

Simple aggregation Recall the definitions:

```
type Partition a = [a]
type RDD a = [Partition a]
```

```
aggregate :: b → (b → a → b) → (b → b → b) → RDD a → b
aggregate z (⊗) (⊕) = foldl (⊕) z ∘ map (foldl (⊗) z)
```

Aggregation for pair RDDs Extending aggregation to pair RDDs:

```
type PairRDD a b = RDD (a, b)
```

```
key :: (a, b) → a
key = fst
value :: (a, b) → b
value = snd
```

The functional call `hasValue k v ps` locates the rightmost occurrence of (k, u) in `ps`, if any, and returns `u`. Otherwise it returns the default value `v`.

```
hasValue :: Eq a ⇒ a → b → [(a, b)] → b
hasValue k v ps = last_v (map value (filter ((≡ k) ∘ key) ps))
```

where $\text{last}_v = \text{last} \circ (v:)$ (thus $\text{last}_v [] = v$). The definition above is equivalent to $\text{foldl } (\lambda r p \rightarrow \text{if } k \equiv \text{key } p \text{ then value } p \text{ else } r) v$. Also, let $\text{keyEq } k = (\equiv k) \circ \text{key}$, since we use it often

The function `addTo` updates a key-value list by a given key-value pair.

```
addTo :: Eq a ⇒ a → b → [(a, b)] → [(a, b)]
addTo k v ps = (k, v) : filter (¬ ∘ (≡ k) ∘ key) ps
```

Given $(\odot) :: c \rightarrow b \rightarrow c$, $(\hat{\odot}_z) :: [(a, c)] \rightarrow (a, b) \rightarrow [(a, c)]$ is (\odot) lifted to list of key-value pairs:

```
∘ :: Eq a ⇒ (c → b → c) → c → [(a, c)] → (a, b) → [(a, c)]
r ∘_z (k, v) = addTo k (hasValue k z r ∘ v) r
```

It is helpful to give $(\hat{\odot}_z)$ an inductive characterization:

```
[ ]      ∘_z (k, v) = [(k, z ∘ v)]
(xs ++ [(j, u)]) ∘_z (k, v) | k ≡ j      = xs ++ [(k, u ∘ v)]
                        | otherwise = (xs ∘_z (k, v)) ++ [(j, u)]
```

The function `repartition` is unspecified. The following dummy definition is merely there to allow the code to type check.

$\text{repartition} :: [(a, b)] \rightarrow \text{PairRDD } a \ b$
 $\text{repartition } xs = [xs]$

Finally, `aggregateByKey` may be defined by:

$\text{aggregateByKey} :: \text{Eq } a \Rightarrow c \rightarrow (c \rightarrow b \rightarrow c) \rightarrow (c \rightarrow c \rightarrow c) \rightarrow$
 $\text{PairRDD } a \ b \rightarrow \text{PairRDD } a \ c$
 $\text{aggregateByKey } z \ (\otimes) \ (\oplus) =$
 $\text{repartition} \circ \text{foldl } (\hat{\oplus}_z) \ [] \circ \text{concat} \circ \text{map } (\text{foldl } (\hat{\otimes}_z) \ [])$

Aggregation with a given key The function `aggregateWithKey`

$\text{aggregateWithKey} :: \text{Eq } a \Rightarrow$
 $a \rightarrow c \rightarrow (c \rightarrow b \rightarrow c) \rightarrow (c \rightarrow c \rightarrow c) \rightarrow \text{PairRDD } a \ b \rightarrow c$
 $\text{aggregateWithKey } k \ z \ (\otimes) \ (\oplus) =$
 $\text{aggregate } z \ (\otimes) \ (\oplus) \circ$
 $\text{filter } (\neg \circ \text{null}) \circ$
 $\text{map } (\text{map value} \circ \text{filter } (\text{keyEq } k))$

$\text{lookUp} :: \text{Eq } a \Rightarrow a \rightarrow b \rightarrow \text{PairRDD } a \ b \rightarrow b$
 $\text{lookUp } k \ z = \text{last}_z \circ \text{concat} \circ \text{map } (\text{map value} \circ \text{filter } ((\equiv k) \circ \text{key}))$

$\text{lookUpL} :: \text{Eq } a \Rightarrow a \rightarrow b \rightarrow [(a, b)] \rightarrow b$
 $\text{lookUpL } k \ z = \text{last}_z \circ \text{map value} \circ \text{filter } ((\equiv k) \circ \text{key})$

Lemma 1. *For all j, k, xs, v, z , and (\odot) , we have:*

1. $\text{filter } (\text{keyEq } k) (xs \hat{\odot}_z (k, v)) = \text{filter } (\text{keyEq } k) xs \hat{\odot}_z (k, v).$
2. $\text{filter } (\text{keyEq } k) (xs \hat{\odot}_z (j, v)) = \text{filter } (\text{keyEq } k) xs$, if $k \neq j$.

Proof. We prove 1. only. **Case** $xs := []$.

$\text{filter } (\text{keyEq } k) ([] \hat{\odot}_z (k, v))$
 $= \text{filter } (\text{keyEq } k) [(k, z \odot v)]$
 $= [(k, z \odot v)]$
 $= [] \hat{\odot}_z (k, v)$
 $= \text{filter } (\text{keyEq } k) [] \hat{\odot}_z (k, v) .$

Case $xs := xs \uparrow [(k, u)]$.

$\text{filter } (\text{keyEq } k) ((xs \uparrow [(k, u)]) \hat{\odot}_z (k, v))$
 $= \{ \text{inductive definition of } (\hat{\odot}_z) \}$
 $\text{filter } (\text{keyEq } k) (xs \uparrow [(k, u \odot v)])$
 $= \text{filter } (\text{keyEq } k) xs \uparrow [(k, u \odot v)]$

$$\begin{aligned}
&= \{ \text{inductive definition of } (\hat{\odot}_z) \} \\
&\quad (\text{filter } (\text{keyEq } k) \text{ } xs \mathbin{++} [(k, u)]) \hat{\odot}_z (k, v) \\
&= \text{filter } (\text{keyEq } k) (xs \mathbin{++} [(k, u)]) \hat{\odot}_z (k, v) .
\end{aligned}$$

Case $xs := xs \mathbin{++} [(j, u)]$, where $j \neq k$.

$$\begin{aligned}
&\quad \text{filter } (\text{keyEq } k) ((xs \mathbin{++} [(j, u)]) \hat{\odot}_z (k, v)) \\
&= \{ \text{inductive definition of } (\hat{\odot}_z) \} \\
&\quad \text{filter } (\text{keyEq } k) ((xs \hat{\odot}_z (k, v)) \mathbin{++} [(j, u)]) \\
&= \text{filter } (\text{keyEq } k) (xs \hat{\odot}_z (k, v)) \\
&= \{ \text{induction} \} \\
&\quad \text{filter } (\text{keyEq } l) xs \hat{\odot}_z (k, v) \\
&= \text{filter } (\text{keyEq } k) (xs \mathbin{++} [(j, u)]) \hat{\odot}_z (k, v) .
\end{aligned}$$

□

Lemma 2. For all k , (\oplus) , and z , we have

$$\text{filter } (\text{keyEq } k) \circ \text{foldl } (\hat{\oplus}_z) [] = \text{foldl } (\hat{\oplus}_z) [] \circ \text{filter } (\text{keyEq } k) .$$

Proof. Induction on the input.

Case $xs := []$. Both sides reduces to $[]$.

Case $xs := xs \mathbin{++} [(k, v)]$

$$\begin{aligned}
&\quad \text{filter } (\text{keyEq } k) (\text{foldl } (\hat{\oplus}_z) [] (xs \mathbin{++} [(k, v)])) \\
&= \text{filter } (\text{keyEq } k) (\text{foldl } (\hat{\oplus}_z) [] xs \hat{\oplus}_z (k, v)) \\
&= \{ \text{Lemma 1} \} \\
&\quad \text{filter } (\text{keyEq } k) (\text{foldl } (\hat{\oplus}_z) [] xs) \hat{\oplus}_z (k, v) \\
&= \{ \text{induction} \} \\
&\quad \text{foldl } (\hat{\oplus}_z) [] (\text{filter } (\text{keyEq } k) xs) \hat{\oplus}_z (k, v) \\
&= \text{foldl } (\hat{\oplus}_z) [] (\text{filter } (\text{keyEq } k) (xs \mathbin{++} [(k, v)]))
\end{aligned}$$

Case $xs := xs \mathbin{++} [(j, v)]$, where $j \neq k$.

$$\begin{aligned}
&\quad \text{filter } (\text{keyEq } k) (\text{foldl } (\hat{\oplus}_z) [] (xs \mathbin{++} [(j, v)])) \\
&= \text{filter } (\text{keyEq } k) (\text{foldl } (\hat{\oplus}_z) [] xs \hat{\oplus}_z (j, v)) \\
&= \{ \text{Lemma 1} \} \\
&= \text{filter } (\text{keyEq } k) (\text{foldl } (\hat{\oplus}_z) [] xs) \\
&= \{ \text{induction} \} \\
&\quad \text{foldl } (\hat{\oplus}_z) [] (\text{filter } (\text{keyEq } k) xs) \\
&= \text{foldl } (\hat{\oplus}_z) [] (\text{filter } (\text{keyEq } k) (xs \mathbin{++} [(j, v)]))
\end{aligned}$$

□

Lemma 3. For all k , (\odot) , and z , we have

$$\text{foldl } (\hat{\odot}_z) [] \circ \text{filter } (\text{keyEq } k) \$ xs = \\ [(k, \text{foldl } (\odot) z \circ \text{map value} \circ \text{filter } (\text{keyEq } k) \$ xs)]$$

if there exists at least one (k, v) in xs . Otherwise $\text{foldl } (\hat{\odot}_z) [] \circ \text{filter } (\text{keyEq } k) \$ xs = []$.

Proof. The "otherwise" part is trivially true. For the first part, we perform an induction on the input.

Case $xs := [(k, v)]$. Both sides reduce to $[(k, z \odot v)]$.

Case $xs := xs \# [(k, v)]$.

- **Case** there exists at least one $[(k, u)]$ in xs .

$$\begin{aligned} & \text{foldl } (\hat{\odot}_z) [] (\text{filter } (\text{keyEq } k) (xs \# [(k, v)])) \\ &= \text{foldl } (\hat{\odot}_z) [] (\text{filter } (\text{keyEq } k) xs \# [(k, v)]) \\ &= \text{foldl } (\hat{\odot}_z) [] (\text{filter } (\text{keyEq } k) xs) \hat{\odot}_z (k, v) \\ &= \{ \text{induction} \} \\ & \quad [(k, \text{foldl } (\odot) z \circ \text{map value} \circ \text{filter } (\text{keyEq } k) \$ xs)] \hat{\odot}_z (k, v) \\ &= \{ \text{definition of } (\hat{\oplus}_z) \} \\ & \quad [(k, (\text{foldl } (\odot) z \circ \text{map value} \circ \text{filter } (\text{keyEq } k) \$ xs) \odot v)] \\ &= [(k, \text{foldl } (\odot) z \circ \text{map value} \circ \text{filter } (\text{keyEq } k) \$ (xs \# [(k, v)]))] . \end{aligned}$$

- **Case** $[(k, u)]$ does not occur in xs .

$$\begin{aligned} & \text{foldl } (\hat{\odot}_z) [] (\text{filter } (\text{keyEq } k) (xs \# [(k, v)])) \\ &= \text{foldl } (\hat{\odot}_z) [] (\text{filter } (\text{keyEq } k) xs \# [(k, v)]) \\ &= \text{foldl } (\hat{\odot}_z) [] (\text{filter } (\text{keyEq } k) xs) \hat{\odot}_z (k, v) \\ &= \{ \text{filter } (\text{keyEq } k) xs = [] \} \\ &= [] \hat{\odot}_z (k, v) \\ &= [(k, z \odot v)] \\ &= [(k, \text{foldl } (\odot) z \circ \text{map value} \circ \text{filter } (\text{keyEq } k) \$ (xs \# [(k, v)]))] . \end{aligned}$$

Case $xs := xs \# [(j, v)]$ where $j \neq k$. In this case there must exists some $[(k, u)]$ in xs .

$$\begin{aligned} & \text{foldl } (\hat{\oplus}_z) [] (\text{filter } (\text{keyEq } k) (xs \# [(j, v)])) \\ &= \text{foldl } (\hat{\oplus}_z) [] (\text{filter } (\text{keyEq } k) xs) \\ &= \{ \text{induction} \} \\ & \quad [(k, \text{foldl } (\odot) z \circ \text{map value} \circ \text{filter } (\text{keyEq } k) \$ xs)] \\ &= [(k, \text{foldl } (\odot) z \circ \text{map value} \circ \text{filter } (\text{keyEq } k) \$ (xs \# [(k, v)]))] . \end{aligned}$$

□

Corollary 4. *As a corollary, we have that for all k , (\odot) , and z ,*

$$\text{last_z} \circ \text{map value} \circ \text{foldl } (\hat{\odot}_z) [] \circ \text{filter (keyEq } k) = \\ \text{foldl } (\odot) z \circ \text{map value} \circ \text{filter (keyEq } k) \quad .$$

When the input is empty or does not contain entries having key k , both sides reduce to z .

Corollary 5. *For all k , z and (\odot) we have:*

$$\text{concat} \circ \text{map (map value} \circ \text{filter (keyEq } k) \circ \text{foldl } (\hat{\odot}_z) []) = \\ \text{map (foldl } (\odot) z) \circ \text{filter } (\neg \circ \text{null}) \circ \text{map filter (keyEq } k) \quad .$$

Proof. We reason:

$$\begin{aligned} & \text{concat} \circ \text{map (map value} \circ \text{filter (keyEq } k) \circ \text{foldl } (\hat{\odot}_z) []) \\ &= \{ \text{Lemma 2} \} \\ & \text{concat} \circ \text{map (map value} \circ \text{foldl } (\hat{\oplus}_z) [] \circ \text{filter (keyEq } k)) \end{aligned}$$

Let the input be xss . Induction.

Case $xss := []$. Both sides reduce to $[]$

Case $xss := xss \mathbin{++} [xs]$

- there exists at least one (k, v) in xs

$$\begin{aligned} & \text{concat} \circ \text{map (map value} \circ \text{foldl } (\hat{\oplus}_z) [] \circ \text{filter (keyEq } k)) \$ (xss \mathbin{++} [xs]) \\ &= (\text{concat} \circ \text{map (map value} \circ \text{foldl } (\hat{\oplus}_z) [] \circ \text{filter (keyEq } k)) \$ xss) \mathbin{++} \\ & \quad (\text{map value} \circ \text{foldl } (\hat{\oplus}_z) \circ \text{filter (keyEq } k) \$ xs) \\ &= \{ \text{induction} \} \\ & \quad (\text{map (foldl } (\odot) z) \circ \text{filter } (\neg \circ \text{null}) \circ \text{map (map value} \circ \text{filter (keyEq } k)) \$ xss) \mathbin{++} \\ & \quad (\text{map value} \circ \text{foldl } (\hat{\oplus}_z) \circ \text{filter (keyEq } k) \$ xs) \\ &= \{ \text{Lemma 3} \} \\ & \quad (\text{map (foldl } (\odot) z) \circ \text{filter } (\neg \circ \text{null}) \circ \text{map (map value} \circ \text{filter (keyEq } k)) \$ xss) \mathbin{++} \\ & \quad [\text{foldl } (\odot) z \circ \text{map value} \circ \text{filter (keyEq } k) \$ xs] \\ &= \text{map (foldl } (\odot) z) \circ \text{filter } (\neg \circ \text{null}) \circ \text{map (map value} \circ \text{filter (keyEq } k)) \$ (xss \mathbin{++} [xs]) \end{aligned}$$

- there exists no (k, v) in xs .

$$\begin{aligned} & \text{concat} \circ \text{map (map value} \circ \text{foldl } (\hat{\oplus}_z) [] \circ \text{filter (keyEq } k)) \$ (xss \mathbin{++} [xs]) \\ &= \{ \text{same as above} \} \\ & \quad (\text{map (foldl } (\odot) z) \circ \text{filter } (\neg \circ \text{null}) \circ \text{map (map value} \circ \text{filter (keyEq } k)) \$ xss) \mathbin{++} \\ & \quad (\text{map value} \circ \text{foldl } (\hat{\oplus}_z) \circ \text{filter (keyEq } k) \$ xs) \\ &= \{ \text{Lemma 3} \} \\ & \quad \text{map (foldl } (\odot) z) \circ \text{filter } (\neg \circ \text{null}) \circ \text{map (map value} \circ \text{filter (keyEq } k)) \$ xss \\ &= \text{map (foldl } (\odot) z) \circ \text{filter } (\neg \circ \text{null}) \circ \text{map (map value} \circ \text{filter (keyEq } k)) \$ (xss \mathbin{++} [xs]) \quad . \end{aligned}$$

□

Finally, the main theorem:

Theorem 6. *For all $k, z, (\otimes)$ and (\oplus) , we have:*

$$\text{lookUp } k \circ \text{aggregateByKey } z (\otimes) (\oplus) = \text{aggregateWithKey } k z (\otimes) (\oplus) .$$

Proof. We reason:

$$\begin{aligned}
& \text{lookUp } k \circ \text{aggregateByKey } z (\otimes) (\oplus) \\
= & \quad \{ \text{definition of lookUp and aggregateByKey} \} \\
& \text{last}_z \circ \text{concat} \circ \text{map} (\text{map value} \circ \text{filter} (\text{keyEq } k)) \circ \\
& \quad \text{repartition} \circ \text{foldl} (\hat{\oplus}_z) [] \circ \text{concat} \circ \text{map} (\text{foldl} (\hat{\otimes}_z) []) \\
= & \quad \{ \text{routine naturality laws. See below.} \} \\
& \text{last}_z \circ \text{map value} \circ \text{filter} (\text{keyEq } k) \circ \text{foldl} (\hat{\oplus}_z) [] \circ \\
& \quad \text{concat} \circ \text{map} (\text{foldl} (\hat{\otimes}_z) []) \\
= & \quad \{ \text{Lemma 2} \} \\
& \text{last}_z \circ \text{map value} \circ \text{foldl} (\hat{\oplus}_z) [] \circ \text{filter} (\text{keyEq } k) \circ \\
& \quad \text{concat} \circ \text{map} (\text{foldl} (\hat{\otimes}_z) []) \\
= & \quad \{ \text{corollary of Lemma 3} \} \\
& \text{foldl} (\oplus) z \circ \text{map value} \circ \text{filter} (\text{keyEq } k) \circ \text{concat} \circ \text{map} (\text{foldl} (\hat{\otimes}_z) []) \\
= & \quad \{ \text{naturality, filter } p \circ \text{concat} = \text{concat} \circ \text{map} (\text{filter } p) \} \\
& \text{foldl} (\oplus) z \circ \text{concat} \circ \text{map} (\text{map value} \circ \text{filter} (\text{keyEq } k) \circ \text{foldl} (\hat{\otimes}_z) []) \\
= & \quad \{ \text{Corollary 5} \} \\
& \text{foldl} (\oplus) z \circ \text{map} (\text{foldl} (\otimes) z) \circ \text{filter} (\neg \circ \text{null}) \circ \text{map filter} (\text{keyEq } k) \\
= & \quad \{ \text{definition} \} \\
& \text{aggregateWithKey } k z (\otimes) (\oplus) .
\end{aligned}$$

The "routine naturality laws" in the second step we need are:

- $\text{map} (\text{filter } p) \circ \text{repartition} = \text{repartition} \circ \text{filter } p$,
- $\text{map} (\text{map } f) \circ \text{repartition} = \text{repartition} \circ \text{map } f$,
- $\text{concat} \circ \text{partition} = \text{id}$.

□