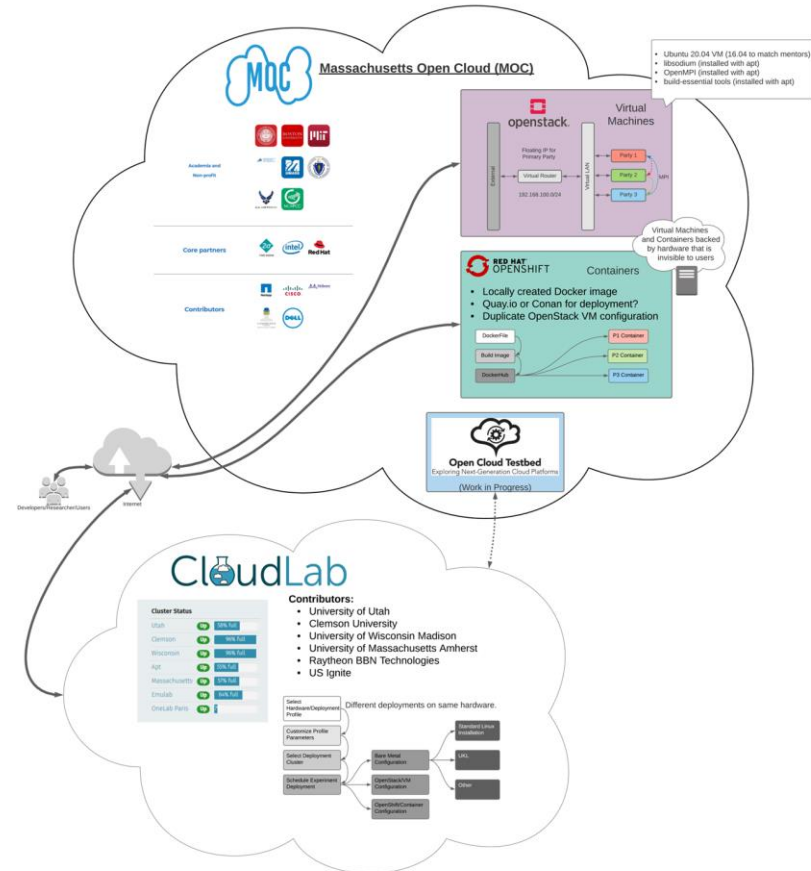


Secure Multiparty Computation Final Report

Developer | Hasnain Abdur Rehman | hasnain@bu.edu
Developer | Pierre-François Wolfe | pwolfe@bu.edu
Developer | Samyak Jain | samyakj@bu.edu
Developer | Suli Hu | sulihu@bu.edu
Developer | Yufeng Lin | yflin@bu.edu
Mentor/Client | John Liagouris | liagos@bu.edu
Mentor/Client | Vasiliki Kalavri | vkalavri@bu.edu
Subject-Matter Expert | Mayank Varia | varia@bu.edu

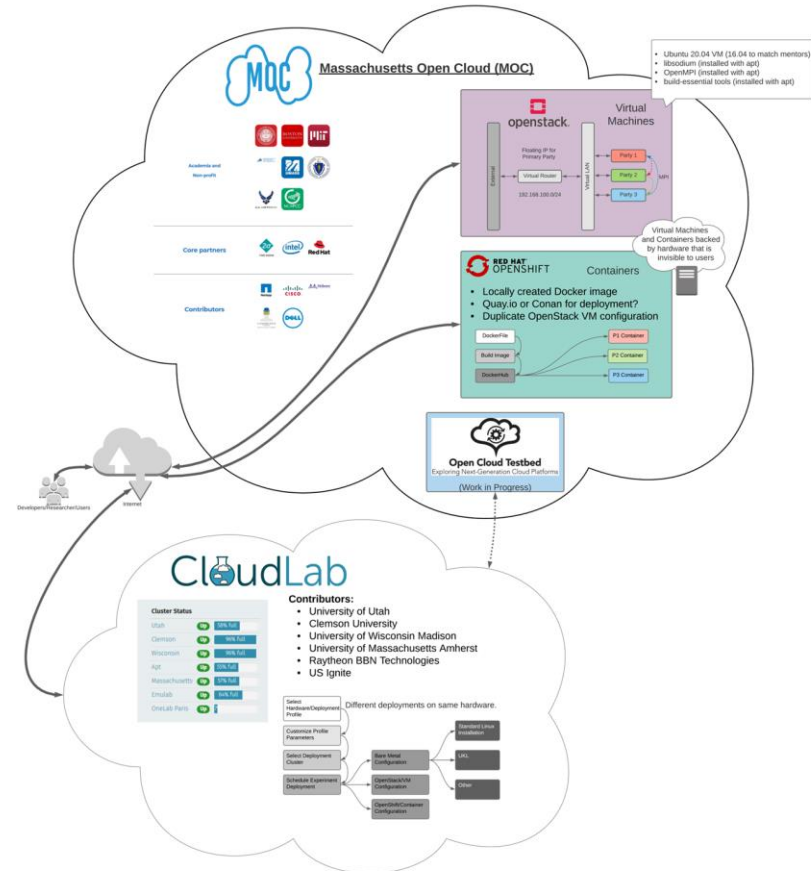
Presentation Outline

- Overview
- Deployment Environments
- Testing Instrumentation
- Data Analysis
- Automation/Replicability
- Conclusions & Future Work



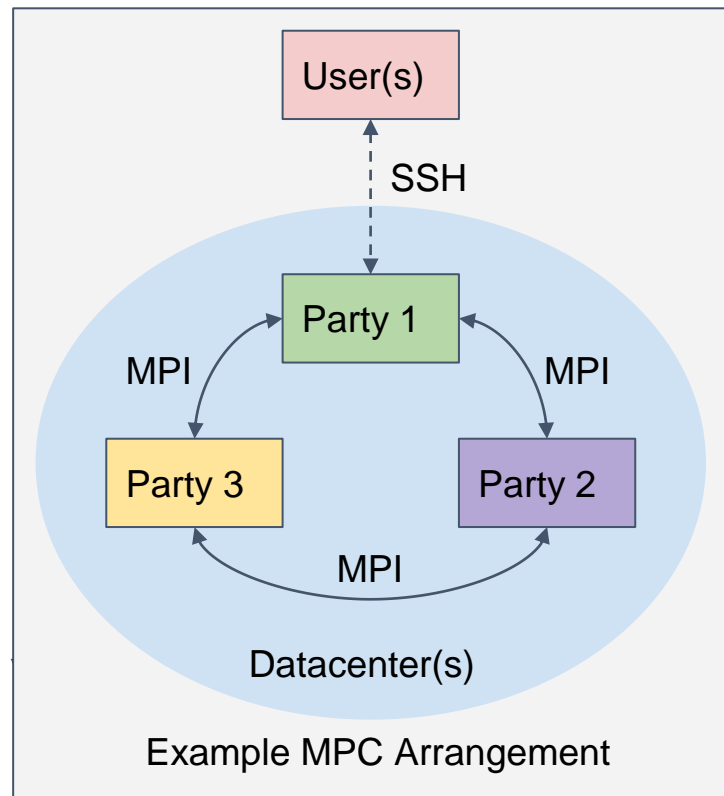
Presentation Outline

- Overview
- Deployment Environments
- Testing Instrumentation
- Data Analysis
- Automation/Replicability
- Conclusions & Future Work




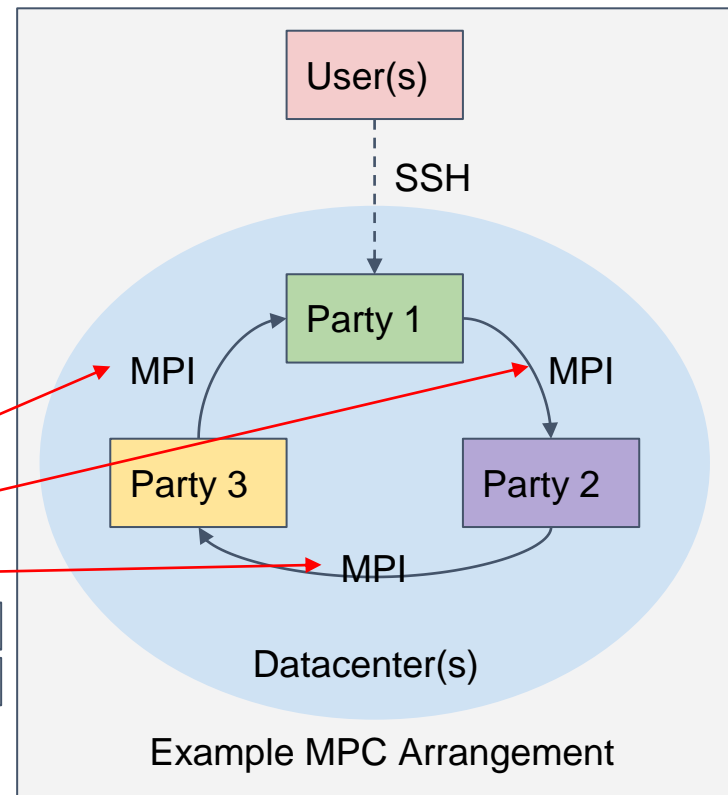
Motivation

- MPC enables...
 - Shared Computation on Private Data
 - Protects the Privacy of Data
 - Mutually Agreed Computation
- Our mentors...
 - Use three party Secret Sharing MPC
 - Perform Database Queries with MPC
 - Keep all operations secure (under MPC) insecure steps



MPC Data Exchange Primitives

- Secret Sharing MPC
 - Latency Bounded
 - Ring communication pattern
 - Each party passes data clockwise
- MPI Communication
 - Sync/Async?
 - Serial/Batch? 
 - Quantity of messages?
 - Other?...



Project Goals

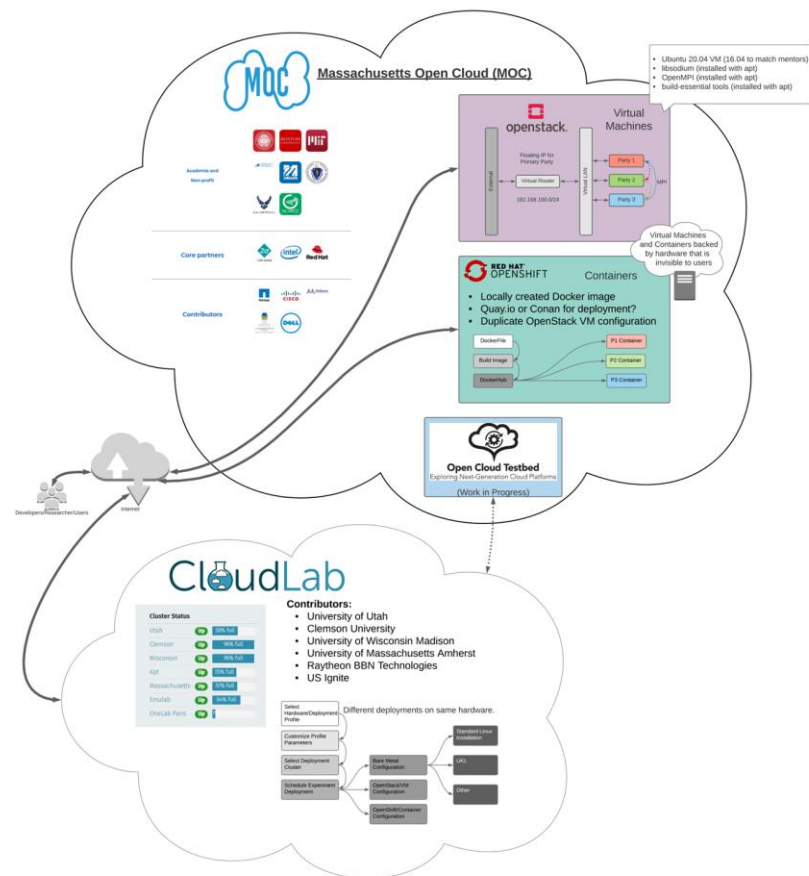
- Deployment
 - Run MPC library in multiple environments
- Testing
 - Create test scenarios conducive to analysis
 - Improve instrumentation of MPC experiments
 - Benchmark the networking layer
 - Analyze and gain initial insights
- Automation
 - Make ongoing testing easy and repeatable

Accomplishments

- Deploy to:
 - CloudLab bare-metal (various hardware, clusters, topologies)
 - MOC OpenStack (VMs)
 - MOC OpenShift (containers)
- Testing (Profiling)
 - Improve exp-exchange testbench (in source code)
 - Collect benchmark and trace data with Score-P
 - Inspect with CUBE GUI (and others)
- Automation
 - Manual config → Shell scripts → Ansible

Presentation Outline

- Overview
- **Deployment Environments**
- Testing Instrumentation
- Data Analysis
- Automation/Replicability
- Conclusions & Future Work



OpenStack

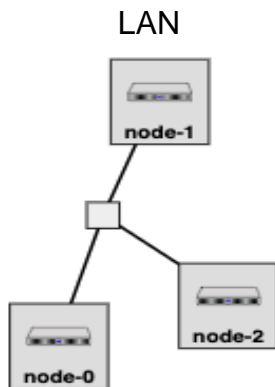
- Single VM to determine setup manually
- Some shell scripts to help automate setup
- One reference VM (semi-automated)
 - Snapshot and duplicate for other parties
- Troubleshoot SSH communication between VMs

CloudLab

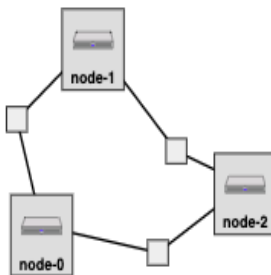
- Initial test with predefined scripts
 - Testing with the provided openstack image
- Custom scripts
 - What went wrong, then what did we learn?
- Compare some different topologies
 - LAN vs RING
 - Single vs multi cluster
 - Keep a local or VM result for reference?

Cloudlab Topologies

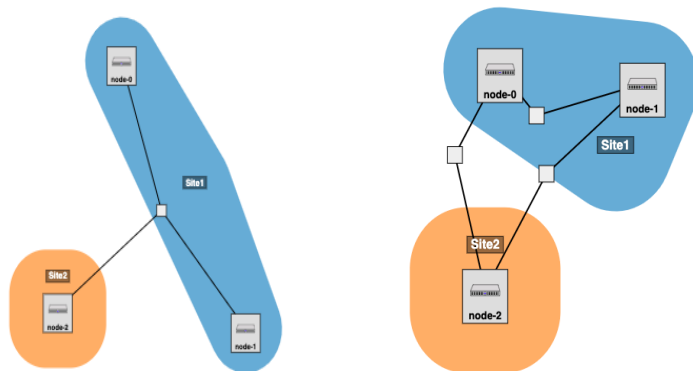
Single Cluster



Ring



Multi Cluster



```

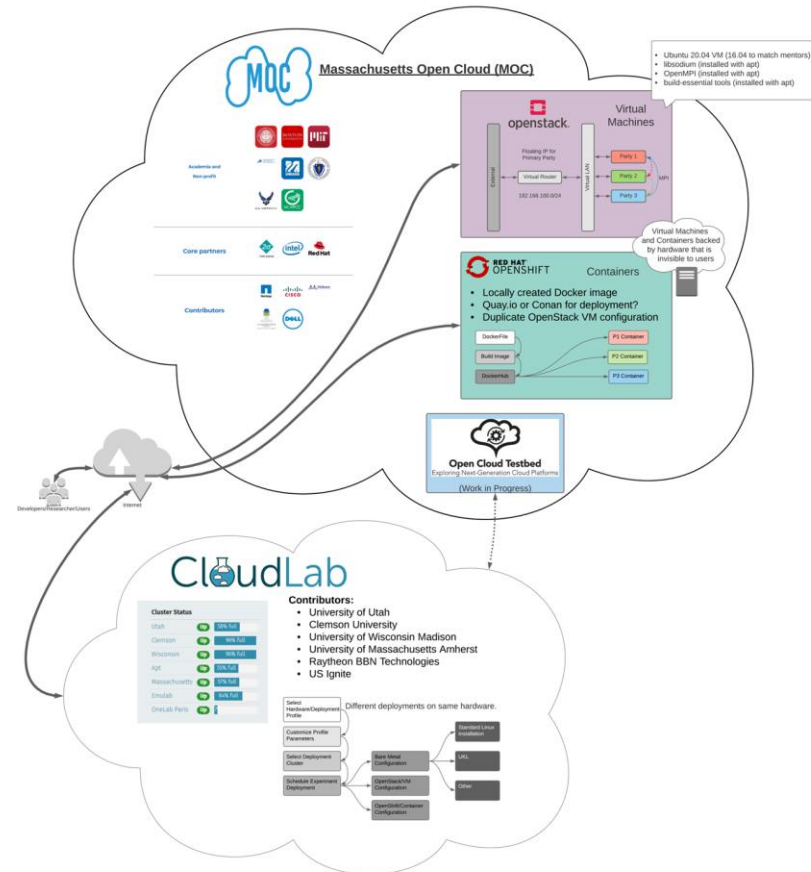
1  ""ubuntu baremetal multi-site LAN""
2
3  # Import the Portal object.
4  import geni.portal as portal
5  # Import the ProtoGENI library.
6  import geni.rspec.pg as pg
7  # Import the Emulab specific extensions.
8  import geni.rspec.emulab as emulab
9
10 pc = portal.Context()
11
12 pc.defineParameter("node_type_1", "Hardware Type for Site 1",
13     portal.ParameterType.NODETYPE, "any")
14 pc.defineParameter("node_type_2", "Hardware Type for Site 2",
15     portal.ParameterType.NODETYPE, "any")
16 pc.defineParameter("node_count", "Number of Machines",
17     portal.ParameterType.INTEGER, 3)
18
19 params = pc.bindParameters()
20 request = portal.context.makeRequestRSpec()
21
22 node = []
23
24 # Create selected number of nodes
25 for i in range(params.node_count):
26     node.append(request.RawPC('node-%d' % i))
27     node[i].disk_image = 'urn:publicid:IDN+emulab.net+image+emulab-ops:UBUNTU16-64-STD'
28     if (i < params.node_count - 1): #Condition can be changed based on requirement
29         node[i].Site("Site1")
30         node[i].hardware_type = params.node_type_1
31     else:
32         node[i].Site("Site2")
33         node[i].hardware_type = params.node_type_2
34
35 # Create a LAN for all the connections
36 lan = request.LAN("lan")
37 lan.bandwidth = 100000
38
39 # Create a link between each of the nodes to make a ring
40 for i in range(params.node_count):
41     iface = node[i].addInterface("eth1")
42     iface.addAddress(pg.IPv4Address("192.168.1."+str(i+1), "255.255.255.0"))
43     lan.addInterface(iface)
44
45 # Print the generated rspec
46 pc.printRequestRSpec(request)
  
```

Containers

- Initial attempt: Using a C/C++ package manager to build the app, and containerize using OpenShift s2i tool .
- Docker deployment using Dockerfiles + Docker Compose.
 - Set up SSH
 - Install and test MPI + MPC codebase.
- OpenShift
 - Unexpected behavior, unlike Docker containers
 - SSH issue
 - Permissions issues: MPI failing to launch required dependencies, create new files.

Presentation Outline

- Overview
- Deployment Environments
- **Testing Instrumentation**
- Data Analysis
- Automation/Replicability
- Conclusions & Future Work



Instrumenting the Exchange Primitive Test

Before

```
11 int main(int argc, char** argv) {
12
13     if (argc < 2) {
14         printf("\nUsage: %s [INPUT_SIZE]\n\n", argv[0]);
15         return -1;
16     }
17 }
```

```
40 /* =====
41 | 1. Measure exchange_array
42 | ===== */
43 // start timer
44 gettimeofday(&begin, 0);
45
46 exchange_shares_array(r1s1, r1s2, ROWS);
47
48 // stop timer
49 gettimeofday(&end, 0);
50 seconds = end.tv_sec - begin.tv_sec;
51 micro = end.tv_usec - begin.tv_usec;
52 elapsed = seconds + micro*1e-6;
53
54 if (rank == 0) {
55     printf("BATCHED\t%d\t%.3f\n", ROWS, elapsed);
56 }
```

Input Arguments

Experiment Timing

```
167 if (rank == 0) {
168     // Output collected data to csv file
169     time_t rawtime;
170     struct tm * timeinfo;
171     char fname[36];
172
173     // YYYYMMDD_hhmmss_exp-exchange_log.csv
174     time(&rawtime);
175     timeinfo = localtime(&rawtime);
176     strftime(fname, 36, "%Y%m%d_%H%M%S_exp-exchange_log.csv", timeinfo);
177     FILE * logfile;
178     logfile = fopen (fname, "w+");
179
180     fprintf(logfile, "ROWS,GENSHR,SEEDS,BATCHED,SYNC,ASYNC,\n");
181     for (int nstep = 0; nstep < NSTEP; nstep++){
182         if (nstep == 0)
183             CURR_ROW = ROWS;
184         else
185             CURR_ROW += STEP;
186         for (int iter = 0; iter < ITER; iter++){
187             fprintf(logfile, "%d,", CURR_ROW);
188             for (int meas = 0; meas < MEAS; meas++){
189                 fprintf(logfile, "%ld.%09ld,", (long long)time_stamp[nstep][iter][meas].tv_sec, time_stamp[nstep][iter][meas].tv_nsec);
190             }
191             fprintf(logfile, "\n");
192         }
193     }
194     fclose(logfile);
195     // Store system details to another file
196 }
```

After

```
29 int main(int argc, char** argv) {
30
31     if (argc < 5) {
32         printf("\nUsage: %s [INPUT_SIZE_START] [STEP_SIZE] [NUM_STEPS] [NUM_ITER]\n\n", argv[0]);
33         return -1;
34     }
35 }
```

```
102
103 /* =====
104 | 1. Measure exchange_array
105 | ===== */
106 // start timer
107 // CLOCK_PROCESS_CPUTIME_ID or CLOCK_MONOTONIC perhaps (check man clock_gettime)
108 clock_gettime(CLOCK_MONOTONIC, &begin);
109
110 exchange_shares_array(r1s1, r1s2, CURR_ROW);
111
112 // stop timer
113 clock_gettime(CLOCK_MONOTONIC, &end);
114 elapsed = diff(begin, end);
115 time_stamp[nstep][iter][2] = elapsed;
116
117 if (rank == 0) {
118     printf("BATCHED\t%d\t%ld.%09ld\n", CURR_ROW, (long long)elapsed.tv_sec, elapsed.tv_nsec);
119 }
```

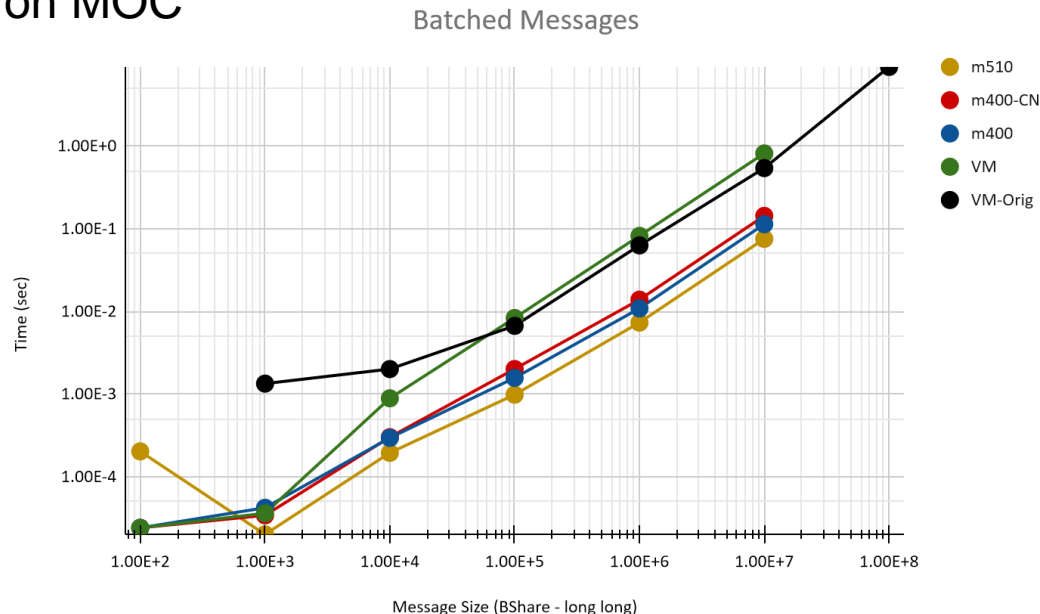
Output Format

Instrumenting the Exchange Primitive Test

- Before:
 - 3 measurements
 - gettimeofday
 - Call multiple times
 - For different sizes
 - For multiple runs
 - txt output & python import
 - No benchmark/trace collection
- After:
 - 5 measurements
 - clock_gettime
 - Call once
 - For range of sizes
 - For multiple runs
 - csv output (easier import)
 - Benchmark/trace collection with Score-P

Collected data from exp-exchange

- Overlaid batched runs from different environments
- Compare data collected from
 - Openstack on MOC
 - Cloudlab



Benchmark/Profiling



- Scorep:
 - Benchmarking/Profiling framework
 - Allows for data collection from MPI (and other sources)
 - Modify Makefile in order to build with Score-P
 - Setup overview:
 - Installed PPA package on ubuntu system
 - Install additional package (sudo apt-get install libz-dev)
 - Specify “scorep” in makefile, update all binary executable on all vms
- Analysis
 - CUBE GUI: Inspecting *.cubex benchmark file

Running Exp-exchange with Score-P

```
ylm@cc-mpc-main:~$ ls scorep-20201207_1014_3385102085294411/
MANIFEST.md  profile.cubex  scorep.cfg
```

- Execute new binary file → Scorep folder including *.cubex file to evaluate → analysis of data
- Textual output supported when no Gui tool available
- Use cube to analyze data in *.cubex (details to be continued)

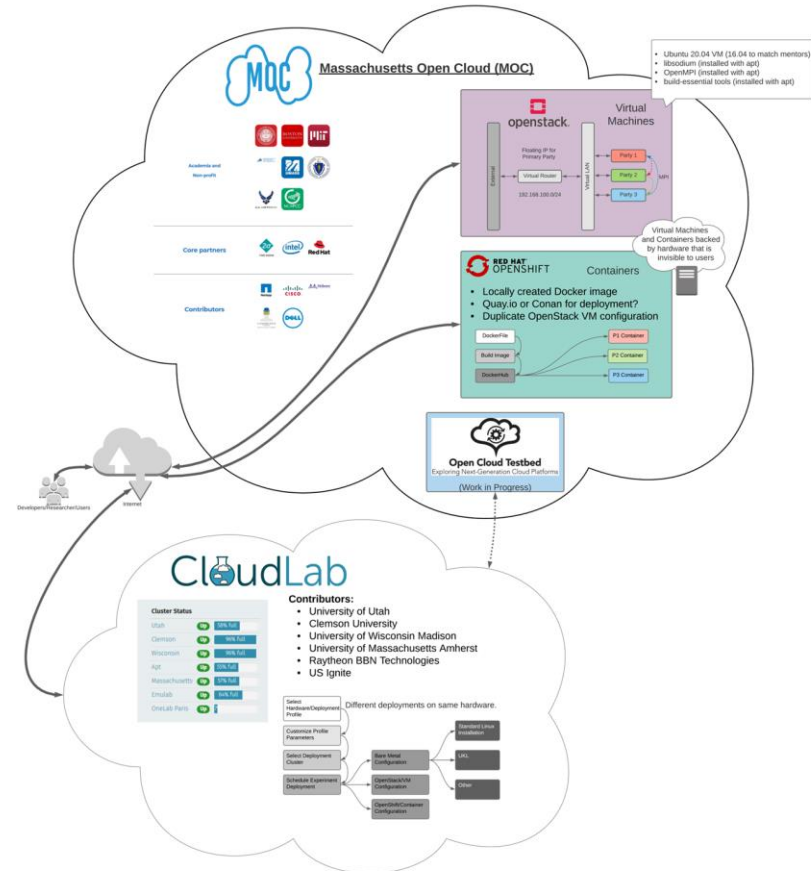
```
ylm@cc-mpc-main:~/mpc_shared/ccproject/experiments$ cat Makefile
# CC= gcc -std=c99
# CFLAGS= -O3 -Wall -lsodium      #Use this CFLAGS when deploying in OSX
CFLAGS= -O3 -Wall                #Use this CFLAGS & DEP when deploying in
DEP= -lsodium -lm                #Use this along with CFLAGS for Linux de
MPI=scorep mpicc
SRC=../src
PRIMITIVES= $(SRC)/comm.c $(SRC)/party.c $(SRC)/primitives.c $(SRC)/shar
RELATIONAL= $(SRC)/relational.c
```

```
ylm@cc-mpc-main:~/scorep-20201207_1014_3385102085294411$ scorep-score -r profile.cubex
Estimated aggregate size of event trace:          2432kB
Estimated requirements for largest trace buffer (max_buf): 828kB
Estimated memory requirements (SCOREP_TOTAL_MEMORY): 4097kB
(hint: When tracing set SCOREP_TOTAL_MEMORY=4097kB to avoid intermediate flushes
or reduce requirements using USR regions filters.)
```

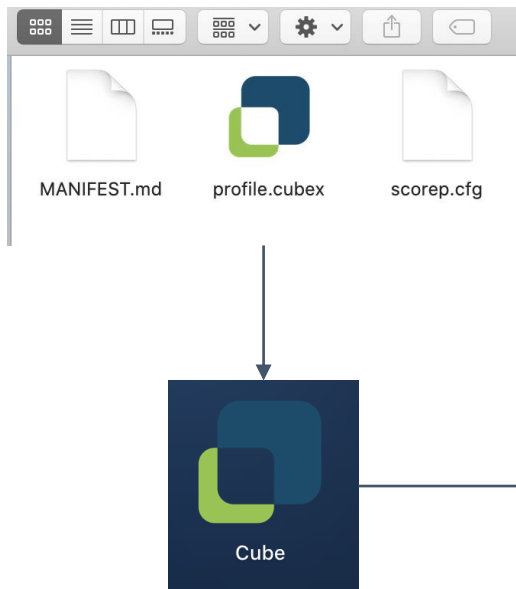
flt	type	max_buf[B]	visits	time[s]	time[%]	time/visit[us]	region
	ALL	847,365	73,111	1.41	100.0	19.28	ALL
	USR	442,494	49,055	0.01	1.0	0.28	USR
	MPI	352,700	18,038	1.38	97.8	76.47	MPI
	COM	52,130	6,015	0.02	1.2	2.76	COM
	SCOREP	41	3	0.00	0.0	21.55	SCOREP
	USR	208,234	24,027	0.00	0.2	0.13	check_init
	USR	104,130	12,015	0.00	0.2	0.25	get_rank
	MPI	89,089	3,003	0.01	0.7	3.24	MPI_Irecv
	MPI	89,089	3,003	0.07	4.9	23.13	MPI_Isend
	MPI	61,305	3,007	0.06	3.9	18.39	MPI_Send
	MPI	61,183	3,007	0.16	11.5	53.78	MPI_Recv
	USR	52,052	6,006	0.00	0.2	0.42	get_pred
	MPI	52,052	6,006	0.15	10.7	25.12	MPI_Wait
	USR	52,052	6,006	0.00	0.2	0.47	get_succ
	COM	26,000	3,000	0.01	0.6	2.69	exchange_shares_async
	COM	26,000	3,000	0.01	0.4	2.07	exchange_shares
	USR	26,000	1,000	0.00	0.2	2.24	generate_bool_share
	SCOREP	41	3	0.00	0.0	21.55	exp-exchange
	USR	26	1	0.00	0.0	59.69	init_sharing
	COM	26	3	0.00	0.0	54.74	exchange_rsz_seeds
	COM	26	3	0.00	0.0	56.31	generate_and_share_random_data
	COM	26	3	0.00	0.0	30.60	init
	MPI	26	3	0.00	0.0	14.96	MPI_Comm_rank
	MPI	26	3	0.00	0.0	1.59	MPI_Comm_size
	COM	26	3	0.00	0.0	25.63	exchange_shares_array
	MPI	26	3	0.00	0.1	268.18	MPI_Finalize
	MPI	26	3	0.93	66.1	310484.80	MPI_Init
	COM	26	3	0.00	0.1	603.82	main

Presentation Outline

- Overview
- Deployment Environments
- Testing Instrumentation
- **Data Analysis**
- Automation/Replicability
- Conclusions & Future Work



Why CUBE Gui



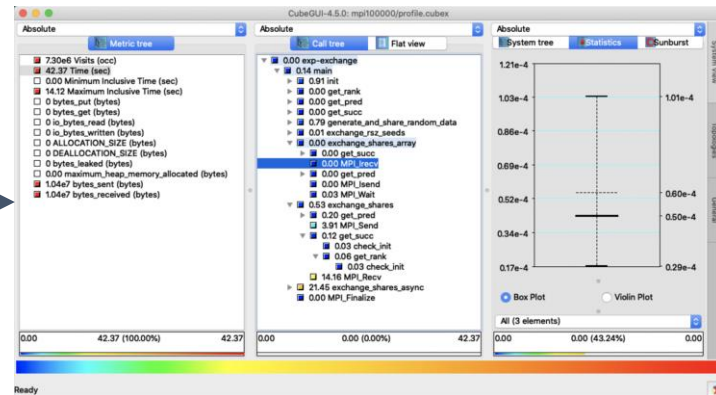
Experiment Row Data (.txt):

1. Limited to instrumented tests
2. No Secondary Statistic Info

Score-P Data (.cubex):

1. Inconvenient to view in text
2. No derived metrics when viewed directly

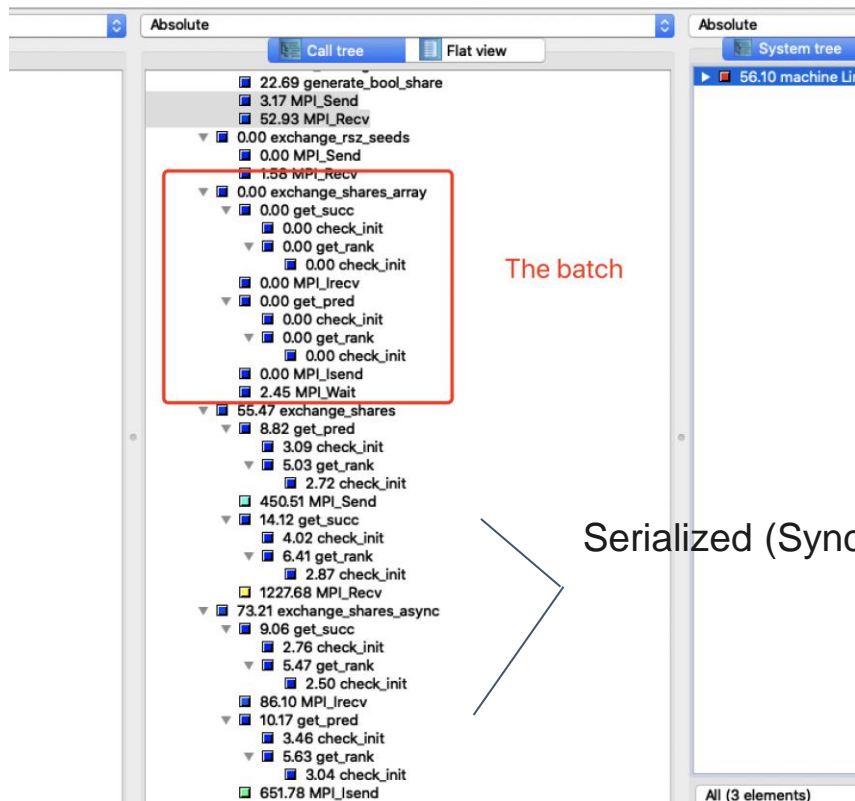
CUBE can help explore benchmark data



Data Insights

1. **Batch MPI** communication is **more efficient** than **Serialized (Sync/Async)**
2. Increasing batch sizes always helped in our tests
3. MPI Initialization cost is constant,
communication efficiency reaches a steady point for larger messages
1. The standard deviation of **Sending bytes** is always **2x** of the sd of **Receiving bytes**.

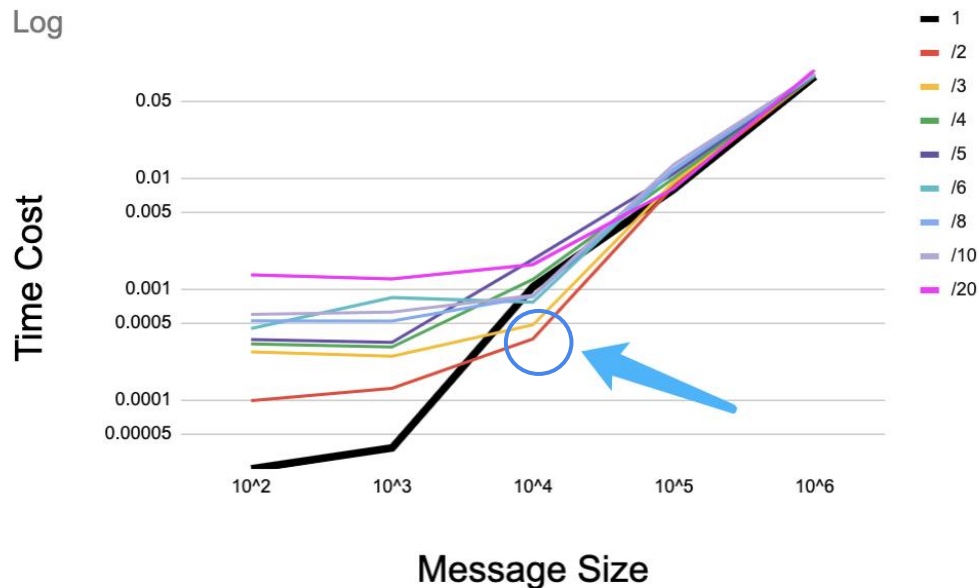
Batch MPI communication is **more efficient** than Serialized (Sync/Async)



1/100 of the time cost
for every message size

Serialized (Sync/Async)

Increasing batch sizes always helped in our tests.

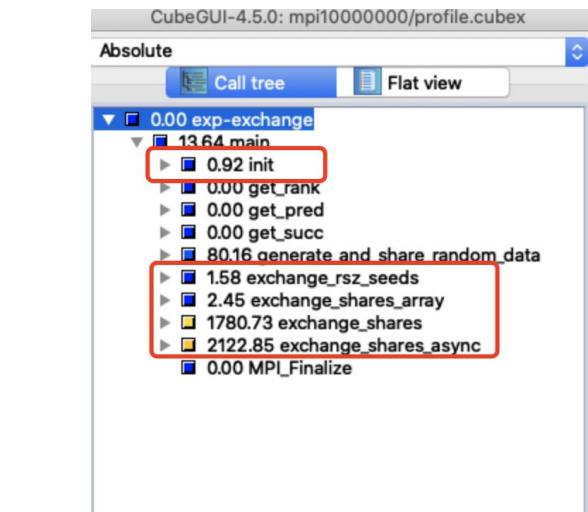


```

102
103  /* =====
104  1. Measure exchange_array
105  ===== */
106  // start timer
107  // CLOCK_PROCESS_CPUTIME_ID or CLOCK_MONOTONIC perhaps (check man
108  clock_gettime)
109  clock_gettime(CLOCK_MONOTONIC, &begin);
110
111  long part_CURR_ROW = CURR_ROW/6;
112  exchange_shares_array(r1s1, r1s2, part_CURR_ROW);
113  exchange_shares_array(r1s1, r1s2, part_CURR_ROW);
114  exchange_shares_array(r1s1, r1s2, part_CURR_ROW);
115  exchange_shares_array(r1s1, r1s2, part_CURR_ROW);
116  exchange_shares_array(r1s1, r1s2, part_CURR_ROW);
117  exchange_shares_array(r1s1, r1s2, part_CURR_ROW);
118  exchange_shares_array(r1s1, r1s2, part_CURR_ROW);
119
120  // stop timer
121  clock_gettime(CLOCK_MONOTONIC, &end);
122  elapsed = diff(begin, end);
123  time_stamp[nstep][iter][2] = elapsed;
124
125  if (rank == 0) {
126    printf("BATCHED\t%d\t%d\t%09d\n", CURR_ROW, (long long)elapsed.
127          tv_sec, elapsed.tv_nsec);
128  }
129
130  /* =====
131  2. Measure synchronous element-wise exchange
132  ===== */
133  // start timer
134  clock_gettime(CLOCK_MONOTONIC, &begin);
135
136  for (long i=0; i<CURR_ROW; i++) {

```

Initialization cost is constant,
communication efficiency: We have low efficiency for small messages which increases with message size until hitting reaching a steady efficiency value.



POP Assessment : exp-exchange

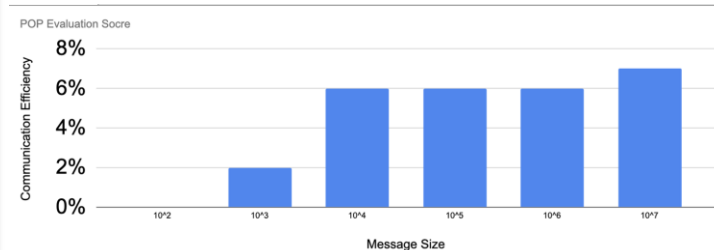
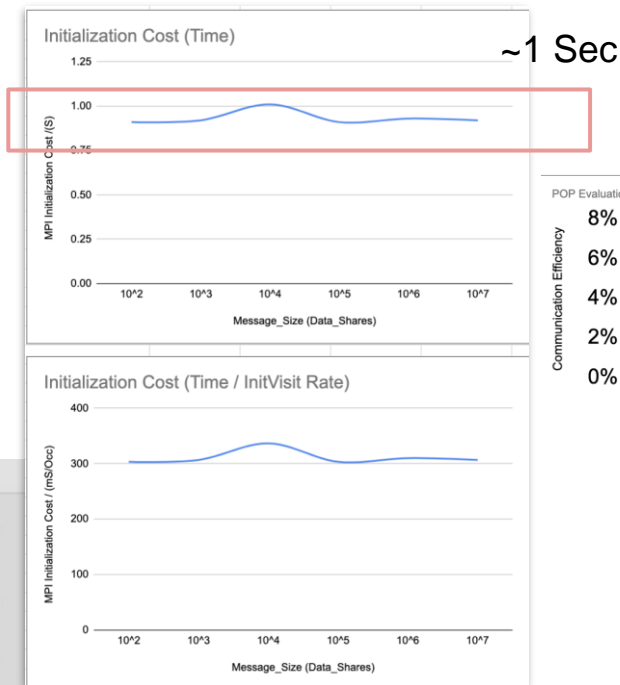
Parallel Efficiency 6% very poor 0.06

Load balance 83% very good 0.84

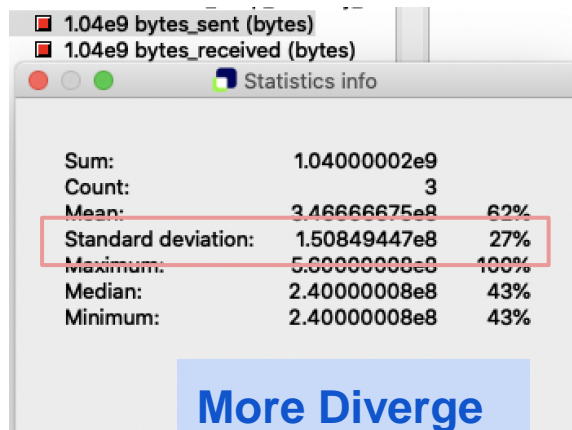
Communication Efficiency 7% very poor 0.08

Serialisation Efficiency

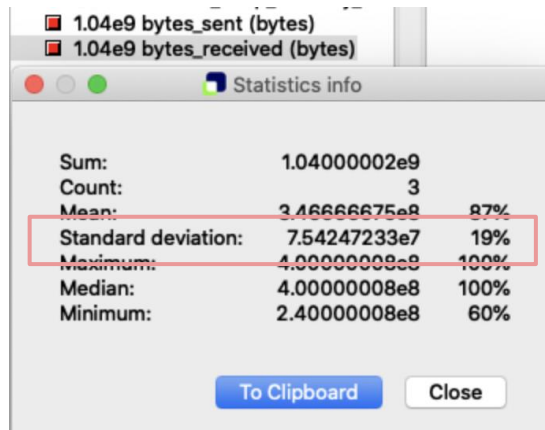
Transfer efficiency



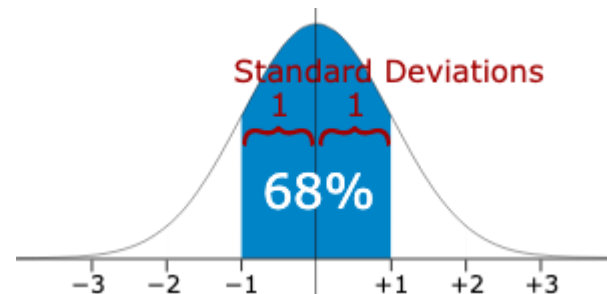
The standard deviation of **Sending bytes** is always **2x** of the SD of **Receiving bytes**.



Sending



Receiving

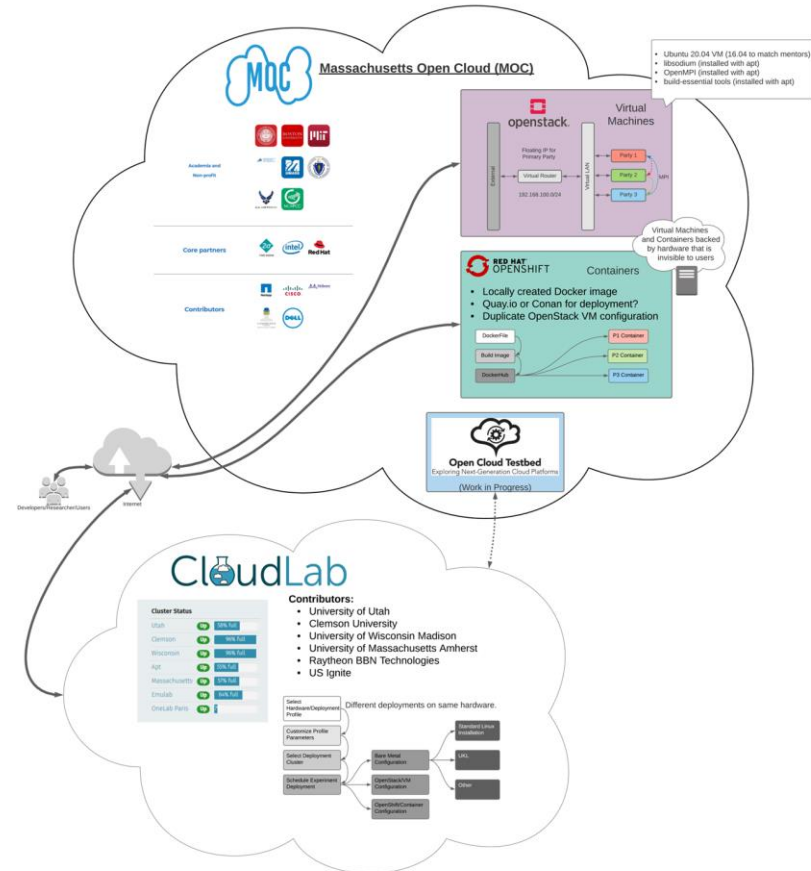


This is interesting to note and we hypothesize that this could indicate the opportunity for further optimization in the **sending** process.

Independently testing Sync and Async MPI transmissions might be a good follow-on experiment.

Presentation Outline

- Overview
- Deployment Environments
- Testing Instrumentation
- Data Analysis
- **Automation/Replicability**
- Conclusions & Future Work



Improved Automation

- Manual setup and testing
- Scripts - semi-automated
 - Shell Scripts
 - Packaged source *.tar.gz
- Ansible - advanced automation
 - MOC
 - Cloudlab

```

1  #!/bin/bash
2
3  # Bundle MPC Source code and Installer script for CloudLab
4  # MPC tests and experiments on Ubuntu 20.04 setup
5
6  # -c : creates a new archive
7  # -z : uses gzip compression
8  # -f : specifies files to be added
9  # In this case I only want *.c, *.h, and Makefile
10 # I use some regex in the paths specified to achieve this
11 tar -czf cloudlab_setup.tar.gz ../src/*.[ch] ../tests/*.[ch],Makefile} ../experiments/*.[ch],
    Makefile} setup.sh

```

Boston University CS & ECE

```

1  #!/bin/bash
2
3  # Installer script for CloudLab
4
5  # Check Linux Version
6  OS=$(lsb_release -i | cut -f2)
7  RELEASE=$(lsb_release -r | cut -f2)
8
9  # Alternate approach...use this to support CentOS eventually...
10 #grep ^NAME /etc/os-release
11
12 # Ubuntu
13 if [[ $OS == "Ubuntu" ]]; then
14   # Perform update
15   apt update -y && apt upgrade -y
16   # Install dependencies
17   if [[ $RELEASE == "16.04" ]]; then
18     apt install -y make gcc libopenmpi-dev openmpi-bin libsodium18 libsodium-dev
19   elif [[ $RELEASE == "18.04" ]]; then
20     apt install -y make gcc libopenmpi-dev openmpi-bin libsodium23 libsodium-dev
21   elif [[ $RELEASE == "20.04" ]]; then
22     apt install -y make gcc libopenmpi-dev openmpi-bin libsodium23 libsodium-dev
23   else
24     echo "Some other Ubuntu version"
25   fi
26
27 # Other
28 else
29   echo "Not Ubuntu"
30 fi
31
32 # Build MPC codebase (This will only work if the dependencies were successfully installed)
33 cd experiments
34 make all
35
36 # Now we're done with setup and can run something like this:
37 # mpirun --host localhost,192.168.1.2,192.168.1.3 -np 3 hostname
38 # mpirun --host localhost,192.168.1.2,192.168.1.3 -np 3 exp-exchange 1000

```

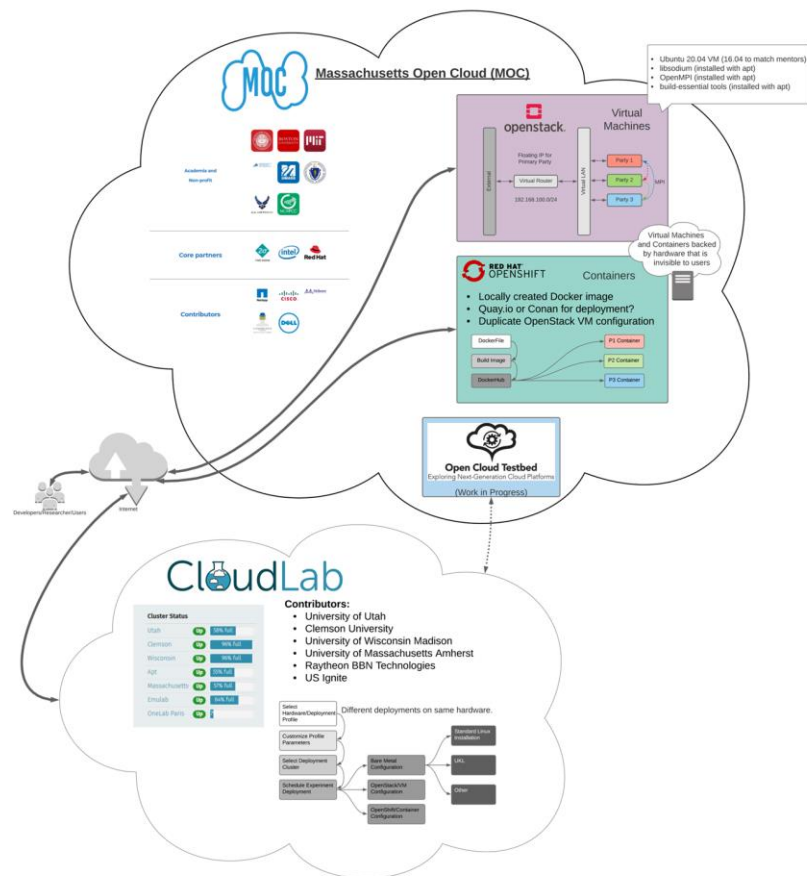
BOSTON
UNIVERSITY

From a user perspective...

1. Instantiate three nodes with and OS installed.
 - a. VMs on OpenStack (GUI web interface)
 - b. Bare-metal nodes on CloudLab (geni-lib script)
2. Clone project git repository
3. Modify ansible environment files
 - a. Specify SSH key for accessing nodes
 - b. Modify IP address or name specifying nodes
4. Conditional: If using a different OS (not Ubuntu 20.04)
 - a. Create a file specifying OS version of packages
5. Run ad-hoc ansible commands or playbook
 - a. Playbook will execute an entire experiment scenario

Presentation Outline

- Overview
- Deployment Environments
- Testing Instrumentation
- Data Analysis
- Automation/Replicability
- Conclusions & Future Work



Conclusion - Lessons Learned

- OpenStack experience
- Docker/OpenShift experience
- Linux system administration/configuration
- MPI instrumenting and testing
- Geni-lib scripting
- Advanced Shell Scripting
- Ansible first time experience

Links to documentation and Repo here

- Github repository
 - Public: https://github.com/BU-CLOUD-F20/Secure_MultiParty
 - Private: <https://github.com/jliagouris/ccproject>
- Google Drive (data collection)
 - https://drive.google.com/drive/folders/143M9s_VxxjImq6VWSjr6qNxkYMCzSz9j?usp=sharing
 - Request access to other resources from project team
- Additional Resources:
 - Massachusetts Open Cloud (MOC): <https://massopen.cloud/>
 - Cloudlab (Bare-metal): <https://www.cloudlab.us/>
 - Open-MPI: <https://www.open-mpi.org/>
 - Score-P: <https://www.vi-hps.org/projects/score-p>
 - Scalasca (CUBE GUI and more...): <https://www.scalasca.org/>
 - Ansible: <https://docs.ansible.com/ansible/latest/index.html>

Future Work - Ongoing Efforts

- End-to-end application on this framework
- What is the best deployment (performance)?
 - Openstack/OpenShift on Bare-Metal (follow-on test)
- Build user-friendly front-end UI to perform steps
- More extensive data analysis for MPI using Score-P
- Deploy the codebase on Linux Unikernel

Thank you

...any questions?