

<http://poloclub.gatech.edu/cse6242>

CSE6242: **Data** & **Visual** Analytics

# Data Collection

## Duen Horng (Polo) Chau

Associate Professor, College of Computing

Associate Director, MS Analytics

Georgia Tech

## Mahdi Roozbahani

Lecturer, Computational Science & Engineering, Georgia Tech

Founder of **Filio**, a visual asset management platform

Partly based on materials by

Professors Guy Lebanon, Jeffrey Heer, John Stasko, Christos Faloutsos

# How to Collect Data?

## Method

## Effort

---

Download

Low



---

API

(Application program interface)

Medium



---

Scrape/Crawl

High



# Data you can just download

NYC Taxi data: Trip (11GB), Fare (7.7GB)

StackOverflow (xml)

Wikipedia (data dump)

Atlanta crime data (csv)

Soccer statistics

Data.gov

...

# Data you can just download

If you have leads, let us know on Piazza!

More datasets on course website:

CSE6242A,Q/CX4242A Schedule Homework Project Warnings Policies **Datasets** Resources

There are [multiple CSE6242 sections](#). This is the course homepage for **campus CSE6242A,Q/CX4242A**.

CSE6242A,Q/CX4242A Fall 2021

**Data** and Visual Analytics

Georgia Tech, College of Computing

"DVA Live" on Tue & Thu, 3:30pm-4:45pm (Atlanta time) via [BlueJeans Event](#)

# Collect Data via APIs

## Google Data API (e.g., Google Maps Directions API)

<https://developers.google.com/gdata/docs/directory>

## Twitter (small subset)

<https://dev.twitter.com/streaming/overview>

## Last.fm (Pandora has unofficial API)

## Flickr

## data.nasa.gov

## data.gov

## Facebook (your friends only)

API	GData Status	See Also
 Google Analytics Data Export API	Replaced by <a href="#">Google Analytics Core Reporting API</a> (starting at version 2.4).	<a href="#">Migration Guide: M APIs to v2.4 &amp; v3.0</a>
 Google Apps Provisioning API	Shut down. Replaced by the <a href="#">Admin SDK Directory API</a> .	<a href="#">Current Google Ap</a>
 Google Base Data API	Not available since June 1, 2011. Replaced by the <a href="#">Content API for Shopping</a> .	<a href="#">New Shopping API of the Base API</a>
 Blogger Data API	Replaced by the <a href="#">latest Blogger API</a> .	
 Google Book Search API	Shut down. Replaced by <a href="#">Google Books API Family</a> .	<a href="#">Google books API (on Stack Overflow)</a>
 Google Calendar API v2	Shut down. Replaced by latest <a href="#">Google Calendar API</a> .	
 Google Code Search Data API	Shut down in Jan 15, 2012. No replacement API.	<a href="#">A fall sweep (Goog</a>
 Google Contacts API	GData version is still live. Replaced by <a href="#">Google People API</a> for read-only access.	<a href="#">Google Contacts A Google People API</a>
 Google Documents List Data API	Shut down. Replaced by <a href="#">Google Drive API</a> .	
 Google Finance Portfolio Data API	Shut down. No replacement API.	<a href="#">Spring cleaning fo (Google blog post)</a>
 Google Health Data API	The product was <a href="#">discontinued</a> as of January 1, 2013. No replacement API.	<a href="#">An update on Goog Google PowerMet</a>

# Data that needs scraping

Amazon (reviews, product info)

ESPN

eBay

Google Play

Google Scholar

...

# How to Scrape?

## Google Play example

*Goal: collect the network of similar apps*



### Shazam - Discover Music

Shazam Entertainment Limited Music & Audio ★★★★★ 3,102,253

Teen

Contains ads

Add to Wishlist

Install

Similar

See more



### SoundHound Music

The popular music app with 300 million+ downloads globally!

★★★★★ FREE



### TrackID™ - Music

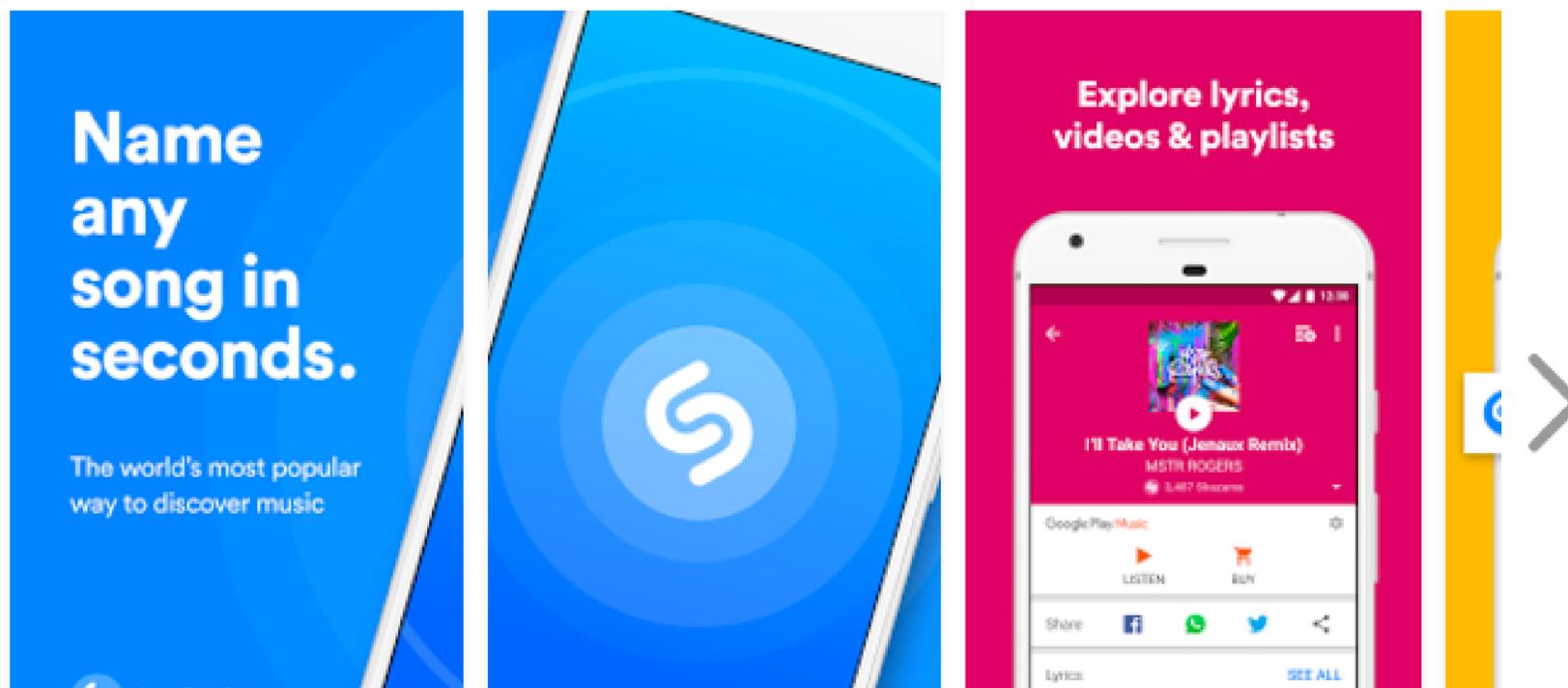
TrackID™ is the best way to identify the music playing around you.

★★★★★ FREE



### Musixmatch Lyrics

Enjoy lyrics for Spotify, Youtube and many other players





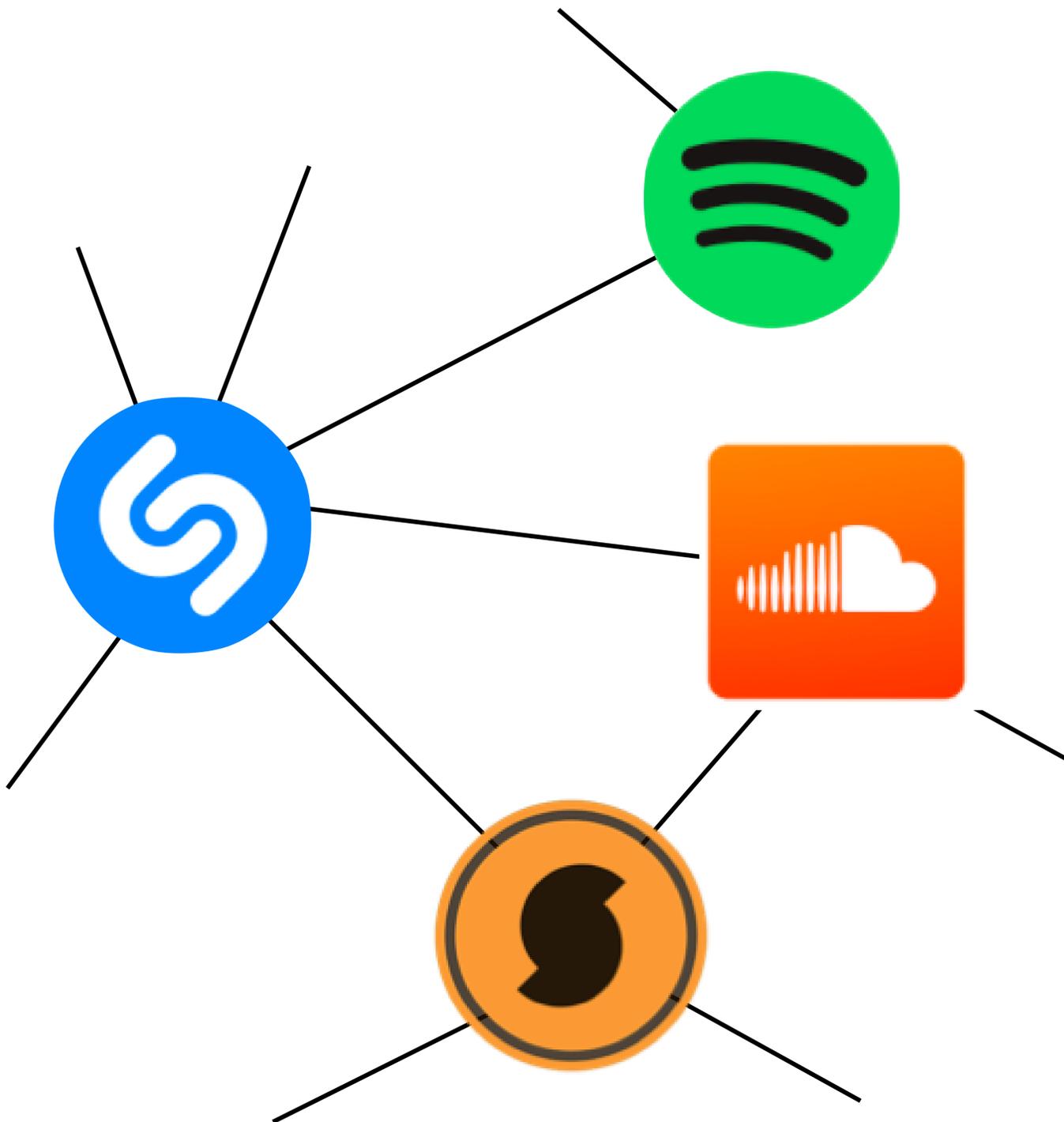
# Name any song in seconds

Shazam will identify any music playing around you.

**GET IT NOW**

# How to Scrape?

Goal: Write a **program/algorithm** to scrape Google Play to **collect a million-node network** of similar apps



Each **node** is an app

An **edge** connects two similar apps

Hint: start with some apps (e.g., Shazam), and go from there.

# How to Scrape?

## Google Play example

*Goal: collect the network of similar apps*

<https://play.google.com/store/apps/details?id=com.shazam.android>



<https://play.google.com/store/apps/details?id=com.spotify.music>

# Popular Scraping Libraries

**Selenium.** Supports multiple languages. <http://www.seleniumhq.org>

**Beautiful Soup.** Python. <https://www.crummy.com/software/BeautifulSoup>

**Scrapy.** Python. <https://scrapy.org>

**JSoup.** Java. <https://jsoup.org>

## Important considerations:

**Different web content shows up depending on web browsers used**  
Scraper may need different “web driver” (e.g., in Selenium), or browser “user agent”

**Data may show up after certain user interaction (e.g., click a button)**

- Scraper may need to simulate the actions.
- Selenium supports more actions than beautiful soup:  
<http://www.discoversdk.com/blog/web-scraping-with-selenium>