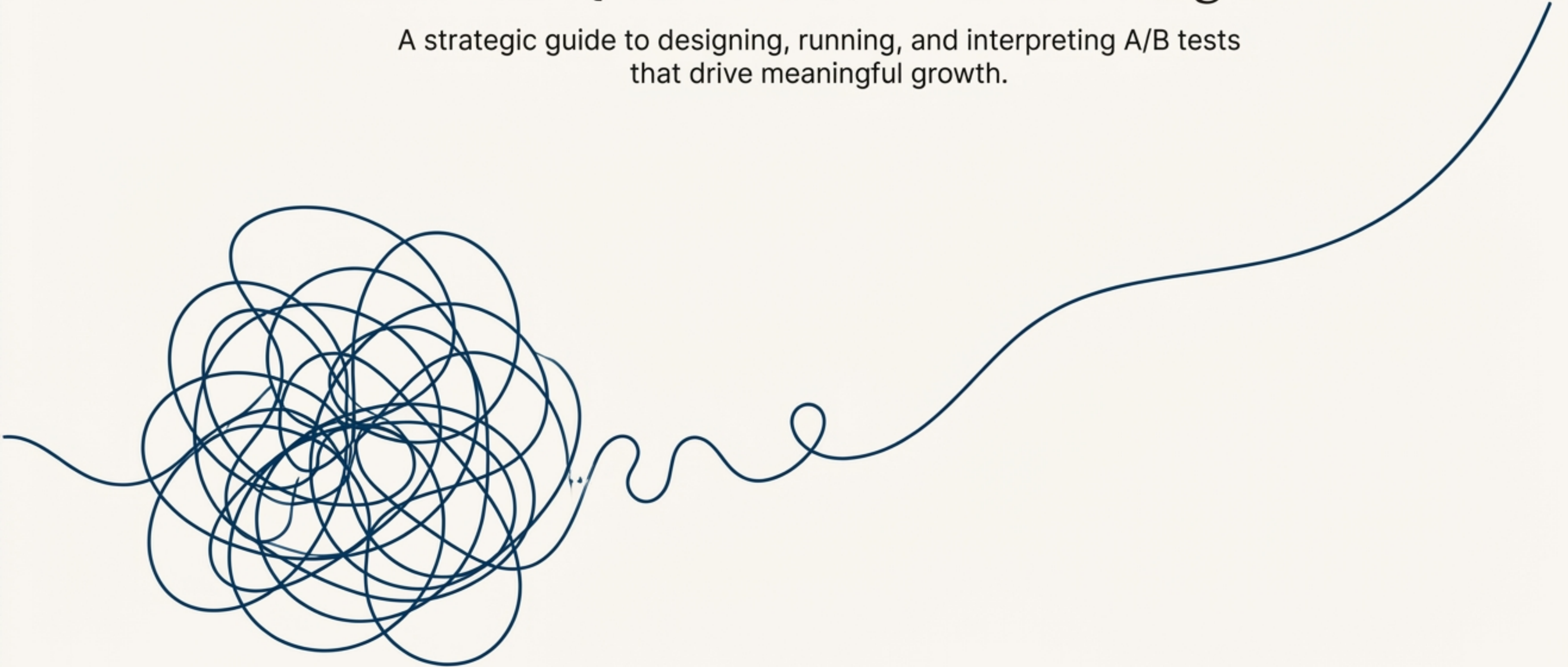


# The Experimentation Playbook: From Business Questions to Validated Insights

A strategic guide to designing, running, and interpreting A/B tests that drive meaningful growth.





# Experimentation Replaces Guesswork with Rigour

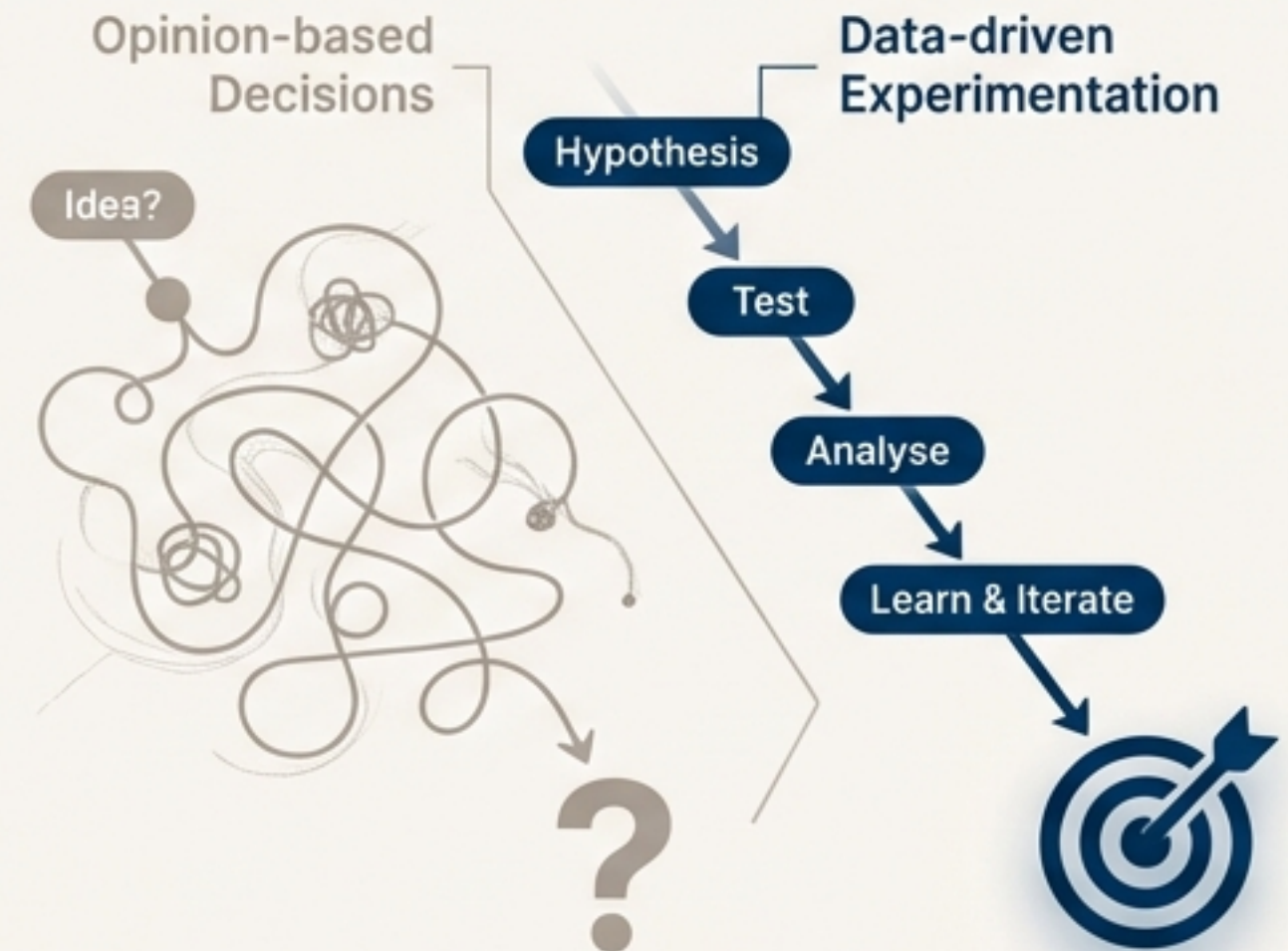
A mature experimentation programme eliminates subjective decision-making ('we think') and replaces it with objective, quantitative evidence ('we know'). It is a core business function for continuous product improvement and maximising ROI.

A/B testing is a randomised experimentation process used to empirically determine which version of a variable drives the maximum positive impact on defined business metrics.

The goal is to move beyond optimising for a 'local maximum' (a small, isolated win like a button colour) to pursuing a 'global maximum' by optimising the entire user journey (e.g., onboarding flows, pricing displays).

---

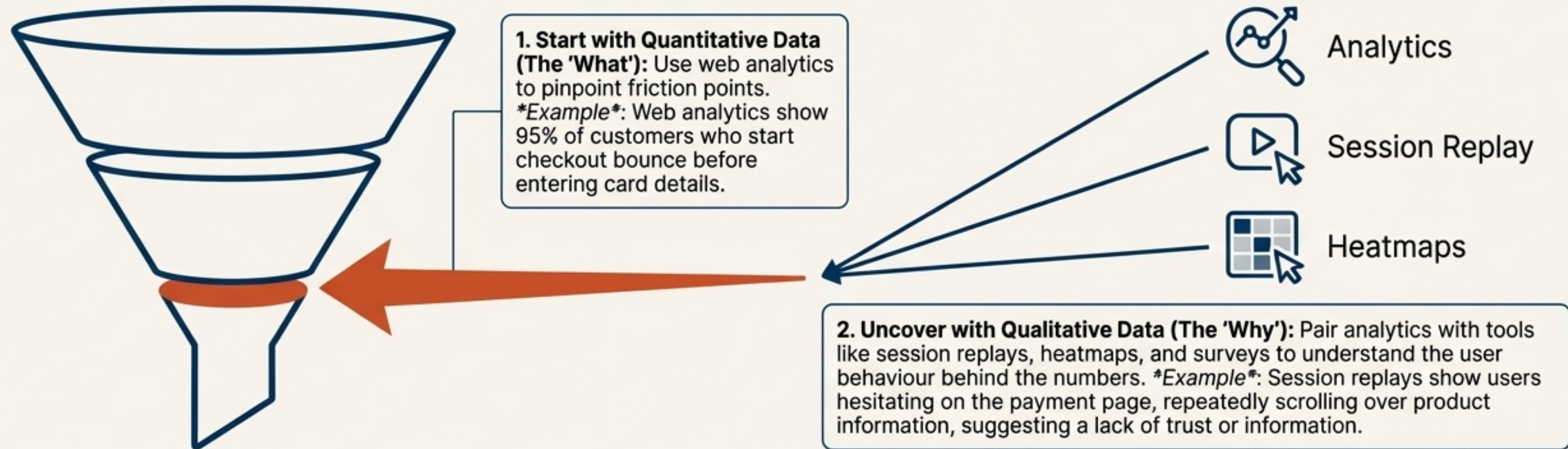
This playbook provides a structured, three-part framework to transform testing from a tactic into a strategic engine for growth.





# High-Impact Opportunities Are Found, Not Guessed

**Key Idea:** The experimentation process begins not with design changes, but with rigorous data collection to identify and validate user problems.



## Case in Point: Bannersnack

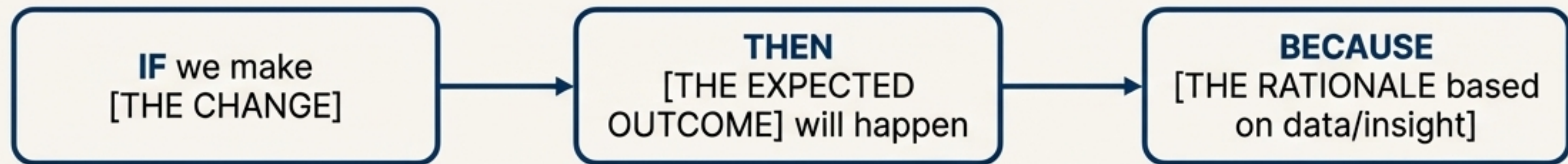
Bannersnack noticed low adoption of a new 'timeline' feature. Session replays revealed users were completely ignoring the button. Their hypothesis—"If we make the button bigger, more users will try it"—was validated with an A/B test, increasing feature adoption by 12%.



# A Strong Hypothesis Is the Compass for Your Experiment

An A/B test can only answer a **clear, closed-ended question**. A well-formed **hypothesis** translates a user problem into a testable prediction.

## The Anatomy of a Hypothesis



Example: IF we add a 'secure payment' icon to the checkout page, THEN more users will convert, BECAUSE we observed user hesitation and believe this addresses trust concerns.

## The Statistical View

### Null Hypothesis ( $H_0$ )

The status quo. It asserts there is **no difference** between the control and the variant.

e.g., "The icon will have no effect on conversion."



### Alternative Hypothesis ( $H_a$ )

The claim you want to prove. It asserts there **is a difference**.

e.g., "The icon will have a positive effect on conversion."



# Rigorous Metrics Define Success and Protect the Business

**Key Idea:** A successful experiment requires more than one success metric. A strategic set of metrics provides a holistic view of an experiment's impact.



## Primary Metric

The single KPI that determines if the hypothesis is validated. It must directly measure the behaviour you are trying to influence.

**\*Example\*:** Click-through rate on the 'Add to Cart' button.



## Secondary Metrics

Provide supporting context and help understand the "why" behind the primary metric's movement. They track broader user journey impacts.

**\*Example\*:** Average order value, time on page, bounce rate.



## Guardrail Metrics

The essential safety net. They monitor critical aspects of product health to prevent unintended negative consequences.

Ensure a win in one area doesn't cause a critical loss elsewhere.


**\*Real-World Example\*:** Netflix tests new recommendation algorithms (Primary: engagement) but uses stream start times and buffering ratios as guardrails to prevent degrading the core viewing experience.





# The Statistical Foundations of a Trustworthy Test

Before collecting data, you must define the statistical parameters that determine sample size, test duration, and how results are interpreted.

## The Three Levers of Experimental Design

 **Significance Level ( $\alpha$ )**  
The risk of a False Positive (Type I Error)—declaring a winner when one doesn't exist. Conventionally set at 5% (or 0.05), meaning you accept a 5% chance of being wrong. This corresponds to a 95% confidence level.

 **Statistical Power ( $1-\beta$ )**  
The probability of detecting a true effect if it exists, avoiding a False Negative (Type II Error). Conventionally set at 80%, meaning you have an 80% chance of finding a real winner.

 **Minimum Detectable Effect (MDE)**  
The smallest lift in the primary metric that the experiment is designed to detect. This is the most critical practical input, as it directly governs the required sample size.

## The Core Trade-off

Parameter	Impact on Required Sample Size (N)
Decrease Significance Level $\alpha$ (e.g., 5% -> 1%)	Sharply Increases N $\uparrow\uparrow$
Increase Statistical Power (e.g., 80% -> 90%)	Increases N $\uparrow$
Decrease MDE (Detect smaller effects)	Dramatically Increases N $\uparrow\uparrow\uparrow$



# Mastering the MDE: Your Lever for Balancing Precision and Practicality

**Key Idea:** MDE measures the sensitivity of your experiment. A low MDE detects small changes but requires a massive sample size; a high MDE is faster but may miss subtle, valuable improvements.

## MDE as a Resource Allocation Tool

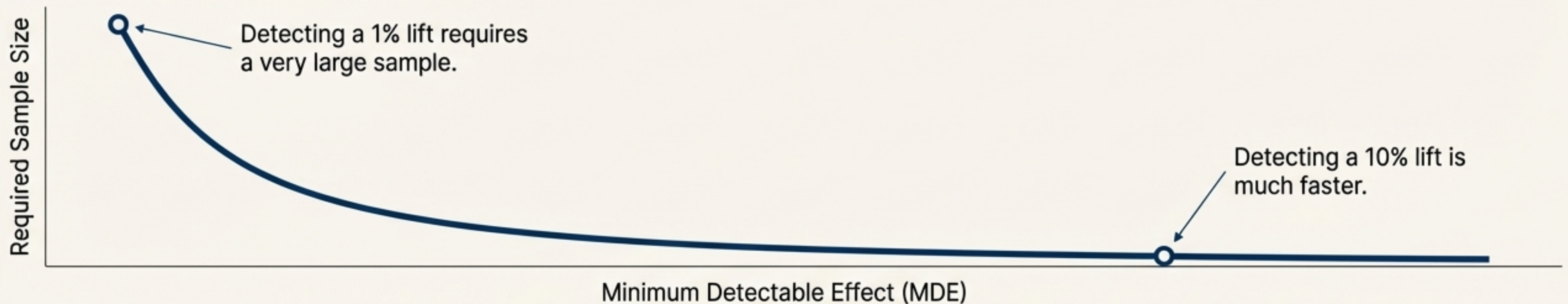
### Low MDE (High Sensitivity)

- Detects small lifts (e.g., 1-2%).
- Requires very high traffic and/or a long duration.
- Use for optimising critical, high-volume flows like checkout.

### High MDE (Low Sensitivity)

- Only detects large lifts (e.g., 10%+).
- Requires less traffic and runs faster.
- Use for testing bold changes or on lower-traffic pages.

## The Relationship Visualised



## Expert Prescription

Your MDE choice must align with your traffic volume and the cost of traffic acquisition. If a test cannot achieve 80% power to detect a practically relevant MDE (a common rule of thumb is 2-5%), the results are unlikely to be trustworthy or worth the development effort.



# Two Non-Negotiables for Data Integrity: Duration and Randomisation

## Pillar 1: Run Your Test for the Right Duration

**It's not just about hitting a sample size number.**



**Account for Business Cycles:** If your typical purchase cycle is one week, a three-day test is unrepresentative. Run tests for at least one full business cycle, ideally two.



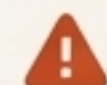
**Avoid Day-of-Week Effects:** User behaviour on a Monday is different from a Saturday. **Always run tests for full weeks** (e.g., 7, 14, or 21 days) to capture these patterns and avoid skewing results. Low-traffic sites may need up to eight weeks.

		1	2	3	4	5
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21					

Two Full Cycles

## Pillar 2: Get Randomisation Right

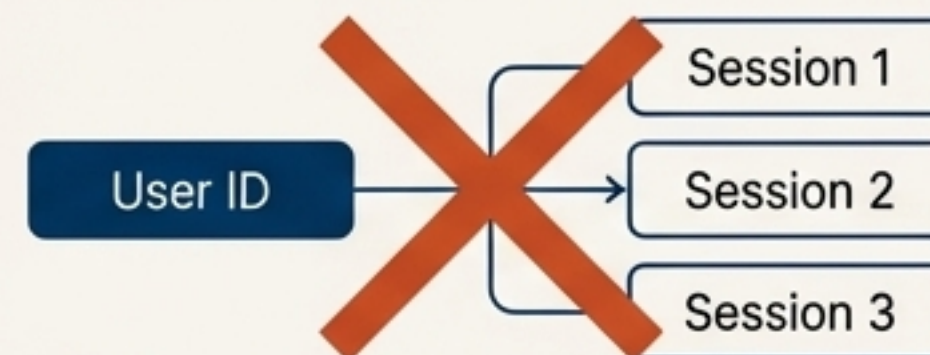
**Randomisation Unit:** The identifier used to assign users to groups (e.g., User ID, Device ID).



### Pitfall to Avoid - The Mismatch Problem

A critical error occurs when the **randomisation unit** does not match the **analysis unit**.

*Example:* You randomise by **User ID** but analyse conversion rate **per session**. A single user can have multiple sessions, so these data points are not independent.



This inflates statistical significance and dramatically increases the risk of a false positive.

**Expert Prescription:** Always randomise at the user level for consistency. Ensure your analysis unit is not more granular than your randomisation unit.



# Avoid These Pitfalls That Corrupt Your Results

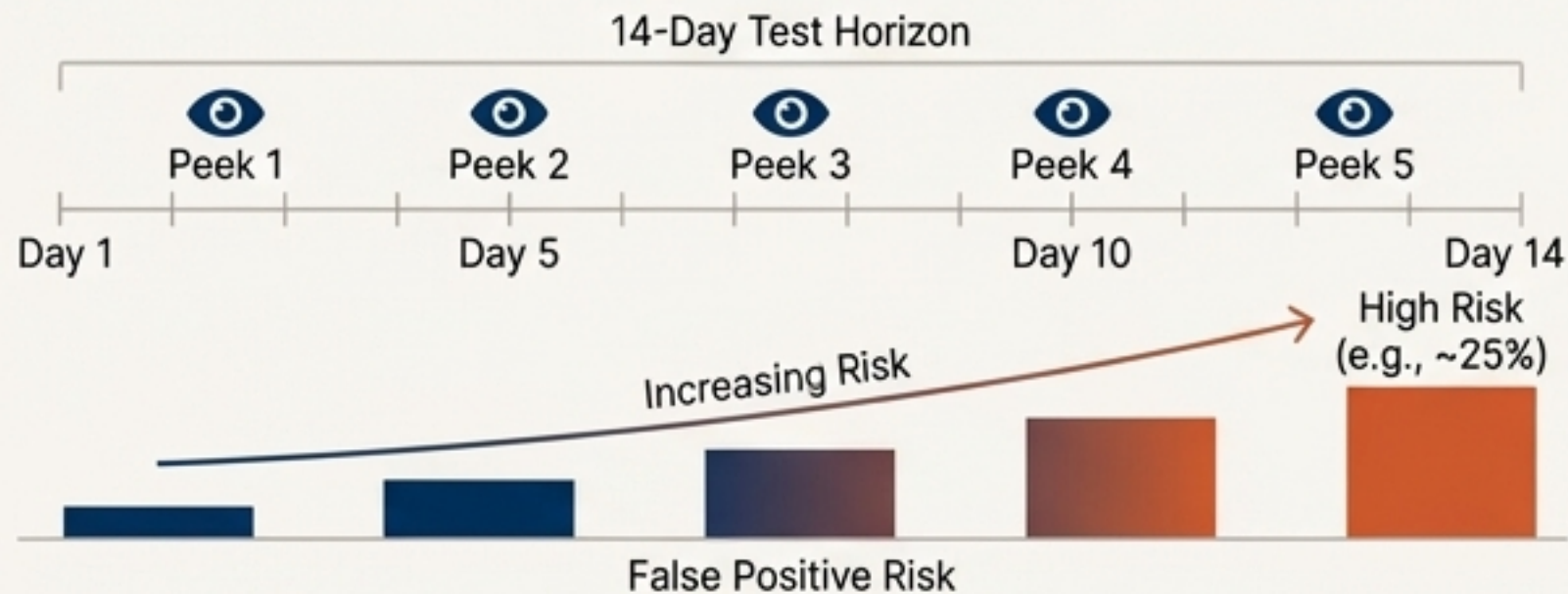
## Pitfall 1: 'Peeking' at Results Prematurely

### What it is

Checking results before the pre-determined sample size is reached.

### The Danger

Every "peek" is a statistical test that invalidates the initial significance calculation. It dramatically inflates the cumulative probability of a Type I error (False Positive). Peeking just 10 times can turn a result that looks 99% significant into one that is only 95% significant.



### The Prescription

Use a fixed-horizon test and do not analyse data until the end.

Alternatively, use a testing platform with a sequential testing engine that is specifically designed to manage error rates during continuous monitoring.

## Pitfall 2: Testing Too Many Variables at Once

### What it is

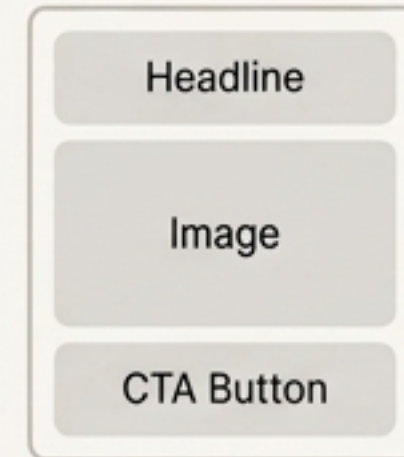
Changing the headline, image, and CTA button all in one variant.

### The Danger

It becomes impossible to isolate which specific change was responsible for the observed result.

You learn nothing about why it won or lost.

### Original Variant (A)



### Variant (B)



### The Prescription

Test one change at a time to understand its causal impact.

If you must test multiple elements, use a Multivariate Test (MVT), but be aware it requires substantially more traffic.



# Interpreting Results with Confidence Intervals

**Key Idea:** Move beyond a simple p-value. Confidence Intervals (CIs) provide a range of plausible values for the true lift, giving you a clearer picture of the magnitude and uncertainty of the effect.

## How to Read a Confidence Interval

- A 95% CI means that if you repeated the experiment many times, the interval would contain the true difference 95% of the time.
- **The Zero Rule:** If the CI range *does not* include zero, the result is statistically significant.
- **The Width Rule:** A **narrow** interval indicates high precision (more confidence). A wide interval indicates high uncertainty (less confidence), often due to small sample size or noisy data.





# Statistical Significance Is Not Business Significance

## Statistical Significance

Confirms that an observed effect is unlikely due to random chance. It tells you an effect *exists*.

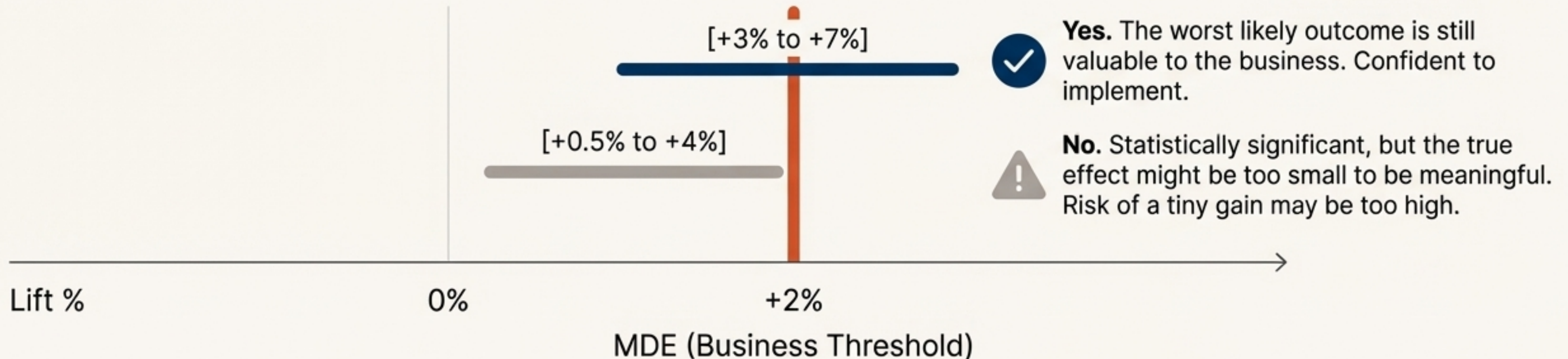
## Practical Significance

Refers to the *magnitude* of the effect. Is the lift large enough to justify the cost of implementation and the opportunity cost?

## Using Confidence Intervals to Judge Practical Significance

Don't just check if the CI **excludes zero**. Examine the **entire range**.

**The MDE Test:** Is the *lower bound* of your confidence interval above your pre-defined MDE?





# The Go/No-Go Decision Framework for Mixed Outcomes

## The Common Dilemma

Your test shows a significant win on the primary metric, but a critical guardrail metric (like page load time, error rates, or user churn) is significantly harmed.

## The Rule for Moving Forward

Ship the change **if and only if** the treatment is significantly superior on at least one success metric **AND** significantly non-inferior (or safe) on all guardrail metrics.

## The Decision Matrix

Primary Metric Result	Guardrail Metric Result	Interpretation	Recommended Action
Significant Positive Lift	No significant change (Safe)	Ideal outcome. Effective and robust.	✔ <b>GO:</b> Implement the variation.
Significant Positive Lift	Significant Negative Impact	Success achieved at the cost of product health. A hidden loss.	✖ <b>NO-GO / INVESTIGATE:</b> Do not ship. Redesign to mitigate the harm and re-test.
Not Statistically Significant	No significant change (Safe)	Change had no measurable impact.	↺ <b>LEARN / ITERATE:</b> Document learnings and test a bolder hypothesis.
Significant Negative Impact	No significant change (Safe)	Change actively harmed the key metric.	✖ <b>NO-GO:</b> Document the failure and pivot to a new hypothesis.

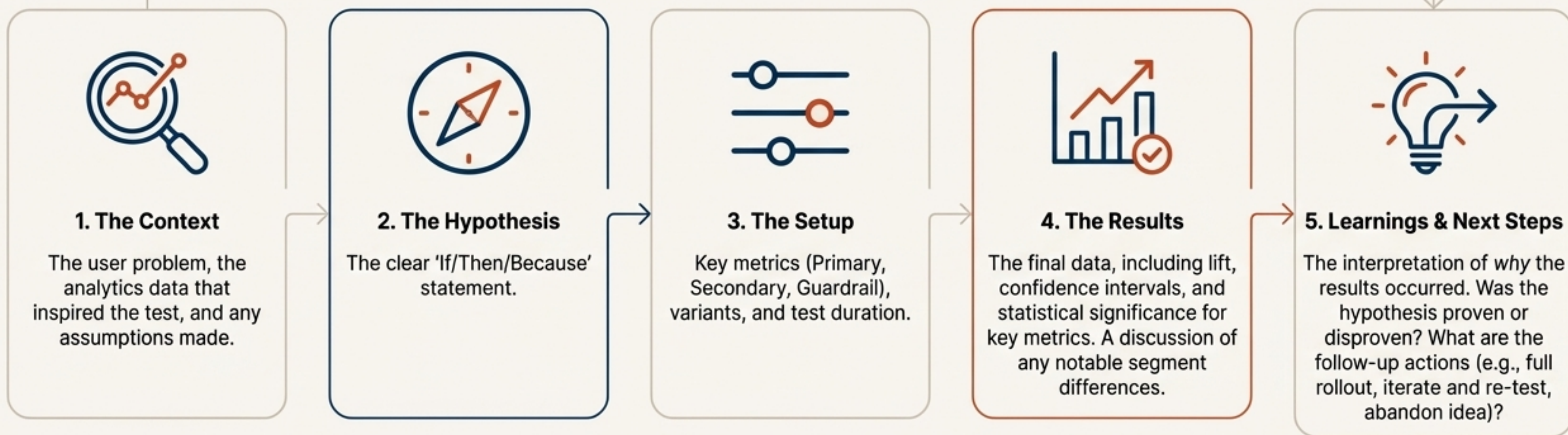


# Operationalising Learnings to Build Institutional Knowledge

**Key Idea:** Every test—win, loss, or inconclusive—is an opportunity to learn. Thorough documentation prevents repeating mistakes and builds a deeper understanding of your users.

## Anatomy of an Experiment Report

Create a simple, standardised template that everyone sticks to. It should include:



**A documented failure is valuable institutional knowledge.**  
**An undocumented test is a wasted resource.**



# From A/B Testing to a Continuous Growth Engine

A disciplined approach transforms A/B testing from a series of one-off tactics into a strategic, repeatable process that drives sustainable growth.



Mature experimentation is not about finding the perfect button colour. It's about **building a culture of curiosity and discipline** that **systematically reduces uncertainty and validates the path to achieving your most important business goals.**