

PDF Search Engine

1.1

Generated by Doxygen 1.8.11

Contents

1	Module Index	1
1.1	Modules	1
2	Namespace Index	3
2.1	Namespace List	3
3	Hierarchical Index	5
3.1	Class Hierarchy	5
4	Class Index	7
4.1	Class List	7
5	Module Documentation	9
5.1	Indexing	9
5.1.1	Detailed Description	9
5.2	Debugging	10
5.2.1	Detailed Description	10
5.2.2	Macro Definition Documentation	10
5.2.2.1	CASSERT	10
5.2.2.2	DUMP_TO_FILE	10
5.2.2.3	NON_UNIT_TEST_ASSERT	10
5.3	Stemming	11
5.3.1	Detailed Description	11
5.4	Mathematics	12
5.4.1	Detailed Description	12

5.4.2	Function Documentation	12
5.4.2.1	compare_doubles(const double actual, const double expected, const double delta=1e-6)	12
5.4.2.2	compare_doubles_greater(const double left, const double right, const double delta=1e-6)	13
5.4.2.3	compare_doubles_less(const double left, const double right, const double delta=1e-6)	13
5.4.2.4	compare_doubles_less_or_equal(const double left, const double right, const dou- ble delta=1e-6)	13
5.4.2.5	int_to_bool(const T intVal)	14
5.4.2.6	safe_divide(const T dividend, const T divisor)	14
5.4.2.7	safe_modulus(const T dividend, const T divisor)	14
5.5	StringOperations	16
5.5.1	Detailed Description	16
5.6	Utilities	17
5.6.1	Detailed Description	18
5.6.2	Macro Definition Documentation	18
5.6.2.1	size_of_array	18
5.6.3	Function Documentation	18
5.6.3.1	copy_member(inT begin, inT end, outT dest, member_extract_functorT get_value)	18
5.6.3.2	is_within(const T value, const T first, const T second)	18
5.6.3.3	operator()(T value) const	18
5.6.3.4	within(T range_begin, T range_end)	19
5.6.3.5	within_range(const T start, const T end, const T value)	19

6 Namespace Documentation	21
6.1 stemming Namespace Reference	21
6.1.1 Detailed Description	21
6.2 string_util Namespace Reference	22
6.2.1 Detailed Description	24
6.2.2 Function Documentation	25
6.2.2.1 axtoi(const T *hexStr, size_t length=-1)	25
6.2.2.2 find_last_not_of(const T *string, const T *search, size_t offset=std::basic_string< T >::npos)	25
6.2.2.3 find_last_of(const T *string, const T ch, size_t offset=-1)	25
6.2.2.4 find_matching_close_tag(const T *string, const T openSymbol, const T closeSymbol, const bool fail_on_overlapping_open_symbol=false)	26
6.2.2.5 has_suffix(const T *text, const size_t text_length, const T *suffix, const size_t suffix_length)	26
6.2.2.6 is_hex_digit(const T ch)	26
6.2.2.7 is_space(const T ch)	26
6.2.2.8 remove_blank_lines(string_typeT &Text)	26
6.2.2.9 remove_extra_spaces(string_typeT &Text)	27
6.2.2.10 strcspn_pointer(const T *string1, const T *string2, const size_t string2Length)	27
6.2.2.11 strnatordcmp(const T *first_string, const T *second_string, bool case_insensitive=false)	27
6.2.2.12 strnchr(const T *string, const T ch, size_t numberOfCharacters)	27
6.2.2.13 strncspn(const T *stringToSearch, const size_t stringToSearchLength, const T *searchString, const size_t searchStringLength)	27
6.2.2.14 strnistr(const T *string, const T *strSearch, const size_t string_len)	28
6.2.2.15 strnlen(const T *str, const size_t maxlen)	28
6.2.2.16 strrstr(const T *string, const T *search, size_t offset)	28
6.2.2.17 strtod_ex(const Tchar_type *nptr, Tchar_type **endptr)	28
6.2.2.18 strtol(const char *str, char **strend, int radix)	28

7	Class Documentation	29
7.1	AVLTreeIndex Class Reference	29
7.1.1	Detailed Description	30
7.1.2	Member Function Documentation	30
7.1.2.1	addIndex(string &word, string &doc)	30
7.1.2.2	findIndex(string key)	30
7.1.2.3	insert(string x, string docname)	30
7.1.2.4	insert(string x, string docname, int count)	31
7.1.2.5	split(const std::string &s, char delim)	31
7.1.2.6	topwords()	31
7.1.2.7	totalWordsIndexed()	31
7.2	backup_variable< T > Class Template Reference	32
7.2.1	Detailed Description	32
7.3	bookmark Struct Reference	32
7.3.1	Detailed Description	33
7.4	comparable_first_pair< T1, T2 > Class Template Reference	33
7.4.1	Detailed Description	33
7.5	document Struct Reference	34
7.5.1	Detailed Description	34
7.6	DocumentParser Class Reference	34
7.6.1	Detailed Description	35
7.6.2	Constructor & Destructor Documentation	36
7.6.2.1	DocumentParser(IndexHandler *)	36
7.6.3	Member Function Documentation	36
7.6.3.1	extract(string)	36
7.6.3.2	getStemmed(string &)	36
7.6.3.3	getTotalPages()	37
7.6.3.4	isStopWord(string &)	37
7.6.3.5	numOfDocs()	37
7.6.3.6	rawTextExtract(string fileStream)	37

7.7	double_less Class Reference	38
7.7.1	Detailed Description	38
7.8	stemming::english_stem< string_typeT > Class Template Reference	38
7.8.1	Detailed Description	39
7.8.2	Member Function Documentation	41
7.8.2.1	operator()(string_typeT &text)	41
7.9	string_util::equal_basic_string_i_compare< T > Class Template Reference	41
7.10	string_util::equal_string_compare< T > Class Template Reference	42
7.11	string_util::equal_string_i_compare< T > Class Template Reference	42
7.12	hashy Class Reference	43
7.12.1	Detailed Description	43
7.12.2	Constructor & Destructor Documentation	44
7.12.2.1	hashy()	44
7.12.3	Member Function Documentation	44
7.12.3.1	addIndex(string &word, string &name)	44
7.12.3.2	displayIndex(int index)	44
7.12.3.3	findIndex(string word)	44
7.12.3.4	Hashy(string key)	45
7.12.3.5	insertFromFile(string &word, string &doc, int &count)	45
7.12.3.6	NumItemsInIndex(int index)	45
7.12.3.7	split(const std::string &s, char delim)	45
7.12.3.8	topwords()	46
7.12.3.9	totalWordsIndexed()	46
7.13	index_node Struct Reference	46
7.13.1	Detailed Description	47
7.14	indexextractor Class Reference	47
7.14.1	Detailed Description	47
7.14.2	Constructor & Destructor Documentation	47
7.14.2.1	indexextractor(string)	47
7.15	IndexHandler Class Reference	48

7.15.1 Detailed Description	48
7.15.2 Constructor & Destructor Documentation	48
7.15.2.1 IndexHandler(char type)	48
7.15.3 Member Function Documentation	49
7.15.3.1 addIndex(string, string)	49
7.15.3.2 getDocs(string)	49
7.15.3.3 topFifty()	49
7.15.3.4 totalWordsIndexed()	49
7.16 IndexInterface Class Reference	50
7.16.1 Detailed Description	50
7.17 item Struct Reference	50
7.17.1 Detailed Description	51
7.18 string_util::less_basic_string_compare< T > Class Template Reference	51
7.19 string_util::less_string_compare< T > Class Template Reference	51
7.20 string_util::less_string_i_compare< T > Class Template Reference	52
7.21 string_util::less_string_n_compare< T > Class Template Reference	52
7.22 string_util::less_string_natural_order_i_compare< T > Class Template Reference	53
7.23 string_util::less_string_ni_compare< T > Class Template Reference	53
7.24 stemming::no_op_stem< string_typeT > Class Template Reference	53
7.25 QueryEngine Class Reference	54
7.25.1 Detailed Description	54
7.25.2 Constructor & Destructor Documentation	55
7.25.2.1 QueryEngine(IndexHandler *, DocumentParser *)	55
7.25.3 Member Function Documentation	55
7.25.3.1 findTDIFs(vector< document > &currdocs)	55
7.25.3.2 getDifference(vector< document > r, vector< document > l)	55
7.25.3.3 getIntersection(vector< document > l, vector< document > r)	56
7.25.3.4 getUnion(vector< document > l, vector< document > r)	56
7.25.3.5 querySearch(string query)	56
7.25.3.6 relevancySort(vector< document > &currdocs)	56

7.25.3.7	search(string word)	57
7.26	rawOutputExtractor Class Reference	57
7.26.1	Detailed Description	58
7.26.2	Constructor & Destructor Documentation	58
7.26.2.1	rawOutputExtractor(IndexHandler *, indexextractor *)	58
7.26.3	Member Function Documentation	58
7.26.3.1	getDocName()	58
7.26.3.2	getTotalPages()	58
7.26.3.3	getWordCount()	59
7.26.3.4	Init(const char *pszInput)	59
7.27	SearchEngine Class Reference	59
7.27.1	Detailed Description	61
7.27.2	Constructor & Destructor Documentation	61
7.27.2.1	SearchEngine()	61
7.27.3	Member Function Documentation	61
7.27.3.1	addBookmark(string book, string query)	61
7.27.3.2	addDocumentsToIndex(string docpath)	61
7.27.3.3	addToHistory(string query)	62
7.27.3.4	chooseStructure(char type)	62
7.27.3.5	display_search_results(string word)	62
7.27.3.6	displayBookmarks()	62
7.27.3.7	displayRawFile(string filePath)	63
7.27.3.8	getTotalPages()	63
7.27.3.9	numWordsIndexed()	63
7.27.3.10	readFilePathFromTXTFile()	63
7.27.3.11	readTotalPages()	63
7.27.3.12	topfifty()	64
7.28	stemming::stem< string_typeT > Class Template Reference	64
7.28.1	Detailed Description	69
7.28.2	Member Function Documentation	69

7.28.2.1	hash_french_yui(string_typeT &text, const wchar_t *vowel_string)	69
7.28.2.2	is_apostrophe(const wchar_t &ch) const	69
7.28.2.3	is_suffix_in_r1(const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U)	69
7.28.2.4	is_suffix_in_r2(const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U)	69
7.28.2.5	is_suffix_in_rv(const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U)	70
7.29	string_util::string_tokenize< T > Class Template Reference	70
7.29.1	Detailed Description	70
7.29.2	Constructor & Destructor Documentation	70
7.29.2.1	string_tokenize(const T &val, const T &delim)	70
7.29.3	Member Function Documentation	71
7.29.3.1	get_next_token()	71
7.29.3.2	has_more_delimiters() const	71
7.29.3.3	has_more_tokens() const	71
7.30	string_util::string_trim< char_typeT > Class Template Reference	71
7.30.1	Detailed Description	71
7.31	TextExtractor Class Reference	72
7.31.1	Detailed Description	72
7.31.2	Constructor & Destructor Documentation	72
7.31.2.1	TextExtractor(IndexHandler *, indexextractor *)	72
7.31.3	Member Function Documentation	73
7.31.3.1	getDocName()	73
7.31.3.2	getTotalPages()	73
7.31.3.3	getWordCount()	73
7.31.3.4	Init(const char *pszInput, string docName)	73
7.32	userinterface Class Reference	73
7.32.1	Detailed Description	74
7.32.2	Constructor & Destructor Documentation	74
7.32.2.1	userinterface()	74
7.32.3	Member Function Documentation	75
7.32.3.1	use()	75
7.33	within< T > Class Template Reference	75
7.33.1	Detailed Description	75

Chapter 1

Module Index

1.1 Modules

Here is a list of all modules:

Indexing	9
Debugging	10
Stemming	11
Mathematics	12
StringOperations	16
Utilities	17

Chapter 2

Namespace Index

2.1 Namespace List

Here is a list of all documented namespaces with brief descriptions:

stemming	
Namespace for stemming classes	21
string_util	22

Chapter 3

Hierarchical Index

3.1 Class Hierarchy

This inheritance list is sorted roughly, but not completely, alphabetically:

backup_variable< T >	32
binary_function	
double_less	38
string_util::equal_basic_string_i_compare< T >	41
string_util::equal_string_compare< T >	42
string_util::equal_string_i_compare< T >	42
string_util::less_basic_string_compare< T >	51
string_util::less_string_compare< T >	51
string_util::less_string_i_compare< T >	52
string_util::less_string_n_compare< T >	52
string_util::less_string_natural_order_i_compare< T >	53
string_util::less_string_ni_compare< T >	53
bookmark	32
document	34
DocumentParser	34
index_node	46
indexextractor	47
IndexHandler	48
IndexInterface	50
AVLTreeIndex	29
hashy	43
item	50
stemming::no_op_stem< string_typeT >	53
pair	
comparable_first_pair< T1, T2 >	33
QueryEngine	54
rawOutputExtractor	57
SearchEngine	59
stemming::stem< string_typeT >	64
stemming::english_stem< string_typeT >	38
string_util::string_tokenize< T >	70
string_util::string_trim< char_typeT >	71
TextExtractor	72
unary_function	
within< T >	75
userinterface	73

Chapter 4

Class Index

4.1 Class List

Here are the classes, structs, unions and interfaces with brief descriptions:

AVLTreeIndex

The **AVLTreeIndex** (p. 29) class has multiple custom insert methods to allow it to either insert from inverted index or directly from parser 29

backup_variable< T >

Class that remembers its original value from construction 32

bookmark

The bookmark struct This struct uses two strings, one for the PDF name, the other for the query 32

comparable_first_pair< T1, T2 >

Pair interface that compares on the first item 33

document

The document struct This struct contains a string used as the word or the PDF name, and an int used as the count, and a double for the tf/idf. This struct is used in many different ways depending on the function it is used in 34

DocumentParser

Parses the PDFs. It extracts the text using the **TextExtractor** (p. 72) class, stems the words and removes stop words using the **IndexExtractor** class, and computes the term frequency using the custom AVLTree or the custom Hash table 34

double_less

"less" interface for double values 38

stemming::english_stem< string_typeT >

English stemmer 38

string_util::equal_basic_string_i_compare< T > 41

string_util::equal_string_compare< T > 42

string_util::equal_string_i_compare< T > 42

hashy

The hashy class has multiple custom insert methods to allow it to either insert from inverted index or directly from parser 43

index_node

The **index_node** (p. 46) struct allows for storage of inverted indices with other data in tree; most importantly the documents it is found in 46

indexextractor

The indexextractor class 47

IndexHandler

The **IndexHandler** (p. 48) class 48

IndexInterface

The **IndexInterface** (p. 50) class 50

item	
The hashy class is the custom implementation of the HashTable. The basic hash parts of the code was based off of the series of videos by Paul Programming	50
string_util::less_basic_string_compare< T >	51
string_util::less_string_compare< T >	51
string_util::less_string_i_compare< T >	52
string_util::less_string_n_compare< T >	52
string_util::less_string_natural_order_i_compare< T >	53
string_util::less_string_ni_compare< T >	53
stemming::no_op_stem< string_typeT >	53
QueryEngine	
The QueryEngine (p. 54) uses the inverted file index and searches it given a properly formatted Boolean query. It also ranks the results using tf/idf statistics	54
rawOutputExtractor	
The rawOutputExtractor (p. 57) class	57
SearchEngine	
The SearchEngine (p. 59) class	59
stemming::stem< string_typeT >	
The base class for language-specific stemmers	64
string_util::string_tokenize< T >	
Tokenizes a string using a set of delimiters	70
string_util::string_trim< char_typeT >	
Trims whitespace from around a string	71
TextExtractor	
The TextExtractor (p. 72) class	72
userinterface	
The UserInterface is a command line menu driven class that makes use of the SearchEngine (p. 59). It allows the user to enter into two modes, Maintenance and Query, allows the user many options including adding a new pdf to the inverted index, clearing the index, searching the PDF, outputting total pages, outputting total words indexed, outputting the top fifty words, outputting the corpus paths, storing and clearing the search history, storing and clearing bookmarks and outputting the raw text from a selected pdf	73
within< T >	
Determines if a value is within a given range	75

Chapter 5

Module Documentation

5.1 Indexing

Library for stemming words down to their root words.

5.1.1 Detailed Description

Library for stemming words down to their root words.

Date

2003-2016

Copyright

Oleander Software, Ltd.

Author

Oleander Software, Ltd.

This program is free software; you can redistribute it and/or modify it under the terms of the BSD License.

5.2 Debugging

Functions used for debugging.

Macros

- `#define CASSERT(x) typedef char __C_ASSERT__[(x) ? 1 : -1]`
- `#define NON_UNIT_TEST_ASSERT(x) assert(x)`
- `#define DUMP_TO_FILE(x, file) __debug::__dump_to_file((x), (file))`

5.2.1 Detailed Description

Functions used for debugging.

Date

2008-2015

Copyright

Oleander Software, Ltd.

Author

Oleander Software, Ltd.

This program is free software; you can redistribute it and/or modify it under the terms of the BSD License.

5.2.2 Macro Definition Documentation

5.2.2.1 `#define CASSERT(x) typedef char __C_ASSERT__[(x) ? 1 : -1]`

Validates that an expression is true at compile time. If the expression is false then compilation will fail.

5.2.2.2 `#define DUMP_TO_FILE(x, file) __debug::__dump_to_file((x), (file))`

Prints data stream to a specified file.

5.2.2.3 `#define NON_UNIT_TEST_ASSERT(x) assert(x)`

If unit test symbol (`__UNITTEST`) is defined then does nothing; otherwise asserts. This is useful for suppressing asserts when unit testing.

5.3 Stemming

Library for stemming words down to their root words.

Namespaces

- **stemming**

Namespace for stemming classes.

5.3.1 Detailed Description

Library for stemming words down to their root words.

Date

2003-2015

Copyright

Oleander Software, Ltd.

Author

Oleander Software, Ltd.

This program is free software; you can redistribute it and/or modify it under the terms of the BSD License.

5.4 Mathematics

Math and statistics classes.

Classes

- class **double_less**
"less" interface for double values.

Functions

- template<typename T >
T **safe_modulus** (const T dividend, const T divisor)
Modulus operation that checks for modulus by zero or into zero (returns zero for those situations).
- template<typename T >
T **safe_divide** (const T dividend, const T divisor)
Division operation that checks for division by zero or into zero (returns zero for those situations).
- bool **compare_doubles** (const double actual, const double expected, const double delta=1e-6)
Compares two double values (given the specified precision).
- bool **compare_doubles_less** (const double left, const double right, const double delta=1e-6)
Compares two double values for less than (given the specified precision).
- bool **compare_doubles_less_or_equal** (const double left, const double right, const double delta=1e-6)
Compares two double values for less than or equal to (given the specified precision).
- bool **compare_doubles_greater** (const double left, const double right, const double delta=1e-6)
Compares two double values for greater than (given the specified precision).
- bool **double_less::operator()** (const double &left, const double &right) const
- template<typename T >
bool **int_to_bool** (const T intVal)
Converts an integral type to a boolean. Compilers complain about directly assigning an int to a bool (casting doesn't help either), so this works around that.

5.4.1 Detailed Description

Math and statistics classes.

Date

2004-2015

Copyright

Oleander Software, Ltd.

Author

Oleander Software, Ltd.

This program is free software; you can redistribute it and/or modify it under the terms of the BSD License.

5.4.2 Function Documentation

5.4.2.1 bool **compare_doubles** (const double *actual*, const double *expected*, const double *delta* = 1e-6) [inline]

Compares two double values (given the specified precision).

Parameters

<i>actual</i>	The value being reviewed.
<i>expected</i>	The expected value to compare against.
<i>delta</i>	The tolerance of how different the values can be. The larger the delta, the higher precision used in the comparison.

Returns

True if the value matches the expected value.

5.4.2.2 `bool compare_doubles_greater (const double left, const double right, const double delta = 1e-6) [inline]`

Compares two double values for greater than (given the specified precision).

Parameters

<i>left</i>	The value being reviewed.
<i>right</i>	The other value to compare against.
<i>delta</i>	The tolerance of how different the values can be. The larger the delta, the higher precision used in the comparison.

Returns

True if the value is greater than the other value.

5.4.2.3 `bool compare_doubles_less (const double left, const double right, const double delta = 1e-6) [inline]`

Compares two double values for less than (given the specified precision).

Parameters

<i>left</i>	The value being reviewed.
<i>right</i>	The other value to compare against.
<i>delta</i>	The tolerance of how different the values can be. The larger the delta, the higher precision used in the comparison.

Returns

True if the value is less than the other value.

5.4.2.4 `bool compare_doubles_less_or_equal (const double left, const double right, const double delta = 1e-6) [inline]`

Compares two double values for less than or equal to (given the specified precision).

Parameters

<i>left</i>	The value being reviewed.
<i>right</i>	The other value to compare against.
<i>delta</i>	The tolerance of how different the values can be. The larger the delta, the higher precision used in the comparison.

Returns

True if the value is less than or equal to the other value.

5.4.2.5 `template<typename T > bool int_to_bool (const T intVal) [inline]`

Converts an integral type to a boolean. Compilers complain about directly assigning an int to a bool (casting doesn't help either), so this works around that.

Parameters

<i>intVal</i>	The integer value to convert to a boolean.
---------------	--

Returns

The boolean equivalent of the integer.

5.4.2.6 `template<typename T > T safe_divide (const T dividend, const T divisor) [inline]`

Division operation that checks for division by zero or into zero (returns zero for those situations).

Parameters

<i>dividend</i>	The dividend (i.e., the value being divided).
<i>divisor</i>	The divisor (i.e., the value dividing by).

Returns

The quotient of the division operation, or zero if one of the values was invalid.

Note

If the template type has floating point precision, then the result will retain its precision.

5.4.2.7 `template<typename T > T safe_modulus (const T dividend, const T divisor) [inline]`

Modulus operation that checks for modulus by zero or into zero (returns zero for those situations).

Parameters

<i>dividend</i>	The dividend (i.e., the value being divided).
<i>divisor</i>	The divisor (i.e., the value dividing by).

Returns

The remainder of the modulus operation, or zero if one of the values was invalid.

5.5 StringOperations

Classes

- class **string_util::string_tokenize**< T >
Tokenizes a string using a set of delimiters.

5.5.1 Detailed Description

Classes for string operations.

5.6 Utilities

Utility classes.

Classes

- class **within**< T >
Determines if a value is within a given range.
- class **comparable_first_pair**< T1, T2 >
pair interface that compares on the first item
- class **backup_variable**< T >
class that remembers its original value from construction.

Macros

- #define **size_of_array**(x) (sizeof(x)/sizeof(x[0]))

Functions

- template<typename T >
T **within_range** (const T start, const T end, const T value)
- template<typename T >
bool **is_within** (const T value, const T first, const T second)
- **within**< T >::**within** (T range_begin, T range_end)
- bool **within**< T >::**operator**() (T value) const
- **comparable_first_pair**< T1, T2 >::**comparable_first_pair** (const T1 &t1, const T2 &t2)
- bool **comparable_first_pair**< T1, T2 >::**operator**< (const **comparable_first_pair**< T1, T2 > &that) const
- bool **comparable_first_pair**< T1, T2 >::**operator**== (const **comparable_first_pair**< T1, T2 > &that) const
- **backup_variable**< T >::**backup_variable** (const T &value)
- void **backup_variable**< T >::**operator**= (const T &value)
- bool **backup_variable**< T >::**operator**== (const T &value) const
- bool **backup_variable**< T >::**operator**< (const T &value) const
- bool **backup_variable**< T >::**operator**<= (const T &value) const
- bool **backup_variable**< T >::**operator**> (const T &value) const
- bool **backup_variable**< T >::**operator**>= (const T &value) const
- void **backup_variable**< T >::**operator**+ (const T &value)
- void **backup_variable**< T >::**operator**+= (const T &value)
- void **backup_variable**< T >::**operator**- (const T &value)
- void **backup_variable**< T >::**operator**-= (const T &value)
- **backup_variable**< T >::**operator** const T () const
- T * **backup_variable**< T >::**operator**& ()
- const T & **backup_variable**< T >::**get_value** () const
- T & **backup_variable**< T >::**get_value** ()
- bool **backup_variable**< T >::**has_changed** () const
- template<typename T >
bool **is_either** (const T value, const T first, const T second)
Determines if a given value is either of two other given values.
- template<typename T >
bool **is_neither** (const T value, const T first, const T second)

Determines if a given value is neither of two other given values.

- `template<typename inT , typename outT , typename member_extract_functorT >`
`outT copy_member (inT begin, inT end, outT dest, member_extract_functorT get_value)`
- `template<typename inT , typename outT , typename _Pr , typename member_extract_functorT >`
`outT copy_member_if (inT begin, inT end, outT dest, _Pr meets_criteria, member_extract_functorT get_value)`

Copies a member value between objects based on specified criteria.

5.6.1 Detailed Description

Utility classes.

Date

2003-2015

Copyright

Oleander Software, Ltd.

Author

Oleander Software, Ltd.

This program is free software; you can redistribute it and/or modify it under the terms of the BSD License.

5.6.2 Macro Definition Documentation

5.6.2.1 `#define size_of_array(x) (sizeof(x)/sizeof(x[0]))`

Returns

The item count of an array.

Note

Do not call this on an empty array. Also, this is meant for arrays of intrinsic types only.

5.6.3 Function Documentation

5.6.3.1 `template<typename inT , typename outT , typename member_extract_functorT > outT copy_member (inT begin, inT end, outT dest, member_extract_functorT get_value) [inline]`

calls a member function of elements in a container for each element in another container

5.6.3.2 `template<typename T > bool is_within (const T value, const T first, const T second) [inline]`

Returns

True if a value is within a given range.

5.6.3.3 `template<typename T > bool within< T >::operator() (T value) const [inline]`

Returns

True if value is within the valid range of values.

Parameters

<i>value</i>	The value to review.
--------------	----------------------

5.6.3.4 `template<typename T > within< T >::within (T range_begin, T range_end) [inline]`

Constructor.

Parameters

<i>range_begin</i>	The beginning of the valid range.
<i>range_end</i>	The end of the valid range.

5.6.3.5 `template<typename T > T within_range (const T start, const T end, const T value) [inline]`

Range checks a given value and truncates it if it is too high or low.

Parameters

<i>start</i>	The start of the valid range.
<i>end</i>	The end of the valid range.
<i>value</i>	The value to be range checked.

Returns

The value if within the valid range. If it was too large, then the end of the range is returned. If too low, then the start of the range is returned.

Chapter 6

Namespace Documentation

6.1 stemming Namespace Reference

Namespace for stemming classes.

Classes

- class **english_stem**
English stemmer.
- class **no_op_stem**
- class **stem**
The base class for language-specific stemmers.

Enumerations

- enum **stemming_type** {
 no_stemming, **danish**, **dutch**, **english**,
 finnish, **french**, **german**, **italian**,
 norwegian, **portuguese**, **spanish**, **swedish**,
 STEMMING_TYPE_COUNT }

Variables

- const wchar_t **UPPER_Y_HASH** = 7
- const wchar_t **LOWER_Y_HASH** = 9
- const wchar_t **UPPER_I_HASH** = 10
- const wchar_t **LOWER_I_HASH** = 11
- const wchar_t **UPPER_U_HASH** = 12
- const wchar_t **LOWER_U_HASH** = 13

6.1.1 Detailed Description

Namespace for stemming classes.

6.2 string_util Namespace Reference

Classes

- class **equal_basic_string_i_compare**
- class **equal_string_compare**
- class **equal_string_i_compare**
- class **less_basic_string_compare**
- class **less_string_compare**
- class **less_string_i_compare**
- class **less_string_n_compare**
- class **less_string_natural_order_i_compare**
- class **less_string_ni_compare**
- class **string_tokenize**
Tokenizes a string using a set of delimiters.
- class **string_trim**
trims whitespace from around a string

Functions

- wchar_t **tolower_western** (const wchar_t c)
lowercases any Western European alphabetic characters
- double **strtod** (const char *str, char **strend, int radix)
- double **strtod** (const wchar_t *str, wchar_t **strend, int radix)
- double **strtod** (const char *str, char **strend)
strtod
- double **strtod** (const wchar_t *str, wchar_t **strend)
- int **atoi** (const char *str)
atoi
- int **atoi** (const wchar_t *str)
- long **atol** (const char *str)
atol
- long **atol** (const wchar_t *str)
- int **tolower** (char c)
tolower
- wchar_t **tolower** (wchar_t c)
- int **toupper** (char c)
toupper
- wchar_t **toupper** (wchar_t c)
- template<typename T >
T * **memset** (T *dest, int c, size_t count)
memset
- char * **memset** (char *dest, int c, size_t count)
- wchar_t * **memset** (wchar_t *dest, int c, size_t count)
- const char * **strchr** (const char *s, int ch)
strchr
- const wchar_t * **strchr** (const wchar_t *s, wchar_t ch)
- const char * **strstr** (const char *s1, const char *s2)
strstr
- const wchar_t * **strstr** (const wchar_t *s1, const wchar_t *s2)
- size_t **strcspn** (const char *string1, const char *string2)

strcspn

- `size_t strcspn (const wchar_t *string1, const wchar_t *string2)`
- `char * strncat (char *strDest, const char *strSource, size_t count)`

strncat

- `wchar_t * strncat (wchar_t *strDest, const wchar_t *strSource, size_t count)`
- `int wctomb (wchar_t *s, wchar_t wc)`

wctomb

- `int wctomb (char *s, wchar_t wc)`
- `size_t strlen (const char *text)`
- `size_t strlen (const wchar_t *text)`
- `int strcmp (const char *string1, const char *string2)`

strcmp

- `int strcmp (const wchar_t *string1, const wchar_t *string2)`
- `int strncmp (const char *string1, const char *string2, size_t count)`

strncmp

- `int strncmp (const wchar_t *string1, const wchar_t *string2, size_t count)`
- `char * strncpy (char *strDest, const char *strSource, size_t count)`

strncpy

- `wchar_t * strncpy (wchar_t *strDest, const wchar_t *strSource, size_t count)`
- `template<typename charT >`
`int itoa (long value, charT *out, const size_t length)`

functions not available in ANSI C

- `template<typename T >`
`bool is_space (const T ch)`
- `template<typename T >`
`bool is_hex_digit (const T ch)`
- `template<typename T >`
`int atoi (const T *hexStr, size_t length=-1)`
- `template<typename T >`
`size_t strlen (const T *str, const size_t maxlen)`
- `template<typename T >`
`const T * stristr (const T *string, const T *strSearch)`

search for substring in string (case-insensitive)

- `template<typename T >`
`const T * strnistr (const T *string, const T *strSearch, const size_t string_len)`
- `template<typename T >`
`const T * strrstr (const T *string, const T *search, size_t offset)`
- `template<typename T >`
`int strnicmp (const T *first, const T *last, size_t count)`

Case-insensitive comparison by character count.

- `template<typename T >`
`int stricmp (const T *first, const T *last)`

Case-insensitive comparison.

- `template<typename T >`
`int strnatordcmp (const T *first_string, const T *second_string, bool case_insensitive=false)`
- `template<typename T >`
`int strnatordncasecmp (const T *a, const T *b)`

Compare, recognizing numeric strings and ignoring case.

- `template<typename T >`
`bool has_suffix (const T *text, const size_t text_length, const T *suffix, const size_t suffix_length)`
- `template<typename T >`
`const T * find_matching_close_tag (const T *string, const T openSymbol, const T closeSymbol, const bool fail_on_overlapping_open_symbol=false)`

- `template<typename T >`
`const T * find_matching_close_tag (const T *string, const T *openSymbol, const T *closeSymbol)`
Searches for a matching tag, skipping an extra open/close pairs of symbols in between.
- `template<typename T >`
`const T * strnchr (const T *string, const T ch, size_t numberOfCharacters)`
- `template<typename T >`
`const T * strcspn_pointer (const T *string1, const T *string2, const size_t string2Length)`
- `template<typename T >`
`size_t strncspn (const T *stringToSearch, const size_t stringToSearchLength, const T *searchString, const size_t searchStringLength)`
- `template<typename T >`
`size_t find_last_not_of (const T *string, const T *search, size_t offset=std::basic_string< T >::npos)`
- `template<typename T >`
`size_t find_last_of (const T *string, const T ch, size_t offset=-1)`
- `template<typename T >`
`size_t find_first_not_of (const T *stringToSearch, const size_t stringToSearchLength, const T *searchString, const size_t searchStringLength)`
- `template<typename T >`
`T remove_all_whitespace (const T &text)`
Removes all whitespace from a string.
- `template<typename Tchar_type , typename T >`
`void remove_all (T &text, Tchar_type char_to_replace)`
Removes all instances of a character from a string.
- `template<typename Tchar_type , typename T >`
`void replace_all (T &text, Tchar_type text_to_replace, Tchar_type replacement_text)`
helper functions for stemmers
- `template<typename T , typename Tchar_type >`
`void replace_all (T &text, const Tchar_type *text_to_replace, const Tchar_type *replacement_text)`
- `template<typename T >`
`void replace_all (T &text, const T &text_to_replace, const T &replacement_text)`
- `template<typename string_typeT >`
`size_t remove_extra_spaces (string_typeT &Text)`
- `template<typename string_typeT >`
`size_t remove_blank_lines (string_typeT &Text)`
- `template<typename Tchar_type >`
`double strtod_ex (const Tchar_type *nptr, Tchar_type **endptr)`
- `template<typename Tchar_type >`
`bool is_one_of (const Tchar_type character, const Tchar_type *char_string)`

6.2.1 Detailed Description

Date

2003-2015

Copyright

Oleander Software, Ltd.

Author

Oleander Software, Ltd.

This program is free software; you can redistribute it and/or modify it under the terms of the BSD License.

6.2.2 Function Documentation

6.2.2.1 `template<typename T > int string_util::atoi (const T * hexStr, size_t length = -1) [inline]`

Converts string in hex format to int. Default figures out how much of the string is a valid hex string, but passing a value to the second parameter overrides this and allows you to indicate how much of the string to try to convert.

Parameters

<i>hexStr</i>	The string to convert. How much of the string to analyze. The value -1 (the default) will tell the function to read until there are no more valid hexadecimal digits.
---------------	---

Returns

The value of the string as an integer.

6.2.2.2 `template<typename T > size_t string_util::find_last_not_of (const T * string, const T * search, size_t offset = std::basic_string<T>::npos) [inline]`

Searches for a single character not from a sequence in a string in reverse.

Parameters

<i>string</i>	The string to search in.
<i>search</i>	The sequence of characters to skip.
<i>offset</i>	Where to begin the search. If -1, then the reverse search will begin at the end of the string.

Returns

The position of where the last non-matching character is at, or -1 if it can't be found.

6.2.2.3 `template<typename T > size_t string_util::find_last_of (const T * string, const T ch, size_t offset = -1) [inline]`

Searches for the last instance of a character in a string in reverse.

Parameters

<i>string</i>	The string to search.
<i>ch</i>	The character to search for.
<i>offset</i>	The offset in the string to begin the search from. The default (-1) will begin the search at the end of the string.

Returns

The offset of the found character, or -1 if not found.

6.2.2.4 `template<typename T > const T* string_util::find_matching_close_tag (const T * string, const T openSymbol, const T closeSymbol, const bool fail_on_overlapping_open_symbol = false) [inline]`

Searches for a matching tag, skipping an extra open/close pairs of symbols in between.

Parameters

<i>openSymbol</i>	The opening symbol.
<i>closeSymbol</i>	The closing symbol that we are looking for.
<i>fail_on_overlapping_open_symbol</i>	Whether it should immediately return failure if an open symbol is found before a matching close symbol.

Returns

A pointer to where the closing tag is, or NULL if one can't be found.

6.2.2.5 `template<typename T > bool string_util::has_suffix (const T * text, const size_t text_length, const T * suffix, const size_t suffix_length) [inline]`

Indicates whether a larger strings ends with the specified suffix. Lengths are provided by the caller for efficiency. This function is case sensitive.

6.2.2.6 `template<typename T > bool string_util::is_hex_digit (const T ch) [inline]`

Determines whether a character is a hexademical digit (0-9,A-F,a-f).

Parameters

<i>ch</i>	The letter to be analyzed.
-----------	----------------------------

6.2.2.7 `template<typename T > bool string_util::is_space (const T ch) [inline]`

Determines whether a character is a space, tab, or newline. Also includes double-width and no break spaces.

Parameters

<i>ch</i>	The letter to be analyzed.
-----------	----------------------------

6.2.2.8 `template<typename string_typeT > size_t string_util::remove_blank_lines (string_typeT & Text)`

Removes blank lines from a block of text

Parameters

<i>Text</i>	The text to have blank lines removed from.
-------------	--

Returns

Number of characters (not lines) removed from the block.

6.2.2.9 `template<typename string_typeT > size_t string_util::remove_extra_spaces (string_typeT & Text)`

Strips extraneous spaces/tabs/carriage returns from a block of text so that there isn't more than one space consecutively.

6.2.2.10 `template<typename T > const T* string_util::strcspn_pointer (const T * string1, const T * string2, const size_t string2Length) [inline]`

Searches for a single character from a sequence in a string and return a pointer if found.

6.2.2.11 `template<typename T > int string_util::strnatordcmp (const T * first_string, const T * second_string, bool case_insensitive = false) [inline]`

Natural order comparison functions. Compare, recognizing numeric strings.

6.2.2.12 `template<typename T > const T* string_util::strnchr (const T * string, const T ch, size_t numberOfCharacters) [inline]`

Searches for a single character in a string for n number of characters. Size argument should be less than or equal to the length of the string being searched.

Parameters

<i>string</i>	The string to search in.
<i>ch</i>	The character to search for.
<i>numberOfCharacters</i>	The number of characters to search through in the string.

Returns

A pointer in the string where the character was found, or NULL if not found.

6.2.2.13 `template<typename T > size_t string_util::strncspn (const T * stringToSearch, const size_t stringToSearchLength, const T * searchString, const size_t searchStringLength) [inline]`

Searches for a single character from a sequence in a string for n number of characters.

Parameters

<i>stringToSearch</i>	The string to search.
<i>stringToSearchLength</i>	The length of the string being searched.
<i>searchString</i>	The sequence of characters to search for.
<i>searchStringLength</i>	The length of the sequence string.

Returns

The index into the string that the character was found. Returns the length of the string if not found.

6.2.2.14 `template<typename T> const T* string_util::strnistr (const T * string, const T * strSearch, const size_t string_len) [inline]`

Searches for substring in a larger string (case-insensitively), limiting the search to a specified number of characters.

6.2.2.15 `template<typename T> size_t string_util::strlen (const T * str, const size_t maxlen) [inline]`

Returns

The number of characters in the string pointed to by *str*, not including the terminating '\0' character, but at most *maxlen*. In doing this, *strlen* looks only at the first *maxlen* characters at *str* and never beyond *str*+*maxlen*. This function should be used for input that may not be NULL terminated.

Parameters

<i>str</i>	The string to review.
<i>maxlen</i>	The maximum length of the string to scan.

Returns

The valid length of the string or *maxlen*, whichever is shorter.

6.2.2.16 `template<typename T> const T* string_util::strrstr (const T * string, const T * search, size_t offset) [inline]`

Search string in reverse for substring. "offset" is how far we are in the source string already and how far to go back.

6.2.2.17 `template<typename Tchar_type> double string_util::strtod_ex (const Tchar_type * nptr, Tchar_type ** endptr) [inline]`

Converts strings to double values, but also takes into account ranges (returning the average). For example, a string like "5-8" will return 6.5. Hyphens and colons are seen as range separators.

6.2.2.18 `double string_util::strtol (const char * str, char ** strend, int radix) [inline]`

ANSI C decorators strtol

Chapter 7

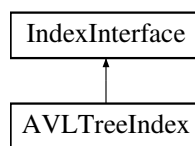
Class Documentation

7.1 AVLTreeIndex Class Reference

The **AVLTreeIndex** (p. 29) class The **AVLTreeIndex** (p. 29) class has multiple custom insert methods to allow it to either insert from inverted index or directly from parser.

```
#include <avltreeindex.h>
```

Inheritance diagram for AVLTreeIndex:



Public Member Functions

- **AVLTreeIndex** ()
AVLTreeIndex::AVLTreeIndex (p. 29) The AVLTree constructor.
- virtual vector< **document** > * **findIndex** (string key)
AVLTreeIndex::findIndex (p. 30) This finds the index of a specified word.
- virtual void **addIndex** (string &word, string &doc)
AVLTreeIndex::addIndex (p. 30) This adds the index given a word and PDF name.
- void **insert** (string x, string docname)
AVLTreeIndex::insert This inserts a word and PDF name into the AVL tree.
- void **insert** (string x, string docname, int count)
AVLTreeIndex::insert This inserts the data from the inverted file index and into the hash table.
- void **remove** (string x)
AVLTreeIndex::remove This removes a word from the AVLTree.
- virtual void **display** ()
AVLTreeIndex::display (p. 29) This displays the AVLTree.
- virtual void **readIndex** ()
AVLTreeIndex::readIndex (p. 29) This reads the inverted file index and inserts the data into the AVLTree.
- virtual void **writelIndex** ()
AVLTreeIndex::writelIndex (p. 29) This writes the AVLTree to the inverted file index .txt file.

- virtual vector< **document** > **topwords** ()
***AVLTreeIndex::topwords** (p. 31) This returns the top fifty words from the AVLTree.*
- virtual int **totalWordsIndexed** ()
***AVLTreeIndex::totalWordsIndexed** (p. 31) This returns the total words indexed.*
- vector< std::string > **split** (const std::string &s, char delim)
***AVLTreeIndex::split** (p. 31) This splits the supplied string by a given delimiter.*

7.1.1 Detailed Description

The **AVLTreeIndex** (p. 29) class The **AVLTreeIndex** (p. 29) class has multiple custom insert methods to allow it to either insert from inverted index or directly from parser.

7.1.2 Member Function Documentation

7.1.2.1 void **AVLTreeIndex::addIndex** (string & *word*, string & *docname*) [virtual]

AVLTreeIndex::addIndex (p. 30) This adds the index given a word and PDF name.

Parameters

<i>word</i>	The given word.
<i>docname</i>	The given PDF name.

Implements **IndexInterface** (p. 50).

7.1.2.2 vector< **document** > * **AVLTreeIndex::findIndex** (string *word_key*) [virtual]

AVLTreeIndex::findIndex (p. 30) This finds the index of a specified word.

Parameters

<i>word_key</i>	The specified word.
-----------------	---------------------

Returns

vector<document> The vector of document that contains all off the PDF names and counts of the specified word.

Implements **IndexInterface** (p. 50).

7.1.2.3 void **AVLTreeIndex::insert** (string *x*, string *docname*)

AVLTreeIndex::insert This inserts a word and PDF name into the AVL tree.

Parameters

<i>x</i>	The word to be inserted.
<i>docname</i>	The PDF name to be inserted.

7.1.2.4 void AVLTreeIndex::insert (string *x*, string *docname*, int *count*)

AVLTreeIndex::insert This inserts the data from the inverted file index and into the hash table.

Parameters

<i>x</i>	The word to be inserted.
<i>docname</i>	The PDF name to be inserted.
<i>count</i>	The count to be inserted.

7.1.2.5 vector< string > AVLTreeIndex::split (const std::string & *s*, char *delim*)

AVLTreeIndex::split (p. 31) This splits the supplied string by a given delimiter.

Parameters

<i>s</i>	The string to be split.
<i>delim</i>	The character to be used as the delimiter.

Returns

vector<string> The vector of delimited strings.

7.1.2.6 vector< document > AVLTreeIndex::topwords () [virtual]

AVLTreeIndex::topwords (p. 31) This returns the top fifty words from the AVLTree.

Returns

vector<document> The vector of documents objects of the top fifty words and their counts.

Implements **IndexInterface** (p. 50).

7.1.2.7 int AVLTreeIndex::totalWordsIndexed () [virtual]

AVLTreeIndex::totalWordsIndexed (p. 31) This returns the total words indexed.

Returns

int The total words indexed.

Implements **IndexInterface** (p. 50).

The documentation for this class was generated from the following files:

- avltreeindex.h
- avltreeindex.cpp

7.2 backup_variable< T > Class Template Reference

class that remembers its original value from construction.

```
#include <utilities.h>
```

Public Member Functions

- **backup_variable** (const T &value)
- void **operator=** (const T &value)
- bool **operator==** (const T &value) const
- bool **operator<** (const T &value) const
- bool **operator<=** (const T &value) const
- bool **operator>** (const T &value) const
- bool **operator>=** (const T &value) const
- void **operator+** (const T &value)
- void **operator+=** (const T &value)
- void **operator-** (const T &value)
- void **operator-=** (const T &value)
- **operator const T** () const
- T * **operator&** ()
- const T & **get_value** () const
- T & **get_value** ()
- bool **has_changed** () const

7.2.1 Detailed Description

```
template<typename T>
class backup_variable< T >
```

class that remembers its original value from construction.

The documentation for this class was generated from the following file:

- utilities.h

7.3 bookmark Struct Reference

The bookmark struct This struct uses two strings, one for the PDF name, the other for the query.

```
#include <searchengine.h>
```

Public Member Functions

- **bookmark** (string m, string q)

Public Attributes

- string **mark**
- string **query**

7.3.1 Detailed Description

The bookmark struct This struct uses two strings, one for the PDF name, the other for the query.

The documentation for this struct was generated from the following file:

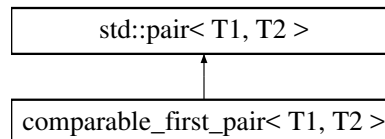
- searchengine.h

7.4 comparable_first_pair< T1, T2 > Class Template Reference

pair interface that compares on the first item

```
#include <utilities.h>
```

Inheritance diagram for comparable_first_pair< T1, T2 >:



Public Member Functions

- **comparable_first_pair** (const T1 &t1, const T2 &t2)
- bool **operator<** (const **comparable_first_pair**< T1, T2 > &that) const
- bool **operator==** (const **comparable_first_pair**< T1, T2 > &that) const

7.4.1 Detailed Description

```
template<typename T1, typename T2>
class comparable_first_pair< T1, T2 >
```

pair interface that compares on the first item

The documentation for this class was generated from the following file:

- utilities.h

7.5 document Struct Reference

The document struct This struct contains a string used as the word or the PDF name, and an int used as the count, and a double for the tf/idf. This struct is used in many different ways depending on the function it is used in.

```
#include <indexinterface.h>
```

Public Member Functions

- **document** (string name)
- **document** (string name, int count)
- string **getName** ()
- void **setName** (string name)
- void **setCount** (int num)
- int **getCount** ()

Public Attributes

- string **docname**
- int **count** =1
- double **tdif** =0

7.5.1 Detailed Description

The document struct This struct contains a string used as the word or the PDF name, and an int used as the count, and a double for the tf/idf. This struct is used in many different ways depending on the function it is used in.

The documentation for this struct was generated from the following file:

- indexinterface.h

7.6 DocumentParser Class Reference

The **DocumentParser** (p. 34) class parses the PDFs. It extracts the text using the **TextExtractor** (p. 72) class, stems the words and removes stop words using the IndexExtractor class, and computes the term frequency using the custom AVLTree or the custom Hash table.

```
#include <DocumentParser.h>
```

Public Member Functions

- **DocumentParser** (**IndexHandler** *)

The **DocumentParser** (p. 34) class parses the PDFs. It extracts the text using the **TextExtractor** (p. 72) class, stems the words and removes stop words using the **IndexExtractor** class, and computes the term frequency using the custom AVLTree or the custom Hash table.

- **DocumentParser** ()

DocumentParser::DocumentParser (p. 36) The overloaded default **DocumentParser** (p. 34) constructor.

- string **getStemmed** (string &)

DocumentParser::getStemmed (p. 36) This gets the stemmed version of the word.

- bool **isStopWord** (string &)

DocumentParser::isStopWord (p. 37) This checks if the word is a stop word.

- bool **extract** (string)

DocumentParser::extract (p. 36) This extracts the text from the PDF. This also stores the total pages and the total word indexed count in the appropriate .txt file.

- int **getTotalPages** ()

DocumentParser::getTotalPages (p. 37) This gets the total number of pages.

- void **readInWordyMap** ()

DocumentParser::readInWordyMap (p. 35) This reads the word count and associated PDF name to a .txt file.

- int **numOfDocs** ()

DocumentParser::numOfDocs (p. 37) This returns the number of documents in the corpus.

- void **clearWordTxt** ()

DocumentParser::clearWordTxt (p. 35) This clears the word count and associated PDF name from the .txt file.

- bool **rawTextExtract** (string fileStream)

DocumentParser::rawTextExtract (p. 37) This outputs the raw text from a supplied PDF.

Public Attributes

- map< string, int > **wordy**
- vector< **TextExtractor** > **m**
- vector< string > **k**

7.6.1 Detailed Description

The **DocumentParser** (p. 34) class parses the PDFs. It extracts the text using the **TextExtractor** (p. 72) class, stems the words and removes stop words using the **IndexExtractor** class, and computes the term frequency using the custom AVLTree or the custom Hash table.

CSE 2341 **DocumentParser.h** (p. ??)

Author

Aviraj Sinha (owner)
Patrick Yienger

Version

1.0 05/07/17 The **DocumentParser** (p. 34) class

7.6.2 Constructor & Destructor Documentation

7.6.2.1 `DocumentParser::DocumentParser (IndexHandler * i)`

The **DocumentParser** (p. 34) class parses the PDFs. It extracts the text using the **TextExtractor** (p. 72) class, stems the words and removes stop words using the `IndexExtractor` class, and computes the term frequency using the custom AVLTree or the custom Hash table.

CSE 2341 `DocumentParser.cpp`

Author

Aviraj Sinha (owner)
Patrick Yienger

Version

1.0 05/07/17 **DocumentParser::DocumentParser** (p. 36) The **DocumentParser** (p. 34) constructor.

Parameters

<i>i</i>	The IndexHandler (p. 48) object.
----------	---

7.6.3 Member Function Documentation

7.6.3.1 `bool DocumentParser::extract (string fileStream)`

DocumentParser::extract (p. 36) This extracts the text from the PDF. This also stores the total pages and the total word indexed count in the appropriate .txt file.

Parameters

<i>fileStream</i>	The file path to the corpus of PDFs.
-------------------	--------------------------------------

Returns

`bool` The flag for a successful text extract.

7.6.3.2 `string DocumentParser::getStemmed (string & word)`

DocumentParser::getStemmed (p. 36) This gets the stemmed version of the word.

Parameters

<i>word</i>	The word to be stemmed
-------------	------------------------

Returns

string The stemmed word

7.6.3.3 int DocumentParser::getTotalPages ()

DocumentParser::getTotalPages (p. 37) This gets the total number of pages.

Returns

totalPages The total number of pages.

7.6.3.4 bool DocumentParser::isStopWord (string & word)

DocumentParser::isStopWord (p. 37) This checks if the word is a stop word.

Parameters

<i>word</i>	The word to be checked
-------------	------------------------

Returns

bool The flag of whether the word is a stop word

7.6.3.5 int DocumentParser::numOfDocs ()

DocumentParser::numOfDocs (p. 37) This returns the number of documents in the corpus.

Returns

int The number of documents in the corpus.

7.6.3.6 bool DocumentParser::rawTextExtract (string *fileStream*)

DocumentParser::rawTextExtract (p. 37) This outputs the raw text from a supplied PDF.

Parameters

<i>fileStream</i>	The file path to the supplied PDF.
-------------------	------------------------------------

Returns

bool The flag for a successful text extract.

The documentation for this class was generated from the following files:

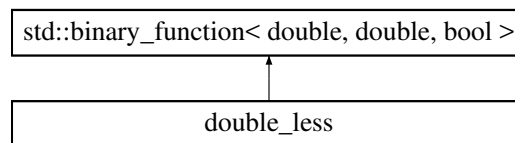
- DocumentParser.h
- DocumentParser.cpp

7.7 `double_less` Class Reference

"less" interface for double values.

```
#include <safe_math.h>
```

Inheritance diagram for `double_less`:



Public Member Functions

- `bool operator()` (`const double &left`, `const double &right`) `const`

7.7.1 Detailed Description

"less" interface for double values.

The documentation for this class was generated from the following file:

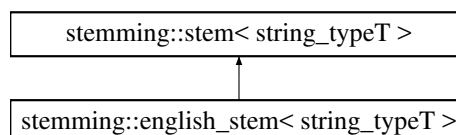
- `safe_math.h`

7.8 `stemming::english_stem< string_typeT >` Class Template Reference

English stemmer.

```
#include <english_stem.h>
```

Inheritance diagram for `stemming::english_stem< string_typeT >`:



Public Member Functions

- `void operator()` (`string_typeT &text`)

Additional Inherited Members

7.8.1 Detailed Description

```
template<typename string_typeT = std::wstring>
class stemming::english_stem< string_typeT >
```

English stemmer.

Overview:

I have made more than one attempt to improve the structure of the Porter algorithm by making it follow the pattern of ending removal of the Romance language stemmers. It is not hard to see why one should want to do this: step 1b of the Porter stemmer removes ed and ing, which are i-suffixes (*) attached to verbs. If these suffixes are removed, there should be no need to remove d-suffixes which are not verbal, although it will try to do so. This seems to be a deficiency in the Porter stemmer, not shared by the Romance stemmers. Again, the divisions between steps 2, 3 and 4 seem rather arbitrary, and are not found in the Romance stemmers. Nevertheless, these attempts at improvement have been abandoned. They seem to lead to a more complicated algorithm with no very obvious improvements. A reason for not taking note of the outcome of step 1b may be that English endings do not determine word categories quite as strongly as endings in the Romance languages. For example, condition and position in French have to be nouns, but in English they can be verbs as well as nouns. We are all conditioned by advertising. They are positioning themselves differently today. A possible reason for having separate steps 2, 3 and 4 is that d-suffix combinations in English are quite complex, a point which has been made elsewhere. But it is hardly surprising that after twenty years of use of the Porter stemmer, certain improvements do suggest themselves, and a new algorithm for English is therefore offered here. (It could be called the 'Porter2' stemmer to distinguish it from the Porter stemmer, from which it derives.) The changes are not so very extensive: (1) terminating y is changed to i rather less often, (2) suffix us does not lose its s, (3) a few additional suffixes are included for removal, including (4) suffix ly. In addition, a small list of exceptional forms is included. In December 2001 there were two further adjustments: (5) Steps 5a and 5b of the old Porter stemmer were combined into a single step. This means that undoubling final ll is not done with removal of final e. (6) In Step 3 ative is removed only when in region R2. To begin with, here is the basic algorithm without reference to the exceptional forms. An exact comparison with the Porter algorithm needs to be done quite carefully if done at all. Here we indicate by * points of departure, and by + additional features. In the sample vocabulary, Porter and Porter2 stem slightly under 5% of words to different forms. Dr. Martin Porter Define a vowel as one of

- a e i o u y Define a double as one of
- bb dd ff gg mm nn pp rr tt Define a valid li-ending as one of
- c d e g h k m n r t Define a short syllable in a word as either (a) a vowel followed by a non-vowel other than w, x or Y and preceded by a non-vowel, or * (b) a vowel at the beginning of the word followed by a non-vowel. So rap, trap, entrap end with a short syllable, and ow, on, at are classed as short syllables. But uproot, bestow, disturb do not end with a short syllable. A word is called short if it consists of a short syllable preceded by zero or more consonants. R1 is the region after the first non-vowel following a vowel, or the end of the word if there is no such non-vowel. R2 is the region after the first non-vowel following a vowel in R1, or the end of the word if there is no such non-vowel. If the word has two letters or less, leave it as it is. Otherwise, do each of the following operations, Set initial y, or y after a vowel, to Y, and then establish the regions R1 and R2.

Algorithm:

Step 1a: Search for the longest among the following suffixes, and perform the action indicated:

- sses
 - Replace by ss.
- ied+ ies*
 - Replace by i if preceded by just one letter, otherwise by ie (so ties -> tie, cries -> cri).

- s
 - Delete if the preceding word part contains a vowel not immediately before the s (so gas and this retain the s, gaps and kiwis lose it).
- us+ ss
 - Do nothing. **Step 1b:** Search for the longest among the following suffixes, and perform the action indicated:
- eed eedly+
 - Replace by ee if in R1.
- ed edly+ ing ingly+
 - Delete if the preceding word part contains a vowel, and then
 - If the word ends at, bl or iz add e (so luxuriat -> luxuriate), or
 - If the word ends with a double remove the last letter (so hopp -> hop), or
 - If the word is short, add e (so hop -> hope). **Step 1c:** Replace suffix y or Y by i if preceded by a non-vowel which is not the first letter of the word (so cry -> cri, by -> by, say -> say) **Step 2:** Search for the longest among the following suffixes, and, if found and in R1, perform the action indicated:
- tional
 - Replace by tion.
- enci
 - Replace by ence.
- anci
 - Replace by ance
- abli
 - Replace by able.
- entli
 - Replace by ent.
- izer ization
 - Replace by ize.
- ational ation ator
 - Replace by ate.
- alism aliti alli
 - Replace by al.
- fulness
 - Replace by ful.
- ousli ousness
 - Replace by ous.
- iveness iviti
 - Replace by ive.
- biliti bli+
 - Replace by ble.
- ogi+
 - Replace by og if preceded by l.
- fulli+
 - Replace by ful.
- lessli+
 - Replace by less.
- li+

- Delete if preceded by a valid li-ending. **Step 3:** Search for the longest among the following suffixes, and, if found and in R1, perform the action indicated:
 - tional+
 - Replace by tion.
 - ational+
 - Replace by ate.
 - alize
 - Replace by al.
 - icate iciti ical
 - Replace by ic.
 - ful ness
 - Delete.
 - ative*
 - Delete if in R2. **Step 4:** Search for the longest among the following suffixes, and, if found and in R2, perform the action indicated:
 - al ance ence er ic able ible ant ement ment ent ism ate iti ous ive ize
 - Delete
 - ion
 - Delete if preceded by s or t. **Step 5:** Search for the following suffixes, and, if found, perform the action indicated:
 - e
 - Delete if in R2, or in R1 and not preceded by a short syllable.
 - l
 - Delete if in R2 and preceded by l.

7.8.2 Member Function Documentation

7.8.2.1 `template<typename string_typeT = std::wstring> void stemming::english_stem< string_typeT >::operator() (string_typeT & text) [inline]`

Parameters

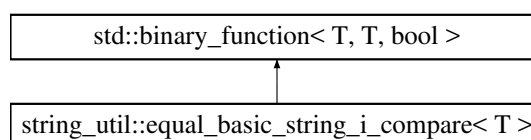
<i>in, out</i>	<i>text</i>	English string to stem.
----------------	-------------	-------------------------

The documentation for this class was generated from the following file:

- english_stem.h

7.9 string_util::equal_basic_string_i_compare< T > Class Template Reference

Inheritance diagram for string_util::equal_basic_string_i_compare< T >:



Public Member Functions

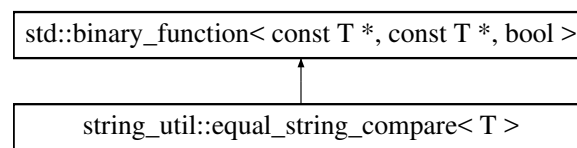
- bool **operator()** (const T &a_, const T &b_) const

The documentation for this class was generated from the following file:

- string_util.h

7.10 string_util::equal_string_compare< T > Class Template Reference

Inheritance diagram for string_util::equal_string_compare< T >:



Public Member Functions

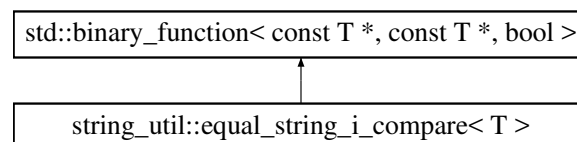
- bool **operator()** (const T *a_, const T *b_) const

The documentation for this class was generated from the following file:

- string_util.h

7.11 string_util::equal_string_i_compare< T > Class Template Reference

Inheritance diagram for string_util::equal_string_i_compare< T >:



Public Member Functions

- bool **operator()** (const T *a_, const T *b_) const

The documentation for this class was generated from the following file:

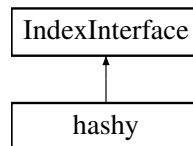
- string_util.h

7.12 hashy Class Reference

The hashy class The hashy class has multiple custom insert methods to allow it to either insert from inverted index or directly from parser.

```
#include <hash.h>
```

Inheritance diagram for hashy:



Public Member Functions

- **hashy ()**
The hashy class is the custom implementation of the HashTable. The basic hash parts of the code was based off of the series of videos by Paul Programming.
- **int Hashy (string key)**
***hashy::Hashy** (p. 45) The hash function that returns the int value of the supplied string. This function is based off of the hash function for strings developed by Paul Larson.*
- **void addIndex (string &word, string &name)**
***hashy::addIndex** (p. 44) This adds the supplied word and name to the hash table. It also updates the count.*
- **int NumItemsInIndex (int index)**
***hashy::NumItemsInIndex** (p. 45) This returns the number of items in a specified index.*
- **void display ()**
***hashy::display** (p. 43) This displays the contents of the hash table.*
- **void displayIndex (int index)**
***hashy::displayIndex** (p. 44) This displays the contents of the hash table at a specified index.*
- **vector< document > * findIndex (string word)**
***hashy::findIndex** (p. 44) This finds the index of a specified word.*
- **virtual void readIndex ()**
***hashy::readIndex** (p. 43) This reads the data from the inverted file index and into the hash table.*
- **virtual void writeIndex ()**
***hashy::writeIndex** (p. 43) This writes the inverted file index from the hash table and into the index .txt.*
- **virtual void insertFromFile (string &word, string &doc, int &count)**
***hashy::insertFromFile** (p. 45) This inserts the data from the inverted file index and into the hash table.*
- **vector< std::string > split (const std::string &s, char delim)**
***hashy::split** (p. 45) This splits the supplied string by a given delimiter.*
- **virtual vector< document > topwords ()**
***hashy::topwords** (p. 46) This finds and returns the top fifty words and their counts.*
- **virtual int totalWordsIndexed ()**
***hashy::totalWordsIndexed** (p. 46) This returns the total words indexed.*

7.12.1 Detailed Description

The hashy class The hashy class has multiple custom insert methods to allow it to either insert from inverted index or directly from parser.

7.12.2 Constructor & Destructor Documentation

7.12.2.1 hashy::hashy ()

The hashy class is the custom implementation of the HashTable. The basic hash parts of the code was based off of the series of videos by Paul Programming.

CSE 2341 hash.cpp

Author

Parick Yienger (owner)

Version

1.0 05/07/17 **hashy::hashy** (p. 44) This is the hashy constructor. It sets all of the elements of the HashTable to default values.

7.12.3 Member Function Documentation

7.12.3.1 void hashy::addIndex (string & word, string & name) [virtual]

hashy::addIndex (p. 44) This adds the supplied word and name to the hash table. It also updates the count.

Parameters

<i>word</i>	The word to be added.
<i>name</i>	The pdf the word came from.

Implements **IndexInterface** (p. 50).

7.12.3.2 void hashy::displayIndex (int index)

hashy::displayIndex (p. 44) This displays the cotents oh the hash table at a specified index.

Parameters

<i>index</i>	The specified index.
--------------	----------------------

7.12.3.3 vector< document > * hashy::findIndex (string word) [virtual]

hashy::findIndex (p. 44) This finds the index of a specified word.

Parameters

<i>word</i>	The specified word.
-------------	---------------------

Returns

vector<document> The vector of document that contains all off the PDF names and counts of the specified word.

Implements **IndexInterface** (p. 50).

7.12.3.4 int hashy::Hashy (string *key*)

hashy::Hashy (p. 45) The hash function that returns the int value of the supplied string. This function is based off of the hash function for strings developed by Paul Larson.

Parameters

<i>key</i>	The string to be hashed.
------------	--------------------------

Returns

int The integer hash value of the string.

7.12.3.5 void hashy::insertFromFile (string & *word*, string & *doc*, int & *count*) [virtual]

hashy::insertFromFile (p. 45) This inserts the data from the inverted file index and into the hash table.

Parameters

<i>word</i>	The word to be inserted.
<i>doc</i>	The PDF name to be inserted.
<i>count</i>	The count to be inserted.

7.12.3.6 int hashy::NumItemsInIndex (int *index*)

hashy::NumItemsInIndex (p. 45) This returns the number of items in a specified index.

Parameters

<i>index</i>	The specified index.
--------------	----------------------

Returns

int The number of words in the index.

7.12.3.7 vector< string > hashy::split (const std::string & *s*, char *delim*)

hashy::split (p. 45) This splits the supplied string by a given delimiter.

Parameters

<i>s</i>	The string to be split.
<i>delim</i>	The character to be used as the delimiter.

Returns

vector<string> The vector of delimited strings.

7.12.3.8 vector< document > hashy::topwords () [virtual]

hashy::topwords (p. 46) This finds and returns the top fifty words and their counts.

Returns

vector<document> The vector of documents objects of the top fifty words and their counts.

Implements **IndexInterface** (p. 50).

7.12.3.9 int hashy::totalWordsIndexed () [virtual]

hashy::totalWordsIndexed (p. 46) This returns the total words indexed.

Returns

The total words indexed.

Implements **IndexInterface** (p. 50).

The documentation for this class was generated from the following files:

- hash.h
- hash.cpp

7.13 index_node Struct Reference

The **index_node** (p. 46) struct allows for storage of inverted indices with other data in tree; most importantly the documents it is found in.

```
#include <avltreeindex.h>
```

Public Attributes

- string **word_key**
- **index_node** * **left**
- **index_node** * **right**
- vector< **document** > **documents**
- int **height**

7.13.1 Detailed Description

The **index_node** (p. 46) struct allows for storage of inverted indices with other data in tree; most importantly the documents it is found in.

The documentation for this struct was generated from the following file:

- avltreeindex.h

7.14 indexextractor Class Reference

The indexextractor class.

```
#include <indexextractor.h>
```

Public Member Functions

- **indexextractor** (string)
The IndexExtractor class stems the words and removes stopwords. The code makes use of the
- void **useStopWords** (string)
indexextractor::useStopWords (p. 47) adds the stop words from a file
- bool **isStopWord** (string &)
indexextractor::isStopWord (p. 47) This checks if it is in the stop words set
- string **getStemmed** (string &)
indexextractor::getStemmed (p. 47) This gets the stemmed version of the word

7.14.1 Detailed Description

The indexextractor class.

7.14.2 Constructor & Destructor Documentation

7.14.2.1 indexextractor::indexextractor (string *stoptxt*)

The IndexExtractor class stems the words and removes stopwords. The code makes use of the

CSE 2341 IndexExtractor.cpp

Author

Aviraj Shina (owner)

Version

1.0 05/07/17 **indexextractor::indexextractor** (p. 47) This class is responsible for stemming and stop words

The documentation for this class was generated from the following files:

- indexextractor.h
- indexextractor.cpp

7.15 IndexHandler Class Reference

The **IndexHandler** (p. 48) class.

```
#include <indexhandler.h>
```

Public Member Functions

- **IndexHandler** (char type)
*The **IndexHandler** (p. 48) class creates and stores the inverted index file. It uses the custom data structures to write and read the inverted index.*
- void **writelIndex** ()
***IndexHandler::writelIndex** (p. 48) This writes the inverted index file to the .txt file from the appropriate data structure.*
- void **addIndex** (string, string)
***IndexHandler::addIndex** (p. 49) This adds the word and the PDf name to the chosen data structure.*
- vector< **document** > * **getDocs** (string)
***IndexHandler::getDocs** (p. 49) This returns the name of all of the PDFs containing the requested word.*
- int **getDocumentFrequency** (string word, string docname)
- int **getCorpusFrequency** (string word)
- void **displayIndices** ()
***IndexHandler::displayIndices** (p. 48) This displays the indices of the appropriate data structure.*
- void **readIndex** ()
***IndexHandler::readIndex** (p. 48) This reads the data from the inverted index file and into the appropriate data structure.*
- vector< **document** > **topFifty** ()
***IndexHandler::topFifty** (p. 49) This returns the top fifty words and their count.*
- int **totalWordsIndexed** ()
***IndexHandler::totalWordsIndexed** (p. 49) This returns the total number of words indexed.*

7.15.1 Detailed Description

The **IndexHandler** (p. 48) class.

7.15.2 Constructor & Destructor Documentation

7.15.2.1 IndexHandler::IndexHandler (char type)

The **IndexHandler** (p. 48) class creates and stores the inverted index file. It uses the custom data structures to write and read the inverted index.

CSE 2341 IndexHandler.h

Author

Aviraj Shina (owner)
Patrick Yienger

Version

1.0 05/07/17 **IndexHandler::IndexHandler** (p. 48) The index handler constructor. This initialaizes the chosen data structure.

Parameters

<i>type</i>	The data structure to be initialized.
-------------	---------------------------------------

7.15.3 Member Function Documentation

7.15.3.1 void IndexHandler::addIndex (string word, string docname)

IndexHandler::addIndex (p. 49) This adds the word and the PDF name to the chosen data structure.

Parameters

<i>word</i>	The word to be added.
<i>docname</i>	The PDF name to be added.

7.15.3.2 vector< document > * IndexHandler::getDocs (string word)

IndexHandler::getDocs (p. 49) This returns the name of all of the PDFs containing the requested word.

Parameters

<i>word</i>	The word requested.
-------------	---------------------

Returns

vector<document> The vector of document object that stores the count and PDF name of the requested word.

7.15.3.3 vector< document > IndexHandler::topFifty ()

IndexHandler::topFifty (p. 49) This returns the top fifty words and their count.

Returns

vector<document> This returns the vector of document objects containing the word and the count.

7.15.3.4 int IndexHandler::totalWordsIndexed ()

IndexHandler::totalWordsIndexed (p. 49) This returns the total number of words indexed.

Returns

totalWordsIndexed

The documentation for this class was generated from the following files:

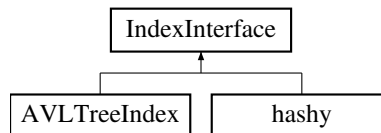
- indexhandler.h
- indexhandler.cpp

7.16 IndexInterface Class Reference

The **IndexInterface** (p. 50) class.

```
#include <indexinterface.h>
```

Inheritance diagram for IndexInterface:



Public Member Functions

- virtual vector< **document** > * **findIndex** (string key)=0
- virtual void **addIndex** (string &word, string &doc)=0
- virtual void **display** ()=0
- virtual void **readIndex** ()=0
- virtual void **writeIndex** ()=0
- virtual vector< **document** > **topwords** ()=0
- virtual int **totalWordsIndexed** ()=0

7.16.1 Detailed Description

The **IndexInterface** (p. 50) class.

The documentation for this class was generated from the following file:

- indexinterface.h

7.17 item Struct Reference

The hashy class is the custom implementation of the HashTable. The basic hash parts of the code was based off of the series of videos by Paul Programming.

```
#include <hash.h>
```

Public Attributes

- string **word_key**
- int **count**
- string **docName**
- vector< **item** > **vec**

7.17.1 Detailed Description

The hashy class is the custom implementation of the HashTable. The basic hash parts of the code was based off of the series of videos by Paul Programming.

CSE 2341 **hash.h** (p. ??)

Author

Parick Yienger (owner)

Version

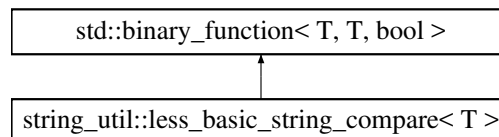
1.0 05/07/17 The item struct This is a struct that contains the word, count, and PDFName. It also uses a vector in the top index on the hash table to handle collisions.

The documentation for this struct was generated from the following file:

- hash.h

7.18 string_util::less_basic_string_compare< T > Class Template Reference

Inheritance diagram for string_util::less_basic_string_compare< T >:



Public Member Functions

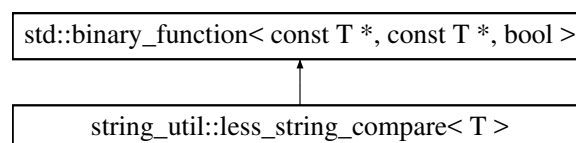
- bool **operator()** (const T &a_, const T &b_) const

The documentation for this class was generated from the following file:

- string_util.h

7.19 string_util::less_string_compare< T > Class Template Reference

Inheritance diagram for string_util::less_string_compare< T >:



Public Member Functions

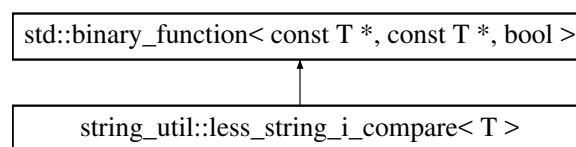
- bool **operator()** (const T *a_, const T *b_) const

The documentation for this class was generated from the following file:

- string_util.h

7.20 string_util::less_string_i_compare< T > Class Template Reference

Inheritance diagram for string_util::less_string_i_compare< T >:



Public Member Functions

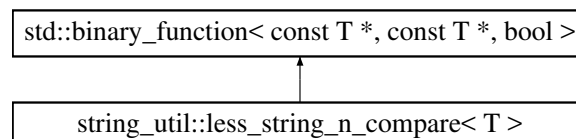
- bool **operator()** (const T *a_, const T *b_) const

The documentation for this class was generated from the following file:

- string_util.h

7.21 string_util::less_string_n_compare< T > Class Template Reference

Inheritance diagram for string_util::less_string_n_compare< T >:



Public Member Functions

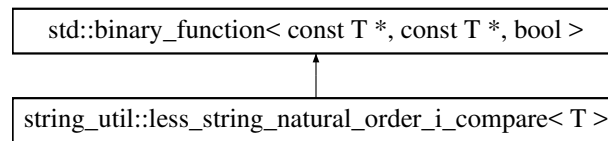
- **less_string_n_compare** (size_t comparison_size)
- bool **operator()** (const T *a_, const T *b_) const

The documentation for this class was generated from the following file:

- string_util.h

7.22 string_util::less_string_natural_order_i_compare< T > Class Template Reference

Inheritance diagram for string_util::less_string_natural_order_i_compare< T >:



Public Member Functions

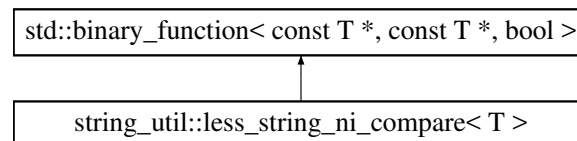
- bool **operator()** (const T *a_, const T *b_) const

The documentation for this class was generated from the following file:

- string_util.h

7.23 string_util::less_string_ni_compare< T > Class Template Reference

Inheritance diagram for string_util::less_string_ni_compare< T >:



Public Member Functions

- **less_string_ni_compare** (size_t comparison_size)
- bool **operator()** (const T *a_, const T *b_) const

The documentation for this class was generated from the following file:

- string_util.h

7.24 stemming::no_op_stem< string_typeT > Class Template Reference

Public Member Functions

- void **operator()** (const string_typeT &) const
No-op stemming of declared string type.
- template<typename T >
void **operator()** (const T &) const
No-op stemming of flexible string type.

The documentation for this class was generated from the following file:

- stemming.h

7.25 QueryEngine Class Reference

The **QueryEngine** (p. 54) uses the inverted file index and searches it given a properly formatted Boolean query. It also ranks the results using tf/idf statistics.

```
#include <queryengine.h>
```

Public Member Functions

- **QueryEngine** (**IndexHandler** *, **DocumentParser** *)

The **QueryEngine** (p. 54) uses the inverted file index and searches it given a properly formatted Boolean query. It also ranks the results using tf/idf statistics.

- vector< **document** > **querySearch** (string query)

QueryEngine::querySearch (p. 56) This searches the inverted index file for a given boolean query. It calls the search method to search for the words and returning the proper results based on whether a union, intersection or difference was required.

- vector< **document** > **search** (string word)

QueryEngine::search (p. 57) This searches the inverted index file for a word, returning the results.

- vector< **document** > **getIntersection** (vector< **document** > l, vector< **document** > r)

QueryEngine::getIntersection (p. 56) This returns the intersection of two queries.

- vector< **document** > **getUnion** (vector< **document** > l, vector< **document** > r)

QueryEngine::getUnion (p. 56) This returns the union of two queries.

- vector< **document** > **getDifference** (vector< **document** > r, vector< **document** > l)

QueryEngine::getDifference (p. 55) This returns the difference of two queries.

- void **findTDIFs** (vector< **document** > &currdocs)

QueryEngine::findTDIFs (p. 55) This calculates the tf/idf values of the query results.

- void **relevancySort** (vector< **document** > &currdocs)

QueryEngine::relevancySort (p. 56) This sorts the vector<document> by tf/idf value.

Public Attributes

- int **fileSize** =0

7.25.1 Detailed Description

The **QueryEngine** (p. 54) uses the inverted file index and searches it given a properly formatted Boolean query. It also ranks the results using tf/idf statistics.

CSE 2341 QueryEngine.h

Author

Aviraj Shina (owner)

Version

1.0 05/07/17 The **QueryEngine** (p. 54) class

7.25.2 Constructor & Destructor Documentation

7.25.2.1 QueryEngine::QueryEngine (IndexHandler * *i*, DocumentParser * *d*)

The **QueryEngine** (p. 54) uses the inverted file index and searches it given a properly formatted Boolean query. It also ranks the results using tf/idf statistics.

CSE 2341 QueryEngine.h

Author

Aviraj Shina (owner)

Version

1.0 05/07/17 **QueryEngine::QueryEngine** (p. 55) The constructor for the query engine.

Parameters

<i>i</i>	The index handler object.
<i>d</i>	The documen parser object.

7.25.3 Member Function Documentation

7.25.3.1 void QueryEngine::findTDIFs (vector< document > & *currdocs*)

QueryEngine::findTDIFs (p. 55) This caculates the tf/idf values of the query results.

Parameters

<i>currdocs</i>	The vector<document> to caculate the tf/idf for.
-----------------	--

7.25.3.2 vector< document > QueryEngine::getDifference (vector< document > *r*, vector< document > *l*)

QueryEngine::getDifference (p. 55) This returns the difference of two queries.

Parameters

<i>l</i>	The vector<document> of the first query.
<i>r</i>	The vector<document> of the secnd query.

Returns

vector<document> The difference of the first and second query.

7.25.3.3 `vector< document > QueryEngine::getIntersection (vector< document > l, vector< document > r)`

QueryEngine::getIntersection (p. 56) This returns the intersection of two queries.

Parameters

<i>l</i>	The vector<document> of the first query.
<i>r</i>	The vector<document> of the second query.

Returns

vector<document> The intersection of the first and second query.

7.25.3.4 `vector< document > QueryEngine::getUnion (vector< document > l, vector< document > r)`

QueryEngine::getUnion (p. 56) This returns the union of two queries.

Parameters

<i>l</i>	The vector<document> of the first query.
<i>r</i>	The vector<document> of the second query.

Returns

vector<document> The union of the first and second query.

7.25.3.5 `vector< document > QueryEngine::querySearch (string query)`

QueryEngine::querySearch (p. 56) This searches the inverted index file for a given boolean query. It calls the search method to search for the words and returning the proper results based on whether a union, intersection or difference was required.

Parameters

<i>query</i>	The query to search for.
--------------	--------------------------

Returns

vector<document> A vector of document objects containing the results of the search sorted by tf/idf values.

7.25.3.6 `void QueryEngine::relevancySort (vector< document > & currdocs)`

QueryEngine::relevancySort (p. 56) This sorts the vector<document> by tf/idf value.

Parameters

<i>currdocs</i>	The vector<document> to be sorted.
-----------------	------------------------------------

7.25.3.7 vector< document > QueryEngine::search (string word)

QueryEngine::search (p. 57) This searches the inverted index file for a word, returning the results.

Parameters

<i>word</i>	The word to search for.
-------------	-------------------------

Returns

vector<document> A vector of document objects containing the results of the search.

The documentation for this class was generated from the following files:

- queryengine.h
- queryengine.cpp

7.26 rawOutputExtractor Class Reference

The **rawOutputExtractor** (p. 57) class.

```
#include <rawoutputextractor.h>
```

Public Member Functions

- **rawOutputExtractor** (IndexHandler *, indexextractor *)
rawOutputExtractor::rawOutputExtractor (p. 58) The **rawOutputExtractor** (p. 57) constructor.
- **rawOutputExtractor** ()
rawOutputExtractor::rawOutputExtractor (p. 58) The **rawOutputExtractor** (p. 57) constructor.
- virtual ~**rawOutputExtractor** ()
rawOutputExtractor::~~rawOutputExtractor (p. 57) The **rawOutputExtractor** (p. 57) destructor.
- int **getTotalPages** ()
rawOutputExtractor::getTotalPages (p. 58) Returns the total pages.
- string **getDocName** ()
rawOutputExtractor::getDocName (p. 58) Returns the PDF name.
- int **getWordCount** ()
rawOutputExtractor::getWordCount (p. 59) Returns the word count.
- void **Init** (const char *pszInput)
rawOutputExtractor::Init (p. 59) Uses Podofo to extract the text from the given PDF.

Public Attributes

- **IndexHandler** * **ih**
- **indexextractor** * **ie**
- string **mm**
- vector< string > **a**
- int **wordCount** =0
- string **docsName**

7.26.1 Detailed Description

The **rawOutputExtractor** (p. 57) class.

7.26.2 Constructor & Destructor Documentation

7.26.2.1 **rawOutputExtractor::rawOutputExtractor** (**IndexHandler** * *ih*, **indexextractor** * *ie*)

rawOutputExtractor::rawOutputExtractor (p. 58) The **rawOutputExtractor** (p. 57) constructor.

Parameters

<i>ih</i>	The IndexHandler (p. 48) object.
<i>ie</i>	The IndexExtractor object.

7.26.3 Member Function Documentation

7.26.3.1 string **rawOutputExtractor::getDocName** ()

rawOutputExtractor::getDocName (p. 58) Returns the PDF name.

Returns

docsName The PDF name.

7.26.3.2 int **rawOutputExtractor::getTotalPages** ()

rawOutputExtractor::getTotalPages (p. 58) Returns the total pages.

Returns

totalPages

7.26.3.3 int rawOutputExtractor::getWordCount ()

rawOutputExtractor::getWordCount (p. 59) Returns the word count.

Returns

wordCount The word count.

7.26.3.4 void rawOutputExtractor::Init (const char * pszInput)

rawOutputExtractor::Init (p. 59) Uses Podofo to extract the text from the given PDF.

Parameters

<i>pszInput</i>	The path to the corpus of PDFs
<i>docName</i>	The name of the pdf to extract text from.

The documentation for this class was generated from the following files:

- rawoutputextractor.h
- rawoutputextractor.cpp

7.27 SearchEngine Class Reference

The **SearchEngine** (p. 59) class.

```
#include <searchengine.h>
```

Public Member Functions

- **SearchEngine** (string docpath, char)
- **SearchEngine** ()

*The **SearchEngine** (p. 59) class uses the various class to create a PDF Search Engine. The class makes use of the **DocumentParser** (p. 34), **IndexHandler** (p. 48), and **QueryEngie** class to perform all of the various tasks of the Mantience mode and Query mode of the user interface.*
- vector< **document** > **display_search_results** (string word)

***SearchEngine::display_search_results** (p. 62) This searches the inverted file index for a query and displays the results.*
- void **relevencySort** (vector< **document** > &)
- void **writelIndex** ()

***SearchEngine::writelIndex** (p. 59) This writes the index from the appropriate data structure and into the inverted file index.*
- void **readIndex** ()

***SearchEngine::readIndex** (p. 59) This reads the index from the inverted index file and into the appropriate data structure.*
- void **clearIndex** ()

***SearchEngine::clearIndex** (p. 59) This clears the inverted file index (.txt file).*

- int **getTotalPages** ()
SearchEngine::getTotalPages (p. 63) Returns the total page count.
- int **numWordsIndexed** ()
SearchEngine::numWordsIndexed (p. 63) Returns the total indexed word count.
- vector< **document** > **topfifty** ()
SearchEngine::topfifty (p. 64) Returns the top fifty most common words.
- bool **addDocumentsToIndex** (string docpath)
SearchEngine::addDocumentsToIndex (p. 61) Allows the user to choose between using a Hash table or AVL Tree to create the inverted index file. This also initializes the **IndexHandler** (p. 48), **Documentparser**, and **QueryEngine** (p. 54) objects to allow the Search Engine access to their functions. It also updates the tf/idf, inverted file index, page count, word count, and top fifty words.
- void **chooseStructure** (char type)
SearchEngine::chooseStructure (p. 62) Allows the user to choose between using a Hash table or AVL Tree to create the inverted index file. This also initializes the **IndexHandler** (p. 48), **Documentparser**, and **QueryEngine** (p. 54) objects to allow the Search Engine access to their functions. It also updates the tf/idf, inverted file index, page count, word count, and top fifty words.
- void **findTDIFs** (vector< **document** > &currdocs)
- int **readTotalPages** ()
SearchEngine::readTotalPages (p. 63) This outputs the total number of pages stored from the PDFs on the inverted index.
- void **clearTotalPages** ()
SearchEngine::clearTotalPages (p. 60) This clears the total number of pages .txt file.
- void **clearWordTxt** ()
SearchEngine::clearWordTxt (p. 60) This clears the total indexed word count .txt file.
- void **clear** ()
SearchEngine::clear (p. 60) Allows the user to clear all of the stored data. It clears the inverted file index, the total page count, the total indexed word count, the file paths, the bookmarks and the search history.
- bool **displayRawFile** (string filePath)
SearchEngine::displayRawFile (p. 63) This displays the raw text of a requested PDF.
- void **displayTop50** ()
SearchEngine::displayTop50 (p. 60) This displays the top fifty words from the PDFs on the inverted index.
- void **writeFilePathToTXTFile** ()
SearchEngine::writeFilePathToTXTFile (p. 60) This writes the file paths to a text file allowing the path to be stored.
- vector< string > **readFilePathFromTXTFile** ()
SearchEngine::readFilePathFromTXTFile (p. 63) This reads the file paths from the .txt file of file paths.
- void **clearFilePath** ()
SearchEngine::clearFilePath (p. 60) This clears the file path .txt file.
- void **writeBookmark** (string book, string query)
- void **readBookmarks** ()
SearchEngine::readBookmarks (p. 60) This reads the bookmarks from the .txt file and adds them to the bookmark vector.
- void **addBookmark** (string book, string query)
SearchEngine::addBookmark (p. 61) This adds bookmarks to the .txt file.
- void **addToHistory** (string query)
SearchEngine::addToHistory (p. 62) This adds the query to the search history.
- void **clearHistory** ()
SearchEngine::clearHistory (p. 60) This clears the history .txt file.
- void **displayHistory** ()
SearchEngine::displayHistory (p. 60) This displays the search history.
- void **clearBookmarks** ()
SearchEngine::clearBookmarks This clears the bookmarks .txt file.
- vector< **bookmark** > **displayBookmarks** ()
SearchEngine::displayBookmarks (p. 62) This displays the bookmarks.

Public Attributes

- int **totalPages** =0

7.27.1 Detailed Description

The **SearchEngine** (p. 59) class.

7.27.2 Constructor & Destructor Documentation

7.27.2.1 SearchEngine::SearchEngine ()

The **SearchEngine** (p. 59) class uses the various class to create a PDF Search Engine. The class makes use of the **DocumentParser** (p. 34), **IndexHandler** (p. 48), and **QueryEngine** class to perform all of the various tasks of the Maintenance mode and Query mode of the user interface.

CSE 2341 SearchEngine.cpp

Author

Aviraj Shina (owner)
Patrick Yienger

Version

1.0 05/07/17 SearchEngine::SearchEngine The search engine constructor.

7.27.3 Member Function Documentation

7.27.3.1 void SearchEngine::addBookmark (string *book*, string *query*)

SearchEngine::addBookmark (p. 61) This adds bookmarks to the .txt file.

Parameters

<i>book</i>	The bookmark to be added.
<i>query</i>	The query the bookmark was from.

7.27.3.2 bool SearchEngine::addDocumentsToIndex (string *docpath*)

SearchEngine::addDocumentsToIndex (p. 61) Allows the user to choose between using a Hash table or AVLTree to create the inverted index file. This also initializes the **IndexHandler** (p. 48), **DocumentParser**, and **QueryEngine** (p. 54) objects to allow the Search Engine access to their functions. It also updates the tf/idf, inverted file index, page count, word count, and top fifty words.

Parameters

<i>type</i>	The user choice
-------------	-----------------

Returns

bool The flag for a successful text extract.

7.27.3.3 void SearchEngine::addToHistory (string *query*)

SearchEngine::addToHistory (p. 62) This adds the query to the search history.

Parameters

<i>query</i>	The query to be added.
--------------	------------------------

7.27.3.4 void SearchEngine::chooseStructure (char *type*)

SearchEngine::chooseStructure (p. 62) Allows the user to choose between using a Hash table or AVLTree to create the inverted index file. This also initializes the **IndexHandler** (p. 48), Documentparser, and **QueryEngine** (p. 54) objects to allow the Search Engine access to their functions. It also updates the tf/idf, inverted file index, page count, word count, and top fifty words.

Parameters

<i>type</i>	The user choice
-------------	-----------------

7.27.3.5 vector< document > SearchEngine::display_search_results (string *word*)

SearchEngine::display_search_results (p. 62) This searches the inverted file index for a query and displays the results.

Parameters

<i>word</i>	The query to search for.
-------------	--------------------------

Returns

vector<document> The vector containing the document objects from the query search.

7.27.3.6 vector< bookmark > SearchEngine::displayBookmarks ()

SearchEngine::displayBookmarks (p. 62) This displays the bookmarks.

Returns

vector<bookmarks> The vector of bookmark objects that were displayed.

7.27.3.7 bool SearchEngine::displayRawFile (string *filePath*)

SearchEngine::displayRawFile (p. 63) This displays the raw text of a requested PDF.

Parameters

<i>filePath</i>	The path to the requested PDF
-----------------	-------------------------------

Returns

bool The flag for a successful text extract.

7.27.3.8 int SearchEngine::getTotalPages ()

SearchEngine::getTotalPages (p. 63) Returns the total page count.

Returns

pages The total page count.

7.27.3.9 int SearchEngine::numWordsIndexed ()

SearchEngine::numWordsIndexed (p. 63) Returns the total indexed word count.

Returns

totalWordsIndexed

7.27.3.10 vector< string > SearchEngine::readFilePathFromTXTFile ()

SearchEngine::readFilePathFromTXTFile (p. 63) This reads the file paths from the .txt file of file paths.

Returns

vector<string> The vector of the file paths.

7.27.3.11 int SearchEngine::readTotalPages ()

SearchEngine::readTotalPages (p. 63) This outputs the total number of pages stored from the PDFs on the inverted index.

Returns

totalPages

7.27.3.12 `vector< document > SearchEngine::topfifty ()`

SearchEngine::topfifty (p. 64) Returns the top fifty most common words.

Returns

top50 The vector<document> containing the top fifty words and their count.

The documentation for this class was generated from the following files:

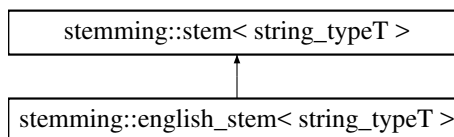
- searchengine.h
- searchengine.cpp

7.28 `stemming::stem< string_typeT >` Class Template Reference

The base class for language-specific stemmers.

```
#include <stemming.h>
```

Inheritance diagram for `stemming::stem< string_typeT >`:



Protected Member Functions

- void **find_r1** (const string_typeT &text, const wchar_t *vowel_list)
- void **find_r2** (const string_typeT &text, const wchar_t *vowel_list)
- void **find_spanish_rv** (const string_typeT &text, const wchar_t *vowel_list)
- void **find_french_rv** (const string_typeT &text, const wchar_t *vowel_list)
- void **find_russian_rv** (const string_typeT &text, const wchar_t *vowel_list)
- void **update_r_sections** (const string_typeT &text)
- bool **is_apostrophe** (const wchar_t &ch) const
- void **trim_western_punctuation** (string_typeT &text) const
- bool **is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U) const
is_suffix for one character
- bool **is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U) const
is_suffix for two characters
- bool **is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U) const
is_suffix for three characters
- bool **is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U) const
is_suffix for four characters

- **bool is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U) const
is_suffix for five characters
- **bool is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const wchar_t suffix6L, const wchar_t suffix6U) const
is_suffix for six characters
- **bool is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const wchar_t suffix6L, const wchar_t suffix6U, const wchar_t suffix7L, const wchar_t suffix7U) const
is_suffix for seven characters
- **bool is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const wchar_t suffix6L, const wchar_t suffix6U, const wchar_t suffix7L, const wchar_t suffix7U, const wchar_t suffix8L, const wchar_t suffix8U) const
is_suffix for eight characters
- **bool is_suffix** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const wchar_t suffix6L, const wchar_t suffix6U, const wchar_t suffix7L, const wchar_t suffix7U, const wchar_t suffix8L, const wchar_t suffix8U, const wchar_t suffix9L, const wchar_t suffix9U) const
is_suffix for nine characters
- **bool is_partial_suffix** (const string_typeT &text, const size_t start_index, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U)
comparison for two characters
- **bool is_partial_suffix** (const string_typeT &text, const size_t start_index, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U)
comparison for three characters
- **bool is_suffix_in_rv** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U)
RV suffix functions.
- **bool is_suffix_in_rv** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U)
RV suffix comparison for two characters.
- **bool is_suffix_in_rv** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U)
RV suffix comparison for three characters.
- **bool is_suffix_in_rv** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U)
RV suffix comparison for four characters.
- **bool is_suffix_in_rv** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U)
RV suffix comparison for five characters.
- **bool is_suffix_in_rv** (const string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const wchar_t suffix6L, const wchar_t suffix6U)
RV suffix comparison for six characters.

- R2 suffix comparison for seven characters.*

- R2 deletion for five character suffix.*

- R2 deletion for six character suffix.*

- R2 deletion for seven character suffix.*

- R2 deletion for eight character suffix.*

- bool **delete_if_is_in_rv** (string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const bool success_on_find=true)
- bool **delete_if_is_in_rv** (string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const bool success_on_find=true)
- bool **delete_if_is_in_rv** (string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const bool success_on_↵ find=true)
- bool **delete_if_is_in_rv** (string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const bool success_on_find=true)
- bool **delete_if_is_in_rv** (string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const bool success_on_find=true)
- bool **delete_if_is_in_rv** (string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_↵_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const wchar_t suffix6L, const wchar_↵_t suffix6U, const bool success_on_find=true)
- bool **delete_if_is_in_rv** (string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_↵_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const wchar_t suffix6L, const wchar_↵_t suffix6U, const wchar_t suffix7L, const wchar_t suffix7U, const bool success_on_find=true)
- bool **delete_if_is_in_rv** (string_typeT &text, const wchar_t suffix1L, const wchar_t suffix1U, const wchar_↵_t suffix2L, const wchar_t suffix2U, const wchar_t suffix3L, const wchar_t suffix3U, const wchar_t suffix4L, const wchar_t suffix4U, const wchar_t suffix5L, const wchar_t suffix5U, const wchar_t suffix6L, const wchar_↵_t suffix6U, const wchar_t suffix7L, const wchar_t suffix7U, const wchar_t suffix8L, const wchar_t suffix8U, const bool success_on_find=true)
- void **remove_german_umlauts** (string_typeT &text)
- void **italian_acutes_to_graves** (string_typeT &text)
- void **hash_dutch_yi** (string_typeT &text, const wchar_t *vowel_string)
Hash initial y, y after a vowel, and i between vowels into hashed character.
- void **unhash_dutch_yi** (string_typeT &text)
- void **hash_german_yu** (string_typeT &text, const wchar_t *vowel_string)
Hash 'u' and 'y' between vowels.
- void **unhash_german_yu** (string_typeT &text)
- void **hash_french_yui** (string_typeT &text, const wchar_t *vowel_string)
- void **unhash_french_yui** (string_typeT &text)
- void **hash_y** (string_typeT &text, const wchar_t *vowel_string)
- void **unhash_y** (string_typeT &text)
- void **hash_italian_ui** (string_typeT &text, const wchar_t *vowel_string)
Hash u after q, and u, i between vowels.
- void **unhash_italian_ui** (string_typeT &text)
- void **remove_dutch_umlauts** (string_typeT &text)
- void **remove_dutch_acutes** (string_typeT &text)
- void **remove_spanish_acutes** (string_typeT &text)
- size_t **get_r1** () const
- void **set_r1** (const size_t val)
- size_t **get_r2** () const
- void **set_r2** (const size_t val)
- size_t **get_rv** () const
- void **set_rv** (const size_t val)
- void **reset_r_values** ()

7.28.1 Detailed Description

```
template<typename string_typeT = std::wstring>
class stemming::stem< string_typeT >
```

The base class for language-specific stemmers.

The template argument for the stemmers are the type of std::basic_string that you are trying to stem, by default std::wstring (Unicode strings). As long as the char type of your basic_string is wchar_t, then you can use any type of basic_string. This is to say, if your basic_string has a custom char_traits or allocator, then just specify it in your template argument to the stemmer.

Example:

```
typedef std::basic_string<wchar_t, myTraits, myAllocator> myString;
myString word(L"documentation");
stemming::english_stem<myString> StemEnglish;
StemEnglish(word);
```

7.28.2 Member Function Documentation

7.28.2.1 `template<typename string_typeT = std::wstring> void stemming::stem< string_typeT >::hash_french_yui (string_typeT & text, const wchar_t * vowel_string) [inline], [protected]`

Hash u or i preceded and followed by a vowel, and y preceded or followed by a vowel. u after q is also hashed. For example, jouer -> joUer ennue -> ennule yeux -> Yeux quand -> qUand

7.28.2.2 `template<typename string_typeT = std::wstring> bool stemming::stem< string_typeT >::is_apostrophe (const wchar_t & ch) const [inline], [protected]`

Determines if a character is an apostrophe (includes straight single quotes).

Parameters

<i>ch</i>	The letter to be analyzed.
-----------	----------------------------

7.28.2.3 `template<typename string_typeT = std::wstring> bool stemming::stem< string_typeT >::is_suffix_in_r1 (const string_typeT & text, const wchar_t suffix1L, const wchar_t suffix1U) [inline], [protected]`

R1 suffix functions.

R1 suffix comparison for one character

7.28.2.4 `template<typename string_typeT = std::wstring> bool stemming::stem< string_typeT >::is_suffix_in_r2 (const string_typeT & text, const wchar_t suffix1L, const wchar_t suffix1U) [inline], [protected]`

R2 suffix functions.

R2 suffix comparison for one character

7.28.2.5 `template<typename string_typeT = std::wstring> bool stemming::stem< string_typeT >::is_suffix_in_rv (const string_typeT & text, const wchar_t suffix1L, const wchar_t suffix1U) [inline], [protected]`

RV suffix functions.

RV suffix comparison for one character

The documentation for this class was generated from the following file:

- stemming.h

7.29 string_util::string_tokenize< T > Class Template Reference

Tokenizes a string using a set of delimiters.

```
#include <string_util.h>
```

Public Member Functions

- **string_tokenize** (const T &val, const T &delim)
- bool **has_more_tokens** () const
- bool **has_more_delimiters** () const
- T **get_next_token** ()

7.29.1 Detailed Description

```
template<typename T>
class string_util::string_tokenize< T >
```

Tokenizes a string using a set of delimiters.

Date

2010

7.29.2 Constructor & Destructor Documentation

7.29.2.1 `template<typename T > string_util::string_tokenize< T >::string_tokenize (const T & val, const T & delim) [inline]`

Constructor which takes the string to parse and the delimiters to use.

Parameters

<i>val</i>	The string to parse.
<i>delim</i>	The set of delimiters to separate the string.

7.29.3 Member Function Documentation

7.29.3.1 `template<typename T> T string_util::string_tokenize<T>::get_next_token () [inline]`

Returns

The next token from the original string as a string object Note that empty tokens can be returned if there is proceeding or trailing delimiters in the string, or if there are repeated delimiters next to each other.

7.29.3.2 `template<typename T> bool string_util::string_tokenize<T>::has_more_delimiters () const [inline]`

Returns

Whether or not there are more delimiters in the string. This is useful for seeing if there are any delimiters at all when first loading the string.

7.29.3.3 `template<typename T> bool string_util::string_tokenize<T>::has_more_tokens () const [inline]`

Returns

Whether or not there are more tokens in the string.

The documentation for this class was generated from the following file:

- string_util.h

7.30 string_util::string_trim< char_typeT > Class Template Reference

trims whitespace from around a string

```
#include <string_util.h>
```

Public Member Functions

- `const char_typeT * operator() (const char_typeT *value, size_t length=std::basic_string< char_typeT >::npos)`
- `size_t get_trimmed_string_length () const`

7.30.1 Detailed Description

```
template<typename char_typeT>
class string_util::string_trim< char_typeT >
```

trims whitespace from around a string

The documentation for this class was generated from the following file:

- string_util.h

7.31 TextExtractor Class Reference

The **TextExtractor** (p. 72) class.

```
#include <textextractor.h>
```

Public Member Functions

- **TextExtractor (IndexHandler *, indexextractor *)**
TextExtractor::TextExtractor (p. 72) The **TextExtractor** (p. 72) constructor.
- virtual **~TextExtractor ()**
TextExtractor::~~TextExtractor (p. 72) The **TextExtractor** (p. 72) destructor.
- int **getTotalPages ()**
TextExtractor::getTotalPages (p. 73) Returns the total number of pages in the pdf.
- string **getDocName ()**
TextExtractor::getDocName (p. 73) Returns the PDF name.
- int **getWordCount ()**
TextExtractor::getWordCount (p. 73) Returns the word count.
- void **Init** (const char *pszInput, string docName)
TextExtractor::Init (p. 73) Uses Podofo to extract the text from the given PDF.

Public Attributes

- **IndexHandler * ih**
- **indexextractor * ie**
- string **mm**
- vector< string > **a**
- int **wordCount** =0
- string **docsName**

7.31.1 Detailed Description

The **TextExtractor** (p. 72) class.

7.31.2 Constructor & Destructor Documentation

7.31.2.1 TextExtractor::TextExtractor (IndexHandler * ih, indexextractor * ie)

TextExtractor::TextExtractor (p. 72) The **TextExtractor** (p. 72) constructor.

Parameters

<i>ih</i>	The IndexHandler (p. 48) object.
<i>ie</i>	The IndexInterface (p. 50) object.

7.31.3 Member Function Documentation

7.31.3.1 string TextExtractor::getDocName ()

TextExtractor::getDocName (p. 73) Returns the PDF name.

Returns

docsName The PDF name.

7.31.3.2 int TextExtractor::getTotalPages ()

TextExtractor::getTotalPages (p. 73) Returns the total number of pages in the pdf.

Returns

totalPages The total number of pages in the pdf.

7.31.3.3 int TextExtractor::getWordCount ()

TextExtractor::getWordCount (p. 73) Returns the word count.

Returns

wordCount The word count.

7.31.3.4 void TextExtractor::Init (const char * pszInput, string docName)

TextExtractor::Init (p. 73) Uses Podofo to extract the text from the given PDF.

Parameters

<i>pszInput</i>	The path to the corpus of PDFs
<i>docName</i>	The name of the pdf to extract text from.

The documentation for this class was generated from the following files:

- textextractor.h
- textextractor.cpp

7.32 userinterface Class Reference

The UserInterface is a command line menu driven class that makes use of the **SearchEngine** (p. 59). It allows the user to enter into two modes, Maintenance and Query, allows the user many options including adding a new pdf

to the inverted index, clearing the index, searching the PDF, outputting total pages, outputting total words indexed, outputting the top fifty words, outputting the corpus paths, storing and clearing the search history, storing and clearing bookmarks and outputting the raw text from a selected pdf.

```
#include <userinterface.h>
```

Public Member Functions

- **userinterface** ()

*The UserInterface is a command line menu driven class that makes use of the **SearchEngine** (p. 59). It allows the user to enter into two modes, Maintenance and Query, allows the user many options including adding a new pdf to the inverted index, clearing the index, searching the PDF, outputting total pages, outputting total words indexed, outputting the top fifty words, outputting the corpus paths, storing and clearing the search history, storing and clearing bookmarks and outputting the raw text from a selected pdf.*

- void **use** ()

use method runs the user interface. It

7.32.1 Detailed Description

The UserInterface is a command line menu driven class that makes use of the **SearchEngine** (p. 59). It allows the user to enter into two modes, Maintenance and Query, allows the user many options including adding a new pdf to the inverted index, clearing the index, searching the PDF, outputting total pages, outputting total words indexed, outputting the top fifty words, outputting the corpus paths, storing and clearing the search history, storing and clearing bookmarks and outputting the raw text from a selected pdf.

CSE 2341 UserInterface.h

Author

Patrick Yienger (owner)

Version

1.0 05/07/17 The userinterface class

7.32.2 Constructor & Destructor Documentation

7.32.2.1 **userinterface::userinterface** ()

The UserInterface is a command line menu driven class that makes use of the **SearchEngine** (p. 59). It allows the user to enter into two modes, Maintenance and Query, allows the user many options including adding a new pdf to the inverted index, clearing the index, searching the PDF, outputting total pages, outputting total words indexed, outputting the top fifty words, outputting the corpus paths, storing and clearing the search history, storing and clearing bookmarks and outputting the raw text from a selected pdf.

CSE 2341 UserInterface.cpp

Author

Patrick Yienger (owner)

Version

1.0 05/07/17 **userinterface::userinterface** (p. 74) The constructor for the user interface.

7.32.3 Member Function Documentation

7.32.3.1 void userinterface::use ()

use method runs the user interface. It

userinterface::use (p. 75) The method that runs the user interface menu. It allows the user to enter into two modes, Maintenance and Query, allows the user many options including adding a new pdf to the inverted index, clearing the index, searching the PDF, outputting total pages, outputting total words indexed, outputting the top fifty words, outputting the corpus paths, storing and clearing the search history, storing and clearing bookmarks and outputting the raw text from a selected pdf.

The documentation for this class was generated from the following files:

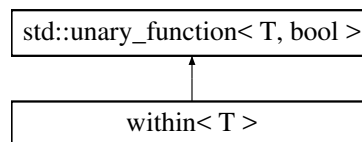
- userinterface.h
- userinterface.cpp

7.33 within< T > Class Template Reference

Determines if a value is within a given range.

```
#include <utilities.h>
```

Inheritance diagram for within< T >:



Public Member Functions

- **within** (T range_begin, T range_end)
- bool **operator()** (T value) const

7.33.1 Detailed Description

```
template<typename T>
class within< T >
```

Determines if a value is within a given range.

The documentation for this class was generated from the following file:

- utilities.h

Index

- AVLTreeIndex, 29
 - addIndex, 30
 - findIndex, 30
 - insert, 30, 31
 - split, 31
 - topwords, 31
 - totalWordsIndexed, 31
- addBookmark
 - SearchEngine, 61
- addDocumentsToIndex
 - SearchEngine, 61
- addIndex
 - AVLTreeIndex, 30
 - hashy, 44
 - IndexHandler, 49
- addToHistory
 - SearchEngine, 62
- axtoi
 - string_util, 25
- backup_variable< T >, 32
- bookmark, 32
- CASSERT
 - Debugging, 10
- chooseStructure
 - SearchEngine, 62
- comparable_first_pair< T1, T2 >, 33
- compare_doubles
 - Mathematics, 12
- compare_doubles_greater
 - Mathematics, 13
- compare_doubles_less
 - Mathematics, 13
- compare_doubles_less_or_equal
 - Mathematics, 13
- copy_member
 - Utilities, 18
- DUMP_TO_FILE
 - Debugging, 10
- Debugging, 10
 - CASSERT, 10
 - DUMP_TO_FILE, 10
 - NON_UNIT_TEST_ASSERT, 10
- display_search_results
 - SearchEngine, 62
- displayBookmarks
 - SearchEngine, 62
- displayIndex
 - hashy, 44
- displayRawFile
 - SearchEngine, 63
- document, 34
- DocumentParser, 34
 - DocumentParser, 36
 - extract, 36
 - getStemmed, 36
 - getTotalPages, 37
 - isStopWord, 37
 - numOfDocs, 37
 - rawTextExtract, 37
- double_less, 38
- extract
 - DocumentParser, 36
- find_last_not_of
 - string_util, 25
- find_last_of
 - string_util, 25
- find_matching_close_tag
 - string_util, 25
- findIndex
 - AVLTreeIndex, 30
 - hashy, 44
- findTDIFs
 - QueryEngine, 55
- get_next_token
 - string_util::string_tokenize, 71
- getDifference
 - QueryEngine, 55
- getDocName
 - rawOutputExtractor, 58
 - TextExtractor, 73
- getDocs
 - IndexHandler, 49
- getIntersection
 - QueryEngine, 55
- getStemmed
 - DocumentParser, 36
- getTotalPages
 - DocumentParser, 37
 - rawOutputExtractor, 58
 - SearchEngine, 63
 - TextExtractor, 73
- getUnion
 - QueryEngine, 56
- getWordCount

- rawOutputExtractor, 58
- TextExtractor, 73
- has_more_delimiters
 - string_util::string_tokenize, 71
- has_more_tokens
 - string_util::string_tokenize, 71
- has_suffix
 - string_util, 26
- hash_french_yui
 - stemming::stem, 69
- Hashy
 - hashy, 45
- hashy, 43
 - addIndex, 44
 - displayIndex, 44
 - findIndex, 44
 - Hashy, 45
 - hashy, 44
 - insertFromFile, 45
 - NumItemsInIndex, 45
 - split, 45
 - topwords, 46
 - totalWordsIndexed, 46
- index_node, 46
- IndexHandler, 48
 - addIndex, 49
 - getDocs, 49
 - IndexHandler, 48
 - topFifty, 49
 - totalWordsIndexed, 49
- IndexInterface, 50
- indexextractor, 47
 - indexextractor, 47
- Indexing, 9
- Init
 - rawOutputExtractor, 59
 - TextExtractor, 73
- insert
 - AVLTreeIndex, 30, 31
- insertFromFile
 - hashy, 45
- int_to_bool
 - Mathematics, 14
- is_apostrophe
 - stemming::stem, 69
- is_hex_digit
 - string_util, 26
- is_space
 - string_util, 26
- is_suffix_in_r1
 - stemming::stem, 69
- is_suffix_in_r2
 - stemming::stem, 69
- is_suffix_in_rv
 - stemming::stem, 69
- is_within
 - Utilities, 18
- isStopWord
 - DocumentParser, 37
- item, 50
- Mathematics, 12
 - compare_doubles, 12
 - compare_doubles_greater, 13
 - compare_doubles_less, 13
 - compare_doubles_less_or_equal, 13
 - int_to_bool, 14
 - safe_divide, 14
 - safe_modulus, 14
- NON_UNIT_TEST_ASSERT
 - Debugging, 10
- NumItemsInIndex
 - hashy, 45
- numOfDocs
 - DocumentParser, 37
- numWordsIndexed
 - SearchEngine, 63
- operator()
 - stemming::english_stem, 41
 - Utilities, 18
- QueryEngine, 54
 - findTDIFs, 55
 - getDifference, 55
 - getIntersection, 55
 - getUnion, 56
 - QueryEngine, 55
 - querySearch, 56
 - relevancySort, 56
 - search, 57
- querySearch
 - QueryEngine, 56
- rawOutputExtractor, 57
 - getDocName, 58
 - getTotalPages, 58
 - getWordCount, 58
 - Init, 59
 - rawOutputExtractor, 58
- rawTextExtract
 - DocumentParser, 37
- readFilePathFromTXTFile
 - SearchEngine, 63
- readTotalPages
 - SearchEngine, 63
- relevancySort
 - QueryEngine, 56
- remove_blank_lines
 - string_util, 26
- remove_extra_spaces
 - string_util, 27
- safe_divide
 - Mathematics, 14
- safe_modulus

- Mathematics, 14
- search
 - QueryEngine, 57
- SearchEngine, 59
 - addBookmark, 61
 - addDocumentsToIndex, 61
 - addToHistory, 62
 - chooseStructure, 62
 - display_search_results, 62
 - displayBookmarks, 62
 - displayRawFile, 63
 - getTotalPages, 63
 - numWordsIndexed, 63
 - readFilePathFromTXTFile, 63
 - readTotalPages, 63
 - SearchEngine, 61
 - topfifty, 63
- size_of_array
 - Utilities, 18
- split
 - AVLTreeIndex, 31
 - hashy, 45
- Stemming, 11
- stemming, 21
- stemming::english_stem
 - operator(), 41
- stemming::english_stem< string_typeT >, 38
- stemming::no_op_stem< string_typeT >, 53
- stemming::stem
 - hash_french_yui, 69
 - is_apostrophe, 69
 - is_suffix_in_r1, 69
 - is_suffix_in_r2, 69
 - is_suffix_in_rv, 69
- stemming::stem< string_typeT >, 64
- strcspn_pointer
 - string_util, 27
- string_tokenize
 - string_util::string_tokenize, 70
- string_util, 22
 - axtoi, 25
 - find_last_not_of, 25
 - find_last_of, 25
 - find_matching_close_tag, 25
 - has_suffix, 26
 - is_hex_digit, 26
 - is_space, 26
 - remove_blank_lines, 26
 - remove_extra_spaces, 27
 - strcspn_pointer, 27
 - strnatordcmp, 27
 - strnchr, 27
 - strncspn, 27
 - strnistr, 28
 - strnlen, 28
 - strrstr, 28
 - strtod_ex, 28
 - strtol, 28
- string_util::equal_basic_string_i_compare< T >, 41
- string_util::equal_string_compare< T >, 42
- string_util::equal_string_i_compare< T >, 42
- string_util::less_basic_string_compare< T >, 51
- string_util::less_string_compare< T >, 51
- string_util::less_string_i_compare< T >, 52
- string_util::less_string_n_compare< T >, 52
- string_util::less_string_natural_order_i_compare< T >, 53
- string_util::less_string_ni_compare< T >, 53
- string_util::string_tokenize
 - get_next_token, 71
 - has_more_delimiters, 71
 - has_more_tokens, 71
 - string_tokenize, 70
- string_util::string_tokenize< T >, 70
- string_util::string_trim< char_typeT >, 71
- StringOperations, 16
- strnatordcmp
 - string_util, 27
- strnchr
 - string_util, 27
- strncspn
 - string_util, 27
- strnistr
 - string_util, 28
- strnlen
 - string_util, 28
- strrstr
 - string_util, 28
- strtod_ex
 - string_util, 28
- strtol
 - string_util, 28
- TextExtractor, 72
 - getDocName, 73
 - getTotalPages, 73
 - getWordCount, 73
 - Init, 73
 - TextExtractor, 72
- topFifty
 - IndexHandler, 49
- topfifty
 - SearchEngine, 63
- topwords
 - AVLTreeIndex, 31
 - hashy, 46
- totalWordsIndexed
 - AVLTreeIndex, 31
 - hashy, 46
 - IndexHandler, 49
- use
 - userinterface, 75
- userinterface, 73
 - use, 75
 - userinterface, 74
- Utilities, 17

- copy_member, 18
- is_within, 18
- operator(), 18
- size_of_array, 18
- within, 19
- within_range, 19

within

- Utilities, 19

within< T >, 75

within_range

- Utilities, 19