

# PLCP 配列に基づくテキスト圧縮

クップル ドミニク (九州大学 / JSPS)

共著 :

Patrick Dinklage

Jonas Ellert

Johannes Fischer

Manuel Penschuck

# greedy 圧縮

名前	方法	参考
LZ77	unidirectional	Lempel,Ziv '77
longest first	grammar	中村 + '09
lcpcomp	bidirectional	Dinklage+ '17
LZ-LFS	hybrid	Mauer+ '17 西 + '18
LZRR	bidirectional	西本 +'19

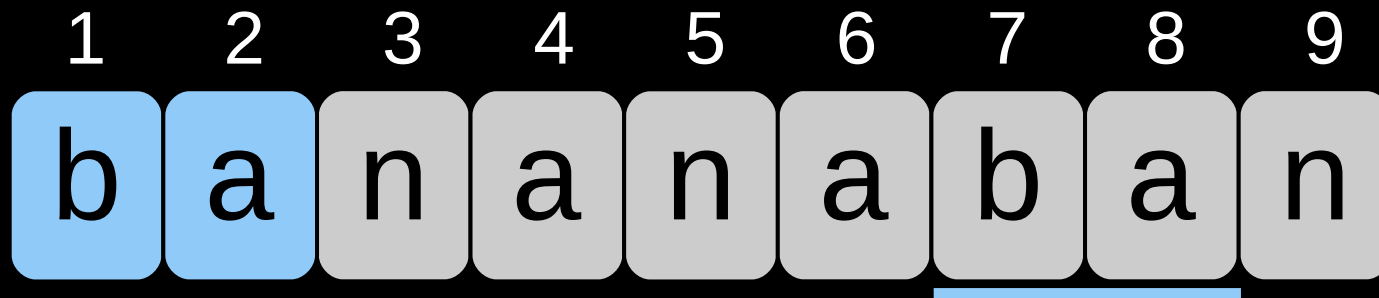
# bidirectional 圧縮

- 部分文字列を参照に置換

1 2 3 4 5 6 7 8 9  
b a n a n a b a n

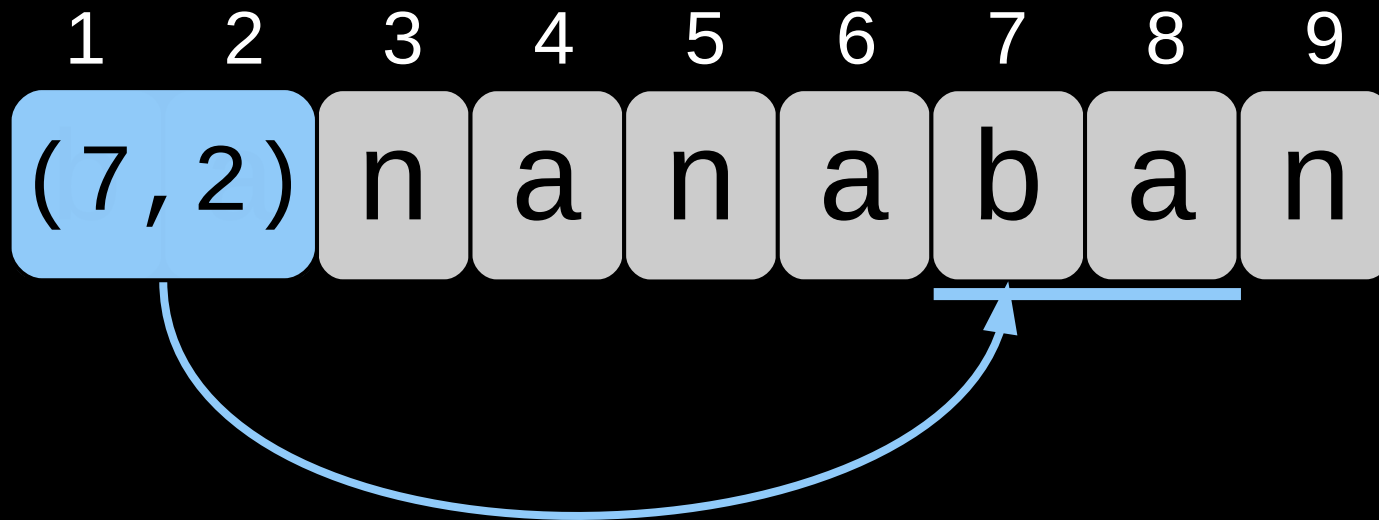
# bidirectional 圧縮

- 部分文字列を参照に置換
- 任意な順序



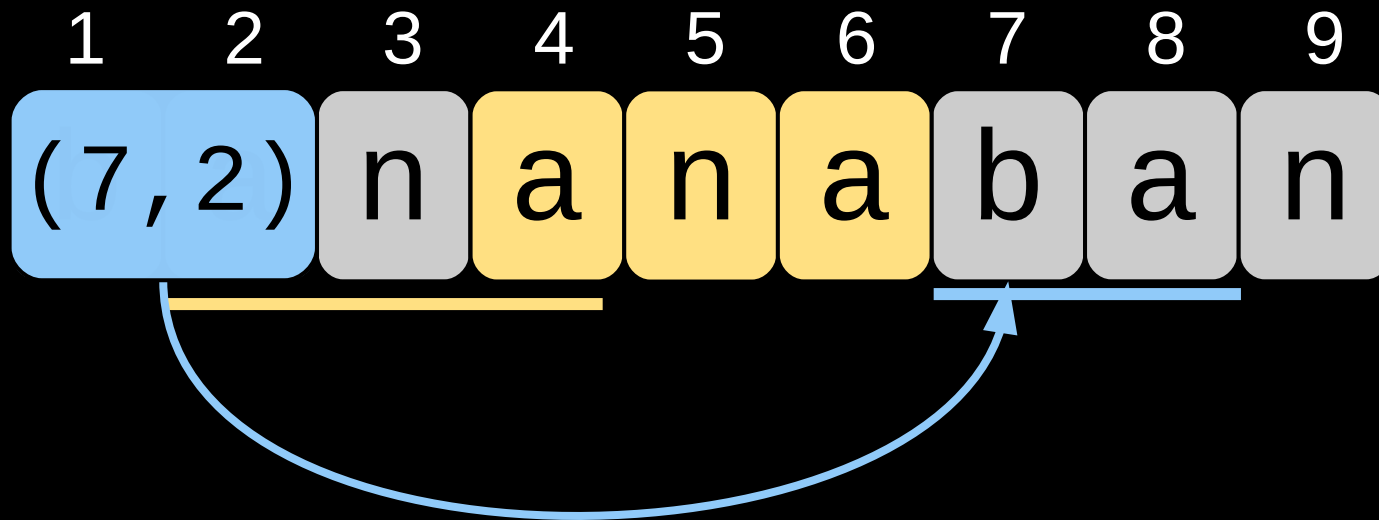
# bidirectional 圧縮

- 部分文字列を参照に置換
- 任意な順序



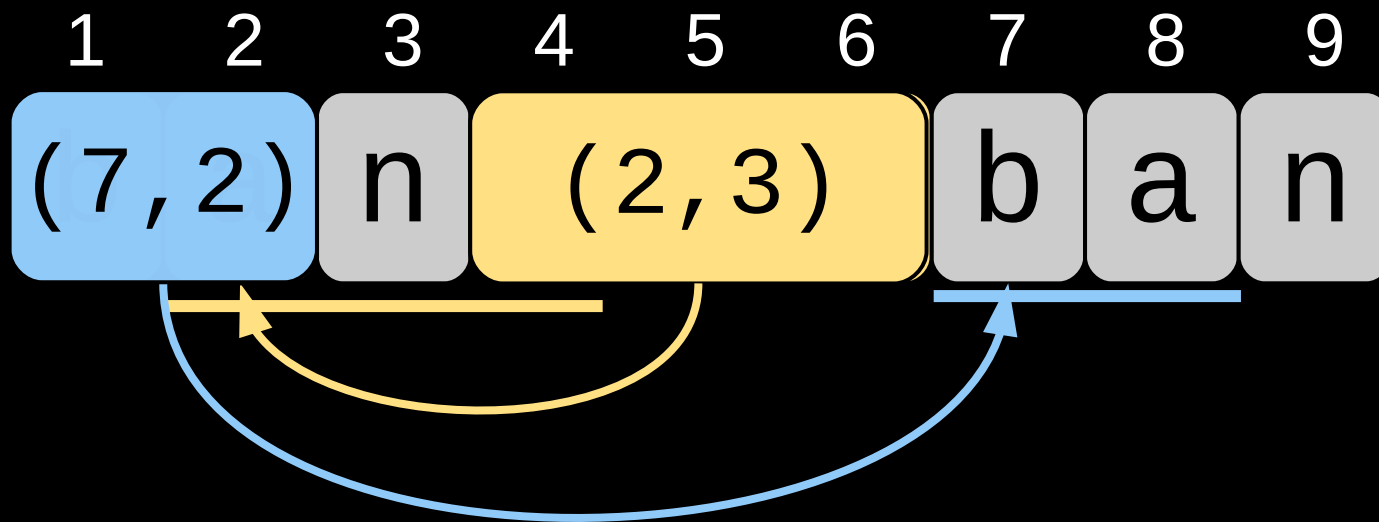
# bidirectional 圧縮

- 部分文字列を参照に置換
- 任意な順序
- 自己参照も OK



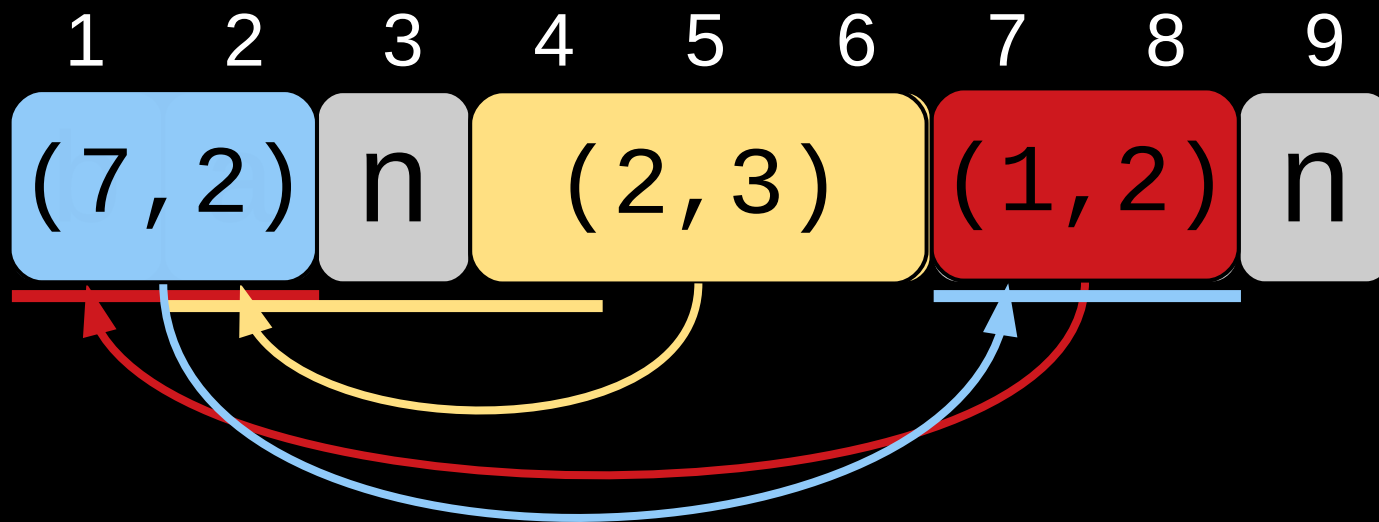
# bidirectional 圧縮

- 部分文字列を参照に置換
- 任意な順序
- 自己参照も OK



# bidirectional 圧縮

- 部分文字列を参照に置換
- 任意な順序
- 自己参照も OK
- 復元のため、cycle は禁止





# lcpcomp

- 最長の繰り返し部分文字列を置換、
- 繰り返す

1	2	3	4	5	6	7	8	9
b	a	n	a	n	a	b	a	n

# lcpcomp

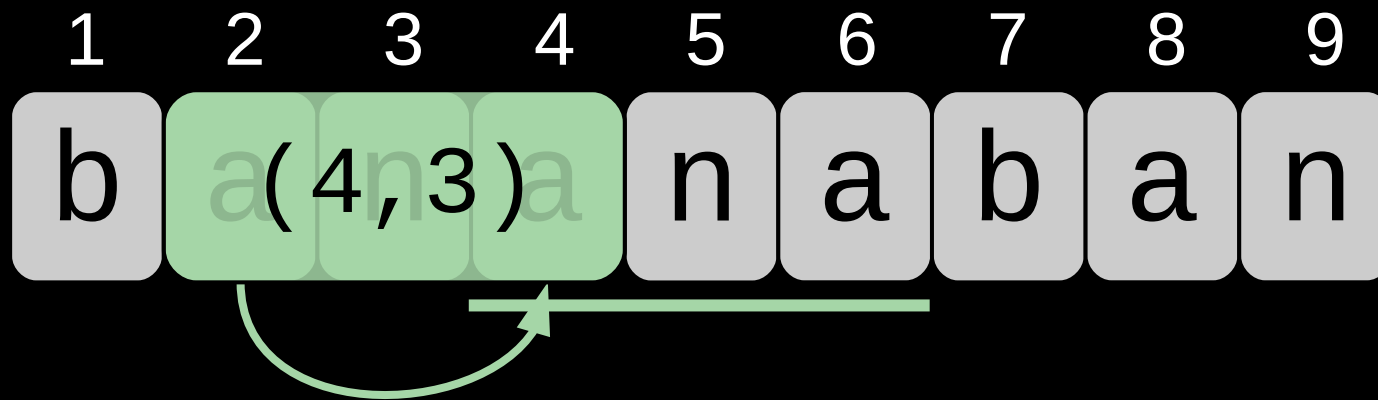
- 最長の繰り返し部分文字列を置換、
- 繰り返す



位置 4 から 3 文字をコピーし、

# lcpcomp

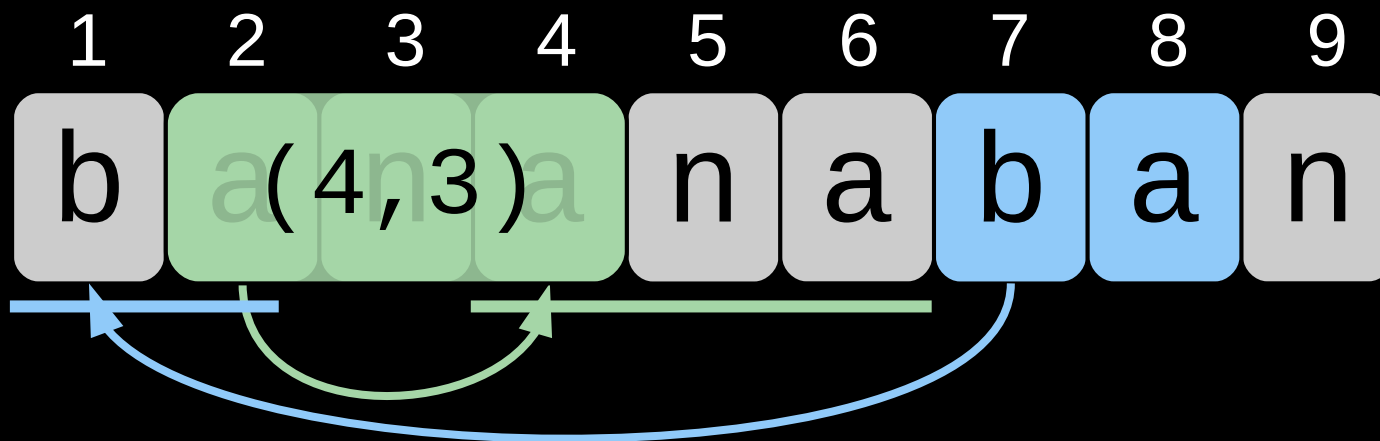
- 最長の繰り返し部分文字列を置換、
- 繰り返す



位置 4 から 3 文字をコピーし、

# lcpcomp

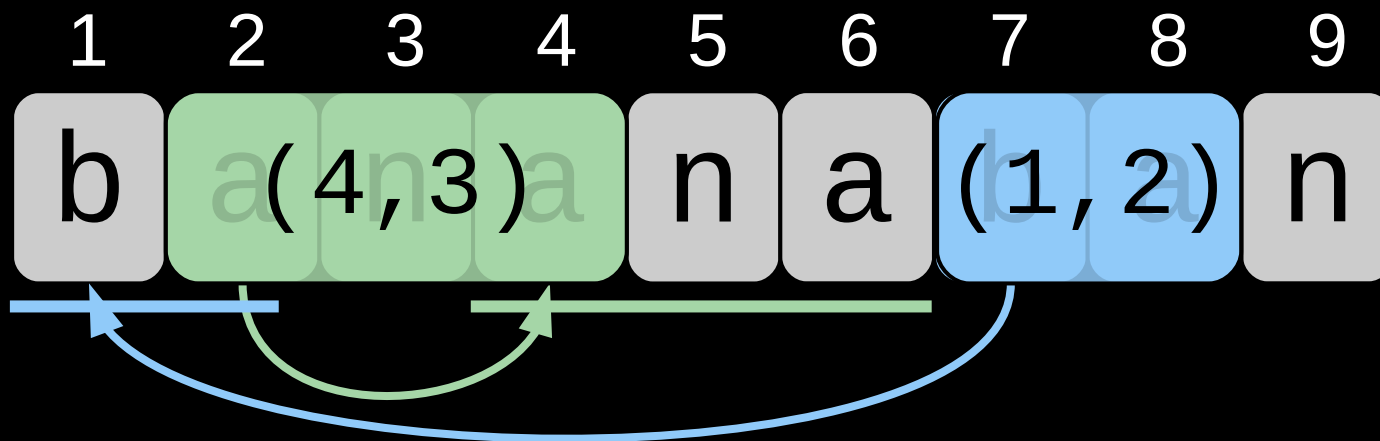
- 最長の繰り返し部分文字列を置換、
- 繰り返す



位置 4 から 3 文字をコピーし、  
位置 1 から 2 文字をコピーする

# lcpcomp

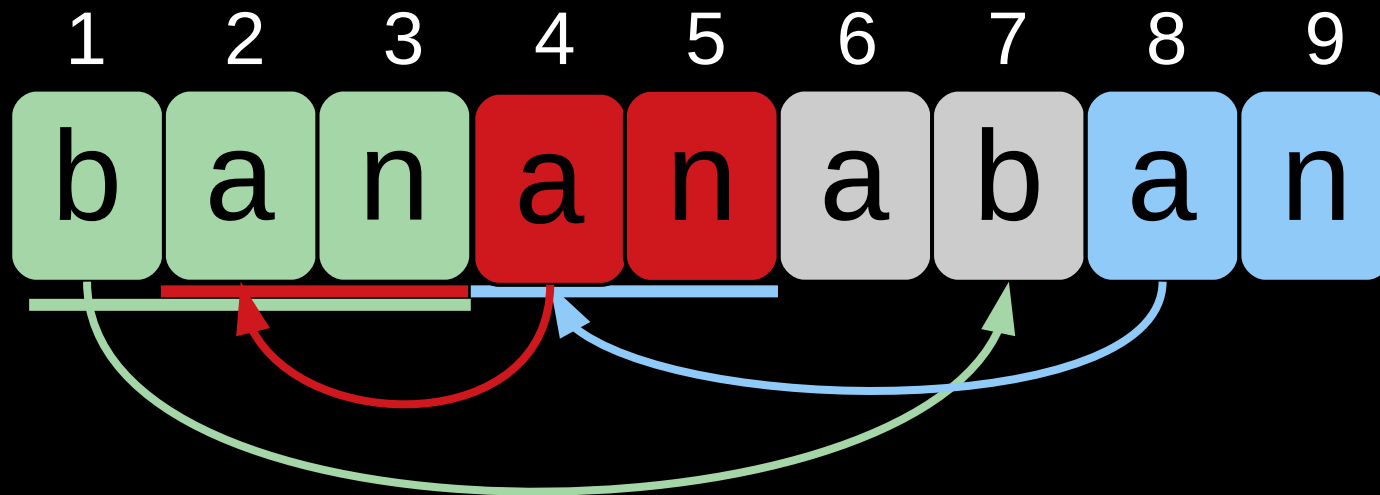
- 最長の繰り返し部分文字列を置換、
- 繰り返す



位置 4 から 3 文字をコピーし、  
位置 1 から 2 文字をコピーする

# 復元

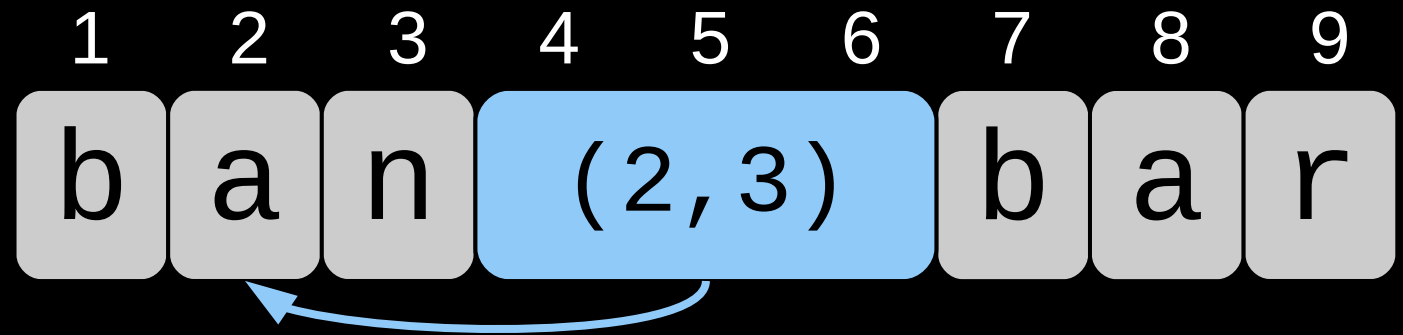
- cycle を防ぐためのルール
- 参照は辞書式順序の小さい接尾辞



# LZ77 に比べて

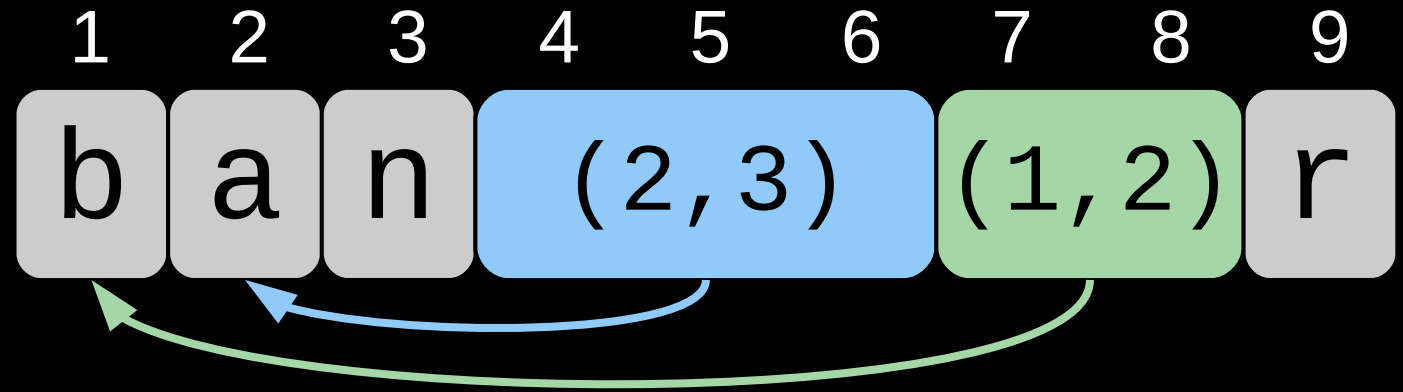
1 2 3 4 5 6 7 8 9  
b a n a n a b a r

# LZ77 に比べて

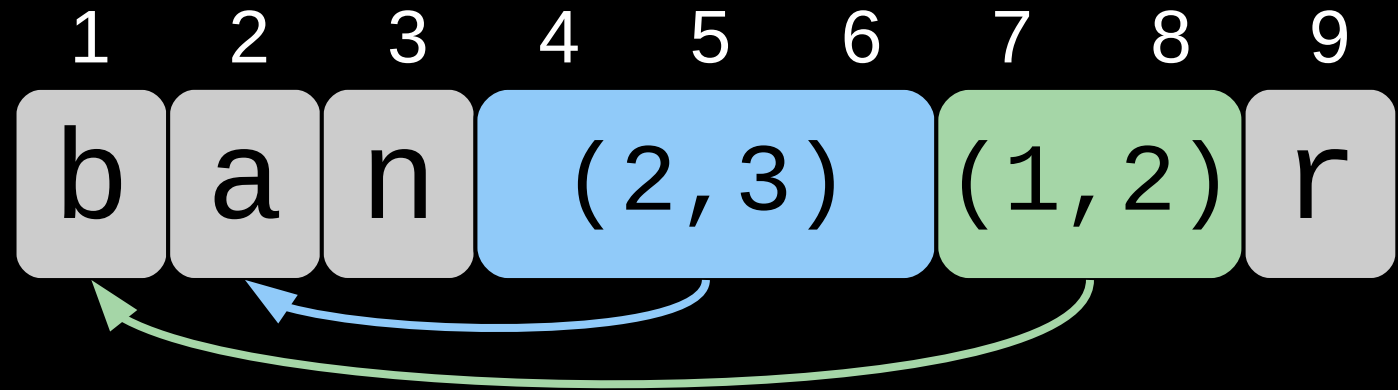




# LZ77 に比べて



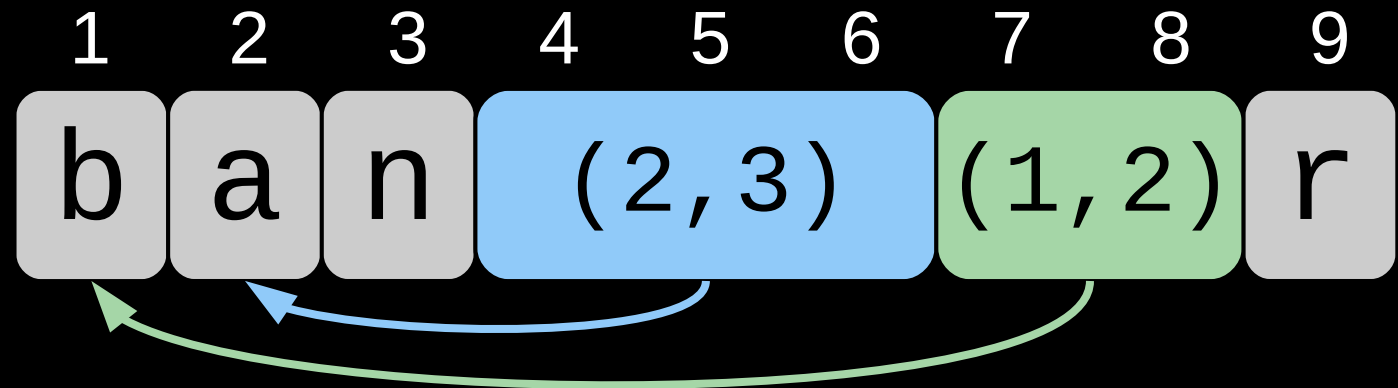
# LZ77 に比べて



- LZ77

参照はテキスト順序で前の位置

# LZ77 に比べて



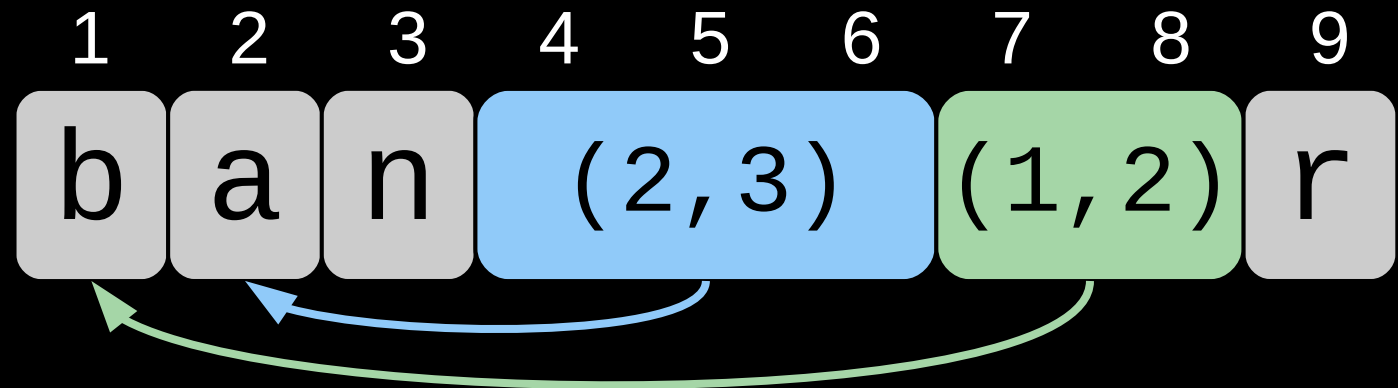
- LZ77

参照はテキスト順序で前の位置

- lcpcomp

参照は辞書式順序で前の接尾辞

# LZ77 に比べて



- LZ77

参照はテキスト順序で前の位置

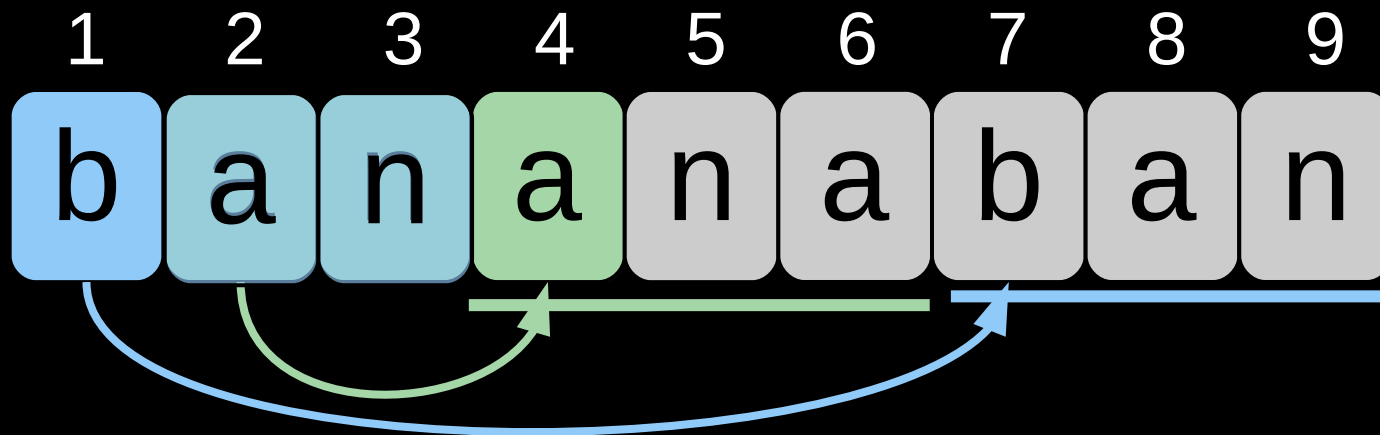
- lcpcomp

参照は辞書式順序で前の接尾辞

⇒ どちらも greedy かつ cycle がない

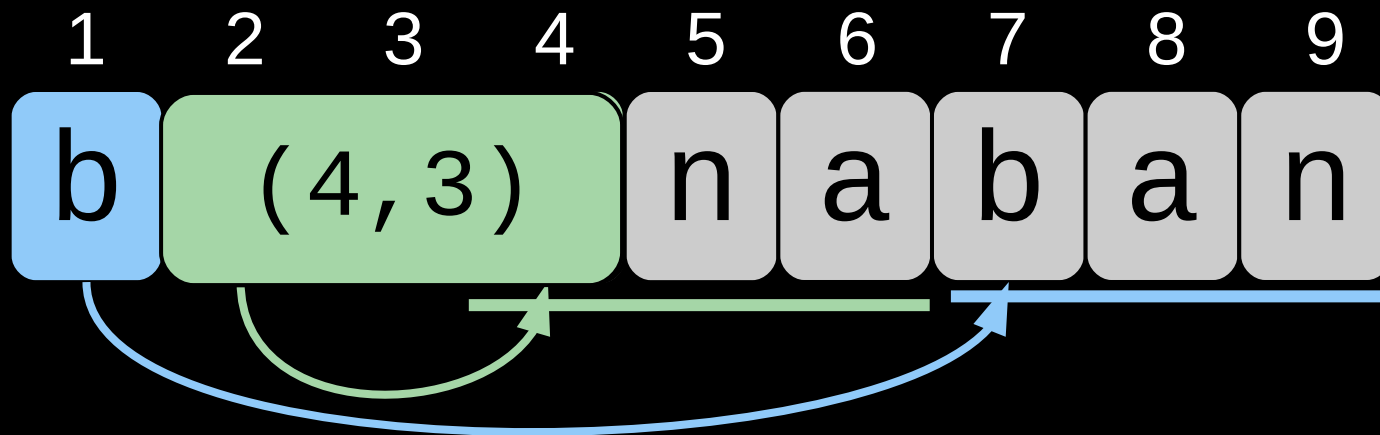
# lcpcomp

- どの参照を作るか選択できる



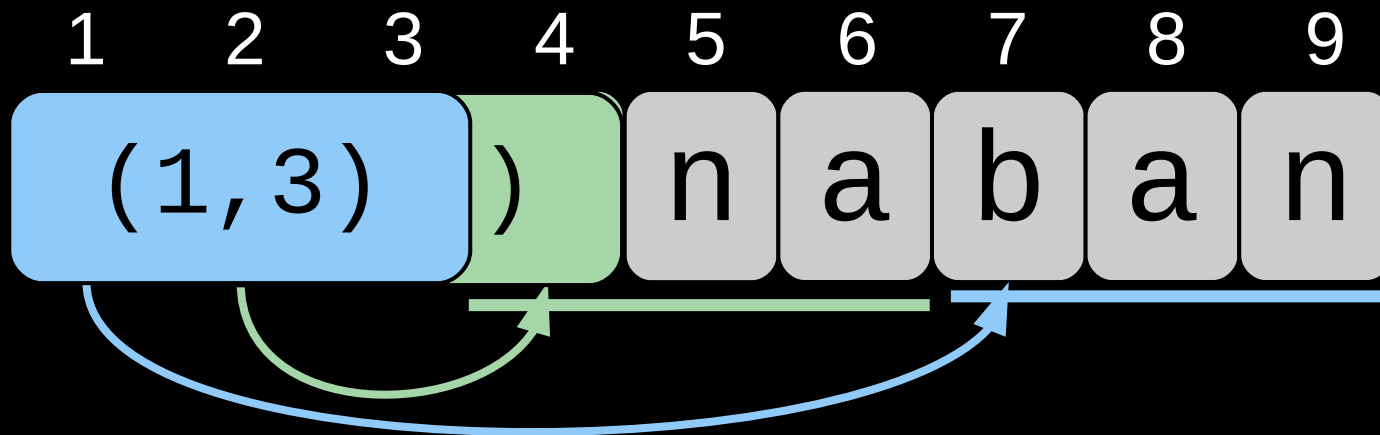
# lcpcomp

- どの参照を作るか選択できる



# lcpcomp

- どの参照を作るか選択できる



# plcpcomp

plcpcomp は今回の新しい提案方法

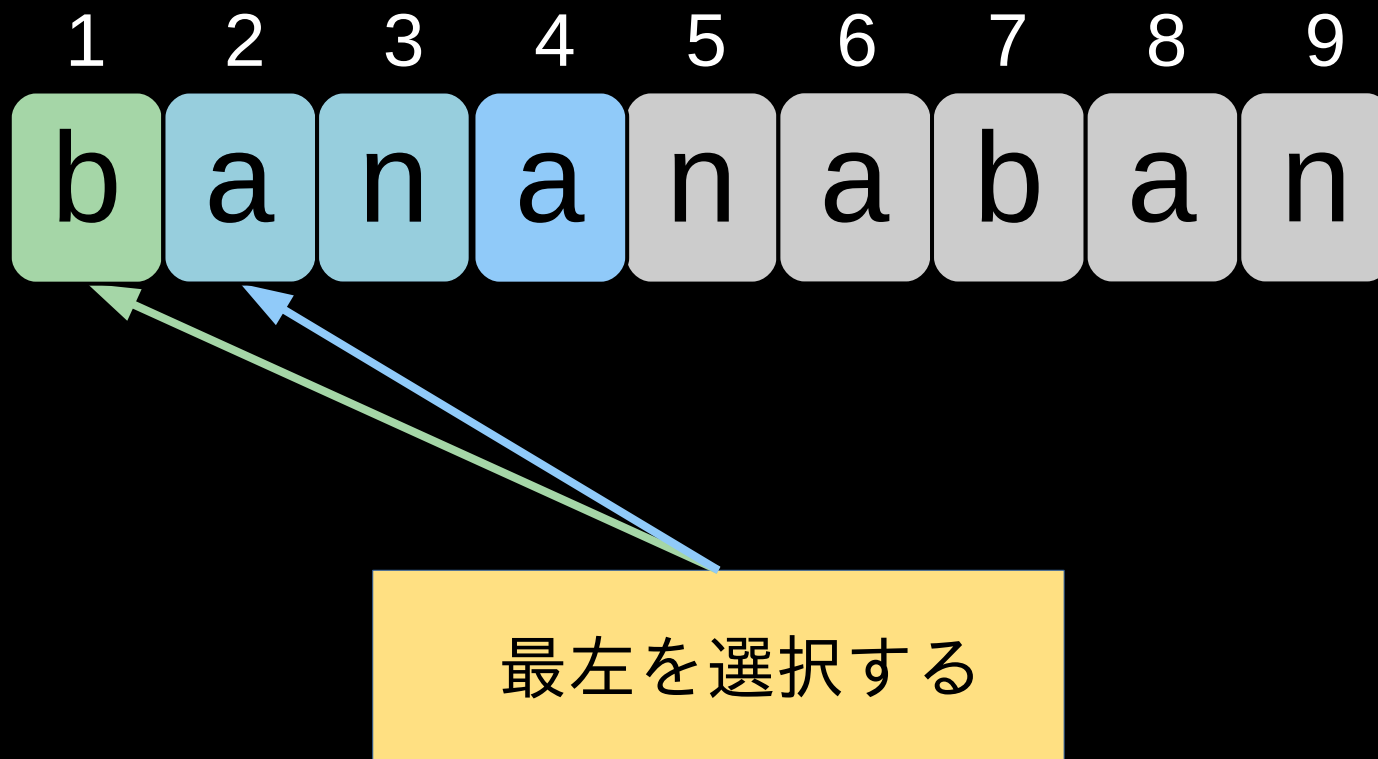
1	2	3	4	5	6	7	8	9
b	a	n	a	n	a	b	a	n

最左を選択する



# plcpcomp

plcpcomp は今回の新しい提案方法



# plcpcomp

plcpcomp は今回の新しい提案方法

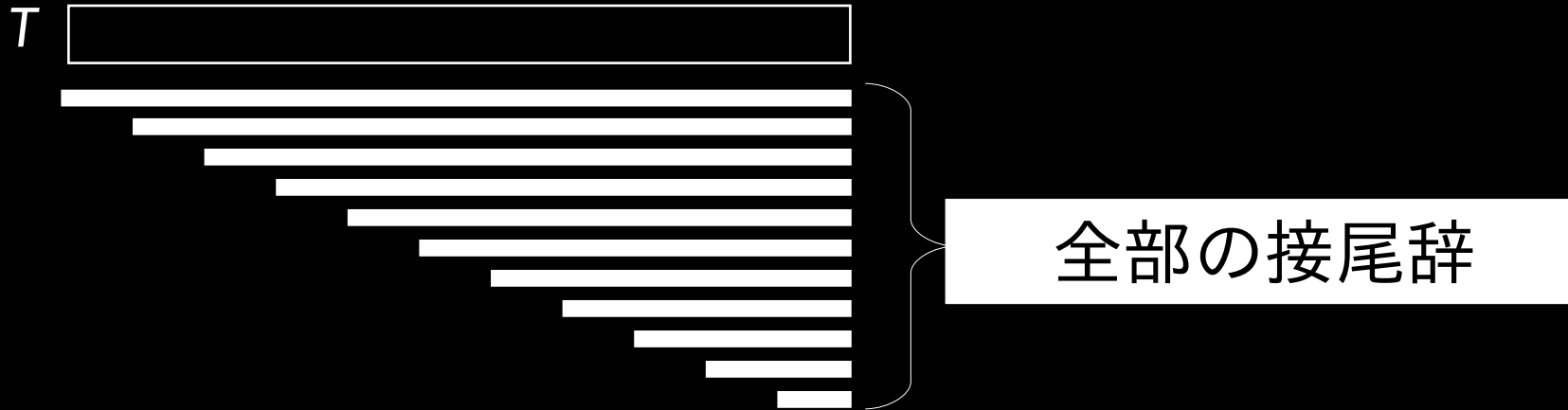


最左を選択する

# テキストデータ構造

$T$

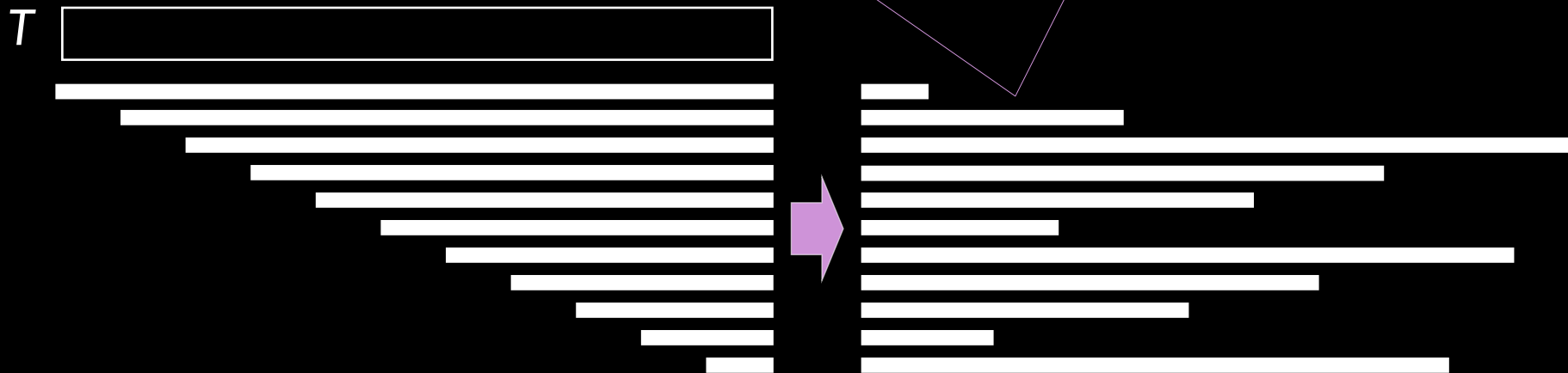
# テキストデータ構造



# テキストデータ構造

- SA: 接尾辞配列

辞書順で接尾辞を並べる



# テキストデータ構造

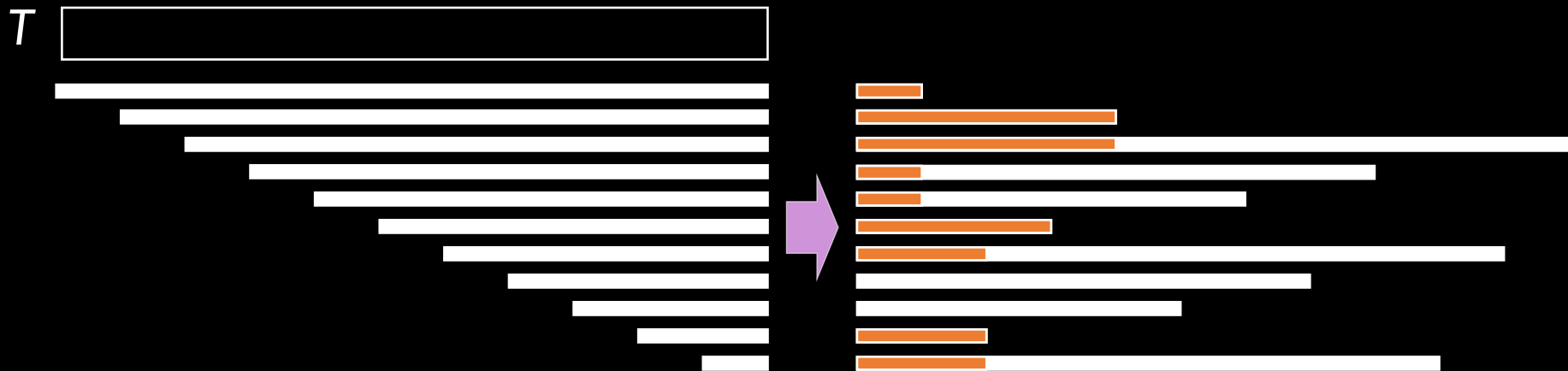
- SA: 接尾辞配列
- LCP: 隣の接尾辞の最長一致の接頭辞

辞書順で接尾辞を並べる



# テキストデータ構造

- SA: 接尾辞配列
- LCP: 隣の接尾辞の最長一致の接頭辞
- PLCP[i]:  $T[i..]$  という接尾辞の LCP の長さ



# どうやって？

## 目的

- 最長の繰り返し部分文字列  $S$  を見つける
- $S =$  2つの接尾辞の  $LCP$

## そのために

- $|S| = PLCP[i]$
- $S$  は  $\operatorname{argmax} PLCP[i]$  から出現する
- $PLCP[i]$  は  $i$  段目の接尾辞と  $\Phi[i]$  段目の接尾辞の  $LCP$  の長さ



# Φ 配列

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

Φ	7	4	5	8	9	-	2	6	1
SA	6	8	4	2	7	1	9	5	3
ISA	6	4	9	3	8	1	5	2	7
	1	2	3	4	5	6	7	8	9
	b	a	n	a	n	a	b	a	n

# Φ 配列

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

Φ	7	4	5	8	9	-	2	6	1
SA	6	8	4	2	7	1	9	5	3
ISA	6	4	9	3	8	1	5	2	7
	↑								
	1	2	3	4	5	6	7	8	9
	b	a	n	a	n	a	b	a	n

# Φ 配列

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

Φ	7	4	5	8	9	-	2	6	1
SA	6	8	4	2	7	1	9	5	3
ISA	6	4	9	3	8	1	5	2	7
	1	2	3	4	5	6	7	8	9
	b	a	n	a	n	a	b	a	n

# Φ 配列

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

Φ	7	4	5	8	9	-	2	6	1
SA	6	8	4	2	7	1	9	5	3
ISA	6	4	9	3	8	1	5	2	7
	1	2	3	4	5	6	7	8	9
	b	a	n	a	n	a	b	a	n

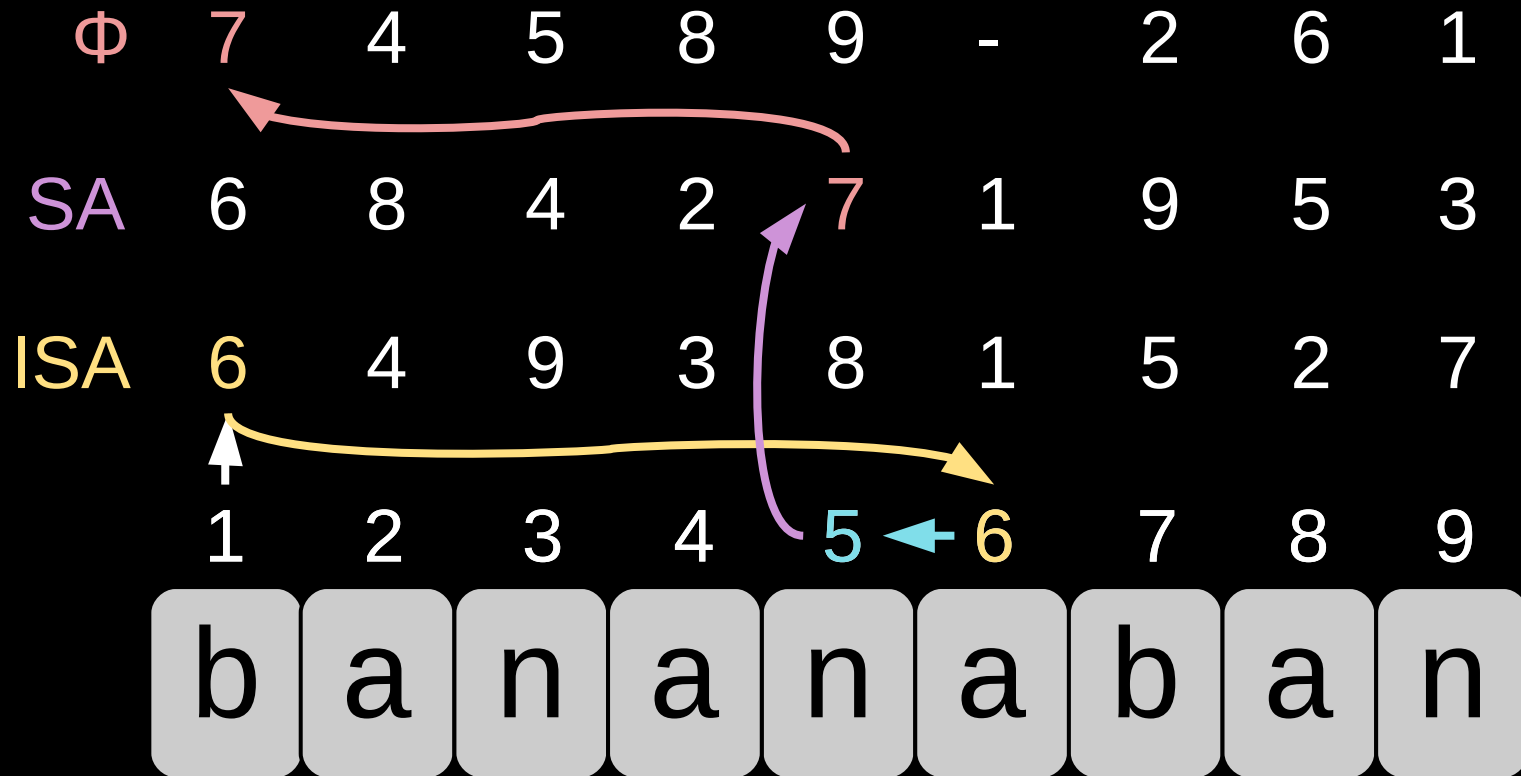
# Φ 配列

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

Φ	7	4	5	8	9	-	2	6	1
SA	6	8	4	2	7	1	9	5	3
ISA	6	4	9	3	8	1	5	2	7
	1	2	3	4	5	6	7	8	9
	b	a	n	a	n	a	b	a	n

# Φ 配列

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



# アルゴリズムの方針

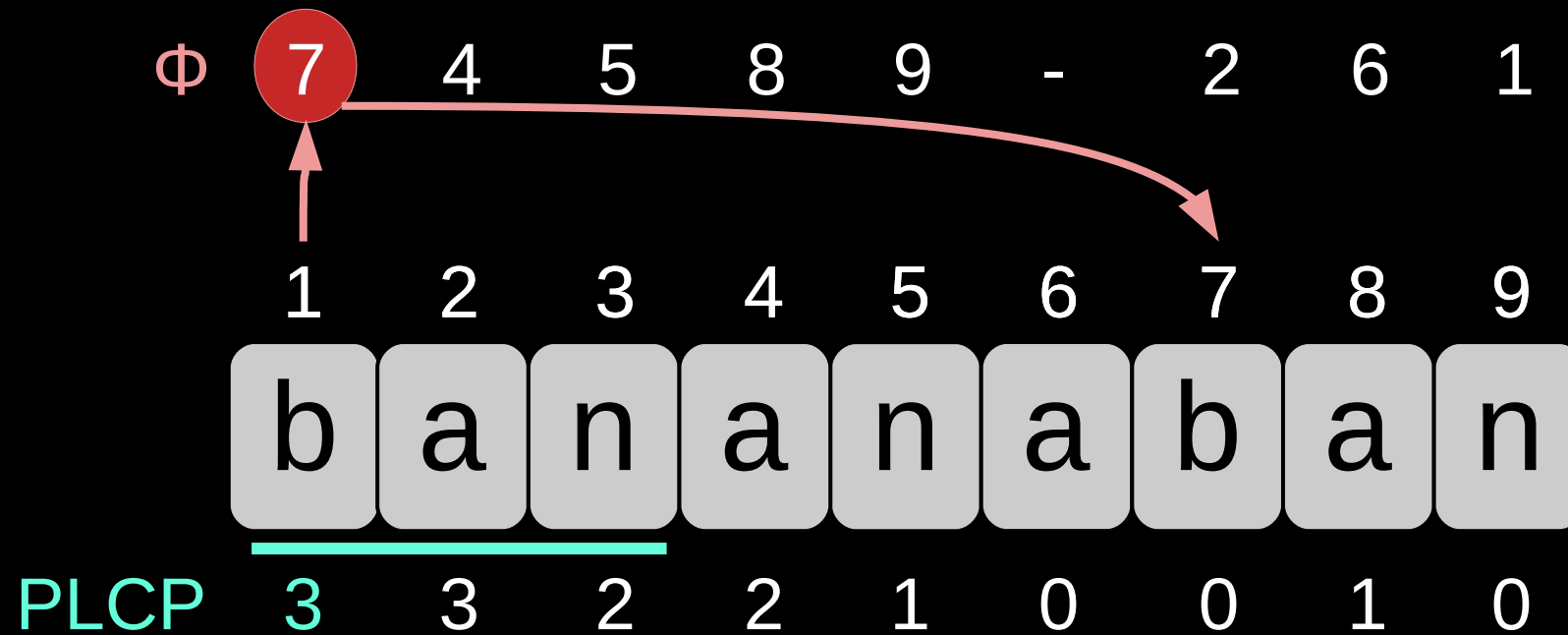
$\Phi$  7 4 5 8 9 - 2 6 1

1 2 3 4 5 6 7 8 9  
b a n a n a b a n

PLCP 3 3 2 2 1 0 0 1 0

- PLCP: 参照の長さ
- $\Phi$ : 参照がどこから始まるか

# アルゴリズムの方針

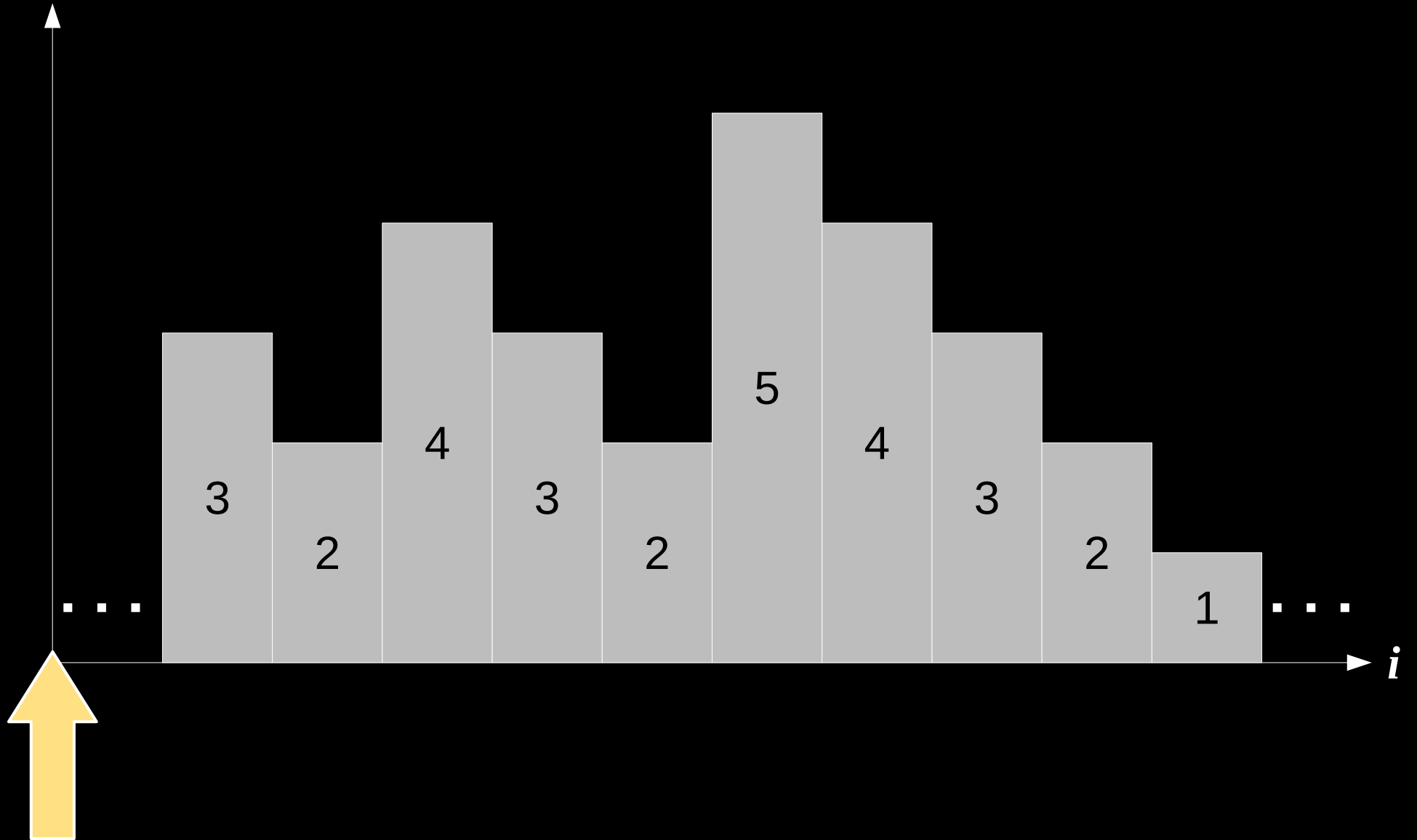


- **PLCP**: 参照の長さ
- $\Phi$ : 参照がどこから始まるか

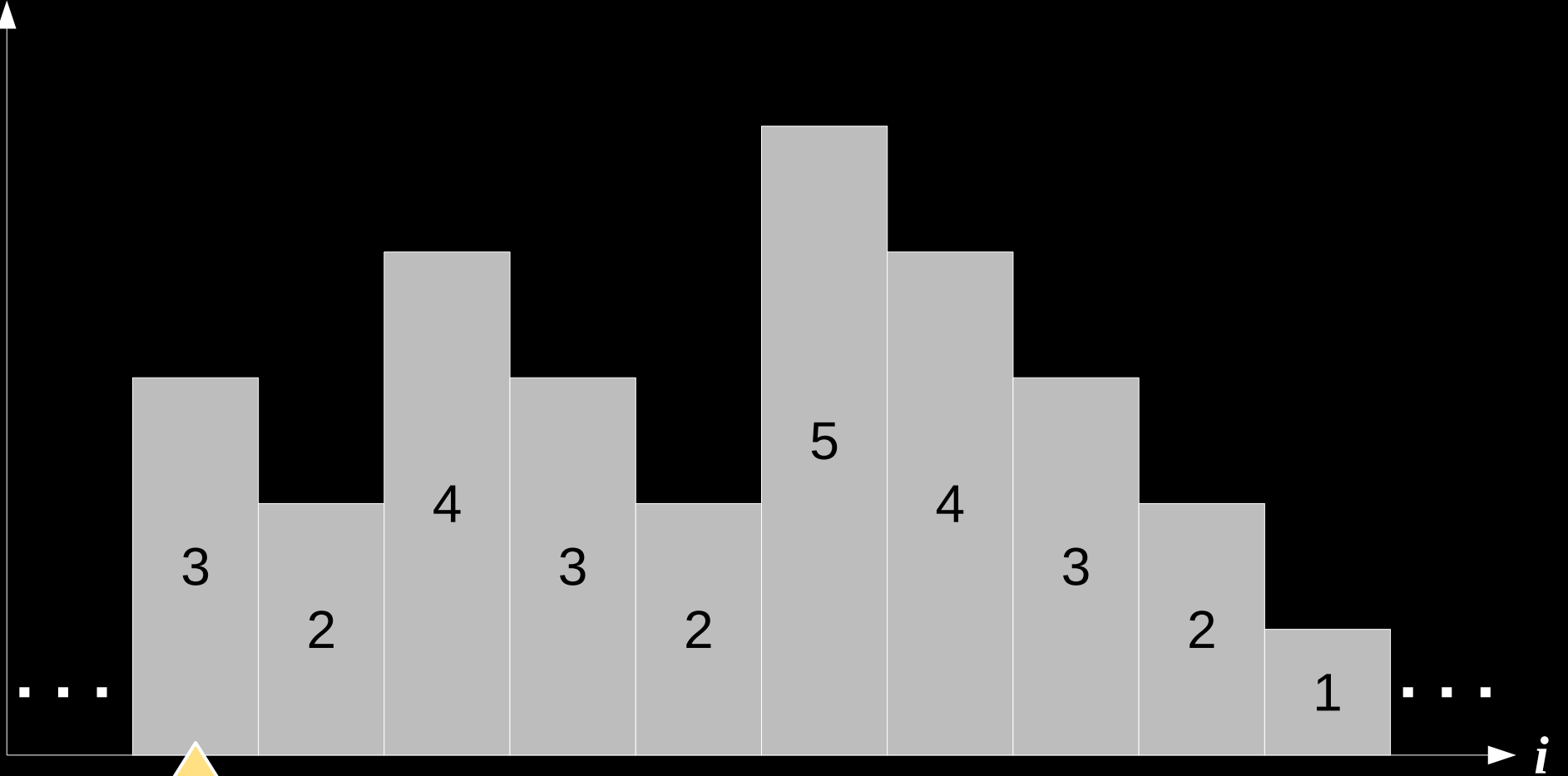


# 算出

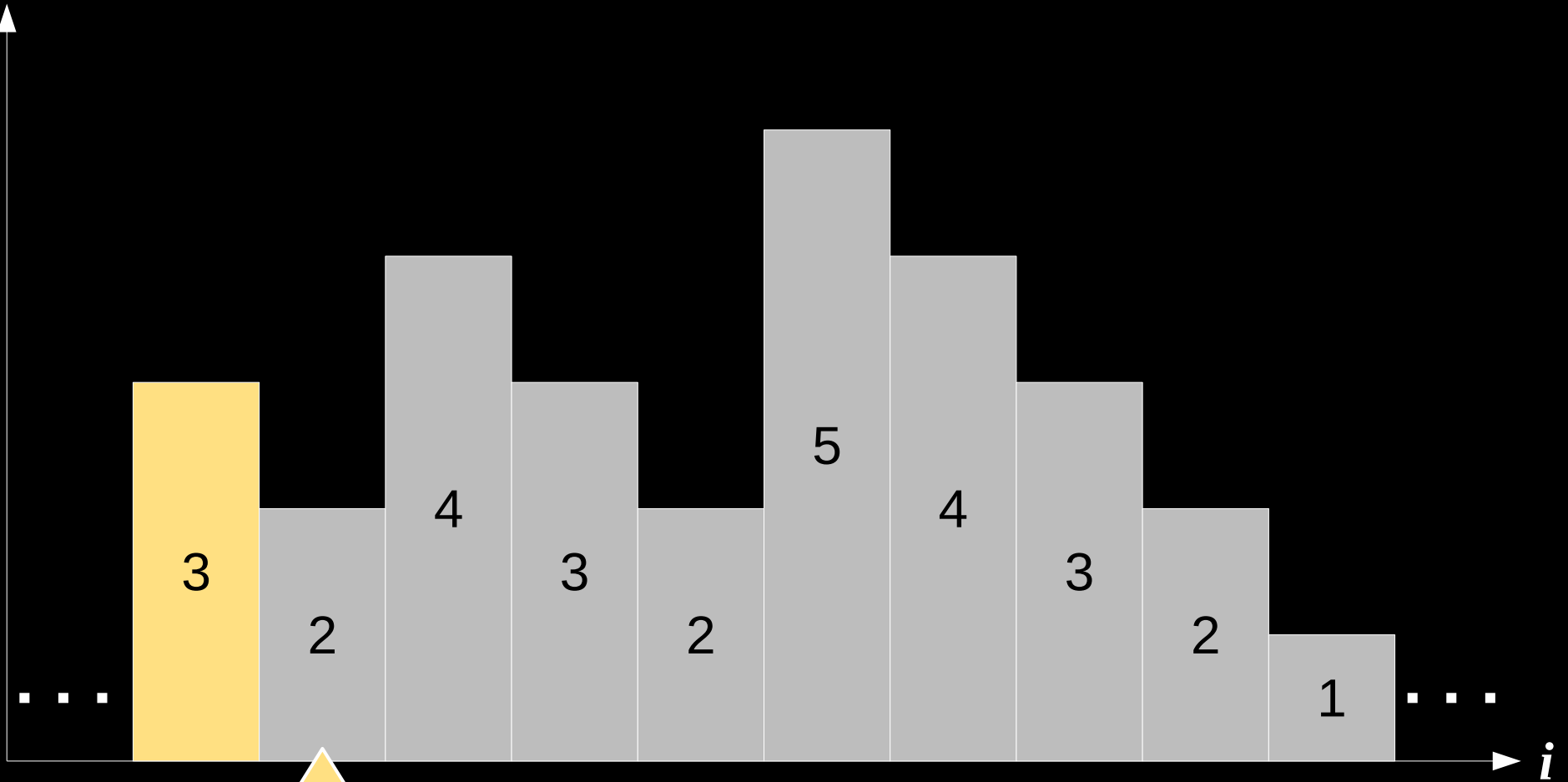
PLCP[ $i$ ]



PLCP[ $i$ ]

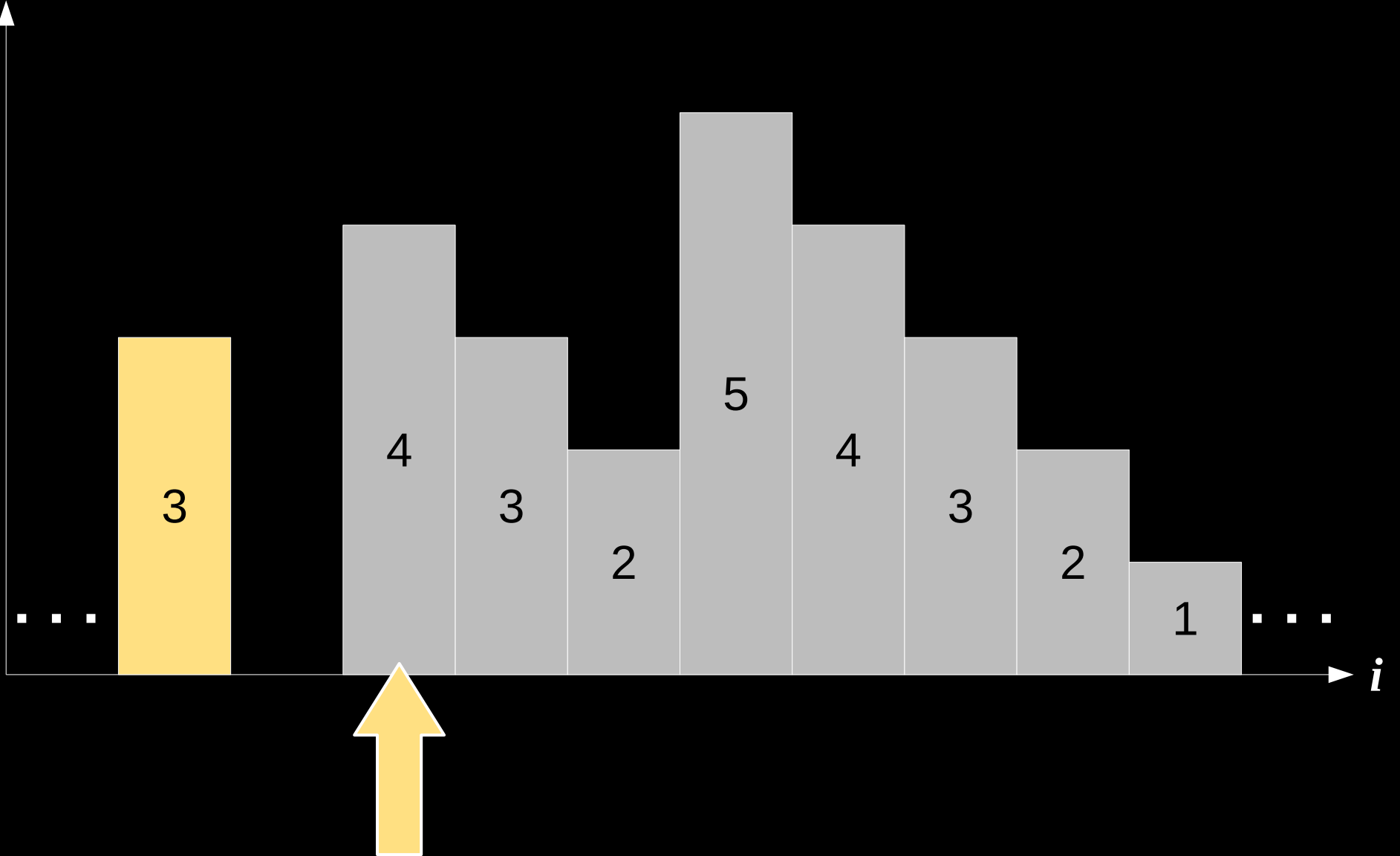


PLCP[ $i$ ]

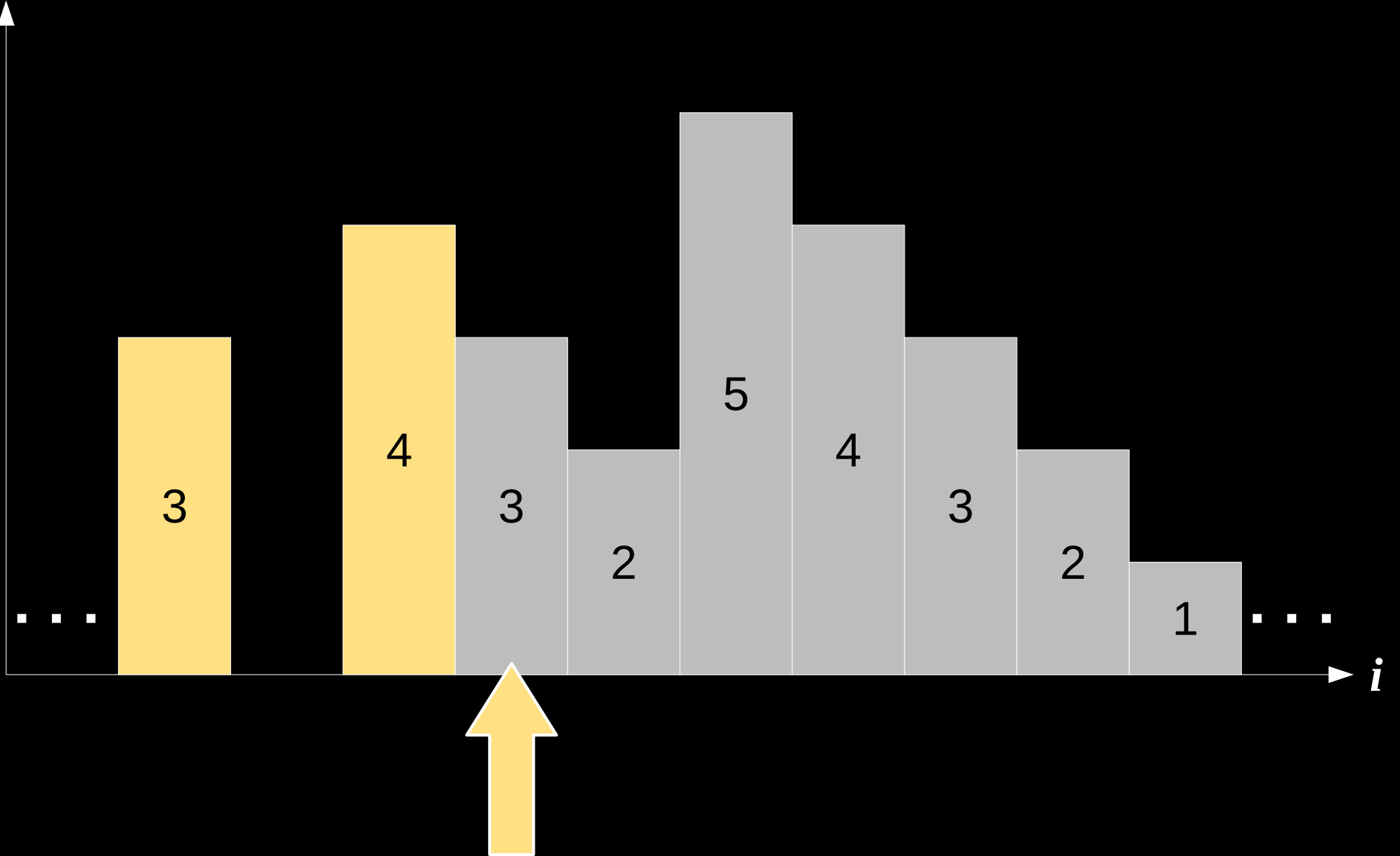


前の PLCP 値のほうが大きい  
⇒ 捨てる

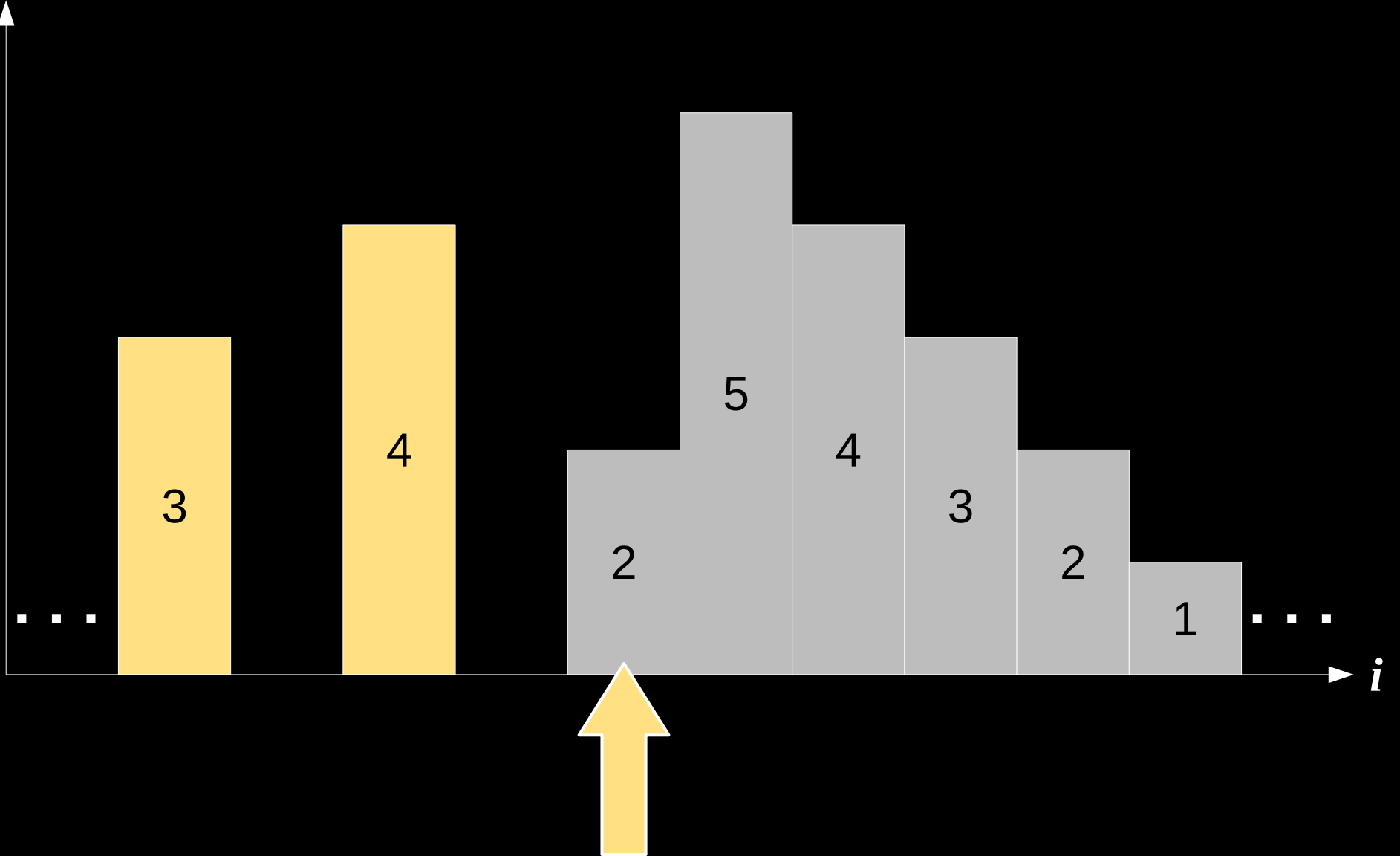
PLCP[ $i$ ]



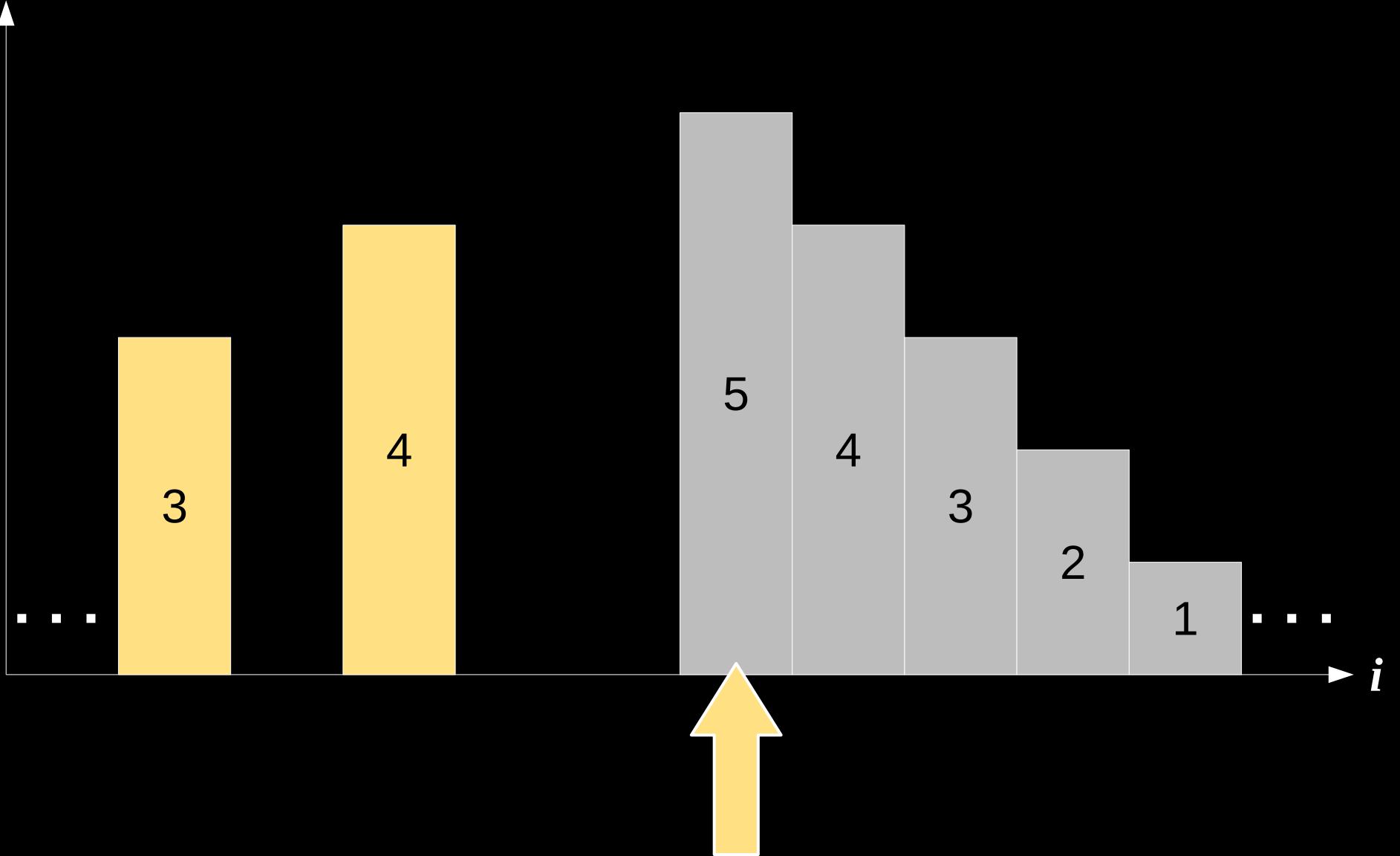
PLCP[ $i$ ]



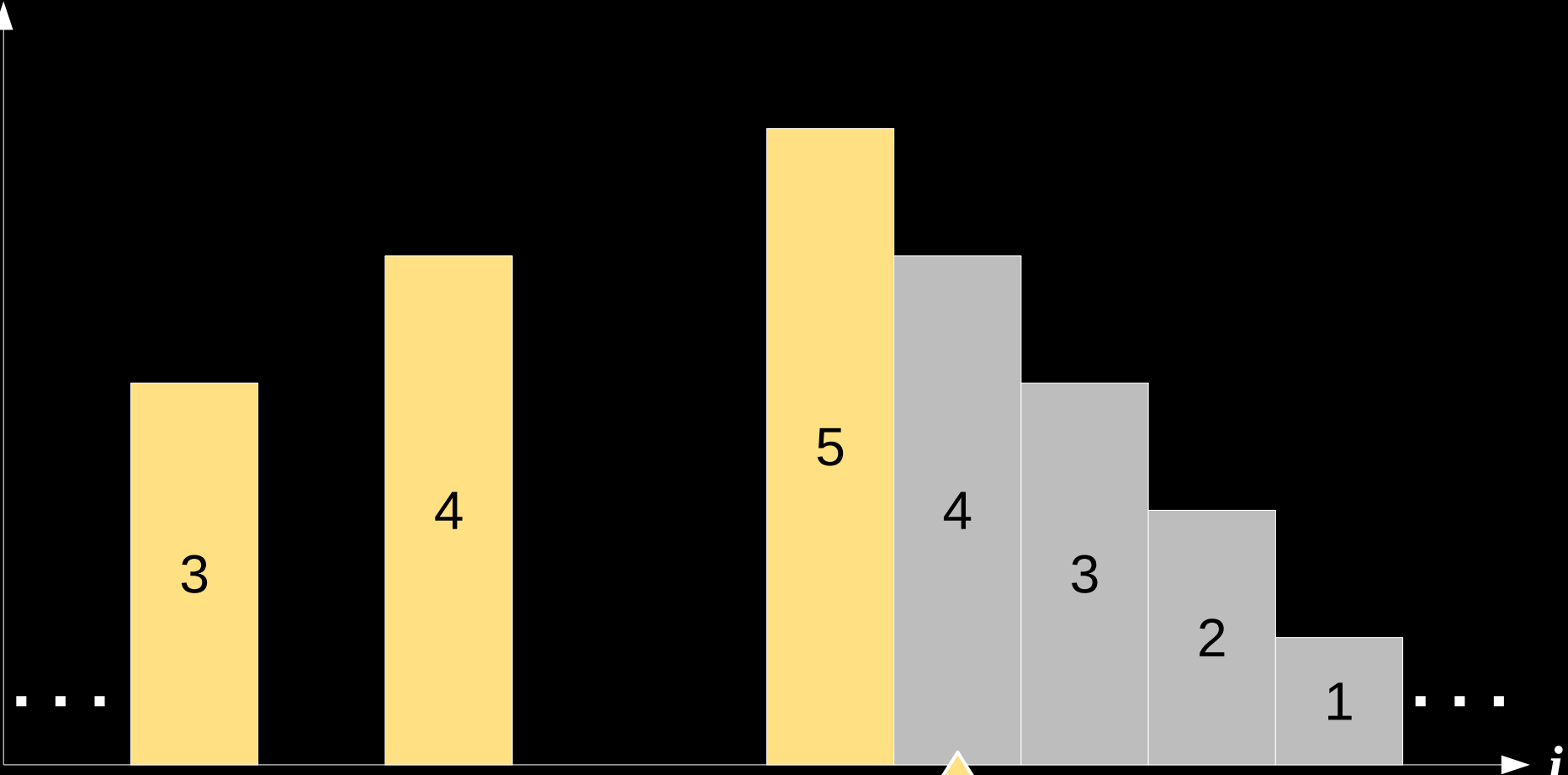
PLCP[*i*]



PLCP[ $i$ ]

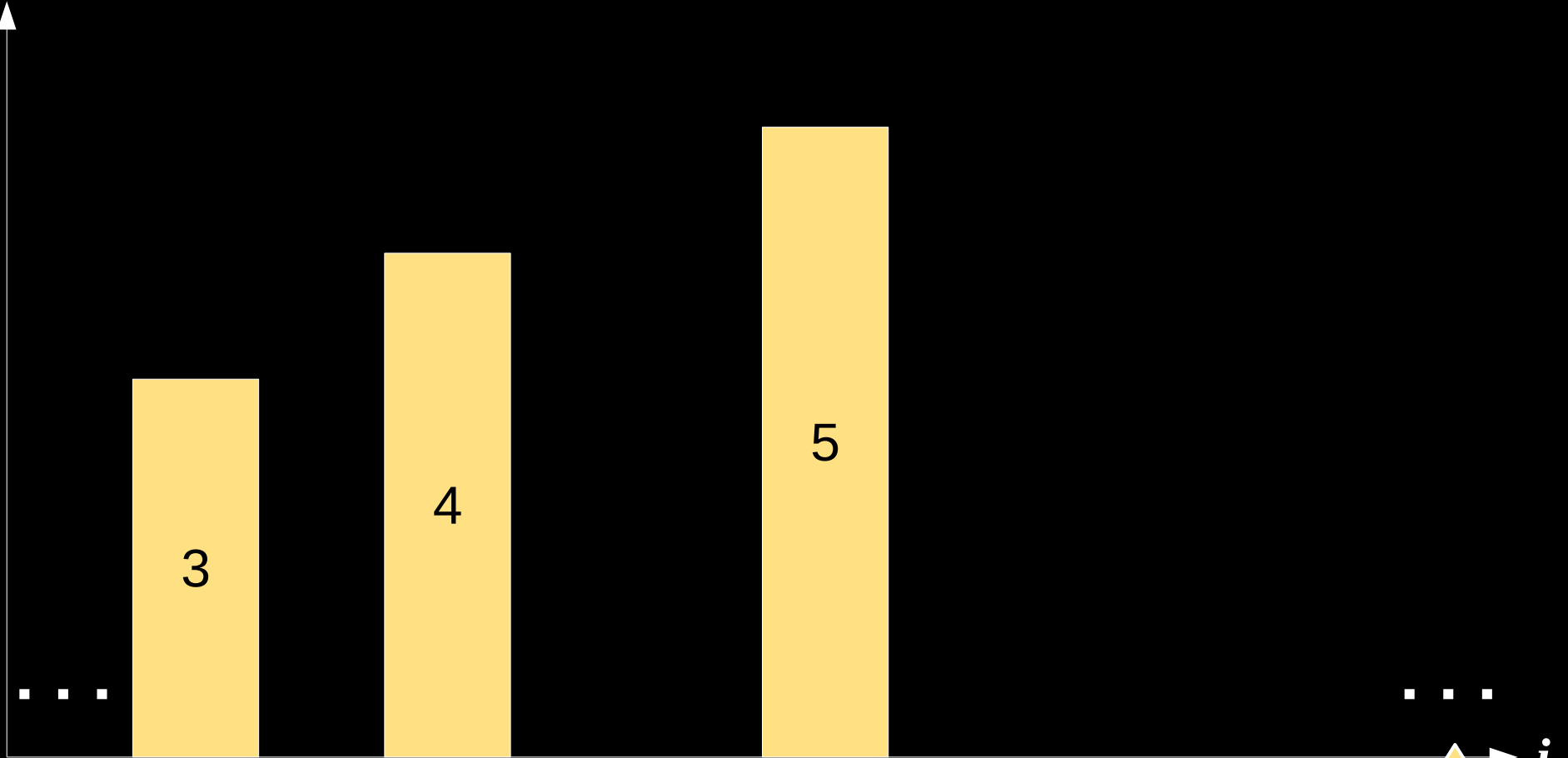


PLCP[ $i$ ]





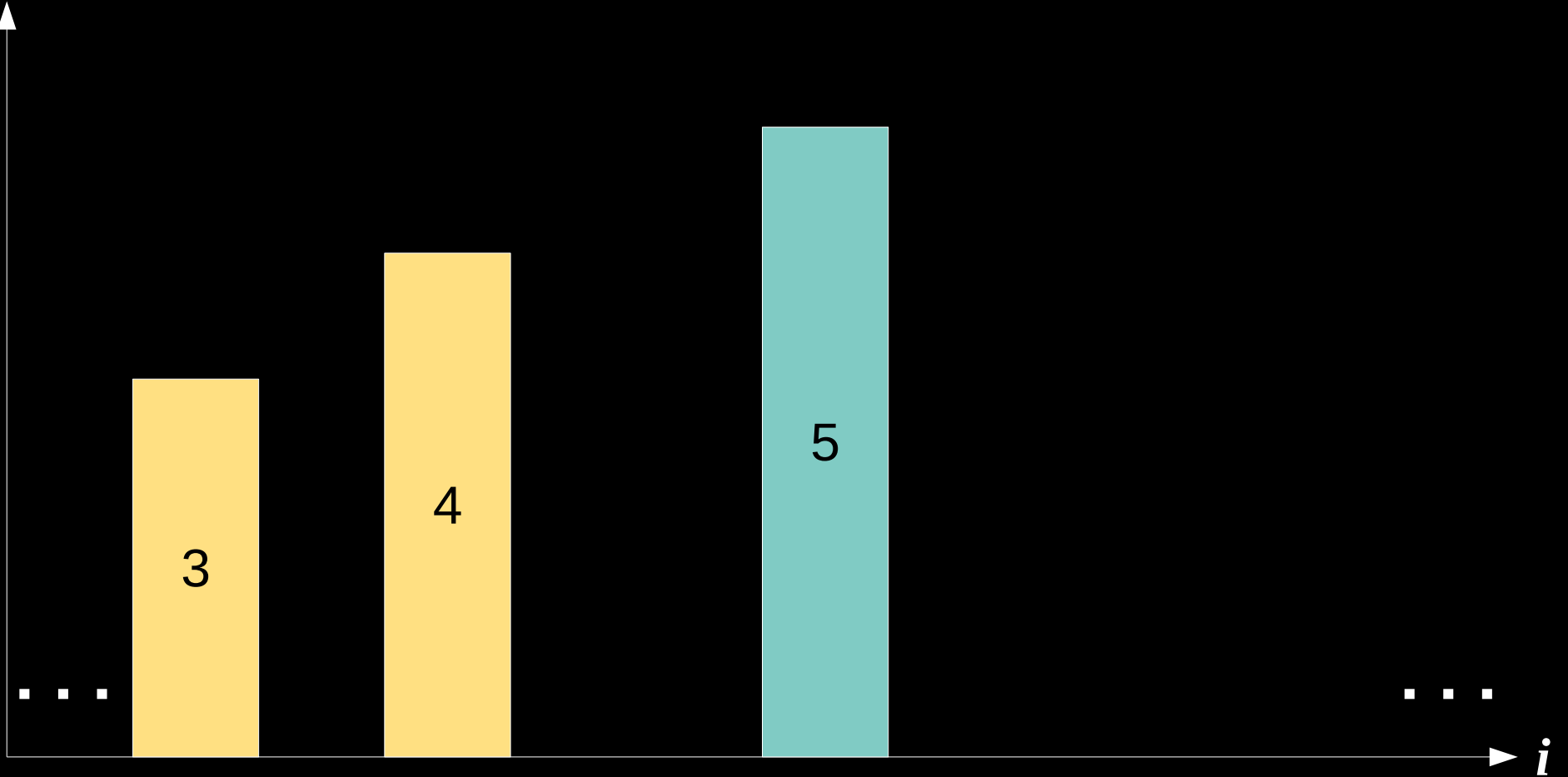
PLCP[ $i$ ]



山  $> 5$  がない

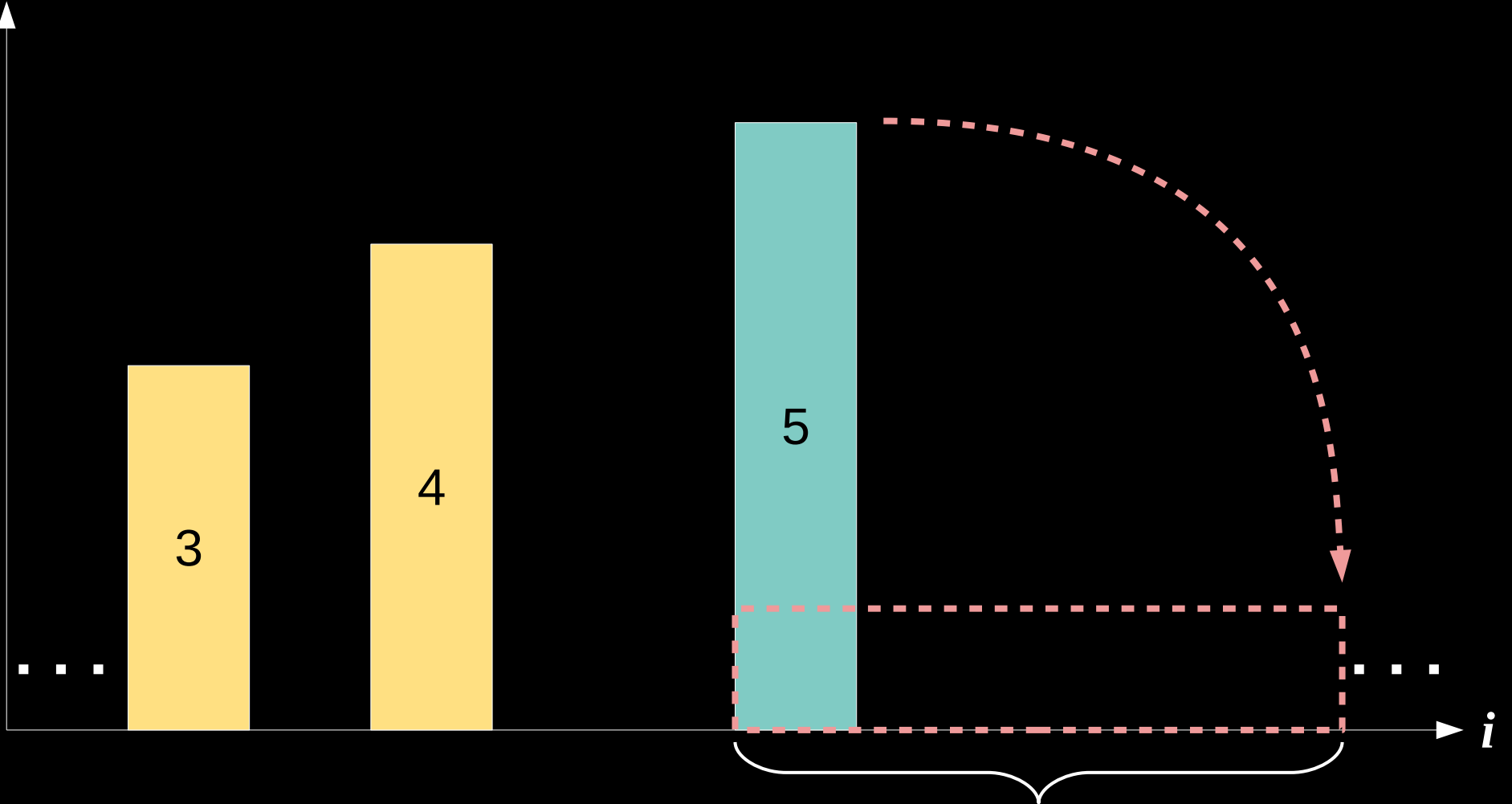


PLCP[ $i$ ]



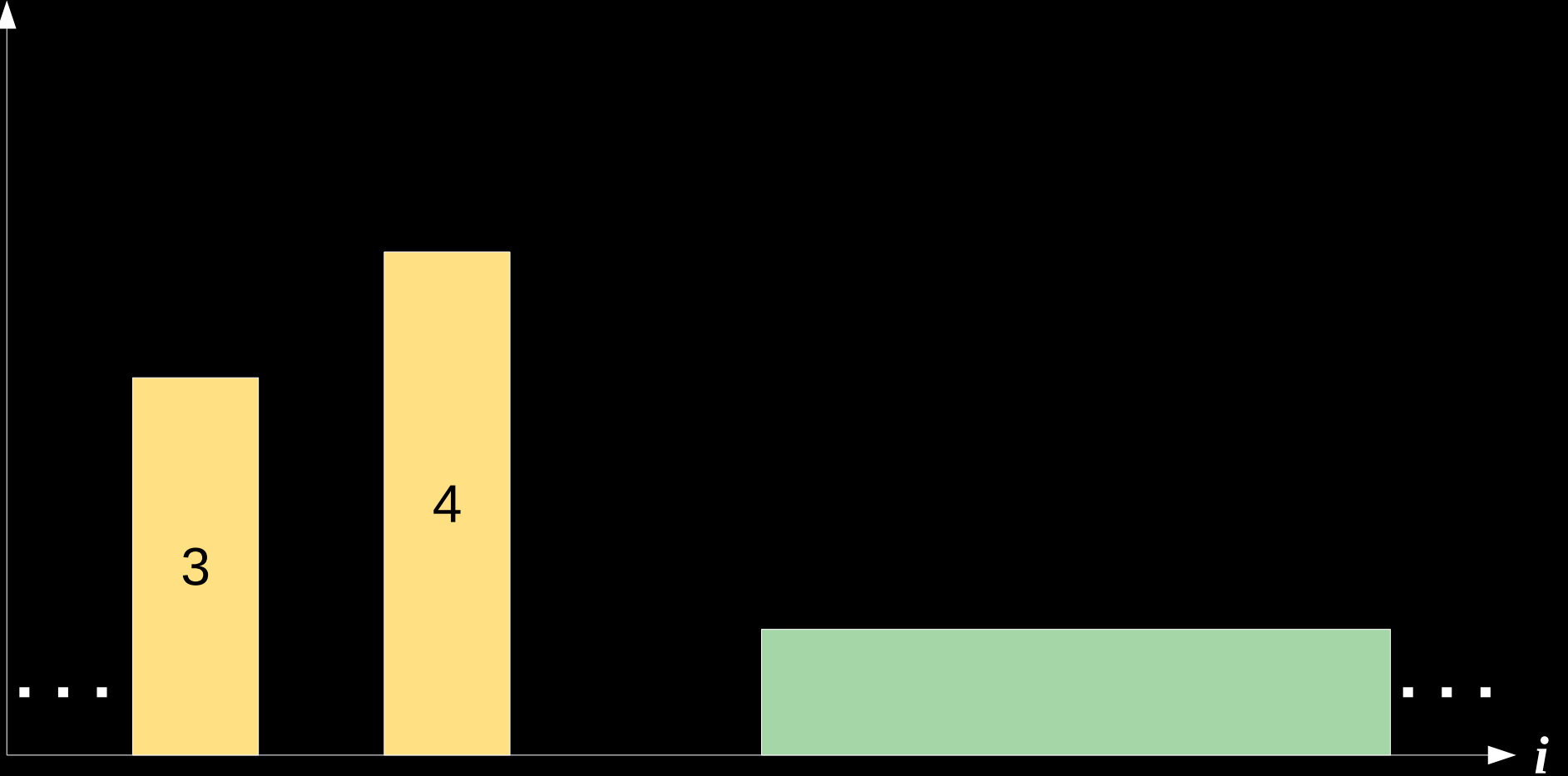
一番高い山

PLCP[ $i$ ]

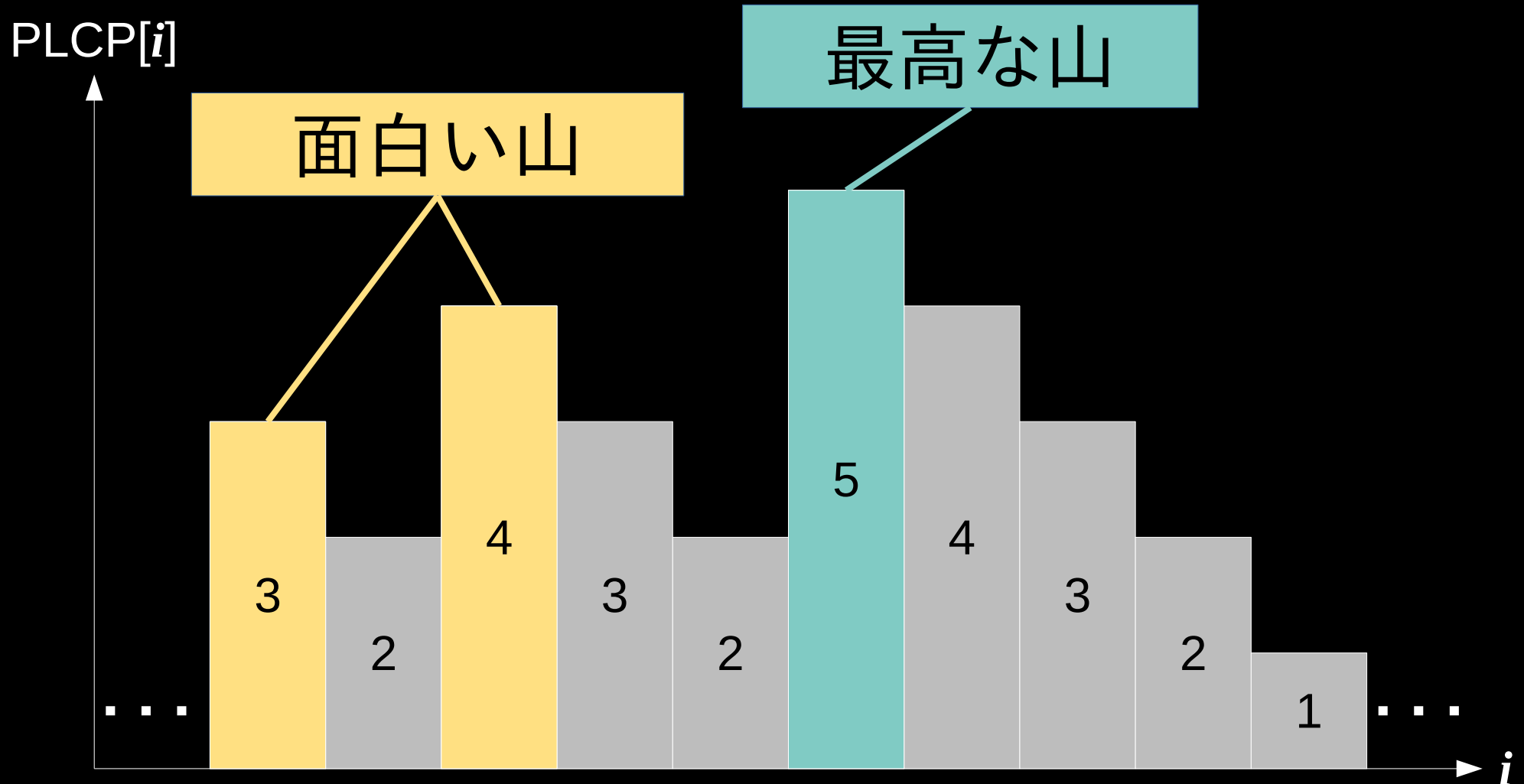


参照の長さは 5

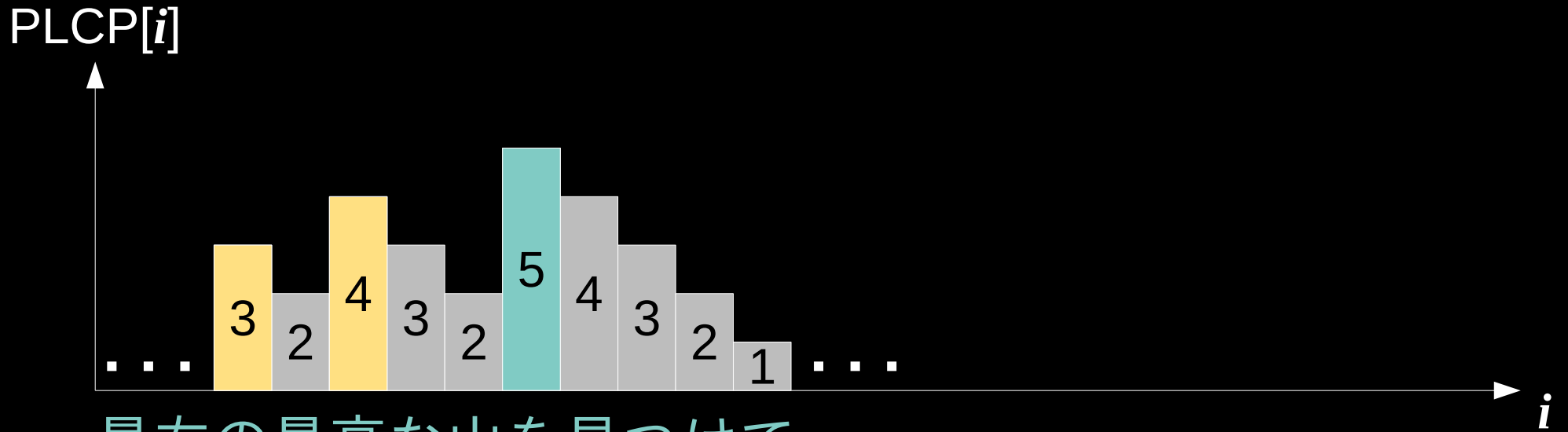
PLCP[ $i$ ]



# 2つの定義



# 繰り返して



最左の最高な山を見つけて

この接頭辞で plcpcomp を応用して

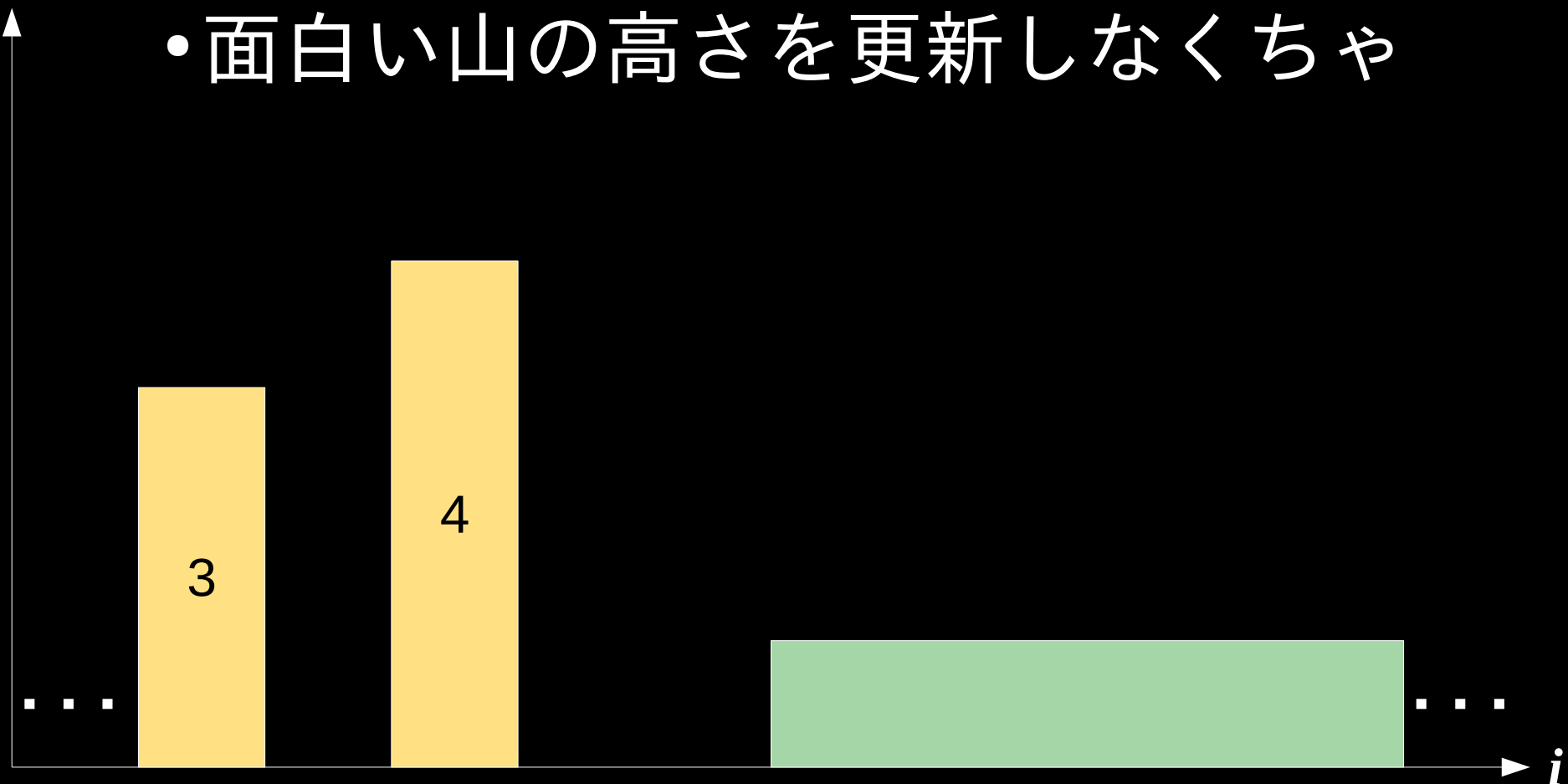
線形時間？

この接尾辞で plcpcomp を応用して

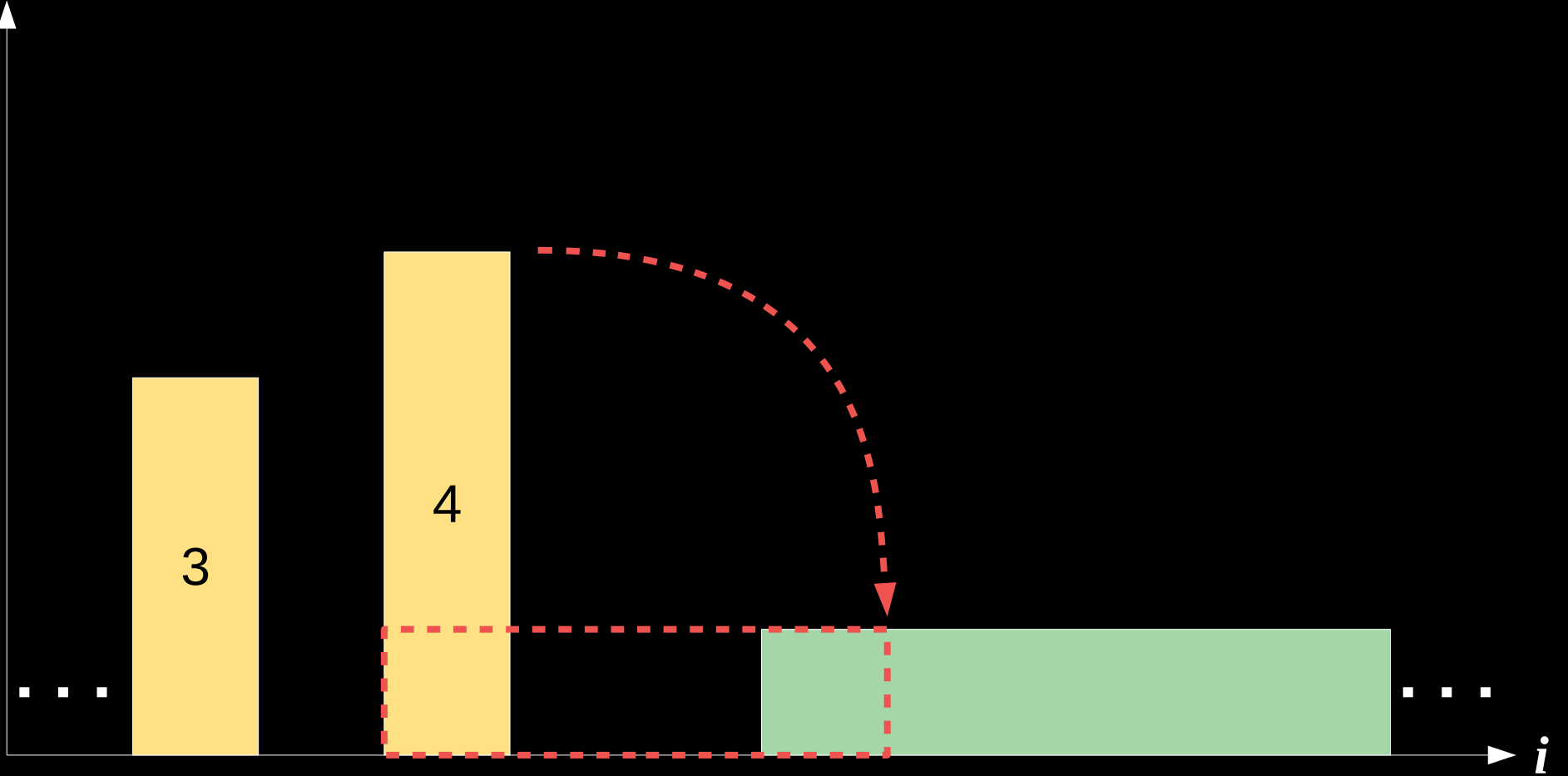
面白い山しか見ないで！

PLCP[ $i$ ]

- 新しい最高な山は面白い山
- 面白い山の高さを更新しなくちゃ

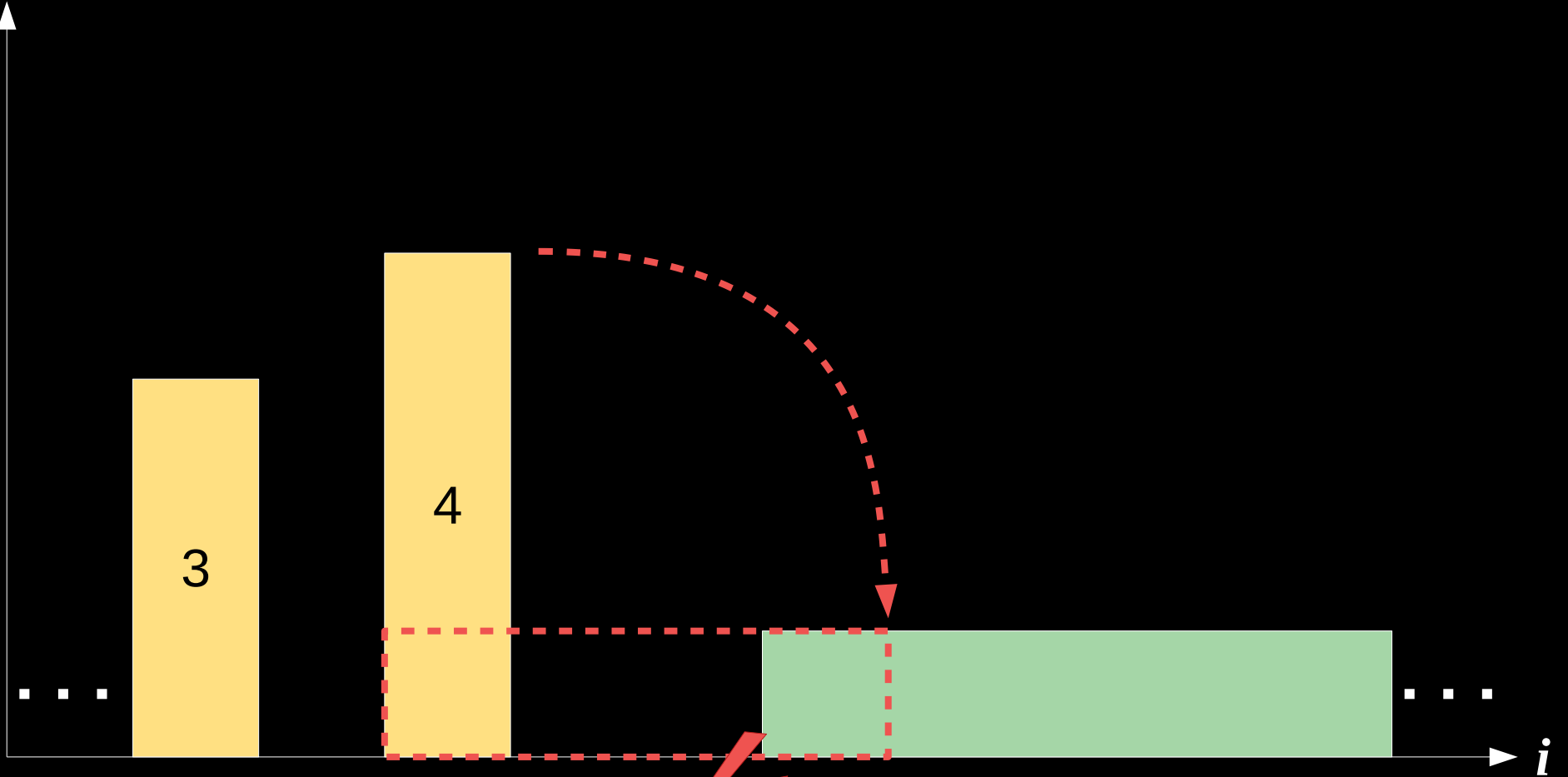


PLCP[ $i$ ]



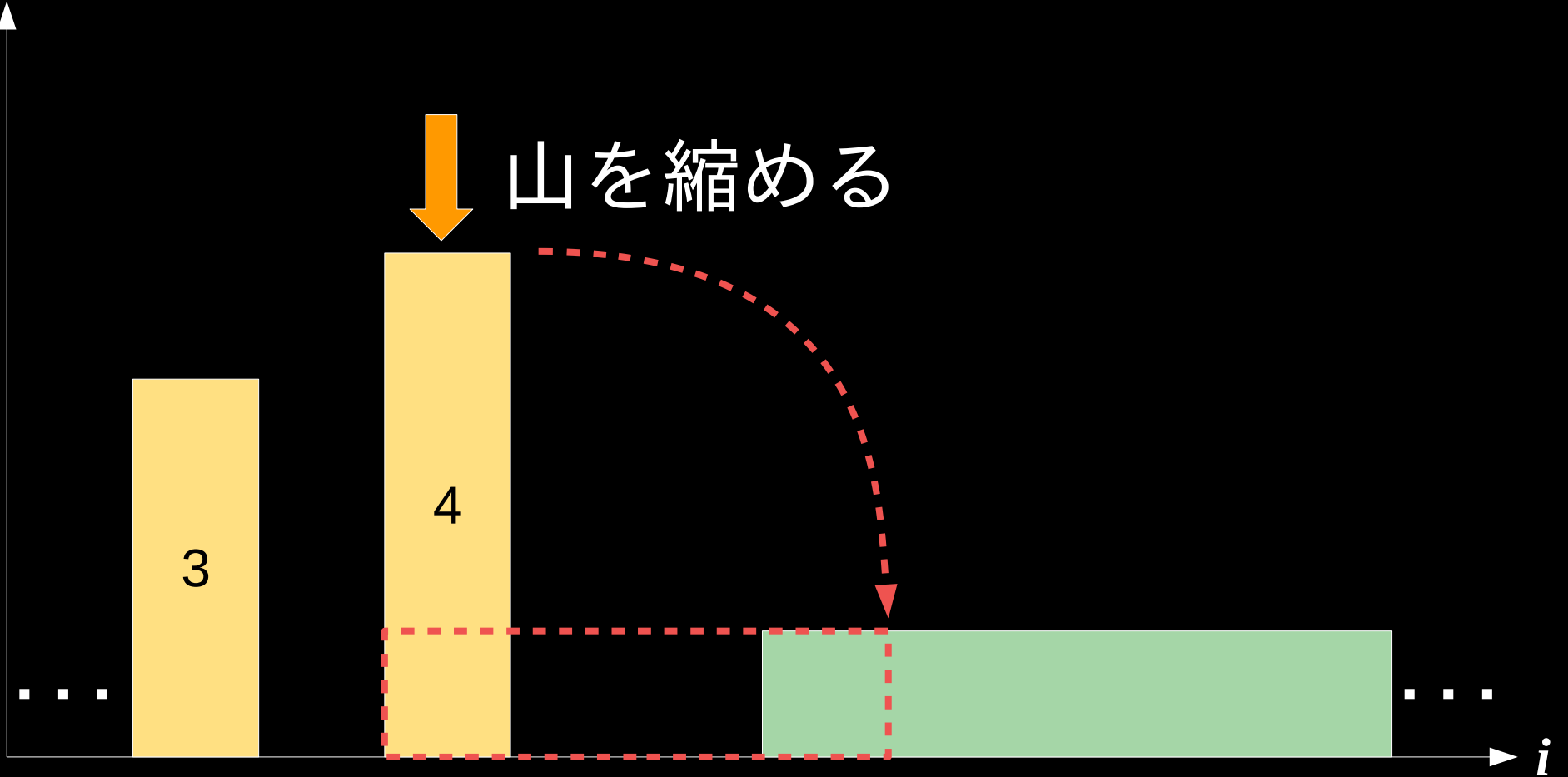


PLCP[ $i$ ]

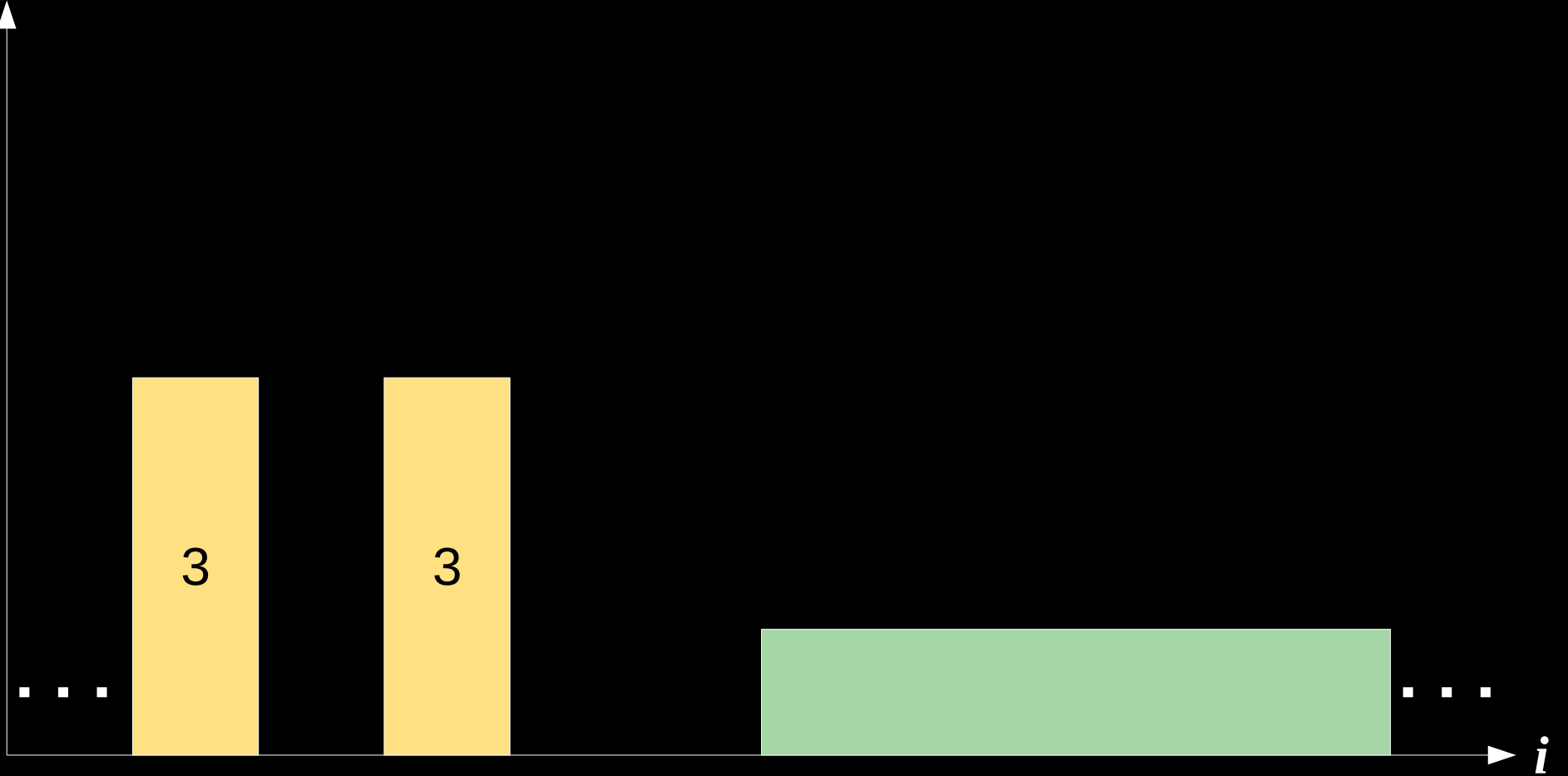


⚡ 1文字の重なり

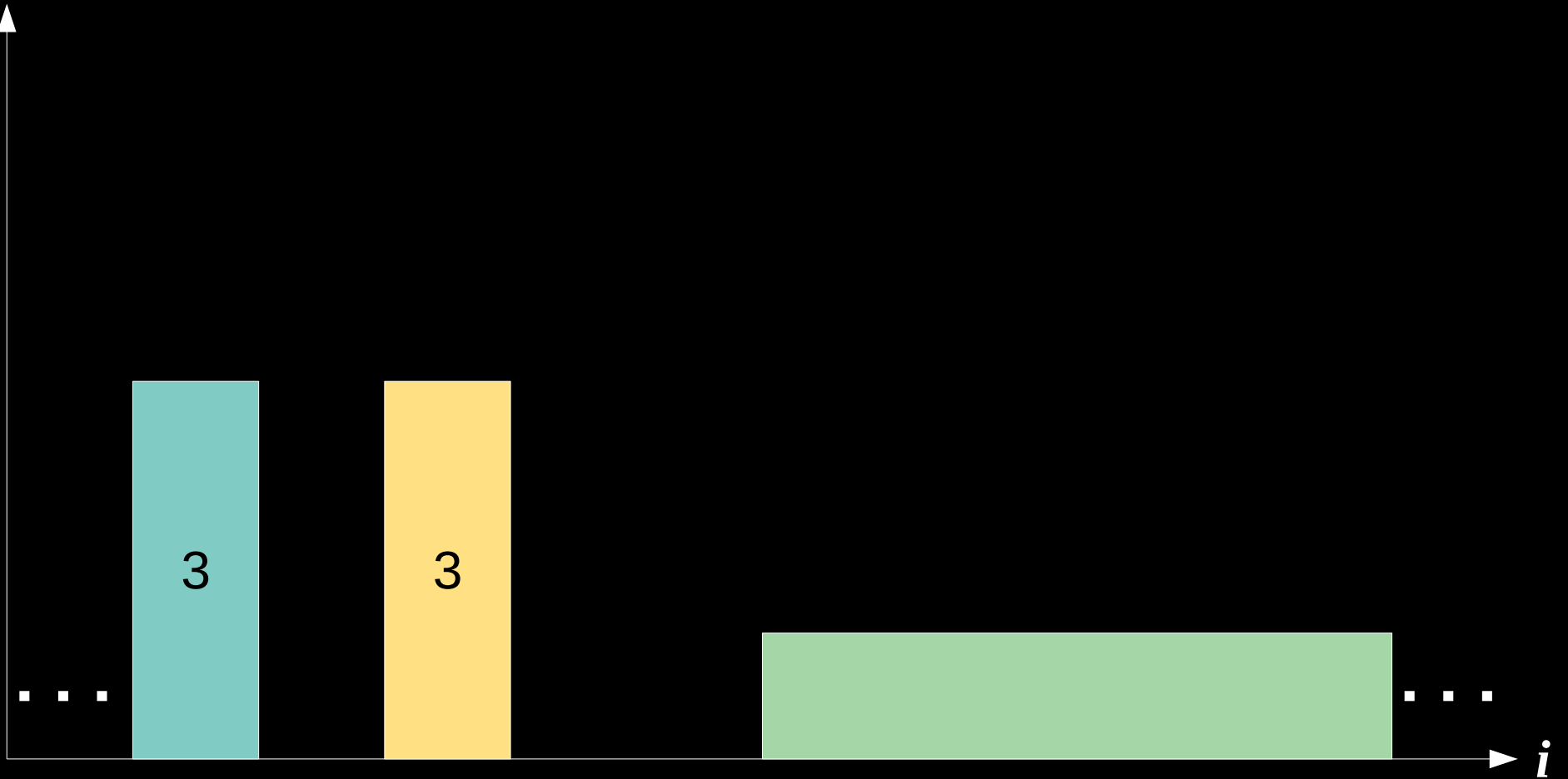
PLCP[ $i$ ]



PLCP[ $i$ ]

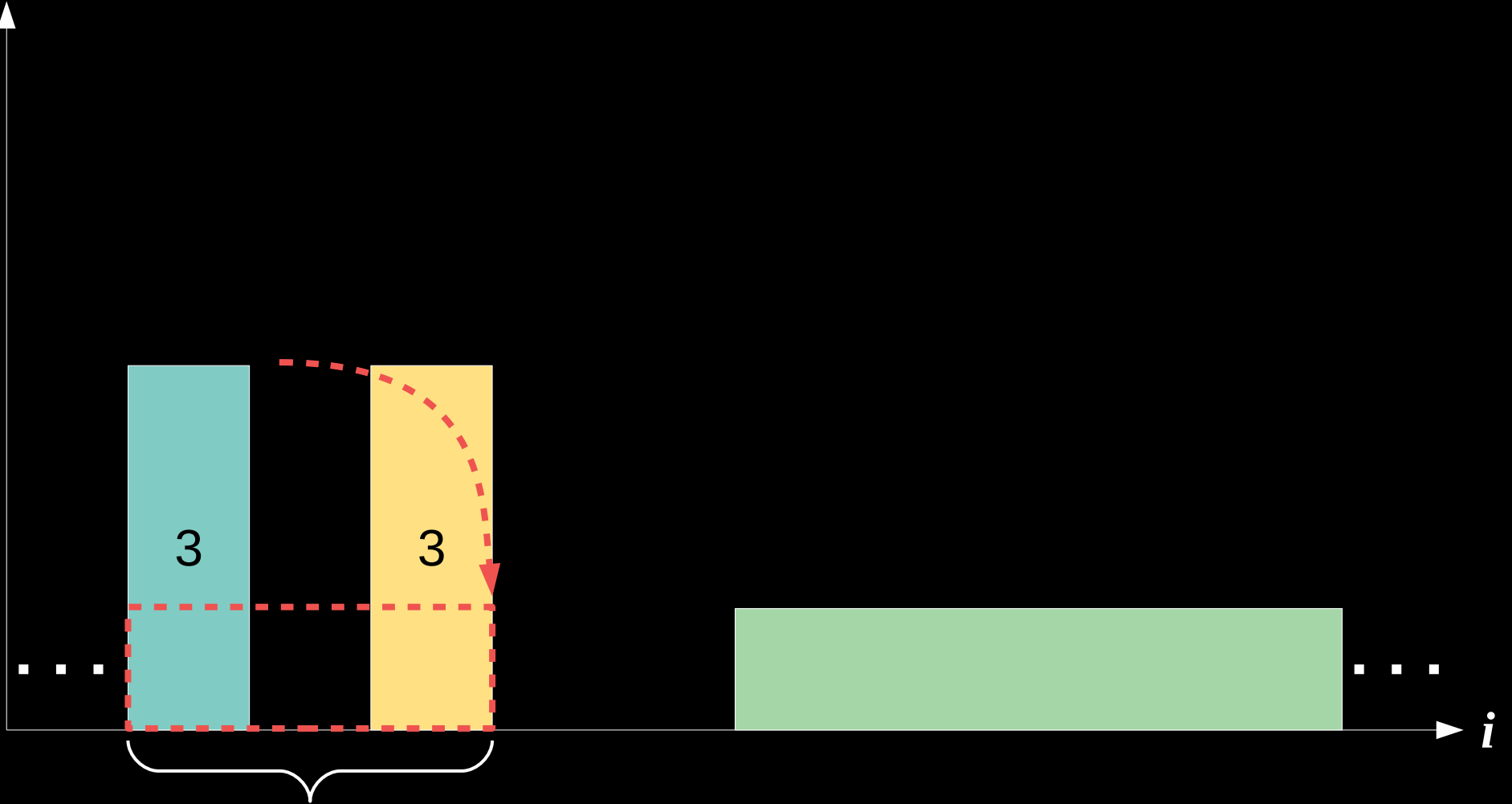


PLCP[ $i$ ]



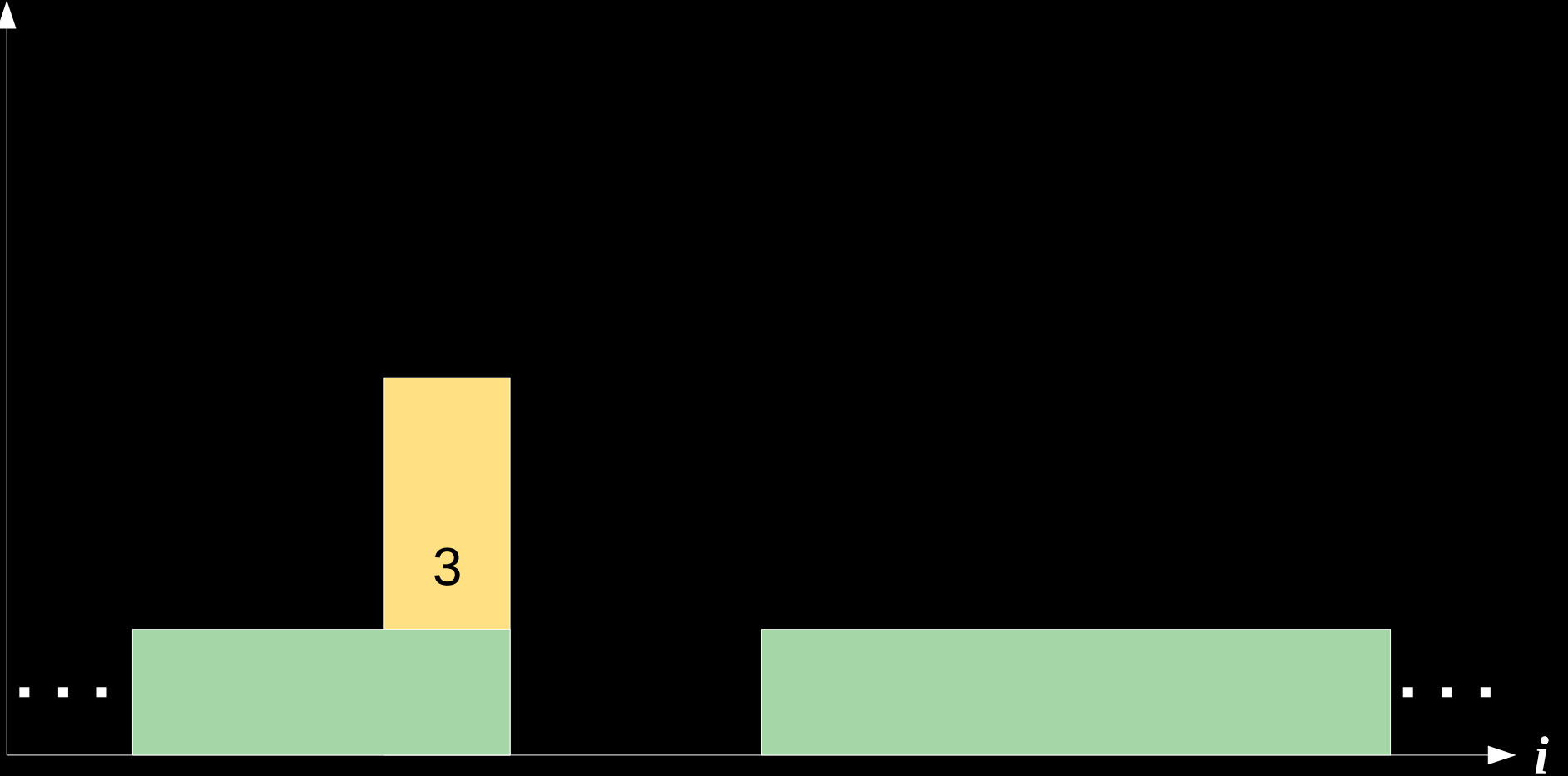
新しい最左の最高な山

PLCP[ $i$ ]

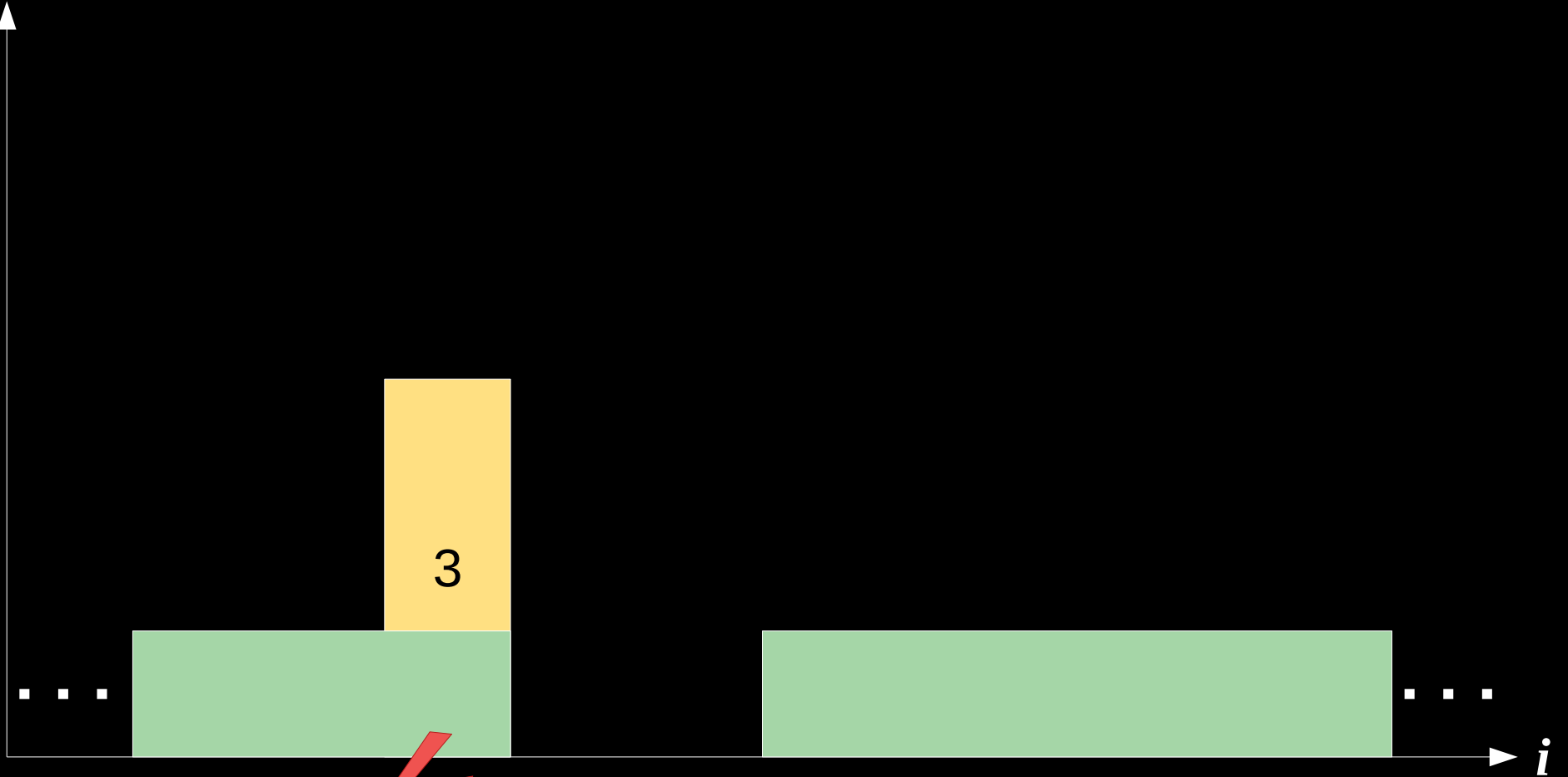


参照の長さは 3

PLCP[*i*]

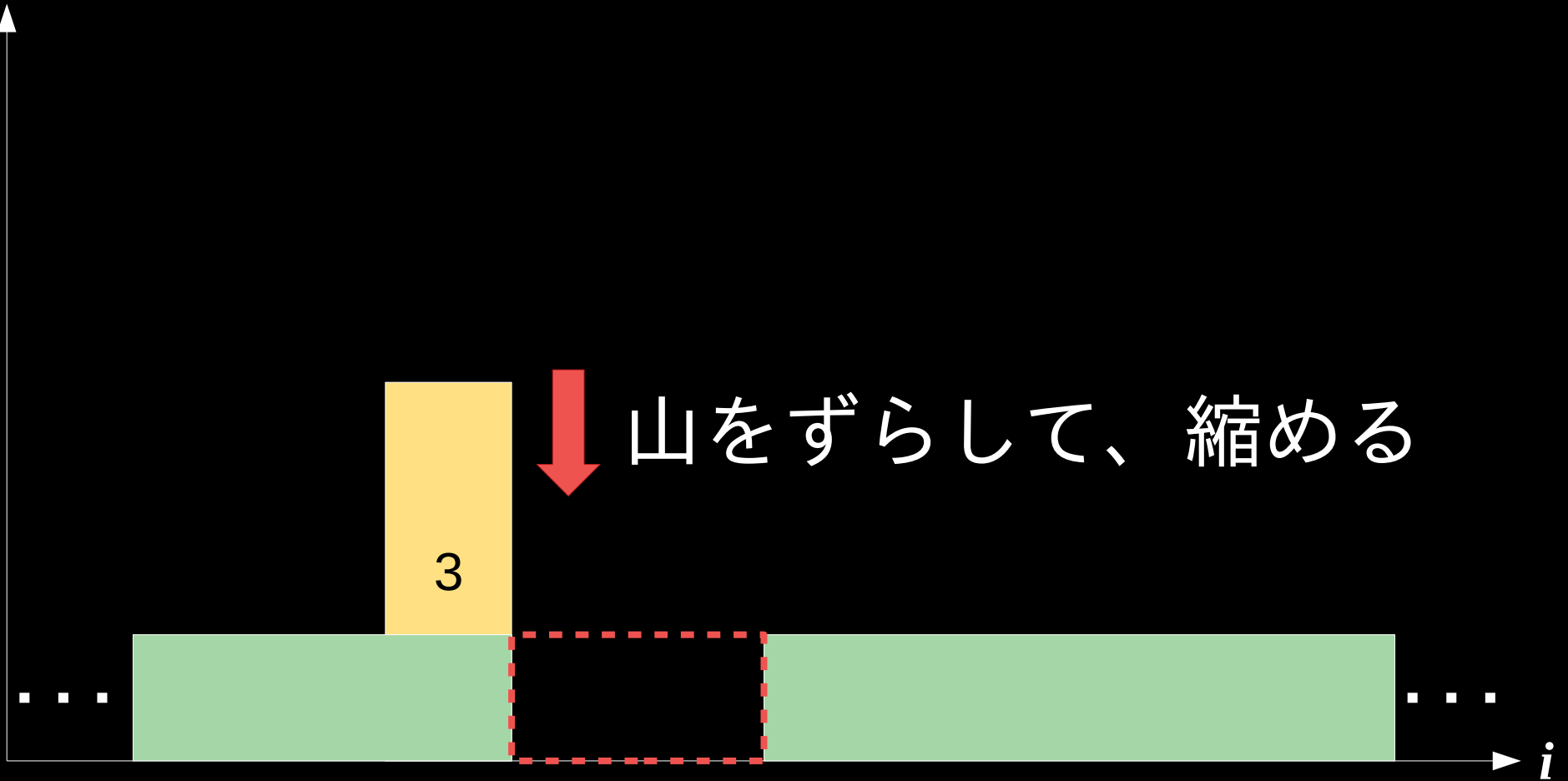


PLCP[ $i$ ]



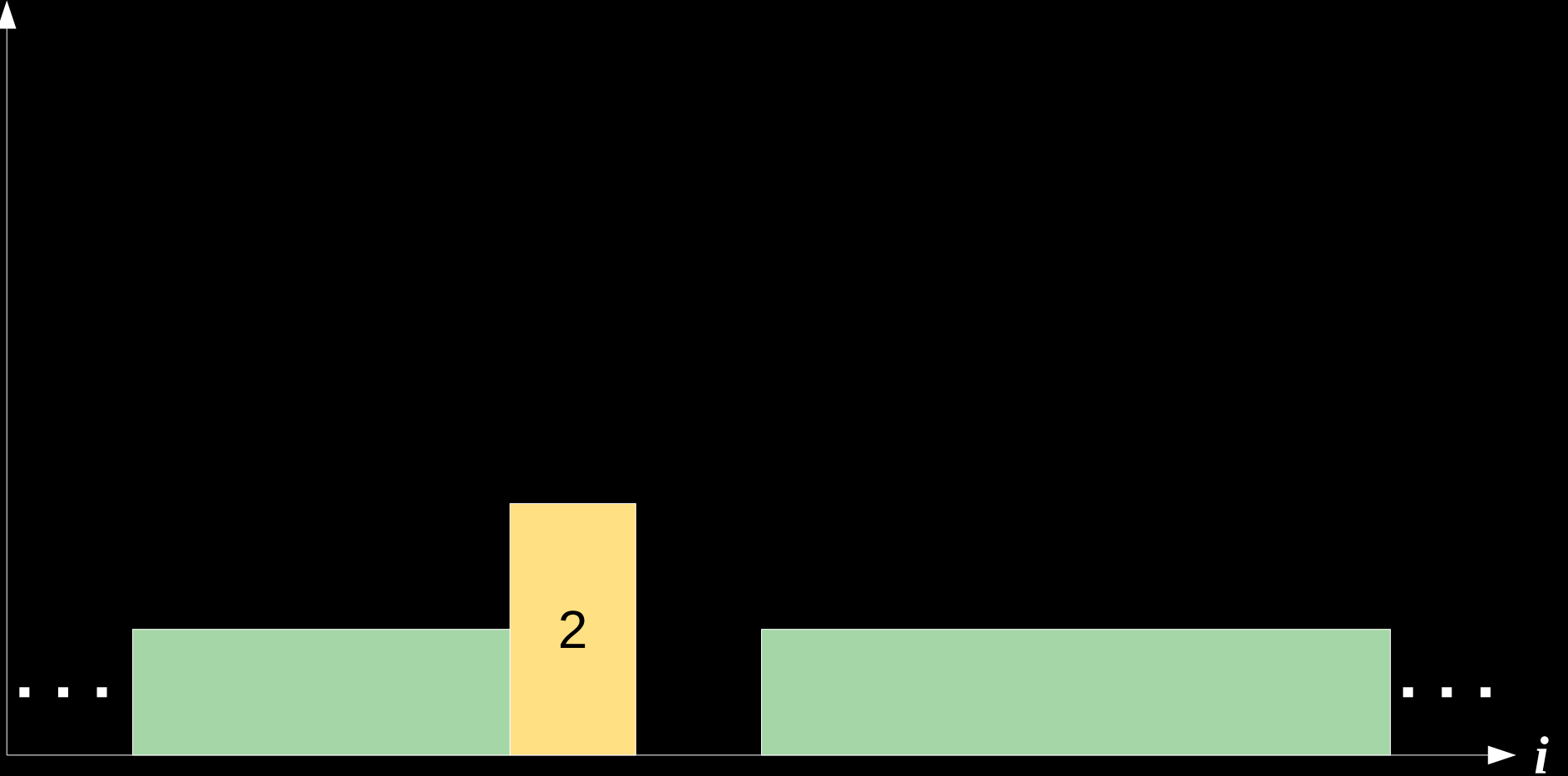
⚡ 右側の重なり

PLCP[ $i$ ]

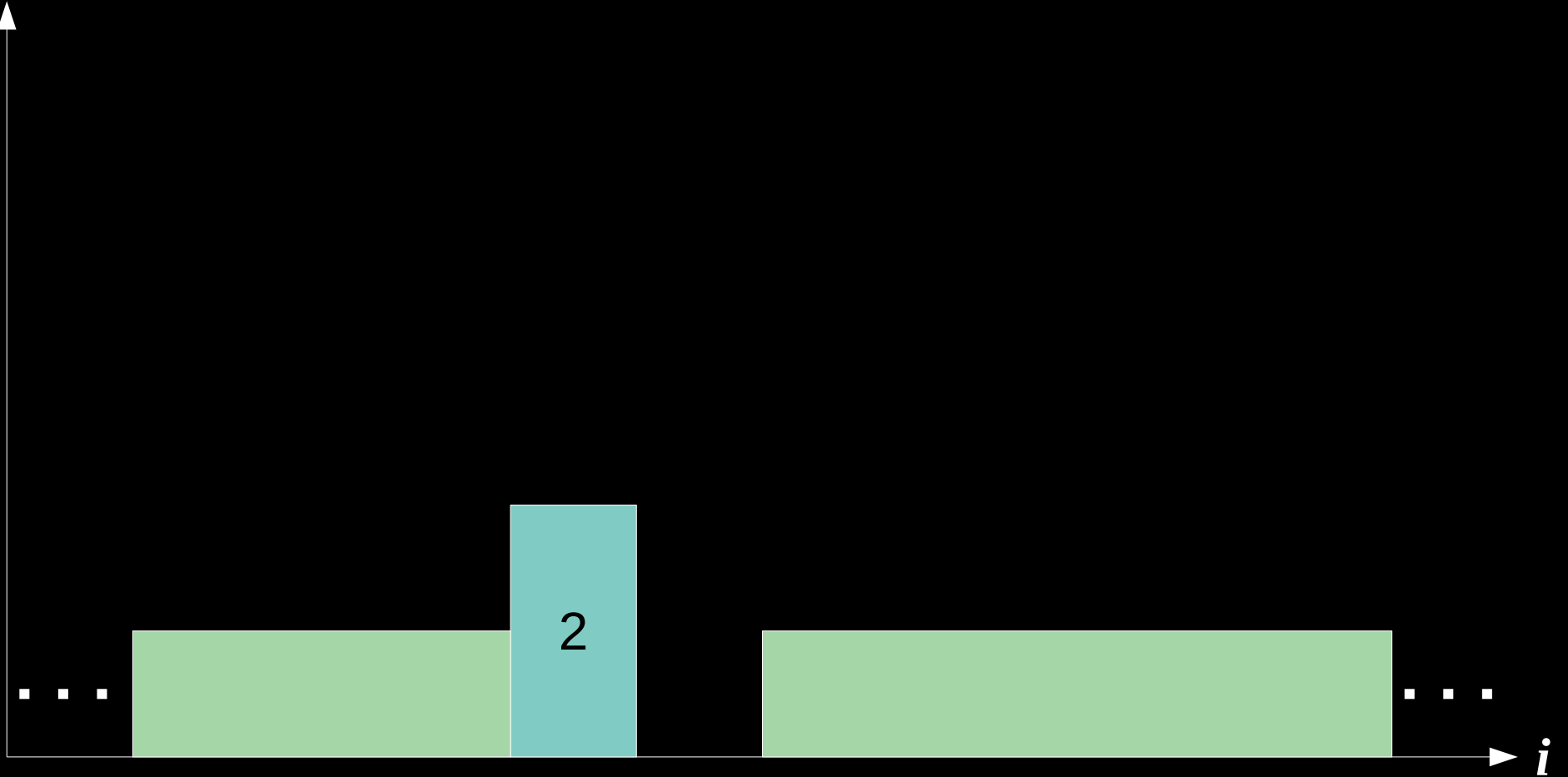




PLCP[*i*]

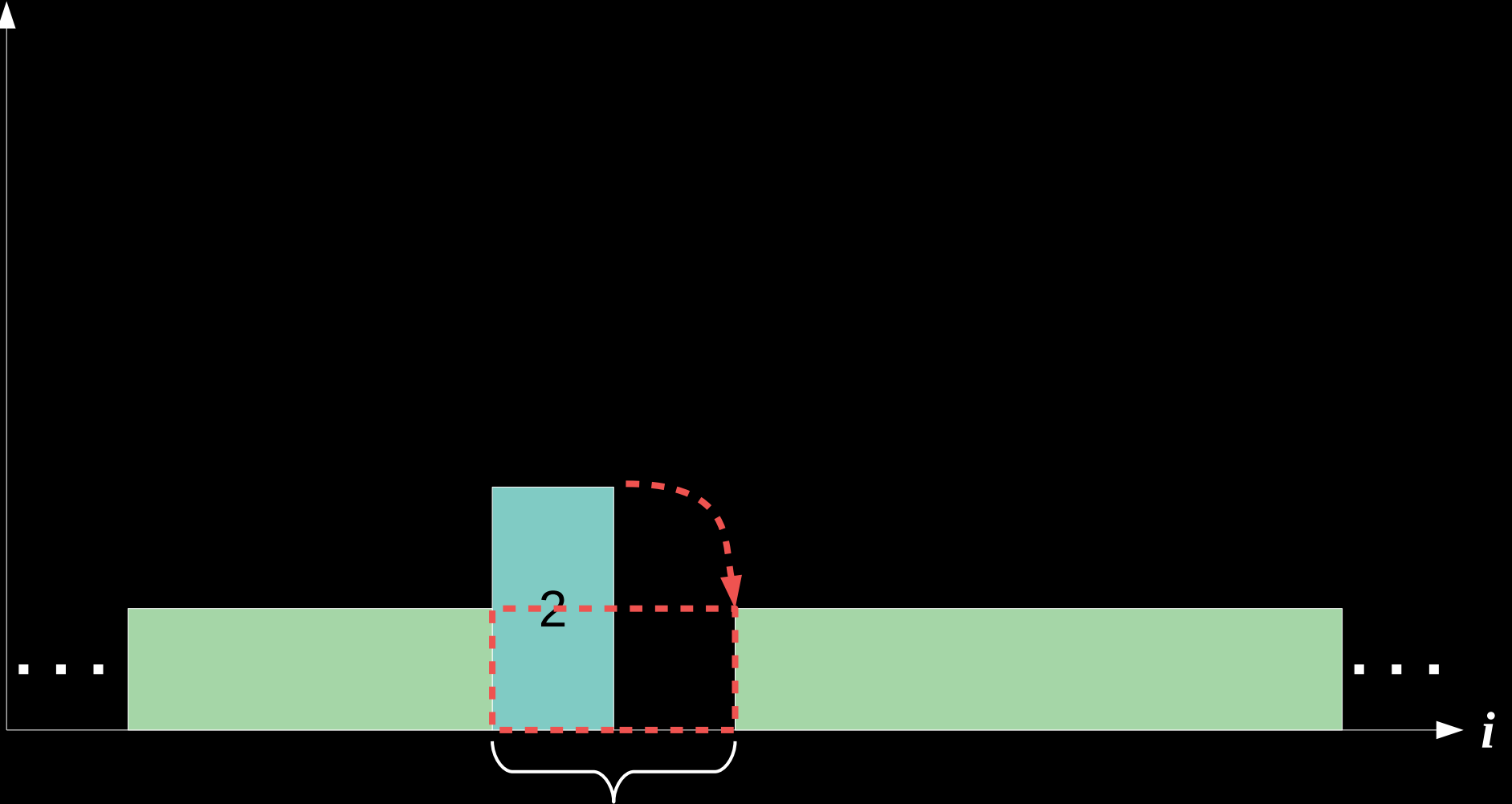


PLCP[ $i$ ]



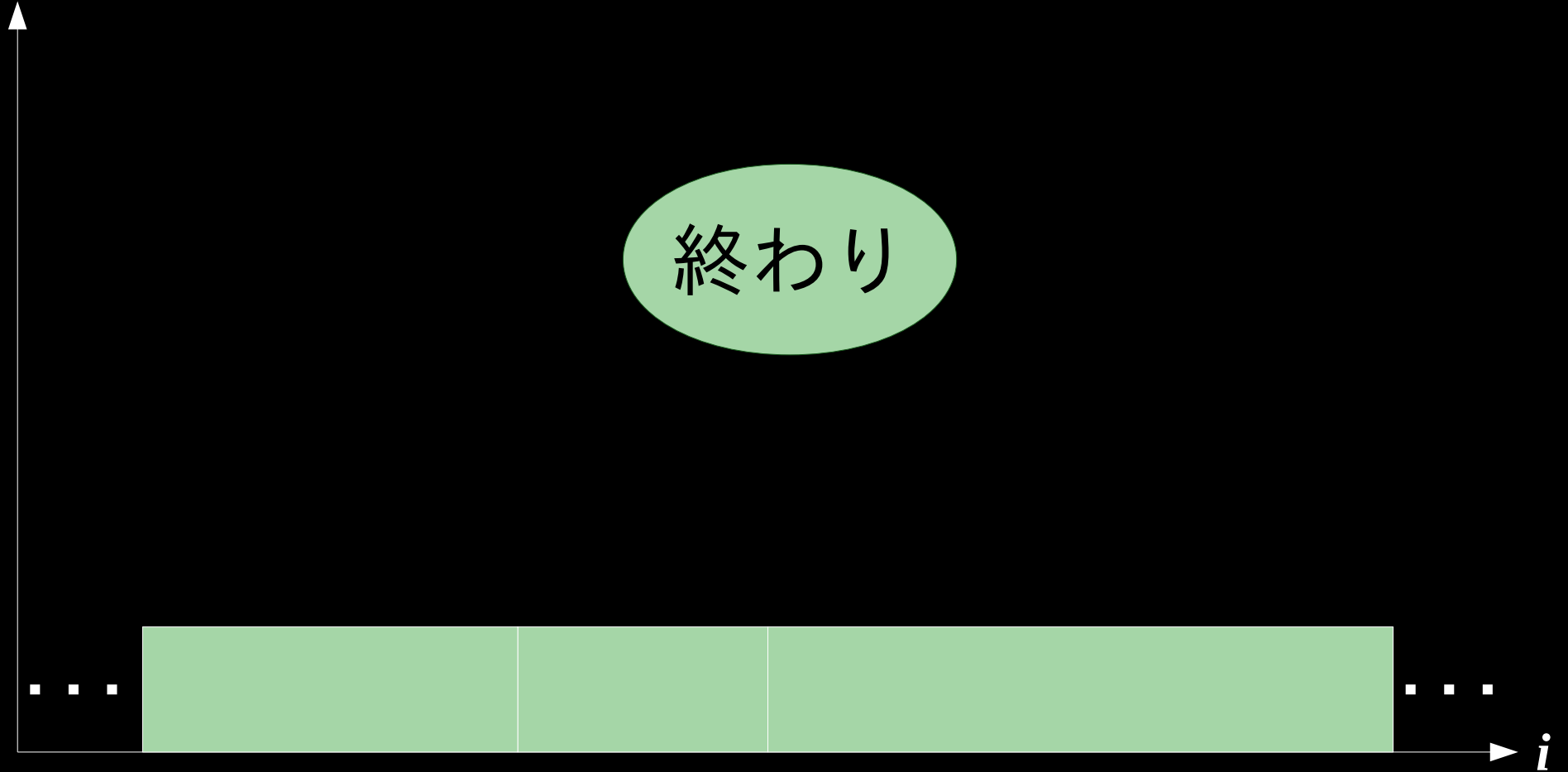
新しい最左の最高な山

PLCP[ $i$ ]

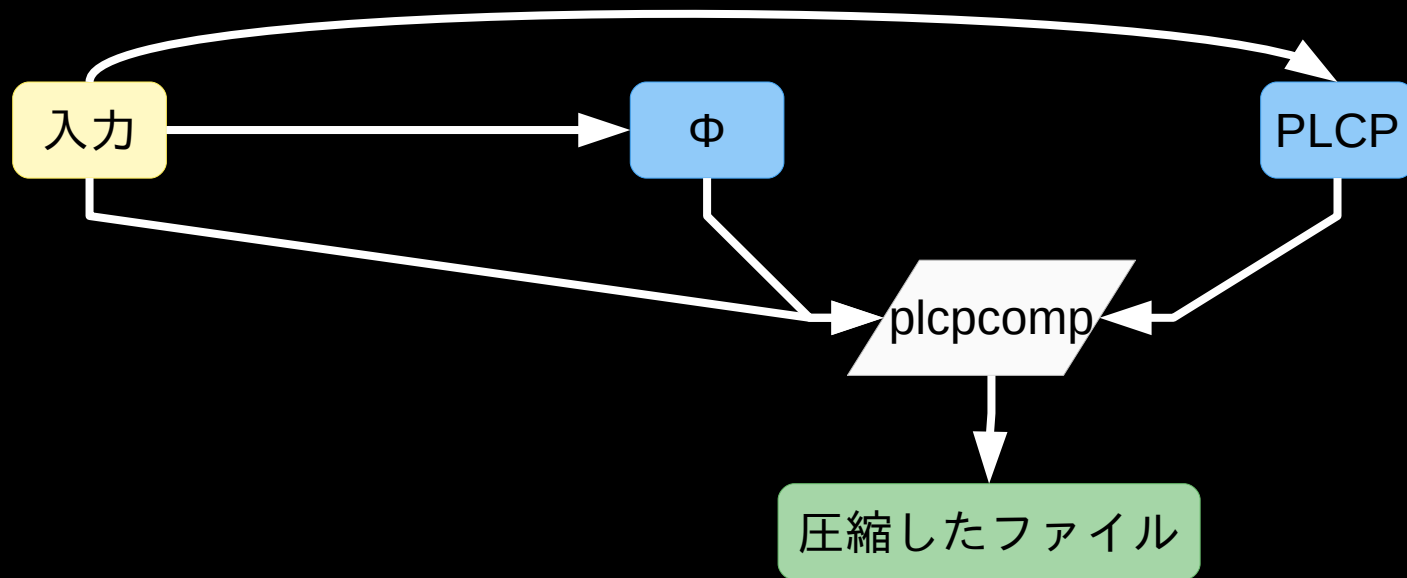


参照の長さは 2

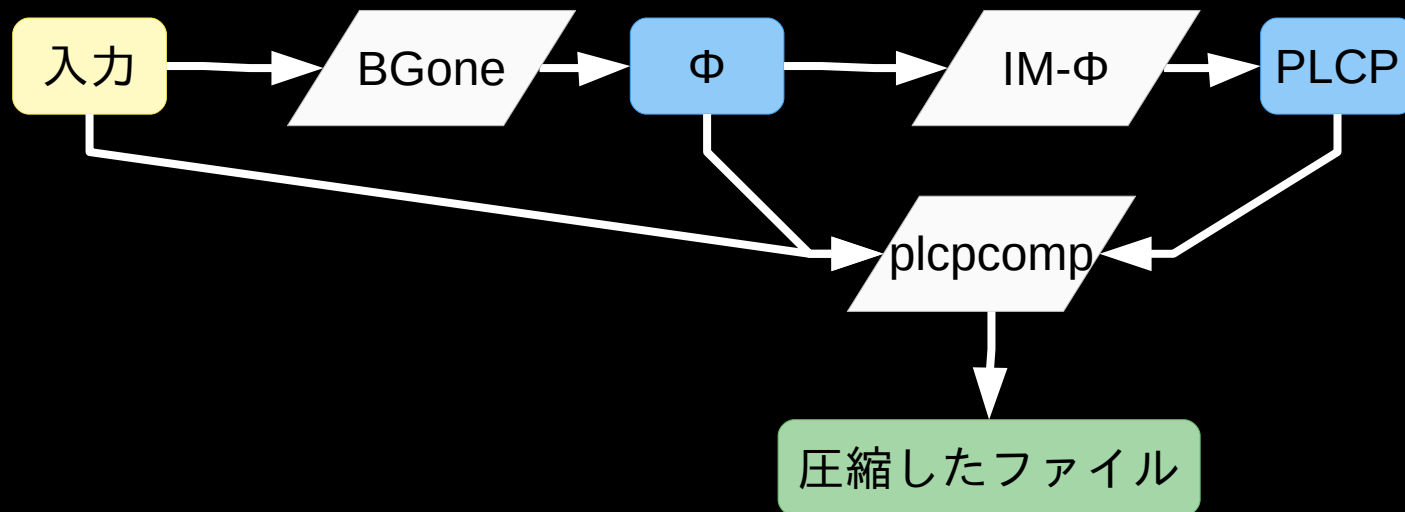
PLCP[ $i$ ]



# データ構造

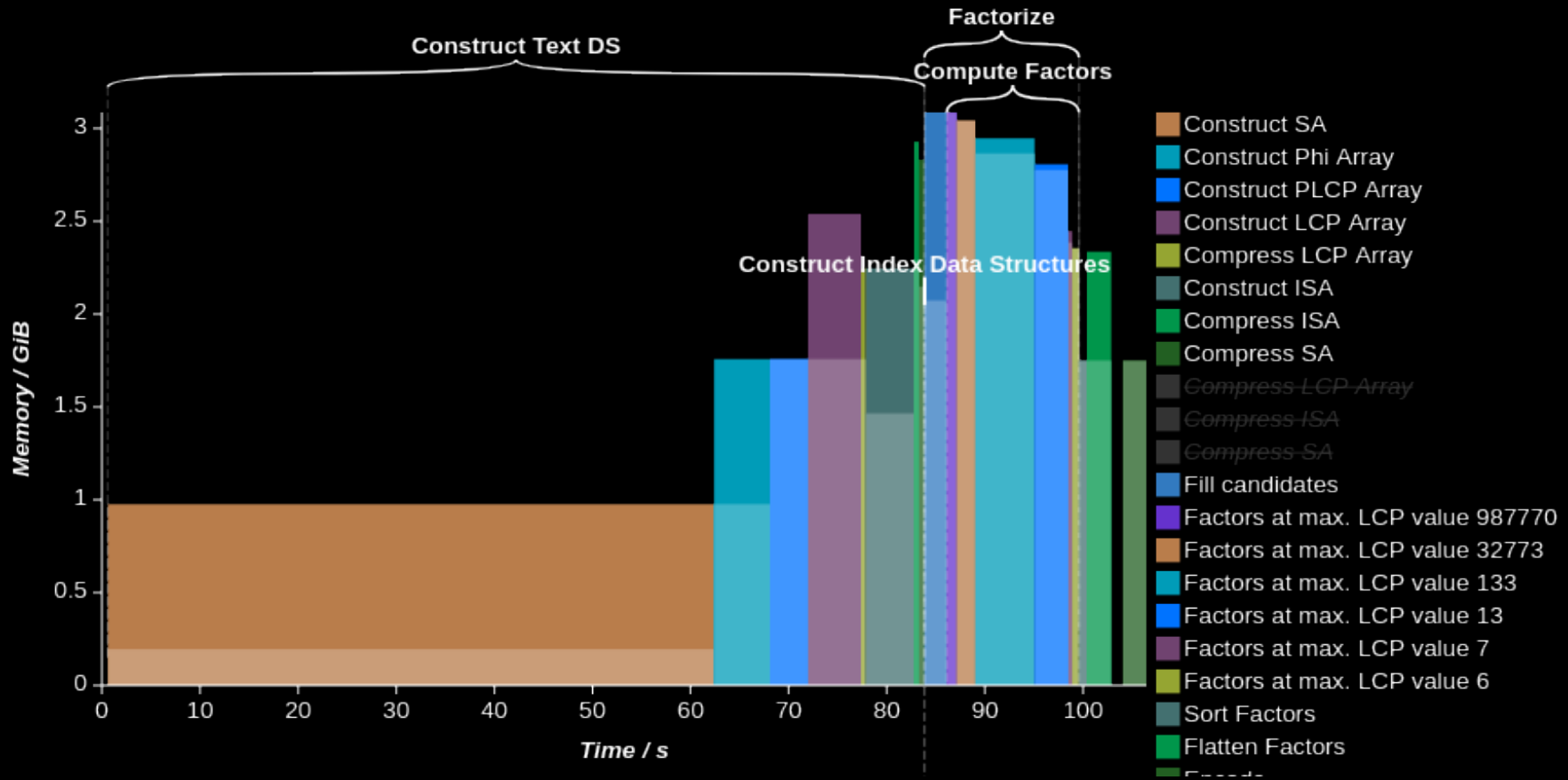


# データ構造とアルゴリズム

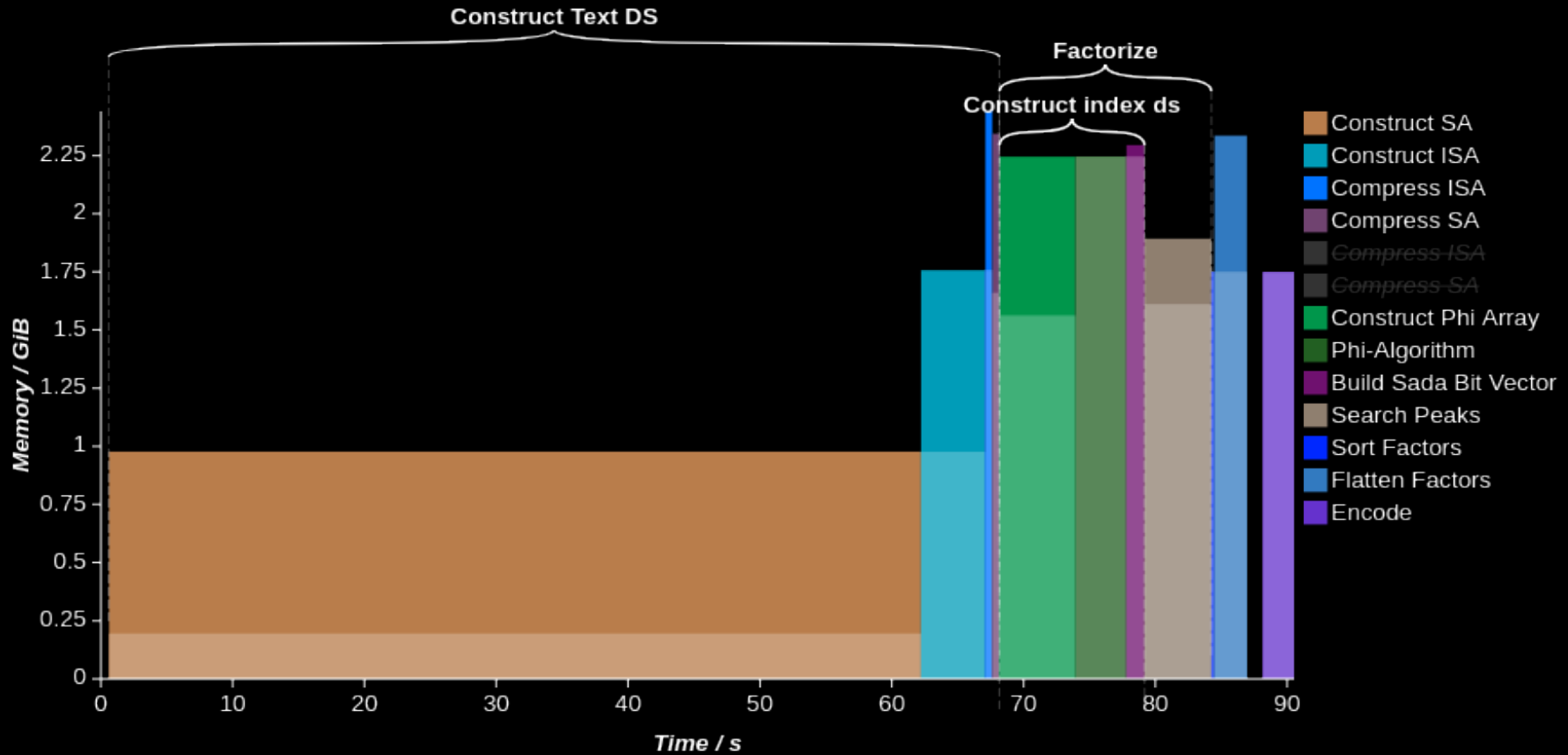


- BGone: 後藤, 坂内 '14
  - IM-Φ: Kärkkäinen+'09
- } 線形時間

# lcpcomp



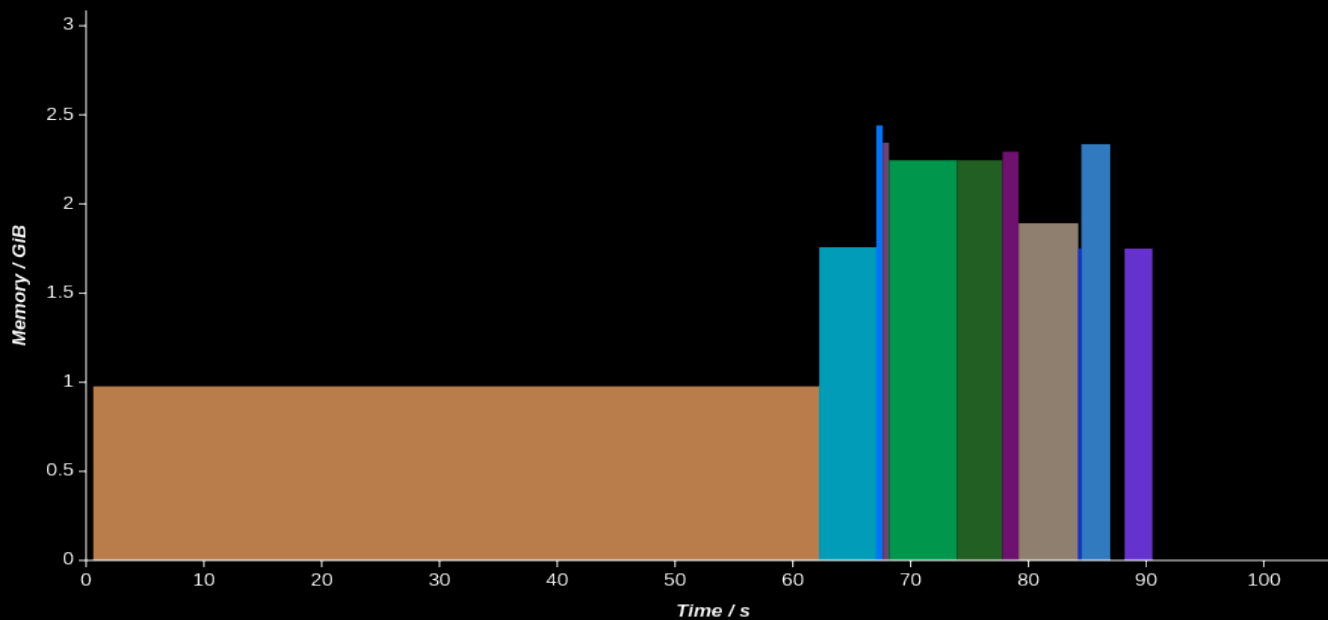
# plcpcomp



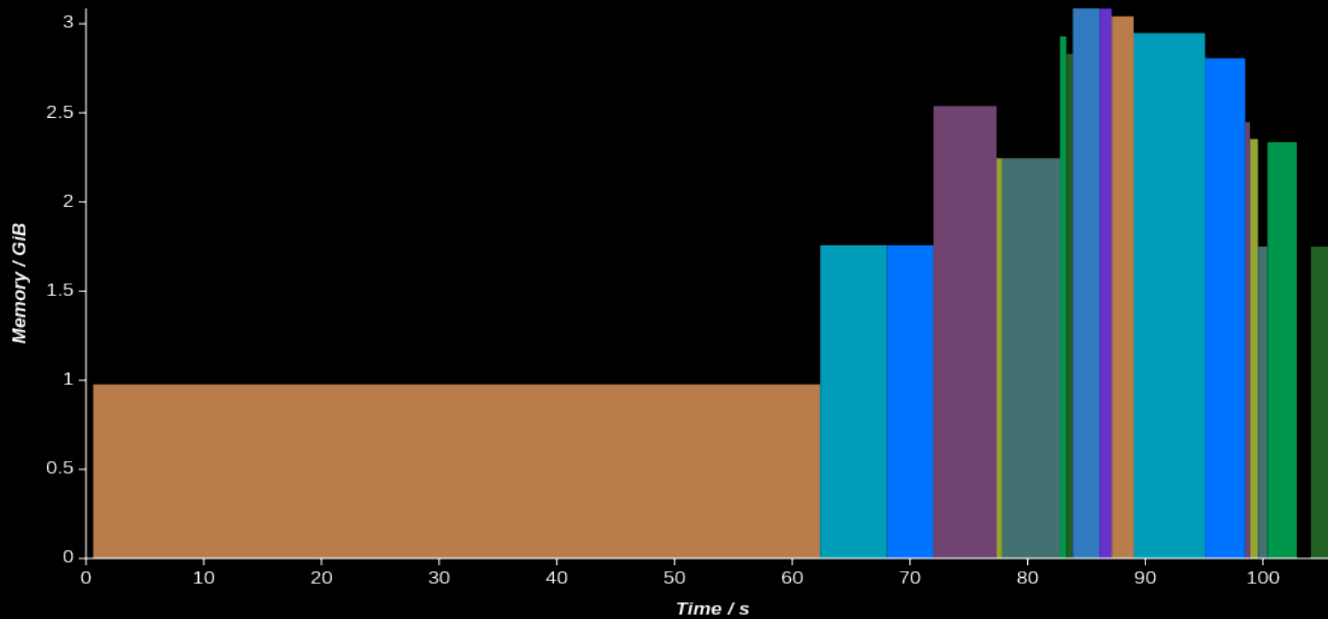


# 見た目の比較

plcpcomp



lcpcomp



# まとめ

## plcpcomp

- bidirectional 圧縮
- 線形時間
- lcpcomp より実際に早い
  - テキスト、PLCP、 $\Phi$  を線形に走査できる

まだ説明していないこと：

- 外部メモリアルゴリズム (EM) もある
- EM で LZ77 のような方法の中では一番早い