# Buckets

Group by Operation

=> Combined locally

hyd,
hyd,
hyd,
:

P1

P1

P2

groupby

P2

540 mb
customers.csv

P3

P4

P5

P200

only 8 will be filled

# Join in Spark

☐→ hyd
☐→ delhi

Customers

Orders

Customers    10    2 ———————→    5↑    50    Orders
              8 ⇒  8 ———————→    8|    8
              2    10 ——————→   10  ⇐ 10
             50    50            50    5

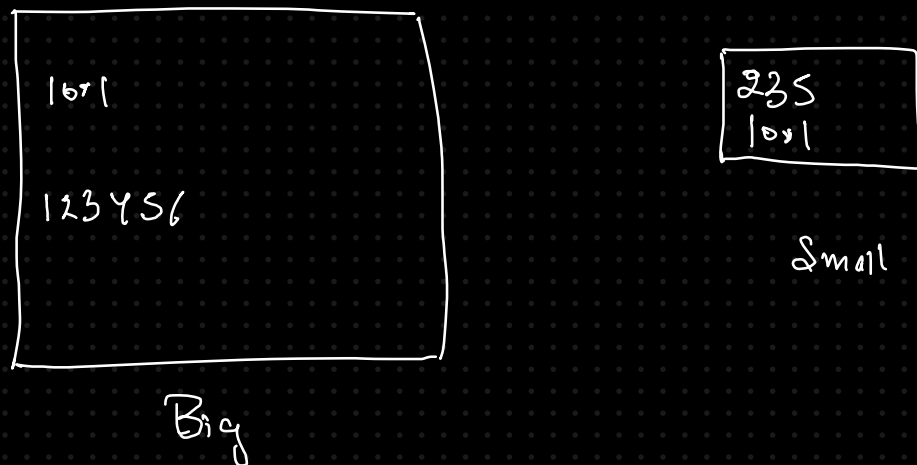# Broadcast Join

```
┌──────────────────────┐                    ┌──────────────┐
│                      │                    │  235         │
│   16ı1               │                    │  10ı1        │
│                      │                    └──────────────┘
│   123456             │                         Small
│                      │
│                      │
└──────────────────────┘
         Big
```

| Feature | Shuffle Sort Merge Join | Broadcast Join |
|---|---|---|
| When Used? | Default for large tables | When one table is small (≤ 50MB) |
| Shuffling? | **Yes** (expensive) | **No** (avoids shuffle) |
| Execution Strategy | **Sort & Merge** | **Hash-based lookup** |
| Performance | **Slower** due to shuffle | **Faster** (O(1) hash lookup) |
| Memory Usage | Low | **Higher** (broadcasted table stored in memory) |
| When to Use? | Large datasets with no small table | Small table + large dataset |
| Can be Forced? | No | **Yes** ( `.hint("broadcast")` ) |

# Skewness in the Data

Skewness in the Data