# Inverted Indexing
## for Text Retrieval

# Three components of the web search problem

- Gathering web content
  - ❖ Web crawling

- Construction of the inverted index
  - ❖ Indexing

- Ranking documents given a query
  - ❖ Retrieval

- First two steps are typically carried out off-line
- The retrieval step needs to be operated in real time

# What is inverted index?

- First, what is index?

## Index

# Example of inverted index

**one fish, two fish**

**red fish, blue fish**

**cat in the hat**

**green eggs and ham**

| | 1 | 2 | 3 | 4 |
|-------|---|---|---|---|
| blue | | 1 | | |
| cat | | | 1 | |
| egg | | | | 1 |
| fish | 1 | 1 | | |
| green | | | | 1 |
| ham | | | | 1 |
| hat | | | 1 | |
| one | 1 | | | |
| red | | 1 | | |
| two | 1 | | | |

| | |
|-------|-------|
| blue | 2 |
| cat | 3 |
| egg | 4 |
| fish | 1 → 2 |
| green | 4 |
| ham | 4 |
| hat | 3 |
| one | 1 |
| red | 2 |
| two | 1 |

# More abstract view of inverted index

terms     postings

| $term_1$ | $d_1$ | $p$ | → | $d_5$ | $p$ | → | $d_6$ | $p$ | → | $d_{11}$ | $p$ | ... |

| $term_2$ | $d_{11}$ | $p$ | → | $d_{23}$ | $p$ | → | $d_{59}$ | $p$ | → | $d_{84}$ | $p$ | ... |

| $term_3$ | $d_1$ | $p$ | → | $d_4$ | $p$ | → | $d_{11}$ | $p$ | → | $d_{19}$ | $p$ | ... |

...   ...

- An inverted index consists of posting lists
- A posting list is comprised of individual postings
  - ❖ Each posting consists of a document id and a payload
    - ▪ Payload example: the occurrence frequency of the term in the corresponding document
  - ❖ Generally, postings are sorted by document id

# Baseline implementation of inverted indexing

```
1: class MAPPER
2:     procedure MAP(docid n, doc d)
3:         H ← new ASSOCIATIVEARRAY
4:         for all term t ∈ doc d do
5:             H{t} ← H{t} + 1
6:         for all term t ∈ H do
7:             EMIT(term t, posting ⟨n, H{t}⟩)
```
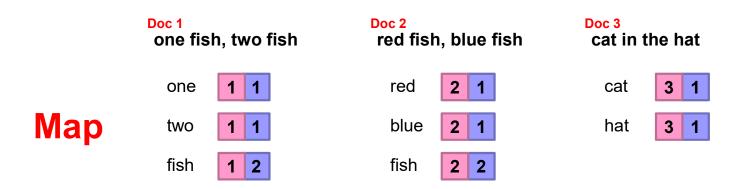
```
1: class REDUCER
2:     procedure REDUCE(term t, postings [⟨n₁, f₁⟩, ⟨n₂, f₂⟩ ...])
3:         P ← new LIST
4:         for all posting ⟨a, f⟩ ∈ postings [⟨n₁, f₁⟩, ⟨n₂, f₂⟩ ...] do
5:             P.ADD(⟨a, f⟩)
6:         P.SORT()
7:         EMIT(term t, postings P)
```

# Illustration of the baseline algorithm

**Doc 1**
**one fish, two fish**

**Doc 2**
**red fish, blue fish**

**Doc 3**
**cat in the hat**

**Map**

| | | |
|---|---|---|
| one | 1 | 1 |
| two | 1 | 1 |
| fish | 1 | 2 |

| | | |
|---|---|---|
| red | 2 | 1 |
| blue | 2 | 1 |
| fish | 2 | 2 |

| | | |
|---|---|---|
| cat | 3 | 1 |
| hat | 3 | 1 |

**Shuffle and Sort:** aggregate values by keys

**Reduce**

| | | | | |
|---|---|---|---|---|
| cat | 3 | 1 | | |
| fish | 1 | 2 | 2 | 2 |
| one | 1 | 1 | | |
| red | 2 | 1 | | |

| | | |
|---|---|---|
| blue | 2 | 1 |
| hat | 3 | 1 |
| two | 1 | 1 |

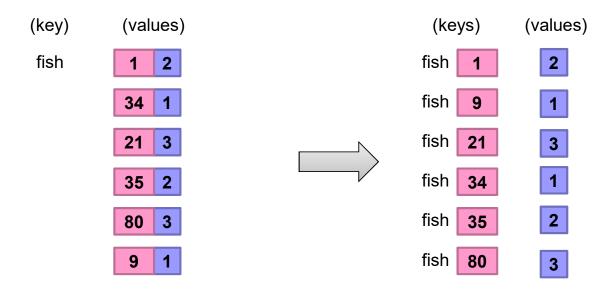# Inverted Indexing: Pseudo-Code

```
1: class MAPPER
2:     procedure MAP(docid n, doc d)
3:         H ← new ASSOCIATIVEARRAY
4:         for all term t ∈ doc d do
5:             H{t} ← H{t} + 1
6:         for all term t ∈ H do
7:             EMIT(term t, posting ⟨n, H{t}⟩)
```

$$H\{t\} \leftarrow H\{t\} + 1$$

$$\text{EMIT}(\text{term } t, \text{posting } \langle n, H\{t\}\rangle)$$

```
1: class REDUCER
2:     procedure REDUCE(term t, postings [⟨a₁, f₁⟩, ⟨a₂, f₂⟩ ...])
3:         P ← new LIST
4:         for all posting ⟨a, f⟩ ∈ postings [⟨a₁, f₁⟩, ⟨a₂, f₂⟩ ...] do
5:             APPEND(P, ⟨a, f⟩)
6:         SORT(P)
7:         EMIT(term t, postings P)
```

$$\text{REDUCE}(\text{term } t, \text{postings } [\langle a_1, f_1\rangle, \langle a_2, f_2\rangle \ldots])$$

$$\text{for all posting } \langle a, f\rangle \in \text{postings } [\langle a_1, f_1\rangle, \langle a_2, f_2\rangle \ldots] \text{ do}$$

$$\text{APPEND}(P, \langle a, f\rangle)$$

What's the problem?

# Scalability issue of the baseline implementation

- Initial implementation: terms as keys, postings as values
  - ❖ Reducers must buffer all postings associated with key (to sort)
  - ❖ What if we run out of memory to buffer postings?

# Another try

| (key) | (values) | | | (keys) | (values) |
|-------|----------|--|--|--------|----------|
| fish | 1 | 2 | | fish 1 | 2 |
| | 34 | 1 | | fish 9 | 1 |
| | 21 | 3 | | fish 21 | 3 |
| | 35 | 2 | | fish 34 | 1 |
| | 80 | 3 | | fish 35 | 2 |
| | 9 | 1 | | fish 80 | 3 |

- Value-to-key conversion

# Revised implementation

```
1: class MAPPER
2:     method MAP(docid n, doc d)
3:         H ← new ASSOCIATIVEARRAY
4:         for all term t ∈ doc d do
5:             H{t} ← H{t} + 1
6:         for all term t ∈ H do
7:             EMIT(tuple ⟨t, n⟩, tf H{t})
```

```
1:  class REDUCER
2:      method INITIALIZE
3:          t_prev ← Ø
4:          P ← new POSTINGSLIST
5:      method REDUCE(tuple ⟨t, n⟩, tf [f])
6:          if t ≠ t_prev ∧ t_prev ≠ Ø then
7:              EMIT(term t, postings P)
8:              P.RESET()
9:          P.ADD(⟨n, f⟩)
10:         t_prev ← t
11:     method CLOSE
12:         EMIT(term t, postings P)
```