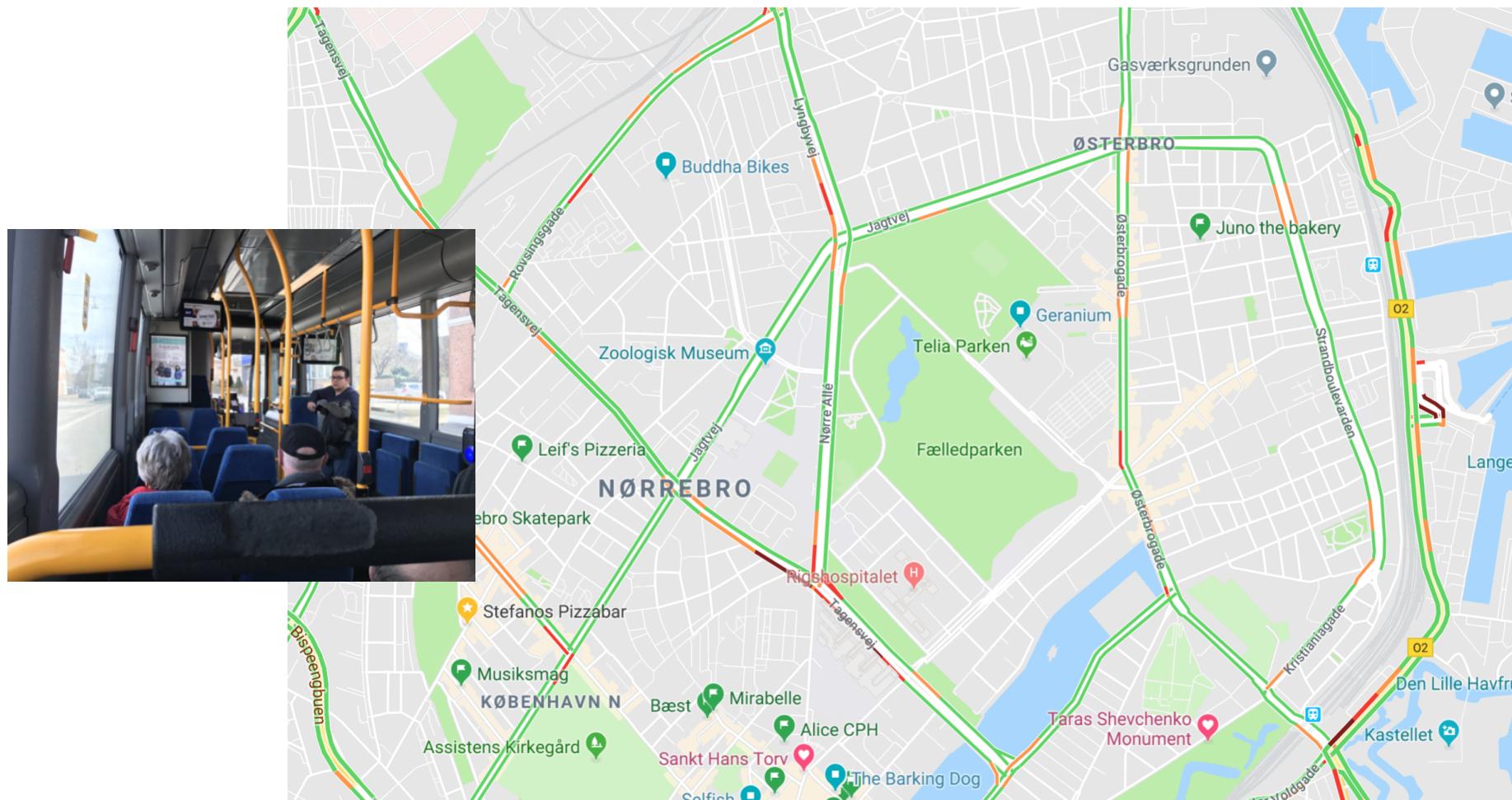


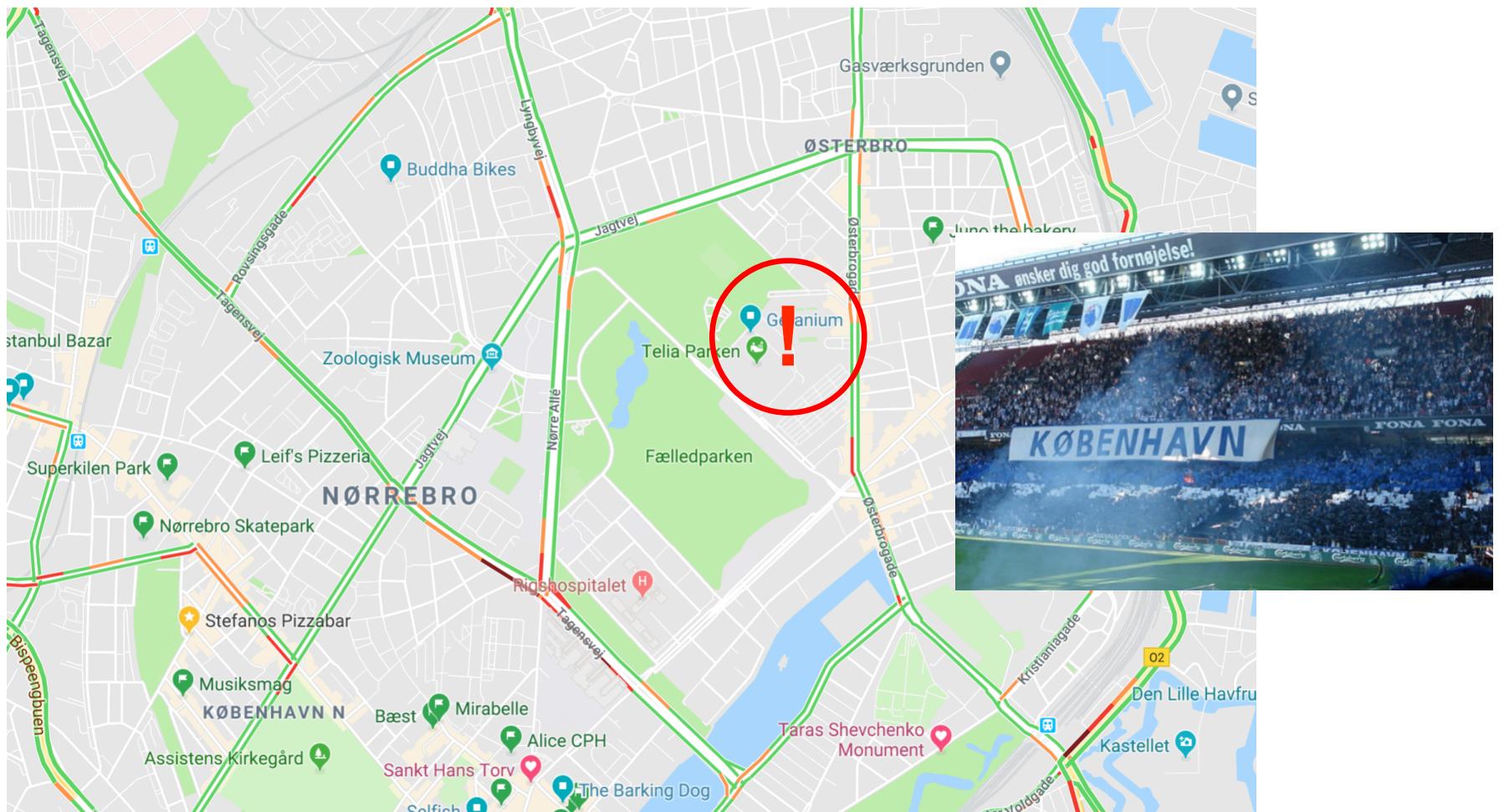
42578 Advanced Business Analytics

Web Data Mining

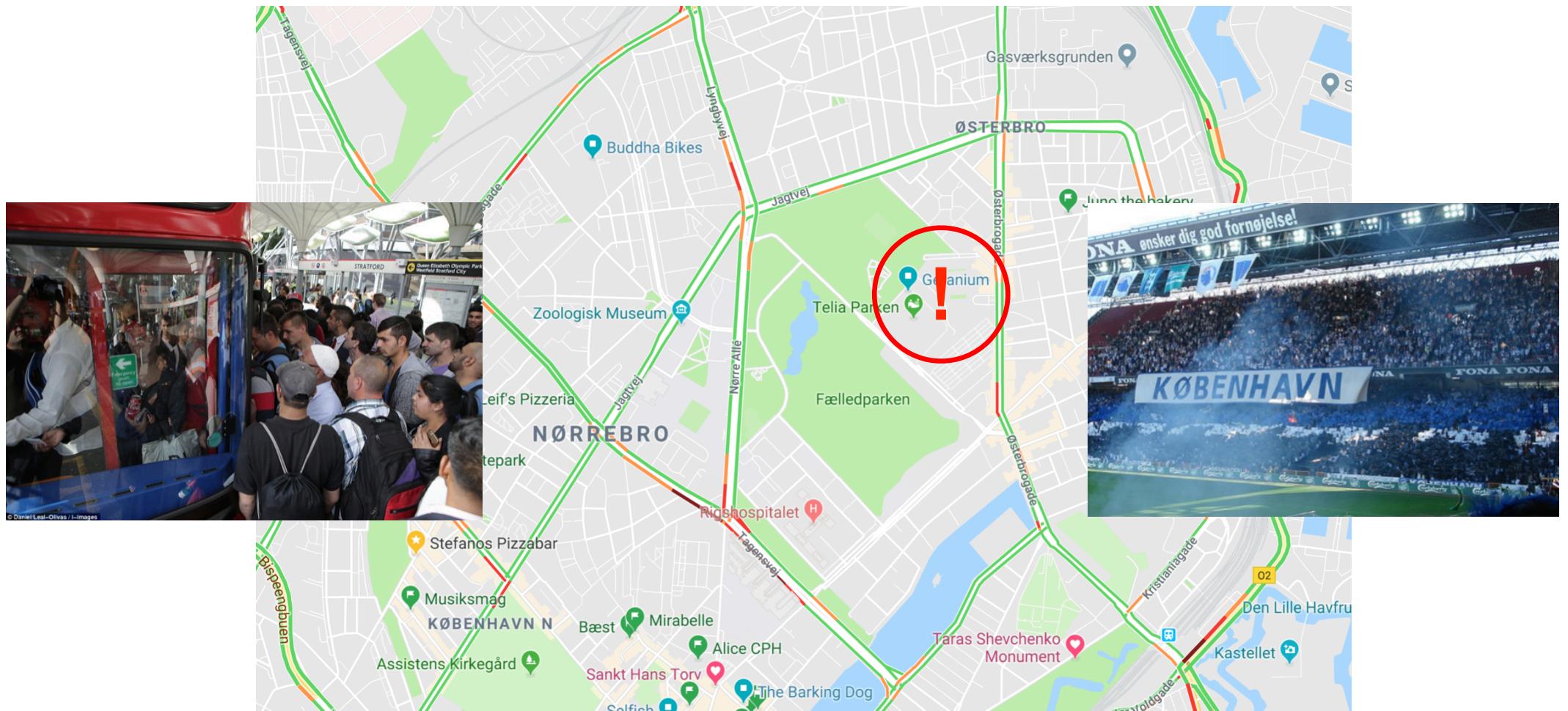
Transport demand problem: Average day



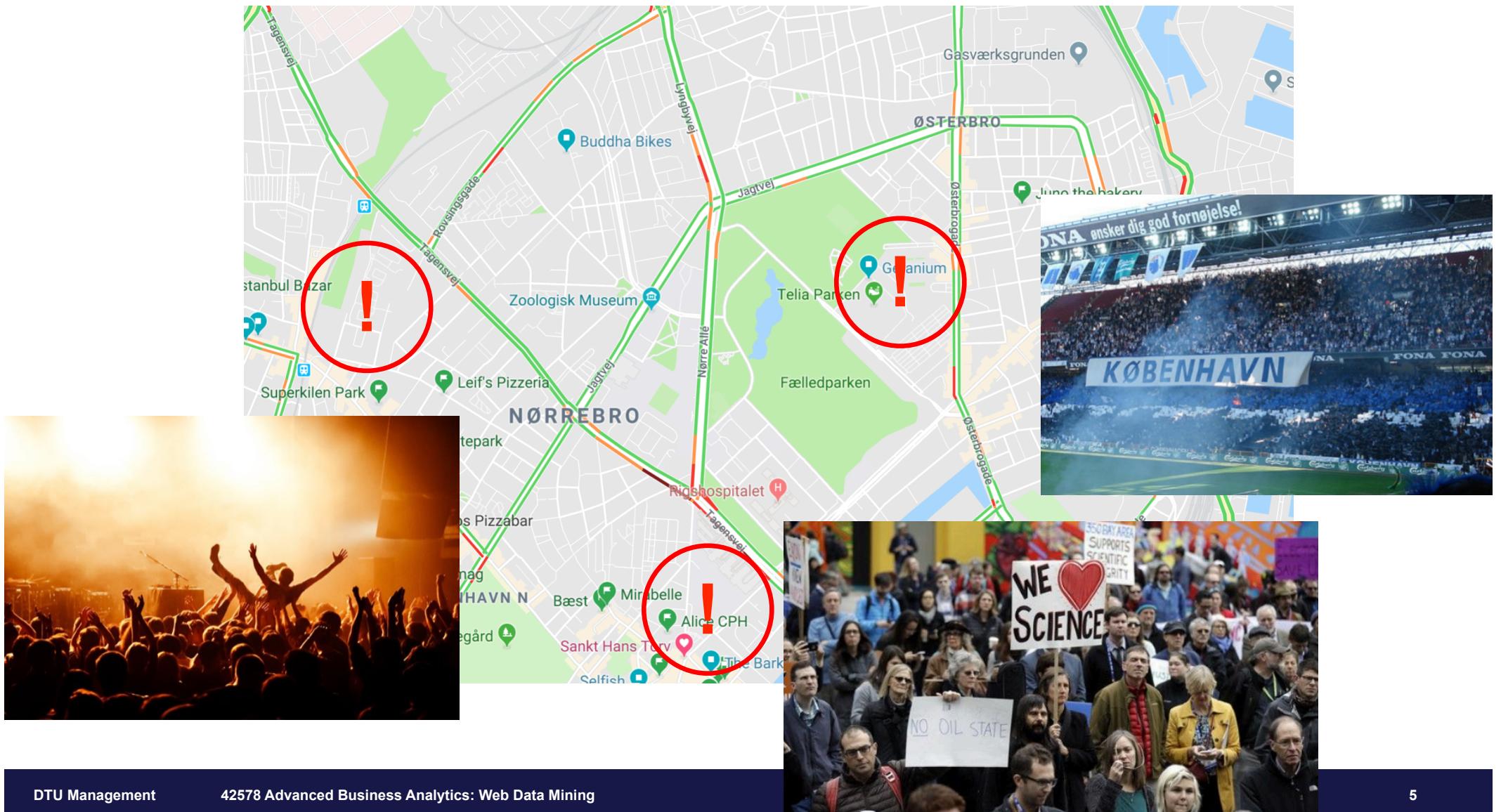
Transport demand problem: Large event



Transport demand problem: Large event



Transport demand problem: A lot of stuff...

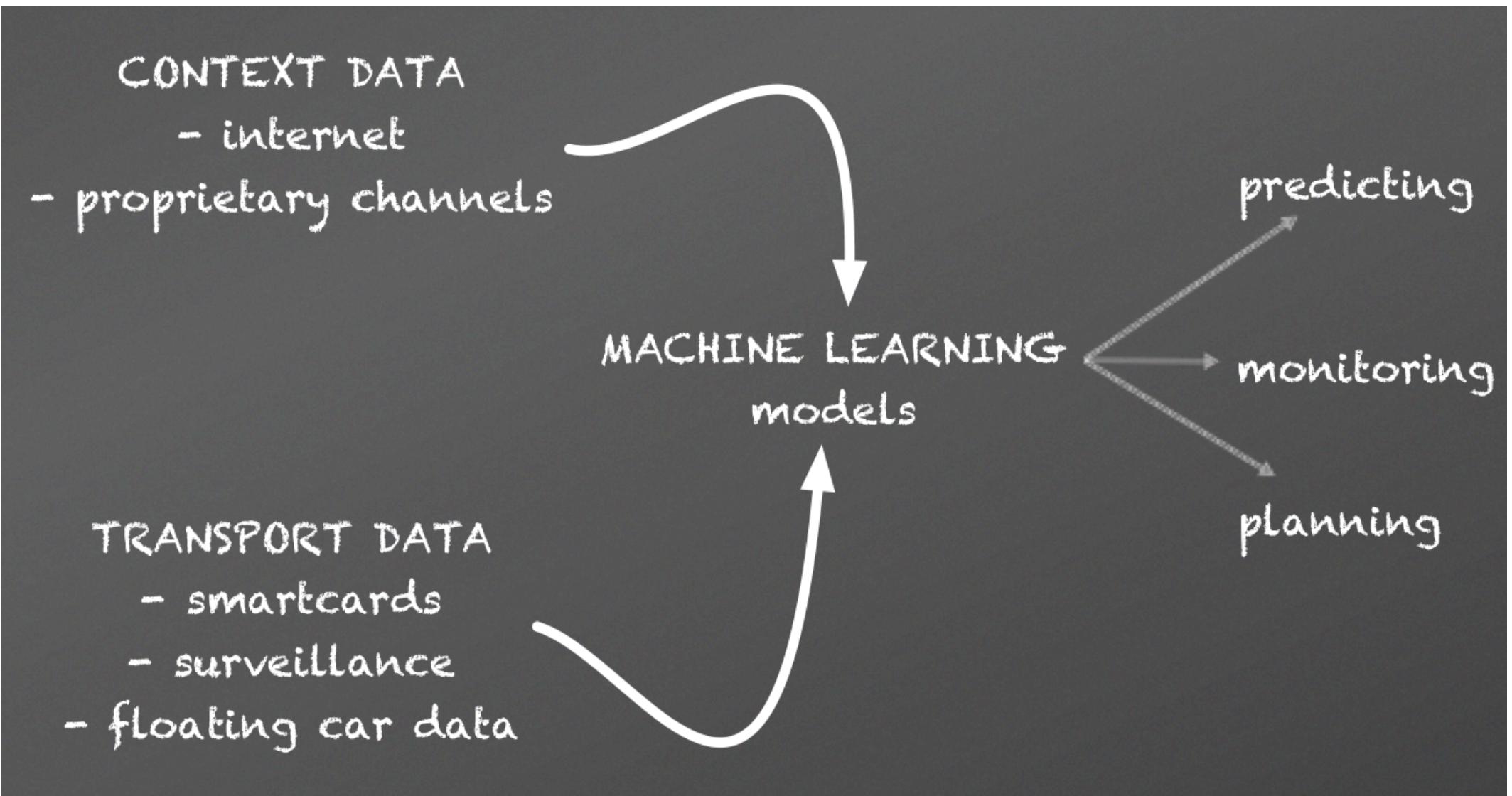




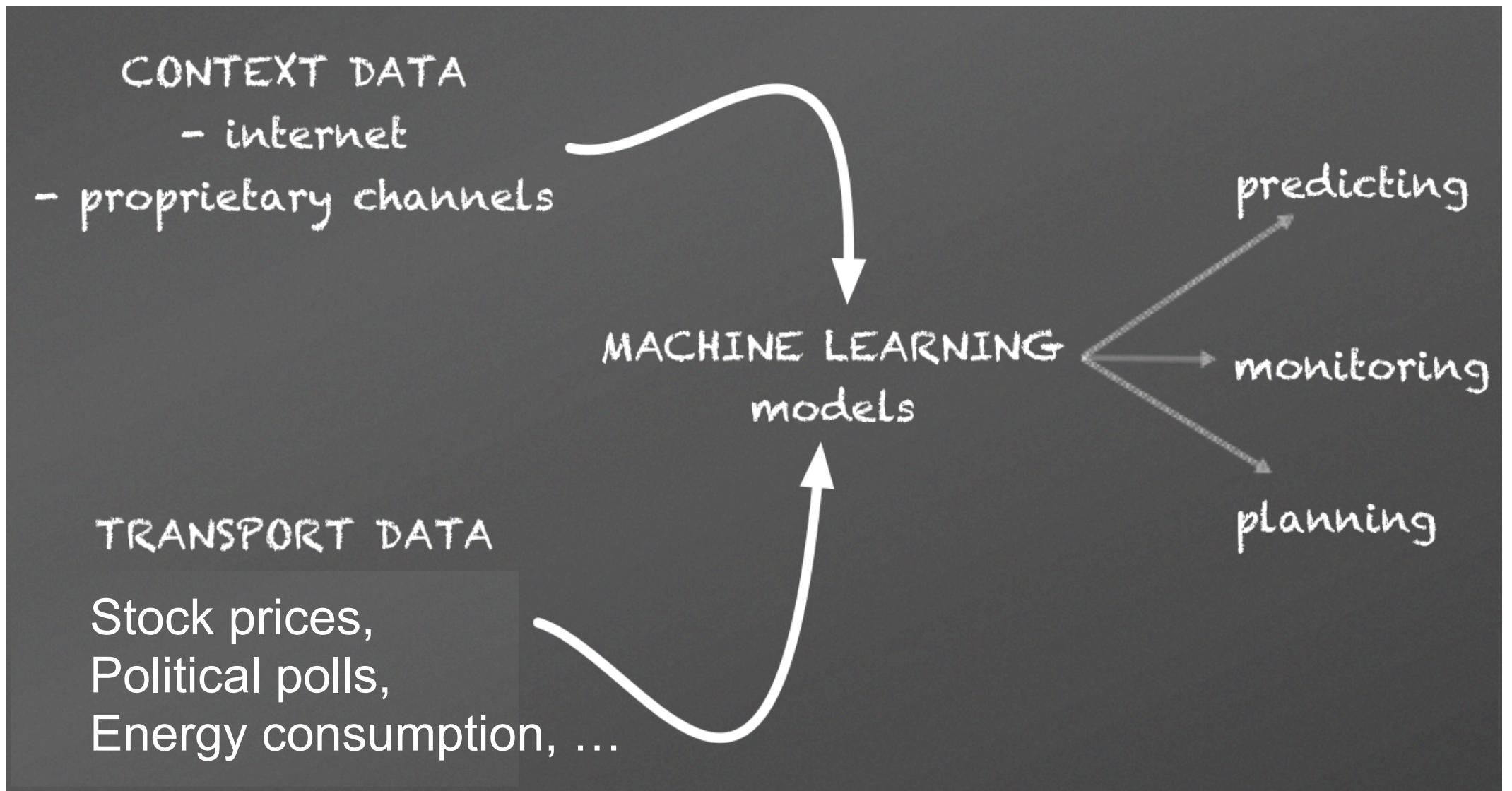
Where to get the information?

- Internet (obviously)
- ...

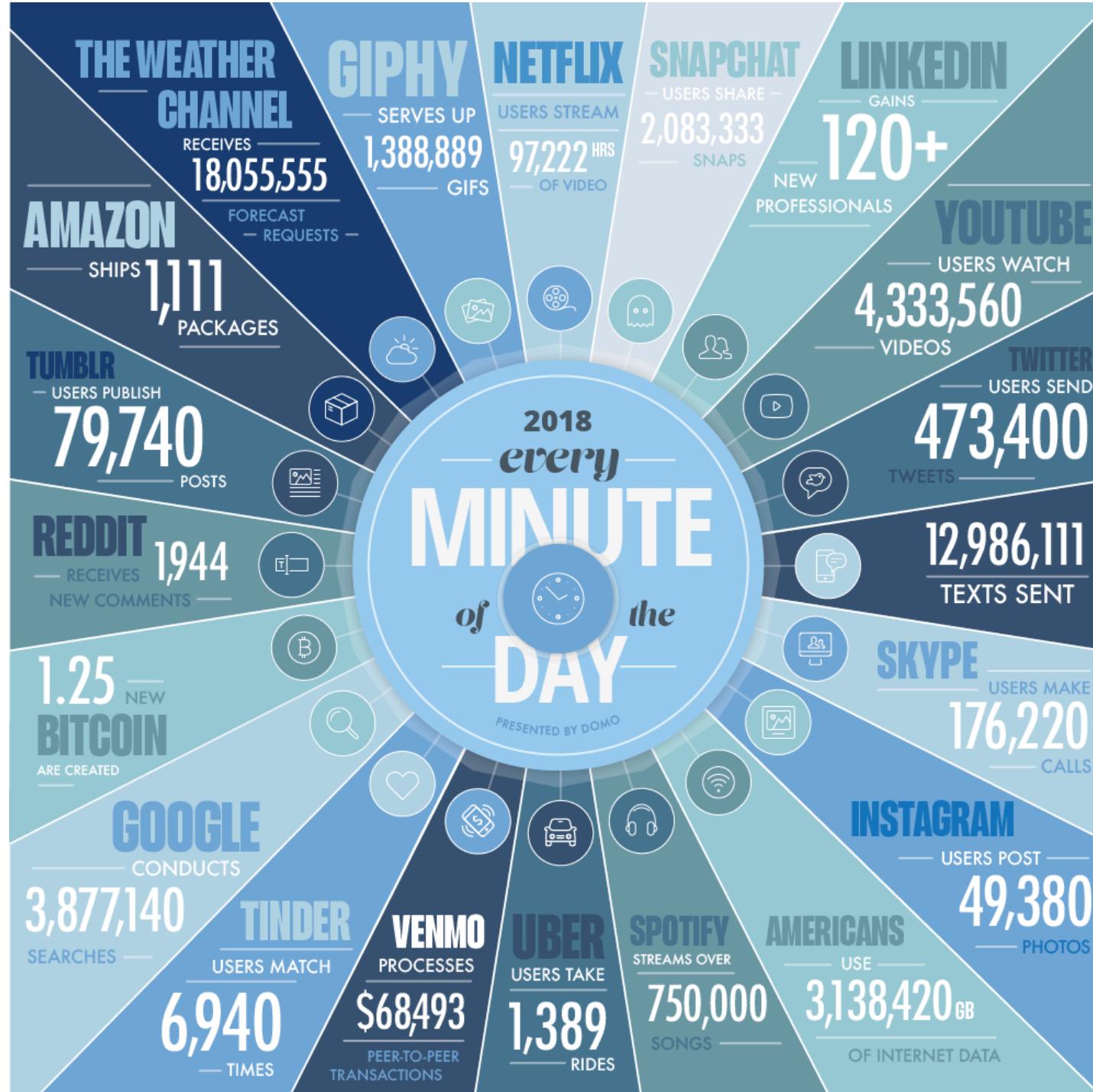
Predictions in Context

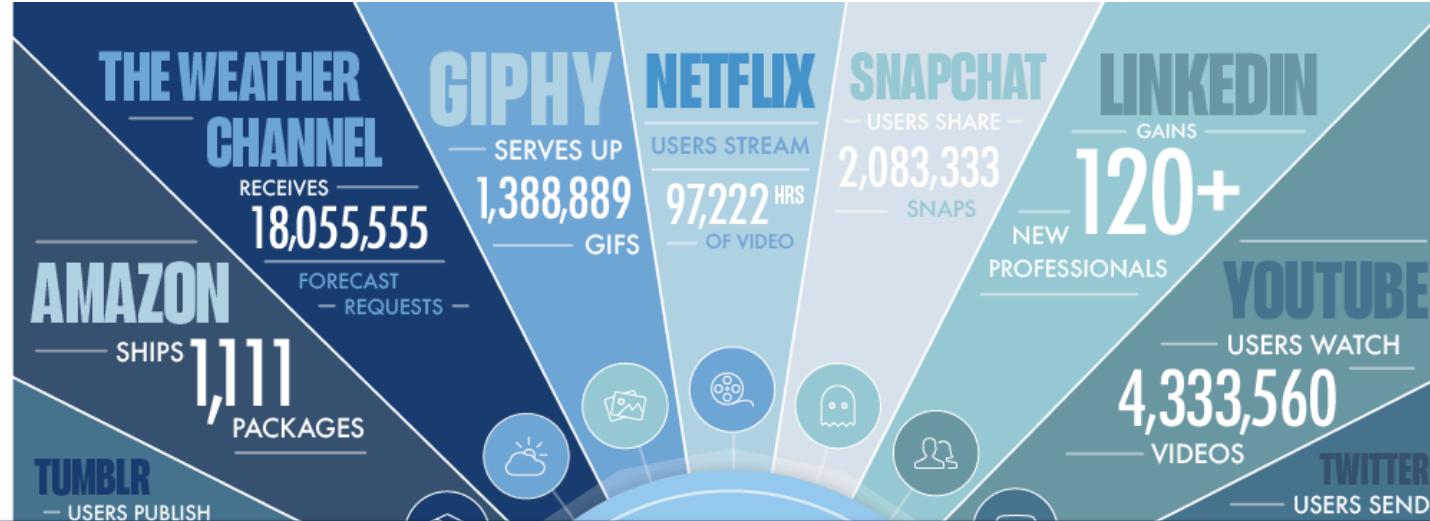


Predictions in Context



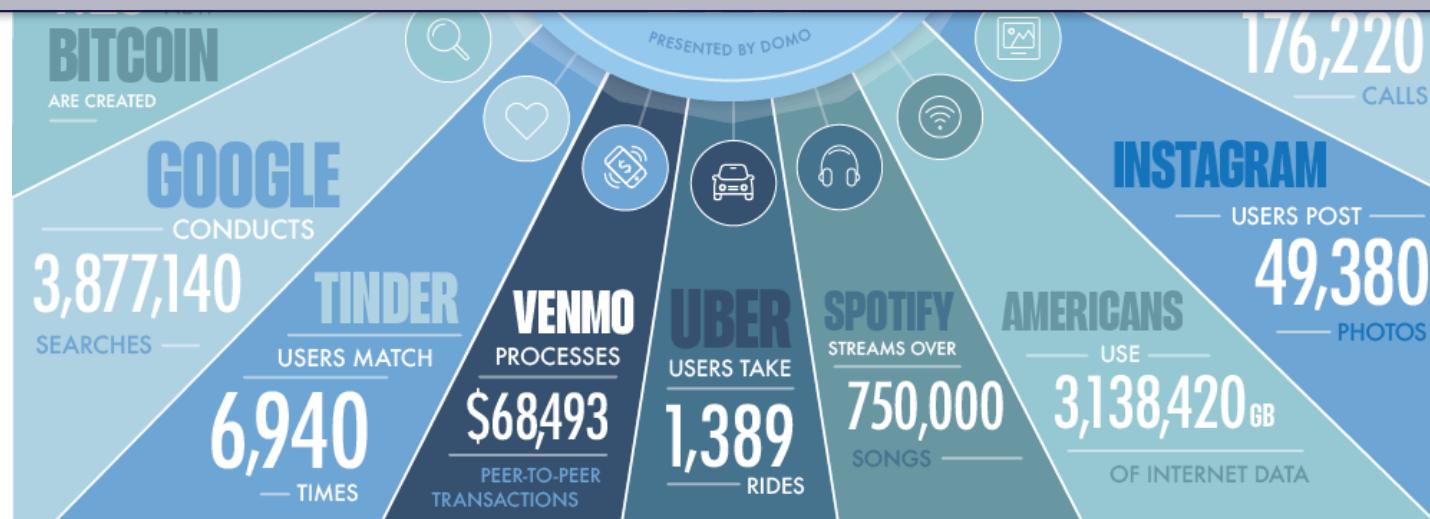
Let's focus on the Internet





A digital footprint of the humanity

(past, present, future,
facts, emotions, interpretations,
expectations, attitudes, plans, ...)



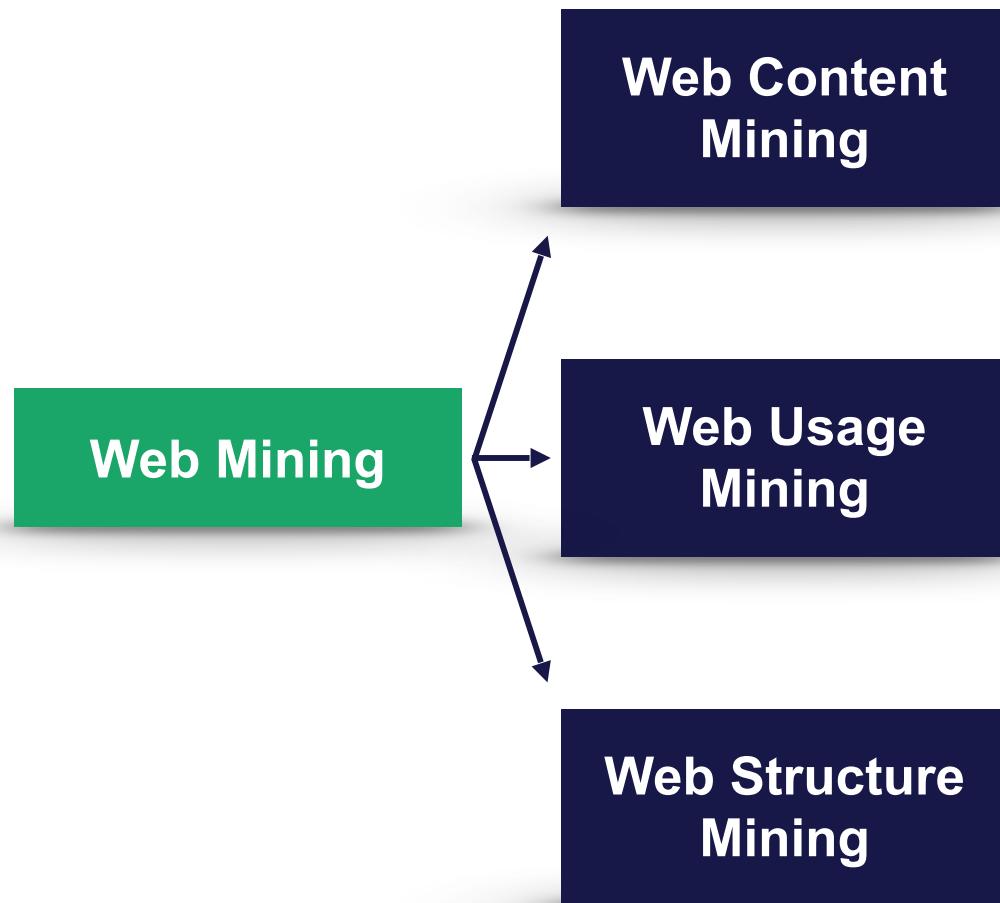
What is Web mining

- The potential of extracting valuable knowledge from the Web is quite evident
- Web mining is the collection of methods to fulfil this potential

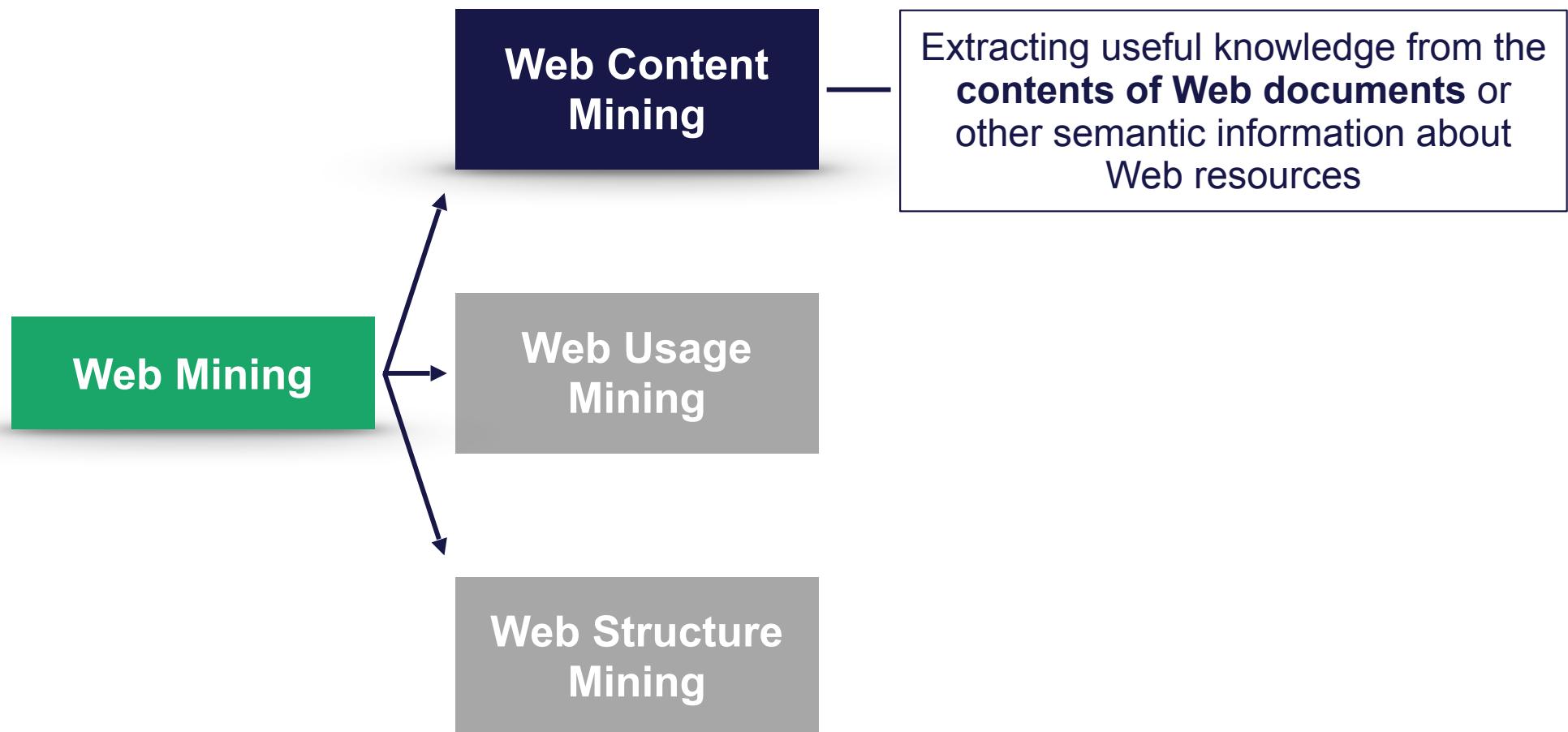
Web mining is application of *data mining* and *machine learning* techniques to extract useful knowledge from the content, structure, and usage of Web resources

- Nowadays, it is more relevant than at any other time during the history of the Web

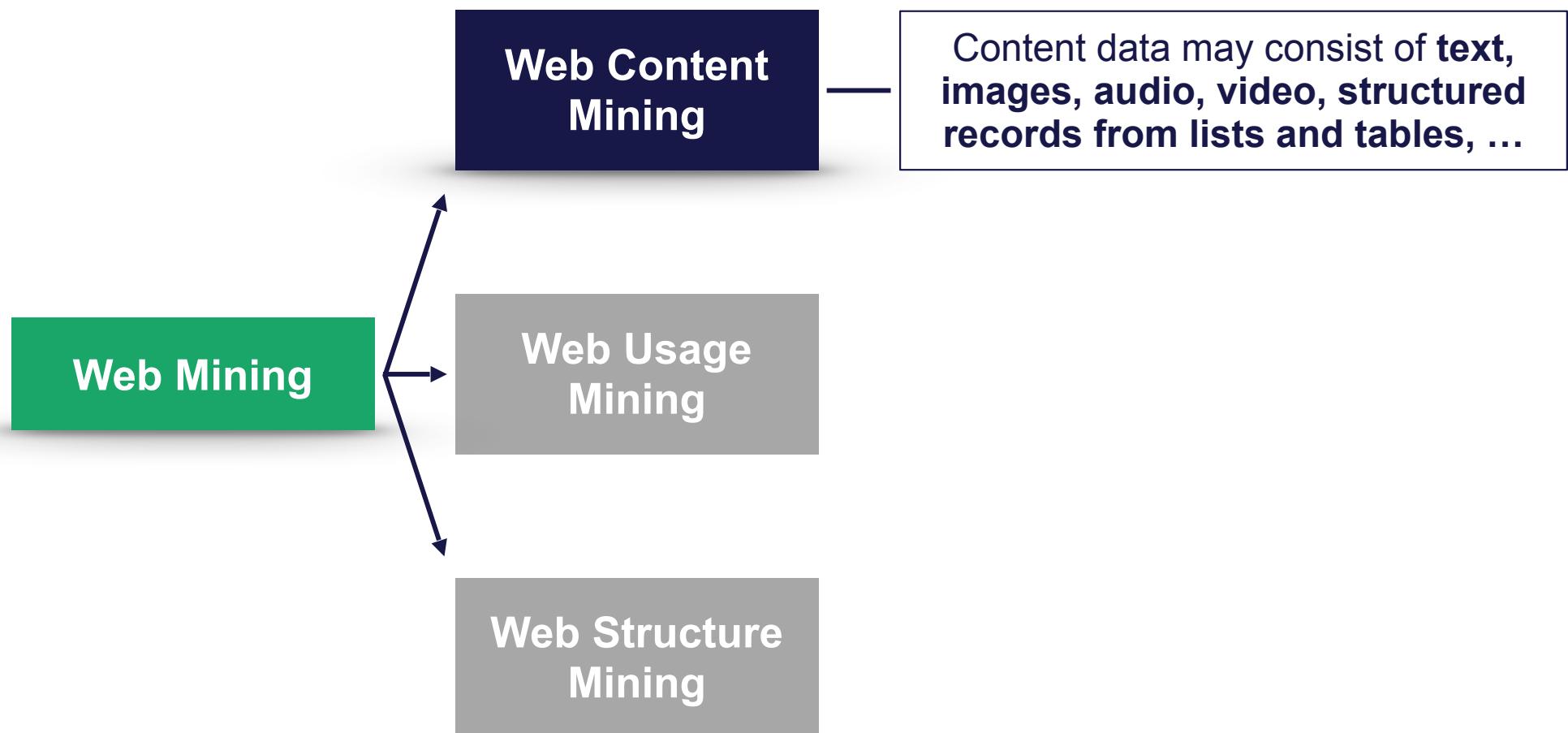
Types of Web Mining



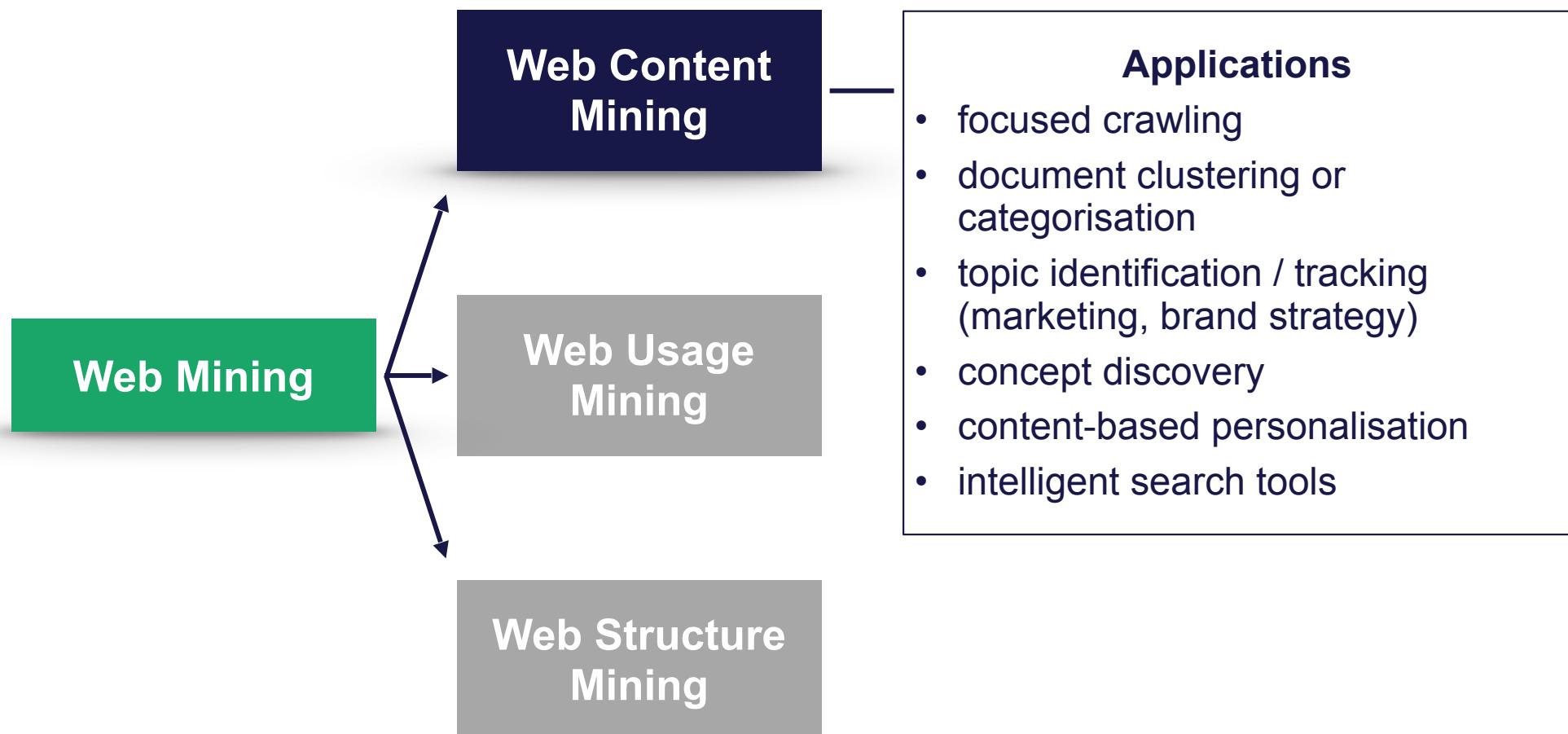
Types of Web Mining



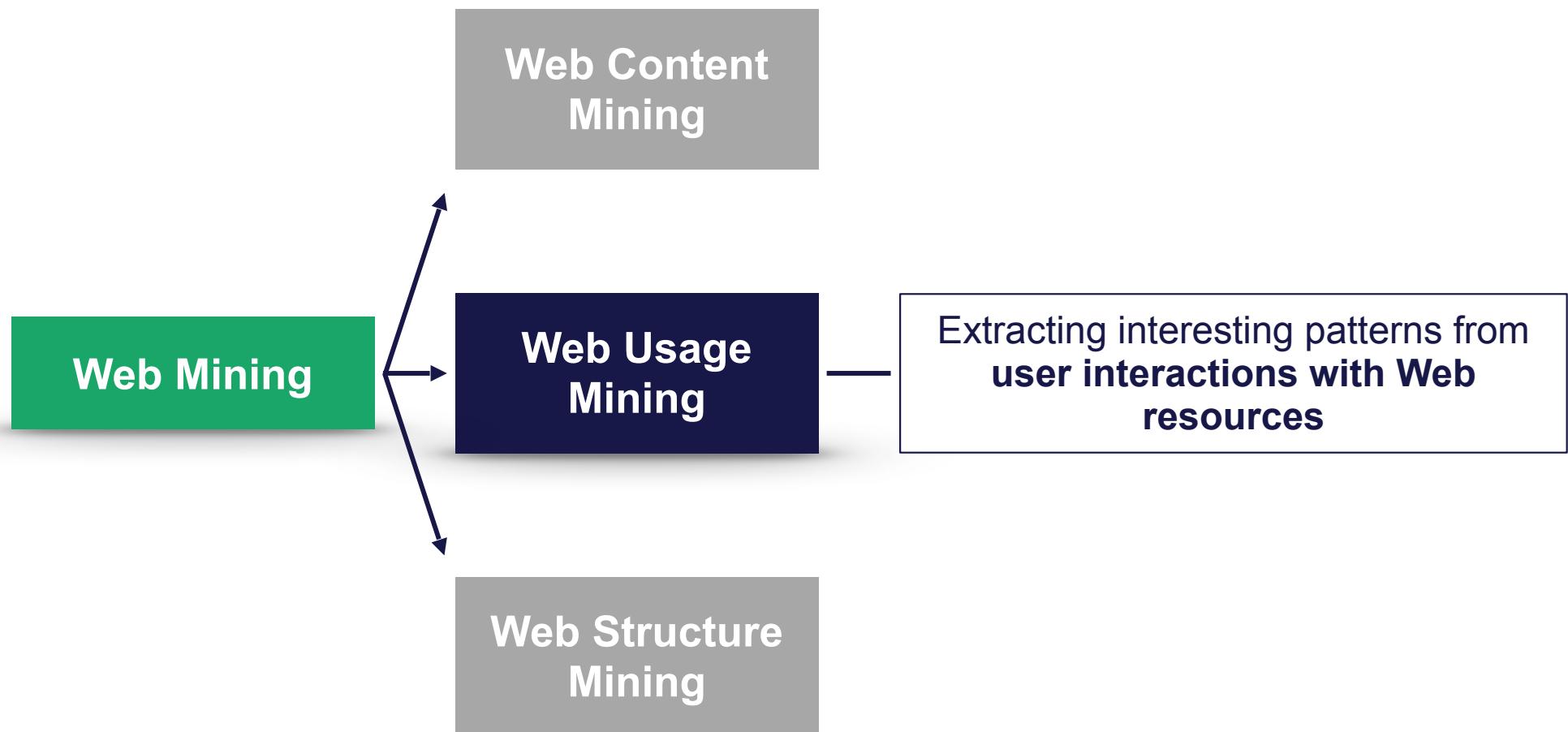
Types of Web Mining



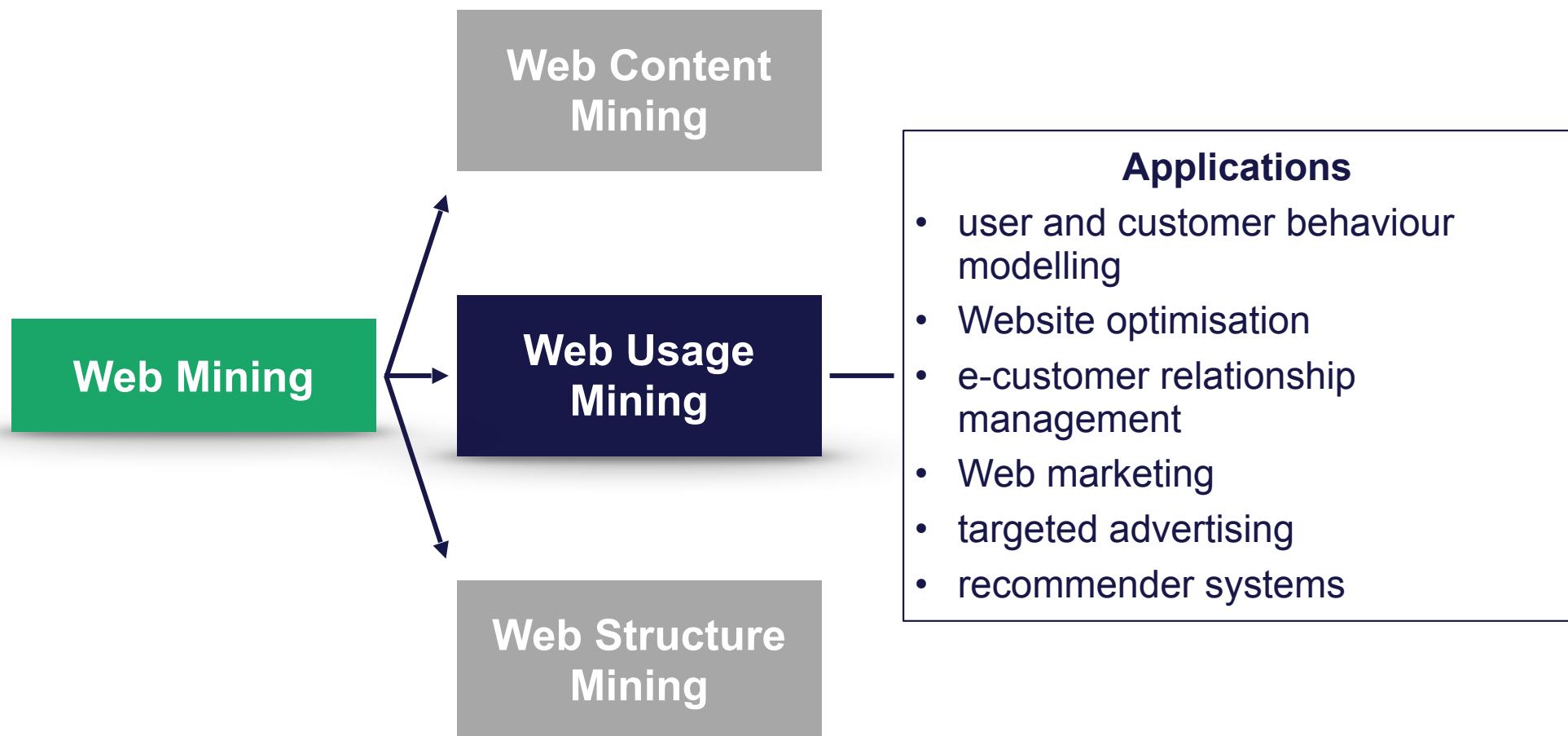
Types of Web Mining



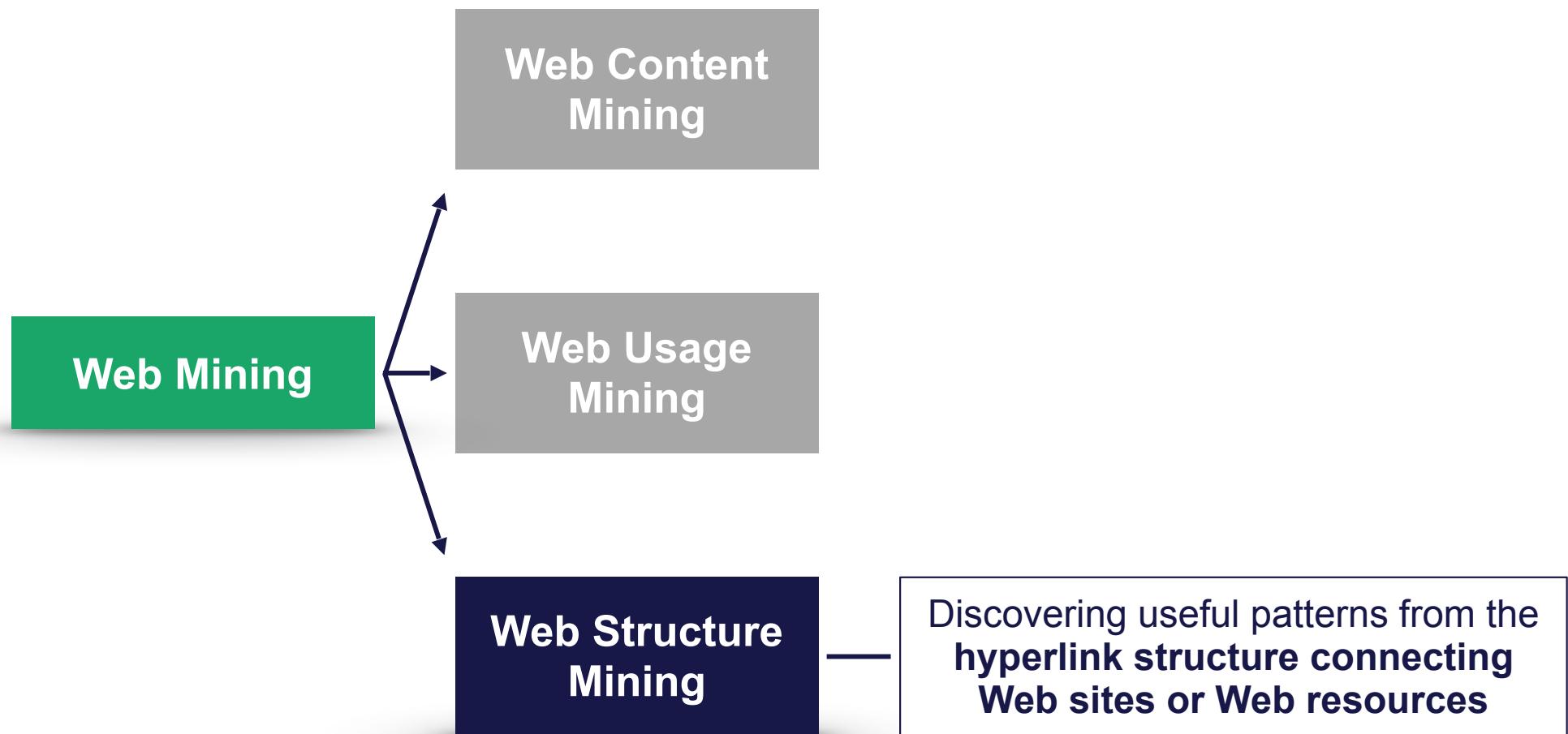
Types of Web Mining



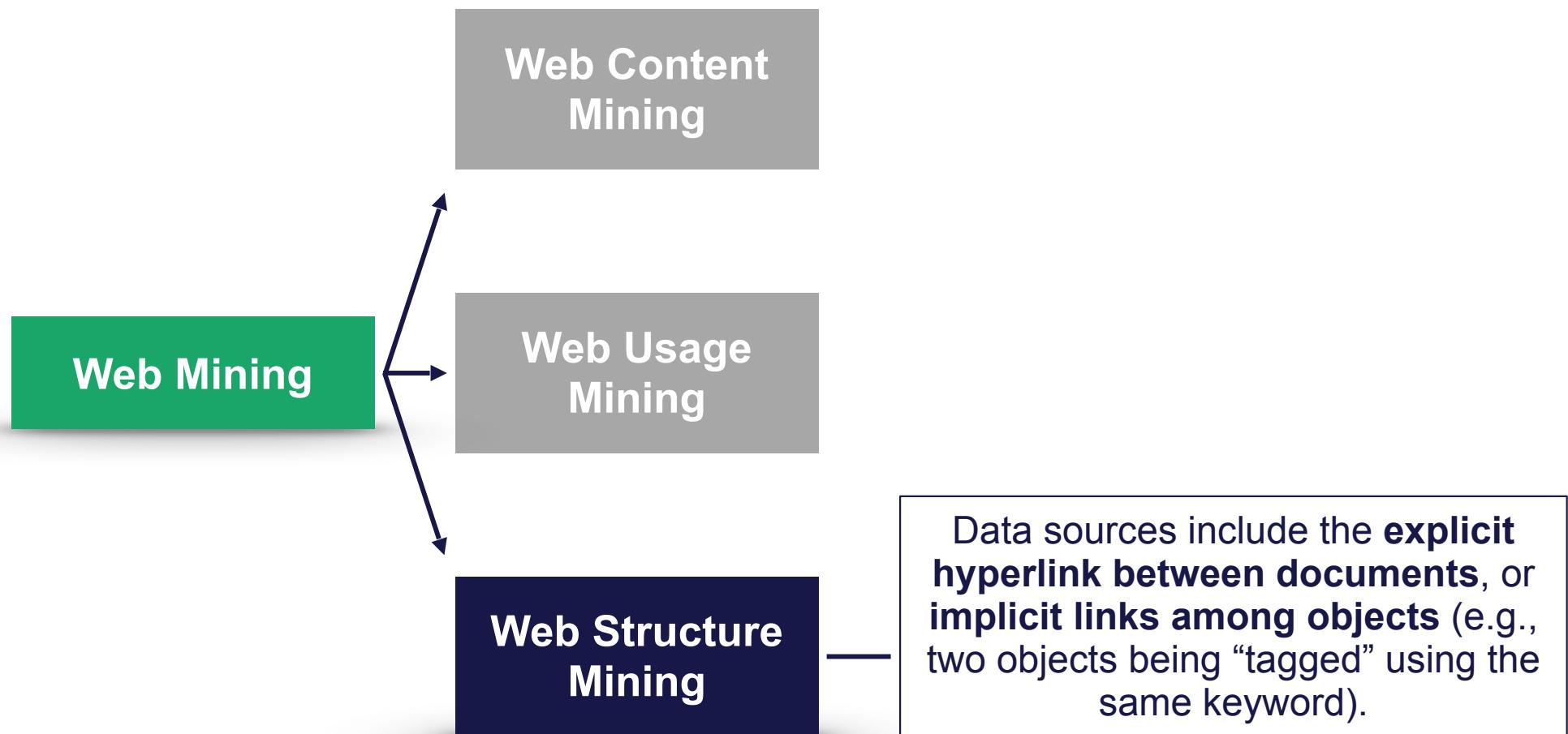
Types of Web Mining



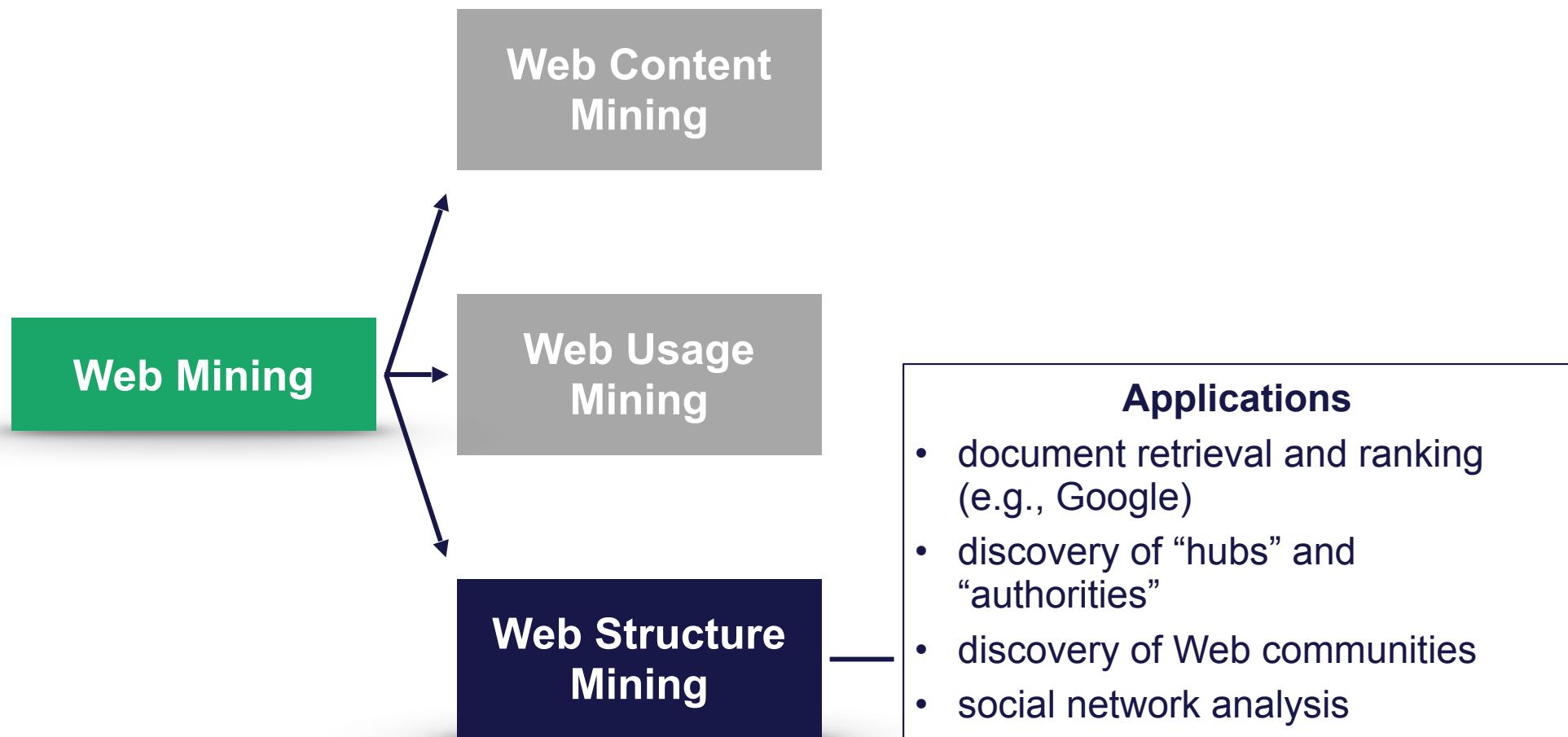
Types of Web Mining



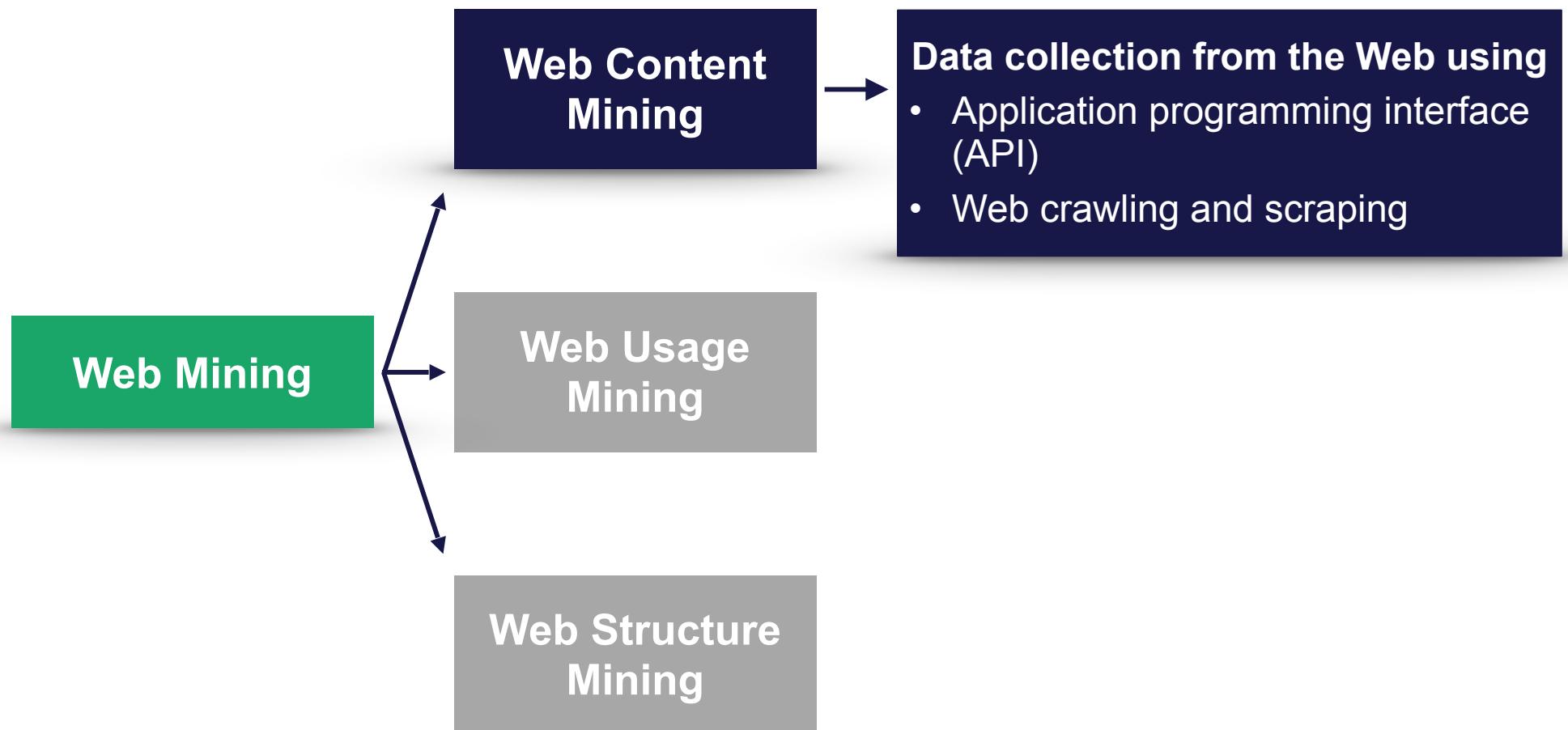
Types of Web Mining



Types of Web Mining



Today's focus: Web content mining



Transport demand problem: Context data

- Examples of the useful online sources:
 - ?

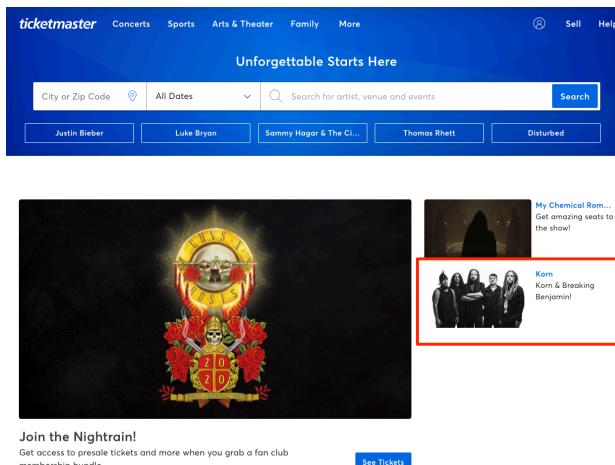
Transport demand problem: Context data

- **Examples of the useful online sources:**
 - Events announcements
 - Social networks
 - News
 - Waze, Google, Inrix
 - Weather
 - ...

Example: ticketmaster.com

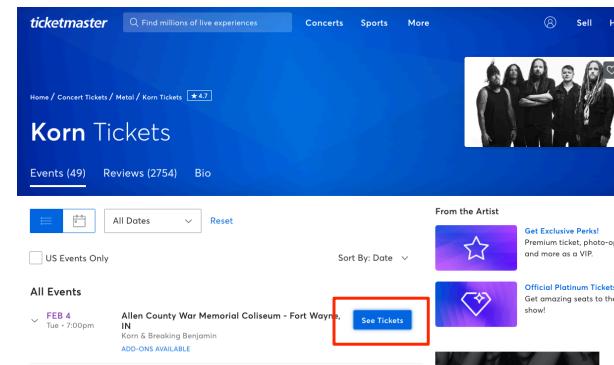
- Go to <https://www.ticketmaster.com/>
- Click on one of the events listed
- Click “see tickets”
- Click “more info” next to the event title in the top left corner

1



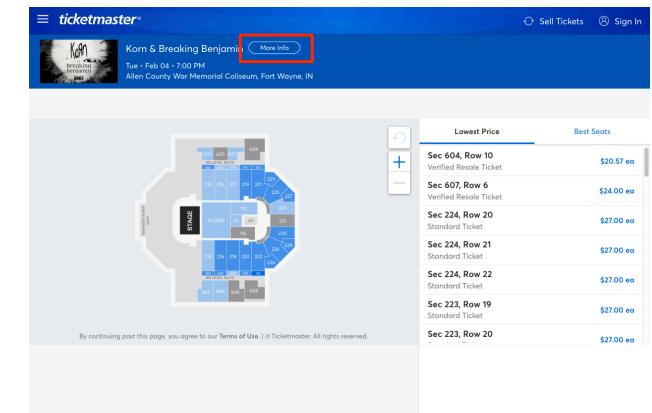
The screenshot shows the Ticketmaster homepage. At the top, there's a navigation bar with links for Concerts, Sports, Arts & Theater, Family, More, Sell, and Help. Below the navigation is a search bar with placeholder text "Unforgettable Starts Here". The search bar includes fields for "City or Zip Code" and "All Dates", and a "Search" button. Below the search bar is a grid of featured events. One event, "Korn & Breaking Benjamin", is highlighted with a red box around its thumbnail image.

2



The screenshot shows the "Korn Tickets" page. At the top, it says "Home / Concert Tickets / Metal / Korn Tickets" with a rating of 4.7. There's a thumbnail of the band Korn. Below the header, there are tabs for "Events (49)", "Reviews (2754)", and "Bio". The main content area shows a list of events. One event, "FEB 4 Tue • 7:00pm Allen County War Memorial Coliseum - Fort Wayne, IN Korn & Breaking Benjamin ADD-ONS AVAILABLE", has a "See Tickets" button highlighted with a red box. To the right, there are sections for "From the Artist" (with options for "Get Exclusive Perks!" and "Official Platinum Tickets") and "Sort By: Date".

3



The screenshot shows the event details page for "Korn & Breaking Benjamin" on February 4th at the Allen County War Memorial Coliseum. At the top, it shows the event title, date, and location. Below that is a seating chart of the venue. To the right, there's a table of ticket prices for different sections and rows. The table includes columns for "Lowest Price", "Best Seats", and ticket descriptions like "Verified Resale Ticket" and "Standard Ticket".

Example: ticketmaster.com

- A lot of useful information... How to get (“extract”) it?

The screenshot shows a modal window titled "Event Info" for an event at the Allen County War Memorial Coliseum. The window is divided into several sections:

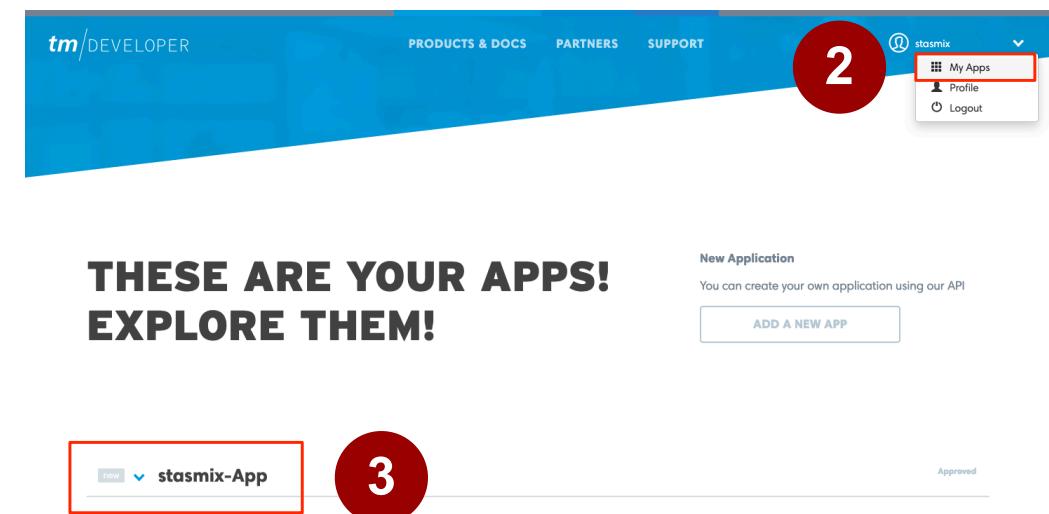
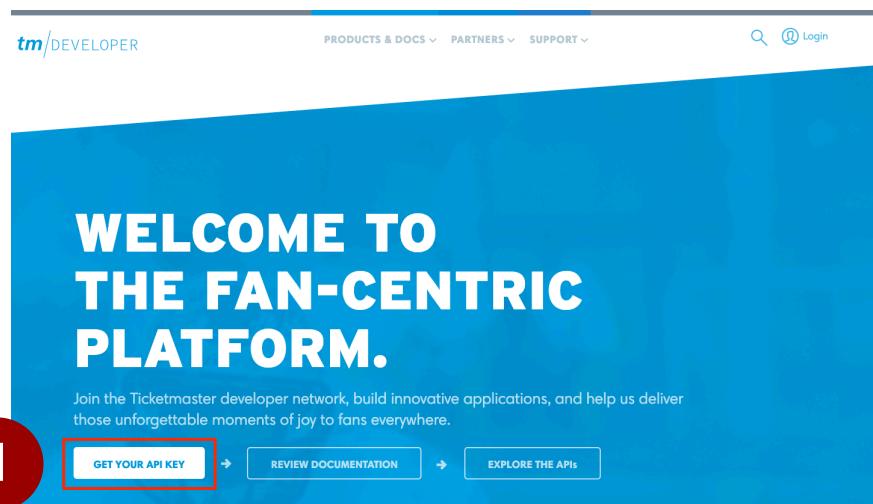
- Type:** Hard Rock/Metal (highlighted by a red box)
- Title:** Korn & Breaking Benjamin (highlighted by a red box)
- Date / time:** Tue • Feb 04 • 7:00 PM (highlighted by a red box)
- Venue:** Allen County War Memorial Coliseum
Fort Wayne, IN (highlighted by a red box)
- Performers:** Lineup (highlighted by a red box) featuring Korn, Breaking Benjamin, and Bones UK.
- Ticket Limits:** An 8 ticket limit for this event.
- Accessible Tickets:** Information about accessible seating tickets.

Application programming interface (API)

- The easiest way to get data from the Internet 😊
- Data is structured or semi-structured
- **Steps:**
 1. Make sure the data source provide APIs for data collection.
 2. Obtain API key or other forms of authorisation.
 3. Read documentation
 4. Coding

API Example: ticketmaster.com

- Go to <https://developer.ticketmaster.com/>
- Click on “GET YOUR API KEY”
- Sign up
(create a new email address if you don’t want to use your own one)
- Log in, go to “My Apps” and click on your first automatically created app



API Example: ticketmaster.com

- Here you can see your credentials to make API requests and other information such as usage quotas
- **Always keep your credentials private!**

new ▾ stasmix-App Approved

Keys Details Affiliate IDs ANALYTICS DELETE EDIT

Below are keys you can use to access the API products associated with this application (stasmix-App). The actual keys need to be approved and approved for an API product to be capable of accessing any of the URLs defined in the API product.

Consumer Key	AN24mrb1iTUGU [REDACTED]
Consumer Secret	Wko1s [REDACTED]
Key Issued	Tue, 02/04/2020 - 14:26
Key Expires	Never
Product	Status
Public APIs 5000 requests every 1 day	Approved
OAuth Product 100 requests every 1 minute	Approved

API Example: ticketmaster.com

- Go to the documentation page (we will use **Discovery API**)
<https://developer.ticketmaster.com/products-and-docs/apis/discovery-api/v2/>
- There are many ways to authorise yourself and use API
- We will use the simplest “GET” method which corresponds to a plain URL
- A GET request URL has the simple structure:
url_base ? parameter1=value1 & parameter2=value2 ...
- Scroll down to “Examples” below the map and copy-paste it to a new tab
(your API key is automatically appended in the examples)

url_base

Get a list of all events in the United States [https://app.ticketmaster.com/discovery/v2/events.json?
countryCode=US&apikey=AN24mrN6](https://app.ticketmaster.com/discovery/v2/events.json?countryCode=US&apikey=AN24mrN6)

parameter1 parameter2 (your API customer secret)

API Example: ticketmaster.com

- Congratulations! You will get a structured response in the JSON format

```
{
  "embedded": {
    "events": [
      {
        "name": "93 ENT Birthday Bash",
        "type": "event",
        "id": "C5eV2pdR-BGMM",
        "test": false,
        "url": "https://www.ticketmaster.com/t3-ent-birthday-bash-columbia-south-carolina-03-2020/event/2000583B818C8CD7",
        "locale": "en-us",
        "images": [
          {
            "ratio": "16_9",
            "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_RETINA_PORTAIT_16_9.jpg",
            "width": 640,
            "height": 160,
            "fallback": false
          },
          {
            "ratio": "4_3",
            "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_RETINA_CUSTOM.jpg",
            "width": 305,
            "height": 225,
            "fallback": false
          },
          {
            "ratio": "16_9",
            "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_RETINA_LANDSCAPE_16_9.jpg",
            "width": 1136,
            "height": 639,
            "fallback": false
          }
        ],
        "ratio": "16_9",
        "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_RECOMMENDATION_16_9.jpg",
        "width": 100,
        "height": 56,
        "fallback": false
      },
      {
        "name": "EVENT_DETAIL PAGE_16_9",
        "type": "event_detail",
        "id": "495231_ARTIST_PAGE_16_9",
        "test": true,
        "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_ARTIST_PAGE_16_9.jpg",
        "locale": "en-us",
        "images": [
          {
            "ratio": "3_2",
            "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_ARTIST_PAGE_3_2.jpg",
            "width": 140,
            "height": 115,
            "fallback": false
          }
        ],
        "ratio": "3_2",
        "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_ARTIST_PAGE_3_2.jpg",
        "width": 1427,
        "height": 1152,
        "fallback": false
      },
      {
        "name": "ARTIST_PAGE_3_2",
        "type": "event_detail",
        "id": "495231_ARTIST_PAGE_3_2",
        "test": false,
        "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_ARTIST_PAGE_3_2.jpg",
        "locale": "en-us",
        "images": [
          {
            "ratio": "3_2",
            "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_ARTIST_PAGE_3_2.jpg"
          }
        ],
        "ratio": "3_2",
        "url": "https://sl.ticketmaster.net/dam/a/497/bdbfdcl1-f05d-437f-a6fe-e35c8d587497_495231_ARTIST_PAGE_3_2.jpg",
        "width": 203,
        "height": 160,
        "fallback": false
      }
    ]
  }
}
```

The JSON response contains a single event object with various details such as name, type, id, test status, URL, locale, images, and ratios. The images section lists multiple image URLs for different aspect ratios (16:9, 4:3, and 3:2). The event also includes a detailed page (ARTIST_PAGE_16_9) and a specific page for the artist (ARTIST_PAGE_3_2).

- To make it more human readable, you can either use JSON formatter such as <https://jsonformatter.curiousconcept.com/> or directly use Python (next slide)

API Example: ticketmaster.com

```
In [1]: import requests

api_key = "AN241...gN6"
# It's a very bad programming practice
# to put passwords or secret keys into the code...
url_base = "https://app.ticketmaster.com/discovery/v2/events.json"
params = ["countryCode=US", "apikey="+api_key]

request_url = url_base + "?" + "&".join(params)

response = requests.get(request_url)
data = response.json()
data
```

```
Out[1]: {'_embedded': {'events': [{'_embedded': {'attractions': [{'_links': {'self': {'href': '/discovery/v2/attractions/K8vZ917pkBV?locale=en-us'}}}, 'aliases': ['NBA Youngboy'], 'classifications': [{"family": False, 'genre': {'id': 'KnvZfZ7vAv1', 'name': 'Hip-Hop/Rap'}, 'primary': True, 'segment': {'id': 'KZFzniwnSyZfZ7v7nJ', 'name': 'Music'}, 'subGenre': {'id': 'KZazBEonSMnZfZ7vkvl', 'name': 'Hip-Hop/Rap'}, 'subType': {'id': 'KZFzBERXgnZfZ7vAd7', 'name': 'Musician'}, 'type': {'id': 'KZAyXgnZfZ7v7la', 'name': 'Individual'}}], 'externalLinks': {'homepage': [{'url': 'http://youngboynba.com/'}, {'musicbrainz': [{"id": "127d8614-0876-4408-886c-6f110b047a78"}], 'twitter': [{"url": "https://twitter.com/ggyoungboy"}]}], 'id': 'K8vZ917pkBV', 'images': [{"fallback": False, 'height': 360, 'ratio': '16_9', '...': '...', 'venues': [{"_links": {"self": {"href": '/discovery/v2/venues/KovZpZAE7dtA?locale=en-us'}}}, '...': '...', 'address': {"line1": "801 Lincoln St."}, 'aliases': ['carolina center', 'colonial center'], '...': '...']}]}
```

JSON format

One of the most popular web data formats

```
{  
    "var_1": "string_val",  
    "var_2": 1.56,  
    "var_3": true,  
    "var_4": [1, 2, 3],  
    "var_5": {  
        "var_6": ["a", "b", "c"]  
    }  
}
```

Can be parsed to a Python dictionary using json library

```
import json
```

APIs: Summary



- parameters
- authorization

- JSON
- XML
- ...

Playtime!

- Take a look at the “api_example.ipynb”
- Do Part 1 - Twitter.ipynb

and/or

- Do Part 2 - Weather.ipynb

APIs: Problems

- No API is provided
- Free access is limited
- Response does not contain all the information

API Problems Example: [eventbrite.com](#)

“Personal anecdote” from Stanislav Borysov (2020 spring semester):

- Last semester, we used [eventbrite.com](#) API as an example
- However, when I checked it again last week, I found that the event search API was shut down on Thursday, December 12, 2019.
- There are still some API functionality left but maybe it is not enough for our purposes.
- Also, maybe [eventbrite.com](#) has more events and provides better event description, so we would like to use it instead of [ticketmaster.com](#)
- **What to do?**

API Problems Example: [eventbrite.com](#)

- Last semester, we used [eventbrite.com](#) API as an example
- However, when I checked it again last week, I found that the event search API was shut down on Thursday, December 12, 2019.
- There are still some API functionality left but maybe it is not enough for our purposes.
- Also, maybe [eventbrite.com](#) has more events and provides better event description, so we would like to use it instead of [ticketmaster.com](#)
- **What to do?**

Generic web crawler / scraper!

Web scraping: Overall idea

- After browsing the [eventbrite.com](https://www.eventbrite.com) website, you find out that each event's information can be found at https://www.eventbrite.dk/e/event_id where each event has its own unique id *event_id*
- **Pseudo code:**

get a list of event ids {ids}, for example, from the main page [eventbrite.com](https://www.eventbrite.com)

for each **event_id** in {ids}

 access https://www.eventbrite.dk/e/event_id, store page content in **C**

 extract event name **n**, date **d**, time **t**, etc from **C**

 store (**event_id**, **n**, **d**, **t**) in database

Web crawling

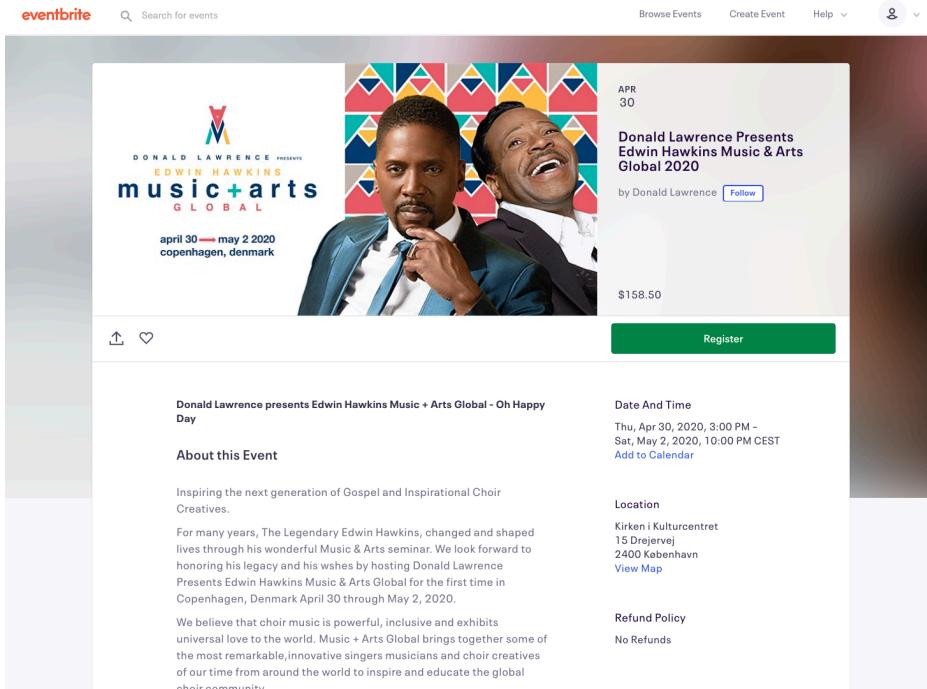
Web scraping

Example: Eventbrite.com

- Go to eventbrite.com
- Click on one of the events listed
- View the webpage HTML source
 - **Chrome**
 - *Right click → View page source*
 - *View → Developer → View source*
 - **Safari**
 1. *Preferences → Advanced → Show develop menu*
 2. *Develop → Show page source*
 - **Firefox**
Tools → Web developer → Page source
 - **Edge**
Right click → View source

Example: Eventbrite.com

What you see



What your browser sees

XML/HTML format

```
<!DOCTYPE html >

  <head>
    <meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
  </head>
  <body>
    <div id="header">
      <ul id="A">
        <li><a href="#">A link</a></li>
      </ul>
      <div id="B"><span class="about">About something</span></div>
    </div>

    <div id="main">Main container...</div>

    <p id="footer">This is the footer. <a href="#">Yet another link</a></p>
  </body>
</html>
```

Web crawlers: Selected useful Python methods

- Read html files:

```
import urllib.request
contents = urllib.request.urlopen("http://example.com/foo/bar").read()
#returns a string containing the source html content for the url
```

- Regex that finds specific patterns in a text:

```
import re
results = re.findall(pattern, contents)
for result in results:
    # do something
#from contents, find string(s) that matches the pattern specified by regex
```

- Wait for several seconds to reduce the risks of being detected and banned

```
import time
time.sleep(5)
#wait for 5 seconds.
```

Regular Expression

- Regex - An advanced search.
 - “Normal search” only deals with finding fixed character sequences.
 - Regex can handle various patterns.
- An interactive tutorial:
 - <http://regexone.com/>
- A place to quickly test a written regex against a source text:
 - <http://regexpal.com/>

Regular Expression

Regular- Expression Patterns

- ^ Matches beginning of line.
- \$ Matches end of line.
- . Matches any single char except newline.
- [...] Matches any single char in brackets.
- [^...] Matches any single char not in brackets.
- \w Matches word characters.
- \W Matches nonword characters.
- \s Matches whitespace.
- \S Matches nonwhitespace.
- \d Matches digits.
- \D Matches nondigits.
- \A Matches beginning of string.
- \Z Matches end of string.
- \z Matches end of string.
- \G Matches point where last match finished.
- x|y Matches either x or y.

- [0-9] Match any digit; same as [0123456789]
- [a-z] Match any lowercase ASCII letter
- [A-Z] Match any uppercase ASCII letter
- [a-zA-Z0-9] Match any of the above
- [^aeiou] Match any other than a lowercase vowel
- [^0-9] Match anything other than a digit.

Special symbols:

- * 0 or more repetitions of the previous symbol
- + 1 or more repetitions of the previous symbol
- ? Optional character

The most useful one for web crawlers

<tag>(.*)?</tag>

match everything surrounded by
<tag><tags>

Regular Expression: Example

HTML content:

```
<div class="txt-block" itemprop="actors" itemscope itemtype="http://schema.org/Person">
  <h4 class="inline">Stars:</h4>
  <name size=3>Ben Ziegler</name>,
  <name size=5>Glenna Hill</name>,
  <name size=4>Jason Woolfolk</name>
  <span class="ghost">|</span>
  <span class="see-more inline nobr">
    <a href="fullcredits?ref_=tt_ov_st_sm" itemprop='url'> See full cast and crew</a> &raquo;
  </span>
</div>
```

Regular Expression: Example

- Match the three names surrounded by <name> tags:
 - <name size=\d>(.*)</name>

The screenshot shows a regular expression testing interface. At the top, there is a search bar containing the regular expression pattern: <name size=\d>(.*)</name>. Below the search bar, the results panel displays a portion of an HTML document. The names 'Ben Ziegler', 'Glenna Hill', and 'Jason Woolfolk' are highlighted in yellow, indicating they are matches for the regular expression. The rest of the HTML code is shown in a monospaced font.

```
<name size=\d>(.*)</name>

<div class="txt-block" itemprop="actors" itemscope itemtype="http://schema.org/Person">
    <h4 class="inline">Stars:</h4>
    <name size=3>Ben Ziegler</name>,
    <name size=5>Glenna Hill</name>,
    <name size=4>Jason Woolfolk</name>
        <span class="ghost">|</span>
        <span class="see-more inline nobr">
            <a href="fullcredits?ref_=tt_ov_st_sm" itemprop='url'> See full cast and crew</a> &raquo;
        </span>
    </div>
```

Regular Expression: Example

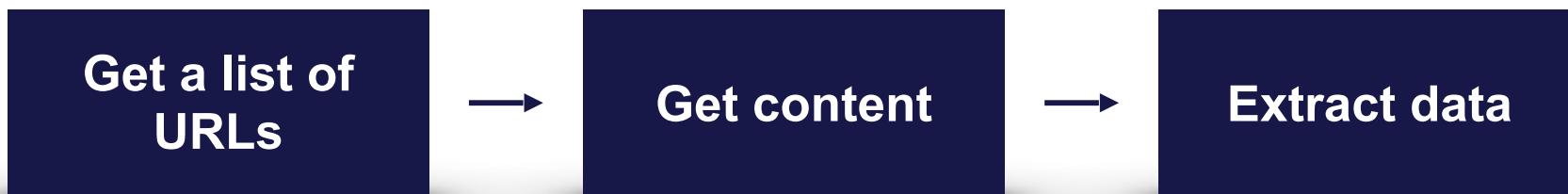
```
import re

html_content = """
<div class="txt-block" itemprop="actors" itemscope>
    <h4 class="inline">Stars:</h4>
    <name size=3>Ben Ziegler</name>,
    <name size=5>Glenna Hill</name>,
    <name size=4>Jason Woolfolk</name>
    <span class="ghost">|</span>
    <span class="see-more inline nobr">
        <a href="fullcredits?ref_=tt_ov_st_sm" itemprop='url'>
            See full cast and crew
        </a> &raquo;
    </span>
</div>
"""

names = re.findall(r'<name size=\d>(.*?)</name>', html_content)
for name in names:
    print(name)
```

Ben Ziegler
Glenna Hill
Jason Woolfolk

Web crawlers / scrapers: Summary



Playtime!

- Do Part 3 - Web scraping.ipynb
- Do Part 4 - HTML and web scraping.ipynb (optional)

Web crawlers / scrapers: Summary

- Try not to (almost never) implement yourself
- A lot of frameworks are available, e.g., Scrapy (Python)
 - Tricks to imitate a real user not to be blocked
- A lot of HTML parsers are available (PyQuery, BeautifulSoup)
Easy parsing functionality (e.g., “*get all blocks where class=event*”)
- Read Terms & Conditions carefully, obey “robots.txt” rules

Web crawlers / scrapers: Summary

- Try not to (almost never) implement yourself
- A lot of frameworks are available, e.g., Scrapy (Python)
 - Tricks to imitate a real user not to be blocked
- A lot of HTML parsers are available (PyQuery, BeautifulSoup)
Easy parsing functionality (e.g., “*get all blocks where class=event*”)
- **Read Terms & Conditions carefully, obey “robots.txt” rules!**



APIs vs Web crawlers

	Pros	Cons
Third party APIs	<ul style="list-style-type: none">• Convenient, easy to use• Safe, won't be blocked• Fast	<ul style="list-style-type: none">• Need to manage API keys• Inflexible• Limit on access
Web crawlers / scrapers	<ul style="list-style-type: none">• Very flexible. Theoretically, you can collect anything you find.	<ul style="list-style-type: none">• Data structure can change• May be blocked• Terms & Conditions