

Overview

Last week

- The web consists of servers and clients communicating using **HTTP requests** and **HTTP responses** (that contain **HTML** which may contain **JavaScript**)
- HTTP requests are usually **GET** or **POST** requests
 - GET: parameters in URL
 - POST: parameters in HTTP body

Today

- different **languages & encodings** of data in HTTP traffic
- two notions of **sessions** (incl. for security):
 - **TLS / HTTPS** at transport level
 - **cookies** at application level

Web Security

Languages & encodings

Languages/formats in web pages

```
<html><title>The various languages and formats used inside web pages</title>
<body>
  <h1 style="color:blue;">Sample exam question</h1> is 3 &lt; 4?
  <a href="https://duckduckgo.com/?q=how+to+encode+<+in+HTML%3F">A link with special characters</a>
  <a href="https://duckduckgo.com/?q=how+to+encode+%2F+in+a+URL%3F">And another one</a>
  <script> var x = 'a string with a single quote \' and double quote ".";
            alert(x);
  </script>
</body>
</html>
```

How many languages (or formats) do you see in the text above?

Which encoded characters do you see?

Homework exercise: If you are not sure how the browser would display this,
copy this text into an .html file and open it in a browser

Web pages contain HTML, CSS, JavaScript, and URLs

```
<html><title>The various languages and formats used inside web pages</title>
```

```
<body>
```

```
<h1 style="color:blue;">Sample exam question<h1> is 3 &lt; 4?
```

```
<a href="https://duckduckgo.com/?q=how+to+encode+<+in+HTML%3F">A link with special characters</a>
```

```
<a href="https://duckduckgo.com/?q=how+to+encode+%2F+in+a+URL%3F">And another one</a>
```

```
<script> var x = 'a string with a single quote \' and double quote "';
```

```
    alert(x);
```

```
</script>
```

```
</body>
```

```
</html>
```

- Special characters may need to be *encoded* aka *escaped* to prevent unintended effects
- Which characters have to be encoded, and how, depends on the **context**.
 - Eg < is a special character in HTML, but not in a URL
- Within a single language there can be several contexts. Eg
 - / is a special character in URLs, but not in the query string (i.e after the ?)
 - For a **JavaScript strings** inside **JavaScript** the outer quotes (' or ") determine which quotes inside the string need to be escaped.

URL encoding aka %-encoding

Replaces reserved characters that have a special meaning in URLs

`/?!*' ; : @ & = + $, # () []`

with their ASCII value in hex preceded with escape character `%`

<code>/</code>	<code>#</code>	<code>space</code>	<code>=</code>	<code>?</code>	<code>%</code>	<code>...</code>
<code>%2F</code>	<code>%23</code>	<code>%20</code> or <code>+</code>	<code>%3D</code>	<code>%3F</code>	<code>%25</code>	<code>...</code>

Try this out with eg `https://duckduckgo.com/?q=%3F`

Possible sources of confusion (and bugs or security issues?)

- Encoding space as `+` comes from older `x-www-form-urlencoded` format
- The reserved characters are different for different parts of the URL.
Eg `/` in the path of a URL must be encoded, in the query it need not be
- What happens if you URL-encode unreserved characters? eg `A` \rightarrow `%41`
- What happens if you double URL-encode? eg `%` \rightarrow `%25` \rightarrow `%2525`

HTML encoding

Replaces HTML special characters with similar looking ones

<	>	&	“
<	>	&	"

- HTML encoding and URL encoding are *very* different things, used for *very* different **contexts**
 - *Things can get confusing: what about URLs inside HTML or vv?*
- HTML also has the notion of **character encoding**: which character set is used, eg ASCII or UTF-8 (default)
- Some browser engines are sloppy or forgiving, and will let you get away with *not* encoding e.g. & as **&** in webpages
 - <http://validator.w3.org> checks if a page is correct HTML
- On top of HTML-encoding, websites may apply additional input **sanitisation** to remove or replace tags it wants to disallow in user input;
 - eg `<script>` tags are commonly stripped from user input

base64 encoding

HTTP is text-based, so all data transmitted has to be **text**
– ie. **printable, displayable characters**

Base64 encoding turns ‘raw’ **binary data - bytes** - into **text**
so that it can be transferred via HTTP

- 6 bits coded up as one of the 64 standard characters
a-z A-Z 0-9 + /
- Groups of 3 bytes (ie 24 bits) represented as 4 characters
- Padding with **=** or **==** to make sure results is multiple of 4 characters long

base64 encoding

Bits		0	1	0	0	1	1	0	1	0	1	1	0	0	0	0	1	0	0						
Base64 encoded	Sextets	19						22						4				Padding							
	Character	T						W						E				=							
	Octets	84 (0x54)						87 (0x57)						69 (0x45)				61 (0x3D)							

- groups of 6 bits coded up as one of the standard characters
a-z A-Z 0-9 + /
- So 3 bytes represented as 4 characters
- Padding with zeroes to make the input a multiple of 6 bits
- Padding with = or == to make sure results is multiple of 4 characters long

Details not that important for this course, but you may come across base64-encoded data

Encoding user content for security

If web pages contain **user-supplied content** this may need to be **encoded** or **sanitised** to prevent malicious user content from triggering unwanted effects.

```
<html><title>Mallory's Radboud Student Homepage</title>
<body>
  <h1 color ="color:red;">Welcome to Mallory's homepage</h1>
  Some text that Mallory provided.
  <a href="https://ru.nl/Mallory">My contact information</a>
  <a href="https://brightspace.ru.nl/d2l/home/427025">My favourite course</a>
  <a href="https://ru.osiris-student.nl/grades.html?uid=s123456">My grades</a>
  <script> someJavaScriptFunction(someOtherString+'Mallory'); </script>
</body>
</html>
```

Beware of confusion

encoding

- **changing the representation of data**
- **no information is lost or changed**
- **eg HTML or URL encoding**

sanitisation and validation

- **removing or changing 'problematic' data**
- **some information is lost; possibly an entire request is rejected as invalid**
- **eg removing `<script>` tags in input or rejecting incorrect date 31/2/2024**

Sometimes both operations are combined; this combination is then often called sanitisation.

Usually sanitisation removes or changes problematic data whereas validation simply rejects the entire input.

More synonyms exist: eg encoding can also be called escaping

Exercise to hand in this week

- Figure out how Brightspace encodes & sanitises user input in Discussion Forums
 - **client-side** in the browser and/or **server-side**
 - for **header** and **body** of forum posts

Web Security

Authentication & Session Management

Security shortcomings of internet

“No security built into the internet”

But what does that mean?

- *No way of knowing who you are communicating with, apart from an IP address*
 - *ie no authentication*
- *Any party along the way (wifi router, ISP, ...) can read or modify the communication*
 - *ie no integrity & confidentiality of communication*

Adding security

Two security requirements we want to add

1. Authentication of parties involved
of a) web site by user and b) vice versa

How?

For a) TLS certificates aka X509 certificates

For b) username/password or more secure solutions, eg MFA (Multi-factor Authentication)

2. Integrity & confidentiality of communication

How?

TLS: not just by *encryption* but also by *integrity protection* with

- MACs (Message Authentication Codes)
- digital signatures

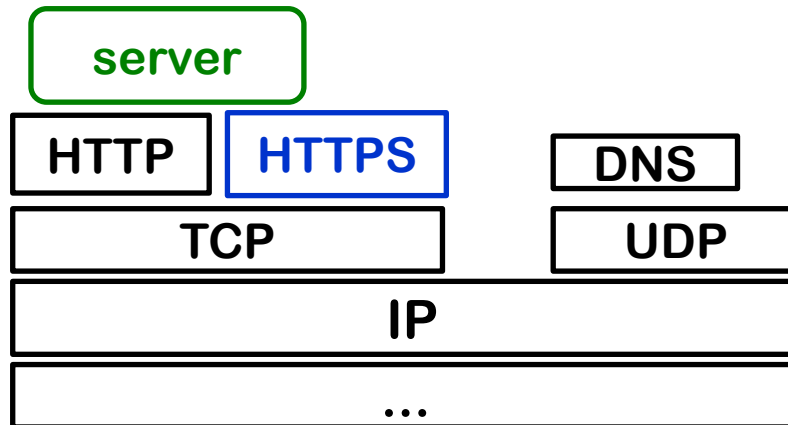
Today: two notions of sessions

1. HTTPS at the network layer

- by TLS, on top of TCP or inside QUIC
- includes authentication of the server

2. Session management at the application layer

- by web application using sessions IDs and/or cookies
- includes authentication of the user



HTTPS

Attacker models for the internet & the web



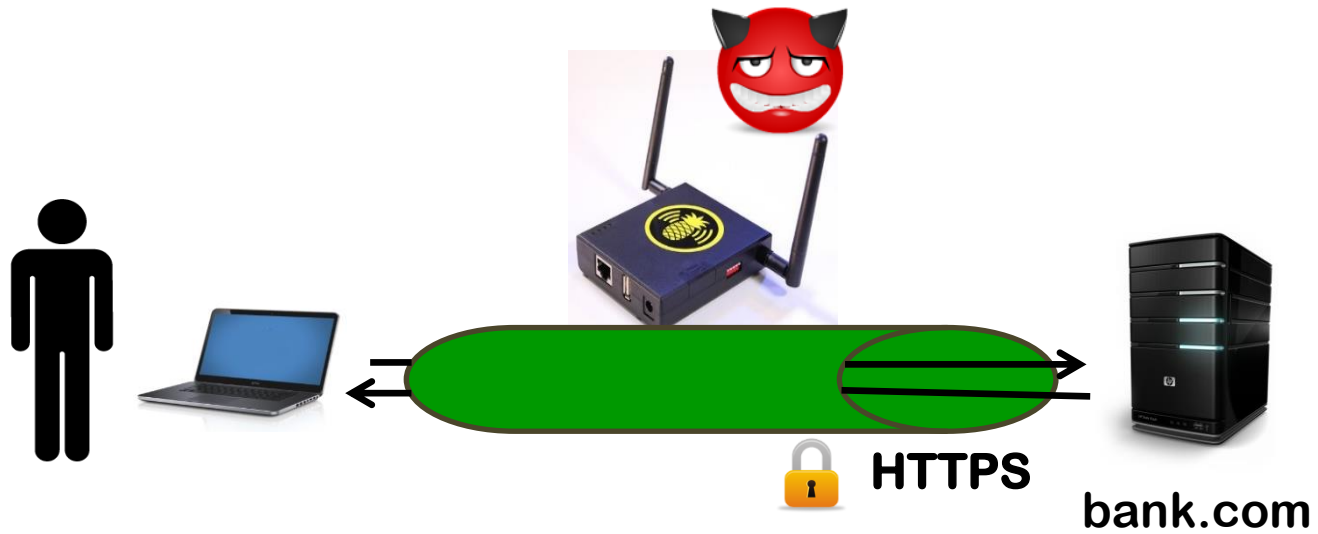
Eg a malicious or compromised **ISP, router, or WiFi access point**

WiFi security (eg WPA2) should prevent attackers eavesdropping on Wifi traffic


c) **Malicious or vulnerable end points (browser or server)**

A malicious server (eg fakebank.com) can act as MitM by relaying traffic to real website bank.com

(partial) security solution: TLS



TLS

1. Server sends **X509 server certificate** to client
 - Signed by a **Certificate Authority (CA)** or self-signed
 - Browsers come pre-configured with list of trusted CAs
2. Client checks that certificate has not been revoked
 - **by requesting Certificate Revocation List (CRL) from CA**
3. Client authenticates the server, with a challenge-response protocol
4. Client and server then agree a **session key**
5. Subsequent HTTP traffic in **a secure tunnel** 

Does a self-signed certificate provide any security guarantee?

Yes, because at least clients knows they keep talking to the same server

Does it prevent Man-in-the-Middle attack?

Not if MitM intercepts the first communication,
but it does prevent MitM breaking a session after it has been started

TLS – crypto details; not important for this course

1. Server sends **X509 server certificate** to client
*includes server's **public key PK***
2. Client checks that certificate has not been revoked
3. Client authenticates the server, with a challenge-response protocol
*Client sends **nonce n** encrypted with public key PK, and checks if server response includes n which proves knowledge of corresponding private key*
4. Client and server then agree a **session key**
Typically an AES key
5. Subsequent HTTP traffic in a **secure tunnel**
Traffic encrypted and MACed with session key
 - *encryption for confidentiality, MACing for integrity**Periodically the session key is refreshed*

HTTPS: HTTP over TLS

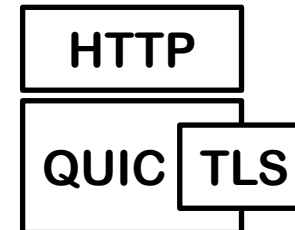
HTTP

TLS

Security guarantees :

- **Confidentiality & integrity of the session**
 - Attacker on the network can still see *that* two IP addresses communicate (an example of **meta-data**), but not *what*
 - All HTTP content, incl. headers & URL parameters are protected inside TLS tunnel
 - Attacker cannot change any traffic or replay it
- **Server authentication**, using certificates
- Possibly, but uncommon: **Client authentication** with a client certificate
 - Usually servers use another means to authenticate clients: often **passwords** 😞

The same holds if TLS is used as part of QUIC



Aside: name confusion TLS vs SSL

TLS (Transport Layer Security) used to be called **SSL** (Secure Sockets Layer)

- TLS version 1.0 is SSL version 3.1
- Latest TLS version is 1.3

X509 certificates still sometimes called **SSL certificates** and a well-known open source TLS implementation is **OpenSSL**

We'll come back to TLS later in this lecture to discuss its limitations.

Sessions (at application level)

Functional shortcoming of HTTP

HTTP is **stateless** and has **no notion of session**, ie

- No state is recorded about history of previous requests
- (Hence) no notion of a sequence of requests belonging together in one conversation between client and server

This is very clumsy for interaction between a client and server

- *Has this user logged in?*
- *Did the user select English or Dutch as language for the site?*
- *Has the user put items in their online shopping basket?*
- *Did the user already agree to our privacy policy?*

Why can't we use IP address for this?

- **Different clients may share the same IP address**
Eg different browsers & apps on the same device,
different users on lilo.science.ru.nl,
or different users on a local wifi network (esp for IPv4)
- **Multiple web applications can share the same IP address**
 - especially web applications hosted in the cloud
- **Clients and servers can change IP address**
 - eg. clients on mobile devices, when switching from mobile network to WiFi or v.v.
 - also: web applications hosted in the cloud, if they are migrated to other server

Session & session data

There is usually **session data** associated with a session that needs to be remembered.

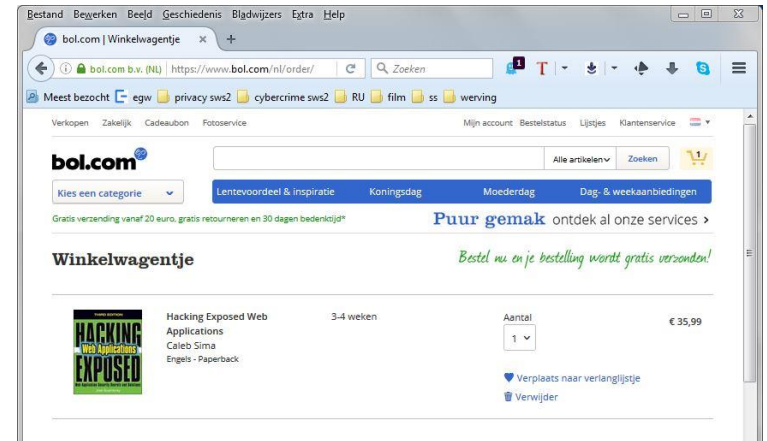
- Eg: content of online shopping basket

Ways to keep track of such data:

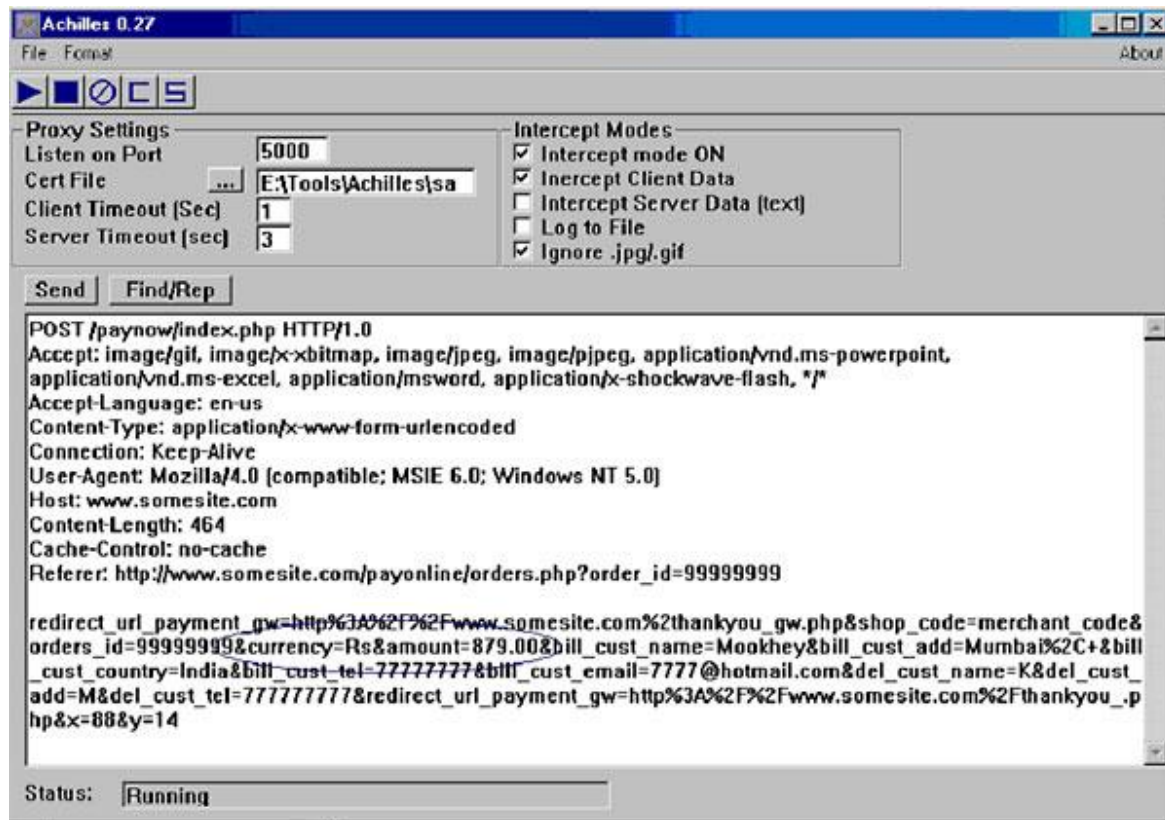
1. **send it back & forth between server and browser** with each request and response eg using hidden parameters
2. **record it at the client side** using HTML local storage
3. **record it at the server side** and just send back & forth a unique identifier

Pros & Cons ?

- Con 3: server has to record lots of info for many sessions
- Con of 1 & 2: client could mess with this data



Things that can go wrong with session data



Classic security flaw: the price is recorded in a hidden form field, as shown in the proxy output above.

The client can change this...

Misplaced trust in the client

For data for which integrity is important (eg prices)

the server should never trust the client

to provide this data or to return this data unaltered

Instead, the server should

- store such data server-side

or

- add a cryptographic integrity check
 - eg using a MAC (Message Authentication Code) or Digital Signature
 - Such a check should also include a **time stamp** or some **session id** that frequently changes, to avoid replay or roll-back attacks.

Session data for authentication

Authentication often involves a notion of session, and then goes in two steps

1. **Actual authentication**, say with a username/password plus the response of an MFA token
2. **Creating a session**, with **session identifier** aka **session token** for fast & easy (re)authentication without repeating step 1



Most web applications use **session cookies** for this purpose

- Such cookies provide **identity** and **proof that this identity has verified (aka authenticated)**
- These cookies are just as valuable for attacker as original credentials used to authenticate, eg username/password *plus* MFA response

Sessions managed by the web application

Typical steps

1. Web application creates & manages sessions
 - **Session data** is stored at server and associated with a unique **session ID**
2. Client is informed of session ID
 - and client attaches session ID to subsequent requests so server knows about previous requests

Web application frameworks usually provide built-in support for session management, but web application developers can implement their own

- NB it is better to use existing solutions than inventing your own
- Still, don't underestimate the complexity of using these correctly

Solution 1: session ID in URL

Web page returned by the server contains links with **session ID as extra parameter**

```
<html>
```

```
Example web page with session IDs in the URL.
```

```
The user can now click
```

```
<a href="http://demo.net/nextpage.php?sid=1234">here</a>
```

```
or
```

```
<a href="http://demo.net/anotherpage.php?sid=1234">here</a>
```

```
passing on its session id back to the server
```

```
wherever he goes next.
```

```
</html>
```

Hence: every user gets their own unique copy of a web page.

Solution 2: session ID in hidden parameter

```
<htm>
```

The form below uses a hidden field

```
<form method="POST" action= "http://ru.nl/register.php">  
  Email: <input type="text" name="Your email address">  
  <input type="hidden" name="sid" value="s1234">  
  <input type="submit" value="Click here to submit">  
</form>
```

Hidden means hidden from the user by browser,
not hidden from a proxy like ZAP.

A hidden form field could also be used to track user preferences, eg

```
<input type="hidden" name="Language" value="Dutch">
```


Session ID in URL vs hidden parameter

Can you think of a downside of a session ID in the URL?

If you give a link with your session ID to someone else, then that person might continue with your session!

Also, bookmarking a URL incl. the session ID does not (or should not) make sense, as the next time you use the bookmark you should start a different session

Solution 3: sessionID in a cookie

Standard solution built into HTTP and browser

- Cookie is piece of information that is **set by the server** and **stored by the browser**
 - namely when HTTP response includes `Set-Cookie` field in header
 - It belongs to some domain, eg `www.test.com`
 - It includes **expiry date**, **domain name**, optional **path**, optional **flags**
 - eg `secure`, `HttpOnly`, and `SameSite` flags
- Cookie is **automatically** included in any HTTP request by the browser, for any request to that domain
 - in the `Cookie` field of HTTP request
- Cookie can include any type of information
 - sensitive information, such as **session ID**
 - less sensitive information, such as language preferences

Example cookie traffic

- **Setting** a cookie set with an HTTP response

```
HTTP/1.0 200 OK
Content-type text/html
Set-Cookie: language=Dutch
Set-Cookie: sessionID=123; Expires=Tue, 26 Apr 2021 11:30:00 GMT
...
```

- **Sending** a cookie in an HTTP request

```
GET someurl.html HTTP/1.0 200 OK
Host: example.com
Cookie: language=Dutch, sessionID=123
```

Different types of cookies

- **non-persistent** aka **session cookies**
 - only stored while current browser session lasts
- **persistent cookies**
 - preserved between browser sessions
 - useful for maintaining login and user preferences across sessions
 - bad for privacy

Domains, subdomain, and top level domains

The domain in a cookie can be a **subdomain** of a website (eg `cs.ru.nl` is a subdomain of `ru.nl`) which raises questions, such as

Are cookies for `cs.ru.nl` sent with requests to `ru.nl`? Or v.v.?

Can `ru.nl` set a cookie for `cs.ru.nl` ?

Complex rules restrict cookie access across (sub)domains [RFC 6265]

Overall rationale: **subdomains need not trust their superdomain**

- Subdomains can *access* cookie for domain, but not vice versa
- Subdomains can *set* cookie for direct superdomain, but not vv
- With the `HostOnly` flag, cookies can further restrict access

For details, check [RFC6265] and hope browsers do not still implement outdated parts of [RFC 2109] or [RFC 2965].

For **top level domains**, eg `.nl`, there are additional rules, to prevent say `ru.nl` from setting a cookie for `.nl`

But does this work as intended for countries using 3 level domain names?
Eg for `somecompany.co.uk`, where `co.uk` is not a top level domain

Different ways to provide session ID

1. Encoding it in the URL

Downsides: 1) stored in logs (eg browser history), 2) can be cached & bookmarked, 3) visible in the browser location bar.

2. Hidden form field

Better: won't appear in URLs, so cannot be bookmarked, and less likely to be logged

3. Cookies

Best choice: automatically handled by browser; easier & more flexible.

But such automation has downsides, as we'll see: **CSRF**

Now: attacking this!

Some session attacks

Aim of attacker: **get the session ID**

- This can be session cookie, or other form of session ID
- If the victim is logged in, this is just as good as stealing his username and password!
- **Stealing the session ID**
 - eg sniffing network traffic
(or XSS attacks discussed later)
- **Session Prediction**
 - try to guess a session ID
- **Brute Force**
 - try many guesses for the session ID, until you get lucky
- **Session Fixation**
 - make the victim use a known session ID

Sessions & authentication

*How can a web site mitigate the risk of **session token** aka **session ID** being stolen and abused?*

- by having a time-out to terminate inactive sessions
- by having a prominent LogOut button on every webpage
- by tracking meta-information about the user (eg location, browser, operation system) to spot suspicious use

Session ID prediction attack

Suppose you can check your grades in blackboard on page
`brightspace.ru.nl/grades.php?s=s776823`

Is this a security problem?

If s776823 is your student number and also the session id
(in the URL in this case) then it is!

Attacker could try other student IDs or – better still – the
university employee number of a teacher.

Session fixation attack

If the sessionID is in the URL, an attacker can

1. start a session with say bank.com and obtain a session ID;
2. craft a link with that session ID and gets victims to click it, by
 - a) emailing victims with that link in the email; *or*
 - b) luring victims to a webpage with that link
3. The victim now goes to the website using a known session ID;
4. If victim logs in, and *session ID is not changed*, then attacker can join the session & abuse the user's rights!

Therefore: web server should **change session ID on login actions**

If the session-ID is a hidden form field, it does not end up in URL, attacker cannot email a link, but option 2b) is still possible, with a POST request.

Variant: attacker has already logged in, so victim joins the attacker's session and may enter confidential data (eg credit card number) for attacker's account.

[aka **Login CSRF**]

Making attacks on sessions harder

- Use long enough, random session IDs – ie with enough entropy
 - prevents session prediction and brute forcing
- Change session ID after any change in privilege level
 - eg after logging in
 - prevents session fixations
- Expire sessions
 - eg by setting expiration on cookies
 - reduces the attack surface in time
- Use HTTPS
 - for *all* requests & responses that include session ID, not just the login
 - prevents networking sniffing of session ID
- Let clients re-authenticate before important actions
 - reduces the value of any stolen session ID

CSRF – only if session IDs are stored in cookies

Interaction of two features:

- 1) any web page can link to another other web page
- 2) browser automatically attaches the cookies of A.com to any requests to A.com

can be abused:

if attacker creates webpage or HTML email with links to bank.com, the browser will automatically attach the bank's cookies with the correct value to authenticate these requests (assuming victim is logged on)

Why us this a security problem?

Attacker can trick user in performing action at bank.com

Do session ID in URL or hidden form field also come with this risk?

No, as the attacker would need to guess the value of the session ID; the browser will not supply the right value for the attacker.

Standard solution to prevent CSRF

Use two special numbers to identify a session

1. a fixed session ID stored in a cookie
2. a changing **CSRF token**, as URL parameter or hidden form field, that changes to a new random value for each request

For any malicious cross-site requests, say from mafia.com to bank.com, the browser will attach the right session ID cookie, but these requests will not have the right CSRF token.

CSRF token aka **anti-CSRF token** aka **tokenization**

Exercises for this week

A. Checking input sanitisation in Brightspace

How is input encoded & sanitised in Discussion Forums, at the client side and/or at the server side?

B. Checking security settings for some sites where you have a login: 1) whether it support HTTP(S), 2) which cookie flags it uses for session, and 3) whether it supports HSTS and Certificate Transparency.

C. One more WebGoat lesson

Authentication Flaws - Authentication Bypasses

**NB A & B to be handed in (in pairs) via Brightspace
Deadline Thursday Feb 8, 23:59**