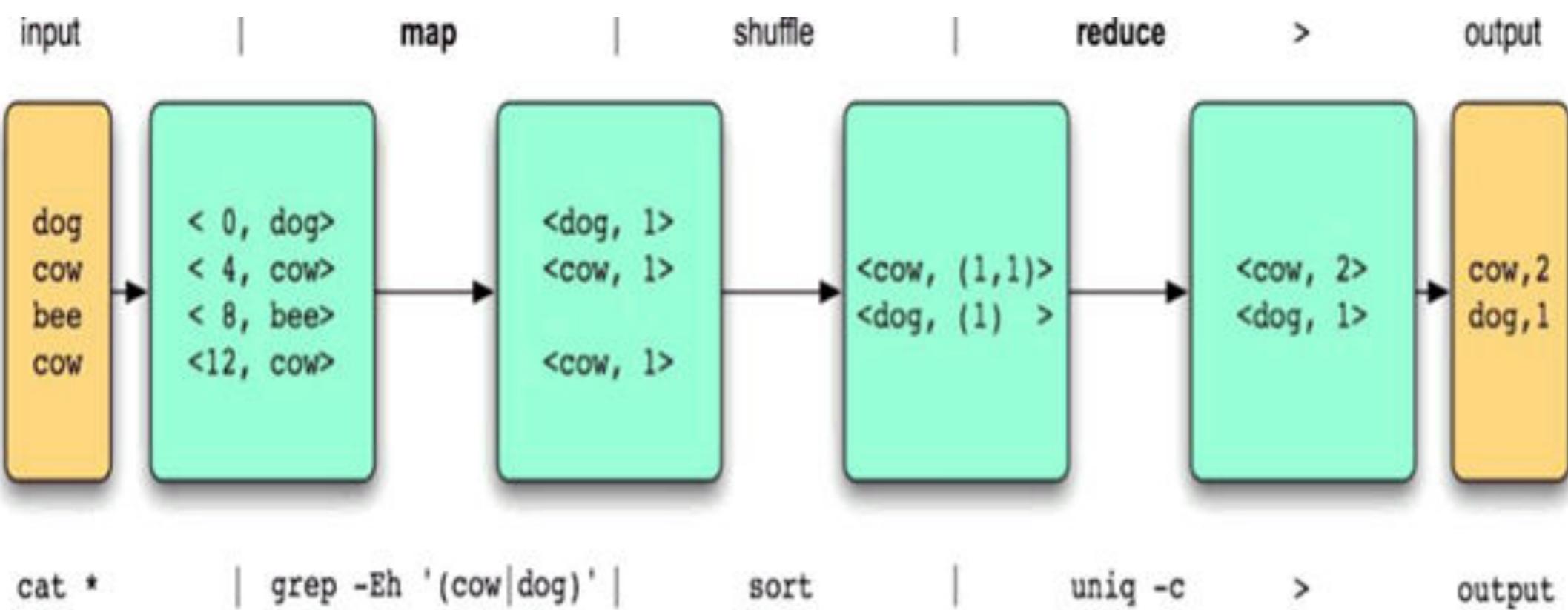


MapReduce Patterns



Intermediate Data

- **Written locally**
- **Transferred from mappers to reducers over network**
- **Issue**
 - Performance bottleneck
- **Solution**
 - Use combiners
 - Use **In-Mapper Combining**

Original Word Count

```
1: class MAPPER
2:   method MAP(docid a, doc d)
3:     for all term t  $\in$  doc d do
4:       EMIT(term t, count 1)

1: class REDUCER
2:   method REDUCE(term t, counts [c1, c2, ...])
3:     sum  $\leftarrow$  0
4:     for all count c  $\in$  counts [c1, c2, ...] do
5:       sum  $\leftarrow$  sum + c
6:     EMIT(term t, count sum)
```

- How many intermediate keys per mapper?
- How can we improve this?
- Is it a “real” improvement?

Stateless In-Mapper Combining

```
1: class MAPPER
2:   method MAP(docid a, doc d)
3:     H  $\leftarrow$  new ASSOCIATIVEARRAY
4:     for all term t  $\in$  doc d do
5:       H{t}  $\leftarrow$  H{t} + 1
6:     for all term t  $\in$  H do
7:       EMIT(term t, count H{t})
```

- Custom local aggregator
- Coding overhead
- Is it a “real” improvement?

Stateful In-Mapper Combining

```
1: class MAPPER
2:   method INITIALIZE
3:      $H \leftarrow \text{new ASSOCIATIVEARRAY}$ 
4:   method MAP(docid  $a$ , doc  $d$ )
5:     for all term  $t \in \text{doc } d$  do
6:        $H\{t\} \leftarrow H\{t\} + 1$ 
7:   method CLOSE
8:     for all term  $t \in H$  do
9:       EMIT(term  $t$ , count  $H\{t\}$ )
```

- Custom local aggregator
- Coding overhead
- Is it a “real” improvement?

In-Mapper Combining Analysis

- **Advantages:**

- Complete local aggregation control (how and when)
- Guaranteed to execute
- Direct efficiency control on intermediate data creation
- Avoid unnecessary objects creation and destruction (before combiners)

- **Disadvantages:**

- Breaks the functional programming background (state)
- Potential ordering-dependent bugs
- Memory scalability bottleneck (solved by memory foot-printing and flushing)

Matrix Generation

- **Common problem:**
 - Given an input of size N, generate an output matrix of size $N \times N$
- **Example: word co-occurrence matrix**
 - Given a document collection, emit the bigram frequencies

“Pairs”

```
1: class MAPPER
2:   method MAP(docid a, doc d)
3:     for all term w ∈ doc d do
4:       for all term u ∈ NEIGHBORS(w) do
5:         EMIT(pair (w, u), count 1)           ▷ Emit count for each co-occurrence

1: class REDUCER
2:   method REDUCE(pair p, counts [c1, c2, ...])
3:     s ← 0
4:     for all count c ∈ counts [c1, c2, ...] do
5:       s ← s + c                         ▷ Sum co-occurrence counts
6:     EMIT(pair p, count s)
```

- We must use custom key type
- Intermediate overhead? Bottlenecks?
- Can we use the reducer as a combiner?
- Keys space?

“Stripes”

```
1: class MAPPER
2:   method MAP(docid  $a$ , doc  $d$ )
3:     for all term  $w \in \text{doc } d$  do
4:        $H \leftarrow \text{new ASSOCIATIVEARRAY}$ 
5:       for all term  $u \in \text{NEIGHBORS}(w)$  do
6:          $H\{u\} \leftarrow H\{u\} + 1$                                  $\triangleright$  Tally words co-occurring with  $w$ 
7:       EMIT(Term  $w$ , Stripe  $H$ )
8:
9: class REDUCER
10:   method REDUCE(term  $w$ , stripes [ $H_1, H_2, H_3, \dots$ ])
11:      $H_f \leftarrow \text{new ASSOCIATIVEARRAY}$ 
12:     for all stripe  $H \in \text{stripes } [H_1, H_2, H_3, \dots]$  do
13:       SUM( $H_f, H$ )                                          $\triangleright$  Element-wise sum
14:     EMIT(term  $w$ , stripe  $H_f$ )
```

- We must use custom key and value types
- Intermediate overhead? Bottlenecks?
- Can we use the reducer as a combiner?
- Keys space?