

15-210

PARALLEL AND SEQUENTIAL
ALGORITHMS AND DATA
STRUCTURES

LECTURE 8

SETS AND TABLES-II

SYNOPSIS

- How search engines work
- Single-threaded sequences

BUILDING A SEARCH ENGINE

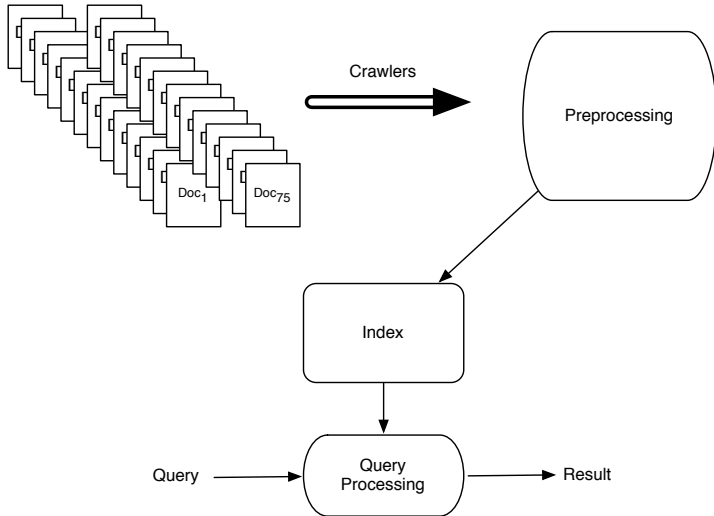
How do search engines work?

- What are the inputs?
 - ▶ (Billions and billions of) documents consisting of “words”.
- How do we interact with the search engine
 - ▶ (Boolean) Keyword queries
 - ▶ List of matching documents (URLS)

HOW DOES THE SEARCH REALLY WORK?

- User inputs a query (say a couple of words)
- SE starts searching for the words in each document one-by-one
- Returns documents when they match.
- Not really!
 - ▶ Not scalable (even for one user)
- Preprocessing

PREPROCESSING



PLAN

- What kinds of queries we want to have.
- What is the interface we want to have.
- How do we implement all these

QUERIES

- **Bingle** (:-) supports logical queries on words involving
 - ▶ And: “15210” And “course” And “slides”
 - ▶ Or: “15210” Or “15150”
 - ▶ AndNot: “15210” AndNot “Pittsburgh”
- “CMU” And “fun” And (“courses Or “clubs”)
- Result would be a list of webpages/documents that match.

THE INTERFACE

```
signature INDEX = sig
  type word = string
  type docId = string
  type 'a seq
  type index
  type docList
```

```
val makeIndex : (docId * string) seq -> index
```

```
val find : index -> word -> docList
```

```
val And : docList * docList -> docList
```

```
val AndNot : docList * docList -> docList
```

```
val Or : docList * docList -> docList
```

```
val size : docList -> int
```

```
val toSeq : docList -> docId seq
```


DOCUMENTS

- Indexing a tweet database.

$$T = \langle (\text{"jack", "chess club was fun"}, \\ (\text{"mary", "I had a fun time in 210 class today"}, \\ (\text{"nick", "food at the cafeteria sucks"}, \\ (\text{"sue", "In 217 class today I had fun reading my email"}, \\ (\text{"peter", "I had fun at nick's party"}, \\ (\text{"john", "tiddlywinks club was no fun, but more fun than 218"} \\ \rangle$$

- “jack” is a document id
- “chess club was fun” is a document

USING THE INTERFACE

$T = \langle$ ("jack", "chess club was fun"),
("mary", "I had a fun time in 210 class today"),
("nick", "food at the cafeteria sucks"),
("sue", "In 217 class today I had fun reading my email"),
("peter", "I had fun at nick's party"),
("john", "tiddlywinks club was no fun, but more fun than 218"),
 \rangle

$f = (\text{find } (\text{makeIndex}(T))) : \text{word} \rightarrow \text{doclist}$

$\text{toSeq}(\text{And}(f \text{ "fun"}, \text{Or}(f \text{ "class"}, f \text{ "club"})))$
 $\Rightarrow \langle \text{"jack"}, \text{"mary"}, \text{"sue"}, \text{"john"} \rangle$

USING THE INTERFACE

```
T = ⟨ ("jack", "chess club was fun"),  
      ("mary", "I had a fun time in 210 class today"),  
      ("nick", "food at the cafeteria sucks"),  
      ("sue", "In 217 class today I had fun reading my email"),  
      ("peter", "I had fun at nick's party"),  
      ("john", "tiddlywinks club was no fun, but more fun than 218"),  
      ⟩  
      size(AndNot(f "fun", f "tiddlywinks"))  
⇒ 4
```

THE MAKEINDEX FUNCTION

```
1 fun makeIndex(docs) =  
2 let  
3   fun tagWords(id, str) =  $\langle (w, id) : w \in \text{tokens}(str) \rangle$   
4   Pairs = flatten  $\langle \text{tagWords}(d) : d \in \text{docs} \rangle$   
5   Words = Table.collect(Pairs)  
6 in  
7    $\{w \mapsto \text{Set.fromSeq}(d) : (w \mapsto d) \in \text{Words}\}$   
8 end
```

- What does tagWords do?

$\text{tagWords}(\text{"jack"}, \text{"chess club was fun"})$
 $\Rightarrow \langle (\text{"chess"}, \text{"jack"}), (\text{"club"}, \text{"jack"}), (\text{"was"}, \text{"jack"}), (\text{"fun"}, \text{"jack"}) \rangle$

THE PAIRS FUNCTION

```
1 fun makeIndex(docs) =  
2 let  
3   fun tagWords(id, str) =  $\langle (w, id) : w \in \text{tokens}(str) \rangle$   
4     Pairs = flatten  $\langle \text{tagWords}(d) : d \in \text{docs} \rangle$   
5     Words = Table.collect(Pairs)  
6 in  
7    $\{w \mapsto \text{Set.fromSeq}(d) : (w \mapsto d) \in \text{Words}\}$   
8 end
```

- What does Pairs do?

$$\text{Pairs} = \langle ("chess", "jack"), ("club", "jack"), ("was", "jack"),$$

$$("fun", "jack"), ("I", "mary"), ("had", "mary"),$$

$$("fun", "mary"), \dots \rangle$$

THE COLLECT FUNCTION

```
1 fun makeIndex(docs) =  
2 let  
3   fun tagWords(id, str) =  $\langle (w, id) : w \in \text{tokens}(str) \rangle$   
4     Pairs = flatten  $\langle \text{tagWords}(d) : d \in \text{docs} \rangle$   
5     Words = Table.collect(Pairs)  
6 in  
7    $\{w \mapsto \text{Set.fromSeq}(d) : (w \mapsto d) \in \text{Words}\}$   
8 end
```

- What does collect do?

```
Words = {("a"  $\mapsto$   $\langle$  "mary"  $\rangle$ ),  
         ("at"  $\mapsto$   $\langle$  "mary", "peter"  $\rangle$ ),  
         ...  
         ("fun"  $\mapsto$   $\langle$  "jack", "mary", "sue", "peter", "john"  $\rangle$ ),  
         ...
```

FINAL TOUCHES

```
1 fun makeIndex(docs) =  
2 let  
3   fun tagWords(id, str) =  $\langle (w, id) : w \in \text{tokens}(str) \rangle$   
4     Pairs = flatten  $\langle \text{tagWords}(d) : d \in \text{docs} \rangle$   
5     Words = Table.collect(Pairs)  
6 in  
7    $\{w \mapsto \text{Set.fromSeq}(d) : (w \mapsto d) \in \text{Words}\}$   
8 end
```

- What is happening here?
- Sequences are converted to tables.

MAKEINDEX COSTS

```
1 fun makeIndex(docs) =  
2 let  
3   fun tagWords(id, str) =  $\langle (w, id) : w \in \text{tokens}(str) \rangle$   
4   Pairs = flatten  $\langle \text{tagWords}(d) : d \in \text{docs} \rangle$   
5   Words = Table.collect(Pairs)  
6 in  
7    $\{w \mapsto \text{Set.fromSeq}(d) : (w \mapsto d) \in \text{Words}\}$   
8 end
```

- Assuming tokens have a upper bound on length
 - ▶ $W_{\text{makeIndex}}(n) \in O(n \log n)$, $S_{\text{makeIndex}} \in O(\log^2 n)$
 - ▶ What does n represent?

REST OF THE INTERFACE

fun *find* T $v = Table.find\ T\ v$

fun *And*(s_1, s_2) = $s_1 \cap s_2$

fun *Or*(s_1, s_2) = $s_1 \cup s_2$

fun *AndNot*(s_1, s_2) = $s_1 \setminus s_2$

fun *size*(s) = $|s|$

fun *toSeq*(s) = *Set.toSeq*(s)

SINGLE-THREADED ARRAY SEQUENCES

- Updating an array sequence in an imperative language takes $O(1)$ work.
- In a functional setting, everything is **persistent**.
- An update to a sequence of n elements needs
 - ▶ $O(n)$ work for `arraySequence` implementation to copy and update.
 - ▶ $O(\log n)$ work for `treeSequence` implementation (because of substructure sharing)
- Interfaces do not provide functions for updating a single position.
 - ▶ to discourage sequential (and expensive) computation.

SINGLE-THREADED ARRAY SEQUENCES

- A *map* can be implemented as follows

```
fun map f S =  
  iter (fn ((i, S'), v)  $\Rightarrow$  (i + 1, update (i, f(v)) S'))  
    (0, S)  
  S
```

- Iterates n times ($i = 0, \dots, n - 1$)
- and updates the value S_i with $f(S_i)$.
- `arraySequence`: Each update will do $O(n)$ work for a total $O(n^2)$ work
- `treeSequence`: Each update will do $O(\log n)$ work for a total $O(n \log n)$ work.

SINGLE-THREADED SEQUENCES

- A new ADT: **Single Threaded Sequence**: `stseq`
- Useful when a bunch of items need to be updated.
- Straightforward interface
- Cost specification imply non-functional stuff under the hood!

STSEQ INTERFACE AND COSTS

	Work	Span
$\text{fromSeq}(S) : \alpha \text{ seq} \rightarrow \alpha \text{ stseq}$ Converts from a regular sequence to a stseq.	$O(S)$	$O(1)$
$\text{toSeq}(ST) : \alpha \text{ stseq} \rightarrow \alpha \text{ seq}$ Converts from a stseq to a regular sequence.	$O(S)$	$O(1)$
$\text{nth } ST \ i : \alpha \text{ stseq} \rightarrow \text{int} \rightarrow \alpha$ Returns the i^{th} element of ST. Same as for seq.	$O(1)$	$O(1)$
$\text{update } (i, v) \ S : (\text{int} \times \alpha) \rightarrow \alpha \text{ stseq} \rightarrow \alpha \text{ stseq}$ Replaces the i^{th} element of S with v .	$O(1)$	$O(1)$
$\text{inject } I \ S : (\text{int} \times \alpha) \text{ seq} \rightarrow \alpha \text{ stseq} \rightarrow \alpha \text{ stseq}$ For each $(i, v) \in I$ replaces the i^{th} element of S with v .	$O(I)$	$O(1)$

- Cost bounds for `nth` and `update` only valid for the “current” version of the sequence.

MAP WITH STSEQ

```
1  fun map f S = let  
2      S' = StSeq.fromSeq(S)  
3      R = iter  
4          (fn ((i, S''), v) => (i + 1, StSeq.update (i, f(v)) S''))  
5          (0, S')  
6          S'  
7  in  
8      StSeq.toSeq(R)  
9  end
```

- Overall work and span is $O(n)$ (Why?)
- Multiple updates can be done in $O(n)$ time.

IMPLEMENTING STSEQ

- Keep two full copies of the sequence
 - ▶ **Original** and **Current**
 - ▶ We keep a **change log**: updates to the original to get **Current**.
- When **Current** is updated
 - ▶ We create a “new” **Current** with the update, and update change log.
 - ▶ Mark the previous version as old, remove its **Current** and but keep the old change log.
- Any item from the current version is accessible in constant work.
- Any item from the any previous version is accessible but needs more work.

IMPLEMENTING STSEQ

Change Log

Original

()

Current

IMPLEMENTING STSEQ

Change Log

Original

()

Current

update(3, 5)

Original

()

Old Version1

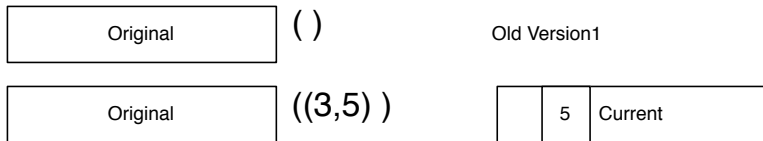
Original

((3,5))

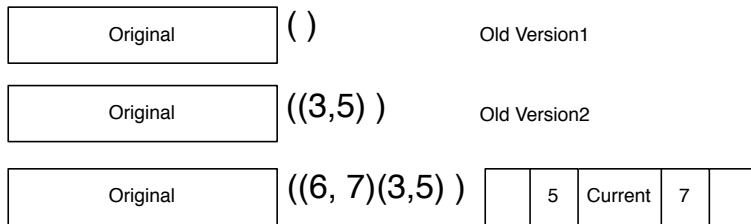
	5	Current
--	---	---------

- There really is only one copy of the **Original**.
- All point to that copy.

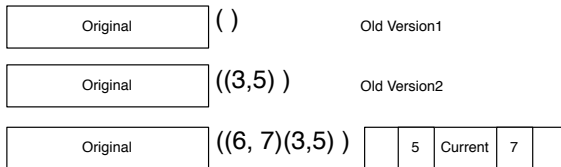
IMPLEMENTING STSEQ



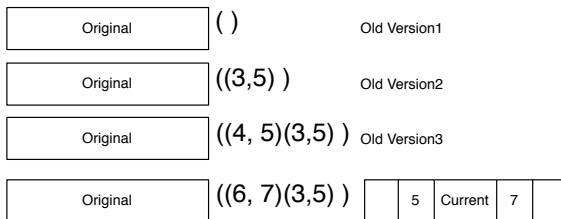
update(6, 7)



IMPLEMENTING STSEQ



update_{Oldversion2}(4, 5)



SUMMARY

- How search engines work
- Single-threaded sequences