

# Distributed Systems

The second half of *Concurrent and Distributed Systems*

<https://www.cl.cam.ac.uk/teaching/current/ConcDisSys>

Dr. Martin Kleppmann (mk428@cam)

University of Cambridge

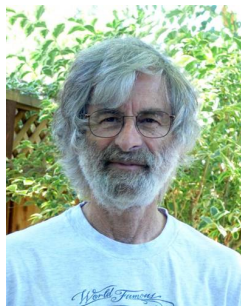
Computer Science Tripos, Part IB



This work is published under a  
Creative Commons BY-SA license.

# A distributed system is...

- ▶ *"... a system in which the failure of a computer you didn't even know existed can render your own computer unusable."* — Leslie Lamport



# A distributed system is...

- ▶ *“...a system in which the failure of a computer you didn't even know existed can render your own computer unusable.”* — Leslie Lamport
- ▶ ...multiple computers communicating via a network...
- ▶ ...trying to achieve some task together
- ▶ Consists of “nodes” (computer, phone, car, robot, ...)

# Recommended reading

- ▶ van Steen & Tanenbaum.  
**“Distributed Systems”**  
(any ed), free ebook available
- ▶ Cachin, Guerraoui & Rodrigues.  
**“Introduction to Reliable and Secure Distributed Programming”** (2nd ed), Springer 2011
- ▶ Kleppmann.  
**“Designing Data-Intensive Applications”**,  
O'Reilly 2017
- ▶ Bacon & Harris.  
**“Operating Systems: Concurrent and Distributed Software Design”**, Addison-Wesley 2003

# Relationships with other courses

- ▶ **Concurrent Systems** – Part IB  
(every distributed system is also concurrent)
- ▶ **Operating Systems** – Part IA  
(inter-process communication, scheduling)
- ▶ **Databases** – Part IA  
(many modern databases are distributed)
- ▶ **Computer Networking** – Part IB Lent term  
(distributed systems involve network communication)
- ▶ **Further Java** – Part IB Michaelmas  
(distributed programming practical exercises)
- ▶ **Security** – Part IB Easter term  
(network protocols with encryption & authentication)
- ▶ **Cloud Computing** – Part II  
(distributed systems for processing large amounts of data)

# Why make a system distributed?

# Why make a system distributed?

- ▶ **It's inherently distributed:**

e.g. sending a message from your mobile phone to your friend's phone

# Why make a system distributed?

- ▶ **It's inherently distributed:**  
e.g. sending a message from your mobile phone to your friend's phone
- ▶ **For better reliability:**  
even if one node fails, the system as a whole keeps functioning

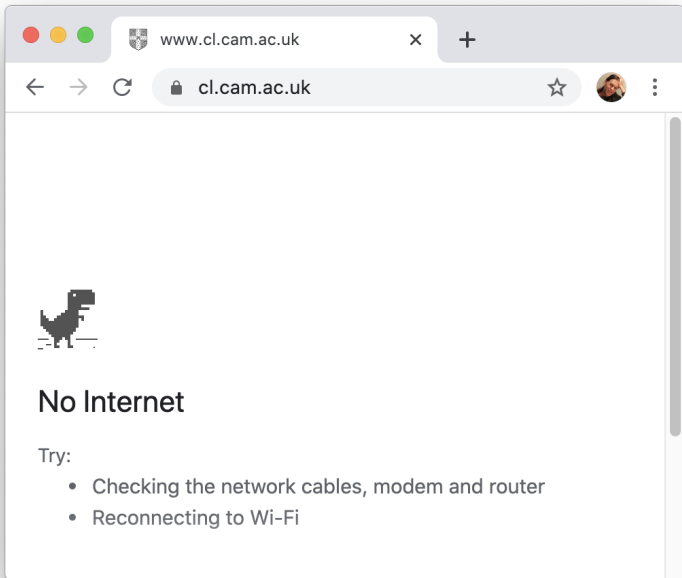


# Why make a system distributed?

- ▶ **It's inherently distributed:**  
e.g. sending a message from your mobile phone to your friend's phone
- ▶ **For better reliability:**  
even if one node fails, the system as a whole keeps functioning
- ▶ **For better performance:**  
get data from a nearby node rather than one halfway round the world

# Why make a system distributed?

- ▶ **It's inherently distributed:**  
e.g. sending a message from your mobile phone to your friend's phone
- ▶ **For better reliability:**  
even if one node fails, the system as a whole keeps functioning
- ▶ **For better performance:**  
get data from a nearby node rather than one halfway round the world
- ▶ **To solve bigger problems:**  
e.g. huge amounts of data, can't fit on one machine



# Why NOT make a system distributed?

The trouble with distributed systems:

- ▶ Communication may fail (and we might not even know it has failed).
- ▶ Processes may crash (and we might not know).
- ▶ All of this may happen nondeterministically.

# Why NOT make a system distributed?

The trouble with distributed systems:

- ▶ Communication may fail (and we might not even know it has failed).
- ▶ Processes may crash (and we might not know).
- ▶ All of this may happen nondeterministically.

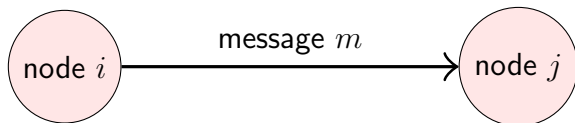
**Fault tolerance:** we want the system as a whole to continue working, even when some parts are faulty.

This is hard.

Writing a program to run on a single computer is comparatively easy?!

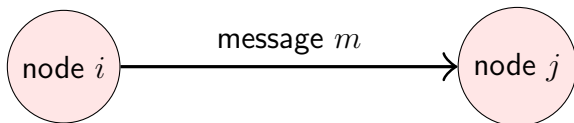
# Distributed Systems and Computer Networking

We use a simple abstraction of communication:



# Distributed Systems and Computer Networking

We use a simple abstraction of communication:



Reality is much more complex:

- ▶ **Various network operators:**  
eduroam, home DSL, cellular data, coffee shop wifi, submarine cable, satellite. . .
- ▶ **Physical communication:**  
electric current, radio waves, laser, hard drives in a van. . .

# Hard drives in a van?!



<https://docs.aws.amazon.com/snowball/latest/ug/using-device.html>

High latency, high bandwidth!



# Latency and bandwidth

**Latency:** time until message arrives

- ▶ In the same building/datacenter:  $\approx 1$  ms
- ▶ One continent to another:  $\approx 100$  ms
- ▶ Hard drives in a van:  $\approx 1$  day

# Latency and bandwidth

**Latency:** time until message arrives

- ▶ In the same building/datacenter:  $\approx 1$  ms
- ▶ One continent to another:  $\approx 100$  ms
- ▶ Hard drives in a van:  $\approx 1$  day

**Bandwidth:** data volume per unit time

- ▶ 3G cellular data:  $\approx 1$  Mbit/s
- ▶ Home broadband:  $\approx 10$  Mbit/s
- ▶ Hard drives in a van: 50 TB/box  $\approx 1$  Gbit/s

(Very rough numbers, vary hugely in practice!)


Concurrent and Distributed Sys

+

← → ↻

Not Secure | cst.cam.ac.uk/teaching/2021/ConcDisSys

☆ ⚙️ 👤 ⋮

 UNIVERSITY OF CAMBRIDGE

[Study at Cambridge](#)

[About the University](#)

[Research at Cambridge](#)

Quick links

▼

Search

🔍

Department of Computer Science and Technology

Log in with Raven

Home

The department ▼

Initiatives ▼

Research ▼

Admissions ▼

Current students ▼

Job vacancies ➡️

Intranet

Concurrent and Distributed Systems

**Principal lecturer:** Dr David Greaves  
Martin Kleppmann

**Students:** Part IB CST 50%, Part IB CST 75%

**Course code:** ConcDisSys

**Prerequisite course:** [Object-Oriented Programming](#)  
[Operating Systems](#)

**This course is a prerequisite for:** [Cloud Computing](#)  
[Distributed Ledger Technologies: Foundations and Applications](#)  
[Mobile and Sensor Systems](#)

Related links

[Course materials](#)

[Information for supervisors](#)

# Client-server example: the web

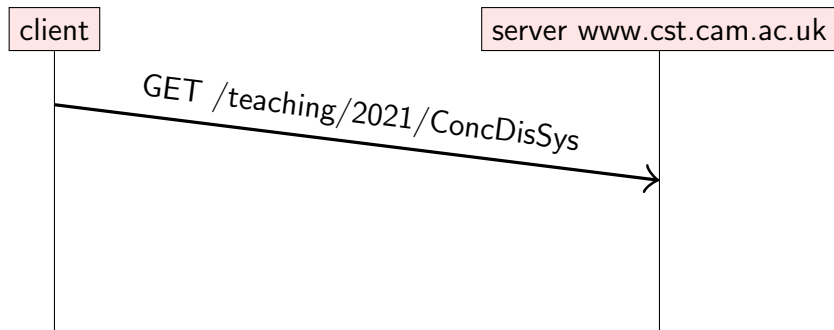
Time flows from top to bottom.

client

server [www.cst.cam.ac.uk](http://www.cst.cam.ac.uk)

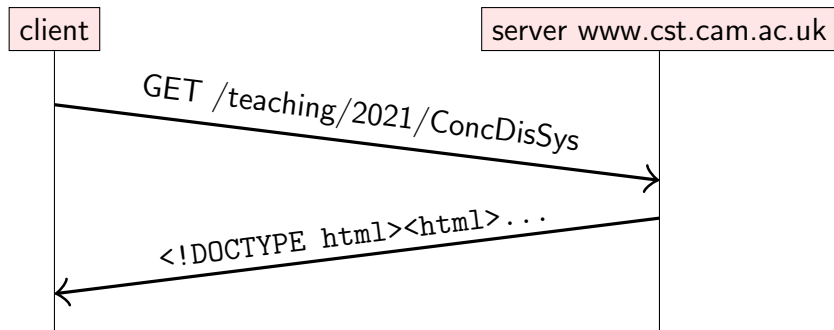
# Client-server example: the web

Time flows from top to bottom.

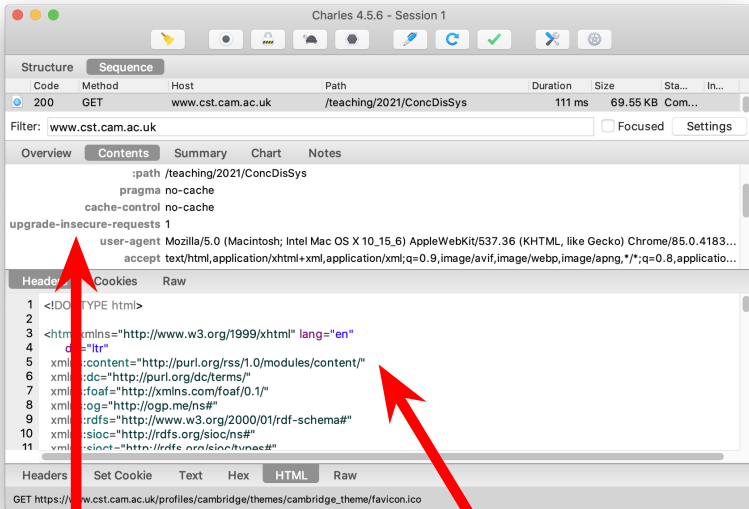


# Client-server example: the web

Time flows from top to bottom.







request message

response message



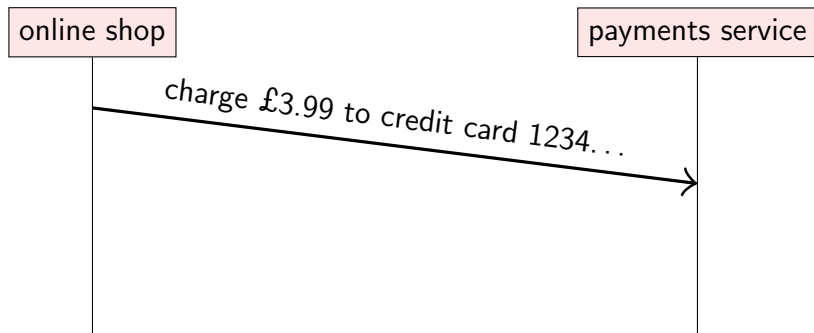


# Client-server example: online payments

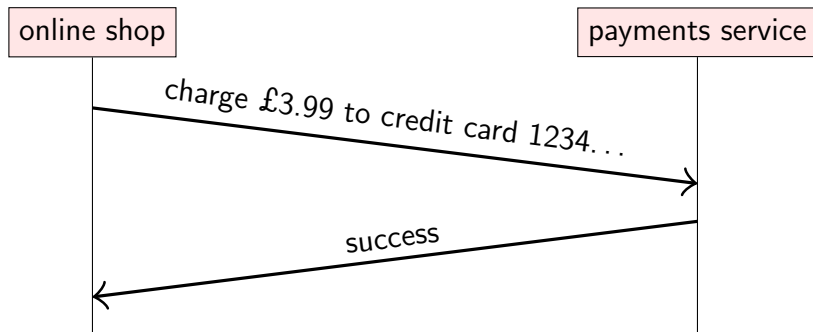
online shop

payments service

# Client-server example: online payments



# Client-server example: online payments



# Remote Procedure Call (RPC) example

*// Online shop handling customer's card details*

```
Card card = new Card();  
card.setCardNumber("1234 5678 8765 4321");  
card.setExpiryDate("10/2024");  
card.setCVC("123");
```

```
Result result = paymentsService.processPayment(card,  
    3.99, Currency.GBP);
```

```
if (result.isSuccess()) {  
    fulfilOrder();  
}
```

# Remote Procedure Call (RPC) example

*// Online shop handling customer's card details*

```
Card card = new Card();  
card.setCardNumber("1234 5678 8765 4321");  
card.setExpiryDate("10/2024");  
card.setCVC("123");
```

```
Result result = paymentsService.processPayment(card,  
    3.99, Currency.GBP);
```

```
if (result.isSuccess()) {  
    fulfilOrder();  
}
```



Implementation of this function is on another node!

online shop

RPC client

RPC server

payment service



processPayment() stub

waiting



online shop

RPC client

RPC server

payment service

processPayment() stub

marshal args

$m_1$

unmarshal args

waiting

```
{
  "request": "processPayment",
  "card": {
    "number": "1234567887654321",
    "expiryDate": "10/2024",
    "CVC": "123"
  },
  "amount": 3.99,
  "currency": "GBP"
}
```

$m_1 =$



online shop

RPC client

RPC server

payment service

processPayment() stub

marshal args

$m_1$

unmarshal args

processPayment()  
implementation

waiting

```
{  
  "request": "processPayment",  
  "card": {  
    "number": "1234567887654321",  
    "expiryDate": "10/2024",  
    "CVC": "123"  
  },  
  "amount": 3.99,  
  "currency": "GBP"  
}
```

$m_1 =$

online shop

RPC client

RPC server

payment service

processPayment() stub

waiting

marshal args

$m_1$

unmarshal args

processPayment()  
implementation

unmarshal result

$m_2$

marshal result

$m_1 =$

```
{
  "request": "processPayment",
  "card": {
    "number": "1234567887654321",
    "expiryDate": "10/2024",
    "CVC": "123"
  },
  "amount": 3.99,
  "currency": "GBP"
}
```

$m_2 =$

```
{
  "result": "success",
  "id": "XP61hHw2Rvo"
}
```

online shop

RPC client

RPC server

payment service

processPayment() stub

waiting

function returns

marshal args

$m_1$

unmarshal args

processPayment()  
implementation

unmarshal result

$m_2$

marshal result

```
{
  "request": "processPayment",
  "card": {
    "number": "1234567887654321",
    "expiryDate": "10/2024",
    "CVC": "123"
  },
  "amount": 3.99,
  "currency": "GBP"
}
```

$m_1 =$

```
{
  "result": "success",
  "id": "XP61hHw2Rvo"
}
```

$m_2 =$

# Remote Procedure Call (RPC)

Ideally, RPC makes a call to a remote function look the same as a local function call.

**“Location transparency”:**  
system hides where a resource is located.

# Remote Procedure Call (RPC)

Ideally, RPC makes a call to a remote function look the same as a local function call.

**“Location transparency”:**  
system hides where a resource is located.

In practice. . .

- ▶ what if the service crashes during the function call?
- ▶ what if a message is lost?
- ▶ what if a message is delayed?
- ▶ if something goes wrong, is it safe to retry?

# RPC history

- ▶ SunRPC/ONC RPC (1980s, basis for NFS)
- ▶ CORBA: object-oriented middleware, hot in the 1990s
- ▶ Microsoft's DCOM and Java RMI (similar to CORBA)
- ▶ SOAP/XML-RPC: RPC using XML and HTTP (1998)
- ▶ Thrift (Facebook, 2007)
- ▶ gRPC (Google, 2015)
- ▶ REST (often with JSON)
- ▶ Ajax in web browsers

# RPC/REST in JavaScript

```
let args = {amount: 3.99, currency: 'GBP', /*...*/};
let request = {
  method: 'POST',
  body:    JSON.stringify(args),
  headers: {'Content-Type': 'application/json'}
};
```

```
fetch('https://example.com/payments', request)
  .then((response) => {
    if (response.ok) success(response.json());
    else failure(response.status); // server error
  })
  .catch((error) => {
    failure(error); // network error
  });
```

# RPC in enterprise systems

**“Service-oriented architecture” (SOA) / “microservices”:**

splitting a large software application into multiple services (on multiple nodes) that communicate via RPC.



# RPC in enterprise systems

“**Service-oriented architecture**” (SOA) / “microservices”:

splitting a large software application into multiple services (on multiple nodes) that communicate via RPC.

Different services implemented in different languages:

- ▶ interoperability: datatype conversions
- ▶ **Interface Definition Language (IDL)**:  
language-independent API specification

# gRPC IDL example

```
message PaymentRequest {  
    message Card {  
        required string cardNumber = 1;  
        optional int32 expiryMonth = 2;  
        optional int32 expiryYear = 3;  
        optional int32 CVC          = 4;  
    }  
    enum Currency { GBP = 1; USD = 2; }  
  
    required Card      card      = 1;  
    required int64      amount    = 2;  
    required Currency   currency = 3;  
}  
  
message PaymentStatus {  
    required bool    success      = 1;  
    optional string  errorMessage = 2;  
}  
  
service PaymentService {  
    rpc ProcessPayment(PaymentRequest) returns (PaymentStatus) {}  
}
```