

.htaccess and other oddities

Website Planning

What Are those files?

On the right is the file listing from the root directory of a website as seen in a FTP client. You may recognise *index.php* as being the website homepage, but what are all the other files?

This presentation aims to explain what they are and how they're used.

Name	Size	Type
accessibility		File folder
contact		File folder
core-competencies		File folder
core-courses		File folder
design-principles		File folder
error-files		File folder
faq		File folder
forum		File folder
includes		File folder
our-philosophy		File folder
our-students		File folder
preparing-for-study		File folder
programme-details		File folder
site-map		File folder
style		File folder
teaching-team		File folder
web-design-books		File folder
webteachingday		File folder
.htaccess	1 KB	HTACCESS File
favicon.ico	23 KB	Icon
google1abe8c03c06acc43.html	1 KB	Firefox HTML Document
index.php	7 KB	PHP Script
robots.txt	2 KB	Text Document
sitemap.xml	5 KB	XML Document

Summary

- ▶ [.htaccess](#) (hypertext access)
 - ▶ • custom error pages
 - ▶ • password protection
 - ▶ • redirects from one file to another
 - ▶ • rewriting URLs
 - ▶ • hot link prevention
 - ▶ • deny access
- ▶ [sitemap.xml](#) (Google sitemap)
- ▶ [robots.txt](#) (disallow crawling)
- ▶ [humans.txt](#) (credit the makers)
- ▶ [favicon.ico](#) (favourites icon)

Website Planning

THE .htaccess FILE

What is a .htaccess file?

- .htaccess is a localised server configuration file that can be used to override default server configuration settings.
- Originally, the file's primary purpose was to facilitate password protection to web folders; hence the name (hypertext access).
- On modern servers, .htaccess can be used to perform a range of tasks, including...

What can .htaccess do?

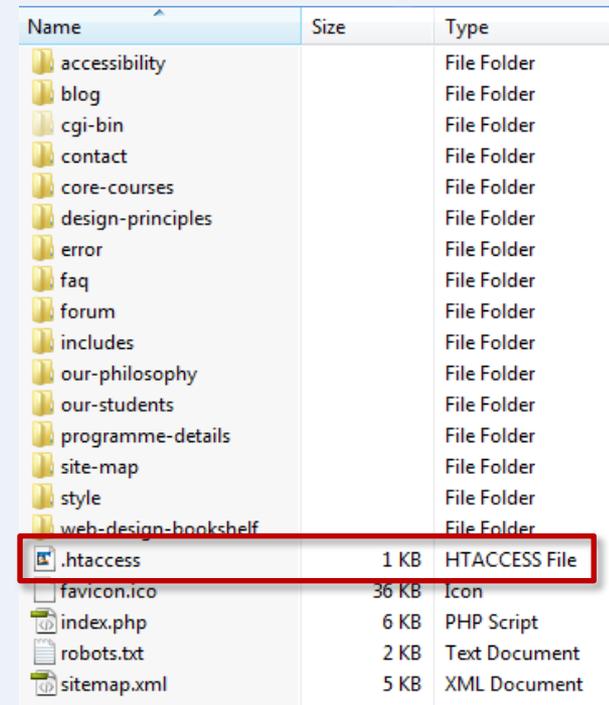
- **Custom Error Pages** – configure the use of custom error pages (e.g. 404 “page not found”).
- **Password Protection** – in combination with a .htpasswd file (containing encrypted username and password).
- **Redirection** – can redirect requests for one page or one folder to another (useful if your site changes).

What can .htaccess do?

- **Rewrite URLs** – for consistency and for the benefit of search engines you can decide whether your site uses “www” or not. This is known as *URL Canonicalization*.
- **Prevent Hotlinking** – can prevent your web content (usually images) from being embedded in sites outside of your server.
- **Deny access** – block access to your website from specific IP addresses.
- And a great deal more.

Where does .htaccess live?

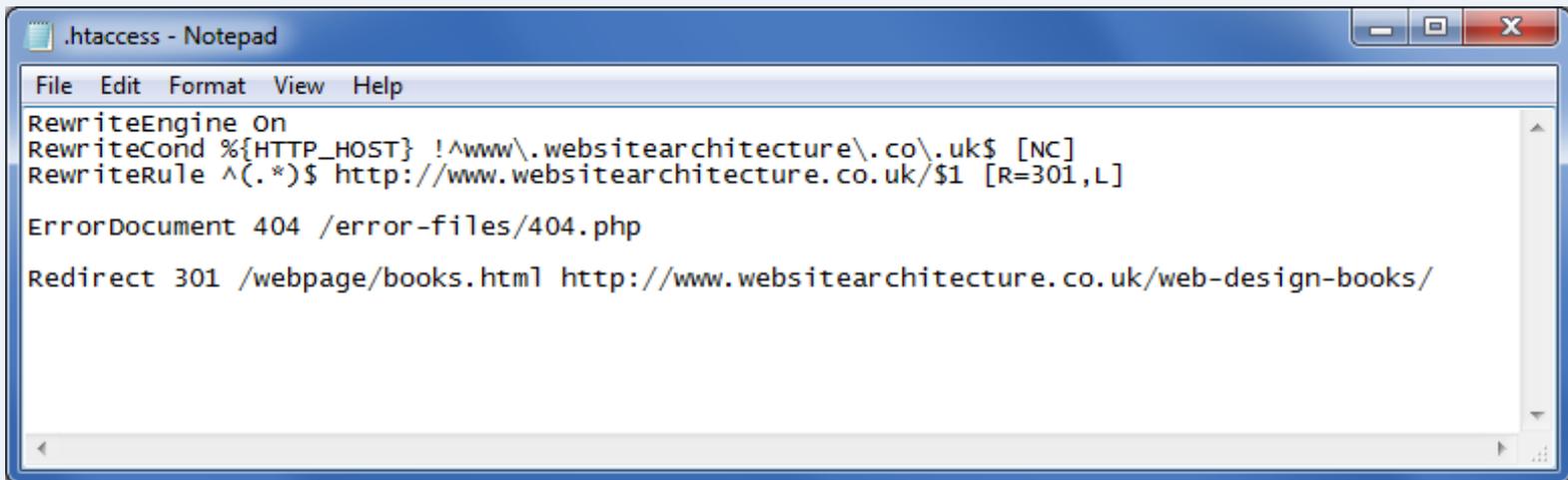
- Websites do not need a .htaccess file but if they exist, they are placed in the root folder (using FTP).
- There may be additional .htaccess files if password protection is used. Each secure folder will have its own .htaccess file.
- The leading dot tells the web server that this is a hidden file, so you may need to tell your FTP client to display hidden files before you can see it.



Name	Size	Type
accessibility		File Folder
blog		File Folder
cgi-bin		File Folder
contact		File Folder
core-courses		File Folder
design-principles		File Folder
error		File Folder
faq		File Folder
forum		File Folder
includes		File Folder
our-philosophy		File Folder
our-students		File Folder
programme-details		File Folder
site-map		File Folder
style		File Folder
web-design-bookshelf		File Folder
.htaccess	1 KB	HTACCESS File
favicon.ico	36 KB	Icon
index.php	6 KB	PHP Script
robots.txt	2 KB	Text Document
sitemap.xml	5 KB	XML Document

What does .htaccess look like?

- .htaccess files are simple ASCII text files and can be viewed and edited in any text editor, even Notepad.

A screenshot of a Notepad window titled ".htaccess - Notepad". The window contains the following text:

```
File Edit Format View Help
RewriteEngine On
RewriteCond %{HTTP_HOST} !^www\.websitearchitecture\.co\.uk$ [NC]
RewriteRule ^(.*)$ http://www.websitearchitecture.co.uk/$1 [R=301,L]

ErrorDocument 404 /error-files/404.php

Redirect 301 /webpage/books.html http://www.websitearchitecture.co.uk/web-design-books/
```

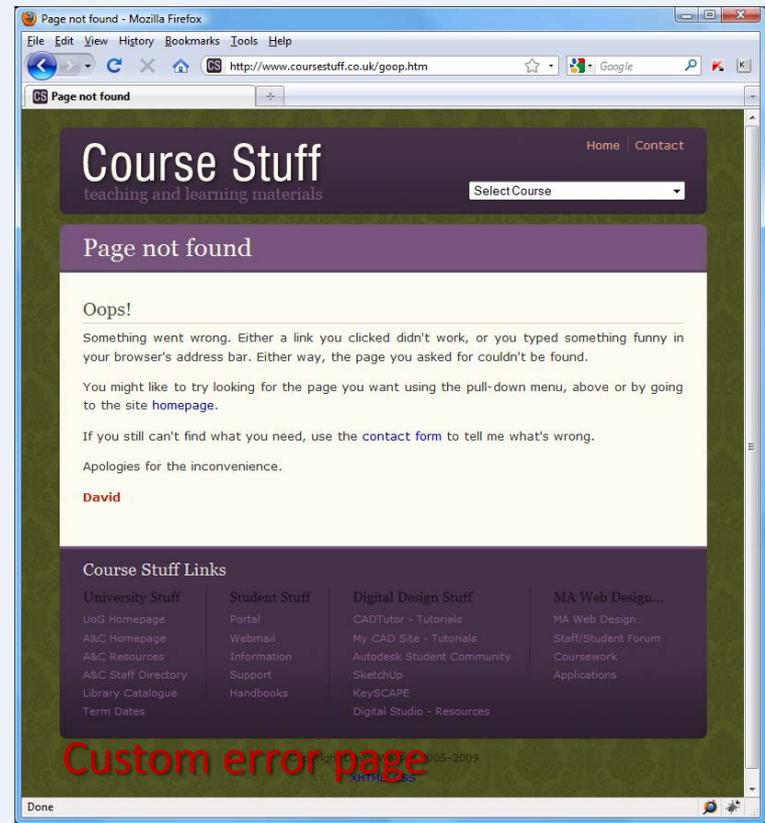
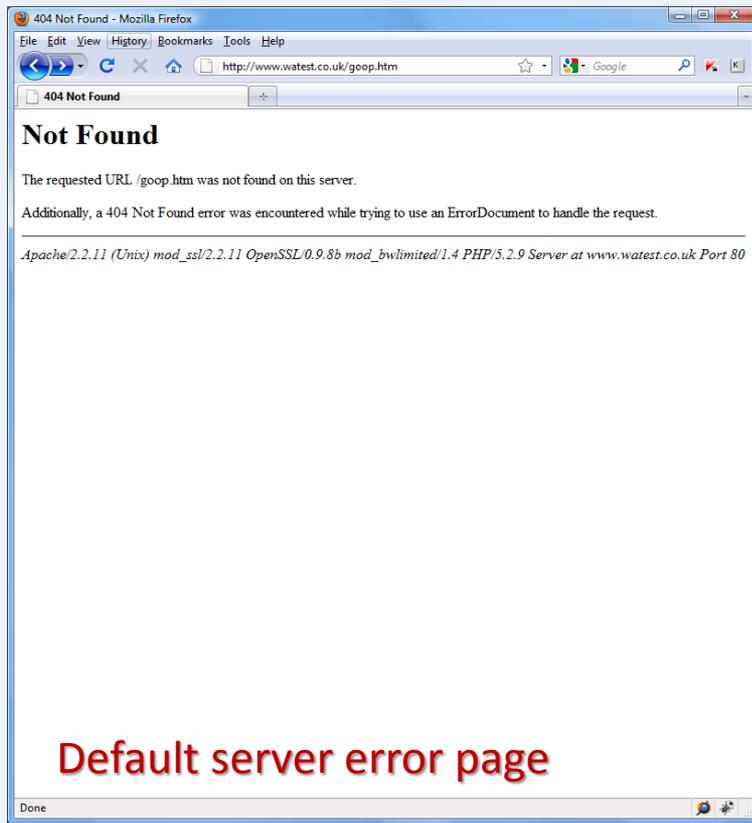
- The file contains one or more lines, known as “configuration directives”.

Website Planning

.htaccess: CUSTOM ERROR PAGES

Custom Error Pages

- All good websites make use of custom error pages; they are an excellent usability tool.



- The most common error is the 404, “page not found”.

Server Errors

- When a hypertext request fails, the server determines the reason and allocates an error code.
- If a requested page cannot be found, the error code is 404.
- However, such codes are meaningless to the normal user and should be avoided.
- Far better to use a useful custom error page to help the user recover from the error.

Creating a custom error page

- Custom error pages are no different to any other web page – they are built using HTML and CSS (and optionally PHP).
- The custom error page should look and feel like part of your site and should include plenty of navigation options – but not too many.
- You tell the server to serve your custom error page by adding a directive to `.htaccess`.

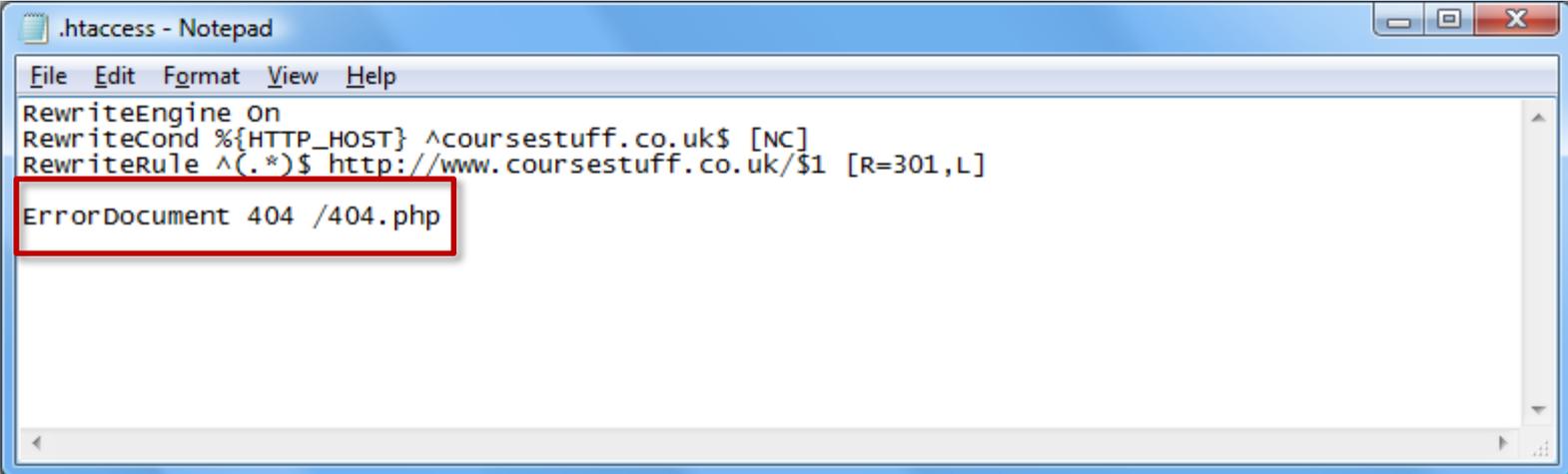
The ErrorDocument directive

```
ErrorDocument 404 /error/404.html
```

- *ErrorDocument* = the directive
- *404* = the error type code
- */error/404.html* = the path from the web root to the page that should be served in the event of this particular error. In this case, a file called *404.html* in a folder called *error*.
- Each of the above elements is separated by a space.

The ErrorDocument directive

- Below is the .htaccess file at coursestuff.co.uk and you can see that in this case, the error file is in the root folder and is a PHP file (*404.php*).



```
.htaccess - Notepad
File Edit Format View Help
RewriteEngine On
RewriteCond %{HTTP_HOST} ^coursestuff.co.uk$ [NC]
RewriteRule ^(.*)$ http://www.coursestuff.co.uk/$1 [R=301,L]
ErrorDocument 404 /404.php
```

Hosting control panel

The screenshot shows a web hosting control panel interface. At the top, there is a red navigation bar with a 'Log Out' button on the right and menu items for 'Home', 'Hosting', 'Domains', 'Internet', 'Telecoms', and 'Help'. Below the navigation bar, the page title is 'Custom Error Pages'. A breadcrumb trail shows the path: 'Home / Account Manager / dwatson / Custom Error Pages'. A paragraph explains that users can choose the nature of the response presented to a customer, replacing the default response generated by Apache. This is ideal for redirecting viewers to a certain page should an application error or Page Not Found message be generated. Users may specify any HTTP error code for the custom error and supply a relative path to a local error document, an external address or a plain text message.

Examples:

- 500 /500.html
- 401 Permission Denied
- 404 http://domain.com/script.php

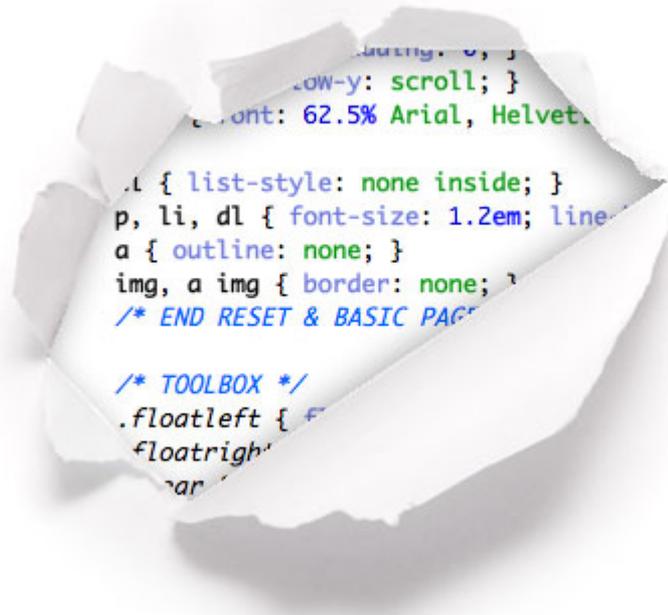
Please note that external redirects are not possible for 401 errors.

Type	Document	Action
<input type="text"/>	<input type="text"/>	<input type="button" value="Create"/>

Some web hosting control panels allow you to set up error directives via a simple form. Pentangle have such a form which automatically creates the .htaccess file for you.

Humour?

- It has become somewhat of a tradition to inject some humour into your custom 404 error page – there are plenty of good examples...



Take a look at the [404 Research Lab](#) or [50 Creative and Inspiring 404 Pages](#) for inspiration

clearleft.com



+44 845 838 6163 @clearleft info@clearleft.com

WHAT WE DO / OUR WORK / WHO WE ARE / OUR BLOG

Page not found

There are known knowns. These are things we know that we know.
There are known unknowns. That is to say, there are things that
we know we don't know. But there are also unknown unknowns.
There are things we don't know we don't know.

We don't know what you were looking for and we don't know we don't know. [Let us know.](#)

Work with us?

Got a great idea or a problem that you need help solving?
We'd love to chat about your business.

Give Andy a call on +44 845 838 6163

Tweet us @clearleft

Email us info@clearleft.com

Or come and visit our office

acromediainc.com



HOME // PORTFOLIO // SERVICES // ABOUT US // CONTACT

Dude, we can't find that page!



We can't find the page you are looking for, we're terribly sorry, but these things happen.

Luckily our site has a lot of other pages that are really awesome, you should check them out:

- [Home](#)
- [Portfolio](#)
- [Services](#)
- [About Us](#)
- [Contact US](#)

Or, you can try searching (little orange search tab, top right), or browse through our [sitemap](#).

We thank you for visiting Acro Media, have a great day.



NEXT STEP → **We'd Love to Discuss Your Online Needs!**
Please provide a few basic details, and we'll get in touch

Name

Email

Phone Optional

Give Us a Shout Now!
Toll Free: **1.877.763.8844**
Local Phone: **1.250.763.8884**
Email: solutions@acromediainc.com

Question? ×
We're here to help

Copyright © 1998-2013 Acro Media Inc. All Rights Reserved. [Privacy Policy](#) | [Sitemap](#) | [Drupal](#)  Offline - Leave a message here!

smashingmagazine.com

SMASHING
MAGAZINE

Books eBooks Job Board Shop

Search

e.g. JavaScript

CODING

CSS
HTML
JavaScript
Techniques

DESIGN

Web Design
Typography
Inspiration
Business

MOBILE

Responsive
iPhone & iPad
Android
Design Patterns

We Couldn't Find Your Page! (404 Error)

By [Smashing Magazine's Server](#)

🕒 Right now 📄 404, Errors 💬 No comments

Unfortunately, the page you've requested **cannot be displayed**. It appears that you've lost your way either through an outdated link or a typo on the page you were trying to reach.

Please feel free to [return to the front page](#) or use the **search box** in the upper area of the page to find the information you were looking for. We are very sorry for any inconvenience.

WHAT CAN YOU DO?

- Use the search box at the top of the page to find the information you were looking for.
- Have a look at our navigation menu, as well as the list of tags in the right sidebar of the page. There, you might find the information you were looking for, as well as other useful articles that you might be interested in.
- You can also have a look at Smashing Magazine's [articles on defensive design](#).
- Report this error using our [contact form](#). We'll be very grateful.

Website Planning

.htaccess: PASSWORD PROTECTION

Password protection

- Password protection requires a `.htaccess` file in the folder to be protected and a `.htpasswd` file located anywhere on the domain (ideally in a secure location).
- In many cases, the `.htpasswd` file is located in the same folder as `.htaccess` but if you have access to folders above the web root, it should be placed there as it is more secure.

How it works...

1. User requests access to folder by entering address in browser.

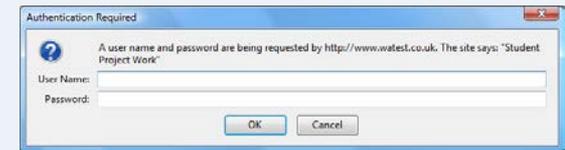


2. Server checks if folder contains .htaccess. If authentication is required...

Authorization Required

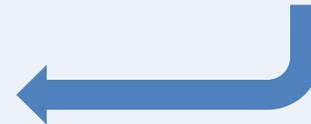
This server could not verify that you are authorized to access the document requested. Either you supplied the wrong credentials (e.g., bad password), or your browser doesn't understand how to supply the credentials required.

Apache/2.2.3 (CentOS) Server at www.davidwatson.info Port 80

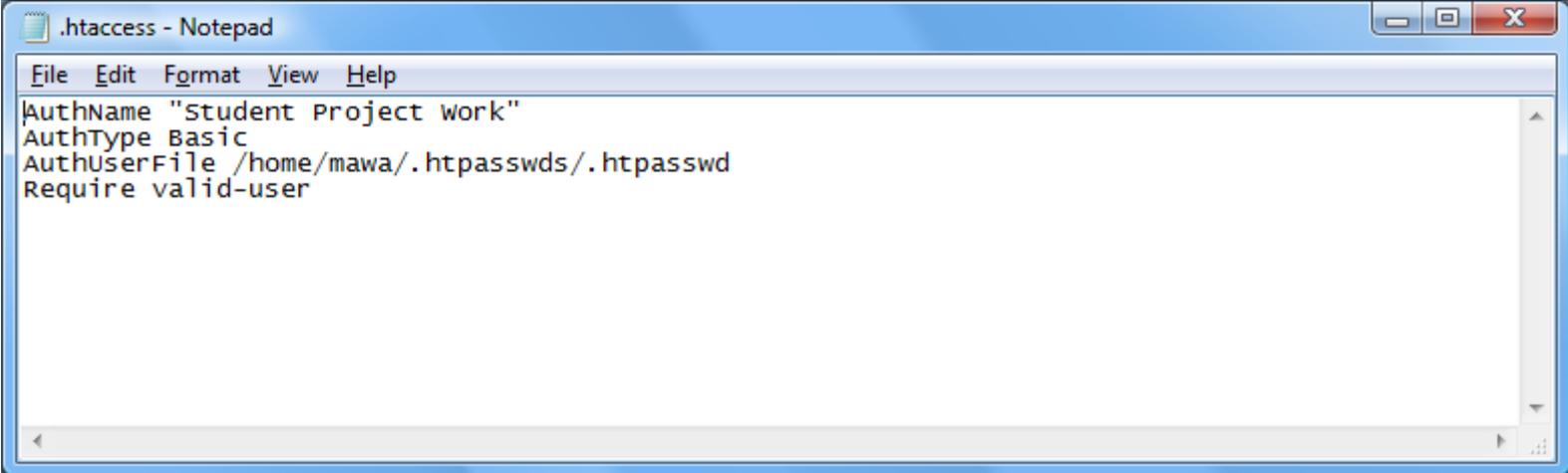


...user is asked to enter User Name and Password.

3. Server checks details against .htpasswd file. If correct, access is granted, if incorrect a 401 error is issued and error page displayed.



Password protection .htaccess

A screenshot of a Notepad window titled ".htaccess - Notepad". The window contains the following text:

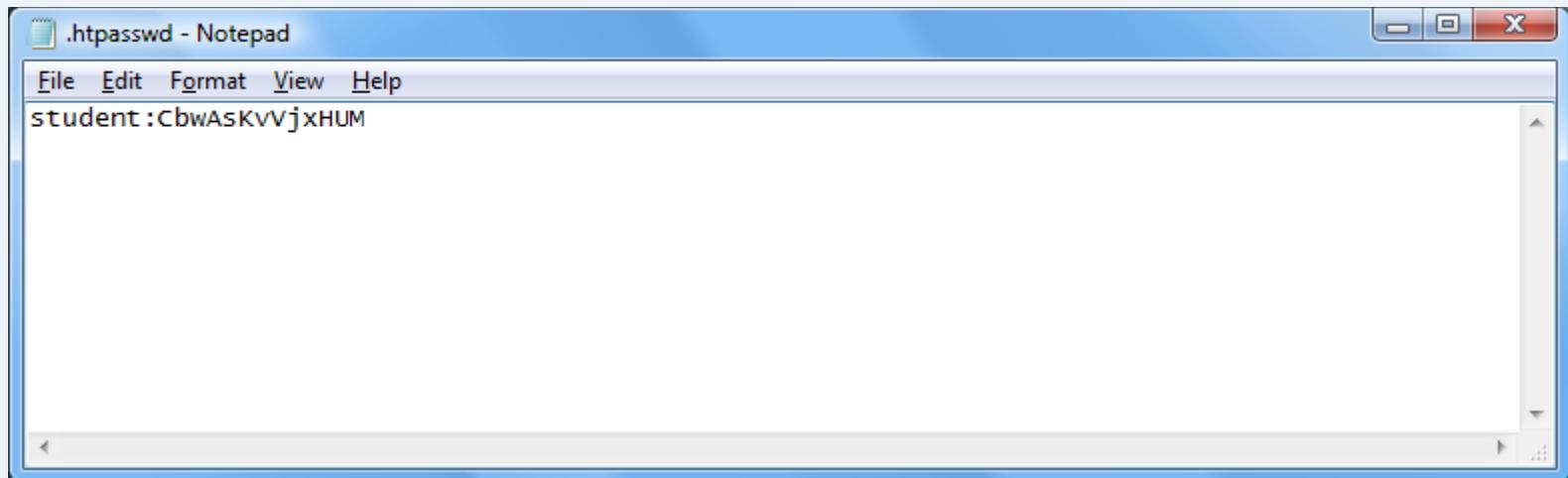
```
File Edit Format View Help
AuthName "Student Project work"
AuthType Basic
AuthUserFile /home/mawa/.htpasswd/.htpasswd
Require valid-user
```

- *AuthName* = text that will display on the authentication dialogue box.
- *AuthType* = method used, *Basic* is the default.
- *AuthUserFile* = server path to the password file.
- *Require* = type of access (e.g. group access can be specified)

Take a look at [Authentication, Authorization and Access Control](#) for more information

Password protection .htpasswd

- The .htpasswd file contains a list of all the valid User Name/Password combinations, one on each line.
- The User Name is plain text but the Password is encrypted using the MD5 algorithm.

A screenshot of a Notepad window titled ".htpasswd - Notepad". The window has a menu bar with "File", "Edit", "Format", "View", and "Help". The text area contains a single line: "student:cbwAsKvVjxHUM". The window has standard Windows window controls (minimize, maximize, close) in the top right corner.

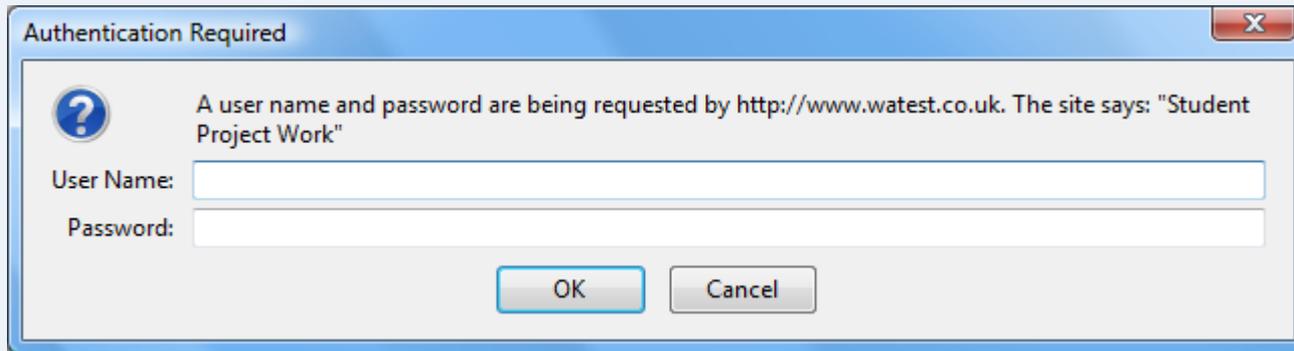
```
.htpasswd - Notepad
File Edit Format View Help
student:cbwAsKvVjxHUM
```

How to make .htpasswd

- There are plenty of free online tools that will automatically create .htpasswd files for you.
- Use Notepad to save your .htpasswd file and then upload to your site using FTP.
- Once both .htaccess and .htpasswd are in place, the folder is protected and accessible only by entering the correct authentication details.

Authentication

- The authentication dialogue box varies depending on browser. FireFox is shown below:



- Notice that "Student Project Work" is the text defined in the *AuthName* directive.

401 Error



- If the authentication is unsuccessful (User Name or Password are incorrect), a 401 error is issued.
- If you wanted, you could make a custom error page for 401 errors.

Hosting control panel

Scripting

	SSL Install and Configure Shared SSL		FrontPage Extensions Install and configure FrontPage extensions
	Password Protection Configure Password protected directories for this web site 		Custom Error Pages Configure Custom Error Pages
	Error Log View this site's error log		One-Click Installs Install popular applications on this site
	Permissions Reset the file permissions throughout the account		php.ini Editor Edit the php.ini configuration & view the phpinfo output

Setting up password protection manually can be a bit of a faff, so most hosting control panels have a tool you can use to do it more easily. Part of the Pentangle control panel is shown above.

Website Planning

.htaccess: REDIRECTION

Websites change

- Websites change: FACT
- In some cases you may want to rename a file or even rename your folders for SEO or for consistency as a site expands.
- So what happens when that popular page has to move or is renamed?
- All the inbound links will be broken, including those from search engines – disaster!

Inbound links

- So, you need to make some major changes to your site...
- ...how can this be done without breaking all the inbound links?
- You can use a 301 redirect to tell search engines where the content has moved to.
- Furthermore, a 301 redirect tells search engines that this is a *permanent* move, so they can update their index accordingly.

The 301 Redirect

- You can use a 301 “permanent” redirect in `.htaccess`.
- This does 2 things:
 - it serves a new page when an old page is requested.
 - it tells search engines to change their index and replace the old page with the new one.

Directive syntax:

`Redirect[space]301[space]old path from root[space]new absolute path`

The example below redirects any request for the folder `/acad` to the new folder `/tutorials/autocad`, for example:

a request for `/acad/index.html` is redirected to `/tutorials/autocad/index.html`

```
Redirect 301 /acad/ http://www.cadtutor.net/tutorials/autocad/
```

Continue redirecting

- Although search engines will learn the new location of content very quickly via your 301 redirect, inbound links are not usually updated in any systematic way, so it's a good idea to keep the redirect in place for as many years as you think appropriate.
- Most webmasters want their content to be correct and a quick email asking them to update their link usually works.

Temporary moves

- It's less common that you may need to move content temporarily...
- ...but if you do, there's a way to do that too.
- Simply use a 302 redirect directive.
- This redirects user requests in the same way as a 301 but it tells search engines not to update their index.

`Redirect 302 /existing/ http://www.temporary.co.uk/mystuff/`

Website Planning

.htaccess: REWRITING URLS

Rewriting URLs

- .htaccess allows you to rewrite any URL and change its form using a Rewrite Engine module in the Apache server, called *mod_rewrite*.
- Common uses:
 - to change <http://www.mydomain.com> to <http://mydomain.com> or vice versa.
 - to change mydomain.co.uk to mydomain.com
 - to change difficult URLs (generated by blogs etc.) to search engine friendly ones.

Canonicalization

- Canonicalization is an SEO issue.
- Search engines may consider <http://www.mysite.com> and <http://mysite.com> to be different websites when, in fact, they are the same.
- The following directive forces all URLs to be rewritten with the “www” even if the request was made without it.

```
RewriteEngine On
```

```
RewriteCond %{HTTP_HOST} ^mysite.com$ [NC]
```

```
RewriteRule ^(.*)$ http://www.mysite.com/$1 [R=301,L]
```

Regular Expressions

RewriteEngine On

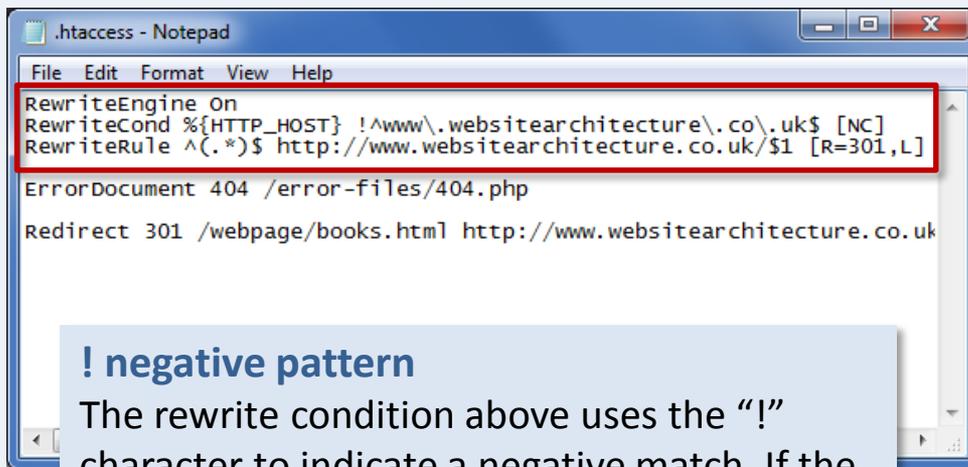
RewriteCond %{HTTP_HOST} ^mysite.com\$ [NC]

RewriteRule ^(.*)\$ http://www.mysite.com/\$1 [R=301,L]

- The directive strings for RewriteCond and RewriteRule look a bit odd.
- They use *regular expressions* (regex) to match URL patterns.
- There's no need to craft your own regex, just use those that others have designed and substitute your own domain details.

Normalising TLDs

- If you have a number of Top Level Domains (e.g. [.com](#), [.net](#), [.co.uk](#)) for the same name, mod_rewrite can be used to change them all to one preferred TLD.



```
.htaccess - Notepad
File Edit Format View Help
RewriteEngine On
RewriteCond %{HTTP_HOST} !^www\.websitearchitecture\.co\.uk$ [NC]
RewriteRule ^(.*)$ http://www.websitearchitecture.co.uk/$1 [R=301,L]
ErrorDocument 404 /error-files/404.php
Redirect 301 /webpage/books.html http://www.websitearchitecture.co.uk
```

! negative pattern

The rewrite condition above uses the “!” character to indicate a negative match. If the requested URL does **not** match this pattern, it will be rewritten so that it matches the pattern defined in the rewrite rule.

On the left is the .htaccess file used at the websitearchitecture website. The directive changes all TLD variations, with or without the “www” to the preferred URL.

For example,

<http://websitearchitecture.net>

will be rewritten as:

<http://www.websitearchitecture.co.uk>

and that’s what will appear in the address bar.

Tidy URL parameters

- URLs with parameters look untidy and may look suspicious to users who don't understand how they work. They may also be bad for SEO.
- The RewriteEngine can be used to tidy such URLs.

RewriteEngine On

RewriteRule `^([0-9]+)\/?$ index.php?id=$1` [NC]

`http://interaction.gallery/dream/index.php?id=25`

becomes

`http://interaction.gallery/dream/25`

Website Planning

.htaccess: PREVENT HOTLINKING

Stop Hotlinking!

- `mod_rewrite` can also be used to prevent people hotlinking (or inline linking) to your content and stealing your bandwidth.
- The directives below (added to `.htaccess`) will cause a “failed request” when `.GIF`, `.JPG`, `.JS` or `.CSS` files are requested from outside the server.

RewriteEngine on

```
RewriteCond %{HTTP_REFERER} !^$
```

```
RewriteCond %{HTTP_REFERER} !^http://(www\.)?mydomain.com/.*$ [NC]
```

```
RewriteRule \.(gif|jpg|js|css)$ - [F]
```

Serving Alternate Content

- `mod_rewrite` can even be used to serve alternate content in response to a hot linking request.
- The directives below serve an image called *angryman.gif* every time a .GIF or .JPG file is requested from outside the server.

RewriteEngine on

RewriteCond %{HTTP_REFERER} !^\$

RewriteCond %{HTTP_REFERER} !^http://(www\.)?mydomain.com/.*\$ [NC]

RewriteRule \.(gif|jpg)\$ http://www.mydomain.com/angryman.gif [R,L]

Website Planning

.htaccess: DENY ACCESS

Deny access by IP address

```
order allow,deny
deny from 123.16.14.245
deny from 41.251.66.32
deny from 105.238.0.
allow from all
```

deny from...

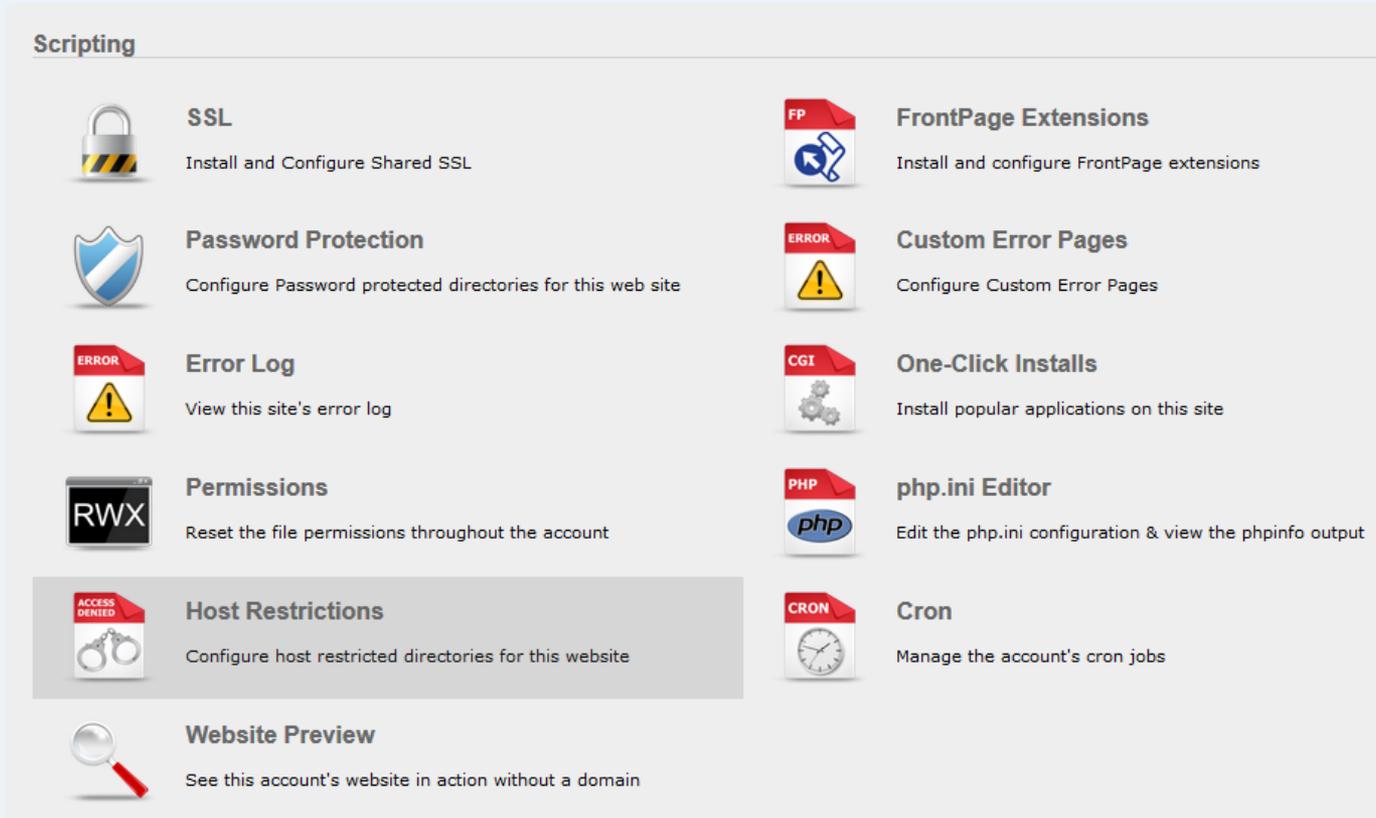
You can deny access from any specific IP address by adding a “deny from” directive and adding the explicit IP address, e.g. **123.16.14.245**. But you can also deny access from an IP range by omitting one or more sets of digits. So, **105.238.0.** means all IP addresses between 105.238.0.0 and 105.238.0.225.

There may be times when you want to prevent access to your website from certain IP addresses. Say you suspect a hacking attempt and you have the user IP address from your server logs or you just want to stop a bandwidth-hogging bot.

Simply, add any IP addresses you want to deny access to in your .htaccess file using the syntax shown above.

This can also be used to deny access to specific folders – just add a .htaccess file to that folder with the appropriate deny/allow directives.

Host restriction from control panel



Scripting

-  **SSL**
Install and Configure Shared SSL
-  **Password Protection**
Configure Password protected directories for this web site
-  **Error Log**
View this site's error log
-  **Permissions**
Reset the file permissions throughout the account
-  **Host Restrictions**
Configure host restricted directories for this website
-  **Website Preview**
See this account's website in action without a domain
-  **FrontPage Extensions**
Install and configure FrontPage extensions
-  **Custom Error Pages**
Configure Custom Error Pages
-  **One-Click Installs**
Install popular applications on this site
-  **php.ini Editor**
Edit the php.ini configuration & view the phpinfo output
-  **Cron**
Manage the account's cron jobs

Just like many of the other .htaccess functions, denying access by IP address (or *host restriction*) can be implemented from your hosting control panel.

.htaccess is your friend

- There's more to .htaccess than we've covered here, there are a number of security functions that can be implemented for example.
- However, you should at least be aware of the functions covered because you will need to use them from time-to-time and although some of the syntax looks like gobbledygook, .htaccess can be a very powerful friend.

.htaccess made easy



[.htaccess made easy](#) the book by Jeff Starr

Website Planning

sitemap.xml

sitemap.xml

```
<?xml version="1.0" encoding="UTF-8"?>
<urlset xmlns="http://www.google.com/schemas/sitemap/0.84">
  <url>
    <loc>http://www.websitearchitecture.co.uk/</loc>
    <changefreq>weekly</changefreq>
    <priority>0.5</priority>
  </url>
  <url>
    <loc>http://www.websitearchitecture.co.uk/programme-details</loc>
    <changefreq>weekly</changefreq>
    <priority>0.5</priority>
  </url>
  <url>
    <loc>http://www.websitearchitecture.co.uk/core-courses</loc>
    <changefreq>weekly</changefreq>
    <priority>0.5</priority>
  </url>
</urlset>
```

As its name suggests, sitemap.xml is an XML file that lists all the important content on your website. It tells Google and other search engine spiders which content you would like them to index. It also includes options that allow you to specify how often the content changes and its relative importance.

Element Definitions

Element	Required?	Description
<urlset>	Yes	The document-level element for the Sitemap. The rest of the document after the '<?xml version>' element must be contained in this.
<url>	Yes	Parent element for each entry. The remaining elements are children of this.
<loc>	Yes	Provides the full URL of the page, including the protocol (e.g. http, https) and a trailing slash, if required by the site's hosting server. This value must be less than 2,048 characters.
<lastmod>	No	The date that the file was last modified, in ISO 8601 format. This can display the full date and time or, if desired, may simply be the date in the format YYYY-MM-DD.
<changefreq>	No	<p>How frequently the page may change:</p> <ul style="list-style-type: none">• always• hourly• daily• weekly• monthly• yearly• never <p>'Always' is used to denote documents that change each time that they are accessed. 'Never' is used to denote archived URLs (i.e. files that will not be changed again).</p> <p>This is used only as a guide for crawlers, and is not used to determine how frequently pages are indexed.</p>
<priority>	No	<p>The priority of that URL relative to other URLs on the site. This allows webmasters to suggest to crawlers which pages are considered more important.</p> <p>The valid range is from 0.0 to 1.0, with 1.0 being the most important. The default value is 0.5.</p> <p>Rating all pages on a site with a high priority does not affect search listings, as it is only used to suggest to the crawlers how important pages in the site are to one another.</p>

The sitemap protocol is recognised by Google, Yahoo! And Microsoft.

Wikipedia: [Sitemaps](#)

Building sitemaps

- You can easily build your own sitemaps if you have a simple site with a few pages. All the information you need is available at sitemaps.org.
- If you have a site with many 100s or 1000s of pages, what should you do then?
- Fortunately, there are a number of free services that will crawl your site and build sitemap.xml for you. For example: XML-Sitemaps.com.
- However, always check that you get what you want. These services do not discriminate and you may want to edit the result before using it.
- Google Webmaster Tools recommends you use sitemap.xml for all your sites – that's a pretty good hint that you should have one!

- Site Dashboard
- Site Messages (2)
- Search Appearance **i**
- Search Traffic
- Google Index
- Crawl
 - Crawl Errors
 - Crawl Stats
 - Fetch as Google
 - Blocked URLs
 - Sitemaps**
 - URL Parameters
 - Security Issues
 - Other Resources
- Labs

Sitemaps

ADD/TEST SITEMAP

By me (1) All (1)

Sitemaps content



Google Webmaster Tools
 Once you have created and uploaded your sitemap.xml file, you should submit it to Google using Webmaster Tools. This ensures that Google knows it exists and how to find it. Once submitted and indexed, you can keep track of its use by Google.

Sitemaps (All content types)

Download All Resubmit Delete Show 25 rows 1-1 of 1

<input type="checkbox"/>	#	Sitemap ▲	Type	Processed	Issues	Items	Submitted	Indexed
<input checked="" type="checkbox"/>	1	/sitemap.xml	Sitemap	Jan 30, 2014	-	Web	32	1

1-1 of 1

Website Planning

robots.txt

robots.txt

```
User-agent: *
Disallow: /error/
Disallow: /includes/
Disallow: /forum/clientscript/
Disallow: /forum/cpstyles/
Disallow: /forum/customavatars/
Disallow: /forum/customgroupicons/
Disallow: /forum/customprofilepics/
Disallow: /forum/images/
Disallow: /forum/includes/
Disallow: /forum/install/
Disallow: /forum/signaturepics/

Sitemap: http://www.websitearchitecture.co.uk/sitemap.xml
```

The purpose of robots.txt is to tell crawlers/spiders where they should not go. In other words, it lists any content that you **do not** want indexed. By default, spiders will index any content they find.

In the example above, robots.txt is also used to alert spiders to the fact that sitemap.xml is available. Essentially, that file tells spiders what you **do** want them to index.

Building robots.txt

- As its name suggests, robots.txt is just a simple text file and you can easily write your own following the protocol at robotstxt.org.
- All spiders request robots.txt when they first access a website. If the file is not found, a 404 error is issued and the spider continues with crawling your site.
- Even if you have no content to hide, having a robots.txt file avoids the 404 error and the serving of your custom error page, if you have one.

Empty robots.txt file

```
=====  
User-agent: *  
Disallow:  
  
=====
```

It's probably a good idea to include a robots.txt file in your web root in order to avoid 404 errors. Something like the text above is all you need (note the 2 blank lines after "Disallow:"). Don't forget to add your sitemap when you have one in place.

Note: this is not a substitute for password protection because not all spiders play by the rules!

Webmaster Central: [Do I need a robots.txt file?](#)



- Site Dashboard
- Site Messages (2)
- ▶ Search Appearance ⓘ
- ▶ Search Traffic
- ▶ Google Index
- ▼ Crawl
 - Crawl Errors
 - Crawl Stats
 - Fetch as Google
 - Blocked URLs**
 - Sitemaps
 - URL Parameters
- Security Issues
- Other Resources
- ▶ Labs

Blocked URLs

If your site has content you don't want Google or other search engines to access, use a robots.txt file to specify how search engines should crawl your site's content.

Check to see that your robots.txt is working as expected. (Any changes you make to the robots.txt content below will not be saved.)

robots.txt file	Blocked URLs [?]	Downloaded	Status
http://www.websitearchitecture.co.uk/robots.txt	322	8 hours ago	200 (Success)

robots.txt analysis

Value	Result
Line 56: Sitemap: http://www.websitearchitecture.co.uk/sitemap.xml	Valid Sitemap reference detected

<http://www.websitearchitecture.co.uk/robots.txt> content - edit to test changes

```
User-agent: *
Disallow: /error/
Disallow: /includes/
Disallow: /forum/clientscript/
Disallow: /forum/cpstyles/
Disallow: /forum/customavatars/
Disallow: /forum/customgroupicons/
Disallow: /forum/customprofilepics/
Disallow: /forum/images/
Disallow: /forum/includes/
Disallow: /forum/install/
Disallow: /forum/signaturepics/
Disallow: /forum/ajax.php
Disallow: /forum/announcement.php
Disallow: /forum/attachment.php
Disallow: /forum/calendar.php
Disallow: /forum/converse.php
Disallow: /forum/cron.php
Disallow: /forum/editpost.php
Disallow: /forum/faq.php
Disallow: /forum/global.php
```

URLs Specify the URLs and user-agents to test against.

<http://www.websitearchitecture.co.uk/>

Google Webmaster Tools

You can check the effectiveness of robots.txt and to see whether it is being correctly interpreted using Google Webmaster Tools. You can also see the last time robots.txt was downloaded (by Google) and whether the request was completed successfully.

Website Planning

humans.txt

the idea

standard

H-team

Friends

submit!

Humans!

About humans.txt

What is humans.txt?



It's an initiative for knowing the people behind a website. It's a TXT file that contains information about the different people who have contributed to building the website.

Why a TXT?



Because it's something *simple and fast* to create. Because it's **not intrusive with the code**. More often than not, the owners of the site don't like the authors signing it; they claim that doing so may make the site less efficient. By adding a txt file, you can prove your authorship (not your property) in an external, fast, easy and accessible way.



Where is it located?

In the site root. Just next to the robots.txt file.

If possible, you can also add an author tag to the <head> of the site:

```
<link type="text/plain"
rel="author" href="http://domain
/humans.txt" />
```



Why should I?

You don't have to if you don't want. The only aim of this initiative is to know **who the authors of the sites** we visit are.



Who should I mention

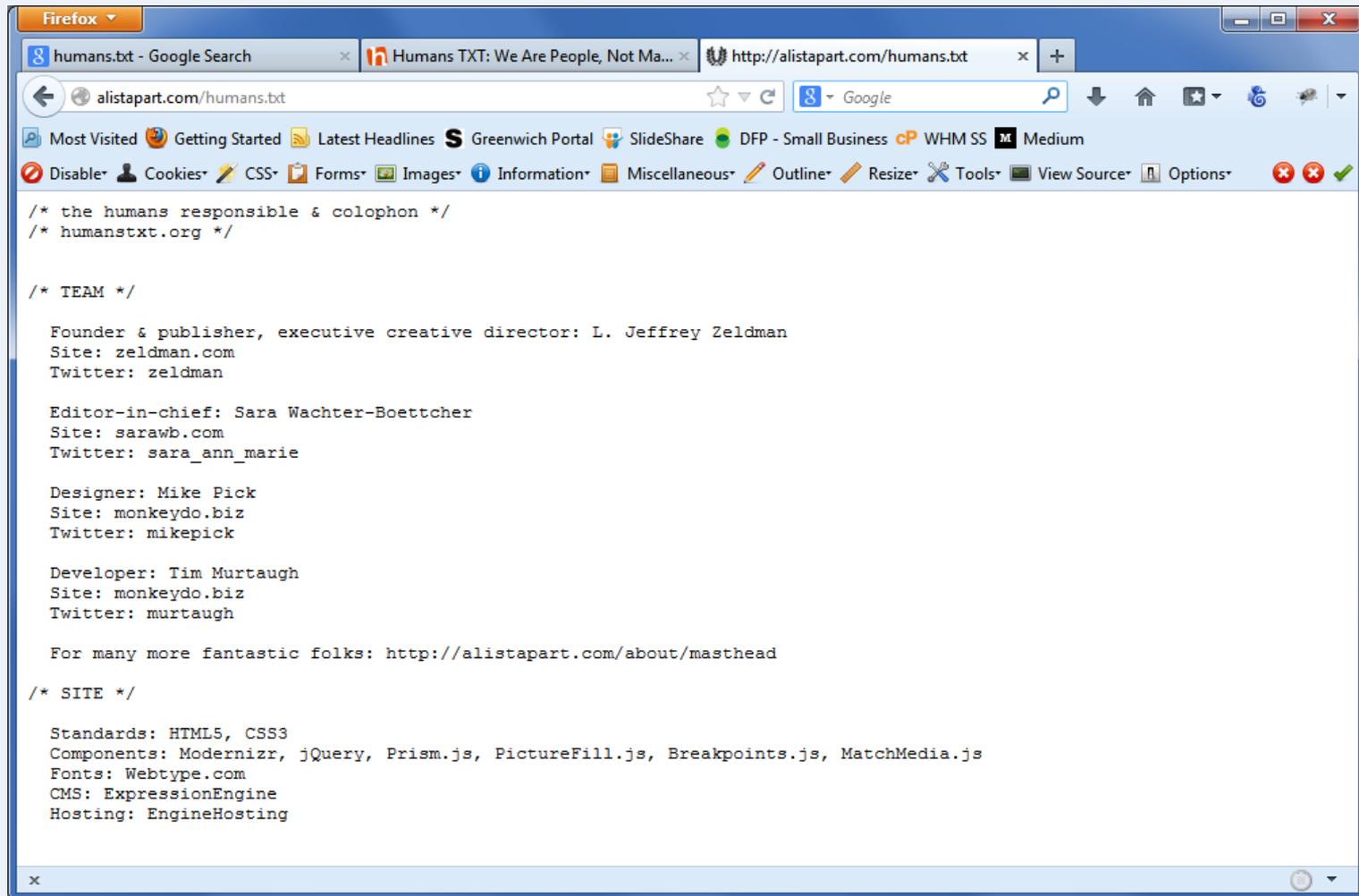
Whoever you want to, provided they wish you to do so. You can mention the developer, the designer, the copywriter, the webmaster, the SEO, SEM or SMO...

As you can see, the number of people who may take part of the creation of a site can be big, so the list is almost endless.

humans.txt

Optionally, you may add a humans.txt file to the root folder of your website. This file is for humans to read (hence the name) and should contain information about the authors of the website and details of the technologies and methods used in its construction as well as any other relevant information.

Unlike robots.txt, this file has no practical function and is not commonly used but it does demonstrate good attention to detail and it's a nice way to give credit to those involved in a design project.



The screenshot shows a Firefox browser window with the address bar displaying `alistapart.com/humans.txt`. The page content is a plain text file with the following text:

```
/* the humans responsible & colophon */
/* humanstxt.org */

/* TEAM */

Founder & publisher, executive creative director: L. Jeffrey Zeldman
Site: zeldman.com
Twitter: zeldman

Editor-in-chief: Sara Wachter-Boettcher
Site: sarawb.com
Twitter: sara_ann_marie

Designer: Mike Pick
Site: monkeydo.biz
Twitter: mikepick

Developer: Tim Murtaugh
Site: monkeydo.biz
Twitter: murtaugh

For many more fantastic folks: http://alistapart.com/about/masthead

/* SITE */

Standards: HTML5, CSS3
Components: Modernizr, jQuery, Prism.js, PictureFill.js, Breakpoints.js, MatchMedia.js
Fonts: Webtype.com
CMS: ExpressionEngine
Hosting: EngineHosting
```

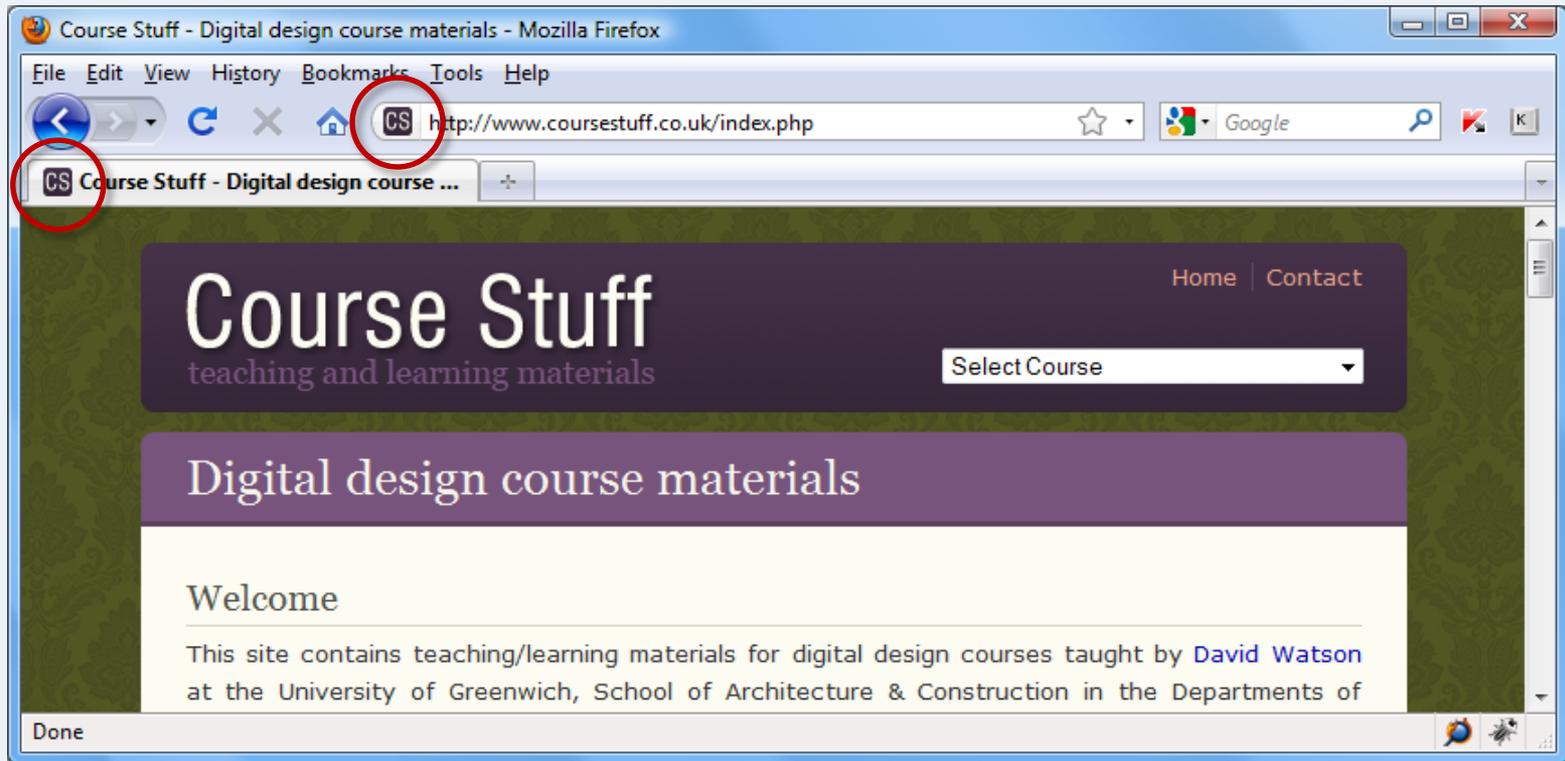
alistapart.com/humans.txt is a good example of a typical humans.txt file it contains brief details of those involved and the technologies used.

Website Planning

favicon.ico

What is a Favicon?

- A Favicon is a small graphic image that appears in the address bar and in other places when a website is viewed in a browser.

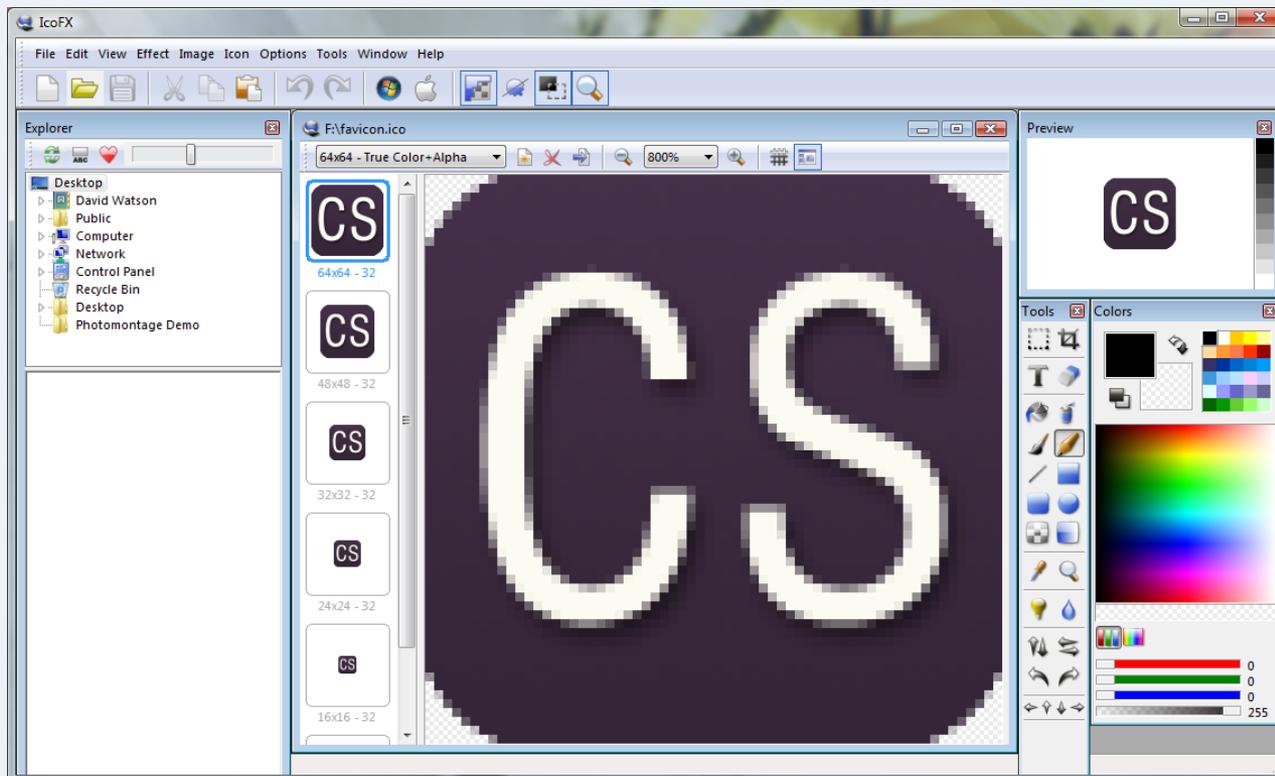


How do I create a Favicon?

- A Favicon is a special type of image file (.ico) that is not commonly supported by mainstream applications – Photoshop has no native support, Fireworks CS4 and above does.
- Fortunately, there are plenty of free and low-cost options for creating favicons.
- Plugins are available for Photoshop and Fireworks.
- There are many online image converters and editors like [x-icon editor](#).
- There are some great free icon editors like [Icon Editor](#) and [Icon Editor Pro](#) (a portable app.)

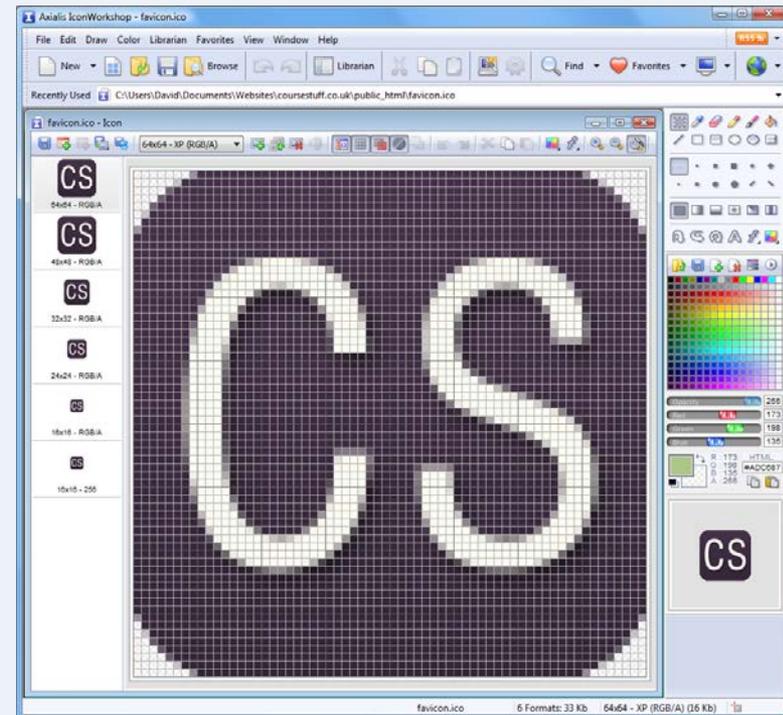
Can't I just use a PNG?

- Most browsers support GIF, JPG and PNG file formats for Favicons.
- Internet Explorer 10 and below support only ICO files.



Axialis IconWorkshop

- If you create a lot of icons, it may be worth spending a bit of money (\$49) on an application like [IconWorkshop](#) or [IcoFX](#).
- This includes a Photoshop plugin that allows you to design the graphic in Photoshop and then export to IconWorkshop for completion.



Adding the Favicon to your site

- When you save your icon, it should be called **favicon.ico**, this is the default filename the server will look for, just as it looks for index.html as a default homepage.
- Use FTP to upload favicon.ico to the *root* folder of your website.
- There is no need to add a link tag to the `<head>` of your HTML files if you use the default filename and place it in the root folder.

When do I need a link tag?

- You only need to point to a Favicon using a `<link>` tag if:
 - Your icon file is called something other than `favicon.ico` or is in a sub-folder.
 - You want to use different icons for different parts of your site.
 - You want to conform to W3C preferences!

```
<link rel="icon" href="/folder/favicon.ico" />
```

All change!

- With the advent of HTML5, favicon.ico is effectively deprecated (we shouldn't really use it) but it still works perfectly well.
- There are also a wider range of contexts where icons are used – desktop, tablet, phone...
- In principle, we should use the ,PNG format, create one file for each image size and link to them from the <head>.
- See this [useful article](#) at CSS Tricks for details.

Redirect 301 start end