# Upgrade data architecture by making six foundational shifts.



Systems of engagement

Mobile | Web | CRM | Physical channels

Application databases

④ APIs and management platform

Unified data, analytics core

⑥ Data lake

⑤ Curated data vault (by domain)

Analytics

Raw data vault

② Real-time streaming ③

Systems of record

Core processing systems

① Cloud-enabled data platforms and services (serverless, noOps)

① From on-premise to cloud-based data platforms

② From batch to real-time data processing

③ From pre-integrated commercial solutions to modular, best-of-breed platforms

④ From point-to-point to decoupled data access

⑤ From an enterprise warehouse to domain-based architecture

⑥ From rigid data models to flexible, extensible data schemas

dm

**Practical approaches about how to make your data platform a success**

# About us

**Jelle De Vleminck**

Senior data engineer focused on making other engineers more productive

**Kristof Martens**

Senior data engineer focused on building secure, at scale data platforms in the cloud for organizations across different industries.

dm

# Business and IT are often acting like rivaling superpowers during the cold war



**Business (decentralised)**



**IT department (centralised)**

# A difference in values

**Business**

**IT department**

*Freedom*

*Speed*
*Agility*
*Experiments*
*User Experience*
*Ease of use*
*Self Service*

CONTROL

Cost
Process & Procedures
Reuse
Standardization
Operations
Security, legal, compliance
Oversight

# They carry dangerous weapons



**Business**



**IT department**

# Mutually Assured Destruction



**Business**



**IT department**

# They build walls



**Business**



**IT department**

CONTROL

Yearly budget exercise
Maintain what we built
More functionality
Priority calls
Escalations



Freedom

Tickets
Queues
Inflated estimates (time / budget)
Compartmentalisation
No (budget/manpower)
Committees

dm

# As a consequence

**Business**

**IT department**

Slow and expensive

Cowboys

Let's build it ourselves!

Why is nobody using our data platform?

# In the end, we all want the same



**Business**



**IT department**



Peace Love and Harmony



WE SHARE THE SAME BIOLOGY

REGARDLESS OF IDEOLOGY

# How to break down walls and create mutual understanding (in data)

# Instructions

Go to

# www.menti.com

Enter the code

# 5932 4620

Or use QR code

# Practical approaches about how to make your data platform a success

**1.** **Tie your platform to the strategic goals of the company**
Look for reasons to do things instead of finding reasons not to do it

**2.** **Embrace a platform thinking approach**
The platform mindset and a roadmap on how to get there

**3.** **How do you get started?**
Practical tips on how to approach this

**4.** **Balancing the act of centralization and decentralization**
Values of Freedom and Control are not conflicting

# Link data platform design to your strategic goals

# Leverage your current company goals



**Improve existing operations**

**Create innovative services**

**Become more data driven**

**Get closer to customers**

**All take a huge advantage of data analytics and data products**

**BUT...**

How can we **leverage all data** that we have available to us?

How can we make sure that data projects are developed in a **technical sustainable and scalable** way?

How can the business **focus on business challenges** instead of engineering and IT operational issues inherent to the data worlds?

# Every company is different

| | TECH | 🏢 | 🏛️ |
|---|---|---|---|
| Budget | ✅ | ✅✅✅✅ | ✅✅ |
| Synergy and reuse | ✅ | ✅✅✅✅ | ✅✅✅ |
| Agility | ✅✅✅✅ | ✅✅ | ✅ |
| Tech skills across organisation | ✅✅✅✅ | ✅✅ | ✅ |
| Legal / Governance requirements | ✅ | ✅✅✅ | ✅✅✅✅ |
| Data platform | ❌ | ✅ (Shared) | ✅ (# Instances) |

dm

# Change how business and IT interact with each other

From delegating responsibilities to other teams

|  | Use case 1 | Use case 2 | Use case 3 | Use case 4 |
|---|---|---|---|---|
|  | App team 1 | App team 2 | App team 3 | App team 4 |
|  | DE team | Data ingest team | Data science team | CI/CD |
|  | Cloud | Infra | DB | Security |
|  |  | Ops |  |  |

# Change how business and IT interact with each other

To offering capabilities that enable teams to work effectively

IT as enabler instead of bottleneck

| Use case 1 | Use case 2 | Use case 3 | Use case 4 |
|---|---|---|---|
| App team 1 | App team 2 | App team 3 | App team 4 |

Data platform team
Offering integrated E2E workflow

Web platform team

Cloud

Make sure you are building the right capabilities

dm

# Beware Big Design Upfront & Analysis Paralysis

"No architecture survives first contact with developers"



This is a wish list, not an architecture!



Accept from the beginning that your architecture will evolve, based on what you will learn. Start small, show value, lure in others!

*"Build the plane while flying it!"*

# Embrace a platform thinking approach

# Classic IT departments often operate in silos with different teams owning the different key layers of the end product



**Ops team**

Operations layer

**Frontend team**

Presentation layer

**Development team**

Business logic layer

**Database Admin team**

Data access layer

🤨 **Forgetting about the end-users.**

🧭 **Limited ownership and accountability on decisions.**

🚢 **Lack of agility & business adaptability.**

# Moving towards long-lived product-oriented teams where people with different expertise work together towards a shared goal



Creating value by refining data products

Team Apple
use case team

Team Pear
use case team

Team Orange
use case team

Operations layer

Presentation layer

Business logic layer

Data access layer

🎁 **A stronger feeling of ownership.**

🚀 **Optimised for flow of change.**

🏃‍♀️ **More motivation in day-to-day tasks**

💪 **Developing skills enabling end-to-end value creation.**

**The main issue with autonomous product-oriented teams is that teams will spend a lot of their day-to-day work on tasks not directly impacting business value. We call that <u>undifferentiated heavy lifting</u>.**

**This is where the idea of a <u>platform</u> and enabling teams comes in. They support the product-oriented teams in doing their job.**

# Make the easy things the rights things by providing paved roads



## Data product lifecycle

A data platform is not a bunch of tools randomly thrown together

A data platform should offer an *integrated experience* that guides them through all steps of the data product lifecycle by offering *paved roads*

# Don't force people to use the platform



**But make them want to use it.
Involve them in how the platform needs to evolve.**

# How do you get started?

# From close collaboration to limited collaboration (discovery) through to X-as-a-Service for established, predictable delivery

# An example: Going from documenting API calls to experimentation as a service

```bash
bash                                    Copy code

aws sagemaker create-notebook-instance \
    --notebook-instance-name my-notebook-instance \
    --instance-type ml.t2.medium \
    --role-arn arn:aws:iam::123456789012:role/SageMaker-Execution-Role \
    --subnet-id subnet-12345 \
    --security-group-ids sg-12345 \
    --lifecycle-config-name my-lifecycle-config \
    --tags Key=Environment,Value=Development
```

**Create a new notebook**                                      ✕

\* 🏷 Notebook name : `meaningful-ostrich`

\* 📁 Project : `imbalance_price_forecast` ⌄

\* 🔲 Environment : `dev` ⌄

\* 🛡 IAM identity ⓘ : `neo-iam-imbalancepriceforecast-dev` ⟳

⌄ Advanced options

🐍 Python version : `3.10` ⌄

🌐 Instance type : `mx.2xlarge` ⌄

🔋 Instance lifecycle : `spot` ⌄

🕐 Max idle time (60 minutes) : ——————◯——————

🗐 Disk size (10 Gb) : —◯————————

[Cancel]  [Create]

# Put yourself in your users shoes

| Journey Steps — Which step of the experience are you describing? | Post a question — Why do they even start the journey? | | Triage and help — How can they feel successful? | | | User gets an answer — How can we make them feel satisfied? | | |
|---|---|---|---|---|---|---|---|---|
| **Actions** — What does the customer do? What information do they look for? What is their context? | Post a question in support channel | Fill in a form in the workflow | Provide context | Read relevant links | Try suggested solution | Open document | | |
| **Needs and Pains** — What does the customer want to achieve or avoid? *Tip: Reduce ambiguity, e.g. by using the first person narrator.* | I want my problem resolved quickly / I don't fiddle with unfamiliar controls | I worry about asking a silly question / I don't waste time reading manuals | I understand how this can help me get my job done | I want to be able to track progress | I worry about delays | I can start creating right away | My problem is fixed | I learn from how others do it |
| **Touchpoint** — What part of the service do they interact with? | support channel | support workflow | Jira Service Desk ticket | Support Engineer | | [the / support channel] | Support Engineer | Support Engineer / PM / TL |
| **Customer Feeling** — What is the customer feeling? *Tip: Use the **emoji app** to express more emotions* | 🤯 | | 🤔 | | | 🥳 | | |
| *Backstage* | | | | | | | | |
| **Opportunities** — What could we improve or introduce? | Send a welcome message / reassure | If the request is not relevant: give the link to the responsibilities doc | Post all updates from the ticket into a thread | Auto-track the question and the conversation in a ticket, post the ticket number right | | Make sure we provided all the relevant links to enable the customer to solve the problem | | |

# Self-service documentation is a <u>first class deliverable</u>.

**Platform Documentation**

- 🚀 Get started
  - › Development tools setup.
- 👷 Purpose-Based Access Control
- › 🌊 Storing & querying data from the data lake
- 📥 Data I/O with ingress & egress buckets
- 🧨 DAGs monitoring & alerting
- › ❄️ Data warehousing with Snowflake
- › 🦊 Version control & CI/CD with GitLab
- › 🕊️ Realtime data layer with Kafka.
- › 📚 Data catalog & schema registry
- › 💂 Data Quality
- › 🧪 Development & experimentation environments
- 🗄️ Self-service database
- › 🐳 Hosting containerized applications on AWS EKS
- › 🧱 Shared resources

🔍 **Audit regularly.**

💡 **Make it simple.**

🧩 **Have an onboarding exercise.**

🖌️ **Have consistent formats.**

📥 **Collect feedback on every page.**

dm

# Balancing the act of centralization and decentralization

# Data Mesh or Data Mess?



## What Is Data Mesh?



datamesh-architecture.com

## Data Mesh Architecture



datamesh-architecture.com

**Data mesh is not a religion**: Don't be dogmatic in your principles and keep a flexible mindset.

**Take a practical approach**: See what works and skip what doesn't work for you

**Do not decentralise everything**: Data mesh principles allow for shared capabilities

**?** **Be careful of ambiguous terminology**: What is a data domain and how granular is it?

**Processing vs Data ownership**: Who "owns" a derived data product using data from other domains?

# Introduce simple concepts that are repeatable and easy to understand

Unit of deploy ment

Has own git repo

Has data inputs & outputs

Data Product

SLA's owned by team

Has container artefact

Compute & Schedulin g

Where DP's are deployed to

Fine grained access to DP's

DEV, ACC & PRD separated from each other

Data Environ ment

Multiple storage techs

Identical *worst case* security setup

Read from higher envs?

# Infrastructure environment

dm

# Decide upfront how you want to offer the capabilities of your data platform

| Data platform capability | Centralised | Self service | Decentralised |
|---|---|---|---|
| Ingestion | | ● | |
| Processing | ● | | |
| Access & sharing | ● | | |
| Security | ● | | |
| ML / AI capabilities | | ● | |
| BI capabilities | | ● | |
| Data modelling | | | ● |
| Dashboards | | ● | |
| Build/release/deploy | | | ● |
| API hosting | | ● | |
| Data product serving | | ● | |
| Governance | ● | | |

# An example (AWS Focussed)



**Decentralised Capabilities**

**Self Service Capabilities**

**Centralised Capabilities**

Can access Data Product specific IAM permissions with scoped data access

Shared Data Platform AWS Account

| Data storage | Data Ingest | Data Processing | Data Science | CICD | DP registration | Data & User access | Data security | Infra / compute | Costing |
|---|---|---|---|---|---|---|---|---|---|

Operational Domain AWS Account

API

Data Product 1

Operational Domain On-Premise

salesforce   SAP

Data Product 2

DP 3

Data engineer / scientists working on data product

**Domain or Business team**

**Paved roads provided by Data Platform (or Enabling) team**

**Data Platform team**

dm

# Pro tip: Leverage Kubernetes to manage separation



**Decentralised Capabilities**

Can access Data Product specific IAM permissions with scoped data access

Operational Domain AWS Account

Data Product 1

Operational Domain On-Premise

Data Product 2

DP 3

Data engineer / scientists working on data product

**Domain or Business team**

**Self Service Capabilities**

Shared Data Platform AWS Account

| Data storage | Data Ingest | Data Process-ing | Data Science | CICD |

**Paved roads provided by Data Platform (or Enabling) team**

**Centralised Capabilities**

| DP regis tration | Data & User access | Data security | Infra / compute | Costing |

**Data Platform team**

# Pro tip: Leverage automation to do this at scale and integrate all components

**Decentralised Capabilities**

**Self Service Capabilities**

**Centralised Capabilities**

Can access Data Product specific IAM permissions with scoped data access

**aws** Shared Data Platform AWS Account

| Data storage | Data Ingest | Data Process-ing | Data Science | CICD | DP regis tration | Data & User access | Data security | Infra / compute | Costing |
|---|---|---|---|---|---|---|---|---|---|

**aws** Operational Domain AWS Account

API

Data Product 1

Operational Domain On-Premise

salesforce

SAP

Data Product 2

DP 3

Data engineer / scientists working on data product

**Domain or Business team**

**Paved roads provided by Data Platform (or Enabling) team**

**Data Platform team**

dm

# If done properly you can strike the right balance...

**Centralised**



**Self Service**

**Decentralised**

# Example of a paved road

Idea　　Experiment　　Build　　Deploy　　Run　　Scale　　Product

Browse available data sources in the data catalog

Register a new potential data product purpose & request data access.



Data owners validate the purpose and whether data can be used for that.
After approval everything is set up automatically.

Ideally offer an integrated user experience for data scientists/engineers/analysts

Experiment with data via one-click access to notebooks scoped to that project

Industrialise via Cloud IDE's that are scoped to that project

Build and deploy via standardized tooling and templates offered via paved roads

Schedule your data pipelines via Airflow or a different scheduler

Follow up on alerts and offer easy access to logs and metrics if needed

# It is possible to build this yourselves, or give yourselves a head start

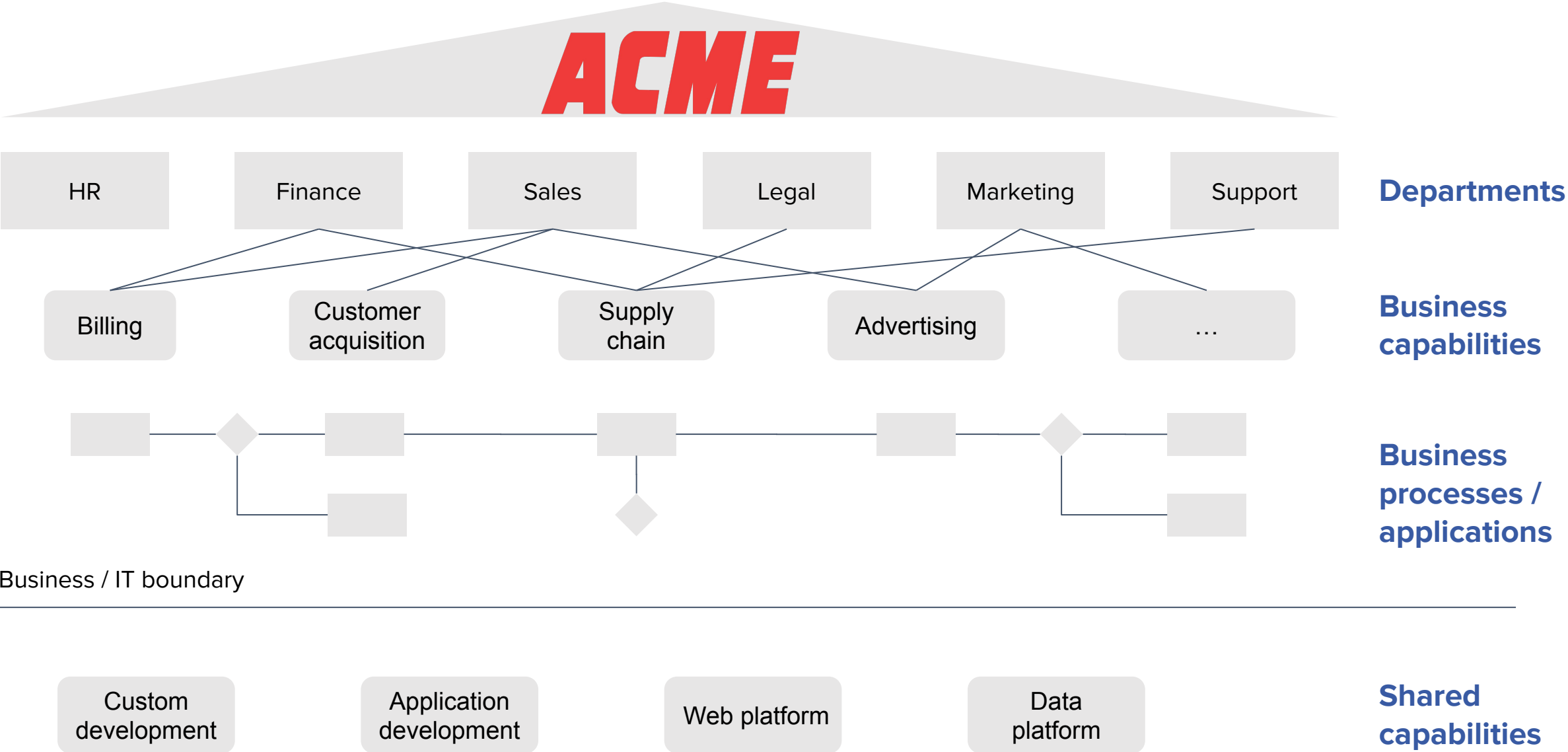# Practical approaches about how to make your data platform a success

**1.**  **Tie your platform to the strategic goals of the company**
Look for reasons to do things instead of finding reasons not to do it

**2.**  **Embrace a platform thinking approach**
The platform mindset and a roadmap on how to get there

**3.**  **How do you get started?**
Practical tips on how to approach this

**4.**  **Balancing the act of centralization and decentralization**
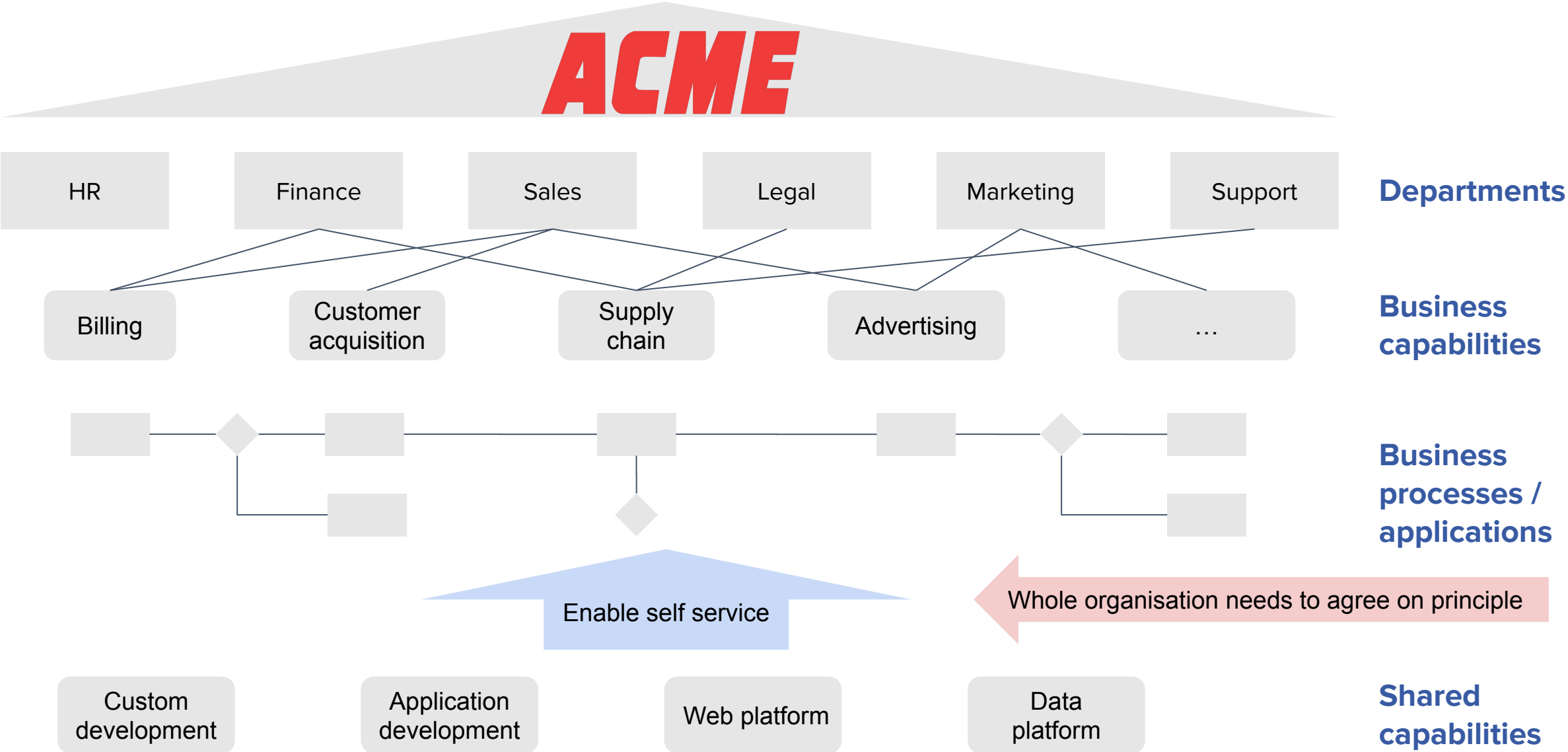Values of Freedom and Control are not conflicting
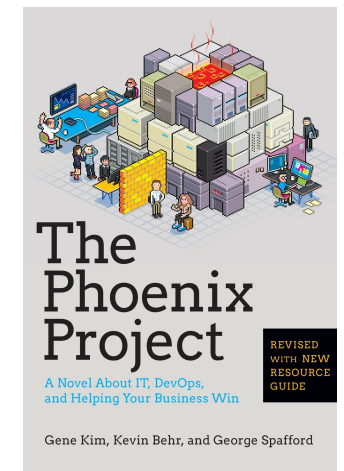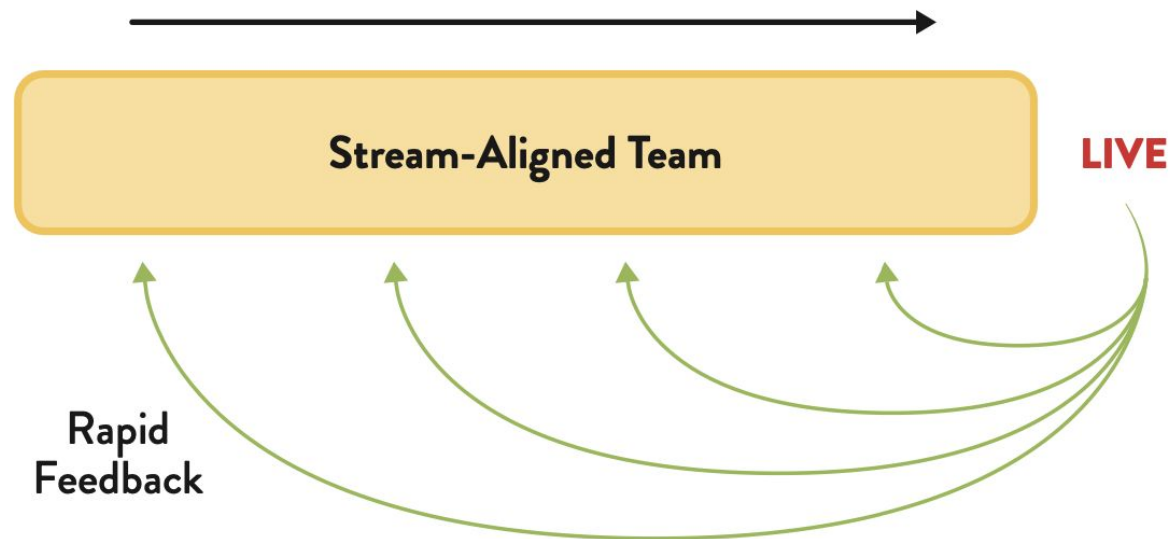
dm

The end

# Backup slides

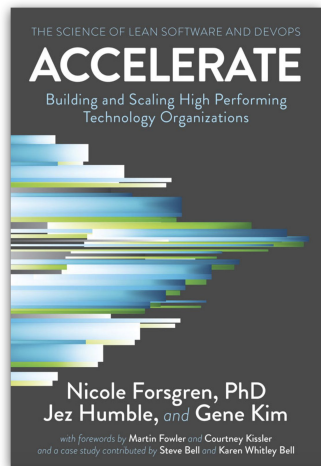# Make sure you are building capabilities the business needs

# Make sure you are building capabilities the business needs

"we must ensure delivery teams are cross-functional, with all the skills necessary to design, develop, test, deploy, and operate the system on the same team."

Stream-Aligned Team

LIVE

Rapid Feedback

# Pro tip: own the perception of how teams look at you as a platform team

**Support your users by defining your commitments, support channels and support workflow for each type of request**
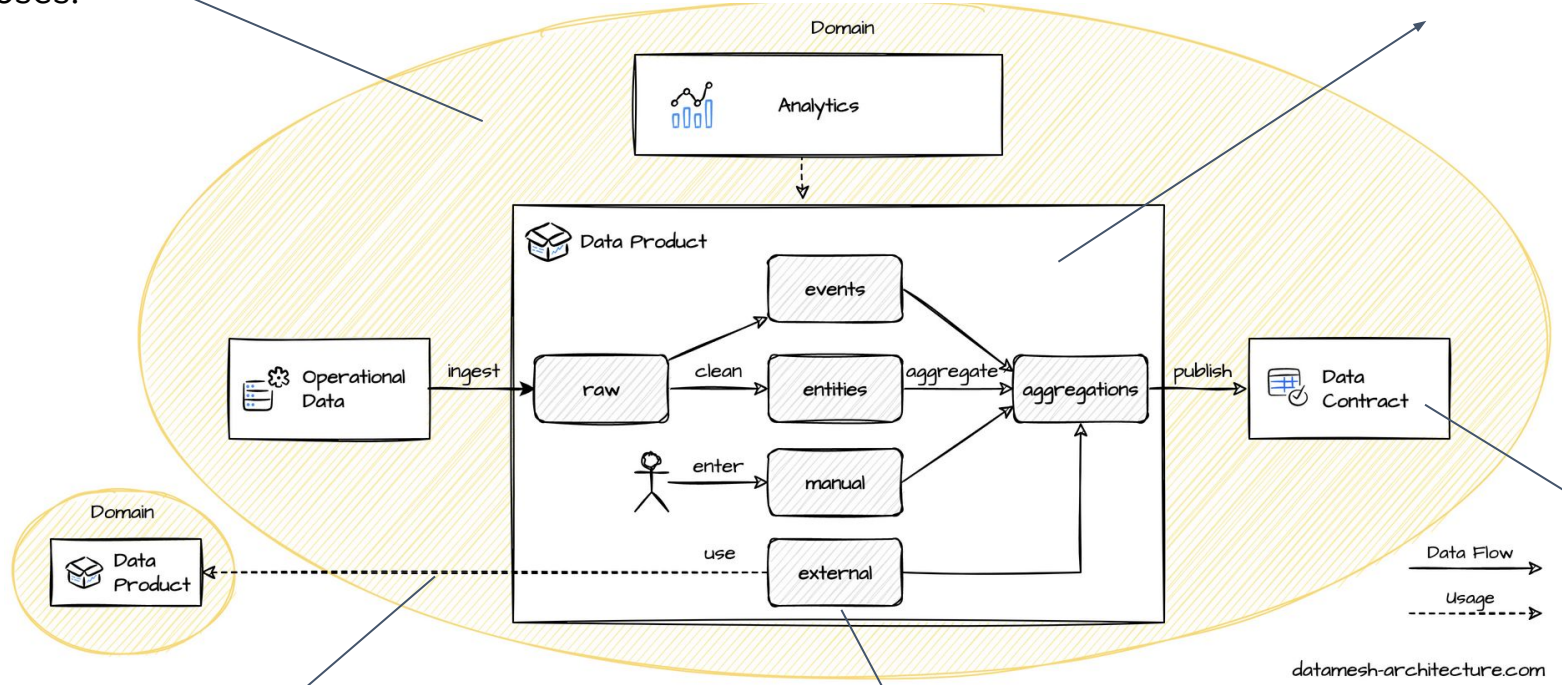
# Data Mesh is not an excuse to decentralise everything (AWS focused)



Domain is not the equivalent of AWS account, especially for analytical purposes!

Capabilities to build this can be offered by the data platform team via paved roads.

There is no standardized way to share data contracts

Please don't use dedicated API's for all X-domain interactions! S3/Glue, Athena/DWH SQL,... are also API's.

Data can be stored in shared data stores (S3/DWH's). This makes it easier to share data across data products/domains and to standardise data access patterns

datamesh-architecture.com