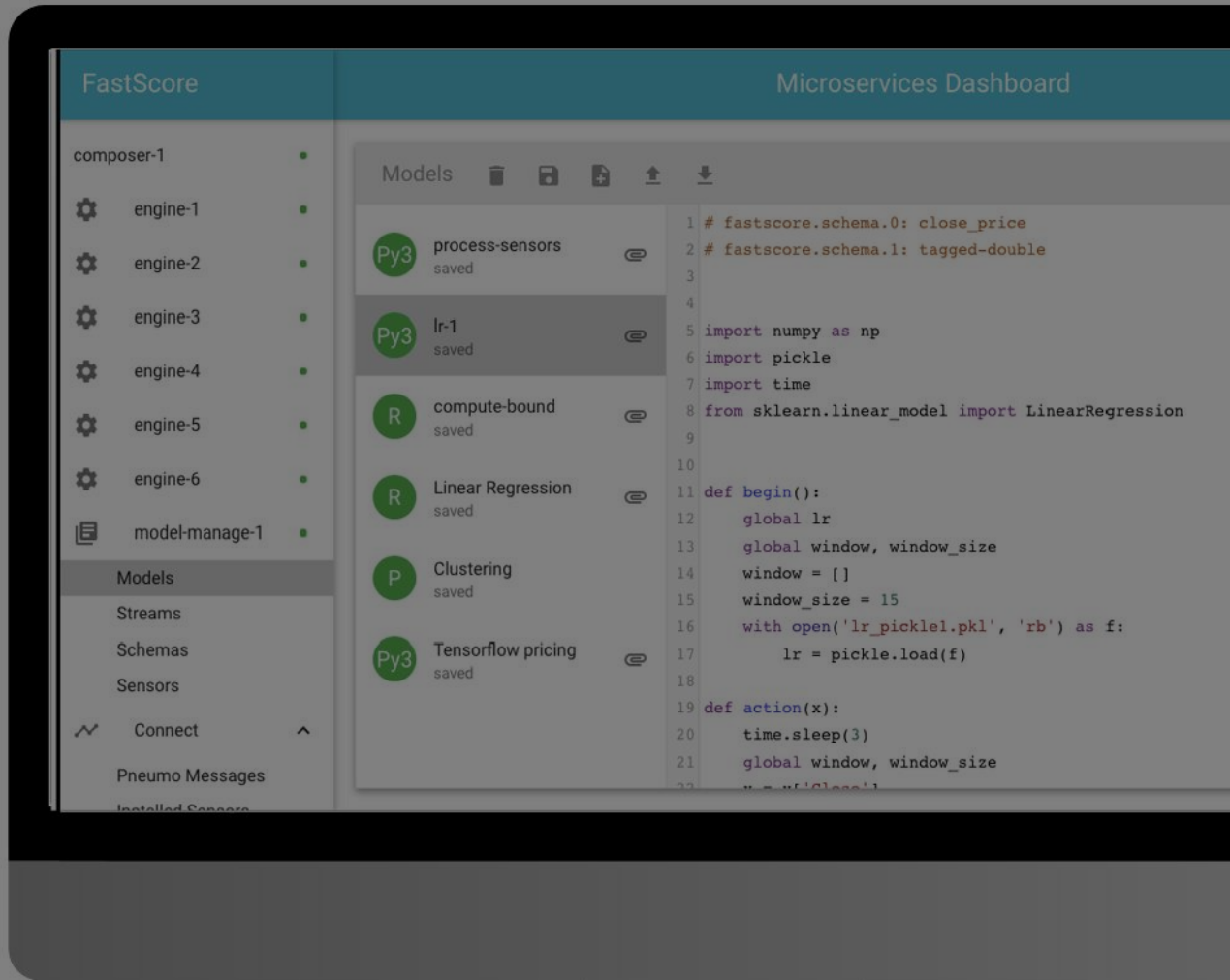# Achieving AI and ML Operational Excellence

Presented by Rehgan Avon, Product Manager

open data

# Featured Speaker

- Studied Integrated Systems Engineering specializing in Data Analytics from Ohio State University

- Launched the Big Data and Analytics Association at OSU

- Founded and coordinates the Women in Analytics Conference in Columbus, Ohio

- Product Manager at Open Data Group

open data

*"We learn more by looking for the answer to a question and not finding it than we do from learning the answer itself."* – Lloyd Alexander

*"If I had an hour to solve a problem and my life depended on the solution, I would spend the first 55 minutes determining the proper question to ask for once I know the proper question, I could solve the problem in less than 5 minutes."* – Albert Einstein

# Market Trend & Drivers

- Applied Machine Learning and AI gain traction in a majority of Fortune 500 companies, across all lines of business
- Open Source adoption continues to slowly erode traditional monolithic platform dominance of enterprise analytics
- Select analytic workloads will be transitioned to cloud (or hybrid) infrastructure, creating cost and application opportunities
- Achieving Operational Excellence on deployed ML and AI models will become as important as creating new models
- Commoditization (e.g. cloud) is allowing the enterprise to benefit from unprecedented available compute and storage at historically low costs.
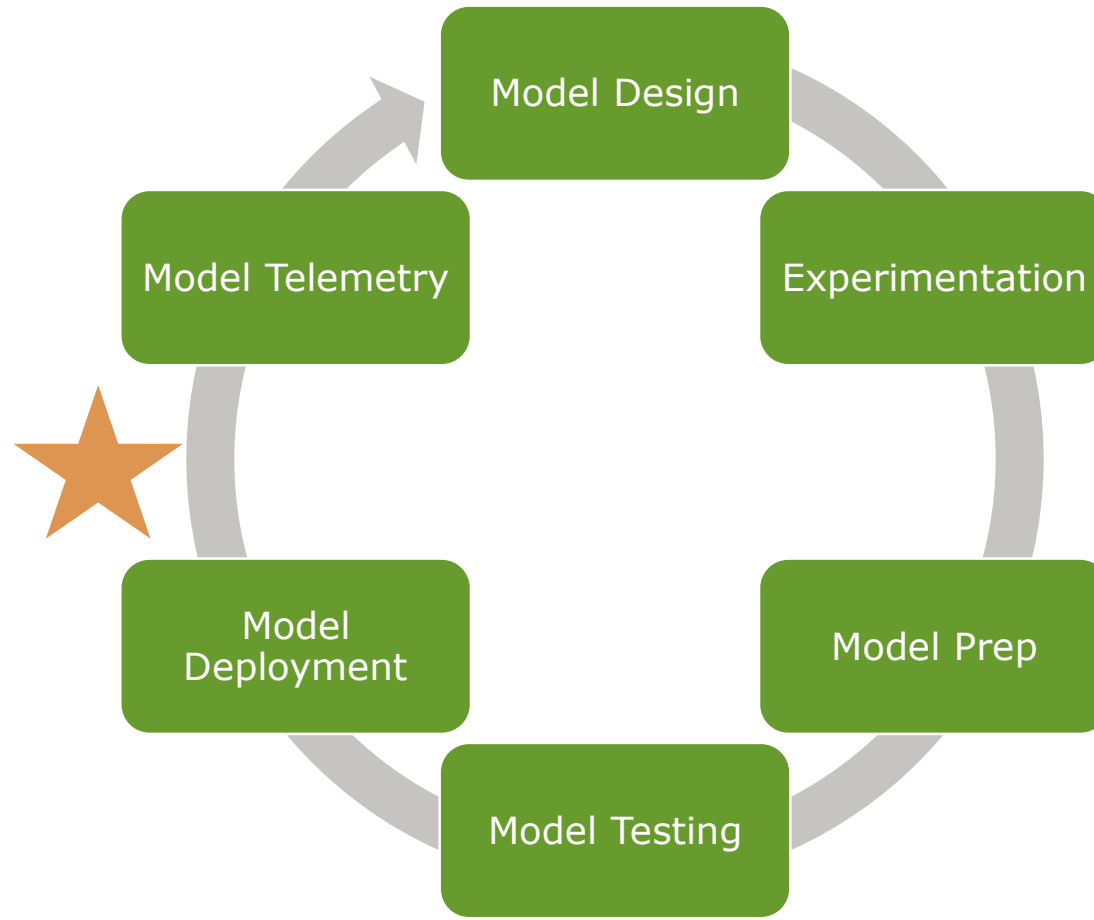
# Gartner

Market trends and insights

- *By 2020, more than 40% of data science tasks will be automated, resulting in increased productivity and broader usage by citizen data scientists.[1]*

- *The challenge now is to deploy and operationalize at scale.[2]*
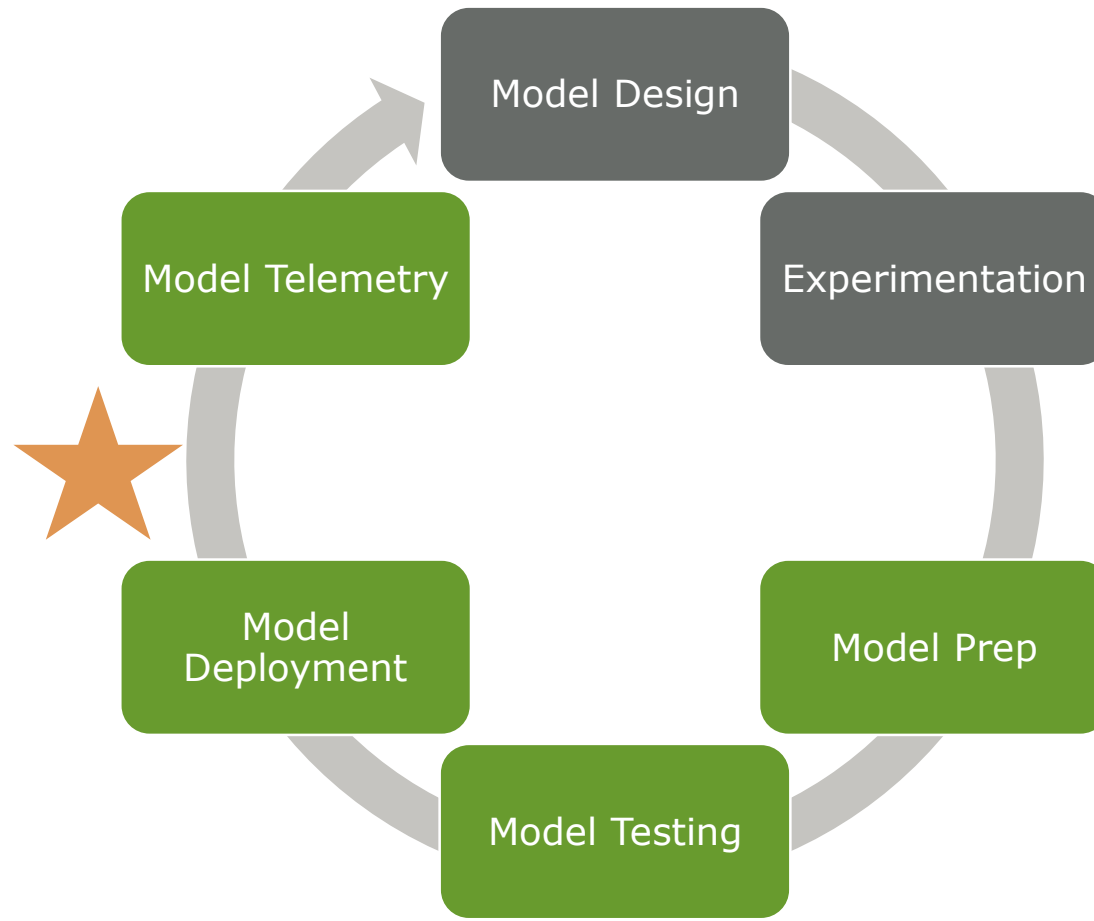
[1] *"Gartner Says More Than 40 Percent of Data Science Tasks Will Be Automated by 2020", Susan Moore et al January 2017*

[2] *"How to Operationalize Machine Learning and Data Science Projects", E Brethenoux et al July 2018*

# Model Development Life Cycle

# Model Development Life Cycle

# Top Process Challenges

Volume of Changes

Diversity of Teams Involved

Rapid Iteration

What is the solution?

Goal:  Build for high quality improvements, quickly

# Top Process Challenges

**Volume of Changes**

Diversity of Teams Involved

Rapid Iteration

Goal:  Build for high quality improvements, quickly
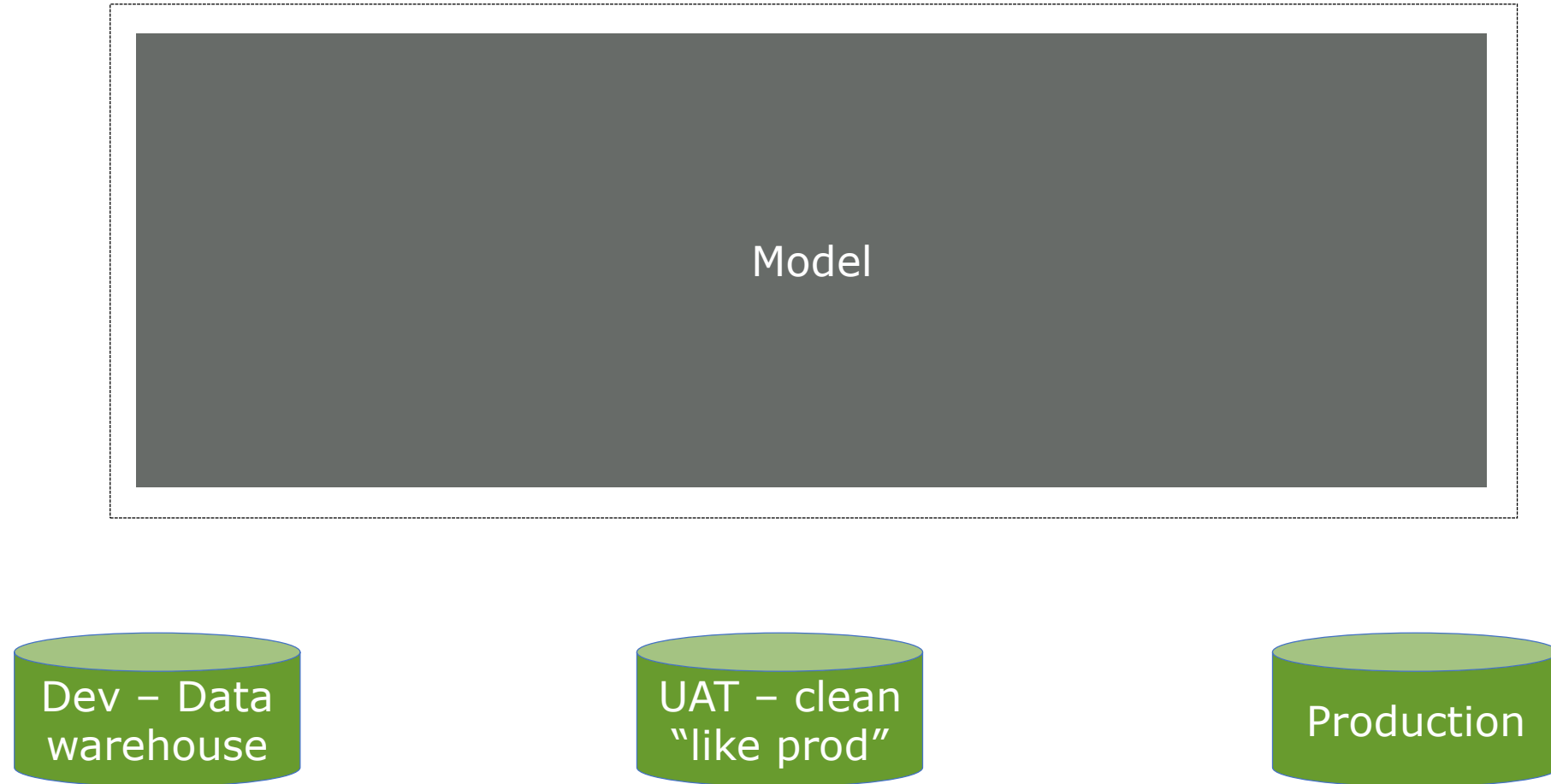
open data

# Top Process Challenges

## Volume of Changes

Data Locations
Infrastructure
Surrounding Systems
Security Requirements
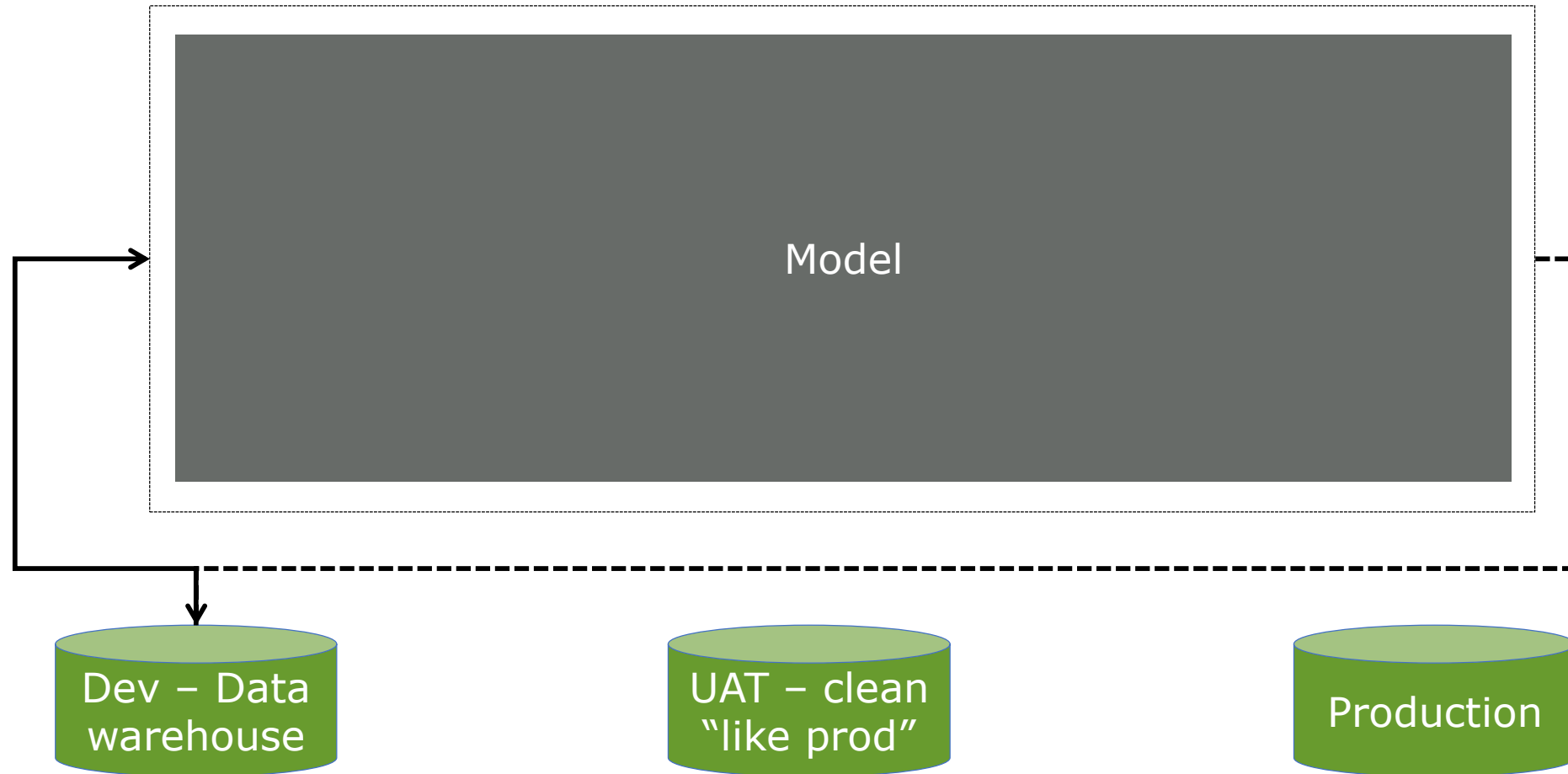
Goal:  Build for high quality improvements, quickly
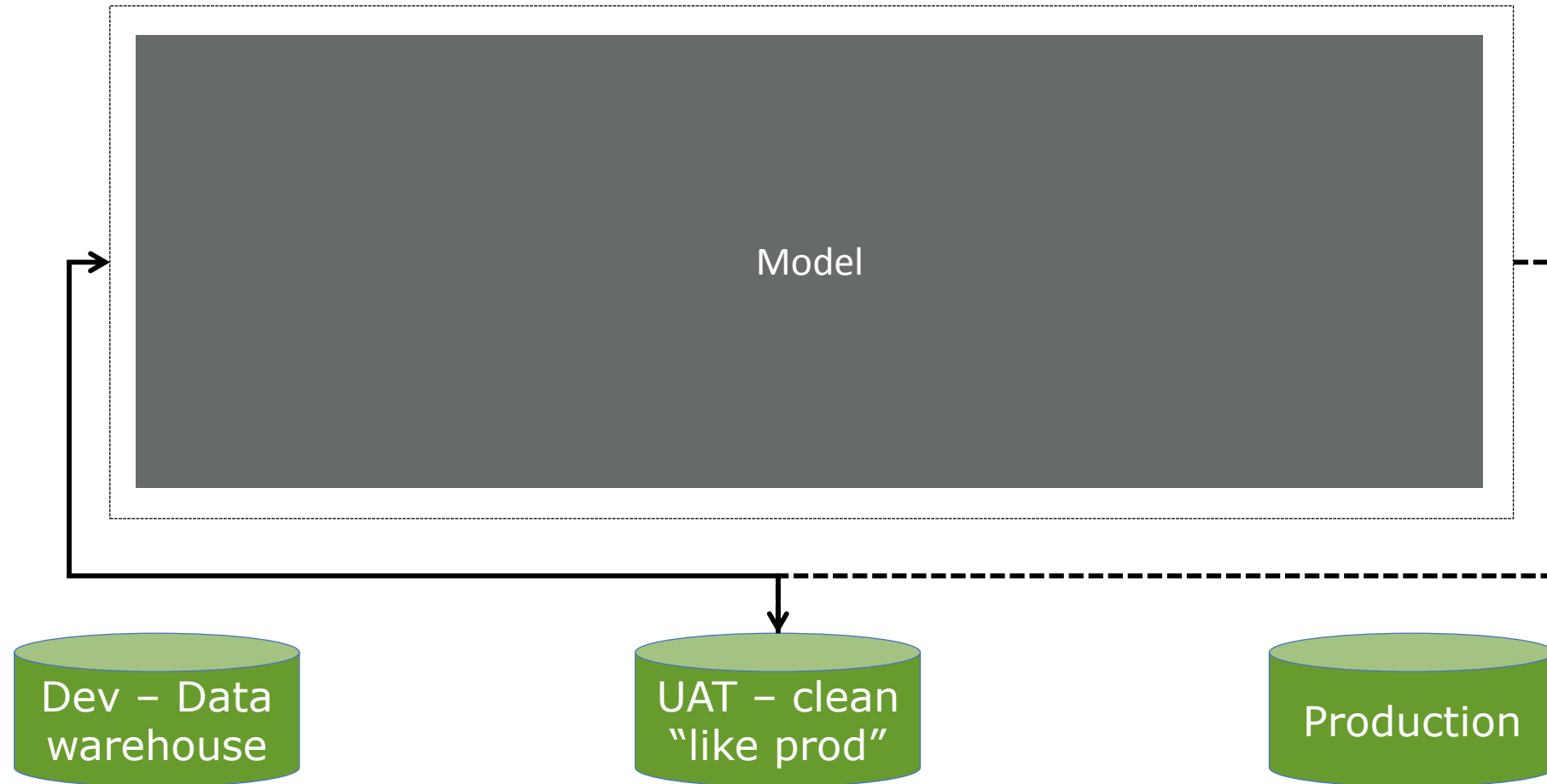
# What Can Change to Deployment?

## Data Locations

Model

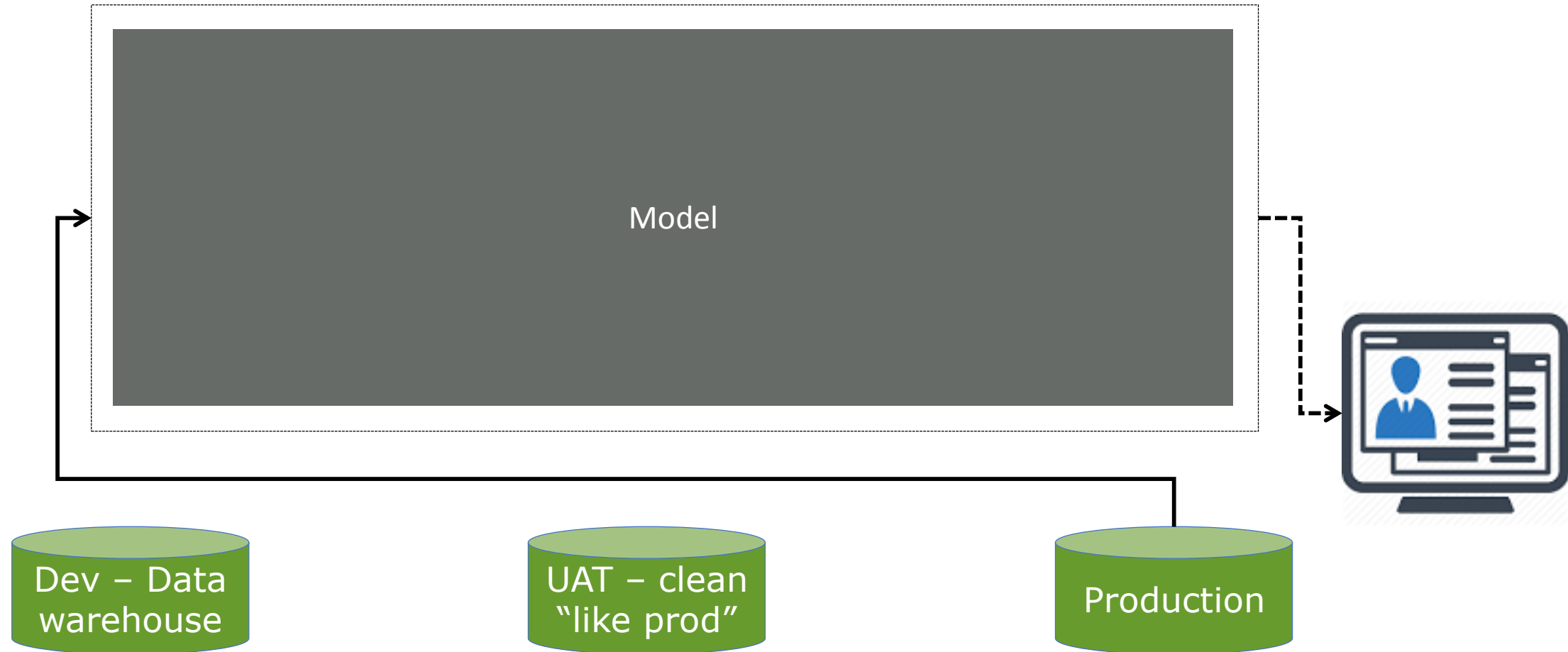Dev – Data warehouse

UAT – clean "like prod"

Production
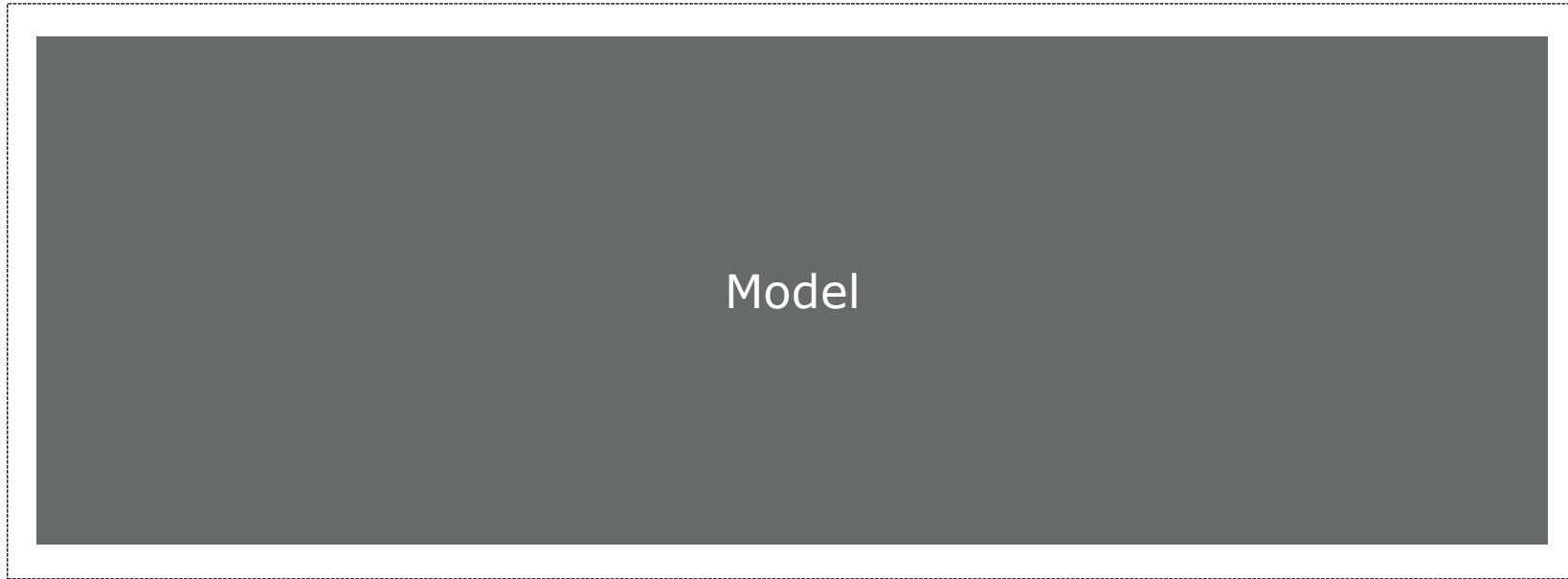
# What Can Change to Deployment?

## Data Locations - Dev

open data

# What Can Change to Deployment?

## Data Locations - Production

# What Can Change to Deployment?

Model Environment

Model

AWS-Dev

AWS-UAT

Azure-Prod

# What Can Change to Deployment?

Model Environment - Dev

open data

# What Can Change to Deployment?

Model Environment - Prod

Model

AWS-Dev

AWS-UAT

Azure-Prod

# What Can Change to Deployment?

Model Environment

Model

*Must support model language and libraries

open data

# What Can Change to Deployment?

| | Dev | UAT | Production |
|---|---|---|---|
| SCM | Open | Restricted | Read only |
| Meta data | All versions | Test results | All results |
| Security | DS test against sandbox | AI eng test against prod data | No access |

# Top Process Challenges

Volume of Changes

Diversity of Teams Involved

Rapid Iteration

What is the solution?

Goal:  Build for high quality improvements, quickly

open data

# Top Process Challenges

Volume of Changes

**Diversity of Teams Involved**

Rapid Iteration

What is the solution?

Goal:  Build for high quality improvements, quickly

# Who is Involved and What do they Care About?

Business
SLAs and Acceptance Criteria

IT
Security, costs, sustaining

Data Science
Build an impactful model

Engineering
Automation, scale, tooling

# Top Process Challenges

Volume of Changes

Diversity of Teams Involved

Rapid Iteration

What is the solution?

Goal:  Build for high quality improvements, quickly

open data

# Top Process Challenges

Volume of Changes

Diversity of Teams Involved

Rapid Iteration

What is the solution?

Goal: Build for high quality improvements, quickly

open data

# Top Process Challenges
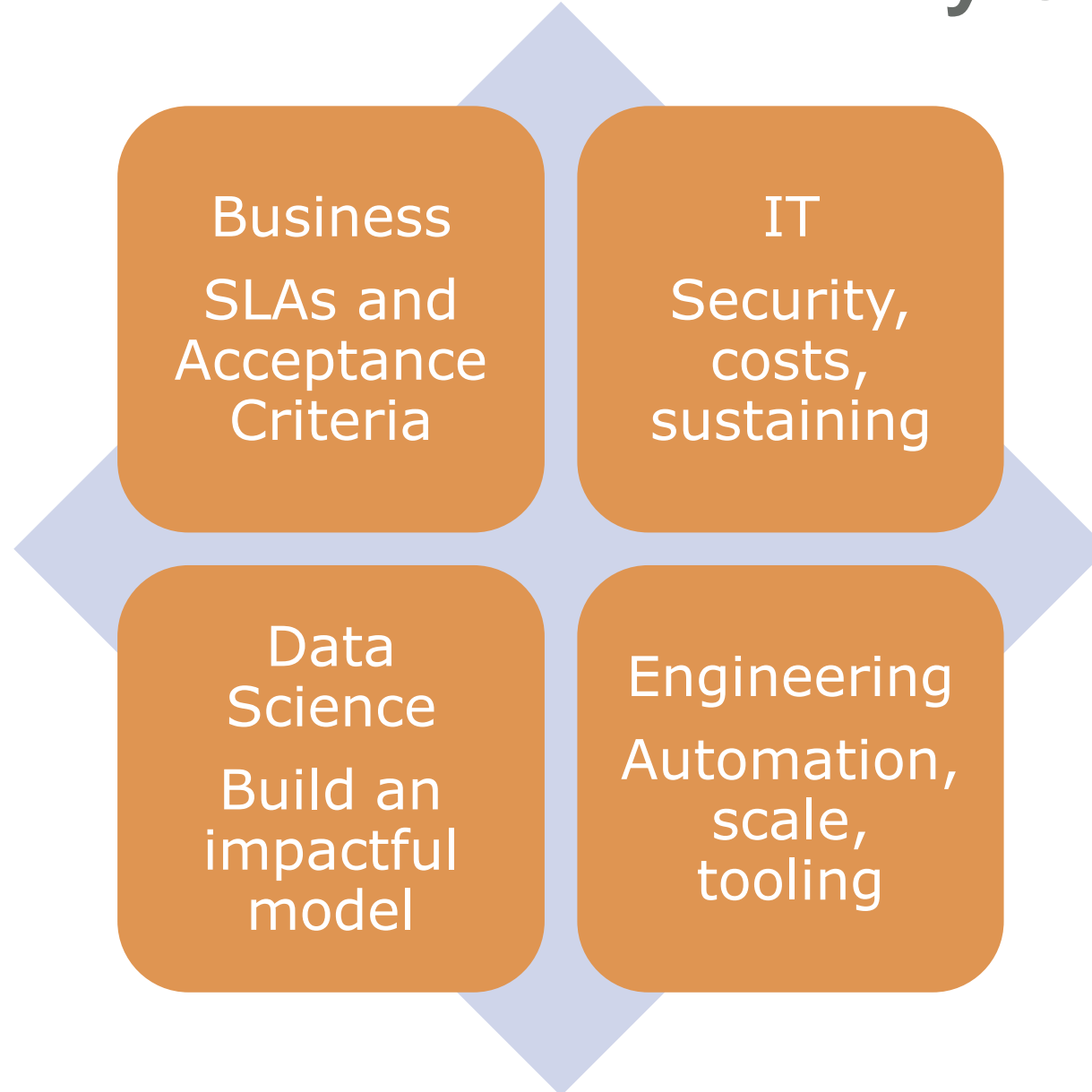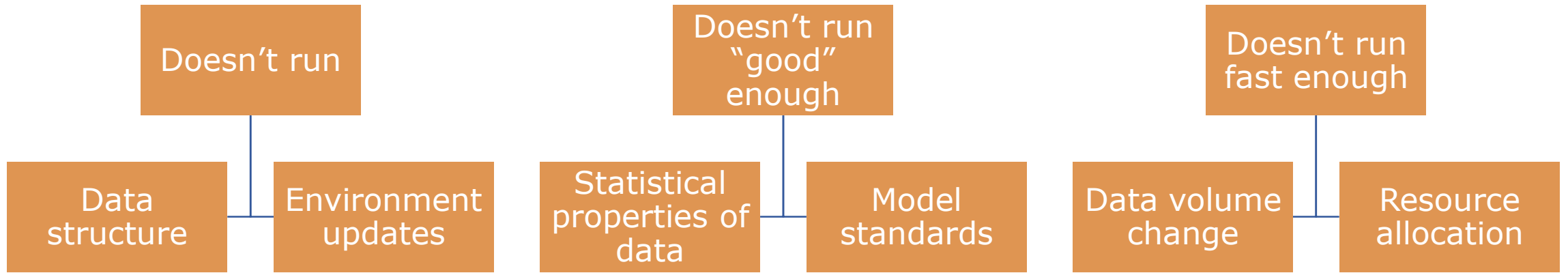
Volume of Changes

Diversity of Teams Involved

Rapid Iteration

## What is the solution?

Goal:  Build for high quality improvements, quickly

open data

# Creating Meaningful Abstractions to Supplement the Process

| Abstraction | Description |
|---|---|
| Model Artifacts – Coefficients | Referenceable statistics on trained data |
| Execution Ready Model | Scoring Code |
| Model Data Definitions | Avro Schemas |
| Data Transport Descriptors | JSON file with necessary information |
| Model Environment Requirements | Mode library/code dependencies |
| Performance Monitoring Models (Telemetry) | Models custom to report performance statistics |
| Sensors | Computational performance tracking |

# Core Concepts to Enable MDLC

Model

open data

# Core Concepts to Enable MDLC

```python
1   # fastscore.input: gbm_input
2   # fastscore.output: gbm_output
3
4   import cPickle
5   import numpy as np
6   import pandas as pd
7
8   from sklearn.ensemble import GradientBoostingRegressor
9   from sklearn.pipeline import Pipeline
10
11  import imp
12
13  def begin():
14      FeatureTransformer = imp.load_source("FeatureTransformer",
15                                  "score_auto_gbm/FeatureTransformer.py")
16      global gbmFit
17      with open("score_auto_gbm/gbmFit.pkl", "rb") as pickle_file:
18          gbmFit = cPickle.load(pickle_file)
19
20  def action(datum):
21      score = list(gbmFit.predict(pd.DataFrame([datum])))[0]
22      yield score
```

# Core Concepts to Enable MDLC

31 lines (30 sloc) | 1.09 KB

```
1    {
2      "type": "record",
3      "name": "CarRecord",
4      "fields": [
5        {"name": "make", "type": "string"},
6        {"name": "fuelType", "type": "string"},
7        {"name": "aspiration", "type": "string"},
8        {"name": "numDoors", "type": "string"},
9        {"name": "bodyStyle", "type": "string"},
10       {"name": "driveWheels", "type": "string"},
11       {"name": "engineLocation", "type": "string"},
12       {"name": "wheelBase", "type": "double"},
13       {"name": "length", "type": "double"},
14       {"name": "width", "type": "double"},
15       {"name": "height", "type": "double"},
16       {"name": "curbWeight", "type": "int"},
17       {"name": "engineType", "type": "string"},
18       {"name": "numCylinders", "type": "string"},
19       {"name": "engineSize", "type": "int"},
20       {"name": "fuelSystem", "type": "string"},
21       {"name": "bore", "type": "double"},
22       {"name": "stroke", "type": "double"},
23       {"name": "compressionRatio", "type": "double"},
24       {"name": "horsepower", "type": "int"},
25       {"name": "peakRPM", "type": "int"},
26       {"name": "cityMPG", "type": "int"},
27       {"name": "highwayMPG", "type": "int"},
28       {"name": "price", "type": "int"}
29     ]
30   }
```

2 lines (1 sloc) | 19 Bytes

```
1    {"type": "double"}
```

# Core Concepts to Enable MDLC

14 lines (13 sloc) | 195 Bytes

```
 1  {
 2    "Transport": {
 3      "Type": "file",
 4      "Path": "/root/data/input_data.jsons"
 5    },
 6    "Envelope": {
 7      "Type":"delimited"
 8    },
 9    "Encoding": "json",
10    "Schema": {
11      "$ref":"gbm_input"
12    }
13  }
```
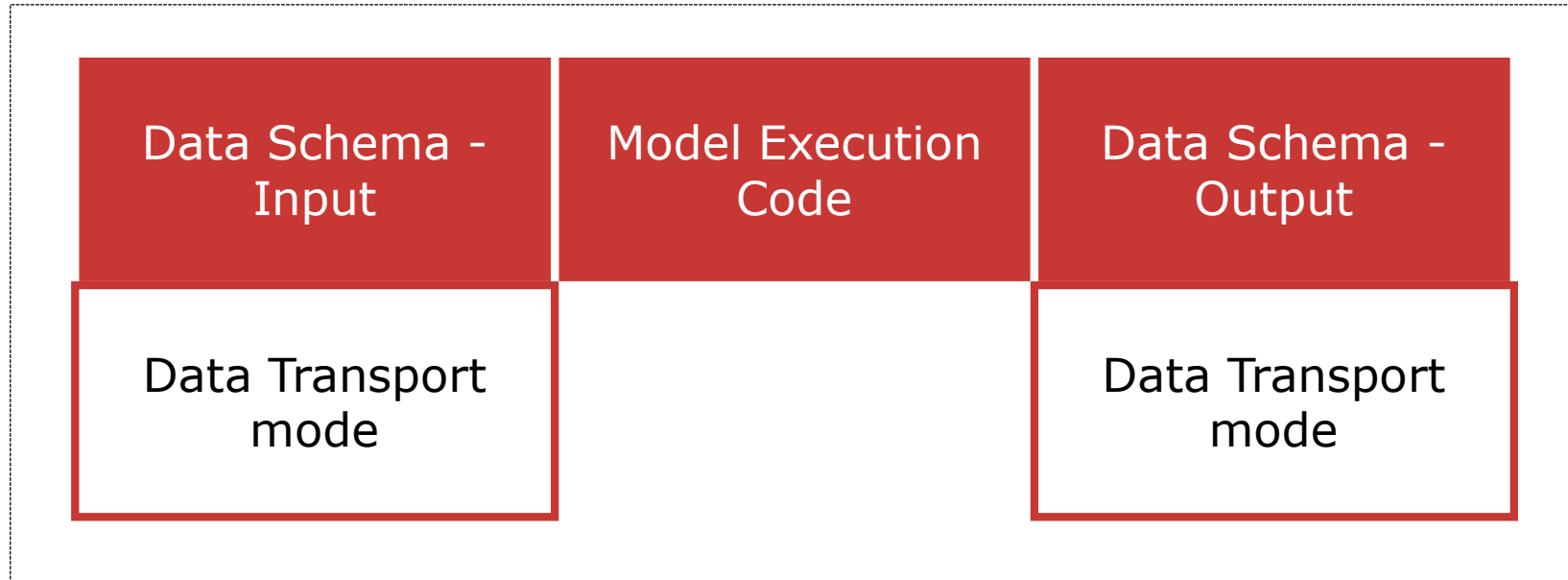
13 lines (12 sloc) | 225 Bytes

```
 1  {
 2    "Version": "1.2",
 3    "Loop": false,
 4    "Transport": {
 5      "Type": "kafka",
 6      "BootstrapServers": ["172.17.0.1:9092"],
 7      "Topic": "input"
 8    },
 9    "Envelope": null,
10    "Encoding": "json",
11    "Schema": {"$ref":"gbm_input"}
12  }
```
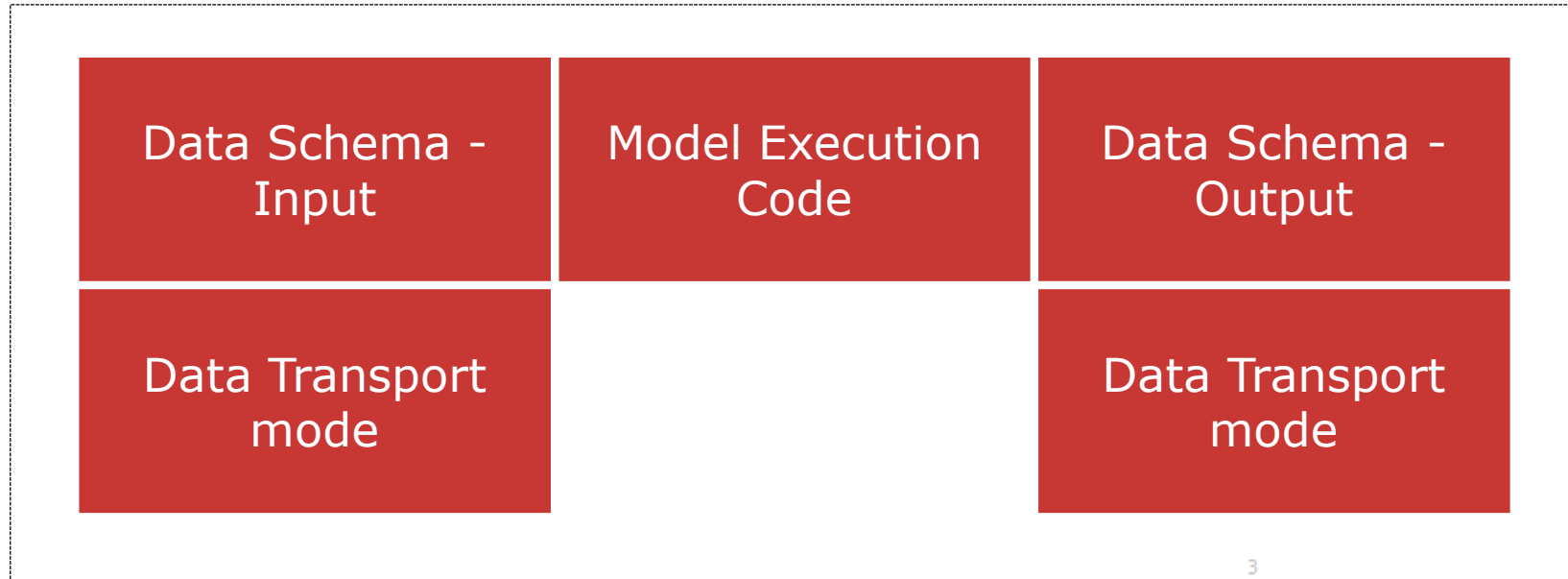
# What Can Change?

Transport Modes

| Data Schema - Input | Model Execution Code | Data Schema - Output |
|---|---|---|
| Data Transport mode | | Data Transport mode |

Will depend on use case, phase, and application needs: Batch, on-demand or "streaming"

# Core Concepts to Enable MDLC

## Language Agnostic

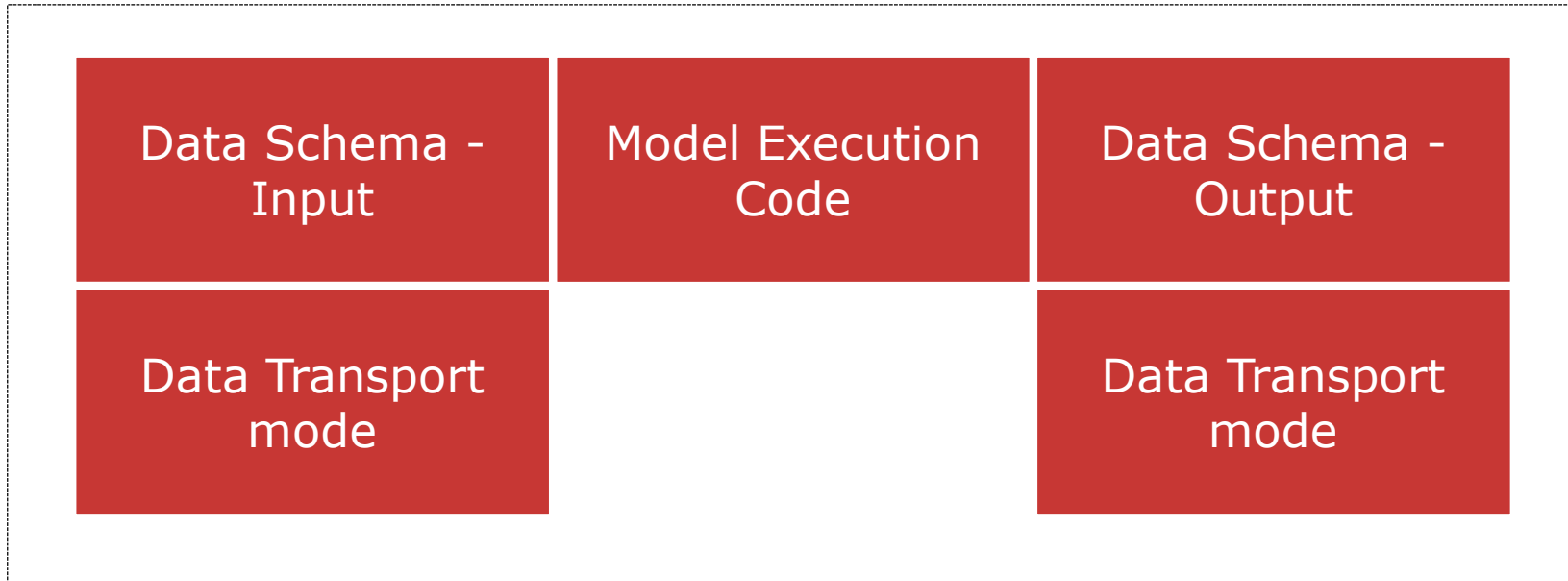| | | |
|---|---|---|
| Data Schema - Input | Model Execution Code | Data Schema - Output |
| Data Transport mode | | Data Transport mode |

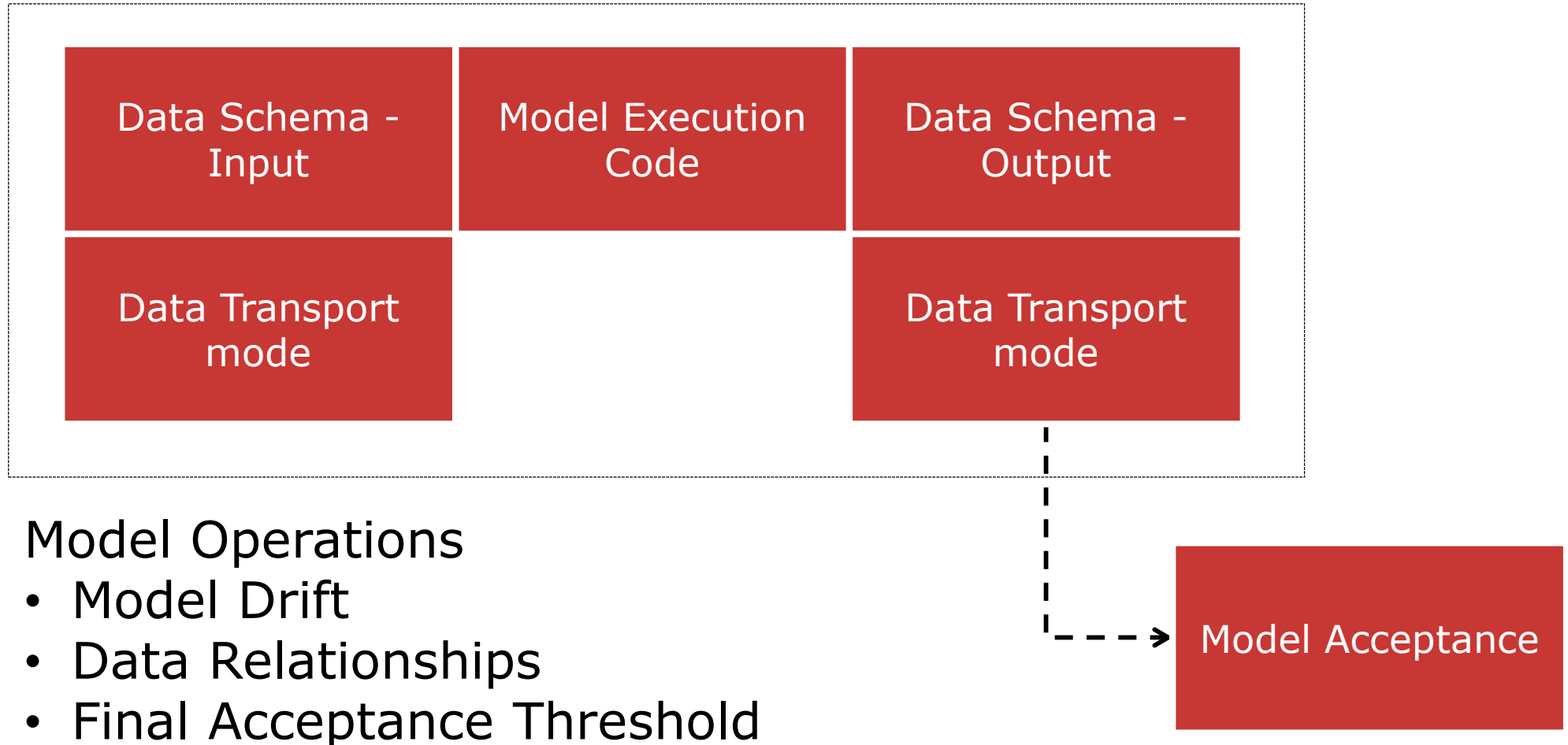Infrastructure Independence (Microservice)

```
3
4    import cPickle
5    import numpy as np
6    import pandas as pd
7
8    from sklearn.ensemble import GradientBoostingRegressor
9    from sklearn.pipeline import Pipeline
10
11   import imp
```
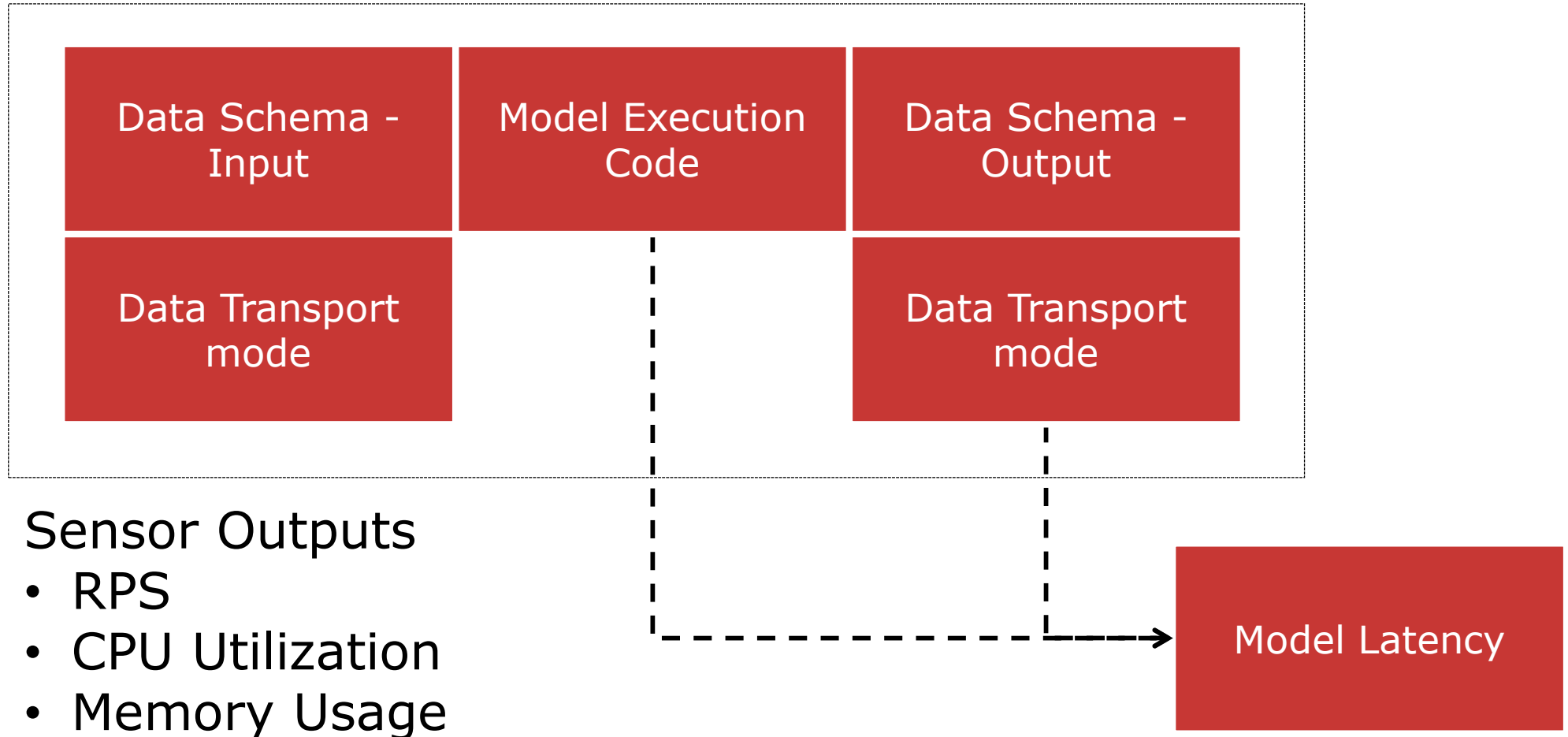
open data

# Core Concepts to Enable MDLC

| Data Schema - Input | Model Execution Code | Data Schema - Output |
|---|---|---|
| Data Transport mode | | Data Transport mode |

## The Unit of Execution for a Model

# Model Telemetry

| Data Schema - Input | Model Execution Code | Data Schema - Output |
| --- | --- | --- |
| Data Transport mode | | Data Transport mode |

Model Operations
- Model Drift
- Data Relationships
- Final Acceptance Threshold
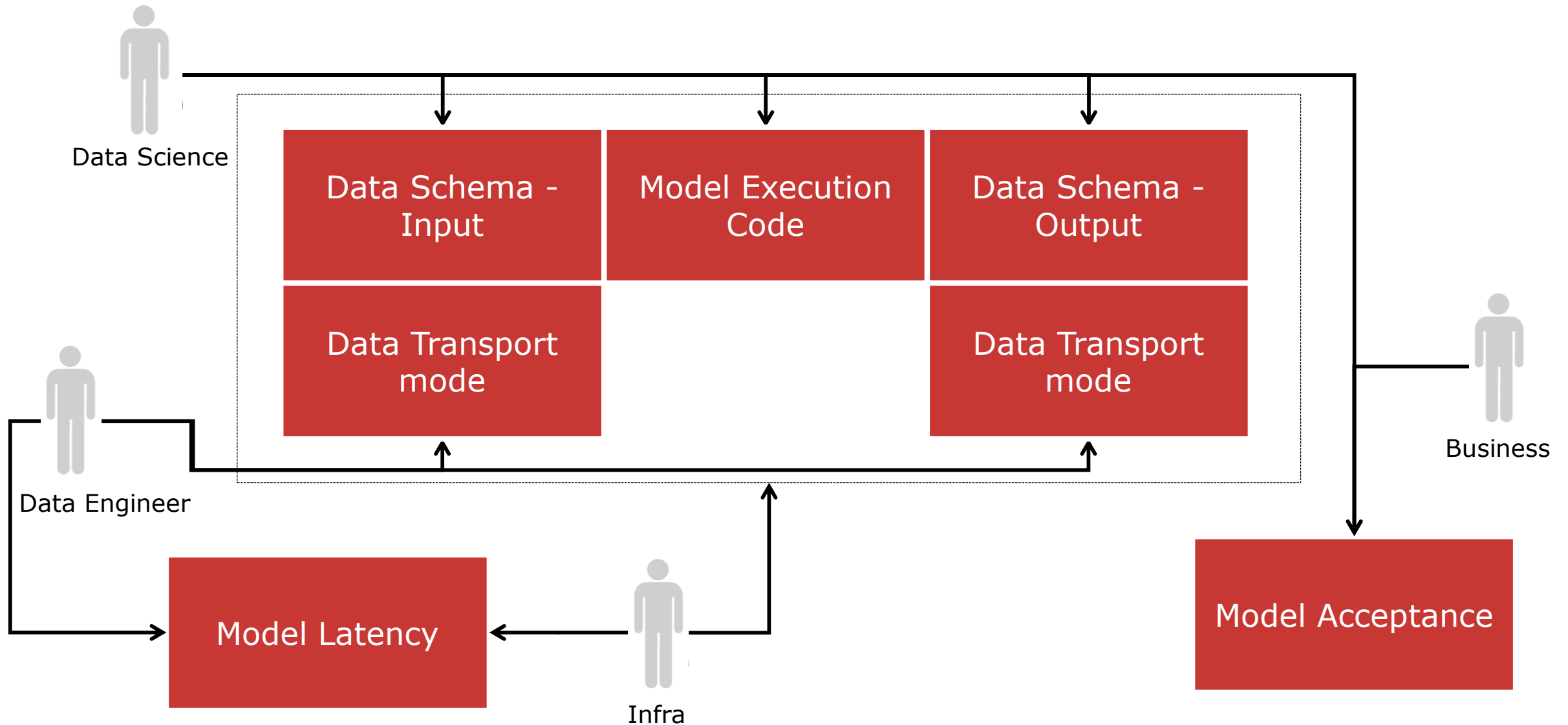
Model Acceptance

open data

# Model Telemetry

# Creating Meaningful Abstractions to Supplement the Process

| Abstraction | Description |
| --- | --- |
| Model Artifacts – Coefficients | Referenceable statistics on trained data |
| Execution Ready Model | Scoring Code |
| Model Data Definitions | Avro Schemas |
| Data Transport Descriptors | JSON file with necessary information |
| Model Environment Requirements | Mode library/code dependencies |
| Performance Monitoring Models (Telemetry) | Models custom to report performance statistics |
| Sensors | Computational performance tracking |

# Who Owns What?

# Model Development Life Cycle

# Model Prep

What assumptions are made during creation?

- Model execution environment will be consistent

- Model will receive and produce the same type of data

- Model will perform to SLAs
  - Latency
  - Accuracy (or other performance metric)

# Model Testing

- Unit Testing – does the model execution properly
- End-to-End Testing - does the analytic workflow execute properly
- Acceptance Testing – does the model perform to SLAs
- Comparison Testing – A/B Testing

open data

# Model Deployment

Pre-Production (Technical Production)

- Production data
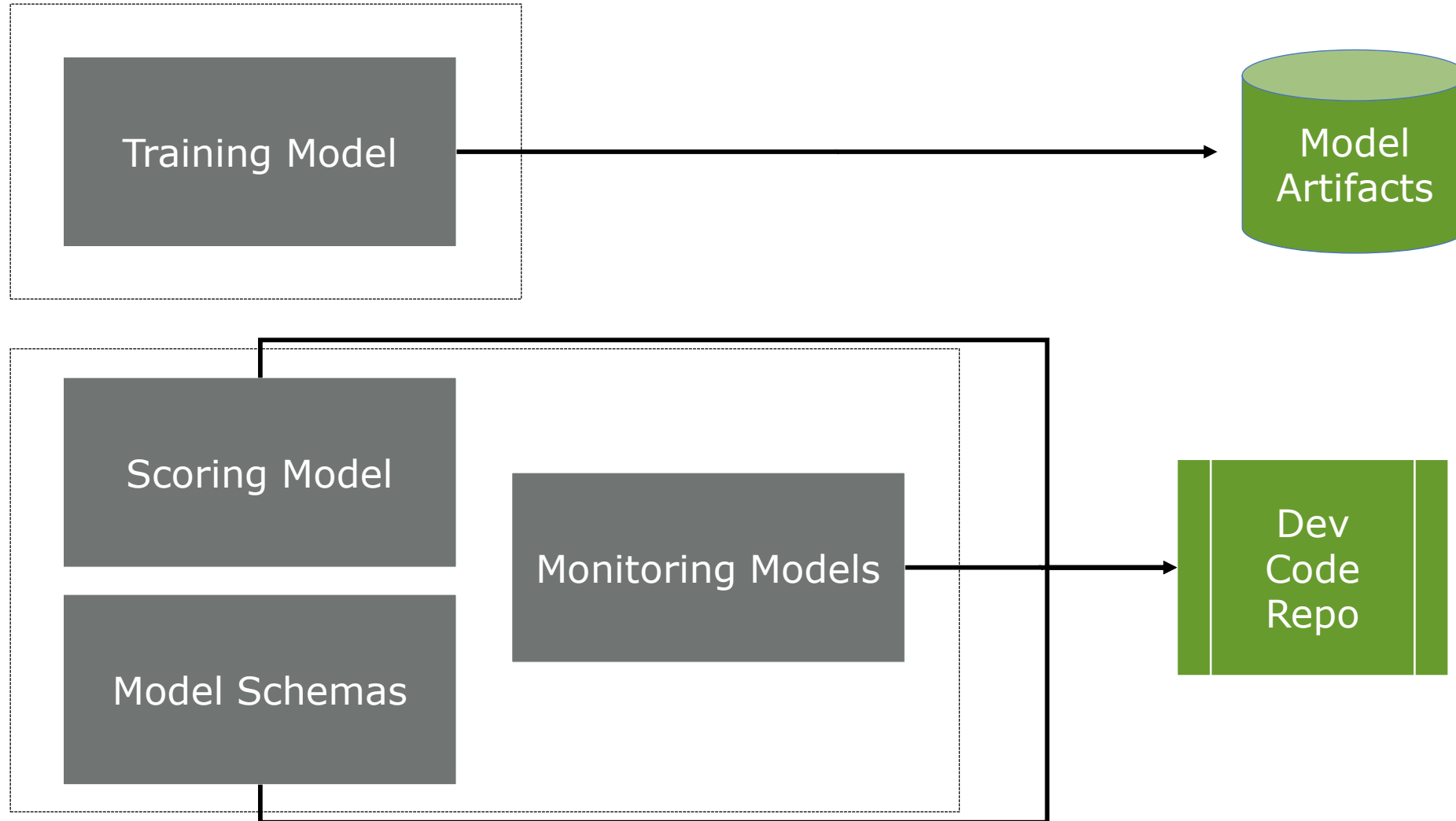
- Production Infrastructure

- Outputs to review

Production (Business Production)

- Production data

- Production Infrastructure

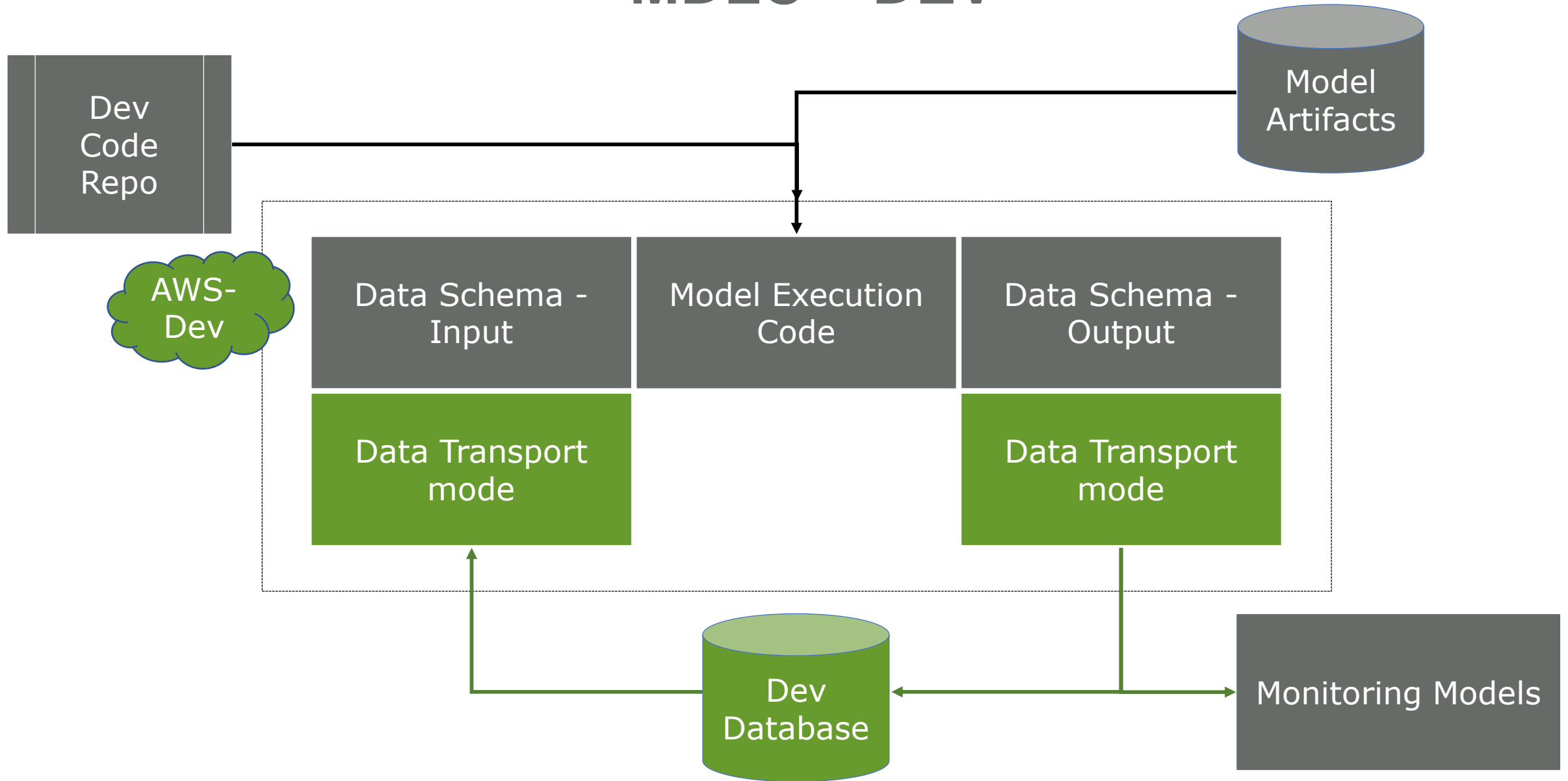- Outputs to production systems/operations for usage

# Model Performance Monitoring

- What aspects of the model are important to watch?
- What statistic identifies the quality of outputs?
- What is the calculated limit that would describe an "unacceptable" output?
- How do you define an "unacceptable" output?
- How many times does the model reach the threshold before the limit is hit?
- What actions are taken when the limit is hit?
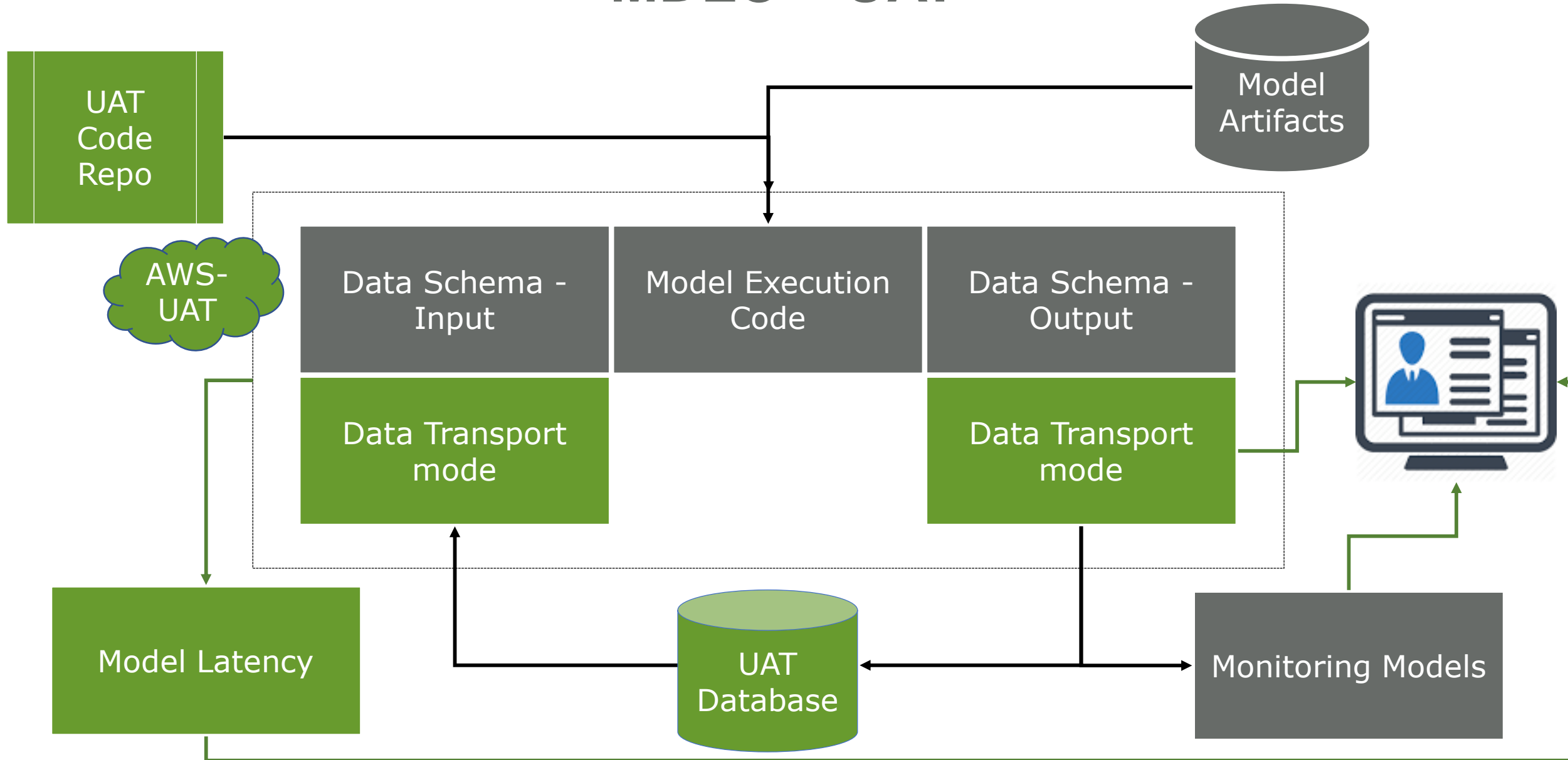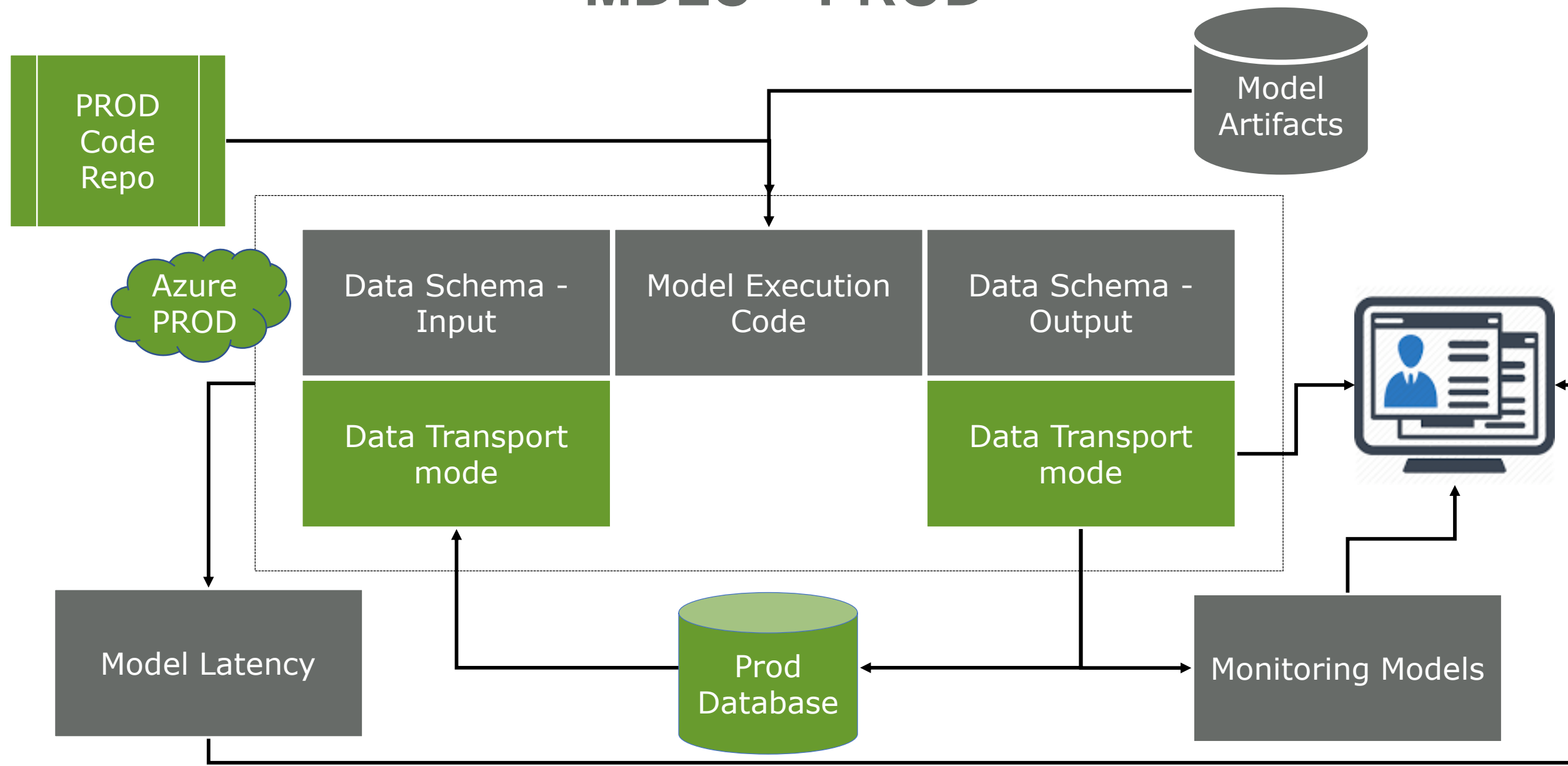- Who is responsible for taking action?

# MDLC – LAB DEV

# MDLC – DEV

open data
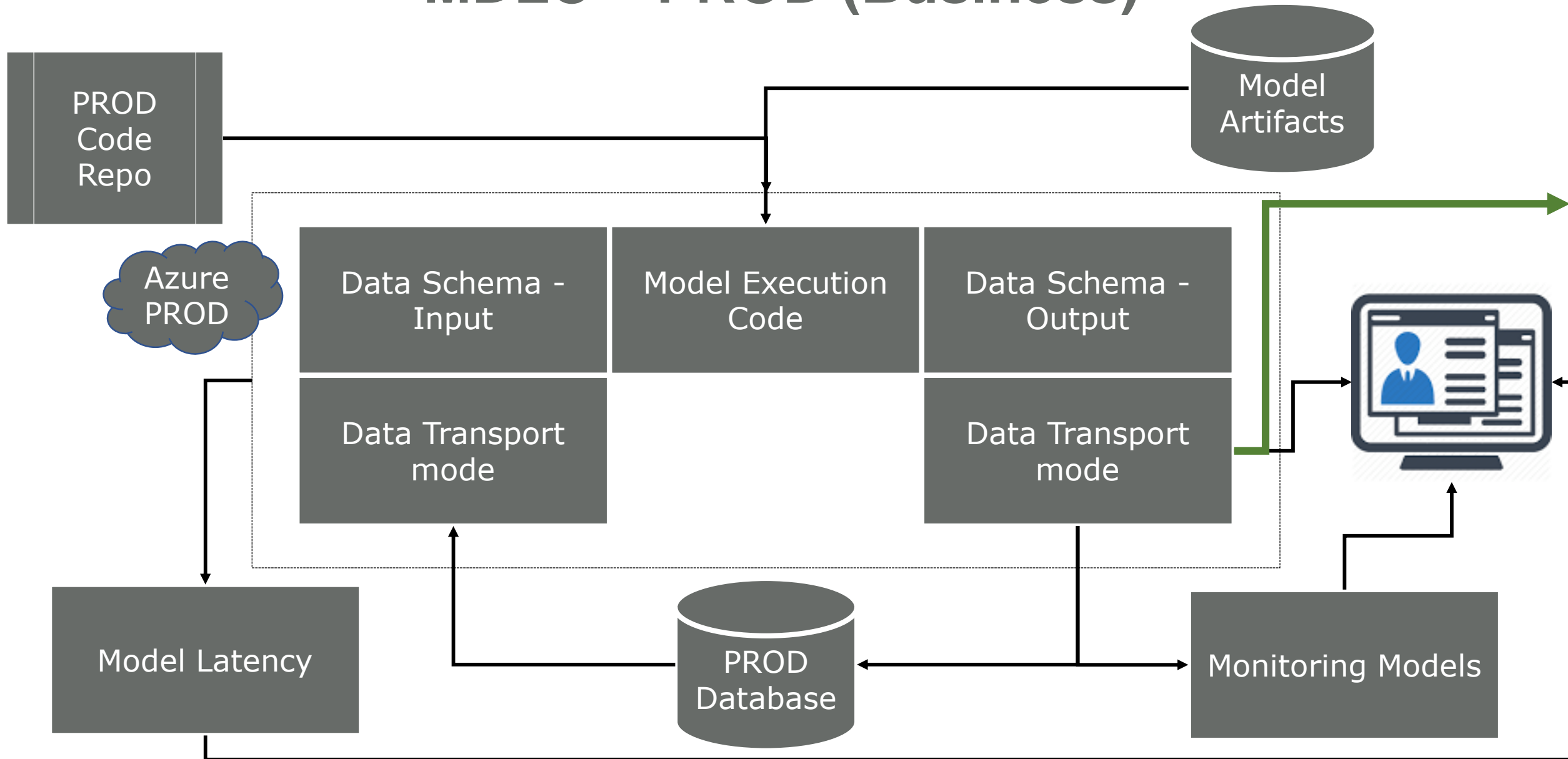
# MDLC – UAT



UAT Code Repo

Model Artifacts

AWS-UAT

| Data Schema - Input | Model Execution Code | Data Schema - Output |

Data Transport mode

Data Transport mode

Model Latency

UAT Database

Monitoring Models

# MDLC – PROD

# MDLC – PROD (Business)



PROD Code Repo

Model Artifacts

Azure PROD

Data Schema - Input

Model Execution Code

Data Schema - Output

Data Transport mode

Data Transport mode

Model Latency

PROD Database

Monitoring Models

**Goal:  Build for high quality improvements, quickly**

# Thank You