# AutoComplete*

Vishak Srikanth
Final Project for OCS15: AP Computer Science A
Stanford Online High School
May 2021
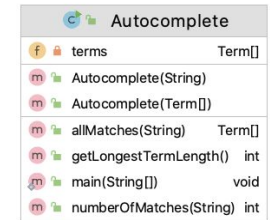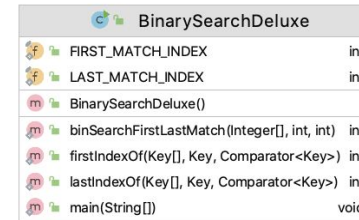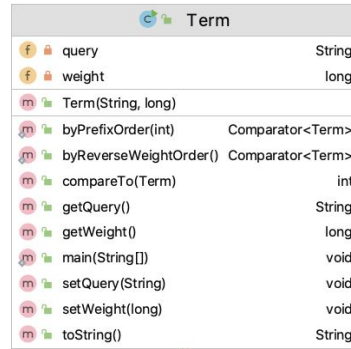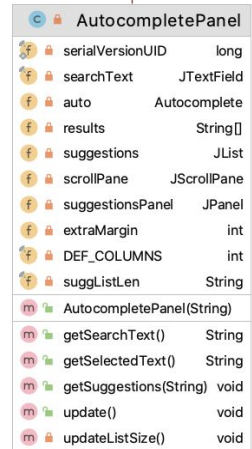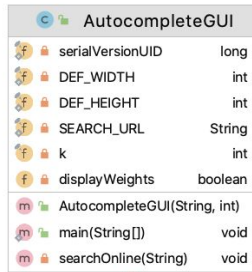
# AutoComplete: Project Background

Goal

- Given a **prefix**, find all queries that **<u>start</u>** with the given prefix string, in descending order of weight (i.e priority of results)

    => Implement *Autocomplete* for a given set of *N terms*, where a term is a query string with associated non-negative weight

- Autocomplete is ubiquitous!
    - Google: Search search suggestions as user types search terms
    - Facebook: Find suggestions for Friends as user types
    - LinkedIn: People Search filters people as you type names
    - IMDB, Netflix: Show movies that match as the user types
    - Mobile phones: Intelligent text completions as user types

# UML Diagram of Project Structure



**AutocompleteGUI**
- serialVersionUID : long
- DEF_WIDTH : int
- DEF_HEIGHT : int
- SEARCH_URL : String
- k : int
- displayWeights : boolean
- AutocompleteGUI(String, int)
- main(String[]) : void
- searchOnline(String) : void

**Term**
- query : String
- weight : long
- Term(String, long)
- byPrefixOrder(int) : Comparator<Term>
- byReverseWeightOrder() : Comparator<Term>
- compareTo(Term) : int
- getQuery() : String
- getWeight() : long
- main(String[]) : void
- setQuery(String) : void
- setWeight(long) : void
- toString() : String

**BinarySearchDeluxe**
- FIRST_MATCH_INDEX : int
- LAST_MATCH_INDEX : int
- BinarySearchDeluxe()
- binSearchFirstLastMatch(Integer[], int, int) : int
- firstIndexOf(Key[], Key, Comparator<Key>) : int
- lastIndexOf(Key[], Key, Comparator<Key>) : int
- main(String[]) : void

**Autocomplete**
- terms : Term[]
- Autocomplete(String)
- Autocomplete(Term[])
- allMatches(String) : Term[]
- getLongestTermLength() : int
- main(String[]) : void
- numberOfMatches(String) : int

**AutocompletePanel**
- serialVersionUID : long
- searchText : JTextField
- auto : Autocomplete
- results : String[]
- suggestions : JList
- scrollPane : JScrollPane
- suggestionsPanel : JPanel
- extraMargin : int
- DEF_COLUMNS : int
- suggListLen : String
- AutocompletePanel(String)
- getSearchText() : String
- getSelectedText() : String
- getSuggestions(String) : void
- update() : void
- updateListSize() : void

**ByReverseWightOrder**
- ByReverseWightOrder()
- compare(Term, Term) : int

**ByPrefixOrder**
- rCharsforComparison : int
- ByPrefixOrder(int)
- compare(Term, Term) : int

# Key Classes & Their Use

- Term: Data type that has a query string and an associated integer weight
  - Supports comparison of terms by three different orders:
    - lexicographic order by query string (natural order: alphabetic with ties (identical strings)ordered in their ordinal position in original array)
    - descending order by weight (i.e. priority of search results such as highest to lowest votes or box office collections for movies, or highest to lowest city population, highest to lowest word frequency in a corpus such as wikipedia)
    - lexicographic order by query string but using only the first r characters (used to query prefix string from a database)
- BinarySearchDeluxe: Enhanced Binary Search algorithm that has methods to return position of **first** and **last match** of a key in an array of sorted objects to be searched
  - Makes $1 + \lceil \log_2 N \rceil$ comparisons (worst case) where N is the length of the array. [ compare is a single call to Comparator<Key>.compare() ]

# Key Classes & Their Use

- Autocomplete: Main data type that provides autocomplete functionality
  - Has a field that holds an array of Term objects
  - Uses static methods from BinarySearchDeluxe class
  - Flow:
    - Sort terms in lexicographic order
    - Use binary search to find the all query strings that start with a given prefix (for e.g as user types each character the prefix to search will change)
    - Sort the matching terms in descending order by weight
  - Methods:
    - Returns all terms that start with the given prefix, in descending order of weight [uses log N + M log M comparisons (worst case) if N is the number of terms and M is the number of matching terms].
    - Returns the number of terms that start with the given prefix [makes log N comparisons (worst case)
- AutocompleteGUI: Swing GUI reference implementation provided by Princeton CS department as part of CS226 course
  - Used as is except for one small change to remove the dependency on a custom input file reader jar and replaced with other generic classes and methods that are part of any standard Java 8 install

# Testing

- All databases used are in the subfolder called 'data' in the working directory of IntelliJ project
- Variety of database files used for testing: movies (250K entries), cities (93K entries), google n-gram words (1m+ for 5-grams for example) etc.
- Each Database text file has following structure:
  - Header line with # of entries in DB
  - Each line has a weight, a tab, and then the query term for the string as shown
- A term Array is created by parsing the input file using standard Java buffered reader and writer interfaces
- The Autocomplete data structure can be created once the terms array is created
- The GUI app is run against a specific database with a max hit count (i.e. limits matching results to a number specified)



```
229447
 760507625  Avatar (2009)
 658672302  Titanic (1997)
 623357910  The Avengers (2012)
 534858444  The Dark Knight (2008)
 460935665  Star Wars (1977)
 448139099  The Dark Knight Rises (2012)
 436471036  Shrek 2 (2004)
 423315812  Pirates of the Caribbean: Dead Man's Chest (2006)
 422783777  The Lion King (1994)
 415004880  Toy Story 3 (2010)
 408992272  Iron Man Three (2013)
 408010692  The Hunger Games (2012)
 403706375  Spider-Man (2002)
 402348347  Jurassic Park (1993)
 402111870  Transformers: Revenge of the Fallen (2009)
 381011219  Harry Potter and the Deathly Hallows: Part 2 (2011)
 380262555  Star Wars: Episode III — Revenge of the Sith (2005)
 377845905  The Lord of the Rings: The Return of the King (2003)
```



Autocomplete Me

Search query: San                                    Search Google

☑ Show weights
| Santiago, Chile | 4837295 |
| Santo Domingo, Dominican Republic | 2201941 |
| Sanaa, Yemen | 1937451 |
| Santa Cruz de la Sierra, Bolivia | 1364389 |
| San Antonio, Texas, United States | 1327407 |

# Unit Testing Term Class

- A term Array is created by parsing cities-medium.txt
- The terms are sorted based on a prefix sort of 1st 3 characters
- The main method requires 3 command line arguments:
  - Database text input file (cities-medium.txt in demo)
  - Prefix length to sort terms entries by (3 in demo)
  - Output file (term-unittest-prefix-sorted-cities-medium.out in demo)
- Output file shows cities in sorted order based on 1st 3 characters

# Unit Testing BinarySearchDeluxe Class

- Main method creates an array of Integers having duplicates and an array of Terms having duplicate strings
- Boxed Integer array is sorted in natural order (ascending int value)
- Terms are sorted based on default comparator which sorts terms in lexicographic order (alphabetic)
- All output is printed locally to show the sorting and the binary search results to find the 1st and last match of the search key in both cases

# Unit Testing AutoComplete Class

Main method takes 3 command line parameters:

- Reads a database input file (imdb-votes.txt) and creates an Autocomplete object loading the content into an array of Term objects
- Reads a query file (query.in) that contains # of results to be returned followed by each query to be against the loaded database returning the specified number of results
- Writes output to the specified file autocomplete-unit-test-results.out

# Testing AutoCompleteGUI Class

Demo shows 2 different scenario:

1. Autocomplete that returns top 5 matches (search results are weighted in descending order of population) as user searches the cities database. When user double clicks on their selection it puts them on a Google search results page with selected term with search already performed.

2. Autocomplete that returns upto 10 top matches (search results are weighted in descending order of IMDB votes) as user searches the IMDB ratings database and then puts user on a Google search results page with selected term once the user selects a movie and presses Search on Google button in UI.