CS 744: MAPREDUCE

Shivaram Venkataraman Fall 2021

ANNOUNCEMENTS

- Assignment I deliverables
 - Code (comments, formatting)
 - Report
 - Partitioning analysis (graphs, tables, figures etc.)
 - Persistence analysis (graphs, tables, figures etc.)
 - Fault-tolerance analysis (graphs, tables, figures etc.)

• See Piazza for Tutorial



BACKGROUND: PTHREADS

```
void *myThreadFun(void *vargp)
                                    - by thread
{
    sleep(1);
    printf("Hello World\n");
    return NULL;
                                     - Threads there memory
file descriptors
}
int main()
{
    pthread_t thread_id_1, thread_id_2;
    pthread_create(&thread_id_1, NULL, myThreadFun, NULL); ____ Oreating a new pthread_create(&thread_id_2, NULL, myThreadFun, NULL); ______ thread
    pthread join(thread id 1, NULL);
    pthread join(thread id 2, NULL);
    exit(0);
```

holks Conditional Variables

BACKGROUND: MPI machines to land processes

int main(int argc, char** argv) { MPI Init(NULL, NULL);

> // Get the number of processes int world size; MPI Comm size(MPI COMM WORLD, &world_size);

> // Get the rank of the process int world_rank; ____ MPI Comm rank(MPI COMM WORLD, &world rank);

```
// Print off a hello world message
printf("Hello world from rank %d out of %d processors\n",
      world rank, world size);
                                                Hello world
```

```
// Finalize the MPI environment.
MPI Finalize();
```

}

mpirun -n 4 -f host_file ./mpi_hello_world

Hello world

rank 0 rank 1

MOTIVATION

Build Google Web Search

- Crawl documents, build inverted indexes etc.

Need for

- automatic parallelization Automate how many processes and network, disk optimization
- network, disk optimization
- handling of machine failures

OUTLINE

- Programming Model
- Execution Overview
- Fault Tolerance
- Optimizations

GFS **PROGRAMMING MODEL** Data type: Each record is (key, value)

Map function:

$$(K_{in}, V_{in}) \rightarrow \mathsf{list}(K_{inter}, V_{inter})$$

Reduce function:

$$(\mathsf{K}_{\mathsf{inter}},\mathsf{list}(\mathsf{V}_{\mathsf{inter}})) \rightarrow \mathsf{list}(\mathsf{K}_{\mathsf{out}},\mathsf{V}_{\mathsf{out}})$$



EXAMPLE: WORD COUNT

each line is input to map (def mapper(line): for word in line.split(): output(word, 1) (cs, 1) (cours, 1)

(CS, 3)

def reducer(key, values):
 output(key, sum(values))



CS -> key [1,1,1...] -> Values



WORD COUNT EXECUTION: PART2



ASSUMPTIONS



ASSUMPTIONS

- I. Commodity networking, less bisection bandwidth
- 2. Failures are common
- 3. Local storage is cheap
- 4. Replicated FS
- 5. Input is splittable



FAULT RECOVERY Inputs are replicated Tasks are deterministic

If a task crashes:

- Retry on another node
- If the same task repeatedly fails, end the job



FAULT RECOVERY

If a node crashes:

Relaunch its current tasks on other nodes
 M/het ab aut task insute ? File system realized

What about task inputs ? File system replication



FAULT RECOVERY

If a task is going slowly (straggler):

- some tasks make very Now progress er node (even for some data) - Launch second copy of task on another node
- Take the output of whichever finishes first





MORE DESIGN

Master failure

MAPREDUCE: SUMMARY

- Simplify programming on large clusters with frequent failures
- Limited but general functional API
 - Map, Reduce, Sort
 - No other synchronization / communication
- Fault recovery, straggler mitigation through retries

DISCUSSION

https://forms.gle/prXWmR97A3xozHkH9

DISCUSSION

Indexing pipeline where you start with HTML documents. You want to index the documents after removing the most commonly occurring words.

- Compute most common words.
- Remove them and build the index. 2

What are the main shortcomings of using MapReduce to do this?

One MapReduce job to compute most common words Second MR job to build index I can you comb (1) Atore the most common words (2) Read the input data multiple times one MR?





(a) 4 (b) loste very similar- but back up tasks do speed up execution

Input rates are much higher

MapReduce Usage Statistics Over Time

	Aug, '04	Mar, '06	Sep, '07	Sep, '09
Number of jobs	29K	171K	2,217K	3,467K
Average completion time (secs)	634	874	395	475
Machine years used	217	2,002	11,081	25,562
Input data read (TB)	3,288	52,254	403,152	544,130
Intermediate data (TB)	758	6,743	34,774	90,120
Output data written (TB)	193	2,970	14,018	57,520
Average worker machines	157	268	394	488

Jeff Dean, LADIS 2009

NEXT STEPS

- Next lecture: Spark
- Assignment I: Use Piazza!