# Map Reduce

# Typical application

# What if…

INPUT
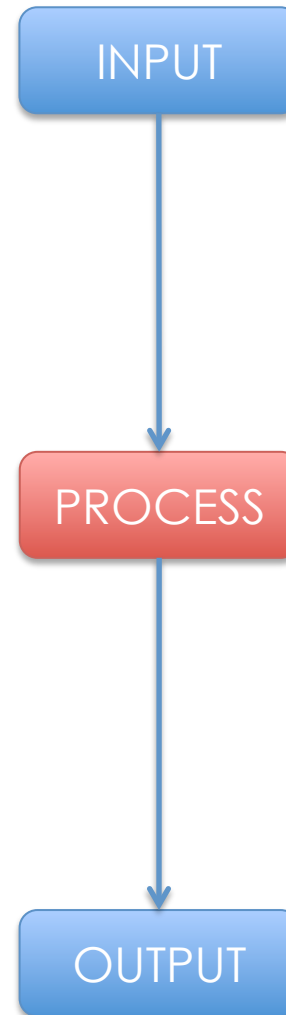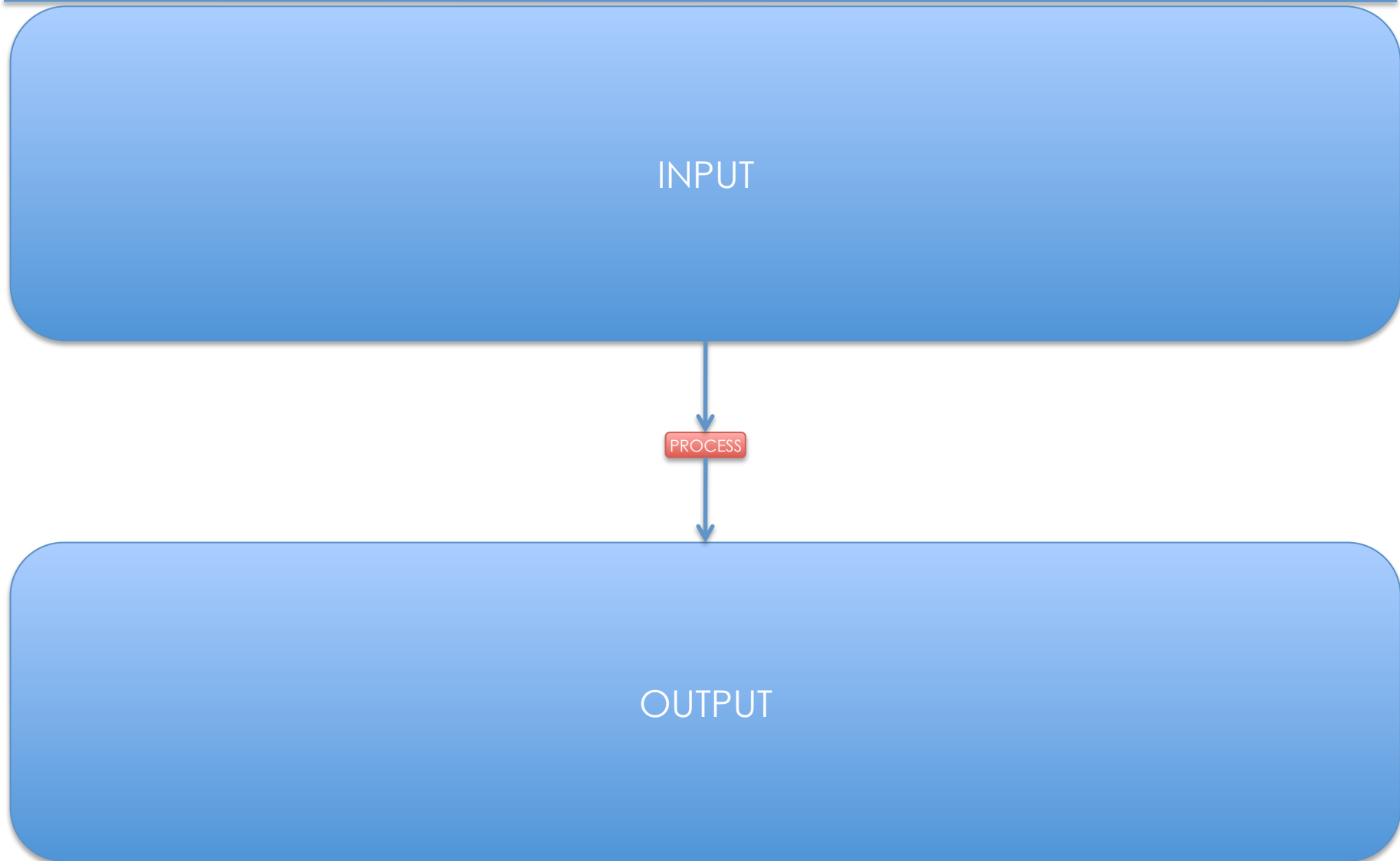
PROCESS

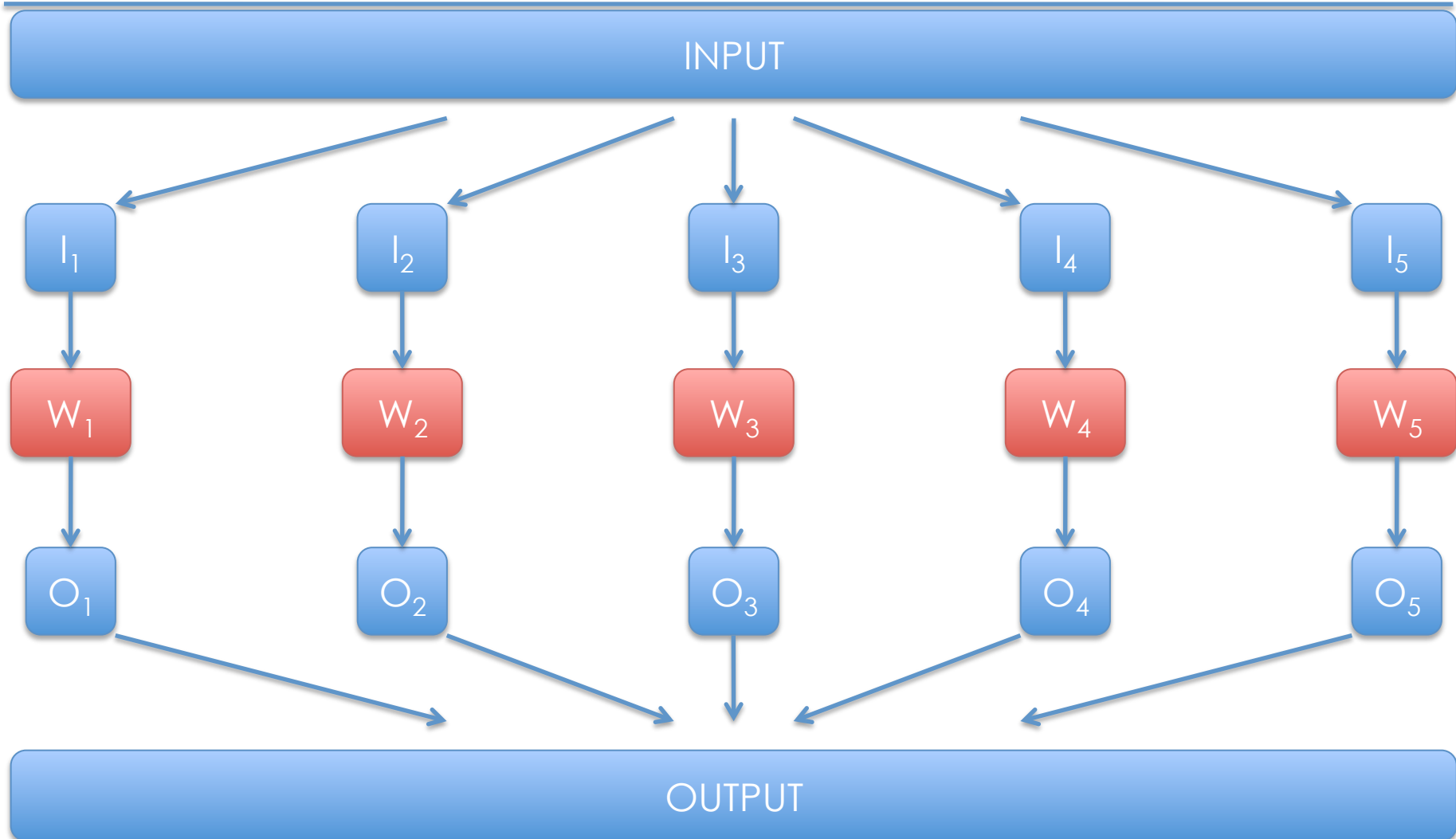OUTPUT

# Divide & Conquer

# Questions

- How do we split the input?

- How do we distribute the input splits?

- How do we collect the output splits?

- How do we aggregate the output?

- How do we coordinate the work?

- What if input splits > num workers?

- What if workers need to share input/output splits?
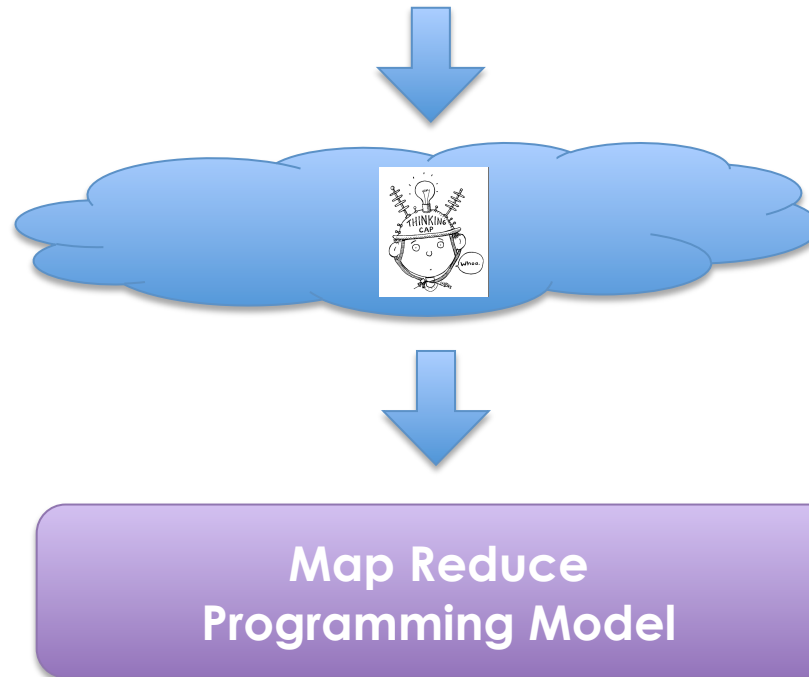
- What if a worker dies?

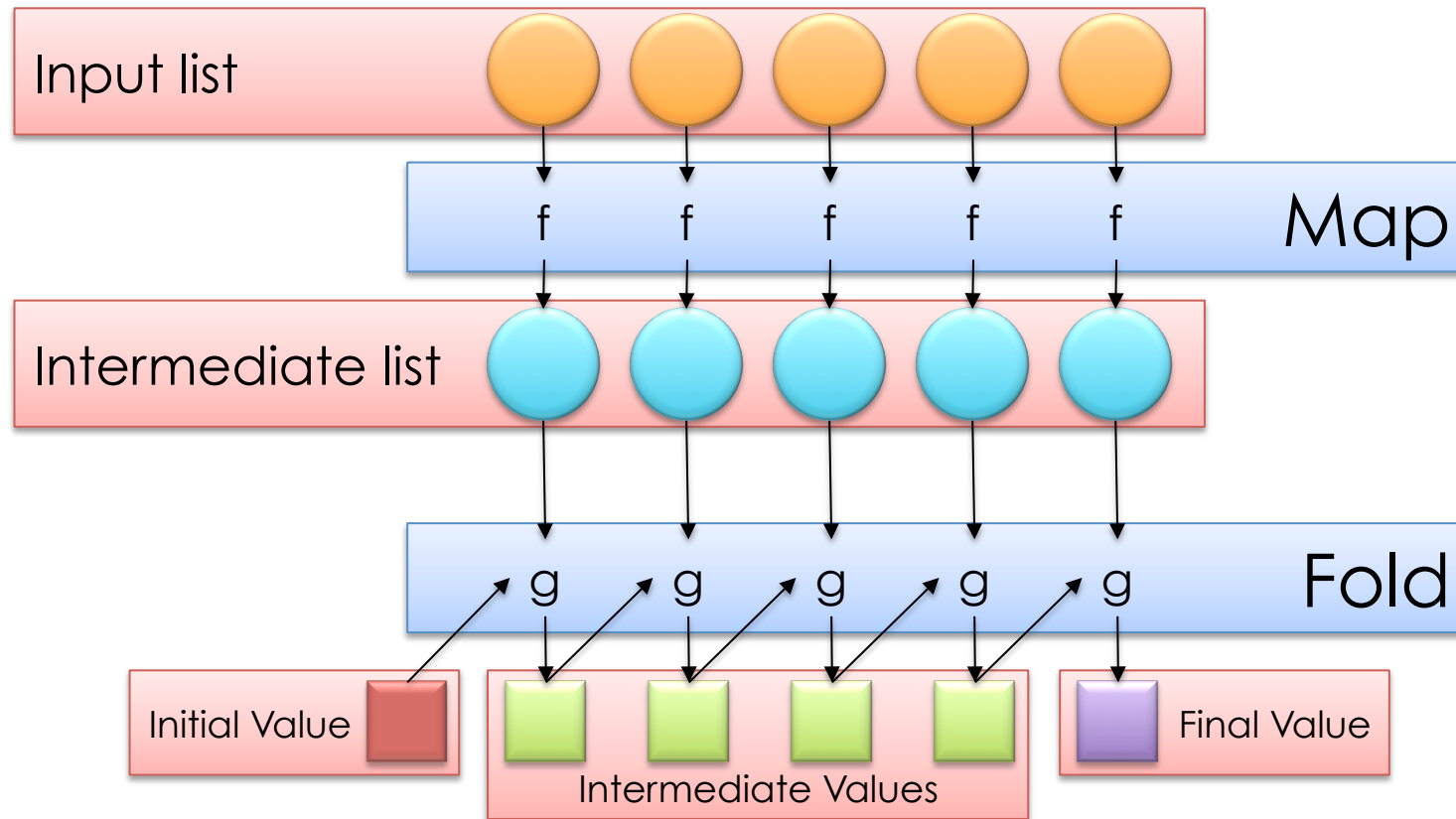- What if we have a new input?

# Design ideas

- Scale "out", not "up"
  - Low end machines

- Move processing to the data
  - Network bandwidth bottleneck

- Process data sequentially, avoid random access
  - Huge data files
  - Write once, read many

- Seamless scalability
  - Strive for the unobtainable

- Right level of abstraction
  - Hide implementation details from applications development

# Typical Large-Data Problem

- Iterate over a large number of records
- Extract something of interest from each
- Shuffle and sort intermediate results
- Aggregate intermediate results
- Generate final output

**Map Reduce Programming Model**

# From functional programming...
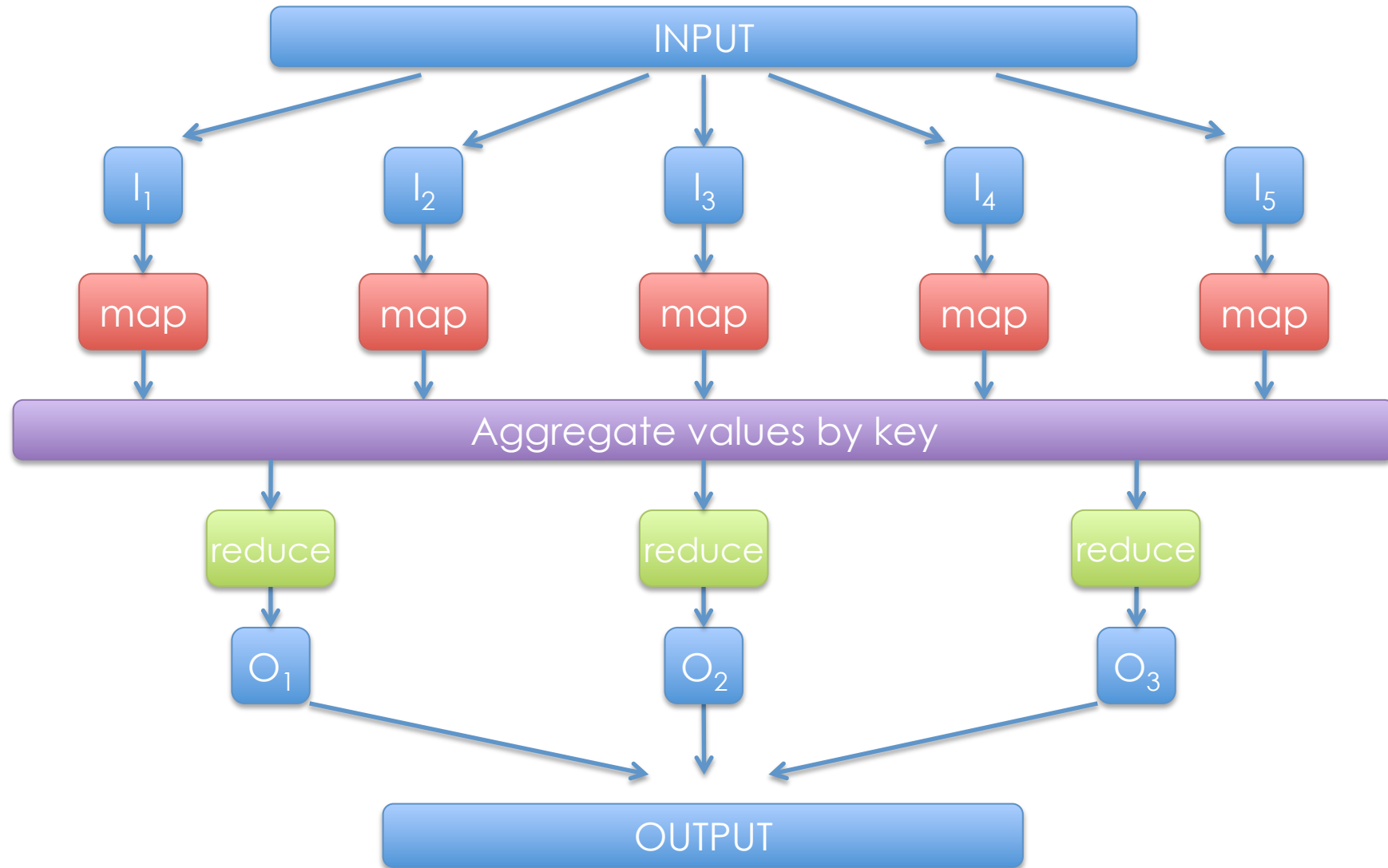
# …To MapReduce

- Programmers specify two functions:

  **map** $(k_1, v_1) \rightarrow [(k_2, v_2)]$

  **reduce** $(k_2, [v_2]) \rightarrow [(k_3, v_3)]$

- All values with the same key are sent to the same reducer

- Input keys and values $(k_1, v_1)$ are drawn from different domain than output keys and values $(k_3, v_3)$

- Intermediate keys $(k_2, v_2)$ and values are from the same domain as the output keys and values $(k_3, v_3)$

- The runtime handles everything else…

# Programming Model (simple)
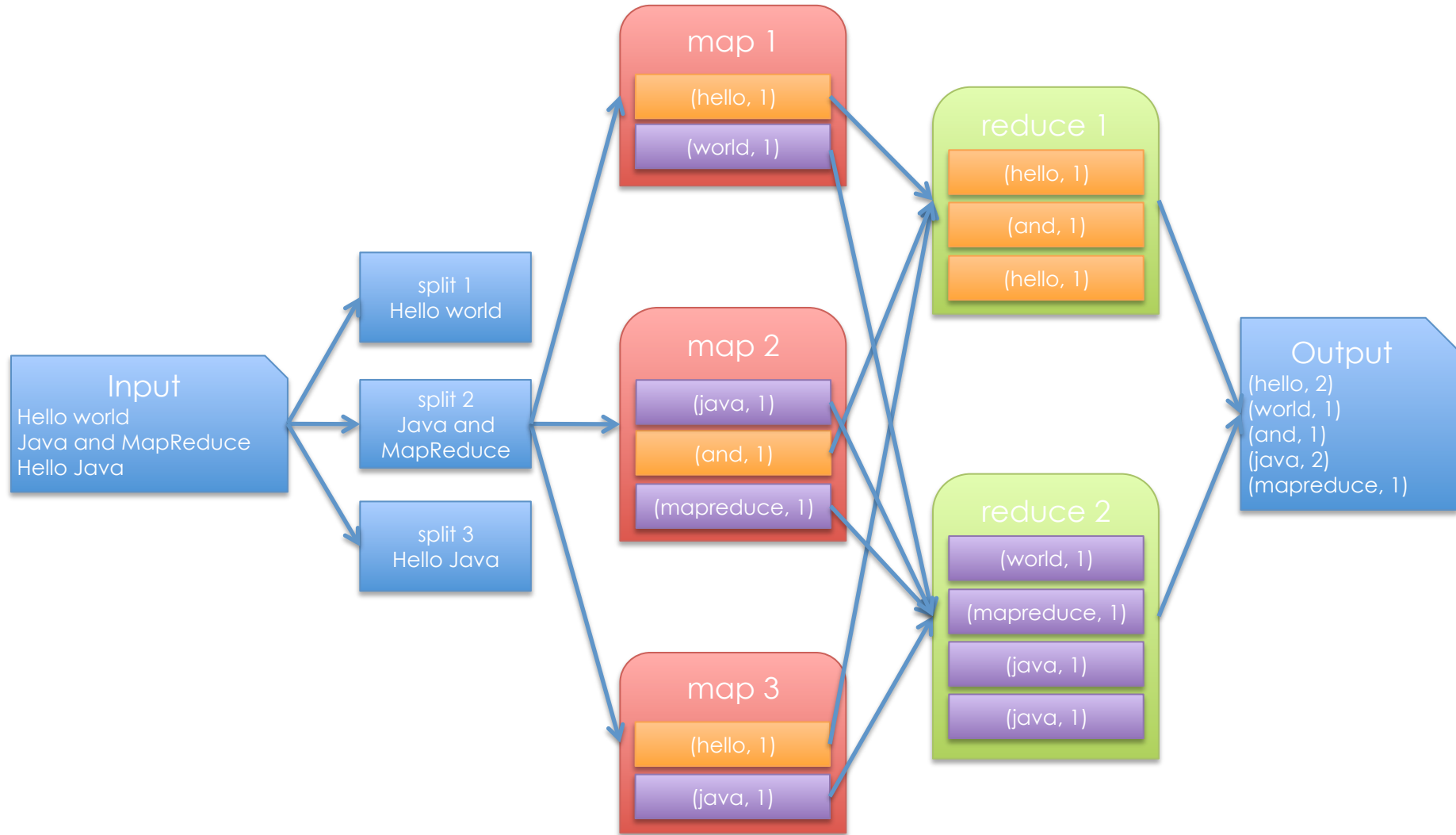
# Example (I)

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

1:  class REDUCER
2:      method REDUCE(term t, counts [c_1, c_2, ...])
3:          sum ← 0
4:          for all count c ∈ counts [c_1, c_2, ...] do
5:              sum ← sum + c
6:          EMIT(term t, count sum)
```
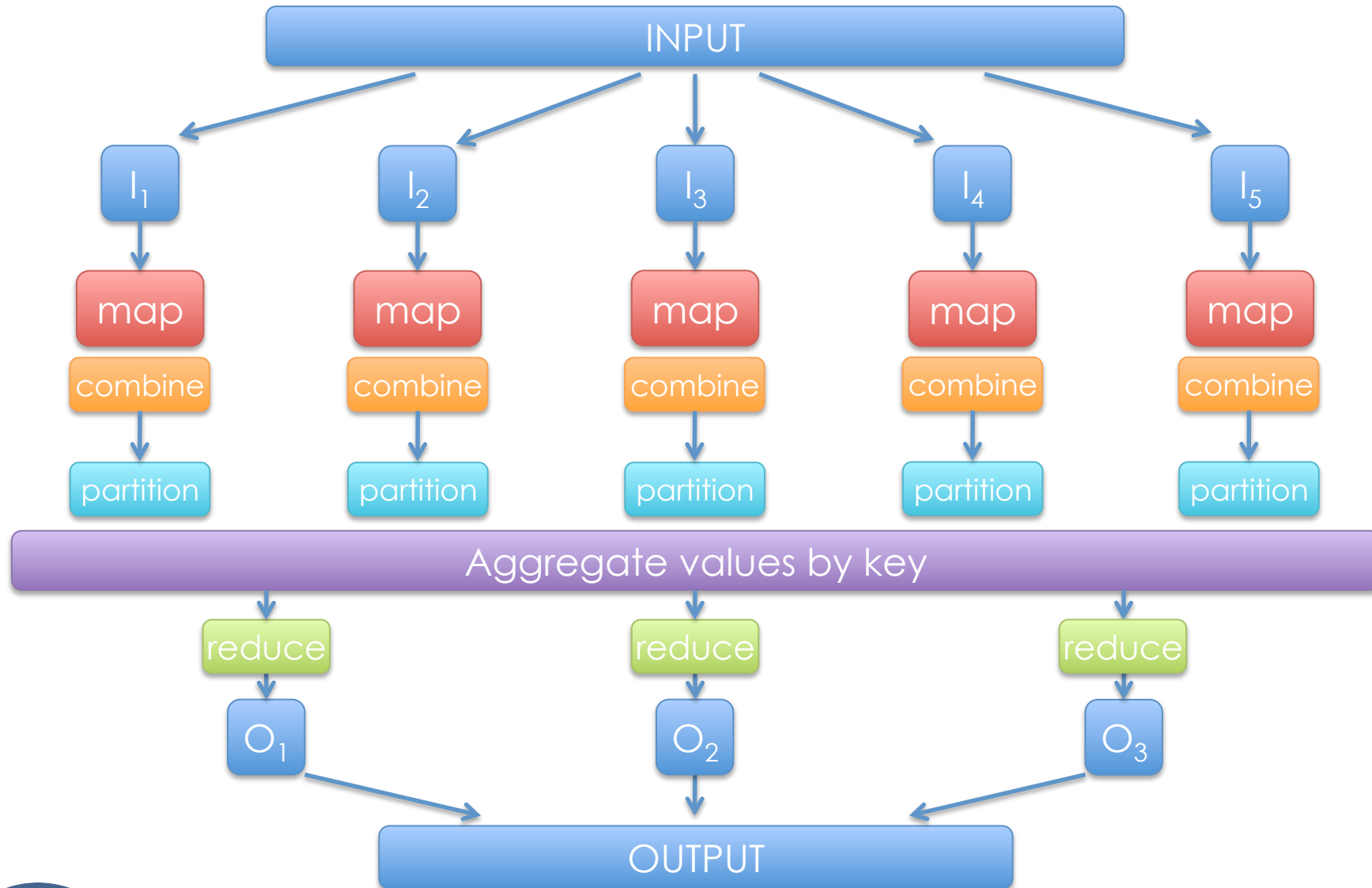
# Example (II)

# Runtime

- ## Handles scheduling
  - Assigns workers to map and reduce tasks
- ## Handles "data distribution"
  - Moves processes to data
- ## Handles synchronization
  - Gathers, sorts, and shuffles intermediate data
- ## Handles errors and faults
  - Detects worker failures and restarts
- ## Everything happens on top of a distributed FS

# Partitioners and combiners

- Programmers specify two functions:
  **map** $(k_1, v_1) \rightarrow [(k_2, v_2)]$
  **reduce** $(k_2, [v_2]) \rightarrow [(k_3, v_3)]$
  – All values with the same key are reduced together

- The execution framework handles everything else…

- Not quite…usually, programmers also specify:

  **partition** $(k_2, \text{number of partitions}) \rightarrow \text{partition for } k_2$
  – Often a simple hash of the key, e.g., hash(k') mod n
  – Divides up key space for parallel reduce operations

  **combine** $(k_2, v_2) \rightarrow [(k_2, v_2)]$
  – Mini-reducers that run in memory after the map phase
  – Used as an optimization to reduce network traffic

# Programming Model (complete)

# MapReduce can refer to…

- The programming model
- The execution framework (aka "runtime")
- The specific implementation

# MapReduce Implementations

- Google has a proprietary implementation in C++
  - Bindings in Java, Python

- Hadoop is an open-source implementation in Java
  - Development led by Yahoo, used in production
  - Now an Apache project
  - Rapidly expanding software ecosystem

- Lots of custom research implementations
  - For GPUs, cell processors, etc.