

# データ構造とアルゴリズム

## 第9週

掛下 哲郎

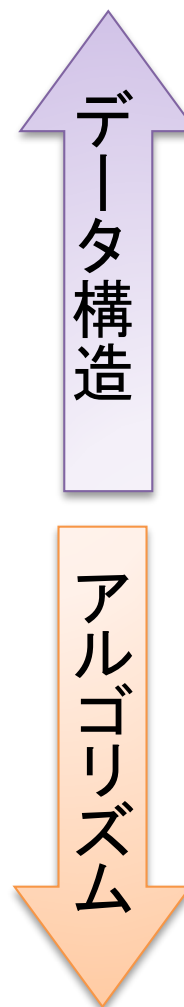
kake@is.saga-u.ac.jp

# 前回のまとめ

- データの整列
  - データを順番に並べるアルゴリズム
- トピック
  1. バブルソート(Bubble Sort)
  2. 選択ソート(Selection Sort)
  3. 挿入ソート(Insertion Sort)
  4. シェルソート(Shell Sort)
  5. ヒープソート(Heap Sort)

# 講義スケジュール

週	講義計画
1-2	導入
3	探索問題
4-5	基本的なデータ構造
6	動的探索問題とデータ構造
7	アルゴリズム演習(第1回)
8-9	データの整列
10-11	グラフアルゴリズム
12	文字列照合のアルゴリズム
13	アルゴリズム演習(第2回)
14	アルゴリズムの設計手法
15	計算困難な問題への対応

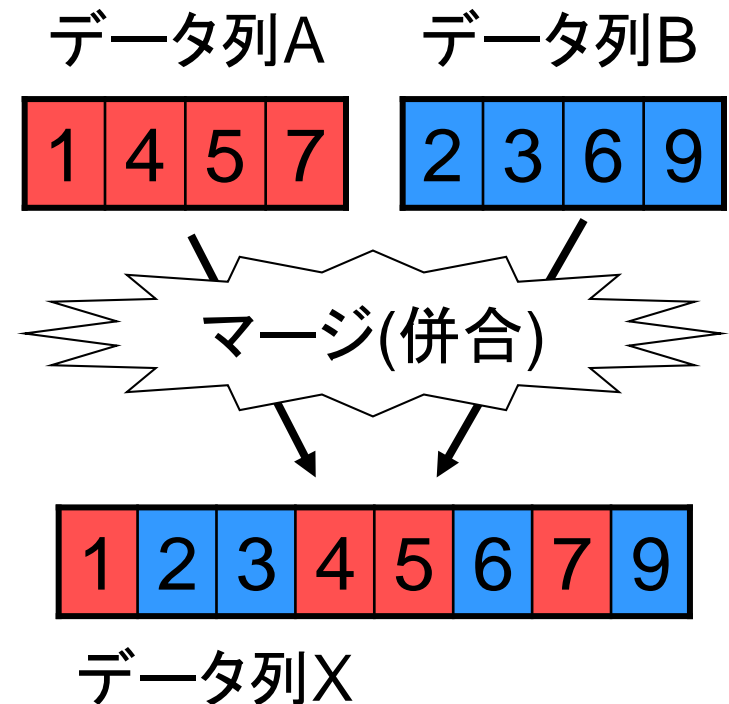


# 今日の内容

- データの整列 (sorting; ソーティング) の続き
- トピック
  - 高度なソート方法
    - ・ マージソート (Merge Sort)
    - ・ クイックソート (Quick Sort)
  - 整列に必要な最小限の比較回数
  - 比較を用いないソート方法
    - ・ バケットソート (Bucket Sort)

# マージソート (merge sort)

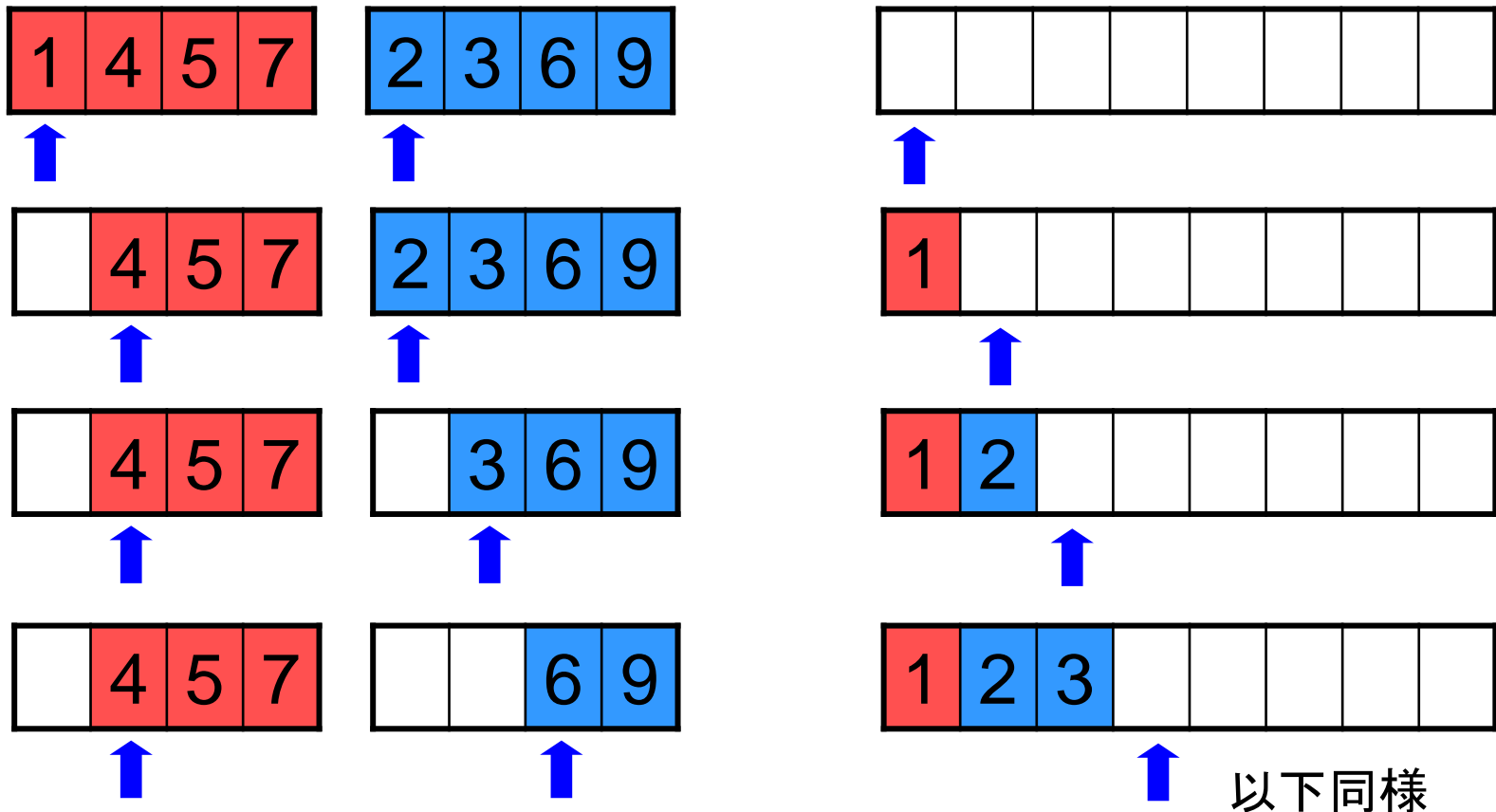
- 「併合ソート」ともよぶ
- 基本手順
  - データ列をA, Bに分割
  - それぞれのデータ列をソート
  - 2つのソート済みデータ列A, Bを1つのソート済みデータ列Xにまとめる(マージ)
- アイデア
  - マージ(併合)処理が高速にできることを利用



# 併合(マージ)の計算量

- 併合は, データ数に比例した手間で済む
- 問題:** マージしたいデータと同じサイズの領域が必要

$$O(|A|+|B|)$$

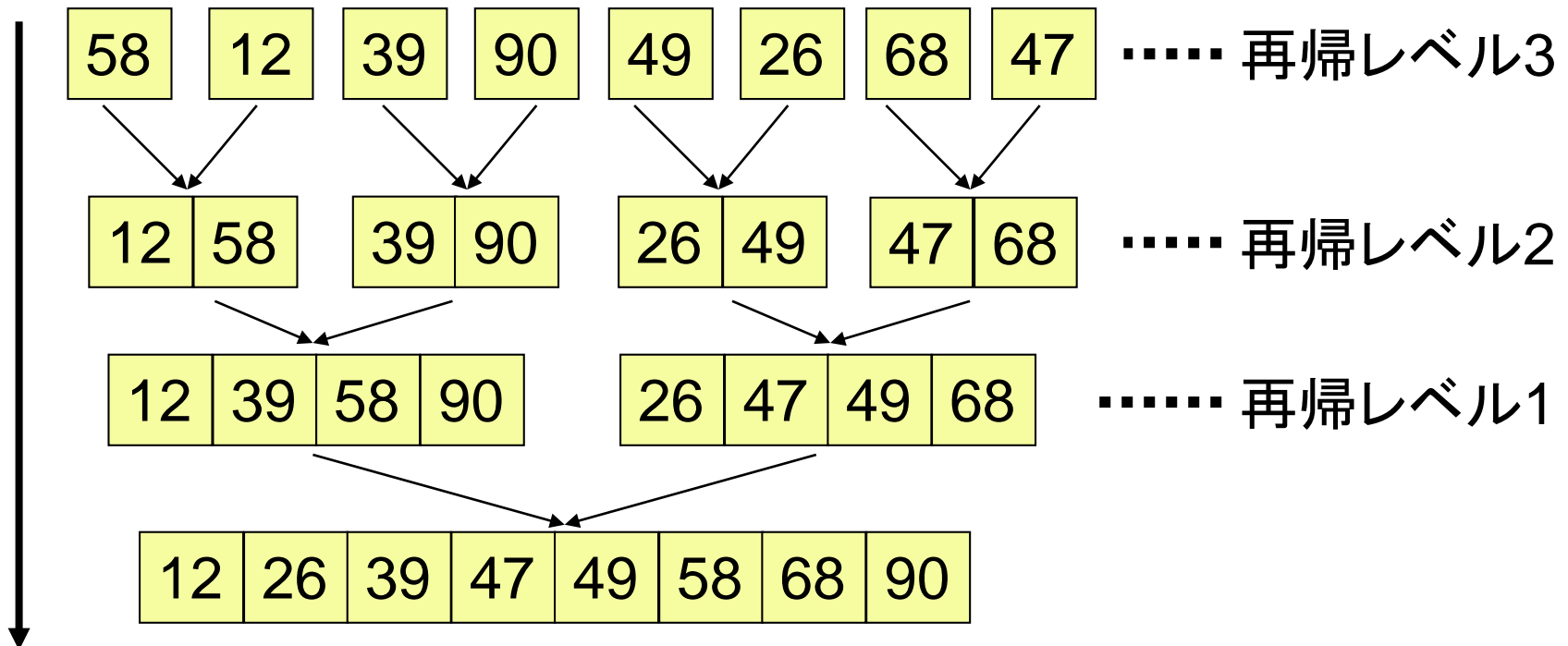


# マージアルゴリズム

- **問題**: 2つのソート済み列AとBを, 1つのソート済み列Xにまとめる.
- **アルゴリズム**
  1. A, B, Xの現在位置を, それぞれの先頭とする.
  2. AかBの現在位置に要素がある限り, 以下の処理を繰り返す.
    - 2.1 AとBの現在位置の要素のうち, 前者が小さいならば以下の処理を行う.
      - 2.1.1 Aの現在位置の要素をXの現在位置に移動する.
      - 2.1.2 Aの現在位置を1つ進める.
    - 2.2 そうでない場合, 以下の処理を行う.
      - 2.2.1 Bの現在位置の要素をXの現在位置に移動する.
      - 2.2.2 Bの現在位置を1つ進める.
    - 2.3 Xの現在位置を1つ進める.

# マージソートの計算手順

- ボトムアップ方式のマージソート
  - 長さ1のデータ列からスタートし、隣同士のデータ列の併合を繰り返す。





# マージソートのアルゴリズム

1. データ列をAとBに分割する.
2. それぞれのデータ列を, マージソートアルゴリズムを用いて再帰的にソートする.
3. 2つのソート済みデータ列A, Bをマージして, 1つのソート済みデータ列Xにまとめる.

再帰呼び出し: ある関数から,  
その関数自身を呼び出す

# マージソートの効率

- データ数 $n$ の場合のマージソートの効率を $S(n)$ とする.
- ステップ1:  $O(n)$
- ステップ2:  $2 \times S(n/2)$
- ステップ3:  $O(n)$

$$S(n) = 2 \times S(n/2) + n$$

$$n = 2^k \text{と置く.} \Rightarrow S(2^k) = 2 \times S(2^{k-1}) + 2^k$$

$$\text{両辺を} 2^k \text{で割る.} \Rightarrow S(2^k)/2^k = S(2^{k-1})/2^{k-1} + 1$$

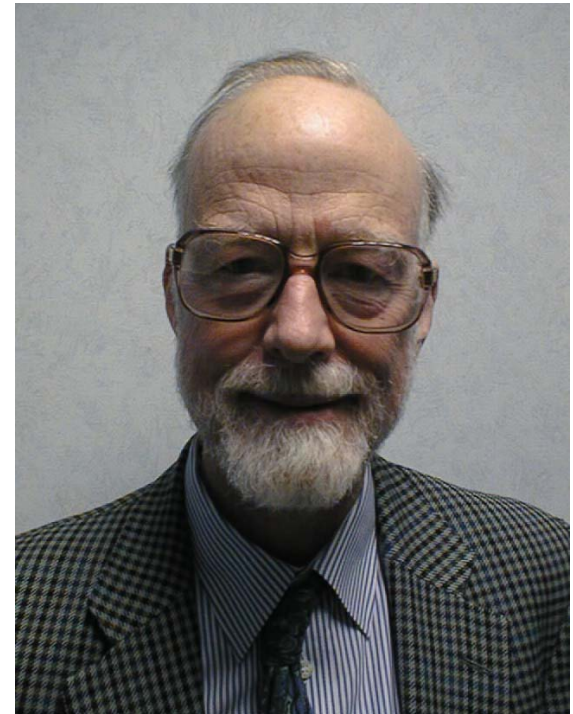
$$S(2^k)/2^k \text{を} T(k) \text{と置く.} \Rightarrow T(k) = T(k-1) + 1$$

$$T(k) = O(k) \Rightarrow S(2^k) = O(k \times 2^k)$$

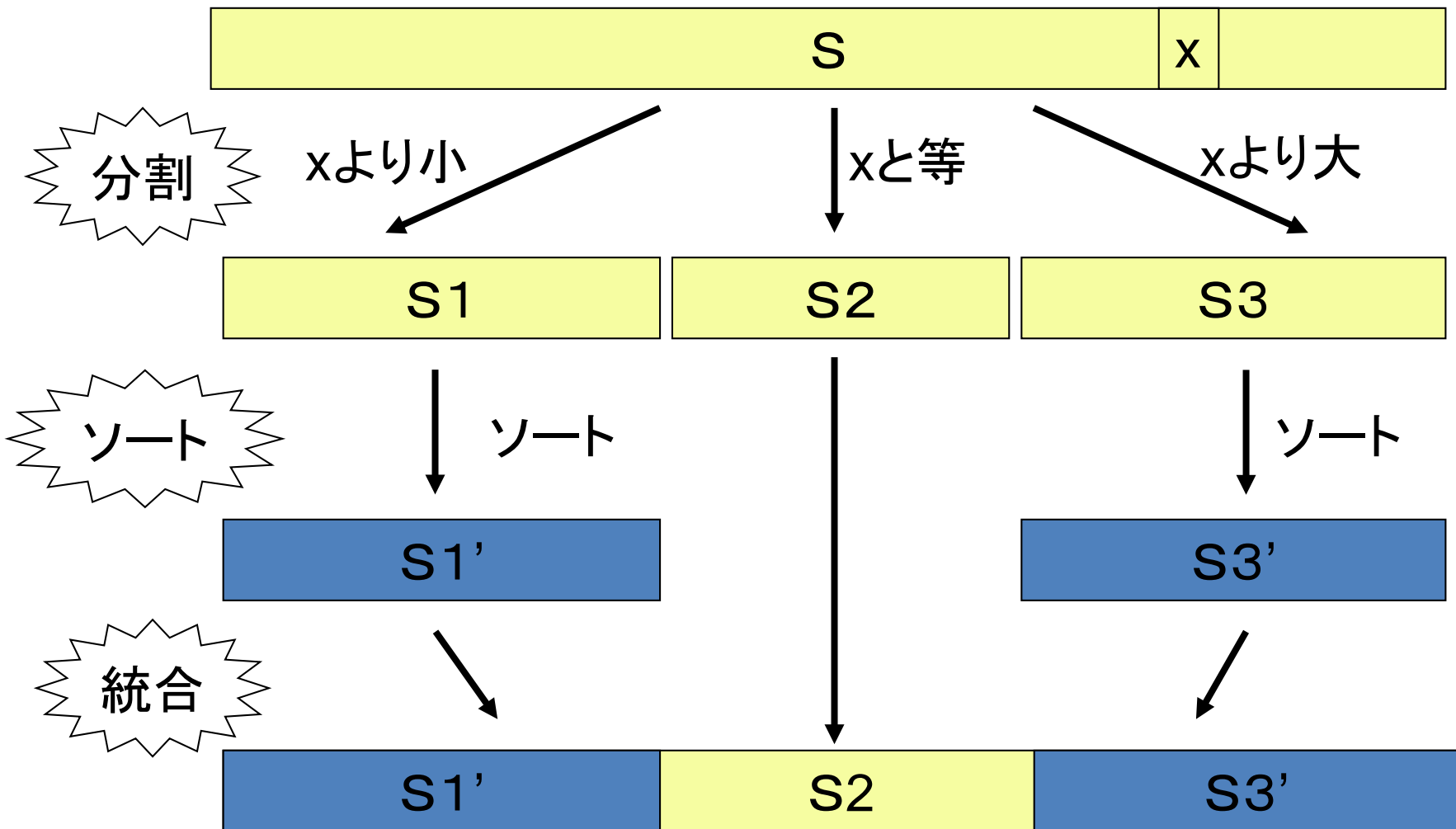
$$n = 2^k \text{より } k = \log_2 n \Rightarrow S(n) = O(n \log n)$$

# クイックソート(Quick Sort)

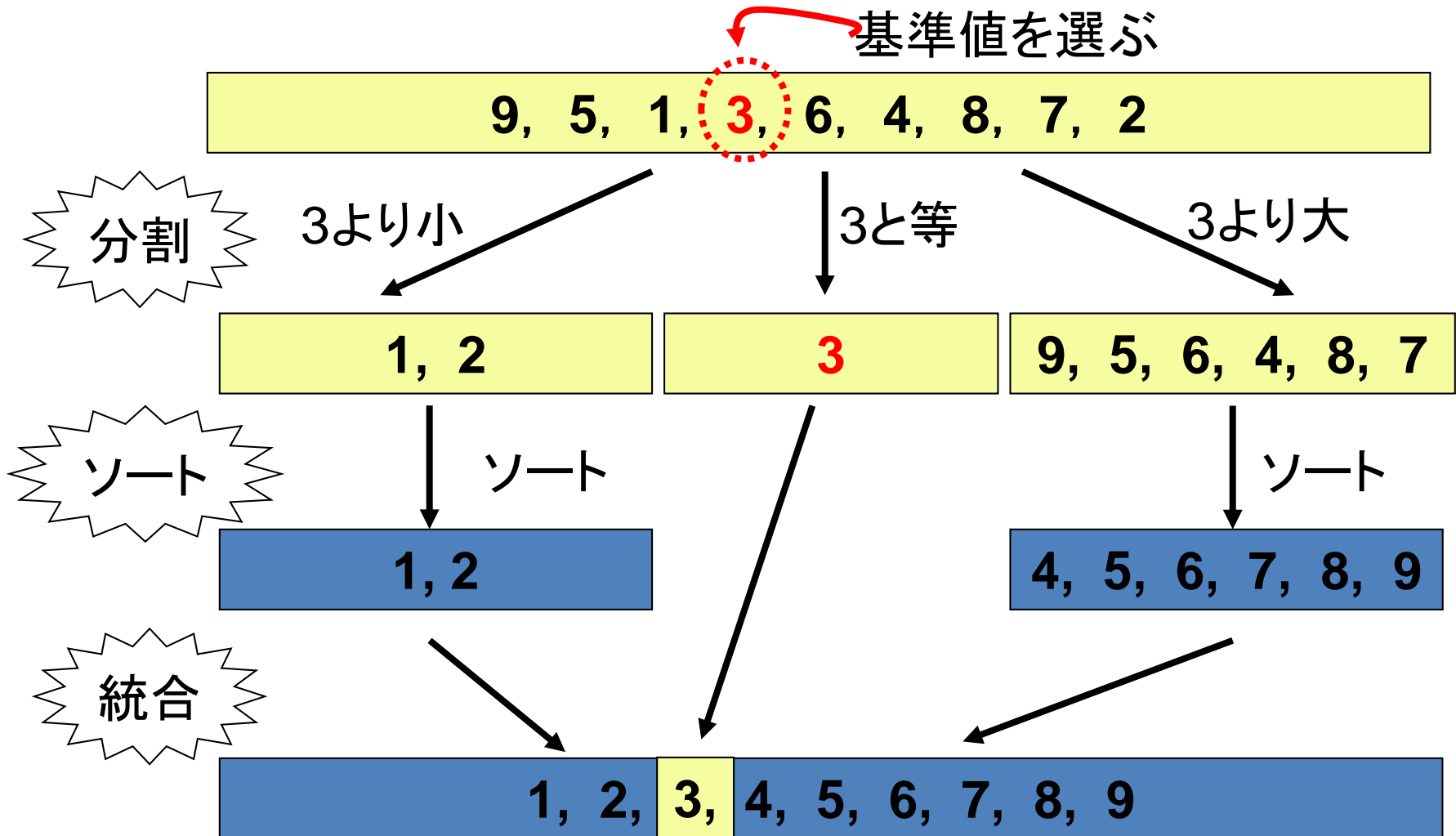
- とても有名. 実際にも広く使用されている.
- **分割統治** (divide and conquer) 法の一つ
- 1960年, Hoare (ホーア)
  - オックスフォード大名誉教授
  - チューリング賞受賞者
  - Hoare理論
  - CSP(並行処理の理論)
  - その他多数の業績



# クイックソートの手順(イメージ)



# クイックソートの手順（具体例）



# クイックソートの手順

(※ 並べたいデータ列をSとする)

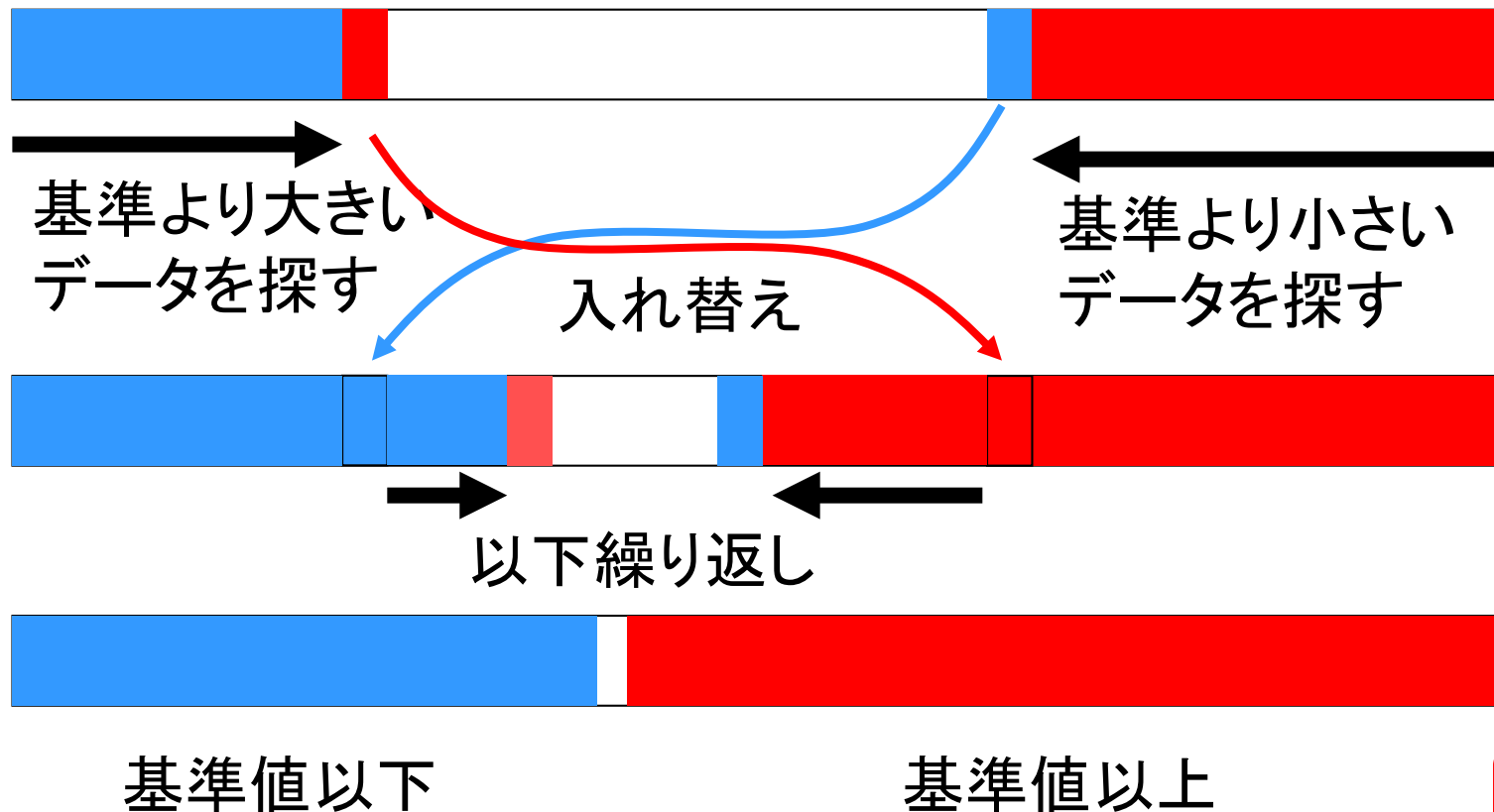
1.  $|S| \leq 1 \rightarrow S$ を返して終了
2. Sから要素を1つ選んでxとする.  
xを基準に, Sを3つのデータ列に分割  
    S1 ..... x より小さいデータが入る  
    S2 ..... x と等しいデータが入る  
    S3 ..... x より大きいデータが入る
3. S1とS3を, クイックソートで再帰的にソート
4. S1, S2, S3 の順に連結して返す.

# 再帰呼び出しの停止性

- 分割後のS1とS3を, クイックソートを使ってそれぞれソート(再帰呼び出し)する.
  - 再帰呼び出し: ある関数から, その関数自身を呼び出す
- $|S2| \geq 1$ より, S1やS3は, Sよりはサイズが小さくなる
  - 分割を続けていけば, 最終的にはサイズが1以下となり, それ以上再帰呼出が続かない.
  - 停止

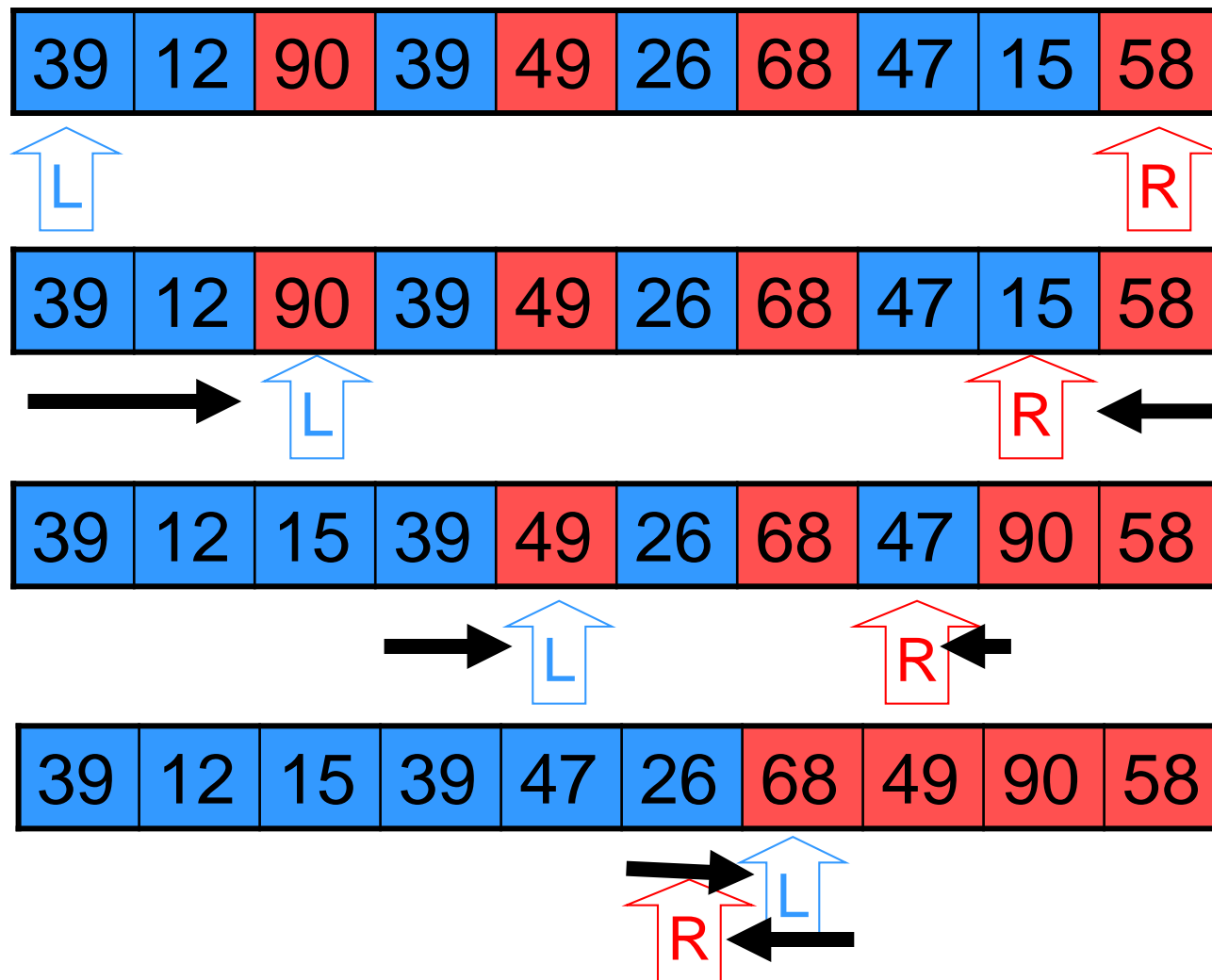
# 基準値による分割 $\rightarrow O(n)$

- アイデア: (※2つに分割)
  - 両端から、真ん中へ向かって処理を進める





# 分割の例



(基準値49  
で分割)

# 分割のアルゴリズム

[問題] 配列に格納された $n$ 個の整数と分割値 $x$ を与える. 配列の先頭部分には $x$ より小さい整数だけを, その後には $x$ より大きい整数だけを格納せよ. ただし,  $x$ より小さい整数と $x$ より大きい整数が少なくとも1つずつは存在すると仮定する.

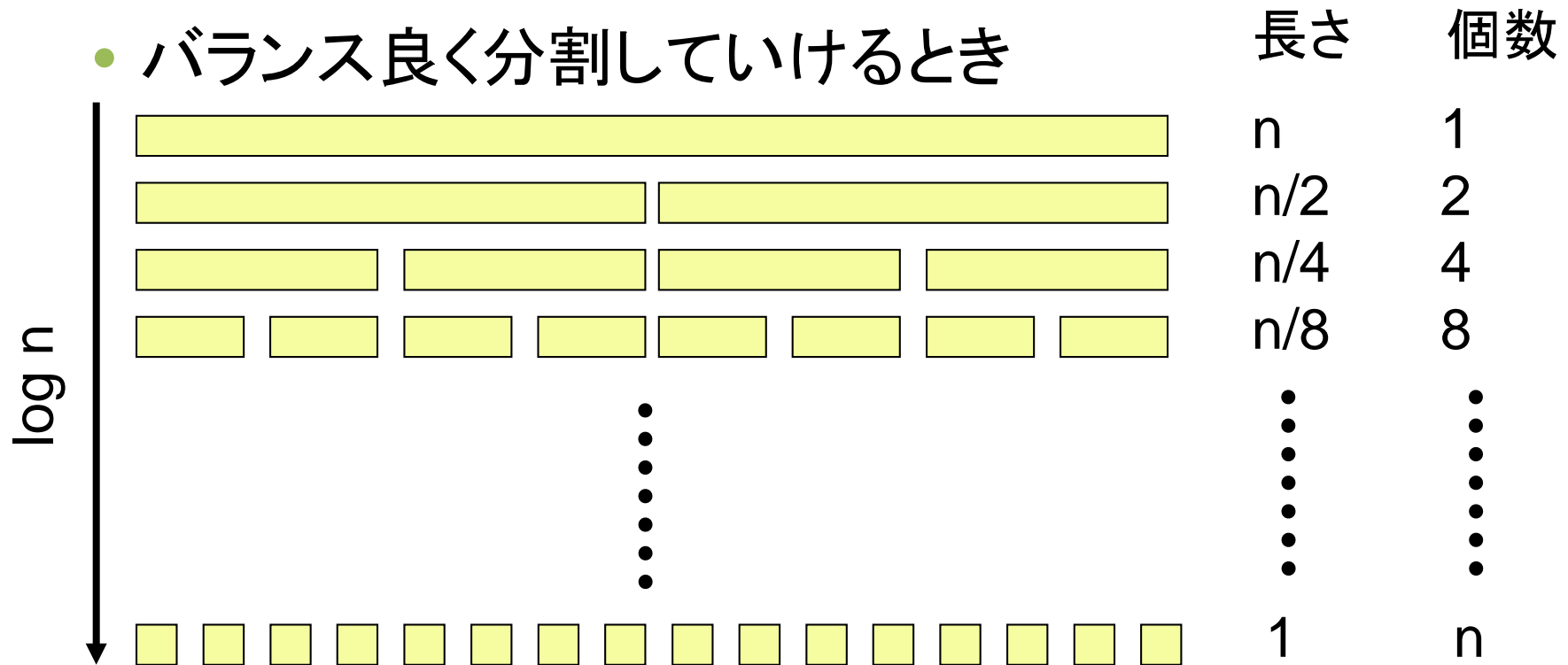
1. 場所1を配列の先頭, 場所2を配列の末尾とする.
2. 場所1が場所2より前にある限り, 以下の処理を繰り返す.
  - 2.1 場所1の要素が $x$ よりも小さい限り, 場所1を1つ進める.
  - 2.2 場所2の要素が $x$ よりも大きい限り, 場所2を1つ戻す.
  - 2.3 場所1の要素と場所2の要素を交換する.
  - 2.4 場所1を1つ進め, 場所2を1つ戻す.
3. 場所1を境界として返す.

# クイックソートのアルゴリズム

1. ソート範囲の要素数が2以下の場合には, 必要に応じて要素を入れ替え, 終了する.
2. そうでなければ以下の処理を実行する.
  - 2.1 分割値 $x$ を, ソート範囲の左端, 中央(小数点以下は切り捨て), 右端の位置にある要素の中央値とする.
  - 2.2  $x$ に基づいて下限と上限の間の要素に分割アルゴリズムを適用し, 境界を求める.
  - 2.3 下限と境界の間にある要素を再帰的にソートする.
  - 2.4 境界+1と上限の間にある要素を再帰的にソートする.

# クイックソートの計算量

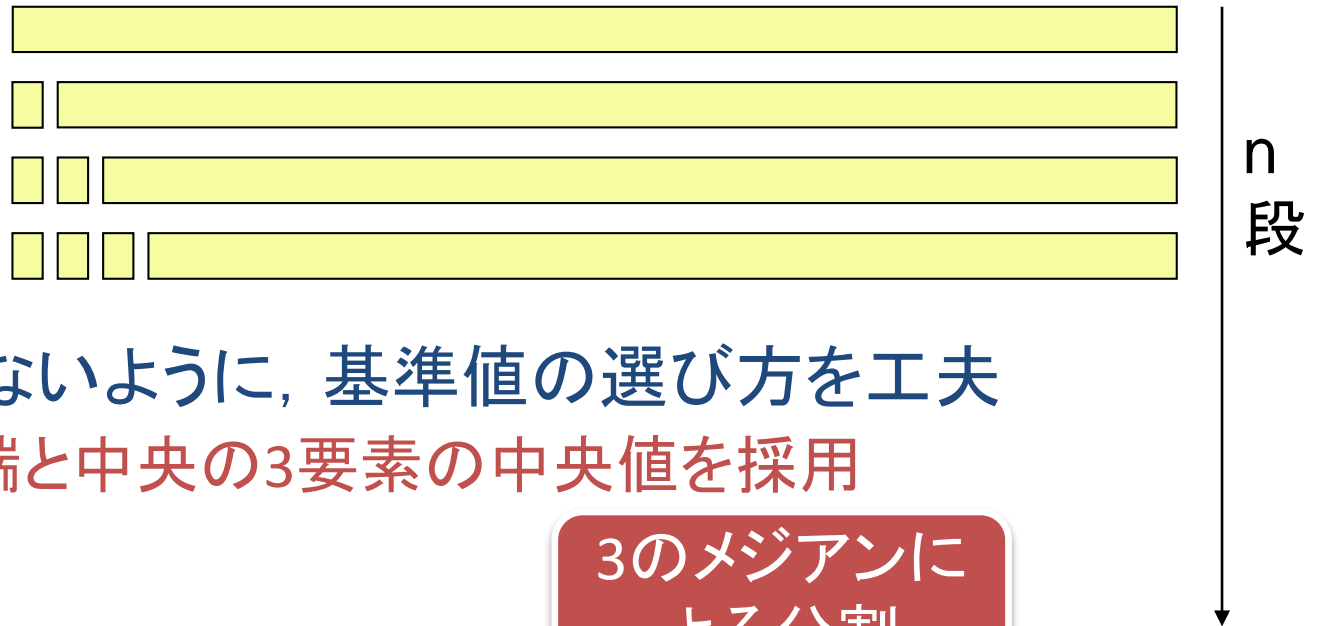
- バランス良く分割していけるとき



分割に必要な手間の合計  $O(n \log n)$  ステップ

# クイックソートの計算量

- 平均的にも,  $O(n \log n)$  ステップ
- ただし, 最悪時は  $O(n^2)$  ステップ
  - 各分割で, 基準値が偶然, 最大値/最小値だった場合



- そうならないように, 基準値の選び方を工夫
  - 例: 両端と中央の3要素の中央値を採用

3のメジアンに  
よる分割

# 整理

- マージソート(Merge Sort)
  - マージ(併合)を利用
  - 常に  $O(n \log n)$  の手間
  - □ マージのために、余分な記憶領域が必要
- クイックソート(Quick Sort)
  - 基準値で分割 → 分割した部分列を再帰的にソート
  - 平均  $O(n \log n)$ , 最悪  $O(n^2)$



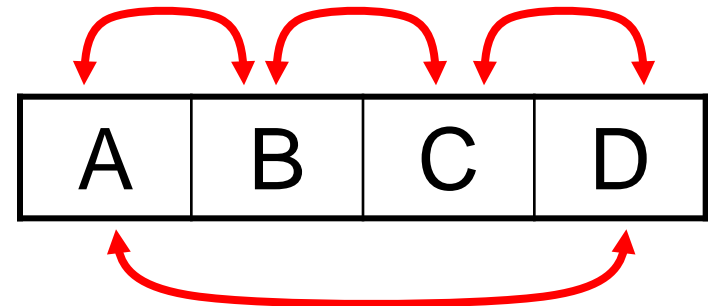
# 整列に必要な手間

- $n$ 個のデータを並べ替えるとき、何回の比較が必要か？
- 正しく並べ替えるには「**必要な情報**」を得なければならない。
- **必要な情報**を得るのに、最悪どれだけの手間がかかるのか？

比較を用いたソーティングアルゴリズム  
のオーダーの最小値は？

# n回で充分？

- n=4で考察
- A,B,C,Dの大小関係を比較によって特定
- 右の例なら  
 $A < B$ 、 $B > C$ 、 $C < D$ 、 $A < D$
- 同じ結果を持つ別のデータ系列が存在
- n回の比較では一意に特定できないことがある  
→ n回の比較だけでは不十分



2	3	1	4
---	---	---	---

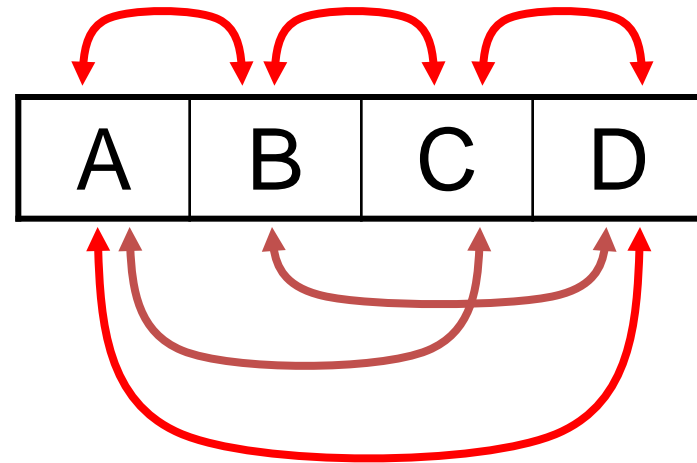
2	4	1	3
---	---	---	---



# $n^2$ 回比較すれば充分

- $n$ データの、全てのペアを比較すれば完全
- $n$ データ中のペアの数

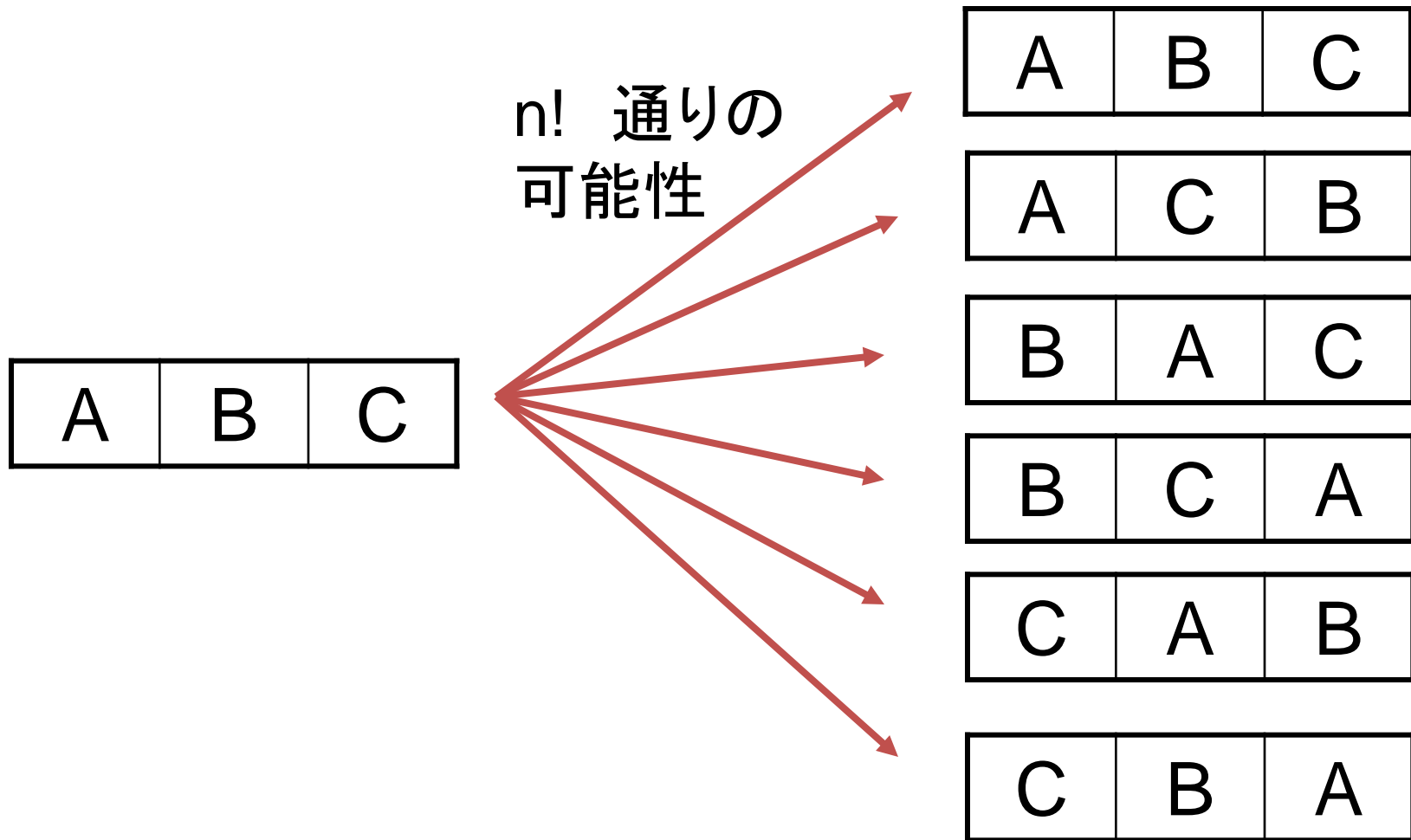
$${}_nC_2 = n(n-1)/2$$



# 整列に必要な手間

- 整列に必要な情報
  - $n$  回の比較では取得しきれない
  - $n(n-1)/2$  回の比較だと取得できる
- 比較に基づいた整列アルゴリズム
  - $O(n)$ のアルゴリズムは理論的に不可能
  - $O(n^2)$ のアルゴリズムは可能。実在。
- 疑問
  - $O(n^2)$ のアルゴリズムは最適？
  - これ以上オーダーを改善できない？

# 比較回数の下界 $\rightarrow O(n \log n)$



# 比較回数の下界 $\rightarrow O(n \log n)$

A	B	C
---	---	---

A	C	B
---	---	---

B	A	C
---	---	---

B	C	A
---	---	---

C	A	B
---	---	---

C	B	A
---	---	---

$A < B$

AとBの大小関係  
が分かったら、大体  
半分に絞り込める

$B < A$

A	B	C
---	---	---

A	C	B
---	---	---

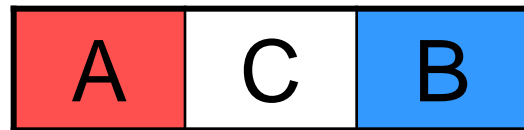
C	A	B
---	---	---

B	A	C
---	---	---

B	C	A
---	---	---

C	B	A
---	---	---

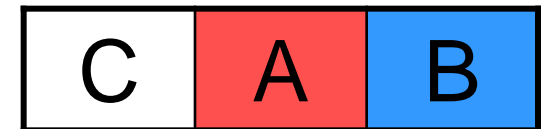
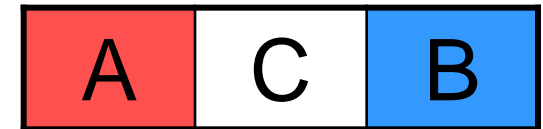
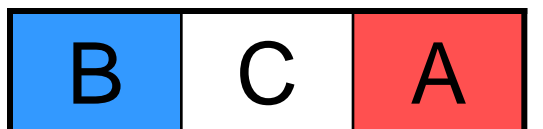
# 比較回数の下界 $\rightarrow O(n \log n)$

 $B < C$ 

BとCの大小  
関係が分か

ると、更に  $C < B$

絞り込める

 $B < C$  $C < B$ 

# 比較回数の下界 $\rightarrow O(n \log n)$

- 最初は  $n!$  パターンの可能性がある
  - $n$ 個の相異なる数字を並べる並べ方
- 1回の比較で、可能性を大体半分に絞り込める
- 最終的に1パターンに絞り込むには  $\log_2 (n!)$  回比較しないとダメな場合がある
  - つまり、最悪時には  $\log_2 (n!)$  回の比較が必要
- $\log_2 (n!) \doteq n \log_2 n$  (スターリングの公式)
  - $\log_2 (n!) = \log_2 1 + \log_2 2 + \log_2 3 + \cdots + \log_2 n$

# バケットソート

比較を行わなければ, より  
高効率にソートできる

- 適用条件
  - データ範囲が1～Nに限定されている
- 手順
  1. 1～Nの番号をつけた箱(バケット)を用意
  2. 各データを対応する箱に入れる
  3. 1番目の箱から順番にデータを取り出し, 並べる
- 時間計算量:  $O(n+N) \Rightarrow$  比較を行わないので効率が良い
- 空間計算量:  $O(N) \Rightarrow$  大容量のメモリが必要  
eg. 「32ビットデータ」  $\rightarrow N=2^{32}=4[G]$

# バケットソートのアルゴリズム

1. 最初のデータを読む.
2. データ $x$ がある限り, 以下の処理を繰り返す.
  - 2.1 データ $x$ を $x$ 番目のバケット(一次元リスト)の先頭に挿入する.
  - 2.2 次のデータを読む.
3. バケットの各要素について, 当該バケットの要素を表示する.



# 本日のまとめ

- マージソート(Merge Sort)
  - マージ(併合)を利用
  - 常に  $O(n \log n)$  の手間
  - マージのために、余分な記憶領域が必要
- クイックソート(Quick Sort)
  - 基準値で分割 → 分割した部分列を再帰的にソート
  - 平均  $O(n \log n)$ , 最悪  $O(n^2)$
- 比較回数の下界 →  $O(n \log n)$
- バケットソート(Bucket Sort)
  - データ範囲が1～Nに限定
  - 効率は良いが、大容量メモリが必要.

# ソートアルゴリズムの特徴(その1)

ソートアルゴリズム	特徴
バブルソート	<ul style="list-style-type: none"><li>• アルゴリズムが最も単純で実現が容易</li><li>• 比較回数, コピー回数が<math>O(n^2)</math>必要.</li></ul>
選択ソート	<ul style="list-style-type: none"><li>• アルゴリズムが比較的単純</li><li>• コピー回数が<math>O(n)</math>で済むため, 各要素のサイズが大きい場合に効率が良い.</li></ul>
挿入ソート	<ul style="list-style-type: none"><li>• データの交換(代入3回)が不要. 移動(代入1回)で済む.</li><li>• 平均比較回数がバブルソート等の半分で済む.</li><li>• データがほぼソートされている場合には, ほぼ<math>O(n)</math>時間で実行できる.</li></ul>
シェルソート	<ul style="list-style-type: none"><li>• 高々<math>O(n^{3/2})</math>時間でソートできる.</li><li>• 多くの場合, 挿入ソートより高速. ヒープソートより単純.</li></ul>
ヒープソート	<ul style="list-style-type: none"><li>• 最悪の場合でも<math>O(n \log n)</math>時間でソートできる.</li><li>• ヒープを保存するために, <math>O(n)</math>の余分な作業領域が必要</li></ul>

# ソートアルゴリズムの特徴(その2)

ソートアルゴリズム	特徴
マージソート	<ul style="list-style-type: none"><li>• 最悪の場合でも<math>O(n \log n)</math>時間でソートできる.</li><li>• マージ処理を実行するために, <math>O(n)</math>の余分な作業領域が必要</li></ul>
クイックソート	<ul style="list-style-type: none"><li>• 平均<math>O(n \log n)</math>時間だが, 最悪<math>O(n^2)</math>時間が必要.</li><li>• 高速ソートアルゴリズムとして有名.</li><li>• 実際にも良く使われている.</li><li>• 余分な作業領域は, 最悪の場合でも<math>O(\log n)</math>で済む.</li></ul>
バケットソート	<ul style="list-style-type: none"><li>• データ範囲が<math>1 \sim N</math>に限定されている.</li><li>• バケットを保持するために, <math>O(N)</math>の余分な作業領域が必要.</li><li>• データの相互比較を行わないため, <math>O(n+N)</math>時間で実行できる.</li></ul>

# 確認テスト(第9回)

- 以下のデータ列を、各種ソーティングアルゴリズムでソートした場合の変更過程を示せ。

7	3	6	1	2	5	4	8
---	---	---	---	---	---	---	---

- バケットソート
- ヒープソート
- マージソート
- クイックソート