



"Controlled" Distributed Systems

- 1 responsible organization
- Low churn
- Examples:
 - Amazon DynamoDB
 - Client/server
- "Secure environment"
- High availability
- Can be homogeneous / heterogeneous

"Fully" Decentralized Systems

- N responsible organizations
- High churn
- Examples:
 - BitTorrent
 - Blockchain
- "Hostile environment"
- Unpredictable availability
- Is heterogeneous



"Controlled" Distributed Systems

- Mechanisms that work well:
 - Consistent hashing (DynamoDB, Cassandra)
 - Master nodes, central coordinator

 Network is under control or client/server → no NAT issues

"Fully" Decentralized Systems

- Mechanisms that work well:
 - Consistent hashing (DHTs)
 - Flooding/broadcasting Bitcoin
- NAT and direct connectivity huge problem



"Controlled" Distributed Systems

- Consistency
 - Leader election (Zookeeper, Paxos, Raft)

- Replication principles
 - More replicas: higher availability, higher reliability, higher performance, better scalability, but: requires maintaining consistency in replicas
- Transparency principles apply

"Fully" Decentralized Systems

- Consistency
 - Weak consistency: DHTs
 - Nakamoto consensus (aka proof of work)
 - Proof of stake Leader election, PBFT protocols Is Bitcoin eventually consistent?
 - Some argue no, some argue it has even stronger guarantees
- Replication principles apply to fully decentralized systems as well
- Transparency principles apply



- Spring Term Distributed Systems (DSy)
 - Tightly/loosely coupled
 - Heterogeneous systems
 - Small-scale systems
 - Distributed systems

(we will also talk about blockchains in this lecture)

- Fall Term Blockchain (BICh)
 - Loosely coupled
 - Heterogeneous systems
 - Large-scale systems
 - Decentralized systems

(we will also talk about distributed systems in this lecture, but DSy is highly recommended)





Access Token / Refresh Token

- Access Token only short lifetime, e.g., 10min.
 - If public key / secret is known, the content in the token can be trusted, e.g., in the serivce
 - Can have userId, role, etc.
 - No need to query DB for those information, e.g.:
 type TokenClaims struct {
 MailFrom string `json:"mail_from,omitempty"`
 MailTo string `json:"mail_to,omitempty"`
 jwt.Claims
 }
- Refresh Token longer lifetime, e.g., 6 month
 - A refresh token is used to get a new access token
 - IAM / Auth server creates access tokens

- Only access token, with long lifetime
 - If a user credential is revoked how to inform every service?
- Only refresh token
 - Tightly coupled Service/Auth, every request to Service, Auth needs to be involved for every access
- Access + Refresh token
 - If a user credential is revoked, user has max.
 10min more to access service
 - Auth only involved if access token is expired





Networking: Layers

- Networking: Each vendor had its own proprietary solution not compatible with another solution
 - IPX/SPX 1983, AppleTalk 1985, DECnet 1975, XNS 1977
- Nowadays most vendors build compatible networks hardware/software from different vendors
 - Cisco, Dell, HP, Huawei, Juniper, Lenovo, Linksys, Netgear, MicroTik, Siemens, Ubiquiti, etc.
- Goal of layers: interoperability
 - 1984: ISO 7498 The Basic Reference Model for Open Systems Interconnection

OSI model	"Internet model"
Application	Application
Presentation	
Session	
Transport	Transport
Network	Internet
Data link	Link
Pysical	

			Data
		TCP Header	Data
	IP Header	TCP Header	Data
Ethernet Header	IP Header	TCP Header	Data



TCP/IP from an Application Developer View

- Server in golang (repo)
 - git clone https://github.com/tboce k/DSy
 - Download GoLand, or others
 - go run server.go → server
- Listening on TCP port 8081
 - Return string in uppercase
- Node.js version
 - Download WebStorm, or other
- Client:
 - nc localhost 8081

```
const net = require('net');
const server = new net.Server();
server.listen(8081, function() {
    console.log('Launching server...');
});

server.on('connection', function(socket) {
    socket.on('data', function(chunk) {
        console.log(`Data received from client: $
    {chunk.toString()}`);

socket.write(chunk.toString().toUpperCase() +
"\n");
    });
});
```

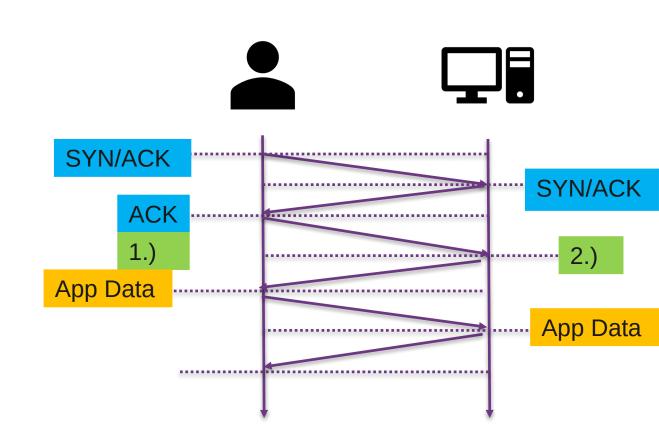
```
package main
import ("bufio"
    "fmt"
    "net"
    "strings")
func main() {
    fmt.Println("Launching server...")
    ln, _ := net.Listen("tcp", ":8081") // listen
on all interfaces
    for {
        conn, _ := ln.Accept() // accept connection
on port
        message, _ :=
bufio.NewReader(conn).ReadString('\n') //read line
        fmt.Print("Message Received:",
string(message))
        newMessage := strings.ToUpper(message)
//change to upper
        conn.Write([]byte(newMessage + "\n"))
//send upper string back
```



```
PING sydney.edu.au (129.78.5.8) 56(84) bytes of data.
64 bytes from scilearn.sydney.edu.au (129.78.5.8): icmp_seq=1 ttl=233 time=307 ms
64 bytes from scilearn.sydney.edu.au (129.78.5.8): icmp_seq=2 ttl=233 time=305 ms
64 bytes from scilearn.sydney.edu.au (129.78.5.8): icmp_seq=3 ttl=233 time=305 ms
```

Layer 4 – TCP + TLS

- Ping to Australia: 329ms
 - One way ~ 165ms
- TCP + TLS handshake:
 - 3RTT = 987ms! No data sent yet
- TLS 1.3, finished Aug 2018
 - <u>1 RTT</u> instead of 2
 - 1.) Client Hello, Key Share
 - 2.) Server Hello, key Share, Verify Certificate, Finished
 - 0 RTT possible, for previous connections, loosing perfect forward secrecy
 - 90% of browsers used already support it





QUIC / HTTP/3

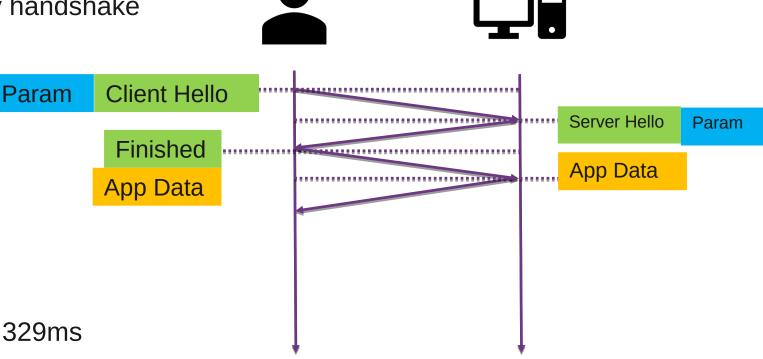
• QUIC: 1RTT connection + security handshake

For known connections: 0RTT

Built in security

- "Google's 'QUIC' TCP alternative slow to excite anyone outside Google" [link] (9%, 25%, 75%)
 - Facebook
 - Cloudflare, state of HTTP

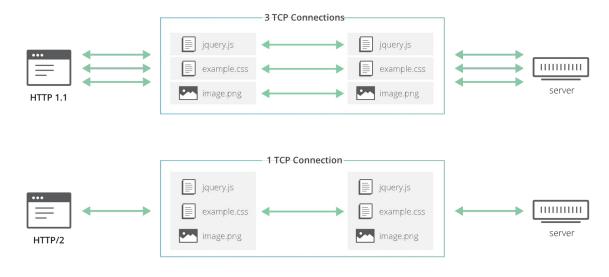
Example Australia: from 987ms to 329ms



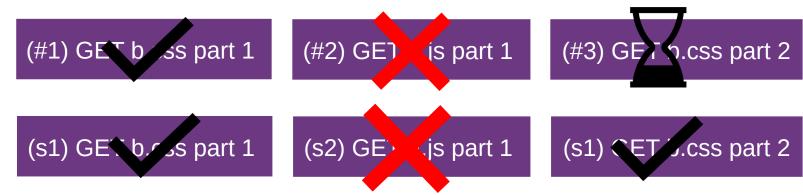


QUIC / HTTP3

- Multiplexing in HTTP/2
 - HTTP/1 → HTTP/2
- HTTP/2: Head-of-line blocking
 - One packet loss, TCP needs to be ordered
 - QUIC can multiplex requests: one stream does not affect others
- HTTP/3 is great, but...
 - NAT → SYN, ACK, FIN, conntrack knows when connection ends, not with QUIC, timeouts, new entries, many entries
 - HTTP header compression, referencing previous headers
 - Many TCP <u>optimizations</u>



source: https://blog.cloudflare.com/the-road-to-quic/







Protocols

- Custom encoding/decoding
 - You control every aspect
 - You send more time on it
- Little-endian / Big-endian
 - sequential order where bytes are converted into numbers

115

118

121

123

124

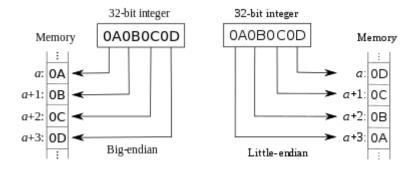
126

131

134

- Networking, e.g. TCP headers: Big-endian
- Most CPUs e.g., x86:
 Little-endian, RISC-V: Bi-endianness

```
public static boolean decodeHeader(final ByteBuf buffer, final InetSocketAddress recipientSocket,
       final InetSocketAddress senderSocket, final Message message) {
   LOG.debug("Decode message. Recipient: {}, Sender:{}.", recipientSocket, senderSocket);
   final int versionAndType = buffer.readInt();
   message.version(versionAndType >>> 4);
   message.type(Type.values()[(versionAndType & Utils.MASK_0F)]);
   message.protocolType(ProtocolType.values()[versionAndType >>> 30]);
   message.messageId(buffer.readInt());
   final int command = buffer.readUnsignedByte();
   message.command((byte) command);
   final Number160 recipientID = Number160.decode(buffer);
   //we only get the id for the recipient, the rest we already know
   final PeerAddress recipient = PeerAddress.builder().peerId(recipientID).build();
   message.recipient(recipient);
   final int contentTypes = buffer.readInt();
   message.hasContent(contentTypes != 0);
    message.contentTypes(decodeContentTypes(contentTypes, message));
```





Application Protocol: HTTP

- HTTP (HyperText Transfer Protocol): foundation of data communication for www
- Started in 1989 by Tim Berners-Lee
 - HTTP/1.1 published in 1997
 - HTTP/2 published in 2015
 - More efficient, header compression, multiplexing
 - HTTP/3 published in 2022
- Request / response (resource)
- HTTP resources identified by URL
 - https://dsl.i.ost.ch/design/ost_logo.svg

```
Scheme
            User info
```

Host

Text-based protocol

```
openssl s_client -connect dsl.i.ost.ch:443 -showcerts
... TLS handshake ...
GET /
```

- HTTP is a stateless protocol
 - Server maintains no state
- Browser sends a bit more...

```
▼ Request Headers (359 B)

                                                                   Host: dsl.hsr.ch
                                                                   User-Agent: Mozilla/5.0 (X11; Linux x86_64; rv:73.0) Gecko/20100101 Firefox/73.0
                                                                   Accept: text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,*/*;q=0.8
                                                                   Accept-Language: en-US, en; q=0.5
                                                                   Accept-Encoding: gzip, deflate, br
                                                                   DNT: 1
                                                                   Connection: keep-alive
                                                                   Upgrade-Insecure-Requests: 1
                                                                   Cache-Control: max-age=0
                                                                   TE: Trailers
                                                                                               Query
                                                                                                                  Fragment
                                                              Port
                                                                           Path
http://tbocek:password@dsl.i.ost.ch:443/lect/fs21?id=1234&lang=de#topj
```



Deployment Strategies

- Many strategies and variations [link, link, link]
- Rolling Deployment
 - New version is gradually deployed to replace the old version - without taking the entire system down at once
 - + Minimal downtime, low risk
 - Complexity, longer deployment times
- Blue-Green Deployment
 - 2 environments, current prod (blue), current prod with new release (green). Test, then switch
 - + Instant rollback, 0 downtime
 - 2 prod environments, keep data in sync

Canary Releases

- Canary in a coal mine new version to a small group of users or servers first, if all goes well, more users
 - + Risk reduction, user feedback
 - Complexity, inconsistencies
- Feature Toggle
 - Fine grained canary, set feature for specific users
 - + More risk reduction, specific user feedback
 - Increase complexity of codebase, config management
- Big Bang
 - Deploy everything at once
 - + Simple
 - High risk, limited rollback



Latency Numbers Every Programmer Should Know

- Interactive [link] from 1990 2020
 - Network stays ~ 150ms
 - L1: 1ns / branch miss 3ns example
- HDD / SSD / NVMe (Non-Volatile Memory Express) - comparison, 2

- Latency and throughput important
- Napkin Math [link]
 - Cost

	_			
NVMe	Latency vs	Other	Storage	Protocols
I 4 A IAIC	Latelley vs	. Other	Otol age	1 10000013

Feature	NVMe	SATA SSD	HDD
Read Latency	~20 µs	~100 µs	2-5 ms
Write Latency	~30 µs	~200 µs	5-10 ms
Queue Depth	64K queues × 64K commands	32 commands	1 command
Bandwidth	Up to 16GB/s (PCIe 4.0)	600MB/s	~150MB/s

Operation	Latency	Throughput	1 MiB	1 GiB
Sequential Memory R/W (64 bytes)	0.5 ns			
- Single Thread, No SIMD		10 GiB/s	100 μs	100 ms
- Single Thread, SIMD		20 GiB/s	50 μs	50 ms
- Threaded, No SIMD		30 GiB/s	35 μs	35 ms
⊢ Threaded, SIMD		35 GiB/s	30 µs	30 ms
Network Same-Zone		10 GiB/s	100 μs	100 ms
- Inside VPC		10 GiB/s	100 μs	100 ms
- Outside VPC		3 GiB/s	300 μs	300 ms
Hashing, not crypto-safe (64 bytes)	25 ns	2 GiB/s	500 μs	500 ms
Random Memory R/W (64 bytes)	50 ns	1 GiB/s	1 ms	1s
Fast Serialization [8] [9] †	N/A	1 GiB/s	1 ms	1s
Fast Deserialization [8] [9] †	N/A	1 GiB/s	1 ms	1s
System Call	500 ns	N/A	N/A	N/A
Hashing, crypto-safe (64 bytes)	100 ns	1 GiB/s	1 ms	1s
Sequential SSD read (8 KiB)	1 μs	4 GiB/s	200 μs	200 m
Context Switch [1] [2]	10 μs	N/A	N/A	N/A
Sequential SSD write, -fsync (8KiB)	10 μs	1 GiB/s	1 ms	1s
TCP Echo Server (32 KiB)	10 μs	4 GiB/s	200 μs	200 m
Decompression [11]	N/A	1 GiB/s	1 ms	1s
Compression [11]	N/A	500 MiB/s	2 ms	2s
Sequential SSD write, +fsync (8KiB)	1 ms	10 MiB/s	100 ms	2 min
Sorting (64-bit integers)	N/A	200 MiB/s	5 ms	5s
Sequential HDD Read (8 KiB)	10 ms	250 MiB/s	2 ms	25



Introduction

- Bitcoin is an <u>experimental</u> digital currency
 - Bitcoin is fully peer-2-peer (no central entity)
 - 1st Bitcoin issued on January 3, 2009
 - Smallest unit: 0.00000001 BTC (1 satoshi)
- Key characteristics
 - Maximum of ~21 million BTC
 - Every transaction broadcast to all peers
 - Every peers knows all transactions (~660 GByte as of today)
 - Validation by proof-of-work (partial hash collision)
 - Difficult to fake proof-of-work
 - No double-spending
- The initiator is unknown so far





Bitcoin - Introduction

- Not relying on trust, but on strong cryptography
- Weak anonymity (pseudonimity)
 - All peers know all transactions
 - Clustering: e.g. if a transaction has multiple input addresses, assume those addresses belong to the same wallet. (example)
- Not controlled by a single entity
 - Development community, no central bank forks Bitcoin Cash, SV
- BIP: Bitcoin Improvement Proposals
- Bitcoins can be exchange for real currencies
 - Several companies allow to exchange BTC for Dollar, Euro, ...
- US, CH considered Bitcoin friendly, China (energy) not that much



Mechanism

- A wallet has public-private keys (wallet.dat)
 - Public key, ECDSA 256 bit → Bitcoin address (can receive bitcoins)
 - Simple address ~ base58(RIPEM160(Sha256(ecdsa public key)))
 - E.g. 1GCeaKuhDYnNLNR6LGmBtKhPqEJD4KeEtF
 - Private key used for signing transactions



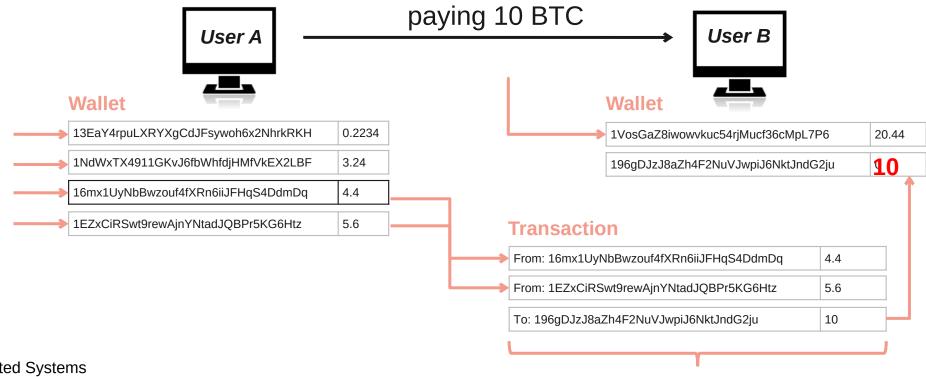
- Peer A wants to send BTC to peer B → creates transaction message
- Transaction contains input / output
 - where the BTC came from and where it goes
- Peer A broadcasts the transaction to all the peers in the network
- Transaction stored in blocks → block is created / verified ~10min





Key Bitcoin Operations

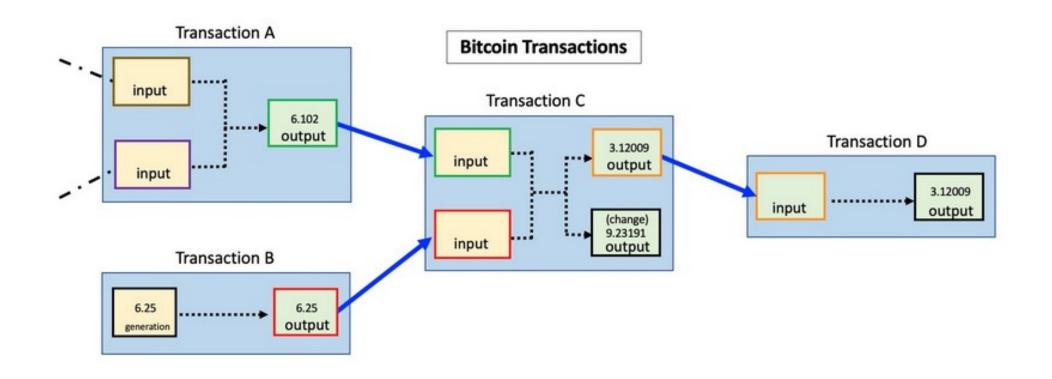
- Private key authorizes the transaction ("access")
 - If keys are stolen, thief may use "your" coins
 - If keys are lost, coins are lost
 - In UTXO (unspent transaction output) systems, complete output is spent



Sign with Private Key of User A



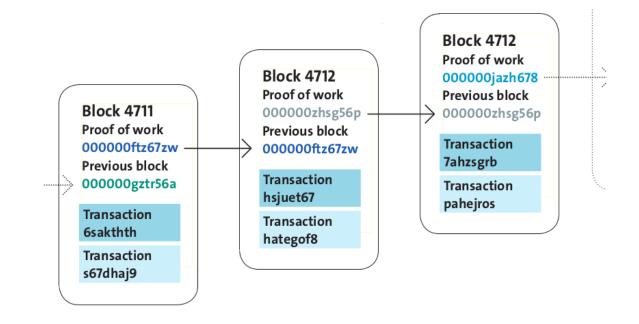
Transactions





Blockchain

- Transactions are collected in blocks
 - New block created approximately every 10 min
- Blocks contain solved crypto puzzles
 - In the form of partial <u>hash collisions</u> (SHA256)
- A block has a pointer to previous <u>block</u> → Blockchain
- Creation of blocks is called mining (reward)
 - Mining / creating blocks → Miner get currently
 3.125 BTC per creation
 - adjustable difficulty 6 blocks / h
 - Sometime in 2028 reward will be 1.5625





Discussion (1)

- Disadvantages
 - Power consumption
 - ~ as much as Poland
 - Not scalable
 - Bitcoin with ~7 tps vs. VISA 57,000 tps (23.12) [tps: transactions per sec]



- Anonymity
 - Can be used for illegal activities

- Advantages
 - Low (fixed) tx fees
 - \sim 1.2 satoshi per byte / 0.25USD (\sim 200bytes tx)
 - Scalable
 - Hardware/storage gets faster



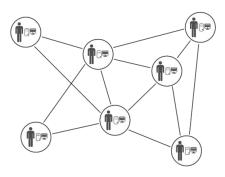
- Anonymity
 - Preserving privacy



Discussion (2)

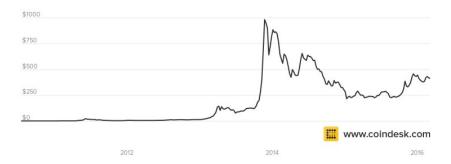
- Advantages
 - No major "crashes"
 - Mt.Gox / FTX was exchange site!
 - Decentralized
 - Open protocol
 - Forks





- Many other blockchain use cases
 - Smart contracts

- Disadvantages
 - Volatile exchange rate



- Central elemements
 - Core developers
 - Mining farms [link]

