

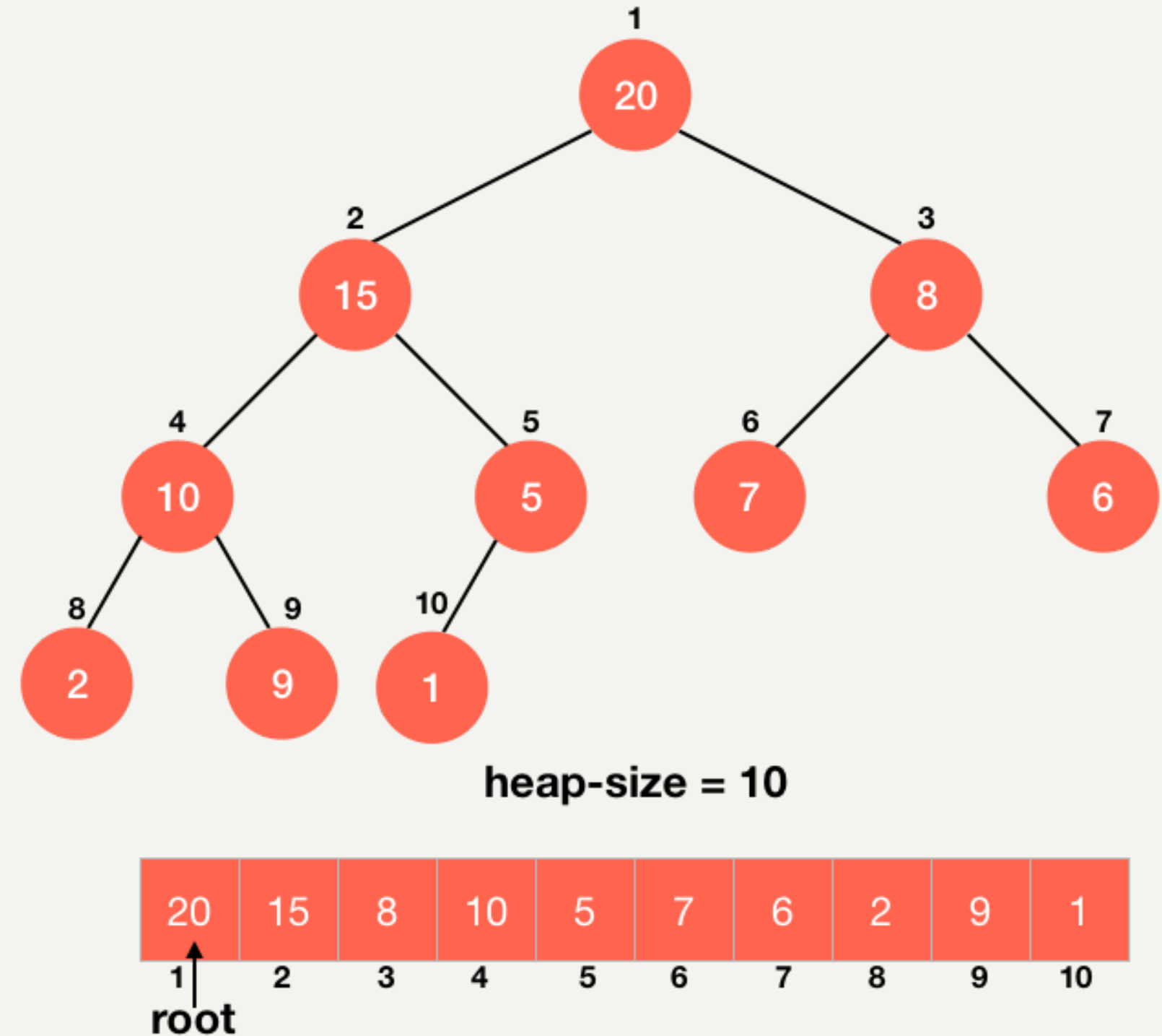
HEAPSORT

(пирамидальная сортировка)

Подготовил: Волков Артемий

ИСТОРИЯ ПИРАМИДАЛЬНОЙ СОРТИРОВКИ

HeapSort является одним из лучших и эффективных методов сортировки в компьютерном программировании и была изобретена Дж. У. Дж. Уильямсом в 1964 году. Сортировка кучей включает создание структуры данных кучи из заданного массива, а затем использование кучи для сортировки массива.



КРАТКАЯ ИСТОРИЯ

ОСНОВНЫЕ ПРИНЦИПЫ РАБОТЫ

Пирамидальная сортировка (или сортировка кучей, HeapSort) — это метод сортировки сравнением, основанный на такой структуре данных как двоичная куча. Она похожа на сортировку выбором, где мы сначала ищем максимальный элемент и помещаем его в конец. Далее мы повторяем ту же операцию для оставшихся элементов.

Двоичная куча — это законченное двоичное дерево, в котором элементы хранятся в особом порядке: значение в родительском узле больше (или меньше) значений в его двух дочерних узлах. Первый вариант называется max-heap, а второй — min-heap. Куча может быть представлена двоичным деревом или массивом.

Поскольку двоичная куча — это законченное двоичное дерево, ее можно легко представить в виде массива, а представление на основе массива является эффективным с точки зрения расхода памяти. Если родительский узел хранится в индексе i , левый дочерний элемент может быть вычислен как $2i + 1$, а правый дочерний элемент — как $2i + 2$ (при условии, что индексирование начинается с 0).

Алгоритм пирамидальной сортировки в порядке по возрастанию:

- 1) Постройте max-heap из входных данных.
- 2) На данном этапе самый большой элемент хранится в корне кучи. Замените его на последний элемент кучи, а затем уменьшите ее размер на 1. Наконец, преобразуйте полученное дерево в max-heap с новым корнем.
- 3) Повторяйте вышеуказанные шаги, пока размер кучи больше 1.

```
1 usage
public void sort(int arr[])
{
    int n = arr.length;

    // Построение кучи (перегруппируем массив)
    for (int i = n / 2 - 1; i ≥ 0; i--)
        heapify(arr, n, i);
    count++;

    // Один за другим извлекаем элементы из кучи
    for (int i=n-1; i≥0; i--)
    {
        // Перемещаем текущий корень в конец
        int temp = arr[0];
        arr[0] = arr[i];
        arr[i] = temp;
        count++;
        // Вызываем процедуру heapify на уменьшенной куче
        heapify(arr, i, 0);
    }
}

// Процедура для преобразования в двоичную кучу поддерева с корневым узлом i, что является
// индексом в arr[]. n - размер кучи
3 usages
void heapify(int arr[], int n, int i)
{
    int largest = i; // Инициализируем наибольший элемент как корень
    int l = 2*i + 1; // левый = 2*i + 1
    int r = 2*i + 2; // правый = 2*i + 2

    // Если левый дочерний элемент больше корня
    if (l < n && arr[l] > arr[largest])
        largest = l;
    count++;

    // Если правый дочерний элемент больше, чем самый большой элемент на данный момент
    if (r < n && arr[r] > arr[largest])
        largest = r;
    count++;
    // Если самый большой элемент не корень
    if (largest ≠ i)
    {
        int swap = arr[i];
        arr[i] = arr[largest];
        arr[largest] = swap;
        count++;
        // Рекурсивно преобразуем в двоичную кучу затронутое поддерево
        heapify(arr, n, largest);
    }
}
```

ОЦЕНКА ВРЕМЕННОЙ СЛОЖНОСТИ

В методе sort, мы используем цикл перегруппирующий массив в кучу, который запускается один раз для создания кучи и имеет время выполнения $O(n)$. Затем, мы вызываем `heapify` для каждого узла, чтобы поддерживать свойство max-heap всякий раз, когда мы удаляем или вставляем узел в кучу. Поскольку узлов 'n', общее время выполнения алгоритма оказывается равным $O(n \log(n))$, и мы используем функцию `max-heapify` для каждого узла.

Математически мы видим, что:

- 1) Первое удаление узла занимает $\log(n)$ времени
- 2) Второе удаление занимает $\log(n-1)$ времени
- 3) Третье удаление занимает $\log(n-2)$ времени
- 4) И так далее до последнего узла, что займет $\log(1)$ времени

Таким образом, суммируя, мы получаем:

$$\log(n) + \log(n-1) + \log(n-2) + \dots + \log(1)$$

так как $\log(x) + \log(y) = \log(x * y)$, мы получаем

$$= \log(n * (n-1) * (n-2) * \dots * 2 * 1)$$

$$= \log(n!)$$

При дальнейшем упрощении (используя приближение Стирлинга) $\log(n!)$ оказывается

$$= n * \log(n) - n + O(\log(n)) = O(n \log(n))$$

приближение Стирлинга:

$$\ln \Gamma(n + 1) = \ln n! = n \ln n - n + O(\ln n)$$

ИНФОРМАЦИЯ О ВХОДНЫХ ДАННЫХ

Входные данные задаются как
(int) (Math.random() * (i+1))
Всего 1000 наборов данных,
начиная с 1100 элементов в одном
наборе, заканчивая 101000
элементами в одном наборе (в
каждом последующем наборе
данных на 100 элементов больше
чем в предыдущем)

```
for (int j = 0; j < 1000; j++) {  
    for (int i = 0; i < (1000 + count); i++) {  
        x = String.valueOf((int) (Math.random() * (i+1)) + ";");  
        fileOutputStream.write(x.getBytes());  
    }  
    count = count + 100;  
}
```

1373; 953; 493; 3140; 818; 38; 1580; 2537; 2299; 2388; 2421; 1533; 2138; 622; 644; 2595; 4711; 285; 4937; 5342; 96; 318; 4727; 4988; 690; 6037; 3181; 385; 3401; 773; 5732; 3264;
0; 1373; 953; 493; 3140; 818; 38; 1580; 2537; 2299; 2388; 2421; 1533; 2138; 622; 644; 2595; 4711; 285; 4937; 5342; 96; 318; 4727; 4988; 690; 6037; 3181; 385; 3401; 773; 5732; 3264;
7; 36; 2792; 2304; 2267; 594; 902; 2395; 708; 129; 1404; 1533; 1441; 647; 1219; 231; 4058; 2323; 2533; 2215; 3728; 3647; 5176; 2541; 5338; 5558; 5796; 1592; 4826; 4332; 3486; 4981; 644;
5; 510; 925; 349; 1891; 2224; 3192; 1448; 2992; 179; 3455; 2752; 2469; 4344; 1470; 995; 2143; 1466; 3001; 3245; 2290; 4407; 112; 4371; 1638; 1190; 1829; 252; 2809; 3007; 2373; 4327; 28;
8; 600; 1073; 9; 2682; 684; 3029; 1134; 2466; 2934; 3325; 1226; 3771; 1595; 1615; 2467; 666; 2269; 4158; 3979; 1269; 3516; 1167; 507; 737; 2124; 33; 4608; 228; 1917; 555; 882; 2487; 5722;
11; 1873; 1012; 711; 381; 821; 354; 2541; 126; 185; 3356; 3080; 3912; 3043; 76; 2708; 743; 2657; 1318; 4698; 535; 3382; 4260; 1966; 16; 4643; 5804; 2088; 1359; 3375; 3775; 4951; 960; 65;
1392; 756; 947; 1833; 1435; 1952; 352; 2818; 3261; 2645; 4037; 2072; 3466; 2328; 415; 382; 3420; 2017; 4280; 1985; 3743; 1163; 3796; 3461; 865; 534; 2841; 5291; 5009; 963; 199; 4417; 57;
1992; 2054; 1494; 1468; 3005; 852; 1562; 1021; 788; 610; 1114; 2644; 3925; 343; 3106; 138; 2226; 2992; 1277; 347; 3366; 4205; 3356; 3872; 4878; 4919; 4265; 3611; 5963; 1989; 733; 4885;
631; 2330; 1083; 1996; 3043; 867; 3412; 3393; 547; 1663; 1718; 3369; 3159; 625; 3699; 4117; 3362; 4304; 3889; 2113; 2293; 3643; 2674; 1527; 3877; 2167; 1189; 3482; 4337; 1870; 398; 27;
749; 2475; 2763; 2468; 638; 3089; 48; 2575; 1704; 1176; 3622; 948; 3134; 2785; 3277; 323; 1975; 239; 3888; 4240; 4422; 1878; 3198; 1370; 2273; 3954; 2027; 5239; 6164; 5672; 319; 1742;
2139; 1134; 439; 1452; 2298; 2209; 562; 1252; 3235; 645; 2349; 985; 1027; 1661; 2645; 3125; 147; 4613; 2000; 1985; 3743; 1163; 3796; 3461; 865; 534; 2841; 5291; 5009; 963; 199; 4417; 57;
8; 211; 2468; 2887; 1479; 2664; 2785; 1879; 1574; 2175; 3504; 2140; 2151; 366; 3957; 2638; 2802; 2813; 3778; 4128; 2377; 2225; 2258; 4254; 2432; 3434; 722; 2391; 3467; 1034; 1061; 4722; 2400;
2067; 1068; 934; 1977; 2201; 3035; 1366; 676; 81; 3122; 825; 1688; 3185; 1186; 1175; 2378; 4353; 3917; 3817; 5329; 1953; 2560; 2485; 339; 563; 5700; 1127; 2400;
2744; 113; 845; 245; 1195; 1412; 2243; 466; 126; 3318; 212; 388; 2164; 1043; 2895; 4149; 3199; 4903; 3552; 3755; 2307; 4253; 2955; 1388; 223; 513; 1184; 5836; 4730; 6035; 2779; 316; 3;
1767; 1993; 1129; 438; 1981; 2505; 2955; 1093; 3640; 2674; 3221; 2183; 2546; 3982; 1027; 2122; 596; 3176; 4283; 5120; 4739; 1684; 2619; 5236; 4919; 1027; 4614; 1866; 2420; 3351; 3060;
2339; 1883; 3006; 2955; 3272; 1969; 2729; 126; 3451; 3153; 3441; 2402; 2495; 455; 2450; 860; 3875; 173; 4627; 5181; 1302; 189; 789; 3115; 4480; 453; 4890; 6125; 734; 5085; 5887; 3830;
534; 86; 2142; 593; 3211; 2916; 1999; 3482; 1736; 2340; 1238; 2533; 2774; 2973; 3751; 737; 3044; 936; 4731; 1107; 2511; 2881; 2056; 5281; 2902; 14; 1611; 1005; 2172; 5536; 3939; 3094;
1101; 555; 1304; 1037; 1329; 1825; 348; 217; 1302; 804; 3256; 2820; 3373; 4039; 223; 587; 3925; 4120; 2607; 2350; 1932; 96; 2691; 1539; 3519; 764; 1327; 3336; 5655; 4793; 4564; 4071;
2754; 1165; 690; 423; 2816; 1783; 1641; 215; 523; 1462; 701; 238; 1281; 288; 4061; 731; 2027; 2185; 2671; 337; 2554; 3855; 5513; 3415; 4898; 4174; 4833; 5989; 1002; 3291; 5552; 834; 1177; 3057; 39;
1387; 2389; 1757; 2959; 1660; 727; 1899; 2771; 494; 1640; 3616; 3058; 4126; 2578; 4243; 4201; 2661; 286; 836; 1572; 4412; 4658; 3738; 1616; 4723; 1594; 5105; 5303; 310; 4046; 5910; 45;
8; 2571; 1756; 84; 2416; 1024; 3403; 2713; 1454; 2453; 646; 2945; 570; 652; 519; 4519; 816; 4591; 2043; 4527; 3956; 5126; 2167; 477; 456; 441; 2894; 4664; 5129; 3763; 3067; 463; 1587; 3;
29; 125; 2780; 2346; 2020; 7; 2745; 1329; 1703; 470; 919; 3217; 3067; 3792; 3172; 948; 4597; 2851; 4216; 757; 5117; 1362; 4835; 2327; 1932; 263; 2180; 1732; 299; 6140; 4525; 673;
2434; 2432; 1578; 659; 2427; 612; 3048; 54; 2876; 1616; 1401; 484; 167; 2356; 492; 1055; 622; 2045; 4162; 966; 2777; 3592; 4176; 1723; 5281; 2143; 5005; 4890; 6051; 1822; 6047; 5405;
64; 939; 293; 2736; 1403; 1484; 470; 145; 3134; 2099; 3888; 2012; 1789; 1638; 4050; 1728; 4195; 3199; 4500; 4479; 148; 2569; 5356; 4140; 5076; 1425; 687; 5786; 320; 3559; 2065; 797; 44;
62; 1412; 2159; 966; 2951; 3086; 861; 195; 864; 2957; 1526; 3552; 3528; 2190; 2726; 48; 2288; 4104; 1743; 3187; 4247; 2204; 2563; 3683; 1852; 4453; 4030; 5013; 2503; 592; 5461; 3040; 2;
4; 592; 1917; 2953; 3301; 2958; 2325; 3651; 2517; 209; 3469; 1533; 2910; 1812; 3609; 3285; 1589; 1446; 2412; 2322; 2781; 5122; 618; 2577; 3558; 947; 25; 3857; 4188; 2968; 1028; 6615; 2;
53; 2626; 2051; 3068; 801; 983; 891; 2874; 1667; 2235; 2275; 3640; 499; 739; 4004; 2112; 4293; 2137; 1845; 4235; 3833; 2306; 2638; 2425; 4835; 2624; 5050; 4358; 4322; 1934; 1397; 1859;
2; 1182; 228; 788; 1907; 2413; 2440; 3403; 318; 3113; 2288; 3747; 1926; 3652; 3707; 609; 3468; 4226; 2451; 696; 3142; 1352; 1420; 4322; 742; 4594; 2059; 2223; 5317; 5465; 3767; 5231; 3;
772; 1922; 2184; 1880; 1784; 613; 3255; 1357; 1260; 3517; 1031; 1088; 4216; 1526; 40; 1350; 3583; 2325; 4542; 4388; 4675; 4138; 2418; 2765; 3385; 2136; 3844; 2319; 5765; 4969; 4417;
55; 790; 1232; 1222; 2564; 1006; 3425; 1054; 3631; 274; 143; 438; 4067; 1580; 3254; 1853; 1901; 4432; 1487; 93; 993; 198; 4635; 787; 3939; 5207; 364; 580; 4755; 4688; 1884; 3276; 2914;

ТАБЛИЦА ПОЛУЧЕННЫХ ЗНАЧЕНИЙ КОЛИЧЕСТВА ИТЕРАЦИЙ И ВРЕМЕНИ РАБОТЫ В ЗАВИСИМОСТИ ОТ РАЗМЕРА ДАННЫХ

количество итераций	время (нс)	длина набора
32352	820200	1100
68149	183400	1200
107406	151900	1300
150054	159700	1400
196363	174500	1500
246186	189500	1600
299541	197700	1700
356371	212200	1800
416751	234000	1900
480630	232400	2000
548134	240300	2100
619320	287500	2200
694272	269700	2300
772981	285100	2400
855549	296400	2500
941886	306600	2600
1031872	328100	2700
1125783	335800	2800
1223403	303100	2900
1324840	340600	3000
1430226	261400	3100
1539396	224700	3200
1652362	207700	3300
1769208	207500	3400
1889733	216800	3500
2014201	220400	3600
2142261	230200	3700
2274210	245900	3800
2410162	242600	3900
2549922	249100	4000
2693310	337800	4100
2841007	260800	4200
2992563	270800	4300
3148149	336200	4400
3308005	291500	4500
3471978	335800	4600
3639837	340000	4700
3811606	397400	4800
3987588	413500	4900
4167696	484800	5000
4351789	342100	5100
4540092	343400	5200
4732506	346000	5300
4929001	355300	5400
5129559	362300	5500
5334369	370600	5600
5543113	377700	5700
5755962	371600	5800
5972664	375700	5900
6193483	378200	6000

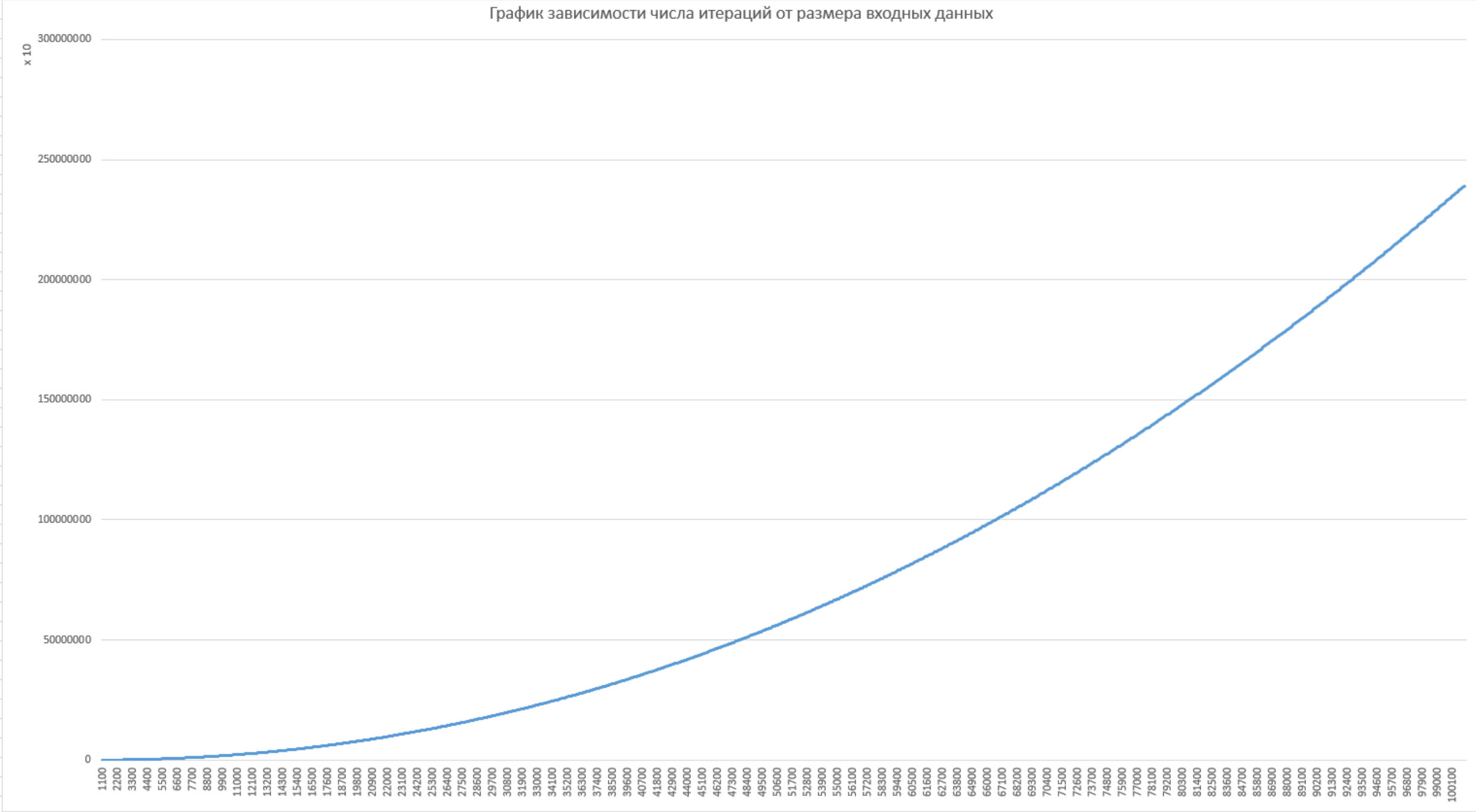


График зависимости количества итераций от объема входных данных

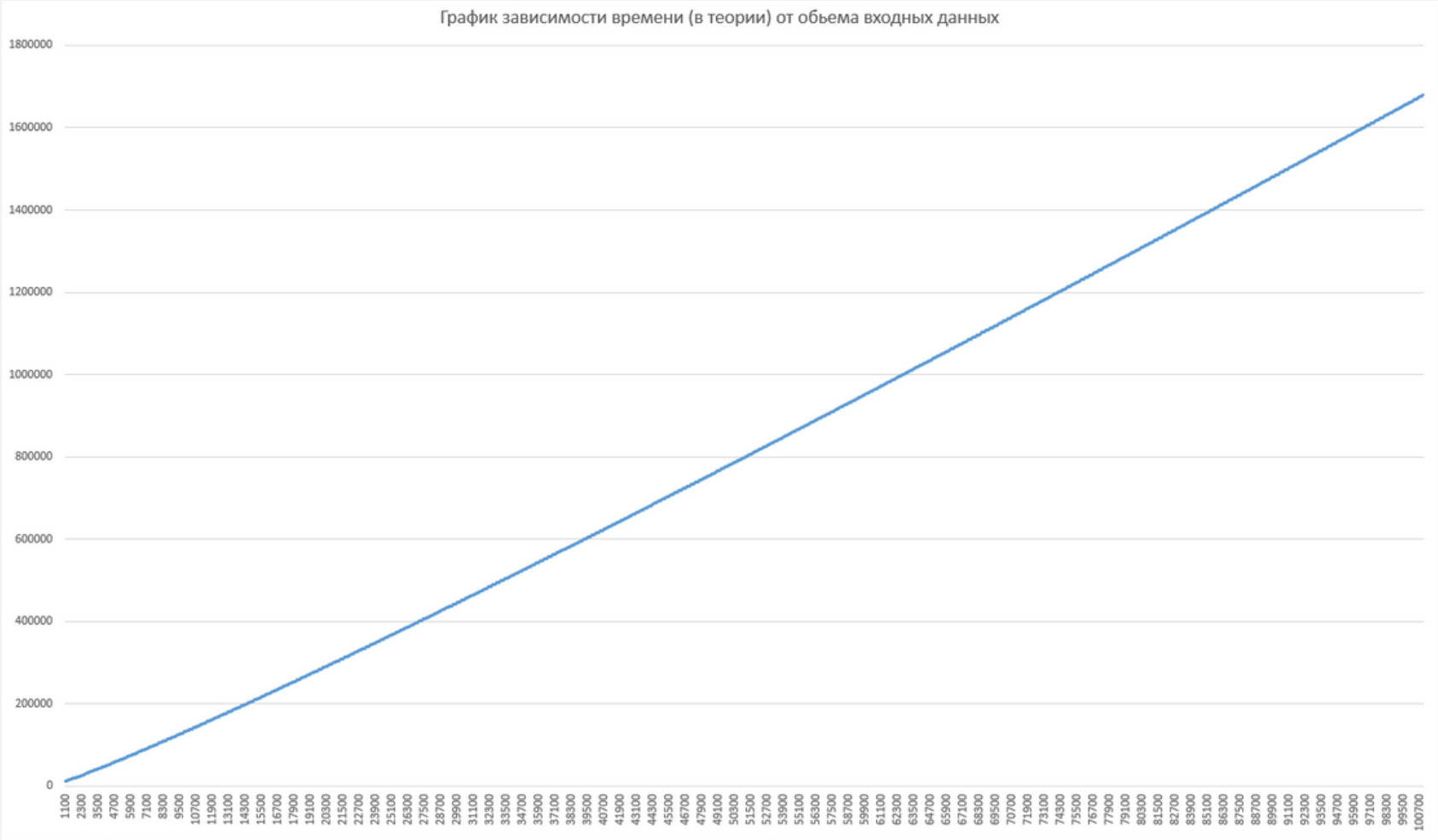
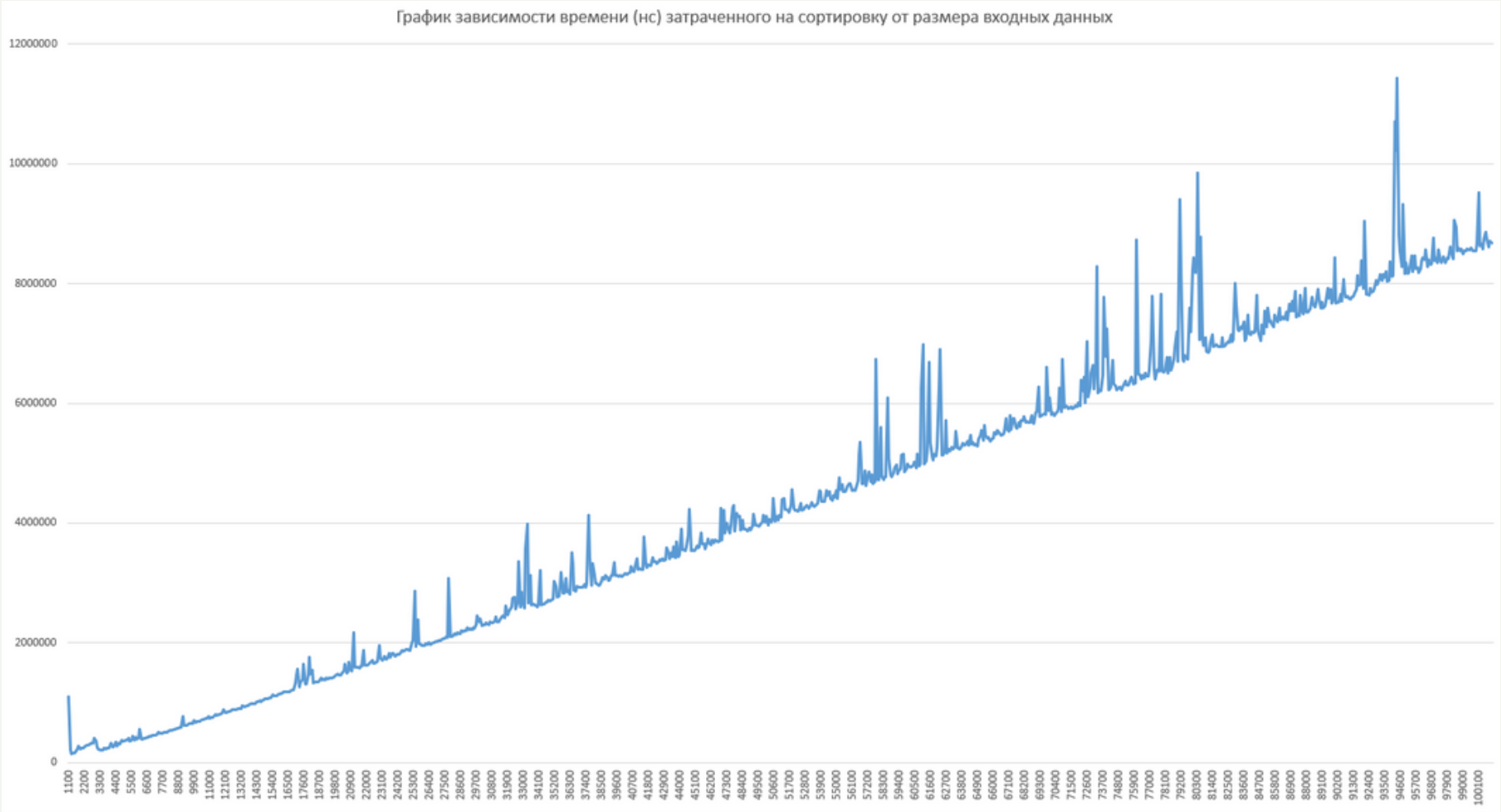


График зависимости времени (нс) от объема входных данных (на практике)

График зависимости времени от объема входных данных (в теории)

ПОЧЕМУ ГРАФИКИ НА ПРАКТИКЕ И В ТЕОРИИ РАЗЛИЧАЮТСЯ?

Теоретически скорость пирамидальной сортировки определяется количеством сравнений и перестановок, необходимых для сортировки входных данных, что обычно выражается как временная сложность $O(n \cdot \log(n))$. Эта временная сложность основана на предположении, что все операции занимают постоянное время, независимо от размера входных данных.

Однако на практике на скорость пирамидальной сортировки может влиять ряд факторов, не учитываемых при теоретическом анализе. Некоторые из этих факторов включают:

Эффекты кэширования. Если данные не хранятся в кеше, то время, необходимое для их извлечения из памяти, может быть намного больше, чем время, необходимое для выполнения сравнения или подкачки. В результате фактическое время работы пирамидальной сортировки может быть больше, чем предсказывает теоретический анализ.

Детали реализации. Фактическая реализация пирамидальной сортировки также может существенно повлиять на ее производительность. Например, способ построения кучи или порядок, в котором элементы меняются местами, могут повлиять на количество требуемых сравнений и свопов, что, в свою очередь, может повлиять на время выполнения.

Аппаратные ограничения. На производительность пирамидальной сортировки также может влиять аппаратное обеспечение, на котором она выполняется. Например, если процессор имеет небольшой кэш, то производительность пирамидальной сортировки может быть ниже ожидаемой.

Характеристики входных данных. Производительность пирамидальной сортировки также может зависеть от характеристик входных данных. Например, если данные уже частично отсортированы или содержат много дубликатов, то количество требуемых сравнений и перестановок может быть меньше, чем предсказано теоретическим анализом, что приводит к ускорению времени выполнения.

Следовательно, хотя теоретическая временная сложность `heapsort` составляет $O(n \cdot \log(n))$, на ее фактическую производительность может повлиять ряд практических факторов. На практике фактическое время выполнения пирамидальной сортировки может быть медленнее или быстрее, чем время в теории, в зависимости от реализации и характеристик входных данных и оборудования.

ПЛЮСЫ И МИНУСЫ HEAPSORT

Плюсы сортировки кучей:

- 1) Он имеет временную сложность $O(n \log n)$, что делает его эффективным для сортировки больших наборов данных.
- 2) Это алгоритм сортировки на месте, что означает, что для выполнения операции сортировки не требуется дополнительная память.
- 3) Он стабилен, то есть сохраняет порядок одинаковых элементов во входных данных.
- 5) Его относительно легко реализовать и понять.

Минусы сортировки кучей:

- 1) Он имеет большой постоянный коэффициент, что означает, что это может быть не самый быстрый алгоритм сортировки для небольших наборов данных.
- 2) Это не дружественный к кэшу алгоритм, а это означает, что он может страдать от низкой производительности при использовании с большими наборами данных, которые не помещаются в кэш-память.

ПРИМЕНИМОСТЬ HEAPSORT

Эта сортировка особенно полезна для сортировки больших наборов данных, где память является ограничением. Свойство сортировки на месте делает его популярным выбором для встраиваемых систем и других сред с ограниченным объемом памяти. Кроме того, его свойство стабильности делает его подходящим для приложений, где важно сохранять порядок одинаковых элементов, например, в запросах к базе данных. Однако низкая производительность кэш-памяти может ограничить его полезность в некоторых приложениях, где часто осуществляется доступ к данным из памяти. В целом, сортировка кучей — хороший выбор для сортировки больших наборов данных с ограниченным объемом памяти, где важна стабильность, а первоначальный порядок данных не имеет значения.