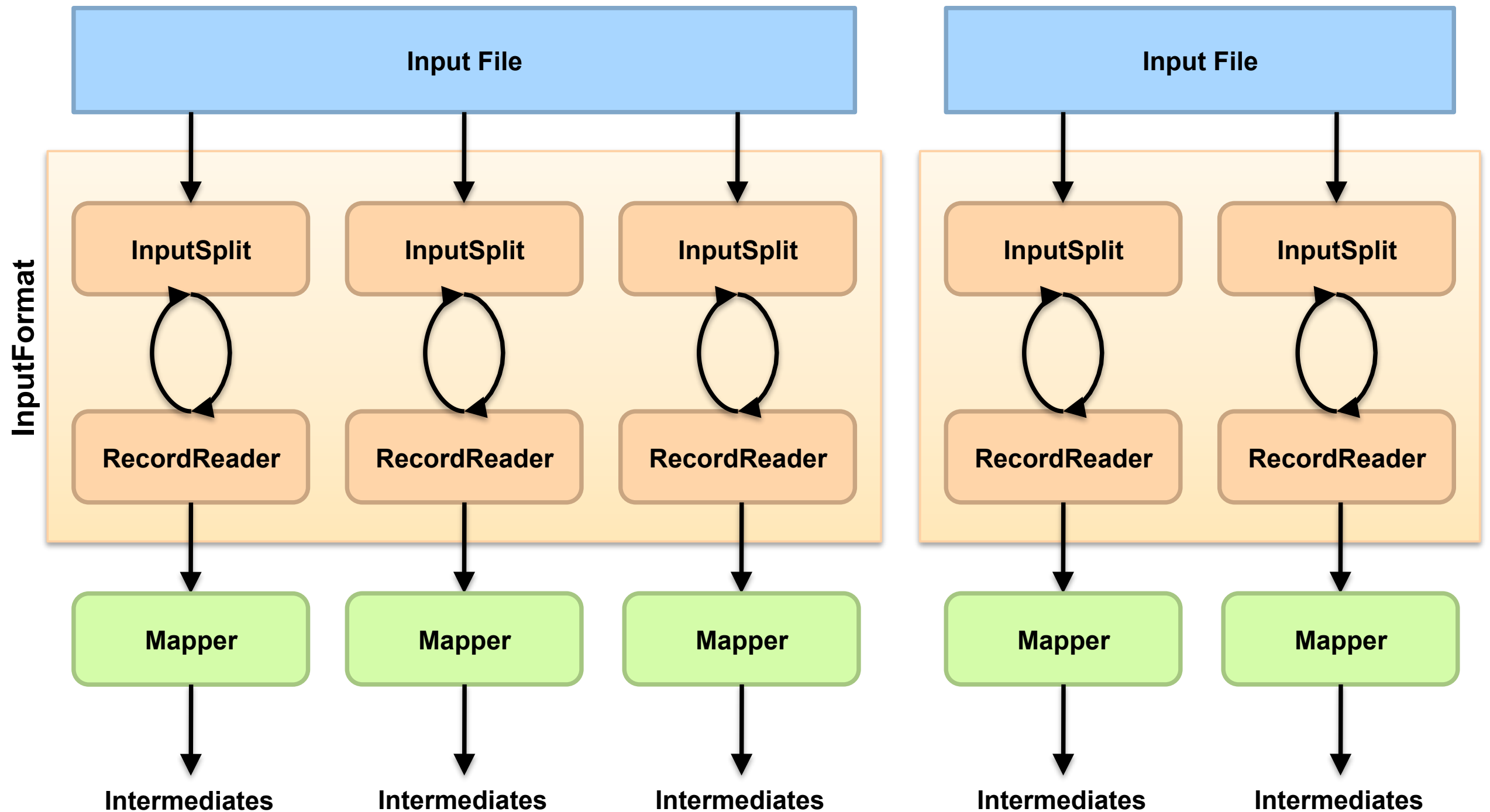


Hadoop Dataflow (I)



- Data sets are specified by **InputFormats**
 - Defines input data (e.g., a directory)
 - Identifies partitions of the data that form an **InputSplit**, each of which will be assigned to a mapper
 - Provide the **RecordReader** implementation to extract (k, v) records from the input source
- Base class implementation is **FileInputFormat**
 - Will read all files out of a specified directory and send them to the mappers
 - **TextInputFormat** – Treats each ‘\n’-terminated line of a file as a value
 - **KeyValueTextInputFormat** – Maps ‘\n’- terminated text lines of “k SEP v”
 - **SequenceFileInputFormat** – Binary file of (k, v) pairs with some add’l metadata
 - **SequenceFileAsTextInputFormat** – Same, but maps (k.toString(), v.toString())



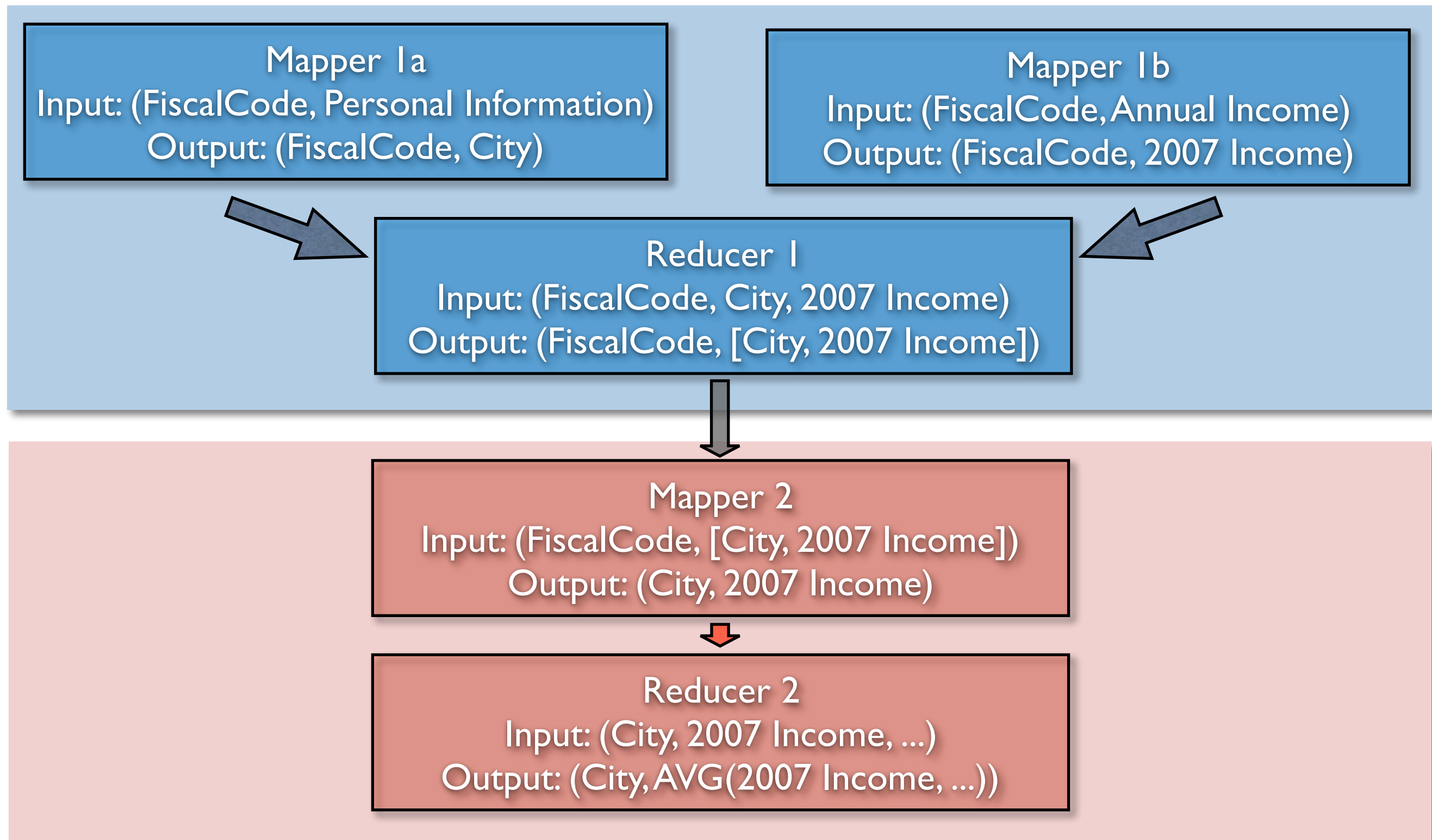
Computing average income (I)

- FiscalTable 1: (FiscalCode, {Personal Information})
 - ABCDEF123: (Alice Rossi; Pisa, Toscana)
 - ABCDEF124: (Bebo Verdi; Firenze, Toscana)
 - ABCDEF125: (Carlo Bianchi; Genova, Liguria)
- FiscalTable 2: (FiscalCode, {year, income})
 - ABCDEF123: (2007, € 70,000), (2006, € 65,000), (2005, € 60,000),...
 - ABCDEF124: (2007, € 72,000), (2006, € 70,000), (2005, € 60,000),...
 - ABCDEF125: (2007, € 80,000), (2006, € 85,000), (2005, € 75,000),...
- **Task:** Compute average income in each city in 2007
- **Note:** Both inputs sorted by FiscalCode

taken from: <http://research.google.com/>



Computing average income (II)

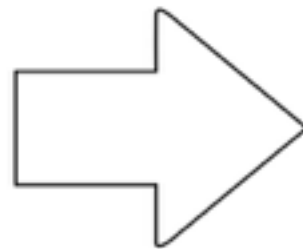


taken from: <http://research.google.com/>



Overlaying satellite images (I)

- Stitch Imagery Data for Google Maps (simplified)
 - Imagery data from different content providers
 - Different formats
 - Different coverages
 - Different timestamps
 - Different resolutions
 - Different exposures/tones
- Large amount of data to be processed
- **Goal:** produce data to serve a "satellite" view to users

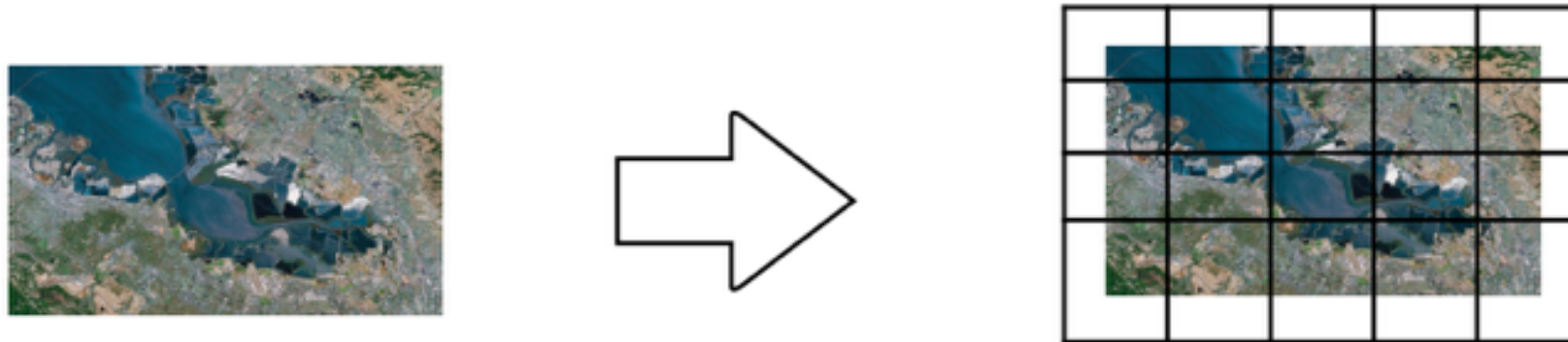


taken from: <http://research.google.com/>

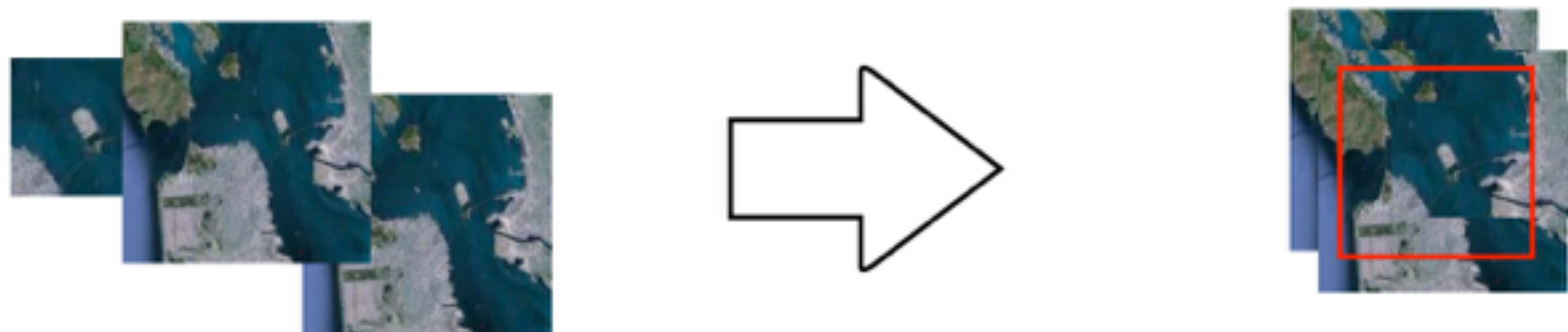


Overlaying satellite images (II)

1. Split the whole territory into "tiles" with fixed location IDs
2. Split each source image according to the tiles it covers



3. For a given tile, stitch contributions from different sources, based on its freshness and resolution, or other preference



4. Serve the merged imagery data for each tile, so they can be loaded into and served from a image server farm.

taken from: <http://research.google.com/>



Overlaying satellite images (III)

map(String key, Image value):

- ▶ // key: image file name
- ▶ // value: image data
- ▶ Tile whole_image;
- ▶ switch (file_type(key)):
 - JPEG: Convert_JPEG(value, whole_image);
 - GIF: Convert_GIF(value, whole_image);
 - ...
- ▶ // split whole_image according to the grid into tiles
- ▶ List<Tile> tile_images = Split_Image(whole_image);

- ▶ for (Tile t: tile_images):
 - **emit**(t.getLocationId(), t);

taken from: <http://research.google.com/>



Overlaying satellite images (IV)

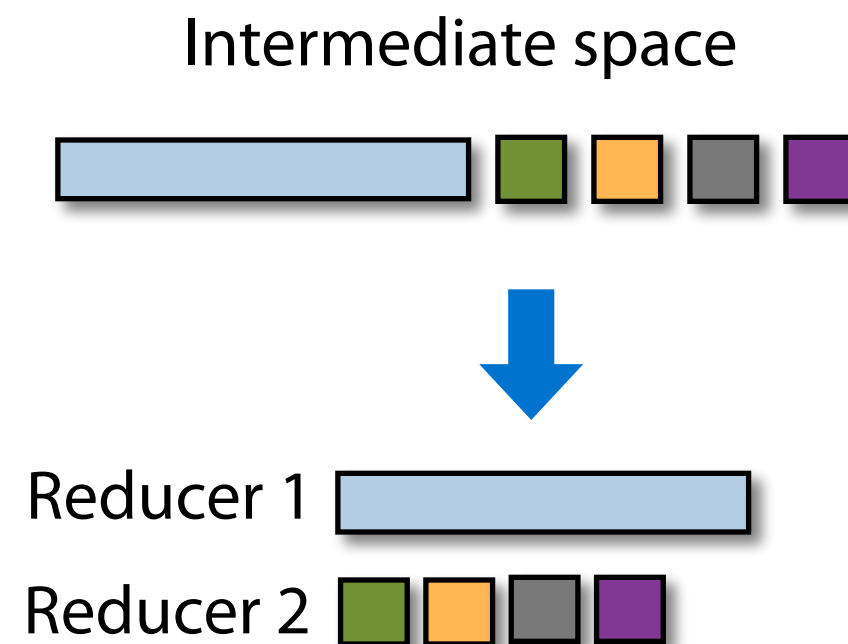
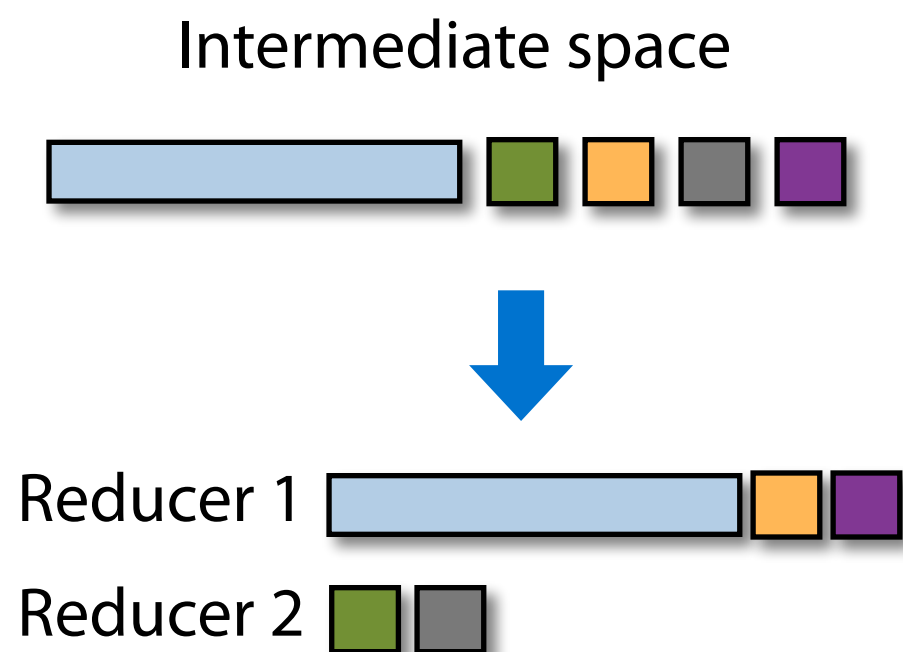
reduce(int key, List<Tile> values):

- ▶ // key: locationId,
- ▶ // values: tiles from different sources
- ▶ // sort values according to resolution and timestamp;
- ▶ Collection.sort(values, ...)
- ▶ Tile mergedTile;
- ▶ for (Tile v: values):
 - // overlay pixels in v to mergedTile based on coverage;
 - mergedTile.overlay(v);
- ▶ // Normalize mergedTile to be the serve tile size;
- ▶ mergedTile.normalize();
- ▶ **emit**(key, mergedTile);

taken from: <http://research.google.com/>

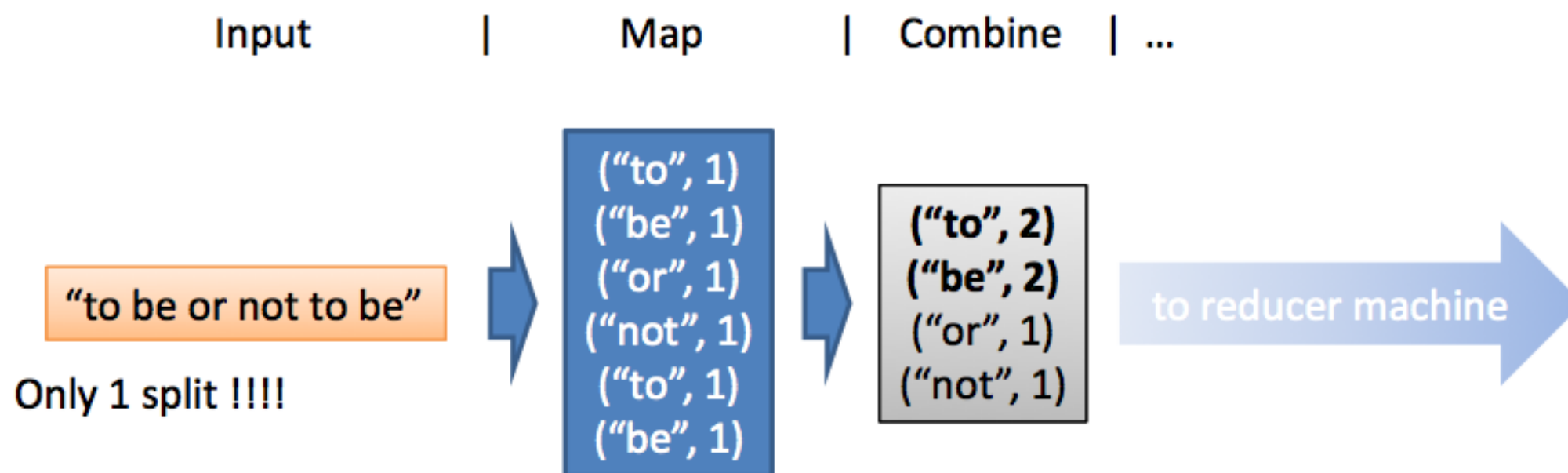
Partitioners

- **Balance the key assignments to reducers**
 - By default, intermediate keys are hashed to reducers
 - Partitioner specifies the node to which an intermediate key-value pair must be copied
 - Divides up key space for parallel reduce operations
 - Partitioner only considers the key and ignores the value

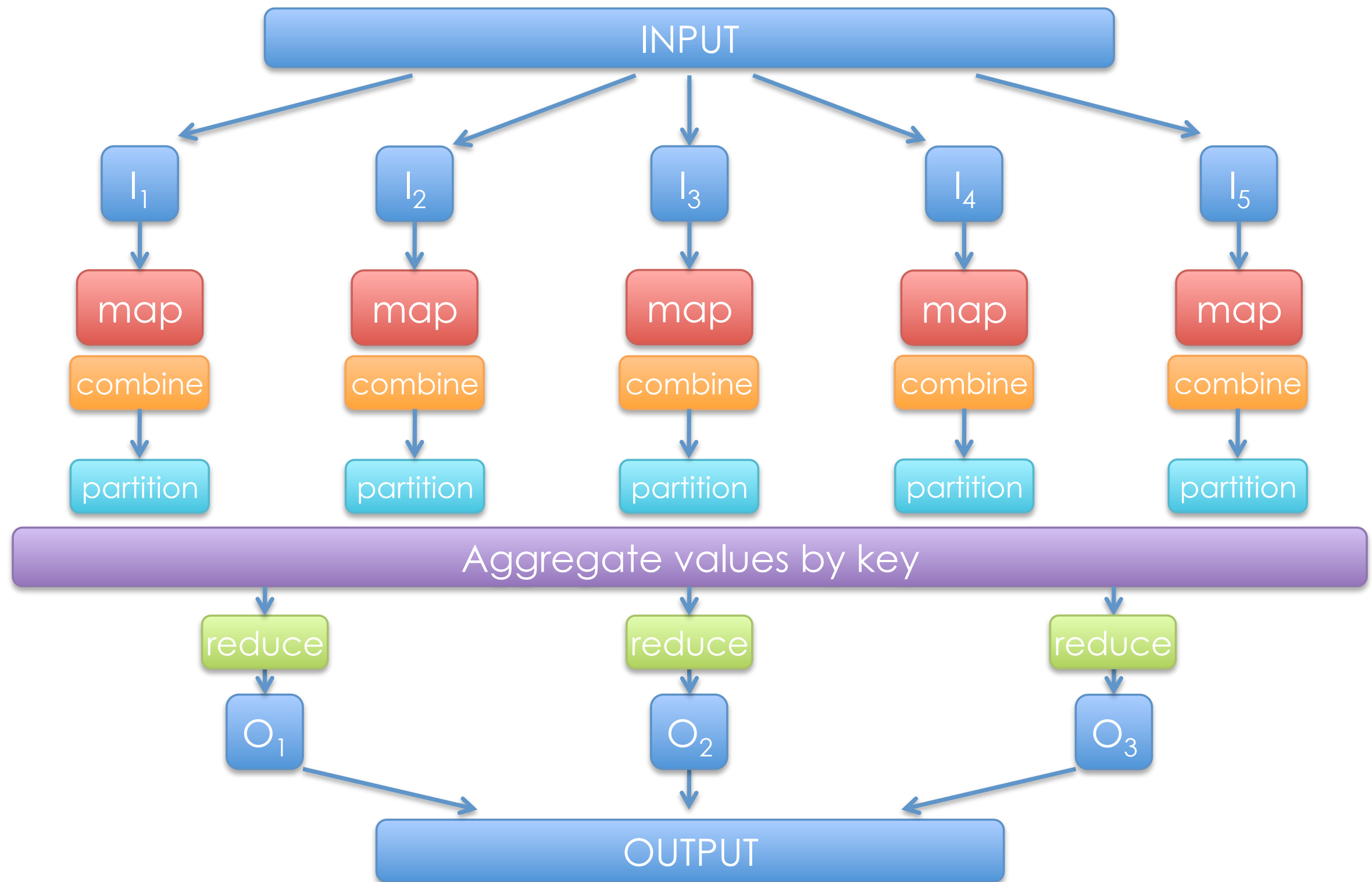


Combiners

- **Local aggregation before the shuffle**
 - All the key-value pairs from mappers need to be copied across the network
 - The amount of intermediate data may be larger than the input collection itself
 - Perform local aggregation on the output of each mapper (same machine)
 - Typically, a combiner is a (local) copy of the reducer



Programming Model (complete)

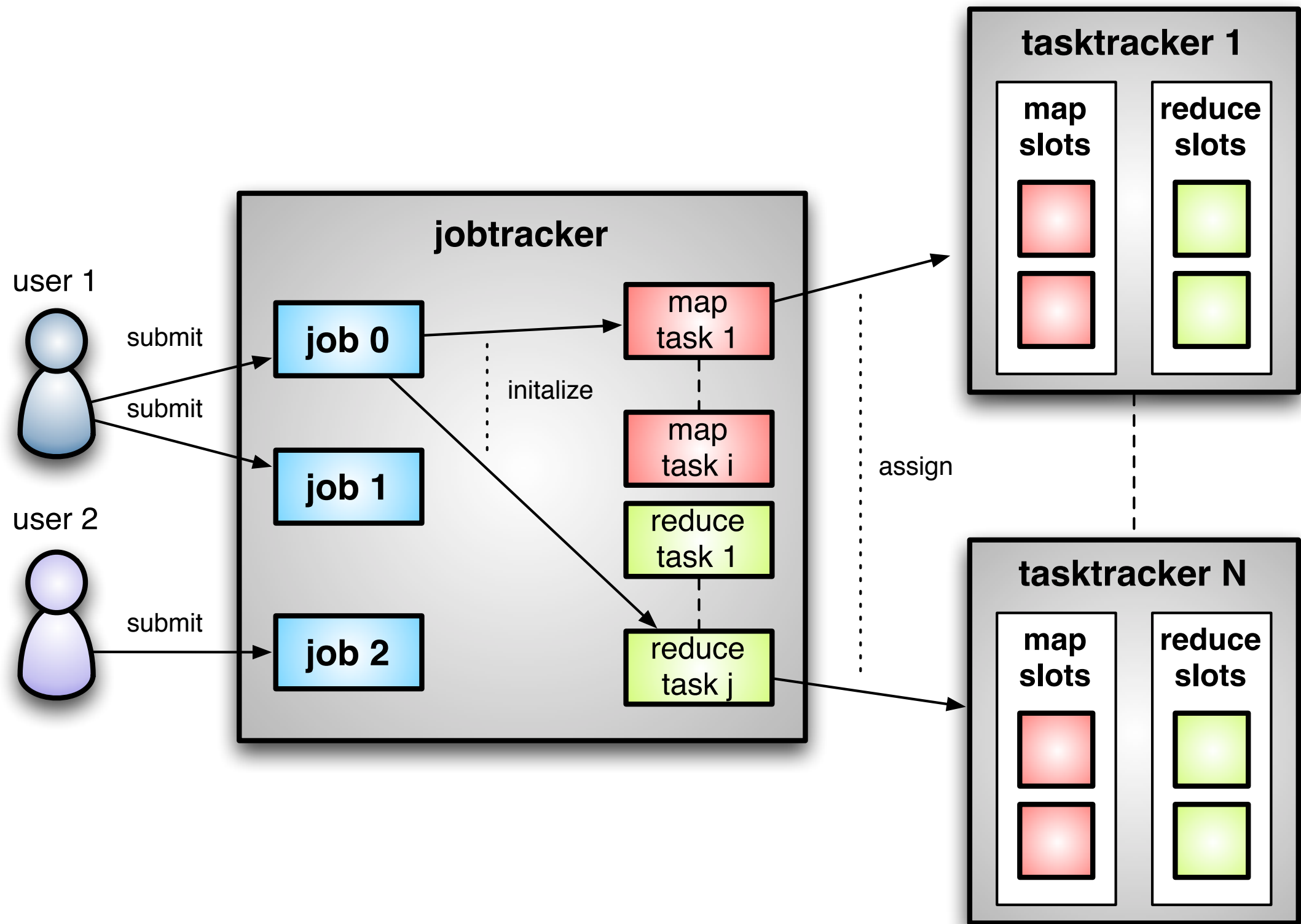


- **Handles scheduling**
 - Assigns workers to map and reduce tasks
- **Handles “data distribution”**
 - Moves processes to data
- **Handles synchronization**
 - Gathers, sorts, and shuffles intermediate data
- **Handles errors and faults**
 - Detects worker failures and restarts
- **Everything happens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

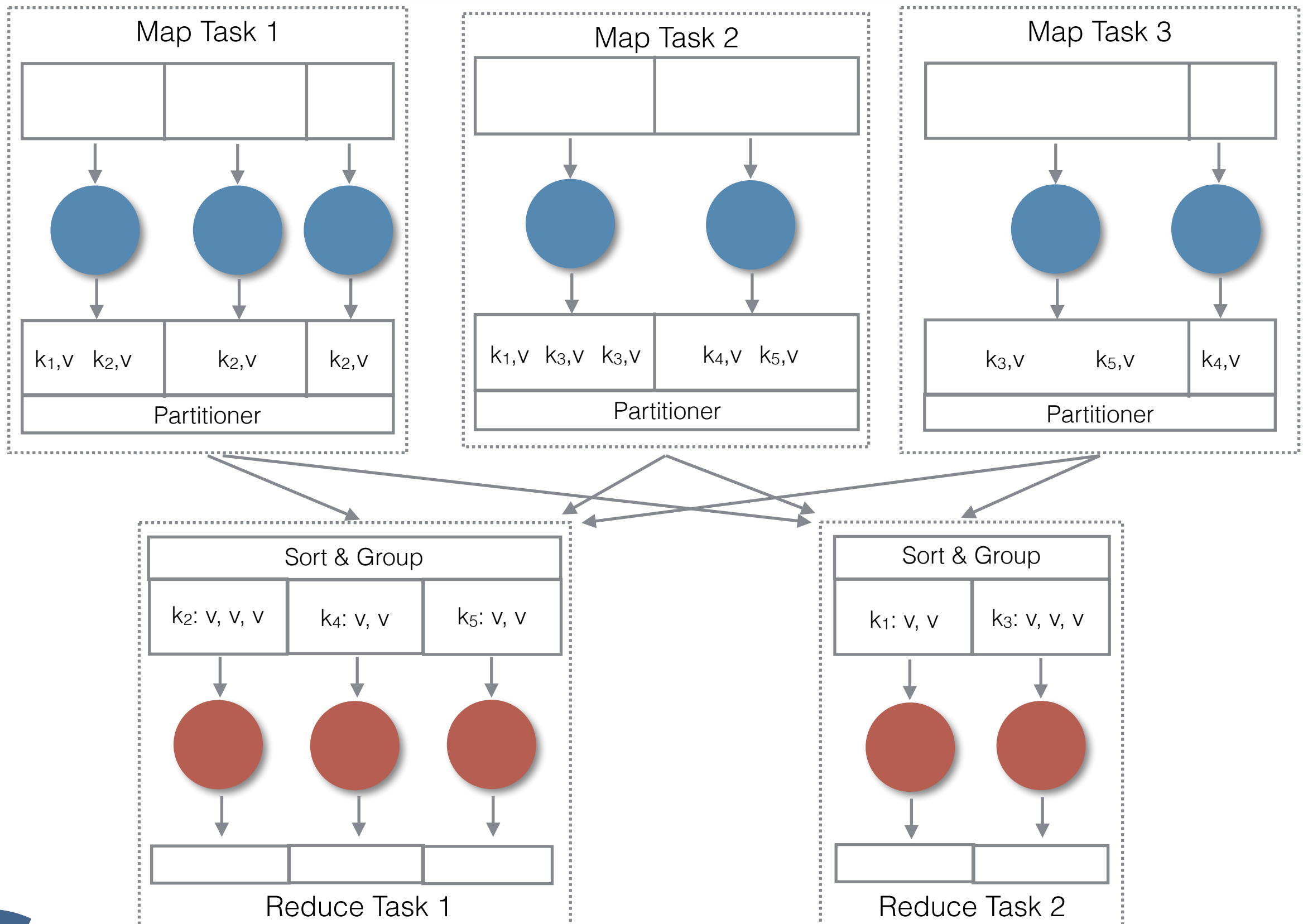
Terminology

- **Job**
- **Task**
- **Slot**
- **JobTracker**
 - Accepts Map/Reduce jobs submitted by users
 - Assigns Map and Reduce tasks to Task Trackers
 - Monitors task and Task Tracker status, re-executes tasks upon failure
- **TaskTracker**
 - Run Map and Reduce tasks upon instruction from the Job Tracker
 - Manage storage and transmission of intermediate output
- **Splits**
 - Data locality optimization

Runtime



Diagram



Scheduling

- **One master, many workers**

- Input data split into M map tasks (typically 64 MB in size)
- Reduce phase partitioned into R reduce tasks ($\text{hash}(k) \bmod R$)
- Tasks are assigned to workers dynamically
- Often: $M=200,000$; $R=4000$; workers=2000

- **Master assigns each map task to a free worker**

- Considers locality of data to worker when assigning a task
- Worker reads task input (often from local disk)
- Worker produces R local files containing intermediate k/v pairs

- **Master assigns each reduce task to a free worker**

- Worker reads intermediate k/v pairs from map workers
- Worker sorts & applies user's reduce operation to produce the output

Parallelism

- **Map functions run in parallel, create intermediate values from each input data set**
 - The programmer must specify a proper input split (chunk) between mappers to enable parallelism
- **Reduce functions also run in parallel, each will work on different output keys**
 - Number of reducers is a key parameter which determines map-reduce performance

Speculative Execution

- **Problem: Stragglers (i.e., slow workers) significantly lengthen the completion time**
 - Other jobs may be consuming resources on machine
 - Bad disks with soft (i.e., correctable) errors transfer data very slowly
 - Other weird things: processor caches disabled at machine init
- **Solution: Close to completion, spawn backup copies of the remaining in-progress tasks.**
 - Whichever one finishes first, "wins"
- **Additional cost: a few percent more resource usage**
- **Example: A sort program without backup = 44% longer.**

Fault Tolerance (I)

Failures in Literature

- **LANL data** (DSN 2006)
 - Data collected over 9 years
 - Covered 4,750 machines and 24,101 CPUs
 - Distribution of failures
 - Hardware ~ 60%, Software ~ 20%, Network/Environment/Humans ~ 5%, Aliens ~ 25%*
 - Depending on a system, failures occurred between once a day to once a month
 - Most of the systems in the survey were the cream of the crop at their time
- **PlanetLab** (SIGMETRICS 2008 HotMetrics Workshop)
 - Average frequency of failures per node in a 3-months period
 - Hard failures: 2.1
 - Soft failures: 41
 - Approximately failure every 4 days

taken from: <http://research.google.com/>

Fault Tolerance (II)

Failures in Google Data Centers

- **DRAM errors analysis** (SIGMETRICS 2009)
 - Data collected over 2.5 years
 - 25,000 to 70,000 errors per billion device hours per Mbit
 - Order of magnitude more than under lab conditions
 - 8% of DIMMs affected by errors
 - Hard errors are dominant cause of failure
- **Disk drive failure analysis** (FAST 2007)
 - Annualized Failure Rates vary from 1.7% for one year old drives to over 8.6% in three year old ones
 - Utilization affects failure rates only in very old and very old disk drive populations
 - Temperature change can cause increase in failure rates but mostly for old drives

taken from: <http://research.google.com/>

MapReduce Fault Tolerance

- **Master keeps track of progress of each task and worker nodes**
 - If a node fails, it re-executes the completed as well as in-progress map tasks on other nodes that are alive
 - It also executes in-progress reduce tasks.
- **If particular input key/value pairs keep crashing**
 - Master blacklists them and skips them from re-execution
- **Tolerate small failures, allow the job to run in best-effort basis**
 - For large datasets containing potentially millions of records, we don't want to stop computation for a few records not processing correctly
 - User can set the failure tolerance level

Performance

- **Maximizing Map input transfer rate**
 - Input Locality
 - Minimal deserialization overhead
- **Small intermediate output**
 - $M \times R$ transfers over the network
 - Minimize/compress transfers
 - Avoid shuffling/sorting if possible (e.g. map-only computations)
 - Use combiners and/or partitioners!!!
 - Compress everything (automatic)
- **Opportunity to Load Balance**
- **Changing algorithm to suit architecture yields best implementation**

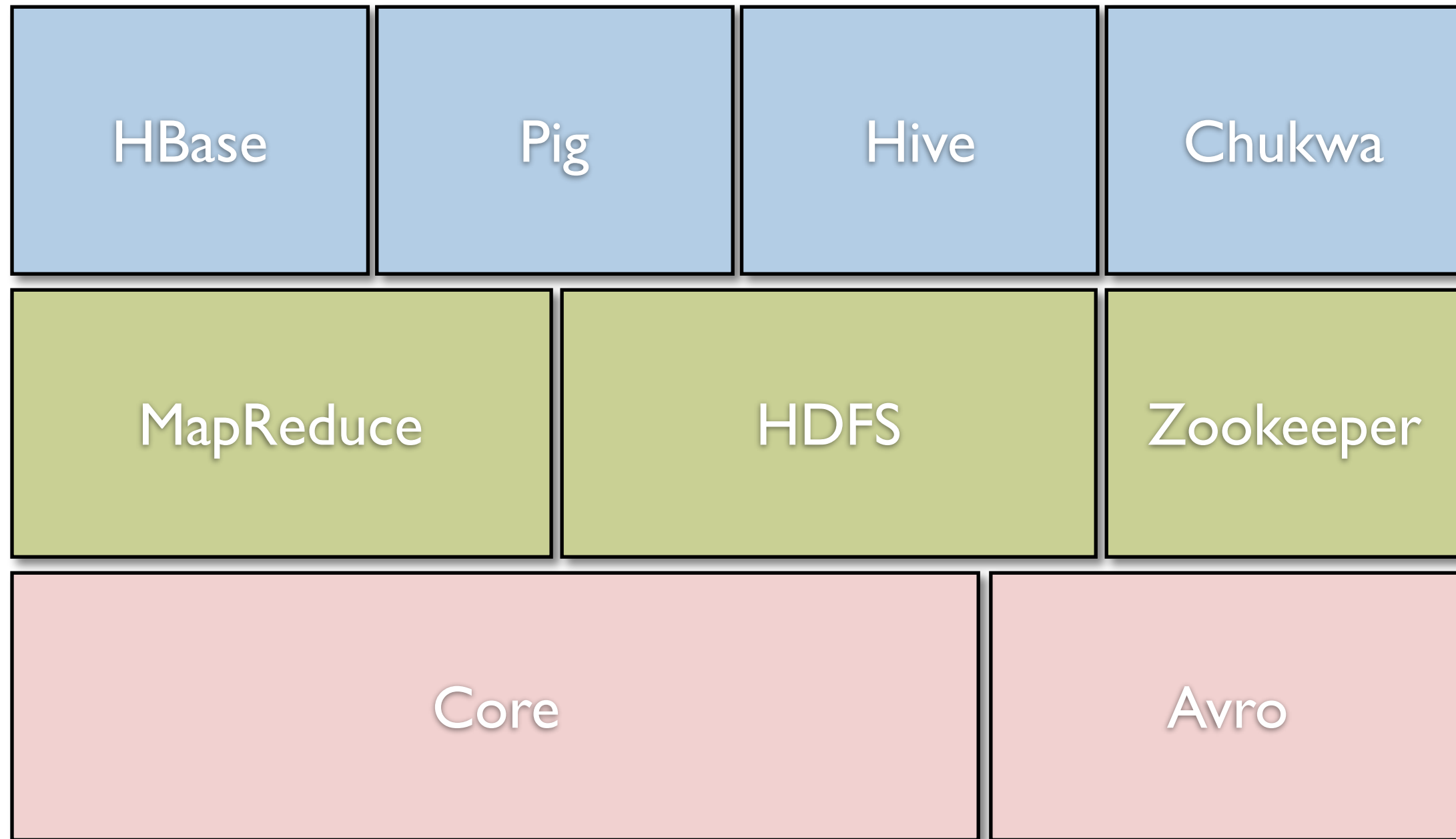


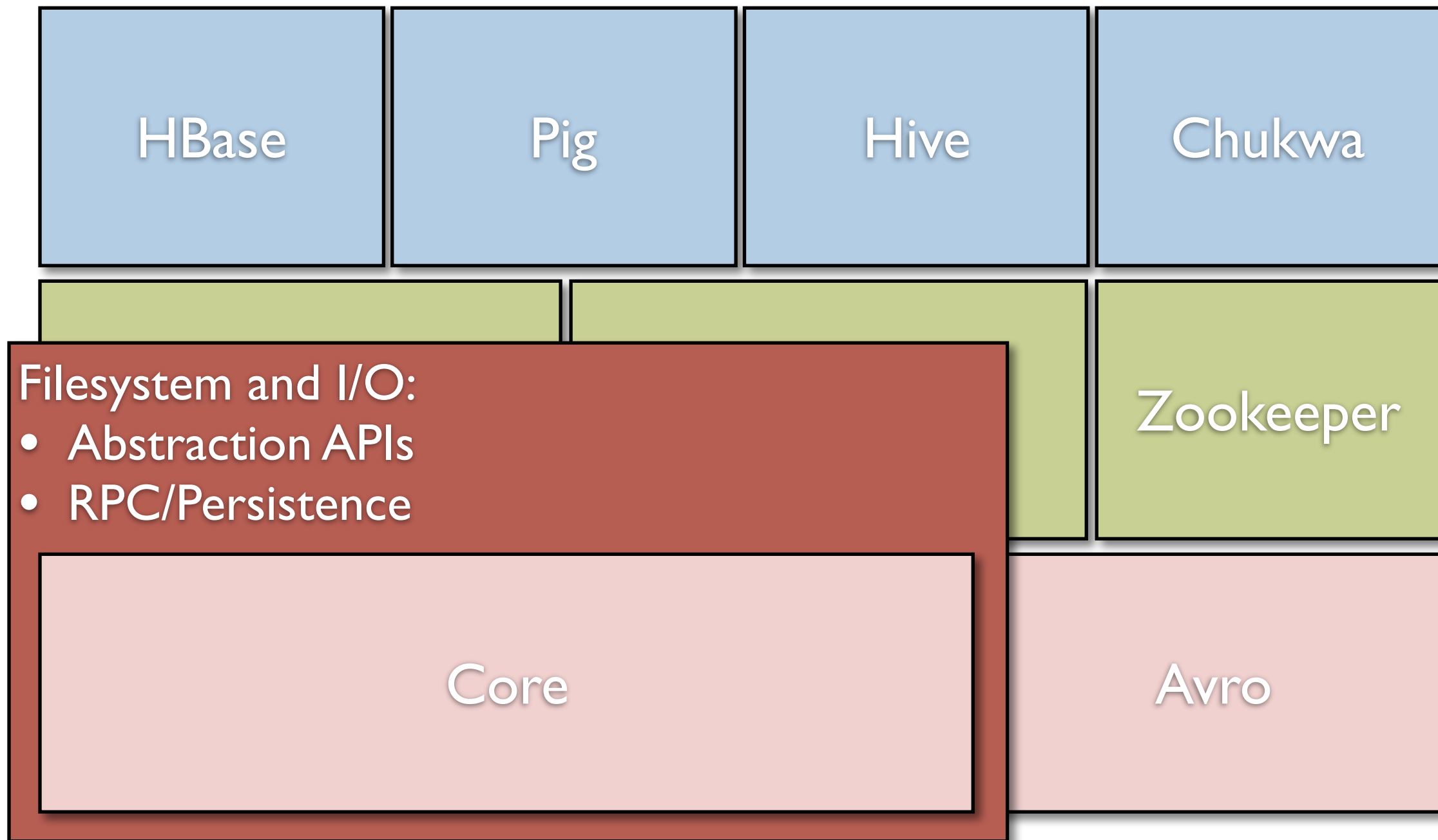
The right tool for the right job

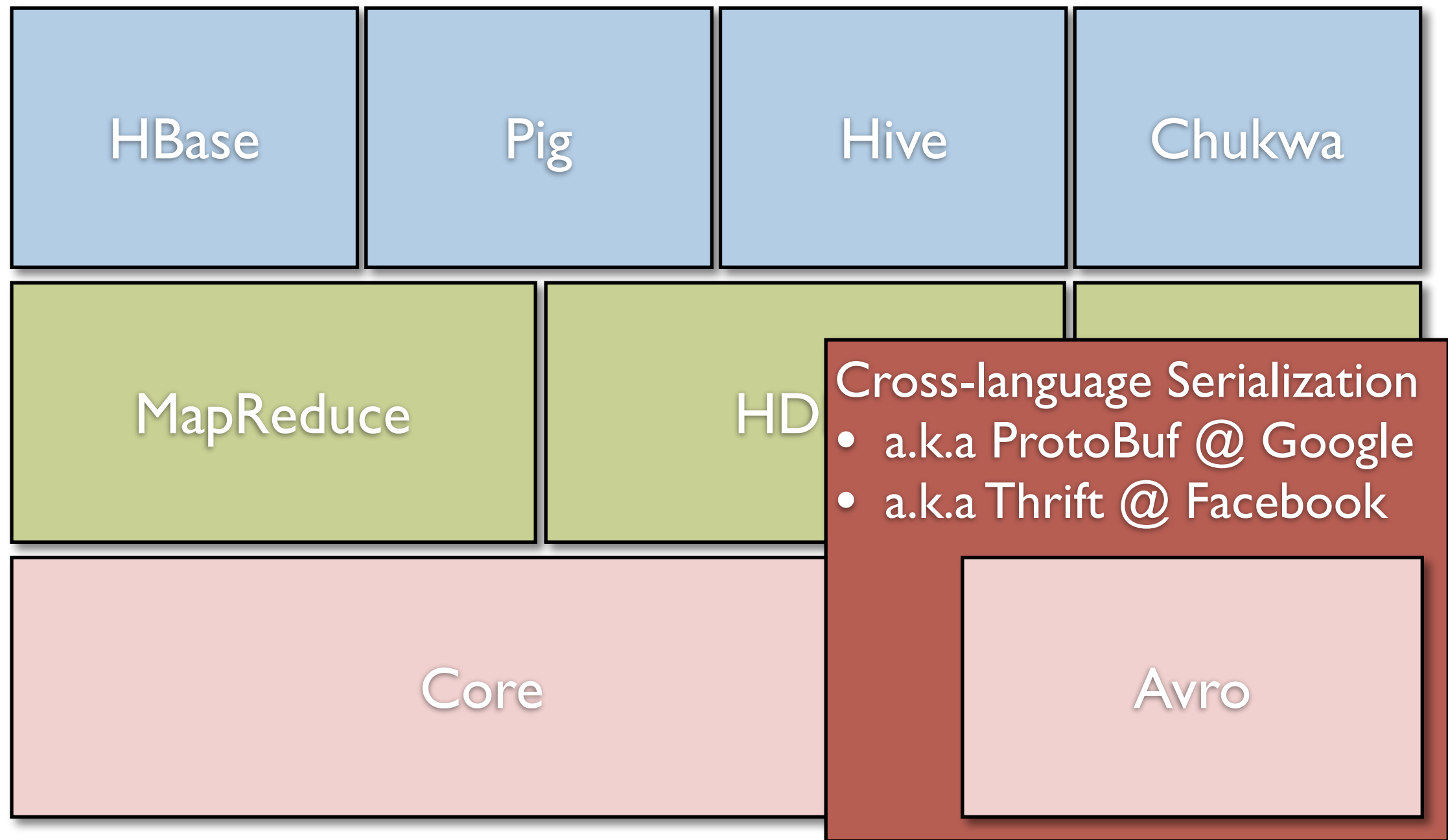
| | MPI | MapReduce | DBMS/SQL |
|------------------------|--|--|--|
| What it is | A general parallel programming model | A programming paradigm and its associated execution system | A system to store, manipulate and serve data |
| Programming Model | Message passing between nodes | Restricted to map and reduce operations | Declarative on data query Stored procedures |
| Data Organization | No assumptions | Shared files | Organized data structures |
| Data to be manipulated | Any | key-value pairs | Tables with rich attributes |
| Execution Model | Nodes are independent | Map/Shuffle/Reduce Checkpointing/Backup Physical data locality | Transaction Query optimization Materialized view |
| Usability | Steep learning curve difficult to debug | Simple concept Could be hard to optimize | Declarative interface Could be hard to debug at runtime |
| Key Selling Point | Flexible to accommodate various applications | Process large amount of data with commodity hardware | Interactive querying Maintain a consistent view across client |

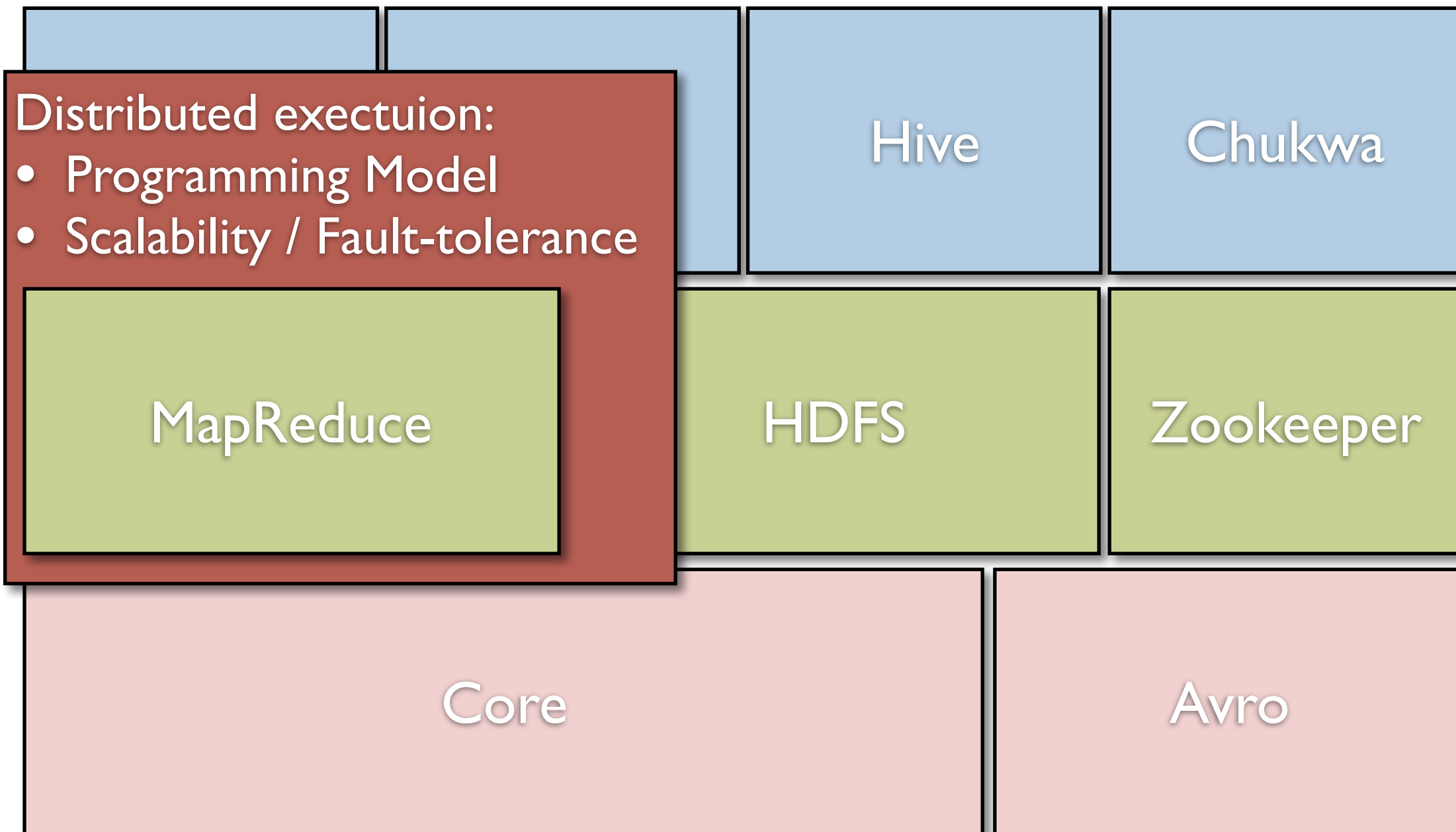
taken from: <http://research.google.com/>

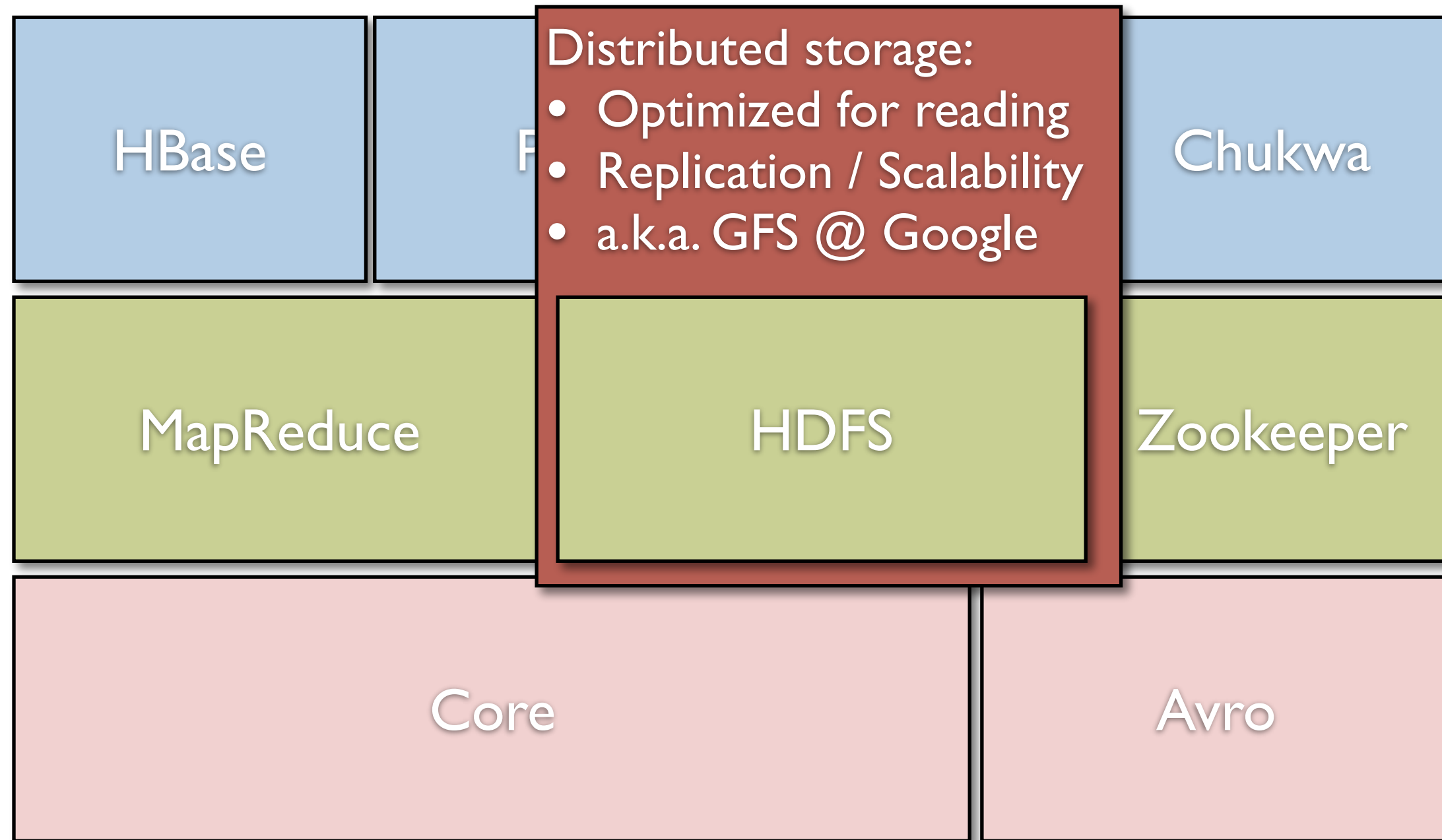
HADOOP

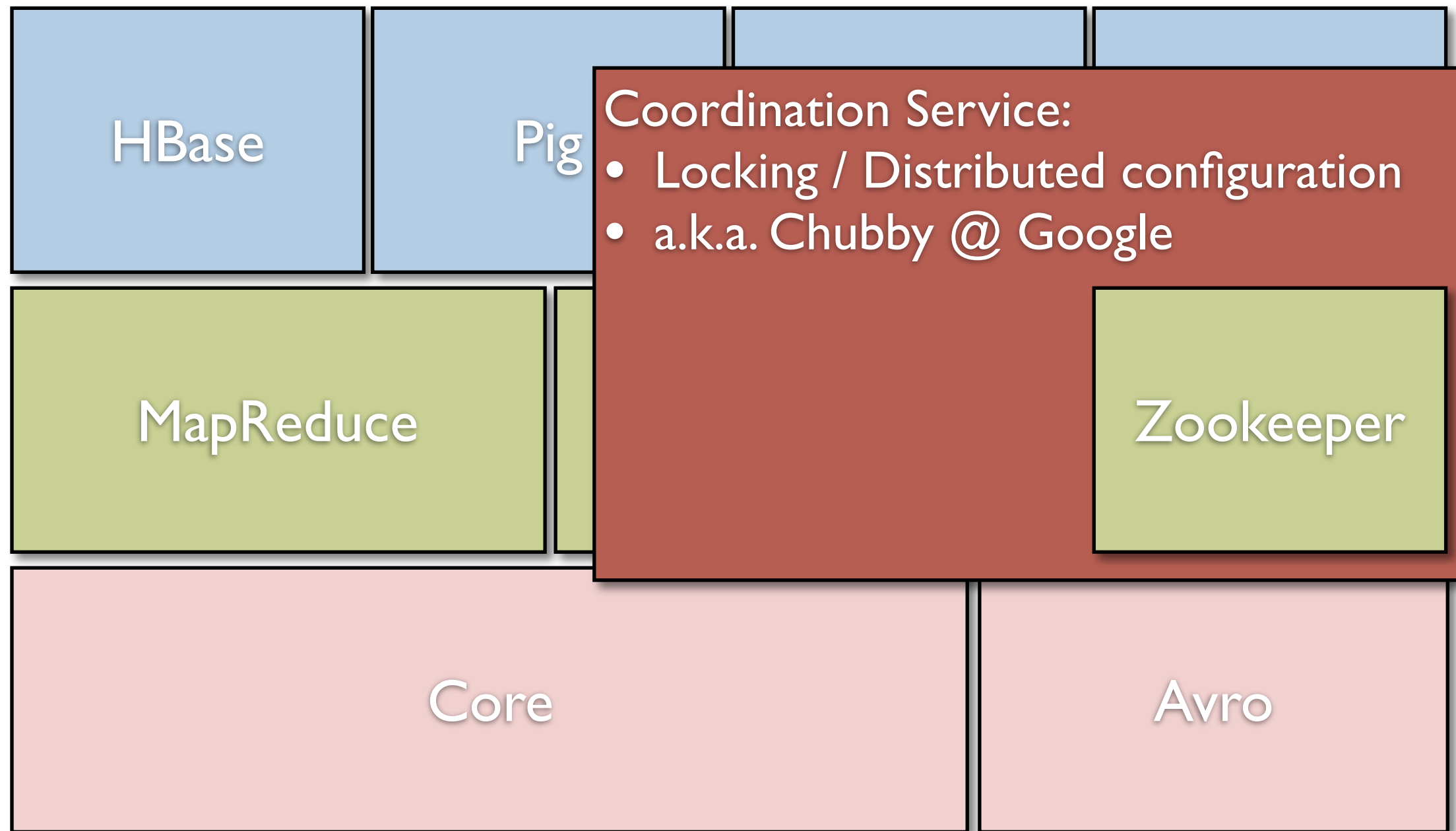


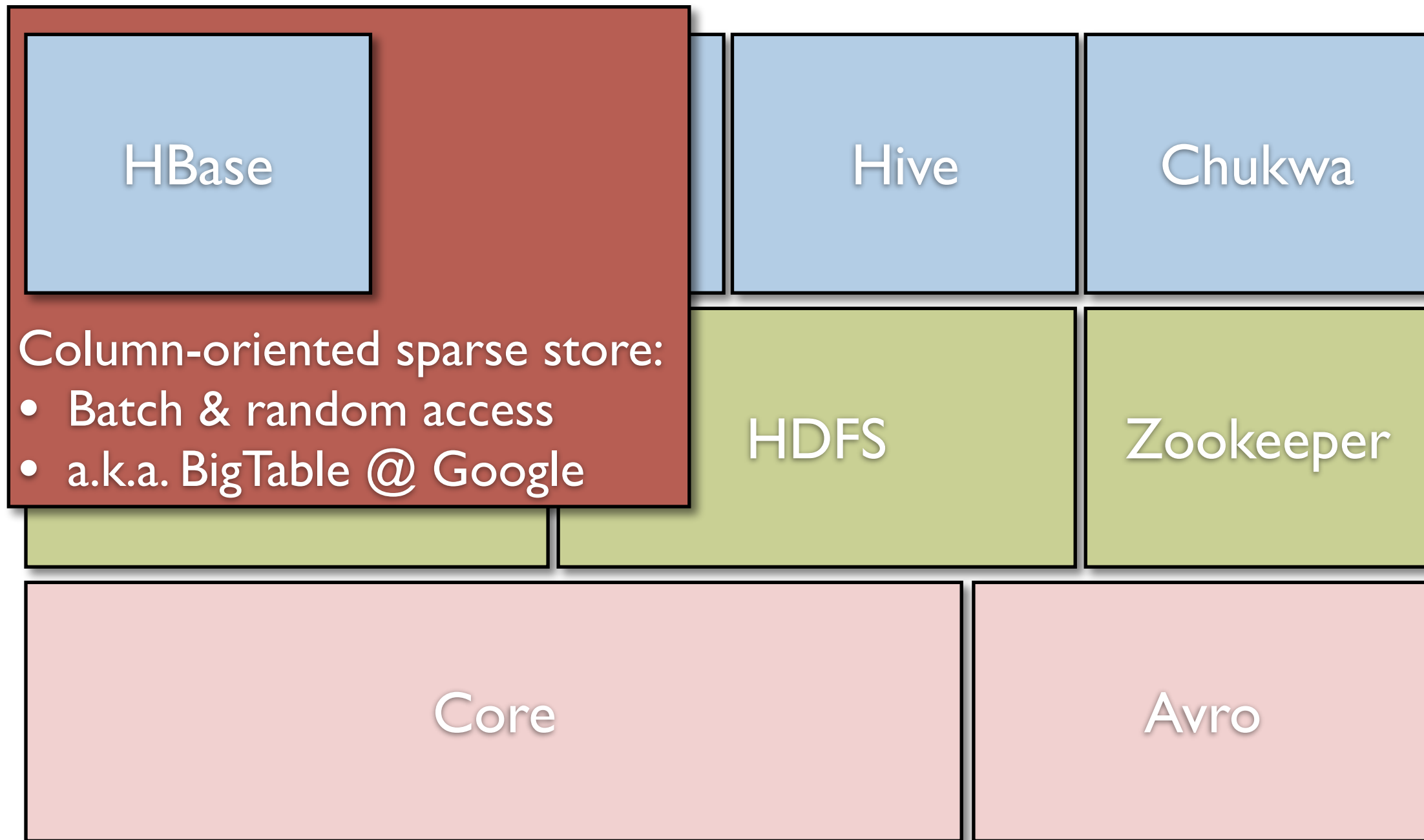


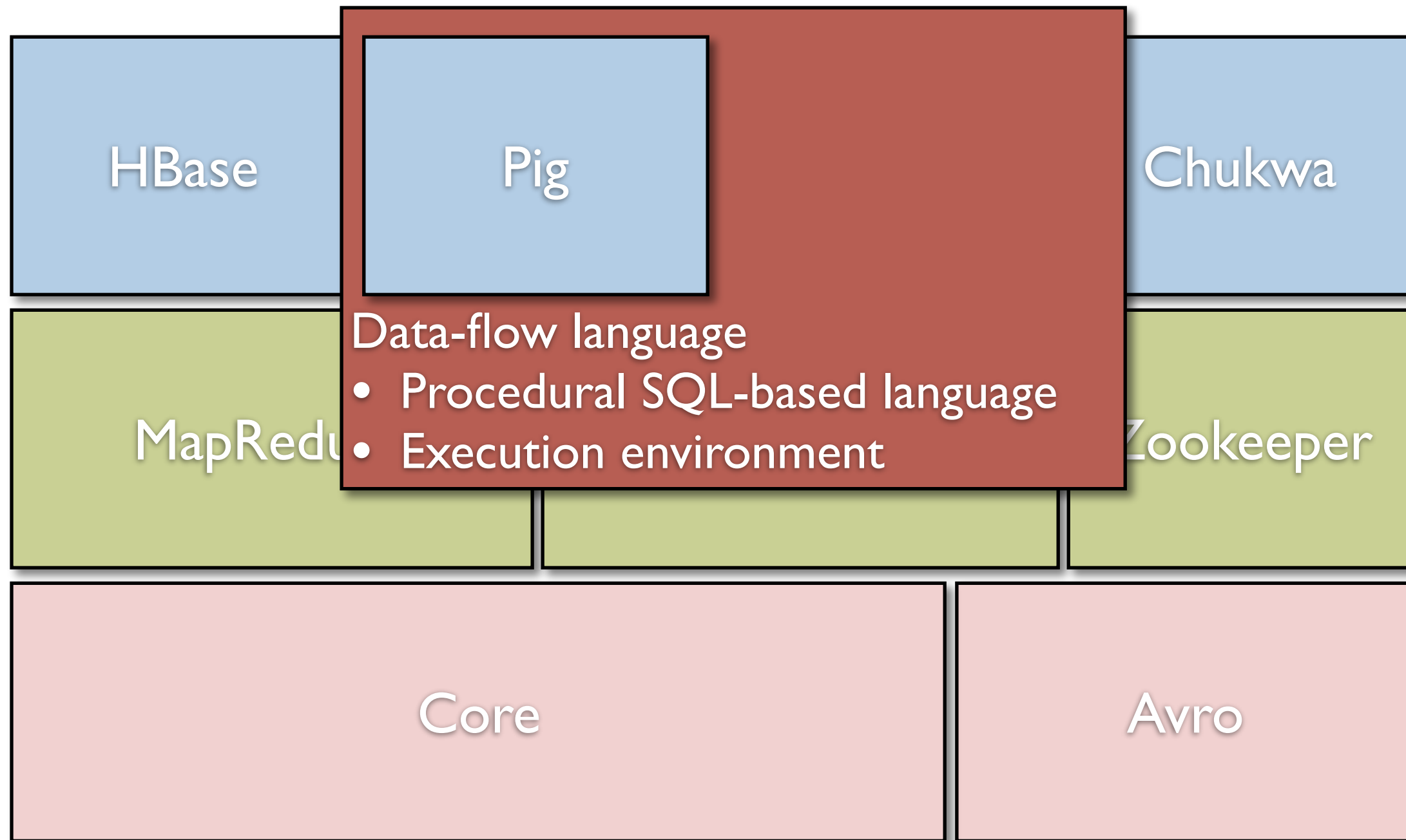


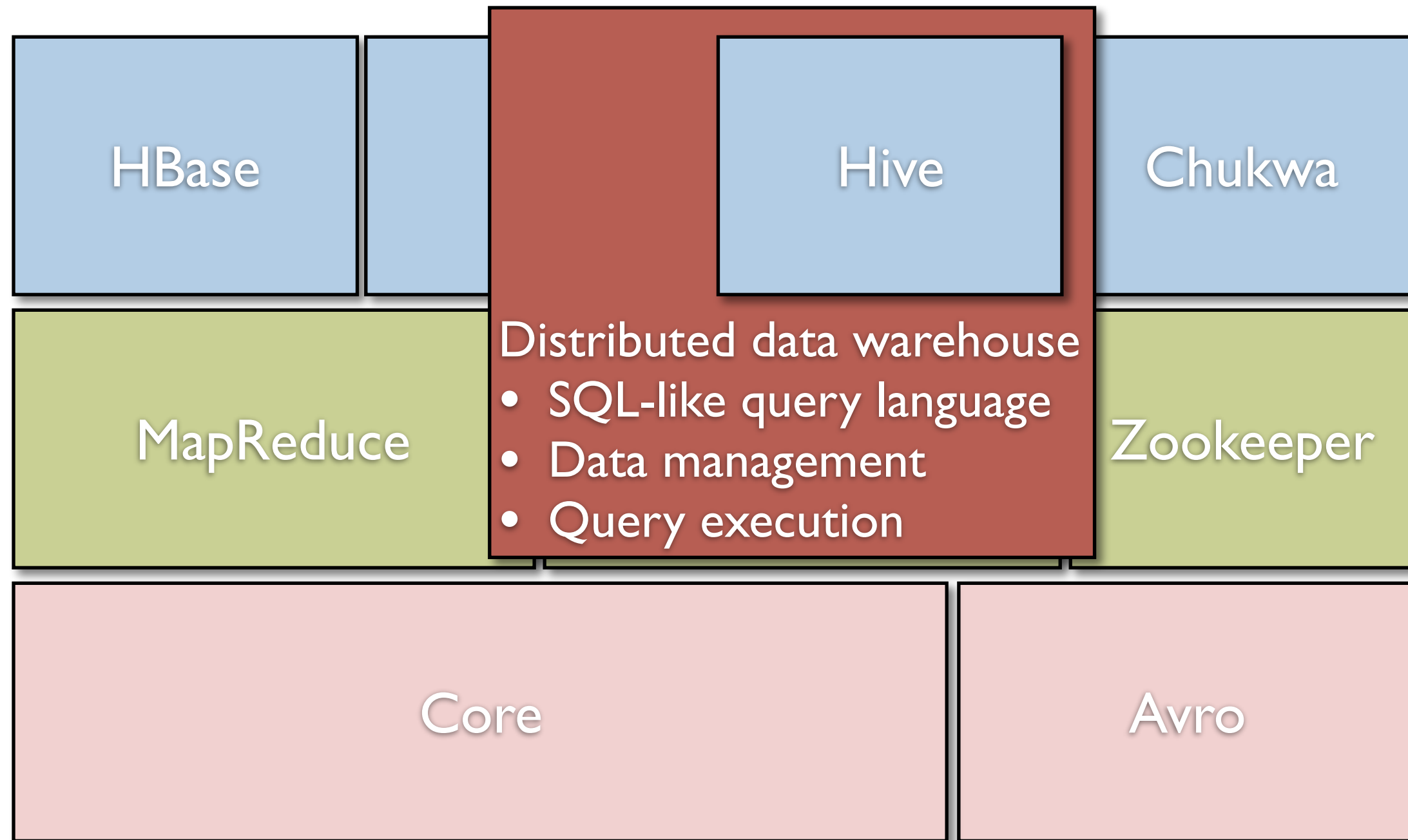




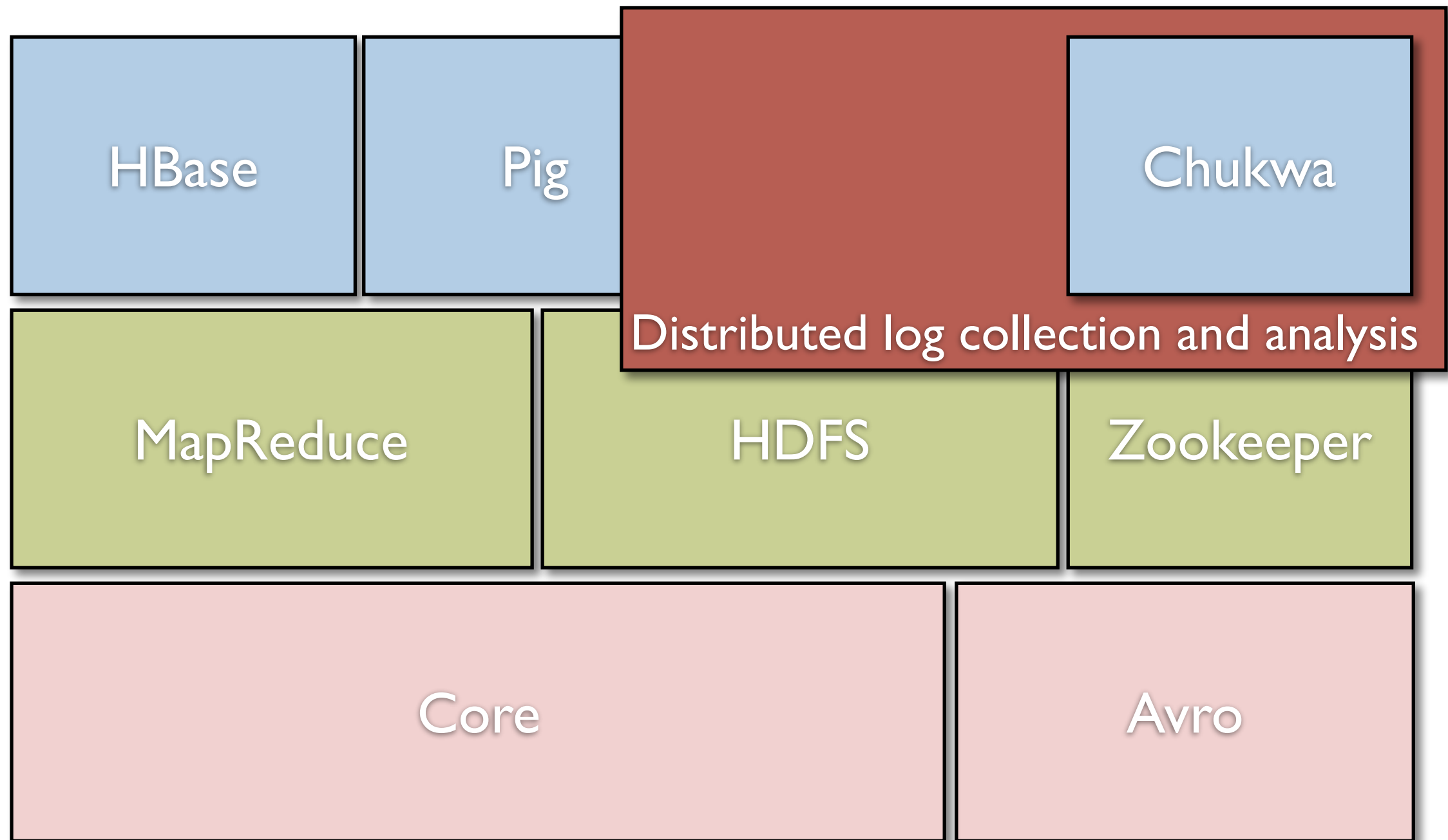




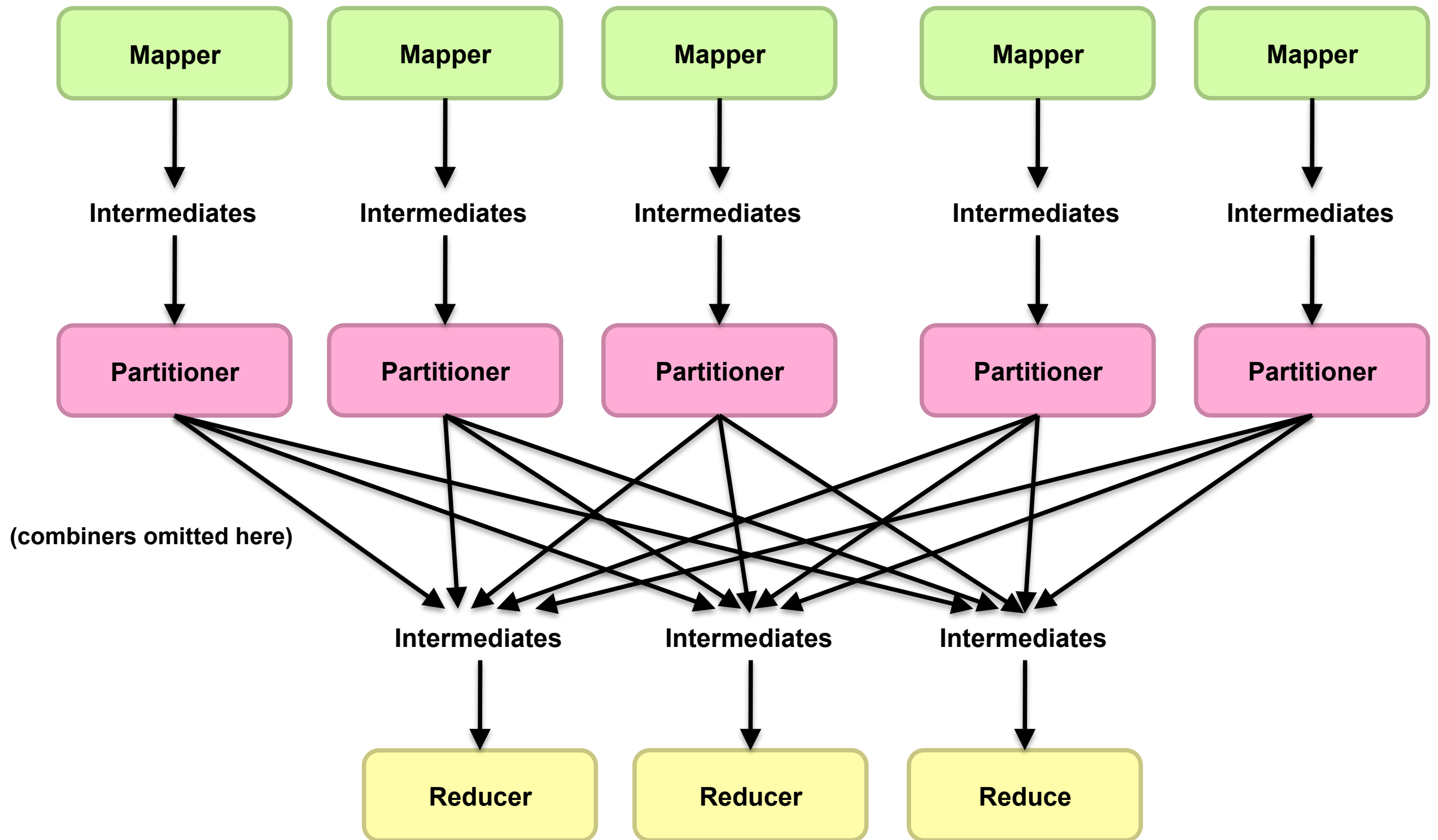


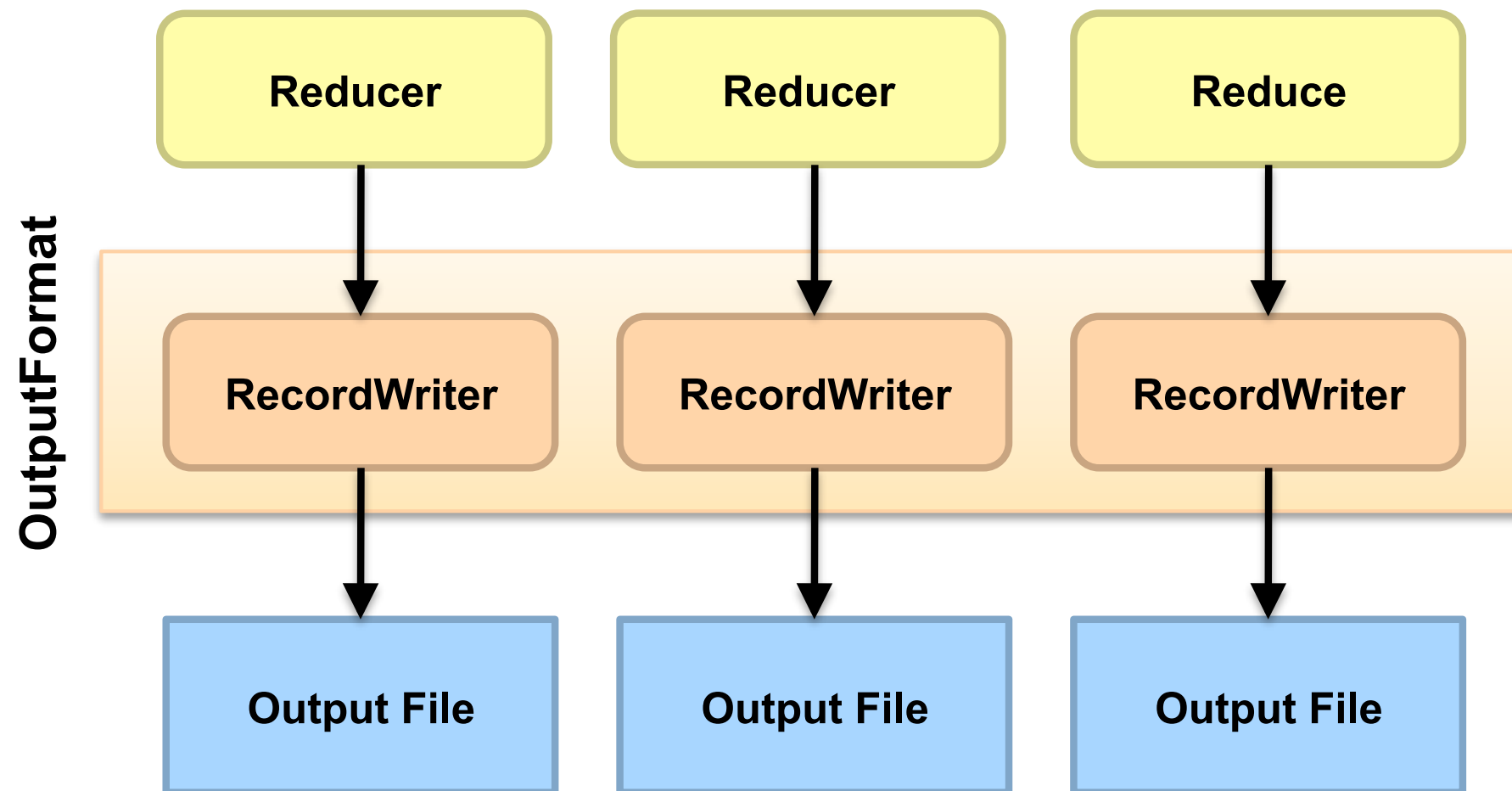


HADOOP



Hadoop Dataflow (II)





- Data sets are specified by **OutputFormats**
 - Analogous to InputFormat
- Base class implementation is **FileOutputFormat**
 - TextOutputFormat – Writes “key val\n” strings to output file
 - SequenceFileOutputFormat – Uses a binary format to pack (k, v) pairs
- Other implementation is **NullOutputFormat**
 - Discards output to /dev/null

Hadoop Shuffle & Sort

- **Map Side**

- Mapper outputs are buffered in memory in a circular buffer
- When buffer reaches threshold, contents are “spilled” to disk
- Spills are merged in a single partitioned file (sorted within each partition)
- Combiners run here

- **Reduce Side**

- Firstly, mapper outputs are copied over to the reducer machine
- “Sort” is a multi-pass merge of map outputs (in memory and on disk)
- Combiners run here
- Final merge pass goes directly into reducer

- **Probably the most complex aspect of the framework!**