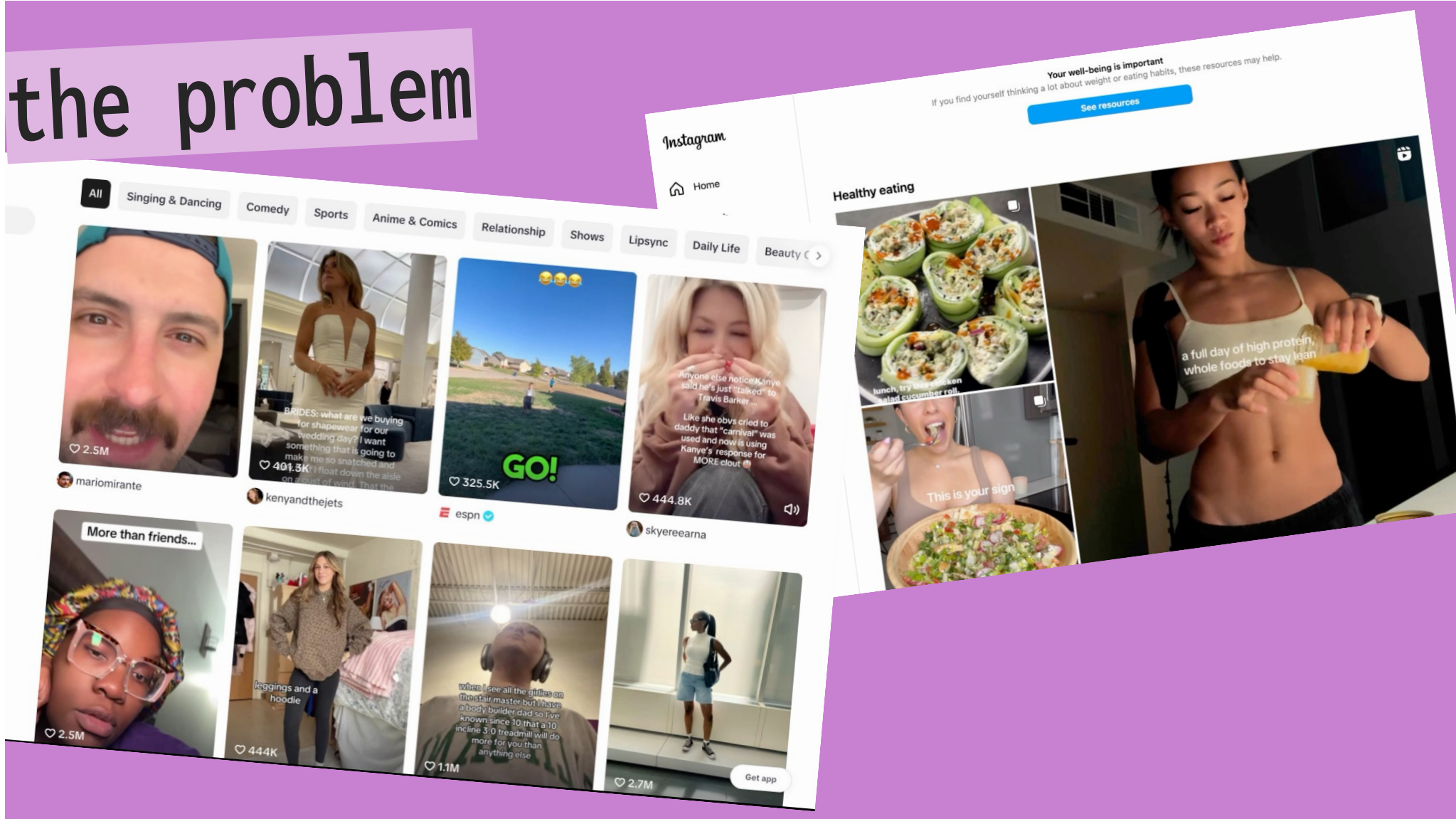# Passive Scraping for social media

(and everything else)

Jonathan Soma
Columbia Journalism School
js4571@columbia.edu
@dangerscarf

# the problem

# tiers of problem-solving

- Using a tool
- Writing a scraper
- Undocumented APIs
- Intercepting browser requests
- Pack-ratting with HAR and WARC/WACZ files

# A solution: instaloader

```python
import instaloader

L = instaloader.Instaloader()
L.login(username, password)

profile = instaloader.Profile.from_username(L.context, "catrepublic")

for index, post in enumerate(profile.get_posts()):
    if index >= 10:
        break
    L.download_post(post, target=username)
```

catrepublic

Q Search

2025-02-26_22-09-20_UTC_1.jpg
2025-02-26_22-09-20_UTC_2.jpg
2025-02-26_22-09-20_UTC.txt
2025-02-28_01-51-06_UTC_1.jpg
2025-02-28_01-51-06_UTC_2.jpg
2025-02-25_23-20-32_UTC.jpg
2025-02-25_02-26-49_UTC_1.jpg
2025-02-25_02-26-49_UTC_2.jpg

2025-02-28_01-51-06_UTC_comments.json
2025-02-28_01-51-06_UTC.json
2025-02-28_01-51-06_UTC.txt
2025-03-01_23-30-30_UTC.jpg
2025-03-01_23-30-30_UTC.json
2025-02-25_23-20-32_UTC.txt
2025-02-25_16-23-25_UTC_comments.json
2025-02-25_16-23-25_UTC.json

2025-03-01_23-30-30_UTC.txt
2025-03-03_22-35-58_UTC.jpg
2025-03-03_22-35-58_UTC.json
2025-03-03_22-35-58_UTC.txt
2025-02-26_22-09-20_UTC_comments.json
2025-02-26_22-09-20_UTC.json
2025-02-25_02-26-49_UTC.jpg
2025-02-24_08-11-28_UTC.txt

2025-02-25_23-20-32_UTC.mp4
2025-02-25_16-23-25_UTC_1.jpg
2025-02-25_16-23-25_UTC_2.jpg
2025-02-25_16-23-25_UTC_3.jpg
2025-02-25_16-23-25_UTC_4.jpg
2025-02-25_23-20-32_UTC.json
2025-02-24_08-11-28_UTC_comment

Dropbox > Soma > Curriculum > 2025-nicar > passive-scraping > catrepublic

37 items, 750.19 GB available

2025-02-25_16-23-25_UTC_comments.json

[
  {
    "id": "18068495113697485",
    "created_at": 1740582827,
    "text": "Beautiful Mango..... Has Oreo Bear, and Princess Peach been adopted.??????",
    "owner": {
      "id": "45571346482",
      "is_verified": false,
      "profile_pic_url": "https://scontent-lga3-1.cdninstagram.com/v/
t51.2885-19/476916367_2952187921610547_3712381710209805776_n.jpg?stp=dst-jpg_s150x150_tt6&_nc_cat=103&_nc_oc=Q6cZ2AFL1q_r5e6UWK5WggcZ5xRMWx8XzJTxCS7shJnAFw9mcd
lga3-1.cdninstagram.com&_nc_ht=scontent-
QcfvOw8fmfRE8U9bvnZltcc&_nc_ohc=AWVoaATp76YQ7kNvgG_gS1C&_nc_gid=74e0f069cd894b2194bbae626b257
urVMDIZA1tjbM3nUpFoT5aJDwWkLod1b-hy_lG_xow&oe=67CD7941&_nc_sid=d885a2",
      "username": "neo_a_cat60"
    },
    "likes_count": 0,
    "answers": []
  },
  {
    "id": "18053046008144998",
    "created_at": 1740527874,
    "text": "Handsome",
    "owner": {
      "id": "1533405592",
      "is_verified": false,
      "profile_pic_url": "https://scontent-lga3-3.cdninstagram.com/v/t51.2885-19/10735036_70
lga3-3.cdninstagram.com&_nc_cat=106&_nc_ht=scontent-
QcfvOw8fmfRE8U9bvnZltcc&_nc_ohc=vXa2qs8w5ZwQ7kNvgH8QEyW&_nc_gid=74e0f069cd894b2194bbae626b25771e&
k04_XMfovwKGYGPzUnctgitanXdtPk9eMIRu-CXr7AA&oe=67CD758B&_nc_sid=d885a2",
      "username": "ali_cat_tom_cat"
    },
    "likes_count": 1,
    "answers": []
  },
  {
    "id": "18053122958147007",
    "created_at": 1740520750,
    "text": "\u2764MANGO!",
    "owner": {
      "id": "3583278215",
      "is_verified": false,
      "profile_pic_url": "https://scontent-lga3-2.cdninstagram.com/v/t51.2885-19/13712147_1087141904698593_1898202199_
stp=dst-jpg_s150x150_tt6&_nc_ht=scontent-
lga3-2.cdninstagram.com&_nc_cat=109&_nc_oc=Q6cZ2AFL1q_r5e6UWK5WggcZ5xRMWx8XzJTxCS7shJnAFw9mcdkDd
QcfvOw8fmfRE8U9bvnZltcc&_nc_ohc=pzcQV3iSR8cQ7kNvgEpDfSs&_nc_gid=74e0f069cd894b2194b
6bt8FGtita45oxJEgIBmc51dAp7uR_cbhvhWVJAcQcg&oe=67CD5F3A&_nc_sid=d885a2",
      "username": "rickyrealrogue"
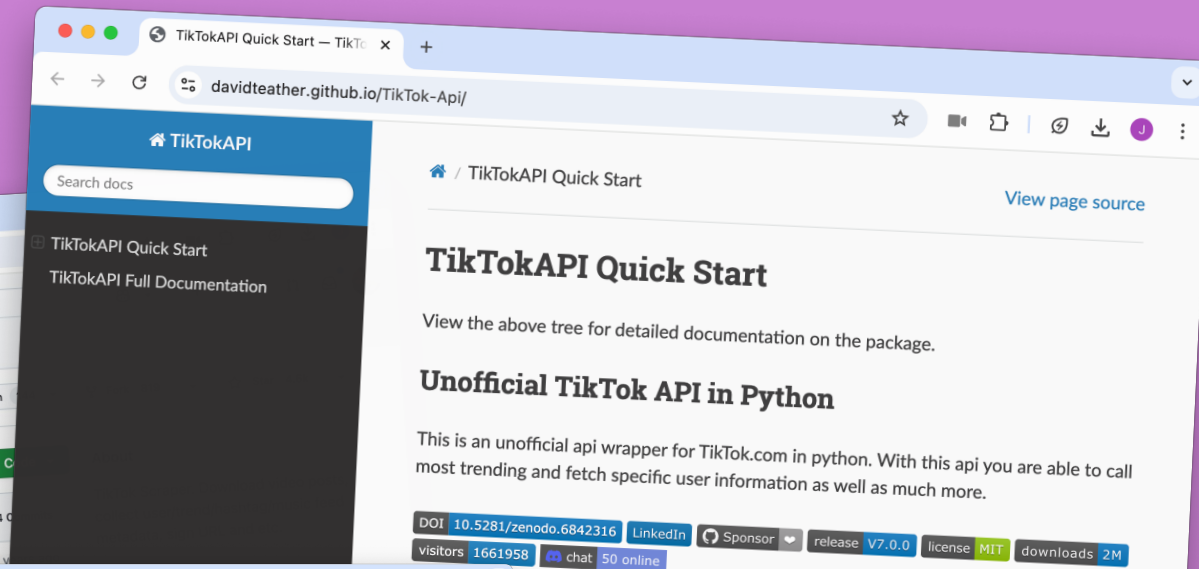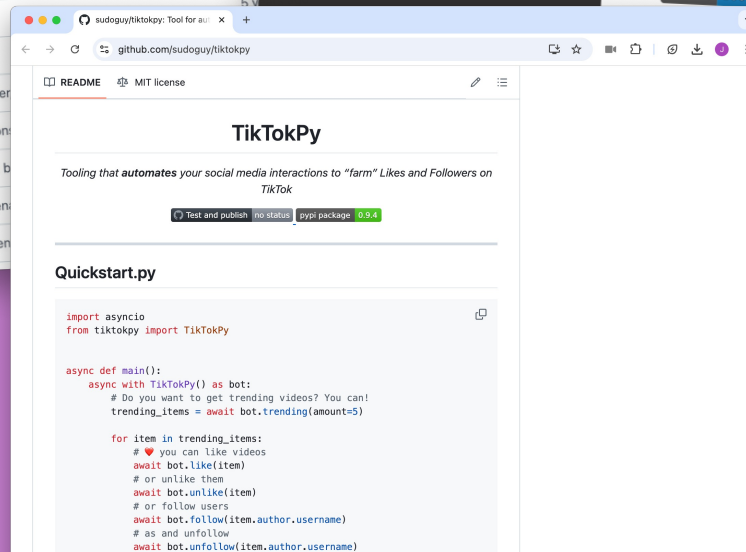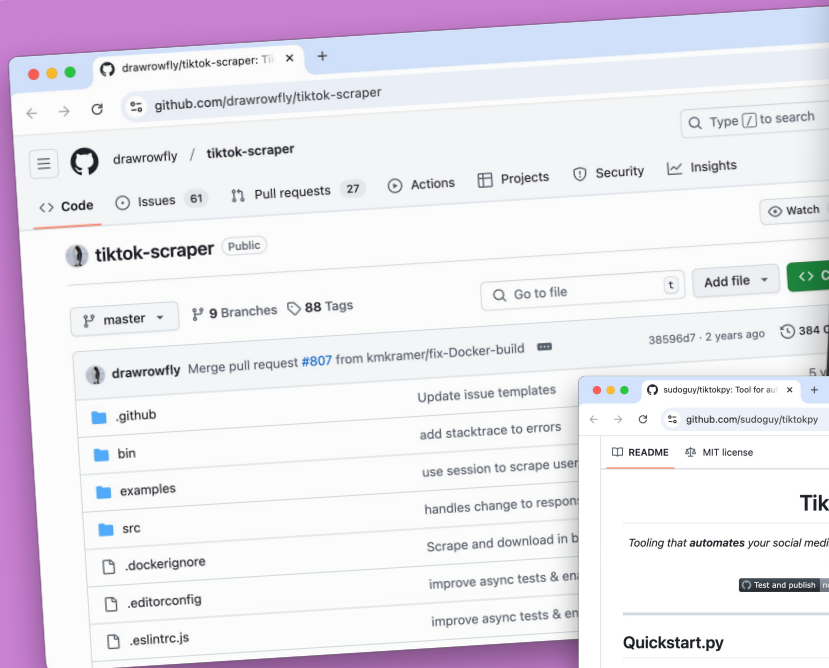    },
    "likes_count":

```python
import instaloader
import getpass

L = instaloader.Instaloader(
    download_pictures=True,
    download_videos=True,
    download_video_thumbnails=True,
    download_geotags=False,
    download_comments=True,
    save_metadata=True,
    compress_json=False
)

try:
    L.load_session_from_file('dangerscarf')
    print("Session loaded")
except:
    username = input("Enter your Instagram username: ")
    password = getpass.getpass("Enter your Instagram password: ")

    try:
        L.login(username, password)
        print("Login successful!")
    except instaloader.exceptions.TwoFactorAuthRequiredException:
        print("Two-factor authentication required.")
        code = input("Enter the 2FA code from your authenticator app or SMS: ")
        try:
            L.two_factor_login(code)
            print("Two-factor authentication successful!")
            L.save_session_to_file()
        except instaloader.exceptions.InstaloaderException as e:
            print(f"Two-factor authentication failed: {e}")
    except instaloader.exceptions.InstaloaderException as e:
        print(f"Login failed: {e}")
```

# solution.....?



**Browser — drawrowfly/tiktok-scraper**

github.com/drawrowfly/tiktok-scraper

drawrowfly / tiktok-scraper

Type [/] to search

<> Code · ⊙ Issues 61 · ⥊ Pull requests 27 · ▷ Actions · ▦ Projects · ⊘ Security · ⩘ Insights

tiktok-scraper  Public

⊙ Watch

master · ⎇ 9 Branches · ⊙ 88 Tags

Go to file · Add file · <> C

drawrowfly Merge pull request #807 from kmkramer/fix-Docker-build · 38596d7 · 2 years ago · 384

- .github — Update issue templates
- bin — add stacktrace to errors
- examples — use session to scrape user
- src — handles change to respon
- .dockerignore — Scrape and download in b
- .editorconfig — improve async tests & en
- .eslintrc.js — improve async tests & en

**Browser — TikTokAPI Quick Start**

davidteather.github.io/TikTok-Api/

🏠 TikTokAPI

Search docs

- ▣ TikTokAPI Quick Start
- TikTokAPI Full Documentation

🏠 / TikTokAPI Quick Start

View page source

## TikTokAPI Quick Start

View the above tree for detailed documentation on the package.

## Unofficial TikTok API in Python

This is an unofficial api wrapper for TikTok.com in python. With this api you are able to call most trending and fetch specific user information as well as much more.

DOI 10.5281/zenodo.6842316 · LinkedIn · ♡ Sponsor · release V7.0.0 · license MIT · downloads 2M · visitors 1661958 · 💬 chat 50 online

...signed to **retrieve data** TikTok. It **can not be used post or upload** content to ...e behalf of a user. It has **no support for any user-authenticated routes**, if you ...t while being logged out on their website you can't access it here.

...rs have paid to be placed here or are my own affiliate links which I may earn ...from, and beyond that I do not have any affiliation with them. The ...kage will always be free and open source. If you wish to be a sponsor of this

**Browser — sudoguy/tiktokpy**

github.com/sudoguy/tiktokpy

📖 README · ⚖ MIT license

## TikTokPy

Tooling that **automates** your social media interactions to "farm" Likes and Followers on TikTok

Test and publish no status · pypi package 0.9.4

### Quickstart.py

```python
import asyncio
from tiktokpy import TikTokPy

async def main():
    async with TikTokPy() as bot:
        # Do you want to get trending videos? You can!
        trending_items = await bot.trending(amount=5)

        for item in trending_items:
            # ❤ you can like videos
            await bot.like(item)
            # or unlike them
            await bot.unlike(item)
            # or follow users
            await bot.follow(item.author.username)
            # as and unfollow
            await bot.unfollow(item.author.username)
```

# solution: yt-dlp

```python
import yt_dlp

# List of YouTube URLs
urls = [
    "https://www.youtube.com/watch?v=N5wvtYRYbfA",
    "https://www.youtube.com/watch?v=5MHEMedLWeI",
    "https://www.youtube.com/watch?v=d6ipSIy7MEk",
    "https://www.youtube.com/watch?v=ZsVlvsxfnXw",
    "https://www.youtube.com/watch?v=b8e9YgvkKgc",
    "https://www.youtube.com/watch?v=Thc0vtnWYOo"
]

# Configure yt-dlp options
ydl_opts = {
    'format': 'bestvideo[height<=720]+bestaudio/best[height<=720]',
    'merge_output_format': 'mp4',
    'outtmpl': 'downloads/%(title)s.%(ext)s',
}

with yt_dlp.YoutubeDL(ydl_opts) as ydl:
    ydl.download(urls)
```

downloads

Name

AI, Hugging Face and non-chatbot models (Practical AI for Journa

Local models/private AI (Practical AI for Investigative Journalism, Ses

Sorting documents (Practical AI for Investigative Journalism, Se

Structured, validated data from LLMs (Practical AI for Investigati

Transcription and audio models (Practical AI for Investigative Jou

Why generative AI is a dead end for responsible journalism (Prac

# solution: writing a scraper



```python
from playwright.async_api import async_playwright
playwright = await async_playwright().start()
browser = await playwright.chromium.launch(headless=False)
page = await browser.new_page()
```

the problem with scrapers

```
<div id="product-list">
    <div class="product" id="product-1">
        <h2 class="product-title">Oak Wood Stain</h2>
        <p class="product-description">A rich, warm oak stain that
enhances the natural grain.</p>
        <span class="product-price">$19.99</span>
    </div>

    <div class="product" id="product-2">
        <h2 class="product-title">Walnut Wood Stain</
        <p class="product-description">A deep, luxuri
for a classic finish.</p>
        <span class="product-price">$29.99</span>
    </div>

    <div class="product" id="product-3">
        <h2 class="product-title">Mahogany Wood Stain<
        <p class="product-description">A rich mahogany s
brings out vibrant reddish tones.</p>
        <span class="product-price">$39.99</span>
    </div>
</div>
```

```
<section class="products">
    <article class="product" data-id="oak-stain">
        <header>
            <h2>Oak Wood Stain</h2>
        </header>
        <p>A rich, warm oak stain that enhances the natural grain.</p>
        <footer>
            <strong>$1
        </footer>
    </article>

    <article class="p
        <header>
            <h2>Walnut
        </header>
            lu
```

Drag the puzzle piece into place

Drag the puzzle piece into place

Select all squares with

**traffic lights**

If there are none, click skip

SKIP

# tiers of problem-solving

- ~~Using a tool~~
- ~~Writing a scraper~~
- Undocumented APIs
- Intercepting browser requests
- Pack-ratting with HAR and WARC/WACZ files

# solution: undocumented APIs

HANDS-ON   INTERMEDIATE

**Findings and using undocumented APIs**

**Time**: Friday, March 7, 3:30 – 4:30 p.m. (1h)
**Location**: Gray's Bay, 8th fl, eighth floor (BYO)

SHOW FEWER DETAILS

This tutorial will introduce reporters to an exciting and often overlooked data source found on every website. You will learn how to find and use hidden APIs as a reporting resource, and hear about how this data source has been used in past reporting. We'll be working off this scripted documented: https://inspectelement.org/apis

This session is for reporters who want to diversify their data sources. You don't need to write code: we'll teach participants to find hidden APIs in your web browser, but knowing some coding will let you to unlock detailed and rich datasets hidden in plain sight. Laptops will be provided.

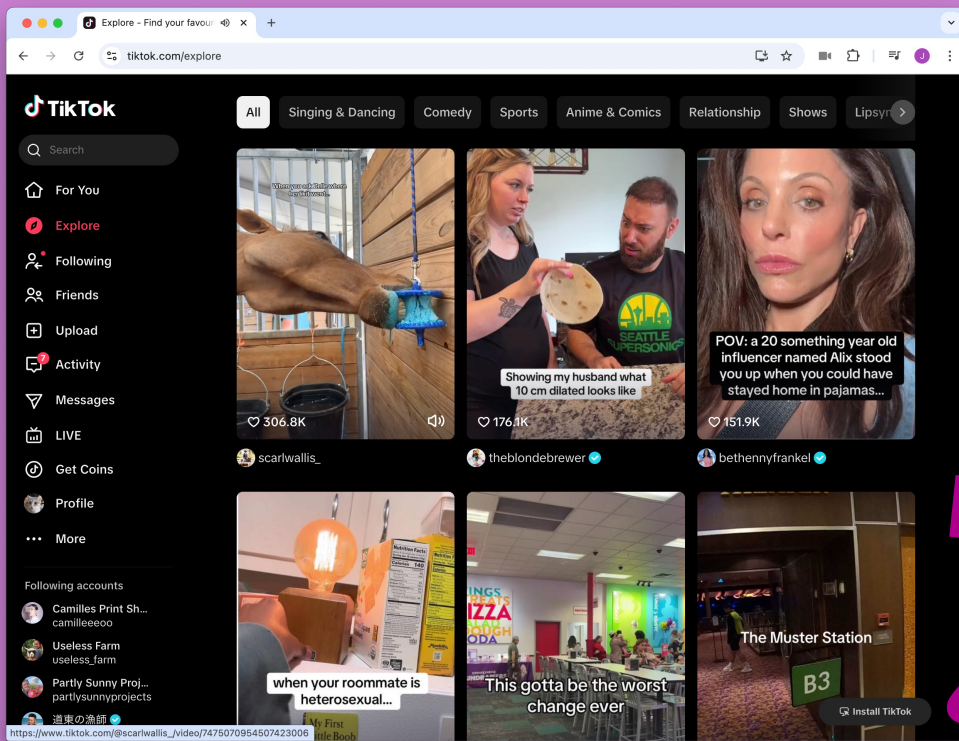# the problem: they change (and we're lazy!)

/products/list.json

```json
{
  "products": [
    {
      "id": "product-1",
      "title": "Oak Wood Stain",
      "description": "A rich, warm oak stain that enhances the natural grain.",
      "price": "$19.99"
    },
    {
      "id": "product-2",
      "title": "Walnut Wood Stain",
      "description": "A deep, luxurious walnut stain for a classic finish.",
      "price": "$29.99"
    },
    {
      "id": "product-3",
      "title": "Mahogany Wood Stain",
      "description": "A rich mahogany stain that brings out vibrant reddish tones.",
      "price": "$39.99"
    }
  ]
}
```

/api/v2/products

```json
{
  "products": [
    {
      "data_id": "oak-stain",
      "name": "Oak Wood Stain",
      "details": "A rich, warm oak stain that enhances the natural grain.",
      "price": 19.99,
      "currency": "USD"
    },
    {
      "data_id": "walnut-stain",
      "name": "Walnut Wood Stain",
      "details": "A deep, luxurious walnut stain for a classic finish.",
      "price": 29.99,
      "currency": "USD"
    },
    {
      "data_id": "mahogany-stain",
      "name": "Mahogany Wood Stain",
      "details": "A rich mahogany stain that brings out vibrant reddish tones.",
      "price": 39.99,
      "currency": "USD"
    }
  ]
}
```

solution: intercepting API calls

give me a thing!

us being sneaky

here you go!

```python
async def process_response(response):
    if response.ok and response.url.startswith("https://www.tikt
        try:
            parsed = urlparse(response.url)
            m = hashlib.md5()
            m.update(parsed.query.encode('utf-8'))
            m.update((await response.text()).encode('utf-8'))
            filename = m.hexdigest()

            end = parsed.path.lstrip("/")
            path = Path(end).joinpath(filename).with_suffix("
            path.parent.mkdir(parents=True, exist_ok=True)

            print("Writing content to ", path)

            query = parse_qs(parsed.query)
            query = {k: v[0] for k, v in query.items()}

            content = {
                'url': response.url,
                'query': query,
                'headers': await response.all_headers(),
                'data': await response.json(),
            }

            path.write_text(json.dumps(content, indent=2
        except Exception as e:
            print("Error processing response", e)

page.on("response", process_response)
```
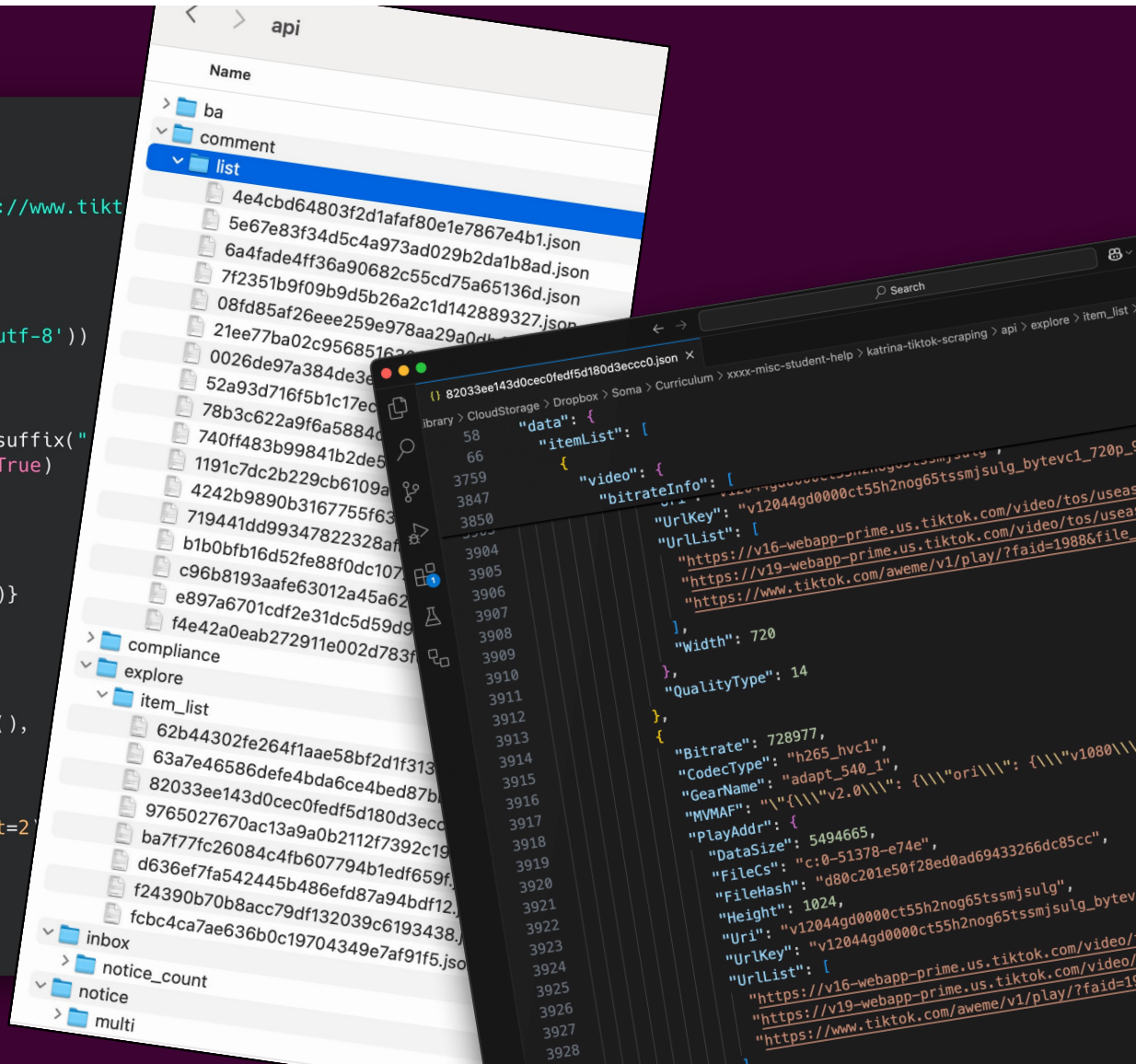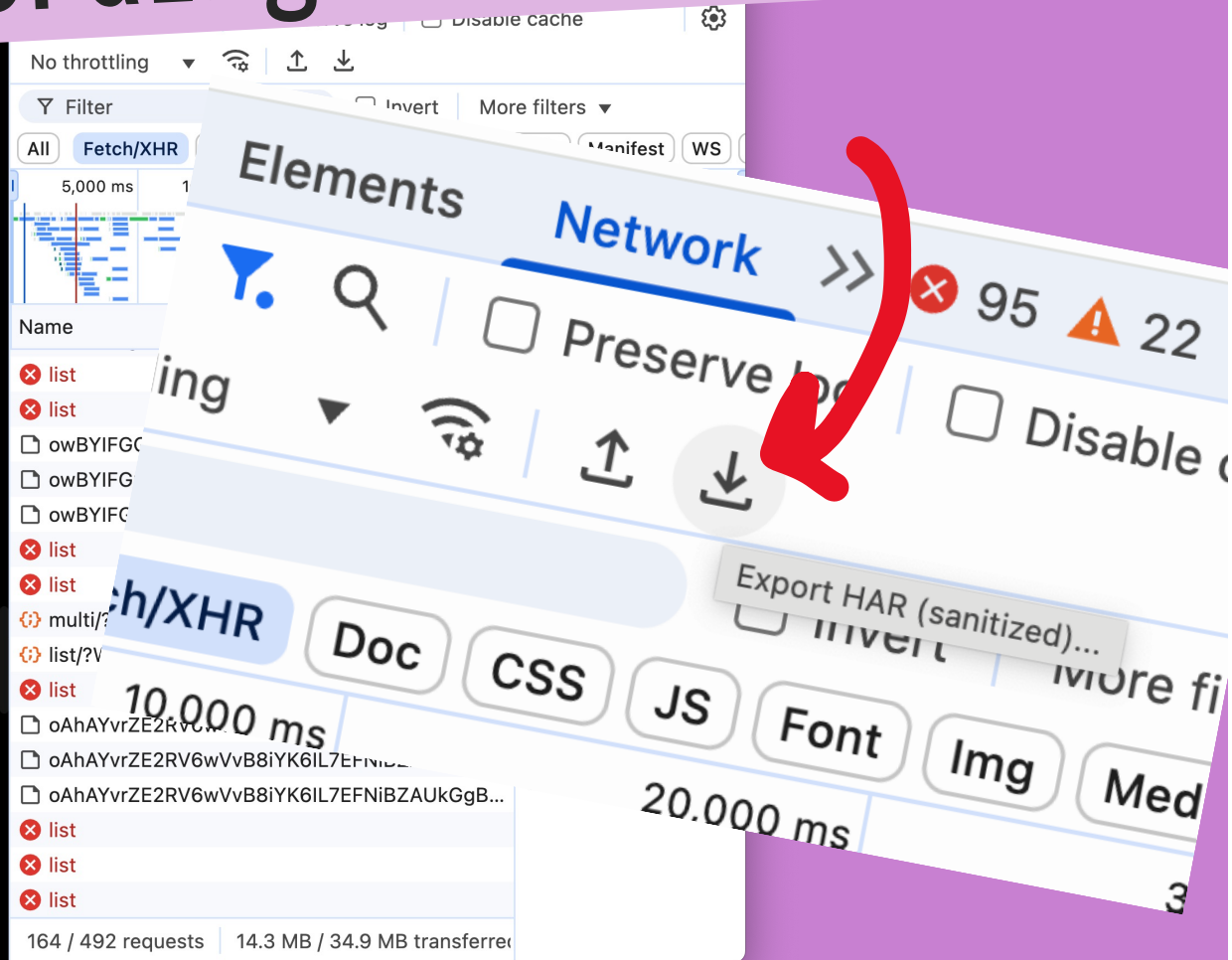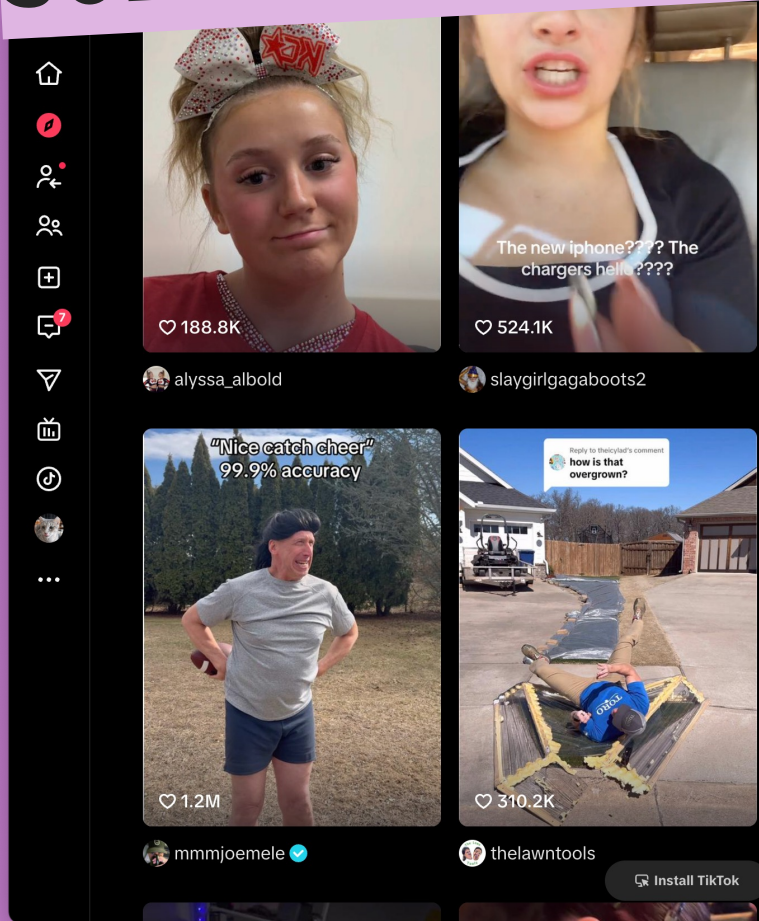
api

Name

- ba
- comment
  - list
    - 4e4cbd64803f2d1afaf80e1e7867e4b1.json
    - 5e67e83f34d5c4a973ad029b2da1b8ad.json
    - 6a4fade4ff36a90682c55cd75a65136d.json
    - 7f2351b9f09b9d5b26a2c1d142889327.json
    - 08fd85af26eee259e978aa29a0d
    - 21ee77ba02c956851620
    - 0026de97a384de3e
    - 52a93d716f5b1c17ec
    - 78b3c622a9f6a5884d
    - 740ff483b99841b2de5
    - 1191c7dc2b229cb6109a
    - 4242b9890b3167755f63
    - 719441dd99347822328af
    - b1b0bfb16d52fe88f0dc107
    - c96b8193aafe63012a45a62
    - e897a6701cdf2e31dc5d59d9
    - f4e42a0eab272911e002d783f
- compliance
- explore
  - item_list
    - 62b44302fe264f1aae58bf2d1f313
    - 63a7e46586defe4bda6ce4bed87b
    - 82033ee143d0cec0fedf5d180d3ecc
    - 9765027670ac13a9a0b2112f7392c19
    - ba7f77fc26084c4fb607794b1edf659f
    - d636ef7fa542445b486efd87a94bdf12.
    - f24390b70b8acc79df132039c6193438.
    - fcbc4ca7ae636b0c19704349e7af91f5.json
- inbox
  - notice_count
- notice
  - multi

82033ee143d0cec0fedf5d180d3eccc0.json
Library > CloudStorage > Dropbox > Soma > Curriculum > xxxx-misc-student-help > katrina-tiktok-scraping > api > explore > item_list

```json
58        "data": {
66          "itemList": [
3759          {
3847            "video": {
3850              "bitrateInfo": [
                    ...gd0000ct55h2nog65tssmjsulg_bytevc1_720p_9
                  "UrlKey": "v12044gd0000ct55h2nog65tssmjsulg/video/tos/useas
3904          "UrlList": [
3905            "https://v16-webapp-prime.us.tiktok.com/video/tos/useas
3906            "https://v19-webapp-prime.us.tiktok.com/?faid=1988&file
3907            "https://www.tiktok.com/aweme/v1/play/?faid=1988&file
3908          ],
3909          "Width": 720
3910        },
3911        "QualityType": 14
3912      },
3913      {
3914        "Bitrate": 728977,
3915        "CodecType": "h265_hvc1",
3916        "GearName": "adapt_540_1",
3917        "MVMAF": "\"{\\\"v2.0\\\": {\\\"ori\\\": {\\\"v1080\\\"
3918        "PlayAddr": {
3919          "DataSize": 5494665,
3920          "FileCs": "c:0-51378-e74e",
3921          "FileHash": "d80c201e50f28ed0ad69433266dc85cc",
3922          "Height": 1024,
3923          "Uri": "v12044gd0000ct55h2nog65tssmjsulg_bytevc
3924          "UrlKey": "v12044gd0000ct55h2nog65tssmjsulg/video/
3925          "UrlList": [
3926            "https://v16-webapp-prime.us.tiktok.com/video/
3927            "https://v19-webapp-prime.us.tiktok.com/aweme/v1/play/?faid=19
3928            "https://www.tiktok.com/aweme/v1/play/?faid=19
```

# the problem: they change (and we're lazy!)

**/products/list.json**

```json
{
  "products": [
    {
      "id": "product-1",
      "title": "Oak Wood Stain",
      "description": "A rich, warm oak stain that enhances the natural grain.",
      "price": "$19.99"
    },
    {
      "id": "product-2",
      "title": "Walnut Wood Stain",
      "description": "A deep, luxurious walnut stain for a classic finish.",
      "price": "$29.99"
    },
    {
      "id": "product-3",
      "title": "Mahogany Wood Stain",
      "description": "A rich mahogany stain that brings out vibrant reddish tones.",
      "price": "$39.99"
    }
  ]
}
```

**/api/v2/products**

```json
{
  "products": [
    {
      "data_id": "oak-stain",
      "name": "Oak Wood Stain",
      "details": "A rich, warm oak stain that enhances the natural grain.",
      "price": 19.99,
      "currency": "USD"
    },
    {
      "data_id": "walnut-stain",
      "name": "Walnut Wood Stain",
      "details": "A deep, luxurious walnut stain for a classic finish.",
      "price": 29.99,
      "currency": "USD"
    },
    {
      "data_id": "mahogany-stain",
      "name": "Mahogany Wood Stain",
      "details": "A rich mahogany stain that brings out vibrant reddish tones.",
      "price": 39.99,
      "currency": "USD"
    }
  ]
}
```

# tiers of problem-solving

- ~~Using a tool~~
- ~~Writing a scraper~~
- ~~Undocumented APIs~~
- ~~Intercepting browser requests~~
- Pack-ratting with HAR and WARC/WACZ files

solution: recording with HAR

www.tiktok.com.har

```
                                    }
                                }
                            }
                        },
                        "_priority": "High",
                        "_resourceType": "fetch",
                        "cache": {},
                        "connection": "443",
                        "pageref": "page_1",
                        "request": {
                            "method": "GET",
                            "url": "https://www.tiktok.com/api/recommend/item_list/?WebIdLastTime=0&aid=1988&app_language=en&ap
                            "httpVersion": "http/2.0",
                            "headers": [
                                {
                                    "name": ":authority",
                                    "value": "www.tiktok.com"
                                },
                                {
                                    "name": ":method",
                                    "value": "GET"
                                },
                                {
                                    "name": ":path",
                                    "value": "/api/recommend/item_list/?WebIdLastTime=0&aid=1988&app_language=en&app_name=tiktok_we
                                },
                                {
                                    "name": ":scheme",
                                    "value": "https"
```

```python
from haralyzer import HarParser
import json
import base64
import hashlib


def make_filepath(response):
    parsed = urlparse(response.url)

    # Hash the query parameters and the response.text to make a uni
    m = hashlib.md5()
    m.update(json.dumps(entry.request.raw_entry).encode('utf-8'))
    m.update(json.dumps(entry.response.raw_entry).encode('utf-8'))
    hash = m.hexdigest()

    # the full path is based on the URL
    path = Path(parsed.netloc).joinpath(parsed.p
    if len(path.name) > 200:
        path = path.with_name(path.name[-50:])
    path = path.with_suffix(f".{hash}.jsor"`

    return path

def decode_response(response):
    if response.textEncoding == 'base64':
        content = base64.b64decode(response.text).decode('utf-8')
    else:
        content = response.text

    if isinstance(content, str):
        try:
            # If it's a JSON string, parse it to an object
            content = json.loads(content)
        except json.JSONDecodeError:
            pass  # Keep as string if not valid JSON
    return content
```

```python
filename = "www.instagram.com.har"
output_folder = Path("output").joinpath(filename)

har = HarParser.from_file(filename)

for page in har.pages:
    entries = page.filter_entries(status_code='200', content_type='application/json
    for entry in entries:
        content = decode_response(response)

        # Find out where it should go
        filepath = make_filepath(response)
        output_folder.joinpath(filepath)
```

```python
import jmespath
import pandas as pd

result = jmespath.search("sectional_items[*].layout_content.fill_items[]", content)
len(result)
```

```python
pd.options.display.max_columns = None

df = pd.json_normalize(result)
df
```

```python
    entries = page.filter_entries(status_code='200', content_type='application/
    for entry in entries:
        try:
            content = decode_response(entry.response)
            if content is None:
                continue

            filepath = make_filepath(entry.response, guess_extension=False)
            path = output_folder.joinpath(filepath)

            # Save based on content type
            path.parent.mkdir(parents=True, exist_ok=True)
            if isinstance(content, (dict, list)):
                path.write_text(json.dumps(content, indent=2))
            elif isinstance(content, str):
                path.write_text(content)
            elif isinstance(content, bytes):
                path.write_bytes(content)

            print(f"Writing {path}")
        except Exception as e:
```

# HAR Data Extractor

Transform JSON API responses from HAR files into CSV format

Drag and drop HAR file here

By your buddy Jonathan Soma. See the code at jsoma/har2data

# the problem

- Giant files!!!
- Chunking!!!
- Base64 encoding!!!
- Other things I probably don't even know about yet!!!
- (But they do also sometimes work, too, it just depends)

# solution: saving WARC files

ty jeremy merrill

# Web Archive JSON Extractor

Extract and analyze API requests from web archives (WARC/WACZ files). This tool finds JSON API requests made by web pages and allows you to explore and export the data.

Drag and drop WARC/WACZ files here, or click to select files

## API Requests

Filter by path...

Sort by: Default

Select visible

No API requests found. Upload a WARC/WACZ file to begin.

## JSON Content (Click on a field to select its path)

Select a JSON file to view its content

JQ Path: Enter JQ path (e.g., .data.items[0].name) ?
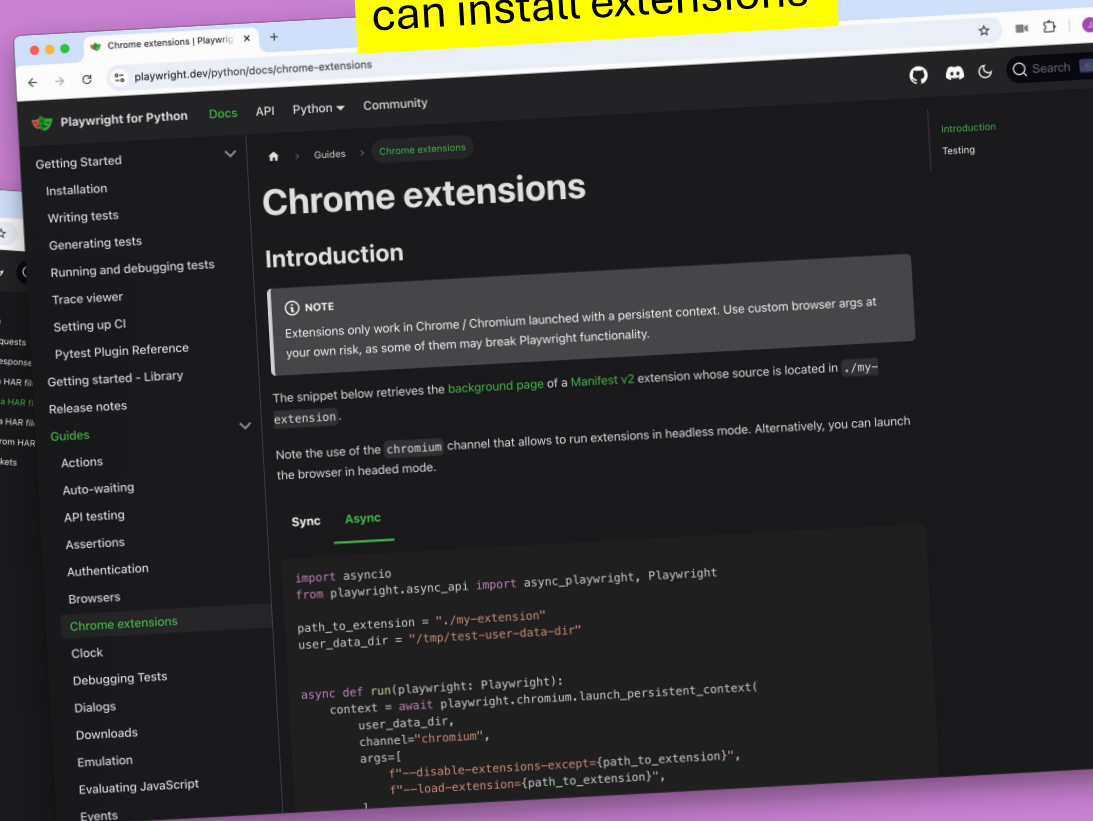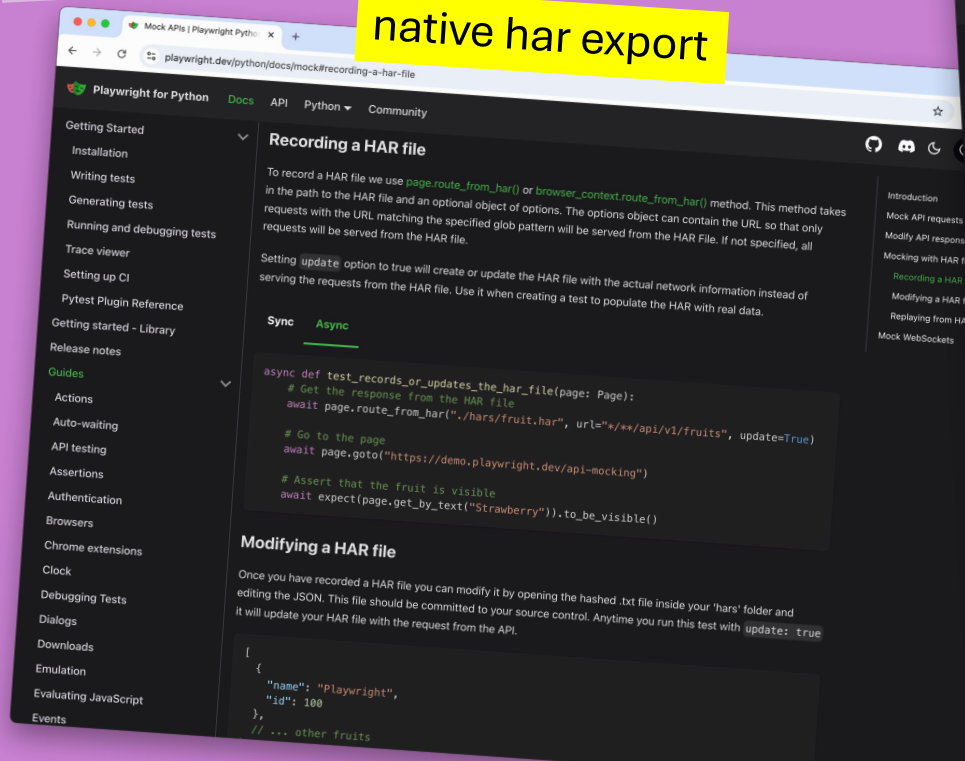
# the problem: automation..?



```python
from playwright.async_api import async_playwright
playwright = await async_playwright().start()
browser = await playwright.chromium.launch(headless=False)
page = await browser.new_page()
```

# just kidding! Automate with Playwright!

**can install extensions**

**native har export**



## Recording a HAR file

To record a HAR file we use `page.route_from_har()` or `browser_context.route_from_har()` method. This method takes in the path to the HAR file and an optional object of options. The options object can contain the URL so that only requests with the URL matching the specified glob pattern will be served from the HAR File. If not specified, all requests will be served from the HAR file.

Setting `update` option to true will create or update the HAR file with the actual network information instead of serving the requests from the HAR file. Use it when creating a test to populate the HAR with real data.

**Sync** **Async**

```
async def test_records_or_updates_the_har_file(page: Page):
    # Get the response from the HAR file
    await page.route_from_har("./hars/fruit.har", url="*/**/api/v1/fruits", update=True)

    # Go to the page
    await page.goto("https://demo.playwright.dev/api-mocking")

    # Assert that the fruit is visible
    await expect(page.get_by_text("Strawberry")).to_be_visible()
```

## Modifying a HAR file

Once you have recorded a HAR file you can modify it by opening the hashed .txt file inside your 'hars' folder and editing the JSON. This file should be committed to your source control. Anytime you run this test with `update: true` it will update your HAR file with the request from the API.

```
[
    {
        "name": "Playwright",
        "id": 100
    },
    // ... other fruits
```

## Chrome extensions

### Introduction

> ⓘ **NOTE**
>
> Extensions only work in Chrome / Chromium launched with a persistent context. Use custom browser args at your own risk, as some of them may break Playwright functionality.

The snippet below retrieves the background page of a Manifest v2 extension whose source is located in `./my-extension`.

Note the use of the `chromium` channel that allows to run extensions in headless mode. Alternatively, you can launch the browser in headed mode.

**Sync** **Async**

```
import asyncio
from playwright.async_api import async_playwright, Playwright

path_to_extension = "./my-extension"
user_data_dir = "/tmp/test-user-data-dir"

async def run(playwright: Playwright):
    context = await playwright.chromium.launch_persistent_context(
        user_data_dir,
        channel="chromium",
        args=[
            f"--disable-extensions-except={path_to_extension}",
            f"--load-extension={path_to_extension}",
```

https://bit.ly/nicar-passive

# Passive Scraping for social media

(and everything else)

Jonathan Soma
Columbia Journalism School
js4571@columbia.edu
@dangerscarf