# Improving Scalability and Fault Tolerance in an Application Management Infrastructure

Nikolay Topilski, Jeannie Albrecht, and Amin Vahdat

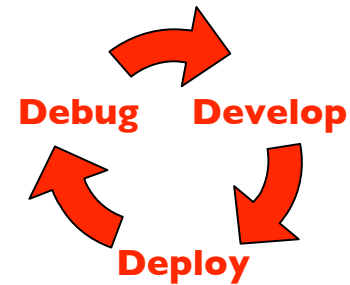Williams College & UC San Diego

# Large-Scale Computing

- Large-scale computing has many advantages

  - Increased computing power leads to improved performance, scalability, and fault tolerance

- Also introduces many new challenges

  - Building and managing distributed applications to leverage full potential of large-scale environments is difficult
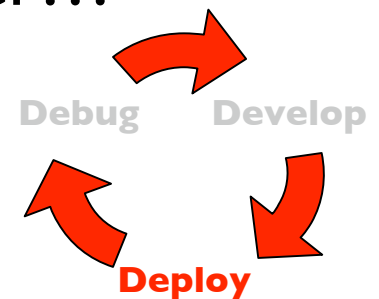
# Distributed Application Management

- Develop-Deploy-Debug cycle
  - Develop software
  - Deploy on distributed machines
  - Debug code when problems arise
- Management challenges in large-scale environments
  - Configuring resources
  - Detecting and recovering from failures
  - Achieving scalability and fault tolerance

- Research goal: Build an application management infrastructure that addresses these challenges

# Deploying an Application

- Steps required to deploy an application
  1. Connect to each resource
  2. Download software
  3. Install software
  4. Run application
  5. Check for errors on each machine
  6. When we find an error, we start all over…
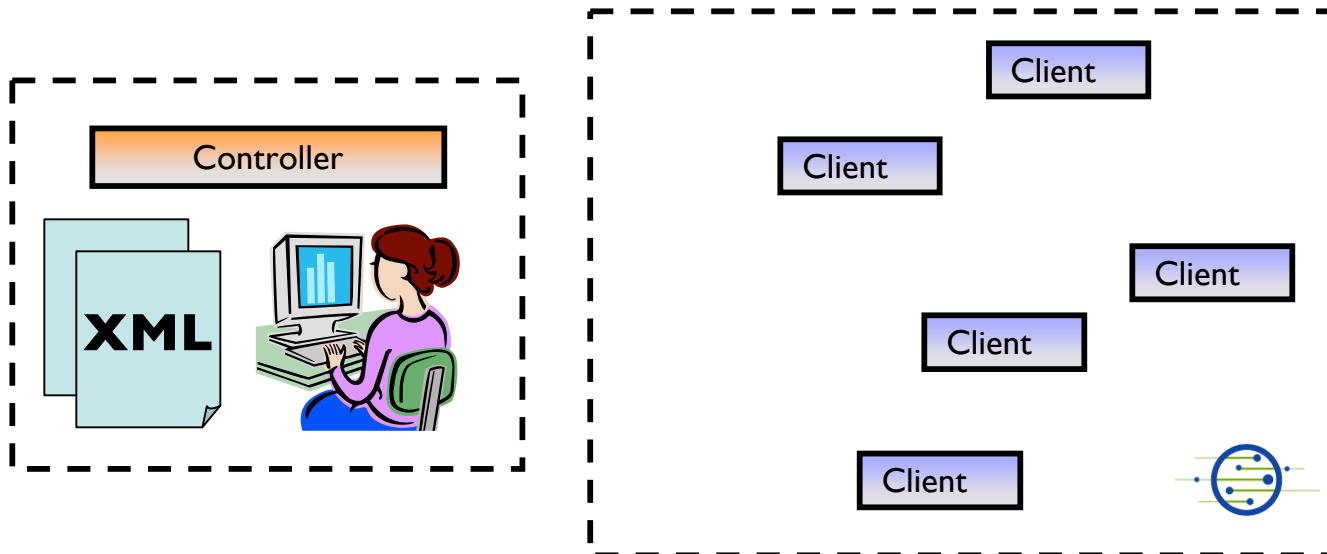- A better alternative: Plush

Debug    Develop

Deploy

# Plush

- Distributed application management infrastructure
  - Designed to simplify management of distributed applications
  - Help software developers cope with the challenges of large-scale computing
  - Support most applications in most environments

- Talk overview
  - Give brief overview of Plush architecture
  - Discuss scalability and fault tolerance limitations in original design
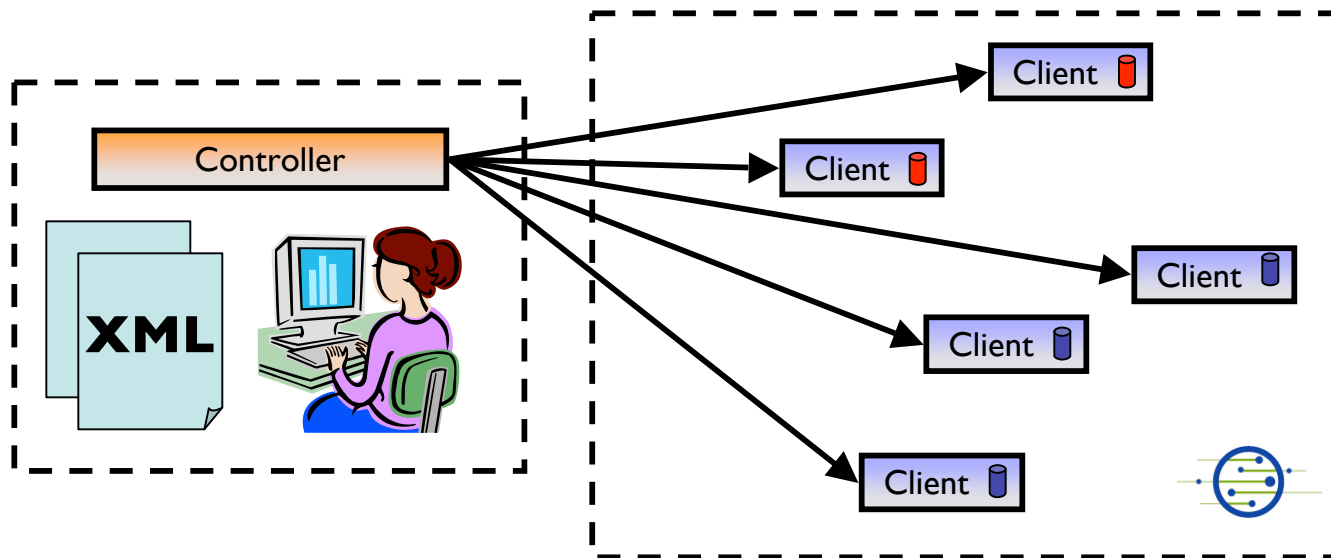  - Investigate ways to improve limitations

# Plush Overview

- Plush consists of two main components:
  - Controller - runs on user's Desktop
  - Client - runs on distributed resources
- To start application, user provides controller with application specification and resource directory (XML)

# Plush Overview

- Controller makes direct TCP connection to each client process running remotely
  - Communication mesh forms star topology
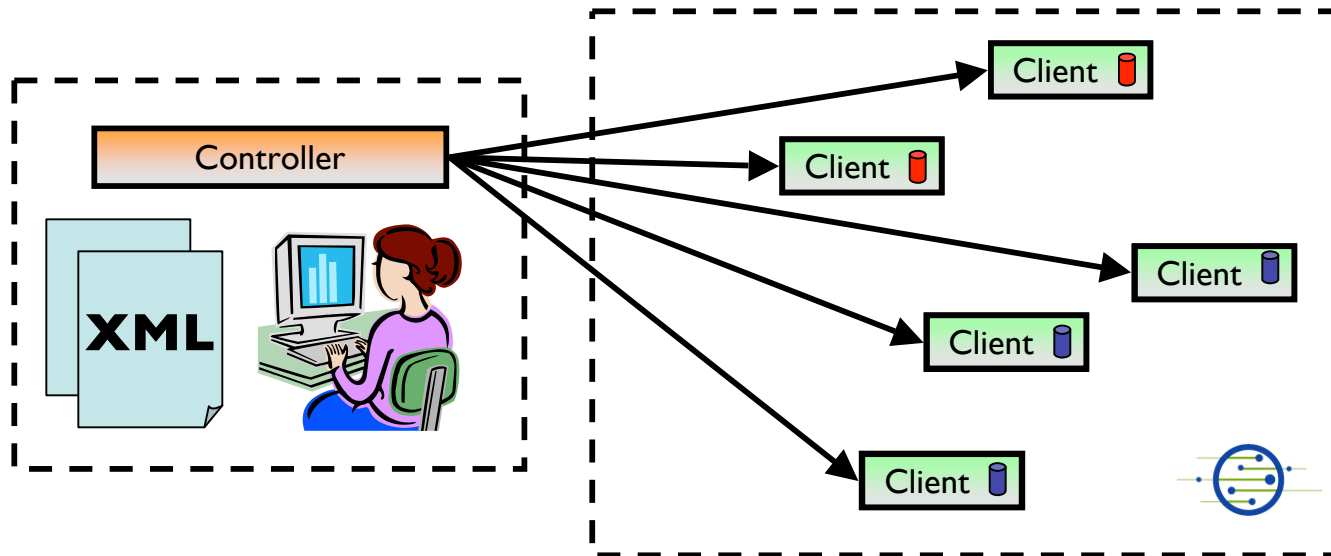- Controller instructs clients to download and install software (described in app spec)

# Plush Overview

- When all resources have been configured, controller instructs clients to begin execution

- Clients monitor processes for errors
  - Notify controller if failure occurs

# Plush Overview

- Once execution completes, controller instructs clients to "clean up"
  - Stop any remaining processes
  - Remove log files
  - Disconnect TCP connections

# Plush User Interfaces

- Command-line interface used to interact with applications
- Nebula (GUI) allows users to describe, run, & visualize applications
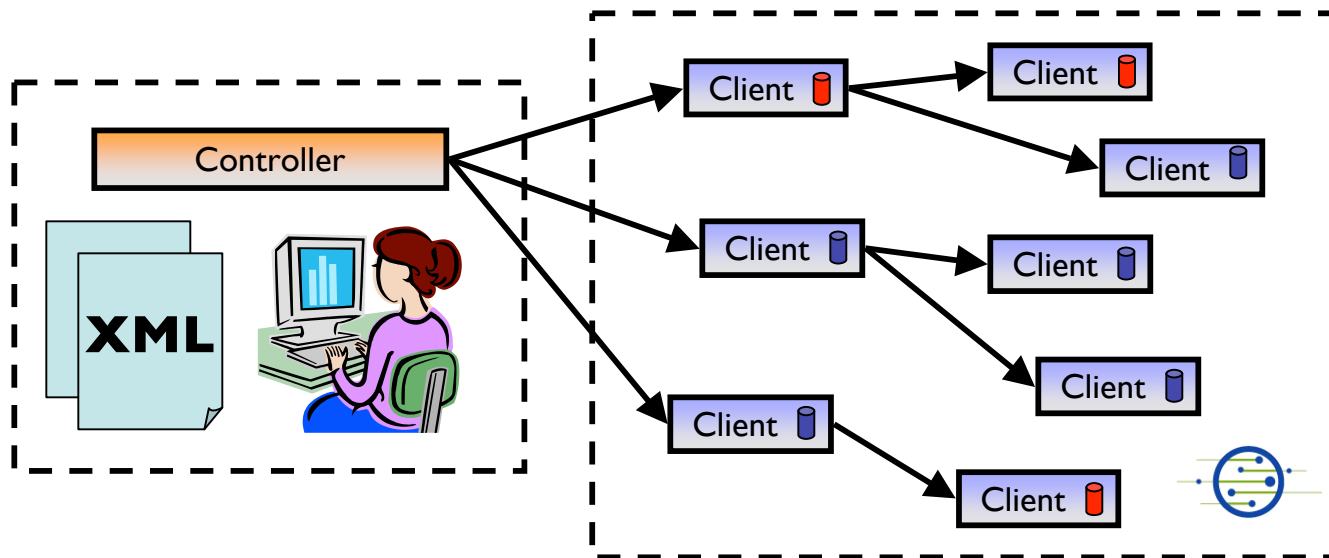- XML-RPC interface for managing applications programmatically

# Limitations

- Plush was designed with PlanetLab in mind…
  - … in 2004!
  - PlanetLab grew from 300 machines to 800+
- Plush now supports execution in a variety of environments in addition to PlanetLab
  - Some have 1000+ resources

- Problems
  - Star topology does not scale beyond ~300 resources
  - Tree topology scales but is not resilient to failure

# Insights

- We need a *resilient* overlay tree in place of the star
- Lots of people have already studied overlay tree building algorithms
- Mace is a framework for building overlays
  - Developed at UCSD
  - Simplifies development through code reuse
- Solution: Combine Plush with overlay tree provided by Mace!
  - Allow us to explore different tree building protocols
  - Leverage existing research in overlay networks without "reinventing the wheel"
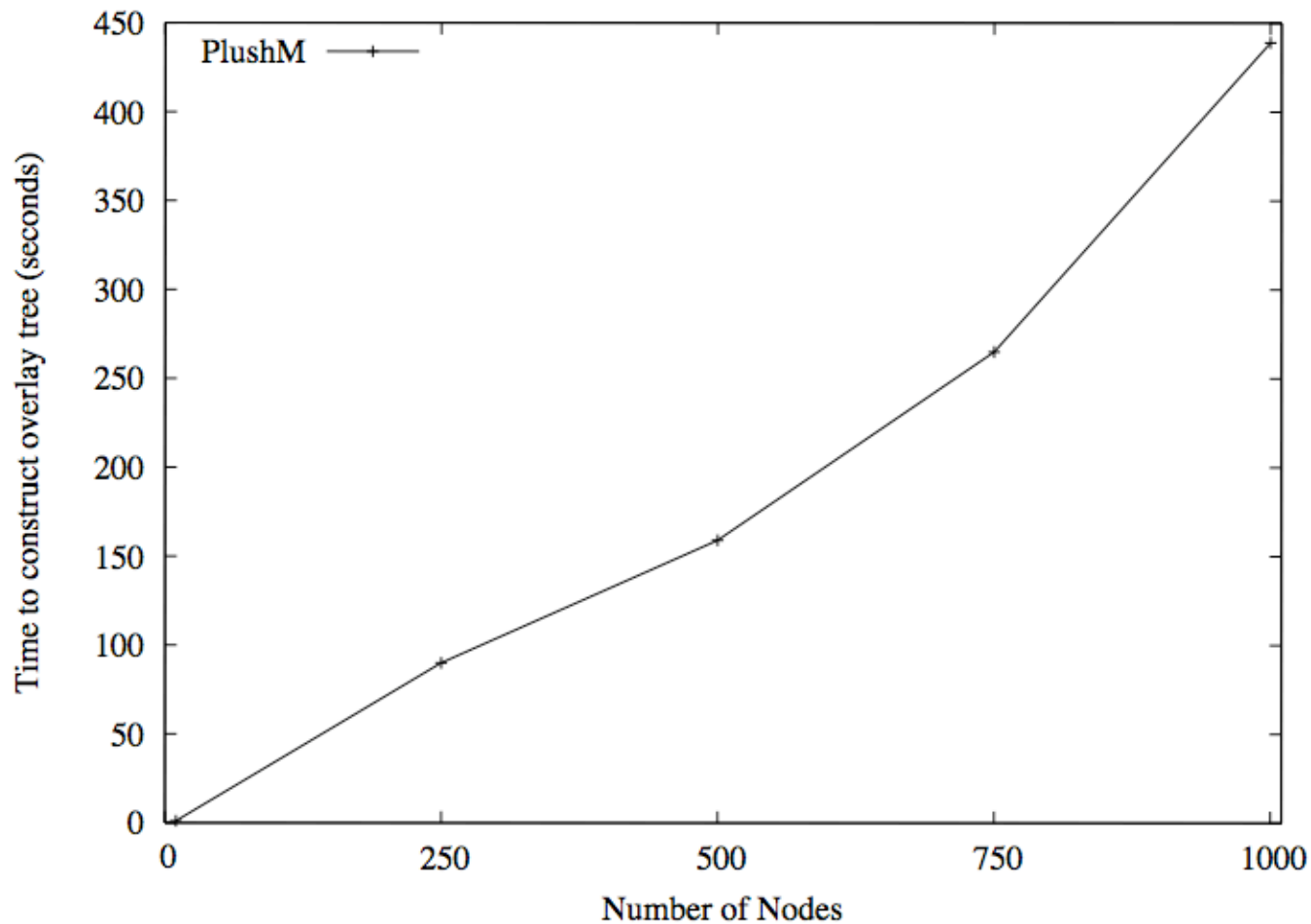  - Improve scalability and fault tolerance of Plush

# Introducing PlushM

- We extended the existing communication fabric in Plush to allow interaction with Mace ($\Rightarrow$ PlushM)

- PlushM still uses same abstractions for application management as Plush

- We chose RandTree as our initial overlay topology
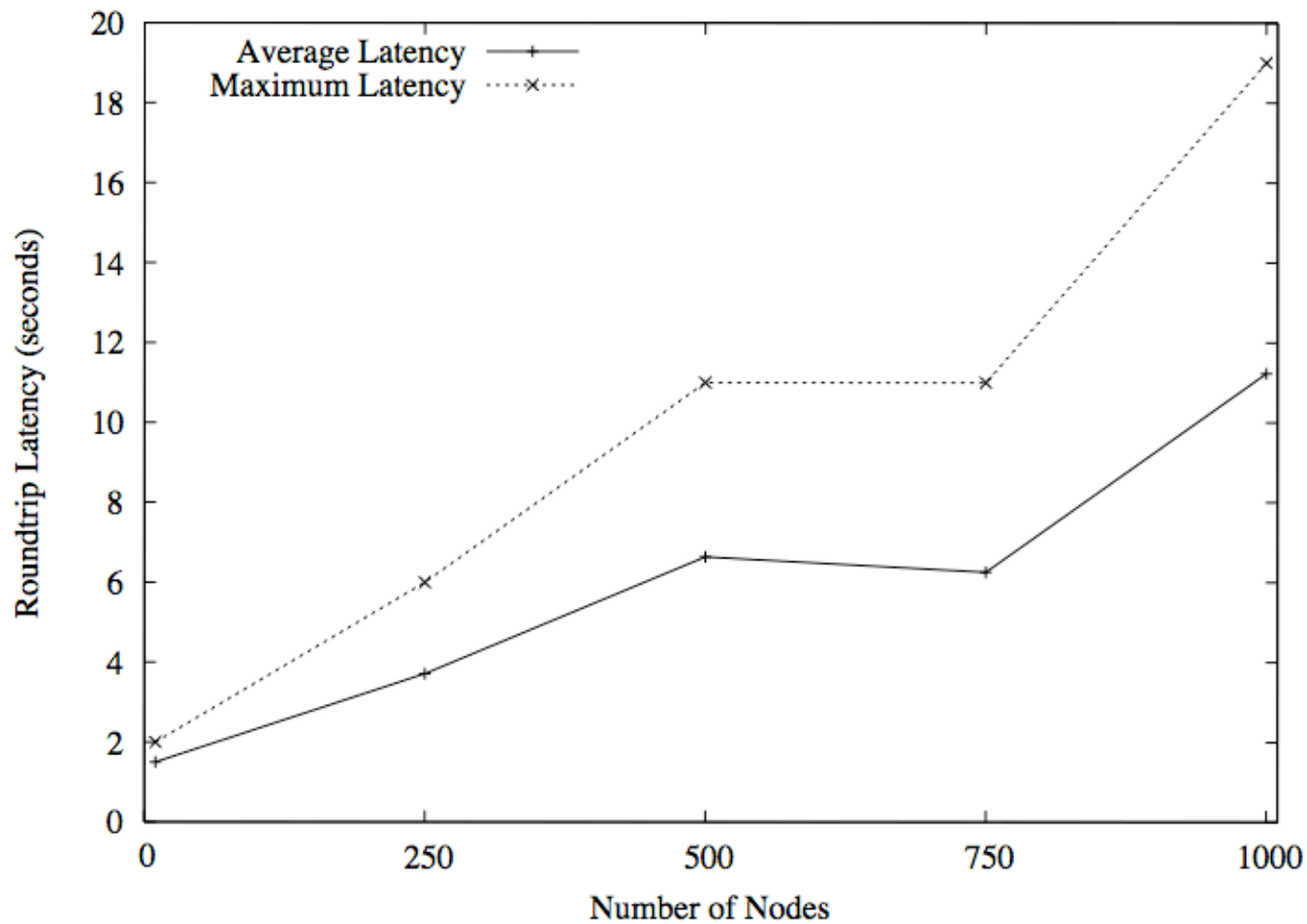  - Random overlay tree that reconfigures when failure occurs

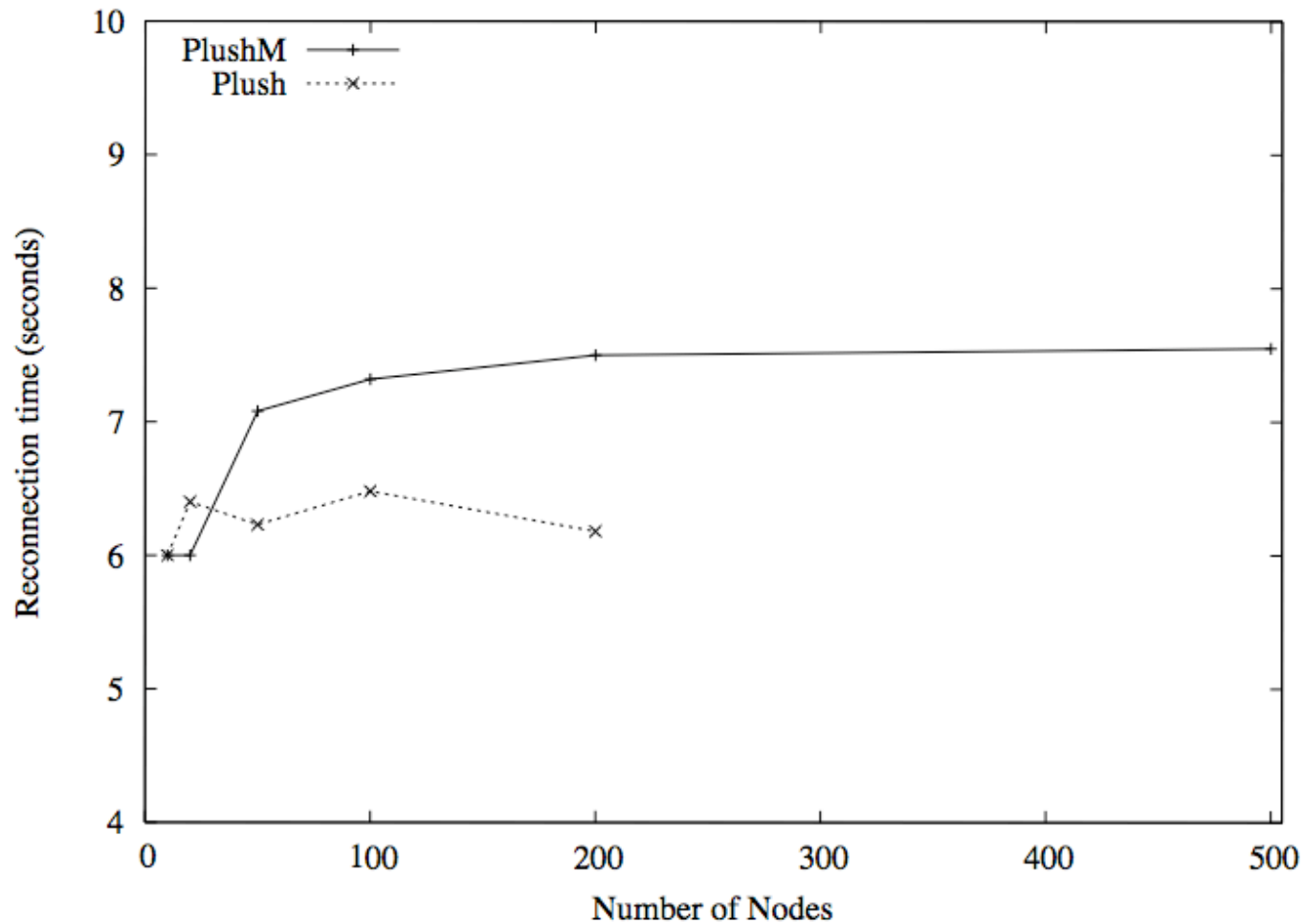# Evaluating Scalability

- Overlay tree construction time

# Evaluating Scalability

- Message propagation time

# Evaluating Fault Tolerance

- Reconfiguration time after disconnect (ModelNet)

# Conclusions and Future Work

- Plush provides distributed application management in a variety of environments
  - Original design has scalability/fault tolerance limitations in large-scale clusters
- PlushM replaces Plush's communication infrastructure with Mace overlay to provide better scalability (1000 resources) and fault tolerance
- Future work
  - Evaluate PlushM on larger topologies
  - Investigate the user of other Mace overlays in addition to RandTree
  - Explore ways to improve PlushM performance

# Thank you!

Plush http://plush.cs.williams.edu

Mace http://mace.ucsd.edu

## Email

ntopilsk@cs.ucsd.edu

jeannie@cs.williams.edu

vahdat@cs.ucsd.edu