

# Design doc.

Alexader Yermakov [OUI 20407 DSA Course]{26/06/2019}

## MajellaMap

- Open addressing HashMap
- Can use single values as input to behave simply like HashSet *(i.e. hash table, as was requested)*
- Uses MurMur3 hashing algorithm
- Uses Robin Hood hashing *(trying to fit items as close as possible to their best position)*
- Uses Fibbonachi mapping of hash values *(to get more even distribution and reduce preliminary expansions)*
- Limits the amount of maximal scans to  $\log_2 n$  of current size, if reached expands
- *Currently cannot shrink if already expanded*

## This project can...

- Use text and dictionary as input *(via command line args)*
- Load text to internal RedBlack Tree data structure
- Load dictionary to HashTable styled HashMap via MajellaMap
- Run though Mapped dictionary and remove all the occurences of its word in the Tree
- Print the result in readable format *(examples in sample\_output.pdf)*

# Known limitations

- Although has a default dictionary and demo file structure, I couldn't make it run when packaged to Jar. :(
- That means it HAS to get two arguments as input [text\_file] and [dictionary\_file]
- Uses regex to remove punctuation from splitted singular words that mean something like "*(apple*" will get transformed to *apple* and pass the check, technically in the scope of the assignment it doesn't matter. *(it was an attempt to generify the pattern matching, so I can split it to multiple stages in the future).*

## INSERT/DELETE Complexities

On all **put()** operations, we get  **$\log_2 n$**  worst case scenario.

So we get  **$\log_2 n$**  amortized complexity, if we don't count the map reallocation when it expands.

1. By using Knuth's multiplicative method (*Fibonacci mapping*) we get more or less even distribution of hash values, when they get reduced to lower range (*current map size*).
2. On each addition we calculate item position and see if someone already occupies this place.
3. If yes, we compare its *povertyLevel* to our, that is basically a characteristic that represents how far away from the best position this item is currently sitting.
4. If item is wealthier than the one we trying to insert, we steal its

position and continue inserting the item we stole the place from.

5. At most we will do it  $\log_2 n$  cycles of this loop, at which point we will force expansion and rehashing.

On all **remove()** operations, we get  $\log_2 n$  average case scenario.

1. The to actually find the item we want to remove is still  $\log_2 n$ .
2. However after finding we will proceed to *fix* the Map, from the point we removed the item, to its end, up until we hit some item that is already slotted in its best place, or an empty cell.
3. Current MAX\_LOAD\_FACTOR is 0.53, so it is practically impossible for us to get in to situation, when we need to readjust  $n$  items from the first cell to the last, for that we would have to have MAX\_LOAD\_FACTOR of 1.
4. Since we readjust the items and try to distribute them evenly even with bad hashing function we will expand before any such scenario can happen.

## Implementation details

- Our MajellaMap size is growing in powers of two +  $\log_2$  of its size, however from mapping function point of view its size is just two in some power, without adding its power to final result.
- For that reason it will never spill out of bound.
- We store the temporarily stolen items if last slot of the Map, it can never be mapped, for reasons above.
- Checking if this slot is empty or not by the end of insertion cycle, tells us if we need to expand or not, however LOAD\_FACTOR can

affect this decision too.

- I used [12Dict](#) collection of dictionaries to test the spell checking, they are included with rest of the files.
- As a secondary text for tests, I've used first three chapters of [The History of the Peloponnesian War](#), its also included.

## Used tech.

- IntelliJ Community ed.
- Maven build tool
- Git