# Secure Multiparty Computation Sprint 4

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

**Boston University** CS & ECE

**BOSTON UNIVERSITY**
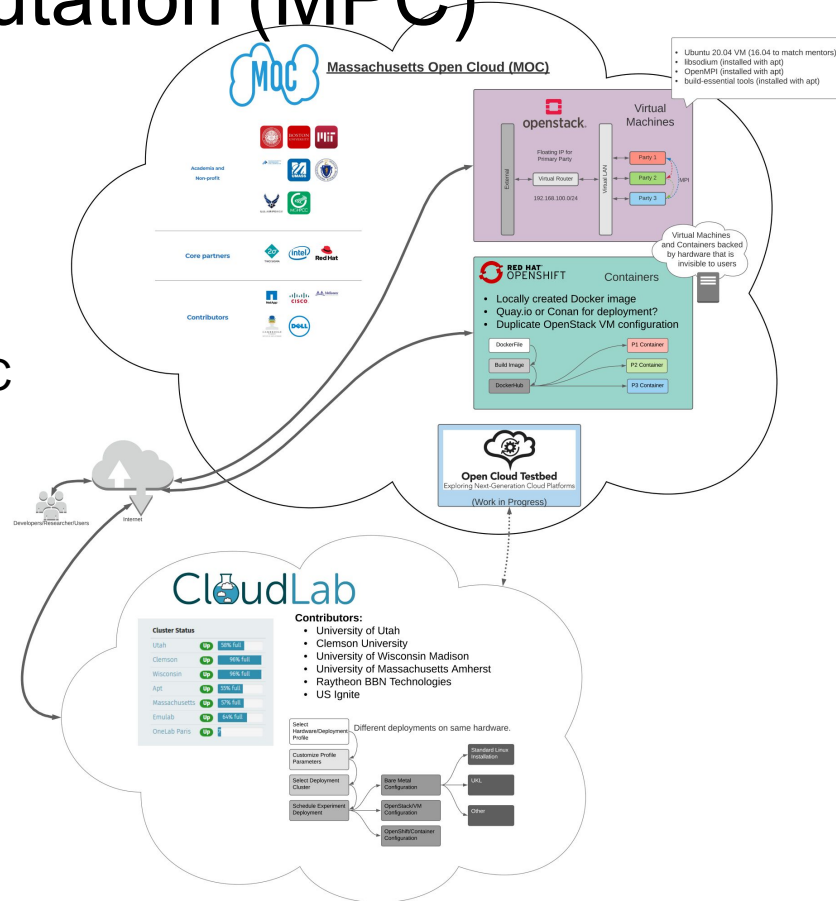
# Presentation Outline

- Project Recap
- Project Goals & Sprint 4 Stories/Tasks
- Work Accomplished & Information Learned
  - Bare-Metal → CloudLab Testing
  - Containers → Docker/OpenShift
- Project Organization Assessment (Burndown)
- Sprint 5 goals (Mentor priorities)

**Boston University** CS & ECE

BOSTON
UNIVERSITY

# Recap of Multi-Party Computation (MPC)

- MPC enables...
  - Shared Computation on Private Data
  - Protects the Privacy of Data
  - Mutually Agreed Computation
- Our mentors…
  - Are using three party Secret Sharing MPC
  - Perform Relational Queries with MPC
  - Keep all parts secure vs. splitting into secure and insecure steps
- Our mission…
  - Profile this new MPC library
  - Identify bottlenecks
  - Compare deployment scenarios and find the best performance

**Boston University** CS & ECE

# Project Goals & Sprint 4 Stories/Tasks

- Presentation
  - Finalize SmartNIC presentation
- CloudLab
  - Geni-lib automation
  - Bare metal testing
- Containers
  - Debug OpenMPI on Docker
  - Docker/Docker-compose work
  - Status with OpenShift (kompose)

**Boston University** CS & ECE

---

### ∨ Sprint 4

|  |  |
| --- | --- |
|  | 60 closed |
| 29 Oct 2020-12 Nov 2020 | 88 total |

#117 As a student in the CC course, I want prepare and deliver a practice presentation about the selected article to be ready for the class presentation. 🕐     40

#11 As a team member, I want to build a containerized MPC environment on OpenShift, like one on the VMs.     12

#162 As a team member, I want to benchmark the results obtained by running the codebase on multi bare-metal nodes     16

#211 As a team member, I want to create a demo summarizing accomplishments in order to show progress to the clients     20

**BOSTON UNIVERSITY**

# Bare-Metal Testing

CloudLab, geni-lib scripts, data collection, ...

**Boston University** CS & ECE

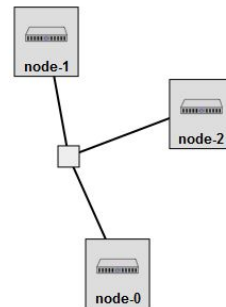# CloudLab: Geni-lib Automation

- LAN vs Link
- Static IP
- Setup
  - Dependencies
  - MPC Code duplicate
    - geni.rspec.pg.Install
- Remaining Issue:
  - geni.rspec.pg.Execute

**Boston University** CS & ECE

```python
1   """ubuntu baremetal ring of nodes"""
2
3   #
4   # NOTE: This code was machine converted. An actual human would not
5   #       write code like this!
6   #
7
8   # Import the Portal object.
9   import geni.portal as portal
10  # Import the ProtoGENI library.
11  import geni.rspec.pg as pg
12  # Import the Emulab specific extensions.
13  import geni.rspec.emulab as emulab
14
15  pc = portal.Context()
16
17  pc.defineParameter("node_type", "Hardware Type",
18                      portal.ParameterType.NODETYPE, "any")
19  pc.defineParameter("node_count", "Number of Machines",
20                      portal.ParameterType.INTEGER, 3)
21
22  params = pc.bindParameters()
23
24  request = portal.context.makeRequestRSpec()
25
26  node = []
27  link = []
28
29  # Create selected number of nodes
30  for i in range(params.node_count):
31      node.append(request.RawPC('node-%d' % i))
32      node[-1].disk_image = 'urn:publicid:IDN+emulab.net+image+emulab-ops:UBUNTU16-64-STD'
33      node[-1].hardware_type = params.node_type
34
35  # Create a LAN for all the connections
36  lan = request.LAN("lan")
37
38  # Create a link between each of the nodes to make a ring
39  for i in range(params.node_count):
40      iface = node[i].addInterface("if1")
41      iface.component_id = "eth1"
42      iface.addAddress(pg.IPv4Address("192.168.1."+str(i+1), "255.255.255.0"))
43      lan.addInterface(iface)
44
45  # Install and execute scripts on each node
46  for i in range(params.node_count):
47      node[i].addService(pg.Install(url="https://www.dropbox.com/s/7t91cf0ugt66ypl/cloudlab_setup.tar.gz", path="/home/mpc"))
48      node[i].addService(pg.Execute(shell="bash", command="/home/mpc/setup.sh"))
49
50  # Print the generated rspec
51  pc.printRequestRSpec(request)
```
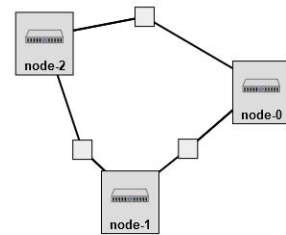
# CloudLab: Geni-lib Further Tweaks

- **Modifications**
  - Link multiplexing
  - Best Effort
- **Issues with Ring Topology**
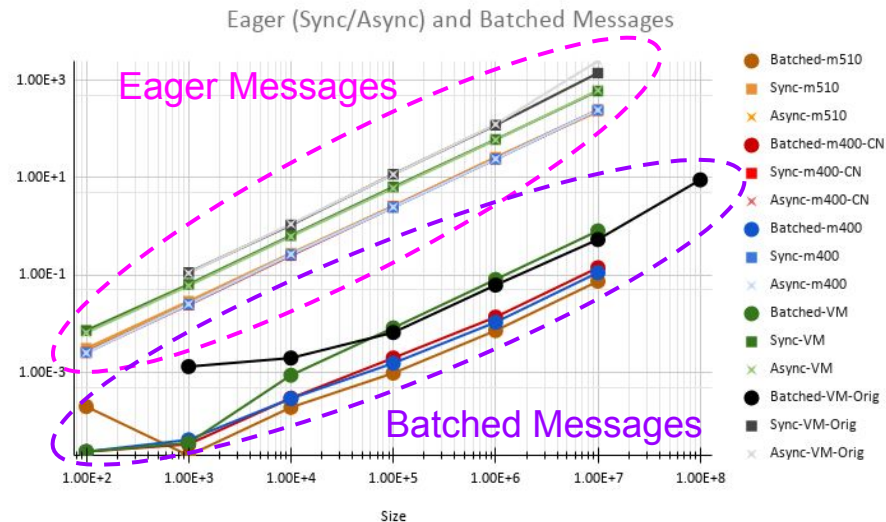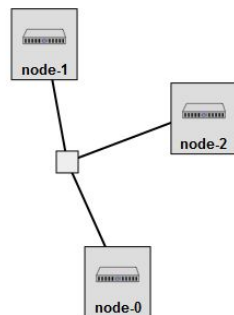  - Couldn't map to hardware
- Current Status: Still WIP

**Boston University** CS & ECE



```python
"""ubuntu baremetal ring of nodes"""

#
# NOTE: This code was machine converted. An actual human would not
#       write code like this!
#

# Import the Portal object.
import geni.portal as portal
# Import the ProtoGENI library.
import geni.rspec.pg as pg
# Import the Emulab specific extensions.
import geni.rspec.emulab as emulab

pc = portal.Context()

pc.defineParameter("node_type", "Hardware Type",
                   portal.ParameterType.NODETYPE, "any")
pc.defineParameter("node_count", "Number of Machines",
                   portal.ParameterType.INTEGER, 3)

params = pc.bindParameters()

request = portal.context.makeRequestRSpec()

node = []
link = []

# Create selected number of nodes
for i in range(params.node_count):
    node.append(request.RawPC('node-%d' % i))
    node[-1].disk_image = 'urn:publicid:IDN+emulab.net+image+emulab-ops:UBUNTU16-64-STD'
    node[-1].hardware_type = params.node_type

# Create a LAN for all the connections
#lan = request.LAN("lan")

# Create a link between each of the nodes to make a ring
#for i in range(params.node_count):
#    iface = node[i].addInterface("if1")
#    iface.component_id = "eth1"
#    iface.addAddress(pg.IPv4Address("192.168.1."+str(i+1), "255.255.255.0"))
#    lan.addInterface(iface)

# Create two links between them
link1 = request.Link(members = [node[0], node[1]])
iface1 = node[0].addInterface("if1")
iface1.component_id = "eth1"
iface1.addAddress(pg.IPv4Address("192.168.1."+str(1), "255.255.255.0"))
link1.addInterface(iface1)

link2 = request.Link(members = [node[1], node[2]])
iface2 = node[1].addInterface("if2")
iface2.component_id = "eth1"
iface2.addAddress(pg.IPv4Address("192.168.1."+str(2), "255.255.255.0"))
link2.addInterface(iface2)

link3 = request.Link(members = [node[2], node[0]])
iface3 = node[2].addInterface("if3")
iface3.component_id = "eth1"
iface3.addAddress(pg.IPv4Address("192.168.1."+str(3), "255.255.255.0"))
link3.addInterface(iface3)

# Turn on link multiplexing. Note that this also turns on vlan encapsulation
# You have to set this both links.
link1.link_multiplexing = True
link2.link_multiplexing = True
link3.link_multiplexing = True

# But the resource mapper is going to try to prevent the two links from oversubscribing
# the physical link. For example, trying to create two 1Gb multiplexed links on top of a 1Gb
# physical link. Sometimes this is the correct behaviour. But if not, do this to turn
# off the checks.
link1.best_effort = True
link2.best_effort = True
link3.best_effort = True

# Install and execute scripts on each node
for i in range(params.node_count):
    node[i].addService(pg.Install(url="https://www.dropbox.com/s/7t91cf0ugt66ypl/cloudlab_setup.tar.gz", path="/home/mpc"))
    node[i].addService(pg.Execute(shell="bash", command="/home/mpc/setup.sh"))

# Print the generated rspec
pc.printRequestRSpec(request)
```

# CloudLab Testing

- ## Shared Control Network
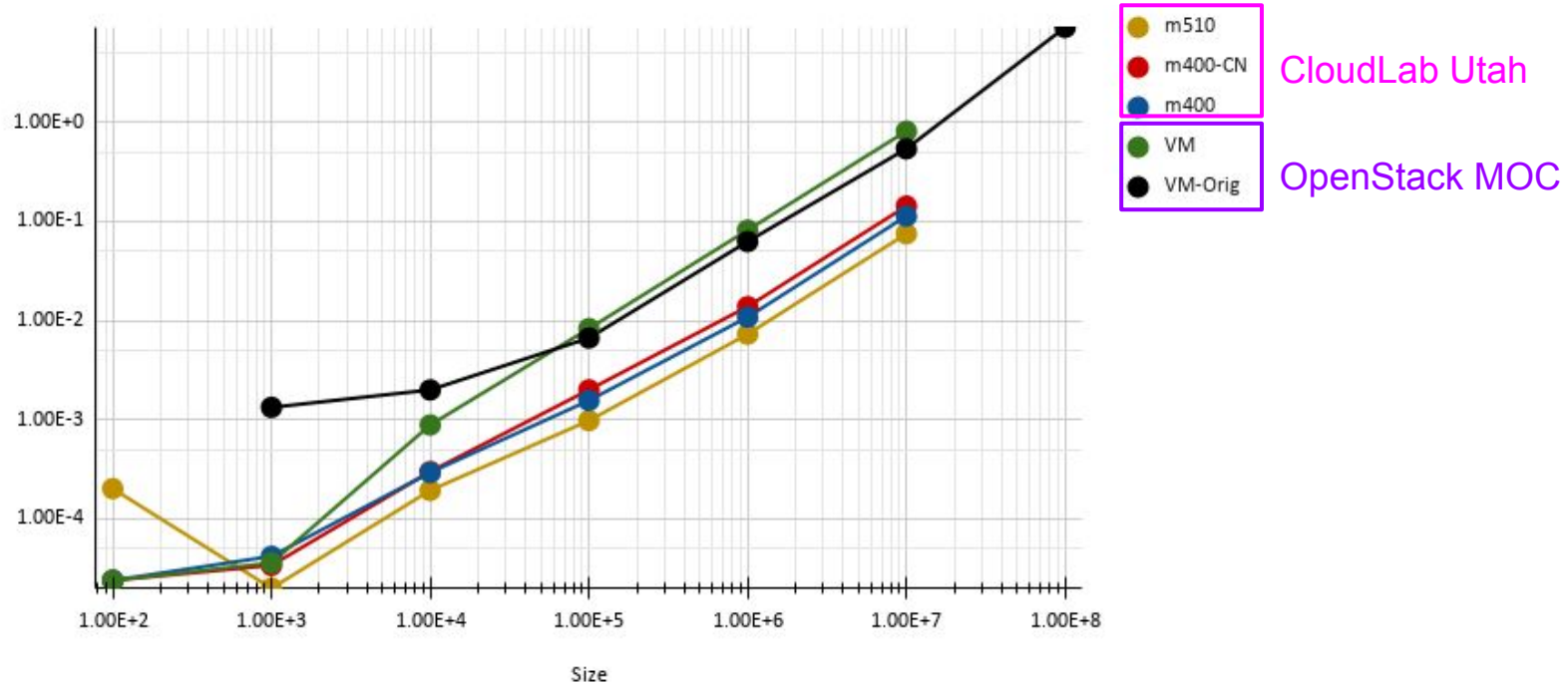  - ### Hardware
    - #### m400 (ARM)
- ## LAN topology
  - ### Hardware
    - #### m510  (x86_64)
    - #### m400 (ARM)





**Boston University** CS & ECE

# Batched Messages - Different Deployments



**Boston University** CS & ECE
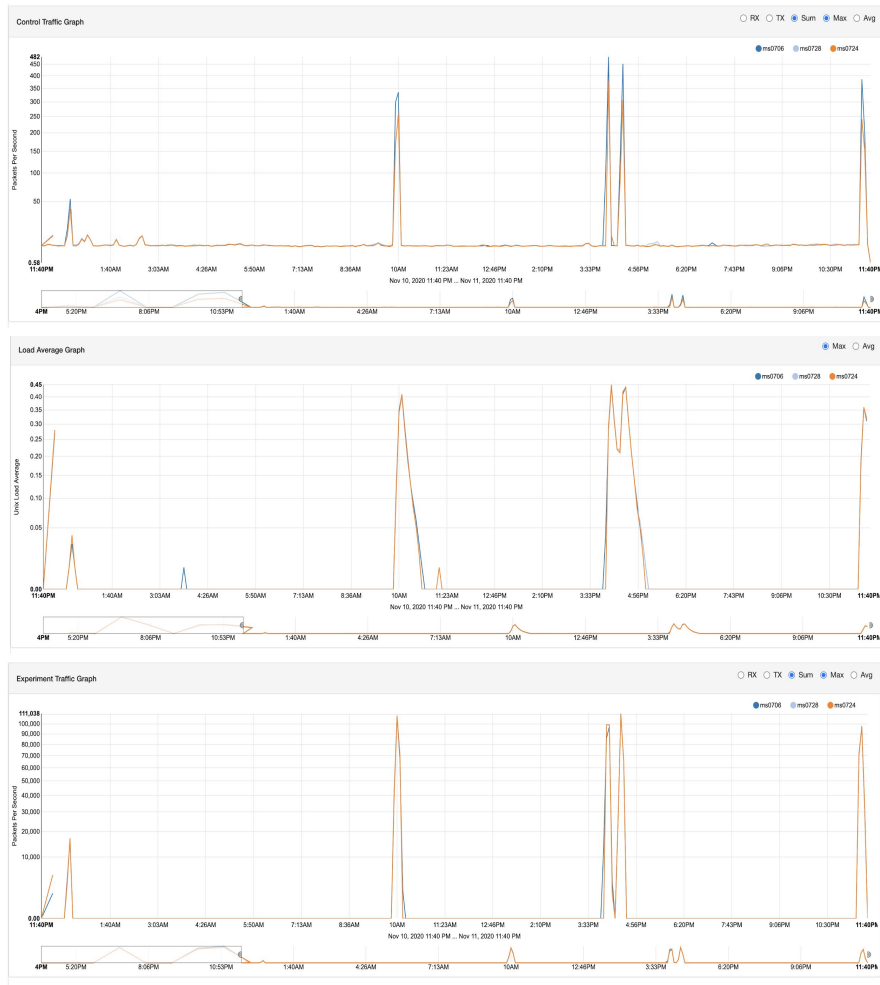
# Observations

- News tests are using new Ubuntu image (14.04.1 → 16.04.01)
- When running 10M size test, performance varies might because of the time period. Test gets stuck in peak-hours and function normally in early morning EST.
- Experiments of smaller sizes are not restricted by time period.

**Boston University** CS & ECE

# Working with Containers

Docker, docker-compose, OpenShift, ...

**Boston University** CS & ECE

# MPC in Docker Container Debugging...

- Running tests.sh → OK
- Running exp-exchange
  - Issue with size greater than 505… ex: with 1000
  - Cryptic message…
- Determine source…
  - Some clues but overall meaning still unclear

```
mpirun -np 3 exp-exchange1000

# which produces the following:

root@ebd1c7f24dfe:~/experiments#

[ebd1c7f24dfe:00046] Read -1, expected 8000, errno = 1

[ebd1c7f24dfe:00046] Read -1, expected 8000, errno = 1

[ebd1c7f24dfe:00047] Read -1, expected 8000, errno = 1

[ebd1c7f24dfe:00047] Read -1, expected 8000, errno = 1

[ebd1c7f24dfe:00045] Read -1, expected 8000, errno = 1

[ebd1c7f24dfe:00047] Read -1, expected 8000, errno = 1

[ebd1c7f24dfe:00046] Read -1, expected 8000, errno = 1

BATCHED 1000     0.00005

SYNC    1000     0.00047

ASYNC   1000     0.00042
```

Size changes value

[Hostname:PID]

**Boston University** CS & ECE

BOSTON
UNIVERSITY

# OpenMPI Debugging Continued...

```
$12 = {MPI_SOURCE = 0, MPI_TAG = 193, MPI_ERROR = 0, _cancelled = 0, _ucount = 8000}
(gdb) p result1
$13  0
(gdb) p status2
$14 = {MPI_SOURCE = 0, MPI_TAG = 193, MPI_ERROR = 0, _cancelled = 0, _ucount = 8000}
(gdb) p result2
$15  0
(gdb) l
220        }
221      else { //P3
222        result1 = MPI_Recv(r1s1, ROWS, MPI_LONG_LONG, 0, SHARE_TAG, MPI_COMM_WORLD, &status1);
223        result2 = MPI_Recv(r1s2, ROWS, MPI_LONG_LONG, 0, SHARE_TAG, MPI_COMM_WORLD, &status2);
224        }
225      }
(gdb)
```

MPI_Recv → No error directly returned...

```
mpirun -np 3 exp-exchange 1000
# which produces the following:
root@ebd1c7f24dfe:~/experiments#
[ebd1c7f24dfe:00046] Read -1, expected 8000, errno = 1
[ebd1c7f24dfe:00046] Read -1, expected 8000, errno = 1
[ebd1c7f24dfe:00047] Read -1, expected 8000, errno = 1
[ebd1c7f24dfe:00046] Read -1, expected 8000, errno = 1
[ebd1c7f24dfe:00045] Read -1, expected 8000, errno = 1
[ebd1c7f24dfe:00047] Read -1, expected 8000, errno = 1
[ebd1c7f24dfe:00046] Read -1, expected 8000, errno = 1
BATCHED 1000    0.00005
SYNC    1000    0.00047
ASYNC   1000    0.00042
```

Looking for matching message pattern on github...

```
68  #if OPAL_BTL_SM_HAVE_CMA
69  int mca_btl_sm_get_cma (mca_btl_base_module_t *btl, mca_btl_base_endpoint_t *endpoint, void *local_address,
70                          uint64_t remote_address, mca_btl_base_registration_handle_t *local_handle,
71                          mca_btl_base_registration_handle_t *remote_handle, size_t size, int flags,
72                          int order, mca_btl_base_rdma_completion_fn_t cbfunc, void *cbcontext, void *cbdata)
73  {
74      struct iovec src_iov = {.iov_base = (void *)(intptr_t) remote_address, .iov_len = size};
75      struct iovec dst_iov = {.iov_base = local_address, .iov_len = size};
76      ssize_t ret;
77
78      /*
79       * According to the man page :
80       * "On success, process_vm_readv() returns the number of bytes read and
81       * process_vm_writev() returns the number of bytes written.  This return
82       * value may be less than the total number of requested bytes, if a
83       * partial read/write occurred.  (Partial transfers apply at the
84       * granularity of iovec elements.  These system calls won't perform a
85       * partial transfer that splits a single iovec element.)".
86       * So since we use a single iovec element, the returned size should either
87       * be 0 or size, and the do loop should not be needed here.
88       * We tried on various Linux kernels with size > 2 GB, and surprisingly,
89       * the returned value is always 0x7ffff000 (fwiw, it happens to be the size
90       * of the larger number of pages that fits a signed 32 bits integer).
91       * We do not know whether this is a bug from the kernel, the libc or even
92       * the man page, but for the time being, we do as is process_vm_readv() could
93       * return any value.
94       */
95      do {
96          ret = process_vm_readv (endpoint->segment_data.other.seg_ds->seg_cpid, &dst_iov, 1, &src_iov, 1, 0);
97          if (0 > ret) {
98              opal_output(0, "Read %ld, expected %lu, errno = %d\n", (long)ret, (unsigned long)size, errno);
99              return OPAL_ERROR;
100         }
101         src_iov.iov_base = (void *)((char *)src_iov.iov_base + ret);
102         src_iov.iov_len -= ret;
103         dst_iov.iov_base = (void *)((char *)dst_iov.iov_base + ret);
104         dst_iov.iov_len -= ret;
105     } while (0 < src_iov.iov_len);
106
107     /* always call the callback function */
108     cbfunc (btl, endpoint, local_address, local_handle, cbcontext, cbdata, OPAL_SUCCESS);
109
110     return OPAL_SUCCESS;
111 }
112 #endif
```

# OpenMPI Debugging Resolution

- **Issue for messages larger than ~1k**
- **Shared memory Byte Transport Layer (BTL)**
  - "Sm" was the original version
  - "Vader" is the current version
- **CMA (Cross Memory Attach)**
  - Kernel support required for "Zero copy" mechanism
- **Bypass with parameter**
  - mpirun --mca btl_vader_single_copy_mechanism none -np 3 exp-exchange 1000

disable CMA in vader #3270

Closed   hunsa opened this issue on Apr 2, 2017 · 9 comments

hunsa commented on Apr 2, 2017

Hello all,

We experienced problems with Open MPI when communication larger messages, for example with `MPI_Gather`, `MPI_Allgather`, etc. Large means messages larger than 1k (for very small messages the problem did not occur).

We received error messages like this
Read -1, expected 8000, errno = 38

I suspected that there is something wrong with the CMA support.

And indeed, until Open MPI 2.0.2 the configure script would have the following result
`#define OPAL_BTL_SM_HAVE_CMA 0`

Now, the new CMA detection method in 2.1.0 leads to
`#define OPAL_BTL_SM_HAVE_CMA 1`

That wouldn't be a problem if we could disable CMA in vader during configure, but it does not seem to be possible. So, `--with-cma=no` or `--without-cma` will have no effect and we will end up with
`#define OPAL_BTL_SM_HAVE_CMA 1`

Currently, we can set
`OMPI_MCA_btl_vader_single_copy_mechanism=none`
and the error messages will not show up.

It would be great if we could manually disable CMA support during configure.

Thank you

# Working with Docker and docker-compose

- docker-compose.yml
  - Define topology
  - Multiple "services" each a container based on Dockerfile
  - Virtual network for communication
- Build Configuration
  - docker-compose build
- Launch Three Containers
  - docker-compose up

- Dockerfile
  - Install dependencies
  - Configure non-root user
  - Configure ssh
  - Build MPC code
- Build Image
  - docker build -t mpc .
- Launch Single Container
  - docker run --name mpc -it mpc /bin/bash

**Boston University** CS & ECE

# Container Configuration Hierarchy

**docker-compose**

**Three containers**

**Network Config.**

**Dockerfile**

**Dependencies**

**SSH Setup**

**Code Build**

**Boston University** CS & ECE

```yaml
version: '2'

services:
  party-0:
    hostname: party-0
    build:
      dockerfile: Dockerfile
      context: .
    tty: true
    networks:
      mpc_net:
        ipv4_address: 192.168.1.11
    ports:
      - "22"
  party-1:
    hostname: party-1
    build:
      dockerfile: Dockerfile
      context: .
    tty: true
    networks:
      mpc_net:
        ipv4_address: 192.168.1.12
    ports:
      - "22"
  party-2:
    hostname: party-2
    build:
      dockerfile: Dockerfile
      context: .
    tty: true
    networks:
      mpc_net:
        ipv4_address: 192.168.1.13
    ports:
      - "22"

networks:
  mpc_net:
    driver: bridge
    ipam:
      driver: default
      config:
        - subnet: 192.168.1.0/24
          gateway: 192.168.1.1
```

```dockerfile
# Dockerfile based on: https://github.com/oweidner/docker.openmpi/blob/master/Dockerfile
# https://codeburst.io/direct-connection-to-a-docker-container-with-ssh-56e1d2744ee5
# Build this image: docker build -t mpc .

FROM ubuntu:20.04
MAINTAINER Pierre-Francois Wolfe <pwolfe@bu.edu>

ENV USER mpc
ENV HOME=/home/${USER}
ARG DEBIAN_FRONTEND=noninteractive

RUN apt update -y && \
    apt-get install -y --no-install-recommends sudo apt-utils && \
    apt-get install -y --no-install-recommends openssh-server \
    make ssh gcc libopenmpi-dev openmpi-bin libsodium23 libsodium-dev && \
    apt clean && \
    apt purge && \
    rm -rf /var/lib/apt/lists/* /tmp/* /var/tmp/*

RUN mkdir /var/run/sshd
RUN echo 'root:${USER}' | chpasswd
RUN sed -i 's/PermitRootLogin without-password/PermitRootLogin yes/' /etc/ssh/sshd_config

# SSH Login fix. Otherwise user is kicked off after login
RUN sed 's@session\s*required\s*pam_loginuid.so@session optional pam_loginuid.so@g' -i /etc/pam.d/sshd

ENV NOTVISIBLE "in users profile"
RUN echo "export VISIBLE=now" >> /etc/profile

# ---------------------------------------------
# Add an 'mpc' user

RUN adduser --disabled-password --gecos "" ${USER} && \
    echo "${USER} ALL=(ALL) NOPASSWD:ALL" >> /etc/sudoers
```

```dockerfile
# Set-up                          key
# ---------------------------------------------
ENV SSHDIR ${HOME}/.ssh/

RUN mkdir -p ${SSHDIR}

ADD ssh/config ${SSHDIR}/config
ADD ssh/id_rsa.mpi ${SSHDIR}/id_rsa
ADD ssh/id_rsa.mpi.pub ${SSHDIR}/id_rsa.pub
ADD ssh/id_rsa.mpi.pub ${SSHDIR}/authorized_keys

RUN chmod -R 600 ${SSHDIR}* && \
    chown -R ${USER}:${USER} ${SSHDIR}

# ---------------------------------------------
# Copy MPC code
# ---------------------------------------------
COPY src/* /home/${USER}/code/src/
COPY experiments/* /home/${USER}/code/experiments/
COPY tests/* /home/${USER}/code/tests/
COPY launch.sh /home/${USER}/code/experiments/

WORKDIR /home/${USER}/code/experiments
RUN make exp-exchange

ENV OMPI_ALLOW_RUN_AS_ROOT=1
ENV OMPI_ALLOW_RUN_AS_ROOT_CONFIRM=1

RUN chown -R ${USER}:${USER} ${HOME}/code

RUN /usr/bin/ssh-keygen -A

EXPOSE 22

CMD ["/usr/sbin/sshd", "-D"]
```

# Docker-compose to OpenShift

- ## Current (local)
  - ### docker-compose
    - #### Docker Containers
- ## In Progress
  - ### OpenShift (MOC)
    - #### Kubernetes
      - ##### Docker Containers
- ## Conversion Tool
  - ### Kubernetes Kompose

**Boston University** CS & ECE

```
docker-compose.yml
```

```
kompose --provider openshift --file
docker-compose.yml convert
```

```
party-0-service.yaml
party-1-service.yaml
party-2-service.yaml
party-0-deploymentconfig.yaml
party-0-imagestream.yaml
party-1-deploymentconfig.yaml
party-1-imagestream.yaml
party-2-deploymentconfig.yaml
party-2-imagestream.yaml
```

```
oc new-app scripts_party-0
scripts_party-1 scripts_party-2
```

Launching on OpenShift WIP

APPLICATION
scriptsparty-0

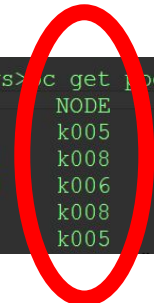| > | DEPLOYMENT CONFIG scriptsparty-0 | ⊗ 1 Error | No deployments for scriptsparty-0 | ⋮ |
| > | DEPLOYMENT CONFIG scriptsparty-1 | ⊗ 1 Error | No deployments for scriptsparty-1 | ⋮ |
| > | DEPLOYMENT CONFIG scriptsparty-2 | ⊗ 1 Error | No deployments for scriptsparty-2 | ⋮ |

BOSTON UNIVERSITY

# OpenShift - Deploying Pods on Specific Nodes

When deploying multiple pods of same application, pods get assigned to different worker nodes randomly.

```
B:\Study\Boston University\Fall 2020\EC 528 - Cloud Computing\Project\openshift-client-windows>oc get pods -o wide
NAME                            READY   STATUS               RESTARTS   AGE   IP            NODE
docker101tutorial-794f8f8dd8-c5mlb   1/1   #Cre Running       1          34d   10.128.1.34   k005
docker101tutorial-794f8f8dd8-dlgm8   1/1     Running          1          34d   10.128.8.101  k008
docker101tutorial-794f8f8dd8-js4ww   1/1   oc app Running     1          34d   10.128.4.189  k006
docker101tutorial-794f8f8dd8-wljmd   1/1     Running          1          34d   10.128.8.93   k008
docker101tutorial-794f8f8dd8-zkzkh   1/1   #Upd Running       1          34d   10.128.1.33   k005
```

Running all MPC parties pods on the same nodes may have different latencies/computation delays as is the case otherwise.

**Boston University** CS & ECE

**BOSTON UNIVERSITY**

OpenShift - Deploying Pods on Specific Nodes

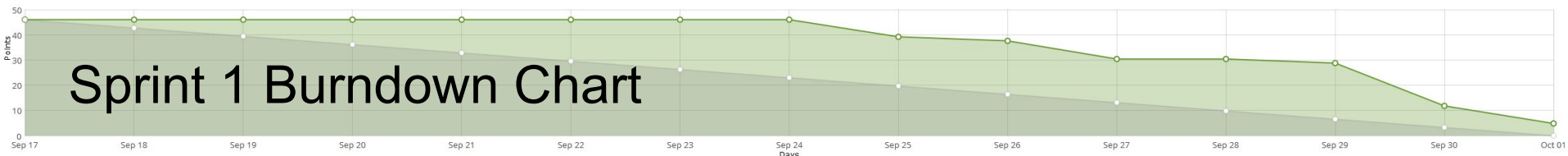Tried specifying what nodes to place pods on.

Failed → Most documentation for this is for OpenShift 3, while the latest version we're using is OpenShift 4.

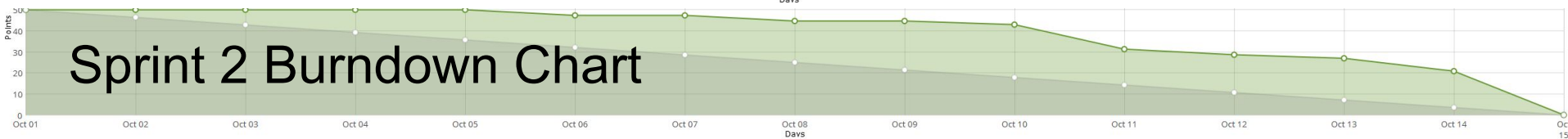An authorization error is encountered: Admin access / elevated rights needed !

```
B:\Study\Boston University\Fall 2020\EC 528 - Cloud Computing\Project\openshift-client-windows>oc get nodes
Error from server (Forbidden): nodes is forbidden: User "hasnain@bu.edu" cannot list nodes at the cluster scope: no RBAC policy matched
B:\Study\Boston University\Fall 2020\EC 528 - Cloud Computing\Project\openshift-client-windows>oc edit namespace
Error from server (Forbidden): namespaces is forbidden: User "hasnain@bu.edu" cannot list namespaces at the cluster scope: no RBAC policy matched
```
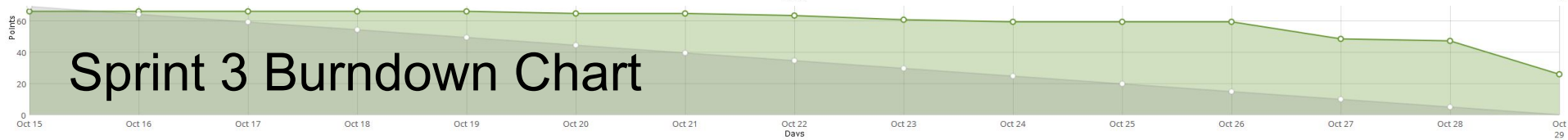
**Boston University** CS & ECE
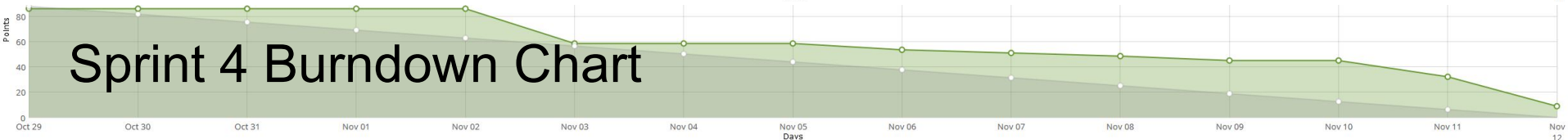
Sprint 1 Burndown Chart

Sprint 2 Burndown Chart

Sprint 3 Burndown Chart

Sprint 4 Burndown Chart

**Boston University** CS & ECE

# Sprint 5 - Some Known Stories

- As a researcher, I want to use a comprehensive test suite across all deployments...
  - Exp-exchange: iterate over message size, transaction size, other MPI options...
  - Add score-p wrapper toggle (and use some MPI tools for insights)
  - Push all changes to test setups (CloudLab *.tar.gz, rebuild Docker, rsync to VMs)
- As a researcher, I want to deploy final tests to each environment of interest.
  - Docker containers on OpenShift (finish debugging)
  - Revisit existing VM setup on OpenStack
  - Ring topology on CloudLab (if possible, otherwise LAN)
- As a researcher, I want to compare the same tests run on different platforms.
  - Conclusions about best performance of tested methods
  - Documentation and any other insights

**Boston University** CS & ECE

**BOSTON UNIVERSITY**

# Thank you

...any questions?

**Boston University** CS & ECE