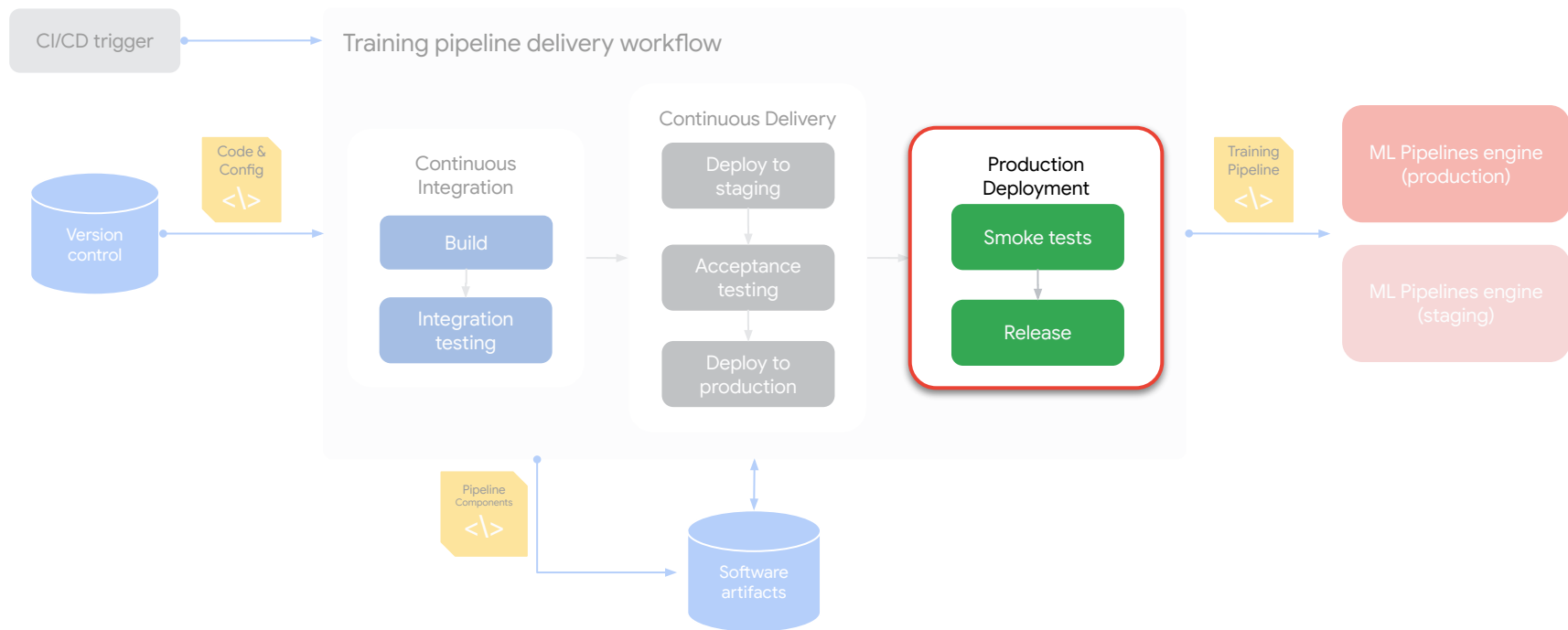# Production Deployment

# **MLOps:** Training Operationalization

# The MLOps **Personas**

ML
Engineer

ML
Researcher

Data
Scientist

Data
Engineer

Software
Engineer

**DevOps**

Business
Analyst

# Things go wrong even *after* testing?

# Things go wrong even *after* testing?

1. **Assumed** testing environments is the **same** as the production environment

# Things go wrong even *after* testing?

1. **Assumed** testing environments is the **same** as the production environment

2. **Hard to** make the test environment **match** the production environment

# Things go wrong even *after* testing?

1. **Assumed** testing environments is the **same** as the production environment

2. **Hard to** make the test environment **match** the production environment
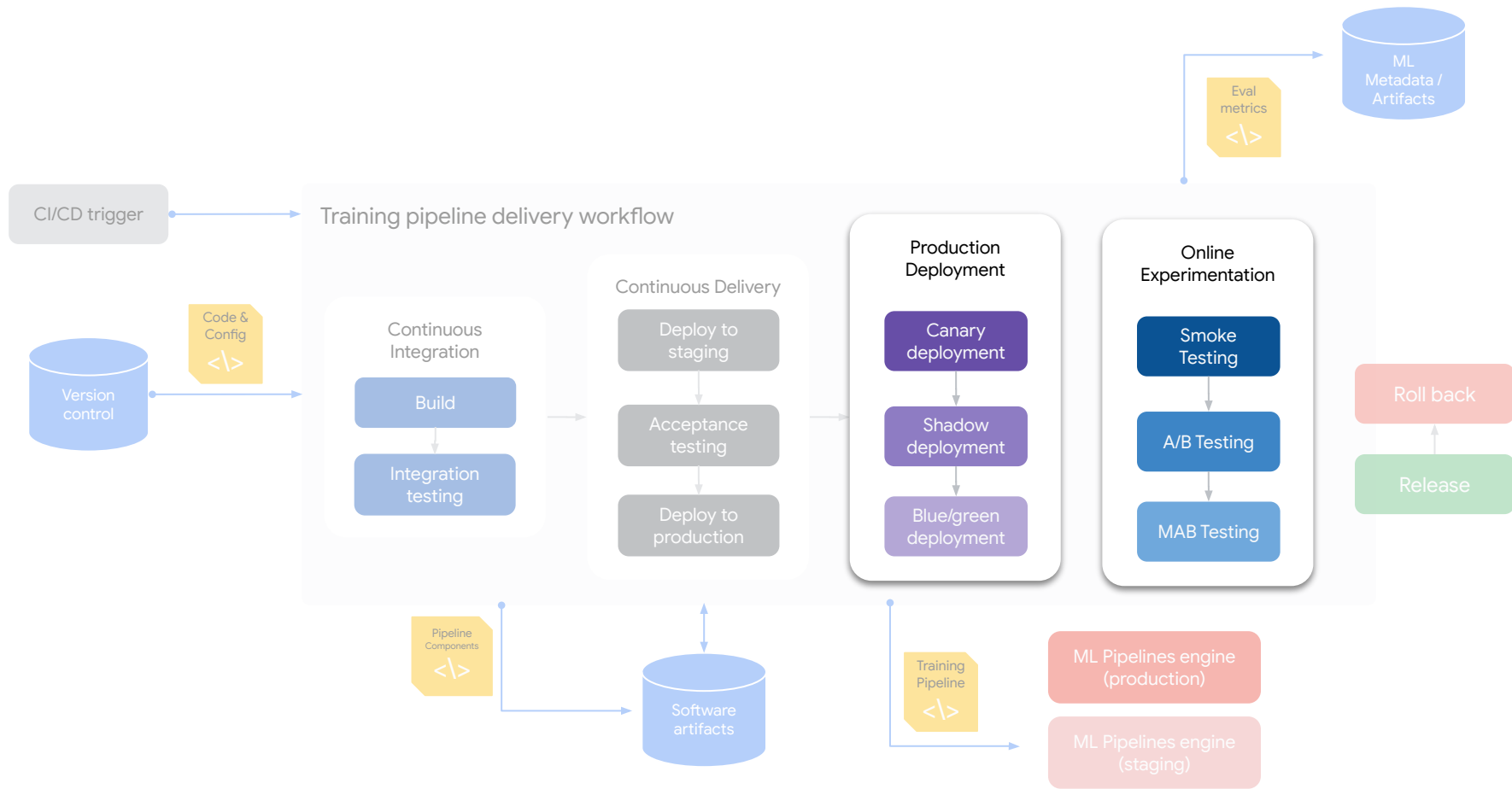
3. **Impossible to test all** the different **conditions** and scenarios in practice
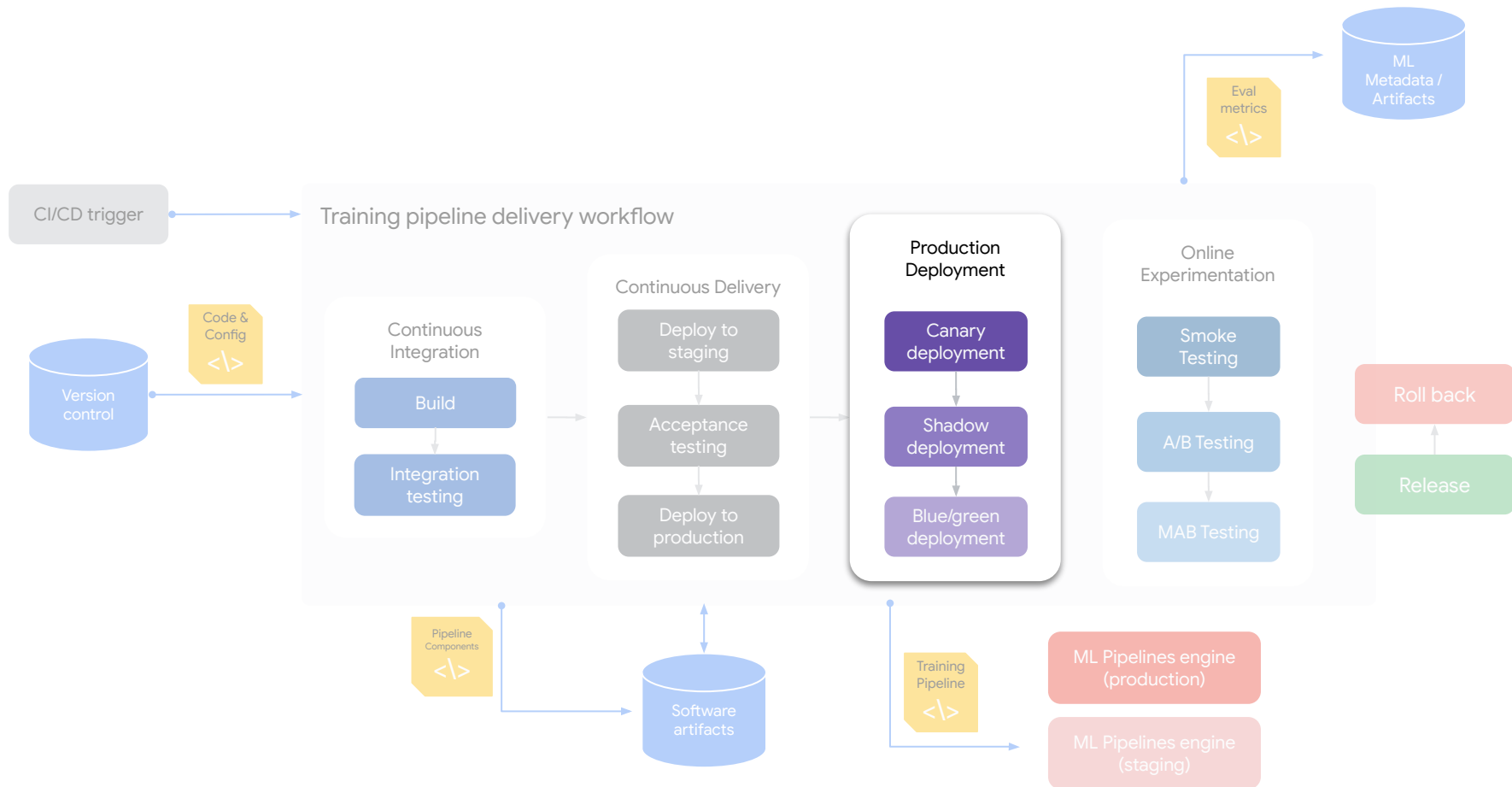
# Things go wrong even *after* testing?

1. **Assumed** testing environments is the **same** as the production environment

2. **Hard to** make the test environment **match** the production environment

3. **Impossible to test all** the different **conditions** and scenarios in practice
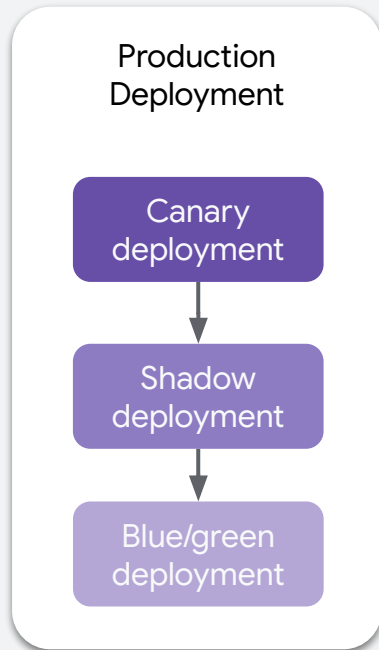
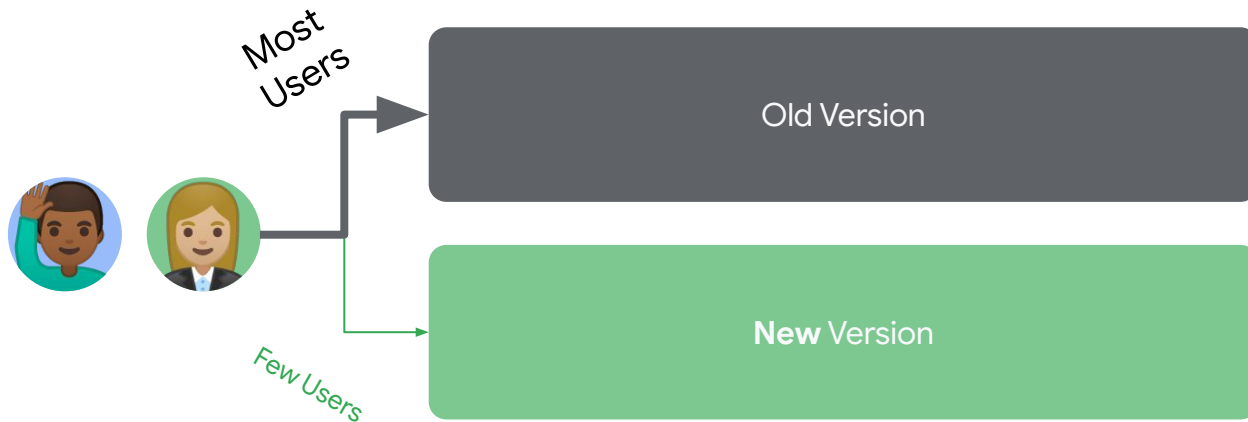⇒ *"Test" live in production environment with real-traffic*

ML Metadata / Artifacts

Eval metrics
</>

CI/CD trigger

**Training pipeline delivery workflow**

Code & Config
</>

Version control

**Continuous Integration**

Build

Integration testing

**Continuous Delivery**

Deploy to staging

Acceptance testing

Deploy to production

**Production Deployment**

Canary deployment

Shadow deployment

Blue/green deployment

**Online Experimentation**

Smoke Testing

A/B Testing

MAB Testing

Roll back

Release

Pipeline Components
</>

Software artifacts

Training Pipeline
</>

ML Pipelines engine (production)

ML Pipelines engine (staging)

CI/CD trigger

Training pipeline delivery workflow

ML Metadata / Artifacts

Eval metrics
</>

Code & Config
</>

Version control

Continuous Integration

Build

Integration testing

Continuous Delivery

Deploy to staging

Acceptance testing

Deploy to production

Production Deployment

Canary deployment

Shadow deployment

Blue/green deployment

Online Experimentation

Smoke Testing

A/B Testing

MAB Testing

Roll back

Release

Pipeline Components
</>

Software artifacts

Training Pipeline
</>

ML Pipelines engine (production)

ML Pipelines engine (staging)

# Production Deployment Strategies

1. Canary Deployment
2. Shadow Deployment
3. Blue/Green Development

Production Deployment

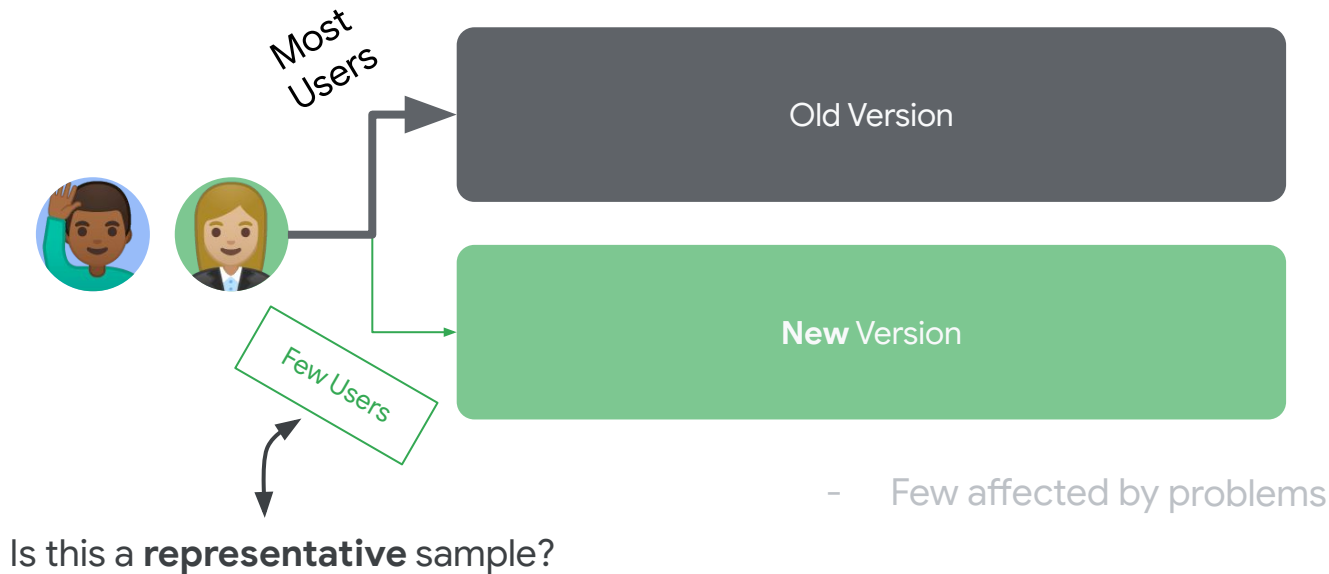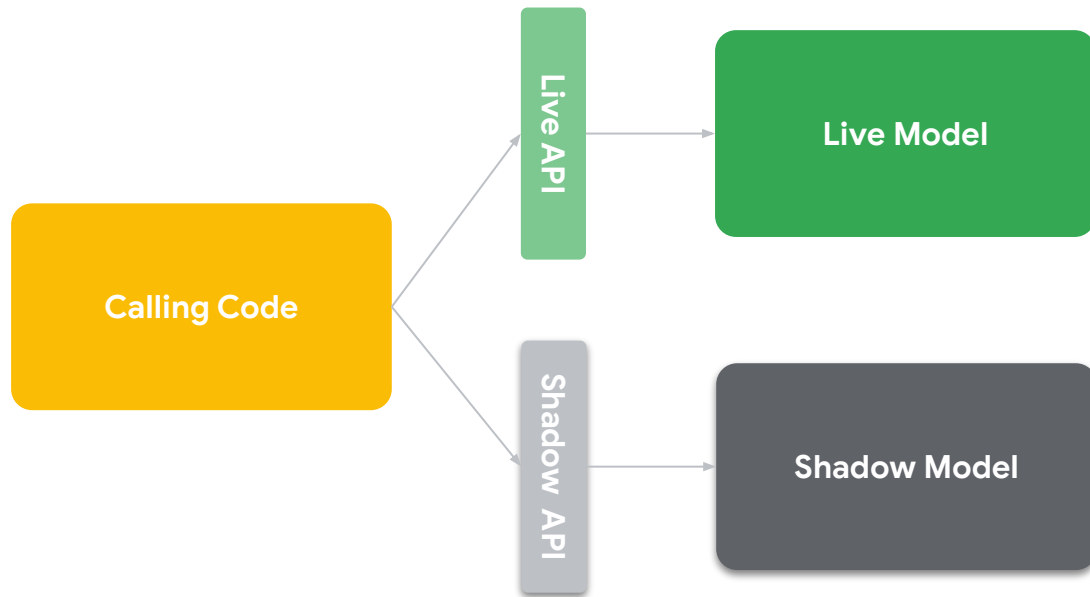Canary deployment

Shadow deployment

Blue/green deployment

# Canary Deployment

# Canary Deployment

Most Users

Old Version

New Version

Few Users

- Few affected by problems

# Canary Deployment



Most
Users

Old Version

Few Users

**New** Version

- Few affected by problems
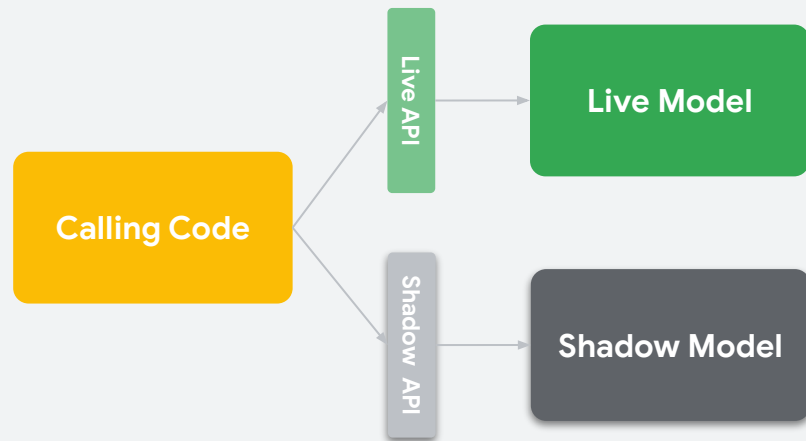
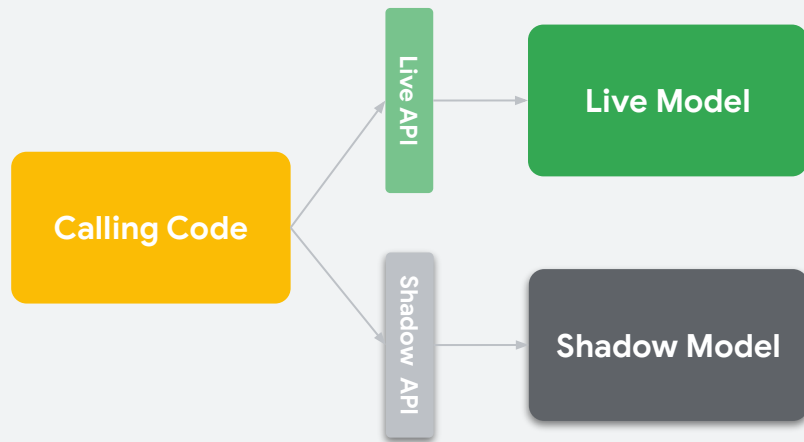Is this a **representative** sample?

# Shadow Deployment

# Shadow Deployment

+ **No impact to production** traffic (customer)
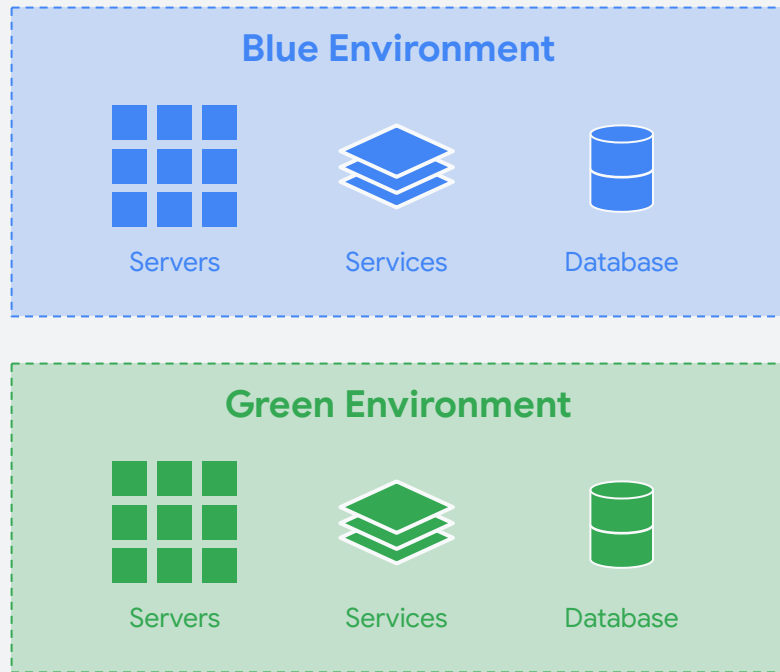
+ **Check stability** (requirements)

# Shadow Deployment

\+   **No impact to production** traffic
(customer)

\+   **Check stability**
(requirements)

\-   Increased **cost** (parallel)

\-   Complex **update-in-place**

\-   Need to carefully design
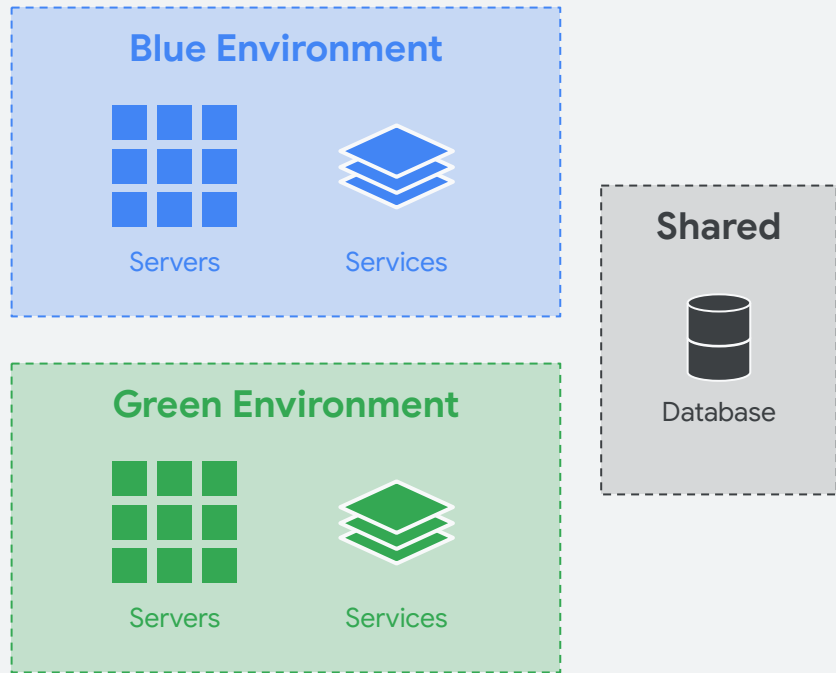**performance testing**

# Blue/Green Deployment

In its purest form, blue-green asks us to **duplicate** every resource our application leverages.



**Blue Environment**

Servers      Services      Database

**Green Environment**
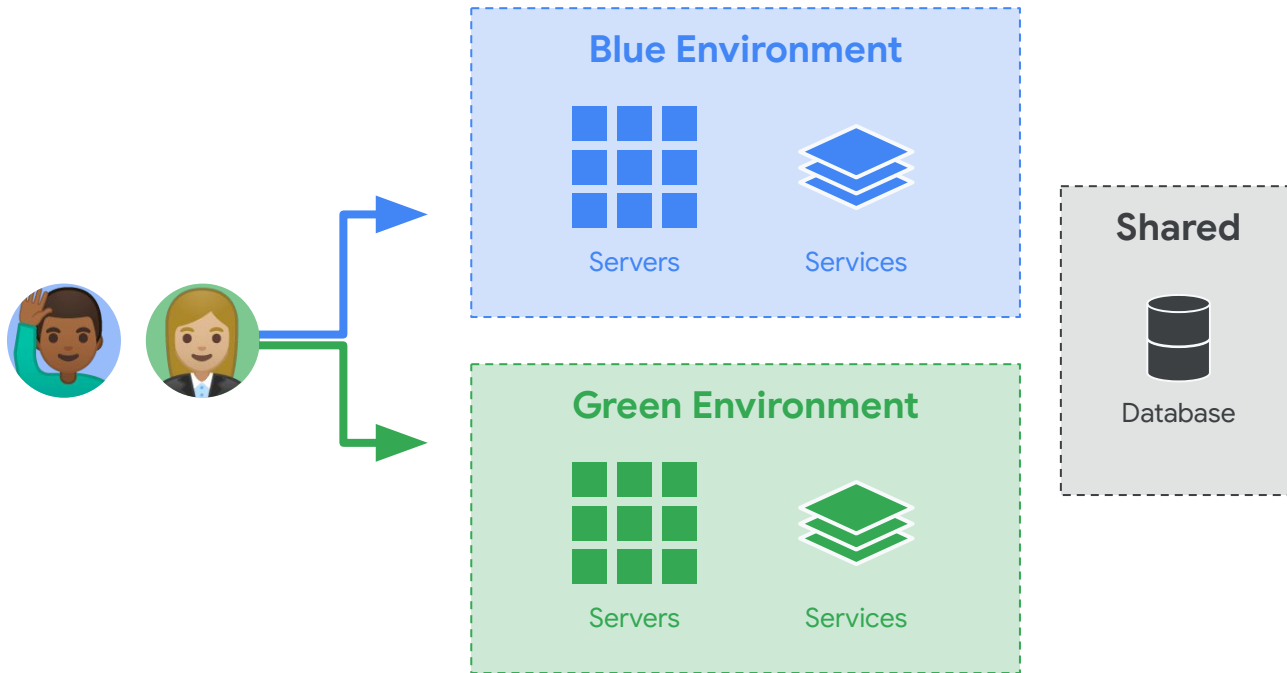
Servers      Services      Database

# Blue/Green Deployment

In its purest form, blue-green asks us to **duplicate** every resource our application leverages.

But, we often see blue-green deployments with **shared** components.

**Blue Environment**

Servers

Services

**Green Environment**

Servers

Services
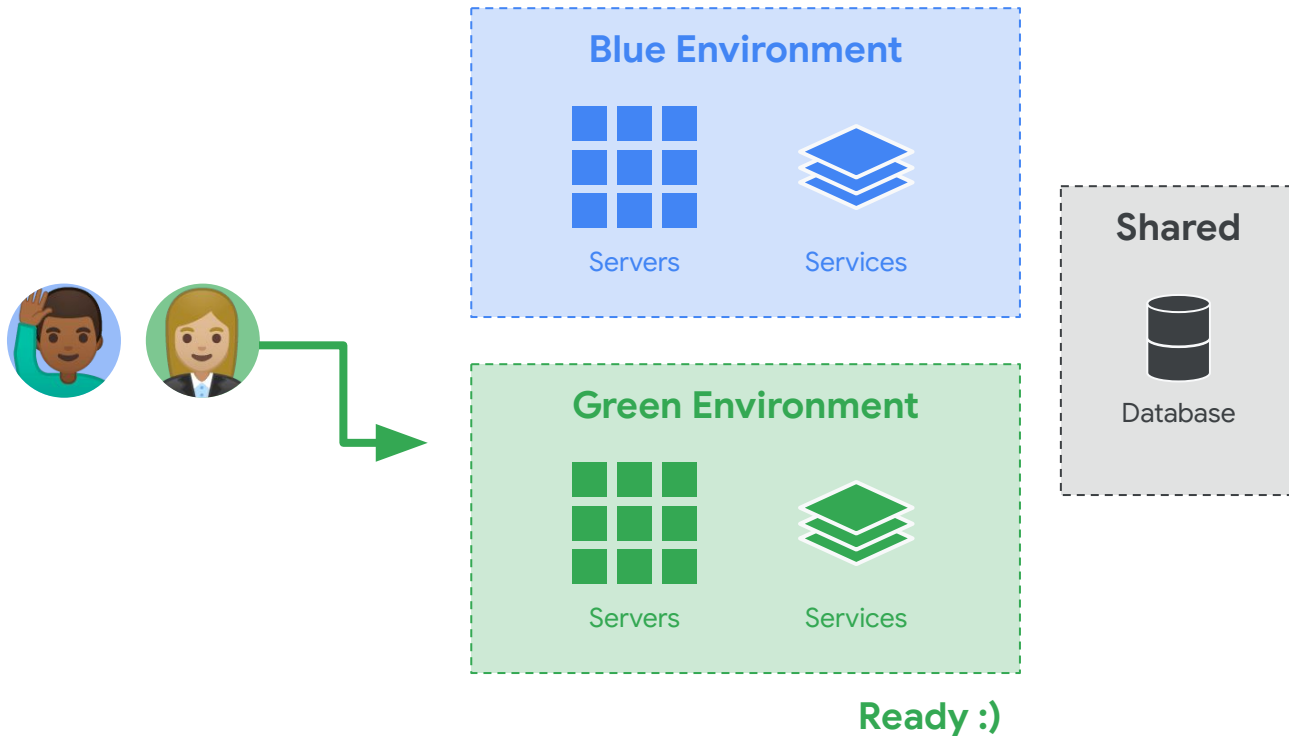
**Shared**

Database

# Blue/Green Deployment

# Blue/Green Deployment:
# **Alternate Producers**

# Blue/Green Deployment:
# **Alternate Producers**

# Blue/Green Deployment

+ No system downtime

+ Accurate testing

+ Reliability testing of deployments



**Blue Environment**

Servers    Services

**Green Environment**

Servers    Services

**Shared**

Database

# Blue/Green Deployment

For **TinyML**:

+ **When we can't make continuous updates**

**Blue Environment**

Servers

Services

**Green Environment**

Servers

Services

**Shared**

Database

# Blue/Green Deployment

For **TinyML**:

+   When we can't make continuous updates

**+   When we can't support two identical environments**

**Blue Environment**

Servers          Services

**Shared**

Database

**Green Environment**

Servers          Services

# Blue/Green Deployment

For **TinyML**:

+ When we can't make continuous updates

+ When we can't support two identical environments

+ **When infrastructure does not allow a router to direct users**

## Blue Environment

Servers          Services

## Green Environment

Servers          Services

## Shared

Database

ML Metadata / Artifacts

Eval metrics
</\>

CI/CD trigger

Training pipeline delivery workflow

Production Deployment

Online Experimentation

Code & Config
</\>

Continuous Integration

Continuous Delivery

Deploy to staging

Acceptance testing

Deploy to production

Canary deployment

Shadow deployment

Blue/green deployment

Smoke Testing

A/B Testing

MAB Testing

Roll back

Release

Version control

Build

Integration testing

Pipeline Components
</\>

Software artifacts

Training Pipeline
</\>

ML Pipelines engine (production)

ML Pipelines engine (staging)