

# Specyfikacja implementacyjna „gentex”

Maciej Skarbek  
Nr albumu: 271088

23 marca 2015

# Wprowadzenie

Celem projektu jest stworzenie aplikacji w języku C, która będzie generować teksty wyjściowe na podstawie analizy innych tekstów wykorzystując przy tym łańcuchy Markova. Przy tworzeniu programu korzystać będziemy z systemu kontroli wersji git.

## Opis modułów

### **managment**

Rozpoznaje argumenty wywołań i steruje innymi modułami. Zawiera funkcje "main".

### **store**

Przechowuje macierz przejść z podziałem na prefiksy i sufiksy. Implementacja drzewa.

### **generation**

Generuje tekst wynikowy i zapisuje go do pliku.

### **reading**

Odczytuje teksty podane przez użytkownika. Analizuje je dzieląc na prefiksy i sufiksy a następnie przekazuje do modułu "store".

### **backup**

Tworzy i umożliwia wczytanie plików pośrednich z których w przyszłości można generować teksty wynikowe.

### **stats**

Tworzy statystyki łączone dla tekstów bazowych i oddzielną statystykę dla tekstu wynikowego (prawdopodobieństwo wystąpienia pojedynczego słowa, n-gramu z jakich został wygenerowany tekst i wskaźnik PMI).

### **error**

Obsługuje błędy zaistniałe podczas działania.

# Opis najważniejszych funkcji i struktur

## managment

- *void add(char \*\*prefix, char \* suffix)*  
prefix - tablica wyrazów z których składa sie prefiks  
suffix - sufiks  
Rozpoznanie argumentów wywołania i przekazanie sterowania do odpowiednich modułów.

## store

- Struktury  

<i>typedef struct{</i>	<i>typedef structnode{</i>	<i>typedef struct{</i>
<i>char **prefix;</i>	<i>ngram * g;</i>	<i>tree_t;</i>
<i>char **suffix;</i>	<i>structnode * left,*right;</i>	<i>intnumber_gram;</i>
<i>intsize_s;</i>	<i>}node_t,*tree_t;</i>	<i>intsize;</i>
<i>intn_s;</i>		<i>ngram **tab;</i>
<i>}ngram;</i>		<i>intn_s_max;</i>
		<i>}store;</i>

Tworzone jest drzewo do przechowywania prefiksów i sufiksów jak również tablica w której będą wskaźniki do komurek drzewa. Tablica będzie używana przy losowaniu prefiksów jak również przy tworzeniu pliku pośredniego (gdybyśmy chcieli stworzyć plik pośredni z drzewa był by on posortowany i przy wczytywaniu zamiast drzewo powstała by nam lista co znacząco wydłużyło by czas pracy programu).

- *tree\_t insert( tree\_t t, char \*\*prefix, char \* suffix )*  
return - funkcja zwraca drzewo z dodanym elementem  
t - drzewo do którego ma zostac dodany element  
prefix - tablica wyrazów z których składa sie prefiks  
suffix - sufiks  
Wstawia do drzewa prefiksy i sufiksy, jeśli prefiks juz istnieje to dopisuje tylko sufiks. Dodaje również wzkaźnik na "ngram" do "tab".
- *void add\_from\_backup(char \*\*prefix, char \*\*suffix, int n\_s)*  
prefix - tablica wyrazów z których składa sie prefiks  
suffix - tablica wszystkich sufiksów dla danego prefiksu  
Dodaje prefiks i całą listę sufiksów do drzewa.
- *ngram\* rand\_prefix()*  
return - zwraca strukturę pokazaną wyżej  
Losuje prefiks.
- *char\* rand\_suffix(char\*\* prefix)*  
return - zwraca sufiks

prefix - tablica wyrazów z których składa się prefiks  
Losuje sufiks dla podanego prefiksu.

## generation

- *void generation()*  
Generuje tekst wynikowy i zapisuje go do pliku.

## reading

- *void reading( char \* name\_file)*  
name\_file - plik który ma zostać zanalizowany i wczytany do "store"  
Czyta podany plik, dzieli na prefiksy i sufiksy a następnie przekazuje do "store" i "stat".

## backup

- *void backup()*  
Tworzy plik wczytuje i tworzy plik pośredni.

## stats

- *void stat\_add\_word(char \* word)*  
word - wyraz przeczytany z pliku poddanego analizie  
Dodaje pojedynczy wyraz do drzewa ze statystykami (zlicza ilość wystąpień).
- *void stat\_add\_ngram(char\*\* prefix, char \* suffix)*  
prefix - tablica wyrazów z których składa się prefiks  
suffix - sufiks  
Dodaje n-gram do drzewa ze statystykami i zlicza ilość wystąpień.
- *double get\_probability(tree\_stat t, char \* word)*  
return - zwraca prawdopodobieństwo  
t - nazwa drzewa w którym mamy szukać wyrazu(może to być kilka wyrazów) i informacji o ilości wystąpień  
word - wyraz dla którego liczone jest prawdopodobieństwo  
Liczy prawdopodobieństwo dla n-gramu.
- *long double get\_pmi(char\* wngram)*  
return - zwraca wartość wskaźnika PMI  
wngram - napis dla którego ma być liczony wskaźnik  
Liczy wskaźnik PMI dla n-gramu.

- *void write\_stat( char \* name\_file\_stat)*  
name\_file\_stat - nazwa pliku do którego mają być zapisywane statystyki Zapisuje statystyki do podanego pliku.

## **error**

- *void fatal(int err, const char \*msg)*  
err - instrukcja logiczna która warunkuje wystąpienie błędu np. uchwyt\_do\_pliku == NULL  
msg - Wiadomość która ma zostać wyświetlona w przypadku spełnienia się warunku Dostaje wyrażenie logiczne "err" i w przypadku potwierdzenia wystąpienia błędu pokazuje wiadomość "msg" i np. zamyka program.

## **Testowanie**

### **Użyte narzędzia**

#### **time**

Polecenie time do pomiaru czasu działania programu aby wybrać optymalny algorytm jaki zaimplementujemy.

#### **valgrind**

Narzędzie do debugowania pamięci i wykrywania wycieków pamięci.

### **Sposób testów**

Testy będą wykonywane dla pojedynczych funkcjonalności programu, dla całych modułów a następnie po przyłączeniu kolejnego modułu zostaną wykonane kompleksowe testy całego programu. Szczególną uwagę należy zwrócić na punkty krytyczne w których możemy spodziewać się błędów.

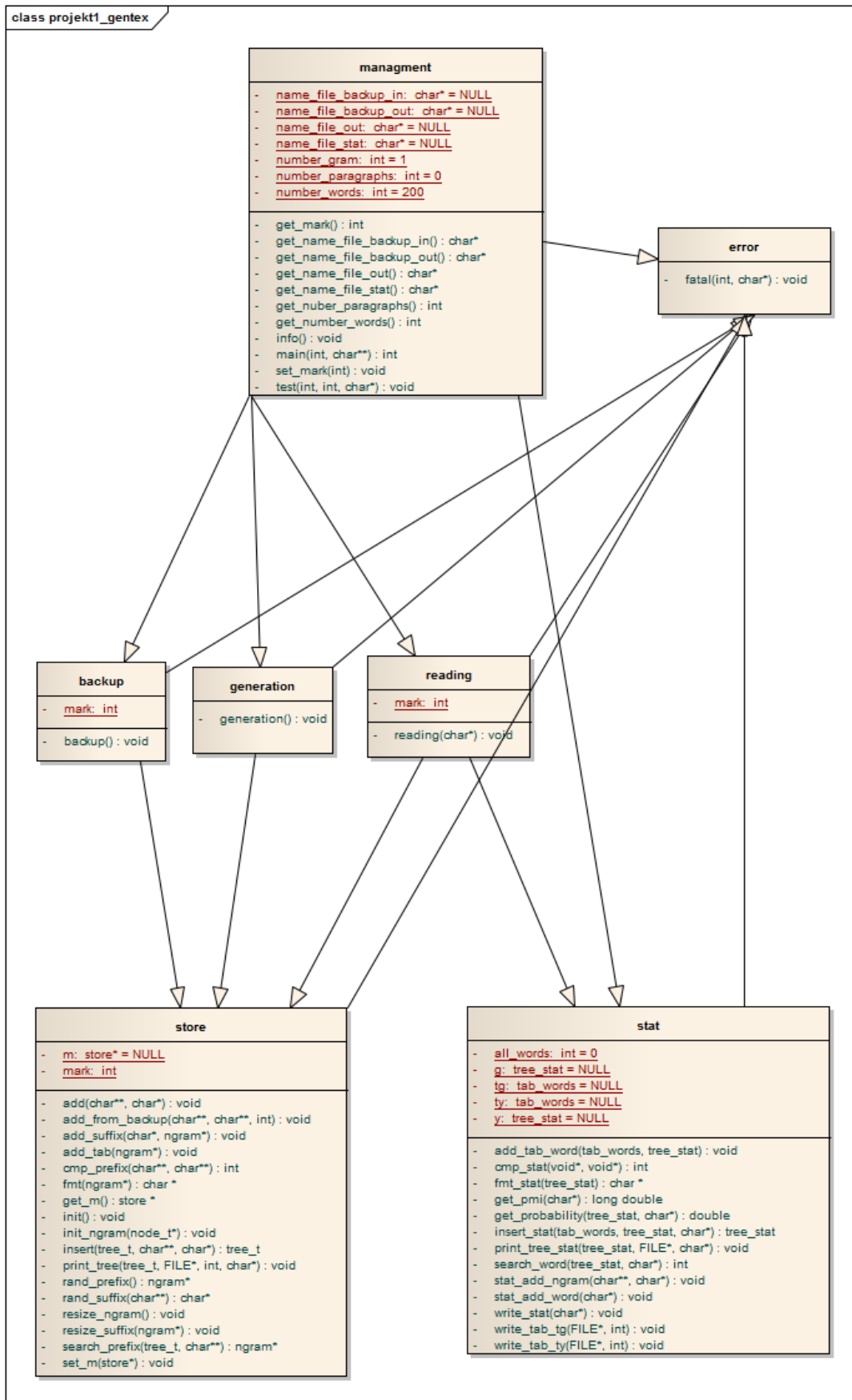
### **Punkty krytyczne**

- Bardzo duże pliki wejściowe
- Puste pliki wejściowe
- Bardzo małe pliki wejściowe (np 2 słowa)
- Błędne wywołanie programu

- Podanie pliku do zapisu który już istnieje
- Mała ilość pamięci urządzenia na którym zostanie uruchomiony program
- Wycieki pamięci

## Diagram modułów

Znaczenie strzałki (moduł1  $\longrightarrow$  moduł2) moduł1 korzysta z jakiejś funkcji w moduł2.



Rysunek 1: