

БАГАТОПЛАТФОРМЕНА СИСТЕМА ПЕРЕВІРКИ ОФОРГАФІЇ НА ОСНОВІ СЛОВНИКОВОЇ ТЕХНОЛОГІЇ HUNSPELL

Основною метою цього дослідження є розробка нового методу підходу до проектування підсистеми перевірки орфографії на основі технології перевірки орфографії за словником Hunspell. Запропоновано та описано структурну і об'єктно-орієнтовну схему підсистеми перевірки орфографії. Проведено аналіз ефективності запропонованих підходів на основі експерименту.

The main purpose of this study is to develop a new method of approach to the design of a subsystem spell checking technology based spell checking dictionary for Hunspell. Proposed and described structural and object-oriented scheme subsystem checker. The analysis of the effectiveness of the proposed approach based on experiment.

Постановка проблеми. Якість інформаційно-пошукових та інтеграційних можливостей у великій мірі залежить від якості даних, що містяться в базі даних. Під час роботи з даними якого можуть виникати ситуації, що призводять до появи різного типу помилок.

Формування цілі статті. В сучасному світі розробники інформаційних систем значну увагу приділяють перевірці орфографії, але багато цих засобів є платними і не є багатоплатформеними.

Підсистема перевірки орфографії, що реалізована на технології Java для операційних систем Windows, MAC OS, Linux, схематично має структуру, представлена автором на рис. 1.

Експерименти з оцінювання швидкодії системи. На основі описаних підходів у статті Ярмолюка Р. С. [1] було розроблено програмне забезпечення, що реалізує можливості пошуку орфографічних помилок у базах даних на технології Java.

Технологія Java була обрана через те, що вона дозволяє створювати незалежні від платформи програми шляхом компіляції в проміжне представлення, яке називається байт кодом [2].

Процес установки з'єднання з SQL-сервером може продовжуватися набагато більший час, ніж безпосередня обробка

запиту до бази після установки з'єднання. Слідуючи логіці, треба уникати зайвих з'єднань до бази, не від'єднуючись від неї там, де це можна зробити, якщо надалі планується продовжити роботу з SQL-сервером [4]. Зокрема технологія Java завдяки багатоплатформеності дозволяє будувати клієнт/сервер програми, запиту яких будуть працювати на різних платформах із нормальною швидкістю.

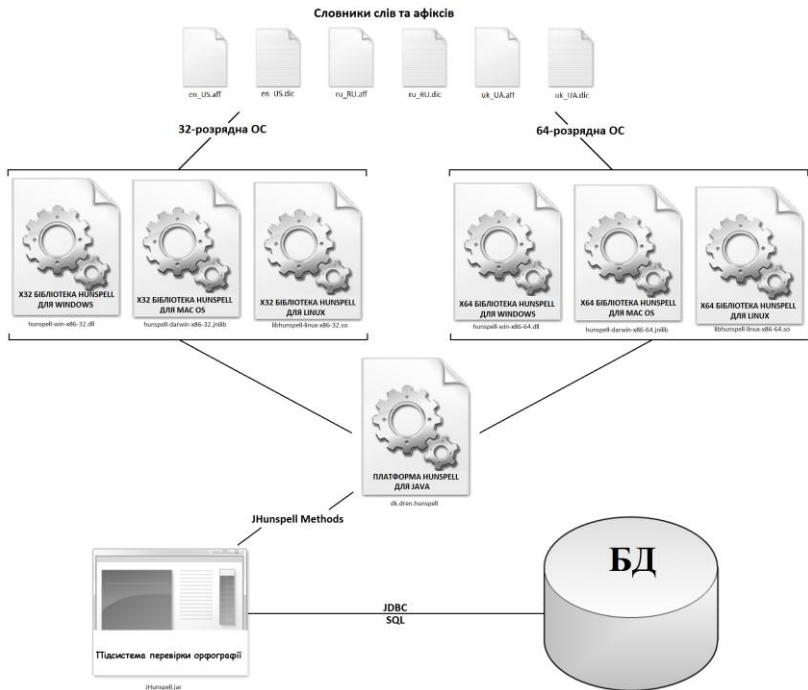


Рис. 1. Схема структурної організації підсистеми перевірки орфографії бази даних

Під час експлуатації системи виникла необхідність у дослідженні швидкодії ядра системи Hunspell [3], а саме оцінку часу затраченого на перевірку великого масиву слів. Для оцінки верхньої часової межі для кожного слова повинен відбуватись певний перебір словника, тобто слово з великою ймовірністю повинно бути відсутнім у словнику.

Технологічні умови проведення експерименту. Для перевірки орфографії було обрано загальнолексичні словники

OpenOffice.org. Причина вибору обумовлюється їх безкоштовністю і постійним оновленням. Параметри кожного із словників наведені у табл. 1.

Тестування проводилось на двох комп'ютерних платформах на базі операційної системи Microsoft Windows. Критичні параметри кожної із систем наведені у табл. 2.

Таблиця 2

Параметри комп'ютерних платформ тестування

	Процесор	Оперативна пам'ять	Операційна система
Перша платформа	Intel(R) Pentium (R) CPU B970 2.30 GHz	4 ГБ	Microsoft Windows 7 Максимальна 64- розрядна
Друга платформа	AMD Athlon(tm) II X2 245 Processor 2.90 GHz	12 ГБ	Microsoft Windows XP SP3 64- розрядна

Дослідження швидкодії ядра системи Hunspell. На першому етапі під час експерименту для створення вхідних даних для перевірки було розроблено клас із методами генерації випадкових слів певної довжини. Метод `randomWord.UA(int len, int nom)`, `randomWord.RU(int len, int nom)`, `randomWord.ENG(int len, int nom)` (Рис. 2) використовуються для генерування `nom` кількості символьних послідовностей UNICOD (тестових слів довжини `len`) для українського, російського та англійського алфавіту відповідно. На другому етапі метод `test.tspell (String[] s, int index)` реалізує перевірку кожного тестового слова із масиву слів `s` за допомогою методу `misspelled(s[i])` та повертає час (у мілісекундах), затрачений на перевірку `nom` кількості тестових слів довжини `len` за словником відповідної мови `index`.

Тестування проводилось для кожного масиву слів відповідної мови. Розмірність масиву складала 100 000 тестових слів відповідної довжини. Максимальна довжина слова 23 символів UNICOD. Оскільки тестові слова генерувались як довільна послідовність символів, то розглядається найгірший тестовий варіант: кожне із 100 000 слів з великою ймовірністю (близько 0,95) є неправильним. Результати оцінки верхньої часової межі швидкодії системи перевірки орфографії на основі технології Hunspell показані на рис. 3 та 4.

```

1  class randomWord {
2      int a;
3      String[] UA (int len,int nom) {
4          Random generator = new Random();
5          String[] ret=new String[nom];
6          char[] word1 = new char[len];
7          for (int i=0; i<nom;i++)
8          { for (int j = 0; j < len; j++)
9              { word1[j] = (char)(generator.nextInt(32) + 1072);
10                 switch (word1[j])
11                 { case 'н': word1[j] = 'i'; break;
12                   case 'ь': word1[j] = 'i'; break;
13                   case 'э': word1[j] = 'e'; break; } }
14                 ret[i]=new String(word1); }
15          return ret; }
16      String[] RU (int len,int nom) {
17          Random generator = new Random();
18          String[] ret=new String[nom];
19          char[] word1 = new char[len];
20          for (int i=0; i<nom;i++)
21          { for (int j = 0; j < len; j++)
22              { word1[j] = (char)(generator.nextInt(32) + 1072); }
23              ret[i]=new String(word1); }
24          return ret; }
25      String[] ENG (int len,int nom) {
26          Random generator = new Random();
27          String[] ret=new String[nom];
28          char[] word1 = new char[len];
29          for (int i=0; i<nom;i++)
30          { for (int j = 0; j < len; j++)
31              { word1[j] = (char)(generator.nextInt(26) + 97); }
32              ret[i]=new String(word1); }
33          return ret; } }

```

Рис. 2. Код класу randomWord

Аналіз результатів експерименту. Проаналізувавши отримані результати, можемо зробити такі висновки:

- Швидкість перевірки слова лінійно залежить від кількості символів, з яких складається слово, і ця тенденція зберігається незалежно від апаратної платформи або параметрів словників.
- Основним фактором, що впливає на швидкодію методу, є розмір словника основ слів.
- Збільшення затрачуваного часу для обробки великих слів на англійській мові може пояснюватись особливістю афіксної конструкції мови.

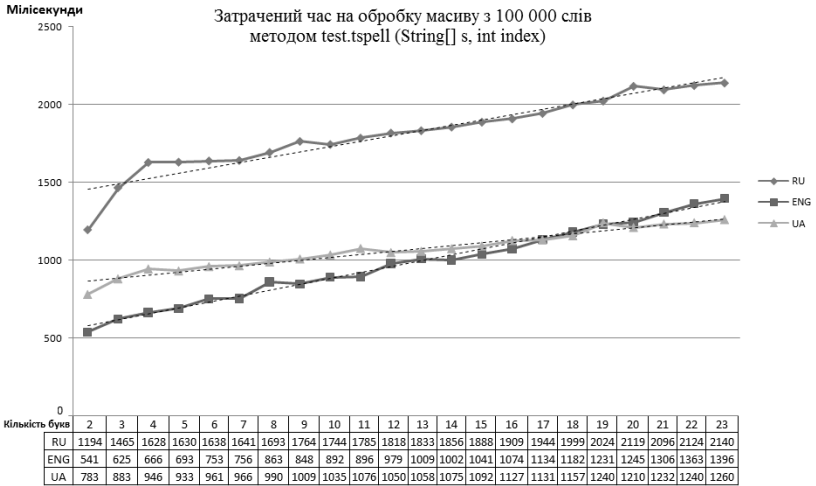


Рис. 3. Результати тестування на першій апаратній платформі

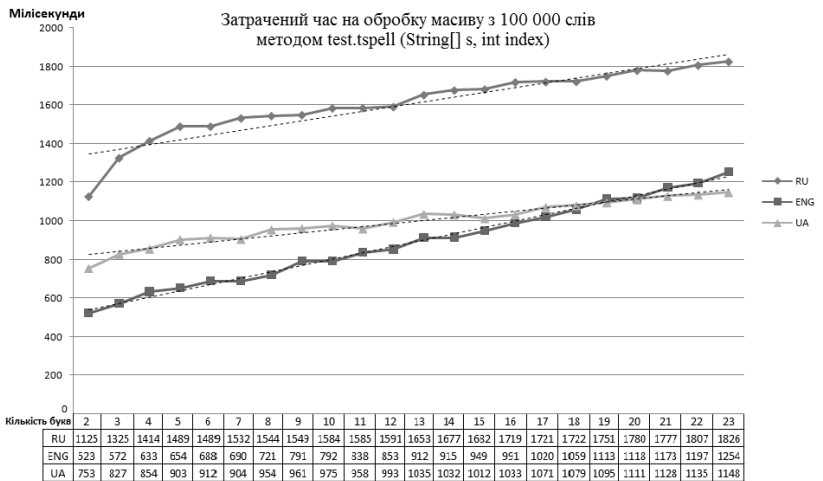


Рис. 4. Результати тестування на другій апаратній платформі

Висновки. У результаті проведеного дослідження було розроблено підхід для проектування підсистеми автоматичної корекції орфографічних помилок у записах бази даних на основі словникової технології. Було обрано відкриту технологію Hunspell для побудови на її основі даної підсистеми. Запропонована схема структурної

організації підсистеми на основі технології Java. Проведено експеримент із тестуванням швидкодії ядра системи.

Література

1. Ярмолюк Р.С. Підсистема перевірки орфографії електронного каталогу бібліотеки на основі технології Hunspell / Р.С. Ярмолюк // Вісник Національного Університету “Львівська Політехніка”, серія Інформаційні системи та мережі, - 2012. - № 743 – С.219-230.
2. Герберт Шилдт Java. Полное руководство. 7-е издание. / Г. Шилдт // Вильямс, 2007. –1034 с.
3. Surhone L.M. HunsPELL / L.M. Surhone, M.T. Tennoe, S.F. Henssonow – Betascript Publishing, 2010 – 116 p.
4. Кузнецов М. С. Самоучитель MySQL 5 / М. С. Кузнецов // Петербург – 2007. – 560 с.