# Captchas & The Net



"Some speak of an Armageddon; A time when humans will build machines they neither understand nor control.

To myself I whisper, 'We already have.'"
– Taylor Swift

**And We Call It "Machine Learning"**

# Bug Of The Day…

- Yet Another Buffer overflow…
  - You think we'd be bored of them by now
- But…
  - The operating system and device drivers are special…
- They need very low level access
  - As they are working in a world where everything *is* just a big pile of bits!
- Perhaps you could use rust…
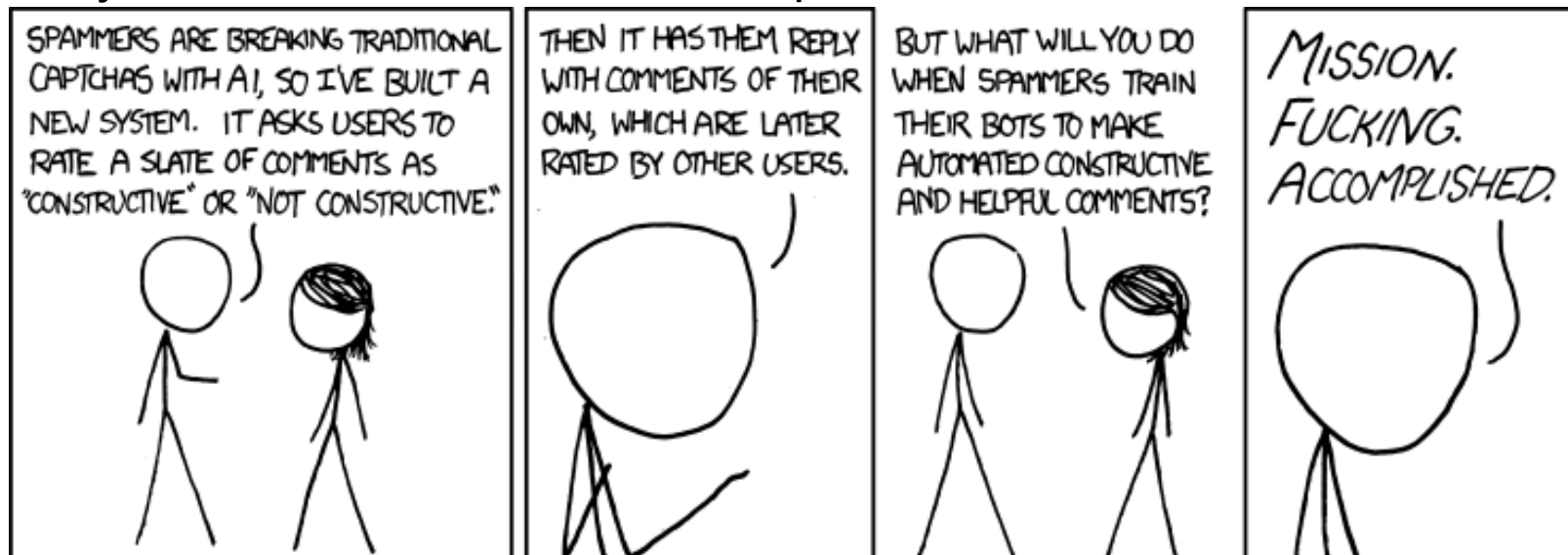  - But you would need to rewrite a huge amount of code

**ars** TECHNICA

*NOTICE OF ABSENCE —*

## Unpatched Linux bug may open devices to serious attacks over Wi-Fi

Buffer overflow can be triggered in Realtek Wi-Fi chips, no user interaction needed.

DAN GOODIN - 10/17/2019, 2:35 PM

2

# The Problem:
# Automation…

- You host some website…

- It is intended for **_human_** usage

  - One person, one mouse, one clickstream of behavior…

- But you want to lock out **_robot_** usage

- Why?

  - Selling something

  - Offering something for free

  - Dealing with load from an attack

- Enter the CAPTCHA:
  A way to go "Is this a human?"

# CAPTCHAs:
# How Lazy Cryptographers Do AI

- The whole point of CAPCHAs is **not** just to solve "is this human"...
  - But leverage bad guys to force them to solve hard problems
  - Primarily focused on machine vision problems



4

Visual code | Audio code                                         Help

yz4ccruAAJ

Type the code shown  [                    ]          ⟳ Try a new code

By clicking the "Create My Account" button below, I certify that I have read and agree to the Yahoo! Terms of Service, Yahoo! Privacy Policy and Communication Terms of Service, and to receive account related communications from Yahoo! electronically. Yahoo! automatically identifies items such as words, links, people, and subjects from your Yahoo! communications services to deliver product features and relevant advertising.
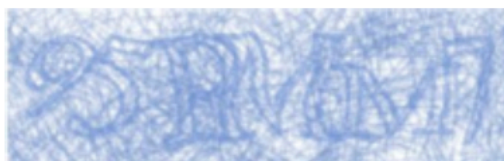
**Create My Account**

# CAPTCHAs

- *Reverse Turing Test*: present "user" a challenge that's easy for a human to solve, hard for a program to solve
- One common approach: distorted text that's difficult for character-recognition algorithms to decipher
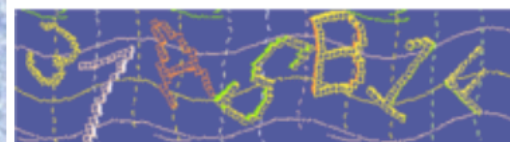


Security Check
Enter **both words** below, separated by a space.
Can't read the words below? Try different words or an audio captcha.

Text in the box:

6

(a) Aol.                          (b) mail.ru                        (c) phpBB 3.0

(d) Simple Machines Forum         (e) Yahoo!                         (f) youku

Figure 1: Examples of CAPTCHAs from various Internet properties.

Problems?

vatinkes πύργους



ReCAPTCHA™

stop spam.
read books.

**Verify Your Registration**

* Enter the code shown: [                    ]   More info

This helps prevent automated registrations.

fuck CloudFlare

**Qualifying question**

Just to prove you are a human, please answer the
following math challenge.

Q: Calculate:

$$\frac{\partial}{\partial x}\left[4 \cdot \sin\left(7 \cdot x - \frac{\pi}{2}\right)\right]\Bigg|_{x=0}$$

A: [                    ]

*mandatory*

Note: If you do not know the answer to this question,
reload the page and you'll get another question.

**Please enter the code you see below.** what's this?
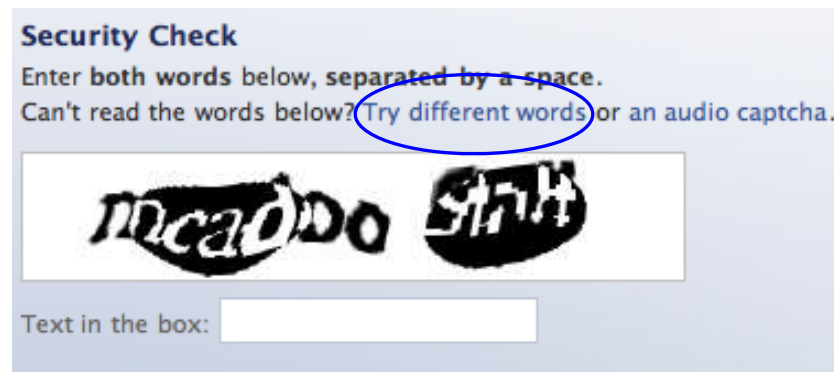
8

# Issues with CAPTCHAs

- Inevitable arms race: as solving algorithms get better, defense erodes

Figure 4: Examples of images from the hard CAPTCHA puzzles dataset.

9

# Issues with CAPTCHAs

- Inevitable arms race: as solving algorithms get better, defense erodes, or gets harder for humans

# Asirra

*Asirra is a human interactive proof that asks users to identify photos of cats and dogs. It's powered by over* **two million photos** *from our unique partnership with* Petfinder.com. *Protect your web site with Asirra — free!*



**Please click on the images that show cats:**

adopt me | adopt me | adopt me | adopt me

adopt me | adopt me | adopt me | adopt me

adopt me | adopt me | adopt me | adopt me

(Score Test)
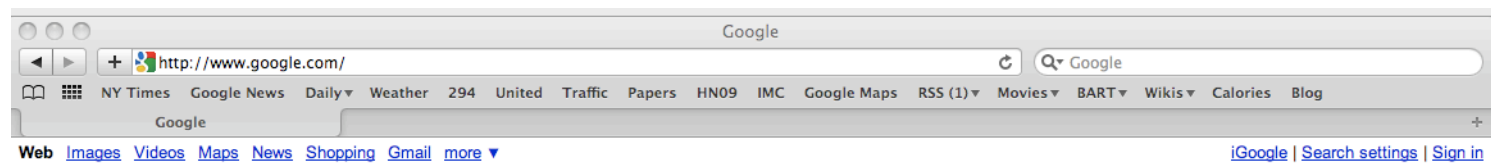
11

# Issues with CAPTCHAs

- Inevitable arms race: as solving algorithms get better, defense erodes, or gets harder for humans



- *Accessibility*: not all humans can see
- *Granularity*: not all bots are bad (e.g., crawlers)

12

# Issues with CAPTCHAs, con't

- Deepest problem: CAPTCHAs are inherently vulnerable to *outsourcing* attacks
  - Attacker gets real humans to solve them

13

## CAPTCHA solving service

✓ **Cheapest price on the market**
   Starting from 0.5USD per 1000 images, depending on your daily upload volume

✓ **Pay as you go**
   Pay-per-captcha payment basis. Minimum refill is 1 USD, no recurring charges

✓ **99.99% uptime since 2007**
   Vast amount of workers and premium infrastructure allows us to provide highly reliable 24/7/365 service

✓ **Solving Google Recaptcha since 2016**
   You may fully rely on our stable solution and forget about browser emulation

**Create Account**

🔒 Customers Area

---

API 📎

# Documentation in English

Created by Administrator
Last updated Oct 17, 2018

## About version 2.0

API version 2 works on address https://api.anti-captcha.com/ and it works only via HTTP methods, data format is JSON.

**To solve a captcha, you need:**

1. Create captcha task via via createTask method which will return task ID.
2. Wait a few seconds to let system assign captcha to an employee and retrieve resu
3. Request captcha solution with getTaskResult method. If captchas is not solved ye step #2.

| createTask | | getTaskResult |
|---|---|---|
| - NoCaptchaTask<br>- NoCaptchaTaskProxyless<br>- ImageToTextTask | ···wait 5-10s··▶ | |

16

| Language | Example | AG | BC | BY | CB | DC | IT | All |
|---|---|---|---|---|---|---|---|---|
| English | one two three | 51.1 | 37.6 | 4.76 | 40.6 | 39.0 | 62.0 | 39.2 |
| Chinese (Simp.) | 一 二 三 | 48.4 | 31.0 | 0.00 | 68.9 | 26.9 | 35.8 | 35.2 |
| Chinese (Trad.) | 一 二 三 | 52.9 | 24.4 | 0.00 | 63.8 | 30.2 | 33.0 | 34.1 |
| Spanish | uno | | | | | | | |
| Italian | uno | | | | | | | |
| Tagalog | isá da | | | | | | | |
| Portuguese | um | | | | | | | |
| Russian | один | | | | | | | |
| Tamil | ஒன்று | | | | | | | |
| Dutch | een t | | | | | | | |
| Hindi | एक | | | | | | | |
| German | eins | | | | | | | |
| Malay | satu dua tiga | 0.00 | 1.42 | 0.00 | 0.00 | 0.55 | 29.4 | 5.23 |
| Vietnamese | một hai ba | 0.46 | 2.07 | 0.00 | 0.00 | 1.74 | 18.1 | 3.72 |
| Korean | 일 이 삼 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 20.2 | 3.37 |
| Greek | ένα δύο τρία | 0.45 | 0.00 | 0.00 | 0.00 | 0.00 | 15.5 | 2.65 |
| Arabic | ثلاثة اثنين واحد | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 15.3 | 2.56 |
| Bengali | এক দুই তিন | 0.45 | 0.00 | 9.89 | 0.00 | 0.00 | 0.00 | 1.72 |
| Kannada | ಒಂದು ಎರಡು ಮೂರು | 0.91 | 0.00 | 0.00 | 0.00 | 0.55 | 6.14 | 1.26 |
| Klingon | ⌐ < ∈ | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 1.12 | 0.19 |
| Farsi | سه دو یک | 0.45 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.08 |

Table 2: Percentage of responses from the services with correct answers for the language CAPTCHAs.

Re: CAPTCHAs – Understanding CAPTCHA-Solving Services in an Economic Context

Marti Motoyama, Kirill Levchenko, Chris Kanich, Damon McCoy, Geoffrey M. Voelker and Stefan Savage
University of California, San Diego
{mmotoyam, klevchen, ckanich, dlmccoy, voelker, savage}@cs.ucsd.edu

17

# These Days:
# CAPTCHAs are ways of *training* AI systems

- Plus are all about an economic protection
  - Even the best CAPTCHA doesn't say "Is this a human or a bot"…
  - but…
  - "Is this a human or a bot willing to spend a couple pennies?"
- Acts as a hard limit on what a CAPTCHA can really protect!

TO COMPLETE YOUR REGISTRATION, PLEASE TELL US WHETHER OR NOT THIS IMAGE CONTAINS A STOP SIGN:

NO   YES

ANSWER QUICKLY—OUR SELF-DRIVING CAR IS ALMOST AT THE INTERSECTION.

SO MUCH OF "AI" IS JUST FIGURING OUT WAYS TO OFFLOAD WORK ONTO RANDOM STRANGERS.

18

# Network Security

- Why study network security?
  - Networking greatly extends our overall attack surface
    - Networking = the Internet
  - Opportunity to see *how large-scale design affects security issues*
  - Protocols a great example of *mindless agents* in action

- This lecture + next: sufficient background in networking to then explore security issues in next ~8 lectures

- Complex topic with many facets
  - We will omit concepts/details that aren't very security-relevant
  - But to no small extent we are speed running about 1/2 a dozen worth of "networking" lectures!
  - By all means, ask questions when things are unclear

19

# Protocols

- A protocol is an agreement on how to communicate

- Includes syntax and semantics
  - How a communication is specified & structured
    - Format, order messages are sent and received
  - What a communication means
    - Actions taken when transmitting, receiving, or timer expires

- E.g.: making a comment in lecture?
  1. Raise your hand.
  2. Wait to be called on.
  3. Or: wait for speaker to **pause** and vocalize
  4. If unrecognized (after timeout): vocalize w/ "excuse me"
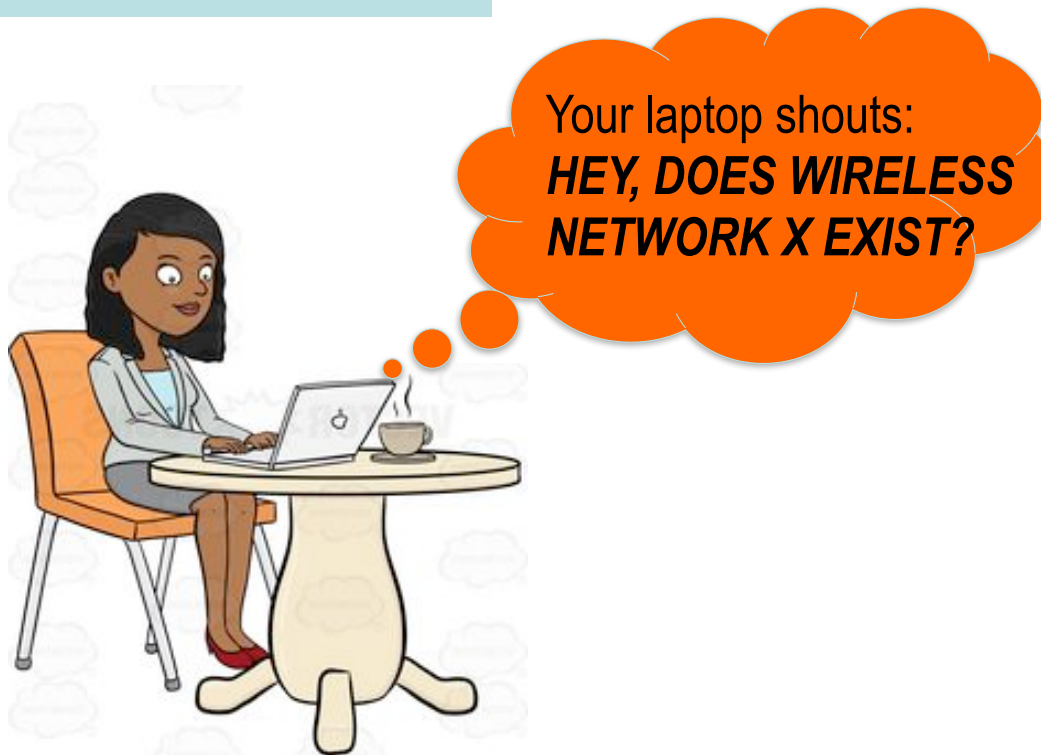
20

# So Let's Do A Google Search...

- Walk into a coffee shop

- Open a laptop

- Search google...

## Coffee Shop

**2. Configure your connection**

192.168.1.14

The configuration includes:
(1) An Internet address (**IP address**) your laptop should use; typ. 32 bits (IPv4). May also include 64b of the 128b IPv6 address
(2) The address of a "**gateway**" system to use to access *hosts* beyond the local network
(3) The address of a **DNS server** ("*resolver*") to map names like `google.com` to IP addresses

Coffee Shop

192.168.1.14

gateway

3. Find the address of google.com

Coffee Shop

3. Find the address of google.com

192.168.1.14

gateway

resolver

router

The Rest of the Internet

Coffee Shop

google.com?

192.168.1.14

gateway

resolver

router

3. Find the address of google.com

The Rest of the Internet

**3. Find the address of google.com**

Coffee Shop

192.168.1.14

google.com?

resolver

gateway

router

*(The resolver now itself uses DNS queries to other DNS servers to figure out the address associated with `google.com`.)*

*The Rest of the Internet*

Coffee Shop

192.168.1.14

3. Find the address of google.com

gateway

resolver

router

google.com's address is
172.217.6.78

*The Rest of the Internet*

Coffee Shop

4. Connect to `google.com` server

192.168.1.14

gateway

resolver

router

*The Rest of the Internet*

*Coffee Shop*

192.168.1.14

gateway

resolver

router

*The Rest o*

*the Interne*

172.217.6.78

Your laptop now *establishes a connection* with the web server at `172.217.6.78`. It uses **TCP** for this rather than UDP, to obtain reliability.

Coffee Shop

TCP SYN

192.168.1.14

gateway

resolver

router

4. Connect to `google.com` server

The Rest of the Internet

172.217.6.78

The first step of establishing the connection is to send a TCP connection request ("SYN") to the server.

Coffee Shop

*The Rest of the Internet*

192.168.1.14

gateway

resolver

router

TCP SYN ACK

172.217.6.78

If the server accepts the connection, it replies with a "SYN ACK".

Coffee Shop

TCP ACK

192.168.1.14

gateway

resolver

router

4. Connect to `google.com` server

The Rest of
the Interne

172.217.6.78

Your laptop completes the connection establishment by likewise sending an acknowledgement.

Coffee Shop

4. Connect to `google.com` server

192.168.1.14

gateway

resolver

router

The Rest of the Internet

172.217.6.78

At this point the connection is established and data can be (reliably) exchanged.

Coffee Shop

192.168.1.14

gateway

resolver

Here's a certificate that vouches for my public key, `google.com`

The Rest of the Internet

172.217.6.78

5. Establish a secure connection using **TLS** (https)

Coffee Shop

192.168.1.14

gateway

resolver

router

Here's your proof

The Rest of the Internet

172.217.6.78

5. Establish a secure connection using **TLS** (https)

Coffee Shop

GET /search?query=
who+is+outis%3F…

192.168.1.14

gateway

resolver

router

The Rest of
the Internet

172.217.6.78

6. Finally, your laptop can send along
your query!
(Using HTTP inside the *TLS channel*)

# Layering

- Internet design is strongly partitioned into layers

  - Each layer relies on services provided by next layer below …

  - … and provides services to layer above it

- Analogy:

  - Consider structure of an application you've written and the "services" each layer relies on / provides

| Code You Write |
|:---:|
| Run-Time Library |
| System Calls |
| Device Drivers |
| Voltage Levels / Magnetic Domains |

} Fully isolated from user programs

47

# Internet Layering ("Protocol Stack"/"OSI Model")

| # | Layer |
|---|-------|
| 8 | **People & Politics** |
| 7 | **Application** |
| 4 | **Transport** |
| 3 | **(Inter)Network** |
| 2 | **Link** |
| 1 | **Physical** |

Note on a point of potential confusion: these diagrams are always drawn with lower layers **below** higher layers …

But diagrams showing the layouts of packets are often the *opposite*, with the lower layers at the **top** since their headers <u>precede</u> those for higher layers

(And nobody remembers what layers 5 and 6 are for ("Session" and "Presentation) for the trivia buffs because they aren't really used)

(also, layer 8 is a "joke", but really is important)

# Packets and The Network

- Modern networks break communications up into packets
  - For our purposes, packets contain a variable amount of data up to a maximum specified by the particular network
- The sending computer breaks up the message and the receiving computer puts it back together
  - So the software doesn't actually see the packets per-se
  - Network itself is *packet switched*: sending each packet on towards its next destination
- Other properties:
  - Packets are received *correctly* or not at all in the face of *random* errors
    - The network does not enforce correctness in the face of adversarial inputs:
      They are checksums not cryptographic MACs.
  - Packets may be *unreliable* and "dropped"
    - Its up to higher-level protocols to make the connection Reliable

49

# Horizontal View of a Single Packet

**First bit transmitted**

| Link Layer Header | (Inter)Network Layer Header (IP) | Transport Layer Header | *Application Data: structure depends on the application …* |

# Vertical View of a Single Packet

**First bit transmitted**

| Link Layer Header |
|---|
| **(Inter)Network Layer Header (IP)** |
| Transport Layer Header |
| *Application Data: structure depends on the application* . . . . . . |

51

# Internet Layering ("Protocol Stack")

| 7 | Application |
| 4 | Transport |
| 3 | (Inter)Network |
| 2 | Link |
| 1 | Physical |

52

# Layer 1: Physical Layer

| | |
|---|---|
| 7 | **Application** |
| 4 | **Transport** |
| 3 | **(Inter)Network** |
| 2 | **Link** |
| 1 | **Physical** |

Encoding bits to send them over a <u>single</u> **physical link** e.g. patterns of
*voltage levels / photon intensities / RF modulation*

53

# Layer 2: Link Layer

| | |
|---|---|
| 7 | **Application** |
| 4 | **Transport** |
| 3 | **(Inter)Network** |
| 2 | **Link** |
| 1 | **Physical** |

Framing and transmission of a collection of bits into individual messages sent across a single "subnetwork" (one physical technology)

Might involve multiple *physical links* (e.g., modern Ethernet)

Often technology supports broadcast transmission (**every** "node" connected to subnet receives)

54

# Layer 3: (Inter)Network Layer *(IP)*

| | |
|---|---|
| 7 | **Application** |
| 4 | **Transport** |
| 3 | **(Inter)Network** |
| 2 | **Link** |
| 1 | **Physical** |

*Different* **for each Internet "hop"**

Bridges multiple "subnets" to provide *end-to-end* internet connectivity between nodes
  • Provides global addressing

Works across different link technologies

55

# Layer 4: Transport Layer

7    **Application**

4    **Transport**

3    **(Inter)Network**

2    **Link**

1    **Physical**

*End-to-end* communication between processes

Different services provided:
  TCP = reliable *byte stream*
  UDP = unreliable *datagrams*

(*Datagram* = single packet message)

56

# Layer 7: Application Layer

| | |
|---|---|
| 7 | **Application** |
| 4 | **Transport** |
| 3 | **(Inter)Network** |
| 2 | **Link** |
| 1 | **Physical** |

Communication of whatever you wish

Can use whatever transport(s) is convenient

Freely structured

E.g.:
  Skype, SMTP (email), HTTP (Web), Halo, BitTorrent

57

# 4.5: Some Crypto…

7 **Application**

4 **Transport**

3 **(Inter)Network**

2 **Link**

1 **Physical**

TLS cryptography
(aka the 's' in HTTPS)

Often basically used as a
"layer 4.5" transport layer to
encrypt otherwise
unencrypted network
connections

Other times crypto may be at
the application layer (e.g. ssh

# Internet Layering ("Protocol Stack")

| | |
|---|---|
| 7 | **Application** |
| 4 | **Transport** |
| 3 | **(Inter)Network** |
| 2 | **Link** |
| 1 | **Physical** |

Implemented only at hosts, not at interior routers ("dumb network")

59

# Internet Layering ("Protocol Stack")

| | |
|---|---|
| 7 | **Application** |
| 4 | **Transport** |
| 3 | **(Inter)Network** |
| 2 | **Link** |
| 1 | **Physical** |

Implemented everywhere

# Internet Layering ("Protocol Stack")

| 7 | **Application** |
|---|---|
| 4 | **Transport** |
| 3 | **(Inter)Network** |

} *~ Same* **for each Internet "hop"**

| 2 | **Link** |
|---|---|
| 1 | **Physical** |

} *Different* **for each Internet "hop"**

61

# Hop-By-Hop vs. End-to-End Layers

Host A communicates with Host D

# Hop-By-Hop vs. End-to-End Layers

Host A communicates with Host D

Host C

Host A

Host D

E.g., Ethernet

Router 1

Router 2

Router 3

E.g., Wi-Fi

Router 5

Host B

Router 6

Router 7

Host E

Router 4

*Different* **Physical & Link Layers (Layers 1 & 2)**

63

# Hop-By-Hop vs. End-to-End Layers

Host A communicates with Host D

Host C

Host A

Router 1

Router 2

Host D

Router 3

Router 5

E.g., HTTP over TCP over IP

Host B

Router 6

Router 7

Host E

Router 4

*Same* Network / Transport / Application Layers (3/4/7)
(Routers **ignore** Transport & Application layers)

64

# Layer 3: (Inter)Network Layer *(IP)*

| | |
|---|---|
| 7 | **Application** |
| 4 | **Transport** |
| 3 | **(Inter)Network** |
| 2 | **Link** |
| 1 | **Physical** |

Bridges multiple "subnets" to provide *end-to-end* internet connectivity between nodes
- Provides global addressing

Works across different link technologies

65

# IPv4 Packet Structure
# (IP version 6 is different)

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

# IP Packet Structure

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 1 | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

Specifies the length of the entire IP packet: bytes in this header plus bytes in the **Payload**

67

# IP Packet Structure

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flag | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 1 | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

Specifies how to interpret the start of the **Payload**, which is the header of a *Transport Protocol* such as **TCP** (6) or **UDP** (17)

68

# IP Packet Structure

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) |
|---|---|---|---|
| 16-bit Identification | | | 3-bit Flag~~s~~ ~~13-bit Fragment Offset~~ |
| 8-bit Time to Live (TTL) | **6** | | 1~~~~ |
| 32-bit Source IP Address | | | |
| 32-bit Destination IP Address | | | |
| Options (if any) | | | |
| Start of TCP Header | | | |

Specifies how to interpret the start of the **Payload**, which is the header of a *Transport Protocol* such as **TCP** (6) or **UDP** (17)

69

# IP Packet Structure

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

# IP Packet Header (Continued)

- Two IP addresses
  - Source IP address (32 bits in main IP version, IPv4)
  - Destination IP address (32 bits, likewise)

- Destination address
  - Unique identifier/locator for the receiving host
  - Allows each node to make forwarding decisions

- Source address
  - Unique identifier/locator for the sending host
  - Recipient can decide whether to accept packet
  - Enables recipient to send reply back to source

71

# The Basic Ethernet Packet: The near-universal Layer 2

- ## An Ethernet Packet contains:
  - A preamble to synchronize data on the wire
    - We normally ignore this when talking about Ethernet
  - 6 bytes of destination MAC address
    - In this case, MAC means media access control address, not message authentication code!
  - 6 bytes of source MAC address
  - Optional 4-byte VLAN tag
  - 2 bytes length/type field
  - 46-1500B of payload

| DST MAC | SRC MAC | VLAN | Type | PAYLOAD |
|---------|---------|------|------|---------|

# The MAC Address

- The MAC acts as a device identifier
  - The upper 3 bytes are assigned to a manufacturer
    - Can usually identify product with just the MAC address
  - The lower 3 bytes are assigned to a specific device
    - Making the MAC a de-facto serial #
- Usually written as 6 bytes in hex:
  - e.g. `13:37:ca:fe:f0:0d`
- A device ***should ignore*** all packets that aren't to itself or to the broadcast address (`ff:ff:ff:ff:ff:ff`)
  - But almost all devices can go into ***promiscuous mode***
    - This is also known as "sniffing traffic"
- A device generally should only send with its own address
  - But this is enforced with software and can be trivially bypassed when you need to write "raw packets"

73