

Notes at <https://github.com/hadley/web-scraping>

# Scraping websites with R

Using rvest and the tidyverse

Hadley Wickham

Chief Scientist, Posit

July 2024



# Introductions

# Getting data off a website

Easy

~~1. Official API~~

On Wednesday

2. Unofficial API

3. Static HTML

4. Dynamic HTML

Hard

**1. HTML structure**

**2. rvest basics**

**3. Extracting data**

**4. Bonus content**

Pagination, Live HTML, Unofficial APIs, LLMs

# HTML structure

---

# HTML is a tree

```
<!doctype html>
<html lang="en-US">
  <head>
    <title>Page title</title>
    <meta ... >
    <script ...></script>
    ...
  </head>
  <body>
    ...
  </body>
</html>
```

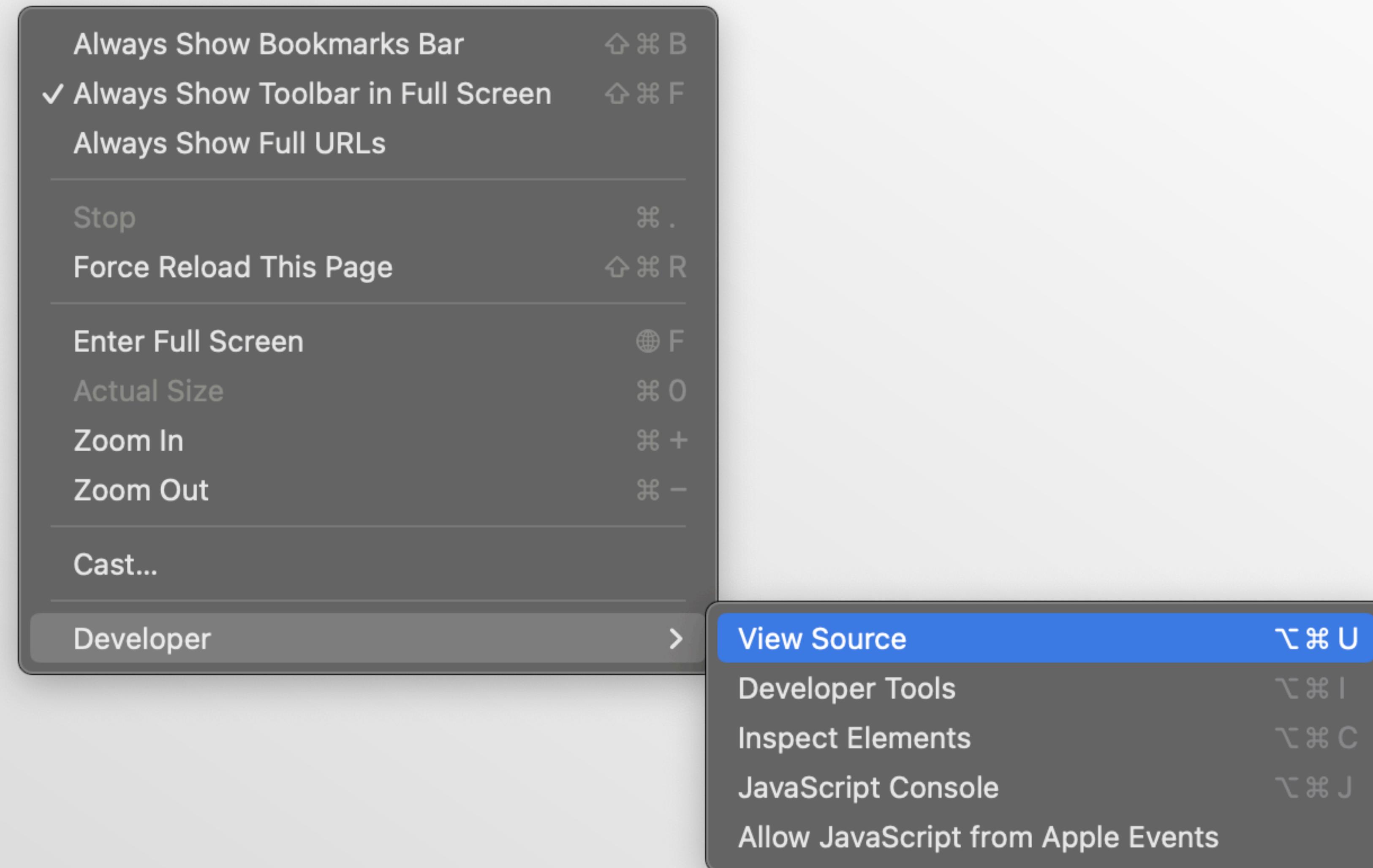
# The tree is made up of elements



# The stuff you see on a page comes from the body

```
<body>
  <h1>Top level heading</h1>
  <p>A paragraph containing text that is <b>bold</b> and
     an image: </p>
  <ul>
    <li>A bulleted list</li>
    <li>Bullet two</li>
  </ul>
  <table>
    <tr><th>A</th><th>B</th></tr>
    <tr><td>1</td><td>2</td></tr>
  </table>
</body>
```

# Best way to see the tree of a real page is to use DevTools



I recommend Chrome for this, even if it isn't your daily driver

# Or right-click and choose inspect

Hoping to resolve the matter with a blockade of deadly battleships, the greedy Trade Federation

Look Up "resolve"

Copy

Copy Link to Highlight

Search Google for "resolve"

Print...

Translate Selection to English

Open in Reading Mode NEW

Attack of the Clones

Released: 2002

Director: George Lucas

debates this alarming chain of events, heeded two Jedi Knights, the guardians of conflict....

Inspect

Speech >

Services >

Elements | Console | Sources | Network | Performance | Memory | Application | Security

```
<!DOCTYPE html>
<!-- Generated by pkgdown: do not edit by hand -->
<html lang="en">
  ><head> ...
  ><body> flex
    <div id="MathJax_Message" style="display: none;"></div>
    <a href="#container" class="visually-hidden-focusable">Skip to content</a>
    <nav class="navbar fixed-top navbar-light navbar-expand-lg bg-none headroom headroom--top
headroom--not-bottom"> flex
      ><div class="container"> ...
      ></div> flex
    </nav>
    ><div class="container template-article" id="container"> ...
    ></div> flex
    ><footer> ...
    ></footer> flex
  </body>
...</html> == $0
```

# Your turn

Go to <<https://rvest.tidyverse.org/articles/starwars.html>>, open the developer tools, and use the "elements" view to answer the following questions:

What element contains the film titles?

What element contains the name of the director? What attributes does this element have?

How many paragraphs does the “crawl” at the start of each movie have?

Where does the table of contents live relative to the film data?

# Static vs dynamic

The HTML displayed in the elements pane is usually generated from a HTML file that you can find in the sources pane.

Sometimes, however, the HTML is generated **dynamically** with Javascript. We'll come back to this later, but I wanted to illustrate the difference.

Let's look at <<https://rvest.tidyverse.org/dev/articles/starwars-dynamic.html>>

# rvest basics

---

# Easiest to get data from HTML if it's already in a table

```
<table>
  <tr>
    <th>Character</th> Header
    <th>Role</th>
    <th>Affiliation</th>
  </tr>
  Row <tr>
    <td>Luke Skywalker</td> Cell  
(d = data)
    <td>Jedi Knight</td>
    <td>Rebel Alliance</td>
  </tr>
  <tr>
    <td>Darth Vader</td>
    <td>Sith Lord</td>
    <td>Galactic Empire</td>
  </tr>
</table>
```

# Can read with html\_table()

```
library(rvest)
html <- minimal_html("<table>
<tr>
  <th>Character</th>
  <th>Role</th>
  <th>Affiliation</th>
</tr>
<tr>
  <td>Luke Skywalker</td>
  <td>Jedi Knight</td>
  <td>Rebel Alliance</td>
</tr>
</table>")
html_table(html)
```

# Let's look at a more realistic example

```
# I want to get the sound track for Star Wars movie  
  
url <- "https://en.wikipedia.org/wiki/Star_Wars_(soundtrack)"  
html <- read_html(url)  
  
html >  
  html_table() >  
  _[5:8]
```

# Can we do better than asking for tables 5 through 8?

```
# I want to get the sound track for Star Wars movie  
  
url <- "https://en.wikipedia.org/wiki/Star_Wars_(soundtrack)"  
html <- read_html(url)  
  
html >  
  html_table() >  
  _[5:8]
```

# Your turn

Open <[https://en.wikipedia.org/wiki/Star Wars \(soundtrack\)](https://en.wikipedia.org/wiki/Star_Wars_(soundtrack))>

Using the developer tools, can you find something in the structure of the HTML that uniquely identifies these tables?

# We can use html\_elements() with CSS selectors

```
url <- "https://en.wikipedia.org/wiki/Star_Wars_(soundtrack)"
```

```
html <- read_html(url)
```

```
html >
```

```
  html_elements(".tracklist") >  
  html_table()
```

# .tracklist means all elements with class tracklist

# CSS selectors

**CSS** = cascading style sheets

Primary purpose is to separate the visual appearance (style) from its underlying semantics.

Used to say (e.g.) “make this box blue” or “make all links green”.

It's a domain specific language for selecting elements in the HTML tree. We'll use it to identify the HTML elements that contain the data we care about.

# Most important selectors

- `.brown` = all elements with class "brown"
- `#abc` = single element with id "abc"
- `p` = all paragraphs
- `p.important` = all paragraphs with "important" class
- `p b` = all bold elements that are descendants of a paragraph
- `p > b` = all bold elements that are children of a paragraph

# Your turn

Use <https://flukeout.github.io/> to learn and practice the most important selectors.

# Extracting data

---

1. Find the “rows” with `html_elements()`
2. Find the “columns” with `html_element()`
3. Extract the data with `html_text2()` or `html_attr()`
4. Make a tibble
5. Clean it up

# Your turn

- Head back to [https://rvest.tidyverse.org/articles/  
starwars.html](https://rvest.tidyverse.org/articles/starwars.html)
- What are the rows? How can you identify them with a css selector?
- What are the columns? How can you identify them with a css selector?

# Solution

starwars.R

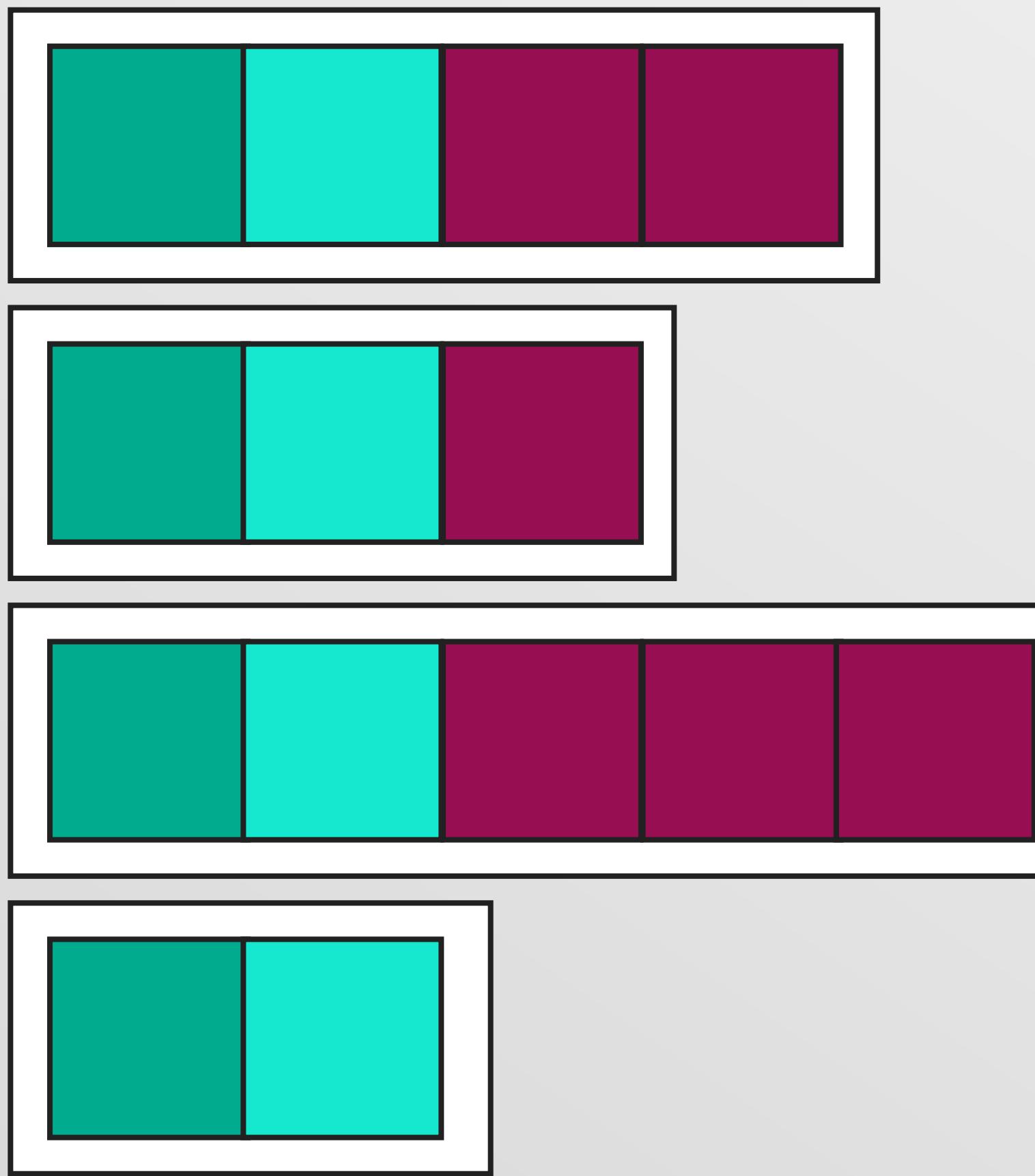
# Your turn

- Go to <https://quotes.toscrape.com/>
- Identify the rows and columns (including the URL to the author page), and the selectors that will identify them.
- (Challenge: how might store the tags in a list column? Why might this be useful? What makes it hard?)
- Scrape into a tibble

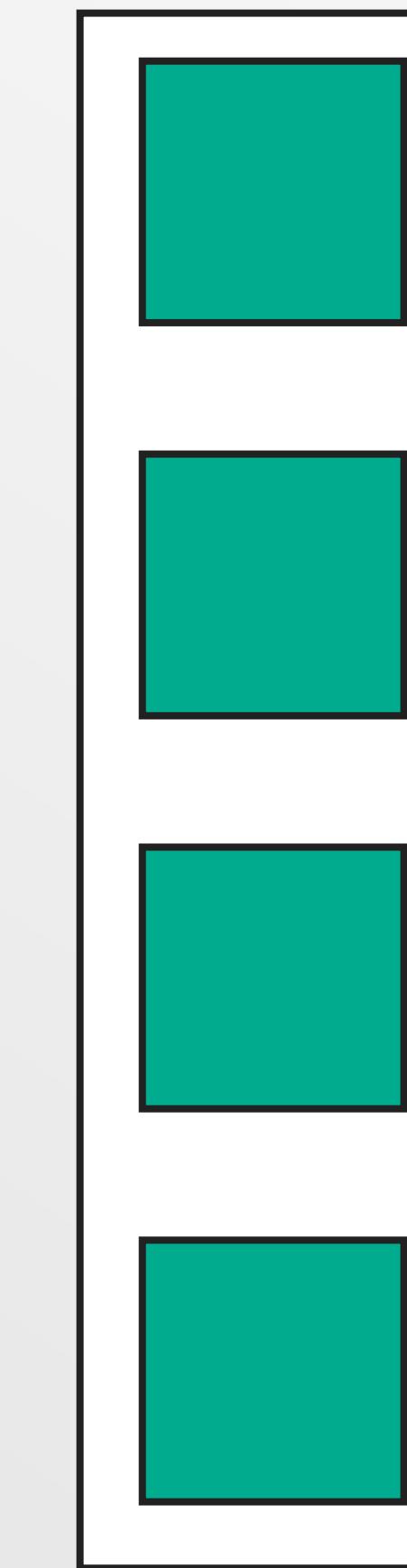
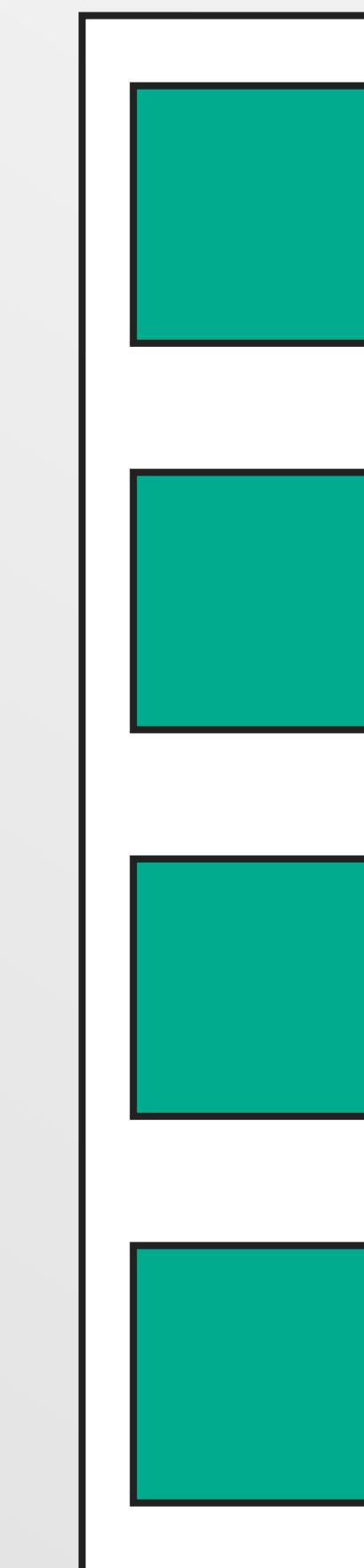
# Solution

quotes.R

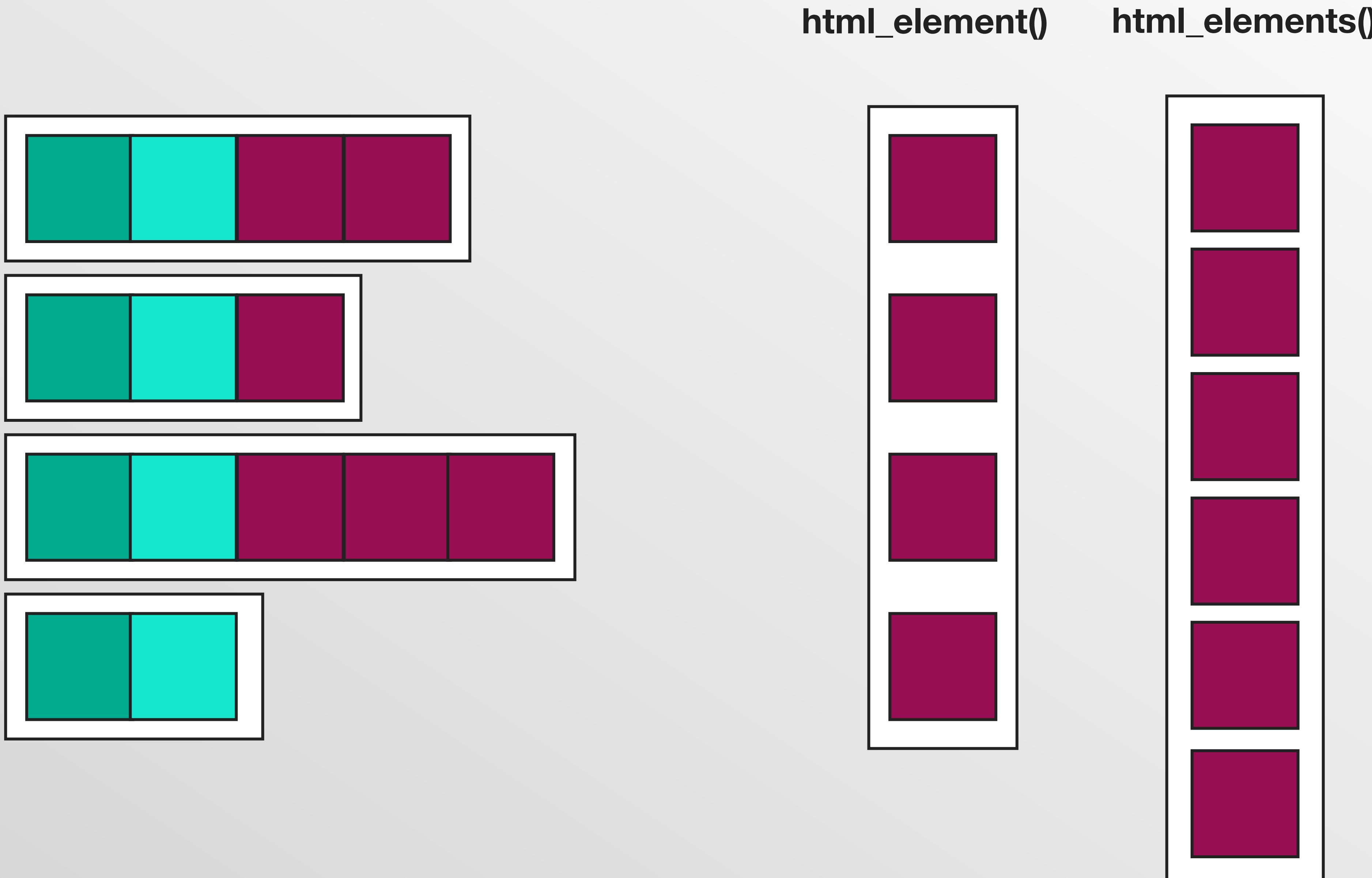
# There's no difference when there's one match



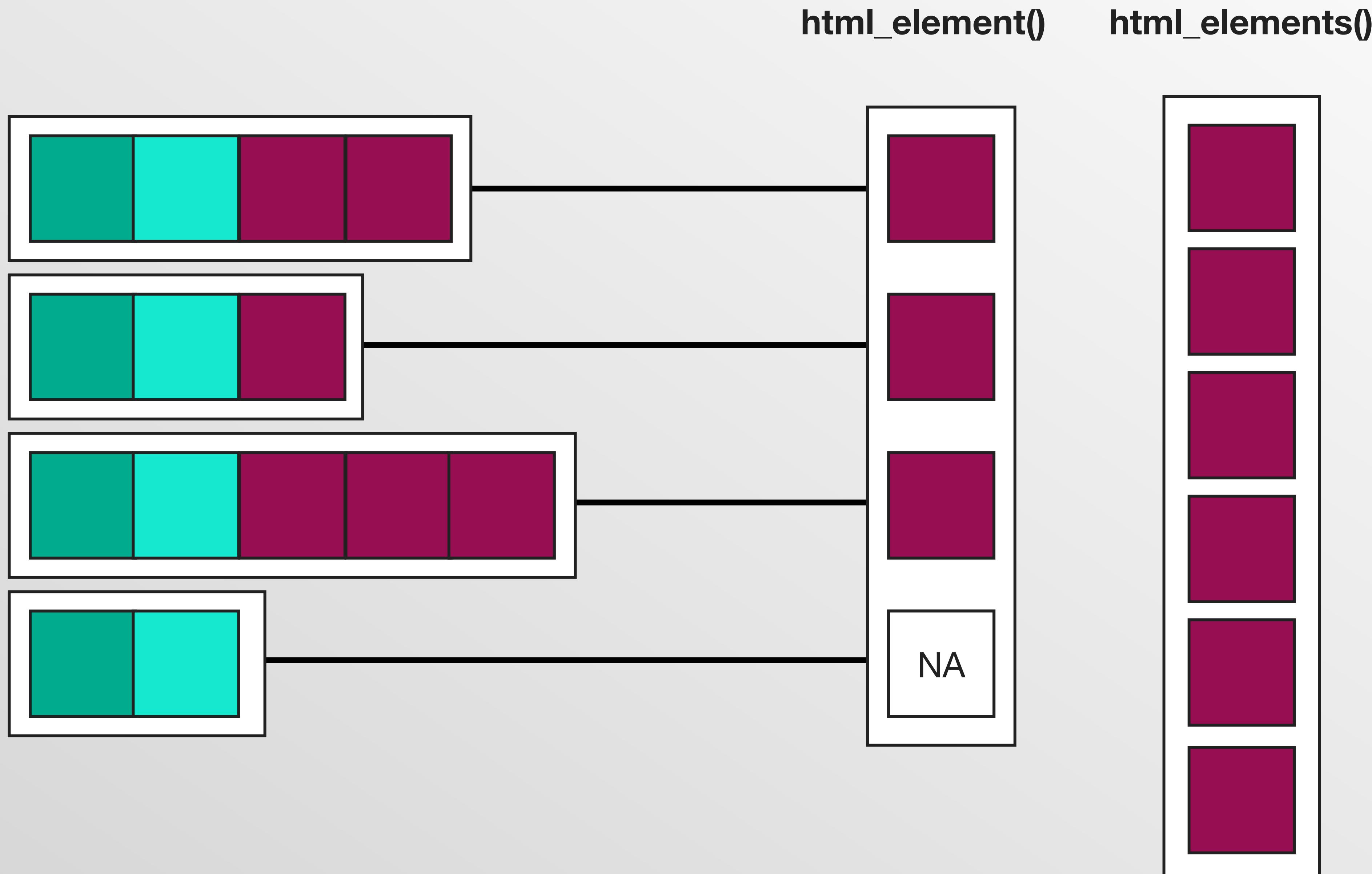
`html_element()`      `html_elements()`



# There's a big difference when numbers of matches varies



# html\_element() is always one-to-one



# `html_elements()` vs `html_element()`

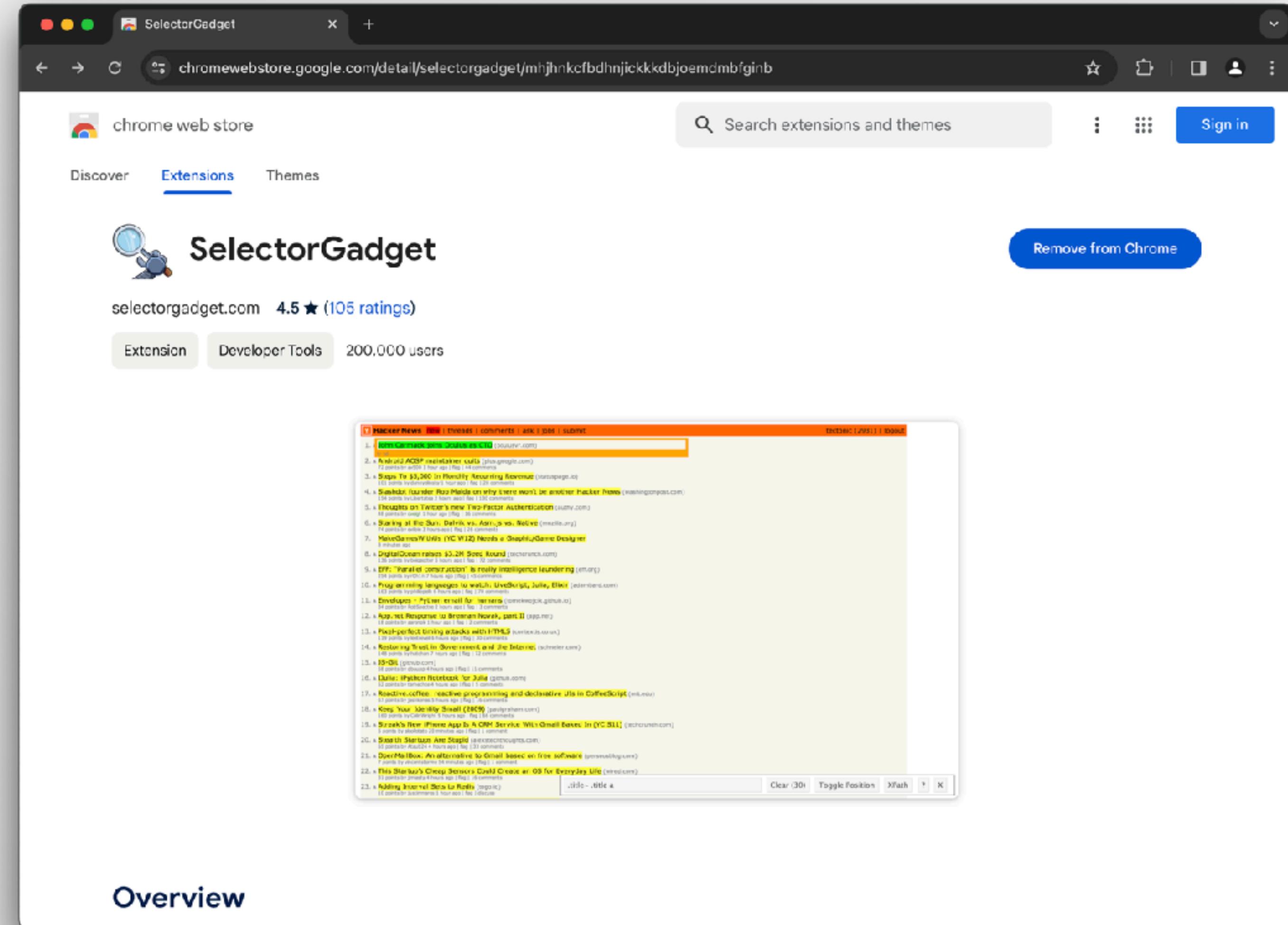
<code>html_elements()</code>	<code>html_element()</code>
$n \rightarrow m$	$n \rightarrow n$
<code>length(0)</code>	<code>NA</code>
Find rows	Find column in each row

# Ways to find the selector

1. Directly inspecting the HTML
2. In DevTools, right-click & choose “Copy selector” (then simplify)
3. SelectorGadget

# Google for selector gadget

<https://chrome.google.com/webstore/detail/selectorgadget>



## More practice

Scrape all the books off <<http://books.toscrape.com/>>

Make sure to capture the (full) title, the rating, the path to the cover image and the price.

(Don't worry about the pagination yet)

# Pagination

---

Demo

pagination.R

# Live HTML

---

F Forbes America's Top College +

forbes.com/top-colleges/    

≡    

# Forbes

# America's Top Colleges

## 2023

**Princeton University tops the 2023 list, while California public colleges continue to show their strength.**

Edited by Emma Whitford with Janet Novack



Illustration by Will Cordell For Forbes

**Forbes' annual list showcases 500 of the finest U.S. colleges, ranked using data on student success, return on investment and alumni influence.**

---

**E**very year, Forbes spotlights the 500 U.S. colleges that check all the boxes: impressive graduation rates, high graduate salaries, and great outcomes for low-income students, to name a few. Princeton University (#1)—with median early-career salaries north of \$88,000 and a 98% six-year graduation rate—aces these criteria. But it's not just the Ivy League that delivers on the promise of a quality education.

F Forbes America's Top College +

forbes.com/top-colleges/    

# Forbes

## America's Top Colleges List

FILTER LIST BY:

OVERALL STATE PUBLIC/PRIVATE SCHOOL SIZE CAMPUS SETTING

SEARCH 

RANK ^	NAME	STATE	TYPE	AV. GRANT AID	AV. DEBT	MEDIAN 10-YEAR SALARY	FINANCIAL GRADE
1	Princeton University	NJ	Private not-for-profit	\$47,136	\$7,216	\$177,300	A+
2	Yale University	CT	Private not-for-profit	\$58,715	\$4,968	\$163,900	A+
3	Stanford University	CA	Private not-for-profit	\$56,211	\$8,868	\$173,800	A+
4	Massachusetts Institute of Technology	MA	Private not-for-profit	\$32,562	\$7,235	\$182,800	A+
5	University of California, Berkeley	CA	Public	\$21,406	\$7,202	\$161,300	
6	Columbia University	NY	Private not-for-profit	\$57,726	\$13,338	\$150,900	A+
7	University of California, Los Angeles	CA	Public	\$17,592	\$5,965	\$140,300	
8	University of Pennsylvania	PA	Private not-for-profit	\$50,778	\$10,510	\$165,700	A+
9	Harvard	MA	Private not-for-profit	\$50,840	\$8,700	\$169,520	A+

  
**Top Colleges 2023: Why California's Public Universities Are So Good**  
It's no surprise that Princeton and Harvard made the top 25 colleges, but despite California's woes, four of its public universities did so too.  
[READ MORE >](#)

# Your turn

<<https://www.forbes.com/top-colleges/>>

Where does the data live? What defines the rows and the columns? (Hint: it looks like a table, but it's not)

What happens if you try to read this data into R?

# Solution

forbes-live.R

# Unofficial API

---

# Unofficial API

- Most websites that dynamically generate HTML do so from a JSON file.
- If you can find that JSON file you can work with it directly, making your life much easier.
- Find it with network pane in the browser developer tools.
- It's often obvious, but even when not it's worth spending 30 minutes on because it'll easily save that much time.
- (I think this Forbes site is an outlier; most of the time it will be a bit easier.)
- Another useful resource is <<http://inspectelement.org/apis.html>>



# Two useful heuristics to start with

- Filter to “Fetch/XHR” + Sort by size (decreasing)
- Use search to find text that you know must be in the data

DevTools - www.forbes.com/top-colleges/

Elements Console Sources Network Performance Memory Application Security Lighthouse Recorder > 18 ▲ 2794 1914 ⚙ :

Preserve log Disable cache No throttling Filter Invert Hide data URLs Hide extension URLs All Fetch/XHR Doc CSS JS Font Img Media Manifest WS Wasm Other

Blocked response cookies Blocked requests 3rd-party requests

20000 ms 40000 ms 60000 ms 80000 ms 100000 ms 120000 ms 140000 ms 160000 ms 180000 ms 200000 ms 220 ms

Name	Status	Type	Initiator	Size	Time	Waterfall
true.json?limit=1000&fields=organizationName,acade...u...	200	fetch	index-47e544c657c997b4.js:1	939 kB	160 ms	
auction	200	fetch	clientag.js?cid=8CU4ZAPFY...	76.7 kB	332 ms	
feature-decisions	200	fetch	444-a3070e0912093a5d.js:1	53.5 kB	264 ms	
auction	200	fetch	clientag.js?cid=8CU4ZAPFY...	53.1 kB	405 ms	
auction	200	fetch	clientag.js?cid=8CU4ZAPFY...	44.1 kB	278 ms	
auction	200	fetch	clientag.js?cid=8CU4ZAPFY...	41.0 kB	434 ms	
auction	200	fetch	clientag.js?cid=8CU4ZAPFY...	39.8 kB	410 ms	
TX_Ribbon_728x90.json	200	xhr	lottie.min.js:1	27.0 kB	28 ms	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	15.3 kB	1.48 s	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	13.8 kB	466 ms	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	13.5 kB	1.22 s	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	13.3 kB	2.98 s	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	13.3 kB	2.76 s	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	13.3 kB	2.42 s	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	13.1 kB	1.60 s	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	13.1 kB	1.87 s	
sodar?sv=200&tid=gda&tv=r20240221&st=env	200	xhr	show_ads_impl_fy2021.js:102	12.5 kB	126 ms	
sodar?sv=200&tid=gpt&tv=m202402200101&st=env	200	xhr	pubads_impl.js?cb=310813...	12.5 kB	66 ms	
sodar?sv=200&tid=gda&tv=r20240221&st=env	200	xhr	show_ads_impl_fy2021.js:102	12.2 kB	188 ms	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	11.7 kB	2.08 s	
pbjs?s=1039601	200	fetch	clientag.js?cid=8CU4ZAPFY...	8.1 kB	249 ms	
pbjs?s=1039599	200	fetch	clientag.js?cid=8CU4ZAPFY...	8.1 kB	372 ms	
ads?pvrid=1448819261062600&correlator=2040876848...	200	fetch	pubads_impl.js?cb=310813...	7.6 kB	847 ms	
auction	200	fetch	clientag.js?cid=8CU4ZAPFY...	7.5 kB	397 ms	
327?referer=https%3A%2F%2Fwww.forbes.com%2Fto...	200	xhr	script.js:147	6.8 kB	70 ms	
sodar?sv=200&tid=xfad&tv=01_247&st=int	200	xhr	Enabler_01_247.js:119	5.9 kB	58 ms	

210 / 1197 requests | 935 kB / 4.6 MB transferred | 12.1 MB / 79.0 MB resources | Finish: 3.0 min | DOMContentLoaded: 174 ms | Load: 181 ms

DevTools - www.forbes.com/top-colleges/

Elements Console Sources Network **Performance** Memory Application Security Lighthouse Recorder > 8 2149 2066 ⚡

Search  

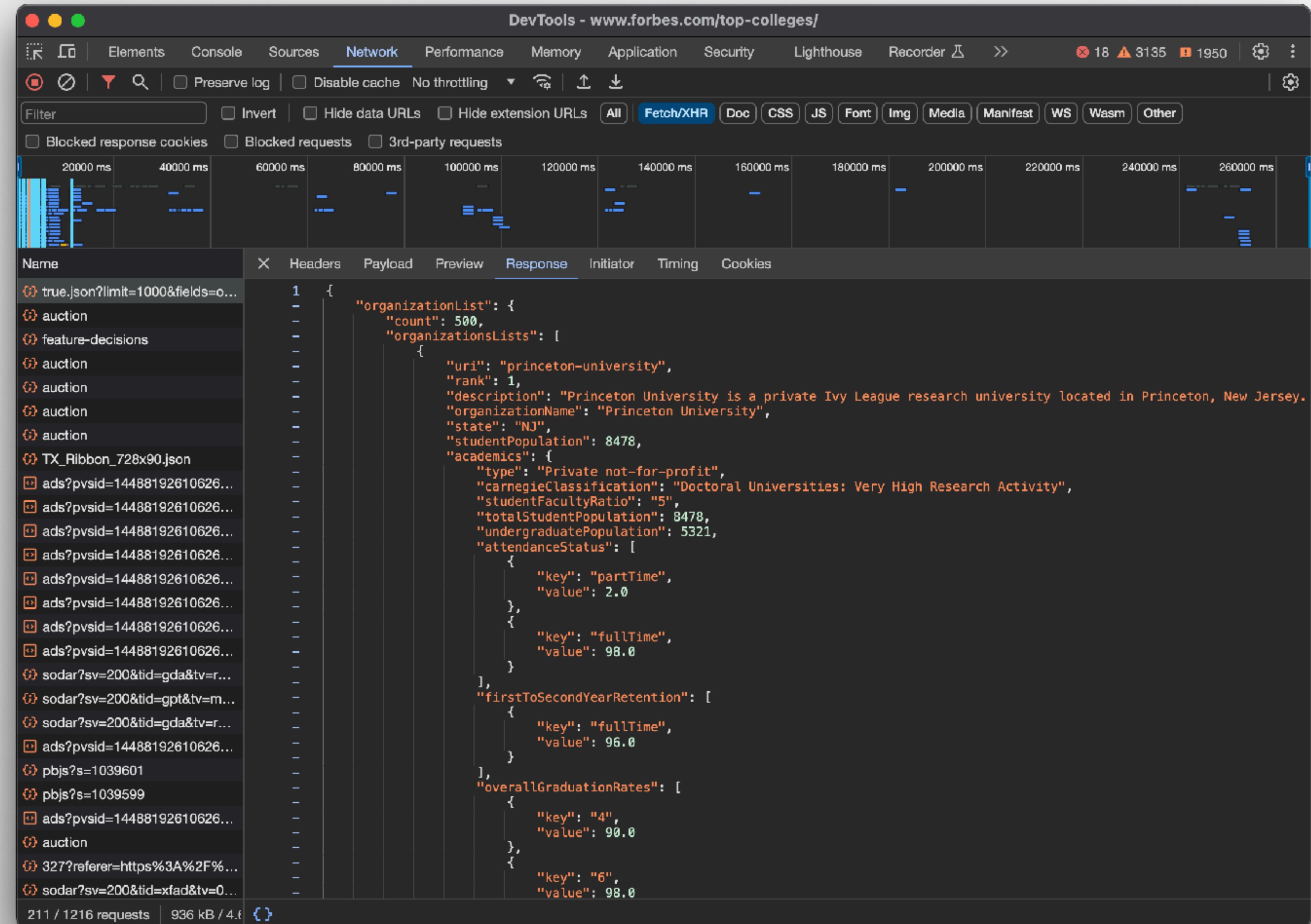
Aa payload.json bacon.forbes.com/bacon-forbes-prd/... 5000 ms 10000 ms 15000 ms 20000 ms 25000 ms 30000 ms 35000 ms

payload.json bacon.forbes.com/bacon-forbes-prd/... 5000 ms 10000 ms 15000 ms 20000 ms 25000 ms 30000 ms 35000 ms

1 ...It's no surprise that Princeton and Harvard ma...  
1 ...nce.", "IntroSubheading": "Princeton University to...  
1 ...students, to name a few. Princeton University (#1...  
  
true.json — www.forbes.com/forbesapi/org/top-colle...  
1 ...nizationsLists": [{"url": "princeton-university", "ran...  
1 ...rank": 1, "description": "Princeton University is a...  
1 ...ch university located in Princeton, New Jersey. A...  
1 ...ge in the United States, Princeton has a deep his...  
1 ...cial service professions. Princeton provides gener...  
1 ... class of 2026. In 2022, Princeton reported an en...  
1 ...021 freshman semester at Princeton University. "...  
1 ...y.", "organizationName": "Princeton University", "s...

Name	Status	Type	Initiator	Size	Time	Waterfall
▶ play.png	200	png	web_video.js:614	(disk cac...	1 ms	
CSI?v=2&s=ima&dmc=8&puid=o~ltln...	204	ping	web_video.js:290	17 B	383 ms	
interaction/?ai=BNCH_J8vtZbm3G_6...	200	gif	web_video.js:438	64 B	91 ms	
pixel.gif?e=25&q=2&hp=1&zMoatGNl...	200	gif	VM374:3	265 B	21 ms	
CSI?v=2&s=ima&dmc=8&puid=e~ltln...	204	ping	web_video.js:290	17 B	337 ms	
▶ play.png	200	png	web_video.js:614	(memory ...	0 ms	
CSI?v=2&s=ima&dmc=8&puid=f~ltln6...	204	ping	web_video.js:290	17 B	399 ms	
CSI?v=2&s=ima&dmc=8&puid=g~ltln...	204	ping	web_video.js:290	17 B	409 ms	
CSI?v=2&s=ima&dmc=8&puid=h~ltln...	204	ping	web_video.js:290	17 B	410 ms	
CSI?v=2&s=ima&dmc=8&puid=i~ltln6...	204	ping	web_video.js:290	17 B	410 ms	
CSI?v=2&s=ima&dmc=8&puid=j~ltln6...	204	ping	web_video.js:290	17 B	410 ms	
CSI?v=2&s=ima&dmc=8&puid=k~ltln...	204	ping	web_video.js:290	17 B	410 ms	
CSI?v=2&s=ima&dmc=8&puid=l~ltln6...	204	ping	web_video.js:290	17 B	409 ms	
CSI?v=2&s=ima&dmc=8&puid=m~ltln...	204	ping	web_video.js:290	17 B	409 ms	
CSI?v=2&s=ima&dmc=8&puid=n~ltln...	204	ping	web_video.js:290	17 B	409 ms	
interaction/?ai=B0TPYJ8vtZe29J8DZ...	200	gif	web_video.js:438	64 B	92 ms	
pixel.gif?e=9&q=1&hp=1&sst=1&ra=...	200	gif	VM39:2	265 B	170 ms	
dc_oe=ChMlyfCMivvphAMVIVIJCR3...	200	gif	express_html_inp...	63 B	50 ms	
dc_oe=ChMlqb_KivvphAMVtkcJCR1...	200	gif	express_html_inp...	63 B	50 ms	
dc_oe=ChMI--3bivvphAMVuI4JCR1B...	200	gif	express_html_inp...	63 B	53 ms	
interaction/?ai=B9a8ZJ8vtZdSiGsrj...	200	gif	web_video.js:438	64 B	89 ms	
dc_oe=ChMlwCDi_vphAMVb_EYAh2...	200	gif	express_html_inp...	63 B	50 ms	
dc_oe=ChMI-M30lvvphAMVl90YAh3...	200	gif	express_html_inp...	63 B	51 ms	
log?logid=kfk&evtid=pbad&itype=MA...	200	gif	clientag.js?cid=8C...	164 B	20 ms	
imsync.ashx?pl=3641587759298641...	200	script	tag.aspx?102:3	29 B	74 ms	
ping?h=forbes.com&p=%2Ftop-colle...	200	gif	chartbeat.js:29	200 B	43 ms	
pixel.gif?e=25&q=2&hp=1&kq=2&lo=...	200	gif	VM494:2	265 B	21 ms	
event.png?impid=8078c17e9221445...	204	ping	dv-measurements...	295 B	121 ms	
event.png?impid=64cf2de838b64f15...	204	ping	dv-measurements...	295 B	125 ms	
event.png?impid=0cd3a824c91747e...	204	ping	dv-measurements...	295 B	125 ms	
h?a=657665248&u=4194874624626...	200	gif	heap-657665248.j...	260 B	45 ms	

Search finished. Found 11 matching lines in 2 files. 1091 requests | 6.1 MB transferred | 83.1 MB resources | Finish: 32.68 s | DOMContentLoaded: 107 ms | Load: 110 ms



DevTools - www.forbes.com/top-colleges/

Elements Console Sources Network Performance Memory Application Security Lighthouse Recorder > > 18 3135 1950 | : | Preset log | Disable cache No throttling | Filter Invert Hide data URLs Hide extension URLs All Fetch/XHR Doc CSS JS Font Img Media Manifest WS Wasm Other Blocked response cookies Blocked requests 3rd-party requests

Name X Headers Payload Preview Response Initiator Timing Cookies

true.json?limit=1000&fields=o...  
auction  
feature-decisions  
auction  
auction  
auction  
TX\_Ribbon\_728x90.json  
ads?pvrid=14488192610626...  
ads?pvrid=14488192610626...  
ads?pvrid=14488192610626...  
ads?pvrid=14488192610626...  
ads?pvrid=14488192610626...  
ads?pvrid=14488192610626...  
ads?pvrid=14488192610626...  
ads?pvrid=14488192610626...  
sodar?sv=200&tid=gda&tv=r...  
sodar?sv=200&tid=gpt&tv=m...  
sodar?sv=200&tid=gda&tv=r...  
ads?pvrid=14488192610626...  
pbjs?s=1039601  
pbjs?s=1039599  
ads?pvrid=14488192610626...  
auction  
327?referer=https%3A%2F%...  
sodar?sv=200&tid=xfad&tv=0...

Open in Sources panel  
Open in new tab  
Clear browser cache  
Clear browser cookies

Copy >

Block request URL  
Block request domain

Sort By >  
Header Options >

Override headers  
Override content  
Show all overrides

Save all as HAR with content

Copy URL

Copy as cURL (highlighted)  
Copy as PowerShell  
Copy as fetch  
Copy as fetch (Node.js)

Copy response  
Copy stack trace

Copy all URLs  
Copy all as cURL  
Copy all as PowerShell  
Copy all as fetch  
Copy all as fetch (Node.js)

Copy all as HAR

212 / 1218 requests | 936 kB / 4.6 {}

# This gives a giant curl call

```
curl 'https://www.forbes.com/forbesapi/org/top-colleges/2023/position/true.json?  
limit=1000&fields=organizationName,academics,state,financialAid,rank,medianBaseSalary,campusSetting,studentPopulation,squareImage,uri,description,grade' \  
-X 'GET' \  
-H 'Accept: */*' \  
-H 'Sec-Fetch-Site: same-origin' \  
-H 'Cookie: _ga_DLD85VJ5QY=GS1.1.1706542437.3.1.1706542642.60.0.0; VWO=27.700;  
_ketch_consent_v1_=eyJzWhhdmlvcmFsX2FkdmVydGlzaW5nIjp7InN0YXR1cyI6ImdyYW50ZWQiLCJjYW5vbmljYWxQdXJwb3NlcI6WyJzWhhdmlvcmFsX2FkdmVydGlzaW5nIl19LCJhbhFseXRpY3MiOnsic3RhdHVzIjoiz3JhbnRlZCIsImNhbm9uaWNhb  
FB1cnBvc2VzIjpbImFuYWx5dGljcyJdfSwiZnVuY3Rpb25hbCI6eyJzdGF0dXMi0iJncmFudGVKIiwiY2Fub25pY2FsUHVycG9zZXMiOlsichJvZF9lbmhbmNlbWVudCIsInBlcnNvbmFsaXphdGlvbiJdfSwicmVxdWlyZWQiOnsic3RhdHVzIjoiz3JhbnRlZCIsI  
mNhbm9uaWNhbFB1cnBvc2VzIjpbImVzc2VudGlhbF9zZXJ2aNlcI6eyJdfX0%3D;  
_swb_consent_=eyJlbnZpcm9ubWVudENvZGUiOijwcm9kdWN0aW9uIiwiAWRlbnRpdGllcyI6eyJfZ29vZ2xlQW5hbHl0aNlcI6eyJfZ29vZ2xlQW5hbHl0aWNzQ2xpZW50SUQiOijHQTEuMi4zNTcxMjI4NjguMTcwNTI1OTE40CIsInN3Yl93ZWJzaXRLX3NtYXJ0X3RhZyI6IjU3YzY4NDNhLWQ4ZT  
ktNDBjMy05YmMxLThjYWNkZWU3ZDE20SJ9LCJqdXJpc2RpY3Rpb25Db2RlIjoidXNfZ2VuZXJhbCIsInByb3BlcnR5Q29kZSI6IndlYnNpdGVfc21hcnRfdGFnIiwigHvycG9zZXMiOnsiYW5hbHl0aWNzIjp7ImFsbG93ZWQiOij0cnVliiwibGVnYWxCYXNpc0NvZG  
UiOijkaXNjbG9zdXJlIn0sImJlaGF2aW9yYWxfYWR2ZXJ0aNpbmcioOnsiYWxsb3dlZCIsInRydWUiLCJsZdhbEJhc2lzQ29kZSI6ImRpc2NsB3N1cmUifSwiZnVuY3Rpb25hbCI6eyJhbGxvd2VkJoidHJ1ZSIsImxlZ2FsQmFzaXNDb2RlIjoiZGlzY2xvc3VyzS  
J9LCJyZXF1aXJlZCIsImxlZ2FsQmFzaXNDb2RlIjoiZGlzY2xvc3VyzSj9fSwiY29sbGVjdGVkQXQiOje3MDY1NDI10Dl9; _fbp=fb.1.1706121217898.164324328; _ga=GA1.2.357122868.1705259188;  
_gcl_au=1.1.1366914331.1705344922; _gid=GA1.2.412386523.1706542438; us_privacy=1---; usprivacy=1---; AWSALB=tnC2CAJGdokK0IsOKzBarRIkeT8r4YS5/  
wR6+n+A2VQX2Gfxa3lgP2KrEv6bGPRZyGIhGFvSCrNCftxLR8EXyMI3eDVjI5kBMLcIv9BthsZsyM9Vp0eTKR5yeR/H; AWSALBCORS=tnC2CAJGdokK0IsOKzBarRIkeT8r4YS5/  
wR6+n+A2VQX2Gfxa3lgP2KrEv6bGPRZyGIhGFvSCrNCftxLR8EXyMI3eDVjI5kBMLcIv9BthsZsyM9Vp0eTKR5yeR/H; BCSessionID=2ae07b14-70be-40b1-9e90-ed5d91b1be4f; ki_r=;  
ki_t=1706121218201%3B1706542438105%3B1706542587089%3B2%3B12; client_id=14d486f22e5f65b3107348cd0ffea50923; lux_uid=170654243670875530; AMP_TOKEN=%24NOT_FOUND; _swb=57c6843a-  
d8e9-40c3-9bc1-8cacdee7d169; rbzid=CXwv4IPKOa5WMsN9vGFFx/FfHDJXwSL2pHJ/swUFnuvHRpgd1qGWgsSHDvkfxbe6A3W7IDkaJbmIeiPvd/wC7NLIMs6nZ4N6B7HWA1lvWFqWciQ/+GZ1Bm7YCvrGGhX059tvv6mZYAXbx6MHil6+/  
BVKF18m2Z5gx0M2r0M2j9D/QiW7bH4TR8S/oJmWPQi7TJT5F+3SMGf6SjtfG64f74hLDwf+zEy5y95JETChlw70g8eg5Q05AALhKK+rhpV4z7F29XqrM2aGXgKH//+Q=; rbzsessionid=da3b925ab2cf238acdfa18f6e699045b;  
_ga_HY3LZWHH6W=GS1.1.1705344921.1.1.1705345703.0.0.0; _uetvid=a2f10940b3d711ee8c1881f6a3138e42; amp_9c5697=N1292876525...1hk77ksom.1hk77kvfo.2.2.4; notice_behavior=implied,us; fadve2etidvcnt=2;  
_clk=10sfr1j%7C2%7Cfif%7C0%7C1474; fadve2etid=N1292876525; fadvfpuid=FA77ce891e24882d9a03e9d2bc5bf16cf3; _ga_0Y2Y7WWQP1=GS1.1.1705259187.1.1.1705259215.0.0.0;  
_ga_JFZ3B3QM86=GS1.1.1705259187.1.1.1705259215.0.0.0; cmapi_cookie_privacy=permit 1,2,3; fadvuke2etid=N255829421; blaize_session=d72df655-7d08-4673-bbd3-beadbe316b40;  
blaize_tracking_id=418dd2e2-817f-48c7-9792-e36d6ce05bf9' \  
-H 'Sec-Fetch-Dest: empty' \  
-H 'Accept-Language: en-GB,en-US;q=0.9,en;q=0.8' \  
-H 'Sec-Fetch-Mode: cors' \  
-H 'Host: www.forbes.com' \  
-H 'User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/17.2.1 Safari/605.1.15' \  
-H 'Referer: https://www.forbes.com/top-colleges/' \  
-H 'Accept-Encoding: gzip, deflate, br' \
```

# And you can translate to R with `httr2::curl_translate()`

```
request("https://www.forbes.com/forbesapi/org/top-colleges/2023/position/true.json") ▷  
req_method("GET") ▷  
req_url_query(  
  limit = "1000",  
  fields = "organizationName,academics,state,financialAid,rank,medianBaseSalary,campusSetting,studentPopulation,squareImage,uri,description,grade",  
) ▷  
req_headers(  
  Accept = "*/*",  
  Cookie = "_ga_DLD85VJ5QY=GS1.1.1706542437.3.1.1706542642.60.0.0; VWO=27.700;  
_ketch_consent_v1=eyJzWhhdmIvcmFsX2FkdmVydGlzaW5nIjp7InN0YXR1cyI6ImdyYW50ZWQiLCJjYW5vbmljYWxQdXJwb3NlcI6WyJzWhhdmIvcmFsX2FkdmVydGlzaW5nIl19LCJhbhFseXRpY3MiOnsic3RhdHVzIjoiz3JhbnRlZCIsImNhbm9uaWNhbFB1cnBvc2VzIjpbImFuYWx5dGljcyJdfSwiZnVuY3Rpb25hbCI6eyJzdGF0dXMiOijncmFudGVkIiwiY2Fub25pY2FsUHVycG9zZXMiOlsichJvZF9lbmhbmNlbWVudCIsInBlcnNvbmFsaXphdGlvbijdfSwicmVxdWlyZWQiOnsic3RhdHVzIjoiz3JhbnRlZCIsImNhbm9uaWNhbFB1cnBvc2VzIjpbImVzc2VudGlhbF9zZXJ2aWNlcI6eyJfZ29vZ2xlQW5hbHl0aWNzQ2xpZW50SUQiOijHQTeuMi4zNTcxMjI4NjguMTcwNTI1OTE40CIsInN3Yl93ZWJzaXrlx3NtYXJ0X3RhZyI6IjU3YzY4NDNhLWQ4ZTktNDBjMy05YmMxLThjYWNkZWU3ZDE20SJ9LCJqdXJpc2RpY3Rpb25Db2RlIjoidXNFZ2VuZXJhbCIsInByb3BlcnR5Q29kZSI6IndlYnNpdGVfc21hcnRfdGFniIwicHVycG9zZXMiOnsiYW5hbHl0aWNzIjp7ImFsbG93ZWQiOij0cnVlIiwibGVnYWxCYXNpc0NvZGUioiJkaXNjbG9zdXJlIn0sImJlaGF2aW9yYWxfYWR2ZXJ0aNpbmcioNSiYWxs3dlZCI6InRydWUiLCJsZWdhbEJhc2lzQ29kZSI6ImRpc2Ns3N1cmUifSwiZnVuY3Rpb25hbCI6eyJhbGxvd2VkJoidHJ1ZSiImxlZ2FsQmFzaXNDb2RlIjoiZGlzY2xvc3VyzSj9LCJyZXF1aXJlZCI6eyJhbGxvd2VkJoidHJ1ZSiImxlZ2FsQmFzaXNDb2RlIjoiZGlzY2xvc3VyzSj9fSwiY29sbGVjdGVkQXQiOje3MDY1NDI10Dl9; _fbp=fb.1.1706121217898.164324328; _ga=GA1.2.357122868.1705259188;  
_gcl_au=1.1.1366914331.1705344922; _gid=GA1.2.412386523.1706542438; us_privacy=1---; usprivacy=1---; AWSALB=tnC2CAJGdokK0IsOKzBarRIkeT8r4YS5/wR6+n+A2VQX2Gfxa3lgP2KrEv6bGPRZyGIhGFvSCrNCftxLR8EXyMI3eDVjI5kBMLcIv9BthsZsyM9Vp0eTKR5yeR/H; AWSALBCORS=tnC2CAJGdokK0IsOKzBarRIkeT8r4YS5/wR6+n+A2VQX2Gfxa3lgP2KrEv6bGPRZyGIhGFvSCrNCftxLR8EXyMI3eDVjI5kBMLcIv9BthsZsyM9Vp0eTKR5yeR/H; BCSessionID=2ae07b14-70be-40b1-9e90-ed5d91b1be4f; ki_r=; ki_t=1706121218201%3B1706542438105%3B1706542587089%3B2%3B12; client_id=14d486f22e5f65b3107348cd0ffea50923; lux_uid=170654243670875530; AMP_TOKEN=%24NOT_FOUND; _swb=57c6843a-d8e9-40c3-9bc1-8cacdee7d169; rbzid=CXwv4IPKOa5WMsN9vGFFx/FfHDJXwSL2pHJ/swUFnuvHRpgd1qGWgsSHDvkfxbe6A3W7IDkaJbmIeiPvd/wC7NLDIMs6nZ4N6B7HWA1lvWFqWciQ/+GZ1Bm7YCvrGGhX059tvv6mZYAXbx6MHil6+/BVKF18m2Z5gx0M2r0M2j9D/QiW7bh4TR8S/oJmWPQi7JT5F+3SMGf6SjtfG64f74hLDwf+zEy5y95JETChlw70g8eg5Q05AALhKK+rhpV4z7F29XqrM2aGXgKh//+Q==; rbzsessionid=da3b925ab2cf238acdfa18f6e699045b; _ga_HY3LZWHH6W=GS1.1.1705344921.1.1.1705345703.0.0.0; _uetvid=a2f10940b3d711ee8c1881f6a3138e42; amp_9c5697=N1292876525 ... 1hk77ksom.1hk77kvfo.2.2.4; notice_behavior=implied,us; fadve2etidvcnt=2; _clck=10sfr1j%7C2%7Cfif%7C0%7C1474; fadve2etid=N1292876525; fadvfpuid=FA77ce891e24882d9a03e9d2bc5bf16cf3; _ga_OY2Y7WWQP1=GS1.1.1705259187.1.1.1705259215.0.0.0; _ga_JFZ3B3QM86=GS1.1.1705259187.1.1.1705259215.0.0.0; cmapi_cookie_privacy=permit 1,2,3; fadvuke2etid=N255829421; blaize_session=d72df655-7d08-4673-bbd3-beadbe316b40; blaize_tracking_id=418dd2e2-817f-48c7-9792-e36d6ce05bf9",  
`Accept-Language` = "en-GB,en-US;q=0.9,en;q=0.8",  
Host = "www.forbes.com",  
`User-Agent` = "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/605.1.15 (KHTML, like Gecko) Version/17.2.1 Safari/605.1.15",  
`Accept-Encoding` = "gzip, deflate, br",  
) ▷
```

# Then simplify to the essentials

```
url ← "https://www.forbes.com/forbesapi/org/top-colleges/2023/  
position/true.json"
```

```
req ← request(url) ▷  
  req_url_query(limit = "1000") ▷  
  req_perform()
```

# Demo

forbes-api.R

# If the data is stored as JSON in the HTML

This is a pain, but is fortunately relatively rare.

I have had some luck in the past with using the V8 R package to run the javascript code and then extract the JSON object back into R.

# LLMs

---

# Houston pollen and mold counts

The screenshot shows a web browser window displaying the Houston Health Department's website. The page title is "Houston Pollen and Mold Count - Friday, July 5, 2024". The main content area features four categories of pollen and mold counts, each with an associated image:

- TREE POLLEN**: NONE (0)
- WEED POLLEN**: MEDIUM (24)
- GRASS POLLEN**: MEDIUM (6)
- MOLD SPORES**: MEDIUM (7,604)

Below this, a section titled "Major tree pollen counted" lists two entries:

- Acer (Maple): 0
- Morus (Mulberry): 0

The browser's address bar shows the URL: <https://www.houstonhealth.org/services/pollen-mold>.

<https://www.houstonhealth.org/services/pollen-mold>

The following ranges are based on the data reported by all counting stations for weeds, grasses, trees and mold according to the [National Allergy Bureau](#).

**Weed pollen is measured per cubic meter of air.**

- 1-9 pollen per cubic meter of air or less = Low
- 10 to 49 = Medium
- 50 to 499 = Heavy
- Greater than 500 = Extremely Heavy

**Grass pollen is measured per cubic meter of air.**

- 1-4 pollen per cubic meter of air or less = Low
- 5 to 19 = Medium
- 20 to 199 = Heavy
- Greater than 200 = Extremely Heavy

**Tree pollen is measured during the tree pollen season of mid-January through mid-April.** Cedar Elm tree pollen is measured in September and October. All tree pollen types are reported by HDHHS laboratory.

- 1 to 14 pollen per cubic meter of air = Low
- 15 to 89 = Medium
- 90 to 1499 = Heavy
- Greater than 1500 = Extremely Heavy

**Mold spores are measured per cubic meter of air.**

- 1 to 6499 = Low
- 6500 to 12999 = Medium
- 13000 to 49999 = Heavy
- Over 50000 = Extremely Heavy

**The HHD laboratory reports all major types of mold spore irritants.**

# Prompt

Take the following data and convert it into a call to tibble::tribble, like the following two example rows for pollen

```
thresholds <- tibble(  
  ~category, ~low, ~high, ~value,  
  "weed", 1, 9, "low",  
  "weed", 500, Inf, "extremely high"  
)
```

Here's the data ...