

Lexicon case study

- **Special case of Set/Map**
 - ◆ Keys always strings, in fact, English words
 - ◆ No associated value
- **Core operations**
 - ◆ add
 - ◆ containsWord
 - ◆ containsPrefix

Lexicon as sorted vector

```
private:  
    Vector<string> words; abacus abate ... zygote
```

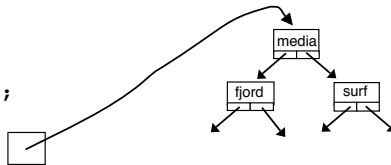
- containsWord?**
 - ◆ Binary search using ==, $O(\lg N)$ time
- containsPrefix?**
 - ◆ Binary search using substr ==, $O(\lg N)$
- add?**
 - ◆ Binary search to find position, shuffle to insert, $O(N)$
- space usage?**
 - ◆ Per-entry = sizeof(string) \sim length* sizeof(char)
 - ◆ \sim 8 bytes per word (assuming words average len 8 chars)

1

2

Lexicon as binary search tree

```
private:  
    struct node {  
        string word;  
        node *left, *right;  
    };  
    node *root;
```

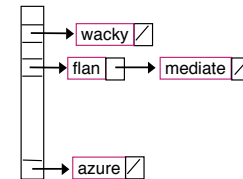


- containsWord?**
 - ◆ Tree search using ==, $O(\lg N)$ time
- containsPrefix?**
 - ◆ Tree search using substr ==, $O(\lg N)$
- add?**
 - ◆ Tree search to find position, insert new leaf, $O(\lg N)$
- space usage?**
 - ◆ Per-entry = string + left/right ptrs + balance factor
 - ◆ \sim 18 bytes per word

3

Lexicon as hash table

```
private:  
    struct cell {  
        string word;  
        cell *next;  
    };  
    cell *buckets[NBuckets];
```



- containsWord?**
 - ◆ Hash to bucket, search bucket, $O(N/B)$
- containsPrefix?**
 - ◆ Search all buckets, $O(N)$
- add?**
 - ◆ Hash to bucket, search/add to bucket, $O(N/B)$
- space usage?**
 - ◆ Per-entry = string + 4-byte pointer + 4-byte bucket ptr
 - ◆ \sim 16 bytes per word

4

Summary so far

	Sorted vector	BST	Hashtable
containsWord	$O(\lg N)$	$O(\lg N)$	$O(1)$
containsPrefix	$O(\lg N)$	$O(\lg N)$	$O(N)$
add	$O(N)$	$O(\lg N)$	$O(1)$
bytes per word	8	18	16

- **OSPD2.txt file has 125,000 words**
 - ♦ Average length 8 characters, file is 1.1 MB
- **At 8-18 bytes per word, total is 1MB to 2.3 MB**
- **Boggle lexicon actually only takes 350K (.3 MB)....**
 - ♦ How???

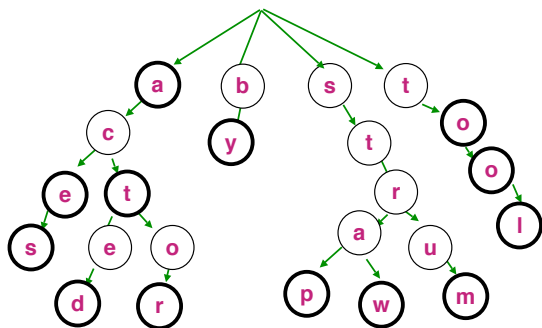
5

Notice any patterns?

straddle	straightener	straitjacket	strangled	strategically
straddled	straighteners	straitjackets	stranglehold	strategies
straddler	straightening	straitlaced	strangleholds	strategist
straddlers	straightens	straitly	strangler	strategists
straddles	straighter	straitness	stranglers	strategy
straddling	straightest	straits	strangles	stratification
strafe	straightforward	strake	strangling	stratifications
strafed	straightforwardly	straked	strangulate	stratified
straffer	straightforwardness	strakes	strangulated	stratifies
strafers	straighting	stramash	strangulates	stratify
strafes	straightjacket	stramashes	strangulating	stratifying
strafing	straightjackets	stramonies	strangulation	stratosphere
straggle	straightlaced	stramony	strap	stratospheres
straggled	straightly	strand	strapless	stratospheric
straggler	straights	stranded	strapped	stratum
straggles	straightway	strander	strapper	stratums
straggles	strain	stranders	strappers	stratus
stragglier	strained	stranding	strapping	stravage
straggliest	strainer	strands	strappings	stravaged
stragglings	strainers	strang	straps	stravages
straggly	straining	strange	strass	stravaging
straight	strains	strangely	strasses	stravaig
straightaway	strait	strangeness	strata	stravaiged
straightaways	straiten	stranger	stratagem	stravaiging
straighted	straitened	strangered	stratagems	stravaigs
straightedge	straitening	strangering	stratal	straw
straightedges	straitens	strangers	stratas	strawberries
straighten	straiter	strangest	strategic	strawberry
straightened	straitest	strangle	strategical	

6

Letter trie

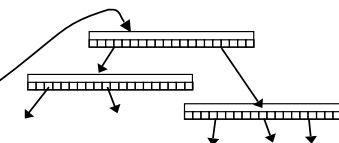


- Trie from retrieval (but pronounced "try")
- 26-way branching tree
- Paths are prefixes
- Path that ends at thick circle is word

7

Lexicon as trie

```
private:
struct node {
    char letter;
    bool isWord;
    node *children[26];
};
node *root;
```



containsWord?

- ♦ Trace path, check isWord, $O(\text{length})$

containsPrefix?

- ♦ Trace path $O(\text{length})$

add?

- ♦ Trace path, add new nodes as needed $O(\text{length})$

space usage?

- ♦ Per-node = 106 bytes
- ♦ Trie for OSPD2 has ~250,000 nodes = 26.5 MB!!!

8

Dynamic array of children

```
private:
struct node {
    char letter;
    bool isWord;
    node **children;
    int nChildren;
};
node *root;
```

Replace static array of children with dynamically-sized array

- ◆ Most of 26 entries are NULL anyway

Space usage?

- ◆ Per-node = 10 bytes + NumChildren*4 bytes
- ◆ Each node averages 6 children, so total 34 bytes
- ◆ Trie for OSPD2 has ~250,000 nodes = **8.5 MB...**

Flatten tree into array

```
private:
struct node {
    char letter;
    bool isWord;
    int childIndex;
    int nChildren;
};
node *root;
```

'F'	'A'	'B'	'C'	'D'	'S'
1	27	39	45	55	T-1
26	12	6	20	8	0

Flatten tree into array

- ◆ Like we did for heap, but not fixed location
- ◆ Saves all space used for pointers
- ◆ Makes data structure much less flexible (insert/delete words)

Space usage?

- ◆ Per-node = 10 bytes
- ◆ Trie for OSPD2 has ~250,000 nodes = **2.5 MB**
- ◆ (Now same ballpark as bst/hash)

Where to go next?

adapted	basted	depicted	impacted	interrupted
adapter	baster	depicter	impacter	interrupter
adapters	basters	depicters	impacters	interrupters
adapting	basting	depicting	impacting	interrupting
adaption	bastion	depiction	impaction	interruption
adapptions	bastions	depictions	impactions	interruptions
adapts	basts	depicts	impacts	interrupts
adopted	billed	deserted	indented	invented
adopter	billier	deserter	indenter	inventer
adopters	billiers	deserters	indenters	inventers
adopting	billing	deserting	indenting	inventing
adoption	billion	desertion	indention	invention
adoptions	billions	desertions	indentions	inventions
adopts	bills	deserts	indents	invents
affected	camped	detected	inserted	perfected
affecter	camper	detector	inserter	perfecter
affecters	campers	detectors	inserters	perfecters
affecting	camping	detecting	inserting	perfecting
affectation	campion	detection	insertion	perfection
affectations	campions	detections	insertions	perfections
affects	campus	detects	inserts	perfects
asserted	corrupted	fruited	intercepted	ported
asserter	corrupter	fruiter	interceptor	porter
asserters	corrupters	fruiters	interceptors	porters
asserting	corrupting	fruiting	intercepting	porting

Dawg: directed acyclic word graph

Unify suffixes as well as prefixes!

Directed

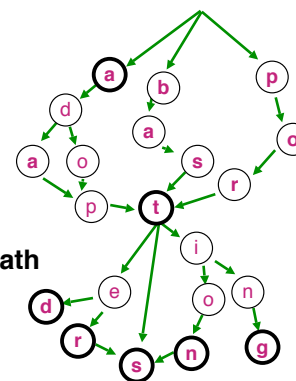
- ◆ Arcs go one-way only

Acyclic

- ◆ No cycles
 - banana, bananan...a

Graph

- ◆ Generalized tree
- ◆ Can be more than one path between nodes



Lexicon as dawg

```
private:
struct node {
    char letter;
    bool isWord;
    int childIndex;
    int nChildren;
};
node *root;
```

F	'A'	'B'	'C'	'D'	'S'
1	27	39	45	55	-1
26	12	6	20	8	0

Unify suffixes as well as prefixes

- ◆ Trie had 250,000 nodes, dawg has just 80,000
- ◆ 125,000 words — so many words have no unique nodes

Space usage?

- ◆ Per-node = 10 bytes
- ◆ Dawg for OSPD2 has ~80,000 nodes = **.8MB**

13

The final result

```
private:
struct node {
    5 bits letter;
    1 bit isWord;
    1 bit for lastChild
    25 bits childIndex;
};
node *root;
```

F	'A'	'B'	'C'	'D'	'S'
1	27	39	45	55	-1
F	F	F	F	F	T

Bitpack node into smallest size

- ◆ Only 26 letters, don't need full ASCII, squeeze in 5 bits
- ◆ isWord/last each 1 bit
- ◆ Remaining 25 bits used for index
- ◆ Each node shrinks from 10 bytes down to 4

Space usage?

- ◆ Dawg for OSPD2 has ~80,000 nodes = **.32 MB**
- ◆ Just ~30% of the size of word file!

14

Cool facts about the dawg

Easily written/read to disk

- ◆ Just take the entire array and write it out as is
- ◆ Indexes don't change, no pointers!
- ◆ Easy to restore
- ◆ That's what the binary "lexicon.dat" file is

Data structure facilitates other interesting operations

- ◆ Tweak width of "beam" as you traverse path to explore just near neighbors
 - Regex, spelling corrections, slight permutations
- ◆ Solving puzzles
 - Anagrams, hangman, scrabble, etc.

Original reference

- ◆ "The World's Fastest Scrabble Player", Appel & Jacobsen, CACM May 1988

15