

# plcpcmp

**Patrick Dinklage**

**Jonas Ellert**

**Johannes Fischer**

***Dominik Köppl***

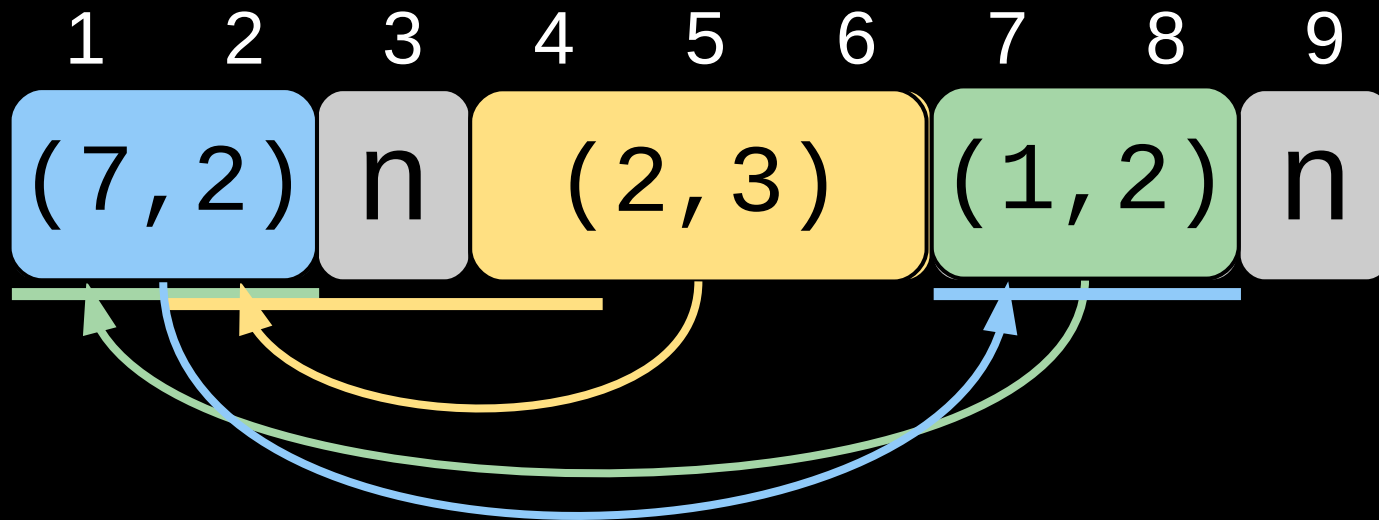
**Manuel Penschuck**

# greedy compression

scheme	method	ref
LZ76	unidirectional	Lempel,Ziv '76
longest first	grammar	中村 + '09
lcpcomp	bidirectional	Dinklage+ '17
LZ-LFS	hybrid	Mauer+ '17 西 + '18
自己参照なし	without self-ref	藤原 + '19
unicomp	unidirectional	峰松 + '19

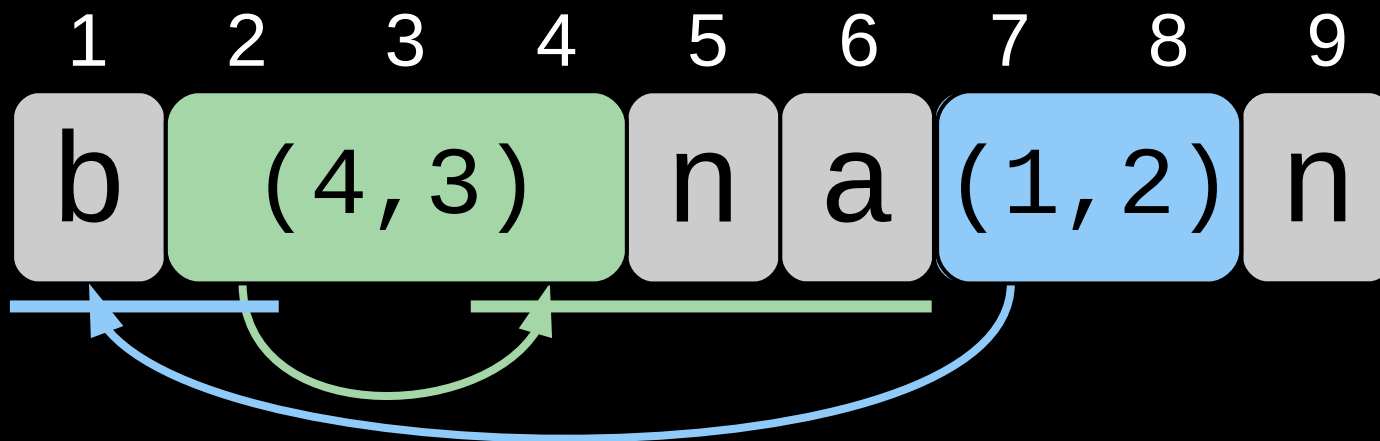
# bidirectional compression

- replace substrings by references
- in any direction
- allow self-references
- no cycles



# lcpcomp

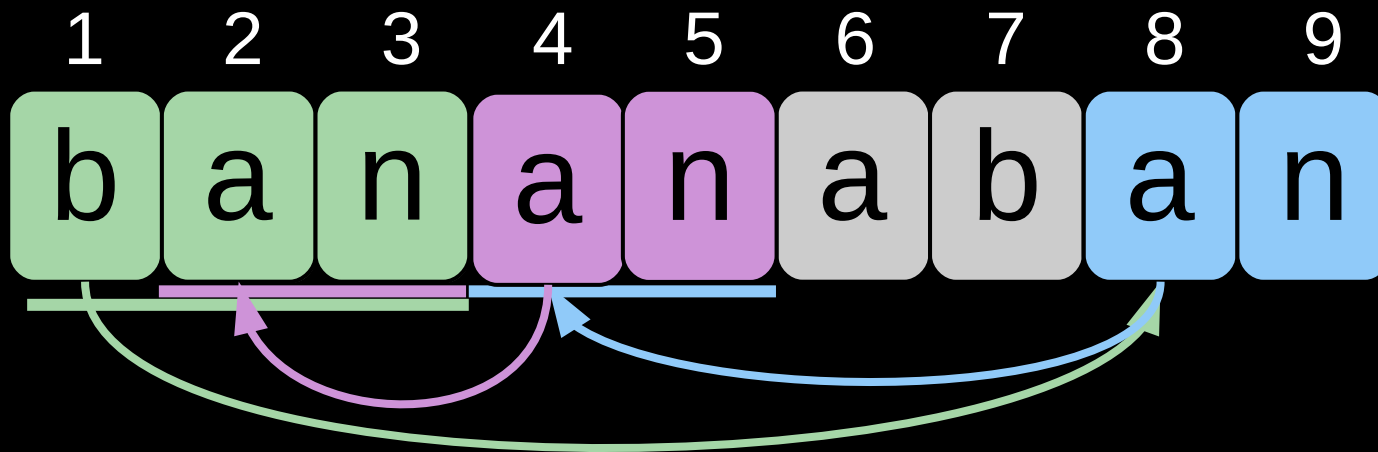
- replace longest reoccurring substring
- recurse



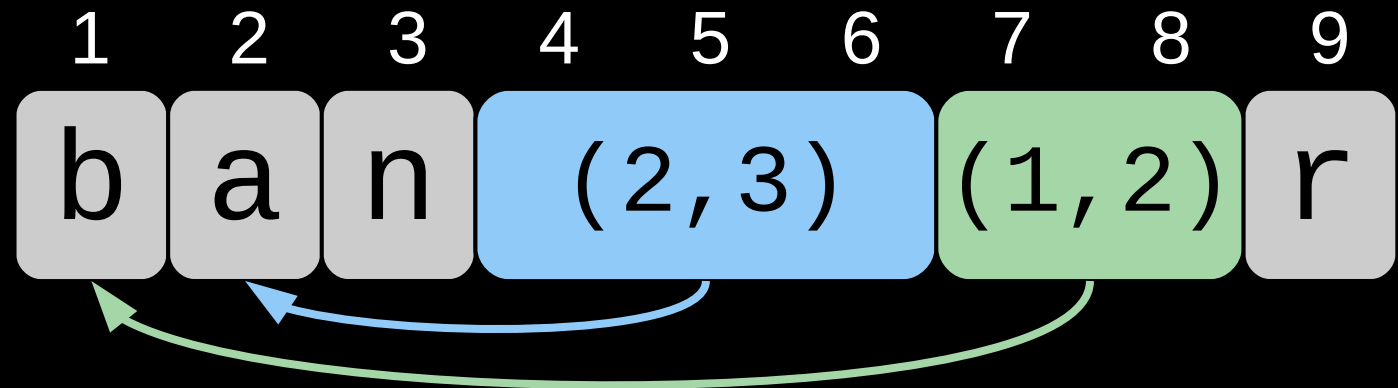
*Starting from position 4, copy 3 symbols.  
Starting from position 1, copy 2 symbols.*

# cycles

- need restriction to prevent cycles
- rule: reference to *lexicographically smaller* suffix



# ~ LZ77



- LZ77

- reference to *previous text position*

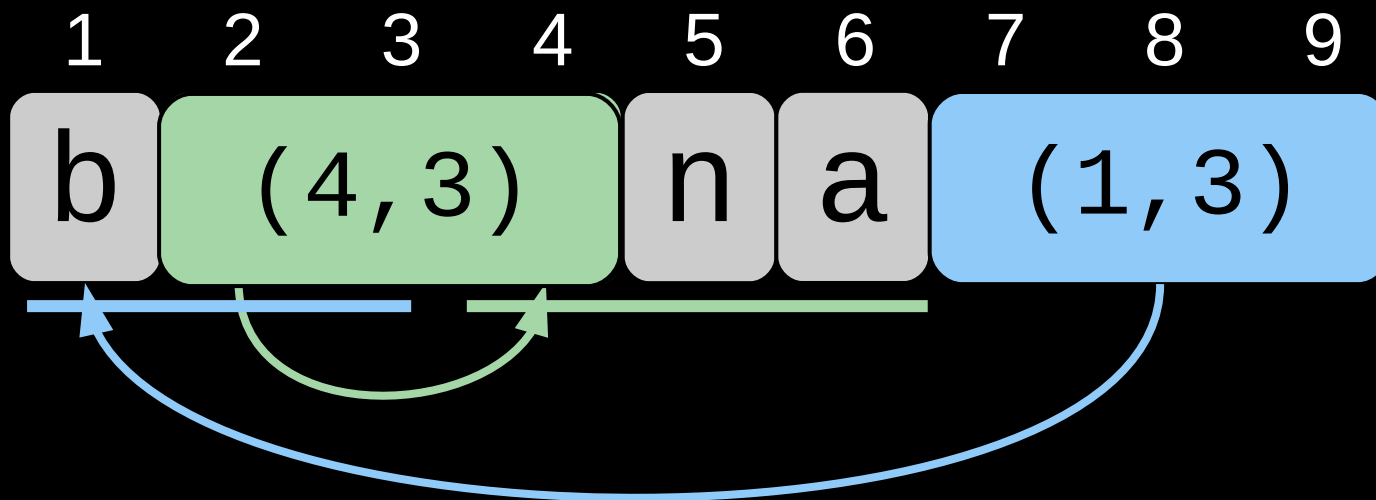
- lcpcomp

- reference to *lexicographically smaller suffix*

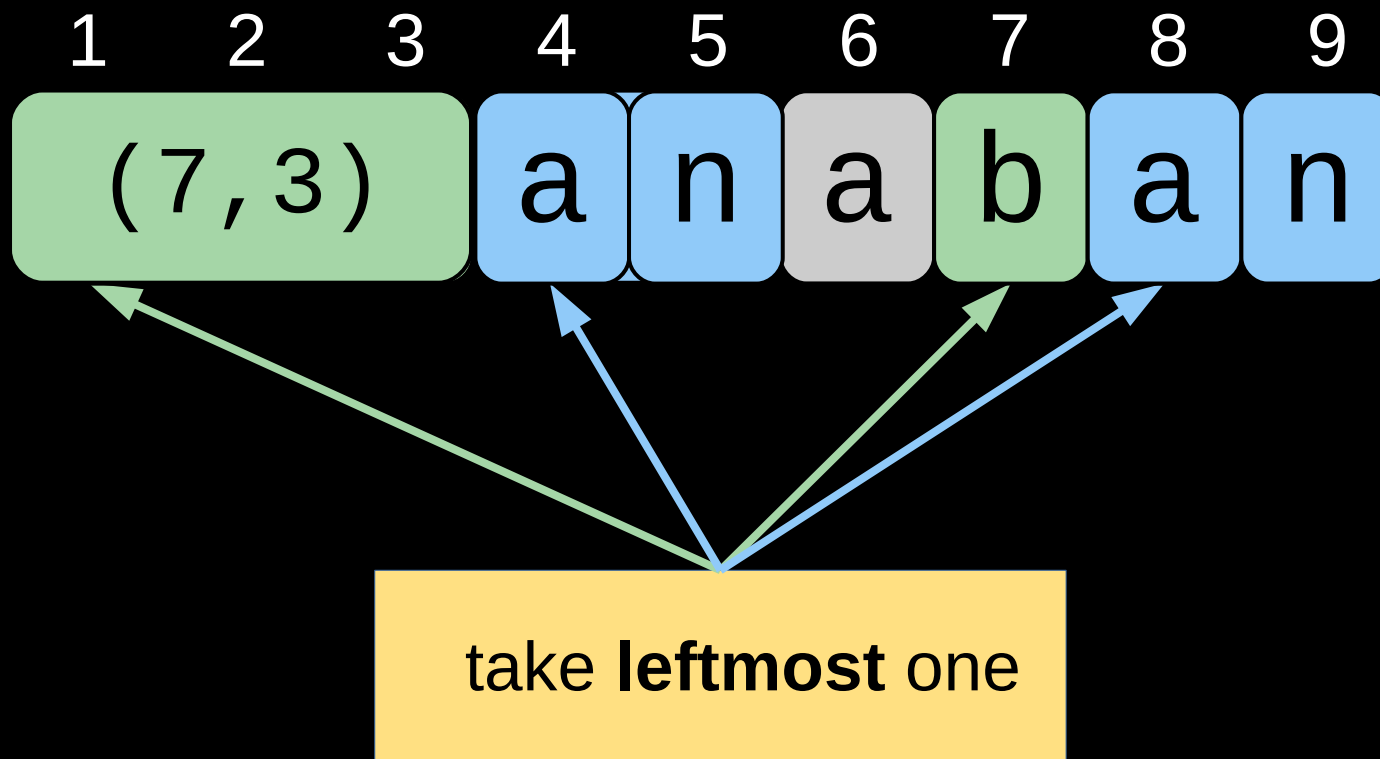
=> both greedy + reference order constraint

# lcpcomp

- multiple choices
- order not clear



# plcpcomp: tie break





# how?

## find

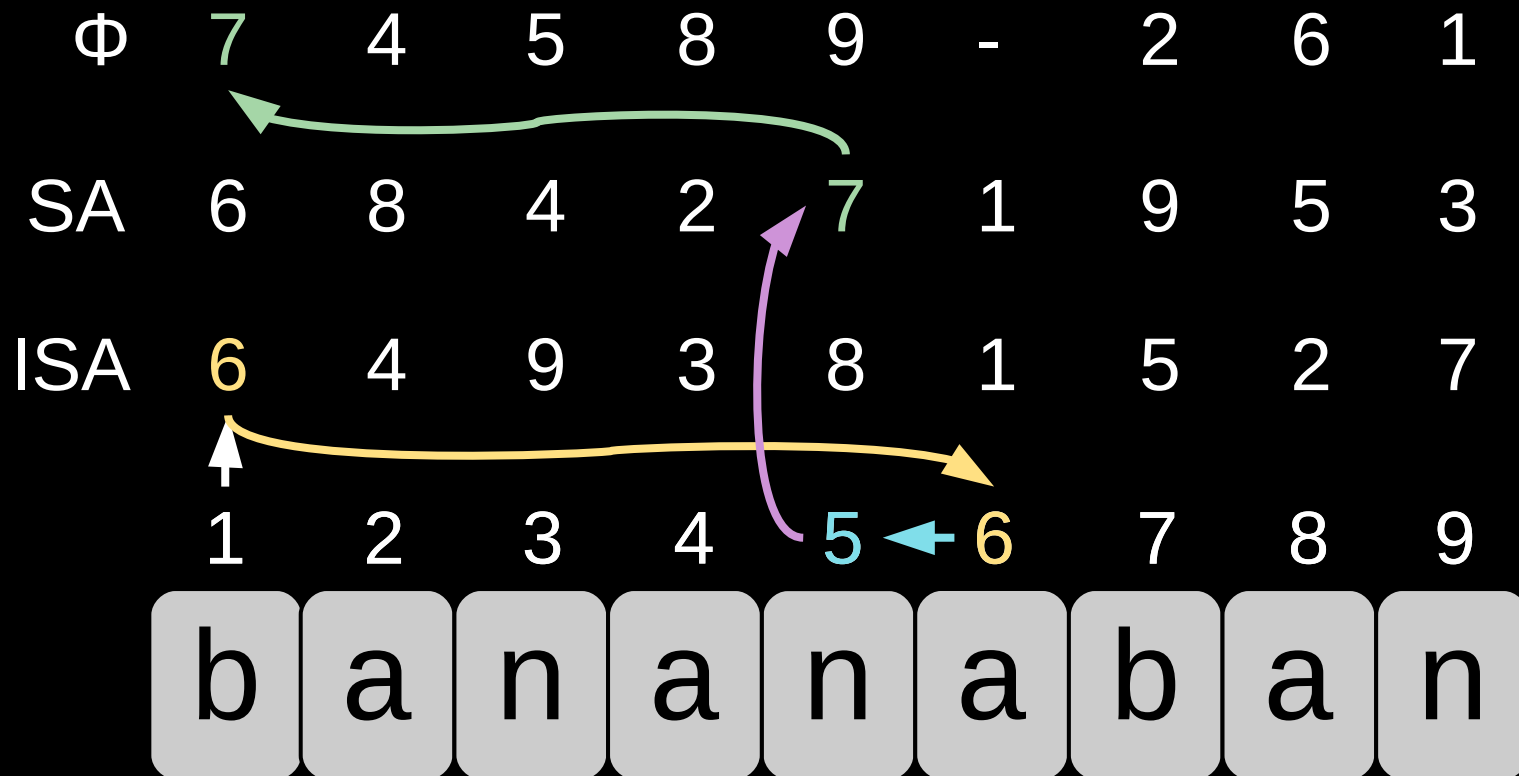
- longest re-occurring substring  $S$
- = longest common prefix (LCP) of two suffixes

## by

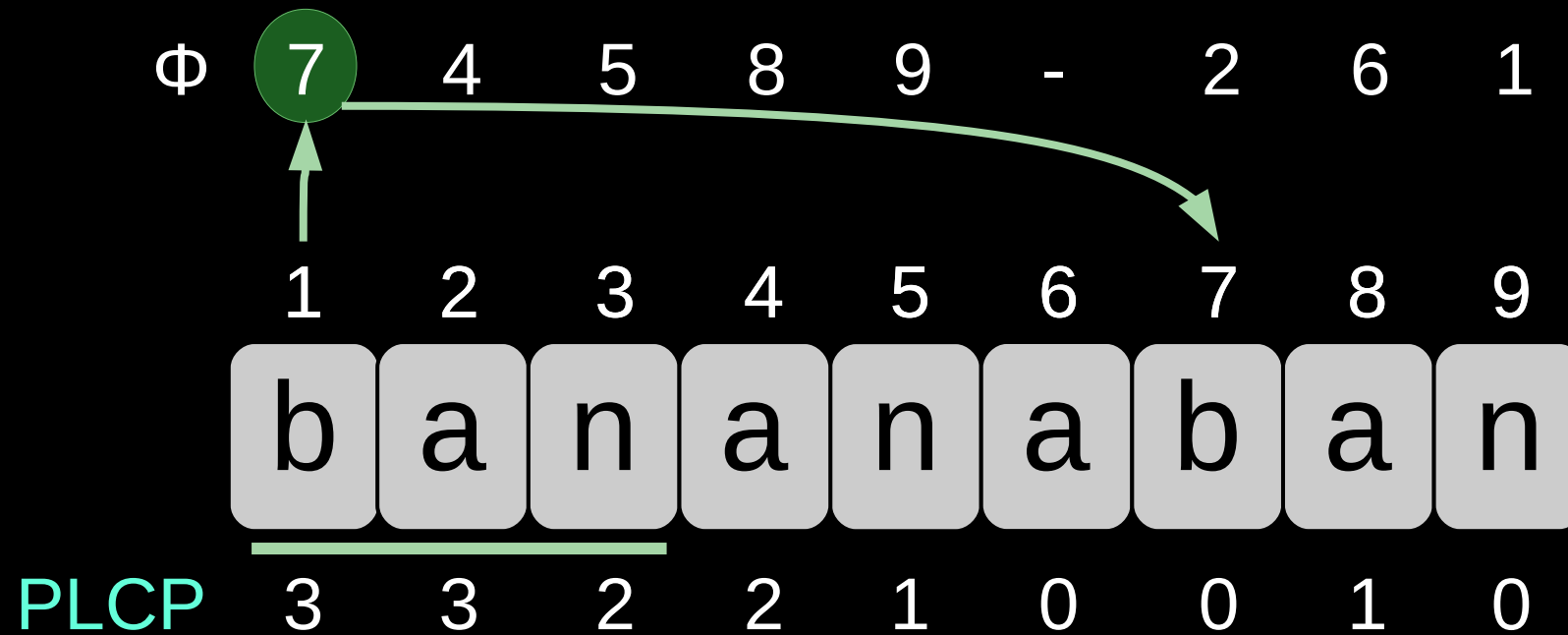
- $|S| = \text{PLCP}[i]$
- $S$  starts at  $\text{argmax } \text{PLCP}[i]$
- $\text{PLCP}[i] = \text{LCP}$  of  $i$ -th suffix and  $\Phi[i]$ -th suffix

# $\Phi$ array

$$\Phi[i] := SA[ISA[i] - 1]$$



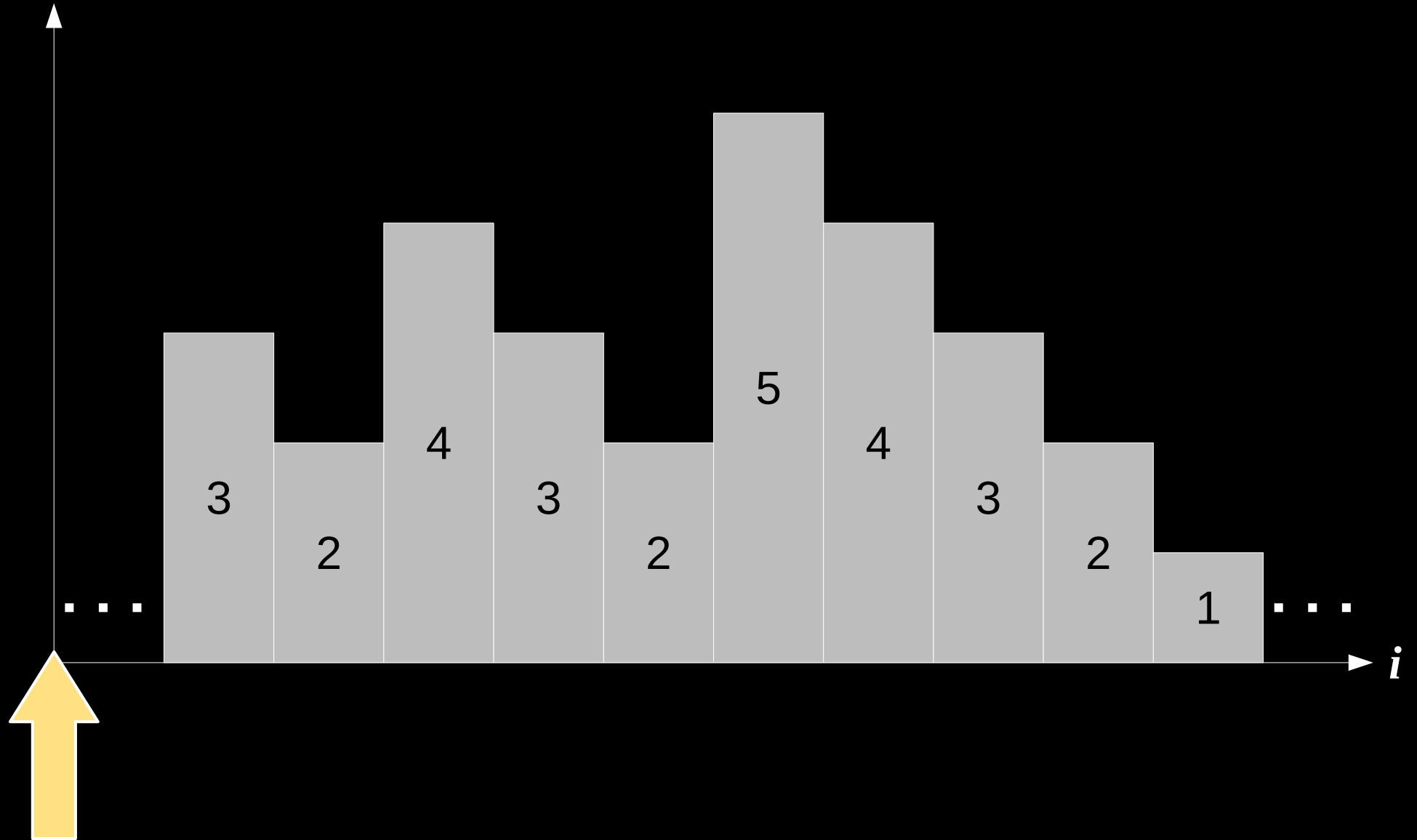
# main idea



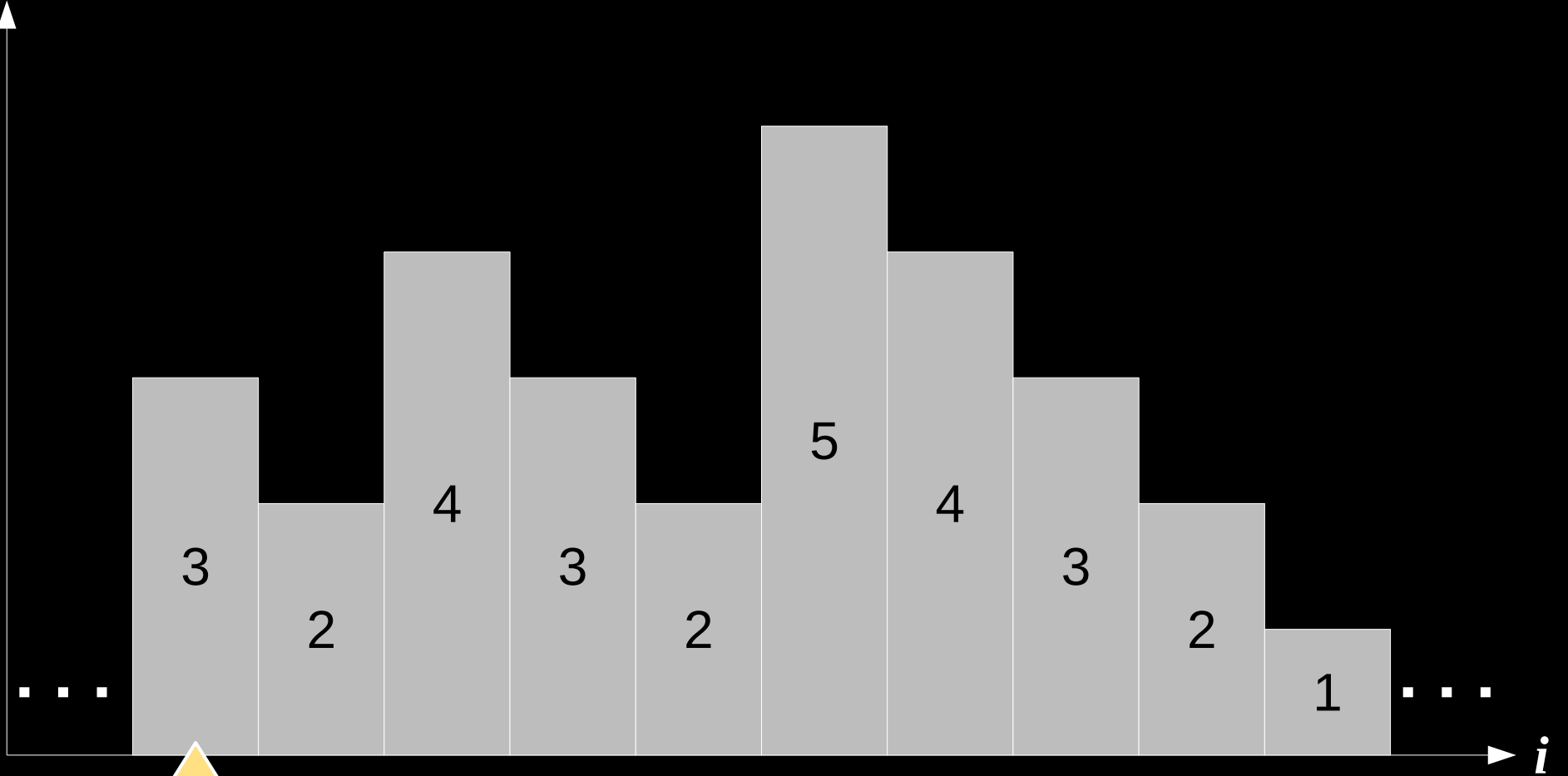
- **PLCP**: length of factor
- $\Phi$ : reference

# computation

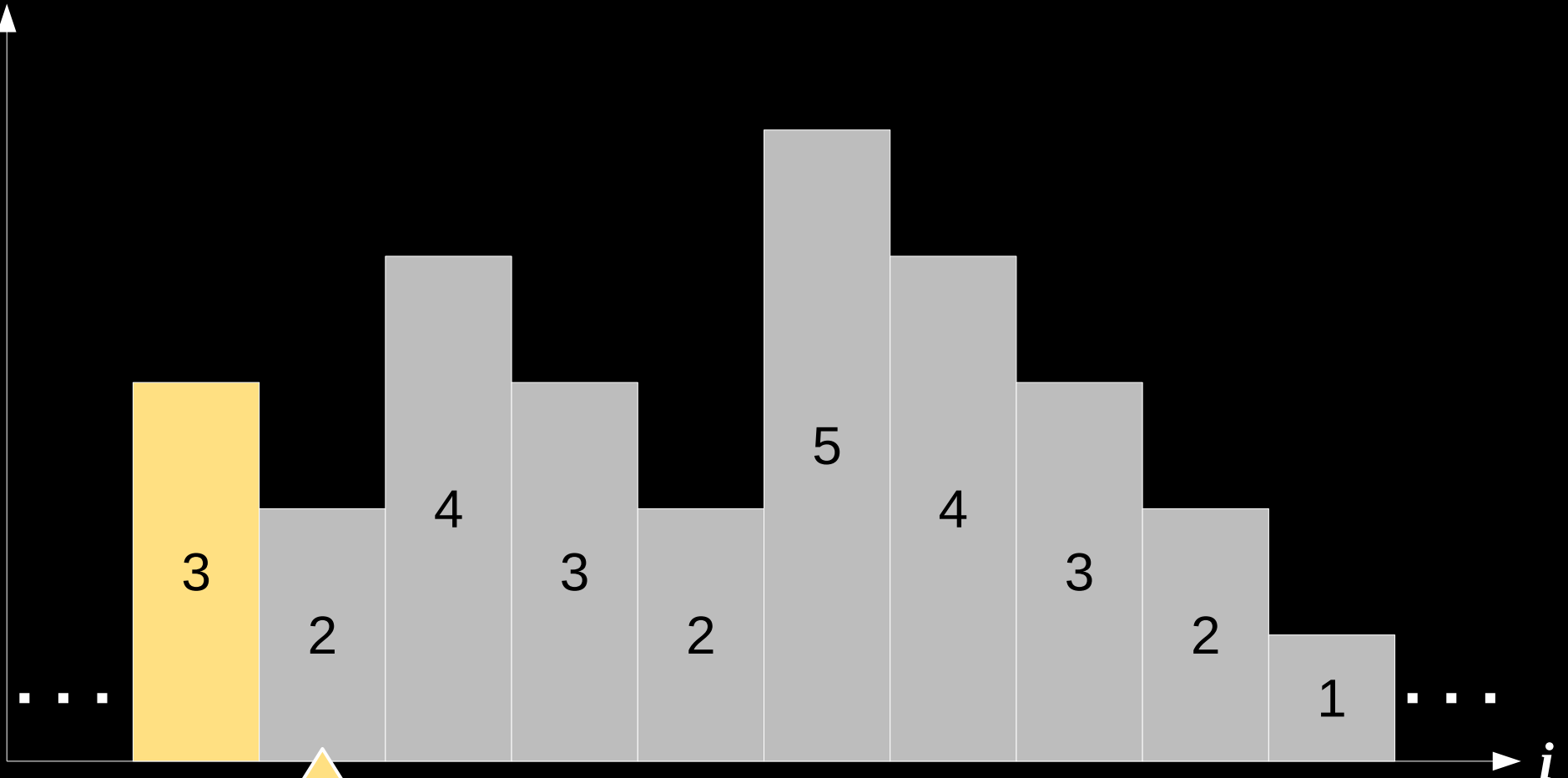
PLCP[ $i$ ]



PLCP[ $i$ ]

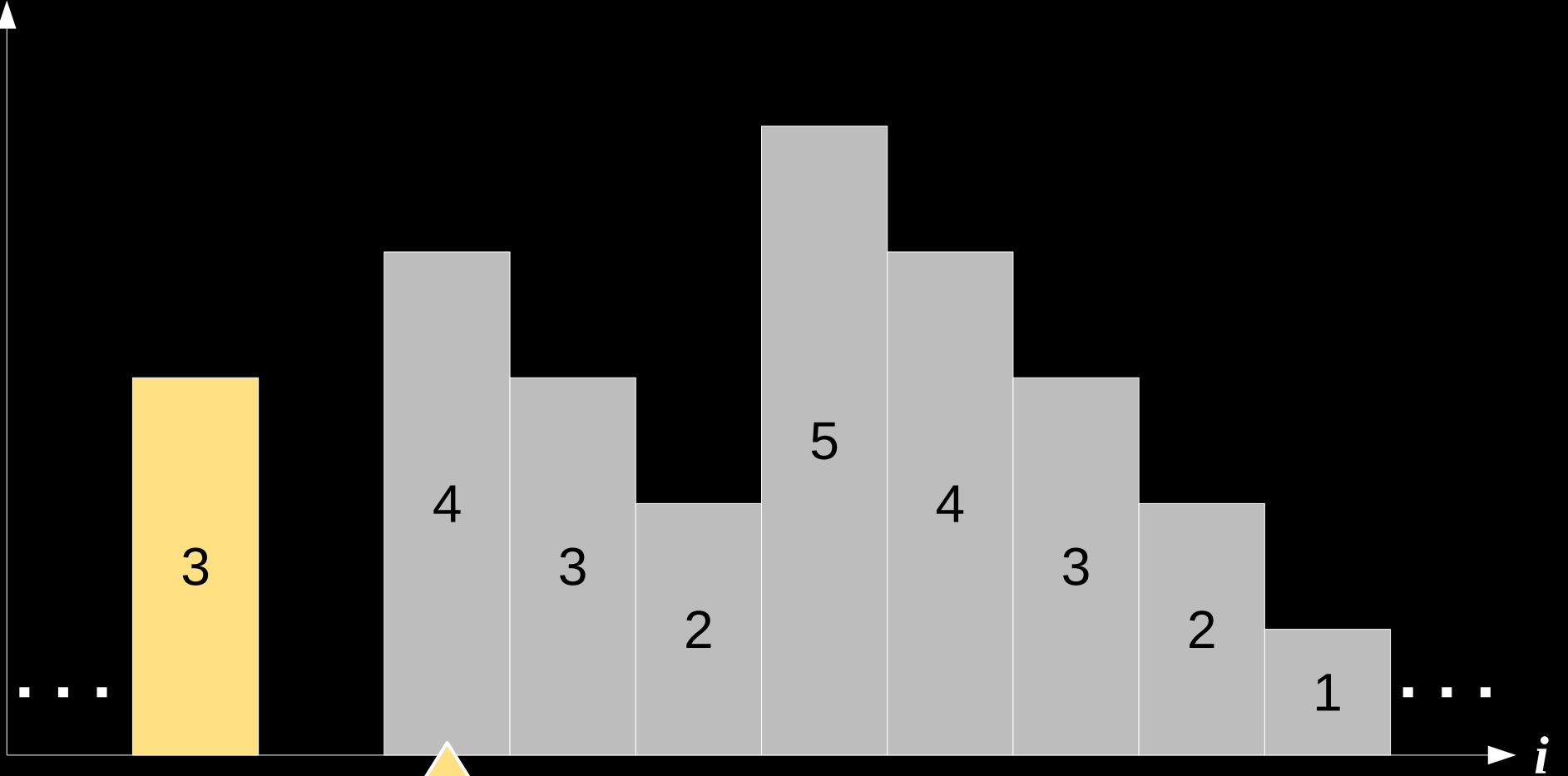


PLCP[ $i$ ]

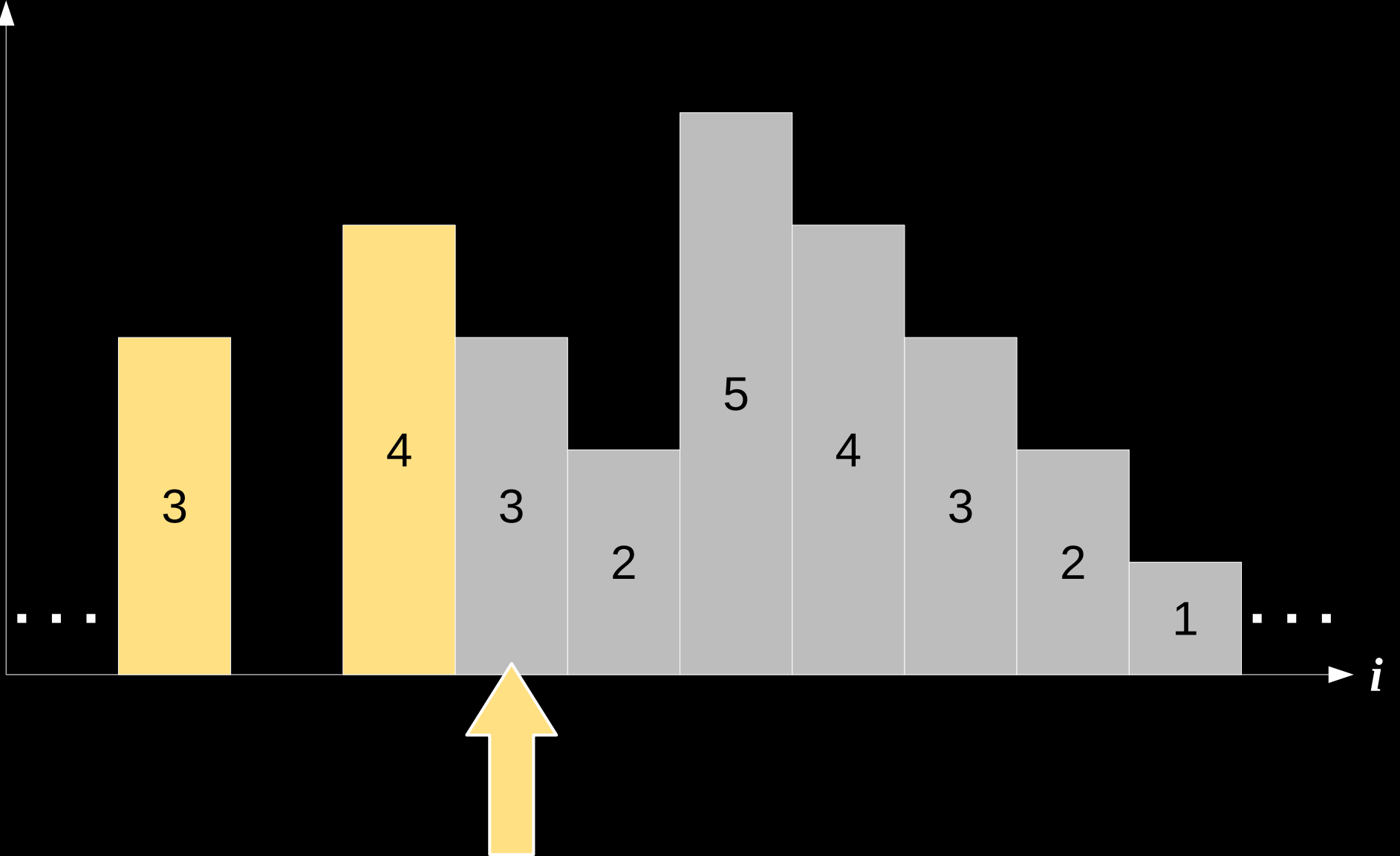


discard, since previous position's PLCP is higher

PLCP[ $i$ ]

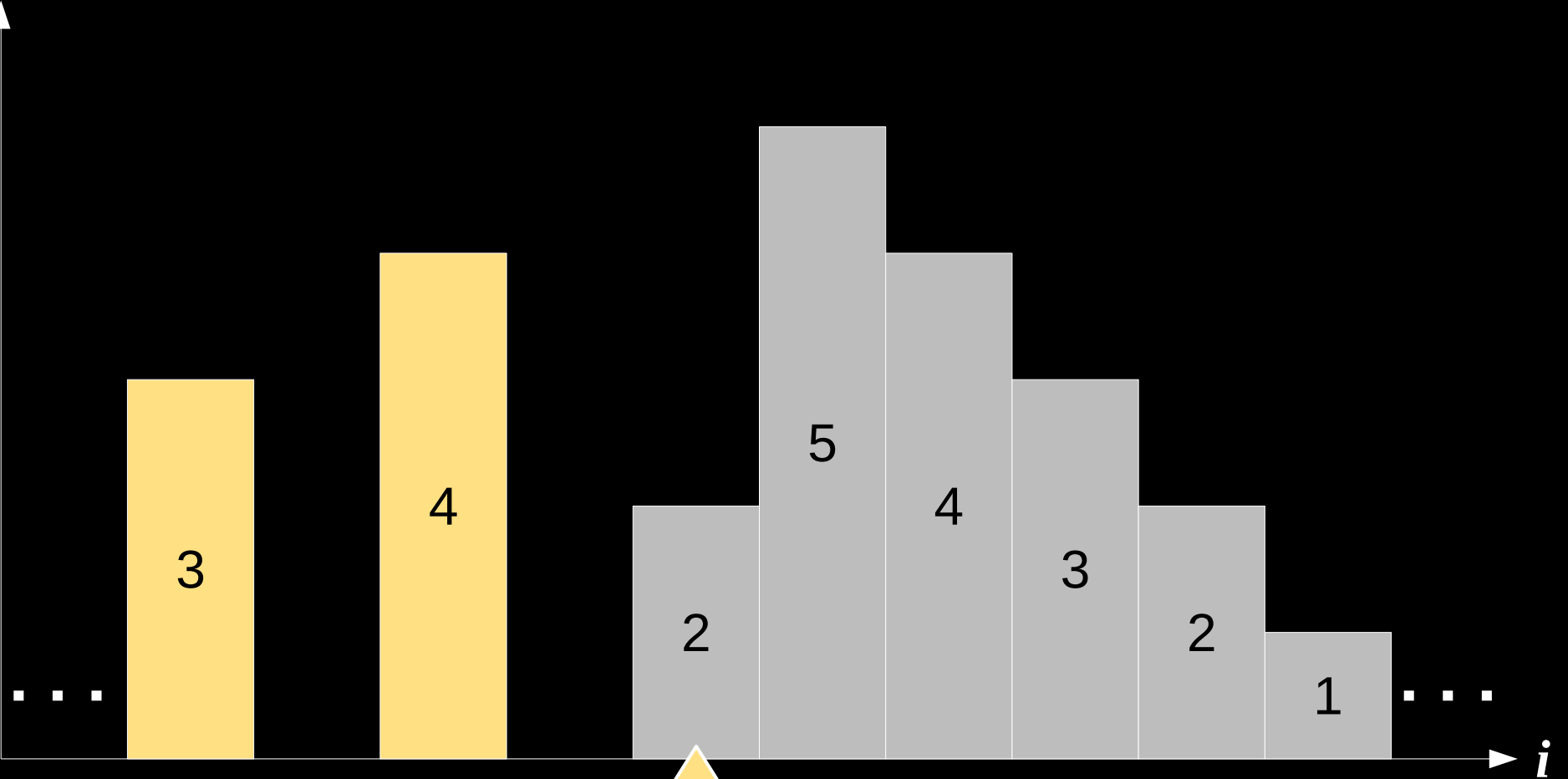


PLCP[ $i$ ]

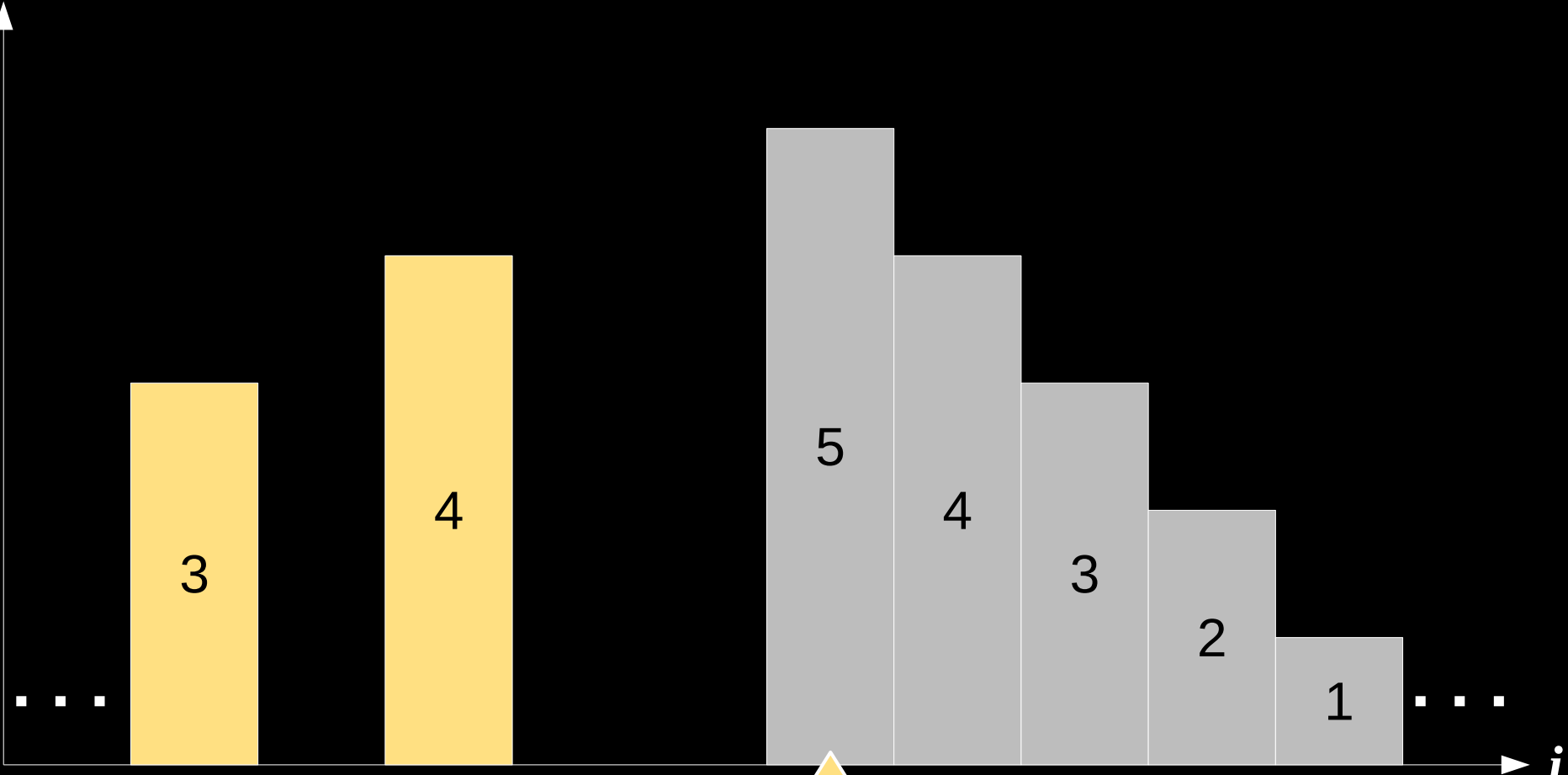




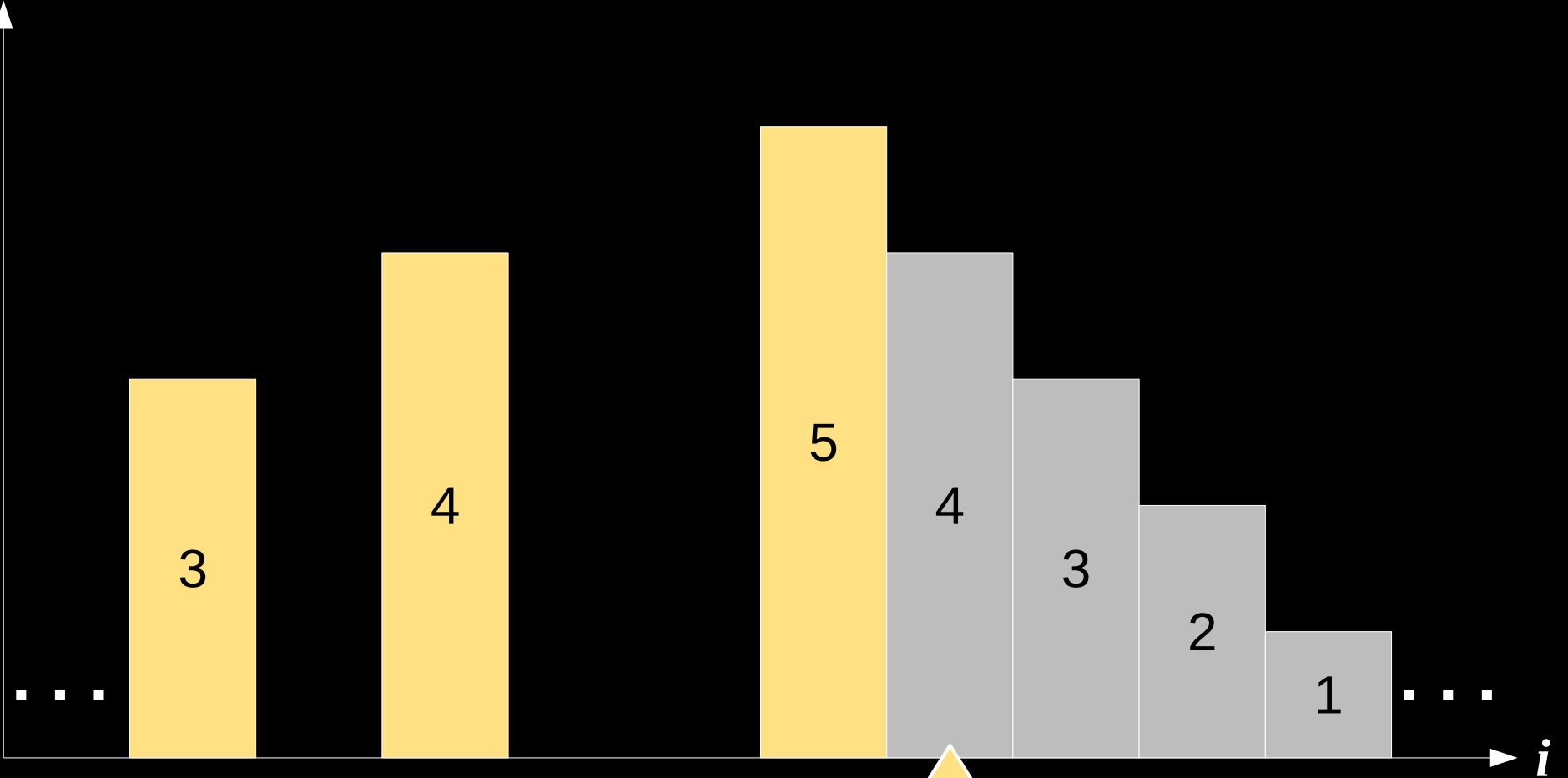
PLCP[*i*]



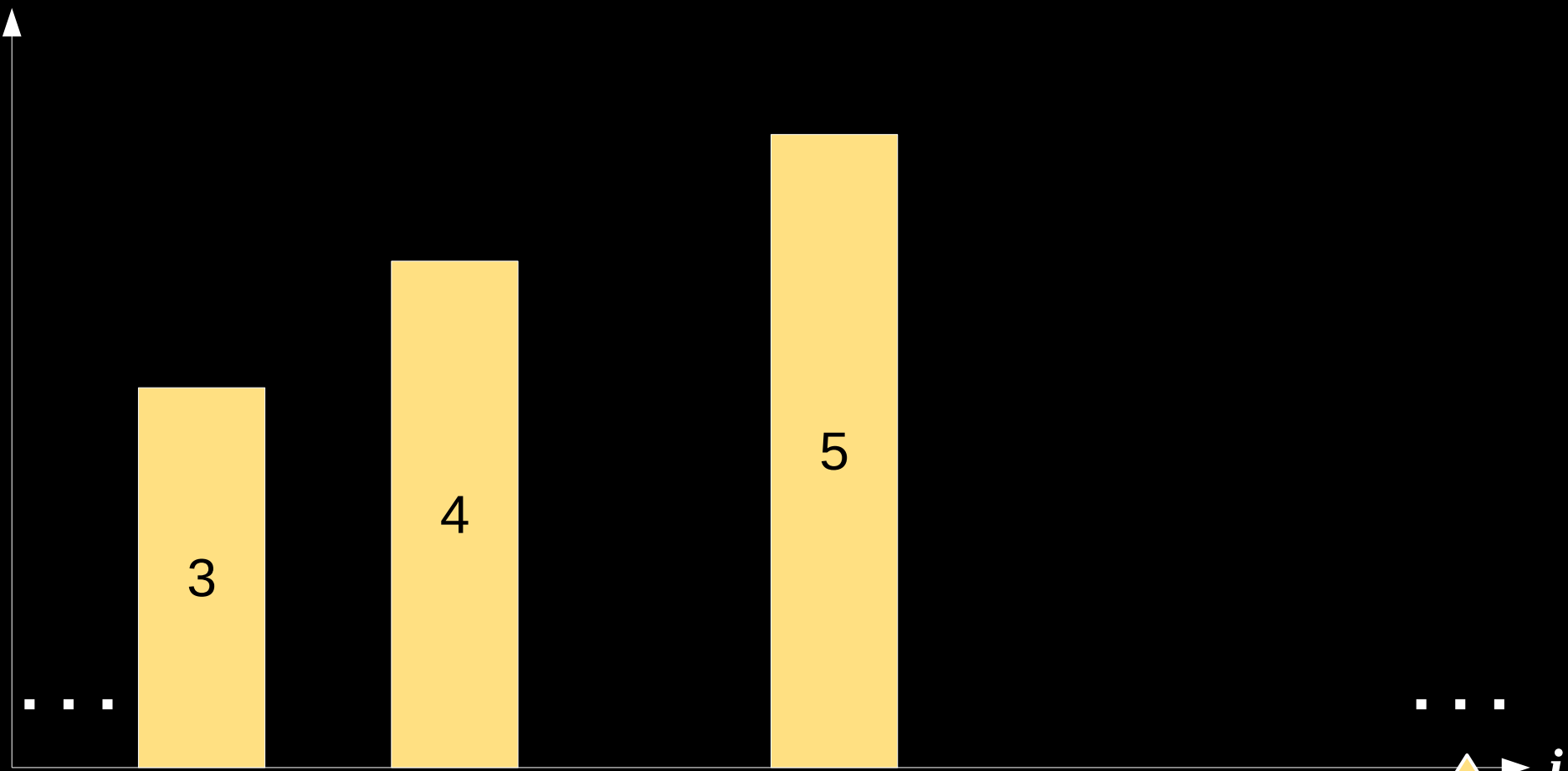
PLCP[ $i$ ]



PLCP[ $i$ ]



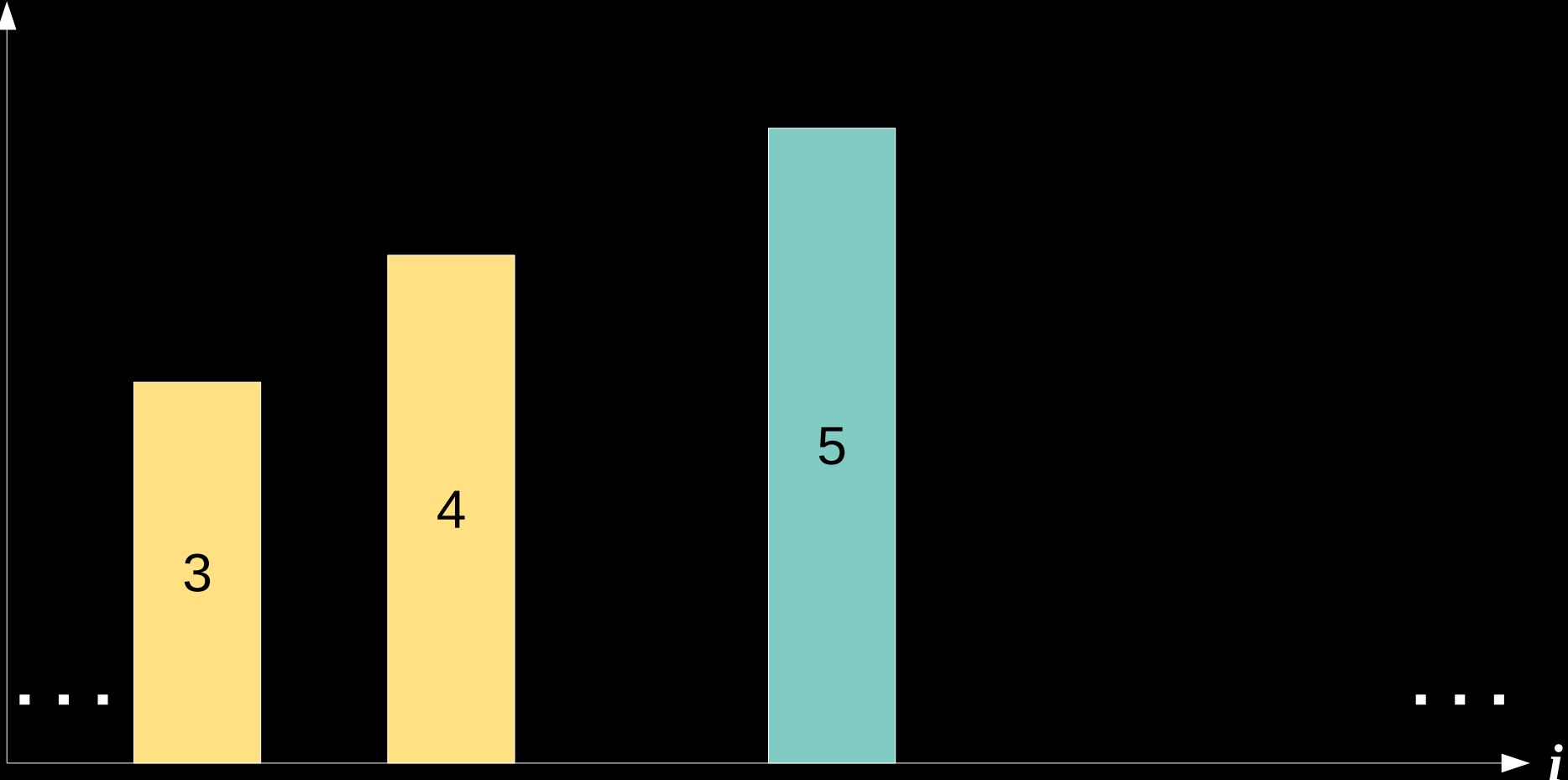
PLCP[*i*]



no peak > 5

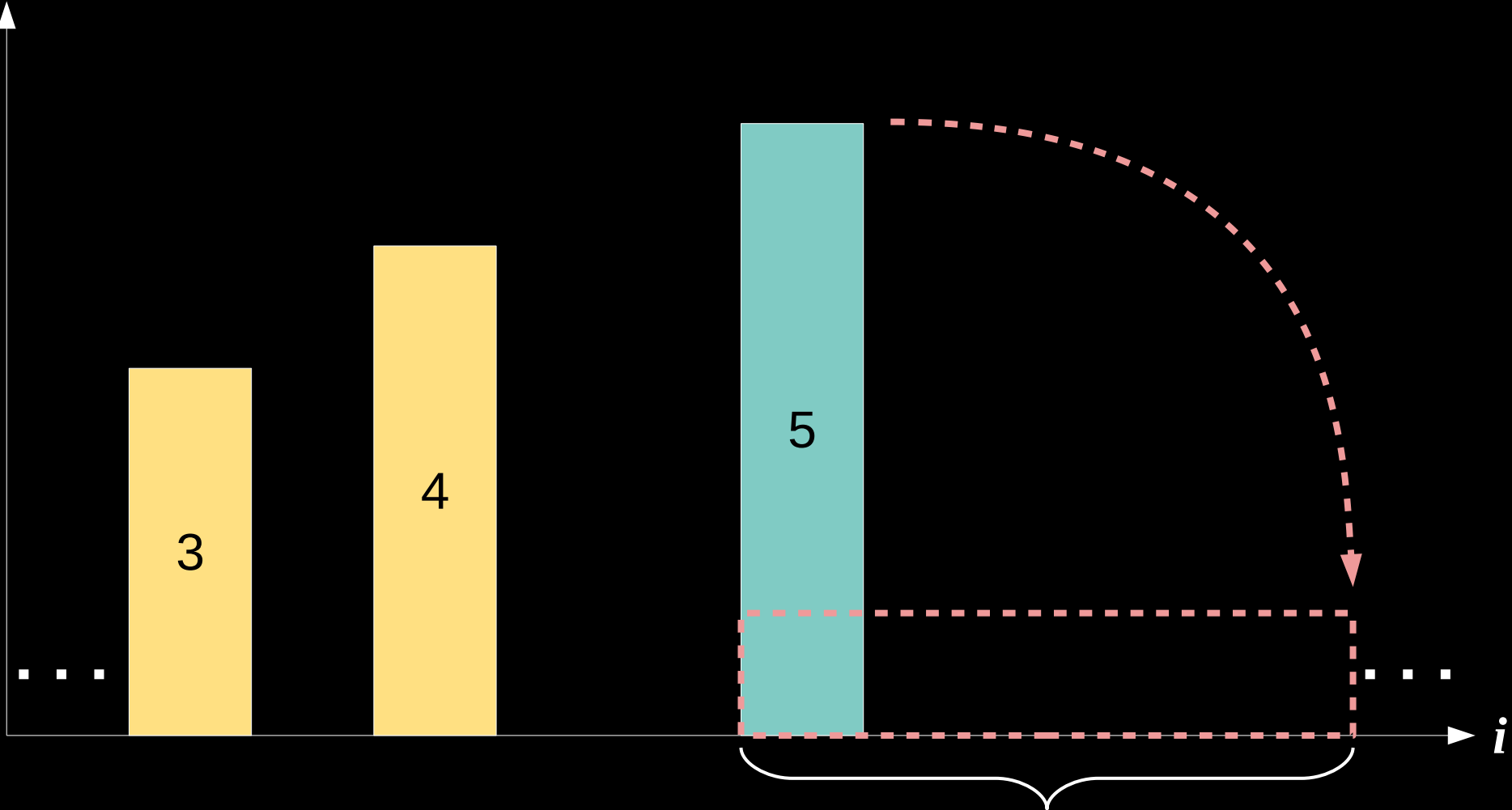


PLCP[ $i$ ]



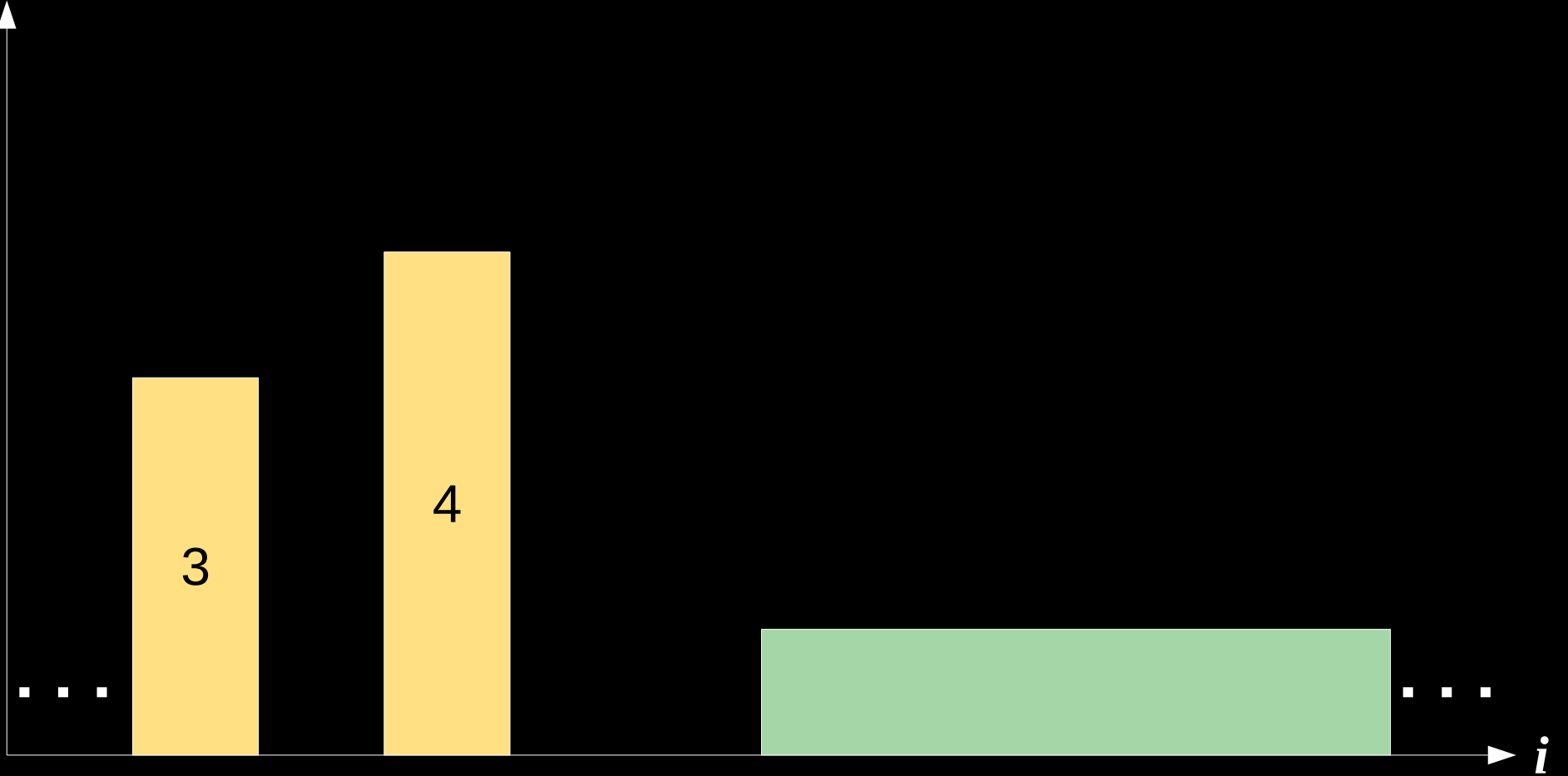
*maximum peak*

PLCP[ $i$ ]

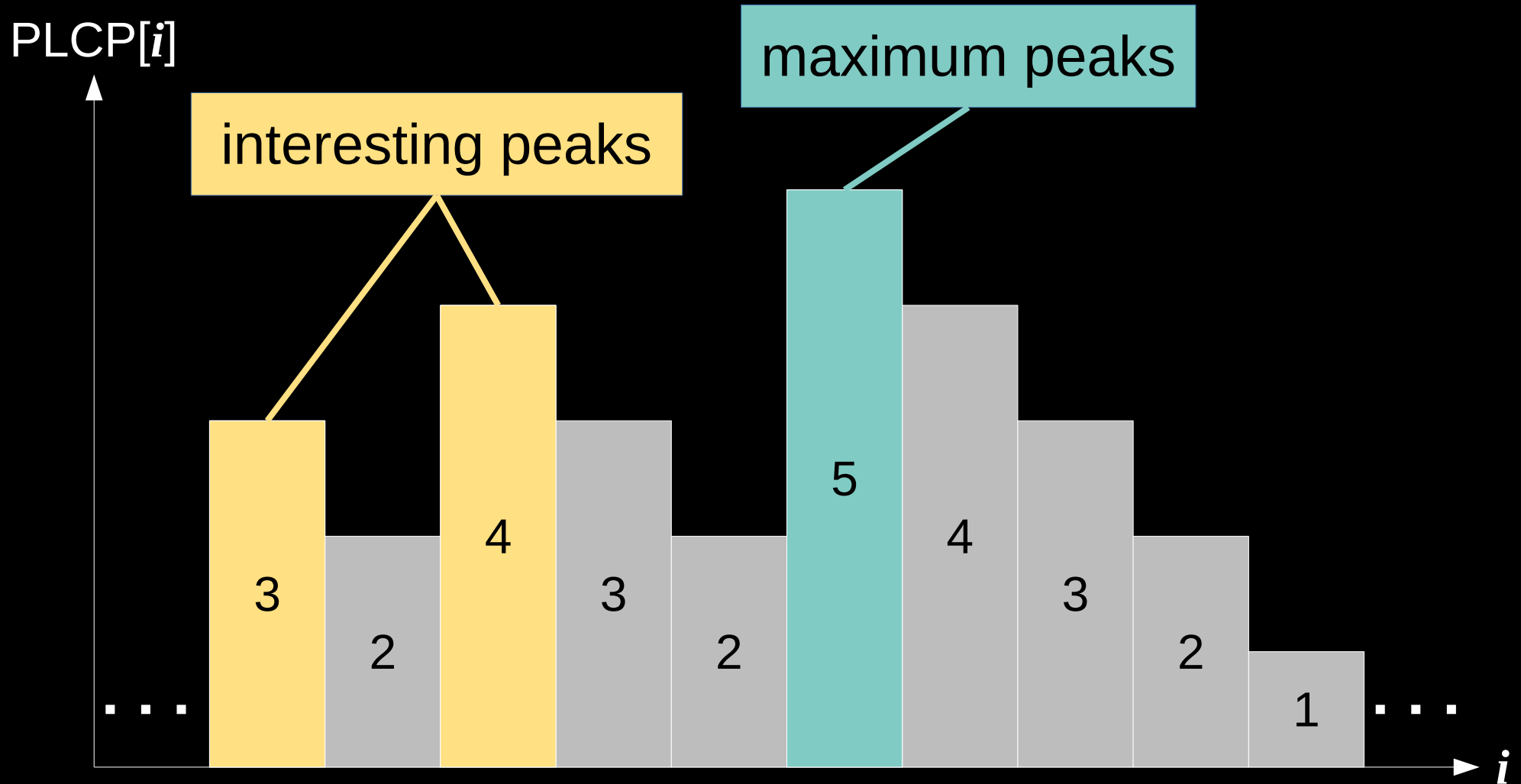


*factor of length 5*

PLCP[ $i$ ]

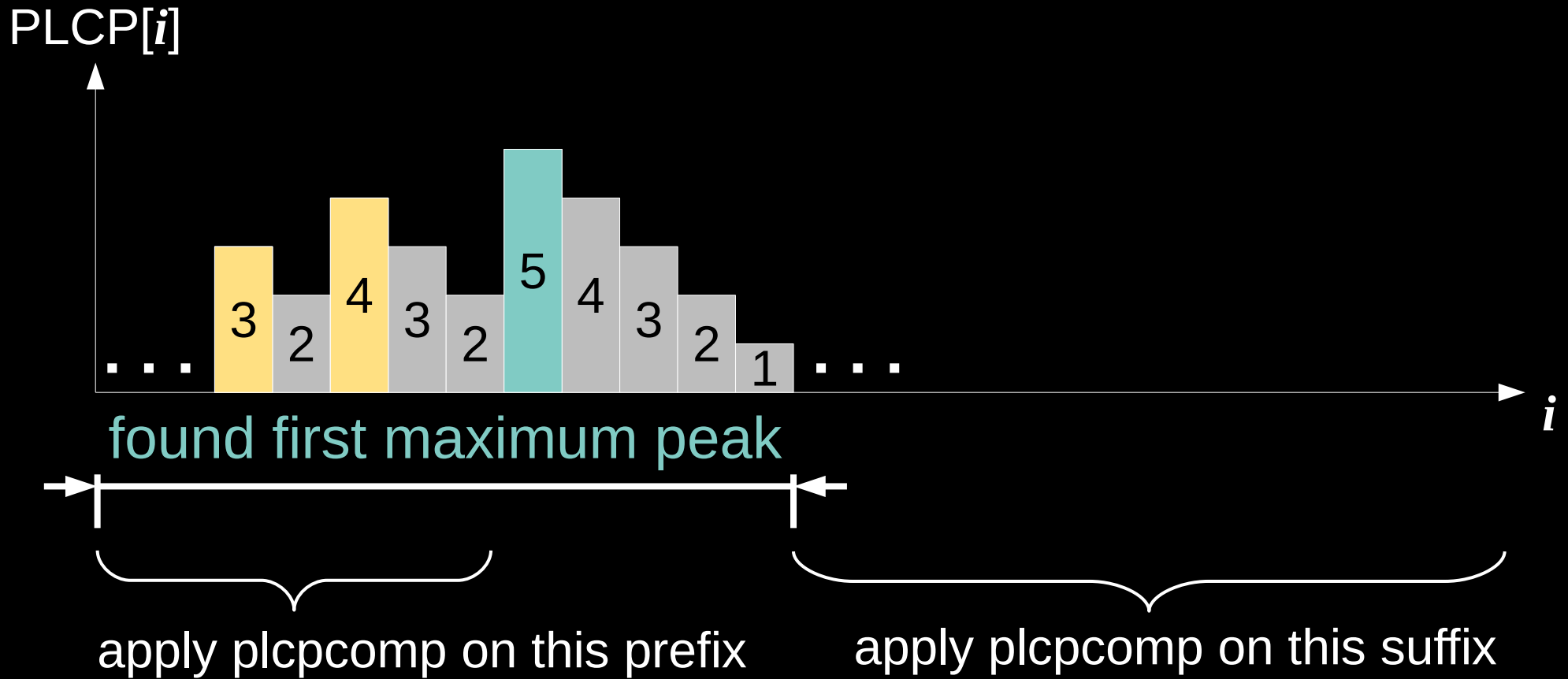



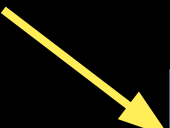
# 2 definitions





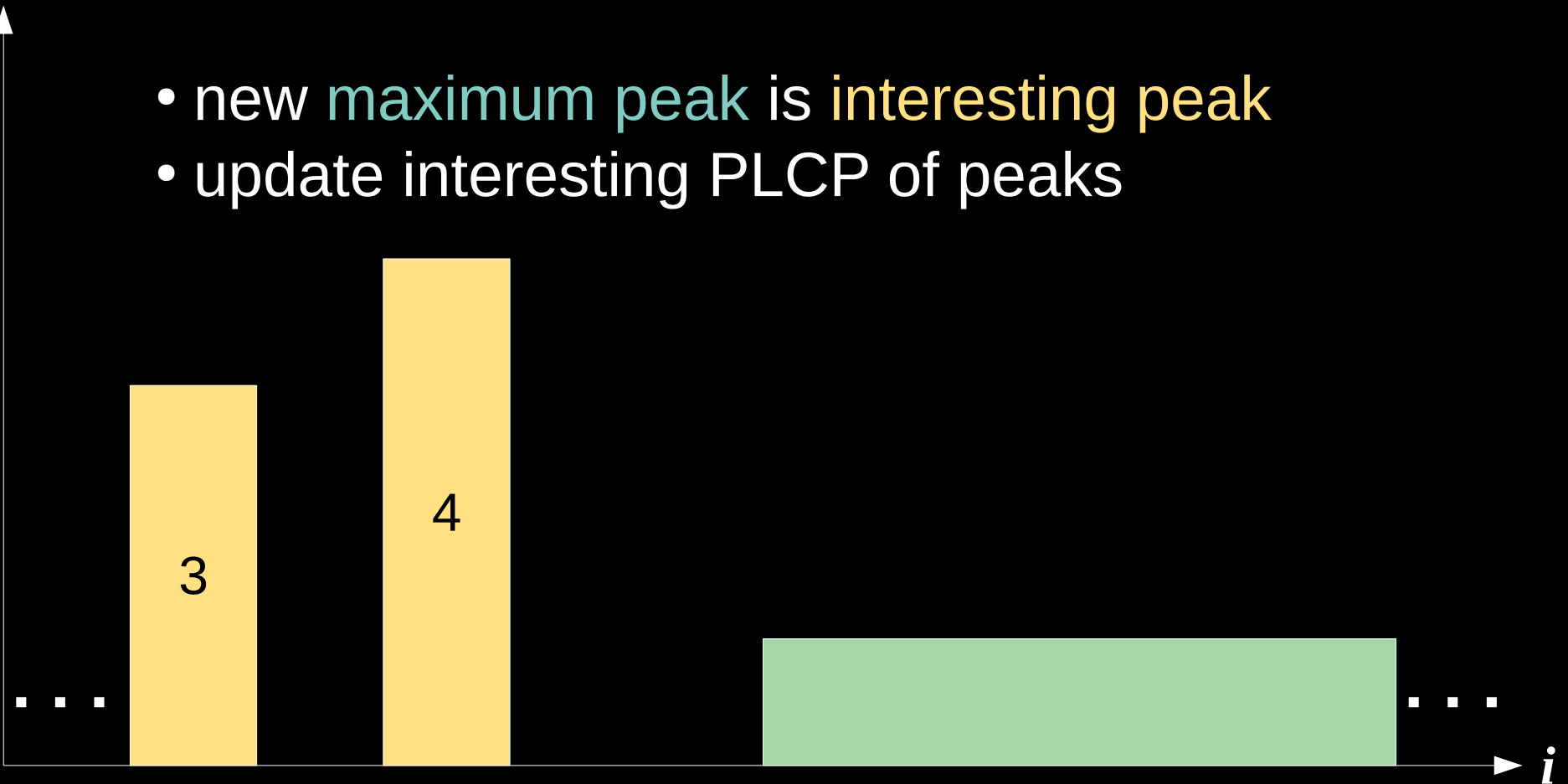
# recurse



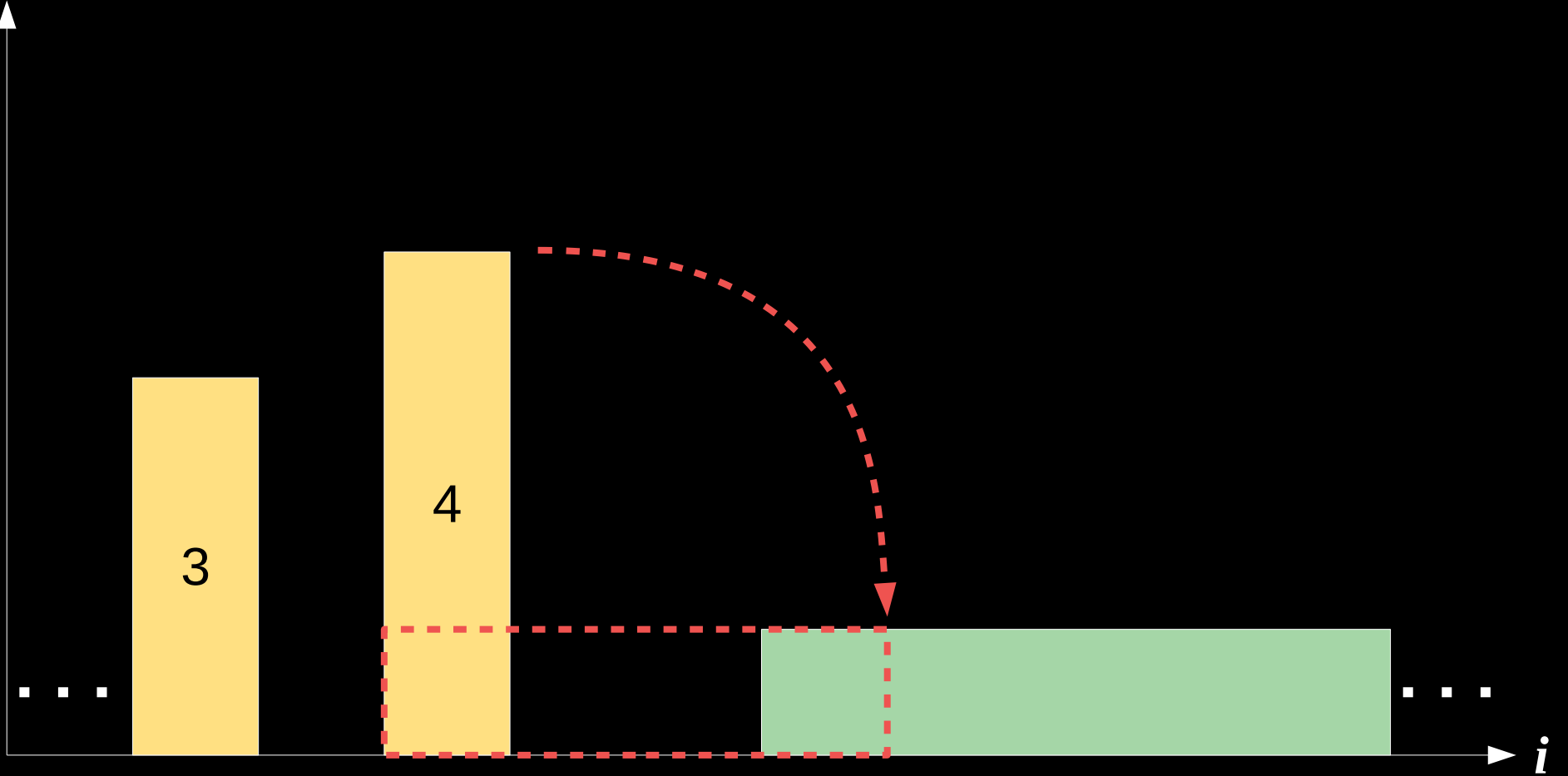
linear time?   only examine interesting peaks!

PLCP[ $i$ ]

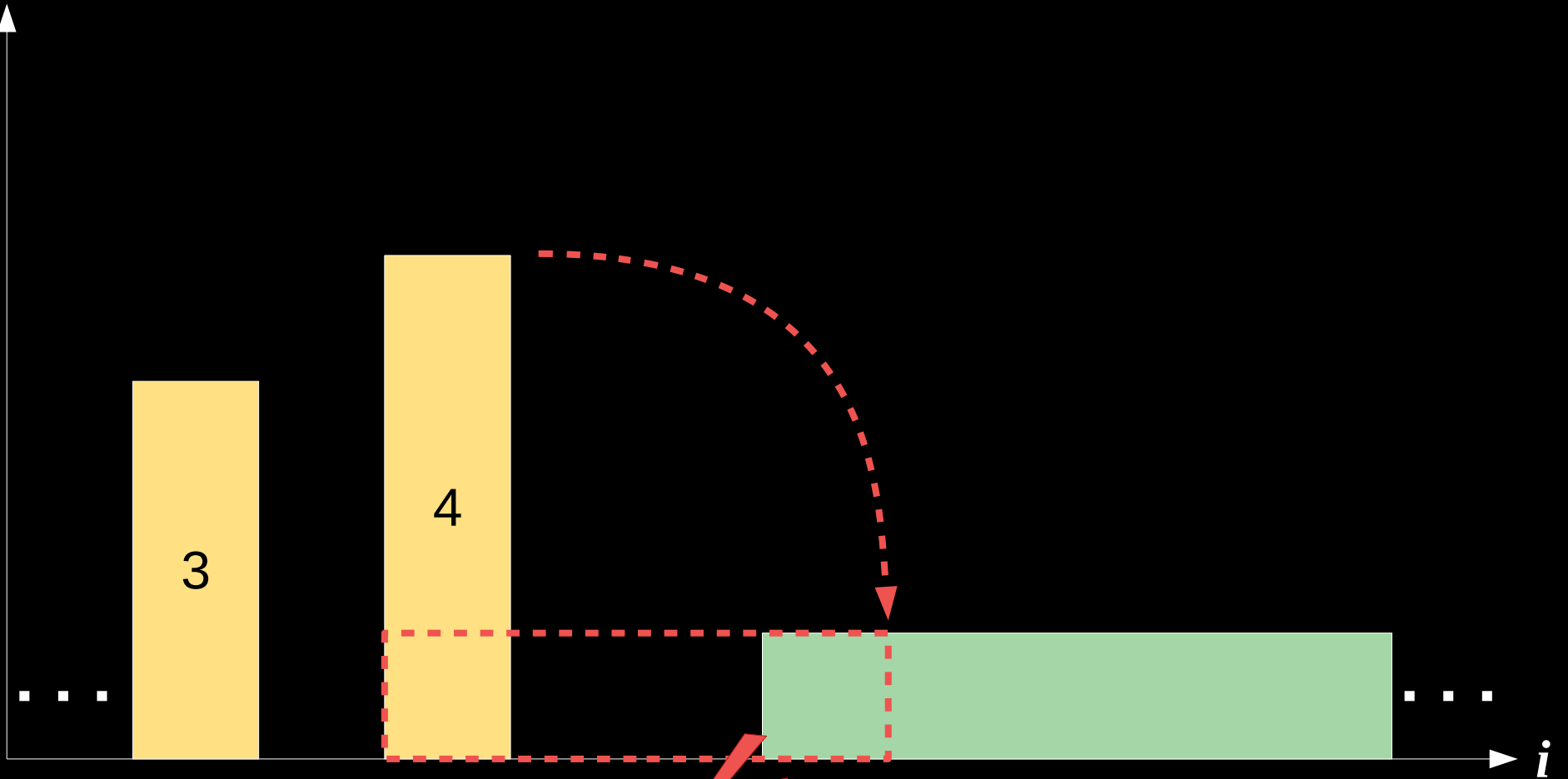
- new maximum peak is interesting peak
- update interesting PLCP of peaks



PLCP[ $i$ ]

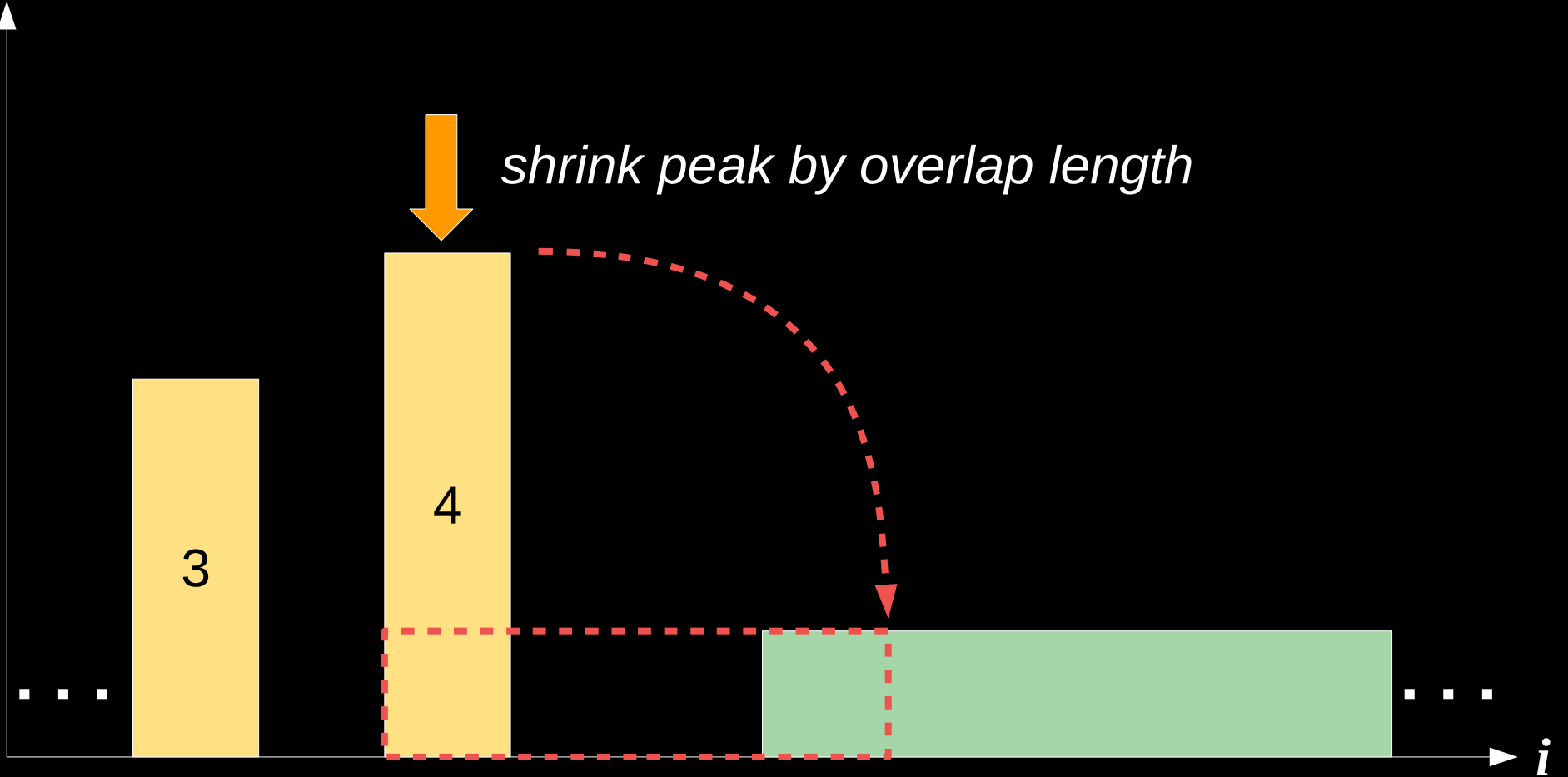


PLCP[ $i$ ]

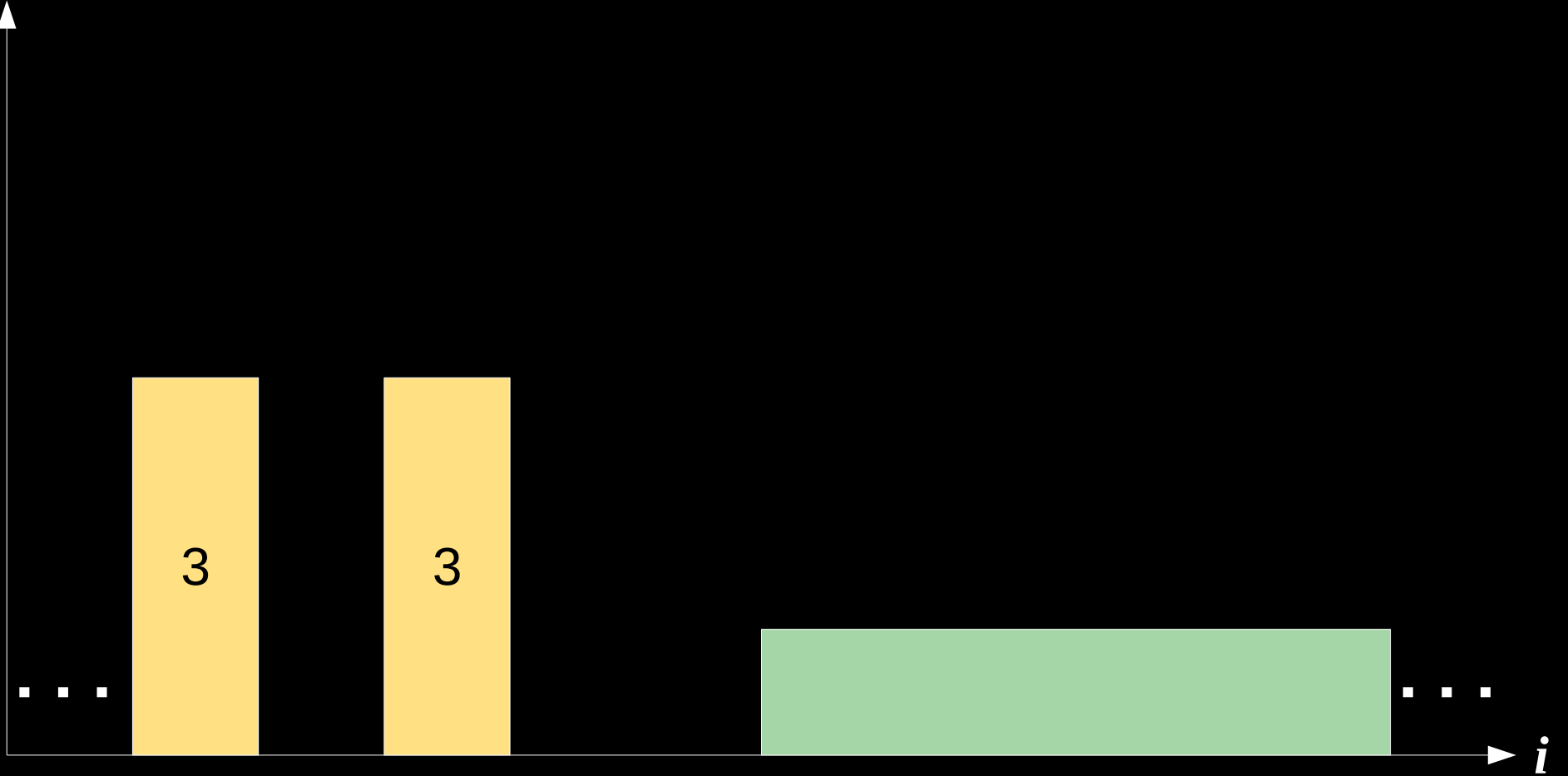


*left overlap of length 1*

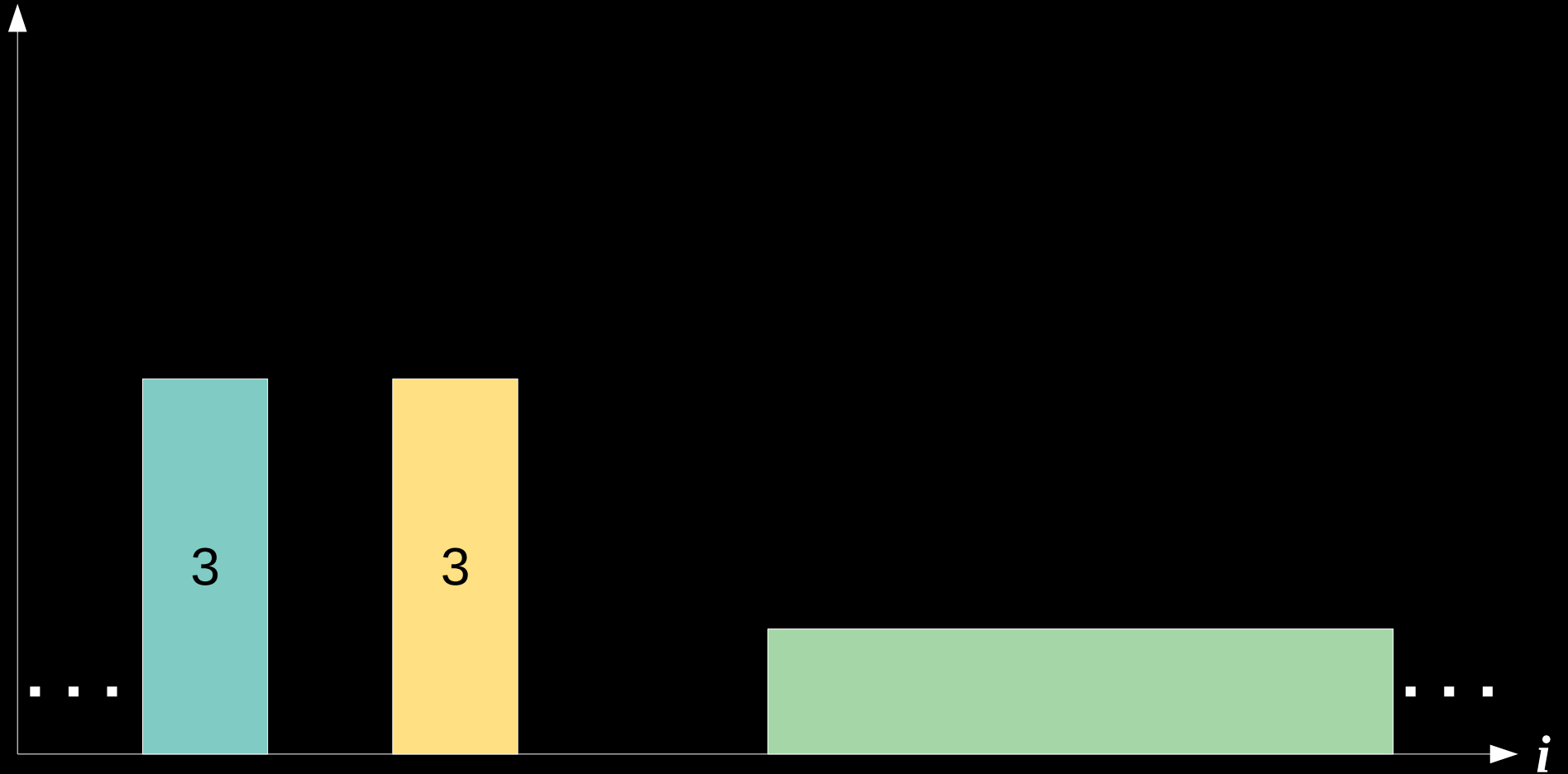
PLCP[ $i$ ]



PLCP[*i*]

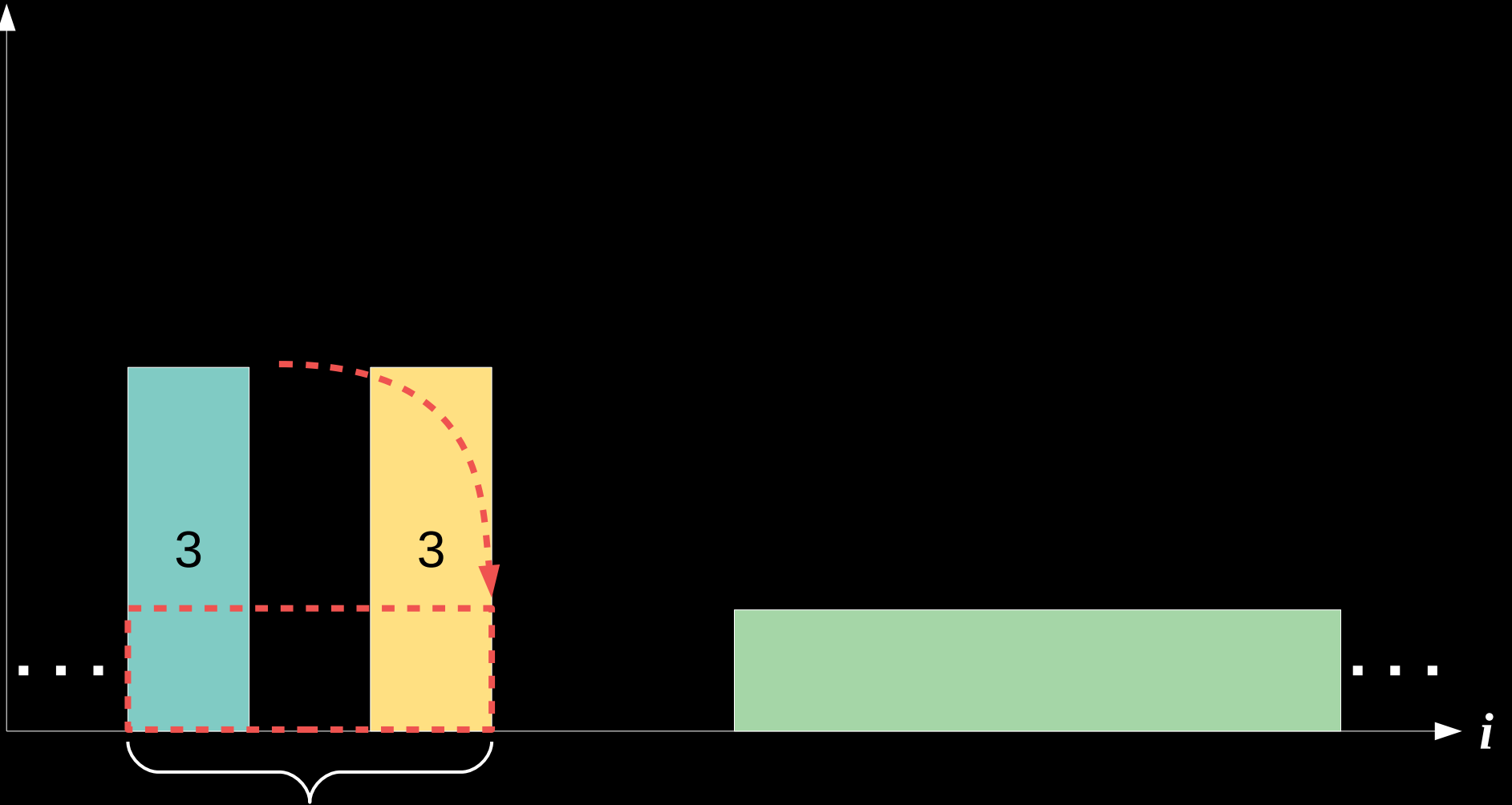


PLCP[ $i$ ]



*new highest peak*

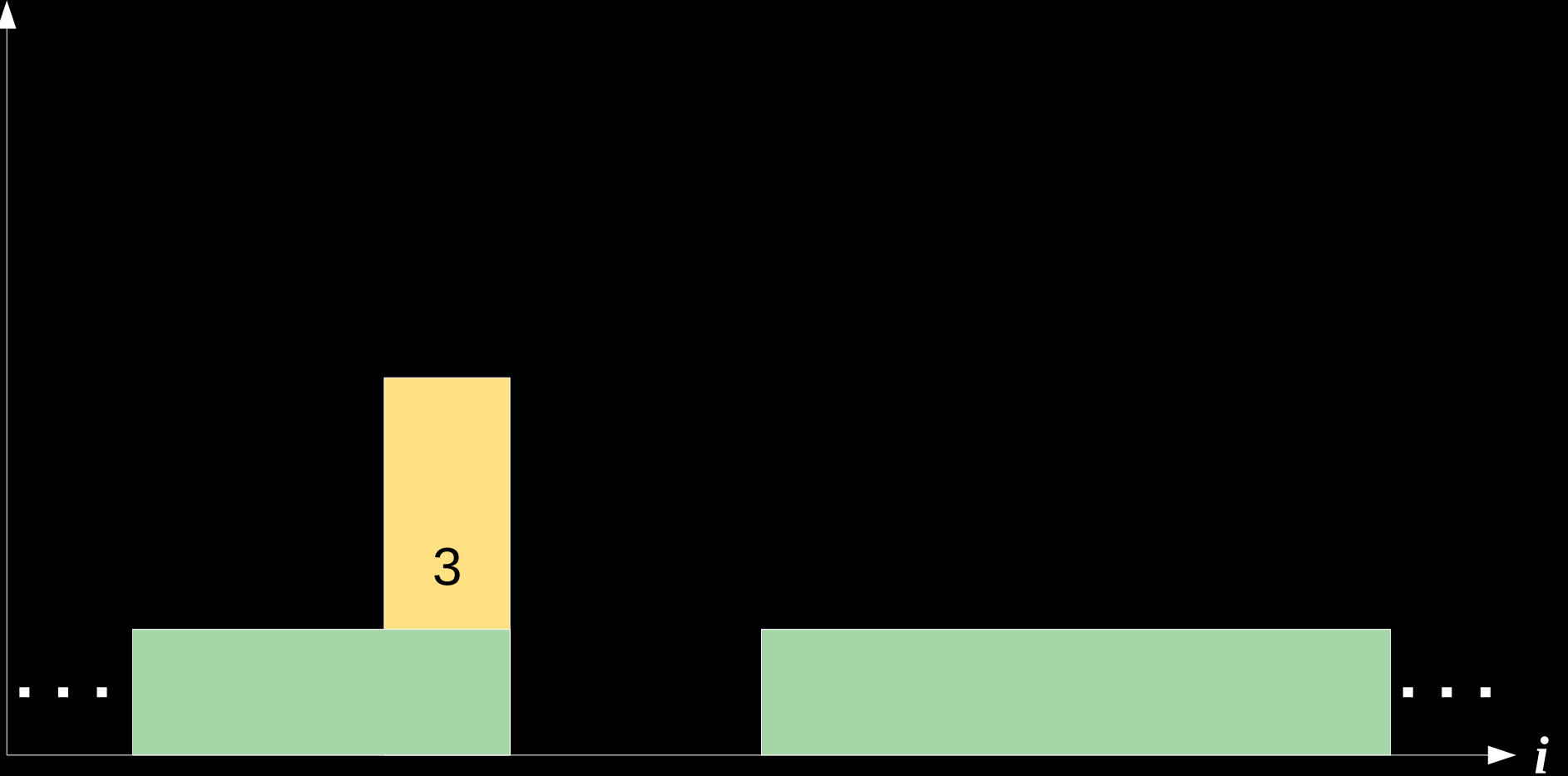
PLCP[ $i$ ]



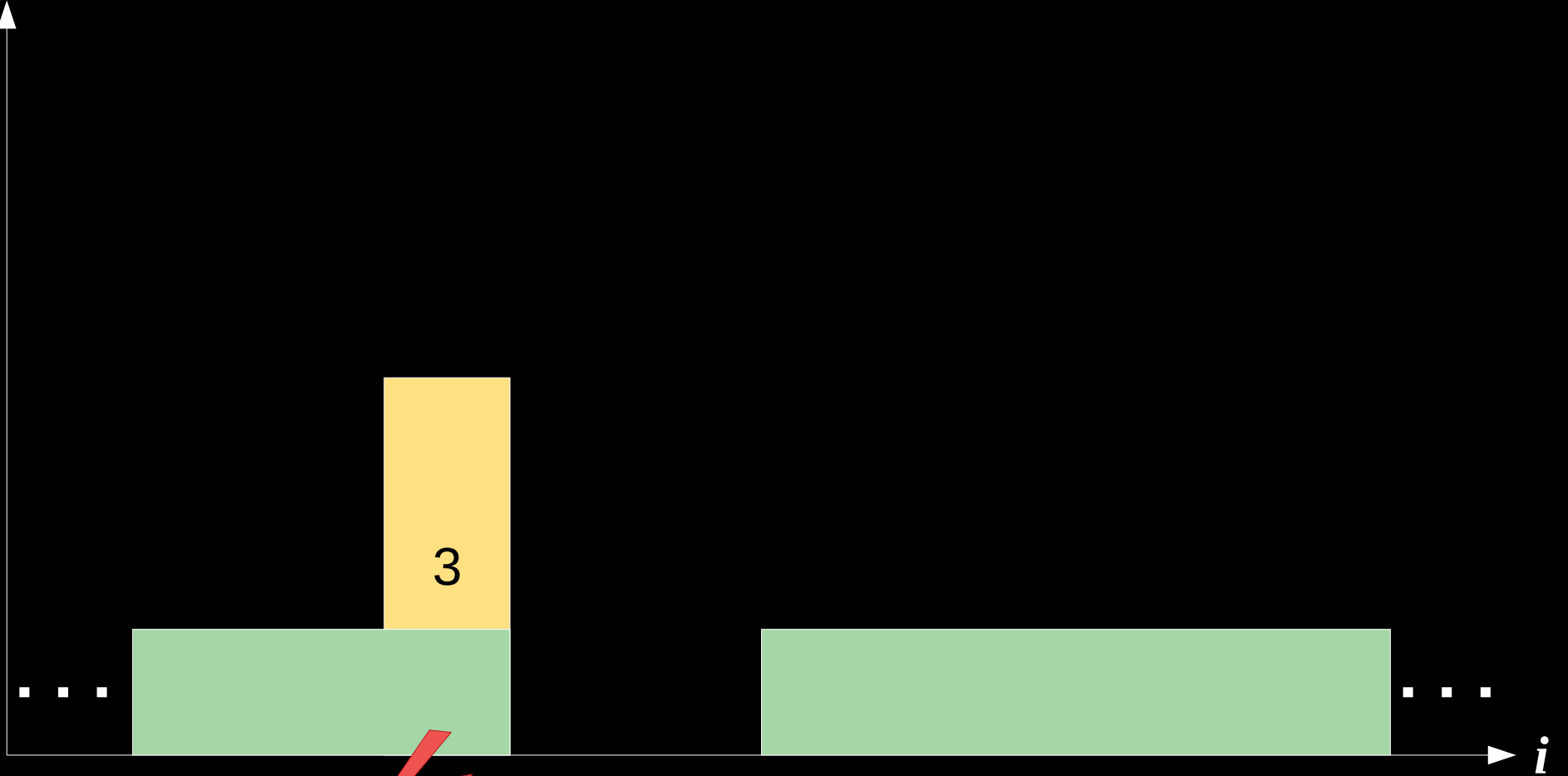
*factor of length 3*



PLCP[ $i$ ]

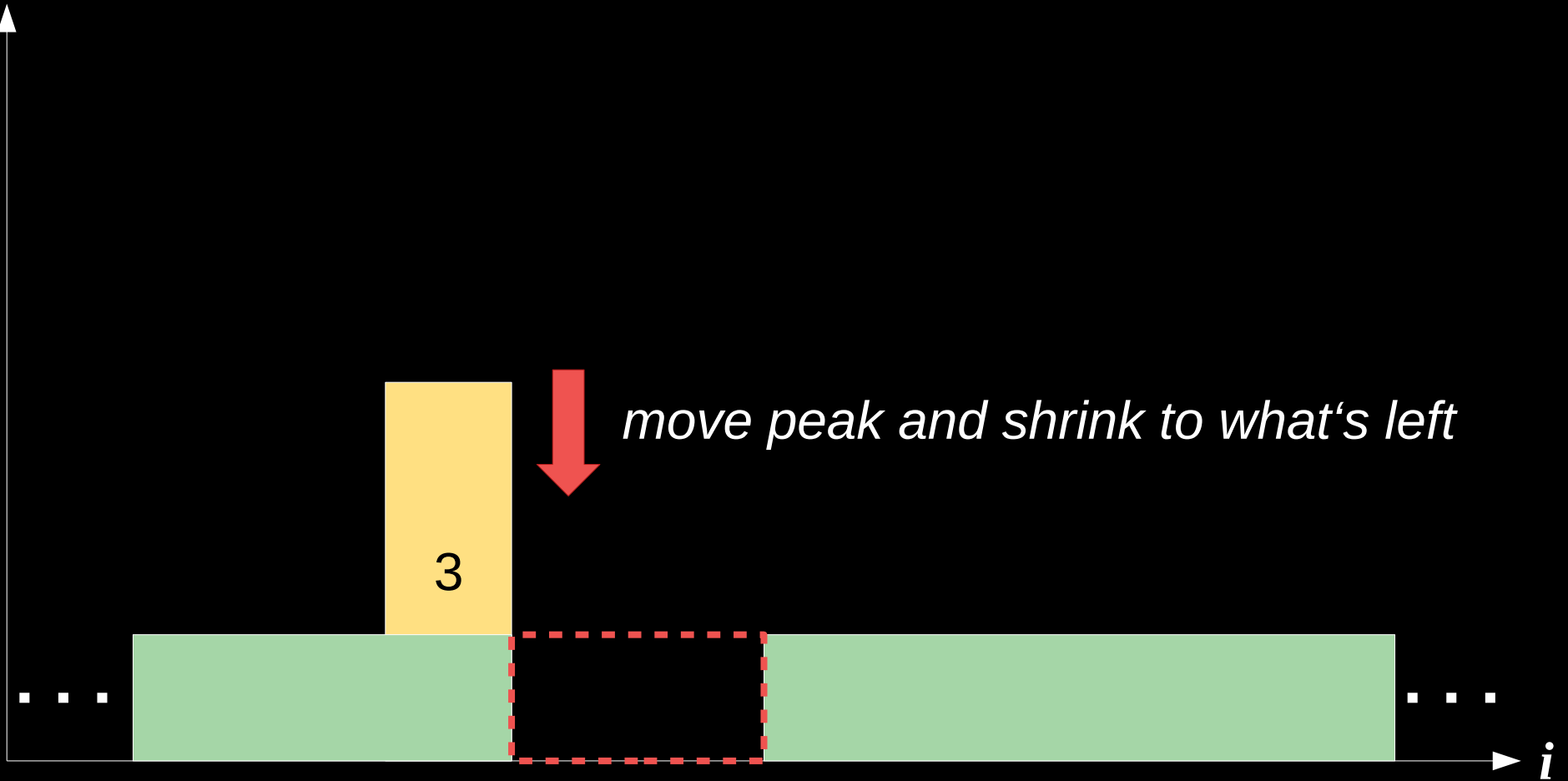


PLCP[ $i$ ]

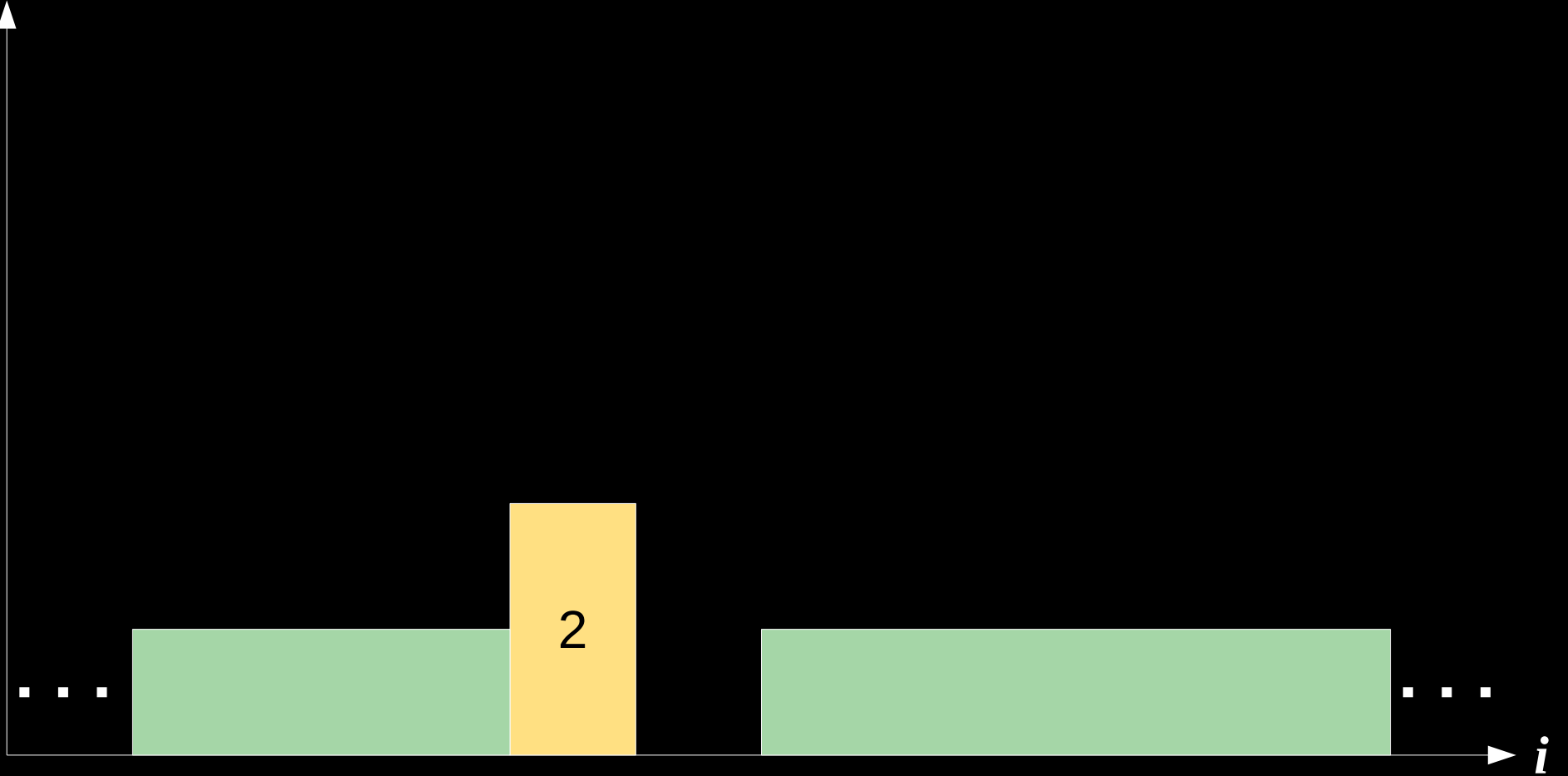


 *right overlap*

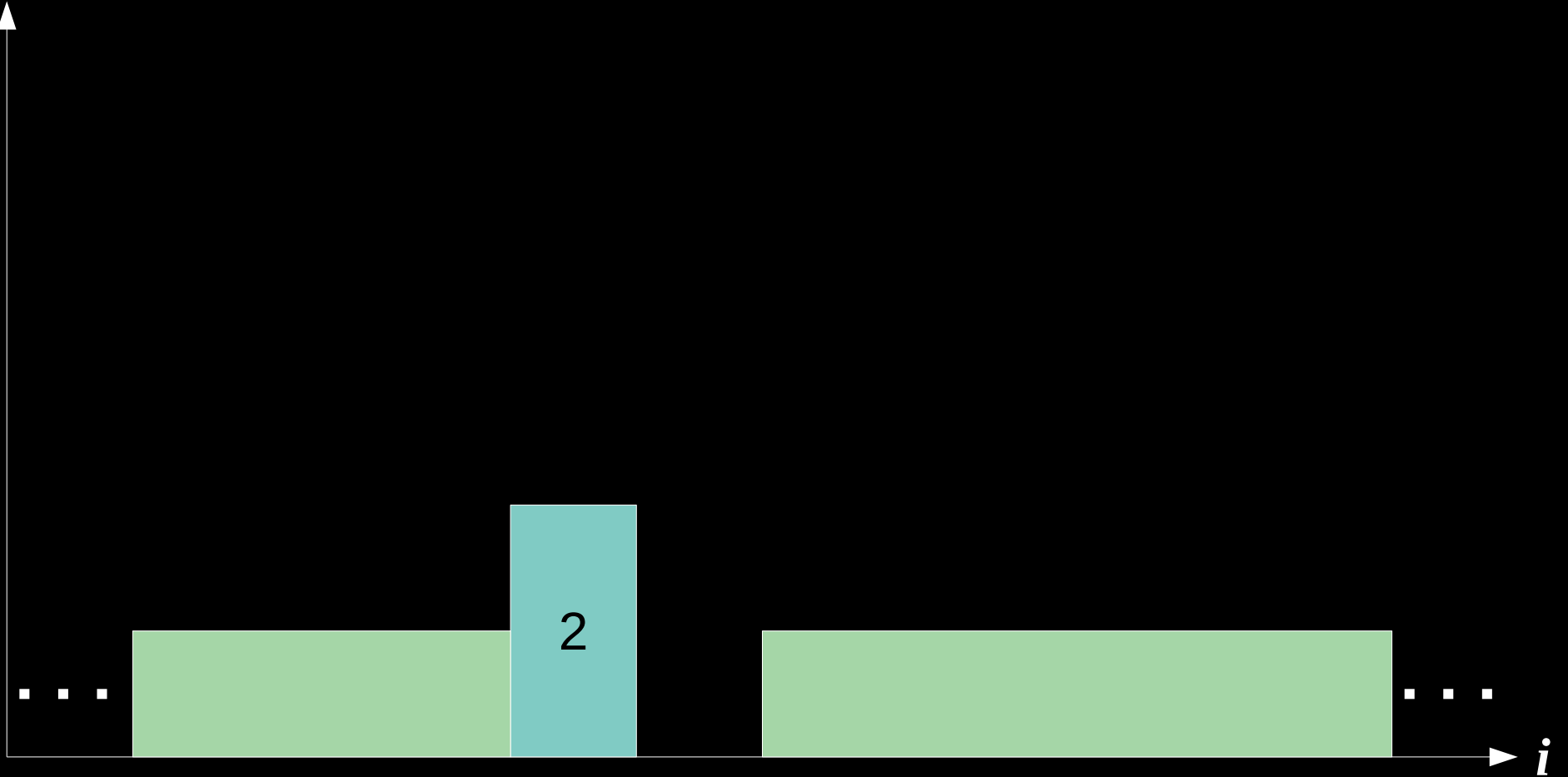
PLCP[ $i$ ]



PLCP[*i*]

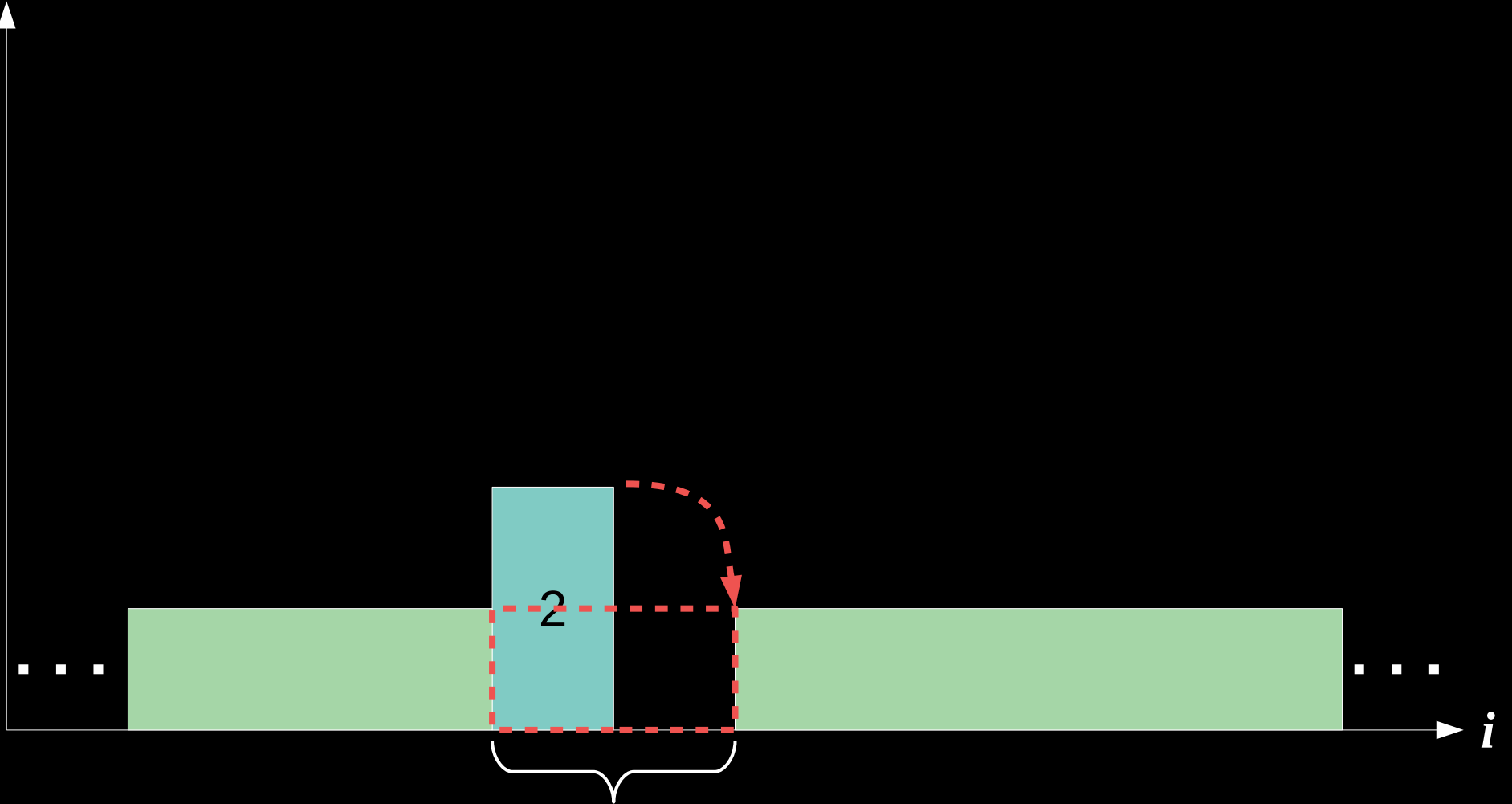


PLCP[ $i$ ]



*new maximum peak*

PLCP[ $i$ ]

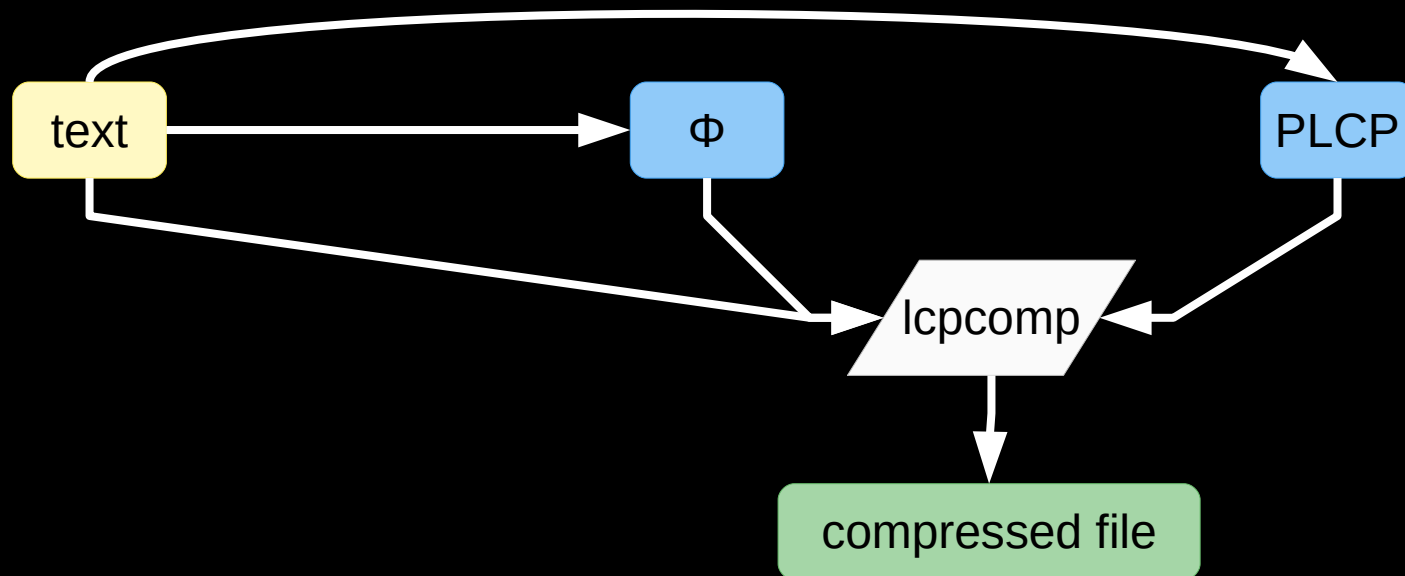


*factor of length 2*

PLCP[*i*]

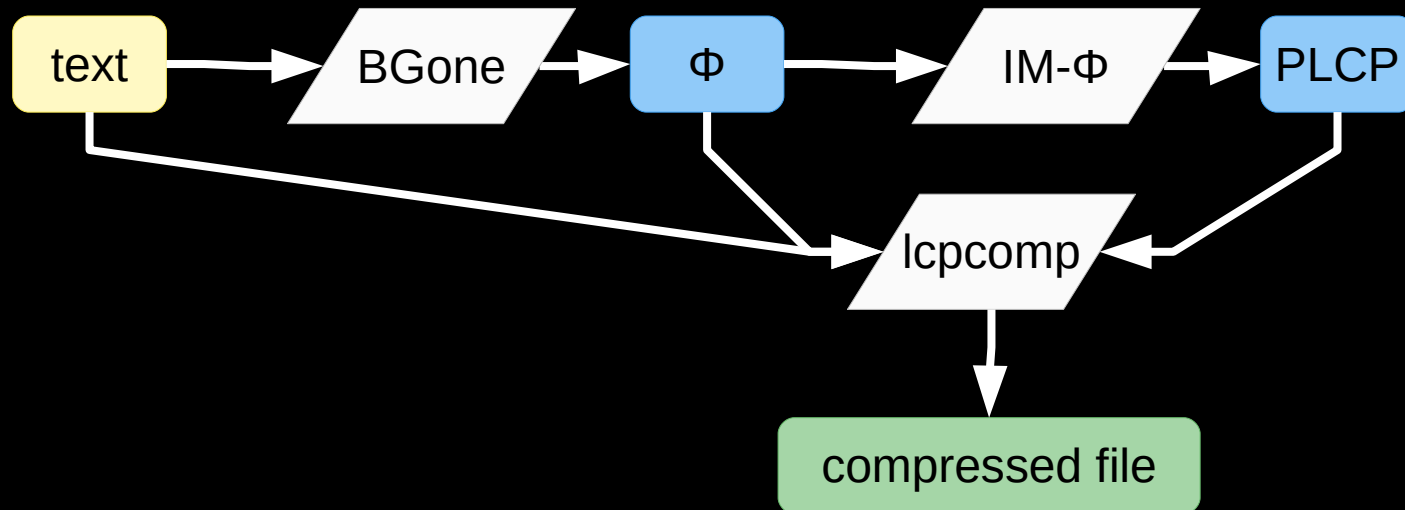


# data structures



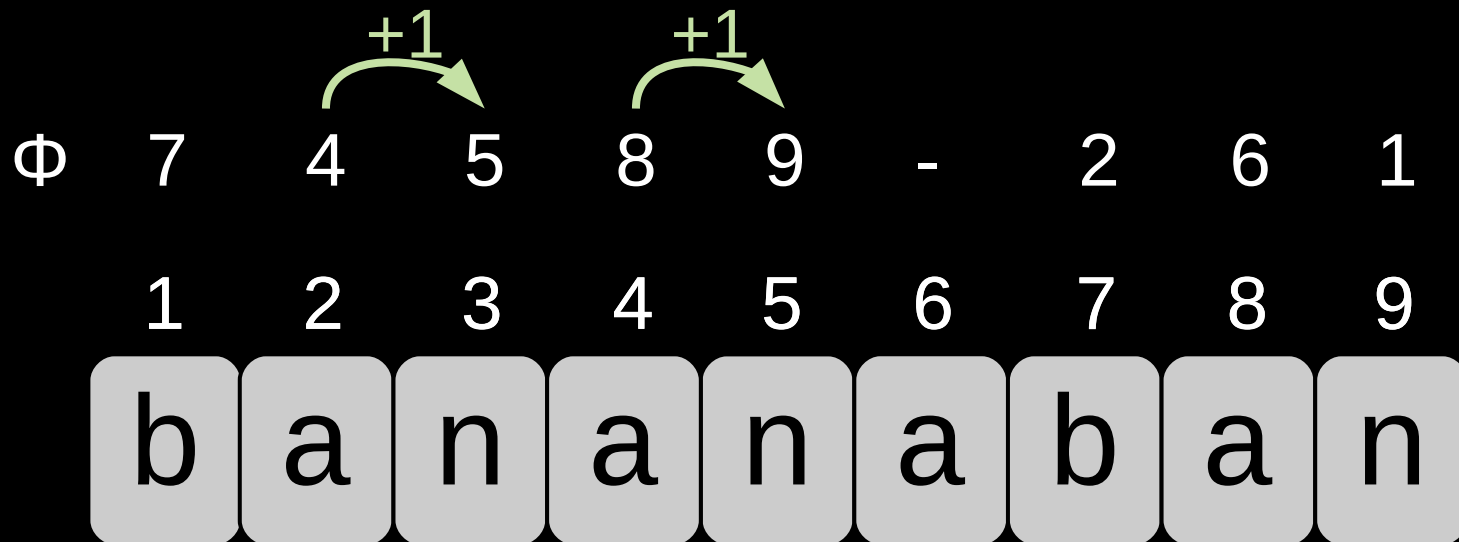


# data structures & algorithms



- BGone: 後藤, 坂内 '14
  - IM- $\Phi$ : Kärkkäinen+'09
- }  $O(n)$  time

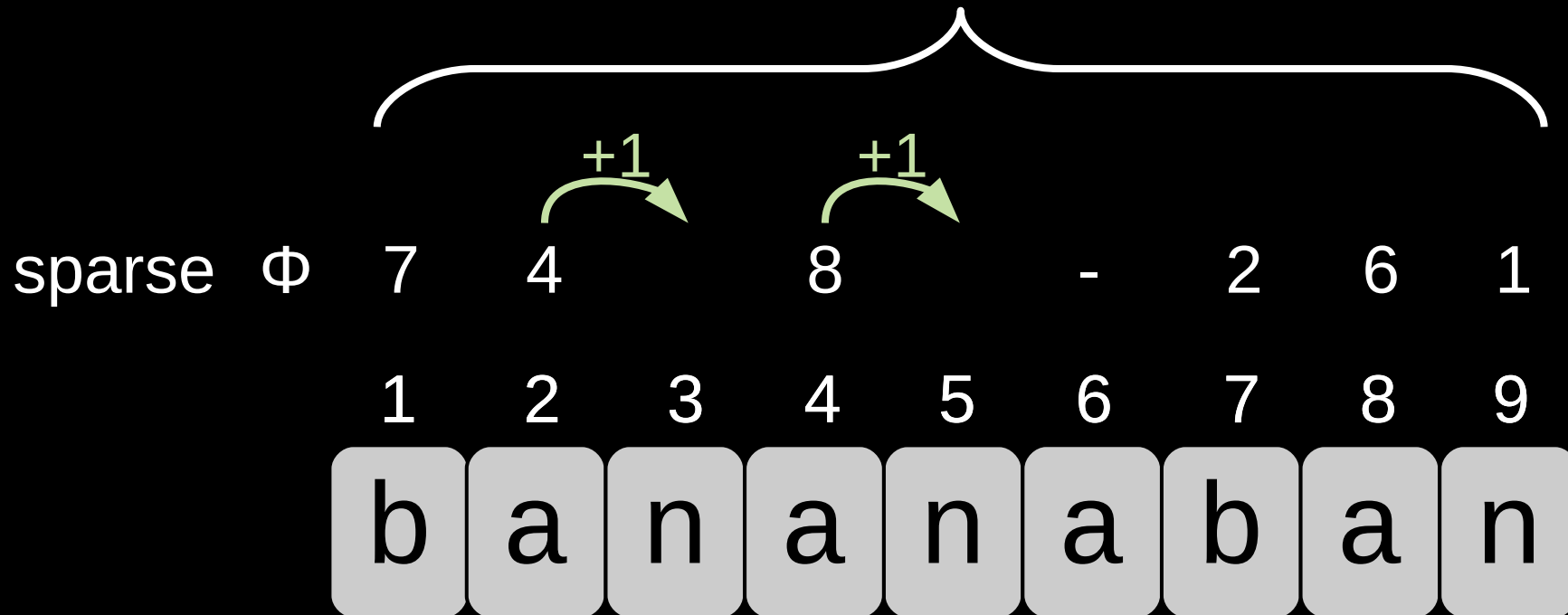
# sparse $\Phi$ /PLCP



# sparse $\Phi$ /PLCP

$|\text{sparse } \Phi| = \#\text{runs in BWT} =: r$

[Kärkkäinen+ '16]



# complexity

• time :  $O(n)$

• space: PLCP + bit vector construct  $\Phi$

$$\max(\underbrace{r \lg n}_{\text{sparse } \Phi} + \overbrace{3n+o(n)}^{\text{PLCP + bit vector}} + \underbrace{(O(n \lg n) + 1)^{0.5} \lg n}_{\text{maintain interesting peaks}}, \overbrace{n \lg n}^{\text{construct } \Phi})$$

sparse  $\Phi$

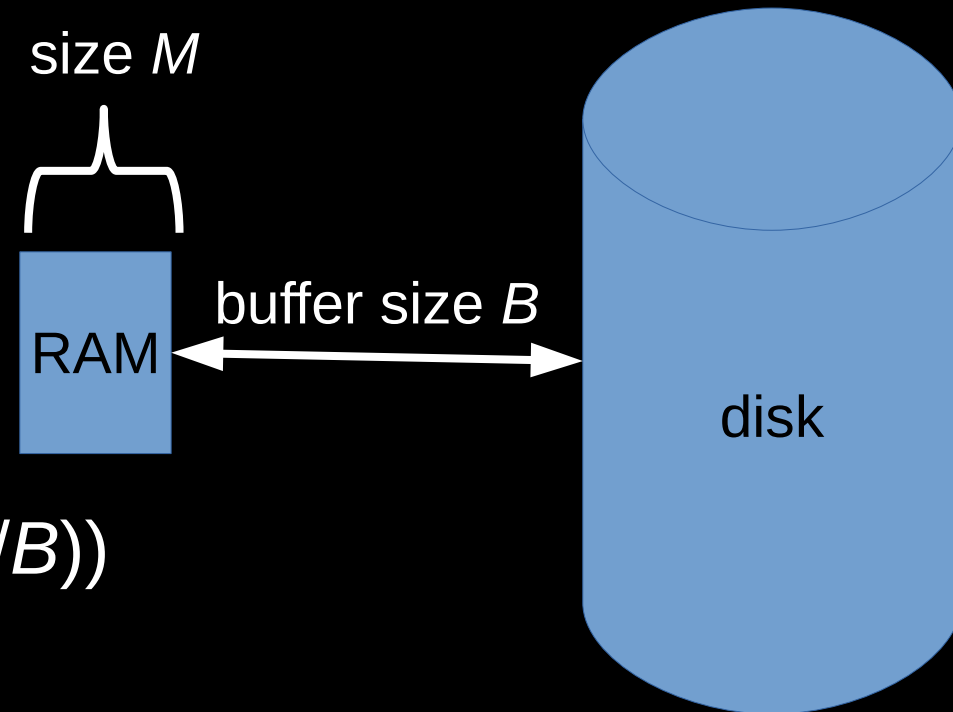
maintain interesting peaks

# external memory

- # of I/Os

- $\text{scan}(n) = \Theta(n/B)$

- $\text{sort}(n) = \Theta\left(\frac{n}{B} \log_{M/B} \frac{n}{B}\right)$



# execution

given text, PLCP,  $\phi$

- scan PLCP
  - create tuples  $\langle j, \text{PLCP}[j] \rangle$
  - sort
- scan  $\phi$ 
  - create tuples  $\langle j, \text{PLCP}[j], \phi[j] \rangle$
- scan text
  - create output



3 scan + 1 sort

# summary

## plcpcomp

- bidirectional compression scheme
- algorithmic improvement to lcpcomp
- linear scan of text, PLCP,  $\Phi$
- $\Rightarrow$  works in EM
- **fastest** solution to compress file in EM with ratio  $\sim$  LZ77

## Outlook

- bounds for #factors?
- bounds for maintaining interesting peaks
- algorithmic optimization
- more variants
  - using  $\Phi^{-1}$
  - unidirectional