

# Using LinkController

---

Michael De La Rue

---

Copyright © 1997-2001 Michael De La Rue

Published by ...

Permission is granted to distribute and change this manual under the terms of the GNU public license.

This is the alpha version of this manual and is very incomplete.

# Introduction

LinkController is a system for checking links.

Most HTML pages contain references to other HTML pages (links). These allow the readers of those pages to locate other related resources (web pages etc.). Unfortunately the location of ‘resources’<sup>1</sup> can change, resources can disappear completely or the system providing the resource can break. When this happens, the link which used to find them will no longer work. The only reliable way to detect this problem is to periodically check over the resources and take corrective action for the ones that have gone missing.

LinkController is designed to make that task much more efficient. It automates the task of checking which links have been broken for a period of time and then of finding which documents they occur in.

LinkController is copyrighted software and is distributed under the GNU Public License which should have been included as the file ‘COPYING’ in the distribution.

---

<sup>1</sup> resource is a general term for HTML pages and all of the other things that can be referenced by URLs



# 1 Getting Started

This section of the manual assumes that the programs that make up LinkController are already installed and working on your computer. If not, then See [Section B.2 \[URLs\], page 47](#). We assume the standard setup where your system administrator runs the link checking for you. Other setups will need slightly different behaviour. Speak to the person who set up link controller.

The first thing to do is to run `configure-link-control` to configure the system. This will ask a series of questions about your configuration and create a configuration file which will be used by the various programs which make up LinkController.

Next you have to work out which links you are interested in. Do this by extracting the links from your web pages (see [Chapter 5 \[Extracting Links\], page 13](#)). The output file from this with the list of links found will be stored in the location you gave during configuration.

Assuming that you have the default install, your links will be automatically copied and checked over time following that.

After a short time (about a day) you will begin to get information about links which didn't work with `link-report --not-perfect`. After some more time (a week or so) you can use `link-report` to find out which links are really broken.



## 2 Configuration

### 2.1 Setting Configuration Variables

The `configure-link-control` program can be used by users to configure LinkController. This will ask you a series of questions and then generate a configuration in your home directory.

The configuration that is controlled by this program is related to reporting and fixing links. For other configuration see See [Chapter 14 \[Administration\]](#), page 31.

### 2.2 Setting Configuration Variables

All of the variable information is stored in the file `‘.link-control.pl’` in your home directory or `‘/etc/link-control.pl’` for system wide configuration. The configuration files are written directly in Perl (the programming language LinkController is written in). You can set the configuration variables by putting lines like this.

```
$::links='/var/lib/link_database.bdbm';
```

Please note the semi colon at the end of the line and the use of single quotes so that Perl doesn't do anything strange to your values.

### 2.3 Link Control Configuration Variables

`$::user_address`

is the email address which the robot declares to the world as it goes around checking links. If you want to check links yourself, you must set this to a valid email address, because if something goes badly wrong, it is the only way for a user at another site to know how to contact you.

`$::base_dir`

This is the base directory for all of the configuration files. If this variable is defined then the other variables will default as given below and do not need to be set individually.

`$::links` tells you what file is being used to store information about links. This could easily be a shared database used by everyone on your system. Defaults to `‘$::base_dir/links.bdbm’`.

`$::schedule`

tells the system where to find the schedule file used to decide which links should be checked next and when that should be. You will need to set this and create the file in order to do link checking. Defaults to `‘$::base_dir/schedule.bdbm’`.

**`$$:page_index`**

tells the system where to find the schedule file used to decide which links should be checked next and when that should be. You will need to set this and create the file in order to do link checking. Defaults to '`$$:base_dir/page_has_link.cdb`'.

**`$$:link_index`**

tells the system where to find the schedule file used to decide which links should be checked next and when that should be. You will need to set this and create the file in order to do link checking. Defaults to '`$$:base_dir/link_on_page.cdb`'.

**`$$:infostrucs`**

This variable points to the configuration file where definitions of infostructures are should be put See [Section 2.4 \[Infostructure Configuration\], page 6](#). Defaults to '`$$:base_dir/infostrucs`'.

**`$$:link_stat_log`**

This variable is the name of a file where important link status changes will be logged. The current definition is links which have just been discovered to be broken. This can be used in email notification. See [Section 7.1 \[Email Reporting\], page 18](#).

## 2.4 Configuring Infostructures

The infostructure configuration is used to find links within the pages when we are building our databases. It is kept in a separate file defined by the `$$:infostrucs` configuration variable.

The format of the file is one line for each infostructure with configuration directives separated by spaces. For example

```
directory http://example.com/manual /var/www/html/manual
www http://example.com/strange_database
```

The first directive describes how `extract-links` program should extract the links. It currently has three possible values. The value `www` means to actually use the given URL to download the web pages. The value `directory` means that `extract-links` should assume that all of the files are stored in a directory and that the directory structure matches the structure of the infostructure. The final value "advanced"

In the case where we use the `directory` directive, a third directive is present on each line with the full path to the base directory of the infostructure.

More advanced configuration is possible by defining the information directly in Perl.

## 3 Advanced Configuration

There are various advanced ways to configure LinkController. These are mostly not needed for simple checking of a small collection of web pages. For larger sites and special situations however, they may well make life much easier.

### 3.1 Advanced Infostructure Configuration

Using more advanced configuration it is possible to skip over certain resources when we are doing link extraction and to ignore some of the links. You may want to skip over this section initially and come back to it only when you find that there are links or pages being checked that you would rather avoid.

For this section, we assume that you already know how to make basic Perl code. If not, then please read through the Perl manual pages ‘perl’, ‘perlsyn’ and ‘perldata’. You may find that the examples given below are sufficient to get you started.

In order to get `extract-links` to extract links using an advanced infostructure, you must use the `advanced` keyword. In the infostructure file. Infostructures not listed there will be ignored, but won’t cause any harm.

Advanced configuration is in the ‘.link-controller.pl’ configuration file by making definitions into the `%::infostrucs` hash. These look like the following

```
$::infostrucs{http://www.mypages.org/} = {
    mode => "directory";
    file_base => "/home/myself/www",
    prune_re => "^(/home/myself/www/statistics)" #ignore referrals
    . "(cgi-bin)", #do CGIs separately
    exclude_re => "\.secret$", #secrets shouldn't get into link database
};

$::infostrucs{http://www.mypages.org/cgi-bin/} = {
    mode => "www";
    exclude_re => "query", #query space is infinite!!
};
```

There are a number of keywords that can be used.

‘mode’ This decides how to download the links. Either ‘www’ or ‘directory’.

‘file\_base’  
If defined, this defines the directory which matches the URL where the infostructure is based. This must be defined if the mode is set to directory.

‘resource\_include\_re’  
If defined, this regular expression must be matched by the *URL* for every resource before links will be extracted from it.

‘resource\_exclude\_re’  
If defined, this regular expression must *not* be matched by the *URL* for every resource before links will be extracted from it.

**'link\_include\_re'**

If defined, this regular expression must be matched by every *URL* found before it will be extracted and saved.

**'link\_exclude\_re'**

If defined, this regular expression must *not* be matched by every *URL* found before it will be extracted and saved.

**'prune\_re'**

Used only in directory mode, this will completely exclude all files and sub-directories of directories matched by the regular expression.

N.B. the exclude and include regular expression can be used together. For a match, the include regular expression must match and the exclude must not match. In other words excludes override includes.

In order for the infostructure to be used by `extract-links` an entry must still be made in the `'infostrucs'` file. For this use the `advanced` keyword. The second argument is a URL used to look up the definition in the `::$infostrucs` hash.

```
advanced  http://www.mypages.org/
advanced  http://www.mypages.org/cgi-bin/
```

The URL used here must match *exactly* the one used in the hash. It is important to note that `'directory'` and `'www'` definitions in the `'infostrucs'` file will override any advanced configuration given.

## 3.2 Authorisation Configuration

One problem when checking links, especially within an intranet situation is that some pages can be protected with basic authentication. In order to extract links from those pages or to simply know that they are there, we have to get through that authentication. By using the advanced Authorisation Configuration we can give LinkController authority to access these pages and allow link checking to work as normal.

*Using this method to allow LinkController to work in an environment with authentication is inherently a security issue since authentication tokens must be stored, effectively in plaintext, in files. This risk may, however, not be much higher than the one that you currently accept, so this can be useful*

We can store the authentication tokens simply in the `%::credentials` hash which we can create in the `'link-controller.pl'` configuration file. The keys in the hash are the exact realm string which will be sent by the web server. Each value of this hash is a hash with a pair of keys. The `'credentials'` key should be associated to the authentication token. The `'uri_re'` key should be a regular expression which matches the web pages you want to visit. For security reasons it shouldn't match any others.

```
::$credentials = {
  my_realm => { uri_re => "https://myhost.example.com",
                credential => "my_secret" }
} );
```

As a minor sanity check, every `uri_re` will be tested against the strings `http://3133t3hax0rs.rhere.` and `http://3133t3hax0rs.rhere.com/secretstuff/www.goodplace.com/`. If they match then those credentials will be disallowed. The owners of `3133t3hax0rs.rhere.com` will just have to hack the code..

For more discussion about the security risks and how to mitigate them see the file `authorisation.pod` included with the LinkController distribution. If you didn't understand the security risk from the above description then probably you should consider avoiding using this mechanism.



## 4 Configuring CGI Programs

The CGI programs have separate configuration variables and configuration.

`'$WWW::Link::Repair::infostrucbase'`

The URL that gets to the base directory of the infostructure.

`'$WWW::Link::Repair::filebase'`

The file-name of the directory which is equivalent to the URL

**FIXME:** this section needs to be rewritten.



## 5 Extracting Links

This section is written assuming that you are using a standard HTML infostructure in a directory or on the World Wide Web

The first part of using link controller is to extract the links. When doing this, a pair of index files is built which list which URLs happen on which pages along with a file listing all of the URLs in the infostructure.

**FIXME:** compare and contrast multi-user configuration with single user

The first stage of the process is done by `extract-links`<sup>1</sup>.

There are two modes for extract links `directory` and `www`. The key difference between them is that the latter actually downloads from a server so it is less efficient but will work in more circumstances and is more likely to represent your site as seen by users. This is assuming that all of your WWW pages are interconnected so it can find them.

**FIXME :** need to describe modes of operation of extract link

`extract-links` creates three files. The first two files (`*.cdb`) are the index files for your infostructure and are located wherever you have configured them to by default they are called `link_on_page.cdb`, `page_has_link.cdb`. The third file is the database file `links.db`. `extract-links` can also optionally create a text file which lists all of the URLs in the infostructure, one per line.

---

<sup>1</sup> the command names in link controller are quite long.. you might want to make your life easier by using command completion which will finish what you have started.. once you've typed a little of the command use `escape escape` or `tab` depending on your shell.. if this doesn't work then you may like to upgrade to a newer shell such as bash or zsh.



## 6 Testing Links

If you are using someone else's link information then you may be able to skip this part and go straight on to the next one on generating reports.

Testing links takes a long time. To begin to detect broken ones will take about ten days. This is a deliberate feature of LinkController. It is designed this way firstly to give people at the other end of links a chance to repair their resources and secondly to reduce the amount of network bandwidth LinkController uses at a given time and so its impact on other people's Internet usage.

The key program which you want to use is `test-link`. I run this from a shell script which directs its output to a log file

**FIXME** actually I now just use a cron job.

```
#!/bin/sh
#this is just a little sample script of how I run the program.

LOGDIR=$HOME/log
test-link >> \
    $LOGDIR/runlog-`/bin/date +%Y-%m-%d`.log 2>&1
#assumes the use of a gnu style date command which can print
#out full dates.
```

And I run this shell script from my 'crontab' with a command like this

```
42 02 * * *    /..directories./run-daily-test.sh
```

The string `/..directories./` should be replaced with the directory where you have the script. Remember to make the script executable.

This will now run until completion each night. However, you should make sure that it does actually finish. If you have too many links to check in the given time, then you can end up with a backlog and the system will take a long time to stop. To avoid this, either make testing less frequent or make checking run faster. This will have to be done by editing the program itself at present.



## 7 Reporting Problems

The easiest way to find out which links are broken is to use the command line interface. The simplest report you can generate is just a list of all the known broken links. Do this like so:

```
link-report
```

On the system I'm testing on right now, this gives:

```
broken:-      file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/cgi
              http://www.ippt.gov.pl/docs-1.4/cgi/examples.html
broken:-      file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/curr
ent/httpd_1.4_irix5.2.Z
              http://www.ippt.gov.pl/docs-1.4/setup/PreExec.html
broken:-      file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/curr
ent/httpd_1.4_linux.Z
              http://www.ippt.gov.pl/docs-1.4/setup/PreExec.html
broken:-      file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/curr
ent/httpd_1.4_osf3.0.Z
              http://www.ippt.gov.pl/docs-1.4/setup/PreExec.html
broken:-      file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/curr
ent/httpd_1.4_solaris2.4.Z
              http://www.ippt.gov.pl/docs-1.4/setup/PreExec.html
broken:-      file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/curr
ent/httpd_1.4_solaris2.4.tar.Z
              http://www.ippt.gov.pl/docs-1.4/setup/PreCompiled.html
Sorry, couldn't find info for url file://ftp.ncsa.uiuc.edu/Web/httpd/U
nix/ncsa_httpd/current/httpd_1.4_source.tar.Z
please remember to check you have put it in full format
broken:-      file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/docu
ments/usage.ps
              http://www.ippt.gov.pl/docs-1.4/postscript-docs/Overview.html
..etc...
```

Which just tells you which links are broken. (In this chapter examples are folded at the 70th row, so that they fit well on narrow screens and in  $\text{\TeX}$ . I wanted to use real text rather than making something up.)

We also know which page they are broken on and can go and look at that on the World Wide Web or directly as a file on the server.

You can get a complete list of options for `link-report` (or any other program which is part of LinkController) using

```
link-report -h
```

For more advanced reporting and editing of documents with broken links you may want to use the Emacs interface (see [Chapter 12 \[Emacs\], page 27](#)).

### 7.1 Email Reporting of Newly Broken Links

It's possible to arrange automatic reporting by email of links which have become newly broken. This is done by getting `test-link` to make a list of links that become broken using the `link_stat_log` variable (see [Section 2.3 \[Link Variables\], page 5](#)) and calling `link-report` to report on those links.

Typically, you may don't want to have a report every time that `test-link` runs, but probably once a day instead. In this case, run a script like the following from your crontab.

```
#!/bin/sh
STAT_LOG=$HOME/link-data/stat-log
WORK=$STAT_LOG.work
EMAIL=me@example.com
mv $STAT_LOG $WORK
if [ -s $WORK ]
then
    link-report --broken --url-file=$STAT_LOG --infostructure |
        mail -s "link-report for 'date'" $EMAIL
fi
```

Every time that this script is run, it will rename the status change log file and then mail a report with all of the new broken links to the specified email address.

In future this feature may be folded into `link-report` directly.

## 8 Examining Individual Files

When you have just written an HTML page, you often want to check it before you put it up for use. You can do this immediately using the `check-page` program. Simply run something like

```
check-page filename.html
```

And it will list all of the links that it is unsure about along with the line number the problem occurred on. This program works particularly well when you editing with Emacs (see [Section 12.2 \[check-page in Emacs\], page 27](#)).



## 9 Repairing Links

The program responsible for repairing links is `fix-link`. It simply accepts two URLs and changes all of the occurrences of the first link in your documents into the second link. It assumes that you have permission to edit all of the problem files and that there is a replacement link. For example

```
fix-link http://www.ed.ac.uk/~mikedlr/climbing/ \  
        http://www.tardis.ed.ac.uk/~mikedlr/climbing/
```

Typed at the shell prompt would have updated the location of my Climbing pages when they moved some while ago and

```
fix-link http://www.tardis.ed.ac.uk/~mikedlr/climbing/ \  
        http://www.tardis.ed.ac.uk/climb/
```

Will change them to the very latest location.

At present, there's no facility for automatically updating the databases when you do this. Instead, you have to run `extract-links` regularly so that new links are noticed. Maybe a later version of LinkController will change this.



## 10 Making Suggestions

A link in the database can have suggestions associated with it. These are normally alternative URLs which somebody or something has decided would make a good replacement for the URL of the Link. Humans can add to the database with the `suggest` program. For example use:

```
suggest file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/current/htt
pd_1.4_linux.Z \
    http://delete.me.org/
Link suggestion accepted. Thank you
```

If you try the same thing again you get

```
suggest file://ftp.ncsa.uiuc.edu/Web/httpd/Unix/ncsa_httpd/current/htt
pd_1.4_linux.Z \
    http://delete.me.org/
Already knew about that suggestion. Thanks though.
```

These suggestions will make it easier for others to repair links, especially if they are using the CGI interface.



## 11 CGI Interface

The CGI interface is not fully developed and has a number of issues related to security to be considered. I have however used it and shown that it can work, so if you want to you could try the same. The two programs `fix-link.cgi` and `link-report.cgi` replace the normal ones `fix-link` and `link-report`. They should be interfaced through an HTML page which feeds the needed information to `link-report.cgi`.

The main security question is how to do authentication of the user. This will have to be set up using the features of the web server.



## 12 The Emacs Interface

LinkController's reporting system is designed to be independent of the interface to it, and often the shell interface will be all that is needed. However another convenient interface is through `emacs`. There are two parts to this integration.

### 12.1 Finding Files with Broken Links

There is a special Emacs mode called `link-report-dired` written for locating files with broken links. The mode is based on `find-dired` and works very similarly. It runs the program `link-report` with an option which makes it list file names in the same way as the `ls` program does. The user can then move around the buffer as normal in Emacs and enter files using a single key press (normally `f`).

### 12.2 Finding Broken Links in Files Within Emacs

The program `check-page` was specially designed so that it outputs in a format which can be read by Emacs' `compile` mode. You can use it within Emacs and then step from error to error correcting them.

To do this, after you have set up your system and completed link checking use the command `M-x compile check-page filename`. You will now see another buffer open up with all of the errors shown there. You can use the key `M-'` (that's a real back quote, not an apostrophe) to step between errors.



## 13 Setting up LinkController

This chapter is aimed at administrators setting up LinkController or who want to have a better understanding of the way that their installation is set up.

The first stage is to actually build and install the programs. This is covered in the document 'INSTALL' which is included with the distribution.

Once you have installed the software, the next step is to configure LinkController so that it knows where you have all of your data. The program `default-install` provides one model of this.



## 14 Administration

There are various aspects of administration. This is mostly related to testing links.

### 14.1 Default Installation

Running `default-install -all` should set everything up correctly. There are various variations on this command which do different things, but the summary is

- Create a linkcont user-id and group which will be used for running programs
- Create a working directory where LinkController will keep it's data
- Create configuration files, especially `/etc/link-control.pl`
- Create cron scripts which will run LinkController automatically.

Using this command it is also possible to activate users and groups e.g. `default-install -user username` or `default-install -group groupname` in which case the specified users will become a member of the linkcont group.

### 14.2 User Administration

User administration is really only needed if you are running link testing centrally for your users. This makes sense since it means that if several users have a link to the same place (likely in any given site) then you will only have to check that link once.

In this case, the important question is which links are copied into the checking database. This is controlled by the program `copy-links-from-users` and decides copies data from users which are in the `lntusr` group.

The command `default-install` can be used to manipulate which users are in the group e.g. `default-install --user username` or `default-install -group groupname` in which case the specified users will become a member of the `lntusr` group.

Another form of user administration is limitation on which users have access to the database. This can be done with normal file permissions. There isn't any specific control to stop users from seeing which links other users have put into the database.

### 14.3 Cron Scripts

In order to be effective, link testing should be done every day. Furthermore, it is a good idea to do the testing at low usage times, which normally means at night. For this reason normally a cron script will be used.

- copy the links from each user with `copy-links-from-users`
- add them to the database with `extract-links --in-url-list`
- build a schedule for testing the links with `build-schedule`

- test the links with `test-link`

The other thing which is done is to remove old links from the database. This only needs to be done weekly.

The program `default-install` can create these scripts.

## 15 Robots and Sensible Behaviour

The most important thing about a program like this is to realise that if you set it up incorrectly and used it in the wrong way, you could upset a large number of people who have set up their web servers in the assumption that they would be used normally by human beings browsing through on Netscape.

Probably it's true that the only way forward is for every WWW site to begin to set up robot defences and detect when someone starts to download from at an unreasonable rate and then cut off the person doing the downloading. I suggest that you don't do this for at least two reasons.

- respect for the person's time
- a wish not to be the person who is cut off

There are probably many other reasons, but that's one for the good side in you and one for the selfish. What more do you need.

For suggestions about what constitutes 'correct' behaviour, it's worth seeing the Robots World Wide Web page. <http://info.webcrawler.com/mak/projects/robots/robots.html>

There are a number of points which make LinkController relatively safe as a link. These are all related to the design and limitations on `test-link`.

- `test-link` does *not* recurs. It only tests links that are specifically listed in the schedule database.
- There is a limit to the number of links that will be tested in one run. This defaults to 1000, but can be configured.
- The schedule for link testing is designed to spread the testing of links across time
- The testing system will not test links at a given site faster than a certain rate.

The last limitation is inherited from the `LWP::RobotUA` module and the documentation for that covers the details of how it works. `test-link` tries to re-order testing of links as needed so that a limit on the rate of visits to one site does not cause a limit on overall testing speed.



## 16 Link Database Maintenance

For the most part the link database shouldn't need much maintenance. There are a number of cases where it might, however. If it becomes corrupt, you may try the `db_recover` command. Probably, however, it's just better to recover the database and link checking schedule from a recent backup. Time is not really critical since the work is normally easy to regenerate. You should make sure that link checking doesn't run in at the time you do backups, however.

The other thing is that occasionally you may want to recover space in the link database by dumping and un-dumping it. See the Berkeley database documentation for more details.

### 16.1 Link Ageing

Sometimes we have a link which is no longer in use within our infostructure. However, it's not a good idea to throw away information about it immediately. It could be that certain files have been temporarily deleted and will come back. Alternatively, a link could have been found broken and been corrected, but someone has a copy of the old page. When they re-install the copy, we will have to deal with that link again.

If we had not kept that link around and they immediately do a link check on their document they may see nothing wrong, and, because of the nature of link-controller, we won't start reporting the link as broken until we have confirmed that it really is, some days later.

If, on the other hand, we keep checking all of the links which are aged, we will cause ourselves considerable extra work.

To handle this, links are aged. Once a link has reached greater than a certain age, the link will not be checked any more. Once the link has reached a much larger age, it will be completely deleted from the link database. The age is reset each time the link is extracted from the infostructure, so links which are still in use will continue to be checked.

In the meantime, the link will still be repeatedly scheduled based to it's normal checking time. This causes us to examine it quite regularly, but that is okay, since we will do almost no work when we look at it.

By default links which have not been refreshed will be ignored after one week and will be deleted from the database after two months.



## 17 Bugs and bug reporting

This version of LinkController is still in early development. There are many changes to come. Undoubtedly there are many bugs in the software already and will soon be more.

A bug is when

- the software doesn't do something the documentation says it should
- the software does something the documentation says it shouldn't
- the software does something surprising and that isn't documented
- the software does something strange but the documentation doesn't explain why
- it is difficult or impossible to understand what the documentation is trying to say

some of these mean fixing the documentation and some the software. All of them are bugs and should be reported and fixed.

If you find a bug, I will be grateful to hear about it. Even if you don't know how to fix it or anything, it is useful to know what is wrong so that other people don't get caught out but *read the BUGS file first* please. If the bug is listed there then the only useful thing that you can do is fix it. If you do this and contribute it to me then that is very useful.

When you report a bug, please tell me what release of link controller you were using. This is the number which was in the name of the file that LinkController came in. If your problem was with a specific program, please also run '`program --version`' and send the output. This tells me exactly which version of that program you were running.

Since this is a developers release, I'd hope most users would be able to make some level of fixes. If you do this, send me context differences (use '`diff -u`' if it works or try '`diff -c`' otherwise). I use CVS, so as long as I know which version you have I will be able to find the original file and see your changes. However it's also important to explain them because I won't be able to use them unless I (relatively stupid computer type) understand them.

Send bug reports to the address you get by changing words into punctuation in the following.

```
link minus controller at scotclimb dot org dot uk
```

This mailing address is sent only to me right now, but may become a list in future. Use my (Michael De La Rue) personal address to contact me please. N.B. I am **extremely** inefficient about answering email. Don't worry if you don't get a reply.

The ideas, and history



## 18 History

LinkController was originally inspired by MOMspider and having the MOMspider code available was very useful when starting the creation of this kit, but, it shares almost no code with MOMspider, other than what has come to it from the LibWWW-Perl library.

Philosophically, the MOMspider heritage is obvious in the wish to handle big jobs efficiently. In the working practice there are far more differences than similarities, partly caused by Perl language changes.

I decided to completely separate the exploration of the local infostructure, looking for links to be checked, from the actual checking process. This means that checking can be spread over a large number of days and still run efficiently.

The basic aim of this link checking kit is to be able to efficiently handle any size of link checking job. At the bottom end we have checking new pages as they are written. Here we want to use information from previous checks to avoid having to check all of each page every time. At the other end we have massive info structures (sites) which deal in many thousands of links and could not possibly all be checked in one day. For this latter case the aim is to be able to efficiently spread the link checking load into all available low usage periods.

My primary aim in writing this was not to write very efficient code for the small scale case (takes minimum time to do everything), but rather code which would scale well. If your system can check 1000 links in two days, it will hopefully be able to check almost 7000 links in two weeks. I'm trying to make sure all data structures which grow with the number of links are kept on disk.



## 19 Acknowledgements

Although I wrote this system by myself, this would not have been nearly as easy and almost certainly wouldn't have been finished without the help of the following people and organisations.

### Esoterica Internet Portugal

Esoterica provided me with full access to the Internet in Portugal and use of their computers for free which allowed me to keep up on both this software and the Linux Access HOWTO. In particular I'd like to thank all of the members of staff who helped me very much. These people include Mario Francisco Valente (the instigator of Mini Linux) who first agreed to me using their kit, set me up to use their machines, and along with Luis Sequeira provided a sounding board for some ideas. Luis also provided the odd lift home in the evening. Also Martim de Magalhaes Pereira and Mr Mendes. See them all on

<http://www.esoterica.pt/esoterica/quemsomos.html>

For more about esoterica (Internet Services in Portugal) see:

<http://www.esoterica.pt/esoterica/>

These pages are in Portugese<sup>1</sup> of course.

### IPPT PAN Poland

Thanks go to IPPT PAN (part of PAN - Polska Akademia Naukowa) in Poland and in particular Piotr Pogorzelski who allowed me use of facilities for testing this software, provided a willing victim for having his web pages tested and made a number of suggestions which have been incorporated into the software.

### The Tardis Project

Supported by the Computing Science department of the University of Edinburgh, the Tardis project provides an experimental framework in which students, former students and other related people to do their own work on fully Internet connected Unix and Linux hosts.

The use of the facilities of the Tardis Project has made it much easier for me to develop software like this. In particular, the large amount of disk space the administrators have allow me to use is very useful.

---

<sup>1</sup> Whilst the above names are mangled here. See the correct versions in the original texinfo or on the Web pages.

## Other Free Software Authors

It is through the software provided by the Free Software Foundation (such as the `gcc` C compiler, Emacs, the file utilities), the authors of the various packages which make up a working Linux System (Linux by Linus Torvalds, Alan Cox, etc.... filesystems and support by Theodore Tytso, Stefan Tweedie etc.. Linux-Libc by HJ Lu, based on GNU `glibc` from the FSF.. the list is indefinite) and the authors of Perl and its modules, especially Gisle Aas and Martijn Kostler for LibWWW-Perl that I was able to set this up.

I'd particularly like to thank Tim Goodwin the author of the Perl CDB module who made and accepted a number of alterations to that, at my request. These alterations made this package simpler to write and easier to maintain.

The Free Software Foundation web pages are at

<http://www.gnu.ai.mit.edu/>

## Appendix A Invoking the LinkController Programs

Because they use the Perl `Getopt::Mixed` module, all of the LinkController command line programs respond to the standard POSIX style command line options. At least the following two options will be implemented.

`--help` This option will give a list of all of the options understood by the program along with brief explanations of what they do. At present it may lie a little, but that is being corrected.

`--version` This option will give some version information for the program.

You can use the `--help` option to get help on each program, for example:

```
extract-links --help
```

will give something like

```
extract-links [arguments] [url-base [file-base]]
```

```
-V --version          Give version information for this program
-h --help --usage    Describe usage of this program.
  --help-opt=OPTION  Give help information for a given option
-v --verbose[=VERBOSITY] Give information about what the program is doing. Set
value to control what information is given.
```

```
-e --exclude-regex=REGEX Exclude expression for excluding files.
-p --prune-regex=REGEX Regular expression for excluding entire directories.
-d --default-infostrucs handle all default infostrucs (as well as ones listed
on command line)
```

```
-l --link-database=FILENAME Database to create link records into.
-c --config-file=FILENAME Load in an additional configuration file
```

```
-o --out-url-list=FILENAME File to output the url of each link found to
-i --in-url-list=FILENAME File to input urls from to create links
```

Extract the link and index information from a directory containing HTML files or from a set of WWW pages with URLs which begin with the given URL and which can be found by starting from that URL and searching other such pages.

You can then use that information to get the program to do what you want.



## 20 Packages Which Can Be Useful with LinkController

### 20.1 The CDB utilities

In order to have LinkController working you must have installed these. It is worth looking at the utilities that are provided, especially `cdbdump` which will let you look at the contents of the file. You should be aware that `cdbget` program which is provided *won't* be able to get at the full contents of the index files since they contain repeated keys.

More information on cdb and new releases can be got from the www page.

<http://pobox.com/~djb/cdb.html>

### 20.2 The Tie-Transact-Hash Perl Module and Programmes

This is a Perl module written by myself which includes a program which allows direct examination and editing of Berkeley databases. It can be useful for debugging and correcting problems in the LinkController Link database or schedule file.

Tie::TransactHash can be downloaded from CPAN, the Comprehensive Perl Archive Network get there via:

<http://www.perl.com/perl/info/software.html>



## Appendix B Terms

### B.1 Resource

A resource is almost anything. ‘It’ can range from a person to an HTML file to a computer to a database or presumably eventually to phone numbers, possibly physical hardware. This generality is a very important concept for the World Wide Web. Really the key thing about a resource is that it can be ‘identified’. See [Section B.2 \[URLs\], page 47](#), for more details.

### B.2 URLs

A URL or ‘Uniform Resource Locator’ are the essence of the World Wide Web. Approximately, they are addresses through which ‘resources’ can be located. The idea is that almost anything can be given some kind of address in a form that a machine can work with. By defining a set of rules, this can then be converted into a URL. A URL has two parts. The first tells us what rules to use and the second tells us what the address is.

### B.3 Infostructure

An infostructure is a concept which was introduced in Link Checking in the MOMspider package. It is a collection of related resources.

### B.4 Link

The term link in LinkController is used for a connection between two resources. It’s existence really comes from the ‘class’ or piece of type of computer data which is used to store information about ‘links’. Properties of a link include:

- Knowing what the URL of the target resource of the connection is.
- Knowing whether the connection to the target resource has been working recently.
- Knowing when the connection to the target resource was last checked.

Within the programs, a link is different from a URL in that it is specifically aimed at checking connections, where a URL just specifies what the connection should be if it is working.



## Program and Variable Name Index

Perl variables, which use `::` to separate the different parts, are separated here using `/` instead to make it easier to work in info with this file.

### \$

<code>\$base_dir</code> .....	5
<code>\$infostrucs</code> .....	6
<code>\$Link/Repair/filebase</code> .....	11
<code>\$Link/Repair/infostrucbase</code> .....	11
<code>\$link_index</code> .....	6
<code>\$links</code> .....	5
<code>\$page_index</code> .....	6
<code>\$schedule</code> .....	5
<code>\$user_address</code> .....	5

### C

check-page, in emacs.....	27
---------------------------	----

### E

extract-links.....	13
--------------------	----

### F

fix-link.....	21
fix-link.cgi.....	25

### L

link-report.....	17
link-report-dired.....	18, 27
link-report.cgi.....	25

### S

suggest.....	23
--------------	----

### T

test-link, using.....	15
-----------------------	----



# Concept Index

## A

acknowledgements .....	41
authentication .....	25
authentication, basic .....	8
authorisation .....	8

## B

bandwidth .....	33
basic authentication .....	8
broken links, finding .....	17
broken links, finding in Emacs .....	18, 27

## C

CGI interface .....	25
CGI, configuration .....	11
checking individual pages .....	19
configuration .....	5
configuration variables .....	5
configuration, infostructure .....	6
crontab, example .....	15

## D

dangers .....	33
---------------	----

## E

Emacs interface .....	27
extracting links .....	13

## F

file, individual, checking .....	19
filtering links .....	7

## H

history MOMspider .....	39
HTML pages, groups of .....	47

## I

infostructure .....	47
infostructure, configuration .....	7
interface, CGI .....	25
interface, Emacs .....	27

## L

link index, creating .....	13
links, examining .....	17
links, examining, in Emacs .....	18, 27
links, extracting .....	13
links, repairing .....	21

## P

page index, creating .....	13
page, checking .....	19

## R

regular expression, exclude .....	7
regular expression, include .....	7
repairing links .....	21
reports .....	17
resource .....	47
robots .....	33

## S

security .....	8
setting variables .....	5
suggestions .....	23

## U

URL .....	47
-----------	----

## V

variables, configuration .....	5
variables, setting .....	5

## W

WWW .....	47
-----------	----

