



# **CS5412/LECTURE 11**

## **THE GEOSCALE CLOUD**

**Ken Birman**  
**CS5412 Spring 2022**

# CLOUD SYSTEMS HAVE GLOBAL SCALE



... and so they need information from global sources, consolidated to wherever a query occurs

For example, with Blockchain, anyone anywhere can see every transaction and validate it. With social networks, whether I am at home or visiting in Seattle, I get updates from all my family and friends.

To work well, global cloud systems need global data replication!

# THESE ARE NOT ONLY ABOUT REPLICATION!

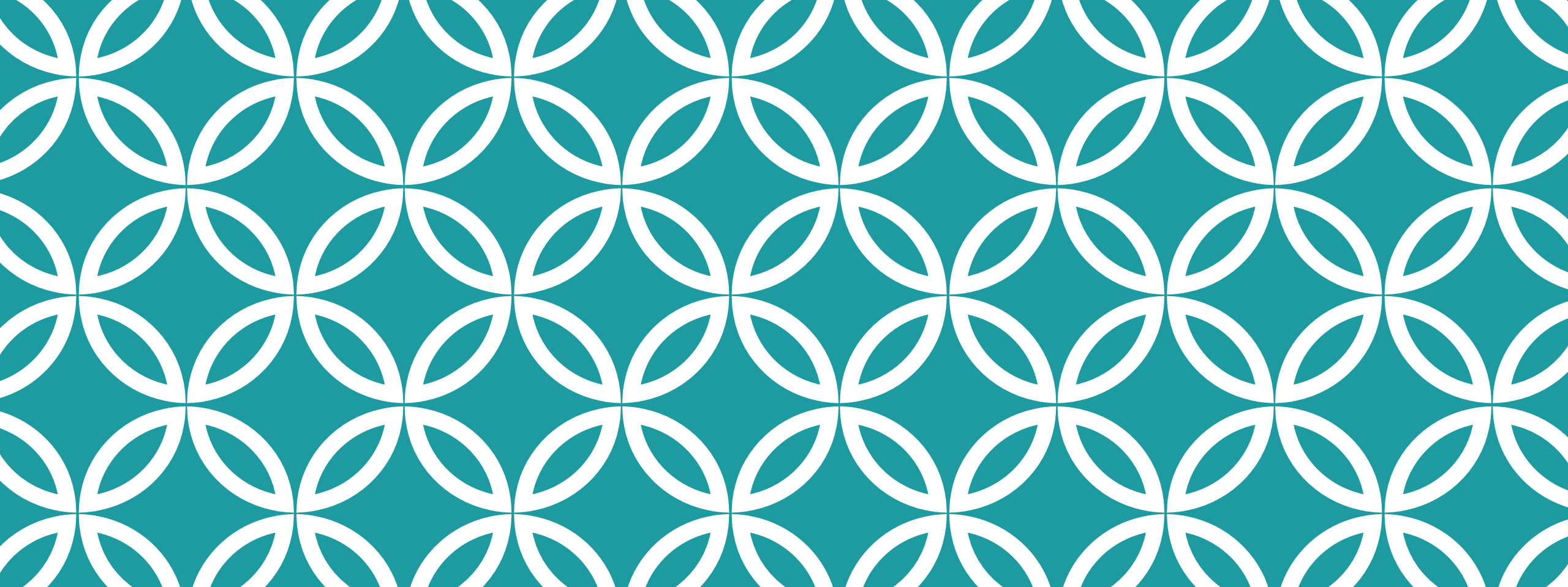
We have point-to-point, subset, and true global replication patterns. Moreover, replication is one of a few dimensions of global state sharing.

Questions of consistency also arise: “how strong should the guarantees be?”

Today we will see a series of topics on this broad issue

# TODAY'S LECTURE: TOPIC BREAKDOWN

- Concept: Regions with Availability Zones (data centers) in them.
- Constraints: Why we can't just use TCP or RDMA or other standards
- Point to point solutions via "Zone aware" message queuing and apps
- Replication, including Google's unique Spanner atomic multicast
- 5G and how this will leverage but also push beyond the limits



# AVAILABILITY ZONES

Basic building block

# BEYOND THE DATACENTER

Early in the semester we discussed Facebook's global blob cache, but since then talked primarily about technologies used inside a single datacenter.

How do cloud developers approach global-scale application design?

Today we will discuss “georeplication” and look at some solutions.

# AVAILABILITY ZONES

Companies like Amazon and Microsoft faced a problem early in the cloud build-out: servicing a data center can require turning it off!

Why?

- Some hardware components are too critical to service while active, like the datacenter power and cooling systems, or the “spine” of routers.
- Some software components can’t easily be upgraded while running, like the datacenter management system.
- In fact there is a very long list of cases like these

# AVAILABILITY ZONES

So... they decided that instead of building one massive datacenter, they would put two or even three side by side.

When all are active, they spread load over them, so everything stays busy.

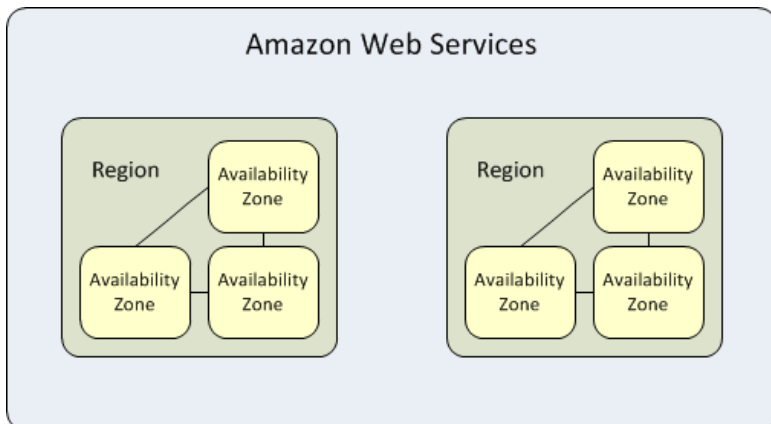
But this also gives them an option for shutting one down entirely to do upgrades (and with three, they would still be able to “tolerate” a major equipment failure in one of the two others).



# AVAILABILITY ZONES – AMAZON AWS

AWS Edge is less capable

The AWS “region” has an availability zone structure

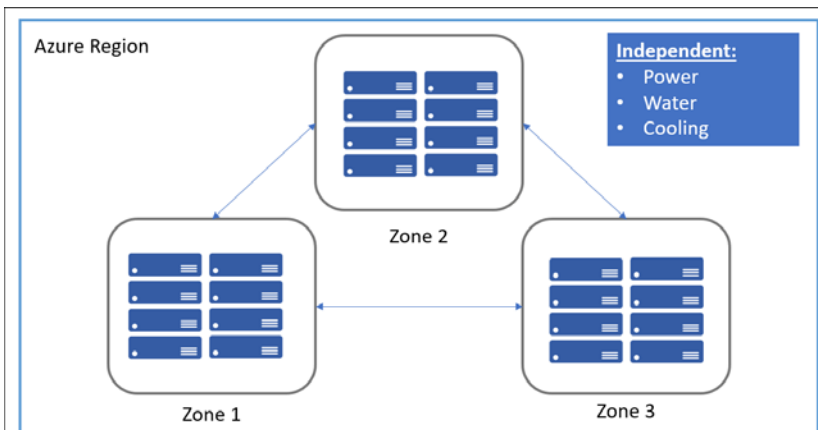


## AWS Regions

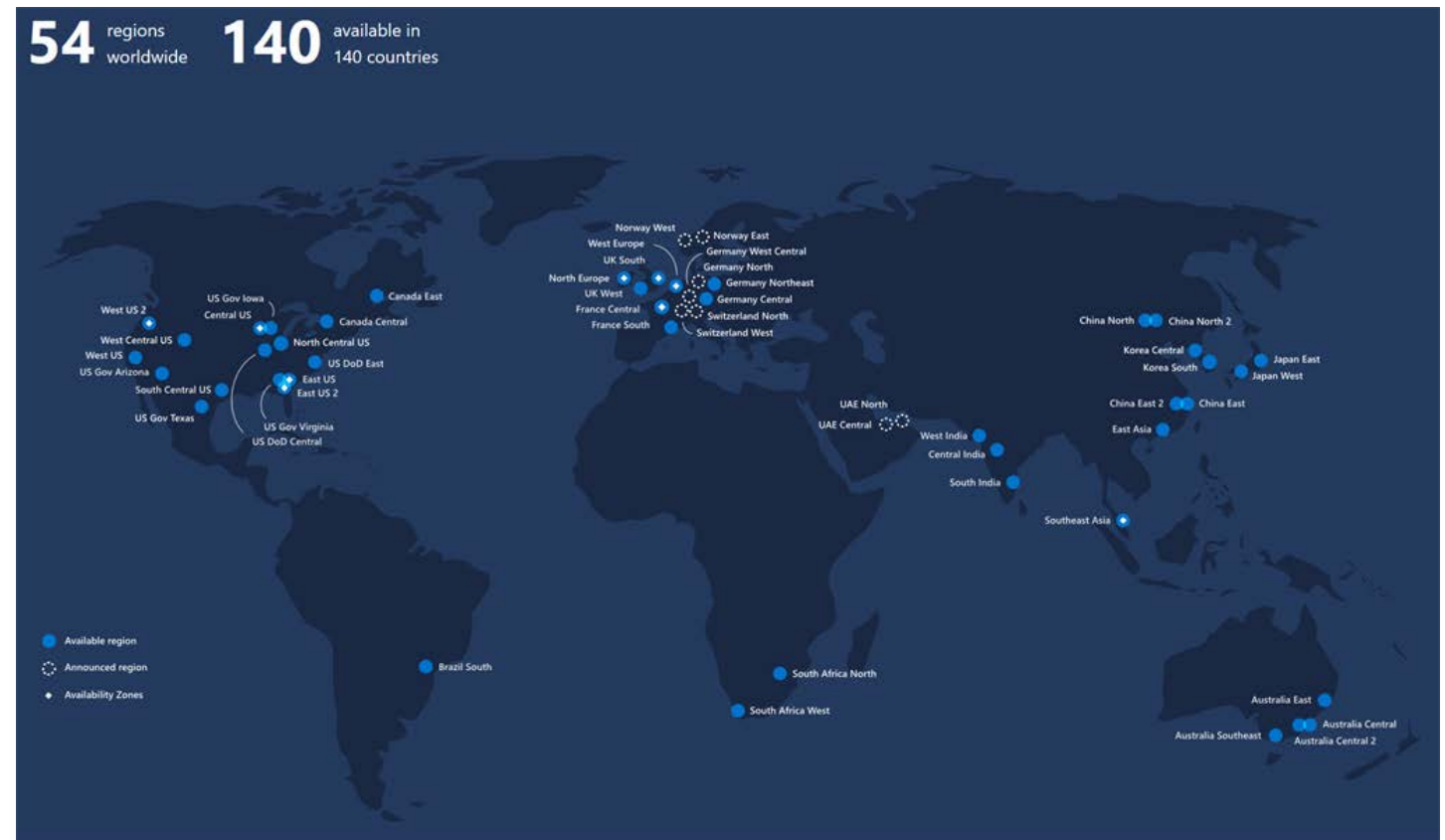


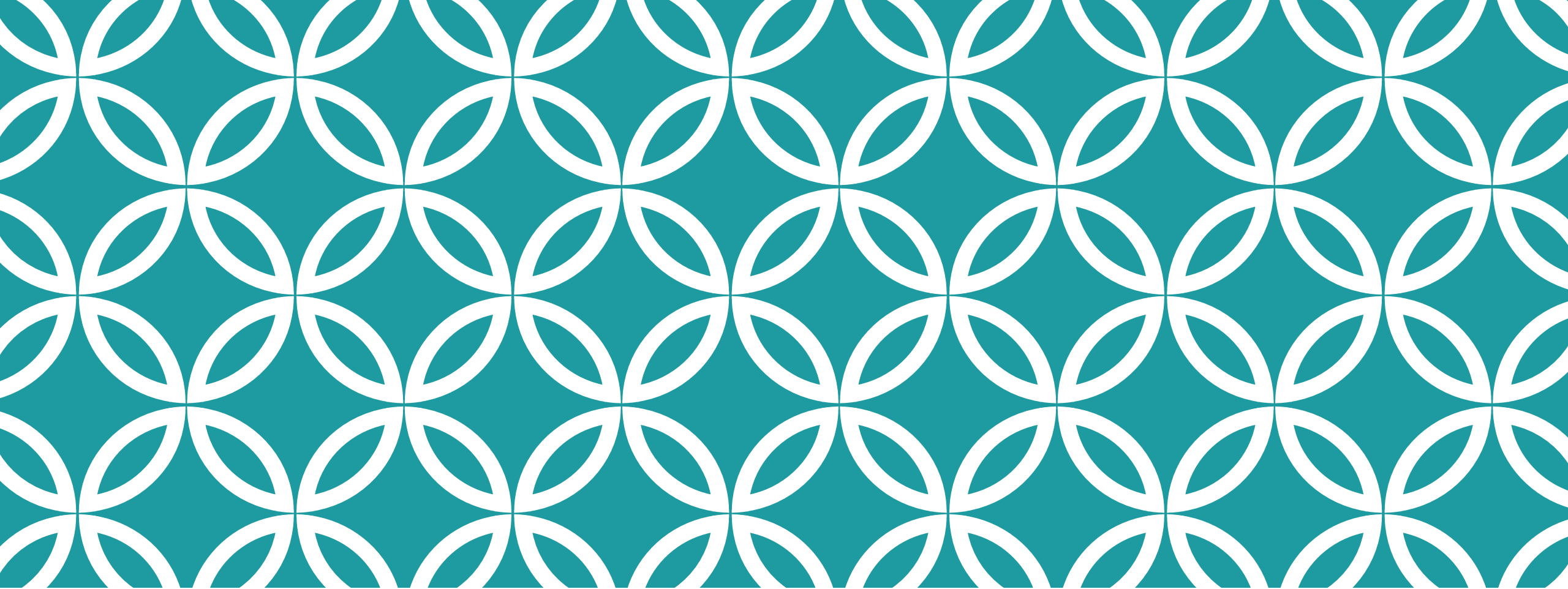
# AVAILABILITY REGIONS AND ZONES – MICROSOFT AZURE

Notice the blue circles.  
Those are regions  
with three data centers.  
A “zone” additionally is  
physically spread out.



Azure Availability Zones are separate buildings 10's of miles apart within a Region





# CONSTRAINTS

**TCP/IP will be blocked!**

# CONNECTIVITY LIMITATIONS

Every datacenter has a firewall forming a very secure perimeter: applications cannot pass data through it without following the proper rules.

Zone-redundant services are provided by AWS and Azure and others to help you mirror data across zones, or even communicate from your service in Zone A to a “sibling” in Zone B.

Direct connections via TCP would probably be blocked: it is easy to connect into a datacenter but hard to connect *out* from inside!

# WHY DO THEY RESTRICT OUTGOING TCP?

The modern datacenter network can have millions of IP addresses inside each single datacenter.

But these won't actually be unique IP addresses if you compare across different data centers: the addresses *only make sense within a data center*. In fact, many IP addresses only make sense *within your own "private cloud"*!

Thus a computer in datacenter A often would not have an IP address visible to a computer in datacenter B, blocking connectivity!

# HOW DOES A BROWSER OVERCOME THIS?

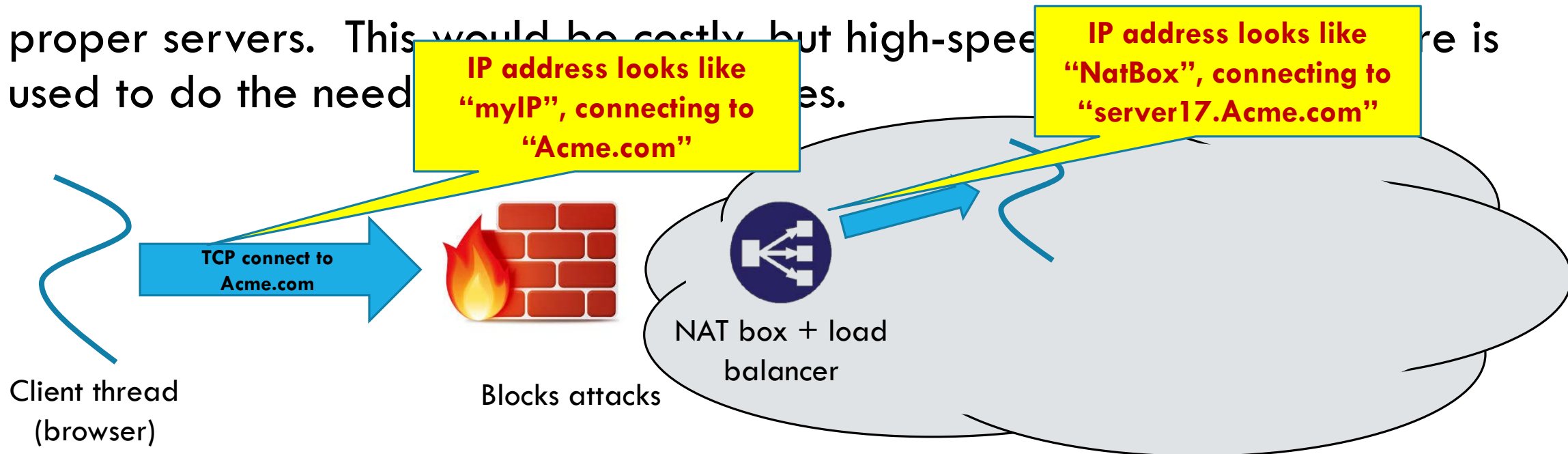
Here at Cornell, your browser is on the *public* internet, not internal to a datacenter.

So it sends a TCP connect request to the datacenter over one of a small set of datacenter IP addresses covering the full datacenter.

- AWS, which hosts for many other sites, has a few IP addresses per site.
- Some systems mimic this approach, others have their own ways of ensuring that traffic to Acme.com gets to Acme's servers, even if hosted on their framework.

# ARRIVING TRAFFIC: NAT BOX

On arrival, packets are scanned for possible BOT traffic or DDoS attacks, then the IP addresses are translated to “load balance” work over the proper servers. This would be costly, but high-speed hardware is used to do the need.



# HOW DOES IT PULL THIS TRICK OFF?

The NAT box maintains a table of “internal IP addresses (and port numbers) and the matching “external” ones.

As messages arrive, if they are TCP traffic, it does a table lookup and substitution, then adjusts the packet header to correct the checksum.

Thus “server17.Acme.com” cannot be accessed directly and yet your messages are routed to it, and vice versa.



# BUT WITH GEOREPLICATION, THIS BLOCKS YOUR CODE FROM CONNECTING TO ITSELF

If you have machine A.Acme.com inside AWS or Azure, and then try to connect to B.Acme.com, it works inside a single datacenter.

- In fact you will be running in an “enterprise VLAN” or “VPC”: what seems to be a private cloud.
- If you were to launch Ethereal or a similar sniffer you only see traffic from your own machines, not traffic from other datacenter tenants.

But if B was in a different datacenter, the connection simply won't work.

- Both A and B are behind NAT boxes, so neither can see the other!

# COULD THIS BE SOLVED?

Actually, yes.

Because two NAT boxes are employed here, Amazon or Azure actually could allow connectivity using NUTSS, a special way of tunnelling

But they don't do so because they don't want uncontrolled connections.

# OPTIONS?

Some vendors (not AWS or Azure) do offer ways to make an A-B connection across datacenters in the same region (availability zone).

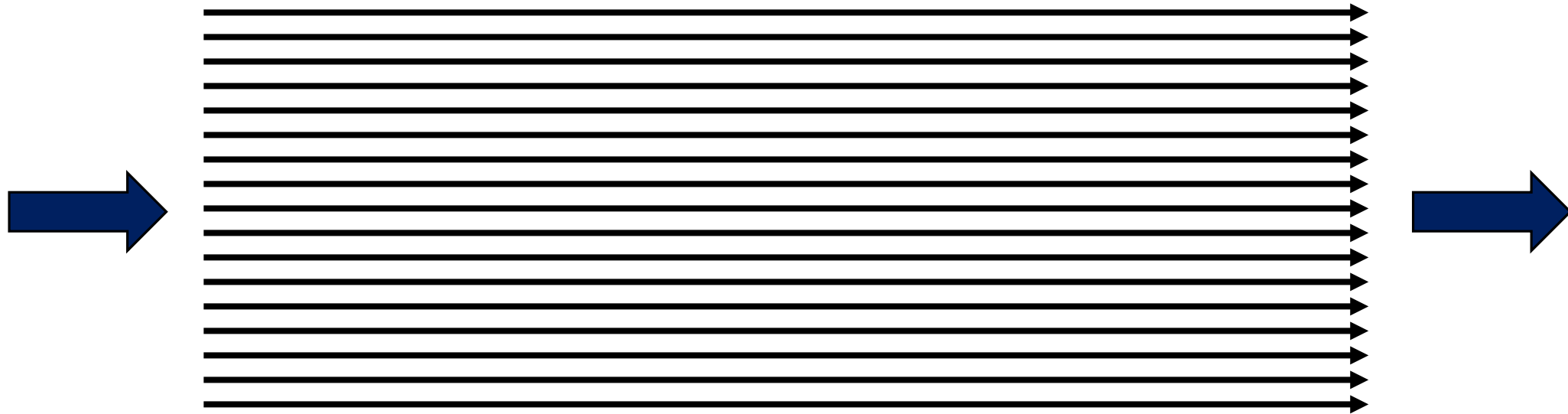
You need to use a special library they provide and otherwise, the connection would fail. And they charge for this service.

With AWS and Azure, you use an existing “Zone-aware” service

# THE BASIC IDEA



Many side-by-side TCP links, on a dedicated optical network



# ZONE AWARE SERVICES ON AZURE

Linux Virtual Machines

Windows Virtual Machines

Virtual Machine Scale Sets

Managed Disks

Load Balancer

Public IP address

Zone-redundant storage

SQL Database

Event Hubs

Service Bus

VPN Gateway

ExpressRoute

Application Gateway

# A FEW WORTH NOTING

**Zone-aware storage** is a storage mirroring option.

Files written in zone A would be available as read-only replicas in zone B. B sees them as a read-only file volume under path `/Zone/A/...`

This is a very common way to share information between data centers.

# A FEW WORTH NOTING



The **Azure Service Bus** in its “Premium” configuration

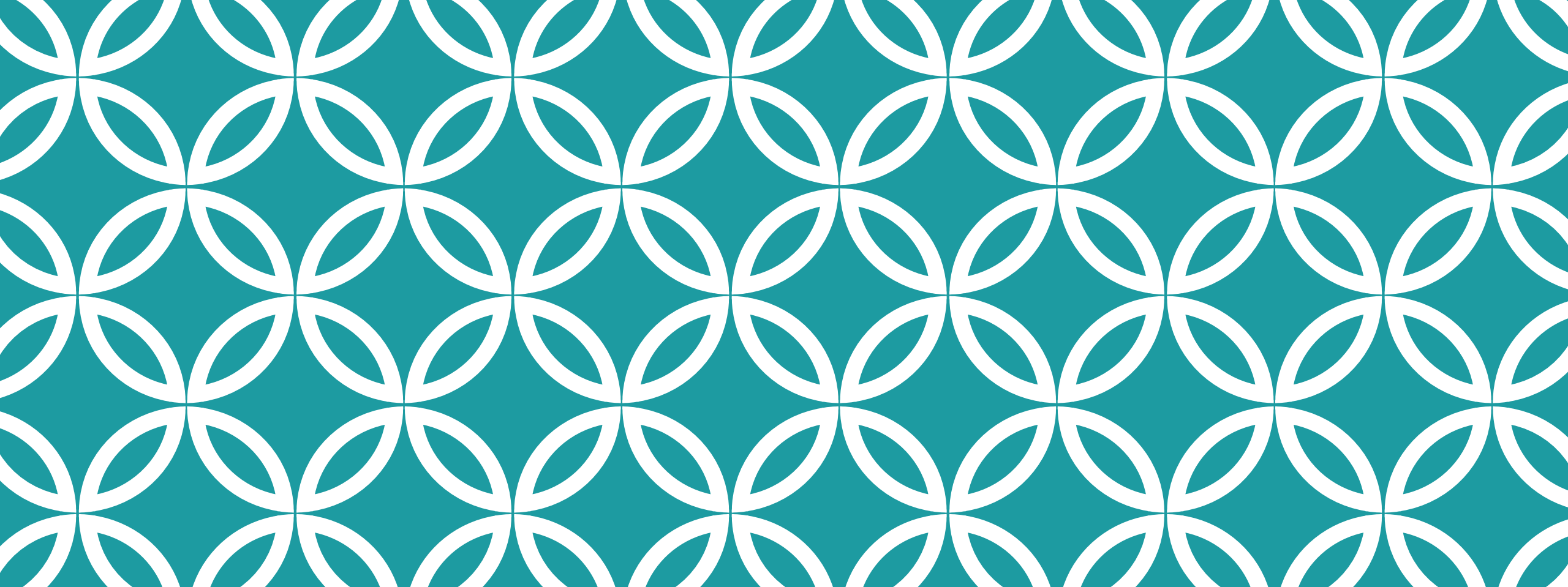
- This is a message bus used by services within your VPC to communicate.
- Azure offers a configuration that automatically transfers data across zones under your control.
- Again, you need to follow their instructions to set it up. They charge but the performance and rate have historically favored this model.
- Basically, two queues hold messages, and then they use a set of side by side TCP connections to shuttle data in both directions. Very efficient!

# A FEW WORTH NOTING

## Azure's **zone-redundant virtual network gateway**

- Used when you really do want a connection of your own, via TCP
- Setup is fairly sophisticated but they walk you through it.
- In effect, creates a special “pseudo-public” IP address for your services, which can then connect to each other. Not visible to other external users
- But performance might be balky: this isn't their most performant option. And they charge by the megabyte transferred over the links.





## **REMINDER: MESSAGE BUS (DDS) VS. MESSAGE QUEUE (KAFKA)**

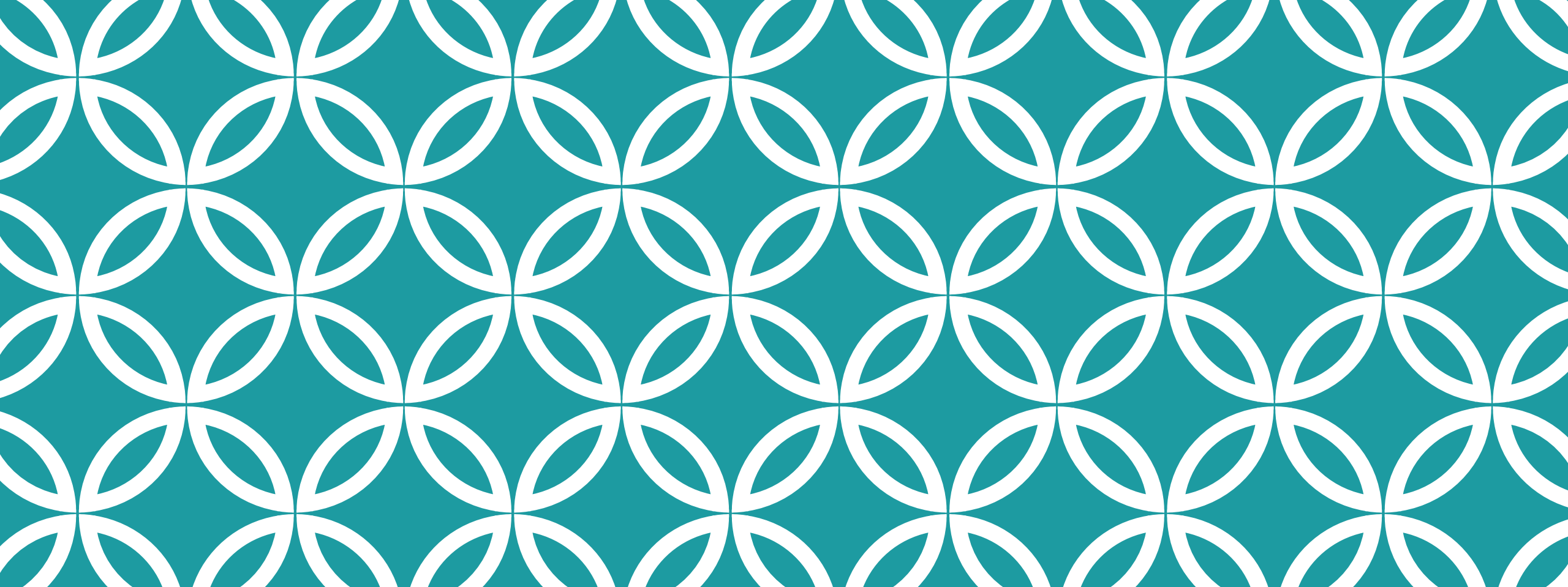
**A message queue like Kafka is a form of sharded key-value email system, for applications**

# QUEUING VERSUS MESSAGE BUS MODEL

We mentioned these earlier.

A message queue is a store-and-forward scheme, like email. Useful when doing batched processing (“give me all the emails from the past hour”)

A message bus is used when you want minimal latency on a per-event basis. Data is sent without any persistent storage, like a text message.



# LONG DELAYS

Heavy-tailed latency: Big issue

# HEAVY TAILED LATENCY

A big concern with georeplication is erratic delays.

Within an availability zone, the issue is minimal: the networks are short (maybe blocks, maybe a few miles), so latencies are tiny.

But with global WAN links, latencies can be huge and variation grows.

- Mean delay from New York to London: 90ms
- Mean delay from New York to Tokyo: 103ms



The number of events in the  $i$ 'th popularity bin is proportional to  $1/i^\alpha$

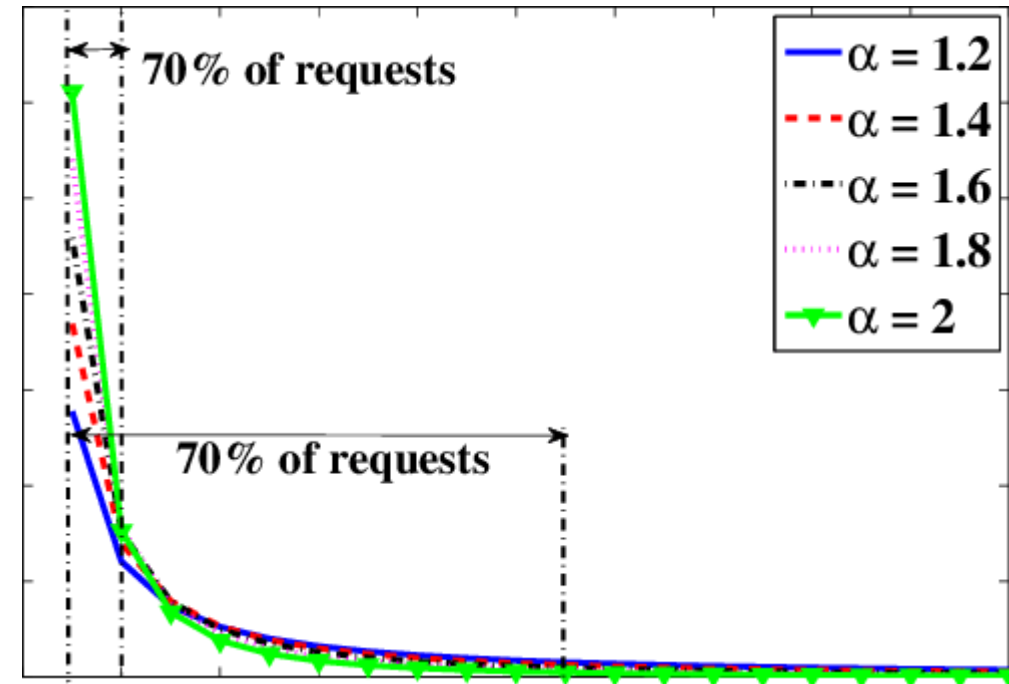
Zipf's law

# HEAVY TAILED LATENCY

Studies reveal “Zipf” latency distributions

Most traffic gets through very quickly

But some traffic takes extremely long.  
We say that these distributions are heavy-tailed if the “very slow” traffic adds up to a substantial portion of the total traffic.



# CAUSES FOR HEAVY TAILS?

Packet drops. When a WAN link gets noisy, it might suddenly drop a slew of packets, then recover. The effect is that some packets rush through, but others show up very late – TCP will need to delay the early ones to deliver in order.

Routing or link speed changes: At WAN scale these events are infrequent but they do happen and can cause a few seconds disruption.

Node crashes, especially in the machines handling the links or queuing logic

# ISSUE THIS RAISES

Should we want to wait for *all* sites to respond, or just a quorum?

Quorum: a subset of the sites large enough to include a majority. Any two quorums are certain to overlap. Call this  $Q$ : we require  $Q + Q > N$

- If task A updates a quorum and task B later reads the quorum, then B is certain to overlap with A on some server. So A “sees” B’s update.
- One can get fancier by introducing separate quorums for writes and for reads, such that  $Q_w + Q_w > N$ , and  $Q_w + Q_r > N$

# ISSUE THIS RAISES

Seemingly, a quorum is far better.

On the other hand, perhaps we would want all sites within an availability zone, but then wouldn't need to wait for geo-replicas to respond?

Leads to a three-level concept of Paxos stability.

- Locally stable in datacenter... Availability-zone stable... WAN stable



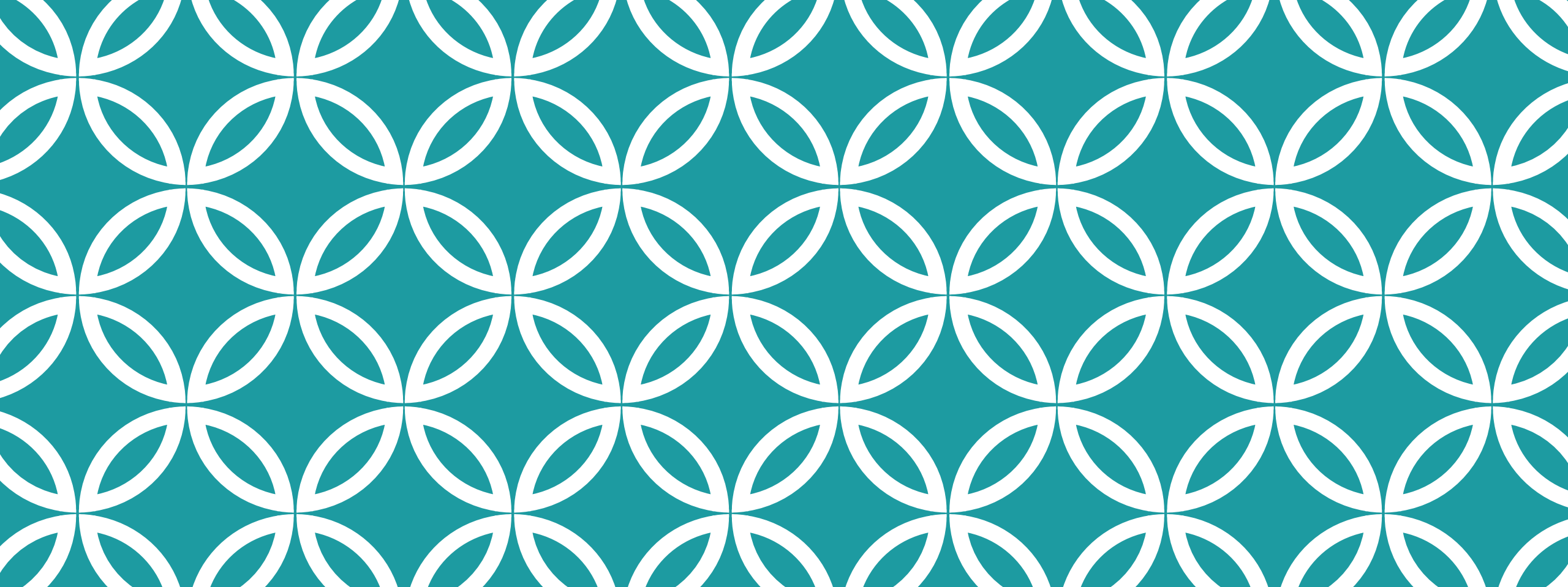
# WHAT HAVE EXPERIMENTS SHOWN?

Basically, Paxos performs poorly with high, erratic latencies.

A big issue is that delay isn't symmetric:

- The path from Zone A to Zone B might be slow
- Yet the path from Zone B to Zone A could be fast at that same instant.

So the outgoing proposals experience one set of delays, and replies from Paxos members experience different delays. You end up waiting “for everyone”



# **ATOMIC MULTICAST VIA SPANNER (GOOGLE IDEA)**

**Yet another way to create state  
machine replication solutions!**

# GOOGLE SPANNER: BACKGROUND ON CHUBBY

Google was one of the first to use a service built from Paxos in real datacenter settings. It was called Chubby, and was created by Mike Burrows with a Cornell PhD graduate, Tushar Chandra.

In a single data center, Chubby worked well, although it was dramatically slower than Derecho.

But in a WAN setting, Google struggled to try and build a global version of Chubby. Basically, they couldn't pull it off!

# GOOGLE SPANNER INNOVATION: TRUETIME

Google uses actual time as a way to build an asynchronous totally ordered data replication solution called Spanner.

Google **TrueTime** is a global time service that uses atomic clocks together with GPS-based clocks to issue timestamps with the following guarantee:

- For any two timestamps,  $T^0$  and  $T^1$
- If  $T^1$  was generated after generation of  $T^0$  finished,
- ... then  $T^0 < T^1$

# WHY SUCH A TORTURED DEFINITION?

With the guarantee they offer, if some operation B was timestamped T, and could possibly have seen the effects of operation A with timestamp T', then we can order B and A so that A occurs before B.

The full API is very elegant and simple: `TT.now()`, `TT.before()`, `TT.after()`. `TT.now()` returns a range from `TT.before()` to `TT.after()`.

The basic guarantee is that the true time is definitely in the range `TT.now()`. `TT.before()` is definitely a past time, and `TT.after()` is definitely in the future.

# IMPLEMENTING TRUETIME

Google starts with a basic observation:

- Suppose clocks are synchronized to precision  $\delta$
- It is 10:00.000 on my clock, and someone wants to run transaction X.
- What is the largest possible timestamp any zone could have issued?

My clock could be slow, and some other clock could be fast.

So the largest (worst case) possible will be  $T+2\delta$

# MINIMIZING $\delta$ IN TRUETIME

Google uses a mix of atomic clocks and GPS synchronization for this. They synchronize against the mix every 30s, then might drift in between.

GPS can be corrected for various factors, including Einstein's relativistic time dilation, and they take those steps.

In the end their value of  $\delta$  is quite small (0-6ms). So at actual time 10am, transaction X might get a `TT.after()` timestamp like **(10:00.006, zone-id, uid)**.

➤ The zone id and uid are to break ties.

# NOTE: TWO DISTINCT NOTIONS OF $<$

With  $T^0$  and  $T^1$  both true-time timestamps, we can look at

- $T^0 < T^1$ : This is evaluated as  $T^0.\text{after}() < T^1.\text{before}()$ . If  $T^0$  and  $T^1$  are concurrent, then neither  $T^0 < T^1$  nor  $T^0 > T^1$
- But each timestamp also has an actual time value. We can order timestamps by this value, e.g.  $T^0.\text{now}() < T^1.\text{now}()$ , even if the times are concurrent. This ordering is used for sorting.

The first of these genuinely implies that  $T^0$  happened before  $T^1$

The second is a total ordering but in fact  $T^0$  could be concurrent with  $T^1$



# USES OF EACH KIND OF ORDERING

Given a set of transactions, we can sort them by `tt.now()`. Spanner does this.

If applications hold “leases” (locks with timeouts) on resources, we can use the  $T^0 < T^1$  form of ordering to know, with certainty, when a lease has expired – and this works even without A being able to talk to B.

Another nice use case is for partitioning: If we track when processes last touched bases with a membership service, we can ensure that if B loses connectivity, it will shut down – and A can deduce this, too. Example: if B was in control of a drone, but loses contact to the main system, B shuts down and A can safely take over.

# SPANNER'S USE OF TRUE TIME

Spanner is a transactional database system (in fact using a key-value structure, but that won't matter today).

Any application can generate a timestamped Spanner transaction. It won't be executed right away – they use stored procedures triggered by a message (*the message holds arguments for the triggered procedure*).

These are relayed over the Google version message services. Their service ***delivers messages from Zone X to Zone Y in timestamp order.***

# HOW TO USE GOOGLE SPANNER

First, get a job at Google – this is an internal-only product!

But now you register your transactional logic as stored procedures, on all the Spanner instances.

When you wish to run a transaction you encode the arguments to your procedure into a message. Spanner will deliver it globally, causing your stored logic to run in a fixed order on all Google's datacenters worldwide.

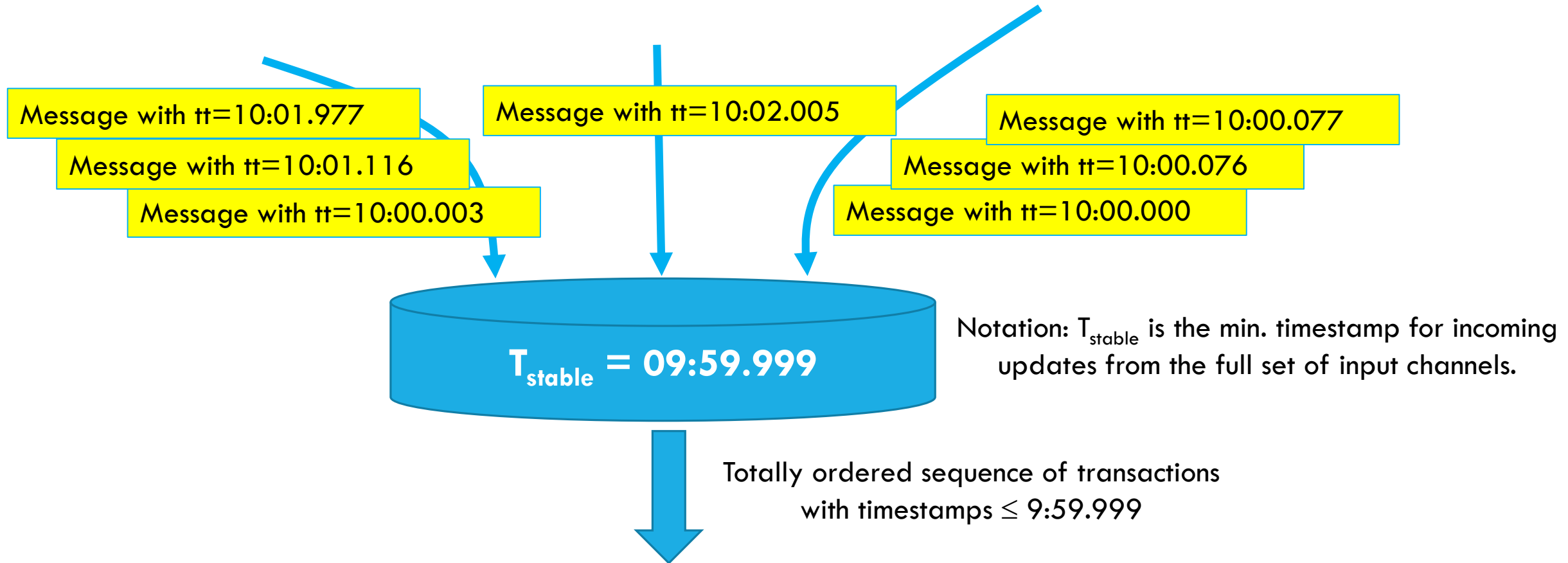
# LIFE OF A SPANNER INSTANCE

Spanner has connections to perhaps hundreds of remote zones.

Transactions flow in on these connections, and it stores them until every zone has sent some transaction with a timestamp of value  $T_{\text{stable}}$  or larger.

Then it can apply all transactions with timestamps  $\leq T_{\text{stable}}$ , in order.

# LIFE OF A SPANNER INSTANCE



# WHY DOES THIS WORK?

If Spanner has received messages with timestamps  $> T_{\text{stable}}$  from every zone, *it has seen every transaction with timestamp  $\leq T_{\text{stable}}$ !*

- This is because the connections deliver messages in order, by timestamp.
- If an earlier transaction were to show up, it would violate this rule.

So, if it now places those transactions into timestamp order (breaking ties by (zone-id, uid)), they can be applied to the global database in total order.

# HOW LONG COULD AN UPDATE BE DELAYED?

If incoming messages from any single datacenter pause, Spanner will have to wait for them (because they could have timestamps that order those messages “next”). So when an incoming stream pauses, Spanner pauses.

Eventually they declare the datacenter itself to be down – like with virtual synchrony membership management. This step is tricky and is the complex part of the Spanner technology. They need a form of atomic multicast so that if *any* datacenter executes an update, every datacenter does so.

# A NICE FEATURE OF SPANNER

Spanner can safely assume that the links to all its data centers are working, and it just waits to hear from all of them. If a data center is taken offline, Spanner is told, so then it won't wait for that link.

Thus Spanner can make ordering decisions “as soon as possible”.

This is leading to interesting work on using RAID-style erasure coding to overcome the delay impact of a slow TCP link side by side with other TCP links that aren't so slow...



# WHAT LIMITS SPANNER?



One limiting factor is that although the zone-to-zone data transfers run over large numbers of parallel TCP connections, the messages need to be put into order to obtain this “monotonicity” property.

A second limiting factor is Zipf-like latency with heavy tails: Spanner will often have to wait for “laggards”.

At global scale the effect can be significant (many seconds).

# PAXOS ALL OVER AGAIN?

With the classic implementation of Paxos, we have a back-and-forth interaction *that traverses every link several times in both directions*. Paxos experiences delays repeatedly.

With Google Spanner, there are global asynchronous flows but no back-and-forth coordination of this kind. So Spanner experiences delays once.

Intriguing observation: Derecho's version of Paxos "behaves" like Spanner.

# WHAT LIMITS SPANNER?

Google researchers report that the most frustrating issue is that a transaction cannot even be processed locally without waiting this way!

In some situations, a speculative result may be acceptable: “If there are no conflicts, my transaction would run and you would win the auction!”

But in other settings, consistency is absolutely needed, so there is no choice.

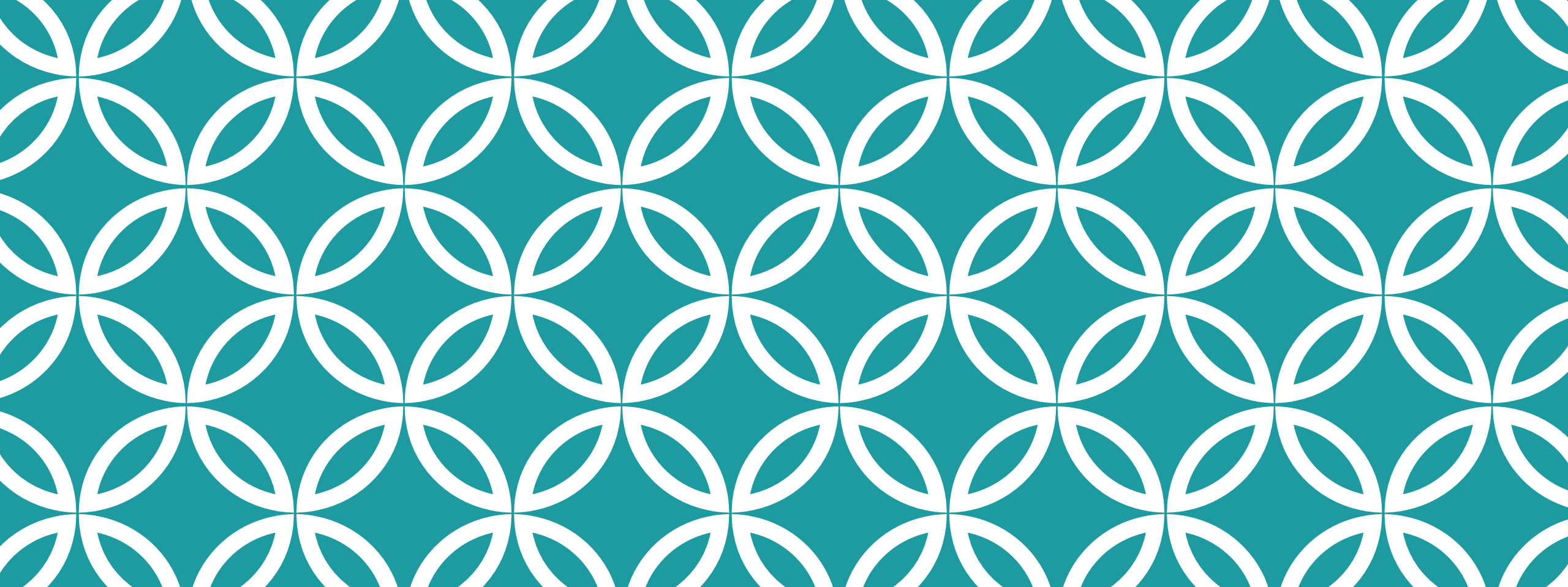
# TRICKS TO MINIMIZE IMPACT

If a zone hasn't been sending any transactions, it should "pause".

- Send a special "End of sequence" message
- Cease to send new transactions
- Other zones no longer need to wait.

Later, to resume, it will have to get permission to restart:

- Send a "resume request"
- Every other zone must acknowledge this before new transactions can start.



# OTHER OPTIONS

**Sometimes you just have to improvise**

# OTHER OPTIONS?

With a “primary zone” model, we can shard our database and assign each shard a primary owner (only the owner zone can update that shard).

Then you can make the rule that the primary owner can always do transactions on shards that it owns, without waiting.

But if *any* transaction would need to access shards for which it isn't primary, than *all* must use the Spanner ordering mechanism (you can't mix methods).

# WHAT LIMITS THE PRIMARY ZONE METHOD?

In some situations it is hard to know what shards a transaction would access before the transaction code actually executes.

For example, the keys a transaction will touch might be a function of the data it reads in some initial step. So until it executes, we don't know the key set, and can't know if those will all be local shards.

Spanner is really aimed at cases like these.

# MORE REMARKS ABOUT IP ADDRESSES

We have seen that with modern cloud systems, you could connect to Acme.com and be routed to Amazon.com instead.

In fact the modern cloud has considerable control over this, and you can influence that layer.

For each external client, when it initially connects, you can program the cloud DNS to route that specific request to a particular data center, and control whether or not the DNS record will be cached.

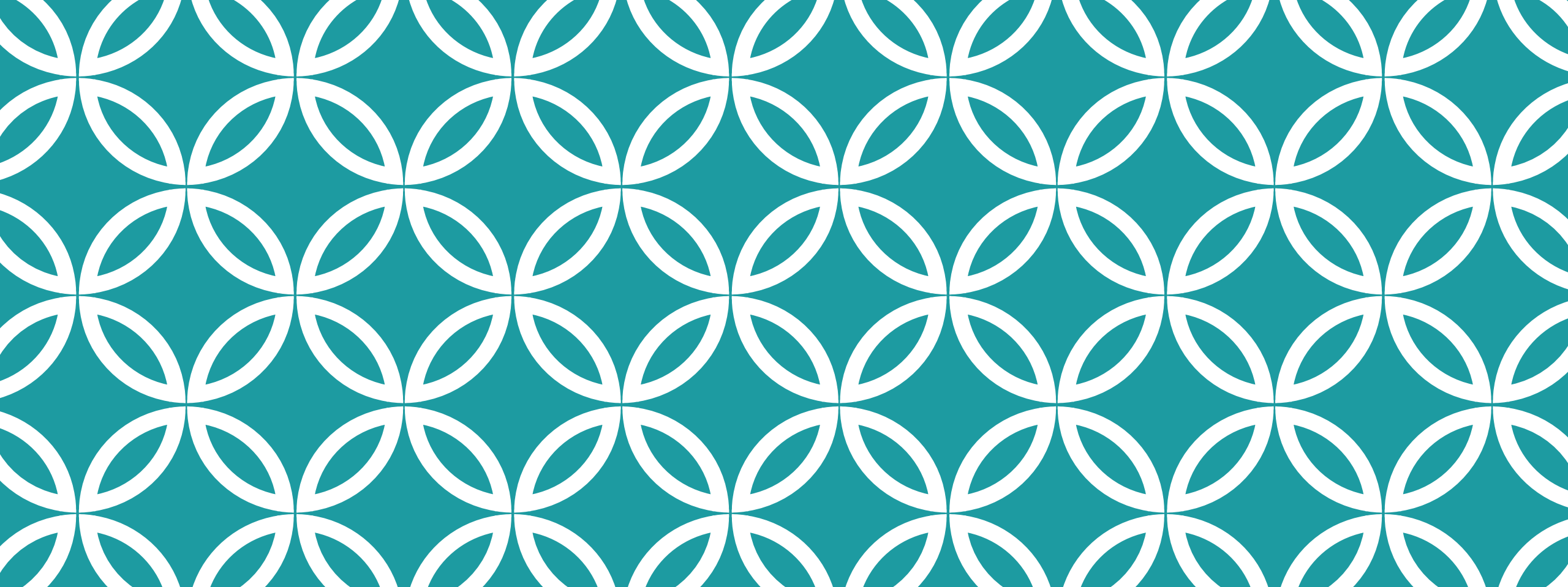


# MORE CONTROL

You can use information in a cookie left on the client to advise the cloud on which of your servers would be the best one to handle a connection.

And you can select between a wide range of elasticity and routing and load-balancing “recipes” that are offered by the vendor.

In the next generation of routers (SDN routers) you’ll be able to even provide packet inspection rules that could route based on values in specific fields of the incoming request.



# **5G: MORE OF THE SAME, PLUS SOME NEW ISSUES**

**In the US, 5G has been around  
the corner for almost a decade**

# MOBILITY AND GEOREPLICATION

With the advent of 5G, more and more IoT systems will include mobile clients in which the sensors are actually inside a vehicle like a car or plane.

5G hubs are very much like Azure IoT Edge, and in fact Azure IoT Edge may be a leading platform option for 5G services.

With mobile users, the user itself may have a dynamically changing IP address, and might be using multipath TCP to maintain connectivity.

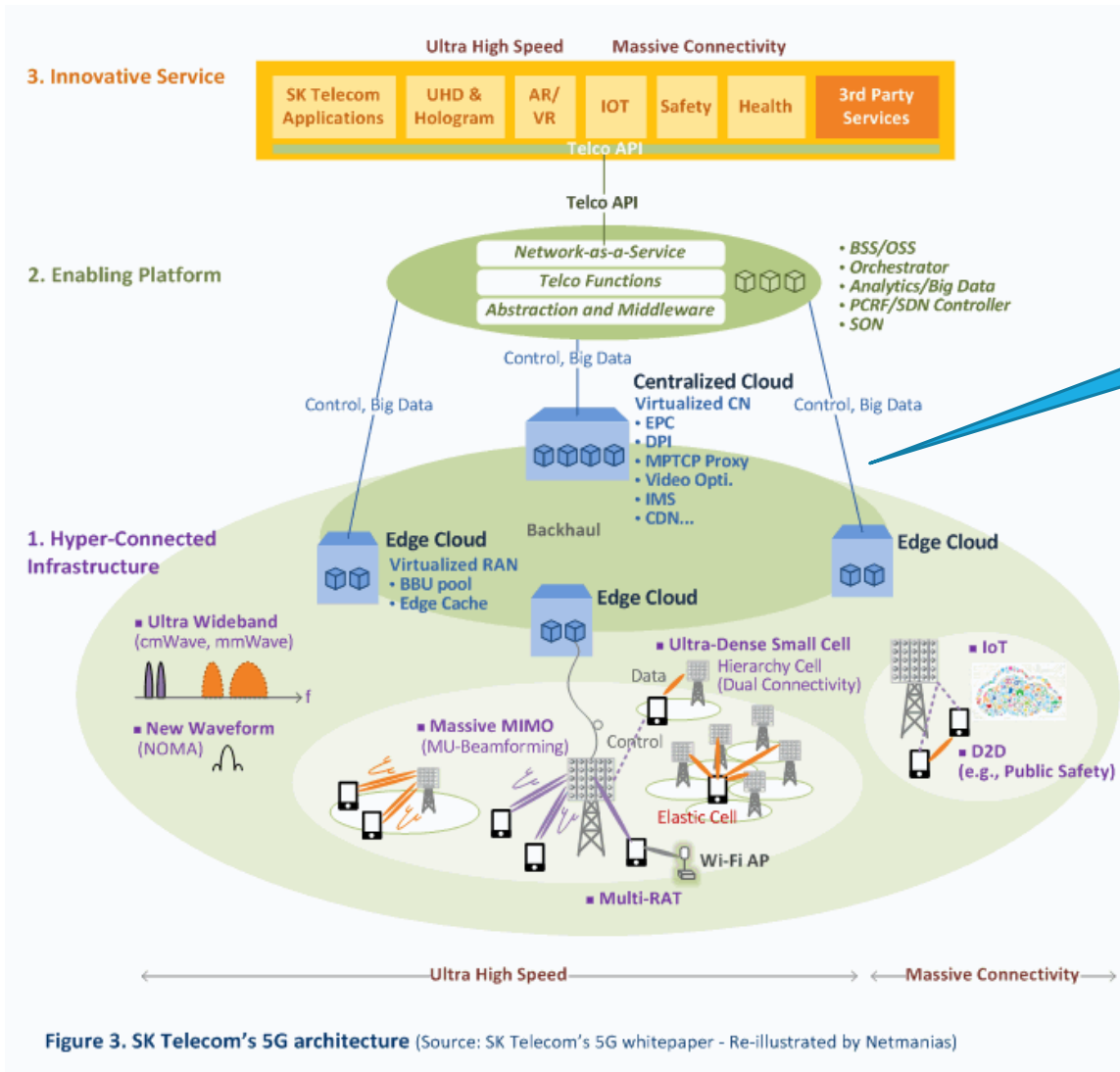
# 5G WILL ENABLE A NEW EDGE CLOUD!

Many experts believe that as these 5G point of presence systems roll out, they will be so interconnected at the edge that they will rapidly become a new form of cloud.

Just like the standard cloud, we'll be able to run apps on it!

The application will talk to the local 5G PoP, but those systems will need to use tools like messages queues or replication to share data.

# THE PATH TO 5G



They envision a service infrastructure of their own living on Azure IoT Edge or similar!

*Is 5G just an IoT cloud integrated more closely to telecommunications?*

Figure 3. SK Telecom's 5G architecture (Source: SK Telecom's 5G whitepaper - Re-illustrated by Netmanias)

# 5G WILL ALSO ENABLE “QoS OPTIONS”

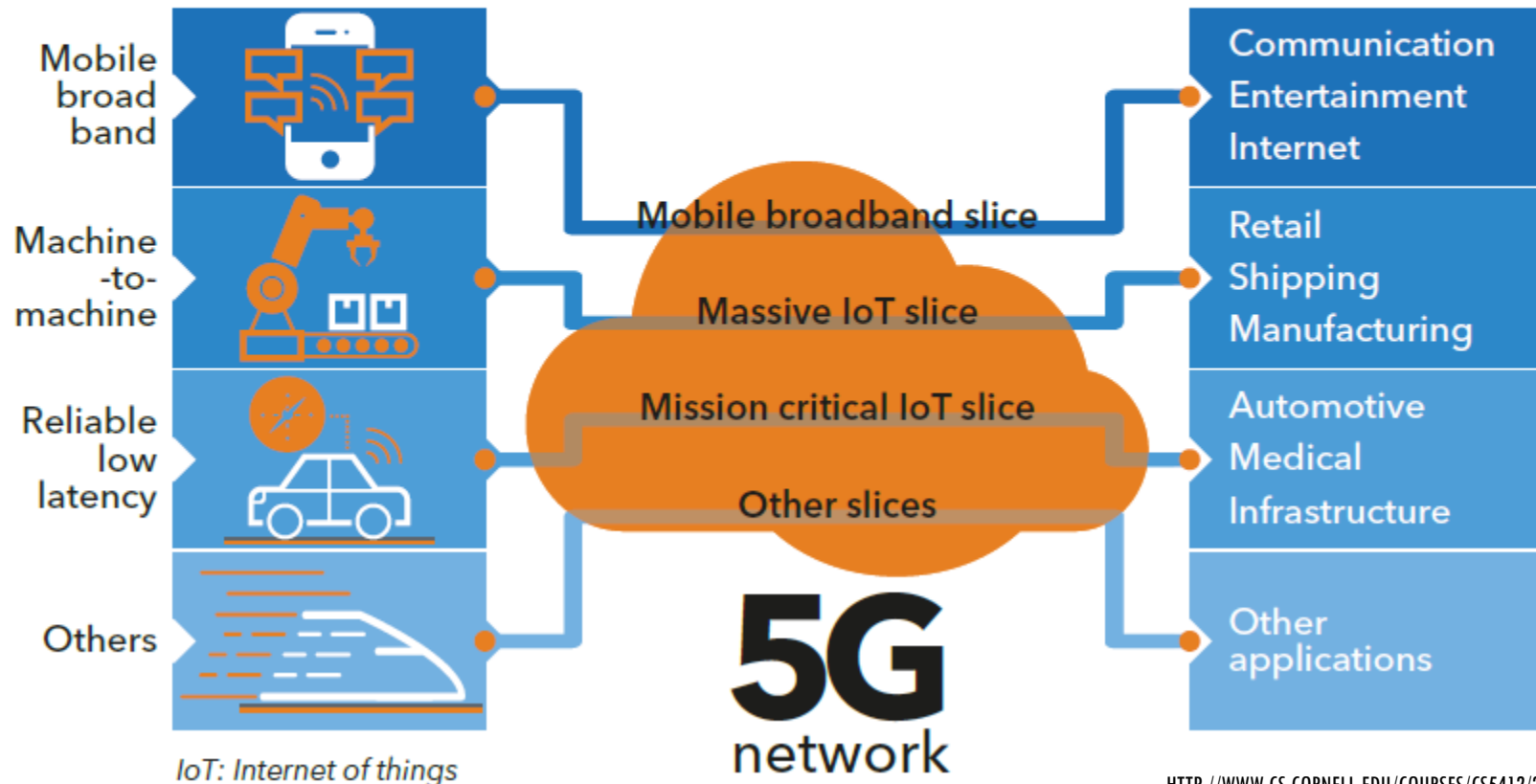


4G networks do not enable the range of services that the future requires. 5G will be faster and more flexible.

**4G**  
network

## 5G network slicing

5G network slicing enables service providers to build virtual end-to-end networks tailored to application requirements.



# KEY INSIGHTS

5G is just another georeplication scenario, and the cloud can evolve to support it while retaining its main features!

Main areas where change will occur:

- Very large number of cellular regions (“5G edge cloudlets”)
- Mobility  $\Rightarrow$  Growing need to move data in anticipation of demand
- Opportunity create a new hosting infrastructure for 5G “apps”

# SUMMARY



“...he doth bestride the narrow world/Like a Colossus...”

GeoReplication is best viewed as having two scales:

- Availability Zone: Just neighboring data centers. With some cost, you can use TCP and build “normal” protocols. Derecho should work this way.
- True global scope: Here, techniques like Google Spanner are best. Researchers have experimented with Paxos at global scope, but it performs poorly due to high-latency links.
- The notion of ***wide-area stability*** is useful and might be worth using in other contexts.
- 5G mobility will introduce an additional layer of services