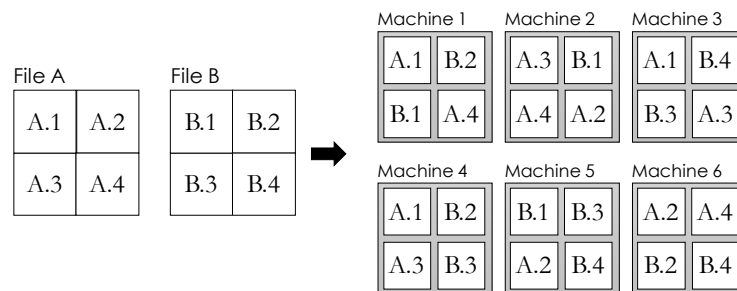


Big Data & Ray

Name:

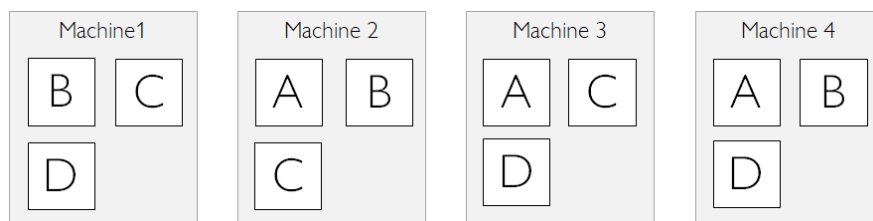
1 Big Data & Distributed File System

1. Consider the following layout of the files A and B onto a distributed file-system of 6 machines.



Assume that all blocks have the same file size and computation takes the same amount of time.

- (a) If we were to lose machines $M1$, $M2$, and $M3$ which of the following file or files would we lose (select all that apply).
 A. File A B. File B **C. We would still be able to load both files.**
- (b) If each of the six machines fail with probability p , what is the probability that we will lose block $B.1$ of file B?
 A. $3p$ **B. p^3** C. $(1 - p)^3$ D. $1 - p^3$



2. The figure above shows four distinct file blocks labeled A, B, C, and D spread across four machines, where each machine holds exactly 3 blocks.
- (a) For the figure above, at most, how many of our machines can fail without any data loss?

2

- (b) Suppose that instead of 4 machines, we have only 3 machines that can store 3 blocks each. Suppose we want to be able to recover our data even if two machines fail. What is the maximum total number of distinct blocks we can store? 3
- (c) Same as part b, but now suppose we only need to be able to recover our data if one machine fails. What is the maximum total number of distinct blocks we can store? 4
3. Suppose we use the map reduce paradigm to compute the total lab scores for each student in DS100. Suppose there are 800 students and 12 labs, and exactly 1 submission per student. Each execution of the map operation is an execution of the autograder, i.e. it will compute the score for a single lab for a single student. The reduce operation computes the total score for each student by adding up all of the lab scores.
- (a) How many key value pairs will be generated in total after all map operations have completed execution? 9600
- (b) How many distinct keys will there be? 800
- (c) How many final key value pairs will remain after all reduce operations have completed? 800
4. As described in class, the traditional data warehouse is a large tabular database that is periodically updated through the ETL process, which combines data from several smaller data sources into a common tabular format. The alternative is a data lake, where data is stored in its original natural form. Which of the following are good reasons to use a data lake approach?
- ☐ A. The data is sensitive, e.g. medical data or government secrets.
 - ☐ B. To maximize compatibility with commercial data analysis and visualization tools.
 - ☒ C. When there is no natural way to store the data in tabular format.
 - ☐ D. To ensure that the data is clean.
5. We want to store a big file on a distributed file system by splitting it into smaller fragments. Assuming the file divides evenly into 800 fragments and we use **4-way replication** answer the following questions.
- (a) If the distributed file system contains 8 separate nodes, how many fragments of the file will be stored on each node?
- ☐ A. 100 ☒ B. 400 ☐ C. 800 ☐ D. 3200
- (b) What is the maximum number of machines that can fail and still guarantee that we can read the entire file.
- ☐ A. 1 ☐ B. 2 ☒ C. 3 ☐ D. 4 ☐ E. 5

6. Which of the following statements are correct?
- (a) In the star schema, the *dimension table* contains the relationships between different facts in the separate *fact tables*.
☐ A. True ☒ B. False
 - (b) Star schemas can help eliminate update errors by reducing duplication of data.
☒ A. True ☐ B. False
 - (c) During the *reduce phase* of MapReduce all the records associated with a given key are sent to the same machine.
☒ A. True ☐ B. False
 - (d) Because files are spread across multiple machines, reading a large file from a distributed file system is usually slower than reading the same large file from a single drive.
☐ A. True ☒ B. False
 - (e) When using MapReduce, we need to have a memory buffer that is big enough to load all the data from disk to memory.
☐ A. True ☒ B. False

2 Distributed/Parallel Computing & Ray

2.1 Primer on ray

2.1.1 Overview

Ray is a distributed execution engine. The same code can be run on a single machine to achieve efficient multiprocessing, and it can be used on a cluster for large computations.

When using Ray, several processes are involved.

1. Multiple **worker** processes execute tasks and store results in object stores. Each worker is a separate process.
2. **Task** is a *stateless* function that can be executed on a remote worker. **Actor** is a *stateful* object that lives in a remote process.
3. One **object store** per node stores immutable objects in shared memory and allows workers to efficiently share objects on the same node with minimal copying and deserialization.

4. One **raylet** per node assigns tasks to workers on the same node.
5. A **driver** is the Python process that the user controls. For example, if the user is running a script or using a Python shell, then the driver is the Python process that runs the script or the shell. A driver is similar to a worker in that it can submit tasks to its raylet and get objects from the object store, but it is different in that the raylet will not assign tasks to the driver to be executed.
6. A **Redis server** maintains much of the system's state. For example, it keeps track of which objects live on which machines and of the task specifications (but not data). It can also be queried directly for debugging purposes.

2.1.2 Asynchronous Computation in Ray

Ray enables arbitrary Python functions to be executed asynchronously. This is done by designating a Python function as a **remote function**.

For example, a normal Python function looks like this.

```
1 def add1(a, b):  
2     return a + b
```

A remote function looks like this.

```
1 @ray.remote  
2 def add2(a, b):  
3     return a + b
```

2.1.3 Remote functions

Whereas calling `add1(1, 2)` returns 3 and causes the Python interpreter to block until the computation has finished, calling `add2.remote(1, 2)` immediately returns an object ID and creates a **task**. The task will be scheduled by the system and executed asynchronously (potentially on a different machine). When the task finishes executing, its return value will be stored in the object store.

```
1 x_id = add2.remote(1, 2)  
2 ray.get(x_id) # 3
```

The following simple example demonstrates how asynchronous tasks can be used to parallelize computation.

```
1 import time  
2  
3 def f1():  
4     time.sleep(1)
```

```

5
6 @ray.remote
7 def f2():
8     time.sleep(1)
9
10 # The following takes ten seconds.
11 [f1() for _ in range(10)]
12
13 # The following takes one second (assuming the system has at least
14   ten CPUs).
15 ray.get([f2.remote() for _ in range(10)])

```

2.2 Parameter Server

Say we are training a *large* logistic regression model that have 500,000,000 rows of data. We use gradient descent to find the optimal parameters. Each worker can compute gradient for the model for only 10,000,000 rows of data, it takes 10ms to compute the gradient. We want to train for 10 iteration to achieve good result.

(a) What's the ideal duration to finish compute gradients for one iteration?

☐ A. 1s ☒ B. 5s ☐ C. 10s

(b) Why can't we achieve the ideal duration? **There will be synchronization and communication cost.**

(c) We use the parameter server pattern, given the code:

```

1 @ray.remote
2 class ParameterServer(object):
3     def __init__(self, learning_rate):
4         self.net = model.SimpleLR(learning_rate=learning_rate)
5
6     def apply_gradients(self, *gradients):
7         self.net.apply_gradients(np.mean(gradients, axis=0))
8         return self.net.variables.get()
9
10    def get_weights(self):
11        return self.net.variables.get()
12
13 @ray.remote
14 class Worker(object):
15     def __init__(self, worker_index, batch_size=50):
16         self.worker_index = worker_index
17         self.batch_size = batch_size
18         self.data = input_data.read_data_sets()
19         self.net = model.SimpleLR()

```

```
20
21     def compute_gradients(self, weights):
22         self.net.variables.set(weights)
23         xs, ys = self.data.train.next_batch(self.batch_size)
24         return self.net.compute_gradients(xs, ys)
```

Fill in the code below, how would we complete the training loop?

```
1 current_weights = randomly_initialized_weight
2 while True:
3     gradients = [worker.___.remote(current_weights)
4                   for worker in workers]
5     current_weights = ps.___.remote(*gradients)
```

```
1 current_weights = randomly_initialized_weight
2 while True:
3     gradients = [worker.compute_gradients.remote(current_weights)
4                   for worker in workers]
5     current_weights = ps.apply_gradients.remote(*gradients)
```
