

Harp Tutorial

Bingjing Zhang



INDIANA UNIVERSITY BLOOMINGTON

SCHOOL OF INFORMATICS AND COMPUTING

Outline

Concepts

Data Types

APIs

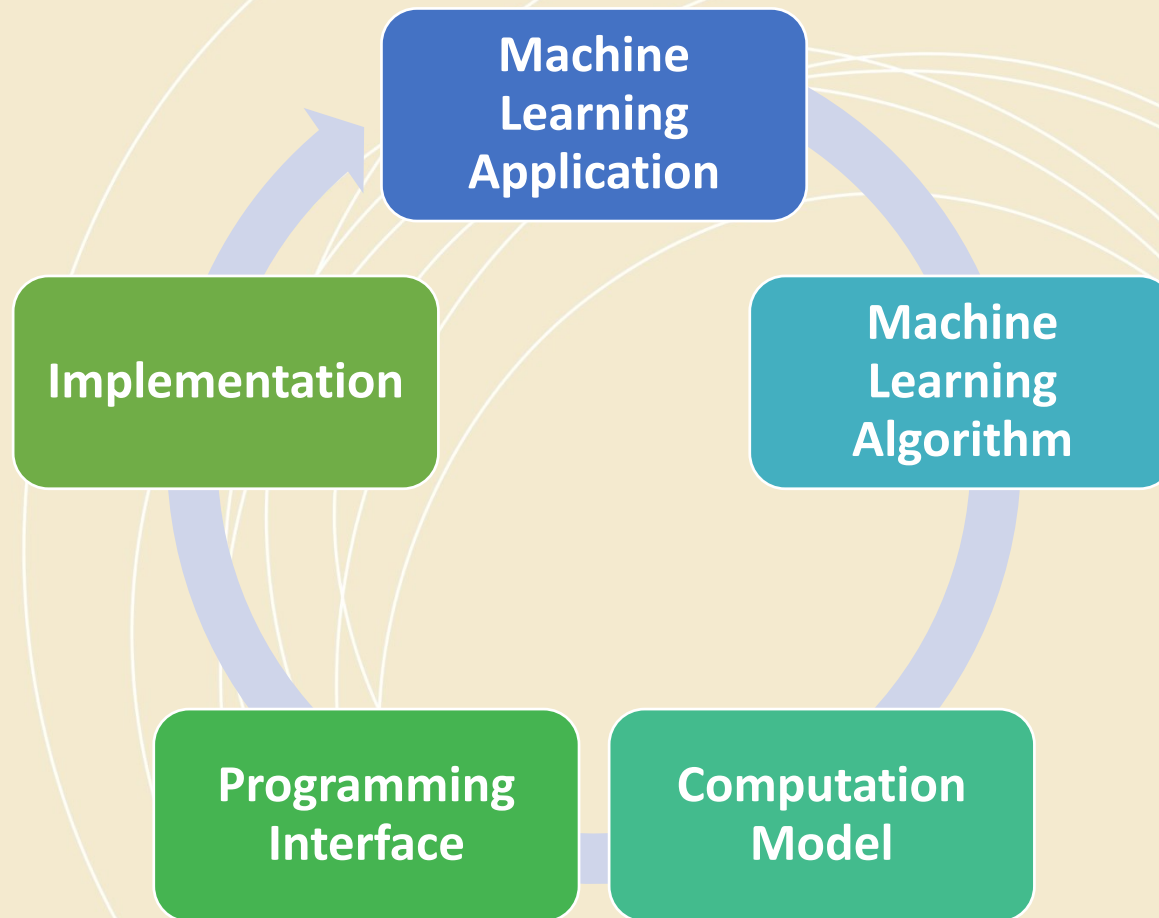
Application Examples



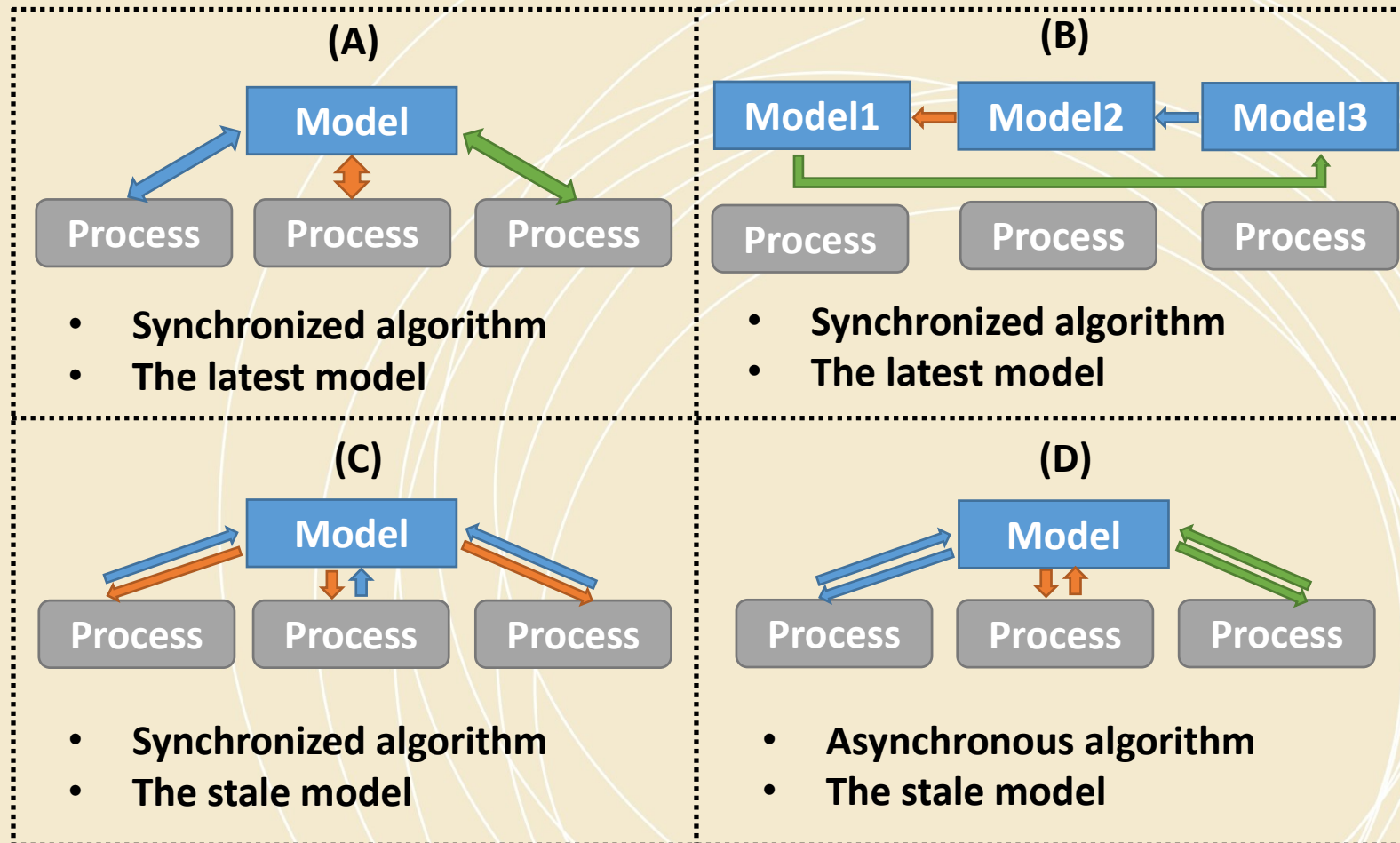
INDIANA UNIVERSITY BLOOMINGTON

SCHOOL OF INFORMATICS AND COMPUTING

Parallelization of Machine Learning Applications



Inter-node Computation Models



Inter-node Computation Models

(Training Data Items Are Partitioned to Each Process)

Computation Model A

- Once a process trains a data item, it locks the related model parameters and prevents other processes from accessing them. When the related model parameters are updated, the process unlocks the parameters. Thus the model parameters used in local computation is always the latest.

Computation Model B

- Each process first takes a part of the shared model and performs training. Afterwards, the model is shifted between processes. Through model rotation, each model parameters are updated by one process at a time so that the model is consistent.

Computation Model C

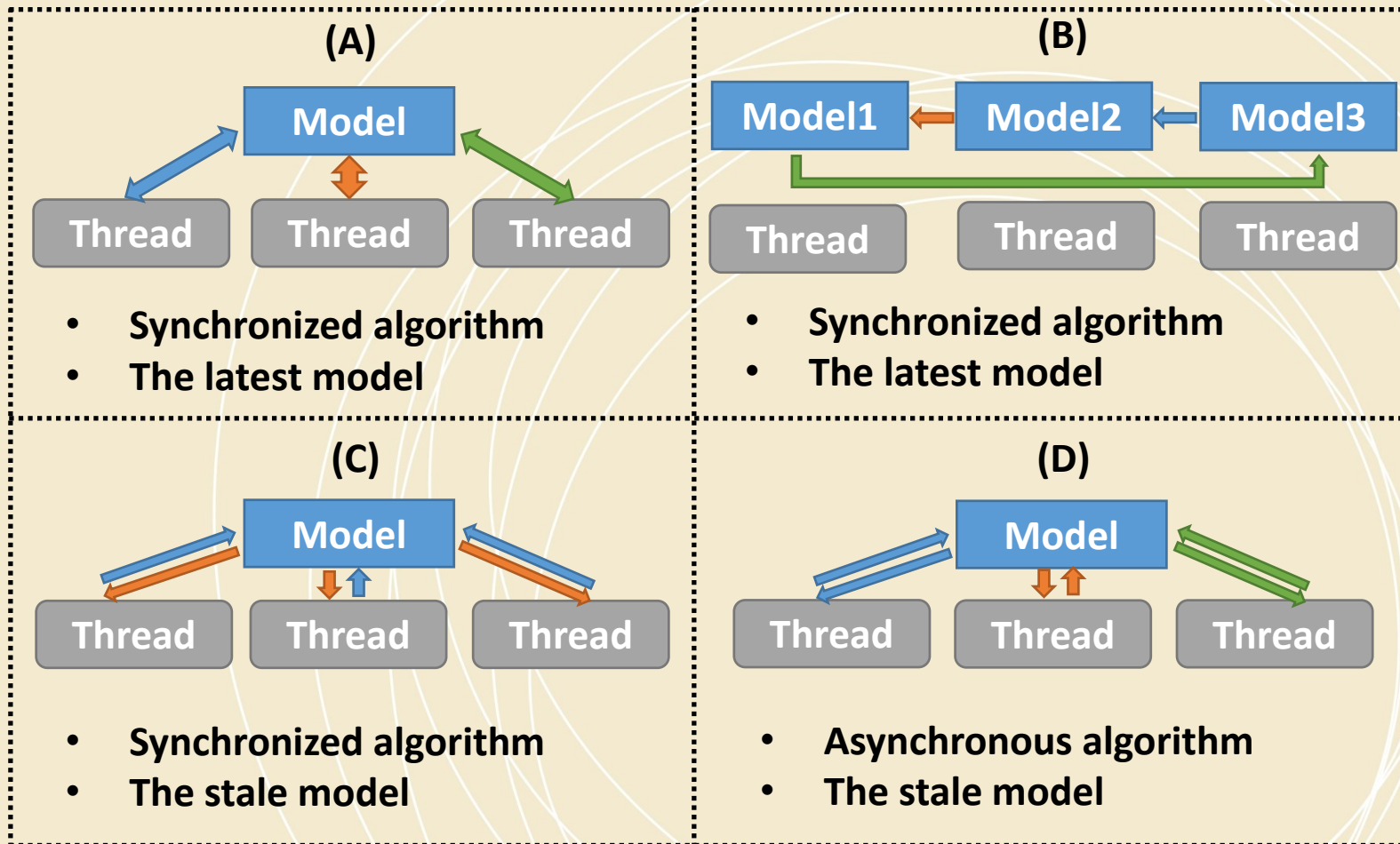
- Each process first fetches all the model parameters required by local computation. When the local computation is completed, modifications of the local model from all processes are gathered to update the model.

Computation Model D

- Each process independently fetches related model parameters, performs local computation, and returns model modifications. Unlike A, workers are allowed to fetch or update the same model parameters in parallel. In contrast to B and C, there is no synchronization barrier.



Intra-node Computation Models



Intra-node Computation Models

(Training Data Items Are Scheduled to Each Thread)

Computation Model A

- Once a thread trains a data item, it locks the related model parameter and prevents other threads from accessing it. When the model parameter is updated, the thread unlocks the parameter. Thus the model parameters used per thread is always the latest.

Computation Model B

- Each thread first takes a part of the shared model and performs training. Afterwards, the model part is sent or scheduled to another thread. Model parameters in each model part are updated by one thread at a time so that the model is consistent.

Computation Model C

- Each thread has a copy of the related model parameters. When the local computation is completed, modifications of the local model from all threads are combined to update the model.

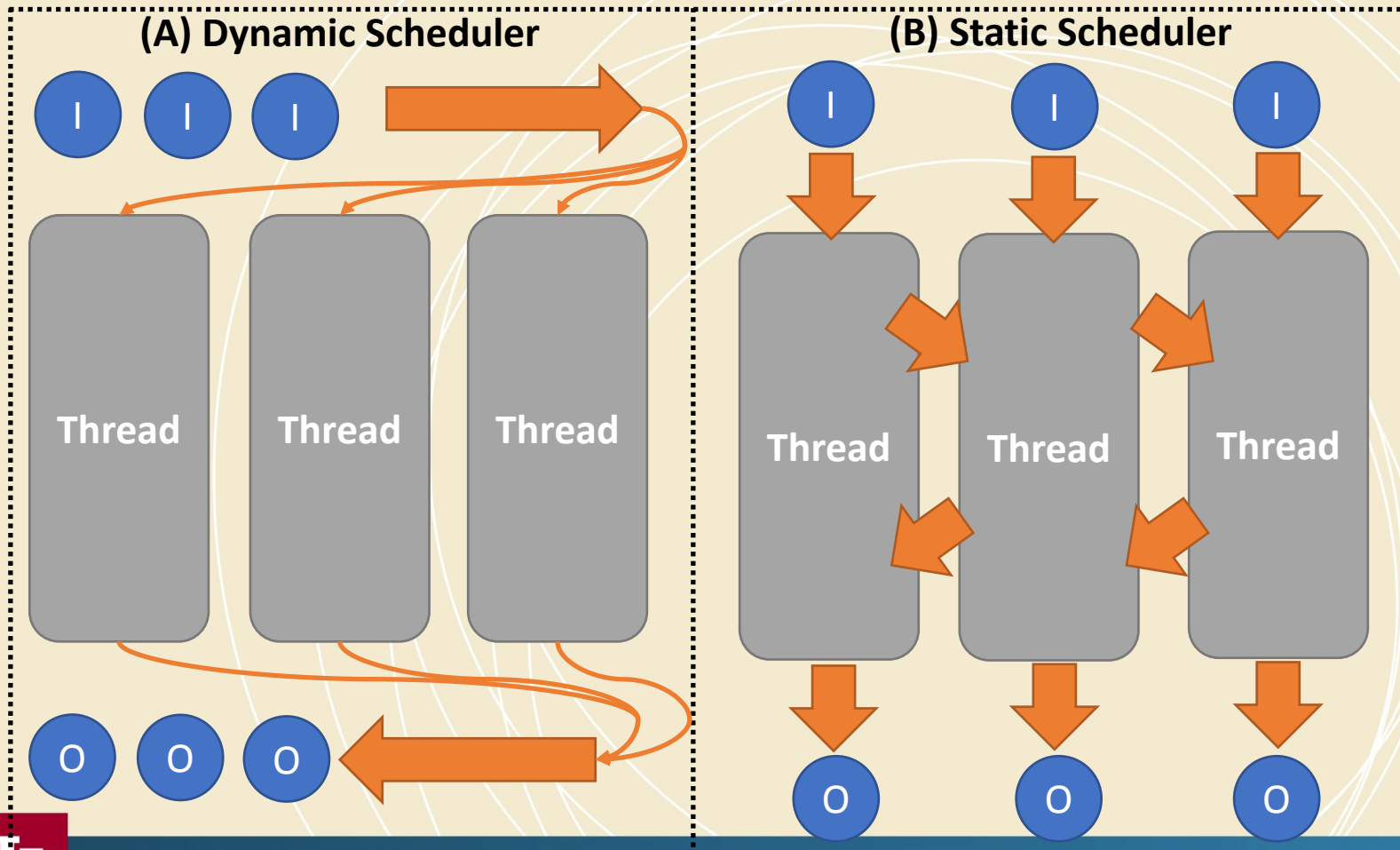
Computation Model D

- Each thread independently reads the related model parameters, performs computation, and update the parameters. Unlike A, threads are allowed to read or write the same model parameter in parallel.



Schedule Training Data Partitions to Threads

(only Data Partitions in Computation Model A, C, D;
Data and/or Model Partitions in B)



Intra-node Schedulers

(A) Dynamic Scheduler

- All computation models can use this scheduler.
- All the inputs are submitted to one queue.
- Threads dynamically fetch inputs from the queue.
- The main thread can retrieve the outputs from the output queue.

(B) Static Scheduler

- All computation models can use this scheduler.
- Each thread has its own input queue and output queue.
- Each thread can submit inputs to another thread .
- The main thread can retrieve outputs from each task's output queue.



Harp K-means

Computation
Model B

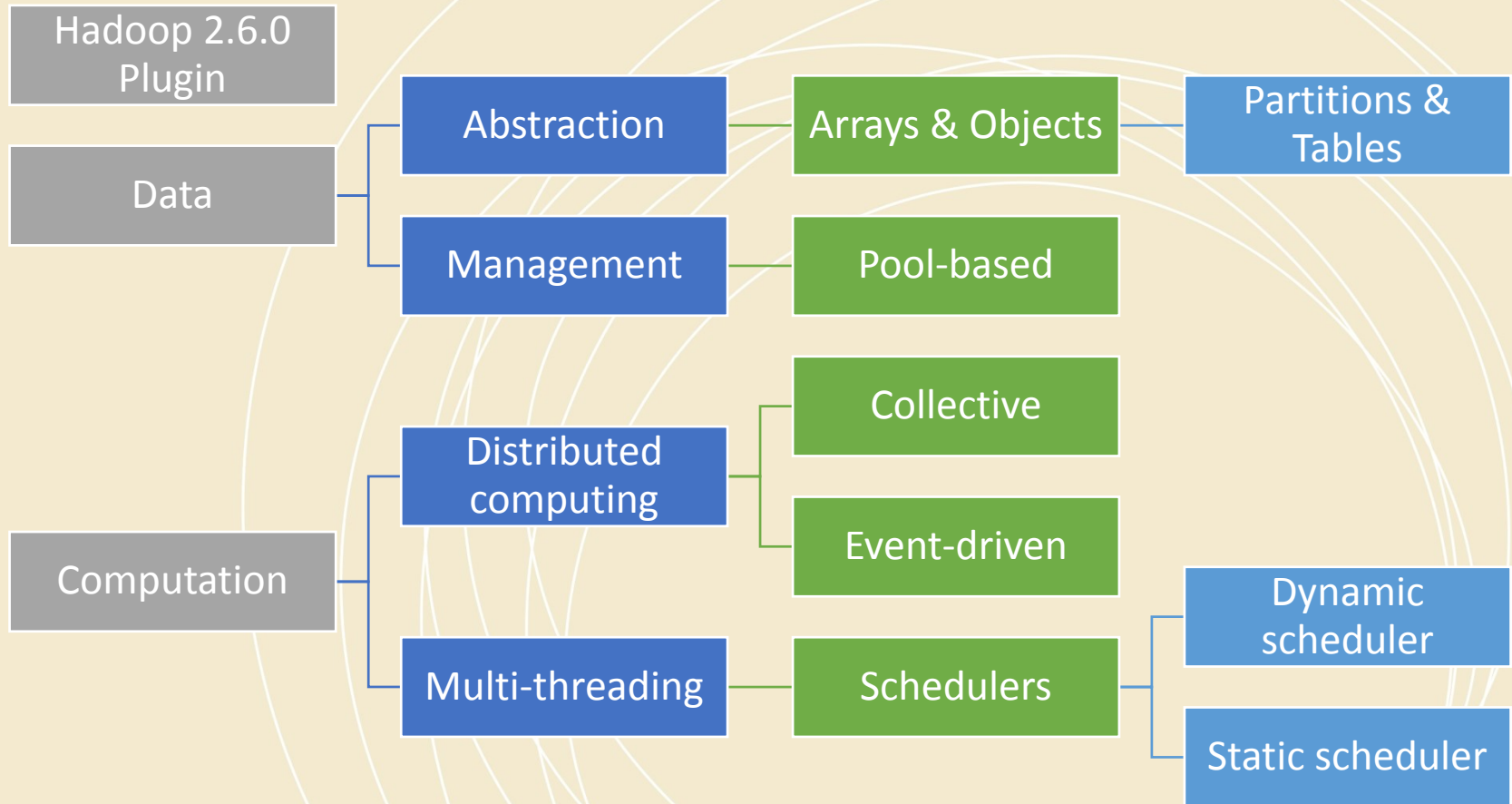
- uses rotation to do partial centroid comparison/update each time

Computation
Model C

- uses allreduce to sum all the local versions of intermediate centroids



Harp Features



Data Types

Arrays & Objects

Primitive Arrays

- ByteArray, ShortArray, IntArray, FloatArray, LongArray, DoubleArray

Serializable Objects

- Writable



Partitions & Tables

Partition

- An array/object with partition ID

Table

- The container to organize partitions

Key-value Table

- Automatic partitioning based on keys



APIs

Scheduler

- DynamicScheduler
- StaticScheduler

Collective

- broadcast
- reduce
- allgather
- allreduce
- regroup
- pull
- push
- rotate

Event Driven

- getEvent
- waitEvent
- sendEvent



mapCollective interface

protected void setup(Context context)

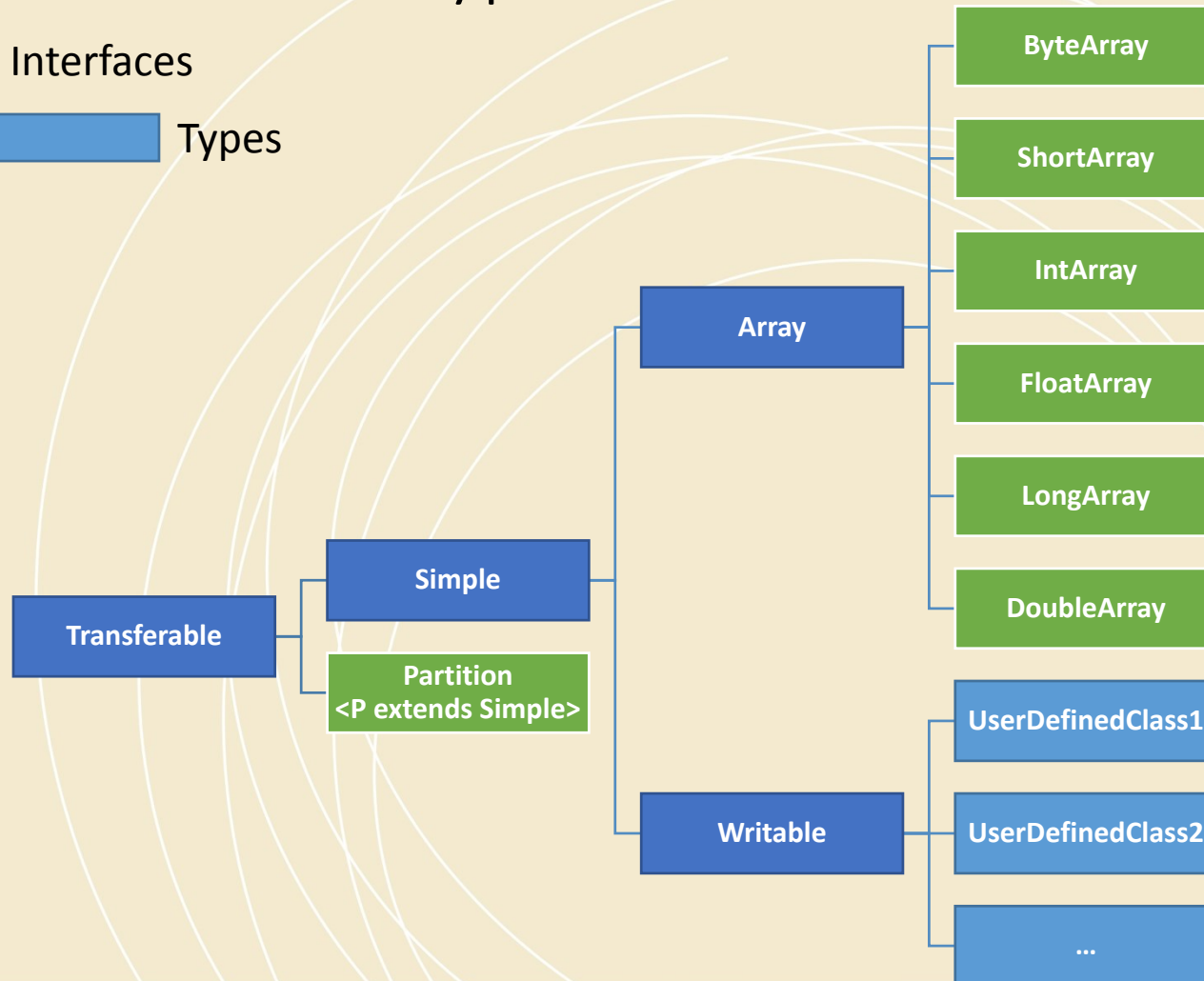
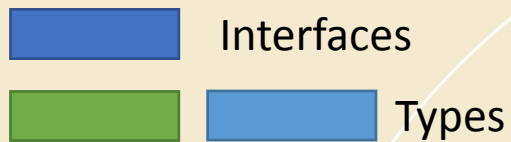
- The interface invoked before running the task, used for fetching job configurations to the Map task.

protected void
mapCollective(KeyValReader reader, Context context)

- The main interface to process key-value pairs
- KeyValReader is used to read all the key-value pairs to the task



Data Interfaces & Types



Array & Writable

Array

- **create(int len, boolean approximate)** - fetch array allocation from the pool. The boolean parameter indicates if the real allocation size can be padded to increase the chance for reuse.
- **release()** – release the array back to the pool
- **free()** – free the allocation to GC
- **get()** – get the array body
- **start()** – get the start index of the array
- **size()** – get the size of array

Writable

- Common methods
 - **create(Class<W> clazz)** – create a object based on the real class
 - **release()** – release the object to the pool
 - **free()** – free the object to GC
- Interfaces should be implemented
 - **getNumWriteBytes()** – calculate the number of bytes to be serialized
 - **write(DataOutput out) & read(DataInput in)** – interfaces for serialization / deserialization
 - **clear()** – clean the fields of the object is before releasing to the pool



An Example of Writable

```
public class Barrier extends Writable {
    private boolean status;

    public Barrier() {
        status = false;
    }
    public void setStatus(boolean st) {
        status = st;
    }

    public boolean getStatus() {
        return status;
    }

    @Override
    public void write(DataOutput out) throws IOException {
        out.writeBoolean(status);
    }

    @Override
    public void read(DataInput in) throws IOException {
        status = in.readBoolean();
    }

    @Override
    public int getNumWriteBytes() {
        return 1;
    }

    @Override
    public void clear() {
    }
}
```



Partition & Table

Partition<P extends Simple> extends Transferable

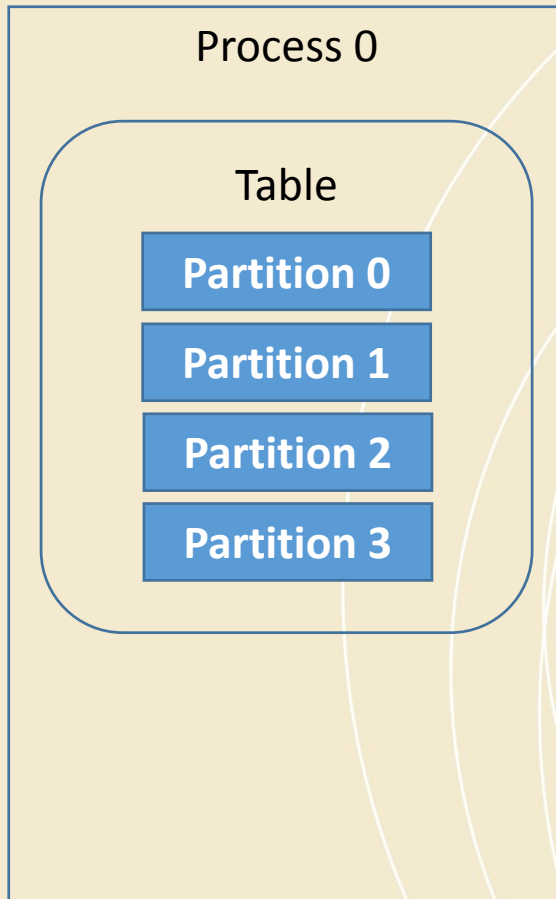
- **Partition(int partitionID, P partition)**
- **id()** - get the partition ID
- **get()** - get the content object of the partition

Table<P extends Simple>

- **Table(int tableID, PartitionCombiner<P> combiner)**
- **tableID** - user defined ID, default ID is allowed
- **combiner** - combiner can combine partitions with the same ID in the table



An Example of Table & Partition



```
Table<DoubleArray> table = new Table<>(0, new DoubleArrPlus());
for (int i = 0; i < numPartitions; i++) {
    DoubleArray array = DoubleArray.create(size, false);
    table.addPartition(new Partition<>(i, array));
}
```

```
public class DoubleArrPlus extends PartitionCombiner<DoubleArray> {
    public PartitionStatus combine(DoubleArray curPar, DoubleArray
newPar) {
        double[] doubles1 = curPar.get(); int size1 = curPar.size();
        double[] doubles2 = newPar.get(); int size2 = newPar.size();
        if (size1 != size2) {
            return PartitionStatus.COMBINE_FAILED;
        }
        for (int i = 0; i < size2; i++) {
            doubles1[i] = doubles1[i] + doubles2[i];
        }
        return PartitionStatus.COMBINED;
    }
}
```



broadcast

```
public <P extends Simple> boolean broadcast(String contextName,  
String operationName, Table<P> table, int bcastWorkerID, boolean  
useMSTBcast)
```

- **contextName** – user defined name to separate operations in different groups
- **operationName** – user defined name to separate operations
- **Table<P> table** – the data structure for broadcasting/receiving
- **bcastWorkerID** – the worker to broadcast the data
- **useMSTBcast** – default broadcast algorithm is pipelining, set this option to true to enable Minimum Spanning Tree (MST) algorithm
- e.g., `broadcast("main", "broadcast-centroids", table, 0, false);`



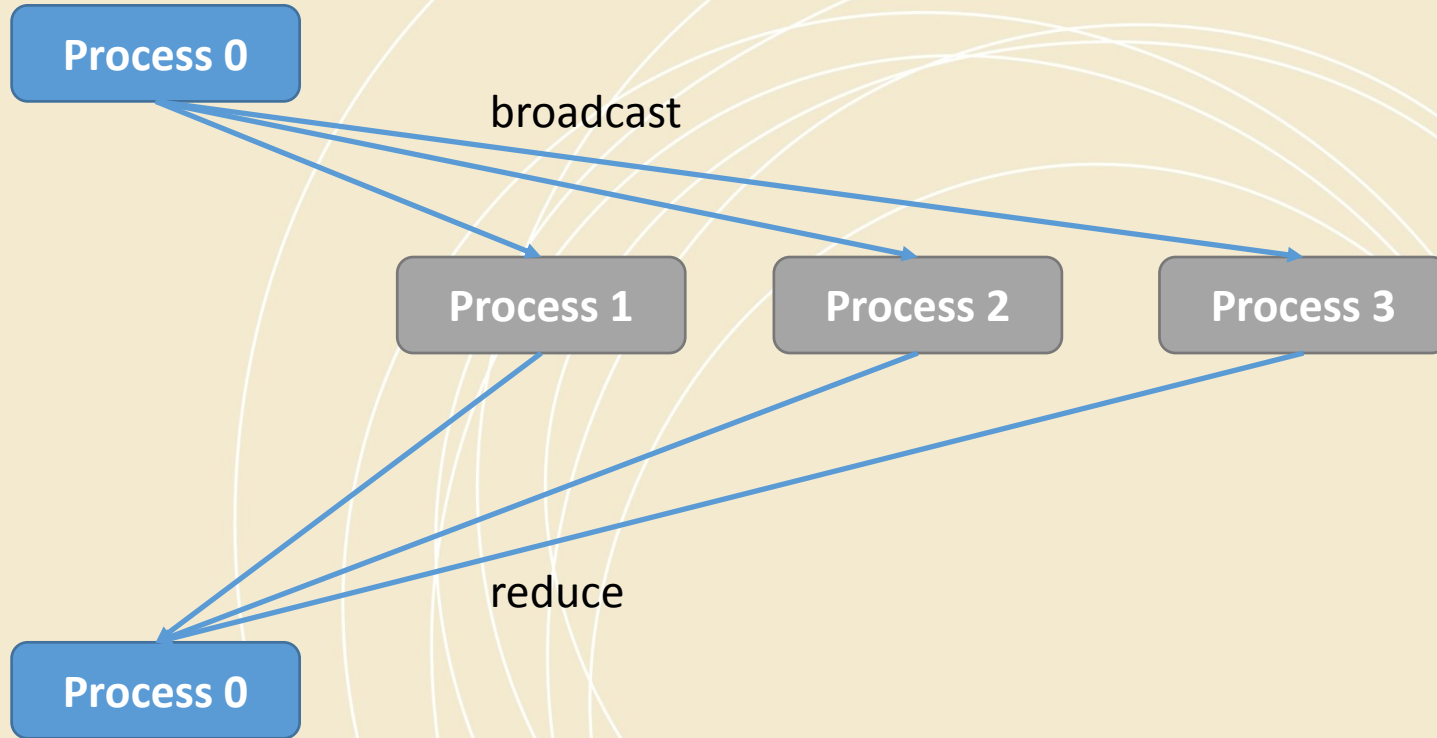
reduce

```
public <P extends Simple> boolean reduce(String contextName,  
String operationName, Table<P> table, int reduceWorkerID)
```

- **contextName** - user defined name to separate operations in groups
- **operationName** - user defined name to separate each operation
- **table** - the data structure for reducing data
- **reduceWorkerID** - the worker ID to receive the reduced data
- e.g., `reduce("main", "reduce-centroids", table, 0);`



Broadcast & Reduce



allreduce

Normally...

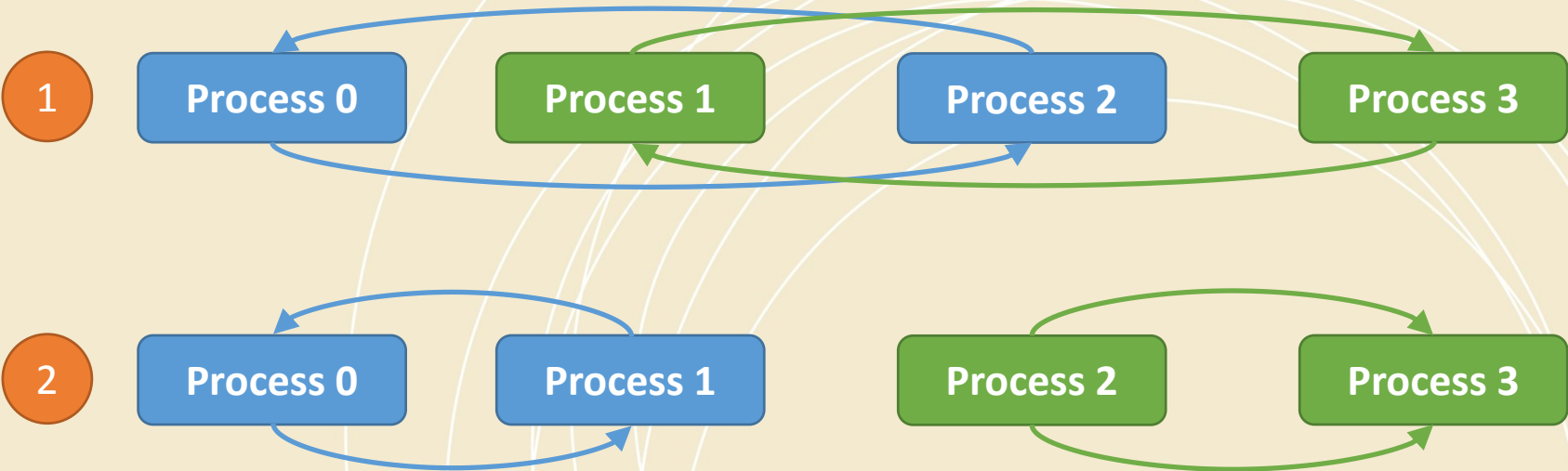
- `public <P extends Simple> boolean allreduce(String contextName, String operationName, Table<P> table)`
- e.g., `allreduce("main", "allreduce", table);`

An alternative way... (if the dataset is large)

- `regroup("main", "regroup", table, new Partitioner(getNumWorkers()));`
- `allgather("main", "allgather", table);`



Allreduce



regroup

```
public <P extends Simple, PT extends Partitioner> boolean  
regroup(String contextName, String operationName, Table<P>  
table, PT partitioner)
```

- **contextName** – user defined name to separate operations in different groups
- **operationName** - user defined name to separate operations
- **table** - the data structure for regrouping data
- **partitioner** – tells which partition to go to which worker for regrouping, e.g. **new Partitioner(numWorkers)**
- e.g., `regroup("main", "regroup", table, new Partitioner(getNumWorkers()));`

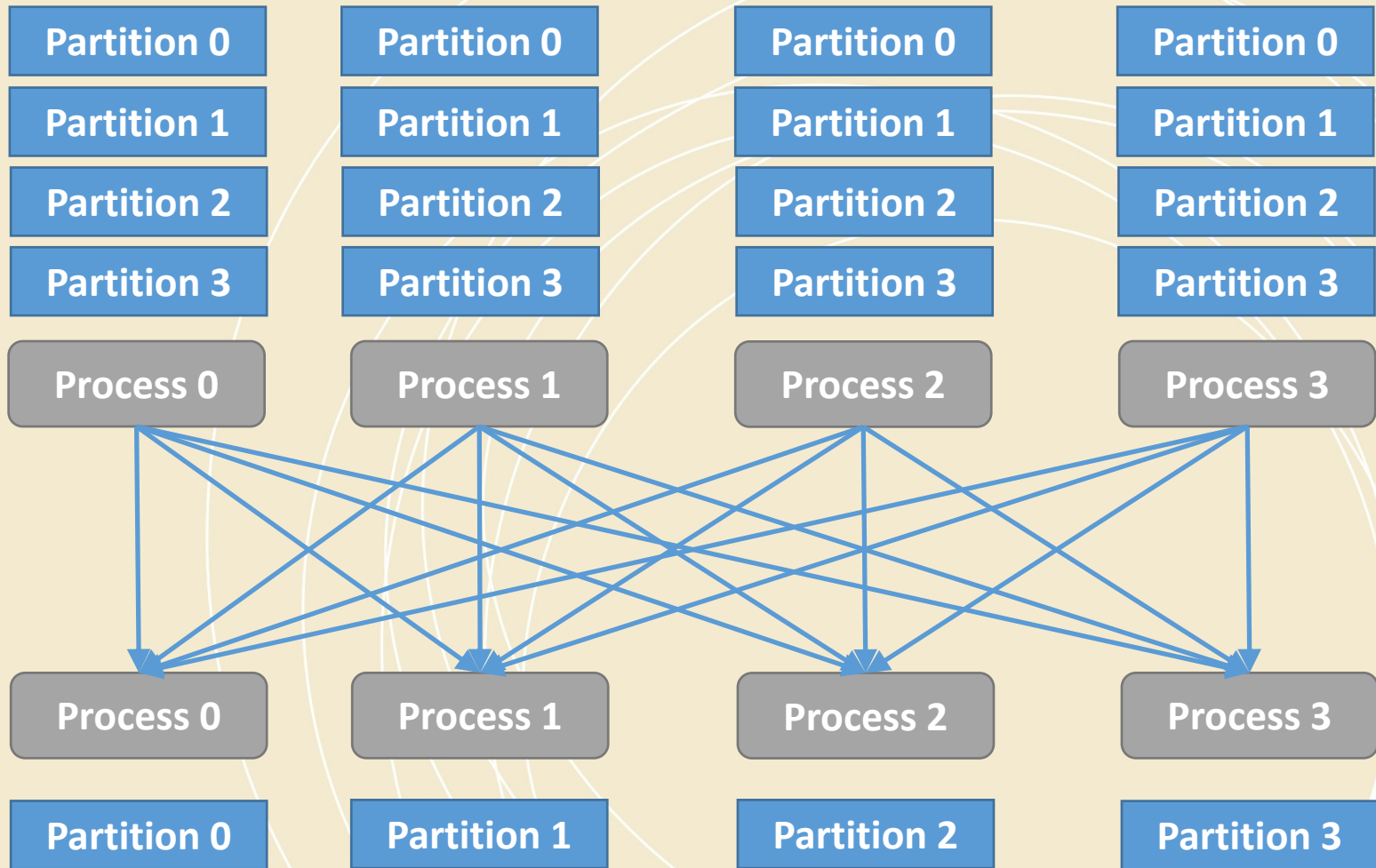


Default Partitioner

```
public class Partitioner {  
  
    private int numWorkers;  
  
    public Partitioner(int numWorkers) {  
        this.numWorkers = numWorkers;  
    }  
  
    public int getWorkerID(int partitionID) {  
        int workerID = partitionID % numWorkers;  
        return workerID;  
    }  
}
```



Regroup



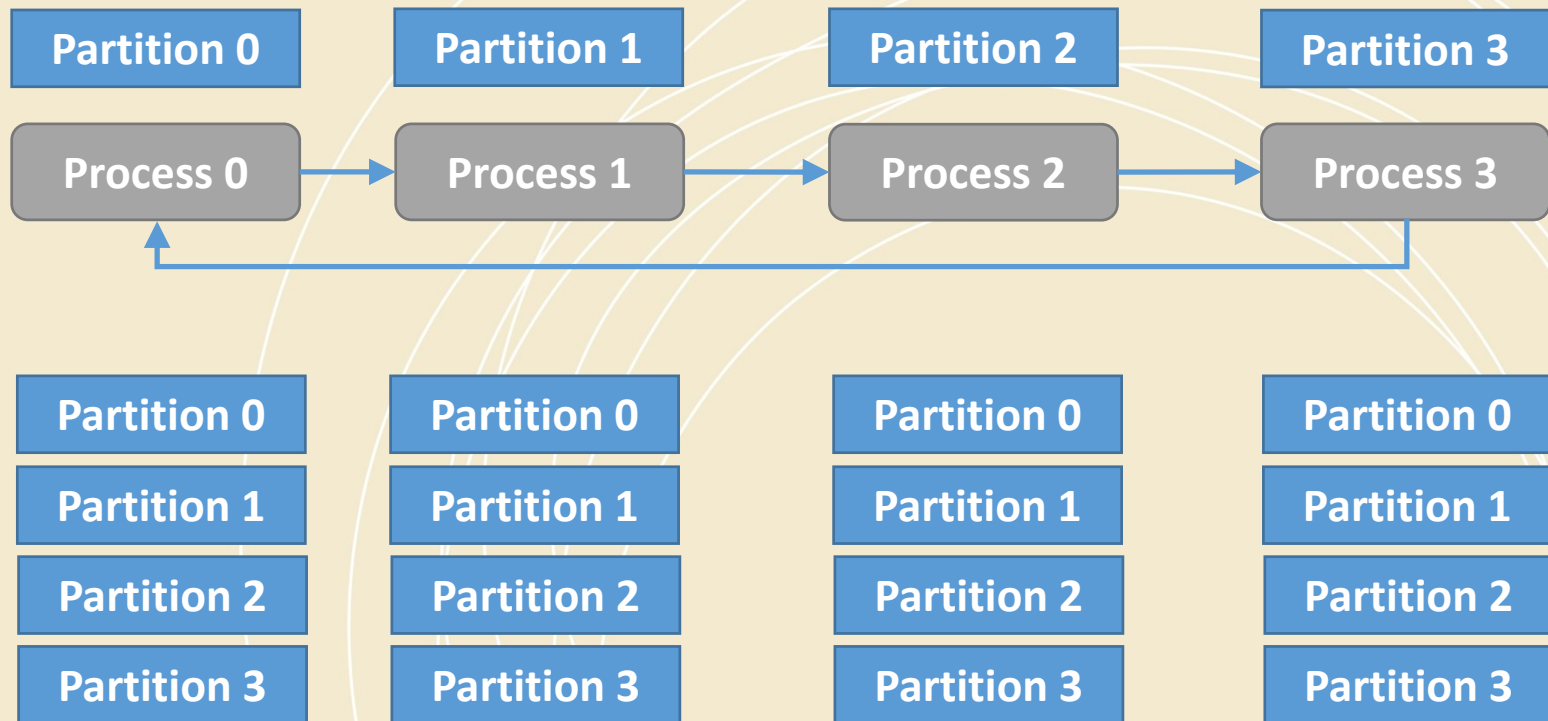
allgather

```
public <P extends Simple> boolean allgather(String contextName,  
String operationName, Table<P> table)
```

- **contextName** – user defined name to separate operations in different groups
- **operationName** - user defined name to separate operations
- **table** - the data structure for allgather data
- e.g., allgather("main", "allgather", table);



Allgather



push & pull

```
public <P extends Simple, PT extends  
Partitioner> boolean push(String  
contextName, String operationName,  
Table<P> localTable, Table<P>  
globalTable, PT partitioner)
```

- Send the partitions from localTable to globalTable based on the partition ID matching
- **localTable** - contains temporary local partitions
- **globalTable** - is viewed as a distributed dataset where each partition ID is unique across processes
- **partitioner** – if some local partitions is not shown in the globalTable, a partitioner can be used to decide where partitions with this partition ID go

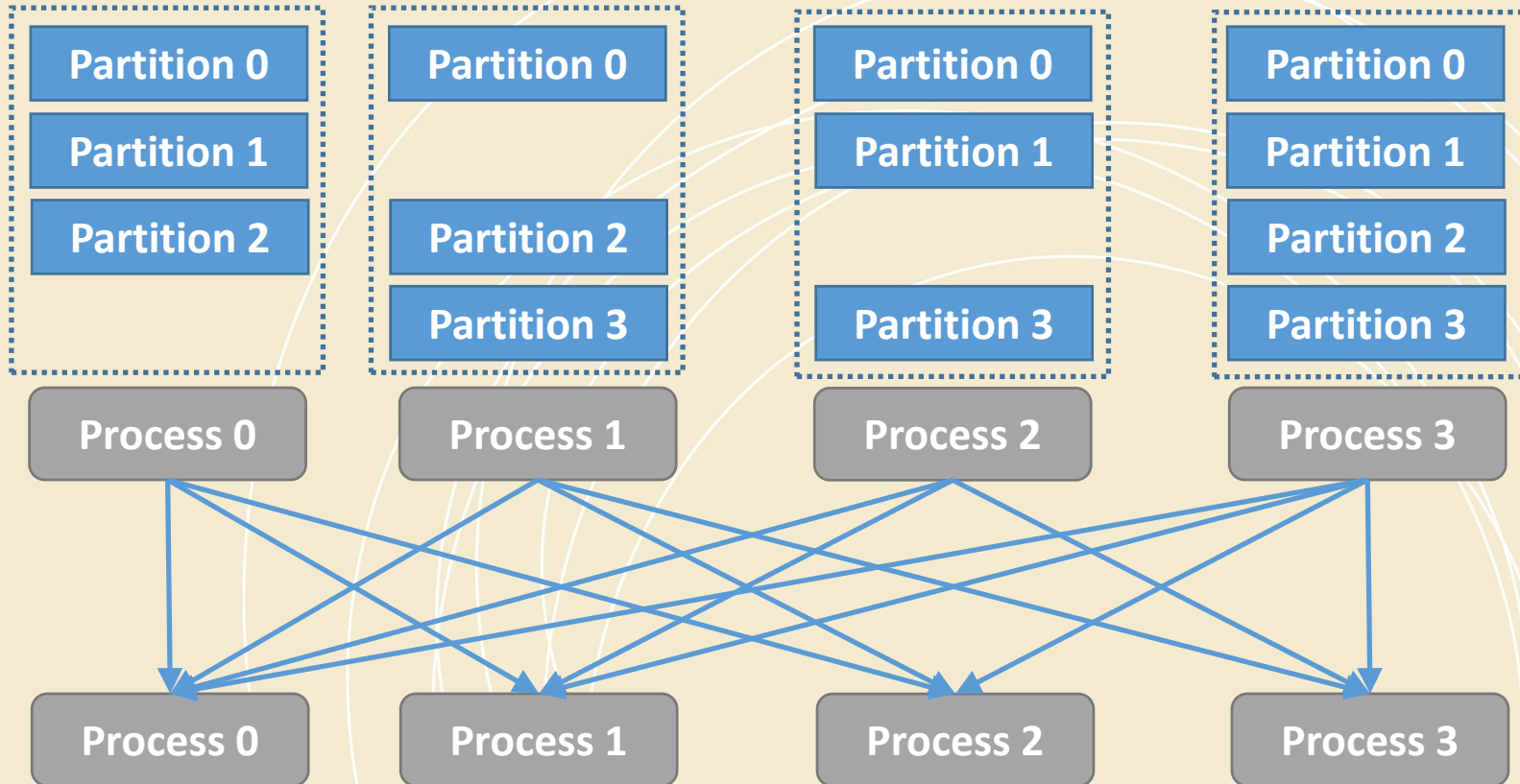
```
public <P extends Simple> boolean  
pull(String contextName, String  
operationName, Table<P> localTable,  
Table<P> globalTable, boolean useBcast)
```

- Retrieve the partitions from globalTable to localTable based on partition ID matching
- **useBcast** – if broadcasting is used when a partition is required to send to all the processes.
- e.g, pull("main", "global-local", localTable, globalTable, true);
- e.g., push("main", "local-global", localTable, globalTable, new Partitioner(getNumWorkers()));



Push

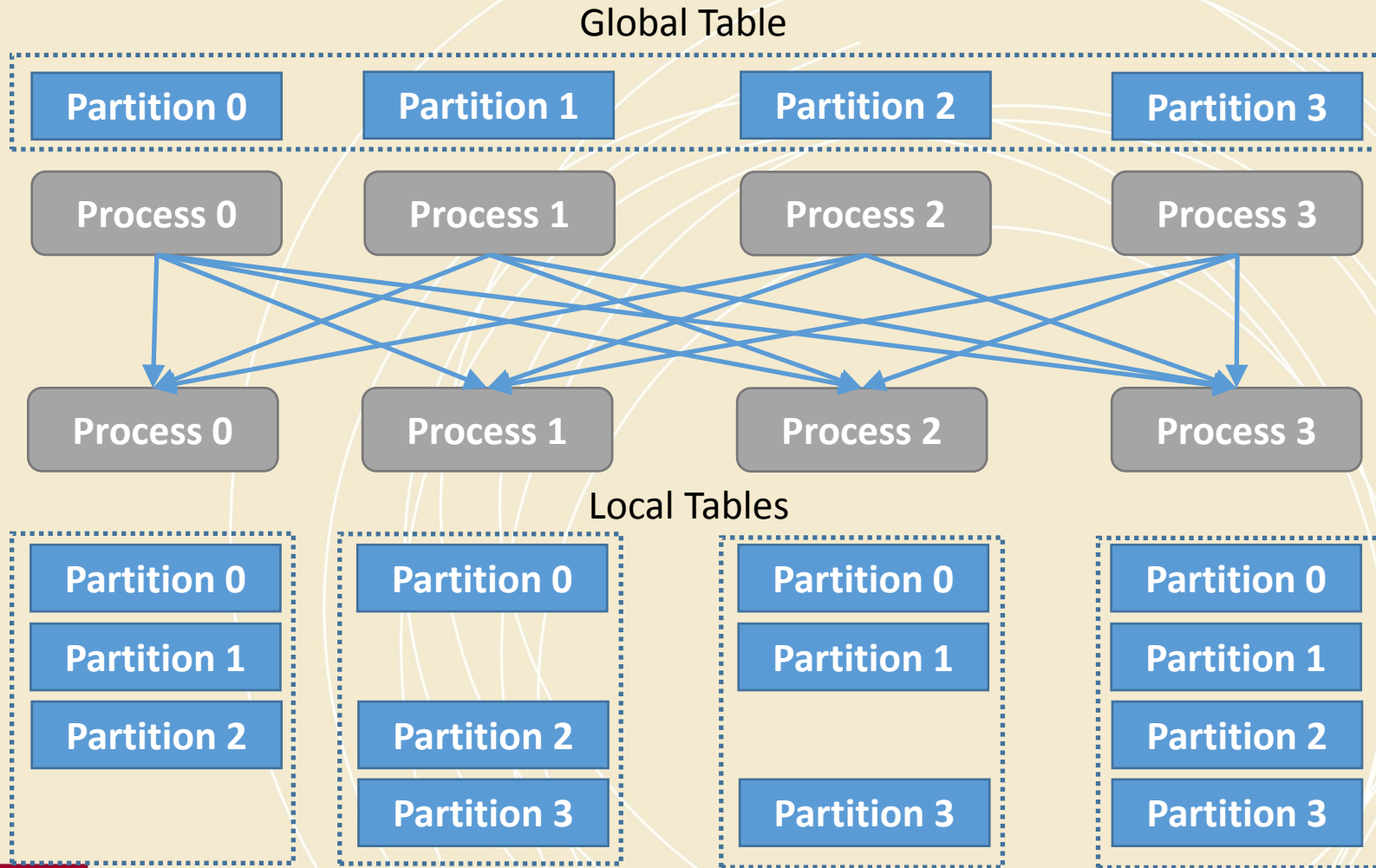
Local Tables



Global Table



Pull



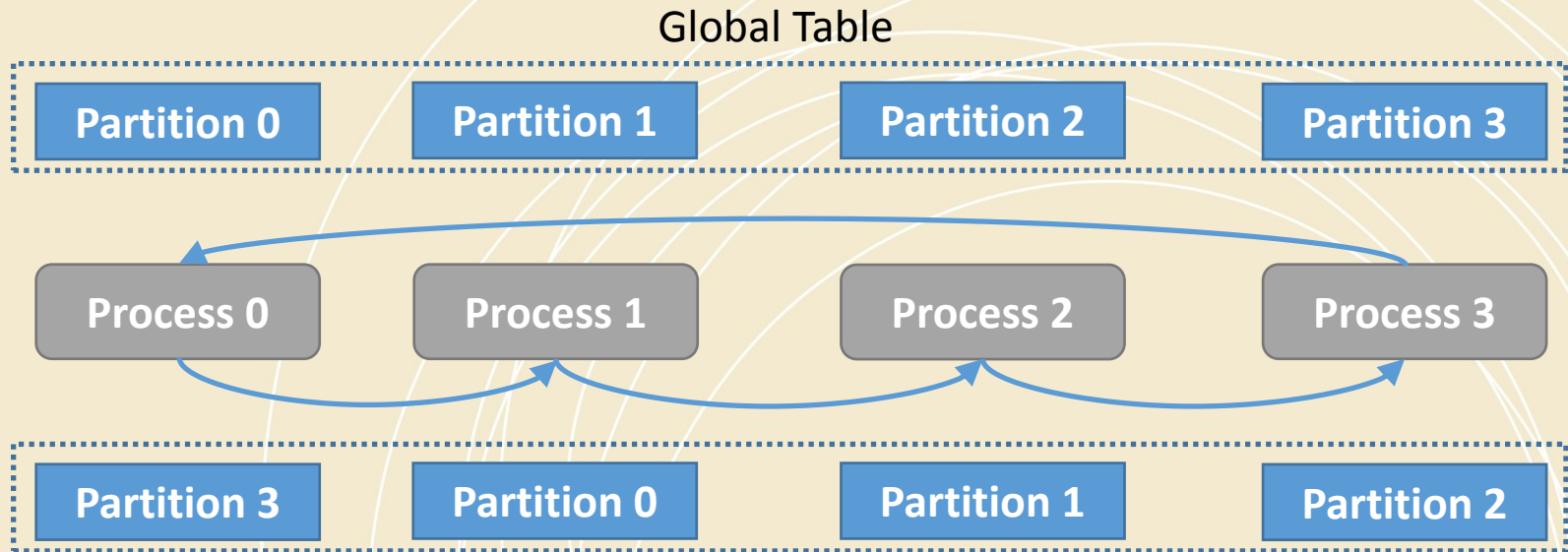
rotate

```
rotate(String contextName, String operationName, Table<P>  
globalTable, Int2IntMap rotateMap)
```

- **contextName** - user defined name to separate operations in different groups
- **operationName** - user defined name to separate operations
- **globalTable** - the data structure for reducing data
- **rotateMap** – the mapping between source worker and target worker
- e.g., rotate("main", "rotate", table, null);



Rotate



DynamicScheduler

public class **DynamicScheduler**<I, O, T extends Task<I, O>> - class declaration

public **DynamicScheduler**(List<T> tasks) - constructor

public synchronized void **submit**(I input) – submit an input

public synchronized void **start**() – start the threads

public synchronized void **pause**() – pause the threads after the current inputs are processed

public synchronized void **pauseNow**() - pause the threads immediately

public synchronized void **cleanInputQueue**() – clean the input queue when the execution is paused or stopped

public synchronized void **stop**() – join the threads

public synchronized boolean **hasOutput**() – check if there is an output in the output queue

public synchronized O **waitForOutput**() – fetch an output from the output queue



StaticScheduler

public class **StaticScheduler**<I, O, T extends Task<I, O>> - class declaration

public **StaticScheduler**(List<T> tasks) - constructor

public synchronized void **submit**(int taskID, I input) – submit an input to a task

public synchronized void **start**() – start threads

public synchronized void **pause**() – pause the threads

public synchronized void **cleanInputQueue**() – clean the input queue

public synchronized void **stop**() – stop the threads

public O **waitForOutput**(int taskID) – wait for the output from a task

public boolean **hasOutput**(int taskID) – check if a task has an output



Adjust Memory Settings in YARN

yarn-site.xml

```
<property>
  <name>yarn.scheduler.minimum-allocation-mb</name>
  <value>512</value>
</property>

<property>
  <name>yarn.scheduler.maximum-allocation-mb</name>
  <value>2048</value>
</property>

<property>
  <name>yarn.nodemanager.resource.memory-mb</name>
  <value>4096</value>
</property>

<property>
  <name>yarn.nodemanager.vmem-check-enabled</name>
  <value>false</value>
</property>
```



Adjust Memory Settings in YARN cont'd

mapred-site.xml

```
<!-- set the Map process memory size -->
<property>
  <name>mapreduce.map.collective.memory.mb</name>
  <value>512</value>
</property>

<property>
  <name>mapreduce.map.collective.java.opts</name>
  <value>-d64 -server -Xmx256m -Xms256m</value>
</property>

<!-- minimize the application master memory size -->
<property>
  <name>yarn.app.mapreduce.am.resource.mb</name>
  <value>512</value>
</property>

<property>
  <name>yarn.app.mapreduce.am.command-opts</name>
  <value>-Xmx256m -Xms256m</value>
</property>
```



Harp MapCollective Job Settings at the Client Side

```
Job job = Job.getInstance(configuration, "kmeans_job");
FileInputFormat.setInputPaths(job, inputDir);
FileOutputFormat.setOutputPath(job, outputDir);
job.setInputFormatClass(MultiFileInputFormat.class);
job.setJarByClass(KMeansLauncher.class);
job.setMapperClass(KMeansCollectiveMapper.class);
org.apache.hadoop.mapred.JobConf jobConf =
    (JobConf) job.getConfiguration();
jobConf.set("mapreduce.framework.name", "map-collective");
jobConf.setNumMapTasks(numMapTasks);
jobConf.setInt("mapreduce.job.max.split.locations", 10000);
job.setNumReduceTasks(0);
```

