



SWCON104
Web & Python Programming

Web Crawling and Web Server Development

Department of Software Convergence

Today

- Let's Crawl and Build a Web Server using Python

Practice

- Practice_21_WebServerUsingPython

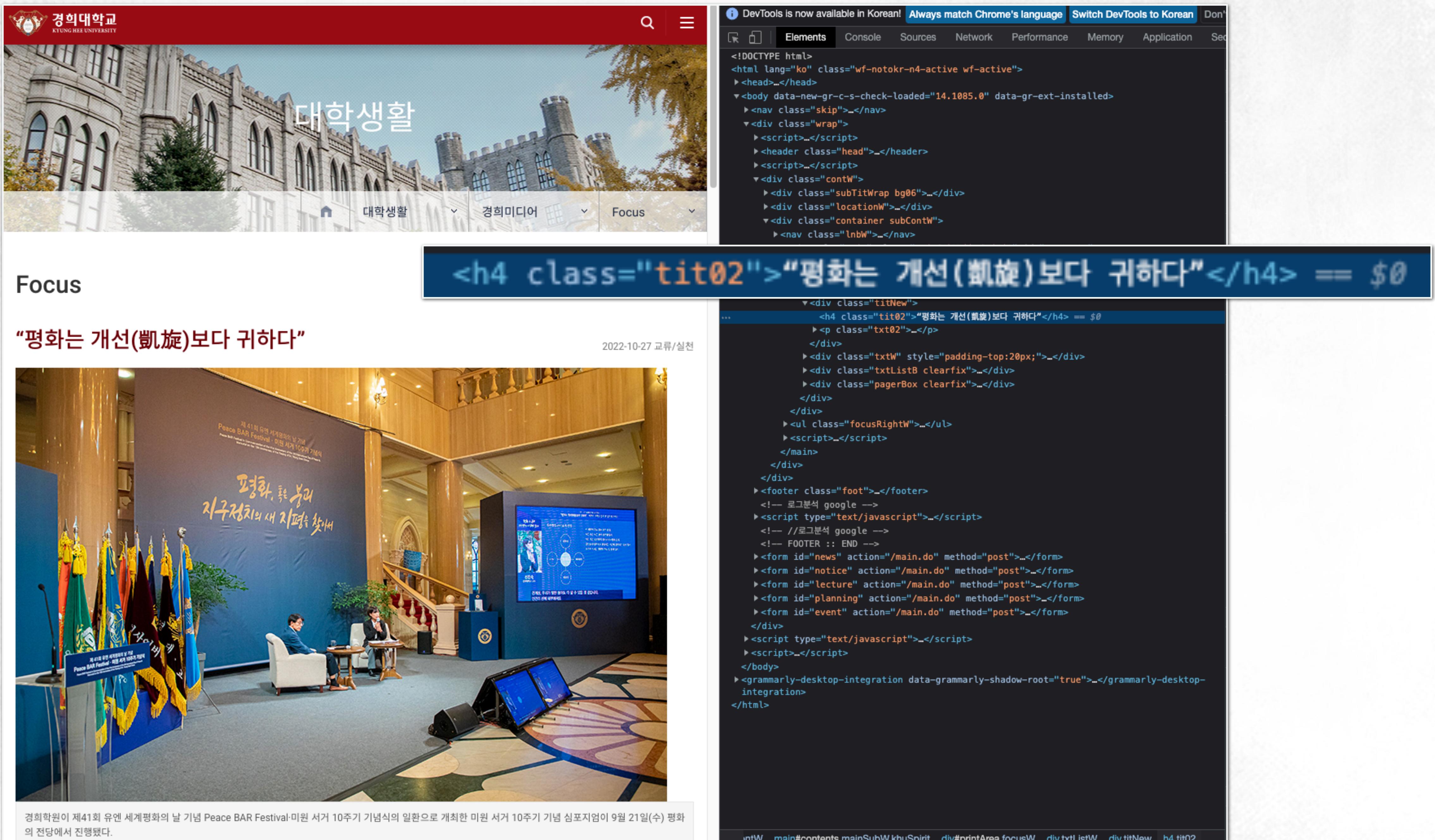
Crawling a web server

- Tools Installation
 - Requests: pip install requests
 - » <https://requests.readthedocs.io/en/latest/>
 - BeautifulSoup: pip install bs4
 - » <https://www.crummy.com/software/BeautifulSoup/>

경희대학교 대학생활 Focus 기사 제목 크롤링

Site access and check source code

<https://khu.ac.kr/kor/focus/detail.do?seq=2164725&page=1&pageSize=5>



The screenshot shows a news article from Kyung Hee University's website. The article is titled "평화는 개선(凱旋)보다 귀하다" (Peace is more valuable than victory) and was published on October 27, 2022. The article features a photograph of two people in a formal setting, likely a lecture or panel discussion, with flags in the background. The DevTools console on the right side of the browser window highlights the title element in the source code.

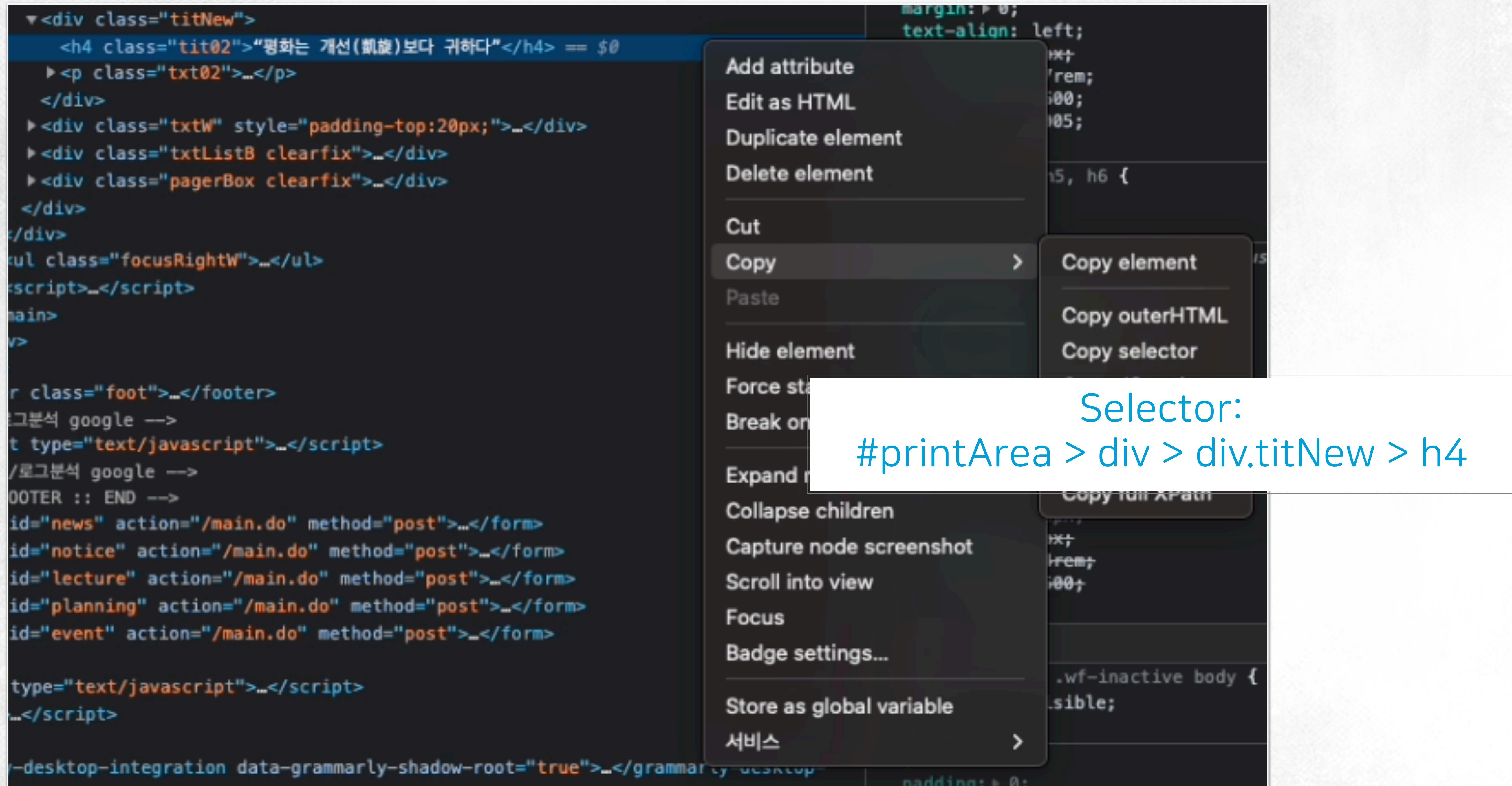
```
<h4 class="tit02">"평화는 개선(凱旋)보다 귀하다"</h4> == $0
```

```
<!DOCTYPE html>
<html lang="ko" class="wf-notokr-n4-active wf-active">
  <head>...
    <body data-new-gr-c-s-check-loaded="14.1085.0" data-gr-ext-installed>
      <nav class="skip"></nav>
      <div class="wrap">
        <script>...</script>
        <header class="head">...</header>
        <script>...</script>
      </div>
      <div class="contW">
        <div class="subTitWrap bg06">...</div>
        <div class="locationW">...</div>
        <div class="container subContW">
          <nav class="lnbW">...</nav>
          ...
        </div>
      </div>
    </body>
  </html>
```

경희대학교 대학생활 Focus 기사 제목 크롤링

Check selectors & tags

- https://khu.ac.kr/kor/focus/detail.do?seq=2164725&page=1&pageSize=5



경희대학교 대학생활 Focus 기사 제목 크롤링

- Create code and crawl target information
 - <https://khu.ac.kr/kor/focus/detail.do?seq=2164725&page=1&pageSize=5>

```
import requests
from bs4 import BeautifulSoup

req = requests.get('https://khu.ac.kr/kor/focus/detail.do?seq=2164725&page=1&pageSize=5')
html = req.text
soup = BeautifulSoup(html, 'html.parser')
my_titles = soup.select(
    '#printArea > div > div.titNew > h4'
)
for title in my_titles:
    print(title.text)
```

Code
webCrawling_1.py

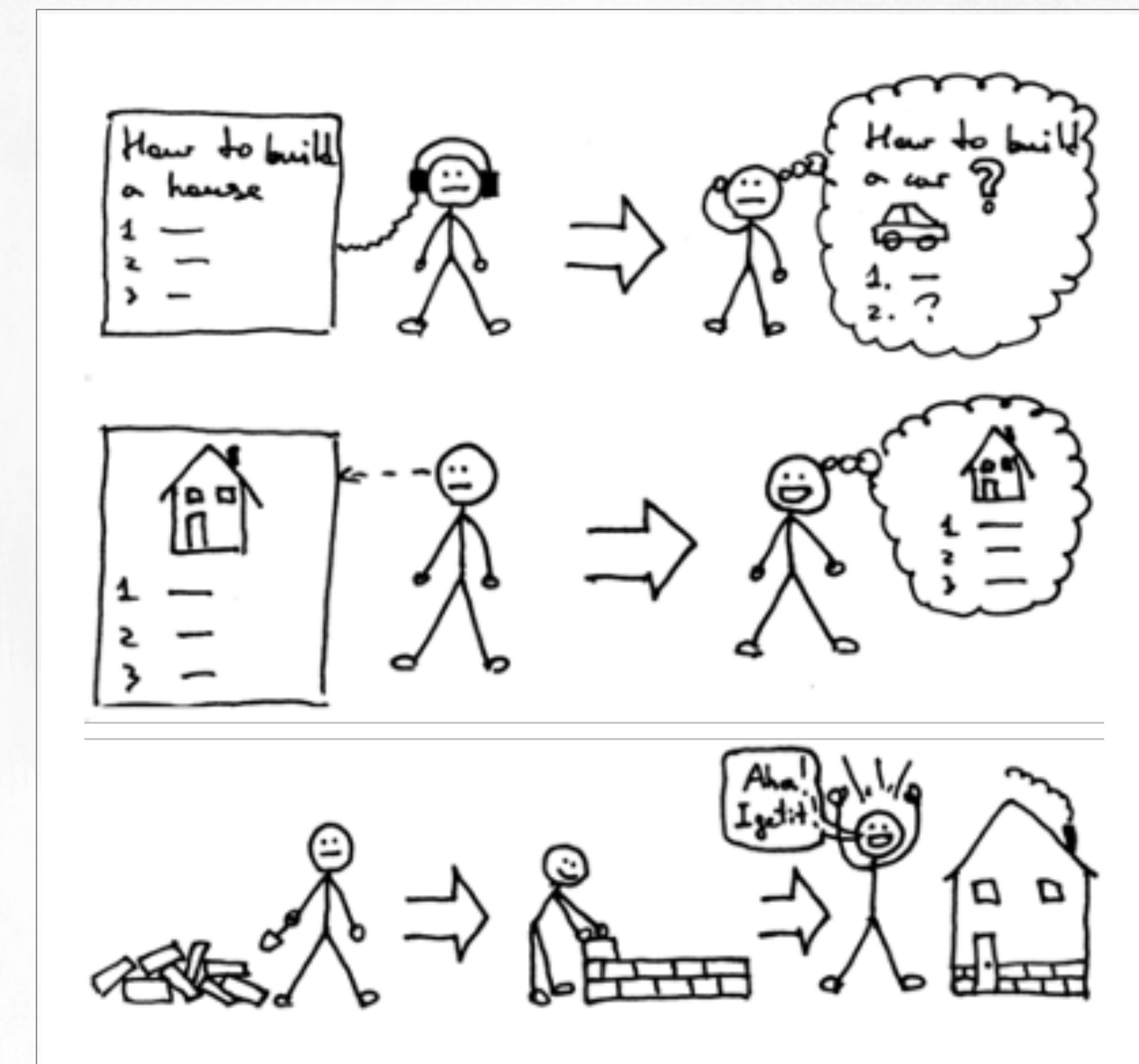
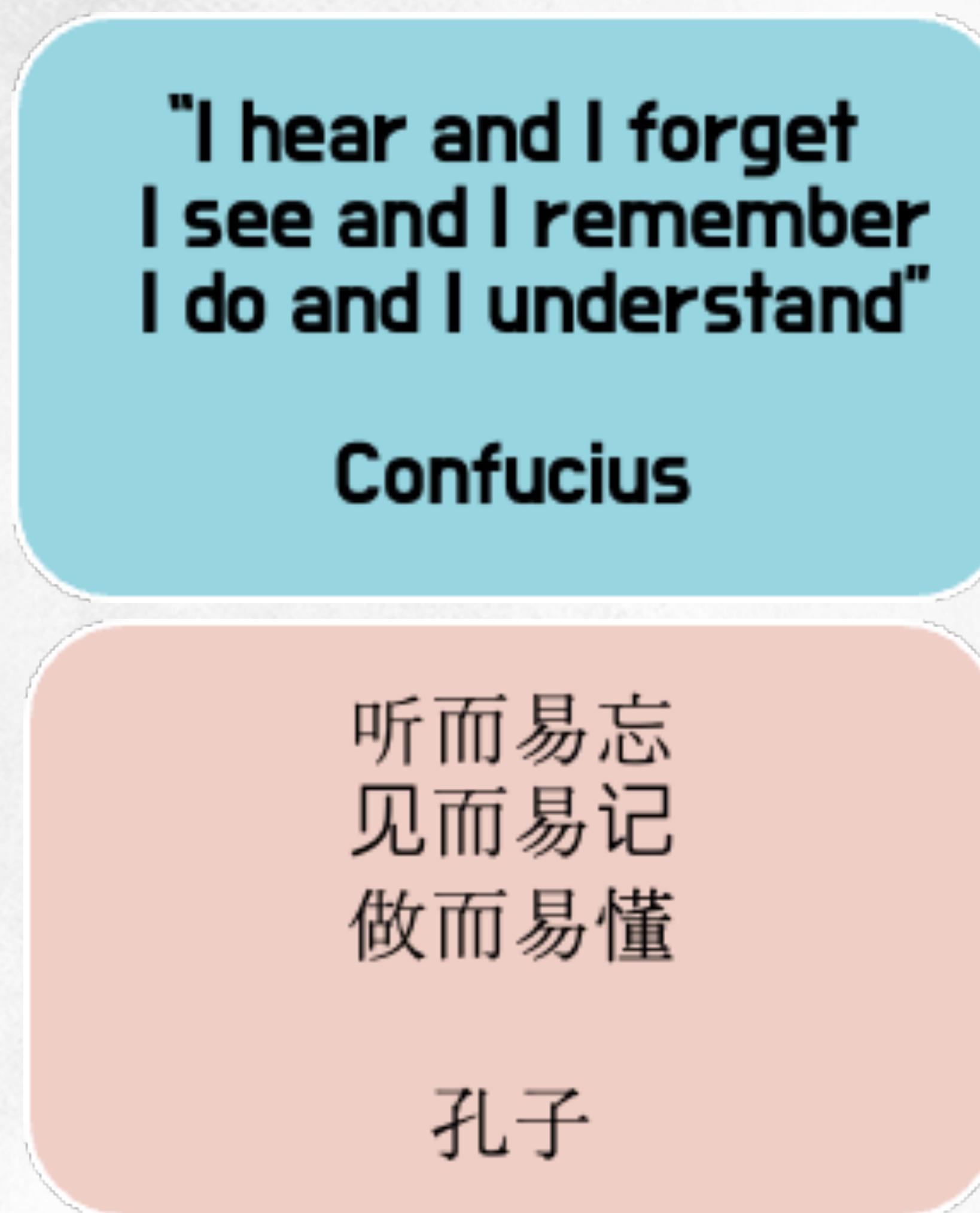
Output

```
[drsungwon~$ python webCrawling_1.py
"평화는 개선(凱旋)보다 귀하다"
```

Building a web server

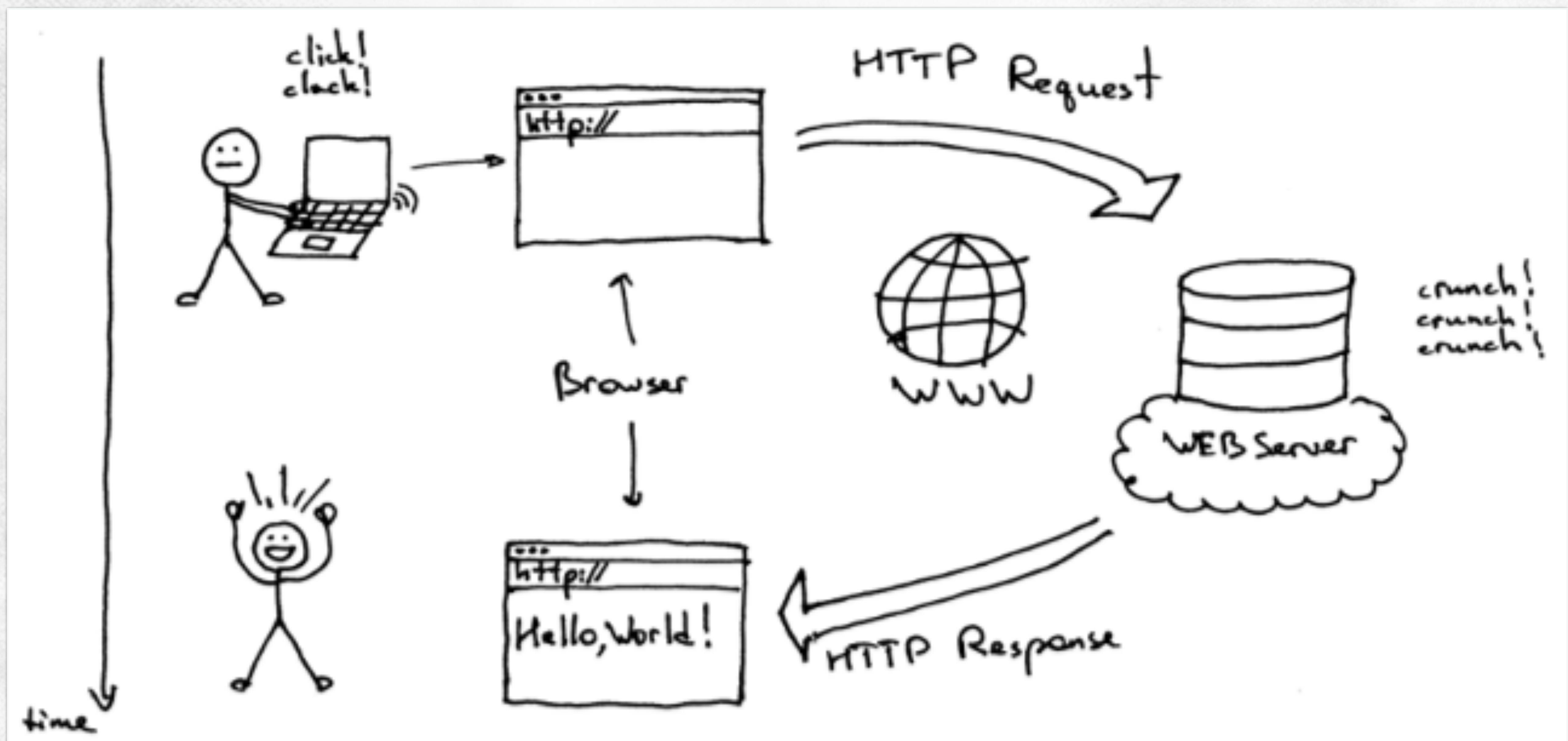
- Why?

- Understands conventional web services in view of developer
- Develops HTML/CSS/JS web service server using python
- Develops server for backend and micro services



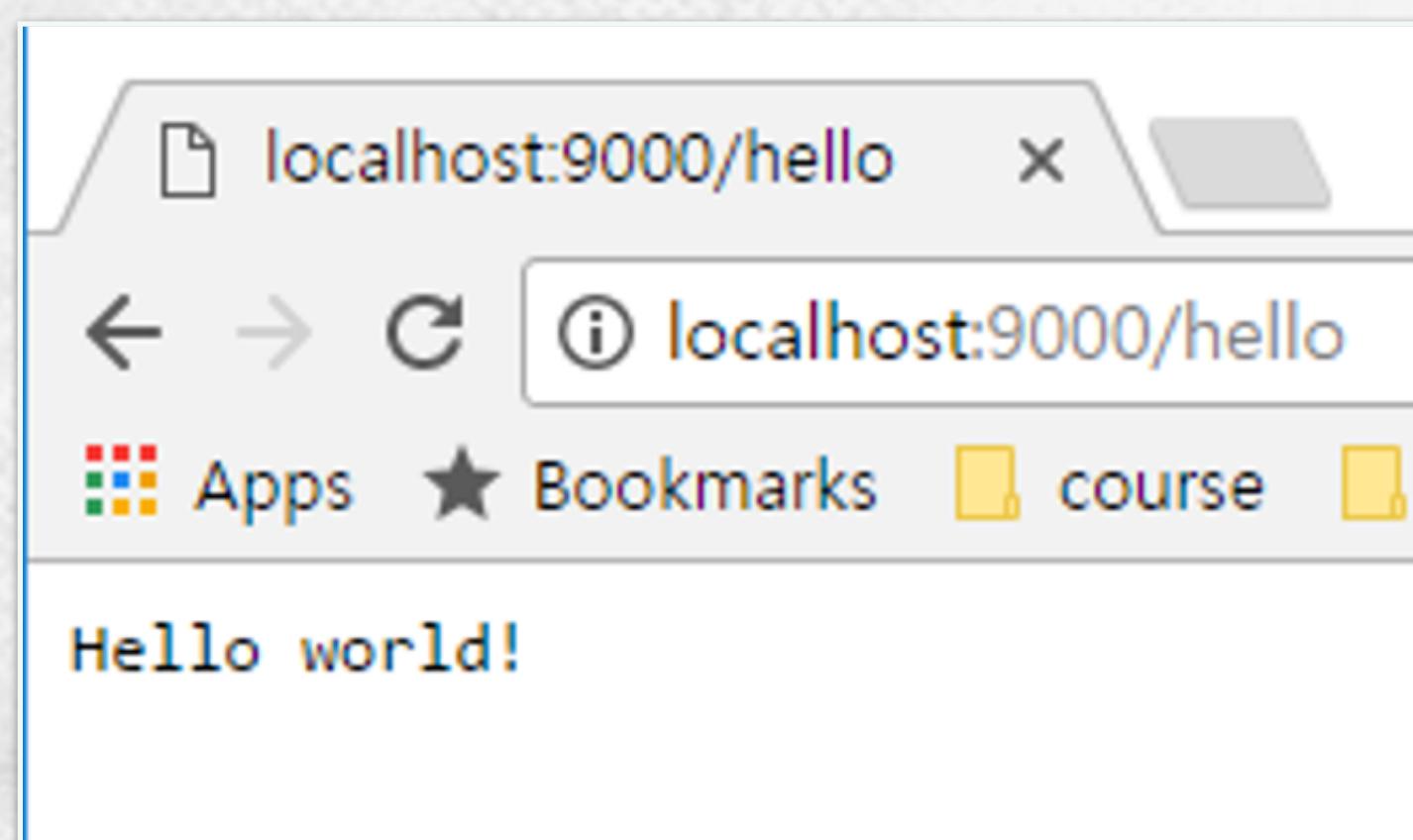
What is a web server?

- Waits for a client to send a request
- When it receives a request, it generates a response and sends it back to the client



Simple web server: webserver1.py

- A very simple implementation of a web server in Python



```
import socket

HOST, PORT = '', 9000

listen_socket = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
listen_socket.setsockopt(socket.SOL_SOCKET, socket.SO_REUSEADDR, 1)
listen_socket.bind((HOST, PORT))
listen_socket.listen(1)
print ('Serving HTTP on port', PORT, '...')
while True:
    client_connection, client_address = listen_socket.accept()
    request = str(client_connection.recv(1024), 'utf-8')
    print (request)

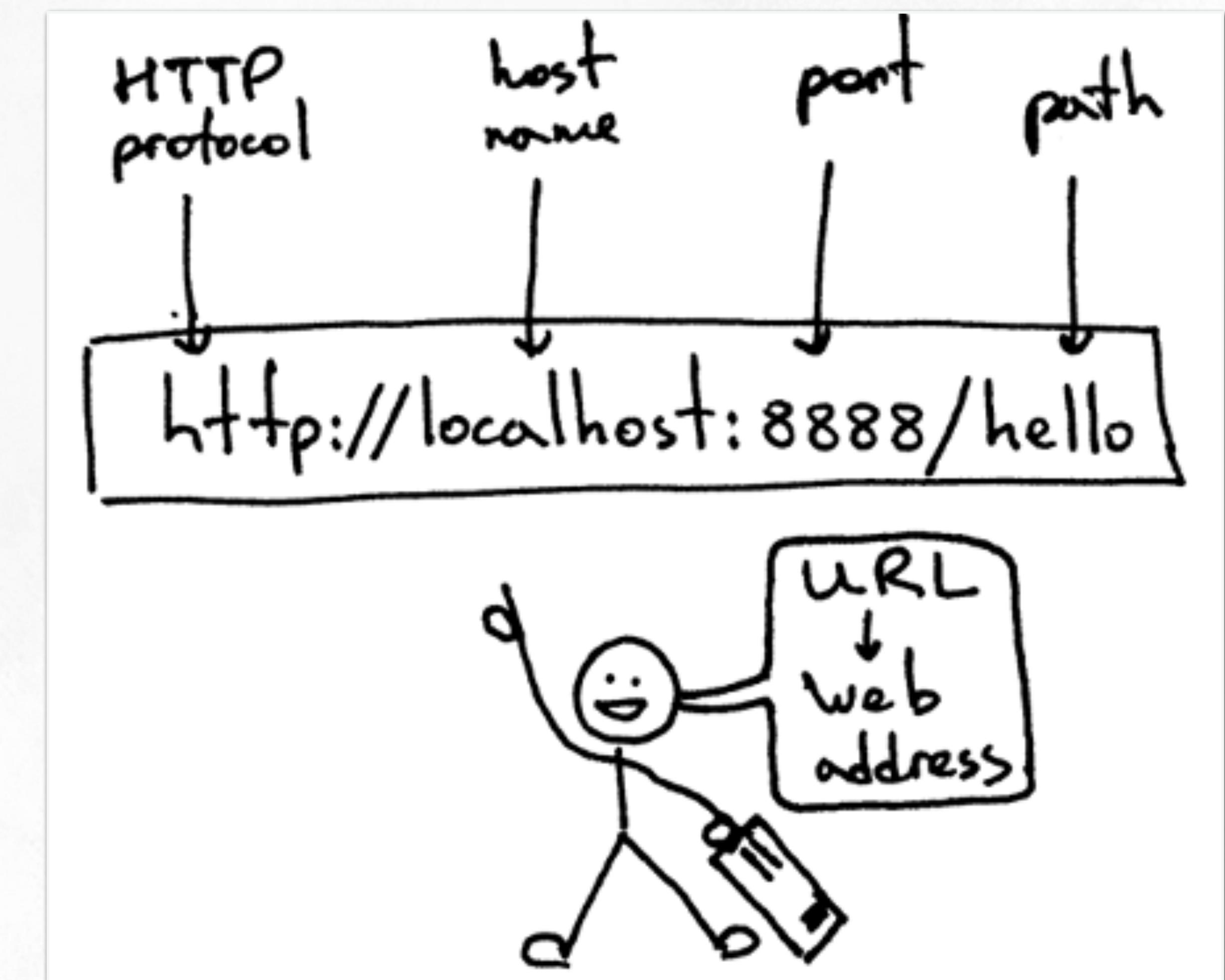
    http_response = """\
HTTP/1.1 200 OK

Hello world!
"""

    client_connection.sendall(bytes(http_response, 'utf-8'))
    client_connection.close()
```

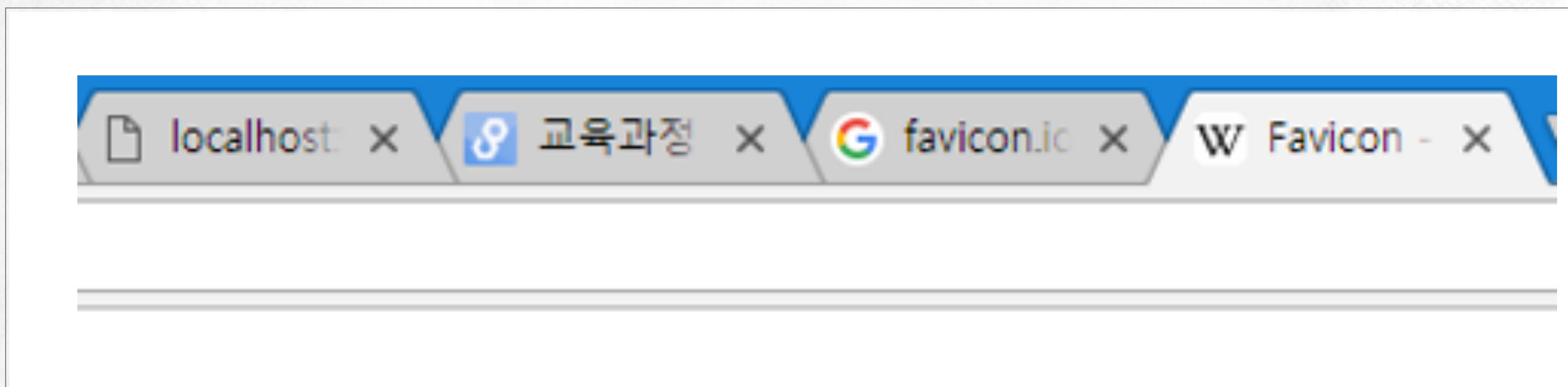
webserver1: How it works?

- The browser:
 - Establishes a TCP connection with the web server
 - Sends an HTTP request over the TCP connection to the server
 - Waits for the server to send an HTTP response back
 - Receives the response and displays it



webserver1: favicon

- A favicon (short for favorite icon)
 - a shortcut icon, website icon, tab icon, URL icon, bookmark icon
 - a file containing one or more small icons, associated with a particular website or webpage



webserver1: favicon

```
Serving HTTP on port 9000 ...
```

```
GET /hello HTTP/1.1
Host: localhost:9000
Connection: keep-alive
Upgrade-Insecure-Requests: 1
User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML,
like Gecko) Chrome/66.0.3359.181 Safari/537.36
Accept: text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,image/a
png,*/*;q=0.8
Accept-Encoding: gzip, deflate, br
Accept-Language: ko-KR,ko;q=0.9,en-US;q=0.8,en;q=0.7
```

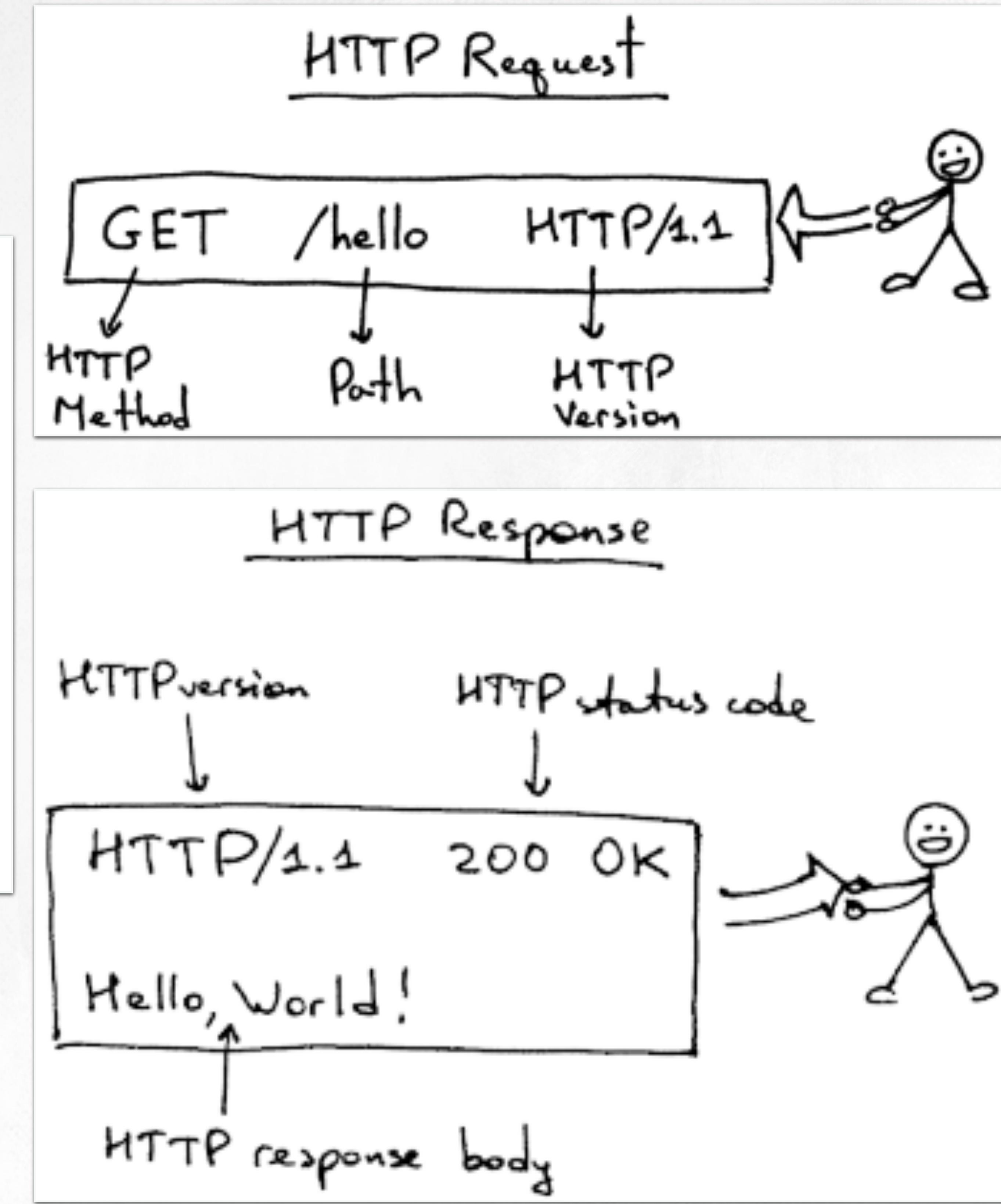
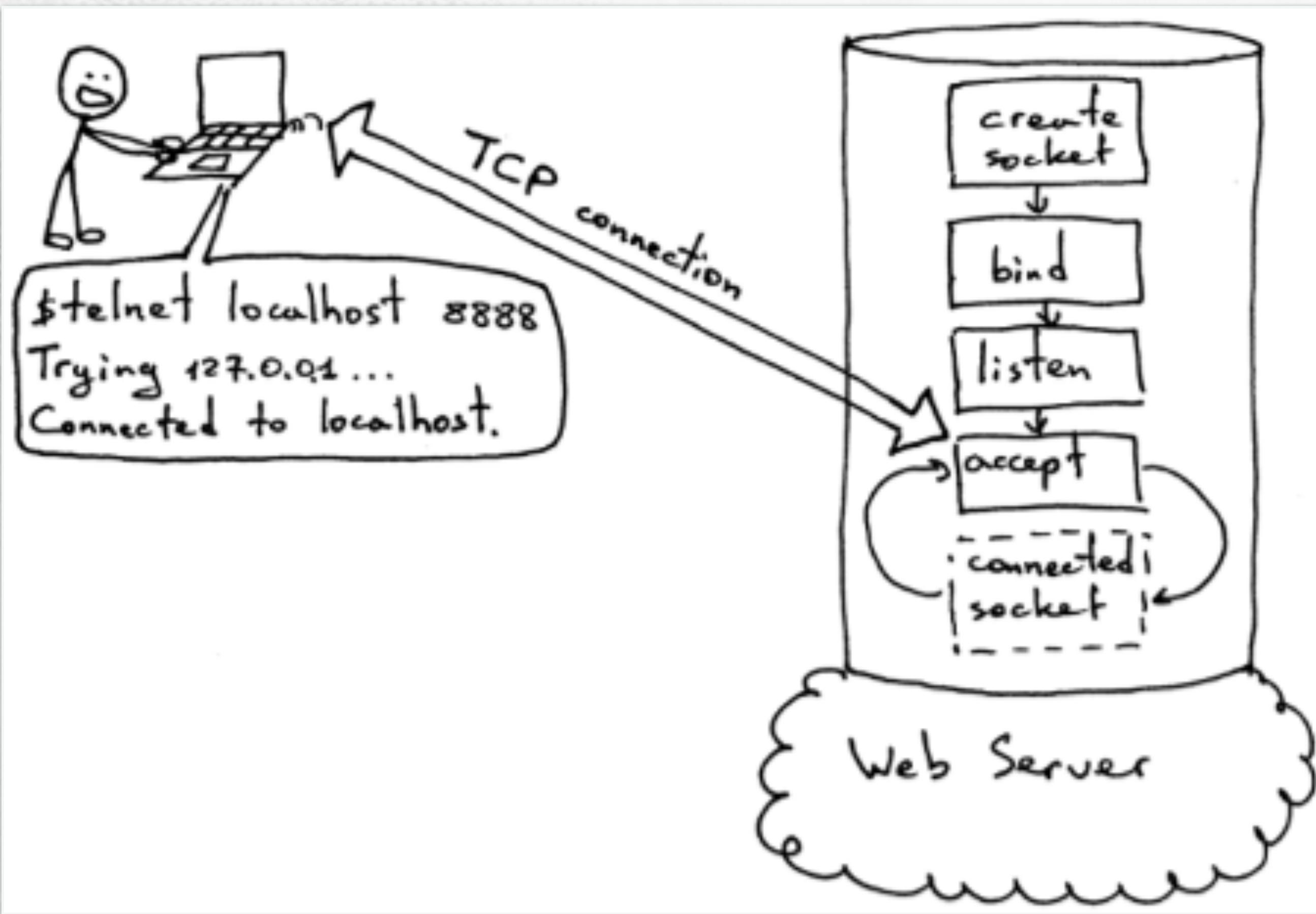
```
GET /favicon.ico HTTP/1.1
Host: localhost:9000
Connection: keep-alive
User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML,
like Gecko) Chrome/66.0.3359.181 Safari/537.36
Accept: image/webp,image/apng,image/*,*/*;q=0.8
Referer: http://localhost:9000/hello
Accept-Encoding: gzip, deflate, br
Accept-Language: ko-KR,ko;q=0.9,en-US;q=0.8,en;q=0.7
```

**GET request comes in twice:
once for /hello and another for /favicon.ico**

webserver1: TCP connection

- TCP is Transmission Control Protocol
- Client and the server establish a TCP connection under HTTP

webserver1: stage by stage



webserver1: code explanation

- Import socket module
 - HOST = "" # service apply to any host name
 - PORT = 9000 # open port, 80 or 8080 are common to avoid conflict
 - <variable> = **socket.socket(<family>, <type>)**
-
- **Socket**: programming interface to TCP and UDP
 - Socket **families**:
 - AF_INET: IPv4 protocols
 - AF_INET6: IPv6 protocols
 - AF_UNIX: UNIX domain protocols

webserver1: code explanation

- Setting socket options
 - <socket_object>.setsockopt(**level**, option_name, value)
- **Level**: the categories of options.
 - SOL_SOCKET: socket-level options
 - IPPROTO_IP: protocol numbers
 - Available options are determined by your OS and IP version
- Binding the port to the socket
- Tell the computer to wait and to listen on that port

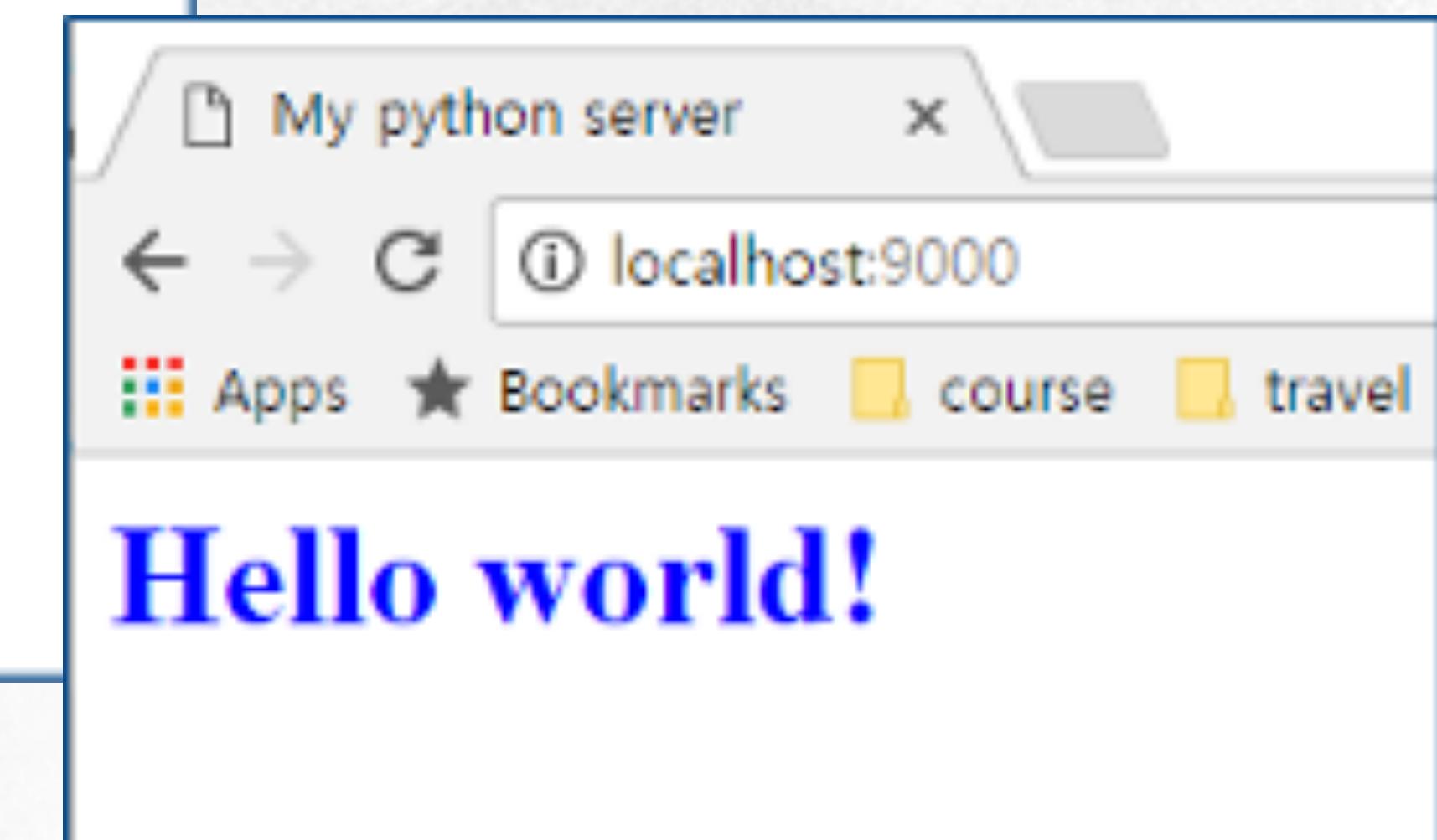
webserver1: handling a server request

- When request is made,
- The server accepts the request, and
- Sends data (a response)
 - Data
 - » First line: a status line (protocol, protocol version, message number, and status)
 - » Two new line characters ("\\n\\n") to distinguish the protocol information from the page content
 - » Then the rest of the data
- Close the server socket

```
while True:  
    client_connection, client_address = listen_socket.accept()  
    request = str(client_connection.recv(1024), 'utf-8')  
    print (request)  
  
    http_response = """\nHTTP/1.1 200 OK  
  
Hello world!  
"""  
  
    client_connection.sendall(bytes(http_response, 'utf-8'))  
    client_connection.close()
```

webserver1: handling a server request

```
while True:  
    client_connection, client_address = listen_socket.accept()  
    request = str(client_connection.recv(1024), 'utf-8')  
    print (request)  
  
    http_response = """\nHTTP/1.1 200 OK  
  
Hello world!  
"""  
  
    client_connection.sendall(bytes(http_response, 'utf-8'))  
    client_connection.close()
```



```
http_response = """\nHTTP/1.1 200 OK  
  
<html><head><title>My python server</title></head><body><H1 style="color:blue">Hello world!</H1></body></html>  
"""
```

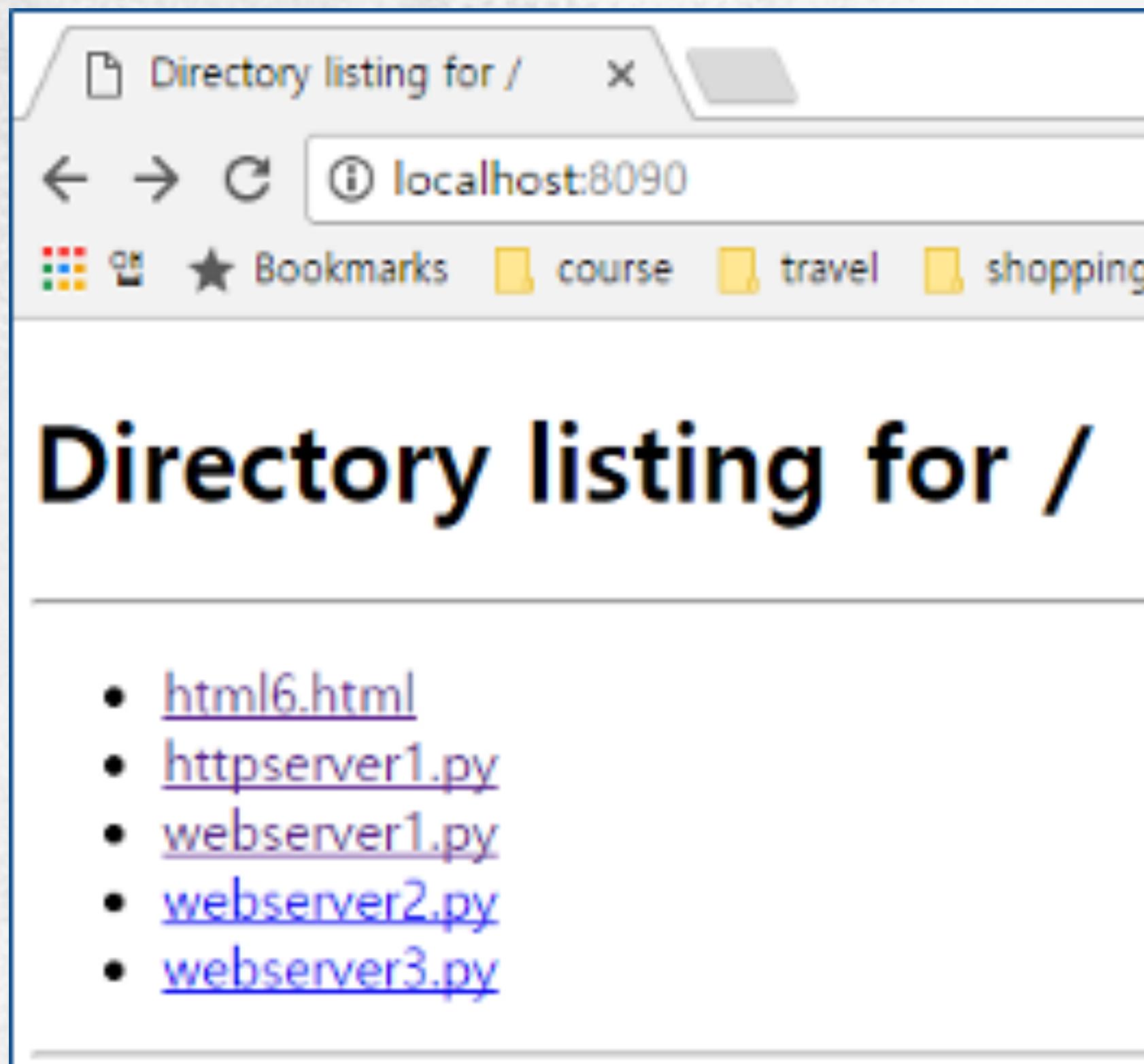
webserver1: send HTML files

```
while True:  
    client_connection, client_address = listen_socket.accept()  
    request = str(client_connection.recv(1024), 'utf-8')  
    print (request)  
  
    http_response = """\nHTTP/1.1 200 OK  
  
Hello world!  
"""  
    client_connection.sendall(bytes(http_response, 'utf-8'))  
    client_connection.close()
```

```
while True:  
    client_connection, client_address = listen_socket.accept()  
    request = str(client_connection.recv(1024), 'utf-8')  
    print (request)  
  
    http_response = "HTTP/1.1 200 OK\n\n"  
    file = open('html6.html', 'r+b')  
    client_connection.sendall(bytes(http_response, 'utf-8'))  
    client_connection.sendfile(file)  
  
    file.close()  
    client_connection.close()
```

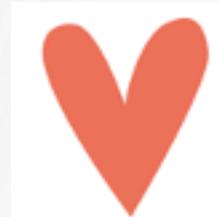
httpserver1: simple HTTP server

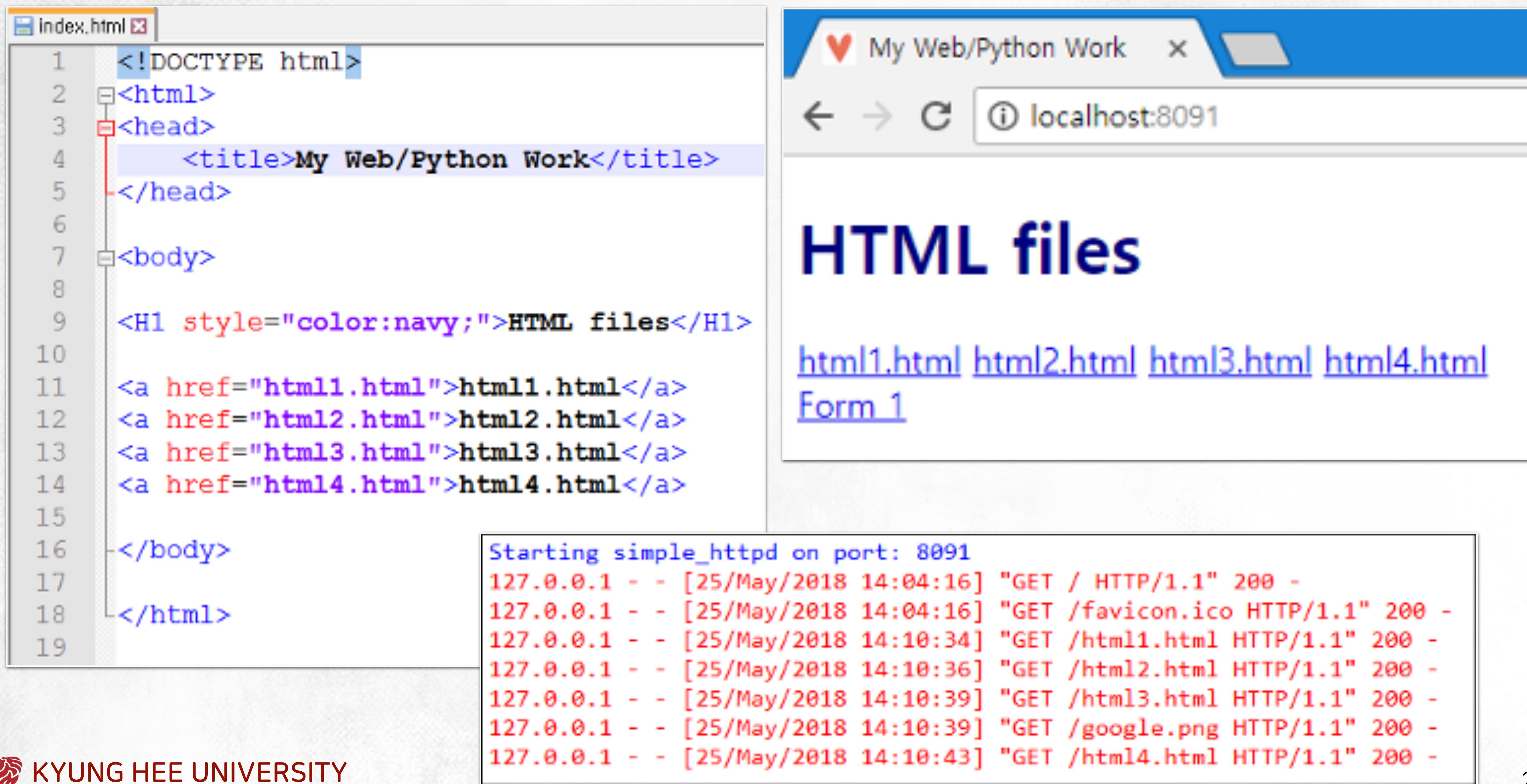
- A very simple implementation of a HTTP server in Python



```
from http.server import SimpleHTTPRequestHandler, HTTPServer
port = 8090
server_address = ('', port)
httpd = HTTPServer(server_address, SimpleHTTPRequestHandler)
print("Starting simple_httpd on port: " + str(httpd.server_port))
httpd.serve_forever()
```

httpserver2: make an index page

- index.html
- Link to other pages
- Save an favicon.ico  to the same folder



The image shows a development environment with three main components:

- Code Editor:** An IDE window titled "index.html" displays the following HTML code:

```
1 <!DOCTYPE html>
2 <html>
3 <head>
4   <title>My Web/Python Work</title>
5 </head>
6
7 <body>
8
9   <H1 style="color:navy;">HTML files</H1>
10
11  <a href="html1.html">html1.html</a>
12  <a href="html2.html">html2.html</a>
13  <a href="html3.html">html3.html</a>
14  <a href="html4.html">html4.html</a>
15
16 </body>
17
18 </html>
```
- Browser:** A browser window titled "My Web/Python Work" shows the rendered HTML. The title bar includes a red heart icon. The address bar shows "localhost:8091". The content area displays "**HTML files**" in large blue font, followed by four blue links: "html1.html", "html2.html", "html3.html", and "html4.html", and a link to "Form 1".
- Terminal:** A terminal window at the bottom shows the log output of the "simple_httpd" server:

```
Starting simple_httpd on port: 8091
127.0.0.1 - - [25/May/2018 14:04:16] "GET / HTTP/1.1" 200 -
127.0.0.1 - - [25/May/2018 14:04:16] "GET /favicon.ico HTTP/1.1" 200 -
127.0.0.1 - - [25/May/2018 14:10:34] "GET /html1.html HTTP/1.1" 200 -
127.0.0.1 - - [25/May/2018 14:10:36] "GET /html2.html HTTP/1.1" 200 -
127.0.0.1 - - [25/May/2018 14:10:39] "GET /html3.html HTTP/1.1" 200 -
127.0.0.1 - - [25/May/2018 14:10:39] "GET /google.png HTTP/1.1" 200 -
127.0.0.1 - - [25/May/2018 14:10:43] "GET /html4.html HTTP/1.1" 200 -
```

httpserver2: explanation using inheritance

- A very simple implementation of a HTTP server in Python

```
from http.server import HTTPServer, SimpleHTTPRequestHandler

class testHTTPServer_RequestHandler(SimpleHTTPRequestHandler):
    def do_GET(self):
        super().do_GET()
        print("do_get")

port = 8095
httpd = HTTPServer(('', port), testHTTPServer_RequestHandler)
print("Starting simple_httpd on port: " + str(httpd.server_port))
httpd.serve_forever()
```

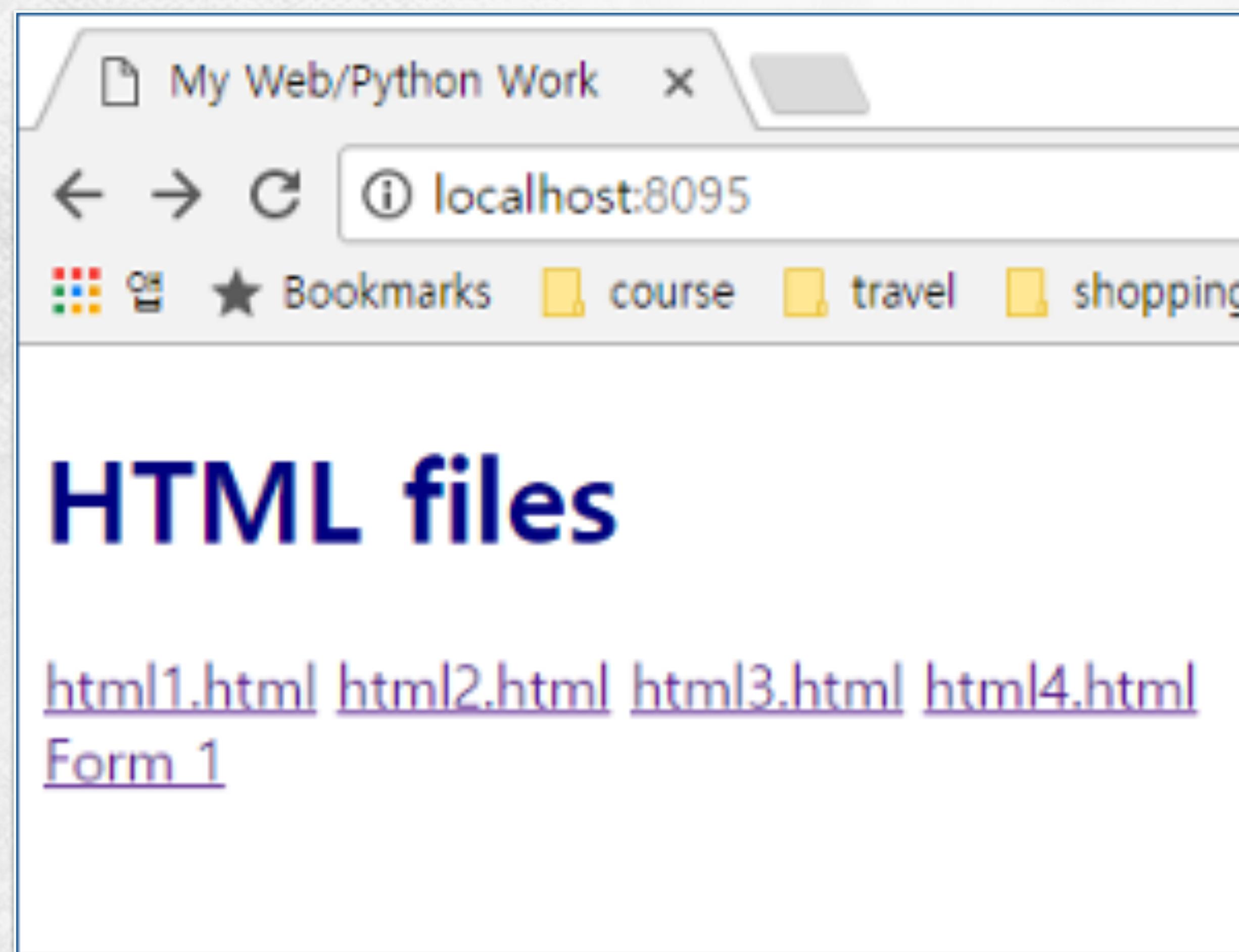
httpserver2: explanation using inheritance

Standard HTTP handler in Python

```
servertest01.py x server.py x
523     v: (v.phrase, v.description)
524     for v in HTTPStatus.__members__.values()
525 }
526
527
528 class SimpleHTTPRequestHandler(BaseHTTPRequestHandler):
529
530     """Simple HTTP request handler with GET and HEAD commands.
531
532     This serves files from the current directory and any of its
533     subdirectories. The MIME type for files is determined by
534     calling the .guess_type() method.
535
536     The GET and HEAD requests are identical except that the HEAD
537     request omits the actual contents of the file.
538
539     """
540
541     server_version = "SimpleHTTP/" + __version__
542
543     def do_GET(self):
544         """Serve a GET request."""
545         f = self.send_head()
546         if f:
547             try:
548                 self.copyfile(f, self.wfile)
549             finally:
550                 f.close()
551
552     def do_HEAD(self):
553         """Serve a HEAD request."""
554         f = self.send_head()
555         if f:
556             f.close()
557
558
559     def send_head(self):
560         """Common code for GET and HEAD commands.
561
562         This sends the response code and MIME headers.
563
564         Return value is either a file object (which has to be copied
565         to the outputfile by the caller unless the command was HEAD,
566         and must be closed by the caller under all circumstances), or
567         None, in which case the caller has nothing further to do.
568
569         """
570         path = self.translate_path(self.path)
571         f = None
572         if os.path.isdir(path):
573             parts = urllib.parse.urlsplit(self.path)
574             if not parts.path.endswith('/'):
575                 # redirect browser - doing basically what apache does
576                 self.send_response(HTTPStatus.MOVED_PERMANENTLY)
577                 new_parts = (parts[0], parts[1], parts[2] + '/',
578                             parts[3], parts[4])
579                 new_url = urllib.parse.urlunsplit(new_parts)
580                 self.send_header("Location", new_url)
581                 self.end_headers()
582             return None
583
584             for index in "index.html", "index.htm":
585                 index = os.path.join(path, index)
586                 if os.path.exists(index):
587                     path = index
588                     break
589
590             else:
591                 return self.list_directory(path)
592
593         ctype = self.guess_type(path)
```

httpserver3: HTML form exercise

- Edit index.html to link a form page (formsInput1.html)
- Check how the server receives user input values



The screenshot shows a web browser window titled "HTML Forms 1". The address bar displays "localhost:8095/formsInput1.html". The main content area is a form with the following fields:

- First name:
- Last name:
- User password:
- Food preferences:
 - Pizza
 - Tacos
 - Salad
-

httpserver3: fetch and process user input

The image shows a split-screen environment. On the left, a code editor displays the Python script `httpserver3.py`. On the right, a web browser window shows an HTML form titled "HTML Forms 1". Below the browser is a terminal window displaying the server's log.

Code Editor (httpserver3.py):

```
httpserver3.py - D:\course\2017s_python\myserver\httpserver3.py (3.5.3)
File Edit Format Run Options Window Help
from http.server import HTTPServer, SimpleHTTPRequestHandler
from urllib.parse import parse_qs, urlparse

class testHTTPServer_RequestHandler(SimpleHTTPRequestHandler):

    def do_GET(self):
        url = self.path
        form = parse_qs(urlparse(url).query)
        if (form != {}):
            self.process_form(form)

        super().do_GET()
        print("do_get")

    def process_form(self,form):
        if 'food' in form:
            if form['food'][0] == 'Pizza':
                print(form['firstname'][0] + ", call Dominos tonight!")
            elif form['food'][0] == 'Tacos':
                print(form['firstname'][0] + ", go to TacoBell tonight!")
            elif form['food'][0] == 'Salad':
                print(form['firstname'][0] + ", have a Caesar Salad tonight!")

port = 9095
httpd = HTTPServer(('', port), testHTTPServer_RequestHandler)
print("Starting simple_httpd on port: " + str(httpd.server_port))
httpd.serve_forever()
```

Browser (HTML Forms 1):

localhost:9095/formsInput1.html?firstname=John&lastname=Smith&password=dddd&food=Tacos

First name: John
Last name: Smith
User password:
 Pizza
 Tacos
 Salad

Terminal (Log):

```
Starting simple_httpd on port: 9095
John, go to TacoBell tonight!
127.0.0.1 - - [25/May/2018 14:36:27] "GET /formsInput1.html?firstname=John&lastname=Smith&password=dddd&food=Tacos HTTP/1.1" 200
-
do_get
```

Frameworks (3rd party softwares)



Flask

web development,
one drop at a time

Welcome to Flask's documentation. Get started with [Installation](#) and then get an overview with the [Quickstart](#). There is also a more detailed [Tutorial](#) that shows how to create a small but complete application with Flask. Common patterns are described in the [Patterns for Flask](#) section. The rest of the docs describe each component of Flask in detail, with a full reference in the [API](#) section.

Flask depends on the [Jinja](#) template engine and the [Werkzeug](#) WSGI toolkit. The documentation for these libraries can be found at:

```
from flask import Flask

app = Flask(__name__)

@app.route("/")
def hello_world():
    return "<p>Hello, World!</p>"
```

Frameworks (3rd party softwares)

Latest version: [20.1.0](#)



gunicorn

Gunicorn 'Green Unicorn' is a Python WSGI HTTP Server for UNIX. It's a pre-fork worker model. The Gunicorn server is broadly compatible with various web frameworks, simply implemented, light on server resources, and fairly speedy.

[View source](#) [Download](#)



QUICKSTART
Read the quickstart guide to get started using Gunicorn.



DEPLOYMENT
Learn how to deploy the Gunicorn server.



COMMUNITY
Get in touch with the community.



DOCUMENTATION
Read the documentation to learn more about Gunicorn.

Frameworks (3rd party softwares)

Latest version: [20.1.0](#)

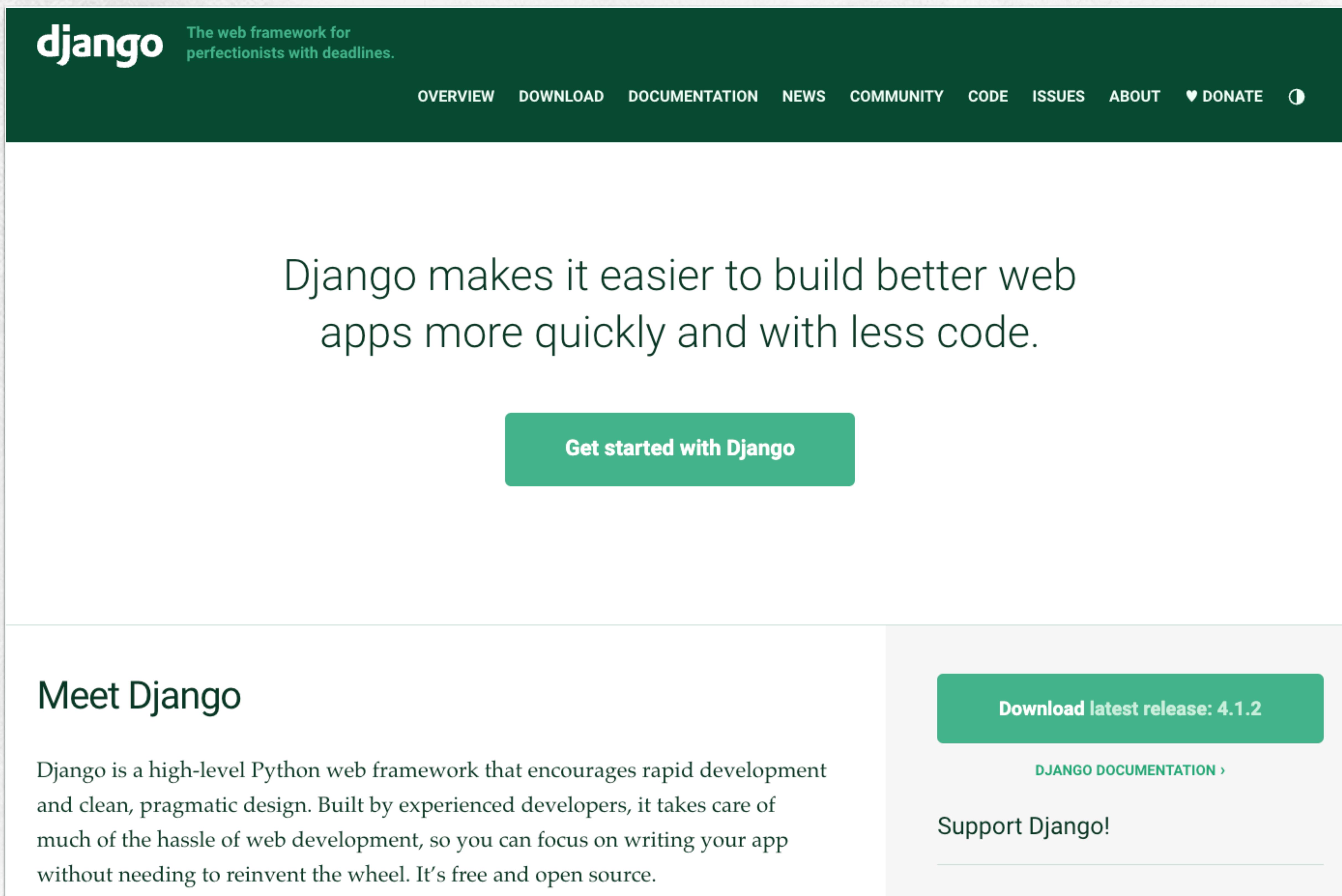


gunicorn

Gunicorn 'Green Unicorn' is a Python WSGI HTTP Server for UNIX. It's a pre-fork worker model. The Gunicorn server is broadly compatible with various web frameworks, simply implemented, light on server resources, and fairly speedy.

```
$ pip install gunicorn
$ cat myapp.py
def app(environ, start_response):
    data = b"Hello, World!\n"
    start_response("200 OK", [
        ("Content-Type", "text/plain"),
        ("Content-Length", str(len(data)))
    ])
    return iter([data])
$ gunicorn -w 4 myapp:app
[2014-09-10 10:22:28 +0000] [30869] [INFO] Listening at: http://127.0.0.1:8000 (30869)
[2014-09-10 10:22:28 +0000] [30869] [INFO] Using worker: sync
[2014-09-10 10:22:28 +0000] [30874] [INFO] Booting worker with pid: 30874
[2014-09-10 10:22:28 +0000] [30875] [INFO] Booting worker with pid: 30875
[2014-09-10 10:22:28 +0000] [30876] [INFO] Booting worker with pid: 30876
[2014-09-10 10:22:28 +0000] [30877] [INFO] Booting worker with pid: 30877
```

Frameworks (3rd party softwares)



The screenshot shows the official Django website. At the top, there's a dark green header bar with the word "django" in white lowercase letters and a tagline "The web framework for perfectionists with deadlines." Below the header are navigation links: OVERVIEW, DOWNLOAD, DOCUMENTATION, NEWS, COMMUNITY, CODE, ISSUES, ABOUT, a "DONATE" button with a heart icon, and a user icon.

The main content area features a large, bold text statement: "Django makes it easier to build better web apps more quickly and with less code." Below this is a teal button with the text "Get started with Django".

On the left side of the main content area, there's a section titled "Meet Django" with a paragraph describing the framework. On the right side, there are two teal buttons: one for "Download latest release: 4.1.2" and another for "DJANGO DOCUMENTATION >". Below these buttons is a link to "Support Django!".

Meet Django

Django is a high-level Python web framework that encourages rapid development and clean, pragmatic design. Built by experienced developers, it takes care of much of the hassle of web development, so you can focus on writing your app without needing to reinvent the wheel. It's free and open source.

Download latest release: 4.1.2

DJANGO DOCUMENTATION >

Support Django!

Frameworks (3rd party softwares)



The web framework for
perfectionists with deadlines.

OVERVIEW DOWNLOAD DOCUMENTATION NEWS COMMUNITY CODE ISSUES ABOUT ♥ DONATE

members/templates/myfirst.html :

```
<!DOCTYPE html>
<html>
<body>

<h1>Hello World!</h1>
<p>Welcome to my first Django project!</p>

</body>
</html>
```

ter web
code.

members/views.py :

```
from django.http import HttpResponse
from django.template import loader

def index(request):
    template = loader.get_template('myfirst.html')
    return HttpResponse(template.render())
```

Meet Django

Django is a high-level Python web framework that encourages rapid development and clean, pragmatic design. It takes much of the hassle of web development, so you can focus on writing your app without needing to reinvent the wheel. It's free and open source.



KYUNG HEE UNIVERSITY



Thank you