

WYDZIAŁ INFORMATYKI I TELEKOMUNIKACJI
POLITECHNIKA WROCŁAWSKA

ANALIZA EKSPERYMENTALNA NOWYCH ALGORYTMÓW SORTOWANIA W MIEJSCU

DAMIAN BALIŃSKI
NR INDEKSU: 250332

Praca inżynierska napisana
pod kierunkiem
dr inż. Zbigniewa Gołębiewskiego



Politechnika
Wrocławska

WROCŁAW 2021

Spis treści

1	Wstęp	1
1.1	Cel pracy	1
1.2	Zakres pracy	1
1.3	Przegląd literatury	1
1.4	Zawartości pracy	1
2	Analiza problemu	3
2.1	Model matematyczny przypadku średniego	3
2.2	Model matematyczny przypadku pesymistycznego	3
2.3	Założenia	3
2.3.1	Założenia odnośnie testowanych parametrów	3
2.3.2	Założenia odnośnie danych wejściowych	3
3	Przegląd podstawowych algorytmów sortujących	5
3.1	Quick Sort	5
3.1.1	Analiza algorytmu Quick Sort	5
3.1.2	Problemy związane z algorytmem Quick Sort	7
3.2	Merge Sort	7
3.2.1	Analiza algorytmu Merge Sort	7
3.2.2	Problemy związane z algorytmem Merge Sort	7
4	Przegląd hybrydowych algorytmów sortujących	9
4.1	Główne sposoby modyfikacji algorytmów	9
4.2	Rodzina algorytmów Quick Sort z deterministycznym algorytmem wyboru pivota	9
4.3	Rodzina algorytmów Quick Sort z niedeterministycznym algorytmem wyboru pivota	9
4.4	QuickMerge Sort	9
4.4.1	Pseudokod	9
4.4.2	Analiza algorytmu	9
4.4.3	Wnioski	10
4.5	Intro Sort	10
4.5.1	Pseudokod	10
4.5.2	Analiza algorytmu	10
4.5.3	Wnioski	10
5	Implementacja systemu	11
5.1	Struktura systemu	11
5.2	Koncepcje architektury silnika testującego	11
5.2.1	Wstrzykiwanie zależności	11
5.2.2	Obiektowość	11
5.2.3	Bezstanowość	12
5.3	Model aplikacji	12
5.3.1	Diagram klas	12
5.3.2	Diagram aktywności	13
5.4	Wzorce projektowe	13

5.4.1	Fasada	13
5.4.2	Obiekt-Wartość	14
5.4.3	Strategia	14
5.4.4	Budowniczy	14
6	Podsumowanie	15
	Bibliografia	17
A	Słownik pojęć	19
A.1	Notacja $O()$	19
A.2	Algorytm działający w miejscu	19
A.3	Algorytm stabilny	19
B	Środowisko uruchomieniowe aplikacji	21
B.1	Zmienne środowiskowe	21
B.2	Biblioteki zewnętrzne	21
B.3	Instalowanie aplikacji	21
B.4	Przykładowy plik konfiguracyjny	21

Wstęp

Ogólna historia algorytmów sortujących. Motywacja tworzenie algorytmów

1.1 Cel pracy

TODO: Motywem przewodnim pracy jest analiza nowoczesnych algorytmów sortujących w miejscu, takich jak koncepcja QuickMerge Sort. W tym celu przygotowano analizę porównawczą podstawowych algorytmów sortujących oraz dokonano przeglądu zmodyfikowanych wersji tych algorytmów oraz przeanalizowano nowoczesne algorytmy hybrydowe, będące połączeniem dwóch lub wielu algorytmów podstawowych.

1.2 Zakres pracy

TODO: Aby ułatwić analizę algorytmów przygotowany został silnik testujący oraz silnik graficzny, które w oparciu o plik konfiguracyjny przeprowadzają testy oraz tworzą wizualizację wyników tych testów. Wykorzystując podane narzędzia została przeprowadzona analiza podstawowych oraz hybrydowych algorytmów.

1.3 Przegląd literatury

TODO: Ogólny opis pracy Sebastiana Wilda. Pomysł na algorytm QuickMerge Sort.

1.4 Zawartości pracy

TODO: Ogólny sposób organizacji dokumentu. Rozdział pierwszy - analiza matematyczna problemu. Rozdział drugi - przegląd podstawowych algorytmów sortujących. Rozdział trzeci - przegląd hybrydowych algorytmów sortujących.



Analiza problemu

2.1 Model matematyczny przypadku średniego

2.2 Model matematyczny przypadku pesymistycznego

2.3 Założenia

2.3.1 Założenia odnośnie testowanych parametrów

TODO: Testowana będzie złożoność czasowa, w tym celu analizowane są takie parametry jak: liczba porównań, liczba swapów oraz liczba przypisań, z pominięciem operacji wykonywanych na iteratorach pętli - wyjaśnienie dlaczego.

TODO: Ponieważ rzeczywisty czas wykonywania algorytmu różni się w zależności od maszyny oraz architektury systemu na którym przeprowadzany test, zostało przyjęte następujące założenie: Złożoność czasowa została określona wzorem:

$$T = n_c + 3 \cdot n_s + n_a$$

Gdzie:

T - czas trwania algorytmu

n_c - liczba operacji porównania

n_s - liczba operacji zamiany miejsc

n_a - liczba operacji przypisania

2.3.2 Założenia odnośnie danych wejściowych

TODO: Wiele algorytmów sortujących bazuje na pewnych założeniach odnośnie wejściowego zbioru danych. Np. algorytm ... doskonale radzi sobie ze zbiorem danych prawie posortowanym, tzn. takim w którym W tej pracy zakładamy że dane wejściowe będą losowym ciągiem liczb powtórzeń.



Przegląd podstawowych algorytmów sortujących

3.1 Quick Sort

Historia algorytmu Quick Sort sięga drugiej połowy XX wieku. W roku 1959 brytyjski naukowiec Tony Hoare opracował, a dwa lata później opublikował pierwszą wersję tego algorytmu. Od tamtego czasu powstało wiele udoskonaleń tego algorytmu, jednak jego koncepcja nadal jest widoczna we współczesnych językach programowania ¹. Na cześć algorytmu Quick Sort standardowa funkcja sortująca w języku C++ nosi nazwę `qsort` ².

Algorytm Quick Sort składa się z dwóch etapów. Pierwszym z nich jest partycjonowanie zbioru wejściowego. Po tym kroku tablica wejściowa jest rozbita na dwa rozłączne zbiory, w których wszystkie elementy pierwszego zbioru są skumulowane po lewej stronie tablicy oraz każdy z tych elementów jest większy od dowolnego elementu z drugiej tablicy. Drugim etapem jest rekurencyjne sortowanie lewej oraz prawej podtablicy. Algorytm Quick Sort wykorzystuje technikę dziel i zwyciężaj, ponieważ problem sortowania tablicy wejściowej rozбивa na sortowanie dwóch podtablic.

3.1.1 Analiza algorytmu Quick Sort

Liczba operacji wykonywanych przez algorytm Quick Sort została przeanalizowana pod kątem trzech przypadków: optymistycznego, średniego oraz pesymistycznego.

Dla algorytmu Quick Sort przypadek optymistyczny (3.1) następuje wówczas, gdy algorytm partycjonowania przy każdym wywołaniu dzieli tablicę wejściową na dwie równe części. Efekt ten uzyskano, wprowadzając dane już posortowane oraz stosując algorytm wybierający pivot dokładnie w połowie tablicy. Z analizy eksperymentalnej wynika, że w przypadku optymistycznym algorytm działa ze złożonością czasową $O(n \log n)$.

Przypadek pesymistyczny (3.2) zachodzi, gdy drzewo wywołań rekurencyjnych jest możliwie najgłębsze. Efekt ten uzyskano, wprowadzając dane już posortowane oraz stosując algorytm wybierający pivot jako ostatni element tablicy. W tej sytuacji w kolejnych iteracjach rozpatrywana jest tablica z rozmiarem o jeden mniejszy od poprzedniej, a więc drzewo wywołań rekurencyjnych ma głębokość n . Z analizy wynika, że złożoność czasowa algorytmu w przypadku pesymistycznym wynosi $O(n^2)$.

Przypadek średni (3.1) został zbadany wprowadzając losowe dane z powtórzeniami oraz stosując algorytm wybierający pivot jako ostatni element tablicy. Analiza eksperymentalna wykazała, że w przypadku średnim algorytm Quick Sort ma złożoność czasową równą $O(n \log n)$, a więc jest tego samego rzędu co dla przypadku optymistycznego.

Porównując liczbę wykonywanych operacji można zauważyć, że algorytm Quick Sort wykonuje prawie dwa więcej operacji porównania niż operacji zamiany miejsc. Liczba pojedynczych operacji przypisania ro-

¹Dokumentacja biblioteki sortującej w języku java: <https://docs.oracle.com/en/java/javase/11/docs/api/java.base/java/util/Arrays.html>

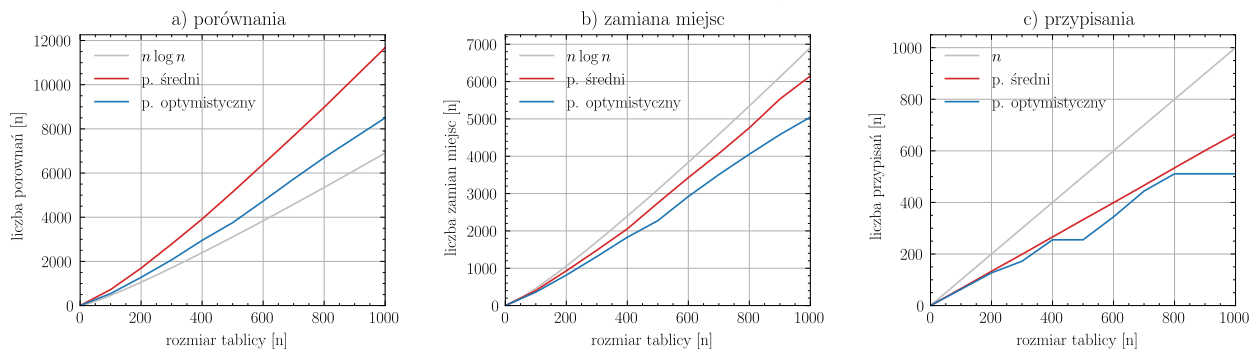
²Dokumentacja funkcji sortującej `qsort`: <https://en.cppreference.com/w/cpp/algorithm/qsort>



śnie liniowo, a więc jest znikoma w porównaniu z liczbą pozostałych operacji.

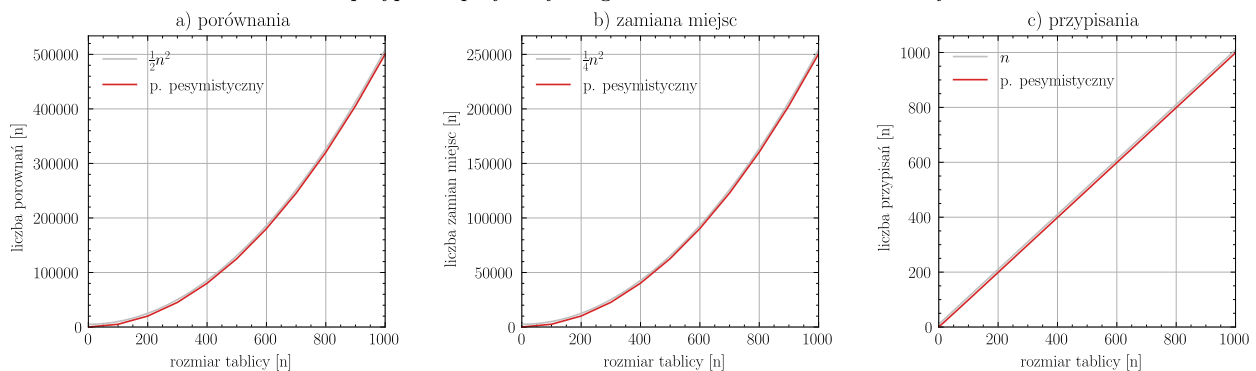
Analizując rozkład prawdopodobieństwa liczby wykonanych operacji (3.3) użyto tablicy losowych danych o stałym rozmiarze $n = 10000$. Można zauważyć, że liczba operacji porównania oraz liczba operacji zamiany miejsc nie są przedstawiane za pomocą rozkładu normalnego. Bardziej prawdopodobne jest wykonanie większej liczby tych operacji w stosunku do wartości średniej. Z kolei rozkład liczby wykonanych operacji przypisania przedstawia się za pomocą rozkładu normalnego, z jednakowym prawdopodobieństwem liczba ta może być większa lub mniejsza od wartości średniej.

Liczba wykonanych operacji w algorytmie Quick Sort dla przypadku średniego oraz optymistycznego w zależności od rozmiaru tablicy



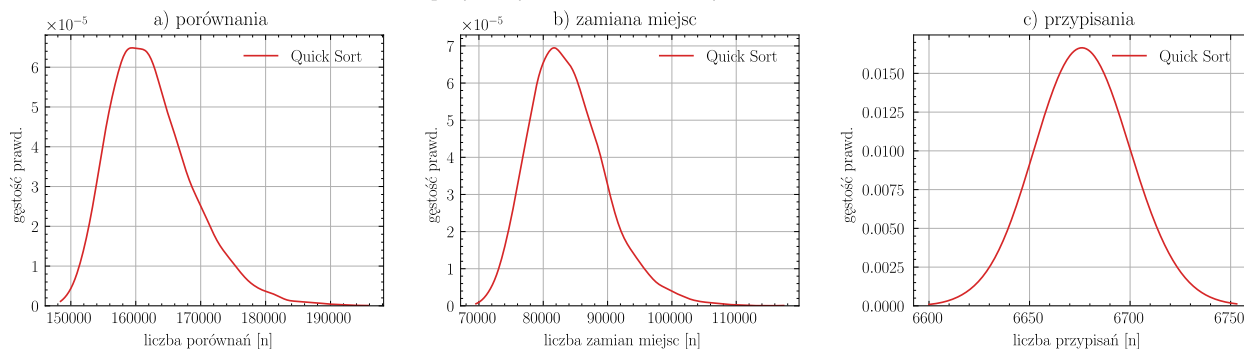
Rysunek 3.1

Liczba wykonanych operacji w algorytmie Quick Sort dla przypadku pesymistycznego w zależności od rozmiaru tablicy



Rysunek 3.2

Rozkład prawdopodobieństwa liczby wykonanych operacji
w algorytmie Quick Sort dla losowych danych
przy stałym rozmiarze tablicy $n = 10000$



Rysunek 3.3

3.1.2 Problemy związane z algorytmem Quick Sort

Głównym problemem algorytmu Quick Sort jest jego słaba pesymistyczna złożoność czasowa. Ponieważ algorytm działa rekurencyjnie, w przypadku pesymistycznym głębokość drzewa wywołań rekurencyjnych może przekroczyć maksymalną liczbę ramek stosu, powodując awaryjne zatrzymanie programu.

Kolejnym problemem tego algorytmu jest stosunkowo duża liczba wykonywanych operacji porównania w stosunku do liczby pozostałych operacji. Punkt ten jest szczególnie istotny w sytuacji, gdy sortowane są złożone struktury, dla których wykonanie pojedynczej operacji porównania jest znacznie kosztowniejsze od pozostałych operacji. W tym przypadku bardziej wskazany wydaje się użycie algorytmu Merge Sort, którego analizę przeprowadzono w kolejnym rozdziale.

3.2 Merge Sort

Ogólna koncepcja algorytmu, autor, rok powstania, bez pseudokodu.

3.2.1 Analiza algorytmu Merge Sort

Wykresy liczby operacji porównania, zamiany miejsc, przypisania. Eksperymentalna analiza wartości oczekiwanej liczby operacji.

3.2.2 Problemy związane z algorytmem Merge Sort

Konieczność alokacji dodatkowych zasobów.



Przegląd hybrydowych algorytmów sortujących

4.1 Główne sposoby modyfikacji algorytmów

1. Podmiana algorytmów składowych, np. inny algorytm partycjonowania.
2. Łączenie wielu algorytmów w jeden.

4.2 Rodzina algorytmów Quick Sort z deterministycznym algorytmem wyboru pivota

Algorytmy Quick Sort z różnymi algorytmami partycjonowania oraz różnymi algorytmami wyboru pivota
Wykresy porównujące algorytmy.

1. Partycjonowanie metodą Lemuto (domyślne)
2. Partycjonowanie metodą Hoare
3. Wybór pivota metodą median of three - pierwszy, środkowy, ostatni.
4. Wybór pivota jako mediana of medians
5. Wybór pivota metodą pseudomedian of nine
6. Wybór pivota algorytmem quick select

4.3 Rodzina algorytmów Quick Sort z niedeterministycznym algorytmem wyboru pivota

1. QuickSort z losowaniem pivota
2. QuickSort z losowaniem trzech liczb, wybór mediany (power of three choices)

4.4 QuickMerge Sort

Koncepcja algorytmu, mocne strony (Merge Sort bez konieczności alokacji pamięci)

4.4.1 Pseudokod

4.4.2 Analiza algorytmu

Wykresy liczby wykonywanych operacji w porównaniu do algorytmów bazowych. Wykresy gęstości liczby wykonywanych operacji dla stałej liczby n , np $n = 10000$.



4.4.3 Wnioski

Wyniki analizy porównawczej

4.5 Intro Sort

Ogólny opis algorytmu, gdzie jest wykorzystywany (`std::sort` w `g++`), zalety.

4.5.1 Pseudokod

4.5.2 Analiza algorytmu

Wykresy liczby wykonywanych operacji w porównaniu do algorytmów bazowych. Wykresy gęstości liczby wykonywanych operacji dla stałej liczby n , np $n = 10000$. Wykresy dla różnych algorytmów partycjonowania.

4.5.3 Wnioski

Wyniki analizy porównawczej

Implementacja systemu

5.1 Struktura systemu

Aplikacja składa się z dwóch modułów: silnika testującego oraz silnika graficznego. Działanie systemu jest określone na podstawie wspólnego pliku konfiguracyjnego. W pliku konfiguracyjnym określone są rodzaje testów jakie należy przeprowadzić oraz metadane potrzebne do wygenerowania wizualizacji.

Silnik testujący to generyczna biblioteka algorytmów sortujących oraz narzędzie przetwarzające te algorytmy. Aplikacja w oparciu o plik konfiguracyjny generuje zestaw testowy oraz utrzuła wyniki przeprowadzonych testów na dysku. W zależności od konfiguracji, silnik testujący może sumować, zliczać lub uśredniać liczbę wykonywanych operacji takich jak: liczba porównań, liczba operacji zamiany miejsc, liczba operacji przypisania oraz czas trwania algorytmu. Ta część aplikacji została napisana w języku C++¹ z wykorzystaniem technik programowania obiektowego.

Silnik graficzny to zbiór skryptów przetwarzających wyniki z silnika testującego. Na podstawie pliku konfiguracyjnego oraz danych testowych generowane są wizualizacje graficzne w postaci wykresów, dzięki czemu użytkownik końcowy może w łatwy sposób analizować oraz porównywać badane algorytmy. Ta część systemu została napisana w języku Python² przy użyciu biblioteki matplotlib³.

5.2 Koncepcje architektury silnika testującego

5.2.1 Wstrzykiwanie zależności

Większość algorytmów sortujących składa się z kilku odrębnych kroków. Niektóre z tych kroków są na tyle złożone, że stanowią osobne algorytmy. Dla przykładu jednym etapów sortowania metodą Quick Sort jest partycjonowanie danych wejściowych na rozłączne zbiory. Aby w łatwy sposób umożliwić modyfikację testowanych algorytmów, bez konieczności ponownej implementacji całego procesu, zastosowano technikę wstrzykiwania zależności. Jeżeli algorytm testujący korzysta z innego algorytmu, to algorytm składowy jest wstrzykiwany w trakcie działania programu. Dzięki temu lekka modyfikacja testowanego algorytmu ogranicza się do podmiany jego algorytmów składowych, bez konieczności ingerowania w strukturę bazową.

5.2.2 Obiektowość

Aby uprościć organizację kodu zastosowano model obiektowy. Każdy z algorytmów wykorzystywanych w systemie został zamodelowany za pomocą odrębnej klasy. Dla każdej rodziny algorytmów tego samego typu istnieje nadrzędna klasa bazowa określająca interfejs dla tej rodziny. Korzyści wynikające z zastosowanego modelu, takie jak statyczny polimorfizm oraz dziedziczenie, gwarantują bardziej wiarygodne działanie programu oraz umożliwiają wykrywanie błędów strukturalnych już na etapie kompilacji projektu.

¹Dokumentacja języka C++: <https://en.cppreference.com>

²Dokumentacja języka Python: <https://docs.python.org/3/>

³Dokumentacja biblioteki matplotlib: <https://matplotlib.org/>



5.2.3 Bezstanowość

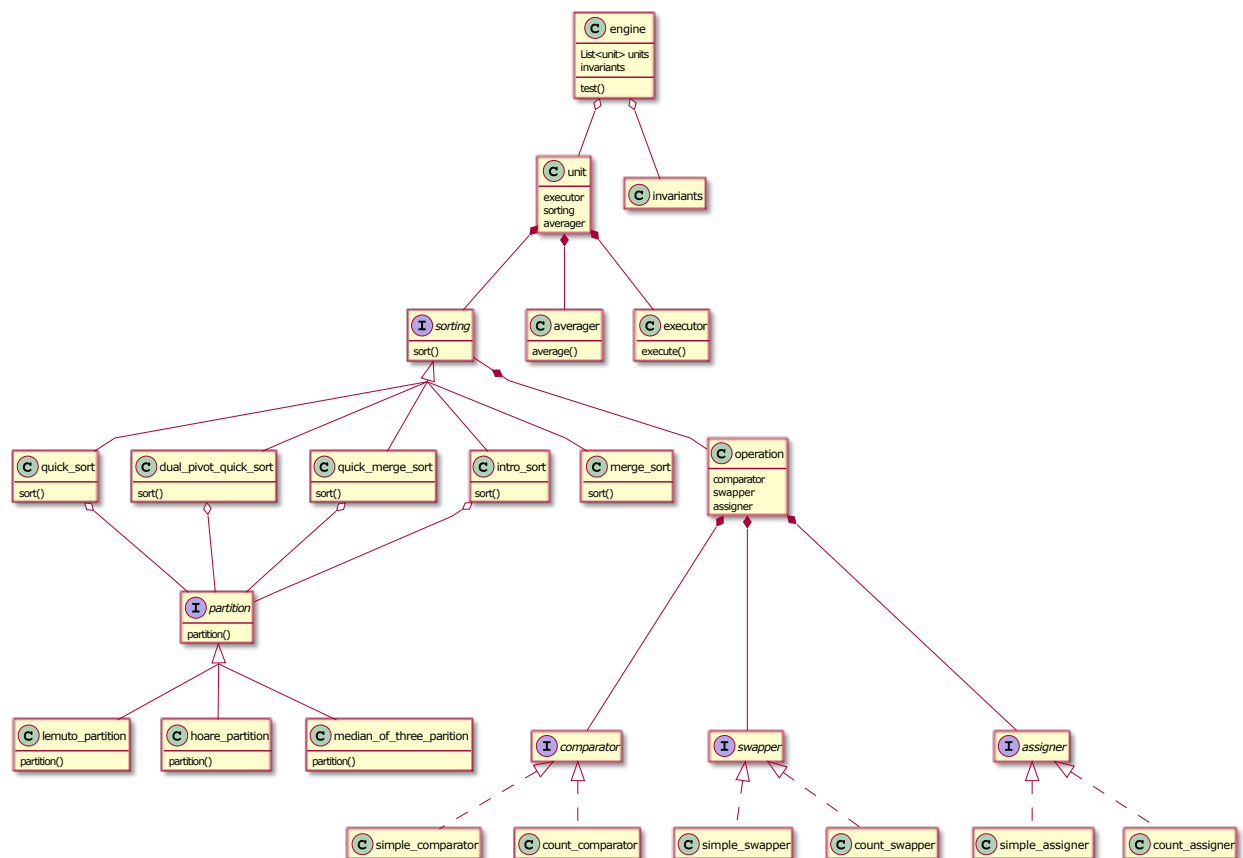
Powszechnym problemem w programowaniu obiektowym jest przechowywanie stanu. Problem ten wynika po części z praktyki hermetyzacji danych wewnątrz obiektowej abstrakcji. Użytkownik zewnętrzny korzystając z interfejsu danego komponentu nie ma dostępu do procesów zachodzących w jego wnętrzu. Może to prowadzić do tzw. efektów ubocznych (ang. side effects), przez co wyniki zwracane przez program stają się niewiarygodne.

Aby tego uniknąć zastosowano model bezstanowy. Żaden z algorytmów sortujących w zaimplementowanym systemie nie posiada zmiennych składowych, które mogłyby zostać zmodyfikowane w trakcie działania programu. Podczas testowania dane są przekazywane poprzez sygnatury metod, wzorując się na technice programowania funkcyjnego. Dzięki temu wszystkie algorytmy sortujące wykorzystane w implementacji są oznaczone jako niemodyfikowalne, nie mogą zmienić stanu aktualnie testowanego algorytmu. Gwarantuje to całkowitą separację poszczególnych testów.

5.3 Model aplikacji

5.3.1 Diagram klas

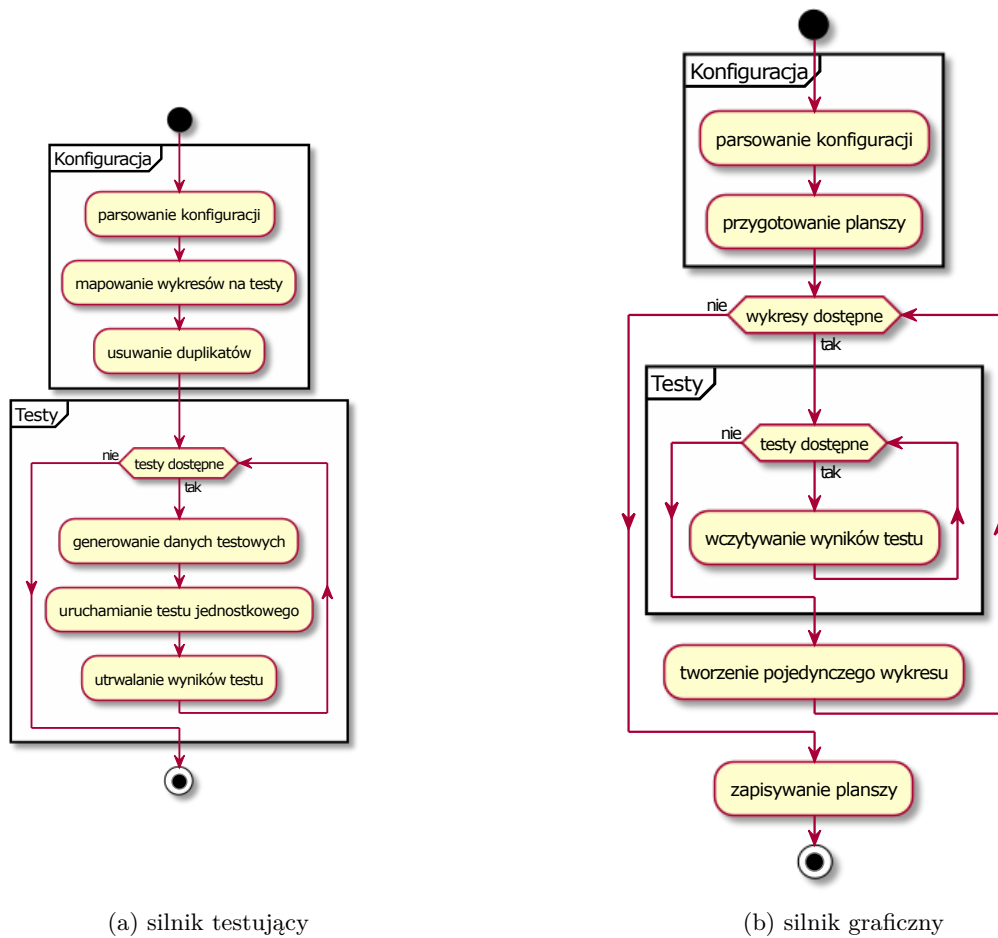
TODO: opis diagramu



Rysunek 5.1: Diagram klas silnika testującego

5.3.2 Diagram aktywności

TODO: opis diagramu

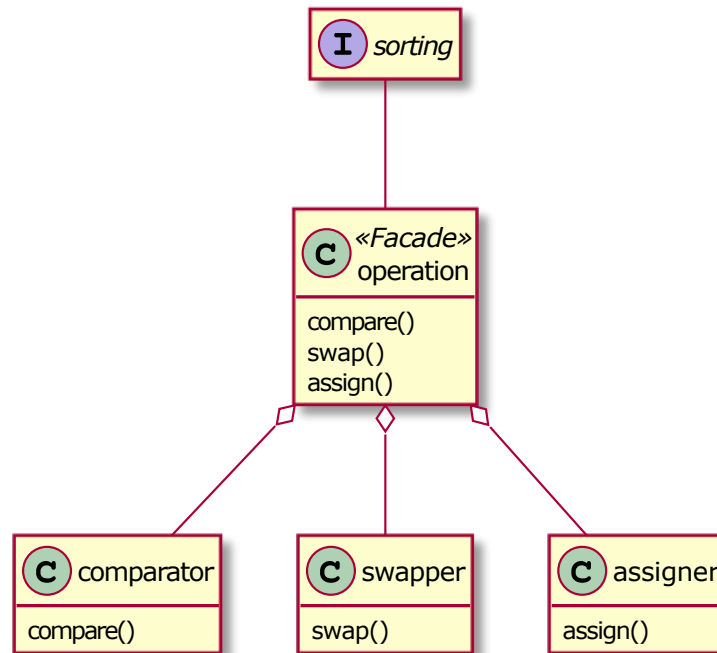


Rysunek 5.2: Diagram aktywności projektowanego systemu

5.4 Wzorce projektowe

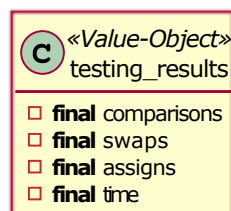
5.4.1 Fasada

W trakcie działania algorytm sortujący wykonuje wiele operacji atomowych, takich jak porównywanie elementów, zamiana elementów miejscami oraz operacje przypisania. Aby uniknąć nadmiaru odpowiedzialności dla klas sortujących zastosowano obiekt pośredniczący **operation** będący równocześnie **fasadą**. Fasada zapewnia jednolity interfejs dla wszystkich operacji atomowych oraz przekierowuje ich działanie do obiektów bezpośrednio odpowiedzialnych za ich wykonanie.



5.4.2 Obiekt-Wartość

Proces testowania algorytmu składa się z wielu iteracji. Każdy z atomowych testów wchodzących w skład iteracji powinien być całkowicie niezależny i odseparowana od innych testów. Aby to zapewnić, dane pochodzące z osobnych testów są przekazywane za pomocą **obiektów-wartości**. Pola w takim obiekcie po inicjalizacji stają się niemodyfikowalne. Użytkownik może jedynie odczytać ich wartość, bez możliwości ich modyfikacji. **Obiekt-wartość** jest gwarancją, że wyniki pochodzące z testu są rzetelne oraz nie zostały zmodyfikowane w trakcie przepływu danych pomiędzy procesami.



5.4.3 Strategia

5.4.4 Budowniczy

Podsumowanie

Podsumowanie wyników testowania algorytmów. Wnioski z analizy algorytmów hybrydowych.



Bibliografia



Słownik pojęć

A.1 Notacja $O()$

A.2 Algorytm działający w miejscu

A.3 Algorytm stabilny

Algorytm nie zamienia kolejnością elementów o tej samej wartości.



Środowisko uruchomieniowe aplikacji

Platforma uruchomieniowa aplikacji - Windows. Wykorzystane języki programowania - C++, Python.

B.1 Zmienne środowiskowe

Opis zmiennych środowiskowych TEST-DIRECTORY, CONFIG-DIRECTORY, PLOT-DIRECTORY.

B.2 Biblioteki zewnętrzne

Biblioteki w C++ oraz Pythonie potrzebne do uruchomienia aplikacji wraz z numerami wersji.

B.3 Instalowanie aplikacji

Uruchamianie skryptu zaciągającego potrzebne zależności. Uruchamianie pliku makefile kompilującego i instalującego aplikację. Cykl pracy programu - tworzenie konfiguracji, testowanie silnikiem testującym, wizualizacja wyników przy użyciu silnika graficznego.

B.4 Przykładowy plik konfiguracyjny

Plik konfiguracyjny w jsonie. Omówienie pliku, na początku są dane współdzielone przez wszystkie testy. Potem plik zawiera listę elementów typu plot, czyli listę osobnych testów.

