# Autonomic Condor Clouds

David Wolinsky
ACIS P2P Group
University of Florida

# So What's the Big Deal

- Support connectivity across the Internet, in constrained locations, and with clouds

- Simplify packaging

- Minimize Condor configuration

- Reduce downtime

- Let's try to make this easy ....

# Discussion Goals

- The High-Level Overview

- Self-Configurable Condor Components

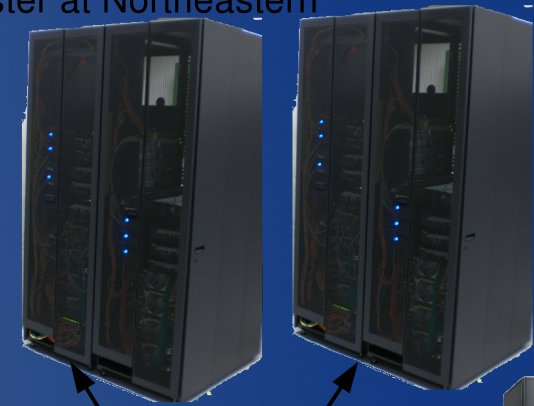- Virtual Networking

- The User Experience (Time allowing)

# High level Overview

- Self-contained VM Appliance

- Configuration through virtual floppy

- Fully distributed, decentralized Virtual Private Network via P2P Overlay

- Job scheduling via Condor

- Customization through Debian and Stacked File system

**Archer**

Cluster at Texas at Austin

Cluster at Northeastern

New user: startd, schedd, etc

Virtual network over the Internet

Cluster at Florida State

Future clusters at Florida, Minnesota, Cornell, and Northwestern

UF UNIVERSITY of FLORIDA ACIS

# Everything But Virtual Networking

- Virtual Machines (VMs)
    - Support for VMware, VirtualBox, KVM, and Xen
    - Homogenous environment on heterogenous resources
- Configuration
    - Provided by Virtual floppy, 25 KB download
    - P2P overlay info
    - Users identification certificates
    - Condor machine type
- Customization
    - Based upon Debian 4.0, access to apt repositories
    - Stacked file systems allow users to modify image and share only the changes
- Condor .... :-)

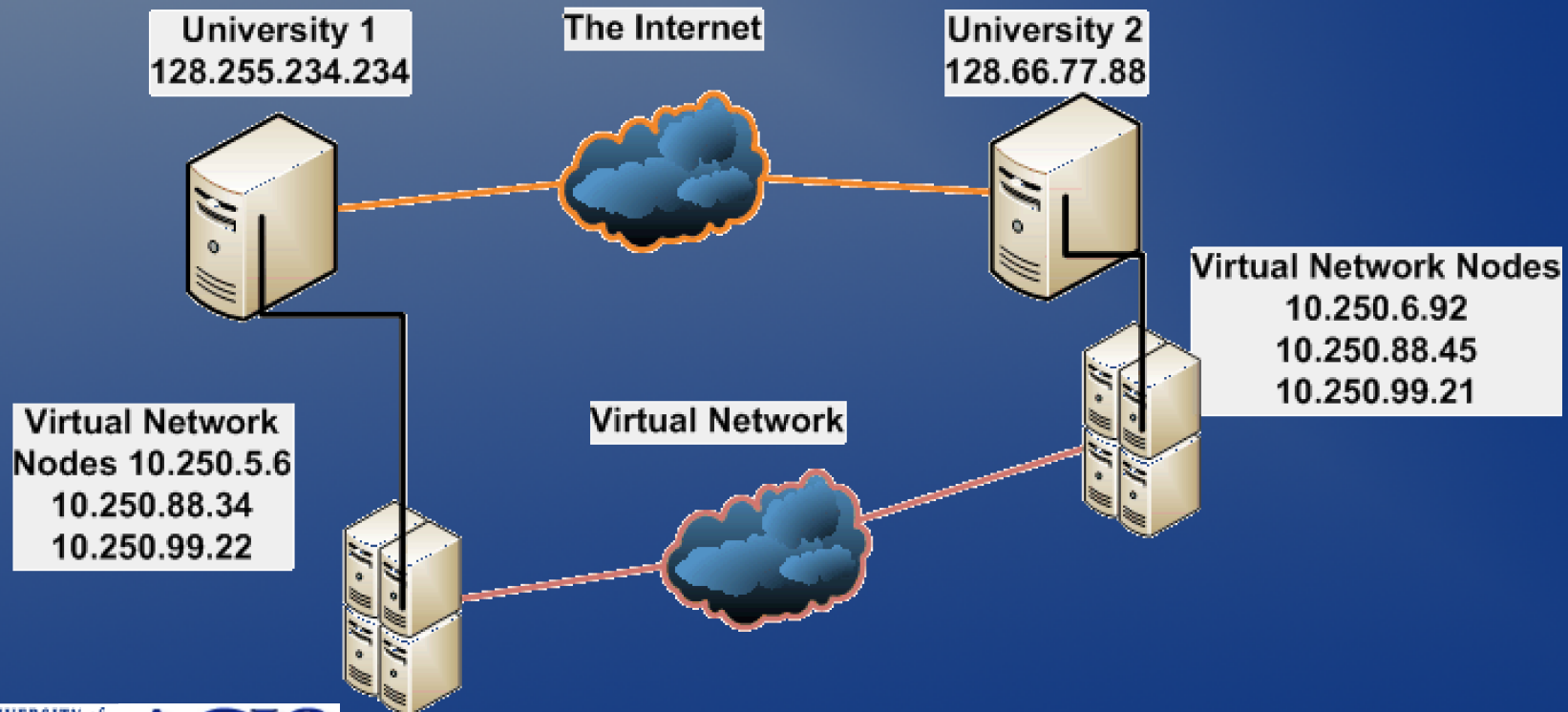# Configurable Components

- Daemons to start
    - Worker – execution only - startd
    - Client - submission and execution – startd and schedd
    - Master / Server – negotiator and collector
- Which Master / Server to connect or flock to
- User / Group Resource ownership and preemption
- Client can share files via autofs enabled NFS
- Monitoring and binding to a dynamic IP address

# More Opportunities

- Support for VM Universe
- Distributed data storage
- Web portal front-end
- Self-policing security system
- Self-sustaining condor cluster
- Portable (decentralized) Condor System Configurer
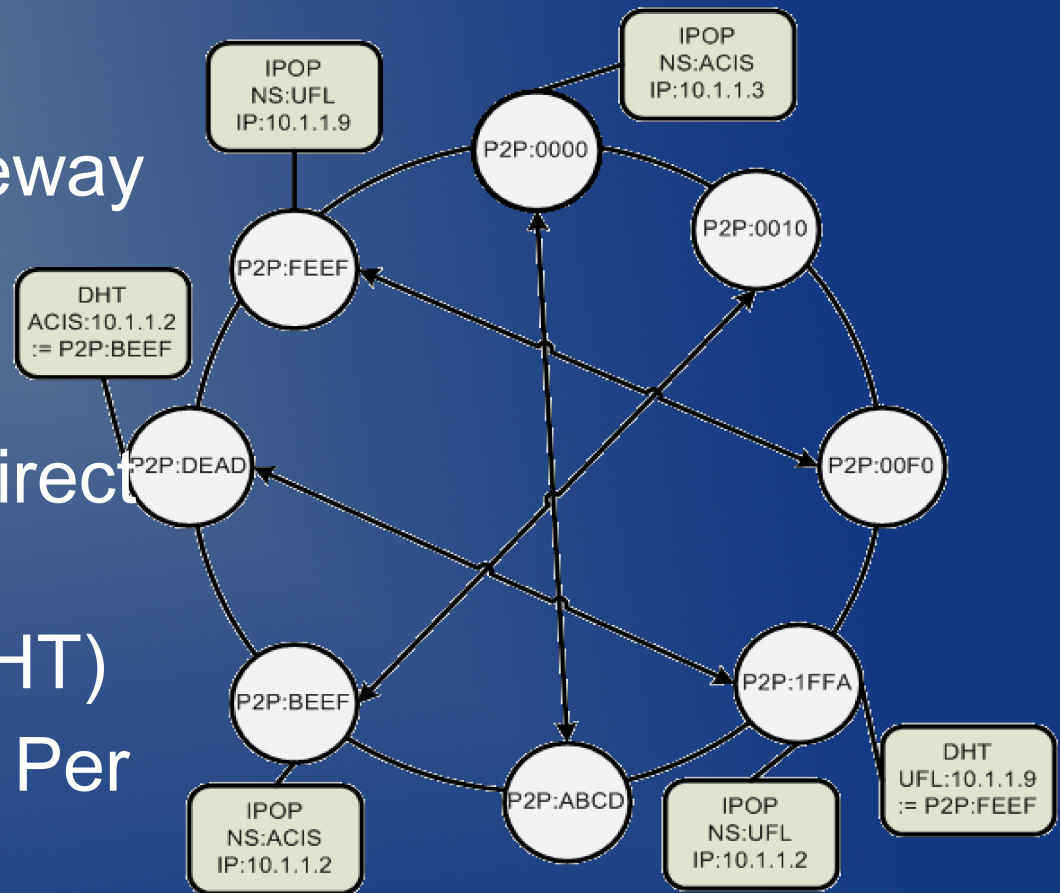- Configurable through Social Networks and User Portals

# Virtual Networking

- Unified layer 3 (IP) network for all machines
- Cross-site communication without a middleware broker

# IPOP

- Open Source
- NAT Traversal (STUN)
- Transparent Subnet Gateway
- Structured P2P Network Overlay
- Provides tunneling and direct shortcuts
- Distributed data store (DHT)
- Multiple Virtual Networks Per Overlay

# P2P Overlay

- Several hundred well distributed nodes
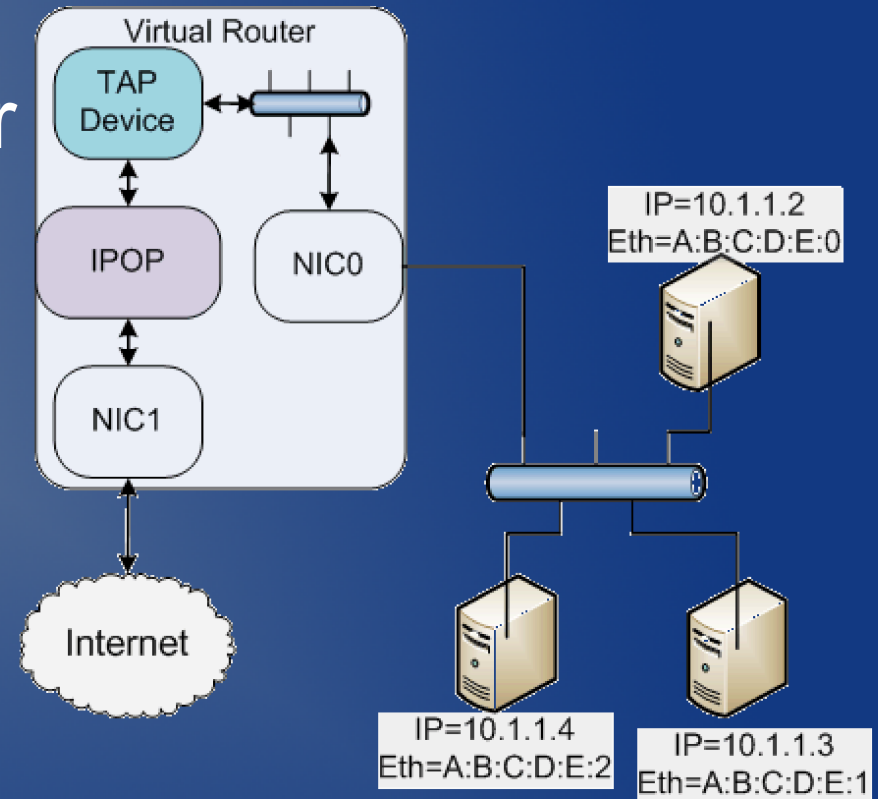- Assist in bootstraping and NAT traversal
- Runs on top of Planetlab

# VN Interfaces

- Each machine has VN Interface running locally
- Machine has VN and "Internet" connectivity

# VN Routers

- Single VN instance (Router) for entire cluster

- Limited to no resource configuration

- Isolated VN overhead

- May have "Internet" connectivity if there is an "Internet" router

# The User Experience

- Time permitting video - http://www.youtube.com/watch?v=1XDvITdayhs

- Otherwise slides –

  - Boot VM and obtain IP Addresses

  - Condor access

  - Direct connectivity (i.e. low ping overheads)

```
griduser@C240195038: /home/griduser

slot1@C185000234.i LINUX        INTEL   Unclaimed Idle      0.000   1485  0+03:15:05
slot2@C185000234.i LINUX        INTEL   Unclaimed Idle      0.000   1485 17+11:44:48
slot1@C185025227.i LINUX        INTEL   Unclaimed Idle      0.000    620  0+00:55:05
slot2@C185025227.i LINUX        INTEL   Unclaimed Idle      0.000    620  1+08:41:03
C186098058.ipop    LINUX        INTEL   Unclaimed Idle      0.000   2971  0+00:55:04
slot1@C186185058.i LINUX        INTEL   Unclaimed Idle      0.000   1485  0+00:50:04
slot2@C186185058.i LINUX        INTEL   Unclaimed Idle      0.000   1485  6+08:40:33
slot1@C191158136.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+02:15:05
slot2@C191158136.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+18:07:09
slot1@C193080211.i LINUX        INTEL   Unclaimed Idle      0.000    620 17+11:10:58
slot2@C193080211.i LINUX        INTEL   Unclaimed Idle      0.000    620  0+00:10:05
C194165156.ipop    LINUX        INTEL   Unclaimed Idle      0.000   3225  0+03:25:05
slot1@C195096083.i LINUX        INTEL   Unclaimed Idle      0.000   1485  0+03:40:04
slot2@C195096083.i LINUX        INTEL   Unclaimed Idle      0.000   1485  3+19:27:53
slot1@C206214018.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+14:35:52
slot2@C206214018.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+02:45:06
slot1@C211145230.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+17:56:31
slot2@C211145230.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+02:25:06
slot1@C223073125.i LINUX        INTEL   Unclaimed Idle      0.000    620  0+00:50:04
slot2@C223073125.i LINUX        INTEL   Unclaimed Idle      0.000    620  3+08:37:05
slot1@C224255233.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+02:20:04
slot2@C224255233.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+18:08:52
slot1@C228022109.i LINUX        INTEL   Unclaimed Idle      0.030   1485  0+00:25:04
slot2@C228022109.i LINUX        INTEL   Unclaimed Idle      0.000   1485  9+00:37:36
slot1@C232105165.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+02:30:04
slot2@C232105165.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+14:35:25
slot1@C235052143.i LINUX        INTEL   Unclaimed Idle      0.000   1485  0+00:45:04
slot2@C235052143.i LINUX        INTEL   Unclaimed Idle      0.000   1485  7+16:51:09
slot1@C235252250.i LINUX        INTEL   Unclaimed Idle      0.000   1485  4+08:51:42
slot2@C235252250.i LINUX        INTEL   Unclaimed Idle      0.000   1485  0+00:45:05
C240195038.ipop    LINUX        INTEL   Owner     Idle      0.320    249  0+00:05:12
slot1@C245091047.i LINUX        INTEL   Unclaimed Idle      0.000   1485  0+02:00:04
slot2@C245091047.i LINUX        INTEL   Unclaimed Idle      0.000   1485  9+09:59:04
slot1@C254063065.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+00:40:07
slot2@C254063065.i LINUX        INTEL   Unclaimed Idle      0.000   1775  0+16:07:10

            Total Owner Claimed Unclaimed Matched Preempting Backfill

  INTEL/LINUX   136     2       1       133       0        0         0

        Total   136     2       1       133       0        0         0
griduser@C240195038:~$
```
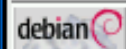
debian | 0 | X xmessage | griduser@C240... | griduser@C240195... | 10:24:03 AM

# Where We Are Now

- CondorView
- Tutorials
- Hadoop Appliance (Uses Condor to Organize)
- Archer
  - Computer Architecture research/education
  - Hundreds of cores distributed over 6 universities in the US over 3 years
- Over 100,000 CPU Hours used since October



UF UNIVERSITY of FLORIDA ACIS

# Projects Using IPOP

- Purdue – BoilerGrid

- Clemson – Campus Grids

- FCCN in Portugal – Hadoop-based Web Indexing - GaPPa

# Questions

- Our Projects
  - Grid Appliance, Archer – grid-appliance.org
  - IPOP – ipop-project.org
- Acknowledgements
  - ACIS P2P Group
  - Condor Group – Ben Burnett and Alain Roy
  - Effort sponsored by the NSF under grants OCI-0721867 and CNS-0751112.
  - Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF.

UF UNIVERSITY of FLORIDA ACIS

# Simics w/ NFS

- Simics requires license server – runs on condor master – Time sync issue :(

- Can use LARGE (many GB) input files

- Expensive to transfer over the wide area!

- NFS allows block level access, only transfer what you need, and data cache, so if a job comes again, we reduce our transfer overhead

- From http://www.grid-appliance.org/wiki/index.php/Archer:Simics

# User Feedback

The best thing we like about Archer is its computing power and easy accessibility and usability. We can easily access 50 or more Intel Xeon cores remotely for our computing needs to accomplish something we cannot do before. For example, in our recent data prefetching and cache management studies, we used more than 10000 simulation hours on Archer.

Besides the huge computational resources, Archer is unbelievably easy to use. All the needed software is packaged in a virtual machine, including the IPOP, which provides communication among the grid. It is all transparent to users. Users just need to download the Grid Appliance package, and is ready to go. Archer even provides a couple of popular CPU simulators by default, like Simics, SimpleScalar and PTLsim, etc. Our group uses Simics frequently and glad to see it is available on Archer. Besides, Archer employs Condor to manage all the tasks and resources, which makes it easy to deploy/monitor the tasks and need not worry about the resources. With the newly added feature of NFS interface, we can do more in a customized way. It allows mounting a local virtual disk to the grid, and sharing user-specific files, i.e., large Simics checkpoints.  All the grid nodes can access the files shared in the NFS file system. This feature helps build our own simulation environments efficiently.

- Jih-Kwon Peir