# Inverted Indexing
## for Text Retrieval

1

---

## 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

2

---

## What is inverted index?

- First, what is index?



3

---

## Example of inverted index



4

---

## More abstract view of inverted index



- 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

5

---

## 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)
```

6

## 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 |

7

## 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}⟩)
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)
```
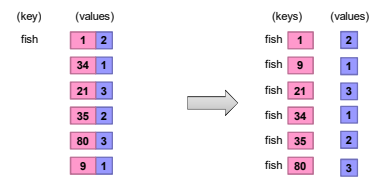
**What's the problem?**

8

## 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?

9

## Another try

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

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

- Value-to-key conversion

10

## 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 ← ∅
4:        P ← new POSTINGSLIST
5:    method REDUCE(tuple ⟨t, n⟩, tf [f])
6:        if t ≠ t_prev ∧ t_prev ≠ ∅ 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)
```

11