

Cluster Computing

Weijie Zhao

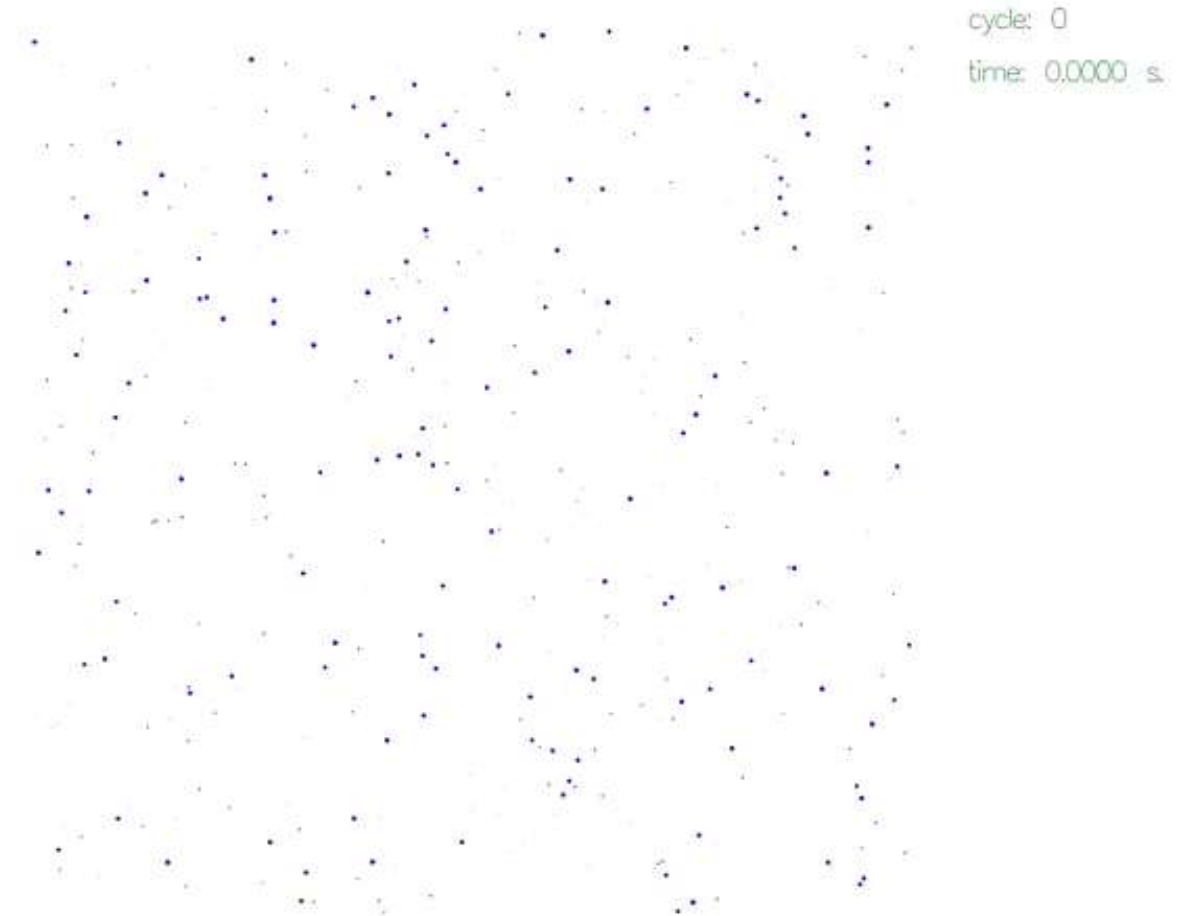
09/21/2023

HW1 Review

- 19/19 submissions
 - 7/19 correct solutions
 - Fastest solution:
 - Quinn Tucker 0.35s (7.89s momentum)
 - Runner-ups:
 - Kevin Penkowski 17.93s gapped evaluation
 - Karamcheti Pritham 30.83s parallel: accuracy, grad, update
 - Solutions no slower than 15.78s will get 15 pts
- xxxtargzlog 8 245.64 [0.01, 0.01, 0.02, 0.04, 0.04, 0.04, 0.07, 120.0, 120.0, 5.41] [8, 9]
- Random seed for generator 22566789
 - All grades will be finalized at the end of 9/28

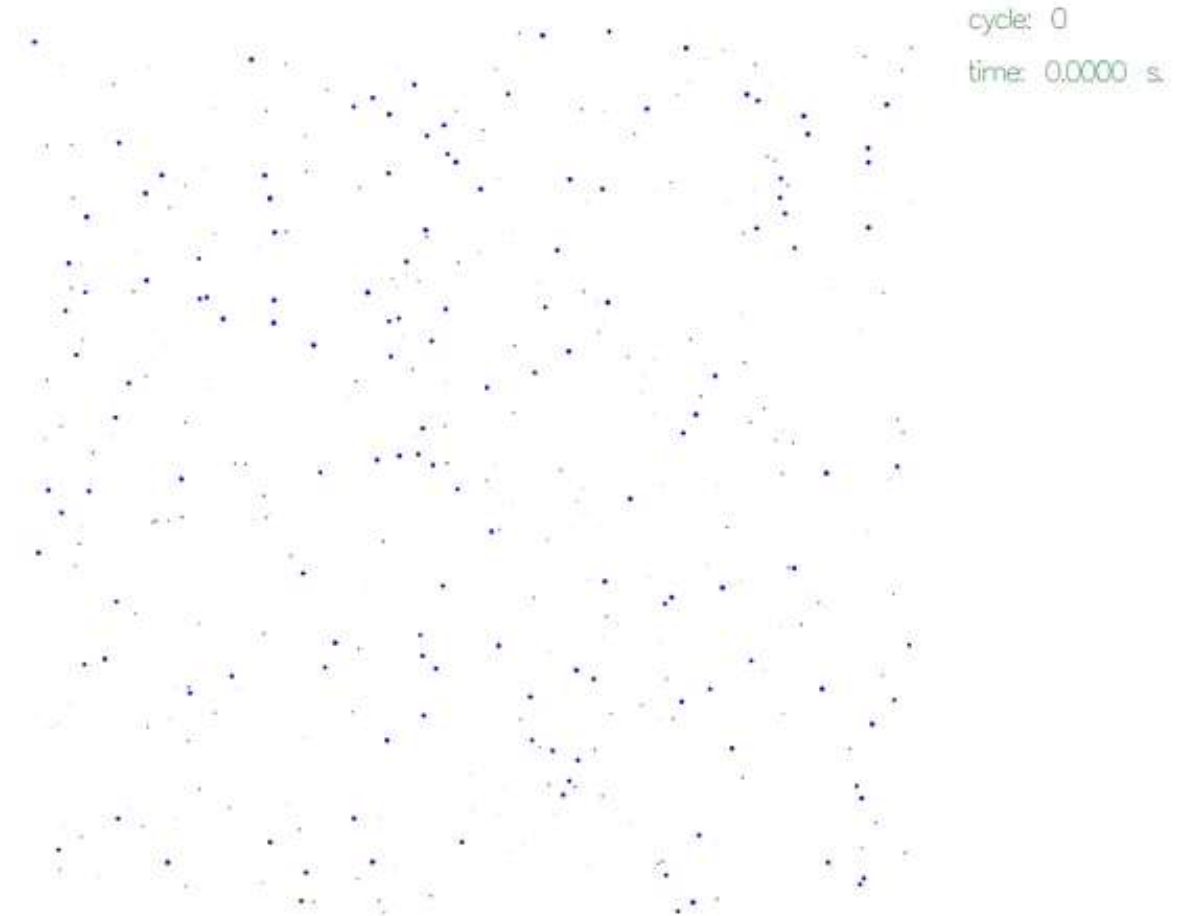
N-Body Problem

- Given N objects
 - Mass
 - Velocity
- Compute the status of each object
- Universal Gravitation
 - $O(N^2)$ forces



N-Body Problem

- Given N objects
 - Mass
 - Velocity
- Compute the status of each object
- Universal Gravitation
 - $O(N^2)$ forces
- We need to scale!
- (Or have a better algorithm)



Cluster Computing

- Putting many (cheap) computers in a cluster
 - The computers in the same cluster do not have to be the same
 - Communication topology
 - All nodes are fully connected
 - Hubs
- Eventually, the communication will be the bottleneck
- For most cases, a network filesystem is employed

Message Passing Interface (MPI)

- Introduced in early 90's
- Each process may have multiple threads
- Each process has its own address space
- Inter-process communication

MPI Example

```
#include <stdio.h>
#include <mpi.h>
int main(int argc, char *argv[])
{
    MPI_Init(&argc, &argv);
    printf("hello world!\n");
    MPI_Finalize();
    return 0;
}
```

MPI Example

```
#include <stdio.h>
#include <mpi.h>
int main(int argc, char *argv[])
{
    int rank, size;
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    printf("hello world from %d of %d!\n", rank, size);
    MPI_Finalize();
    return 0;
}
```


MPI Communications

```
int MPI_Send(  
    void* data,  
    int count,  
    MPI_Datatype datatype,  
    int destination,  
    int tag,  
    MPI_Comm communicator)
```

```
int MPI_Recv(  
    void* data,  
    int count,  
    MPI_Datatype datatype,  
    int source,  
    int tag,  
    MPI_Comm communicator,  
    MPI_Status* status)
```

MPI Communications

```
int MPI_Probe(  
    int source,  
    int tag,  
    MPI_Comm comm,  
    MPI_Status* status)
```

```
int MPI_Get_count(  
    MPI_Status* status,  
    MPI_Datatype datatype,  
    int* count)
```

MPI Communications

```
int MPI_Isend(  
    const void *buf,  
    int count,  
    MPI_Datatype datatype,  
    int dest,  
    int tag,  
    MPI_Comm comm,  
    MPI_Request *request)
```

MPI Communications

```
int MPI_Wait(  
    MPI_Request *request,  
    MPI_Status *status)
```

```
int MPI_Test(  
    MPI_Request *request,  
    int *flag,  
    MPI_Status *status)
```

Communicator

```
int MPI_Comm_split(  
    MPI_Comm comm,  
    int color,  
    int key,  
    MPI_Comm * newcomm)
```

```
int MPI_Comm_free(MPI_Comm *comm)
```

Compilation and Execution

- MPICH, OpenMPI
- mpicc, mpiCC, mpic++
- mpiexec, mpirun
- mpiexec -np 4 ./a.out
- mpiexec --showme

- SLURM
 - sbatch
 - srun

MPI Configuration

- For each node, create a user that can ssh to all other nodes
- Install MPICH/OpenMPI
- `mpirun -np 4 --hostfile myhost_file ./a.out`
 - `node1 slots=2 max_slots=10`
 - `node2 slots=2 max_slots=10`
- `mpirun -np 4 --hostfile myhost_file --byslot ./a.out`
- `mpirun -np 4 --hostfile myhost_file --bynode ./a.out`

MPI Collective Communications

- MPI_Barrier(MPI_Comm communicator)
- MPI_Bcast(void* data, int count, MPI_Datatype datatype, int root, MPI_Comm communicator)
- MPI_Reduce(const void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)
 - MPI_MIN, MPI_MAX, MPI_MINLOC, MPI_MAXLOC, MPI_BOR, MPI_BXOR, MPI_LOR, MPI_LXOR, MPI_BAND, MPI_LAND, MPI_SUM and MPI_PROD
- MPI_Allreduce(const void* send_buffer, void* receive_buffer, int count, MPI_Datatype datatype, MPI_Op operation, MPI_Comm communicator)

Cluster Computing

- MPI
 - Inter-node communication
 - High-performance computing
- Node failure
 - Broken hardware
 - Software bugs
 - Insufficient resources
- Node failure happens commonly for clusters with 1,000+ nodes
 - $(1 - p)^{1000}$

Cluster Computing

- MPI
 - Inter-node communication
 - High-performance computing
- Node failure
 - Broken hardware
 - Software bugs
 - Insufficient resources
- Node failure happens commonly for clusters with 1,000+ nodes
 - $(1 - p)^{1000}$

We need a system to
handle these failures!

Distributed File System

- Decouple data and computing resources
- Replication to take care of node/disk failures
- HDFS
 - Name node
 - Data node

Common Data Analysis Tasks

- Given a large data, find some statistics
- Given a page view log, find the number of users
- Given a page view log, find the number of users group by browser
- Given a page view log, find the number of users from NY state group by browser

Map and Reduce

- PageRank
- $PR(x) = \sum_{y \text{ links to } (x)} (PR(y) / \text{out_degree}(y))$
- Iterative disk I/O

Spark

- Spark
- Resilient Distributed Dataset (RDD)
 - Immutable
 - Transformations
 - map
 - filter
 - reduceByKey
 - join
 - ...
 - Actions
 - count
 - collect
 - ...