

Visit <https://bit.ly/lede-ai-playwright> for material

Magical Scraping!

stealing (??) data from
the internet with the magic
of Playwright + AI

HTML

the language that all web
sites are built with

```
<h1>This is a headline</h1>
<h2>This is a smaller headline</h2>
<h3>And an even smaller headline</h3>
```

This is a headline

This is a smaller headline

And an even smaller headline

```
<h1>This is a headline</h1>
<h2>This is a smaller headline</h2>
<h3>And an even smaller headline</h3>
<p>This is a paragraph</p>

<p>This is a paragraph with
<a href="google.com">a link</a></p>
<div>This is ANYTHING</div>
<div>Anything! Anything in the world</div>
```

This is a headline

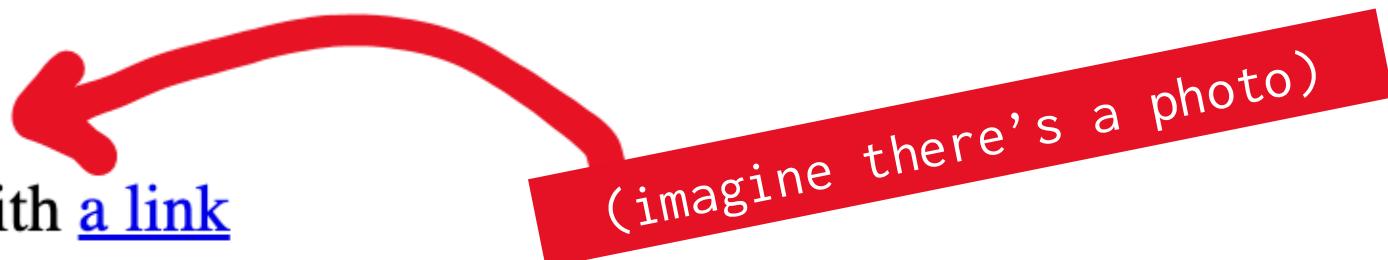
This is a smaller headline

And an even smaller headline

This is a paragraph

This is a paragraph with a link

This is ANYTHING
Anything! Anything in the world



(imagine there's a photo)

<h1>An incredible story</h1>

<p>By J. Soma</p>

<p>This is the start of the story.</p>

<p>It is an amazing story!</p>

<p>”It’s incredible,” said the source.</p>

<p>An editor agreed: “it’s true!”</p>

<p>Additional reporting by Mulberry the fat
mean cat</p>

<p>Edit: An earlier version said “Jonathan”
Soma</p>

An incredible story

By J. Soma

This is the start of the story.

It is an amazing story!

"It's incredible," said the source.

An editor agreed: "it's true!"

Additional reporting by Mulberry the fat mean cat

Edit: An earlier version said "Jonathan" Soma

<h1>An incredible story</h1>

<p>By J. Soma</p>

<p>This is the start of the story.</p>

<p>It is an amazing story!</p>

<p>”It’s incredible,” said the source.</p>

<p>An editor agreed: “it’s true!”</p>

<p>Additional reporting by Mulberry the fat
mean cat</p>

<p>Edit: An earlier version said “Jonathan”
Soma</p>

```
<h1>An incredible story</h1>
<p id="byline">By J. Soma</p>
<p>This is the start of the story.</p>
<p>It is an amazing story!</p>
<p>"It's incredible," said the source.</p>
<p>An editor agreed: "it's true!"</p>
<p class="additional">Additional reporting by
Mulberry the fat mean cat</p>
<p class="correction">Edit: An earlier version
said "Jonathan" Soma</p>
```

An incredible story

By J. Soma

This is the start of the story.

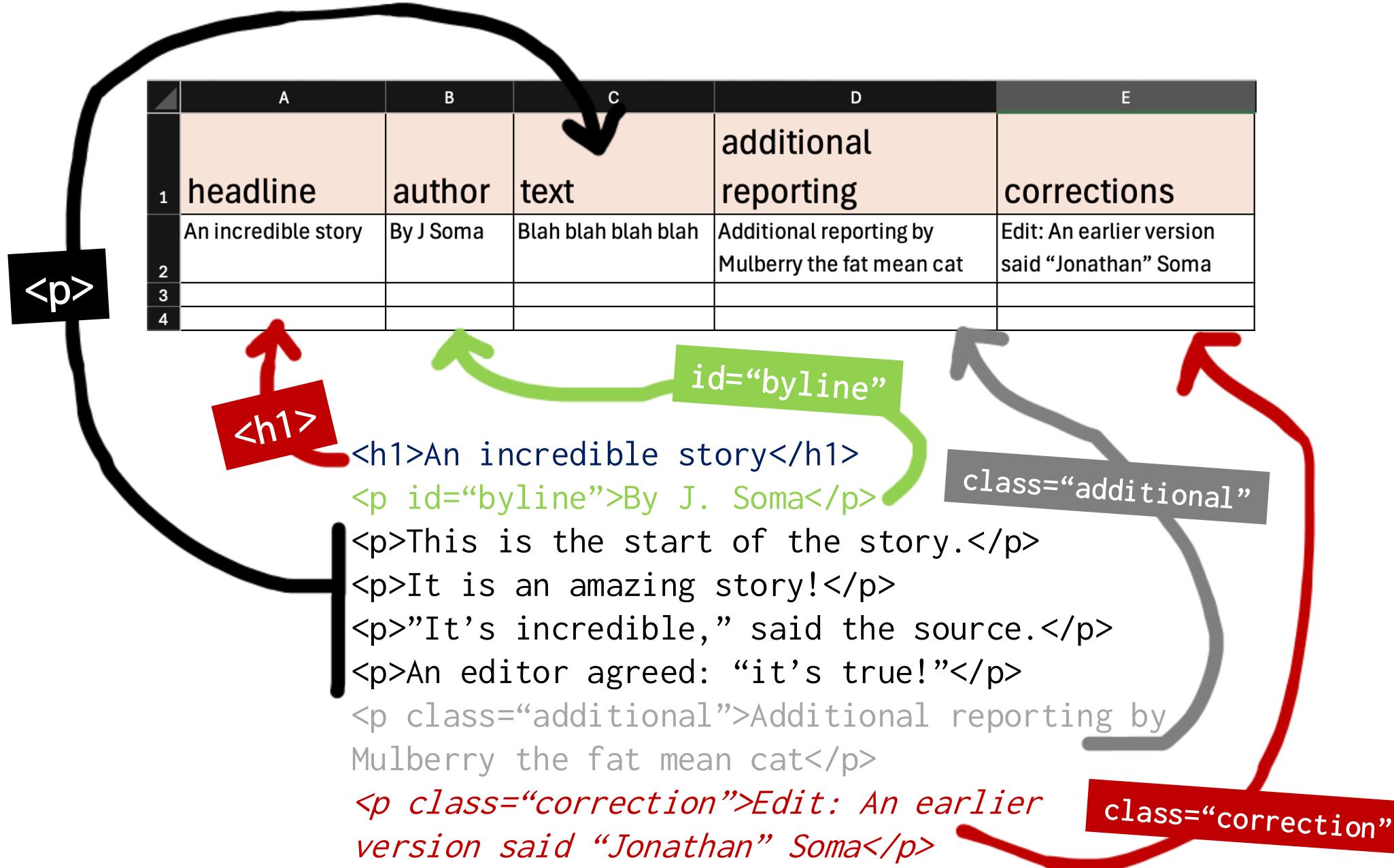
It is an amazing story!

”It’s incredible,” said the source.

An editor agreed: “it’s true!”

Additional reporting by
Mulberry the fat mean cat

*Edit: An earlier version said “Jonathan”
Soma*



Let's see it on
the internet!

bbc.com

BBC Home - Breaking News

Home News Sport Business Innovation Culture Travel Earth Video Live

Register Sign In

What will be a row in our spreadsheet?

 **LIVE VP pick Walz to speak as Democrats deploy star guests**

The Minnesota governor will run alongside presidential nominee Kamala Harris for November's US election.

 **Three things the Democrats have avoided so far at the DNC**

What they have tried to avoid says as much about their weaknesses as what they choose to highlight, writes Anthony Zurcher.

2 hrs ago | US & Canada

- Day three of Democratic Convention features party's rising stars
- Obamas mock Trump over 'black jobs' and crowd sizes

 **Divers find five bodies in wreck of Sicily yacht**

 **Gaza nurse says whole family, including quadruplets, killed in air strike**

5 hrs ago | Middle East

 **Andrew Tate held overnight after police raid homes in Romania**

The internet personality and his brother have been remanded in custody as part of a probe into new allegations.

44 mins ago | Europe

 **Ancient ocean of magma found on Moon south pole**

The findings are from India's historic Chandrayaan-3 mission that landed on the Moon's south pole.

6 hrs ago | Science & Environment

► **German Navy blasts out Darth Vader theme on Thames**

BBC Home - Breaking News, +

bbc.com

Home News Sport Business Innovation Culture Travel Earth Video Live

Register Sign In

 **LIVE VP pick Walz to speak as Democrats deploy star guests**

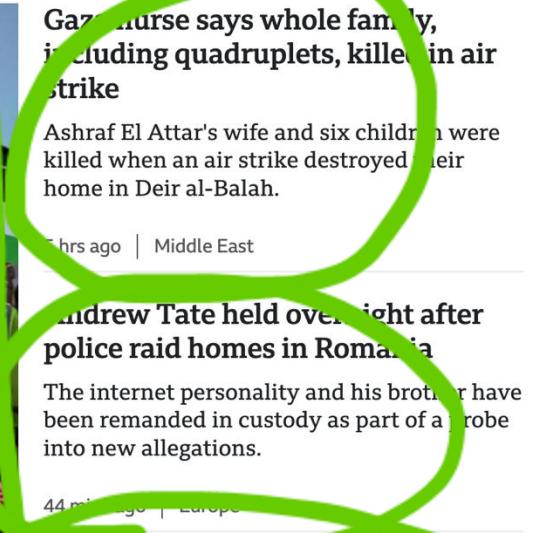
The Minnesota governor will run alongside presidential nominee Kamala Harris for November's US election.

 **Three things the Democrats have avoided so far at the DNC**

What they have tried to avoid says as much about their weaknesses as what they choose to highlight, writes Anthony Zurcher.

- Day three of Democratic Convention features party's rising stars
- Obamas mock Trump over 'black jobs' and crowd sizes

 **Divers find five bodies in wreck of Saily yacht**

 **Gaza nurse says whole family, including quadruplets, killed in air strike**

Ashraf El Attar's wife and six children were killed when an air strike destroyed their home in Deir al-Balah.

 **Ancient ocean of magma found on Moon south pole**

The findings are from India's historic Chandrayaan-3 mission that landed on the Moon's south pole.

 **► German Navy blasts out Darth Vader theme on Thames**

What will be a column of data?

Andrew Tate held overnight after police raid homes in Romania

The internet personality and his brother have been remanded in custody as part of a probe into new allegations.

44 mins ago

| Europe



LIVE VP pick Walz to speak as
Democrats deploy star guests

The Minnesota governor will run alongside
presidential nominee Kamala Harris for
November's US election.

メルカリ - 日本最大のフリマサー × +

jp.mercari.com/search?category_id=79

楽器・機材 ×

ログイン 会員登録 ベル 出品

絞り込み クリア 楽器・機材 の検索結果 ↑ おすすめ順 この検索条件を保存する

除外キーワード カテゴリー ホビー・楽器・アート 楽器・機材 すべて ブランド サイズ 価格 価格なし出品 あんしん鑑定 割引オプション 商品の状態

What will be a row in our spreadsheet?

商品名	価格
korg kross61 電池駆動軽量化シンセ	¥28,000
iQ7 ステレオマイク Lightning接続 (修理品)	¥3,800
ギターストラップ 未使用	¥990
音楽之友社 うたとピアノの絵本 1 みぎで	¥700
バイオリン用 頸当て 黒檀 訳あり	¥1,500
YUCKY様専用 Pioneer DDJ-FLX4	¥40,000
ACアダプター JH35P1200150D	¥880
GoogleChromecast グーグル クロームキャスト	¥4,700
日本（フジゲン）製 ibanez SR1000 フレット ...	¥55,000
SE ELECTRONICS DM1 DYNAMITE	¥8,500

<https://jp.mercari.com/item/m97630098490>

メルカリ - 日本最大のフリマサー × +

jp.mercari.com/search?category_id=79

楽器・機材 ×

ログイン 会員登録 ベル 出品

絞り込み クリア

除外キーワード

カテゴリ

ホビー・楽器・アート

楽器・機材

すべて

ブランド

サイズ

価格

価格なし出品

あんしん鑑定

割引オプション

商品の状態

検索結果

販売中のみ表示

↑↓ おすすめ順 この検索条件を保存する

エフェクター ジャンク ベース bare fender chase パワーサプライ オカリナ kemper ホルン

¥28,000 Korg kross61 電池駆動軽量化シンセ

¥3,800 iQ7 ステレオマイク Lightning接続 (修理品)

¥990 ギターストラップ 未使用

¥700 音楽之友社 うたとピアノの絵本 1みぎで

¥1,500 バイオリン用 頸当て 檜 訳あり

¥40,000 YUCKY様専用 Pioneer DDJ-FLX4

¥880 ACアダプター JH35P1200150D

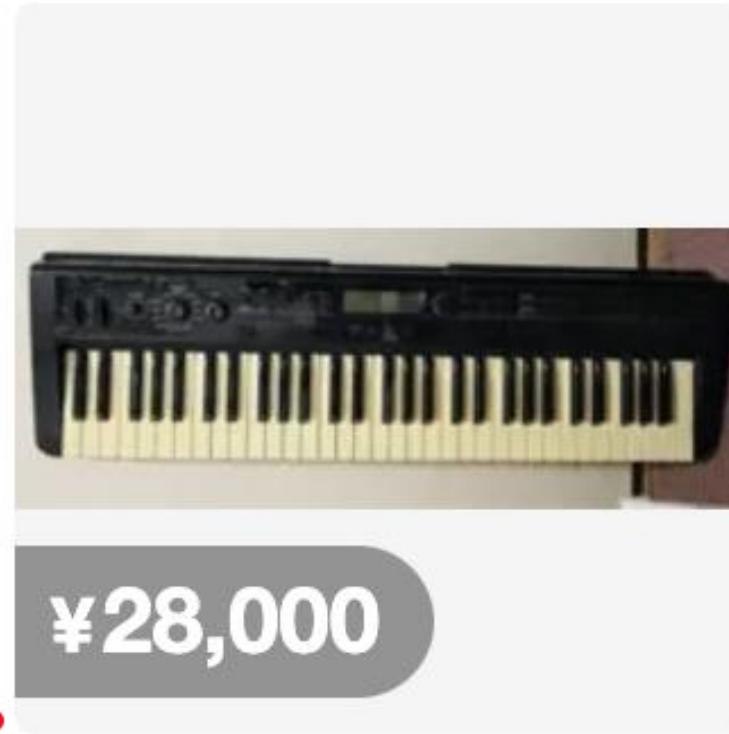
¥4,700 GoogleChromecast グーグル クロームキャスト ...

¥55,000 日本(フジゲン)製 ibanez SR1000 フレット ...

¥8,500 SE ELECTRONICS DM1 DYNAMITE

<https://jp.mercari.com/item/m97630098490>

What will be a column of data?



price

¥28,000

name

korg kross61 電池駆動軽
量化シンセ

image

Computers don't see like us
(usually), we need to see

HTML code

Visit a page on Mercari

Do this now

【2024年最新】楽器・機材の人

ログイン 会員登録 ベル 出品

楽器・機材 の検索結果

クリア

除外キーワード

カテゴリー

ホビー・楽器・アート

楽器・機材

ブランド

サイズ

価格

価格なし出品

あんしん鑑定

割引オプション

¥1,250 新品 D'Addario ダダリオ アコースティックギター弦

¥5,500 新品 D'Addario ダダリオ アコースティックギター弦

¥15,000 Pearl スネアドラム MUS1455M

¥2,400 新品 D'Addario ダダリオ アコースティックギター弦

¥399 バタフライフィンガーピック 4個セット ゴ...
バタフライフィンガーピック 4個セット ゴ...

¥579 オイル漬け牛骨製ナット

¥580 YOASOBI 群青 ピアノ楽譜 ソロ

¥1,800 新品 D'Addario ダダリオ アコースティックギター弦

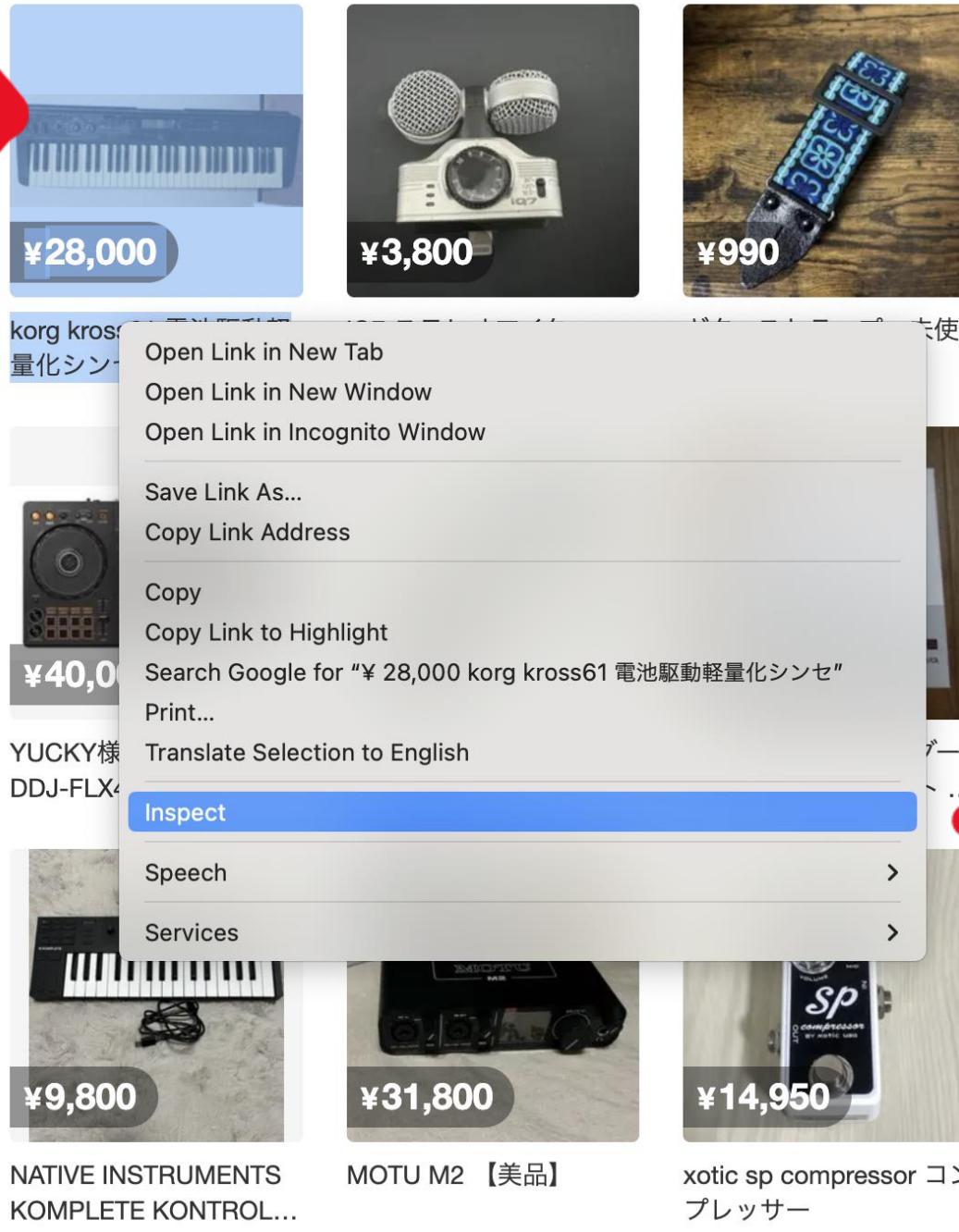
¥680 新品 D'Addario ダダリオ アコースティックギター弦

¥31,000 ♪森の工房♪フルートの店 おかげさまで100本超
YAMAHA 211SII 銀メッキ
美品 分解磨き 調整済み
部活応援!!銀メッキ!!美品!!調整済!!ヤマハフル...

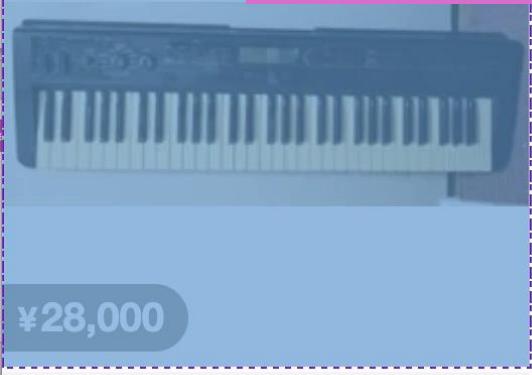
↑↓ おすすめ順 この検索条件を保存する

エフェクター ジャンク ベース bare fender chase パワーサプライ オカリナ kemper ホルン

Right click (or
command+click)
and Inspect



Find the HTML
for our “row”



METHOD
ONE

この検索条件を保存する

Mouse/click
around



会員登録 ログイン



```
<div id="item-grid" data-testid="search-item-grid">
  <ul class="sc-50e0525c-0 sc-bcd1c877-0 JKB0Z ipWzwL"> grid
    <li data-testid="item-cell" class="sc-bcd1c877-2 cvAXgx">
      <div>
        <a href="/item/m48062282876" data-location="search_result:best_match:body:item_list:item_thumbnail" class="sc-bcd1c877-1 lpjZwE" data-testid="thumbnail-link"> flex
          <div class="merItemThumbnail fluid_a6f874a2" role="img" aria-label="kross61 電池駆動軽量化シンセの画像" data-bbox="212 480 418 780" data-testid="itemThumbnail"> 28,000円" id="m48062282876" itemtype="ITEM_TYPE_MERCARI">
            <figure class="itemThumbnail_a6f874a2">
              <div class="before_a6f874a2" aria-hidden="true"></div>
              <div.imageContainer_f8ddf3a2 picture img>
                <img alt="A small thumbnail image of a DJ-style audio mixer." data-bbox="422 568 612 756" data-testid="itemThumbnailImage"/>
              </div>
            </figure>
          </div>
        </a>
      </div>
    </li>
  </ul>
</div>
```

Find the HTML
for our “row”



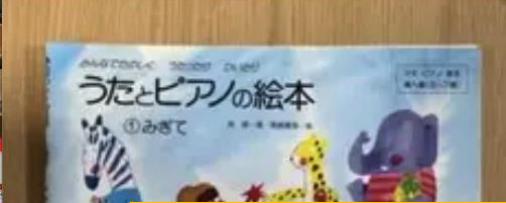
¥28,000



¥3,800



¥990



¥700



¥40,000

METHOD
TWO

この検索条件を保存する



会員登録

ログイン



Click
this

Search + click
here

```
<div id="item-grid" data-testid="search-item-grid">
  <ul class="sc-50e0525c-0 sc-bcd1c877-0 JKBOZ ipWzwL"> grid
    <li data-testid="item-cell" class="sc-bcd1c877-1 lpjZwE" data-testid="thumbnail-link"> flex
      <div class="merItemThumbnail fluid_a6f874a2" role="img" aria-label="korg kross61 電池駆動式シンセの画像
        28,000円"
      </div>
    </li>
  </ul>
</div>
```

◀ 7-0.JKBOZ.ipWzwL li.sc-bcd1c877-2.cvAXgx ▶

Styles Computed Layout Event Listeners >

Filter

:hov .cls +, □

```
element.style {  
}
```

amazon.com
search results

Sales at Lidl

List of amphibians
from Germany's Red
List Center

Headlines from
Nikkei Asia

items on
Mercari

Now try it.

Anywhere!!!!

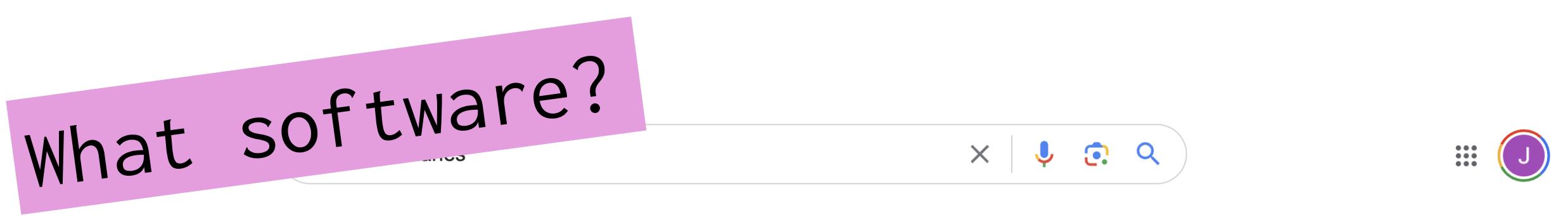
Wikipedia's
Fictional Big
Cats

School Board Minutes
from Grand Island Public
Schools in Nebraska

Scraping is just connecting to a website
and pulling data into your spreadsheet.
But we need some software
to do that for us!

Open up your browser and go to

your spreadsheet



All Images Videos News Books Web Finance

Tools



Scraping libraries

From sources across the web



Selenium



Lxml



Playwright



ZenRows S.L.



Beautiful Soup



Requests



Puppeteer



Cheerio



Scrapy



MechanicalSoup



Urllib3



[\[–\]](#) **Swingbiter**  **70 points** 2 years ago

Learn the basic html elements that build up a website.

 [-] **coventous**  22 points 2 years ago

 I recommend checkin

W 1 point 2 years ago

Look into learning and understanding HTML, then, so Google for that. Oh and bonus tip, there's (official? chrome/etc) which made me very happy (I like Docker!)

2 ITEM DATA RESPONSE.000

[-] luizv4z  **12 points** 2 years ago

From my own research, run away from Selenium. The right direction is CDP (Chrome Developers Protocol).

```
2 all_links = soup.find_all(name=
```

~~Do python on them until~~

Bea  [-] riisen  3 points

for small projects use built-in tools, its am

Beginner projects go with scrapy, its an

for bigger projects
add save save-RES repo

[permalink](#) [source](#) [embed](#) [save](#)

F.S. permit, unless

[permalink](#) [source](#) [embed](#) [save](#) [save-RES](#) [report abuse](#)

permalink source embed save save-RES report

[permalink](#) [source](#) [embed](#) [save](#) [save-RES](#) [report](#) [reply](#) [hide child comments](#)

-] ned334 5 points 2 years ago

Google "Selenium find_element(By.XPATH, '/XPATH/')"

All elements have an XPath that you can copy from chrome by Inspect -> right click on code block -> copy full Xpath.

Scraping solved

Selenium is garbage!

...but BeautifulSoup
can't scrape all sites.

Top discounts for you

Page 1 of 3



Continental Grand Prix
5000 Rennrad Faltreifen //
28-622 (700x28C)

85

400+ viewed in past month

\$66.00

FREE Delivery Friday,
Aug 23



Wink Super X-Large Ultra
Thin Lubricated Latex
Condoms, Premium Latex
for Smooth and Natural...

76

100+ viewed in past month

\$14.99 (\$0.30/Count)

FREE Delivery Friday,
Aug 23



Assortment Rubber Ducks
in Bulk, 50-Pack Assorted
Mini Duckies Toy for
Ducking Cruise Ships, 2"

328

5K+ viewed in past month

\$24.99

Overnight by 8:00
AM



50 Pack Rubber Ducks in
Bulk, Jeep Ducks for
Ducking, Assorted Rubber
Ducks Jeep Ducking, Bab...

569

4K+ viewed in past month

\$21.89

Overnight by 8:00
AM



ArtCreativity Assorted
Rubber Ducks Jeep Ducking
(100Pack) - Rubber Ducki...

982

700+ viewed in past month

Limited time deal

\$39.98

List: \$49.99

Today by 10:00 PM



Kicko Assorted Rubber
Ducks with Mesh Bag - 50
Ducklings, 2 Inch – Jeep
Ducks for Kids, Baby Bath...

2,543

2K+ viewed in past month

\$28.99

FREE Delivery
Saturday, Aug 24



Glitter Rubber Ducks in
Bulk - (Pack of 50) Assorted
2-inch Duck Toys for Baby
Shower Rubber Duckies,...

308

1K+ viewed in past month

\$26.79

Today by 10:00 PM



Bulk savings to consider

Page 1 of 7



Chochkees Asso

Ducks Toy Duckies for Kids
and Toddlers, Bath Birthday

Baby Showers Classroom,...

1,155



Rubber Ducks Jeep Ducking
(50 Pack) - Rubber Duckie...

100+ viewed in past month



Set, Mini Colorful Rubber
Duckies Bath Toy for
Child,Float & Squeak Tiny...

763



Bulk, Jeep Ducks for
Ducking, Assorted Rubber
Ducks Jeep Ducking, Bab...

569



Ducks: Fun Unique Military-
Inspired Bath Toys for Jeep
Ducking or Play - 2 inches

800+ viewed in past month



Bulk, Jeep Ducks for
Ducking, Assorted Rubber
Ducks for Jeeps, Bath Toy...

569



Glitter Rubber Duck Toy
Assortment Duckies for
Kids, Bath Birthday Gifts...



Some sites need interaction

Company



TEXAS DEPARTMENT OF LICENSING & REGULATION

TDLR License Data Search (Active Licenses only)

[Search Help](#) | [Download License files](#) | [Download Other](#) | [Questions/Comments](#)

Inquire by License Type	Inquire by License #
Choose One (Optional)	<input type="text"/> (Numeric only)
Inquire by Expiration Date	
<input type="text"/> (mmddyyyy)	
Inquire by Name (Last, First) or by Business Name	
<input type="text"/>	
Inquire by Location (City)	
Choose One (Optional)	Type the first letter to scroll down.
Inquire by County	
Choose One (Optional)	Type the first letter to scroll down.
Inquire by Zip Code	
<input type="text"/>	
<input type="button" value="Search"/> <input type="button" value="Reset"/>	

If license not found, please contact Customer Service at 800-803-9202

Data last updated: 8/21/2024 06:01

[Bookmark This Page](#)

Some sites need interaction

Fast and reliable end-to-end 🐍

playwright.dev/python/

Playwright for Python Docs API Python ▾ Community

GitHub Discord Moon Search

Playwright enables reliable end-to-end testing for modern web apps.

GET STARTED ⚡ Star 64k+

I love Playwright!



Playwright is perfect!

But! It's new, so

ChatGPT isn't very good
at it. *But!* We try.

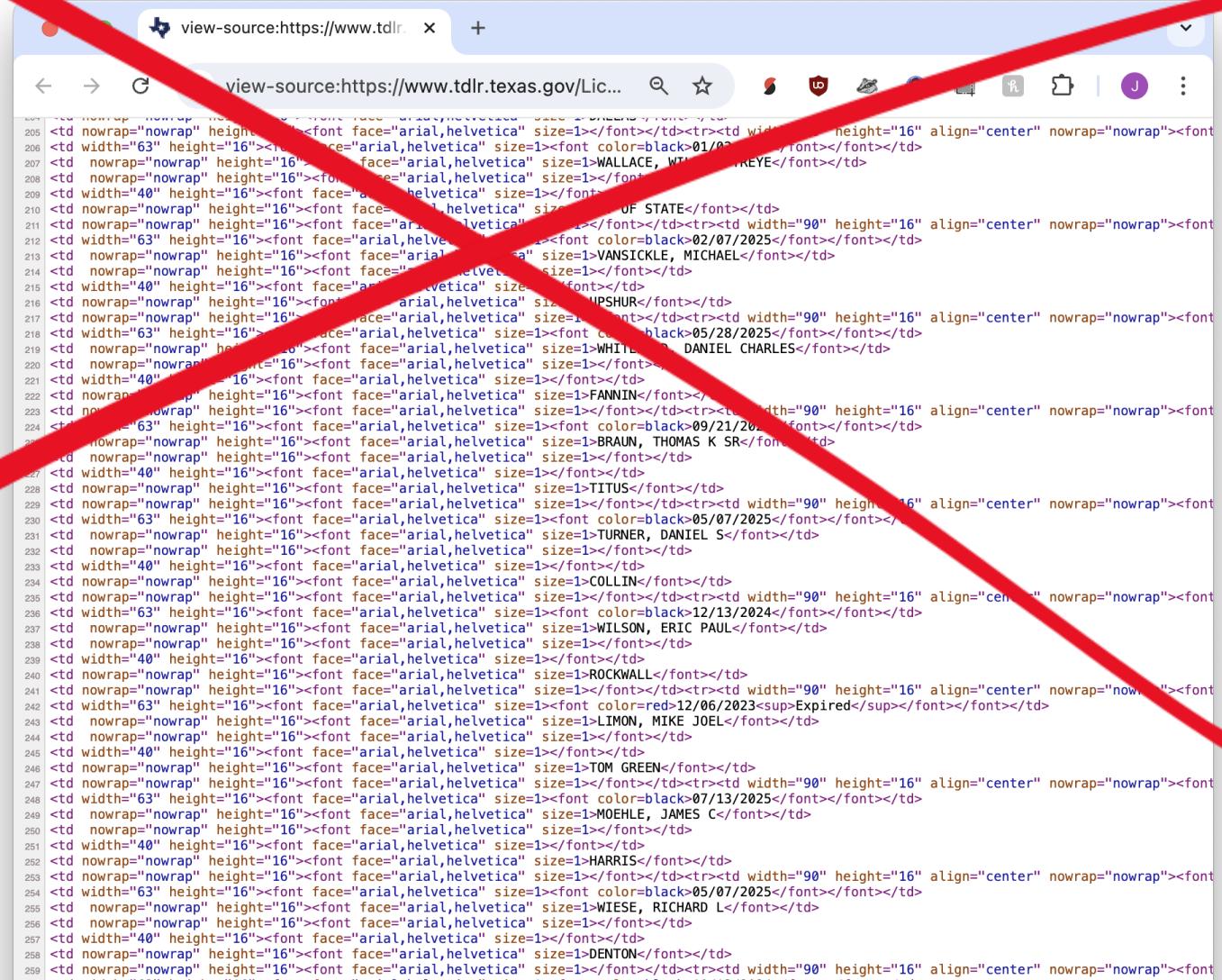
Playwright + ChatGPT +
pasting samples of HTML

=

infinite scrapers

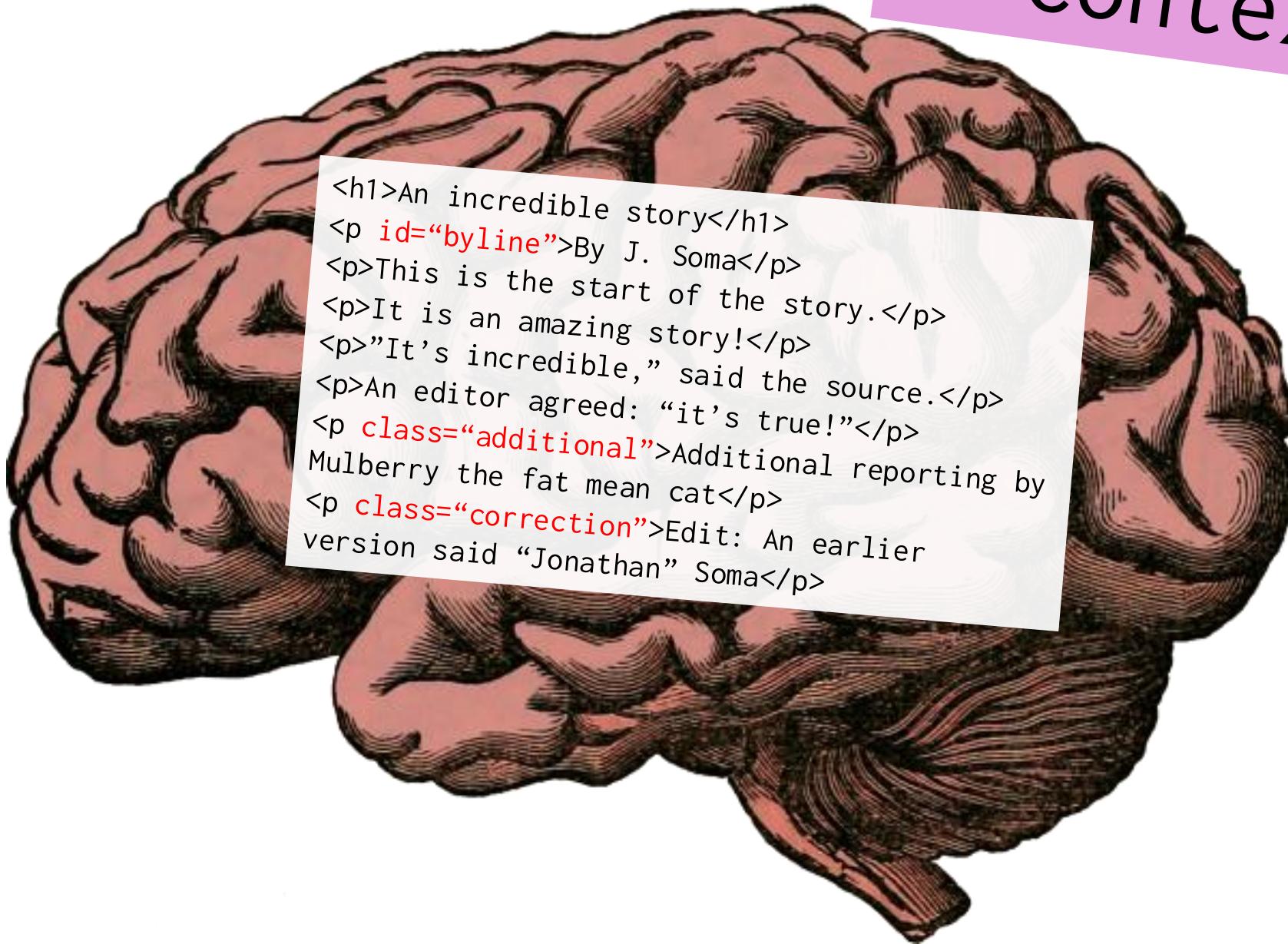


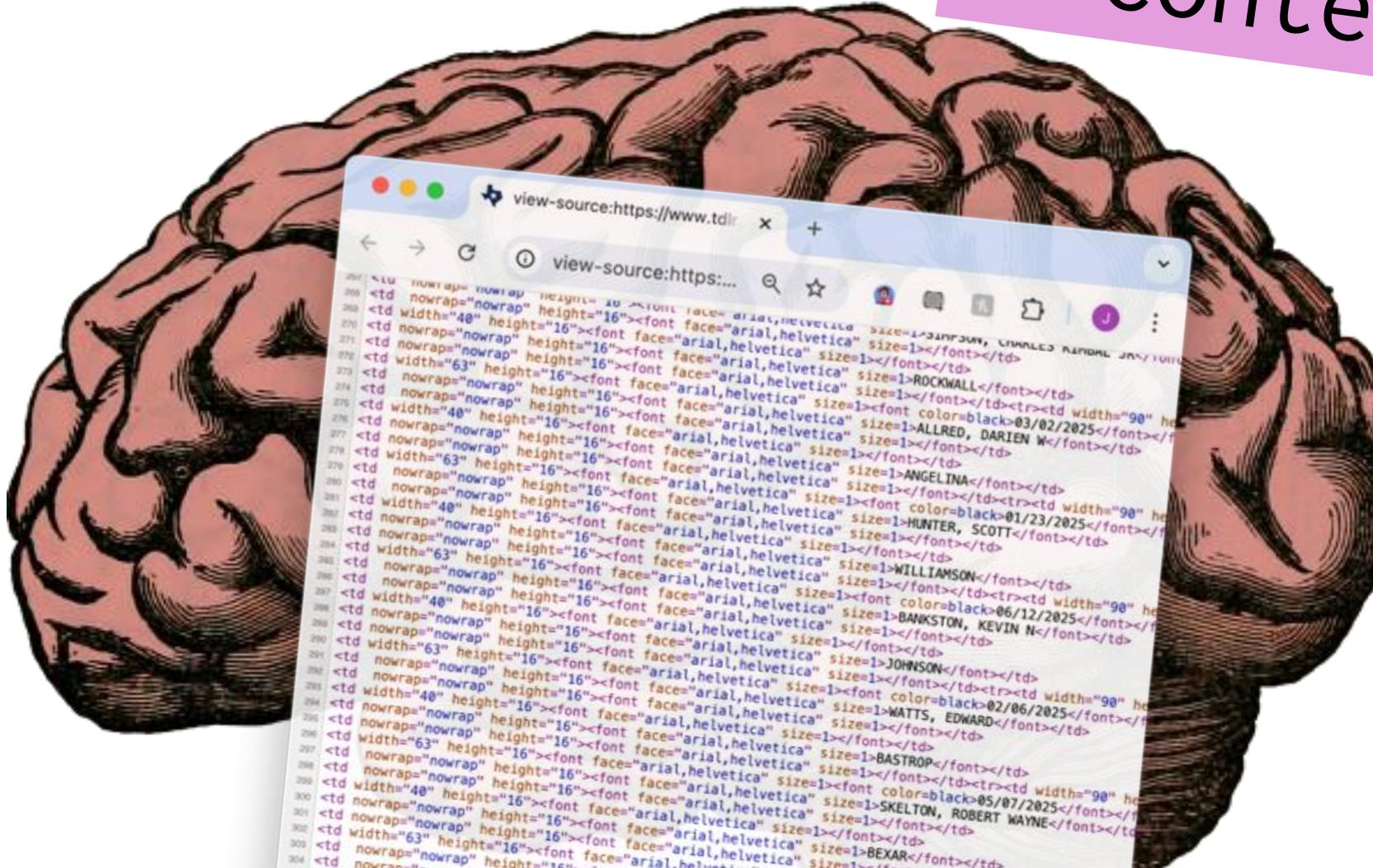
“Write a scraper for these search results:”



The image shows a browser window with the URL `view-source:https://www.tdlr.texas.gov/Lic...` in the address bar. The main content area displays the raw HTML source code of a search results page. The code consists of numerous `<td>` and `<tr>` elements, primarily in black and white, with some blue and purple text appearing as hyperlinks. A large, thick red 'X' is drawn across the entire content area, indicating that the task of writing a scraper for this specific page is no longer applicable.

“Context window”





“Context window”

What we might need:

- Start URL
- Forms to fill out (*optional*)
- “Rows” of our spreadsheet
- Pagination/“next” pages (*optional*)
- A magical prompt

Visit <https://bit.ly/lede-ai-playwright> for material

Magical Scraping!

stealing (??) data from
the internet with the magic
of Playwright + AI