

Stork: Package Management for Distributed VM Environments

Justin Cappos, Scott Baker, Jeremy Plichta, Duy Nyugen, Jason Hardies, Matt Borgard, Jeffry Johnston, and John H. Hartman – University of Arizona

ABSTRACT

In virtual machine environments each application is often run in its own virtual machine (VM), isolating it from other applications running on the same physical machine. Contention for memory, disk space, and network bandwidth among virtual machines, coupled with an inability to share due to the isolation virtual machines provide, leads to heavy resource utilization. Additionally, VMs increase management overhead as each is essentially a separate system.

Stork is a package management tool for virtual machine environments that is designed to alleviate these problems. Stork securely and efficiently downloads packages to physical machines and shares packages between VMs. Disk space and memory requirements are reduced because shared files, such as libraries and binaries, require only one persistent copy per physical machine. Experiments show that Stork reduces the disk space required to install additional copies of a package by over an order of magnitude, and memory by about 50%. Stork downloads each package once per physical machine no matter how many VMs install it. The transfer protocols used during download improve elapsed time by 7X and reduce repository traffic by an order of magnitude. Stork users can manage groups of VMs with the ease of managing a single machine – even groups that consist of machines distributed around the world. Stork is a real service that has run on PlanetLab for over four years and has managed thousands of VMs.

Introduction

The growing popularity of virtual machine (VM) environments such as Xen [3], VMWare [31], and Vservers [17, 18], has placed new demands on package management systems (e.g., apt [2], yum [36], RPM [27]). Traditionally, package management systems deal with installing and maintaining software on a single machine whether virtual or physical. There are no provisions for inter-VM sharing, so that multiple VMs on the same physical machine individually download and maintain separate copies of the same package. There are also no provisions for inter-machine package management, centralized administration of which packages should be installed on which machines, or allowing multiple machines to download the same package efficiently. Finally, current package management systems have relatively inflexible security mechanisms that are either based on implicit trust of the repository, or public/private key signatures on individual packages.

Stork is a package management system designed for distributed VM environments. Stork has several advantages over existing package management systems: it provides *secure* and *efficient* inter-VM package sharing on the same physical machine; it provides centralized package management that allows users to determine which packages should be installed on which VMs without configuring each VM individually; it allows multiple physical machines to download the same package efficiently; it ensures that package

updates are propagated to the VMs in a timely fashion; and it provides a flexible security mechanism that allows users to specify which packages they trust as well as delegate that decision on a per-package basis to other (trusted) users.

Stork's inter-VM sharing facility is important for reducing resource consumption caused by package management in VM environments. VMs are excellent for isolation, but this very isolation can increase the disk, memory, and network bandwidth requirements of package management. It is very inefficient to have each VM install its own copy of each package's files. The same is true of memory: if each VM has its own copy of a package's files then it will have its own copy of the executable files in memory. Memory is often more of a limiting factor than disk, so Stork's ability to share package files between VMs is particularly important for increasing the number of VMs a single physical machine can support. In addition, Stork reduces network traffic by only downloading a package to a physical machine once, even if multiple VMs on the physical machine install it.

Stork's inter-machine package management facility enables centralized package management and efficient, reliable, and timely package downloads. Stork provides package management utilities and configuration files that allow the user to specify which packages are to be installed on which VMs. Machines download packages using efficient transfer mechanisms such as BitTorrent [9] and CoBlitz [22], making downloads

efficient and reducing the load on the repository. Stork uses fail-over mechanisms to improve the reliability of downloads, even if the underlying content distribution systems fail. Stork also makes use of publish/subscribe technology to ensure that VMs are notified of package updates in a timely fashion.

Stork provides all of these performance benefits without compromising security; in fact, Stork has additional security benefits over existing package management systems. First, Stork shares files securely between VMs. Although a VM can delete its link to a file, it cannot modify the file itself. Second, a user can securely specify which packages he or she trusts and may delegate this decision for a subset of packages to another user. Users may also trust other users to know which packages *not* to install, such as those with security holes. Each VM makes package installation decisions based on a user's trust assumptions and will not install packages that are not trusted. While this paper touches on the security aspects of the system that are necessary to understand the design, a more rigorous and detailed analysis of security is available through documentation on our website [29].

In addition, Stork is flexible and modular, allowing the same Stork code base to run on a desktop PC, a Vserver-based virtual environment, and a PlanetLab node. This is achieved via pluggable modules that isolate the platform-specific functionality. Stork accesses these modules through a well-defined API. This approach makes it easy to port Stork to different environments and allows the flexibility of different implementations for common operations such as file retrieval.

Stork has managed many thousands of VMs and has been deployed on PlanetLab [23, 24] for over four years. Stork is currently running on hundreds of PlanetLab nodes and its package repository receives a request roughly every ten seconds. Packages installed in multiple

VMs by Stork typically use over an order of magnitude less space and 50% the memory of packages installed by other tools. Stork also reduces the repository load by over an order of magnitude compared to HTTP-based tools. Stork is also used in the Vserver [18] environment and can also be used in non-VM environments (such as on a home system) as an efficient and secure package installation system. The source code for Stork is available at <http://www.cs.arizona.edu/stork>.

Stork

Stork provides manual management of packages on individual VMs using command-line tools that have a syntax similar to `apt` [2] or `yum` [36]. Stork also provides centralized management of groups of VMs. This section describes an example involving package management, the configuration files needed to manage VMs with Stork, and the primary components of Stork.

An Example

Consider a system administrator that manages thousands of machines at several sites around the globe. The company's servers run VM software that allow different production groups more flexible use of the hardware resources. In addition, the company's employees have desktop machines that have different software installed depending on their use.

The system administrator has just finished testing a new security release for a fictional package `foobar` and she decides to have all of the desktop machines used for development update to the latest version along with any testing VMs that are used by the coding group. The administrator modifies a few files on her local machine, signs them using her private key, and uploads them to a repository. Within minutes all of the desired machines that are online have the updated `foobar` package installed. As offline machines come online or new VMs

File Type	Repository	Client	Central Mgmt	Signed and Embedded
User Private Key	No	No	Yes	No
User Public Key	No †	Yes	Yes	No
Master Configuration File	No †	Yes	Yes	No
Trusted Packages (TP)	Yes	Yes	Yes	Yes
Pacman Packages	Yes	No	Yes	Yes
Pacman Groups	Yes	No	Yes	Yes
Packages (RPM, tar.gz)	Yes	Yes	Yes	Secure Hash
Package Metadata	Yes	Yes	Yes	No
Repository Metahash	Yes	Yes	No	Signed Only

Table 1: Stork File Types: This table shows the different types of files used by Stork. The repository column indicates whether or not the file is obtained from the repository by the clients. The client column indicates whether or not the file is used for installing packages or determining which packages should be installed locally based upon the files provided by the centralized management system. The centralized management column indicates if the files are created by the management tools. The signed/embed column indicates which files are signed and have a public key embedded in their name.

† In order to automatically deploy Stork on PlanetLab, this restriction is relaxed. See the PlanetLab section for more details.

are created, they automatically update their copies of foobar as instructed.

The subsequent sections describe the mechanisms Stork uses to provide this functionality to its users. The walkthrough section revisits this example and explains in detail how Stork provides the functionality described in this scenario.

File Types

Stork uses several types of files that contain different information and are protected in different ways (Table 1). The user creates a public/private key pair that authenticates the user to the VMs he or she controls. The *public key* is distributed to all of the VMs and the *private key* is used to sign the configuration files. In our previous example, the administrator's public key is distributed to all of the VMs under her control. When files signed by her private key were added to the repository, the authenticity of these files was independently verified by each VM using the public key.

The *master configuration file* is similar to those found in other package management tools and indicates things such as the transfer method, repository name, user name, etc. It also indicates the location of the public key that should be used to verify signatures.

The user's trusted packages file (*TP file*) indicates which packages the user considers valid. The TP file does not cause those packages to be installed, but instead indicates trust that the packages have valid contents and are candidates for installation. For example, while the administrator was testing the latest release of foobar she could add it to her trusted packages file because she believes the file is valid.

There are two pacman files used for centralized management. The *groups.pacman* file allows VMs to be categorized into convenient groups. For example, the administrator could configure her pacman groups file to create separate groups for VMs that perform different tasks. VMs can belong in multiple groups such as ALPHA and ACCOUNTING for an alpha test version of accounting software. Any package management instructions for either the ALPHA group or the ACCOUNTING group would be followed by this VM.

The *packages.pacman* file specifies what actions should be done on a VM or a group of VMs. Packages can be installed, updated, or removed. Installation is different from updating in that installation will do nothing if there is a package that meets the criteria already installed while update ensures that the preferred version of the package is installed. For example, when asked to install foobar, if any version of the package is currently installed then no operation will occur. If asked to update foobar, Stork checks to see if the administrator's TP file specifies a different version of foobar and if so, replaces the current version with the new version.

The *packages* (for example, the foobar RPM itself) contain the software that is of interest to the user. The *package metadata* is extracted from packages and is published by the repository to describe the packages that are available. The *repository metahash* is a special file that is provided by the repository to indicate the current repository state.

Architecture

Stork consists of four main components:

- a *repository* that stores configuration files, packages, and associated metadata;
- a set of *client tools* that are used in each Stork client VM to manage its packages by interacting either directly with the repository or through the nest when it is available;
- a *nest* process that runs on physical machines and coordinates sharing between VMs as well as providing repository metadata updates to its client VMs and downloading packages;
- and *centralized management tools* that allows a user to control many VMs concurrently, create and sign packages, upload packages to the repository, etc.

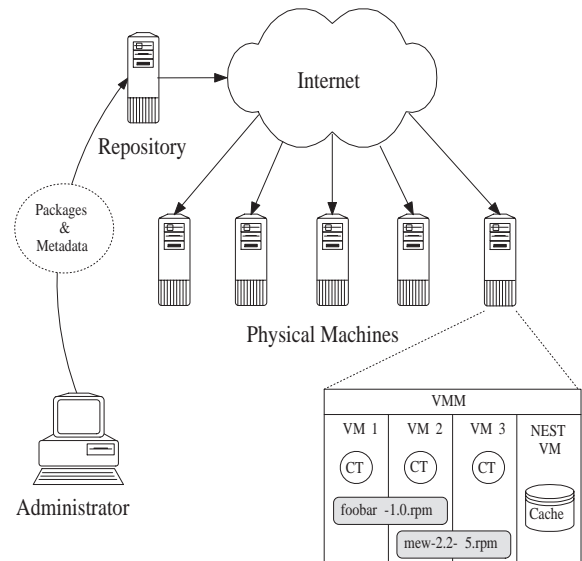


Figure 1: Stork Overview. Stork allows centralized administration and sharing of packages. The administrator publishes packages and metadata on the repository. Updates are propagated to VMs running on distributed physical machines. Each physical machine contains a single nest VM, and one or more client VMs that run the Stork client tools.

The client tools consist of the *stork command-line tool* (referred to simply as *stork*), which allows users to install packages manually, and *pacman*, which supports centralized administration and automated package installation and upgrade. While a client VM may communicate with the repository directly, it is far more

efficient for client VMs to interact with their local nest process, who interacts with the repository on their behalf.

Repository

The Stork repository's main task is to serve files much like a normal web server. However, the repository is optimized to efficiently provide packages to Stork client VMs. First, the repository provides secure user upload of packages, trusted packages files, and pacman packages and groups files. Second, the repository pushes notifications of new content to interested VMs. Third, the repository makes packages available via different efficient transfer mechanisms such as BitTorrent.

Handling Uploaded Data The Stork repository allows multiple users to upload files while retaining security. TP, groups.pacman, and packages.pacman files must be signed by the user that uploads them. Every signed file has a timestamp for the signature embedded in the portion of the file protected by the signature. The public key of the user is embedded in the file name of the signed file (similar to self-certifying path names [19]). This avoids naming conflicts and allows the repository to verify the signature of an uploaded file. The repository will only store a signed file with a valid signature that is newer than any existing signed file of the same name. This prevents replay attacks and allows clients to request files that match a public key directly.

Packages and package metadata are treated differently than configuration files. These files are not signed, but instead incorporate a secure hash of their contents in their names. This prevents name collisions and allows clients to request packages directly by secure hash. In all cases, the integrity of a file is verified by the recipient before it is used (either by checking the signature or the secure hash, as appropriate). The repository only performs these checks itself to prevent pollution of the repository and unnecessary downloads, rather than to ensure security on the clients.

Pushing Notifications The repository notifies interested clients when the repository contents have changed. The repository provides this functionality by pushing an updated repository metahash whenever data has been added to the repository. However, this does not address the important question of *what* data has been updated. This is especially difficult to address when VMs may miss messages or suffer other failures.

One solution is for the repository to push out hashes of all files on the repository. As there are many thousands of metadata files on the repository, it is too costly to publish the individual hashes of all of them and have the client VMs download each metadata file separately. Instead, the repository groups metadata files together in a tarball organized by type. For example, one tarball contains all of the trusted packages files, another with all of the pacman files, etc. The hashes of these tarballs are put into the repository

metahash which is pushed to each interested client VM. No matter how many updates the client VM misses, it can examine the hash of the local tarballs and the hashes provided by the repository and determine what needs to be retrieved.

Efficient Transfers The repository makes all of its files available for download through HTTP. However, having each client download its files via separate HTTP connections is prohibitively expensive. The repository therefore supports different transfer mechanisms for better scalability, efficiency, and performance. Some transfer mechanisms are simple (like CoBlitz and Coral) which require no special handling by the repository and others (like BitTorrent) which do.

To support BitTorrent [9] downloads the repository runs a BitTorrent tracker and a modified version of the btlaunchmany daemon provided by BitTorrent. The btlaunchmany daemon monitors a directory for any new or updated files. When a new file is uploaded to the repository it is placed in the monitored directory. When the daemon notices the new file it creates a torrent file that is later seeded. Unique naming is achieved by appending the computed hash of the shared file to the name of the torrent. The torrent file is placed in a public location on the repository for subsequent download by the clients through HTTP.

Client Tools

The client tools are used to manage packages in a client VM and include the stork, pacman, and stork_receive_update commands. The stork tool uses command-line arguments to install, update, and remove packages. Its syntax is similar to apt [2] or yum [36]. The stork tool resolves dependencies and installs additional packages as necessary. It also upgrades and removes packages. The stork tool downloads the latest metadata from package repositories, verifies that packages are trusted by the user's TP file, and only installs trusted files.

Package management with the stork tool is a complex process involving multiple steps including *dependency resolution*, *trust verification*, *download*, and *installation*. For example, consider the installation of the foobar package. Assume foobar depends on a few other packages, such as emacs and glibc, before foobar itself can be installed. In order to perform the installation of foobar, the stork tool must determine whether foobar, emacs, and glibc are already installed on the client and if not, locate candidate versions that satisfy the dependencies. These steps are similar to those performed by other package managers [2, 36, 27]. Finally Stork ensures that those candidates satisfy the trust requirements that the user has specified.

Figure 2 shows a TP file example. This file specifically allows emacs-2.2-5.i386.rpm, several versions of foobar, and customapp-1.0.tar.gz to be installed. Each package listed in the TP file includes the hash of the package, and only packages that match the hashes

may be installed. It trusts the planetlab-v4 user to know the validity of any package it says (this user has a list of hashes of all of the Fedora Core 4 packages). It also trusts the stork user to know the validity of any packages that start with “stork”.

Once satisfactory trusted candidates have been found, Stork downloads the packages from the repository and verifies that the packages it downloaded match the entries in the TP file, including the secure hashes. Finally, the packages themselves are installed.

Package removal is much less complex than installation. Before removing a package, the stork command first checks to see if other packages depend upon the package to be removed. For RPM packages, stork leverages the rpm command and its internal database to check dependencies. Tar packages do not support dependencies at this time and can always be removed. If there are dependencies that would be broken by removal of the package, then stork reports the conflict and exits. Stork removes an installed package by deleting the package’s files and running the uninstall scripts for the package.

The pacman (“package manager”) tool is the entity in a VM that locally enacts centralized administration decisions. The pacman tool invokes the appropriate stork commands based on two configuration files: groups.pacman (Figure 3) and packages.pacman (Figure 4). The groups.pacman file is optional and defines VM groups that can be used by an administrator to manage a set of VMs collectively. The groups.pacman syntax supports basic set operations such as union, intersection, compliment, and difference. For example, an administrator for a service may break their VMs into alpha VMs, beta VMs, and production VMs. This

allows developers to test a new release on alpha VMs (where there are perhaps only internal users) before moving it to the beta VMs group (with beta testers) and finally the the production servers.

```
<GROUPS>
<GROUP NAME="ALPHA">
<INCLUDE NAME="planetlab1.arizona.net"/>
<INCLUDE NAME="planetlab2.arizona.net"/>
</GROUP>
<GROUP NAME="ACCOUNTING">
<INCLUDE NAME="ALPHA"/>
<INCLUDE NAME="p11.unm.edu"/>
</GROUP>
</GROUPS>
```

Figure 3: Example groups.pacman. The “ALPHA” group consists of two machines in Arizona. The “ACCOUNTING” group also includes a machine at the University of New Mexico.

The packages.pacman file specifies which packages should be installed, updated, or removed in the current VM based on a combination of VM name, group, and physical machine. This makes it easy, for example, to specify that a particular package should be installed on all VMs on a physical machine, while another package should only be installed on alpha VMs, etc.

Although pacman can be run manually, typically it is run automatically via one of several mechanisms. First, pacman establishes a connection to the stork_receive_update daemon. This daemon receives the repository metahashes that are pushed by the repository whenever there is an update. Upon receiving this notification, stork_receive_update alerts pacman to the new

```
<?xml version="1.0" encoding="ISO-8859-1" standalone="yes" ?>
<TRUSTEDPACKAGES>
<!-- Trust some packages that the user specifically allows -->
<FILE PATTERN="emacs-2.2-5.i386.rpm" HASH="aed4959915ad09a2b02f384d140c4\
626b0eba732" ACTION="ALLOW"/>
<FILE PATTERN="foobar-1.01.i386.rpm" HASH="16b6d22332963d54e0a034c11376a\
2066005c470" ACTION="ALLOW"/>
<FILE PATTERN="foobar-1.0.i386.rpm" HASH="3945fd48567738a28374c3b238473\
09634ee37fd" ACTION="ALLOW"/>
<FILE PATTERN="simple-1.0.tar.gz" HASH="23434850ba2934c39485d293403e3\
293510fd341" ACTION="ALLOW"/>
<!-- Allow access to the planetlab Fedora Core 4 packages -->
<USER PATTERN="*" USERNAME="planetlab-v4" PUBLICKEY="MFwwDQYJKoZIhvcNAQEB\
BQADSwAwSAJBALtGteQPdLa0kYv+k1FWTk1H9Y7frYh15JV1hgJa5P1GI3yK+R22UsD65_J4P\
V92RUgVd_uJMuB8Q4b1lw4o6JMCAwEAAQ" ACTION="ALLOW"/>
<!-- Allowing the 'stork' user lets stork packages be installed -->
<USER PATTERN="stork*" USERNAME="stork" PUBLICKEY="MFwwDQYJKoZIhvcNAQEBBQADSwAw\
SAJBAKgZCjFKD19ISoclFbuZsQze6bXtu+QYF64TLQ1I9fgEg2CDyGQVOsZ2CaX1ZEZ_O69AYZ\
p8nj+YJLIJM3+W3DMCAwEAAQ" ACTION="ALLOW"/>
</TRUSTEDPACKAGES>
```

Figure 2: Example TP File. This file specifies what packages and users are trusted. Only packages allowed by a TP file may be installed. FILE actions are used to trust individual packages. USER actions allow hierarchical trust by specifying a user whose TP file is included. The signature, timestamp, and duration are not shown and are contained in an XML layer that encapsulates this file.

information. A change to the repository metahash indicates that the repository contents have changed which in turn may change which packages are installed, etc. Second, when `stork_receive_update` is unavailable `pacman` wakes up every 5 minutes and polls the repository for the repository metahash. As before, if there is a discrepancy between the stored data and the described data, `pacman` downloads the updated files. Third, `pacman` also runs when its configuration files change.

The `stork_receive_update` daemon runs in each client VM and keeps the repository's metahash up-to-date. Metadata is received from the repositories using both push and pull. Pushing is the preferred method because it reduces server load, and is accomplished using a multicast tree or publish/subscribe system such as PsEPR [5]. Heartbeats are pushed if no new metahash is available. If `stork_receive_update` doesn't receive a regular heartbeat it polls the repository and downloads new repository metahash if necessary. This download is accomplished using an efficient transfer mechanism from one of Stork's *transfer modules* (discussed further in the transfer modules section). This combination of push and pull provides an efficient, scalable, fault tolerant way of keeping repository information up-to-date in the VMs.

Nest

The Stork nest process enables secure file-sharing between VMs, prevents multiple downloads of the same content by different VMs, and maintains up-to-date repository metadata. It accomplishes these in two ways. First, it operates as a shared cache for its client VMs, allowing metadata and packages to be downloaded once and used by many VMs. Second, it performs package installation on behalf of the VMs, securely sharing read-only package files between multiple VMs that install the package (discussed further in the sharing section). The nest functionality is implemented by the `stork_nest` daemon.

The `stork_nest` daemon is responsible for maintaining connections with its client VMs and processing requests that arrive over those connections (typically via a socket, although this is configurable). A client must first authenticate itself to `stork_nest`. The authentication persists for as long as the connection is established. Once authenticated, the daemon then fields

requests for file transfer and sharing. File transfer operations use the shared cache feature of the repository to provide cached copies of files to the clients. Sharing operations allow the clients to share the contents of packages using the *prepare interface* (discussed further in the section on prepare modules).

Typically, the nest runs on each machine that runs Stork; however, there may be cases where the nest is not run, such as in a desktop machine or a server that does not use VMs. In the case where no nest is running or the nest process fails, the client tools communicate directly with the repository.

Centralized Management Tools

The centralized management tools allow Stork users to manage their VMs without needing to contact the VMs directly. In our example the administrator wanted to install `foobar` automatically on applicable systems under her control rather than logging into them individually. Unlike the client tools that are run in Stork client VMs, the centralized management tools are typically run on the user's desktop machine. They are used to create TP files, `pacman` packages and groups files, the master configuration file, public/private keypairs, etc. These files are used by the client tools to decide what actions to perform on the VM. In addition to managing these files, the centralized management tools also upload metadata and/or packages to the repository, and assist the user in building packages.

The main tool used for centralized management is `storkutil`, a command-line tool that has many different functions including creating public/private key pairs, signing files, extracting metadata from packages, and editing trusted packages, `pacman` packages and groups files. Administrators use this tool to create and modify the files that control the systems under their control. While files can be edited by other tools and then resigned, `storkutil` has the advantage of automatically resigning updated files. After updating these files they are then uploaded to the repository.

Stork on PlanetLab

Stork currently supports the Vserver environment, non-VM machines, and PlanetLab [23, 24]. The PlanetLab environment is significantly different from

```
<PACKAGES>
  <CONFIG SLICE="stork" GROUP="ACCOUNTING">
    <INSTALL PACKAGE="foobar" VERSION="2.2"/>
    <REMOVE PACKAGE="vi"/>
  </CONFIG>
  <CONFIG>
    <UPDATE PACKAGE = "firefox"/>
  </CONFIG>
</PACKAGES>
```

Figure 4: Example packages.pacman. VMs in the slice (a term used to mean a VM on PlanetLab) “stork” and in the group “ACCOUNTING” will have `foobar` 2.2 installed and `vi` removed. All VMs in this user's control will have `firefox` installed and kept up-to-date with the newest version.

the other two, so several extensions to Stork have been provided to better support it.

PlanetLab Overview

PlanetLab consists of over 750 nodes spread around the world that are used for distributed system and network research. Each PlanetLab node runs a custom kernel that superficially resembles the Vserver [18] version of Linux. However there are many isolation, performance, and functionality differences.

The common management unit in PlanetLab is the *slice*, which is a collection of VMs on different nodes that allow the same user(s) to control them. A node typically contains many different VMs from many different slices, and slices typically span many different nodes. The common PlanetLab (mis)usage of the word “slice” means both the collection of similarly managed VMs and an individual VM.

Typical usage patterns on PlanetLab consist of an authorized user creating a new slice and then adding it to one or more nodes. Many slices are used for relatively short periods of time (a week or two) and then removed from nodes (which tears down the VMs on those nodes). It is not uncommon for a group that wants to run an experiment to create and delete a slice that spans hundreds of nodes in the same day. There are relatively loose restrictions as to the number of nodes slices may use and the types of slices that a node may run so it is not uncommon for slices to span all PlanetLab nodes.

Bootstrapping Slices on PlanetLab

New slices on PlanetLab do not have the Stork client tools installed. Since slices are often short-lived and span many nodes, requiring the user to log in and install the Stork client tools on every node in a slice is impractical. Stork makes use of a special initscript to automatically install the Stork client tools in a slice. The initscript is run whenever the VMM software instantiates a VM for the slice on a node. The Stork initscript communicates with the nest on the node and asks the nest to share the Stork client tools with it. If the nest process is not working, the initscript instead retrieves the relevant RPMs securely from the Stork repository.

Centralized Management

Once the Stork client tools are running they need the master configuration file and public key for the slice. Unfortunately the ssh keys that are used by PlanetLab to control slice access are not visible within the slice, so Stork needs to obtain the keys through a different mechanism. Even if the PlanetLab keys were available it is difficult to know which key to use because many users may be able to access the same VM. Even worse, often a different user may want to take control of a slice that was previously managed by another user. Stork’s solution is to store the public key and master configuration file on the Stork repository.

The repository uses PlanetLab Central’s API to validate that users have access to the slices they claim and stores the files in a area accessible by https. The client tools come with the certificate for the Stork repository which pacman and stork use to securely download the public key and master configuration file for the slice. This allows users to change the master configuration file or public key on all nodes by simply adding the appropriate file to the Stork repository.

Modularity

Stork is highly modular and uses several interfaces that allow its functionality to be extended to accommodate new protocols and package types:

Transfer A transfer module implements a transport protocol. It is responsible for retrieving a particular object given the identifier for that object. Transfer protocols currently supported by Stork include CoBlitz [21], BitTorrent [9], Coral [12], HTTP, and FTP.

Share A share module is used by the Stork nest to share files between VMs. It protects files from modification, maps content between slices, and authenticates client slices. Currently Stork supports PlanetLab and Linux VServers. Using an extensible interface allows Stork to be customized to support new VM environments.

Package A package module provides routines that the Stork client tools use to install, remove, and interact with packages. It understands several package formats (RPM, tar) and how to install them in the current system.

Prepare A prepare module prepares packages for sharing. Preparing a package typically involves extracting the files from the package. The Prepare interface differs from the Package interface in that package install scripts are not run and databases (such as the RPM database) are not updated. The nest process uses the prepare module to ready the package files for sharing.

Transfer Modules

Transfer modules are used to download files from the Stork repository. Transfer modules encapsulate the necessary functionality of a particular transfer protocol without having to involve the remainder of Stork with the details.

Each transfer module implements a `retrieve_files` function that takes several parameters including the name of the repository, source directory on the repository, a list of files, and a target directory to place the files in. The transfer module is responsible for opening and managing any connections that it requires to the repositories. A successful call to `retrieve_files` returns a list of the files that were successfully retrieved.

Transfer modules are specified to Stork via an ordered list in the main Stork configuration file. Stork always starts by trying the first transfer module in the

list. If this transfer module should fail or return a file that is old, then Stork moves on to the next module in the list.

Content Retrieval Modules

CoBlitz uses a content distribution network (CDN) called CoDeeN [33] to support large files transfers without modifying the client or server. Each node in the CDN runs a service that is responsible for splitting large files into chunks and reassembling them. This approach not only reduces infrastructure and the need for resource provisioning between services, but can also improve reliability by leveraging the stability of the existing CDN. CoBlitz demonstrates that this approach can be implemented at low cost, and provides efficient transfers even under heavy load.

Similarly, the Coral module uses a peer-to-peer content distribution network that consists of volunteer sites that run CoralCDN. The CoralCDN sites automatically replicate content as a side effect of users accessing it. A file is retrieved via CoralCDN simply by making a small change to the hostname in an object's URL. Then a peer-to-peer DNS layer transparently redirects browsers to nearby participating cache nodes, which in turn cooperate to minimize load on the origin web server. One of the system's key goals is to avoid creating hot spots. It achieves this through Coral [12], a latency-optimized hierarchical indexing infrastructure based on a novel abstraction called a *distributed sloppy hash table* (DSHT).

BitTorrent is a protocol for distributing files. It identifies content by URL and is designed to integrate seamlessly with the web. Its advantage over HTTP is that nodes that download the same file simultaneously also upload portions of the file to each other. This greatly reduces the load on the server and increases scalability. Nodes that upload portions of a file are called *seeds*. BitTorrent employs a *tracker* process to track which portions each seed has and helps clients locate seeds with the portions they need. BitTorrent balances seed loads by having its clients preferentially retrieve unpopular portions, thus creating new seeds for those portions.

Stork also supports traditional protocols such as HTTP and FTP. These protocols contact the repository directly to retrieve the desired data object. It is preferable to use one of the content distribution networks instead of HTTP or FTP as it reduces the repository load.

Stork supports all of these transfer mechanisms with performance data presented in the results section. One key observation is that although these transfer methods are efficient, the uncertainties of the Internet make failure a common case. For this reason the transfer module tries a different transfer mechanism when one fails. For example, if a BitTorrent transfer fails, Stork will attempt CoBlitz, HTTP, or another mechanism until the transfer succeeds or gives up. This

provides efficiency in the common case, and correct handling when there is an error.

Nest Transfer

In addition to the transfer modules listed above, Stork supports a nest transfer module. The nest transfer module provides an additional level of indirection so that the client asks the nest to perform the transfer on its behalf rather than performing the transfer directly. If the nest has a current copy of the requested item in its cache, then it can provide the item directly from the cache. Otherwise, the nest will invoke a transfer module (such as BitTorrent, HTTP, etc.) to retrieve the item, which it will then provide to the client and cache for later use.

Push

Stork supports metadata distribution to the nests using a publish/subscribe system [11]. In a publish/subscribe system, subscribers register their interest in an event and are subsequently notified of events generated by publishers. One such publish/subscribe system is PsEPR [5]. The messaging infrastructure for PsEPR is built on a collection of off-the-shelf instant messaging servers running on PlanetLab. PsEPR publishes events (XML fragments) on channels to which clients subscribe. Behind the scenes PsEPR uses overlay routing to route events among subscribers.

The Stork repository pushes out metadata updates through PsEPR. It also pushes out the repository's metahash file that contains the hashes of the metadata files; this serves as a heartbeat that allows nodes to detect missed updates. In this manner nodes only receive metadata changes as necessary and there is no burden on the repository from unnecessary polling.

Directory Synchronization

In addition to pushing data, Stork also supports a mechanism for pulling the current state from a repository. There are several reasons why this might be necessary, with the most obvious being that the publish/subscribe system is unavailable or has not published data in a timely enough manner. Stork builds upon the transfer modules to create an interface that supports the synchronization of entire directories.

Directory synchronization mirrors a directory hierarchy from the repository to the client. It first downloads the repository's metahash file (the same file that the repository publishes periodically using PsEPR). This file contains a list of all files that comprise the repository's current state and the hashes for those files. Stork compares the hashes to the those of the most recent copies of these files that it has on disk. If a hash does not match, then the file must be re-downloaded using a transfer module.

Share Modules

Virtual machines are a double-edged sword: the isolation they provide can come at the expense of sharing between them. Sharing is used in conventional

systems to provide performance and resource utilization improvements. One example is sharing common application programs and libraries. They are typically installed in a common directory and shared by all users. Only a single copy of each application and library exists on disk and in memory, greatly reducing the demand on these resources. Supporting different versions of the same software is an issue, however. Typically multiple versions cannot be installed in the same common directory without conflicts. Users may have to resort to installing their own private copies, increasing the amount of disk and memory used.

Stork enables sharing in a VM environment by weakening the isolation between VMs to allow file sharing under the control of the nest. Specifically, read-only files can be shared such that individual slices cannot modify the files, although they can be unlinked. This reduces disk and memory consumption. These benefits are gained by all slices that install the same version of a package. It also allows slices to install different package versions in the standard location in their file systems without conflict.

In Stork, sharing is provided via Share modules that hide the details of sharing on different VM platforms. This interface is used by the nest and provides five routines: `init_client`, `authenticate_client`, `share`, `protect`, and `copy`. `init_client` is called when a client binds to the nest, and initializes the per-client state. `authenticate_client` is used by the nest to authenticate the client that has sent a bind request. This is done by mapping a randomly named file into the client's filesystem and asking it to modify the file in a particular way. Only a legitimate client can modify its local file system, and therefore if the client succeeds in modifying the file the nest requested, the nest knows that it is talking to a legitimate client. The `share` routine shares (or unshares) a file or directory between the client and nest, `protect` protects (or unprotects) a file from modification by the client, and `copy` copies a file between the nest and a client.

The implementation of the Share module depends on the underlying platform. On PlanetLab the Share module communicates with a component of the VMM called Proper [20] to perform its operations. The nest runs in an unprivileged slice – all privileged operations, such as sharing, copying, and protecting files, are done via Proper.

On the Vserver platform the nest runs in the root context, giving it full access to all VM file systems and allowing it to do all of its operations directly. Hard links are used to share files between VMs. The immutable bits are used to protect shared files from modification. Directories are shared using `mount --bind`. Copying is easily done because the root context has access to all VM filesystems.

Package Modules

Stork supports the popular package formats RPM and tar. In the future, other package formats such as

Debian may be added. Each type of package is encapsulated in a package module. Each package module implements the following interfaces:

`is_package_understood`. Returns true if this package module understands the specified package type. Stork uses this function to query each package module until a suitable match is found.

`get_package_provides`. Returns a list of dependencies that are provided by a package. This function is used to generate the metadata that is then used to resolve dependencies when installing packages.

`get_packages_requires`. Returns a list of packages that this package requires. This function is used along with `get_package_provides` to generate the package metadata.

`get_package_files`. Returns a list of the files that are contained in a package. This function is also used when generating package metadata.

`get_package_info`. Returns the name, version, release, and size of a package. This information allows the user to install a specific version of a package.

`get_installed_versions`. Given the name of a package, returns a list of the versions of the package that are installed. This function is used to determine when a package is already installed, so that an installation can be aborted, or an upgrade can be performed if the user has requested upgrades.

`execute_transactions`. Stork uses a transaction-based interface to perform package installation, upgrade, and removal. A transaction list is an ordered list of package actions. Each action consists of a type (install, upgrade, remove) and a package name.

Supported Package Types

`stork_rpm`. Stork currently supports RPM and tar packages. The RPM database is maintained internally by the `rpm` command-line tool, and Stork's RPM package module uses this tool to query the database and to execute the install, update, and remove operations,

`stork_tar`. Tar packages are treated differently because Linux does not maintain a database of installed tar packages, nor is there a provision in tar packages for executing install and uninstall scripts. Stork allows users to bundle four scripts, `.preinstall`, `.postinstall`, `.preremove`, `.postremove` that are executed by Stork at the appropriate times during package installation and removal. Stork does not currently support dependency resolution for tar packages, but this would be a straightforward addition. Stork maintains a database that contains the names and versions of tar packages that are installed that mimics the RPM database provided by the `rpm` tool.

Nest Package Installation

A special package manager, `stork_nest_rpm`, is responsible for performing shared installation of RPM packages. Shared installation of tar packages is not supported at this time. Performing a share operation is a three-phase process.

In the first phase, `stork_nest_rpm` calls `stork_rpm` to perform a private installation of the package. This allows the package to be installed atomically using the protections provided by RPM, including executing any install scripts. In the second phase, `stork_nest_rpm` contacts the Stork nest and asks it to prepare the package for sharing. The prepare module is discussed in the following section. Finally, in the third phase `stork_nest_rpm` contacts the nest and instructs it to share the prepared package. The nest uses the applicable share module to perform the sharing. The private versions of files that were installed by `stork_rpm` are replaced by shared versions. Stork does not attempt to share configuration files because these files are often changed by the client installation. Stork also examines files to make sure they are identical prior to replacing a private copy with a shared copy.

Removal of packages that were installed using `stork_nest_rpm` requires no special processing. `stork_nest_rpm` merely submits the appropriate remove actions to `stork_rpm`. The `stork_rpm` module uses the `rpm` tool to uninstall the package, which unlinks the package's files. The link count of the shared files is decremented, but is still nonzero. The shared files persist on the nest and in any other clients that are linked to them.

Prepare Modules

Prepare modules are used by the nest to prepare a package for sharing. In order to share a package, the nest must extract the files in the package. This extraction differs from package installation in that no installation scripts are run, no databases are updated, and the files are not moved to their proper locations. Instead, files are extracted to a sharing directory.

Prepare modules only implement one interface, the `prepare` function. This function takes the name of a package and the destination directory in which to extract the package.

RPM is the only package format that Stork currently shares. The first step of the `stork_rpm_prepare` module is to see if the package has already been prepared. If it has, then nothing needs to be done. If the package has not been prepared, then `stork_rpm_prepare` uses `rpm2cpio` to convert the RPM package into a `cpio` archive that is then extracted. `stork_rpm_prepare` queries the `rpm` tool to determine which files are configuration files and moves the configuration files to a special location so they will not be shared. Finally, `stork_rpm_prepare` sets the appropriate permissions on the files that it has extracted.

Stork Walkthrough

This section illustrates how the Stork components work together to manage packages using the earlier example in which an administrator installs an updated version of the `foobar` package on the VMs the company uses for testing and on the non-VM desktop machines used by the company's developers.

1. The administrator uses `storkutil` to add the new version of the `foobar` package to her TP file.
2. She uses `storkutil` to add the groups `Devel` and `Test` to her `groups.pacman` file, representing the developer's end systems and the testing VMs, respectively. Since groups can be reused, this step most likely would have been done previously.
3. The administrator uses `storkutil` to add a line to her `packages.pacman` file instructing the `Test` group to update `foobar`. She does the same for the `Devel` group.
4. `Storkutil` automatically signed these files with her private key. She now uploads these files to a Stork repository. If the new version of the `foobar` package is not already on the repository she uploads this as well.
5. The repository treats the TP and `pacman` files similarly. The signatures are verified using the administrator's public key that is embedded in the file name. The new files replace the old if their signatures are valid and their timestamps newer. The `foobar` package is stored in a directory whose name is its secure hash. The package metadata is extracted and made available for download.
6. The repository uses the `publish/subscribe` system `PsEPR` to push out a new repository metahash to the VMs.
7. The VMs are running `stork_receive_update` and obtain the new repository metahash. The `stork_receive_update` daemon wakes up the `pacman` daemon.
8. The `pacman` daemon updates its metadata. On non-VM platforms, the files are downloaded efficiently using whatever transfer method is listed in the Stork configuration file. On VM platforms, `pacman` retrieves the files through the nest (which means the files are downloaded only once per physical machine).
9. `Pacman` processes its metadata and if the current VM is in either the `Test` or `Devel` groups it calls `stork` to update the `foobar` package.
10. The `stork` tool verifies that it has the current metadata and configuration files. This is useful because it is not uncommon for several files to be uploaded in short succession. If this is not the case it retrieves the updated files in the same manner as `pacman`.
11. `Stork` verifies that the specified version of `foobar` is not already installed; if it is, `Stork` simply exits.
12. `Stork` searches the package metadata for the specified package. If no candidate is found then it exits with an error message that the package cannot be found. Multiple candidates may be returned if the metadata database contains several versions of `foobar`.
13. `Stork` verifies that the user trusts the candidate versions of `foobar`. It does this by applying the

rules from the user's TP file one at a time until a rule is found that matches each candidate. If the rule is a DENY rule, then the candidate is rejected. If the rule is an ACCEPT rule, then the candidate is deemed trustworthy. The result of trust verification is an ordered list of package candidates.

14. Stork now has one or more possible candidates for foobar. However, if foobar depends on other packages stork repeats steps 13-17 for the dependencies to determine if those dependencies can be satisfied.
15. Stork now has a list of packages that are to be updated, including foobar and its missing dependencies. Stork uses a transfer module to retrieve foobar and dependent packages. The highest priority transfer method is to contact the repository, which is via the nest in VM environments.
16. In a VM environment the nest receives the requests for foobar and its dependencies from the client VM. If these files are already cached on the nest, then the nest provides those local copies. If not, then the nest invokes the transfer modules (BitTorrent, CoBlitz, etc.) to retrieve the files. When retrieval is complete, the nest shares the package with the client VM.
17. Stork now has local copies of foobar and its dependent packages. The client queries the package modules to find one that can install the package. In non-VM environments the `stork_rpm` module installs the packages using RPM and returns to stork which exits. In VM environments the `stork_nest_rpm` module is tried first (stork will fail over and use `stork_rpm` if this module fails). Because foobar is an RPM package, `stork_nest_rpm` can process it. Stork builds a transaction list and passes it to the `execute_transactions` function of `stork_nest_rpm`.
18. In a VM environment the `stork_nest_rpm` module passes the transaction list to `stork_rpm` in order to install a private non-shared copy of the foobar package.
19. In a VM environment the `stork_nest_rpm` module then contacts the nest and issues a request to prepare and share foobar. The nest uses the appropriate prepare module to extract the files contained in foobar. The nest uses the appropriate share module to share the extracted files with the client VM. Sharing overwrites the private versions of the files in the client's VM with shared versions from the foobar package.

In some cases there will be systems that do not receive the PsEPR update. This could occur because PsEPR failed to deliver the message or perhaps because the system is down. If PsEPR failed then `pacman` check for updates every five minutes. If the system was down then when it restarts `pacman` will run. Either way `pacman` will start and obtain a new repository

metahash and the system will continue the process from Step 8.

If nest or module failures happen, stork fails over to other modules that might be able to service the request. For example, if the packages cannot be downloaded by BitTorrent, the tool will instead try another transfer method like CoBlitz as specified in the master configuration file.

Results

Stork was evaluated via several experiments on PlanetLab. The first measures the effectiveness of Stork in conserving disk space when installing packages in VM environments. The second experiment measures the memory savings Stork provides to packages installed in multiple VMs. The final set of experiments measure the impact Stork has on package downloads both in performance and in repository load.

Disk Usage

The first experiment measured the amount of disk space saved by installing packages using Stork versus installing them in client slices individually (Figure 5). These measurements were collected using the 10 most popular packages on a sample of 11 PlanetLab nodes. Some applications consist of two packages: one containing the application and one containing a library used exclusively by the application. For the purpose of this experiment they are treated as a single package.

Rank	Package Name	Disk Space (KB) Standard	Stork	Percent Savings
1	scriptroute	8644	600	93%
2	undns	13240	652	95%
3	chord	64972	1216	98%
4	j2re	61876	34280	45%
5	stork	320	32	90%
6	bind	6884	200	97%
7	file	1288	36	97%
8	make	808	32	96%
9	cpp	3220	44	99%
10	binutils	6892	60	99%

Figure 5: Disk Used by Popular Packages. This table shows the disk space required to install the 10 most popular packages installed by the slices on a sampling of PlanetLab nodes. The *Standard* column shows how much per-slice space the package consumes if nothing is shared. The *Stork* column shows how much per-slice space the package requires when installed by Stork.

For all but one package, Stork reduced the per-client disk space required to install a package by over 90%. It should be noted that the nest stores an entire

copy of the package to which the clients link; Stork's total space savings is therefore a function of the total number of clients sharing a package.

One package, `j2re`, had savings of only 45%. This was because many of the files within the package were themselves inside of archives. The post-install scripts extract these files from the archives. Since the post-install scripts are run by the client, the nest cannot share the extracted files between slices. By repackaging the files so that the extracted files are part of the package, this issue can be avoided.

Memory Usage

Stork also allows processes running in different slices to share memory because they share the underlying executables and libraries (Figure 6). The primary application was run from each package and its memory usage was analyzed. It was not possible to get memory sharing numbers directly from the Linux kernel running on the PlanetLab nodes. Since the PlanetLab kernel shares free memory pages between VMs and there are many VMs being used by different users on each PlanetLab node, this increases the difficulty of gathering accurate memory usage information.

To obtain approximate results the `pmmap` command was used to dump the processes' address spaces. Using the page map data, it is possible to classify memory regions as shared or private. The results are only approximate, however, because the amount of address space shared does not directly correspond to the amount of memory shared as some pages in the address space may not be resident in memory. More accurate measurements require changes to the Linux kernel that are not currently feasible.

Another difficulty in measuring memory use is that it changes as the program runs. Daemon programs

were simply started and measured. Applications that process input files (such as `java` and `make`) were started with a minimal file that goes into an infinite loop. The remaining applications printed their usage information and were measured before they exited.

The resulting measurements show that Stork typically reduces the memory required by additional processes by 50% to 60%. There are two notable exceptions: `named` and `java`. These programs allocate huge data areas that are much larger than their text segments and libraries. Data segments are private, so this shadows any benefits Stork provides in sharing text and libraries.

Package Retrieval

Stork downloads packages to the nest efficiently, in terms of the amount of network bandwidth required, server load, and elapsed time. This was measured by retrieving a 10 MB package simultaneously from 300 nodes (Figure 7), simulating what happens when a new package is stored on the repository. Obviously faulty nodes were not included in the experiments, and a new randomly-generated 10 MB file was used for each test. Each test was run three times and the results averaged. It proved impossible to get all 300 nodes to complete the tests successfully; in some cases some nodes never even started the test. Faulty and unresponsive nodes are not unusual on PlanetLab. This is dealt with by simply reporting the number of nodes that started and completed each test.

Repository load is important to system scalability, represented as the total amount of network traffic generated by the repository. This includes retransmissions, protocol headers, and any other data. For BitTorrent, this includes the traffic for both the tracker and the initial seed as they were run on the same node;

Rank	Package Name	Application Name	Memory (MB)		Percent Savings
			Standard	Stork	
1	<code>scriptroute</code>	<code>srinterpreter</code>	5.8	3.2	45%
2	<code>undns</code>	<code>undns_decode</code>	4.2	2.0	53%
3	<code>chord</code>	<code>adbd</code>	7.6	2.3	70%
3	<code>chord</code>	<code>lsd</code>	7.5	1.1	86%
4	<code>j2re</code>	<code>java</code>	206.8	169.5	18%
5	<code>stork</code>	<code>stork</code>	3.4	1.2	64%
6	<code>bind</code>	<code>named</code>	36.7	32.1	12%
7	<code>file</code>	<code>file</code>	2.6	1.3	50%
8	<code>make</code>	<code>make</code>	2.5	1.1	54%
9	<code>cpp</code>	<code>cpp</code>	2.5	1.2	52%
10	<code>binutils</code>	<code>objdump</code>	3.3	1.4	59%
10	<code>binutils</code>	<code>strip</code>	2.9	1.0	65%
10	<code>binutils</code>	<code>strings</code>	3.4	1.7	50%

Figure 6: Memory Used by Popular Packages. Packages installed by Stork allow slices to share process memory. The *Standard* column shows how much memory is consumed by each process when nothing is shared. With Stork the first process will consume the same amount as the *Standard* column, but additional processes only require the amount shown in the *Stork* column.

running them on different nodes made negligible difference. At a minimum the repository must send 10 MB, since the clients are downloading a 10 MB file. CoBlitz generated the least network traffic, sending 7.8 times the minimum. BitTorrent sent 3.3 times as much data as CoBlitz and Coral sent 5.5 times as much as CoBlitz. HTTP was by far the worst, sending 39.5 times more than CoBlitz. In fact, HTTP exceeded the product of the number of clients and the file size because of protocol headers and retransmissions.

For each test the amount of useful bandwidth each client received (file data exclusive of network protocol headers) is reported, including both the median and mean, as well as the 25th and 75th percentiles. BitTorrent's mean bandwidth is 2.8 times that of CoBlitz, 3.3 times that of HTTP, and 4.2 times that of Coral. HTTP does surprisingly well, which is a result of a relatively high-speed connection from the repository to the PlanetLab nodes.

Figure 8 shows the cumulative distribution of client completion times. More than 78% of the nodes completed the transfer within 90 seconds using BitTorrent, compared to only 40% of the CoBlitz and 23% of the Coral nodes. None of the HTTP nodes finished within 90 seconds.

The distribution of client completion times also varied greatly among the protocols. The time of HTTP varied little between the nodes: there is only an 18% difference between the completion time of the 25th and 75th percentiles. The BitTorrent clients in the 25th percentile finished in 48% the time of clients in the 75th percentile, while Coral clients differed by 64%. CoBlitz had the highest variance, so that the clients in the 25th percentile finished in 14% of the time of the clients in the 75th percentile, meaning that the slowest nodes took 7.3 times as long to download the file as the fastest.

These results reflect how the different protocols download the file. All the nodes begin retrieving the file at the same time. Clients in BitTorrent favor downloading rare portions of the file first, which leads to most of the nodes downloading from each other, rather than from the repository. The CoBlitz and Coral CDN nodes download pieces of the file sequentially. This causes the clients to progress lock-step through the file, all waiting for the CDN node with the next piece of the

file. This places the current CDN node under a heavy load while the other CDN nodes are idle.

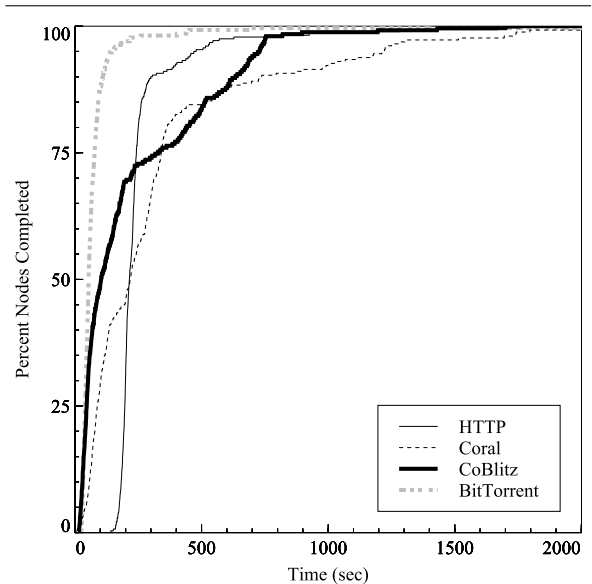


Figure 8: Elapsed Time. This graph shows the cumulative distribution of client completion times. Only nodes that successfully completed are included.

Based on these results Stork uses BitTorrent as its first choice when performing package retrievals, switching to other protocols if it fails. BitTorrent decreased the transfer time by 70% over HTTP and reduces the amount of data that the repository needs to send by 92%.

Related Work

Prior work to address the problem of software management can be roughly classified into three categories: (a) traditional package management systems which resolve package dependencies and retrieve packages from remote systems, (b) techniques to reduce the cost of duplicate installs, and (c) distributed file systems that are used for software distribution.

Traditional Package Management Systems

Popular package management systems [2, 10, 27, 34, 36] typically retrieve packages via HTTP or FTP,

Transfer Protocol	Effective Client Bandwidth (Kbps)				Nodes Completed	Server MB Sent
	25%	Median	Mean	75%		
HTTP	413.9	380.6	321.2	338.1	280/286	3080.3
Coral	651.3	468.9	253.6	234.1	259/281	424.7
CoBlitz	1703.5	737.2	381.1	234.0	255/292	77.9
BitTorrent	2011.8	1482.2	1066.9	1044.0	270/284	255.8

Figure 7: Package Download Performance. This table shows the results of downloading a 10 MB file to 300 nodes. Each result is the average of three tests. The client bandwidth is measured with respect to the amount of file data received, and the mean, median, 25th percentile, and 75th percentile results given. The *Nodes Completed* column shows the number of nodes that started and finished the transfer. The *Server MB Sent* is the amount of network traffic sent to the clients, including protocol headers and retransmissions.

resolve dependencies, and manage packages on the local system. They do not manage packages across multiple machines. This leads to inefficiencies in a distributed VM environment because a service spans multiple physical machines, and each physical machine has multiple VMs. The package management system must span nodes and VMs, otherwise VMs will individually download and install packages, consuming excessive network bandwidth and disk space.

Most package management systems have support for security. In general, however, the repository is trusted to contain valid packages. RPM and Debian packages can be signed by the developer and the signature is verified before the package is installed. This requires the user to have the keys of all developers. In many cases package signatures are not checked by default because of this difficulty. The `trustedpackages` file mechanism in Stork effectively allows multiple signatures per package so that users require fewer keys.

Reducing the Cost of Duplication

Most VMMs focus on providing isolation between VMs, not sharing. However different techniques have been devised to mitigate the disk, memory, and network costs installing duplicate packages.

Disk A good deal of research has gone into preventing duplicate data from consuming additional disk space. For example, many file systems use copy-on-write techniques [6, 8, 14, 15, 16, 30] which allow data to be shared but copied if modified. This allows different “snapshots” of a file system to be taken where the unchanged areas will be shared amongst the “snapshots”. However, this does not combine identical files that were written at different locations (as would happen with multiple VMs downloading the same package).

Some filesystem tools [4] and VMMs [17, 18] share files that have already been created on a system. They unify common files or blocks to reduce the disk space required. This unification happens after the package has been installed; each VM must download and install the package, only to have its copies of the files subsequently replaced with links. Stork avoids this overhead and complexity by linking the files in the first place.

Another technique for reducing the amount of storage space consumed by identical components detects duplicate files and combines them as they are written [25]. This is typically done by using a hash of the file blocks to quickly detect duplicates. Stork avoids the overhead of needing to check file blocks for duplicates on insertion and avoids the need to download the block multiple times in the first place.

Memory There are many proposals that try to reduce the memory overhead of duplicate memory pages. Disco [6] implements copy-on-write memory sharing between VMs which allows not only a process’ memory pages to be shared but also allows duplicate buffer cache pages to be shared. The sharing

provided by Stork is much less effective than Disco, but at a much lower cost.

Stork allows VMs to share the memory used by shared applications and libraries. VMware ESX Server [32] also allows VMs to share memory, but does so based on page content. A background process scans memory looking for multiple copies of the same page. Any redundant copies are eliminated by replacing them with a single copy-on-write page. This allows for more potential sharing than Stork, as any identical pages can be shared, but at the cost of having processes create duplicate pages only to have them culled.

Network Bandwidth A common technique to mitigate the network costs of duplicate data retrieval is to use a proxy server [7, 26, 28, 35]. Proxy servers minimize the load on the server providing the data and also increase the performance of the clients. However, the data still must be transferred multiple times over the network, while the Stork nest provides the data to the client VMs without incurring network traffic (much like each system running its own proxy server for packages). Stork uses techniques such as P2P file dissemination [9] along with proxy based content retrieval [22, 12] to minimize repository load.

Distributed File Systems

Stork uses content distribution mechanisms to download packages to nodes. Alternatively, a distributed file system such as NFS could be used. For example, the relevant software package files could be copied onto a file system that is shared via NFS. There are many drawbacks to this technique including poor performance and the difficulty in supporting different (and existing) packages on separate machines.

Among the numerous distributed files systems Shark [1] and SFS-RO [13] are two that have been promoted as a way to distribute software. Clients can either mount applications and libraries directly, or use the file system to access packages that are installed locally. The former has performance, reliability, and conflict issues; the latter only uses the distributed file system to download packages, which may not be superior to using an efficient content distribution mechanism and does not provide centralized control and management.

Conclusion

Stork provides both efficient inter-VM package sharing and centralized inter-machine package management. When sharing packages between VMs it typically provides over an order of magnitude in disk savings, and about 50% of the memory costs. Additionally, each node needs only download a package once no matter how many VMs install it. This reduces the package transfer time by 70% and reduces the repository load by 92%.

Stork allows groups of VMs to be centrally administered. The `pacman` tool and its configuration

files allow administrators to define groups of VMs and specify which packages are to be installed on which groups. Changes are pushed to the VMs in a timely fashion, and packages are downloaded to the VMs efficiently. Stork has been in use on PlanetLab for over four years and has managed thousands of virtual machines. The source code for Stork may be downloaded from <http://www.cs.arizona.edu/stork>

Acknowledgments

First and foremost, we would like to thank all of the undergraduates who were not coauthors but assisted with the development of Stork including Mario Gonzalez, Thomas Harris, Seth Hollyman, Petr Moravsky, Peter Peterson, Justin Samuel, and Byung Suk Yang. We would also like to thank all of the Stork users. A special thanks goes out to the developers of the services we use including Vivek Pai, KyoungSoo Park, Sean Rhea, Ryan Huebsch, and Robert Adams for their efforts in answering our countless questions. We would especially like to thank Steve Muir at PlanetLab Central for his efforts on our behalf throughout the development of Stork and Proper.

Biographies

Justin Cappos is a Ph. D. student in the Computer Science Department at the University of Arizona. He has been working on projects involving large, real world distributed systems for the past four years. His other research interests include resource allocation, content aggregation, and tools for building distributed systems. He can be reached electronically at justin@cs.arizona.edu.

Scott Baker received a B.S., M.S., and Ph.D. in Computer Science at the University of Arizona. He now works as a software consultant, with a focus in Linux systems programming. He can be reached at bakers@cs.arizona.edu.

Jeremy Plichta is a senior at the University of Arizona majoring in Computer Science, with a minor in Mathematics. After graduating, he plans to pursue a career in industry with the possibility of graduate study at a later date. He designed and maintained the Stork Repository as well as some aspects of the Stork GUI. He can be reached electronically at jplichta@arizona.edu.

Duy Nguyen is currently an undergraduate student at the University of Arizona. He has been working on the Stork Project for a year. His other interests include programming languages, networking, web design, animation, and instructional applications. He can be reached electronically at dqn@email.arizona.edu.

Jason Hardies received a BA in linguistics at the University of Arizona. While a student he worked on the Stork project. After leaving the university in 2006, he joined the healthcare software company Epic Systems, Corp. in Madison, WI where he is a software

developer. He can be reached electronically at jhardies@epicsystems.com.

Matt Borgard is currently an undergraduate at the University of Arizona, studying Computer Science and Creative Writing. His interests include storage, computational linguistics and interactive storytelling. He can be reached electronically at mborgard@email.arizona.edu.

Jeffrey Johnston received a B.S. in Computer Science at the University of Arizona in 2007. He is currently employed at IBM in Tucson, Arizona where he is a software engineer in the z/OS Storage DFSMSHsm department. He can be reached electronically at stork@kidsquid.com.

John H. Hartman is an Associate Professor in the Department of Computer Science at the University of Arizona, which he joined in 1995. He received his Ph.D. in Computer Science from the University of California at Berkeley in 1994. His research interests include distributed file systems, networks, distributed operating systems, and mobile computing. He can be reached electronically at jhh@cs.arizona.edu.

Bibliography

- [1] Annapureddy, S., M. J. Freedman, and D. Mazières, "Shark: Scaling File Servers via Cooperative Caching," *Proceedings 2nd NSDI* Boston, MA, May, 2005.
- [2] *Debian APT tool ported to RedHat Linux*, <http://www.apt-get.org/>.
- [3] Barham, P., B. Dragovic, K. Fraser, S. Hand, T. Harris, A. Ho, R. Neugebauer, I. Pratt, and A. Warfield, "Xen and the Art of Virtualization," *Proceedings 19th SOSP*, Lake George, NY, Oct, 2003.
- [4] Bolosky, W. J., S. Corbin, D. Goebel, and J. R. Douceur, "Single Instance Storage in Windows 2000," *Proceedings 4th USENIX Windows Systems Symposium*, pp. 13-24, Seattle, WA, Aug, 2000.
- [5] Brett, P., R. Knauerhase, M. Bowman, R. Adams, A. Nataraj, J. Sedayao, and M. Spindel, "A Shared Global Event Propagation System to Enable Next Generation Distributed Services," *Proceedings of the 1st Workshop on Real, Large Distributed Systems*, San Francisco, CA, Dec, 2004.
- [6] Bugnion, E., S. Devine, K. Govil, and M. Rosenblum, "Disco: Running Commodity Operating Systems On Scalable Multiprocessors," *ACM Transactions on Computer Systems*, Vol. 15, Num. 4, pp. 412-447, Nov, 1997.
- [7] Chankhunthod, A., P. B. Danzig, C. Neerdaels, M. F. Schwartz, and K. J. Worrell, "A Hierarchical Internet Object Cache," *USENIX Annual Technical Conference*, pp. 153-164, 1996.
- [8] Chutani, S., O. T. Anderson, M. L. Kazar, B. W. Leverett, W. A. Mason, and R. N. Sidebotham,

- “The Episode File System,” *Proceedings of the USENIX Winter 1992 Technical Conference*, pp. 43-60, San Francisco, CA, USA, 1992.
- [9] Cohen, B., “Incentives Build Robustness in BitTorrent,” *Workshop on Economics of Peer-to-Peer Systems*, 2003.
- [10] *Debian – dpkg*, <http://packages.debian.org/stable/base/dpkg>.
- [11] Eugster, P. T., P. Felber, R. Guerraoui, and A.-M. Kermarrec, “The Many Faces of Publish/Subscribe,” *ACM Computing Surveys*, Vol. 35, Num. 2, pp. 114-131, Jun, 2003.
- [12] Freedman, M. J., E. Freudenthal, and D. Mazières, “Democratizing Content Publication with Coral,” *Proceedings 1st NSDI*, San Francisco, CA, Mar., 2004.
- [13] Fu, K., M. F. Kaashoek, and D. Mazières, “The Click Modular Router,” *ACM Transactions on Computer Systems*, Vol. 20, Num. 1, pp. 1-24, Feb, 2002.
- [14] Ghemawat, S., H. Gobioff, and S.-T. Leung, “The Google File System,” *Proceedings 19th SOSP*, Lake George, NY, Oct 2003.
- [15] Hitz, D., J. Lau, and M. Malcolm, “File System Design for an NFS File Server Appliance,” *Proceedings of the USENIX Winter 1994 Technical Conference*, pp. 235-246, San Francisco, CA, USA, 1994.
- [16] Howard, J. H., M. L. Kazar, S. G. Menees, D. A. Nichols, M. Satyanarayanan, R. N. Sidebotham, N., and M. J. West, “Scale and Performance in a Distributed File System,” *ACM Transactions on Computing Systems*, Vol. 6, Num. 1, pp. 51-81, 1988.
- [17] Kamp, P.-H., and R. N. M. Watson, “Jails: Confining the Omnipotent Root,” *Proceedings 2nd International SANE Conference*, Maastricht, The Netherlands, May, 2000.
- [18] *Linux VServers Project*, <http://linux-vserver.org/>.
- [19] Mazières, D., M. Kaminsky, M. F. Kaashoek, and E. Witchel, “Separating Key Management From File System Security,” *Proceedings 17th SOSP*, pp. 124-139, Kiawah Island Resort, SC, Dec, 1999.
- [20] Muir, S., L. Peterson, M. Fiuczynski, J. Cappos, and J. Hartman, “Proper: Privileged Operations in a Virtualised System Environment,” *Proceedings USENIX '05*, Anaheim, CA, Apr, 2005.
- [21] Park, K., and V. S. Pai, “Deploying Large File Transfer on an HTTP Content Distribution Network,” *Proceedings of the 1st Workshop on Real, Large Distributed Systems*, San Francisco, CA, Dec, 2004.
- [22] Park, K., and V. S. Pai, “Scale and Performance in the CoBlitz Large-File Distribution Service,” *Proceedings 3rd NSDI*, San Jose, CA, May, 2005.
- [23] Peterson, L., T. Anderson, D. Culler, and T. Roscoe, “A Blueprint for Introducing Disruptive Technology into the Internet,” *Proceedings HotNets-I*, Princeton, NJ, Oct, 2002.
- [24] *PlanetLab*, <http://www.planet-lab.org>.
- [25] Quinlan, S., and S. Dorward, “Venti: A New Approach to Archival Storage,” *First USENIX Conference on File and Storage Technologies (FAST)*, Monterey, CA, 2002.
- [26] Rabinovich, M., J. Chase, and S. Gadde, “Not All Hits Are Created Equal: Cooperative Proxy Caching Over a Wide-Area Network,” *Computer Networking ISDN Systems*, Vol. 30, pp. 2253-2259, 1998.
- [27] *RPM Package Manager*, <http://www.rpm.org/>.
- [28] *squid: Optimising Web Delivery*, <http://www.squid-cache.org/>.
- [29] *Stork Project*, <http://www.cs.arizona.edu/stork/>.
- [30] Thekkath, C. A., T. Mann, and E. K. Lee, “Frangipani: A Scalable Distributed File System,” *SOSP '97: Proceedings of the Sixteenth ACM Symposium on Operating Systems Principles*, pp. 224-237, ACM Press, New York, NY, USA, 1997.
- [31] *VMWare Workstation*, <http://www.vmware.com/>.
- [32] Waldspurger, C. A., “Memory Resource Management in VMware ESX Server,” *Operating Systems Review*, Vol. 36, pp. 181-194, 2002.
- [33] Wang, L., K. Park, R. Pang, V. Pai., and L. Peterson, “Reliability and Security in the CoDeeN Content Distribution Network,” *Proceedings USENIX '02*, San Francisco, CA, Aug, 2002.
- [34] *Windows Update*, <http://update.windows.com/>.
- [35] *WinProxy*, <http://www.winproxy.com/index.asp>.
- [36] *Yum: Yellow Dog Updater Modified*, <http://linux.duke.edu/projects/yum/>.