

Distributed Systems

Systems are not Applications

connectivity → the network matters

- reliability data can be lost when sent over the wire
- latency it takes time to cross the network
- bandwidth network congestion may interfere
- security bad actors on the Net
- topologies resources can go down or move to another subnet
- administration updates, dependencies, ...
- transport costs serialization / deserialization
- diversity lots of different systems, languages, architectures, ... working together

"If you have to cross the network,
take all the data you might need
with you."

Projects are a poor model for software development,
long-lived products are better

acknowledges, retries, caches, ...
wait, expire, ...
discarded packets, time-critical data, ...
threat model analysis, risk analysis
DNS, dynamic routing, ...
logging, monitoring, ...
auditing, journaling, ...

- (several steps), non-atomic operations
- (continuous) evolution
- (decentralized), distributed logic

Coupling

- platform A works with B
 - time processing time of B affects that of A
 - spatial where is B?

Decoupling

use text based standard representations and protocols

JSON, XML, ... HTTP,
 SNMP

use messages , publish/subscribe, ... WIDP,

HTTP,
SMTP

100 P.

delegate connections to lower layers, load balancers,
proxies, ...,
queues, ...

Messaging

Services coupled to messages, not to each other

- asynchronous messaging → one-way, fire & forget messages
 - strongly-typed messages
 - represent methods as messages
 - include id, version, timestamp in each message → ordering issues
 - include relating info, in order to correlate... what had happened?

Processes

process: a set of activities that are performed in a certain sequence, as a result of internal or external triggers

long running processes : are stateful

use Sagas to manage the flow

use Adapters to manage integrations

① The network is reliable?

it works most of the time ...
... except when it don't work!

- our code must be aware of timeouts, broken links, congestion, ... on the sending way and on the returning way
- in asynchronous scenarios we can try to retry the operation or to recover by other means in the background
- but in synchronous scenarios we must take into account that the user is waiting for the response (response time matters)
(in Request/Response scenarios)
- in distributed systems, with many actors involved, with several steps, there are too many edge cases when we try to recover from a failure in one of the steps

transactions integrity
is tricky in distributed systems

can we build a reliable system on top of a fundamentally unreliable network?

The response to this question is "yes", but we are to going to build it differently than the traditional Request/Response pattern we are used to use on desktop applications

better to use a MQ
(Message Queue),

use a reliable messaging infrastructure

and adapt ^(*) the architecture of the application / system to work with this paradigm

(*) we lose Request/Response kind of operation

② Latency is not a problem ($= \Theta$) ?

time to cross the network
in one direction

- is small for a LAN
- can be large for a WAN

it is many orders of magnitude
(sec to year)
 10^{-9} to 10^6
(ms to ms)
 10^{-3} to 10^3

slower than in-memory access

Don't do use remote objects!
avoid them as much as possible

don't chat back and forth
with a remote object

Don't cross the network
if you don't have to.

Use DTOs, and take all data
you need once, and cache it
locally for further use, whenever
it is feasible to do so

send & forget style
of data interchanges

Inter-object chit-chat shouldn't
cross the network, do it locally,

but beware with how this affects to
the bandwidth limitations
on your networks

"If you have to cross the network,
take all the data you might need
with you."

DTO (Data Transfer Object)

See next page back to
③ back to

we need to have more than one
domain model to resolve forces
of bandwidth (lazy load) and latency (eagerly
fetch)

③ Bandwidth is not a problem (oo) ?

how much data you can move through the network, is going to have severe ramifications on the performance of your application

$$1 \text{ gigabit} = 125 \text{ Mbytes/sec}$$

Tcp/ip headers and coordination take half of it

$$\approx 64 \text{ Mbytes/sec}$$

Data serialization format can take another half

$$\approx 32 \text{ Mbytes/sec}$$

If you have network congestion problems, putting more servers will only worsen them

You need a more efficient way to access to the data
(or a network with more bandwidth ;-)

or a staggered access to the data

or several, different, independent, routes to the data

don't underestimate the bandwidth of a truck full of hard drives going by highway (moving bytes across a network is much much slower)

Although bandwidth keeps growing,
the amount of data to transfer grows faster.

2000

1 gigabit

1 core

2 GB RAM

512 GB disk

2024

10 gigabit

8 cores

64 GB RAM

2 TB disk

congestion problems lead to latency problems, and latency problems lead to timeout problems

problems \Rightarrow retries \Rightarrow more problems



⑦ The network is secure?

The only machine/system that is secure
is one that is disconnected from anything else...
turned off and buried in 3m of concrete

we, human beings are the weakest link in any
security chain. Social engineering is the
best tool of a hacker.

You can't be 100% safe from everything

So, perform a threat model analysis
and balance costs against risks

- what is the cost of protecting something?

- what is the risk (damage) if this something is compromised?

- or what is the advantage (benefit) someone will gain if compromises this something?

But, in first place... is this way of connection necessary?
(attack vector) (convenient)

The best hacker proof door is the door that doesn't exist.

nowadays, with practically all things running
in the cloud (Internet), and
or connected to the
the systems becoming more and more
complex and interconnected,
the ability to build systems that
are secure is more difficult

⑤ The network topology don't change ?

Programming models that can run into issues in environments where the machines and networks are more volatile than you assume.

for example - WCF callbacks contracts

- config file values for machine ids
- clients or servers that disconnect and disappear from the net after subscribing or register themselves with a service or pool
- etc

Be extra aware of fast-dirty-deco code that nevertheless ends-up deployed in production.

- Don't hard-code addresses
- Consider using resilient protocols (multicast)
- (Auto) Discovering can be cool, but it is hard to get it right.

Test your system regularly to know how it performs under changing topologies (for example one or several instances of a server goes down, or reboots at different times, or try to reboot before some other service is available, or...)

chaos monkey tests

⑥ The admin will know what to do ?

The all-knowing admin that maintains flawlessly the system is only possible in small networks. And if so, only until he/she got promoted and leaves...

Will the system survive to changes in working procedures, upgrading policies, new versions, new rules, ...

Make clear any dependency between different parts of the system. Or better, make each part as independent possible

The more configuration options we create, the less likely that anyone knows what a given set of specific settings is going to do

Do invest in some kind of automated configuration management

Take time to slow-down and document

Take snapshots of well-known points and be able to return to them if something goes wrong

Do invest in some kind of logging/monitoring management

Know how looks your system when it is performing well, to detect problems when they start to show-up

The more moving parts you get, the more likely that some go off of sync with each other and something not going to work

The more and more code we write, the more changes we make, bigger the chances of failure at the next upgrade/deployment to production

CI/CD

and automate as much as possible of the pipeline

⑦ Transport cost is not a problem?

Serialization before crossing the network takes time^(*). And the same deserialization on the other side.

(*) it takes a bunch of CPU time

And usually goes undiced, behind the frameworks we are using in our code.

Each architecture is better suited for a certain kind of application.
Is our work to find the most ^{optimal} efficient architecture style given our requirements and constraints.

All is a trade-off balancing the cost perspective, the complexity perspective, the maintainability perspective, the scalability perspective, ...

Physical architecture

It is cheaper/faster as more readability is our system. (for example, reading all from the same machine or ~~all~~ is cheaper/faster than retrieve all over the network from different machines)

Logical architecture

It is more flexible/scalable as more modular and more independent those modules are.

connections costs tend to grow exponentially as we go to larger and larger links.

Some architecture can be feasible inside one unique datacenter, but can be infeasible if distributed all around the world.

⑧ The network is homogenous?
 (intersperelle)

(2005)
 Some years ago, there was only the interoperability between major languages
 (aka .NET ↔ Java)

But nowadays⁽²⁰²⁰⁾ the decentralization trend forces us to manage between different languages, different database styles, different API styles, different encoding formats

go	Relational	REST	JSON
Ruby	NoSQL	RPC	XML
JavaScript	Graph	Webhooks	YAML
:	:	:	:

on different systems hooked to work together

Semantic interoperability will always be hard, budget for it.

for example
 We don't know what to put in place of missing mandatory fields for a system that requires them when the data reaches us from another system that doesn't require them.

be aware that almost nothing works out-of-the-box in the interoperability world

You always need to make adjustments, from small tweaks to big integration rewrites

- Take into account that things can fail on a network or change
- Watch for latency and bandwidth issues
- watch for security issues
- Automate as much as possible - configuration
and don't
 - integration
 - deployment CI/CD
 - monitorization
- Avoid as much as possible interoperability scenarios.
- Avoid as much as possible network communications specially the long-distance ones.
- Allow as much as possible logical modularity and minimize as much as possible the dependencies between modules

The system is atomic ?

a program that had evolved into
BBOM (Big Ball of Mud) → a logical monolithic monster
can be worse → if it is physical distributed!
(different parts of the logical monolith
running in different machines)

scalability { vertical : bigger, more powerfull, worker
horizontal : more workers

Beware of false "services". A service is an indepedent part
of our system that performs a substantial work
on its own, without need of other parts to do so.

If a "service" needs collaboration from a bunch of
other "services" to get its work done, all them are
one logical service (one logical monolith)

logical monolith : all parts of the code
tightly coupled with each other

logical distributed : code organized in
well defined modules
with clear interfaces
between each other

physical monolith : all parts running
in the same machine

physical distributed : different parts running
in different machines
(linked together trough)
(network connections)

A logical distributed program can
be deployed on a physical monolith manner
or on a physical distributed manner.

But a logical monolith program can
only be deployed in a physical monolith manner
(can only scale vertically)
(try to scale horizontally = distributed = leads)
(to worsen performance = network penalty =

The system is atomic (cont)

How do a program turn into a BBoM ?

Some forces that can lead to this outcome:

- integration through a unique database
(all code is tied to that database)
(coupled)
- patch other supplementary data sources
on top of that unique database
(usually because we cannot change the schema
for fear to broke something else)

If you are saying "we don't do this...",
"we must fix this later... when we got time..."

You are onto the pernicious path toward a BBoM

Please, stop! and take time to refactor
the system to maintain it's architectural
stability

A monolithic data architecture cannot scale horizontally!

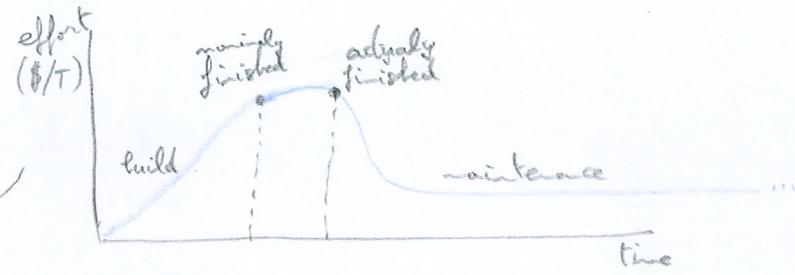
Horizontal scalability is only possible if
our system is designed to allow it.

Taking into account the logical boundaries (if they existed, of course)
in our model, and leveraging the system
to take advantage of them.

It is not an easy task!

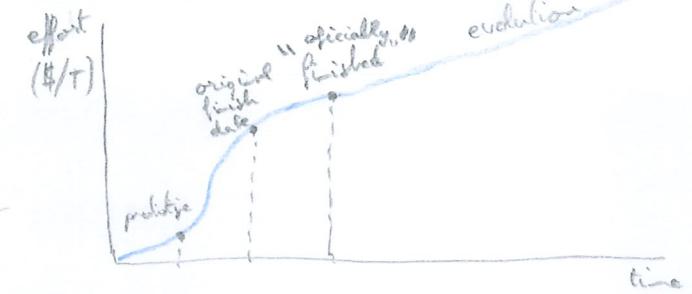
Independent services can run on independent machines. We can raise several nodes with a service
as we need more power for this service

The system is finished ?



There is not "maintenance" phase in a software project.

A software system is never "finished", it is more like a long-lived product in continuous construction/evolution.

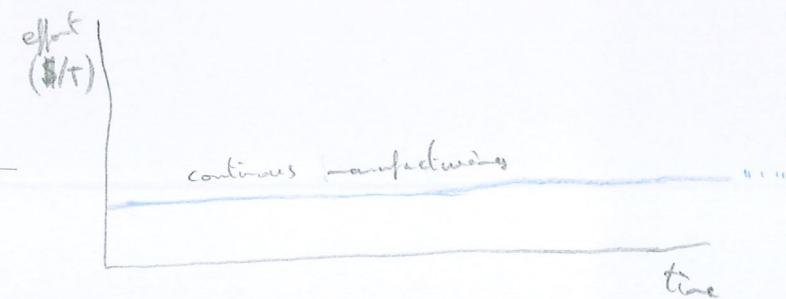


In other kind of projects and products,
the kind of activities performed and skillsets
needed to perform are quite different on
design/build phase than in maintain phase.

But in software systems they are the same,
always design/build ones... all lifetime of
the system.

To talk to the business people : **software is a product
not a project**

Launch a product takes care of the long-term ^(viability) manufacturability
of this product



A better development process ?

When asking for new software
Business people don't give requirements.
 they give workarounds that they thought will allow
 the system to do something they ^{want} need to do.

IT Business Analyst must have sufficient
 "gravitas" or "political capital" to be able
 to say "lets have a conversation" to discover
 and analyse what is actually needed,
 why it is needed and how to achieve it.

IT Architect must have sufficient
 experience to be able to give
 some estimates, try some PoCs (Proof of Concept),
 and refine the estimates, until all
 stakeholders agreed to some requirements
 and associated costs. (requirements are an investment decision)

"Given a well-formed team of N persons,
 that is not working on anything else,

I'm $C\%$ confident work will take between T_1 and T_2

the Holy Grail ?
collaborative communication between
 all people involved,
 business ad IT,
 until they understand each other
 and agree on something
 realistic enough to do

months
millions €
...

well-formed
 A team is a group of people
 that had worked together
 in something similar.
 It is not a bunch of people
 just assembled for this tasks

Business logic can and should be centralized ?

When we implement something in three or five places all around the system ..

What happens when business rules change ?

we risk to forget to change one of these places !

But, if it is implemented only once, in a given component. Then we have other two or four components that depend on this centralized component.

What happens when business rules change, but they change differently for each of those components ?

we need to walk all those components to broke its dependency and rewrite for the new rules

Architecting requires more than one point of view to grasp all its complexity; all of its different aspects of a software system

the pull request or commit that implements/changes the feature
You can tag source code by feature implemented, in a manner that enables you to find later by feature, even if its in multiple files/places. So you don't need to implement it on a single place, avoiding dependencies.

has references/tags
to all places
this feature touches

Duplication isn't the devil!

in more loosely coupled codebases is less likely that a change will going to break something elsewhere

DRY $\xrightarrow{\text{has evolved to}}$ don't be DRY,
we WET

Don't
Repeat
Yourself

Write
Easily
Twice

Coupling between two parts of a system:

- **Afferent** (incoming) : who depends on me? → to many signals (but a change in this component can potentially break a lot of things elsewhere)
- **Efferent** (outgoing) : what are the things I depend on? → too many puts in doubt that this component fulfills only a single responsibility
- **Platform** : specific or proprietary libraries, protocols, ... that are only available on one platform.
- **Temporal** : processing time of one component affects processing time of other component
- **Spatial** : only one specific machine can handle the ^{calls} ~~process~~ function once it is initiated on that machine

Take into account when business requirements demands consistency

Beware of shared resources (DB, API, ...) that hide coupling between components that share them

To minimize platform coupling, we can use

(and improve interoperability)

don't reinvent the wheel!

- text based representations of data (xml, json, ...)
- standards based transfer protocols (HTTP, SMTP, UDP, ...)
- SOAP services, REST APIs, ...

To minimize temporal coupling,

To minimize spatial coupling

- first be sure that the process is not temporally coupled
(for ex: if component A can't proceed until component B provides data)
(cannot do any sensible work)

- avoid repetitive task for some kind of data, cache it when you receive it in the first ask

If you can use stale data

if you publish data,
state its validity

- publisher / subscribers

↳ only one logical publisher for a given kind of event

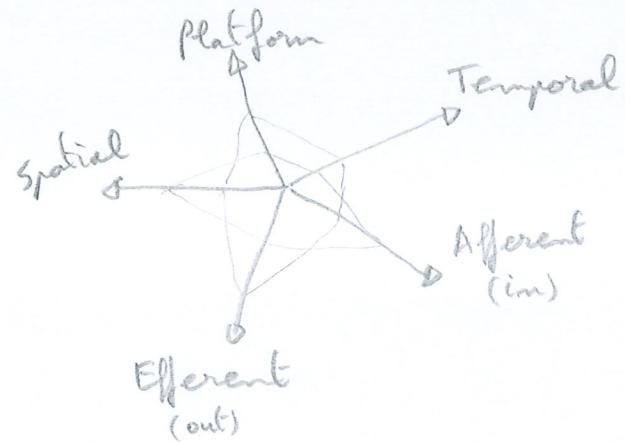
don't publish commands
publish only facts

Don't apply technical rules based on metrics. Persons have a tendency to hide things to make the metrics look good, instead of actually solve the problem that the rules try to catch

measuring degree of coupling is a very subjective matter

- what kind of coupling?
- how much of each kind?

Coupling is a function of 5 different dimensions:



Each technological or architectural decision we take, position us on different parts of this diagram.

Coupling is not loose or tight,
the question we must ask ourselves is

"is this coupling situation reasonable
in this context?"

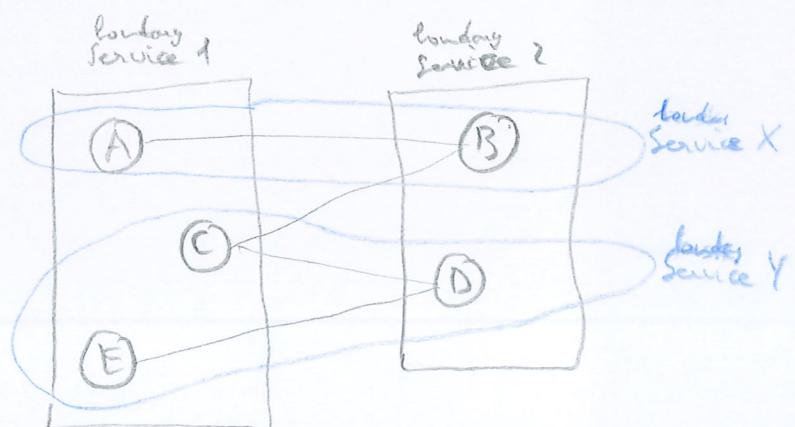
There is no "silver-bullet" for all,

There is no "free-lunch".

All has trade-offs and equilibries

Clear boundaries with
clear relationships between them.

Beware of tasks that cross boundaries at several points. They are a warning signal that some boundaries need rethinking.



Asynchronous messaging \Rightarrow { one-way
fire & forget } messages

- (A) launch something
- (B) continue doing sensible work
- (C) when response arrives (latter)
or when it is convenient for us,
continue with the work that
needs sending from (A)

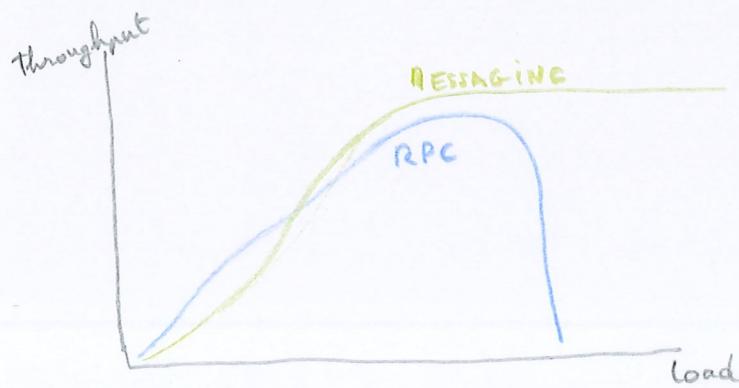
If our task at hand doesn't have (B)
(we must wait hand-on-hand until
doing nothing useful until (C))

The task is synchronous and not asynchronous!
and blocking

Each message has an id

Store & Forward adds reliability

RPC vs Messaging



RPC (Remote Procedure Call)
has blocking calls and
have no time to free memory
With MESSAGING, threads are
independent

The remote server starts rejecting requests ↗
and the system goes down temporarily
to process much latter ↗
The Queue can cache storing messages to disk ↗

Service Interfaces

- public functions in a library or similar
- public remote calls to an API
- queues or topics where to send messages

- or
- command-centric
 - document-centric (resource)

Messaging interfaces

- messages of certain types, with a certain bunch of attributes / properties / fields
(they are as DTOs)
- handlers for each type of message

In order to share message definitions,
distribute
[package managers] become an important tool
to manage artifacts in a distributed
system. (packaging and versioning)

The cross-cutting stuff (security, logging, ...)
usually goes transitively on the [headers]
of the transport system and not in the
body (the payload) of the messages

On top of the messaging delivery system,
you articulate the business logic architecture:

- Event-driven
- or
- CQRS (Command-Query Responsibility Segregation)
- or
- Microservices
- or
- etc

also
DDD (Domain-Driven Design)

Fault tolerance

- when servers crash
- when deadlocks happen
- when network are saturated/down
- when databases are saturated/down
- ...

If failed, you need to log/store the data or you will lose that data

direct calls to APIs need retry, acknowledge, timeouts, ...

message systems simply stockpile the messages to be processed latter, when the database server is available again

store & forward mechanisms receive and store messages in a place to be sent to another place when possible (same machine) (across the network)

^{log}
Write a message in a queue is a small / fast / reliable operation.

No business logic, database transactions, ... are involved!

The messages stay in the queue until the system can process it.

There is a record of the data / action waiting to be performed!

It is easy to transfer messages from one machine to another.

or from one queue (for exple: pending) to another (for exple: done)

And also ^{servers} a log that can be monitored or audited

After some number of retries, if message processing fails, usually, it is tagged as ''poison message''

When the problem is solved, you can retry the failed messages. → you don't lose data!

and put in an error queue for "failed messages"
(adding the appropriate information about the failure reason)
(with the appropriate notification to someone in charge)

→ you have rich monitoring data to help troubleshooting!

Auditing / Journaling

Messages can be saved to a long-term storage after they are processed.

A message queuing system is also a "running log" of what had happened and what is happening in the system

So, you can have a central log of everything that happened.

But it can be difficult to interpret by itself.

"Related to" fields are very useful when tracing what had happened. Helping to correlate one message to another ^{in order} to trace the flow

usually they are based on "Message id" fields

of messages related to a given action

"Timestamp" fields are very useful also. Helping to correlate temporally the flow ~~from~~ different queues from different sources.

Calling external web services

If one step fails?

and you need to **rollback!**
the steps already done?

It is easier to do with
"user generated id's"
(for example UUIDs).

It is difficult to do with

"database generated id's" → first steps have a different id,
or nothing can be done until
the database provides an id
for the first step of the process

It is impossible to do if
the ids must be correlative
and there can not be gaps
between ids. (When a process
fails, other processes had advanced
(the id counter, how can you
redesign the failed id, ... !?!) ...)

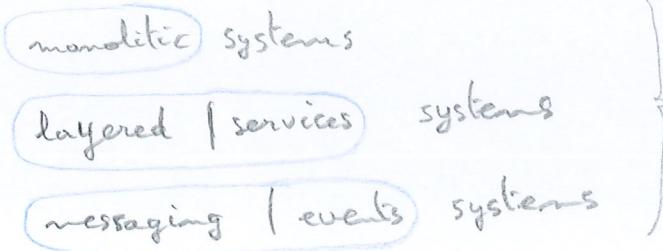
Weird things start happening
(eventual inconsistency) when you
mix different paradigms
(messing, RPC, REST, ...) in
the same endpoint processing
code.

DTC!

Distributed
Transaction
Controller

Distributed transactions
are very very tricky,
very very difficult to
do well. Avoid them
whenever it is feasible
to avoid.

ACID consistency
is very difficult
to attain when
working with
distributed systems



different kinds of architectures (hammers, screwdrivers,..)
for different kinds of needs (nails, screws,..)

organizational inertia is not on your side,
organizations want to keep doing what they always been doing.
Any type of change you are proposing, you have to push it up-hill.

How to reduce ^{do-now} real-time chunk of process (synchronous part)
and expand do-latter chunk of process (asynchronous part)
so you can cope with big load peaks

- Save to database (synchronous part) and batch process (asynchronous part)
- Save messages to a queue (synchronous part) and read from the queue (asynchronous part)

The async-effect is that is cost effective,
not the messaging-effect
(messaging systems are intrinsically async)

With messaging systems, we can separate the do-latter part on several machines.
With batch processes, we need to use the same database; and it can become a bottleneck.

But, the several machines ultimately also need to use the same database.
if you don't include all the data you need to process it in each message.

- don't do value judgments
like "X sues", because they
are easily counter with
"Y sues", ... ".sues", ...
all all sues at the end
- don't appeal to figures of authority, because
anyone of us has his/her
loved authority figures
(instead, we can try
to reach to agree on
who can be a figure
of authority preference
on which all of us agree)

- ground on technical arguments explaining
benefits for the situation
we are talking about
in each moment.

- instead of "- is more -"
try to explain "why - is
more - in this situation",
"how is - more - in this
situation"

Messaging systems forces us to think upfront
on separation of concerns, temporal and spatial coupling, ...
and other aspects of distributed systems.

But this guarantees nothing. We can do a "messaging mess"
as we can do a "call services mess" or a "call functions mess"

A poor designed architecture is a poor designed architecture,
regardless of the technology it is built in.

Start small, break things, do mistakes, on and learn.

Start with applications that are highly async : send email, generate PDFs, ...
- that have basic batch processing processes

little by little - developers are gaining experience with messaging technology
business are gaining confidence with the new technology
- and system admins

The potential benefits only realize themselves if we have enough developers that have gone through enough of the mistakes with the new technology to be able to figure out what is the right way to apply it.

The issue do not reside into the technology, it is into the architectural mentality the developers had. You can do right monoliths right distributed with practically all any technology.
(And so for bad monoliths, bad distributed)

Distributed Communications Patterns

How to manage the state of long-running distributed communication?

It is not about time, is about complexity of state:
 - regular processes → basic process control: if/then
 - longrunning processes →: **state machine**

In long running processes multiple external events/triggers/calls are handled by the same process instance → so, it is **stateful**

multiple steps that change state, until we arrive at the final state → endprocess

each step non deterministic, the process must wait until the step is completed; and then, proceeded with the next step that leads to the next related state

Message Handlers → process a message, and done! (or failed...)

Saga Handlers → a message triggers the process, the saga ^{manages} has an internal state, and works until the process arrives to its final state (or until it fails...)

internal state at each step must also be persisted

a consistent state that can be persisted

↓
and we must rollback the side effects any of the performed work steps had caused

avoid as much as possible distributed transactions !!
(distributed long running processes)

! this is the difficult part of my transactional process

! and more difficult in distributed transactions

to query directly
In a saga , avoid to alter directly
the master data during processing each step.

A saga is about to process all the steps,
persisting its internal state at each one,
until it completes the process and
saves the final result. → usually not directly , but by invoking some other services

If the processing fails ; nothing is altered.

And we can retry the saga .

Rollbacks can be performed :

- reverting the changes
- or
- making new compensatory changes

or sending messages to tell others to do something

The principal responsibility of a Saga
is the orchestration of the work to be done .

In multipart business cases

~~If the business case is complex enough , sometimes~~
we cannot solve/orchestrate them by technology only.

There^{are} also involved business decisions/contracts about
how each part must behave on different situations
and what each part can (or not) expect to do/receive.
(level of service agreements)

doing queries as a part of commands
over history in a highly consistent fashion

→ the "Holy Grail" for distributed transactions

In order to have very encapsulated
objects that do not depend on others.

Do as few request/response outside it
as possible, preferably none.

But for that, this object has to have
inside it all data it needs to work.

The easy part is using the building blocks.

The hard part is analyzing the
business processes to identify
what the steps should be.

And what the boundaries are for each component

In multistep processes external services
When interacting with 3rd party systems
 legacy systems,

- use a saga to manage the flow
- use a separate adapter for each integration