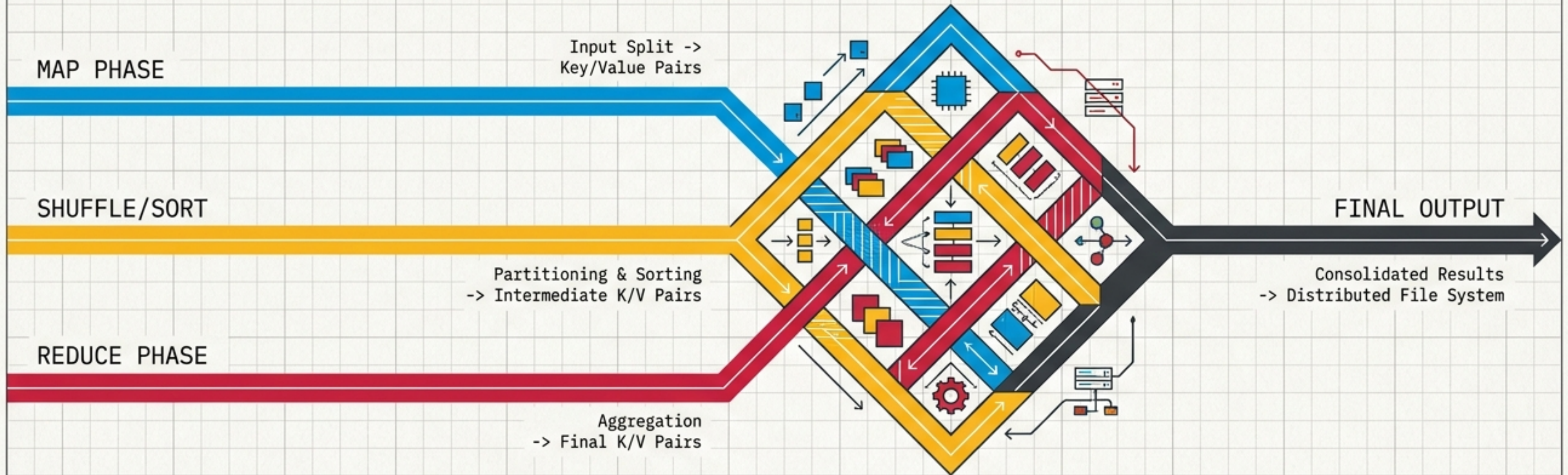
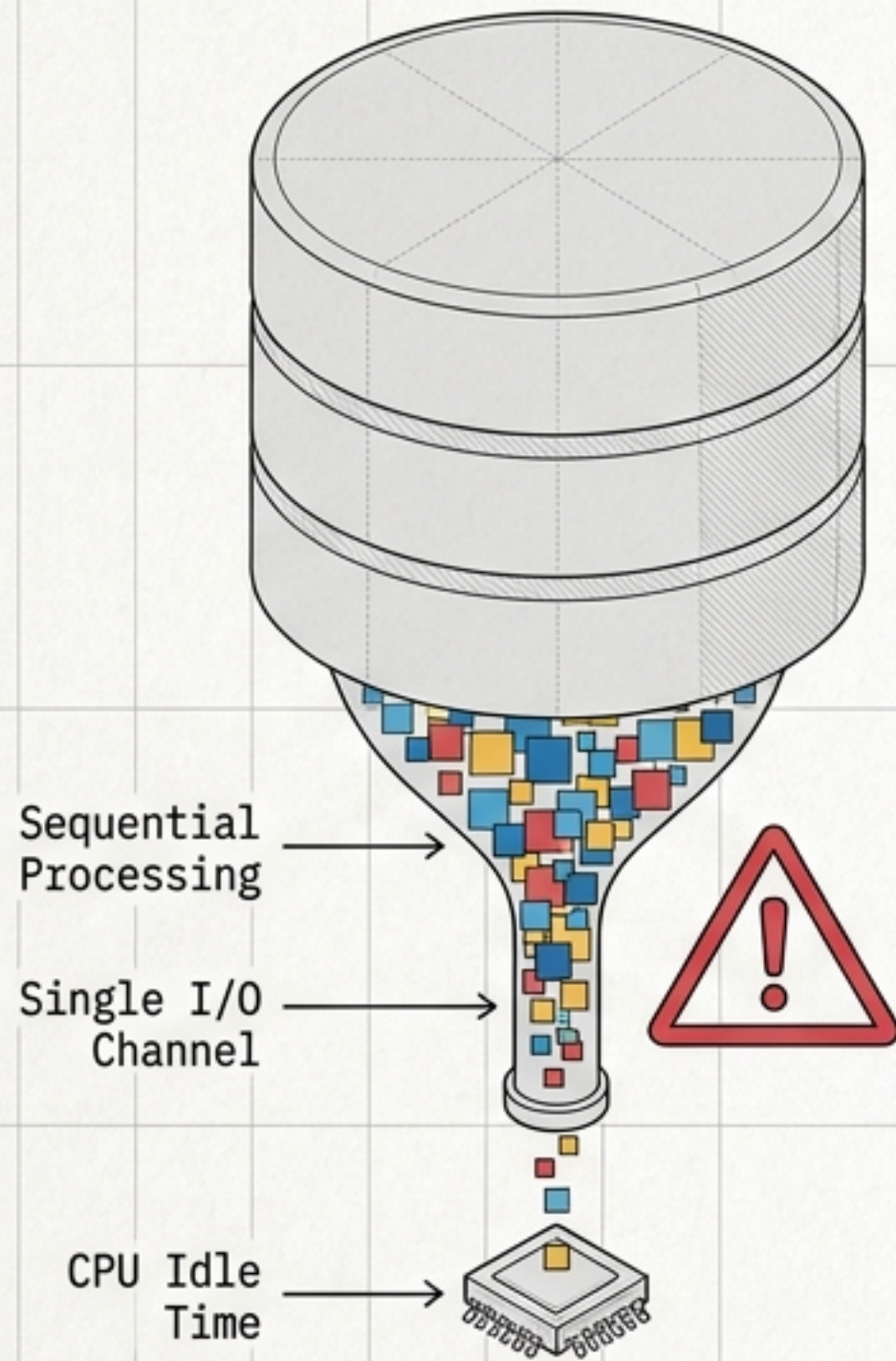


MapReduce: The Anatomy of Distributed Computation

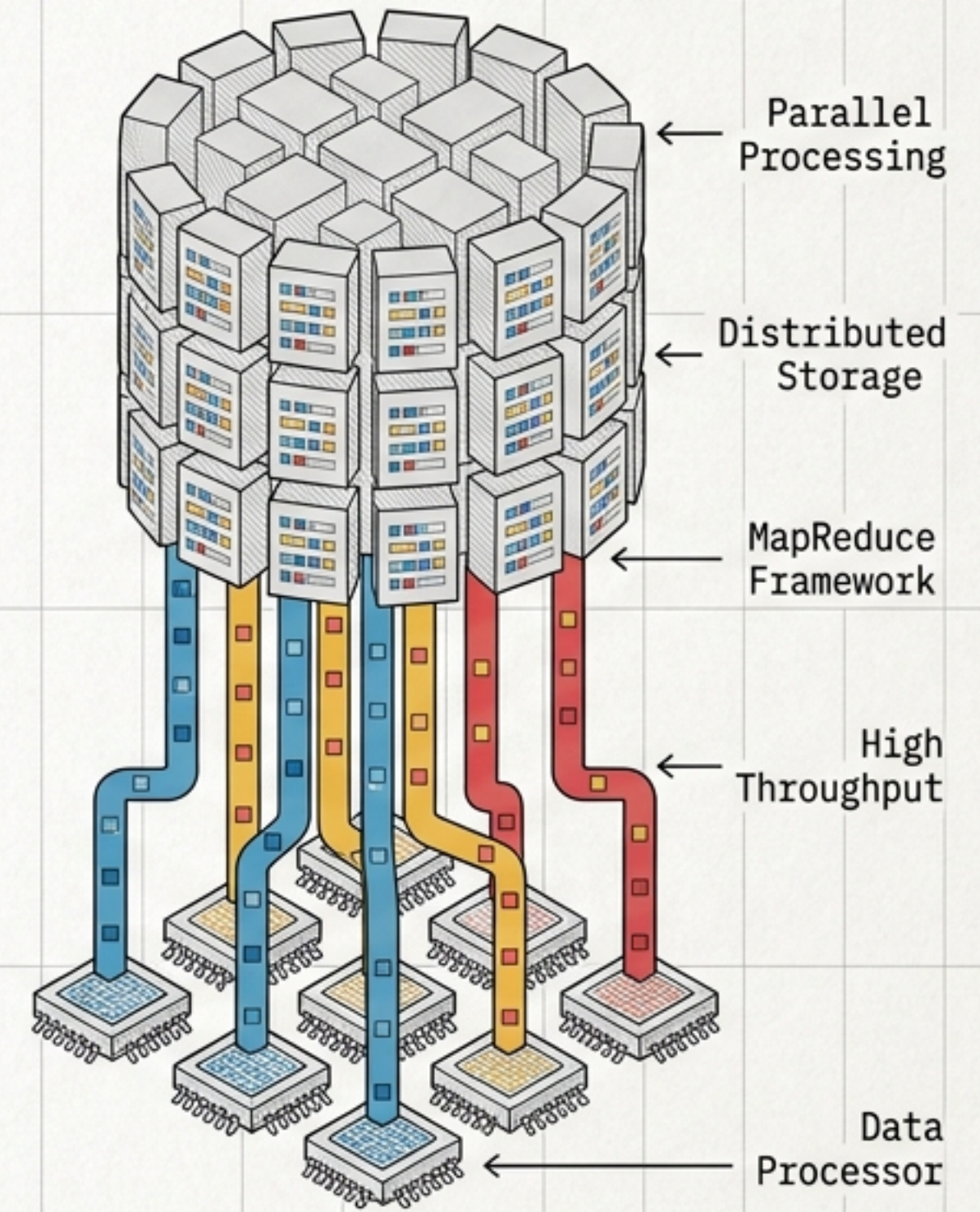
From Abstract Paradigm to Concrete Execution



The Sequential Bottleneck



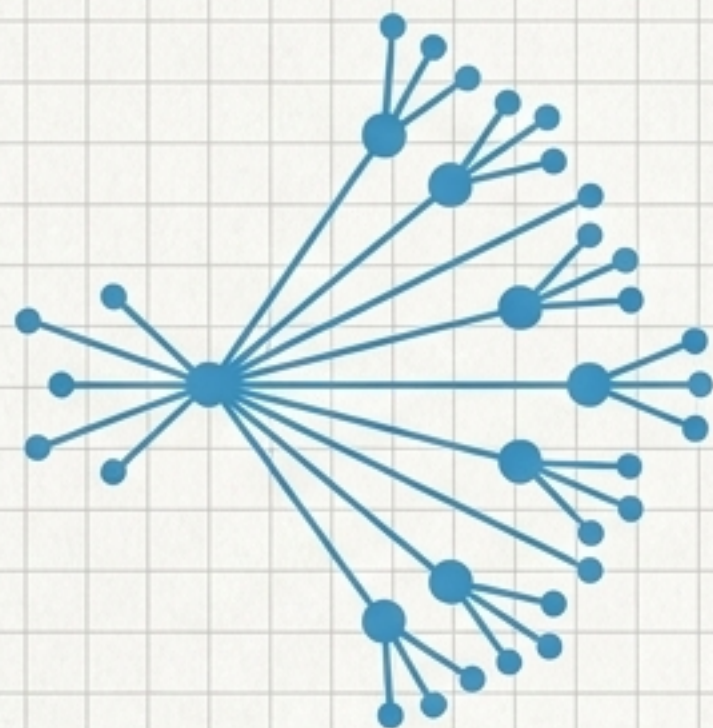
The Distributed Reality



The Big Data scale problem is not a lack of CPU power; it is the physical limitation of moving massive datasets to a single processor.

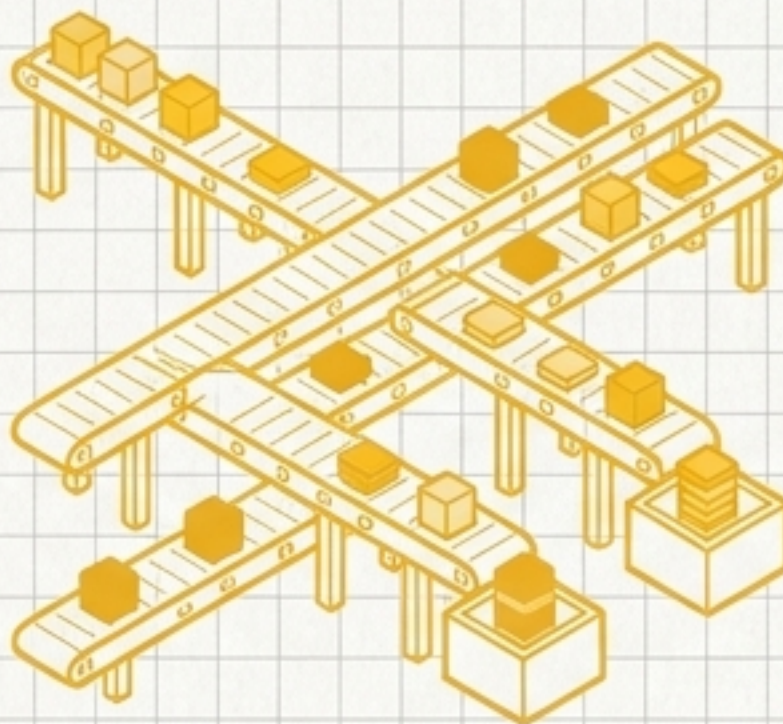
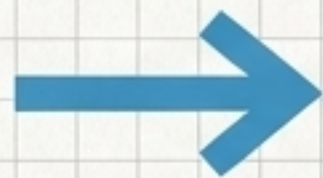
The Paradigm Shift: Traditional vs. MapReduce

Dimension	Traditional Approach	The MapReduce Paradigm
Data Movement	Move data to computation	Move computation to the data
Hardware	Specialized supercomputers	Clusters of commodity hardware
Developer Burden	Developer handles parallelization logic	Framework automates parallelization
Fault Tolerance	Manual checkpointing & recovery	Automatic task reassignment



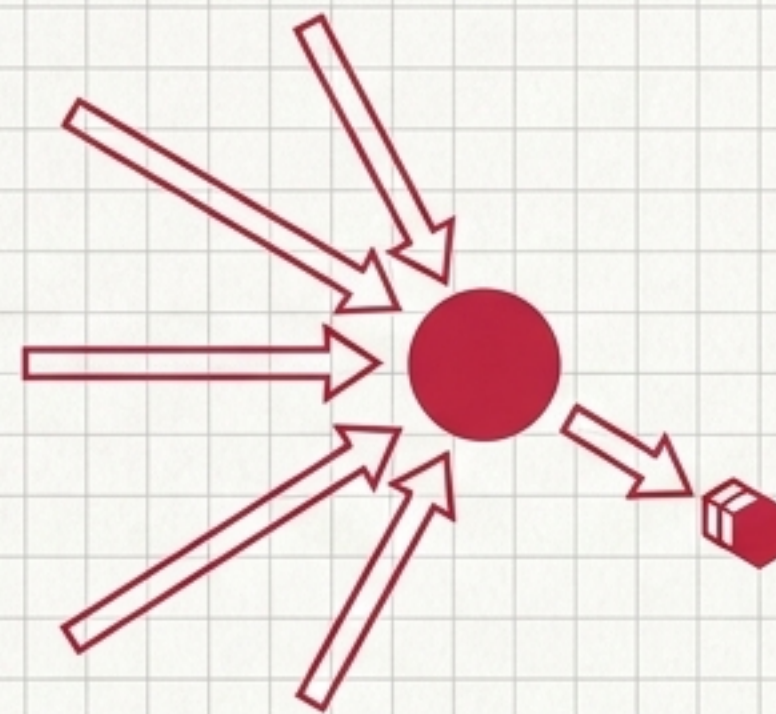
Phase 1: The Map Phase

Divide and Extract.
Filtering and sorting raw input data into structured pairs.



Phase 2: The Shuffle & Sort Phase

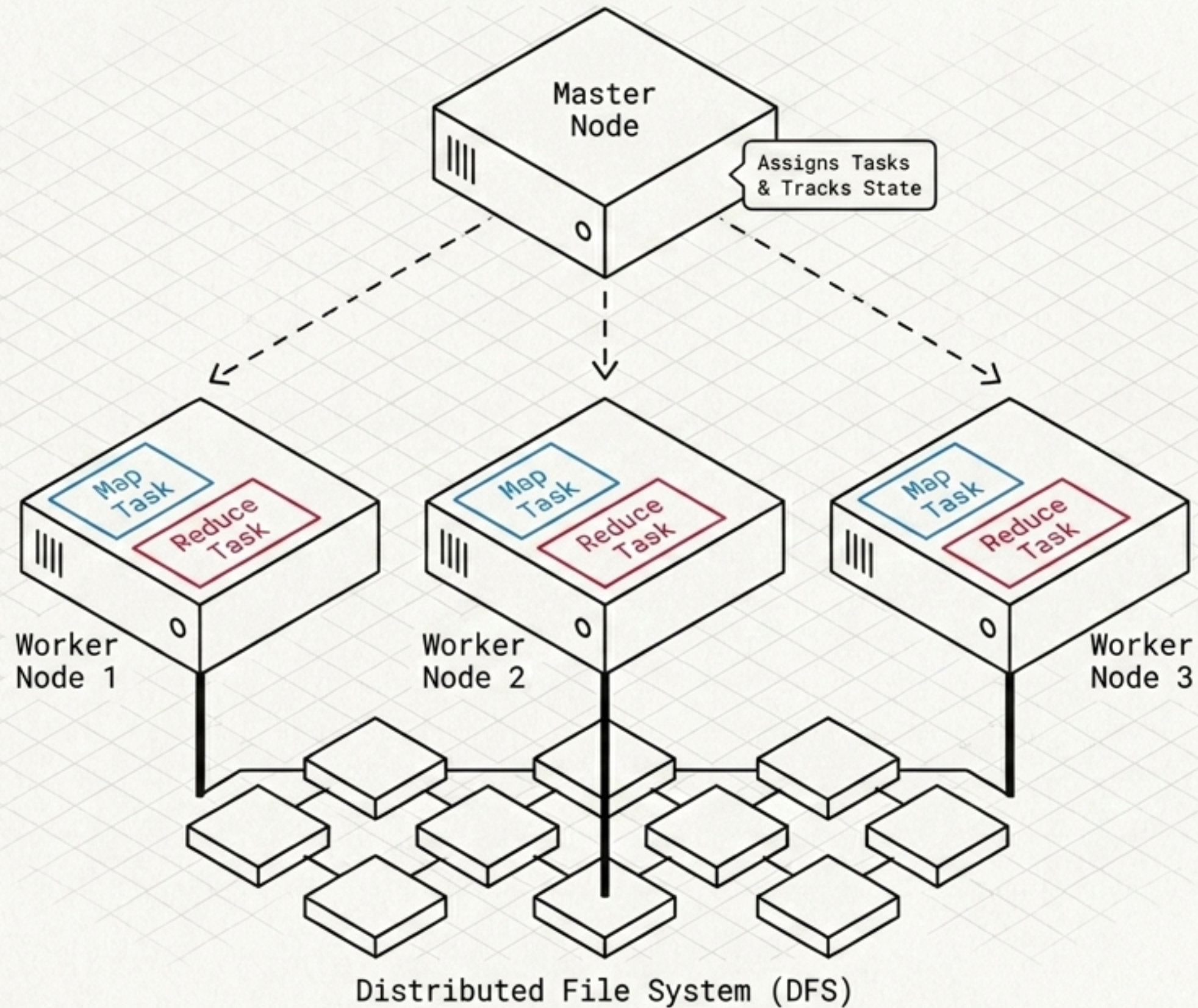
Organize and Route.
Grouping all identical keys together across the network.

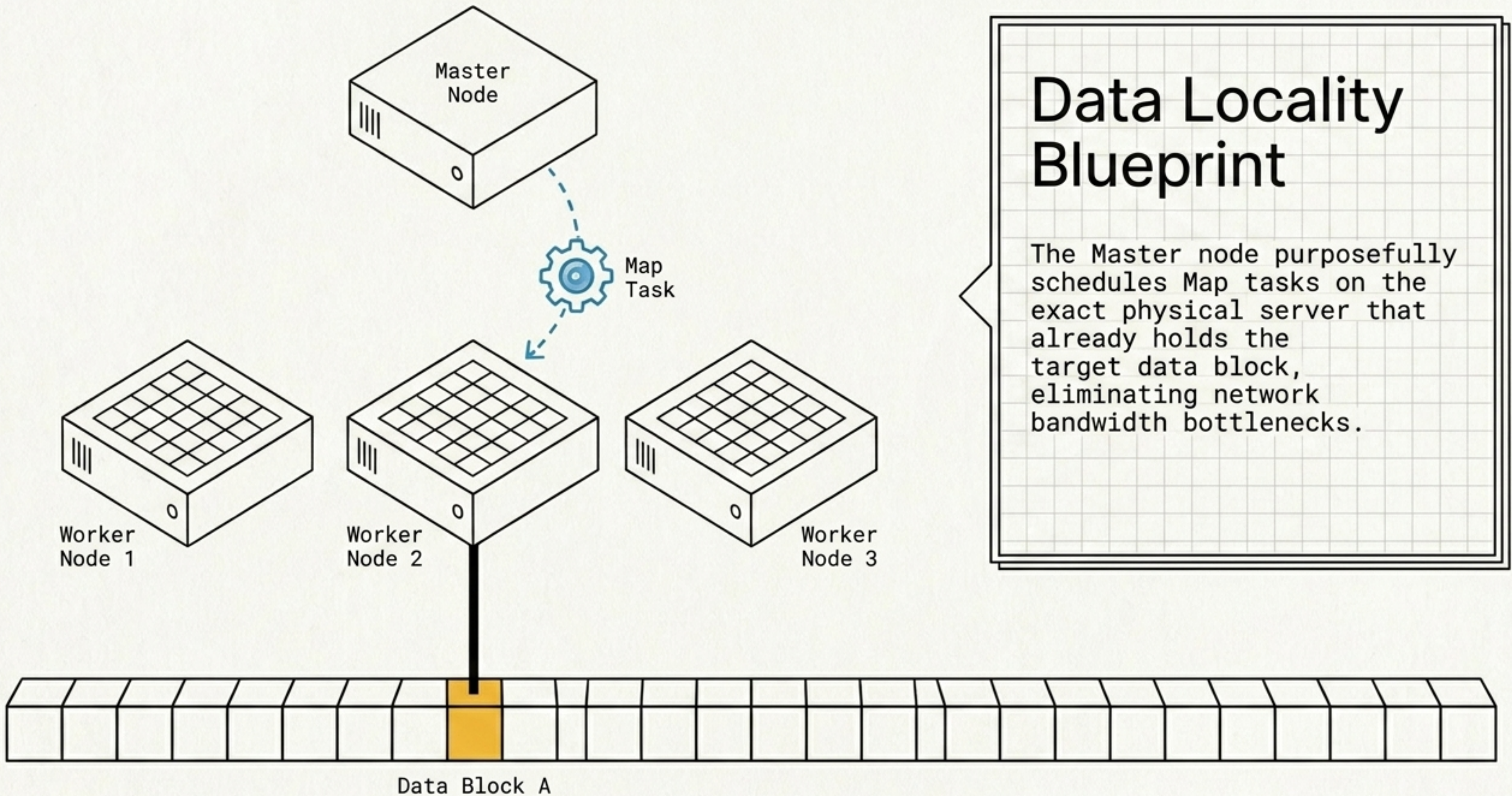


Phase 3: The Reduce Phase

Consolidate and Aggregate.
Summarizing the grouped data into the final output.

MapReduce Cluster Architecture



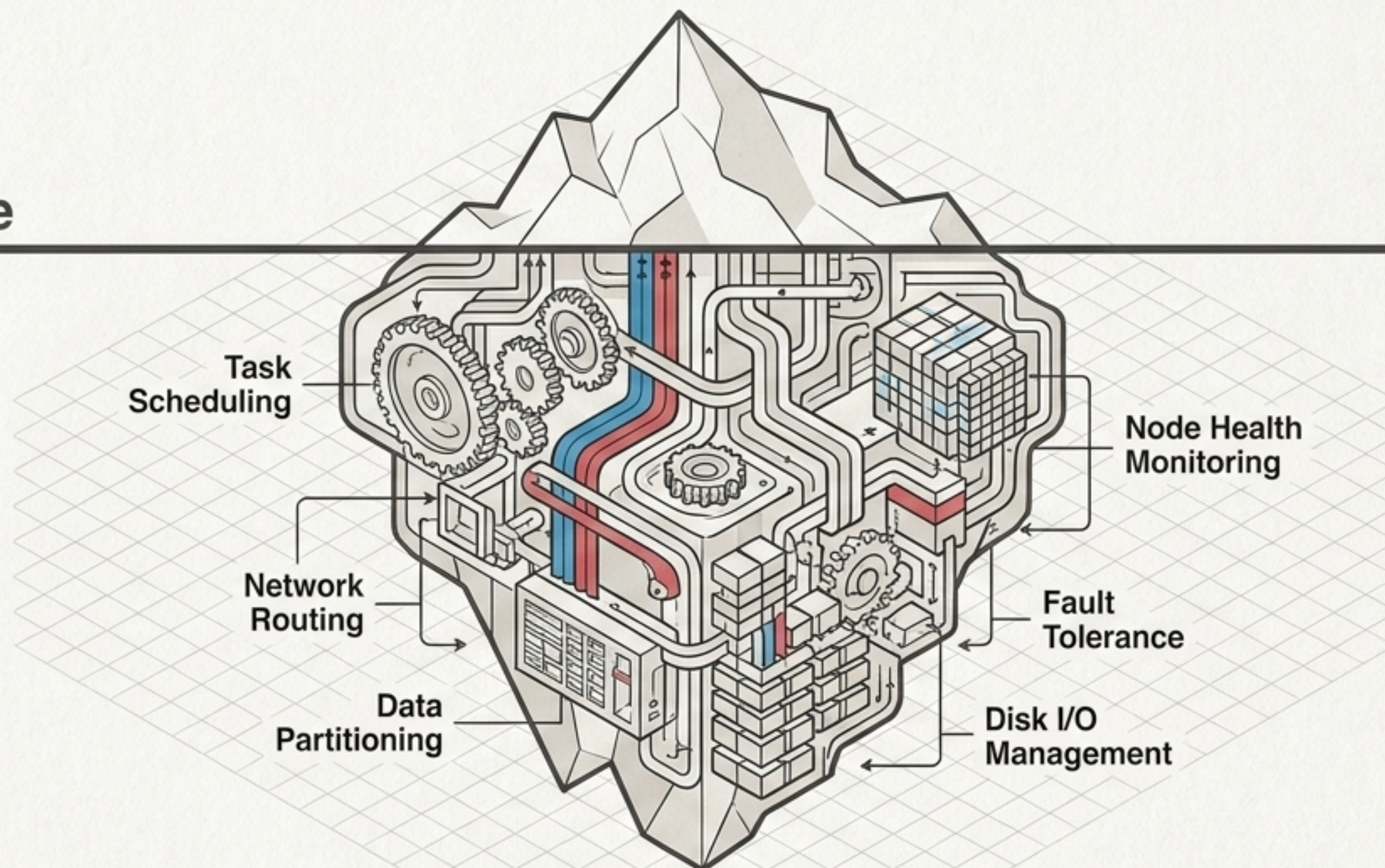


`map(key, value)`

`reduce(key, values)`

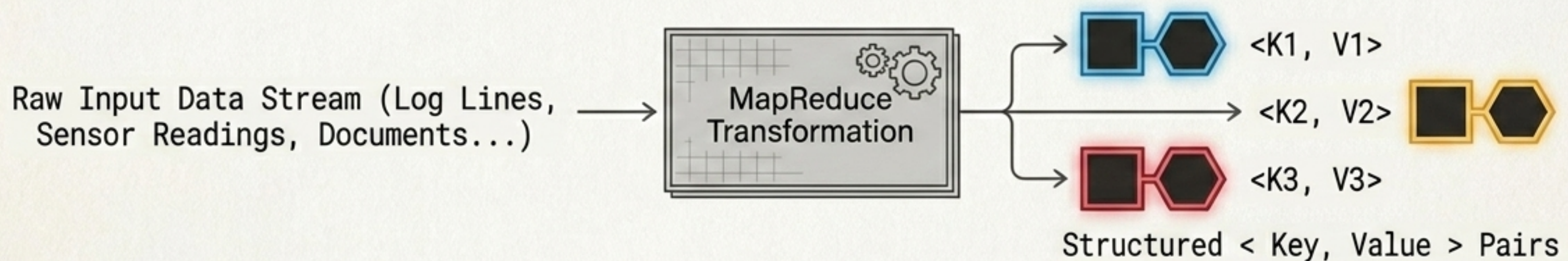
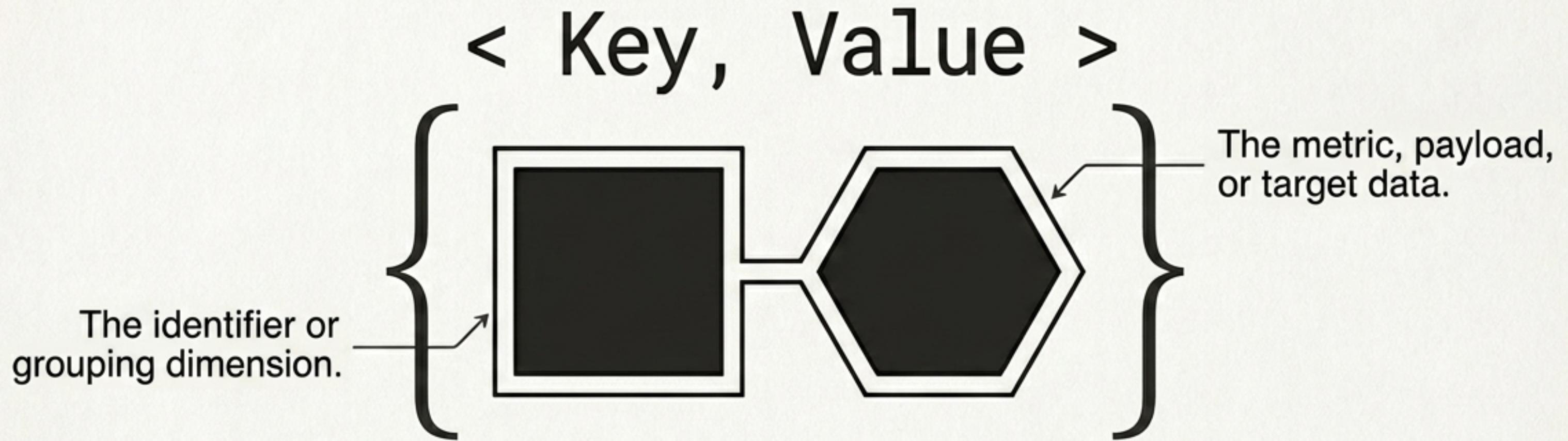
The programmer writes two functions.

The API Surface



“The brilliance of MapReduce is not the computation itself, but the immense complexity it hides behind a simple interface.”

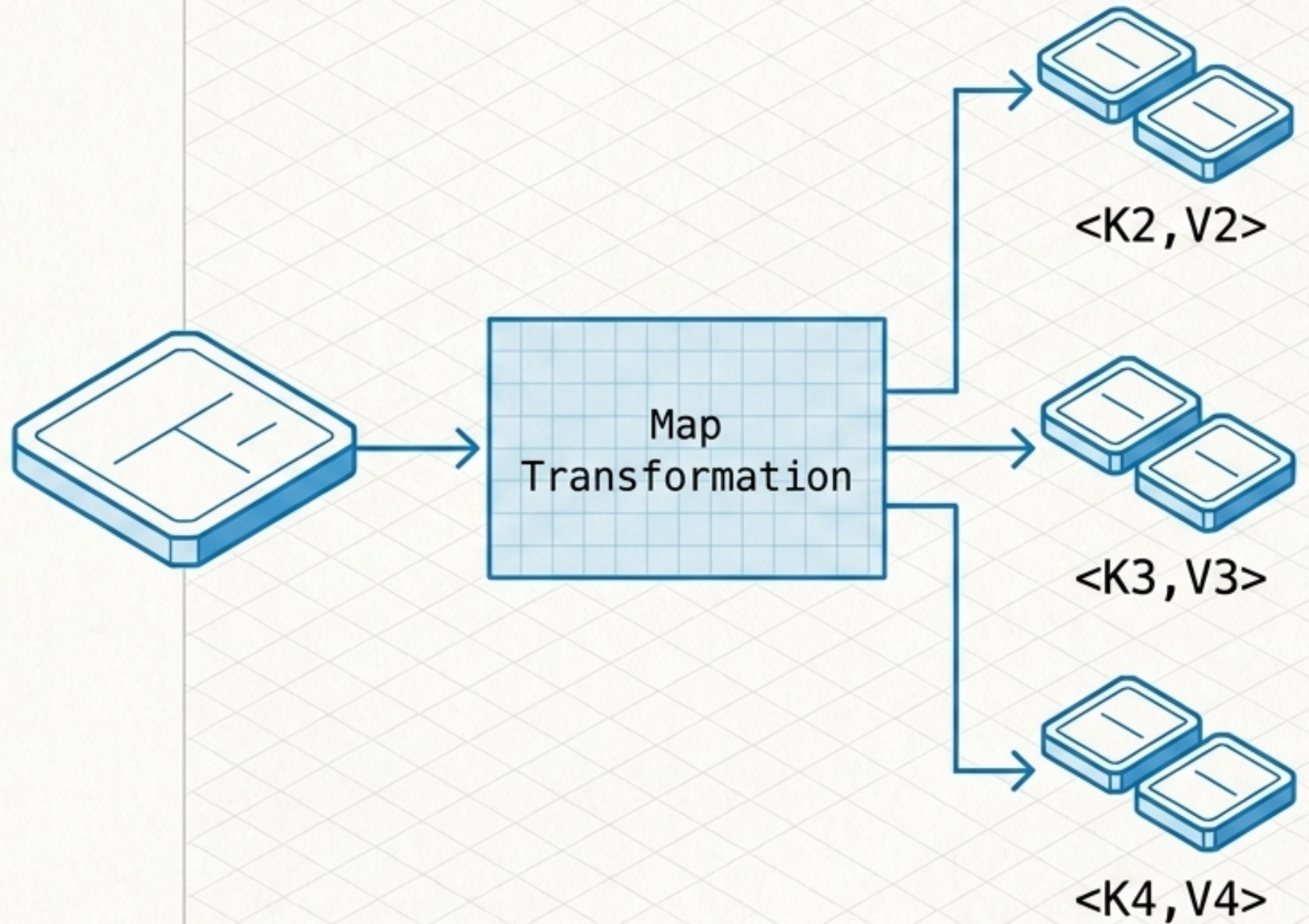
The Universal Currency: Every input and output in a MapReduce pipeline is strictly formatted as a Key-Value pair.



Anatomy of the Map Function

The Mapper takes a single input pair and translates it into zero, one, or multiple intermediate pairs.

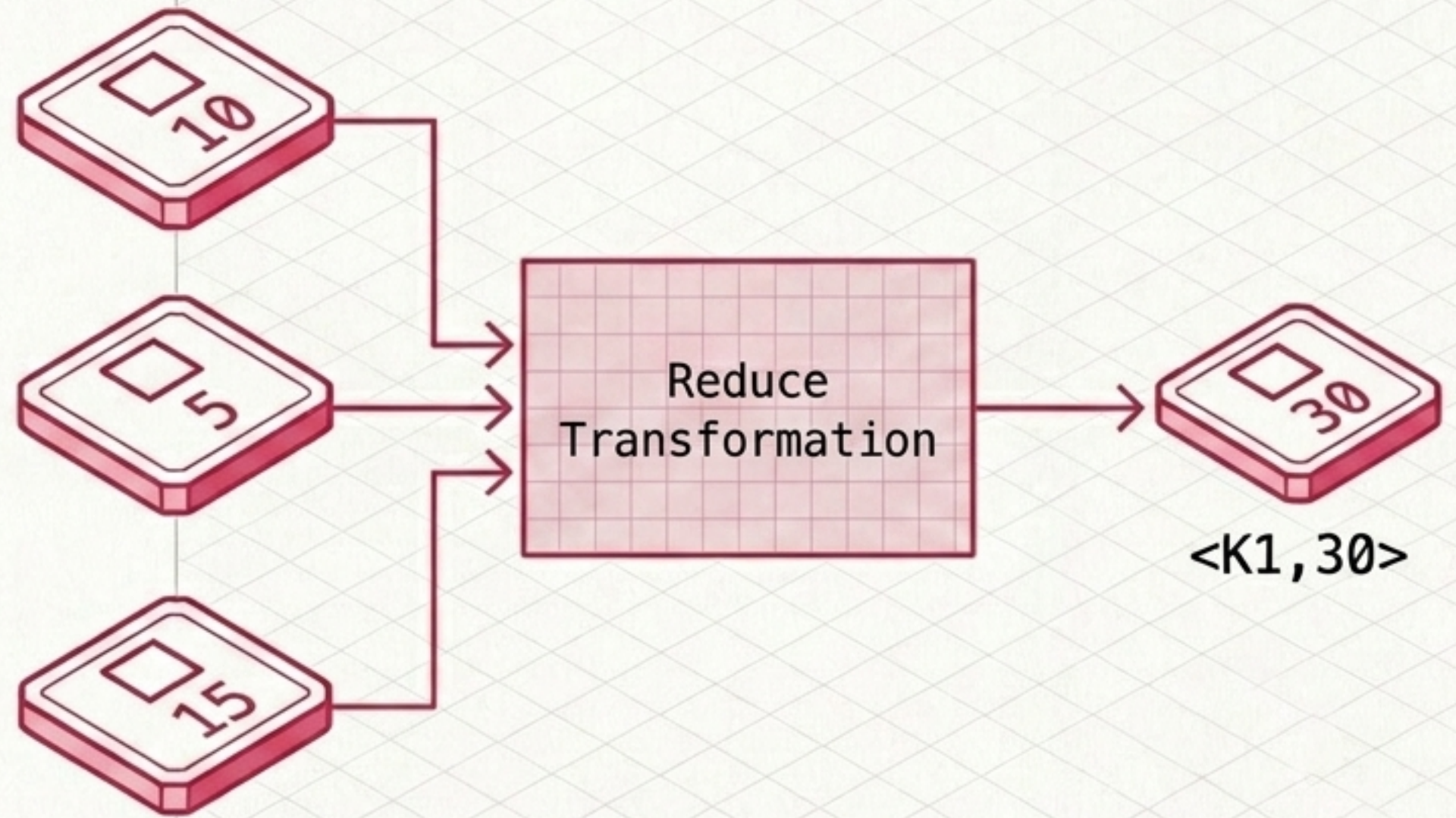
Expansion / Extraction.



Anatomy of the Reduce Function

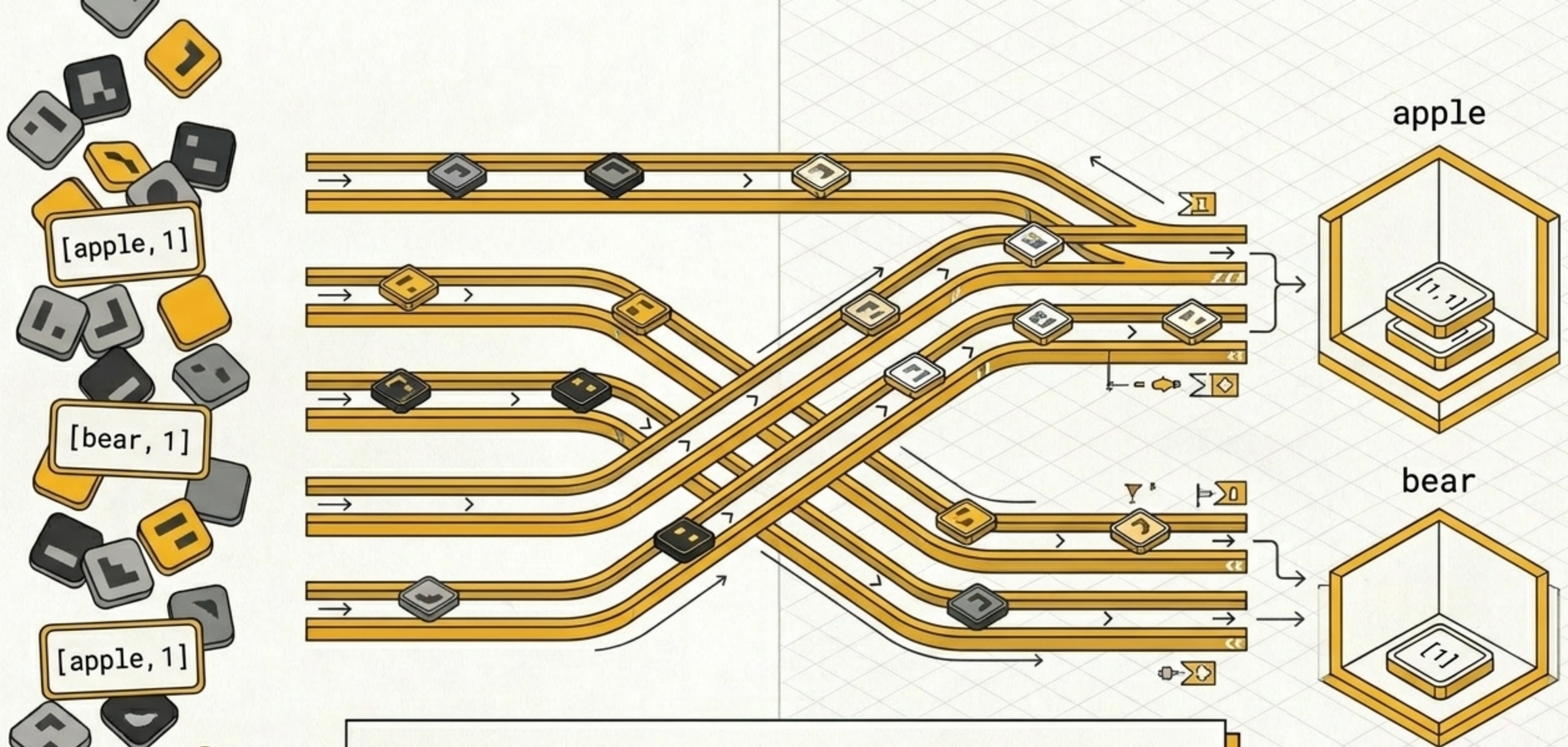
The Reducer accepts a single Key and an iterator of all Values associated with that Key, summarizing them into a final output.

Consolidation / Aggregation.



Anatomy of the Interface

	MAP	REDUCE
Primary Goal	Transform and emit intermediate data.	Aggregate and finalize output data.
Cardinality	1-to-N (One input pair yields many intermediate pairs).	N-to-1 (Many intermediate values yield one final value).
Input Format	<code><k1, v1></code>	<code><k2, list(v2)></code>
Output Format	<code>list(<k2, v2>)</code>	<code>list(<k3, v3>)</code>



The Shuffle & Sort: The framework guarantees that all values associated with a specific key arrive at the same Reducer.