Lecture 31 - Consider the characters / Trie-ing to use data structures

Q: Can you jump to a specific line in a file? — How do you delete a certain line in a file?

```
# Example from HW5
DISK_PATH = pathlib.Path("disk")
                                                           \equiv test.txt
                                                                       def write_file():
                                                           disk \geq \equiv test.txt
   with open(DISK_PATH/"test.txt", "w") as f:
                                                                  hello world 🚺
                                                             1
       f.write("hello world\n")
                                                                  this is some text
                                                             2
       f.write("this is some text\n")
                                                                  on a line∖N
                                                             3
       f.write("on a line\n")
       f.write("four\n")
                                                                  four (N
                                                              4
       f.write("five\n")
                                                             5
                                                                  five
       f.write("six\n")
                                                              6
                                                                  six 🚺
                                                              7
```

A mental model for files in this class:

=> Can think about files as a list of lines, where each line is a string => Each line ends with a newline character (write as "\n"), which ends the line

This is a generally good model for working with text files, but it turns out that things are actually a bit more complicated under the hood, which complicates things like jumping to a specific line, or deleting a specific line.

For more discussion on files (and answers to "can you jump to a specific line?" and "how to delete a line of a file?") see the end of these notes.

Design problems: Try and trie again

Design problem 1: A large-scale internet service assigns an internal ID number to each username. Here are their key criteria:

- These IDs are used for various tasks (creating URLs, storing relationships between users, etc), so lookup needs to be fast.
- The site also adds new users frequently, so assigning new IDs to new users has to be fast
- No two users can have the same ID
- Because the site has so many users, they are concerned about the space being used to store this information

What data structure do you choose? How many many many solution require, in terms OBJECTS of the number U of users?

HASHMAP < 1D, USERNAME >



Classic hashmap structure: each array slot contains a linked list of KVPairs: => Number of objects: 1 HashMap object + underlying array + (1 Node + 1 KVPair)*(U users)

(Also possible to be a bit more compact and merge the Node and KVPair into one object, but structure remains the same.)

Design problem 2: Swype keyboards let people type by passing their fingers over the keys rather than hitting individual keys separately. The software for such a keyboard makes use of a dictionary of available words to help figure out which keys might be relevant, as well as which keys might come "next" based on what a user has typed so far.

What data structure do you use to store this dictionary? \implies HASHSET

How will you use it to predict which keys might come next?



https://www.flickr.com/photos/hahatango/3173033612

For example: How to keep all "qui" words together? => How to group words with a common prefix?

Could use a hashmap: prefix => (words with that prefix) However, lots of potential prefixes: q, qu, qui => hard to scale, because many duplicates

Can we do better???

Tries: Storing prefixes with a tree!

An example dictionary: pie, piece, piano, pet, pail

Idea: Share prefixes of words along branches of the tree

=> some nodes are "endpoints" that store a result (ie, a word in the dictionary, a decision, etc.)

=> This is called a trie (reTRIEval)



How would a trie work for design problem #2?

Can follow the trie to determine if the letters a user types are in the dictionary:

- For each letter the user enters, follow one node on the trie
- Can tell when the path ends at a valid word

- (From there, could possibly do some more computation to figure out the next "closest" word, to help with prediction)

Q: Could you embed a trie in an array?

=> Would only work if we had a defined structure (one slot for each character per branch, so that every node has the same number of child nodes)

Q: Does a trie actually save space?

=> Could become very large due to the number of possible branches. However, there are some tricks to make the trie use fewer nodes (eg. reducing the long chains of nodes with no branching--not important in this course)

Back to Design Problem #1

Some usernames:

What if we need to store a LOT of usernames??

chen chad cami carl cary



Could represent usernames using a trie: time to look up a username would involve searching the tree

=> Even with a lot of users, the search time is based on the <u>number of characters in the</u> <u>username</u> (eg. a small, constant value, like 8 characters)

=> In some applications, this might be faster than using a hashmap, due to the time required to hash the individual characters of each username. (This tradeoff isn't something we'd ask you to make a decision about, but we want you to see the idea and understand why tries are cool.)

What's in a string? Basic version: ASCII Codes

Every character has a corresponding numeric value ("ascii code")



CHARACTE

Every character has some numeric representation

=> A file is just a big, ordered sequence of characters

This is something you'll spend more time learning if you take a systems course (eg. CS 300 or CS 330). For CS 200, it's just helpful to know that files and strings are composed of characters, which are really just numbers.

What does this have to do with jumping to a specific line in a file?

The newline character (\n) is just another character (ie, a number). The only difference is that it has a special meaning: instead of "printing" as a letter or number, it tells the program or terminal reading the file to go to the next line.

```
with open(DISK_PATH/"test.txt", "r") as f:
    while True:
        line = f.readline()
        if line == "":
                                         => readline() just scans the list of characters
            break
                                         until it finds the newline
        print(line.strip())
                                         This is actually just a linear search of all the
                                         characters in the file!
 (View from a <u>hex editor</u>):
  68 65 6C 6C 6F 20 77 6F 72 6C 64 0A 74 68 69 73 hello world this
  20 69 73 20 73 6F 6D 65 20 74 65 78 74 0A 6F 6E is some text.on
  20 61 20 6C 69 6E 65 0A 66 6F 75 72 0A 66 69 76 a line four fiv
  65 0A 73 69 78 <mark>0A</mark> +
                                                 e.six.
```

So: can you jump to a specific line in a file?

=> Not with the tools that we have: don't know how far apart the newline characters are--need to read the file line by line until you get to the line you want

=> To learn how to do better, check out CS 300!

Similarly, if you wanted to delete a line, would need to read the whole file in, and then write again (without the line you want to delete)

=> Essentially, build the big array of characters all over again

Looking closer at the file....



Don't know how far apart the newline characters are => need to read the file line by line until you get to the line you want