

Introduction to Programming for Scientists

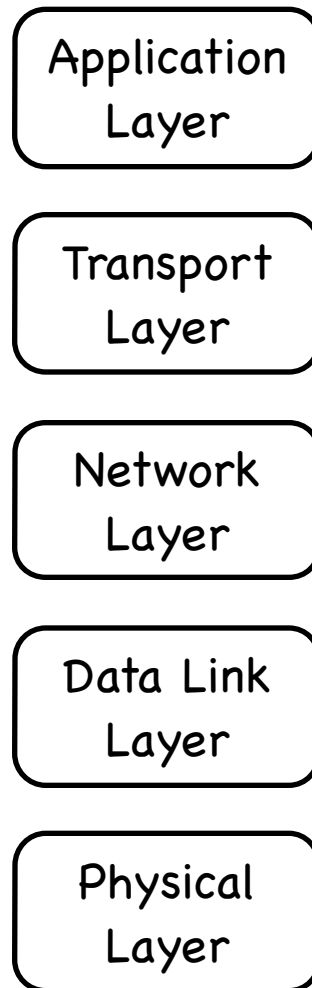
Lecture 6

HTML, XML and the Web
Web Scraping and How to Avoid It

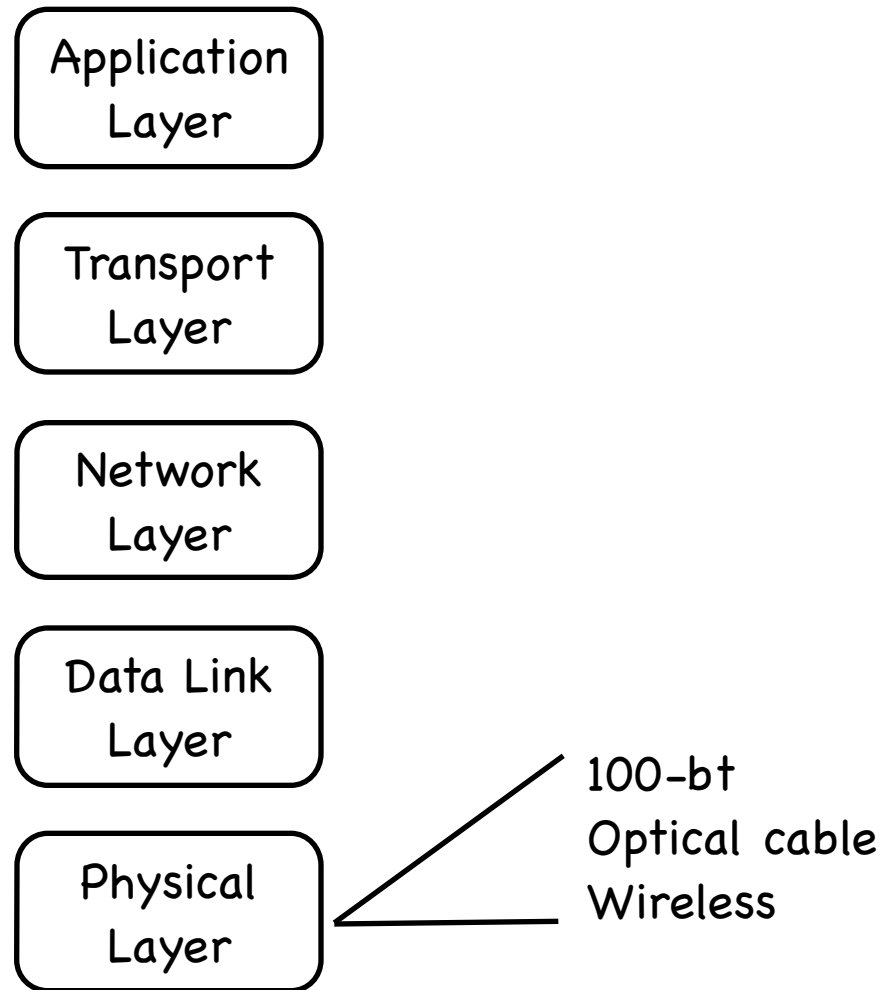
Prof. Steven Ludtke
N410.07, sludtke@bcm.edu

Quick Homework Review

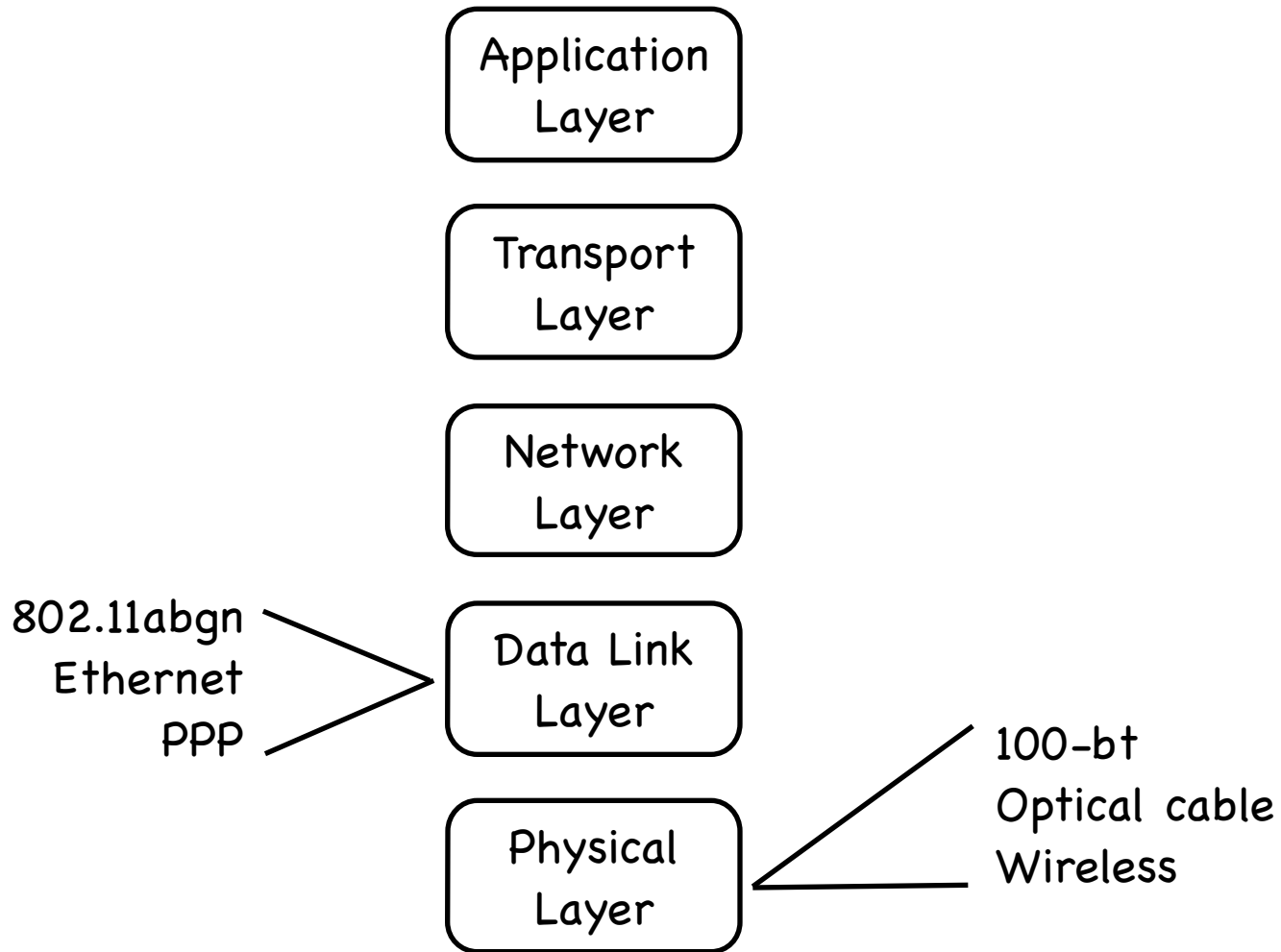
Networking



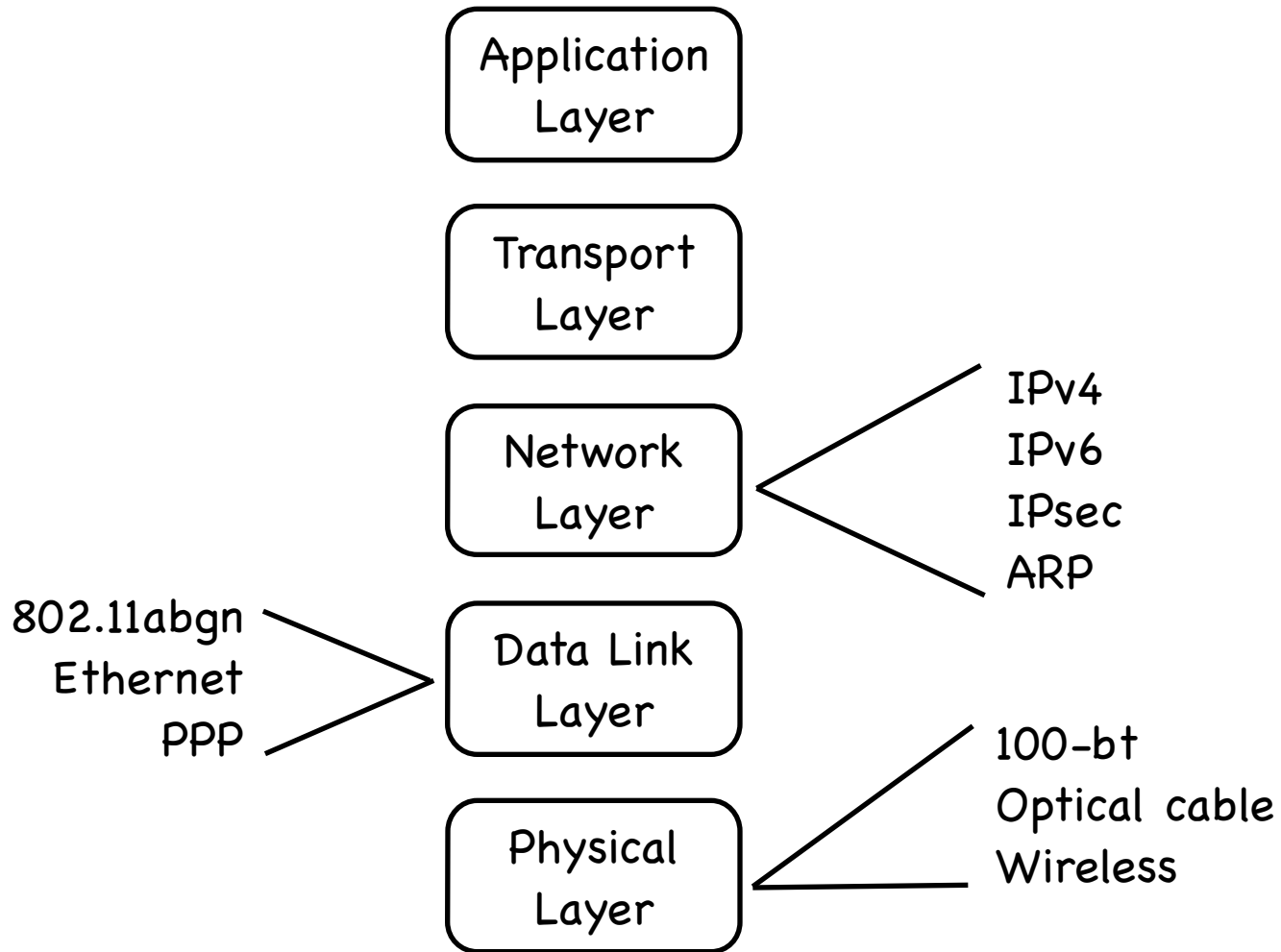
Networking



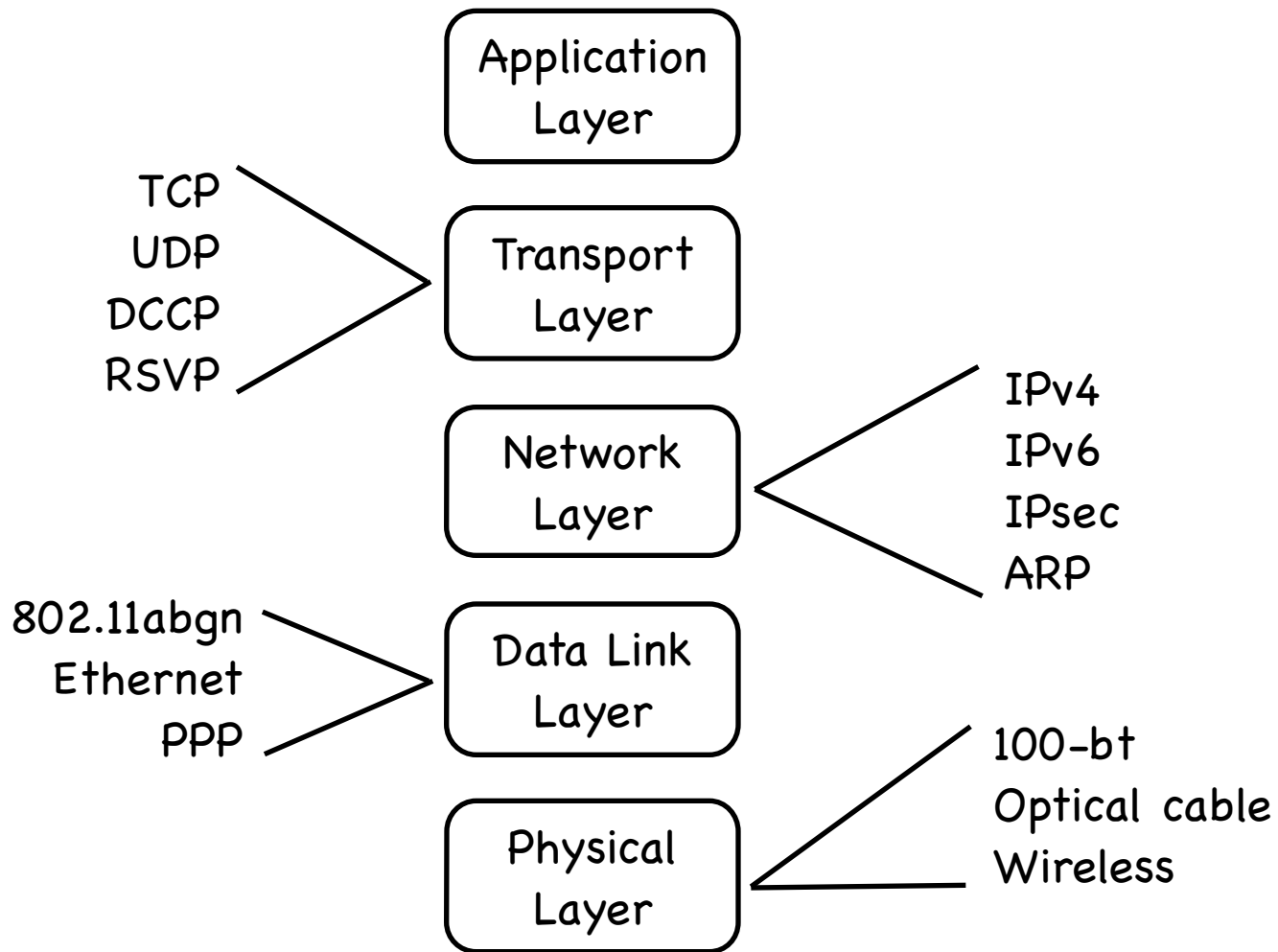
Networking



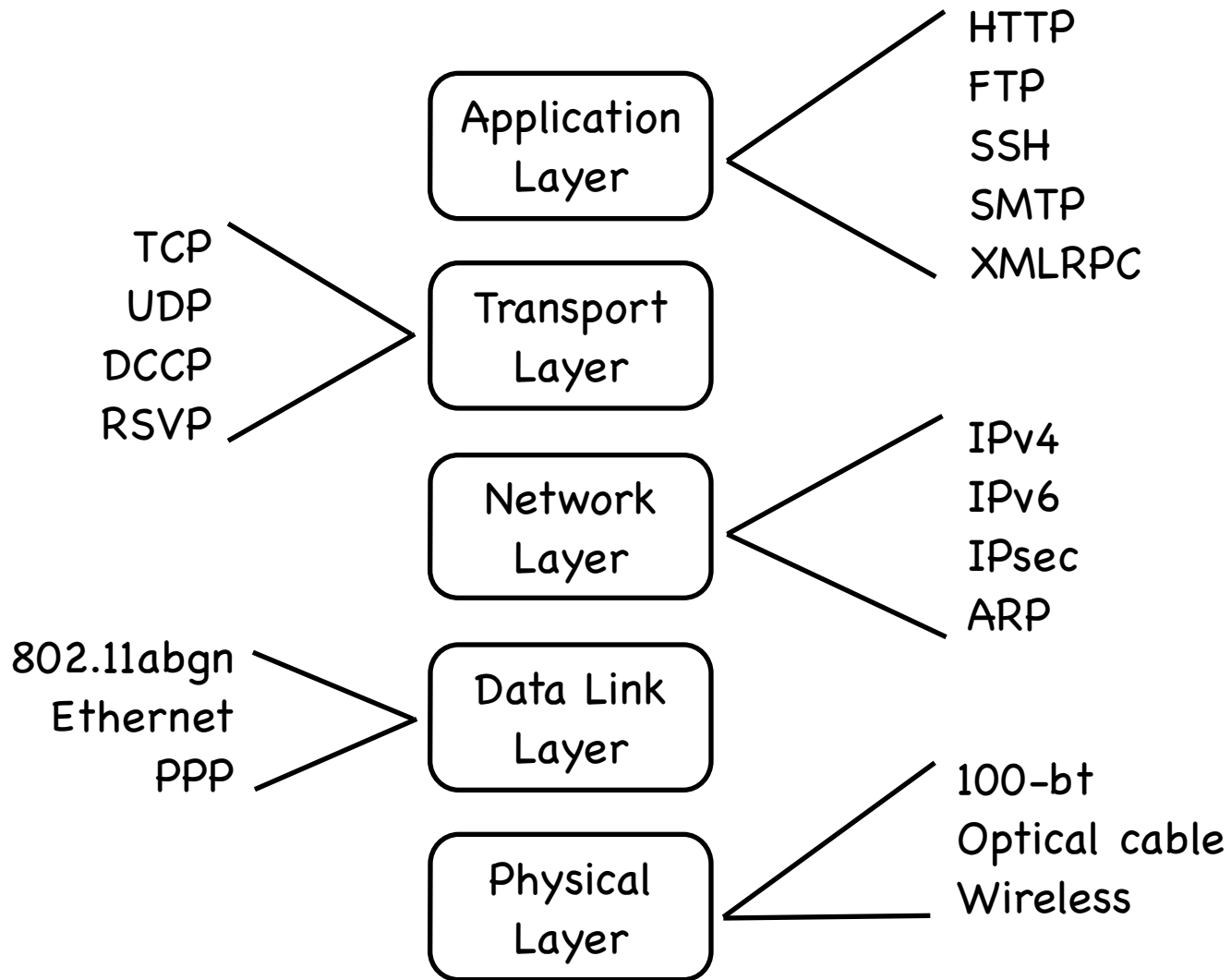
Networking



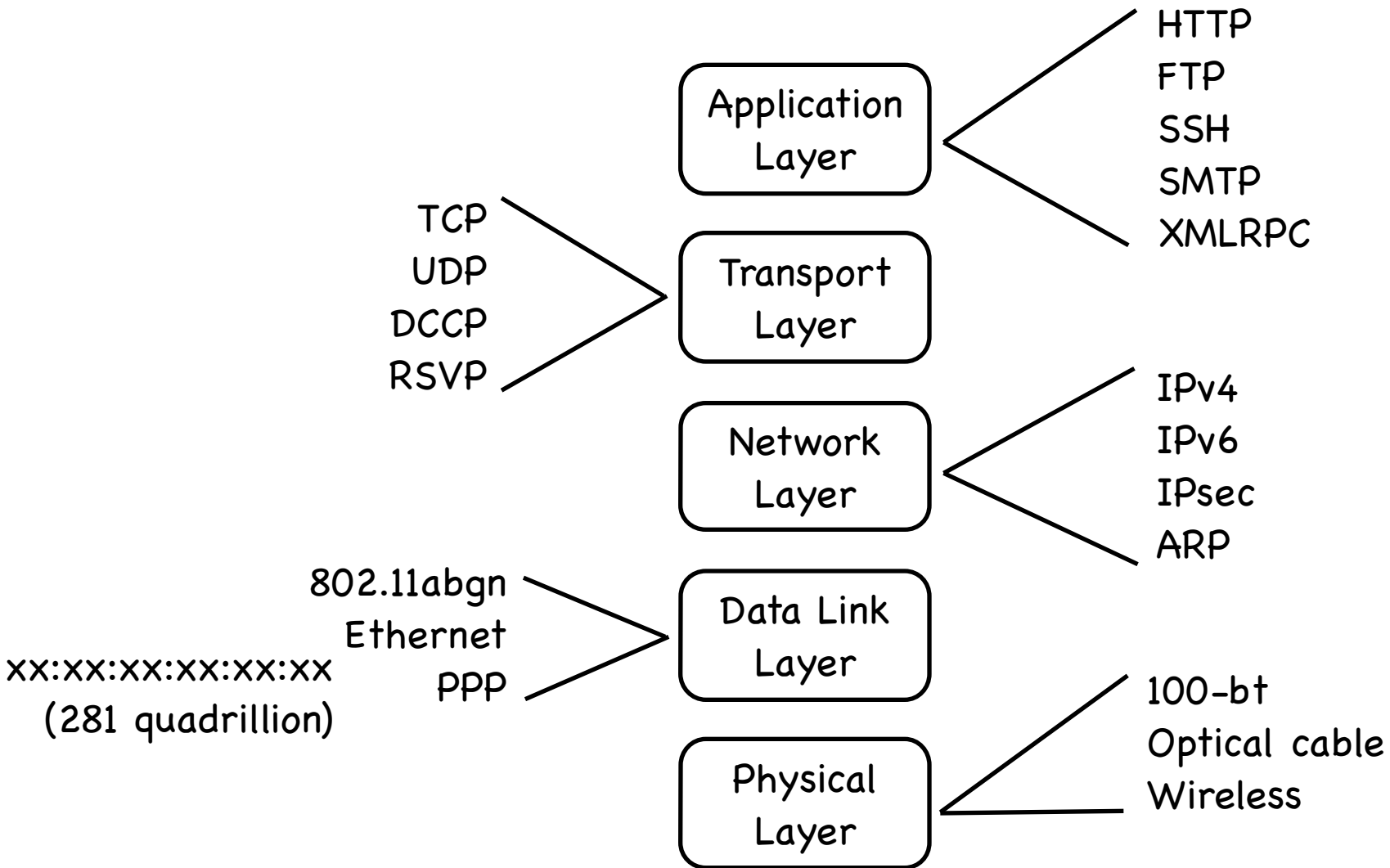
Networking



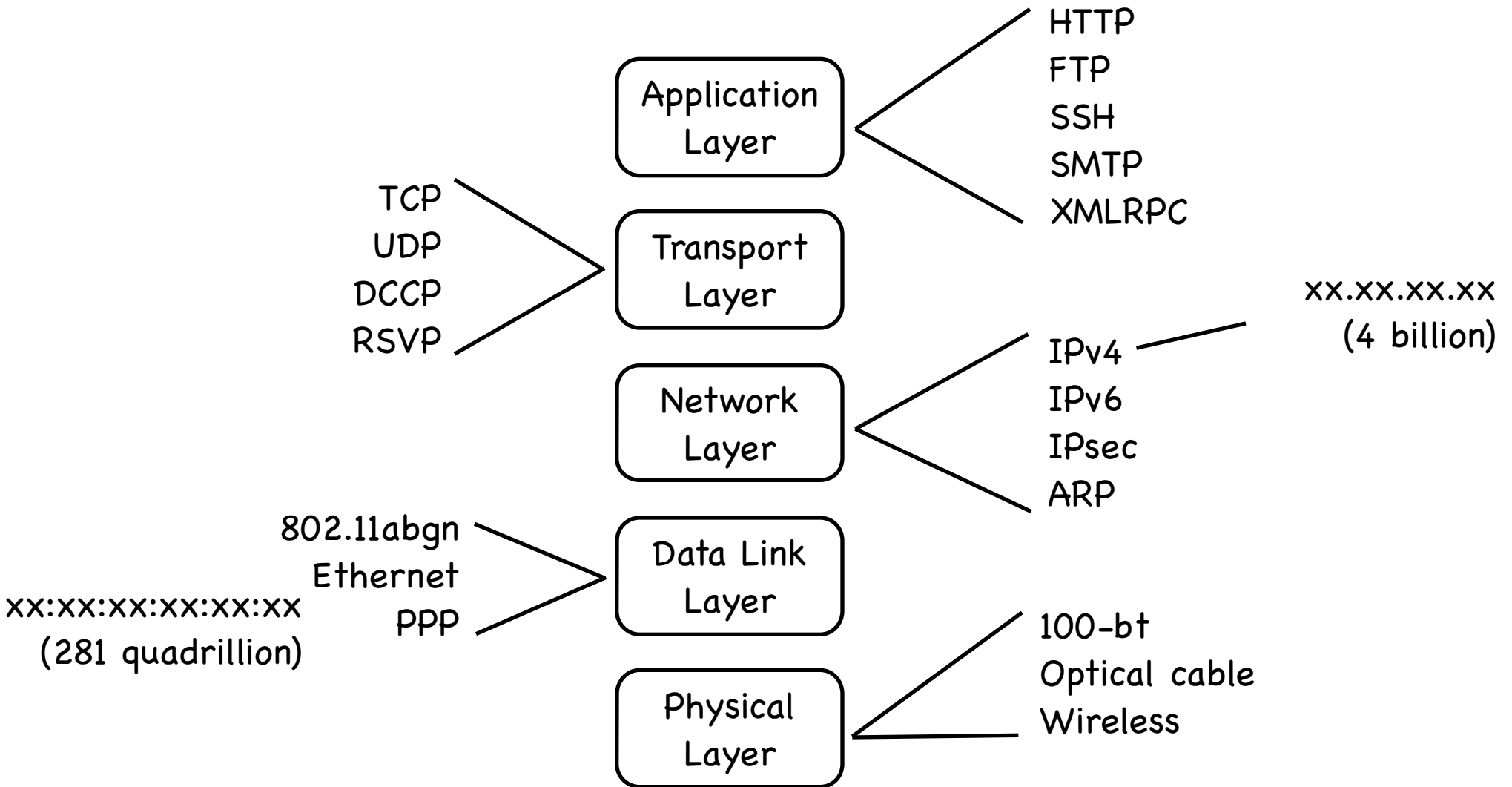
Networking



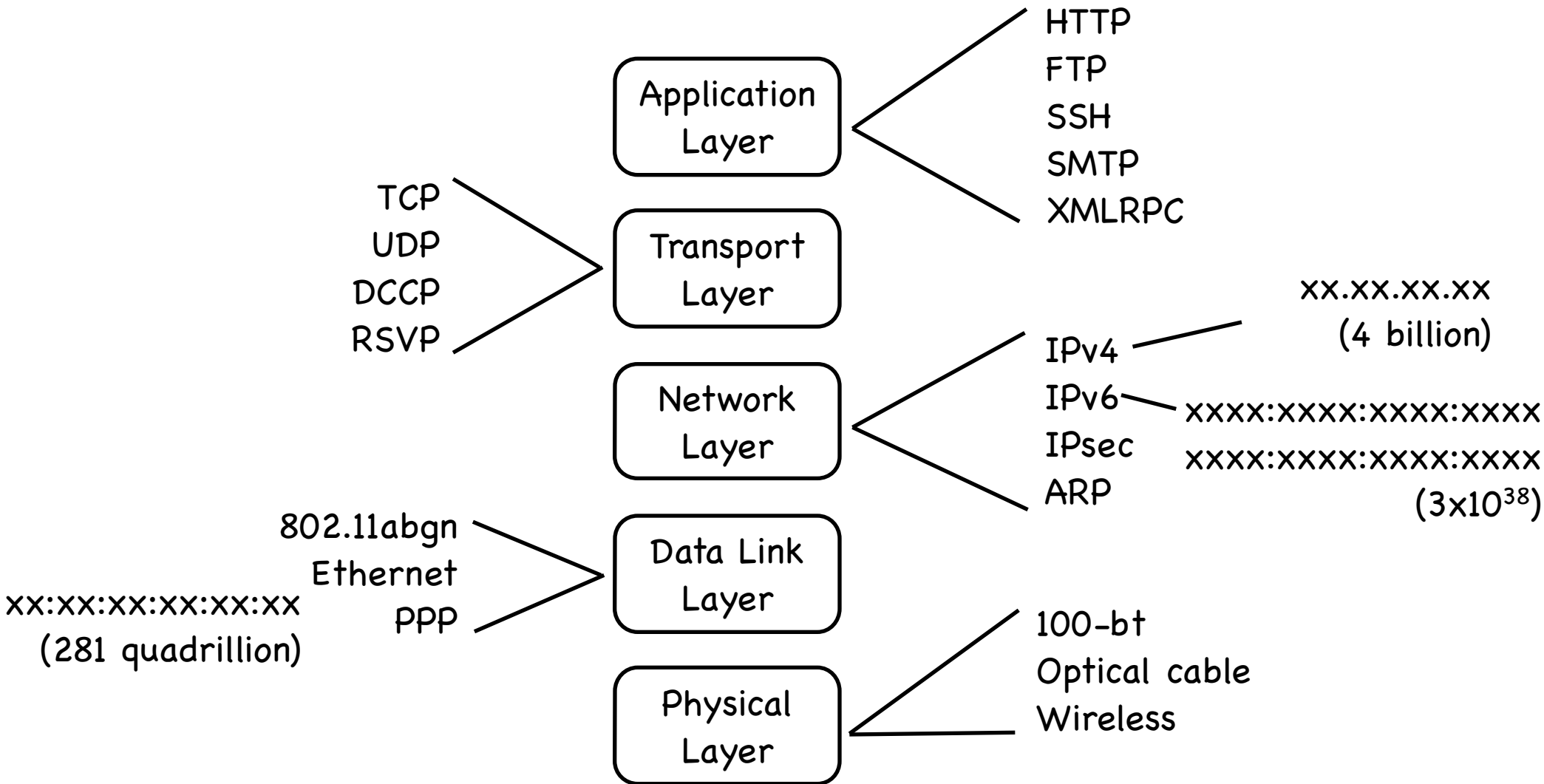
Networking



Networking

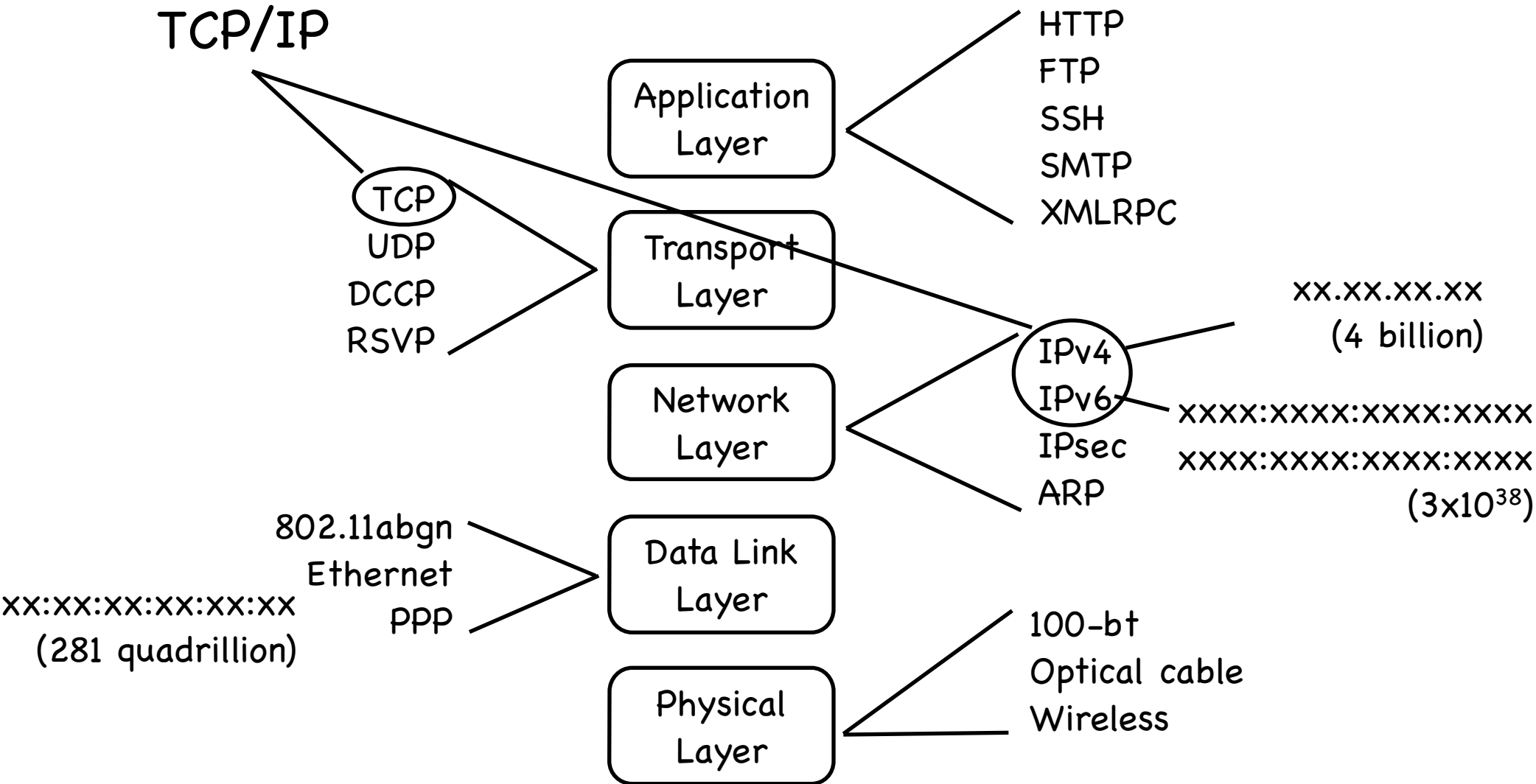


Networking



Networking

TCP/IP



IPv4 Network

- ① IP Address - Computer's unique* address (x.x.x.x)
- ① Netmask - defines local 'subnet', machines the computer can speak to 'directly'
- ① Router - Address used to contact machines outside subnet
- ① DNS Server - Address where names can be mapped to addresses
- ① Port - For a specific connection 0-65535, 0-1023 reserved for system services

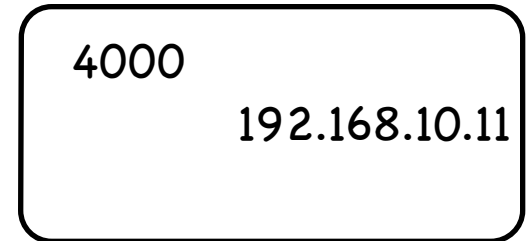
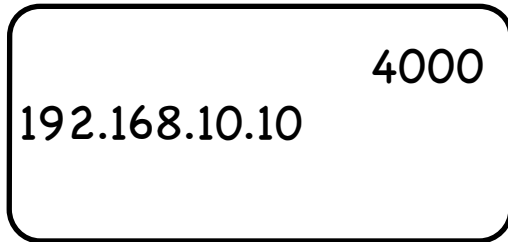
* - Some addresses are for 'private networks'. These are 10.*.*, 172.16.*.*-172.31.*.*, and 192.168.*.*

Common Services

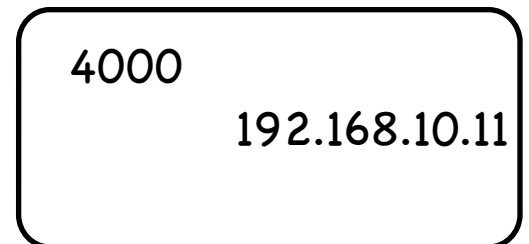
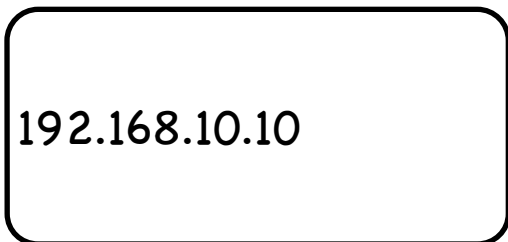
port	service
21	ftp
22	ssh
23	telnet
25	smtp (mail)
79	finger
80	http (web)
110	pop3 (email retrieval)
123	ntp (time)
137-139	Windows file sharing
143	imap (email retrieval)
443	https (secure http)

Sockets (TCP/UDP)

UDP

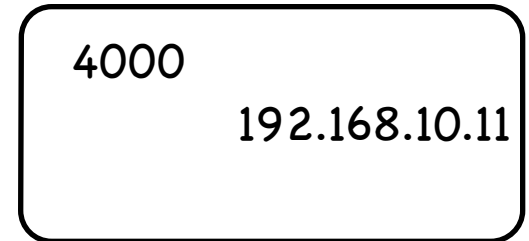
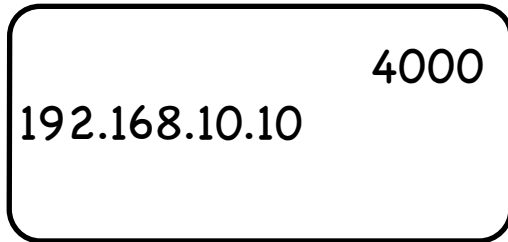


TCP

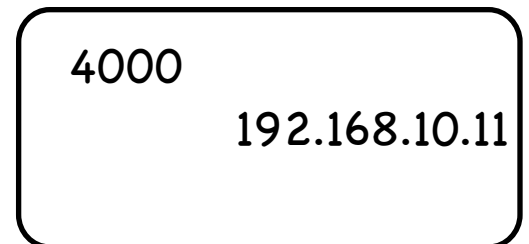
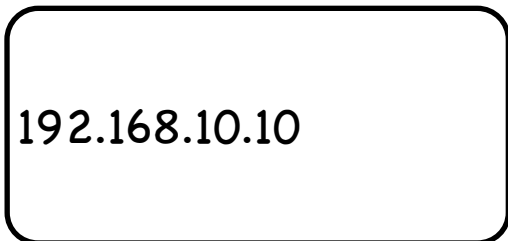


Sockets (TCP/UDP)

UDP

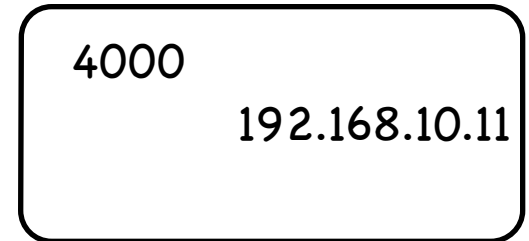
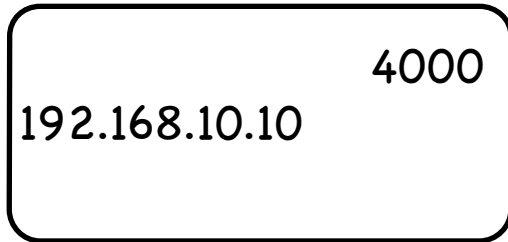


TCP

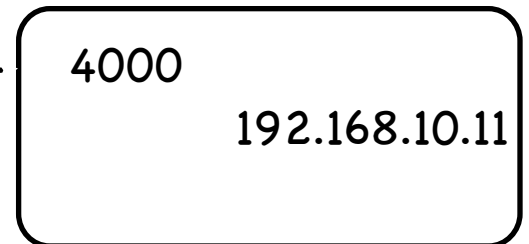
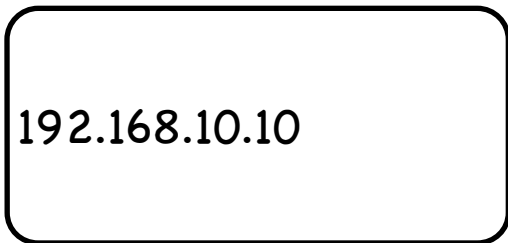


Sockets (TCP/UDP)

UDP

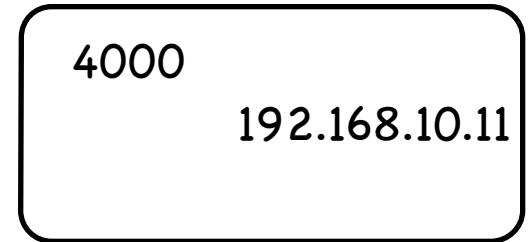
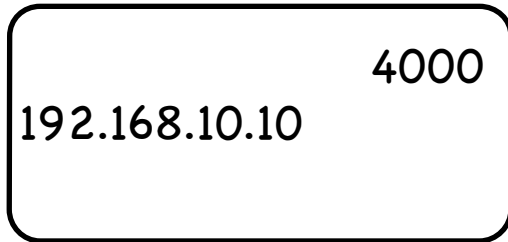


TCP

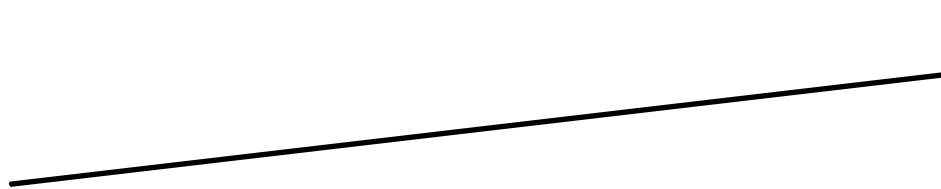
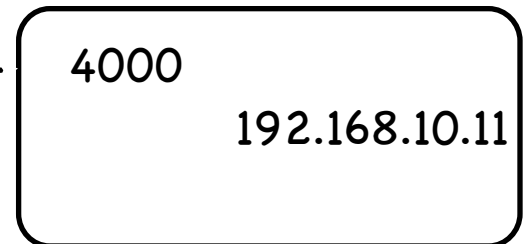
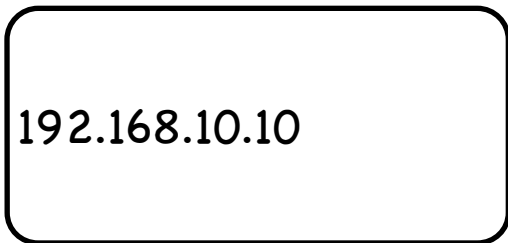


Sockets (TCP/UDP)

UDP



TCP



IPv4 Network Parameters

- ④ IP Address - Computer's unique* address (x.x.x.x)
- ④ Netmask - defines local 'subnet', machines the computer can speak to 'directly'
- ④ Router - Address used to contact machines outside subnet
- ④ DNS Server - Address where names can be mapped to addresses
- ④ Port - For a specific connection 0-65535, 0-1023 reserved for system services

* - Some addresses are for 'private networks'. These are 10.*.**, 172.16.*.*-172.31.*.*, and 192.168.*.*

NAT

TCP

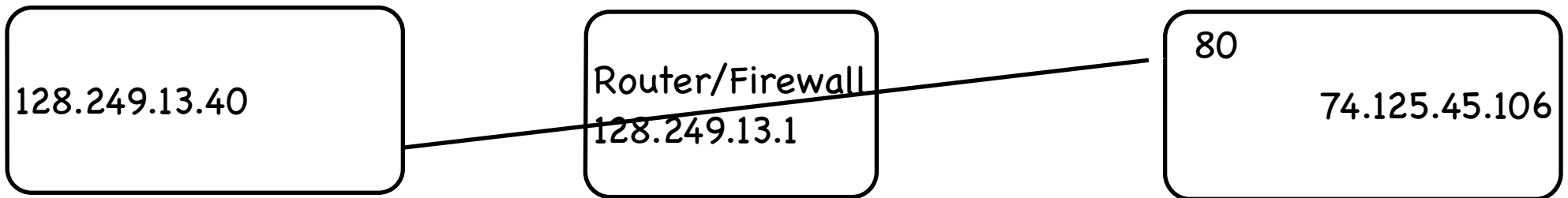
128.249.13.40

Router/Firewall
128.249.13.1

74.125.45.106

NAT

TCP



NAT

TCP

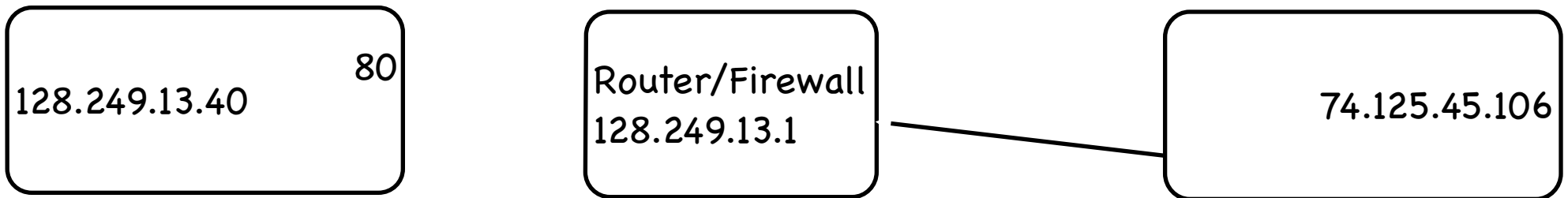
128.249.13.40

Router/Firewall
128.249.13.1

74.125.45.106

NAT

TCP



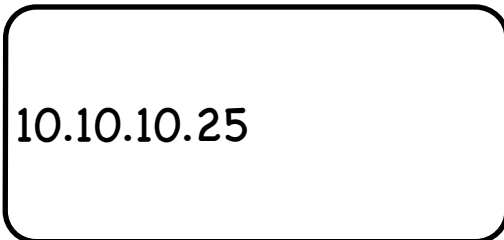
NAT

TCP

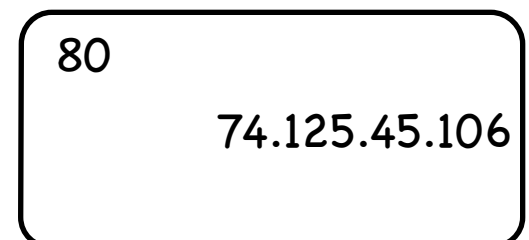
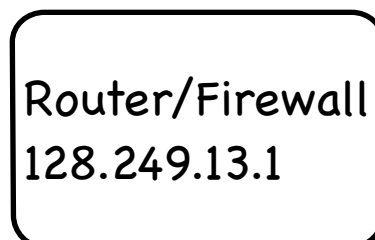
NAT

TCP

Intranet

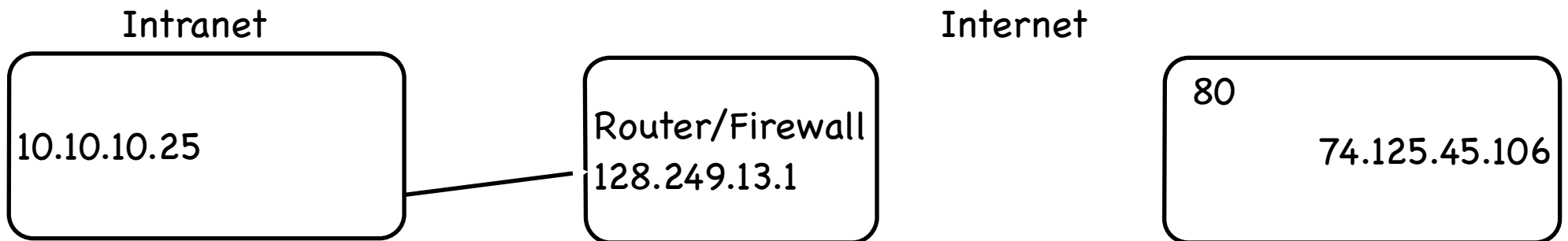


Internet



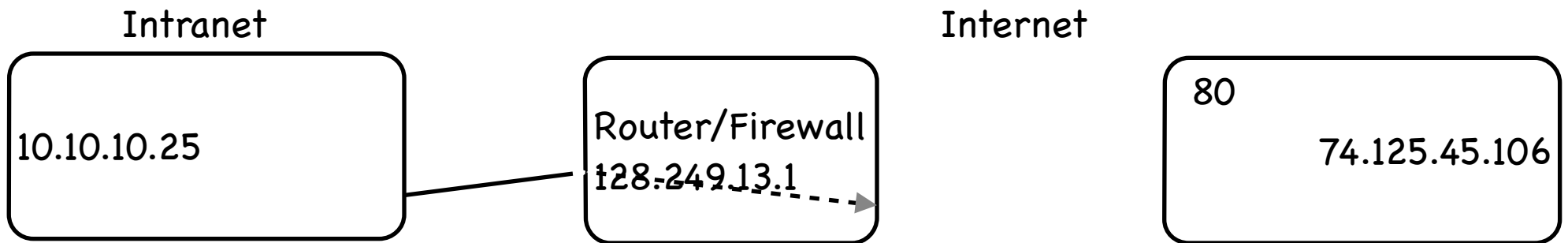
NAT

TCP



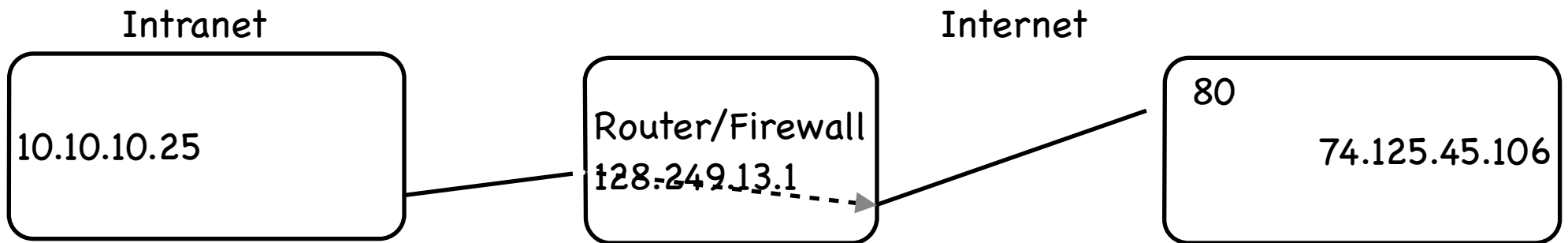
NAT

TCP



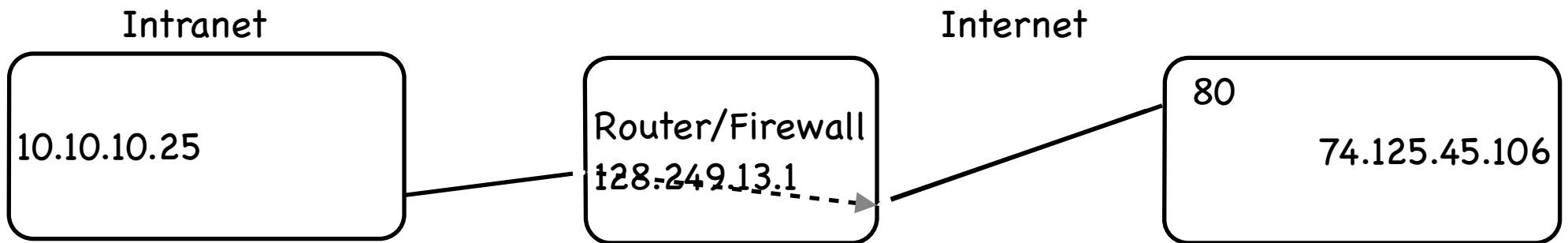
NAT

TCP



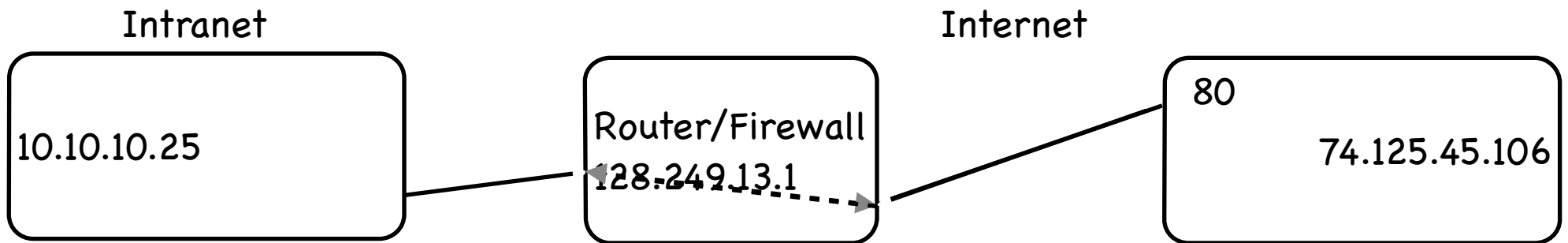
NAT

TCP



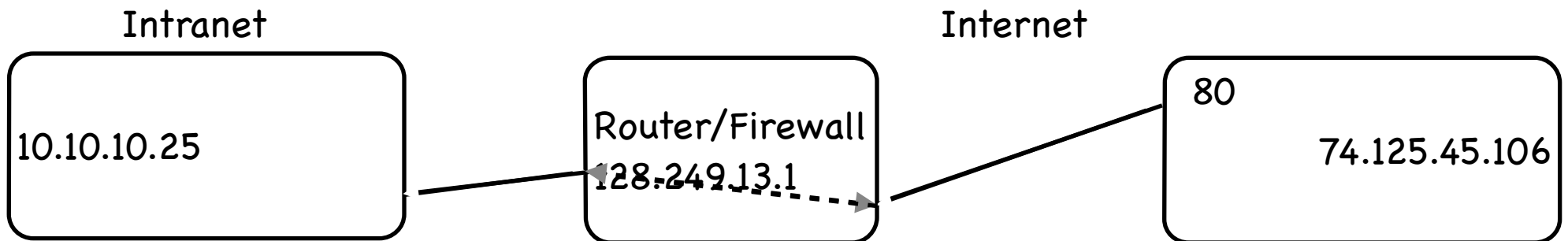
NAT

TCP



NAT

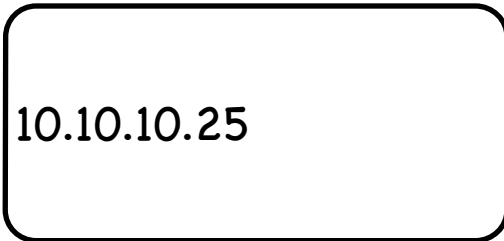
TCP



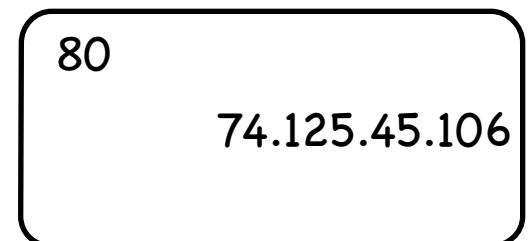
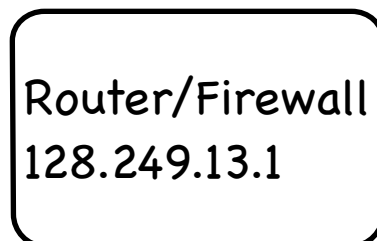
NAT

TCP

Intranet

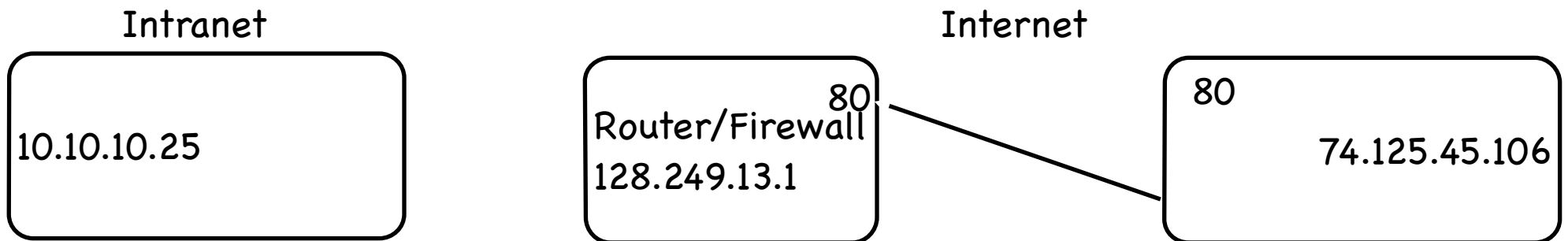


Internet



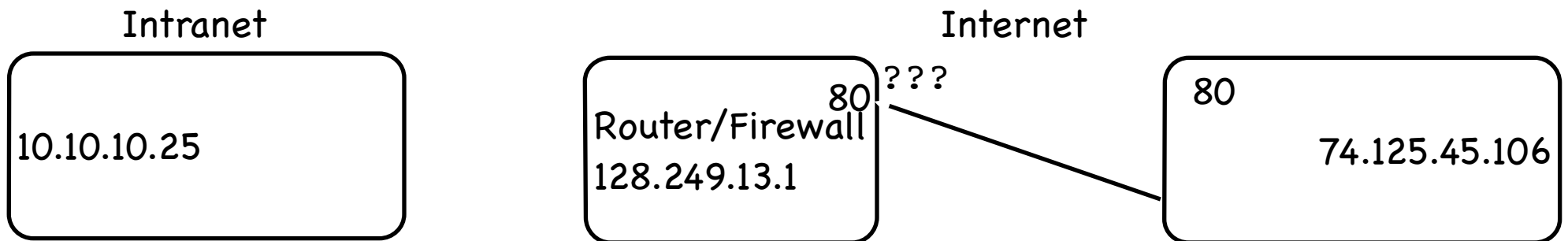
NAT

TCP



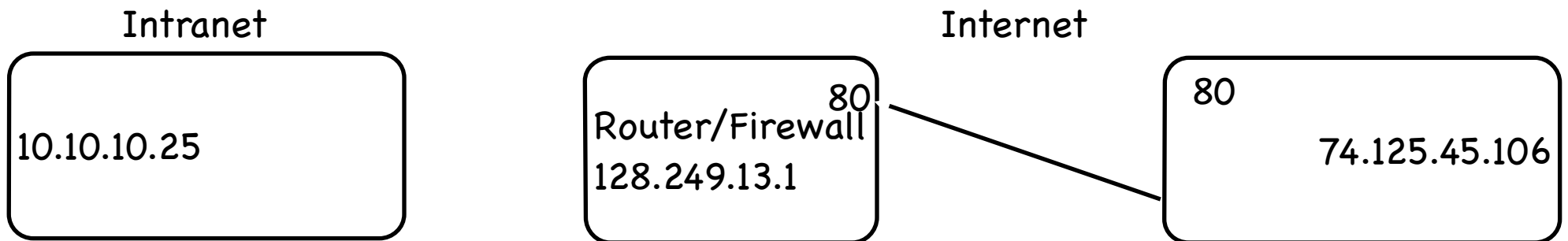
NAT

TCP



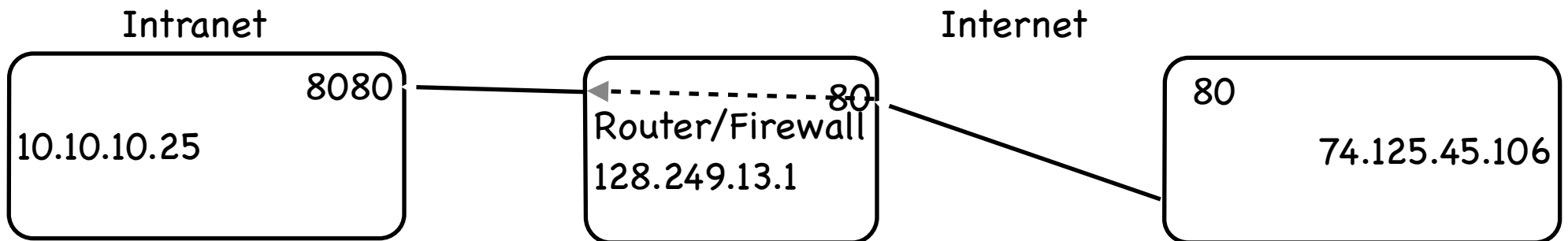
NAT

TCP



NAT

TCP



HISTORY

port	service
ARPANET (1969)	
23	telnet (1969)
21	ftp (1971)
79	finger (1977)
25	smtp (email, 1982)
NSFNET (1986)	
70	gopher (1991)
80	http & HTML (web, 1991)
NCSA Mosaic (1993)	
443	https (secure http, 1994)
NSFNET decommissioned, Internet Commercialized (1995)	
HTML 2.0 (1995)	
22	ssh (1995)
HTML 4 - commercial infighting (1999)	
HTML5/XHTML5 (2014-2016?)	

Glossary

- ④ HTTP - Hypertext Transport Protocol
- ④ HTML - Hypertext Markup Language
- ④ XML - Extensible Markup Language
 - ④ SGML - Parent of XML & HTML
- ④ JavaScript - Python-like language on webpages (NOT Java)
- ④ JSON - JavaScript Object Notation
- ④ AJAX - Asynchronous JavaScript and XML (2005 Google Maps)
- ④ AJAJ - Asynchronous JavaScript and JSON
- ④ CSS - Cascading Style Sheet
- ④ RSS - Really Simple Syndication

urllib2

```
import urllib2
```

```
f=urllib2.urlopen("http://blake.bcm.edu/dl/test.html")
```

```
for i in f: print i
```

HTML

- ④ <http://www.w3.org/History/19921103-hypertext/hypertext/WWW/TheProject.html>

- ④ Declarative language

- ④ HTML is a type of XML, XHTML obeys XML rules more completely

- ④ Python HTMLParser module

- ④ 'commands' in HTML are denoted by
<command option=value option=value>text</command>

- ④ For example:

```
<HTML>
```

```
<HEAD><TITLE>My Page</TITLE></HEAD>
```

```
<BODY>
```

```
<H3>Hi Everyone</H3>
```

```
<P>This is really just some test text to demonstrate how HTML works.
```

```
I can do interesting things like <i>italicize</i> or make text <b>bold</b>,
or even <b><i>both together</i></b>. ta da
```

```
</BODY>
```

CSS

- ④ Cascading Style Sheet

- ④ Used to present websites in a uniform way

- ④ Define how a page looks, then use the definition on many pages

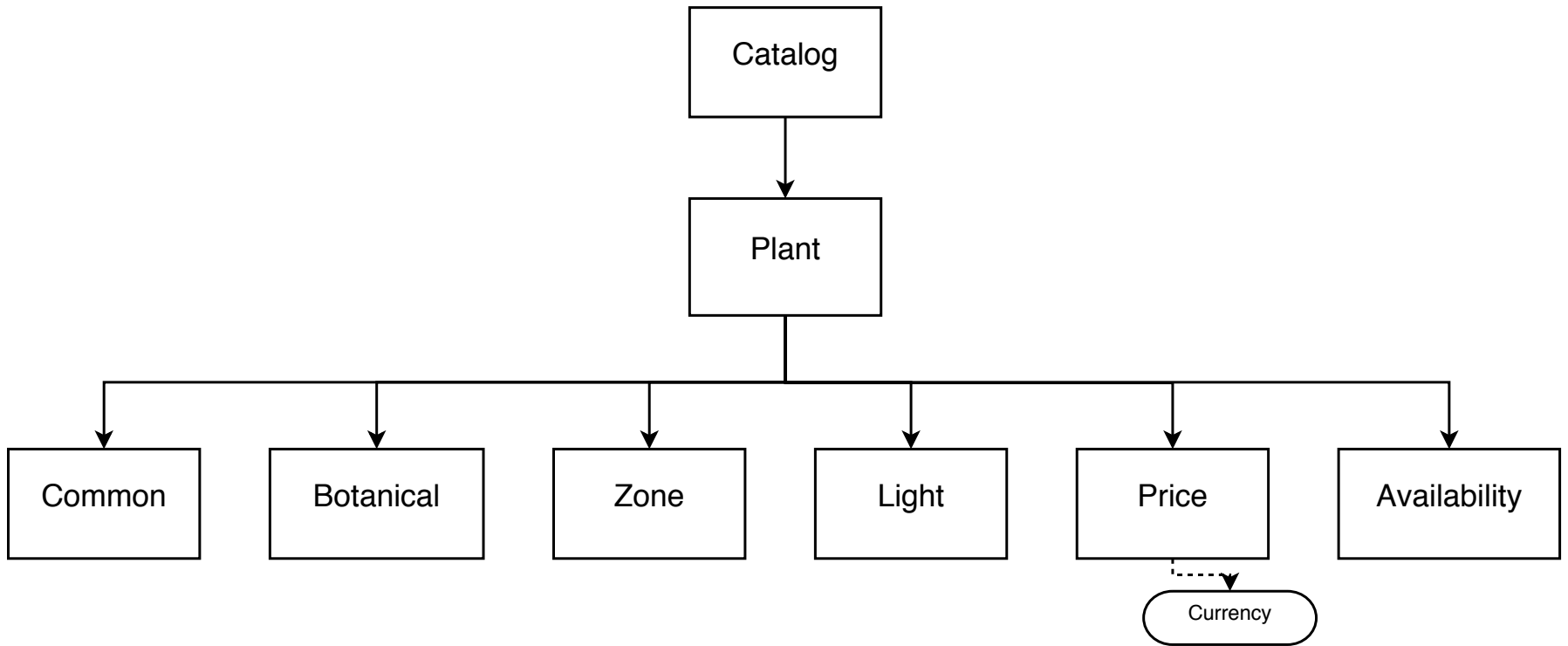
- ④ `<link rel="stylesheet" type="text/css" charset="utf-8" media="all" href="/moin_static185/modern/css/common.css">`

- ④ `<div id="page" lang="en" dir="ltr">`

- ④ ``

XML Basics

- ⑤ `<?xml version="1.0" encoding="UTF-8" ?>`
- ⑤ Tags
 - ⑤ `<tag> content </tag>` or `<tag />`
- ⑤ Attributes
 - ⑤ `<tag attr1="value" attr2="value2"> </tag>`
- ⑤ Nesting
 - ⑤ `<tag 1>content <tag2>nested</tag2></tag1>`
- ⑤ Case sensitive (unlike HTML)



XML Example

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<CATALOG>
  <PLANT>
    <COMMON>Bloodroot</COMMON>
    <BOTANICAL>Sanguinaria canadensis</BOTANICAL>
    <ZONE>4</ZONE>
    <LIGHT>Mostly Shady</LIGHT>
    <PRICE CURRENCY="dollar">2.44</PRICE>
    <AVAILABILITY>031599</AVAILABILITY>
  </PLANT>
  <PLANT>
    <COMMON>Columbine</COMMON>
    <BOTANICAL>Aquilegia canadensis</BOTANICAL>
    <ZONE>3</ZONE>
    <LIGHT>Mostly Shady</LIGHT>
    <PRICE CURRENCY="dollar" >9.37</PRICE>
    <AVAILABILITY>030699</AVAILABILITY>
  </PLANT>
</CATALOG>
```

XML in Python

🌐 xml.sax

- 🌐 Simple API for XML (W3C)
- 🌐 Parse XML files sequentially, callbacks

🌐 xml.dom

- 🌐 Document Object Model (W3C)
- 🌐 View XML as a single hierarchical document

🌐 xml.etree

- 🌐 Python specific, similar to DOM
- 🌐 Easier to use !

Using ElementTree

- ⑤ `import xml.etree.cElementTree` (or `xml.etree.ElementTree`)
- ⑤ `et=xml.etree.cElementTree.parse("xml_example.xml")`
- ⑤ `et=xml.etree.cElementTree.fromstring("XML CODE")`
- ⑤ `e=et.getroot()` - The root object in the XML file
- ⑤ `e[n]` - Children of this element
- ⑤ `e.items()` or `e.attrib` - An element's attributes
- ⑤ `e.text` - Unused text between the start and end tags
- ⑤ `e.tag` - The element's tag as a string

XML Schemas

- ④ Schemas Specifications
 - ④ DTD
 - ④ XML Schema
 - ④ RELAX NG
- ④ Specific Schemas/Ontologies
 - ④ <http://www.bioontology.org>
 - ④ http://en.wikipedia.org/wiki/List_of_XML_markup_languages

RESTful Servers

- ④ Generally return XML
- ④ Client-Server - Servers store data, clients display data
- ④ Stateless - URL uniquely identifies display
- ④ Cacheable - Pages must declare whether their data is

Real World Example

⑤ www.pdb.org/pdb/software/rest.do

⑤ Several interfaces provided

RSS

- ④ Small XML files containing frequently updated information. Link to webpage.
- ④ 2 variants, also Atom.
- ④ Required elements (2.0): title, link, description
- ④ Optional elements: language, copyright, managingeditor, webmaster, pubdate, lastbuilddate, category, generator, docs, cloud, ttl, image, rating, textinput, skiphours, skipdays

Homework 7

- Write a program that retrieves something from the web and processes it in some useful fashion (easier if you find an XML or RSS site). Turn in the program and a brief description of what it retrieves and what it does with it.

BRING YOUR LAPTOPS MONDAY !