
АДРИАНОВ Н.М.
ИВАНОВ А.Б.

АЛГОРИТМЫ И СТРУКТУРЫ ДАННЫХ

TIMSORT

TIMSORT

- Tim Peters (2002)
<http://bugs.python.org/file4451/timsort.txt>
- Python 2.3, Java SE 7, Android 1.5
- В реальных ситуациях массивы данных часто содержат в себе упорядоченные подмассивы.
- Использует (модифицированную) процедуру слияния, суть алгоритма - в каком порядке выполнять слияние. Хочется сливать массивы близкого размера.
- Отсортированные куски храним в стеке, поддерживая некоторые условия (инвариант).
- Гибридный алгоритм: Небольшие подмассивы сортируются вставками.

TIMSORT: КОРРЕКТНОСТЬ АЛГОРИТМА

Февраль 2015: обнаружена ошибка (инвариант алгоритма не соблюдался, что приводило к переполнению стека)

S. de Gouw, J. Rot, F. de Boer, R. Bubel, R. Nähnle:

[Proving that Android's, Java's and Python's sorting algorithm is broken \(and showing how to fix it\)](#)

Предложено два варианта исправления

<https://bugs.openjdk.java.net/browse/JDK-8072909>

- 1) Исправить логику так, чтобы инвариант соблюдался
- 2) Увеличить размер стека

Python: исправлено по варианту 1

Java: исправлено по варианту 2

TIMSORT: КОРРЕКТНОСТЬ И СЛОЖНОСТЬ

Май 2018: Вариант исправления, выбранный в Java - некорректный (в работе 2015 года была ошибка).

N. Auger, C. Nicaud, C. Pivoteau, On the Worst-Case Complexity of TimSort. <https://arxiv.org/pdf/1805.08612.pdf>

<https://bugs.openjdk.java.net/browse/JDK-8203864>

В этой же работе доказано, что сложность алгоритма $O(n \log n)$ и даже $O(n + n \log \rho)$: Timsort – адаптивный алгоритм! Также ведет себя Natural MergeSort (Кнут, т.3).

<http://www-igm.univ-mlv.fr/~juge/slides/poster/ligm-2019.pdf>

TIMSORT: «RUN»

$$a_1 \leq a_2 \leq a_3 \leq \dots \leq a_r$$

$$a_1 > a_2 > a_3 > \dots > a_r$$

Выбираем число $\text{minrun} \in [32, 64)$ так, чтобы N/minrun оказалось степенью двойки или немного меньше.

Пример: $N = 2112$

$\text{minrun} = 32$ – плохой ($2112 = 66 \times 32$),

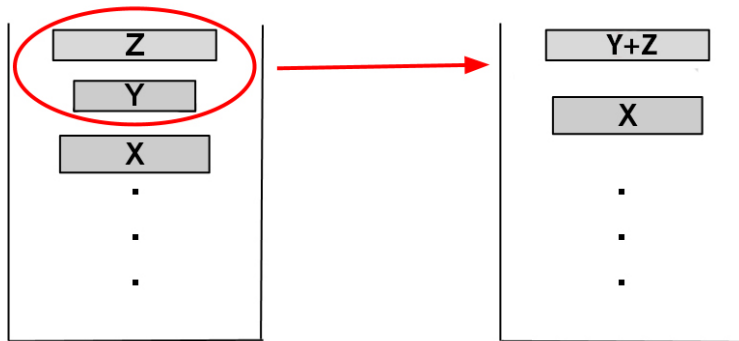
$\text{minrun} = 33$ – хороший ($2112 = 64 \times 33$).

```
int r = 0;    // Becomes 1 if any 1 bits are shifted off
while (n >= minrun) { r |= (n & 1); n >>= 1; }
return n + r;
```

TIMSORT: ОБЩАЯ СХЕМА

1. Определяем `minrun`, текущая позиция = 1
2. Находим `run`, начинающийся с текущей позиции
3. Если `run` убывающий – переворачиваем
4. Если `run` короче `minrun` – дополняем и сортируем вставками
5. Добавляем в стек (начало подмассива, длина)
6. Проверяем условия инварианта, если нарушены – выполняем слияние (см. след. слайд)
7. Если не дошли до конца массива – переходим к 2
8. Сливаем подмассивы, хранящиеся в стеке

TIMSORT: СТЕК



- $X > Y + Z$
- $Y > Z$

Если $X \leq Y + Z$ – сливаем Y с минимальным из X и Z .

Если $Y \leq Z$ – сливаем Y с Z .

TIMSORT: СТЕК

При соблюдении инварианта:

- Размеры кусков, хранящихся в стеке растут: $r_i < r_{i+1}$
- Растут быстрее чисел Фибоначчи: $r_i + r_{i+1} < r_{i+2}$

$$r_2 > r_1$$

$$r_3 > r_1 + r_2 > 2r_1$$

$$r_4 > r_2 + r_3 > 3r_1$$

$$r_5 > r_3 + r_4 > 5r_1$$

или (другое рассуждение):

$$r_{i+2} > r_i + r_{i+1} > 2r_i$$

т.е. скорость роста – экспоненциальная

- Следовательно, высота стека $O(\log n)$

TIMSORT: СЛИЯНИЕ

Слияние всегда выполняется для соседних подмассивов.

Создаем временный массив, размер которого равен размеру меньшего из сливаемых подмассивов. Меньший из подмассивов копируем во временный массив.

TIMSORT: РЕЖИМ «ГАЛОПА»

Still without loss of generality, assume A is the shorter run. In galloping mode, we first look for $A[0]$ in B . We do this via “galloping” comparing $A[0]$ in turn to $B[0]$, $B[1]$, $B[3]$, $B[7]$, ..., $B[2^{*j} - 1]$, ..., until finding the k such that $B[2^{*(k-1)} - 1] < A[0] \leq B[2^{*k} - 1]$. This takes at most roughly $\lg(B)$ comparisons, and, unlike a straight binary search, favors finding the right spot early in B (more on that later).

After finding such a k , the region of uncertainty is reduced to $2^{*(k-1)} - 1$ consecutive elements, and a straight binary search requires exactly $k-1$ additional comparisons to nail it. Then we copy all the B 's up to that point in one chunk, and then copy $A[0]$. Note that no matter where $A[0]$ belongs in B , the combination of galloping + binary search finds it in no more than about $2 * \lg(B)$ comparisons.

If we did a straight binary search, we could find it in no more than $\text{ceiling}(\lg(B+1))$ comparisons - but straight binary search takes that many comparisons no matter where $A[0]$ belongs. Straight binary search thus loses to galloping unless the run is quite long, and we simply can't guess whether it is in advance.

TIMSORT: РЕАЛИЗАЦИЯ В PYTHON 3.6.5

Algorithm 1: TIMSORT

Input: A sequence S to sort

Result: The sequence S is sorted into a single run

```
1 runs ← a run decomposition of  $S$ 
2  $\mathcal{R}$  ← an empty stack
3 while runs  $\neq \emptyset$  do // main loop of TIMSORT
4   remove a run  $r$  from runs and push  $r$  onto  $\mathcal{R}$ 
5   merge_collapse( $\mathcal{R}$ )
6 if  $h \neq 1$  then
7   merge_force_collapse( $\mathcal{R}$ )
```

Algorithm 2: merge_collapse

Input: A stack of runs \mathcal{R}

Result: The invariant is established

```
1 while  $h > 1$  do
2   if ( $h > 2$  and  $r_3 \leq r_2 + r_1$ ) or ( $h > 3$  and  $r_4 \leq r_3 + r_2$ )
3     then
4       if  $r_3 < r_1$  then merge( $R_2, R_3$ )
5       else merge( $R_1, R_2$ )
6     else if  $r_2 \leq r_1$  then merge( $R_1, R_2$ )
7     else break
```

TIMSORT: ПЕРЕФОРМУЛИРОВКА

Algorithm 3: TIMSORT: translation of Alg. 1 and Alg. 2

Input: A sequence to S to sort

Result: The sequence S is sorted into a single run

Note: h = height of the stack \mathcal{R}

R_i = i^{th} top-most run in the stack \mathcal{R}

r_i = the size of R_i .

```
1 runs  $\leftarrow$  the run decomposition of  $S$ 
2  $\mathcal{R} \leftarrow$  an empty stack
3 while runs  $\neq \emptyset$  do // main loop of TIMSORT
4     remove a run  $r$  from runs and push  $r$  onto  $\mathcal{R}$  // #1
5     while true do
6         if  $h \geq 3$  and  $r_1 > r_3$  then merge( $R_2, R_3$ ) // #2
7         else if  $h \geq 2$  and  $r_1 \geq r_2$  then merge( $R_1, R_2$ ) // #3
8         else if  $h \geq 3$  and  $r_1 + r_2 \geq r_3$  then merge( $R_1, R_2$ ) // #4
9         else if  $h \geq 4$  and  $r_2 + r_3 \geq r_4$  then merge( $R_1, R_2$ ) // #5
10        else break
11 while  $h \neq 1$  do merge( $R_1, R_2$ )
```

TIMSORT: СЛОЖНОСТЬ

Теорема. Пусть \mathcal{C} – множество массивов длины n , состоящих из ρ runs длин r_1, \dots, r_ρ . Введем функцию энтропии Шеннона

$$H(p_1, \dots, p_\rho) = - \sum_{i=1}^{\rho} p_i \log p_i$$

и пусть $\mathcal{H} = H(r_1/n, \dots, r_\rho/n)$. Время работы алгоритма Timsort на массивах из \mathcal{C} равно $O(n + n\mathcal{H})$.

Следствие. $\mathcal{H} \leq \log \rho$, т.е. Время работы алгоритма Timsort на массивах из \mathcal{C} равно $O(n + n\rho)$, и следовательно, $O(n \log n)$.

Предложение. Для любого детерменированного алгоритма сортировки *сравнением* найдется массива из множества \mathcal{C} , на котором потребуется, по крайней мере, $n\mathcal{H} - 3n$ сравнений.

Доказательство следствия: надо показать, что $\mathcal{H} \leq \log \rho$.

Функция $f : x \rightarrow -x \ln(x)$ вогнута на $\mathbb{R}_{>0}$, т.к. $f''(x) = -1/x < 0$.

Поскольку $p_i > 0$ и $p_1 + \dots + p_\rho = 1$, то

$$H(p_1, \dots, p_\rho) = \sum_{i=1}^{\rho} f(p_i) / \ln(2) \leq \rho f(1/\rho) / \ln(2) = \log \rho.$$