# Supercomputing InfiniBand Fabric Analysis
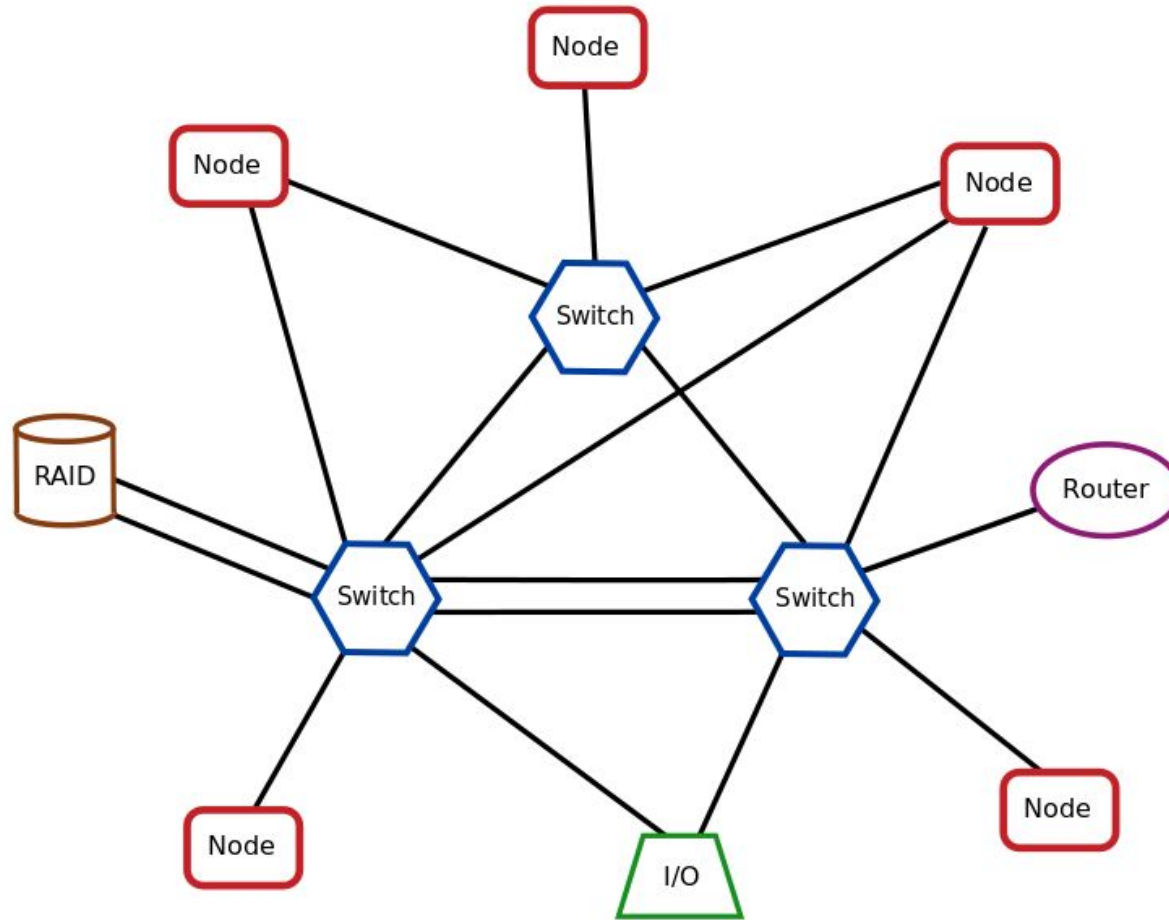
## Todd Yoder

**COLORADO**SCHOOLOF**MINES**
EARTH • ENERGY • ENVIRONMENT

SIParCS — Summer Internships in Parallel Computational Science — BOULDER • CO

# National Center for Atmospheric Research
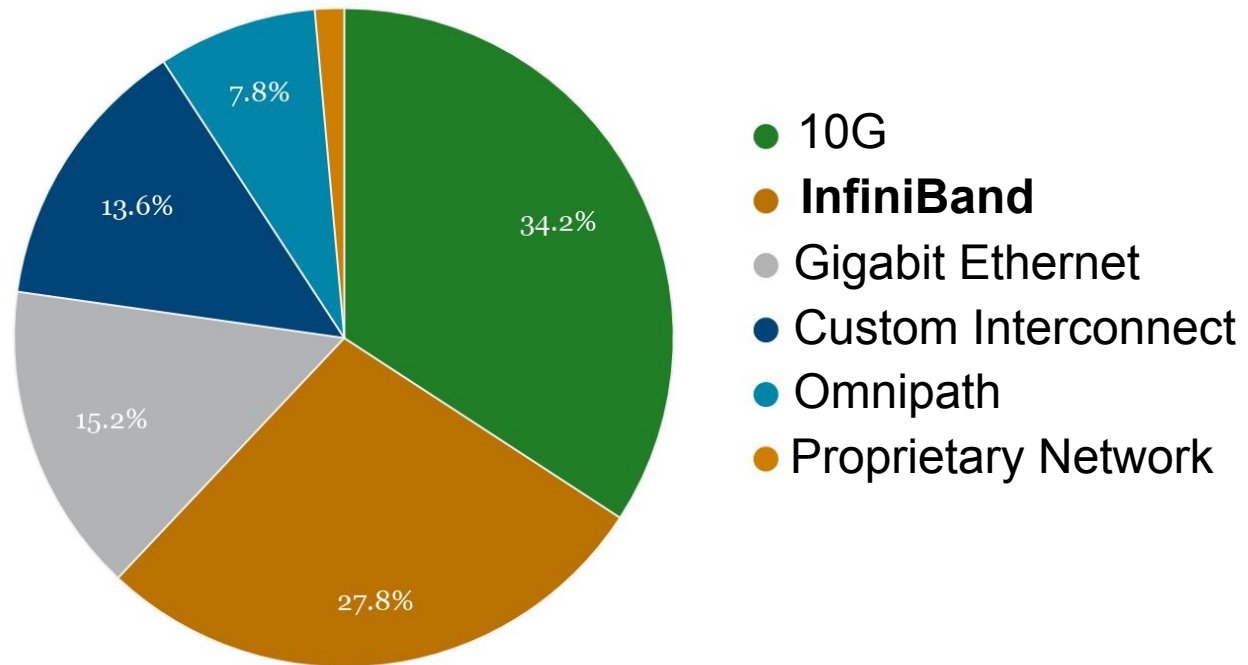
August 3, 2018

**NCAR**

NSF

# Introduction



Simple Supercomputer Fabric

# Introduction

InfiniBand is a computer-networking communications standard for high-performance computing.
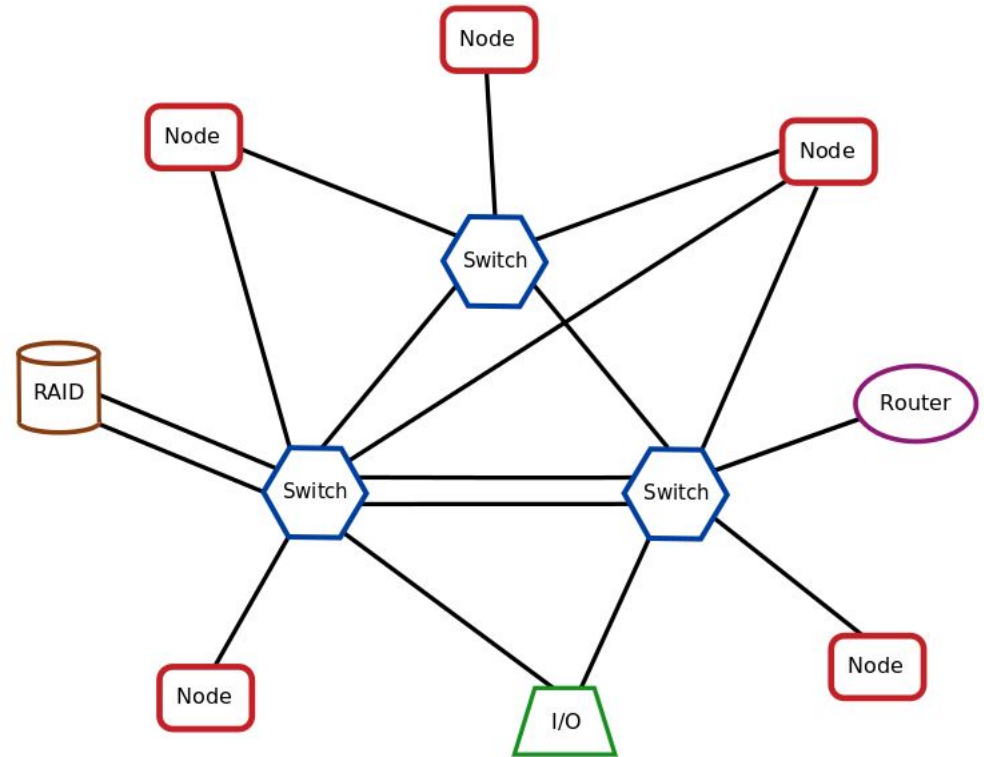
Interconnect Family System Share



- 10G — 34.2%
- **InfiniBand** — 27.8%
- Gigabit Ethernet — 15.2%
- Custom Interconnect — 13.6%
- Omnipath — 7.8%
- Proprietary Network

Interconnects used by the top 500 supercomputers[1]

# Supercomputing InfiniBand Fabric Analysis

## Goal

Develop software tools which analyze basic Graph Theory properties of an InfiniBand graph
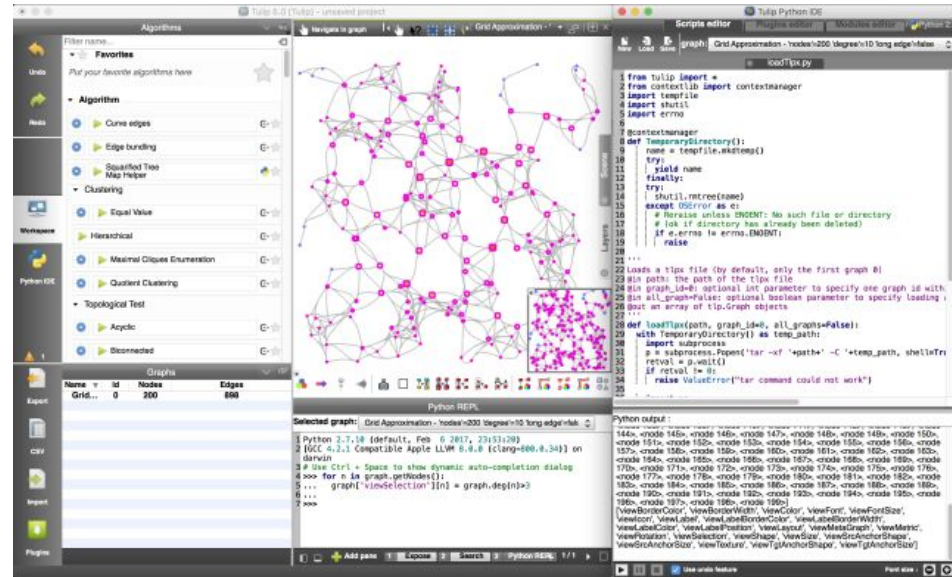


Simple Supercomputer Fabric

# **Tulip**

**Tulip** is a free information visualization framework for analyzing and visualizing relational data. It can be extended with plugins to analyze specific problems.
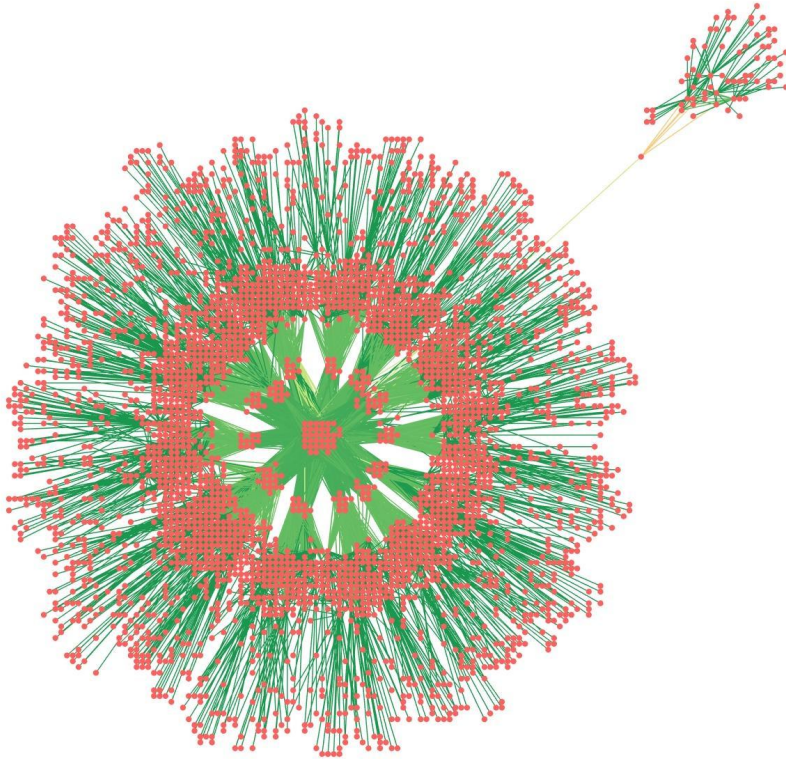
Features:
- 3D visualizations
- Automatic drawing of graphs
- Automatic clustering of graphs
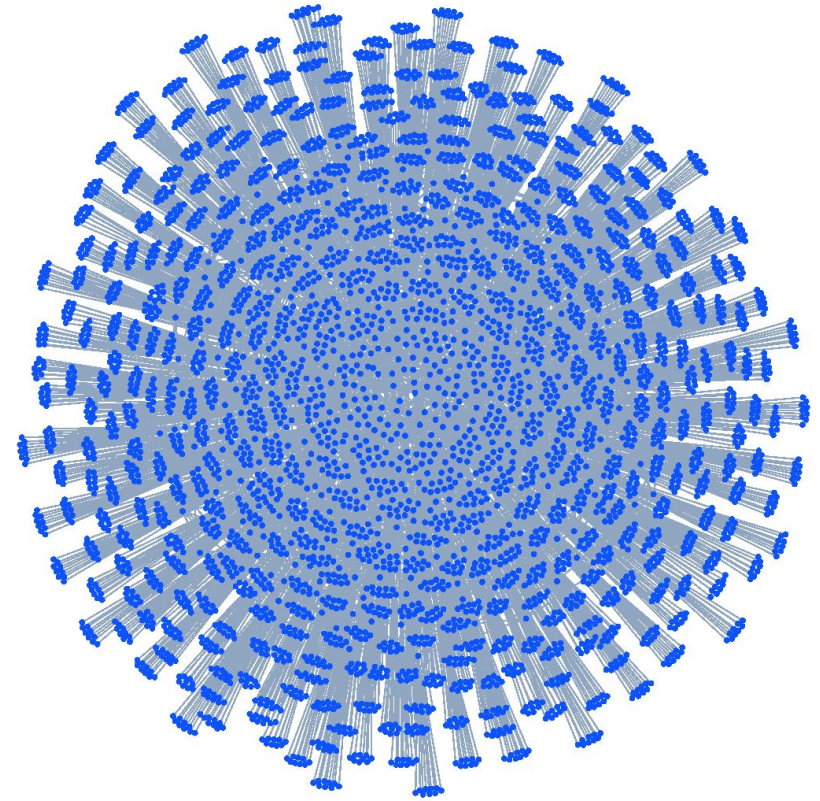- Automatic Metric coloration of graphs
- Open Source
- Free
- Written in C++



Sample screenshot of Tulip's graphic user interface[2]

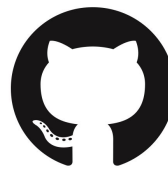# Tulip



NCAR's Yellowstone supercomputer,
a full fat tree[3]

NCAR's Cheyenne supercomputer,
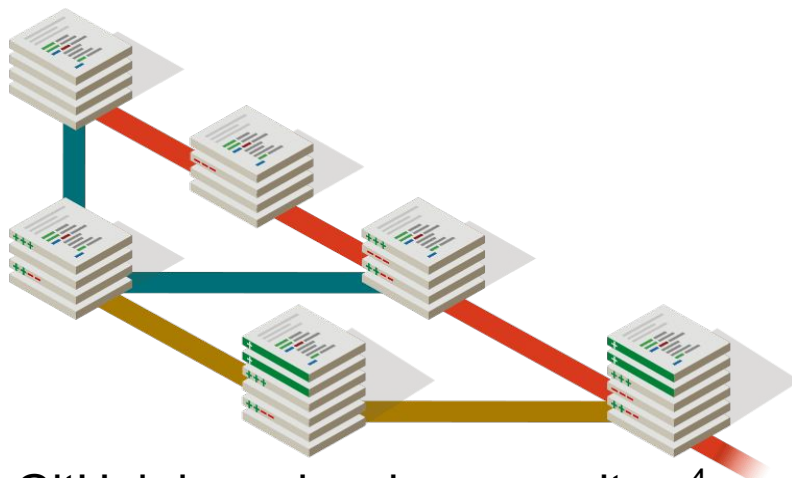a partial 9D enhanced hypercube

# GitHub ⬛ git ◆

https://github.com/NCAR/tulip_infiniband

GitHub: collaboration manager and web-based hosting service for git
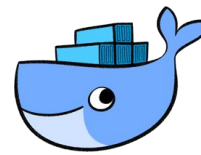
git: version control



GitHub branches in a repository[4]

# Docker

### Makefile  Dockerfile  Readme

Download

1. Install Docker, gcc

2. Download Tulip Infiniband docker folder

3. $ make

## Docker Container

Tulip
Infiniband

Tulip

Prerequisites

ubuntu

# Plugins Developed

## Random Nodes

Selects two random nodes on the graph

*Specific Application:*

Used by other plugins



Laramie: a 3D hypercube test and research supercomputer at NCAR

# Plugins Developed

## Shortest Path

Applies Dijkstra's Algorithm to one of the nodes.

Selects a shortest path between the nodes

*Specific Application:*

● Find routes nodes ought to use to communicate.
● Compare optimal routes with actual routes

A shortest path between two nodes on Laramie

# Plugins Developed

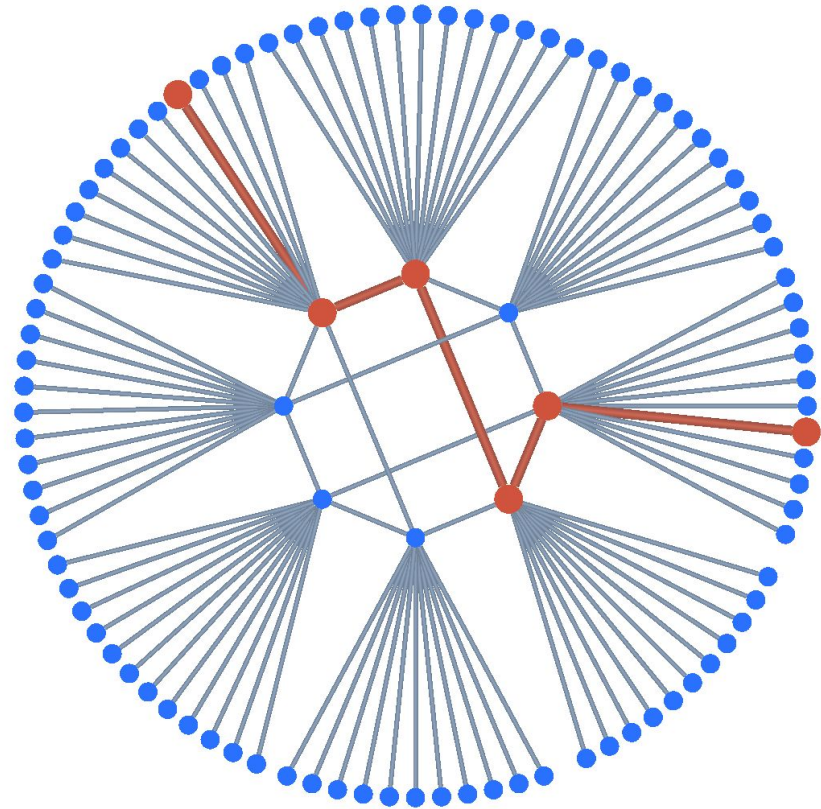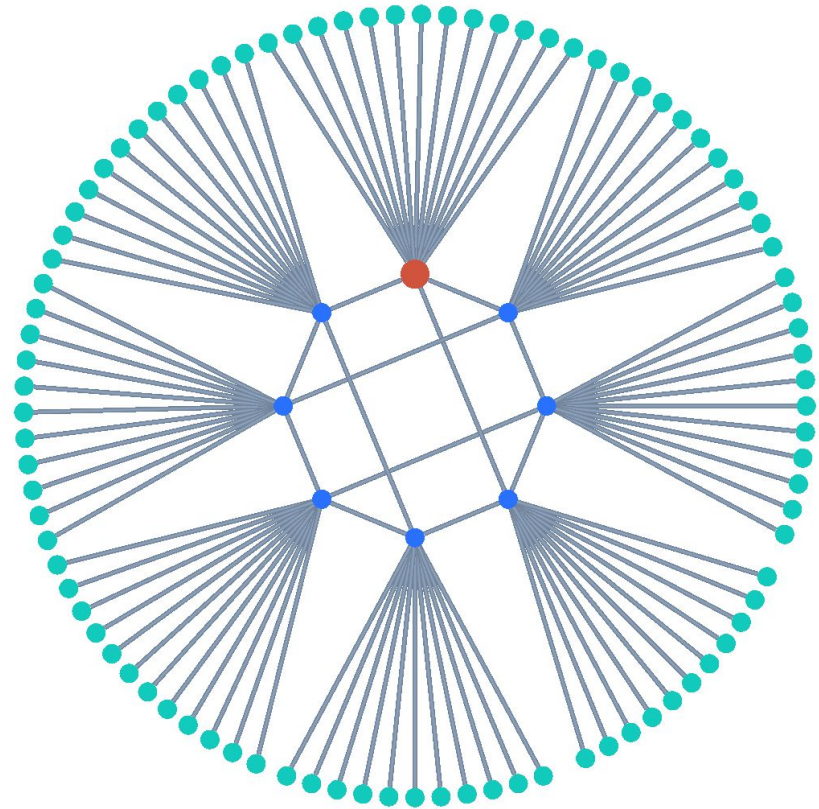## Min Degree and Max Degree

Prints smallest and largest node degrees, respectively

Selects corresponding nodes and prints their node IDs

*Specific Application:*

- Determine where network congestion is likely to occur
- Minimize number of cables in supercomputer while maintaining communication capabilities



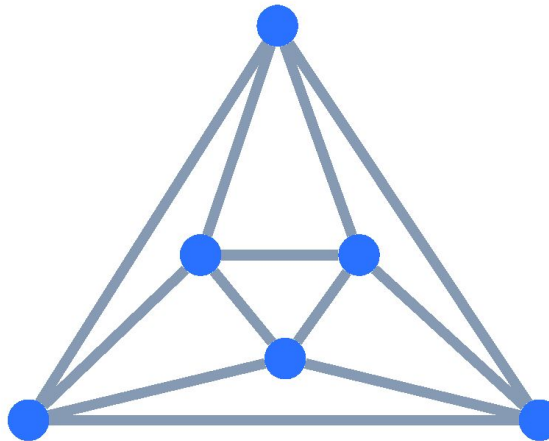Largest degree: 42, Smallest degree: 2

# Plugins Developed

## Regularity Test
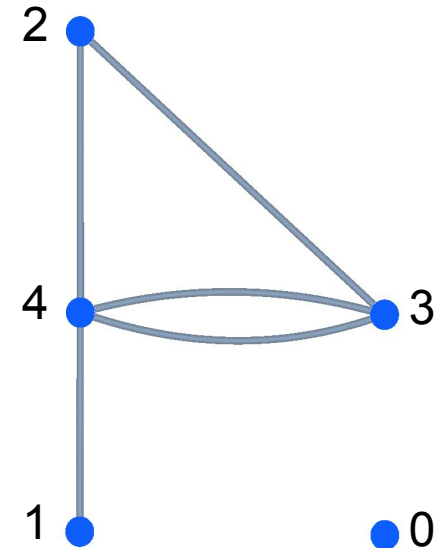
Regular Graph: all nodes have the same degree

Irregular Graph: each node has a unique degree

*Specific Application:*

Determine if switches are not symmetric

A regular graph.
Each node has degree 4

An irregular graph with degrees labeled

# Plugins Developed

## **<u>Bipartite Test</u>**

Bipartite Graph: The nodes can be partitions into two subsets such that every edge connects the two subsets

*Specific Application:*

Enables straightforward full-fabric bandwidth testing
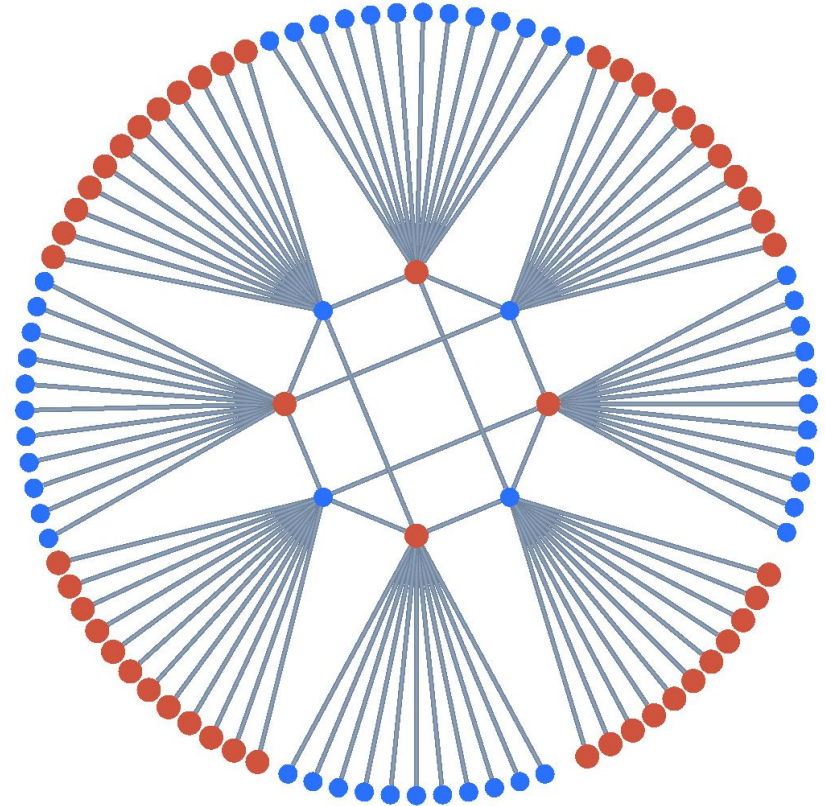


Laramie is bipartite

# Plugins Developed

## **Bipartite Test**

Bipartite Graph: The nodes can be partitions into two subsets such that every edge connects the two subsets

*Specific Application:*

Enables straightforward full-fabric bandwidth testing
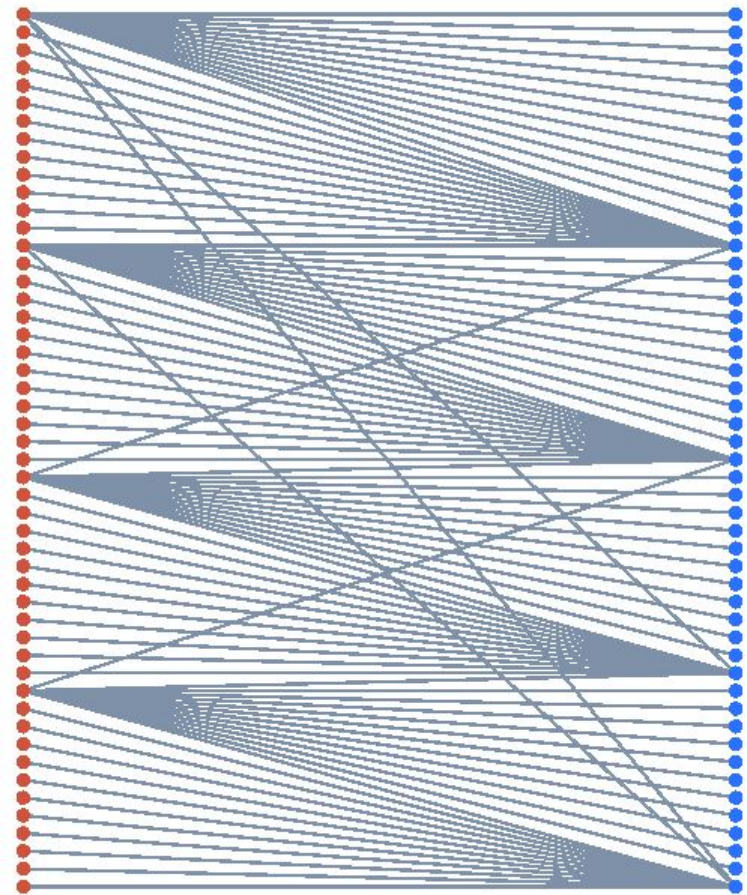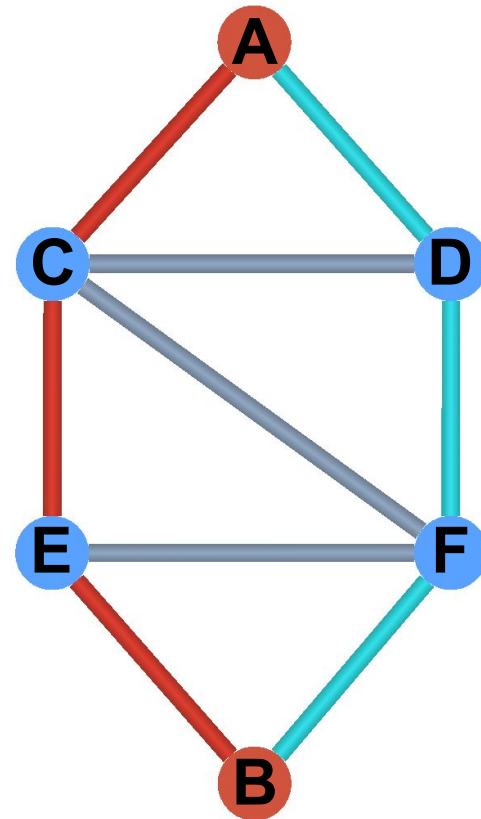


Laramie is bipartite

# Plugins Developed

## **Geodesic Test**

Geodesic Path: path of shortest length between two nodes

*Specific Application:*

- Fabrics need redundancy. It's useful to check that more than one optimal paths exist between nodes
- Helps check for excessive cables



Three geodesic paths from A to B:
red, blue, and ACFB
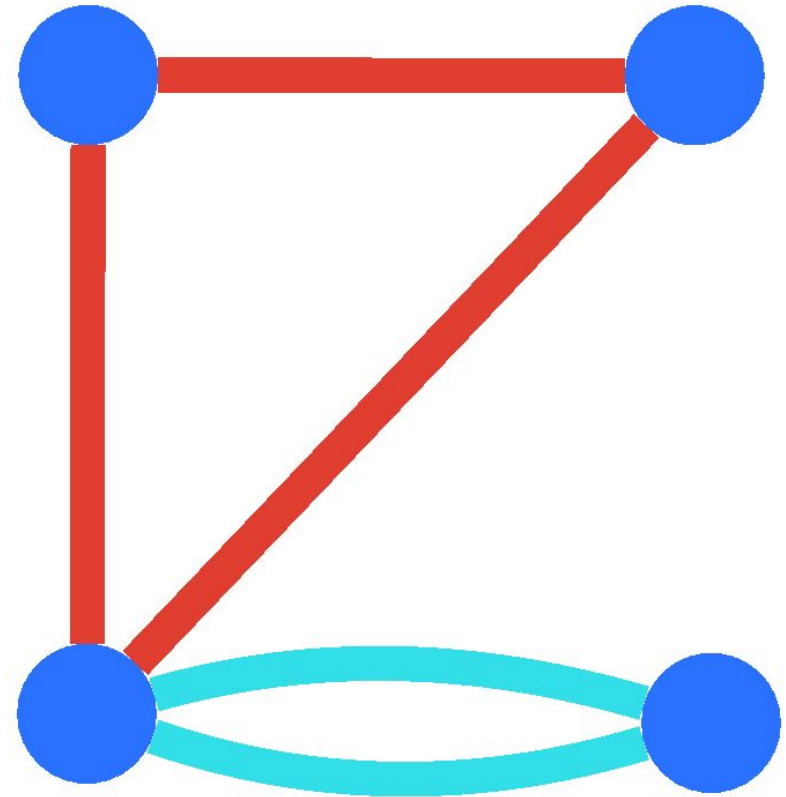
# Plugins Developed

## Node On Cycle Test

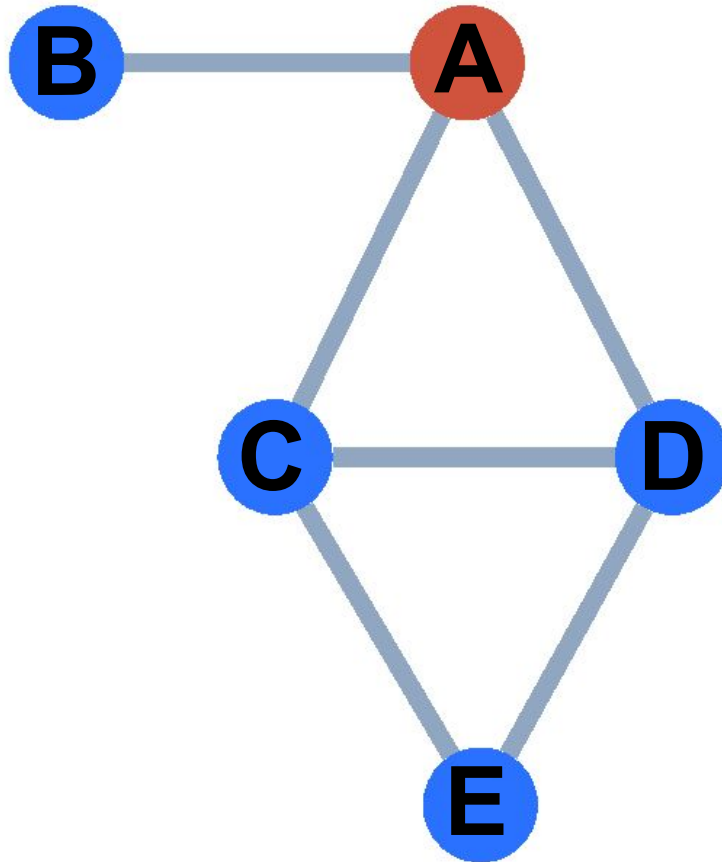Determines if the selected node lies on a cycle

*Specific Application:*

Multicast communications need to be aware of cycles to guard against inefficiencies and infinite loops



Graph with two cycles
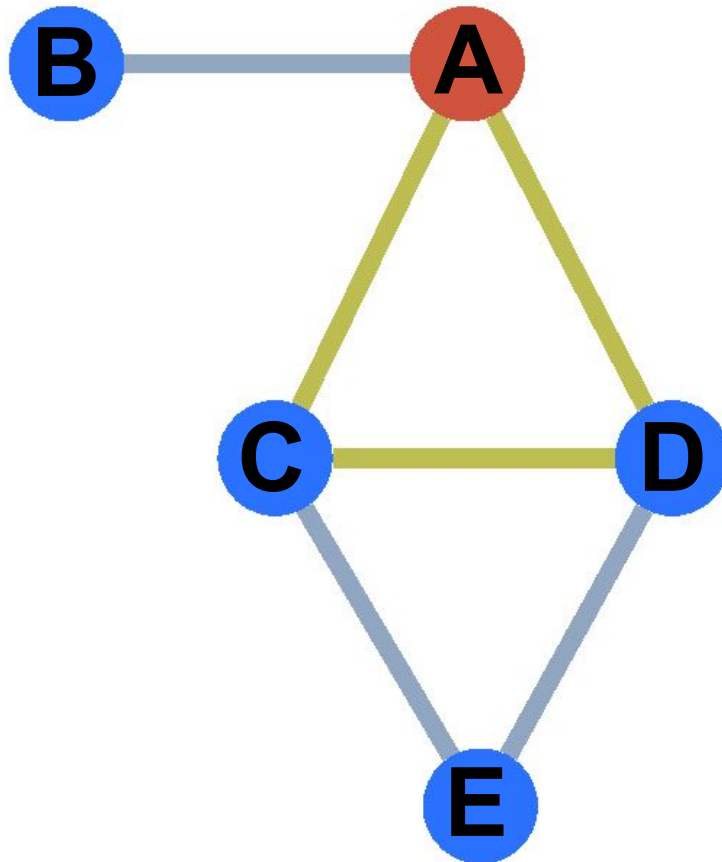
# Plugins Developed

## Node On Cycle Test



Multicast: Send message to multiple nodes, they store and pass on the message
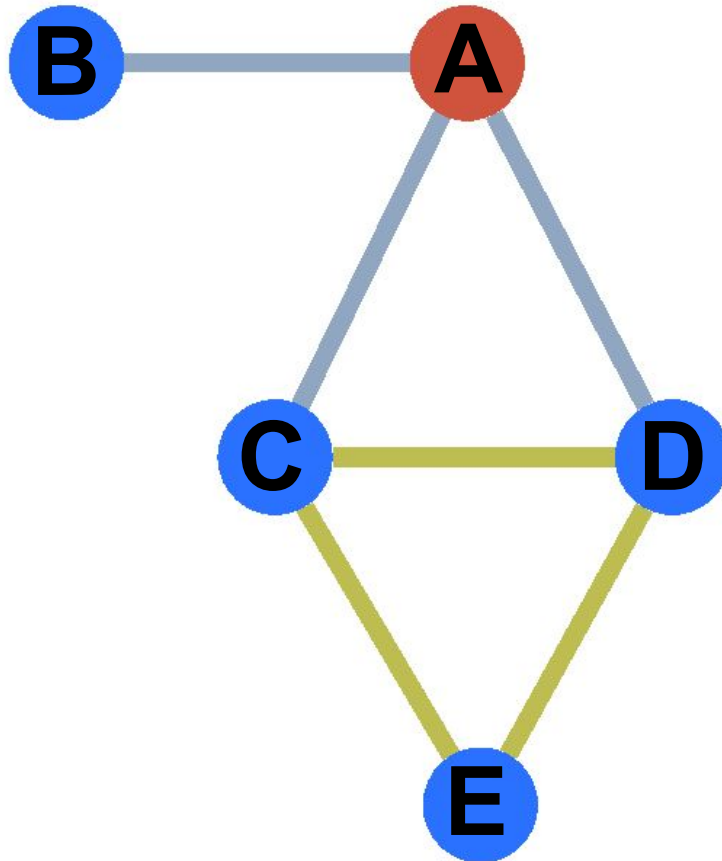
# Plugins Developed

## Node On Cycle Test



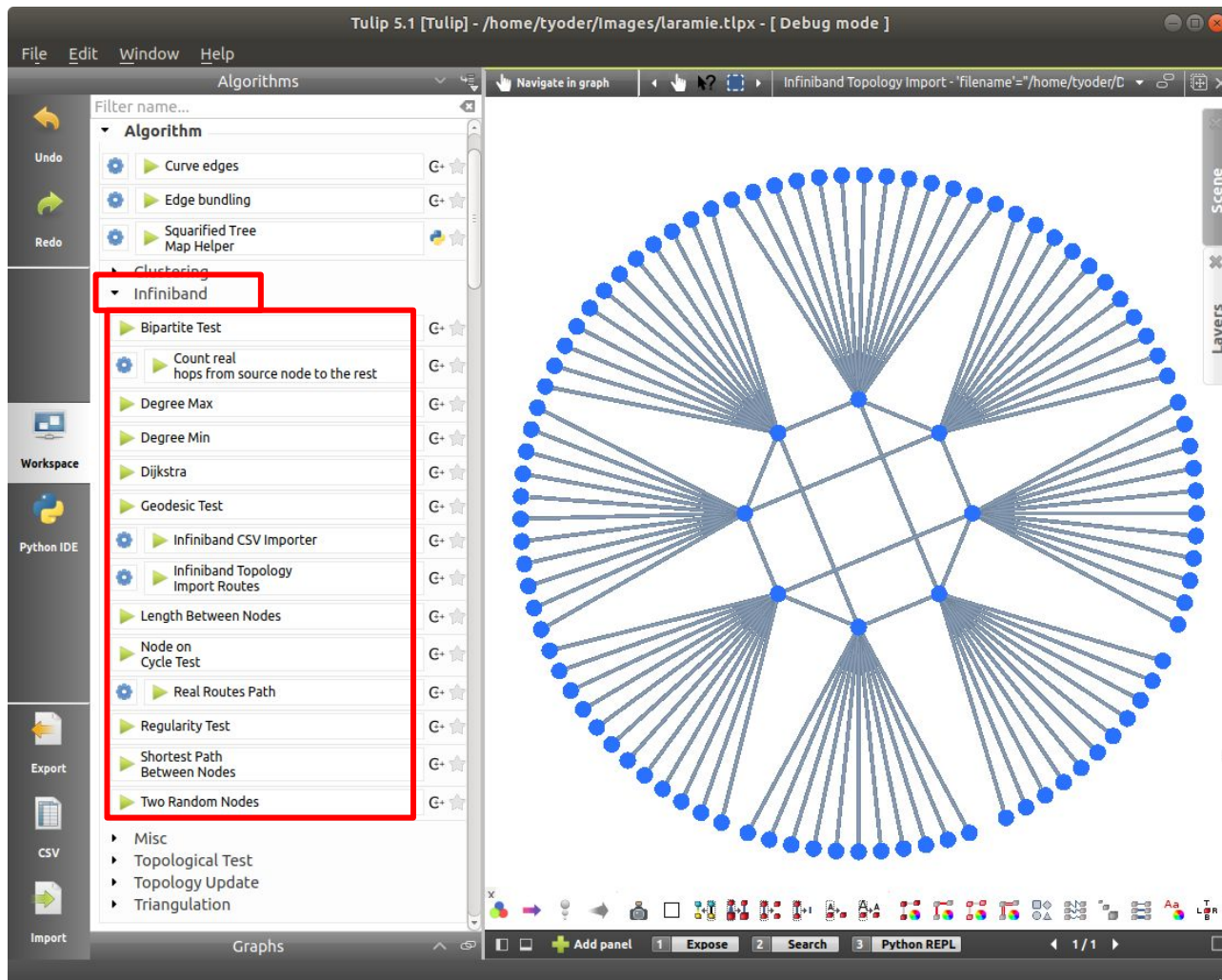Route ACD is inefficient

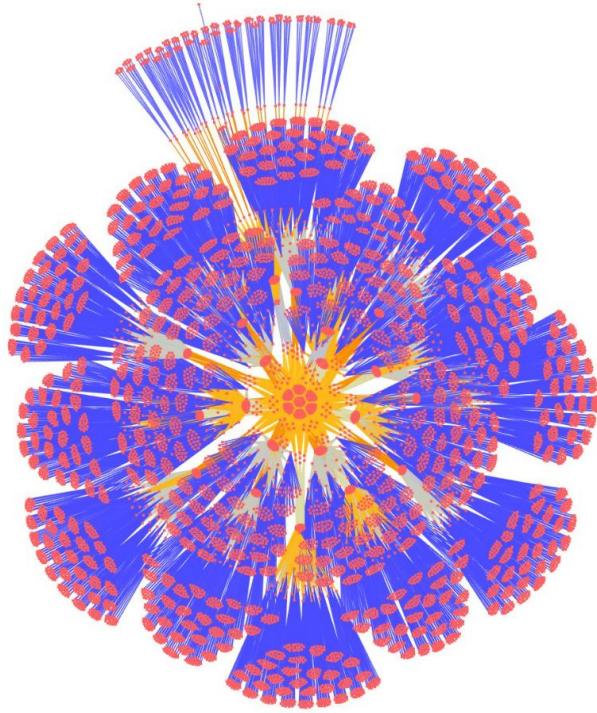# Plugins Developed

## Node On Cycle Test
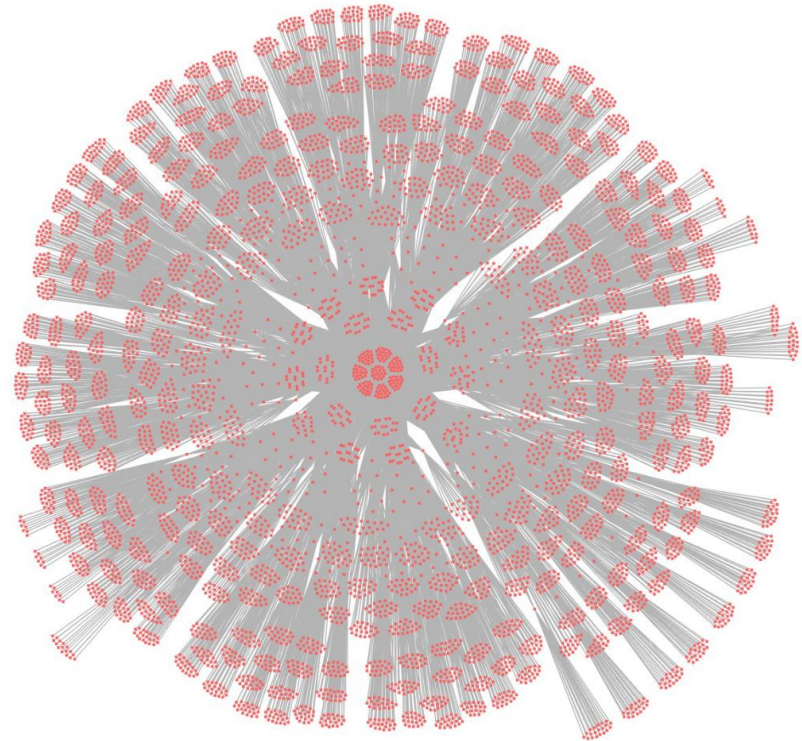


Route CED is an infinite loop!

# Using Tulip Infiniband

# Using Tulip Infiniband



SuperMUC, a supercomputer operated by Leibniz Supercomputing Center



Stampede, a supercomputer operated by Texas Advanced Computing Center until 2017

# Conclusions

Tulip Infiniband can help supercomputer development teams such as SSG make more informed decisions for upgrades, and it provides basic tools for maintenance and performance optimization.

# Future Work

- Write plugins for other graph theory properties

- Convert Dockerfile to Charliecloud or Singularity

- Write plugin which generates a summary of the graph by calling other plugins

# Acknowledgements

- Auber, D., & Mary, P. (2018). Tulip (Version 5.2) [Computer software]. Bordeaux, France: LaBRI, University of Bordeaux I.
- Chartrand, G., & Zhang, P. (2005). *Introduction to graph theory*. Boston: McGraw-Hill Higher Education.
- Futral, W. T. (2002). *InfiniBand architecture development and deployment: A strategic guide to server I/O solutions*. Hillsboro, OR: Intel Press.

# Questions?

## https://github.com/NCAR/tulip_infiniband

# Backup Slides

# Compatibility

Mac doesn't play nice with Graphical User Interfaces in Docker.

XQuartz bridges the gap to provide a GUI through the IP address.

**Linux**

1. Install Docker, gcc

2. Download Tulip Infiniband Docker folder

3. $ make

**Mac**

1. Install Docker, gcc, XQuartz

2. Download Tulip Infiniband Docker folder

3. $ make

# The Dockerfile

**1** **Load Ubuntu image**

Docker provides images with many popular operating systems

**2** **Install Prerequisites**

Tulip and the plugins depend on about two dozen libraries

**3** **Install Tulip**

Tulip is available at https://github.com/Tulip-Dev/tulip

**4** **Install libibautils**

Imports InfiniBand fabric into Tulip. Developed at NCAR
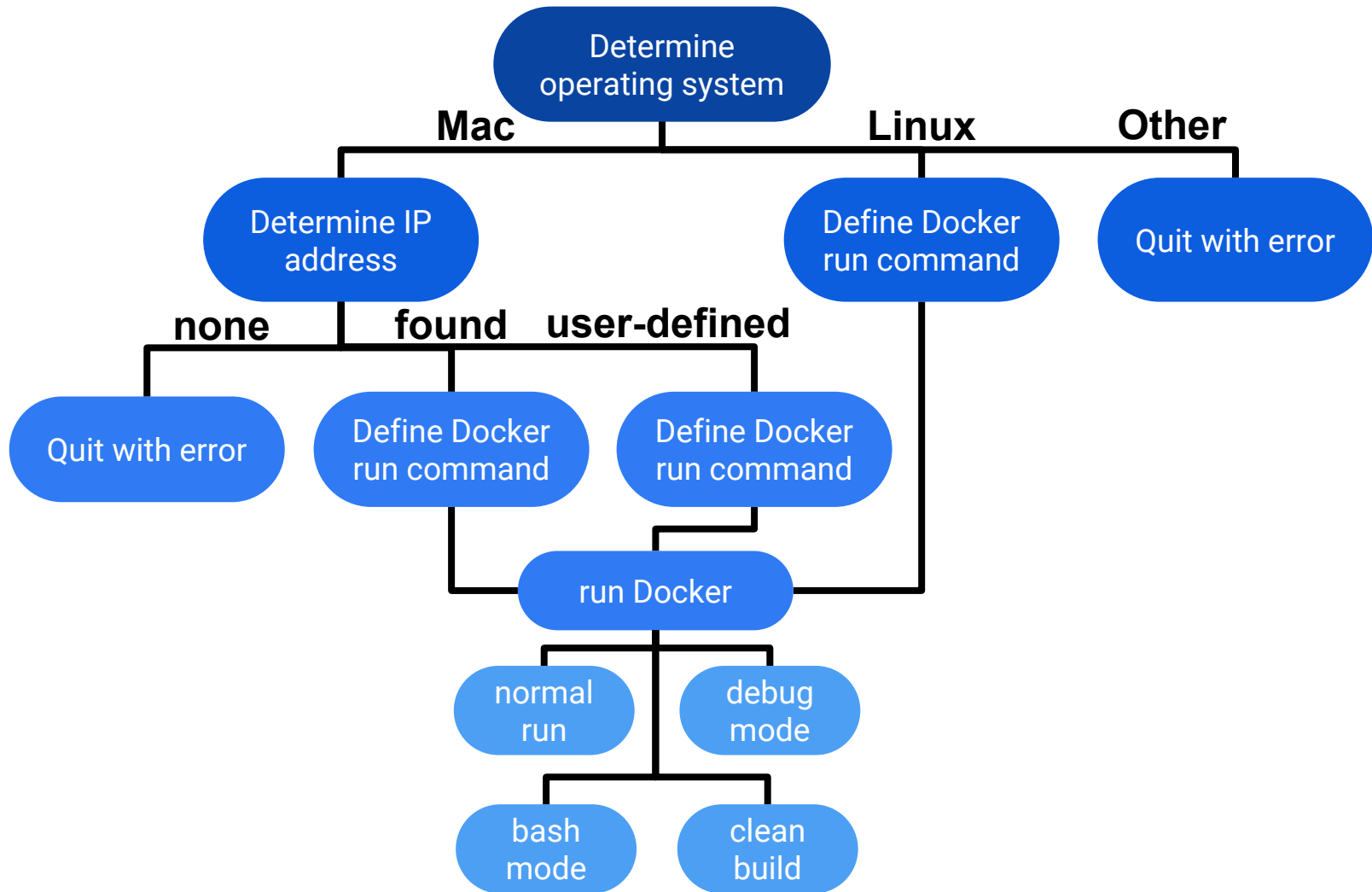
**5** **Install Tulip Infiniband**

Analysis plugins developed for InfiniBand

**6** **Run Tulip**

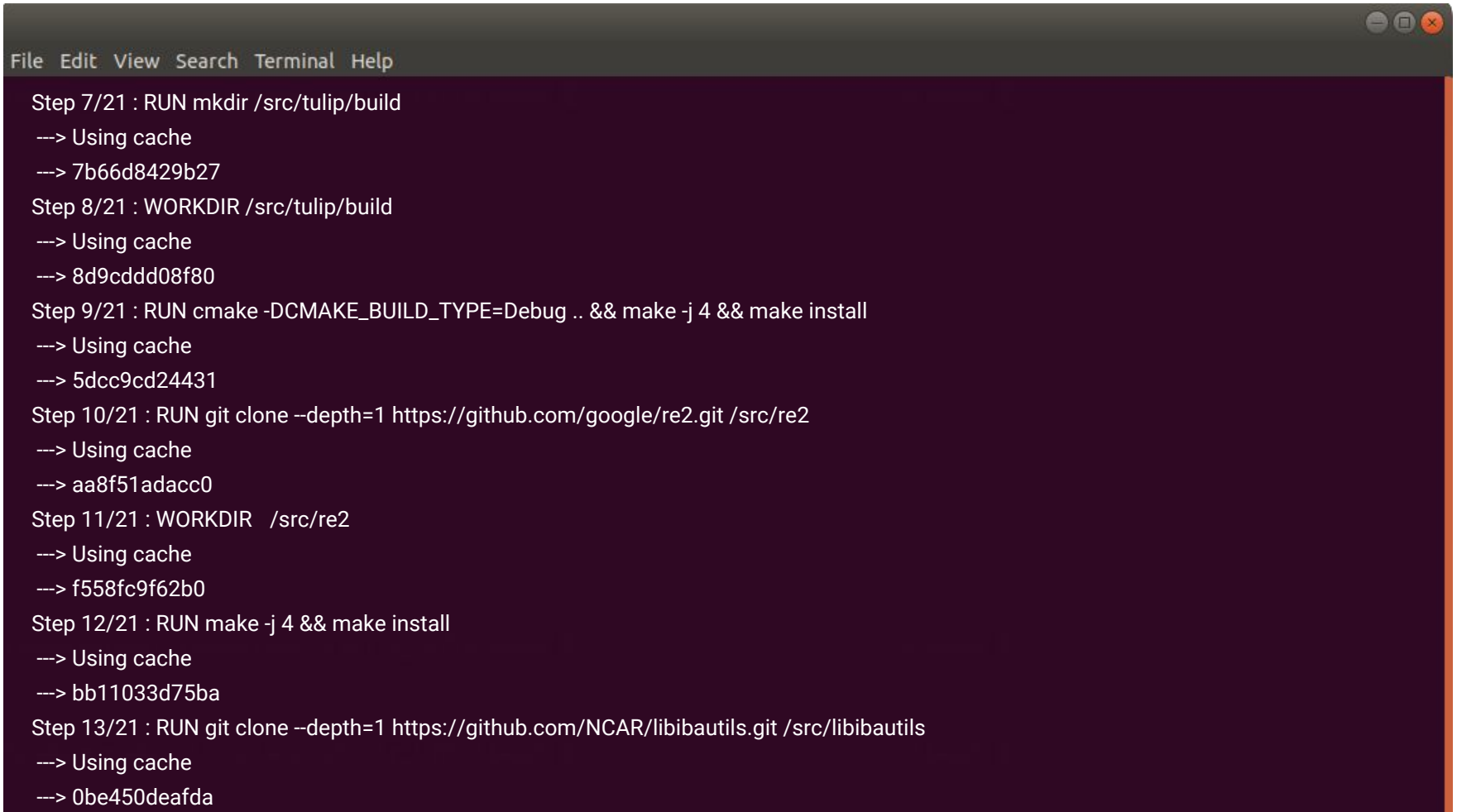Tulip launches with libibautils and Tulip Infiniband plugins

# The Makefile

# Dockerfile

# Dockerfile



```
Step 7/21 : RUN mkdir /src/tulip/build
 ---> Using cache
 ---> 7b66d8429b27
Step 8/21 : WORKDIR /src/tulip/build
 ---> Using cache
 ---> 8d9cddd08f80
Step 9/21 : RUN cmake -DCMAKE_BUILD_TYPE=Debug .. && make -j 4 && make install
 ---> Using cache
 ---> 5dcc9cd24431
Step 10/21 : RUN git clone --depth=1 https://github.com/google/re2.git /src/re2
 ---> Using cache
 ---> aa8f51adacc0
Step 11/21 : WORKDIR   /src/re2
 ---> Using cache
 ---> f558fc9f62b0
Step 12/21 : RUN make -j 4 && make install
 ---> Using cache
 ---> bb11033d75ba
Step 13/21 : RUN git clone --depth=1 https://github.com/NCAR/libibautils.git /src/libibautils
 ---> Using cache
 ---> 0be450deafda
```

NCAR | air • planet • people

# Dockerfile



```
File  Edit  View  Search  Terminal  Help

Step 14/21 : RUN mkdir /src/libibautils/build
 ---> Using cache
 ---> 1894569cc73d
Step 15/21 : WORKDIR /src/libibautils/build
 ---> Using cache
 ---> ce6e3bf7bebc
Step 16/21 : RUN cmake -DCMAKE_BUILD_TYPE=Debug .. && make -j 4 && make install
 ---> Using cache
 ---> 7229aeefadbf
Step 17/21 : RUN git clone --depth=1 https://github.com/NCAR/tulip_infiniband.git /src/tulip_infiniband
 ---> Using cache
 ---> 4577cc6dc57c
Step 18/21 : RUN mkdir /src/tulip_infiniband/build
 ---> Using cache
 ---> 0624b01bc6db
Step 19/21 : WORKDIR /src/tulip_infiniband/build
 ---> Using cache
 ---> 4e323a4a6a0f
Step 20/21 : RUN cmake -DCMAKE_MODULE_PATH="/src/tulip/cmake;/src/tulip_infiniband"        -DCMAKE_BUILD_TYPE=D
/src/tulip_infiniband && make -j 4 && make install
 ---> Using cache
```

# Dockerfile

```
---> dae84f4415c1
Step 21/21 : CMD env LD_LIBRARY_PATH=/usr/local/lib/tulip/:/usr/local/lib:$LD_LIBRARY_PATH tulip
 ---> Using cache
 ---> cffdc673553e
Successfully built cffdc673553e
Successfully tagged tulip:latest
```

## NCAR
*air · planet · people*

# Plugin Algorithms

## Bipartite Test

1. Set selected node as src
2. Place src in Group A
3. Place all neighbors of src in Group B
4. Place the neighbors' neighbors in Group A
5. Continue until all nodes are classified
6. Not bipartite if a node belongs to both groups

## Geodesic Test

1. Verify selected edges form a path
2. Count number of edges selected
3. Call Shortest Path plugin to get length between the end nodes
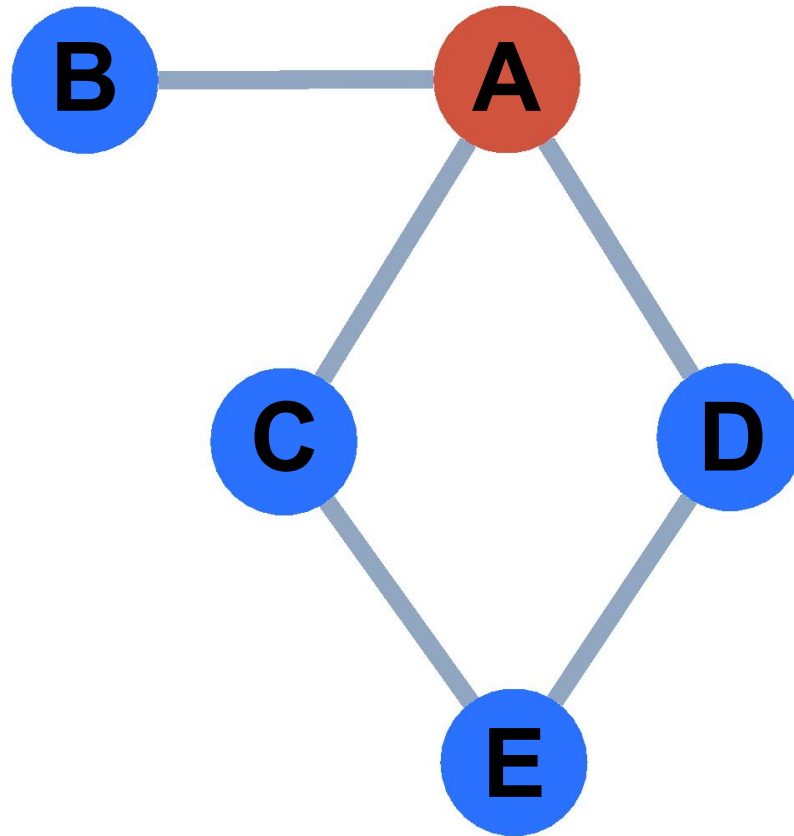4. Compare lengths

# Plugin Algorithms

## **Node On Cycle Test**

1. Store selected node as src
2. Check src degree. False if degree < 2
3. Check for self-loops
4. Check for two edges connecting src to the same neighbor
5. For node N, beginning with src,
   a. Store N as parent of all parentless neighbors, unless src
   b. If a neighbor already has a parent, src is on a cycle if the paths to N and its neighbor are disjoint except for src
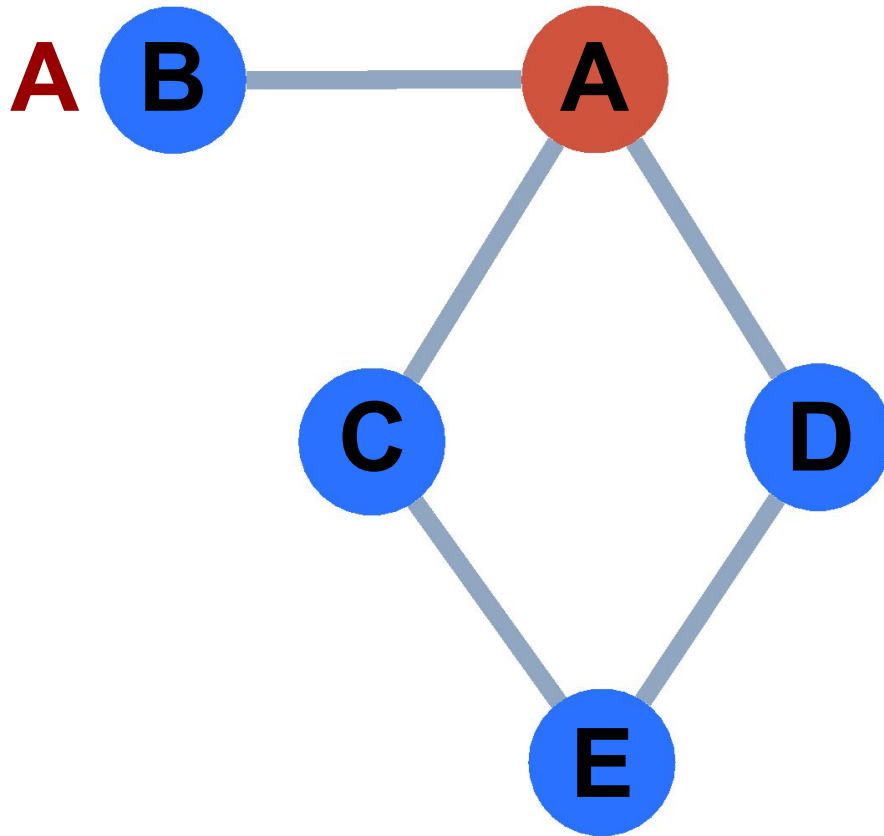   c. Add N's neighbors to queue to be considered

# Plugins Developed

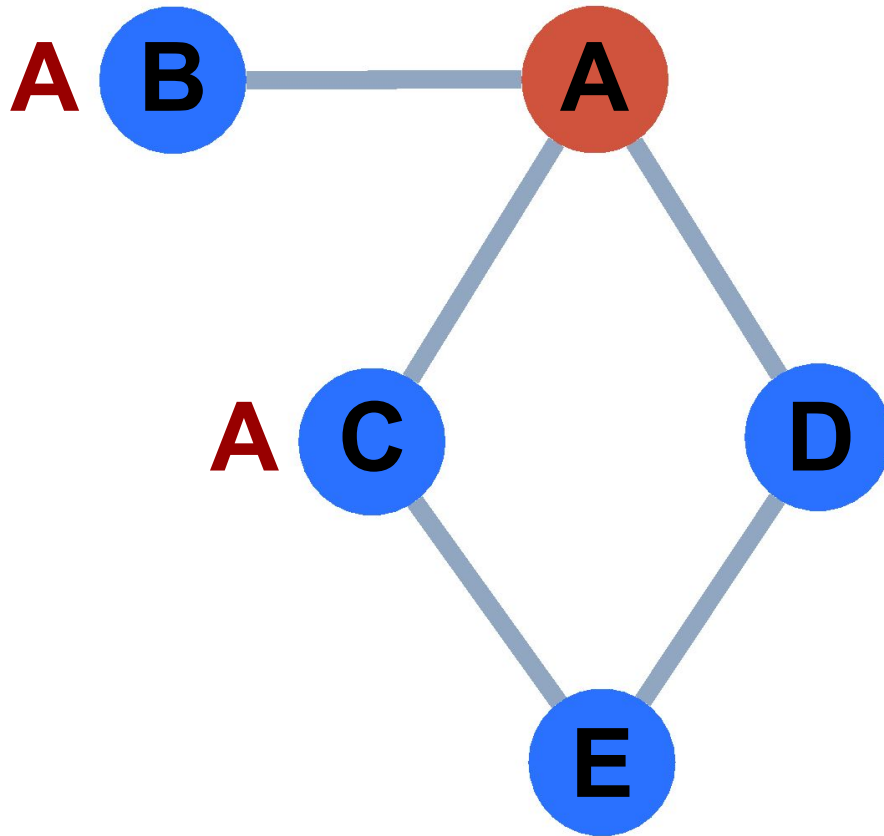## Node On Cycle Test



**Queue**
**A**

# Plugins Developed

## Node On Cycle Test



**Queue**
A
B

# Plugins Developed

## Node On Cycle Test
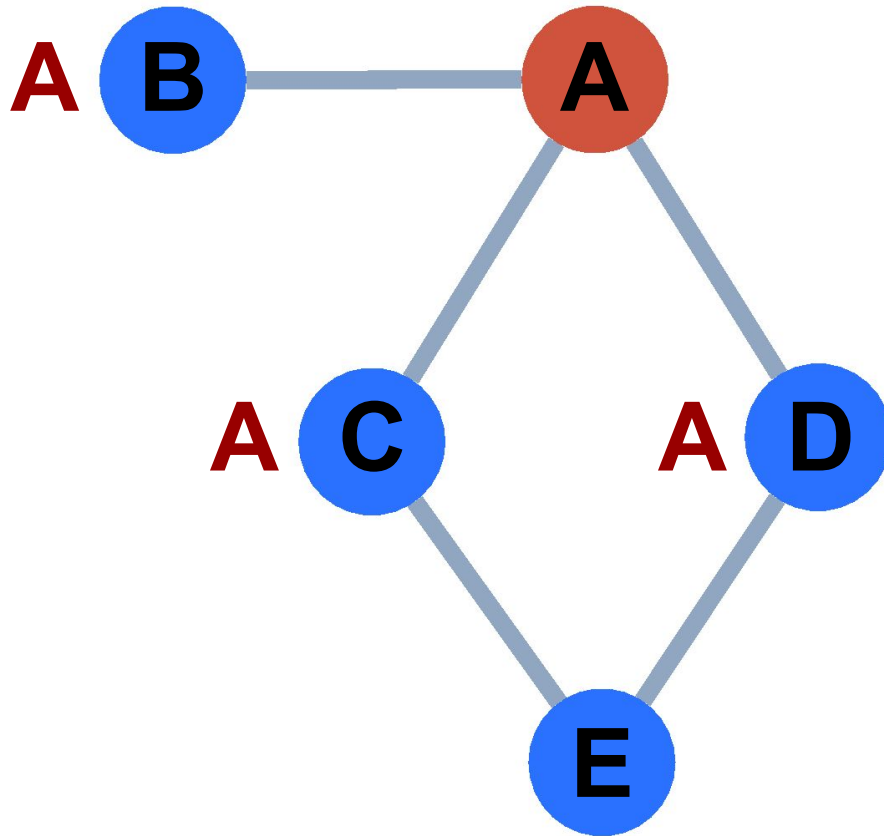


**Queue**
A
B
C

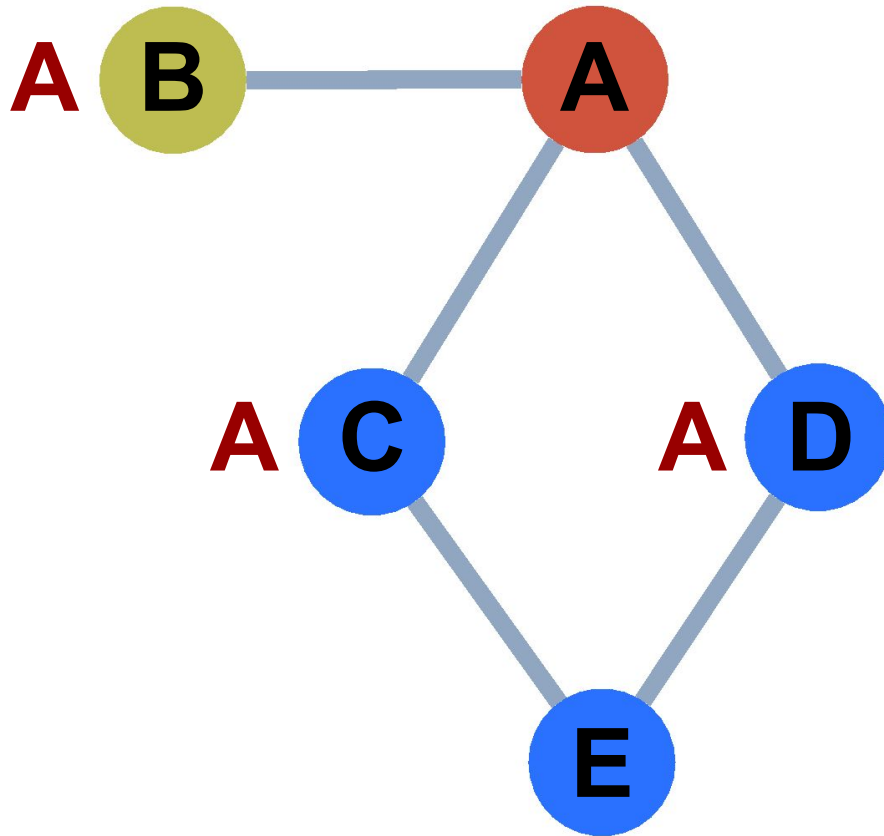# Plugins Developed

## Node On Cycle Test



**Queue**
A
B
C
D

# Plugins Developed

## Node On Cycle Test



Queue
B
C
D
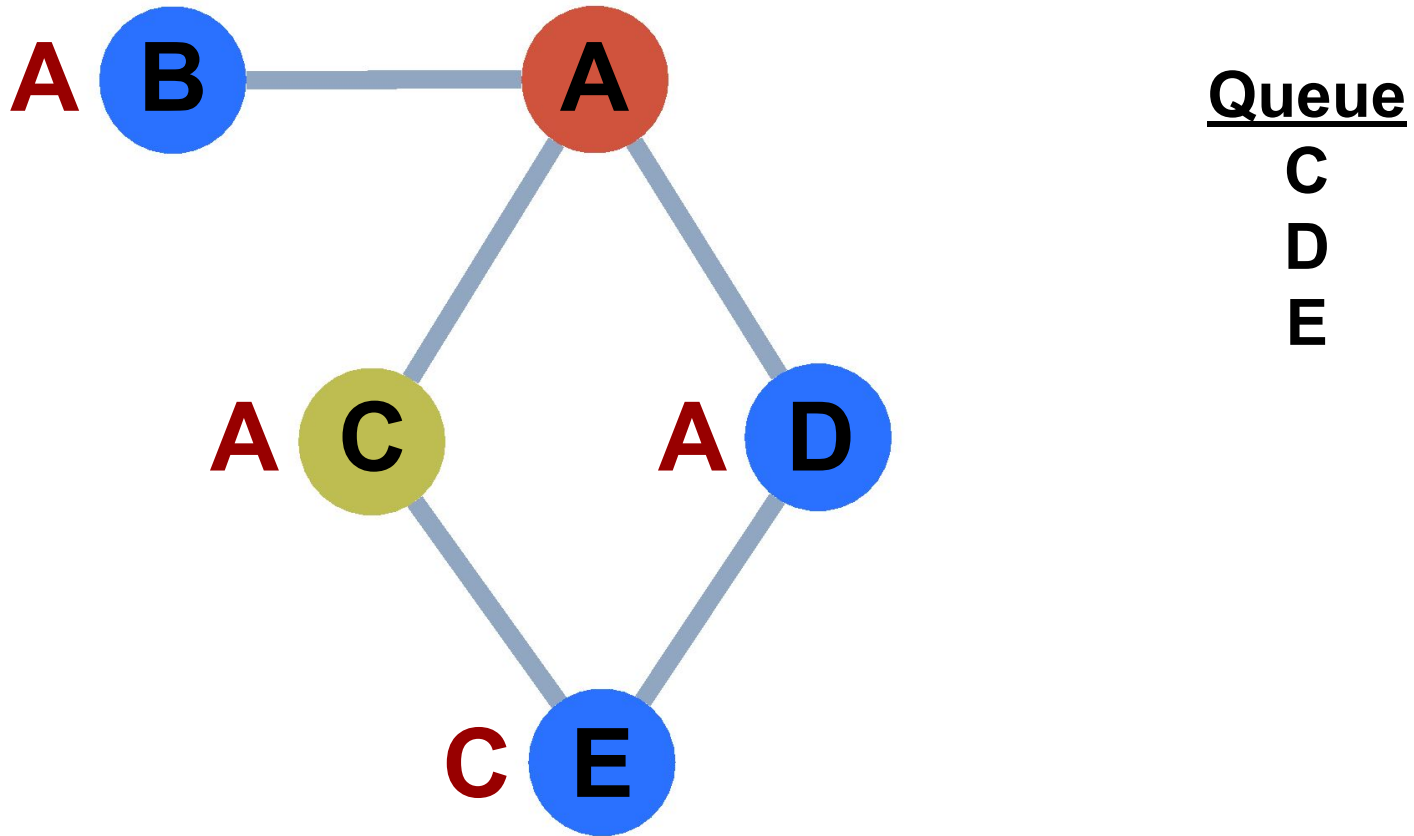
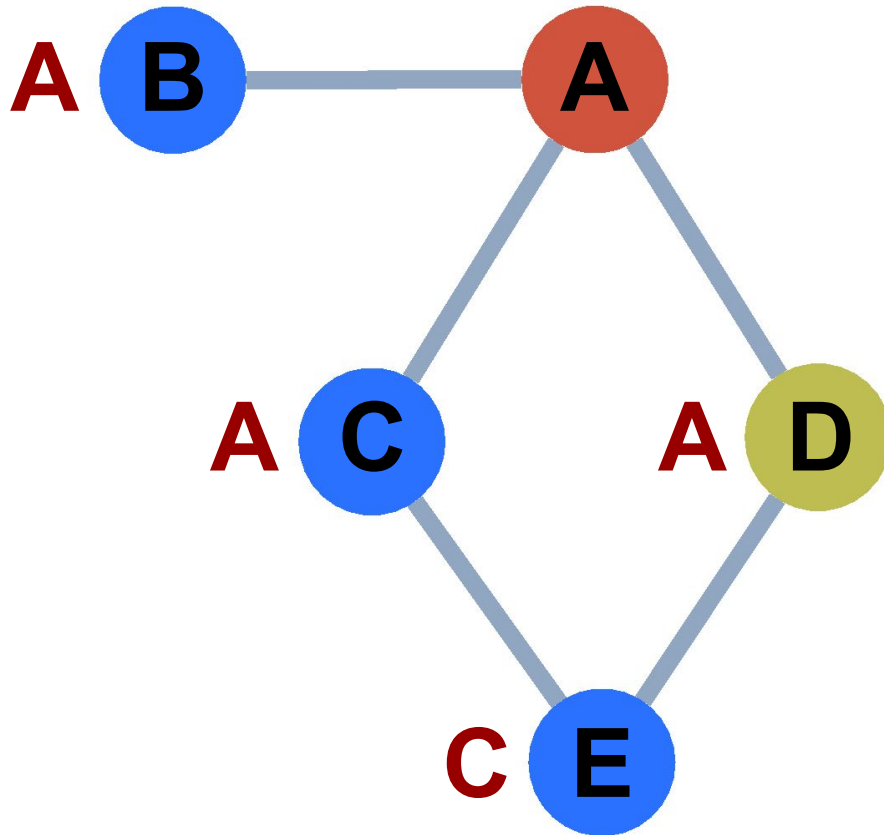# Plugins Developed

## Node On Cycle Test

# Plugins Developed

## Node On Cycle Test



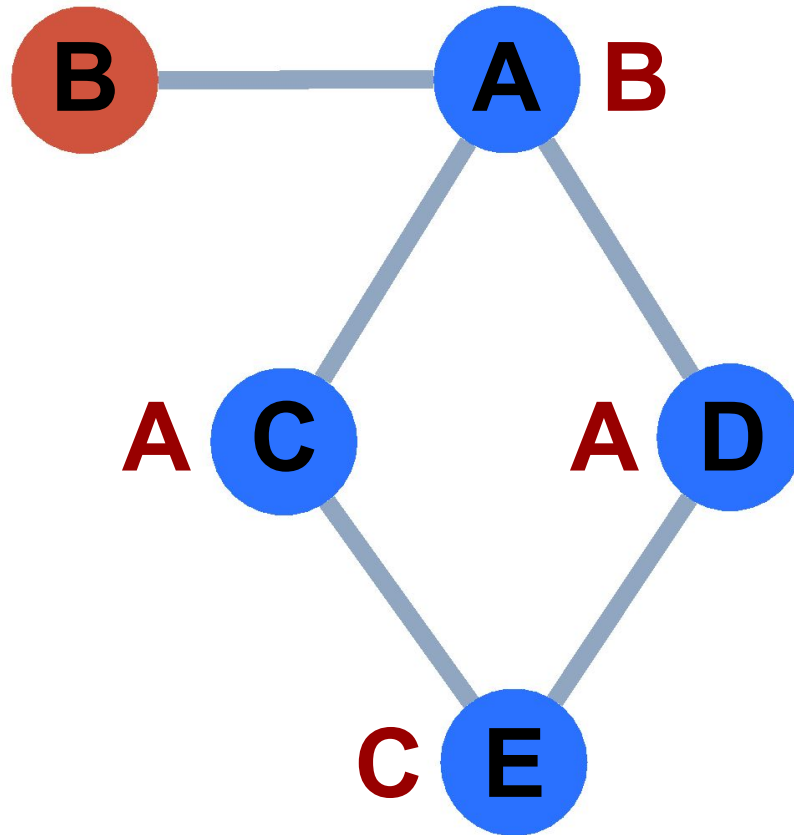**Queue**
D
E

**D-Path**
A

**E-Path**
C
A

Paths disjoint except for A implies
A is on a cycle

# Plugins Developed

## Node On Cycle Test



**D-Path**
A
B

**E-Path**
C
A
B

Paths share a node implies we found a cycle, but it doesn't include B

# Image Sources

[1]https://www.top500.org/statistics/list/

[2]Auber, D., & Mary, P. (2018). Tulip (Version 5.2) [Computer software]. Bordeaux, France: LaBRI, University of Bordeaux I.

[3]https://github.com/NCAR/tulip_infiniband

[4]https://arstechnica.com/gadgets/2017/11/microsoft-and-github-team-up-to-take-git-virtual-file-system-to-macos-linux/

Slide 8

- Makefile image: http://www.iconarchive.com/show/oxygen-icons-by-oxygen-icons.org/Mimetypes-text-x-makefile-icon.html
- Dockerfile image: https://www.iconsdb.com/black-icons/text-file-5-icon.html
- README image: https://findicons.com/search/readme
- Folder image: https://dumielauxepices.net/sites/default/files/folders-clipart-computer-folder-616425-7549368.png
- Puzzle piece: http://autism-works.org/wp-content/uploads/2013/12/2012-puzzle-piece.png

Slide 26:

- Bridge: http://pngimg.com/uploads/bridge/bridge_PNG12.png