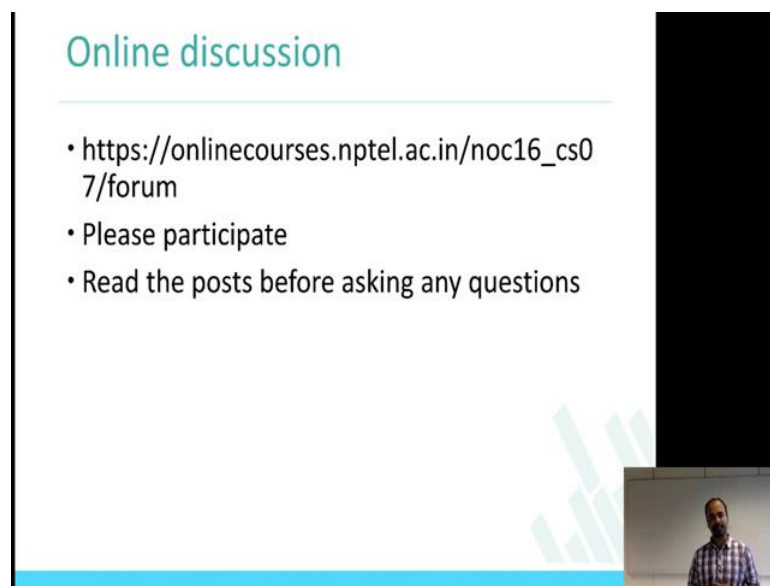


**Privacy and Security in Online Social Networks**  
**Prof. Ponnuram Kumaraguru (“PK”)**  
**Department of Computer Science and Engineering**  
**Indian Institute of Technology, Madras**

**Week - 2.1**  
**Lecture – 06**  
**OSM APIs and tools for data collection**

Welcome back to the course on Privacy and Security in Online Social Media, week 2.

(Refer Slide Time: 00:16)



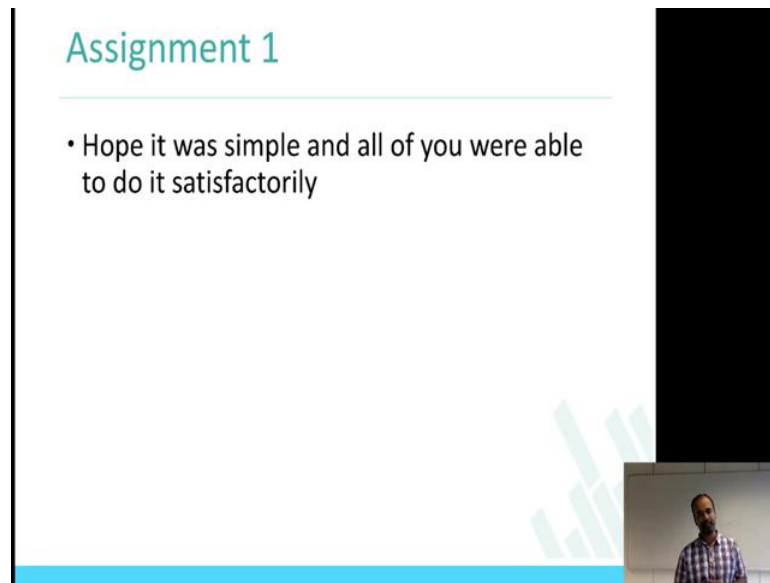
**Online discussion**

- [https://onlinecourses.nptel.ac.in/noc16\\_cs07/forum](https://onlinecourses.nptel.ac.in/noc16_cs07/forum)
- Please participate
- Read the posts before asking any questions

The slide features a light blue header, a thin horizontal line, and a list of three bullet points. A small video inset in the bottom right corner shows the professor speaking. The slide is framed by a black border on the left and top, and a blue bar at the bottom.

I hope you are participating in the online forum that we have in the course. I already see a lot of people asking questions, and trying to answer. My sincere request will be, please, please read the posts before you actually ask the question; that is, read the posts that have been already asked, the questions that have already been asked and the answers that has been already given, before asking the question. And, please participate also in the online forum, not just only asking questions; if you know the answers for the questions that others are asking, please try and answer them also.

(Refer Slide Time: 00:52)



**Assignment 1**

- Hope it was simple and all of you were able to do it satisfactorily

The slide features a light blue bar at the bottom and a small video inset in the bottom right corner showing a man in a plaid shirt. A faint bar chart is visible in the background.

I hope most of you got to see the assignment 1 that we had posted. So, I think, the weekend of the week 2 is the deadline for the assignment 1. Please try to work it out. The assignment 1 is actually pretty simple. We have just captured some questions from the slides that we did, and some from the tutorials.

(Refer Slide Time: 01:15)

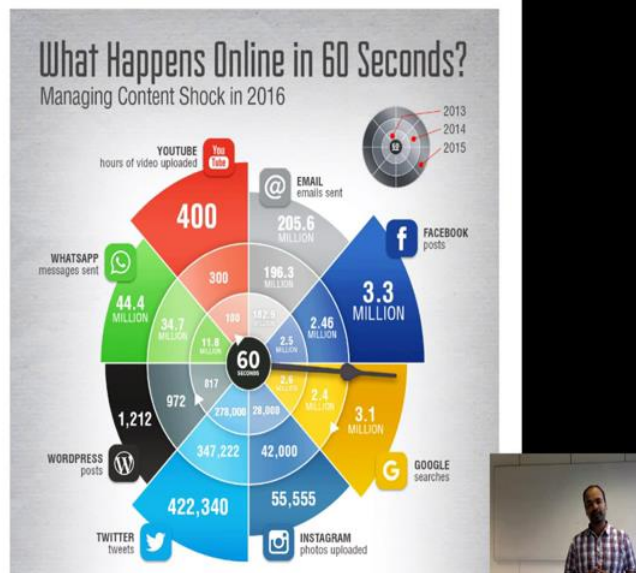


The slide displays a collection of social media logos arranged in a grid-like fashion. The logos include YouTube, flickr, foursquare, tumblr, tinder, facebook, in, whisper, twitter, g+, Instagram, and Pinterest. A small video inset in the bottom right corner shows a man in a plaid shirt.

So, let me just give you a quick summary of what we have seen until now and then, I will go ahead with the topics that I wanted to cover today. So, first, we saw what social media is; different types of social networks, different types of content that gets generated on our

social network; classical online social media services, and then some, which are more like ephemeral social networks and anonymous social networks.

(Refer Slide Time: 01:41)



We also saw what online social media means in 60 seconds; so, how much of data is getting generated on social media in 60 seconds. We saw 400 hours of videos uploaded on YouTube, and 3.3 million posts are done on Facebook and things like that. This basically shows us that, large amount of content that are getting generated on online social media services.

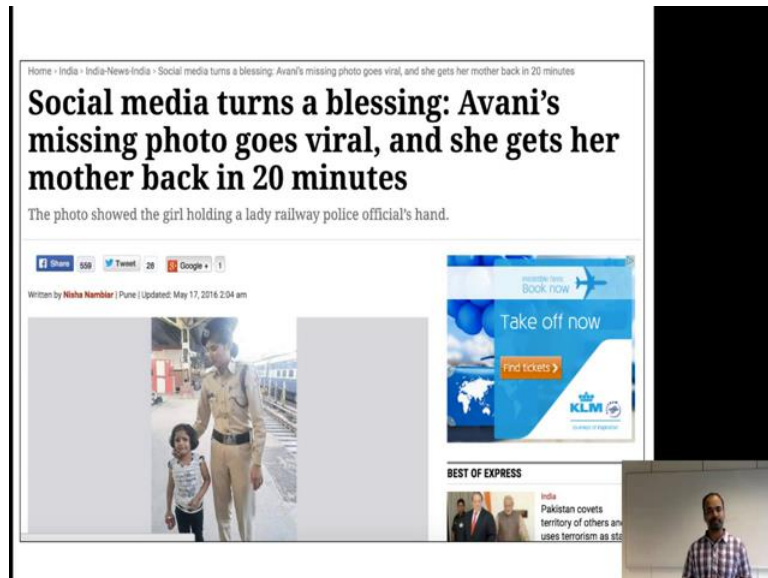
(Refer Slide Time: 02:09)

4 / 5 V's of Online Social Media

The slide features the text "4 / 5 V's of Online Social Media" in a teal font, positioned to the left of a large, bold teal letter "V". A small inset shows a person in a video call.

We also saw what 4 or 5 V's of online social media are, they are volume, velocity, veracity, variety and value - those are the 5 V's of online social media.

(Refer Slide Time: 02:25)



And then, I looked at, I showed you some events, where online social media has played an important role in the real world and in the society also.

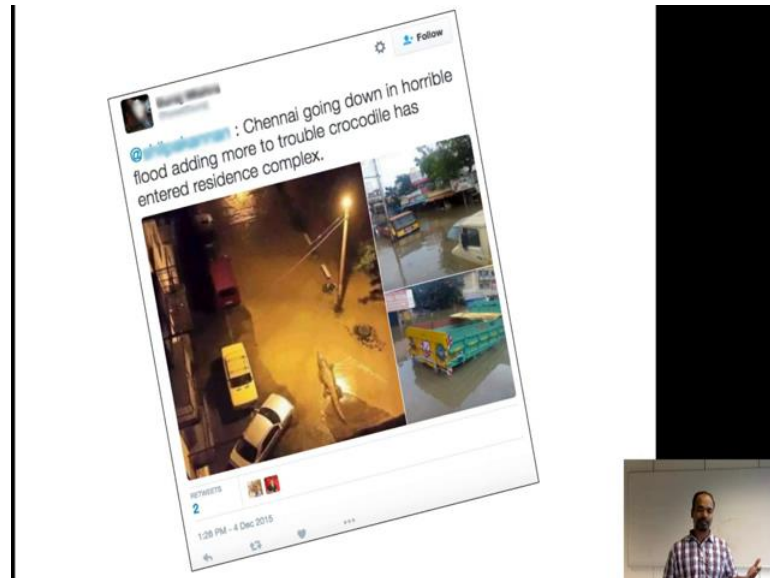
(Refer Slide Time: 02:37)



Telling you about different issues on online social media, for example, in this case, it is compromised account; an account was compromised, where the post said 'Two explosions in White House and Barrack Obama injured', and, there was, there was after

effects of this tweet. So, we looked at different issues that are happening on online social media; compromised account, fake content in this case; and image of a crocodile on the streets of Chennai, while Chennai floods was going on in December 2015, caused panic among citizens.

(Refer Slide Time: 02:59)

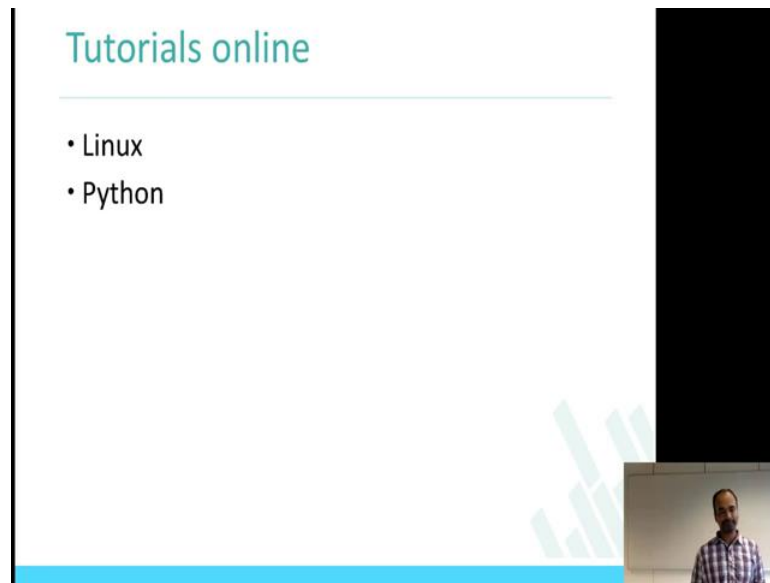


(Refer Slide Time: 03:11)



And, there are also people who lose jobs and others issues, because of the usage of online social media.

(Refer Slide Time: 03:19)



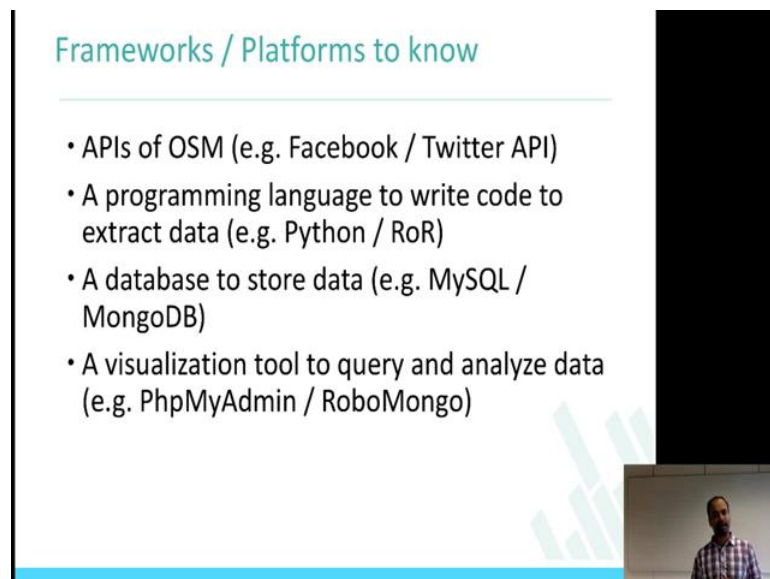
The slide features a title 'Tutorials online' in teal text at the top left. Below the title is a horizontal line. Underneath the line, there is a bulleted list with two items: 'Linux' and 'Python'. The slide has a light blue footer bar. On the right side, there is a vertical black bar and a small video inset showing a man in a plaid shirt.

## Tutorials online

- Linux
- Python

And, in the week 1, we also covered a little bit about Linux and python; hopefully, you are all set, in terms of using the platforms, because, I think, there were some questions about, ‘can we use windows?’ You should be able to use windows, and do programs on python, but it just said our support will be mostly on Linux. And, of course learning Linux will also be good for you.

(Refer Slide Time: 03:44)



The slide features a title 'Frameworks / Platforms to know' in teal text at the top left. Below the title is a horizontal line. Underneath the line, there is a bulleted list with four items: 'APIs of OSM (e.g. Facebook / Twitter API)', 'A programming language to write code to extract data (e.g. Python / RoR)', 'A database to store data (e.g. MySQL / MongoDB)', and 'A visualization tool to query and analyze data (e.g. PhpMyAdmin / RoboMongo)'. The slide has a light blue footer bar. On the right side, there is a vertical black bar and a small video inset showing a man in a plaid shirt.

## Frameworks / Platforms to know

- APIs of OSM (e.g. Facebook / Twitter API)
- A programming language to write code to extract data (e.g. Python / RoR)
- A database to store data (e.g. MySQL / MongoDB)
- A visualization tool to query and analyze data (e.g. PhpMyAdmin / RoboMongo)

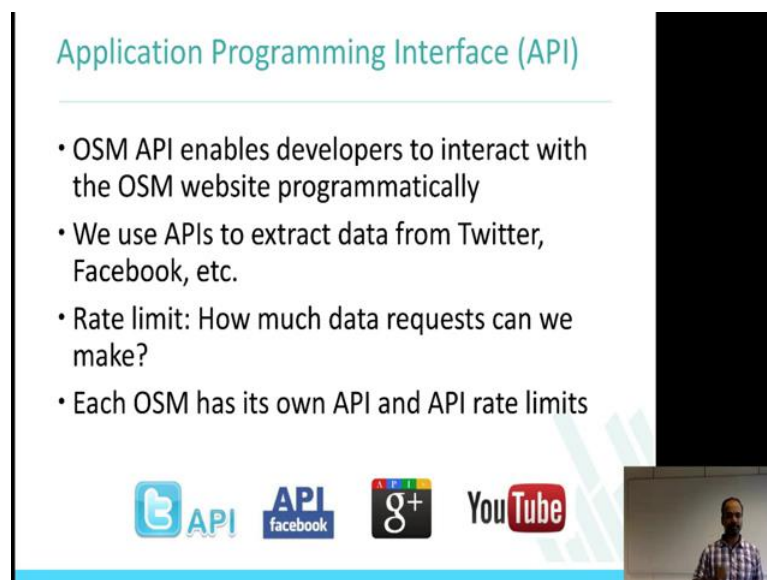
So, what I want to cover today is a couple of things; one, I wanted to actually look at different frameworks or platforms, that you would get to know while doing this course,

or **in another** terms, you should know while doing this course, and collecting data from online social media, analyzing and making inferences.

We will look at what an API is; different kinds of APIs that are available for Facebook and Twitter. Then, we will also look at programming language. There has been a tutorial on python. So, I will just quickly go over. In any case, my work for this week 2.1, about these topics, are only generally, to introduce and then, we will have a tutorials, which are specifically focused on some of them.

Then, we look at programming languages; and then, we will also look at a little bit of database, how this data is stored, what kind of format that the data is coming out; and a little bit about visualization tool.

(Refer Slide Time: 04:45)



**Application Programming Interface (API)**

- OSM API enables developers to interact with the OSM website programmatically
- We use APIs to extract data from Twitter, Facebook, etc.
- Rate limit: How much data requests can we make?
- Each OSM has its own API and API rate limits

Twitter API   Facebook API   Google+   YouTube

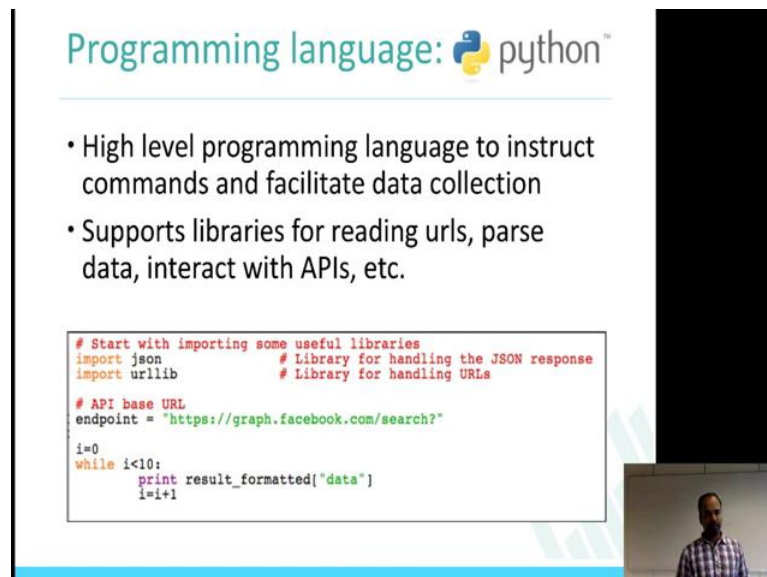
*(A small video inset in the bottom right corner shows a man in a plaid shirt speaking.)*

First, API, which is Application Programming Interface; this basically enables you to interact with the online social media, programmatically, and collect data from there. What does this mean? This basically means that, you can actually have a tunnel that is from your program to the social media services, to collect data. It just creates a tunnel between your program and the online social media services, where you are going to ask some data and then, the social media service is going to respond with saying, here is the data that you asked for, right.


Particularly, in our case, we will actually look at **APIs** for Facebook and Twitter, which will help you to collect data from Facebook and Twitter. There is other APIs also; all other social media services or majority of the social media services provide you with an API. We **can't cover everything** in this course. So, we are going to start looking at only the most popular ones, or the ones that we can actually use for this course, which will help you to understand how APIs work, what data can be collected. So, you can actually do it for other social media services, yourself.

So, **one of the** important thing that you want to also keep in mind is that, about the rate limit, which is that in social media services when we want to collect data you cannot collect the data everything that is available on social, on these services. Because, I am sure, the companies do not want to give you all the data also. They have set it up, you know, by saying that, they have a rate limits for every social media service, and every piece of data that we want to collect from them. So, we will look at something in the tutorials about rate limits, particularly about each **of** the social media API , but I wanted to just give an idea about, there is going to be a rate limit, in terms of the data you can collect from these services.

(Refer Slide Time: 06:33)



The slide features a title 'Programming language: python' with the Python logo. Below the title are two bullet points: 'High level programming language to instruct commands and facilitate data collection' and 'Supports libraries for reading urls, parse data, interact with APIs, etc.'. A code block contains Python code for importing libraries and making a request to a Facebook search endpoint. A small video inset shows a man in a plaid shirt.

```
Programming language:  python™
```

- High level programming language to instruct commands and facilitate data collection
- Supports libraries for reading urls, parse data, interact with APIs, etc.

```
# Start with importing some useful libraries
import json # Library for handling the JSON response
import urllib # Library for handling URLs

# API base URL
endpoint = "https://graph.facebook.com/search?"

i=0
while i<10:
    print result_formatted["data"]
    i=i+1
```

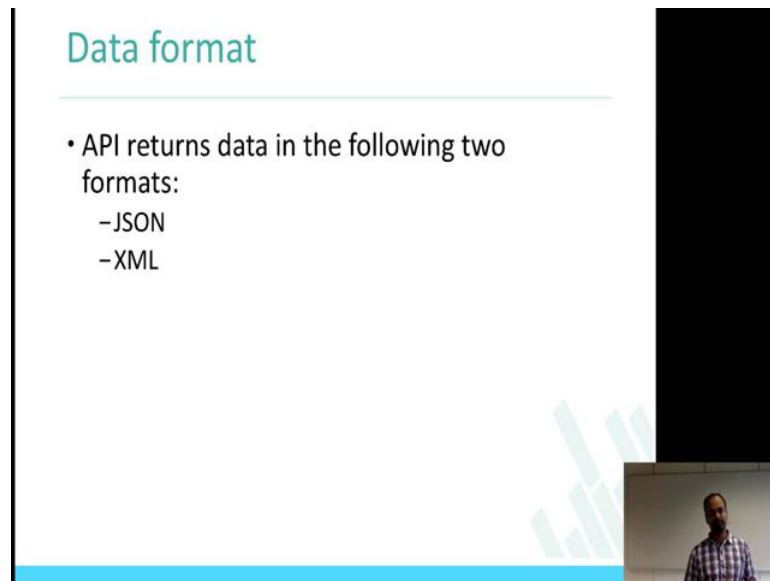
Also then in python, since you have already done a tutorial on python, I will keep it really short. It is basically a programming language, that is used to collect data and is one



of the popular languages currently used in terms of writing API requests to the social media services.

And, it also has a lot of libraries for reading URLs, parsing data, interact with API, and understanding the JSON objects, and things like that.

(Refer Slide Time: 07:00)



The slide is titled "Data format" in a teal color. Below the title is a horizontal line. The main content is a bulleted list:

- API returns data in the following two formats:
  - JSON
  - XML


In the bottom right corner of the slide, there is a small video inset showing a man with a beard wearing a plaid shirt, speaking. The background of the slide features a faint bar chart graphic.

Data format, so particularly the API, when you send the request to Facebook saying that, ‘please give me all the data about friends that PK has’, or, about **the date of birth** of PK, or about my friends’ network. So, what it is going to give you back is actually, it is going to give you in some format. One of the formats that it gives you is actually a JSON format, which we will see in brief what this format means and how we actually interpret the data that is coming back from Facebook, or Twitter. XML, which is also a format with some social media services **give**, or the JSON, is also a little bit like an XML, which is Extended Mark-up Language.

(Refer Slide Time: 07:49)

## JSON

- JSON - JavaScript Object Notation
- Data structuring notation
- Sample:



The screenshot shows the Graph API Explorer interface. At the top, it says 'Graph API Explorer' with 'Application: Graph API Explorer' and 'Locale: English (US)'. Below that, there's an 'Access Token' field with a long token and buttons for 'Debug' and 'Get Access Token'. The 'Graph API' section has an 'FQL Query' input field containing 'GET --> /151831547?fields=id,name' and a 'Submit' button. Below the query, there's a 'Node: 151831547' section with a tree view showing 'id' and 'name'. To the right, the JSON response is displayed: 

```
{
  "id": "151831547",
  "name": "Mehdi Dogra"
}
```

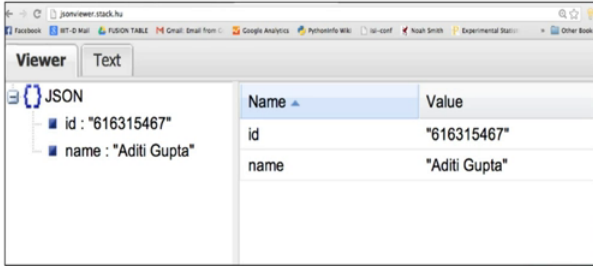
So, here is what a JSON means. JSON means, JavaScript Object Notation, which is a data that you get back from the social media services. So, here is an example that I have in this slide, which just shows you about the JSON object that is returned, when you are asking for id and name of a particular user. So, this is the Graph API Explorer, which you will see in the tutorial in more detail but, it is essentially a through by browser you can actually look at the data, look at the JSON objects of the Facebook data of yourself, or whatever the Facebook API allows, which we will be able to see through this graph API.

So, again, that we emphasize JSON is the JavaScript Object Notation, which is the way that the data is stored in Facebook, data is stored in twitter when you request through the API, for saying, 'give me this data about PK', it is returning the data in JSON format. It is basically the format that most social media services use today.

(Refer Slide Time: 08:53)

# JSON

- Viewing




The screenshot shows a web browser window with the URL 'jsonviewer.stack.hu'. The page has two tabs: 'Viewer' (selected) and 'Text'. The JSON data is displayed in a tree view on the left and a table on the right. The tree view shows a root object with two properties: 'id' with value '616315467' and 'name' with value 'Aditi Gupta'. The table on the right has two columns: 'Name' and 'Value'. It lists the same two properties: 'id' with value '616315467' and 'name' with value 'Aditi Gupta'.

Name	Value
id	"616315467"
name	"Aditi Gupta"

So, when you take the data from JSON, and when you want to interpret the data that is available in this JSON, data that is coming back from Facebook or Twitter, you can actually use JSON dot viewer dot stack dot hu. This is only for you to see visually, what data is coming back; you can take the data that is coming out of Facebook, copy paste it into this JSON viewer, and you will be able to see, what the **fields** are. When you look at the data that is coming back from Facebook, it is generally a block of data; it is just a lot of data that comes back. So, you can actually take it, and put it into the JSON viewer, to see what are the fields that it is actually giving you. We go through this slowly, when we do the tutorials.

(Refer Slide Time: 09:36)




- Relational Database to store data
- Data is stored in rows and columns
- Retrieve using SQL queries
- Sample:

```
mysql> select user_id , user_screen_name from users_data_boston2 limit 2;
```

user_id	user_screen_name
605747286	MacieBliss13
294673452	I_like_HotCake

```
2 rows in set (0.05 sec)

mysql>
```



And, of course when you collect the data, so first is API which is a way by which you want to collect the data, and the data is coming back in JSON. When you collect the data, you have to store it in some format, right. So, the format that majority of the times, the data is stored, is in MySQL. Basically, it is a relational database to store the data, and data is stored in rows and columns, and simple queries, you could use to get the data.

For example, in this case, I am just selecting user id, user screen name from the data that is being collected through Facebook, right. So, that helps **meaning, again** I am emphasizing that, this is not a course on MySQL itself; we will only look at some simple queries on how to look at the data that you have actually stored through the programs that you have written.

(Refer Slide Time: 10:28)

The slide is titled "MongoDB" and shows a terminal window with the following JSON data:

```
db.tweets.find( {text: /'user_screen_name': /}, limit(3), pretty() )
```

```
{
  "_id": "546746b4582d7f72d99994",
  "text": "@Bancroft @CityPolice @P91r will attend to this",
  "user": {
    "screen_name": "@ptrwestbcq"
  }
}
{
  "_id": "546746b4582d7f72d99994",
  "text": "@ptrwestbcq @P91r @d6lqtrffic pls let me clarify to all.Backline Venkatesh was our chief guest who wished the children with flowers.",
  "user": {
    "screen_name": "@ptrwestbcq"
  }
}
{
  "_id": "546746b4582d7f72d99994",
  "text": "Road safety week: Handbills on Traffic awareness given to the public at Yeswestbur by PI Swainathan http://t.co/e2G6jvY5",
  "user": {
    "screen_name": "@ptrwestbcq"
  }
}
{
  "_id": "546746b4582d7f72d99994",
  "text": "Road safety week: Bike rally for Traffic awareness at Jeyanagar today. http://t.co/uA3mcdm",
  "user": {
    "screen_name": "@ptrwestbcq"
  }
}
{
  "_id": "546746b4582d7f72d99994",
  "text": "Road safety week: HDL corner of Police S-I B Demanded with children at a painting competition at Cobben park on Sunday http://t.co/uAM3bhw",
  "user": {
    "screen_name": "@ptrwestbcq"
  }
}
```

A small inset video shows a presenter in a blue shirt.

MongoDB is one of the popular ones, more recently we have started looking at and people are actually using this. So, MongoDB is another way by which the data is stored and the data that is collected from Facebook is actually stored.

(Refer Slide Time: 10:47)

The slide is titled "PhpMyAdmin" and contains the text: "Access MySQL databases and query using browser". A callout box points to the "SQL" tab in the interface with the text "SQL to query databases".

The screenshot shows the PhpMyAdmin interface with a table named "tbl\_whowho" containing the following data:

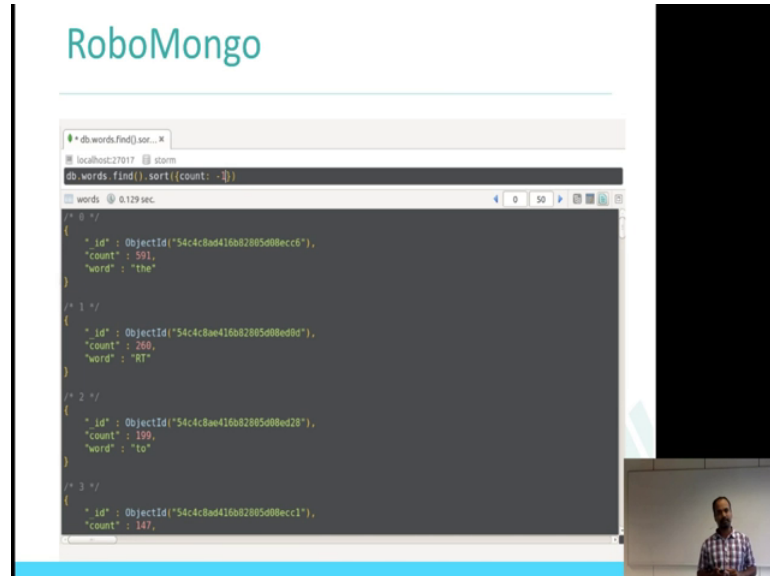
user_id	fname	lname	friends	gender	twitter	facebook	badges	following	photos	checkin
1000388	Tpni	B.	11	female	NULL	1321163198	21	8	2	41
10005512	Tamara	O.	188	female	samaramm	541738782	230	578	73	65
1000980	Maryn	Quaman	522	female	marynquaman	515328219	162	761	65	65
1001189	Abanindo	M.	236	male	imove	517964607	18	4	2	2

A small inset video shows a presenter in a blue shirt.

So, again, let me emphasize which is API; then, there is programming language; then, there is MySQL database or MongoDB, which is data is coming through an API, collected and **dumped into this MySQL or MongoDB**. So, now, we also need a way by which to look at the data that is being stored. So, one of the ways you could use this

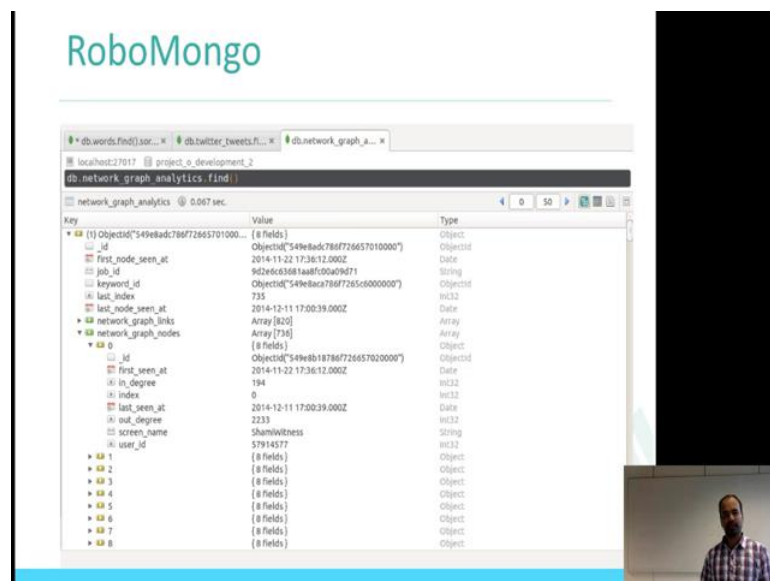
actually phpMyAdmin, which actually allows you to look at the data that you have in your own database.

(Refer Slide Time: 11:20)



So MySQL phpMyAdmin can look at the data from MySQL, and RoboMongo will help you to look at the data from a MongoDB. So, essentially, these are the ways by which you can collect the data, store the data and look at the data that is available with you.

(Refer Slide Time: 11:38)

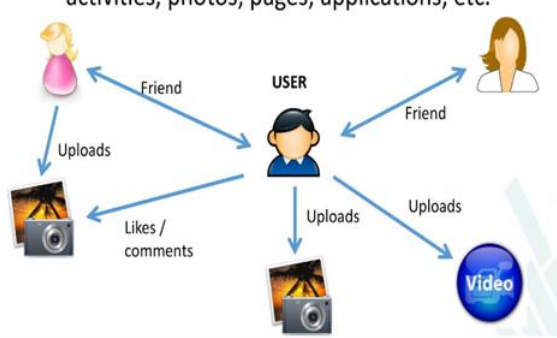


So, this is another view of RoboMongo, which shows you what are the different fields that are available; what data is stored in those fields.

(Refer Slide Time: 11:53)

## All content in graph form

- Graph API
  - Interface to extract data related to User profiles, activities, photos, pages, applications, etc.



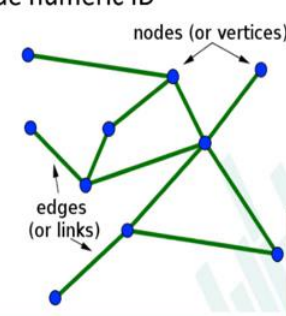
The diagram shows a central 'USER' node (represented by a person icon) connected to two 'Friend' nodes (represented by person icons). The 'USER' node is also connected to three 'Uploads' nodes (represented by photo and video icons) and one 'Likes / comments' node (represented by a photo icon). The connections are labeled 'Friend' and 'Uploads'.

All content on Facebook is actually stored in a graph format; that is, user - the friends that I would have, the pictures that I upload, the videos that I upload, and the status updates that I do, everything is actually a node in the graph. And, every interaction, which is basically like the comments, likes and things like that, becomes edges in this graph. Facebook actually stores all interactions, of all data that they have within the graph format; that is why the API that they have is also called as a graph API.

(Refer Slide Time: 12:27)

## Why is it called the Graph API

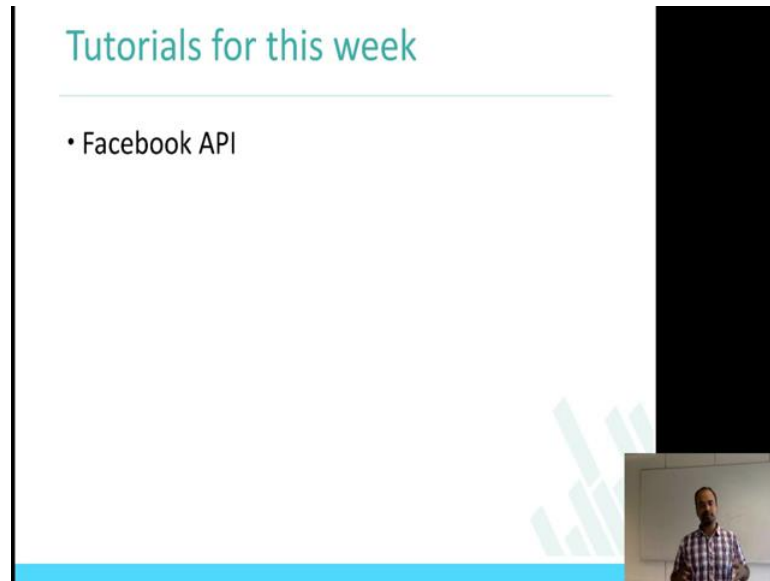
- All objects are stored as nodes of a “graph”
- Connections (likes, friendship etc.) are edges
- All nodes have a unique numeric ID
  - Users
  - Pages
  - Posts
  - ...



The diagram shows a network of nodes (or vertices) connected by edges (or links). The nodes are represented by blue dots and the edges are represented by green lines. The diagram is labeled 'nodes (or vertices)' and 'edges (or links)'.

Here is the another view of the same message, which is, all objects are stored as nodes in the graph; connections like friends, friendships, likes are edges and all nodes have a unique numeric ID, which is users, pages and posts. And, we will be talking mostly about users; we shall later talk about also, **pages**, which is one of the ways by which content can be generated on Facebook.

(Refer Slide Time: 12:52)



In tutorials this week, you will actually look at in detail about what a Facebook API is, how do you actually create the secret key, what kind of authentication that you will have to provide Facebook, in terms of collecting data, what data can be collected and things like that.