

# 情報構造 第五回

ソーティングアルゴリズム

# 今日の予定

- ソーティングとは
- ソーティングアルゴリズム
  - 3つの単純法
  - ヒープソート

# ソートイングとは

- 整列問題

- 入力：n個の数
- 出力：小さい（大きい）順の列

ソートイング例

(4, 5, 2, 1) => (1, 2, 3, 4)

- なぜソートイングアルゴリズムを学ぶのか

- 計算機が遅い頃、ソートを行うために効率の良いアルゴリズムが考えられてきた
  - 現在では、そんなに頑張らなくても良い
  - ライブラリの充実、メモリが大きい、オブジェクト指向
  - 今後、自分でソートアルゴリズムを実装することはないでしょう．．．
- ではなぜ？
- 多くのアルゴリズムが存在する
  - 同じタスクのアルゴリズムの違いを体感できる
  - 計算量の違い

# ソーティングの前準備

- 集合  $A_1, A_2, A_3, \dots, A_n$  を, 項目中のキーの順序で並び替えること
  - 要素  $A_i$ : 型 item

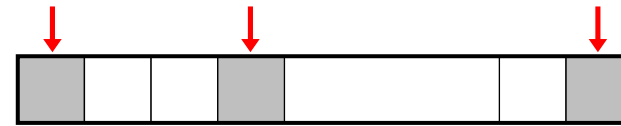
```
typedef struct{
    int key;
    ...
} item;
```
- ソートされる項目は, 要素が item 型, 要素数が  $n$ , すなわち index が 0 から  $n-1$  の配列  $A$  に格納
  - item  $A[i]$ ;
- 第  $i$  項目のキーは,  $A[i].key$  でアクセス
- 時間計算量:
  - ソートする項目の個数 (入力サイズ)  $n$  の関数  $T(n)$  (以下の総和)
    - キーの比較回数
    - 項目の代入回数

# ソーティング法

- 内部ソーティング法

- 高速のランダムアクセスできる内部記憶（主記憶など）上のソート=> 一つの配列

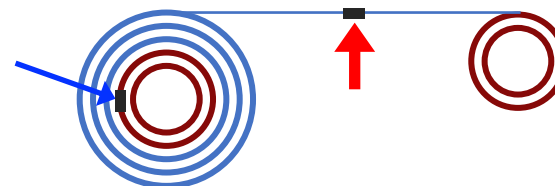
- 単純法
- シェルソート
- ヒープソート
- クイックソート
- 木ソート
- 基数ソート



どこでも同じ時間でアクセス

- 外部ソーティング法

- 大きな領域を持ち、順次アクセスしかできない外部記憶（磁気テープなど）上のソート => 複数の配列（たくさんのメモリを使う）
- マージソート



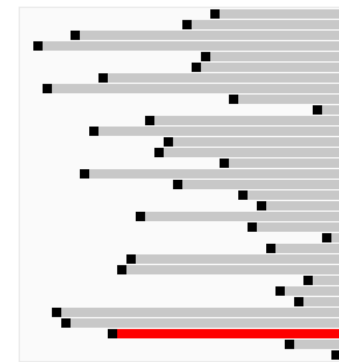
どこかひとつずつ順次アクセス

3つの単純法

# 単純法

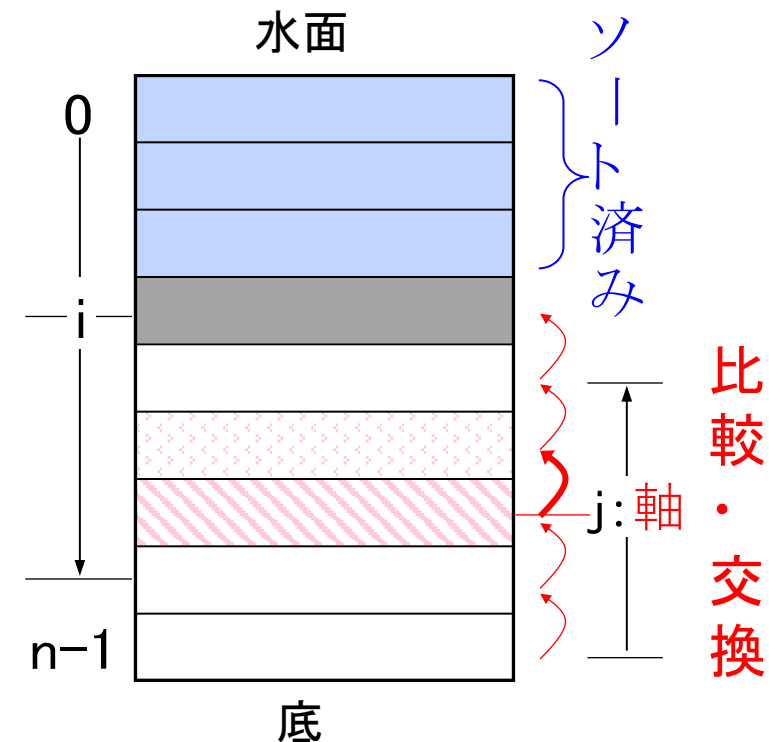
- 3種類の方法
  - 単純交換法（バブルソート）
  - 単純選択法
  - 単純挿入法
- いずれも時間計算量は
  - $O(n^2)$

# 単純交換法（バブルソート）



バブルっぽい図（大きい値を上へ）  
Wikipediaより

- アルゴリズム
  1. 一番下（底）の要素を軸
  2. ひとつ上の要素と比較
    - 大きい時：上下を**交換**する
  3. 軸を一つ上に
    1. 水面ではないなら：2へ
    2. 水面なら：一番上を除いて1.へ
- 泡が上っていくように見える



バブルソートの第iフェーズ



# バブルソート (C言語)

```
void BubbleSort(int A[]) {  
    int i,j,temp;  
    for(i=0; i<(n-1); i++)  
        for(j=n-1; j>i; j--)  
            if( A[j-1]>A[j] ){  
                temp=A[j-1];  
                A[j-1]=A[j];  
                A[j]=temp; }  
}
```

i=0～**n-2**の繰返し (iの**設定・比較: 2回**)  
j=n-1～i+1の**n-i-1**回繰返し (jの**設定・比較: 2回**)  
**比較1回** 必ず行われる  
**代入3回** 最悪時必ず交換  
**最良時**はこれらの代入は行われない

- **入力サイズ**: ソートされる配列要素数 **n**

**最悪時間計算量**  $\sum_{i=0}^{n-2} (2 + 6(n - i - 1)) = 3n^2 - n - 2$

**最良時間計算量**  $\sum_{i=0}^{n-2} (2 + 3(n - i - 1)) = \frac{3}{2}n^2 + \frac{1}{2}n - 2$

時間計算量は、 **$O(n^2)$**

領域計算量は、配列要素数のnと変数 i, j, temp より **n+3:  $O(n)$**

# バブルソートの第0フェーズ

配列 44 55 06 42 94 18 12 67

第0フェーズ

0	i→44	44	44	44	44	44	44	<u>06</u>
1	55	55	55	55	55	55	06 ←j	44
2	06	06	06	06	06	06	06 ←j	55
3	42	42	42	42	12 ←j	12	12	12
4	94	94	94	12 ←j	42	42	42	42
5	18	18	12 ←j	94	94	94	94	94
6	12	12 ←j	18	18	18	18	18	18
7	比較	67 ←j	67	67	67	67	67	67

第1ステップ

第7ステップ 結果

# バブルソートの過程

配列 44 55 06 42 94 18 12 67

0	i→44	→ <u>06</u>	06	06	06	06	06	06
1	55	i→44	→ <u>12</u>	12	12	12	12	12
2	06	55	i→44	→ <u>18</u>	18	18	44	44
3	42	12	55	i→44	→ <u>42</u>	42	55	55
4	94	42	18	55	i→44	<u>44</u>	44	44
5	18	94	42	42	55	i→55	<u>55</u>	55
6	12	18	94	→ <u>67</u>	67	67	i→67	<u>67</u>
7	67	67	67	94	94	94	94	94

第0フェーズ

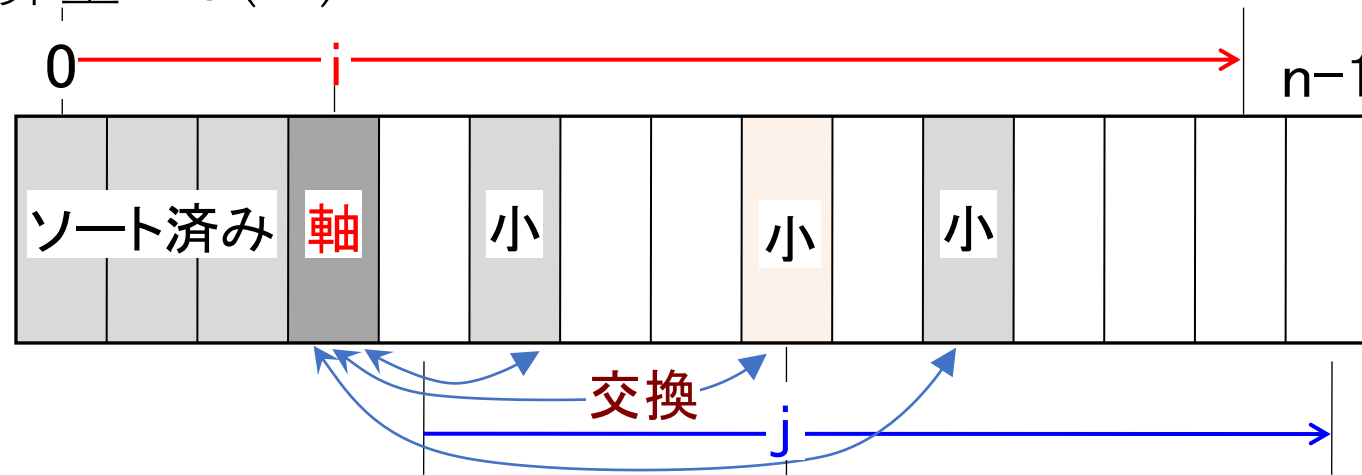
第3フェーズ

第6フェーズ

結果

# 単純選択法

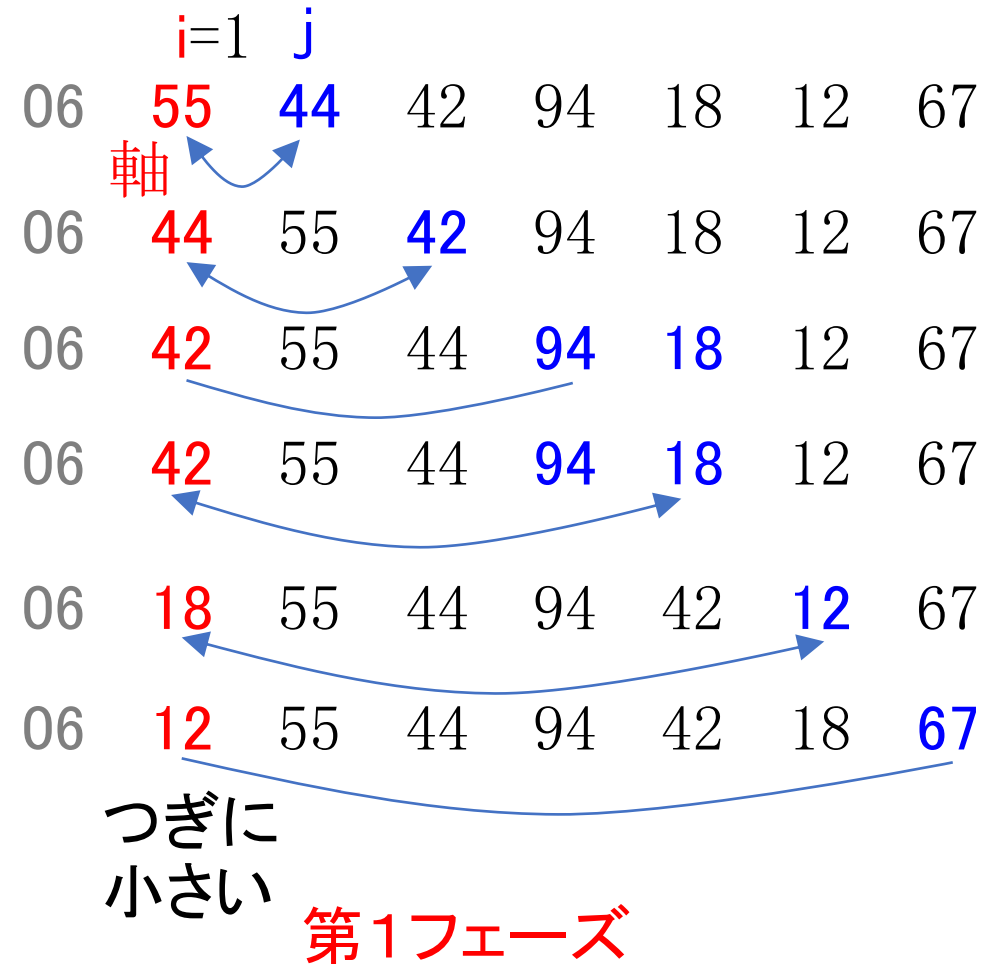
- 軸  $A[i]$  を設定し  $A[0]$  から  $A[n-1]$  まで以下を行う
- 軸  $A[i]$  以降の要素  $A[j]$  として,  $A[i+1]$  から  $A[n-1]$  まで順に調べていき,  $A[j]$  のキーが軸のキー  $A[i].key$  未満のとき,  $A[j]$  を **選択** し軸  $A[i]$  と交換する
- 比較は必ず  $(n^2-n)/2$  回行われる
  - 時間計算量:  $O(n^2)$



# 単純選択法の第 0 フェーズ



# 単純選択法の第1フェーズ



# 単純選択法のプログラムと計算量

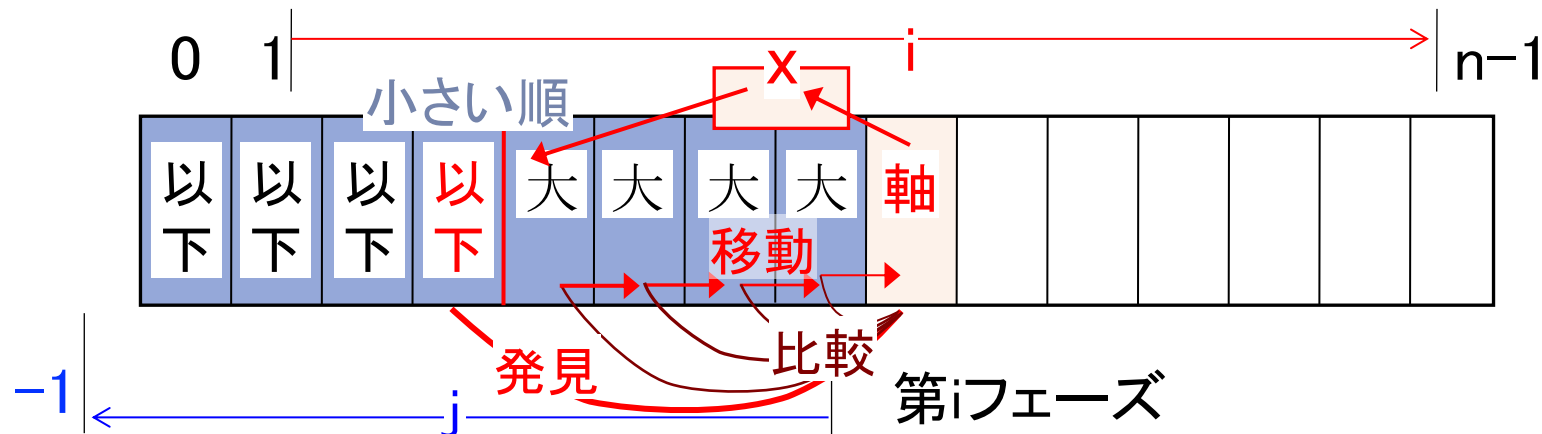
```
for(i=0; i<(n-1); i++){  
    for(j=i; j<=(n-1); j++){  
        if(A[i].key>A[j].key){  
            temp=A[j];  
            A[j]=A[i];  
            A[i]=temp}}}
```

A[i] を軸とする (iの**設定と比較: 2回**)  
n-i-1 回繰り返す (jの**設定と比較: 2回**)  
**比較1回**が必ず行われる  
**最悪3回の交換** (最良は0回)

- 時間計算量
  - 最悪の場合：軸A[i]とすべてのA[j]が交換
  - 最良の場合：軸A[j]とA[i]は交換されない (比較は必ず行われる)
  - $T_{max}(n) = \sum_{i=0}^{n-2} (2 + 6(n-i-1)) = 3n^2 - n - 2$
  - $T_{min}(n) = \sum_{i=0}^{n-2} (2 + 6(n-i-1)) = 3n^2 - n - 2$
  - どちらも  $O(n^2)$
- 領域計算量
  - n+3 (n,temp, i, j)
  - $O(n)$

# 単純挿入法

- 軸  $A[i]$  を設定し  $A[0]$  から  $A[n-1]$  まで以下を行う
- $A[0]$  から 軸  $A[i-1]$  までは **小さい順に並んでいる**
- $j$  を  $i-1$  から一つずつ減らしていき、軸  $A[i].key$  以下の  $A[j].key$  が見つかったら  **$A[j]$  の右に  $A[i]$  を挿入** する
- **$j$  が負 ( $-1$ )** になったら、 $A[0]$  から  $A[i-1]$  の  $key$  がすべて  $A[j].key$  より大きい場合なので  **$A[i]$  を  $A[0]$  先頭に挿入** する。





# 単純挿入法のプログラムと計算量

```
for(i=1; i<=n-1; i++){  
    x=A[i];  
    j=i-1;  
    while(j>=0 && x.key<A[j].key){  
        A[j+1]=A[j];  
        j=j-1;}  
    A[j+1]=x;}
```

iの**設定**, **比較** : 2回

軸の**設定** : 1回

**代入** : 1回

最悪i回のループ (比較2回, 脱出時1回)

代入

代入

- 軸の挿入

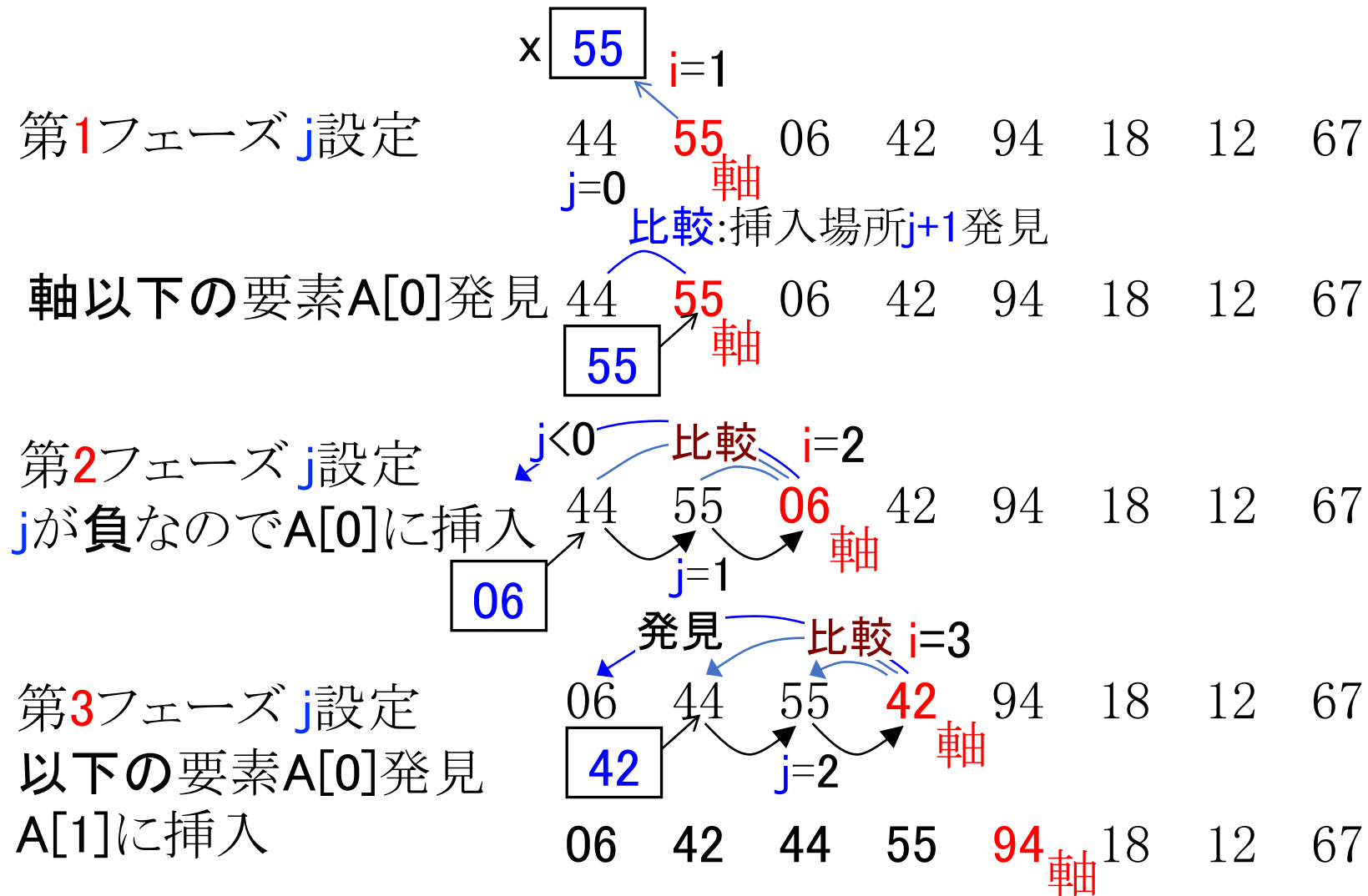
- while文

- A[0]からA[i-1]のkeyがA[i].keyより大きいときは, jが-1となるが, &&は左演算子が不成立のときは右演算子は評価しない (A[-1]は評価しない)
- 最悪 : A[0]までの比較・代入が**i回**, whileループの脱出の判定まで含め4i+1回
- 最良 : whileループはA[i-1]との**比較**ですぐ止まるので脱出判定の**1回**のみ

- 時間計算量

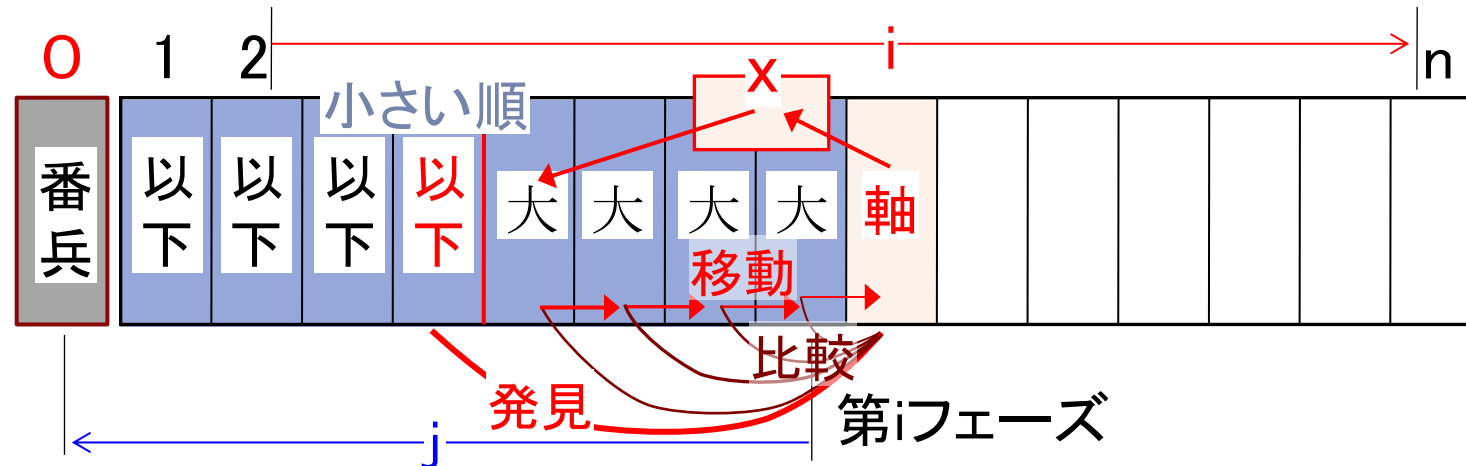
- $T_{max}(n) = \sum_{i=0}^{n-1} (2 + 1 + 1 + 4i + 1 + 1) = 2n^2 + 2n - 4 = O(n^2)$
- $T_{min}(n) = \sum_{i=0}^{n-1} (2 + 1 + 1 + 0 + 1 + 1) = 6(n - 1) = O(n)$

# 単純挿入法の各フェーズ



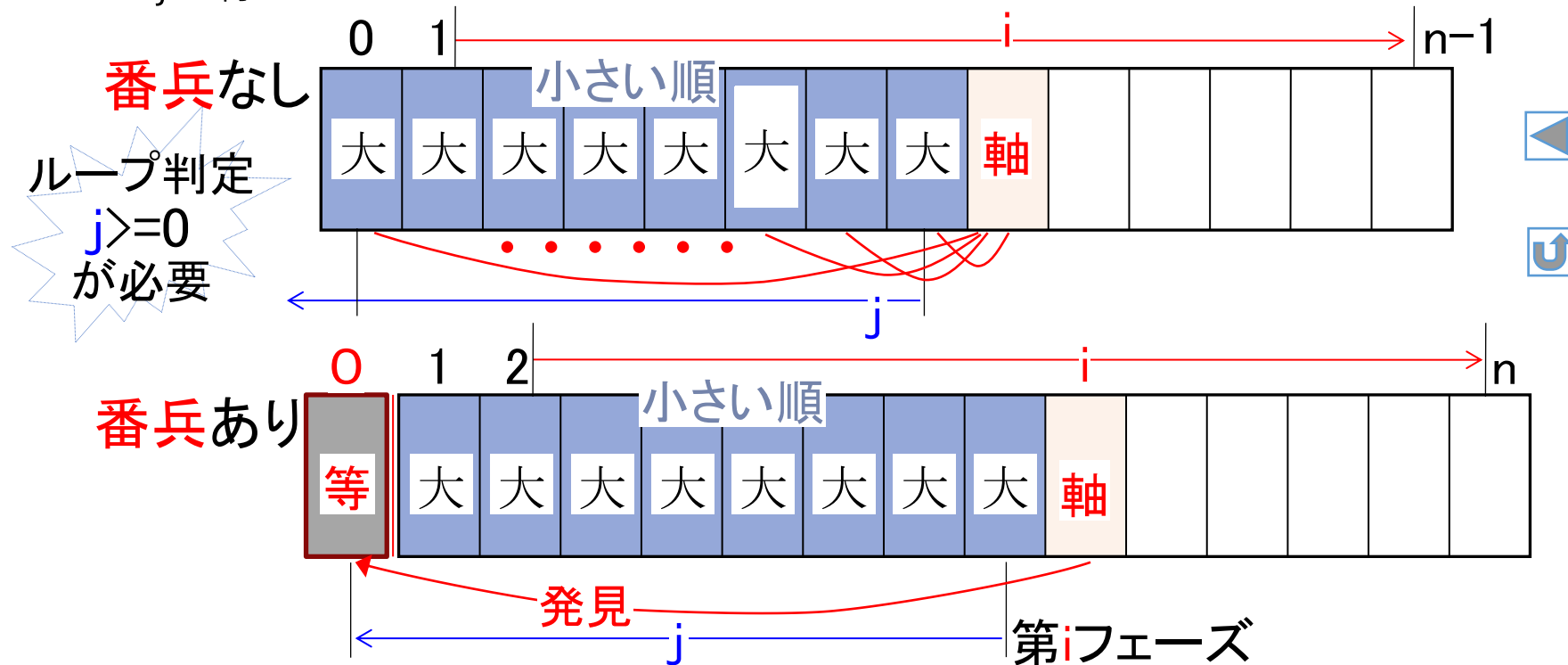
# 番兵を用いた単純挿入法

- 配列Aを $n+1$ 個の要素
  - $A[1]$ から $A[n]$ がソート対象,  $A[0]$ は「番兵」としてindex  $j$ の監視に利用する
- 番兵 $A[0]$ に 軸 $[i]$ を設定
- $A[1]$ から $A[i-1]$ は小さい順に並んでいる
- $j$ を $i-1$ から一つずつ減らしていき, 軸 $A[i].key$ 以下の $A[j].key$ が見つかったら $A[j]$ の右に $A[i]$ を挿入する
- 番兵がいるため $j$ は0より小さくならない ( $j$ の判定がいらない)

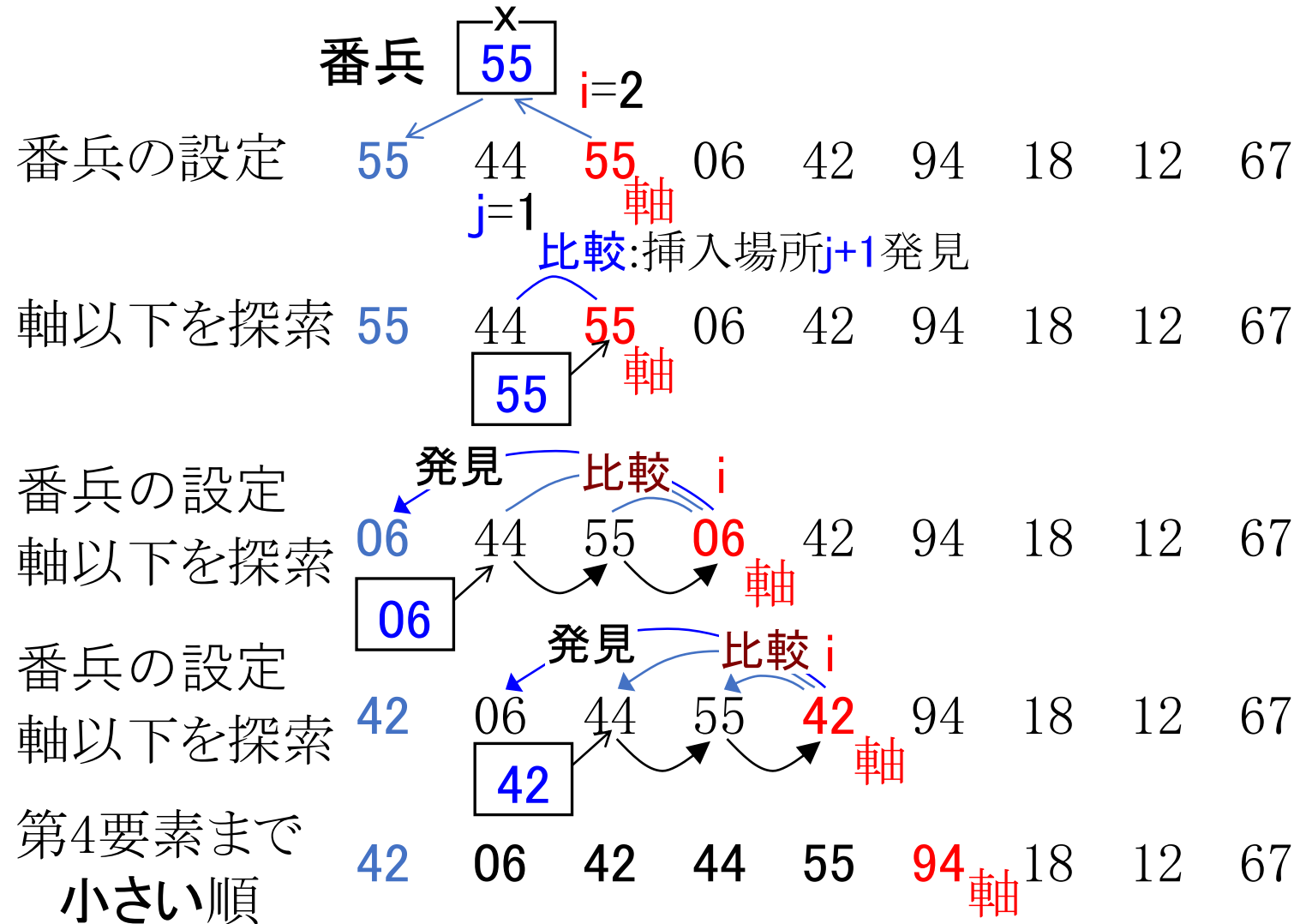


# 番兵の働き

- index  $j$ はどこで止まるか？
  - $j$ は軸のキー以下の要素が見つかったとき停止する
  - $i$ 未満のすべての $A[j]$ が軸のキーより大きいときには
    - $j$ を減らす時にwhileループの条件判定 $j \geq 0$ が必要だが、番兵を使うと最悪番兵で $j$ を停止できる



# 番兵の単純挿入法の各フェーズ



# 番兵あり単純挿入法のプログラムと計算量

for(i=2;i<=n; i++){	iの設定, 比較: 2回
x = A[i];	軸の設定
A[0] = x;	番兵の設定
j= i - 1;	代入
while(x.key<A[j].key){	最悪i回
A[j+1] = A[j];	代入
j=j+1;}	代入
A[j+1]=x;}	軸の挿入

- whileループ

- 最悪: 番兵までの比較がi回より, whileループによるすべての代入は  $2(i-1)$  回
- 最良: whileループは  $j= i-1$  ですぐとまる. 比較1回

- 時間計算量

- $$T_{max}(n) = \sum_{i=2}^n (2 + 3 + i + 2(i-1) + 1) = \frac{3}{2}n^2 + \frac{11}{2}n - 7 = O(n^2)$$
$$T_{min}(n) = \sum_{i=2}^n (2 + 3 + 1 + 1) = 7(n-1) = O(n)$$

# まとめ：二つの単純挿入法の計算量

- 時間計算量
  - 最悪時間計算量： $O(n^2)$
  - 最良時間計算量： $O(n)$
- 領域計算量
  - $O(n)$
- すべての単純法は
  - 時間計算量がよくない： $O(n^2)$
  - 領域計算量は良好： $O(n)$  => 入力配列が小さい場合は有効

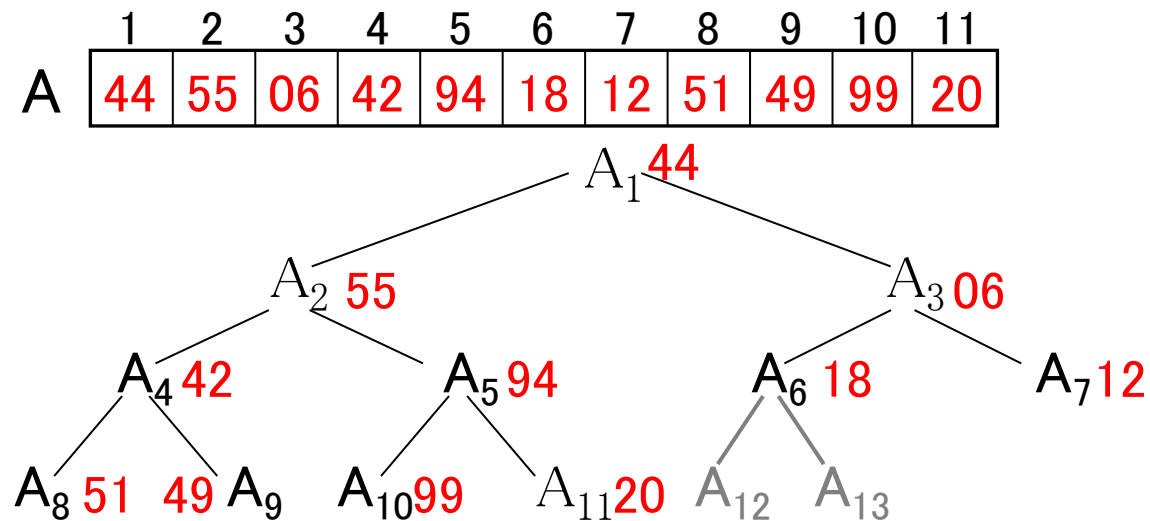
# ヒープソート

ヒープ構造を利用したソート



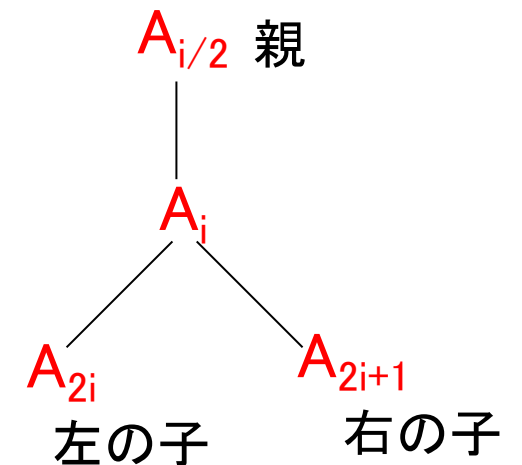
# 完全2分木とは

- 2分木：全長点の子数が最大2個の根付き木
- 完全2分木：葉以外の頂点の子がちょうど2個ですべての葉の高さが等しい2分木
- 配列  $\Rightarrow$  木構造（完全2分木）
- 配列  $A$  は  $A[1]$  から  $A[n]$ （ $A[0]$  は使わない）



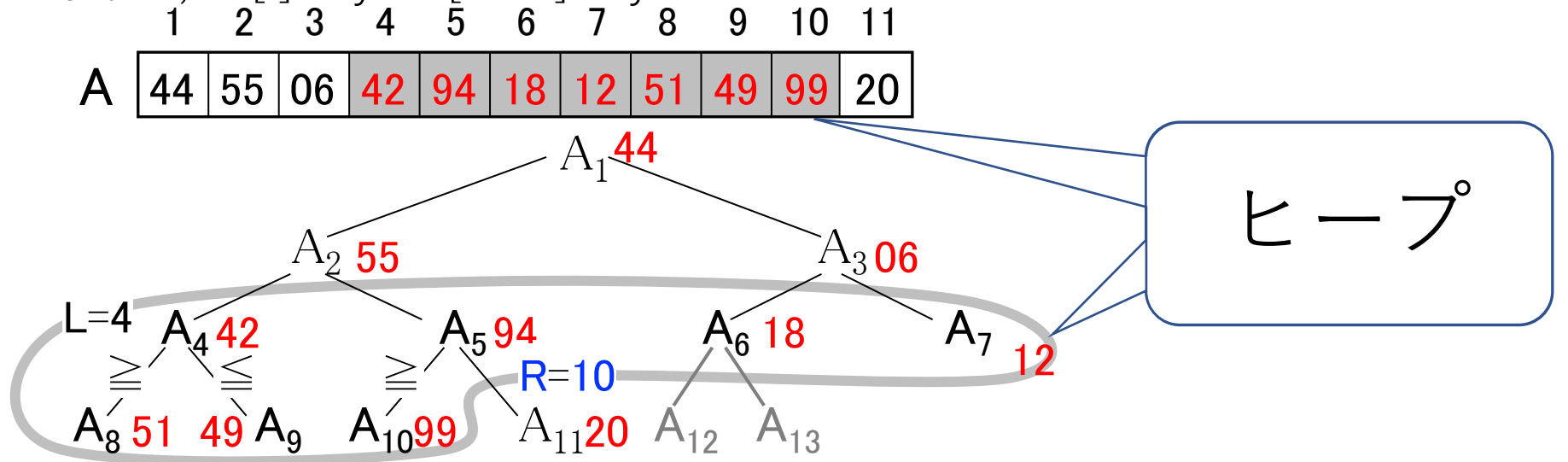
頂点  $i$  において

- 親： $i/2$
- 左の子： $2i$ ,
- 右の子： $2i+1$



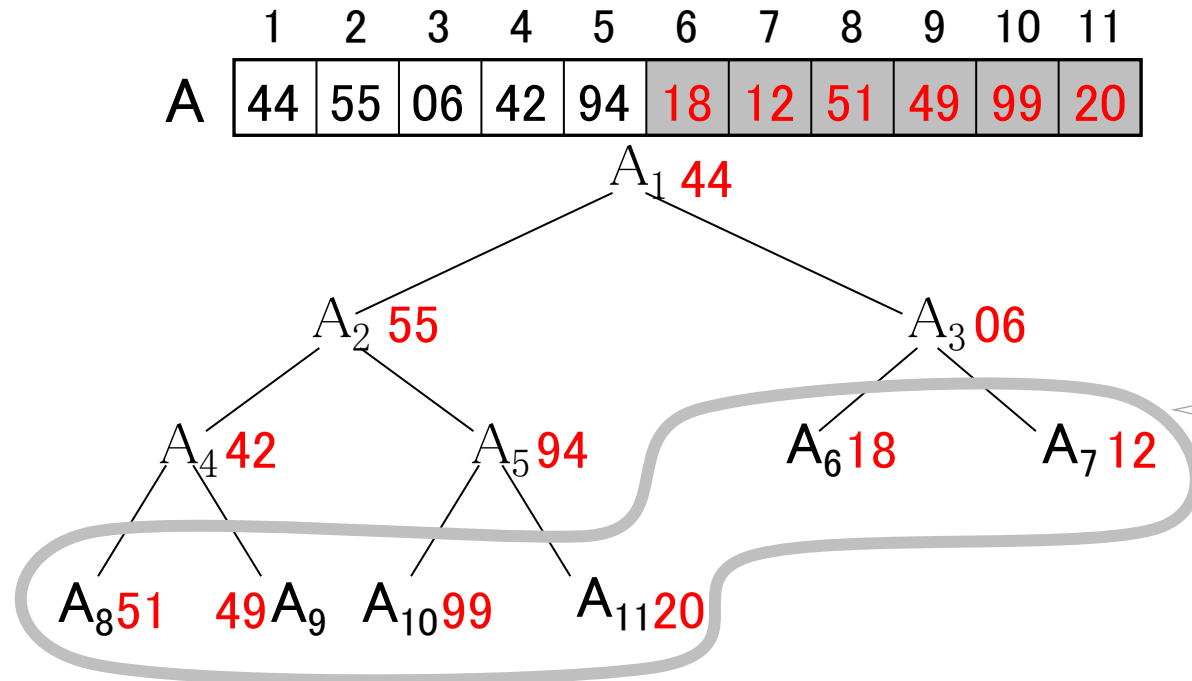
# ヒープとは

- 完全2分木の各頂点はデータをひとつずつ持ち、必ず「ヒープ条件」を満たす
- ヒープ条件（最小ヒープ）
  - ある頂点のデータの値は、その親の持つデータの値以上である
  - 配列A（index 1からn）の部分列A[L], A[L+1], ..., A[R]において
  - すべての  $i=L, \dots, R$  に対して以下を満たす
    - $2i \leq R$  ならば,  $A[i].key \leq A[2i].key$
    - $2i+1 \leq R$  ならば,  $A[i].key \leq A[2i+1].key$



# ヒープ化：最初からヒープ

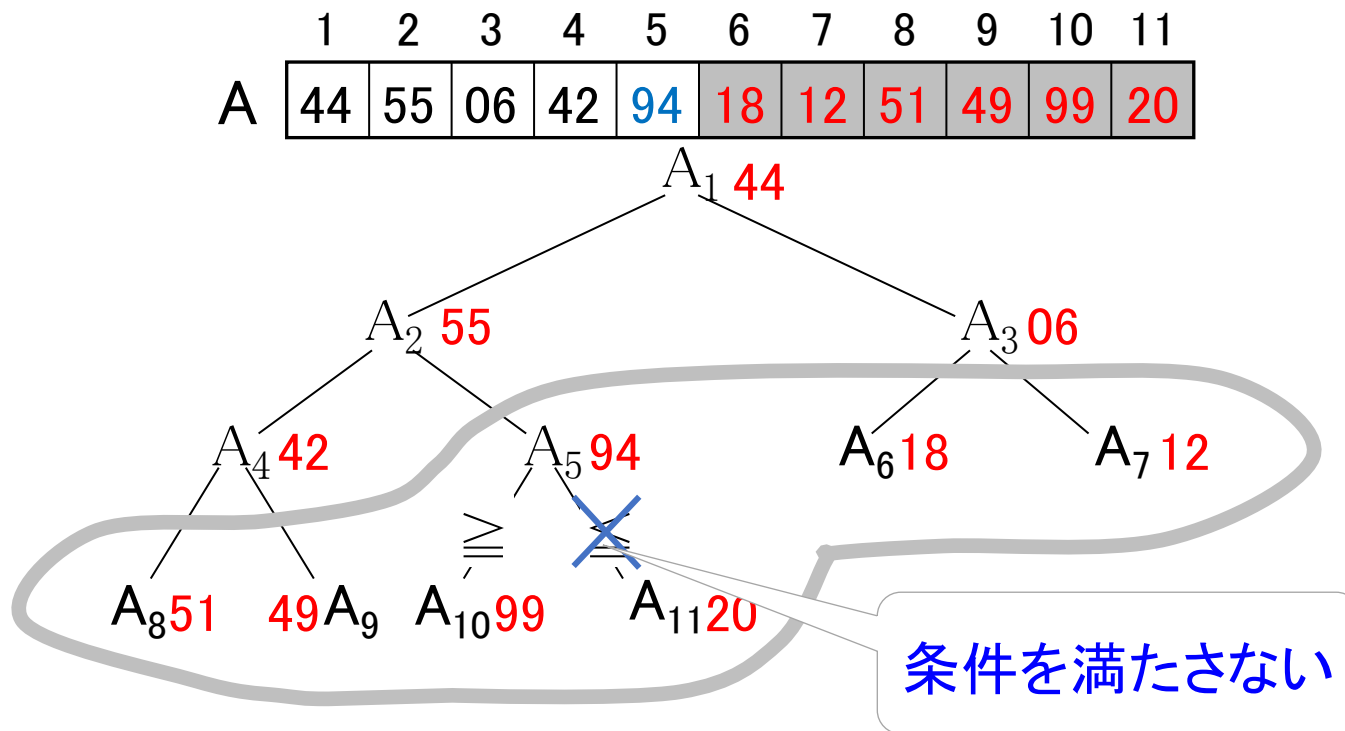
- 部分列 $A[6], \dots, A[11]$ を考える
  - $i=L(6), \dots, R(11)$  のヒープ条件
    - $2i \leq R$  ならば,  $A[i].key \leq A[2i].key$
    - $2i+1 \leq R$  ならば,  $A[i].key \leq A[2i+1].key$
  - $2i \leq R$ ,  $2i+1 \leq R$  が成立しないので, 後半条件は無視しヒープ条件が成立する



配列の**右半分**  
 $A[6] \cdots A[11]$ は  
必ず**ヒープ**

葉はヒープ

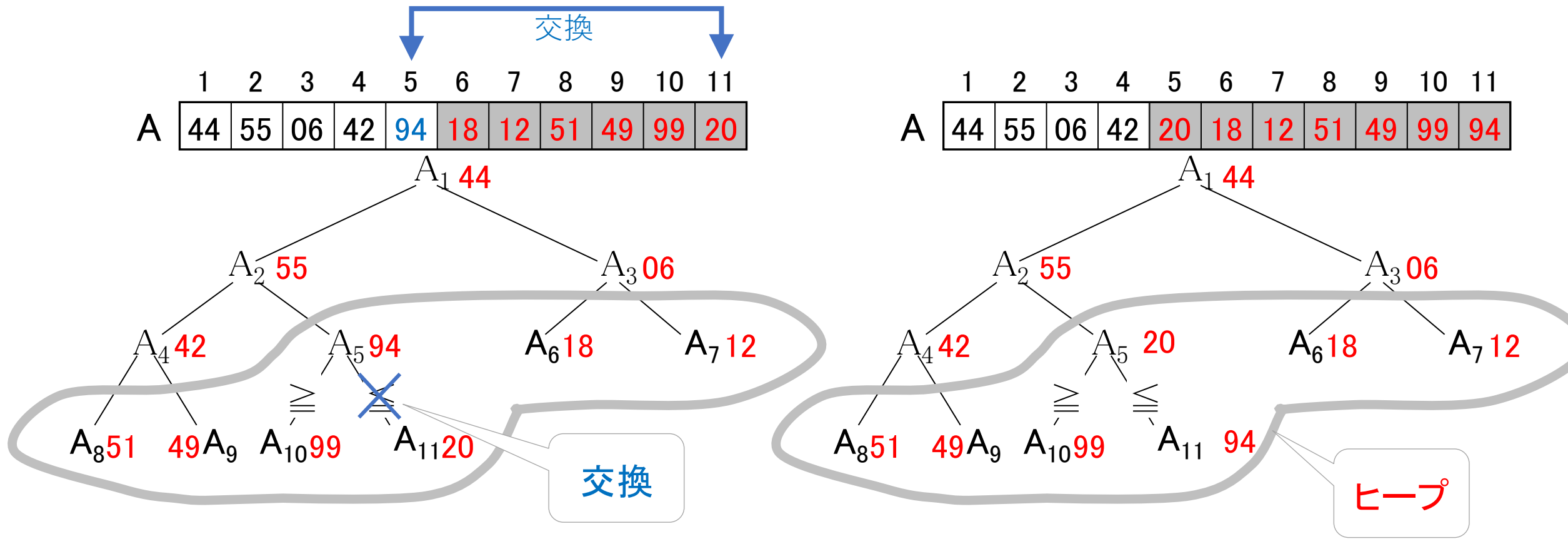
# ヒープ化：ヒープの配列を左に拡大



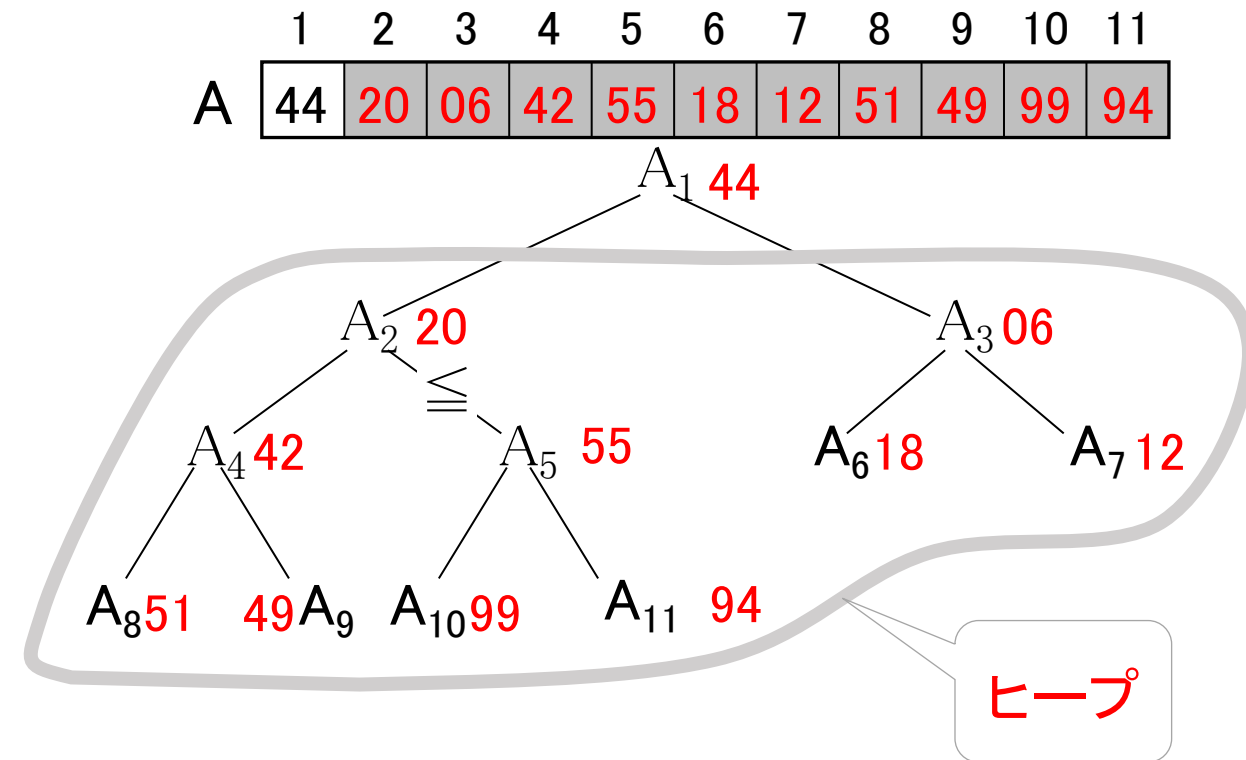
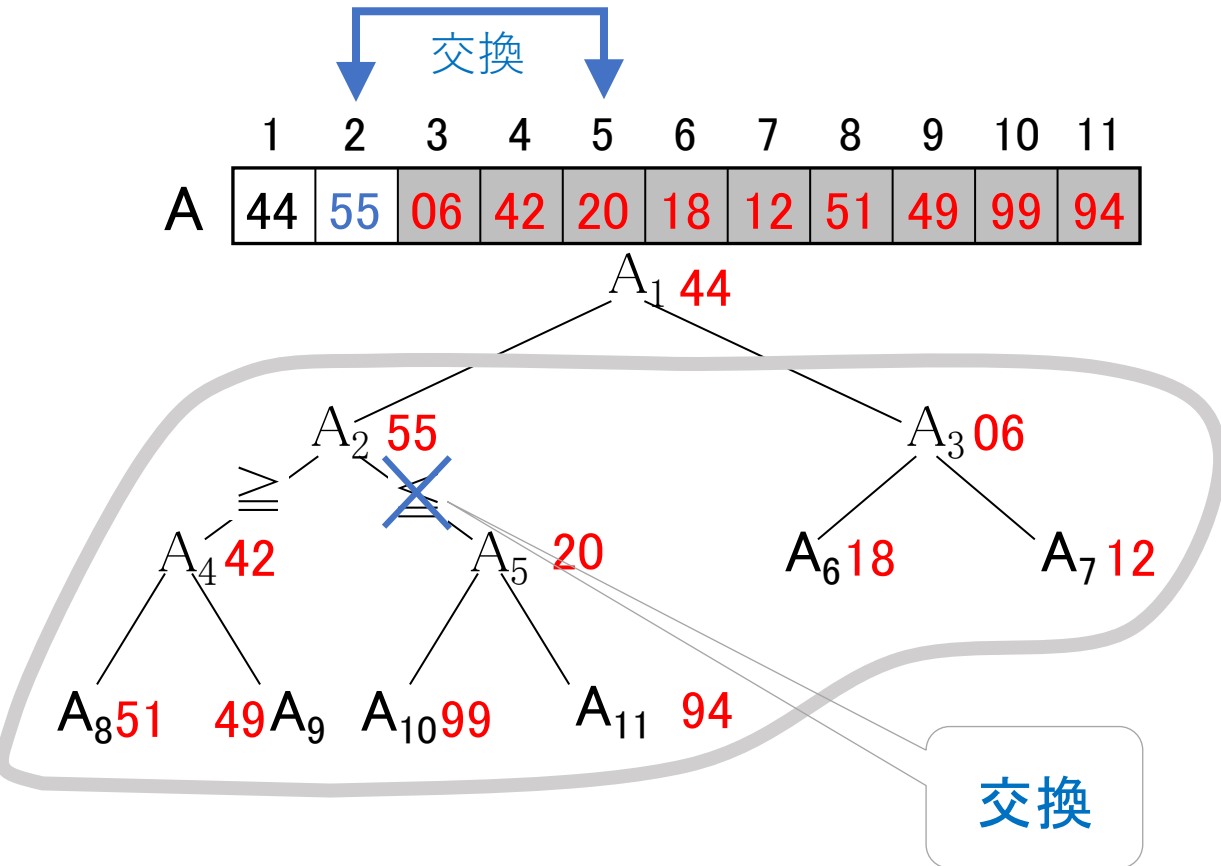
# ヒープ化：ふるい落とし

- ヒープ $A[L+1], \dots, A[R]$ が与えられているとき
- $A[L]$ を追加し,  $A[L], A[L+1], \dots, A[R]$ をヒープになるようにする
- この操作を  $A[L]$ のふるい落とし（シフト）と呼ぶ
  1.  $i=L$
  2. 子の数で場合分け
    - 無し:  $2*i > R$ ならば, おわり
    - 左の子だけ:  $2*i = R$  ならば,  $j=2*i$
    - 左右両方:  $2*i+1 \leq R$ ならば, キーの小さい方のインデックスを $j$ とする
  3.  $A[i].key$ と $A[j].key$ を比較
    - $A[i].key \leq A[j].key$  ならば, おわり
    - そうでなければ,  $A[i]$ と $A[j]$ を交換して, 処理2へ

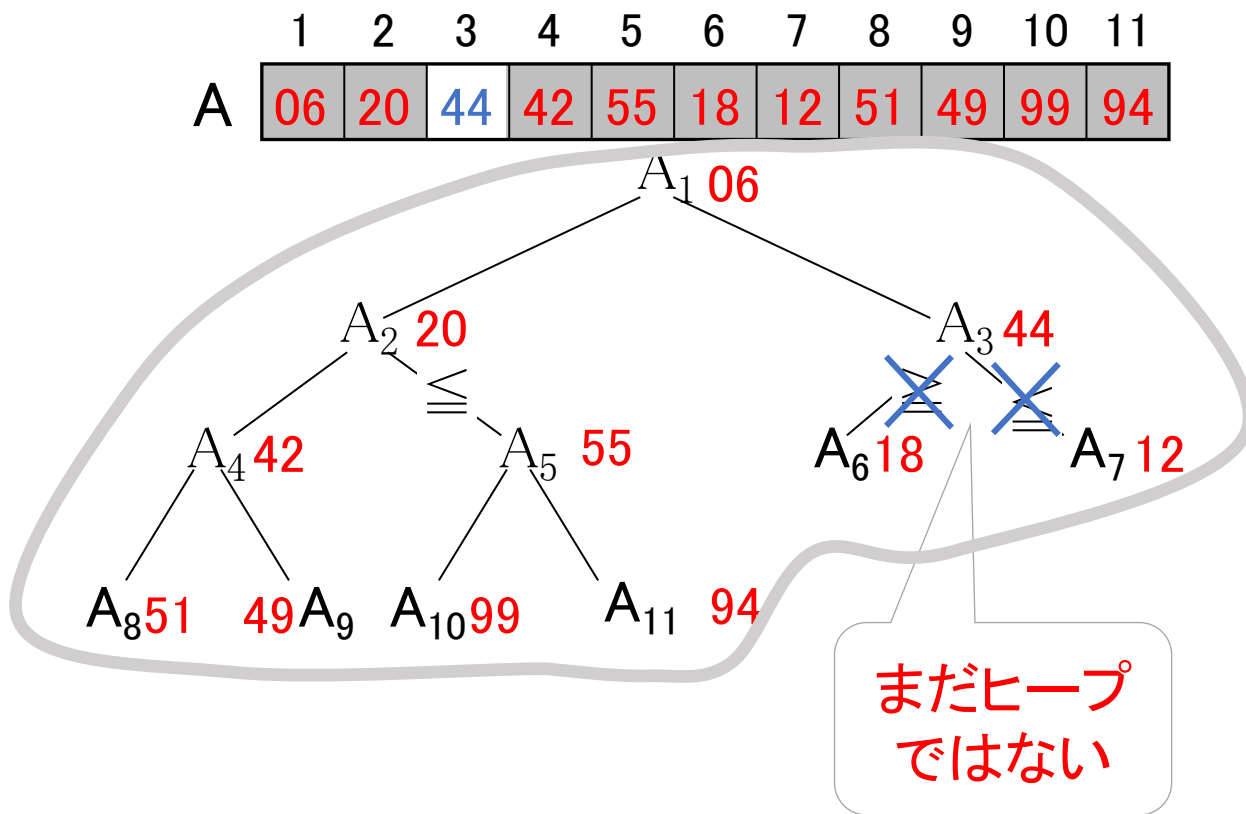
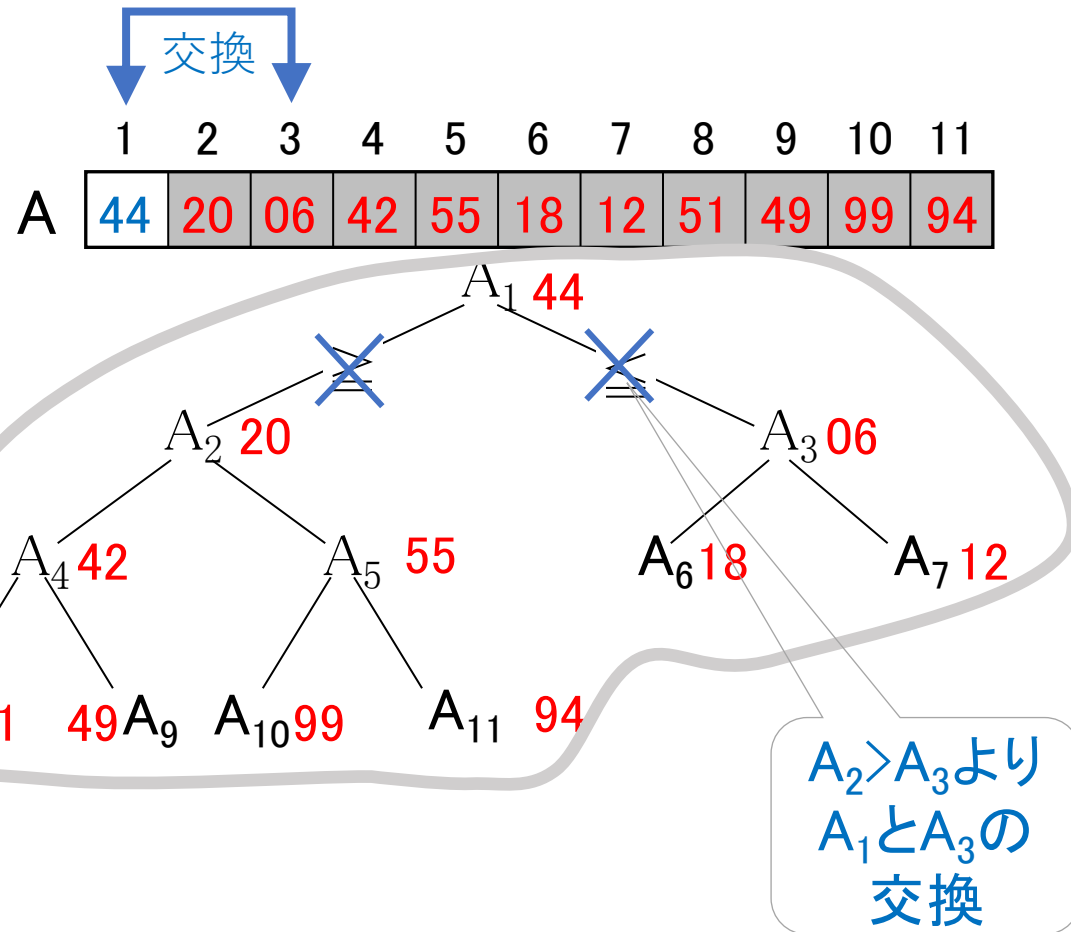
# ヒープ化：ふるい落とし



# ヒープ化：ふるい落とし

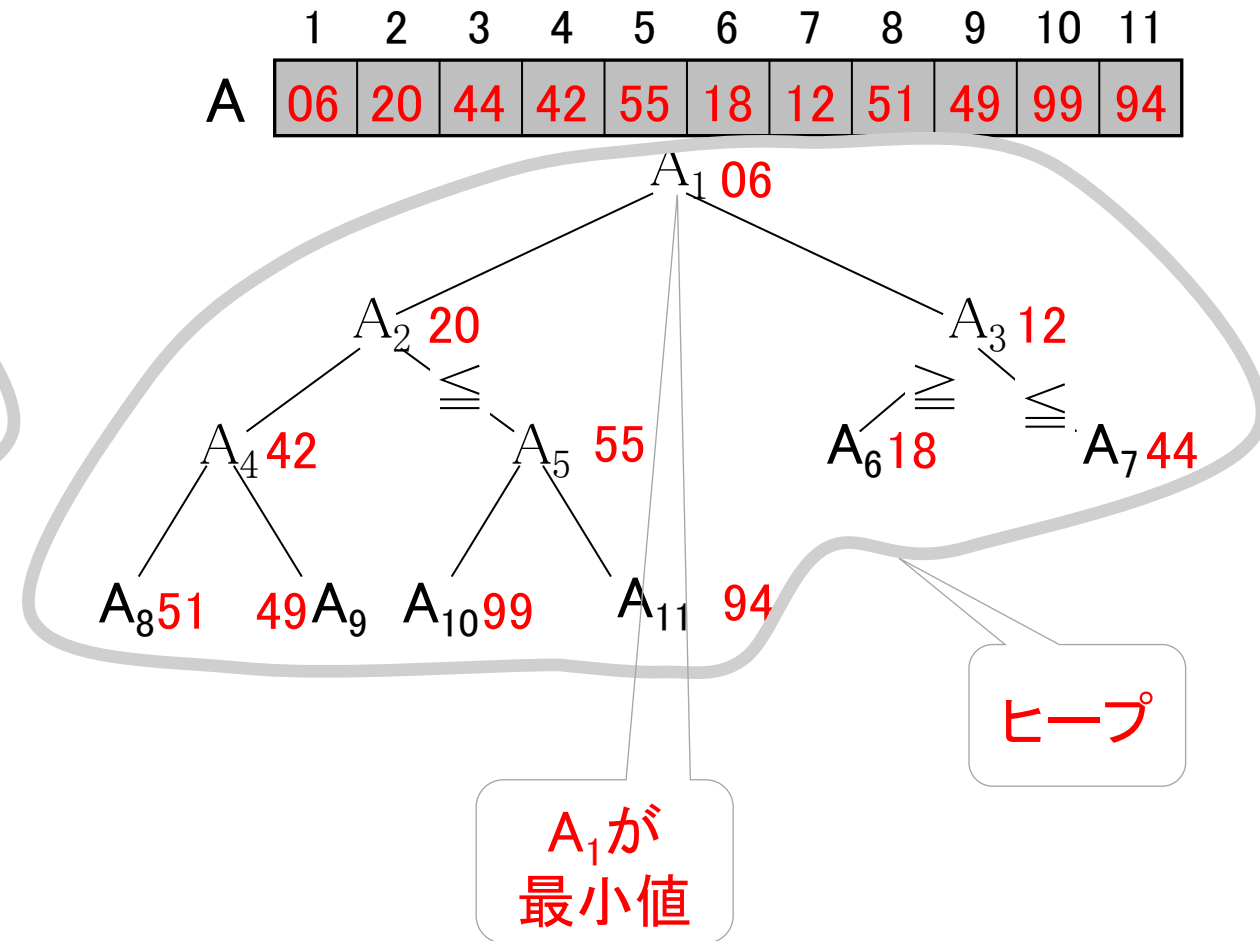
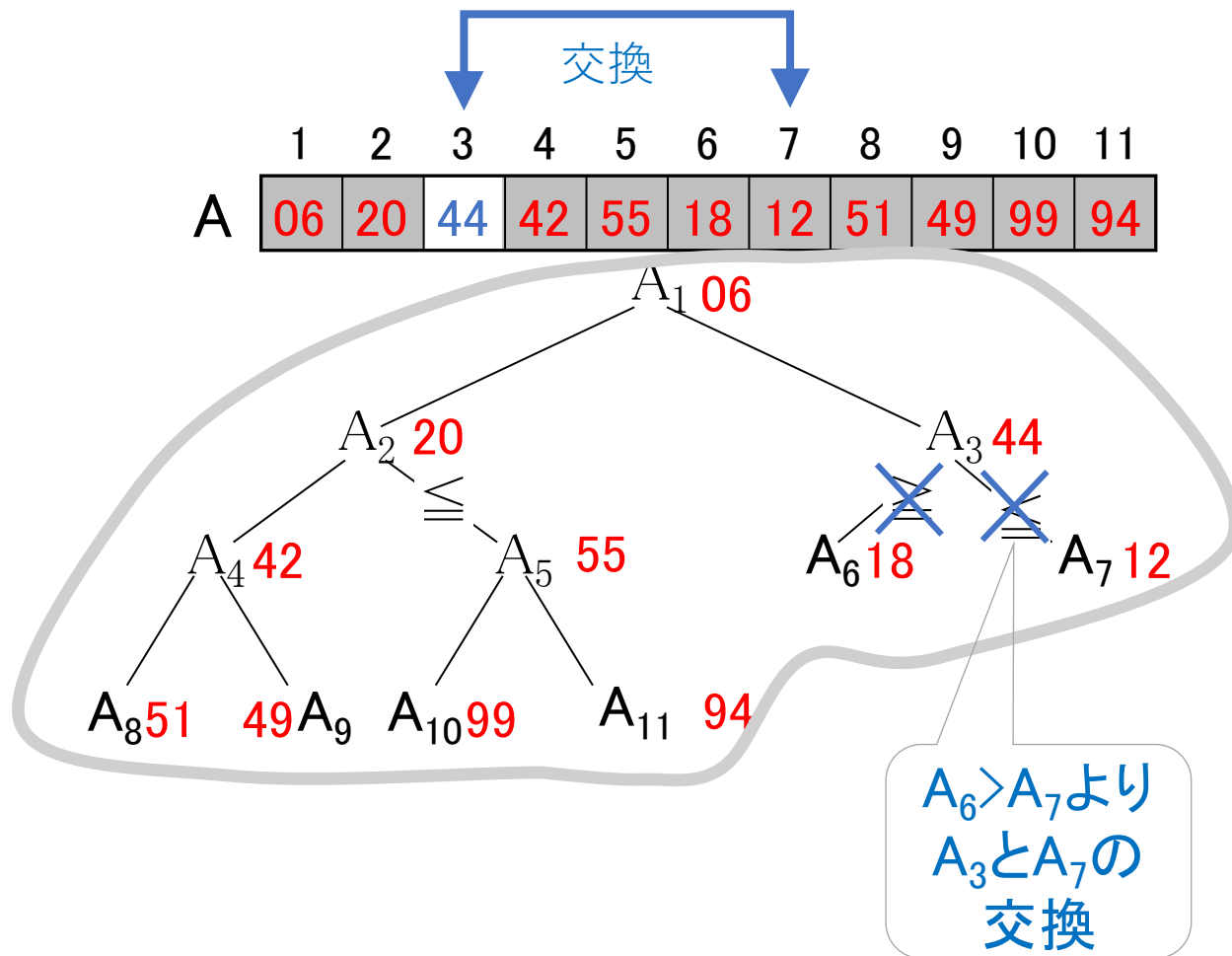


# ヒープ化：ふるい落とし





# ヒープ化：ふるい落とし

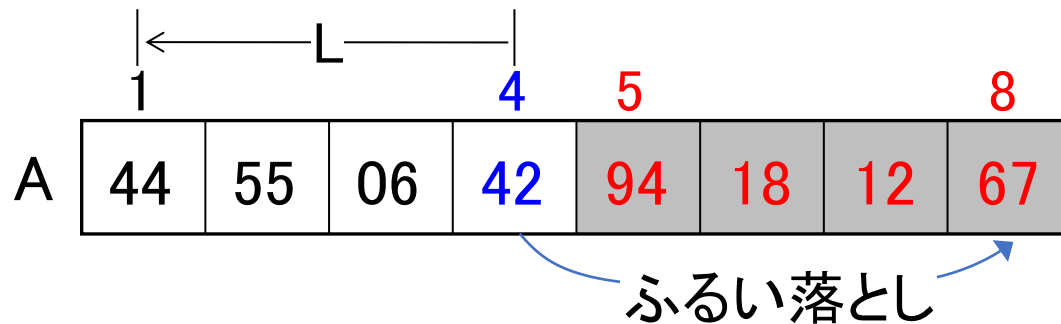


# ヒープソート

- ヒープソートの原理
  1. ヒープの構成
    - 配列 $A[1], \dots, A[n]$  をヒープにする
  2. ソート列の構成
    - ヒープ  $A[1], \dots, A[n]$  をソート列にする

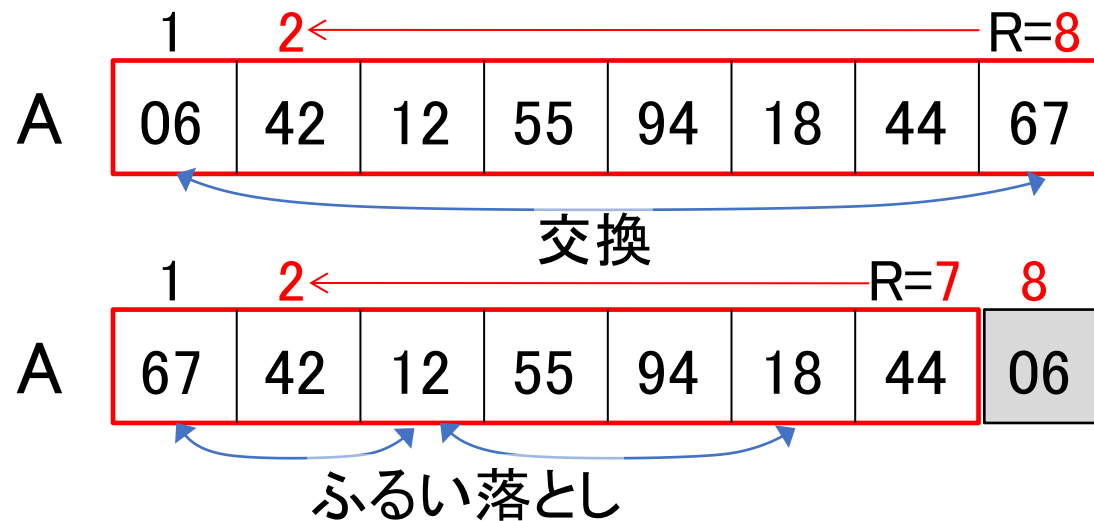
# 1. ヒープの構成

- 配列  $A[1], \dots, A[n]$  をヒープにする
  - $A[n/2 + 1], \dots, A[n]$  は子がないので **始めからヒープ**
  - $A[L], A[L+1], \dots, A[n]$  をヒープにするために  **$A[L]$  のふり落とし** を行う (L は  $n/2$  からデクリメントして 1 まで)



## 2. ソート列の構成

- ヒープ  $A[1], \dots, A[R]$  が与えられたとき,  $A[1]$  と  $A[R]$  を交換する.
- 交換直後は  $A[R]$  が  $A$  の要素で最も小さいキーを持つ
- $A[2], \dots, A[R-1]$  はヒープであるので,  $A[1]$  のみふるい落としを行い  $A[1], \dots, A[R-1]$  をヒープにする



# C言語によるヒープソート

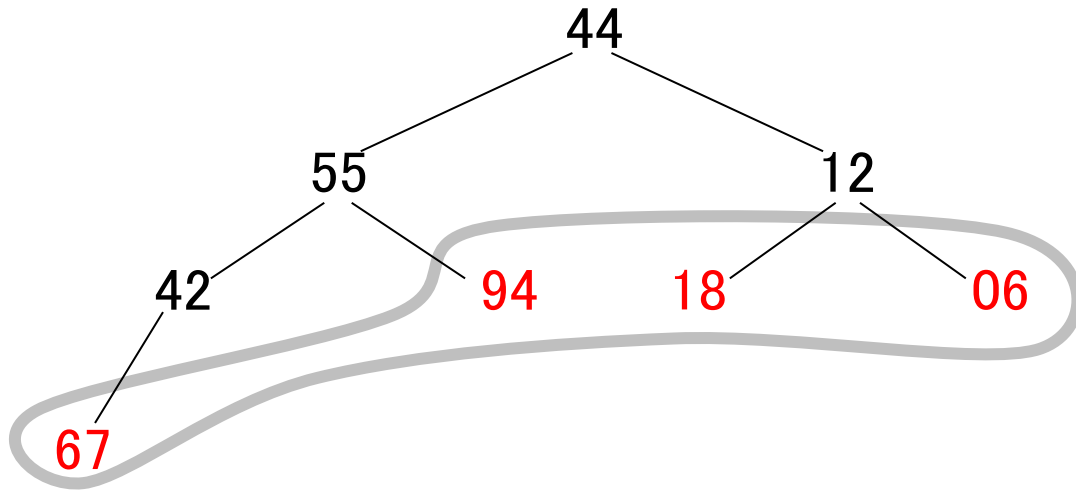
```
/* shift down A[L] into A[L+1],..., A[R] */
```

```
void shift(int L, int R){  
    int i, j; item x;  
    i=L; j=2*i; x=A[i];  
    while(i<=R){  
        if(j<R){  
            if(A[j].key>A[j+1].key){  
                j=j+1;}}  
        if(x.key<=A[j].key){  
            break;  
        }  
        A[i]=A[j]; i=j; j=2*i;  
    }  
    A[i]=x;  
}
```

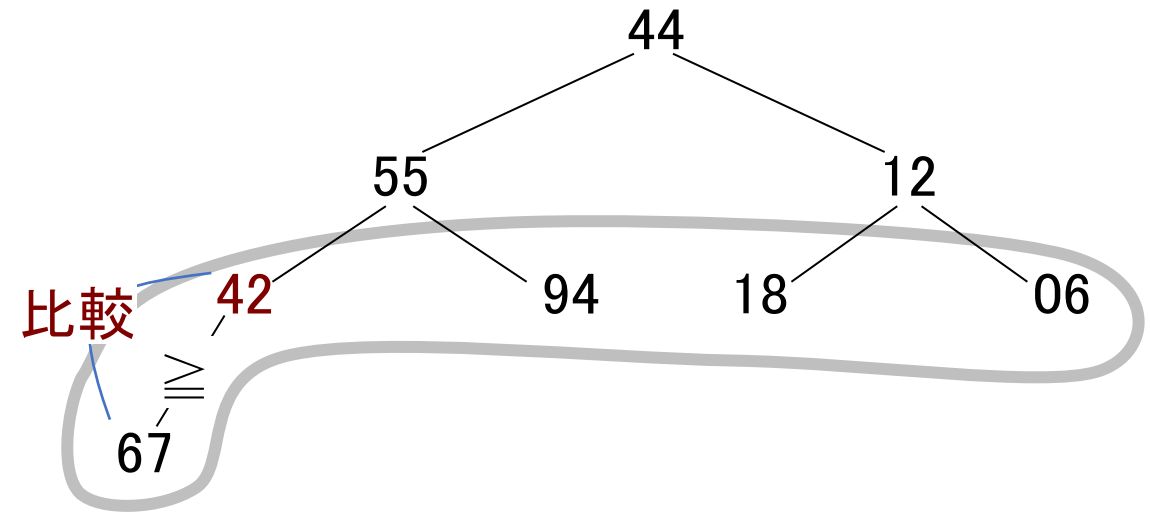
```
void heapsort(){  
    int L, R; item x;  
    L=(n/2)+1;  
    R=n;  
    while(L>1){  
        L=L-1;  
        shift(L, n);} /* A is heaped */  
    while( R>1){  
        x=A[1];  
        A[1]=A[R];  
        A[R]=x;  
        R=R-1;  
        shift(1, R);}  
}
```

動作例：44, 55, 12, 42, 94, 18, 06, 67  
ヒープを構成

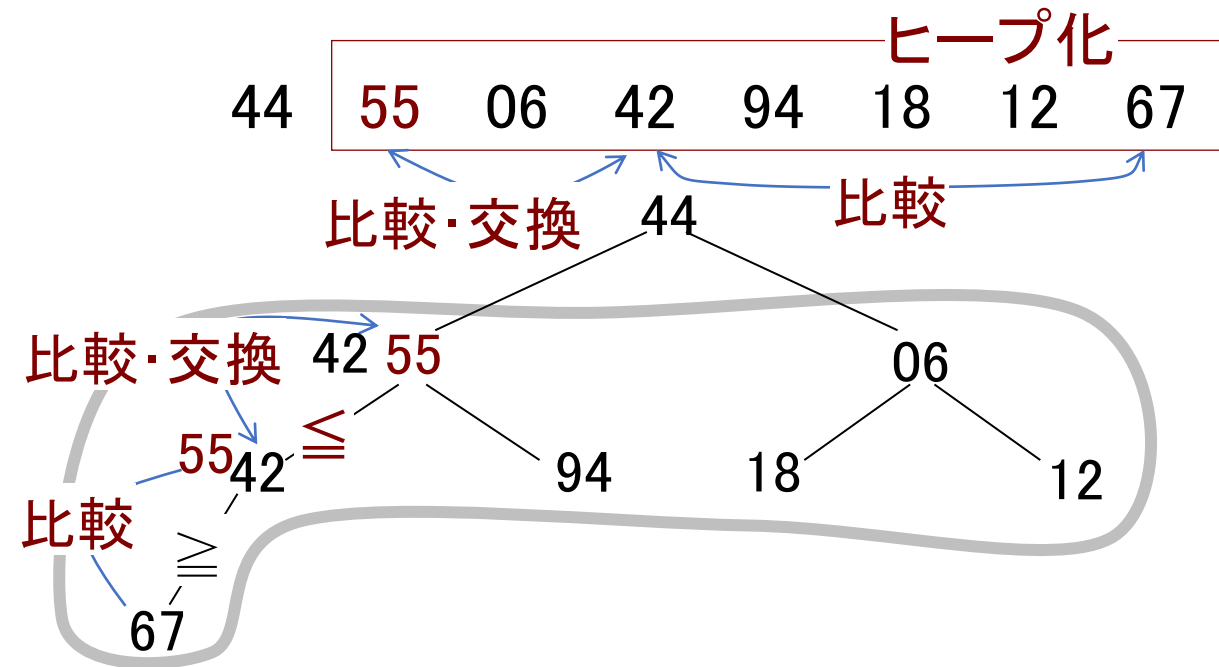
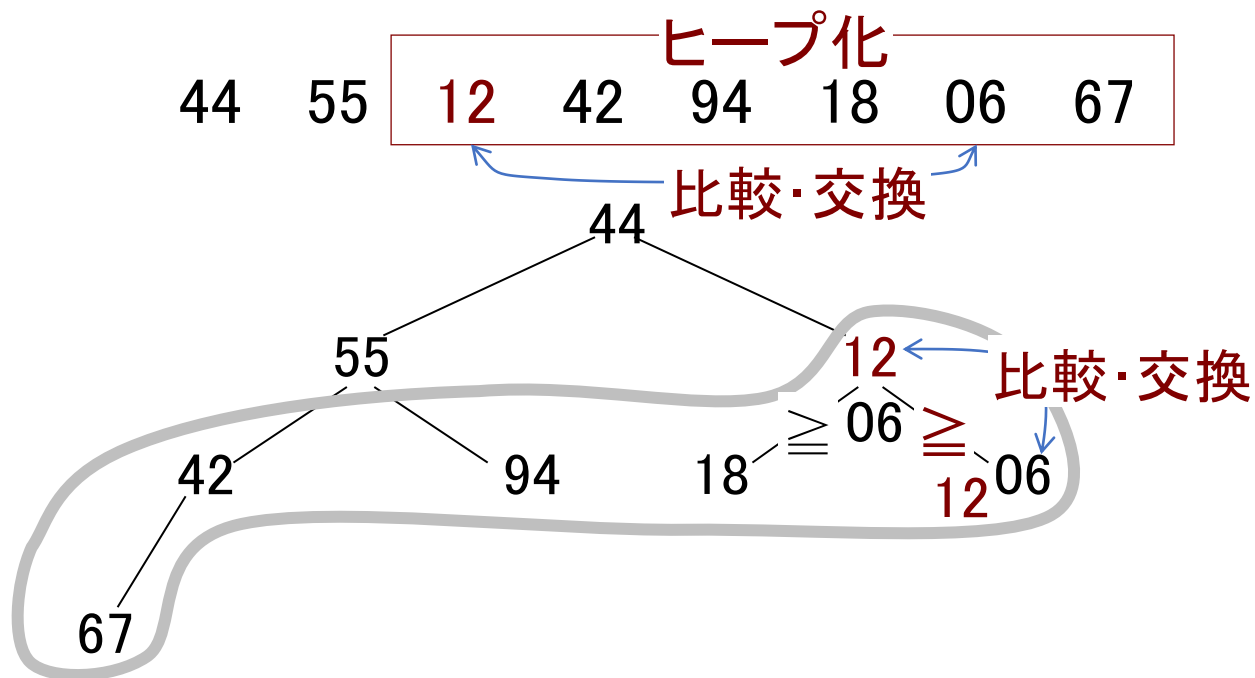
44 55 12 42 94 18 06 67 ヒープ



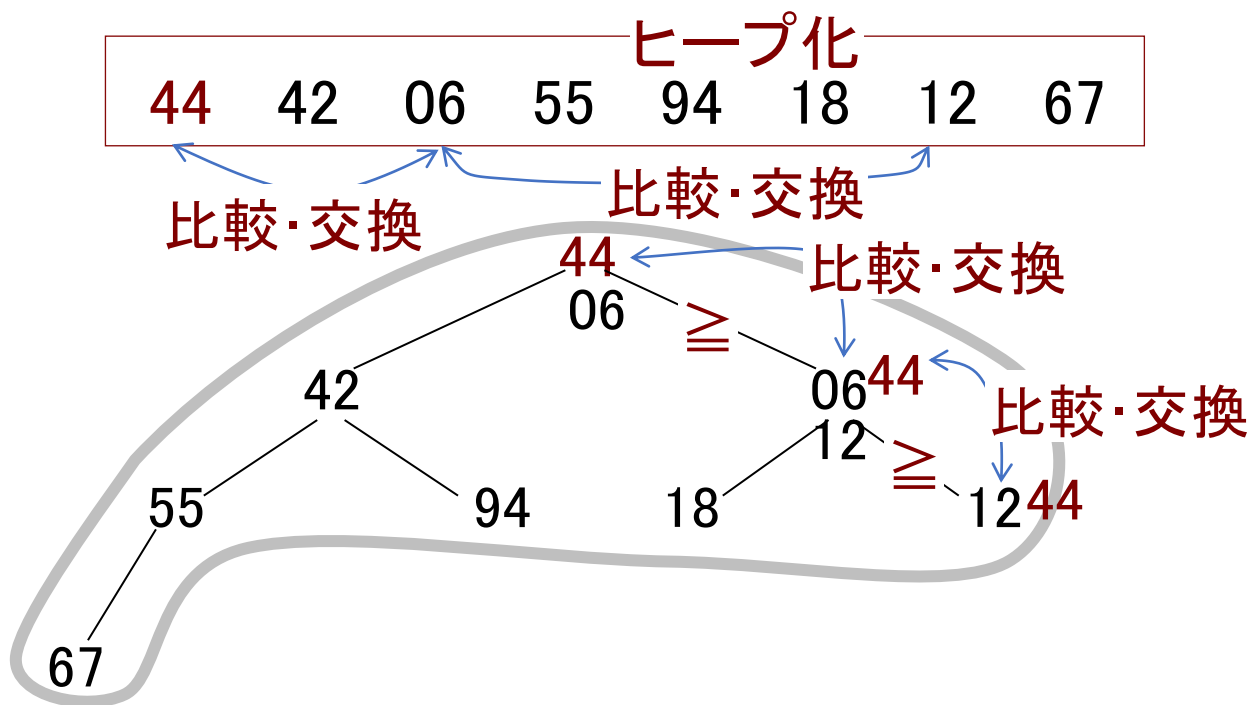
44 55 12 42 94 18 06 67 ヒープ化  
比較



動作例：44, 55, 12, 42, 94, 18, 06, 67  
 ヒープを構成

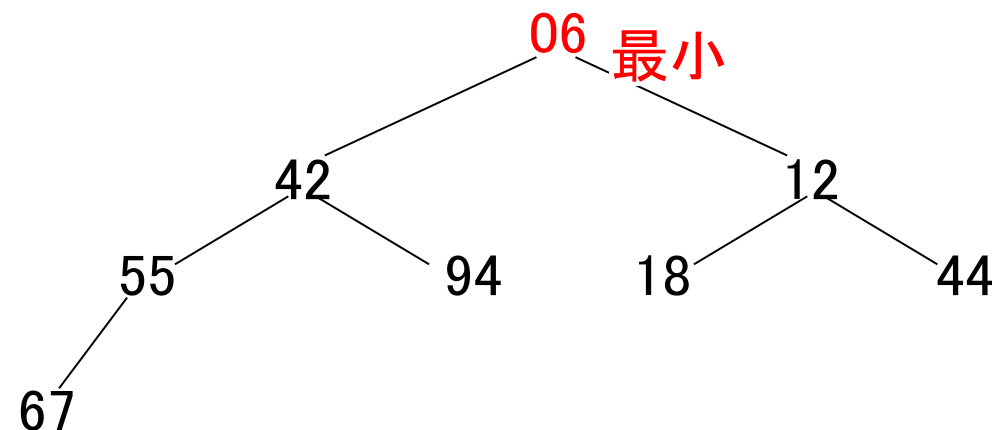


動作例：44, 55, 12, 42, 94, 18, 06, 67  
ヒープを構成



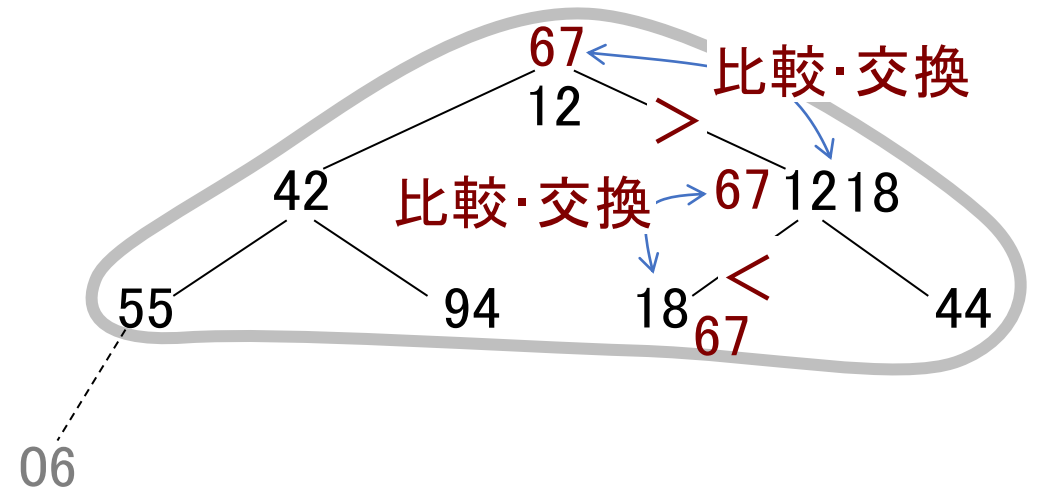
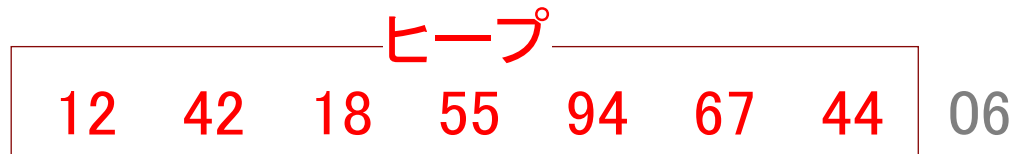
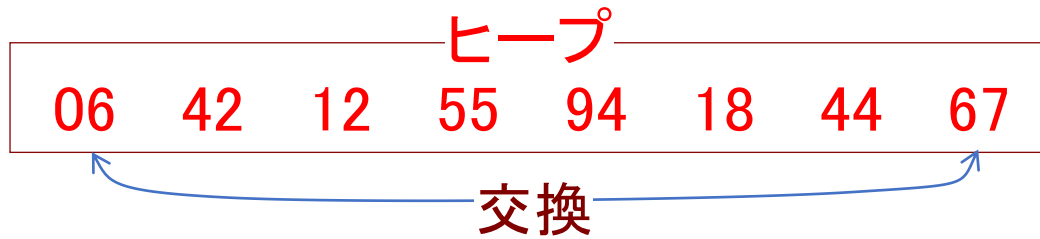
ヒーブ

06 42 12 55 94 18 44 67

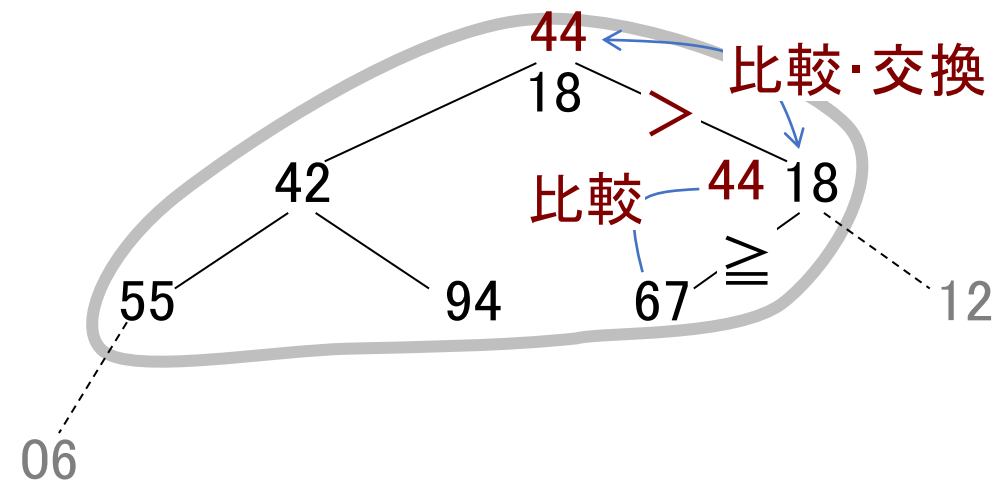
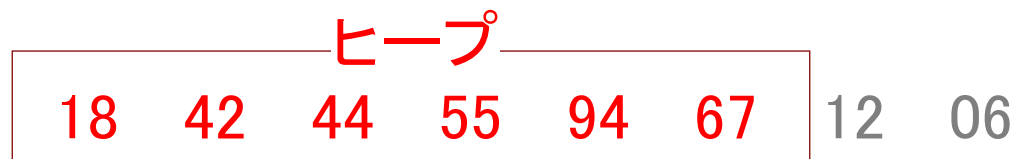
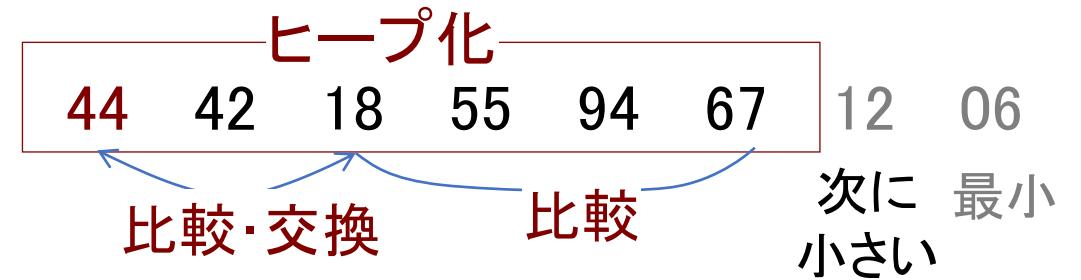
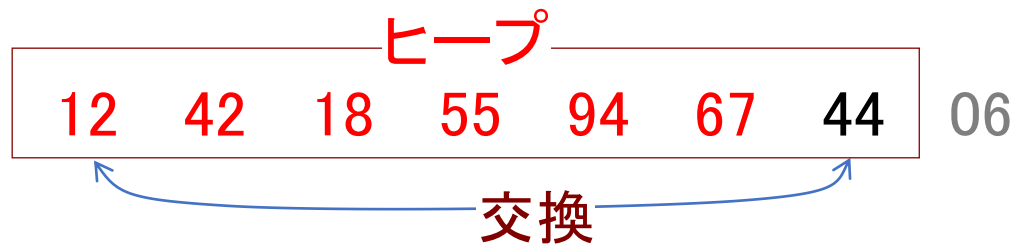




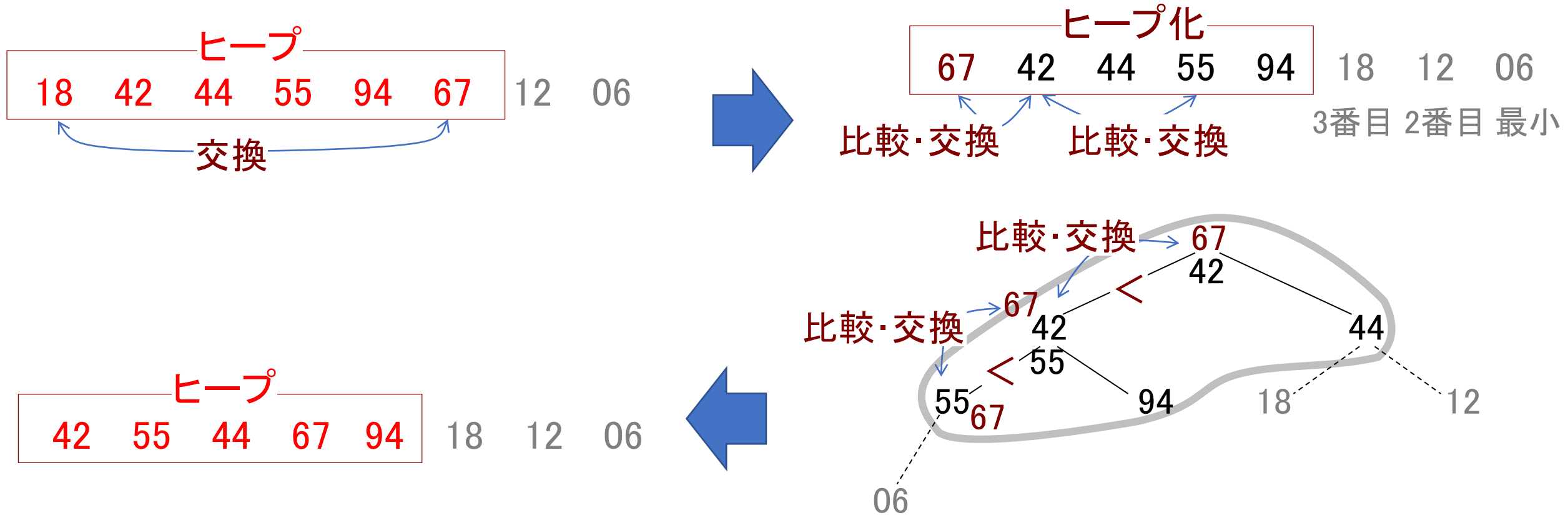
動作例：44, 55, 12, 42, 94, 18, 06, 67  
ソート列を構成



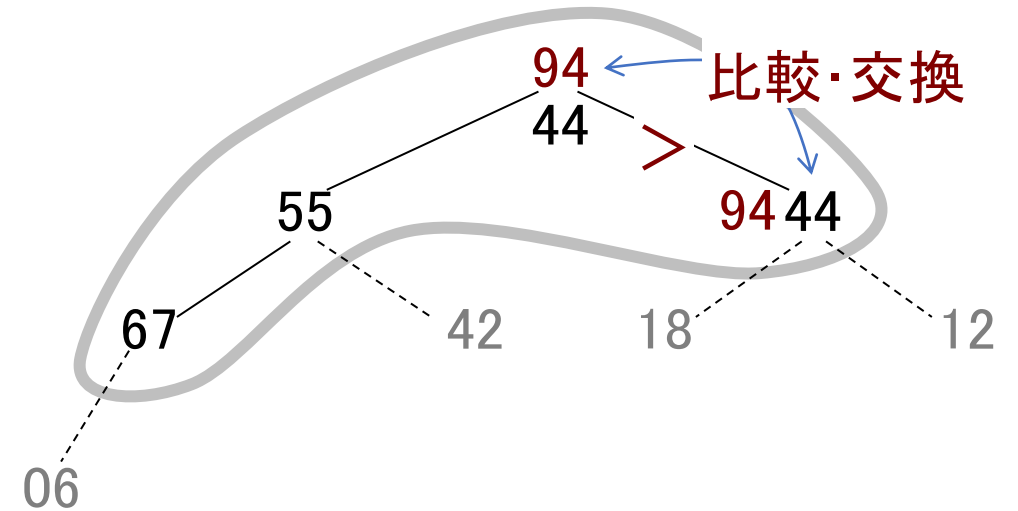
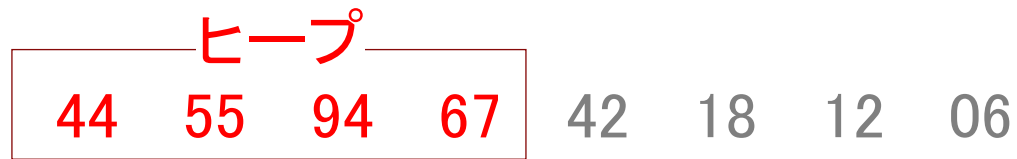
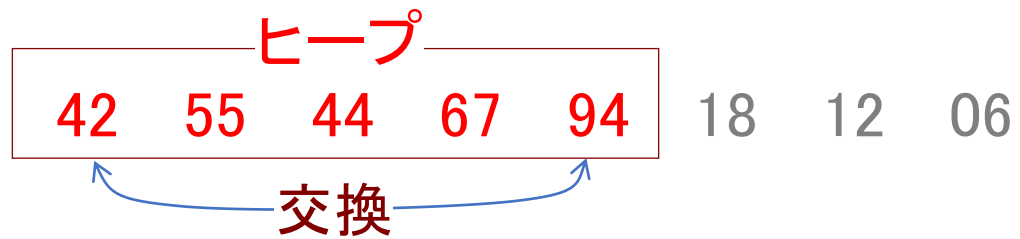
動作例：44, 55, 12, 42, 94, 18, 06, 67  
ソート列を構成



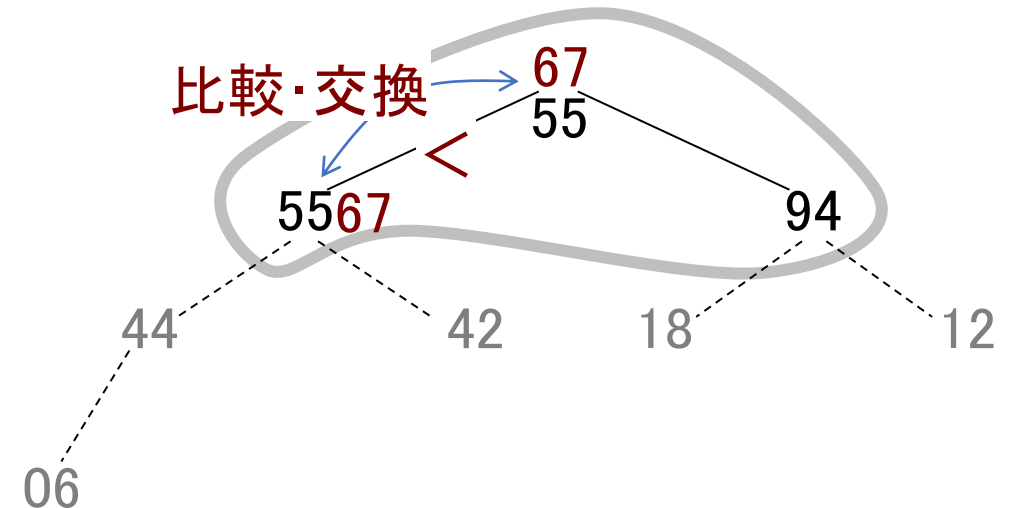
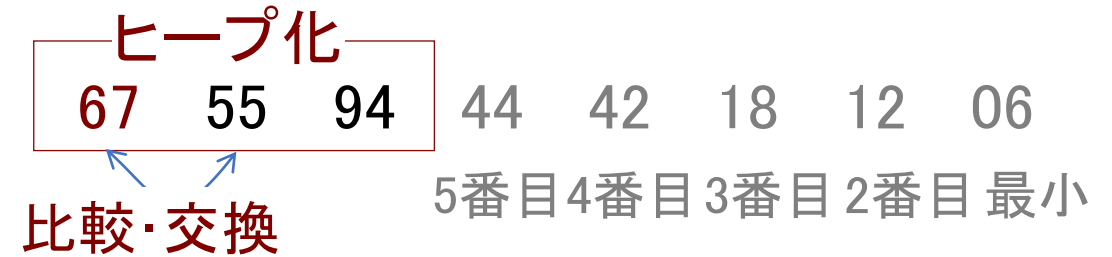
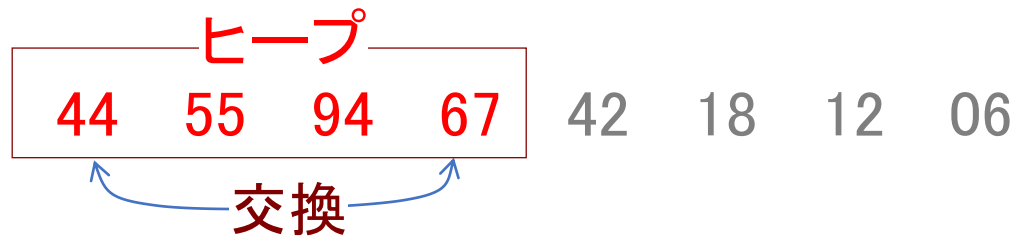
動作例：44, 55, 12, 42, 94, 18, 06, 67  
ソート列を構成



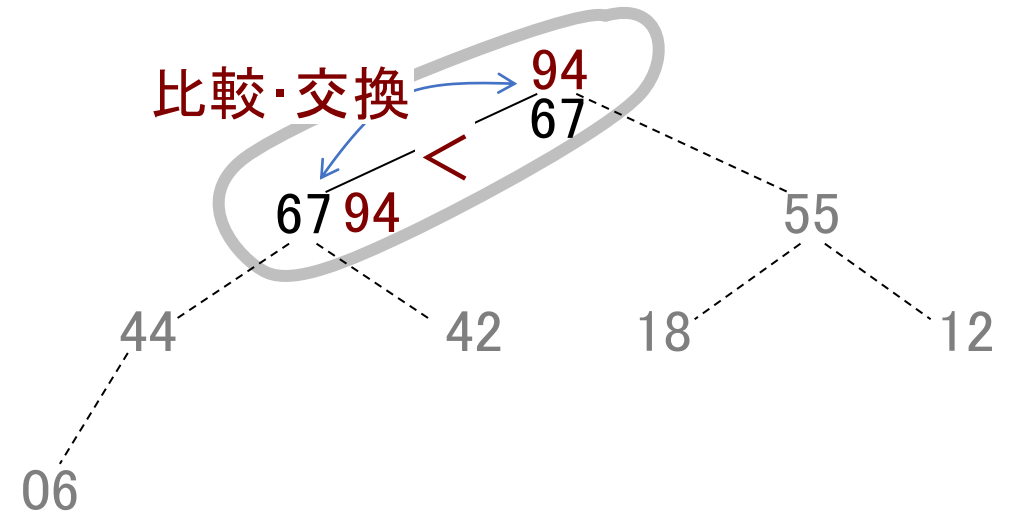
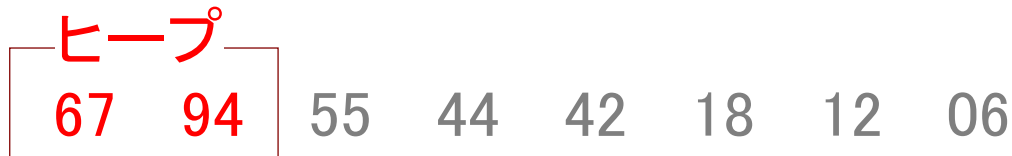
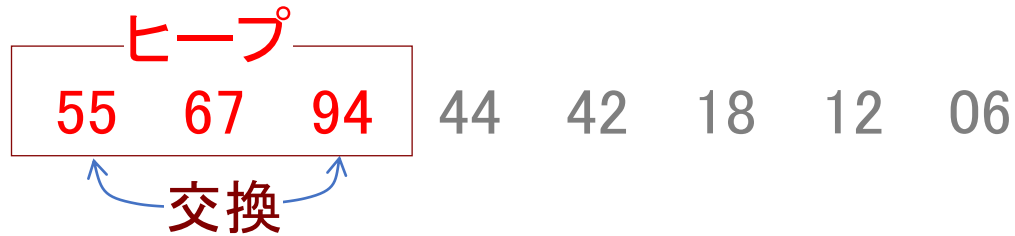
動作例：44, 55, 12, 42, 94, 18, 06, 67  
ソート列を構成



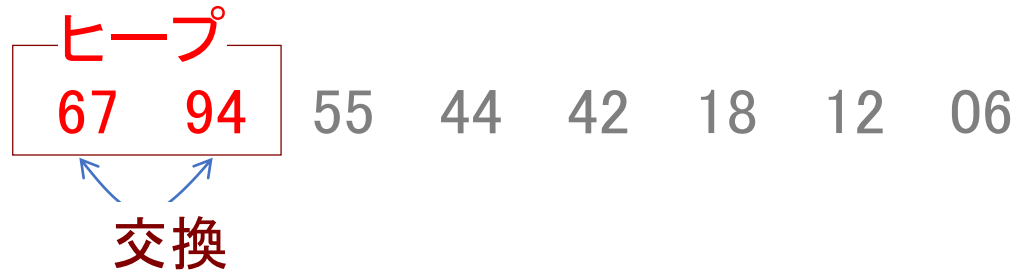
動作例：44, 55, 12, 42, 94, 18, 06, 67  
ソート列を構成



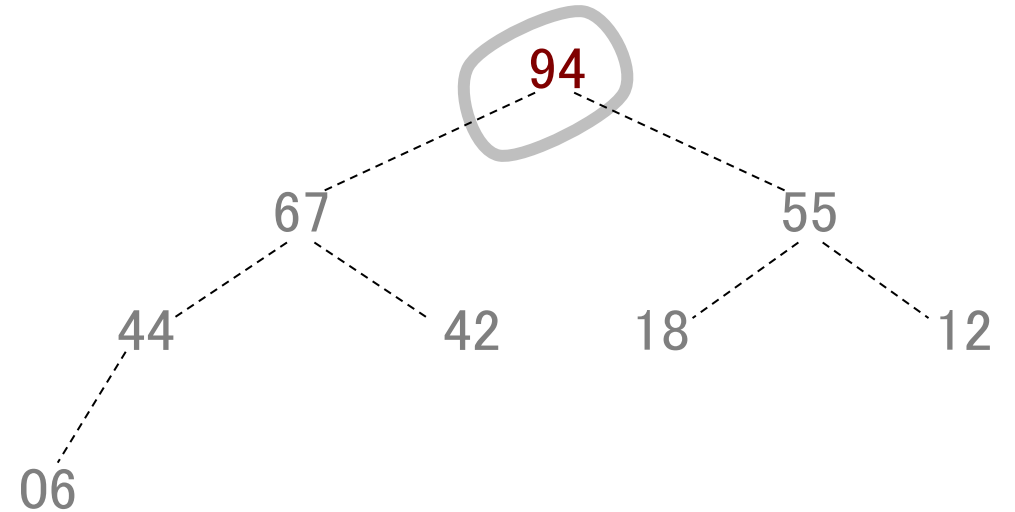
動作例：44, 55, 12, 42, 94, 18, 06, 67  
ソート列を構成



動作例：44, 55, 12, 42, 94, 18, 06, 67  
ソート列を構成



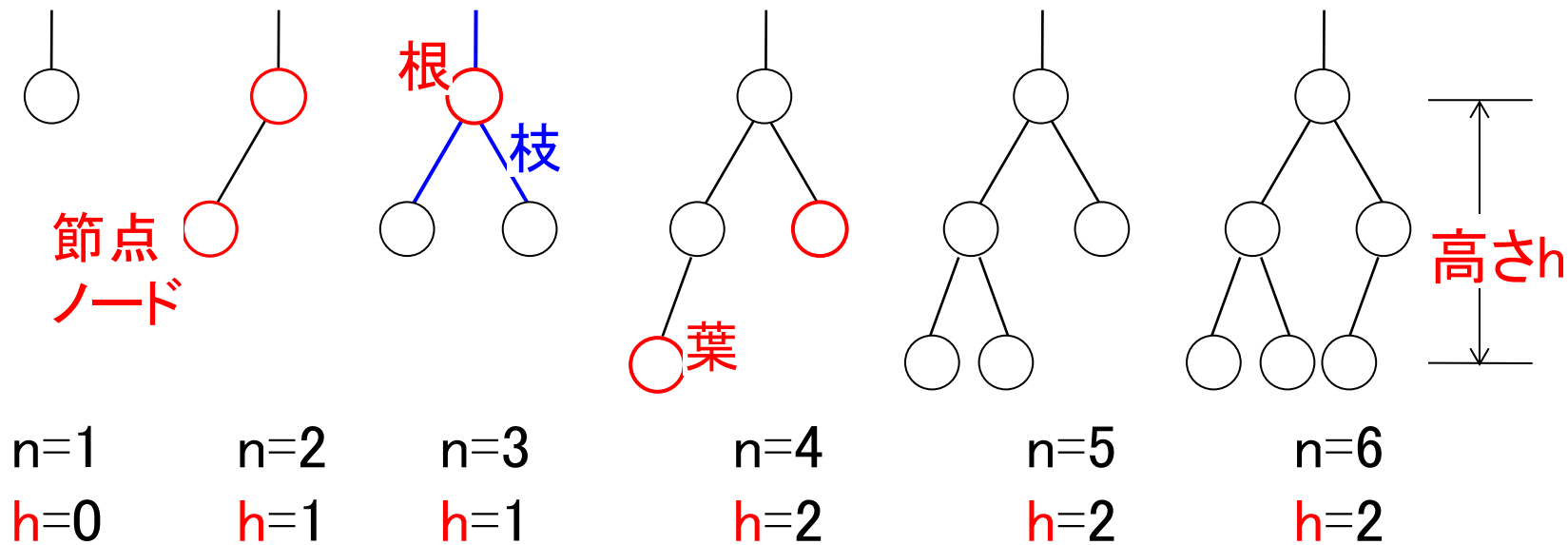
94 67 55 44 42 18 12 06



ソート完了：大きい順

# 2分木の節点数 $n$ と 木の高さ $h$ の関係

- 木の高さ  $h$  とは値から各葉に到る経路中の節点数のうちの「最大数-1」をいう
- 節点を左側から増やしていく2分木を左詰め2分木という



木の高さ  $h = \log_2 n$  ( $n$ は総節点数)

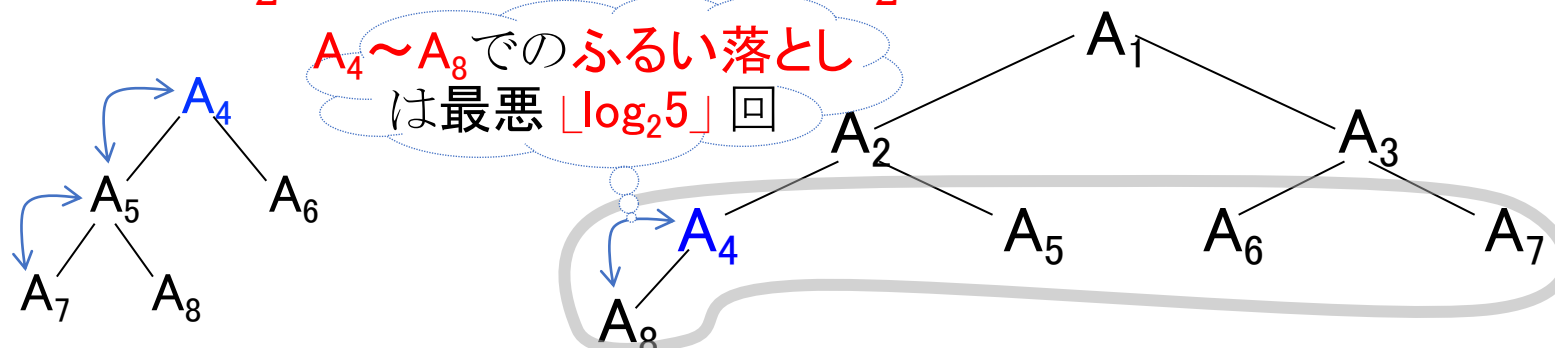


# ヒープソートの時間計算量

- ヒープ構成の時間計算量
  - 最悪で  $O(n\log_2 n)$  かかる
- ソート列構成の時間計算量
  - 最悪で  $O(n\log_2 n)$  かかる
- どちらも同じことから、ヒープソートの時間計算量は
  - 最悪で  $O(n\log_2 n)$  かかる

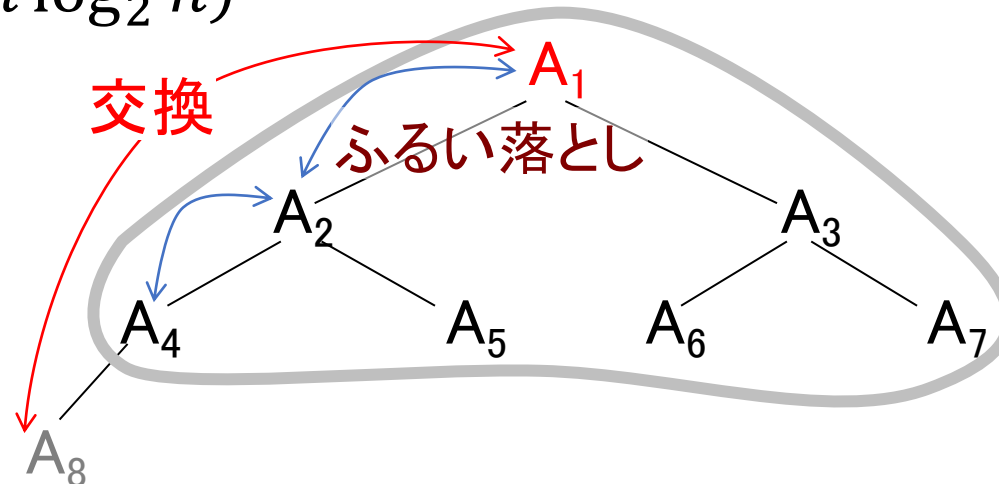
# ヒープ構成の時間計算量

- 最初,  $A_{\frac{n}{2}}$  のふり落とし
  - $A_{\frac{n}{2}}$  から  $A_n$  までの  $n - \frac{n}{2} + 1 = \frac{n}{2} + 1$  個の節点で行われる
  - $\frac{n}{2} + 1$  個からなる左詰め2分木の高さが  $\log_2(\frac{n}{2} + 1)$  なので, 最悪この高さ分に比例する移動が生じる
- 次に  $A_{\frac{n}{2}-1}$  のふり落とし
  - 最悪で  $\log_2(\frac{n}{2} + 2)$  に比例する移動が生じる
- 最後の  $A_1$  のふり落とし
  - 最悪で  $\log_2 n$  に比例する移動が生じる
- $\log_2(\frac{n}{2} + 1) + \log_2(\frac{n}{2} + 2) + \dots + \log_2 n \leq \frac{n}{2} \log_2 n = O(n \log_2 n)$



# ソート列構成の時間計算量

- ヒープ  $A_1, A_2, \dots, A_n$  で,  $A_1$  と  $A_n$  とを交換した後, 再び  $A_1, A_2, \dots, A_{n-1}$  をヒープするために,  $A_1$  のふるい落としを行う.
  - これは  $n-1$  のこの節点の左詰2分木の根から行うので, 最悪  $\log_2(n-1)$  回移動する
- 以降,  $n-2$  個,  $n-3$  個,  $\dots$ , 2 個の節点の左詰2分木での根のふるい落としが行われるので, 最悪で
- $\log_2(n-1) + \log_2(n-2) + \dots + \log_2 2 \leq (n-2) \log_2(n-1)$
- $\leq n \log_2 n = O(n \log_2 n)$



# ヒープソートの領域計算量

- 使用する領域
  - 配列といくつかの変数
- ↓
- ヒープソートの領域計算量
  - $O(n)$

# まとめ

- ソーティングとは
- ソーティングアルゴリズム
  - 3つの単純法
  - ヒープソート

# 演習問題

- 問題1：

- つぎの整数列を3つの単純法を用いてソートせよ。配列の図を描いて、各操作を明示すること。

6   2   7   1   3   5   4

- 問題2：

- つぎの整数列は、10要素の配列の第1要素から第10要素とする。この配列の部分配列でヒープとなる最大長のものを指摘せよ。（ヒント：ヒープになり得る部分配列をすべて列挙する）

50   65   10   17   90   44   12   25   20   33

- 問題3：

- つぎの整数列をヒープソートせよ。ただし、図を用いて、要素の移動・交換を明示すること。

4   5   1   2   6   3   7

# 提出方法

- コンピュータのドローイングソフトなどを利用してもかまいませんが、手書きで結構です。
- 手書きで書いたレポートは、写真に撮って提出してください
  - 単純挿入ソートは、番兵ありとなしのどちらかがあればよいですが、両方あるとなおよいです
  - クイックソートは、第1版と第2版のどちらかがあればよいですが、両方あるとなおよいです
- 提出方法：LETUS
- 締め切り：2023/5/22 10:30まで