



Politechnika Wroclawska

Wydział Elektroniki, Fotoniki i Mikrosystemów

PAMSI

Sprawozdanie nr 2

Sortowanie

Prowadzący:

dr hab. inż. Andrzej Rusiecki

Wykonał:

Jakub Kusz

Wrocław, 24 maja 2022r.

Spis treści

1	Cel projektu	2
2	Zadanie do wykonania	2
3	Filtrowanie danych	2
3.1	Przewidywana złożoność obliczeniowa	2
3.2	Wyznaczona złożoność obliczeniowa	3
4	Quicksort	3
4.1	Złożoność obliczenia	3
4.1.1	Przypadek optymistyczny	3
4.1.2	Przypadek przeciętny	4
4.1.3	Przypadek pesymistyczny	4
4.2	Implementacja	4
4.3	Testy złożoności obliczeniowej	5
5	Mergesort	6
5.1	Złożoność obliczeniowa	7
5.2	Implementacja	7
5.2.1	Scalanie	7
5.2.2	Sortowanie	8
5.3	Testy złożoności obliczeniowej	8
6	Introsort	9
6.1	Złożoność obliczeniowa	9
6.1.1	Przypadek optymistyczny	9
6.1.2	Przypadek przeciętny	9
6.1.3	Przypadek pesymistyczny	9
6.2	Implementacja	10
6.2.1	Podział	10
6.2.2	Sortowanie	10
6.3	Testy złożoności obliczeniowej	10
7	Dyskusja	12
7.1	Czas trwania filtracji	12
7.2	Sortowanie	12
7.2.1	Quicksort	12
7.2.2	Mergesort	12
7.2.3	Introsort	12
8	Podsumowanie	12
8.1	Czas trwania	12
8.2	Średnia arytmetyczna i mediana	13
8.3	Wnioski	14

1. Cel projektu

Celem projektu jest zapoznanie się z różnymi algorytmami sortującymi i ich złożonością obliczeniową zależną od różnych zestawów danych.

2. Zadanie do wykonania

Dla danych w pliku projekt2_dane.csv należy wykonać eksperymenty z sortowaniem danych względem rankingu filmów. Załączony plik jest okrojona baza filmów „IMDb Largest Review Dataset” z kaggle.com. Plik zawiera tylko tytuł oraz ranking. Proszę o wykonanie następujących zadań:

1. Przefiltrowanie danych i usunięcie pustych wpisów w polu ranking (jeśli występują). Proszę zmierzyć i podać w sprawozdaniu czas przeszukiwania. Czy był on zgodny z oczekiwaną złożonością przeszukiwania dla wybranej struktury danych?
2. Przygotować strukturę danych zawierającą odpowiednio: 10 000, 100 000, 500 000, 1 000 000, maksymalną ilość danych z pliku.
3. Przeprowadzić analizę efektywności sortowania na danych z §2 z wykorzystaniem zaimplementowanych algorytmów.
4. Dodatkowo dla każdego zestawu danych proszę podać w tabeli czas sortowania, średnią wartość oraz medianę rankingu.

Zostały przeze mnie wybrane 4 algorytmy sortujące, które poddałem testom:

1. **Quicksort**,
2. **Mergesort**,
3. **Introsort**.

Kod źródłowy programu znajduje się **Pod tym adresem**

3. Filtrowanie danych

Pierwszym krokiem było przefiltrowanie danych z pliku .csv zawierającego tytuł i rankingi. Należało usunąć wpisy, które nie zawierały rankingu. Filtrowanie odbywało się jednocześnie z pobieraniem poszczególnych linii z pliku .csv - jeśli dana linia nie zawierała pola z rankingiem to nie zostawała zapisana do struktury danych.

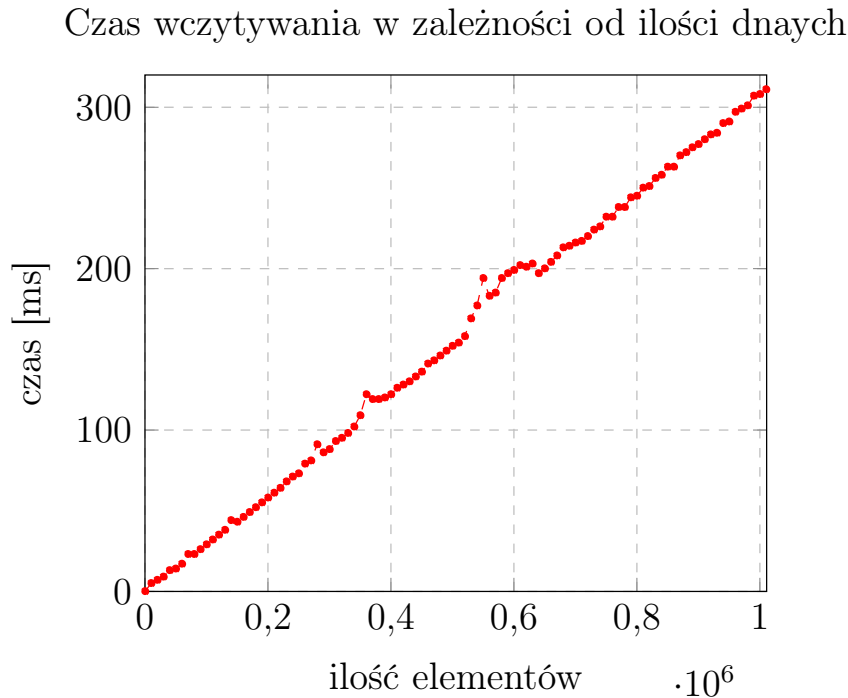
3.1. Przewidywana złożoność obliczeniowa

Przewiduje się złożoność obliczeniową liniową. Złożoność dodania na koniec wykorzystanego w zadaniu kontenera `std::vector` jest stała w czasie, sprawdzenie czy dana linia zawiera pole z rankingiem również, więc czas przefiltrowania danych zależny jest tylko od ich ilości, więc powinien być liniowy. W notacji dużego O:

$$O(n)$$

3.2. Wyznaczona złożoność obliczeniowa

W celu wyznaczenia złożoności obliczeniowej został przeprowadzony test, polegający na mierzeniu czasu wczytywania liczby kolejnych krotności 1000 linii z pliku .csv.



Wykres. 1: Przedstawiający czas filtrowania kolejnych ilości danych

4. Quicksort

Algorytm wykorzystuje technikę "dziel i zwyciężaj". Według ustalonego schematu wybierany jest jeden element w sortowanej tablicy, który będziemy nazywać pivot. Pivot może być elementem środkowym, pierwszym, ostatnim, losowym lub wybranym według jakiegoś innego schematu dostosowanego do zbioru danych. Następnie ustawiamy elementy nie większe na lewo tej wartości, natomiast nie mniejsze na prawo. W ten sposób powstaną nam dwie części tablicy (niekoniecznie równe), gdzie w pierwszej części znajdują się elementy nie większe od drugiej. Następnie każdą z tych podtablic sortujemy osobno według tego samego schematu.

4.1. Złożoność obliczenia

W zależności od rozkładu danych i elementu pivot określa się następujące złożoności obliczeniowe:

4.1.1. Przypadek optymistyczny

W przypadku optymistycznym, jeśli mamy szczęście za każdym razem wybrać medianę z sortowanego fragmentu tablicy, to liczba porównań niezbędnych do uporządkowania n -elementowej tablicy opisana jest rekurencyjnym wzorem:

$$O(n \cdot \log_2 n)$$

4.1.2. Przypadek przeciętny

W przypadku przeciętnym, to jest dla równomiernego rozkładu prawdopodobieństwa wyboru elementu z tablicy:

$$O(1,39 \cdot n \cdot \log_2 n)$$

4.1.3. Przypadek pesymistyczny

W przypadku pesymistycznym, jeśli zawsze wybierzemy element najmniejszy (albo największy) w sortowanym fragmencie tablicy, to:

$$O(n^2)$$

4.2. Implementacja

Poniższy listing przedstawia rekurencyjną implementację algorytmu Quicksort.

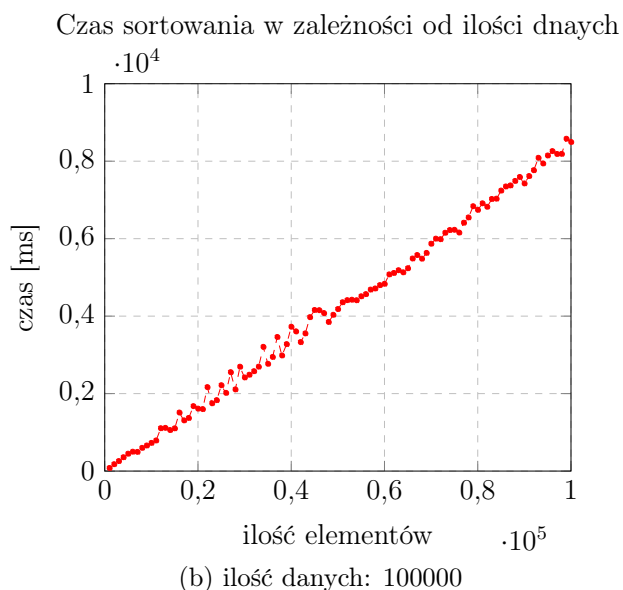
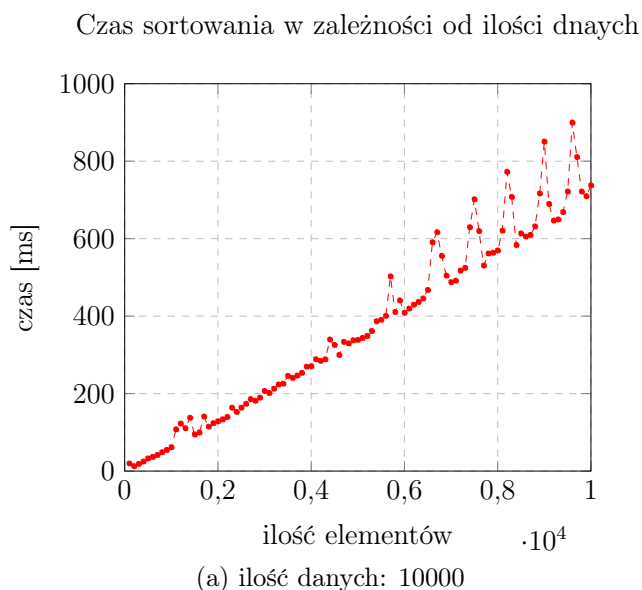
```
1 void quick_sort(double *tab, int left, int right){
2     if(right <= left)
3         return;
4
5     int i = left - 1;
6     int j = right + 1;
7     int pivot = tab[(left+right)/2];
8
9     while(1){
10
11         while(pivot > tab[++i]);
12         while(pivot < tab[--j]);
13
14
15         if( i <= j)
16             std :: swap(tab[i], tab[j]);
17         else
18             break;
19     }
20
21     if(j > left)
22         quick_sort(tab, left, j);
23     if(i < right)
24         quick_sort(tab, i, right);
25 }
```

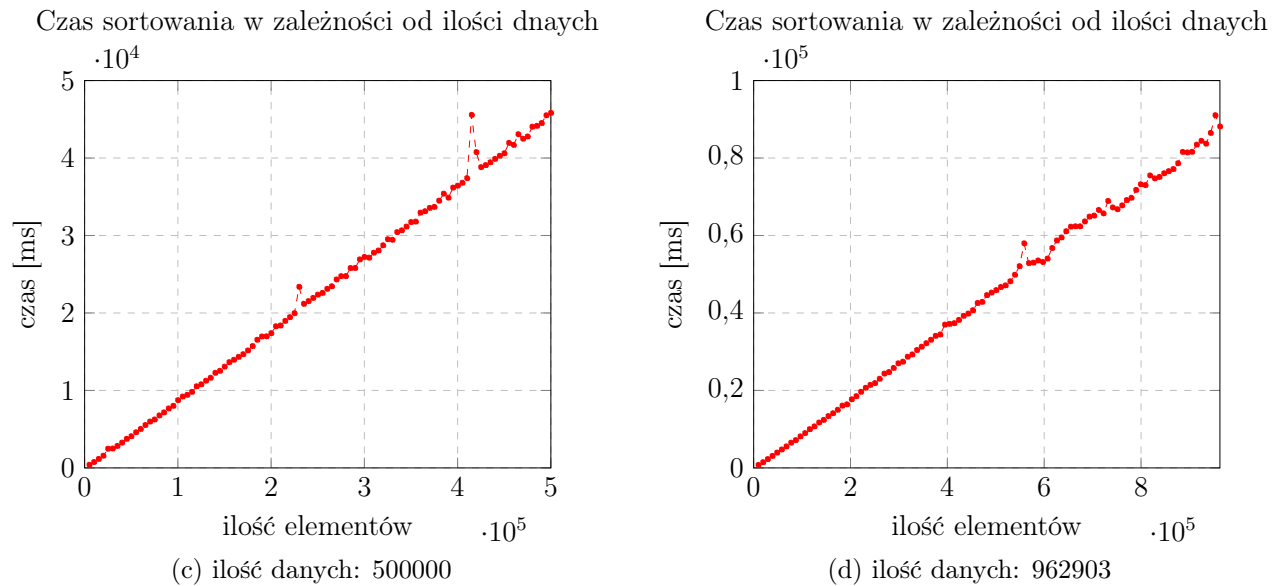
4.3. Testy złożoności obliczeniowej

W celu eksperymentalnego wyznaczenia złożoności obliczeniowej algorytmu Quicksort zostały przeprowadzone testy szybkości dla czterech różniących się ilością danych pakietów:

- 10000,
- 100000,
- 500000,
- 962903.

Dla każdego zestawu danych przeprowadzono 100 pomiarów czasu. Najmniejszą paczką dla poszczególnej ilości danych jest $\frac{x}{100}$ gdzie x to ilość danych. Każdy kolejny pomiar wykonywano dla $n \cdot \frac{x}{100}$ ilości w paczce, gdzie $n \in 1, 2, 3 \dots 100$. Na przykład dla **10000** poszczególne paczki to 100, 200, 300... a dla **500000** to 5000, 10000, 15000... .





Wykres. 2: Wykresy przedstawiające czas sortowania określonej ilości danych

Na Rys. 2 można zauważyć, iż złożoność obliczeniowa prezentuje się jako f. liniowa, co nie jest dużym zaskoczeniem, ponieważ złożoność zarówno optymistyczna(4.1.1) jak i przeciętna(4.1.2) w dużej dziedzinie kształtem również przypominają f. linowe, aczkolwiek nimi nie są.

5. Mergesort

Ideą działania algorytmu jest dzielenie zbioru danych na mniejsze zbiory, aż do uzyskania n zbiorów jednoelementowych, które same z siebie są posortowane, następnie zbiory te są łączone w coraz większe zbiory posortowane, aż do uzyskania jednego, posortowanego zbioru n -elementowego. Etap dzielenia nie jest skomplikowany, dzielenie następuje bez sprawdzania jakichkolwiek warunków. Z kolei łączenie zbiorów posortowanych wymaga odpowiedniego wybierania poszczególnych elementów z łączonych zbiorów z uwzględnieniem faktu, że wielkość zbioru nie musi być równa (parzysta i nieparzysta ilość elementów), oraz tego, iż wybieranie elementów z poszczególnych zbiorów nie musi następować naprzemiennie, przez co jeden zbiór może osiągać swój koniec wcześniej niż drugi. Robi się to w następujący sposób. Kopiujemy zawartość zbioru głównego do struktury pomocniczej. Następnie, operując wyłącznie na kopii, ustawiamy wskaźniki na początki kolejnych zbiorów i porównujemy wskazywane wartości. Mniejszą wartość wpisujemy do zbioru głównego i przesuwamy odpowiedni wskaźnik o 1 i czynności powtarzamy, aż do momentu, gdy jeden ze wskaźników osiągnie koniec zbioru. Wówczas mamy do rozpatrzenia dwa przypadki, gdy zbiór 1 osiągnął koniec i gdy zbiór 2 osiągnął koniec. W przypadku pierwszym nie będzie problemu, elementy w zbiorze głównym są już posortowane i ułożone na właściwych miejscach. W przypadku drugim trzeba skopiować pozostałe elementy zbioru pierwszego pokolei na koniec. Po zakończeniu wszystkich operacji otrzymujemy posortowany zbiór główny.

5.1. Złożoność obliczeniowa

Złożoność obliczeniowa jest niezależna od zestawu danych:

$$O(n \cdot \log_2 n)$$

5.2. Implementacja

Poniższe listingi przedstawiają funkcję funkcję scalającą i funkcję sortującą:

5.2.1. Scalanie

```

1 void merge(double array [], int const left , int const mid, int const right ){
2     int const sub_array_one = mid - left + 1;
3     int const sub_array_two = right - mid;
4     double *left_array = new double[sub_array_one];
5     double *right_array = new double[sub_array_two];
6     for ( int i = 0; i < sub_array_one; i++)
7         left_array [ i ] = array[ left + i ];
8     for ( int j = 0; j < sub_array_two; j++)
9         right_array[ j ] = array[mid + 1 + j];
10    int index_of_sub_array_one = 0;
11    int index_of_sub_array_two = 0;
12    int index_of_merged_array = left;
13    while (index_of_sub_array_one < sub_array_one && index_of_sub_array_two < sub_array_two) {
14        if (left_array [index_of_sub_array_one] <= right_array[index_of_sub_array_two]) {
15            array[index_of_merged_array] = left_array[index_of_sub_array_one];
16            index_of_sub_array_one++;
17        }
18        else {
19            array[index_of_merged_array] = right_array[index_of_sub_array_two];
20            index_of_sub_array_two++;
21        }
22        index_of_merged_array++;
23    }
24    while (index_of_sub_array_one < sub_array_one) {
25        array[index_of_merged_array] = left_array[index_of_sub_array_one];
26        index_of_sub_array_one++;
27        index_of_merged_array++;
28    }
29    while (index_of_sub_array_two < sub_array_two) {
30        array[index_of_merged_array] = right_array[index_of_sub_array_two];
31        index_of_sub_array_two++;
32        index_of_merged_array++;
33    }
34 }
35
```


5.2.2. Sortowanie

```

1
2     void merge_sort(double *arr, int left , int right ){
3         if ( left >= right)
4             return;
5
6         auto mid = left + (right - left ) / 2;
7         merge_sort(arr, left , mid);
8         merge_sort(arr, mid + 1, right);
9         merge(arr, left , mid, right );
10    }
11

```

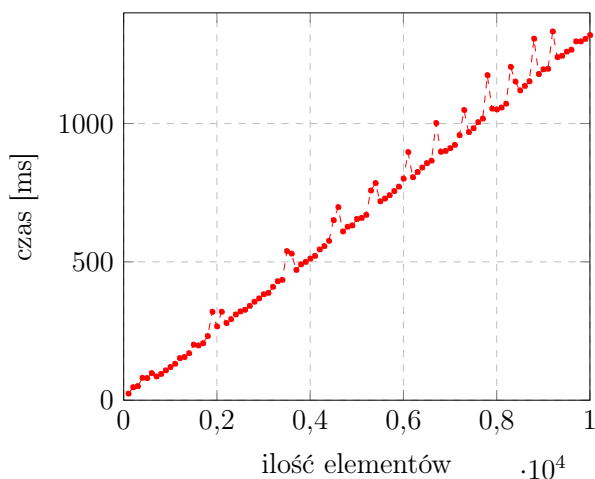
5.3. Testy złożoności obliczeniowej

W celu eksperymentalnego wyznaczenia złożoności obliczeniowej algorytmu Mergesort zostały przeprowadzone testy szybkości dla czterech różniących się ilością danych pakietów:

- 10000,
- 100000,
- 500000,
- 962903.

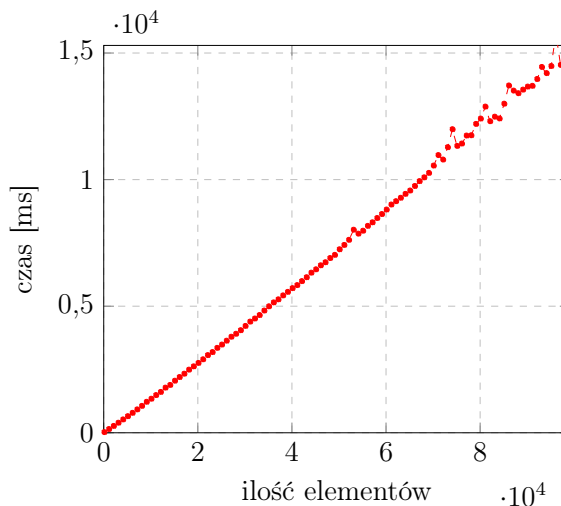
Dla każdego zestawu danych przeprowadzono 100 pomiarów czasu. Najmniejszą paczką dla poszczególnej ilości danych jest $\frac{x}{100}$ gdzie x to ilość danych. Każdy kolejny pomiar wykonywano dla $n \cdot \frac{x}{100}$ ilości w paczce, gdzie $n \in 1, 2, 3 \dots 100$. Na przykład dla **10000** poszczególne paczki to 100, 200, 300... a dla **500000** to 5000, 10000, 15000... .

Czas sortowania w zależności od ilości danych

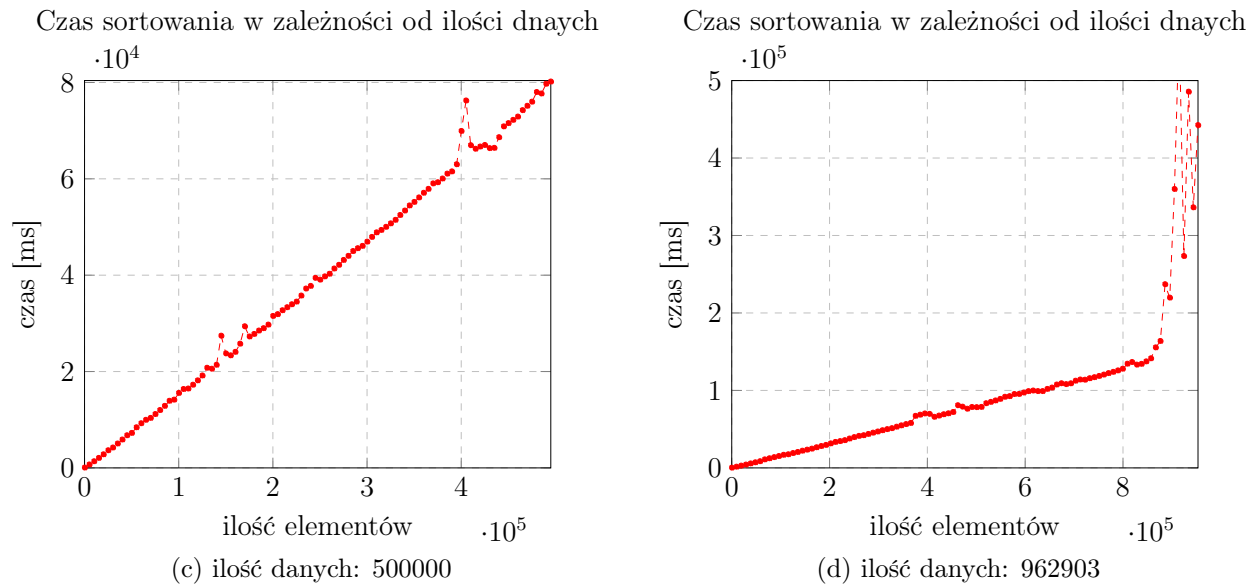


(a) ilość danych: 10000

Czas sortowania w zależności od ilości danych



(b) ilość danych: 100000



Wykres. 3: Wykresy przedstawiające czas sortowania określonej ilości danych

6. Introsort

Sortowanie introspektywne to odmiana sortowania hybrydowego, w której wyeliminowany został problem złożoności $O(n^2)$ występującej w najgorszym przypadku algorytmu sortowania szybkiego. Algorytm łączy w sobie 3 rodzaje sortowań, które wywoływane są dla różnych przypadków. Używa on sortowania przez wstawianie (insertion sort), sortowania przez kopcowanie (heap sort) i sortowania szybkiego (quick sort).

6.1. Złożoność obliczeniowa

6.1.1. Przypadek optymistyczny

$$O(n \cdot \log_2 n)$$

6.1.2. Przypadek przeciętny

$$O(1.39 \cdot n \cdot \log_2 n)$$

6.1.3. Przypadek pesymistyczny

$$O(2 \cdot n \cdot \log_2 n)$$

6.2. Implementacja

Poniższe listingi przedstawiają poszczególne funkcje konieczne do wykonania Introsort

6.2.1. Podział

```

1      int partition (double *data, int left , int right ) {
2          int pivot = data[right ];
3          int temp;
4          int i = left ;
5          for ( int j = left ; j < right; ++j){
6              if (data[j ] <= pivot){
7                  temp = data[j];
8                  data[j ] = data[i ];
9                  data[i ] = temp;
10                 i++;
11             }
12         }
13
14         data[right ] = data[i ];
15         data[i ] = pivot;
16
17         return i;
18     }
19

```

6.2.2. Sortowanie

```

1      void intro_sort(double *arr , int size ) {
2          int partitionSize = partition(arr , 0, size - 1);
3          if ( partitionSize < 16){
4              insertion_sort(arr , size );
5          }
6          else if ( partitionSize >(2 * std :: log( size ))){
7              heap_sort(arr, size );
8          }
9          else{
10             quick_sort(arr, 0, size -1);
11         }
12     }
13

```

6.3. Testy złożoności obliczeniowej

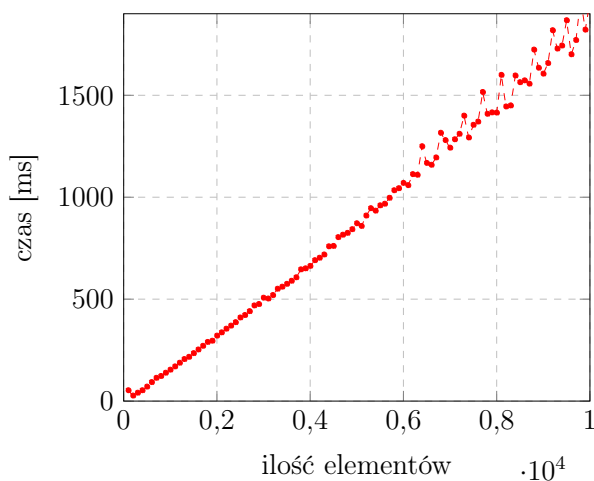
W celu eksperymentalnego wyznaczenia złożoności obliczeniowej algorytmu Introsort zostały przeprowadzone testy szybkości dla czterech różniących się ilością danych pakietów:

- 10000,

- 100000,
- 500000,
- 962903.

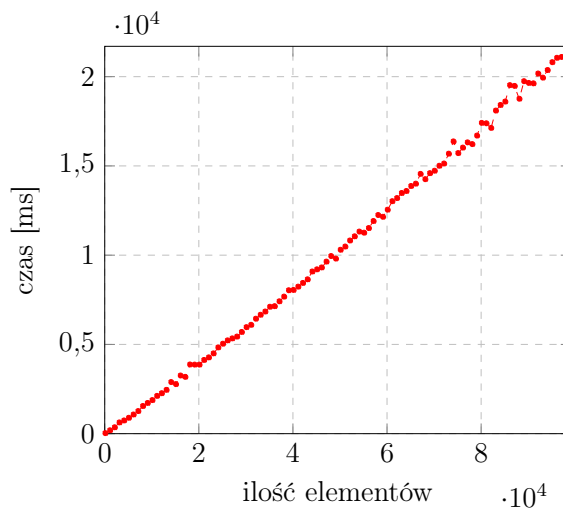
Dla każdego zestawu danych przeprowadzono 100 pomiarów czasu. Najmniejszą paczką dla poszczególnej ilości danych jest $\frac{x}{100}$ gdzie x to ilość danych. Każdy kolejny pomiar wykonywano dla $n \cdot \frac{x}{100}$ ilości w paczce, gdzie $n \in 1, 2, 3 \dots 100$. Na przykład dla **10000** poszczególne paczki to 100, 200, 300... a dla **500000** to 5000, 10000, 15000... .

Czas sortowania w zależności od ilości danych



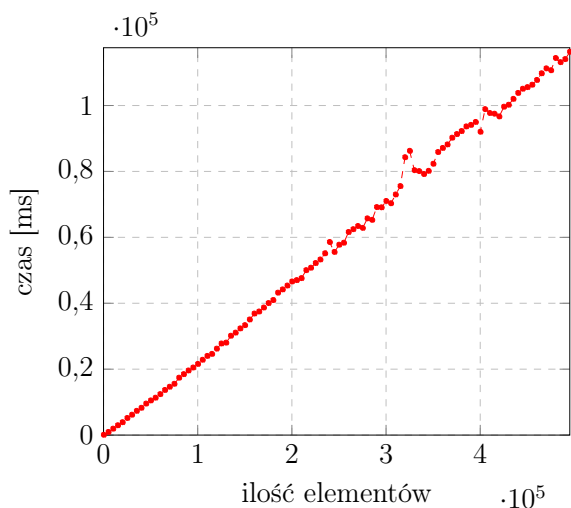
(a) ilość danych: 10000

Czas sortowania w zależności od ilości danych



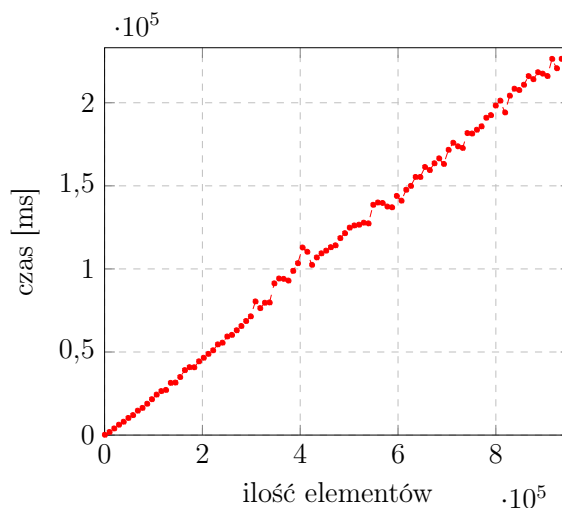
(b) ilość danych: 100000

Czas sortowania w zależności od ilości danych



(c) ilość danych: 500000

Czas sortowania w zależności od ilości danych



(d) ilość danych: 962903

Wykres. 4: Wykresy przedstawiające czas sortowania określonej ilości danych

7. Dyskusja

7.1. Czas trwania filtracji

Wyznaczona doświadczalnie złożoność obliczeniowa filtrowania danych (Wykres 1) jest zgodna ze złożonością prognozowaną (pkt. 3).

Filtrowanie danych trwało \sim : **311 ms**

7.2. Sortowanie

7.2.1. Quicksort

Analizując Wykres 2 można dojść do wniosku, iż złożoność obliczeniowa algorytmu Quicksort jest liniowa, co nie zgadza się z optymistyczną przeciętną jak i pesymistyczną złożonością. (pkt. 4). Można jednak taki rezultat uznać za właściwy, ponieważ przebieg funkcji $n \cdot \log_2 n$ jest kształtem zbliżony do funkcji liniowej.

Sortowanie największego zestawu danych (962903 elementów) trwało \sim : **88,101 ms**

7.2.2. Mergesort

Analizując Wykres 3 można dojść do wniosku, iż złożoność obliczeniowa algorytmu Mergesort jest liniowa, co nie zgadza się z teoretyczną złożonością. (pkt. 5). Można jednak taki rezultat uznać za właściwy, ponieważ przebieg funkcji $n \cdot \log_2 n$ jest kształtem zbliżony do funkcji liniowej. Dla największego zestawu powyżej $9 \cdot 10^5$ ilości danych czas wykonywania sortowania gwałtownie wzrasta, przez co finalny czas sortowania znacznie się wydłużył.

Sortowanie największego zestawu danych (962903 elementów) trwało \sim : **442,203 ms**

7.2.3. Introsort

Analizując Wykres 4 można dojść do wniosku, iż złożoność obliczeniowa algorytmu Introsort jest liniowa, co nie zgadza się z optymistyczną, przeciętną jak i pesymistyczną złożonością. (pkt. 6). Można jednak taki rezultat uznać za właściwy, ponieważ przebieg funkcji $n \cdot \log_2 n$ jest kształtem zbliżony do funkcji liniowej.

Sortowanie największego zestawu danych (962903 elementów) trwało \sim : **232,851 ms**

8. Podsumowanie

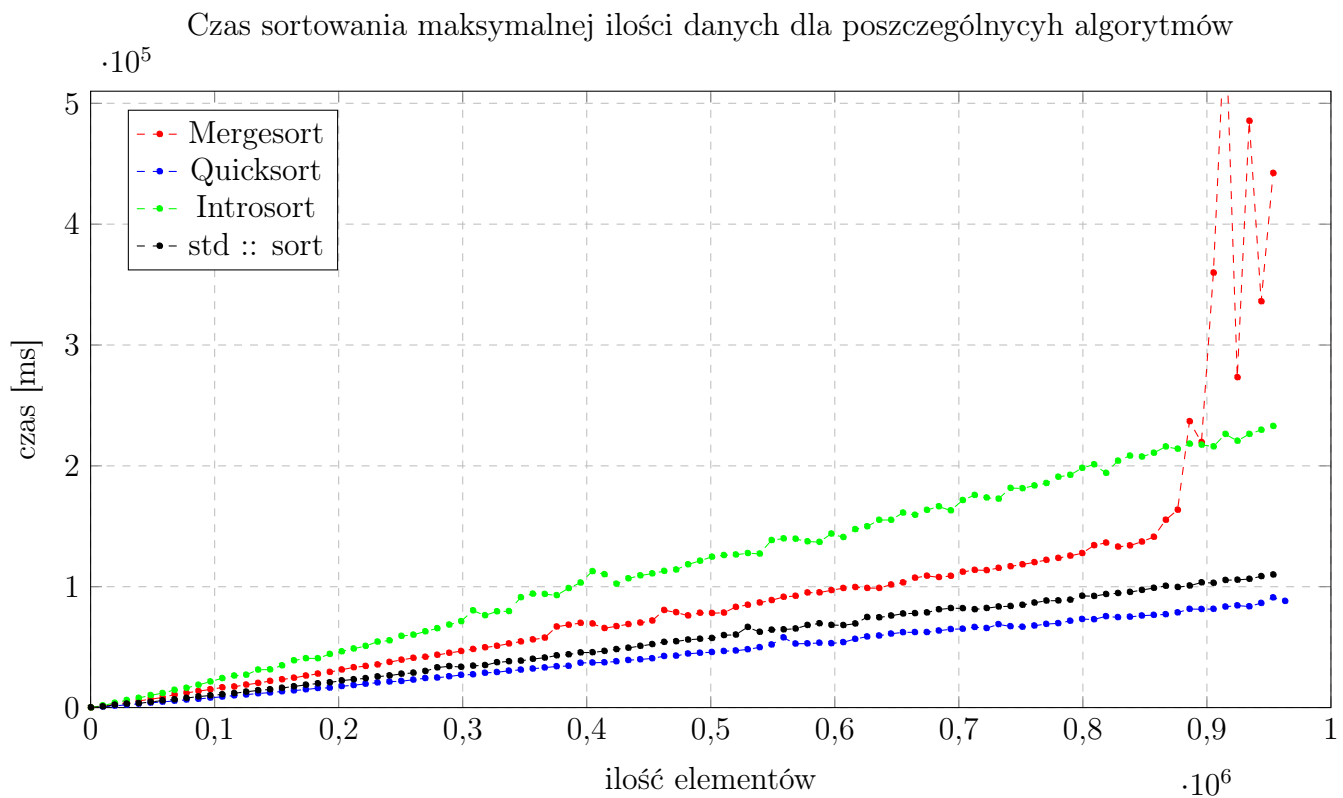
8.1. Czas trwania

Poniższa tabela przedstawia czasy trwania sortowania poszczególnych algorytmów:

Tablica 1: Zawierająca czasy trwania sortowania poszczególnych algorytmów

algorytm	czas trwania [ms]
Quicksort	88,101
Mergesort	442,203
Introsort	232,851
filtrowanie	311

Na poniższym wykresie zostały zaprezentowane charakterystyki złożoności czasowej wszystkich omawianych algorytmów, dodatkowo została zaprezentowana implementacja algorytmu sortującego z biblioteki STL C++ `std::sort()`, w celu porównania szybkości gotowego rozwiązania dostępnego dla programisty, z jego własną implementacją algorytmu sortującego.



Wykres. 6: Przedstawiający czas sortowania maksymalnej ilości danych dla poszczególnych algorytmów

Uwaga. Algorytm `std::sort()` został przedstawiony tylko w celu porównania i zaprezentowania go na wykresie z z napisanymi w ramach zadania algorytmami, nie został on uwzględniony w podsumowaniu i wnioskach, pojawił się jako rodzaj ciekawostki.

8.2. Średnia arytmetyczna i mediana

Dla każdego zestawu danych została obliczona średnia arytmetyczna i mediana rankingu. Zostały one zaprezentowane w poniższej tabeli.

Tablica 2: Zawierająca wartości średniej arytmetycznej i mediany dla poszczególnych zestawów danych

ilość elementów	średnia arytmetyczna	mediana
10000	5.4603	5
100000	6.08993	7
500000	6.66572	7
962903	6.63661	7

8.3. Wnioski

- najszybciej działającym algorytmem sortowania okazał się Quicksort (tab. 1),
- najwolniej działającym algorytmem sortowania okazał się Introsort (tab. 1),
- Mergesort dla danych w ilości większej od ok. $9 \cdot 10^5$ gwałtownie zwalnia przez co dla dużej ilości danych okazuje się najwolniejszy, nie udało mi się ustalić tego przyczyny (Wykres 6),
- w przypadku algorytmów sortujących złożoności czasowe wyznaczone doświadczalnie nie są zgodne z teoretycznymi, aczkolwiek na tyle podobne, iż wyniki doświadczeń można uznać za satysfakcjonujące.
- mimo, iż najszybciej działającym algorytmem jest Quicksort, może on dla niekorzystnego wyznaczenia pivotu okazać się algorytmem o złożoności kwadratowej (pkt. 4), w przeciwieństwie do algorytmów Introsort i Mergesort (pkt. 6,5), które gwarantują taką samą złożoność niezależnie od zestawu danych, lecz działają wolniej.