



# Introduction To Web Scrapping



With David Teather



# Who am I? 🤔

- CS & SWE
- A general clown
- Uses emojis excessively





# TikTokAPI

- Programmatically extract TikTok data from Python
- 800k+ downloads
- 3k+ stars

**About**

The Unofficial TikTok API Wrapper In Python

[davidteather.github.io/TikTok-API/do...](https://github.com/davidteather/tiktok-api)

python api trending hacktoberfest

tiktok tik tok tiktok-scraper

tiktok-api tiktokapi tiktok-python

tiktok-automation tiktok-downloader

tiktok-signature tiktok-trending-page

tiktok-compilations download-tiktoks

Readme

MIT, MIT licenses found

Cite this repository

3k stars

74 watching

723 forks

**Releases** 108

**V5.2.2** Latest  
on Jul 15, 2022

+ 107 releases





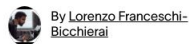
# YikYak “Security Analysis”

- YikYak was exposing GPS coords of posts within 10-15 feet
- [TheResponseTimes.com](https://www.theresponse.com)

**MOTHERBOARD**  
TECH BY VICE

## Anonymous Social Media App Yik Yak Exposed Users' Precise Locations

Privacy researchers have found that it's possible to find out the precise location of Yik Yak posts, potentially exposing users to doxing or stalking.



By Lorenzo Franceschi-  
Bicchieri

May 12, 2022, 11:10am [Share](#) [Tweet](#) [Snap](#)





# GitHub Campus Expert

- Please take more stickers 🥲





# Why should you learn web scraping?





# Why you should learn

- Data is (almost) priceless
- Improves decision making
- Not all companies have their partner's data





# Record Labels

Imagine you're a record label and you want to find artists that will be successful in the future.

Q: What data would be valuable to have?







# Record Labels

- TikTok
- Spotify
- YouTube
- Twitter
- Instagram
- Etc

Especially smaller companies might not have deals with these companies for their artists data





# Record Labels

(Before 2018)

- Hardly any data utilized in music
- Talent relation agents went to multiple concerts a day to scout
- Listened to hundreds of irrelevant hours music





# Record Labels

(After 2018)

- Still listen to music 24/7, but more strategic in choosing what to listen to
- Better predictions of who will be the next big star
- Critics claim the numbers ruin the art in music
- All major labels utilize data heavily



WARNER MUSIC GROUP

Sodatone.





# Tracking.Exposed



**Tracking  
Exposed**



- Trying to research algorithms as a 3rd party
  - Maybe for independent auditing
- Usually don't have research APIs available
  - And some accidentally mislead researchers
- Need to collect tons of data yourself





# Trendpop/Collab

- Collab
  - Signing content creators
- Trendpop
  - Analyzing social media primarily for marketers
    - Talent & Media Agencies
    - Music Streaming Platforms
    - Content Creators

collab

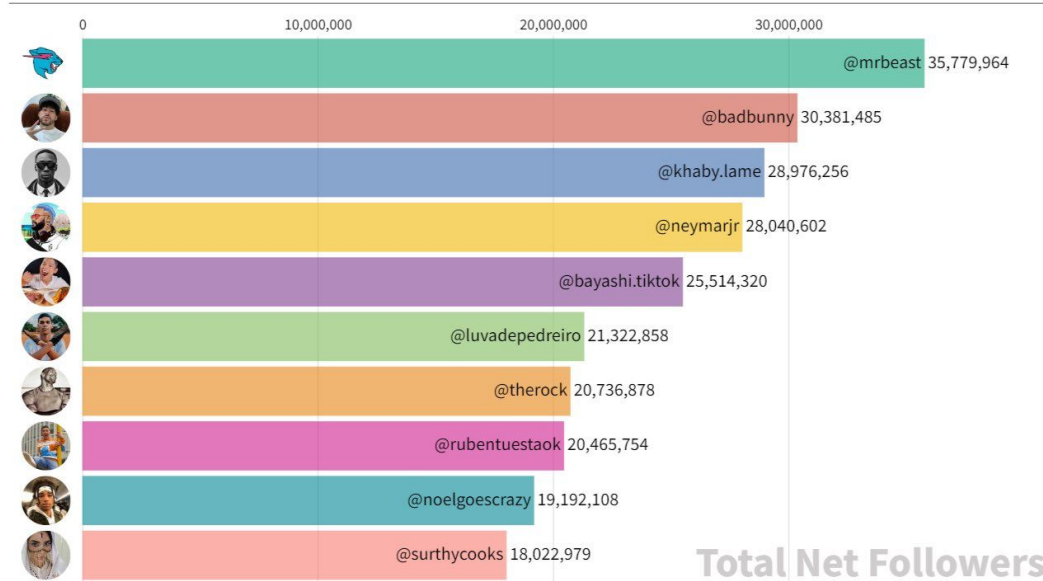
Trendpop





# Trendpop/Collab

## Fastest Growing TikTok Creators Worldwide (2022)



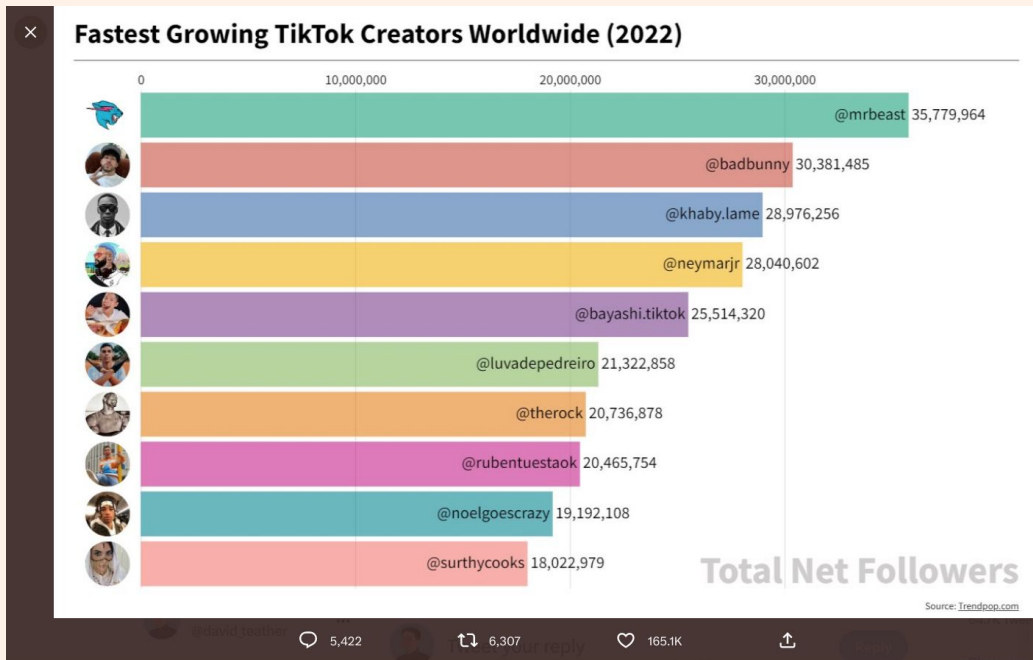
Total Net Followers

Source: Trendpop.com





# Trendpop/Collab





# Possibilities are Endless

(Currently)

- Predominantly hyper-competitive industries
  - Social Media
  - Investment Firms
  - More “Tech” Companies







# Possibilities are Endless

(Future)

- Everyone utilizing data
- Managers in companies will leverage highly accurate data
  - Even data that they might not necessarily own
  - Leverage web scraped data for competing with others
  - etc





# Let's get started!





# Websites Send Files



- Adds the Structure
- The “Base” language





# Websites Send Files



- Adds the Structure
- The “Base” language



- Adds styling
- What makes a website appealing





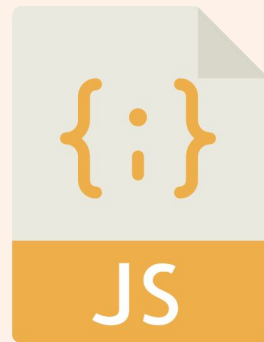
# Websites Send Files



- Adds the Structure
- The “Base” language



- Adds styling
- What makes a website appealing



- Interactive
- “Responsive”





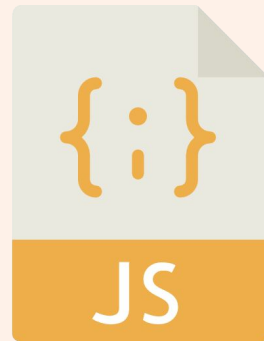
# Websites Send Files



- Adds the Structure
- The “Base” language



- Adds styling
- What makes a website appealing



- Interactive
- “Responsive”





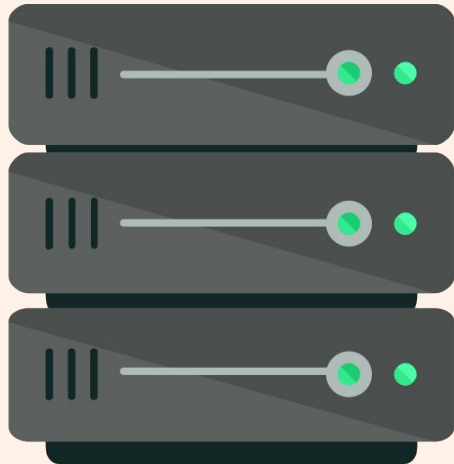
# How Websites Send Data

1. Static Site/Server Side Rendering (SSR)
2. Dynamic / AJAX





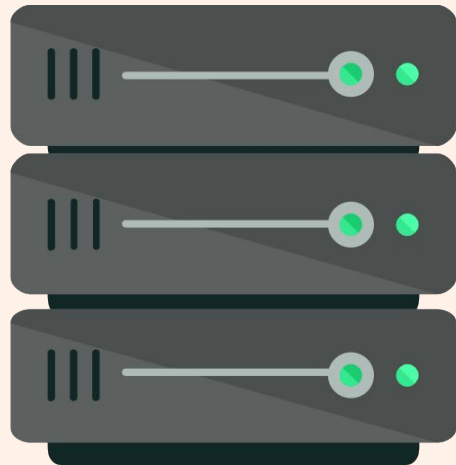
# Static Site / SSR



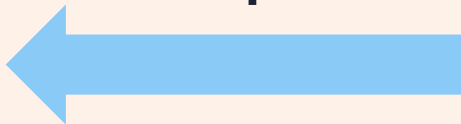




# Static Site / SSR

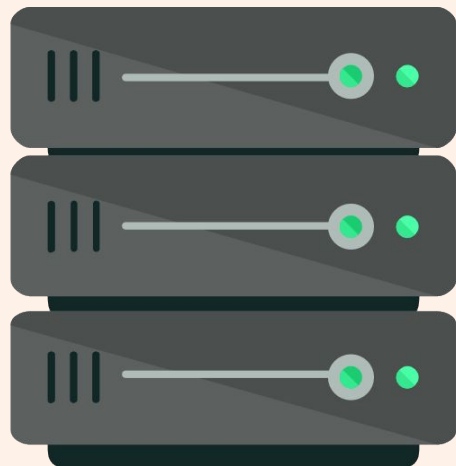


GET example.com

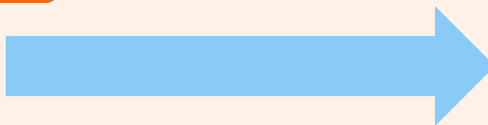
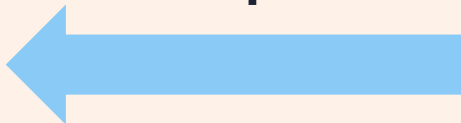




# Static Site / SSR

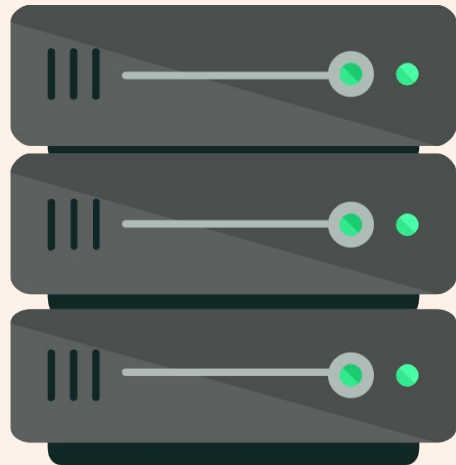


GET example.com

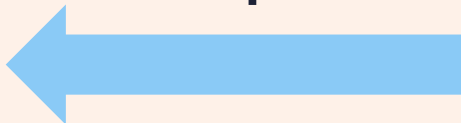




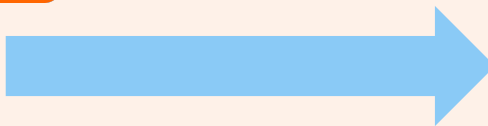
# Static Site / SSR



GET example.com

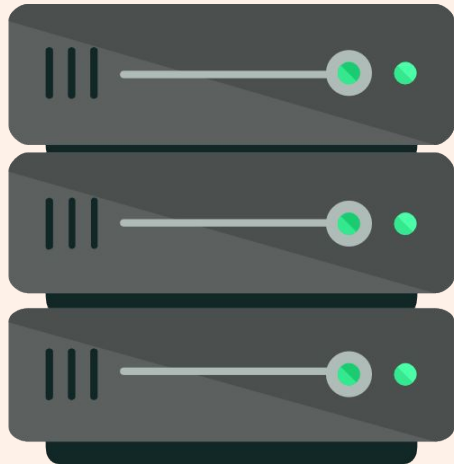


Directly has data



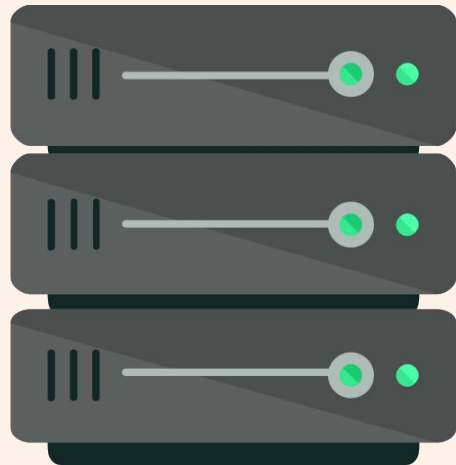


# Dynamic / AJAX





# Dynamic / AJAX

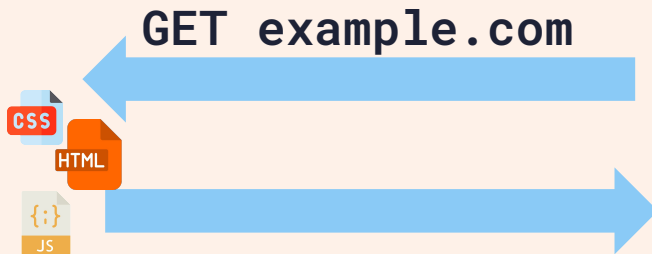
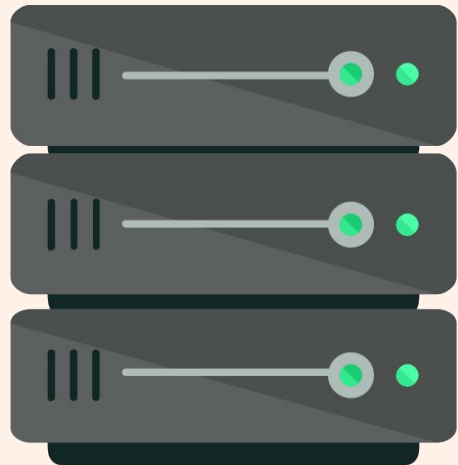


GET example.com



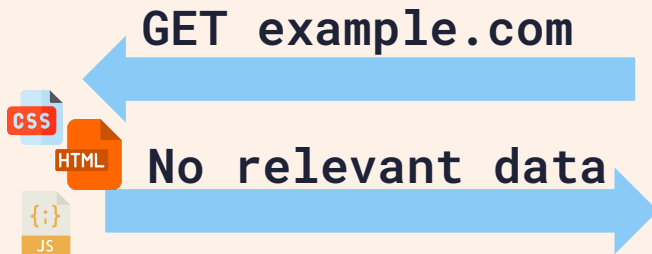
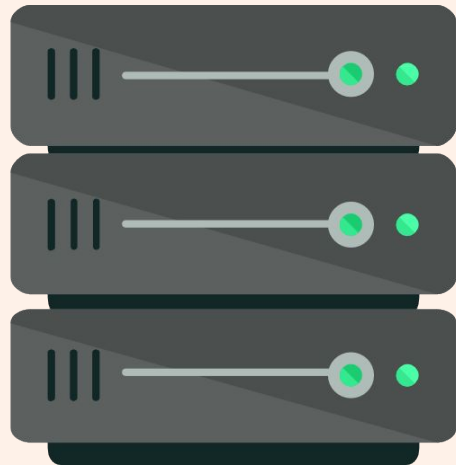


# Dynamic / AJAX



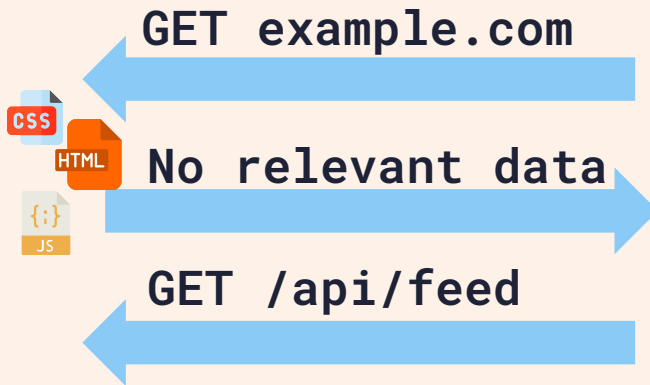
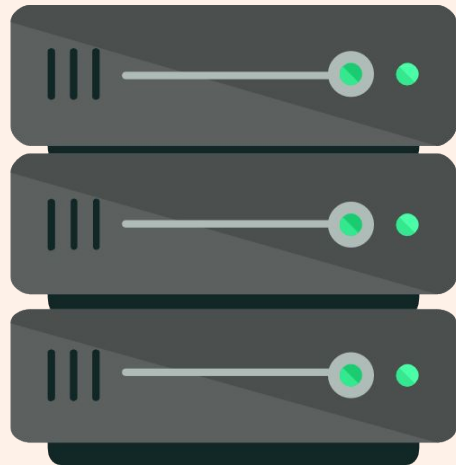


# Dynamic / AJAX





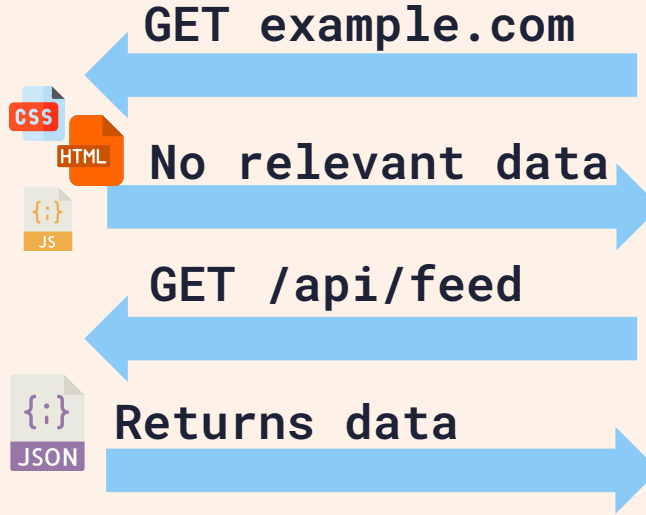
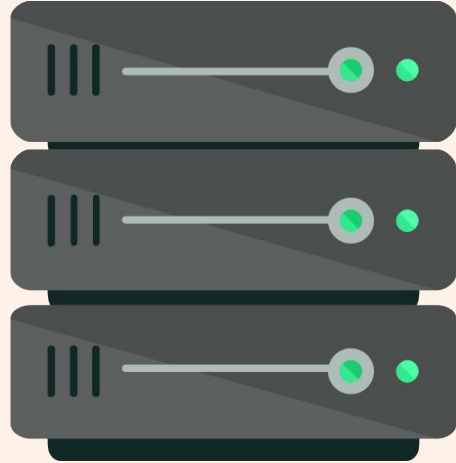
# Dynamic / AJAX







# Dynamic / AJAX





## Some Websites Will Do Both

1. Initially send some data in the HTML
  - For faster load times
2. Then use AJAX
  - For responsive feel

Ex: TikTok, YouTube, etc





# Workshop's Focus

- Extracting data from the HTML
- Dynamic / AJAX sites not covered here
  - [Everything Web Scraping](#)
  - Additional resources at end of slides
  - Talk to me :)



**DEMO TIME!**



[bit.ly/web scraping-madhacks](https://bit.ly/web scraping-madhacks)





# **BONUS:**

## **Advanced Techniques**





# Websites Blocking Requests

- Websites don't want us to web scrape 🤖
- Will block requests that look like a bot
  - But can't block real users or they can't access the site
- Goal: Look identical to real users





# How Bots are Identified

- User Agents
- IPs
- Headers
- Cookies
- Behavior
- Browser Fingerprinting
- Much more





# User Agents

- Identify what is making the request

Browser: Mozilla/5.0 (X11; Linux  
x86\_64) AppleWebKit/537.36 (KHTML,  
like Gecko) Chrome/110.0.0.0  
Safari/537.36

Python: python-requests/2.25.1







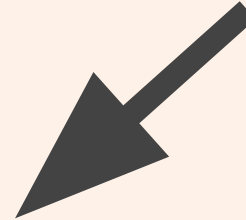
# User Agents

- Identify what is making the request

Browser: Mozilla/5.0 (X11; Linux x86\_64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/110.0.0.0 Safari/537.36

Python: python-requests/2.25.1

Easy to detect  
python requests  
module making  
requests





# User Agents

- Easy to change user-agent

```
requests.get(url, headers={'user-agent':  
'whatever'})
```





# IPs

- Rate limiting
  - Only allow X requests per minute per IP
- Bypassed with proxies
- [Everything Web Scraping Lesson](#)





# Headers

- Metadata passed on each request
- Can identify browser
- User agent is a header





# Cookies



- Another type of header
- Can store session data
  - Random marketing IDs
  - “Your account”
  - Also behavior





# Behavior

- Less applicable to static site web scraping
- Websites could tie your IP/cookies to your history on the site
- Ex: If you're jumping around to pages not accessible from the page you're on could look suspicious





# Browser Fingerprinting

- Uses a combination of all browser attributes to uniquely identify users
- Again, your goal is appear as average as possible
- Cool/Scary Resource: [AmIUnique.org](https://amiunique.org)
  - Even if you don't give it any explicit permissions on your browser usually uniquely identifies me :(





# Browser Automation

- Programmatically control a browser
- Can emulate human behavior
- Can still extract HTML or scrape through that
- Good solution for dynamically generated sites
  - Although [Forging API Requests](#) has some advantages over this







## Connect With Me!



David Teather

[LinkedIn](#)

[GitHub](#)

## Feedback



[bit.ly/david-workshop-feedback](https://bit.ly/david-workshop-feedback)  
(also on the GitHub Repo)

please do this it helps  
me improve :)

