



Lecture 7 Websites

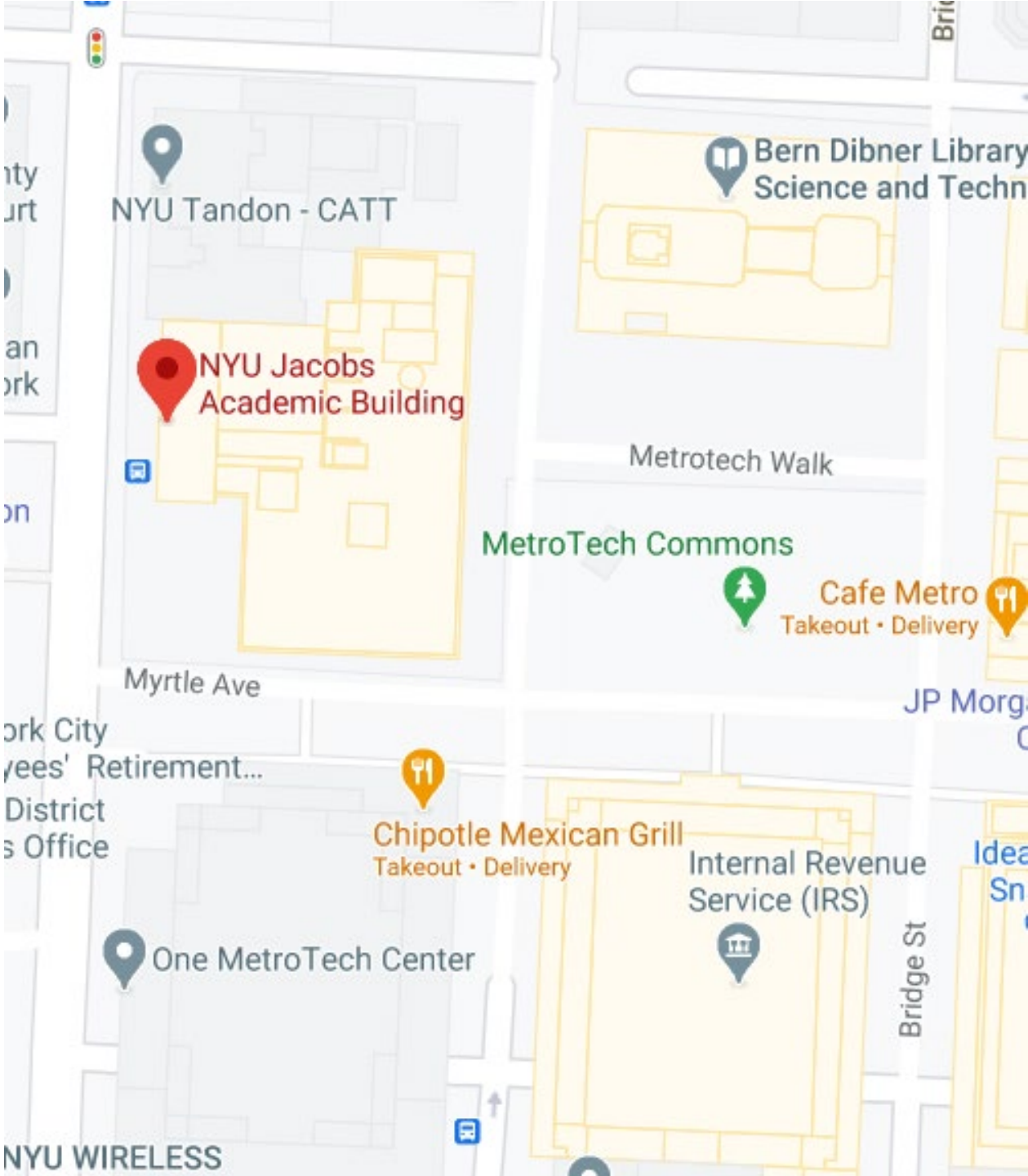
Programming for Business Analytics
MG-GY 8401





Agenda

- Connecting to Websites
- Web-Scraping
- Parsing Text as Data

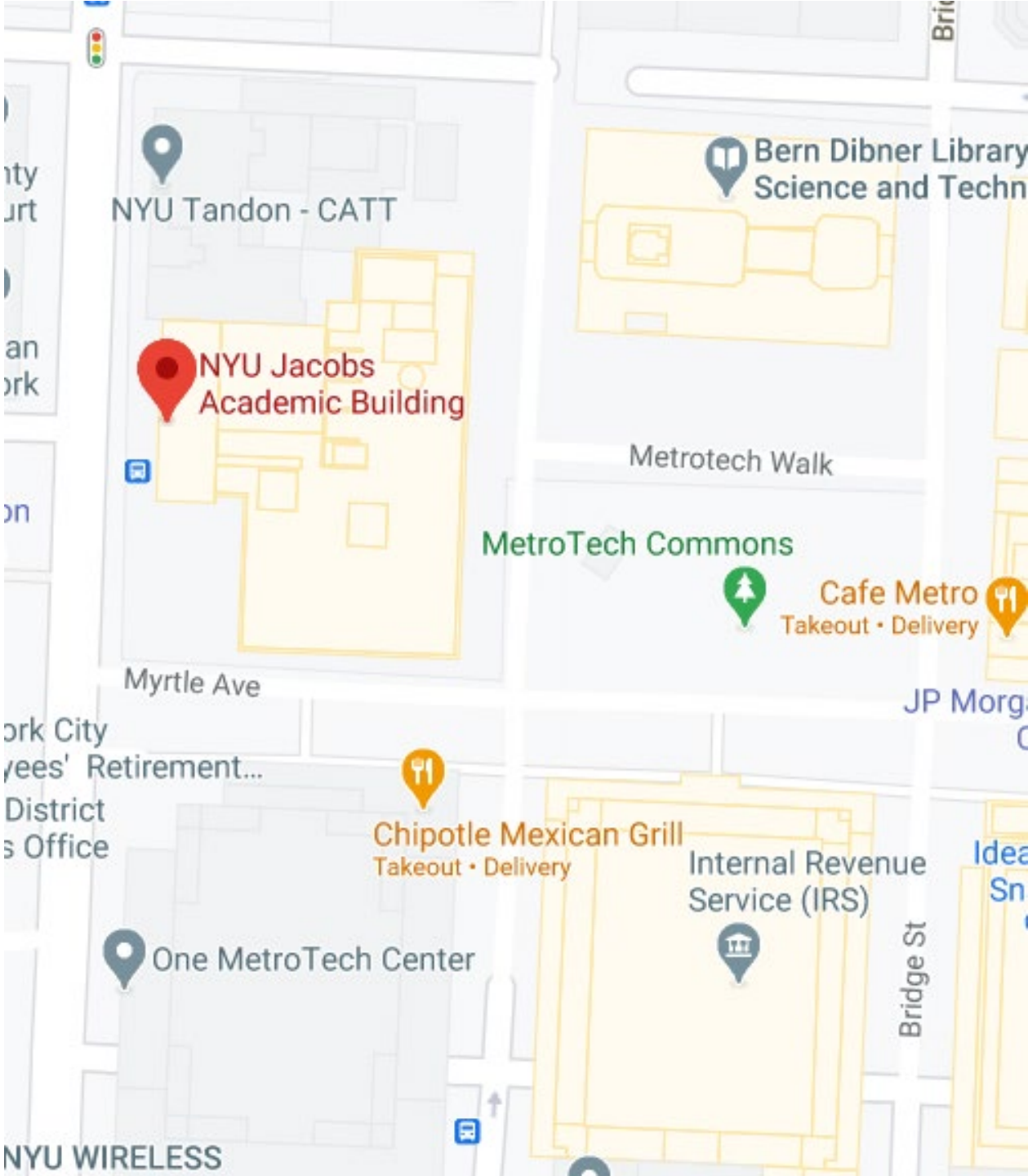


NYU

TANDON SCHOOL
OF ENGINEERING

Logistics

- Homework
 - Homework 6
 - Homework 5



NYU

TANDON SCHOOL
OF ENGINEERING

Logistics

- Homework
 - Homework 7
- Project



Websites

Facebook.com

+27.0%

Netflix.com

+16.0%

YouTube.com

+15.3%

170M

Feb. 29
First U.S. Covid-19 death

Average daily
traffic

26M

200M

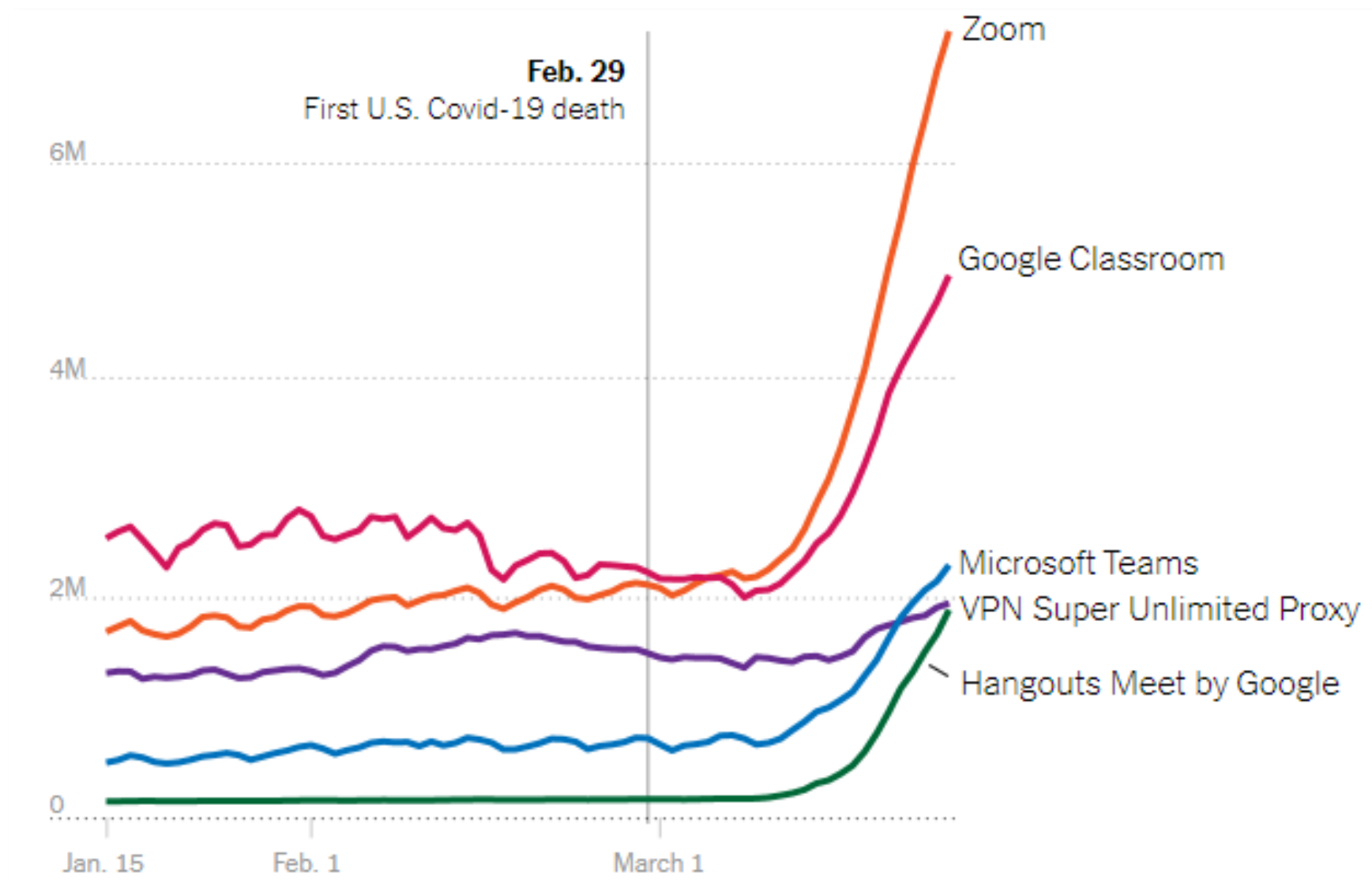
120M

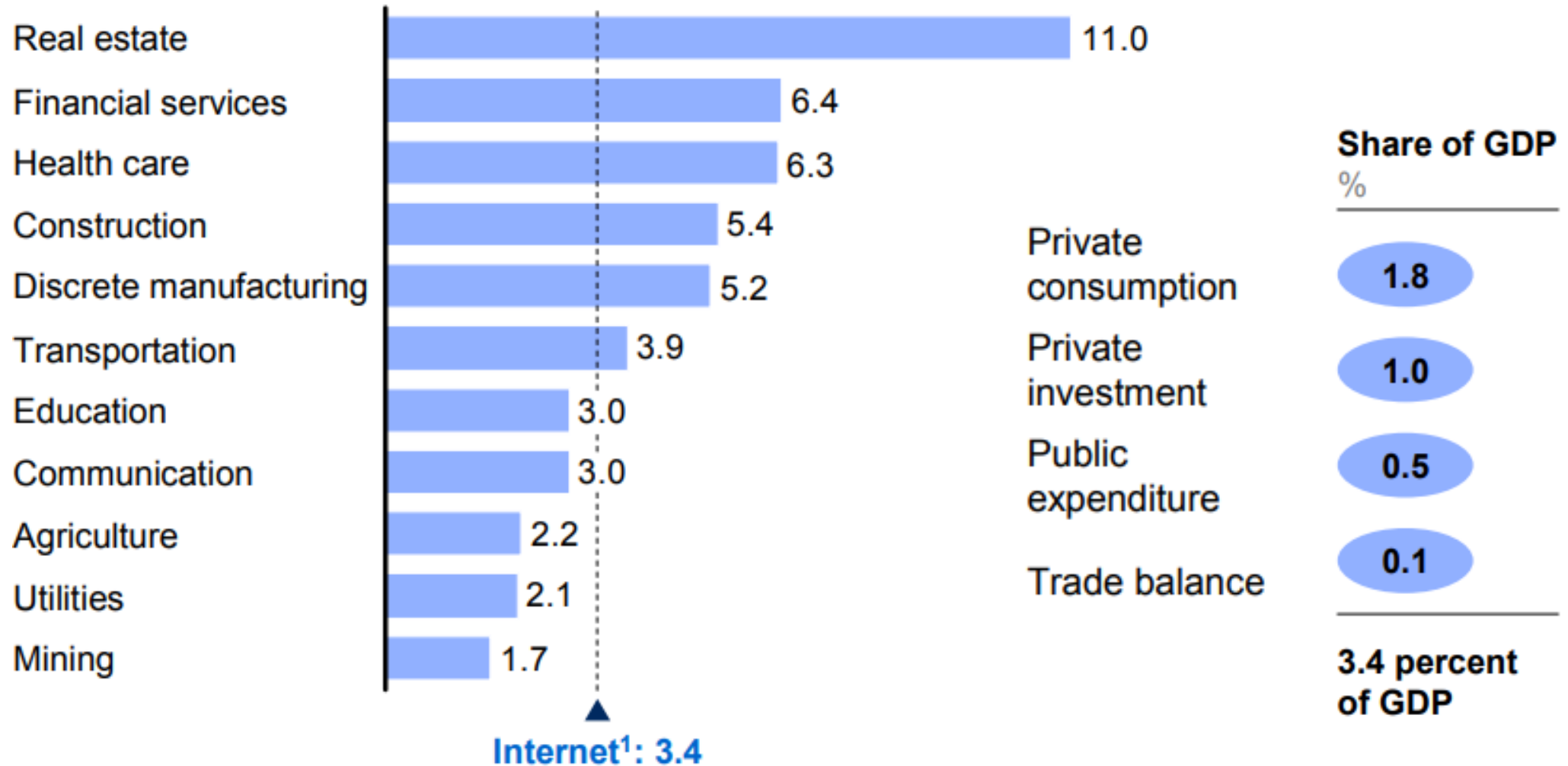
Jan. 15

March 24

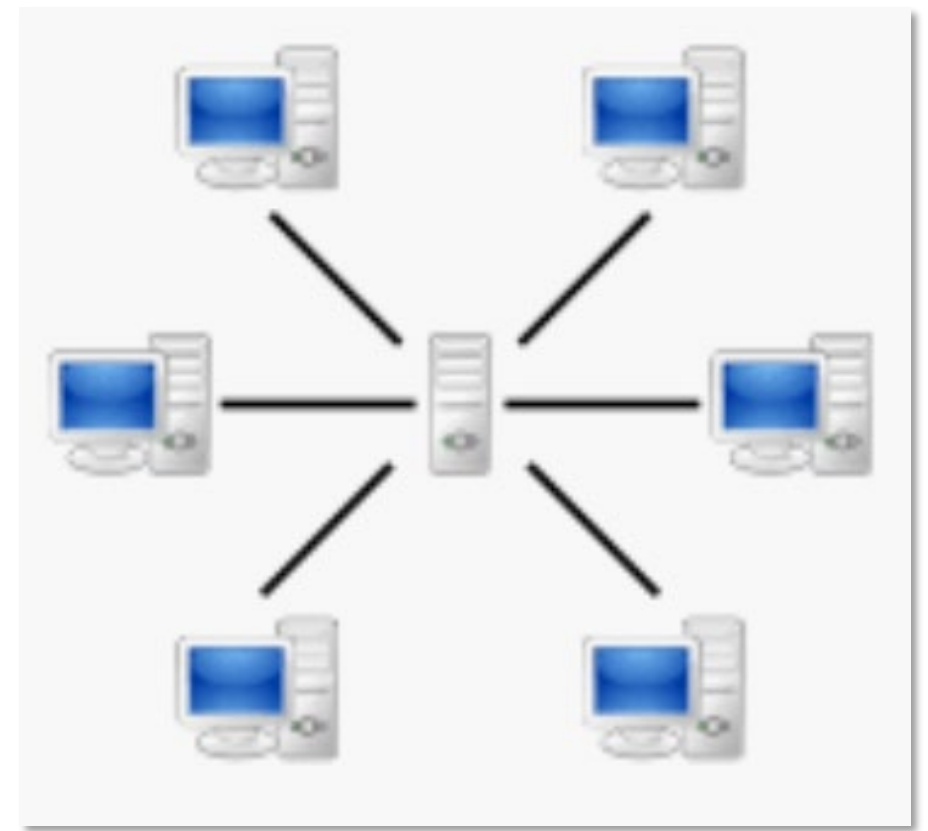
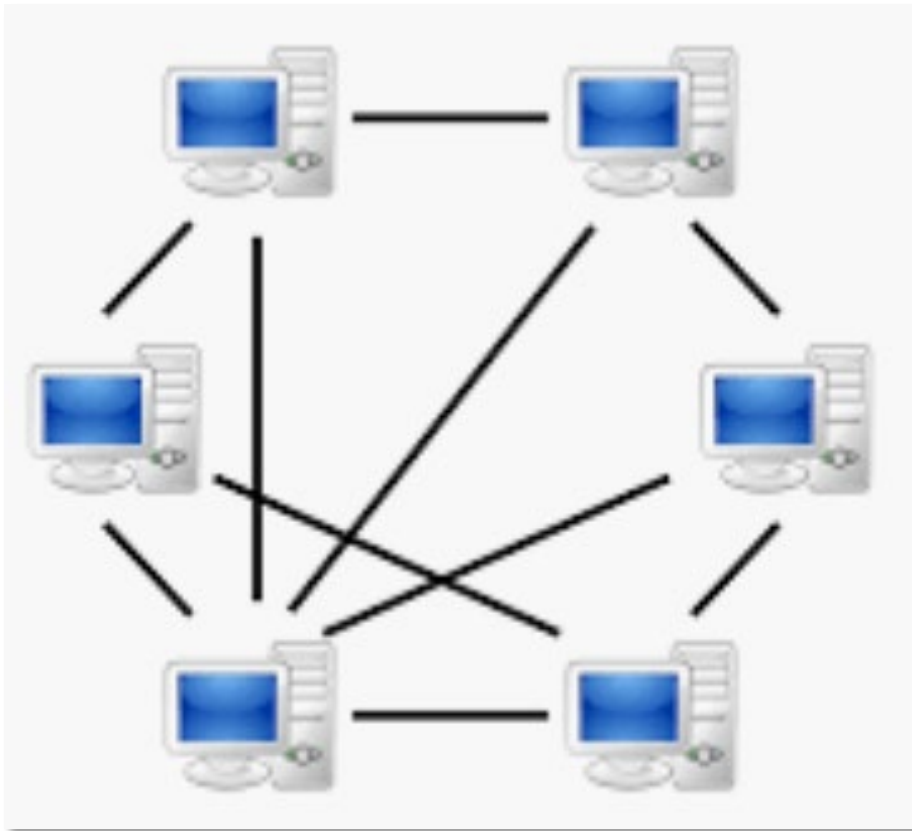
16M

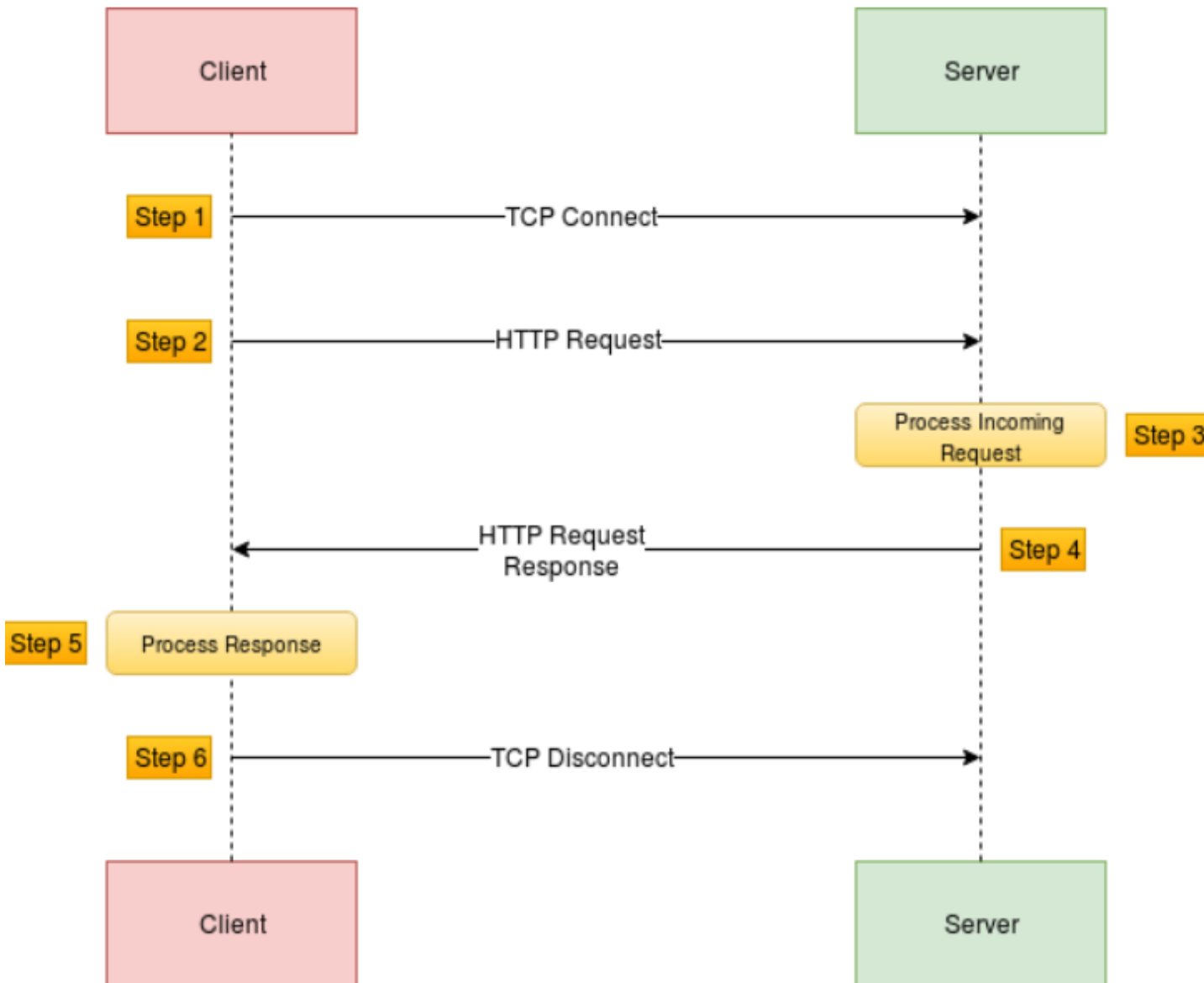
160M





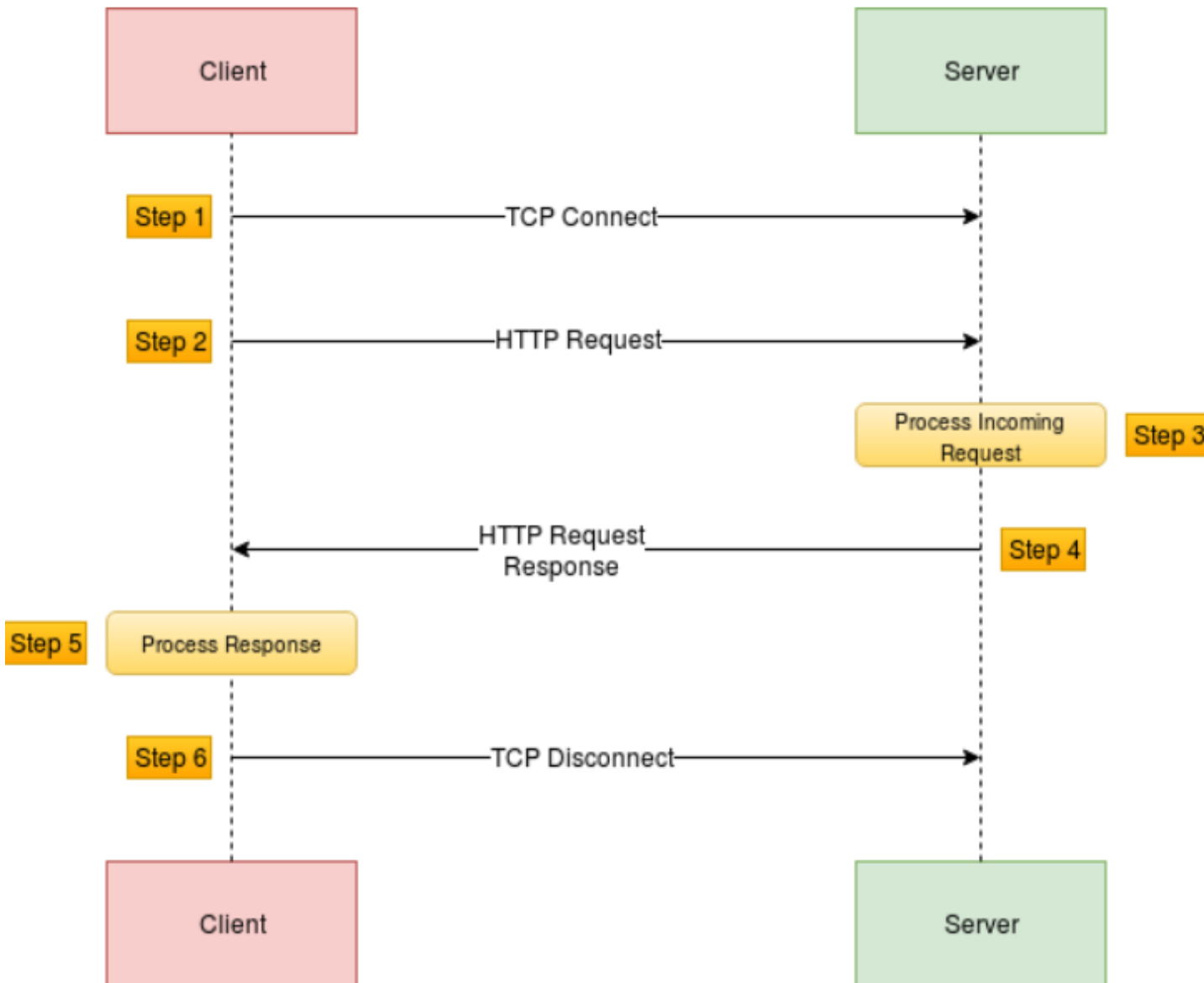
Internet Protocol





Transmission Control Protocol

- ordering of requests and responses
- splitting large packets into small packets
- removal of duplicate packets
- retransmission of lost packets



HyperText Transfer Protocol

- client sends request to server
- server runs an application to process the request
- server sends a response to client



userinfo host port

https://john.doe@www.example.com:123/forum/questions/?tag=networking&order=newest#top

scheme authority path query fragment

Uniform Resource Locator

- credentials for authentication
 - username + password
- location on the website
- query to filter content on website

Which of the following would match the regular expression `jo*hn`?

1. `jooohn`
2. `jon`
3. `jh`
4. `john`
5. `jooooooooohn`

Which of the following would match the regular expression `jo*hn`?

Try it on regex101.com where you can for experimenting with regular expressions

1. jooohn

2. jon

3. jhn

4. john

5. joooooooooohnn

**DELETE****/delete** "The request's DELETE parameters."**GET****/get** The request's query parameters.**PATCH****/patch** The request's PATCH parameters.**POST****/post** The request's POST parameters.**PUT****/put** The request's PUT parameters.



1XX
INFORMATIONAL

2XX
SUCCESS

3XX
REDIRECTION

4XX
CLIENT ERROR

5XX
SERVER ERROR

Which of the following would match **revolution**,
revolutionary, **revolutionaries**?

- A. `revolution[a-z]?`
- B. `revolution[a-z]*`
- C. `revolution[a-z]+`

Which of the following would match **revolution**,
revolutionary, **revolutionaries**?

A. `revolution[a-z]?`

B. `revolution[a-z]*`

C. `revolution[a-z]+`

```
{  
  "firstName": "John",  
  "lastName": "Smith",  
  "age": 27,  
  "address": {  
    "city": "New York",  
    "postalCode": "10021-3100"  
  },  
  "phoneNumbers": [  
    {  
      "type": "home",  
      "number": "212 555-1234"  
    },  
    {  
      "type": "office",  
      "number": "646 555-4567"  
    }  
  ],  
  "children": []  
}
```

- JSON format resembles dictionaries in Python with key-value pairs
- Keys are strings
- Values are
 - Numbers
 - Strings
 - Boolean
 - List
 - Another Dictionary

How could we match social security numbers in text?

Note that social security numbers have the form

###-##-####

1. $[0 \text{ to } 9]\{3\}-[0 \text{ to } 9]\{2\}-[0 \text{ to } 9]\{4\}$
2. $\{0-9\}[3]-\{0-9\}[2]-\{0-9\}[4]$
3. $[0-9]\{3\}-[0-9]\{2\}-[0-9]\{4\}$
4. $[0-9]\{3\}-[0-9]\{3\}-[0-9]\{3\}$

How could we match social security numbers in text?

Note that social security numbers have the form

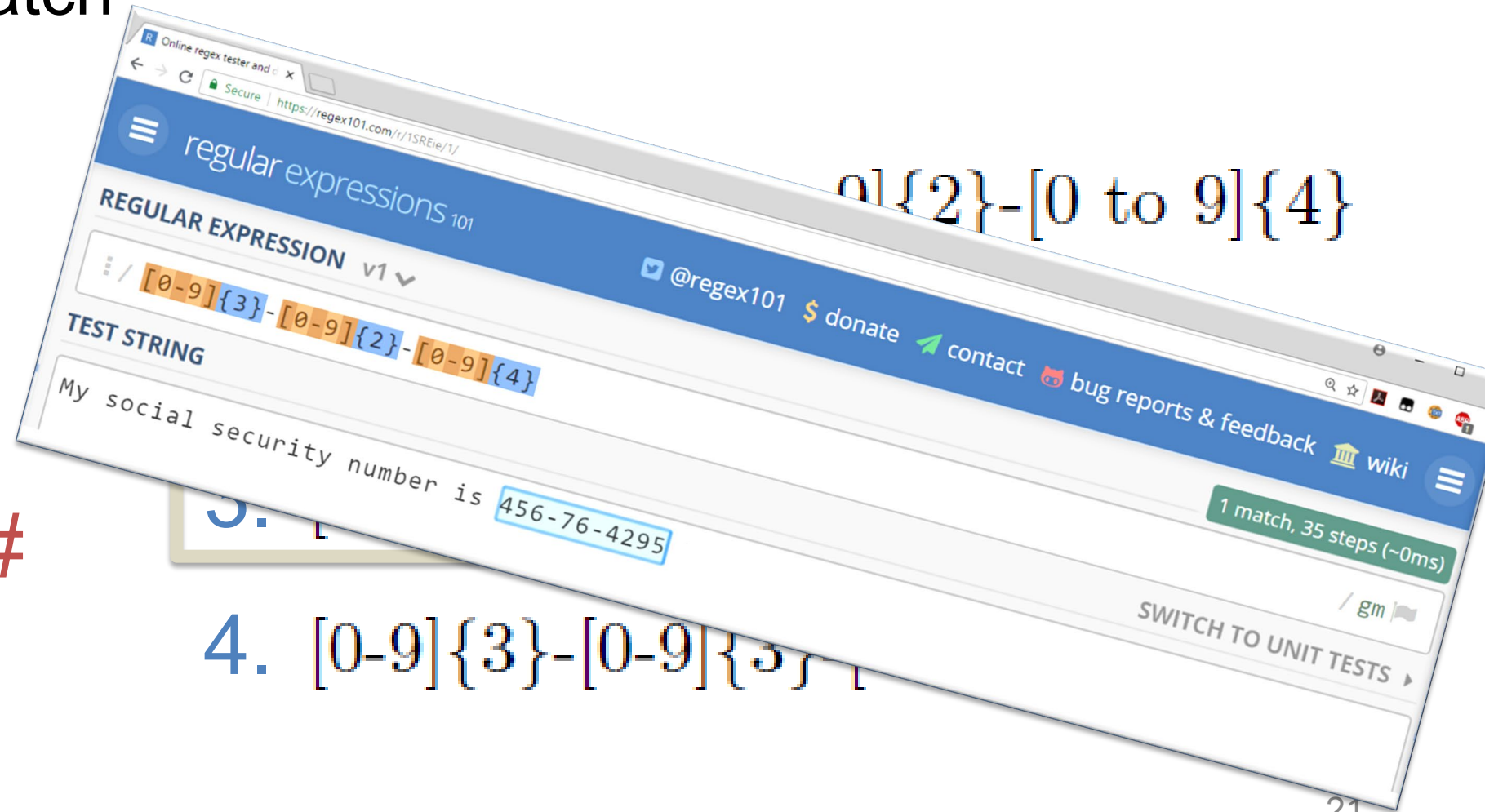
###-##-####

1. $[0 \text{ to } 9]\{3\}-[0 \text{ to } 9]\{2\}-[0 \text{ to } 9]\{4\}$
2. $\{0-9\}[3]-\{0-9\}[2]-\{0-9\}[4]$
3. $[0-9]\{3\}-[0-9]\{2\}-[0-9]\{4\}$
4. $[0-9]\{3\}-[0-9]\{3\}-[0-9]\{3\}$

How could we match social security numbers in text?

Note that social security numbers have the form

###-##-####



4. `[0-9]{3}-[0-9]{2}-[0-9]{4}`



```
<catalog>
  <plant>
    <common>Bloodroot</common>
    <botanical>Sanguinaria canadensis</botanical>
    <zone>4</zone>
    <light>Mostly Shady</light>
    <price currency="USD">$2.44</price>
    <availability>031599</availability>
  </plant>
  <plant type='a'>
    <common>Columbine</common>
    <botanical>Aquilegia canadensis</botanical>
    <zone>3</zone>
    <light>Mostly Shady</light>
    <price currency="USD">$9.37</price>
    <availability>030699</availability>
  </plant>
</catalog>
```



```
<catalog>
  <plant>
    <common>Bloodroot</common>
    <botanical>Sanguinaria canadensis</botanical>
    <zone>4</zone>
    <light>Mostly Shady</light>
    <price currency="USD">$2.44</price>
    <availability>031999</availability>
  </plant>
  <plant type="a">
    <common>Columbine</common>
    <botanical>Aquilegia canadensis</botanical>
    <zone>3</zone>
    <light>Mostly Shady</light>
    <price currency="USD">$9.37</price>
    <availability>030699</availability>
  </plant>
</catalog>
```

Diagram illustrating XML structure and components:

- Start tag:** Points to the opening tag `<plant>`.
- End tag:** Points to the closing tag `</plant>`.
- Content:** Points to the text `Mostly Shady` between the `<light>` and `</light>` tags.



```
<catalog>
  <plant>
    <common>Bloodroot</common>
    <botanical>Sanguinaria canadensis</botanical>
    <zone>4</zone>
    <light>Mostly Shady</light>
    <price currency="USD">$2.44</price>
    <availability>031599</availability>
  </plant>
  <plant type='a'>
    <common>Columbine</common>
    <botanical>Aquilegia canadensis</botanical>
    <zone>3</zone>
    <light>Mostly Shady</light>
    <price currency="USD">$9.37</price>
    <availability>030699</availability>
  </plant>
</catalog>
```

Start tag

Content consists of multiple elements

End tag



```
<catalog>
  <plant>
    <common>Bloodroot</common>
    <botanical>Sanguinaria canadensis</botanical>
    <zone>4</zone>
    <light>Mostly Shady</light>
    <price currency="USD">$2.44</price>
    <availability>031599</availability>
  </plant>
  <plant type='a'>
    <common>Columbine</common>
    <botanical>Aquilegia canadensis</botanical>
    <zone>3</zone>
    <light>Mostly Shady</light>
    <price currency="USD">$9.37</price>
    <availability>030699</availability>
  </plant>
</catalog>
```

Attribute



Review

- HyperText Transfer Protocol
- Javascript Object Notation
- HyperText Markup Language

