# Software Upgrades for Distributed Systems

Sameer Ajmani
Google

Barbara Liskov          Liuba Shrira
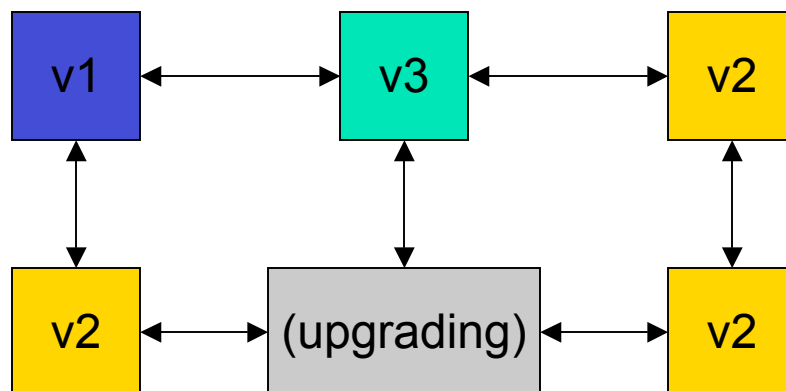MIT                     Brandeis

# Internet Services

- ν Are long-lived, robust
- ν Run on many machines
- ν Must be continuously available
- ν Have persistent state
- ν Face ever-changing requirements

- ν Require software upgrades to
  - ν Fix bugs
  - ν Improve performance
  - ν Add/change/remove features

# Upgrade Requirements

ν Automatic, Controlled Deployment

    ν Ensure continuous availability

    ν Test new software on a few nodes

    ν Upgrade servers before clients

ν **Mixed mode operation**

# Outline

- ν System & Upgrade Model
- ν Specifying Upgrades
- ν Implementation Models

# System Model

- A node is an object of class C
- Different nodes may run different classes
- Nodes communicate via RPCs

# Upgrade Model

- A class upgrade replaces an old class, $C_{old}$, with a new one, $C_{new}$
  - Implements types $T_{old}$, $T_{new}$
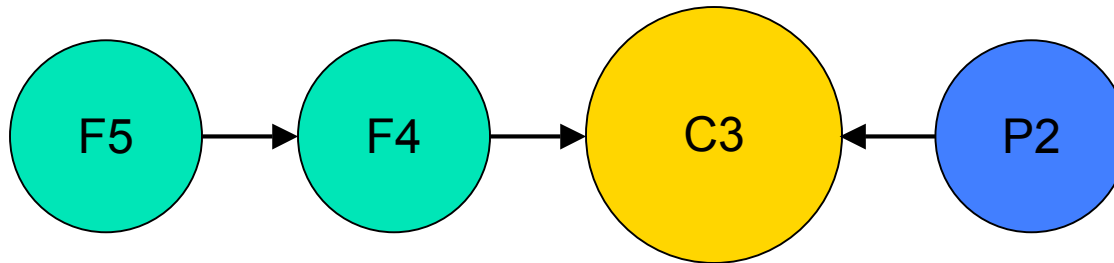  - May be compatible or incompatible
- An upgrade is a set of class upgrades

# Supporting Mixed Mode

ν Each node handles calls to past and future versions of itself

ν Adding support for new versions must be fast

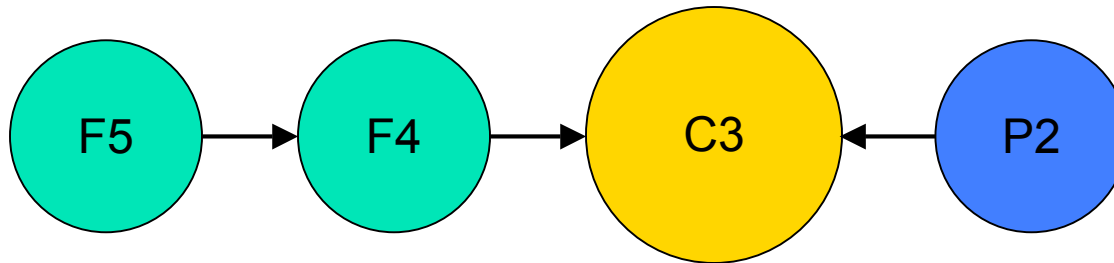ν Removing support for old versions should be easy

# Simulation Objects

ν Each node handles calls to past and future versions of itself

```
( F5 ) → ( F4 ) → ( C3 ) ← ( P2 )
```

ν Future SOs simulate future behavior

ν Past SOs simulate past behavior

ν Adding/removing an SO does not require a restart

# Simulation Objects

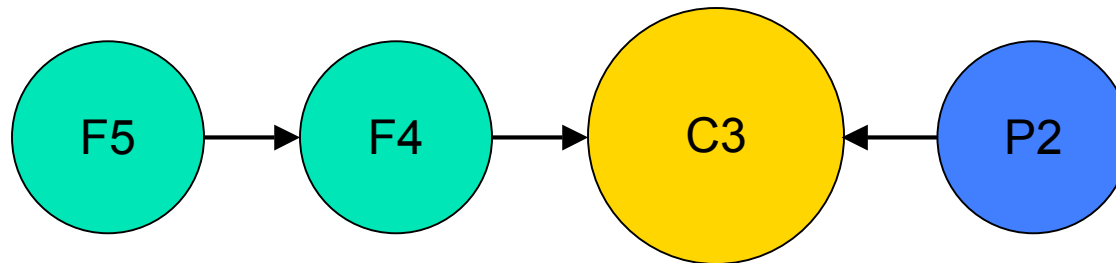- Each node handles calls to past and future versions of itself



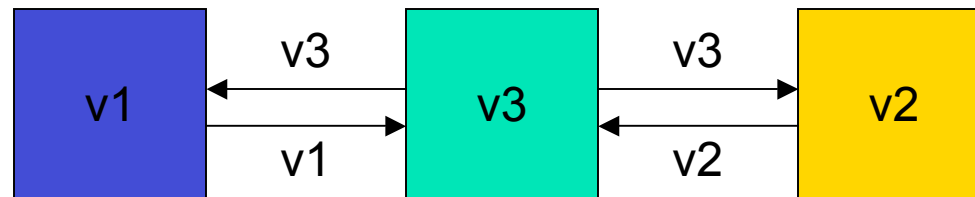- SOs are only required for certain upgrades

# Specifying Upgrades

# Specifying Upgrades

F5 → F4 → C3 ← P2

- ν **Must behave like a single object**
  - ν Even when upgrades are incompatible
- ν **Upgrade specification must define this**
  - ν Goal: no surprises for clients!
  - ν E.g., changing permissions to ACLs

# Constraints on Specifications

- Type requirement
  - A call of version V behaves according to the specification of type $T_V$

# Constraints on Specifications

- ν **Sequence requirement**
  - ν Each event must reflect all earlier ones despite:
    - ν Client upgrades
    - ν Server upgrades
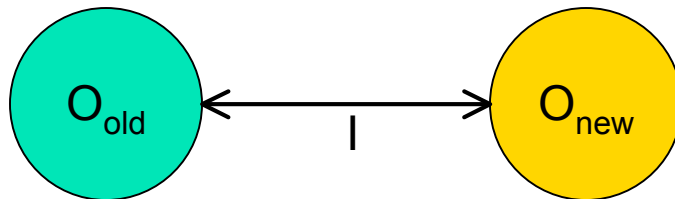    - ν Version introduced
    - ν Version retired

# Example

- ColorSet ◊ FlavorSet
  - Incompatible upgrade
- ColorSet methods: insertColor(x, c), getColor(x), …
  - E.g., { (1, red), (2, blue), (3, red) }
- FlavorSet methods: insertFlavor(x,f), getFlavor(x), …

# Specifications: Invariant

Invariant I relates the object states

$$I(O_{old}, O_{new})$$



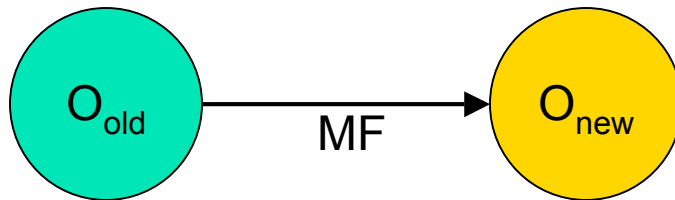$$I: \{x | <x,c> \text{ in } O_{CS}\} = \{x | <x,f> \text{ in } O_{FS}\}$$

# Specifications: Mapping Function

Mapping function MF defines initial state

$$O_{new} = MF(O_{old}) \quad s.t. \ I(O_{old}, O_{new})$$



$$O_{FS} = MF(O_{CS}) = \{<x, grape> | <x, c> \text{ in } O_{CS}\}$$

# Relating Behavior



Only for mutators

# Relating Behavior



**Shadow methods** relate behavior
  Told.m $\lozenge$ Tnew.\$m
  Tnew.p $\lozenge$ Told.\$p

# Shadow Method Specification

void ColorSet.$insertFlavor(x, f)

    Effects: no <x,c> in $\text{this}_{pre} \Rightarrow$

        $\text{this}_{post} = \text{this}_{pre} \cup \{<x,\text{blue}>\}$

- Also ColorSet.$delete, FlavorSet.$insertColor, FlavorSet.$delete

# The Compound Type

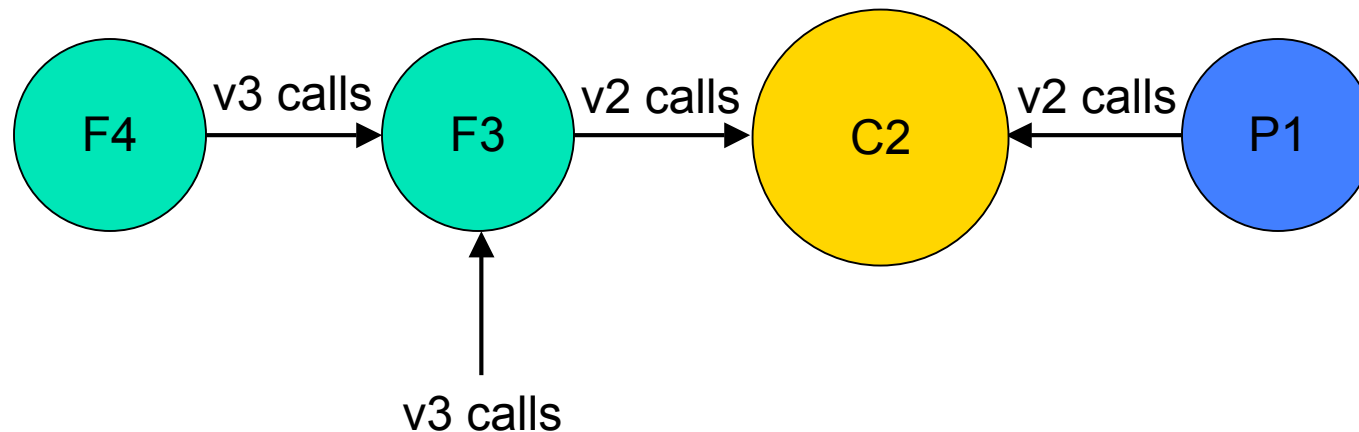- ν I, MF, and the shadow methods define a compound type $T_{old\&new}$
  - ν All the methods, with extended specs for mutators
- ν We would like:
  - ν $T_{old\&new}$ is a subtype of $T_{old}$, $T_{new}$
- ν When this doesn't work:
  - ν Weaken invariant I
  - ν Upgrade scheduling
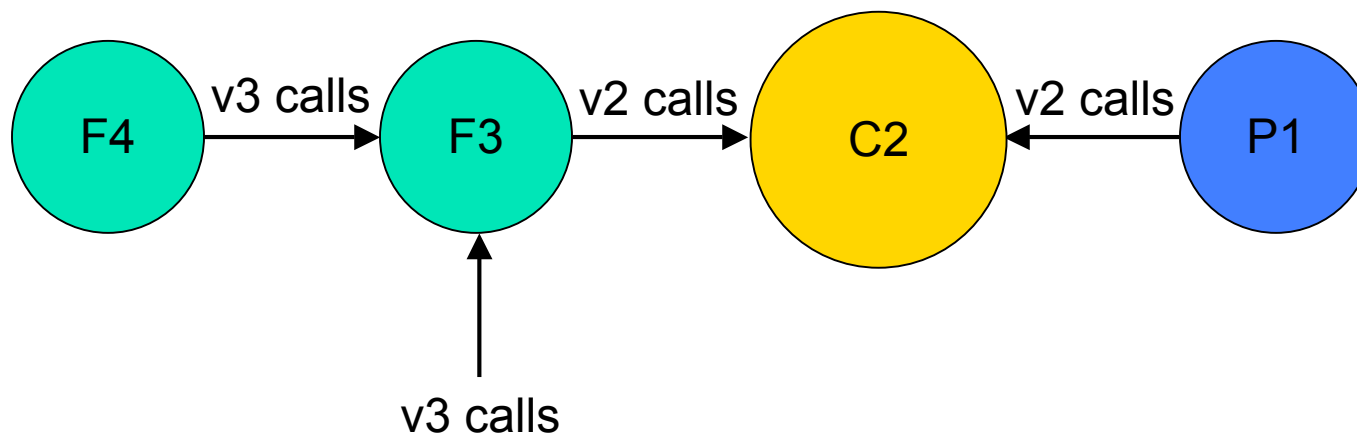  - ν Disallow methods

# Implementation Models

# Direct Model

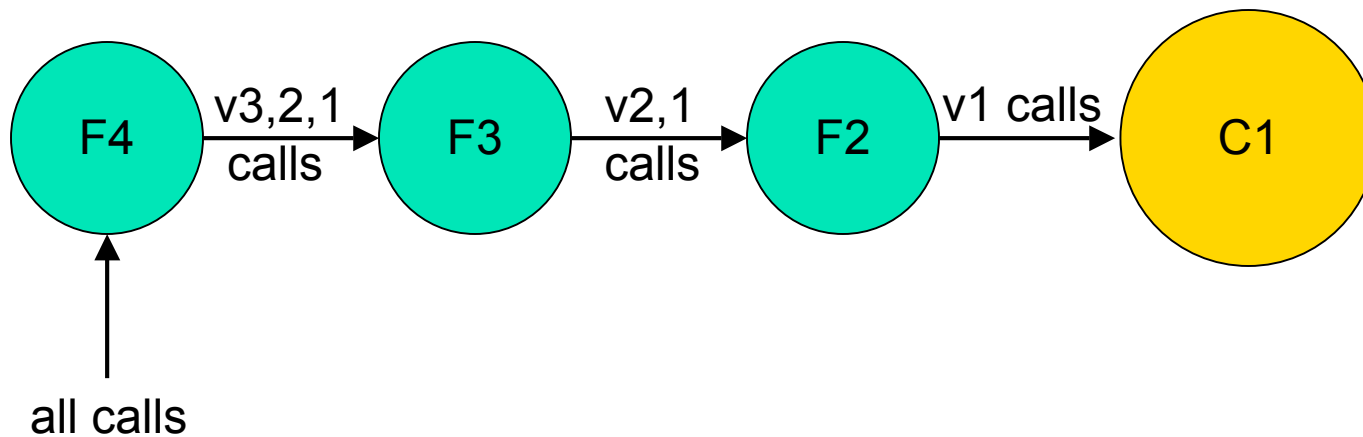- SO handles calls only to its own version
  - And delegates down the chain

# Problems with Direct Model

- Poor expressive power
  - E.g., FlavorSet SO doesn't know about C.insertColor call
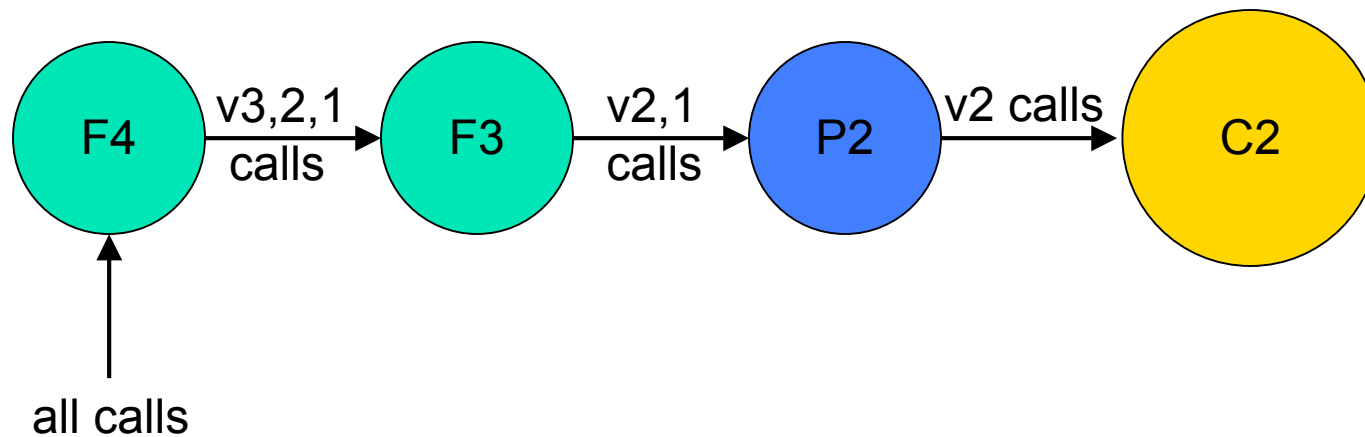- Synchronization

# Interceptor Model

- Newest SO gets all calls
  - And delegates down the chain

# Interceptor Model

ν   After first (incompatible) upgrade is installed

# Interceptor Model

ν After second (possibly compatible)
upgrade is installed

F4 —v3,2,1 calls→ P3 —v3 calls→ C3

all calls → F4

# Interceptor Model Evaluation

ν Excellent expressive power

ν Future and past SO must do more

ν Can reuse code and delegate

F4 —v3,2,1 calls→ P3 —v3 calls→ C3

all calls → F4

# Prototype Implementation: Upstart

- C++ and Sun RPC

- Intercepts socket(), read(), write()

- Imposes minimal overhead

# Summary

- Upstart is the first complete approach
  - Allows mixed mode operation
- The first definition of what must be specified for incompatible upgrades
- A powerful and useful implementation model
- A prototype implementation

# Software Upgrades for Distributed Systems

Sameer Ajmani

Google

Barbara Liskov

Liuba Shrira

MIT

Brandeis

# Disallowing

- Constraint: never disallow methods of the current object

- Future SO may disallow $T_{new}$ methods

- Past SO may disallow $T_{old}$ methods

- In either case, disallow
    - Mutators whose shadows are problem
    - Observers that expose problems

# Some shadows cannot be implemented via delegation

- ν Disallow methods that have unimplementable shadows
- ν Add shadow method to delegate via dynamic updating
  - ν Allowed iff T + shadow is a subtype of T
  - ν E.g., can't add delete() to GrowSet
- ν Implement shadow method in interceptor
  - ν Impacts transform function
  - ν Won't work for past SOs because of retirement

# What does Google do?

- ν Extensible protocols
  - ν Assume defaults for missing fields
  - ν Ignore unexpected fields
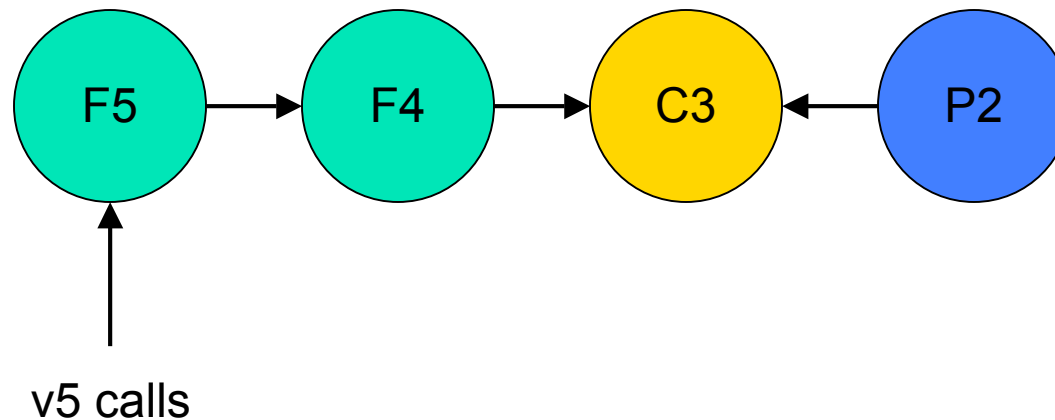- ν Round-robin upgrades among replicas
- ν Datacenter-by-datacenter

# END OF SLIDES

- The remaining slides are leftover from previous talks and may contain stale information

# Code Execution

- Call contains version number
- Called node dispatches



F5 → F4 → C3 ← P2

v5 calls

# Upstart

- A system that supports upgrades
- And a methodology

- Joint work with
  - Barbara Liskov
  - Liuba Shrira

# Class Upgrade

- New and old classes $C_{new}$, $C_{old}$
  - Implement $T_{new}$, $T_{old}$
- Scheduling function SF
- Transform function TF
- Simulation classes $S_{new}$, $S_{old}$

- Might be <span style="color:red">incompatible</span>
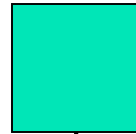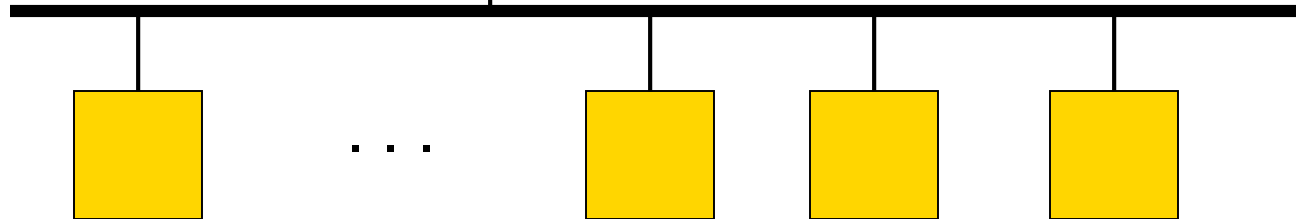  - $T_{new}$ is not a subtype of $T_{old}$

# Class Upgrade

- Replaces an old class, $C_{old}$ with a new one, $C_{new}$
- Every node running old class will switch to new class eventually

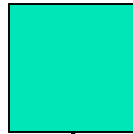- Upgrade is a set of class upgrades

# Defining an Upgrade

UDB

nodes · · ·

- Upgrader enters new upgrade at UDB
- Defines a new version

# Propagating an Upgrade

UDB

nodes . . .

- Nodes query the UDB periodically
- Version numbers flow on all messages

# Executing an Upgrade

ν **If upgrade affects the node**

   ν Runs the SF

      ν And <span style="color:red">simulates the future</span>

   ν Shuts down, restarts, runs TF

   ν Starts up "normally"

      ν And <span style="color:red">simulates the past</span>

# Disallowing Example

- GrowSet ◊ IntSet

- For the future SO:
    - Disallow IntSet.delete

- For the past SO:
    - Disallow GrowSet.isIn


- $T_{old\&new}$ becomes $<T_{future}, T_{past}>$

# What Disallowing Provides

- $T_{future}$ is a subtype of $T_{old}$
- And it implements $T_{new}$

- $T_{past}$ is a subtype of $T_{new}$
- And it implements $T_{old}$

# Transform Functions

- Implement the identity map
- May need to use future SO, create past SO
- Must be restartable
- Cannot make remote calls

# Scheduling Functions

ν Can consult the UDB

ν Examples:

    ν Rolling upgrade

    ν Big flip

    ν Fast reboot

# Implementing Upgrades

- ν Need to provide SOs, TF, SF
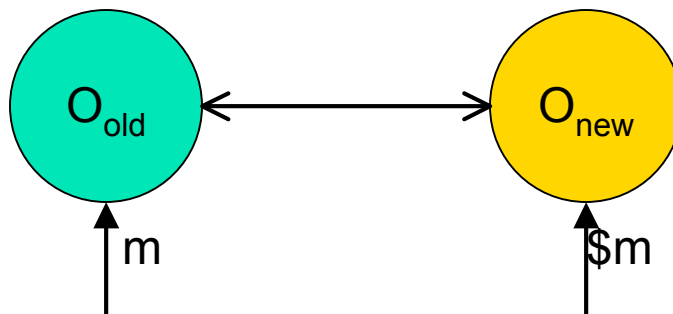- ν For the SOs, need an implementation model

# Summary of Specifications

- Specification defines the compound type $T_{old\&new}$
  - I, MF, and the shadows
- If the compound type isn't a subtype, disallow

# Specifications: Shadow Methods

ν  Shadow methods relate behavior

$$T_{old}.m \lozenge T_{new}.\$m$$
$$T_{new}.p \lozenge T_{old}.\$p$$



e.g., FlavorSet.$insertColor

# Talk Outline

- Upgrade requirements
- Upstart overview
- Specifying upgrades
- Implementing upgrades

# Requirement: Generality

- Support for arbitrary changes
- <span style="color:red">Incompatible upgrades</span>
  - Old features are no longer supported

# Requirement: Continuous Availability

- ν Service is required 24/7
- ν Even when upgrading
- ν Therefore systems upgrade gradually

- ν Implies mixed-mode operation

# Requirement: Controlled Deployment

- ν Systems upgrade gradually
  - ν But with control

- ν Manual control is impractical
  - ν An automatic system
  - ν But upgrader needs control

# Requirement: Persistence

- Systems store important state for users
  - It cannot be lost
  - But may need to be transformed

# Requirement: Ease of Use

- Avoiding feature creep helps
- Upgrader needs to understand only a few recent versions