

XJCO3011: Web Services and Web Data



Session #2 –

HTTP – The Workhorse of the Web

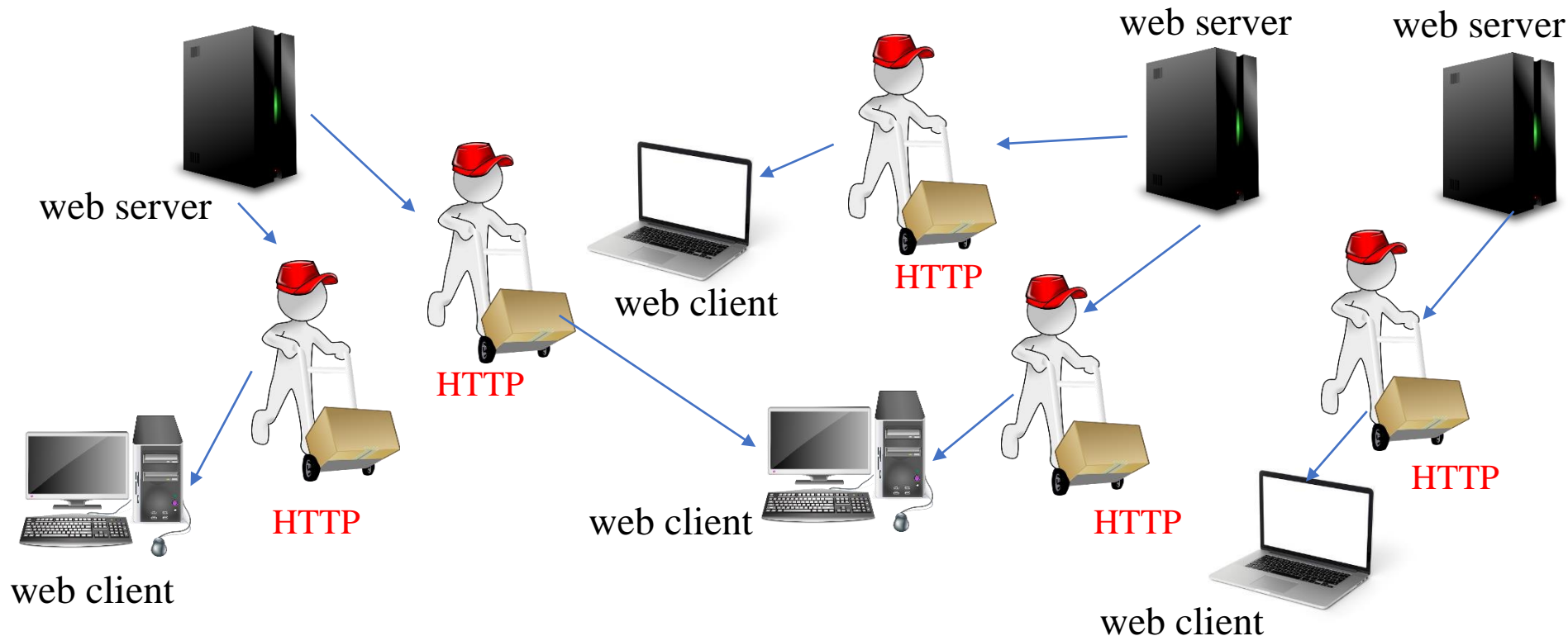
Instructor: Guilin Zhao
Spring 2023

The Hypertext Transfer Protocol (HTTP)

The Multimedia Courier of the Web



- Every day, **billions of multimedia content** (e.g. JPEG images, HTML pages, text files, MPEG movies, WAV audio files, Java applets) **cruise through the Internet** using the HTTP protocol.
- HTTP **moves information quickly** and reliably from **web servers** all around the world to **web browsers** on people's desktops.
- HTTP is an **application layer protocol** that dates back to 1991, but is still the workhorse of the world wide web.
- A **good understanding** of this protocol is **essential** for building web applications, search engines, and RESTful APIs.



A large, abstract blue ink splatter or blotch serves as a background for the title text. The splatter is irregular and textured, with various shades of blue and white, giving it a hand-painted or ink-splashed appearance. It is centered on the slide and occupies a significant portion of the upper half.

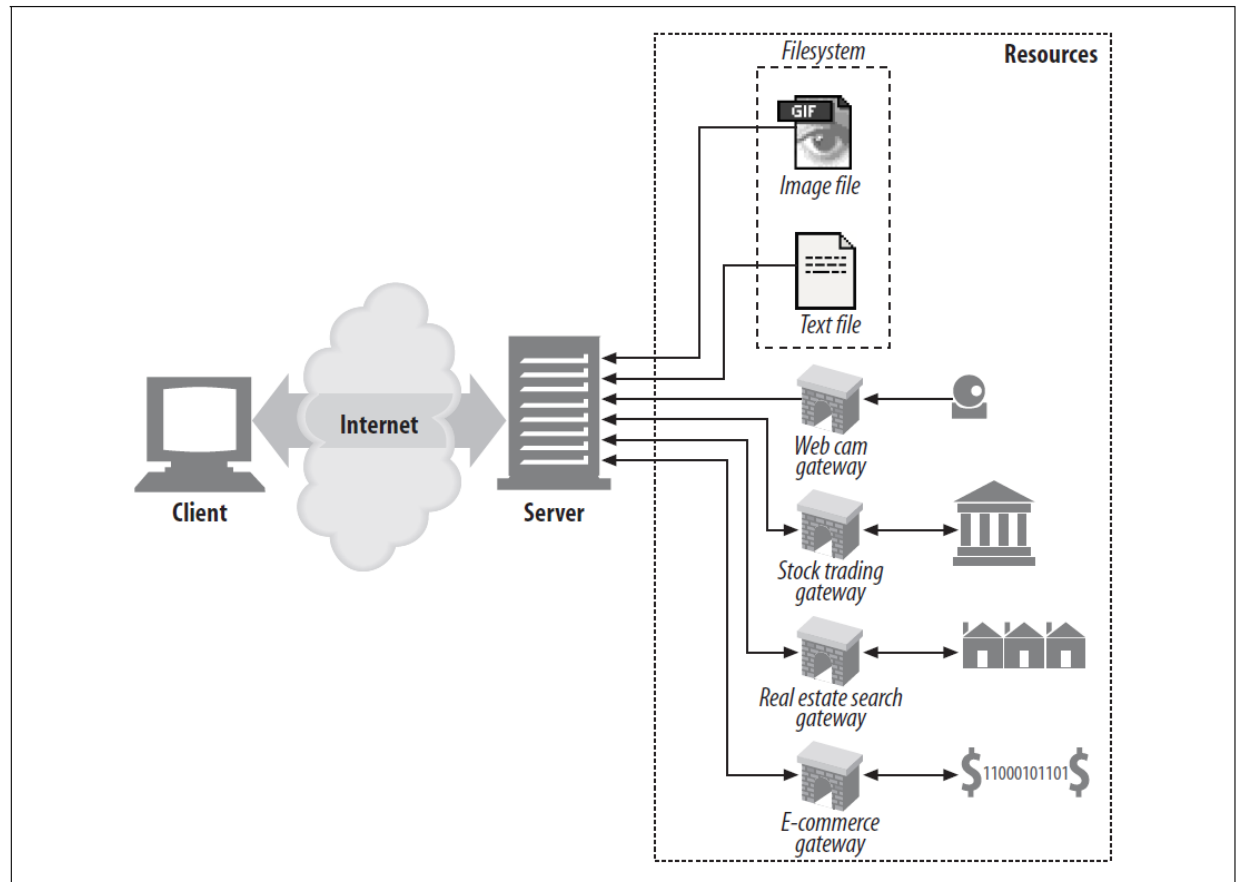
A Quick Introduction to the HTTP Protocol

- A protocol is a set of rules and standards that govern the communication between two or more devices or systems, typically in the context of computer networking.
- HTTP is a protocol used to transfer resources over the web.

- HTTP is a protocol used to transfer resources over the web, and it forms the basis of **web resources**.
- Web resources are various types of files, documents, images, videos, audios, and other files that are transmitted from a web server to a client browser using the HTTP protocol.
- Therefore, HTTP and web resources are closely related, with **HTTP providing a standard way to transfer web resources** so that they can be accessed, shared, and utilized on the web.

Web Resources

- Web servers host *web resources*.
- A web resource is the source of web content.



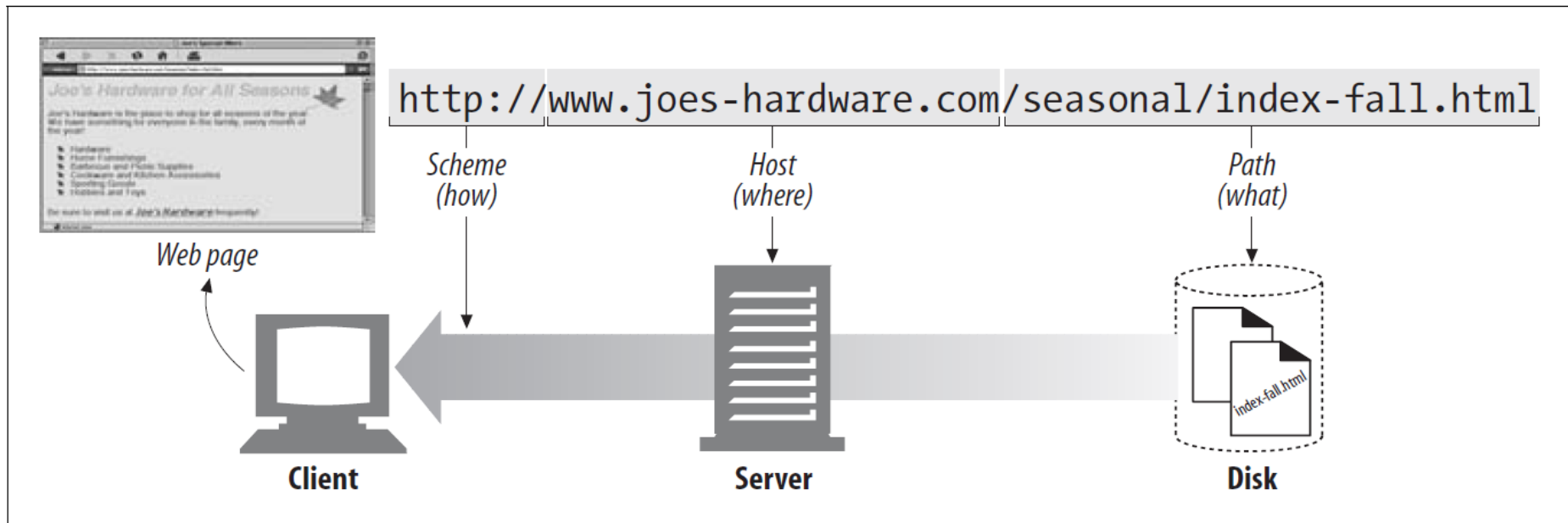
- The simplest kind of web resource is a **static file** on the web server's filesystem (text, HTML, Microsoft Word, Adobe Acrobat, JPEG image, AVI movie, .. etc.).
- However, resources do **NOT** have to be static files. Resources can also be software programs that generate content on demand.

URIs

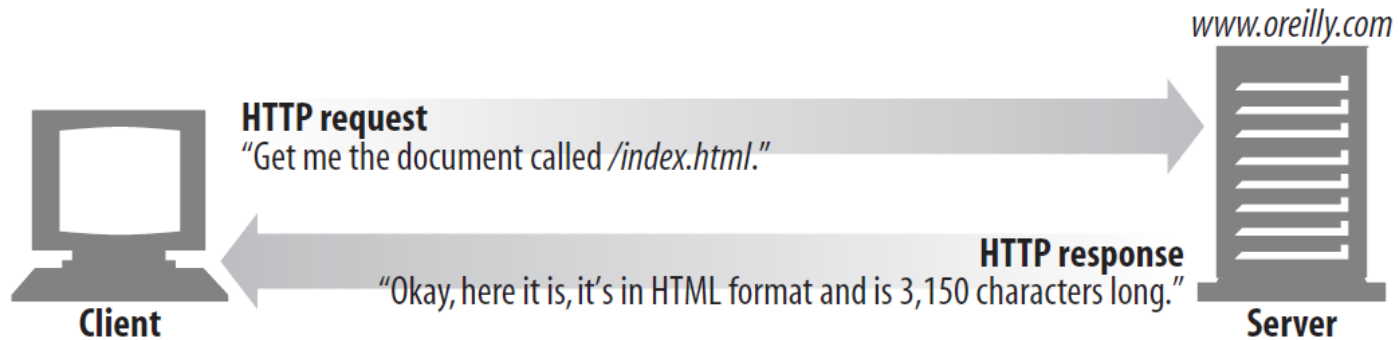
- The **HTTP protocol** defines the rules and format of communication between a client browser and a server, allowing the client browser to request **web resources** and receive them.
- Web resources are typically identified by a Uniform Resource Identifier (URI), and **the client browser uses the URI to send an HTTP request to the server**, which responds by returning the appropriate web resource.

URIs

- Each web resource has a name, so clients can point out what resources they are interested in.
- The resource name is called a *uniform resource identifier*, or **URI**.
- URIs have the general form: “scheme://server location/path”



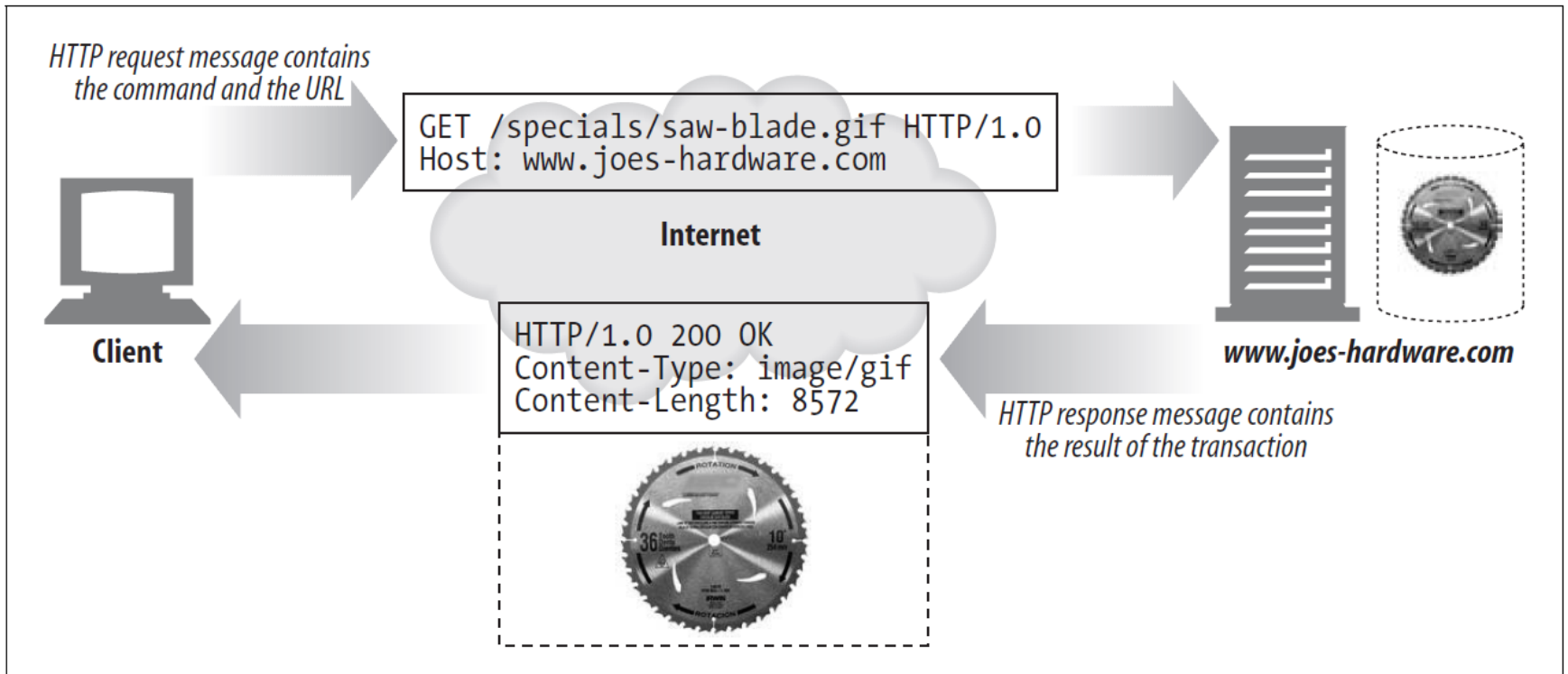
HTTP uses a Client Server Model



HTTP Message Types

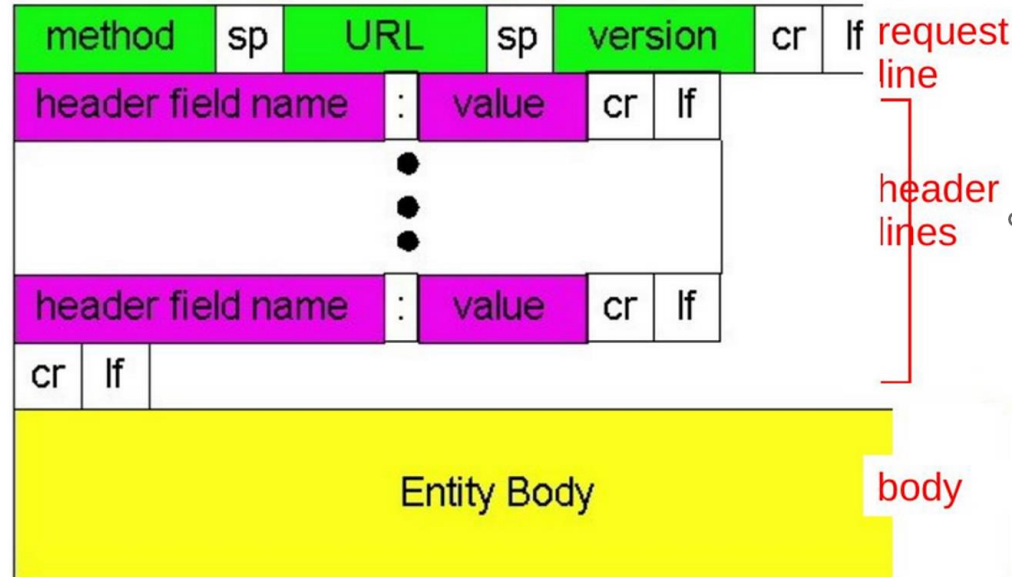
All HTTP messages fall into two types: **request messages** and **response messages**.

- Request messages **request an action from a web server**.
- Response messages **carry results of a request back to a client**.
- Both request and response messages have the same basic message structure

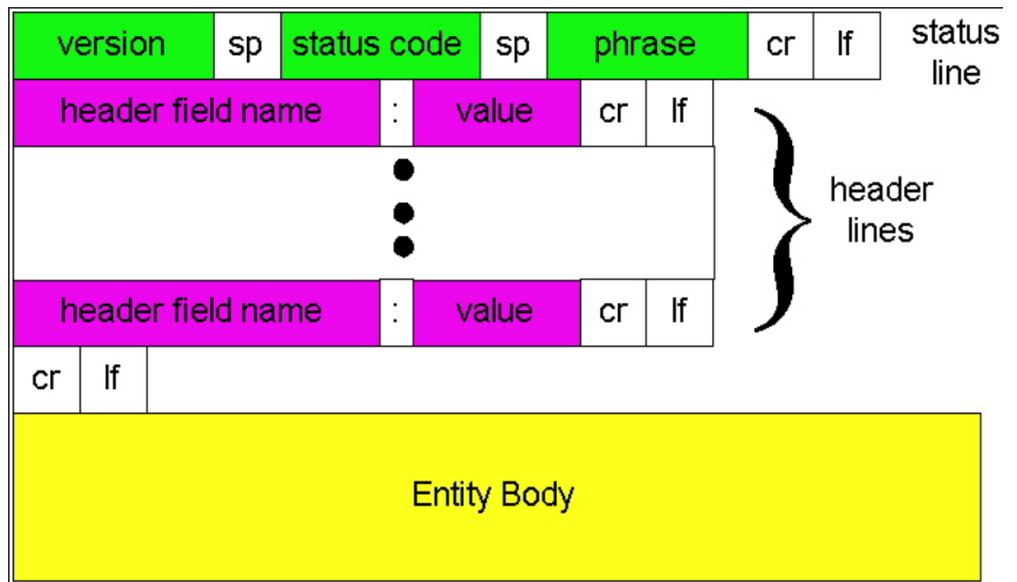


HTTP Message: General Format

- Request messages



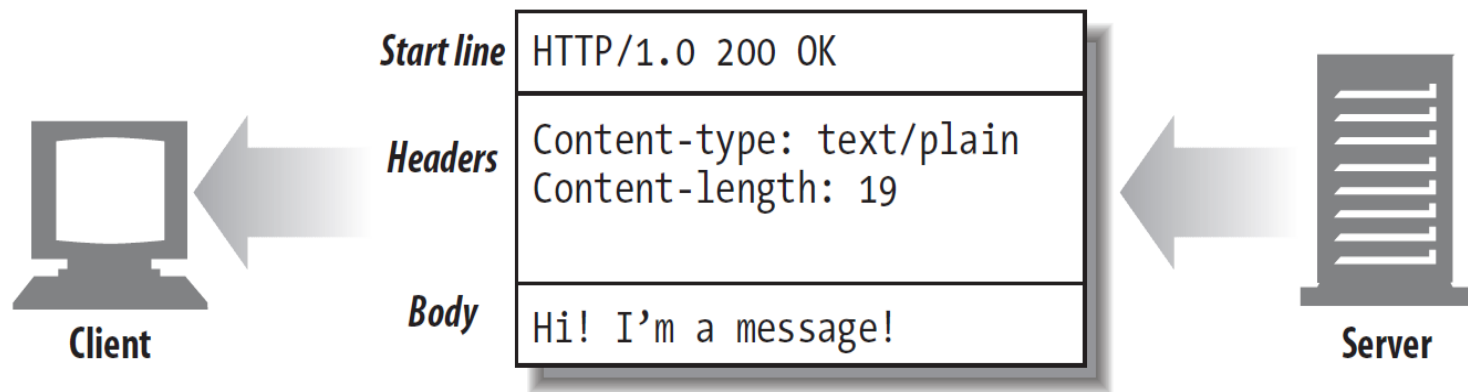
- Response messages



The Parts of an HTTP Message

- A **start line** describing the message.
- A **block of headers** containing some attributes.
- An **optional body** containing data.

Each line ends with a two-character end-of-line sequence, consisting of a carriage return (ASCII 13) and a line-feed character (ASCII 10), called **CRLF** (or `\r\n` in C like languages).



Format of a Request Message

```
<method> <request-URL> <version>  
<headers>  
<entity-body>
```

```
GET /test/hi-there.txt HTTP/1.1
```

```
Accept: text/*
```

```
Host: www.joes-hardware.com
```

Example

- **<method>** The action that the client wants the server to perform on the resource. It is a single verb such as “GET,” “HEAD,” or “POST”
- **<request-URL>** A complete URL naming the requested resource, or the **path component** of the URL.
- **<version>** The version of HTTP that the message is using. Its format looks like: HTTP/<major>.<minor>
- **<headers>** Zero or more headers, each of which is a name, then a colon (:), then a value.
- Headers are terminated by a blank line (CRLF), marking the end of headers and the beginning of the entity body.
- **<entity-body>** The entity body contains a block of arbitrary data. Not all messages contain entity bodies, so sometimes a message terminates with a bare CRLF.

Format of a Response Message

<version> <status> <reason-phrase>
<headers>
<entity-body>

HTTP/1.0 200 OK

Content-type: text/plain
Content-length: 19

Hi! I'm a message!

Example

- **<status>** A **three-digit number** describing what happened during the request. Remember, for example, the **notorious 404** (Not Found) status code.
- The first digit describes the general class of status (“success,” “error,” etc.).
- **<reason-phrase>** A **human-readable version** of the numeric status code, consisting of all the text until the end-of-line sequence.



Status Code Classes

| Overall range | Defined range | Category |
|---------------|---------------|---------------|
| 100-199 | 100-101 | Informational |
| 200-299 | 200-206 | Successful |
| 300-399 | 300-305 | Redirection |
| 400-499 | 400-415 | Client error |
| 500-599 | 500-505 | Server error |

Examples

| Status code | Reason phrase | Meaning |
|-------------|---------------|--|
| 200 | OK | Success! Any requested data is in the response body. |
| 401 | Unauthorized | You need to enter a username and password. |
| 404 | Not Found | The server cannot find a resource for the requested URL. |

The Header Lines

Each HTTP header has a simple syntax: a name, followed by a colon (:), followed by the field value, followed by a CRLF

Examples

| Header example | Description |
|--|---|
| Date: Tue, 3 Oct 1997 02:16:03 GMT | The date the server generated the response |
| Content-length: 15040 | The entity body contains 15,040 bytes of data |
| Content-type: image/gif | The entity body is a GIF image |
| Accept: image/gif, image/jpeg, text/html | The client accepts GIF and JPEG images and HTML |

Header Types

- **General headers** used by both clients and servers. For example, the Date header:

Date: Tue, 3 Oct 1974 02:16:00 GMT

- **Request headers** specific to request messages. They provide extra information to servers, such as **what type of data** the client is willing to receive. For example, to accept any media type we use:

Accept: */*

- **Response headers** provide information to the client. For example, to tell the client that it is talking to a Version 1.0 Tiki-Hut server:

Server: Tiki-Hut/1.0

- **Entity headers** refer to headers that deal with the entity body. For example, the following Content-Type header lets the application know that the data is an HTML document in the iso-latin-1 character set:

Content-Type: text/html; charset=iso-latin-1



Media Types

- In HTTP, Media types are used to indicate the type of content being transferred in the body of an HTTP message. In the HTTP protocol, when a client makes a request to a server, the server responds with an HTTP message that contains a response body, which can be any type of data, such as HTML documents, images, audio, video, and so on. **Media types are used to identify these different types of data.**
- The use of media types in HTTP provides a standardized way to represent data types, which allows different clients and servers to interact with each other. This standardization ensures that **data is correctly parsed and processed, ensuring the reliability and interoperability of data transfer.**
- Additionally, media types can **guide web browsers in displaying different content types**, such as rendering HTML files as web pages and PDF files as documents.

- In HTTP, each response message includes a "**Content-Type**" header field, which specifies the media type of the response message body.
- The Content-Type field value is made up of a **MIME (Multipurpose Internet Mail Extensions)** type and an optional MIME subtype, such as "text/html" or "image/jpeg". The client uses the Content-Type field value to determine how to handle the data in the response message body.
- MIME was **originally designed** for moving messages between different **electronic mail** systems. MIME worked very well, so HTTP adopted it to describe and label its own multimedia content.





MIME types structure

- Each MIME media type consists of a **primary type**, a **subtype**, and a list of optional parameters.
- The type and subtype are separated by a **slash**, and the optional parameters begin with a **semicolon**.
- The primary type can be a predefined type, an IETF-defined (Internet Engineering Task Force) extension token, or an experimental token (beginning with “x-”). Here are a few examples:

| Examples | Type | Description |
|----------|-------------|---|
| | application | Application-specific content format (discrete type) |
| | audio | Audio format (discrete type) |
| | chemical | Chemical data set (discrete IETF extension type) |
| | image | Image format (discrete type) |
| | message | Message format (composite type) |
| | model | 3-D model format (discrete IETF extension type) |
| | multipart | Collection of multiple objects (composite type) |
| | text | Text format (discrete type) |
| | video | Video movie format (discrete type) |

MIME type examples:

- An HTML-formatted text document would be text/html.
- A plain ASCII text document would be text/plain.
- A JPEG version of an image would be image/jpeg.
- An Apple QuickTime movie would be video/quicktime.
- A Microsoft PowerPoint presentation would be application/vnd.ms-powerpoint.
- A GIF-format image would be image/gif.



MIME Type IANA Registration

- MIME types should be registered with **IANA** (Internet Assigned Numbers Authority)
- Registration is simple and open to all.
- MIME type tokens are split into **four classes**, called “**registration trees**,” each with its own registration rules.

| Registration tree | Example | Description |
|----------------------------|---|--|
| IETF | text/html (HTML text) | <p>The IETF tree is intended for types that are of general significance to the Internet community. New IETF tree media types require approval by the Internet Engineering Steering Group (IESG) and an accompanying standards-track RFC.</p> <p>IETF tree types have no periods (.) in tokens.</p> |
| Vendor (vnd.) | image/vnd.fpx (Kodak FlashPix image) | <p>The vendor tree is intended for media types used by commercially available products. Public review of new vendor types is encouraged but not required.</p> <p>Vendor tree types begin with “vnd.”.</p> |
| Personal/Vanity (prs.) | image/prs.btif (internal check-management format used by Nations Bank) | <p>Private, personal, or vanity media types can be registered in the personal tree. These media types will not be distributed commercially.</p> <p>Personal tree types begin with “prs.”.</p> |
| Experimental (x- or x.) | application/x-tar (Unix tar archive) | <p>The experimental tree is for unregistered or experimental media types. Because it's relatively simple to register a new vendor or personal media type, software should not be distributed widely using x- types.</p> <p>Experimental tree types begin with “x.” or “x-”.</p> |

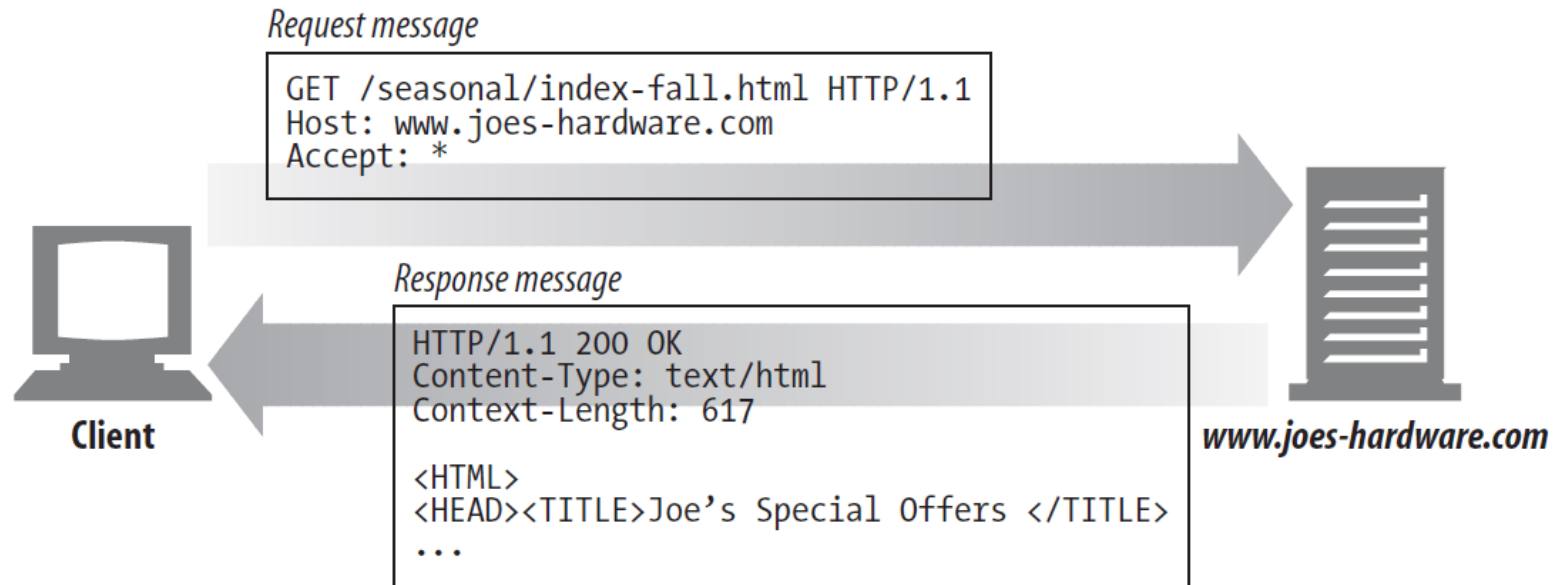
Common HTTP Methods (verbs)

| Method | Description | Message body? |
|---------|--|---------------|
| GET | Get a document from the server. | No |
| HEAD | Get just the headers for a document from the server. | No |
| POST | Send data to the server for processing. | Yes |
| PUT | Store the body of the request on the server. | Yes |
| TRACE | Trace the message through proxy servers to the server. | No |
| OPTIONS | Determine what methods can operate on a server. | No |
| DELETE | Remove a document from the server. | No |

Note that not all methods are implemented by every server. To be compliant with HTTP Version 1.1, a **server need implement only the GET and HEAD** methods for its resources.

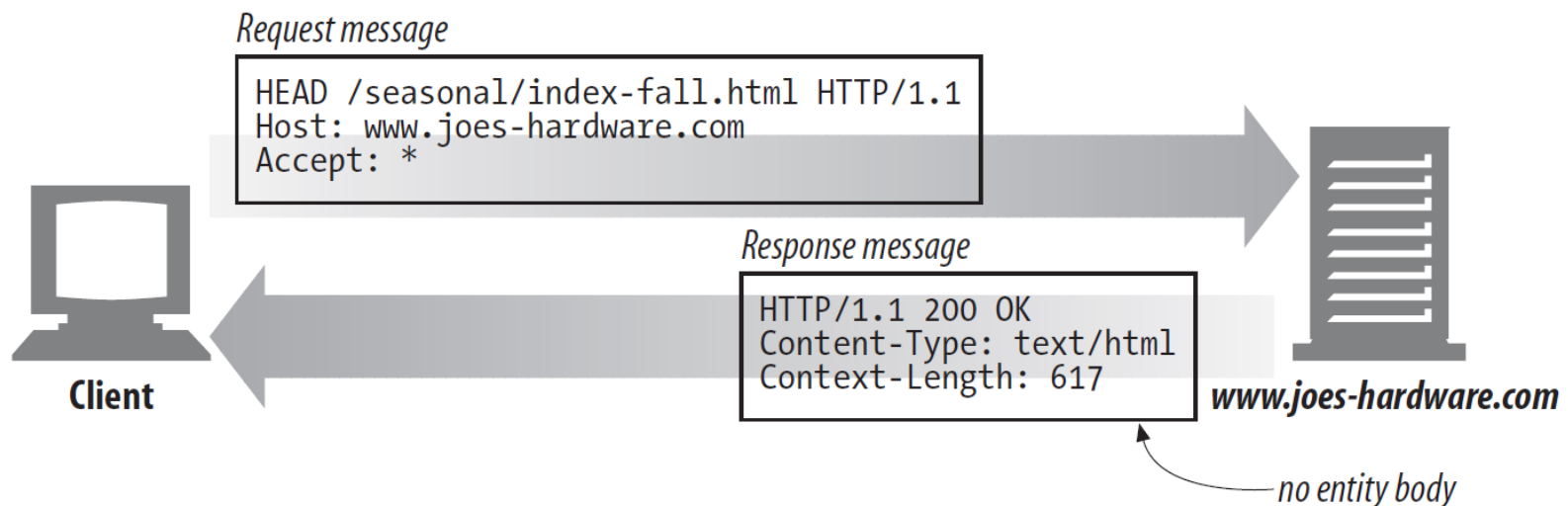
GET

GET is the most common method. It usually is used to ask a server to send a resource.



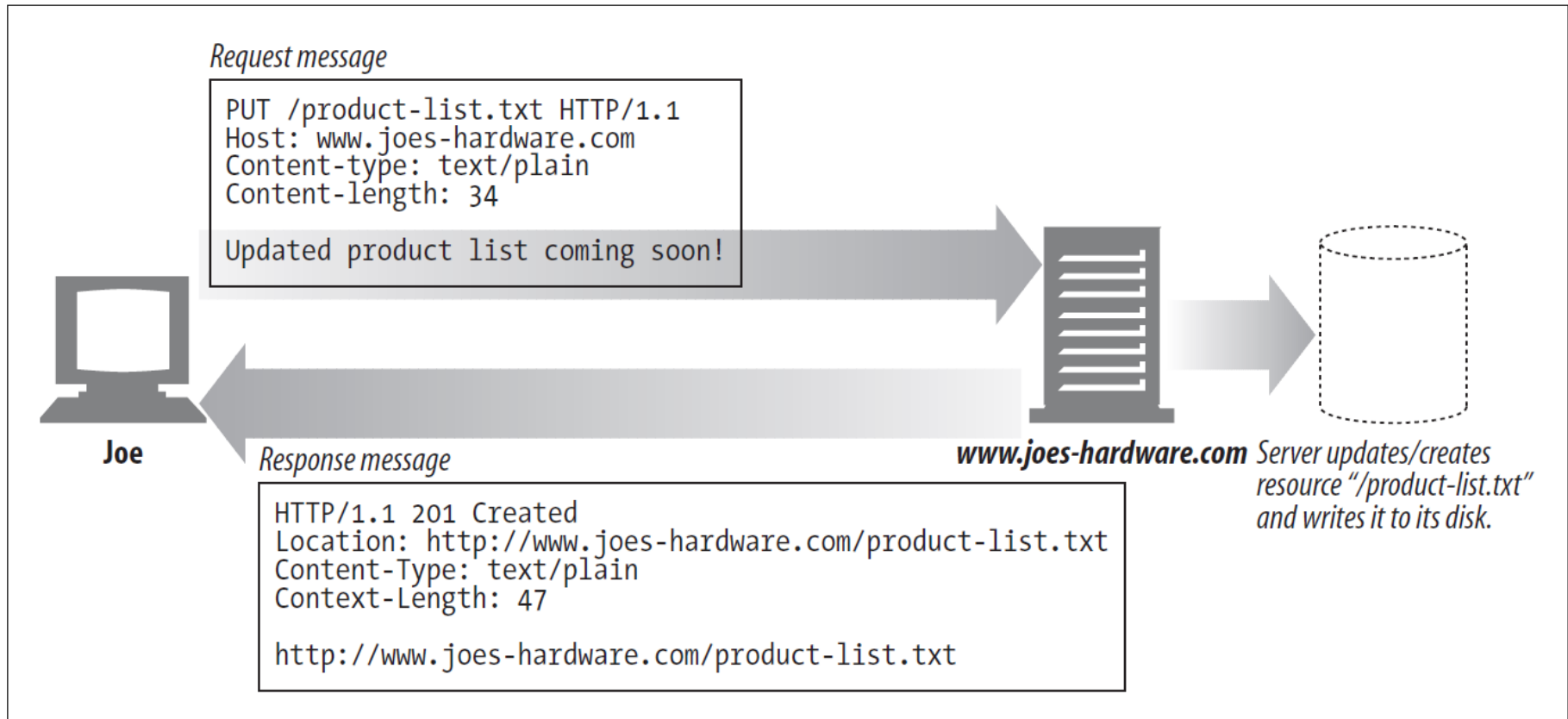
HEAD

- The HEAD method behaves exactly like the GET method, but the server returns only the headers in the response. **No entity body is ever returned.**
- This allows a client to inspect the headers for a resource without having to actually get the resource.
- Using HEAD, you can:
 - Find out about a resource (e.g., determine its type) without getting it.
 - See if an object exists, by looking at the status code of the response.
 - Test if the resource has been modified, by looking at the headers (Last-Modified).



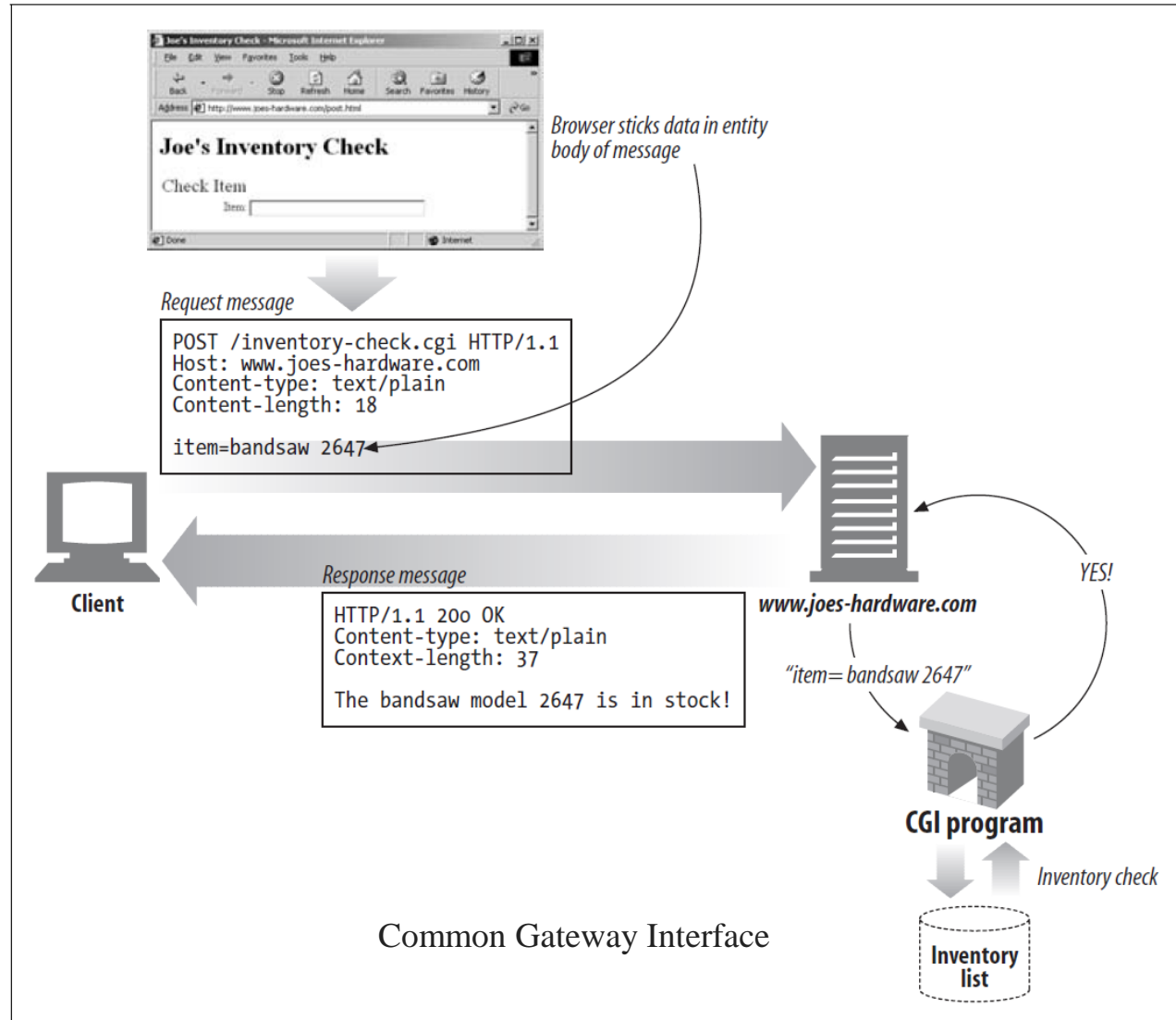
PUT

- The PUT method writes documents to a server, in the inverse of the way that GET reads documents from a server.
- Some publishing systems let you create web pages and install them directly on a web server using PUT



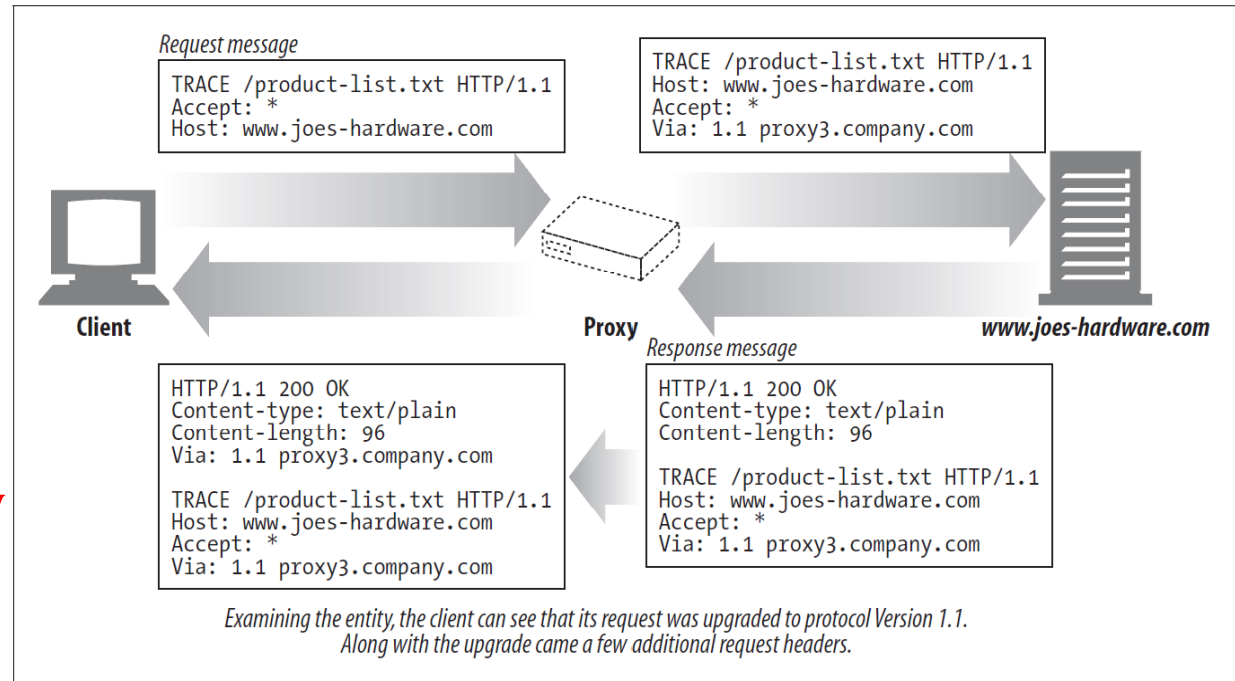
POST

- The POST method **sends input data to the server**. For example, it is often used to support HTML forms. The data from a filled-in form is sent to the server, which then processes it.



TRACE

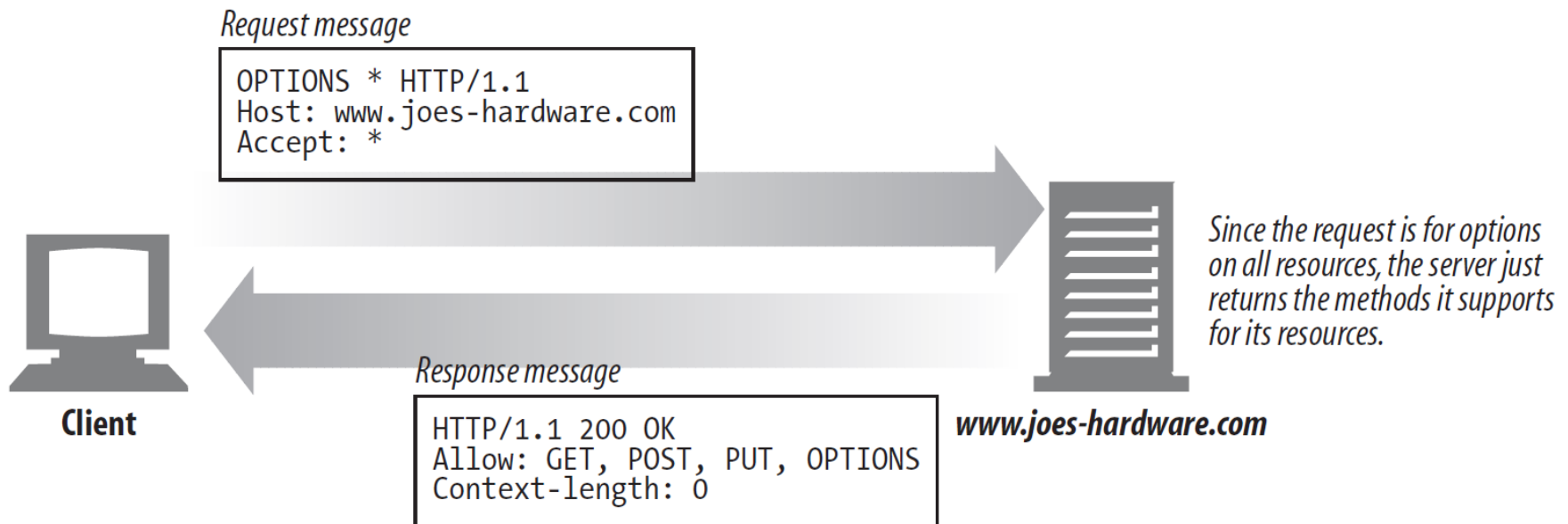
- When a client makes a request, that **request may have to travel through firewalls, proxies, gateways, or other applications.**
- Each of these has the **opportunity to modify the original HTTP request.**



- The TRACE method **allows clients to see how its request looks when it finally makes it to the server.**
- A TRACE request initiates a **“loopback” diagnostic** at the destination server.
- The server at the final leg of the trip bounces back a TRACE response, **with the virgin request message** it received in the body of its response.

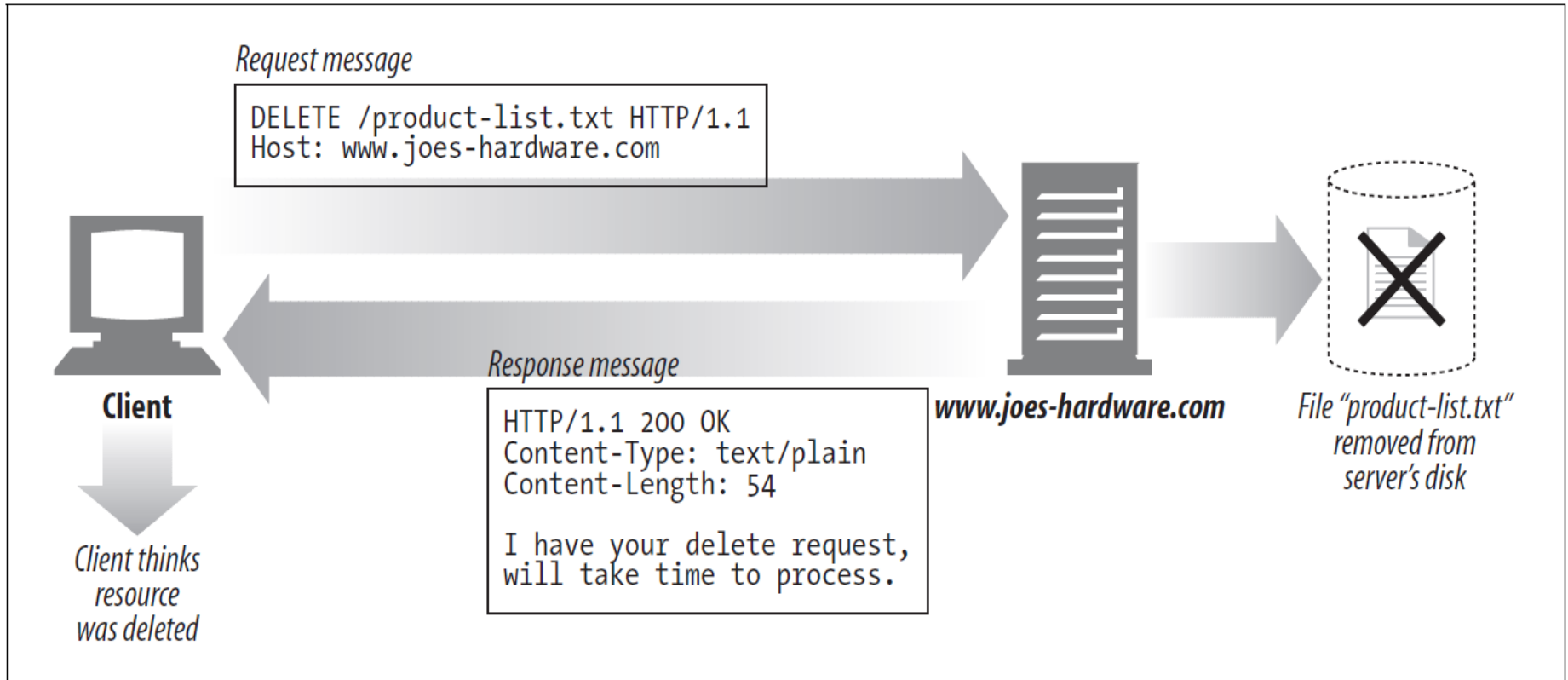
OPTIONS

- The OPTIONS method asks the server to tell us about the **various supported capabilities of the web server**.
- You can ask a server about **what methods it supports in general or for particular resources** (Some servers may support particular operations only on particular kinds of objects).
- This provides a means for client applications to **determine how best to access various resources** without actually having to access them.



DELETE

- The DELETE method asks the server to **delete the resources** specified by the request URL.
- However, the client application is **not guaranteed that the delete is carried out**.
- The HTTP specification allows the server to override the request **without telling the client**.



Extension Methods

- HTTP was designed to be **field-extensible**, so new features wouldn't cause older software to fail.
- Extension methods are **methods that are not defined in the HTTP/1.1** specification.
- They provide developers with a means of **extending the capabilities of the HTTP** services their servers implement on the resources that the servers manage.

Example: the WebDAV HTTP extension that support publishing of web content to web servers over HTTP.

| Method | Description |
|--------|---|
| LOCK | Allows a user to “lock” a resource—for example, you could lock a resource while you are editing it to prevent others from editing it at the same time |
| MKCOL | Allows a user to create a resource |
| COPY | Facilitates copying resources on a server |
| MOVE | Moves a resource on a server |

That's it Folks



Further Reading

HTTP: The Definitive Guide by Brian Totty and David Gourley