# Assignment: -

**1.A developer is assigned a task to scrape 1 lakh website pages from a directory site, while scrapping he is facing such hcaptcha, which are placed to stop people from scrapping**
**As a project Coordinator suggest ways to solve this problem**

As a project co Ordinator I suggest the following ways to slove the issues of facing problem with hcaptcha

**Check for Public Applications:** Look for public APIs that the site may offer for accessing data.

**Use Proxy Servers:** Rotate through a pool of proxy servers to make it appear as if the requests are coming from different IP addresses. This can help bypass IP rate limiting and avoid CAPTCHAs.

**Solve CAPTCHAs:** Utilize CAPTCHA-solving services or libraries that can automatically solve CAPTCHAs when they appear during scraping. Be sure to comply with the site's terms of service and legal requirements.

**Human Intervention:** Employ manual CAPTCHA solving when necessary. Have a human operator ready to solve CAPTCHAs as they occur in real-time.

**Avoid Scraping Too Fast:** Slow down your scraping requests to mimic human browsing behavior. Sites are more likely to impose CAPTCHAs on rapid, automated requests.

**Session Management:** Maintain session state in your scraper so that it appears as if you're a regular user navigating the site.

**Crawl During Off-Peak Hours:** Schedule your scraping tasks during off-peak hours when server loads are lower, reducing the likelihood of encountering CAPTCHAs.

**Continuous Monitoring:** Continuously monitor your scraping process for CAPTCHA challenges and adapt your code to handle them as they arise.

facing hCaptcha while scraping can be a challenge. One way to solve this problem is by using a **headless browser** like **Selenium** to automate the scraping process. This way, the browser can interact with the hCaptcha and solve it programmatically. Another option is to **use a proxy service** to rotate IP addresses and avoid triggering the hCaptcha.

Remember that web scraping may involve legal and ethical considerations, so it's important to be respectful of the site's policies and guidelines.
Additionally, the effectiveness of these strategies may vary depending on the specific website and its security measures.

# 2.Our client has around 10k linkedin people profiles, he wants to know the estimated income range of these profiles. Suggest ways on how to do this?

**To estimate the income range of the LinkedIn profiles,**
you can use a **combination of data scraping and data enrichment techniques**. By scraping the profiles for information such as job titles, industries, and locations, you can then match this data with external sources that provide income data for those specific job titles and locations.
This can give you an estimated income range for each profile.
However, it's important to note that this approach may not provide precise or accurate results for every profile.

Estimating the income range of LinkedIn profiles can be challenging as LinkedIn does not provide specific income information.
However, you can try these general methods:

**Keyword Analysis:** Analyze the profiles for keywords that might indicate income levels. For example, job titles, industries, and the size and type of companies they work for can provide hints about their income range.
**Geographic Data:** Income can vary by location. Use location data on profiles to estimate regional income ranges.
**Education Level:** People with higher education levels may earn more on average. Look at the education section of profiles for insights.
**Connections and Network:** The number and quality of connections can sometimes correlate with income. More connections with influential people might suggest a higher income.
**Industry Benchmarks:** Research industry-specific income data to estimate income ranges for people in certain professions.
**Machine Learning:** You can build a machine learning model if you have a labeled dataset to train on, but this can be resource-intensive and requires quality training data.

# 3.We have a list of 1L company names, need to find linkedin company links of these profiles, how to go about this?

**To find the LinkedIn company links for the list of company names,**
you can use **LinkedIn's API** or a **web scraping tool** to search for each company name and extract the corresponding LinkedIn company profile URLs.
This way, you can gather the LinkedIn company links for the given list of company names.

The another ways is

**Web Scraping Tools:** You can use web scraping tools or libraries such as Scrapy, Beautiful Soup, or Selenium in Python. These tools can help automate the process of searching and collecting LinkedIn company pages.

**LinkedIn Search:** Create a script that goes to LinkedIn's search page and enters each company name as a search query. Then, scrape the search results to find the LinkedIn company pages.

**Parsing Search Results:** Once you get the search results, you'll need to parse the HTML to extract the URLs of the LinkedIn company pages. LinkedIn's website structure may change, so your scraping script may need periodic updates.

**Rate Limiting and Ethical Scraping:** Be mindful of LinkedIn's terms of service and robots.txt file, which may contain information about scraping policies. To avoid overloading LinkedIn servers, implement rate limiting in your scraping script.

**Data Quality:** Keep in mind that not all company names may have LinkedIn company pages, or there may be multiple pages for companies with similar names. You'll need to handle such cases in your script and validate the data.

# 4. How to identify list of companies whose tech stack is built on Python. Give names of 5 companies if possible, by your suggested approach

**To identify companies whose tech stack is built on Python,** you can use various methods.

One approach is to <span style="color:red">**search for job postings or company profiles**</span> that mention Python as a required skill or technology.
This can give you an indication of companies that utilize Python in their tech stack.
As for now I know some well-known companies that heavily use Python include **Google, Facebook, Instagram, Dropbox, and Spotify.**

**Job Postings:** Many companies include the technologies they use in their job postings. You can search for job openings that mention Python as a required skill or part of the tech stack. Websites like LinkedIn, Indeed, and Glassdoor are good places to start.

**GitHub Repositories:** You can search on GitHub for repositories associated with companies. Companies often share their open-source projects, and you can check the programming languages used in their repositories, including Python.

**LinkedIn Company Pages:** Some companies may list the technologies they use on their LinkedIn company pages. You can visit the LinkedIn pages of companies you're interested in and look for "About" or "Tech Stack" sections.

**Social Media :** Companies might discuss their tech stack on platforms like Twitter, Reddit, or tech forums. Searching for mentions of specific companies and Python can provide insights.

# 5. Need to find an API, through which we can send linkedin messages to other linkedin users

There are several APIs available that allow you to send LinkedIn messages to other LinkedIn users programmatically.

One popular option is the **LinkedIn Messaging API**. It provides developers with the ability to send messages to 1st-degree connections on LinkedIn.

**NOTE :** Please note that using the LinkedIn Messaging API requires proper authentication and adherence to LinkedIn's API terms and conditions.

But , Recently LinkedIn has restricted or limited the ability to send unsolicited messages to other LinkedIn users through APIs to prevent spam and ensure user privacy.

They primarily provide the LinkedIn Messaging API to facilitate messaging within approved contexts.

**If you're looking to send messages to LinkedIn users,**
you should consider the following options:

**LinkedIn Messaging API:** If you're running a LinkedIn Ad campaign through LinkedIn Marketing Solutions, you may have access to their messaging features. These can be used for sending messages to users who have interacted with your ads.

**LinkedIn Sponsored InMail:** You can use LinkedIn's Sponsored InMail feature for reaching out to targeted users within the LinkedIn platform. It's a paid advertising feature.

**Manual Outreach:** To send personalized messages to LinkedIn users, you can manually reach out to them through LinkedIn's messaging system. Keep in mind that LinkedIn has messaging limits and policies against spam, so use this method judiciously.

**Third-Party Tools:** There are third-party tools and software available that can automate LinkedIn messaging to some extent. However, be cautious when using these tools, as they may violate LinkedIn's terms of service, and excessive automation could lead to account restrictions.