

# Filtering in Confluence Access Logging

If we want to measure and log "real" pageviews from our Confluence, see [Access Logging in Confluence](#), and we are not doing this from within Confluence, we need to filter out false results.

This is my findings regarding to filtering results:

## Robots and Crawlers

We need to filter out those:

### Kibana / Grafana

```
type: apache AND (request.keyword: \display* OR request.keyword: \pages\viewpage.action*) AND NOT ( agent: *facebook* OR agent: *crawler* OR agent: *bot* OR agent: *Bot* OR agent: *Spider*)
```

### Splunk

```
index=apache AND host=moserver AND (uri="*/display/*" OR uri="*/viewpage.action/*") useragent!="*bot*" useragent!="*spider*" useragent!="*facebookexternalhit*" useragent!="*crawler"
```

## HTTP Codes

Often, if your page is down or under maintenance, apache/nginx logs will return **HTTP Code 502** or **HTTP Code 503** - See <https://www.w3.org/Protocols/rfc2616/rfc2616-sec10.html>

These can be removed from the search results; - or **HTTP Code 200** can be added to give only positive results:

### Kibana / Grafana

```
type: apache AND (request.keyword: \display* OR request.keyword: \pages\viewpage.action*) AND response: 200 AND NOT ( agent: *facebook* OR agent: *crawler* OR agent: *bot* OR agent: *Bot* OR agent: *Spider*)
```

### Splunk

```
index=apache AND host=moserver AND (uri="*/display/*" OR uri="*/viewpage.action/*") useragent!="*bot*" useragent!="*spider*" useragent!="*facebookexternalhit*" useragent!="*crawler*" AND status=200
```

## Reffers

A User request for a Confluence page spawns several request from the Confluence to the Confluence itself, and from the client to the Confluence server, but only one represents the actual page request. In the above, this is covered by:

### Kibana / Grafana

```
request.keyword: \display* OR request.keyword: \pages\viewpage.action*
```

### Splunk

```
uri="*/display/*" OR uri="*/viewpage.action/*"
```

# Monitoring

Monitoring a website from various Sources (Zenoss, Nagios, Datadog etc etc) also represents a large (often the largest) number of hits, and we need to remove these, this can be done in several ways:

- Identifying all agent types of the monitoring
- Often monitoring is on / - giving the agent a **HTTP Code 302 back** code. This is not the case for mysite, where apache makes a rewrite; hence the monitoring gets a **HTTP Code 200 OK**

Looking into Your logs will often make it pretty easy to Identify monitoring:

Here is my datadog monitoring from another server:

i	Time	Event
>	12/03/2017 11:25:40.000	afserver.npn.hosted.netic.dk - - [12/Mar/2017:11:25:40 +0100] "GET / HTTP/1.1" 200 19030 "-" "Datadog Agent/5.10.1" host = moserver   source = /var/log/apache2/www.mos-eisley.dk-access.log   sourcetype = access_combined
>	12/03/2017 11:25:20.000	afserver.npn.hosted.netic.dk - - [12/Mar/2017:11:25:20 +0100] "GET / HTTP/1.1" 200 19029 "-" "Datadog Agent/5.10.1" host = moserver   source = /var/log/apache2/www.mos-eisley.dk-access.log   sourcetype = access_combined
>	12/03/2017 11:25:01.000	afserver.npn.hosted.netic.dk - - [12/Mar/2017:11:25:01 +0100] "GET / HTTP/1.1" 200 19008 "-" "Datadog Agent/5.10.1" host = moserver   source = /var/log/apache2/www.mos-eisley.dk-access.log   sourcetype = access_combined

And from Worldping via [Grafana](#):

>	12/03/2017 11:29:59.000	nyc2.us.collector.raintank.io - - [12/Mar/2017:11:29:59 +0100] "GET / HTTP/1.1" 200 18989 "-" "worldping-api" host = moserver   source = /var/log/apache2/www.mos-eisley.dk-access.log   sourcetype = access_combined
>	12/03/2017 11:29:59.000	ams1.nl.collector.raintank.io - - [12/Mar/2017:11:29:59 +0100] "GET / HTTP/1.1" 200 18993 "-" "worldping-api" host = moserver   source = /var/log/apache2/www.mos-eisley.dk-access.log   sourcetype = access_combined