

Внешняя сортировка

- Принято называть "**внешней**" сортировкой сортировку последовательных файлов, располагающихся во внешней памяти и слишком больших, чтобы можно было целиком переместить их в основную память и применить один из методов внутренней сортировки.

Внешняя сортировка: прямое слияние

- Имеется последовательный файл **A**, состоящий из записей **a1, a2, ..., an**
- Для сортировки используются два вспомогательных файла **B** и **C** (размер каждого из них - **n/2**).
- Сортировка состоит из последовательности шагов, в каждом из которых выполняется распределение состояния файла **A** в файлы **B** и **C**, а затем слияние файлов **B** и **C** в файл **A**.

Внешняя сортировка: прямое слияние

- **Шаг1.** Последовательно читается файл **A**, и записи **$a_1, a_3, \dots, a_{(n-1)}$** пишутся в файл **B**, а записи **a_2, a_4, \dots, a_n** - в файл **C** (начальное распределение)
- **Шаг2.** Начальное слияние производится над парами **$(a_1, a_2), (a_3, a_4), \dots, (a_{(n-1)}, a_n)$** ,
- **Шаг3.** Последовательно читается файл **A**, и в файл **B** записываются последовательные пары с нечетными номерами, а в файл **C** - с четными. (пример)
- **Шаг4.** Слияние **B** и **C** образуются и пишутся в файл **A** упорядоченные четверки записей...

Внешняя сортировка: естественное слияние

- Серией называется подпоследовательность записей

$$a_i, a(i+1), \dots, a_j$$

такая, что

$$a_k \leq a(k+1) \text{ для всех } i \leq k < j,$$

$$a_i < a(i-1)$$

$$\text{и } a_j > a(j+1).$$

Метод естественного слияния основывается на распознавании серий при распределении и их использовании при последующем слиянии.

Внешняя сортировка: естественное слияние

- При распределении распознается первая серия записей и переписывается в файл **B**, вторая - в файл **C** и т.д.
- При слиянии первая серия записей файла **B** сливается с первой серией файла **C**, вторая серия **B** со второй серией **C** и т.д.
- Если просмотр одного файла заканчивается раньше, чем просмотр другого (по причине разного числа серий), то остаток недопросмотренного файла целиком копируется в конец файла **A**.
- Процесс завершается, когда в файле **A** остается только одна серия. (демо)

Сбалансированное многопутевое слияние

- распределение серий исходного файла по m вспомогательным файлам

V_1, V_2, \dots, V_m

и их слияние в m вспомогательных файлов

$C_1, C_2, \dots, C_m.$

- на следующем шаге производится слияние файлов C_1, C_2, \dots, C_m в файлы V_1, V_2, \dots, V_m и т.д., пока в V_1 или C_1 не образуется одна серия (демо)

Многофазная сортировка

- Из имеющихся m вспомогательных файлов $(m-1)$ файл служит для ввода сливаемых последовательностей, а один - для вывода образуемых серий.
- Как только один из файлов ввода становится пустым, его начинают использовать для вывода серий, получаемых при слиянии серий нового набора $(m-1)$ файлов.
- Таким образом, имеется первый шаг, при котором серии исходного файла распределяются по $m-1$ вспомогательному файлу, а затем выполняется многопутевое слияние серий из $(m-1)$ файла, пока в одном из них не образуется одна серия.