



(Photo: deanpalmer.ca)

Pug System Overview

Version: 1.3

February 21, 2016

Copyright © 2013-2016
Packetizer, Inc.



Contents

1	What is Pug?	4
2	Technology Employed.....	4
3	System Overview.....	5
4	Installing Pug.....	7
4.1	Extract the Pug Software	7
4.2	Install MySQL and Create the Pug Database.....	8
4.3	Install AES Crypt	8
4.4	Put the Key File Information in the “Encryption” Table	9
4.5	Define “Locations”	9
4.6	Modify the pug_env File	11
4.7	Modify the config.pl File	11
5	Using Pug.....	14
5.1	Running Pug from cron	14
5.2	Massive Initial Archival	15
5.3	Multiple Instances of Pug.....	15
6	System Operation	15
6.1	Identifying Files to Back Up.....	15
6.2	Marking Files for Future Deletion	16
6.3	Recovering Files	16
6.4	Recovering from a System Failure	17
7	Database Schema.....	19
7.1	“Archive” Table	19
7.2	“ArchivePart” Table.....	20
7.3	“Encryption” Table	20
7.4	“Locations” Table	21
7.5	“Files” Table	22
7.6	“Directories” Table.....	23



Pug System Overview

Legal Information

The Pug software and documentation are copyright © 2013-2014 Packetizer, Inc.

This software is licensed as "freeware". Permission to distribute this software in source and binary forms, including incorporation into other products, is hereby granted without a fee. THIS SOFTWARE IS PROVIDED 'AS IS' AND WITHOUT ANY EXPRESSED OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE. THE AUTHOR SHALL NOT BE HELD LIABLE FOR ANY DAMAGES RESULTING FROM THE USE OF THIS SOFTWARE, EITHER DIRECTLY OR INDIRECTLY, INCLUDING, BUT NOT LIMITED TO, LOSS OF DATA OR DATA BEING RENDERED INACCURATE.

The Pug photo used in this document is copyright © Dean Palmer Photography (deanpalmer.ca) and used with permission. The Pug photo may not be used outside of this document or the Packetizer-owned web site used to distribute this software without explicit written permission from Dean Palmer Photography.

The Pug name, Packetizer name, and Packetizer logo are trademarks of Packetizer, Inc.



1 What is Pug?

Pug™ is cloud-based backup software designed to run on Linux. Inspiration for Pug came from a general disdain for manually handling daily or weekly backup chores, combined with the fact that storage on any external physical medium means that one needs an external storage facility to avoid loss due to fire or other natural disasters.

Many of the existing cloud-backup solutions work by creating a compressed archive of all of the data to be backed up and then pushing that huge file into the cloud. This consumes way more space than is necessary, as files are redundantly stored, and it takes far too long to backup data into the cloud.

Pug takes a very different, more granular approach to the backup problem. With Pug, one defines a set of “locations” (i.e., directories) to archive. Periodically, Pug’s “discover” process will scan those locations looking for new or modified files and will note them in the local Pug database. The Pug “archiver” process also runs periodically to archive any new or modified files in cloud storage, compressing and encrypting each archive before transmission into the cloud. Archive files that are identical are only stored in the cloud one time to avoid wasted space and bandwidth. Thus, from day-to-day, the only files uploaded to cloud storage are new or modified files that are not already stored in the cloud. At any given point in time after the “archiver” runs, all files that exist on your local machine also exist in the cloud.

Pug also allows one to maintain historical copies of files for a specified period of time. Files that are deleted from the local machine remain in cloud storage until Pug is told to expunge the file from cloud storage. As the Pug administrator, you have the ability to see all of the files stored in cloud storage and to retrieve any particular version of the file using the “pget” command.

In the event of a system crash and all local copies of files are lost, one can recover files from Pug using the “precover” command. That command will recover the most recent version of all files that existed on or after the date specified on the command-line.

In short, Pug is an incremental, secure, and efficient cloud archiving solution that eliminates the need for traditional, tedious backup methods employed in the past.

2 Technology Employed

Pug is written entirely in the popular scripting language Perl, meaning that it should be easily ported to any Linux machine without any effort. It was written using Perl 5.14, though it may work with older versions of Perl. There was no intent or desire to use “cutting edge” Perl features, as the primary objective was to ensure portability of code and ease of system recovery in the event of a failure.

Pug relies on MySQL for data storage. It might be possible to utilize other databases, such as SQLite, but no effort was made and no testing has been performed. MySQL was selected for performance reasons and for the fact that tools like phpMyAdmin and MySQL Workbench make it simple to administer.



Pug System Overview

Pug utilizes Amazon S3 (<http://aws.amazon.com/s3/>) for file storage. Amazon offers a truly robust service that virtually guarantees no loss of data. For Pug, this is extremely important. If you prefer a different cloud storage provider, you might be delighted to know that all of the code required to put objects into the cloud and to retrieve objects from the cloud are isolated to a single Perl library that could easily be modified to support any other cloud storage provider. (If you are interested to look at Amazon S3-specific source code, it is found in `lib/cloudutils.pl`.)

Pug also uses AES Crypt™ for file encryption. This software is freely available from <http://www.aescrypt.com> and must be installed, as all files are encrypted with that tool before uploading into the cloud. While Amazon offers secure transmission (using HTTPS) and server-side encryption options, not all cloud services offer the same features. Further, relying on a cloud storage provider to provide encryption means that, should someone gain access to the cloud account, they could access all encrypted files.

3 System Overview

Pug was selected as the name for this software, because there seems to be general consensus among dog owners that pugs are some of the laziest dogs in the world. It was laziness to consistently deal with off-site backups that inspired the creation of Pug. Also, Pug implements a somewhat lazy approach collecting files, archiving them, and is even lazier about getting rid of them once archived. All of these characteristics about Pug are actually good, though, as the end result is a very smooth, transparent, and current archival of all data. Plus, a pug is such a cute dog that it deserved a product named after it.

Pug is comprised of the Pug database (which is document in Section 7) and a set of Perl scripts that look for files to be archived, archive files, retrieve files from storage, and expunge files from storage. The software installation directory is structured as shown below and each file is briefly explained.

Directory	Tool	Description
bin/		
	pget	This is a tool that allows the administrator to get specific files from cloud storage, including previous versions of specific files.
	pharmony	This command will perform a consistency check between cloud storage and the local database to ensure that all files in cloud storage are supposed to be there and that all files in the database are in cloud storage. This command should be used sparingly and only when Pug services are not running. It is a very useful tool for identifying inconsistencies after performing a system recovery or repairing a damaged database.
	plocations	This tool will display the current set of defined “locations” that are to be archived, along with the current scanning and archival status.
	pls	This will provide a listing of all files currently archived by Pug, including multiple versions of archived files (if any).
	precover	This command will allow an administrator to recover all files for a given “location” or a sub-set of those files. This command always retrieves the most recent version of a file from cloud storage.



Pug System Overview

	ptime	This command will display the current Unix time (useful for use with precover) and optionally show the time some number of days prior
sbin/		
	pug_archiver	This program will archive all files marked in the local database by the “discover” process as new or potentially new files. The “archiver” is smart enough to recognize that a given file has already been archived and to make reference to that archived file, versus uploading it again.
	pug_check	This script is called by several other programs to check to make sure that it is safe to perform discovery and archival functions. If this script returns a non-zero status code, the “archiver”, “discover”, and “housekeeper” processes will not run. This is useful to prevent Pug from running when file systems are not mounted or the system time is not synchronized. This script is intended to be modified by the system administrator and it should be observed that the script presently requires ntpd to be running and the time synchronized. If you do not want to force that requirement, comment out that line of code in the installation package.
	pug_db_archiver	This program will archive the local Pug database in cloud storage along with other archived files. In the event of a total system failure, the administrator may retrieve the pug database from cloud storage manually, restore it, and then use the precover command to restore all archived data.
	pug_discover	This program will scan all locations in the database for which scanning is active and look for new or modified files, inserting filenames into the database upon which the “archiver” will take action.
	pug_env	This is a shell script that sets environment variables for the benefit of other Pug scripts. Specifically, it sets the PATH and PERL5LIB environment variables. Use of this is optional, but provided for convenience and intended to be modified for your environment.
	pug_housekeeper	This program should be run periodically to perform various “housekeeping” chores, including expunging files from cloud storage that are no longer needed. This script also performs a number of consistency checks and takes steps to ensure that any errors found in the database are corrected on the next run of the “discover” and “archiver” processes.
	pug_requires	This script exists only to provide verification that all required Perl modules are installed on the system. If any errors are produced when this script is executed, it means some module is missing and must be installed for Pug to operate properly. The script itself contains the complete list of required modules.
docs/		
	Pug System Overview.pdf	This is the document you are reading.
lib/		
	archiveutils.pl	This is a library of functions used primarily by the “archiver”.



	clientutils.pl	This is a library of functions primarily used by the client utilities in the “bin/” directory.
	cloututils.pl	This is a library that contains the functions to interface with the cloud storage service provider (presently, only S3 is supported).
	config.pl	This contains a set of global configuration variables that are used by all Pug software components. This script must be modified to contain values that are appropriate for the environment.
	database.pl	This is a library that contains the functions to connect and disconnect from the database.
	fileutils.pl	This is a library that contains functions for use primary with the “files” table and is primarily used by the “discover” process.
	locations.pl	This is a library that contains functions related to the “locations” table.
	logging.pl	This library contains functions related to logging to syslog.
	utils.pl	This library contains miscellaneous functions, including lock files and string conversions.
schema/		
	pug.sql	This is the database schema file to create an empty “pug” database.

Now, having said so many positive things about Pug, it should be appreciated that Pug is not suitable for all tasks. As a tool for ensuring that all data is permanently archived and that the data can be safely recovered, Pug works well. However, Pug is not suitable for backing up databases of massive size that need to be restored within the hour should there be a system failure. Pug is designed to backup user files. In fact, it *only* backs up user files. Pug is not like “tar” or other commands that backup symbolic links and entire directory structures where those directories do not contain real files. Pug will take note of directory structures, but only directory structures for which it maintains archived files.

4 Installing Pug

4.1 Extract the Pug Software

The first step is to select a suitable location for the Pug software. Given that Pug is intended primarily for use by system administrators and should not be run by any user other than the one authorized to access files related to Pug and the data Pug will archive, it is not advised to install Pug in the usual places like /usr/bin/. Rather, it is recommended to select different location, such as /usr/local/pug and ensure that directory has very restricted permissions (e.g., 0700). In that location, put all of the files in the distribution. For example:

```
$ mkdir /usr/local/pug
$ cd /usr/local/pug
$ tar -xzvf /path/to/pug_distribution.tgz
```

The extraction of the .tgz file will actually create a sub-directory named using the version of the software release. You can use that sub-directory move the contents to /usr/local/pug/ (or whatever directory you choose for Pug).



4.2 Install MySQL and Create the Pug Database

The next step is to create the MySQL database. If MySQL is already in use, this should be a very simple task. If you do not already have MySQL installed and running, you need to do that, first. Most Linux distributions offer MySQL as a package, so the software is readily available. Installing MySQL the first time will require a bit of work, but take notes and save your configuration files. Installing is a second time will be a snap. In any case, we will not go into detail on MySQL installation, maintenance, tuning, etc. There are numerous resources on the Internet that do that already. Documentation abounds here: <http://dev.mysql.com/doc/index.html>. Once you have the MySQL server running, you can create the Pug database with a simple command like this:

```
$ mysql < schema/pug.sql
```

If you do not already have a “user” account defined that Pug can use, be sure to create a database “user” for Pug and give that user full access to the Pug database.

Amazon also has a service called “Amazon RDS” that one could utilize to provide MySQL database services, rather than setting up and managing a database locally. Pug has not been tested with that configuration, but it might be worth exploring if you would prefer to not manage the database installation yourself.

4.3 Install AES Crypt

Download a copy of AES Crypt™ from <http://www.aescrypt.com>. You will need a C compiler like gcc installed to make the AES Crypt binaries. Assuming you are in your home directory and you downloaded “aescrypt-3.0.6b.tar.gz”, here’s what you type with that file in your home directory:

```
$ tar -xzvf aescrypt-3.0.6b.tar.gz
$ cd aescrypt-3.0.6b
$ make
$ make install
```

Only do “make install” if “make” appears to have worked properly. This will install “aescrypt” and “aescrypt_keygen” in /usr/bin so that it is available to all users. To use AES Crypt with Pug, you will need an encryption key file. To create one, do this:

```
$ aescrypt_keygen -g 64 aescrypt.key
```

The number “64” indicates a “password length” and was selected as this generates a password with more than 256-bits of strength. If you wish to have a longer password, you may specify a password length up to 1024 characters long, but it is really overkill. See this page for more information: <http://www.packetizer.com/security/pwgen/>.

Now, copy the key to a secure location (e.g., /usr/local/pug/keys) and ensure that only privileged users have access to this key. You may use this same key for the pug_db_archiver (see discussion later) or you may generate a separate key).

Make sure you do not lose keys. If you lose a key, you will not be able to recover your data!



Keys should be stored in a secure location other than a Pug archive. If someone gains access to the keys, they can decrypt your data (if they can get the encrypted files from cloud storage). And, while you can store copies of key in Pug, you will need to key to get it back. So, keep copies stored elsewhere!

As an example of extreme protection, encrypt your keys using AES Crypt using a really good password (versus the key files which are impossible to memorize) and store the encrypted key file on an SD card stored in a bank vault.

4.4 Put the Key File Information in the “Encryption” Table

Now that you have an encryption key you can use for encrypting files, insert a row into the “encryption” table. There is no tool to do this, so you either have to insert it manually using SQL statements or via a management tool like phpMyAdmin (<http://www.phpmyadmin.net>) or MySQL Workbench (<http://www.mysql.com/products/workbench/>).

The contents of “encryption” table are described in section 7.3. Do make sure you put the current time into the row when creating it. You can use the “ptime” utility that came with Pug to give you the current time value to use.

In general, it’s good to use new keys periodically. When you decide to start using a new key, just generate a new key, put the file in place, and insert a new row. Pug will immediately start using the new key when archiving any new files if that row has a newer timestamp. Note: Old keys cannot be discarded as long as there is a file stored in cloud storage that uses it. Basically, never discard a key.

When putting in the “keyfile” row, ensure you use the fully-qualified pathname to the key file.

4.5 Define “Locations”

Next, you need to define the “locations” you wish to have Pug archive. In Pug, a “location” is basically a “directory” on the system. Examples of “locations” might be /home, /home/bob, and /export/nfs/paris. Pug refers to everything by “location” to help make migration simpler. For example, if one of the “locations” to back up is /export/nfs/paris, where “paris” refers to a machine name, and that remote machine’s name is later changed to “london”, it would be rather frustrating to have to go into the Pug database manually and change all references of “paris” to “london”. With Pug “locations”, making such a change involves 1) stopping Pug (including the “archiver” and “discover” processes) or disabling scanning and archiving for the location (via the “locations” table), 2) unmount “paris”, 3) mount “london” (assumption is this is the same set of files, of course), 4) update the “pathname” field in the “locations” table to the appropriate value, and 5) re-start the Pug software (or re-enabling scanning/archiving on the location).

To define a “location”, insert a row into the “locations” table. The names and meaning of each field in the table is fully documented in 7.4.

Suppose you want to archive the location /home, you want Pug to scan the directory every four hours to look for modified files, you want do not want to archive any files that start with ~ (often used by Microsoft Office products to indicate a temporary office file), you do not want to archive any directories



Pug System Overview

named “tmp” or “temp”, you want Pug to archive any file found 8 hours after it is identified as new or modified, you want to keep the file archived after the user deletes it for 90 days, and you want to allow a maximum of 14 versions to be archived. Let’s also suppose you want to archive a shared directory where users put documents. You do not want to archive any .jpg pictures found there. You want to scan it hourly, since it will likely have more changes than /home. But, since people will work on it all day, you would prefer to not archiving until after 12 hours have past. Also, you want to keep those documents on file for at least one year, but again with a maximum number of versions set at 14. Here’s how you configure a “location” to do just that:

path	scanfreq	fileexcl	pathexcl	archivedelay	expungedelay	maxversions
/home	14400	^~.*	/tmp/,/temp/	28800	7776000	14
/export/shared	3600	^~.*\.jpg\$	/tmp/,/temp/	43200	31536000	14

So what does all of that mean, exactly? Consider the /export/shared location. Scanning the location every hour means that as each hour passes, Pug notes any modified files. Let’s say people in your office work from 8AM to 5PM. Let’s also assume the first change was made at 8:30AM and Pug detects that change at 9AM when it happens to make an hourly scan. That file will not be considered for archival until 9PM. During the 9 hours of the day, perhaps many files are changed. Starting at 9PM, the archiver detects that file and any other files that were detected as changed at the 9AM “discover” run. The “archiver” archives those files and stops. It does this through the night until it has archived the last new or modified file.

Should you prefer the “do it all at once” approach, you can also configure Pug to work that way, too. For example, rather than running the scan each hour, you could run the scanfreq fields to be “1” (every second) and the archivedelay to be “1” (every second), but then only run the “discover” and “archiver” processes once, perhaps having “discover” make a pass over all files at 9PM and then having “archiver” start at 10PM. If you take that approach, do consider that a failure of some sort might mean nothing gets archived for the day. For example, if the “discover” process cannot run because the /export/shared filesystem is offline at 9PM, then the “archiver” will find nothing to archive. If “discover” happened to work, but Amazon is having an outage for an hour or so, then “archiver” might fail. So, at the very least, it would be advisable to let scans happen every hour during the night and let the “archiver” attempt to archive periodically, too.

One thing to note is that you should not have nested “locations”. For example, if you define the “location” /home, that will collect all files in the /home directory and all sub-directories. As such, you should not also define a separate location named /home/bob. Nothing will break if you do this, but it just wastes CPU cycles, requires the “discover” process to scan the /home/bob directory a second time, and inflates the size of the Pug database unnecessarily. It will not result in using more cloud storage space, though, since Pug will recognize that the files found match already-archived files.

An important point to note is that Pug assumes use of UTF-8 in all filenames. It stores all strings in the database in UTF-8 format, too. This is generally a non-issue for those who only use ASCII filenames.



Pug System Overview

However, if you or your users use non-ASCII filenames (including Western European characters that are not ASCII), then it is important to test that those are read by Pug as UTF-8 strings.

If you type “locale” at the Linux command prompt and see some language tag followed by UTF-8, you’re probably not going to have any issues if you see non-ASCII filenames properly with the Linux “ls” command. Even so, it’s worth doing a small test with Pug and also ensuring that Pug has the same Locale settings. Those of particular importance are the environment variables LANG and LC_ALL. Those should be set to “en_US.UTF-8”, for example. These variables are two-part, with the first part indicating the language (which is less important) and the second part indicating the “code page” (UTF-8). It is the UTF-8 “code page” setting that is important.

An easy way to test the behavior of Pug is to create a file that uses non-ASCII characters and ensure that it looks correct in the “files” database and that the “files.pathhash” value is correct for that file. You can use this command on many Linux systems to see the SHA-1 hash of the pathname:

```
echo -n "filename" | shasum
```

For “filename” above, use the value you see in the files.pathname field in the Pug database. The response you get back should match the files.pathhash value. (Note the -n flag on echo means “do not add a trailing newline”. The proper flag used with echo varies by Unix variant, but Linux seems to have settled on this syntax.)

Why is Pug designed to assume only UTF-8? Quite frankly, it is a pain to deal with numerous different character encodings and UTF-8 appears to be the clear winner in terms of universally-agreed character encodings for filenames and the Internet in general.

4.6 Modify the pug_env File

Modify the lib/pug_env file so that it contains the proper values. The PUG_SCRIPTS line should include the PATH to both the “bin/” and “sbin/” directories that contain the Pug software.

The PUG_LIBRARIES line should point to the “lib/” directory where the Pug software is located. The default values are:

```
PUG_SCRIPTS=/usr/local/pug/bin:/usr/local/pug/sbin  
PUG_LIBRARIES=/usr/local/pug/lib
```

4.7 Modify the config.pl File

The lib/config.pl file contains a lot of system-wide configuration parameters that must be set before you try to use Pug. All of the configuration parameters are defined below. Note that all of the variables start with “\$main:”. This is Perl syntax to indicate that the variable is global. Do not change that. Also, if strings are not inside quotes, there may be a reason. Do not add or remove quote marks around parameter values.

Parameter	Description
-----------	-------------



database_user_id	This is the MySQL database “user” identifier that Pug should use.
database_password	This is the MySQL database password associated with the above user ID.
database_name	This is the name of the MySQL database to access.
database_server	This is the hostname of the MySQL database. By default, this value is set to “localhost”, since it is assumed the MySQL database will be local. However, the MySQL database could be installed on an entirely separate machine.
database_use_ssl	If set to 1, then SSL will be used for database connections. A value of 0 indicates that SSL is not used.
aws_access_key_id	This is the Amazon S3 “Key ID” value to use.
aws_secret_access_key	This is the Amazon S3 “secret access key” to use.
aws_s3_bucket_name	This is the name of the Amazon S3 bucket to use.
temp_storage	This is a directory where Pug can create temporary files. This directory should be local (as files are copied to this location, compressed, and encrypted), have sufficient storage space to store at least a few copies of the largest file that will ever be archived, and be secured such that nobody can access this except authorized users. Every file that is archived in cloud storage passes through this directory. Thus, it’s advisable that this directory be owned by “root” with permissions set to 0700. Do not use a directory that might be used for any other purpose.
pug_archive_empty_files	If this value is “Y” then Pug will archive files, even if they have a file size of zero octets. The default is to archive zero-length, but you may set this to “N”.
pug_print_errors_on_stderr	When programs in the sbin directory run, they log messages to syslog. If you would also like to have the output displayed on the screen (useful for when these commands are run manually or when you wish to see output from cron jobs), set this value to “Y”. The default is “N”.
pug_archive_file_prefix	When archiving files to cloud storage, the storage location might be used by other software. This prefix, which MUST be assigned, helps to differentiate Pug files from other files or one server running Pug from another. This may be set to “pug” (default) or the name of a server like “paris”.
pug_archive_file_prefix_separator	Files stored in the cloud are named pug/10353.1, for example. The “pug” at the front is the aforementioned prefix. There is a slash (the default “separator”) followed by numbers that refer to the archive record and file part number. In Amazon S3, this form simulates “folders”. This is purely a matter of preference, but for large sites that use Pug to back up many servers, using the “/” style separator or even a separator like “/<server_name/” might be useful. Using other characters are separators (e.g., “-”) is perfectly valid.



pug_upload_part_size	Uploading to cloud storage can sometimes be challenging when uploading large files. Therefore, Pug breaks files into “parts” that are this size and uploads each “part” separately. That “.1” in the example above is the “part” number. By default, this value is set to 20MB. If Pug uploads a file that is 100MB, it will be done so in 5 parts. This value may be changed at any time, though be mindful of the fact that Pug reads this much data from large files and maintains 2 or 3 copies of that data in RAM. Uploading files using a large “part” size also provide little or no performance gains. For performance reasons, values lower than 8MB should be avoided.
pug_archive_max_sequential_bytes	If this value is greater than zero, it indicates that maximum number of octets to the archiver will store in the cloud before exiting. This value is checked after successfully uploading each complete file, as the archiver will never upload a partial file.
pug_archive_expunge_delay	This parameter tells Pug how long it should wait before expunging an archive file from cloud storage after it detects that the file has been deleted locally and there are no files (including old “deleted” files) referencing it. The default value is one day (86400 seconds), though having a file sit in cloud storage for a week or two is not a bad idea. That said, a longer time period here is not necessary since expunging is better controlled via the “expungedelay” parameter defined for each location in the locations table.
pug_delete_xstatus_delay	As database records related to files, cloud archive files, and archive parts are “deleted”, they are not really deleted. Rather, they are marked with an “X”. This is useful in case the administrator want to look for information related to a specific file or look at the file’s history. However, there is a point where the value diminishes. This parameter indicates how many seconds should pass before the “housekeeper” deletes those records permanently from the database. The default is 31536000 seconds (about one year).
pug_default_directory_create_mask	When Pug needs to re-create a directory as part of the file restoration process, it will first create a directory using this mask and the directory will be owned by the process running the “pget” or “precover” command (e.g., “root”). This value should be restrictive so as to not leak information. The default value is 0700, meaning only the user running pget or precover can access the directory. Proper directory ownership and permissions will be restored if Pug can locate that directory in the Pug database.
pug_db_archive_key	When the database is archived via the pug_db_archiver process, this is the AES Crypt encryption key that will be used to encrypt the file before uploading to cloud storage. Note



	that the encryption key used for all other archive files is extracted from the database. This could be the same key file as used for recovering the Pug data. In any case, it's important to know definitively which key was used for the database file itself. Thus, we declare it to avoid question.
pug_log_app_name	All pug_* processes log to syslog and using the name of the application. This variable exists here only as a reminder and a placeholder. Actual syslog entries would contain things like "pug_archiver", not just "pug".
pug_log_facility	This is the syslog "facility" to use. By default, it is LOG_USER.

5 Using Pug

Once installed and configured, you can start using Pug right away. Since there are several configuration parameters and it is easy to make a mistake the first time, it is recommended that you start using Pug on a small scale (e.g., just one location with a few small/medium size files). Run the "discover" and "archiver" programs manually to see how they work. Delete some local files (do make sure these are not important files!) and test that you can get them back from cloud storage using "pget" or "precover". Use the "pls" command to see the files presently archived in cloud storage.

Modify a file and see that a new version gets archived. The "pls" command should show you two versions of the file.

5.1 Running Pug from cron

Once you feel relatively content that Pug is working for you, consider running these processes via cron by inserting these lines into the crontab file for root:

```
#MIN    HOUR    DOM MON DOW COMMAND
*/15    *       *   *   *   . /path/to/sbin/pug_env; pug_discover
*/15 0-5,18-23 *   *   *   . /path/to/sbin/pug_env; pug_archiver
40      6       *   *   *   . /path/to/pug_env; pug_housekeeper
40      7       *   *   *   . /path/to/pug_env; pug_db_archiver
```

The first command on each line will source the pug_env file to ensure the PATH and PERL5LIB values are set properly, then it will execute pug_discover, pug_archiver, pug_housekeeper, or pug_db_archiver.

The reason for running these commands frequently is so that files are uploaded somewhat evenly during the day, or ideally during the night. When a new or modified file is seen by pug_discover, it inserts a record into the database. However, pug_archiver does not upload that immediately. It is required to wait at least locations.archivedelay seconds before initiating the upload. So, the next time the archiver runs and it sees a file that satisfies this time, it will upload it. If you have a place of business where users create files from 8AM to 6PM, for example, you probably want to insert a 10 hour archivedelay (36000), so that the first file seen at 8AM will not get archived until at least 6PM. Each hour during the night, the archiver selects the next set of files that are eligible for archival. If everything works, all files will be archived before the next business day and uploaded at a steady pace.



Note that directories are not necessarily scanned every 15 minutes. You can control the rate at which a directory is scanned using the `locations.scanfreq` value. So, while “discover” might run every 15 minutes, if there are no directories that are scheduled to be scanned, it will exit immediately.

You can adjust the time cron runs these commands, the `locations.scanfreq`, and `locations.archivedelay` to best suit your needs.

5.2 Massive Initial Archival

Please be advised that if you have tens of thousands or hundreds of thousands of files to archive, the first run of `pug_archiver` might take an incredibly long time to complete. By design, no two instances of `pug_archiver` will run at the same time: it is just a serial process of uploading files to cloud storage.

What you might prefer to do if you have a very large number of files is to “archive” them to local storage. This would be a matter of modifying the `cloututils.pl` library to “store” objects to some local storage device, such as a USB-attached external hard drive. Once all of the files are archived using Pug, the drive can be shipped to Amazon and contents installed in S3. Of course, this would require changes to the Pug code and a good understanding of what Amazon will do with the data, but this is a possible solution to the massive initial archival issue.

5.3 Multiple Instances of Pug

Running multiple, separate instances of Pug are certainly possible. To do that, install the software in some location (e.g., `/usr/local/pug`), but place the `config.pl` file in a different directory. When running any of the Pug programs, you just ensure that the `PERL5LIB` environment variable is set to include the `config.pl` file for the particular instance of Pug you are running. Since the name of the database and all other instance-specific data is isolated to that config file and to the database to which that file refers, then it should be possible to run separate instances.

Note that when the “discover” or “archiver” process run, they create a lock file in the temporary directory specified in the config file. This could cause race conditions for multiple instances of Pug, so you should specify separate temp directories for each. For example `/usr/local/pug/tmp/paris` and `/usr/local/pug/tmp/london` as two different temp directories.

6 System Operation

6.1 Identifying Files to Back Up

The “discover” process (i.e., “`pug_discover`”) is responsible for identifying potential files to be backed up. It does so by searching each defined “location” for new or modified files. It notes any such files in the Pug database and then collects a list of directories and their associated ownership/permissions. The latter is used in restoring proper permissions and ownership should files need to be restored from cloud storage.

Files are inserted into the “files” table with a status of “N”. At some point when the “archiver” (`pug_archiver`) executes, it looks for files with a status of “N” and then checks to see that the current



time is greater than the value of the file's "stime" value (value when the "N" status was applied) plus the `locations.archivedelay` value. If so, then the file is archived.

6.2 Marking Files for Future Deletion

When a file is deleted, the "discover" process notices that, too. When the "discover" process sees a file that is apparently deleted, it marks the file as "deletion scheduled" (status "D") in the database.

However, this does not actually delete the file from cloud storage. If you issue the "pls" command, you will still see the file. It would seem to be available, and that's because it is.

After the file has been scheduled for deletion for more than the length of time specified in `locations.expungedelay`, the "housekeeper" (`pug_housekeeper`) will mark the file as expunged (status "X"). It is only at that point that the file will no longer be visible via "pls".

Note, though, that if there are two identical files on the system, one is deleted and the other is not, there would be only one copy of the file in cloud storage and it would still be marked as archived ("A") in the database. The two files would both share the same "akey" value, which refers to the `archive.skey` field. Only after *all* copies of the same file are deleted will a file be scheduled for deletion from cloud storage.

Once it is determined that there are no longer any files referring to a specific archived file, the archived file will be scheduled for deletion by the "housekeeper". It will not identify files for deletion and delete them in one pass. The next time "housekeeper" runs and after the "`pug_archive_expunge_delay`" (in `config.pl`) delay time has passed, the "housekeeper" will remove the file from cloud storage.

In short, Pug does not get in a hurry to delete files. And, this is the way it should be. One should always have time to recover a file accidentally deleted. And, if you are of the mind that a file should never be deleted, just set the `locations.expungedelay` value to 0 to prevent any files in the specified location from ever being expunged. Likewise, if you wish to prevent any files stored in cloud storage from ever being deleted, set the "`pug_archive_expunge_delay`" value to zero. (Note, though, that if you do this, the only way to map a particular file in the cloud storage to a file on the local system will be through the records marked "X" in the "files" table. So, if you set "`pug_archive_expunge_delay`" to zero, you should also set "`pug_delete_xstatus_delay`" to 0.)

6.3 Recovering Files

If a user deletes a single file or a few files, the easiest way to recover those files is with the "pget" command. The "pget" command has the following syntax:

```
usage: pget { -f <file-id> | -l <location> -p <pathname> } [<new_filename>]
  -f - retrieve a file by File-ID shown in pls
  -l - the "location" from which to retrieve the file
  -p - the "pathname" of the file within that location
```

If you know the File-ID, that's the simplest way to get the file. The File-ID is displayed next to the files shown with the "pls" command. If the target file has a File-ID of "537", then you can recover the file to its original location by issuing this command:



Pug System Overview

```
pget -f 537
```

If you wish to retrieve the file, but place it in a different location, just specify the new pathname as the next item on the command like, like this:

```
pget -f 537 /tmp/somefile.txt
```

If you lose all of the files in a “location” or a subset of those files, you would use the “precover” command. The “precover” command has the following syntax:

```
usage: precover -t <time> -l <location> [<path_prefix>]
      -t - recover all files that were present at this time or later
      -l - the "location" for which files need to be recovered
```

The -t parameter basically means to “recover all files that were present on the system on or after this time”. Do not set this value to zero or some very old time, otherwise the precover command will restore a bunch of files that the user deleted ages ago. The only thing more dissatisfying to a user than losing files is having a bunch of old files reappear that they had long since deleted.

Let’s say you have a location named /export/nfs/paris and the server’s hard drive crashed two days ago. The server is now back online and you want to restore all of the files. It’s reasonable to suggest that we should recover all files that were present on the system as of three days ago. Perhaps we get a few extra “deleted” files, but that is not unusual in any file recovery scenario. But what is the Unix system time “three days ago”? Pug has a utility for that called “ptime”. Just type this:

```
ptime -d 3
```

The program will tell you the current time and what the time was three days ago. Let’s say the time three days ago happened to be “1357118980”. To recover all files in location /export/nfs/paris from three days ago, enter this command:

```
precover -t 1357118980 -l /export/nfs/paris
```

This will cause all archived data files to be pulled out of cloud storage and placed back into their original location. Pug will make an attempt to restore owner, group, permissions, and mtime on each file and owner/group and permissions on each directory.

If only a subset of files, say those in the directory “/export/nfs/paris/documents/contracts/louvre/”, then all of those files could be recovered via this command:

```
precover -t 1357118980 -l /export/nfs/paris documents/contracts/louvre/
```

Note that the “path prefix” part of the command excludes path to the “location”.

6.4 Recovering from a System Failure

To recover from a total system failure, one must re-install the Pug database and restore the data. Since the Pug database is relatively small, it’s strongly advised to archive that at least daily to cloud storage



Pug System Overview

using the `pug_db_archiver` utility. You will note that in the example crontab entries in section 5.1 that this step is taken daily.

Additionally, it's reasonable to have the database replicated to secondary machines. This can be done via using the `mysqldump` command on the Pug database and copying the file from server to server (e.g., via `scp`), employing MySQL's replication functionality, or any other method that works for you.

In the event of a total system disaster, one must restore the database manually. Admittedly, this is perhaps the most painful part of Pug. If you have a dump of the database on a secondary machine, then you can use that file to restore the database.

If you use the `pug_db_archiver`, the archive copy needs to be restored from cloud storage. The location of the archive in cloud storage is part of the configuration parameters `aws_s3_bucket_name`, `pug_archive_file_prefix` and `pug_archive_file_prefix_separator`. Let's assume the bucket name is "mypug", the archive file prefix is "machine1", and the file prefix separator is "/". Then the database file would be stored in Amazon S3 with the following "object" name "mypug/machine1/db.sql.gz.aes.*n*" where *n* is a number indicating the archive part number. Small archives will have a single part, so you will only see the value 1. Larger archives will be broken into a number of parts. Use a tool like `s3cmd` or other to view the files in Amazon S3 to see database archive files and pay close attention to the file timestamps. (If you see an archive part with an older timestamp than the .1 part, then you should ignore those subsequent part numbers as the part is of no value and is leftover from a time when the database was larger in size than the most current backup of the database.) To use `s3cmd` to see the list of files, type this (using the example values):

```
$ s3cmd ls s3://mypug/machine1/db.sql.gz.aes
```

This will show you all of the part numbers. Download each part using a command like this:

```
$ s3cmd get s3://mypug/machine1/db.sql.gz.aes.1
```

As explained above, if you see a higher-numbered part number that has a timestamp older than .1, ignore that part and all subsequent part numbers. You may even safely delete those additional part numbers, but do not do it until you know you have successfully restored the database.

Now, assemble the part numbers in order, like this:

```
$ cat db.sql.gz.aes.1 db.sql.gz.aes.2 > db.sql.gz.aes
```

Now decrypt the file using your archive encryption key (defined by the `pug_db_archive_key` parameter in the configuration file). Let's assume the key is called "mysecret.key". Use "aescrypt" to decrypt and "gunzip" to decompress the file as follows:

```
$ aescrypt -d -k mysecret.key db.sql.gz.aes
$ gunzip db.sql.gz
```

You can use the `db.sql` file to restore the database using the "mysql" command that is part of the MySQL distribution.



Once the database is restored, the Pug software is re-configured and ready to run, you can restore each of the lost locations as described in the previous section.

As a final comment on file recovery, one should not try recovering large numbers of files while at the same Pug is trying to archive files. During any large file recovery operation, it is best to disable Pug. You have several ways to do that:

- Comment out the lines in crontab
- Modify the sbin/pug_check script to always return the value 1
- Modify the locations table of affected locations and set the “scanfreq” and “archivedelay” values to zero

Whatever approach you take, be sure to undo that once the data is recovered so that Pug can again continue archiving files.

7 Database Schema

There are a number of database tables defined to support the Pug software. They are each presented in this section with an explanation of each table and column. Every text field is defined to be a UTF-8 string. All unsigned integers (uint) are 64 bits in length.

7.1 “Archive” Table

The “archive” table is used to store a list of all files that are archived in the cloud. Any two files that have the same hash value are presumed to be identical files and are uploaded only once. Once a new file is added to the “archive” table and uploaded, newly discovered files that match these archived files are merely linked via the relationship “files.akey == archive.skey” and not uploaded again.

Column Name	Type	Description
skey	uint	A unique serial number assigned to each unique file uploaded to the cloud storage
hash	char(40)	The SHA-1 hash of the original file
size	uint	The size of the original file in octets
ekey	uint	Encryption tool and key used to encrypt file (references the “encryption” table)
uhash	char(40)	This SHA-1 hash of the uploaded file, which will be different than the original file due to compression and encryption
usize	uint	Size of the uploaded file, which will be different than the original file due to compression and encryption
status	char(1)	File status, which may be one of: A) Archived D) Deletion scheduled, but the file remains in the archive until the “housekeeper” runs and sees that the archive expunge delay time has passed U) Uploading, not ready for retrieval X) Expunged



stime	uint	Time when the status was last changed (or refreshed in the case where large files are uploaded in parts)
timestamp	uint	Timestamp when this entry was created

7.2 “ArchivePart” Table

If files are of significant size, then uploading an entire file to some cloud storage services in a single unit can fail. Amazon’s S3 service offers a means of breaking object uploads into pieces, but not all services offer this feature and so Pug does not rely on that feature. Pug will handle uploading files in smaller parts that are stored as individual objects in cloud storage. Files larger than the configured “pug_upload_part_size” size will be split into parts. The default value is 20MB, though any reasonable value may be specified (and this value may be changed at any time). When uploading files, the name given to an object part looks like <prefix>-<archive.skey>.<archivepart.part> (e.g., “pug-35565.1”). Note that the entire part is loaded into RAM when uploading or downloading, so the size of “pug_upload_part_size” should not be excessive. Also, the hyphen separator character may be specified in the config file to be something other than a hyphen.

Column Name	Type	Description
skey	uint	A unique serial number assigned to each unique file part uploaded to the cloud storage
akey	uint	The skey of the row in the archive table to which this part belongs
part	uint	The part number (from 1..n) of the uploaded file
hash	char(40)	The SHA-1 hash of this part
partsize	uint	The size of this part
status	char(1)	File status, which may be one of: A) Archived U) Uploading, not ready for retrieval X) Expunged
stime	uint	Time when the status was last changed
timestamp	uint	Timestamp when this entry was created

7.3 “Encryption” Table

This table specifies the encryption tools used and associated keys. Support for a specific tools is built into Pug, with support presently only for AES Crypt. When archiving files, Pug will select the row with the most recent timestamp value. To get the current time from Linux, use the include “ptime” utility.

Column Name	Type	Description
skey	uint	A unique serial number assigned to each combination of encryption tools and encryption keys
tool	char(1)	Encryption tool used to encrypt, which is one of: A) AES Crypt with a 256-bit key (www.aescrypt.com)
keyfile	varchar(256)	The pathname to the key file used for encrypting stored files
timestamp	uint	Timestamp when this row was created



7.4 “Locations” Table

This table specifies all of the locations (directory paths) that should be searched for files to be archived in cloud storage. Pug will periodically make a pass over each location and take note of any new or deleted files. Additions and deletions are noted in the “files” table.

Column Name	Type	Description
key	uint	A unique serial number assigned to each location
path	varchar(256)	The pathname for the location to archive (e.g., “/home”)
scanfreq	uint	Number of seconds that must elapse since lastpass before this location will be searched again (e.g., 3600 for one hour). If this value is zero, the “discover” process will skip this location when looking for files.
lastpass	uint	The time this location was last searched
fileexcl	varchar(256)	<p>This is a comma-separated (“,”) list of regular expressions used to match file names to be excluded from archival. If a comma happens to be in the name of the file to be excluded, it may be preceded by \ to protect it. Since \ is used to protect it, then \\ is required to represent a single \ character. Valid examples include “^Thumbs.db\$” (a specific filename), “.jpg\$” (all files ending with .jpg), “^~.*” (all files that start with a ~ character). If all of these were used together, this field would have the value “^Thumbs.db\$,.jpg\$,^~.*”. These expressions are case-sensitive, as are filenames on Linux systems.</p> <p>Note that it is illegal to use “*.jpg” since “*” is used to match the previous character in a regular expression; syntax as shown previously must be used instead.</p> <p>This column only applies to the file collection process (i.e., determining what files to archive). Pug will produce an error if patterns are invalid.</p> <p>These strings are case-sensitive.</p>
pathexcl	varchar(256)	<p>This is a comma-separated (“,”) list of pathnames to exclude from archival. If a comma happens to be in the name of the file to be excluded, it may be preceded by \ to protect it. Since \ is used to protect it, then \\ is required to represent a single \ character. The pattern match is performed against the relative pathname (including the filename) of the target file relative to the “path” field in this row. If you wish to exclude any pathname that contains “@eaDir” as a subdirectory, then “/@eaDir/” would do that. If you wish to exclude all files ending in .jpg, then “.jpg\$” would do that. To exclude both, this field would contain “/@eaDir/,.jpg\$”.</p>



		<p>Note that it would be illegal to use <code>*.jpg</code> since <code>*</code> is used to match the previous character in a regular expression; syntax as shown previously must be used instead.</p> <p>This column only applies to the file collection process (i.e., determining what files to archive). Pug will produce an error if patterns are invalid.</p> <p>These strings are case-sensitive.</p>
archivedelay	uint	The delay in seconds from the time a file in this location is observed (inserted with an 'N' status) and the time the file will be considered for archival. Having a value equal to a half day or longer will help avoid archiving files that are created temporarily or renamed. If this value is zero, the "archiver" process will skip this location when archiving files.
expungedelay	uint	The delay in seconds that should be imposed after a file has been deleted locally before it is actually expunged from cloud storage (e.g., 86400 for a day, 2592000 for 30 days, 7776000 for 90 days). The purpose for this field is to allow recovery of a deleted file for a period of time from the archive. It should be large enough so as to ensure no file is deleted permanently before recovery can be attempted. If this value is zero, the housekeeper will never expunge files scheduled for deletion files. Rather, they will remain in the "D" state indefinitely.
maxversions	uint	Some files change daily or weekly and, while we may want to retain a deleted file for a certain period of time, we may wish not to retain too many copies. This field indicates the maximum number of "deleted" versions of a file to retain within the "expungedelay" timeframe. If this field is set to 5, for example, then if a sixth file appears with the same name, the oldest copy will be expunged next time the housekeeper operates. If this value is zero, Pug will not enforce a maximum number of versions, thus only the "expungedelay" will be considered when expunging files.
timestamp	uint	Timestamp when this row was created

7.5 "Files" Table

The "files" table is used to record all files that are found and the current status of those files. Details about each file, including user/group ownership, permissions, and modification date are all recorded in this table.

Column Name	Type	Description
skey	uint	A unique serial number assigned to each file
akey	uint	A reference to the skey field in the archive table for this file, or zero if the file has not been archived
location	uint	A reference to the location in which this file is stored



pathname	varchar(256)	This is the relative pathname of the file, relative to the “location” in which the file is stored
pathhash	char(40)	This is an SHA-1 hash of the pathname, which is used to identify all files having the same name
owner	varchar(256)	The name of the owner (converted from the numeric UID)
group	varchar(256)	The name of the group (converted from the numeric GUD)
mode	char(4)	This is a textual representation of the file permissions. It is an octal number where the first digit indicates the SUID(4)/SGID(2)/sticky(1) bits, the second octet is the owner permission for read(4), write(2), execute(1), the third octet is the group permissions, and the last octet is the “other” permissions. For example, “0644” means that the file has owner read/write, group read, and other read permissions.
size	uint	The size of the file in octets
mtime	uint	The file modification time, this along with the file size helps Pug decide when a file on the system may have changed
status	char(1)	File status, which is one of: <ul style="list-style-type: none"> N) New or modified file that needs to be processed by the archiver. If an older version of this file sharing the same pathname exist with in an active (‘A’) state, the archiver will change the status to ‘D’ upon archiving this new version. A) Archived D) Deletion is scheduled, but the file is still available in the archive (files remain in this state until expunged or recovered) R) File has been scheduled for removal, but is not yet expunged. Files scheduled for removal might move to expunged state immediately or might be removed later by the housekeeper. X) File expunged (or otherwise forever gone). If Pug expunges a file permanently from the cloud storage archive, the file status is changed to ‘X’. Likewise, if Pug identifies a new file (‘N’), but it is apparently deleted before the archiver can archive it, the status moves directly from ‘N’ to ‘X’.
stime	uint	Time when the status was last changed, used primarily in deciding when to expunge files (just useful information otherwise)
timestamp	uint	Timestamp when this row was created

7.6 “Directories” Table

The “directories” table is used to record all directories that are along the path of any archived file. This table records the last known owner, group, and permission information so that directories can be re-created with the correct permissions along the path if those directories do not already exist.

Column Name	Type	Description
-------------	------	-------------



Pug System Overview

key	uint	A unique serial number assigned to each directory entry
location	uint	A reference to the location in which this directory exists
pathname	varchar(4096)	This is the relative pathname of the directory, relative to the "location" in which the associated files are stored
pathhash	char(40)	This is an SHA-1 hash of the pathname, which is used to identify all files having the same name
owner	varchar(256)	The name of the owner (converted from the numeric UID)
group	varchar(256)	The name of the group (converted from the numeric GUD)
mode	char(4)	This is a textual representation of the directory permissions. It is an octal number where the first digit indicates the SUID(4)/SGID(2)/sticky(1) bits, the second octet is the owner permission for read(4), write(2), execute(1), the third octet is the group permissions, and the last octet is the "other" permissions. For example, "0755" means that the file has owner read/write/traversal, group read/traversal, and other read/traversal permissions.
timestamp	uint	Timestamp when this row was last updated, which should be each time the "discover" process executes
