



# 在无 X API 授权下实时抓取指定用户推文的方案对比

由于 X（原 Twitter）对未授权 API 的抓取做了严格限制，我们需要另辟蹊径来实时获取某用户的推文数据。下面从开源工具、第三方平台/服务、自主开发三个角度，介绍多种可行方案，并比较它们的实现原理、使用步骤、实时性、数据完整性、优缺点以及被反爬检测的风险。

## 一、使用开源工具进行推文爬取

借助社区维护的开源爬虫工具，我们可以在无需官方API密钥的情况下抓取推文数据。这些工具通常通过解析 Twitter 网页或调用非公开接口来获取数据。

### 1. Snscreape

**方案介绍与原理：** Snscreape 是一个轻量的开源社交媒体爬取工具，支持 Twitter 等平台<sup>1</sup>。它无需 API Key，通过模拟网页请求来抓取公开的推文数据（例如使用 Twitter 未登录状态下的搜索接口）。Snscreape 可以抓取用户时间线、搜索结果、话题标签等，返回包含推文内容和元数据的对象。

**操作步骤：** 安装非常简单：使用 `pip install snscreape` 即可。安装 Python 3.8+ 后，可直接在命令行运行，如：

```
snscreape twitter-user <用户名> > tweets.txt
```

这会抓取指定用户的推文并导出到文本文件。也可以在 Python 中调用 `snscreape.modules.twitter.TwitterUserScraper` 获取推文列表<sup>2</sup>。Snscreape 不需要登录认证，开箱即用。

**实时更新支持：** Snscreape 本身不提供持续监听功能，但可以通过定期运行来近实时获取新推文。例如，每隔几分钟运行一次抓取命令，获取最新的 N 条推文，然后和之前的数据比对来发现新增内容。由于没有推流机制，属于轮询式获取，频率太高可能会碰到速率限制。

**数据完整性：** Snscreape 能获取推文文本、发布时间、推文ID、来源等基本信息<sup>3</sup>。对于带有媒体的推文，可以提取媒体的URL或媒体ID（需额外步骤下载图片/视频）。互动数据（点赞数、回复数、转发数等）在部分版本中也可以获取。对于转发推文，Snscreape会返回转发本身的文本（通常是简单地引用原推文作者和内容链接）以及可能附带的评论文本。总体而言，凡是在网页上公开可见的字段，Snscreape理论上都能抓取<sup>4</sup>。

#### 优点：

- **无需授权：** 不需要官方API密钥或登录凭证，使用门槛低，完全免费开源<sup>5</sup>。
- **历史数据：** 不受官方API时间窗口限制，可抓取很久以前的推文（支持按日期过滤），甚至比某些受限的API获取更多历史记录<sup>6</sup>。
- **使用简便：** 提供CLI和Python库，两三行命令即可拿到数据，支持导出为JSON/CSV等格式。
- **社区维护：** 项目由开发者积极维护，面对 Twitter 界面或接口变化会及时更新，从而恢复抓取功能。

#### 限制：

- **依赖页面结构：** 由于通过解析网页或非公开接口实现，一旦 X 网站改版或反爬策略升级，Snscreape 可能暂时

失效，直到开发者适配更新<sup>5</sup>。使用过程中需要关注项目更新。

- **非实时推送：**无法像流式API那样实时推送新推文，需要轮询获取，实时性略有欠缺。频繁轮询可能导致被暂时禁止访问。
- **媒体处理需手动：**返回数据中媒体内容以链接形式给出，如果需要下载图片或视频，需要自行编写下载逻辑。
- **规模和速率限制：**虽无固定配额，但仍受限于Twitter对未登录访问的限制。大量抓取时可能遇到请求被拒（如HTTP 429）或被要求登陆。因此大规模运行需要控制节奏并可能使用代理IP。

**反爬风险：**总体来说，Snsnscrape尽量模拟正常浏览器行为，但在高频抓取时仍有被检测的可能。X平台具有复杂的反爬机制，若短时间内请求过于频繁，可能触发登录墙或CAPTCHA验证<sup>7</sup>。幸运的是，Snsnscrape不使用账户，因此不存在账号被封的风险；但IP可能会被临时封禁。一些缓解措施包括降低抓取频率、随机化请求间隔、使用不同User-Agent以及必要时使用代理IP以分散流量<sup>8</sup>。

**示例项目：**Snsnscrape项目GitHub上提供了详尽的使用文档和示例<sup>9</sup>。

## 2. Twint (已停维护)

**方案介绍与原理：**Twint是一个曾经非常流行的Twitter爬虫工具，它允许在**无需 API**的情况下抓取推文等数据<sup>10</sup>。Twint的原理是利用Twitter早先的公开搜索接口（例如通过构造查询from:用户名获取用户推文），解析返回的JSON或网页数据，从而爬取内容。当年它因能绕过官方API限制（尤其是可免费获取大量历史推文）而广受欢迎。

**操作步骤：**安装同样简单：`pip install twint`（注意需要Python 3.6-3.8环境）。然后可以通过Python调用或命令行使用，如：

```
import twint
c = twint.Config()
c.Username = "TargetUser"
c.Limit = 500 # 限制抓取条数
c.Store_csv = True
c.Output = "tweets.csv"
twint.run.Search(c)
```

上面示例使用Twint抓取指定用户最近的若干推文并存为CSV。以往Twint支持丰富的搜索过滤参数（关键词、时间段等），使用方式高度可定制。

**实时更新支持：**Twint本身也没有实时推送功能，需要周期性运行来获取新推文。在Twint尚可用的时期，用户可以结合调度脚本让其每隔几分钟抓取一次并追加新数据。但需要注意，由于Twint目前已无法有效绕过新反爬机制，实时更新的意义不大——它已经无法抓取最新数据。

**数据获取能力：**过去Twint能够获取推文文本、发布时间、推文ID，以及作者资料、回复/转发/点赞数量等元数据，甚至可提取推文所含的图片链接、视频链接等（通过解析推文的详细信息）。Twint还支持抓取用户简介、关注列表、粉丝列表等多种信息。然而随着X网站的更新，这些功能如今大多失效或不可靠。

**优点：**（基于历史表现）

- **免费且功能丰富：**不需API密钥，曾支持广泛的数据类型和过滤查询，适合OSINT调查等应用。
- **抓取历史数据：**可以不受时间限制获取大量旧推文，这是官方API免费版难以做到的。
- **输出格式方便：**自带将结果保存为CSV/JSON的功能，易于后续分析。

### 限制：

- **已停止维护：** Twint 项目开发者已停止更新，仓库标记为只读<sup>11</sup>。在2023-2025年间Twitter频繁调整反爬策略后，Twint 很多功能已无法正常使用<sup>12</sup>。
- **需规避登录墙：** 2023年中 Twitter 曾一度强制登录后才能查看内容，这种情况下 Twint 在未登录模式下完全失效。即使后来部分公开访问恢复，Twint依赖的很多公共端点也被关闭或需要认证<sup>13</sup>。
- **不再可靠：** 综合上述原因，2025年的共识是**Twint 已不可用**<sup>13</sup>。任何继续使用Twint的方案都很难稳定获取所需数据，除非对其源码进行大的改造。

**反爬风险：** 由于Twint目前几乎无法抓取数据，谈不上风险。但在其可用时，Twint大批量抓取也曾可能触发Twitter的流量限制。不过当时Twint可以通过使用代理或降低速度来一定程度缓解反爬。在当前环境下，不建议依赖Twint，否则很可能频繁遇到请求被拒甚至IP封禁。

**示例项目：** [Twint GitHub仓库](#) (已归档，仅供参考历史代码)<sup>14</sup>。实际应用中建议考虑替代工具。

## 3. Twscrape

**方案介绍与原理：** Twscrape 是近年出现的一个 Python 工具，它通过**模拟官方Web API调用（GraphQL等）** 来抓取推文数据<sup>15</sup>。简而言之，Twscrape伪装成Twitter网页或App，通过使用真实账户的令牌调用内部接口获取数据，因此能突破未登录限制和部分反爬机制。它支持异步抓取、多账号池轮切、甚至直接调用Twitter隐藏的GraphQL接口获取推文、用户信息等。

**操作步骤：** Twscrape 可以通过 `pip install twscrape` 安装<sup>16</sup>。使用前需要提供至少一个可用的Twitter账户（手机号验证通过的）用于获取数据：

- **添加账号：** 可以在代码中添加账号凭证，或者提供已登录账号的Cookies。<sup>17</sup> 表明，通过cookies登录更稳定，也可用用户名/密码登录（需要能接收验证码的邮箱）。添加多个账号有助于轮换降低单账号的压力。
- **初始化 API 对象：** 调用 `api = twscrape.API()` 初始化后，先登录 `await api.pool.login_all()` (如使用用户名密码方式)。
- **调用抓取：** 然后可使用 `await api.user_by_login("username")` 获取用户信息，或 `await api.user_tweets(user_id, limit=100)` 来获取某用户最近若干推文，等等。Twscrape 提供了丰富的方法，比如搜索推文、获取推文详情、获取转发/回复列表等<sup>18</sup> <sup>19</sup>。

Twscrape 也带有CLI工具，可用命令行执行类似操作，如：

```
twscrape user_tweets <用户ID> --limit=100
```

来抓取用户ID对应账号的100条推文。<sup>20</sup>

**实时更新支持：** 由于Twscrape本质上相当于使用真实账号在调用Twitter接口，它每次抓取都可以获取最新的数据。如果我们维持账号不断线并周期性查询，就能以近实时方式得到新推文。Twscrape自身不具备流式推送，但我们可以通过**轮询GraphQL接口**或利用多线程异步不断查询某用户的最新推文ID，从而达到准实时监控效果。需要注意频率控制，以免触发账号的频繁访问警报。

**数据完整性：** Twscrape 能获取非常全面的数据，因为它实际上调取的是Twitter的内部API。例如：

- **推文内容：** 完整文本（包括带有表情或特殊符号），以及嵌入的媒体URL、GIF等。
- **互动信息：** 点赞、回复、转推、浏览量等计数都可以获取<sup>21</sup> <sup>22</sup>。
- **推文关系：** 能区分推文类型（原创、回复、转推、引述），并获取被引用推文的内容。
- **媒体和附件：** 可以获取媒体列表（如图片链接、视频直链等）以及推文中的链接卡片信息等。
- **用户信息：** Twscrape 也提供获取用户资料、关注者列表、关注列表等功能<sup>23</sup>。

总体而言，Twscrape与官方API拿到的数据差异很小，基本相当于登录账户查看所见即所得的内容。

#### 优点：

- **功能强大**：几乎覆盖所有推文和用户数据（包括互动数和媒体），堪称非官方API的完整替代方案<sup>15</sup>。
- **较稳定**：通过模拟正式接口，避免了解析HTML，受页面结构变化影响较小。如果Twitter改动内部API，Twscrape开发者也会跟进调整。
- **异步高效**：支持Python异步操作，可并发抓取多个用户或多个接口，加快数据收集速度。
- **多账户轮换**：内置账号池机制，可在检测到速率限制时自动切换账号，以平滑通过Twitter的限流规则<sup>24</sup>  
<sup>25</sup>。

#### 限制：

- **需要账号**：必须提供至少一个Twitter账号的凭证（Cookies或密码）。这对某些无法提供账号的场景不适用。而且新注册账号常受限，可能需要购买或借用稳定账号<sup>26</sup>。
- **实现较复杂**：初次配置涉及账号登录、有时要处理验证码或邮件验证。相比Snsccape简单即用，Twscrape在使用上技术门槛略高。
- **潜在维护工作**：尽管Twscrape作者很活跃，但Twitter若改变内部API（如GraphQL查询参数）可能导致工具失效，需要等待更新或自行修改。
- **环境要求**：需要Python 3.10+运行环境<sup>27</sup>，并且熟悉异步编程以充分利用其优势。

**反爬与封禁风险**：Twscrape因使用真实账号，所以风险主要在账号层面：如果抓取频率太高，可能触发账号的反机器人检测。例如出现要求验证手机号或账户暂停。这就是为什么作者建议使用多个账号、轮换代理IP等来降低单账号和单IP压力<sup>28</sup>。另外，要防止使用一个账号连续抓取过量数据——尽管GraphQL接口比网页请求更高效，但Twitter仍有严格的速率限制（每个账号/每IP每15分钟的调用次数等），超限就会返回错误。总的来说，**Twscrape比纯未登录爬虫更不易触发IP封锁**（因为看起来像正常用户行为），但要注意不要让一个账号承担过量爬取任务，否则可能导致账号被临时禁用。

开源项目链接：[Twscrape GitHub仓库](#) 提供了使用文档和示例代码，可供参考学习<sup>15</sup>。

## 4. Twikit

**方案介绍与原理**：Twikit是国内开发者开源的一个Twitter API封装库。它的特点是在**不需要官方API Key**的情况下，实现了类似官方API的大部分功能，包括抓取推文、搜索内容，甚至发布推文等<sup>29</sup><sup>30</sup>。Twikit的实现原理和Twscrape相似，也是通过**网页抓取技术和Twitter内部接口相结合**，来绕过官方API限制<sup>31</sup>。可以把它视作一个“逆向工程的Twitter客户端”。

**操作步骤**：安装方法：`pip install twikit`<sup>32</sup>。使用前同样需要提供一个Twitter账号用于登录（支持邮箱+密码登录，并保存会话Cookies）。基本使用流程：  
- **初始化客户端**：`client = twikit.Client(lang='en-US')` 来创建客户端实例。  
- **登录授权**：调用 `client.login(username, email, password, cookies_file='cookies.json')` 进行账号登录<sup>33</sup><sup>34</sup>。成功后会获取并保存会话cookies以供后续调用使用。  
- **调用功能**：Twikit封装了高层接口，例如 `client.search_tweet('关键词', 'Latest')` 搜索最新推文，<sup>35</sup>；`client.get_user_tweets(user_id, 'Tweets')` 获取某用户的时间线推文列表<sup>36</sup>；`client.create_tweet(text='...')` 发布推文<sup>37</sup>；甚至`client.send_dm(user_id, 'Hello')` 发送私信等<sup>38</sup>。多数方法都是异步的，需要在 `asyncio` 环境下运行。

Twikit提供了类似官方API的调用体验，但实际上背后执行的是模拟Web请求。例如搜索和抓取推文时，它会先获取必要的token然后请求Twitter隐藏接口来拿数据，实现对用户透明。

**实时更新支持：** Twikit 并非为推文实时监测而设计，不过我们可以使用它定期轮询某用户的新推文。通过记录最后抓取的推文ID，然后间隔调用 `get_user_tweets` 查看是否有更新推文。如果需要高实时性，也可以使用 `search_tweet` 按关键词（如某用户名或特定话题）持续搜索最新内容。但需要注意**API速率限制**：Twikit虽无官方限制，但受制于账号行为，如频繁调用可能导致账号被限流。因此实时监控时应适当降低频率并考虑多账号并用。

**数据完整性：** Twikit 致力于提供全面的数据获取能力：

- **推文及互动：** 可以获得推文文本、发布时间、来源，以及回复数、转发数、点赞数等统计<sup>21</sup><sup>22</sup>（这些在Twikit返回的对象里都有对应字段）。
- **媒体内容：** 支持上传媒体并发布推文，也可以获取推文自带的媒体ID/链接，从而下载图片或视频。
- **用户和趋势：** 除了推文，Twikit 还能调用获取当前流行趋势（trending topics）<sup>39</sup>、用户详细资料、关注/粉丝列表等等，可谓功能全面。
- **限制方面：** 由于Twikit本质上调用的是Twitter内部API，只要登录账户有权限看到的数据，它基本都能抓取到。如果目标用户的推文是公开的，就能获取完整内容；若对方加密（protected）账号，则和网页上一样获取不到。

**优点：**

- **功能完整：** 除了抓取推文，还支持发推、私信等操作，相当于非官方的完整API封装<sup>40</sup>。这在需要双向互动（不仅仅读取）的场景下特别有用。
- **免费开源：** 项目使用MIT协议，完全免费<sup>41</sup>。开发者可以自由修改扩展。
- **降低门槛：** 提供类似于官方API的使用方法，对习惯Twitter官方开发的用户来说，上手容易。
- **更新积极：** Twikit在2025年仍保持更新（GitHub已有4k+星），社区反馈活跃。碰到问题可以通过issues讨论，有一定支持力度。

**限制：**

- **账号需求：** 与Twscrape一样，需要账号登录操作，对没有可用账号的用户不友好。而且登录可能触发验证码、人机验证等流程，增加使用复杂度。
- **潜在风险：** Twikit包含发送推文、私信等写操作功能，一旦账号被不当使用甚至可能违反Twitter使用政策而被封号，更需谨慎。
- **操作复杂度：** 对只需读取数据的人来说，引入完整客户端库可能有些繁琐；若无使用发布等功能需求，Twscrape之类更简单。
- **环境依赖：** 需要Python环境支持异步操作，另外因为它涉及账户权限，最好专门准备爬虫账号而非使用主账号，以防万一。

**反爬风险：** Twikit的风险和Twscrape类似，主要集中在账号。因为它等于是用您的账号在运行各种操作，频繁抓取会在Twitter后台留下行为记录。如果按官方API标准，它执行的内部调用很可能超出官方允许频率。因此**极有必要控制好调用频次，最好使用轮换IP和多账号策略**。开发者在文档中也强调了避免单账号高频操作，以保护账号安全<sup>30</sup>。另外，如果使用Twikit的写操作（比如批量发推），更容易引起注意甚至触发封禁，要格外当心。总之，遵循“小批量、多账户、慢频率”的原则可降低被检测的几率<sup>42</sup><sup>43</sup>。

开源项目链接：[Twikit GitHub仓库](#) 提供了英文和中文文档，详细说明了安装配置和各功能用法<sup>29</sup><sup>40</sup>。

## 二、第三方抓取平台/服务

如果不想亲自编写爬虫代码，可以利用一些第三方的平台或服务来抓取推文数据。这些平台通常已经封装好了爬取逻辑，用户通过简单配置即可获取数据，有的甚至提供可视化界面和定时抓取等功能。下面介绍几类代表性的平台：

## 1. Apify 云端爬虫

**方案介绍与原理：** Apify 是一家知名的网页数据抓取云平台。其 **Apify Store** 上有现成的 Twitter (X) 抓取“Actor”（爬虫脚本）。这些 Actor 由开发者编写并公开，供用户直接运行。例如“Twitter Scraper”可以抓取任何用户的推文、主题标签下的推文、推文回复线程、粉丝列表等<sup>44</sup>。Apify 的 Twitter 爬虫被设计为 **无需官方 API**：它利用非官方接口和 Apify 自有的 IP 代理机制，在未登录或模拟登录状态下提取数据<sup>45</sup>。对于用户来说，这相当于一个**托管的爬虫**：不必关注具体实现，只需输入目标用户名等参数即可。

**操作步骤：** 使用 Apify 的一般步骤：

1. **注册 Apify 账号：** 登陆其网页控制台。有免费套餐可用一定配额。
2. **选择 Actor：** 在 Apify Store 中找到所需的 Twitter 爬虫脚本，如“Twitter User Tweets Scraper”<sup>46</sup>。
3. **填写输入参数：** 例如指定 `username: "elonmusk"` 和 `num: 100`（抓取 100 条）<sup>47</sup>。不同 Actor 参数可能略有不同，但通常至少需要用户名或关键词。
4. **运行爬虫：** 点击开始，Apify 后台会分配云浏览器或请求去抓取数据。完成后，可以在控制台查看和下载结果数据。支持 JSON、CSV、Excel 等多种格式导出<sup>48</sup>。
5. **调度与 API：** Apify 支持将任务调度定时运行（如每小时抓取一次），也支持通过 REST API 调用这些爬虫，从而集成到应用程序中。

**实时更新支持：** Apify 平台允许**定时调度**任务，因此可以实现准实时的数据更新。例如设置“每 15 分钟运行一次 Twitter 爬虫 Actor”，就能持续获取新推文。由于 Apify 后端支持并行和 IP 代理，调度频率相对灵活。不过完全的实时（秒级监控）难以做到，但分钟级更新不成问题。此外，Apify 也提供 Webhook 等机制，在任务完成后推送通知，这样数据抓取一完成就可被处理。

**数据完整性：** Apify 的 Twitter 爬虫通常能获取**完整的推文内容和相关数据**，包括：

- 推文文本、发布时间、推文链接及 ID。
  - 作者用户名和显示名，作者认证状态等。
  - 点赞数、回复数、转发数、浏览量等互动指标<sup>49</sup><sup>50</sup>。
  - 推文所含媒体（通常提供媒体的 URL 列表）和外部链接。
  - 对于抓取用户时间线的 Actor，往往还能返回用户简介信息，方便一并记录。
  - 某些 Actor 还支持获取**回复线程**，即某条推文下的所有回复内容；或者**引用转发**（quote tweets）等。
- 例如，Apify 的“Twitter User Tweets Scraper”输出示例中，包含 `full_text`、各类 `_count` 字段，以及 `medias` 字段（无媒体则为 null）等<sup>51</sup>。由此看出，常见数据字段均有覆盖。

**优点：**

- **无需编程：** 图形界面操作，无需自己处理登录、抓取逻辑，适合不擅长编码的人士。
- **功能多样：** Apify 上有不同的现成爬虫，可满足多种用途（用户推文、关键词搜索、粉丝列表、趋势等）。而且这些爬虫经过维护，一般能跟上网站变化。
- **云端可靠：** 抓取任务在云服务器上执行，不依赖本地环境；同时利用 Apify 的代理网络和防封机制，成功率较高。
- **扩展集成：** 提供 API 接口和 Webhook，可以方便地把抓取流程融入自己的应用或数据管道中。
- **调度方便：** 内置计划任务和队列系统，易于实现持续监控和定期抓取，省去了自己搭建调度的工作。

**限制：**

- **成本考虑：** 虽有免费额度，但较大量或高频率抓取需要付费。Apify 按使用的计算时间和流量计费，持续实时监控可能成本不低。
- **第三方依赖：** 数据经过第三方平台处理，可能有合规和隐私顾虑；某些敏感数据抓取可能违反 Apify 政策。
- **黑盒程度：** 虽然省心，但平台封装逻辑不可见。如果 Apify 的爬虫出现错误或数据遗漏，用户无法立即修复，只能等待其开发者更新。

- **速率限制：** Apify虽然有反封机制，但仍需遵守Twitter公开访问限流。一些免费actor可能对**每次运行可抓取的推文数量做了上限**（比如每次最多1万条），如需更多需要特殊方案或购买高级服务<sup>52</sup>。

**反爬风险：** Apify在平台层面做了很多防护。例如，其Twitter爬虫actor大多可以在未登录模式下工作，利用了Twitter匿名接口和Apify的住宅代理池<sup>45</sup>。这使得**单个IP请求量降低**、每次运行都可能分配新IP，从而降低被封风险。另外Apify脚本通常会适当随机等待并模拟正常用户行为。不过，近期X对未授权访问依然敏感，如果目标用户发布频率极高，连续抓取大量分页也可能中途遇到瓶颈。但总体而言，由专业团队维护的Apify爬虫**稳定性和隐蔽性优于个人脚本**。在Bright Data等专业代理加持下，Apify声称其爬虫**没有登录就无配额限制地提取数据**<sup>45</sup>（实际仍受制于网站基础限制，只是通过分布式和代理绕开）。因此使用Apify方案被检测封禁的概率较低，风险主要在于违反平台ToS的法律层面。

**代表项目：** Apify 平台上有很多相关脚本，例如：

- [Twitter User Tweets Scraper](#): 抓取指定用户的推文及回复<sup>46</sup>。
- [Twitter \(X\) Scraper](#): 通用的Twitter数据提取工具<sup>44</sup>。
- [Twitter Followers Scraper](#): 抓取用户关注者列表等。

（以上链接在Apify Store可找到，其README有用法说明。）

## 2. PhantomBuster 自动化服务

**方案介绍与原理：** PhantomBuster 是一种云端的**社交媒体自动化**平台，提供各种“Phantom”小工具来完成特定任务，包括多个Twitter相关的Phantom。例如“Twitter Profile Scraper”可以提取用户简介和最近几条推文，“Twitter Tweet Extractor”可以批量下载指定账户的推文<sup>53</sup>。PhantomBuster 的特点是**使用你的账户模拟人在浏览器上执行操作**：一般通过提供你的Twitter会话cookie，让它在云端浏览器登陆你的账号，然后按脚本要求访问页面并提取数据<sup>54</sup>。相当于一个**云端运行的无头浏览器**，按预定流程抓取数据。

**操作步骤：** 使用PhantomBuster抓取Twitter通常如下：

1. **注册并安装扩展：** 注册PhantomBuster账号，并安装其浏览器插件，用于获取你的Twitter会话cookie<sup>54</sup>。
2. **提供认证信息：** 通过插件“一键连接Twitter”，PhantomBuster会获取你当前登录Twitter的session cookie，并保存到云端<sup>55</sup>。这使Phantom可以以你的身份登录。
3. **选择Phantom并配置：** 在PhantomBuster控制台选择所需的Twitter Phantom，例如“Twitter Tweet Extractor”。配置参数如目标的Twitter用户URL列表、要抓取多少帖子等<sup>56</sup>。通常可输入一个或多个用户名（支持CSV/Google表格批量）<sup>57</sup>。
4. **运行并查看结果：** 启动Phantom后，它会在云端启动浏览器，登录你的账号，然后访问每个目标用户的页面，滚动加载推文并提取需要的数据。完成后，可下载结果数据（JSON/CSV）。
5. **重复和计划：** PhantomBuster支持将任务设置为定时重复运行<sup>58</sup>。例如可以设置每3小时跑一次，以抓取新内容。它也提供API，可通过HTTP请求触发Phantom执行。

**实时更新支持：** PhantomBuster 可以设置**周期性运行**（Scheduling），比如每小时、每天定时运行某个Phantom<sup>58</sup>。这样就能持续更新数据，虽然粒度不是秒级，但频率可以自由设定（注意高频会耗尽账户执行时间配额）。PhantomBuster的**云浏览器**实时访问Twitter页面，理论上只要Twitter前端能显示最新内容，它就能抓取到。因此对于实时性要求较高的场景，可以将间隔调小。但需考虑PhantomBuster每日执行次数或时长限制（免费账户通常每天几十分钟执行时间）。没有原生推送功能，但可以结合其API或Webhook自行实现检测新数据后通知。

**数据完整性：** 由于PhantomBuster采用**浏览器真实渲染**页面并读取数据，它抓取的信息就是网页上呈现的所有公开内容。以“Twitter Tweet Extractor”为例，它会抓取指定用户的**最新数百条推文**，包括：

- 推文文本全文

- 发布时间
  - 每条推文的链接
  - 回复、转推、点赞数量
  - 媒体内容（如包含图片，会获取图片URL；包含视频/GIF，也能提取媒体的链接或缩略图）
  - 如果Phantom配置抓取更多详情，甚至可以点击每条推文查看回复内容。一般Profile Scraper不抓取全部回复，Tweet Extractor只抓主贴内容。
- PhantomBuster还有**Profile Scraper**，可以抓取用户简介（昵称、bio、位置、主页链接、关注数、粉丝数等）。另外有**Search Export**工具可按关键词/话题抓取搜索结果。总之，PhantomBuster能获取的字段与人工打开网页右键复制的内容类似，**数据完整度高于API**（因为部分统计如浏览量也是可见即可得）。例如官方描述中提到“将每个公开信息都提取”<sup>59</sup>。

#### 优点：

- **模拟人工，成功率高：** 使用真实浏览器和用户会话，能避开绝大多数反爬检测。对于JavaScript渲染的动态内容也可以正确加载<sup>60</sup>。
- **丰富的现成工具：** PhantomBuster针对Twitter提供多种Phantom，可抓取推文、粉丝、搜索结果甚至自动关注/取消关注等，功能全面。
- **使用简单：** 非技术用户也能通过其图形界面完成配置。特别是浏览器插件获取cookie的设计，让登录集成非常方便<sup>55</sup>。
- **云端执行：** 不占用本地资源，执行时即使关闭电脑也不影响<sup>61</sup>。而且PhantomBuster托管了解决IP、代理等问题的细节。
- **结果直观：** 支持Excel/CSV导出，适合营销人员直接使用。还可以连接Google Sheets实现自动填充更新。

#### 限制：

- **需要账号授权：** 必须使用你的Twitter账户。这对安全有一定隐患（需信任该服务不会滥用账号）；同时如果账号权限有限（如受年龄地区限制）可能影响抓取范围。
- **受账号可见性限制：** 用你的账号抓取能看到的内容。如果目标账号将你拉黑，你就无法抓取其内容；若目标是受保护（private）账号，你也抓不到。
- **速率与配额：** PhantomBuster收费按执行时间计算。免费试用每天有定额分钟，付费也有上限。同一时间爬太多账户或频繁执行，可能需要购买更高级别。
- **潜在封号风险：** 因为用你的账号做自动化操作，如果Phantom动作不够模拟人（如过于频繁刷新或遍历大量用户），有可能导致你的Twitter账号被标记为机器人行为，甚至触发验证或冻结。这在PhantomBuster官方教程中也有提示，建议稳妥使用。
- **数据量限制：** 由于通过页面加载，单次Phantom运行能抓取的推文条数有限（通常几百条以内，具体取决于Twitter前端允许下拉的量）。大量历史数据抓取效率不如使用API搜索方案。

**反爬风险：** PhantomBuster通过浏览器模拟，大幅降低了反爬被发现的可能性。它使用你的账号cookie，也就绕过了未登录限制。但**风险转嫁到账号**：如果抓取操作频繁，比如短时间查看上千条推文，Twitter可能检测到异常行为。PhantomBuster官方建议通过延时和限制频率来保持安全，同时它的云浏览器会附带正确的浏览器指纹和行为，降低触发Bot检测的概率<sup>62</sup>。很多用户成功使用PhantomBuster进行日常数据提取而账号无恙，但也有报告称如果过度使用（如不停爬取大量用户数据）账号可能会遇到额外验证。因此PhantomBuster总体是在反爬和账号安全间取得平衡：正常使用问题不大，但滥用仍有封号可能<sup>8</sup>。建议使用专门的小号来授权Phantom，以防万一。

**代表服务：** PhantomBuster 官网提供的Phantom包括：

- [Twitter Profile Scraper](#)：提取用户资料和近期推文。
- [Twitter Tweet Extractor](#)：批量导出指定账户的推文<sup>53</sup>。
- [Twitter Search Export](#)：根据搜索词导出推文列表。
- [Twitter Follower Collector](#)：抓取某账户的粉丝名单等。

这些工具大多有对应的教程说明如何使用和注意事项。

### 3. 无代码爬虫工具（Octoparse、Thunderbit 等）

除了上述专门面向社交媒体的数据平台，一些通用的**无代码爬虫工具**也可以用来抓取Twitter。这类工具通常提供图形界面，用户通过点选网页元素或使用内置模板即可定义抓取流程，无需编程。其中具有代表性的是**Octoparse**和近年来出现的**Thunderbit**平台。

#### 方案介绍与原理：

- **Octoparse**: 一款成熟的可视化爬虫软件，支持Windows/Mac客户端和云服务。Octoparse提供了专门的Twitter抓取模板，内置解决了登录和滚动加载等难题，用户只需输入目标账号的URL即可抓取其推文等数据<sup>63</sup>。Octoparse通过其云端代理和反爬技术，实现**无需登录**即可抓取推文<sup>64</sup>。其原理可能是借助自有的IP池不断获取Twitter的匿名访问令牌，并模拟浏览器无限下拉获取数据。对于一些需要登录的操作（如抓取粉丝列表），Octoparse也支持用户提供账号进行登录抓取<sup>65</sup>。
- **Thunderbit**: 国内推出的一款AI驱动的无代码爬虫平台。Thunderbit宣称利用AI自动识别网页结构并适应变化，针对Twitter这样的动态页面也能“自适应”地抓取<sup>66</sup>。Thunderbit与Octoparse类似，也提供云端运行和调度，支持无限下拉、自动点击加载更多等功能<sup>67</sup>。其原理是在后台使用AI分析页面DOM，自动生成抓取规则并执行，因此当Twitter前端变化时，它可以自动调整抓取策略，降低人工维护成本。

#### 操作步骤：

- **Octoparse**: 可以选择使用本地客户端或者Octoparse的云模板。简单起见，使用其**预设模板**方式：登陆Octoparse后，在模板中心搜索“Twitter”。会看到多种模板，例如“Twitter Scraper (按账户URL)”<sup>68</sup>、“Twitter高级搜索Scraper”等。选择相应模板后，输入参数（如目标Twitter用户名或URL，或者搜索关键词）。然后直接运行任务，Octoparse会在云端按照模板流程（模拟浏览器打开Twitter页面->登录（如需要）-> 自动下拉加载推文 -> 提取字段）抓取数据。完成后，可在Octoparse界面下载数据为Excel/CSV/JSON。也可以设置任务**定时**或**循环**执行，以便持续获取新数据。
- **Thunderbit**: Thunderbit则强调**两步提取**：在其界面中，用户打开目标Twitter页面（如用户时间线或搜索结果），点击“AI智能识别字段”，平台会自动检测页面上的推文列表、作者、时间等元素并列出可采集字段<sup>66</sup>。用户确认无误后，再点击“开始抓取”，系统即会云端运行任务并获取数据。Thunderbit能识别无限加载的页面，会自动滚动到底然后抓取全部内容<sup>67</sup>。抓取完成后，可导出数据或通过API获取结果。此外Thunderbit也支持像Octoparse那样定时调度任务、将数据直接导出到Excel/Google Sheets等。总体来说，Thunderbit把复杂的规则配置简化成AI自动识别，提升了易用性。

**实时更新支持：** 两个平台都支持**调度**。在Octoparse云端，可以设置任务定时运行（如每小时、每天）。Thunderbit同样支持多种调度粒度（小时/天/周），用于连续监测某页面变化<sup>69</sup>。例如，可以设定Thunderbit每2小时抓取一次某话题的最新推文<sup>70</sup>。由于它们通过网页抓取新内容，当新的推文出现在页面顶端时，下次运行就会被抓取到。需注意的是，如果抓取频率过高，Twitter可能会出现让你登录或暂时封禁IP的情况，但这些平台通常有应对（例如Octoparse内置了代理池防止频繁访问同一IP<sup>64</sup>）。因此，在合理频率下可以实现接近实时的更新。

#### 数据完整性：

- **Octoparse**: 官方介绍其Twitter Scraper模板可以抓取推文、转推、图片、回复、点赞、浏览量、粉丝/关注数等几乎所有可见信息<sup>63</sup>。从其模板截图看，**推文内容、时间戳、互动数据**都在抓取列表中<sup>71</sup><sup>72</sup>。甚至连**观看数（Views）**和**帖子ID**等也包括在内。Octoparse强调“不需要登录即可抓取全部可见数据”<sup>64</sup>。实际使用中，Octoparse的模板会输出每条推文的一行数据，包含文本、发布日期、点赞/回复/转推数量、媒体链接、推文URL等字段，非常全面。对于图片，会给出图片URL列表；视频可能给视频的封面图和播放链接。
- **Thunderbit**: 由于采用AI自动识别，理论上页面上显示的字段都会被捕捉。例如用户昵称、用户名、认证状态，小蓝标等也可能被识别为字段。Thunderbit提到其能抓取**完整回复线程或子页面**<sup>67</sup>——也就是如果需要，可以点进每条推文抓取评论。这意味着数据获取的深度和广度都能覆盖。不过具体字段取决于AI识别，用

户也可手动调整。总体而言，用Thunderbit抓Twitter，与人眼能看到的信息保持一致，不会遗漏文本或数字内容。同时借助AI，还能自动解析下一页，如点击“显示更多回复”之类的按钮。

#### 优点：

- **零代码：** 无需编写任何程序，通过点击操作即可构建爬虫流程，门槛极低。非常适合没有编程背景的运营、分析人员。
- **全面的数据抓取：** 可见即可抓，无论互动数据还是多媒体内容都能提取<sup>63</sup>。相比官方API筛掉一些内容，这些工具获取的信息更原汁原味。
- **快速上手：** 模板和AI辅助极大减少了设置时间。Octoparse提供现成模板即用，Thunderbit自动生成抓取方案，一两分钟即可启动采集。
- **维护省心：** 平台会处理反爬。例如Octoparse自带智能代理和验证码识别<sup>64</sup>，Thunderbit用AI来适应HTML结构变化<sup>66</sup>。所以当Twitter界面改变或增加限制时，平台会更新策略，用户不必自己修改脚本。
- **云端高并发：** 这些服务通常允许同时跑多个任务，Octoparse云有40+并发进程<sup>73</sup>，可同时抓取多个账户或话题而互不影响，加快数据收集效率。
- **数据交付灵活：** 支持多格式导出和直接对接数据库/API。比如Octoparse能直接把结果推送到Google Sheets、Dropbox等<sup>74</sup>，方便与现有工作流程集成。

#### 限制：

- **费用问题：** 详尽的服务通常需要付费订阅。Octoparse有免费方案但功能有限，云端高并发和大任务需要专业版。Thunderbit提供部分免费抓取，但持续监控和大批量数据也需要购买其套餐。成本相比自己写代码使用免费API会高一些。
- **平台依赖：** 将抓取工作外包给平台，数据安全和隐私需考量。如果抓取敏感信息，上传到第三方服务器处理可能存在风险。
- **灵活性有限：** 模板和AI虽然方便，但面对非常定制化的需求（例如特定交互流程、多步过滤）可能无法实现。无代码工具对**非标准流程**支持不足，需要平台增加新功能。
- **规模限制：** 由于模拟人浏览，抓取特别大的数据集（比如数百万条历史推文）并不现实。无限下拉通常也有边界，Twitter可能在一定深度后停止返回更多旧数据。另外，这些工具受制于各自代理IP池的大小，**并发过高或目标过多**可能超过其能力范围。
- **学习成本：** 虽说无代码，但要熟练运用模板和理解抓取原理，仍需要一点学习。复杂场景下可能需要摸索多次才能拿到想要的数据格式。

**反爬风险：** 从用户视角，这两家工具都宣称“**不会触发封禁**”：

- Octoparse强调其Twitter模板**不需要登录**，从而避免账号被封，并且内置智能代理解决封IP的问题<sup>64</sup>。甚至在FAQ中明确说明无需提供cookies即可安全使用<sup>75</sup><sup>76</sup>。可见Octoparse投入了反检测措施，使得单个用户不会直接暴露在Twitter监控下。风险主要在于平台自己的抓取节点被批量识别（但Octoparse大规模代理池降低了这种风险）。
- Thunderbit由于也是使用未登录模式+AI调整，单次抓取行为看起来和真人浏览相仿，而且AI可能随机化动作，更不易被规则识别。所以其成功率和安全性也较高。Thunderbit官方也把“不用熬夜调整崩溃脚本”作为卖点<sup>7</sup>——意味着它能平稳运行，而不容易踩到反爬雷区。

总体来说，这类**商用无代码工具反爬能力强于个人脚本**，出现封禁情况的概率很低。如果有，也是平台需要解决的问题，用户只需等待修复。因此从个人用户的角度，使用这些平台抓取Twitter几乎**无需担心被封**（当然前提是不抓取登陆后才可见的敏感数据，否则可能违反法律/政策）。例如Octoparse明确表示其方案受Twitter ToS限制较少，但用户也应遵守法律和平台政策。

#### 代表平台：

- [Octoparse Twitter Scraper 模板](#)：官方模板页面列出了能抓取的字段和特点<sup>63</sup><sup>64</sup>。
- [Thunderbit 官方案例](#)：Thunderbit博客有详实的Twitter抓取攻略<sup>77</sup><sup>78</sup>介绍其方法和优势。

- 其他类似平台还有如 **Lobstr** (提供Twitter Search/User Scraper，需要通过浏览器扩展登录账户<sup>79</sup> <sup>80</sup>) 等，可根据需要选择。总体上，国内外都有不少无代码方案可用。

## 三、自主开发方案

对于有编程能力且希望完全控制抓取流程的开发者，可以考虑**自行开发爬虫**。自主开发分为两种主要思路：**直接基于网络请求的爬虫**（利用Twitter未公开的接口或解析网页HTML），以及**基于浏览器自动化的爬虫**（使用Selenium、Playwright等工具驱动一个浏览器完成抓取）。下面分别讨论这两类方案。

### 1. 基于Web请求的爬虫

**实现原理：** 直接使用代码向Twitter网站的接口发送HTTP请求，获取数据并解析。这需要绕过官方API，通常有两种办法：

- **利用公开的隐式接口：** Twitter网页在未登录时加载公开信息时，会调用一些隐式API。例如之前的搜索JSON接口或Twitter的GraphQL接口。这些接口有时不需要OAuth认证，但需要提供特殊的访问令牌（guest token）和正确的HTTP头。开发者可以通过分析网页请求（F12开发者工具）找到这些接口。例如，有文章总结了匿名模式浏览推文的流程：先请求推文页面HTML，提取里面的**guest token**，再带上该token和查询参数请求 `https://api.x.com/graphql/...` 的端点获取数据<sup>81</sup>。这种方法可以获取单条推文详情或用户时间线JSON数据。
- **模拟官方客户端API：** 移动端或Web端的某些API端点可能没有严格授权，如之前的Twitter移动版API用一个全局Bearer Token即可访问搜索结果。开发者可以尝试使用已知的Bearer Token直接请求 `https://api.twitter.com/2/timeline/profile/<user_id>.json` 之类的接口。如果成功，就能得到结构化的推文列表JSON。不过Twitter在近年关闭了很多无授权的端点，所以需要不断尝试更新可用的接口。

总之，基于Web请求的方案需要逆向工程Twitter网站。具体实现包括：获取**guest token**（通常通过请求 `https://x.com/i/api/1.1/guest/activate` 可以拿到一个token），然后将该token添加到后续GraphQL或REST API请求的头部（`Authorization: Bearer <token>` 等）。之后调用例如 `https://api.twitter.com/2/timeline/profile/<uid>.json?count=100` 来获取用户推文。每种接口的参数和返回格式不一，需要仔细分析。**示例：**有人通过Juejin分享了匿名获取数据的方法，核心也是拿到 `ct0` 和 `gt` 两个token，然后请求GraphQL取得推文列表<sup>82</sup>。

#### 操作步骤：

- 分析获取必要参数：** 使用浏览器访问目标内容，例如用户时间线或某条推文，在网络面板查看请求。找到其中data/API请求，提取请求URL、需要的headers如授权token、cookie等。
- 代码实现登录或获取token：** 未登录模式下，先GET一次 `https://twitter.com` 首页，从响应cookie中拿 `guest_id`，或者直接调用激活guest API拿一个token。若需登录模式，则使用程序模拟登录（这很复杂，涉及登录流程和2FA，不推荐匿名完成）。
- 发送数据请求：** 携带上一步获得的认证信息，请求相应API端点。比如抓取用户推文可用之前分析的URL。常用HTTP库如 `requests` 可以设置headers和cookies，然后请求获取JSON。得到数据后解析JSON结构提取需要的字段。Twitter返回的JSON嵌套深，需小心提取推文内容及相关媒体。
- 循环分页获取：** 通常接口会返回 `min_position` 或 `cursor` 用于下一页。需要根据返回指示，不断请求下一页以获取更多推文，直到无新内容。
- 存储处理：** 将抓到的数据存储为本地文件或入库，或者进一步分析。

**实时更新支持：** 通过上述接口自行抓取，可以编写一个循环或者cron任务，每隔一段时间运行一次请求新数据。如果使用GraphQL或v2接口，可以利用since\_id等参数只取新推文。没有现成推送，需要自己轮询对比。理论上也可以长连接订阅，但Twitter已关闭免费流式端点，因此基本还是靠定时查询。频率由自己掌控，但别太频繁以免触发限制。一般1-5分钟一次查询最新推文JSON，可认为是近实时。对于小规模监控这是足够的实时性。

**数据完整性：**直接调用Twitter内部API拿到的通常是**结构化完整数据**。例如timeline接口会返回每条推文的全文、媒体链接、作者信息、统计数据等，与官方API提供的类似，甚至更多。例如GraphQL往往连每条推文的回复thread都嵌在返回中。不过，需要注意某些字段可能需要登录用户权限才能看到（如某些受限内容）。在匿名模式下，获取的数据和游客浏览到的一致。如果想拿到更多（比如看到成人内容或作者受限内容），就必须用登录态cookie。总的来说，通过这种方案获取的数据与Snscreape、Twscrape类似，都属**原始完整**的数据，没有被裁剪。此外，自己解析HTML也是一种办法，但一般不如直接拿JSON可靠。

#### 优点：

- **完全自主：**不依赖第三方库或服务，抓取逻辑可完全掌控，根据需求灵活定制。例如可以特别处理某些字段、过滤内容等。
- **免去官方API限制：**不用申请开发者账号，也没有严格的速率和配额限制（平台有隐形限流，但没有硬性配额）。如能成功，可大批量获取数据。
- **高效率：**直接请求JSON比通过浏览器渲染更轻量级，响应体积小且易解析，适合大规模数据收集。
- **可部署性：**可将代码部署到自己的服务器上，结合代理IP池，实现分布式抓取，扩展性强。

#### 限制：

- **开发复杂度高：**需要对Twitter前端和API有深入了解，调研和调试成本大。稍有不慎请求格式不对就拿不到数据。属于“黑客式”方案，对开发者要求高。
- **不稳定：**Twitter随时可能修改内部API或增加验证。一旦改动，之前的逆向方法可能失效，必须重新研究。例如2023年推特多次调整guest token机制，让很多脚本反复失效。维护工作量大<sup>7</sup>。
- **反爬压力：**直接请求官方接口很容易暴露为自动脚本行为。没有浏览器环境的伪装，Twitter可以通过请求频率、IP地址、User-Agent等检测你。尤其大量并发请求时，触发封禁的风险更高。
- **法律和政策风险：**这种未授权抓取严格说违反Twitter服务条款<sup>83</sup>。虽然很多研究指出抓取公开数据不违法，但仍存在被追责的可能性。Bright Data就曾因类似行为被Twitter起诉（虽然后来因法律原因驳回）<sup>84</sup>。自行抓取需自行承担风险。

**反爬对策及风险：**采用Web请求方式是**最容易被Twitter发现并限制的**。常见风险和应对：

- **IP封禁：**几次请求后可能返回420 Enhance Your Calm或429 Too Many Requests。对此应使用**轮换代理IP**，尤其使用住宅IP或移动IP<sup>85</sup><sup>86</sup>效果更好，不要全部请求从同一IP发出。
- **Token失效：**guest token有有效期且绑定IP，如果频繁更换IP需要重新获取token。可以实现一个池，失败时重新激活。
- **增加Headers：**模拟真实浏览器请求Headers，例如User-Agent设为常见浏览器，带上Accept-Language等，避免默认请求特征。还可带上引用页Referer。这些细节有助于骗过简单检测<sup>62</sup>。
- **降低频率：**切忌高频爆破。要随机暂停，错开请求时间，避免固定间隔<sup>87</sup>。比如每页抓完sleep几秒甚至用噪声调节。
- **监测异常：**如果收到包含验证码的响应或跳转到登录页HTML，说明被识别出是爬虫了。应暂停一段时间或更换IP后再继续<sup>62</sup>。

尽管可以做上述努力，但与浏览器模拟相比，这种方案抗封禁能力始终有限。建议仅在小规模或一次性抓取时采用。如果需要长期稳定服务，还是考虑使用Selenium/无头浏览器以降低风险。

**示例/项目：**由于政策原因，很少有公开的完整代码示例。但一些技术博客和社区帖子提供了片段：

- [掘金文章《免登录twitter官方接口数据解析》](#)详细描述了匿名抓取流程<sup>81</sup>。
- [云+社区文章《Python爬取Twitter数据的挑战与解决方案》](#)也讨论了GraphQL抓取和代理应对<sup>82</sup>。
- 如果需要，可以参考Snscreape的源码（用Python实现了请求Twitter搜索接口）或Twint的旧代码逻辑，作为开发启发。

## 2. 基于浏览器自动化的爬虫

**实现原理：** 使用自动化工具操控一个浏览器实例，像人一样打开Twitter网页、登录（如需要）、导航到目标用户的主页，然后通过脚本提取页面内容。常用工具有 **Selenium**（支持多语言）、**Playwright**、**Puppeteer**（Node.js）等。浏览器自动化的优势是可以**100%模拟用户行为**：执行JavaScript、加载动态内容、处理滚动和点击等。这种方法几乎等同于人在浏览，因此最难被检测为爬虫<sup>60</sup>。缺点是消耗资源大、实现相对繁琐。但对于实时抓取少量用户的推文，浏览器自动化是可靠的方案。

**操作步骤：** 以Python + Selenium为例介绍：

1. **环境配置：** 安装Selenium库以及对应浏览器驱动（如ChromeDriver）。也可以用Selenium自带的webdriver-manager自动下载驱动<sup>88 89</sup>。选择是否**无头模式**运行浏览器（无界面），一般服务器上用无头模式更方便<sup>90</sup>。
2. **模拟登录（可选）：** 目前X平台允许有限未登录访问，但为了防止突然要求登录，可以预先登录一个账号。用自动化脚本打开<https://x.com/login>，填写用户名和密码，然后等待登录完成并保存会话cookie。这个过程可能需要处理验证码或2FA，算是难点。如果不想自动登录，也可以**手工登录**后将cookies导出供脚本使用。
3. **访问目标页面：** 打开浏览器后，让其前往目标用户的推文页面，例如<https://x.com/<用户名>>。页面加载后，可能需要执行一些动作：关闭弹出的登陆提示（有时X会弹窗要求登录，这时如果已登录账号则无此问题）。
4. **滚动加载：** Twitter用户主页采用**无限下拉**加载推文。Selenium可以执行JavaScript或模拟按键下拉页面到底，不断触发加载新推文，直到加载到需要的数量或到底为止<sup>67</sup>。需设定一个停止条件，比如获取到N条推文或者最近的日期够久远。
5. **提取数据：** 使用Selenium的查找功能（如`find_elements`）定位推文元素。Twitter每条推文可能是一个含有`data-testid="tweet"`属性的`<article>`标签，其中包含子元素：文本div（`data-testid="tweetText"`）、时间span（aria-label包含时间）、互动数span等<sup>91 92</sup>。可以通过XPath或CSS选择器提取这些元素的文本内容<sup>93 92</sup>。对于图片，寻找`<img>`标签的`src`属性；对于视频，可能需要获取`<video>`标签的`poster`或`source`属性<sup>94</sup>。Selenium支持获取属性值<sup>94</sup>。依次遍历已加载的推文元素，采集需要的字段并存入数据结构。
6. **保存结果：** 将收集的推文内容、时间、作者等信息保存到文件（CSV/JSON）或数据库。保存后可关闭浏览器释放资源。
7. **循环监控：** 如果需要持续抓取，同一个浏览器会话可以定时刷新页面或重复上述过程。也可以每次启动新会话（开销较大）。Playwright支持浏览器持久会话，可复用登录状态。通过脚本设置定时任务，定期运行抓取流程，获取增量推文。

**实时更新支持：** 浏览器自动化可以做到**准实时**：一种方法是让脚本**常驻运行**，比如每隔60秒自动刷新用户主页，检查是否出现新推文元素；或者不断轮询顶部tweet的时间。如果检测到新推文，则提取新增内容并触发一些自定义逻辑。这接近实时但较复杂。更简单的是**定期重跑**爬虫，每几分钟运行一次，抓取最新N条推文，与上次结果对比发现新内容。Selenium/Playwright的启动耗时较高，不适合秒级频率，但是3-5分钟一次问题不大。用户可以根据新推文发布频率调整监测间隔。总体上，浏览器自动化可以胜任**分钟级实时**抓取。如果需要秒级推文推送，则应该使用官方流API（不在无API讨论范围）。

**数据完整性：** 由于是真实渲染页面，抓到的数据就是**前端显示的一切**。这包括：

- **完整推文文本：** 带表情、换行，甚至如果有“显示更多”的长推文按钮，脚本也可以点击展开确保拿到全文本。
- **媒体链接：** 图片的真实URL（通常是pbs.twimg.com的地址），视频的播放器链接或预览图地址<sup>94</sup>。获取视频源需要进一步网络请求，但至少能拿到Poster图URL。
- **互动统计：** 点赞、回复、转推数可以通过页面元素读取（注意需要滚过它们进入视野才加载数值）。浏览器脚本可确保这些元素出现再读取。
- **元数据：** 推文发布时间（绝对时间或相对时间，需要转换）、作者用户名、是否是置顶推文（页面上第一个tweet可能是置顶的，需标记）、是否包含敏感内容提示等等，都可以读取DOM属性来判断。

- **限制：**唯一缺憾是浏览器看到的就只能是公开信息。如果某条推文被删除或作者隐藏了回复，你也看不到，自然无法抓取。此外，高级信息如每个粉丝具体资料、某话题下全部推文数这些，需要不同页面，单纯用户主页抓取不到。

#### 优点：

- **仿真度最高：**站在Twitter角度来看，这就是一个普通用户在浏览，极难区分<sup>60</sup>。因此**可靠性最佳**，不容易因为反爬而中断。
- **适应动态内容：**无须处理Ajax或GraphQL细节，浏览器会自动执行JS获取内容，适用于网页结构复杂的场景。比如Twitter大量内容通过客户端渲染，直接请求可能难解析，而浏览器能拿到渲染后的结构。
- **开发调试直观：**开发者可以在自动浏览过程中随时观察浏览器界面（可开启有头模式看浏览器动作），便于调试选择器和定位问题。
- **功能扩展灵活：**可以做到一些API无法的事情，如截图推文页面、运行自定义JS代码进行复杂提取、模拟点击用户界面进行互动等。对于需要获取特殊数据或执行操作的场景，浏览器自动化几乎是万能的。
- **多平台支持：**Selenium/Playwright支持Java、Python、Node.js等多语言，选择面广。Playwright还能方便地跨浏览器（Chromium/Firefox/Webkit）测试，增加模拟多样性。

#### 限制：

- **资源消耗大：**启动一个浏览器比发送一个HTTP请求要重得多，占用CPU和内存也高，抓取速度相对慢。一台普通服务器同时开几十个浏览器可能就吃紧。而HTTP请求可以并发上百。
- **实现复杂：**相比调API，编写浏览器脚本步骤多，也要处理页面各种异步行为。要确保等待内容加载完成再提取，处理可能出现的弹窗。需要一定前端知识（HTML/CSS/JS）来正确选取元素和操作。
- **维护成本：**如果Twitter网页结构变化，如类名更新或层级改变，选择器可能失效，需要调整代码。虽然不如直接请求那么频繁失败，但也需要关注页面更新。
- **登录管理：**选择登录还是不登录？如果不登录，有时访问几页就弹登录要求，这时脚本得能处理登录。登录又带来cookie管理和潜在验证码。这个部分编程量不小。
- **无法大规模高速抓取：**浏览器模拟更适合抓少数账户的数据。要抓全量Twitter数据显然不现实。如果要并发抓取很多不同用户，需要集群分布式跑浏览器，这会非常昂贵。对于题目要求的“一次几百条”的抓取，浏览器方案绰绰有余，但如果要求抓几百万推文，就不适合了。

**反爬风险：**如前所述，浏览器自动化被反爬检测到的概率最低。只要**模拟人行为适当**（不要闪电般滚动，避免毫秒级点击），Twitter的系统很难分辨这是脚本还是用户。特别是现在很多无头浏览器可以设置真实浏览器指纹，使得**连headless标志都去除**。当然，如果用登录账号，需要避免**过频操作**导致账号异常。例如频繁刷新同一页面、有规律地24小时不停地scroll，这可能引发怀疑。建议在脚本中加入一些随机等待和人为行为（比如偶尔停顿模拟阅读），这些都进一步降低风险<sup>42</sup><sup>62</sup>。另外IP依然需要考虑：如果都从同一IP自动访问许多不同账户，也可能触发WAF。所以**分布式部署或代理**仍有用，但由于有浏览器Cookies和JS，可使用**更便宜的代理**而不被弹captcha。总之，正确使用浏览器自动化，**基本不会被封禁IP或账号**。Live Proxies的指南也指出，Selenium是应对X反爬的首选方案<sup>95</sup>。

#### 示例项目：

- 官方Selenium 文档和示例（如Live Proxies博客中给出的Selenium抓取推文代码片段<sup>96</sup><sup>97</sup>）。
- Playwright 官方有例子演示登录并抓取社交媒体内容。
- 开源项目方面，可以参考[Twint的GUI版本](#) 或一些GitHub上的Twitter scrape仓库，有些用到了Selenium。

如果需要一个简单上手的例子：用Python Selenium，在登录Twitter账号后执行以下伪代码，可抓取指定用户推文：

```
driver.get(f"https://x.com/{target_user}")
# 等待页面加载若干推文元素
```

```

tweets = driver.find_elements(By.XPATH, '//article[@data-testid="tweet"]')
for tw in tweets:
    text = tw.find_element(By.XPATH, './div[@data-
testid="tweetText"]').text
    time = tw.find_element(By.TAG_NAME, 'time').get_attribute('datetime')
    print(text, time)

```

这个逻辑获取页面上当前加载的推文文本和时间。可扩展加入滚动加载更多及其他数据字段提取。

以上方案各有适用场景和优劣。总结来说：

- 如果追求**快捷简便**且不想编程，首选开源工具如Snsccape（前提是其当前可用）或者第三方平台如Apify、Octoparse等。这些方案支持单次抓取几百条推文且配置简单。
- 如果需要**持续实时监控**并有一定开发能力，使用浏览器自动化或Twsccape/Twikit这样的库会比较稳妥。在避免官方API的前提下，它们能提供接近实时的更新和完整数据。
- 要考虑**反爬**，无论哪种方法，都应遵守平台规则、控制频率并做好代理/IP切换。特别是纯HTTP爬虫，要做好充足的反检测策略 8 87。

最后提醒：抓取Twitter数据需尊重其服务条款和各地法律。上述方案虽然技术上可行，但在商用时请确认不会违反用户隐私和平台政策。在合法合规的前提下，这些工具和方法将有助于我们高效地获取所需的推文数据。

78 98

1 10 10 Best Twitter Scraper Tools in 2025 - Map Lead Scraper

<https://www.mapleadscraper.com/blog/best-twitter-scrappers>

2 3 6 如何在 2026 年抓取 Twitter (X) 数据 [Python + 工具]

<https://www.rapidseedbox.com/zh/blog/mastering-twitter-scraping>

4 5 7 66 67 77 78 2025年用Python抓取Twitter推文全攻略

<https://thunderbit.com/zh-Hans/blog/scrape-tweets-from-twitter-using-python>

8 12 13 15 42 43 60 62 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 How to Scrape

X.com (Twitter) with Python and Without in 2025 | Live Proxies

<https://liveproxies.io/blog/x-twitter-scraping>

9 How to Scrape Tweets With Snsccape | Medium | ScaleReal - Medium

<https://medium.com/scalereal/million-tweets-and-counting-how-snsccape-can-help-you-scrape-big-data-on-twitter-5c0240cab4f3>

11 14 Twint Python Twitter crawler no longer being maintained . Any alternatives? : r/Python

[https://www.reddit.com/r/Python/comments/vb8cmu/twint\\_python\\_twitter\\_crawler\\_no\\_longer\\_being/](https://www.reddit.com/r/Python/comments/vb8cmu/twint_python_twitter_crawler_no_longer_being/)

16 17 18 19 20 23 24 25 26 27 28 twscrape · PyPI

<https://pypi.org/project/twscrape/>

21 22 46 47 48 49 50 51 52 Twitter User Tweets Scraper · Apify

<https://apify.com/twittapi/twitter-user-tweets-scraper>

29 30 32 33 34 35 36 37 38 39 40 GitHub - d60/twikit: Twitter API Scraper | Without an API key |

Twitter Internal API | Free | Twitter scraper | Twitter Bot

<https://github.com/d60/twikit>

31 41 Twikit: 无需API密钥的Twitter API抓取工具-CSDN博客

[https://blog.csdn.net/m0\\_56734068/article/details/142043196](https://blog.csdn.net/m0_56734068/article/details/142043196)

44 45 Twitter (X) Scraper · Apify

<https://apify.com/scrapers/twitter>

53 Twitter Tweet Extractor - PhantomBuster

<https://phantombuster.com/automations/twitter/30442/twitter-tweet-extractor>

54 55 56 57 58 61 Twitter Tweet Extractor tutorial | PhantomBuster

<https://phantombuster.com/automations/twitter/30442/twitter-tweet-extractor/tutorial>

59 Twitter Profile Scraper | PhantomBuster

<https://phantombuster.com/automations/twitter/9375/twitter-profile-scraper>

63 64 71 72 73 74 75 76 Twitter Scraper | Octoparse Templates

<https://www.octoparse.com/scraping-templates/twitter-scraper>

65 Twitter Profile Scraper | Octoparse AI Automation

[https://www.octoparse.ai/appstore/app/dev\\_lq\\_x21/twitter-profile-scraper](https://www.octoparse.ai/appstore/app/dev_lq_x21/twitter-profile-scraper)

68 Octoparse: Web Scraping Tool & Free Web Crawlers

<https://www.octoparse.com/>

69 70 79 80 How to Scrape Twitter Search Results without coding [2024 Edition]

<https://www.lobstr.io/blog/scrape-twitter-search-results>

81 82 免登录twitter官方接口数据解析（针对匿名模式浏览贴文流程逆向）

<https://juejin.cn/post/7389211334641418249>