

15-210

PARALLEL AND SEQUENTIAL  
ALGORITHMS AND DATA  
STRUCTURES

LECTURE 7

COLLECT, SETS AND TABLES

# SYNOPSIS

- The **collect** operation
- The **map-collect-reduce** paradigm
- Sets
- Tables

# THE COLLECT OPERATION

- Group items that share a **common key**.

```
Data = ⟨("jack sprat", "15-210"),  
        ("jack sprat", "15-213"),  
        ("mary contrary", "15-210"),  
        ("mary contrary", "15-251"),  
        ("mary contrary", "15-213"),  
        ("peter piper", "15-150"),  
        ("peter piper", "15-251"),  
        ...⟩
```



```
rosters = ⟨("15-150", ⟨"peter piper", ...⟩)  
          ("15-210", ⟨"jack sprat", "mary contrary", ...⟩)  
          ("15-213", ⟨"jack sprat", ...⟩)  
          ("15-251", ⟨"mary contrary", "peter piper"⟩)  
          ...⟩
```

# THE COLLECT OPERATION

- Very common operation in **Relational Databases**
- Usually called the **Group by** operation.

```
rosters = ⟨("15-150", ⟨ "peter piper", ... ⟩)  
           ("15-210", ⟨ "jack sprat", "mary contrary", ... ⟩)  
           ("15-213", ⟨ "jack sprat", ... ⟩)  
           ("15-251", ⟨ "mary contrary", "peter piper" ⟩)  
           ... ⟩
```

- **Students** are grouped by **Course Numbers**.

# THE COLLECT OPERATION

$$\begin{aligned} collect : (\alpha \times \alpha \rightarrow \textit{order}) &\rightarrow (\alpha \times \beta) \textit{ seq} \\ &\rightarrow (\alpha \times \beta \textit{ seq}) \textit{ seq} \end{aligned}$$

- 1  $\alpha \times \alpha \rightarrow \textit{order}$  is a function for comparing keys of type  $\alpha$
- 2  $(\alpha \times \beta) \textit{ seq}$  is a sequence of **key-value** pairs
- 3  $(\alpha \times \beta \textit{ seq}) \textit{ seq}$  is the resulting sequence:
  - ▶ each unique  $\alpha$  value is paired with **a sequence of all  $\beta$  values it appears with**

# THE COLLECT OPERATION

```
collectStrings = collect String.compare  
rosters = collectStrings( $\langle (n, c) : (c, n) \in \text{Data} \rangle$ )
```

```
rosters =  $\langle$  ("15-150",  $\langle$  "peter piper", ...  $\rangle$ )  
           ("15-210",  $\langle$  "jack sprat", "mary contrary", ...  $\rangle$ )  
           ("15-213",  $\langle$  "jack sprat", ...  $\rangle$ )  
           ("15-251",  $\langle$  "mary contrary", "peter piper"  $\rangle$ )  
           ...  $\rangle$ 
```

- $\langle (n, c) : (c, n) \in \text{Data} \rangle$  arranges the data appropriately.

# THE COLLECT OPERATION

- How would you implement `collect`?
  - ▶ Sort the items on their keys
  - ▶ Partition the resulting sequence
  - ▶ Pull out pairs between each key change

# THE COLLECT OPERATION

- The dominant cost of `collect` is in sorting.
- Work is  $O(W_c n \log n)$ , Span is  $O(S_c \log^2 n)$ 
  - ▶  $W_c$  work bound for the comparison function
  - ▶  $S_c$  span bound for the comparison function
- A  $O(n)$  work can be implemented with **hashing**.
  - ▶ Need a separate hash function
  - ▶ Output not in sorted order



# USING COLLECT IN MAP-REDUCE

- The **map-reduce paradigm** is used to process very large collection of **documents**.
  - ▶ A document is a collection of **words/strings**.
  - ▶ Not the `mapReduce` of 15-150!
- **map-reduce paradigm**  $\equiv$  `map-collect-reduce(s)`.

# USING COLLECT IN MAP-REDUCE

- $f_m$  maps each document to a sequence of **key-value** pairs.
  - ▶  $f_m : (\text{document} \rightarrow (\text{key} \times \alpha) \text{ seq})$
- All key-value pairs in a document are **collected**.
- $f_r$  is applied to the keys to get a single value for a key.
  - ▶  $f_r : \text{key} \times \alpha \text{ seq} \rightarrow \beta$

# AN EXAMPLE

$docs = \langle \text{"this is a document"},$   
 $\text{"this is is another document"},$   
 $\text{"a last document"} \rangle$



$\langle (\text{"this"}, 1), (\text{"is"}, 1), (\text{"a"}, 1), (\text{"document"}, 1),$   
 $(\text{"this"}, 1), (\text{"is"}, 1), (\text{"is"}, 1), (\text{"another"}, 1),$   
 $(\text{"document"}, 1), (\text{"a"}, 1), (\text{"last"}, 1), (\text{"document"}, 1) \rangle$



$\langle (\text{"a"}, 2), (\text{"another"}, 1), (\text{"document"}, 3), (\text{"is"}, 3), (\text{"last"}, 1),$   
 $(\text{"this"}, 2) \rangle$

# MAPREDUCE IN SML

```
1  fun mapCollectReduce  $f_m$   $f_r$  docs =  
2    let  
3      pairs = flatten  $\langle f_m(s) : s \in docs \rangle$   
4      groups = collect String.compare pairs  
5    in  
6       $\langle f_r(g) : g \in groups \rangle$   
7    end
```

- $flatten \langle \langle a, b, c \rangle, \langle d, e \rangle \rangle \Rightarrow \langle a, b, c, d, e \rangle$

# MAPREDUCE IN SML

```
1  fun mapCollectReduce  $f_m$   $f_r$  docs =  
2    let  
3      pairs = flatten  $\langle f_m(s) : s \in docs \rangle$   
4      groups = collect String.compare pairs  
5    in  
6       $\langle f_r(g) : g \in groups \rangle$   
7    end
```

**fun**  $f_m(doc) = \langle (w, 1) : tokens\ doc \rangle$

**fun**  $f_r(w, s) = (w, reduce + 0\ s)$

# MAPREDUCE EXAMPLE IN SML

**fun**  $f_m(doc) = \langle (w, 1) : tokens\ doc \rangle$

**fun**  $f_r(w, s) = (w, reduce + 0\ s)$

$countWords = mapCollectReduce\ f_m\ f_r$

$countWords\ \langle "this\ is\ a\ document",$   
                   $"this\ is\ is\ another\ document",$   
                   $"a\ last\ document" \rangle$

$\Rightarrow \langle ("a", 2), ("another", 1), ("document", 3), ("is", 3),$   
           $("last", 1), ("this", 2) \rangle$

# SETS

- Sets play a very important role in math.
- Often needed in many algorithms.
- Many languages either support sets directly or have libraries for sets.
- In 15-210 we use a purely functional definition for sets:
  - ▶ When updates are done, a new set is returned.

# SETS AS AN ADT

- $\mathbb{U}$  is a universe of elements.
- The SET ADT is a type  $\mathbb{S}$  that represents the power set of  $\mathbb{U}$ .

$$\begin{array}{llll} \text{empty} & : \mathbb{S} & = & \emptyset \\ \text{size}(\mathbf{S}) & : \mathbb{S} \rightarrow \mathbb{Z}_{\geq 0} & = & |\mathbf{S}| \\ \text{singleton}(\mathbf{e}) & : \mathbb{U} \rightarrow \mathbb{S} & = & \{\mathbf{e}\} \\ \text{filter}(\mathbf{f}, \mathbf{S}) & : ((\mathbb{U} \rightarrow \{\mathbf{T}, \mathbf{F}\}) & = & \{\mathbf{s} \in \mathbf{S} \mid \mathbf{f}(\mathbf{s})\} \\ & \times \mathbb{S}) \rightarrow \mathbb{S} & & \end{array}$$



# SETS AS AN ADT

$$\begin{aligned} \text{find}(\mathcal{S}, e) &: \mathcal{S} \times \mathcal{U} &= |\{s \in \mathcal{S} \mid s = e\}| = 1 \\ &\rightarrow \{T, F\} \end{aligned}$$

$$\text{insert}(\mathcal{S}, e) : \mathcal{S} \times \mathcal{U} \rightarrow \mathcal{S} = \mathcal{S} \cup \{e\}$$

$$\text{delete}(\mathcal{S}, e) : \mathcal{S} \times \mathcal{U} \rightarrow \mathcal{S} = \mathcal{S} \setminus \{e\}$$

$$\text{intersection}(\mathcal{S}_1, \mathcal{S}_2) : \mathcal{S} \times \mathcal{S} \rightarrow \mathcal{S} = \mathcal{S}_1 \cap \mathcal{S}_2$$

$$\text{union}(\mathcal{S}_1, \mathcal{S}_2) : \mathcal{S} \times \mathcal{S} \rightarrow \mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$$

$$\text{difference}(\mathcal{S}_1, \mathcal{S}_2) : \mathcal{S} \times \mathcal{S} \rightarrow \mathcal{S} = \mathcal{S}_1 \setminus \mathcal{S}_2$$

- What is the relationship between these two groups?

# SETS AS AN ADT

- We do not really need `find`, `insert`, `delete`!

`find(S, e) = size(intersection(S, singleton(e))) = 1`

`insert(S, e) = union(S, singleton(e))`

`delete(S, e) = difference(S, singleton(e))`

- `intersection`, `union`, **and** `difference`
  - ▶ can operate on multiple elements, and
  - ▶ are suitable for parallelism

# COST MODEL FOR SETS

- Underlying data structure can be
  - ▶ hash-tables
  - ▶ balanced trees
- We will assume a balanced-tree implementation.
- We will assume comparison of two set elements take
  - ▶  $W_c$  work and  $S_c$  span.

# COST MODEL FOR SETS

	<i>Work</i>	<i>Span</i>
size( $S$ ) singleton( $e$ )	$O(1)$	$O(1)$
filter( $f, S$ )	$O\left(\sum_{e \in S} W(f(e))\right)$	$O\left(\log  S  + \max_{e \in S} S(f(e))\right)$
find( $S, e$ ) insert( $S, e$ ) delete( $S, e$ )	$O(W_c \cdot \log  S )$	$O(S_c \cdot \log  S )$

# COST MODELS FOR SETS

$$\begin{array}{ll} \text{intersection}(S_1, S_2) & \text{Work} = O(W_c \cdot m \cdot \log(1 + \frac{n}{m})) \\ \text{union}(S_1, S_2) & \Rightarrow \\ \text{difference}(S_1, S_2) & \text{Span} = O(S_c \cdot \log(n + m)) \end{array}$$

$$n = \max(|S_1|, |S_2|)$$

$$m = \min(|S_1|, |S_2|)$$

- Sets are equal size ( $n = m$ )
  - ▶  $\text{Work} = O(W_c \cdot m \cdot \log(1 + 1)) = O(W_c \cdot n)$
  - ▶  $\text{Span} = O(S_c \cdot \log n)$
- One of the sets is a singleton ( $m = 1$ )
  - ▶  $\text{Work} = O(W_c \cdot \log(1 + n)) = O(W_c \cdot \log n)$
  - ▶  $\text{Span} = O(S_c \cdot \log(n + 1)) = O(S_c \cdot \log n)$

# TABLES

- Table is an ADT for sets of **key-value** pairs.
- $\{(k_1 \mapsto v_1), (k_2 \mapsto v_2), \dots, (k_n \mapsto v_n)\}$
- $\{(k_1, v_1), (k_2, v_2), \dots, (k_n, v_n)\}$
- Each **key** appears only once
- Many languages provide either built-in support or libraries.

# TABLES

- $\mathbb{K}$  is the universe of keys.
- $\mathbb{V}$  is the universe of values.
- $\mathbb{T}$  is a type that represents the power set of  $\mathbb{K} \times \mathbb{V}$ 
  - ▶ restricted so that each key appears at most once.
  - ▶  $\mathbb{S}$  is the power set of  $\mathbb{K}$ .
  - ▶  $\mathbb{Z}_{\geq 0}$  denotes the positive integers.

# TABLE FUNCTIONS

<code>empty</code>	:	$\mathbb{T}$	=	$\emptyset$
<code>size(<math>T</math>)</code>	:	$\mathbb{T} \rightarrow \mathbb{Z}_{\geq 0}$	=	$ T $
<code>singleton(<math>k, v</math>)</code>	:	$\mathbb{K} \times \mathbb{V} \rightarrow \mathbb{T}$	=	$\{(k, v)\}$
<code>filter(<math>f, T</math>)</code>	:	$((\mathbb{V} \rightarrow \{\mathbb{T}, \mathbb{F}\}) \times \mathbb{T})$ $\rightarrow \mathbb{T}$	=	$\{(k, v) \in T \mid f(v)\}$
<code>map(<math>f, T</math>)</code>	:	$((\mathbb{K} \times \mathbb{V} \rightarrow \mathbb{V}) \times \mathbb{T})$ $\rightarrow \mathbb{T}$	=	$\{(k, f(k, v)) \mid ((k, v) \in T)\}$
<code>insert(<math>f, T, (k, v)</math>)</code>	:	$(\mathbb{V} \times \mathbb{V} \rightarrow \mathbb{V}) \times \mathbb{T}$ $\times (\mathbb{K} \times \mathbb{V}) \rightarrow \mathbb{T}$	=	$\begin{aligned} & (T \setminus \{(k, v)\}) \cup \\ & \quad \{(k, f(v, v'))\} \\ & \quad \text{if } (k, v') \in T \\ & T \cup \{(k, v)\} \\ & \quad \text{otherwise} \end{aligned}$
<code>delete(<math>T, k</math>)</code>	:	$\mathbb{T} \times \mathbb{K} \rightarrow \mathbb{T}$	=	$\{(k', v) \in T \mid k \neq k'\}$



# TABLE FUNCTIONS

$$\begin{aligned}\text{find}(T, k) &: \mathbb{T} \times \mathbb{K} \rightarrow (\mathbb{V} \cup \perp) &= \begin{cases} v & (k, v) \in T \\ \perp & \text{otherwise} \end{cases} \\ \text{merge}(f, T_1, T_2) &: (\mathbb{V} \times \mathbb{V} \rightarrow \mathbb{V}) \times \mathbb{T} \times \mathbb{T} \rightarrow \mathbb{T} &= \\ &\bigcup_{k \in \mathbb{K}} \begin{cases} \{(k, f(v_1, v_2))\} & (k, v_1) \in T_1 \wedge (k, v_2) \in T_2 \\ \{(k, v_1)\} & (k, v_1) \in T_1 \wedge (k, v_2) \notin T_2 \\ \{(k, v_2)\} & (k, v_2) \in T_2 \wedge (k, v_1) \notin T_1 \end{cases} \\ \text{extract}(T, S) &: \mathbb{T} \times \mathbb{S} \rightarrow \mathbb{T} &= \{(k, v) \in T \mid k \in S\} \\ \text{erase}(T, S) &: \mathbb{T} \times \mathbb{S} \rightarrow \mathbb{T} &= \{(k, v) \in T \mid k \notin S\} \end{aligned}$$

# TABLE EXAMPLES

- Suppose we have the two tables:
  - ▶  $Summer = \{tree \mapsto green, sky \mapsto blue, cmuq \mapsto tan\}$
  - ▶  $Fall = \{grass \mapsto gray, tree \mapsto brown\}$
- $merge (fn (a, b) \Rightarrow b) Summer Fall$ 
  - ▶  $\{grass \mapsto gray, tree \mapsto brown, sky \mapsto blue, cmuq \mapsto tan\}$

# TABLE EXAMPLES

- Suppose we have the two tables:
  - ▶  $Summer = \{tree \mapsto green, sky \mapsto blue, cmuq \mapsto tan\}$
  - ▶  $Fall = \{grass \mapsto gray, tree \mapsto brown\}$
- $extract(Summer, \{sky, grass\})$ 
  - ▶  $\{sky \mapsto blue\}$

# TABLE EXAMPLES

- Suppose we have the two tables:
  - ▶  $Summer = \{tree \mapsto green, sky \mapsto blue, cmuq \mapsto tan\}$
  - ▶  $Fall = \{grass \mapsto gray, tree \mapsto brown\}$
- $erase(Summer, \{sky, grass\})$ 
  - ▶  $\{tree \mapsto green, cmuq \mapsto tan\}$

# TABLE EXAMPLES

- Other useful functions from the library
- $\text{collect}:(\text{key} \times \alpha) \text{ seq} \rightarrow (\alpha \text{ seq}) \text{ table}$
- $\text{fromSeq}:(\text{key} \times \alpha) \text{ seq} \rightarrow \alpha \text{ table}$ 
  - ▶  $\text{fromSeq}(A) = \{k \mapsto s_0 : (k \mapsto S) \in \text{collect}(A)\}$

# TABLE FUNCTIONS

- Major differences from sets:
  - ▶ `find` returns the value if key is in the table else returns  $\perp$  (NONE).
  - ▶ `insert/merge` need a function to combine if the key is already in the/both table(s).
- Just as with sets, there is symmetry between
  - ▶ `extract` and `find`
  - ▶ `merge` and `insert`
  - ▶ `erase` and `delete`

# COST MODELS FOR TABLES

	<i>Work</i>	<i>Span</i>
$\text{size}(T)$ $\text{singleton}(k, v)$	$O(1)$	$O(1)$
$\text{filter}(f, T)$	$O\left(\sum_{(k,v) \in T} W(f(v))\right)$	$O\left(\log  T  + \max_{(k,v) \in T} S(f(v))\right)$
$\text{map}(f, T)$	$O\left(\sum_{(k,v) \in T} W(f(k, v))\right)$	$O\left(\max_{(k,v) \in T} S(f(k, v))\right)$

# COST MODELS FOR TABLES

	<i>Work</i>	<i>Span</i>
<code>find(<math>S, k</math>)</code>		
<code>insert(<math>T, (k, v)</math>)</code>	$O(W_c \log  T )$	$O(S_c \log  T )$
<code>delete(<math>T, k</math>)</code>		
<hr/>		
<code>extract(<math>T_1, T_2</math>)</code>		
<code>merge(<math>T_1, T_2</math>)</code>	$O(W_c m \log(1 + \frac{n}{m}))$	$O(S_c \log(n + m))$
<code>erase(<math>T_1, T_2</math>)</code>		
<hr/>		
$n = \max( T_1 ,  T_2 ) \quad m = \min( T_1 ,  T_2 )$		



# SUMMARY

- Collect
- Map-Collect-Reduce
- Sets
- Tables