



LUDWIG-
MAXIMILIANS-
UNIVERSITY
MUNICH



DEPARTMENT
INSTITUTE FOR
INFORMATICS



DATABASE
SYSTEMS
GROUP

3.2. Divide-and-Conquer-Methoden



- Einfache Sortieralgorithmen reduzieren die Größe des noch zu sortierenden Arrays pro Schritt lediglich um eins.
- **Ziel der Divide-and-Conquer Sortieralgorithmen:**
Halbierung des zu sortierenden Arrays pro Schritt
- Bei Erreichung dieses Ziels: Laufzeit von $O(n \cdot \log n)$
- Allgemeiner Dive-and-Conquer-Algorithmus:

Allgemeiner Divide-and-Conquer-Algorithmus:

```
ALGORITHMUS DivideAndConquerSort;  
  IF Objektmenge klein genug  
  THEN löse das Problem direkt  
  ELSE  
    Divide: Zerlege die Menge in Teilmengen  
             möglichst gleicher Größe  
  
    Conquer: Löse das Problem für jede der Teilmengen  
  
    Merge: Berechne aus den Teillösungen  
            die Gesamtlösung  
  
  END;
```

MergeSort:

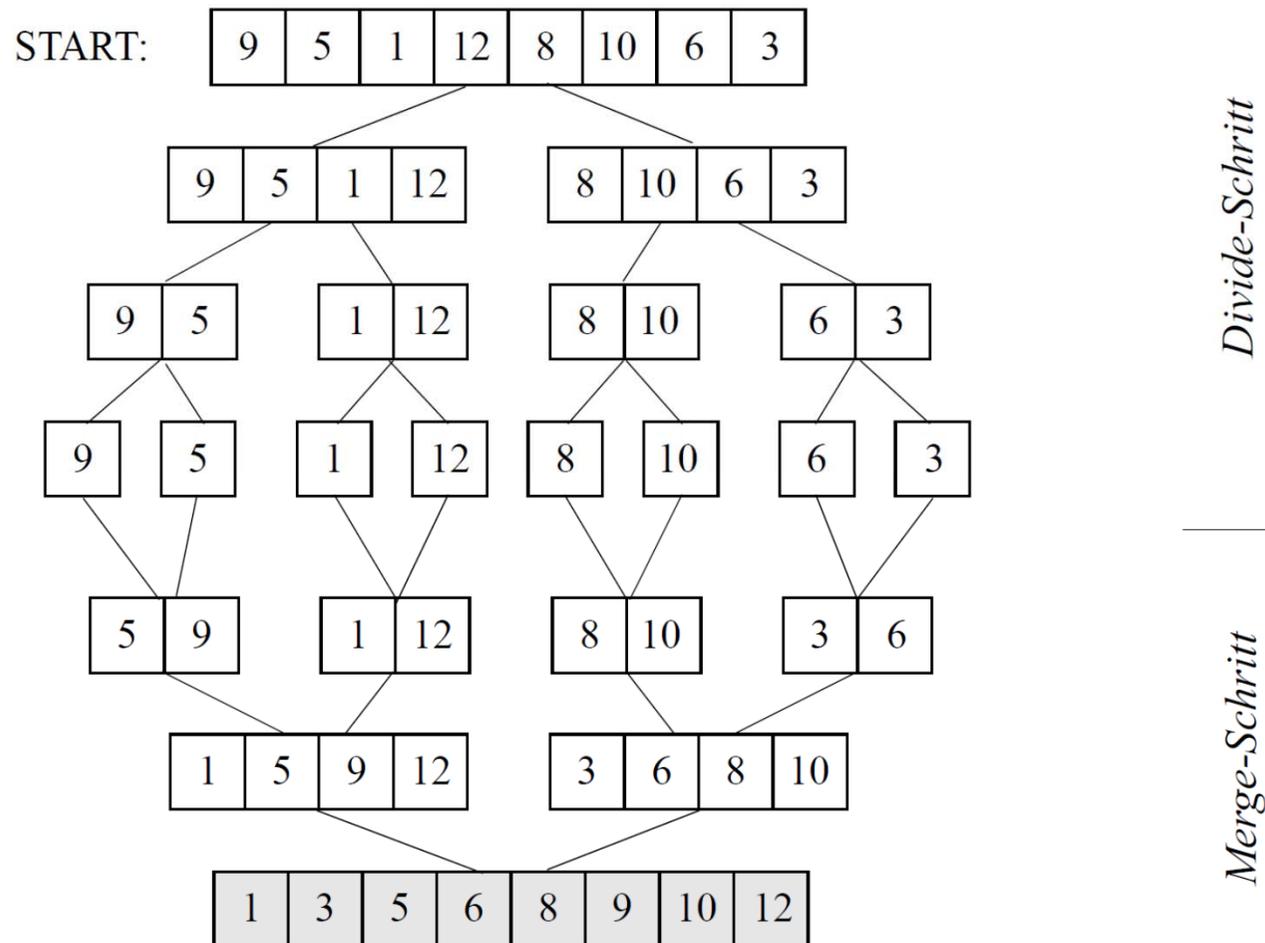
Prinzip:

- Array wird rekursiv bis auf Elementebene zerlegt
- Jedes der Elemente ist trivialerweise sortiert
- Im Anschluss an diese Zerlegung:
Gemäß der Zerlegungshierarchie: Verschmelzung (*Merge*)
von jeweils zwei sortierten Teilarrays, so dass das Ergebnis
wiederum sortiert ist.
- Nach dem letzten Merge-Schritt: gesamtes Array sortiert

Allgemeiner Algorithmus:

```
ALGORITHMUS MergeSort(S);  
  IF |S|=1  
    THEN RETURN S  
  ELSE  
    Divide: S1 :=  $a_1, \dots, a_{\lfloor \frac{n}{2} \rfloor}$ ; S2 :=  $a_{\lfloor \frac{n}{2} \rfloor + 1}, \dots, a_n$ ;  
    Conquer: S1' := MergeSort(S1);  
             S2' := MergeSort(S2);  
    Merge:   RETURN Merge(S1', S2')  
  END;
```

Veranschaulichung:



Laufzeitanalyse:

(Gilt im Wesentlichen für alle Divide-and-Conquer-Verfahren)

Laufzeit $T(n)$ für 1. Rekursionsschritt für n Elemente:

- **Laufzeit des Divide-Schrittes:**
→ $O(1)$ (Aufteilen des Arrays → konstante Zeit)
- **Laufzeit des Conquer-Schrittes:**
→ $2 \cdot T\left(\frac{n}{2}\right)$ (n Elemente → $2 \cdot \frac{n}{2}$ Elemente)
- **Laufzeit des Merge-Schrittes:**
→ $O(n)$ (jedes Element muss betrachtet werden)

Rekursionsgleichung für Zeitbedarf $T(n)$ des MergeSort:

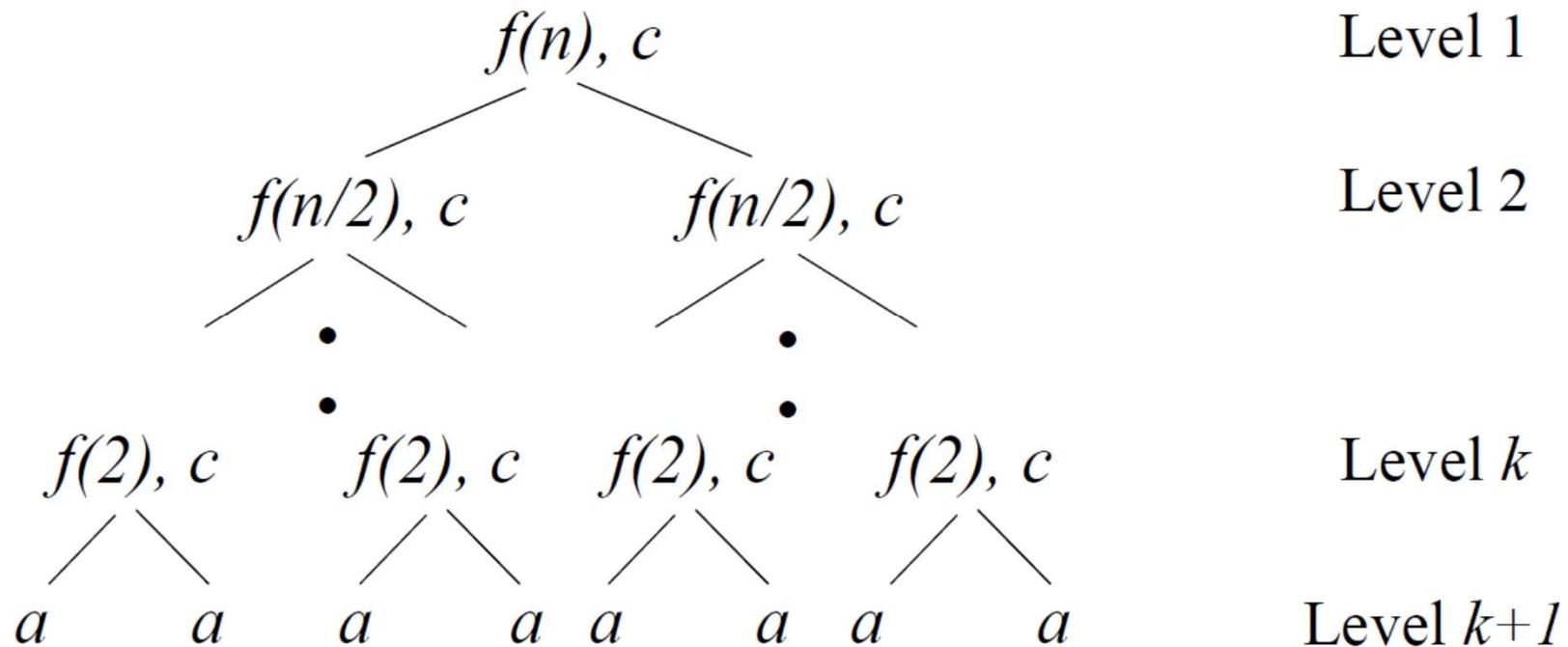
$$T(n) = \begin{cases} O(1) & \text{falls } n = 1 \\ O(1) + 2 \cdot T\left(\frac{n}{2}\right) + O(n) & \text{falls } n > 1 \end{cases}$$

oder anders ausgedrückt (mit Konstanten a und c und einer (mindestens) linear wachsenden Funktion f):

$$T(n) = \begin{cases} a & \text{falls } n = 1 \\ c + 2 \cdot T\left(\frac{n}{2}\right) + f(n) & \text{falls } n > 1 \end{cases}$$

Zur Veranschaulichung:

Baum der Knoten der Rekursionshierarchie für $n = 2^k$:



- Aufwand für alle c 's und a 's: $(n - 1) \cdot c + n \cdot a = O(n)$.
- Da f (mindestens) linear wächst, gilt: $i \cdot f\left(\frac{n}{i}\right) \leq f(n)$
- Hieraus ergibt sich folgende Abschätzung, wobei $k = \log n$:

$$\text{Level 1:} \quad f(n)$$

$$\text{Level 2:} \quad 2 \cdot f\left(\frac{n}{2}\right) \leq f(n)$$

...

$$\text{Level } k - 1: \quad \frac{n}{4} \cdot f(4) \leq f(n)$$

$$\text{Level } k: \quad \frac{n}{2} \cdot f(2) \leq f(n)$$

$$\rightarrow T(n) \leq f(n) \cdot \log n = O(n \cdot \log n)$$

Allgemeines Ergebnis:

- Jedes Divide-and-Conquer-Sortierverfahren, dessen Divide- und Merge-Schritt in $O(n)$ durchgeführt werden kann und das eine balancierte Unterteilung des Problems garantiert, besitzt eine Laufzeit von $O(n \cdot \log n)$.

QuickSort (Hoare 1962):

Prinzip:

- Auswahl eines (beliebigen) Schlüssels x aus dem Array
- **Divide-Schritt:**
Zerlegung des Arrays in Schlüssel $\geq x$ und Schlüssel $< x$
- **Conquer-Schritt:**
Die beiden resultierenden Teilarrays werden rekursiv bis auf Elementebene in gleicher Weise behandelt.
- **Kein Merge-Schritt** durch entsprechende Speicherung der Teilarrays innerhalb des Arrays.

Allgemeiner Algorithmus:

```
ALGORITHMUS QuickSort(S);
  IF |S|=1
    THEN RETURN S
  ELSE
    Divide: Wähle irgendeinen Schlüsselwert  $x=s_j.key$ 
             aus S aus. Berechne eine Teilfolge S1
             aus S mit den Elementen, deren
             Schlüsselwert  $< x$  ist und eine Teilfolge
             S2 mit Elementen  $\geq x$ .

    Conquer: S1' := QuickSort(S1);
              S2' := QuickSort(S2);

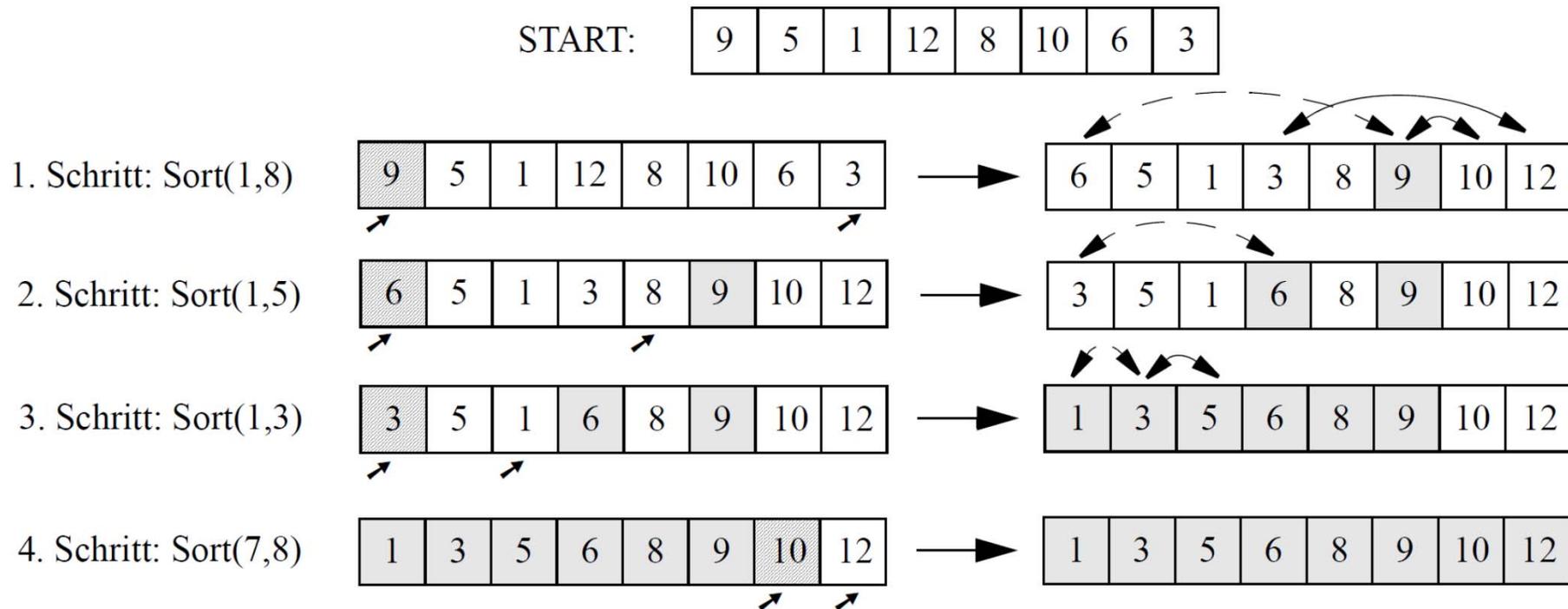
    Merge:   RETURN Concat(S1', S2')
  END;
```

Algorithmus:

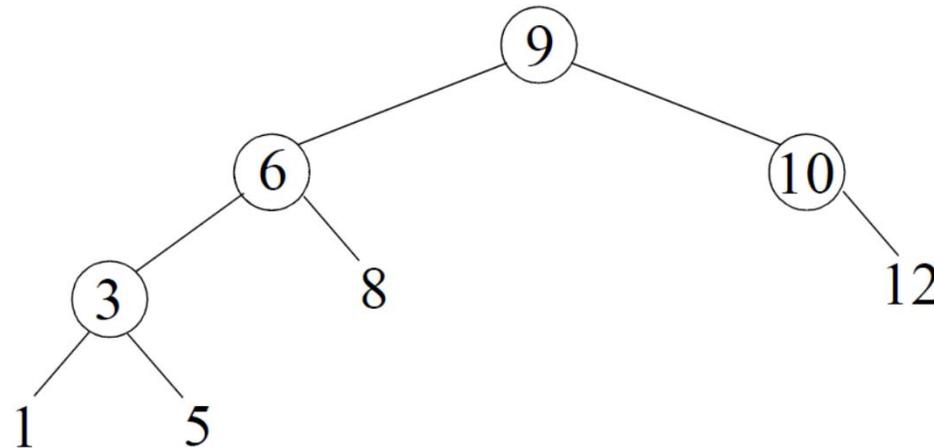
```
void quickSort() {  
    sort(1, n);  
}
```

```
void sort(int li, int r) {  
    int i = li + 1;  
    int j = r;  
    Entry x = a[li]; // hier: Auswahl des linken Randelementes  
    boolean flag = false;  
    do { // Aufteilung in Elemente > x, bzw. Elemente < x  
        while ((a[i].key < x.key) && (i < r)) i = i+1;  
        while (a[j].key > x.key) j = j-1;  
        if (i < j) {  
            vertausche(i, j);  
            i = i+1; j = j-1;  
        } else flag=true;  
    } while (!flag);  
    vertausche(li, j); // x an richtige Stelle bringen  
    if (li < j-1) sort(li, j-1);  
    if (j+1 < r) sort(j+1, r);  
}
```

Veranschaulichung:



- Rekursive Unterteilung des Arrays bezüglich des ausgewählten Elementes
→ binärer Baum über den Arrayelementen



- Höhe des Baumes → Rekursionstiefe des Verfahrens erwünscht:
balancierter binärer Baum → Rekursionstiefe: $O(\log n)$

Laufzeitanalyse:

Zunächst Untersuchung des Zerlegungsschrittes:

Annahme:

Wahl des k -ten Elementes der Sortierungsfolge für den Zerlegungsschritt und Positionierung dieses Elementes an Position k :

- Divide-Schritt:
In jedem Schritt wird jedes Elementes betrachtet
→ n Vergleiche

- Anzahl der Vertauschungsoperationen abhängig von k :
 - (Anzahl der Elemente $< k$) · (Wahrscheinlichkeit für Vertauschung)
 (= Wahrscheinlichkeit, dass das Element im rechten Teil war)
 $= (k - 1) \cdot \frac{n - (k - 1)}{n}$
 - Mittelwert $\overline{M(n)}$ der Vertauschungsoperationen
 über alle $k = 1, 2, \dots, n$ für den ersten Zerlegungsschritt

$$\begin{aligned} \overline{M(n)} &= \frac{1}{n} \cdot \sum_{k=1}^n (k - 1) \cdot \frac{n - (k - 1)}{n} = \frac{1}{n^2} \cdot \sum_{k=0}^{n-1} k \cdot (n - k) = \\ &= \frac{n \cdot (n - 1)}{2 \cdot n} - \frac{2 \cdot n^2 - 3 \cdot n + 1}{6 \cdot n} = \frac{n - \frac{1}{n}}{6} = O(n) \end{aligned}$$

→ der Divide-Schritt benötigt stets $O(n)$ Operationen.

Gesamtzeitbedarf:

- **Günstigster Fall:**

Array wird jeweils in zwei Teilarrays etwa gleicher Größe geteilt

$$\begin{aligned} T_{best}(n) &= c \cdot n + 2T\left(\frac{n}{2}\right) = c \cdot n + 2\left(c \cdot \frac{n}{2}\right) + 4T\left(\frac{n}{4}\right) = \\ &= 2cn + 4T\left(\frac{n}{4}\right) = \dots = \log n \cdot cn + n \cdot T(1) = O(n \cdot \log n) \end{aligned}$$

- **Durchschnittlicher Fall:**

Array wird jeweils bezüglich eines zufällig ausgewählten Elements zerlegt.

Nach Analysen:

Nur um den Faktor $2 \cdot \ln(2) \approx 1,39$ (ca. 40%) schlechter

$$\rightarrow T_{\emptyset}(n) = O(n \cdot \log n)$$

- **Schlechtester Fall:**

Degenerierung (lineare Liste); stets wird das größte/kleinste Element gewählt. (für die obige Implementierung im Falle eines bereits sortierten Elementes)

$$T_{worst}(n) = (n - 1) + (n - 2) + \dots + 1 = \frac{n \cdot (n - 1)}{2} = O(n^2)$$

Beachte hierbei auch die hohe Rekursionstiefe n .

Bemerkungen:

- Austauschen von Schlüsseln über große Distanzen
→ schnell ‚nahezu‘ sortiert
- QuickSort schlecht für kleine n (wie alle komplexen Verfahren)
→ Verbesserung durch direktes Sortieren kleiner Teilarrays
- Methoden, um die Wahrscheinlichkeit des schlechtesten Falls zu vermindern:
 - Zerlegungselement x als Median aus 3 oder 5 Elementen bestimmt
 - ‚Clever QuickSort‘dennoch: QuickSort bleibt immer eine Art ‚Glücksspiel‘, denn es bleibt:
$$T_{worst}(n) = O(n^2)$$
- QuickSort war lange Zeit das im Durchschnittsverhalten beste bekannte Sortierverfahren.



LUDWIG-
MAXIMILIANS-
UNIVERSITY
MUNICH

 DEPARTMENT
INSTITUTE FOR
INFORMATICS

 DATABASE
SYSTEMS
GROUP

3.3. Sortieren mit Hilfe von Bäumen



Sortieren durch Auswählen:

- n -maliges Extrahieren des Minimums/Maximums einer Folge von n Schlüsseln
- Verwenden einer geeigneten Baumstruktur, die diese Operation effizient unterstützt

HeapSort (Williams 1964)

Definition: Heap:

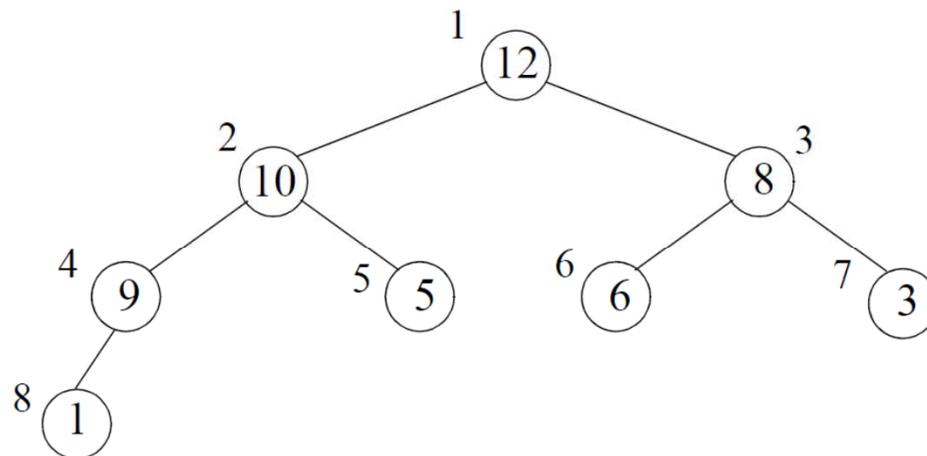
Sei $K = \{k_1, k_2, \dots, k_n\}$ eine Menge von Schlüsseln und T ein binärer Baum, dessen Knoten die Schlüssel aus K speichern.

T ist genau dann ein Heap, wenn gilt:

- T ist vollständig ausgeglichen und
- Der Schlüssel in jedem Knoten ist größer oder gleich den Schlüsseln in den Söhnen des Knotens (*partielle Ordnung*).

Ein Heap ist auf sehr einfache Art mit Hilfe eines sequentiellen Arrays implementierbar.

Beispiel:



1	2	3	4	5	6	7	8
12	10	8	9	5	6	3	1

Vater von Knoten i ist Knoten $\lfloor i/2 \rfloor$

Der linke (rechte) Sohn von Knoten i
ist Knoten $2i$ ($2i+1$).

- Eine Folge von Schlüsseln k_1, k_2, \dots, k_n bildet einen sequentiellen Heap, wenn gilt:

$$k_{\lfloor \frac{j}{2} \rfloor} \geq k_j \quad \text{für} \quad 1 \leq \lfloor \frac{j}{2} \rfloor < j \leq n$$

- In diesem Fall gilt: $k_1 = \max_{1 \leq i \leq n} (k_i)$,
d.h. das Maximum steht in der Wurzel.
- Weiterhin erfüllen alle Blätter $k_{\lfloor \frac{n}{2} \rfloor + 1}, \dots, k_n$ für sich genommen bereits die Heap-Eigenschaft.

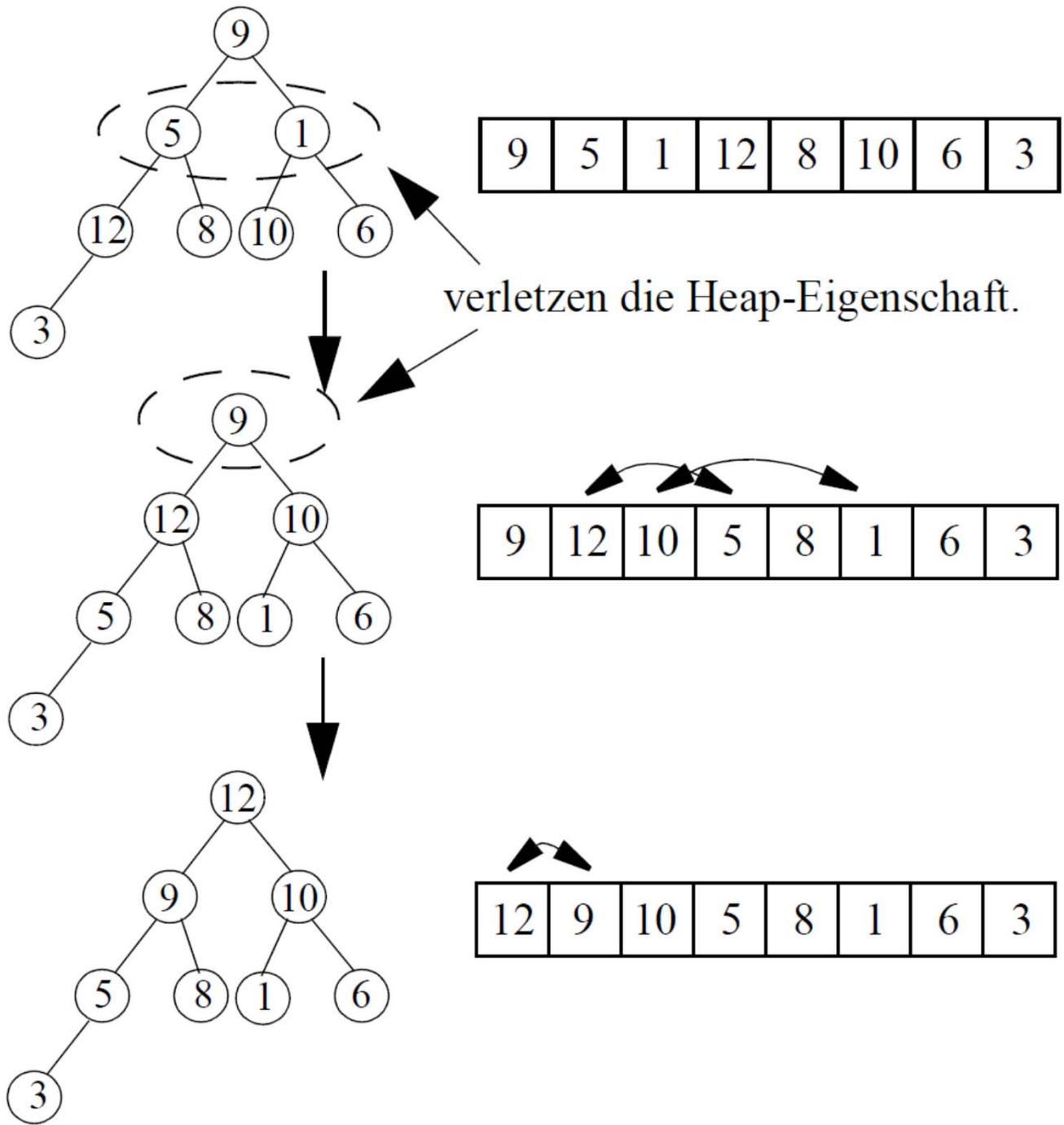
Prinzip:

- Die zu sortierende Schlüssel­folge wird in einem Vorbereitungsschritt in eine sequentielle Heap-Struktur überführt.
- Aus dieser Struktur wird in jedem Schritt das Maximum (die Wurzel) entnommen und der Rest wiederum in Heap-Struktur gebracht.
- Eine sukzessive n -malige Entnahme des Maximums führt zur Sortierung der Folge.

Aufbau der Heap-Struktur:

Ausgangspunkt: unsortiertes Array

- Die $n \text{ DIV } 2$ Elemente in den Blättern erfüllen die Heap-Eigenschaft.
- ‚Absenken‘ aller inneren Schlüssel, die diese Eigenschaft verletzen
 - Erfolgt durch Vertauschen des Schlüssels mit dem größeren seiner Söhne (*sukressive über alle Level*).
- **Ergebnis:** ‚Anfangsheap‘
Der maximale Schlüssel befindet sich in der Wurzel

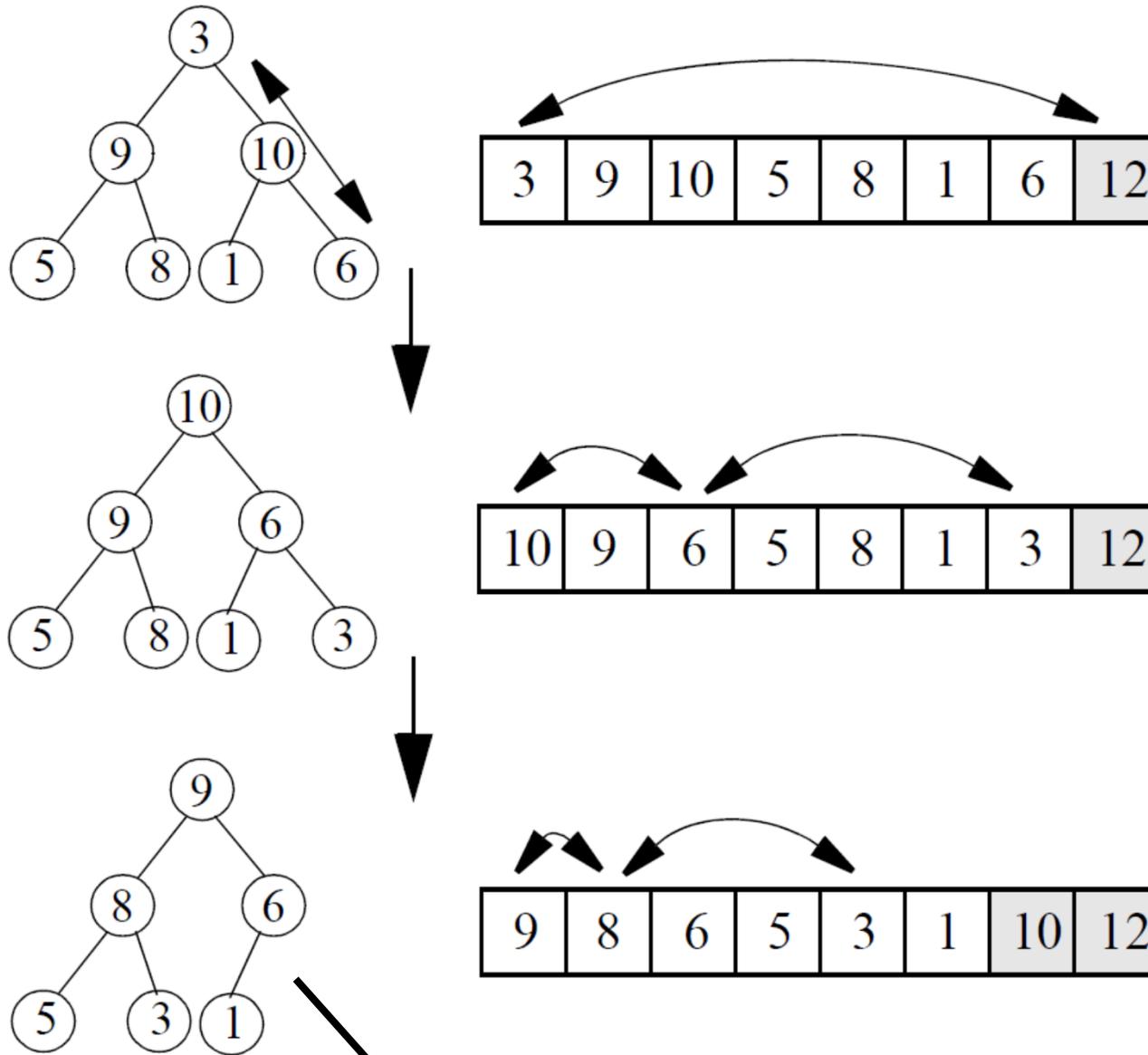


Sortierschritt:

Vertausche: Maximum \leftrightarrow letztes Element

- ‚Absenken‘, um Heap-Eigenschaft zu gewährleisten
- Wiederum steht das Maximum in der Wurzel

Iteration dieses Schrittes, bis das Array sortiert ist.



nächste Folie

