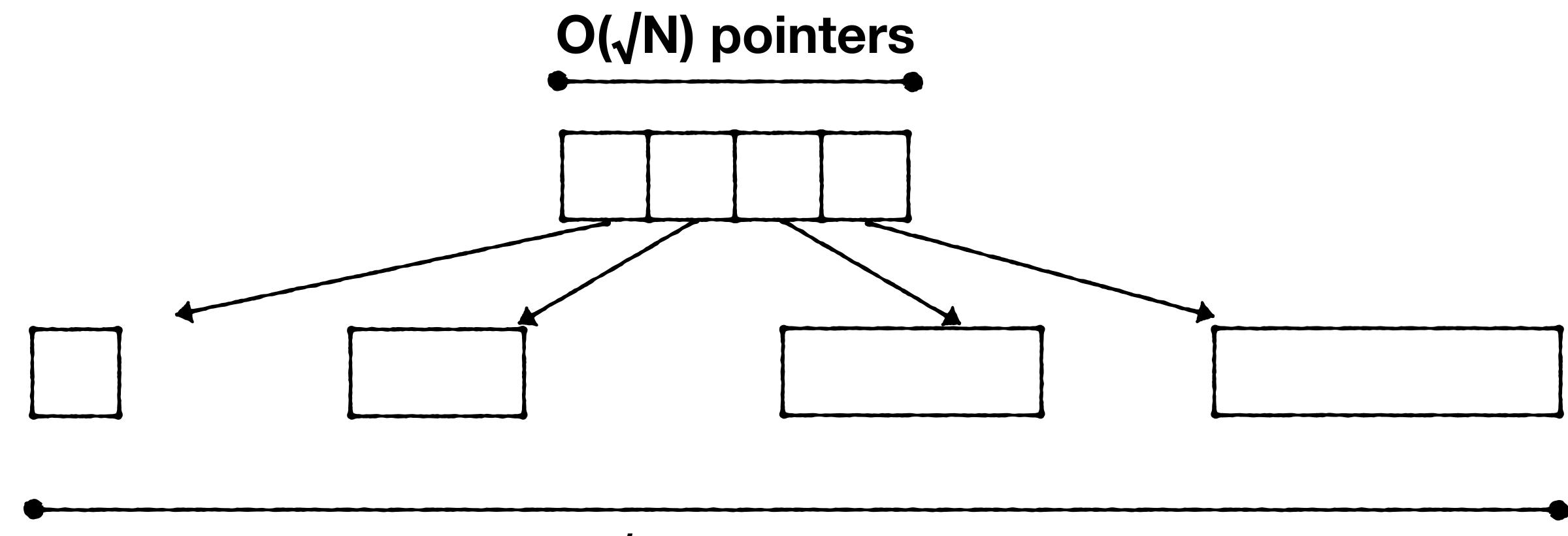
## Resizable Arrays in Optimal Time and Space



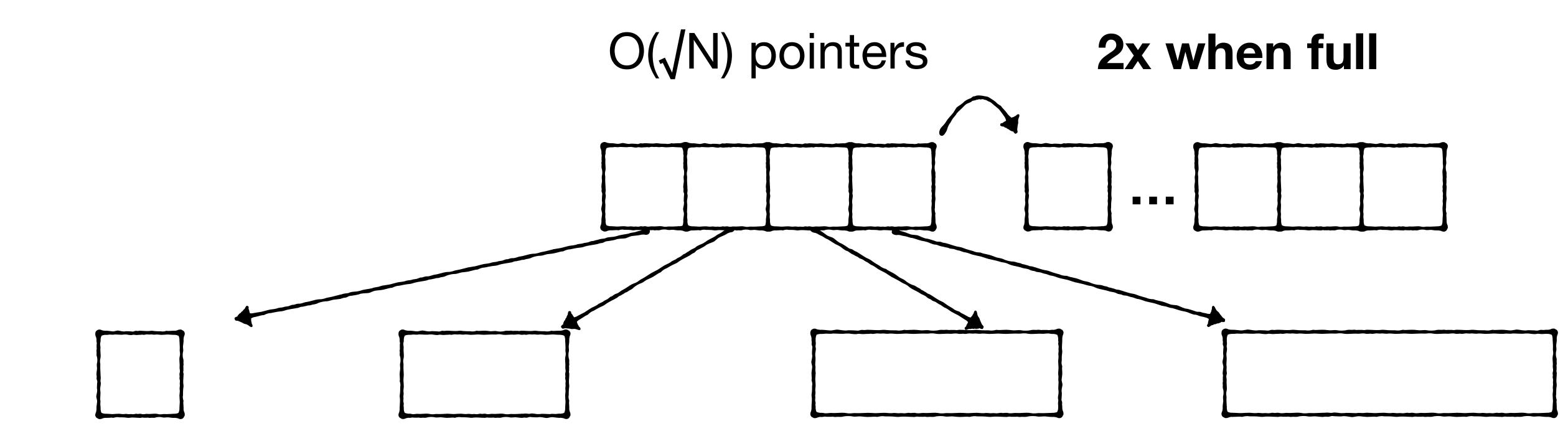
Data blocks should grow in size

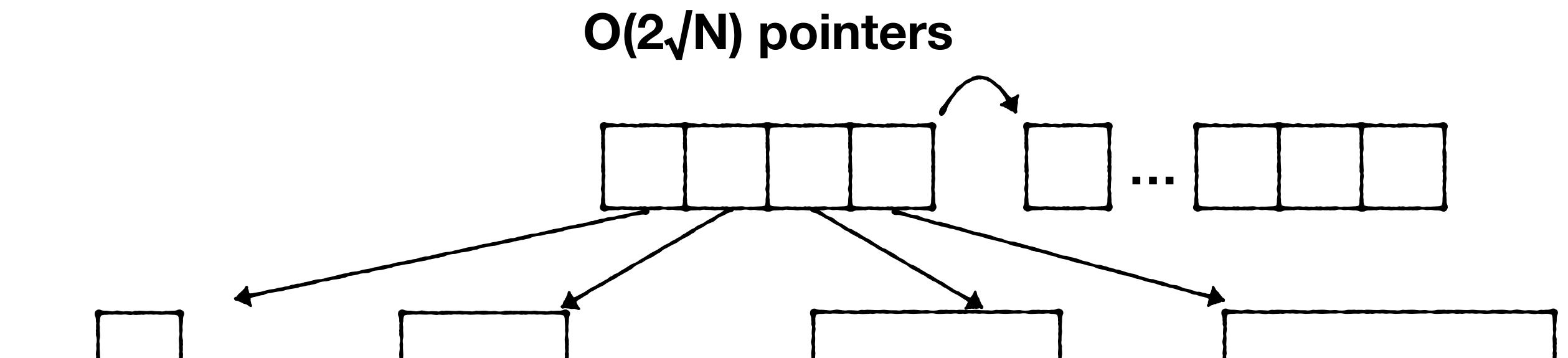
Directory grows more slowly

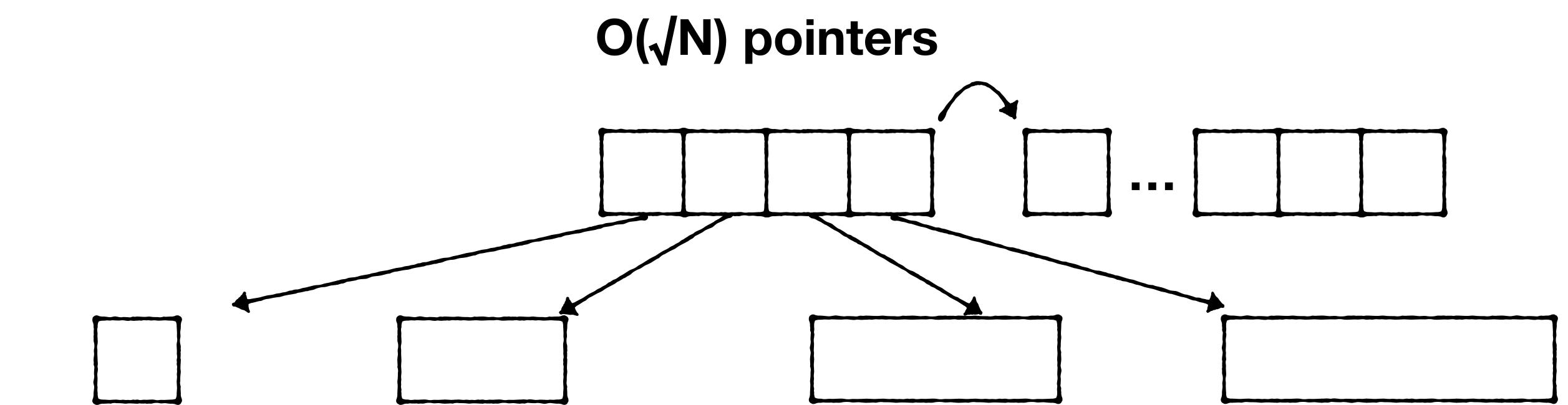




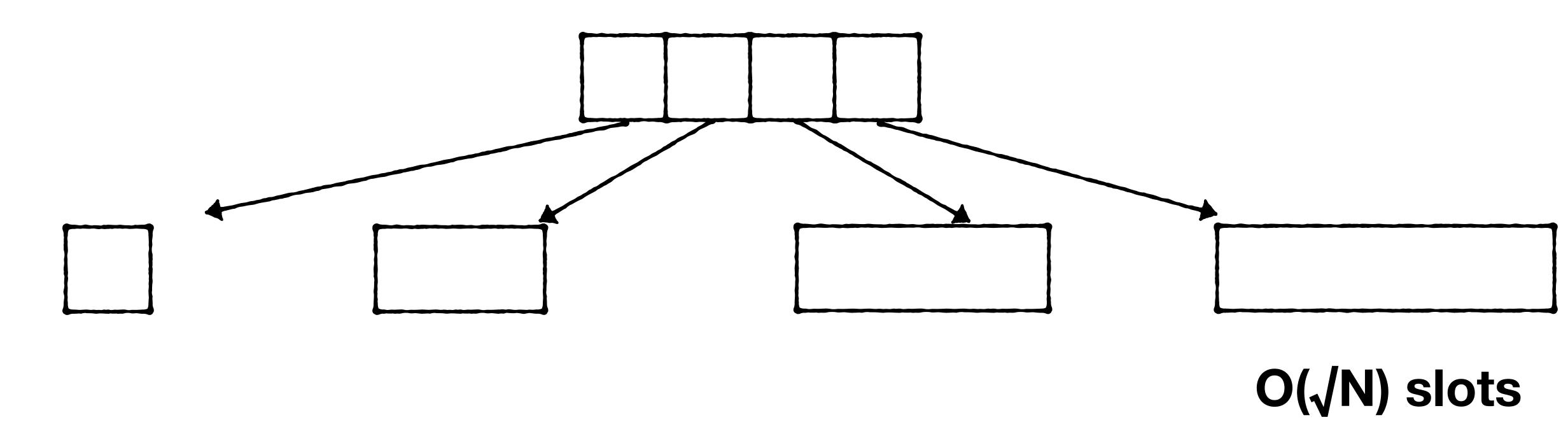
O(\(\sqrt{N}\)) data blocks



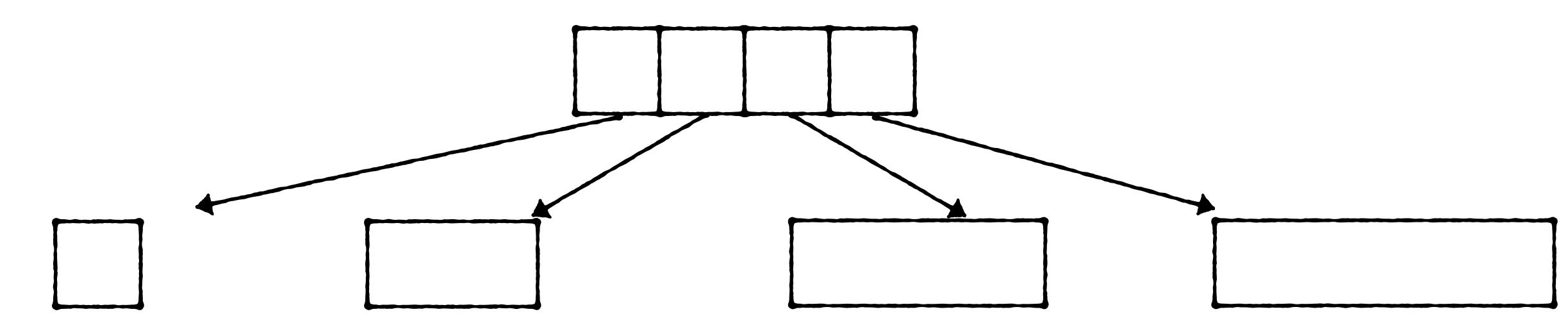




# O(\(\sqrt{N}\)) pointers

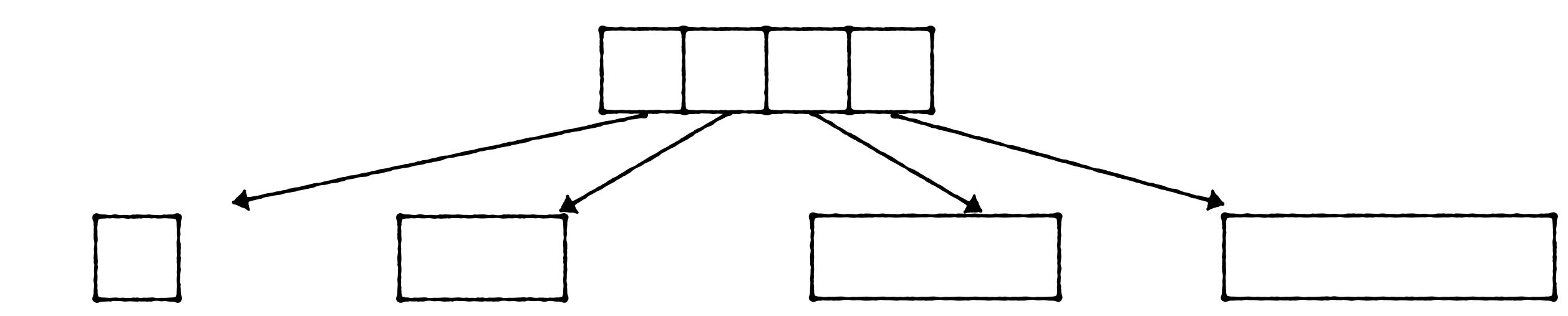


# O(\(\sqrt{N}\)) pointers

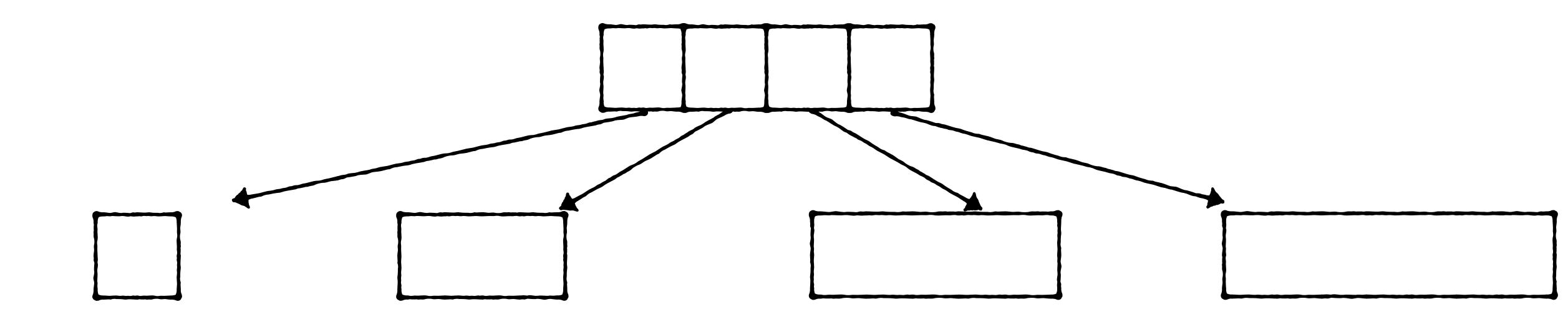


Waste at most O(\langle N) slots

# Max space amp = $O(\sqrt{N}) + O(\sqrt{N}) = O(\sqrt{N})$

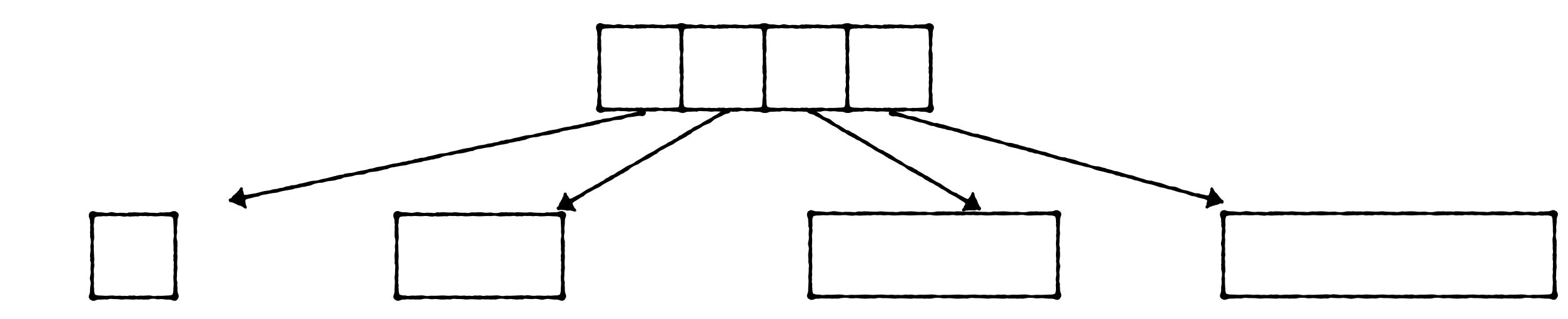


Max space amp =  $O(\sqrt{N})$ 



Challenges: How to grow the block to meet these properties?

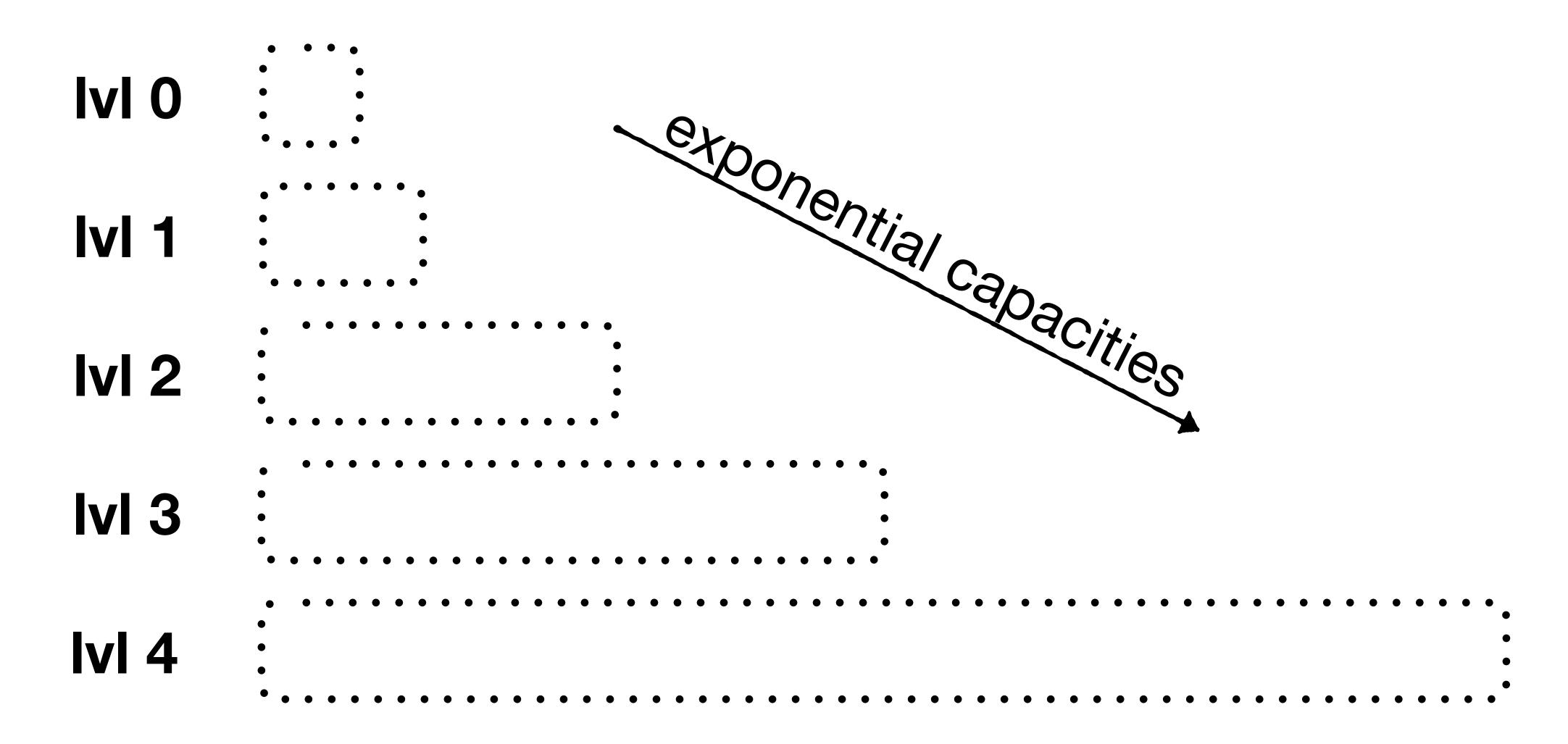
Max space amp =  $O(\sqrt{N})$ 

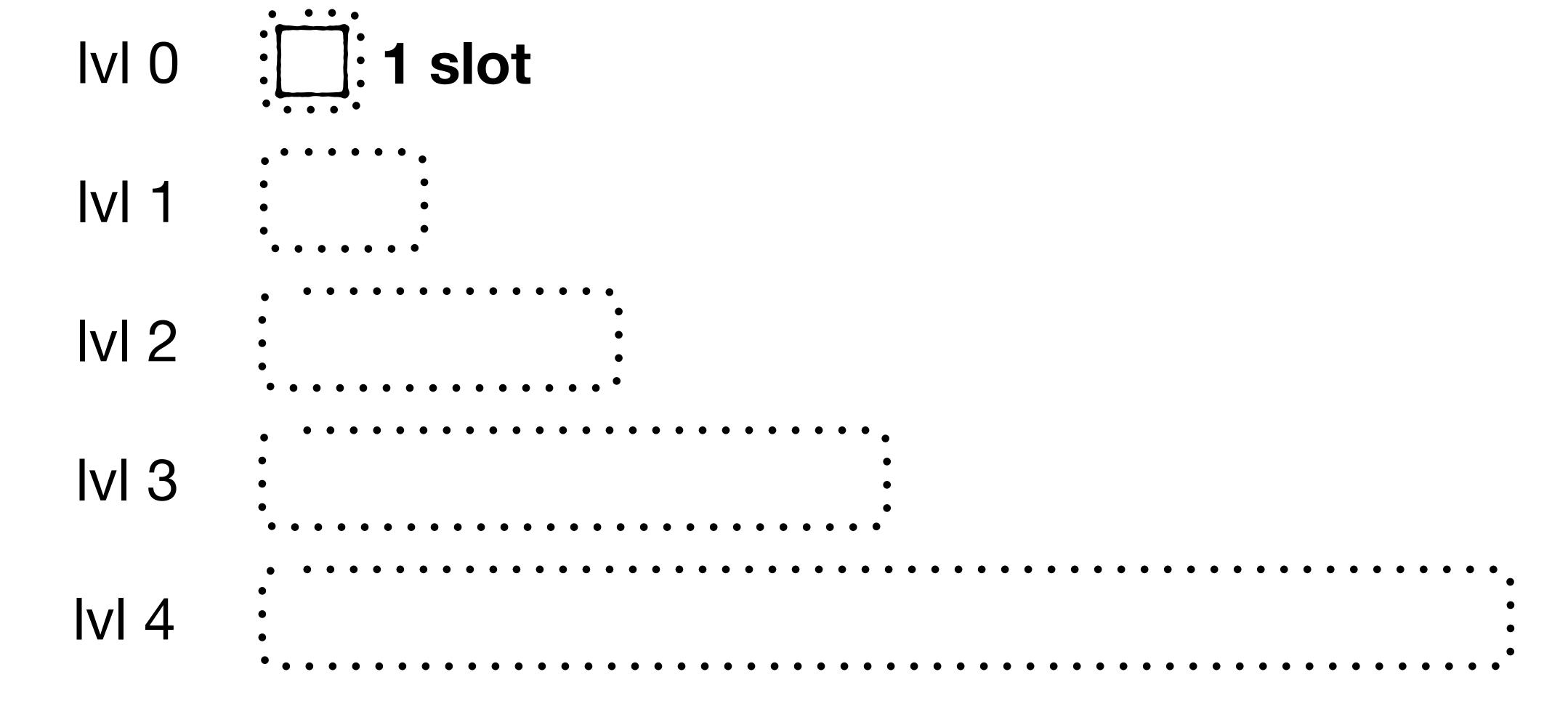


Challenges: How to grow the block to meet these properties?

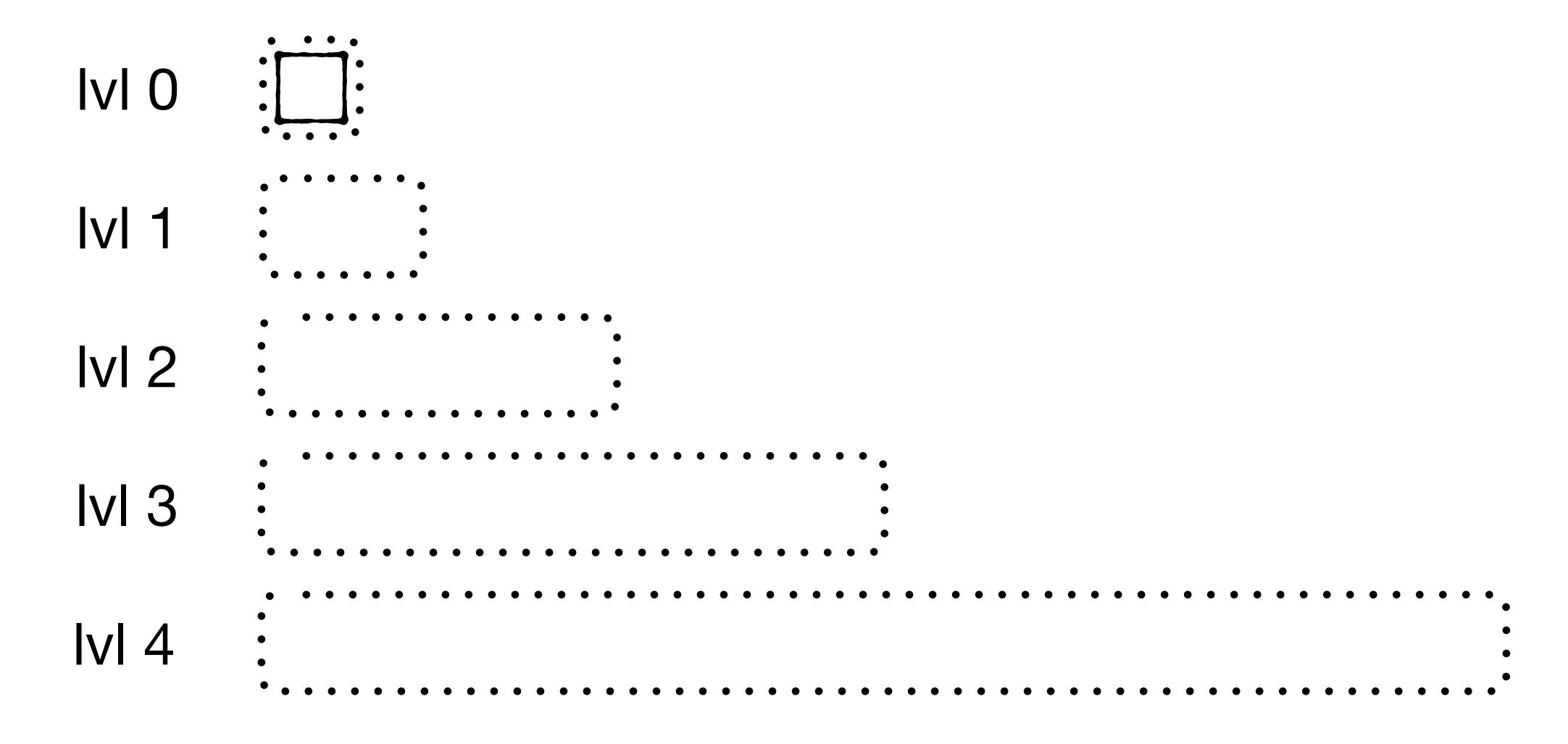
Inferring which block contains which array offset?

## Multiple levels



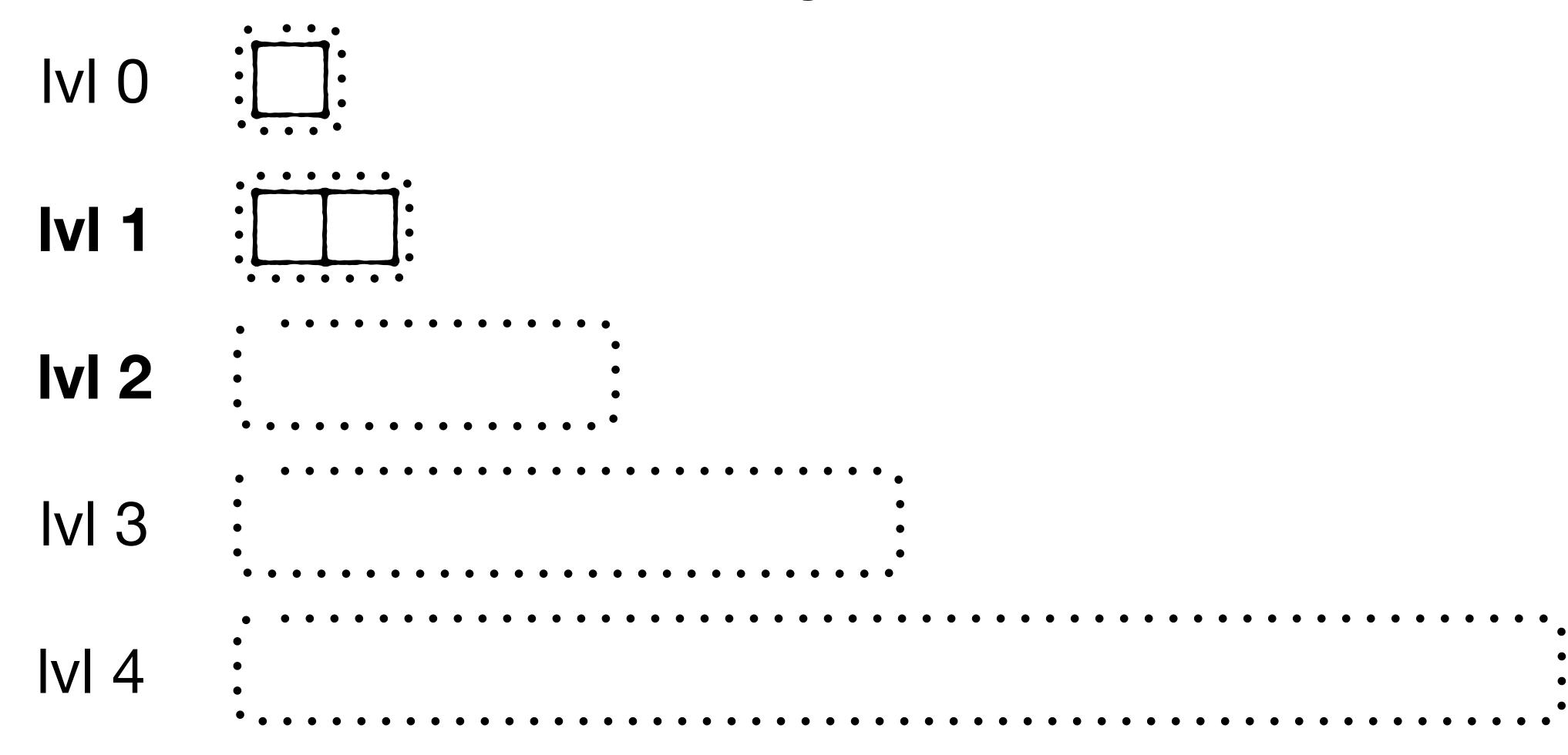


## In every pair of subsequent levels k and k+1

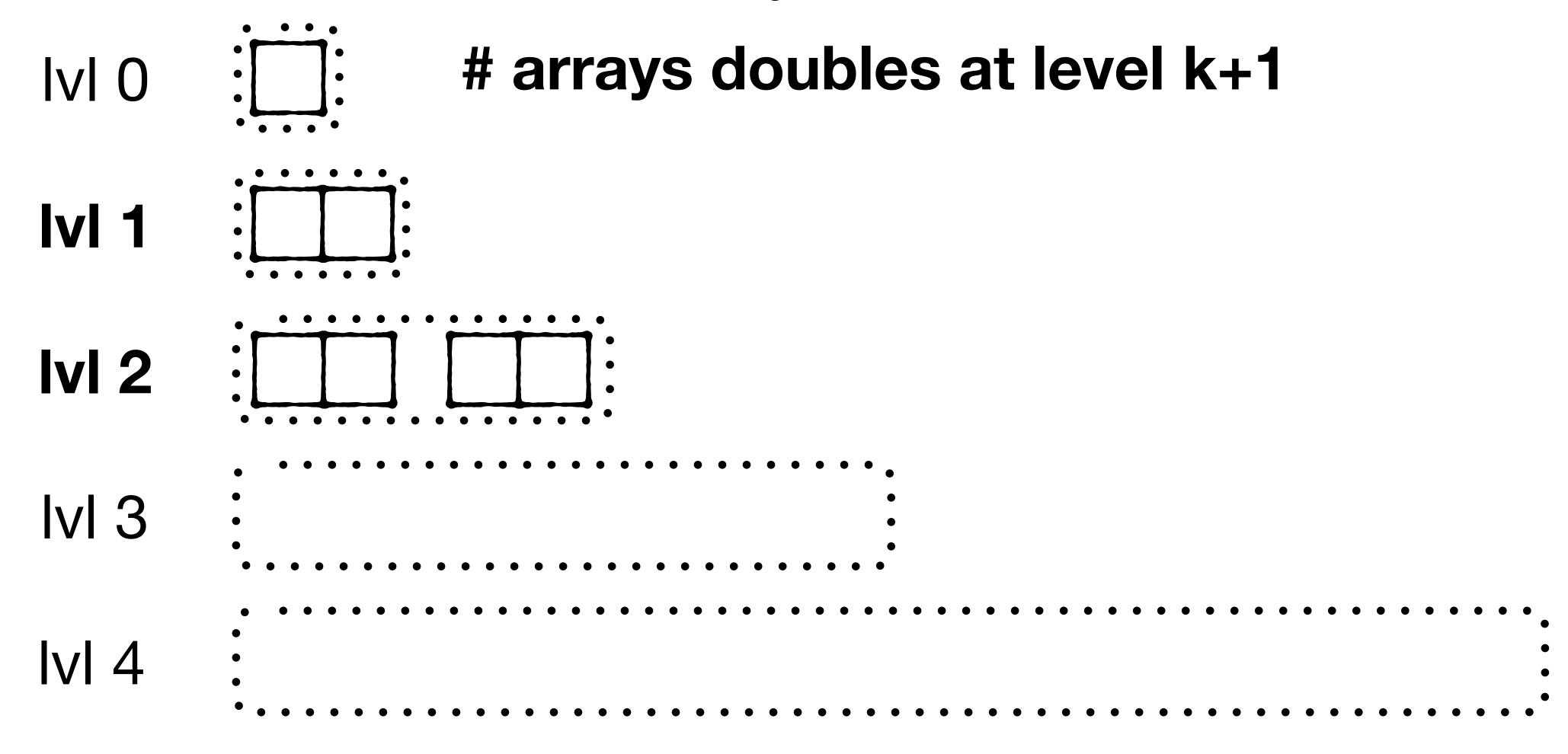


## In every pair of subsequent levels k and k+1

## Size of arrays doubles at level k

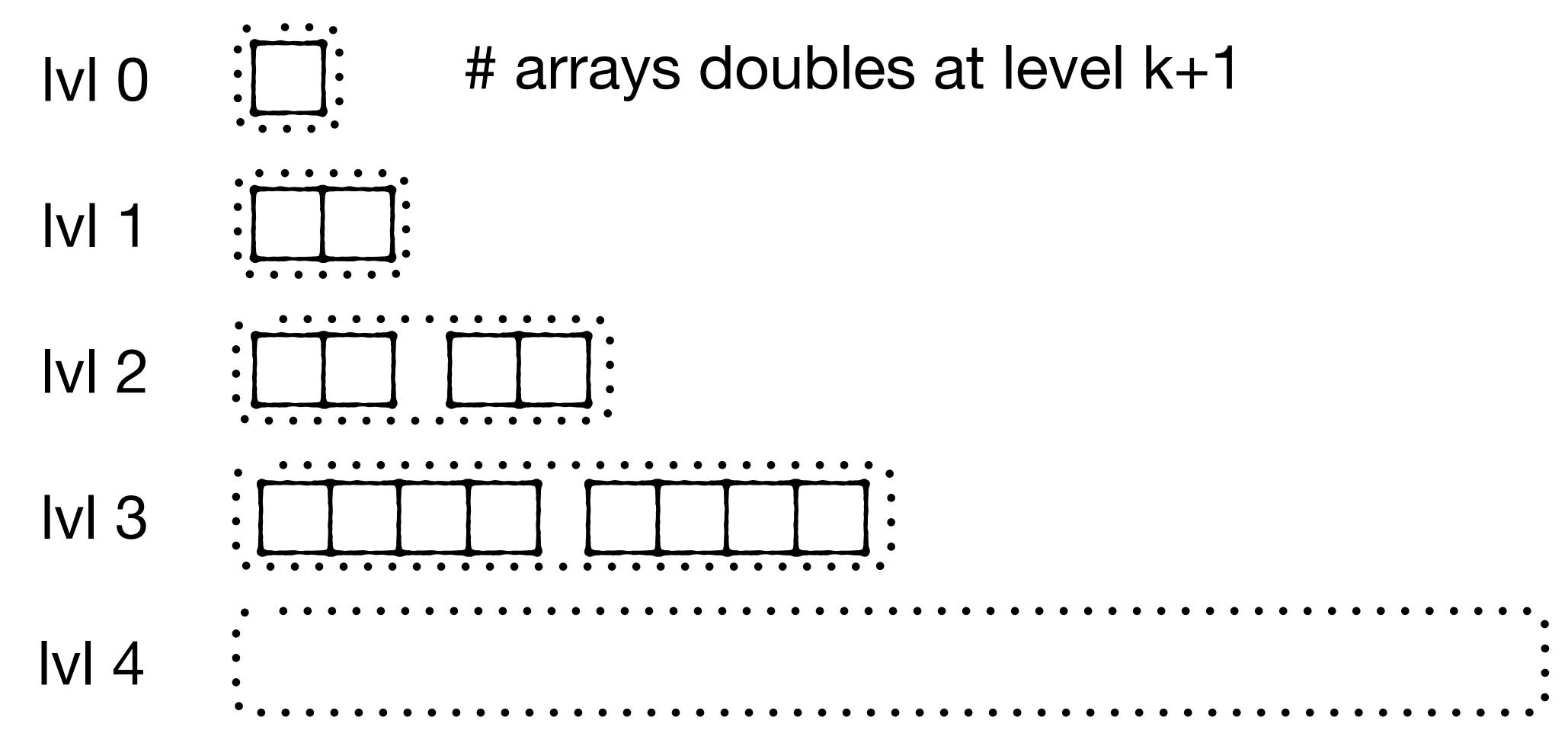


In every pair of subsequent levels k and k+1 Size of arrays doubles at level k

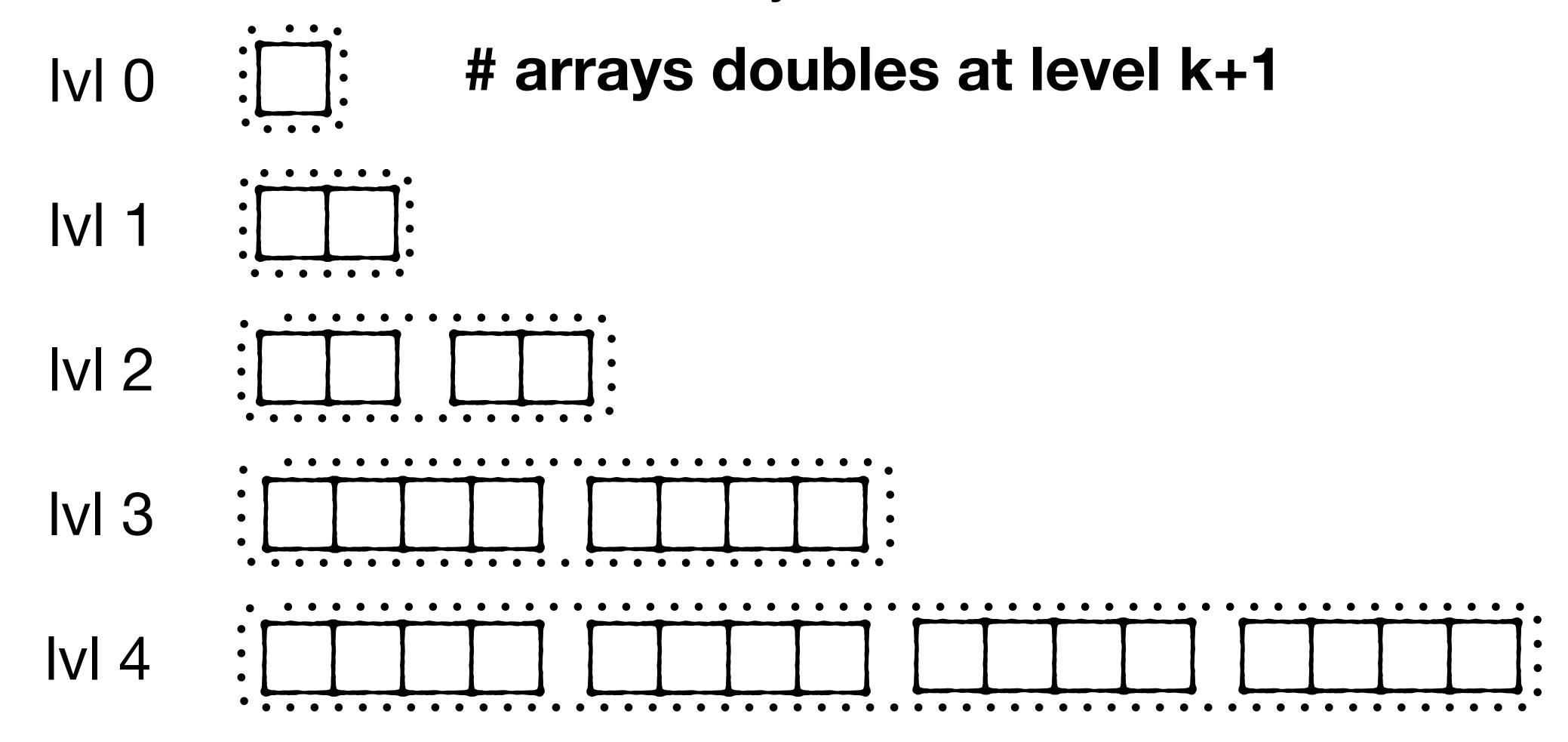


## In every pair of subsequent levels k and k+1

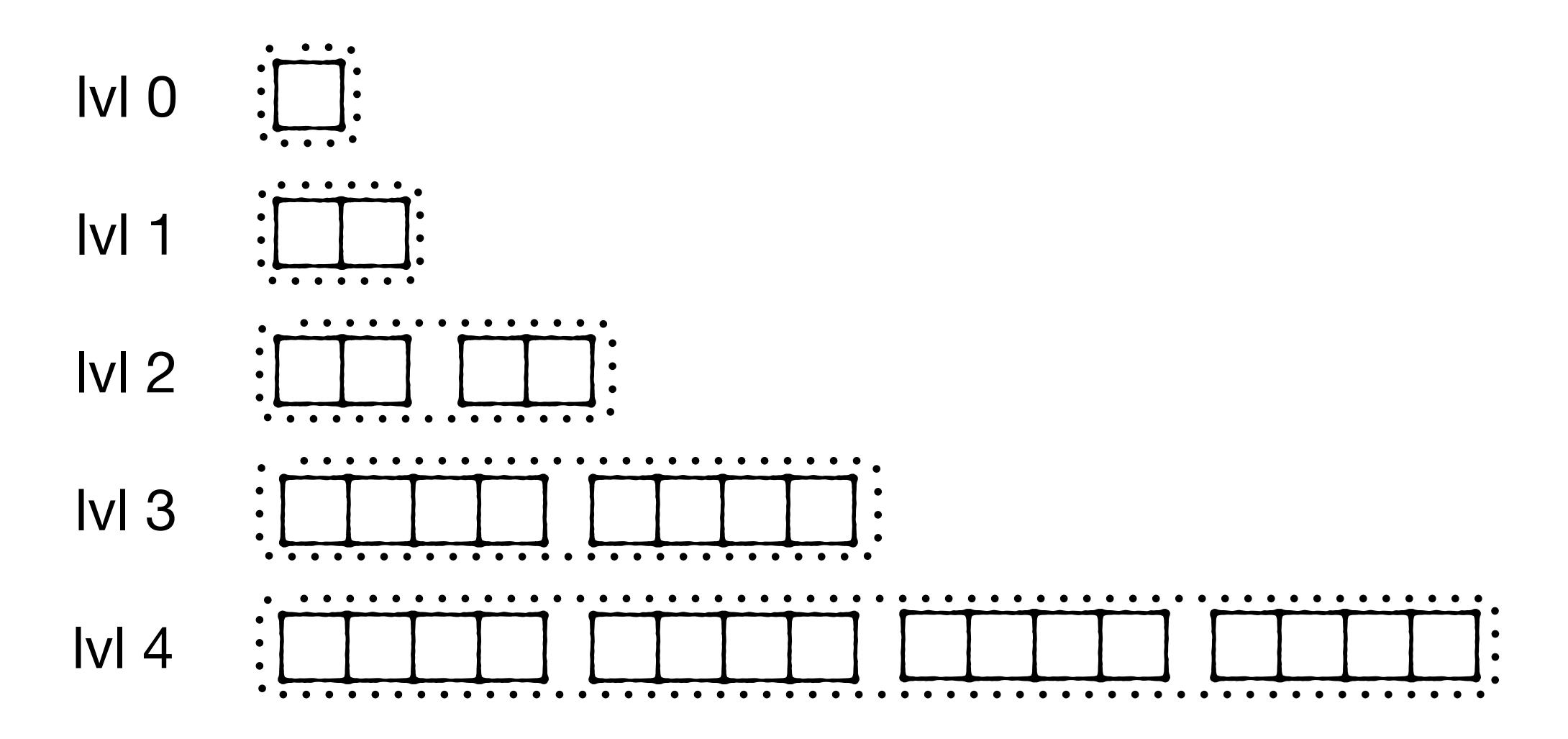
## Size of arrays doubles at level k



In every pair of subsequent levels k and k+1 Size of arrays doubles at level k



## IvI i contains $2^{\lfloor K/2 \rfloor}$ blocks, each with $2^{\lceil K/2 \rceil}$ slots



## IvI i contains $2^{\lfloor K/2 \rfloor}$ blocks, each with $2^{\lceil K/2 \rceil}$ slots

$$|V| 0 \qquad \boxed{ } \qquad \boxed{ } \qquad 2^{\lfloor 0/2 \rfloor} = 1$$

$$|V| 1 \qquad \boxed{ } \qquad \boxed{ } \qquad 2^{\lfloor 1/2 \rfloor} = 1$$

$$|V| 2 \qquad \boxed{ } \qquad \boxed{ } \qquad \boxed{ } \qquad 2^{\lfloor 2/2 \rfloor} = 2$$

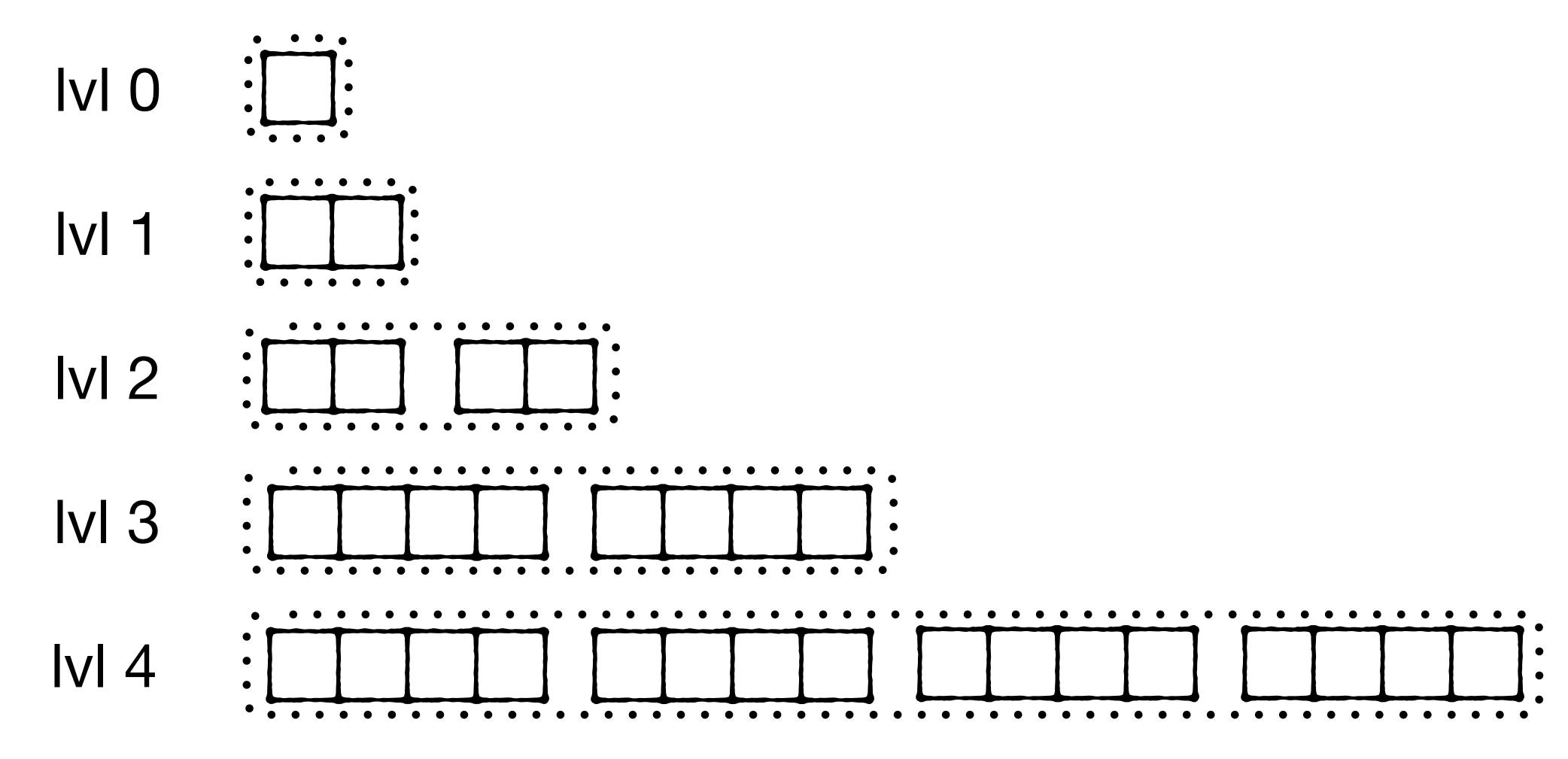
$$|V| 3 \qquad \boxed{ } \qquad \boxed{ } \qquad \boxed{ } \qquad \boxed{ } \qquad 2^{\lfloor 3/2 \rfloor} = 2$$

$$|V| 4 \qquad \boxed{ } \qquad$$

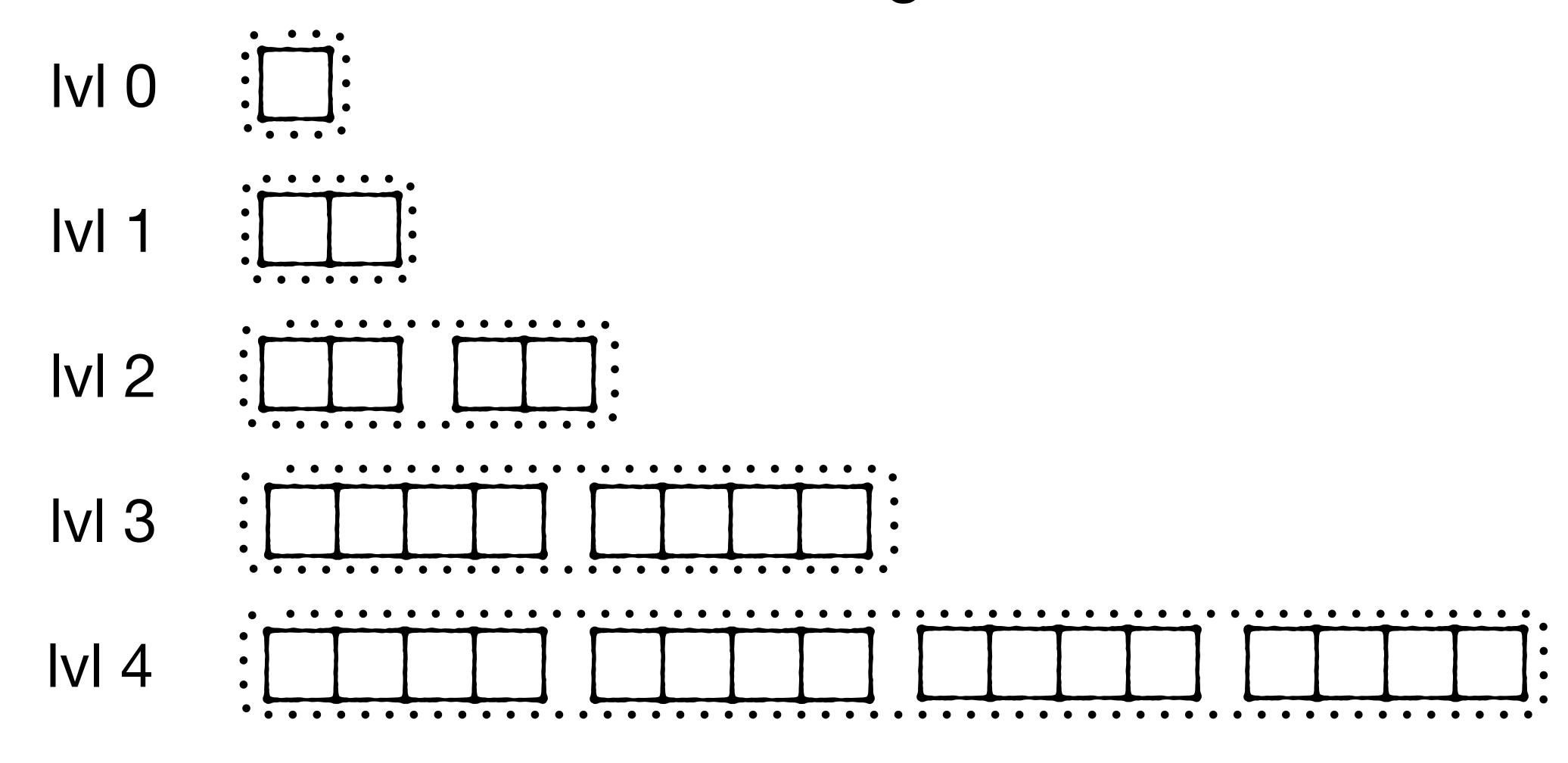
## IvI i contains $2^{\lfloor K/2 \rfloor}$ blocks, each with $2^{\lceil K/2 \rceil}$ slots

| 
$$| V | 0$$
 |  $| 2^{\lceil 0/2 \rceil} = 1$  |  $| V | 1$  |  $| 2^{\lceil 1/2 \rceil} = 2$  |  $| V | 2$  |  $| D | | 2^{\lceil 2/2 \rceil} = 2$  |  $| V | 3$  |  $| D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D | | D$ 

#### # levels?



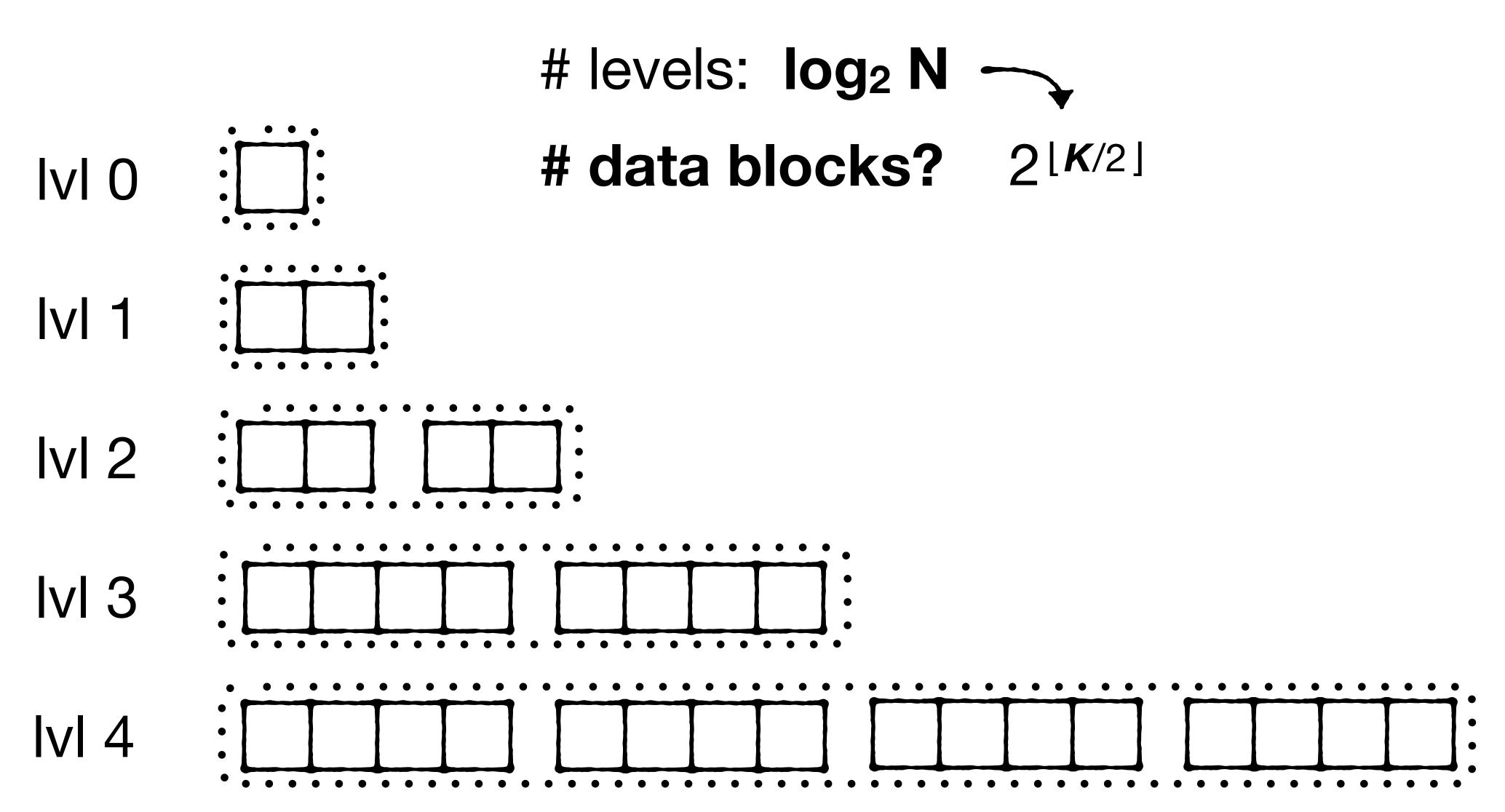
## # levels: log<sub>2</sub> N



IVI 3 : IIII IIII

IVI 4 EMPLIA EMP

# levels: log<sub>2</sub> N # data blocks? Most are here IVI O IVI 1 IVI 2 IVI 3

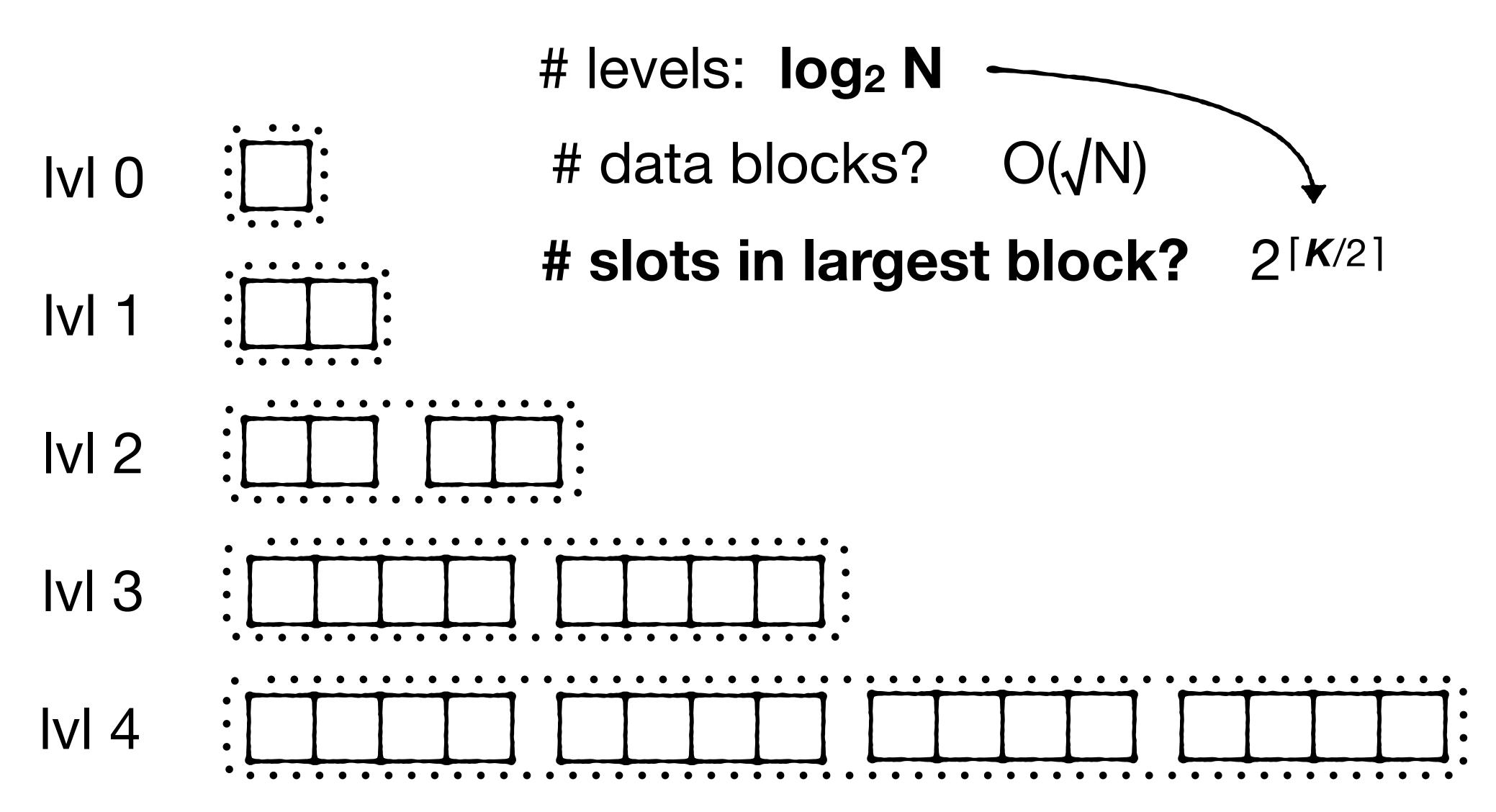


```
# levels: log<sub>2</sub> N —
                       # data blocks? 2 log N/2
IVI O
IVI 1
IVI 2
IVI 3
```

# levels: log<sub>2</sub> N # data blocks? O(\(\sqrt{N}\) IVI 0 IVI 1 IVI 2

IVI 4 EMPLIANTE EN LA ENTRE DE LA ENTRE DELLE DE LA ENTRE DEL ENTRE DE LA ENTRE DELETA DE LA ENTRE DE

# levels: log<sub>2</sub> N # data blocks? IVI O # slots in largest block? IVI 1 IVI 2 IVI 3

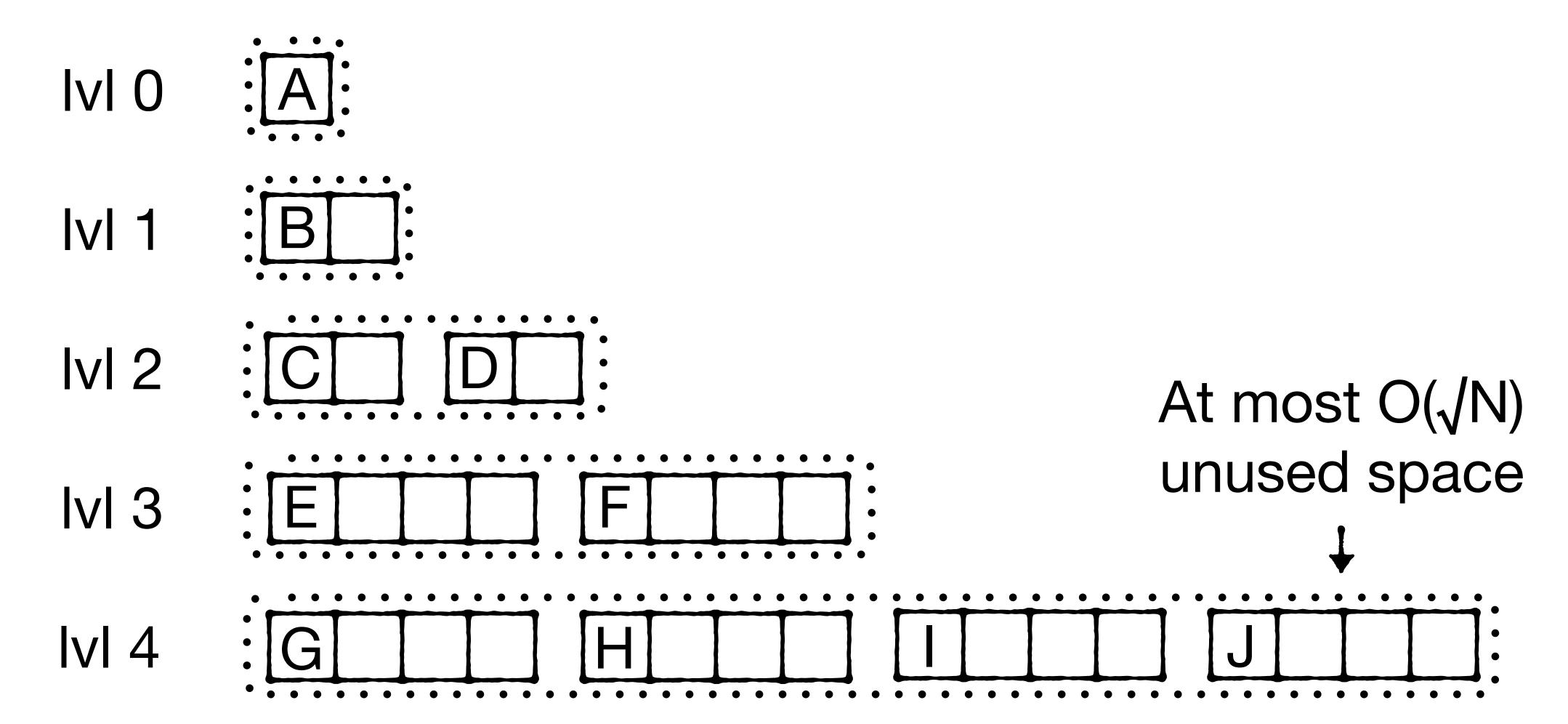


# levels: log<sub>2</sub> N # data blocks? O(√N) IVI O # slots in largest block? O(\lambda N) IVI 1 IVI 2 IVI 3

# levels: log<sub>2</sub> N # data blocks? IVI O IVI 1 IVI 2 At most O(\( \lambda N \right) unused space IVI 3

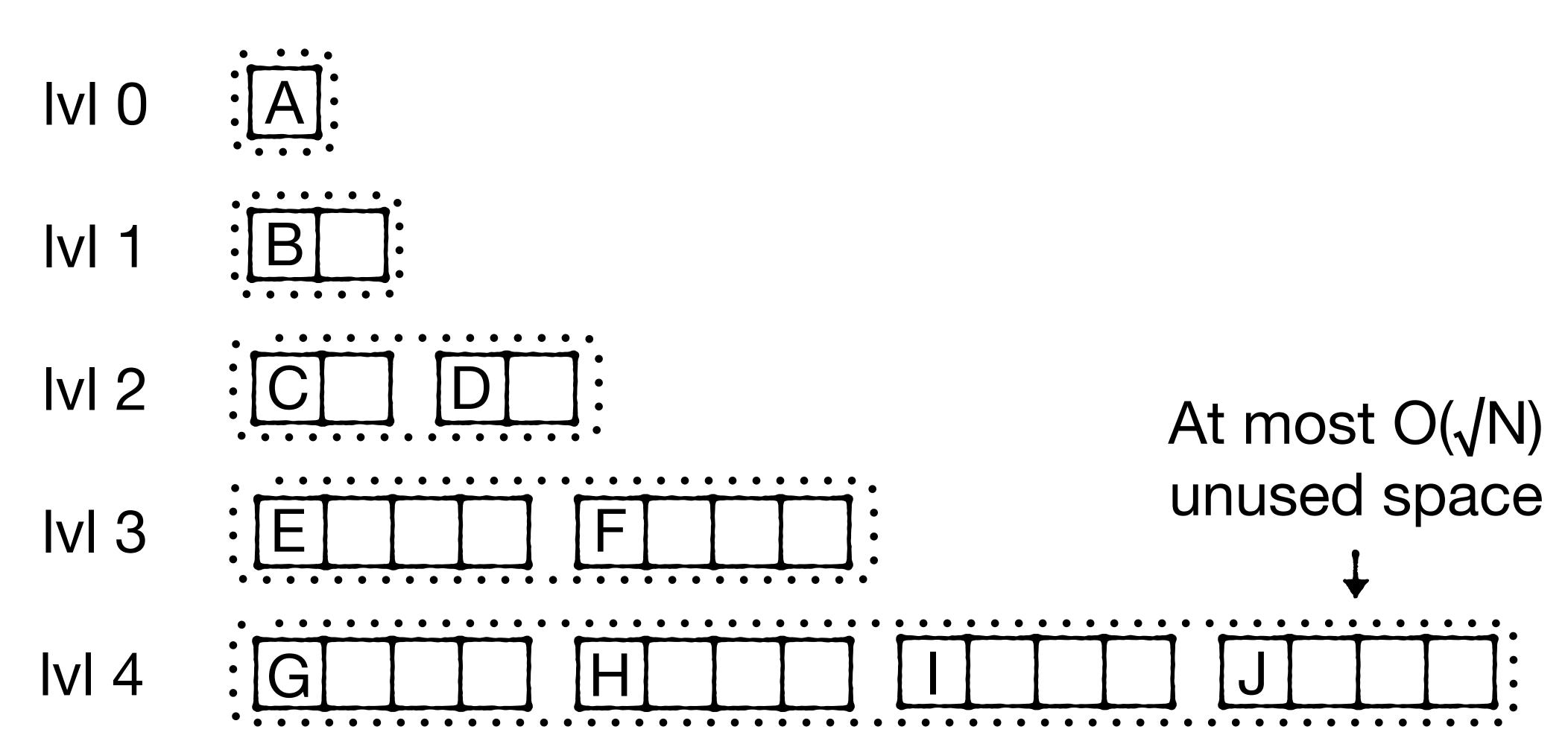
## Directory with O(\lambda N) pointers

ABCDEFGHIJ



## At most half O(\(\sqrt{N}\)) unused space





## ABCDEFGHIJ...

IVI O : A:

IVI 1 B

Max extra space:  $O(\sqrt{N}) + O(\sqrt{N}) = O(\sqrt{N})$ 

IVI 2 : C D

## ABCDEFGHIJ...

IVI O : A:

IVI 1 : B

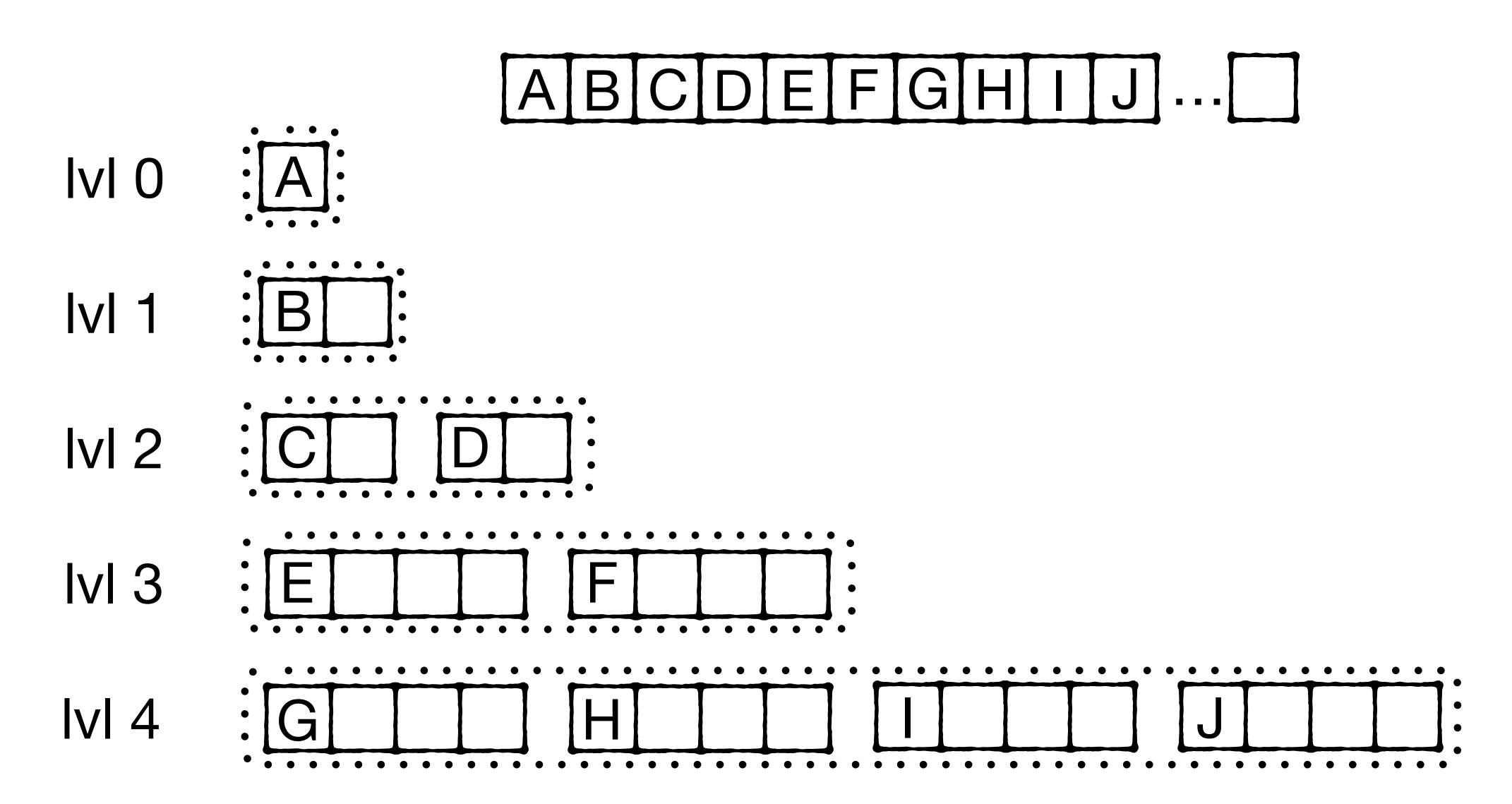
Max space-amp: =  $O(1+N^{-0.5})$ 

IVI 2 : C D

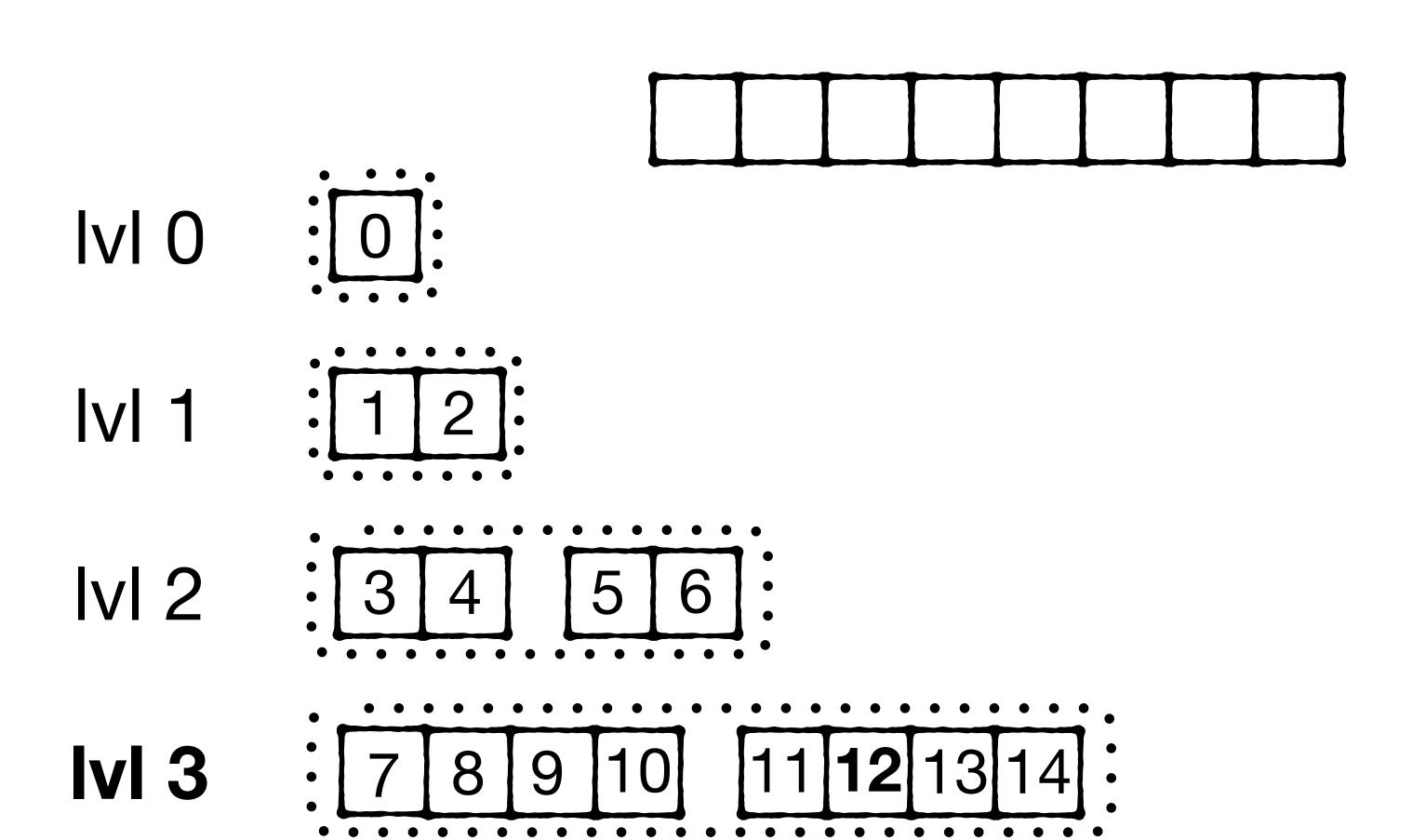
IVI 3 : E F F

IVI 4 : G H H J J J J J J

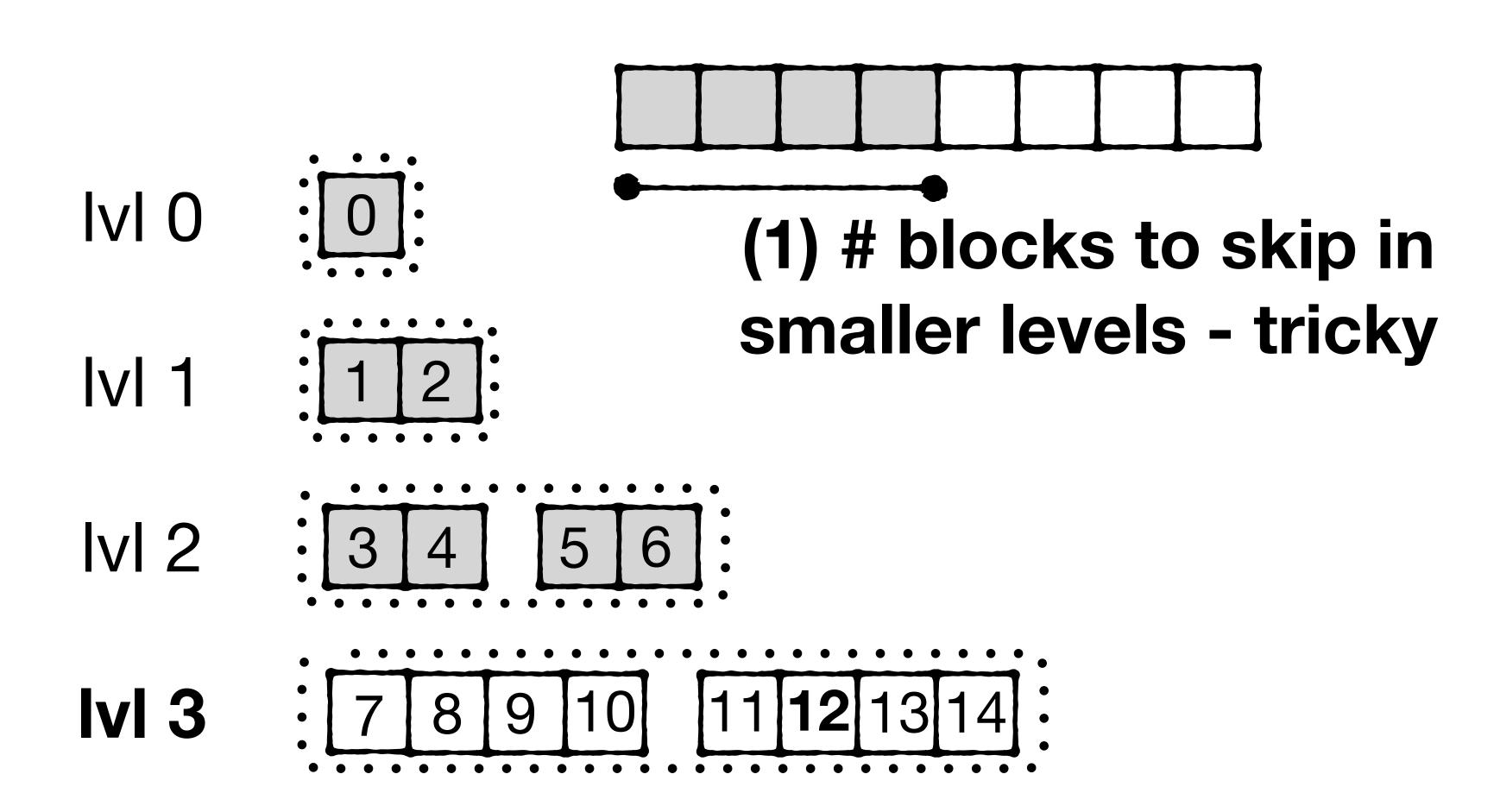
#### How to access slot in O(1) time?



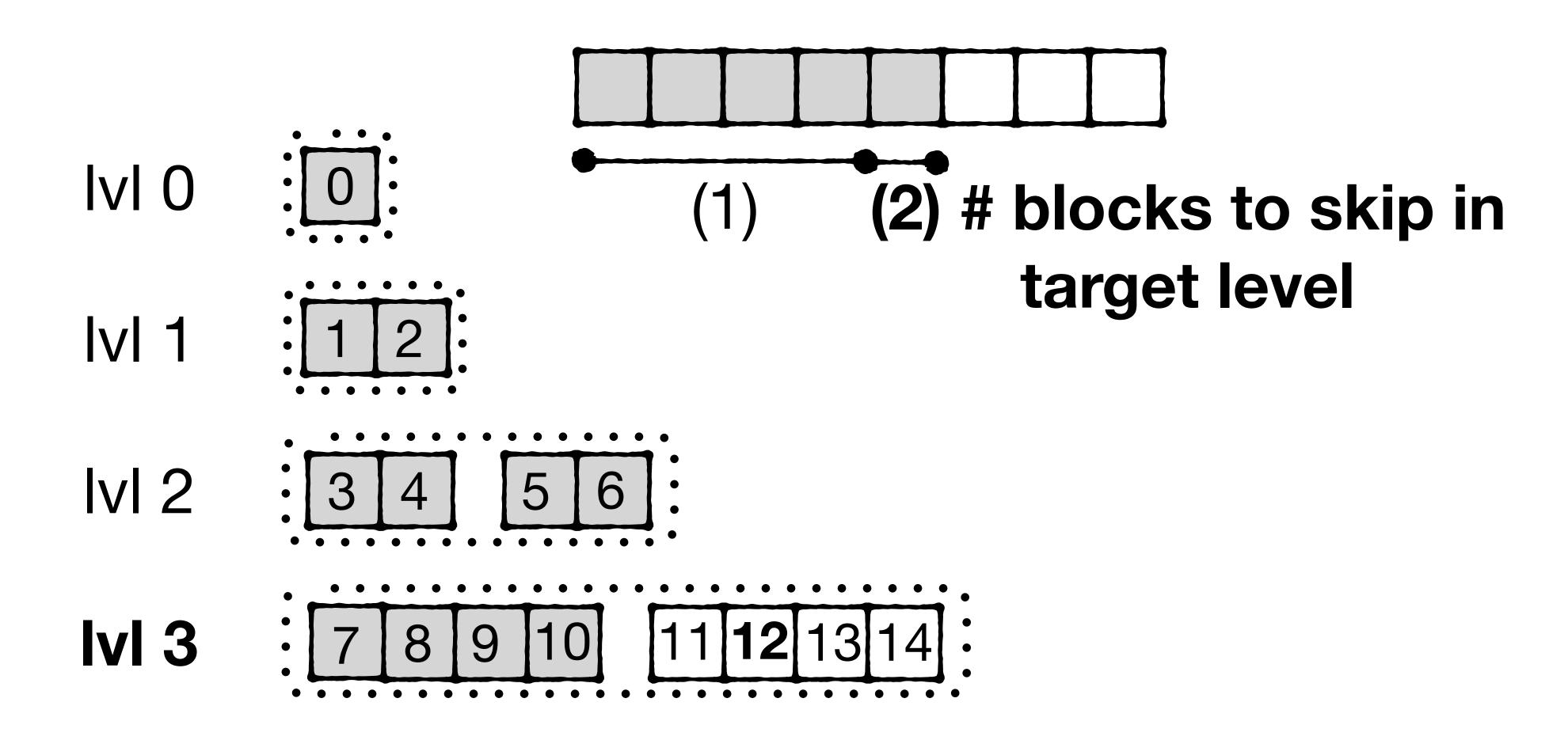
get(12)



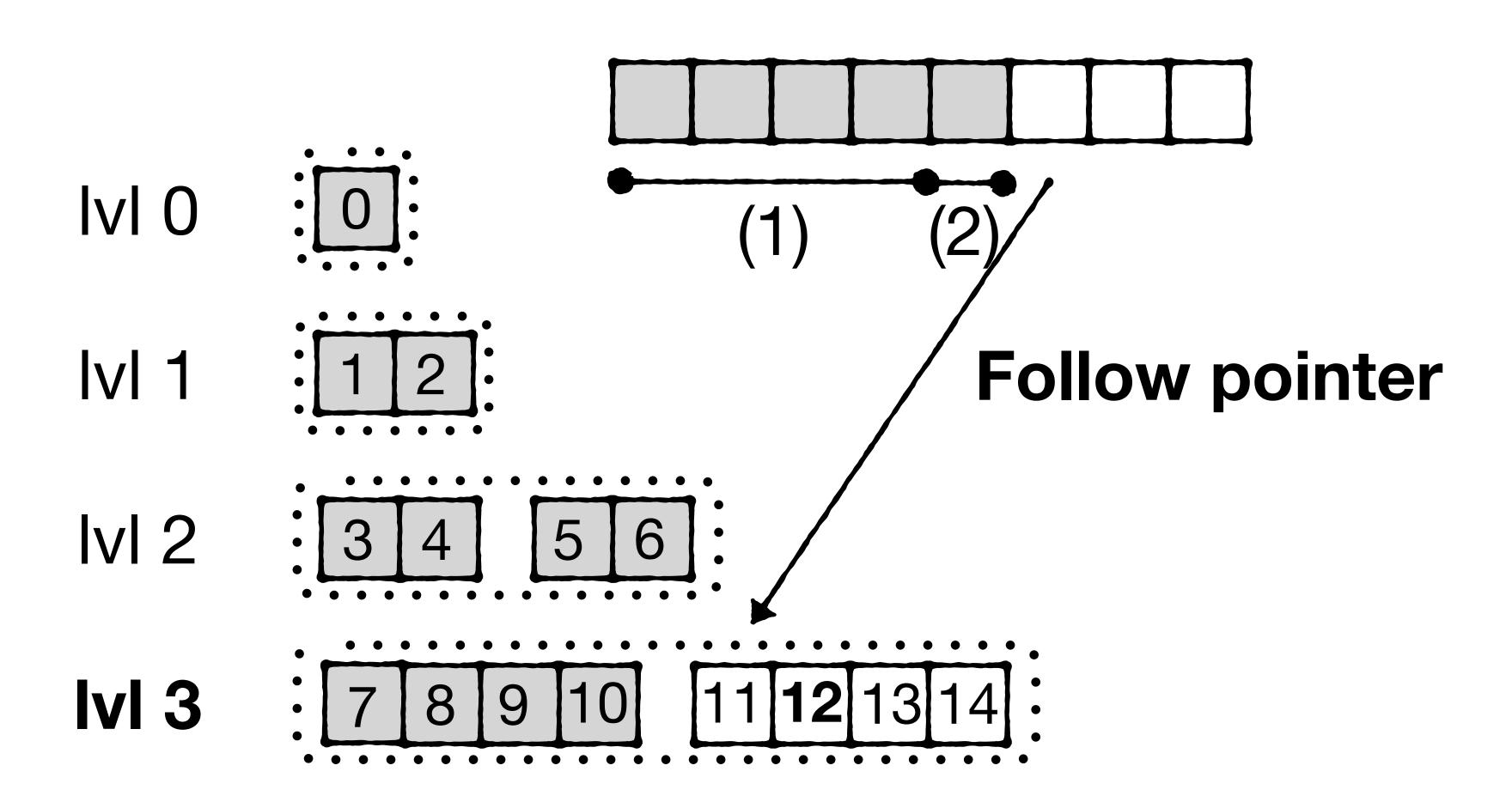
get(12)

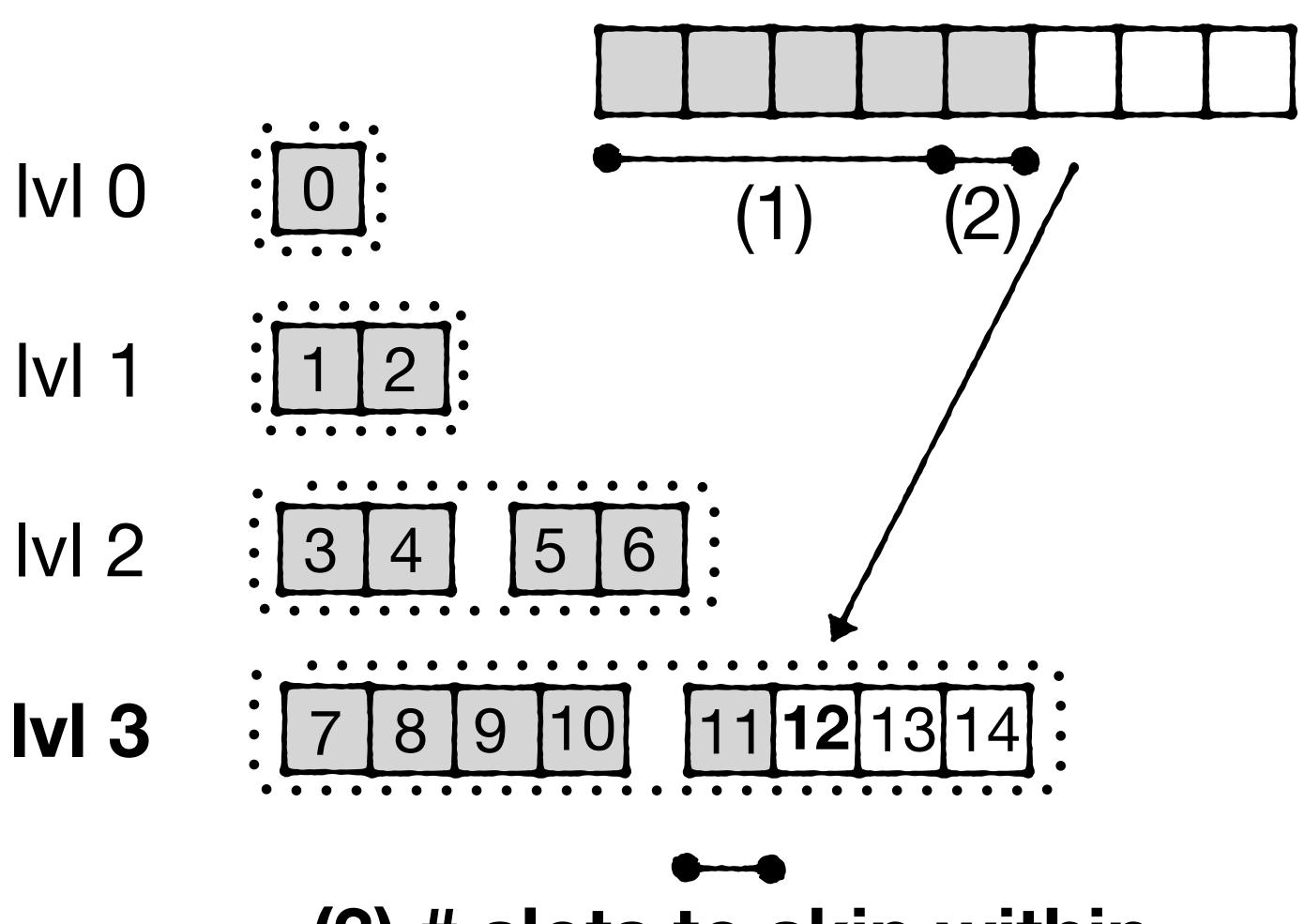


get(12)

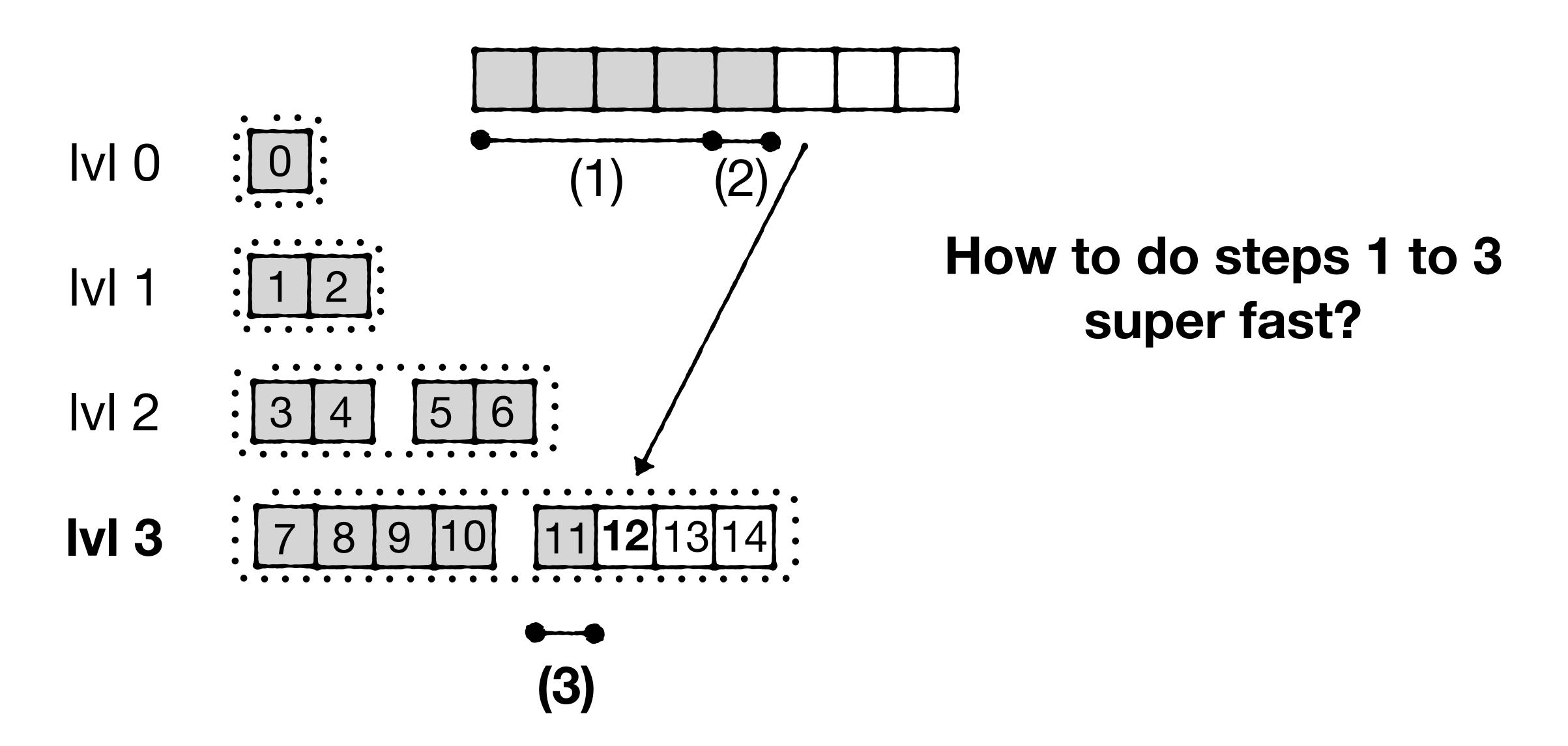


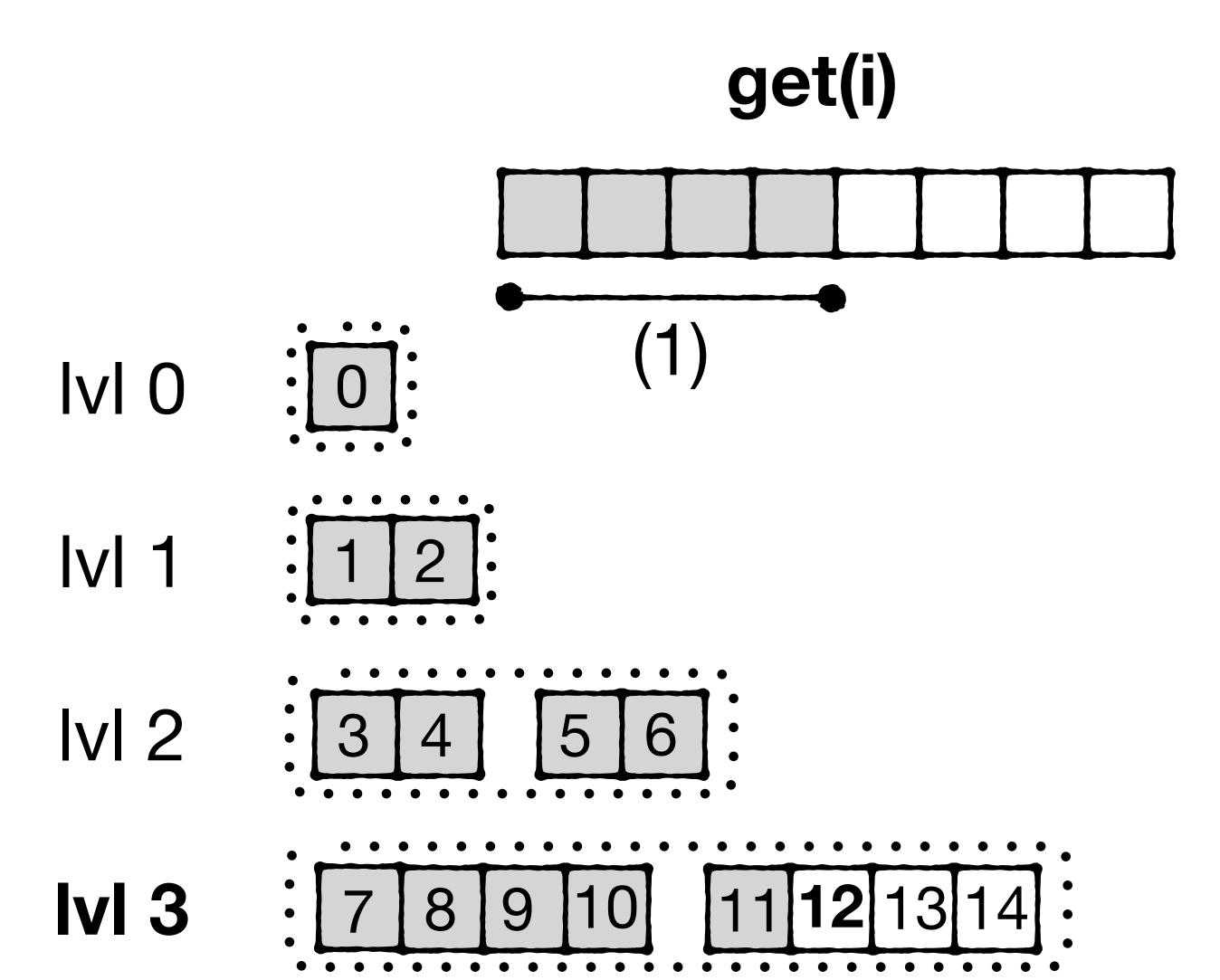
get(12)

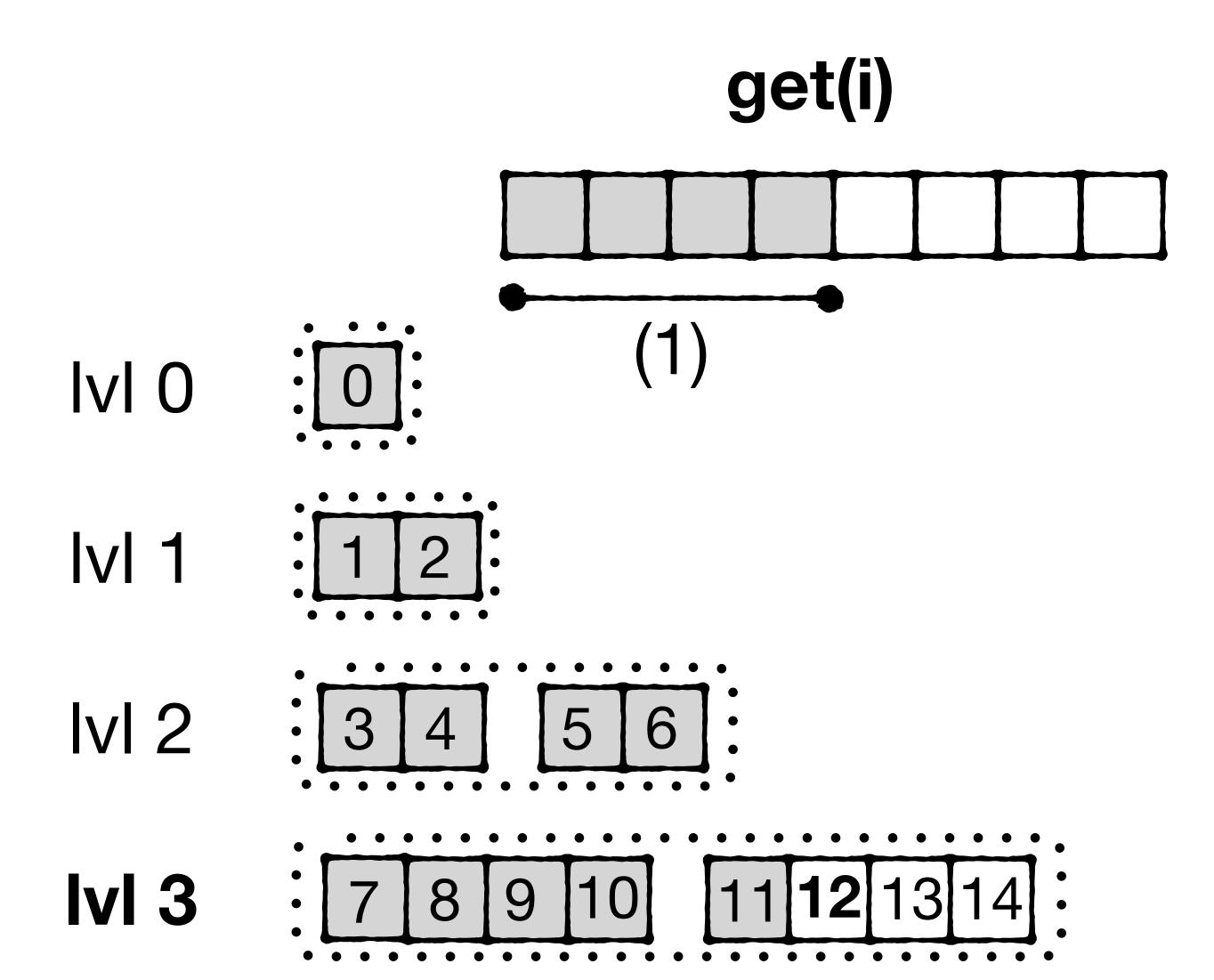




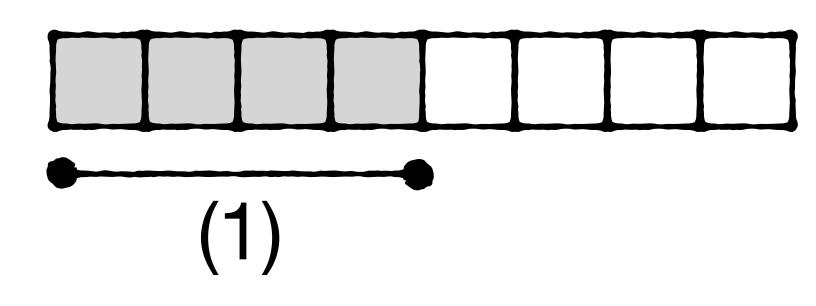
(3) # slots to skip within target block







## get(i)



# Identify target level k $k = \lfloor \log_2(i+1) \rfloor$

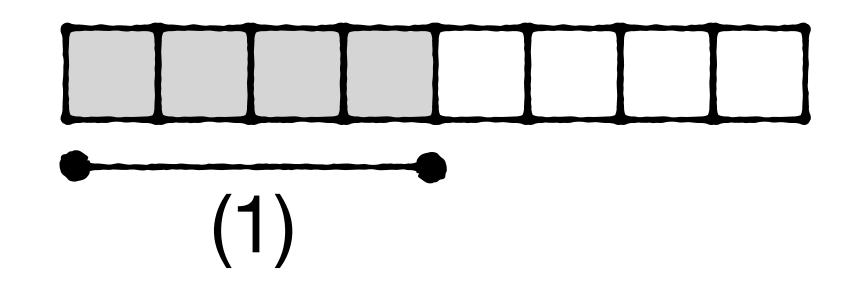
IVI 1 : 1 2

IVI 0

IVI 2 : 3 4 5 6

IVI 3 7 8 9 10 11 12 13 14

## get(12)



Identify target level k  

$$k = \lfloor \log_2(12 + 1) \rfloor = 3$$

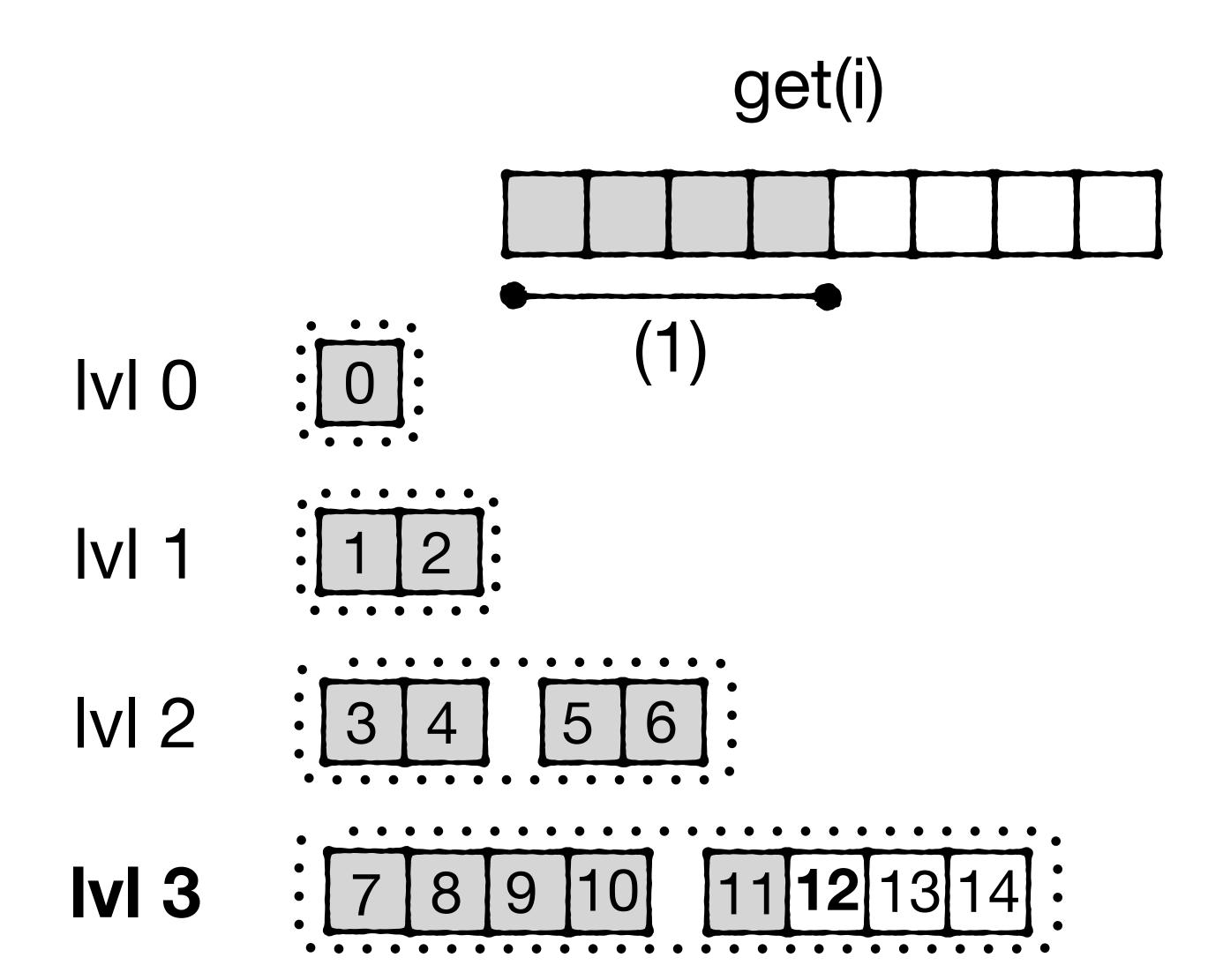
IVI 0

# get(i) IVI 0 IVI 1 IVI 2 IVI 3

Identify target level k

$$k = \lfloor \log_2(i + 1) \rfloor$$

$$\uparrow$$
slow



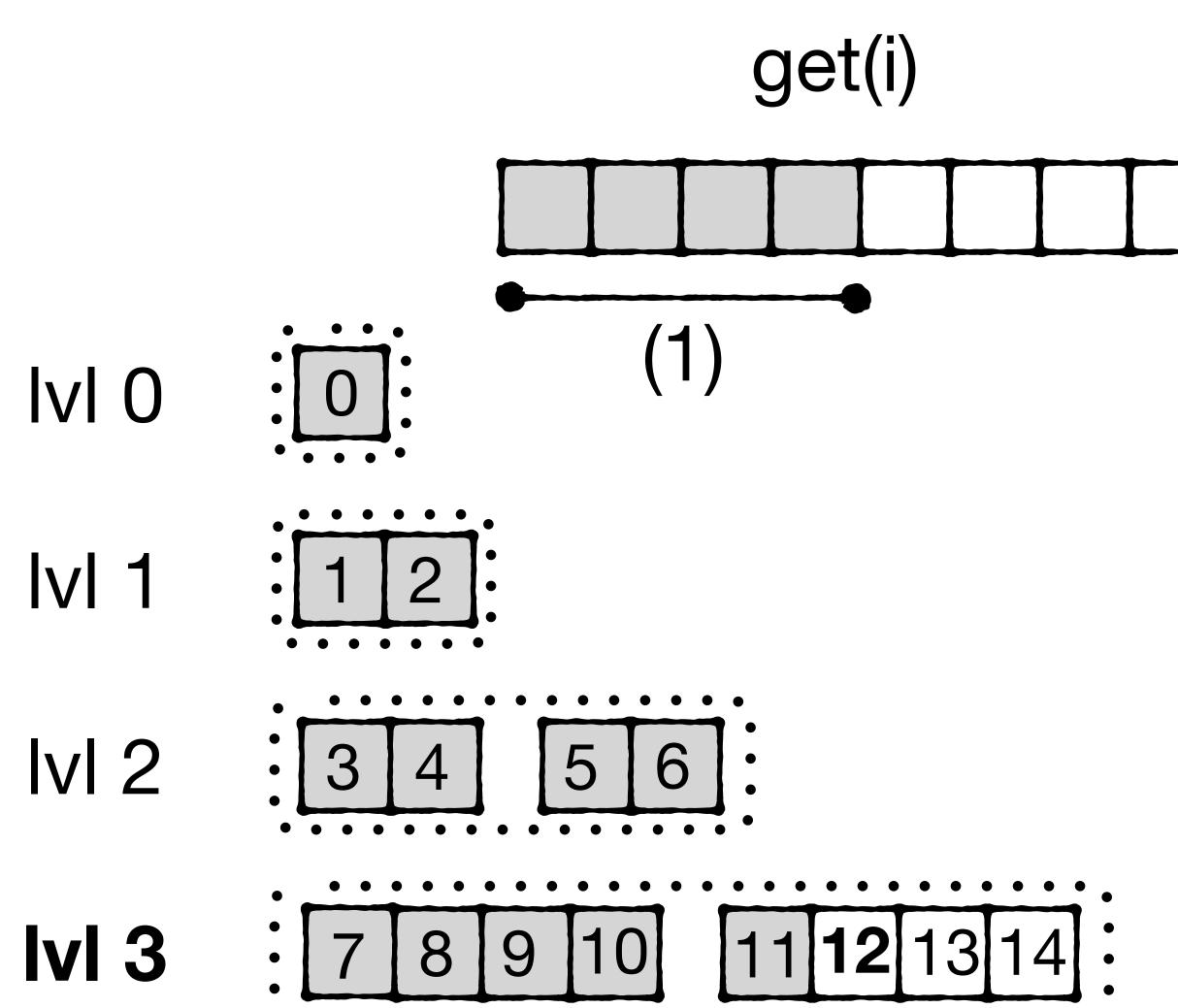
$$k = \lfloor \log_2(i + 1) \rfloor$$

Type casting - also slow

# get(i) IVI 0 IVI 1 IVI 2 IVI 3

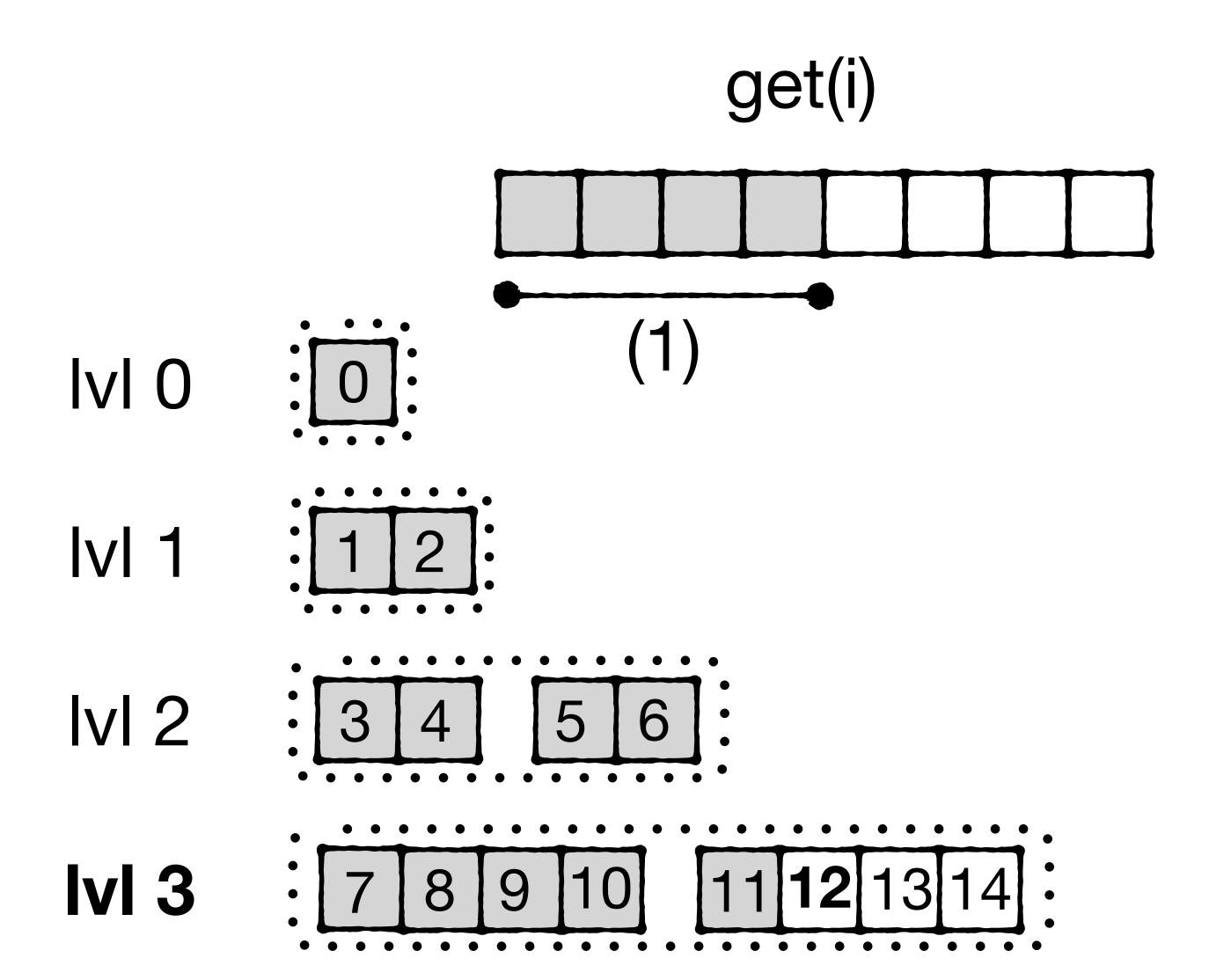
Identify target level k  $k = \lfloor \log_2(i + 1) \rfloor$ 

Insight?



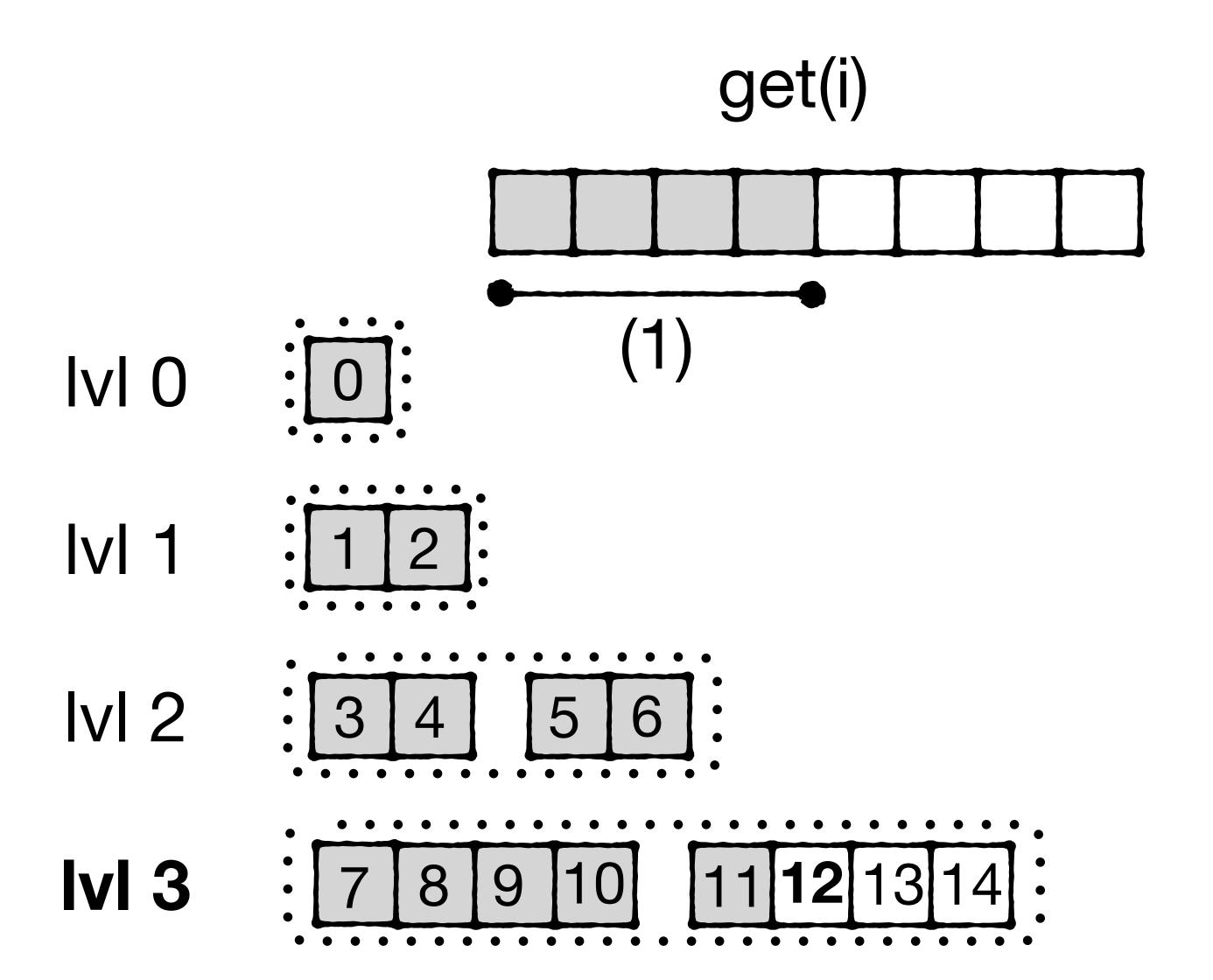
$$k = \lfloor \log_2(i + 1) \rfloor$$

Insight: log<sub>2</sub> amounts to finding index of most significant digit



$$k = \lfloor \log_2(i + 1) \rfloor$$

$$= sizeof(i) - 1 - clz(i+1)$$

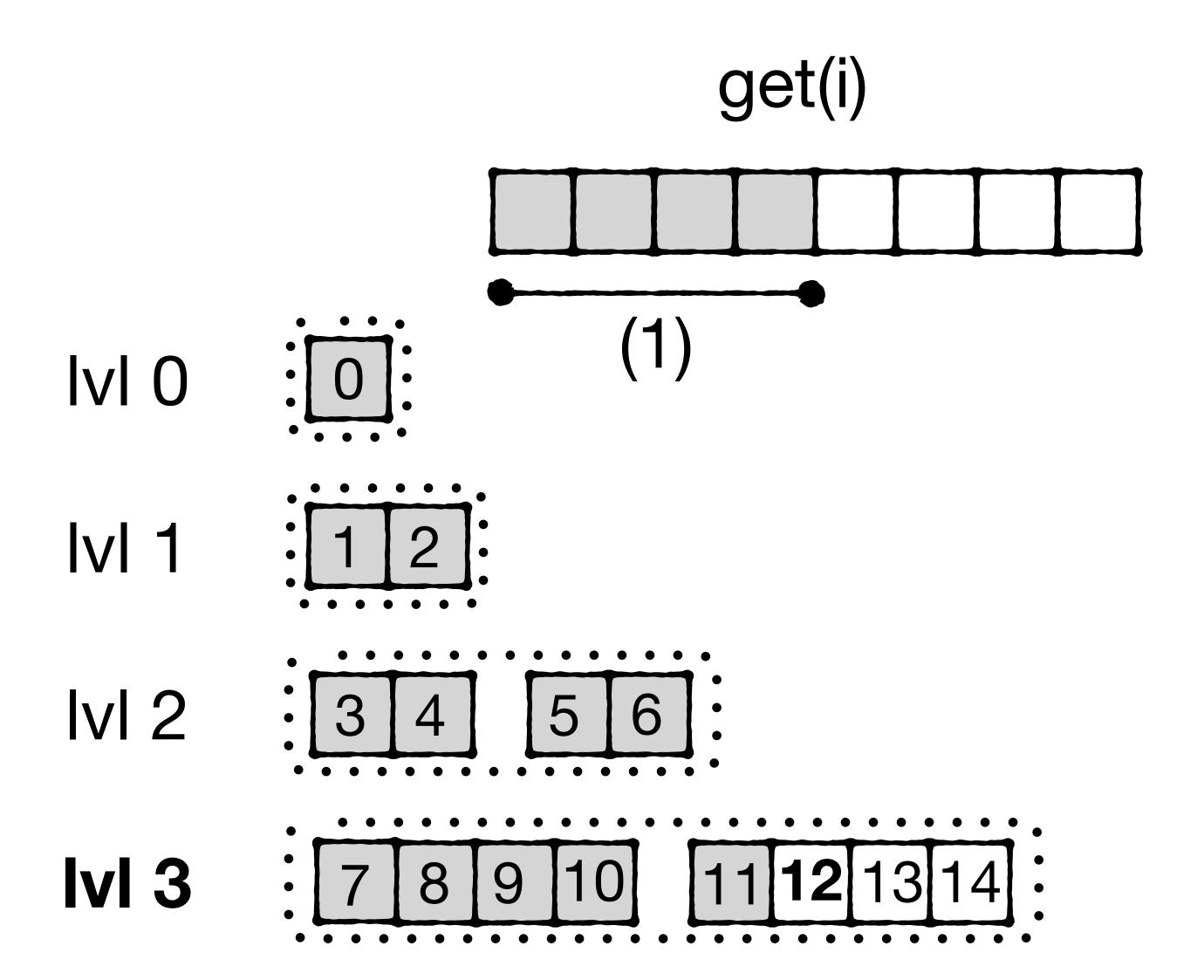


$$k = \lfloor \log_2(i + 1) \rfloor$$

$$= sizeof(i) - 1 - clz(i+1)$$

1

Integer length in bits



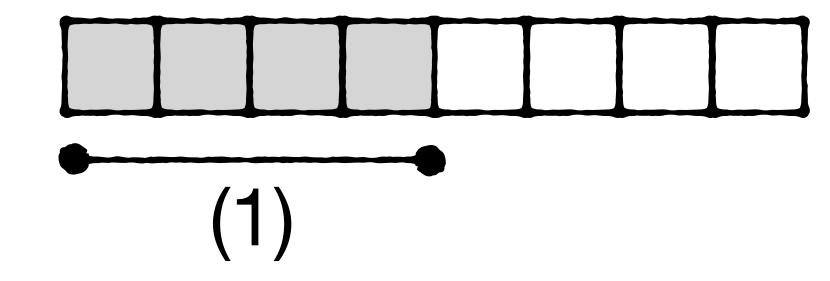
$$k = \lfloor \log_2(i + 1) \rfloor$$

$$= sizeof(i) - 1 - clz(i+1)$$

1

Specialized CPU command for # leading zeros

## get(00001100)



IVI O

IVI 1 : 1 2

IVI 2 : 3 4 5 6

IvI 3

8 9 10 11 12 13 14

Identify target level k

$$k = \lfloor \log_2(i+1) \rfloor$$

$$= sizeof(i) - 1 - clz(i+1)$$

$$8 - 1 - 4 = 3$$

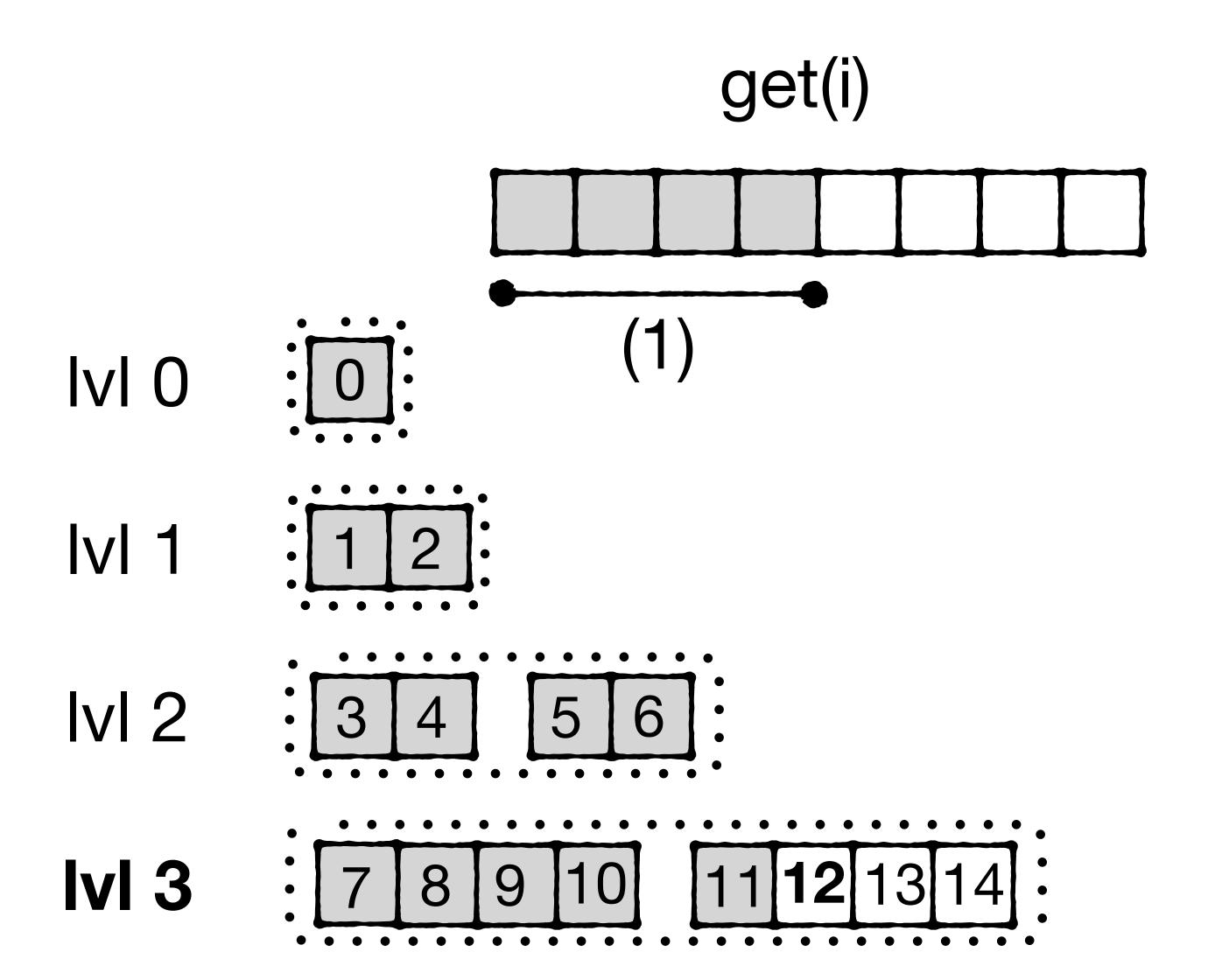
# get(i) IVI O IVI 1 IVI 2 IvI 3

Identify target level k

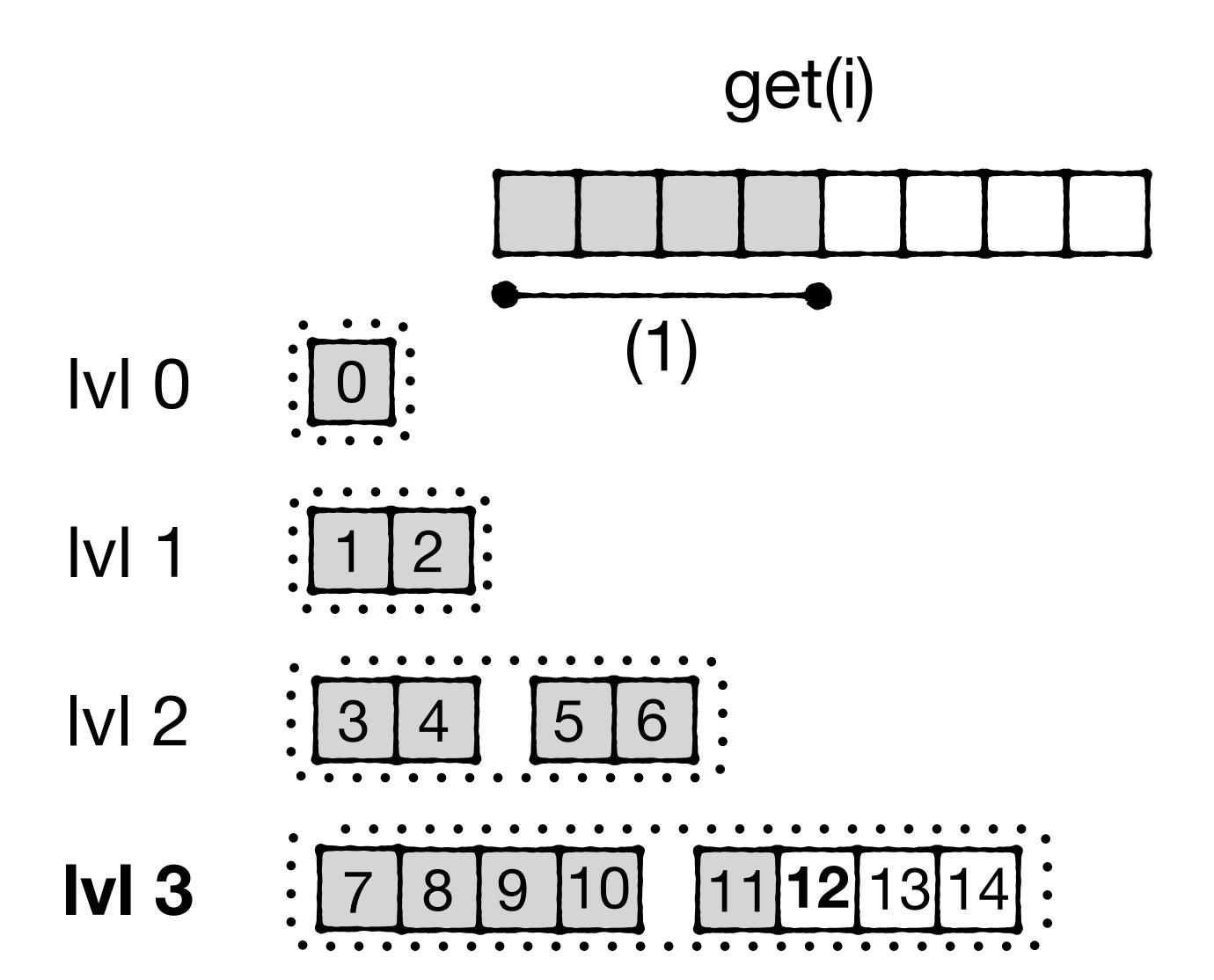
$$k = \lfloor \log_2(i + 1) \rfloor$$

$$= sizeof(i) - 1 - clz(i+1)$$

≈1 ns rather than ≈7ns



# blocks in levels 0 to k-1?



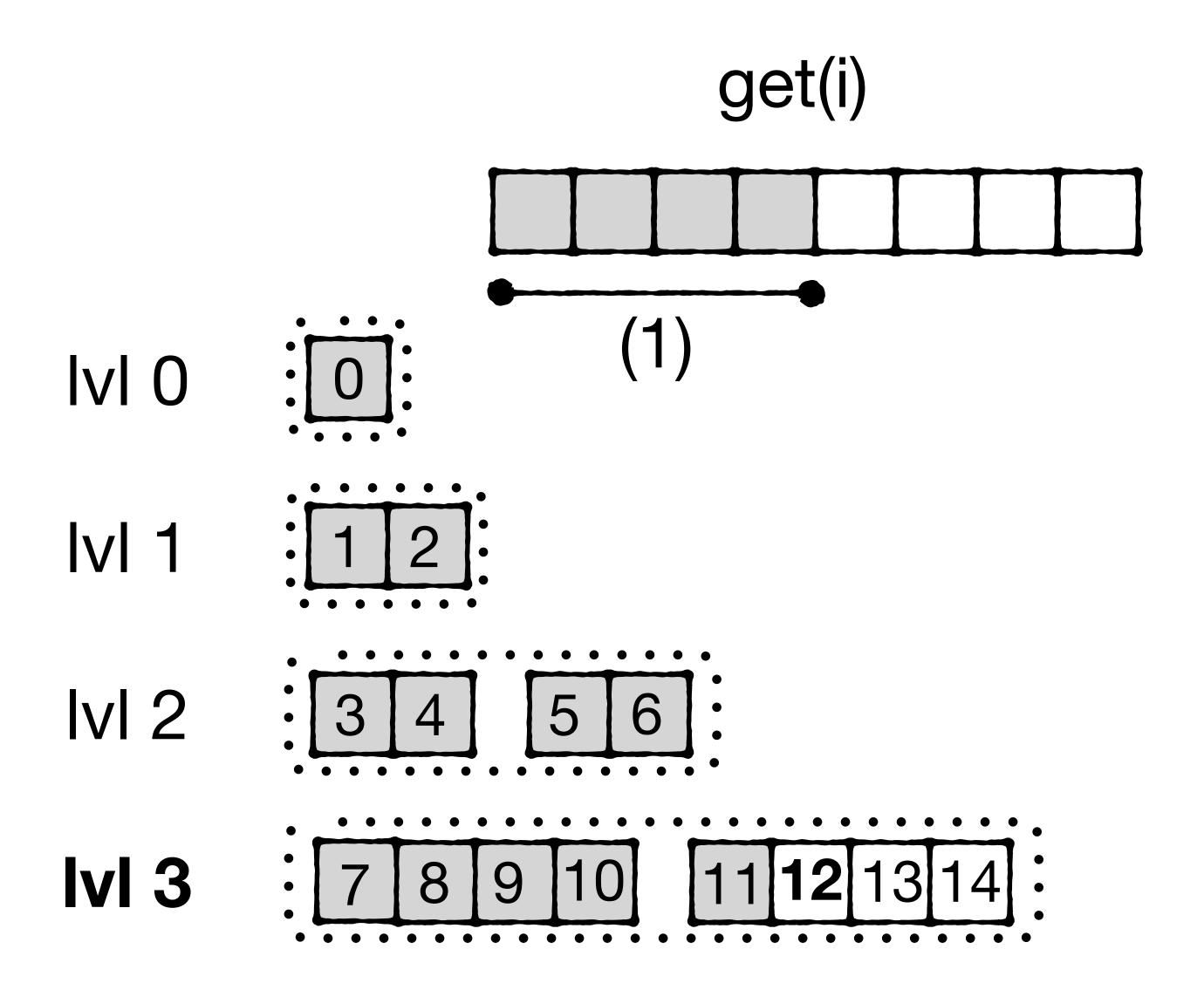
 $= 2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$ 

### get(i) IVI 0 IVI 1 IVI 2 IVI 3

#### # blocks in levels 0 to k-1

$$=2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$

# Original paper gets this wrong, fixed credit to Hyuhng Min



$$= 2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$

Intuition: number of new data blocks grows every other level

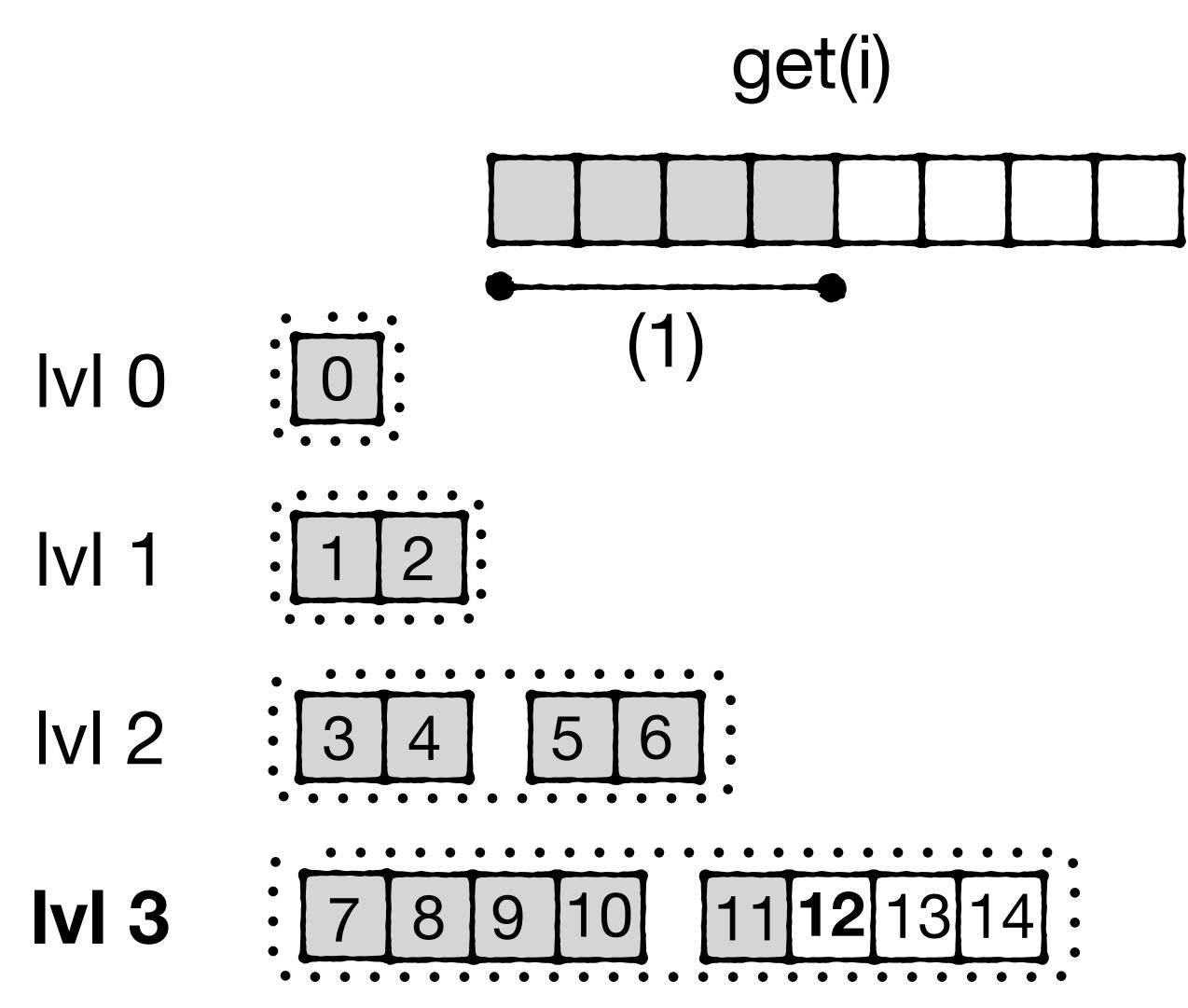
## get(i) IVI 0 IVI 1 IVI 2

IVI 3

# blocks in levels 0 to k-1

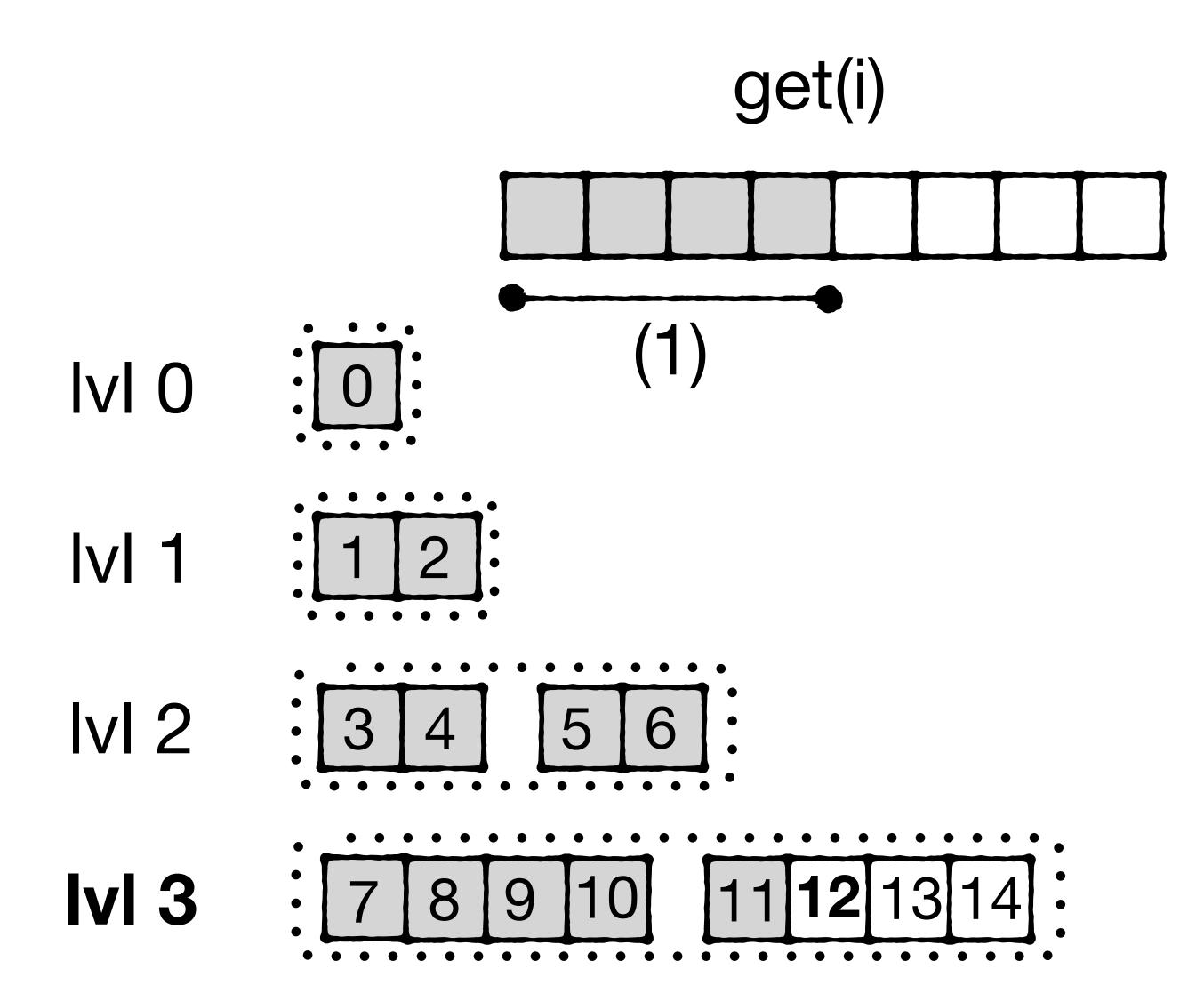
$$=2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$

Level k	# Blocks
0	0
1	1
2	2
3	4
4	6
5	10
6	14
7	22



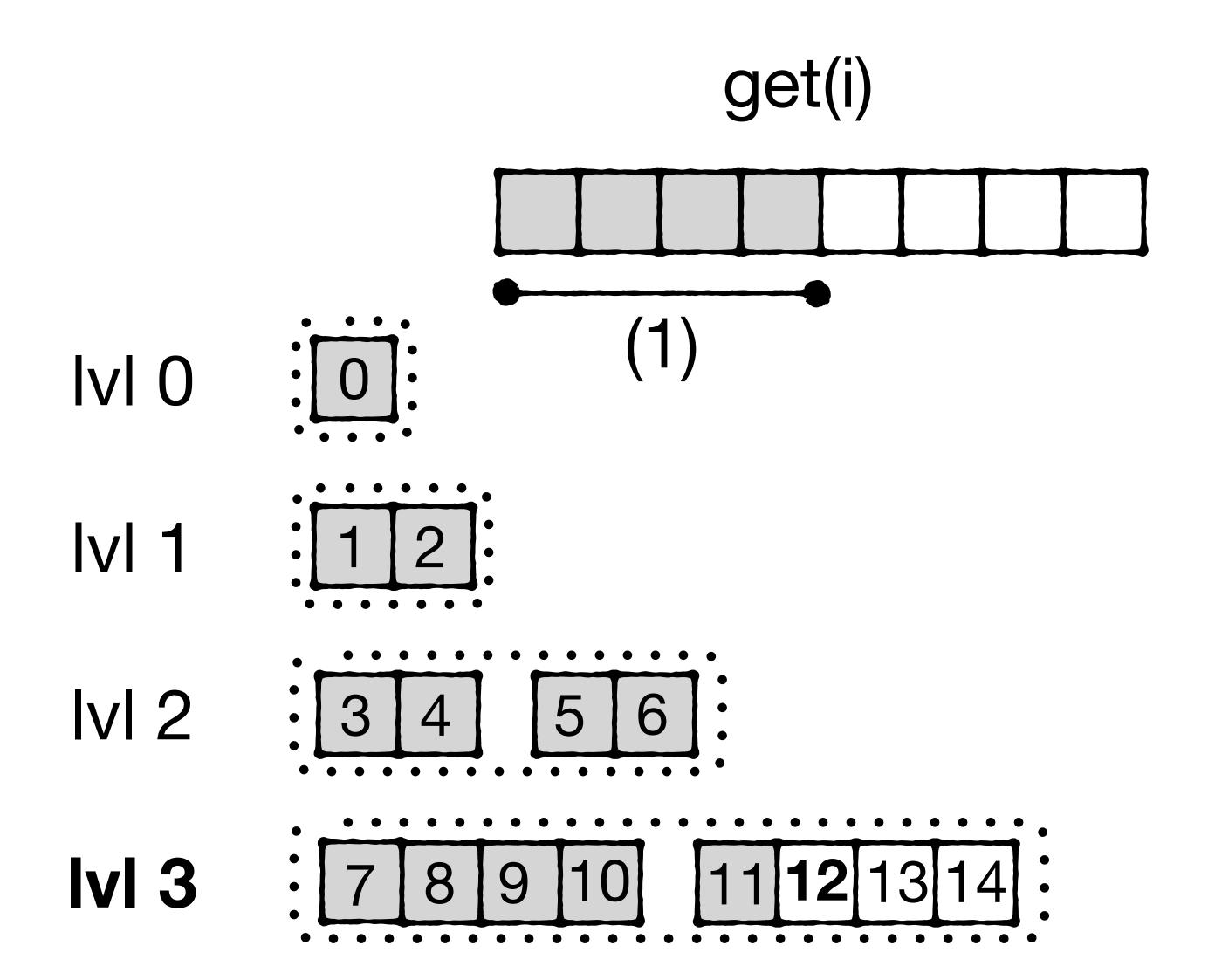
$$= 2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$

Integer division is slow



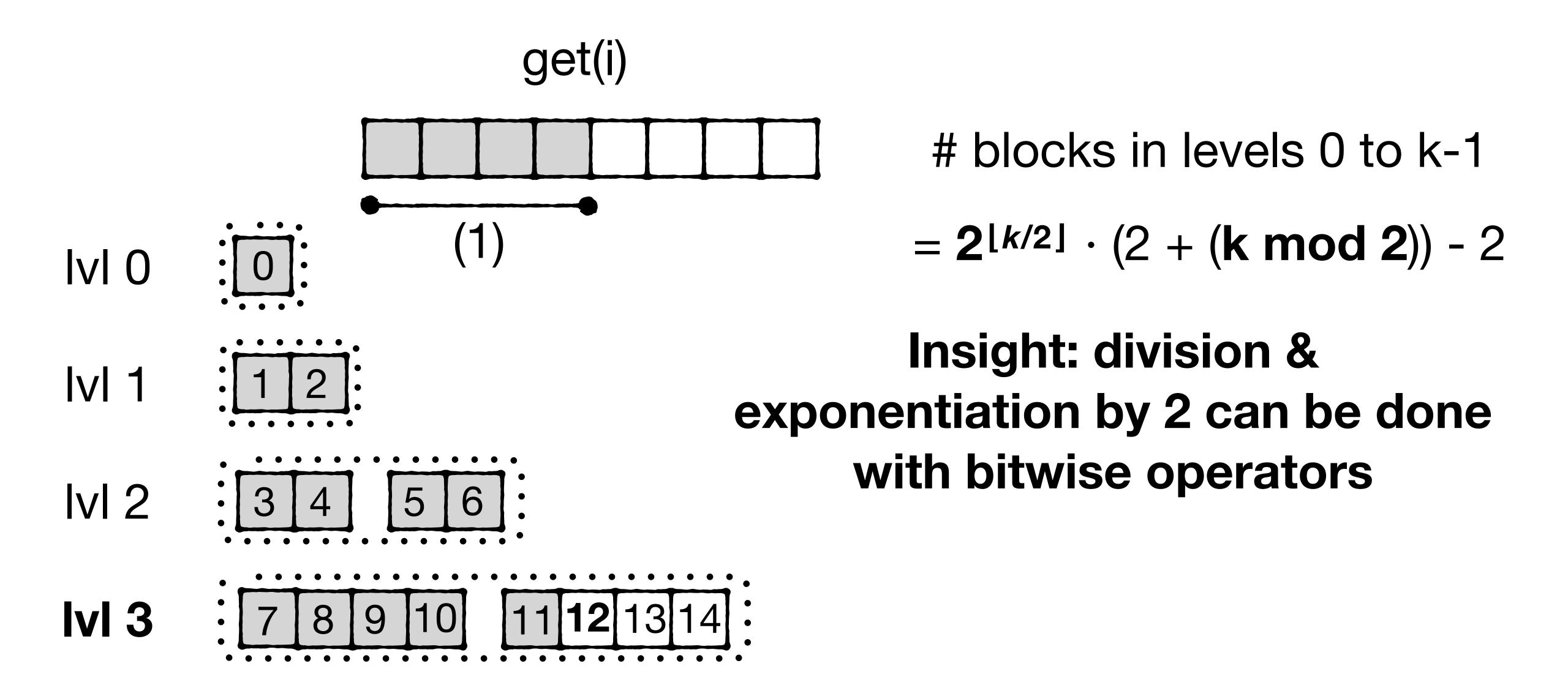
$$= 2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$

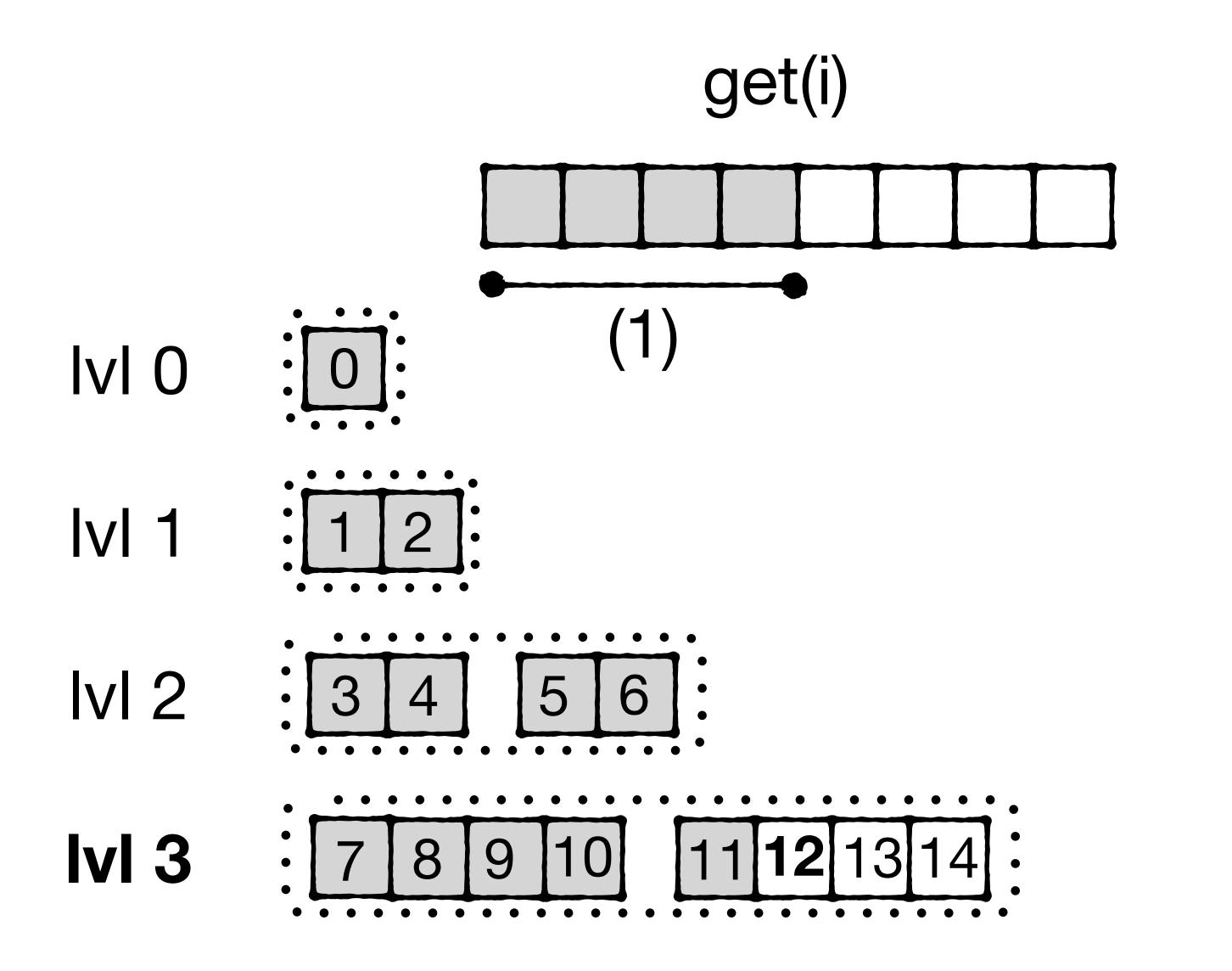
Power is slow



$$=2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$

How to speed up?

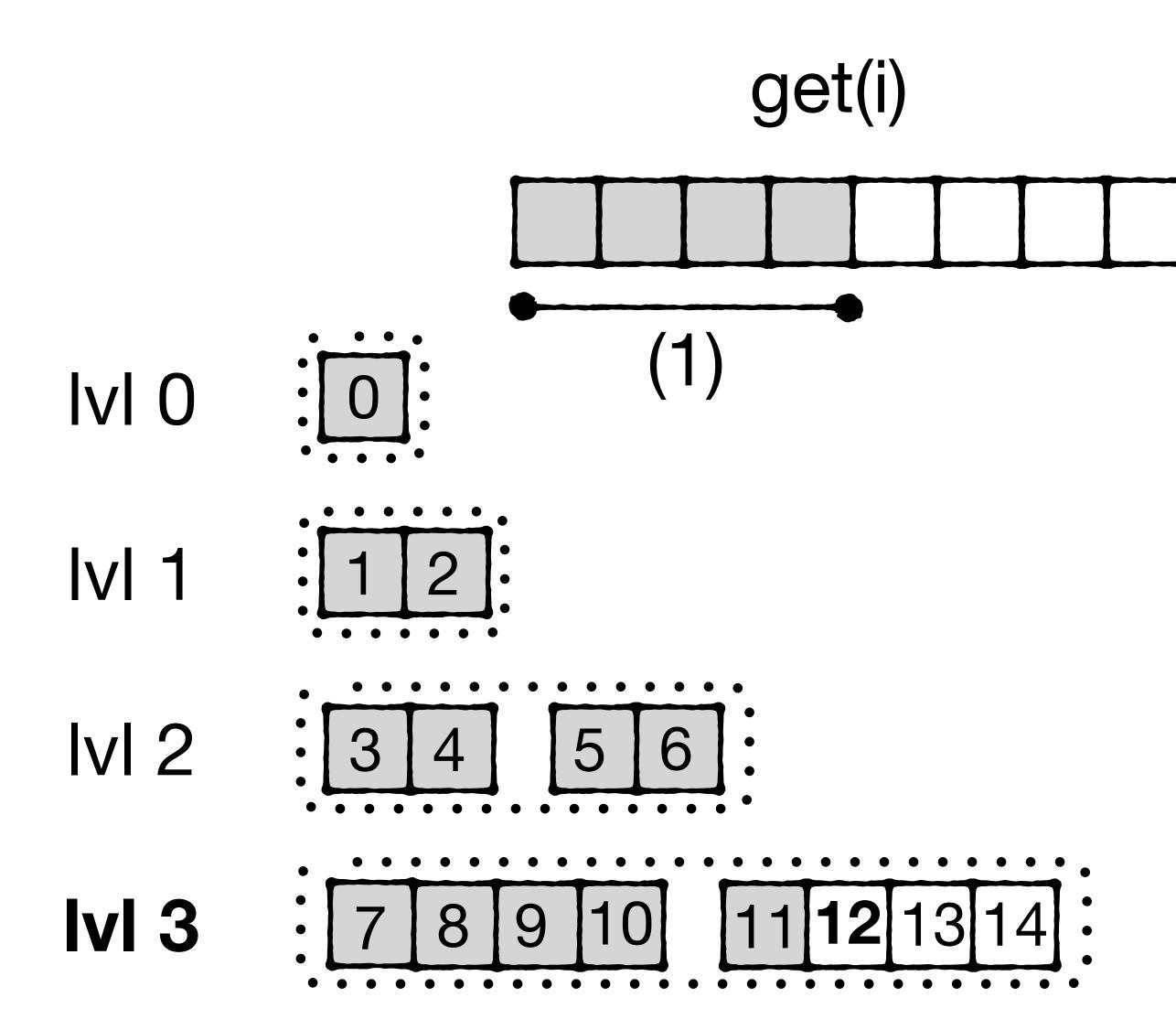




$$=2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$

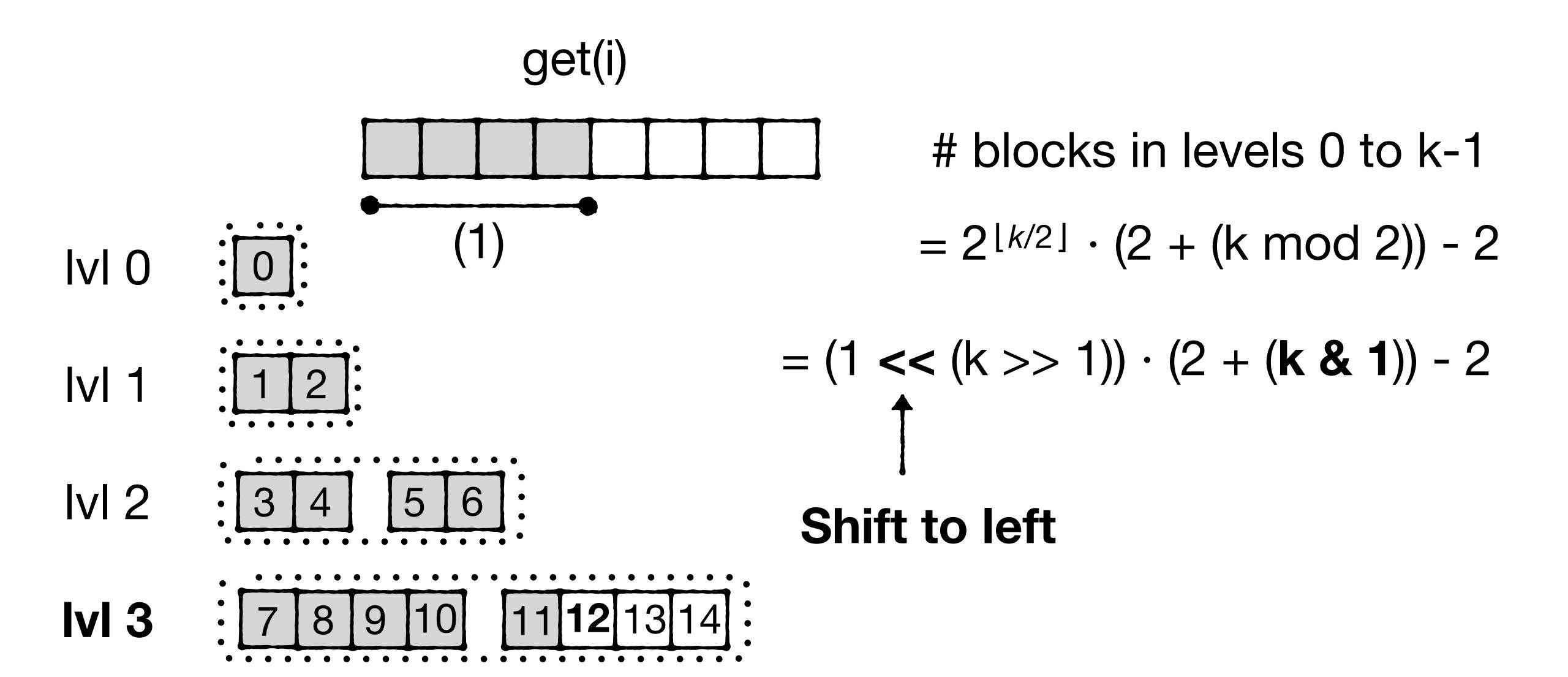
$$= 2^{\lfloor k/2 \rfloor} \cdot (2 + (k \& 1)) - 2$$

"and" with 1



$$= 2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$
$$= 2^{(k >> 1)} \cdot (2 + (k \& 1)) - 2$$

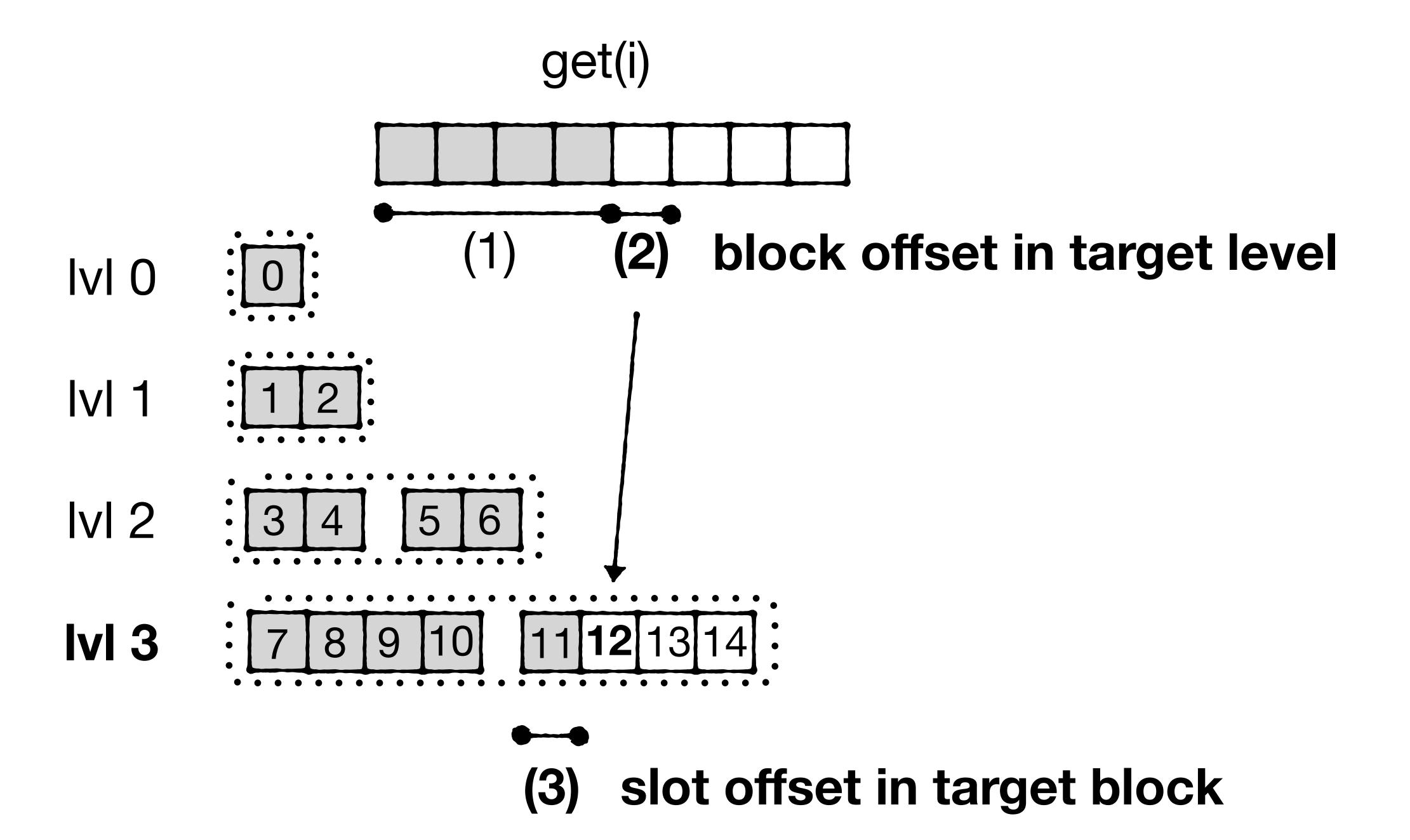
Shift by 1 bit to right

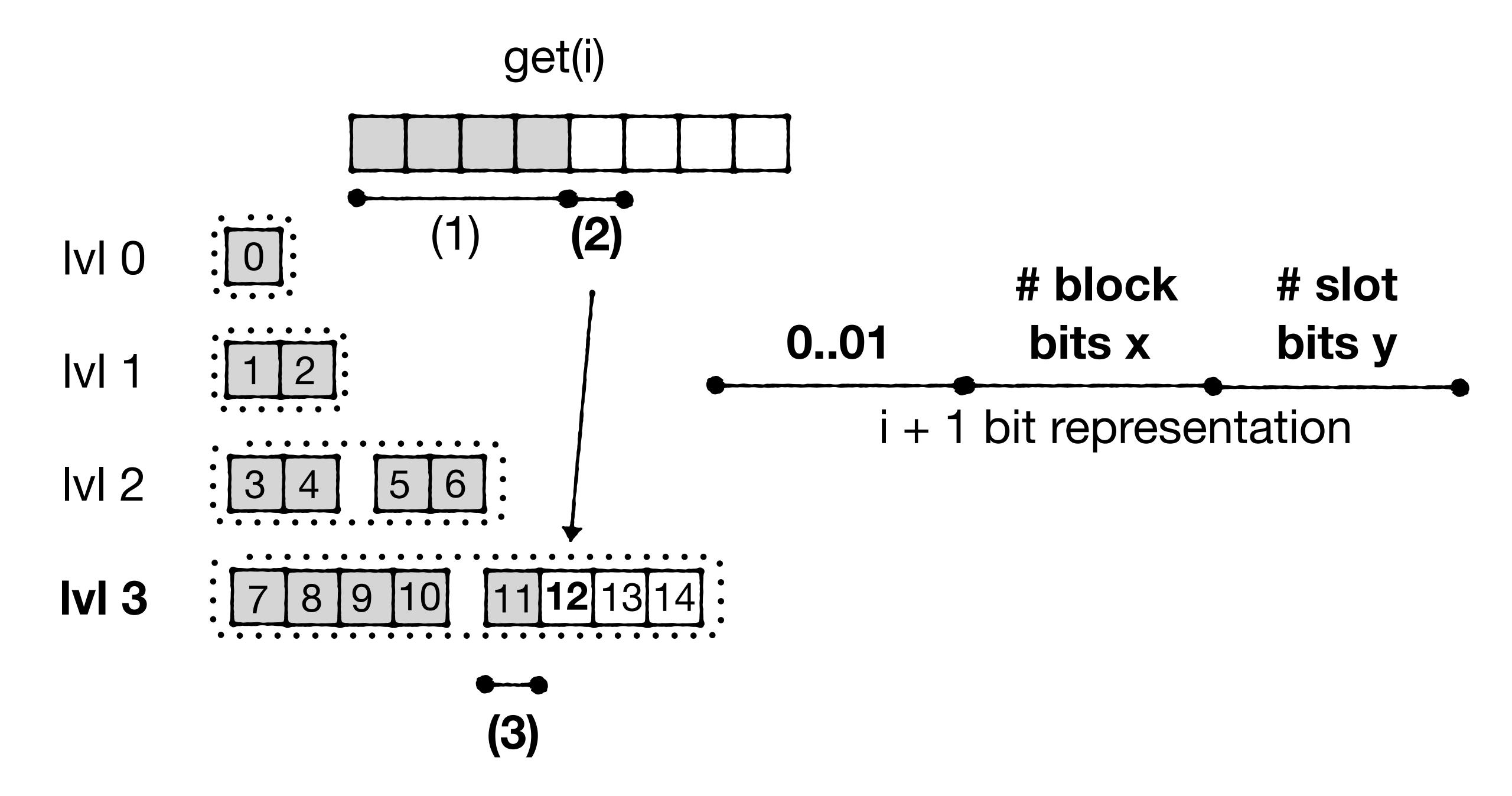


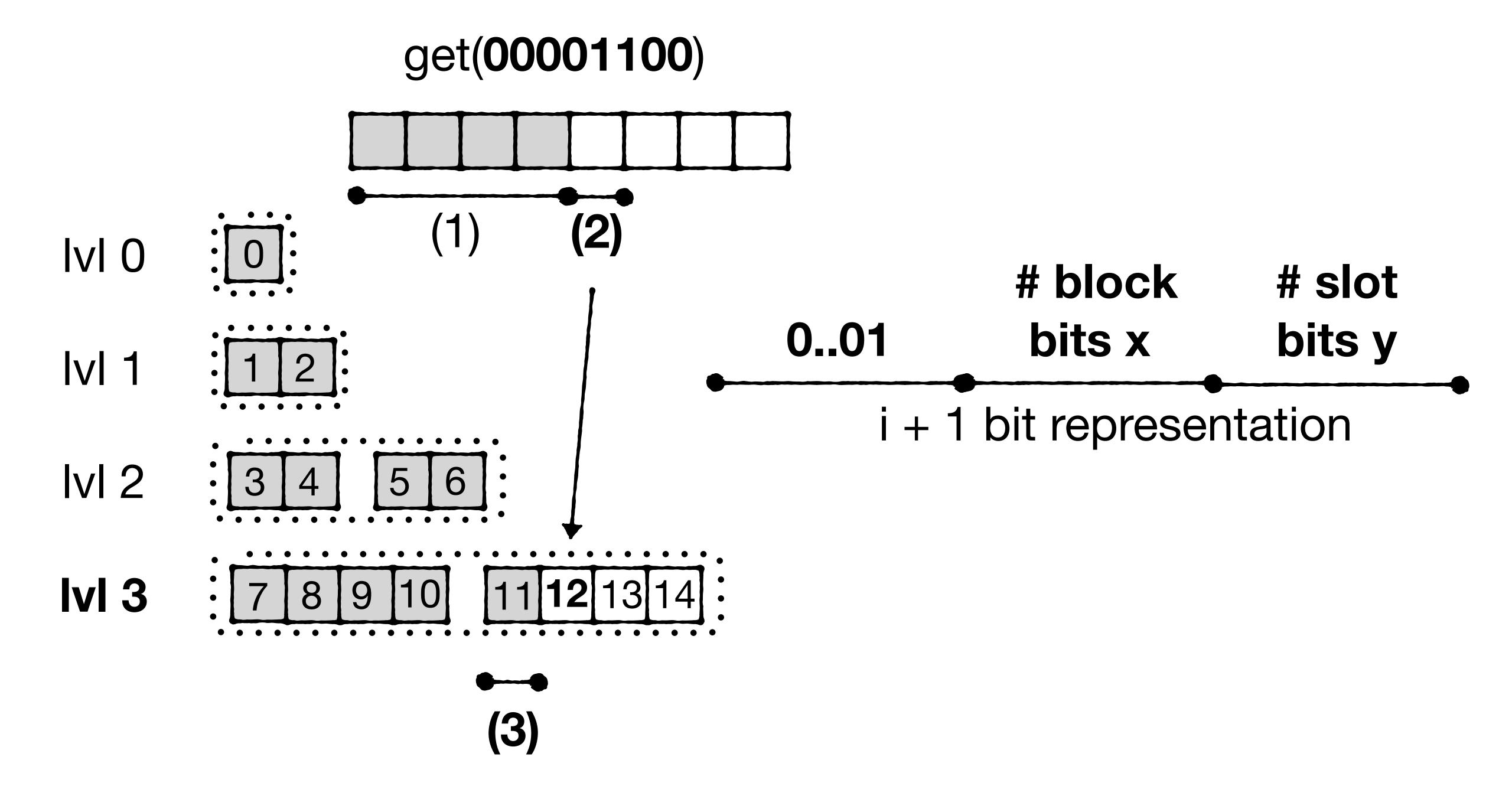
$$= 2^{\lfloor k/2 \rfloor} \cdot (2 + (k \mod 2)) - 2$$

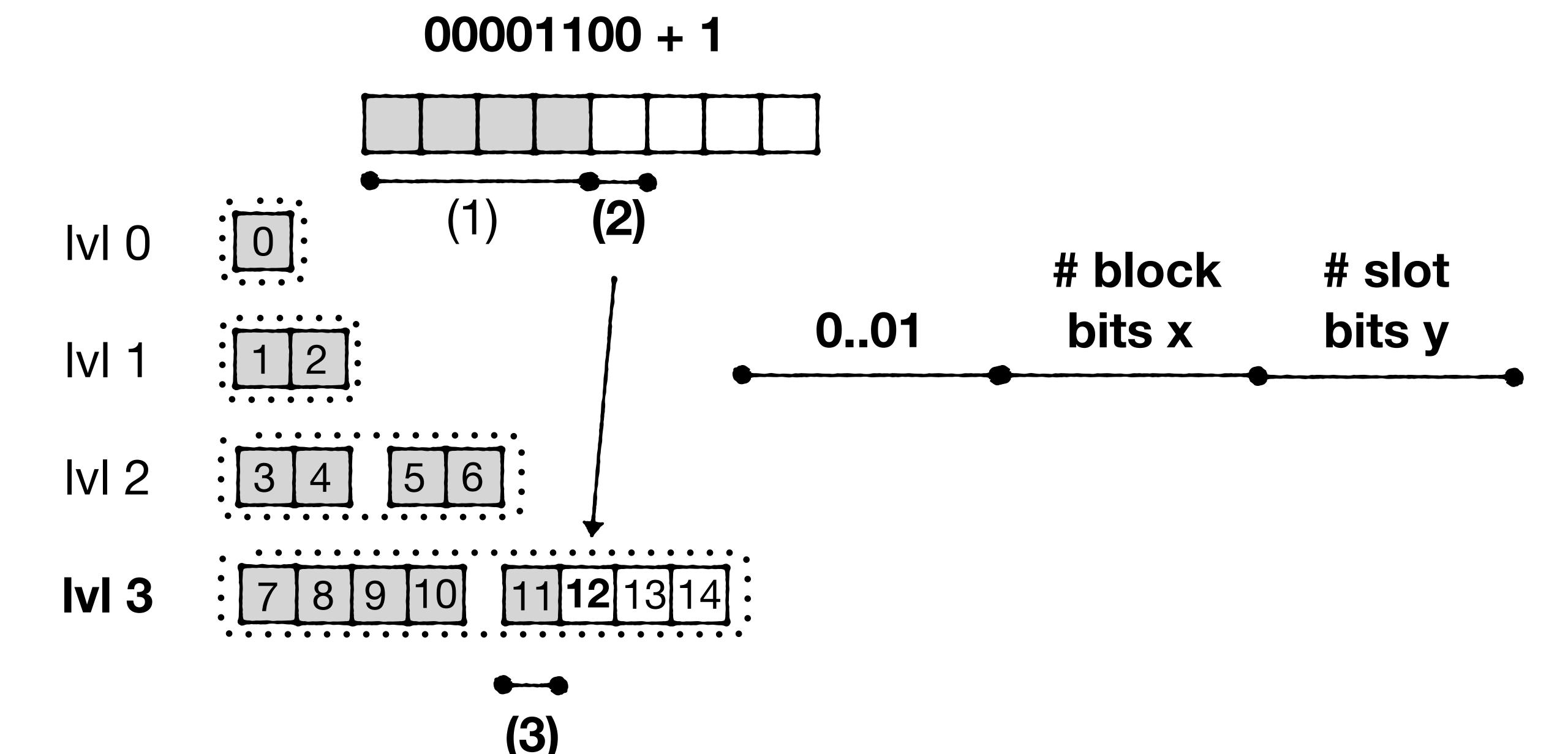
$$= (1 << (k >> 1)) \cdot (2 + (k \& 1)) - 2$$

≈0.6 ns rather than ≈10ns

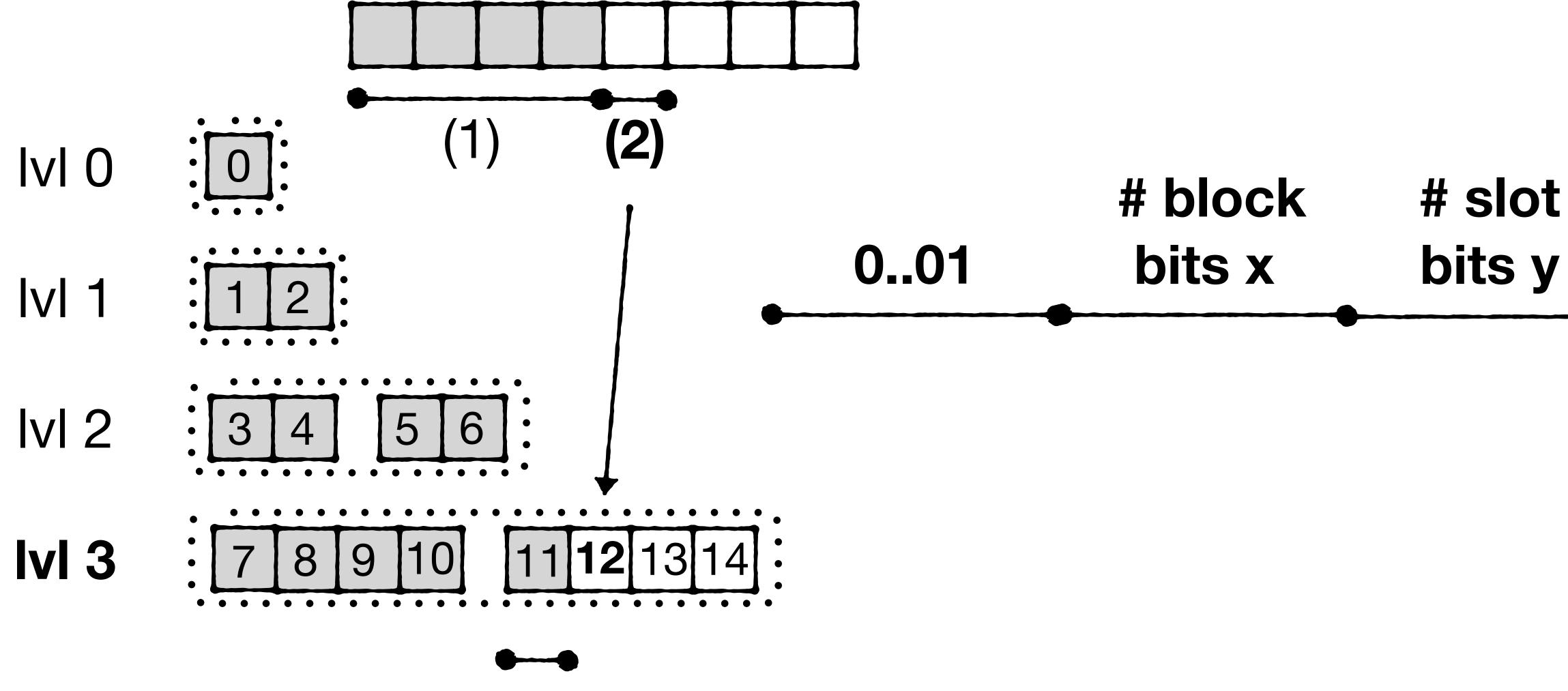


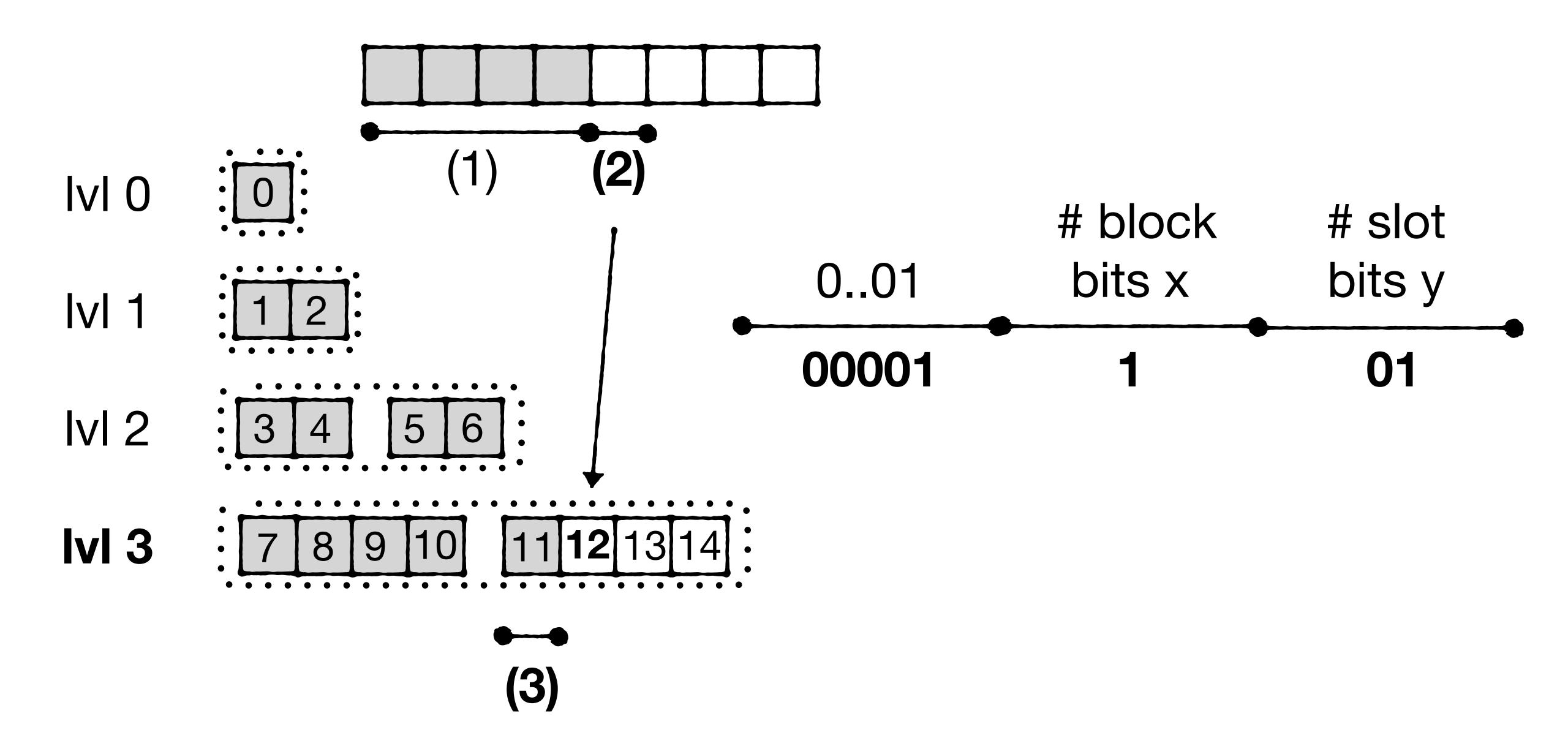


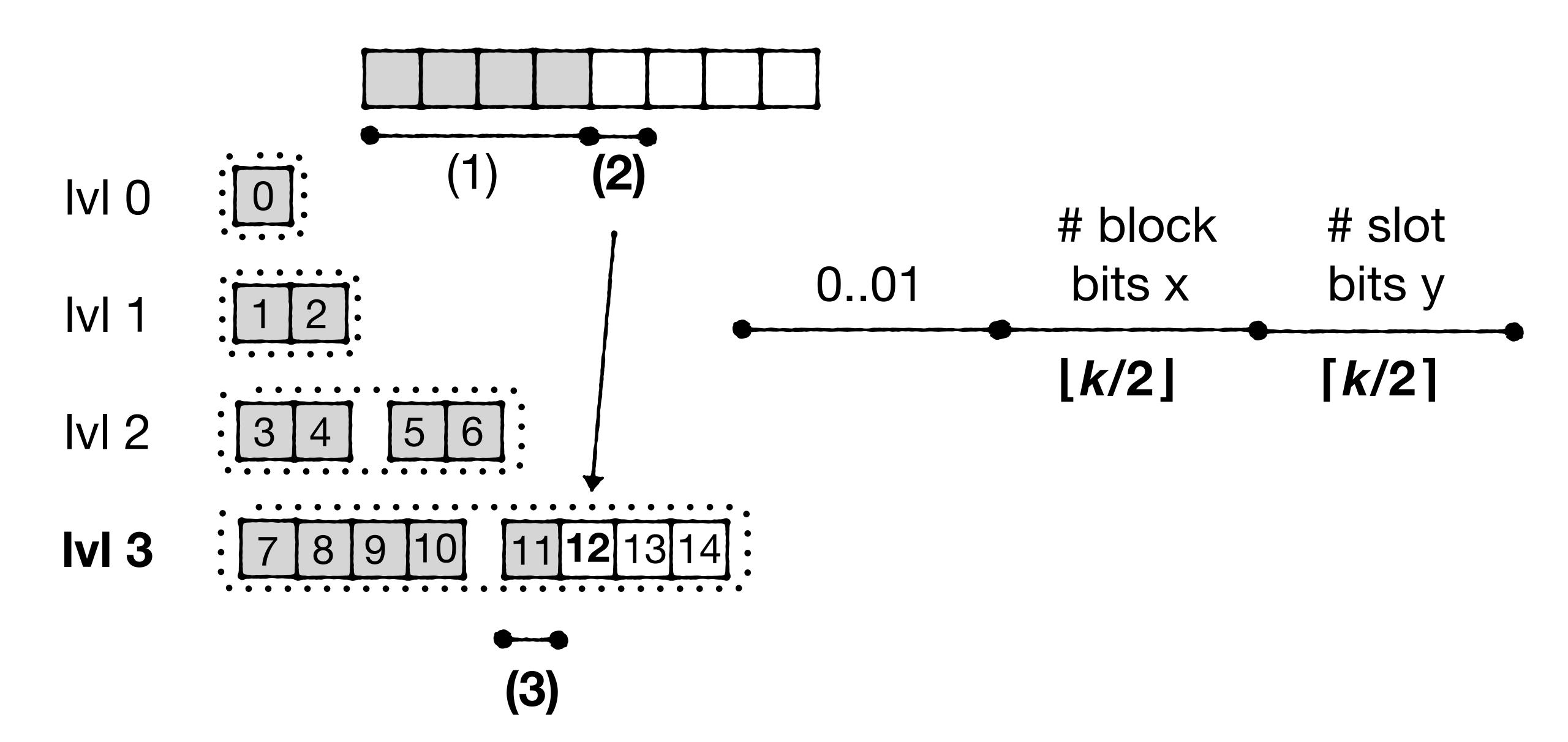


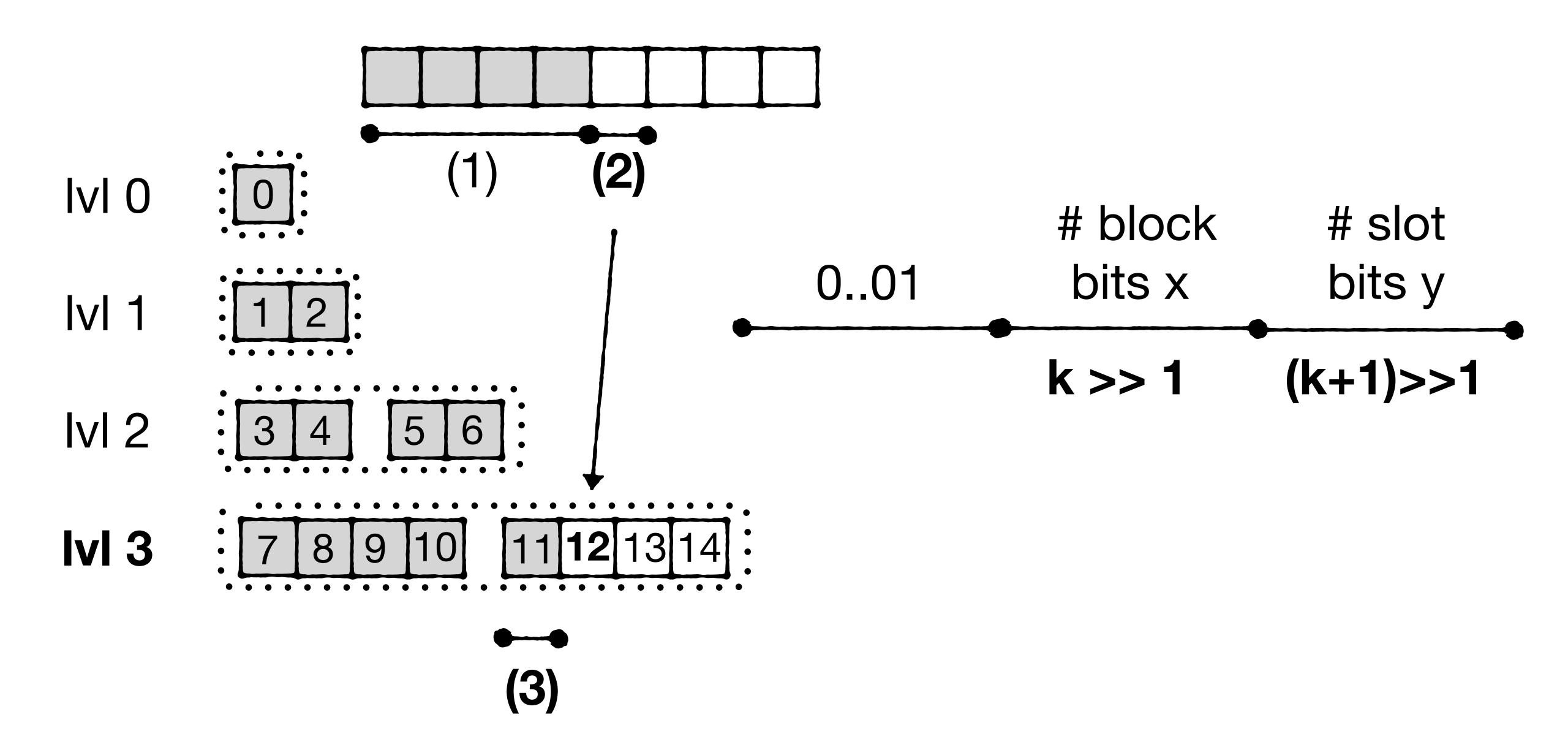


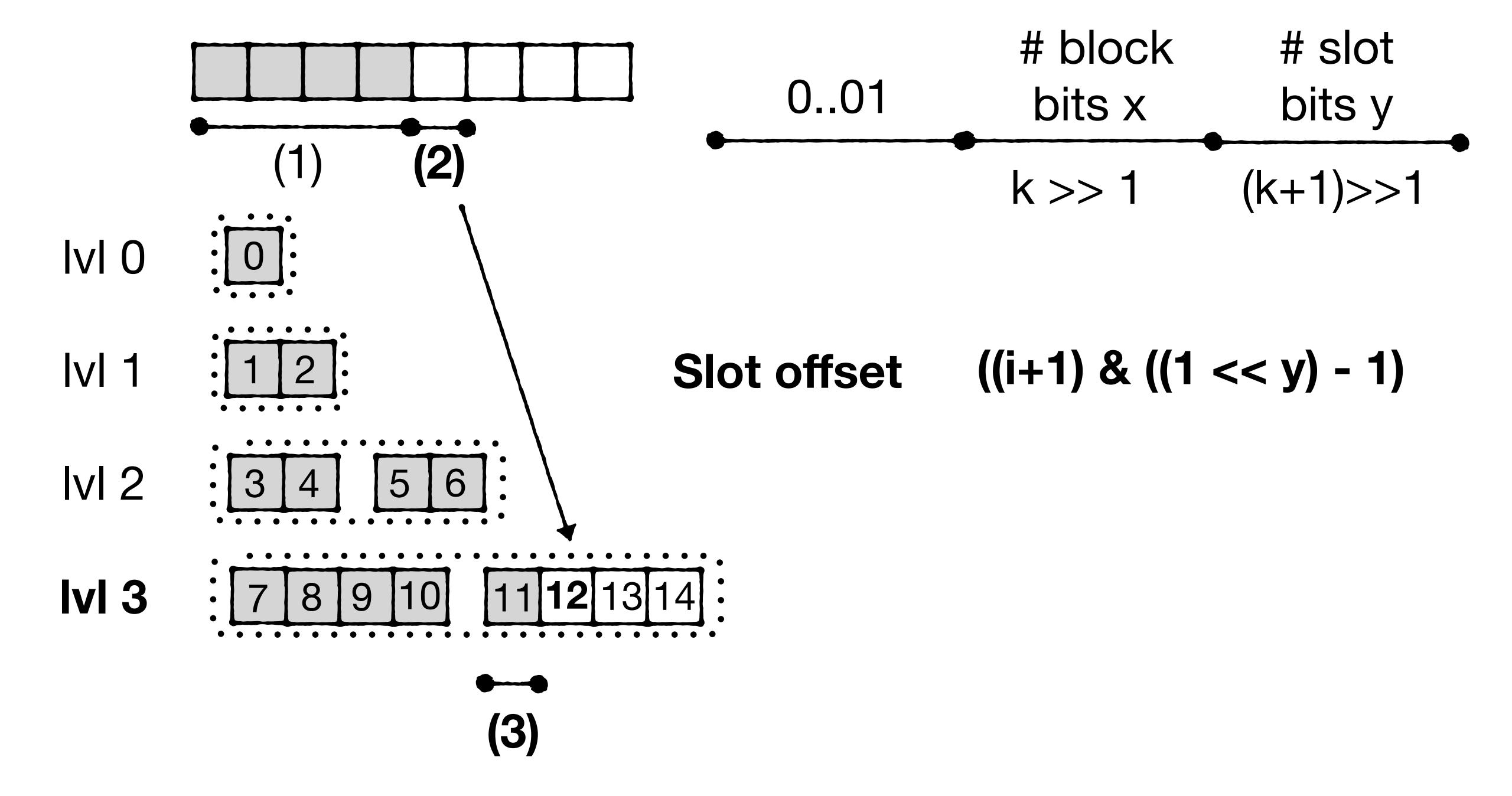
## 

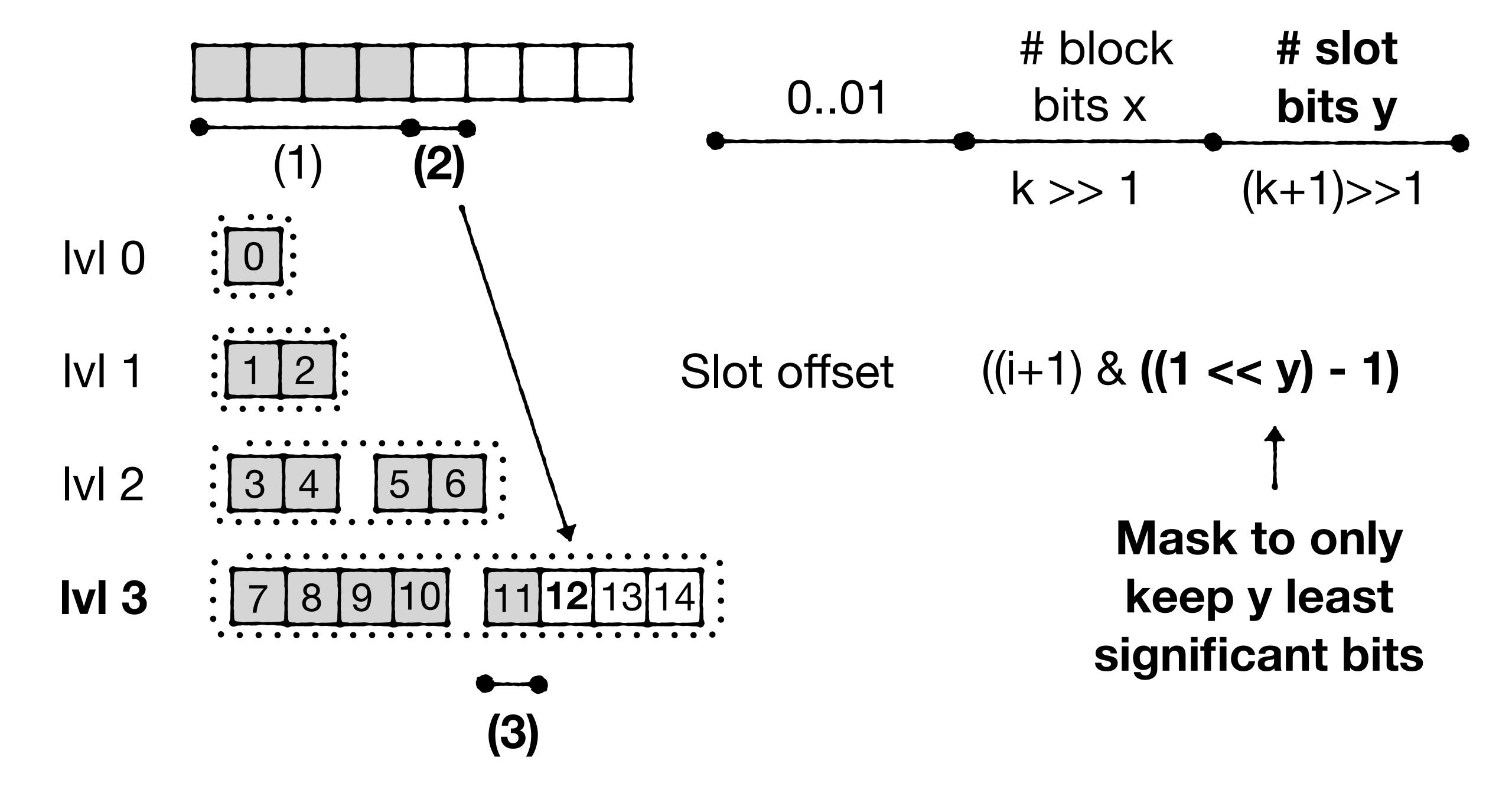


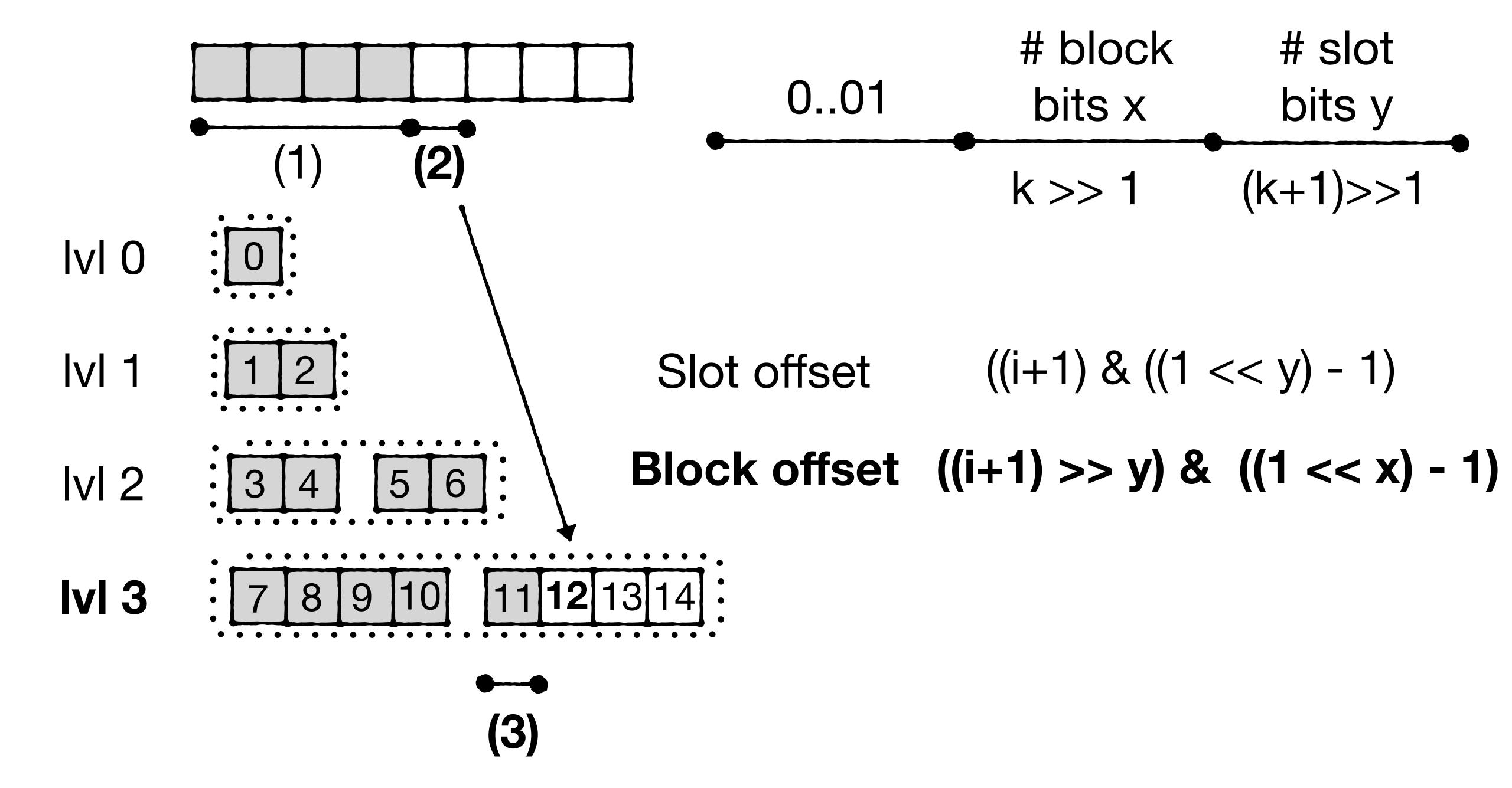


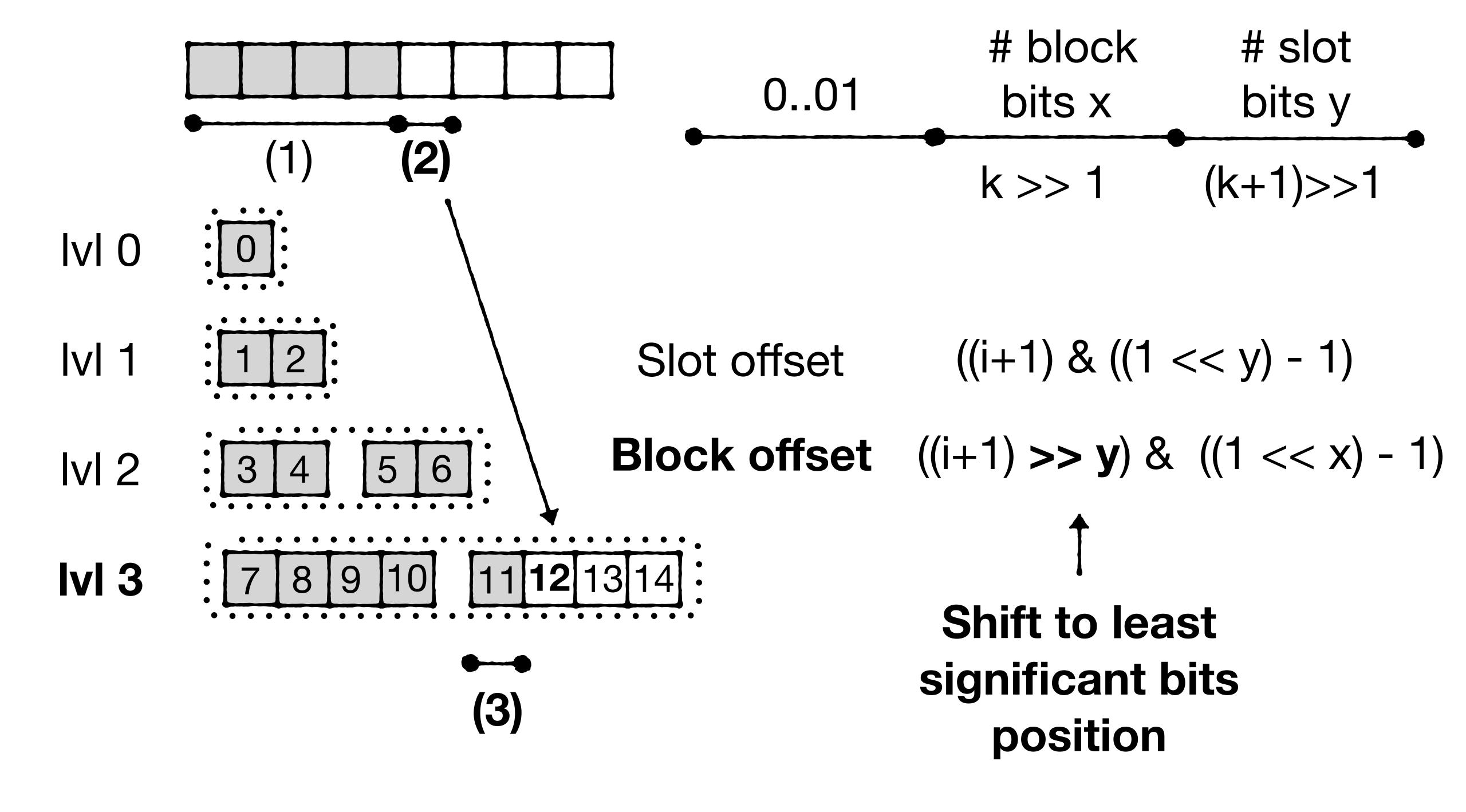


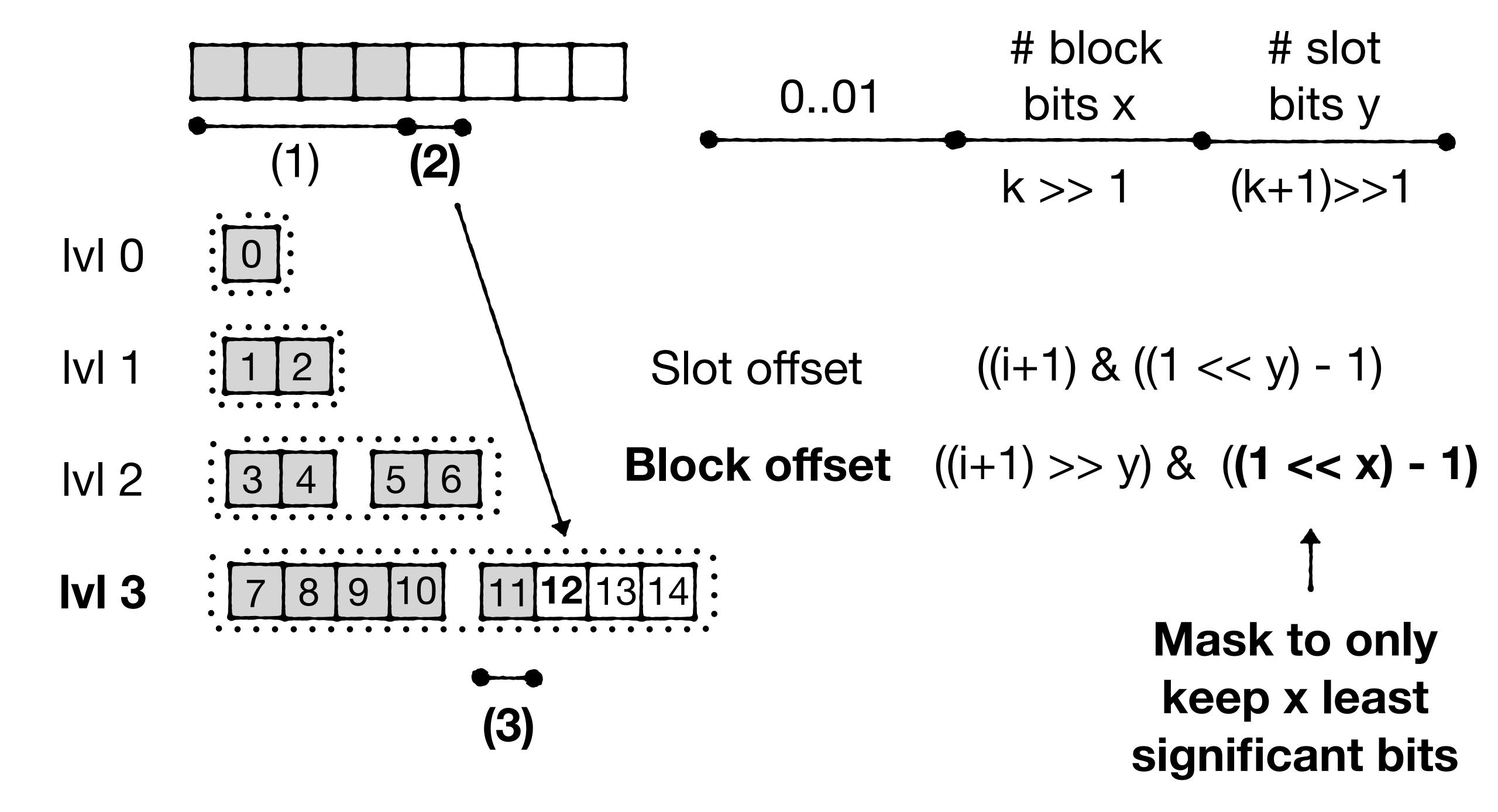


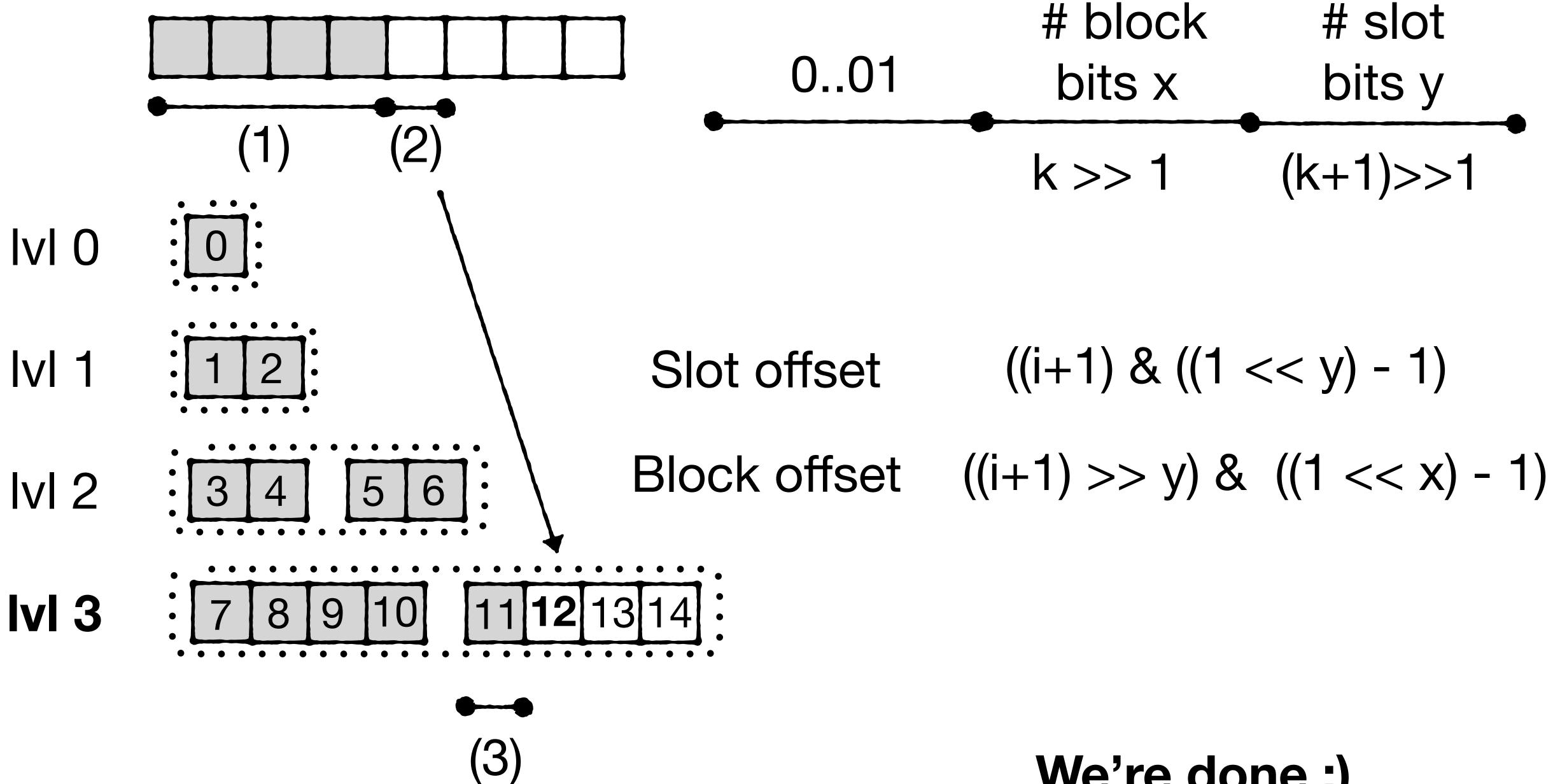












We're done:)

Write-amp Space-amp Read-amp

indirection ≈1 O(1+N<sup>-0.5</sup>) O(1)

Read-amp

indirection

**≈1** 

 $O(1+N^{-0.5})$ 

0(1)

No indirection

G G - 1

O(G)

1

indirection

**≈**1

 $O(1+N^{-0.5})$ 

O(1)

No indirection

**G** > ф

G G - 1 G + G G - 1

1

No indirection

**G** < ф

G - 1

G to G+1

1

Thank you:)