# eppo

# Building A Modern Experimentation Stack

Every ambitious AB testing system has the same architecture, with seven core components. For an experimentation program to generate ROI, they all have to work together. Let's walk through each of these components, and the opportunities and failure modes they present us with.

In a previous era, data teams that wanted to run experiments had to rely on an open-source distributed compute platform like Spark or Hive. This forced teams to build bespoke solutions, since these systems didn't have a canonical method of implementation.

# Nowadays, organizations can build upon a Modern Data Stack, which includes a data warehouse like Snowflake, Databricks, Redshift, or BigQuery, and a transformation layer like dbt or Airflow.

These new tools enable companies to run experiments early and often.

# What is an experiment?

An experiment is a tool for understanding cause and effect. By controlling inputs to a set of experiment subjects and measuring outputs, we gain knowledge about the input/output relationship. We quantify that knowledge using statistics.

### Inputs

**The inputs to experiments can include things like:**

- A new algorithm for recommendations, search, or matching
- A new product feature (or the choice to remove a feature)
- Marketing copy, design alternatives
- Anything you want to test!

### Subjects

The subjects of experiments are generally users who have a User ID and/or a Session ID. Some of these users see version A, and others see version B, randomly, and we compare the two groups.

### Outputs

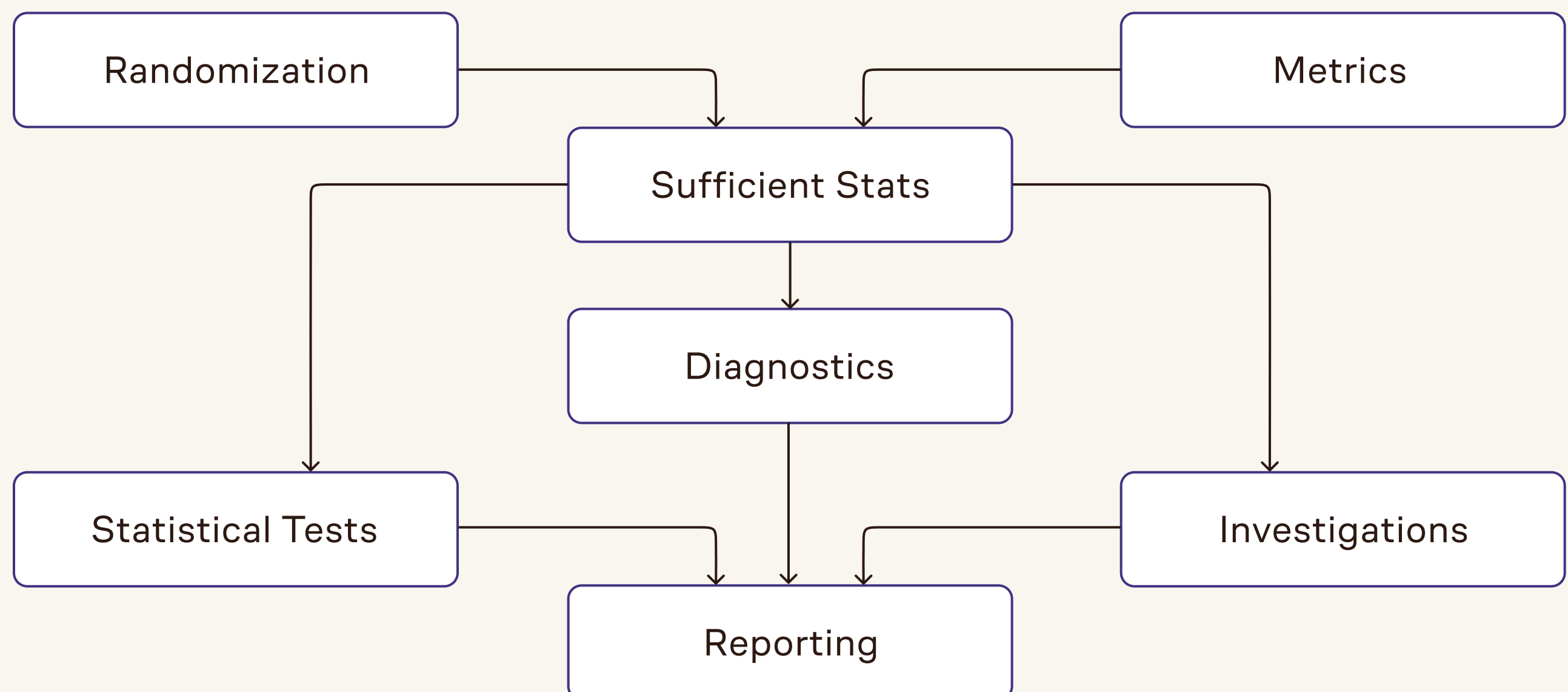**The outputs of experiments will include things like:**

- Conversions (binary outcome)
- Revenue (continuous outcome)
- Any measurable subject behavior

There's extensive evidence that we are all bad at predicting whether new features improve or degrade performance; in fact only about 1/3rd of features lead to positive outcomes. Experimentation is the best tool to efficiently figure out what works, and what does not.

**Companies that have realized this, like Airbnb, Microsoft, Meta, and Netflix, built big in-house experimentation platforms, using their large-scale engineering resources. And every one of these systems has the same architecture.**
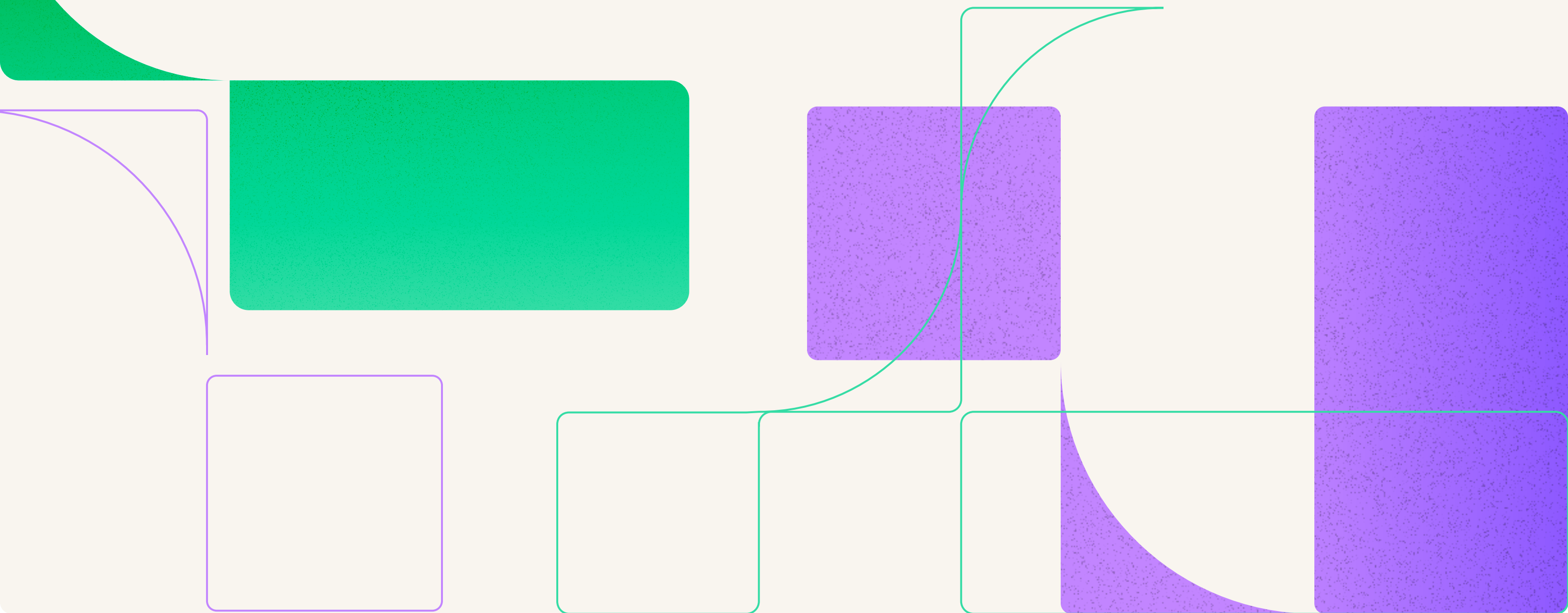
Each of these components has depth and complexity, but for experimentation to generate meaningful ROI, they all have to work together. When even a single component isn't working, it can really set your experimentation practice back.

Let's walk through each of these components, and the opportunities and failure modes they present us with.

```
┌──────────────────┐                    ┌──────────────────┐
│  Randomization   │                    │     Metrics      │
└──────────────────┘                    └──────────────────┘
            │                              │
            └──────────┐        ┌──────────┘
                  ┌──────────────────┐
                  │ Sufficient Stats │
                  └──────────────────┘
        │                 │                 │
        │         ┌──────────────────┐      │
        │         │   Diagnostics    │      │
        │         └──────────────────┘      │
        │                 │                 │
┌──────────────────┐      │      ┌──────────────────┐
│ Statistical Tests│      │      │  Investigations  │
└──────────────────┘      │      └──────────────────┘
        │                 │                 │
        └──────┐   ┌──────────────────┐  ┌──┘
               │   │    Reporting     │  │
               └───└──────────────────┘──┘
```

Every experimentation system has the same architecture
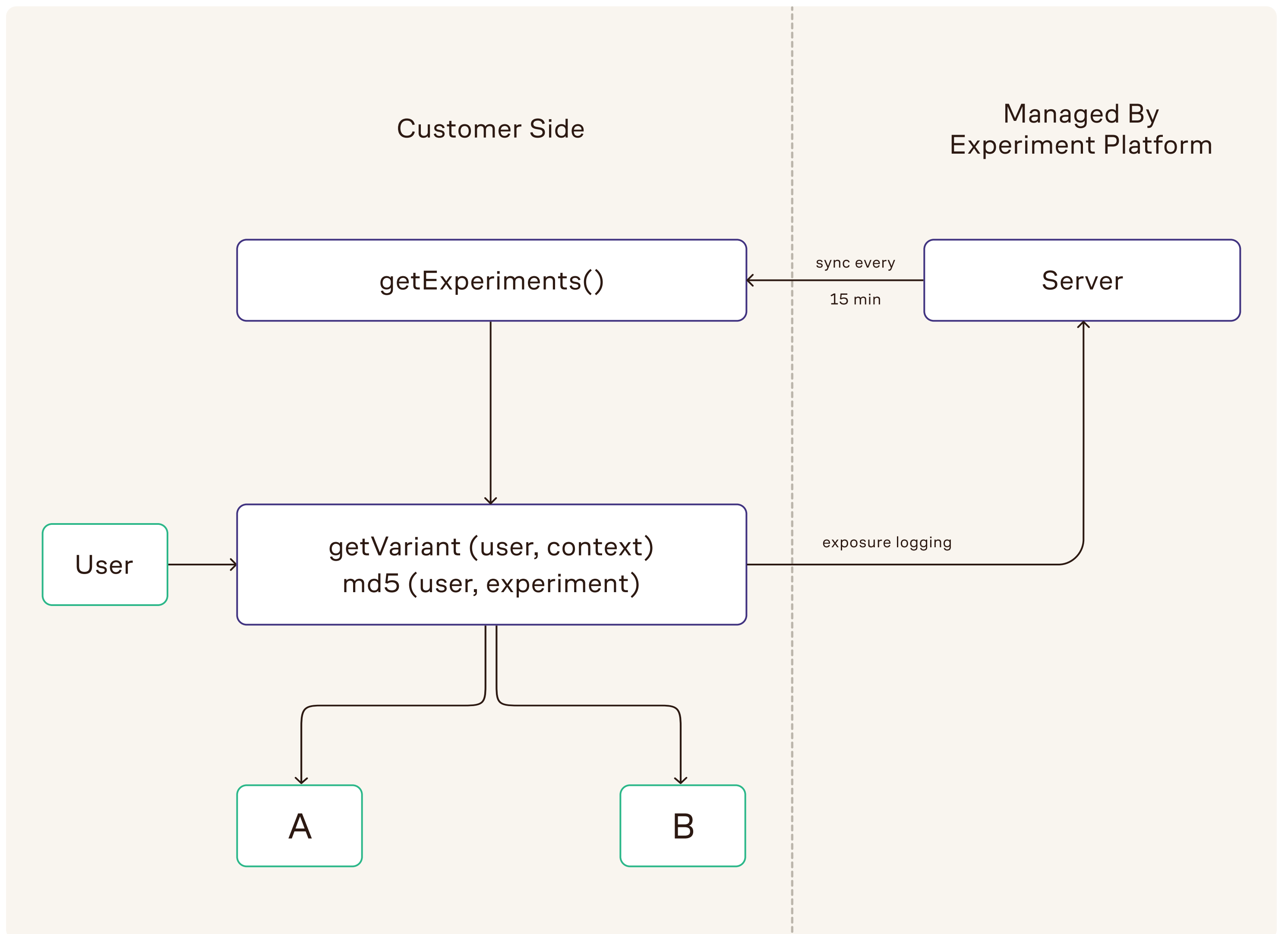
# Randomization

To most people, randomization is probably the most familiar component of an experimentation workflow. It's a process that selects users randomly and independently from any attributes they possess.

There are **three problems** to solve here:

1. How do you serve experiments quickly?
2. How do you randomize properly?
3. How do you instrument what you did?

**Under the hood, all commercial feature-flag tooling looks essentially the same. They all have a way to convey experiment configurations that are synced onto clients every 15 minutes or so.**

Customer Side

Managed By
Experiment Platform

```
getExperiments()
```

sync every
15 min

Server

User → getVariant (user, context)
md5 (user, experiment)

exposure logging

A          B

Pro Tip: Use md5() hashes for assignment

With experiment metadata stored locally, it becomes a quick, low-latency process to split users into different groups. Typically, they have a function called "getVariant," which sends an event back to your warehouse, identifying that one specific user was assigned to one specific experiment.

For randomization, the best practice is to use md5() hashes. Hashing is important, because whether you conduct your randomization offline or online, you can precisely replicate the process. It's idempotent — If the user hits the experiment 10 times, they'll end up in the same treatment group each time.

# Metrics

In experimentation, a metric is a measurement of subject-level activity that occurs after assignment. Metrics are created by aggregating (e.g. summing, counting, etc.) events (facts).

## It's fair to estimate that 80% of a data team's work is to define and serve high-quality metrics to different places, but this is something that a lot of experimentation programs get wrong.
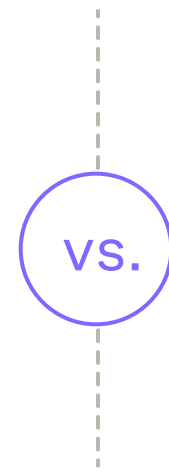
The first thing to emphasize is that in experimentation, you want to define metrics by their underlying facts. A fact table enables you to model atomic actions that happen at a specific moment in time. For experimentation, you need a timestamp to know if a purchase or upgrade happened before or after the user was in the experiment.

Once you turn to the fact table, the best way to make an experimentation program more powerful is to use business metrics. Too many companies just use the metrics that are available via existing tooling, which tends toward shallow metrics like click-conversions. You can boost those metrics, but what you really want to show is that you boosted revenue. Ask yourself the question: If you show your experiment result to your CFO, will they care?

The biggest gap between Airbnb/Netflix/et al. and commercial tools is how easily you can use **business metrics.**

| **Business Metrics** | | **Shallow Metrics** |
|---|---|---|
| • Revenue, Activation, Purchases | | • Signups, "conversions" |
| • What the CFO reads | vs. | • "Directionally accurate" |
| • From databases, Stripe, multiple POS | | • From event streams |

To get real value from experimentation, you want to use financial data, such as from Stripe, and you want to be connected to the data warehouse.

Because teams are so used to shallow click metrics, they usually focus on steps higher up in the funnel, like signups or form submissions.

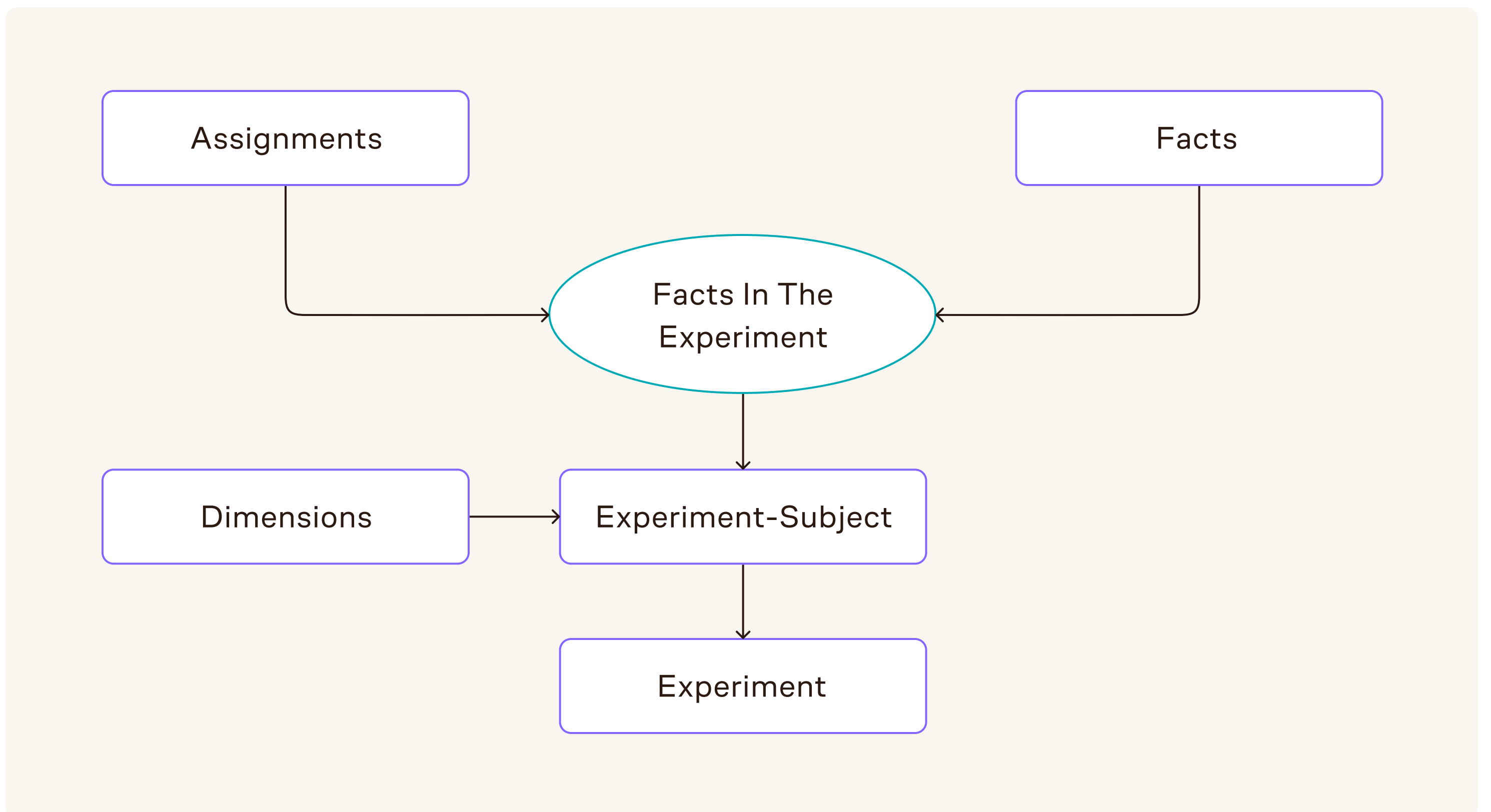| Metric | Δ | p-value |
|---|---|---|
| Search to Book | -0.31% | 0.37 |
| Search to Contact | -1.29% | 0.04 |
| Contact to Book | 0.99% | 0.06 |
| Contact to Accept | 1.58% | 0.00 |
| Accept to Book | 0.58% | 0.11 |

Search → Contact → Accept → Book

Most enterprise experimentation teams find that the vast majority of experiments will boost one part of the funnel, and then depress another, leaving your results completely neutral. If you use business metrics, you'll have more confidence that you're not just playing whack-a-mole.

# Sufficient Statistics

This component is called sufficient statistics because that's what you're calculating to fuel your inference, but it's actually a significant data engineering problem.

Experiment Data Pipelines: One big JOIN and some GROUP-BYs

This is the DAG that most pipelines end up working with for experimentation. On one end, you have the assignment events from the randomization system. On the other end, you have facts, such as a purchase. And then you have to join them together, to determine: of these purchases, how many ended up in this experiment? You group it first at an experiment-subject level, and then at an experiment level.

In this DAG, all the action is primarily in that first join. It's basically a big join and then two group-bys. The group-bys happen very quickly, and most data engineering teams don't use much compute on them. But that experiment join is so large that it can make up a third of all your compute costs, and it's where most of the optimizations will happen.

# Statistical Tests

While there are plenty of conversations around the relative benefits of Frequentist vs. Bayesian statistics, 95% of experimentation programs just run a t-test. While this isn't necessarily wrong, it puts a lot of the onus for expertise on your organization, because there's a bunch of failure modes for t-tests.

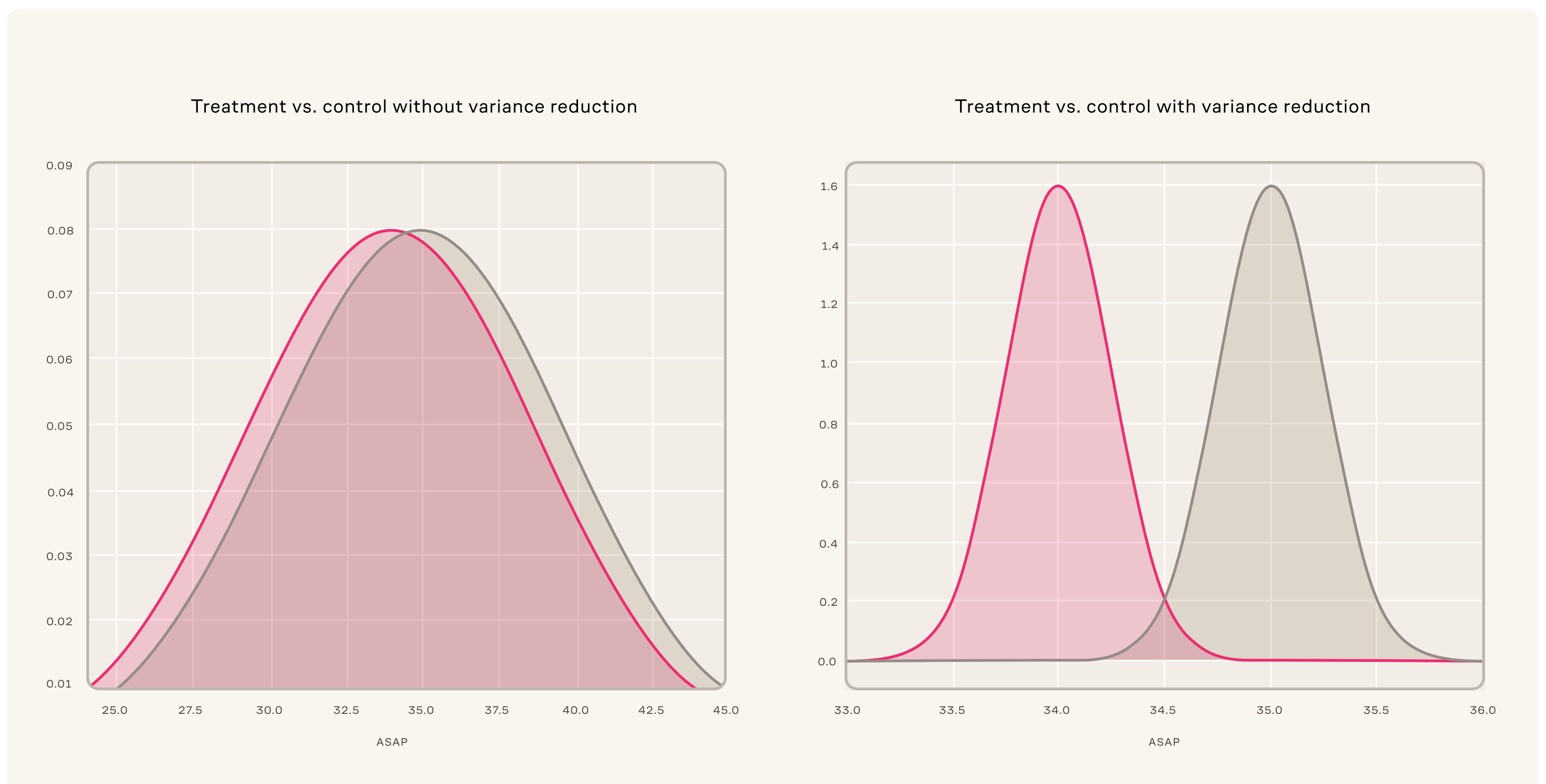## Simple statistical tests stress the organization

**When using t-tests, you need to:**

- Not look at the results until it's done
- Not test multiple variants without a statistical correction
- Not have outliers/power laws

Specifically, you're not supposed to look at the results until they're done. You make a promise saying that you're going to run an experiment, get 1000 samples in both groups, and then look at the results and make a decision. Anyone who runs an experimentation program knows that this never works in practice — peeking is the norm.

There are a bunch of other failure modes. For example, ratio metrics require adjusting for the correlation between the numerator and denominator via the delta method to obtain valid statistical results.

# Because most teams don't realize these points of failure, they end up incurring a bunch of false positives. So simple statistical tests end up stressing the organization.



Advanced statistics can speed things up

But there's a few things that can help avoid statistical problems. To get around the peeking problem, you can use sequential methodology; these methods assume that results of an experiment are monitored during the entire runtime of an experiment. They adjust results in a way that false positives are kept low.

Another way to help your organization is through variance reduction, and a technique like CUPED. By leveraging historical data, you can make your experiments run faster. You can save weeks of product time and increase your learning rate, and there's really no downside.

## CUPED controls for variance that we can predict

### Goal:

Improve # of happy meals purchased

### Insight:

We can predict whether some will purchase a happy meal
- Are they a family?
- Have they purchased a happy meal recently?



Suppose you're McDonald's, and you're running an experiment to boost the number of Happy Meals sold. The easy way to do it is to run some experiment, then see if the purchase number goes up. But suppose every time someone walks in the door, you make a guess as to whether they'll buy a Happy Meal, and then you compare the actual result against your guesses. In that paradigm, if someone walked in with 3 children, you would guess that they'd get a Happy Meal. Suddenly, you are accounting for some natural variation, removing it from the equation, and getting a lot more signal.

**CUPED has become a mainstream method at companies like Airbnb, Microsoft, and Meta, but because it involves a little more complexity from a data engineering and statistics standpoint, you don't yet see it across the commercial landscape.**

# Diagnostics

**The underlying principle of experimentation is trust. The whole process is a leap of faith, where you let a system separate the winners from the losers, so trust is paramount.**

The most common failure mode for experimentation is imbalanced treatment groups. What you think is a 50/50 split may actually be closer to 52/48, and that difference almost always signals an underlying issue. The types of users throwing off this 50/50 split might have low latency, or might be on a broken client. The way to check for this is via the sample ratio mismatch test, which can diagnose the underlying problem. All mainstream experimentation tools build this check into their day-to-day practice.

**Problem:**

Imbalanced treatment groups

**Causes:**

These issues are usually due to
- Latency of experiment delivery
- Bad implementation

**Solution:**

Sample ratio mismatch test (SRM)

Treatment    Control          Treatment    Control

US

Biased implementation          Correct, unbiased implementation

$$x^2 = \sum \frac{(observed - expected)^2}{expected}$$
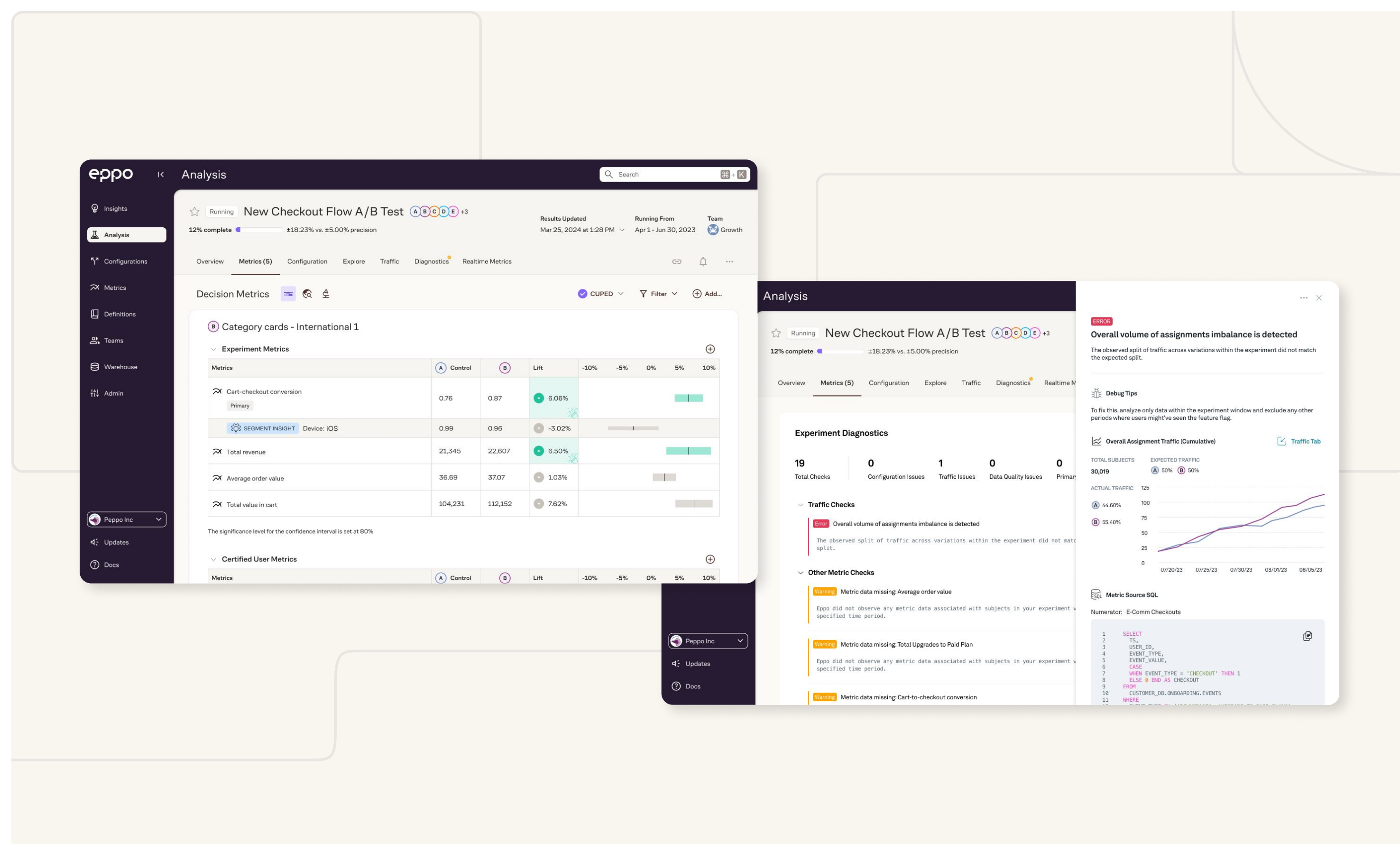
Make sure you have balanced groups

The next diagnostic is outliers. As 'The Black Swan' author Nassim Nicholas Taleb emphasizes, social phenomena metrics follow power laws, and not the finite moment distributions that all of statistics is built upon. There will always be outliers that skew your results, and a small number of observations can wreck the analysis.

To handle outliers, one can use non-parametric tests that make no assumptions about distributions. But in commercial practice, it is more common to see **two things**:

1. Winsorizations, which take the 99th percentile as the maximum for the metric.
2. And CUPED, which accounts for prior history.

# Another challenge with experimentation is that it makes your underlying data quality issues invisible.

In product analytics, it's typically clear that something is broken. If the revenue chart drops to zero, someone will ring alarm bells. However, this is often harder to spot in experiment analyses, as the control and treatment groups are often equally broken, but that does not show up in the difference. So you need a way to monitor data quality before running the experiment.



Monitor your data quality

# Investigations

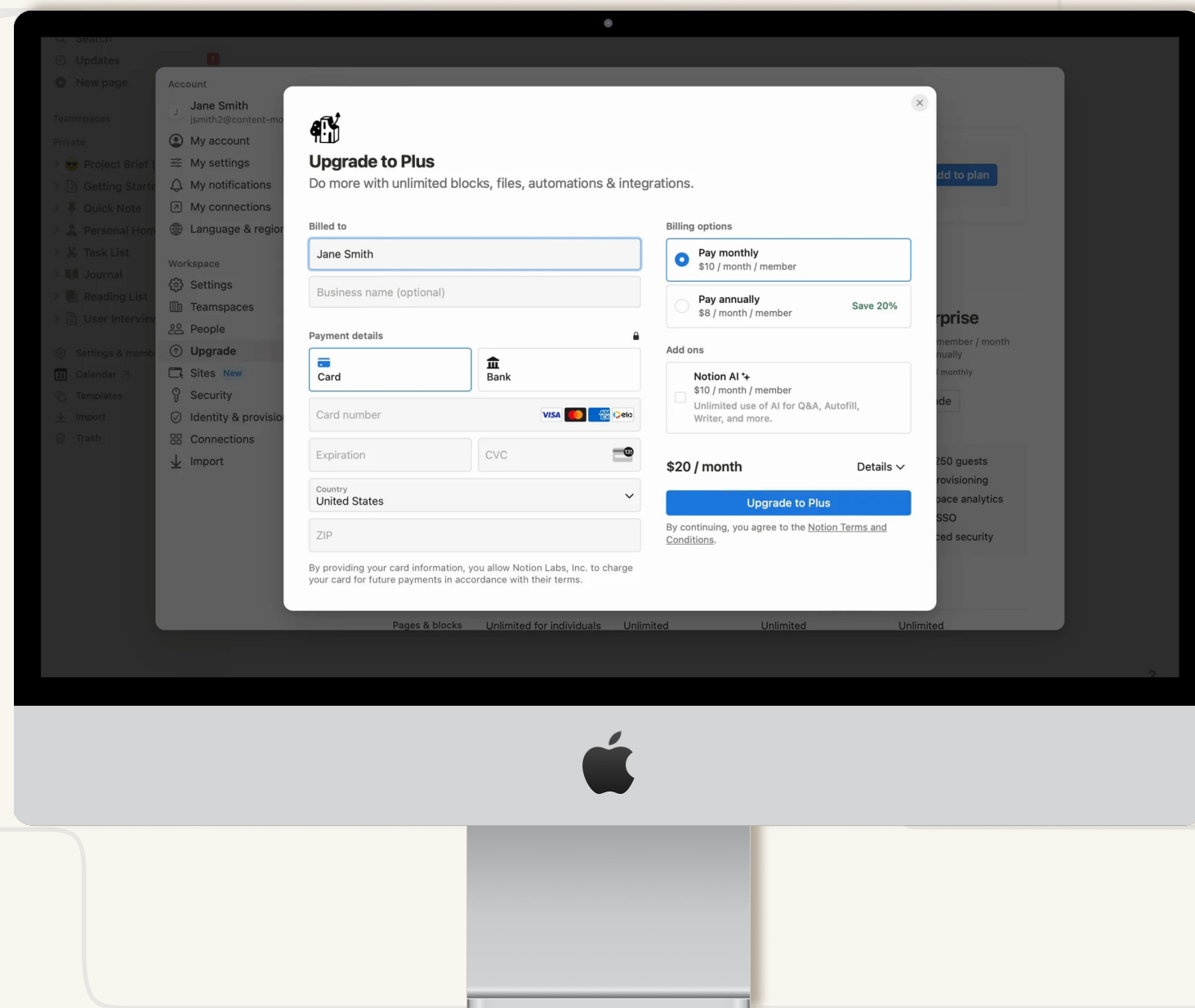Every experimentation practice discovers that delivering wins is a learning process.

**Investigations help you learn**

## "First you must learn to test. Then you learn to learn. Then you learn to win."

**Elena Verna**
Head of Growth & Data @ Dropbox

First you have to be able to run an experiment. And then you have to learn to diagnose problems quickly. And then, you can get to a point where you reliably deliver wins.

| Browser | Δ | p-value |
|---|---|---|
| **All** | -0.27% | 0.29 |
| Chrome | 2.07% | 0.01 |
| Firefox | 2.81% | 0.00 |
| IE | -3.66% | 0.00 |
| Safari | 0.86% | 0.26 |
| **Rest** | -0.74% | 0.33 |

A specific segment of your product's users might particularly dislike your experiment

The most important way to learn is by segmenting your results. In one common scenario, a team sees a result that looks negative, but once they segment it by browser, they discover that it's just a browser-specific bug, and the overall results are positive.

**You need to be able to split your results by user segment, and make the process democratized. That's the only way to reliably understand your experiment results.**

16

# Reporting

Reporting is the most underrated part of data science. You do all this work with such complexity, but in the end, your goal is to drive a decision in your organization. The last mile of data science is communication.



Bad reporting undermines all your hard work on the math and engineering behind getting results

The current state of reporting is suboptimal, and it creates a norm where data scientists on your team have to explain/translate the results in order to make them intelligible.

## Good reporting practices:

- Don't try to teach p-values or stat tests
- Don't list 100 numbers without guidance
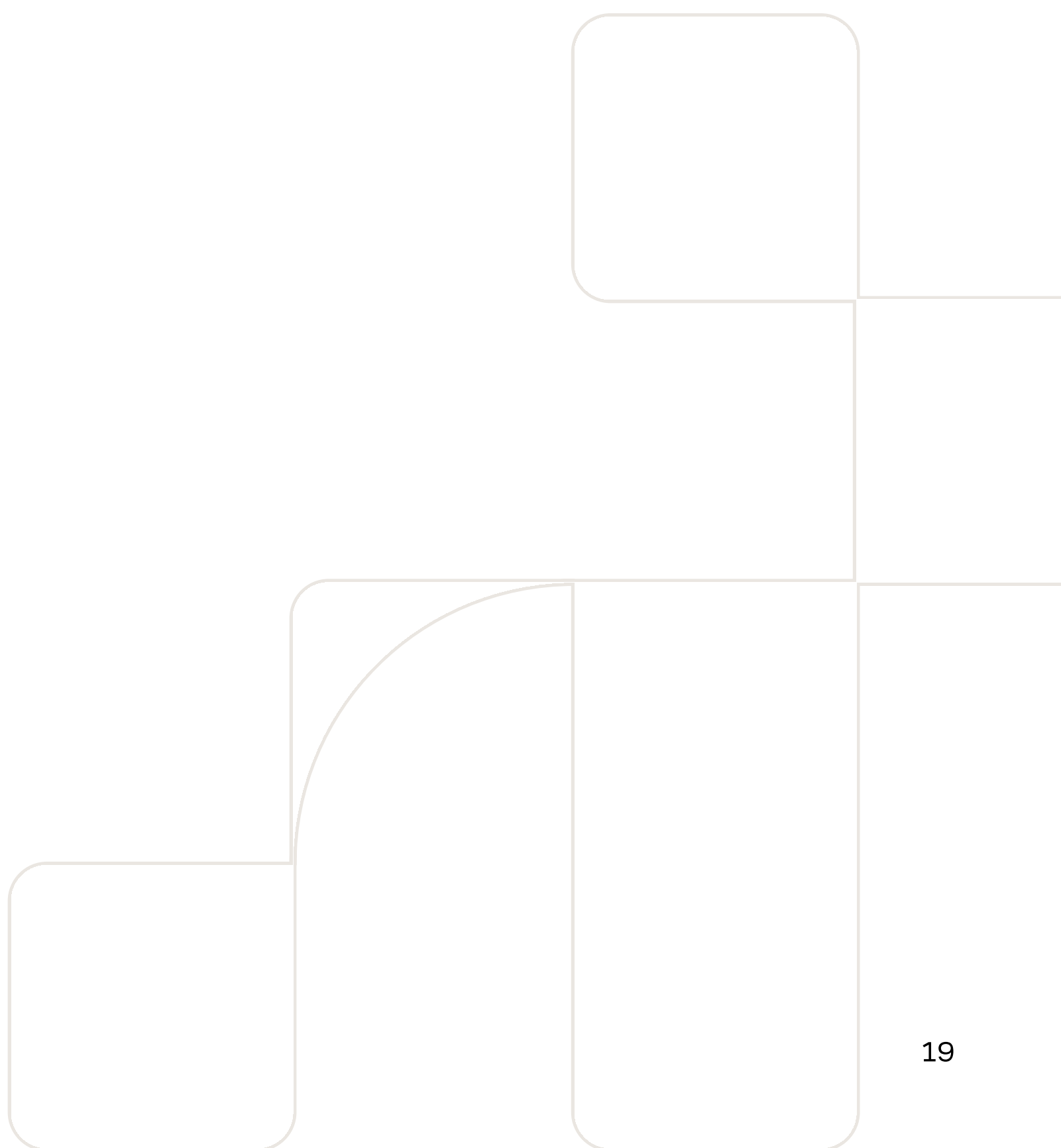- Be opinionated, consistent with choice of numbers

Good reporting assumes no statistics, infrastructure knowledge

The best forms of reporting don't assume context and knowledge. They don't assume you know what a p-value is, or what warehouse tables are. They simplify things, so that a junior Product Manager straight out of college can make a decision based on the results. This is a hugely underrated component of getting experimentation to scale.
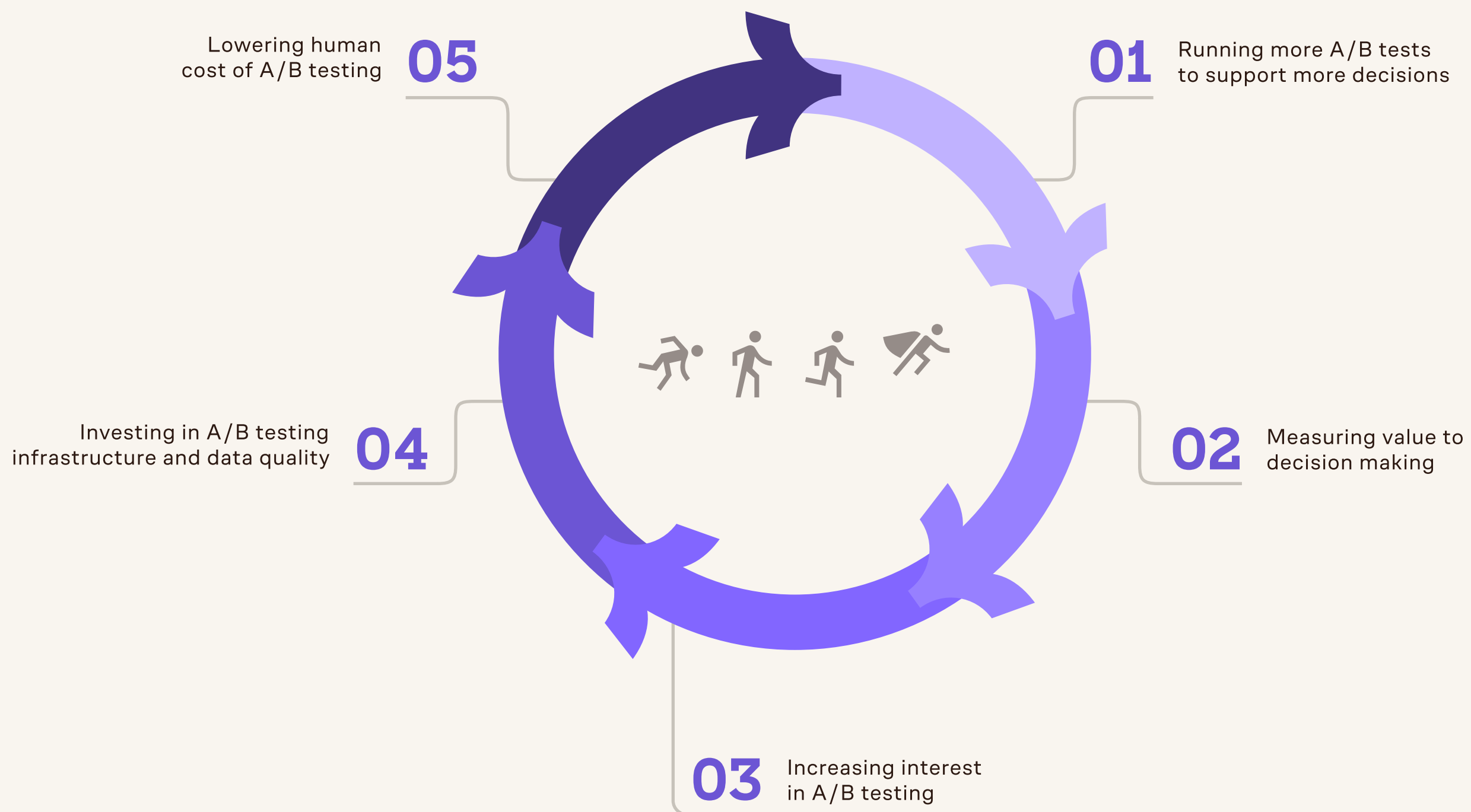
# Conclusion

Ultimately, building an experimentation tool is a complicated endeavor – there are a bunch of different components to get into sync. But even if you start out with something manual and low-fidelity, like Jupyter notebooks, you can begin to establish a flywheel that leads to experimentation success.

# The A/B Testing Flywheel



Crawl, Walk, Run, Fly Progression

When you run your first successful A/B tests, and start to influence decision-making on a small scale, you will start to increase your organization's interest in experimentation. Once you demonstrate that A/B testing leads to better decisions, your team will begin to invest in experimentation infrastructure, which ultimately lowers the labor costs of conducting further tests.

At Eppo, we think this flywheel is still a lot more painful than it needs to be. We built a tool that includes all the necessary components to quickly develop the experimentation culture that the most successful companies have used to win their markets.