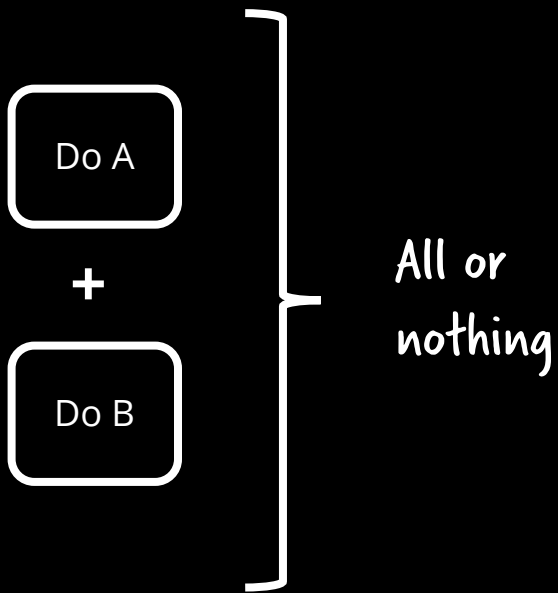


Lost in transaction?

Strategies to deal with
(in)consistency in distributed systems

@berndruecker





once upon a time:

```
try {  
    tx.begin();  
    doA();  
    doB();  
    tx.commit();  
} catch (Exception e) {  
    tx.rollback();  
}
```

or simply:

```
@Transactional  
public void createCustomer(Customer cust) {  
    // ...  
}
```

Atomicity

Consistency

Isolation

Durability

Distributed systems



Distributed systems



WIKIPEDIA
The Free Encyclopedia

[Main page](#)
[Contents](#)
[Featured content](#)
[Current events](#)
[Random article](#)
[Donate to Wikipedia](#)
[Wikipedia store](#)

Interaction

[Help](#)
[About Wikipedia](#)
[Community portal](#)
[Recent changes](#)
[Contact page](#)

Tools

[What links here](#)
[Related changes](#)
[Upload file](#)

Article [Talk](#)

Fallacies of distributed computing

From Wikipedia, the free encyclopedia

The **fallacies of distributed computing** are a set of assertions made by [L Peter Deutsch](#) and others at [Sun Microsystems](#)

Contents [\[hide\]](#)

- [1 The fallacies](#)
- [2 The effects of the fallacies](#)
- [3 History](#)
- [4 See also](#)
- [5 References](#)
- [6 External links](#)

The fallacies [\[edit\]](#)

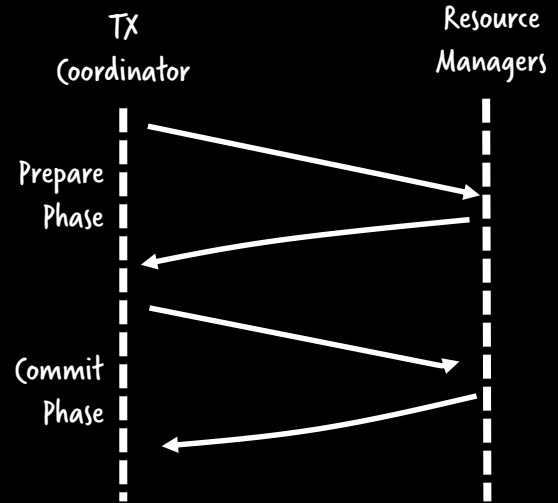
The fallacies are:^[1]

- [1. The network is reliable.](#)
- [2. Latency is zero.](#)
- [3. Bandwidth is infinite.](#)

Distributed systems



But there is two-phase commit (XA)!!





Pat Helland

Distributed Systems Guru
Worked at Amazon,
Microsoft & Salesforce

Life beyond Distributed Transactions: an Apostate's Opinion

Position Paper

Pat Helland

Amazon.Com
705 Fifth Ave South
Seattle, WA 98104
USA

PHelland@Amazon.com

The positions expressed in this paper are personal opinions and do not in any way reflect the positions of my employer Amazon.com.

ABSTRACT

Many decades of work have been invested in the area of distributed transactions including protocols such as 2PC, Paxos, and various approaches to quorum. These protocols provide the application programmer a façade of global serializability. Personally, I have invested a non-trivial portion of my career as a strong advocate for the implementation and use of platforms

Instead, applications are built using different techniques which do not provide the same transactional guarantees but still meet the needs of their businesses.

This paper explores and names some of the practical approaches used in the implementations of large-scale mission-critical applications in a world which rejects distributed transactions. We discuss the management of fine-grained pieces of application data which may be repartitioned over time as the application grows. We also discuss the design patterns used in sending messages between these repartitionable pieces of data.



Pat Helland

Distributed Systems Guru
Worked at Amazon,
Microsoft & Salesforce

“
Grown-Ups Don't Use
Distributed Transactions



Starbucks does not use two phase commit

https://www.enterpriseintegrationpatterns.com/ramblings/18_starbucks.html

Photo by John Ingle

Eric Brewer

Atomicity
Consistency
Isolation
Durability

- ◆ But we forfeit “C” and “I” for availability, graceful degradation, and performance

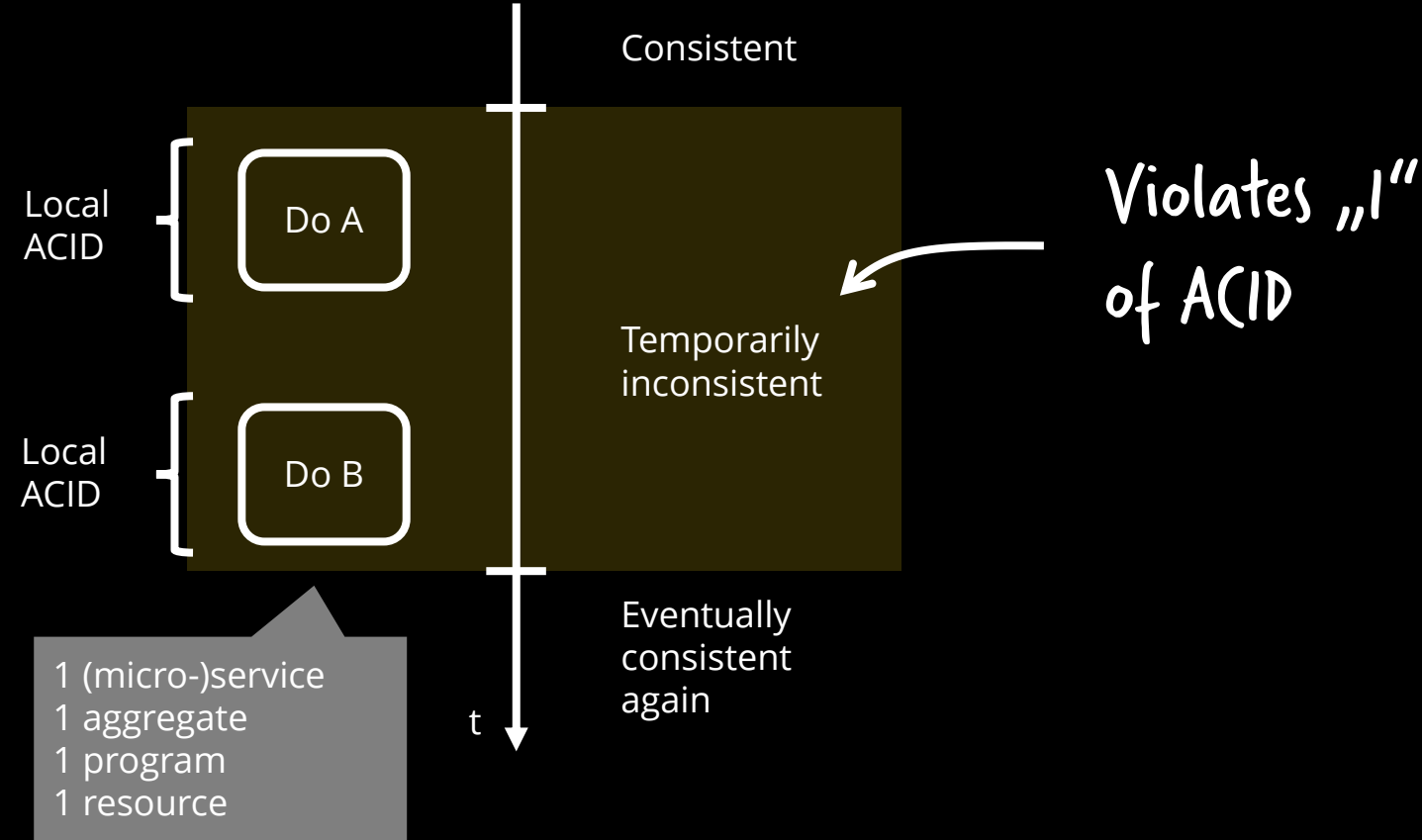
This tradeoff is fundamental.

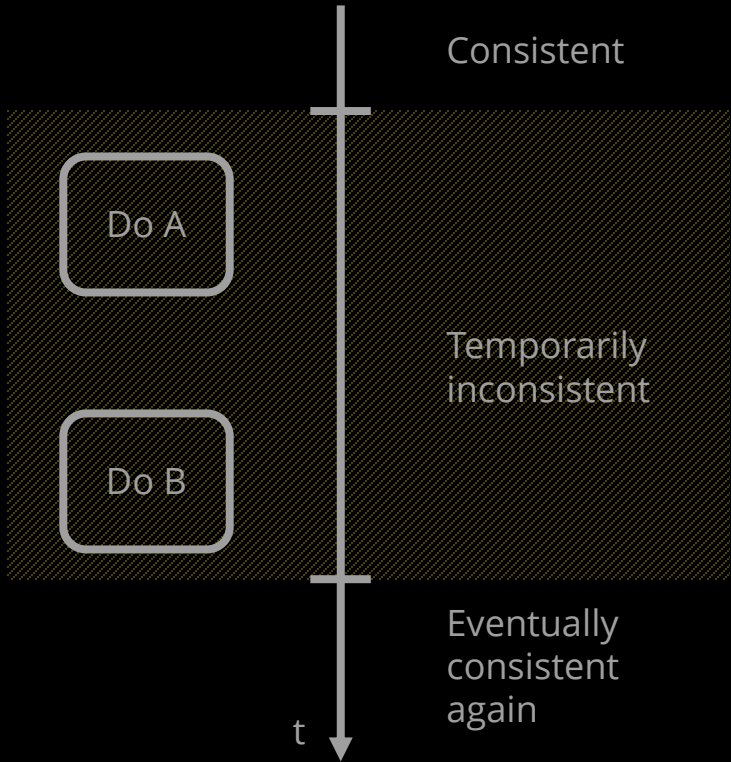
BASE:

- **B**asically **A**vailable
- **S**oft-state
- **E**ventual consistency

PODC Keynote, July 19, 2000

That means





You might know this from:





Pat Helland

„Building on Quicksand“ Paper

A
C
I
D
2.0



Pat Helland

„Building on Quicksand“ Paper

Associative
Commutative
Idempotent
Distributed
2.0

$$(a + b) + c = a + (b + c)$$

$$a + b = b + a$$

$$f(x) = f(f(x))$$



Photo by [pixabay](#), available under [Creative Commons CC0 1.0 license](#).

Requirement: Idempotency of services!



Photo by [pixabay](#), available under [Creative Commons CC0 1.0 license](#).

Requirement: Idempotency of services!

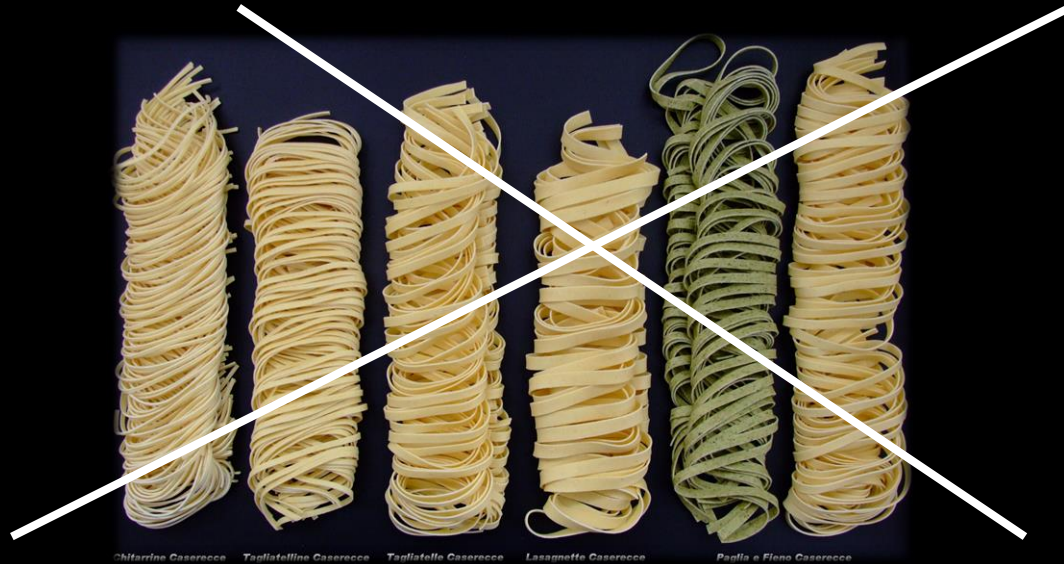
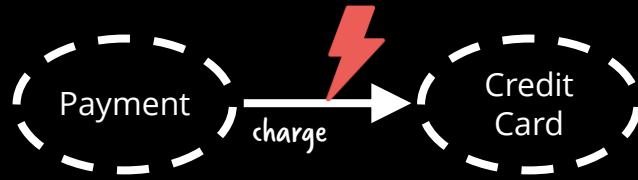
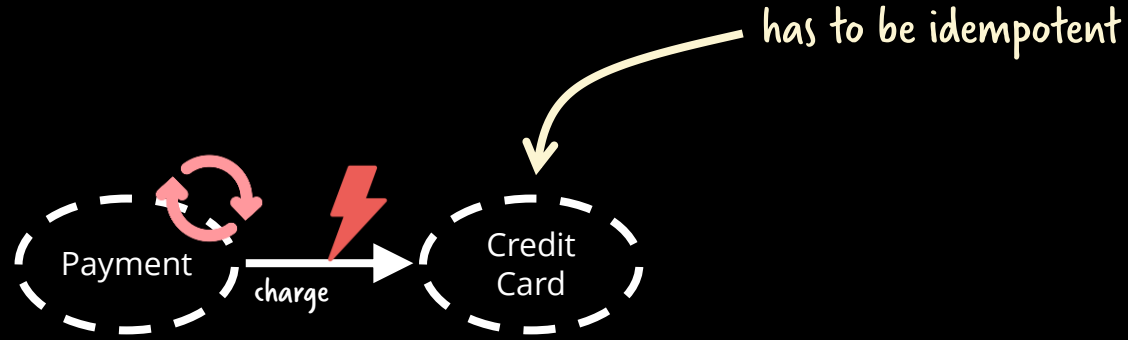


Photo by [Chr.Späth](#), available under [Public Domain](#).

Example



Strategy: retry



Charge Credit Card
cardNumber
amount

← **Not idempotent**

Charge Credit Card
cardNumber
amount
transactionId

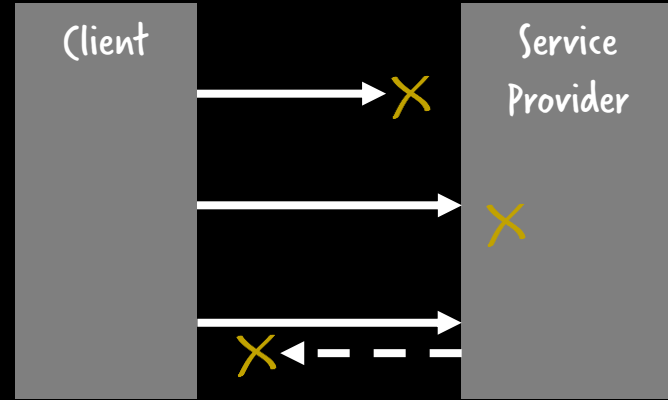
← **Idempotent**

Distributed

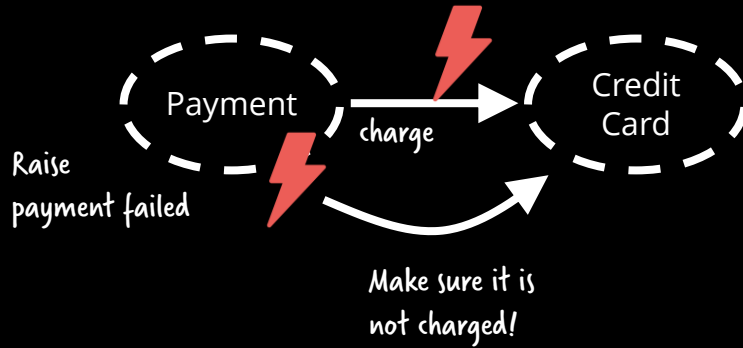


It is impossible to differentiate certain failure scenarios:

Independant of communication style!



Strategy: Cleanup

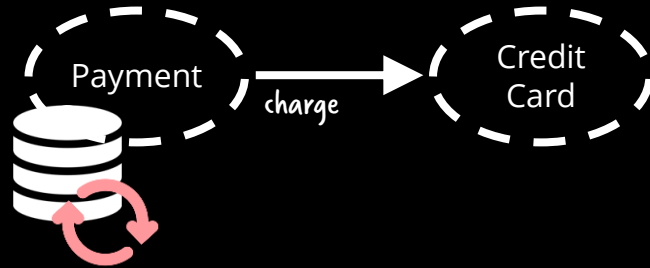


Cancel charge
cardNumber
amount
transactionId

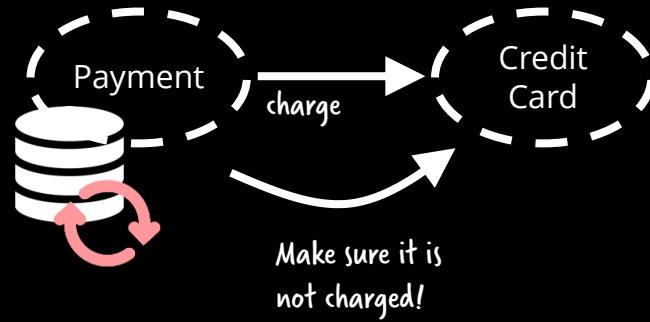


Some communication challenges
require state.

Strategy: Stateful retry



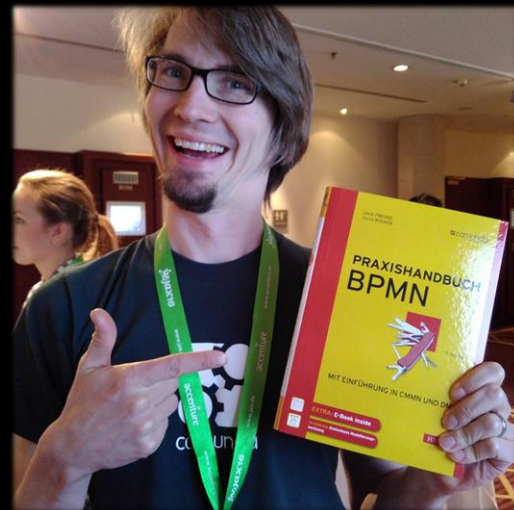
Strategy: Stateful retry



**Warning:
Contains Opinion**



Bernd Ruecker
Co-founder and
Chief Technologist of
Camunda

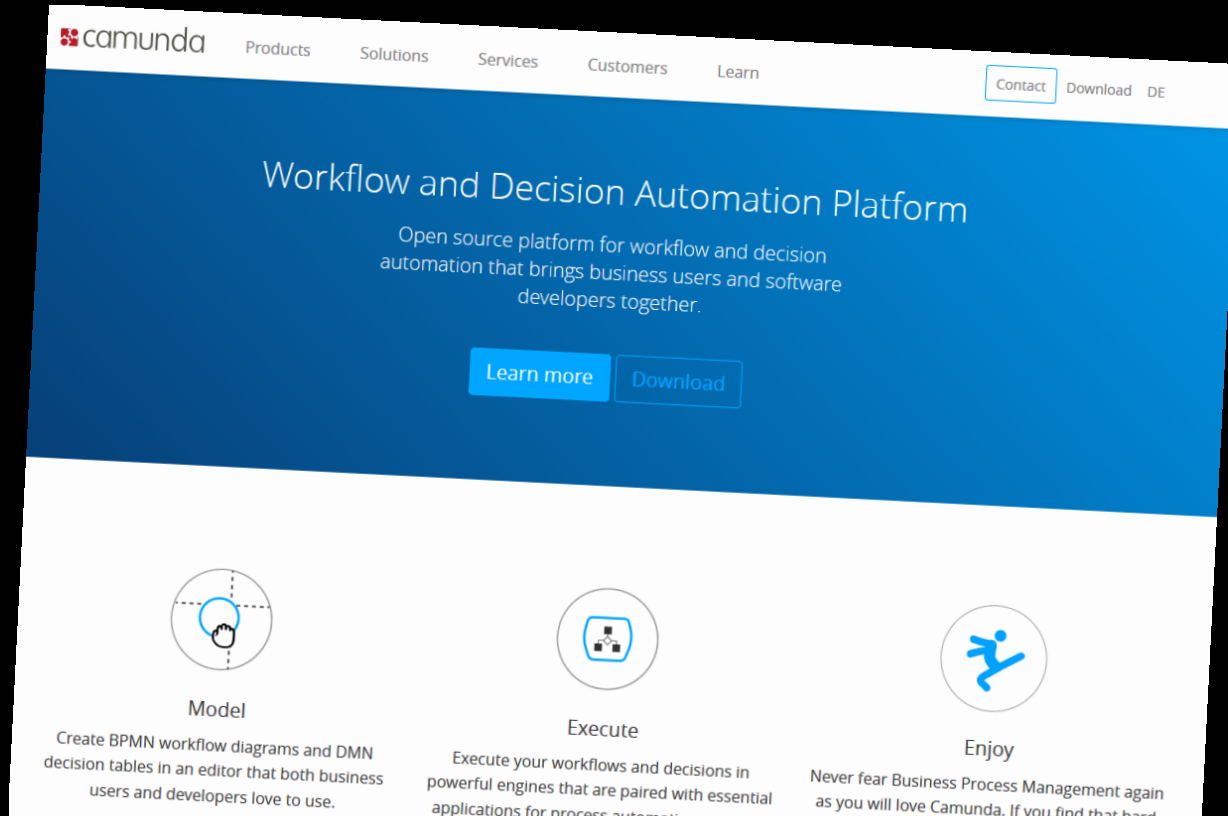


Berlin, Germany



bernd.ruecker@camunda.com
@berndruecker

Let's use a lightweight OSS workflow engine for this:




The image shows a screenshot of the Camunda website homepage. The page has a white header with the Camunda logo and navigation links for Products, Solutions, Services, Customers, and Learn. On the right side of the header, there are buttons for Contact, Download, and DE. The main content area has a blue background with the title "Workflow and Decision Automation Platform" and a sub-headline: "Open source platform for workflow and decision automation that brings business users and software developers together." Below this are two buttons: "Learn more" and "Download". The lower section of the page is white and features three circular icons representing different stages: "Model" (a hand pointing to a lightbulb), "Execute" (a server rack), and "Enjoy" (a person running). Each icon is accompanied by a short paragraph of text.

camunda Products Solutions Services Customers Learn [Contact](#) [Download](#) [DE](#)

Workflow and Decision Automation Platform


Open source platform for workflow and decision automation that brings business users and software developers together.

[Learn more](#) [Download](#)




Model

Create BPMN workflow diagrams and DMN decision tables in an editor that both business users and developers love to use.



Execute

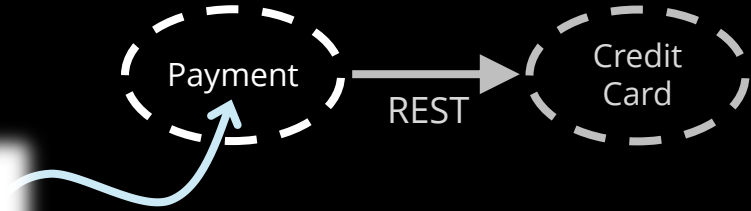
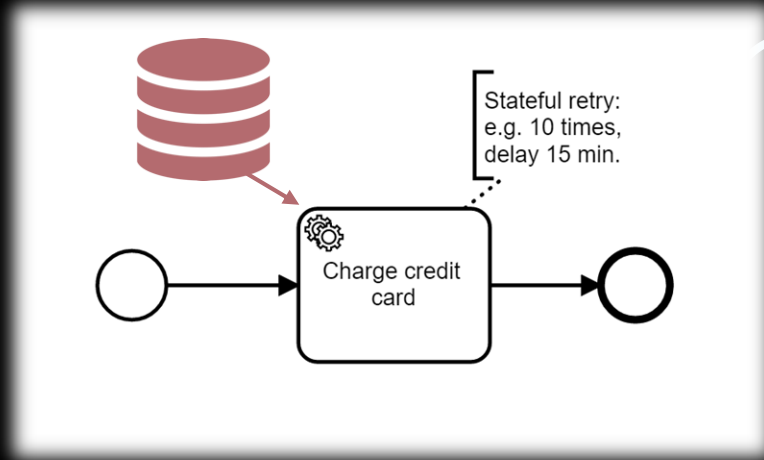
Execute your workflows and decisions in powerful engines that are paired with essential applications for process automation.



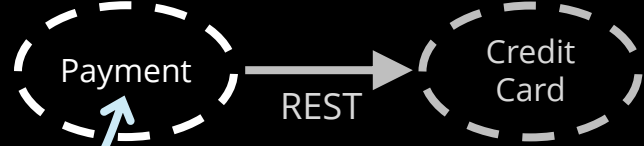
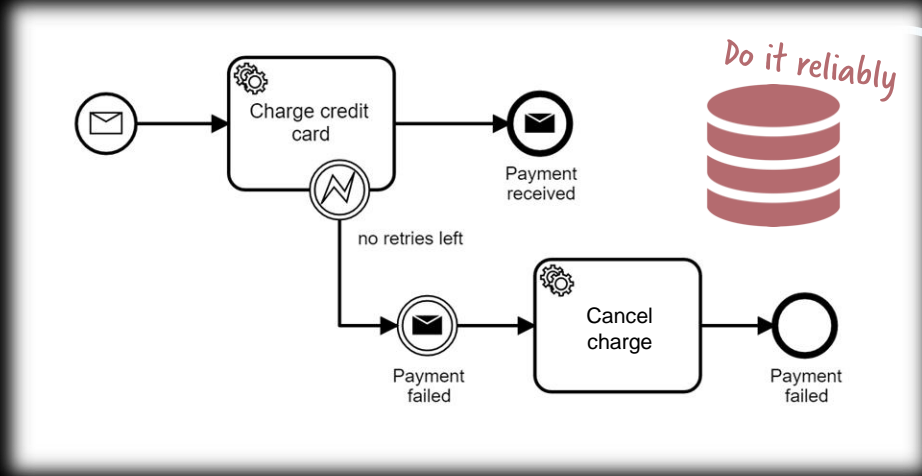
Enjoy

Never fear Business Process Management again as you will love Camunda. If you find that hard

Stateful retry



Stateful retry & cleanup

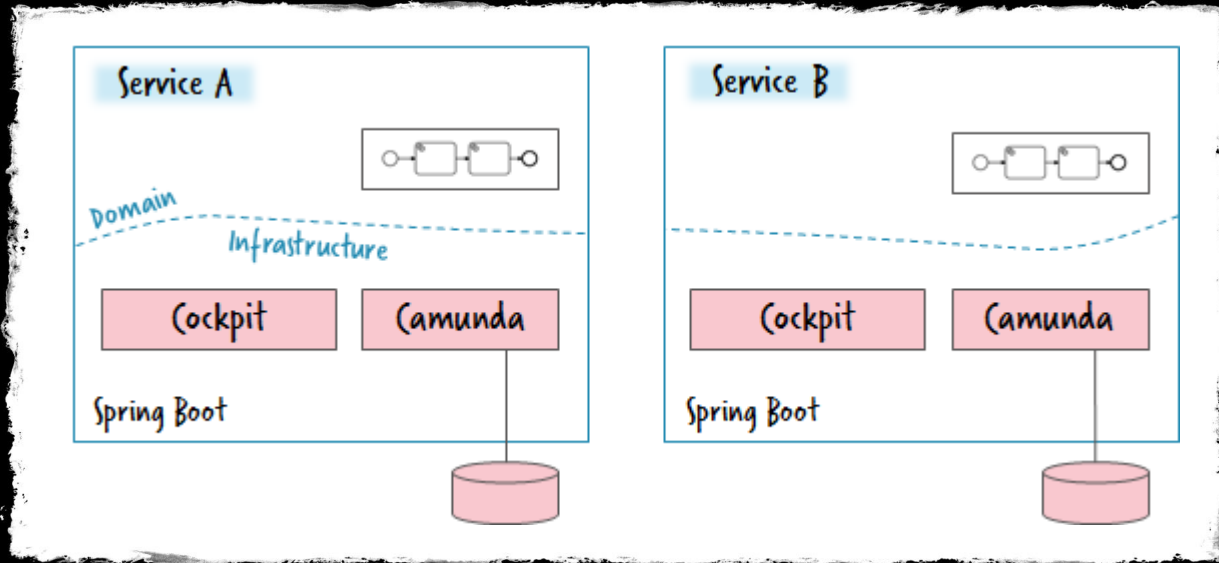


Live hacking

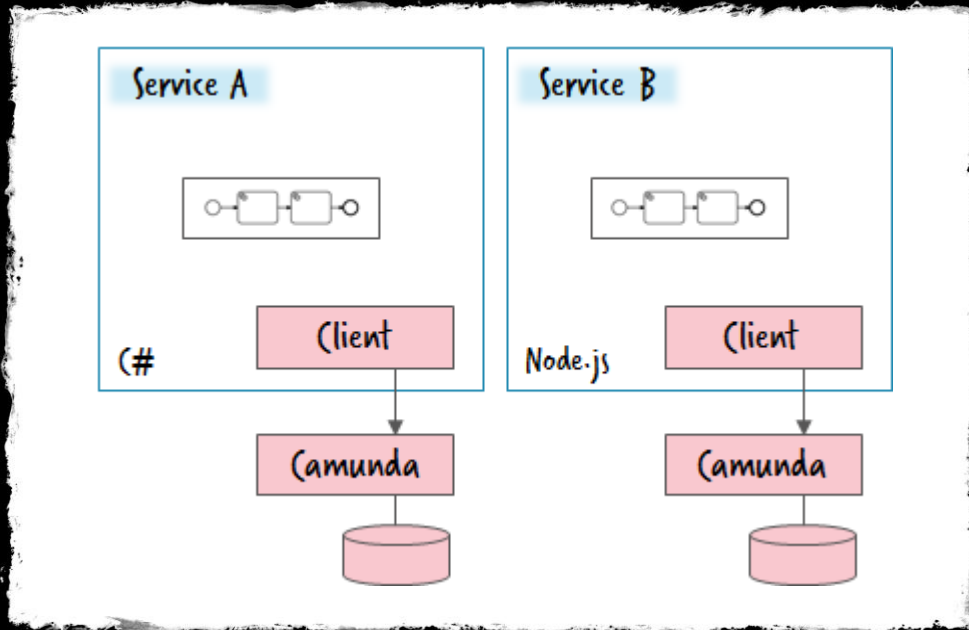


<https://github.com/flowing/flowing-retail/tree/master/rest>

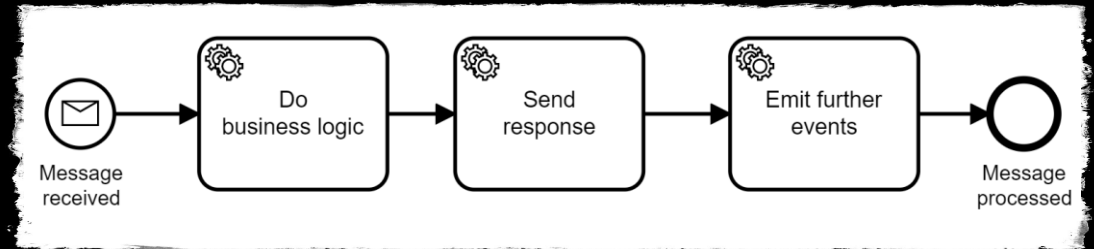
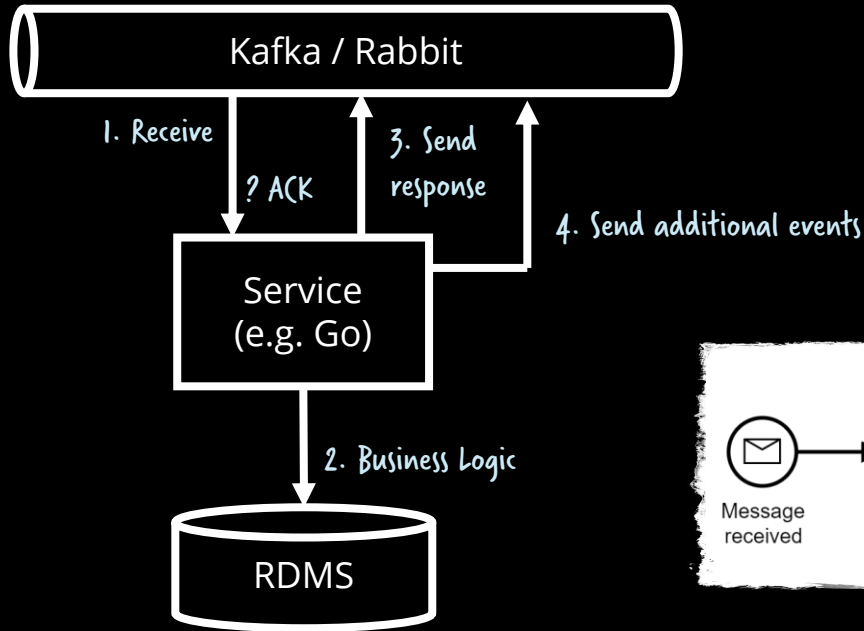
Embedded Engine Example (Java)



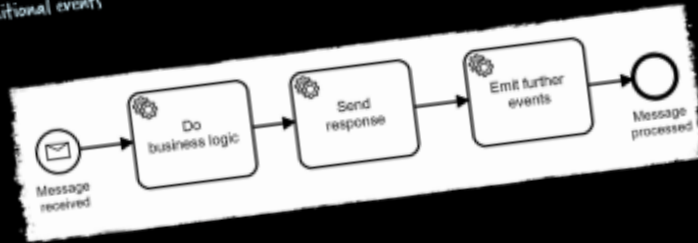
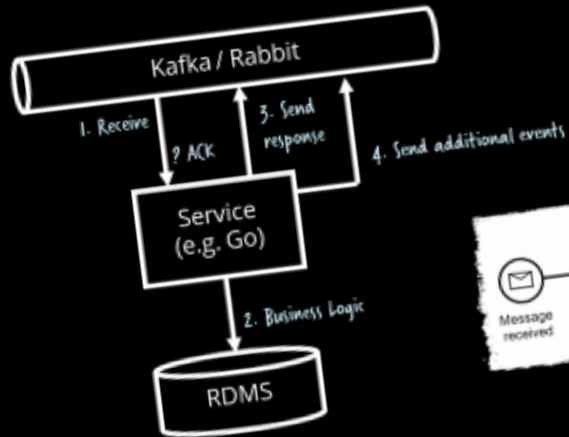
Remote Engine Example (Polyglot)



A relatively common pattern

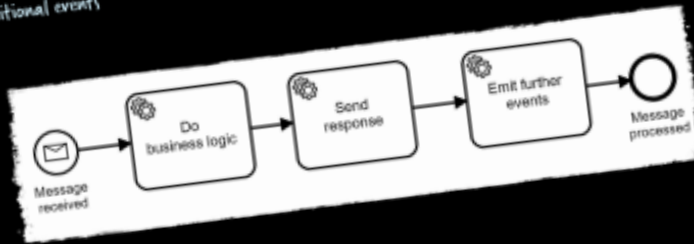
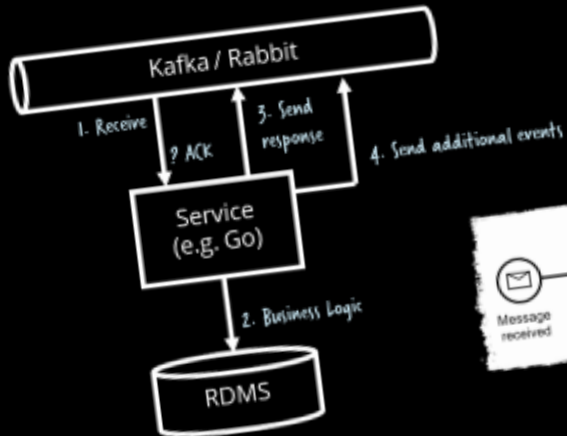


State can solve important basic problems



„Can this handle 15k requests per second?“

State can solve important basic problems

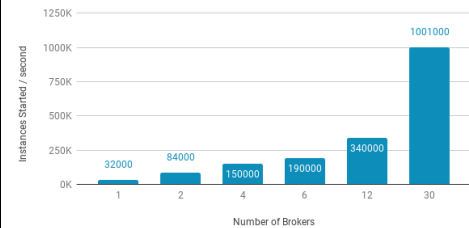


„Yes.“

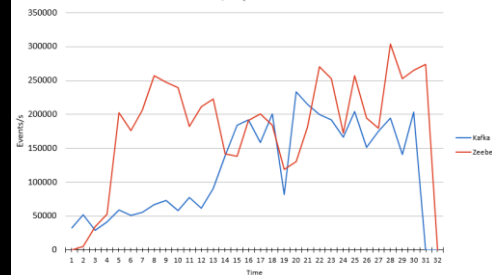


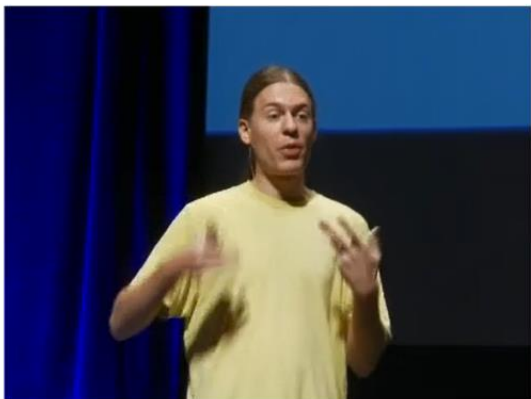
Workflow Instances Started / Second

Single topic, replication factor = 1



Events written/s Apache Kafka vs. Zeebe





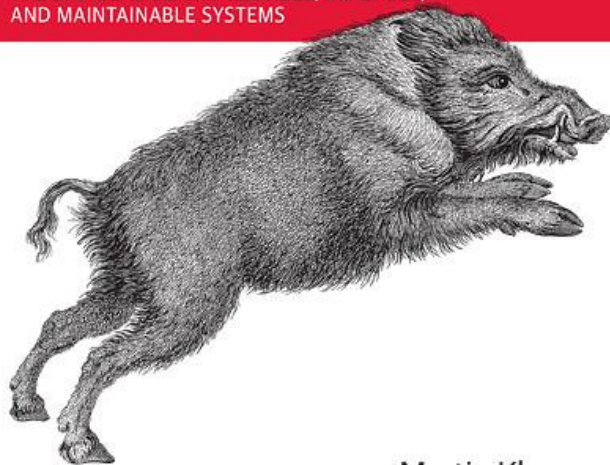
Sept 25-26, 2015

thestrangeloop.com

O'REILLY®

Designing Data-Intensive Applications

THE BIG IDEAS BEHIND RELIABLE, SCALABLE,
AND MAINTAINABLE SYSTEMS



Martin Kleppmann



strangeloop

Sept 25-26, 2015

thestrangeloop.com

Without cross-service transactions:

A

Compensating transactions

≈ abort/rollback at app level

(Garcia-Molina & Salem, 1987)

C

Apologies

(Helland & Campbell, 2009)

detect & fix constraint violations
after the fact, rather than preventing them



Memories, Guesses, and Apologies

Rate this article ★★★★★

 Pat Helland May 15, 2007

 5

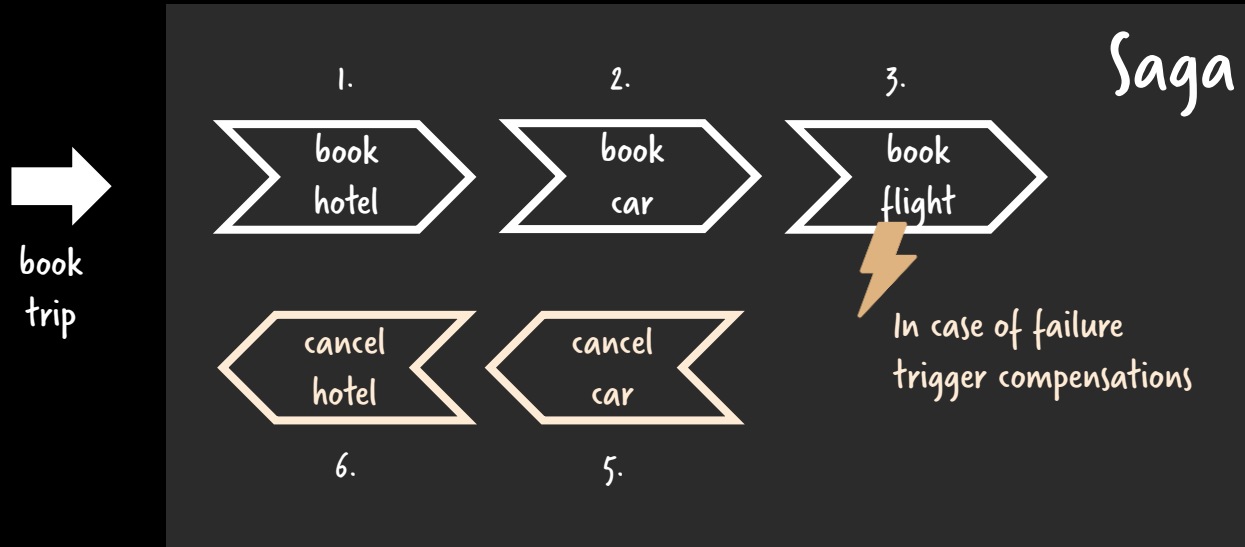
 Share 3  12  0

Well, here I am blogging on the bus with my newly installed Windows Live Writer!!!

This blog is a text version of a five minute "Gong Show" presentation I did at CIDR (Conference on Innovative Database Research) on Jan 8,2007.

All computing can be considered a "M...
...computers suck. Furthermore, it offers additi...
...this is a personal opinion about

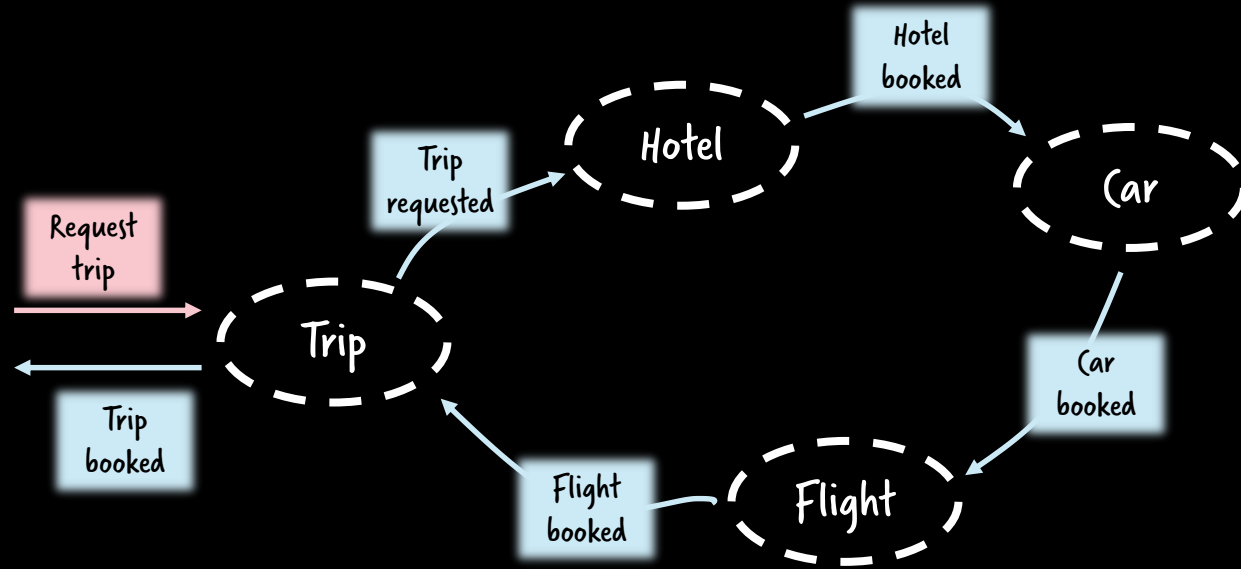
Compensation – the classical example



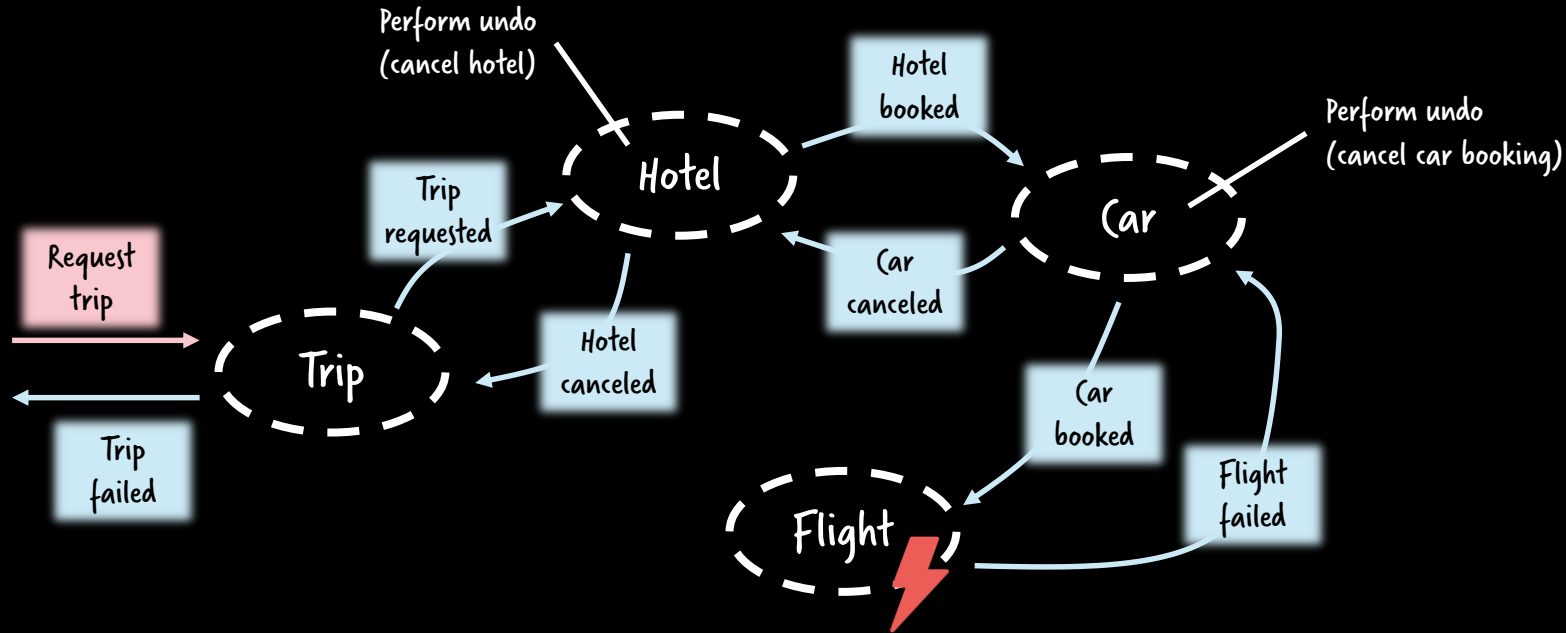
2 alternative approaches: choreography & orchestration



Event-driven choreography



Event-driven choreography





The danger is that it's very easy to make nicely decoupled systems with event notification, without realizing that you're losing sight of that larger-scale flow, and thus set yourself up for trouble in future years.

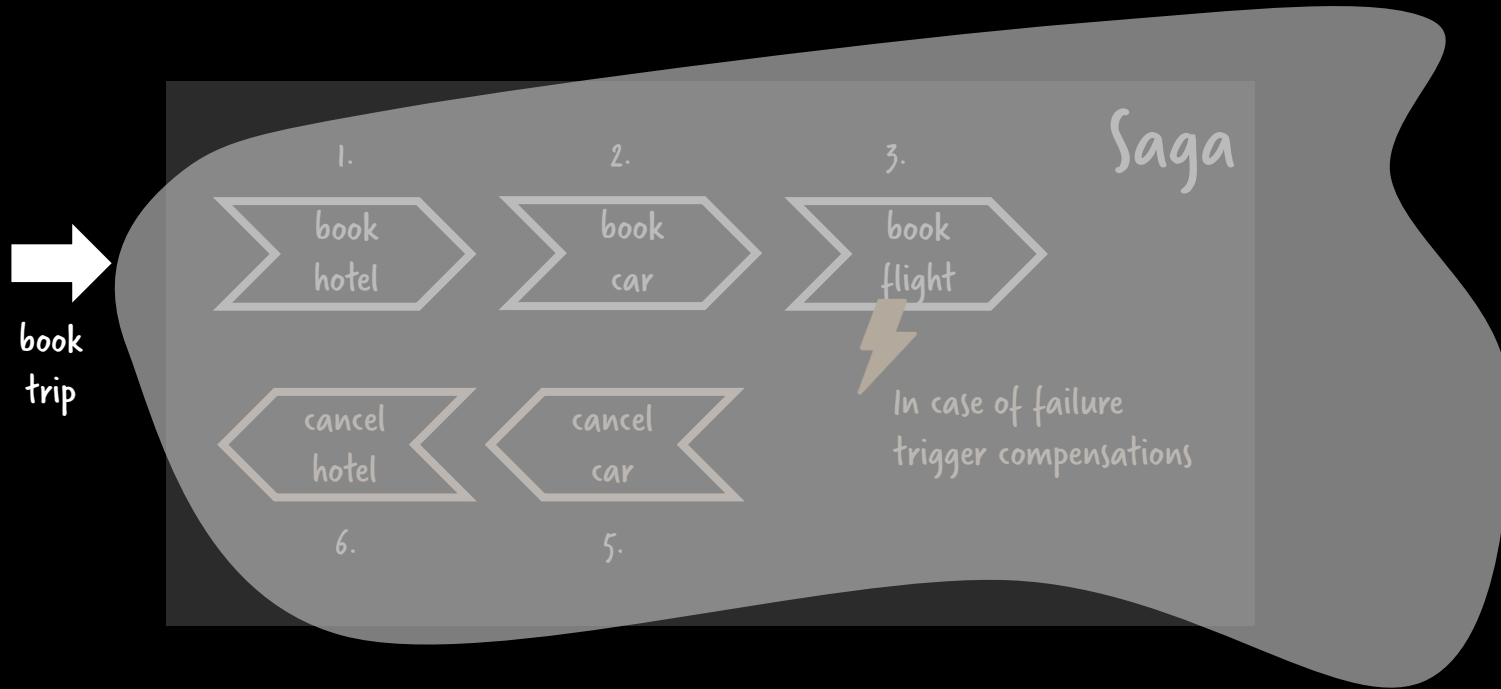


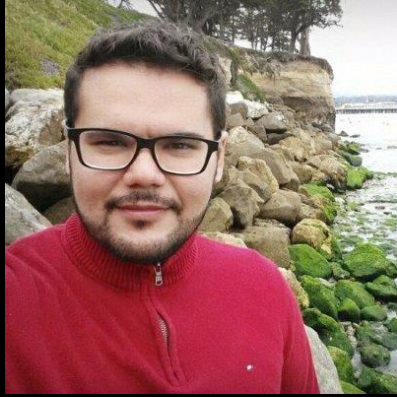
The danger is that it's very easy to make nicely decoupled systems with event notification, without realizing that you're losing sight of that larger-scale flow, and thus set yourself up for trouble in future years.



The danger is that it's very easy to make nicely decoupled systems with event notification, without realizing that you're losing sight of that larger-scale flow, and thus set yourself up for trouble in future years.

Classical example





Denis Rosa
Couchbase

If your transaction involves 2 to 4 steps, choreography might be a very good fit.

However, this approach can rapidly become confusing if you keep adding extra steps in your transaction as it is difficult to track which services listen to which events. Moreover, it also might add a cyclic dependency between services as they have to subscribe to one another's events.

Microservice pioneers have become aware

Traditionally, some of these processes had been orchestrated in an ad-hoc manner using a combination of pub/sub, making direct REST calls, and using a database to manage the state. However, as the number of microservices grow and the complexity of the processes increases, getting visibility into these distributed workflows becomes difficult without a central orchestrator.



Netflix Technology Blog [Follow](#)

Learn more about how Netflix designs, builds, and operates our systems and engineering organizations

Dec 12, 2016 · 7 min read

Netflix Conductor: A microservices orchestrator

The Netflix Content Platform Engineering team runs a number of business processes which are driven by asynchronous orchestration of tasks executing on microservices. Some of these are long running processes spanning several days. These processes play a critical role in getting titles ready for streaming to our viewers across the globe.

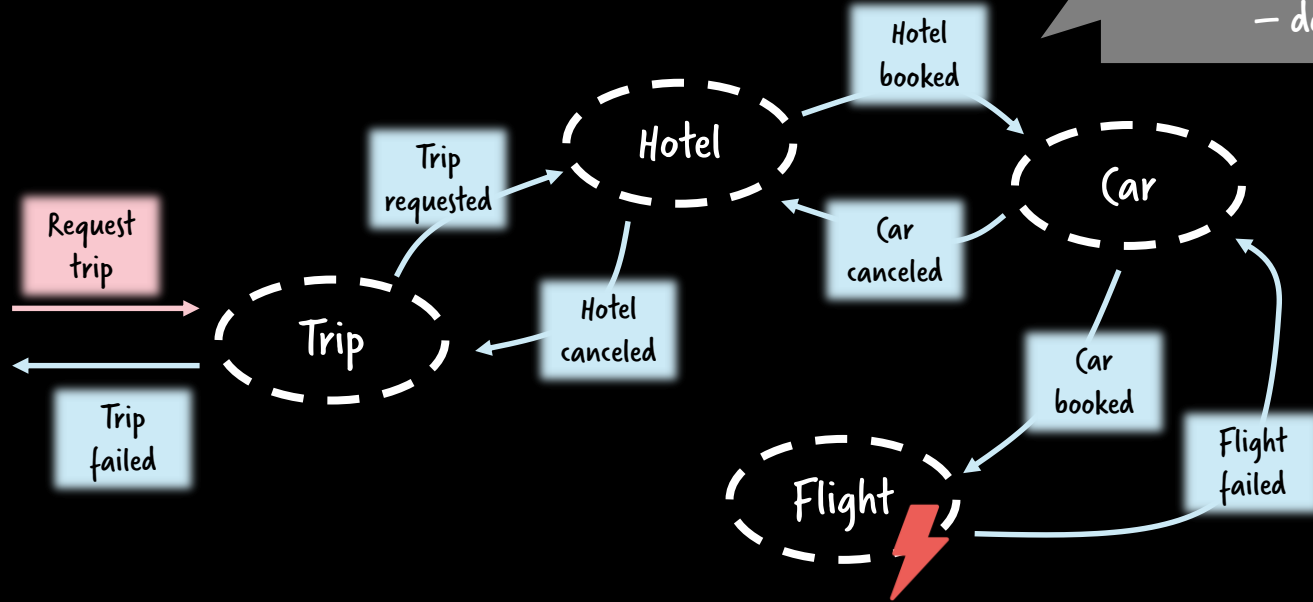
A few examples of these processes are:

- Studio partner integration for content ingestion
- [IMF](#) based content ingestion from our partners
- Process of setting up new titles within Netflix
- Content ingestion, encoding, and deployment to CDN

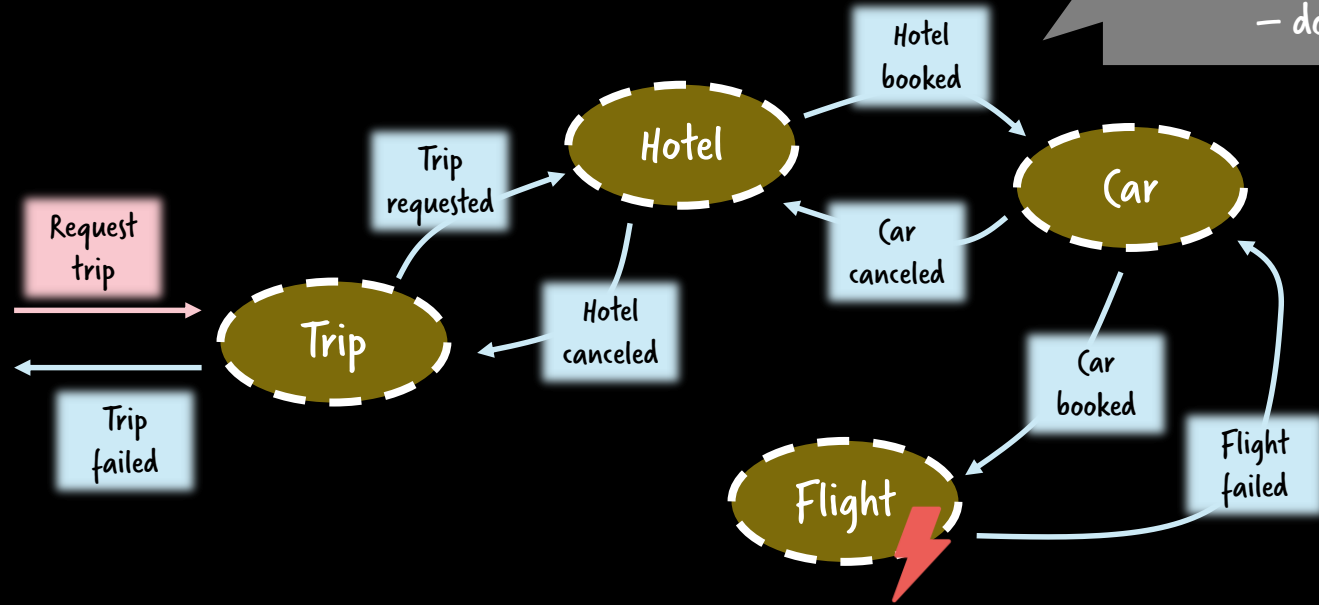
Traditionally, some of these processes had been orchestrated in an ad-hoc manner using a combination of pub/sub, making direct REST calls, and using a database to manage the state. However, as the number of microservices grow and the complexity of the processes increases, getting visibility into these distributed workflows becomes difficult without a central orchestrator.

Implementing changes in the process

We have a new basic agreement with the car rental agency and can cancel for free within 1 hour – do that first!



Implementing changes in the process



We have a new basic agreement with the car rental agency and can cancel for free within 1 hour – do that first!

You have to adjust all services and redeploy at the same time!



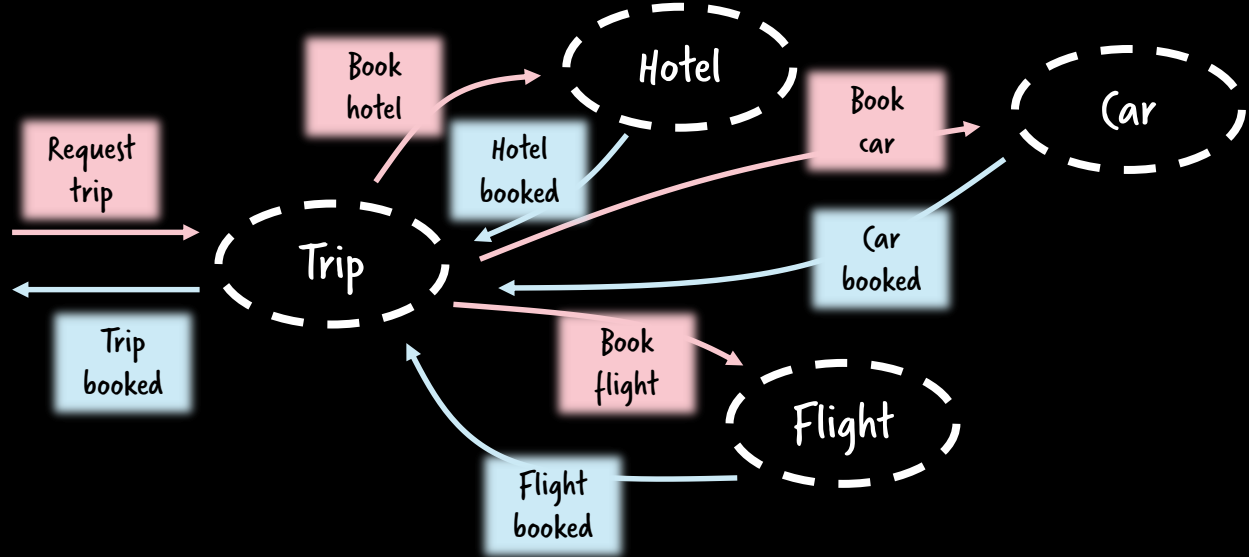
Photo by born1945, available under [Creative Commons BY 2.0 license](https://creativecommons.org/licenses/by/2.0/).

What we wanted



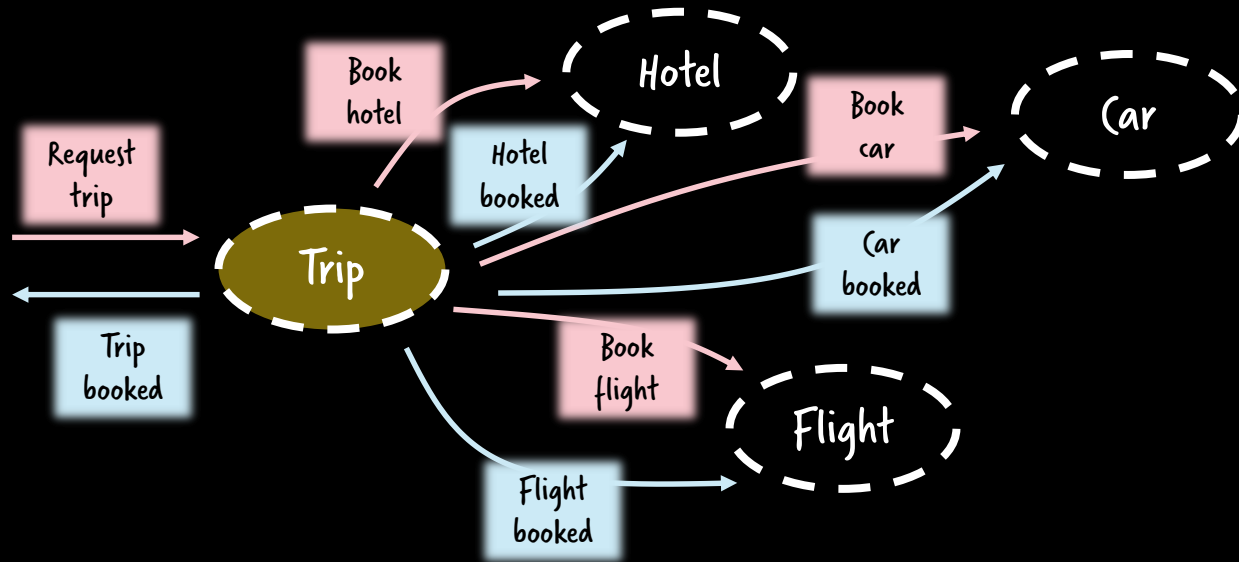
vs. what we got

orchestration



orchestration

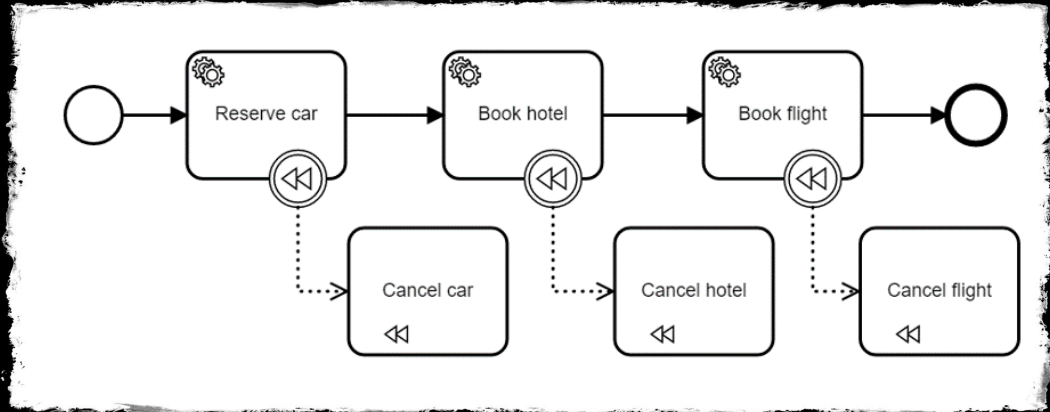
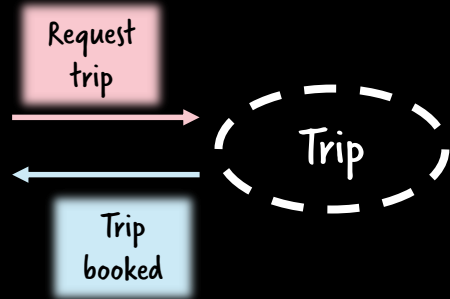
We have a new basic agreement with the car rental agency and can cancel for free within 1 hour – do that first!



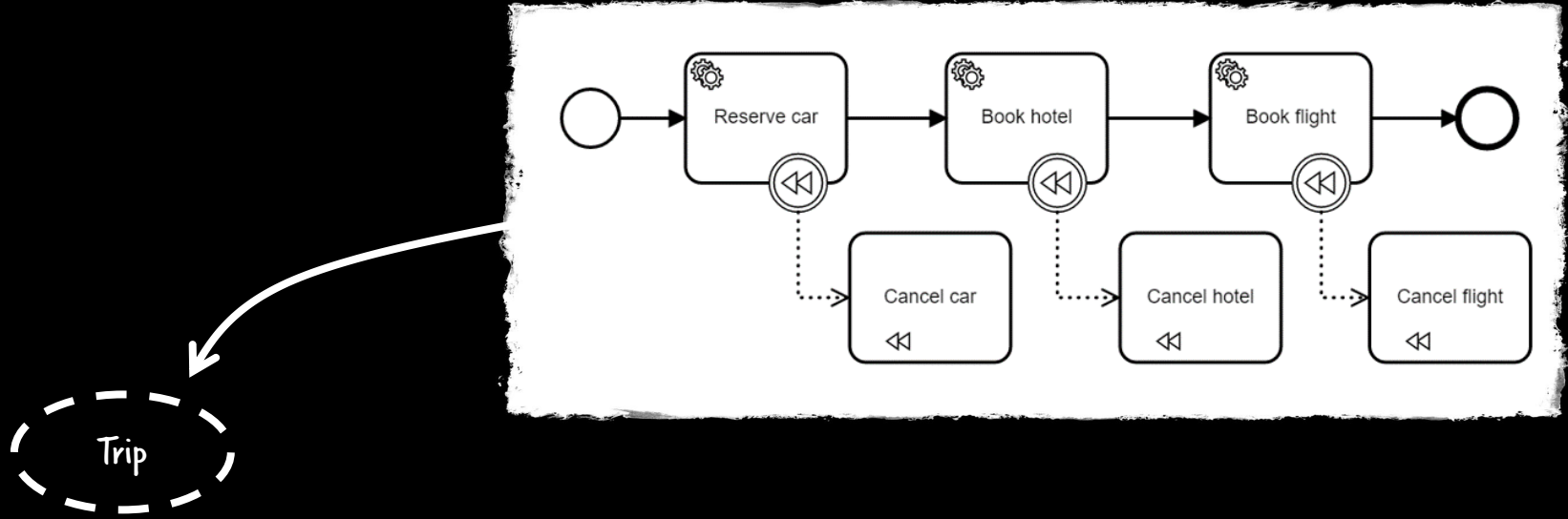
You have to adjust one service and redeploy only this one!

Describe orchestration with BPMN

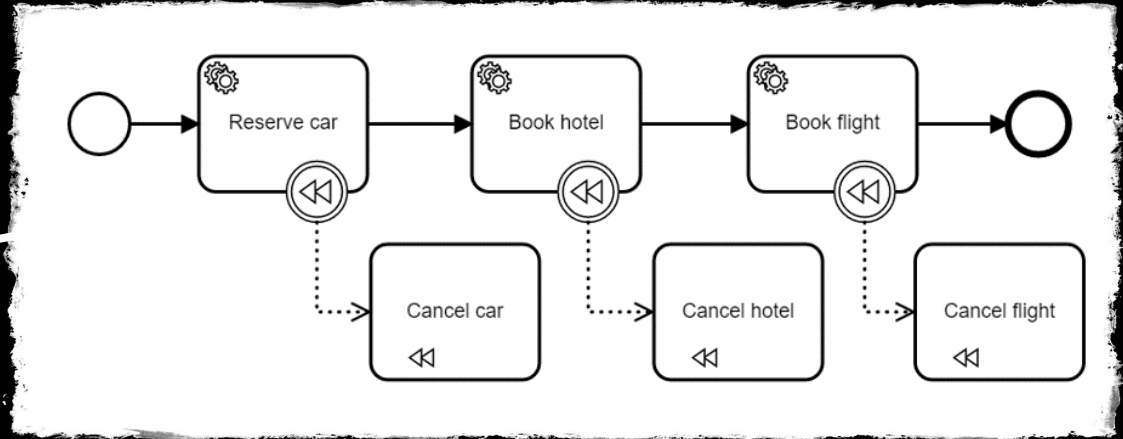
Saga Pattern
(implemented by BPMN
compensation)



The workflow is part of the service

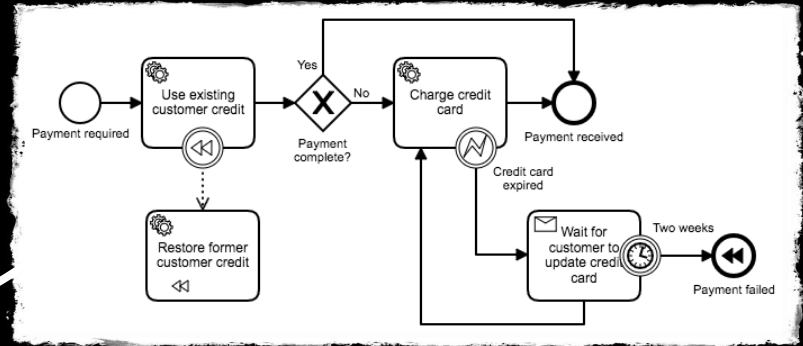


The workflow is part of the service



Trip

Payment





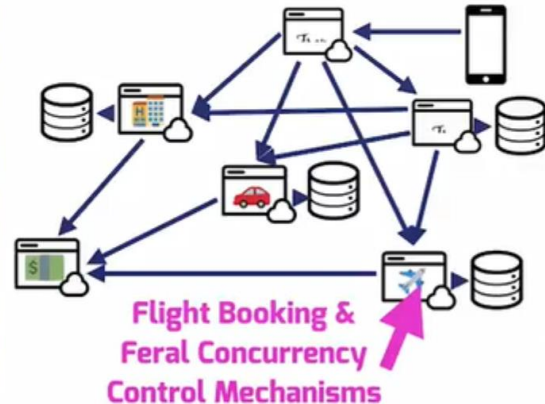
ON THE BEACH

{ IT_Event(); }

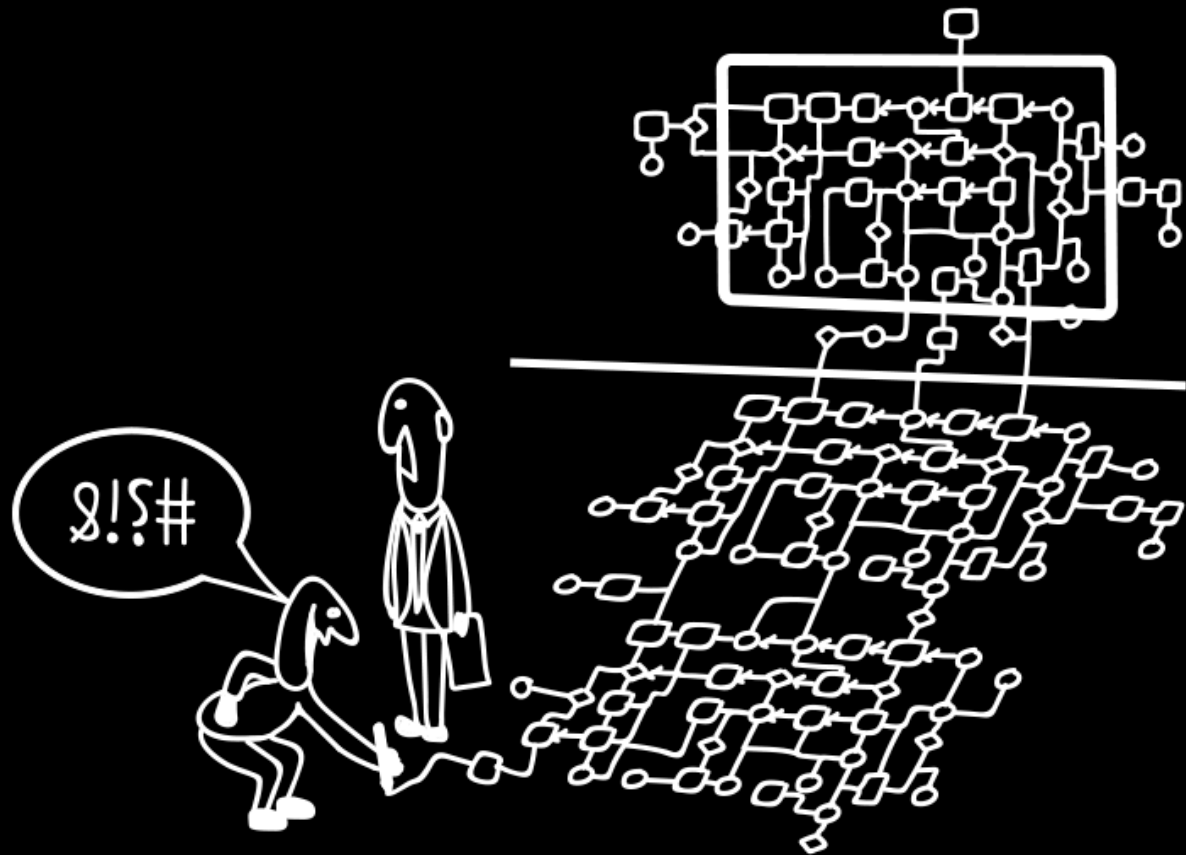


Modular Services

with Distributed Sagas



Graphical models?





Clemens Vasters
Architect at Microsoft

Sagas

Today has been a holiday in some parts of the .NET community debating the Saga pattern. Let's clarify, there are a few frameworks for "Saga" but more that use the term "Saga" for some framework implementation of a state machine or workflow. That's not what a Saga is. Saga is a failure management pattern.

Sagas come out of the realization that particular long-lived transactions (originally even just in the database), but also for distributed transactions across location and/or most importantly can't easily be handled using the classic 2PC model with 2-Phase commit and holding locks for the duration of the work instead. A Saga splits work into individual transactions whose effects can be, optionally, reversed after work has been performed and admitted.

The diagram shows a sequence of three boxes: 'Rental Cars', 'Hotels', and 'Flights'. Each box contains a green arrow pointing right labeled 'Book' and a green arrow pointing left labeled 'Cancel'. A yellow box on the left contains 'Car!', 'Hotel!', and 'Flight!' with an arrow pointing to the 'Book' arrow of the 'Rental Cars' box. A starburst icon is on the 'Book' arrow of the 'Flights' box.

The above shows a single Saga. If you book travel itinerary you want a car and a hotel and a flight. If you can't get all of them, it's probably not worth going. It's also important that you admit one or all of these problems into a distributed 2PC transaction. Instead, you'll have an activity for booking rental cars that knows both how to perform a reservation and also how to admit it and one for a hotel and one for flights.

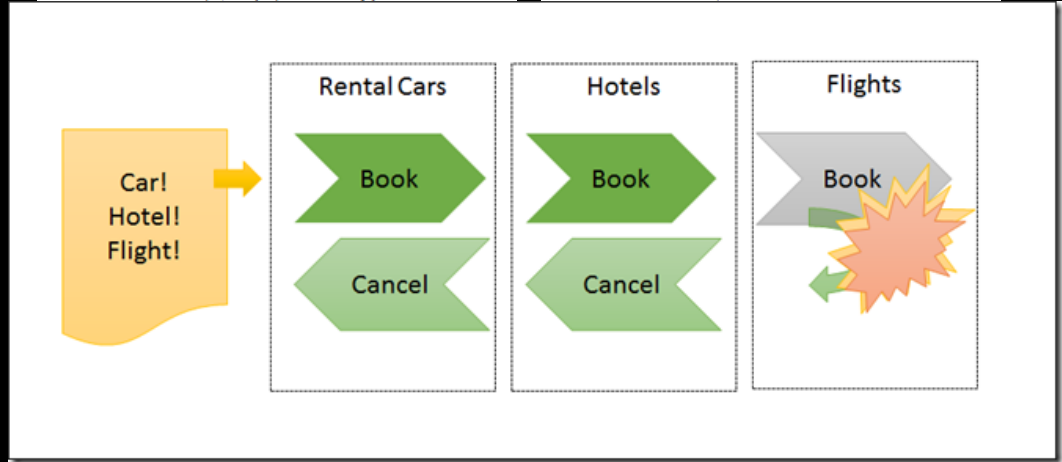
The activities are grouped in a composite job (routing slip) that's handled along the activity chain. If you want you can design the routing slip items so that they can only be understood and manipulated by the intended handler. When an activity completes, it also hands off the completion to the routing slip along with information on where the complementary operation can be reached (e.g. via a Queue). When an activity fails, it drops up (back) and then sends the routing slip back down to the last.

Clemens Vasters
@clemensvasters

```

1: class RoutingSlip
2: {
3:     readonly IEnumerable<ActivityDefinition> activities;
4:     readonly Queue<ActivityDefinition> queue;
5:     public RoutingSlip()
6:     {
7:     }
8:     public RoutingSlip(IEnumerable<ActivityDefinition> activities)
9:     {
10:         queue = new Queue<ActivityDefinition>();
11:     }
12:     public void SubmitAndStart()
13:     {
14:         get { return this.activities.First(); }
15:     }
16:     public void SubmitProgress()
17:     {
18:         get { return this.activities.Last(); }
19:     }
20:     public void Progress()
21:     {
22:         if (this.activities.Count() > 1)
23:             this.activities.Last().SubmitAndStart();
24:     }
25:     public void Complete()
26:     {
27:         var completion = this.activities.Last().Complete();
28:         var activity = completion.ActivityDefinition;
29:         var results = activity.Data.CompletionData;
30:         if (results != null)
31:             queue.Enqueue(results);
32:     }
33: }

```



```

1: //
2: // The activities each implement a reservation step and an undo step. Here's the one for cars
3:
4: class RentalCarActivity : Activity
5: {
6:     public RentalCar() { }
7:     public RentalCar(int days, int location) { }
8:     public override string Description()
9:     {
10:         return string.Format("Rental Car {0} for {1} days at {2}", this.location, this.days, this.location);
11:     }
12:     public override void SubmitAndStart()
13:     {
14:         var reservation = new RentalCarReservation(this);
15:         var completion = reservation.Submit();
16:         var results = completion.Data.CompletionData;
17:         return results;
18:     }
19:     public override void Complete()
20:     {
21:         var reservation = new RentalCarReservation(this);
22:         reservation.Cancel();
23:     }
24: }

```

```

1: //
2: //
3:
4: class FlightActivity : Activity
5: {
6:     public Flight() { }
7:     public Flight(int days, int location) { }
8:     public override string Description()
9:     {
10:         return string.Format("Flight {0} for {1} days to {2}", this.location, this.days, this.location);
11:     }
12:     public override void SubmitAndStart()
13:     {
14:         var reservation = new FlightReservation(this);
15:         var completion = reservation.Submit();
16:         var results = completion.Data.CompletionData;
17:         return results;
18:     }
19:     public override void Complete()
20:     {
21:         var reservation = new FlightReservation(this);
22:         reservation.Cancel();
23:     }
24: }

```

<http://vastars.com/archive/Sagas.html>

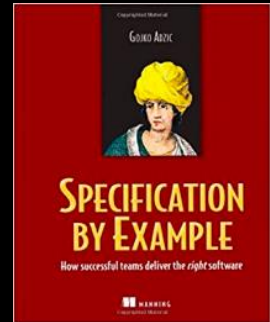
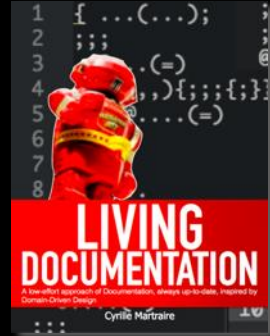
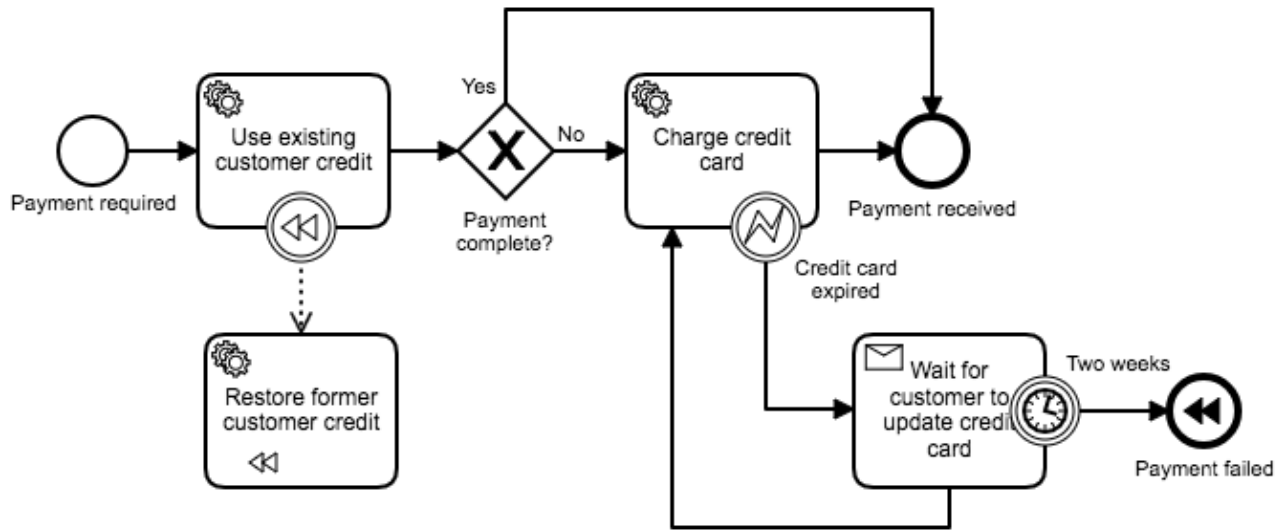
BPMN

Business Process
Model and Notation

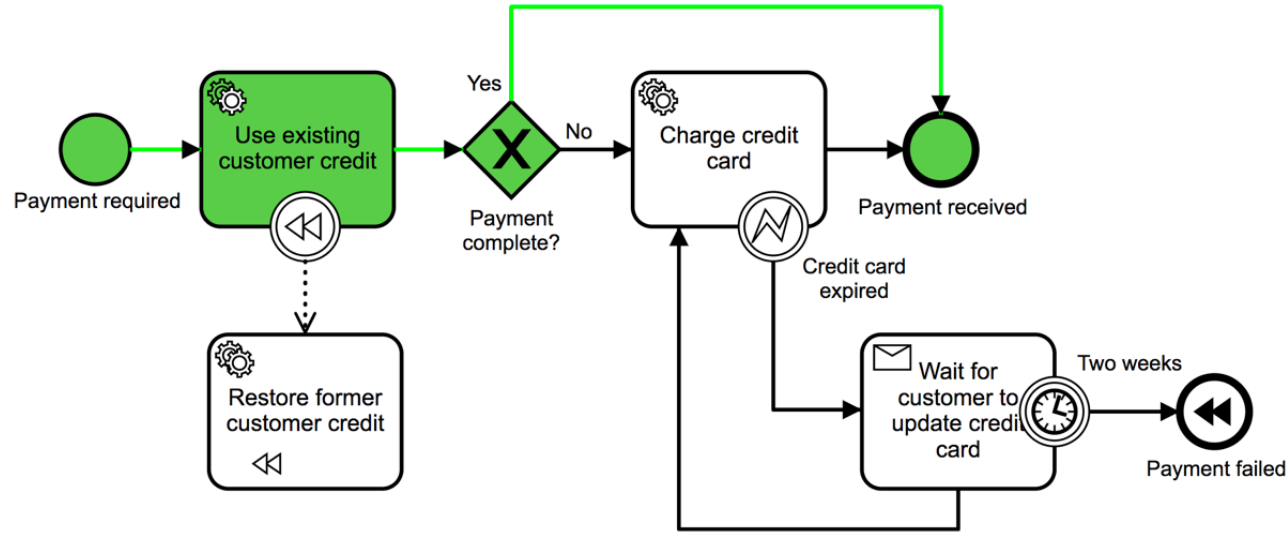
ISO Standard



Living documentation for long-running behaviour

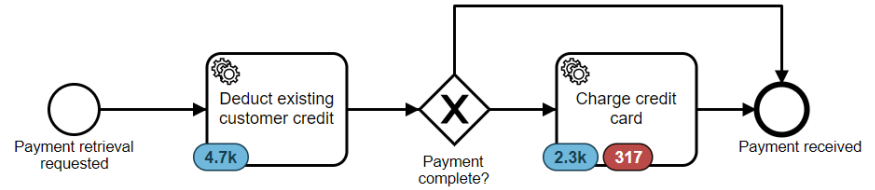


Visual HTML reports for test cases



Dashboard > Processes > paymentV5 : Runtime | History

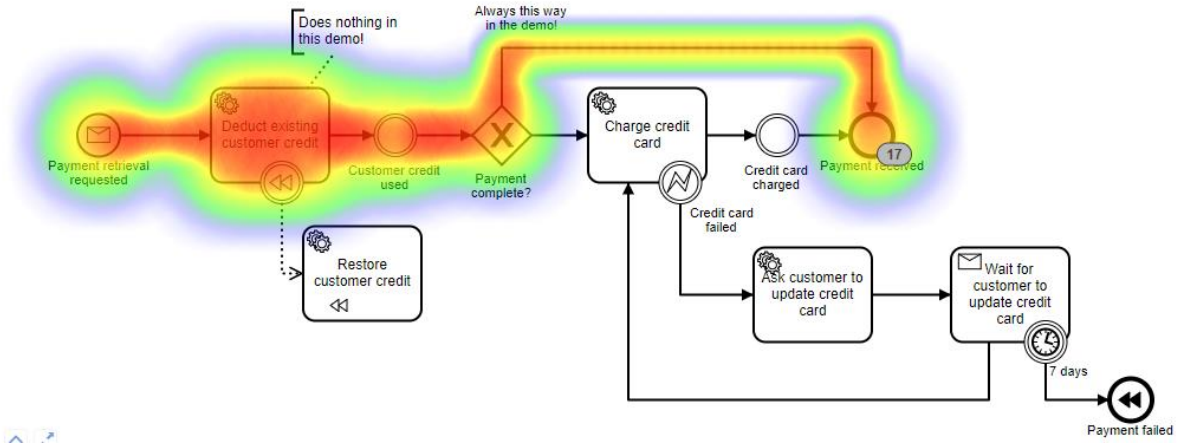
Definition Version: 2
 Version Tag: null
 Definition ID: paymentV5:2:45aea93a-1ad9-11e8-8...
 Definition Key: paymentV5
 Definition Name: null
 History Time To Live: null
 Tenant ID: null



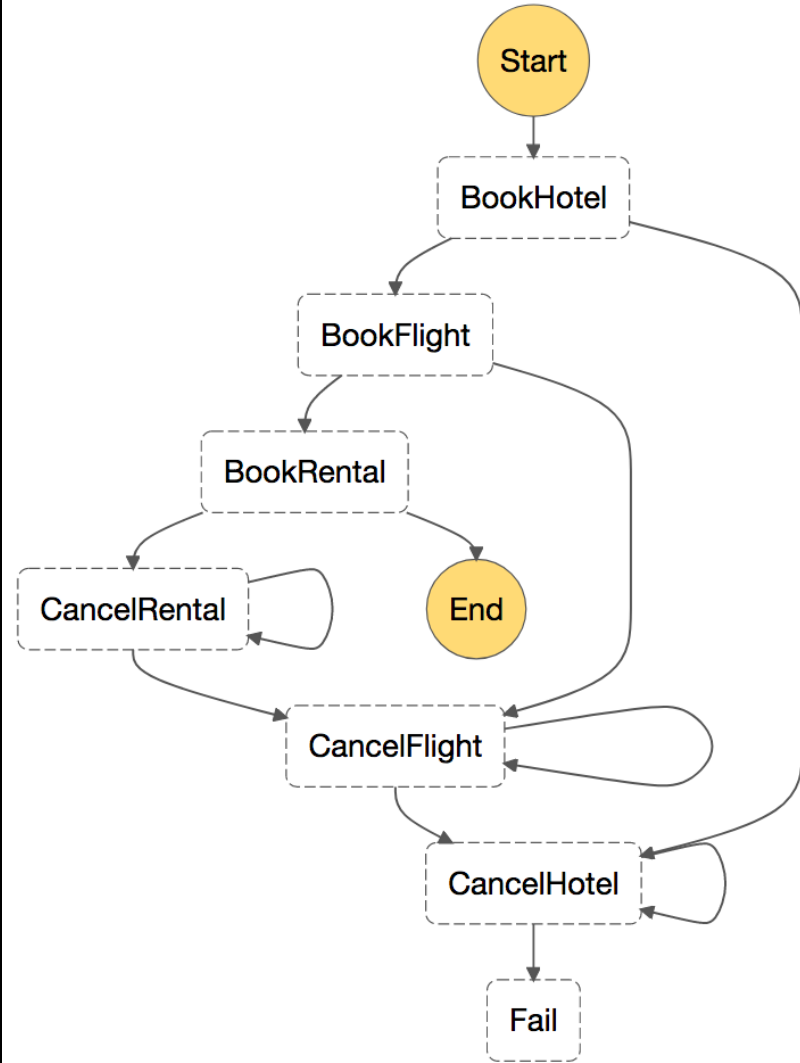
itions Modify

Start Time	Business Key
2018-02-26T10:40:59	
2018-02-26T10:40:18	

Powered by camunda BPM / v7.8.0-ee



Saga with AWS Step Functions



Thoughts on the state machine / workflow engine market



Thoughts on the state machine / workflow engine market

Stack Vendors,
Pure Play BPMS
Low Code Platforms

PEGA, IBM, SAG, ...

Camunda, Zeebe, jBPM,
Activiti, Mistral, ...

oss Workflow or
orchestration Engines

Integration Frameworks

Apache Camel,
Balerina, ...

Homegrown frameworks
to scratch an itch

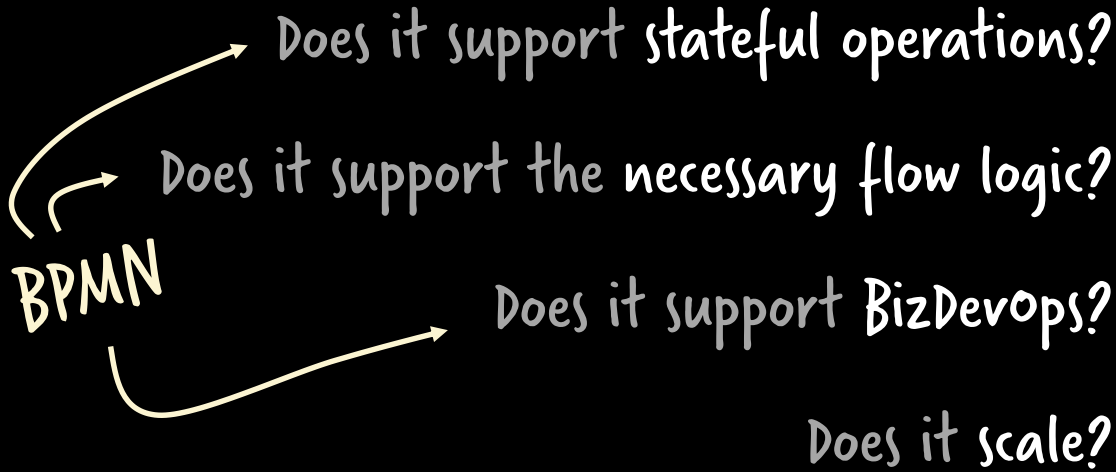
Uber, Netflix, AirBnb, ING, ...

Cloud offerings

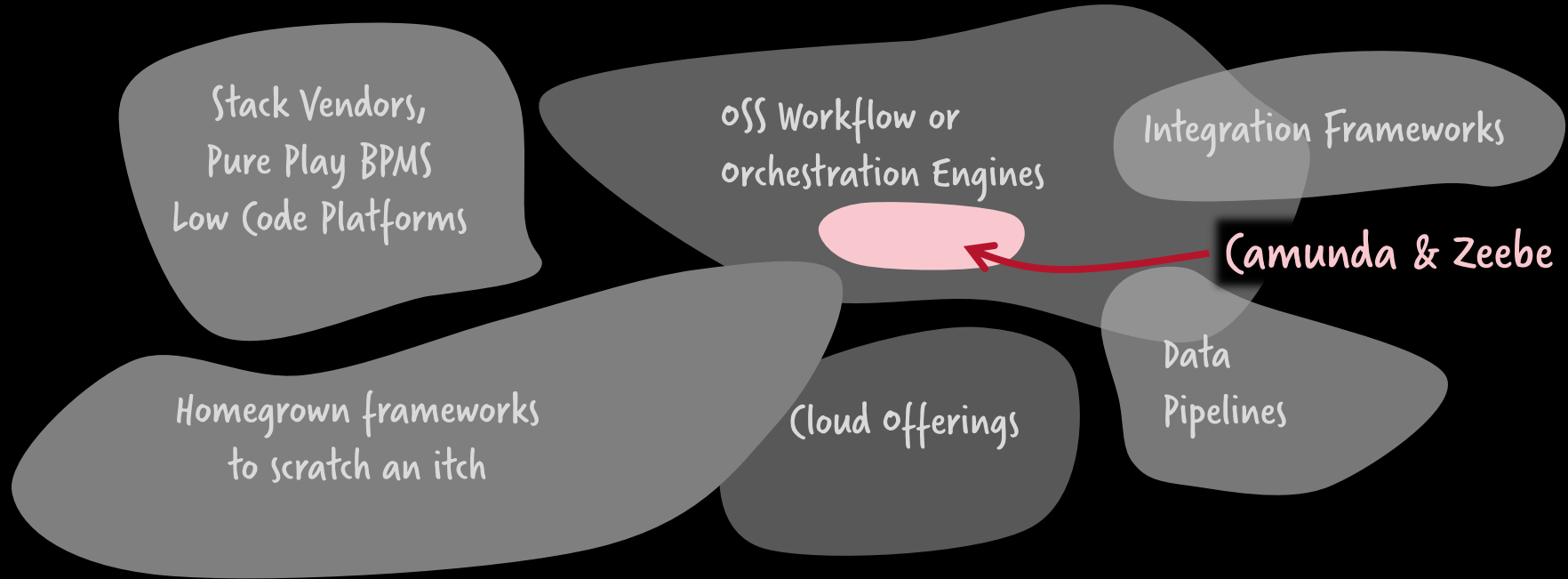
AWS Step Functions,
Azure Durable Functions, ...

Data
Pipelines

Apache Airflow,
Spring Data Flow, ...



My personal pro-tip for a shortlist ;-)



Recap

- Grown ups don't use distributed transactions but eventual consistency
- Idempotency is super important in distributed systems
- Some consistency challenges require state
- Know some strategies
 - Stateful retry & cleanup
 - Saga / Compensation
 - Apologies

Thank you!



Contact: mail@berndruecker.io
[@berndruecker](#)

Slides: <https://berndruecker.io>

Blog: <https://medium.com/berndruecker>

Code: <https://github.com/berndruecker>

InfoWorld
FROM IDG

<https://www.infoworld.com/article/3254777/application-development/3-common-pitfalls-of-microservices-integration-and-how-to-avoid-them.html>

InfoQ
neue

<https://www.infoq.com/articles/events-workflow-automation>

THE NEW STACK

<https://thenewstack.io/5-workflow-automation-use-cases-you-might-not-have-considered/>

