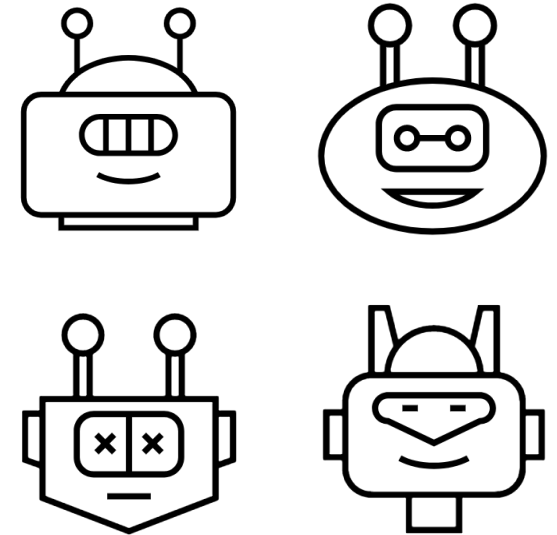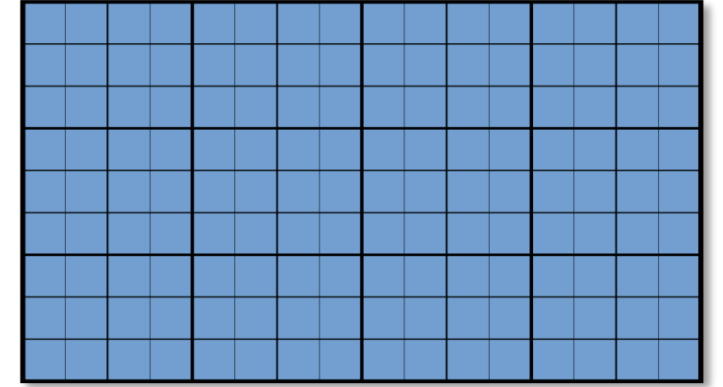# candour

# Advanced robots.txt

Trainer:

Mark Williams-Cook
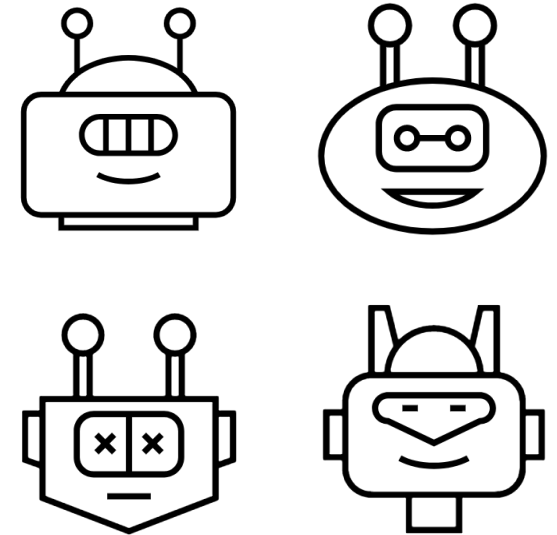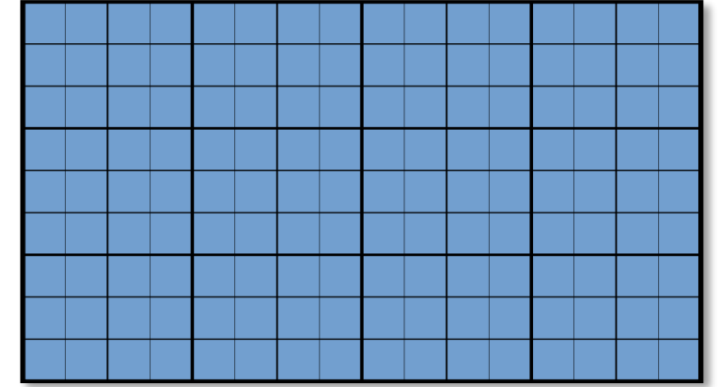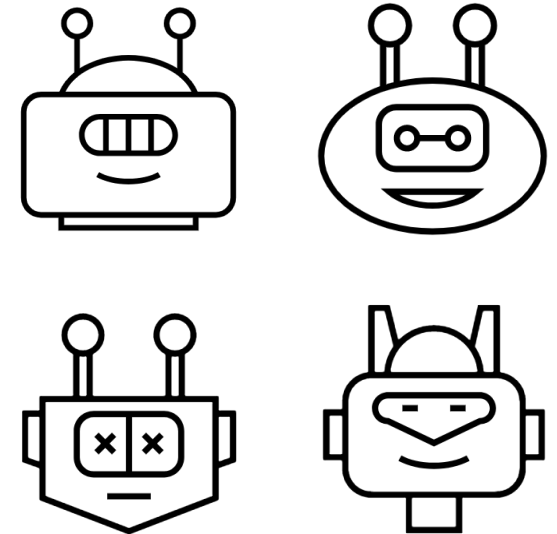Candour

# Filtered navigation
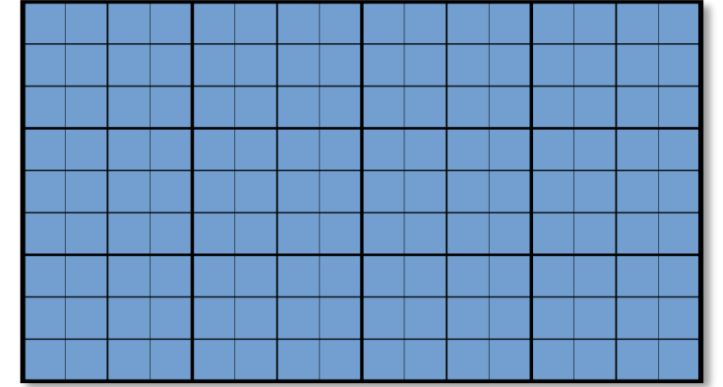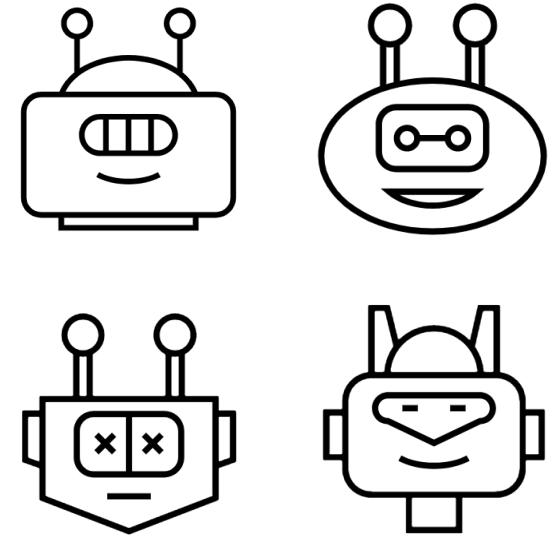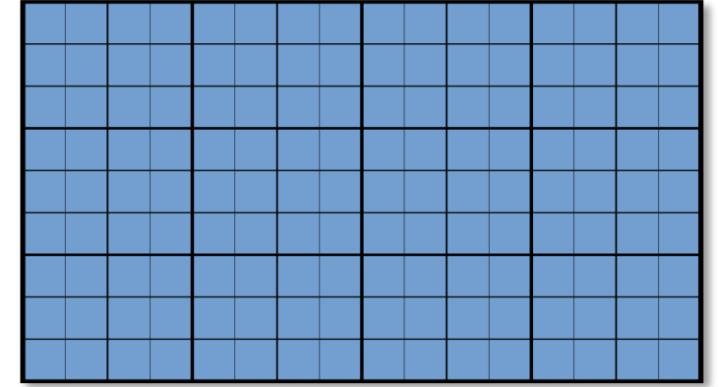
# Filtered navigation

site.com/cat/womens-shoes

# Filtered navigation

site.com/cat/womens-shoes?<span style="color:#e6005c">size=10</span>
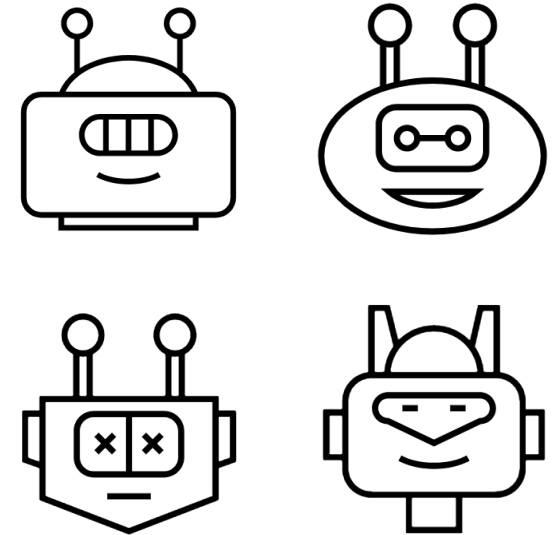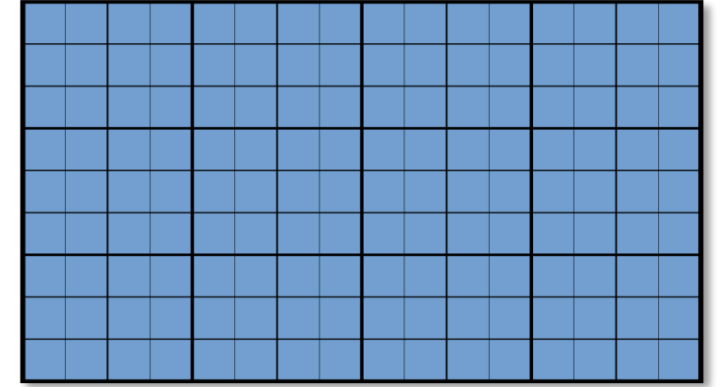
**candour**

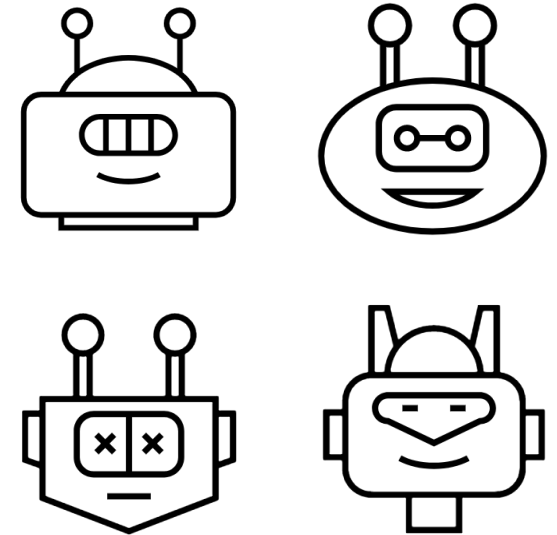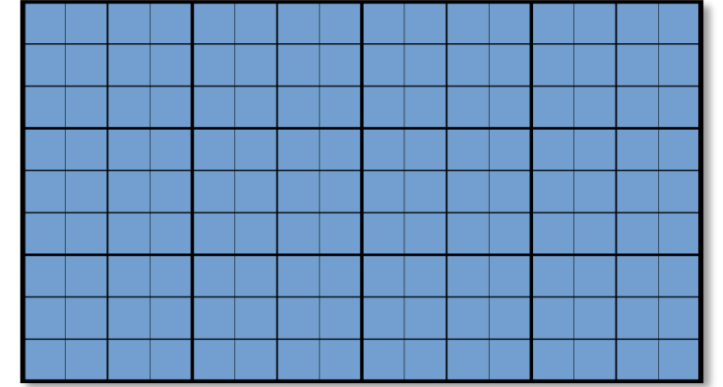# Filtered navigation

site.com/cat/womens-shoes?size=10&colour=red

# Filtered navigation

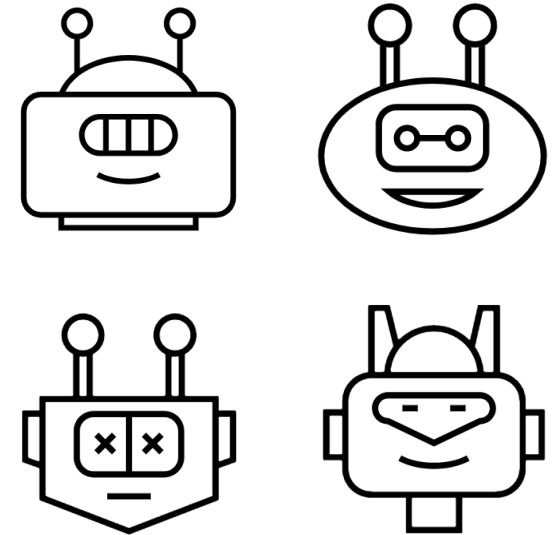site.com/cat/womens-shoes?size=10
&colour=red&price=30

# Filtered navigation
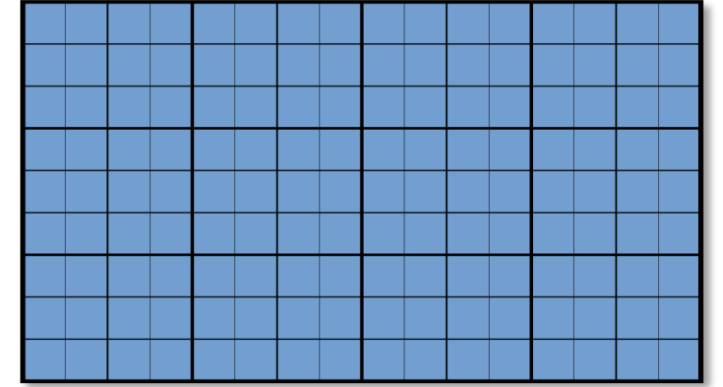
site.com/cat/womens-shoes?size=10
&colour=red&price=30&order=desc

# Filtered navigation

site.com/cat/womens-shoes?size=10
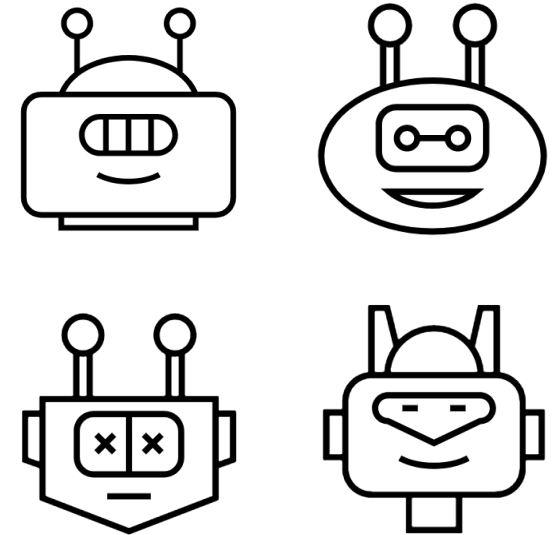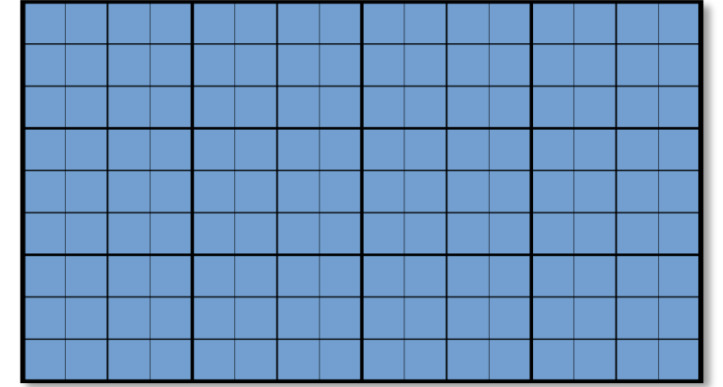&colour=red&price=30&order=desc

12 sizes & 10 colours & 6 price brackets

# Filtered navigation

site.com/cat/womens-shoes?size=10&colour=red&price=30&order=desc

12 sizes & 10 colours & 6 price brackets

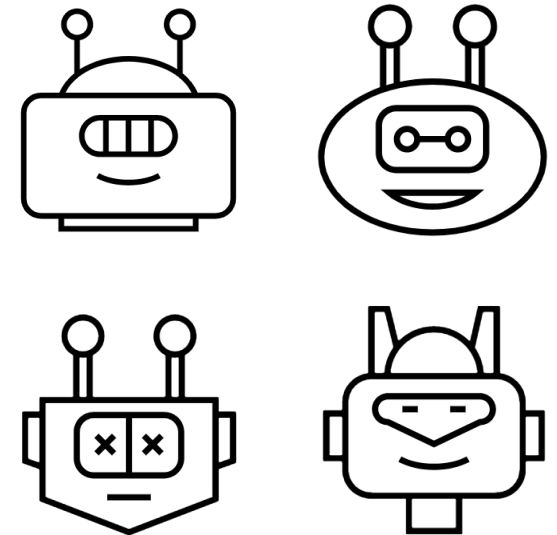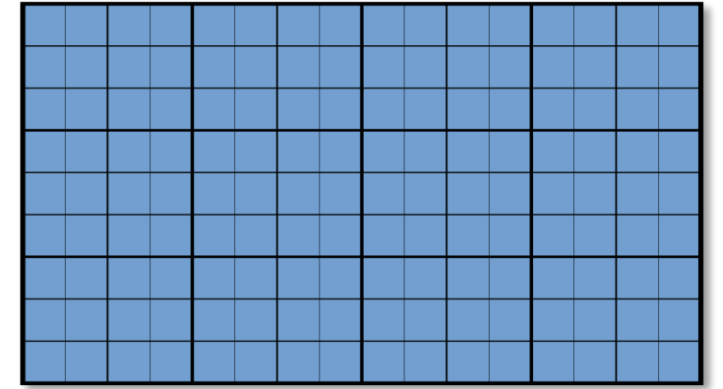This is 10,000s of URL variations!

# Filtered navigation

site.com/cat/womens-shoes?size=10

site.com/cat/womens-shoes?colour=red

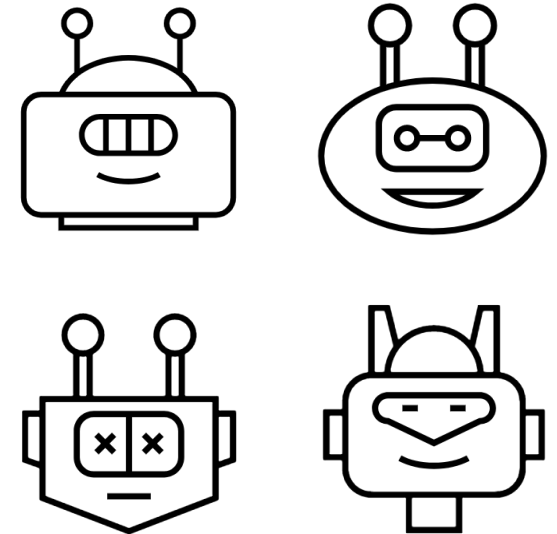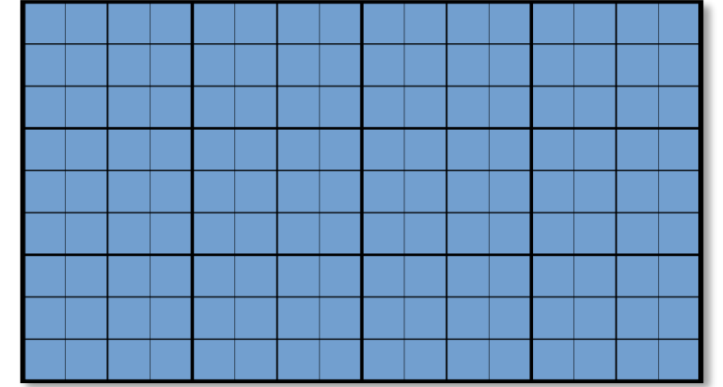site.com/cat/womens-shoes?price=30

# Filtered navigation

site.com/cat/womens-shoes?size=10

site.com/cat/womens-shoes?colour=red

site.com/cat/womens-shoes?price=30

12 sizes & 10 colours & 6 price brackets

= 12 + 10 + 6 = 28 URLs

# Filtered navigation
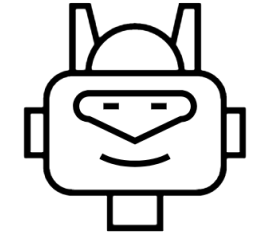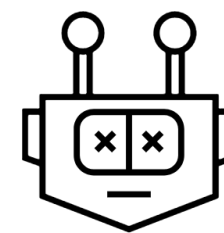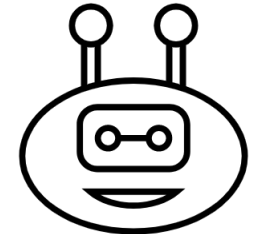
We want to allow: any <u>single</u> query string

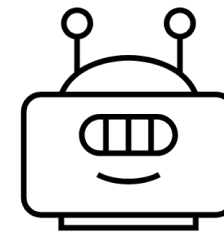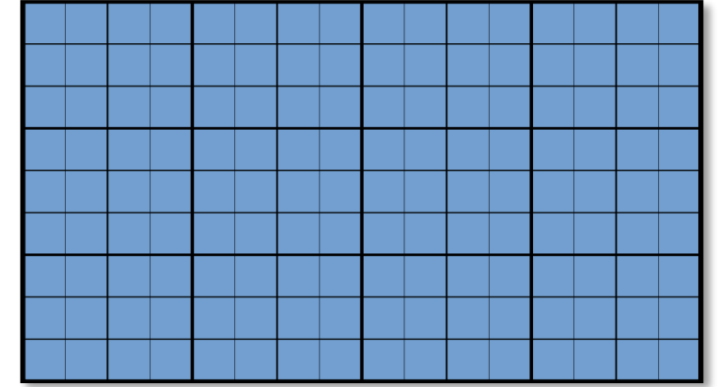site.com/cat/womens-shoes?<span style="color:#e91e63">size=10</span>

site.com/cat/womens-shoes?<span style="color:#e91e63">colour=red</span>

site.com/cat/womens-shoes?<span style="color:#e91e63">price=30</span>

We want to disallow any <u>multiple</u> query string:

site.com/cat/womens-shoes?<span style="color:#e91e63">price=30&size=10</span>
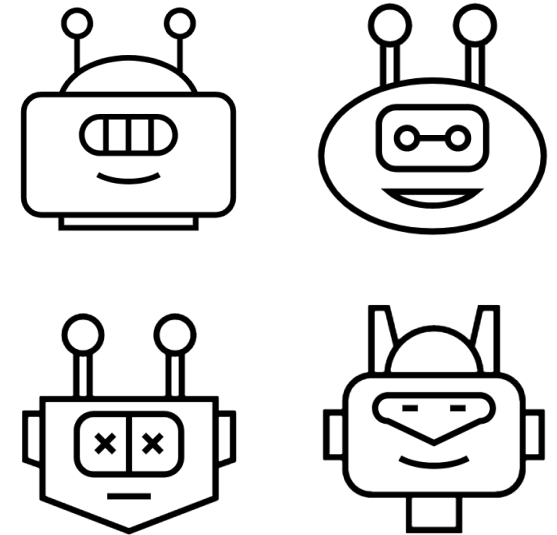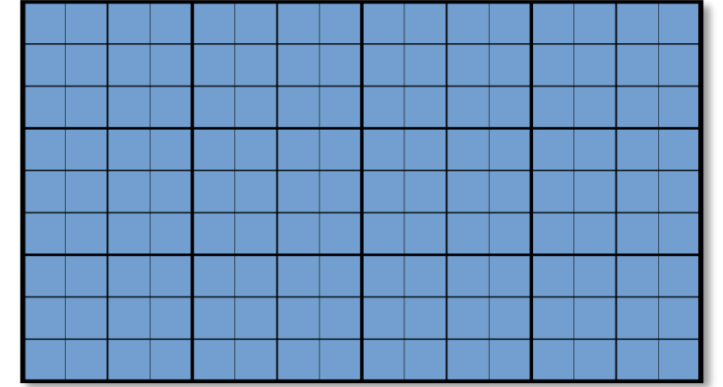
# Filtered navigation

Use wildcards (*) within your robots.txt rules

Wildcards are 0 or more of any valid character

E.g.
Disallow: /cat/*/

Advanced robots.txt
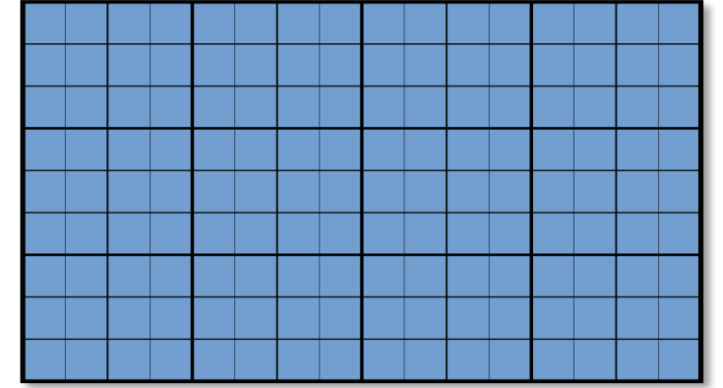
withcandour.co.uk

candour

# Challenge

Write out a robots.txt file that will allow any category URL with a single query string to be crawled but will block crawling of any multiple query string URLs.

candour

# Key concept

Robots.txt is one way to fix such issues – don't rush off to implement anything yet!