

# Collaboration Challenges in Building ML-Enabled Systems:

Communication, Documentation, Engineering, and Process



**Nadia Nahar\***



Shurui Zhou



Grace Lewis



Christian Kästner



# Machine Learning (ML) Component

```
face_detection.ipynb
File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text Copy to Drive

[6] print("[INFO] loading model...")
    prototxt = 'deploy.prototxt'
    model = 'res10_300x300_ssd.xml'
    net = cv2.dnn.readNetFromCaffe(prototxt, model)

    [INFO] loading model...

Use the dnn.blobFromImage function

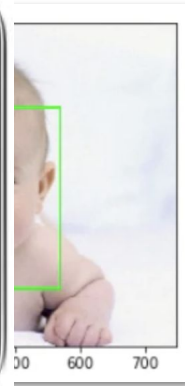
[7] # resize it to have a max dimension of 300 pixels
    image = imutils.resize(image, width=300)
    blob = cv2.dnn.blobFromImage(image, 1.0, (300, 300), (104, 117, 123))

Pass the blob through the neural network

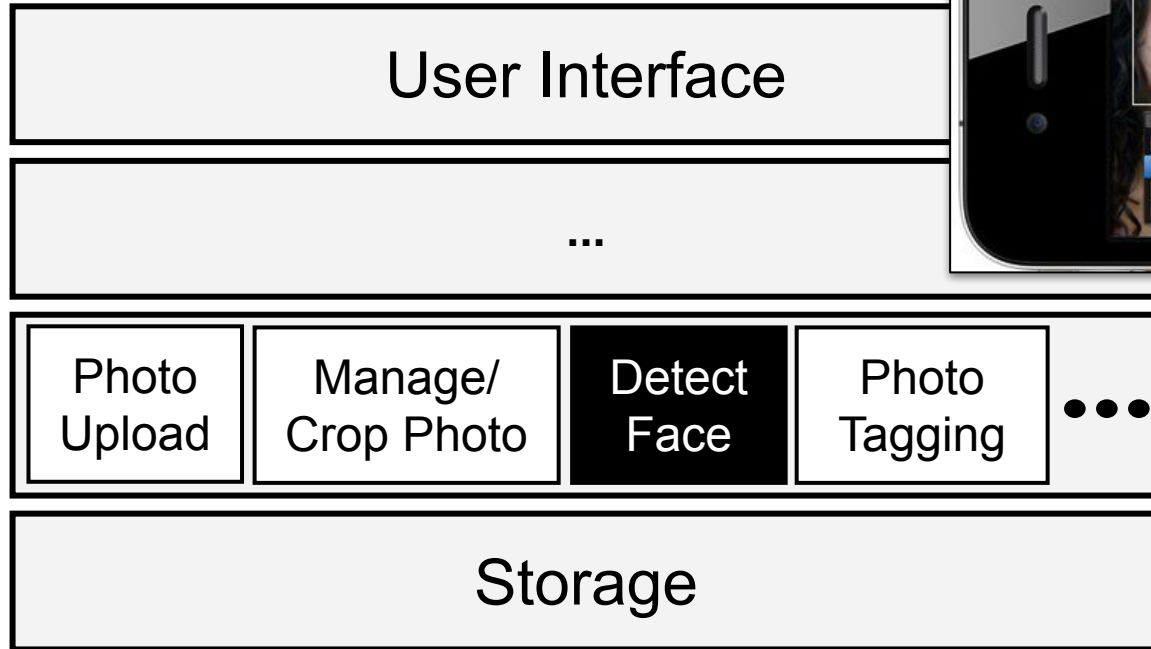
[8] print("[INFO] computing object detections...")
    net.setInput(blob)
    detections = net.forward()

    [INFO] computing object detections...

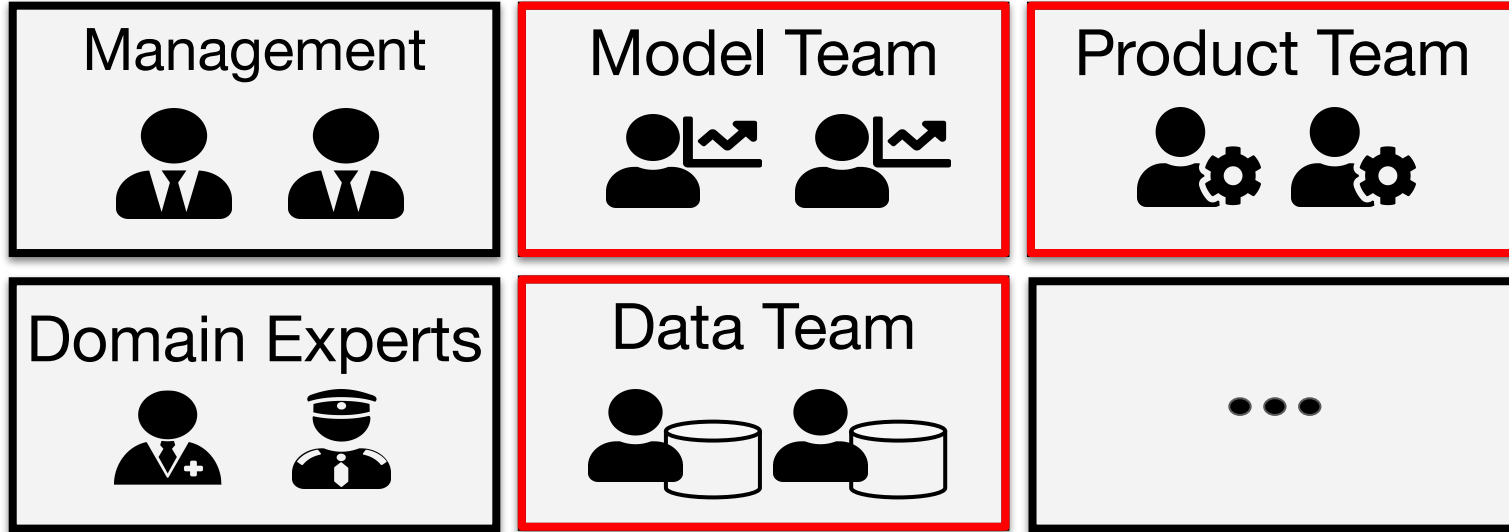
Loop over the detections and draw boxes around the detected faces
```



# ML-Enabled System



# Multiple Teams Collaborating Together



# Research Question

*“What are the collaboration points and corresponding challenges between data scientists and software engineers in building ML-enabled systems?”*

# Why do 87% of data science projects never make it into production?

Collaboration Problems



VB Staff

July 19, 2019 4:

And the third issue, intimately connected to those silos, is the lack of collaboration. Data scientists have been around since the 1950s — and they were individuals sitting in a basement working behind a terminal. But now that it's a team sport, and the importance of that work is now being embedded into the fabric of the company, it's essential that every person on the team is able to collaborate with everyone else: the data engineers, the data stewards, people that understand the data science, or analytics, or BI specialists, all the way up to DevOps and engineering.

“This is a big place that holds companies back because they're not used to collaborating in this way,” Leff says. “Because when they take those insights, and they flip them over the wall, now you're asking an engineer to rewrite a data science model created by a data scientist, how's that work out, usually?”

# WHY DO MACHINE LEARNING PROJECTS FAIL?

Think ahead to production so that you don't let your machine learning project collapse before it even gets started.



Rahul Agarwal  
| Expert Columnist

Agarwal is a senior data scientist currently working with Wal

## 4. YOUR MODEL MIGHT NOT EVEN GO TO PRODUCTION

Let's imagine that you've created this impressive machine learning model. It gives 90 percent accuracy, but it takes around 10 seconds to fetch a prediction. Or maybe it takes a lot of resources to predict.

Is that acceptable? Most likely no.

**Mismatch in Assumptions**

# Top 10 Reasons Why 87% of Machine Learning Projects Fail

In this article, find out why 87% of machine learning projects fail.



by Prajeen MV · Oct. 13, 20 · AI Zone · Opinion

## A Disconnect Between Data Science and Traditional Software Development

A disconnect between Data Science and traditional Software development is another major factor. Traditional software development tends to be more predictable and measurable.

---

**However, Data science is still part-research and part-engineering.**

---




**Different Ways of Working**





# Frustrations shared in Twitter...

All ML projects which turned into a disaster in my career have a single common point:

 I didn't understand the business context first, got over-excited about the tech, and jumped into coding too early.

1:08 PM · Mar 12, 2022 · Twitter Web App

**297** Retweets **39** Quote Tweets **1,786** Likes

Machine Learning lives in an uncanny valley btw Science and Engineering.

It's the worst of both worlds.

We don't care about understanding, just making things "work" (bad science).

We don't care if things work in the real world, just on contrived benchmarks (bad engineering).

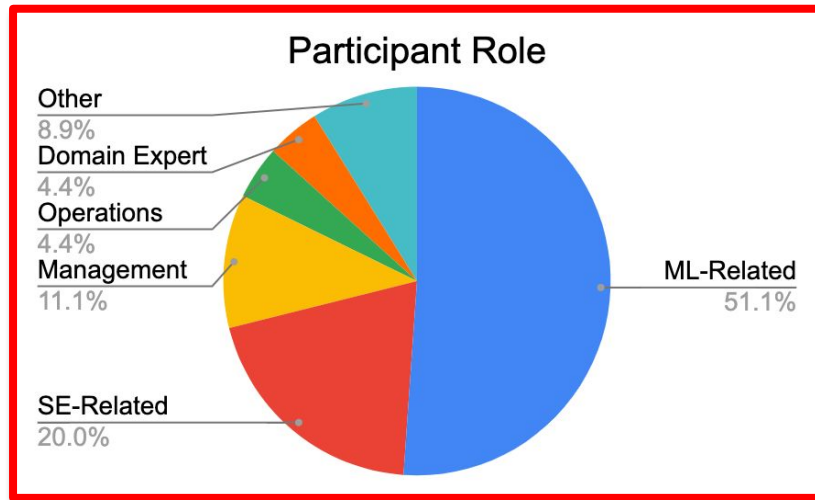
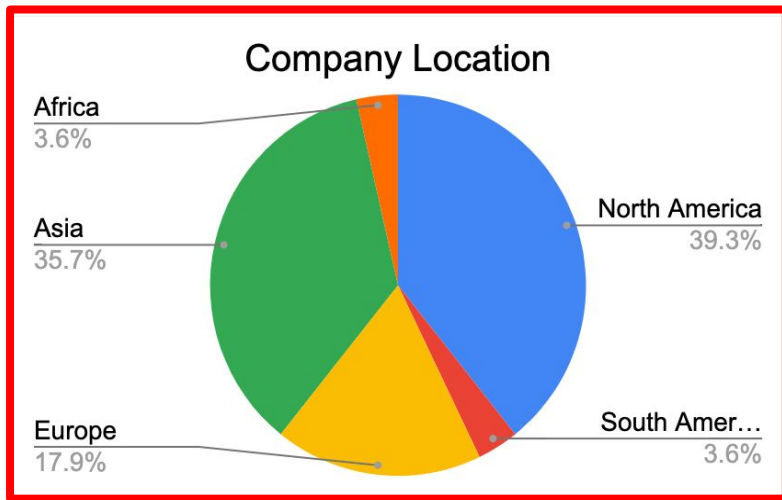
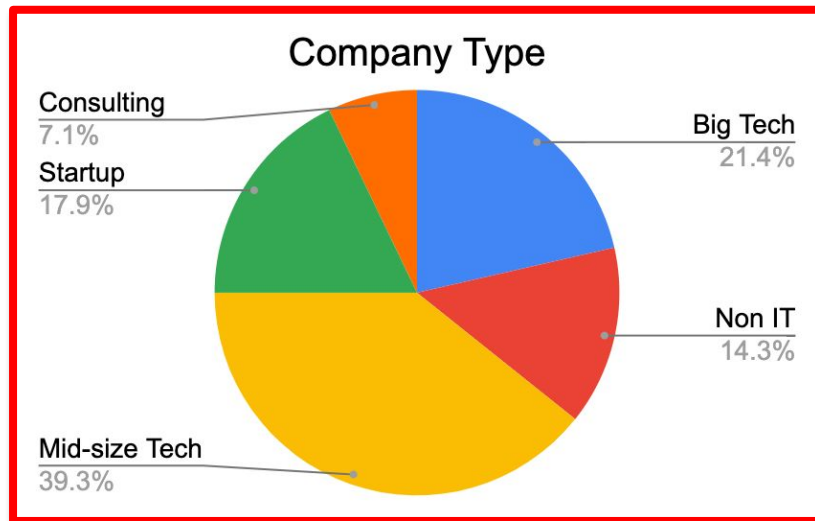
6:45 AM · Jan 29, 2022 · Twitter Web App

**202** Retweets **37** Quote Tweets **1,451** Likes

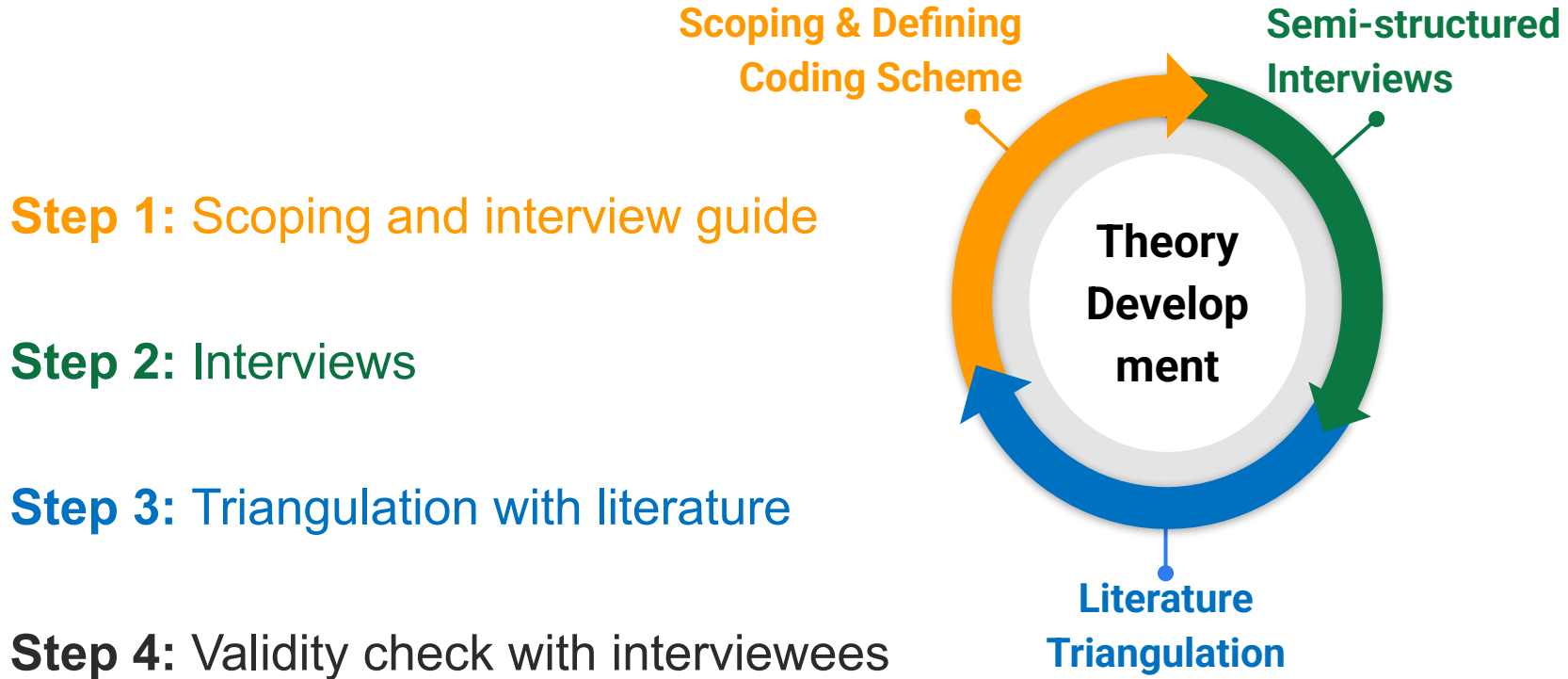
# Research Question

*“What are the **collaboration points** and corresponding **challenges** between data scientists and software engineers in building ML-enabled systems?”*

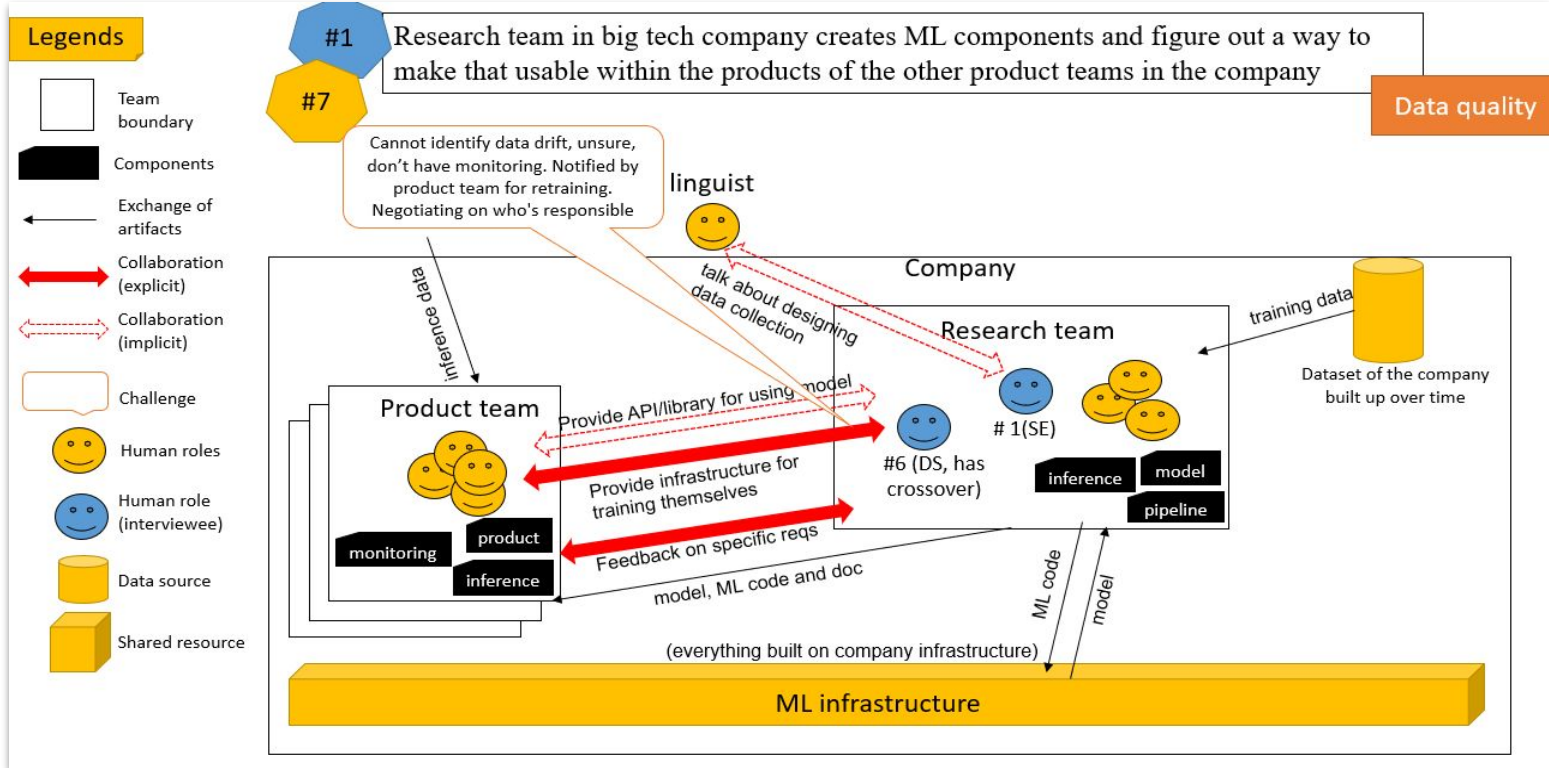
Conducted  
**45** interviews in  
**28** organizations




















# Qualitative Research



# Example Visual Analysis



<b>Requirements and Planning</b>				
-	Product and Model Requirements			 
-	Project Planning		-	- 
<b>Training Data</b>				
-	Negotiating Data Quality and Quantity			 
<b>Product-Model Integration</b>				
-	Responsibility and Cultural Clashes		-	 
-	Quality Assurance for Model and Product			 



Communication



Documentation

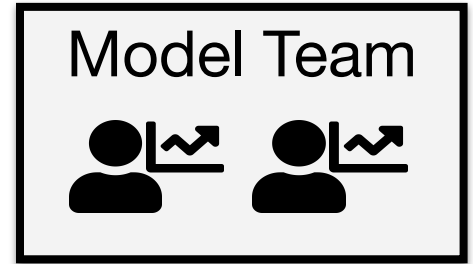
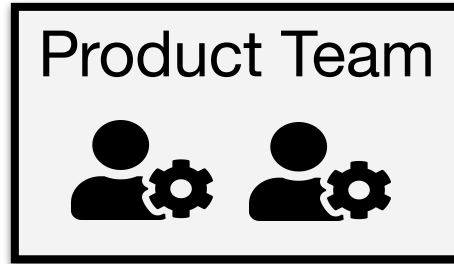


Engineering



Process

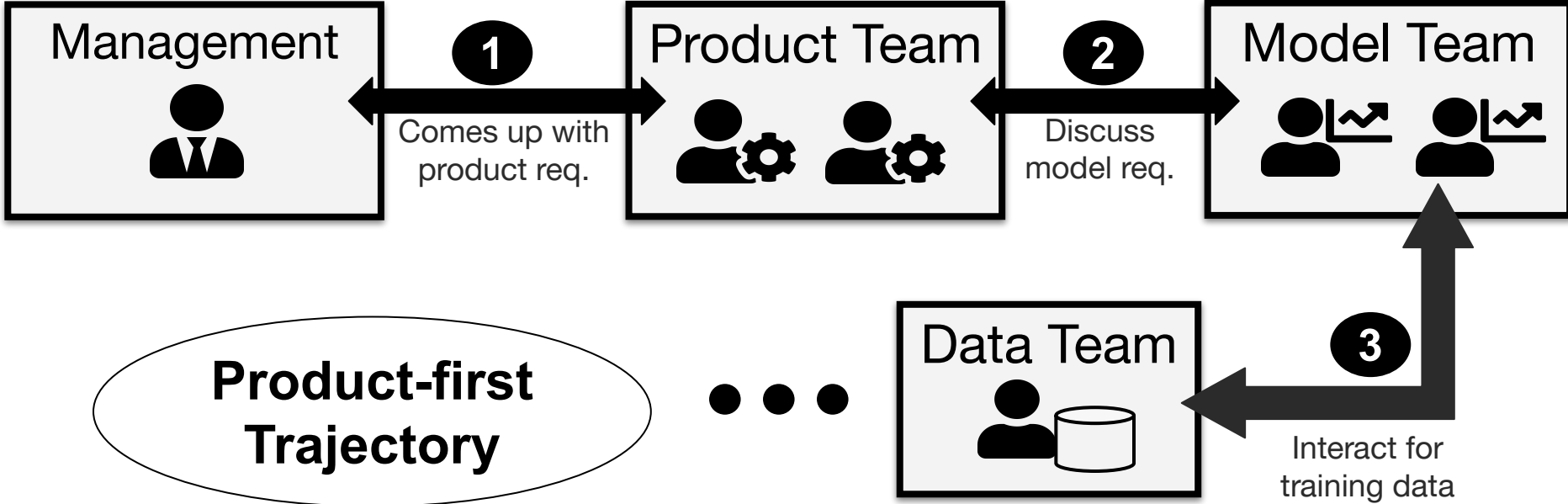
# Collaboration Point: Product and Model Requirements



Different patterns around different organizations.

Let's talk about **two example** orgs.

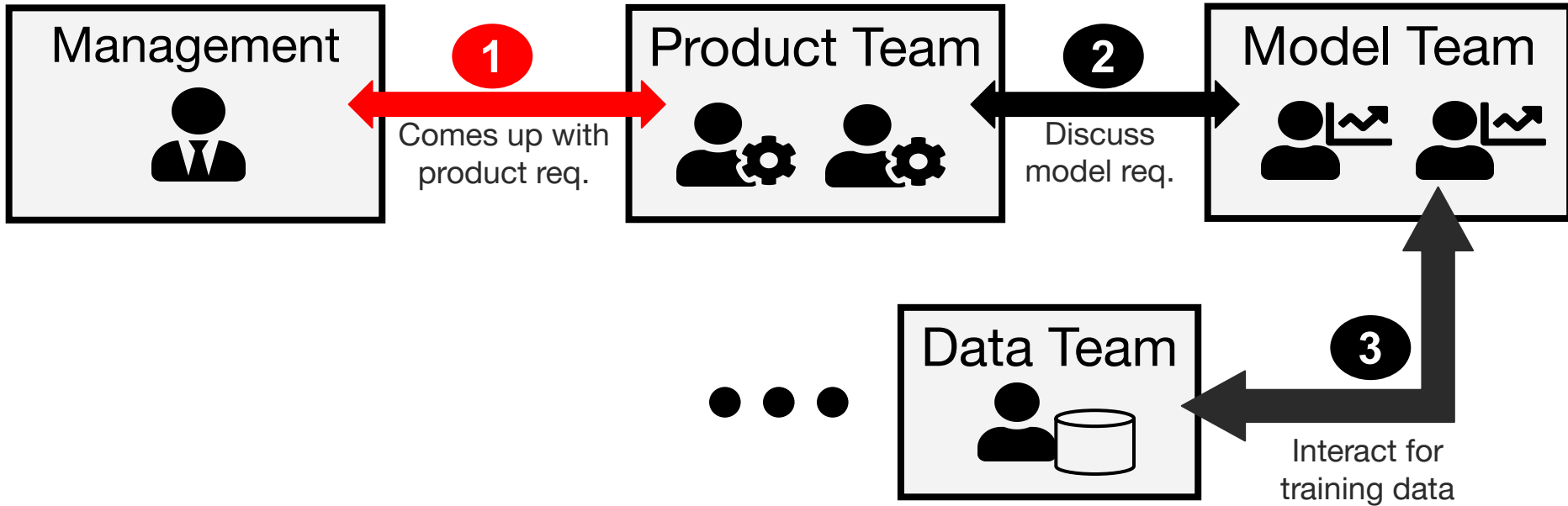
# Org. A: Fraud Detection in Banking Software





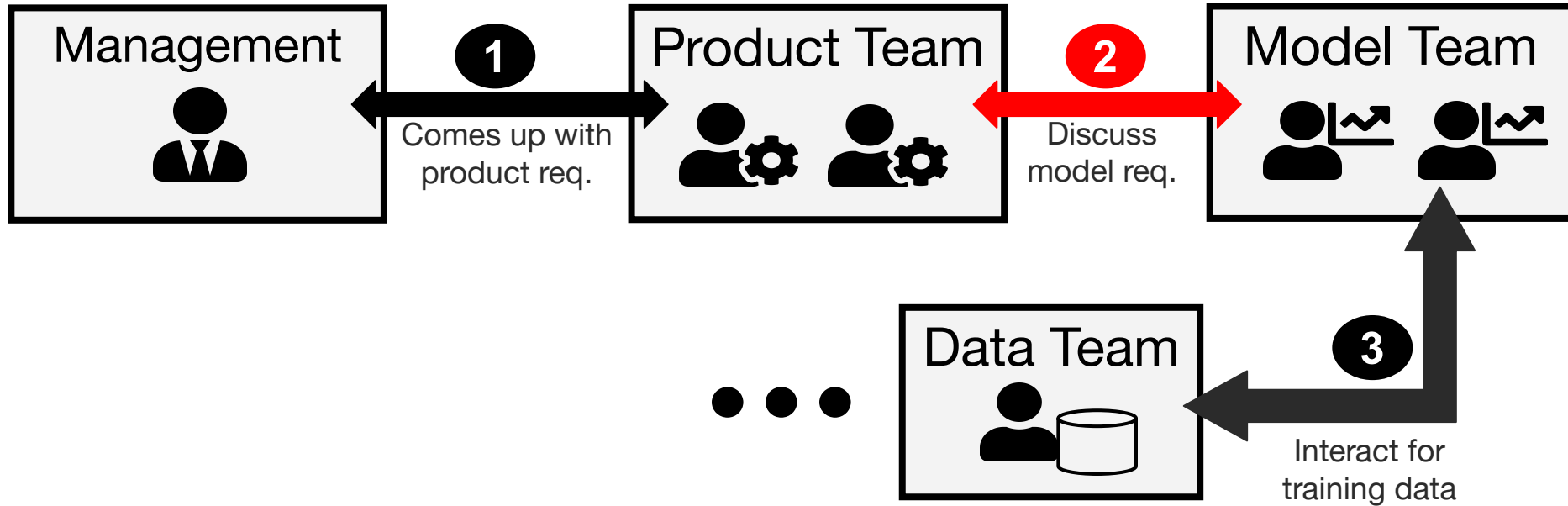


# Problem: Lack of ML Literacy Leads to Unrealistic Requirements





# Problem: Need Data Scientists to Set Correct Expectations





**Communication:** Lack of ML literacy leads to unrealistic requirements



Involving data scientists early when soliciting product requirements



**Documentation:** Product requirements are often not translated into clear model requirements



Adopt more formal requirements documentation for product and model



Communication



Documentation

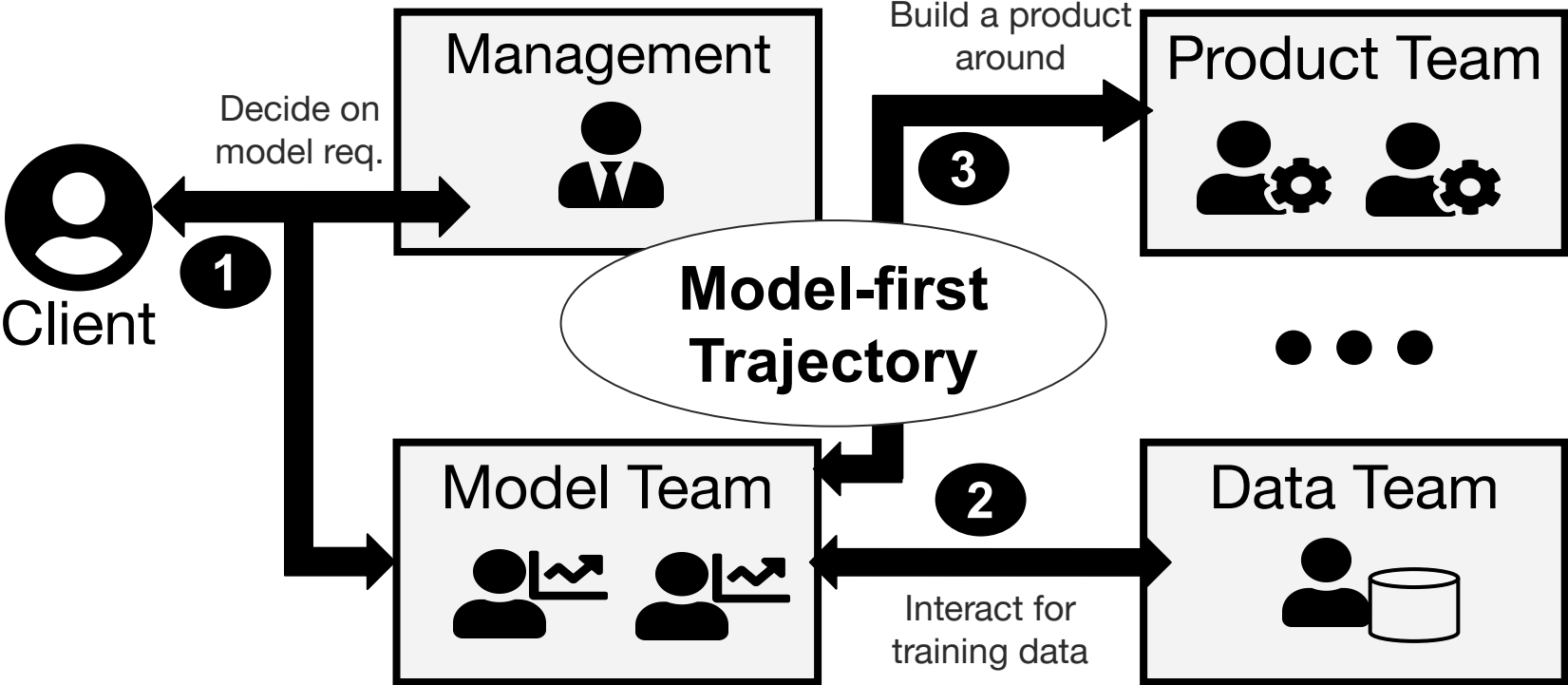


Engineering



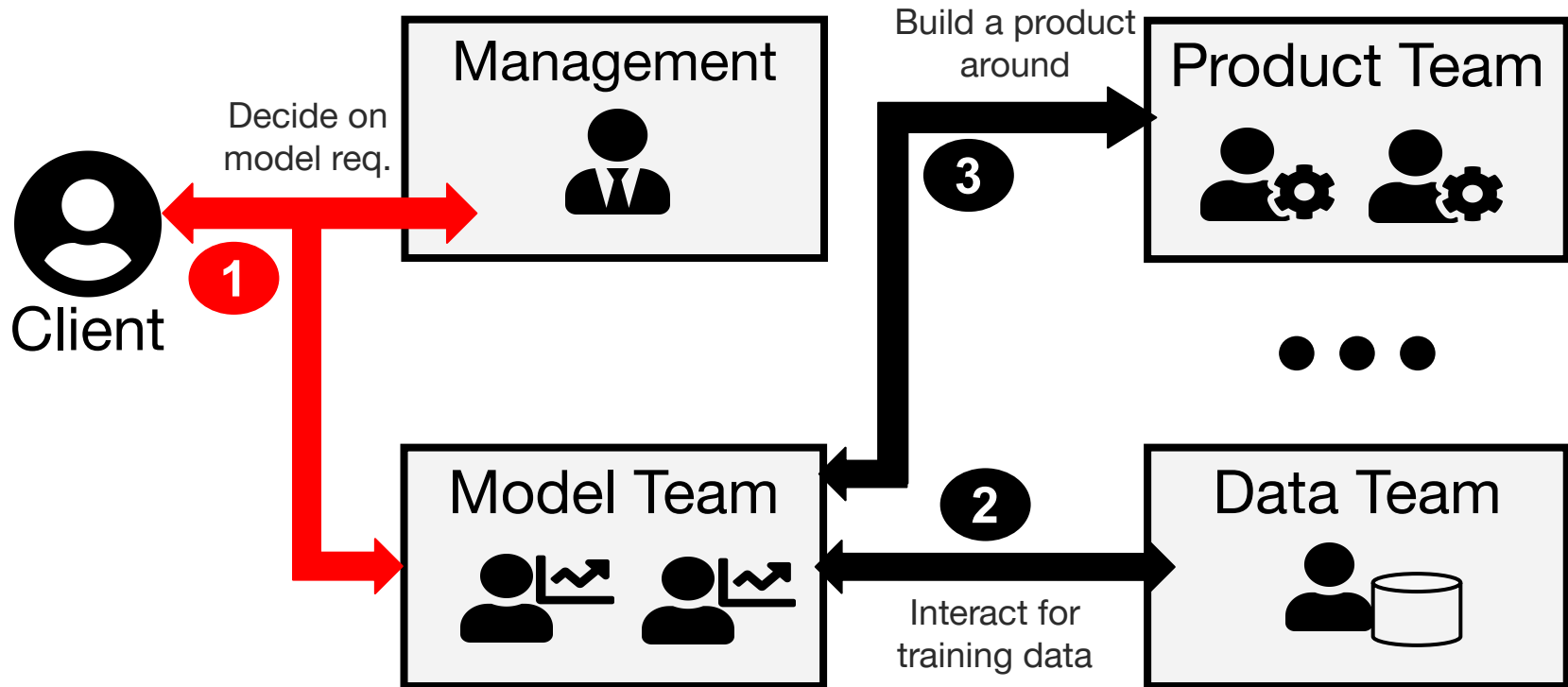
Process

# Org. B: Develop OCR for Local Language

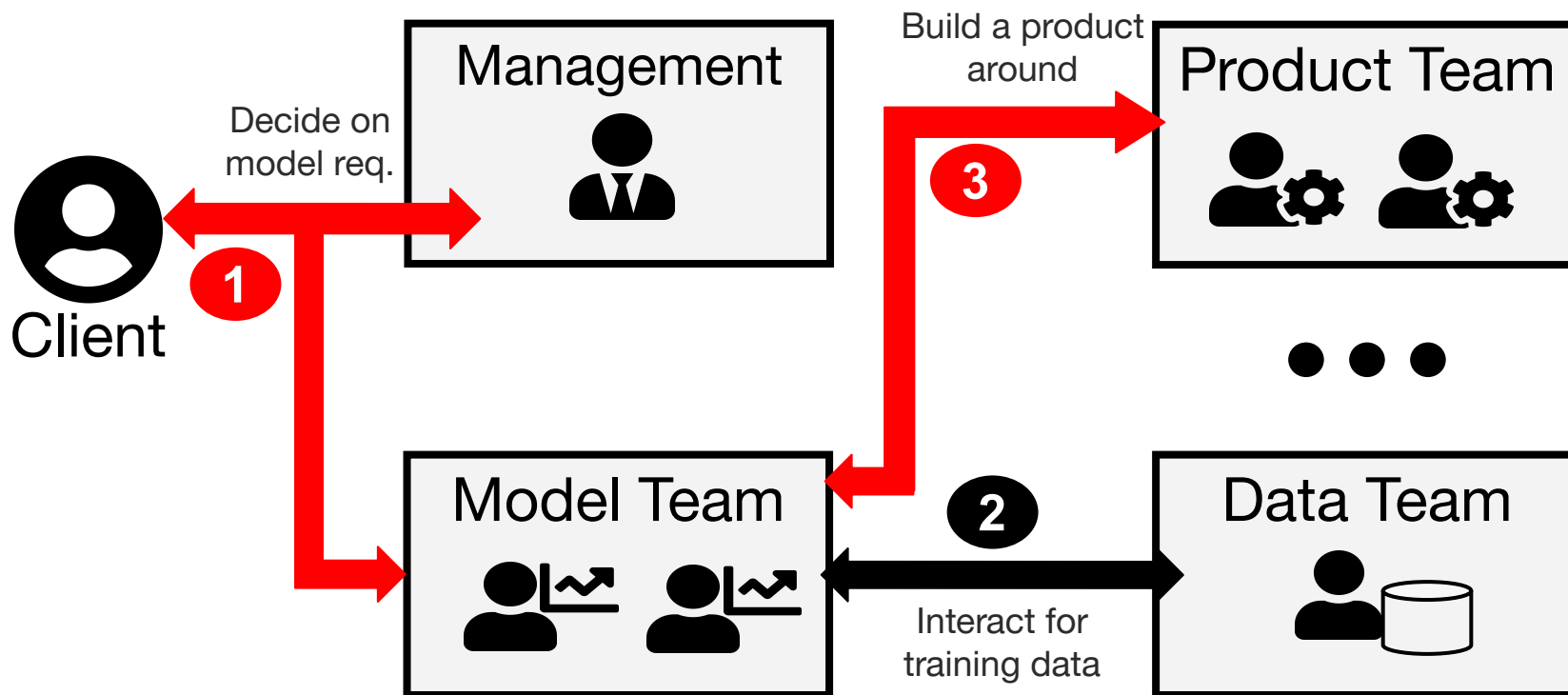




# Problem: Model Team Needs to Educate Client on ML (Less Impact)



# Problem: Less Focus on Entire Product





**Communication:** Model team needs to educate client on ML



ML literacy for customers and product teams: conducting training sessions



**Process:** Pursuing a model-first trajectory entirely without considering product requirements is problematic



Emphasis on collaboration during requirements phase, more research on process needed



Communication




















Documentation



Engineering



Process

Requirements and Planning	
- Product and Model Requirements	   
- Project Planning	 - - 
Training Data	
- Negotiating Data Quality and Quantity	   
Product-Model Integration	
- Responsibility and Cultural Clashes	 -  
- Quality Assurance for Model and Product	   



Communication



Documentation



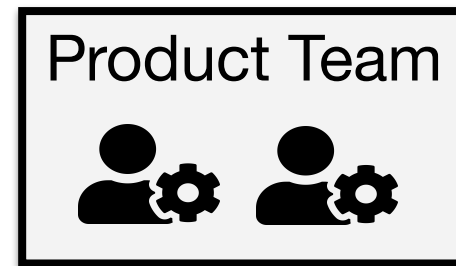
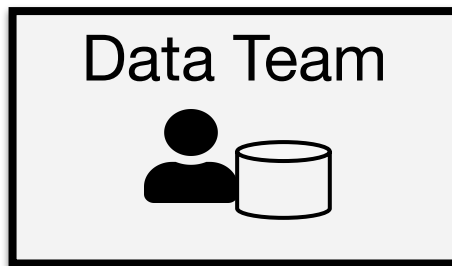
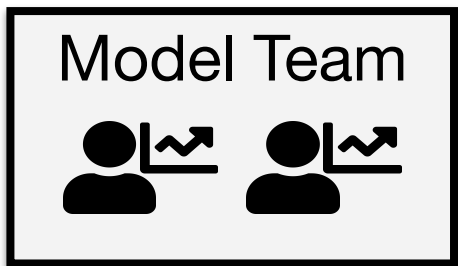
Engineering



Process



# Collaboration Point: Training Data



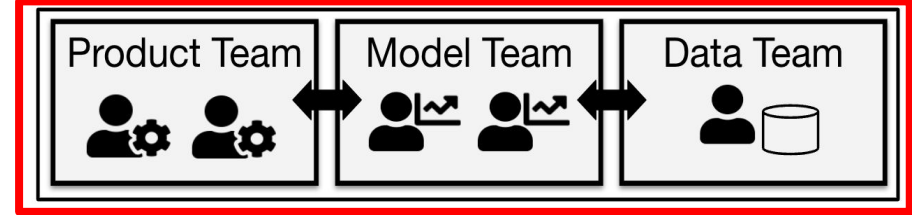
Again different patterns around different organizations.

# Three Collaboration Patterns Around Training Data

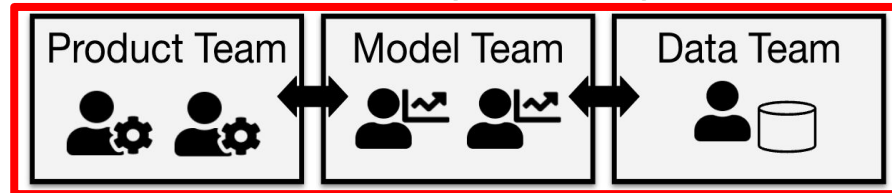
**Provided Data**



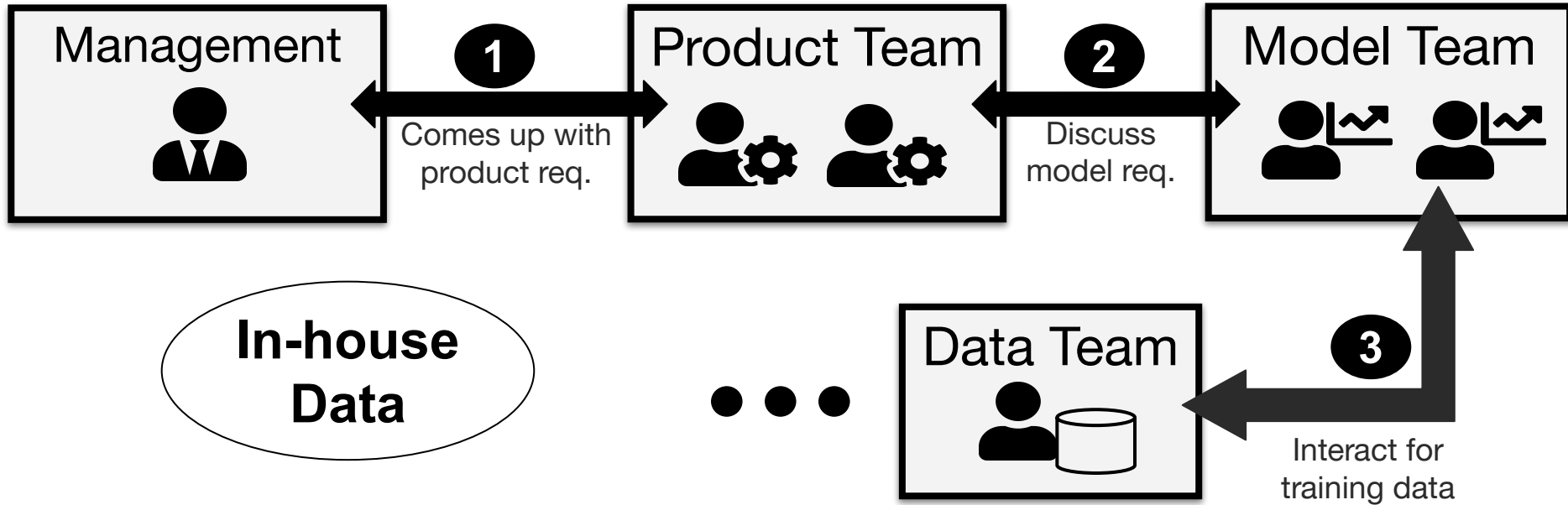
**In-house Data**



**External Data - i) Public, ii) Hired**

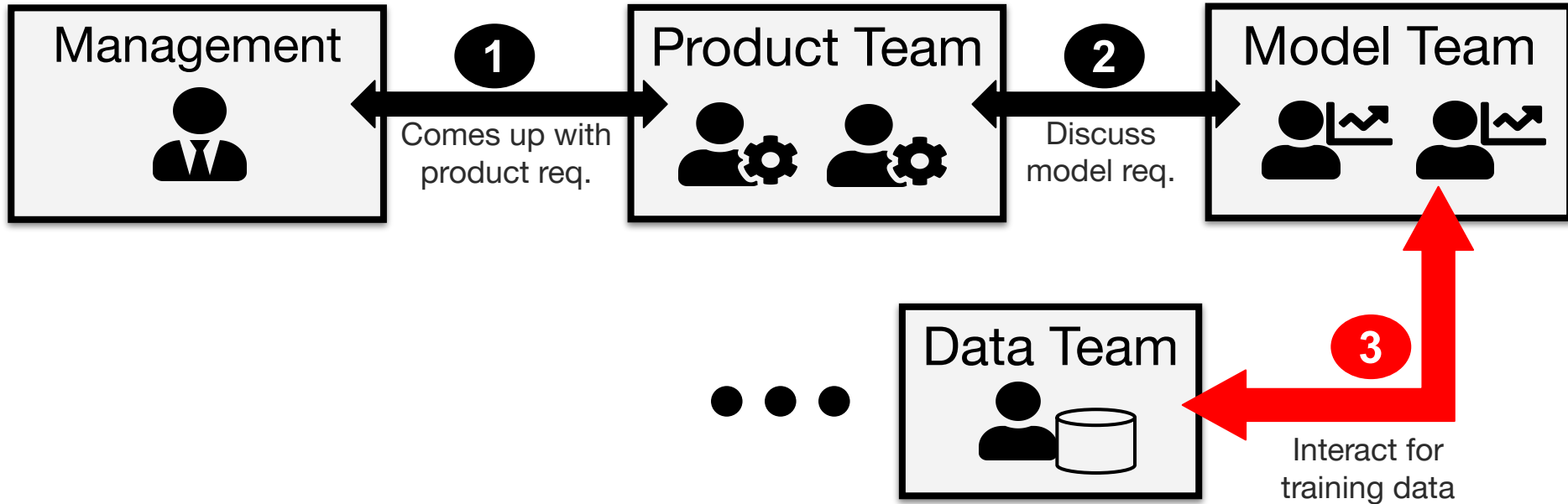


# Org. A: Fraud Detection in Banking Software

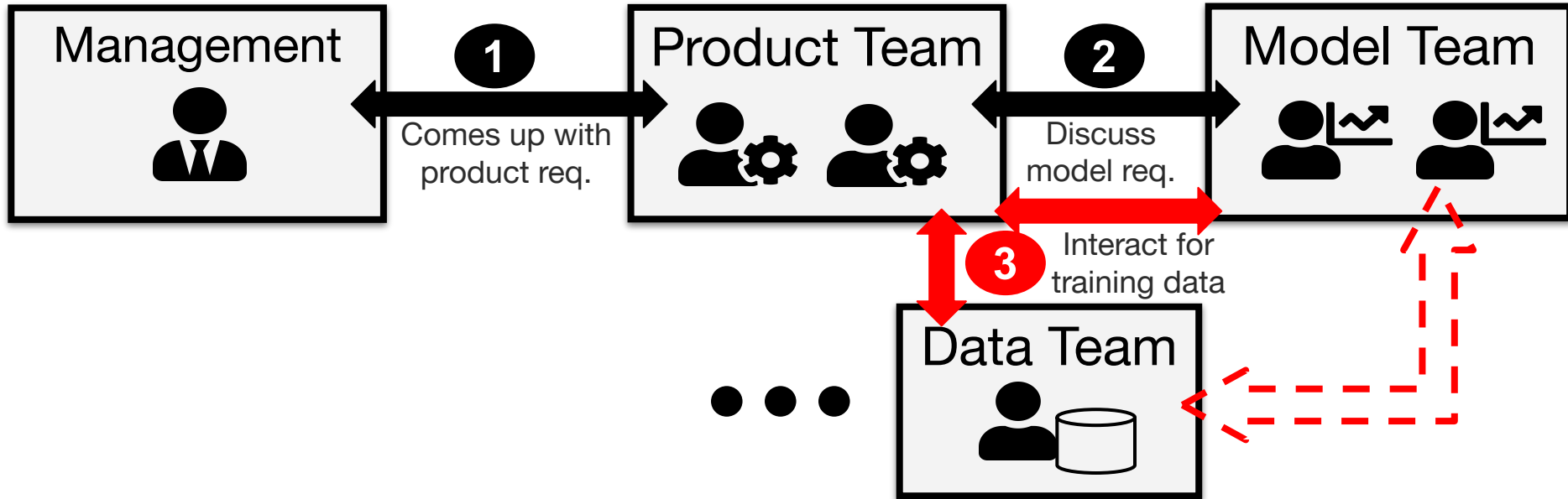




# Problem: Data Access Challenges Due to Power Dynamics



# Problem: Little Help with Data Understanding





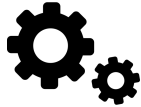
**Communication:** Data Access Challenges Due to Power Dynamics



**Documentation:** Absence of Data Documentation



**Process:** Little Help with Data Understanding



**Engineering:** No Infrastructure to Handle Change in Data



When planning the entire product, it seems important to pay special attention to this collaboration point.



Communication



Documentation



Engineering



Process

## Requirements and Planning

- Product and Model Requirements
- Project Planning



## Training Data

- Negotiating Data Quality and Quantity



## Product-Model Integration

- Responsibility and Cultural Clashes
- Quality Assurance for Model and Product



Communication



Documentation



Engineering



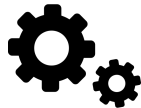
Process



Define processes, responsibilities, and boundaries more carefully



Document APIs at collaboration points between teams



Recruit engineering support for model deployment, monitoring, data validation, etc.



Establish a team culture with mutual understanding and exchange



Communication



Documentation



Engineering



Process



# Summary

## Why do 87% of data science projects never make it into production?

Collaboration Problems

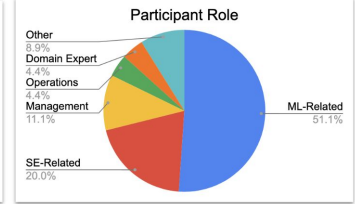
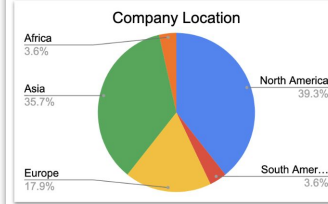
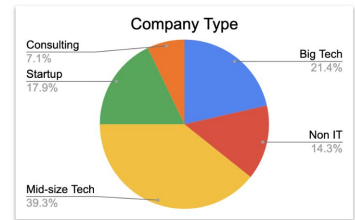
VB Staff

July 19, 2019 4:

And the third issue, intimately connected to those silos, is the lack of collaboration. Data scientists have been around since the 1950s — and they were individuals sitting in a basement working behind a terminal. But now that it's a team sport, and the importance of that work is now being embedded into the fabric of the company, it's essential that every person on the team is able to collaborate with everyone else: the data engineers, the data stewards, people that understand the data science, or analytics, or BI specialists, all the way up to DevOps and engineering.

"This is a big place that holds companies back because they're not used to collaborating in this way," Leff says. "Because when they take those insights, and they flip them over the wall, now you're asking an engineer to rewrite a data science model created by a data scientist, how's that work out, usually?"

Conducted  
45 interviews in  
28 organizations



### Collaboration Points

### Themes

#### Requirements and Planning

- Product and Model Requirements
- Project Planning



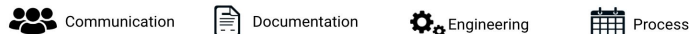
#### Training Data

- Negotiating Data Quality and Quantity



#### Product-Model Integration

- Responsibility and Cultural Clashes
- Quality Assurance for Model and Product



**Communication:** Lack of ML literacy leads to unrealistic requirements



Involving data scientists early when soliciting product requirements



**Documentation:** Product requirements are often not translated into clear model requirements



Adopt more formal requirements documentation for product and model