

Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

Lesson Title: Introduction

Presenter: Max Lotstein

Format: Talking Head

Video Name: T-PSML-0_3_l1_introduction

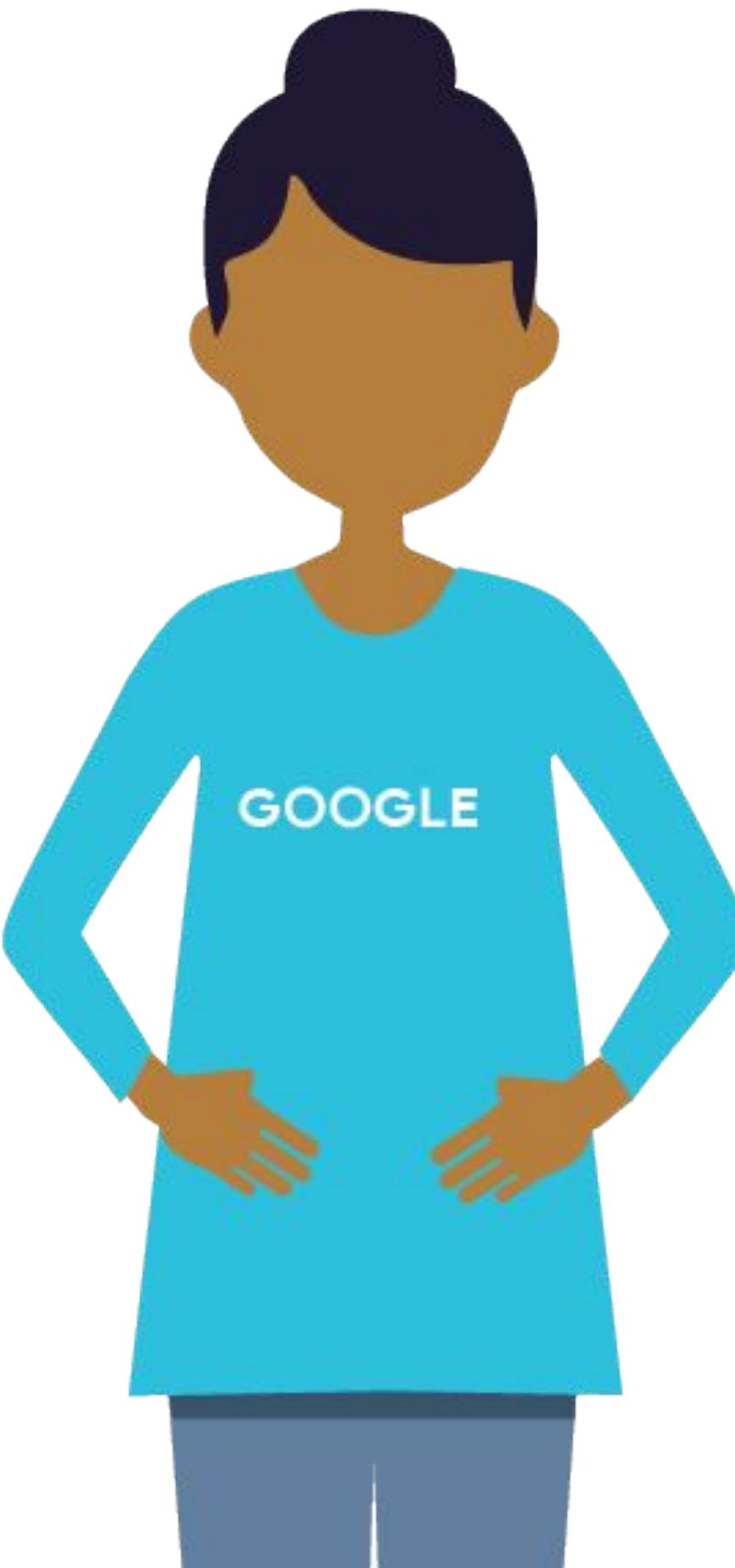


Google Cloud

Designing Adaptable ML Systems

Advanced Machine Learning
on GCP

Max Lotstein



Learn how to...

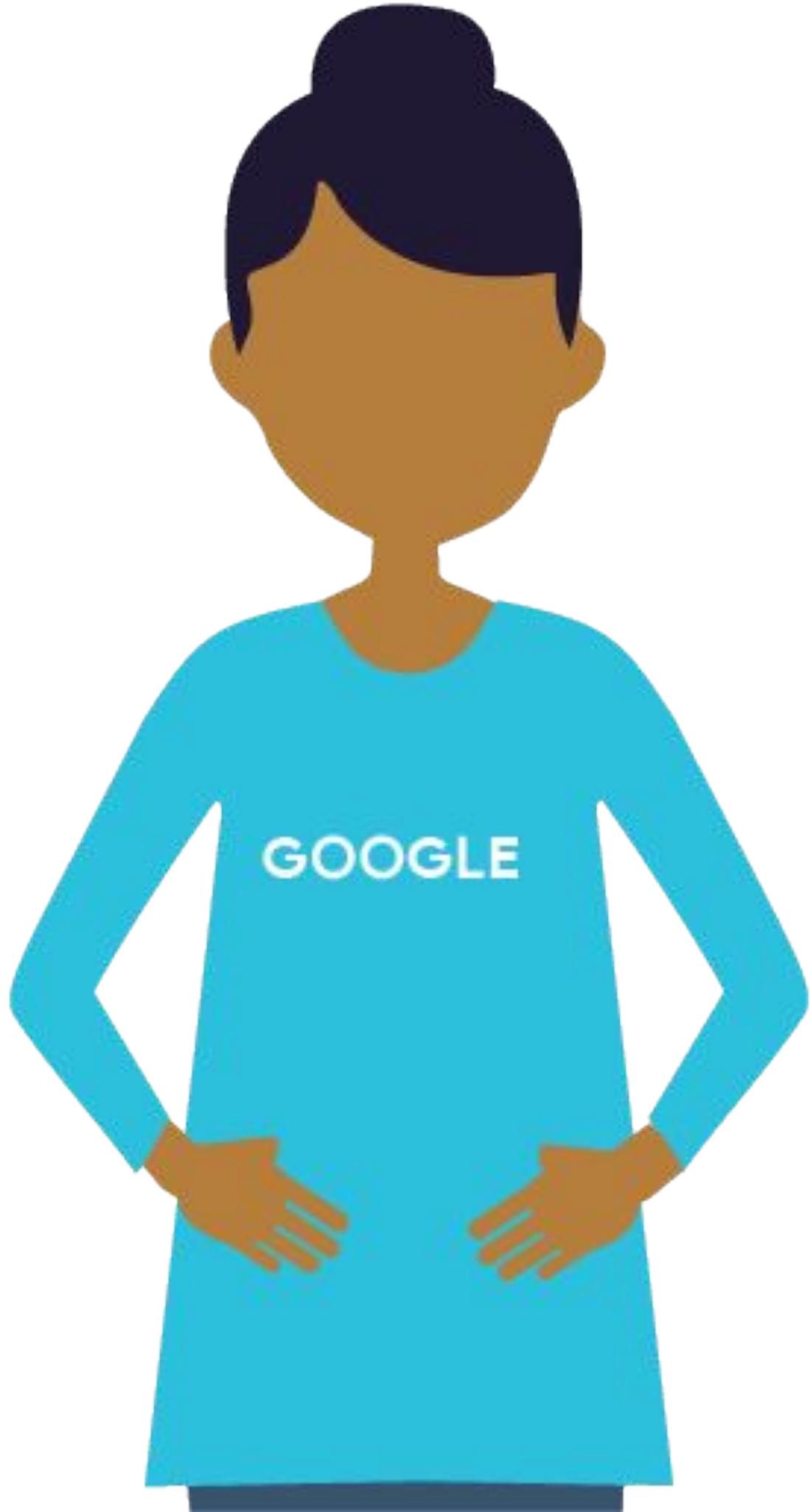
Recognize various data dependencies

Make cost-conscious engineering decisions

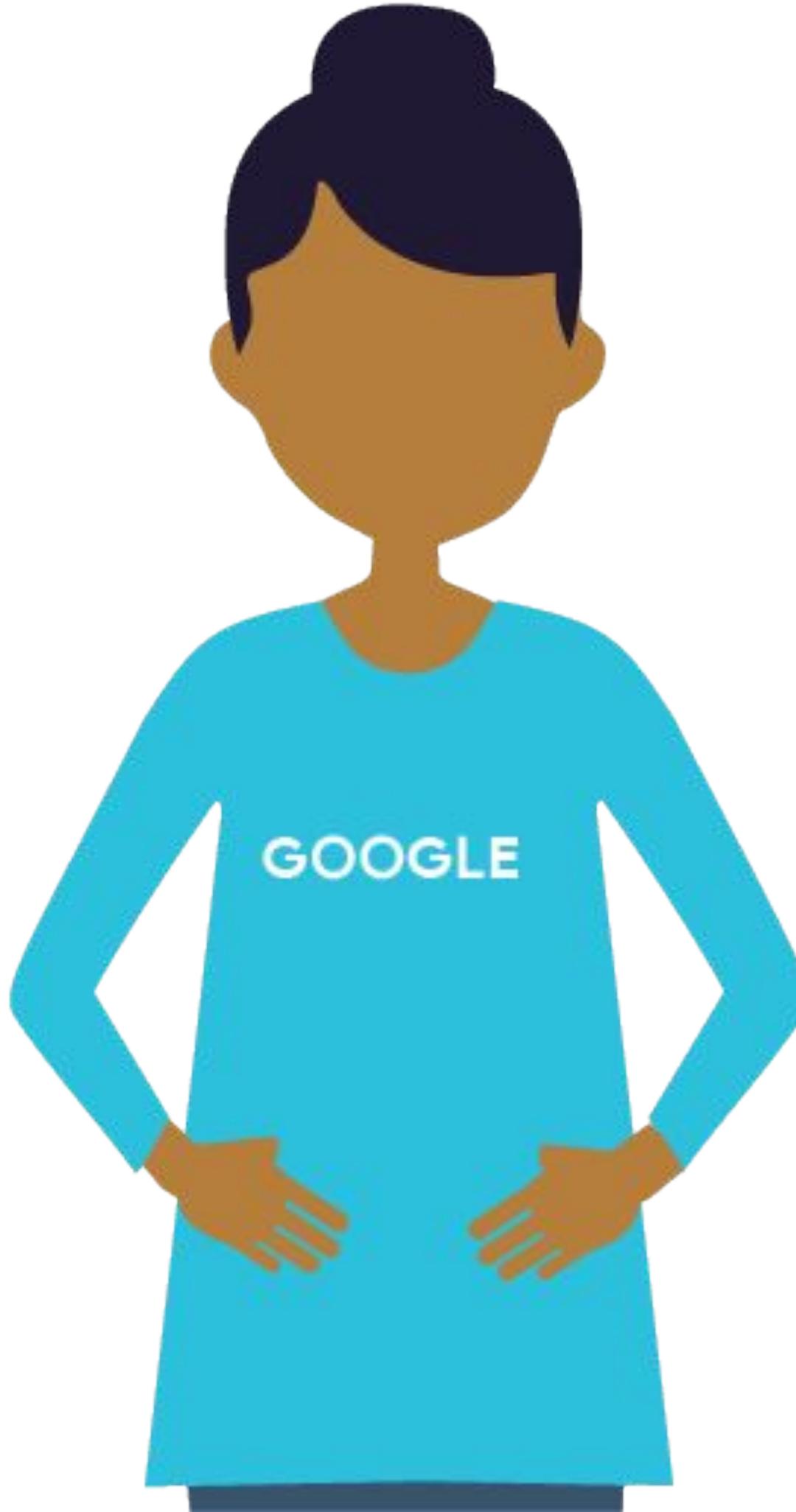
Mitigate model pollution

Implement a pipeline that is immune to one type of dependency

Debug the causes of observed model behavior



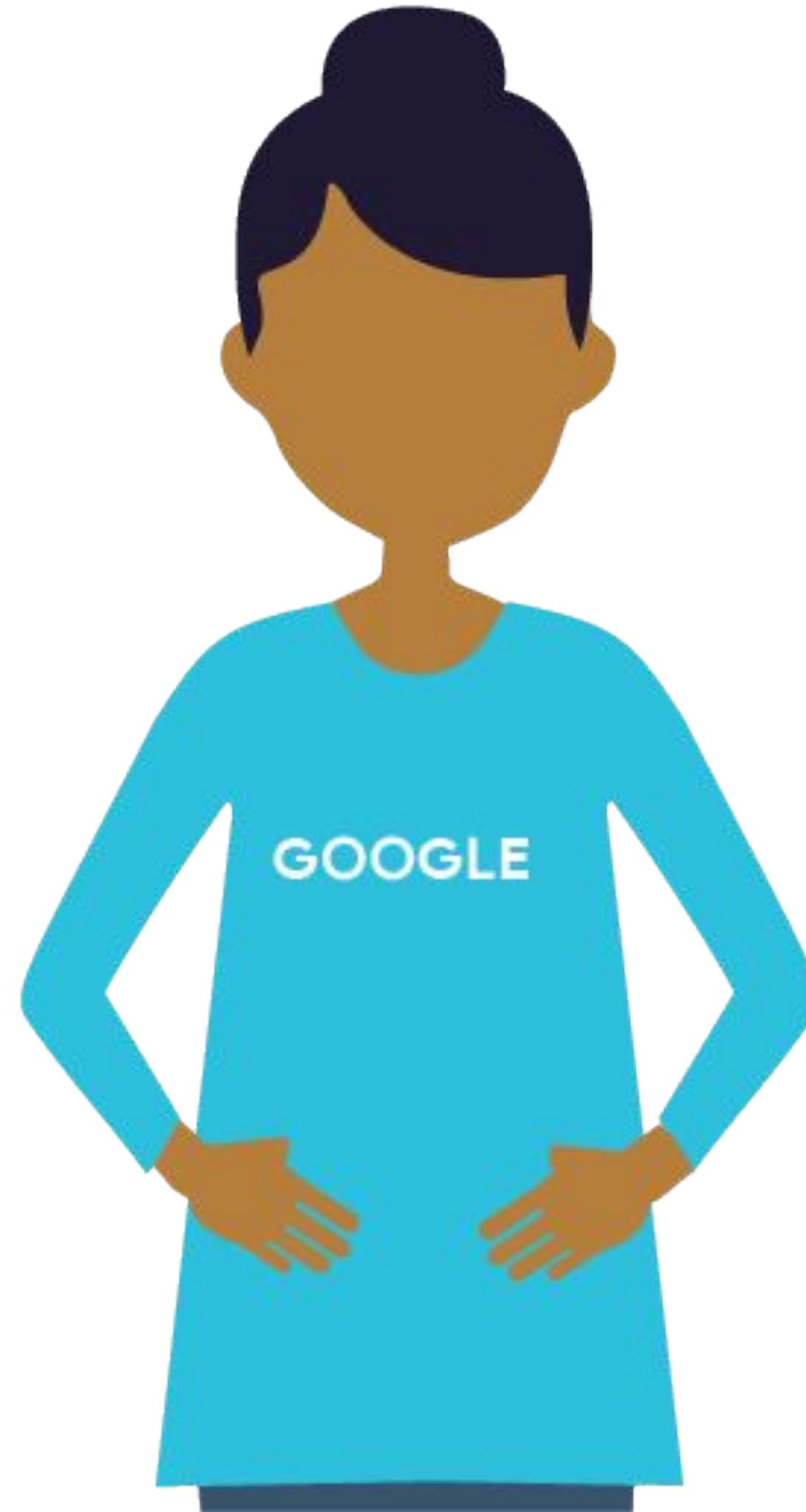
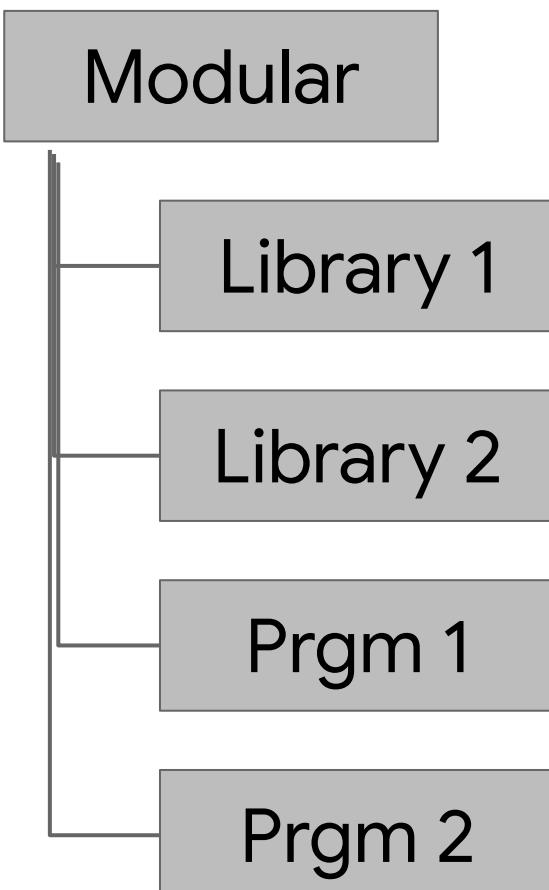
Few Programs are Islands



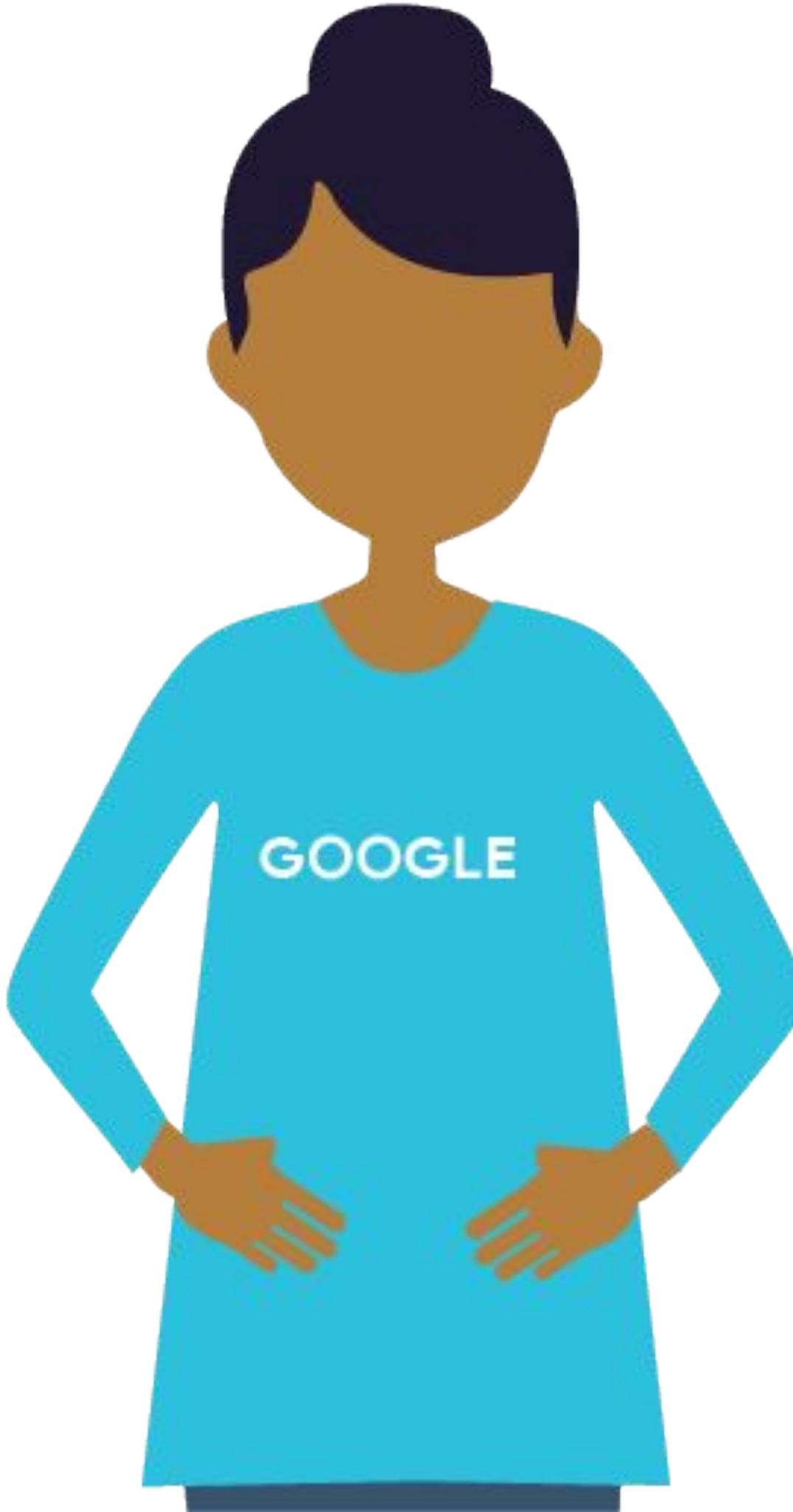
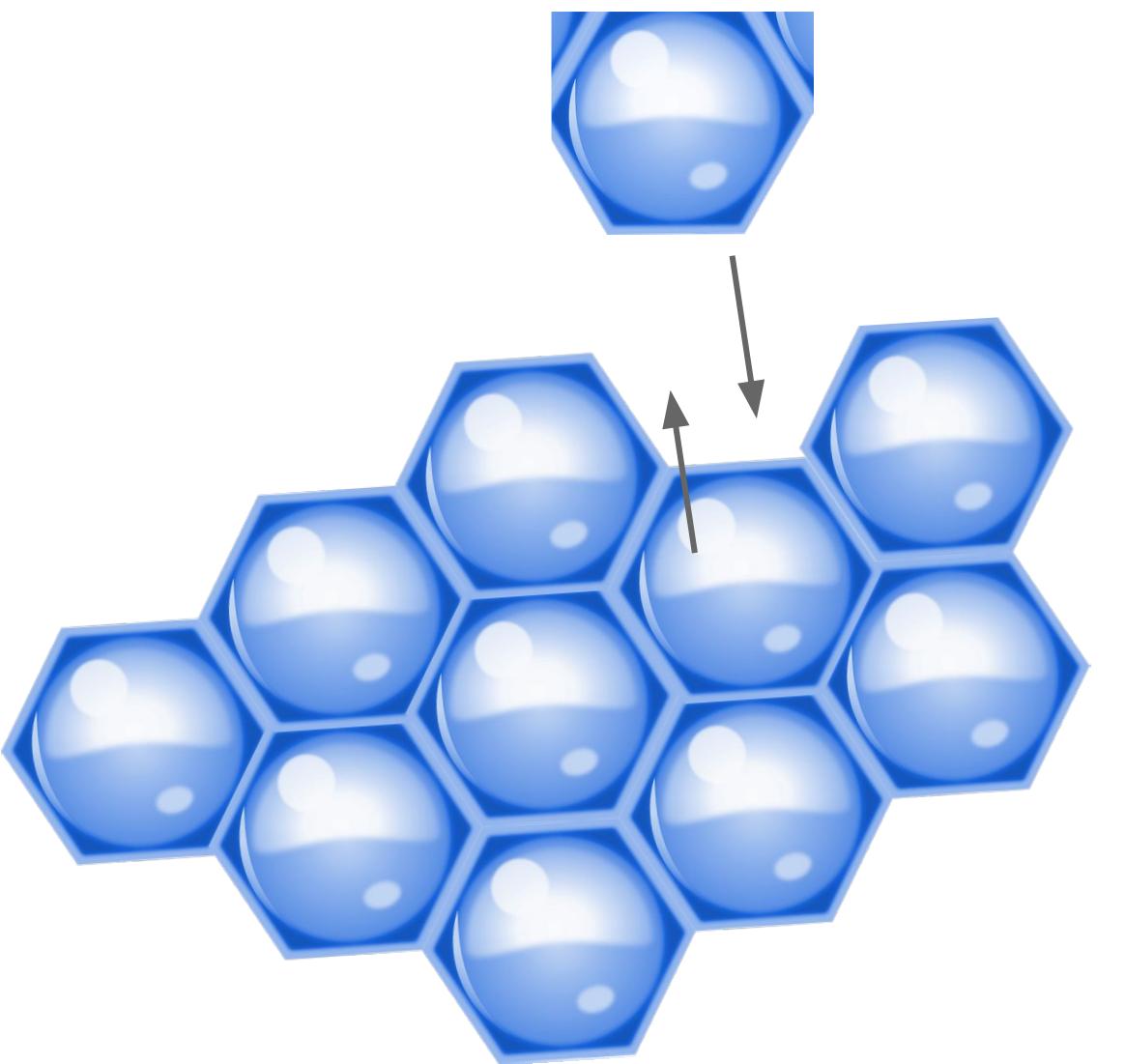
Few Programs are Islands

Monolithic Program

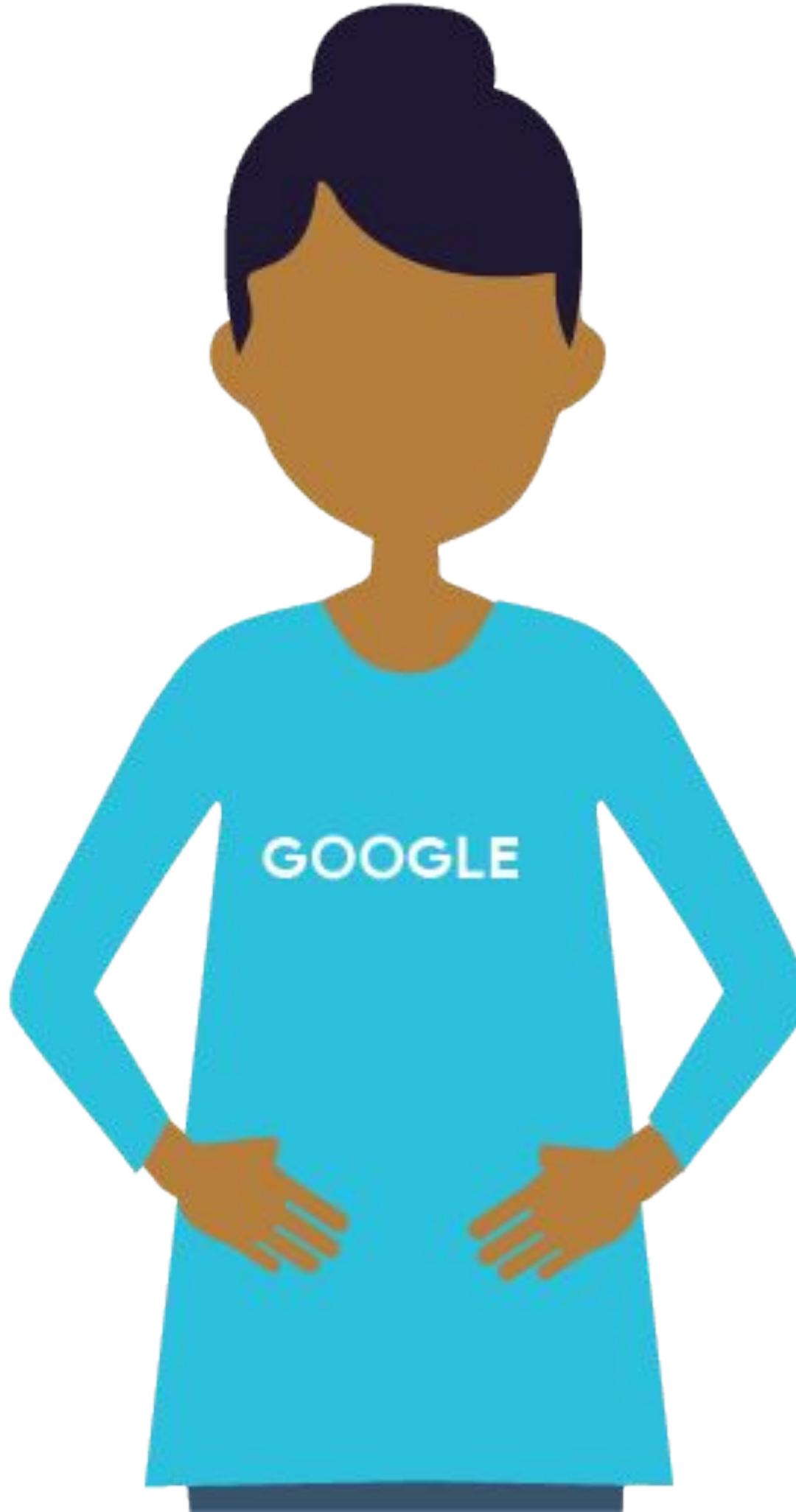
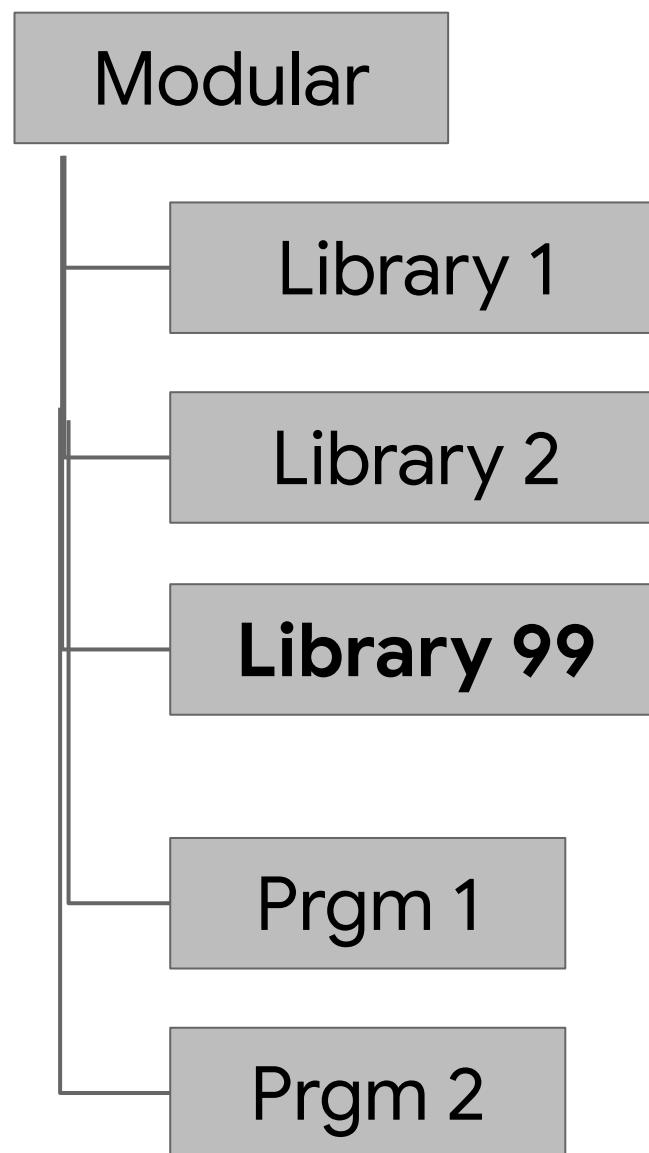
VS



Modular Is More Maintainable



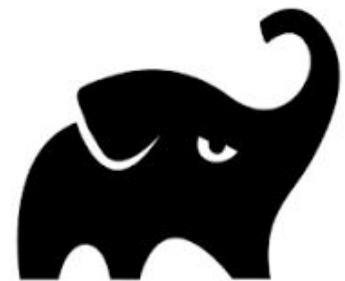
Dependency Management Is Manageable



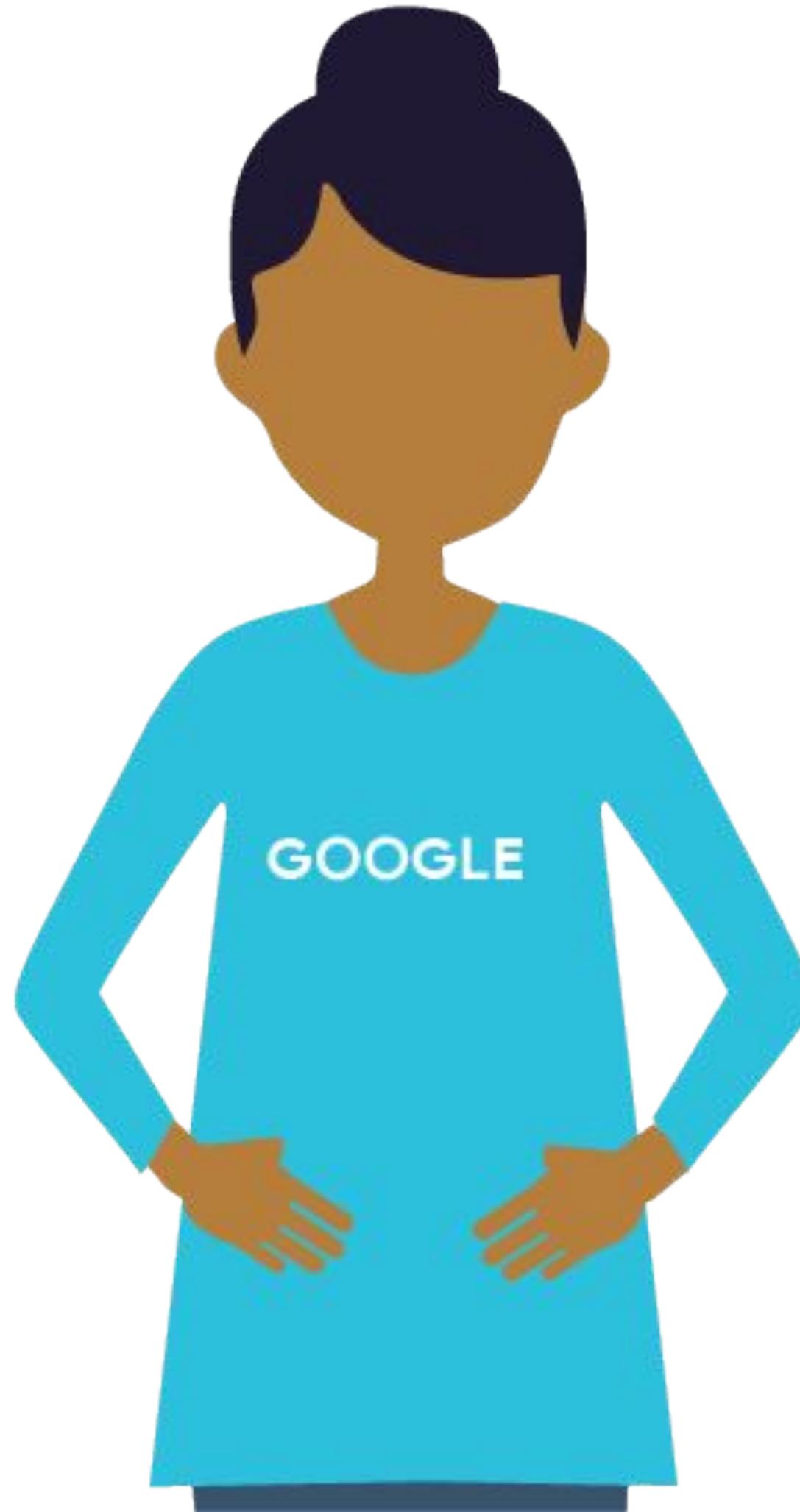
Dependency Management Is Manageable



Maven™

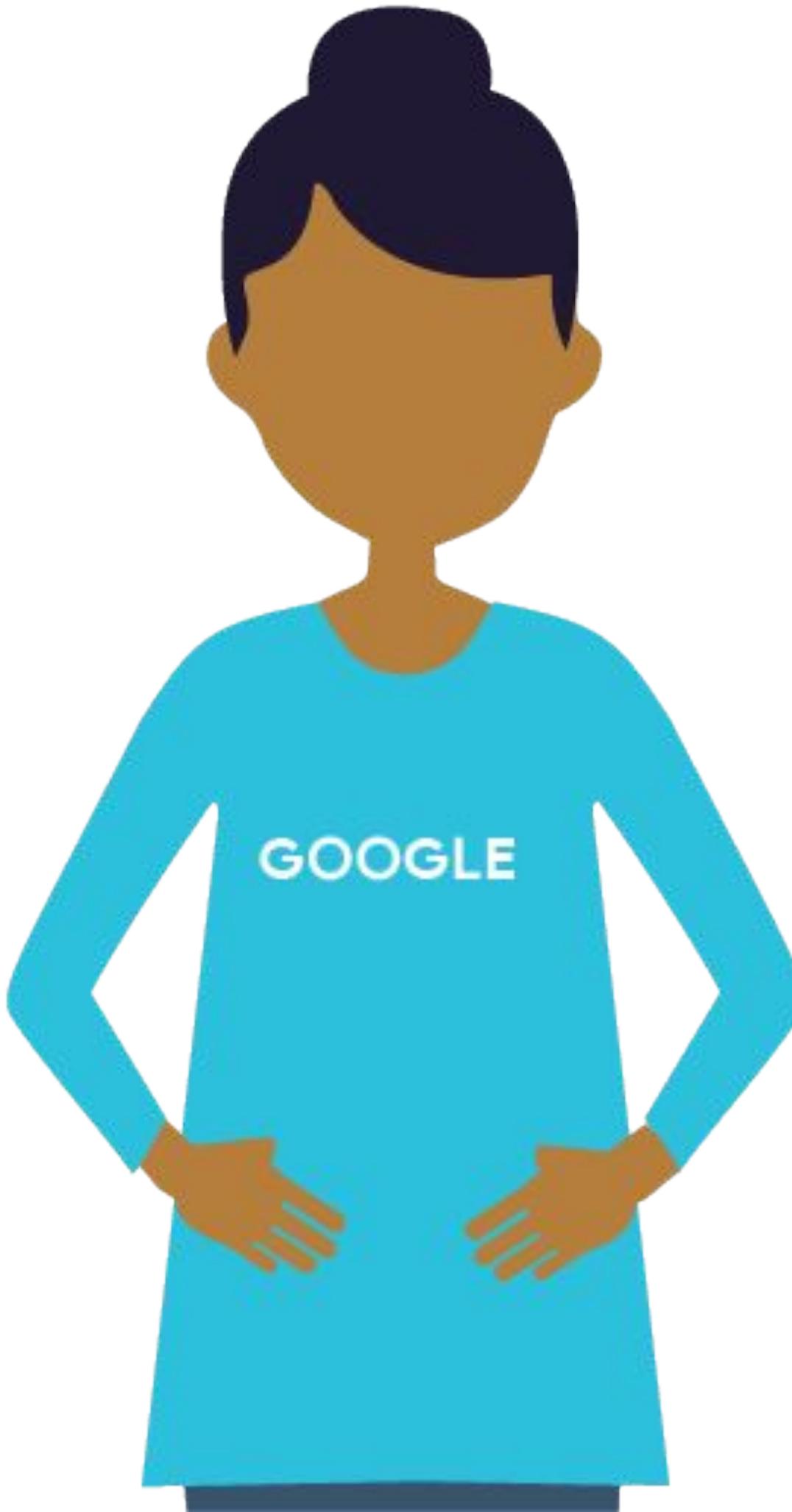


PIP



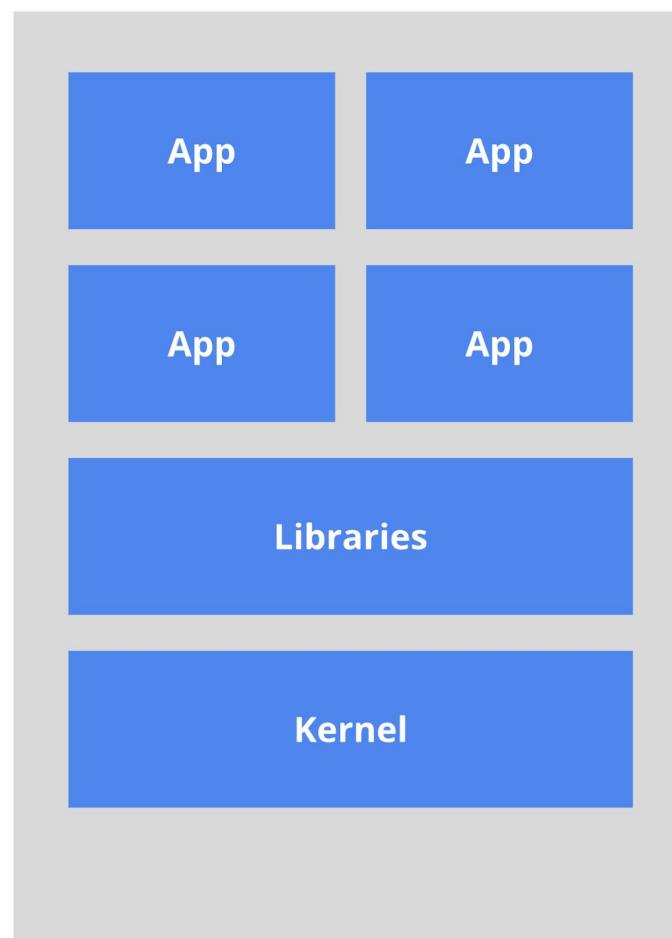
Explicit Dependencies Make Life Easier

```
1. <project xmlns="http://maven.apache.org/POM/4.0.0"
2.   xsi:schemaLocation="http://maven.apache.org/
3.   <modelVersion>4.0.0</modelVersion>
4.
5.   <groupId>com.mycompany.app</groupId>
6.   <artifactId>my-app</artifactId>
7.   <version>1.0-SNAPSHOT</version>
8.   <packaging>jar</packaging>
9.
10.  <name>Maven Quick Start Archetype</name>
11.  <url>http://maven.apache.org</url>
12.
13.  <dependencies>
14.    <dependency>
15.      <groupId>junit</groupId>
16.      <artifactId>junit</artifactId>
17.      <version>4.8.2</version>
18.      <scope>test</scope>
19.    </dependency>
20.  </dependencies>
21. </project>
```



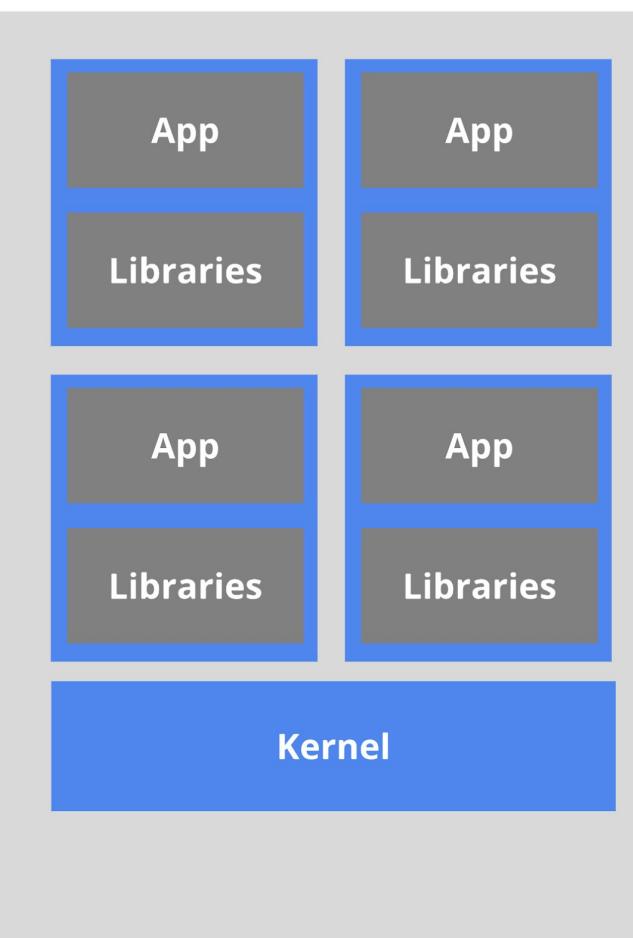
Containers eliminate infrastructure dependencies

The old way: Applications on host

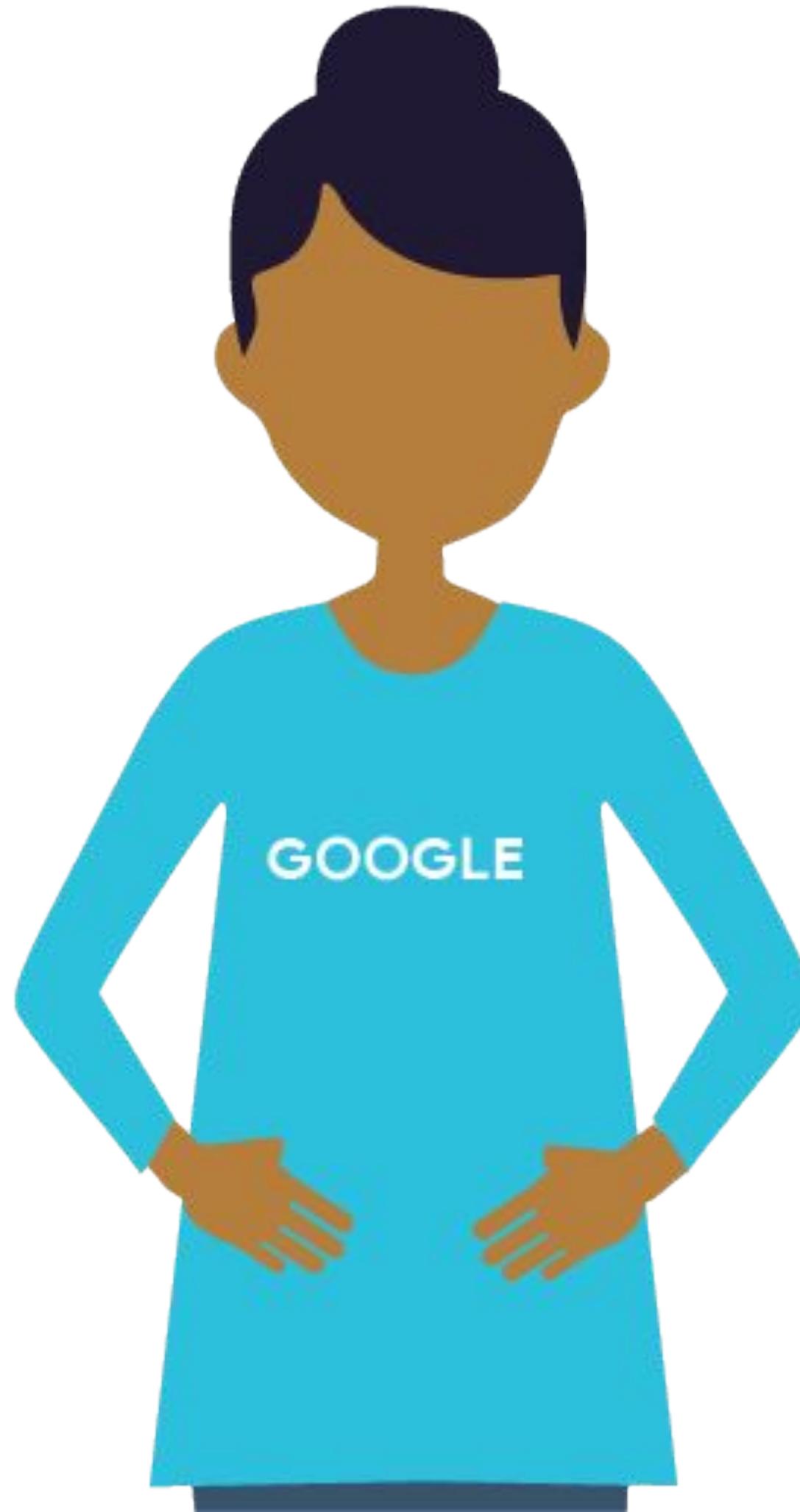


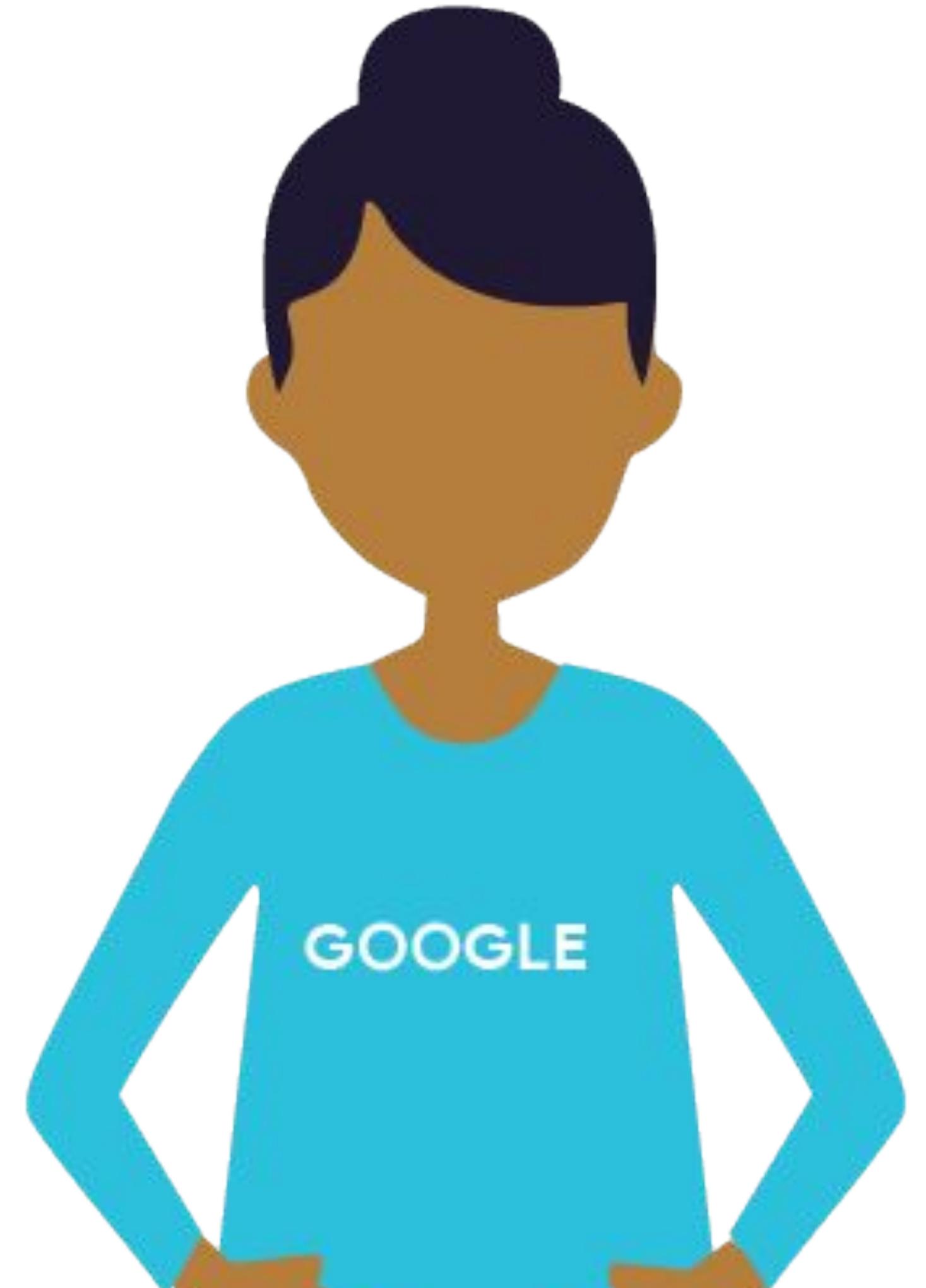
*Heavyweight, non-portable
Relies on OS package manager*

The new way: Deploy containers

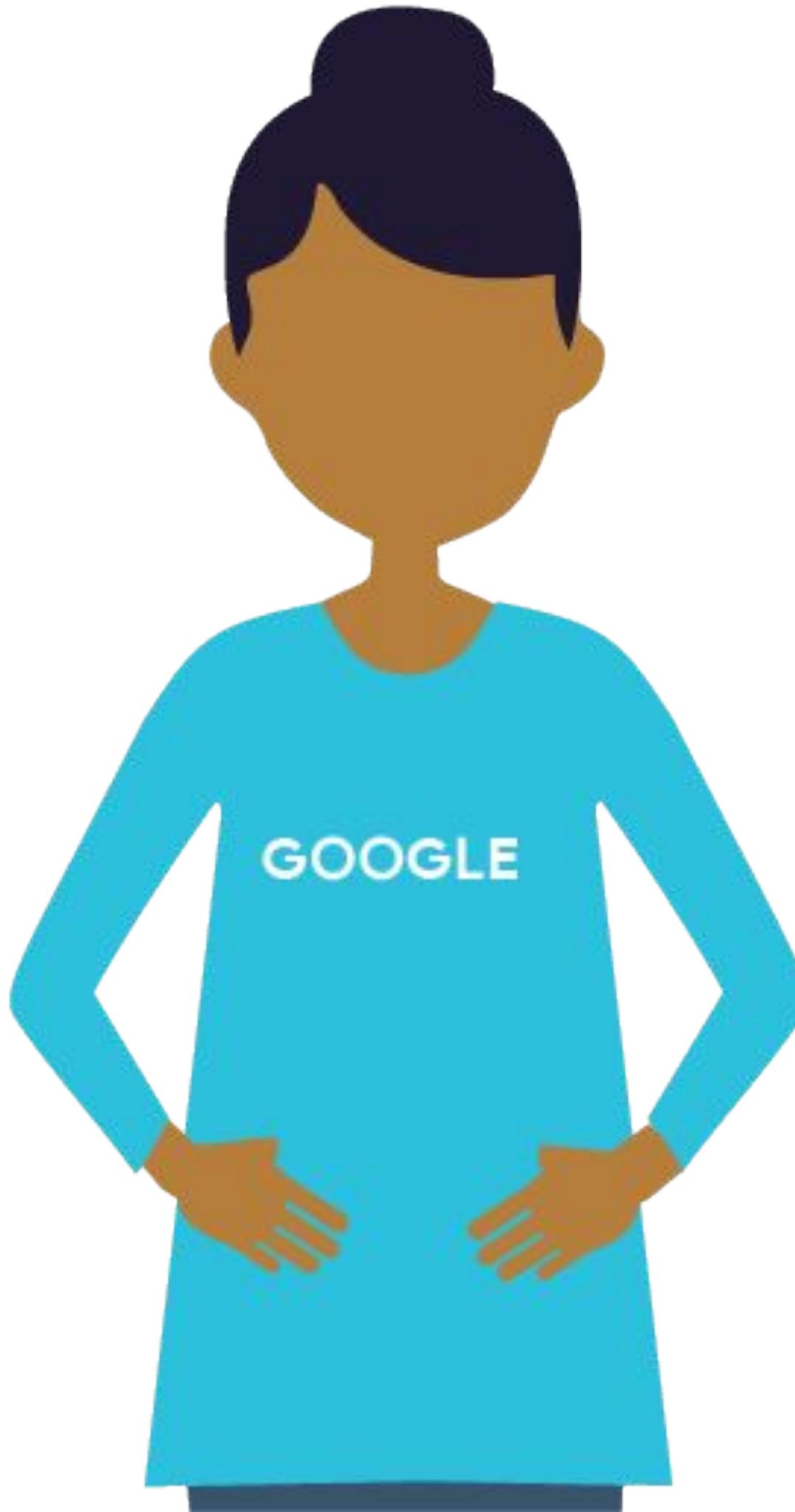
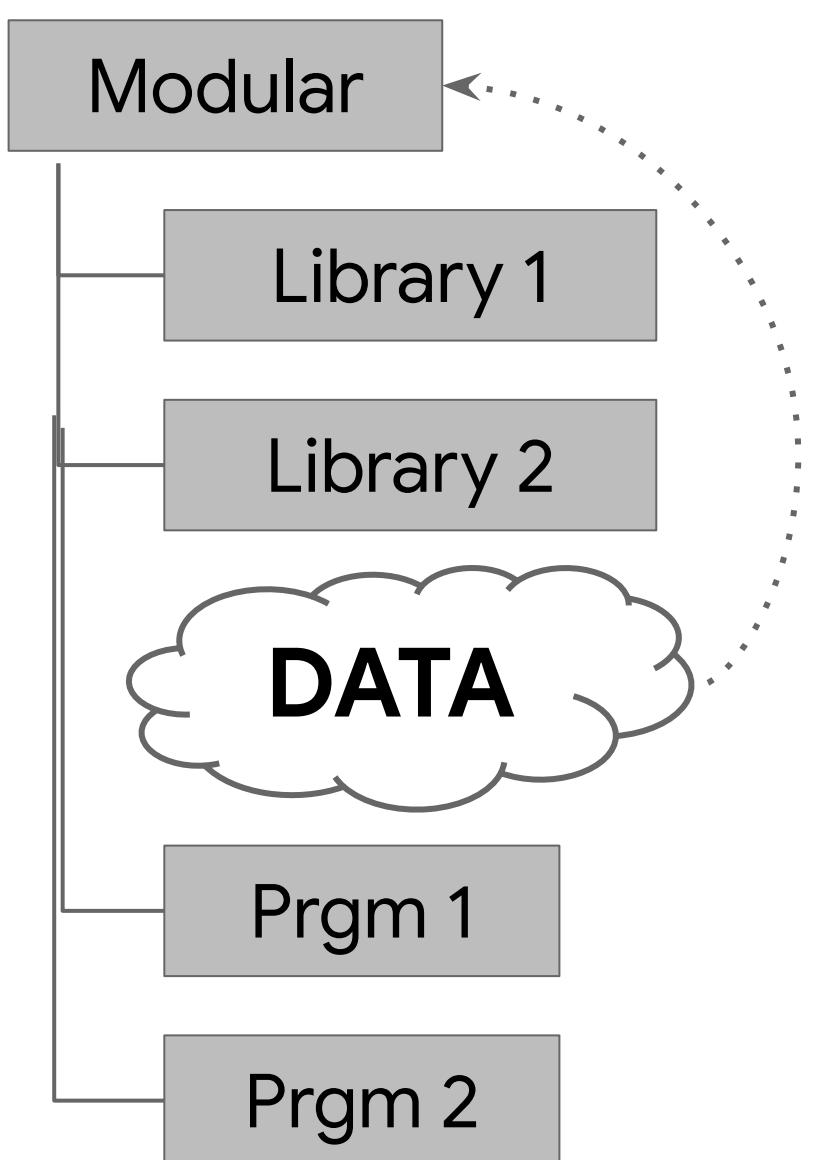


*Small and fast, portable
Uses OS-level virtualization*



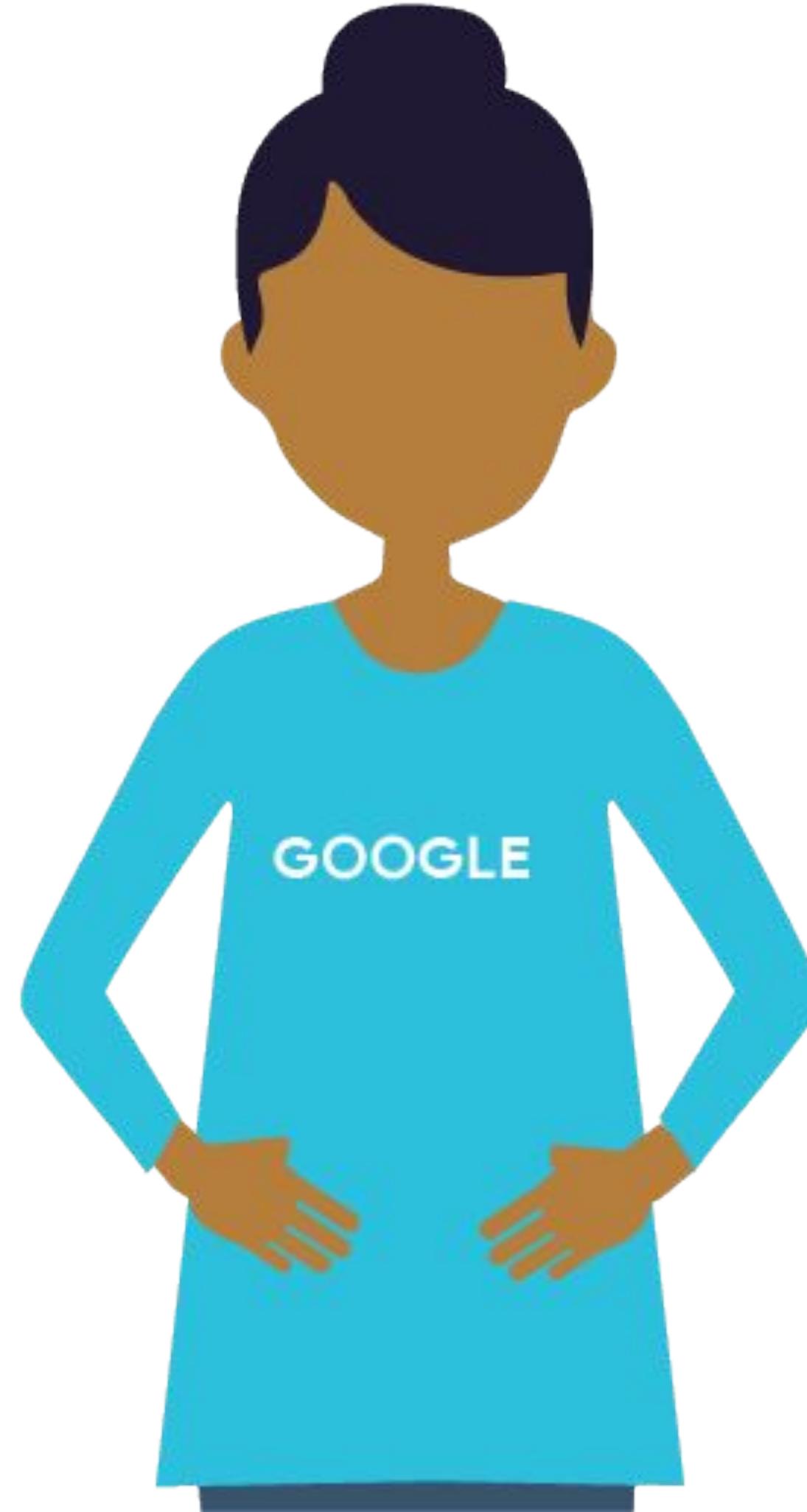


Data: the Dependency Outside the Codebase



Mismanaged Dependencies are Costly

- ✗ Losses in prediction quality
- ✗ Decreases to system stability
- ✗ Decreases in team productivity



Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

Lesson Title: Adapting to Data

Presenter: Max Lotstein

Format: Talking Head

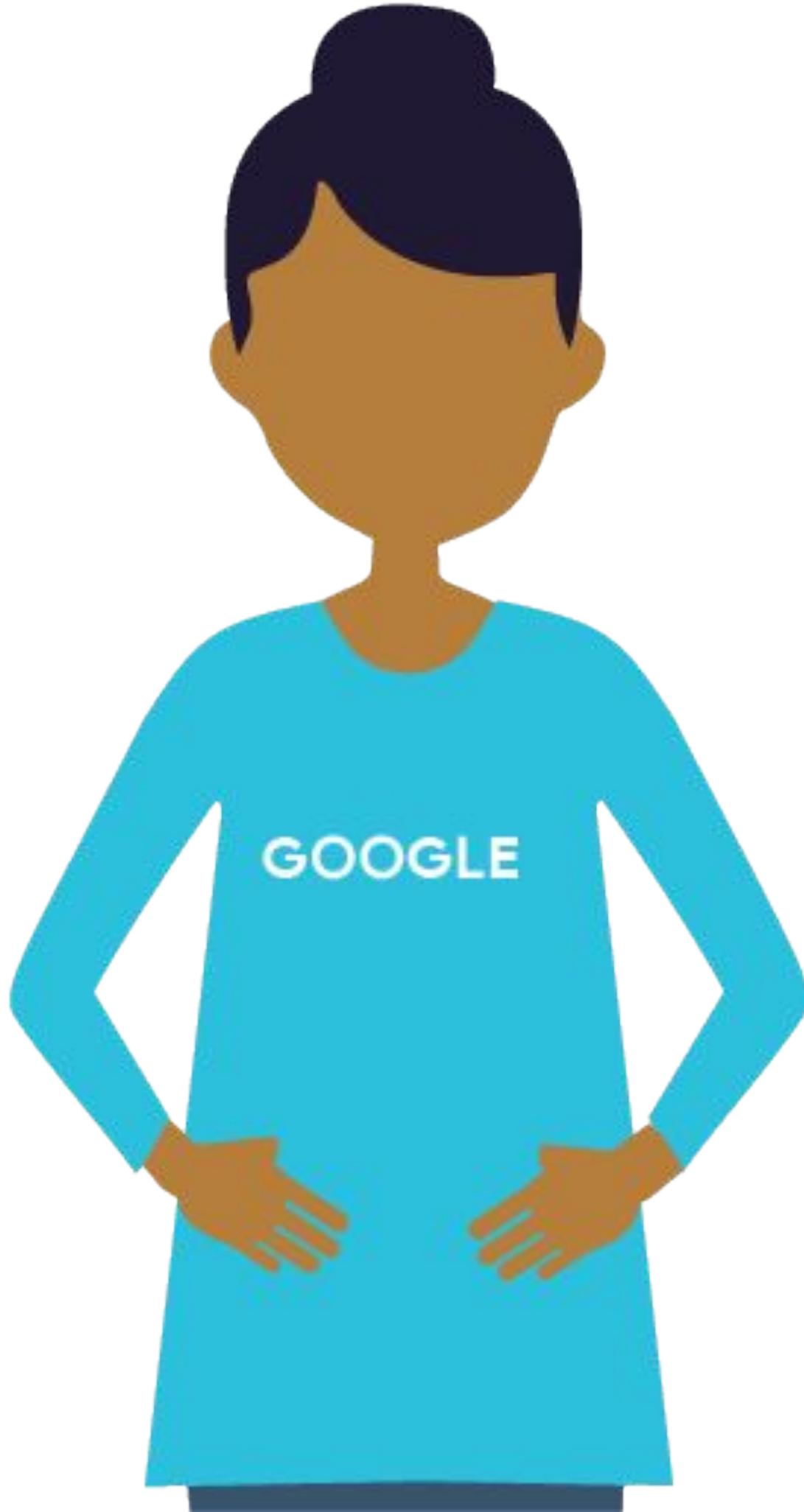
Video Name: T-PSML-0_3_l3_adapting_to_data

Agenda

Adapting to Data

Mitigating Training-Serving
Skew Through Design

Debugging a Production Model

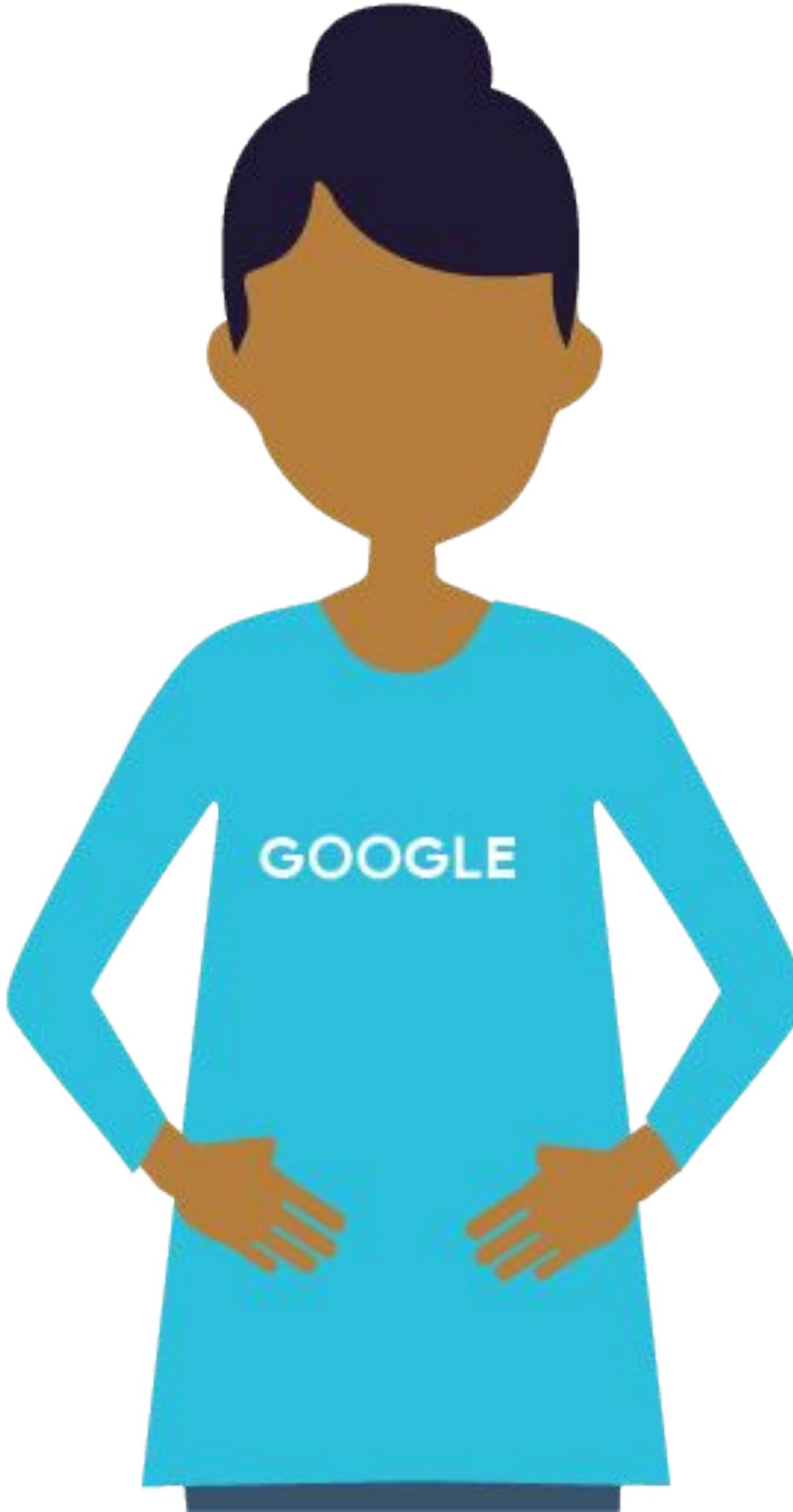


Agenda

Adapting to Data

**Mitigating Training-Serving
Skew Through Design**

Debugging a Production Model

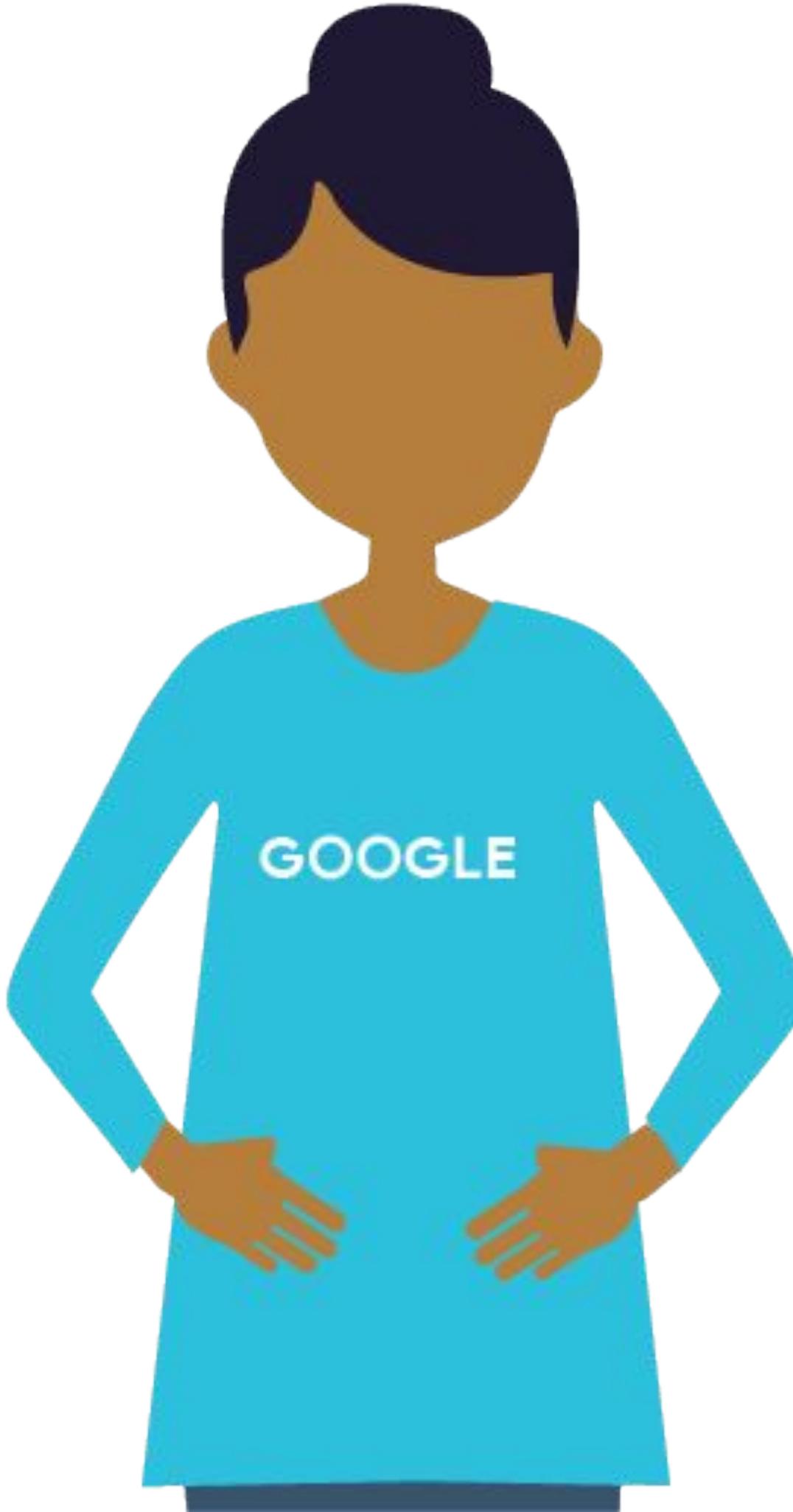


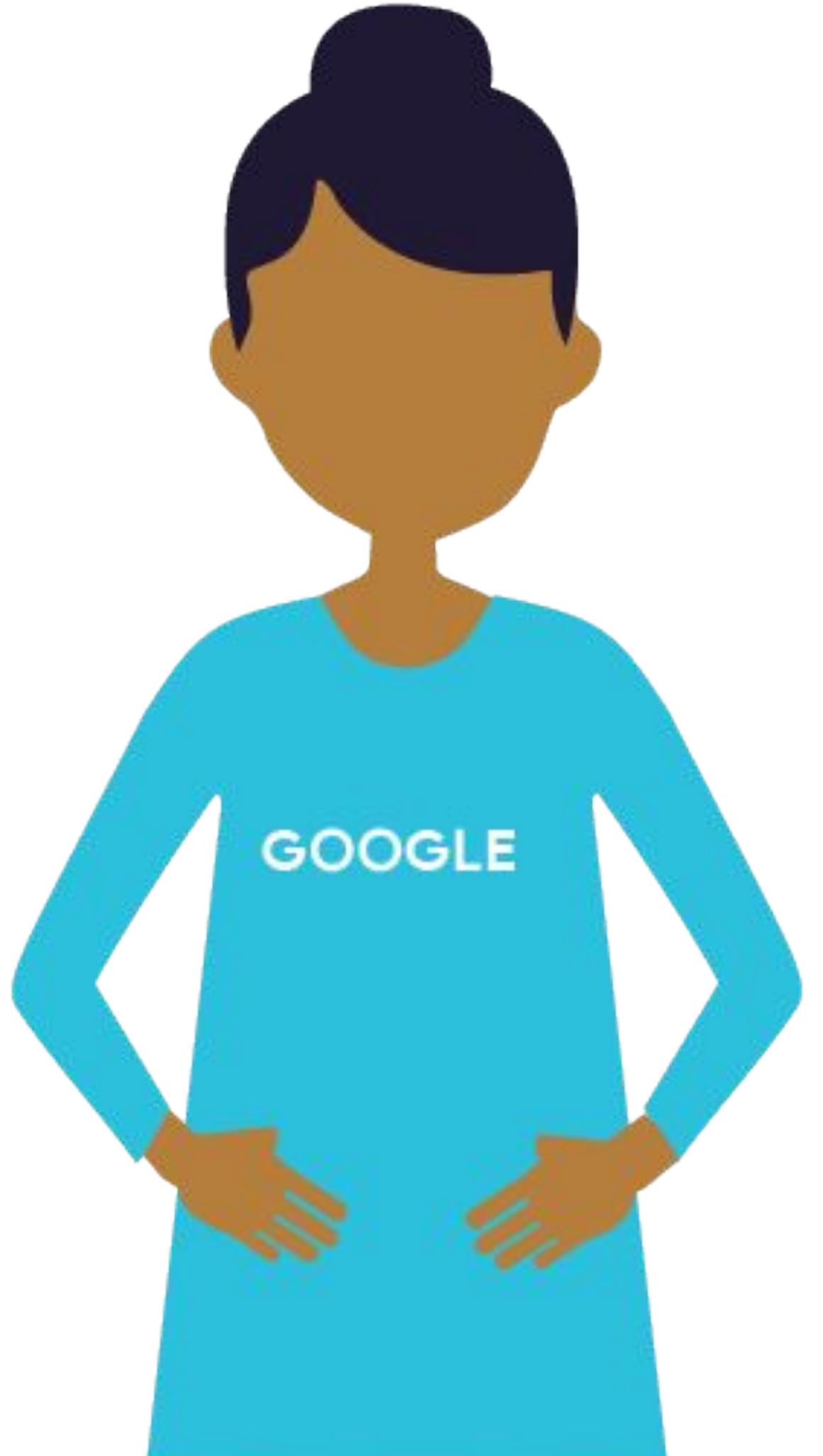
Agenda

Adapting to Data

Mitigating Training-Serving
Skew Through Design

**Debugging a Production
Model**



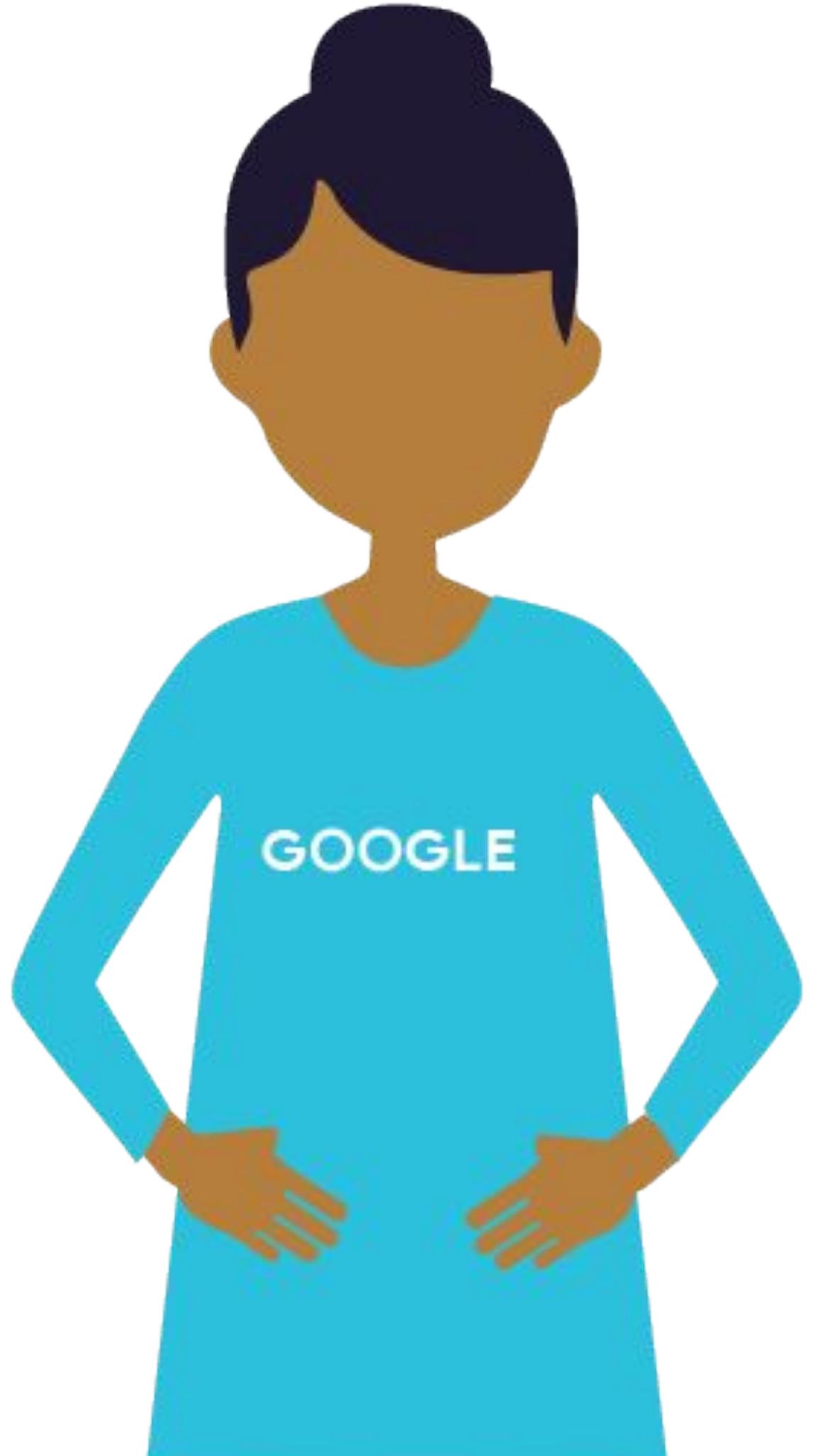


Agenda

Adapting to Data

Mitigating Training-Serving Skew
Through Design

Debugging a Production Model

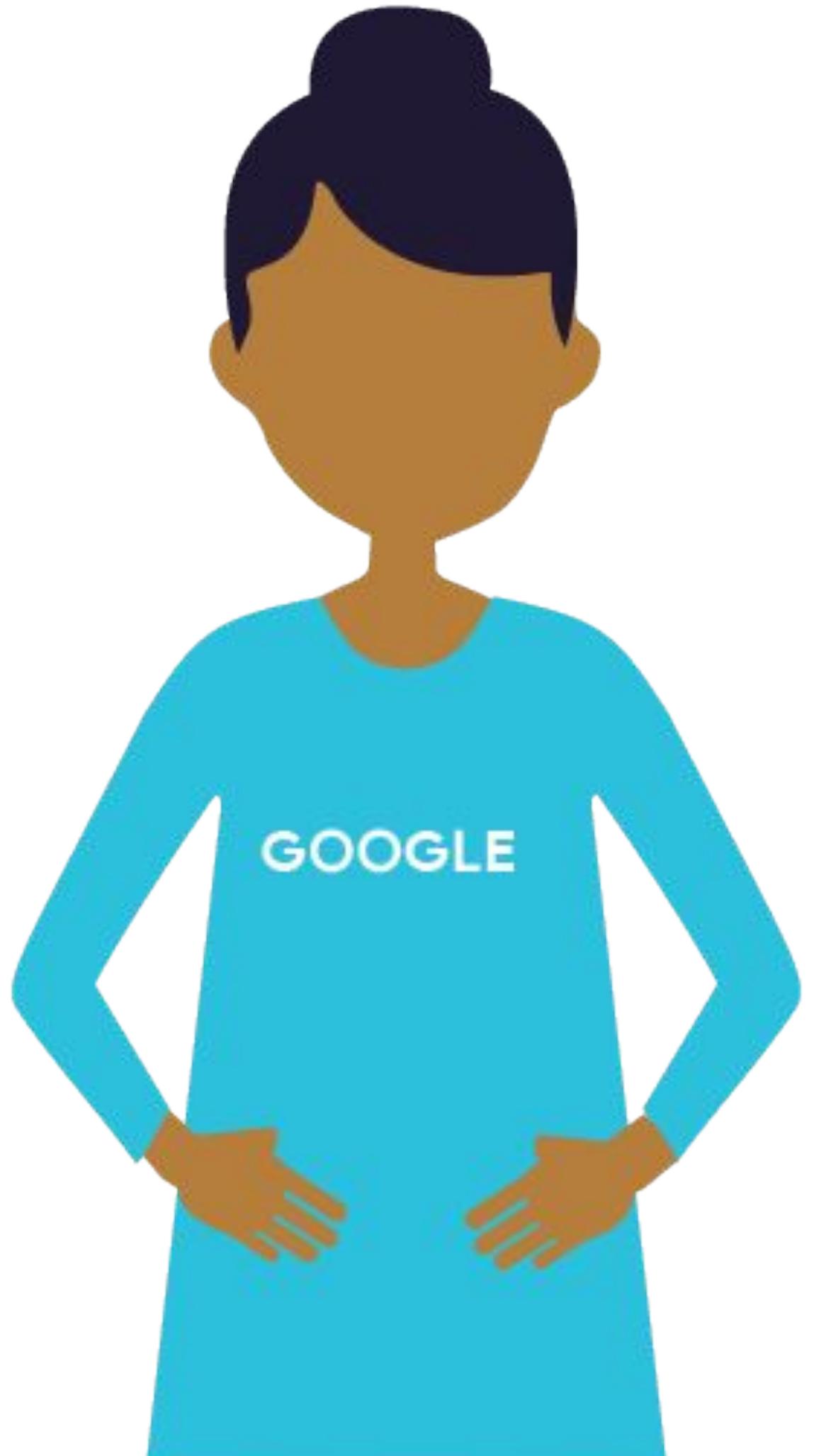


Agenda

Adapting to Data

**Mitigating Training-Serving
Skew Through Design**

Debugging a Production Model

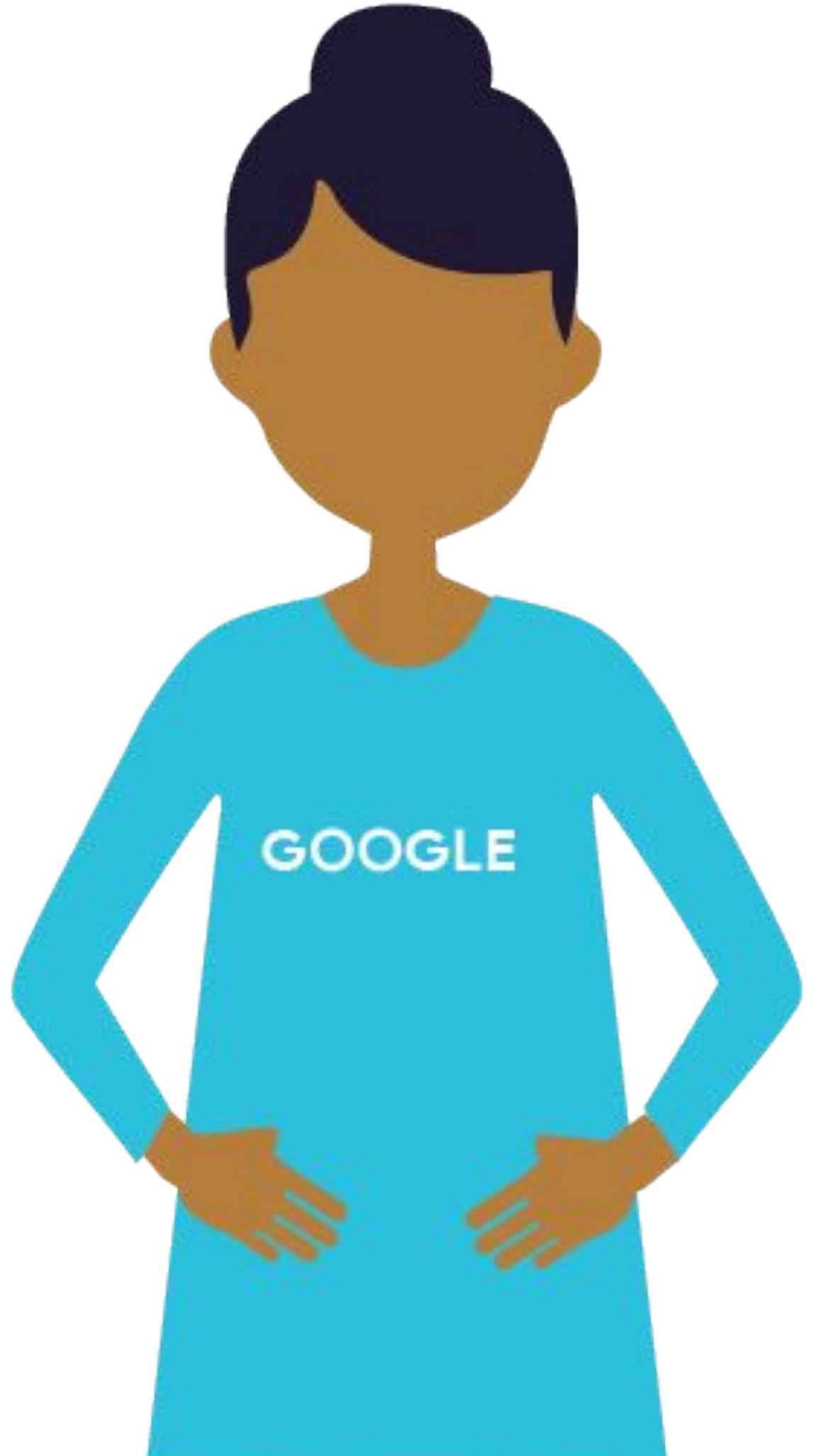


Agenda

Adapting to Data

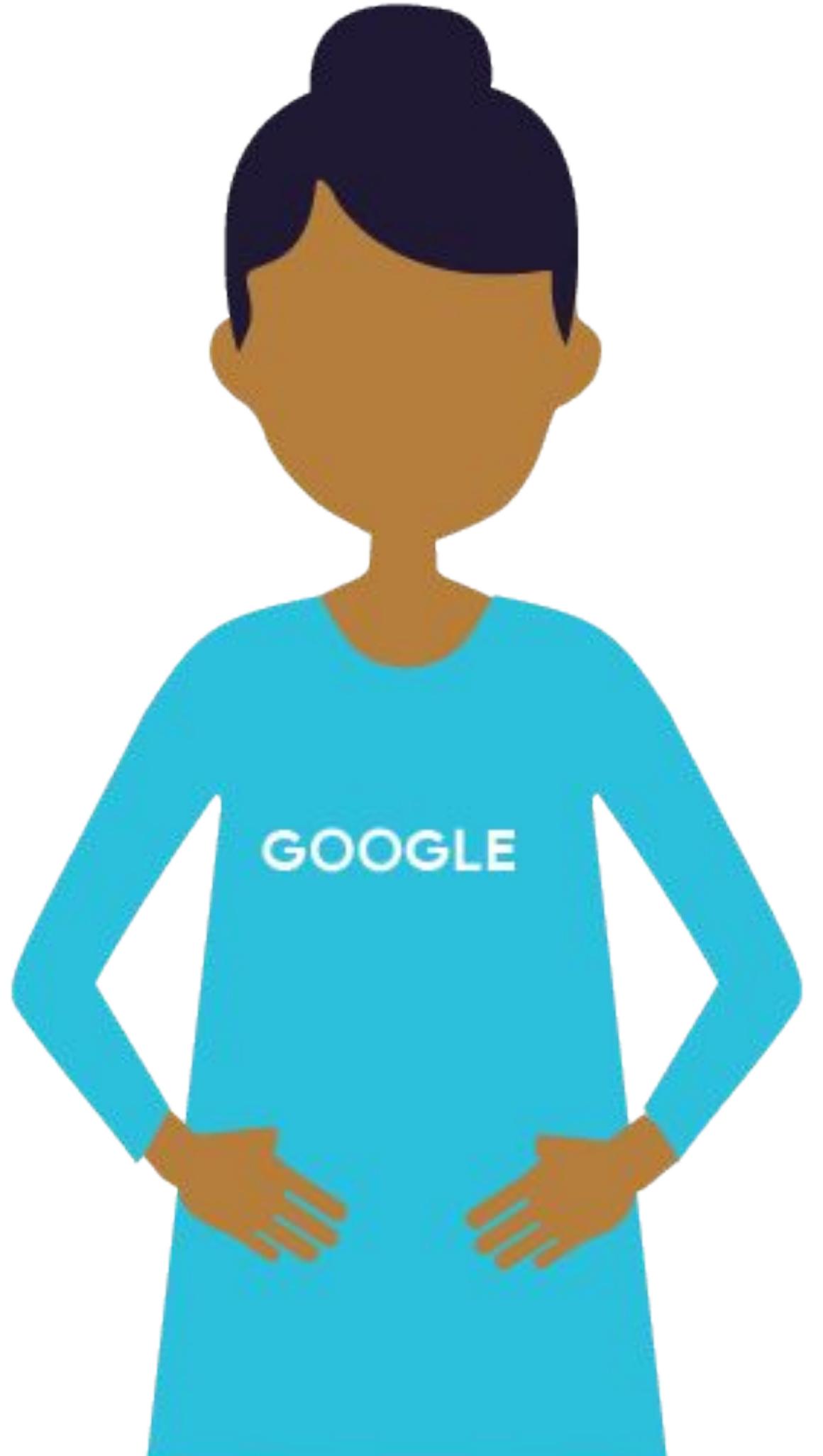
Mitigating Training-Serving Skew
Through Design

Debugging a Production Model



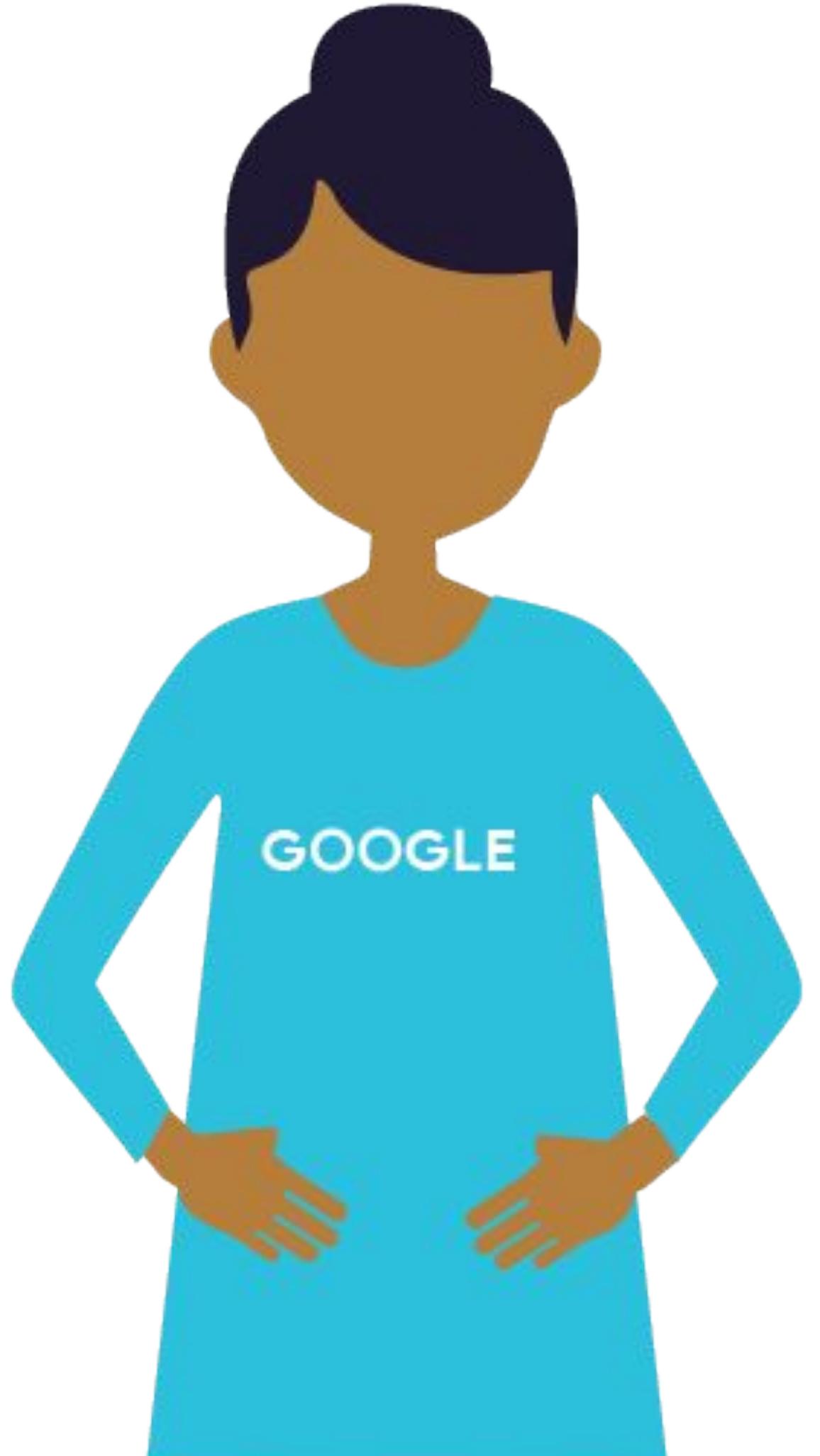
Which of these is least likely to change?

1. An upstream model
2. A data source maintained by another team
3. The relationship between features and labels
4. The distribution of inputs

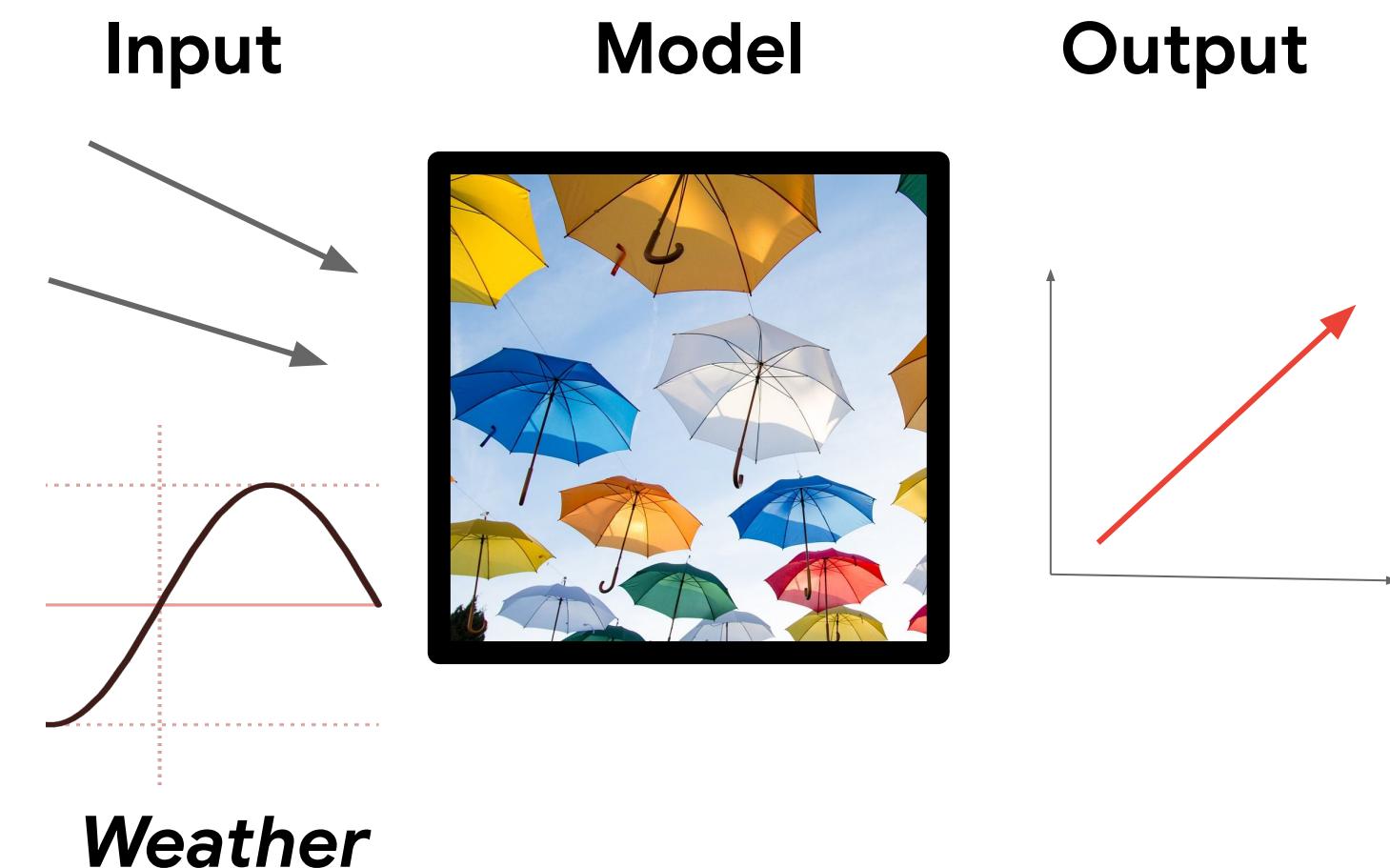


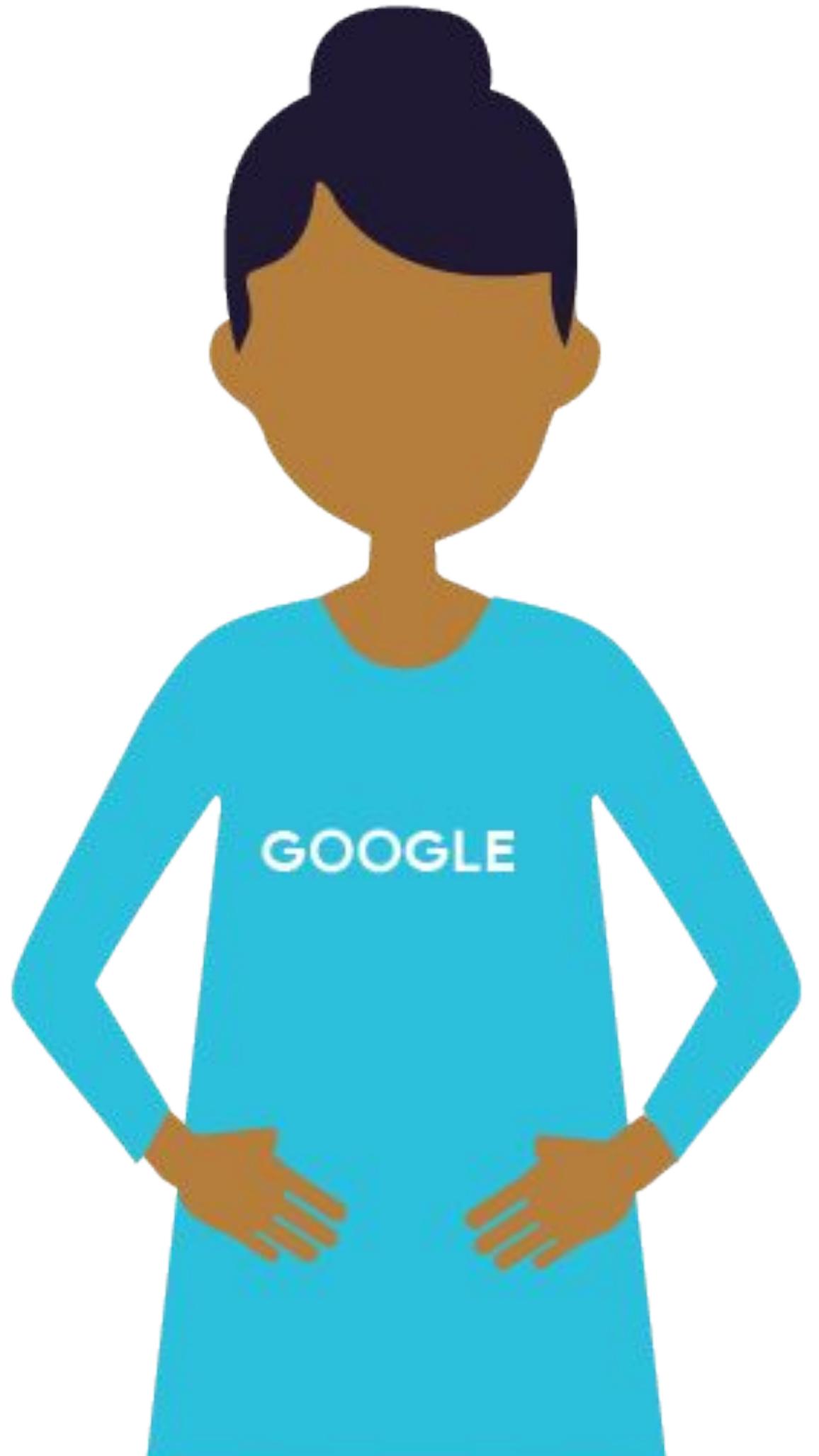
Which of these is least likely to change?

- 1. An upstream model**
- 2. A data source maintained by another team**
- 3. The relationship between features and labels**
- 4. The distribution of inputs**

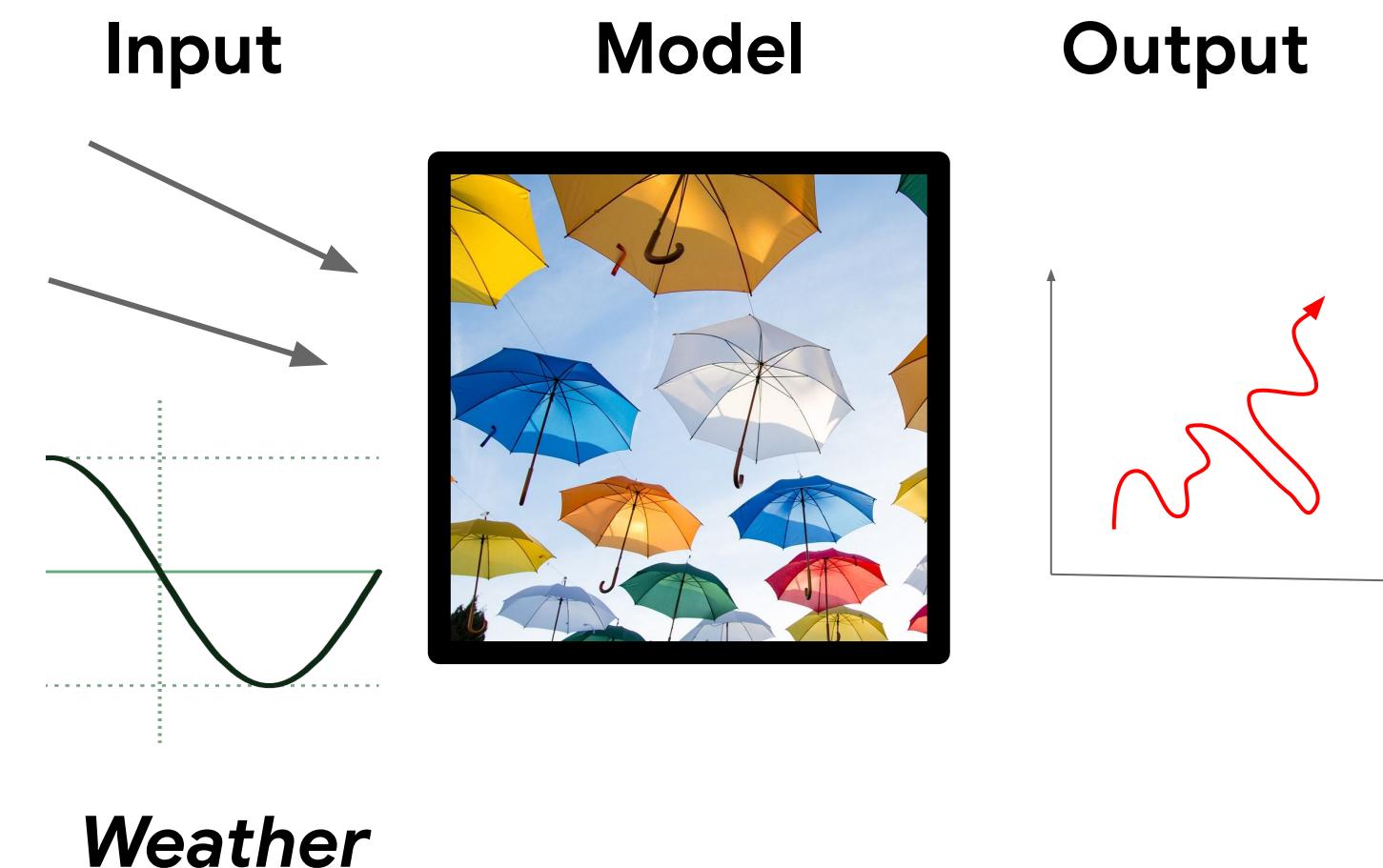


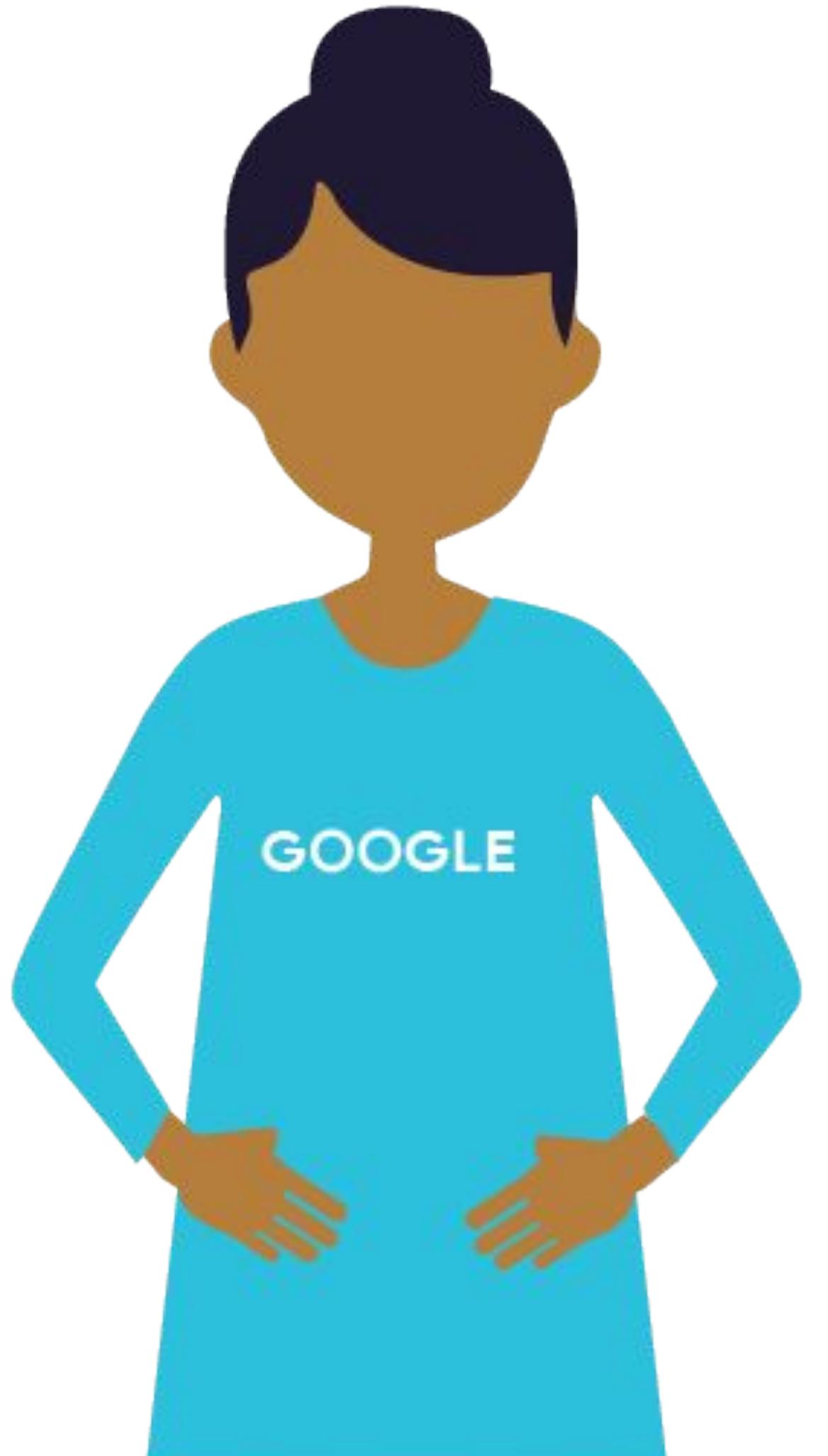
Decoupled upstream data producers



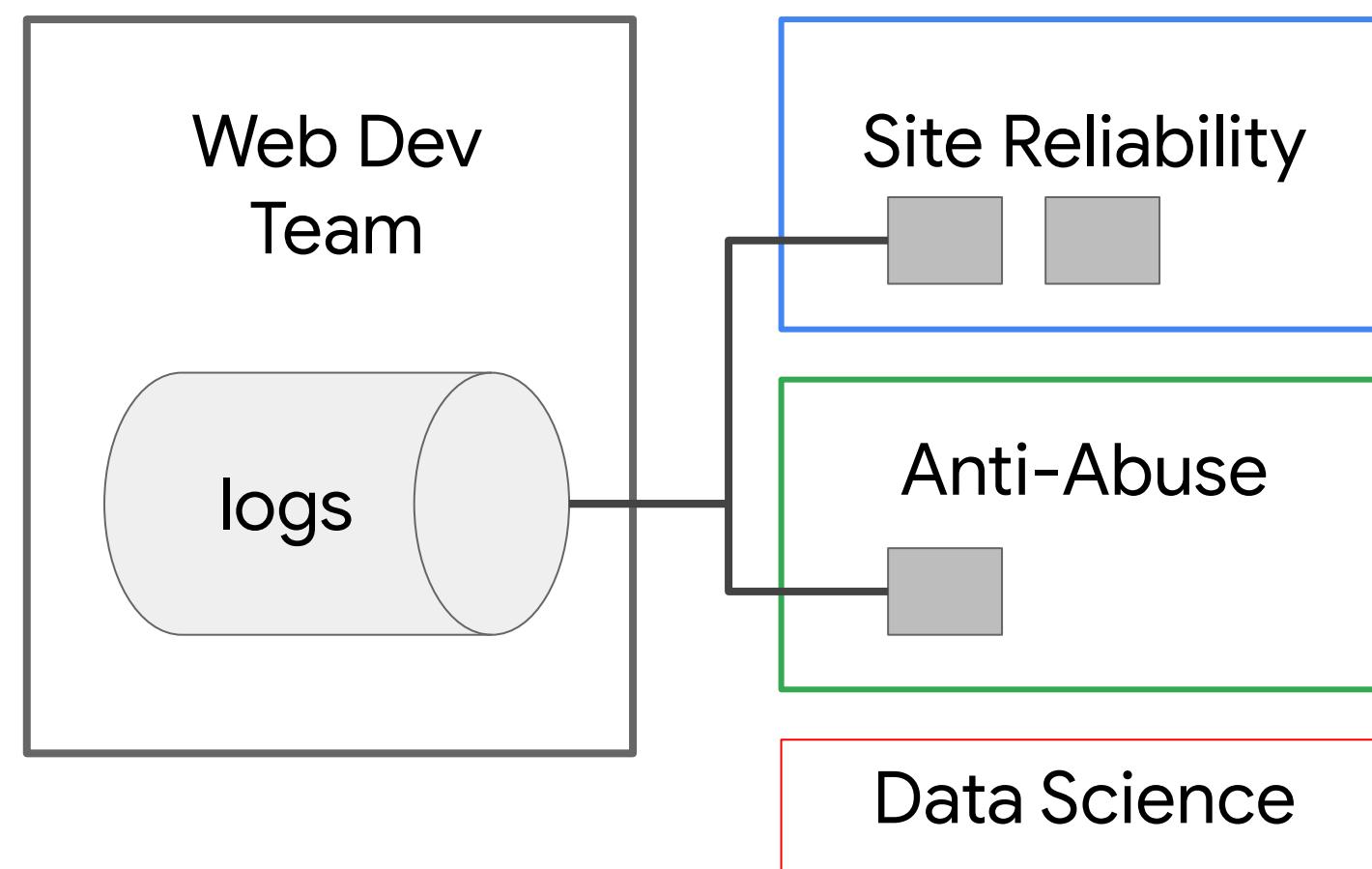


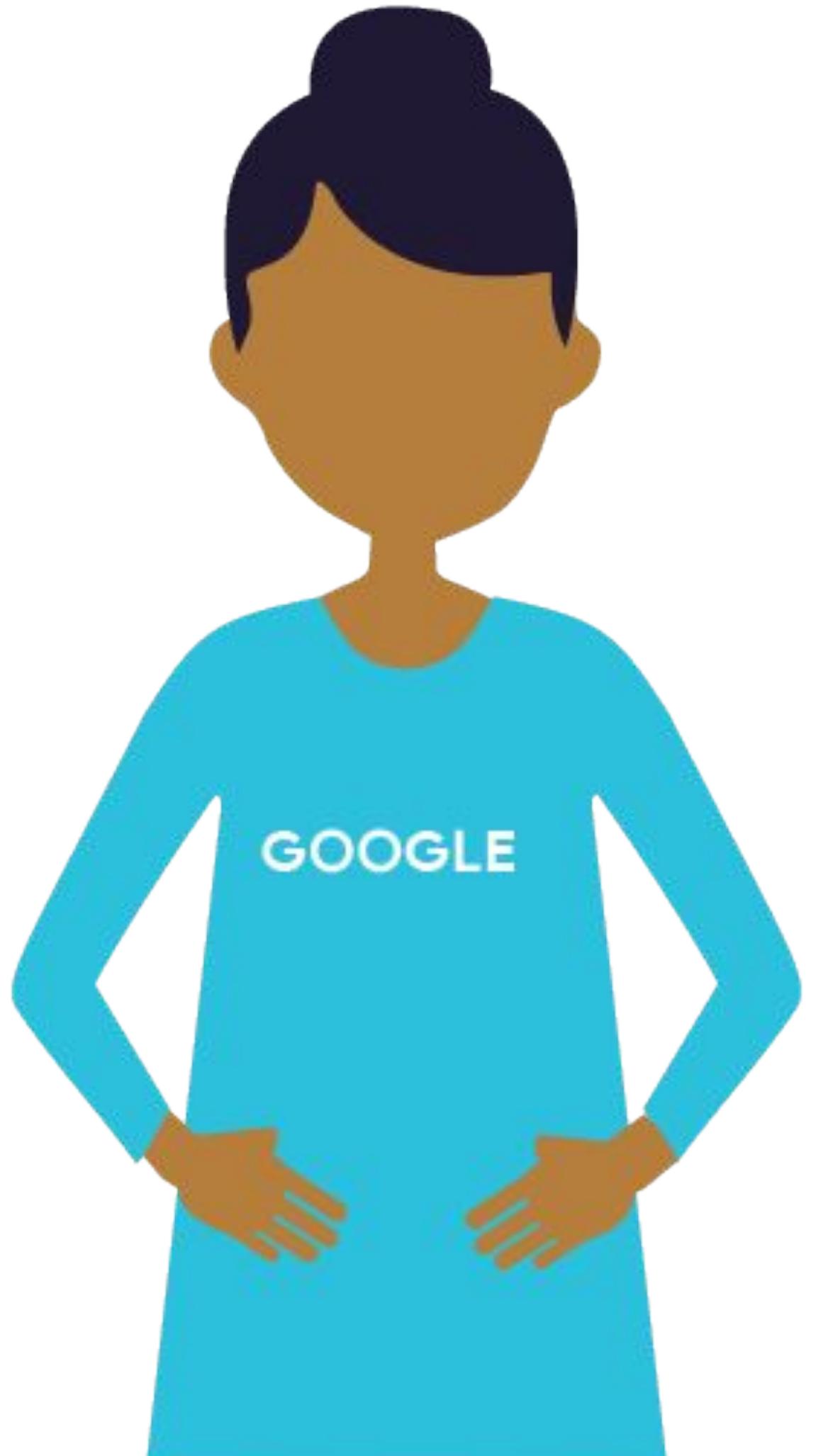
Decoupled upstream data producers



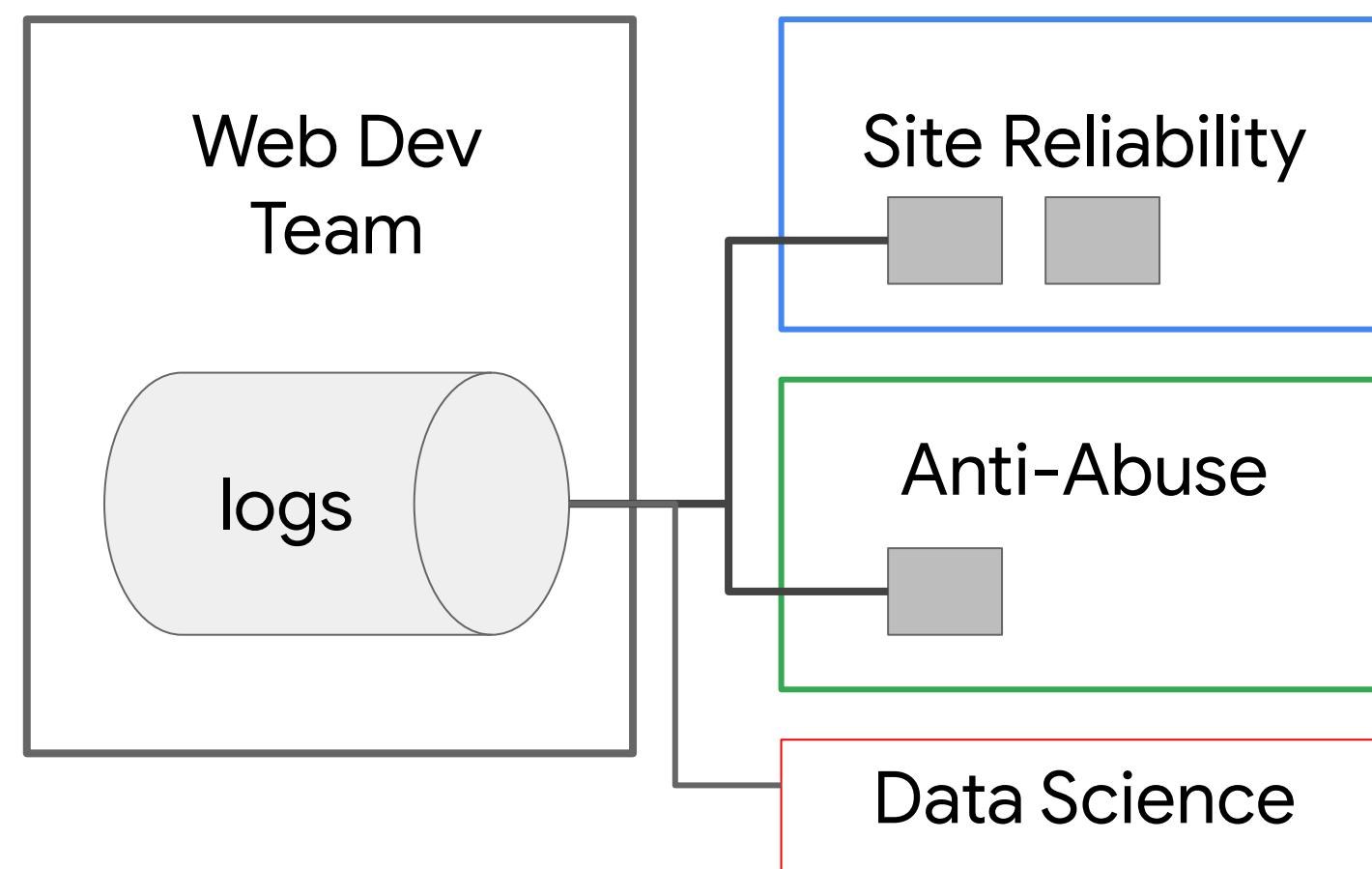


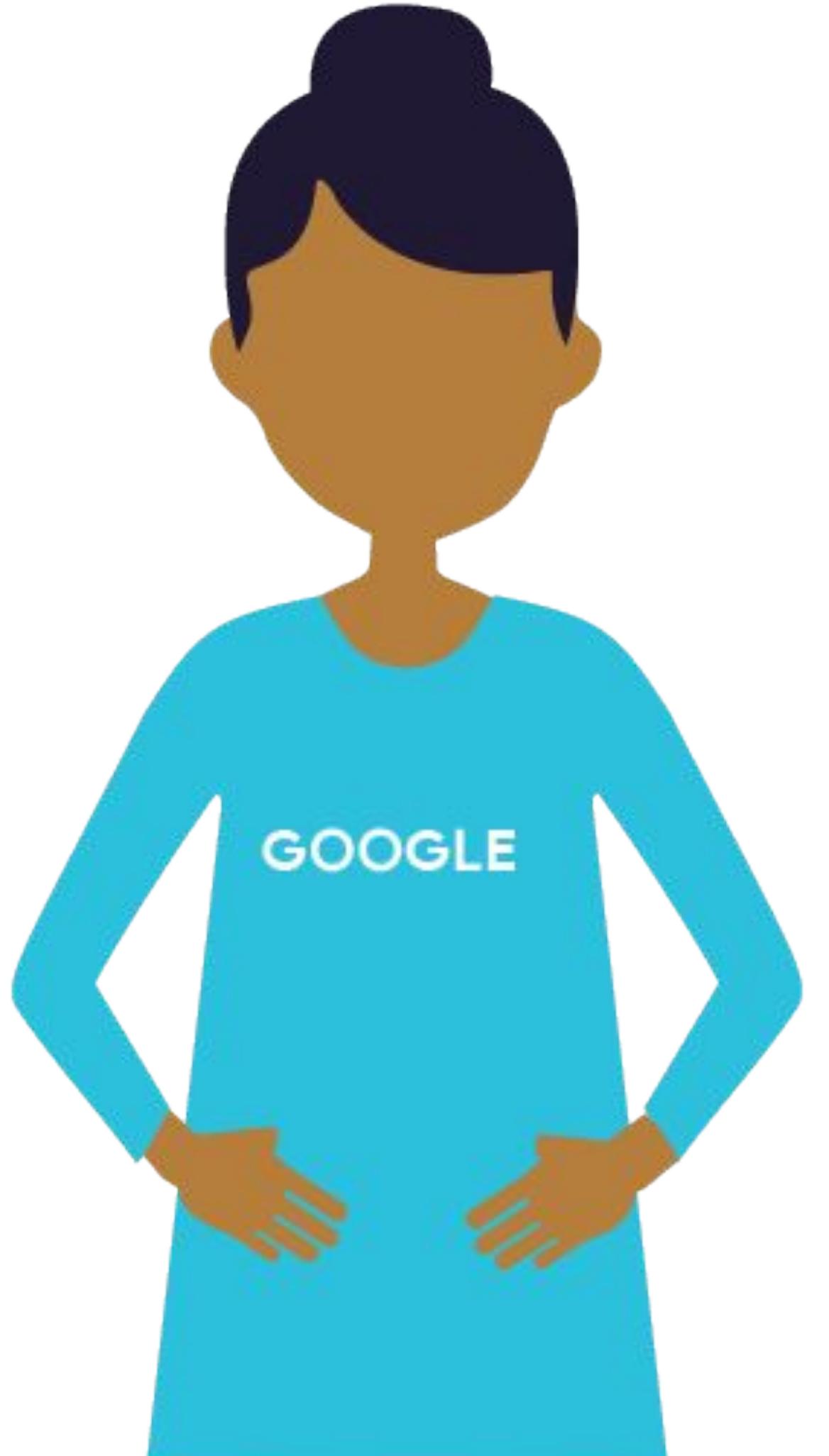
Decoupled upstream data producers



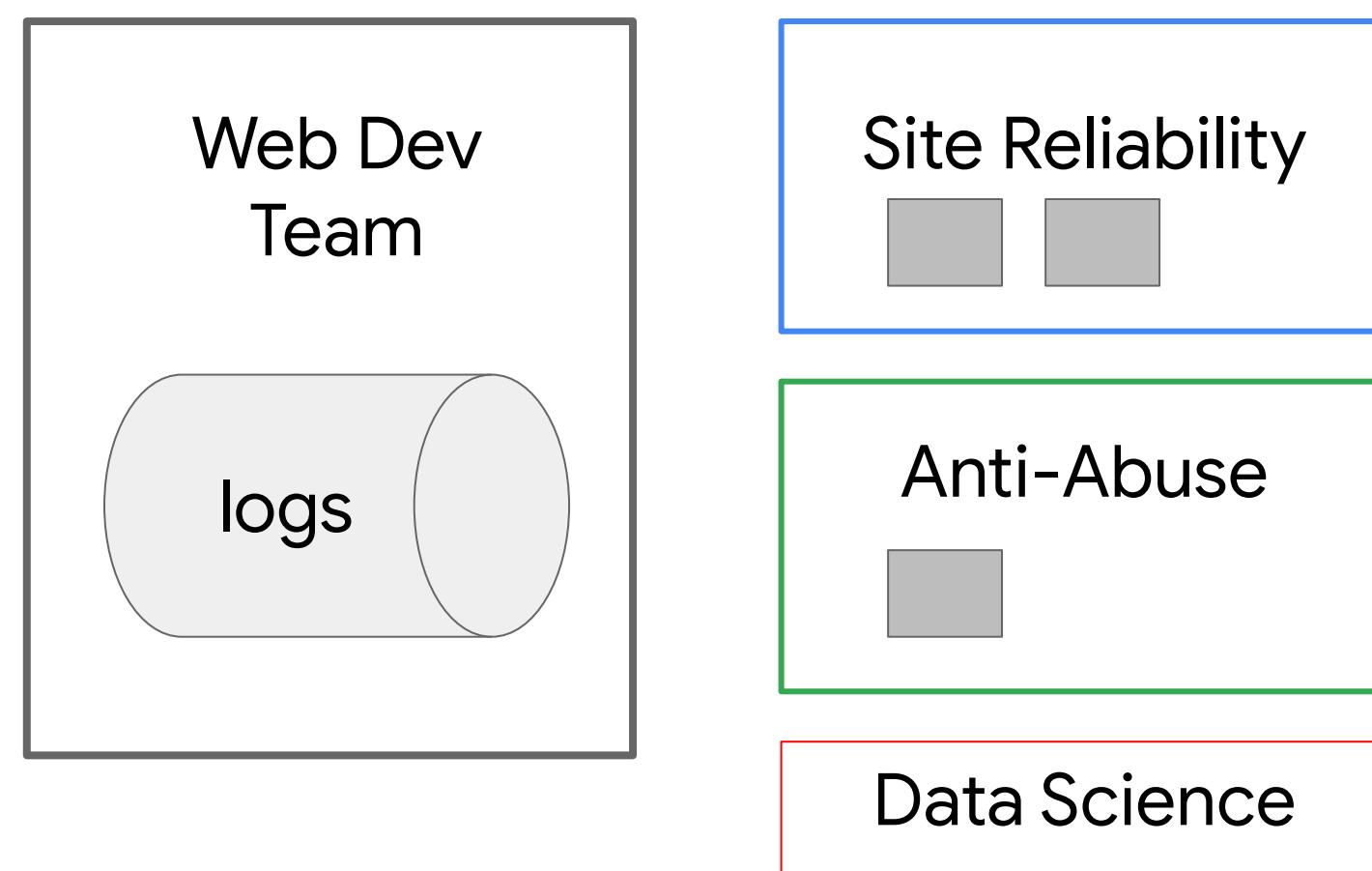


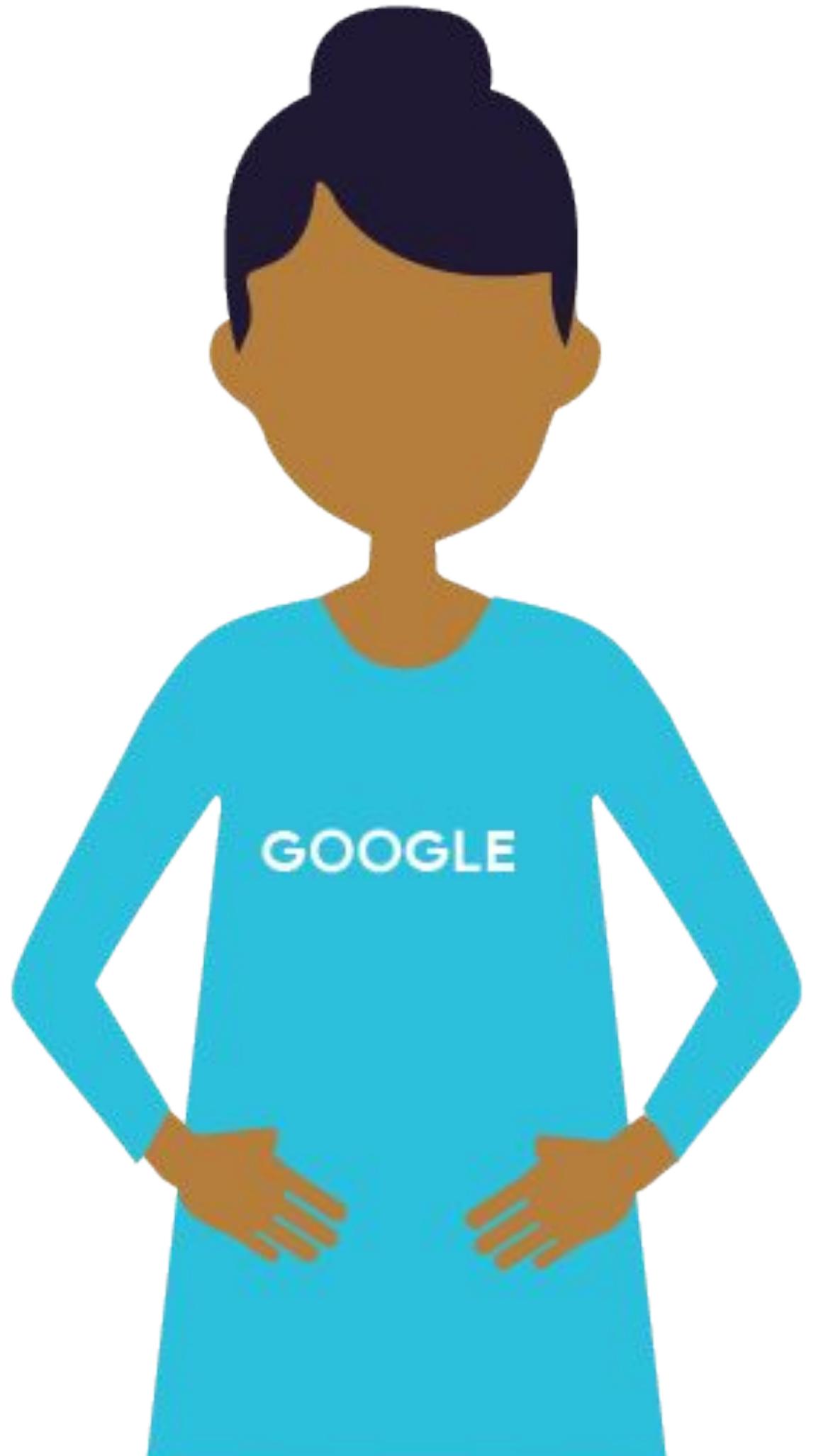
Decoupled upstream data producers



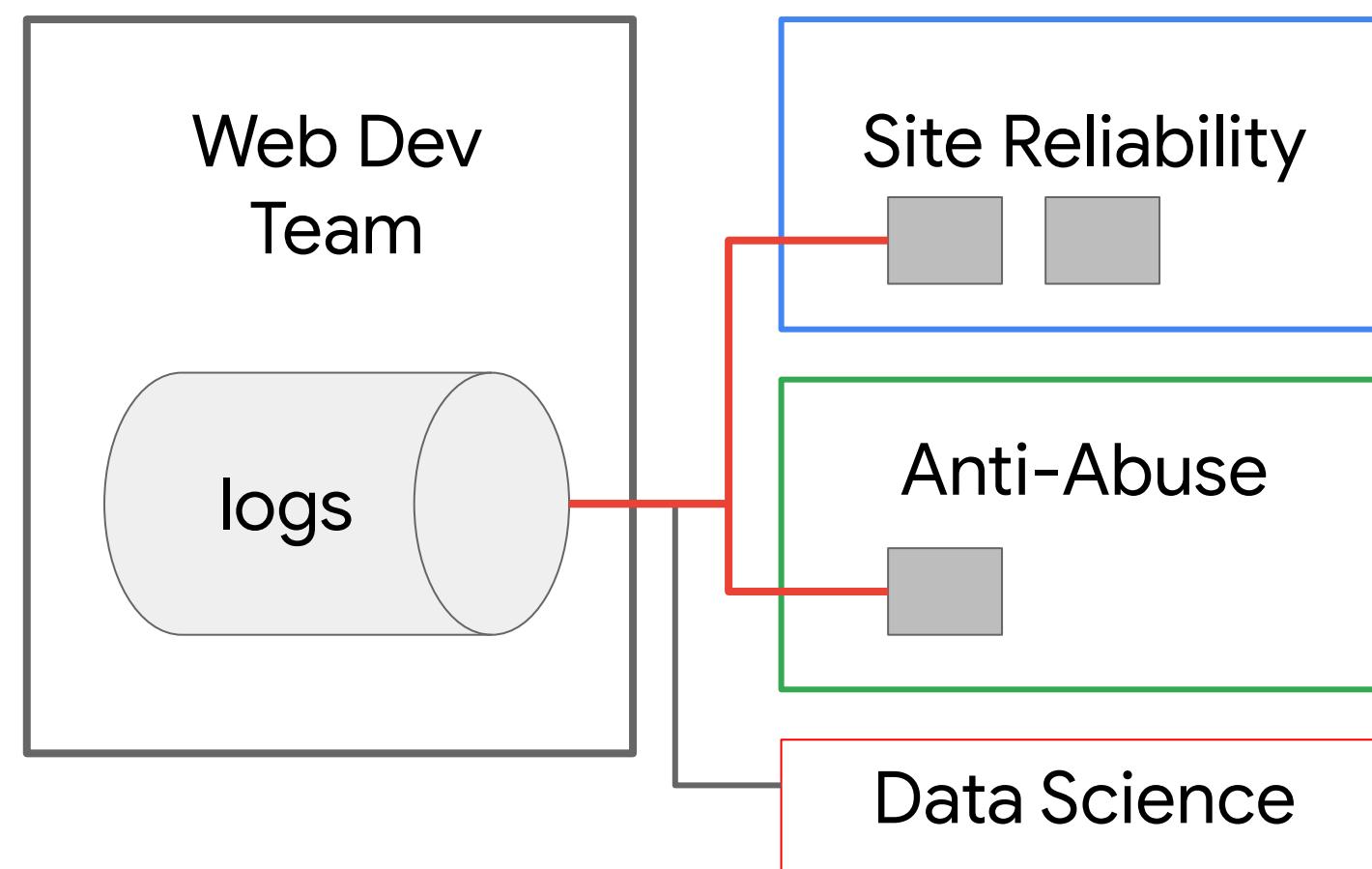


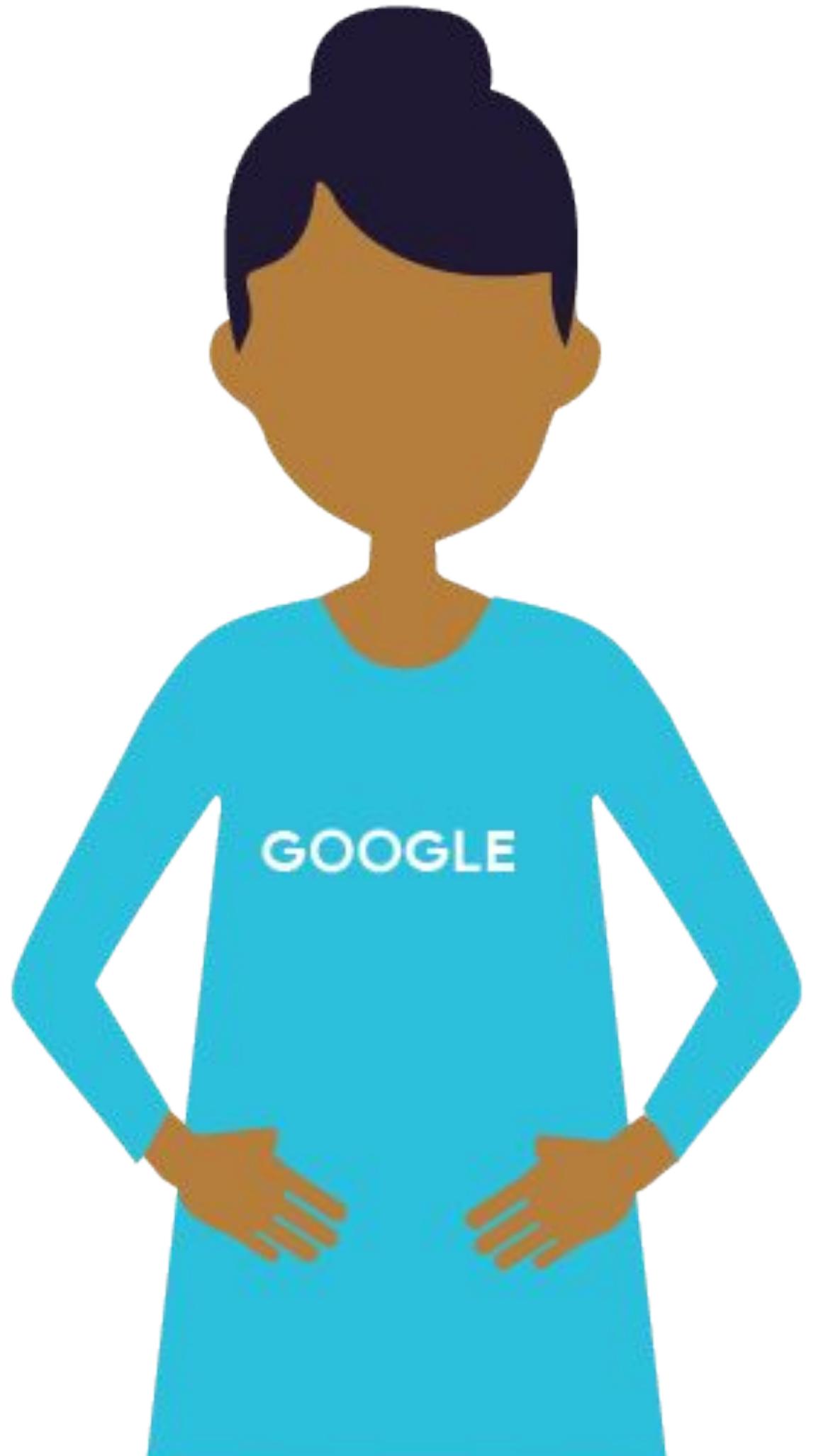
Decoupled upstream data producers



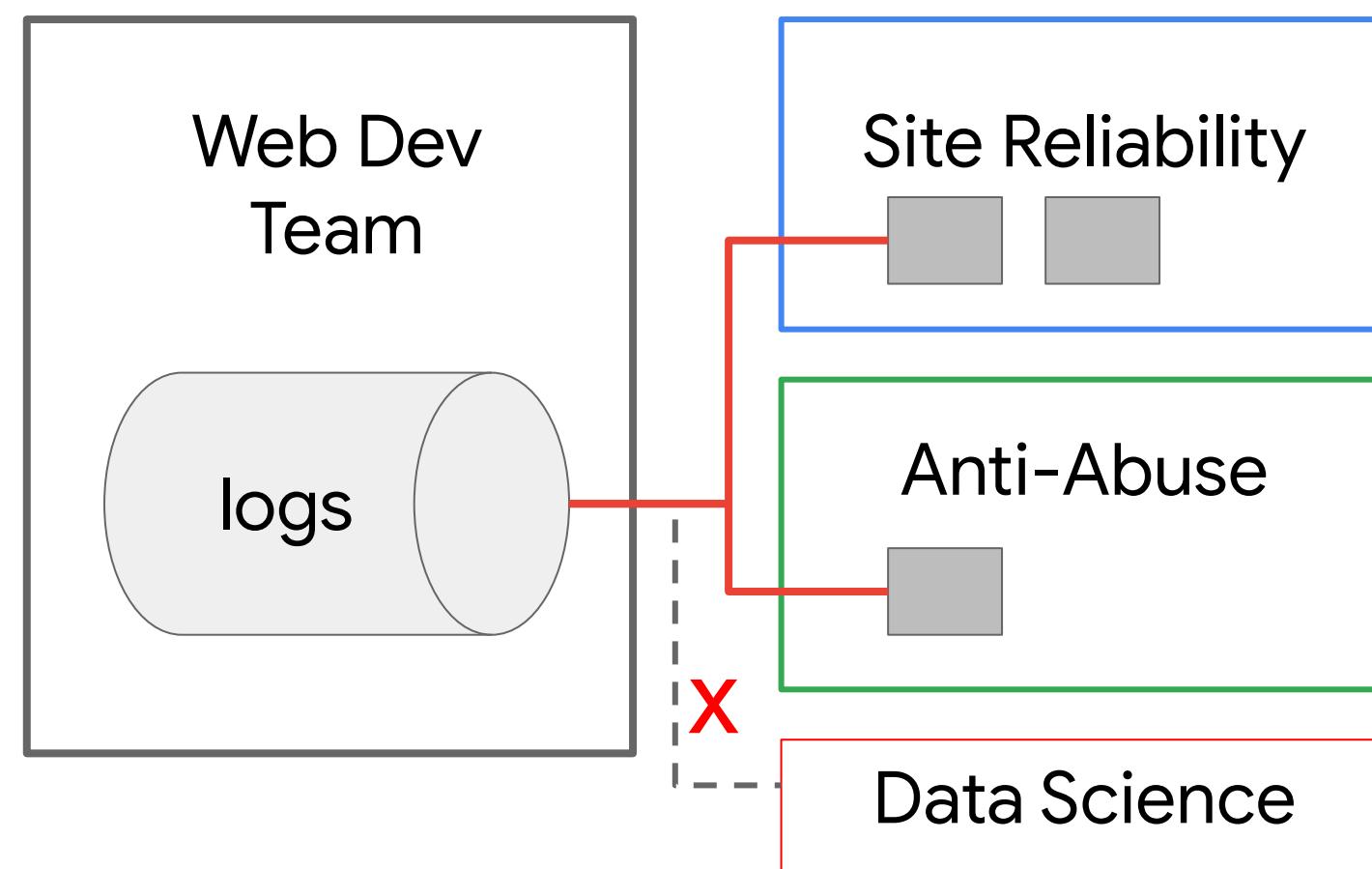


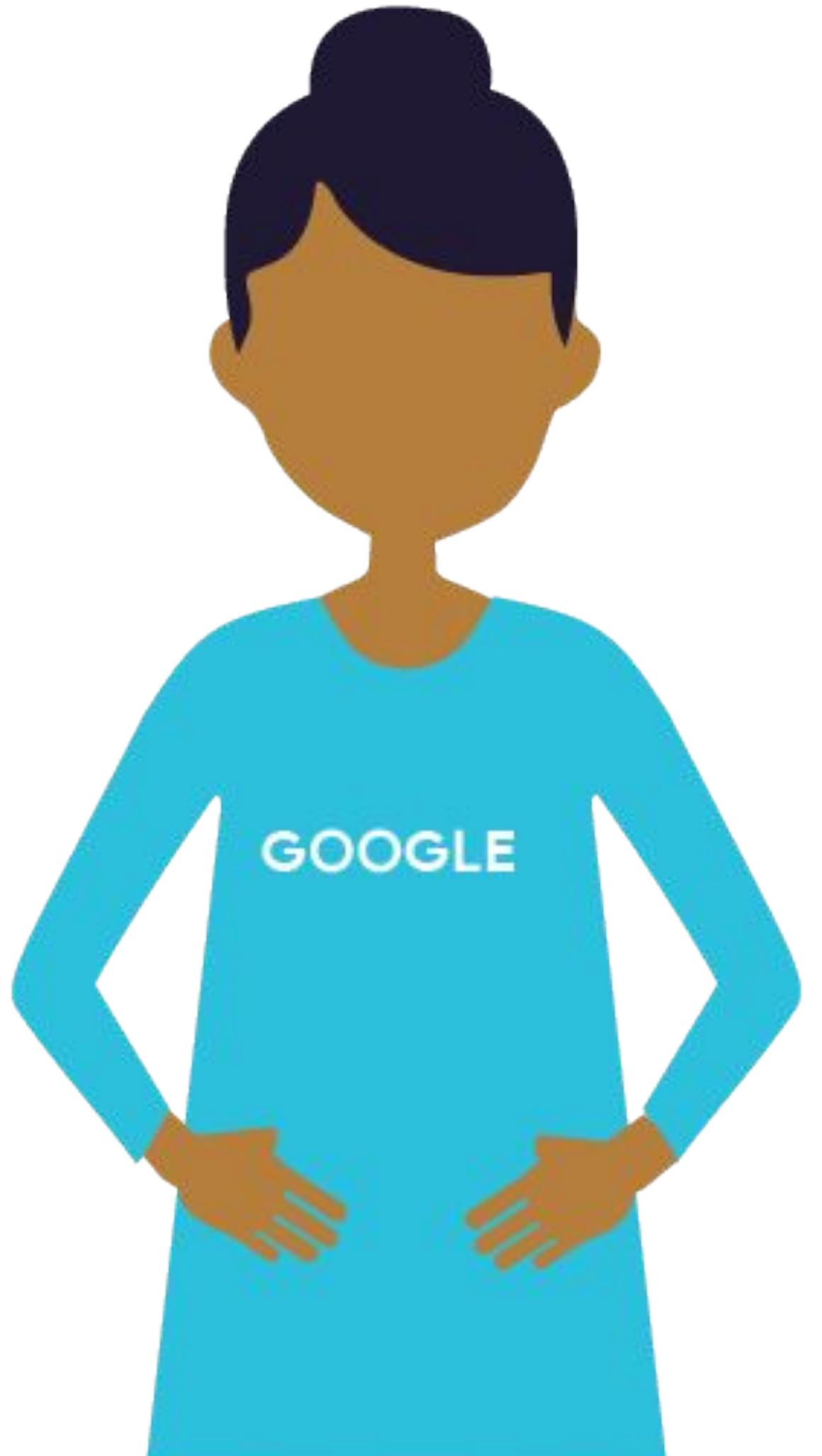
Decoupled upstream data producers



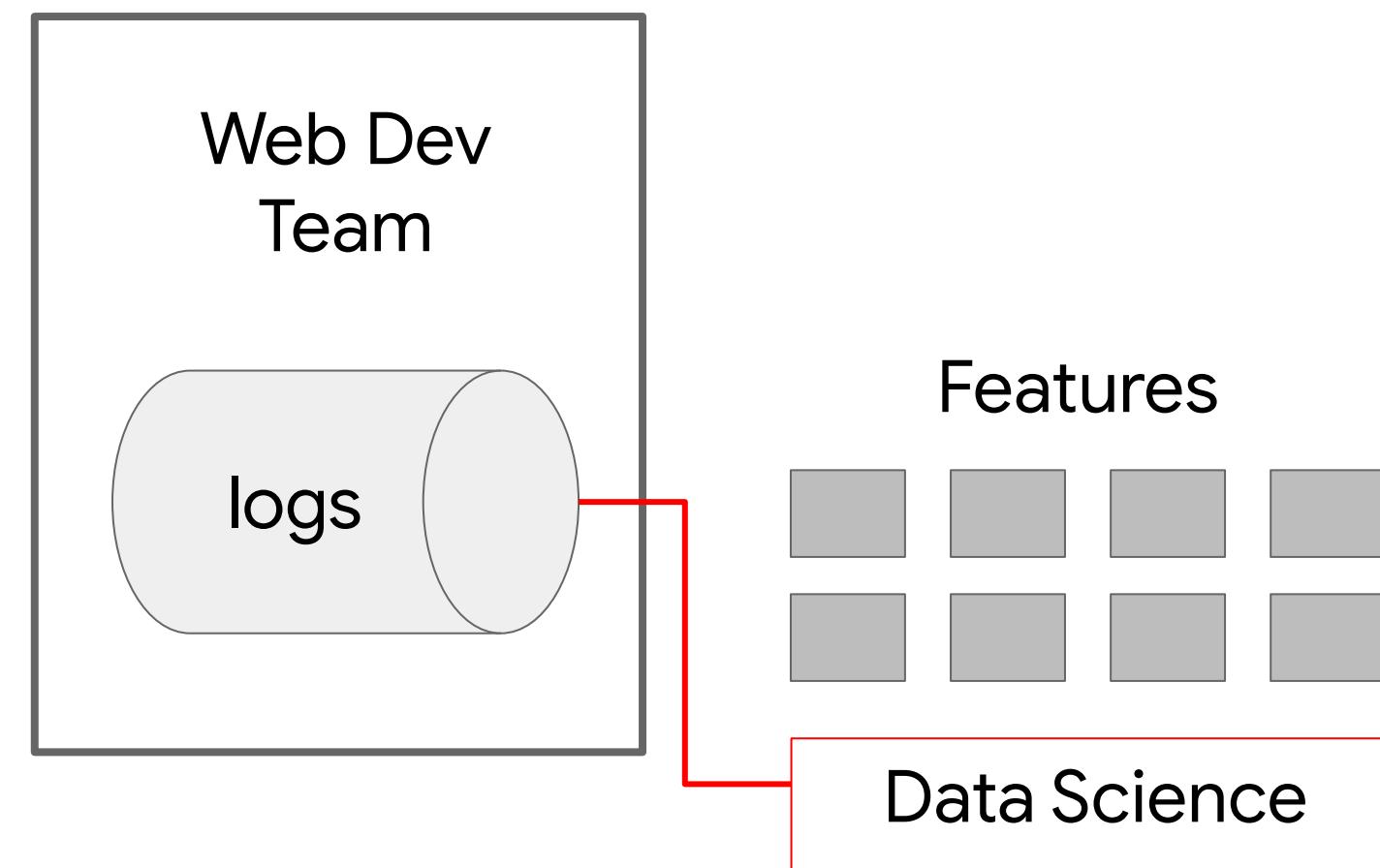


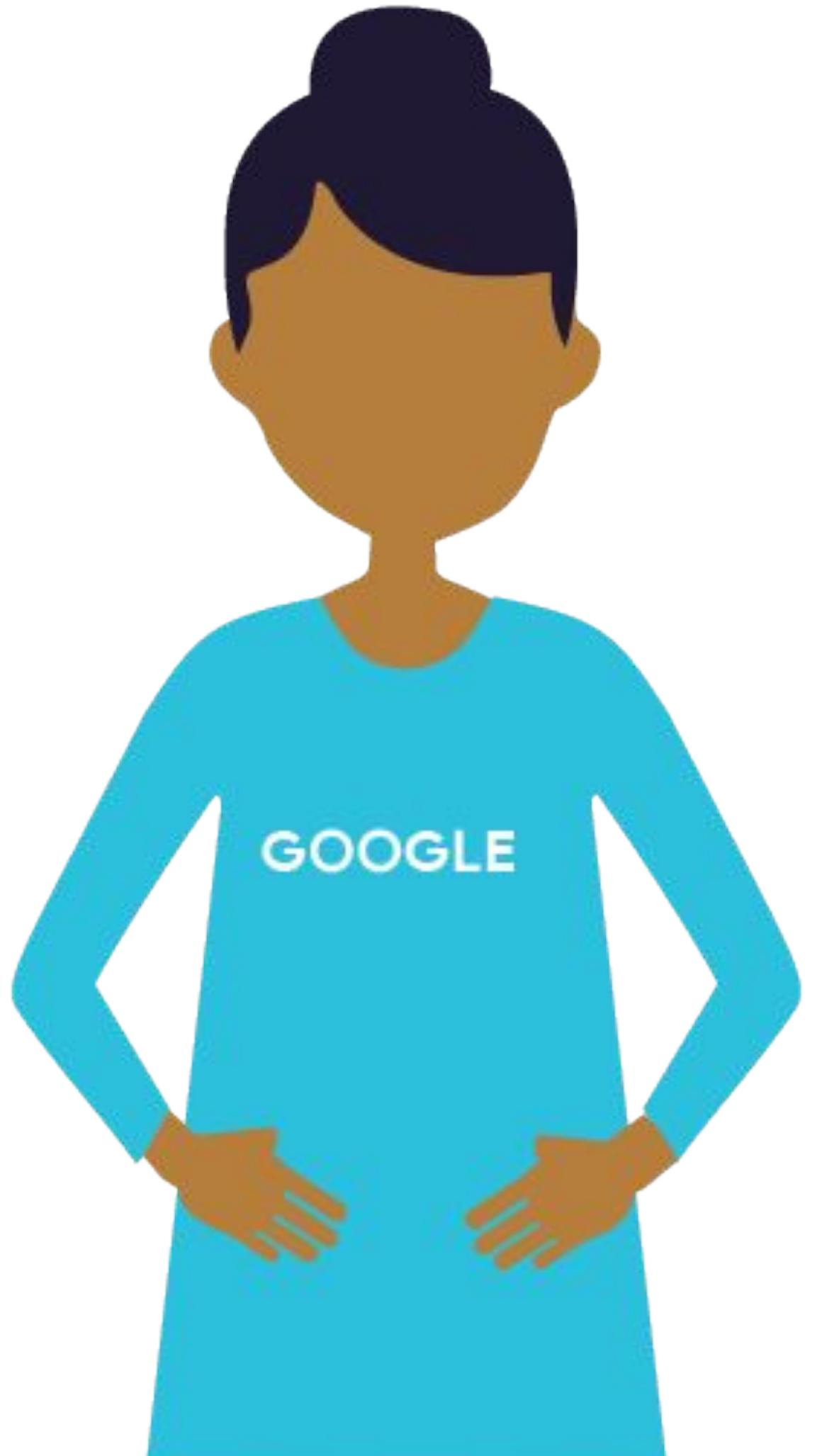
Decoupled upstream data producers



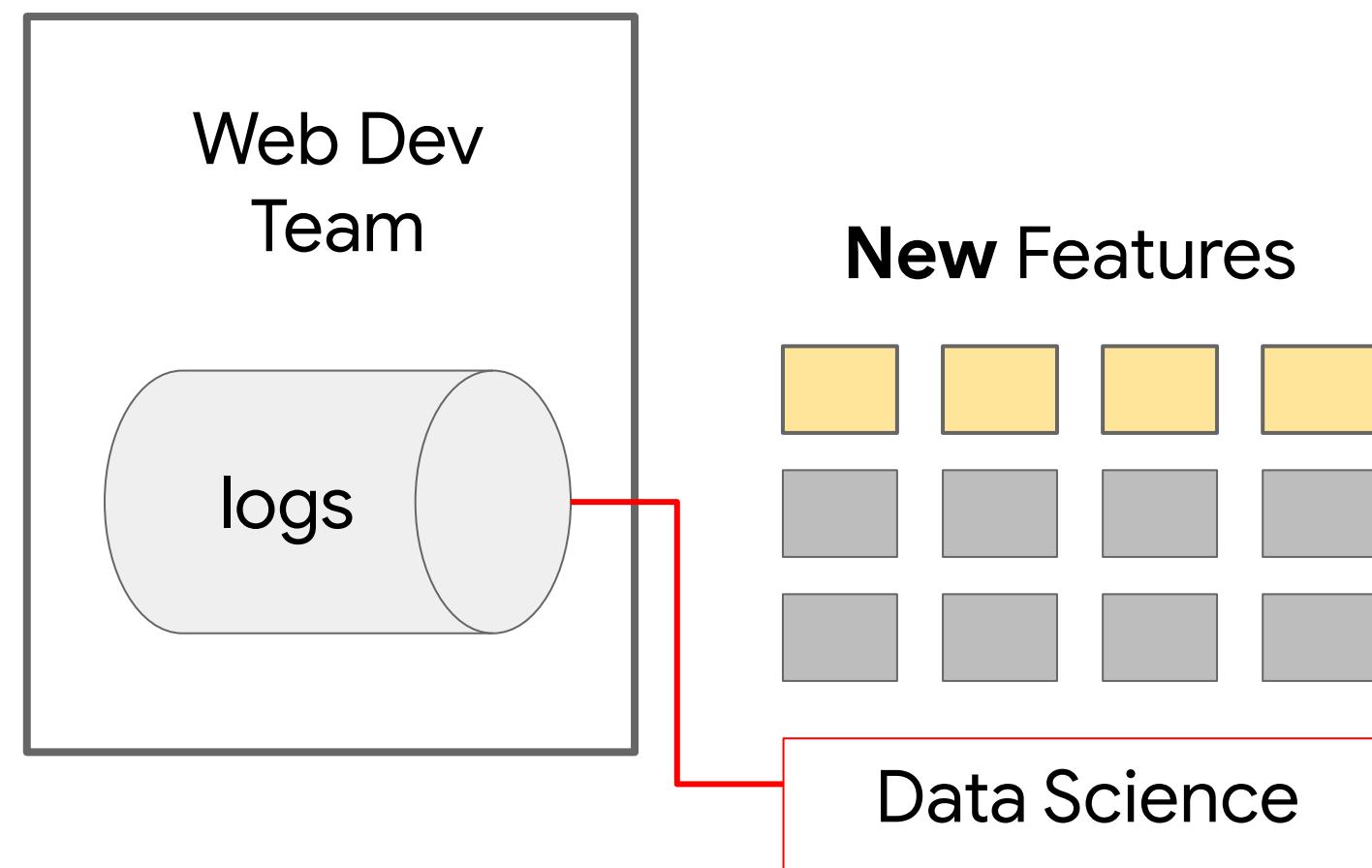


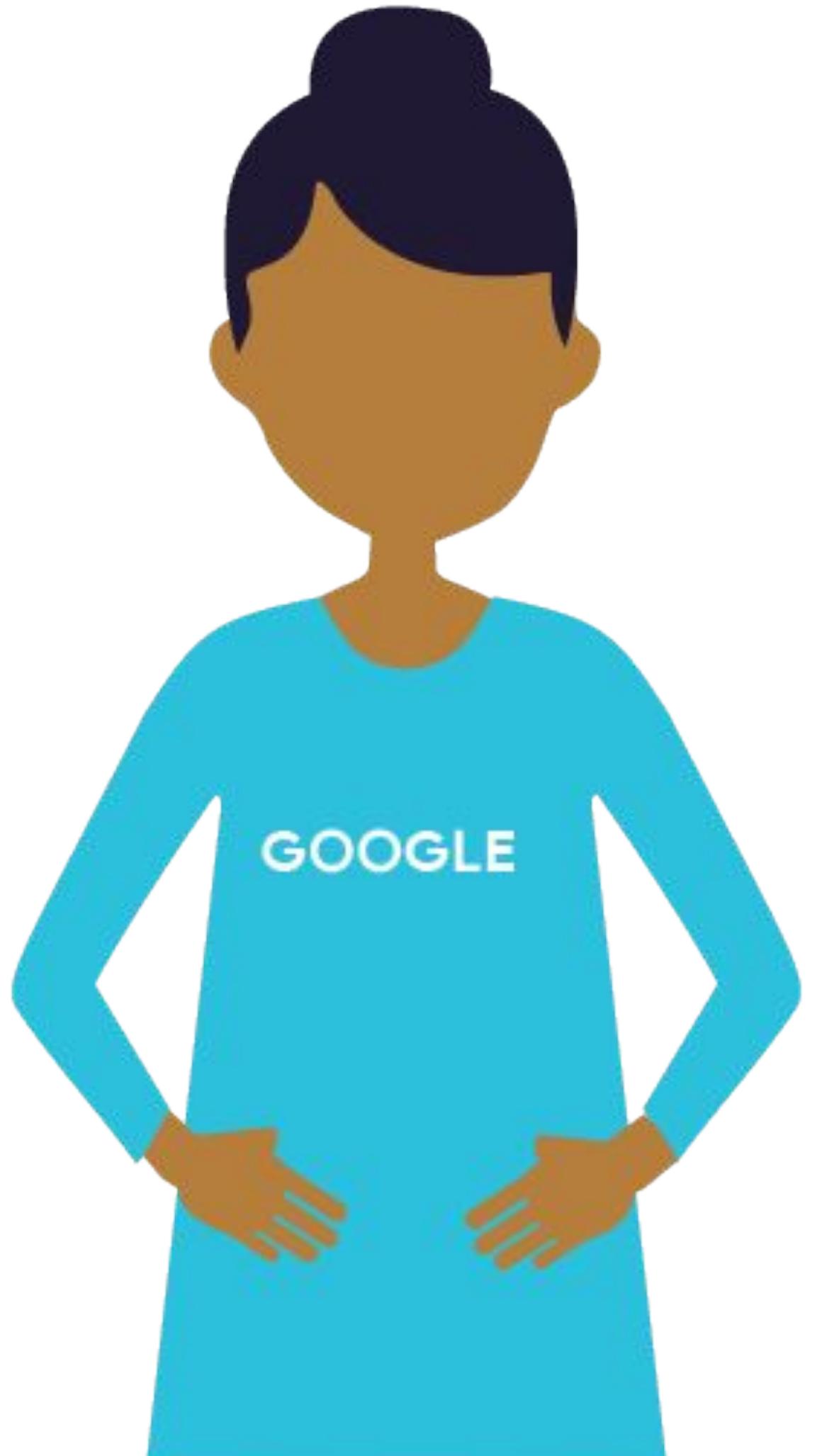
Underutilized data dependencies



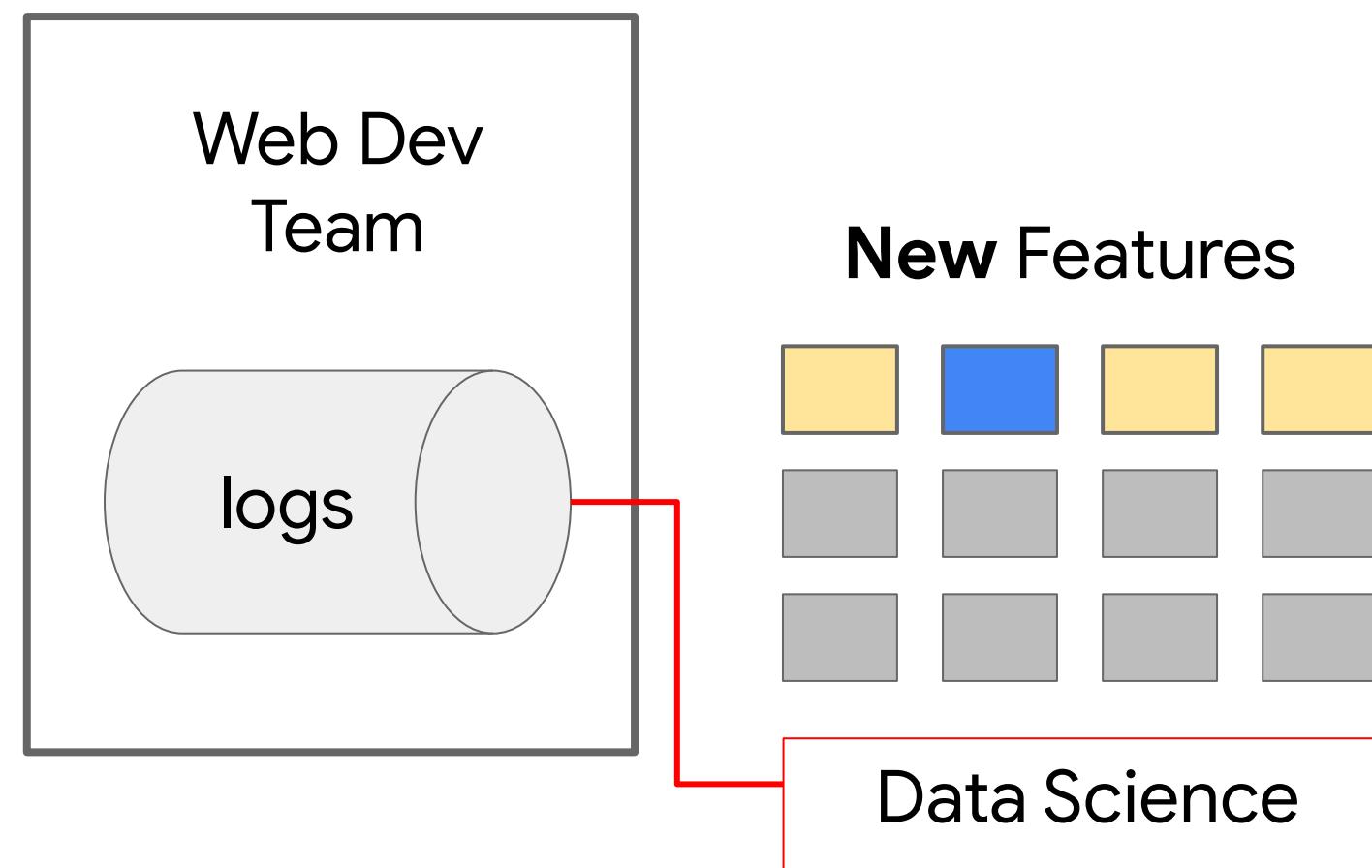


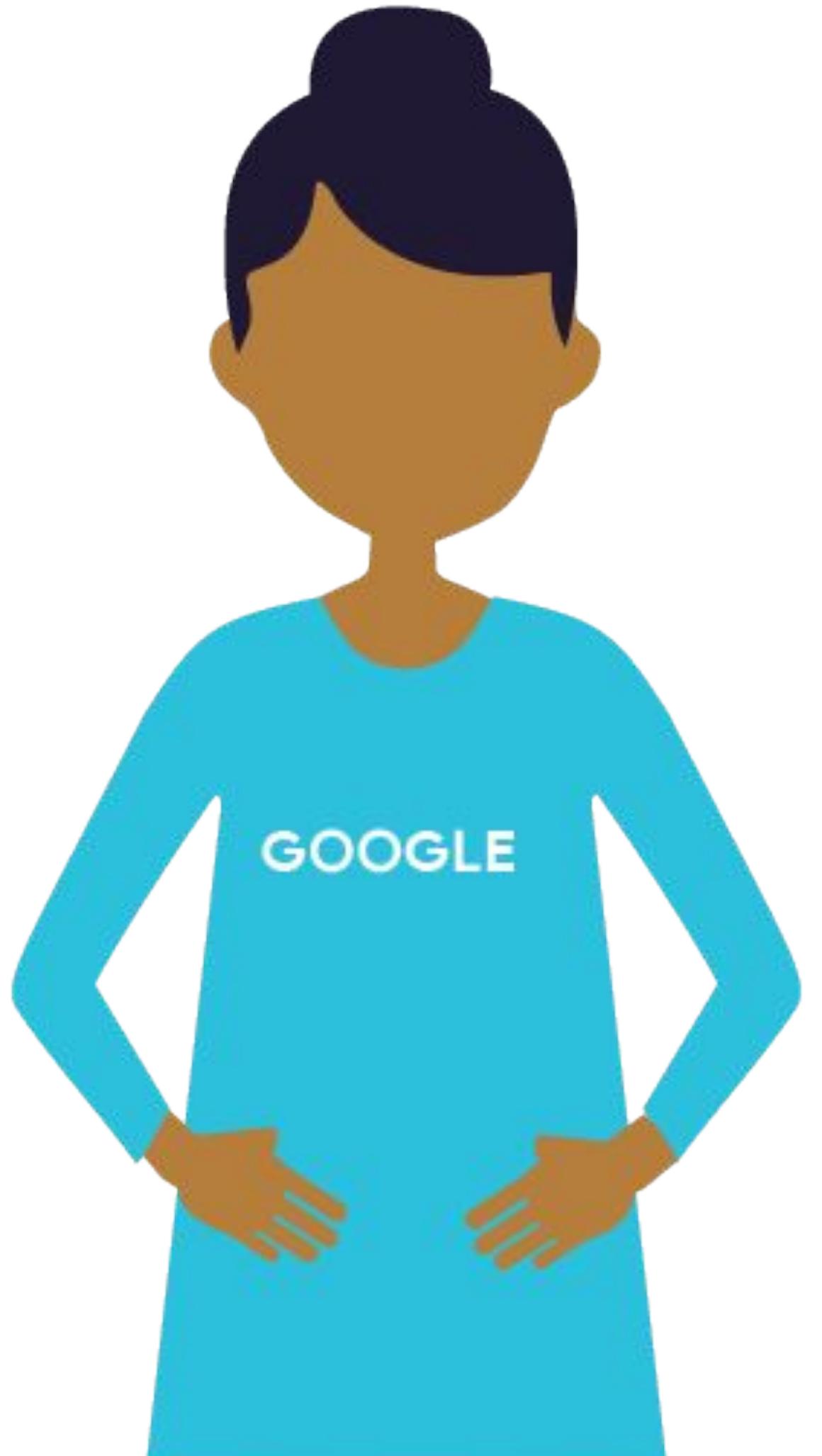
Underutilized data dependencies



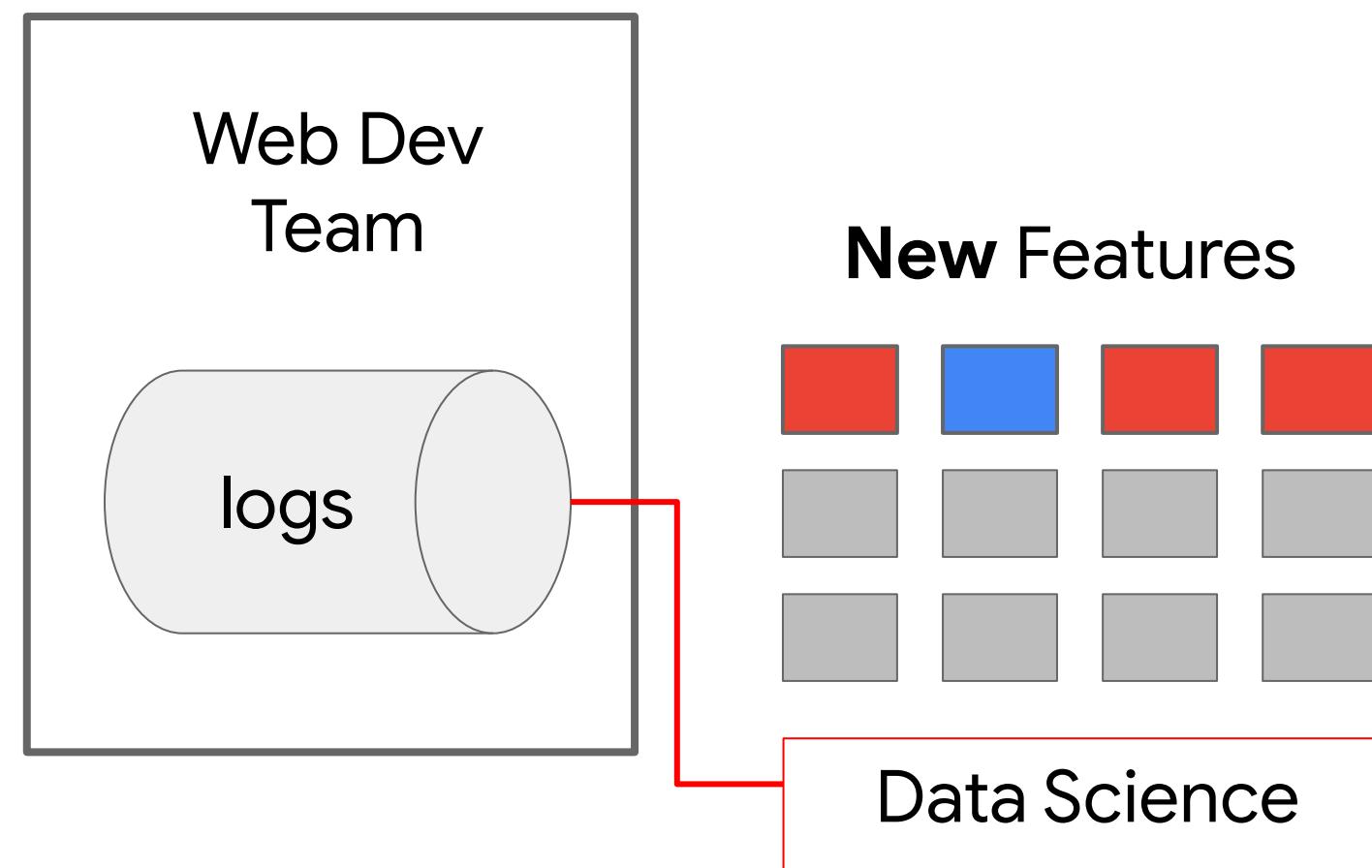


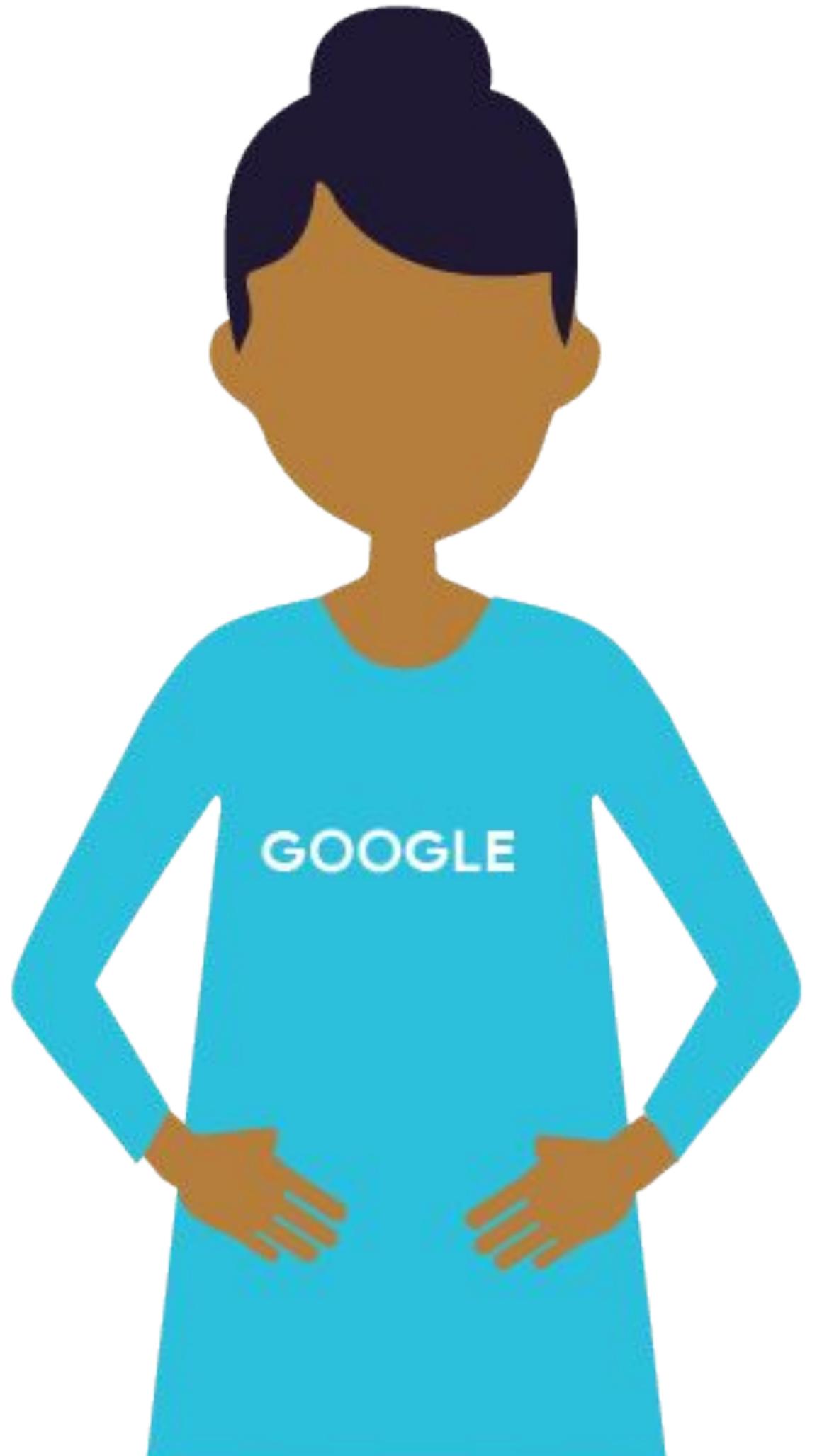
Underutilized data dependencies



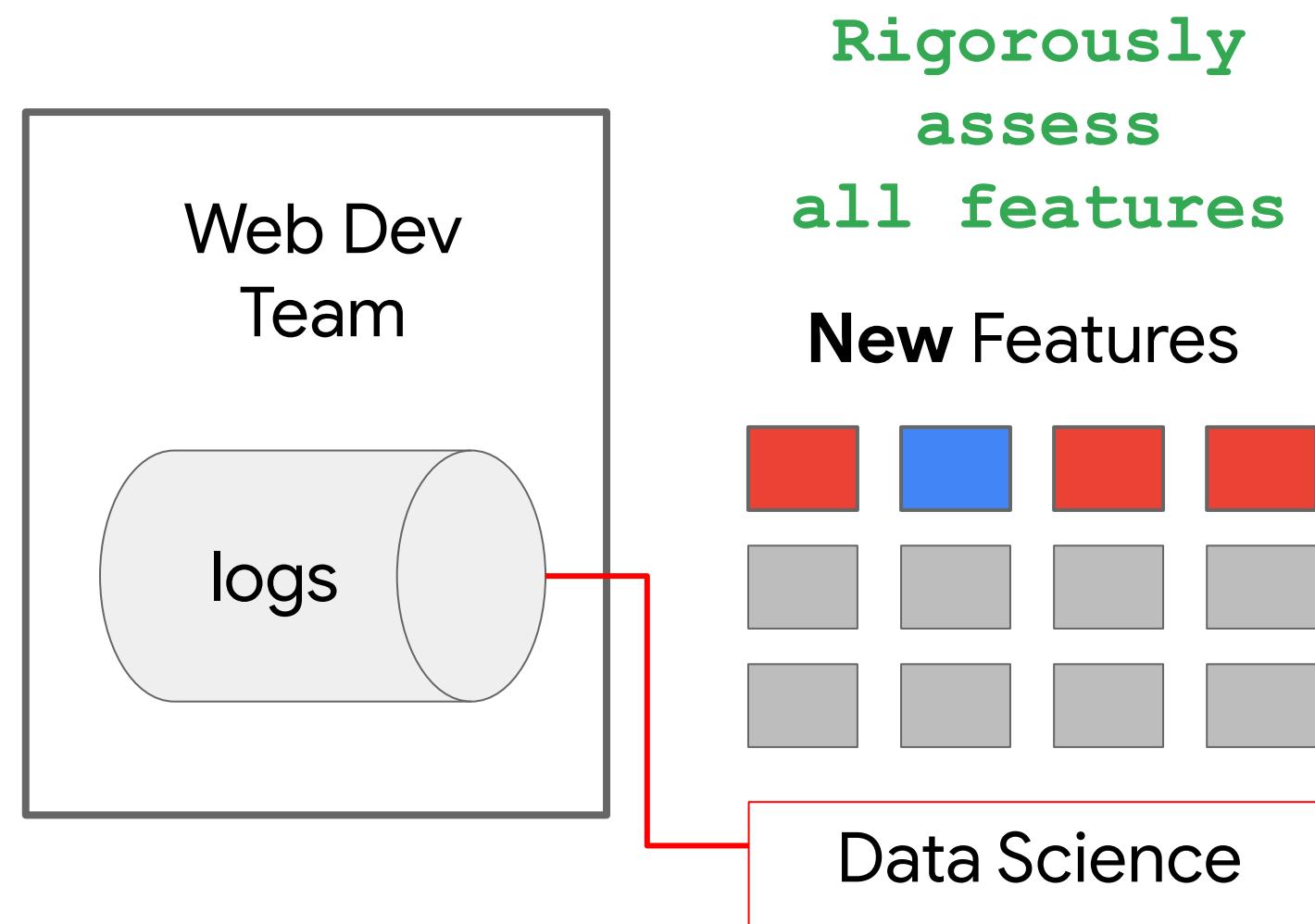


Underutilized data dependencies





Underutilized data dependencies



Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

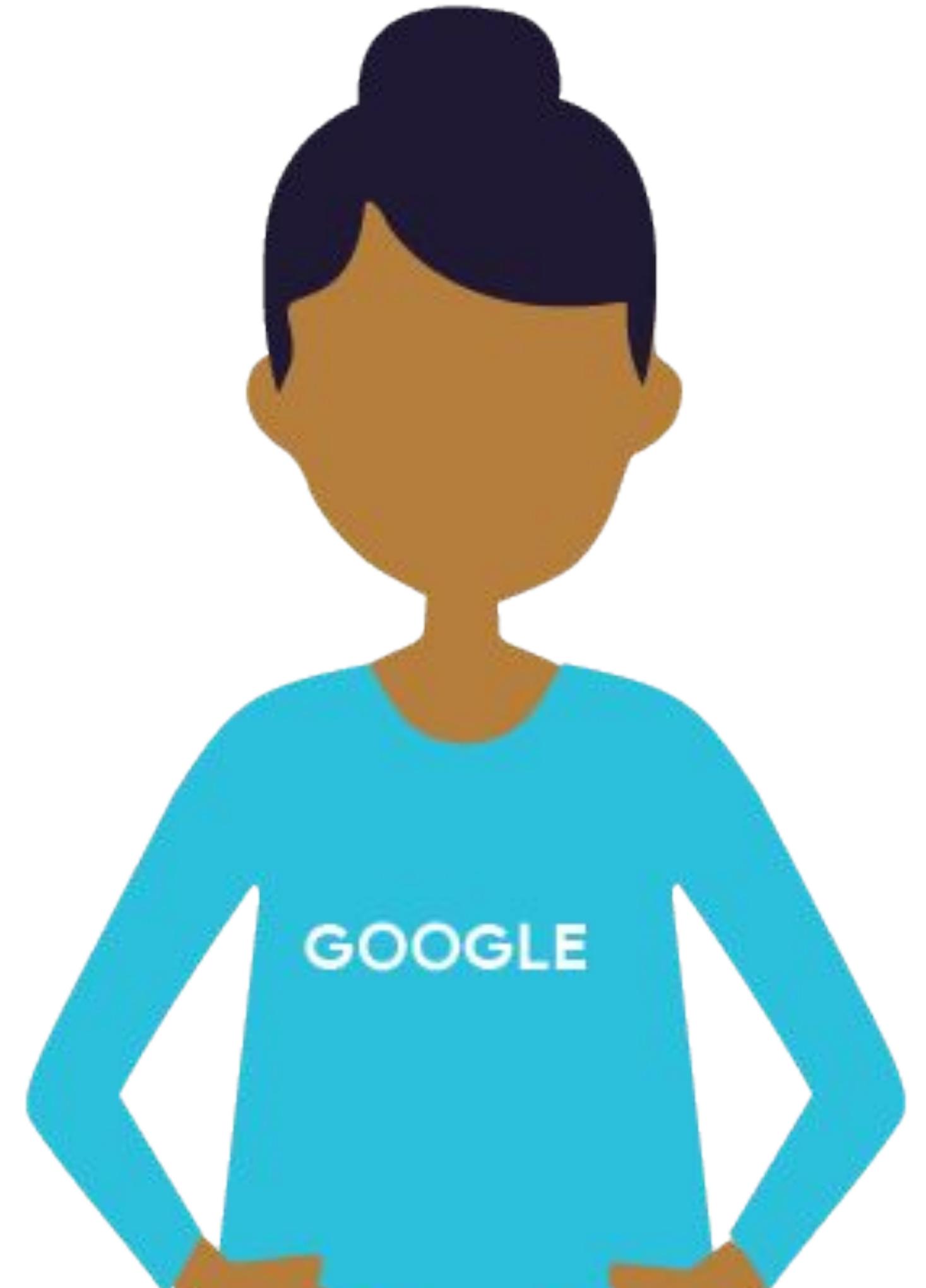
Lesson Title: Adapting to Data: Changing Distributions

Presenter: Max Lotstein

Format: Talking Head

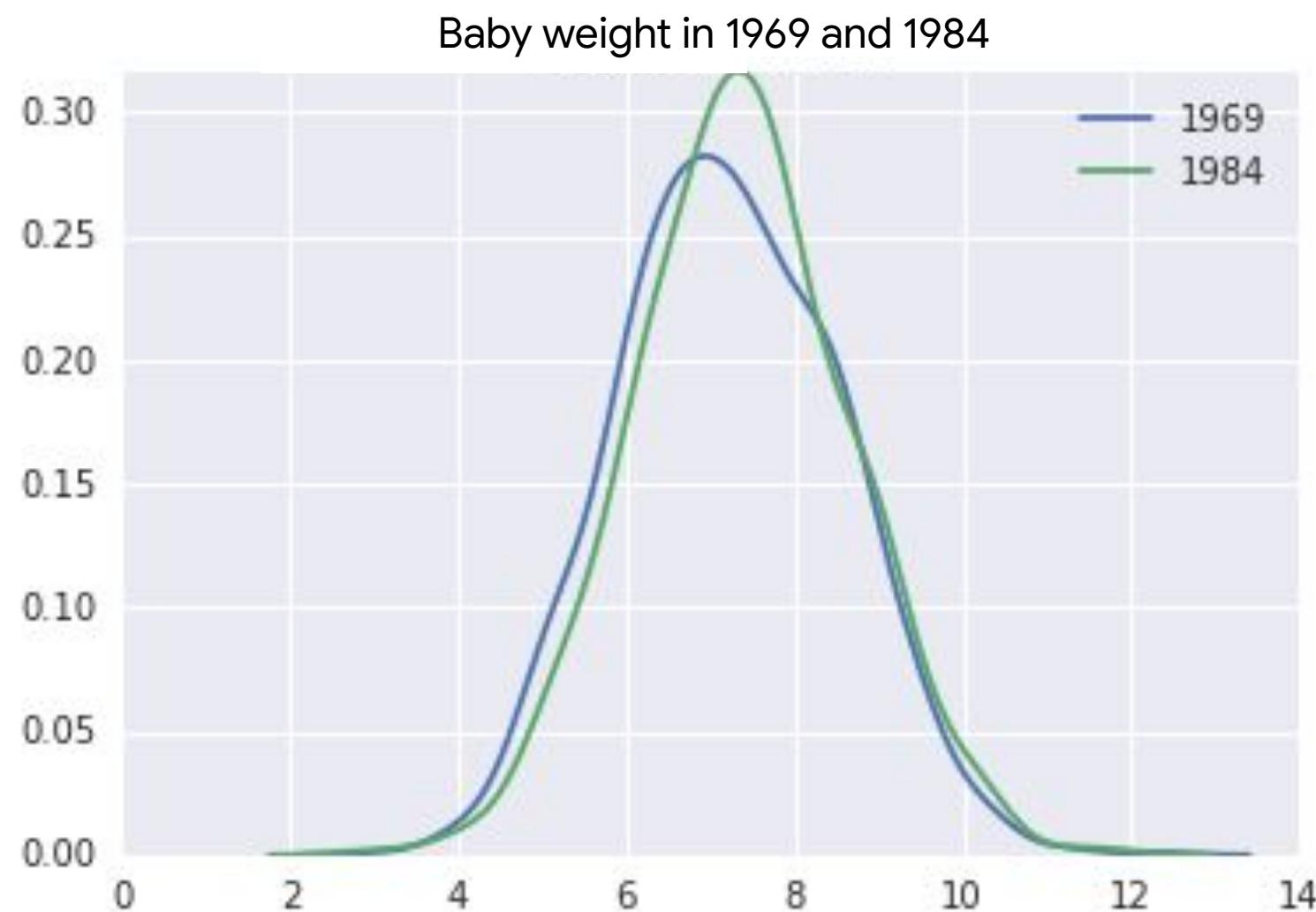
Video Name:

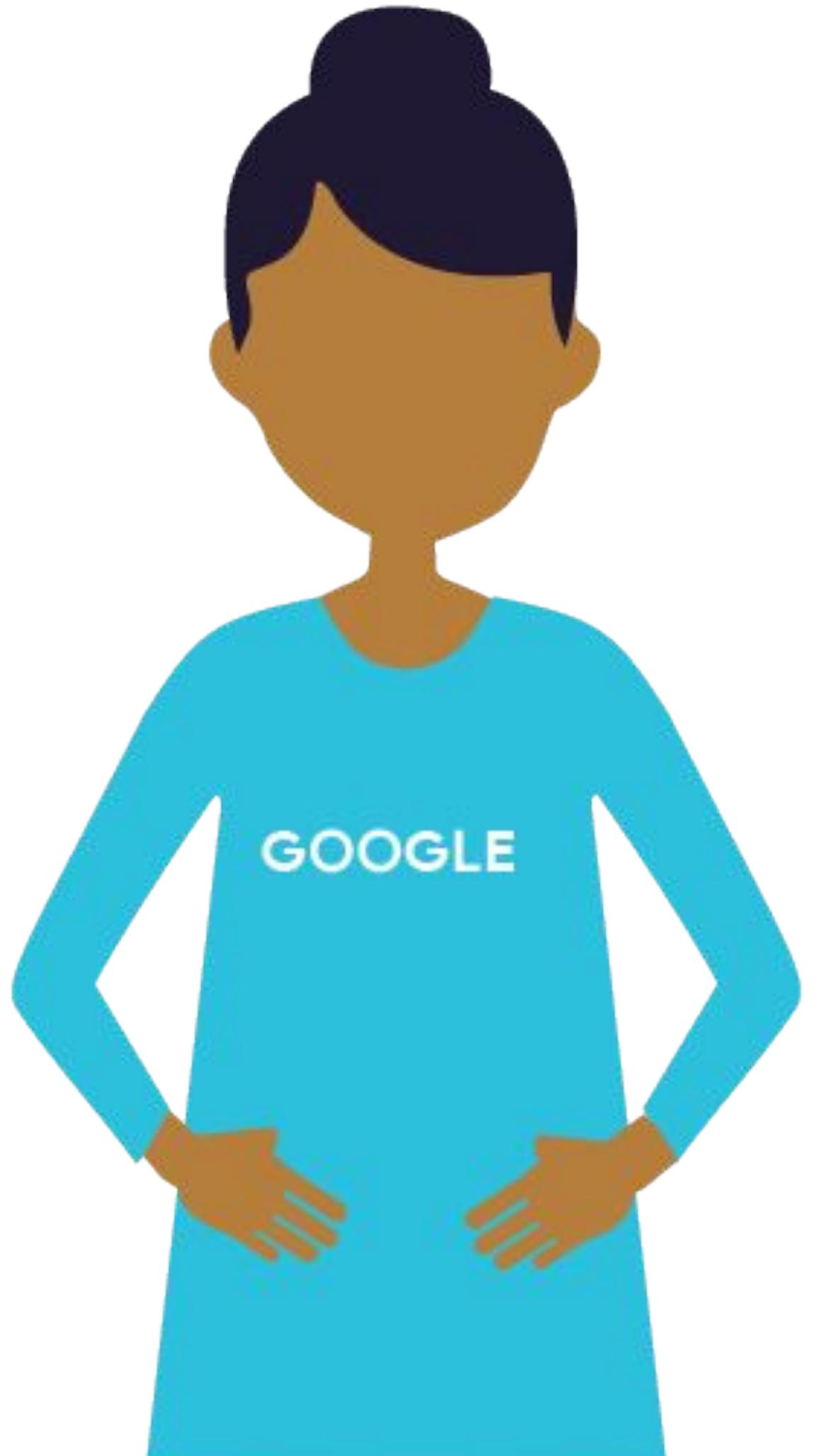
T-PSML-0_3_I4_adapting_to_data:_changing_distributions





Distributions change

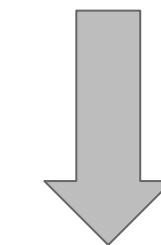




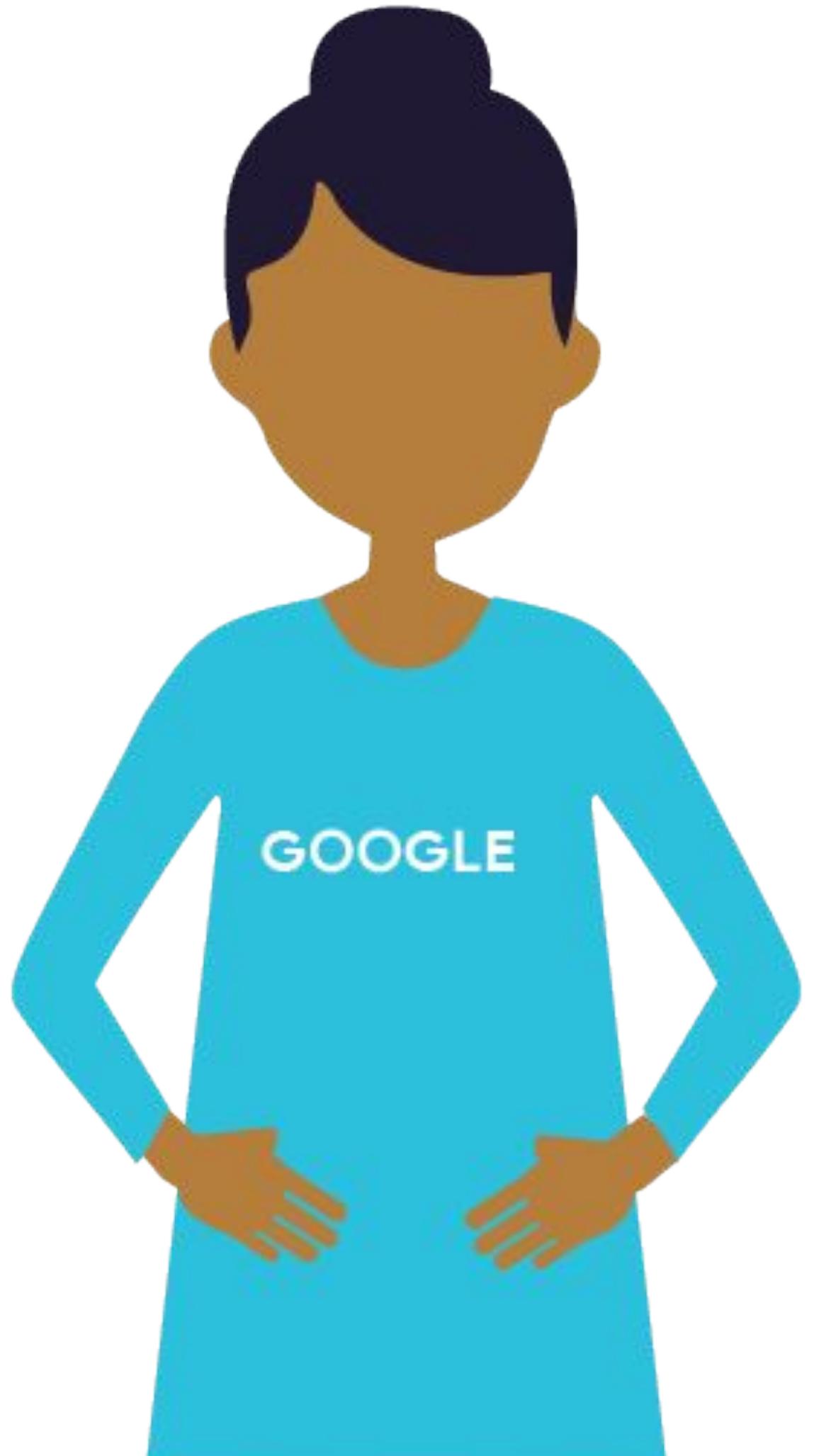
Distributions change

Zip Codes

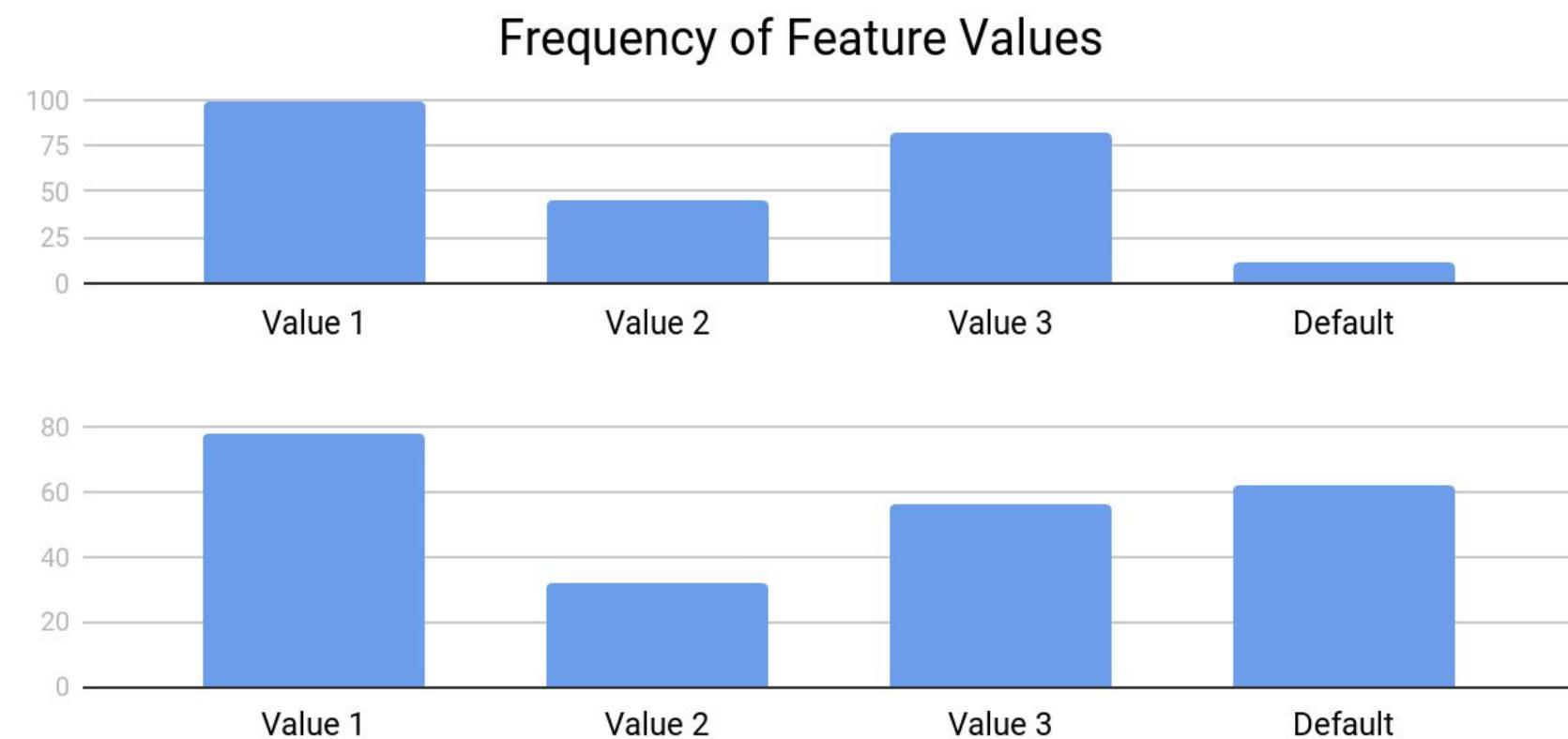
99501
87506
63141
04032
...

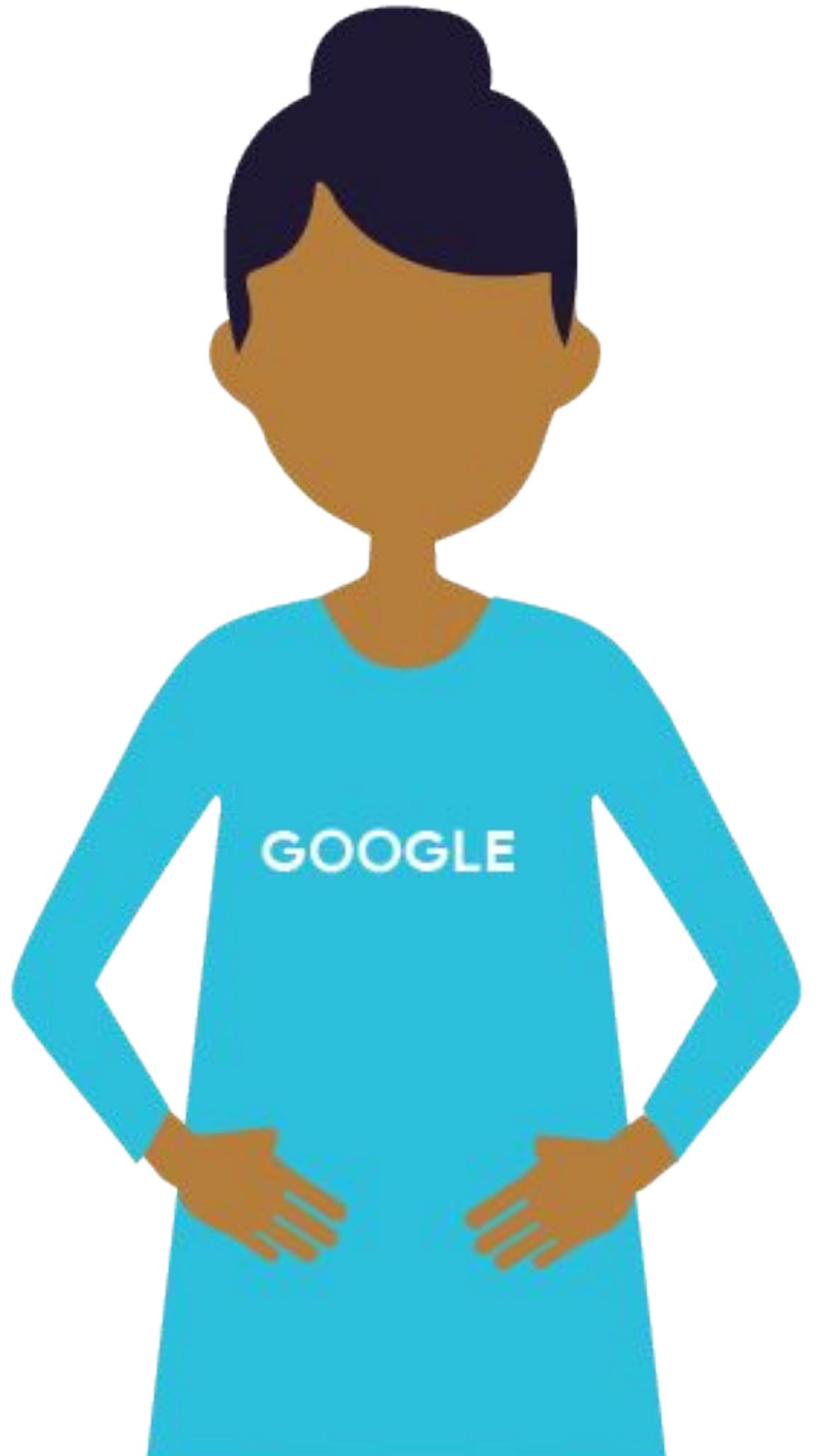


{*element1, element2, element3, ...*}

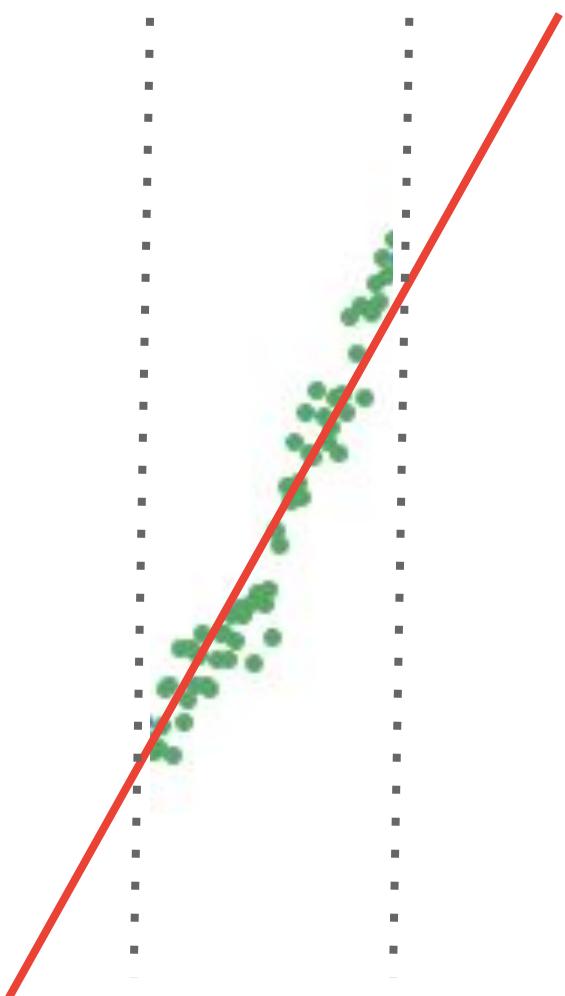


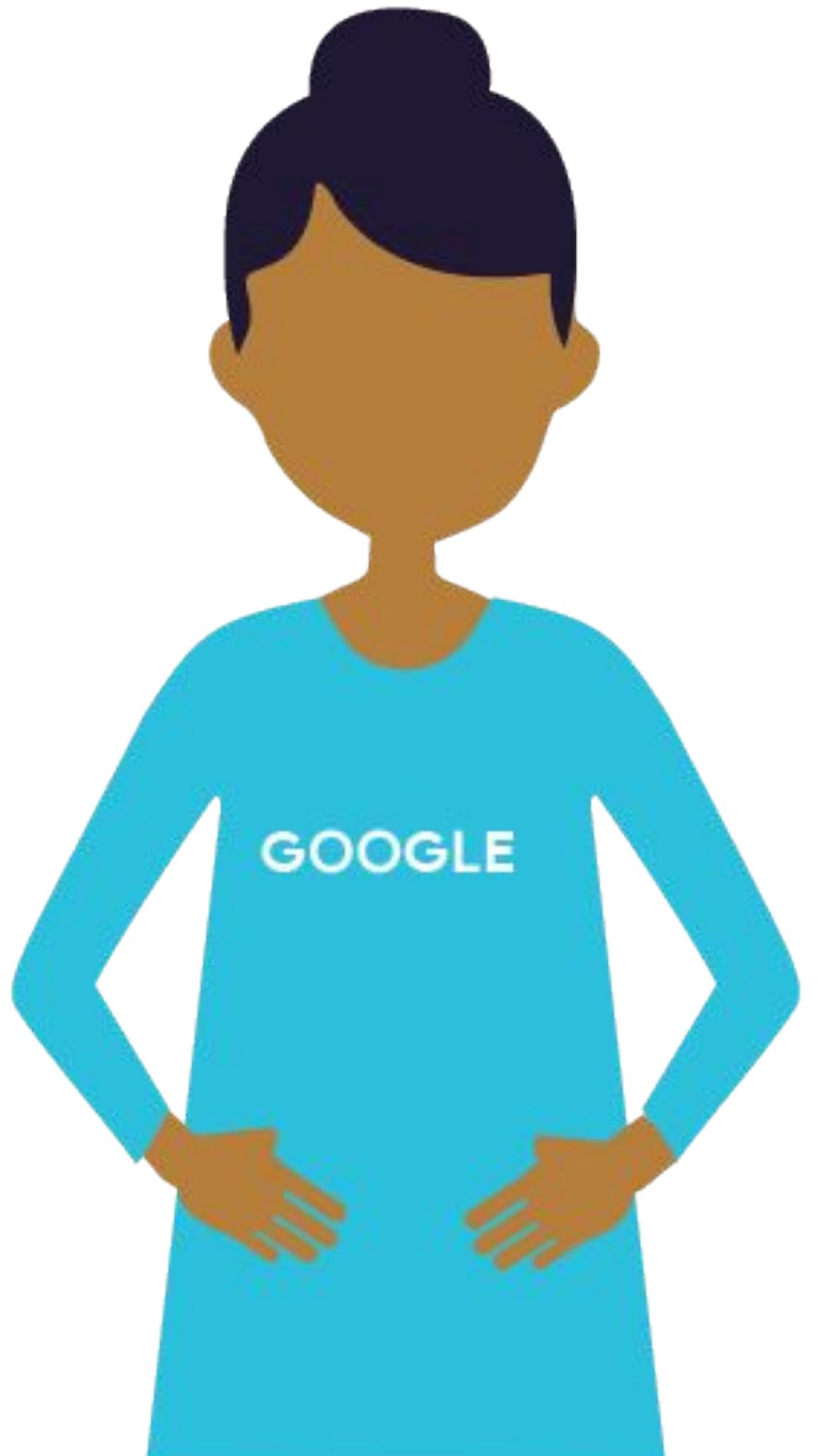
Distributions change



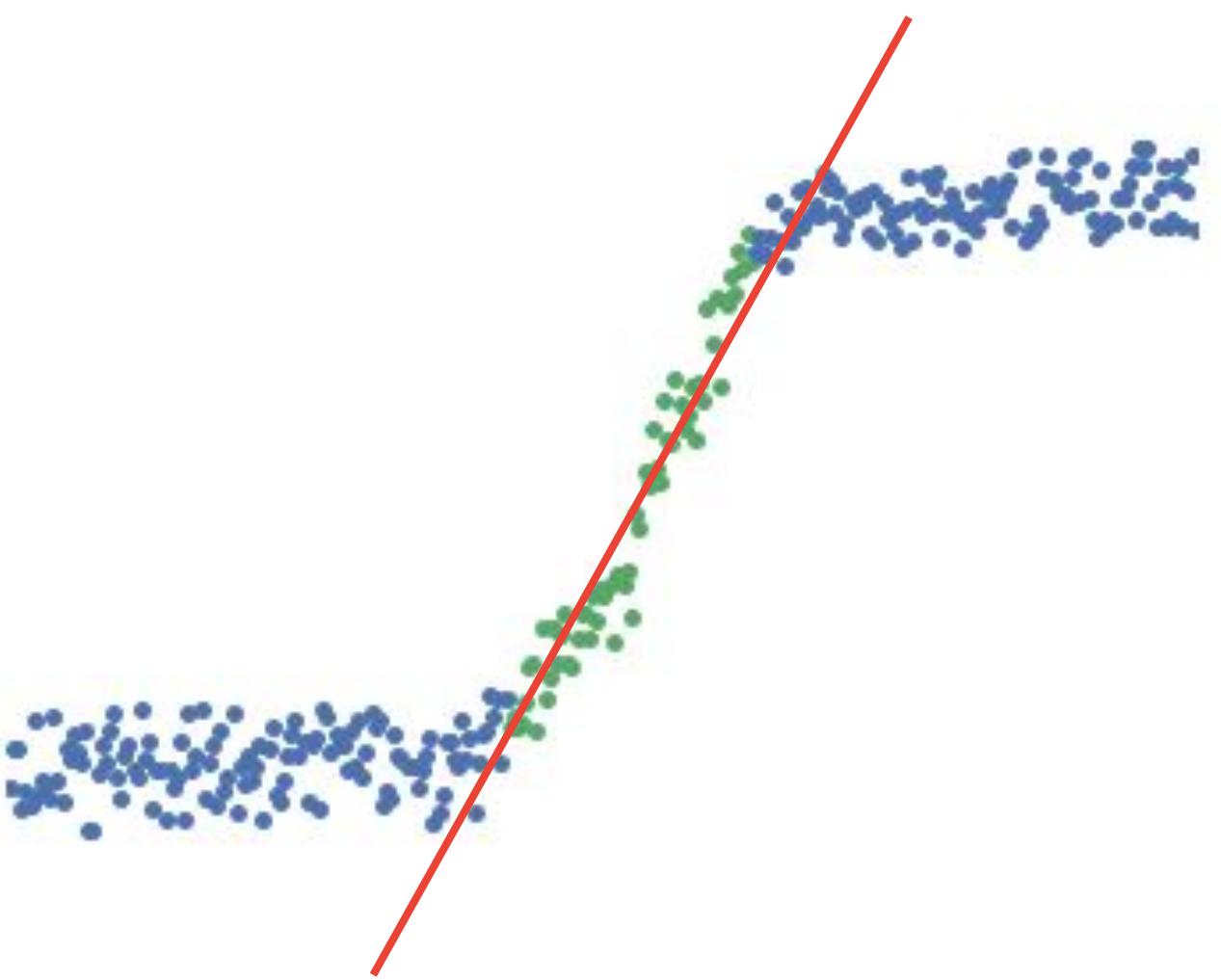


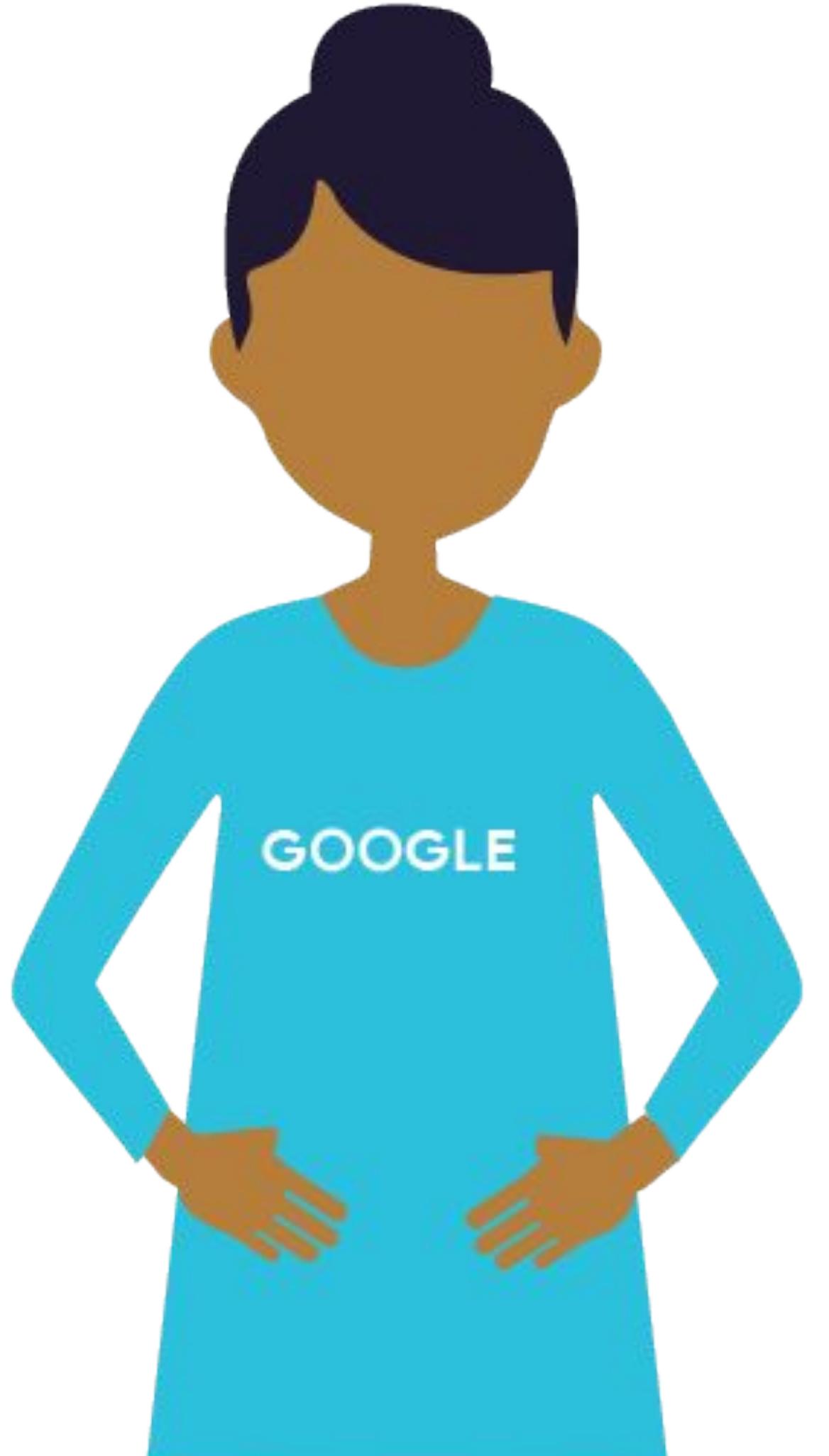
Distributions change





Distributions change





Distributions change

- Monitor descriptive statistics for your inputs and outputs
- Monitor your residuals as a function of your inputs
- Use custom weights in your loss function to emphasize data recency
- Use dynamic training architecture and regularly retrain your model

Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

Lesson Title: Adapting to Data: Lab

Presenter: Max Lotstein

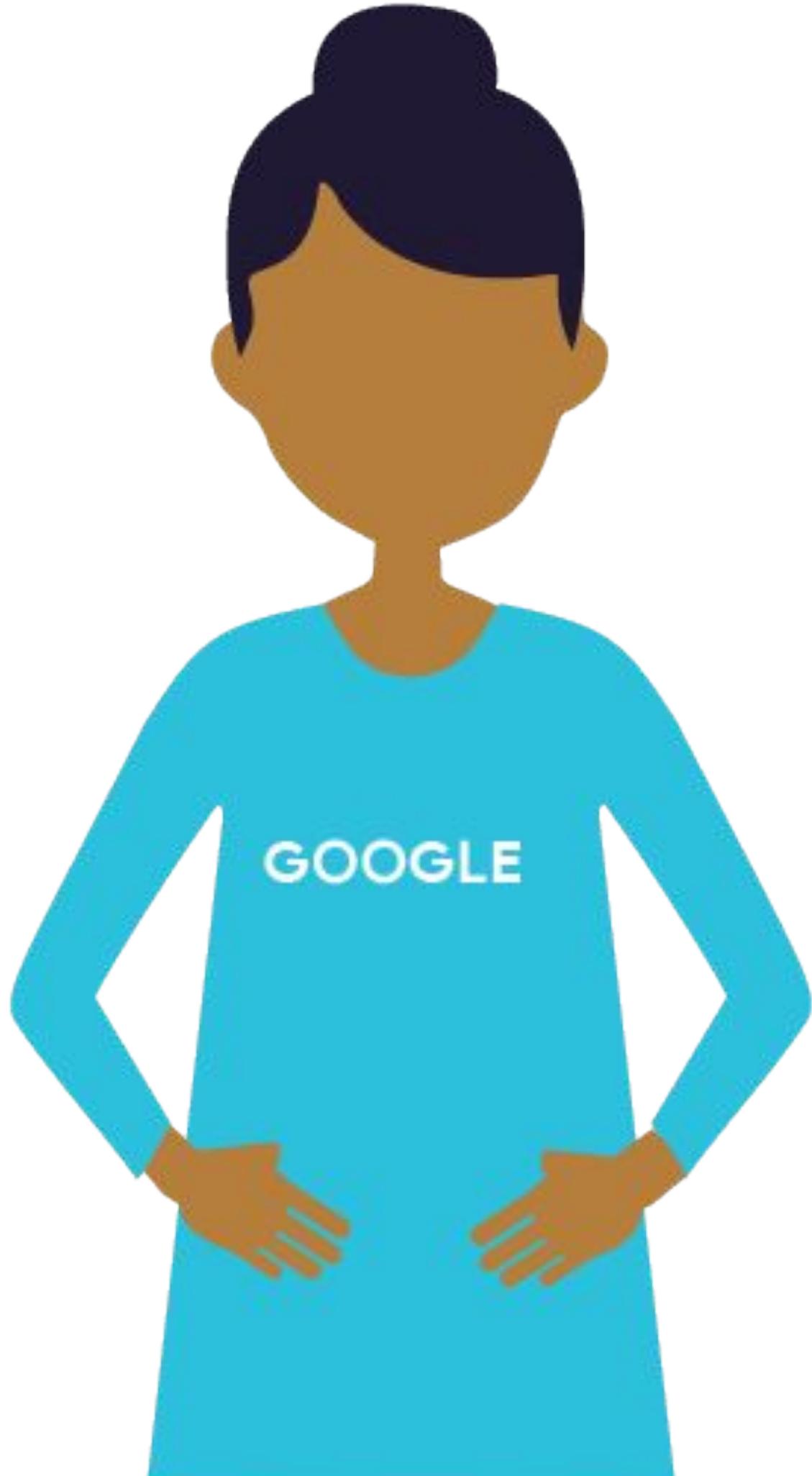
Format: Talking Head

Video Name: T-PSML-0_3_l5_adapting_to_data:_lab_intro

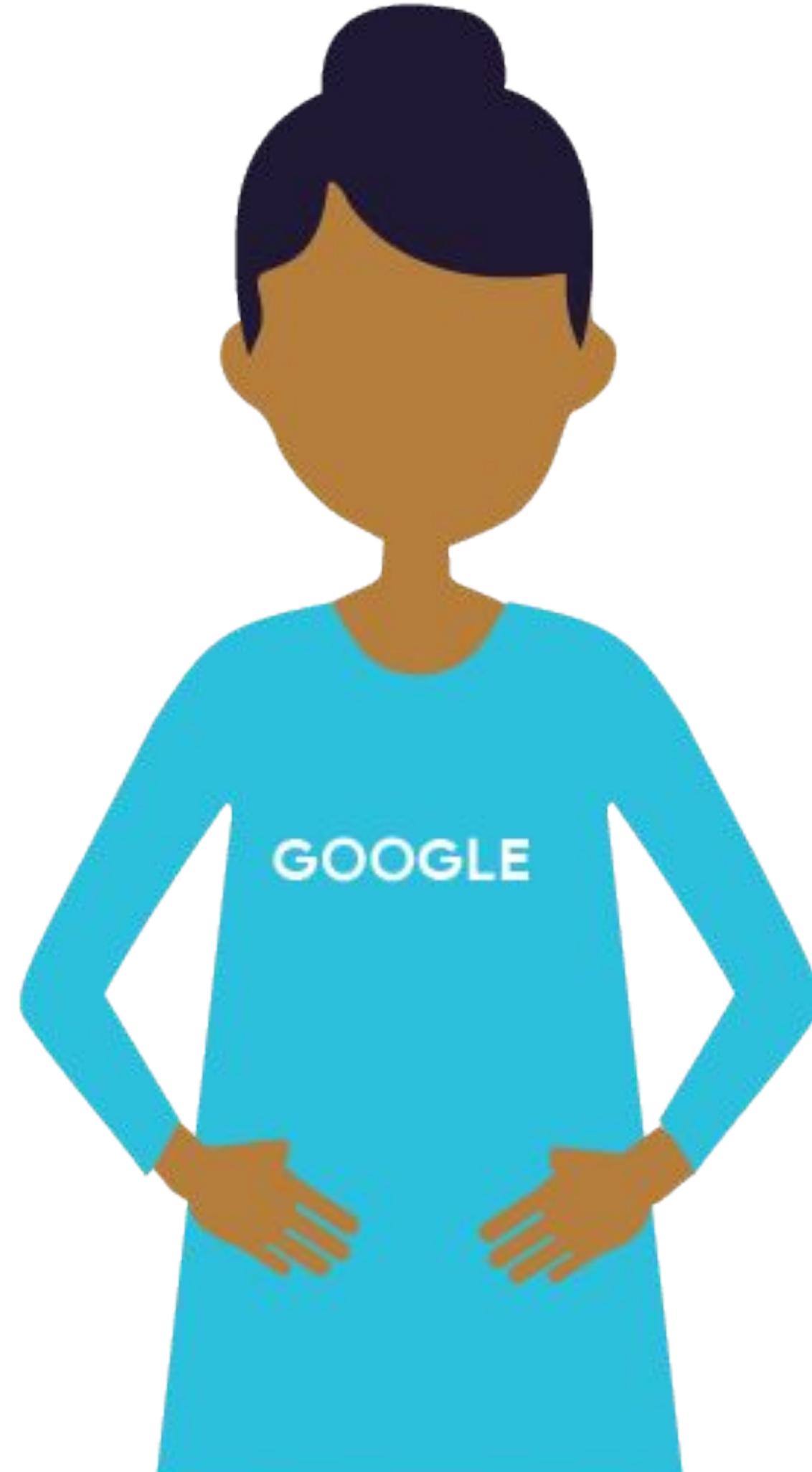
Lab

Making Good ML Engineering
Investments

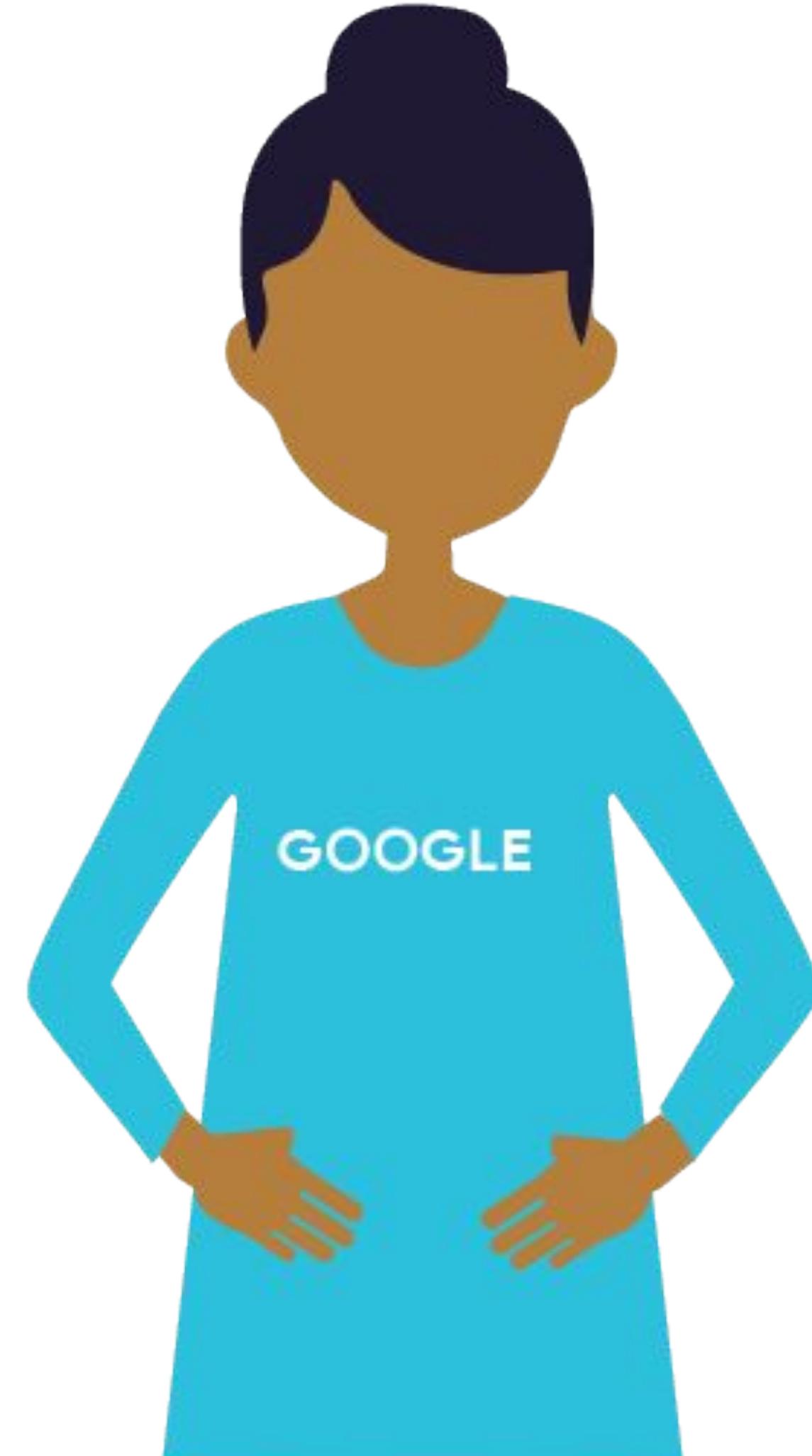
Max Lotstein



Scenario 1: Code Sprint



Scenario 2: A Gift Horse



Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

Lesson Title: Adapting to Data: Right and Wrong Decisions

Presenter: Max Lotstein

Format: Talking Head

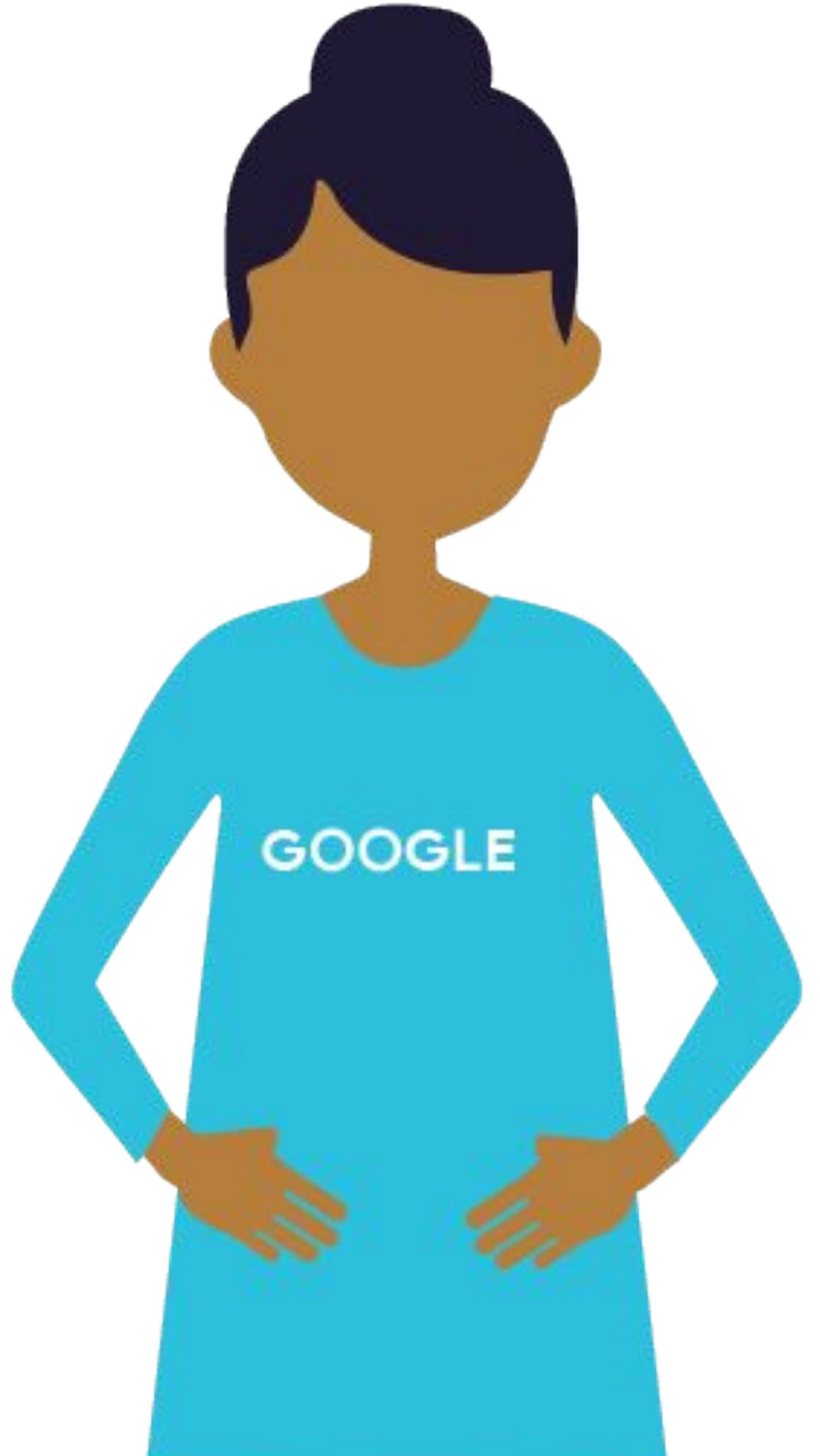
Video Name:

T-PSML-0_3_l7_adapting_to_data:_right_and_wrong_decisions



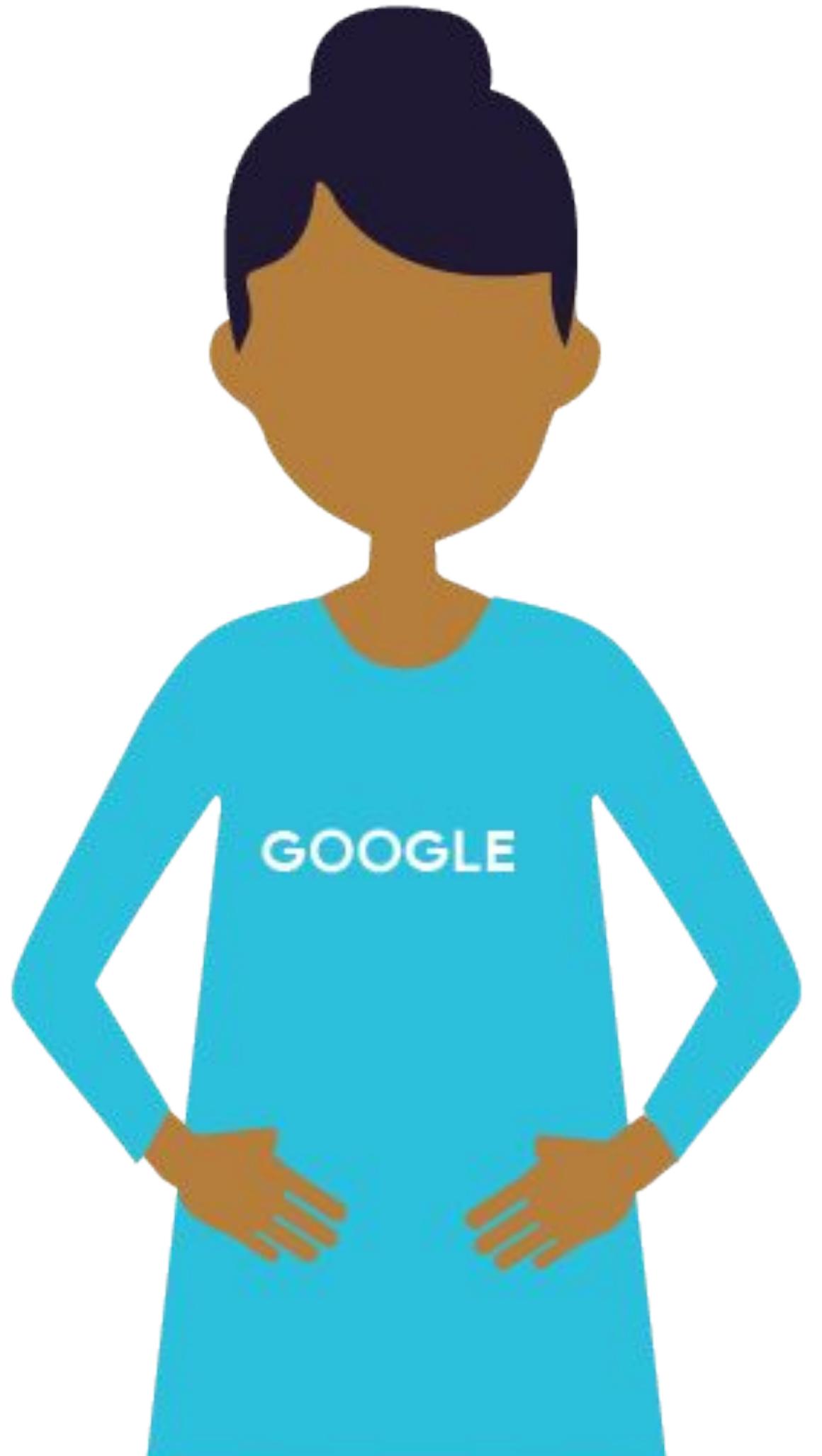
Right and Wrong Data Decisions





Right and Wrong Data Decisions

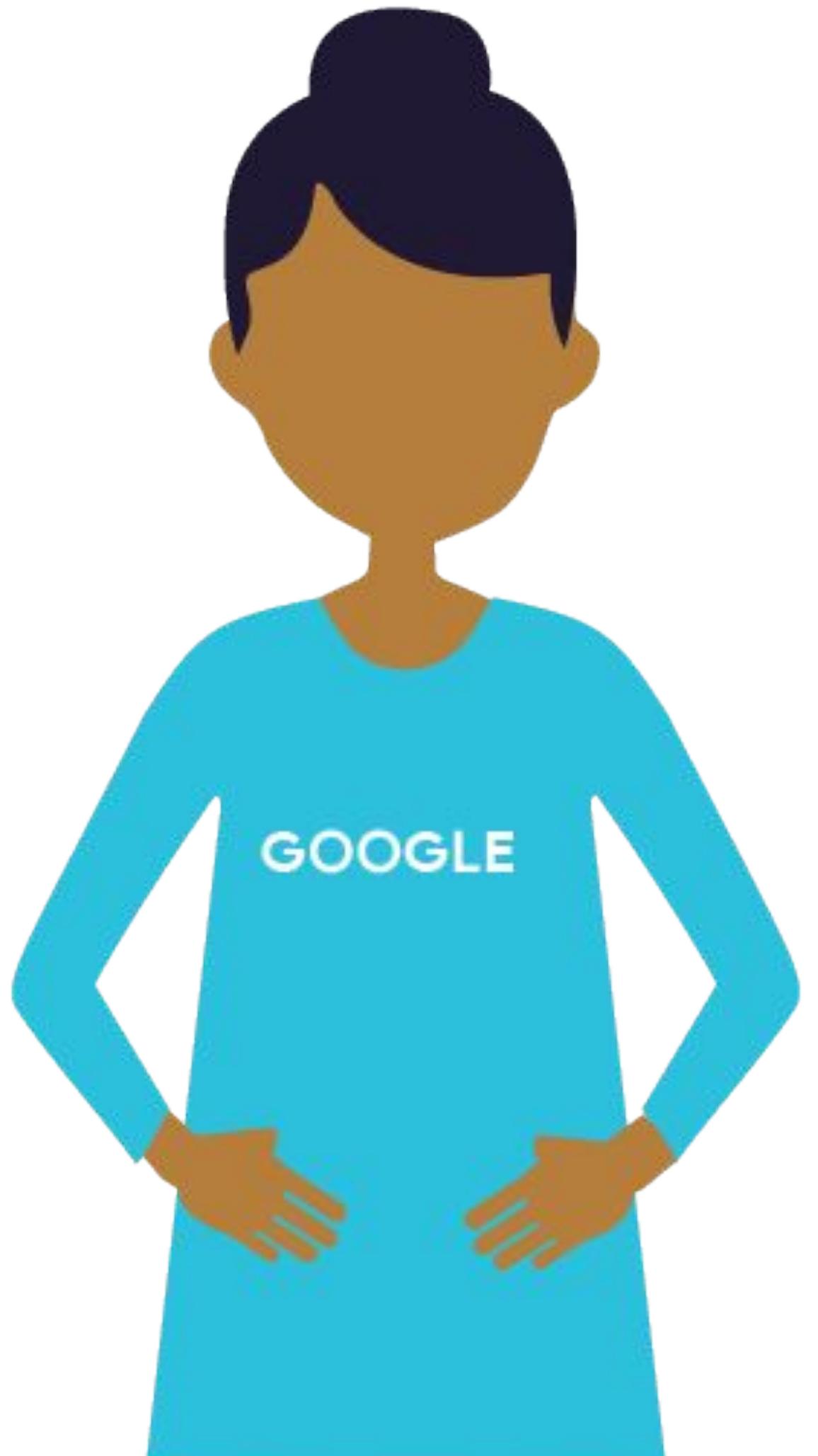
- patient age
- gender
- prior medical conditions
- hospital name
- vital signs
- test results



Data Leakage

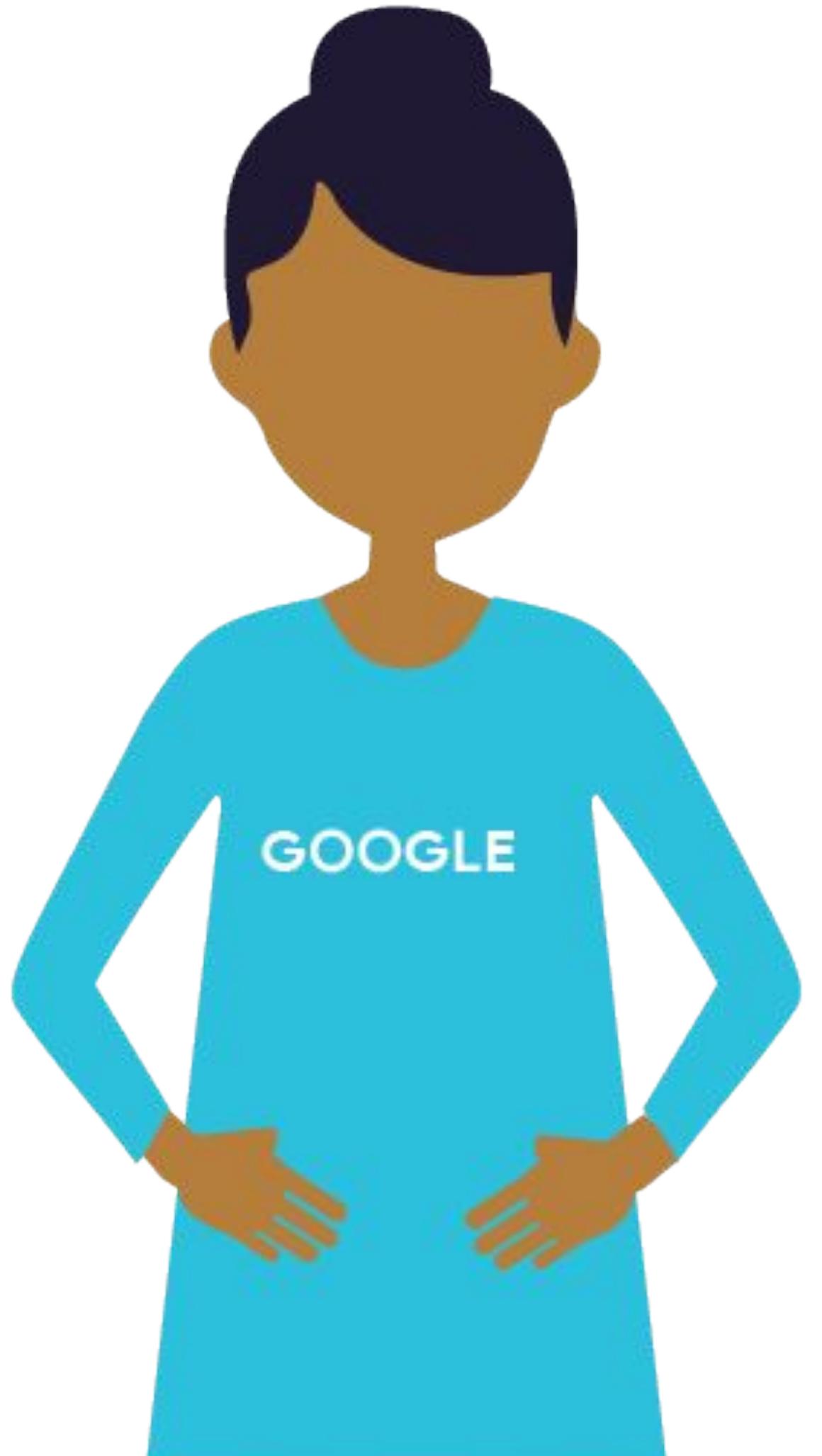


https://upload.wikimedia.org/wikipedia/commons/5/5f/Beth_Israel_Deaconess_Medical_Center_East_Campus.jpg



Predict political affiliation
from metaphors



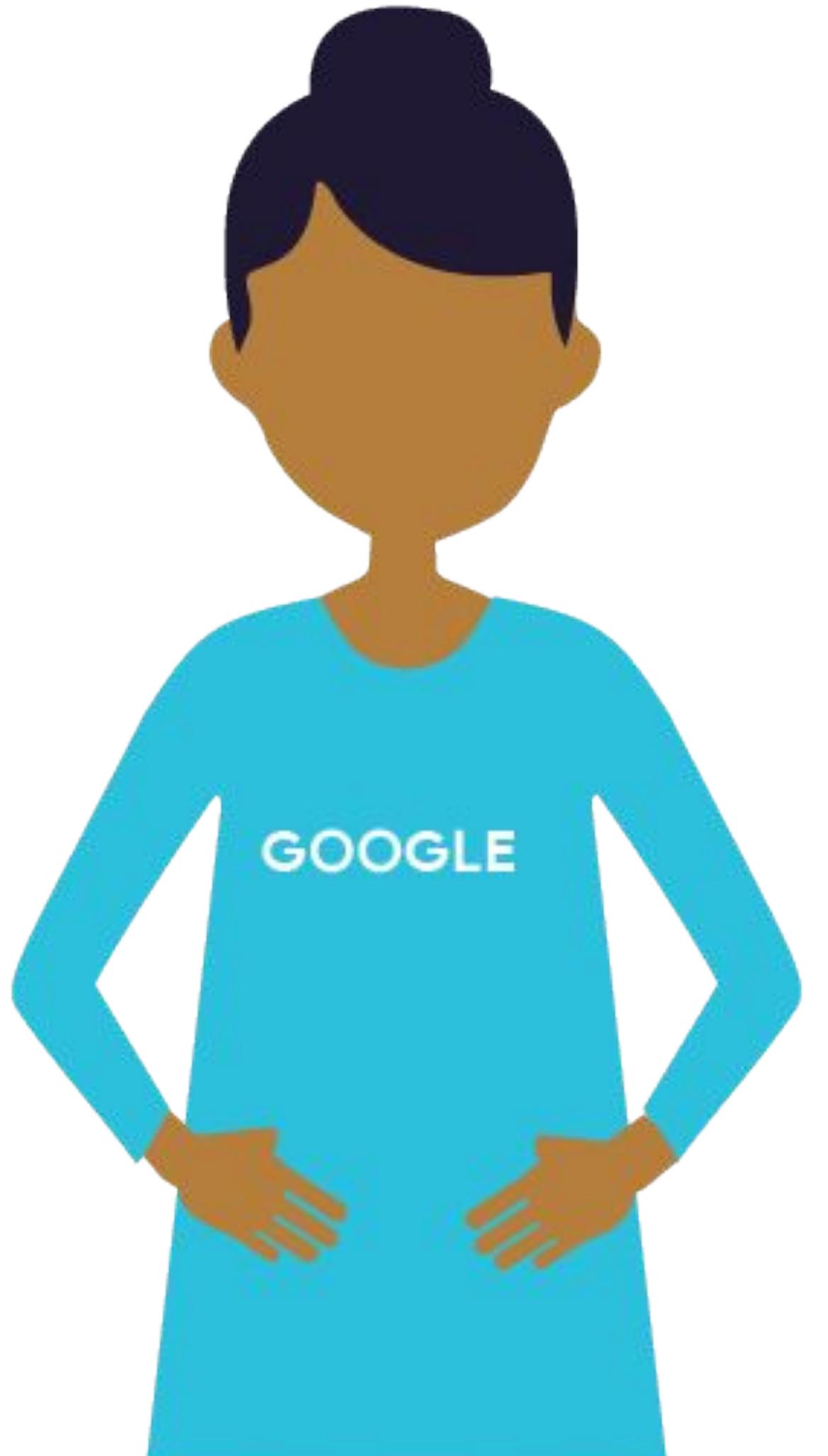


Predict political affiliation from metaphors

Google

the mind is a

- the mind is a **battlefield**
- the mind is a **walled garden**
- the mind is a **muscle**
- the mind is a **powerful tool**
- the mind is a **powerful force**
- the mind is a **powerful**
- the mind is a **prison**
- the mind is a **great servant**
- the mind is a **soft boiled potato**
- the mind is a **beautiful servant**



Predict political affiliation from metaphors

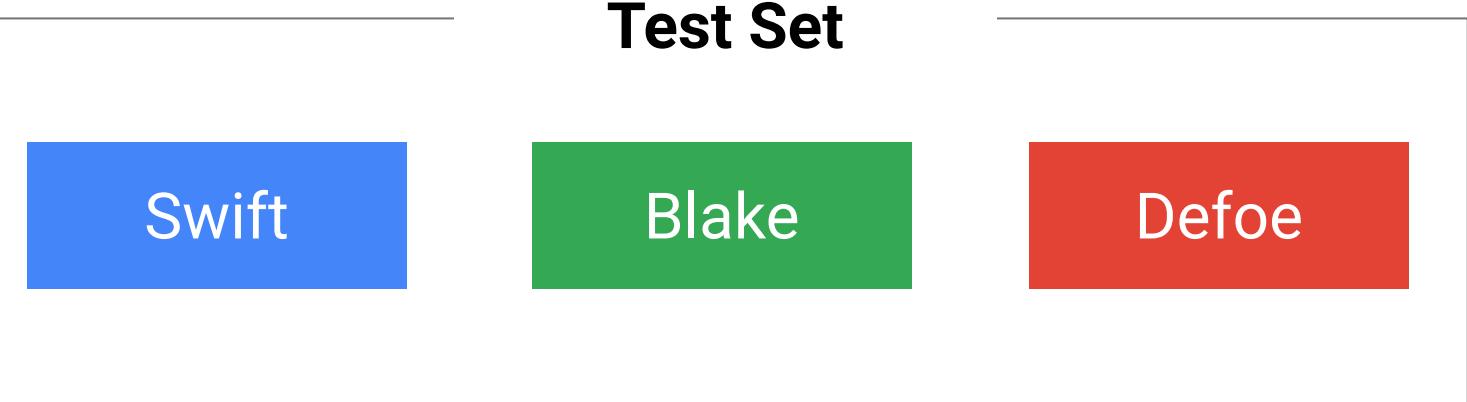
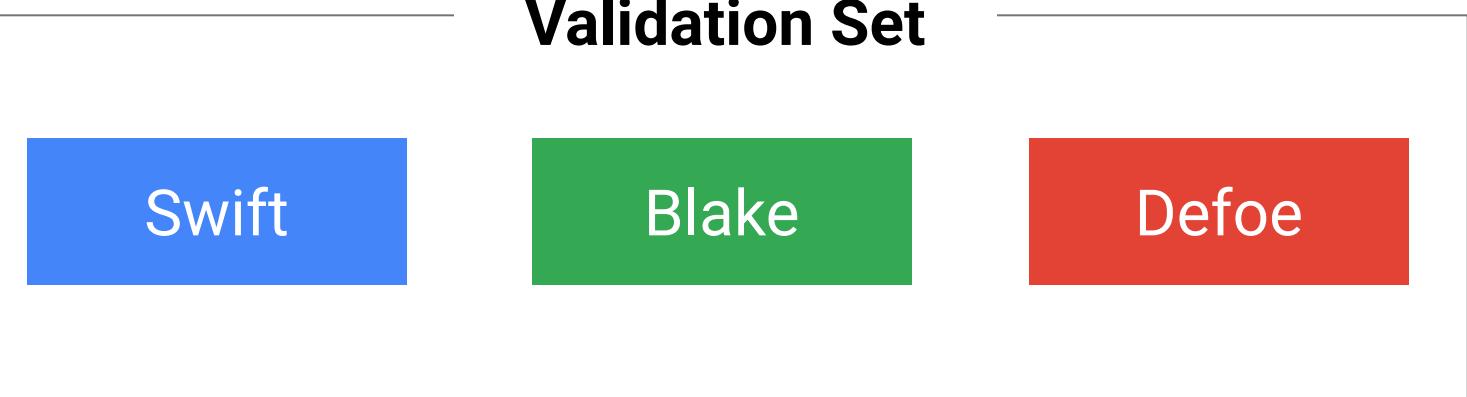
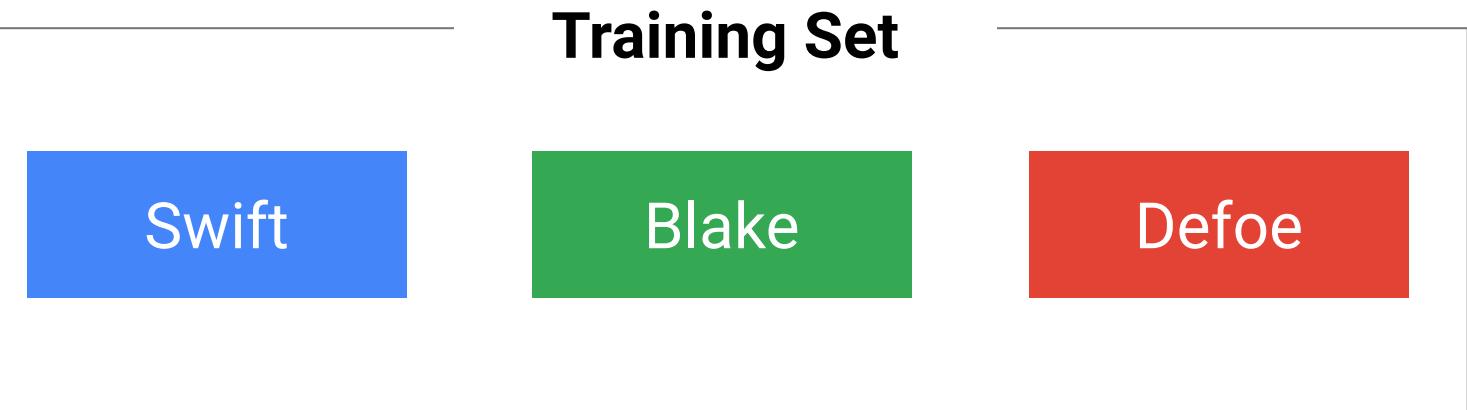
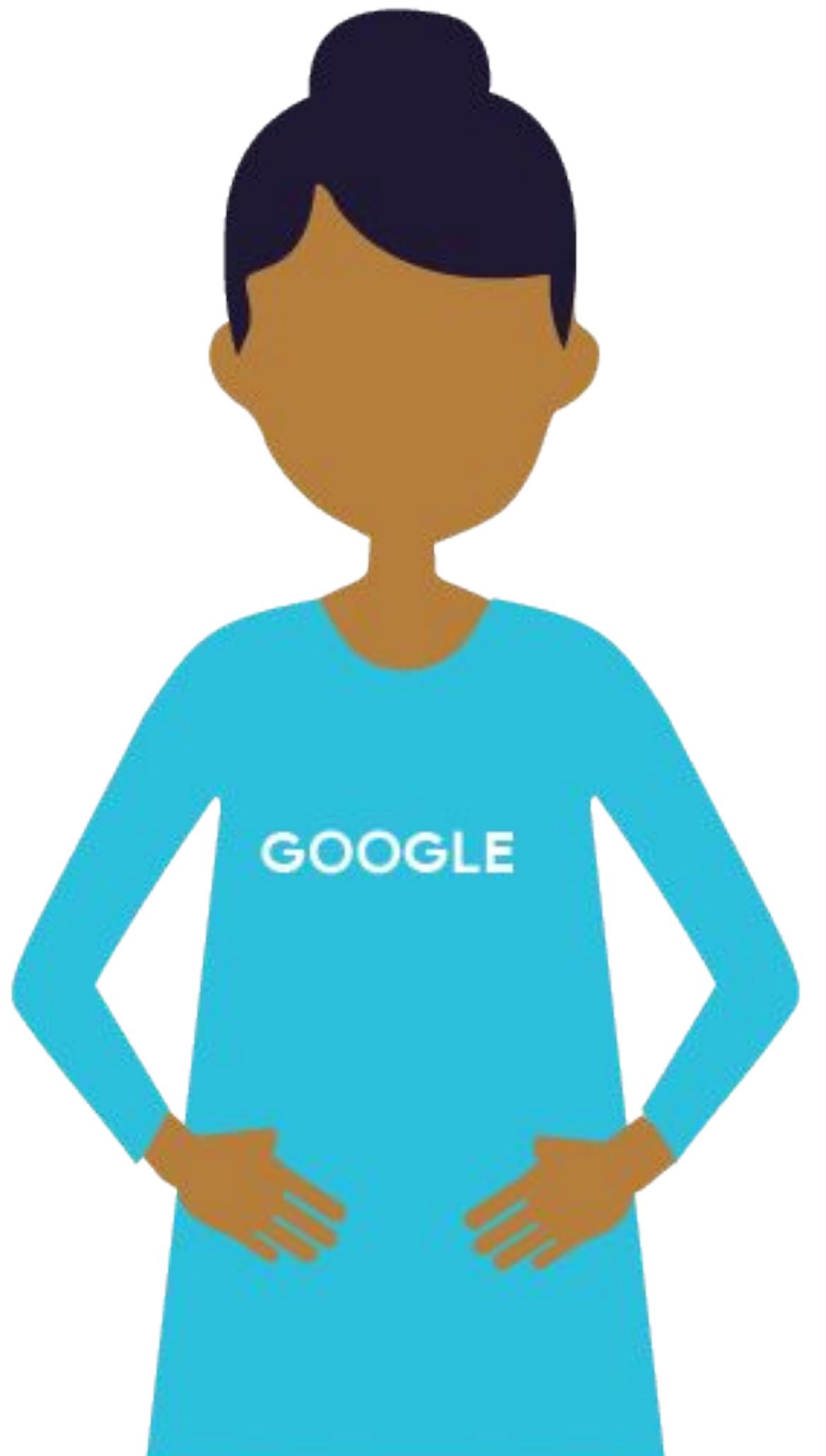
Google

the mind is a

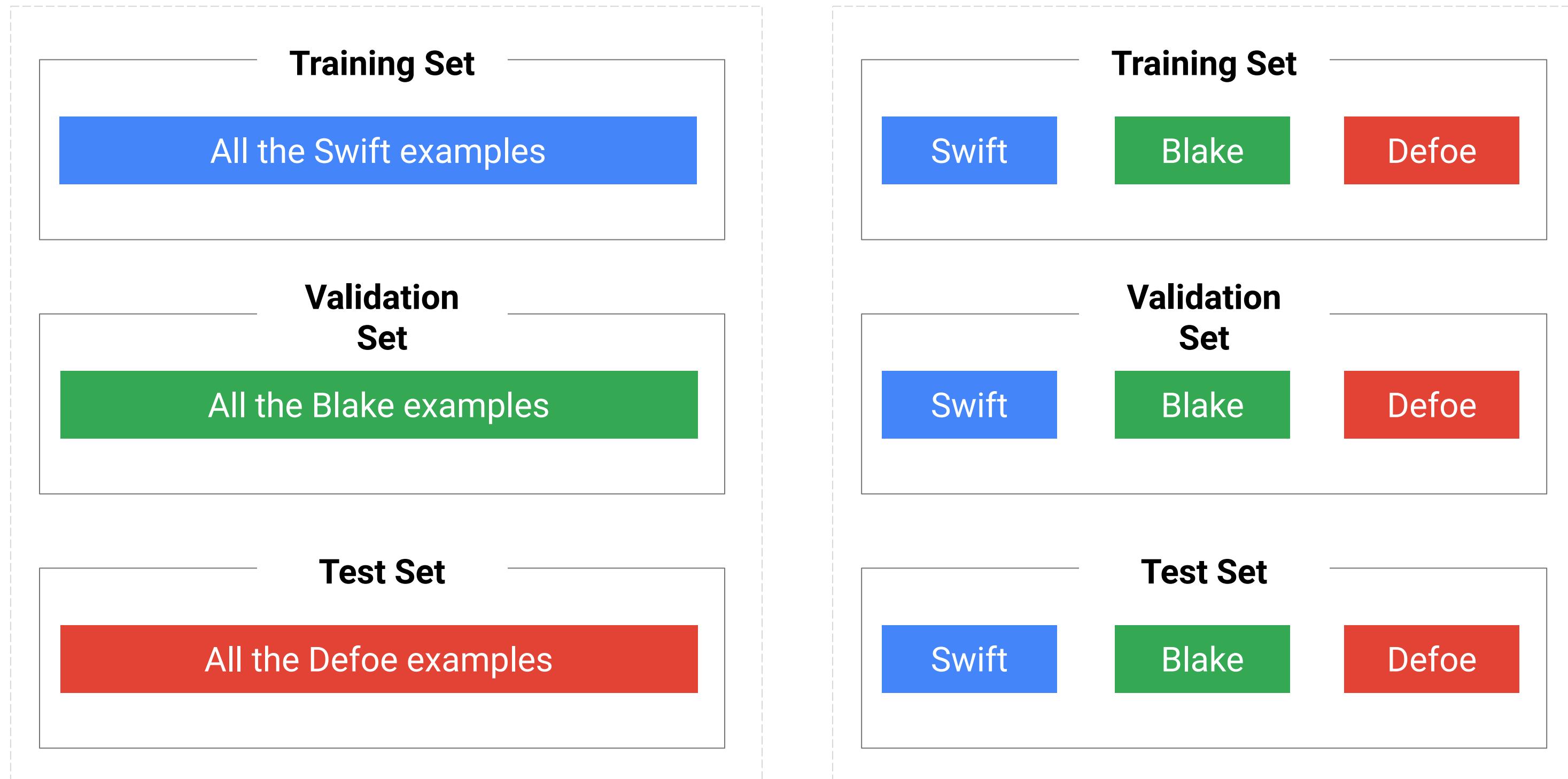
the mind is a **battlefield**
the mind is a **walled garden**
the mind is a **muscle**

the mind is a **powerful tool**
the mind is a **powerful force**
the mind is a **powerful**

the mind is a **prison**
the mind is a **great servant**
the mind is a **soft boiled potato**
the mind is a **beautiful servant**



Solution: Cross-contamination; you have to split by author



Course 2: Production ML Systems

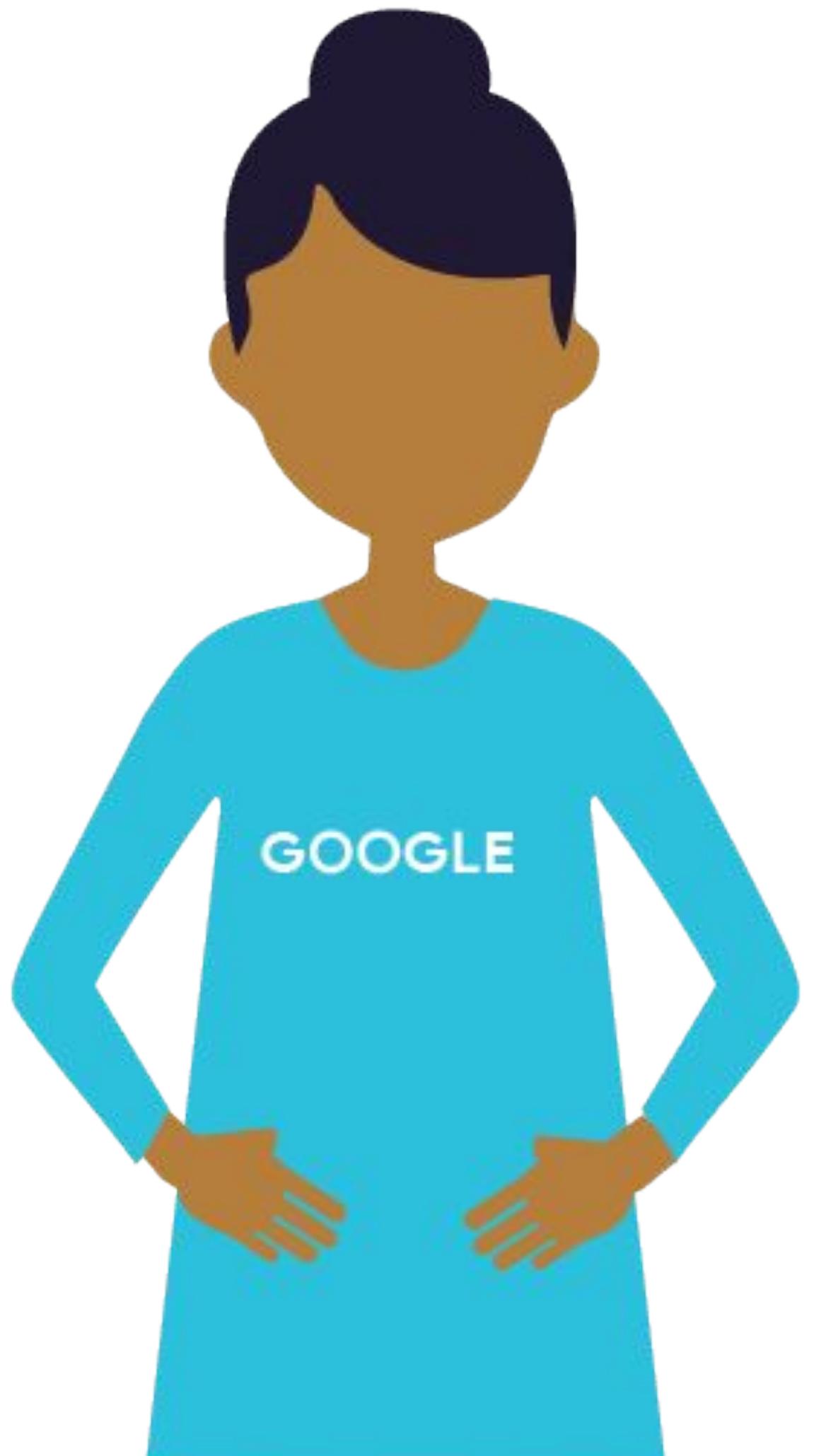
Module 3: Designing Adaptable ML Systems

Lesson Title: **Adapting to Data: System Failure**

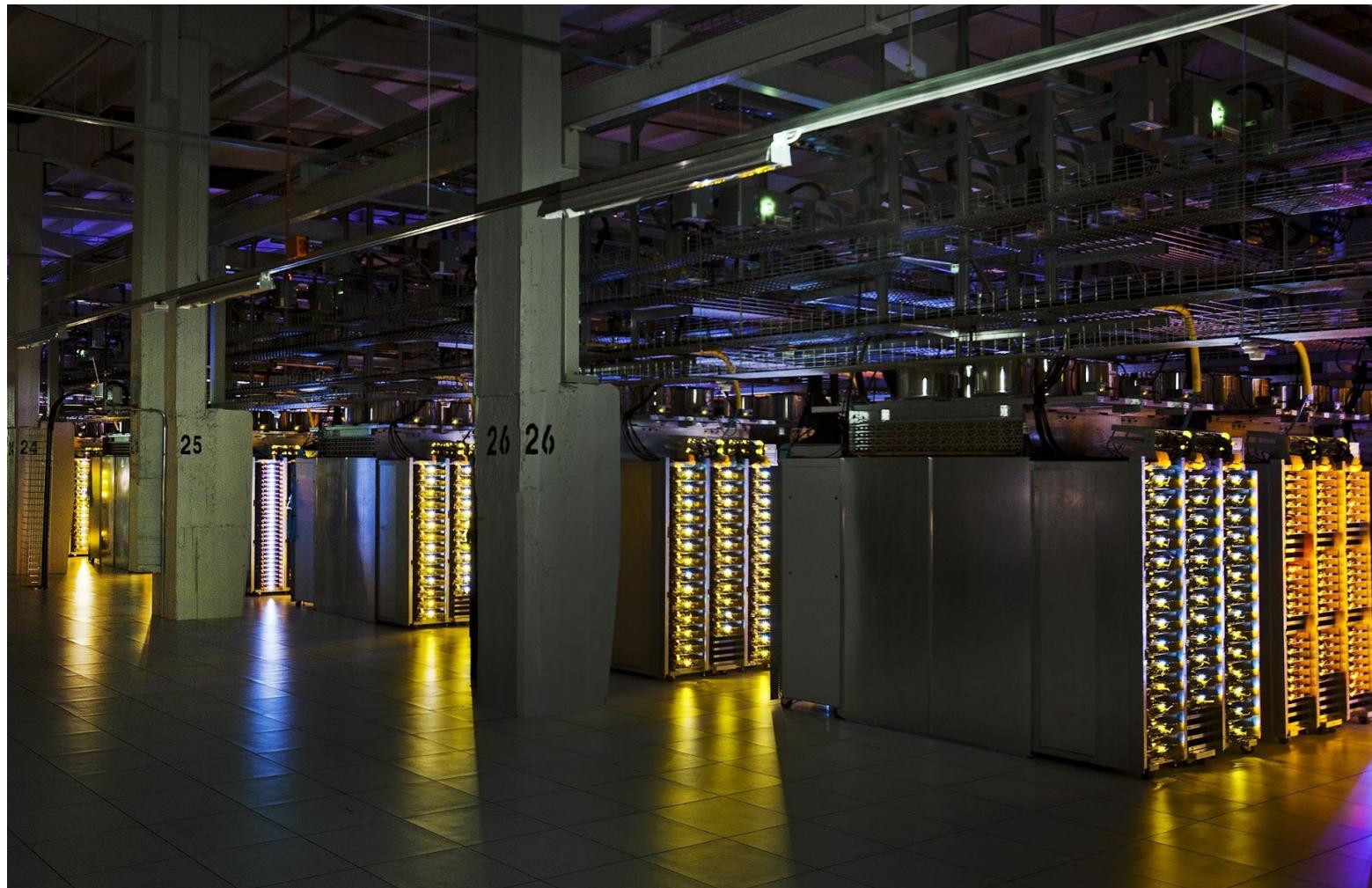
Presenter: Max Lotstein

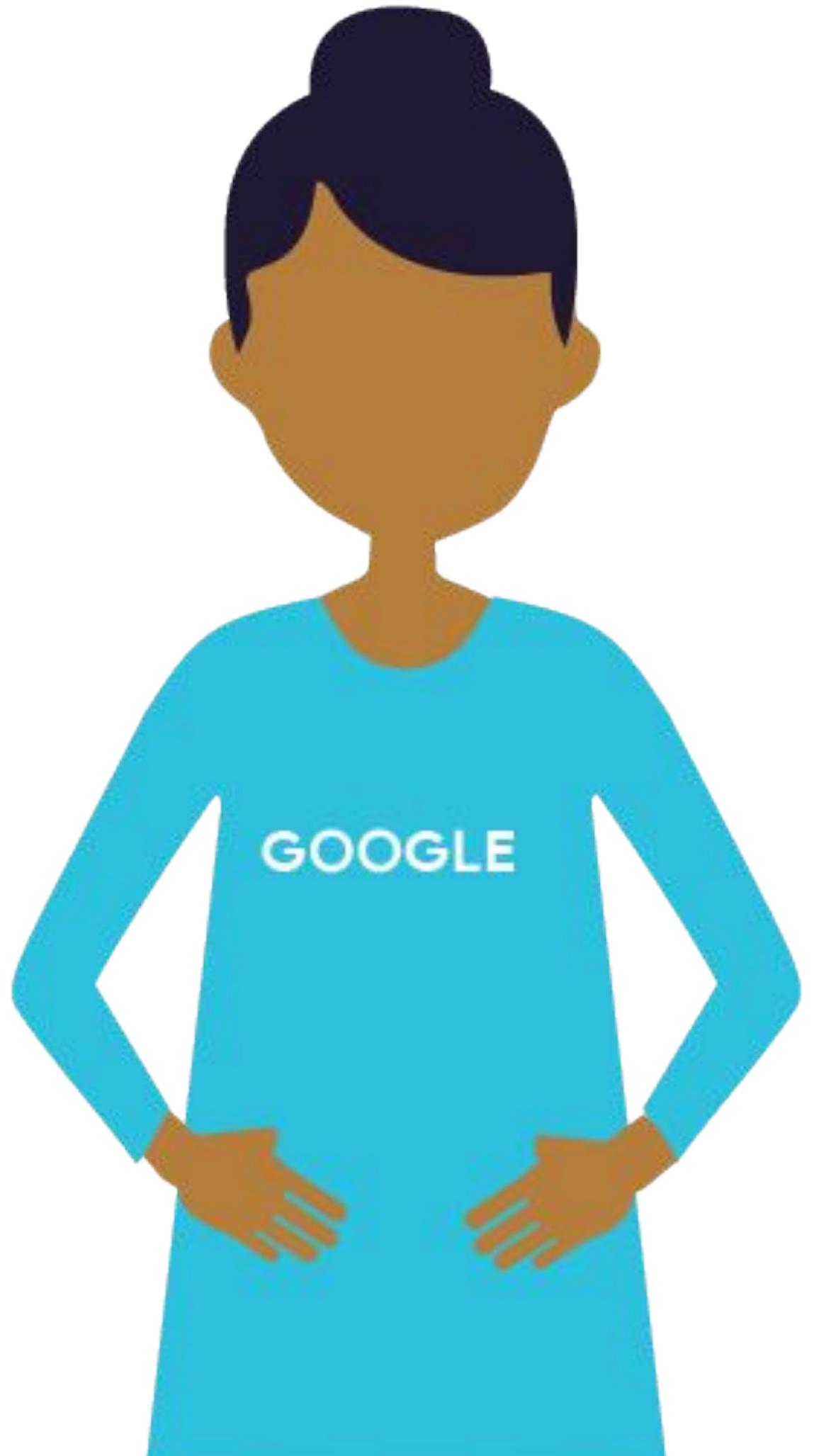
Format: Talking Head

Video Name: T-PSML-0_3_l8_adapting_to_data:_system_failure

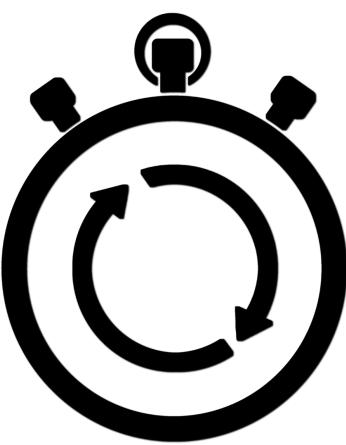
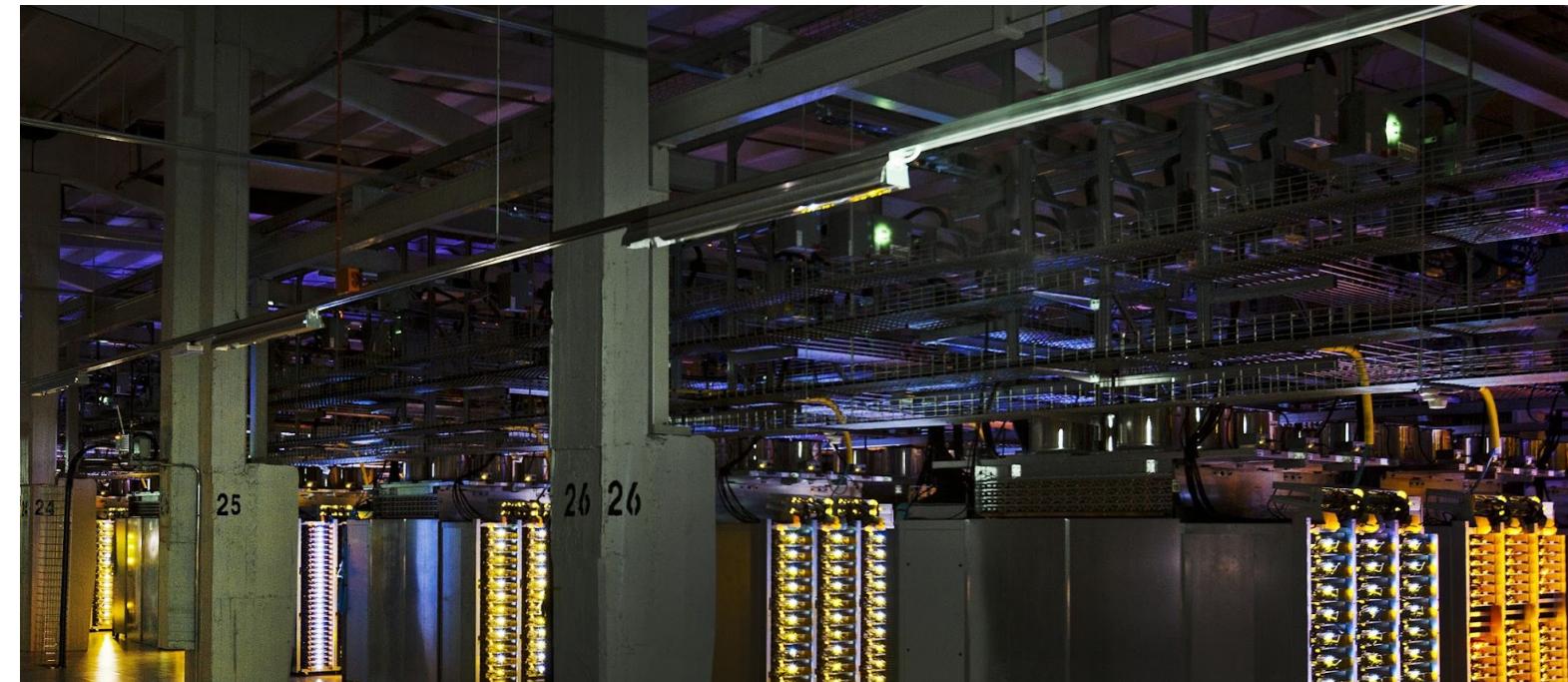


Systems Fail

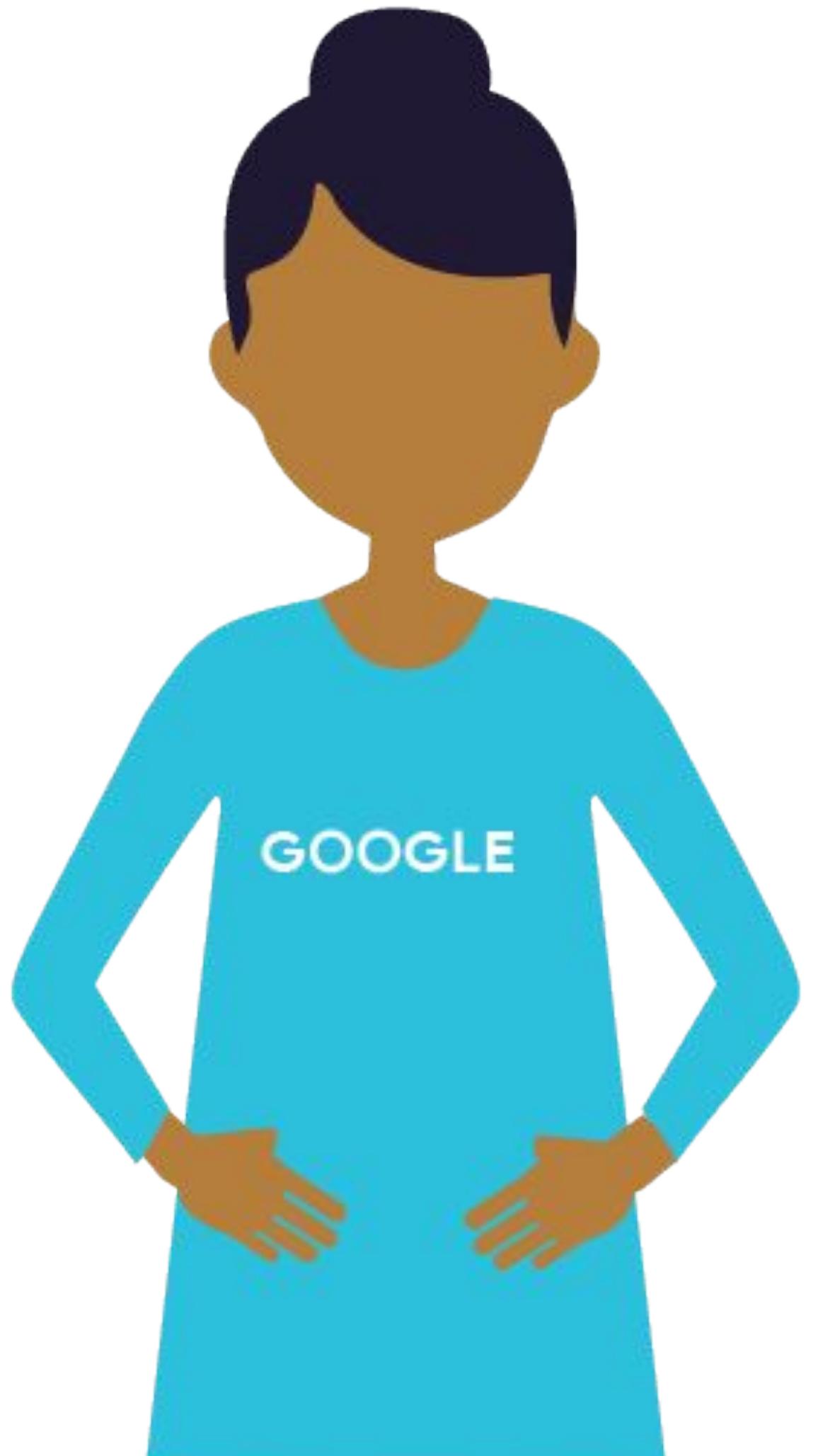




Systems Fail

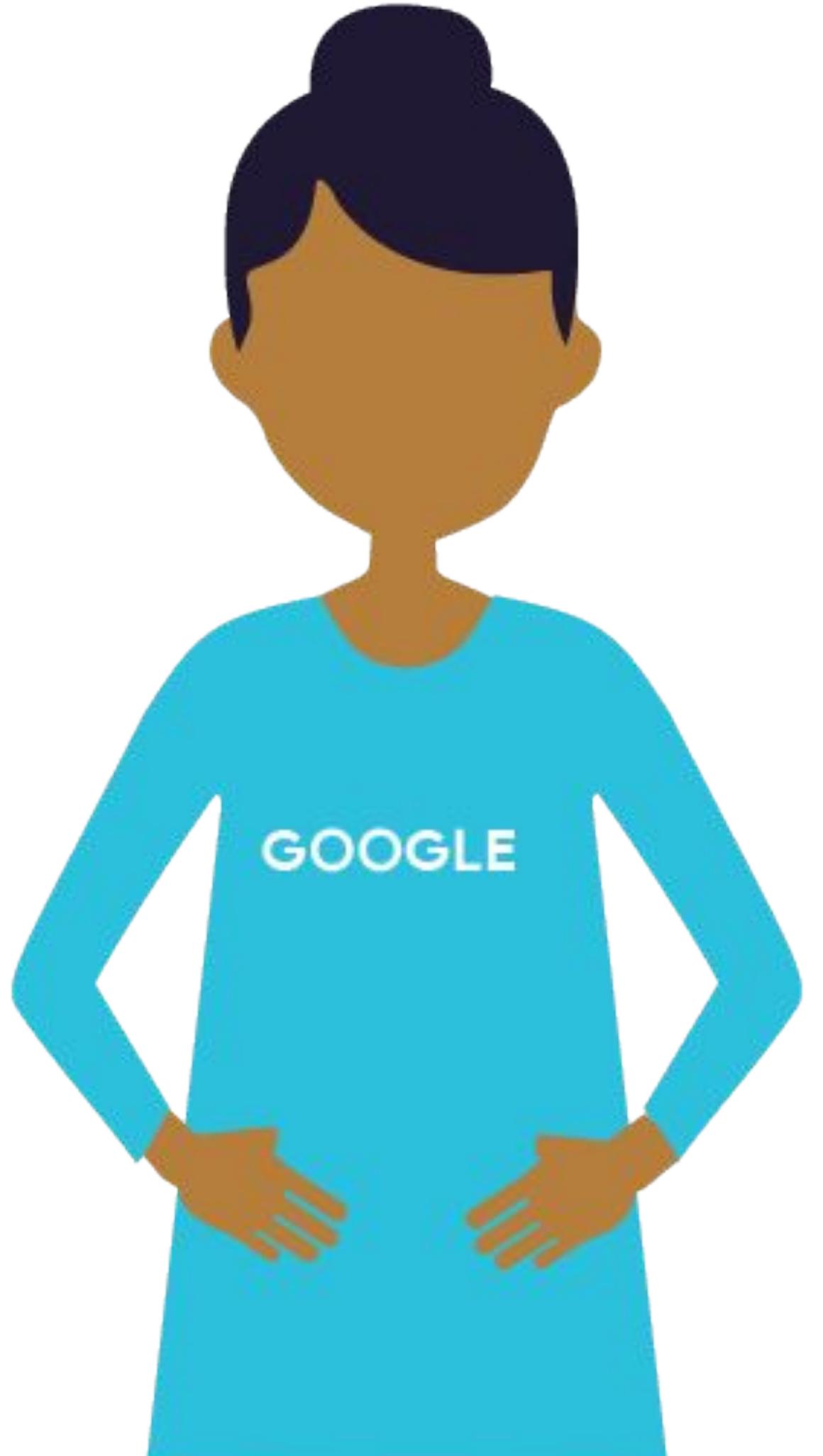


Rollback Initiated
Version 1.0.1
Three Months old

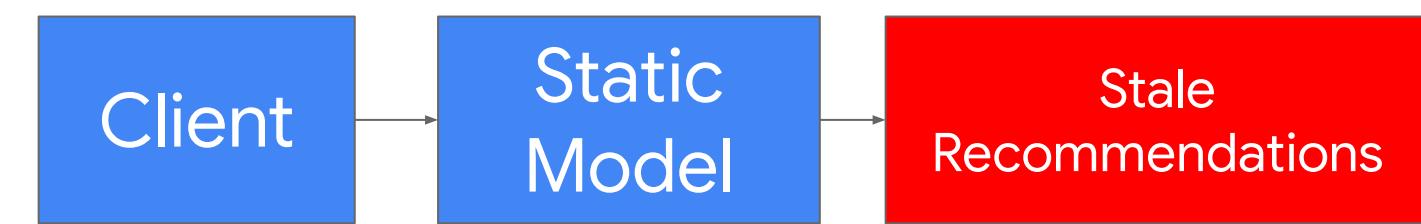


Feedback Loops



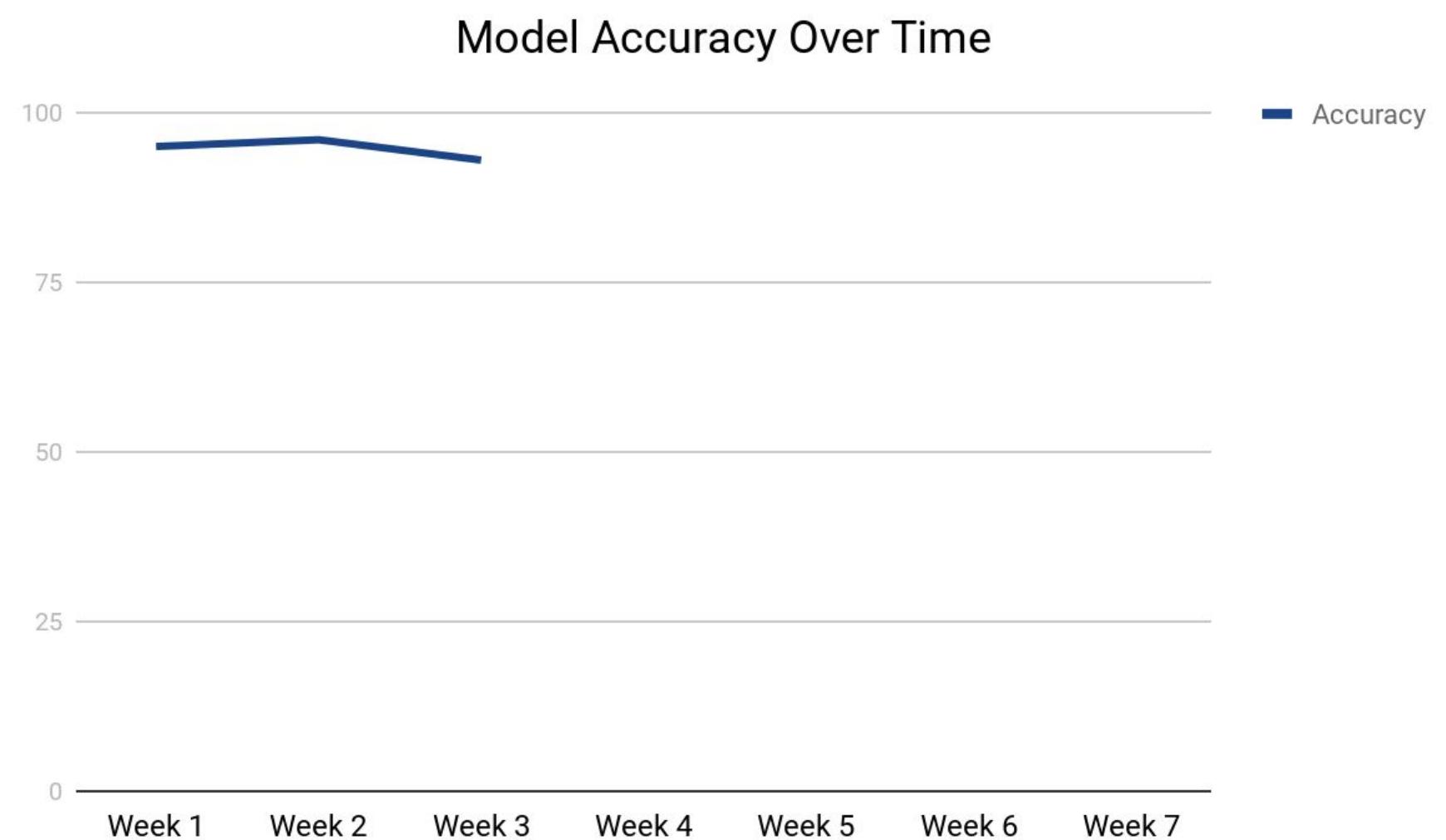


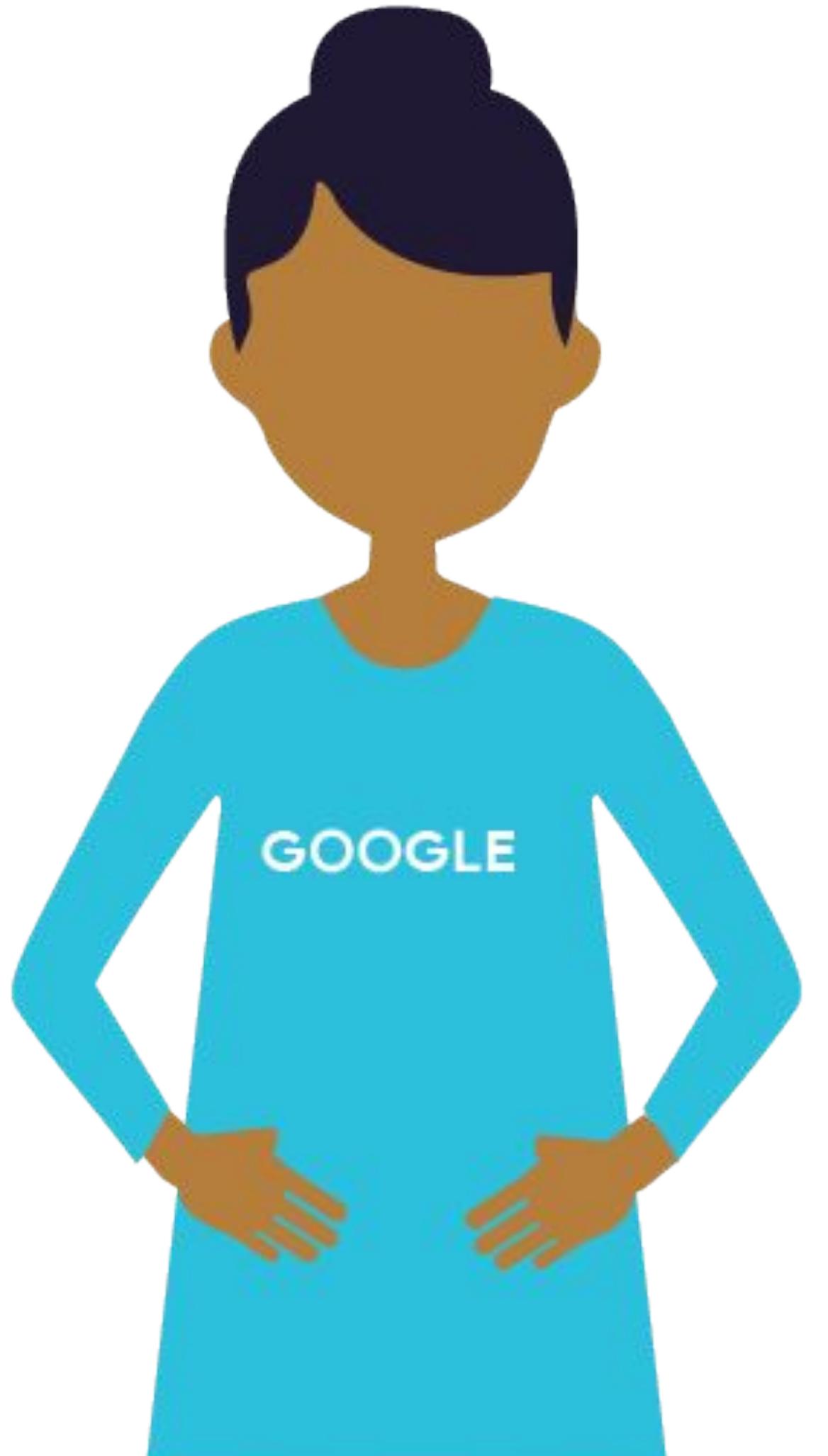
Feedback Loops



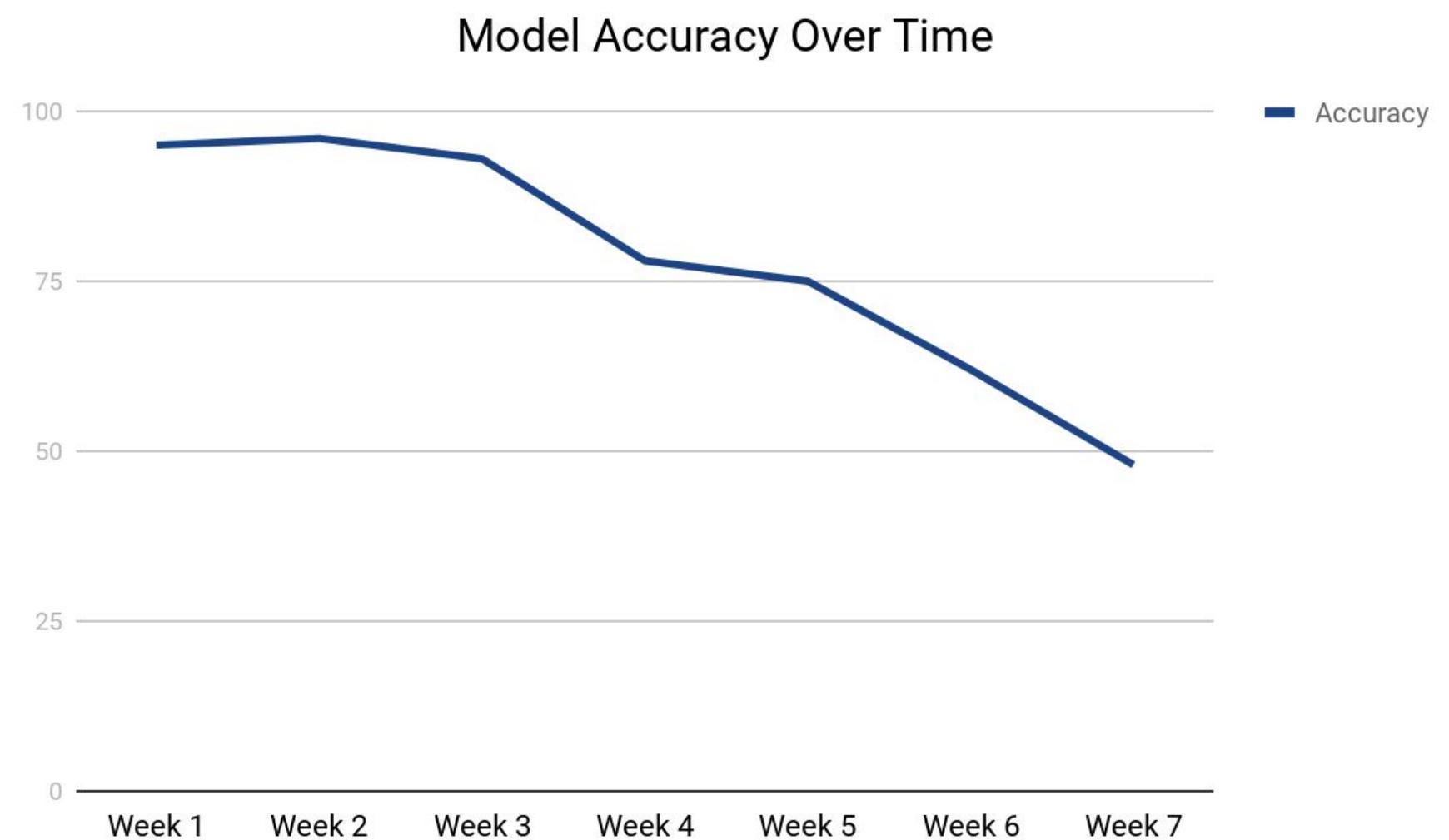


Feedback Loops





Feedback Loops



Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

Lesson Title: **Adapting to Data: Summary**

Presenter: Max Lotstein

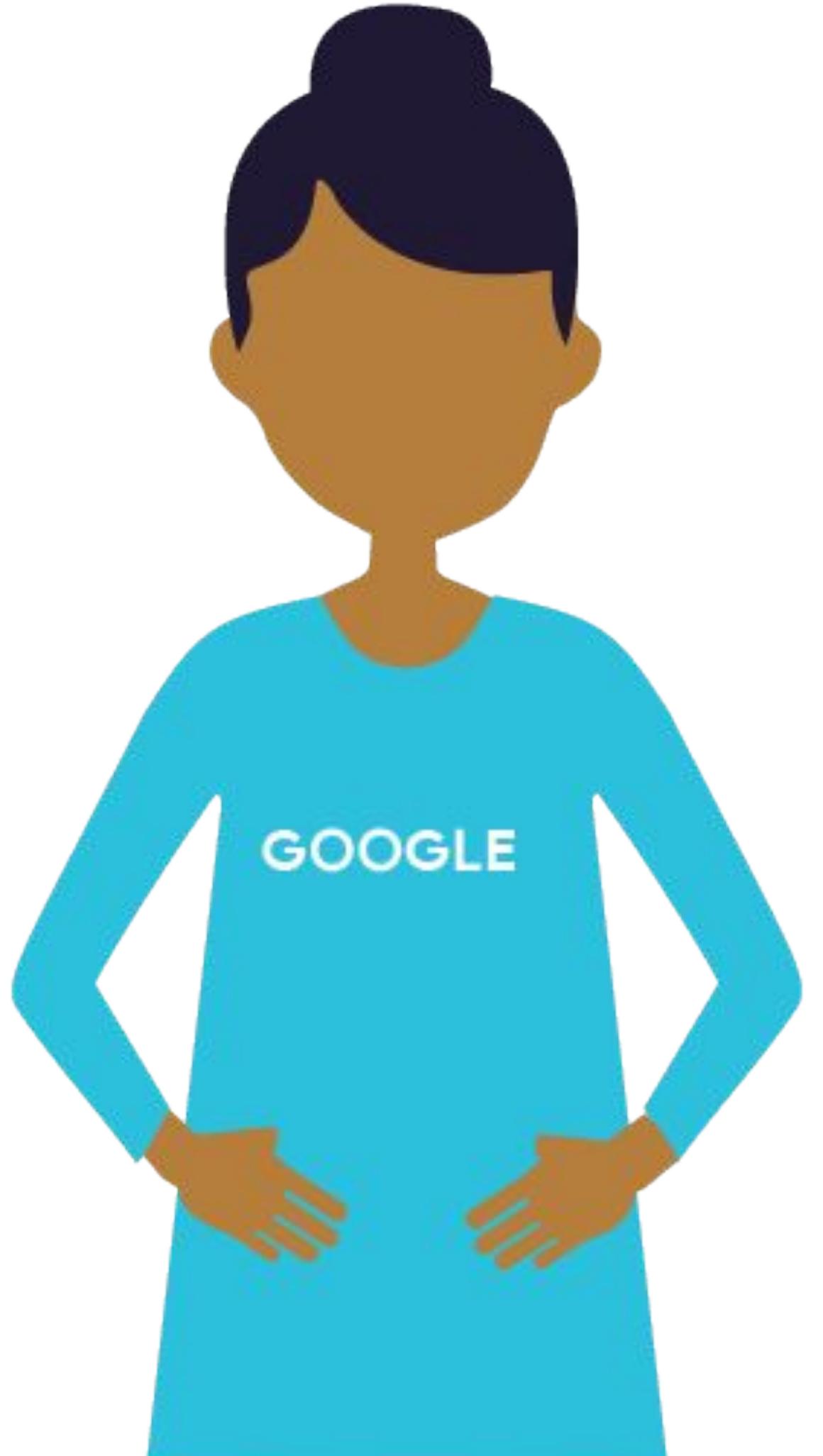
Format: Talking Head

Video Name: T-PSML-0_3_l9_adapting_to_data:_summary



Adapting to Data

- Assess all data sources and features based on both cost and benefit before including into the model
- Communicate with upstream data producers to make your needs known
- Replicate critical data sources
- Monitor descriptive statistics for your inputs and outputs



Adapting to Data

- Monitor your residuals as a function of your inputs
- Use custom weights in your loss function to emphasize data recency
- Use dynamic training architecture and regularly retrain your model
- You get what you optimize for

Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

Lesson Title: **Mitigating Training-Serving Skew Through Design**

Presenter: Max Lotstein

Format: Talking Head

Video Name:

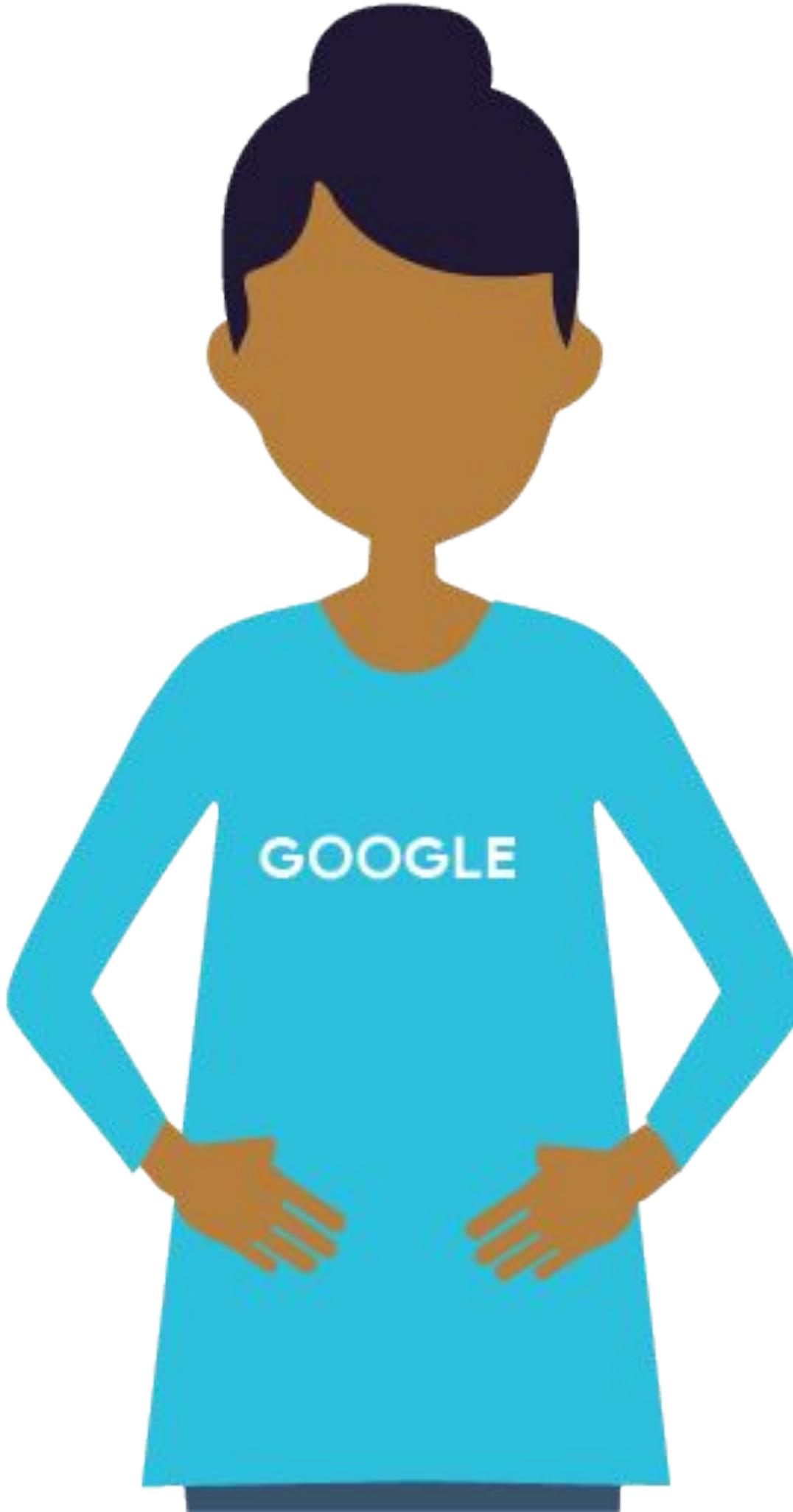
T-PSML-0_3_l10_mitigating_training-serving_skew_through_design

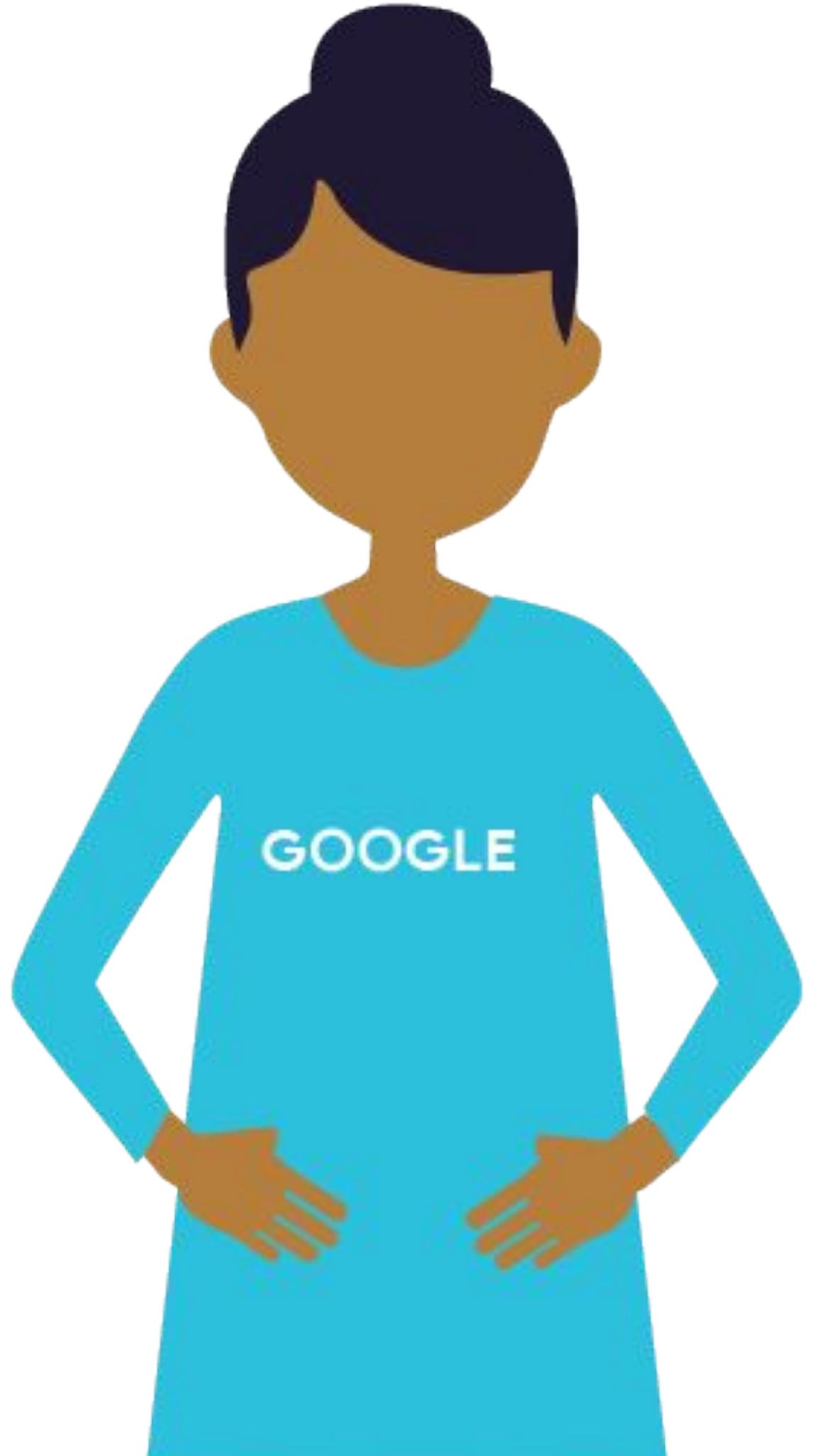
Agenda

Adapting to Data

**Mitigating Training-Serving
Skew Through Design**

Debugging a Production Model



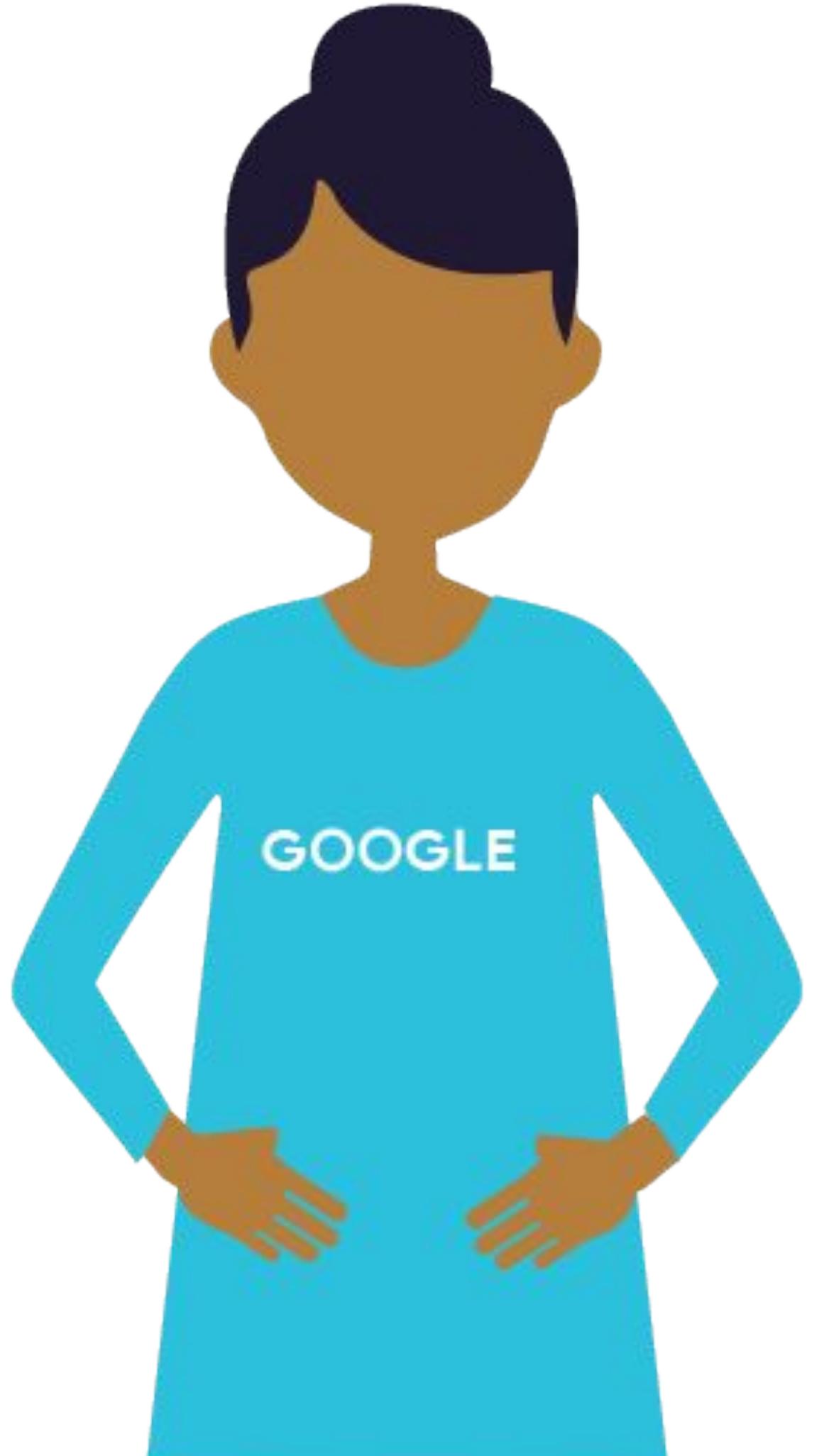


Agenda

Adapting to Data

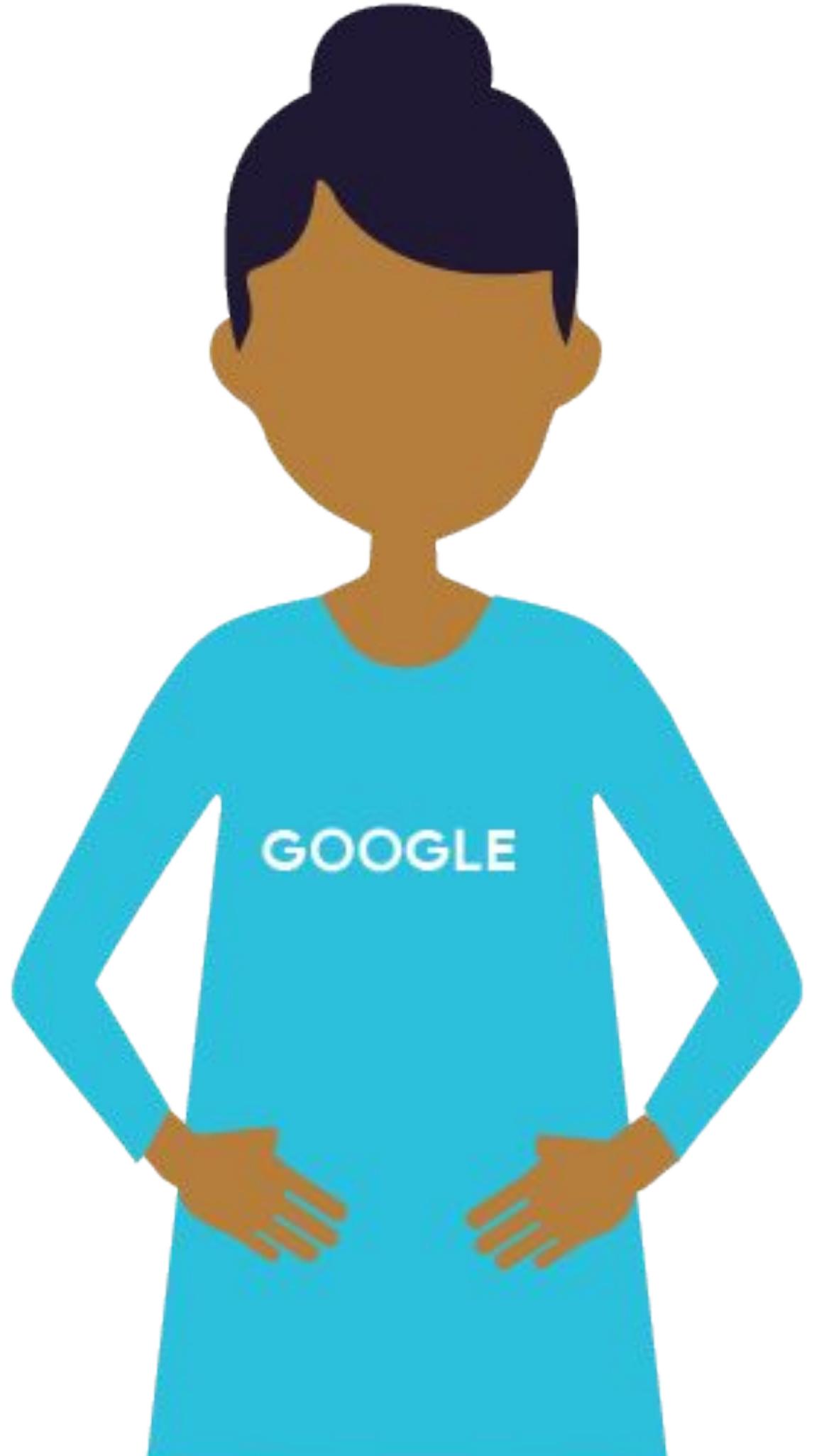
**Mitigating Training-Serving
Skew Through Design**

Debugging a Production Model



Training/Serving Skew

1. A discrepancy between how you handle data in the training and serving pipelines
2. A change in the data between when you train and when you serve.
3. A feedback loop between your model and your algorithm.



How Code Can Create Training/Serving Skew

- Different library versions that are functionally equivalent but optimized differently
- Different library versions that are not functionally equivalent
- Re-implemented functions

Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

Lesson Title: **Lab Intro: Serving ML Predictions in batch and real-time**

Presenter: Max Lotstein

Format: Talking Head

Video Name:

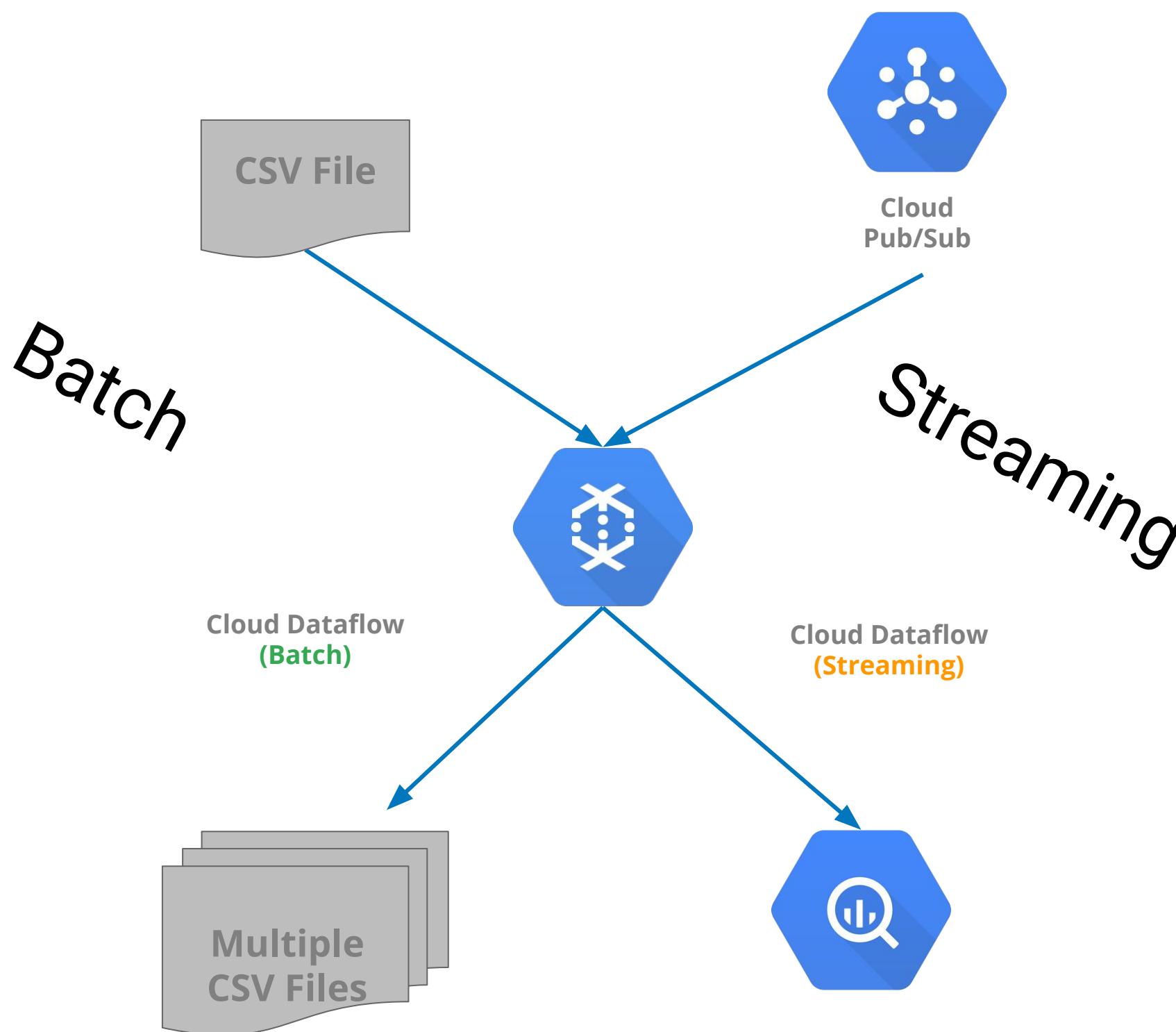
T-PSML-0_3_l11_lab_intro:_serving_ml_predictions_in_batch_and_real-time

Lab

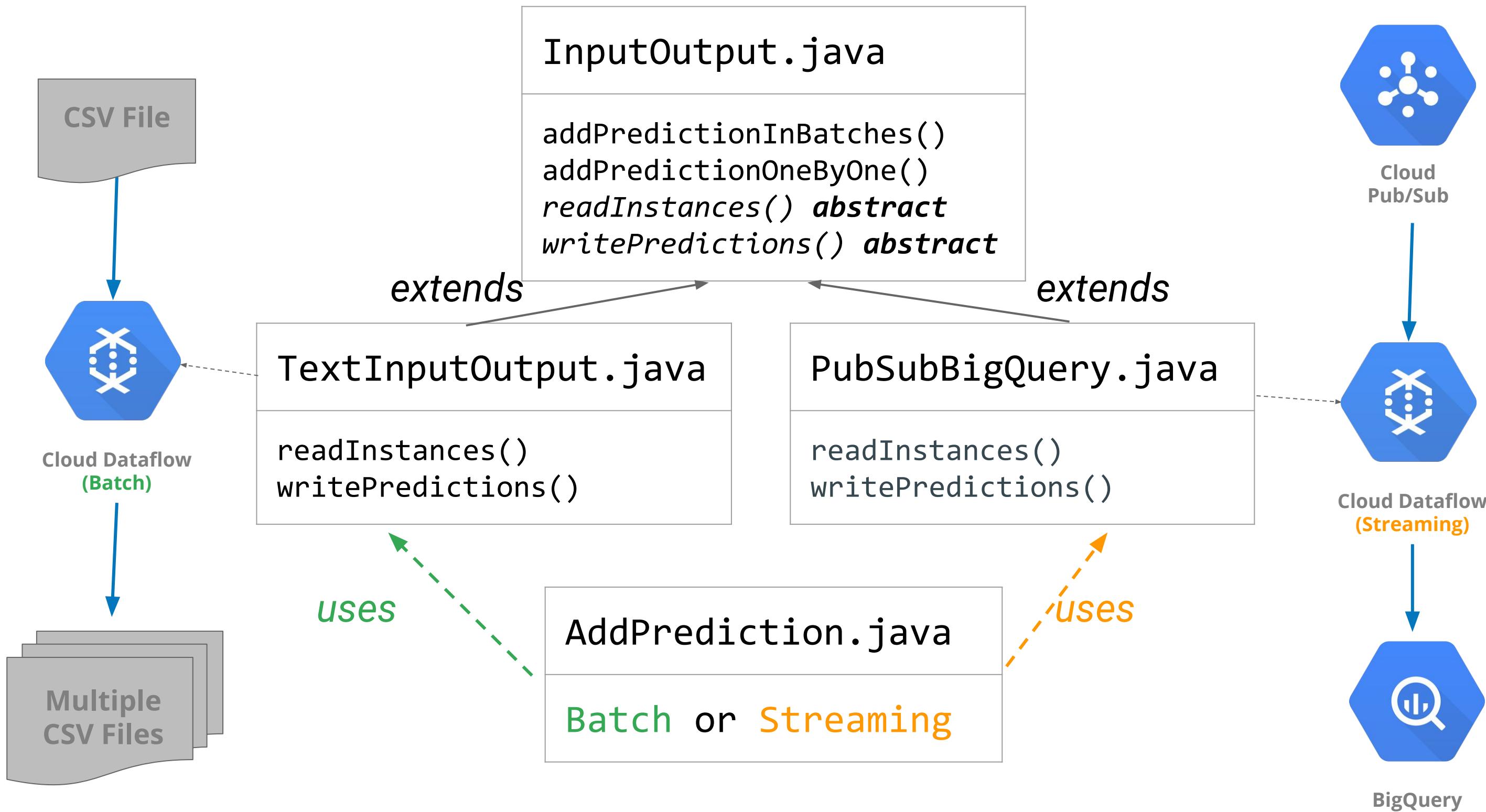
Serving ML Predictions in
batch and real-time

Max Lotstein





Lab: Serving ML Predictions in batch and real-time



Course 2: Production ML Systems

Module 3: Designing Adaptable ML Systems

Lesson Title: **Lab Solution: Serving ML Predictions in batch and real-time**

Presenter: Max Lotstein

Format: Talking Head

Video Name:

T-PSML-0_3_l12_lab_solution:_serving_ml_predictions_in_batch_and_real-time



Title Safe >

< Action Safe

Course 2: Production ML Systems

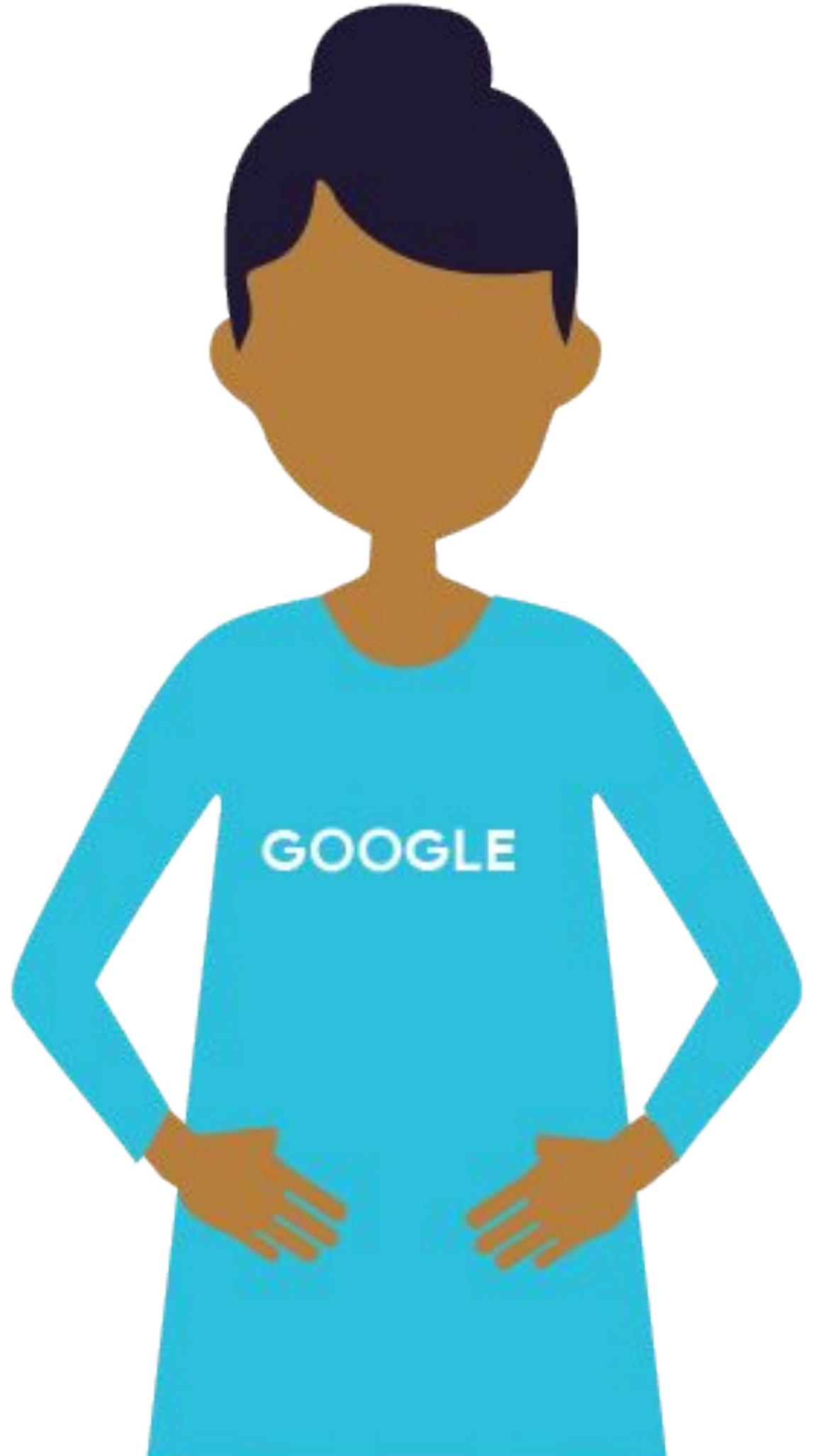
Module 3: Designing Adaptable ML Systems

Lesson Title: Debugging a Production Model

Presenter: Max Lotstein

Format: Talking Head

Video Name: T-PSML-0_3_l13_debugging_a_production_model



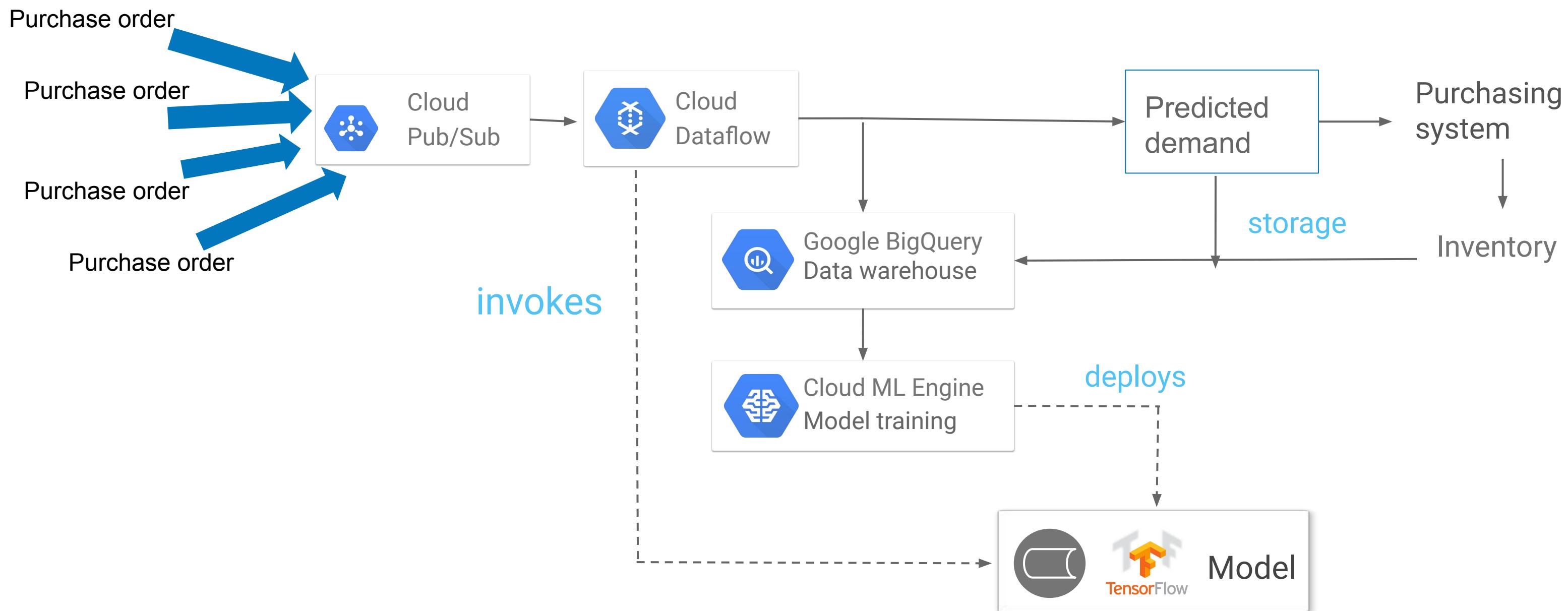
Agenda

Adapting to Data

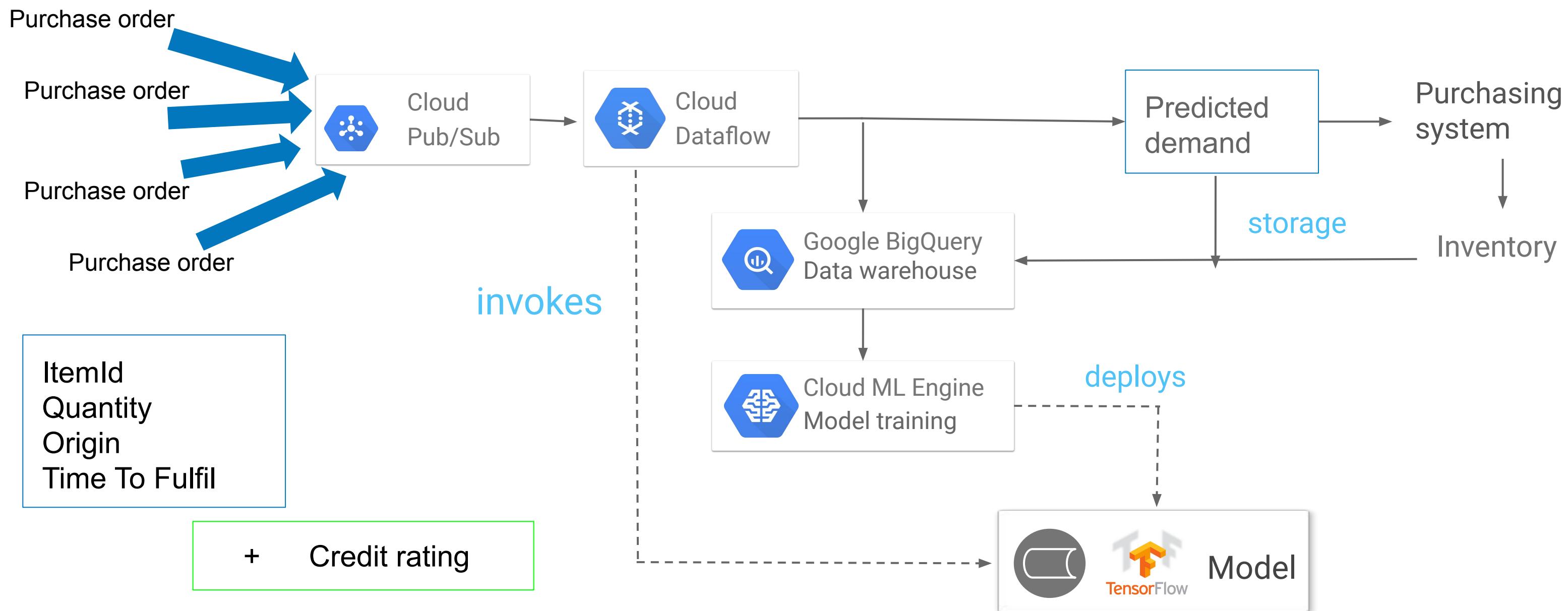
Mitigating Training-Serving Skew
Through Design

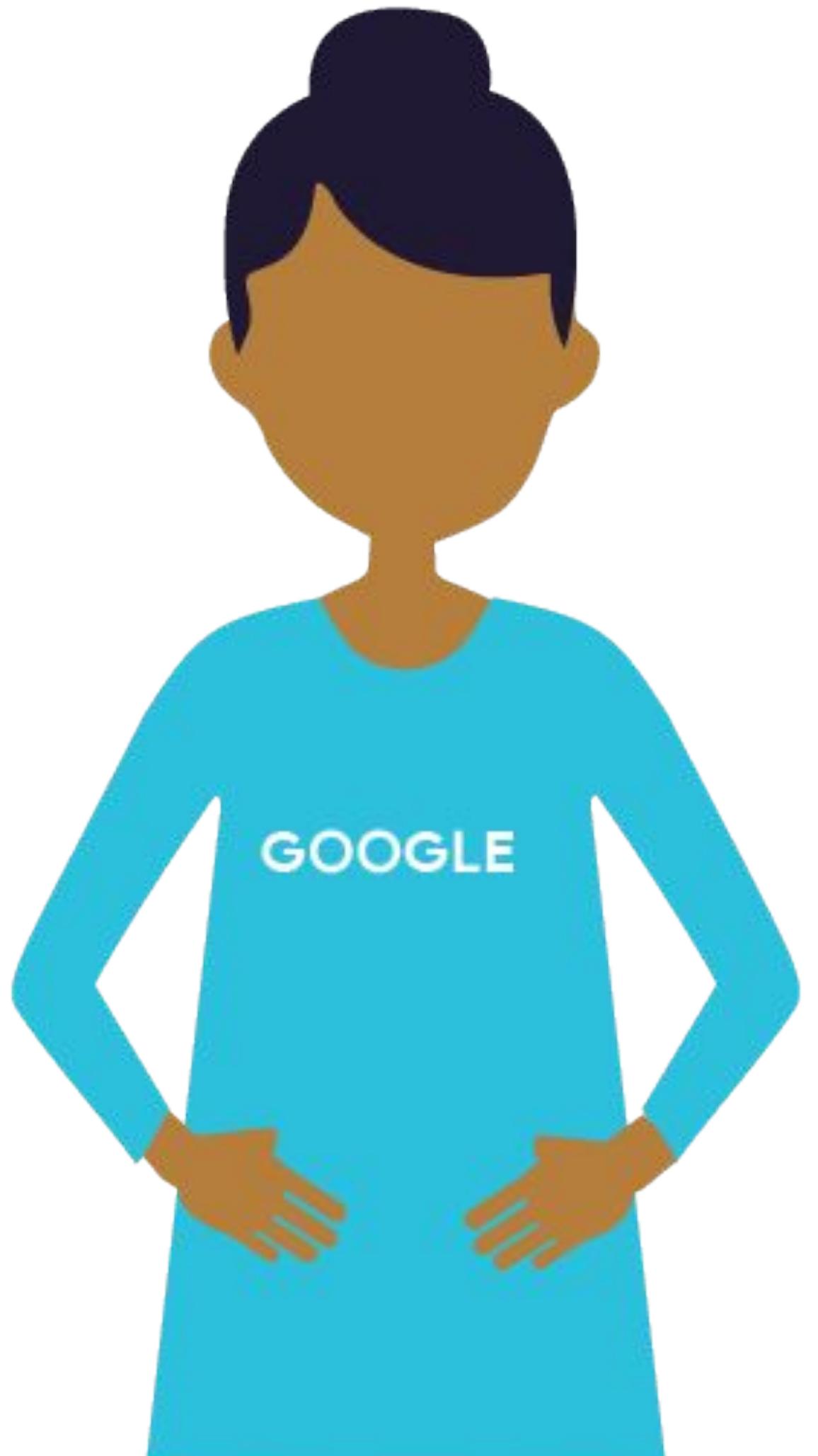
Debugging a Production Model

Predicting Widget Demand

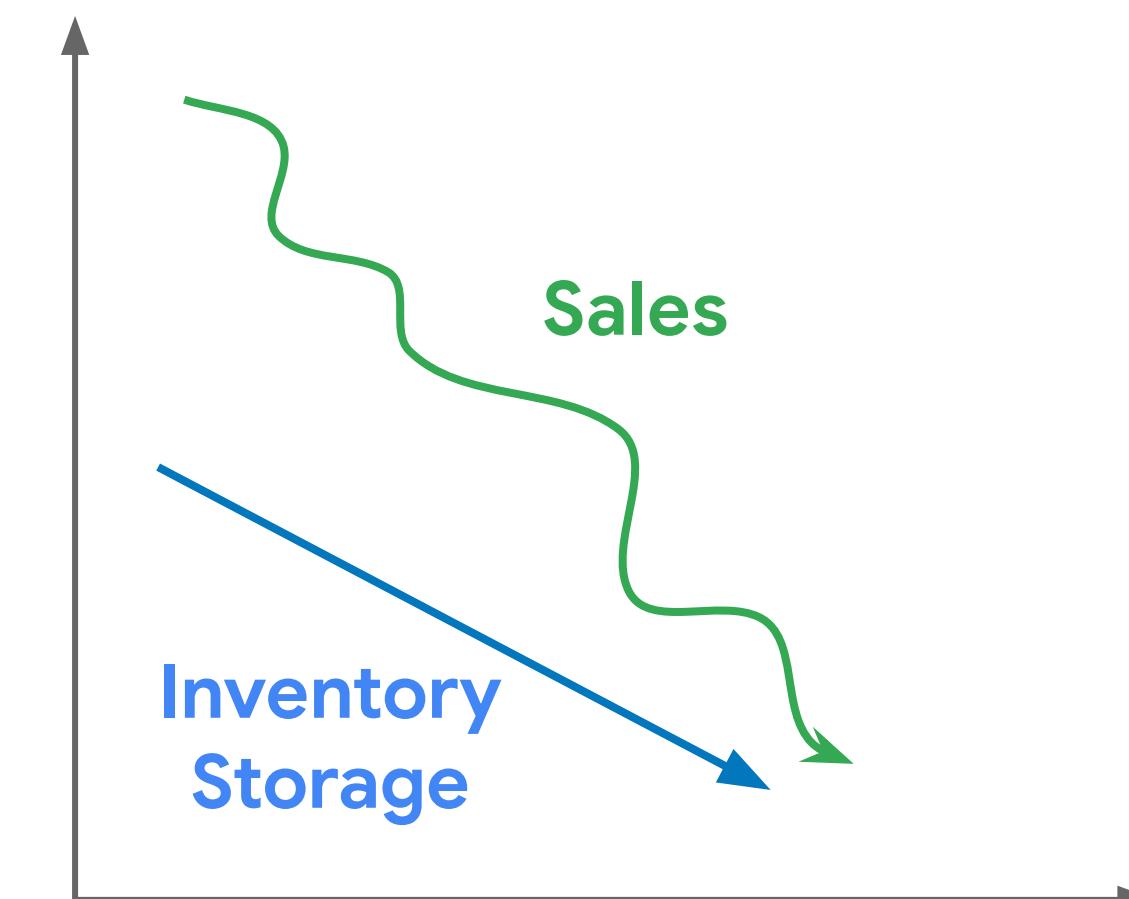


Along comes a new feature

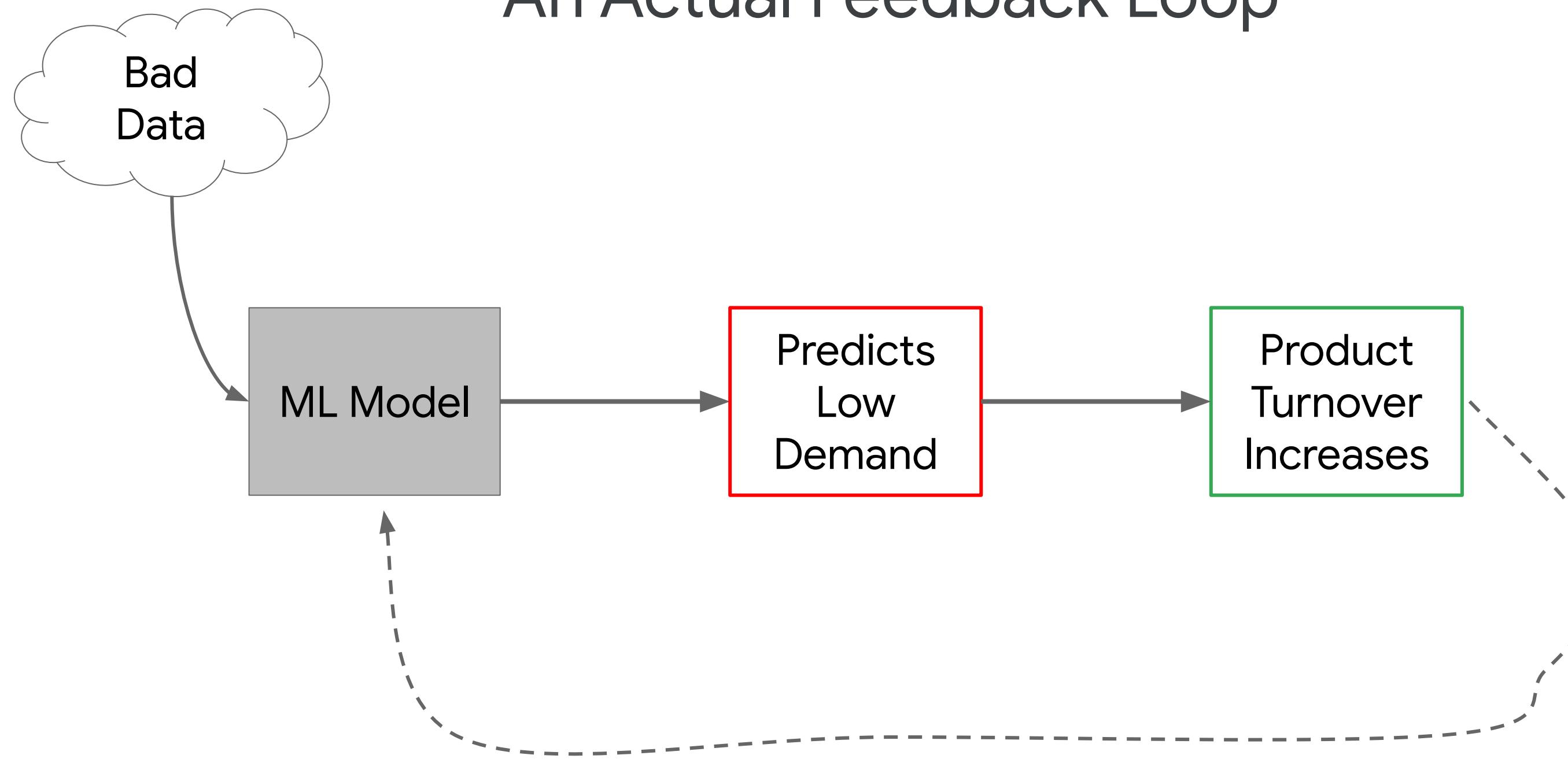




Business Catastrophe 1



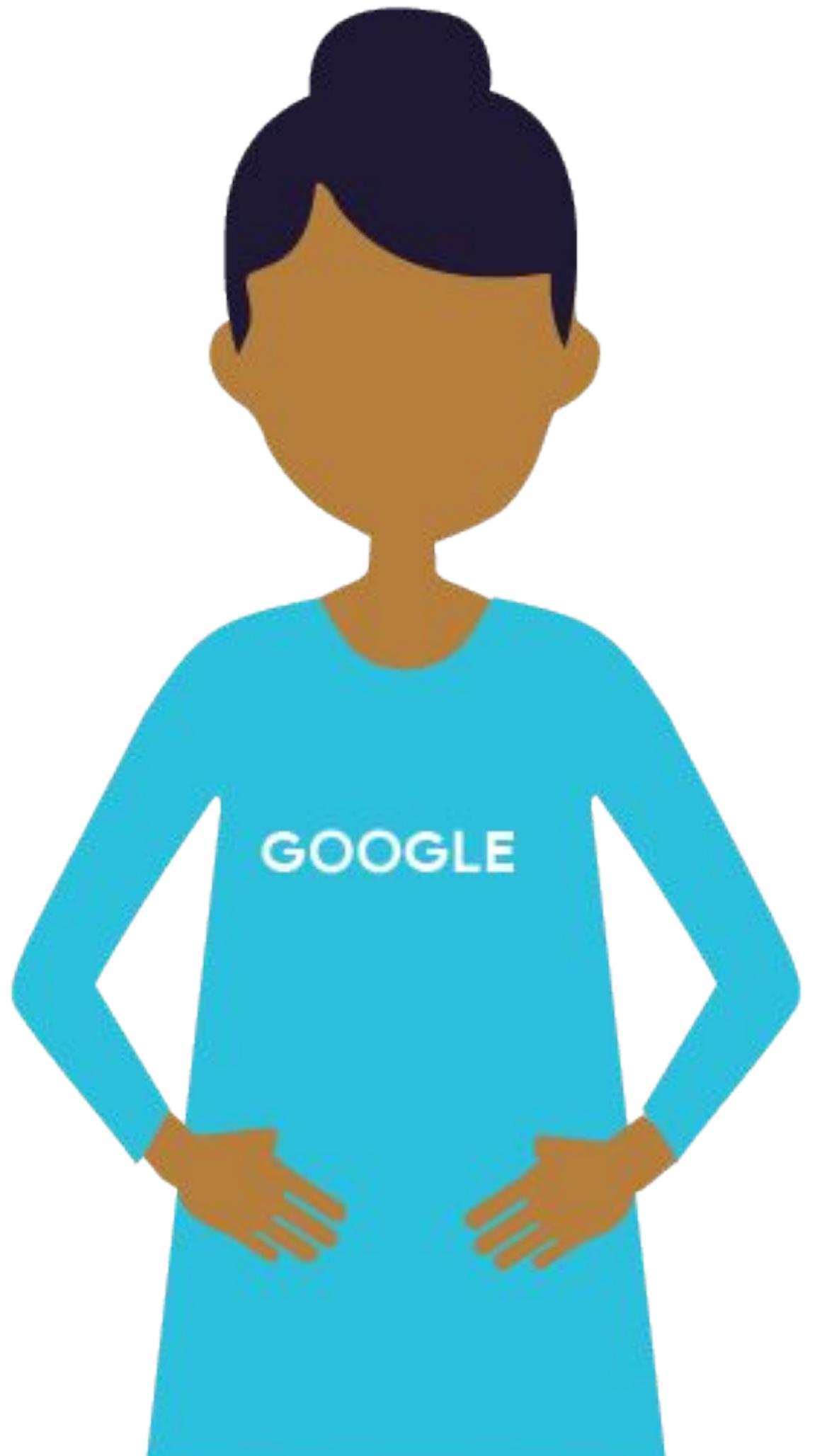
An Actual Feedback Loop



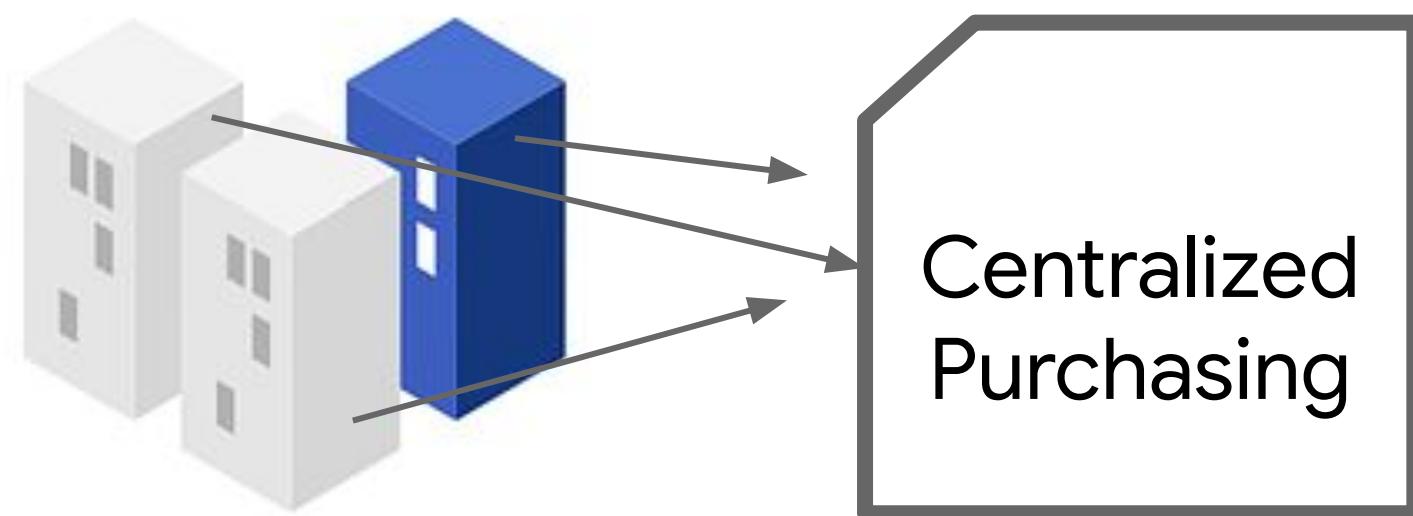


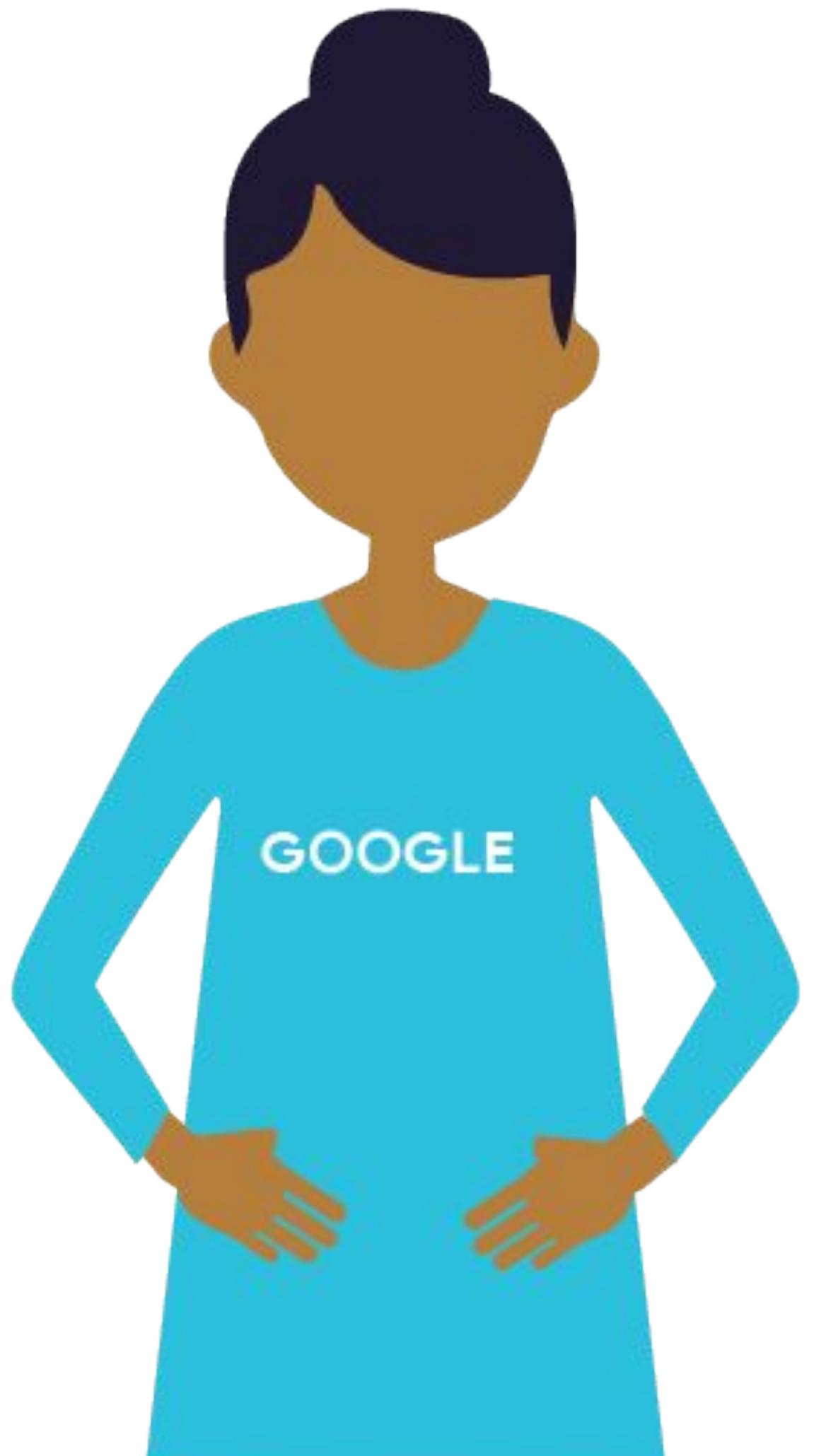
Business Catastrophe 2



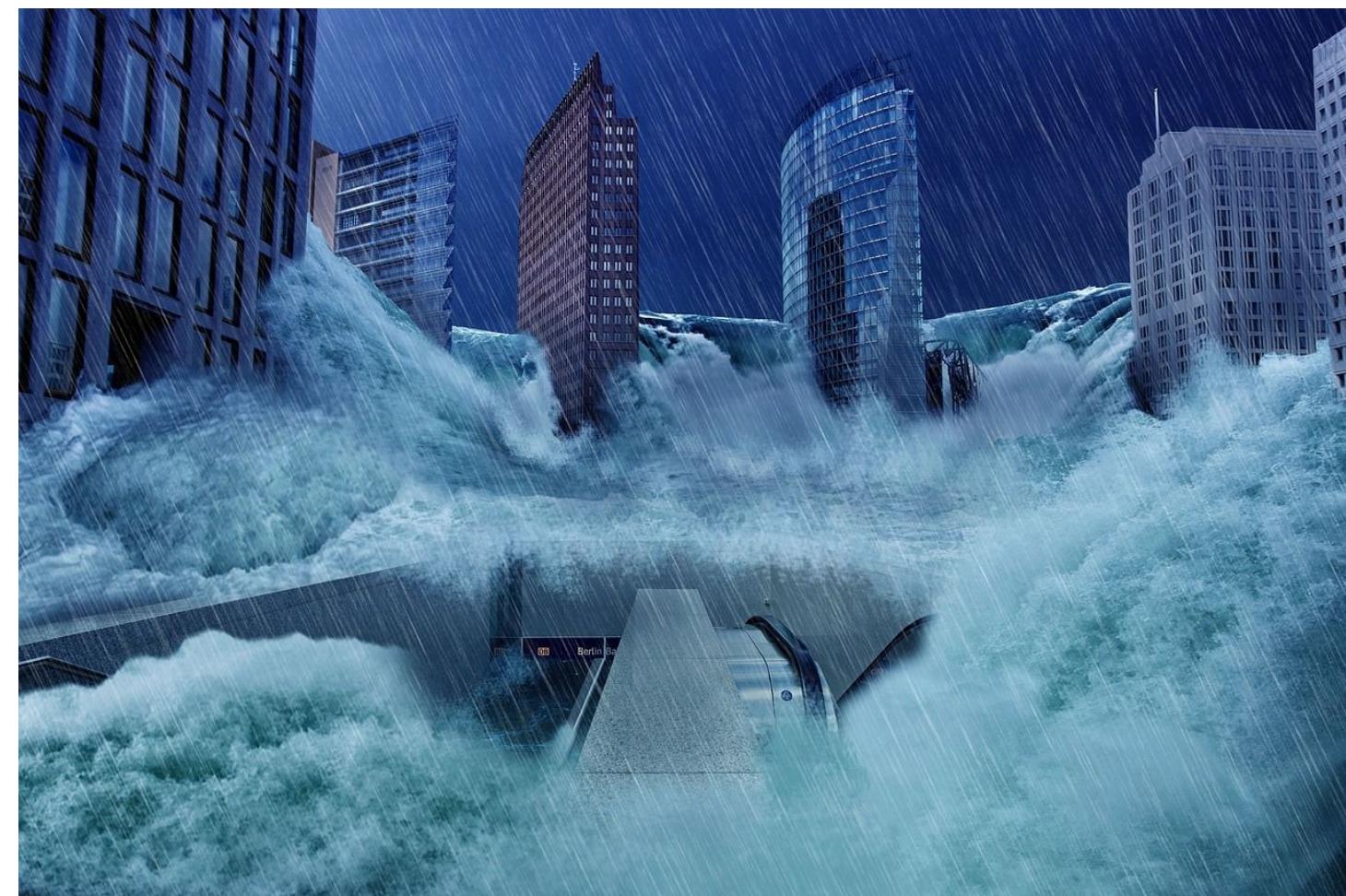


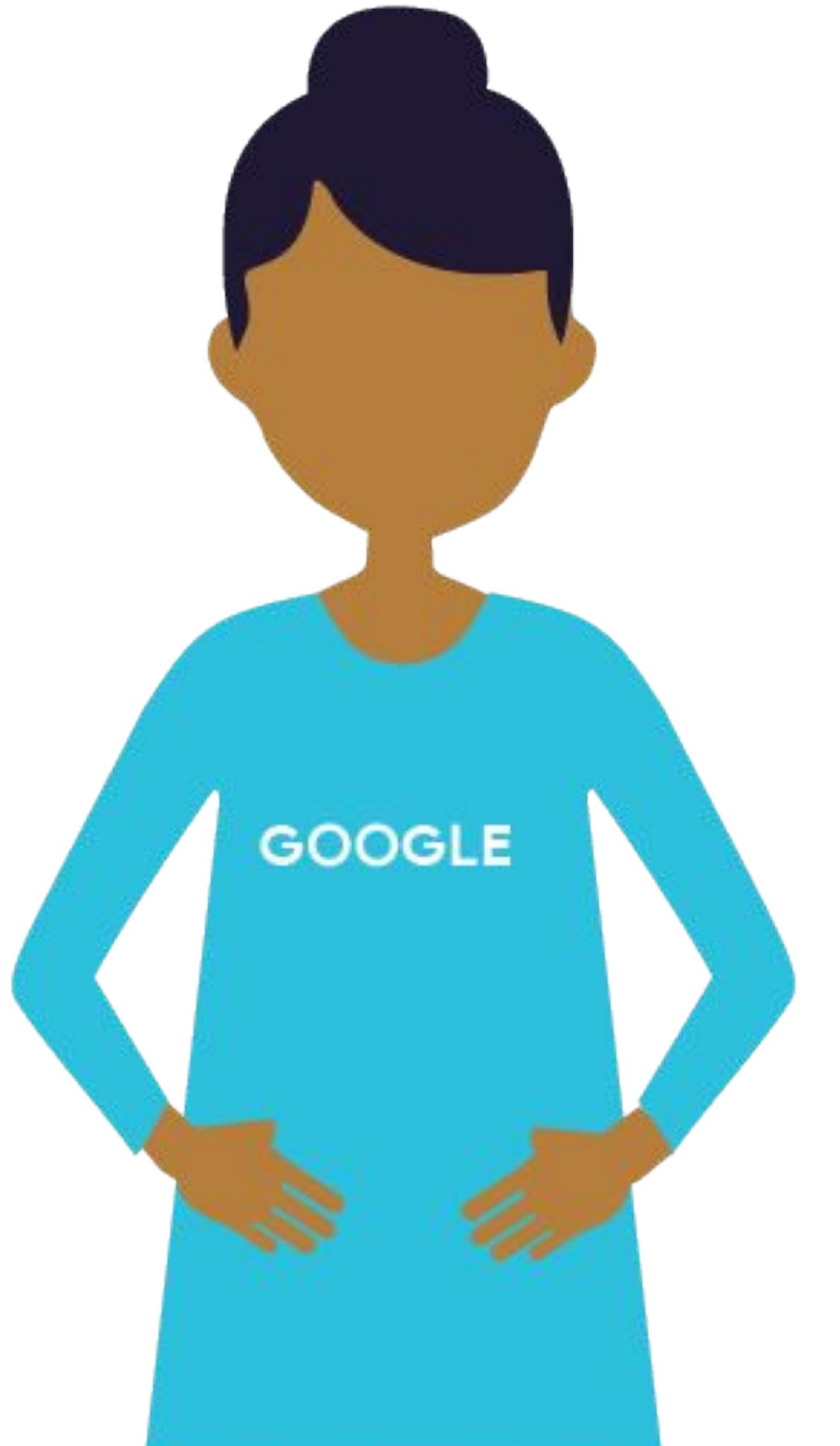
Business Catastrophe 2



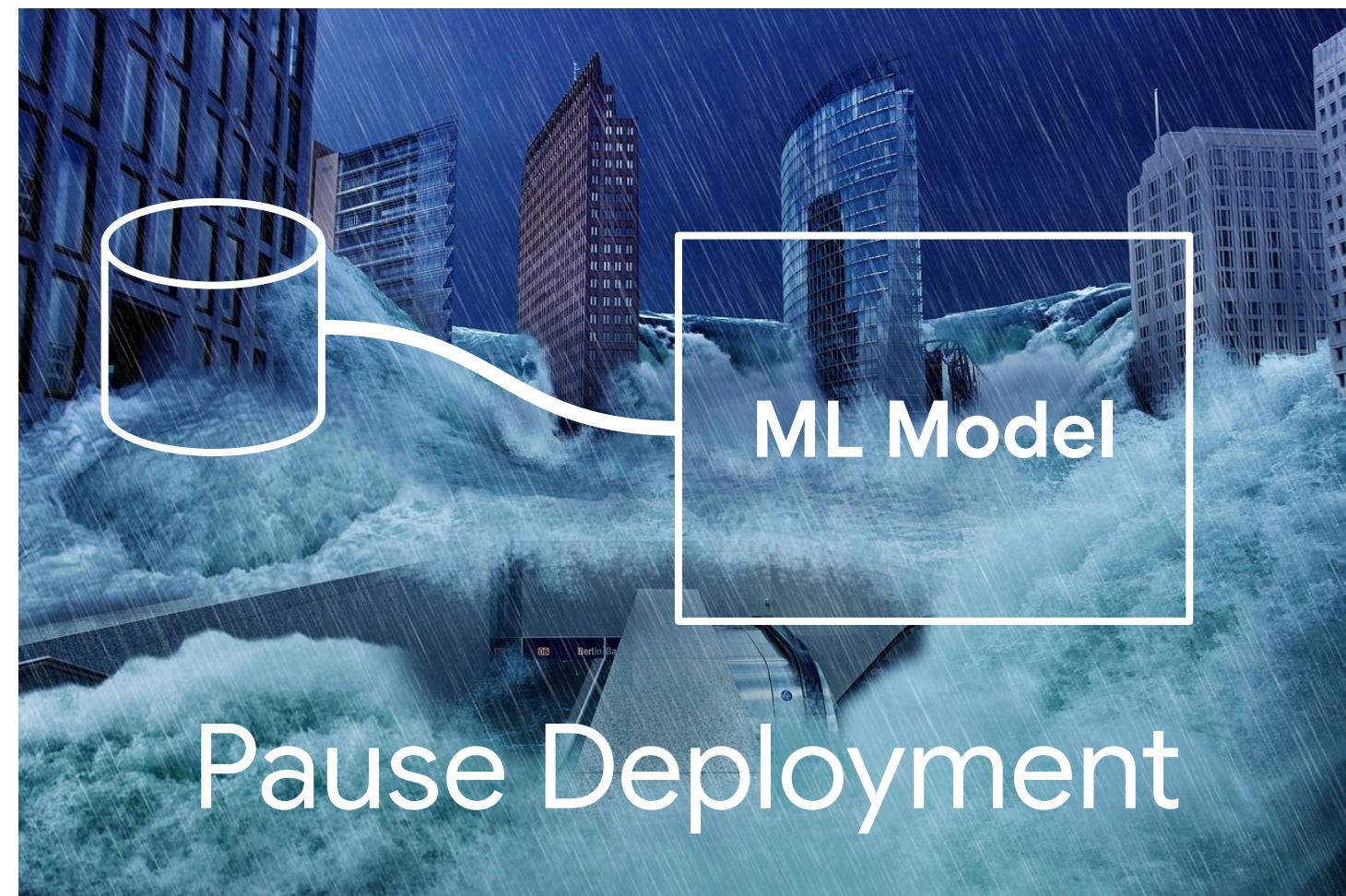


Business Catastrophe 3





Business Catastrophe 3



Course 2: Production ML Systems

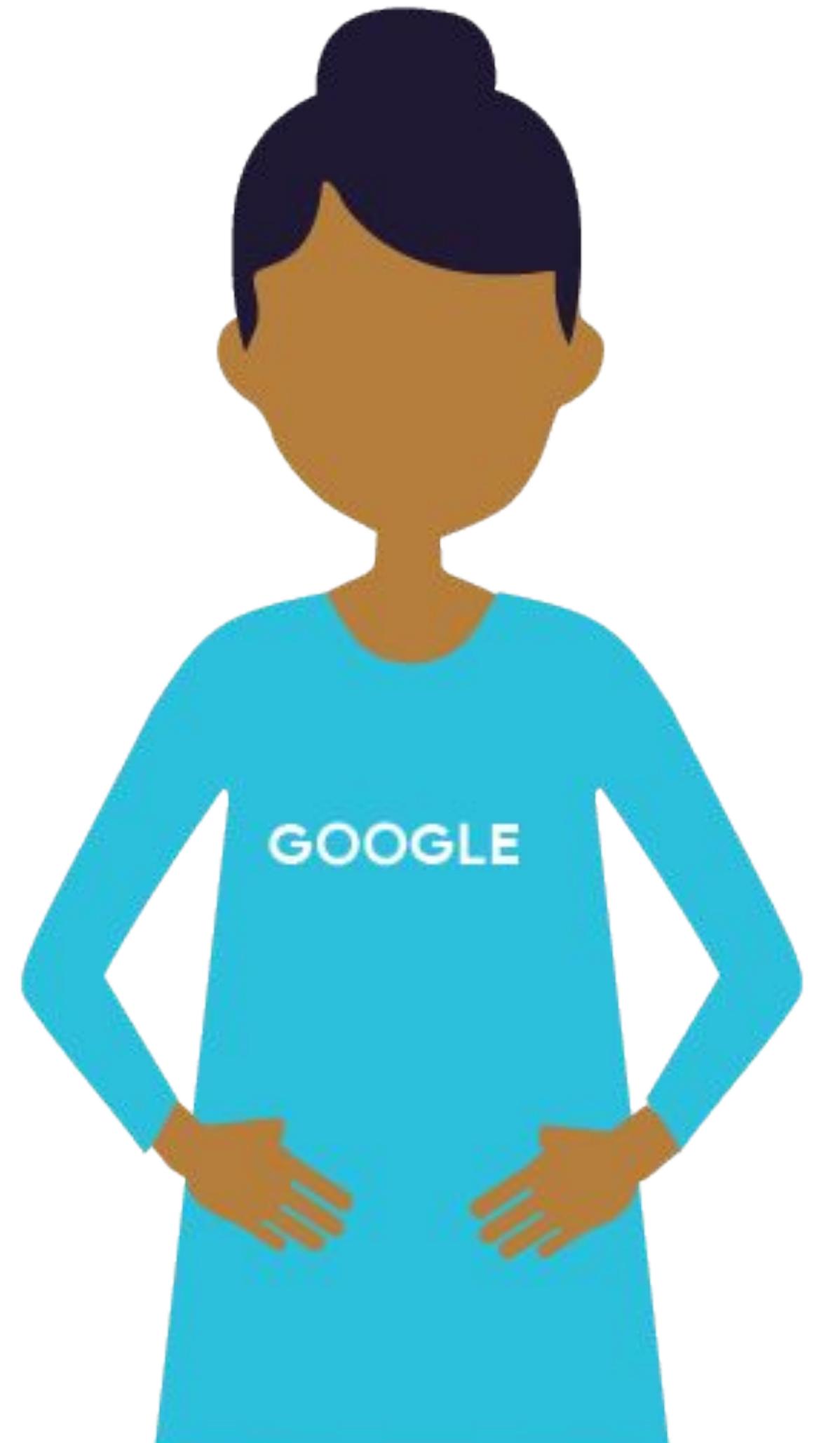
Module 3: Designing Adaptable ML Systems

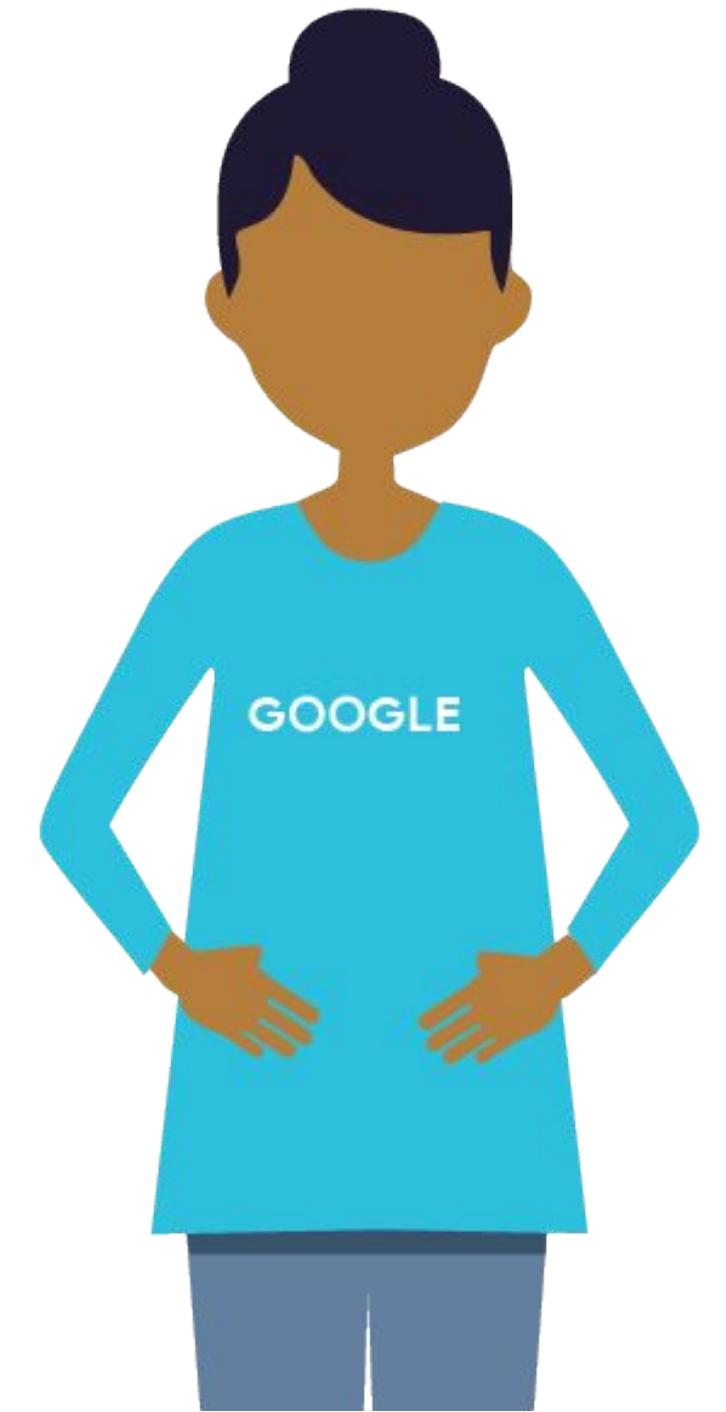
Lesson Title: **Module Summary**

Presenter: Max Lotstein

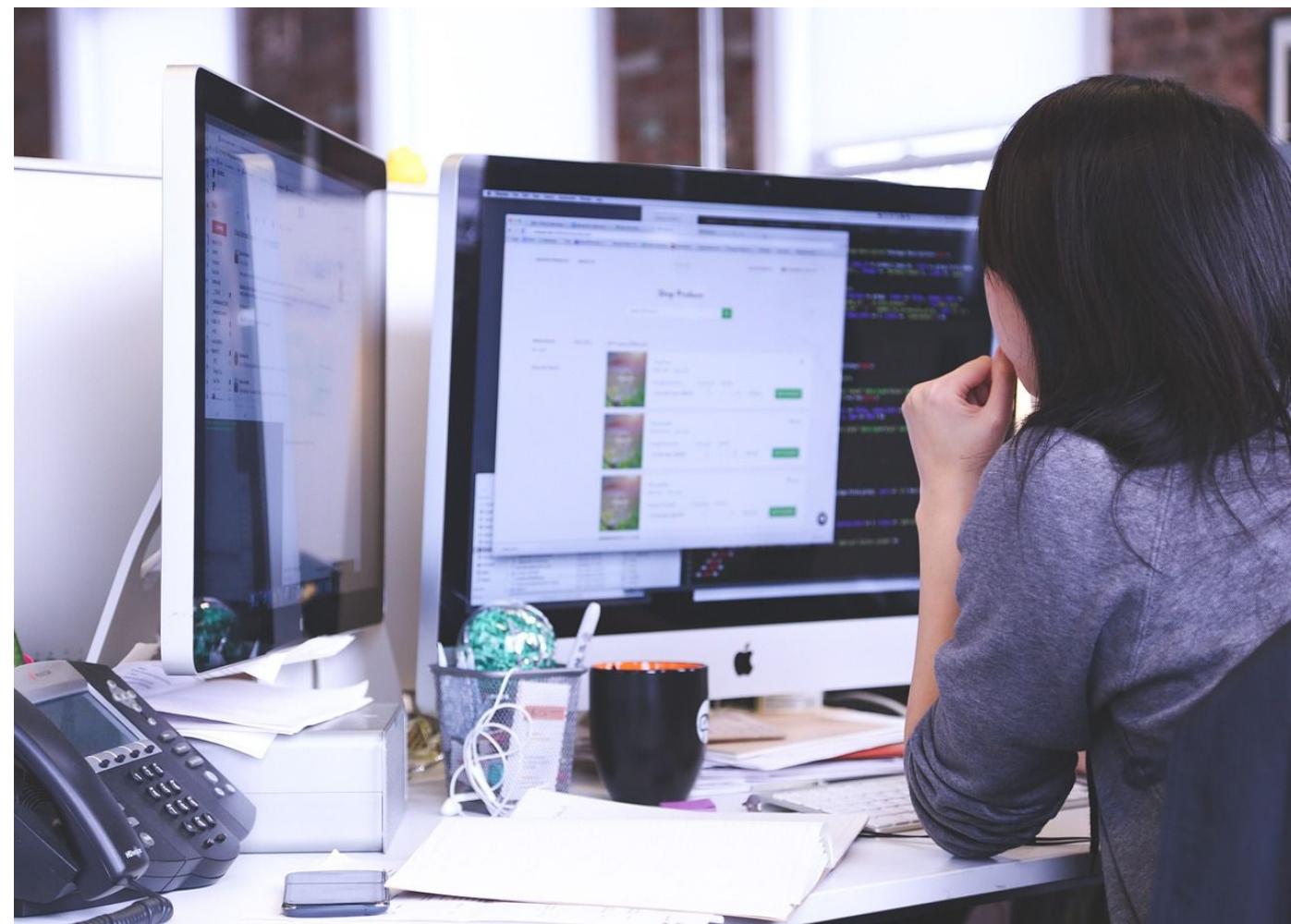
Format: Talking Head

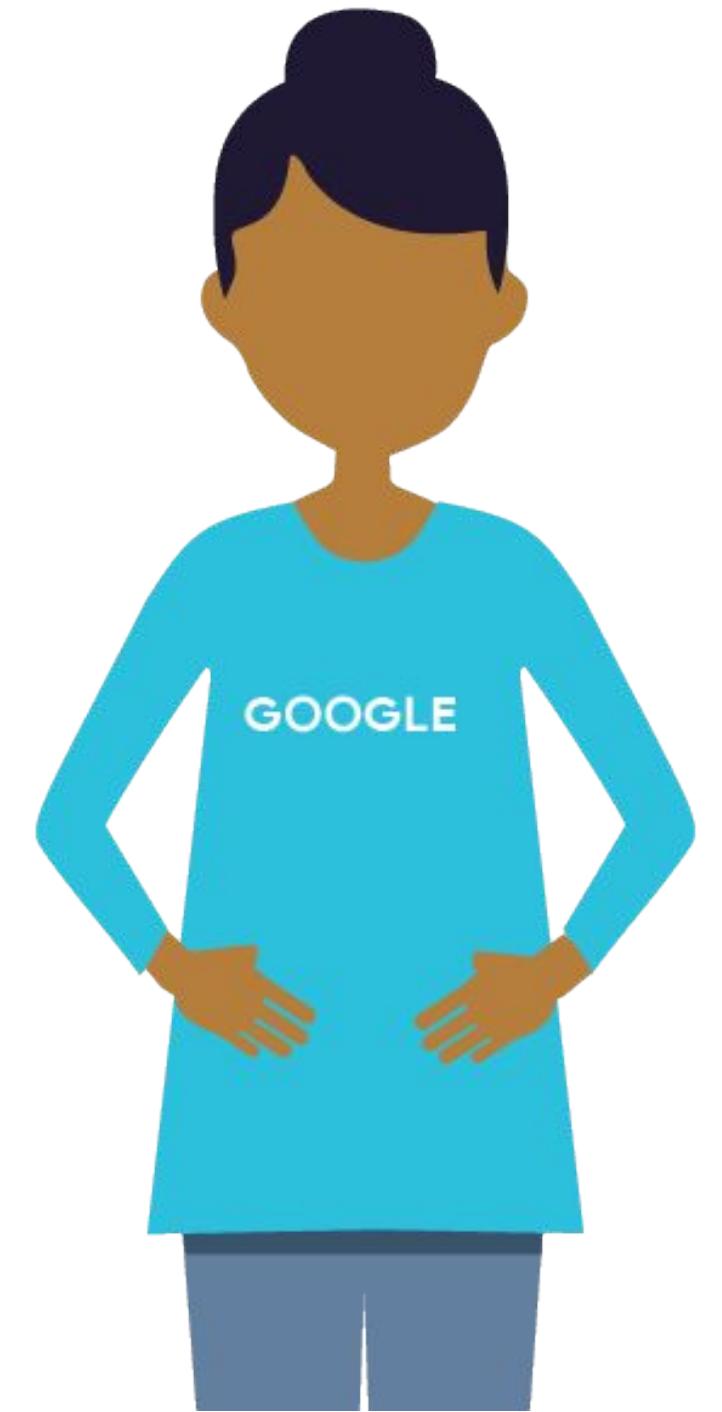
Video Name: T-PSML-0_3_l14_module_summary





Keep humans in the loop





Prioritize maintainability



Get ready to roll back

