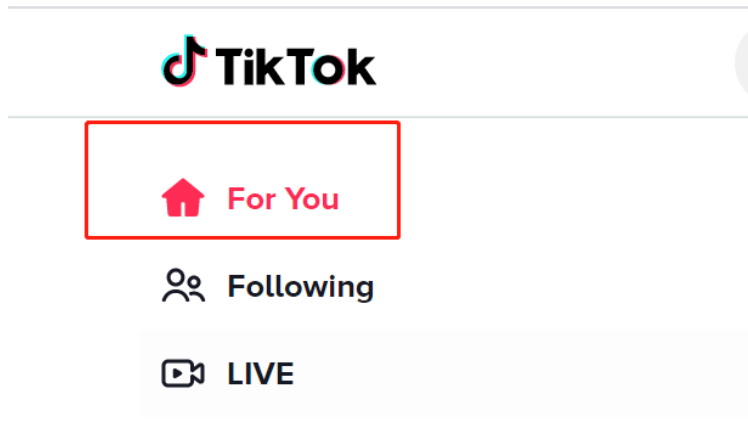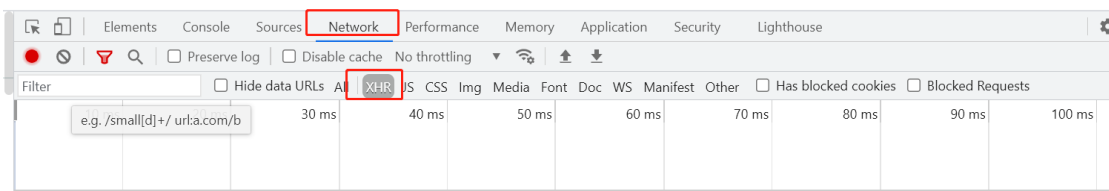**This is the design document for this crawler. I will record the process of how I designed and coded it.**

1. Open TikTok in your browser (Chrome prefer) and go to the **For You** page.



2. Go to Inspect and open the Network tab on top, then click on the XHR.



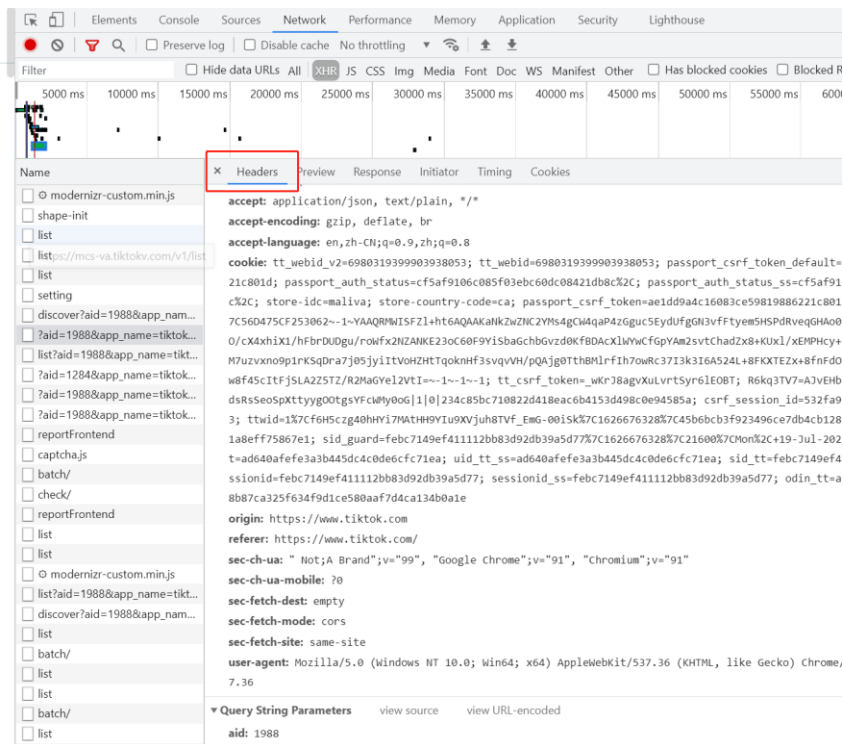3. Refresh the page to reload. Find the API for the video list. The API is called itemList.

4. However, there are only under 30 videos in this list so what we need to do is to scroll down to the next page. To do this, we need to find the variable related to the page number. This variable usually would be saved in the headers.

What we need to do now is to scroll your browser to go to the next page and refresh the browser to see which variable is changed. After doing so, I realized that the only variable changing is the one at the end called _signature

Name

- ⚙ modernizr-custom.min.js
- shape-init
- list
- list
- setting
- list
- discover?aid=1988&app_nam...
- ?aid=1988&app_name=tiktok...
- ?aid=1988&app_name=tiktok...
- list?aid=1988&app_name=tikt...
- ?aid=1284&app_name=tiktok...
- ?aid=1988&app_name=tiktok...
- reportFrontend
- captcha.js
- batch/
- check/
- list
- reportFrontend
- list
- list
- batch/
- list
- batch/
- list
- batch/
- ?aid=1988&app_name=tiktok...
- batch/

27 / 109 requests  63.5 kB / 4.2 M

✕  **Headers**  Preview  Response  Initiator  Timing  Cookies

O/cX4xhiX1/hFbrDUDgu/roWfx2NZANKE23oC60F9YiSbaGchbGvzd0KfBDAcXlWYwCfGpYAm2svtChadZx8+KUxl/xEMPHcy+2UX1GtosEC+duz3uSOFNizCr5
M7uzvxno9p1rKSqDra7j05jyiItVoHZHtTqoknHf3svqvVH/pQAjg0TthBMlrfIh7owRc37I3k3I6A524L+8FKXTEZx+8fnFdOYzY3930YACI3wIO0OtbsCNqCo
w8f45cItFjSLA2Z5TZ/R2MaGYel2VtI=~-1~-1~-1; tt_csrf_token=_wKrJ8agvXuLvrtSyr6lEOBT; R6kq3TV7=AJvEHb16AQAABOEVdMqeoK5z2HQTCzi
dsRsSeoSpXttyyg0OtgsYFcWMy0oG|1|0|234c85bc710822d418eac6b4153d498c0e94585a; csrf_session_id=532fa915a2c0444db62e223ccf56bd5
3; sid_guard=febc7149ef411112bb83d92db39a5d77%7C1626676328%7C21600%7CMon%2C+19-Jul-2021+12%3A32%3A08+GMT; uid_tt=ad640afefe
3a3b445dc4c0de6cfc71ea; uid_tt_ss=ad640afefe3a3b445dc4c0de6cfc71ea; sid_tt=febc7149ef411112bb83d92db39a5d77; sessionid=febc
7149ef411112bb83d92db39a5d77; sessionid_ss=febc7149ef411112bb83d92db39a5d77; odin_tt=adb8cb6e5c6e4da2eab7797a7a8b87ca325f63
4f9d1ce580aaf7d4ca134b0a1e; ttwid=1%7Cf6H5czg40hHYi7MAtHH9YIu9XVjuh8TVf_EmG-00iSk%7C1626676582%7C9219518c3d9244f4c75cc0bd12
3e55ea3e5db303ddfa83ff2c2f852ec1ce4d1a

**origin:** https://www.tiktok.com

**referer:** https://www.tiktok.com/

**sec-ch-ua:** " Not;A Brand";v="99", "Google Chrome";v="91", "Chromium";v="91"

**sec-ch-ua-mobile:** ?0

**sec-fetch-dest:** empty

**sec-fetch-mode:** cors

**sec-fetch-site:** same-site

**user-agent:** Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.124 Safari/537.36

▼ **Query String Parameters**   view source   view URL-encoded

**aid:** 1988

**app_name:** tiktok_web

**device_platform:** web_pc

**device_id:** 6980319399903938053

**region:** CA

**priority_region:**

**os:** windows

**referer:** https://www.tiktok.com/foryou?lang=en&is_copy_url=1&is_from_webapp=v1

**root_referer:** https://www.tiktok.com/logout?redirect_url=https%3A%2F%2Fwww.tiktok.com%2Fforyou%3Flang%3Den%26is_copy_url%3D1%26is_from_webapp%3Dv1&lang=en

**cookie_enabled:** true

**screen_width:** 2560

**screen_height:** 1441

**browser_language:** en

**browser_platform:** Win32

**browser_name:** Mozilla

**browser_version:** 5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91.0.4472.124 Safari/537.36

**browser_online:** true

**verifyFp:** verify_kra5lptq_pfANqQUj_7BJd_41yQ_95Mc_QBzCDnFQqOYw

**app_language:** en

**timezone_name:** America/New_York

**is_page_visible:** true

**focus_state:** true

**is_fullscreen:** false

**history_len:** 8

**battery_info:** 1

**noUser:** 0

**userCount:** 28

**scene:** 17

**from_page:** fyp

**_signature:** _02B4Z6wo00f01M9Q9aAAAIDBGJbY26J0xTzPUPEAAFMmc5

Obviously, this is an encrypted variable. I found the encrypt method under setting->verify->js->sg

Then we could open that file to see how it's encrypted.

Uhm never mind, this is absolutely out of my ability to decrypted it. Therefore, I choose to find some other way to fetch the data.

5. To do so, I first opened an URL with valid video list. (Notice, you can find it under the one has the itemList, just find that one and click on headers)



It should start with https://m.tiktok.com/api/recommend/item_list/…..

6. Then I tried to open this URL in my browser, and I noticed that videos I am getting from this same _signature variable are different. Which means I do not need to decrypt that variable to get different videos. However, this could cause a problem – I would probably get repeat videos since I am getting them all randomly. (I would discuss this problem later)

7. Try                            to                     request              this                    URL

```python
import requests
import re
import csv

url = 'https://m.tiktok.com/api/recommend/item_list/?aid=1988&app_name=tiktok_web&device_platform=web_pc&device_i

headers = {
    'cookie':'tt_webid_v2=6980332468541769218; tt_webid=6980332468541769218; _abck=ED0A47101D299D43655EA43F58
    'user-agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/91
}

response = requests.get(url = url, headers = headers).json()['itemList']
for i in response:
    print(i)
```

As we can see here, this URL is valid, and we could get data from it.

8. Then we just need to extract the info we need from these data and write in our csv file.

9. There is one problem we need to take care. Sometimes, there are emojis in the username or comment and they cannot be written into our csv file because of Unicode error.

10. To solve this problem, we need to 'clean' the data we got.

```python
def string_clean(a):#Tiktok: special characters in the jitter text
    m=''
    f1=open('special.csv','w')
    wf=csv.writer(f1)
    for i in list(a):
        try:
            if i=='\n' or i=='\\' or i=='n'or i=='/':
                continue
            wf.writerow([i])
            m=m+i
        except UnicodeEncodeError:
            continue
    f1.close()
    return m
```

If we are dealing with a special character, we just skip them.

11. Then we need to remove the redundant data(mentioned before in step 6).

```python
videoidlist=[]
for i in range(1,100):
    url = 'https://m.tiktok.com/api/recommend/item_lis
    # url='https://m.tiktok.com/api/recommend/item_lis
    headers = {
        'upgrade-insecure-requests': '1',
        'sec-ch-ua': '" Not;A Brand";v="99", "Microsof
        'cookie':'tt_webid_v2=6980332468541769218; tt_
        'user-agent': 'Mozilla/5.0 (Windows NT 10.0; W
    response = requests.get(url=url, headers=headers,
    print(len(response))
    for i in response:
        print(i)
        video_commentCount = i['stats']['commentCount'
        video_diggCount = i['stats']['diggCount']
        video_playCount=i['stats']['playCount']
        video_shareCount=i['stats']['shareCount']
        video_id = i['id']
        if video_id in videoidlist:
            continue
        videoidlist.append(video_id)
```

If the video is already fetched, we just skip it, otherwise, add this video's id to our list.

12. It could fetch under 30 videos in this process, so we just add a large for loop outside.