# Multimodality Under Governance

## Embeddings, Alignment, Failure Modes, and Drift-Aware Operations

Multimodality is not "more modalities."
It is **learning compatible coordinate systems**—and keeping them stable under change.

---

**Companion Computational Laboratories**

This volume is accompanied by **3 fully executable Google Colab notebooks**, one per chapter, built in **pure Python** and **synthetic-first** style. Each notebook implements a governed laboratory with: deterministic seeds, explicit latent factors, controlled perturbations, diagnostic visualizations, and audit artifacts . These notebooks are not benchmark engines. They are **mechanism labs** designed for reviewability and teaching.

---

Alejandro Reynoso

Chief Scientist, DEFI Capital Research

External Lecturer, Judge Business School, University of Cambridge

February 17, 2026

# Preface: Why Multimodality Is a Frontier Topic Under Governance

Multimodality is often described as a capability upgrade: add images, add audio, add video, and the model becomes "more human." That framing is seductive. It fits the way product demos are shown, it fits the way people experience modern interfaces, and it fits the cultural intuition that intelligence must be multi-sensory. But it is not mechanism-first, and it is not institution-ready. A multimodal model does not become operationally valuable by accumulating modalities. It becomes operationally valuable only if it learns **compatible coordinate systems** across modalities—so that the same latent state can be represented, compared, retrieved, and monitored under a shared geometry.

That single sentence hides the whole frontier. "Compatible coordinate systems" is not a slogan. It is the requirement that makes multimodality difficult in the ways that matter for professional practice. Text and images are not merely different file formats. They are different measurement operators acting on the world. A text description of a chart is not the chart; it is a measurement of the chart through language. An image is not the object; it is a measurement of the object through optics, compression, sensor noise, and cropping. Audio is not the conversation; it is a measurement of the conversation through microphones, rooms, codecs, and diarization. When we build a multimodal model, we are asking the system to discover a latent structure that is invariant enough across these operators that two different observations can be mapped into a single space where "closeness" means something stable.

This is why multimodality belongs in a frontier-oriented sequence of lectures. It is not frontier because it is fashionable. It is frontier because the mechanistic problem is genuinely hard, and because the failure modes are structural rather than cosmetic. In professional domains, the cost of a failure is not that the model makes a strange caption. The cost is that a system that appears coherent becomes unreviewably wrong. A multimodal system can produce a confident narrative because it finds a geometrically convenient alignment, not because it has performed a correct inference. This is exactly the kind of risk amplification that governance-first work exists to control.

To see why the problem is geometric, start from the simplest abstraction. Once a system is embedded, its behavior is vectorial. A text sequence becomes a vector. An image becomes a vector. An audio

segment becomes a vector. "Similarity" becomes a dot product or a cosine. Retrieval becomes neighborhood competition. Clustering becomes density and separation in a high-dimensional space. Alignment becomes a constraint: two encoders, observing different modalities, must map paired observations to nearby points while pushing unpaired observations apart. When we say a multimodal system "understands," what we really mean is that its embedding geometry has become useful: it places things near each other that should be treated as similar for the downstream task.

The frontier is not that we can map multiple modalities into vectors. That part is easy to demonstrate. The frontier is that we can do it *reliably*: under distribution shift, under noisy supervision, under confounders, under pipeline updates, under re-scoring, under fine-tuning, and under adversarial pressure. The frontier is the lifecycle. The frontier is the stability of geometry as a deployable object.

This point is easy to miss because most people meet multimodality through success cases. A retrieval demo finds the right image for a caption. A captioning demo describes a scene. A document assistant answers questions about a PDF. These experiences are impressive, but they hide the question that matters to an institution: *what exactly is the system doing, what can it be trusted to do, and how will we know when it stops doing it?* If we cannot answer those questions, multimodality becomes an uncontrolled surface. It is a capability that can drift silently.

The reason drift is especially dangerous in multimodal systems is that there are more places for drift to occur and more ways for drift to be misinterpreted. In a text-only system, the main measurement operator is tokenization and the distribution of language. In multimodal systems, every modality has its own measurement operator and its own pipeline. Images come with sensor drift, compression drift, cropping drift, and watermark drift. Text comes with style shift, normalization changes, encoding artifacts, and metadata leakage. The pairing relation between modalities can drift through join bugs and ingestion mismatches. The scoring layer can drift through temperature changes, indexing changes, or quantization. Even if the encoders remain frozen, the system can change. And if the system is continuously fine-tuned, then the embedding space itself becomes a moving target.

It is here that governance stops being a philosophical overlay and becomes a technical necessity. Governance-first does not mean "write disclaimers and move on." It means we treat the system as a controlled process with explicit states, explicit update rules, explicit constraints, and explicit audit artifacts. It means we accept that embeddings are not self-explanatory, and we design the monitoring layer that makes the geometry observable. It means we accept that drift is inevitable, and we design bounded interventions and rollback discipline.

Under the *AI 2026* umbrella, this mini-book is intended to sit alongside other frontier-aware topics that share the same structural theme. Long context and long memory are frontier topics not because the window size is large but because increasing the candidate set changes the geometry of competition. Surrogate models are frontier topics not because they are new but because they replace expensive mechanisms with approximations that must be governed under shift. Multimodality

belongs in the same family because it builds a shared geometry across measurement operators and then asks that geometry to remain coherent under change. These are not separate stories. They are three versions of a single question: *how do we operate representation under uncertainty and drift without confusing fluency for truth?*

This is also why the correct professional objective is not persuasion or aesthetic fluency. It is **reviewability**. Reviewability is the ability to show what the system learned, how it changed, and why it can be supervised. In a governed setting, a model output is not accepted because it is eloquent. It is accepted because it is traceable to evidence and because the system that generated it can be audited. Multimodality threatens reviewability precisely because it can create plausible but untraceable couplings. An image can contain signals that a human does not notice. Text can contain metadata that a human treats as irrelevant. The model can learn shortcuts that bind them together. A retrieval system can return the right item for the wrong reason. Without governance, you do not detect this. You celebrate it. That is why this is a frontier topic under governance.

The three chapters in this volume therefore form a single argument that moves from construction to fragility to lifecycle control.

**Chapter 1** treats multimodality as geometry by construction. We build a synthetic world with explicit latent factors and two modalities: small synthetic images and symbolic token sequences. The images are not photographs; they are controlled matrices generated by parameters such as shape, orientation, spatial frequency, phase, and thickness. The text is not natural language; it is a token grammar that encodes those same factors. This design is deliberate. It removes the myth that multimodality is magic. It makes the coordinate-system problem explicit: we know what the latent factors are, we know how each modality measures them, and we can ask whether the learned embedding space respects that structure.

The chapter then trains a two-encoder aligner with a contrastive objective. The objective is simple: paired image and text representations should be close, and non-paired representations should be far. But the point is not the loss value. The point is the geometry. We visualize the embedding space with PCA, we probe factor separability, we compute retrieval accuracy in both directions, and we measure collapse indicators such as mean cosine similarity and embedding variance floors. Students see what "latent space" means in practice: it is a coordinate system in which the same underlying factors become linearly discernible and distance becomes meaningful. They also see what it does *not* mean: it is not a magical container of truth. It is a learned geometry under an objective, and objectives can be satisfied in multiple ways.

This first chapter therefore establishes the correct mental model. Multimodality is not "text plus image." It is **two measurement operators plus one shared coordinate system**. Once students internalize that, a set of natural questions appears: what can go wrong, and how do we detect it?

**Chapter 2** answers by being intentionally critical. Alignment is fragile. We induce failure modes in a controlled setting. We construct modality dominance by making one modality cleaner or by scaling

its encoder outputs so that it overpowers the other. We inject spurious alignment by introducing a confounder feature shared across modalities that is unrelated to the intended latent factors, such as a parity marker or an ID codebook. We corrupt pairings to mimic join bugs and label noise. We induce representation collapse by pushing the temperature too low, pushing the learning rate too high, or applying regularizers that unintentionally crush variance.

The point of Chapter 2 is not to frighten students. The point is to immunize them against a common professional mistake: the belief that high metrics imply semantic alignment. In the confounder experiments, retrieval can improve while factor information declines. This is the most dangerous kind of failure because it looks like success. It teaches students the difference between *task performance* and *mechanism correctness*. A system can perform well by learning the wrong mechanism. In governance-first practice, that is unacceptable, because when the shortcut disappears—as it often does when pipelines change—the system collapses, and the collapse is mysterious unless you have probe metrics that track mechanism.

Therefore, Chapter 2 introduces a diagnostic posture. We treat metrics as evidence, not as verdicts. We implement structural probes such as mutual information proxies between embeddings and true factors, CCA-like correlation proxies between modalities, and sensitivity proxies that approximate how a small training perturbation changes retrieval. We monitor spectra, hubness, and symmetry gaps. The chapter teaches students that a multimodal model is not a single number. It is a geometry that must be interrogated.

By the end of Chapter 2, students have learned the critical taxonomy: dominance, spurious shortcuts, pairing corruption, and collapse. The remaining question is the most important one for real deployments: what happens over time?

**Chapter 3** treats multimodality as a lifecycle system. It formalizes drift taxonomies and implements monitoring, classification, and bounded interventions. The core idea is simple: a multimodal deployment is not one model at one time. It is a sequence of episodes in which the data pipeline, the scoring layer, the pairing relation, and the model itself can change. Therefore, we must operate the system with a monitoring panel and stage gates, not with a one-time evaluation.

The chapter builds a synthetic drift laboratory: a timeline of controlled drift events that mimic production realities. We introduce modality drift (noise or corruption in one modality), scoring drift (temperature shifts), confounder introduction and disappearance, pairing drift (swaps and systematic offsets), and optimization drift (fine-tuning regimes that cross collapse boundaries). For each episode, we compute monitoring deltas relative to baseline: directional retrieval, symmetry gaps, collapse indicators, spectral summaries, and probe metrics. We then classify drift using multi-signal rules that are transparent and allowed to output "unknown." This is a governance lesson: false certainty is worse than uncertainty. If evidence is insufficient, the correct action is to escalate, not to invent a narrative.

Finally, Chapter 3 defines a bounded intervention policy. Not every drift event requires retraining.

Some drift is in scoring and can be corrected by temperature recalibration. Some drift is in data and must be corrected by sanitation filters or confounder ablation. Fine-tuning exists, but it is bounded by stop rules and rollback discipline. This is the operational heart of governance-first: *we do not allow optimization to become a blind response to uncertainty*. We require stage gates that include not only performance but geometric health constraints. We accept interventions only if the system remains reviewable.

All three chapters are paired with synthetic-first Colab laboratories. The labs are designed to be runnable end-to-end, auditable, and teachable. They produce artifacts and logs suitable for review: run manifests, prompts logs, risk logs, deliverables folders with plots and strict JSON summaries. The labs do not claim production readiness. They are training instruments. Their purpose is to train a student's ability to reason about mechanism and governance, not to train a model to win a benchmark.

This synthetic-first stance is not a retreat from reality. It is the only way to teach multimodality without confusing the student. Real multimodal datasets are entangled. They contain unknown confounders, unknown biases, unknown annotation errors, and unknown pipeline effects. In such datasets, a student cannot easily separate "the model failed" from "the dataset was messy" from "the pipeline shifted" from "the metric was mis-specified." Synthetic environments provide ground truth. They allow controlled stress testing. They make causality visible. In a frontier topic, this matters more than scale. A large dataset can produce impressive curves, but it can also hide wrong lessons. A small synthetic dataset can produce correct lessons about mechanism and drift.

The broader pedagogical aim is to give students a clear understanding of what an embedding is in a multimodal context. Students often hear that embeddings are "representations," but the word representation is too vague. In this volume, an embedding is treated as a coordinate assignment. It is a mapping from an observation to a point on a manifold where distances and directions encode the relationships that a downstream task will exploit. Learning a multimodal model is learning two coordinate maps whose images are compatible. That is why geometry matters. That is why spectra and hubness matter. That is why collapse is a geometric pathology. And that is why governance is a geometric necessity: you cannot govern what you cannot measure, and you cannot measure a learned geometry without explicit probes and artifacts.

This is also why multimodality is not merely an engineering topic. It is a systems topic. It sits at the intersection of representation learning, data pipelines, measurement theory, and deployment governance. In a classroom, it should not be taught as a list of architectures. It should be taught as a set of operators and failure modes. The student should leave with a mental model that generalizes beyond any particular model family.

In finance, the analogy is immediate. Markets are multi-modal systems. Price is one modality. Volume is another. Order flow is another. News is another. Fundamental accounting signals are another. In institutional trading and risk, we are always trying to align measurement operators into

a coherent state: we want a regime belief that is consistent with price behavior, consistent with liquidity conditions, consistent with macro releases, and consistent with risk constraints. Multimodal embedding models are therefore not foreign to finance. They are a new implementation of an old problem: constructing a shared state from heterogeneous measurements. But finance also teaches the governance lesson: state that cannot be audited becomes a liability. A risk committee does not accept "the model felt that liquidity was worsening." It accepts a report of measured indicators, thresholds, and versioned logic. The same standard must apply to multimodal AI systems if they are to be used in professional contexts.

Therefore, the thesis of this mini-book can be stated without hype.

We do not treat multimodality as a magical step toward general intelligence. We treat it as a representational engineering problem: learning compatible coordinate systems across measurement operators and keeping those coordinate systems stable under drift. That is frontier work because the lifecycle is difficult. The space is high-dimensional, the supervision is noisy, shortcuts are everywhere, and updates are inevitable. Governance-first does not slow this down; it makes it possible. It provides the instrumentation that lets us know what the system is doing, the controls that bound its failure modes, and the artifact trail that makes it supervisable.

---

**Warning**

**Professional warning.**
In multimodal systems, the most dangerous failure is not a visible error. It is a hidden shortcut that improves metrics while degrading semantics, until a pipeline change removes the shortcut and the system collapses without warning. If you are not monitoring geometry and mechanism, you are not monitoring risk.

---

**Artifact (Save This)**

**Core thesis (one line).**
**Multimodality is shared geometry; governance is what keeps shared geometry usable under drift.**

---

The chapters that follow develop this thesis with explicit constructions, controlled failure inductions, and lifecycle monitoring and intervention logic. If the reader engages with the book as intended—text plus lab—they will not merely learn a vocabulary of multimodal models. They will learn how to reason about multimodal systems as governed geometric objects, which is precisely the skill needed to teach frontier AI responsibly and to deploy it without confusing fluency for truth.

# How to Use This Book

Each chapter is built as a single instructional unit with two inseparable components:

- **The text** formalizes the mechanism, names failure modes, and defines what counts as evidence.
- **The lab** operationalizes the mechanism with synthetic factors, controlled perturbations, and auditable artifacts.

Readers should follow a strict sequence:

1. **Read for objects and operators.** Identify latent factors, measurement operators, embedding maps, and scoring rules.
2. **Run the notebook without modification.** Reproduce baseline behavior before changing knobs.
3. **Inspect diagnostics as primary outputs.** Treat geometry, spectra, retrieval asymmetry, and collapse indicators as the real results.
4. **Stress one axis at a time.** Change noise, confounders, pairing corruption, temperature, or fine-tuning regimes, and record signatures.
5. **Use the audit bundle.** Each run produces manifests, logs, plots, and summaries; write short memos from these artifacts.

This is not a leaderboard exercise. It is a **mechanism and governance** exercise: understand what alignment means, when it fails, and how it can be supervised.

# Contents

# Chapter 1

# Multimodality by Construction

**Abstract.** We treat multimodality as a problem of *compatible coordinate systems* rather than input concatenation. Each modality is a measurement channel: images encode geometric and frequency structure through local spatial interactions, while text encodes symbolic structure through discrete compositional rules. A multimodal model is therefore a machine that learns a shared representation in which different measurements of the same underlying object become comparable. This chapter develops a mechanism-first, synthetic-first account of that claim. We construct a minimal multimodal world with explicit latent factors (shape, orientation, frequency, phase, thickness) and generate paired observations: small synthetic images and symbolic token sequences that encode the same factors. We then train two encoders, one per modality, to map observations into a shared embedding space using a contrastive (InfoNCE) objective with cosine similarity and temperature scaling.

The contribution is pedagogical and operational. Pedagogically, we show how latent factors become geometry: which directions in the embedding space correspond to factors, how clusters emerge, and how distances change under controlled perturbations. Operationally, we insist on reviewability. The companion notebook implements analytic gradients for two-layer MLP encoders, validates them with finite-difference checks, and instruments training with diagnostics that go beyond loss: cross-modal retrieval in both directions, modality symmetry checks, separability scores by factor, and collapse indicators based on embedding variance and average cosine similarity. Finally, we treat stress testing as the scientific method for representation learning. By sweeping noise, pairing quality, and temperature, we map degradation curves and define acceptance criteria as distributions rather than point estimates.

The outcome is a rigorous mental model of multimodal embeddings as governed geometry: alignment is a contract, not a vibe, and the only defensible understanding is one that can be reproduced, inspected, and falsified.

## 1.1 Introduction

### 1.1.1 Motivation: multimodality as coordinate compatibility

Multimodality is often presented as a capability milestone: a model can "see" images, "read" text, and sometimes "hear" audio. For serious students of modern AI systems—especially students trying to reason about frontier behavior rather than repeat fashionable claims—this framing is inadequate. The central difficulty is not merely that multiple input types exist. The deeper difficulty is that each modality is a distinct measurement channel with its own invariances, noise sources, and representational biases, yet a multimodal system must behave as if these channels are coordinated views of one underlying world. The correct mental model is therefore structural: multimodality is the learning of *compatible coordinate systems* under imperfect measurement and limited supervision. The model's job is not to "merge" modalities; it is to impose a geometry in which meaning remains stable when expressed through different forms.

Once pixels and tokens are embedded as vectors, "understanding" becomes geometry. A multimodal model learns at least two maps—one from images to vectors and one from text to vectors—and then tries to place corresponding image–text pairs near each other while pushing non-corresponding pairs

apart. That objective sounds deceptively simple. The geometry it induces is subtle. It is shaped by the density of negatives in each batch, by the temperature parameter that controls similarity sharpness, by normalization and scale choices, and by the relative capacity of each encoder. The resulting space is not merely a convenient intermediate representation; it becomes the system's internal notion of equivalence: what it treats as "the same," what it treats as "close," and what it treats as "irrelevant." If the goal is to teach multimodality at a frontier level, the student must learn to ask the right question: not "does it work," but "what geometry did it learn, what does that geometry mean, and how does it fail?"

The professional motivation is straightforward. Embeddings are not just intermediate computations. In production systems, embeddings become infrastructure. They drive retrieval, clustering, de-duplication, ranking, recommendation, cross-modal search, and conditional generation. A small change in embedding geometry can change system behavior dramatically without changing any explicit business rule. This makes multimodality a governance problem as much as a modeling problem: representation is a policy layer. When the representation changes, the system's effective policy changes. Therefore, teaching multimodality purely as an architecture topic is incomplete. The student must also learn an operational posture: if embeddings are used as infrastructure, then we require stage gates, diagnostics, drift detection, and a reproducible artifact trail.

### 1.1.2 Scope and non-goals: synthetic-first, mechanism-first

This chapter adopts a synthetic-first methodology, not because real data is unimportant, but because real data makes it too easy to confuse statistical success with conceptual understanding. When students learn multimodality from large-scale datasets, they often absorb the wrong lessons: that performance is a function of scale, that alignment "emerges," and that embeddings are mysterious objects one cannot directly interpret. In reality, embeddings can be interpreted, but only if the world that generated them is controllable and the factors that matter are known. Synthetic design provides this controllability. We construct a multimodal world with explicit latent factors—shape type, orientation, frequency content, phase, and thickness—and generate two observations of each latent state: a small image matrix and a symbolic token sequence encoding the same factors. This allows us to test, explicitly and quantitatively, whether the learned embedding geometry reflects the intended semantics.

The phrase *multimodality by construction* is central. It means we are not asking the model to discover what the world is about from messy, confounded data. We are constructing the world so that intended semantics are known. This is a pedagogical strategy: first learn the mechanics of alignment when the truth is available, then study how those mechanics break when the truth is hidden or confounded. In frontier terms, this is analogous to studying controlled physical systems before studying real turbulence. The goal is not to simplify the topic into something trivial; it is to isolate key degrees of freedom so that the student can see what changes what. The moment

a student can predict how temperature will affect retrieval, how noise asymmetry will produce directionally asymmetric retrieval, and how normalization can prevent collapse, the student has moved from superficial knowledge to structural understanding.

This chapter is not a foundation-model survey. We do not attempt to replicate the full complexity of large transformer-based multimodal systems, nor do we depend on external datasets or proprietary training pipelines. We deliberately avoid claims of real-world superiority. The non-goal is to "match production"; the goal is to *understand the mechanism*. This is the only reliable way to keep frontier topics teachable and governable. Once the minimal mechanism is understood, the student is prepared to reason about scaling, architectural complexity, and data realities without collapsing into mysticism.

### 1.1.3 Professional objective: reviewable representations, not persuasion

A second theme is that "more data" and "bigger models" do not solve the conceptual problems of multimodality. They can hide them. Consider a dataset containing a watermark shared across image and text channels, such as an ID embedded in captions or a formatting artifact present in both modalities. A large model can achieve spectacular retrieval by exploiting the shortcut rather than learning the intended semantic relationship. Similarly, if one modality is consistently cleaner than the other, the model can align by allowing one modality to dominate the similarity structure. The resulting system can appear successful while being fragile. If the pedagogical goal is to prepare students for frontier systems, the correct response is not to celebrate high metrics; it is to insist on reviewability: the ability to articulate *why* the model matches what it matches, and *how* that matching changes under controlled perturbations.

This is where governance enters as a methodological necessity rather than an administrative add-on. A professional system is not judged only by performance; it is judged by whether its claims can be reviewed, reproduced, and falsified. In a multimodal setting, that means explicitly instrumenting the representation: logging retrieval@$k$ in both directions, measuring symmetry gaps, tracking embedding variance and mean off-diagonal cosine similarity as collapse indicators, and testing factor separability. The point is not to drown students in metrics. The point is to teach what constitutes evidence about a representation contract. Without evidence, the model is a story. With evidence, the model becomes an object of scientific and professional accountability.

To speak precisely, the professional objective is *reviewable representations, not persuasion*. A persuasive model output is not the same as a correct, stable, and governed system. In a classroom setting, students can be seduced by compelling demos. In professional settings, seduction is dangerous. The output must be reviewable: a third party should be able to run the same notebook, observe the same diagnostics, and validate that the claimed mechanism exists. This is why the companion notebook is structured as an audit-ready laboratory that produces a run manifest, structured logs, plots, and strict JSON reports that separate facts from assumptions and open items.

### 1.1.4 Chapter structure and the notebook contract

This chapter is designed to be read alongside its companion Colab notebook. The reader should not treat the notebook as an appendix. The notebook is the empirical instantiation of the mechanism. The contract is: every conceptual claim in the chapter should be traceable to a reproducible experiment, and every experiment should output artifacts that allow review. The reader is expected to run the notebook end-to-end without modification to reproduce baseline behavior, then to vary specific knobs to stress the mechanism structurally.

The notebook is organized to mirror the conceptual progression. First, it defines a synthetic multimodal world with explicit latent factors and two observation channels. Second, it defines two encoders and trains them with a contrastive objective to produce a shared embedding space. Third, it evaluates alignment through retrieval metrics and geometry probes. Fourth, it stress-tests the representation by sweeping noise, pairing quality, and temperature to map degradation curves. Finally, it exports an audit bundle containing a run manifest, logs, plots, and metrics summaries. This structure is not incidental. It is an implementation of a scientific posture: define the object, train the object, measure the object, stress the object, and record the evidence.

### 1.1.5 What "success" means: geometry, retrieval, and falsifiability

Success in this chapter is not a single high retrieval number. Retrieval is necessary but not sufficient. A model can achieve good retrieval by learning a shortcut that does not reflect the intended factors, or by exploiting a spurious correlation that is stable in the dataset but unstable in production. Therefore, success is defined by a combination of geometric and operational criteria.

First, *geometric success*: the embedding space should reflect the latent factors in interpretable structure. This does not mean there exists a single coordinate axis labeled "orientation." It means that controlled perturbations in orientation induce consistent changes in embedding neighborhoods and distances, and that samples cluster in ways that correspond to factors when visualized through projections such as PCA. It means factor separability scores improve relative to untrained baselines, and the learned representation exhibits meaningful local neighborhoods rather than arbitrary proximity.

Second, *retrieval success*: image→text and text→image retrieval should be strong and symmetric within tolerance. Directional asymmetry is not necessarily a failure, but it must be explainable in terms of modality noise, encoder capacity, or preprocessing. Symmetry is part of the alignment contract: if one direction works but the other fails, the shared space is not truly shared.

Third, *stability success*: the representation should not collapse. Collapse is a geometric failure in which embeddings lose diversity, covariance rank degrades, and mean cosine similarity increases across unrelated samples. Stability is tested not only at the end of training, but across training steps and under perturbations. A representation that is stable only under one precise hyperparameter

setting is not a robust object; it is a brittle artifact.

Fourth, *falsifiability*: the chapter is structured so that claims can be disproven by the same laboratory that supports them. If increasing noise in one modality does not degrade that modality's retrieval direction, we must question whether the modalities are truly coupled or whether the measurement channel is irrelevant. If factor separability does not emerge, we must question whether the objective and architecture are sufficient. If collapse indicators rise while retrieval remains superficially high, we must question whether the representation is being maintained or merely gaming the batch negatives.

In short, success is the student's ability to connect (i) latent factors, (ii) measurement channels, (iii) learning objectives, and (iv) diagnostics into a coherent mechanism. The student should be able to explain why the embedding space looks the way it looks, how it changes under perturbations, and what evidence justifies each claim. Under the broader AI 2026 umbrella, this is the pedagogical endpoint: frontier topics without hype and without mysticism, presented as governed mechanisms that can be reproduced, inspected, and defended.

This introduction therefore frames multimodality as a form of scientific modeling. We specify a world, define measurement channels, write down an objective, train a model, and then test the resulting representation under perturbations. In doing so, we move from "multimodality as capability" to "multimodality as contract." That shift—from capability to contract—is the methodological shift that students need if they are to interact with frontier systems responsibly.

## 1.2 Key Definitions and Mental Model

### 1.2.1 Latent factors, observations, and measurement operators

The objective of this section is to remove ambiguity. Multimodality is a topic where students can speak fluently while holding fuzzy concepts, and fuzzy concepts are incompatible with governed engineering. We therefore define the objects that the rest of the chapter manipulates. The first object is the *latent state* of the world. A latent state is a vector of underlying factors that describe what is present, independent of how it is measured. In our synthetic laboratory these factors are explicit: shape type, orientation, frequency, phase, thickness. In real systems the latent state may include identity, style, lighting, topic, intent, sentiment, or any property that matters to downstream tasks. The crucial point is that "latent" does not mean mystical. It means the variables of interest are not directly provided as structured inputs; they are inferred from observations.

An *observation* is what a modality provides: a matrix of pixel intensities for an image, a sequence of tokens for text, a waveform for audio. The mapping from latent state to observation is the *measurement operator*. In physics this is an instrument model; in multimodal AI it includes both the generative channel (how the world produces pixels or words) and the preprocessing pipeline (resizing, normalization, compression, tokenization). This is not semantic ornamentation. Most practical

failures labeled "drift" are changes in the measurement operator, not changes in the underlying world. A camera pipeline changes exposure or compression. A document pipeline changes OCR or font extraction. A tokenizer changes vocabulary boundaries. The system sees a different observation operator and silently changes its internal geometry.

This chapter adopts the measurement-operator view to prepare students for frontier reality: multimodality is a claim of *invariance*. It claims that meaning should be stable when expressed through different measurement channels. The model is therefore a tool for constructing invariance relations under measurement change. This is why synthetic design is useful: we control the measurement operator and can perturb it deliberately, allowing the student to observe how geometry reacts.

### 1.2.2 Embedding spaces and similarity as operational semantics

An *embedding* is a vector representation produced by an encoder from an observation. The set of all possible embeddings is an *embedding space* or *latent space*. In this book we treat embeddings as geometry, not as compressed text. Once an observation is mapped to a vector, the system operates through dot products, norms, and distances. Similarity functions become operational semantics: if two embeddings are close under the system's similarity metric, the system treats them as related, substitutable, or mutually explanatory.

This implies a critical conceptual shift. The meaning of an embedding is not contained in its coordinates in isolation; it is contained in its *relations* to other embeddings: neighborhoods, clusters, directions of variation, and the structure induced by the similarity metric. The same embedding vector can have different practical meaning if the space around it changes. Therefore, when we talk about "understanding embeddings," we do not mean visualizing raw coordinates. We mean understanding the geometry: which perturbations move points in which directions, whether clusters correspond to factors, and whether neighborhoods remain stable across controlled shocks.

We emphasize *cosine similarity* in the notebook because it operationalizes direction rather than magnitude. But that choice itself is part of the semantics. If cosine similarity is used, normalization is not a minor detail; it is a definitional component of what it means for two items to match. Students must internalize that a multimodal model is not a passive representation; it is a geometric decision system whose semantics are encoded in the metric.

### 1.2.3 Alignment as a geometric contract across modalities

*Alignment* is the property that corresponding observations from different modalities map to nearby points in a shared space. The simplest operational definition is retrieval: given an image, the nearest text embedding should be the paired text; given a text, the nearest image embedding should be the paired image. However, alignment is not merely pair identification. True alignment is a *contract* that the embedding space preserves semantic structure across modalities.

We use "contract" deliberately. A contract has obligations and failure modes. The obligation is: if two observations correspond to the same latent state, they should be close; if they correspond to different latent states, they should be separated. The failure modes are: the model learns shortcuts that satisfy pair identification without learning the intended semantic geometry; one modality dominates the similarity landscape; the embedding space collapses and loses diversity; or the measurement operator changes and the contract silently breaks.

A contract must be testable. The chapter therefore insists on multiple forms of evidence: retrieval metrics, symmetry gaps, separability by latent factors, variance and collapse indicators, and stress-test sweeps. In a governed posture, alignment is never declared by a single metric. It is supported by a set of diagnostics that collectively indicate the contract holds under plausible perturbations.

### 1.2.4 Symmetry and comparability between modality encoders

A multimodal system is often asymmetric by architecture: images may be encoded by convolutional networks or vision transformers, while text may be encoded by transformers or bag-of-token features. Even in simplified notebooks, the two encoders can differ in inductive bias and noise sensitivity. Yet a shared space presupposes comparability: embeddings produced by different encoders must be on compatible scales and must use the shared metric in a way that preserves meaning.

We call this requirement *modality symmetry*. Symmetry does not mean the encoders are identical. It means neither modality is structurally privileged such that the shared space becomes effectively single-modal. Symmetry is tested directionally: image→text retrieval versus text→image retrieval; distributionally: embedding norms and covariance spectra across modalities; and behaviorally: how each direction degrades under asymmetric noise. A robust multimodal system exhibits controlled asymmetry: differences exist but are explainable and bounded.

Symmetry is also a governance issue. In production, teams often monitor global performance and miss directional failures. A model can appear to work because one direction is strong (e.g., text queries retrieving images) while the reverse direction fails, producing systematic user harm in a different workflow. Students should learn that symmetry checks are not optional; they are structural acceptance criteria.

### 1.2.5 Failure pre-view: dominance, shortcuts, collapse

Before building mechanisms, we preview the failure modes that motivate governance. The first is *modality dominance*. If one modality has cleaner signals, higher-dimensional features, or stronger gradients, training can optimize primarily for that modality and treat the other as a weak auxiliary. The symptom is directional asymmetry and hubness: many queries map to the same few candidates, creating a many-to-one geometry.

The second is *spurious alignment* via shortcuts. If both modalities share a feature that correlates

with pair identity but is unrelated to the intended semantics—a watermark, an ID token, formatting artifacts—the model can achieve high retrieval without learning factor geometry. This failure is dangerous because metrics can improve. The correct diagnostic is counterfactual: remove the shortcut and measure the performance drop.

The third is *representation collapse.* Collapse is a geometric degeneracy where embeddings lose diversity. It manifests as rising mean cosine similarity, shrinking covariance rank, reduced effective rank, and degraded neighborhood structure. Collapse can occur due to aggressive optimization, temperature miscalibration, regularization imbalance, or architectural constraints. Collapse can also be induced downstream by bottlenecks or quantization, making it a production risk rather than only a training pathology.

This preview is not pessimism; it is pedagogy. The student should understand that the goal of the chapter is not to present multimodality as magic, but to build a minimal system that reveals how geometry is constructed and how it can fail. These definitions and this mental model are the scaffold that makes the rest of the chapter rigorous: every later claim about embeddings, alignment, or governance will be stated in terms of latent factors, measurement operators, metric geometry, symmetry, and contract failure modes.

## 1.3 Synthetic World Design: Factors, Modalities, and Controls

### 1.3.1 Latent factor design: shape, orientation, frequency, phase, thickness

A synthetic world is only pedagogically valuable if its latent structure is explicit and nontrivial. If the factors are too simple, alignment becomes a toy; if the factors are too complex, the student cannot interpret geometry. We therefore design a factor set that spans multiple representational regimes: discrete categories (shape type), continuous variables (orientation, phase), and mixed discrete-continuous structure (frequency and thickness). This combination forces the shared embedding space to represent both *classification-like* semantics (e.g., "shape A" versus "shape B") and *metric-like* semantics (e.g., orientation distance). In other words, the latent state contains both manifold structure and cluster structure. A competent multimodal representation must accommodate both.

The factor set is also chosen to induce realistic inductive-bias differences between modalities. Images naturally encode orientation and spatial frequency via local structure, while symbolic text can encode these factors only through an engineered grammar. This mismatch is the point. In real multimodal systems, modalities rarely carry information in the same "shape." Images contain rich continuous variation; text contains discrete compositional structure. A shared space must reconcile these differences without allowing one modality to dominate. Our synthetic world makes that tension observable rather than hypothetical.

We also insist on *factor identifiability.* Each factor must be recoverable from each modality (at

least in principle), but not necessarily with equal ease. If a factor is present only in one modality, cross-modal alignment becomes ill-posed: the model can never align on that factor because the other modality lacks the information. Conversely, if all factors are encoded identically in both modalities, the exercise degenerates into learning a redundant mapping. The design goal is balanced asymmetry: each modality carries the same underlying factors, but through different representational channels and noise profiles.

Finally, we treat factor design as a governance choice. In production, teams implicitly choose latent factors by choosing tasks and labels. If the system is trained to retrieve images given text, the latent factors that matter are those that make retrieval succeed. If those factors include confounders, the system will learn confounders. By explicitly defining factors in the synthetic world, we teach the student to see factor choice as an ethical and operational decision: "what should the representation preserve?"

### 1.3.2   Image modality: generative renderer and invariances

The image modality is implemented as small matrices (e.g., $16 \times 16$) that render latent factors into spatial patterns. This is not meant to imitate natural images; it is meant to expose the inductive biases that images naturally support. The renderer must be controllable: for each latent state, we can generate a canonical image and then apply controlled perturbations. Typical rendering primitives include oriented bars, sinusoidal gratings, simple shapes (circles, squares, triangles), and thickness variations. These primitives are sufficiently rich to induce meaningful geometry in pixel space while remaining interpretable.

A key concept is *invariance.* Images naturally support certain invariances: small translations, small rotations, and smooth intensity changes often preserve semantic identity. Our renderer can incorporate or exclude invariances depending on the pedagogical objective. For example, we can design shape identity to be invariant under small orientation changes, while orientation itself remains a factor the model should encode continuously. This teaches the student that "what is invariant" is not a property of the data alone; it is a property of the task and the representation contract.

The renderer also defines the image noise model. Noise is not merely additive Gaussian perturbation. In realistic pipelines, noise includes blur, quantization, compression artifacts, occlusion, and structured distortions. In a synthetic laboratory, we approximate these effects with controlled injections: additive noise, random masks, mild blur via convolution kernels, or structured perturbations that affect certain frequencies more than others. The purpose is to create a family of measurement operators parameterized by noise intensity. This allows us to sweep noise and observe degradation curves in retrieval and geometry.

Importantly, the image renderer is part of the measurement operator. If we change rendering resolution, normalization, or preprocessing, we have changed the operator. In production terms, changing an image resize algorithm can change embedding geometry. In synthetic terms, changing

rendering parameters is an explicit drift knob. We treat this as a first-class object because it is the bridge from classroom understanding to operational reality.

### 1.3.3 Text modality: symbolic grammar, tokenization, feature maps

The text modality is not produced by an LLM. It is produced by a symbolic grammar that encodes the same latent factors as the image. This is essential for two reasons. First, it prevents the student from attributing success to "language understanding." The text channel here is purely structured information. Second, it makes the mapping from latent factors to text explicit and inspectable. In other words, the student can know what information the text carries, and therefore can reason about what the embedding *should* learn.

The grammar is designed to include both discrete tokens (e.g., a token for each shape type) and quantized continuous tokens (e.g., binned orientation or binned phase). This introduces a realistic representational mismatch: the image encodes continuous variables smoothly, while the text encodes them through discretization and order. The student can then observe how this mismatch influences alignment geometry: does the embedding space treat small orientation changes as small distances, or does it treat them as discrete jumps? Does the quantization create artificial clusters? These are not nuisances; they are core lessons about how modality-specific encodings shape the shared space.

Tokenization and feature mapping are treated as explicit engineering choices. In the notebook, text can be represented as bag-of-tokens, bag-of-ngrams, or token-position features. Each choice changes the text measurement operator. A bag-of-tokens representation discards order and treats text as a multiset; a positional representation preserves grammar structure. In production, similar choices arise as tokenization updates, prompt template changes, or the inclusion of metadata fields. The synthetic lab makes these choices visible and measurable.

Finally, we can include controlled text noise: token dropout, token swaps, or grammar corruption. This allows asymmetric noise experiments: degrade text while keeping images fixed, and observe directional retrieval asymmetry. This is one of the clearest demonstrations that multimodal alignment is a coupled system: if one modality loses information, the shared space cannot remain symmetric.

### 1.3.4 Noise models and controllable perturbations

A synthetic world without perturbations is not a laboratory; it is a dataset. The laboratory begins when we define controlled ways to stress the measurement operators and the pairing structure. We organize perturbations into four classes, each corresponding to a distinct production risk.

First, *within-modality noise* alters the observation while leaving latent state unchanged. For images, this may be additive noise, blur, masking, or frequency-dependent distortion. For text, this may be token dropout, token swaps, quantization coarsening, or random insertion of irrelevant tokens. These perturbations test robustness to measurement degradation.

Second, *cross-modal pairing perturbations* alter the correspondence structure. The latent state exists, but the pairing between modalities is partially wrong. In production, this corresponds to mislabeled data, corrupted joins, or asynchronous pipelines that pair the wrong caption to an image. Pairing perturbations are unique because they do not degrade either modality alone; they degrade the coupling. The correct diagnostic is not "image quality" or "text quality" but retrieval inconsistency and training instability under contrastive objectives.

Third, *confounder injection* introduces a shared feature that appears in both modalities and correlates with pair identity but is unrelated to intended semantics. In a synthetic lab, this can be an ID token in text and a watermark pattern in images. Confounders are dangerous because they can increase retrieval while destroying factor geometry. The correct diagnostic is counterfactual: remove the confounder and measure performance drop.

Fourth, *scale and normalization perturbations* alter the embedding mapping indirectly by changing preprocessing or encoder gains. In production, normalization changes can emerge from library updates, quantization, or feature scaling changes. In the lab, we can explicitly scale one modality's embeddings before normalization, or change temperature, and observe dominance or collapse signatures.

These perturbations are not optional extensions. They are the mechanism that turns the chapter into frontier teaching. Frontier risk is rarely "the model cannot learn." Frontier risk is "the model appears to learn, but the representation contract is unstable." Controlled perturbations allow us to measure contract stability directly.

### 1.3.5 Dataset splits and leakage prevention

A governed synthetic laboratory still requires split hygiene. Students often assume leakage is only a concern for real datasets. That is incorrect. Leakage is a concern whenever the experimental design allows information about evaluation to influence training. In synthetic settings, leakage can occur through deterministic generation patterns, reuse of random seeds, or inadvertent coupling of factor assignments across splits. The consequence is that evaluation becomes a reflection of design artifacts rather than representation learning.

We therefore define splits as first-class objects: training, validation, and test. The splits are created deterministically from a global seed, and we assert that indices do not overlap. More importantly, we ensure that any *pairing corruption* or *confounder injection* applied during experiments is applied only within the intended split and is logged as part of the measurement operator configuration. If a confounder is injected into test but not into training, the evaluation answers a different question: generalization under a changed operator. That may be useful, but it must be stated explicitly.

Leakage prevention also includes conceptual hygiene: we must not "tune" based on test results. In a pedagogical notebook, students are tempted to adjust temperature or learning rate until test

retrieval is high. This undermines the lesson. The correct posture is to choose hyperparameters based on validation and to treat test as a final report. The notebook enforces this by logging the selection criterion and by exporting the best checkpoint based on validation metrics rather than test metrics.

Finally, we treat split hygiene as part of governance. A professional system requires that claims be attributable: "This retrieval@1 was achieved under these conditions, with these splits, and under this operator." The audit bundle records these details. Students should leave this section understanding that even in synthetic worlds, experimental integrity is not automatic; it is engineered.

In summary, the synthetic world is not a toy; it is a controlled environment for learning what multimodal geometry is and how it reacts to measurement changes. The factor design creates a rich latent structure, the modalities implement distinct measurement operators, perturbations provide drift knobs, and split hygiene ensures that what we measure is what we claim. This is the foundation for everything that follows: without a governed world, there can be no governed understanding of the embeddings it produces.

## 1.4 Encoders: Two-Layer MLPs as Minimal Function Classes

### 1.4.1 Architecture rationale: capacity, bias, and interpretability

A central design decision in a pedagogical multimodal laboratory is the choice of encoder architectures. If we choose architectures that are too weak, the model cannot represent the mapping and failures become uninformative. If we choose architectures that are too strong, the model becomes difficult to interpret and the student may attribute outcomes to hidden complexity. We therefore use two-layer multilayer perceptrons (MLPs) as a *minimal function class* that remains expressive enough to learn meaningful alignment while being analytically tractable.

The two-layer MLP offers a productive compromise. It can represent nonlinear interactions between input features, allowing the encoder to extract factor-relevant structure from both modalities. Yet it remains simple enough that students can reason about its parameterization: a first affine map to a hidden layer, a nonlinearity, and a second affine map to the embedding dimension. In the notebook, this simplicity enables analytic gradients and finite-difference gradient checking, which is not merely a correctness exercise. It is a governance exercise. A multimodal system that cannot be trusted at the gradient level cannot be trusted at the contract level.

Using two-layer MLPs for both modalities also encourages symmetry. The modalities differ in their input representations, but the encoders share a functional template. This makes it easier to interpret performance differences as arising from the measurement operator and noise profiles rather than from wildly different model classes. When one direction of retrieval degrades, we can ask whether the degradation is due to modality information content, feature mapping, or noise, rather than

blaming an architectural mismatch.

At a deeper level, the two-layer MLP reminds the student that multimodal learning is not fundamentally about exotic architectures. It is about geometry induced by an objective. Even with simple encoders, contrastive objectives can create sophisticated embedding spaces, and those spaces can fail in sophisticated ways. This is a frontier lesson: representational fragility is not a function of architecture complexity alone; it is a function of the incentives created by the objective under finite data and finite negatives.

### 1.4.2 Input normalization and scale control

Input normalization is not a cosmetic preprocessing step; it is part of the measurement operator and therefore part of the semantics of the model. Consider the image modality. If pixel intensities are scaled differently across runs, the encoder sees a different distribution, gradients scale differently, and the resulting embedding geometry can change. In production, changes as small as an update to an image normalization constant can shift retrieval behavior. Therefore, in this chapter, we treat normalization as an explicit component of the experimental contract and we log it as part of the run manifest.

For images, normalization typically includes centering and scaling to a fixed range or variance. For text feature vectors, normalization may include term frequency scaling, clipping, or standardization. The key pedagogical point is that the input scale influences optimization dynamics and can create apparent "dominance" if one modality's inputs produce larger activations and gradients. Without scale control, a modality can dominate simply because its encoder operates in a different numeric regime. This is not an interesting failure mode; it is an engineering artifact. We therefore implement explicit scaling checks and distribution summaries for each modality's input features.

Scale control also includes careful handling of batch statistics. In small synthetic datasets, per-batch normalization can create instability. We prefer deterministic global normalization computed from the training set and then applied consistently to validation and test. This prevents leakage and ensures that evaluation reflects the same measurement operator. The student learns that normalization is a governance obligation: you cannot change it silently without invalidating comparability across runs.

### 1.4.3 Embedding normalization and temperature scaling

Once inputs are encoded, we obtain embeddings. The next design choice is how to compare embeddings. In this chapter, we use cosine similarity, which is implemented by $\ell_2$ normalization of embeddings followed by a dot product. This choice has profound consequences. Cosine similarity emphasizes direction over magnitude. It therefore encourages the model to encode meaning in directions, and it reduces sensitivity to raw embedding scale. However, this is only true if normalization is implemented correctly and consistently across modalities.

Embedding normalization also interacts with *temperature scaling.* In contrastive learning, temperature controls the sharpness of the similarity distribution. A low temperature makes the softmax more peaked, increasing pressure to separate positives from negatives aggressively. This can improve retrieval but can also increase anisotropy and collapse risk. A high temperature smooths the softmax, reducing separation pressure, potentially stabilizing training but also potentially reducing retrieval performance. Temperature is therefore not a mere hyperparameter; it is a structural knob controlling the geometry of competition in the batch.

In a governed posture, temperature must be treated as part of the model contract. The same encoder weights can produce different retrieval behavior under different temperatures, because the similarity matrix is scaled differently. Moreover, temperature interacts with batch size: with more negatives, the required sharpness changes. Students should learn that temperature is a lever on embedding geometry and that its impact must be diagnosed through stress sweeps rather than chosen by intuition.

A further symmetry issue arises: if one modality's embeddings have different norm distributions before normalization, normalization can compress or amplify structure differently across modalities. We therefore compute modality symmetry diagnostics: embedding norm histograms, covariance spectra, and directional retrieval asymmetry. The student learns that "shared space" is not achieved by declaring it; it is achieved by constructing comparable representations and verifying comparability.

### 1.4.4 Initialization and deterministic reproducibility

Initialization defines the starting geometry from which learning proceeds. In contrastive settings, different initializations can lead to different local optima, especially in small or noisy datasets. Therefore, if we want to teach students to reason about geometry rather than to chase accidental outcomes, we must treat initialization as controlled.

We use deterministic random seeds for weight initialization and data ordering. More importantly, we record those seeds and configuration hashes in the audit bundle. This allows exact reproduction of runs. The pedagogical point is that frontier topics require reproducible evidence. A claim about embedding geometry that cannot be reproduced is not a claim; it is an anecdote.

Initialization also has conceptual implications. For example, if initial embeddings are too small or too similar, early training steps can produce large similarity gradients that push the system toward anisotropy. Conversely, if initial embeddings are too large and unnormalized, numeric instability can appear. In a small NumPy-based implementation, these issues are amplified. This is a feature, not a bug. The student sees that geometry is sensitive to numeric regimes, and therefore numeric discipline is part of understanding.

We also encourage the student to view initialization as part of governance: in production, re-training a model without controlling initialization can create drift in embedding space even if the dataset is

unchanged. If embeddings are used as infrastructure, that drift can disrupt downstream systems. Therefore, reproducible initialization is not merely a research convenience; it is a production stability control.

### 1.4.5 Computational complexity and vectorization constraints

Even in a pedagogical notebook, engineering constraints matter. We implement encoders and objectives in NumPy with explicit vectorization to avoid Python loops. This serves two purposes. First, it allows the notebook to run efficiently in Colab, enabling sweeps and diagnostics without excessive runtime. Second, it mirrors the reality that multimodal systems are often bottlenecked by matrix operations, and efficient implementations shape what is feasible.

Computational complexity also influences methodological design. Contrastive objectives scale with batch size because the similarity matrix is $B \times B$. If we increase batch size to improve negative coverage, we increase computation quadratically. Students should understand that in real systems, the negative set is managed through memory banks, queues, or distributed training. In our minimal lab, we cannot replicate those systems, but we can teach the underlying scaling law and show how batch size affects both geometry and compute.

Vectorization constraints also enforce clarity. If we cannot express an operation as a small number of linear algebra operations, we may be implementing a concept incorrectly or inefficiently. The notebook's structure encourages modular functions with clear shapes and explicit assertions. This is not only engineering hygiene; it is interpretability hygiene. When students can see the shapes of tensors and matrices, they can reason about what is being computed and why.

In summary, the encoders in this chapter are deliberately minimal yet rigorous. Two-layer MLPs provide enough expressive power to learn meaningful alignment while remaining interpretable and analytically checkable. Normalization and temperature define the semantics of similarity and must be treated as part of the contract. Initialization and determinism support reproducible geometry. And vectorization constraints keep the laboratory feasible while teaching students the scaling realities that appear at the frontier. This sets the stage for the next section: the contrastive objective that couples the modalities and creates the shared embedding geometry.

## 1.5 Contrastive Alignment Objective

### 1.5.1 InfoNCE as a ranking objective over pairs

The central coupling mechanism in this chapter is a contrastive objective. The purpose of the objective is not to "classify" images or "understand" text in isolation, but to enforce a geometric contract across modalities: paired observations should be more similar than unpaired observations. The most common operationalization is InfoNCE, which can be understood as a differentiable

ranking objective defined over a batch of paired samples.

Let $\{(x_i^{\text{img}}, x_i^{\text{txt}})\}_{i=1}^B$ be a batch of paired observations. Let the encoders produce embeddings

$$z_i^{\text{img}} = f_\theta(x_i^{\text{img}}), \qquad z_i^{\text{txt}} = g_\phi(x_i^{\text{txt}}),$$

and let $\hat{z}$ denote $\ell_2$-normalized embeddings. Define the similarity matrix

$$S_{ij} = \frac{\langle \hat{z}_i^{\text{img}}, \hat{z}_j^{\text{txt}} \rangle}{\tau},$$

where $\tau > 0$ is the temperature. In its symmetric form, InfoNCE trains both directions: image→text and text→image. The image→text loss is

$$\mathcal{L}_{\text{img}\rightarrow\text{txt}} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(S_{ii})}{\sum_{j=1}^B \exp(S_{ij})},$$

and the text→image loss is defined analogously using $S^\top$. The total loss is often the average of the two.

This form reveals the geometric meaning. For each $i$, the model is asked to make the paired text $j = i$ outrank all other texts in the batch under similarity. The batch thus defines a negative set: all other pairs serve as negatives. The objective therefore constructs geometry through competition. In a small batch, the negative set is weak, and the resulting geometry may be under-constrained. In a large batch, the negative set is strong, and the objective creates sharper separation pressure. Students should understand that alignment is not learned "in general"; it is learned relative to the negative distribution induced by batching.

InfoNCE is therefore best understood as an operational surrogate for a mutual information maximization principle: it encourages shared information between the modalities. But this interpretation must be treated carefully. The objective can be satisfied by shared information that is not semantically meaningful (confounders) and can fail even when semantics exist (pair corruption). The rest of the chapter's diagnostics exist precisely because the objective alone does not guarantee the intended contract.

### 1.5.2 Cosine similarity, temperature, and log-sum-exp stability

Implementing InfoNCE correctly requires numerical discipline. The similarity scores can be large in magnitude, especially when temperature is small. The softmax normalization therefore risks overflow if implemented naively. In professional systems, numerical instability is not merely an implementation bug; it can masquerade as a representational phenomenon, such as collapse or dominance, because unstable gradients can distort training dynamics.

We use cosine similarity by normalizing embeddings and taking dot products. This ensures that

similarities lie in $[-1, 1]$ before temperature scaling. Temperature rescales these scores to adjust softmax sharpness. If $\tau$ is small, the effective logits become large, and the softmax becomes extremely peaked, making gradients sparse and potentially unstable. If $\tau$ is large, logits are small, softmax becomes flat, and the model may fail to separate positives from negatives.

To implement the denominator stably, we use the log-sum-exp trick. For each row $i$,

$$\log \sum_{j=1}^{B} \exp(S_{ij}) = m_i + \log \sum_{j=1}^{B} \exp(S_{ij} - m_i),$$

where $m_i = \max_j S_{ij}$. This prevents overflow and produces stable gradients even when similarities are sharp. Students should see that this stability is not optional. If the objective is unstable, the resulting embedding geometry is uninterpretable because it is shaped by numerical artifacts.

Temperature also plays a conceptual role beyond stability. It controls the entropy of the implied matching distribution. Low temperature corresponds to a high-confidence matching policy: the model must concentrate probability mass on the true pair. High temperature corresponds to a softer policy: the model can distribute mass more broadly, potentially preserving diversity but reducing discriminability. This trade-off connects directly to the governance theme: a system that becomes overconfident under drift may appear decisive while being wrong. Temperature is therefore a knob on how the model expresses uncertainty in similarity space.

### 1.5.3 Batch negatives and the geometry of competition

InfoNCE's geometry is batch-dependent. The negative set is the set of mismatched pairs within a batch. This creates a subtle but important consequence: the embedding space is shaped by the distribution of negatives the model sees during training, not by an abstract global notion of "all non-matches." If the batch negatives are too easy, the model may learn superficial separation. If they are too hard or too confounded, the model may exploit shortcuts.

The student should understand two structural effects. First, *negative density* increases with batch size. As batch size grows, each sample competes against more negatives. This increases the chance that some negatives are near-neighbors, forcing the model to learn finer distinctions. But it also increases the risk of interference: semantically similar items can become negatives, encouraging the model to separate items that should arguably be close under the intended factor geometry. In synthetic settings, this can manifest as overly sharp clustering that ignores continuous factor variation. In real settings, it can manifest as brittleness: the model treats near-duplicates as opposites because they co-occurred as negatives.

Second, batch negatives induce *anisotropy pressure*. The model may find it easier to separate by pushing embeddings into a narrow cone of directions and increasing relative angles among a few directions rather than preserving isotropic variance across the sphere. This can create hubness and

degrade interpretability. Monitoring covariance spectra and effective rank helps detect when the space becomes too anisotropic.

A professional point follows: retrieval metrics measured on the same batch structure used in training can be misleading. A model can perform well when negatives match the training distribution but degrade under different negative structures. This motivates stress testing with different batch sizes and negative sampling schemes, even in a minimal notebook. The student learns that geometry is not a static property; it is shaped by competition structure.

### 1.5.4 Analytic gradients for two-layer MLPs

To treat the notebook as a professional instrument, we implement analytic gradients for the two-layer MLP encoders. This serves two purposes. First, it enforces correctness. Second, it teaches the student how the geometry is actually optimized: gradients flow through similarity matrices, softmax normalizations, and nonlinear encoders. Without this understanding, students may treat contrastive learning as an opaque ritual.

At a high level, the gradient of the loss with respect to the similarity matrix $S$ has a simple structure. For the image→text loss, the gradient for row $i$ is

$$\frac{\partial \mathcal{L}}{\partial S_{ij}} = \frac{1}{B}\left(p_{ij} - \mathbf{1}\{j = i\}\right),$$

where $p_{ij} = \frac{\exp(S_{ij})}{\sum_k \exp(S_{ik})}$ is the softmax probability for row $i$. This reveals the ranking nature: the true pair receives negative gradient (increase similarity), and negatives receive positive gradient proportional to how much probability mass they steal.

These gradients then propagate to embeddings through dot products and normalization. The normalization introduces Jacobian terms that couple the embedding coordinates, emphasizing that cosine similarity is not just dot product; it includes a projection onto the unit sphere. Finally, gradients propagate through the two-layer MLP: a first affine transform, a nonlinearity, and a second affine transform. Implementing these gradients carefully in NumPy is an exercise in disciplined linear algebra: explicit shapes, stable computation, and consistent broadcasting.

Why is this in a chapter about multimodality rather than optimization? Because alignment is a contract, and contract enforcement is optimization. If the student does not understand how the objective pushes embeddings around, they cannot interpret why failures occur. For example, dominance can occur because one encoder's gradients are systematically larger. Collapse can occur because temperature makes gradients too sharp and pushes embeddings into degenerate configurations. Analytic gradients make these explanations concrete.

### 1.5.5 Gradient checking and correctness obligations

A professional notebook must not merely implement gradients; it must validate them. Finite-difference gradient checks provide a local correctness test: perturb a parameter slightly, recompute loss, and compare the numerical gradient to the analytic gradient. In a small synthetic setting, we can perform these checks on a tiny batch and a subset of parameters. The goal is not to prove global correctness, but to enforce a minimal correctness obligation: the gradients we use are consistent with the objective we claim to optimize.

This matters for pedagogy because gradient errors can produce apparent learning behavior that is misleading. A sign error can look like "collapse" or "instability." A broadcasting mistake can silently train only a subset of parameters, producing spurious asymmetry. Without gradient checking, these failures become interpretability traps: students may interpret them as conceptual phenomena rather than implementation bugs.

Gradient checking also ties directly to governance. In professional environments, one of the most common sources of model risk is silent implementation error. A governed system therefore uses checks, assertions, and unit-test-like validations. In our notebook, gradient checks play that role. When the check fails, the run should fail loudly, and the artifacts should record the failure. This is the posture we want students to adopt: correctness is not assumed; it is enforced.

In summary, the contrastive objective is the coupling mechanism that creates the shared embedding geometry. InfoNCE imposes a ranking contract through batch competition; cosine similarity and temperature define similarity semantics; stable log-sum-exp implementation ensures numeric correctness; batch negatives shape geometry through competition density; analytic gradients and gradient checks enforce correctness and enable interpretable training dynamics. With this objective in place, we can now describe the training protocol and instrumentation that transform the objective from a formula into a governed laboratory.

## 1.6 Training Protocol and Instrumentation

### 1.6.1 Optimization loop and learning-rate discipline

A contrastive objective specifies what the model should optimize, but the training protocol determines what the model actually becomes. In small synthetic settings, optimization dynamics are especially sensitive: learning rates that are "fine" in large frameworks can be unstable when implemented directly in NumPy, and subtle instabilities can masquerade as representational phenomena. Therefore, we treat the optimization loop not as boilerplate but as part of the scientific instrument.

The training loop in this chapter is intentionally explicit. Each step computes forward embeddings for both modalities, constructs the similarity matrix, evaluates the symmetric InfoNCE loss, computes analytic gradients for both encoders, and applies parameter updates. The loop is deterministic:

batching order is fixed by a seeded permutation, and all sources of randomness are recorded. The learning rate is treated as a structural knob. Too large, and the geometry can oscillate or collapse; too small, and the model may fail to separate positives and negatives, producing a deceptively "smooth" but underfit representation.

Learning-rate discipline is also a pedagogical lesson about frontier systems. Many real-world multimodal failures are not failures of architecture but failures of training stability. Overly aggressive optimization can produce anisotropy, hubness, or collapse; overly conservative optimization can produce weak alignment that looks acceptable on superficial tests but fails under stress. The correct posture is not to hunt for a single learning rate that "works," but to instrument training so that instability is detectable early. This is why we monitor not only loss, but also embedding variance, mean cosine similarity, retrieval@$k$, and gradient norms. Training is treated as a controlled process, not as a single outcome.

Finally, the optimization loop is governed by acceptance conditions. If NaNs appear, if variance collapses below a floor, or if gradient norms explode, the run should fail loudly and record the failure in artifacts. This is professional engineering: a model that silently trains into an unstable state is unacceptable even if it later produces a pretty plot.

### 1.6.2 Metrics logging: loss, retrieval@k, symmetry, variance

A multimodal system is a geometric object, and geometric objects require geometric diagnostics. Logging only the loss is insufficient because loss can decrease while semantic geometry degrades. Therefore, we log a panel of metrics at a frequency that balances visibility and compute cost.

The first metric is the loss itself, both symmetric total loss and directional components (image→text and text→image). This directional decomposition is important because it allows the student to see whether one direction dominates optimization. If one direction improves while the other stagnates, symmetry is already being violated.

The second metric is retrieval@$k$ in both directions. Retrieval@$1$ provides a strict correctness measure; retrieval@$5$ provides a looser measure that reflects neighborhood quality. Directional retrieval asymmetry is a primary indicator of modality imbalance and of asymmetric noise effects. The point is not to treat retrieval as the only goal; the point is to treat retrieval as an operational proxy for whether the contract is being enforced.

The third set of metrics concerns *embedding variance and distribution.* We log per-modality embedding norm statistics (before normalization), embedding coordinate variance, and average pairwise cosine similarity (often approximated efficiently). A rise in mean off-diagonal cosine similarity indicates collapse pressure. A sharp reduction in variance indicates loss of representational capacity. Distributional symmetry between modalities indicates whether the shared space is being used comparably.

Finally, we log auxiliary metrics that support interpretability: separability proxies by latent factor on a validation sample, and covariance spectrum summaries such as effective rank. These metrics are more expensive, so we compute them periodically rather than every step. Their purpose is to prevent false comfort: a model can show good retrieval while embedding geometry becomes degenerate. Periodic structural metrics act as early warnings.

All metrics are recorded as structured logs rather than printed only to the console. Console output is ephemeral; structured logs are auditable. This is why the notebook exports JSON summaries and retains per-step histories in deliverables. The goal is to create evidence that can be reviewed and compared across runs.

### 1.6.3 Early warnings: collapse indicators and instability tests

A core governance requirement is the ability to detect failure early, before a model is "declared trained." In multimodal contrastive learning, two families of early warnings are especially important: collapse and instability.

Collapse warnings include a sustained increase in mean off-diagonal cosine similarity, a sustained decrease in embedding variance, and a contraction of the covariance spectrum (reduced effective rank). Collapse can occur even when loss decreases because the model can satisfy the objective partially by making all embeddings similar and relying on small fluctuations to separate pairs. This may produce a brittle representation that fails under stress. Therefore, collapse indicators are treated as constraints with floors and thresholds.

Instability warnings include gradient norm explosions, oscillatory loss behavior, and sudden changes in retrieval without corresponding changes in loss. Instability is not only a training inconvenience; it is a sign that the representation contract is not being learned robustly. In a controlled lab, we can adjust learning rate or temperature and observe how instability responds. The student learns the causal structure: low temperature can create sharper gradients and instability; high learning rate can amplify it; insufficient normalization can exacerbate it.

Early warning tests also include determinism checks. If rerunning the same configuration yields materially different results, the experiment is not governed. Determinism checks are therefore built into the notebook: fixed seeds, fixed shuffles, and recorded environment fingerprints.

Pedagogically, early warnings teach students to treat training as an observable process. Rather than waiting for final plots, they learn to watch how geometry evolves. This is a frontier mindset: the behavior of a system over time contains information about what it is learning. A system that converges smoothly may be learning a shortcut; a system that oscillates may be too sensitive to noise. The student learns to interpret these patterns as evidence rather than as annoyances.

### 1.6.4   Checkpointing and "best model" selection

Checkpointing is a governance tool. In many educational notebooks, the final model is used by default. This is operationally dangerous because the final model may be worse than earlier states due to overfitting, instability, or late-stage collapse. Therefore, we checkpoint models periodically and select a "best" checkpoint based on validation metrics.

The selection criterion is explicitly defined. For example, we might select the checkpoint that maximizes the average of image→text and text→image retrieval@1 on validation, subject to collapse constraints (variance floors, effective rank floors). This makes selection a constrained optimization problem rather than a naive metric maximization. The student learns that professional selection criteria encode governance: you do not accept a model that achieves higher retrieval by collapsing geometry or by becoming severely asymmetric.

Checkpointing also enables forensic analysis. If later stress tests show a failure, we can compare checkpoints to see whether the failure was present early or emerged late. This supports the chapter's emphasis on interpretability and contract stability. It also mirrors production practice: systems often need rollback options. A checkpointed model is a rollback artifact.

Finally, checkpointing is recorded in the audit bundle. Each checkpoint is saved with metadata: configuration hash, training step, and validation metrics at the time of saving. This ensures that a claim such as "we used the best validation checkpoint" is verifiable.

### 1.6.5   Determinism audits: seeds, ordering, and regression checks

Determinism is a prerequisite for governance. If a model's geometry changes unpredictably across identical runs, monitoring and drift detection become unreliable. Therefore, the notebook includes determinism audits that verify repeatability at multiple levels.

First, we fix seeds for all random number generators used in data construction, batching, and initialization. Second, we fix data ordering by using seeded permutations. Third, we verify that key intermediate outputs (such as a checksum of the first batch embeddings) match expected values when rerun. These checks are not meant to be foolproof; they are meant to catch common sources of non-determinism such as accidental reseeding, shape-dependent randomness, or hidden state.

Regression checks can also be included at a minimal level: run a small number of training steps and ensure that the loss decreases within a known range and that retrieval@1 exceeds a minimal threshold. This is not a benchmark; it is a sanity guard. The purpose is to detect when code changes have broken the objective or the gradient flow.

Determinism audits reinforce the professional posture: a governed multimodal system is not merely one that performs well; it is one whose behavior is stable under the conditions that matter. In a synthetic notebook, we can enforce this stability strongly. In production, we cannot guarantee

full determinism, but the discipline remains: control what you can, record what you cannot, and monitor distributions rather than anecdotes.

In summary, the training protocol in this chapter is treated as an instrument for learning and for governance. The optimization loop is explicit and disciplined, metrics logging goes beyond loss to capture geometry, early warnings detect collapse and instability, checkpointing supports constrained selection and rollback, and determinism audits enforce reproducibility. With training governed in this way, we can move to the next stage: evaluating the learned embedding space as geometry and interpreting what it has learned.

## 1.7 Geometric Diagnostics and Interpretability

### 1.7.1 PCA projections and factor-separated visualization

If multimodality is "learning compatible coordinate systems," then interpretability begins with seeing whether the learned space exhibits structure that corresponds to latent factors. The first tool we use is projection: we take high-dimensional embeddings and map them into a low-dimensional view using principal component analysis (PCA). The goal is not to claim that PCA reveals "the truth" of a representation. The goal is to create an interpretable window into the geometry so that students can reason about clusters, directions, and overlap.

In a synthetic laboratory, factor-separated visualization is especially powerful because we know the true factors for each sample. We can color points by shape type, orientation bin, frequency bin, phase bin, or thickness. We then ask a precise question: does the embedding space organize samples so that factor labels correspond to visible structure? For discrete factors like shape type, we expect cluster-like separations. For continuous factors like orientation, we may expect smooth gradients or curved manifolds. For mixed factors, we may see clusters with internal continuous variation. These patterns are not guaranteed, and that is the point. The student learns that the objective, architecture, and noise determine how factors are represented.

PCA also teaches a subtle lesson about interpretability. The directions PCA finds are directions of maximal variance, not necessarily directions of semantic importance. A factor can be semantically critical but low variance, and therefore invisible in early PCA components. Conversely, a nuisance factor can dominate variance and appear prominent. This is why PCA is a starting point rather than an endpoint. We complement it with quantitative separability scores and controlled perturbation tests. The student learns to treat visualization as evidence, but not as proof.

Finally, PCA projections become a baseline artifact. We save plots for the baseline model and compare them under stress sweeps. If a drift perturbation collapses geometry, the PCA projection will often show contraction into a narrow region or hubness patterns. If a confounder dominates, PCA may show separation that corresponds to the confounder rather than the intended factor. These

visual shifts are not merely illustrative; they provide quick qualitative validation that complements numeric metrics.

### 1.7.2 Similarity heatmaps and neighborhood structure

A shared embedding space is operationally used through similarity computations. Therefore, interpretability must include understanding the similarity matrix structure. A similarity heatmap is a visualization of $S_{ij}$ values between a set of image embeddings and text embeddings (or within-modality embeddings). When the model is well-aligned, we expect a strong diagonal structure in the cross-modal similarity matrix: each image matches its paired text most strongly.

Heatmaps also reveal failure signatures that can be missed by aggregate metrics. For example, *hubness* appears as vertical or horizontal stripes: one text embedding is highly similar to many images, or one image embedding is highly similar to many texts. This can happen even when top-1 retrieval is moderate, and it indicates a many-to-one geometry that is operationally dangerous. A user may experience repetitive retrieval results, and downstream systems may become biased toward hubs.

Heatmaps also reveal *confounding*. If a confounder creates a shared shortcut, the diagonal can become strong even when semantic structure is wrong. In synthetic settings, we can reorder samples by factor and see whether blocks appear. If blocks correspond to confounder groups rather than semantic factors, the model is aligning for the wrong reason. This is why heatmaps are valuable in a governed lab: they provide interpretable signatures that can trigger deeper tests.

Neighborhood structure is a complementary view. For each embedding, we can compute the nearest neighbors in the other modality and inspect whether neighbor factor labels make sense. In synthetic labs, we can quantify neighbor consistency: the fraction of nearest neighbors sharing shape type, or the average orientation difference among neighbors. This moves interpretability from "looks good" to measurable neighborhood semantics.

### 1.7.3 Separability scores: ANOVA-like ratios by factor

Visualization is important but insufficient. We therefore compute quantitative measures of factor separability in embedding space. A principled approach is to measure how much variance in embeddings is explained by factor groupings. For a discrete factor (e.g., shape type), we can compute an ANOVA-like ratio: the between-group dispersion relative to within-group dispersion.

Concretely, let $\mu$ be the global mean embedding, and let $\mu_c$ be the mean embedding of group $c$. The between-group scatter is

$$S_B = \sum_c n_c \|\mu_c - \mu\|^2,$$

and the within-group scatter is

$$S_W = \sum_c \sum_{i \in c} \|z_i - \mu_c\|^2.$$

The ratio $S_B/(S_W + \varepsilon)$ is a separability proxy. A higher ratio indicates that group means are well-separated relative to within-group spread. We compute such ratios for each factor (shape, frequency bin, thickness bin) and for discretized versions of continuous factors (orientation bins, phase bins). This yields a factor-by-factor profile of what the embedding space encodes.

The pedagogical benefit is immediate: students can see that a space may encode shape strongly but orientation weakly, or encode frequency strongly but thickness weakly. They learn that "alignment" is not uniform across semantics. Moreover, separability can improve even when retrieval remains constant, showing that factor geometry can evolve independently of pair identification. Conversely, retrieval can improve while separability degrades, indicating shortcut exploitation. This is precisely why separability metrics are governance-relevant.

Separability metrics also help define acceptance criteria. A professional system may require that certain factors remain encoded above a threshold. For example, if shape identity is critical to a downstream task, we can enforce a minimum separability ratio for shape. In production, such factors are domain-specific; in the lab, we make them explicit.

### 1.7.4 Collapse scores: mean cosine and variance floors

Representation collapse is a geometric degeneracy that must be detected reliably. We therefore compute collapse indicators that are both simple and interpretable. The first is the mean off-diagonal cosine similarity. For a set of normalized embeddings $\{\hat{z}_i\}$, compute the average cosine between different samples:

$$C = \frac{1}{B(B-1)} \sum_{i \neq j} \langle \hat{z}_i, \hat{z}_j \rangle.$$

If $C$ becomes large and positive, embeddings are becoming more aligned with each other, reducing diversity. In a healthy space, unrelated items should have near-zero average cosine (or mildly positive depending on anisotropy), not strongly positive.

The second indicator is the embedding variance floor. We compute coordinate-wise variance across samples and track the mean variance. If variance collapses toward zero, embeddings lose capacity to represent differences. We also compute covariance spectrum summaries: the singular values of the centered embedding matrix. A collapse manifests as most singular values shrinking, reducing effective rank. Effective rank can be computed as

$$r_{\text{eff}} = \exp\left(-\sum_k p_k \log p_k\right),$$

where $p_k$ is the normalized singular value mass. This captures how many directions carry significant

variance.

These collapse indicators serve multiple roles. They are early warnings during training, and they are acceptance criteria during evaluation. They are also diagnostic under stress tests: if lowering temperature increases retrieval but collapses effective rank, we have learned a trade-off that must be governed. The student learns that collapse is not always obvious from retrieval; it is a geometric pathology that can coexist with good headline metrics.

### 1.7.5 Modality symmetry checks: distribution and norm alignment

A shared embedding space must remain shared in a quantitative sense. We therefore compute symmetry checks that compare the distributions of embeddings across modalities. These include:

- *Norm distributions* (pre-normalization): compare mean and variance of embedding norms for image and text encoders. Persistent differences can indicate imbalance in encoder scaling or information content.
- *Covariance spectra*: compare singular value profiles of image embeddings and text embeddings. If one modality has much lower effective rank, it may be under-represented in the shared space.
- *Directional retrieval asymmetry*: compare image→text retrieval and text→image retrieval. Significant gaps indicate dominance or asymmetric noise sensitivity.
- *Neighborhood consistency asymmetry*: compare how factor-consistent neighbors are in each direction.

These checks embody the chapter's contract view: multimodality is not achieved by forcing two encoders into the same vector dimension; it is achieved by maintaining comparability under a shared metric. In a governed posture, symmetry checks are not optional. They are the difference between a system that works for one workflow and a system that works robustly across workflows.

Finally, symmetry checks connect directly to later chapters. In Chapter 2, we study failure modes such as dominance and spurious shortcuts. Symmetry checks are the monitors that reveal those modes. In Chapter 3, we study drift and monitoring. Symmetry checks become stage-gate metrics. Therefore, the student should not treat these diagnostics as "extra plots." They are the instrumentation layer that makes frontier systems professional.

In summary, geometric diagnostics and interpretability are treated as the primary evidence that multimodal alignment has been achieved for the right reasons. PCA projections and heatmaps provide interpretable signatures, separability scores quantify factor encoding, collapse indicators prevent degeneracy from hiding behind retrieval, and symmetry checks enforce that the shared space remains truly shared. With these tools, the chapter can now turn to the methodological posture that defines professional science: stress testing.

## 1.8 Stress Testing as Scientific Method

### 1.8.1 Noise-asymmetry sweeps and degradation curves

A multimodal representation contract is meaningful only if it survives plausible perturbations. In professional settings, those perturbations occur continuously: sensor quality changes, preprocessing pipelines evolve, text distributions shift, and one modality may degrade while the other remains stable. Therefore, stress testing is not an optional "robustness add-on." It is the scientific method for multimodal systems. We treat stress tests as controlled interventions on the measurement operators and we measure how the embedding space and retrieval behavior degrade.

The first stress family is *noise asymmetry*. We add increasing noise to one modality while holding the other fixed, and we measure directional retrieval degradation. If text becomes noisier, we expect text→image retrieval to degrade faster than image→text, because queries from the degraded modality are less informative. Conversely, if images become noisier, image→text should degrade faster. This asymmetry is not a nuisance; it is diagnostic. It provides a structural explanation of how information content differs across modalities and how the shared space behaves when that balance changes.

The output of a noise-asymmetry sweep is a degradation curve: retrieval@$k$ (and auxiliary geometry metrics such as effective rank and mean cosine) as a function of noise intensity. These curves are more informative than point metrics because they reveal *margins*. A representation that degrades smoothly and predictably is more governable than a representation that remains strong until a sudden cliff. In production, cliff behavior is dangerous because it produces unexpected failure in new operating regimes. The student therefore learns to treat degradation curves as stage-gate evidence: we do not ask only whether the model works; we ask how far it can be pushed before it fails.

Noise sweeps also reveal interactions between the objective and the measurement operator. For example, if noise in one modality causes the model to increase hubness (many queries map to a few candidates), the failure is not merely "lower accuracy." It is a geometric contraction that changes user-visible behavior in specific ways. Stress testing therefore links the abstract geometry to operational consequences.

### 1.8.2 Pairing perturbations and robustness margins

The second stress family targets the coupling assumption: that paired observations correspond to the same latent state. In real multimodal pipelines, pairing errors are common: mislabeled captions, incorrect joins across databases, asynchronous content updates, or data ingestion bugs. Contrastive learning is particularly sensitive to pairing quality because the objective assumes that the positive pair is true.

We therefore define a pairing corruption parameter $\rho \in [0, 1]$ representing the fraction of pairs that

are swapped or mispaired. We then train or evaluate under controlled corruption and measure how retrieval, separability, and geometry respond. Pairing perturbations reveal a key conceptual point: multimodal learning is not only about representing factors; it is about enforcing correspondence structure. When correspondence is wrong, the objective optimizes the wrong contract, and the embedding space can become inconsistent or unstable.

Robustness margins are defined as the maximum corruption level at which the system maintains acceptable retrieval and geometric health. In production, such margins correspond to data quality tolerances. A governed system should be designed with explicit tolerances: "we remain stable up to $\rho = 0.05$ corruption," for example. The synthetic lab allows students to observe whether such tolerances are plausible.

Pairing perturbations also teach students to interpret failure signatures. Under corruption, loss may not behave predictably because the objective is being asked to satisfy contradictory constraints. Retrieval can degrade asymmetrically or exhibit non-monotonic behavior. Geometry can show hubness or collapse as the model attempts to resolve contradictions. The student learns that not all performance drops are equal; a drop due to noise is different from a drop due to corrupted correspondence, and monitoring must distinguish them.

### 1.8.3 Temperature sensitivity and embedding anisotropy

Temperature is a structural parameter in contrastive learning. It controls similarity sharpness and therefore the entropy of the implied matching distribution. Stress testing must therefore include temperature sweeps. The goal is not to tune for maximum retrieval; the goal is to understand how temperature changes geometry and failure risk.

At low temperature, the objective pressures embeddings to create very high similarity for true pairs and to strongly suppress negatives. This can improve retrieval but can also increase anisotropy: embeddings may cluster into a narrow cone or create hubs. Collapse indicators may rise. At high temperature, the objective is softer; geometry may remain more diverse but retrieval may weaken. The correct interpretation is that temperature controls a trade-off between discriminability and diversity.

Temperature sensitivity is governance-relevant because production systems often change temperature implicitly. For example, a system might change scaling in similarity computation, adopt a different normalization, or apply quantization that effectively changes similarity sharpness. If the model is temperature-sensitive, these changes can induce drift-like behavior even if the model weights are unchanged. Therefore, the chapter treats temperature sweeps as a simulation of production perturbations: we test whether the contract holds under reasonable similarity scaling changes.

We also measure anisotropy through covariance spectra, effective rank, and mean cosine similarity. If lowering temperature improves retrieval but collapses effective rank, we record this as an explicit

trade-off rather than celebrating the retrieval increase. Students learn the discipline: do not optimize a single metric at the expense of representational health unless you can justify the trade in terms of the downstream application and control mechanisms.

### 1.8.4 Batch size effects and negative-set density

Batch size is not a mere compute parameter; it defines the negative set in InfoNCE, and therefore changes the geometry learned. A larger batch increases negative density, pushing the model to separate more aggressively. This can improve retrieval but can also create interference: semantically similar items become negatives, encouraging the model to separate items that should remain close under factor semantics.

We therefore include batch size as a stress dimension. We run training or evaluation under different batch sizes and observe how retrieval and geometry change. In synthetic settings, batch size effects are instructive because we can control the diversity of latent factors within each batch. If a batch contains many near-duplicates (in factor space), negatives are "hard" and the model must learn fine distinctions. If a batch contains diverse factors, negatives are "easy" and the model may learn superficial separation.

Students learn two frontier-relevant lessons. First, the negative set is a design choice. In large systems, it is managed by sampling strategies and distributed training. Second, "better retrieval" from larger batches may come with a cost: increased anisotropy and reduced factor smoothness. This is why batch size is not only a performance knob; it is a geometry knob.

### 1.8.5 Acceptance criteria: distributions, not point estimates

The final methodological principle is that acceptance must be defined over distributions. A single run can produce a favorable metric by chance, by initialization, or by a transient instability. Professional governance requires that we characterize behavior across seeds, across stress levels, and across plausible perturbations. Therefore, we treat acceptance criteria as distributional stage gates.

Concretely, we define a set of thresholds and margins: minimum retrieval@1 and retrieval@5 in both directions on a test set; maximum allowable symmetry gap; minimum effective rank; maximum mean off-diagonal cosine; minimum separability for critical factors. But we do not check these thresholds at a single operating point. We check them across sweeps. For example, retrieval@1 must remain above a threshold up to a certain noise level; effective rank must remain above a floor across temperature variations; confounder sensitivity must remain below a bound when confounders are absent.

This shifts evaluation from persuasion to evidence. A model that passes only at one point is not robust. A model that passes across a range is governable. The student learns that robustness is not a claim; it is a curve.

Finally, acceptance criteria must be recorded as artifacts. The notebook exports metrics summaries and stress reports as strict JSON. These reports separate what was measured (facts) from what is concluded (analysis) and what remains uncertain (open items). The verification status remains "Not verified" because a synthetic lab does not certify real-world performance. This posture is itself a frontier lesson: the correct stance toward representation learning is disciplined uncertainty coupled with reproducible evidence.

In summary, stress testing is the scientific method for multimodal systems. Noise-asymmetry sweeps reveal information balance and directional robustness, pairing perturbations test the coupling contract, temperature sweeps expose discriminability-diversity trade-offs and anisotropy risk, batch size explores negative-set geometry, and acceptance criteria are defined distributionally rather than as point estimates. With stress testing defined, we can now specify the governance artifact standard that makes these experiments professionally reviewable.

## 1.9   Governance Artifacts and Deliverable Standard

### 1.9.1   Run manifest: configuration hashes and environment fingerprint

A chapter that teaches frontier topics must also teach frontier discipline. In practice, the most common failure in advanced work is not a lack of ideas; it is a lack of traceability. Students and practitioners run experiments, obtain results, and then cannot reliably answer what exactly produced them. This is catastrophic when embeddings become infrastructure, because any downstream behavior change demands forensic reconstruction. Therefore, every run of the companion notebook produces a *run manifest*: a structured record of what was executed.

The run manifest serves as the notebook's primary provenance artifact. It includes (at minimum) a timezone-aware UTC timestamp, a deterministic run identifier, and a complete serialization of the configuration used for data generation, model architecture, training hyperparameters, and evaluation settings. It also includes hashes. The configuration is hashed so that two runs can be compared at the level of exact parameter equality. If code versioning is available, code hashes can be recorded as well. In addition, the environment is fingerprinted: Python version, NumPy version, and any relevant system information that can influence numerical behavior.

This is not bureaucracy. It is the minimum infrastructure for scientific accountability. If a student later asks why a retrieval curve differs from a previous run, the correct answer is not speculation. The correct answer is: compare manifests, identify differences, and rerun if needed. The manifest is the bridge from a notebook experience to production-grade behavior, where every model run must be attributable.

### 1.9.2 Prompts log: redaction, hashing, and provenance

Even in synthetic notebooks, there is value in logging the "prompt" layer: the narrative constraints and instructions that define what the system is supposed to do. In this chapter, prompts are not natural language prompts to an external model; they are the structured experiment directives: "train a shared space with InfoNCE," "sweep noise asymmetry," "compute separability by factor," and so on. Logging them is part of pedagogy: students learn that intentions are part of evidence.

The prompts log is stored as newline-delimited JSON (JSONL), where each entry contains: a timestamp, a short description of the step, and a redacted text field (if any human-readable content is present), plus hashes of the full content. Redaction is included even when the content is harmless because the method should generalize to professional settings where prompts may contain sensitive information. The point is to teach the workflow, not to solve a specific privacy case. Hashes allow a reviewer to confirm consistency without exposing content.

The provenance role is subtle but important. In professional teams, the biggest source of disagreement is often that two people think they ran "the same" experiment but actually ran different experiment specifications. A prompt log provides a second dimension of comparability beyond numeric configuration: it records what was intended, when, and in what sequence. In a governed posture, results are tied not only to parameters but also to procedural steps.

### 1.9.3 Risk log: taxonomy, controls, and open items

A risk log is a formal record of what can go wrong, why it matters, and what controls are in place. In a multimodal setting, the risk log is not optional because the failure modes are often silent. High retrieval can coexist with spurious alignment. Symmetric performance can coexist with collapse onset. Seemingly stable geometry can hide confounding features. Therefore, the notebook produces a risk log each run that includes a taxonomy of risks and an explicit mapping to controls implemented in the code.

The risk taxonomy for this chapter includes at least:

- *Spurious correlation and shortcut learning*: shared artifacts (IDs, watermarks) dominate alignment.
- *Modality dominance*: one modality drives similarity structure, producing asymmetry and hubness.
- *Representation collapse*: embeddings lose diversity; covariance rank contracts; neighborhoods degrade.
- *Train–test leakage*: split hygiene fails; evaluation reflects design artifacts rather than learning.
- *Metric hacking*: optimization or selection overfits to a single metric (e.g., retrieval@1) while degrading geometry.
- *Numerical instability*: overflow/underflow in softmax; NaNs; gradient explosions.

Controls implemented include: deterministic seeds and shuffles; explicit split assertions; stable log-sum-exp; gradient checks; symmetry metrics; separability probes; collapse indicators; stress sweeps with degradation curves; and constrained checkpoint selection. Importantly, the risk log also includes *open items*. These are questions that the notebook cannot answer because the lab is synthetic. For example: "Will these acceptance margins transfer to real image-text data?" or "How will a transformer encoder change anisotropy behavior?" The verification status is explicitly recorded as "Not verified." This is a governance posture: the notebook provides evidence about a mechanism, not certification of real-world performance.

Pedagogically, the risk log is a teaching device. It trains students to articulate failure modes before they occur, to pair them with controls, and to distinguish evidence from confidence. This is exactly the mindset needed for frontier systems.

### 1.9.4   Deliverables folder: plots, metrics summaries, checkpoints

The deliverables folder is the concrete output surface of the notebook. It contains plots, JSON summaries, and model checkpoints. The key requirement is that every plot shown to the student is saved to disk with a deterministic name and recorded in the metrics summary. The deliverables folder is not an afterthought; it is the evidence package.

Typical deliverables include:

- *Plots*: PCA projections colored by factor; similarity heatmaps; degradation curves under noise sweeps; covariance spectrum plots; symmetry gap plots.
- *Metrics summaries*: a strict JSON file that records retrieval@1 and retrieval@5 in both directions, symmetry gaps, collapse indicators, separability scores by factor, effective rank, and stress-test results.
- *Checkpoints*: saved parameter snapshots for the "best" model under the defined selection criteria, and optionally periodic checkpoints for forensic analysis.
- *Stress reports*: JSON files describing sweep configurations and resulting degradation curves.

The deliverables standard enforces a professional habit: do not trust what you cannot store and compare. A plot that exists only in a notebook output cell is not a reliable artifact. A metric printed to console is not a reliable artifact. By saving deliverables systematically, the notebook becomes a reproducible laboratory rather than a one-time demo.

### 1.9.5   Verification posture: "Not verified" by construction

A final governance requirement is epistemic honesty. This notebook is a synthetic laboratory. It demonstrates mechanisms, not guarantees. Therefore, every report produced by the notebook uses a strict JSON schema that explicitly separates:

- **facts_provided**: measured metrics and saved artifact paths,
- **assumptions**: explicit modeling assumptions (e.g., factor encoding sufficiency, noise models),
- **open_items**: questions that remain unresolved or require human review,
- **analysis**: interpretation of facts in light of the mental model,
- **draft_output**: summary statements intended for communication,
- **verification_status**: fixed as `"Not verified"`,
- **questions_to_verify**: explicit next-step questions for real-world translation.

This posture is pedagogically central. Frontier topics are dangerous when they are taught as certainty. The correct student outcome is not blind confidence in multimodal embeddings; it is disciplined ability to reason about them with evidence and to know what has not been proven. "Not verified" is not a weakness. It is a professional safeguard: it prevents synthetic demonstrations from being mistaken for production validation.

In summary, the governance artifacts are the chapter's institutional backbone. The run manifest ensures traceability, the prompts log ensures procedural provenance, the risk log forces explicit recognition of failure modes and controls, the deliverables folder contains auditable outputs, and the verification posture enforces epistemic discipline. Together, these artifacts transform the chapter from a conceptual essay into a supervised laboratory that students can run, inspect, and defend. The final section will then consolidate what was learned and transition to Chapter 2, where we study failure modes explicitly.

## 1.10 Conclusion and Transition to Chapter 2

Multimodality is often treated as a capability label, but this chapter has argued for a stricter interpretation: multimodality is the learning of *compatible coordinate systems* under imperfect measurement. Once observations are embedded as vectors, the system's behavior is governed by geometry. Similarity is not a rhetorical notion; it is an operational rule that determines retrieval, ranking, clustering, and downstream decisions. Under this view, the core question is not whether a model can accept multiple input types, but whether it can maintain a stable representation contract: corresponding observations from different modalities remain near each other for the *right reasons*, across noise, preprocessing shifts, and realistic perturbations.

The synthetic-first methodology was chosen to make this contract visible and testable. By constructing a multimodal world with explicit latent factors (shape type, orientation, frequency, phase, thickness) and two measurement channels (small images and symbolic text), we gained what real-world datasets rarely provide: controlled semantics. This allowed us to treat embeddings as evidence-bearing objects rather than mysterious vectors. Instead of guessing what the model might be encoding, we could measure factor separability, inspect neighborhoods, and visualize geometry via projections. The central pedagogical benefit is that students learn to interpret embeddings

through controlled interventions: change one factor, hold others fixed, and observe how distances and neighborhoods respond. This is interpretability as experimental method, not interpretability as narrative.

A second lesson is that alignment is induced by incentives, and incentives can be satisfied in fragile ways. Contrastive objectives such as InfoNCE enforce cross-modal ranking within a batch. This creates a shared space through competition against batch negatives, but the resulting geometry depends on structural knobs: temperature, batch size, normalization, and encoder capacity. The chapter emphasized that these knobs are not mere hyperparameters. Temperature controls similarity sharpness and the entropy of matching; batch size controls negative density and therefore the geometry of competition; normalization defines the semantics of cosine similarity; encoder capacity defines which invariances and nonlinearities can be represented. Students should leave with the ability to predict qualitative effects: lower temperature typically sharpens retrieval but increases anisotropy risk; larger batches can improve discriminability but can also intensify interference among near-neighbors; asymmetric noise induces directional retrieval asymmetry; and insufficient scale control can create modality dominance that looks like "success" while violating symmetry.

A third lesson is that representation health cannot be inferred from a single headline metric. Retrieval@1 is a useful operational proxy, but it can hide pathologies. A model can maintain retrieval while drifting toward hubness, anisotropy, or partial collapse. Conversely, a model can show moderate retrieval while encoding latent factors cleanly and robustly. This is why we introduced a panel of diagnostics: directional retrieval metrics, symmetry gaps, separability scores by factor, embedding variance floors, mean off-diagonal cosine similarity, and covariance spectrum summaries such as effective rank. The combination is the point. Professional practice requires a set of monitors that together describe the state of the representation contract. Students should internalize that alignment is not a point claim; it is a profile across metrics and perturbations.

The chapter also insisted on stress testing as the scientific method for multimodal systems. Rather than asking whether a model works at one operating point, we swept noise asymmetry, pairing quality, temperature, and batch structure to map degradation curves. These curves reveal margins and cliffs. A robust system degrades smoothly and predictably; a brittle system exhibits sudden failure when a threshold is crossed. In operational settings, cliffs are dangerous because they produce silent failure under small environment changes. Therefore, acceptance criteria should be defined distributionally: not "the model achieved $x$ once," but "the model remains within bounds across plausible perturbations." This is the same posture used in disciplined trading research: mechanisms are validated through surfaces and stress scenarios, not through a single backtest.

Finally, this chapter treated governance artifacts as part of the technical solution. In frontier settings, the inability to reproduce a result is itself a form of failure. The companion notebook's audit bundle—run manifest, prompts log, risk log, deliverables folder, and strict JSON reports separating facts from assumptions and open items—is not overhead. It is the minimum infrastructure for reviewability. It forces epistemic discipline: verification status remains `"Not verified"` because

a synthetic lab demonstrates mechanisms, not production guarantees. Students should learn that the correct response to frontier capability is not hype but instrumented uncertainty: measure what you can, record what you did, state what remains unknown, and define questions to verify before deployment.

These lessons set up the transition to Chapter 2. If Chapter 1 established what alignment is and how embedding geometry can be made interpretable under controlled conditions, Chapter 2 asks the harder question: how does alignment fail in structurally predictable ways? We will study modality dominance, spurious shortcuts, pairing corruption, and collapse as failure modes that can produce misleading success metrics. The purpose is not to become cynical about multimodal systems, but to become professionally competent: to recognize that alignment is fragile, to detect failure signatures early, and to design controls that prevent impressive demos from becoming unreviewable infrastructure. Under the AI 2026 umbrella, the arc is deliberate: the frontier is not defined by what systems can do, but by the gap between capability and controllability. This book exists to make that gap legible.

# Bibliography

[1] A. van den Oord, Y. Li, and O. Vinyals. Representation Learning with Contrastive Predictive Coding. *arXiv preprint* arXiv:1807.03748, 2018.

[2] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A Simple Framework for Contrastive Learning of Visual Representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020.

[3] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum Contrast for Unsupervised Visual Representation Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[4] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning Transferable Visual Models From Natural Language Supervision. *arXiv preprint* arXiv:2103.00020, 2021.

[5] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y.-H. Sung, Z. Li, and T. Duerig. Scaling Up Visual and Vision-Language Representation Learning With Noisy Text Supervision. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021.

[6] X. Zhai, X. Wang, B. Mustafa, A. Steiner, D. Keysers, A. Kolesnikov, and L. Beyer. LiT: Zero-Shot Transfer with Locked-Image Text Tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[7] J. Li, R. Selvaraju, A. Gotmare, S. Joty, C. Xiong, and S. C. H. Hoi. Align before Fuse: Vision and Language Representation Learning with Momentum Distillation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

[8] A. Bardes, J. Ponce, and Y. LeCun. VICReg: Variance-Invariance-Covariance Regularization for Self-Supervised Learning. *arXiv preprint* arXiv:2105.04906, 2021.

[9] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow Twins: Self-Supervised Learning via Redundancy Reduction. *arXiv preprint* arXiv:2103.03230, 2021.

[10] S. Kornblith, M. Norouzi, H. Lee, and G. Hinton. Similarity of Neural Network Representations Revisited. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019.

# Chapter 2

# Multimodal Failure Modes

**Abstract.** Multimodal alignment is often presented as a success story: train with contrastive objectives, obtain high retrieval, and declare a shared embedding space. This chapter argues that such a conclusion is methodologically unsafe. Alignment is a *fragile contract* between measurement channels, and it can fail in structurally predictable ways while headline metrics improve. We develop a mechanism-first taxonomy of four failure modes that dominate practice: modality dominance (one encoder controls similarity through scale or gradient energy), spurious alignment (shared confounders create shortcuts that outperform intended semantics), pairing corruption (mislabeled or mis-joined pairs induce contradictory supervision), and representation collapse (degenerate embeddings with reduced rank and hubness).

Using a synthetic-first laboratory, we re-implement a minimal working aligner (two-layer MLP encoders with symmetric InfoNCE and stable numerics) and then induce failures through controlled interventions: gain scaling and asymmetric noise for dominance; shared ID codebooks and parity markers for confounding; systematic pair swaps for corruption; and low-temperature, high-learning-rate regimes for collapse. The chapter introduces three structural probes for attribution beyond retrieval: (i) a mutual-information proxy between embeddings and true latent factors via discretized empirical entropies, (ii) a CCA-like cross-modal correlation proxy computed by SVD to quantify shared linear structure, and (iii) an influence-style sensitivity proxy measuring how a single update perturbation changes retrieval@1.

The companion notebook operationalizes this framework with governance artifacts: deterministic seeds, split hygiene checks, early-warning indicators (variance trends, mean cosine, covariance spectra, gradient norms), and audit-ready exports (manifest, risk log, strict JSON reports). The chapter's conclusion is operational: multimodal systems must be evaluated as governed mechanisms, with failure signatures and stage gates defined over stress distributions rather than point estimates.

## 2.1 Introduction

### 2.1.1 Why "it works" is not evidence: alignment as a fragile contract

Multimodal systems are routinely demonstrated through a small number of compelling interactions: retrieve an image given a caption, retrieve a caption given an image, or cluster mixed media into coherent groups. When these demos succeed, it is tempting to treat the system as "aligned" in a stable, semantic sense. This chapter begins by rejecting that temptation. In multimodal learning, "it works" is not evidence; it is a hypothesis. The difference matters because alignment is not a property of a single run on a single dataset snapshot. Alignment is a *contract* between measurement channels and a learned geometry, and contracts can fail silently.

To call alignment a contract is to make a precise claim: paired observations—produced by different measurement operators but corresponding to the same latent state—should map to near points in a shared embedding space, and unpaired observations should map to separated points. The contract includes symmetry obligations (image→text should be comparable to text→image), stability obligations (small measurement perturbations should not invert relevance), and semantic obligations (the proximity structure should reflect intended factors rather than incidental artifacts). The

contract can be satisfied for the wrong reasons, and those wrong reasons are not rare corner cases. They are structurally favored by contrastive objectives.

Contrastive objectives are incentive systems. They reward whatever features make positives outrank negatives. If a shared shortcut feature exists, the model is rewarded for exploiting it. If one modality is systematically cleaner or higher signal-to-noise, the model is rewarded for relying on that modality. If pairing is partially wrong, the model is rewarded for satisfying contradictory constraints in whatever degenerate way reduces loss. If temperature is too low or optimization too aggressive, the model is rewarded for creating sharp similarity peaks that can induce anisotropy and collapse. In each case, the model can produce high retrieval under the training distribution while degrading the semantic meaning of the space. The student must learn to treat retrieval not as the conclusion but as the beginning of an investigation.

This is why Chapter 2 is critical. Chapter 1 taught multimodality by construction: how to build a controlled world, train a shared embedding space, and interpret its geometry through diagnostics. Chapter 2 asks a harder question: if we accept the same training mechanism, under what structural conditions does it fail, and how can we detect those failures before they become production incidents? The premise is that failure modes are not accidents. They are stable patterns generated by the interaction of data, objectives, and optimization dynamics. If we can name them and induce them in a laboratory, we can also monitor for them and design controls in production.

The language of fragile contracts also reframes what students should expect from frontier systems. Frontier systems do not primarily fail by refusing to learn. They fail by learning *something else—* a shortcut, a dominance pathway, a degenerate geometry—while producing outputs that look superficially correct. This is why multimodal learning can be more dangerous than unimodal learning: the coupling between modalities creates additional degrees of freedom for shortcuts and silent failures. A model can appear consistent in one direction but not the other. It can appear robust to image noise but not to text noise. It can appear to "understand" semantics while actually matching through IDs, formatting, or dataset artifacts. These are not theoretical worries. They are normal engineering realities.

Therefore, the core methodological shift in this chapter is: we treat multimodal alignment as a *governed mechanism* rather than a capability. A governed mechanism is one whose behavior can be reproduced, instrumented, stressed, and audited. The notebook that accompanies this chapter is designed to produce evidence rather than persuasion. Each experiment induces a failure mode, measures its signatures, and exports artifacts that separate facts from analysis and open items. The goal is not to make students fearful. The goal is to make students competent. Competence means knowing what the system is allowed to claim, what evidence supports that claim, and what controls are required for safe deployment.

## 2.1.2   Failure modes as structural, not anecdotal

It is common to discuss failure through anecdotes: "the model matched the wrong caption" or "the model collapsed to generic embeddings" or "the model used the watermark." Anecdotes are memorable but not actionable. They do not tell you what to monitor, what knob caused the failure, or how likely the failure is under plausible shifts. This chapter insists on a structural framing. A structural failure mode is defined by three components: a causal mechanism, observable signatures, and controllable interventions.

For example, consider modality dominance. The causal mechanism is scale or gradient imbalance: one modality's encoder produces embeddings or gradients with larger magnitude, causing the joint objective to optimize primarily for that modality. The signatures include directional retrieval asymmetry, norm drift, hubness, and covariance spectrum imbalance between modalities. The interventions include gain scaling, asymmetric noise injection, or capacity imbalance. This is structural because it is not a one-off oddity; it is a predictable result of mismatched signal and optimization energy.

Similarly, spurious alignment is structural. The mechanism is the presence of a shared confounder that correlates with pair identity but is unrelated to intended semantics. The signatures include high retrieval combined with poor factor separability, a sharp performance drop when the confounder is removed, and embedding neighborhoods organized by confounder groups rather than semantic factors. The interventions include injecting an ID codebook into both modalities, adding parity markers, or adding a watermark pattern. Again, this is structural: when the objective rewards a shortcut, the shortcut is selected.

Pairing corruption is structural. The mechanism is contradictory supervision: the training objective treats mismatched pairs as positives, forcing the model to satisfy inconsistent constraints. The signatures include instability, degraded symmetry, increased hubness, and sometimes non-monotonic loss behavior as the model searches for degenerate compromises. The interventions include swapping a fraction of pairs or structured swaps that mimic join bugs. The structural lesson is that the multimodal contract is only as good as the pairing pipeline.

Collapse is structural. The mechanism is a degenerate geometry that reduces loss under sharp competition but destroys representational diversity. The signatures include rising mean cosine similarity, shrinking variance, contracting covariance spectra, reduced effective rank, and hubness. The interventions include low temperature, high learning rate, and excessive regularization without variance preservation. Collapse is not merely a bug; it is a solution the optimization can prefer under certain regimes.

By presenting failure modes as structural, we turn them into a curriculum. Students can learn to recognize signatures, understand mechanisms, and reason about controls. This is also a governance curriculum: each failure mode can be mapped to a monitoring plan and acceptance criteria. Anecdotes do not produce governance. Structure does.

This structural framing also prevents a subtle pedagogical error: confusing interpretability with comfort. Students might think that if PCA plots look "nice," the system is safe. In reality, PCA plots can look nice under confounding. Or they can look messy under legitimate continuous factor structure. Therefore, we treat failure modes as objects of experimentation rather than as aesthetic judgments.

Finally, structural framing forces honesty about what the lab can and cannot prove. Our synthetic world can induce and detect failure modes, but it cannot prove that a given production system is safe. It can only teach the student what to look for and how to design controls. This is why reports remain "Not verified." The chapter is a training ground for professional skepticism, not a certification tool.

### 2.1.3 From capability to controllability: what Chapter 2 adds

Chapter 1 established a positive construction: how a shared embedding space can be learned and interpreted when the world is designed to be interpretable. Chapter 2 adds a negative construction: how that space can be made to fail in controlled ways, and how those failures can be diagnosed and governed. This is not a detour. In frontier topics, the boundary between success and failure is where most professional value is created.

The concept of controllability is borrowed from systems thinking. A system is controllable if we can influence its behavior predictably through a set of inputs or controls. In multimodal learning, controllability means we can predict how geometry changes when we change measurement noise, when we change temperature, when pairing quality degrades, or when confounders are introduced. A controllable system is not one that never fails. It is one whose failure modes are legible and bounded. That is the standard we aim to teach.

Chapter 2 therefore introduces three new methodological elements beyond Chapter 1. First, it introduces *induced failures*: we do not wait for the model to fail accidentally; we create the conditions that produce failure. Second, it introduces *structural probes* that go beyond retrieval: mutual-information proxies, CCA-like correlation proxies, and influence-style sensitivity proxies. These probes are designed to help attribute causes, not merely measure outcomes. Third, it introduces *early warning dashboards*: variance trends, mean cosine trends, covariance spectra, and gradient norms tracked over time and across experiments.

These elements are intended to teach a concrete workflow. In professional settings, when a multimodal system changes behavior, engineers must answer: is this due to data drift, measurement drift, pairing drift, confounding, optimization instability, or representation collapse? Chapter 2 offers a laboratory in which these causal categories can be separated and observed. The student learns to distinguish: "retrieval dropped because text noise increased" versus "retrieval remained high but factor separability collapsed because a confounder dominated" versus "retrieval became asymmetric because dominance emerged."

Chapter 2 also forces the student to confront a discomforting truth: high performance can be misleading. A confounder can improve retrieval. A dominance regime can improve retrieval in one direction while ruining the other. A near-collapse anisotropic space can improve top-1 while ruining diversity and increasing hubness. Without a governed diagnostic panel, a team can celebrate improvements while silently degrading the representation contract. This is why we stress that capability without controllability is not progress. Frontier work must be evaluated through controllability, not through spectacle.

### 2.1.4 Operational stakes: embeddings as infrastructure and hidden policy

Why do failure modes matter beyond academic interest? Because embeddings become infrastructure. Once a system produces embeddings, many downstream components treat those embeddings as if they were stable semantics. Search engines rank by nearest neighbor. Recommendation engines cluster by similarity. Deduplication systems drop "redundant" items based on distance thresholds. Monitoring systems detect anomalies by comparing embeddings over time. In each case, the embedding space is a latent policy layer: it determines what the system treats as the same, as similar, as relevant, or as outlier.

If the embedding geometry drifts, the policy drifts. But the drift is hidden, because no explicit rules changed. A user may suddenly see different search results. A compliance pipeline may suddenly group different items. A fraud system may suddenly miss anomalies because embeddings became overly smooth. These are not hypothetical. They are consequences of representation as infrastructure.

Failure modes in Chapter 2 correspond directly to infrastructure incidents. Modality dominance can cause a system to ignore one channel, leading to systematic blind spots. For example, a product search system might over-rely on text metadata and ignore visual signals, or vice versa, causing retrieval errors in cases where the ignored modality is essential. Spurious alignment can cause a system to cluster by irrelevant metadata, creating biased or nonsensical groupings. Pairing corruption can cause a system to learn inconsistent associations, leading to unpredictable retrieval and false matches. Collapse can cause a system to return the same few hubs for many queries, producing repetitive results and reducing diversity.

In regulated or high-stakes settings, these failures can have compliance implications. If embeddings are used to route content for review, a drift-induced change in clustering can alter what gets reviewed. If embeddings are used in audit workflows, a spurious shortcut can create false similarity groups that influence decisions. Therefore, the professional requirement is not only to measure performance but to maintain a reviewable chain of evidence: what changed, why, and what controls were in place.

This is why the notebook produces audit artifacts. The artifacts are not mere teaching props. They model a professional standard: every run generates a manifest, logs, and reports, so that conclusions can be traced back to conditions. Students should learn that frontier systems require this discipline

precisely because failures are silent and policy-like.

### 2.1.5 Notebook contract: controlled failures, probes, and audit artifacts

The companion notebook for Chapter 2 is designed as a controlled failure laboratory. It begins with a working baseline aligner, not because baseline performance is the goal, but because controlled failures require a reference state. Without a baseline, any observed behavior can be blamed on an implementation error or a weak model. With a baseline, we can ask counterfactual questions: what changed when we introduced a confounder? What changed when we scaled one modality? What changed when we corrupted pairings? What changed when we reduced temperature?

The contract of the notebook is therefore:

- **Reproduce the baseline deterministically.** Same seeds, same splits, same results within tolerance.
- **Induce one failure mode at a time.** Change a single structural knob or a controlled small set of knobs, and record them in the manifest.
- **Measure both outcomes and signatures.** Retrieval metrics are outcomes; variance trends, spectra, symmetry gaps, and probe scores are signatures.
- **Use structural probes for attribution.** MI proxy, CCA proxy, and influence proxy help distinguish failure categories.
- **Export audit-ready artifacts.** Every run produces strict JSON reports separating facts, assumptions, open items, analysis, and draft output, with verification status "Not verified."

The notebook also enforces a governance posture. It includes shape assertions, NaN checks, log-sum-exp stability, and gradient checks where applicable. It includes leakage checks for splits. It logs early warnings over time. If a run fails a correctness or health constraint, it should fail loudly and record why. This models professional engineering: silent failure is unacceptable.

Finally, the notebook is designed to be pedagogical without being trivial. It uses two-layer MLP encoders to remain interpretable while still requiring real optimization and diagnostic effort. It uses synthetic worlds with explicit latent factors to allow factor-level interpretation. It uses controlled confounders and corruption to induce realistic failure patterns. The goal is that students can move from reading to doing: they can run the notebook, see failures occur, and learn to recognize signatures. That experience is the core of Chapter 2.

This introduction sets the stage for the rest of the chapter. We have defined why alignment is a fragile contract, why failure modes are structural, what Chapter 2 adds in terms of controllability, why the stakes are operational, and what the notebook contract is. The next section formalizes the key definitions needed for failure analysis so that we can speak precisely about dominance, spurious alignment, pairing corruption, and collapse, and so that we can distinguish signatures from mechanisms in a governed way.

## 2.2 Key Definitions for Failure Analysis

### 2.2.1 Dominance: scale, gradient energy, and directional asymmetry

We begin with modality dominance because it is the most common failure mode that hides behind apparent success. Dominance occurs when one modality effectively determines similarity structure in the shared space. The model may still be called "multimodal" because it accepts two inputs and produces two encodings, but the coupling contract is violated: one modality becomes the primary carrier of meaning, while the other becomes a weak auxiliary. In practice, dominance can be subtle. A system can show acceptable average retrieval while being unusable in one direction or under certain noise conditions. Therefore, we define dominance in a way that can be diagnosed, not merely described.

Let $\hat{z}^{\text{img}}$ and $\hat{z}^{\text{txt}}$ be normalized embeddings. The similarity matrix $S_{ij} = \langle \hat{z}_i^{\text{img}}, \hat{z}_j^{\text{txt}} \rangle / \tau$ induces two retrieval functions: image→text uses row-wise maxima, text→image uses column-wise maxima. In an ideal symmetric system, performance in both directions is comparable, and the embedding distributions are comparable across modalities.

Dominance is a *geometric imbalance* characterized by at least one of the following:

- **Directional performance asymmetry:** retrieval@$k$ differs materially between image→text and text→image in a stable way.
- **Distributional imbalance:** pre-normalization embedding norms, covariance spectra, or effective ranks differ significantly across modalities.
- **Gradient-energy imbalance:** gradient norms with respect to encoder parameters (or intermediate activations) are consistently larger for one modality, indicating that optimization pressure is concentrated there.
- **Hubness asymmetry:** one modality's embeddings become hubs that attract many queries, producing stripes in similarity heatmaps and repetitive retrieval results.

The causal mechanism is usually scale or signal imbalance. If one modality has higher signal-to-noise, the objective can be reduced more efficiently by shaping that modality's encoder. If one modality's features have larger scale, gradients propagate with larger magnitude. If one modality's encoder has higher capacity, it can overfit the pairing contract while the other remains underfit. Dominance can also be induced by temperature and optimization: under sharp competition, the model may compress one modality's representation to a set of anchor points, making it easier to separate positives by mapping the other modality toward those anchors.

Pedagogically, dominance teaches a key frontier lesson: a shared space is not automatically shared. It must be engineered to be symmetric, and symmetry must be monitored. In production, dominance creates hidden policy bias: the system's notion of similarity becomes primarily visual or primarily textual depending on which modality dominates, and that can violate user expectations and

downstream requirements.

### 2.2.2 Spurious alignment: confounders, shortcuts, and leakage channels

Spurious alignment is the failure mode in which the model achieves high apparent alignment by exploiting features that are shared across modalities but unrelated to the intended semantics. These features are confounders. They can be explicit IDs, formatting artifacts, dataset ordering signals, watermarks, or any metadata correlated with pairing. Spurious alignment is especially dangerous because it can improve retrieval metrics, making teams believe the model is becoming "better," while it is becoming less semantically meaningful.

We define a confounder as a variable $c$ that influences both modalities' observations but is not part of the intended latent semantics $f$. In a synthetic world, $f$ might be shape and orientation, while $c$ might be sample index parity. Observations become

$$x^{\text{img}} = \mathcal{M}_{\text{img}}(f, c), \qquad x^{\text{txt}} = \mathcal{M}_{\text{txt}}(f, c).$$

A model that optimizes contrastive loss is rewarded for using any shared information, whether it comes from $f$ or $c$. If $c$ is easier to extract than $f$, the model will prefer $c$. This preference is a structural property of optimization: it selects the easiest signal consistent with the objective.

Spurious alignment has identifiable signatures:

- **Counterfactual fragility:** performance collapses when the confounder is removed or randomized at evaluation.
- **Semantic under-encoding:** factor separability for intended factors remains low even when retrieval is high.
- **Group-structured similarity:** embeddings cluster by confounder groups rather than by semantic factors.
- **Probe inconsistency:** MI proxy with intended factors remains low while retrieval remains high; CCA proxy may remain high due to shared confounder structure.

Leakage channels are a special case of confounding. Leakage occurs when information from the target pairing or split assignment is unintentionally present in the input. For example, if a preprocessing pipeline encodes file paths that contain labels, or if an ID appears in both modalities, the model can learn to match IDs rather than semantics. This is why split hygiene alone is not sufficient; we must also audit feature channels for shared shortcuts.

Pedagogically, spurious alignment forces students to internalize that contrastive learning is not a semantic oracle. It is a matching incentive. Without explicit controls, it will learn whatever makes matching easiest. This is why counterfactual tests are mandatory in a governed evaluation: remove suspected shortcuts and measure the effect.

### 2.2.3  Pairing corruption: label noise, join errors, and contradiction

Pairing is the coupling assumption that defines the multimodal contract. In supervised classification, label noise corrupts the mapping from inputs to labels. In contrastive multimodal learning, pairing corruption corrupts the mapping from one modality to the other. The effect is not merely "noisy supervision"; it is contradiction. The objective treats corrupted pairs as positives and true pairs as negatives (when the true partner appears elsewhere in the batch). This creates a structurally different training problem.

We define pairing corruption through a corruption operator $\pi$ applied to indices. For a fraction $\rho$ of samples, we replace the true pairing $(x_i^{\text{img}}, x_i^{\text{txt}})$ with $(x_i^{\text{img}}, x_{\pi(i)}^{\text{txt}})$ where $\pi(i) \neq i$. Corruption can be random swaps, structured swaps (e.g., within a confounder group), or systematic shifts (e.g., offset by one index). Each corruption pattern induces different failure behavior.

The signatures of pairing corruption include:

- **Training instability:** loss oscillations, non-monotonic behavior, and sensitivity to learning rate.
- **Geometry distortion:** increased hubness, reduced factor separability, and spectrum anomalies as the model tries to satisfy inconsistent constraints.
- **Asymmetry under corruption:** one direction may degrade faster, especially if one modality provides clearer anchors.
- **Probe instability:** influence proxy indicates high sensitivity of retrieval to single-step updates, consistent with contradictory gradients.

The causal mechanism is that the objective is optimizing against an inconsistent matching relation. Under small corruption, the model may tolerate the noise and still learn meaningful structure. Under larger corruption, the model may converge to degenerate compromises: for example, mapping many items to similar embeddings to reduce loss variance, or relying on confounders that remain consistent. This is why pairing corruption is a gateway to other failure modes: it can induce dominance and collapse indirectly.

Operationally, pairing corruption corresponds to data join bugs and pipeline drift, not merely human label mistakes. In production, multimodal datasets are often built by joining logs, metadata tables, and content stores. A small change in join keys or timestamp alignment can introduce systematic mispairing. The model will not warn you; it will learn something else.

### 2.2.4  Collapse: anisotropy, hubness, rank loss, and degeneracy

Collapse is the failure mode where embeddings lose diversity and the representation space degenerates. In contrastive learning, full collapse (all embeddings identical) is often prevented by the presence of negatives, but partial collapse and anisotropic collapse are common: embeddings concentrate in a narrow subspace or form hubs that dominate retrieval. Collapse is therefore not a binary event; it is

a spectrum of degeneracy.

We define collapse through geometric indicators:

- **Mean cosine similarity increase:** off-diagonal cosine similarity becomes significantly positive, indicating over-alignment across unrelated items.
- **Variance collapse:** coordinate-wise variance falls below a floor; embeddings lose dynamic range.
- **Rank collapse:** covariance spectrum contracts; effective rank decreases; most variance is captured by few directions.
- **Hubness increase:** certain embeddings become nearest neighbors for many queries; similarity heatmaps show stripes.
- **Sensitivity spikes:** small perturbations in parameters or temperature produce large changes in retrieval, indicating brittle geometry.

The causal mechanisms include sharp competition (low temperature), overly aggressive optimization (high learning rate), regularization imbalance (strong L2 without variance preservation), and batch composition effects (hard negatives inducing anisotropy). Collapse is also influenced by normalization: if embeddings are normalized, the model may collapse by aligning directions rather than magnitudes. Therefore, collapse indicators must be designed to detect directional degeneracy, not only magnitude shrinkage.

Pedagogically, collapse teaches students that a model can "win" the objective in a way that loses meaning. A highly anisotropic space can separate a small set of categories well but fail to represent continuous factors smoothly. Hubness can produce decent top-1 retrieval on familiar distributions but fail catastrophically under drift, where queries land near hubs. Therefore, collapse is a governance risk: it undermines generalization and stability.

### 2.2.5 Failure signature: observable indicators vs causal mechanisms

The final definition in this section is methodological. A failure mode is not merely a bad outcome; it is a causal pattern. To reason professionally, we must distinguish *signatures* (observable indicators) from *mechanisms* (causal explanations). The same signature can be produced by different mechanisms, and the same mechanism can produce different signatures depending on regime.

For example, directional retrieval asymmetry is a signature. It can be caused by dominance, but it can also be caused by asymmetric noise or by pairing corruption that affects one modality more than the other. Mean cosine increase is a signature. It can be caused by collapse, but it can also be caused by confounders that push embeddings into shared directions. Loss oscillation is a signature. It can be caused by pairing corruption, but it can also be caused by high learning rate under clean data.

Therefore, Chapter 2 emphasizes *attribution.* Attribution requires interventions and probes. We

induce failures in controlled ways, and we compute probes that are sensitive to specific mechanisms: MI proxy to test whether intended factors are encoded; CCA proxy to test whether modalities share linear structure; influence proxy to test sensitivity to training perturbations. We also monitor gradient norms and covariance spectra to distinguish dominance from collapse. This is the governing logic: do not infer cause from a single symptom; use a panel of evidence.

This distinction is the conceptual hinge between Chapter 2 and Chapter 3. Chapter 2 teaches failure taxonomy and signatures under controlled interventions. Chapter 3 will teach drift and monitoring: how to detect when production behavior is entering a failure regime. Without a disciplined separation between signatures and mechanisms, monitoring collapses into noise. With that discipline, monitoring becomes actionable: a signature triggers a specific investigation path and a specific set of controls.

In summary, this section defined the four primary failure modes—dominance, spurious alignment, pairing corruption, and collapse—in terms of causal mechanisms and observable signatures. It also defined the critical methodological distinction between signatures and mechanisms, motivating the need for structural probes and controlled interventions. With these definitions in place, we can now specify the baseline system that will serve as the reference state for all induced-failure experiments.

## 2.3 Baseline System: Minimal Working Aligner as Reference

### 2.3.1 Synthetic world recap: latent factors and measurement operators

A failure analysis is only as strong as its reference point. If we begin with an unstable or poorly specified baseline, every subsequent "failure" can be dismissed as an artifact of the baseline itself. Therefore, Chapter 2 begins by re-establishing a minimal working multimodal aligner under governed conditions. This baseline is not optimized for state-of-the-art performance. It is optimized for interpretability, stability, and reproducibility. Its purpose is to provide a controlled starting geometry against which failures can be induced and attributed.

The synthetic world is defined by latent factors $f$ that represent intended semantics. We reuse a factor set like the one in Chapter 1 because it is sufficiently rich to express both discrete and continuous structure: shape type (discrete), orientation (continuous or binned), spatial frequency (discrete or ordinal), phase (continuous or binned), and thickness (ordinal). These factors define the latent state of each sample, and the experiment's claims are always conditioned on these factors. This is crucial. Without explicit factors, we cannot measure whether the embedding space encodes the semantics we intend. Real-world datasets often lack this clarity; synthetic worlds provide it.

Each modality is a measurement operator applied to the latent factors. For the image modality, the measurement operator $\mathcal{M}_{\text{img}}$ renders factors into a small matrix representation. For example, shape and thickness determine a spatial pattern, orientation determines rotation, frequency and phase

determine grating patterns, and controlled noise is added afterward. For the text modality, the measurement operator $\mathcal{M}_{\text{txt}}$ maps factors into a symbolic token sequence or a structured feature vector derived from that sequence. This mapping is not "language understanding"; it is a designed grammar. The grammar's role is to ensure that text contains the same latent information, but through a qualitatively different encoding: discretized bins, token order, and sparse feature maps.

In this chapter, we treat measurement operators as objects that can be perturbed. Noise intensity, token corruption, rendering distortions, and normalization constants are all part of the operator. This matters because many failure modes can be induced by changing the operator in one modality. Therefore, the baseline world definition includes explicit operator parameters and logs them to the run manifest.

Split hygiene is also part of the world definition. We generate a dataset, then deterministically split indices into train, validation, and test. We assert non-overlap. We ensure that the factor distribution is comparable across splits, or we at least record distribution differences. This prevents accidental leakage and ensures that evaluation reflects generalization rather than memorization of generation artifacts.

The baseline world is therefore defined by: a factor generator $f_i$, two measurement operators producing $(x_i^{\text{img}}, x_i^{\text{txt}})$, a noise model, and a deterministic split. This explicit construction is the foundation for all later claims. When a failure is induced, we will identify whether we perturbed the factor generator, the measurement operators, the pairing relation, or the training regime.

### 2.3.2   Encoders: two-layer MLP symmetry and capacity constraints

The baseline encoders are intentionally minimal but nontrivial: two-layer MLPs for each modality. Each encoder maps its input feature vector to an embedding vector in a shared dimension $d$. The architecture is:

$$h = \sigma(xW_1 + b_1), \qquad z = hW_2 + b_2,$$

followed by $\ell_2$ normalization to produce $\hat{z}$. The nonlinearity $\sigma$ is chosen to be stable and differentiable (e.g., tanh) to support analytic gradient computation. We choose hidden width large enough to allow representation but small enough to maintain interpretability and computational feasibility.

The key design principle is *symmetry*. The modalities differ in their input feature dimensions and distributions, but the encoder class is the same. This reduces confounding. If we observe dominance or asymmetry, we can attribute it more confidently to measurement operator differences, noise, or training dynamics rather than to a mismatch in encoder expressivity. Symmetry does not guarantee symmetric outcomes, but it makes asymmetry interpretable.

Capacity constraints are also deliberate. If the encoders are too powerful, they can memorize pairing structure or confounders even when factor semantics are weak. If they are too weak, they cannot represent the factor mapping, and failures become indistinguishable from underfitting. Two-layer

MLPs sit in a productive regime: they can represent nonlinear interactions but cannot trivially memorize high-dimensional IDs without paying a parameter cost that we can monitor. In synthetic labs, we can also vary capacity as a controlled intervention when studying dominance or confounding. The baseline establishes a capacity regime that produces meaningful alignment without excessive fragility.

Initialization is deterministic and recorded. We use seeded random initialization with scale control (e.g., Xavier-like scaling) to prevent saturating nonlinearities and to maintain stable gradients. We also implement explicit checks: parameter shapes, finite values, and initial embedding variance. These checks are baseline obligations. If the baseline begins in an unstable regime, failure analysis is compromised.

Finally, we log encoder-level diagnostics even in the baseline: pre-normalization norm distributions, gradient norms per modality during training, and covariance spectrum summaries. These metrics are not introduced only when failures appear; they are part of the baseline evidence. This teaches students that monitoring should be continuous, not reactive.

### 2.3.3   Objective: symmetric InfoNCE with stable numerics

The baseline coupling objective is symmetric InfoNCE, as in Chapter 1. The similarity matrix is computed using cosine similarity divided by temperature:

$$S_{ij} = \frac{\langle \hat{z}_i^{\text{img}}, \hat{z}_j^{\text{txt}} \rangle}{\tau}.$$

The loss is the average of two directional cross-entropies: image→text and text→image. Symmetry here is important: it ensures that both modalities contribute to the optimization signal. If we trained only one direction, dominance could be induced trivially by design.

Stable numerics are mandatory. We compute row-wise log-softmax with log-sum-exp stabilization. We include epsilon guards for normalization and log computations. We assert that loss is finite at every step. We also compute intermediate statistics of $S$: mean diagonal similarity, mean off-diagonal similarity, and distribution quantiles. These similarity statistics are early indicators of collapse or confounding.

Temperature is treated as a structural knob. The baseline uses a temperature that produces stable learning without excessive sharpness. Importantly, we do not tune temperature to maximize retrieval. We choose it to achieve a balanced geometry: acceptable retrieval with healthy variance and effective rank. This models professional practice: hyperparameters should be chosen under constrained criteria, not a single metric.

We also treat batch size as part of the objective's operationalization. Because InfoNCE uses batch negatives, batch size influences negative density and therefore geometry. The baseline chooses a

batch size that is feasible in Colab and large enough to create meaningful competition. Batch composition is controlled by deterministic shuffling.

The baseline objective implementation is validated through gradient checking on a small batch and a subset of parameters. This is a correctness obligation. A failure analysis built on incorrect gradients is meaningless. Therefore, the baseline includes analytic gradients for the encoders and a finite-difference check that must pass within tolerance. This check is performed under a numerically stable regime (reasonable epsilon, moderate temperature) to reduce false failures due to floating-point noise.

### 2.3.4 Baseline acceptance: retrieval, symmetry, variance floors

We define baseline acceptance criteria not as claims of real-world performance but as minimal health gates for the synthetic system. The baseline should satisfy:

- **Directional retrieval:** retrieval@1 and retrieval@5 in both directions exceed minimal thresholds on validation (e.g., significantly above chance).
- **Symmetry gap:** the difference between directional retrieval metrics is bounded.
- **Variance floors:** embedding variance and effective rank exceed minimum floors; mean off-diagonal cosine similarity remains below a ceiling.
- **Factor encoding:** separability ratios for key factors (e.g., shape) are above minimal levels; continuous factors show nontrivial structure.
- **Stability:** training does not produce NaNs; gradient norms remain within bounded ranges; loss decreases smoothly in expectation.

These criteria are intentionally modest but non-negotiable. They ensure the baseline is a functioning system with interpretable geometry. They also enforce a core pedagogical message: a model is not accepted because it looks good, but because it satisfies defined constraints.

We also record baseline metrics as a JSON report with strict schema: facts, assumptions, open items, analysis, draft output, verification status, and questions to verify. This report includes artifact paths: plots, checkpoints, and logs. The baseline report becomes the reference for all experiments. Each induced failure experiment will produce its own report and will include a comparison to baseline. This structure teaches students to reason relationally: failures are deviations from a known state, not isolated outcomes.

### 2.3.5 Why a baseline matters: counterfactual control for failure claims

The purpose of a baseline is not to celebrate success; it is to enable causal claims. Without a baseline, when performance degrades under an intervention, we cannot know whether the intervention caused the degradation or whether the system was unstable to begin with. With a baseline, we can make

counterfactual statements: "Under identical conditions except for confounder injection, retrieval increased but factor separability decreased." This is the essence of controlled experimentation.

Baselines also prevent a common methodological trap: moving goalposts. In failure analysis, it is tempting to adjust hyperparameters in each experiment to "make it work." That destroys comparability. The baseline establishes a fixed reference configuration. Induced failures then change one knob at a time relative to that configuration. If additional tuning is performed, it must be logged as a separate experiment, not silently integrated. This is governance: comparability is protected by design.

Another reason baselines matter is that some failures can masquerade as improvements. A confounder can boost retrieval. A dominance regime can boost one direction. A low-temperature setting can sharpen retrieval at the expense of rank. Without baseline comparisons across a panel of metrics, these changes could be interpreted as progress. The baseline provides the anchor against which "improvement" is judged under multi-metric constraints.

Finally, baselines matter because they allow students to learn what "healthy" geometry looks like. In Chapter 1, students saw geometric diagnostics under a stable world. In Chapter 2, they will see those diagnostics under induced failures. The contrast is pedagogically powerful. Students do not learn failure modes in the abstract; they learn them as deviations from a known healthy regime.

With the baseline system defined and accepted, we now have a stable reference for the rest of the chapter. The next section introduces Failure Mode A: modality dominance. We will induce dominance through controlled interventions (gain scaling, asymmetric noise, capacity skew), measure its signatures (directional asymmetry, gradient imbalance, hubness, spectrum differences), and discuss controls that can restore symmetry in both synthetic and production settings.

## 2.4 Failure Mode A: Modality Dominance

### 2.4.1 Mechanism: scale mismatch and gradient imbalance

Modality dominance is the failure mode in which the multimodal contract becomes effectively unimodal: one modality controls similarity structure, and the other becomes a dependent projection. This can happen even when the training objective is symmetric. The reason is that symmetry in the loss does not imply symmetry in optimization. Optimization follows gradient energy, and gradient energy follows signal, scale, and conditioning.

A useful way to formalize dominance is through the notion of *gradient imbalance.* Let $\theta$ denote image encoder parameters and $\phi$ denote text encoder parameters. During training, we can compute gradient norms

$$g_{\text{img}} = \|\nabla_\theta \mathcal{L}\|, \qquad g_{\text{txt}} = \|\nabla_\phi \mathcal{L}\|.$$

If $g_{\text{img}} \gg g_{\text{txt}}$ persistently (or vice versa), the loss landscape is being traversed primarily by changes in one encoder. The system then learns a geometry in which one modality adapts to satisfy matching while the other remains comparatively rigid. In extreme cases, the rigid modality becomes an anchor space and the other becomes a mapper into that space; in other cases, the dominant modality becomes the anchor and the weaker modality is pulled into it.

Why does this gradient imbalance occur? The simplest cause is input scale mismatch. If one modality's input features have larger magnitude or variance, then the activations and gradients of its encoder can also be larger, especially in shallow networks. Even if embeddings are normalized, the pre-normalization geometry influences backpropagation. Another cause is signal-to-noise imbalance. If images are clean and text is noisy, the model can reduce loss faster by shaping the image encoder because image embeddings can become sharp anchors that make positives separable despite noisy text. Conversely, if text is clean and images are noisy, the text encoder can dominate.

A third cause is conditioning and saturation. If one modality's encoder operates in a regime where its nonlinearity saturates (e.g., tanh at large inputs), gradients can vanish, forcing the other modality to carry the optimization. This is a subtle but common issue: a single normalization constant can push one modality into saturation and create dominance without any explicit design decision.

Finally, dominance can be induced by capacity imbalance. If one encoder has a larger hidden width or a more informative feature mapping, it can fit the pairing relation more effectively. The objective does not demand that both modalities contribute equally; it demands only that pairs match. Therefore, the model will use whatever pathway satisfies the objective most efficiently.

Pedagogically, the key lesson is that dominance is not a bug in the loss; it is a property of the coupled optimization problem. Symmetric objectives can produce asymmetric representations if the modalities differ in scale, noise, or capacity. This is why dominance must be monitored explicitly rather than assumed away.

### 2.4.2 Inducing dominance: gain scaling, noise skew, capacity skew

To study dominance as a structural phenomenon, we induce it through controlled interventions that mirror real production changes.

**Gain scaling.** The simplest intervention is to multiply one modality's encoder outputs (pre-normalization) by a gain factor $\alpha$. Although embeddings are normalized for cosine similarity, gain affects gradient flow because it changes the pre-normalization magnitude and therefore the normalization Jacobian. If we apply gain asymmetrically, we can force one modality to produce larger raw activations and gradients. In production, this corresponds to changes in preprocessing scaling, quantization range, or feature normalization constants.

**Noise skew.** We increase noise in one modality's measurement operator while keeping the other fixed. This simulates sensor degradation or pipeline corruption. As noise increases, the noisy

modality becomes less informative, and the objective can be reduced more effectively by relying on the clean modality. Dominance emerges as the system learns to treat the clean modality as the main semantic carrier.

**Capacity skew.** We increase the hidden width of one encoder or enrich its feature mapping (e.g., richer text features) while leaving the other unchanged. This creates an expressivity advantage. The objective then rewards the more expressive modality for fitting the pairing relation and can leave the other modality undertrained. In production, this corresponds to upgrading one encoder architecture or improving one modality's preprocessing pipeline without a corresponding improvement in the other.

Crucially, we induce dominance *without* changing the definition of success. We keep evaluation metrics consistent: directional retrieval, symmetry gaps, variance floors, and factor separability. This ensures that dominance is detected as a deviation from baseline across a panel, not as a subjective impression.

We also emphasize that dominance can look like improvement. If one modality becomes an anchor and the other learns to map into it, retrieval in one direction can improve. For example, if text is clean and images are noisy, text→image retrieval can remain good because the query modality is clean, while image→text may degrade because image queries are noisy. If we only look at average retrieval, we may miss the asymmetry. Therefore, the induced dominance experiments always report directional metrics.

### 2.4.3   Symptoms: retrieval asymmetry, hubness, norm drift

Dominance produces a family of observable signatures that are consistent across induction methods.

**Directional retrieval asymmetry.** The most immediate symptom is a persistent gap between image→text and text→image retrieval. The direction that uses the dominant modality as the query often performs better, because queries from the dominant modality are more informative. The reverse direction often degrades, because queries from the weaker modality land near anchors or hubs without sufficient discrimination.

**Hubness.** Dominance often creates hubness because the weaker modality embeddings become less diverse and collapse toward a set of anchor directions defined by the dominant modality. This produces similarity heatmaps with stripes: many queries from the weaker modality match to a small set of targets in the other modality. Operationally, this corresponds to repetitive retrieval results and reduced diversity.

**Norm drift and distribution mismatch.** Even if we normalize embeddings for similarity, pre-normalization norms can drift. The dominant modality may develop larger raw norms or more structured variance, while the weaker modality may develop smaller norms or saturate. We therefore monitor raw norm distributions, covariance spectra, and effective rank per modality. A persistent

mismatch indicates that the shared space is being used asymmetrically.

**Factor separability imbalance.** Dominance can also appear as factor encoding differences. The dominant modality may encode semantic factors clearly, while the weaker modality may encode them weakly or only through projection. This can be measured by separability ratios computed separately on each modality's embeddings (within-modality structure) and by cross-modal alignment structure.

Pedagogically, students should learn to interpret these symptoms together. A retrieval gap alone does not prove dominance; it could be noise asymmetry. Hubness alone does not prove dominance; it could be partial collapse. Norm drift alone does not prove dominance; it could be scaling. Dominance is inferred through the joint pattern: directional asymmetry plus gradient imbalance plus distribution mismatch.

### 2.4.4   Diagnostics: gradient norms per modality, covariance spectra

Because dominance is primarily an optimization imbalance, diagnostics must include optimization-level measurements, not only geometric outcomes.

**Gradient norms per modality.** We compute per-step gradient norms for each encoder, and we track their ratio over time. A stable ratio far from one indicates dominance pressure. We also compute gradient norms per layer (e.g., for $W_1$ and $W_2$) to detect whether one modality's first layer is saturating (vanishing gradients) while the other remains active.

**Embedding covariance spectra.** We compute the covariance (or centered Gram) matrix of embeddings for each modality and evaluate its singular values. Dominance often manifests as a lower effective rank in the weaker modality. Even when embeddings are normalized, the distribution of directions can become concentrated. Effective rank provides a compact summary of this concentration.

**Cross-modal correlation structure.** A CCA-like proxy can reveal whether shared structure is driven by the dominant modality. If cross-modal principal correlations remain high but one modality's spectrum collapses, the shared structure may be carried by a low-dimensional shortcut.

**Symmetry check suite.** We treat symmetry as an explicit diagnostic package: retrieval gaps, norm distribution divergence, effective rank divergence, and hubness metrics (e.g., how many times each target appears as top-1 across queries). Dominance is identified when multiple symmetry metrics deviate simultaneously.

These diagnostics are not merely computed at the end. They are tracked over time. Dominance often emerges gradually: early training may be balanced, but as the system finds an easier pathway, gradients concentrate in one modality. Time-series monitoring allows students to see dominance as a dynamical process.

### 2.4.5 Controls: normalization discipline, balanced loss weighting, clipping

A failure analysis must end with controls. Controls are not guarantees; they are design responses that reduce risk and increase controllability.

**Normalization discipline.** Ensure that input feature scales are comparable across modalities. Use deterministic training-set-derived normalization, and log it. Ensure that embedding pre-normalization norms are monitored. Avoid saturating nonlinearities through appropriate initialization and input scaling. In production, treat preprocessing as part of the model contract and version it.

**Balanced loss weighting.** Even with symmetric losses, one direction can dominate due to gradient energy. Weighting directional components can partially compensate. For example, if image→text gradients are consistently larger, increase the weight of text→image loss or use adaptive weighting based on gradient norms. However, weighting is not a substitute for scale control; it is a mitigation that must be monitored.

**Gradient clipping and optimizer discipline.** Clip gradients per modality to prevent one modality from dominating updates. Monitor gradient norms and stop or adjust learning rates when imbalance exceeds thresholds. In production, use staged training: pretrain encoders separately to comparable health, then align jointly.

**Capacity matching.** Avoid large capacity upgrades to one modality without corresponding changes to the other, or compensate through regularization and monitoring. If one modality must be upgraded (e.g., improved image encoder), treat the change as

## 2.5 Failure Mode B: Spurious Alignment via Confounders

### 2.5.1 Mechanism: shared shortcuts outperform semantic factors

Spurious alignment is the archetypal failure of contrastive multimodal training because it exploits the very incentive structure that makes contrastive learning powerful. The objective does not ask whether the shared embedding encodes "meaning." It asks whether matched pairs are closer than mismatched pairs. When a shortcut feature is shared across modalities and reliably distinguishes matched from mismatched pairs, the shortcut is an efficient solution. In optimization terms, it is a low-loss pathway with low sample complexity. Unless constrained, it will be selected.

To formalize the mechanism, let $f$ denote intended latent semantics and let $c$ denote a confounder. Both influence the observations:

$$x^{\text{img}} = \mathcal{M}_{\text{img}}(f, c), \qquad x^{\text{txt}} = \mathcal{M}_{\text{txt}}(f, c).$$

The model learns encoders $E_{\text{img}}$ and $E_{\text{txt}}$ such that $\hat{z}^{\text{img}} = E_{\text{img}}(x^{\text{img}})$ and $\hat{z}^{\text{txt}} = E_{\text{txt}}(x^{\text{txt}})$. If $c$ can be extracted more easily than $f$, then a representation that encodes $c$ strongly and $f$ weakly

can still satisfy the InfoNCE ranking constraints. In other words, the objective is satisfied when $\hat{z}^{\text{img}}$ and $\hat{z}^{\text{txt}}$ share any high-mutual-information signal that is stable across pairs. The objective cannot distinguish between "legitimate" shared signal (semantics) and "illegitimate" shared signal (confounders).

This mechanism is not an edge case. Real multimodal datasets are full of confounders: watermarks, camera metadata, timestamps, file naming conventions, web-page template artifacts, OCR artifacts, language or locale patterns, and user-interface chrome. When such artifacts correlate with pairing (or with class labels), the model is rewarded for using them. In production, this can lead to severe brittleness: the moment the confounder distribution shifts (new watermark, new template, new logging format), retrieval and clustering behavior can change abruptly.

The key point is that spurious alignment is not "overfitting" in the classical sense. The model is not memorizing individual training samples. It is learning a generalizable shortcut—generalizable within the dataset artifact regime. That is why it can look strong on validation and test splits if the confounder is present there too. This is why split hygiene alone does not protect against spurious alignment. If the confounder is global, it appears in all splits. The correct defense is counterfactual testing and invariance constraints.

### 2.5.2 Confounder constructions: ID tokens, watermarks, parity markers

To teach spurious alignment as a structural phenomenon, we construct confounders that are easy to control and easy to diagnose.

**ID codebooks.** We assign each sample an index $i$ and construct a random code vector $u_i$ in a shared "ID subspace." We then append $u_i$ to both modalities' raw feature vectors (or inject it into images as a subtle pattern and into text features as an extra token embedding). The codebook is shared across modalities, so the model can match pairs by aligning code vectors rather than by learning factor semantics. This mirrors real settings where filenames, URLs, or database keys leak into both modalities.

**Watermark patterns.** We insert a low-amplitude spatial pattern into images and insert a corresponding token pattern into text. The pattern is unrelated to $f$ but is stable across paired items. This mirrors real watermarks, scanner artifacts, or UI chrome that appear consistently with the content.

**Parity markers and group tags.** We add a binary marker based on sample index parity (even/odd) to both modalities. Although parity is extremely low-information, it can still improve retrieval if the negative set is dense and the marker partitions the dataset. More generally, we can create group tags (e.g., 8 groups) and embed group ID in both modalities. This mirrors domain or source markers: content from the same site, device, or region.

These confounders can be scaled in strength. A strong confounder makes spurious alignment obvious,

but a subtle confounder is pedagogically more valuable because it shows how small shared signals can dominate when semantics are harder. In the notebook, we vary confounder strength and observe how retrieval and factor separability respond. This produces a confounder sensitivity curve analogous to noise sweeps.

We also distinguish two confounder regimes:

- **Confounder-only regime:** the confounder is the only shared signal (or semantics are heavily corrupted). The model must rely on the confounder to succeed.
- **Confounder-plus-semantics regime:** both confounder and semantics are present. The model chooses the easiest mixture. Often, it uses both, but the confounder can dominate early training and lock the model into a shortcut.

This distinction matters because in real systems, confounders rarely fully replace semantics; they compete with semantics. The question is not whether the model can learn semantics at all; it is whether it will prioritize semantics when shortcuts exist.

### 2.5.3   Counterfactual tests: confounder removal and performance collapse

Spurious alignment cannot be diagnosed reliably by observing performance under the confounded distribution alone. If the confounder persists in validation/test, retrieval may look excellent. The diagnosis requires counterfactual evaluation: remove or randomize the confounder and measure what changes.

The counterfactual test suite includes:

**Confounder ablation.** Remove the confounder channel from one or both modalities at evaluation time while keeping the trained model fixed. For example, zero out the ID codebook appendage or remove the watermark pattern. If retrieval drops sharply while factor separability remains low, the model was relying on the confounder.

**Confounder randomization.** Keep confounder magnitude but permute confounder assignments across samples so that the confounder no longer matches pairs. This tests whether the model's geometry depends on confounder identity rather than on semantics. If retrieval collapses, the confounder was decisive.

**Confounder-only diagnostic.** Evaluate MI proxy between embeddings and confounder labels. If confounder MI is high and factor MI is low, the model has encoded the wrong variable strongly.

**Cross-group evaluation.** If confounders define groups (e.g., source tags), evaluate retrieval within-group and cross-group. A confounder-reliant model may perform well within group but poorly across group. This mirrors production: a model trained on one source distribution may fail when content moves to a new source.

Counterfactual evaluation also reveals a subtle governance point: a model can pass standard test splits and still be unsafe. Therefore, the acceptance criteria must include invariance tests. "Not verified" is not only a philosophical label; it reflects the necessity of testing under shifts that remove shortcuts.

Pedagogically, counterfactual tests teach students what it means to claim semantic alignment. Semantic alignment is not "high retrieval"; it is "high retrieval that persists when non-semantic shared signals are removed." This is the structural definition of robustness against spurious alignment.

### 2.5.4   Signatures: high retrieval with low factor separability

Spurious alignment has a characteristic signature profile that distinguishes it from dominance or collapse.

**Retrieval can improve.** Because the confounder can provide a strong matching key, top-1 retrieval can rise. This is why spurious alignment is dangerous: it rewards you with better metrics.

**Factor separability stagnates or degrades.** If the model relies on confounders, it may invest less representational capacity in encoding intended factors. Therefore, separability ratios for shape, orientation, frequency, etc., remain low or even decline relative to baseline. In PCA projections colored by factors, clusters may blur even as retrieval improves.

**Embedding neighborhoods align with confounder groups.** Nearest neighbors share confounder labels rather than factor labels. Similarity heatmaps ordered by confounder group show block structure.

**Symmetry may remain good.** Unlike dominance, spurious alignment can preserve directional symmetry because the confounder is shared in both modalities equally. Therefore, the presence of symmetry does not imply semantic correctness. This is an important student lesson: symmetry is necessary but not sufficient.

**Collapse indicators may remain acceptable.** Spurious alignment does not necessarily cause collapse. The embedding space can remain diverse while being organized by the wrong variable. Therefore, collapse monitors alone do not detect spurious alignment. This motivates the need for factor-based probes and counterfactual tests.

This signature profile shows why Chapter 2 introduces structural probes. Without probes, a team could declare success. With probes, the model is revealed to be misaligned semantically.

### 2.5.5   Controls: invariance tests, ablations, negative sampling design

Controls against spurious alignment must be framed realistically. There is no single trick that eliminates all confounders. The professional response is a layered control strategy.

**Invariance tests as acceptance gates.** Include confounder ablation and randomization tests in evaluation. If performance collapses under plausible shortcut removal, the model fails the stage gate. In production, invariance tests can be approximated by auditing for known artifact features and testing with and without them.

**Data auditing and feature hygiene.** Before training, audit input channels for ID leakage, metadata leakage, and shared preprocessing artifacts. Version preprocessing pipelines and ensure that file paths, URLs, or IDs are not present in features unless explicitly intended. In many failures, the confounder is not subtle; it is a logging convenience that leaked into training.

**Negative sampling and batch design.** The structure of negatives influences shortcut selection. If negatives are too easy, a weak confounder can suffice. If negatives are hard and diverse, the model may need richer semantics. However, hard negatives can also increase collapse risk. Therefore, negative sampling is a control that must be tuned under multi-metric constraints.

**Regularization toward semantic factors.** In synthetic labs, we can add auxiliary losses that encourage factor encoding. In real settings, we can use weak supervision, multi-task objectives, or curated evaluation tasks that require semantic generalization. The key is to create incentives that cannot be satisfied by confounders alone.

**Monitoring in production.** Spurious alignment often reveals itself when a confounder distribution shifts: new watermark, new template, new source. Therefore, monitor embedding behavior by source, device, or preprocessing version. If retrieval or clustering changes sharply when a source changes, spurious shortcuts may be involved.

The core lesson of spurious alignment is that contrastive success is not semantic success. A shared shortcut can dominate and produce better metrics while reducing the intended meaning of the space. Diagnosing this requires counterfactual evaluation and factor-sensitive probes. With spurious alignment established, we proceed to Failure Mode C: pairing corruption, where the coupling relation itself becomes inconsistent and the objective is forced into contradiction.

## 2.6 Failure Mode C: Pairing Corruption and Coupling Breakdown

### 2.6.1 Mechanism: contradictory supervision under contrastive ranking

Pairing is the foundational assumption of multimodal learning: that the paired observations correspond to the same latent state. In supervised learning, label noise corrupts a mapping from inputs to labels; the model can often tolerate modest label noise because the supervision is one-to-many and the model can average. In contrastive multimodal learning, pairing corruption is more severe because the supervision is relational and competitive. A corrupted pair is not just "somewhat wrong"; it directly asserts an incorrect equivalence between two observations, and it simultaneously turns the true partner into an implicit negative when it appears elsewhere in the batch.

To see the contradiction, consider a batch of paired samples indexed by $i \in \{1, \ldots, B\}$. The InfoNCE objective for image→text treats $j = i$ as the positive index for row $i$. If a fraction of pairs are corrupted, then for some $i$, the "positive" $j = i$ is not the true semantic partner. The model is therefore rewarded for bringing together embeddings that correspond to different latent factors, and rewarded for pushing away the true partner if it appears as a negative. This is structurally different from noise in independent labels because the loss is defined over a similarity matrix: every positive choice affects how negatives are weighted.

Under small corruption, the model can sometimes still learn because the clean pairs dominate the objective. Under moderate corruption, the objective becomes inconsistent: there is no embedding geometry that satisfies all pairings simultaneously. Optimization then searches for compromises. These compromises can take degenerate forms: reduce variance (so that mismatches are less costly), create hubs (so that many pairs map to a few anchors), or rely on confounders that remain consistent despite corruption. In this way, pairing corruption often becomes a gateway to other failure modes. It can induce hubness and partial collapse even if the original system was healthy.

The most important conceptual point is that pairing corruption is not merely a problem of "lower data quality." It is a problem of *broken coupling.* Multimodal learning is fundamentally about learning a coupling relation between modalities. When the coupling relation is inconsistent, the objective is not teaching semantics; it is teaching contradiction. Therefore, pairing integrity must be treated as a first-class governance concern.

### 2.6.2 Corruption models: random swaps, structured swaps, partial shuffles

Not all pairing corruption is the same. The pattern of corruption matters, and different patterns correspond to different production risks. Therefore, the chapter distinguishes several corruption models and treats them as separate experimental regimes.

**Random swaps.** For a fraction $\rho$ of samples, we randomly permute text partners. This simulates general join noise or mislabeled captions. Random swaps spread corruption broadly and tend to degrade performance smoothly as $\rho$ increases.

**Structured swaps within factor groups.** We swap partners within a subset defined by a factor group (e.g., within the same shape type). This simulates subtle annotation errors where items are swapped among similar categories. The objective becomes harder to diagnose because retrieval can remain reasonable: swapped items may still share some semantics. However, fine-grained factor encoding may degrade and geometry may become less crisp.

**Systematic offset corruption.** We replace each pairing with the next item in sequence (e.g., $i$ pairs with $i + 1 \mod N$). This simulates pipeline bugs such as off-by-one errors, timestamp misalignment, or sorting inconsistencies. Systematic corruption can create persistent false structure that the model may learn, producing high performance under the corrupted regime but catastrophic

failure under correct pairings.

**Partial shuffles with blocks.** We shuffle partners within contiguous blocks (e.g., per source or per time window). This simulates production joins where misalignment occurs within a time window due to delayed ingestion or caching. The model may learn block-level shortcuts and become sensitive to block boundaries.

These corruption models are pedagogically valuable because they show that "pairing noise" is not one scalar. A system can be robust to random swaps but fragile to systematic offset corruption. It can tolerate within-group swaps but fail under cross-group swaps. Therefore, governance in production must include auditing for specific corruption patterns, not only estimating a global noise rate.

In the notebook, we parameterize corruption by $\rho$ and by corruption type. We sweep $\rho$ and record degradation curves and signature changes. We also run counterfactual evaluations: train under corruption and evaluate under clean pairing, and vice versa. These cross-evaluations reveal whether the model has learned a corrupted coupling relation that generalizes only within the corrupted regime.

### 2.6.3   Instability signatures: loss oscillation, spectrum distortion, hubness

Pairing corruption has a distinctive set of signatures that reflect the contradictory nature of the objective.

**Loss behavior can become non-monotonic.** Under clean data, loss typically decreases smoothly in expectation. Under corruption, the loss can plateau, oscillate, or even decrease while retrieval degrades. This is because the loss is being minimized under an inconsistent target; the model can reduce loss by finding degenerate compromises rather than by learning semantics. Therefore, loss cannot be interpreted as a direct proxy for semantic alignment under corruption.

**Spectrum distortion emerges.** Covariance spectra of embeddings often distort as the model collapses representational degrees of freedom to reduce contradiction cost. Effective rank may decrease. Alternatively, the spectrum may become unusually flat if the model spreads variance to avoid concentrating on any inconsistent direction. Both behaviors can occur depending on regime. Therefore, we track spectrum summaries over time and under corruption sweeps.

**Hubness and many-to-one matching increase.** A common compromise under corruption is to create hubs: map many items to a few anchor embeddings so that mismatched pairs are not too far apart. This reduces the penalty for incorrect positives and reduces the variability of similarity scores. In retrieval heatmaps, hubness appears as stripes. In neighbor counts, some targets become top-1 for many queries. Operationally, this corresponds to repetitive outputs and reduced discriminability.

**Directional asymmetry can appear.** Corruption can induce asymmetry if one modality provides stronger anchors. For example, if text features are cleaner, the model may map image embeddings

to a small set of text anchors to satisfy inconsistent pairings, producing image→text degradation in a different pattern than text→image.

**Sensitivity increases.** Influence-style proxies often show that small perturbations in one batch or one update step can significantly change retrieval. This is consistent with optimization in an inconsistent landscape: small changes can flip which compromise the model prefers.

These signatures are important because pairing corruption may not be obvious from data inspection. In production, join bugs can be subtle. Therefore, the system must monitor for the geometric and optimization signatures that corruption induces.

### 2.6.4   Robustness margins: tolerable corruption levels as stage gates

Because pairing corruption is a realistic risk, a professional system should define tolerable corruption margins. This means: how much pairing noise can the model tolerate before the representation contract becomes unusable? The answer is application-dependent. A retrieval system for casual browsing may tolerate more noise than a medical or financial compliance system. But regardless of the application, the margin must be measured, not assumed.

In the synthetic lab, we define robustness margins by sweep experiments. We increase $\rho$ and measure:

- retrieval@$1$ and retrieval@$5$ in both directions,
- symmetry gaps,
- factor separability ratios,
- collapse indicators (mean cosine, variance floors, effective rank),
- hubness metrics (top-1 frequency distribution),
- gradient norm stability and loss behavior.

We then define a stage gate: the maximum $\rho$ such that all constraints remain within bounds. This is a multi-constraint definition, not a single-metric threshold. The margin is then recorded as part of the deliverables: "This baseline system remains within bounds up to $\rho = 0.05$ random swaps," for example. Importantly, this statement remains "Not verified" beyond the synthetic setting, but it teaches a professional habit: robustness must be quantified as a margin.

Robustness margins also highlight that different constraints fail at different points. Retrieval may remain acceptable up to a higher $\rho$, while factor separability may degrade earlier. Hubness may appear before retrieval collapses. These lead-lag relationships are critical for monitoring: if hubness rises, it may be an early indicator that corruption is increasing even if retrieval remains okay.

In production, robustness margins become operational requirements. Data pipelines can be audited to ensure that expected pairing error rates remain below the margin. When changes occur (new join logic, new ingestion), the pipeline can be validated by re-running synthetic or proxy tests that estimate effective pairing integrity.

### 2.6.5 Controls: data-join audits, filtering, robust objectives, monitoring

Controls against pairing corruption must operate at multiple layers: data pipeline, training objective, and monitoring.

**Data-join audits.** The primary control is to audit pairing construction. This includes verifying join keys, time alignment, and deduplication logic; checking that IDs are consistent; and performing sanity checks on random samples. In production, this is often more effective than any training trick, because the corruption originates upstream.

**Filtering and consistency checks.** Apply heuristics to filter suspicious pairs: extremely low similarity under a pretrained model, inconsistent metadata, or mismatched lengths. In synthetic labs, we can simulate filtering by removing corrupted pairs. The key lesson is that pair integrity is a data governance concern.

**Robust objectives.** There are objective-level mitigations: use multiple positives per anchor, use soft positives, or incorporate clustering-based consistency. However, robust objectives can also introduce new failure modes. Therefore, any objective modification must be evaluated under the same governance panel: do not trade corruption robustness for confounder vulnerability or collapse risk without evidence.

**Monitoring for corruption signatures.** Because corruption can appear due to pipeline drift, the system must monitor for its signatures: loss instability patterns, rising hubness, decreasing effective rank, and increased sensitivity. In production, these monitors can be computed on held-out validation streams or on periodic audits. The key is that corruption is detected indirectly through geometry and optimization behavior.

**Stage gates and rollback discipline.** If a new dataset or pipeline version causes corruption signatures to rise, the system should fail a stage gate and trigger rollback. This is a governance posture: do not allow silent coupling breakdown into production.

The core lesson of pairing corruption is that multimodal learning assumes a coupling relation, and broken coupling creates contradiction. The model will not refuse to train; it will find degenerate compromises. Therefore, pairing integrity must be audited and monitored as a first-class system property. With this failure mode established, we move to Failure Mode D: representation collapse, where geometry degenerates under sharp competition and the embedding space loses diversity, often producing hubness and instability even under clean pairs.

## 2.7 Failure Mode D: Representation Collapse

### 2.7.1 Mechanism: degenerate solutions under sharp competition

Representation collapse is the failure mode in which the embedding space loses effective degrees of freedom and ceases to represent the intended latent structure in a stable way. In the most extreme case, all embeddings become identical, making retrieval impossible. In contrastive learning with negatives, full collapse is often avoided, but partial collapse and anisotropic collapse are common. These forms are more dangerous precisely because they can coexist with decent retrieval on familiar distributions. A system can "work" while its geometry becomes brittle, hub-dominated, and ungovernable.

The mechanism can be understood as a consequence of sharp competition under the InfoNCE objective. Recall that the model is rewarded for making the positive similarity large relative to negatives. When the temperature $\tau$ is low, the softmax distribution becomes sharp: a small difference in similarity creates a large difference in probability mass. This sharpness amplifies gradients and can push the system toward extreme configurations in which a few directions dominate. The model may discover that it can reduce loss by concentrating representations so that positives are always slightly higher than negatives, even if that concentration reduces global diversity.

More formally, consider normalized embeddings on the unit sphere. The objective encourages paired points to align and non-paired points to separate. But separation under high negative density can be achieved by forming clusters around a small number of anchor directions and assigning pairs to the same anchor. This does not yield perfect separation, but it can reduce loss by creating a predictable similarity structure. Such anchoring creates hubness: anchors become nearest neighbors for many queries. In the limit, the space becomes a set of hubs plus small perturbations.

Collapse is also influenced by optimization regime. High learning rates can overshoot and create oscillations that the model resolves by shrinking variance. Strong weight decay can reduce representational capacity and encourage smaller effective subspaces. Saturating nonlinearities can reduce gradient diversity. Batch composition can make negatives too hard or too similar, encouraging the model to adopt brittle configurations. Therefore, collapse is not a single cause; it is an outcome of a regime. This is why the chapter treats collapse induction as a controlled experiment: by varying temperature, learning rate, and regularization, we can push the system into collapse and observe signatures.

The key frontier lesson is that contrastive learning does not only learn semantics; it learns geometry under competition. When competition is too sharp or optimization too aggressive, geometry can degenerate. The correct professional response is not to assume collapse will not happen, but to instrument the system so that collapse is detectable early and preventable through controls.

### 2.7.2 Inducing collapse: low temperature, high LR, strong regularization

To teach collapse as a mechanistic phenomenon, we induce it through three primary knobs, each of which corresponds to realistic production or training changes.

**Low temperature.** We reduce $\tau$ below the baseline regime. This increases similarity sharpness and can push the model toward anisotropic configurations. In some regimes, retrieval improves initially because positives become more distinct. But the cost is that the embedding distribution can concentrate, reducing effective rank. Students learn that "retrieval improvement" can be an early warning, not a success.

**High learning rate.** We increase the learning rate while keeping other settings fixed. This can create unstable updates and can cause the system to settle into degenerate minima or limit cycles. A common stabilizing response of the optimization is to reduce variance and create hubs, because such configurations reduce gradient variability. High learning rates therefore can induce collapse-like geometry even when temperature is moderate.

**Strong regularization.** We increase weight decay or add explicit L2 penalties on parameters or embeddings. Without variance-preserving terms, strong regularization can reduce the ability of the model to represent multiple factors, pushing embeddings into a low-dimensional subspace. This can appear as rank collapse. In real systems, this corresponds to aggressive regularization to prevent overfitting or to meet compression requirements, but it can backfire by inducing degeneracy.

We also consider combined regimes because collapse often emerges from interactions: low temperature plus high learning rate is particularly dangerous. The notebook explores a small factorial grid of these knobs to map where collapse occurs. The important pedagogical point is that collapse is a regime boundary phenomenon. There is often a threshold at which effective rank and hubness change rapidly. Mapping these thresholds teaches students to think in terms of margins and to design training to stay within safe regimes.

### 2.7.3 Indicators: mean cosine rise, variance fall, rank collapse

Collapse must be detected through quantitative indicators that reflect geometry rather than narrative interpretation. The chapter emphasizes a small set of core indicators that are computed per run and tracked over time.

**Mean off-diagonal cosine similarity.** Given normalized embeddings, compute the average cosine similarity between distinct pairs in a batch or dataset. In a healthy representation, this mean should be near zero (or slightly negative) depending on distribution. When collapse or anisotropy emerges, embeddings become more aligned and the mean increases. This is a direct measure of concentration on the unit sphere.

**Embedding variance floors.** Compute coordinate-wise variance of embeddings (after centering)

and track the minimum or average variance. If variance collapses, embeddings lose dynamic range. Variance floors are particularly important because normalization can hide magnitude shrinkage; variance captures directional diversity.

**Covariance spectrum and effective rank.** Compute the eigenvalues (or singular values) of the embedding covariance matrix. A healthy representation has a spectrum that spreads across multiple dimensions. Collapse manifests as a spectrum concentrated in few dimensions. Effective rank provides a scalar summary:

$$r_{\text{eff}} = \exp\Big(-\sum_i p_i \log p_i\Big), \qquad p_i = \frac{\lambda_i}{\sum_j \lambda_j}.$$

A falling effective rank is a strong collapse indicator.

**Hubness metrics.** Count how often each target embedding appears as the top-1 neighbor across queries. In a healthy space, the distribution is relatively uniform. In a collapsed or hub-dominated space, a few targets dominate. Hubness can be quantified by the Gini coefficient of the top-1 counts or by the fraction of queries mapped to the top $m$ hubs.

**Similarity heatmap structure.** While not a numeric indicator alone, heatmap patterns provide a visual diagnostic: stripes indicate hubs; block-diagonal structures may indicate confounders; uniformity indicates collapse. The key is to treat the heatmap as evidence, not as a substitute for metrics.

Students learn that collapse is rarely sudden in loss metrics. Loss can decrease even as rank collapses. Therefore, collapse indicators must be part of acceptance gates. A model that achieves high retrieval but low effective rank fails the representation contract.

### 2.7.4 Spectrum analysis: singular values, effective rank, anisotropy

Because collapse is fundamentally a spectrum phenomenon, the chapter dedicates explicit attention to spectrum interpretation. Students often see eigenvalue plots without understanding what they mean. Here, the meaning is operational: the spectrum describes how many independent directions the model uses to encode variability. If the model encodes only a few directions, it cannot represent multiple factors smoothly and is vulnerable to drift.

We emphasize three spectrum-derived concepts:

**Effective rank vs raw dimension.** An embedding may be $d = 64$ dimensional, but if effective rank is $r_{\text{eff}} = 8$, the representation is effectively low-dimensional. This explains why models can become brittle: many nominal dimensions are unused.

**Anisotropy vs collapse.** Anisotropy means variance is uneven across directions; collapse is extreme anisotropy. Moderate anisotropy is not always bad; some factors may naturally dominate. The question is whether anisotropy becomes so strong that neighborhoods become hub-like and factor

encoding is lost. Therefore, we interpret anisotropy through correlation with factor separability: if anisotropy increases while factor separability improves and remains stable under stress, it may be acceptable. If anisotropy increases while separability stagnates and hubness rises, it is a failure signature.

**Spectrum under stress.** We evaluate how the spectrum changes under noise asymmetry and pairing corruption. A healthy representation should degrade smoothly. A fragile representation will show sharp spectrum contraction under small perturbations. This is a sign that the geometry lacks margins.

This spectrum-focused interpretation prepares students for Chapter 3, where drift detection relies on monitoring distributional changes. In production, covariance spectra can shift due to data drift. If the system is already near a collapse boundary, small drift can push it into hubness. Therefore, spectrum monitoring is not academic; it is operational.

### 2.7.5 Controls: variance/covariance regularizers, temperature discipline, stop rules

Controls for collapse must be treated as a minimum control set for any production-grade contrastive system.

**Temperature discipline.** Choose temperature within a safe regime defined by multi-metric acceptance criteria. Do not chase retrieval improvements by lowering temperature without monitoring effective rank and hubness. In production, treat similarity scaling as part of the contract: changes to normalization, quantization, or scoring must be evaluated as temperature shifts.

**Variance-preserving regularizers.** Use objectives that explicitly preserve variance and reduce redundancy. Examples include variance/covariance regularization approaches (e.g., VICReg-style principles) or redundancy reduction (Barlow Twins-style principles), adapted to multimodal settings. The key conceptual control is: ensure embeddings cannot reduce loss by shrinking diversity. Even if we do not implement these full methods in the minimal notebook, we model the principle through explicit variance floors and stop rules.

**Stop rules and health gates.** Define stop conditions: if effective rank falls below a threshold, or if mean cosine rises above a ceiling, stop training and record failure. This is not optional. In governed systems, training that produces unhealthy geometry is a failed run, not a result to salvage.

**Learning rate and optimizer discipline.** Use stable learning rates, warmup schedules if needed, and gradient clipping. Monitor gradient norms and reduce LR if instability appears. Collapse often emerges under unstable optimization; controlling optimization reduces collapse risk.

**Batch design and hard negative caution.** Hard negatives can improve discriminability but can also induce anisotropy and collapse. Therefore, negative sampling and batch composition must be treated as part of the geometry contract. Evaluate changes through the same spectrum and hubness

monitors.

The core lesson of collapse is that a multimodal system can satisfy a contrastive objective while losing representational meaning through geometric degeneration. Collapse is induced by regimes of sharp competition and unstable optimization; it is detected by mean cosine, variance, spectra, and hubness; and it is controlled through temperature discipline, variance preservation, and explicit stop rules. With dominance, confounding, corruption, and collapse established, we now introduce the structural probes that allow us to attribute mechanisms beyond what retrieval metrics reveal.

## 2.8 Structural Probes for Mechanism Attribution

### 2.8.1 Mutual information proxy: discretization, entropy, and factor dependence

Retrieval metrics tell us whether matched pairs outrank mismatched pairs under a particular distribution. They do not tell us *what* the embeddings encode. In failure analysis, that distinction is critical. Spurious alignment can raise retrieval by encoding confounders; dominance can preserve retrieval by mapping one modality into the other's anchors; collapse can sharpen retrieval while reducing diversity. To attribute these behaviors, we need probes that connect embeddings to latent factors. The first probe is a mutual-information (MI) proxy between embeddings and the true latent factors in the synthetic world.

In principle, mutual information $I(Z; F)$ measures how much information the representation $Z$ contains about a factor $F$. If $I(Z; F)$ is high, then the embeddings are predictive of that factor. In continuous high-dimensional settings, exact MI estimation is difficult. However, in a synthetic pedagogical notebook, we can compute a robust proxy by discretization and empirical entropy estimation.

The method is:

1. Choose a set of embedding statistics $u(Z)$, such as projections onto top principal components or a small set of random directions.
2. Discretize these statistics into bins. For example, quantize each projection into $m$ bins by quantiles.
3. Treat the discretized embedding signature as a discrete variable $\tilde{Z}$.
4. Compute empirical entropies $H(\tilde{Z})$, $H(F)$, and joint entropy $H(\tilde{Z}, F)$ by counting frequencies.
5. Estimate MI proxy:
$$\widehat{I}(\tilde{Z}; F) = H(\tilde{Z}) + H(F) - H(\tilde{Z}, F).$$

Because $F$ can be multivariate (shape, orientation, etc.), we compute MI proxies separately for each factor and optionally for factor subsets. We also compute MI between embeddings and known confounders when present. The comparison is the point: if MI with confounders is high and MI

with intended factors is low, spurious alignment is strongly indicated.

We emphasize that this is a proxy, not a perfect estimator. Discretization loses information. Empirical entropy estimation is biased in small samples. Nevertheless, the proxy is pedagogically powerful because it gives students a quantitative connection between embeddings and semantics. It also supports attribution: collapse tends to reduce MI with all factors; dominance may preserve MI in one modality but reduce it in the other; confounding increases MI with confounders.

In governance terms, MI proxies are used as *diagnostic gates*, not as final proofs. The notebook records them as facts and includes open items: how would MI estimation behave in real data, where factors are unknown and confounders are hidden? The point is to teach students what kind of evidence is needed.

### 2.8.2 CCA-like correlation proxy: principal correlations via SVD

The second probe focuses on *cross-modal shared structure.* Canonical correlation analysis (CCA) measures linear relationships between two sets of variables. In multimodal alignment, CCA-like measures can indicate whether the two modalities share a coherent subspace. This is particularly useful for diagnosing dominance and collapse, and for distinguishing confounding from semantic alignment.

We compute a CCA-like proxy without external libraries using singular value decomposition (SVD). Given centered embedding matrices $Z_{\text{img}} \in \mathbb{R}^{N \times d}$ and $Z_{\text{txt}} \in \mathbb{R}^{N \times d}$, we whiten each:

$$\tilde{Z}_{\text{img}} = Z_{\text{img}} C_{\text{img}}^{-1/2}, \qquad \tilde{Z}_{\text{txt}} = Z_{\text{txt}} C_{\text{txt}}^{-1/2},$$

where $C_{\text{img}} = \frac{1}{N-1} Z_{\text{img}}^\top Z_{\text{img}}$ and similarly for text, and $C^{-1/2}$ is computed via eigen-decomposition with regularization. Then compute the cross-covariance

$$M = \frac{1}{N-1} \tilde{Z}_{\text{img}}^\top \tilde{Z}_{\text{txt}}.$$

The singular values of $M$ are the principal correlations. Large leading singular values indicate strong shared linear structure; a flat spectrum indicates weak coupling.

This probe supports attribution in multiple ways:

- **Dominance:** cross-modal correlations may remain high even when one modality's within-modality spectrum collapses, indicating that coupling is carried by a low-dimensional shared subspace dominated by one modality's anchors.
- **Confounding:** a strong confounder can create strong cross-modal correlations concentrated in a small number of dimensions, even if semantic factor separability is low.
- **Collapse:** as effective rank collapses, whitening becomes unstable and principal correlations can become artificially high or low; tracking the stability of this probe itself becomes informative.

- **Pairing corruption:** correlations may degrade or become noisy as coupling breaks, and sensitivity may increase.

We stress interpretive caution: CCA-like correlation does not prove semantic correctness. It measures shared structure, which can come from confounders. Therefore, the probe must be interpreted jointly with MI proxies and factor separability.

Pedagogically, this probe teaches students to think about multimodal alignment as subspace coupling. Even in nonlinear systems, much of the observable alignment behavior can be explained through how strongly the modalities share a low-dimensional structure. This prepares students for real-world systems, where shared structure is often the first thing that drifts when pipelines change.

### 2.8.3   Influence/sensitivity proxy: one-step perturbation impact on retrieval

The third probe addresses a different question: how sensitive is the system to a small training perturbation? In classical statistics, influence functions estimate how much a parameter estimate changes when a training point is upweighted. In deep learning, full influence function computation is expensive. However, in our minimal NumPy-based system, we can approximate sensitivity through a controlled one-step finite-difference experiment.

The idea is:

1. Start from current parameters $(\theta, \phi)$.
2. Choose a batch $B$ and compute the baseline update $\Delta$ under the optimizer (e.g., one gradient step).
3. Construct a perturbed batch $B'$ that differs from $B$ in a controlled way: swap a fraction of pairs, add confounder strength, or add noise.
4. Compute the perturbed update $\Delta'$.
5. Apply each update to parameters and evaluate retrieval@1 on a fixed validation set:

$$R = \text{retrieval@1}(\theta + \Delta, \phi + \Delta), \quad R' = \text{retrieval@1}(\theta + \Delta', \phi + \Delta').$$

6. Define the sensitivity proxy as $|R' - R|$ (or a vector of metric differences).

This proxy is not a perfect influence function, but it is highly informative in failure regimes. Under clean and stable training, small perturbations of a batch should not drastically change retrieval after one step. Under pairing corruption, the gradient field becomes contradictory, and sensitivity increases. Under near-collapse, the geometry is brittle and small updates can flip hub structure, causing large retrieval changes. Under dominance, the sensitivity may be concentrated in one modality: perturbing the weaker modality may have little effect because it contributes little to the update.

Therefore, this probe connects failure modes to *training fragility.* Training fragility is a critical production risk because it implies that continued training or fine-tuning can produce unstable behavior. If a model is near a failure boundary, minor data shifts can cause sudden behavioral changes. Influence proxies provide an early warning.

Pedagogically, this probe teaches a deep lesson: robustness is not only about test accuracy under fixed parameters; it is also about stability of the learning dynamics. A system that is highly sensitive to small perturbations is difficult to govern, because it can drift quickly under routine updates.

### 2.8.4 Probe interpretation: what each probe can and cannot prove

A governance-first approach requires explicit limits. Probes provide evidence, but they do not provide certainty. Therefore, we state what each probe can and cannot prove.

**MI proxy can:** indicate whether embeddings contain information predictive of known latent factors; compare relative encoding of factors vs confounders; detect loss of factor information under collapse.

**MI proxy cannot:** prove semantic understanding; estimate MI accurately in continuous high-dimensional settings; guarantee invariance under unseen shifts.

**CCA-like proxy can:** quantify strength and dimensionality of shared linear structure; detect subspace coupling changes under dominance or corruption; provide a compact signature for drift.

**CCA-like proxy cannot:** distinguish semantic coupling from confounder coupling; capture non-linear relationships fully; remain stable under extreme rank collapse without careful regularization.

**Influence proxy can:** detect training fragility; distinguish stable from brittle regimes; reveal whether perturbations in batches have outsized effects on retrieval.

**Influence proxy cannot:** provide a full causal attribution; predict long-run training behavior from one-step changes; substitute for full stability analysis.

Students should learn to treat probes as a panel. A robust attribution uses multiple probes and multiple signatures. For example, spurious alignment is supported when retrieval is high, MI with confounders is high, MI with intended factors is low, and counterfactual ablations collapse performance. Dominance is supported when symmetry gaps exist, gradient norms are imbalanced, and one modality's effective rank collapses. Collapse is supported when mean cosine rises, effective rank falls, hubness rises, and MI with factors declines broadly.

### 2.8.5 Probe governance: preventing "probe hacking" and false confidence

Finally, probes themselves can be gamed, intentionally or unintentionally. In research, one can tune hyperparameters to improve a probe metric, just as one can tune to improve retrieval. In

production, a team can over-rely on a single probe and ignore contradictory evidence. Therefore, we include probe governance principles.

**Pre-specify probe thresholds.** Define acceptance criteria for probe metrics before experiments. Do not tune probes post hoc to tell a desired story.

**Use multiple probes.** Require coherence across probes. If a probe conflicts with others, treat it as an open item requiring investigation.

**Audit sensitivity.** For discretization-based MI, test multiple bin counts and report stability. For CCA-like proxies, report sensitivity to regularization strength. Probes that change drastically under minor technical choices are not reliable.

**Record assumptions and open items.** The notebook reports must explicitly label probes as proxies and list questions to verify. This prevents probe metrics from being mistaken for certification.

**Tie probes to interventions.** The strongest probe evidence comes from controlled interventions. For example, if confounder strength increases and MI with confounder increases while MI with factors decreases, that is a causal pattern. Without interventions, probe changes are ambiguous.

With structural probes defined, the chapter now has the tools needed to attribute mechanisms, not just outcomes. The next section specifies the overall experiment suite and reporting standard: how we design sweeps, track early warning dashboards, define acceptance criteria as stage gates, and export strict JSON artifacts that preserve reviewability. This is where the failure taxonomy becomes an operational evaluation workflow.

## 2.9 Experiment Suite and Reporting Standard

### 2.9.1 Experiment design: orthogonal sweeps and factorial structure

A failure taxonomy becomes professionally useful only when it is operationalized into an experiment suite. The experiment suite is the mechanism by which we convert qualitative failure stories into quantitative robustness margins and actionable monitoring signatures. The core design principle is *orthogonality*: we vary one structural knob at a time when we want attribution, and we vary a small factorial set of knobs when we want to map interactions and cliffs.

We begin with a fixed baseline configuration: synthetic world parameters, encoder architectures, optimization settings, temperature, batch size, and split. This baseline is accepted under health gates. Every experiment is then defined as a perturbation of that baseline and is recorded in the run manifest. In a governed lab, the experiment definition is as important as the result, because results without a clear definition are not reproducible evidence.

The suite is organized into four families aligned with the failure modes:

- **Dominance sweeps:** gain scaling factors, noise asymmetry levels, capacity skew levels.
- **Confounding sweeps:** confounder type (ID codebook, parity marker, watermark), confounder strength, confounder prevalence.
- **Pairing corruption sweeps:** corruption rate $\rho$, corruption pattern (random, structured, offset, block), corruption location (train-only, test-only, both).
- **Collapse induction sweeps:** temperature, learning rate, weight decay, and small combinations thereof.

Within each family, we structure experiments to answer two questions:

**(1) Degradation curves.** How does performance and representation health degrade as the knob increases? This yields a curve, not a point. Curves allow margin definitions. If a system degrades smoothly, it is more governable than if it exhibits sudden cliffs.

**(2) Signature evolution.** Which indicators move first? For example, does hubness rise before retrieval drops? Does effective rank fall before MI proxy declines? These lead-lag relationships define early warning dashboards.

We also run cross-evaluations because some failures are distribution-dependent. For confounding, we train with a confounder and evaluate with confounder removed. For pairing corruption, we train with corruption and evaluate under clean pairing, and vice versa. For dominance, we evaluate under symmetric noise and asymmetric noise. These cross-evaluations emulate production shifts: what happens when a watermark disappears, or when pairing quality changes due to pipeline updates?

Finally, we repeat key sweeps across multiple random seeds to avoid a single-seed narrative. Seed variability is not noise; it is part of the system's stability profile. A model that is stable across seeds is more governable. A model that produces divergent geometries across seeds is risky.

### 2.9.2 Early warning dashboard: variance trend, mean cosine, gradient stats

A key output of Chapter 2 is an early warning dashboard: a set of time-series indicators tracked during training and summarized at evaluation. The dashboard is not a visualization preference; it is a governance tool. It provides the basis for stage gates and for post-mortem analysis.

The minimum early warning indicators include:

**Embedding variance trend.** Track per-step (or per-epoch) embedding variance per modality. We record mean variance, minimum variance, and variance distribution summaries. A downward trend signals collapse risk or dominance-induced under-utilization.

**Mean cosine trend.** Track mean off-diagonal cosine similarity for each modality and cross-modal similarities. Rising mean cosine indicates anisotropy/collapse; cross-modal mean can indicate spurious shortcut concentration.

**Covariance spectrum snapshots.** Periodically compute singular values and effective rank. This

can be done less frequently than variance due to cost. A sudden drop in effective rank is a strong warning.

**Gradient norm statistics.** Track gradient norms per modality and per layer. Persistent imbalance indicates dominance. Spikes indicate instability and potential collapse boundary behavior. Under pairing corruption, gradient variability often increases.

**Retrieval trend on validation.** Track directional retrieval@1 and retrieval@5 during training. While retrieval is not sufficient, its trend combined with geometry indicators is informative. For example, retrieval rising while effective rank falling is a classic collapse-like warning.

We record these indicators both as time series (for audit and analysis) and as summary statistics (for quick review). The notebook saves plots of trends and exports a JSON summary that includes threshold violations.

Pedagogically, this dashboard teaches students that representation learning must be monitored like a control system. You do not train and then inspect only the final metric. You watch the system evolve and you detect when it enters an unhealthy regime. This mindset is essential for production: training is a dynamical process, and failure is often visible before it becomes catastrophic.

### 2.9.3 Stress curves: degradation surfaces and robustness margins

Beyond time-series monitoring, Chapter 2 emphasizes stress testing as a surface-mapping exercise. In Chapter 1, we introduced noise sweeps as a pedagogical device. In Chapter 2, stress curves become the core evidence for failure mode characterization.

A stress curve is a mapping from a structural knob $s$ (e.g., confounder strength) to a set of metrics $m(s)$. A stress surface extends this to two knobs $(s_1, s_2)$. The key is that stress tests produce *response functions*. These functions define robustness margins and reveal interactions.

For example, consider dominance under noise asymmetry. We can sweep image noise $\sigma_{\text{img}}$ and text noise $\sigma_{\text{txt}}$ on a grid and measure directional retrieval. The surface reveals whether the system degrades symmetrically or whether one direction collapses quickly when one modality becomes noisy. This is directly interpretable as an operational risk: if text noise increases in production, will image→text retrieval fail rapidly?

Similarly, for confounding, we can sweep confounder strength and measure retrieval and factor MI. A spurious alignment regime appears as high retrieval at high confounder strength with low factor MI. The margin is the confounder strength at which factor MI falls below a floor. This margin is not just academic; it corresponds to the tolerance of the system to shared artifact signals.

For pairing corruption, sweeping corruption rate $\rho$ yields a robustness curve. The margin is the maximum $\rho$ such that constraints hold. The notebook exports this margin as part of a stage gate report. Importantly, different corruption patterns produce different margins. This teaches students

to ask: what corruption pattern is plausible in production? Random swaps may be less plausible than systematic offset bugs.

For collapse induction, sweeping temperature and learning rate yields a collapse boundary. The boundary can be described as a region of safe operation. In production, changes in embedding scaling or scoring can effectively change temperature; changes in fine-tuning regime can change learning rate. The collapse boundary therefore defines acceptable training regimes.

Pedagogically, stress surfaces teach students to think like engineers and scientists: map response functions, not stories. This is how frontier work becomes teachable and governable.

### 2.9.4 Acceptance criteria: distributional gates over seeds and shocks

A governed evaluation requires acceptance criteria that are explicit, multi-metric, and distributional. Chapter 2 therefore defines stage gates not as a single number but as a set of constraints that must hold across seeds and across stress conditions.

A typical acceptance gate includes:

- **Performance constraints:** retrieval@$1$ and retrieval@$5$ above thresholds in both directions under baseline and mild stress.
- **Symmetry constraints:** directional gaps below a threshold; norm distributions comparable; effective rank divergence bounded.
- **Health constraints:** mean cosine below ceiling; variance above floor; effective rank above floor; hubness below ceiling.
- **Semantic constraints:** factor separability above floor; MI proxy with intended factors above floor; MI with known confounders below ceiling when confounders are present.
- **Stability constraints:** training does not produce NaNs; gradient norms within bounds; influence proxy below a sensitivity ceiling under clean regimes.

These constraints are evaluated across multiple seeds. A model that meets constraints for one seed but fails for others is not accepted. This is not pedantry. Seed sensitivity indicates an unstable learning process, which is a production risk.

Constraints are also evaluated under stress. For example, we define a mild noise asymmetry regime and require that directional retrieval remains within bounds. We define a small corruption regime and require that collapse indicators remain healthy. These stress conditions are pre-specified and recorded. The gate is therefore a *distributional* claim: the model remains within bounds across a set of plausible perturbations. This is the appropriate standard for frontier systems that will face drift.

We emphasize that these gates are not meant to certify real-world safety. They are meant to enforce discipline. The notebook's strict JSON reports label verification status as "Not verified" and include questions to verify, such as: how do these gates map to real data shifts? What confounders exist in

production? What corruption patterns are plausible? The gates are a teaching device for governance thinking.

### 2.9.5 Strict reporting: facts vs assumptions vs open items (Not verified)

The final component of the experiment suite is the reporting standard. Without a reporting standard, the notebook produces plots and numbers but not reviewable evidence. Therefore, every experiment exports strict JSON reports that separate epistemic categories:

- **facts_provided:** measurements, metrics, thresholds, plots generated, paths to artifacts, parameter settings from manifest.
- **assumptions:** interpretations that rely on synthetic design choices, proxy assumptions (e.g., discretized MI), and mapping assumptions to production.
- **open_items:** unresolved questions and follow-ups, including potential confounders not tested, additional stress conditions, and production validation needs.
- **analysis:** structured explanation of what happened and which failure mode signatures were observed, referencing facts by name.
- **draft_output:** a narrative summary intended for teaching, written in a way that can be included in course materials.
- **verification_status:** always `"Not verified"` in this synthetic lab context.
- **questions_to_verify:** explicit questions that would need real-data tests or deeper analysis.

This structure prevents a common failure in frontier education: conflating demonstration with proof. The notebook can demonstrate failure modes and diagnostics, but it cannot prove that a given production system is safe. The reporting standard forces that humility into the artifacts.

Moreover, the reporting standard enables reproducibility. Each report references the run manifest and log hashes. This allows a reviewer to reproduce the run and verify that the reported metrics match. This is governance as infrastructure.

With the experiment suite and reporting standard defined, we have turned the failure taxonomy and probes into an operational workflow. The chapter can now conclude by summarizing the core lessons, stating the minimum control set, acknowledging limits, and transitioning to Chapter 3, where the focus shifts from induced failures to lifecycle drift and monitoring under changing environments.

## 2.10 Conclusion and Transition to Chapter 3

### 2.10.1 Core lesson: alignment can fail while metrics improve

The central lesson of Chapter 2 is uncomfortable but professionally essential: multimodal alignment can fail while headline metrics improve. Retrieval@1 can rise because a confounder provides an

easier matching key than the intended semantics. Directional averages can look healthy while one direction quietly collapses under modality dominance. Loss can decrease while effective rank contracts and hubness rises, producing a brittle space that will break under drift. These are not exotic events. They are the natural outcomes of an optimization process that is rewarded for ranking positives above negatives, not for preserving meaning, symmetry, or stability.

In a governed setting, we therefore treat retrieval as a necessary but insufficient indicator. The representation contract is broader: it requires symmetry, diversity, factor dependence, and stability under perturbation. Chapter 2 showed that each of these obligations can be violated without triggering obvious alarms in the usual metrics. This is why the pedagogical posture changes from "train and admire" to "train and interrogate." The student learns to ask: what is the model using, what would break it, and what would I observe if it broke?

## 2.10.2 What the failure taxonomy enables in production monitoring

A taxonomy is valuable only if it enables action. The taxonomy developed here turns vague worries ("shortcuts," "collapse," "data issues") into a structured monitoring plan. Each failure mode maps to a distinctive signature panel: dominance maps to gradient imbalance, norm and spectrum mismatch, and directional retrieval asymmetry; spurious alignment maps to high retrieval with low factor proxies and counterfactual fragility; pairing corruption maps to instability signatures, hubness growth, and sensitivity spikes; collapse maps to rising mean cosine, falling variance floors, contracting spectra, and hub dominance.

In production, we rarely have access to true latent factors, but we do have access to distributional indicators and invariance tests. The taxonomy tells us what to monitor even without semantics labels. For example, if hubness rises and effective rank falls, we have a collapse warning regardless of domain. If directional gaps widen and gradient norms diverge, we have dominance risk regardless of content type. If sensitivity to small fine-tuning changes increases, we have a coupling instability warning that may be driven by pairing quality or by a near-boundary geometry. The taxonomy therefore becomes a lifecycle instrument: it helps teams detect when a model is entering a failure regime and guides the investigation toward plausible mechanisms instead of ad hoc debugging.

## 2.10.3 From diagnosis to control: minimum control set

Diagnosis is not enough. A production-grade posture requires controls: mechanisms that prevent silent failure, or at least force failure to be visible and reviewable. Chapter 2 motivates a minimum control set that should be treated as non-negotiable for multimodal representation systems.

First, enforce reproducibility and artifact discipline: deterministic seeds for controlled tests, versioned preprocessing, and auditable run manifests that record every structural knob. Second, enforce split hygiene and leakage audits that go beyond train/test separation: explicitly search for shared

IDs, metadata, and formatting channels that can become confounders. Third, enforce symmetry constraints as acceptance gates: do not accept a model that improves average retrieval while failing directional parity or distributional parity. Fourth, enforce representation health gates: variance floors, effective rank floors, mean cosine ceilings, and hubness ceilings, with stop rules that terminate training when geometry becomes unhealthy. Fifth, enforce invariance and counterfactual tests for known shortcut families: ablations, randomizations, and cross-source evaluations. Finally, enforce stress sweeps as part of release: measure robustness margins against noise asymmetry, small pairing corruption, and scoring-temperature shifts.

These controls are not burdensome bureaucracy; they are the minimum infrastructure required because embeddings become hidden policy. Without controls, a system can drift into failure while appearing to improve. With controls, failure becomes a visible event with traceable causes.

### 2.10.4 Limits: synthetic evidence and translation risks

This chapter's evidence is synthetic by design. That choice buys interpretability and causal control, but it imposes limits. Synthetic worlds encode a particular factorization of semantics; real worlds do not. Synthetic confounders are explicit; real confounders are subtle, entangled, and often unknown. Synthetic pairing corruption is parameterized; real pairing failures are pipeline-specific and can be correlated with source, time, or user behavior. Even collapse signatures can manifest differently in large-scale models with different normalization and optimization regimes.

Therefore, the correct epistemic label remains "Not verified." The point of the chapter is not to certify any real system. It is to teach a method: how to induce, detect, and interpret failure modes in a controlled setting, and how to translate that method into production monitoring plans. The main translation risk is overconfidence: believing that passing synthetic gates implies safety in the wild. The appropriate stance is the opposite: synthetic success provides a baseline of competence, and real-world deployment requires additional domain-specific validation, threat modeling, and continuous monitoring.

### 2.10.5 Transition: drift, monitoring, and lifecycle governance (Chapter 3)

Chapter 2 treated failures as objects we induce. Chapter 3 treats failures as states we drift into. In real systems, collapse rarely arrives as an explicit hyperparameter mistake; it arrives through gradual distribution changes, scoring changes, data pipeline updates, and incremental fine-tuning. Confounders appear and disappear as sources change. Pairing integrity degrades as joins evolve. Modality dominance emerges when one sensor improves or one metadata pipeline becomes cleaner than the other. These are lifecycle phenomena.

Therefore, Chapter 3 extends the same governance-first posture into monitoring and intervention. It asks: what signals should be tracked continuously, how do we distinguish benign drift from failure

drift, how do we define drift budgets and rollback triggers, and how do we produce audit artifacts that make representational policy changes reviewable? The transition is direct: Chapter 2 provided the taxonomy and signatures; Chapter 3 will provide the lifecycle control loop that keeps multimodal systems within safe operating regimes under change.

# Bibliography

[1] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, vol. 139, pp. 8748–8763, 2021.

[2] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y.-H. Sung, Z. Li, and T. Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, vol. 139, pp. 4904–4916, 2021.

[3] A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

[4] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020.

[5] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9726–9735, 2020.

[6] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, vol. 139, pp. 12310–12320, 2021.

[7] A. Bardes, J. Ponce, and Y. LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations (ICLR)*, 2022.

[8] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2:665–673, 2020.

[9] M. Radovanović, A. Nanopoulos, and M. Ivanović. Hubs in space: Popular nearest neighbors in high-dimensional data. *Journal of Machine Learning Research*, 11(86):2487–2531, 2010.

[10] A. S. Morcos, M. Raghu, and S. Bengio. Insights on representational similarity in neural networks with canonical correlation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.

# Chapter 3

# Multimodal Lifecycle

**Abstract.** We treat multimodal alignment as a lifecycle problem rather than a one-time training outcome. In modern systems, images and text are mapped into a shared embedding space, and downstream behavior is governed by geometry: similarity scores induce retrieval, clustering, ranking, and filtering. This makes embeddings a form of policy. Even when a multimodal model appears to "work" at time $t_0$, its behavior can drift as data pipelines evolve, modality noise changes, scoring scales shift, confounders appear or disappear, or incremental fine-tuning alters representation health. The core risk is silent failure: the system remains numerically stable and may even improve headline metrics while drifting into dominance, spurious alignment, pairing breakdown, or representation collapse.

This chapter develops a governance-first framework for monitoring, detection, and bounded intervention. First, we define a drift taxonomy tailored to multimodal systems: modality drift (noise and corruption asymmetry), scoring drift (temperature and normalization changes), confounder drift (shortcut signals), pairing drift (join bugs and misalignment), and optimization drift (fine-tuning regimes that trigger collapse boundaries). Second, we specify a monitoring panel that goes beyond retrieval accuracy: symmetry gaps, hubness and variance floors, covariance spectra and effective rank, and structural probes that connect embeddings to underlying factors via mutual-information and CCA-like proxies, complemented by a sensitivity proxy that measures training fragility. Third, we propose a multi-signal drift classifier that produces explicit "unknown" outcomes when evidence is insufficient, exporting debug artifacts rather than forcing false certainty.

Finally, we operationalize the framework in a synthetic drift laboratory that generates time-indexed drift episodes, executes monitoring on a frozen baseline, and applies a bounded intervention policy with stage gates and rollback discipline. The result is an auditable lifecycle report: drift events, detected mechanisms, decisions, and artifacts. The chapter's objective is pedagogical and professional: to teach how to keep multimodal embedding systems reviewable under change, without confusing demonstrations with verification.

## 3.1   Introduction

### 3.1.1   Multimodal alignment as a lifecycle problem

Multimodal alignment is often described as a triumph of training: learn a shared embedding space in which an image and its caption are close, and the system becomes capable of retrieval, ranking, and cross-modal reasoning. In practice, that framing is incomplete. It treats alignment as an outcome that can be established once and then assumed. The institutional reality is different. The learned embedding space becomes an operational policy surface: it determines what gets retrieved, what gets grouped, what gets deemed similar, and what gets ignored. The consequences are not static. They evolve because the world feeding the system evolves. Therefore, multimodal alignment is not merely a learning problem; it is a lifecycle problem.

A lifecycle view begins with a simple observation: the production system is not the model alone. It is the model plus the data pipeline, preprocessing, scoring, storage, index, and any fine-tuning processes that continue after the initial training. Each of these components can change. Images

may come from different devices, new compression schemes, or new rendering pipelines. Text may be cleaner or noisier, may include different tokenization artifacts, or may be produced by different annotation processes. Retrieval scoring may be rescaled by a downstream service, or the effective temperature may change because normalization is updated. The dataset may acquire new watermarks or lose old ones. Pairing may degrade because a join key changed, a timestamp alignment shifted, or a batch ingestion process reordered records. Any of these changes can alter the geometry of the embedding space as used in production, even if the encoder weights remain unchanged. And if the system is being incrementally fine-tuned, weight changes can compound these shifts.

In a lifecycle view, the core object of governance is not "the model" but the embedding geometry as experienced by downstream tasks. Geometry can drift in ways that are not immediately visible in headline metrics. A common pattern is the "silent drift" regime: the system remains numerically stable, the loss during fine-tuning decreases, and top-line retrieval on a familiar validation slice remains acceptable, while representation health deteriorates or shortcuts dominate. This happens because the training objective optimizes a particular matching relationship under a particular negative sampling regime. If that regime changes—because batch composition changes, because the distribution of negatives changes, because a confounder shifts—optimization can move the system toward different compromises. The compromises can be operationally damaging even if they produce improvements on narrow evaluation sets.

The lifecycle perspective also forces us to treat monitoring and intervention as part of the design, not as afterthoughts. If embedding geometry is policy, then we need policy monitoring. We need to detect when policy is drifting into a failure regime. We need decision rules about what to do when drift is detected. We need bounded interventions that can restore healthy behavior without creating new risks. We need rollback discipline when interventions fail. And we need artifacts that make these decisions reviewable by humans who are accountable for the system.

This chapter formalizes that posture in a way that is teachable and operational. We do not pretend that a small Colab notebook can replicate the complexity of frontier multimodal systems. Instead, we build a synthetic laboratory where drift can be induced, detected, and classified with causal clarity. The synthetic lab is not a substitute for production validation. It is a didactic device that builds the mental model: how drift emerges, what it looks like in the monitoring panel, why single metrics mislead, and how governance turns detection into controlled action.

### 3.1.2 Why "working at $t_0$" is not evidence of stability

A multimodal model that "works at $t_0$" typically means it passes a snapshot evaluation: retrieval performance is acceptable on a test set drawn from a particular distribution at a particular time. That is a necessary milestone, but it is not evidence of stability. Stability is a dynamic property. It concerns what happens when the distribution changes, when the pipeline changes, when the scoring

changes, when fine-tuning continues, and when the system encounters inputs outside the envelope of the initial evaluation.

The reason snapshot success is weak evidence is geometric. The embedding space defines similarity. Similarity is not absolute; it is relative to the distribution of points that populate the space. If the distribution shifts, neighborhoods shift. If a new source introduces a new cluster of embeddings, nearest neighbors for existing points can change. If the variance in one modality changes, the density of points along certain directions changes. The downstream retrieval experience can change even if the model parameters do not.

Moreover, snapshot evaluation is often narrow. Teams evaluate on data that resembles the training distribution, sometimes inadvertently including the same confounders that made training easy. If a watermark is present in all splits, retrieval can look strong while the model is relying on that watermark rather than semantics. If pairing noise is constant across splits, the model can learn a corrupted coupling that appears consistent within the dataset but fails when pairings are corrected. If scoring is calibrated for one temperature, a downstream rescaling can change effective sharpness and alter hubness. None of these issues is revealed by a single test set.

Stability also involves the training dynamics. Many real systems are updated periodically: new data is added, models are fine-tuned, indices are refreshed. If the representation is near a collapse boundary, small fine-tuning steps can produce large changes in geometry. If the system is sensitive to batch composition, a small change in the data stream can shift learning. If the model is confounder-reliant, removing or changing the confounder distribution can cause sharp performance drops. Therefore, stability must be evaluated as a function of perturbations, not only as a point estimate.

A professional analogy is risk management in finance. A portfolio that performs well under one historical window is not necessarily robust. What matters is how it behaves under stress scenarios, how sensitive it is to parameter changes, and how its exposures evolve as the regime changes. Similarly, a multimodal model that works at $t_0$ is not necessarily robust. What matters is how its geometry behaves under plausible drift scenarios and how it responds to interventions.

This is why Chapter 3 insists on lifecycle governance. We must treat "working at $t_0$" as the beginning of governance, not the end. A system that is not monitored is not stable; it is merely unobserved. A system that is not instrumented cannot be governed. And a system that cannot be governed cannot be responsibly deployed in professional contexts where accountability matters.

### 3.1.3 Drift as geometric policy change in embedding space

To govern drift, we need a precise mental model of what drift is. In this chapter, drift is defined as a change in the distributional and relational structure of embeddings that alters downstream behavior. This is not limited to changes in the encoder weights. Drift can occur purely through changes in

inputs, preprocessing, pairing, or scoring. The embedding space is the mediating representation, and drift is a change in how that space is populated or used.

There are multiple kinds of drift, and not all drift is bad. Some drift is benign: the system sees new types of content, but the geometry remains healthy and the similarity structure remains meaningful. Some drift is harmful: the system's geometry becomes dominated by shortcuts, becomes asymmetric, becomes brittle, or becomes inconsistent. The goal is not to eliminate drift; it is to detect harmful drift and maintain the system within a safe operating regime.

We therefore treat drift as a policy change. In retrieval applications, the embedding space defines what is retrieved. That is a policy decision because it determines which information is surfaced and which is hidden. In ranking applications, embeddings define ordering. That is also a policy. In clustering applications, embeddings define grouping. Again, policy. If the geometry changes, the policy changes. Therefore, drift is not merely a statistical phenomenon; it is a governance phenomenon.

This framing clarifies why we need explicit monitoring panels. We monitor directional retrieval because it is an operational policy outcome. We monitor symmetry gaps because asymmetry implies one modality is being privileged. We monitor collapse indicators because collapse implies the policy is losing representational diversity and becoming hub-dominated. We monitor spectral measures because they indicate whether the representation is becoming low-dimensional and brittle. We monitor structural probes because they indicate whether embeddings continue to encode the intended semantics rather than confounders. And we monitor sensitivity because it indicates how fragile the policy is to incremental updates.

Once drift is framed as policy change, interventions become policy controls. Temperature recalibration is a policy control because it changes similarity sharpness. Loss reweighting is a policy control because it changes which modality the system privileges during training. Pairing sanitation is a policy control because it enforces the coupling contract. Confounder ablation is a policy control because it removes shortcuts. Limited fine-tuning is a policy control because it adjusts the embedding mapping under constraints. And rollback is a policy control because it restores a known policy state when a change fails.

This is the main conceptual innovation of the chapter: to teach students to treat embeddings as policy and drift as policy drift. This makes monitoring and governance feel necessary rather than optional. It also provides a language that can be shared with non-ML stakeholders: the system's behavior is governed by an internal policy surface, and we monitor and control that policy through measurable indicators and explicit decisions.

### 3.1.4   Professional objective: reviewability, not persuasion

This course's broader governance-first posture appears again here, but with a specific emphasis. The objective in professional settings is not to persuade stakeholders that the model is good. It is to make the model's behavior reviewable. Reviewability means that a reviewer can trace what changed, what evidence was observed, what decision was made, and what the system's current operational posture is.

In frontier topics, persuasion is easy. A few impressive retrieval demos can convince an audience that the system is powerful. But a professional system cannot be built on demos. It must be built on artifacts and controls. This is why the chapter emphasizes strict JSON reports, run manifests, redacted prompt logs, risk logs, checkpoints, and plots saved to disk. These artifacts provide a basis for review. They also enforce discipline: you cannot claim that the system is stable without producing evidence across episodes, across metrics, and across interventions.

Reviewability also requires humility about verification. The notebook is synthetic. Therefore, it cannot certify any real production system. The correct label is always "Not verified." But the notebook can still teach reviewability. It can show students how to construct a monitoring panel, how to classify drift, how to justify interventions, and how to export evidence that can be reviewed. This is the pedagogical aim: to teach how to think and work in a way that is compatible with accountability.

The contrast with much of the multimodal literature is instructive. Academic papers often report aggregate metrics on benchmark datasets. Those metrics are valuable, but they are not governance. They do not tell you what happens when the pipeline changes. They do not tell you how to detect confounder drift. They do not tell you what a rollback decision looks like. They rarely include an explicit taxonomy of failure modes with interventions. Chapter 3 complements that literature by focusing on lifecycle governance: the operational layer that is required when systems move from research to institutional contexts.

### 3.1.5   Chapter contributions and deliverables

This chapter's contributions are structured around a synthetic laboratory that makes lifecycle governance concrete. The key deliverables are:

**(1) A drift taxonomy tailored to multimodal systems.** We define modality drift, scoring drift, confounder drift, pairing drift, and optimization drift. Each drift type maps to known failure regimes from Chapter 2: dominance, spurious alignment, pairing corruption, and collapse.

**(2) A monitoring panel with early warning indicators.** The panel includes directional retrieval@$k$, symmetry gaps, norm and spectrum comparisons, collapse indicators (mean cosine, variance floors, hubness), and structural probes (MI proxy, CCA-like proxy, sensitivity proxy). The panel is computed on frozen baseline models and post-intervention models.

**(3) A multi-signal drift classifier.** The classifier maps monitoring deltas to drift labels and allows an explicit "unknown" outcome when evidence is insufficient. It exports debug artifacts in that case rather than forcing false certainty.

**(4) A bounded intervention policy.** Interventions include temperature recalibration, loss reweighting based on gradient norms, renormalization recalibration, pairing sanitation filters, confounder ablation tests, and limited fine-tuning with stop rules. Each intervention is justified with decision records that reference measured facts.

**(5) A lifecycle report and audit bundle.** The notebook exports per-episode tables, stage-gate outcomes, decision records, and artifacts, packaged into an audit bundle with a manifest and logs. The report is structured for reviewability, not for narrative persuasion.

These deliverables are pedagogically aligned with the course's objective: students should leave with a clear understanding that multimodality is not only a modeling trick but a governance challenge. They should understand why a system can drift while seeming to improve, how to detect that drift, and how to respond without self-deception. They should also understand the limits: synthetic evidence is a training ground for thinking, not a substitute for production validation.

With the introduction complete, the next section formalizes the definitions and core concepts that the remainder of the chapter relies upon: latent factors and measurement operators, coupling contracts, drift vs shift, signatures vs mechanisms, and governance vocabulary such as budgets, rollbacks, and stage gates. This establishes the conceptual toolkit needed to interpret the synthetic laboratory's results and to translate them into professional practice.

## 3.2 Definitions and Core Concepts

### 3.2.1 Latent factors, modalities, and measurement operators

A multimodal system begins with a world that is not directly observed. It is observed through modalities. This is not a philosophical statement; it is a modeling statement that allows us to talk precisely about what alignment is supposed to mean and how it can fail. The central construct is a latent state, or latent factor vector, which we denote by $f$. In a synthetic laboratory, $f$ is explicit: it includes discrete factors such as shape type or orientation, and continuous factors such as frequency or phase. In production, $f$ is implicit: it is the underlying content and meaning that the system should preserve across modalities.

A modality is an observation channel. We denote the observations as $x^{\text{img}}$ and $x^{\text{txt}}$ for images and text. These observations are generated from the latent state by measurement operators:

$$x^{\text{img}} = \mathcal{M}_{\text{img}}(f, \epsilon_{\text{img}}), \qquad x^{\text{txt}} = \mathcal{M}_{\text{txt}}(f, \epsilon_{\text{txt}}),$$

where $\epsilon_{\mathrm{img}}$ and $\epsilon_{\mathrm{txt}}$ represent modality-specific noise and corruption. The measurement operators encode the fact that modalities do not observe the latent state directly. They distort it, compress it, and contaminate it with artifacts. This simple formalism is powerful because it makes drift natural: drift can be a change in the measurement operator, not only in the model.

For example, in production, an "image" might be compressed differently, resized differently, or passed through a new preprocessing pipeline. That is a change in $\mathcal{M}_{\mathrm{img}}$. A "text" might be generated by a different annotation pipeline, or tokenized differently, or include new metadata. That is a change in $\mathcal{M}_{\mathrm{txt}}$. These changes can occur without any change in the model. Yet they alter the distribution of observations and therefore the distribution of embeddings.

The measurement operator view also clarifies that multimodal alignment is about learning invariances. We do not want the embedding to encode idiosyncratic artifacts of the measurement operator; we want it to encode stable aspects of $f$. The most important word in that sentence is not "stable" but "aspects." In many real problems, no embedding can preserve all aspects of $f$; the goal is to preserve aspects that are relevant to downstream tasks. That makes alignment a contract: a set of intended invariances and sensitivities.

This contract must be made explicit. If the system is intended for image-caption retrieval, the contract is: embeddings should be close for pairs that share semantic content. If the system is intended for compliance review, the contract is: embeddings should preserve factors related to regulatory meaning and should not be dominated by superficial formatting. In both cases, the measurement operator view forces a disciplined question: which parts of the observation are signal and which are artifact? Drift is often the introduction of new artifacts or the removal of old ones, causing the representation to shift its reliance.

In a synthetic lab, we can instrument this precisely. We can control $f$, and we can control how $\mathcal{M}_{\mathrm{img}}$ and $\mathcal{M}_{\mathrm{txt}}$ introduce noise, corruption, or confounders. In production, we cannot fully specify $f$, but we can still adopt the operator perspective: pipeline changes are changes in $\mathcal{M}$, and these changes must be monitored as potential sources of drift.

### 3.2.2 Shared embedding spaces and coupling contracts

The central engineering construct in multimodal systems is a shared embedding space. We define encoders

$$z^{\mathrm{img}} = E_{\mathrm{img}}(x^{\mathrm{img}}; \theta), \qquad z^{\mathrm{txt}} = E_{\mathrm{txt}}(x^{\mathrm{txt}}; \phi),$$

where $z^{\mathrm{img}}, z^{\mathrm{txt}} \in \mathbb{R}^d$ are embeddings. Typically, we normalize them:

$$\hat{z} = \frac{z}{\|z\| + \varepsilon},$$

and define similarity by cosine:

$$s(i, j) = \hat{z}_i^{\mathrm{img}} \cdot \hat{z}_j^{\mathrm{txt}}.$$

A contrastive objective such as InfoNCE encourages matched pairs $(i, i)$ to have high similarity and mismatched pairs $(i, j)$ to have lower similarity.

This setup is often described as "learning a shared space." But it is more precise to say it learns a *coupling contract.* A coupling is a relationship between two distributions. Here, it is the relationship between the distribution of images and the distribution of texts. The contract says: there exists a mapping such that paired observations occupy nearby points in the space. This contract is relational, not absolute. It does not assert that embeddings represent truth; it asserts that embeddings preserve a particular correspondence relation.

This matters because it clarifies what drift breaks. Drift can break the contract in multiple ways:

- If the measurement operator changes, the mapping from $f$ to observations changes, making the learned encoders miscalibrated.
- If pairing integrity changes (join bugs), the coupling relation itself changes: what counts as a pair is different.
- If scoring changes (temperature scaling), the operational meaning of similarity changes even if embeddings do not.
- If confounders appear, the easiest way to satisfy the contract may shift from semantic coupling to shortcut coupling.

In production, the coupling contract must be stated in measurable terms. This is why we track directional retrieval, symmetry, and structural probes. The contract is not only "pairs are close." It is also "both directions behave similarly," "the space remains diverse," and "the space encodes intended factors rather than shortcuts." These are contractual obligations.

A coupling contract also implies that model updates are contract renegotiations. Fine-tuning changes $\theta$ and $\phi$, therefore it changes the coupling relation. That is not necessarily bad. But it must be governed. A system that fine-tunes without monitoring is renegotiating its contract blindly.

In the synthetic lab, we model this explicitly. We fix a baseline contract at $t_0$. Then we simulate drift episodes that alter the measurement operators or the pairing relation. We monitor whether the contract still holds. If it does not, we apply bounded interventions to restore contract health. This turns an abstract concept into an operational workflow.

### 3.2.3 Drift vs shift: benign change and failure change

The term "drift" is often used loosely to mean any distribution change. That is not sufficient for governance. We need a distinction between benign shift and failure drift. Benign shift means that the input distribution changes but the representation contract remains healthy. Failure drift means that the contract is violated in ways that matter operationally.

Formally, suppose the system at time $t$ produces an embedding distribution $P_t(z^{\text{img}}, z^{\text{txt}})$ and a

similarity distribution over pairs. A shift is any change $P_t \neq P_{t_0}$. A failure drift is a shift that causes one or more contract constraints to be violated: retrieval falls below thresholds; symmetry gaps exceed bounds; effective rank falls; hubness rises; confounder MI exceeds ceilings; sensitivity spikes.

This definition is intentionally operational. It avoids the trap of treating any shift as a failure. In real systems, shift is constant: new topics, new styles, new devices. If governance treats every shift as a problem, it becomes unusable. Governance must instead define safe operating regions. A system can move within that region without triggering interventions. When it leaves the region, governance triggers an action.

The safe region is defined by multi-metric bounds, not a single number. This is because failures are multi-dimensional. A system can maintain retrieval while losing effective rank; it can maintain symmetry while encoding confounders; it can maintain factor separability while becoming sensitive to updates. Therefore, the safe region is a set of constraints.

In the notebook, we explicitly define these constraints and demonstrate how a system can drift into failure modes. Some drift episodes will be benign: increased noise that degrades performance smoothly but does not induce collapse. Others will be failure drift: introduction of a confounder that improves retrieval but violates factor MI constraints. The point is to teach the student that drift detection must be aligned with contractual obligations, not with abstract distributional difference metrics alone.

### 3.2.4 Signatures vs mechanisms: attribution under uncertainty

In Chapter 2, we separated signatures (observable patterns) from mechanisms (causal explanations). Chapter 3 continues this separation because lifecycle governance must operate under uncertainty. In production, we rarely have direct access to mechanisms. We observe symptoms: retrieval gaps, hubness, spectra, probe changes. We must infer plausible mechanisms. That inference can be wrong. Therefore, the system must represent uncertainty explicitly.

A signature is a measurable pattern: for example, directional retrieval gap increases, gradient norms become imbalanced, and one modality's effective rank contracts. This signature suggests dominance. Another signature: retrieval remains high, factor MI drops, confounder MI rises, and ablation collapses performance. This signature suggests spurious alignment. Another signature: loss becomes unstable, hubness rises, sensitivity spikes, and corruption audits fail. This suggests pairing corruption. Another signature: mean cosine rises, variance floors drop, effective rank falls, and the spectrum becomes concentrated. This suggests collapse.

Mechanisms are causal stories: e.g., a new watermark introduced into both modalities created a shortcut; a pipeline join bug introduced systematic pairing offset; an upstream normalization change increased one modality's scale and induced dominance; a downstream scoring scale change effectively lowered temperature and induced anisotropy. Mechanisms are what we want to know to choose

interventions. But we must acknowledge that signature-to-mechanism mapping is imperfect.

Therefore, this chapter introduces an explicit drift classification scheme that includes an "unknown" outcome. The classifier is multi-signal: it uses multiple indicators and requires minimum evidence. If evidence is insufficient, it outputs "unknown" and exports debug artifacts for human review. This is not a weakness. It is governance integrity. Forcing a label when evidence is weak creates false confidence and can trigger the wrong intervention. In professional settings, false confidence is more dangerous than uncertainty because it leads to uncontrolled actions.

The signature-versus-mechanism distinction also clarifies how interventions are justified. We do not justify interventions by narrative. We justify them by measurable facts and by explicit mapping rules. For example, we apply temperature recalibration when scoring drift is detected by similarity scale changes and mean cosine/hubness patterns. We apply loss reweighting when gradient imbalance and symmetry gaps suggest dominance. We apply confounder ablation tests when MI comparisons suggest spurious alignment. We apply pairing sanitation when corruption signatures appear. Each intervention includes open items: what else could explain the signature? What additional tests would reduce uncertainty?

### 3.2.5   Governance vocabulary: stage gates, budgets, rollbacks

Finally, we introduce the governance vocabulary that turns monitoring into action. Three terms are central: stage gates, drift budgets, and rollbacks.

A **stage gate** is an acceptance decision rule. It is defined by a set of constraints that must hold. Stage gates are applied at baseline training (to accept the initial model), after drift detection (to decide whether intervention is needed), and after interventions (to decide whether the intervention succeeded). Stage gates are multi-metric because the contract is multi-dimensional.

A **drift budget** is a tolerance envelope for changes. Rather than reacting to any deviation, we define acceptable ranges of metric changes. For example, retrieval@1 may be allowed to drop slightly under mild noise increases, but effective rank may not be allowed to fall below a floor. Symmetry gaps may be allowed to fluctuate within bounds. Drift budgets prevent governance from being overly reactive, and they create clear triggers for action.

A **rollback** is the disciplined response when an intervention fails. In uncontrolled systems, teams often patch models repeatedly, creating a drift spiral. In a governed system, interventions are bounded, evaluated against stage gates, and rolled back if they fail. Rollback discipline is essential because interventions can introduce new failure modes. For example, fine-tuning to fix retrieval may induce collapse. Reweighting to fix dominance may amplify confounder reliance. Therefore, interventions must be reversible, and the system must maintain a known-good checkpoint.

These governance terms are not administrative. They are engineering primitives for operating a representational policy under change. The synthetic lab implements them explicitly: each episode

produces measurements; classification triggers an intervention plan; interventions are executed under constraints; stage gates decide success; failures trigger rollback; and artifacts record the entire sequence.

With these definitions established, the chapter is prepared to specify the baseline aligner and the reference geometry. The next section defines the synthetic world, the two-encoder architecture, the symmetric InfoNCE objective, the baseline health gates, and the reference artifacts that serve as the anchor for lifecycle monitoring. This anchors the lifecycle discussion in a concrete baseline rather than in abstraction.

## 3.3 Baseline Aligner and Reference Geometry

### 3.3.1 Synthetic multimodal world and factor controls

A lifecycle framework needs an anchor: a reference system whose behavior at $t_0$ is understood and instrumented. In production, that anchor might be a released model version plus its data and pipeline versions. In this chapter, the anchor is a synthetic multimodal world and a baseline aligner trained on it. The purpose of the synthetic world is not to mimic reality; it is to provide causal controllability. We need to be able to say, without ambiguity, what the intended latent factors are, what constitutes semantic coupling, and how drift episodes perturb the data-generation and pairing pipeline.

We define latent factors $f$ as a structured vector with both discrete and continuous components. A typical configuration includes: shape type (discrete), orientation (discrete or continuous), spatial frequency (continuous), phase (continuous), and thickness or amplitude (continuous). These factors are chosen for pedagogical clarity: they can be rendered into small images and encoded into token sequences without requiring large datasets or pretrained models. They also create multiple axes of variation so that the embedding space must represent more than a binary distinction. This matters because many failure modes—collapse and hubness in particular—are only visible when the system must represent multiple degrees of freedom.

The image modality is generated by a rendering operator $\mathcal{M}_{\text{img}}$ that maps factors into a small matrix, such as $16 \times 16$ or $24 \times 24$. Rendering might involve generating simple patterns: a sinusoidal grating with orientation and frequency, a blob with thickness, or a composite of primitives. The exact choice is less important than two structural properties: (i) the mapping is smooth in the continuous factors, and (ii) the mapping is nontrivial across discrete factors. Smoothness ensures that neighborhood relations in factor space can correspond to neighborhood relations in embedding space. Nontriviality ensures that the model cannot solve the task with a trivial binary partition.

The text modality is generated by a symbolic operator $\mathcal{M}_{\text{txt}}$ that maps the same factors into token sequences under a controlled grammar. Tokens represent factor values in discretized bins (e.g.,

orientation bucket, frequency bucket), plus separators and optional positional features. We avoid using any language model. This is critical for pedagogy: the goal is to teach multimodality as coordinate-system compatibility, not to hide complexity inside a pretrained text encoder. The text features are then formed as vector encodings (bag-of-tokens plus positional statistics, or a simple learned token embedding summed across positions). The key is that the text modality is structured and factor-linked, but it remains synthetic and controllable.

The world generator also supports optional channels used for drift experiments: confounder features (shared shortcut signals), modality-specific noise, token corruption, blur-like operations, and pairing corruption. These are not added to the baseline. They are introduced in drift episodes. However, the baseline generator must be designed so that these perturbations can be applied cleanly and reproducibly. That is why the generator is written as a parameterized operator, not as a fixed dataset.

Finally, the world includes explicit split hygiene. We generate a dataset of size $N$, then split indices into train/validation/test deterministically. The factors are generated in a way that avoids leakage: for example, we ensure that the same latent factor combinations do not appear identically across splits unless we intentionally design such a test. In a synthetic setting, leakage can be subtle: if the factor generation uses a fixed random seed per split but with overlapping ranges, identical factor vectors can appear across splits. Therefore, the baseline includes explicit checks: count duplicates of factor vectors across splits and assert that the overlap is below a tiny threshold. This is not merely good practice; it teaches students that leakage exists even in synthetic labs and must be treated explicitly.

### 3.3.2 Two-encoder architecture and symmetry requirements

The baseline aligner consists of two encoders: an image encoder $E_{\text{img}}$ and a text encoder $E_{\text{txt}}$. Each is a two-layer multilayer perceptron (MLP) implemented in NumPy. We choose a 2-layer MLP for two reasons. First, it is the simplest architecture that can exhibit nontrivial representation learning while still being tractable for analytic gradients and gradient checking. Second, it is sufficient to demonstrate the geometric phenomena of interest: alignment, dominance, confounding, and collapse.

Each encoder maps modality-specific features into an embedding vector in $\mathbb{R}^d$. Let $h = \tanh(W_1 x + b_1)$ and $z = W_2 h + b_2$, then normalize $\hat{z} = z/(\|z\| + \varepsilon)$. The tanh nonlinearity is chosen because it is smooth and has bounded outputs, which helps numerical stability. It also allows a clean analytic derivative.

Symmetry is not merely a design aesthetic; it is a contractual requirement. A shared embedding space is used bidirectionally: image-to-text retrieval and text-to-image retrieval. If the system works only in one direction, it is not aligned. Therefore, the baseline enforces symmetry at multiple levels:

- **Objective symmetry:** use a symmetric InfoNCE loss that includes both directions.

- **Architectural symmetry:** both encoders have comparable capacity (similar hidden width), unless we intentionally break this in dominance experiments.
- **Monitoring symmetry:** track directional retrieval metrics and symmetry gaps; track norm and spectrum comparisons across modalities.
- **Optimization symmetry:** ensure learning rates and regularization are applied consistently unless intentionally perturbed.

This symmetry posture is essential because many real multimodal failures are directional. For example, one modality may become the anchor and the other modality may map into it. In such cases, one direction might appear strong while the other direction degrades. Without symmetry monitoring, the system can pass a headline metric while being operationally asymmetric.

The baseline also includes explicit shape and NaN assertions throughout the forward and training pipeline. These are not optional. In governed training, silent NaNs and silent shape mismatches produce irreproducible results and false narratives. The baseline treats these as hard failures.

### 3.3.3 Symmetric InfoNCE objective and stable numerics

The objective used in the baseline is a symmetric InfoNCE loss computed over a batch. Given normalized embeddings $\hat{Z}_{\text{img}} \in \mathbb{R}^{B \times d}$ and $\hat{Z}_{\text{txt}} \in \mathbb{R}^{B \times d}$, we define logits:

$$L = \frac{1}{\tau} \hat{Z}_{\text{img}} \hat{Z}_{\text{txt}}^{\top},$$

where $\tau > 0$ is a temperature parameter. For image→text, the loss is the average cross-entropy over rows, treating the diagonal as the correct class:

$$\ell_{i \to t} = \frac{1}{B} \sum_{i=1}^{B} \left( -L_{ii} + \log \sum_{j=1}^{B} e^{L_{ij}} \right).$$

For text→image, we apply the same formula to columns (equivalently, to $L^{\top}$):

$$\ell_{t \to i} = \frac{1}{B} \sum_{j=1}^{B} \left( -L_{jj} + \log \sum_{i=1}^{B} e^{L_{ij}} \right).$$

The symmetric loss is $\ell = \frac{1}{2}(\ell_{i \to t} + \ell_{t \to i})$.

The key engineering challenge is stable computation. Even in small models, logits can have range that causes overflow in exp. Therefore, we implement log-sum-exp stabilization:

$$\log \sum_j e^{a_j} = m + \log \sum_j e^{a_j - m}, \qquad m = \max_j a_j.$$

We also implement explicit eps guards in normalization. These stability choices are not implementa-

tion details; they affect geometry. For example, changing effective temperature changes sharpness, which is directly linked to collapse risk. Therefore, the baseline treats temperature as a controlled parameter and records it in the manifest.

We also emphasize that InfoNCE is a ranking objective. It does not explicitly enforce variance or decorrelation. Therefore, collapse is possible under certain regimes. This is why baseline training includes collapse monitors even before we induce drift. It teaches students that collapse is not a hypothetical; it is a property of the objective. In practice, the baseline is tuned to avoid collapse, but the monitors remain active.

### 3.3.4   Baseline acceptance: health gates and reproducibility

Before we can talk about drift, we must accept a baseline. Baseline acceptance is not "loss decreased." It is a stage gate decision defined by a multi-metric contract. The baseline must satisfy:

- **Retrieval:** directional retrieval@1 and retrieval@5 exceed minimum thresholds on validation and test.
- **Symmetry:** directional gaps are below a threshold; distributions of norms (pre-normalization) are comparable; effective rank difference is bounded.
- **Representation health:** mean off-diagonal cosine is below a ceiling; variance floors are above a floor; effective rank is above a floor; hubness is below a ceiling.
- **Semantic encoding:** factor separability ratios exceed floors; MI proxy with factors is above floors; MI with confounders is not applicable in the baseline but the framework is prepared.
- **Stability:** no NaNs; gradients are finite; training metrics are reproducible across runs given fixed seeds.

These gates are selected for pedagogical and professional reasons. Pedagogically, they force students to see alignment as multi-dimensional. Professionally, they mirror what real systems require: you cannot accept a model that is strong on retrieval but unhealthy geometrically.

Reproducibility is enforced with deterministic seeds, fixed shuffles, and explicit version logging. The baseline training run writes a run manifest containing configuration, seeds, library versions, and environment fingerprints. It writes metric logs and plots. It saves checkpoints. It packages all artifacts into an audit bundle. This is not an embellishment; it is the foundation of lifecycle governance. Without a baseline artifact set, drift cannot be interpreted because there is no reference state.

### 3.3.5   Reference artifacts: manifests, logs, and checkpoints

The baseline produces reference artifacts that serve as anchors for the remainder of the chapter. These include:

**Run manifest.** A JSON document capturing configuration parameters (data generator settings, model dimensions, training hyperparameters, temperature, seeds), environment metadata (Python version, NumPy version), and file paths to artifacts. The manifest is the primary provenance record.

**Metric logs.** Time-series logs of loss, retrieval metrics, symmetry metrics, collapse indicators, and gradient norms. These logs allow us to compare later runs and to detect deviations.

**Plots.** Visual evidence: PCA projections colored by factors, similarity heatmaps, spectrum plots, trend plots. Visuals are not proofs, but they are diagnostic evidence when coupled with metrics.

**Checkpoints.** Saved encoder parameters with metadata. Checkpoints enable rollback and enable frozen-baseline monitoring. They also allow controlled fine-tuning interventions without losing the original state.

**Reports.** Strict JSON summaries that separate facts, assumptions, open items, analysis, and draft outputs. These reports enforce epistemic discipline and provide a format suitable for review.

The key idea is that the baseline is not merely a model. It is a documented state of the system. Drift is defined relative to this state. Monitoring compares current behavior to baseline behavior. Classification uses deltas. Interventions aim to restore baseline contract properties. Rollback returns to baseline checkpoints.

With the baseline and reference geometry established, we can define the drift taxonomy that the lifecycle loop must handle. The next section specifies drift types that correspond to realistic production changes: modality drift, scoring drift, confounder drift, pairing drift, and optimization drift. Each drift type has characteristic signatures and suggested interventions, which we will later formalize into a classifier and a bounded intervention policy.

## 3.4   Drift Taxonomy for Multimodal Systems

### 3.4.1   Modality drift: noise, corruption, and sensor degradation

Modality drift refers to changes in one observation channel that alter the distribution of inputs without necessarily changing the latent state distribution. In the measurement-operator formalism, modality drift is a change in $\mathcal{M}_{\text{img}}$ or $\mathcal{M}_{\text{txt}}$ that increases noise, changes corruption patterns, or changes the effective resolution of the signal. In production, modality drift is common: camera sensors change, compression schemes update, OCR quality varies, annotation practices evolve, and language style shifts. The key feature is asymmetry: one modality may drift while the other remains stable.

This asymmetry is precisely why modality drift is dangerous. A symmetric multimodal model assumes that both modalities provide compatible coordinate systems for the same latent factors. If one modality degrades, the shared embedding space can become imbalanced. The model may begin

to privilege the cleaner modality, creating dominance. Directional retrieval becomes asymmetric: text→image may remain acceptable while image→text degrades (or vice versa). Alternatively, the model may attempt to maintain alignment by compressing both modalities' variability, increasing hubness and collapse risk. In other words, modality drift can be the upstream cause of multiple downstream failure modes.

In the synthetic drift lab, modality drift is implemented by increasing noise in images (additive noise, blur, occlusion-like corruption) or by corrupting text tokens (drop tokens, substitute tokens, perturb positional features). The magnitude of drift is parameterized and applied over episodes. Importantly, we also include mild stochastic variation so that drift is not perfectly deterministic. This mimics reality: drift is rarely a clean step function; it fluctuates.

The monitoring signatures for modality drift include: directional retrieval degradation concentrated in one direction, increased symmetry gaps, and changes in norm distributions or effective rank that indicate one modality is losing information. Structural probes show a particular pattern: MI proxy with factors declines more for the drifting modality; CCA-like correlations may remain moderate if the stable modality anchors the shared subspace; sensitivity may increase if the system approaches an unstable regime.

The governance posture is to treat modality drift as a budgeted phenomenon. Some degradation is acceptable. The key is to detect when drift pushes the system beyond acceptable asymmetry or beyond health constraints. Interventions may include renormalization recalibration (correcting feature scale), loss reweighting (to prevent dominance), or limited fine-tuning with strict stop rules. But interventions are bounded: we do not chase performance if doing so induces collapse indicators.

### 3.4.2   Scoring drift: temperature, normalization, and calibration shifts

Scoring drift refers to changes in how similarity is computed or used, rather than changes in encoder weights or inputs. This is an underappreciated source of lifecycle failure. In many systems, embeddings are produced by encoders, then passed through a separate service that computes similarity and performs retrieval. That service can change. A downstream team might normalize differently, might rescale scores, might change the temperature used in softmax-based ranking, or might change the indexing algorithm. These changes can alter retrieval behavior even if embeddings do not change.

In our formalism, scoring drift can be represented as a change in the function that maps embeddings to scores:

$$s(i, j) = g(\hat{z}_i^{\text{img}}, \hat{z}_j^{\text{txt}}; \alpha),$$

where $\alpha$ includes scale and temperature-like parameters. For cosine similarity, $g$ is a dot product. But in practice, there may be learned scalers, quantization effects, or approximate nearest neighbor structures that introduce distortions.

A central scoring drift phenomenon is effective temperature change. If scores are multiplied by a factor $c$, then the softmax distribution in contrastive training or downstream ranking becomes sharper as $c$ increases. This can mimic lowering $\tau$. Sharpness interacts with hubness and collapse: a sharper scoring regime makes the system more sensitive to small differences, can cause retrieval to become brittle, and can amplify anisotropy. The system may appear to improve on certain metrics because top-1 choices become more decisive, but it can become less robust under drift and less equitable across modalities.

In the synthetic lab, scoring drift is implemented by changing the scale applied to logits during evaluation (and optionally during fine-tuning) to simulate a downstream calibration change. We then observe how mean cosine, hubness, and effective rank interact with retrieval. The key lesson is that scoring drift can create failure signatures without any model training. Monitoring must therefore include scoring regime checks.

Interventions for scoring drift are often simpler than for other drifts: recalibrate temperature, restore consistent normalization, or adjust scoring thresholds. But governance demands that these changes are logged and evaluated under stage gates. A scoring change is a policy change. It must be reviewable.

### 3.4.3 Confounder drift: shortcut appearance and disappearance

Confounder drift is one of the most damaging lifecycle phenomena because it can cause both silent improvement and silent failure. A confounder is a shared feature present in both modalities that is correlated with pairing but not with intended semantics. In practice, confounders arise from metadata, formatting, watermarks, IDs, source markers, or preprocessing artifacts. They can appear when a new data source is added, when a watermarking pipeline is introduced, or when a metadata join injects an ID into both channels. They can disappear when pipelines are cleaned, when privacy controls remove IDs, or when formatting standards change.

Confounder drift can cause headline metrics to improve dramatically. If the confounder is a strong shared key, the model can align pairs by encoding the confounder rather than semantics. Retrieval improves because the confounder is consistent across splits. This is spurious alignment. The model is not learning the intended latent factors. It is learning a shortcut coupling.

Confounder disappearance is the mirror danger. If the model has relied on the confounder, removing it can cause sudden performance collapse. A system that looked strong in offline evaluation can fail in production when the confounder distribution shifts. This is why confounder drift must be treated as a first-class lifecycle risk.

In the synthetic lab, confounder drift is implemented by introducing a shared codebook feature into both modalities for certain episodes, then removing it. We measure how MI proxy shifts: MI with confounder rises, MI with true factors may stagnate or decline, and retrieval may rise. When

the confounder is removed, retrieval may drop sharply, and factor MI may not recover because the representation never learned factors. This is an explicit demonstration of why "metrics improved" is not sufficient evidence of semantic alignment.

Controls against confounder drift include explicit shortcut audits (MI comparisons, ablation tests), split designs that break confounders, and data governance policies that prevent shared IDs from leaking into both modalities. Interventions include confounder ablation training, filtering out confounded examples, and re-training or fine-tuning on confounder-free data with strict monitoring. In production, the most important control is pipeline governance: prevent confounders from being introduced silently.

### 3.4.4 Pairing drift: join bugs, offsets, and systematic misalignment

Pairing drift refers to changes in the coupling relation: which image is paired with which text. In production, pairing is often constructed through joins: join on ID, join on timestamp, join on content hash. These joins can fail. A key can change. An offset can be introduced. Records can be re-ordered. Deduplication can behave differently. The result is partial misalignment: some fraction of pairs are incorrect.

Pairing drift is especially dangerous because it creates contradictory supervision. If a system is fine-tuned online using new pairs, pairing corruption can push the model into inconsistent regimes. Even without fine-tuning, pairing drift affects evaluation: if test pairs are misaligned, retrieval metrics become noisy and can mislead. Pairing drift can therefore corrupt both the model and the monitoring.

In the synthetic lab, pairing drift is implemented as episodic corruption: random swaps, systematic offset swaps, and block-level shuffles. We apply corruption at different rates and patterns, and we measure signatures from Chapter 2: loss instability, rising sensitivity, hubness, and spectrum distortions. We also test how a frozen baseline behaves under corrupted evaluation, teaching that monitoring must include pairing integrity audits.

Controls include join audits, data validation checks, similarity-based sanity checks, and sanitation filters that drop suspicious pairs during fine-tuning. Interventions are delicate. If pairing drift is detected, the correct first action is often not to retrain but to fix the pipeline. Fine-tuning on corrupted pairs can make the model worse. Therefore, governance rules should include stop conditions: do not fine-tune when pairing integrity is suspect.

### 3.4.5 Optimization drift: fine-tuning regimes and collapse boundaries

Optimization drift refers to changes in the training process after deployment: incremental fine-tuning, adaptation to new data, or changes in hyperparameters that alter the system's geometry. Many systems are updated continually. Fine-tuning is attractive because it can improve performance on

new distributions. But it also introduces collapse risk. A model that is stable under one regime can become unstable under a different fine-tuning regime: lower temperature, higher learning rate, different batch composition, stronger regularization, or different negative sampling.

Optimization drift is not merely "training more." It is changing the dynamical system that updates the parameters. A system can cross a stability boundary and enter a collapse regime. Once collapse begins, it can be self-reinforcing: hubness rises, gradients become biased, and the model converges to a degenerate geometry. This can happen even while loss decreases. Therefore, fine-tuning must be governed with health gates and stop rules.

In the synthetic lab, optimization drift is implemented by episodes that attempt limited fine-tuning under risky settings (e.g., lower temperature and higher LR). We then show how collapse indicators respond. We also show how bounded interventions can prevent collapse: clip gradients, enforce variance floors, abort fine-tuning when effective rank drops, and roll back to the last healthy checkpoint.

This drift type connects the chapter to lifecycle governance most directly. It demonstrates that governance is not only about detecting drift in inputs; it is also about controlling drift induced by our own actions. In production, many catastrophic failures are self-inflicted: aggressive fine-tuning, unreviewed hyperparameter changes, or unmonitored pipeline updates. Optimization drift is the category that covers those self-induced changes, and it motivates the strict intervention discipline of the chapter.

With this taxonomy established, we have a vocabulary for the kinds of lifecycle changes we want to detect. The next section specifies the monitoring panel and early warning indicators that will be computed per episode. These indicators are designed to detect not only performance degradation but also geometric and mechanistic shifts that precede visible failures. They will form the input to the drift classifier and the basis for bounded interventions.

## 3.5   Monitoring Panel and Early Warning Indicators

### 3.5.1   Directional retrieval and symmetry gap metrics

In lifecycle governance, the first monitoring obligation is to measure whether the coupling contract remains operationally valid. Directional retrieval metrics provide the most direct operational signal: when the model is used to retrieve text given an image (image→text) or retrieve images given text (text→image), does it still retrieve correct matches reliably? However, Chapter 2 already showed that retrieval alone is insufficient. In Chapter 3, retrieval is not removed; it is placed in context and augmented with symmetry analysis.

We define retrieval@$k$ as follows. Given a batch or evaluation set of size $N$, compute the similarity

matrix

$$S = \hat{Z}_{\text{img}} \hat{Z}_{\text{txt}}^{\top},$$

where embeddings are normalized and $\hat{Z}_{\text{img}}, \hat{Z}_{\text{txt}} \in \mathbb{R}^{N \times d}$. For each image $i$, rank the text indices $j$ by $S_{ij}$ descending. Retrieval@$k$ for image→text is the fraction of times the correct paired text index appears in the top $k$. For text→image, use the transpose and rank by column.

These metrics have two important properties in a lifecycle setting. First, they are sensitive to scoring drift and indexing drift. Second, they can mask asymmetry. It is common to report an average or a combined score, but that can hide a directional collapse. Therefore, we treat the directional metrics as separate contractual obligations.

The symmetry gap is defined as a difference between directional retrieval metrics:

$$\Delta_k = \text{R@}k(i \to t) - \text{R@}k(t \to i).$$

We monitor $\Delta_1$ and $\Delta_5$. In a healthy symmetric system, these gaps should be small and stable. A sustained increase in gap is an early warning of dominance or modality drift. We do not require the gap to be zero; asymmetries in modalities can create natural differences. But we define budgets: gaps beyond a threshold trigger investigation and potentially intervention.

We also include distributional symmetry checks beyond retrieval. One common trap is to normalize embeddings and then believe the distributions are inherently comparable. But dominance often manifests before normalization: one modality's pre-normalization norms become systematically larger, indicating that one encoder's output scale is different. This matters because pre-normalization scale affects gradients during training and can lead to dominance regimes. Therefore, we track norm statistics (mean, variance, quantiles) before normalization and compare them across modalities.

Finally, we track effective-rank divergence across modalities. Even if both modalities produce $d$-dimensional embeddings, one modality might occupy a lower-dimensional subspace. That is a form of dominance or collapse. Therefore, symmetry monitoring includes geometry, not only performance.

Directional retrieval plus symmetry gaps is the minimal "operational policy" monitor. If these metrics fail, the system is failing its basic contract. But as Chapter 2 emphasized, the system can satisfy these metrics while failing semantically or geometrically. Therefore, we add collapse indicators, spectral monitors, and probes.

### 3.5.2   Collapse indicators: mean cosine, variance floors, hubness

Collapse indicators are early warning signals that the embedding space is losing diversity and becoming dominated by a small set of directions or hubs. The key reason collapse indicators are essential in lifecycle governance is that collapse can emerge without obvious performance degradation. A model can retain retrieval on familiar data while becoming increasingly hub-dominated and brittle,

meaning it will fail under drift.

We monitor three primary collapse indicators:

**Mean off-diagonal cosine similarity.** For a set of normalized embeddings $\{\hat{z}_i\}_{i=1}^N$, compute the cosine similarity matrix $C = \hat{Z}\hat{Z}^\top$. Excluding the diagonal, compute the mean:

$$\mu_{\cos} = \frac{1}{N(N-1)} \sum_{i \neq j} C_{ij}.$$

In an isotropic representation on the sphere, $\mu_{\cos}$ should be near zero. In an anisotropic or collapsed representation, embeddings align and $\mu_{\cos}$ increases. We compute $\mu_{\cos}$ for each modality separately and track it across episodes.

**Variance floors.** Compute the covariance of centered embeddings:

$$\Sigma = \frac{1}{N-1}(Z - \bar{Z})^\top (Z - \bar{Z}),$$

and track summaries of its diagonal (coordinate-wise variances). A healthy space has coordinate variances that are not all near zero. A collapse regime produces small variance in most directions. We therefore track the minimum variance, median variance, and a low quantile. We treat a low quantile dropping below a floor as a stop condition in fine-tuning interventions.

**Hubness.** Hubness captures whether a small set of embeddings become nearest neighbors for many queries. It is computed by counting top-1 retrieval targets: for each query, find the index of the maximum similarity. Then count how many times each target is chosen. A hub-dominated space will have a highly skewed count distribution. We summarize hubness by the fraction of queries mapped to the top $m$ targets (e.g., $m = 1$ or $m = 5$), or by a Gini coefficient over the counts. Hubness matters because it reflects an operational failure mode: retrieval becomes repetitive, and diversity of results collapses even if top-1 accuracy on a narrow set remains.

These indicators are computed per episode on a fixed evaluation set. We interpret them as part of a panel: rising mean cosine, falling variance floors, and rising hubness is a strong collapse signature. If only one indicator moves, it may be benign anisotropy. Therefore, we treat collapse detection as multi-signal.

Importantly, we compute these indicators for both modalities and for the cross-modal similarity matrix. Cross-modal hubness can indicate spurious alignment: a shared confounder can create hubs that dominate cross-modal retrieval. Therefore, hubness is useful beyond collapse; it is a general geometry health monitor.

### 3.5.3 Spectral monitoring: covariance eigenvalues and effective rank

Variance floors and mean cosine are scalar summaries. Spectral monitoring provides a richer view: it tells us how many independent directions the representation uses. This is critical for drift detection because many failure modes manifest as spectrum distortions before performance fails.

Given centered embeddings $Z \in \mathbb{R}^{N \times d}$, we compute covariance $\Sigma$ and its eigenvalues $\lambda_1 \geq \cdots \geq \lambda_d \geq 0$. The eigenvalue spectrum indicates how variance is distributed across directions. A healthy space distributes variance across many dimensions. A collapsed space concentrates variance in a few dimensions. A confounder-dominated space can also concentrate variance, often in a small subspace associated with the shortcut.

We summarize the spectrum with **effective rank**:

$$r_{\text{eff}} = \exp\Big(-\sum_{i=1}^{d} p_i \log p_i\Big), \qquad p_i = \frac{\lambda_i}{\sum_{j=1}^{d} \lambda_j}.$$

Effective rank measures the exponential of the entropy of normalized eigenvalues. If all eigenvalues are equal, $r_{\text{eff}} \approx d$. If only $k$ eigenvalues carry mass, $r_{\text{eff}} \approx k$. In practice, effective rank provides a stable scalar that is sensitive to collapse and to dominance.

We compute effective rank for each modality. We also compute *rank divergence* between modalities:

$$\Delta r_{\text{eff}} = r_{\text{eff}}^{\text{img}} - r_{\text{eff}}^{\text{txt}}.$$

A sustained divergence indicates modality imbalance: one modality is using fewer directions. This can happen under modality drift (one modality loses information) or under dominance (one modality maps into the other's anchors).

Spectral monitoring also supports drift classification. For example, scoring drift (temperature change) can increase apparent sharpness and can change hubness without necessarily changing within-modality spectra. Confounder drift often increases cross-modal coupling and can manifest as a dominant principal component that correlates with confounder. Pairing drift can create noisy spectra and reduce cross-modal correlations. Optimization drift can induce rapid spectral contraction.

In the notebook, spectra are tracked over episodes and plotted. We also export per-episode eigenvalue summaries. However, we emphasize that spectra are not interpreted in isolation. Spectral changes can be benign if the domain becomes more uniform. Therefore, we pair spectral monitoring with probes that connect embeddings to factors and with operational metrics.

### 3.5.4 Structural probes: MI proxy, CCA proxy, sensitivity proxy

Structural probes connect the monitoring panel to mechanism attribution. They are the bridge between "something changed" and "what likely changed." Chapter 2 introduced these probes as tools for diagnosing induced failure modes. In Chapter 3, they become lifecycle instruments.

**MI proxy.** We compute mutual-information proxies between embeddings and true factors in the synthetic world. We also compute MI with known confounders when present. The key lifecycle use is to detect semantic drift. If MI with factors declines, the embedding is losing factor information. If MI with confounder rises relative to factor MI, the system is becoming shortcut-dominated. This can happen even when retrieval improves. Therefore, MI proxies are essential for detecting silent spurious alignment.

**CCA-like proxy.** We compute principal correlations between modality embeddings via whitening and SVD. The lifecycle use is to detect coupling changes. If cross-modal correlations drop sharply, pairing integrity or coupling quality may be failing. If cross-modal correlations become concentrated in a small number of dimensions, confounder coupling may be dominating. If correlations remain high while one modality's effective rank drops, dominance may be emerging. Thus, the CCA proxy is a coupling health monitor.

**Sensitivity proxy.** We measure how a small one-step training perturbation changes retrieval. The lifecycle use is to detect training fragility. If sensitivity spikes, the system may be near a collapse boundary or may be in a contradictory supervision regime due to pairing drift. Sensitivity also supports intervention decisions: if sensitivity is high, fine-tuning interventions are risky and should be bounded or avoided. Conversely, if sensitivity is low, limited fine-tuning may be safe.

These probes are proxies. Therefore, governance requires explicit caveats. We treat probe results as evidence, not as certification. We also require probe stability checks: for example, MI discretization should be tested with multiple bin counts; CCA whitening should be regularized and sensitivity to regularization should be recorded. The notebook is designed to export these stability notes as open items.

### 3.5.5 Operational dashboards: trends, thresholds, and alerts

The monitoring panel becomes operational only when it is turned into dashboards with thresholds and alert logic. This is where governance becomes practical. A dashboard is not merely a plot; it is a set of rules that triggers actions.

We define three layers of monitoring:

**Level 1: Performance alerts.** Trigger when retrieval@1 or retrieval@5 fall below thresholds or when symmetry gaps exceed budgets. These alerts indicate operational degradation.

**Level 2: Health alerts.** Trigger when mean cosine exceeds a ceiling, variance floors drop below a

floor, effective rank drops below a floor, or hubness exceeds a ceiling. These alerts indicate geometric degradation that may precede performance failures.

**Level 3: Mechanism alerts.** Trigger when MI with factors declines sharply, MI with confounders exceeds ceilings, CCA correlations change in signature ways, or sensitivity spikes. These alerts indicate plausible failure mechanisms: spurious alignment, pairing drift, dominance, or collapse risk.

Alerts are evaluated per episode, and the notebook records which alerts triggered. Alerts do not automatically imply an intervention. They feed into the drift classifier, which combines evidence and can output "unknown." The classifier then triggers an intervention plan if the evidence meets minimum requirements.

The chapter emphasizes that thresholds must be pre-specified for reviewability. A common failure in practice is to choose thresholds after seeing results. That is a form of metric hacking. In governed systems, thresholds are part of the contract and must be declared. The notebook records thresholds in the manifest, and it exports alert outcomes in the lifecycle report.

The result is a monitoring system that can be audited: for each episode, we can see what changed, which indicators moved, which alerts triggered, how the drift was classified, what intervention was attempted, and whether stage gates passed. This is the infrastructure needed for lifecycle governance.

With the monitoring panel defined, the next section formalizes drift detection and classification. We define how to construct feature vectors from monitoring deltas, how to score different drift classes with multi-signal rules, how to handle unknown outcomes, and how to calibrate the classifier to avoid false certainty. This prepares the ground for the intervention policy that follows.

## 3.6 Drift Detection and Classification

### 3.6.1 Feature construction from monitoring deltas

Drift detection begins with an apparently simple question: what changed relative to baseline? The answer is not a single number; it is a vector of changes across performance, symmetry, geometry, and probes. We therefore define a feature construction step that maps raw monitoring outputs into a consistent feature vector per episode. The purpose of this step is twofold. First, it standardizes evidence so that classification rules are transparent and reproducible. Second, it separates measurement from interpretation: the feature vector is a factual record of what moved.

Let the baseline episode be $e = 0$ and later episodes be $e = 1, \ldots, E$. For each episode $e$, we compute a monitoring panel $M_e$. From that panel, we construct deltas relative to baseline:

$$\Delta M_e = M_e - M_0,$$

where subtraction is understood component-wise for scalar metrics, and via standardized summaries for distributions.

The feature vector includes at least the following groups:

**Performance deltas.** $\Delta R@1_{i \to t}$, $\Delta R@5_{i \to t}$, $\Delta R@1_{t \to i}$, $\Delta R@5_{t \to i}$, and symmetry gap deltas $\Delta \Delta_1$, $\Delta \Delta_5$. We also include the absolute gap values because the direction of the gap matters for diagnosing dominance.

**Scale and symmetry deltas.** Changes in pre-normalization norm statistics per modality, such as $\Delta$mean norm and $\Delta$norm variance. We include a norm divergence summary, e.g., the absolute difference in means divided by pooled standard deviation. We also include changes in effective rank per modality and the rank divergence.

**Collapse deltas.** Changes in mean off-diagonal cosine similarity per modality, changes in variance floors (e.g., 5th percentile of coordinate variance), changes in hubness summaries (e.g., fraction of queries mapping to top-1 target), and changes in spectrum concentration (e.g., top eigenvalue share).

**Probe deltas.** Changes in MI proxy with each factor, changes in MI proxy with confounder (if present), changes in CCA-like principal correlations (e.g., mean of top-$k$ correlations), and changes in sensitivity proxy.

**Stability deltas.** Changes in gradient norm statistics if fine-tuning occurs; changes in NaN counts (should be zero); changes in any health gate flags.

We explicitly include both raw deltas and normalized deltas. Normalization is needed because different metrics have different scales. For example, retrieval metrics are in $[0, 1]$, while effective rank is in $[1, d]$. Normalized deltas can be computed by dividing by baseline standard deviations estimated via bootstrap or via seed repeats. In a synthetic notebook, we may approximate this using a small number of baseline repeats or by using within-episode variability estimates. The goal is not statistical perfection; the goal is a stable feature representation.

Crucially, feature construction must handle cases where some metrics are undefined. For example, MI with confounder is undefined if no confounder exists. In that case, we record a sentinel value and set a "confounder$_p resent_J flag. This prevents accidental misuse of missing metrics in classification.$

Feature construction is logged as part of the audit bundle. The features are saved per episode. This is governance discipline: if a reviewer disagrees with the classifier, they can inspect the features and propose alternative rules. Without explicit features, classification becomes a black-box narrative.

### 3.6.2 Multi-signal scoring and minimum-evidence rules

Classification is the step where we map signatures to drift labels. Because this mapping is uncertain, we use multi-signal scoring rather than single thresholds. The goal is to reduce false certainty and

to make reasoning reviewable.

We define a set of candidate drift classes:

$$\mathcal{C} = \{\text{benign}, \text{dominance}, \text{spurious\_alignment}, \text{pairing\_corruption}, \text{collapse}, \text{scoring\_drift}, \text{unknown}\}.$$

Depending on the pedagogical design, scoring drift can be treated as its own class or as a subtype that leads to interventions. The key is that the classifier must be explicit about what it recognizes.

For each class $c \in \mathcal{C}$, we define a score function $S_c(\Delta M_e)$, implemented as a weighted sum of evidence terms:

$$S_c = \sum_k w_{c,k}\, \psi_{c,k}(\Delta M_e),$$

where $\psi_{c,k}$ are indicator-like or smoothly increasing functions of features (e.g., ReLU of a thresholded delta), and $w_{c,k}$ are transparent weights. This scoring scheme is intentionally simple. The point is interpretability, not maximal accuracy.

Examples of evidence terms:

**Dominance evidence.** Large increase in symmetry gaps; large divergence in pre-normalization norms; divergence in effective rank; gradient norm imbalance if training occurs. Dominance score increases when multiple of these signals co-occur.

**Spurious alignment evidence.** Retrieval improves or remains high while MI with intended factors declines; MI with confounder increases; CCA correlations concentrate in few dimensions; ablation tests cause sharp performance drop (if performed). Spurious alignment requires confounder presence; otherwise the score is suppressed.

**Pairing corruption evidence.** Retrieval becomes unstable and directionally inconsistent; sensitivity proxy spikes; CCA correlations drop; hubness increases in erratic ways; corruption audits (e.g., similarity consistency checks) fail. Pairing corruption often produces noisy signatures rather than smooth trends.

**Collapse evidence.** Mean cosine increases; variance floors drop; effective rank drops; hubness rises sharply; retrieval may remain stable initially but becomes brittle under mild perturbation. Collapse score requires co-occurrence of multiple geometry indicators.

**Scoring drift evidence.** Changes in similarity scale or effective temperature; shifts in hubness and retrieval sharpness without significant within-modality spectral changes; cross-modal ranking becomes sharper or flatter uniformly. Scoring drift is often detected by comparing distributions of logits and by checking whether embedding distributions themselves changed.

**Benign shift evidence.** Minor changes in retrieval within budget; geometry indicators stable; probes stable; no alerts triggered. Benign shift score increases when changes are small and consistent.

After computing scores, we apply minimum-evidence rules. For example, we do not label spurious

alignment unless confounder MI evidence is present and factor MI evidence shows decline, or unless ablation confirms shortcut reliance. We do not label collapse unless both mean cosine and effective rank show significant movement. We do not label dominance unless symmetry gaps co-occur with either scale divergence or rank divergence. These minimum-evidence rules are the difference between a rigorous classifier and a fragile one. They prevent the system from labeling based on a single noisy metric.

The final label is chosen by maximum score subject to minimum-evidence constraints. If no class meets its minimum evidence, the output is "unknown." This is not an error; it is correct behavior under uncertainty.

### 3.6.3   Unknown-class handling: debug dumps and escalation

The most important governance feature of the drift classifier is the explicit "unknown" state. In many systems, unknown outcomes are treated as failures to be suppressed. In a governance-first system, unknown outcomes are treated as escalation triggers.

When an episode is classified as unknown, the system must do three things:

**(1) Export a debug artifact.** The debug artifact includes the full feature vector, baseline comparisons, raw monitoring outputs, and the per-class scores and which minimum-evidence rules failed. This allows a human reviewer to understand why classification failed and to adjust rules or run additional tests.

**(2) Continue lifecycle evaluation without forcing a label.** The system should still record metrics and monitor trends. Unknown classification should not stop monitoring. However, intervention choices should be conservative: avoid aggressive fine-tuning; prefer diagnostic interventions such as additional audits or counterfactual tests. In a synthetic lab, we can continue to demonstrate later episodes; in production, unknown classification would likely trigger a triage process.

**(3) Record open items and questions to verify.** The report should list hypotheses: e.g., is this a mixture of drift types? Are the thresholds miscalibrated? Did a new failure mode occur? Are probes unstable due to small sample size? This preserves epistemic discipline.

The unknown-handling mechanism prevents a common failure: forcing a plausible-sounding narrative. In lifecycle governance, false narratives are dangerous because they lead to wrong interventions. The correct posture is: when evidence is insufficient, escalate.

In the synthetic lab, unknown classifications are also pedagogically useful. They teach students that drift is not always cleanly classifiable. Real systems exhibit mixed drift: e.g., modality drift plus confounder drift plus scoring drift. Unknown classifications motivate students to run additional tests and to design better monitors.

### 3.6.4 False positives/negatives: calibration and error costs

No drift classifier is perfect. Therefore, we must consider the cost structure of errors. In governance, false positives and false negatives have different costs. A false positive triggers unnecessary intervention, which can be costly and risky if it involves fine-tuning or pipeline changes. A false negative fails to detect harmful drift, which can lead to silent failures and downstream harm.

The appropriate calibration depends on the environment. In high-stakes settings, false negatives may be more costly because they allow harmful behavior to persist. In low-stakes settings, false positives may be more costly because they cause operational churn.

Chapter 3 adopts a conservative principle: when uncertain, do not take aggressive actions. This implies that the classifier should be willing to output unknown rather than forcing a label, and that interventions should be bounded and reversible. This reduces the cost of false positives by limiting intervention severity, and reduces the cost of false negatives by ensuring monitoring continues and by triggering escalation.

Calibration can be performed by evaluating the classifier on synthetic episodes where the true drift type is known. Because the lab controls drift, we can compute confusion patterns: when does the classifier confuse dominance with modality drift, or spurious alignment with scoring drift? These confusion patterns teach students that signatures can overlap. In production, we do not have ground truth, but we can still calibrate by running controlled perturbation tests and observing classifier behavior.

### 3.6.5 Evaluation: classification stability across seeds and regimes

A final requirement for a useful drift classifier is stability. If classification flips across random seeds or small perturbations, it is not governable. Therefore, Chapter 3 evaluates classification stability by repeating key drift episodes across multiple seeds. We then measure how often the classifier outputs the same label and how the scores vary.

Stability evaluation is itself a governance artifact. It informs whether the classifier is reliable enough to be used for automated triggers or whether it should only be advisory. In many systems, the correct decision is to treat drift classification as a human-in-the-loop tool rather than as an automated controller.

We also evaluate classification stability across regimes: for example, under different baseline temperatures or different noise levels. This ensures the classifier is not overfit to one baseline. In a synthetic lab, we can implement this with modest additional computation: run a smaller subset of episodes under alternative baselines and compare classification outcomes.

The output of this evaluation is a report that includes classification tables per episode, per seed, and summary stability metrics. In the notebook, this report is exported in strict JSON. This is

part of the lifecycle audit bundle.

With drift detection and classification defined, we are ready to specify the intervention policy. The next section describes a bounded action space: what interventions are permitted, how they are triggered by classification and evidence, how stage gates evaluate intervention success, and how rollback discipline prevents intervention-induced failure spirals. This is where the chapter moves from diagnosis to control under governance constraints.

## 3.7  Intervention Policy Under Governance Constraints

### 3.7.1  Permitted interventions and bounded action space

A lifecycle framework that detects drift but cannot respond is incomplete. Yet an intervention framework that responds without constraints is dangerous. The purpose of this section is to formalize a bounded intervention policy: a set of allowed actions, a set of prohibited actions, and a disciplined decision procedure that ties actions to evidence and evaluates outcomes via stage gates.

We define the intervention policy as a controller operating on the system state described by monitoring outputs and drift classification. Given an episode $e$, we observe features $\Delta M_e$, produce a classification label $\hat{c}_e$, and select an intervention action $a_e$ from an allowed set $\mathcal{A}$. The policy is bounded in three ways:

**(1) Scope bounds.** Interventions are restricted to actions that can be executed and reversed within the notebook's controlled setting, and analogously in production to actions that can be audited and rolled back. We do not permit "full retraining" as an immediate response; that is an escalation outcome, not an automatic intervention.

**(2) Intensity bounds.** Interventions that modify model parameters (fine-tuning) are limited to a small number of steps $K$ and must respect stop rules. The goal is to demonstrate controlled correction, not unconstrained optimization.

**(3) Evidence bounds.** Actions require minimum evidence. For example, we do not apply confounder ablation unless confounder indicators are strong. We do not apply pairing sanitation unless pairing drift indicators are present. This reduces the risk of applying the wrong intervention.

The allowed intervention set typically includes:

- **Temperature recalibration:** adjust scoring scale to restore baseline effective temperature.
- **Input renormalization recalibration:** adjust modality feature scaling or normalization constants to correct preprocessing drift.
- **Loss reweighting or gradient balancing:** change the relative contribution of modalities in the symmetric loss to counter dominance.
- **Pairing sanitation filter:** detect and drop suspected corrupted pairs prior to fine-tuning, using

similarity-consistency checks.

- **Confounder ablation and retraining:** remove confounder channel or apply counterfactual tests; fine-tune on confounder-free data.
- **Limited fine-tuning with stop rules:** apply a small number of updates under strict health gates and rollback discipline.

We also define explicitly prohibited actions in the pedagogical posture: unconstrained hyperparameter search, threshold changes after observing results, and silent pipeline changes without manifest updates. These prohibitions are governance lessons. In production, many failures come from uncontrolled changes. The synthetic lab instills discipline by preventing those behaviors.

The policy outputs not only an action but a decision record: a structured justification referencing measured facts. Decision records are exported to the audit bundle. This transforms interventions from ad hoc fixes into reviewable governance events.

### 3.7.2 Temperature recalibration and scoring normalization

Temperature recalibration is the archetypal example of a bounded, low-risk intervention. It changes how similarity scores are interpreted without changing embeddings. In systems where scoring drift is detected—e.g., similarity distributions become systematically sharper or flatter relative to baseline—temperature recalibration can restore the operational regime.

In formal terms, if logits are computed as $L = \hat{Z}_{\text{img}}\hat{Z}_{\text{txt}}^{\top}/\tau$, then changing $\tau$ rescales logits. In a retrieval system, this changes the softmax distribution used for probabilistic ranking or training, but even in pure nearest-neighbor retrieval, scaling can affect tie-breaking under approximate search or under thresholding policies. Therefore, temperature is a policy parameter.

The decision rule for recalibration is evidence-based. We look for signatures consistent with scoring drift: large changes in logit scale or similarity distribution without corresponding changes in within-modality spectra, and changes in hubness patterns consistent with sharpness shift. When detected, we compute a recalibration factor by matching a baseline statistic, such as matching the standard deviation of logits or matching a quantile of similarity. We then re-evaluate retrieval and health indicators under the recalibrated temperature and apply stage gates.

The governance lesson is that not all interventions require retraining. Some drift is in the scoring layer. Fixing it requires calibration, not optimization. This is important for production, where retraining is costly and risky. A disciplined system will attempt low-risk calibrations before high-risk fine-tuning.

### 3.7.3 Balancing interventions: gradient norm and loss reweighting

Dominance and modality imbalance often require interventions that affect learning dynamics. A common symptom is that one modality's encoder receives stronger gradients and begins to dominate the shared space. Another symptom is that one modality becomes cleaner, causing the optimizer to rely on it, leaving the other modality under-trained. To address this, we use balancing interventions.

Balancing can be done by reweighting the loss terms. In symmetric InfoNCE, we can assign weights:

$$\ell = \alpha \ell_{i \to t} + (1 - \alpha) \ell_{t \to i},$$

or more generally, we can weight gradients per modality by scaling the backpropagated signals. We choose a rule based on measured gradient norm imbalance. For example, if $\|\nabla_\theta \ell\|$ for the image encoder is consistently larger than $\|\nabla_\phi \ell\|$ for the text encoder, we adjust weights so that effective update magnitudes become comparable. This is not a guarantee of correctness, but it is a bounded control.

Balancing can also include renormalization of embedding outputs before normalization, but because we apply L2 normalization, this is more subtle. Pre-normalization scale affects gradients through the normalization Jacobian. Therefore, balancing can involve adjusting internal scales or applying gradient clipping per modality.

In the synthetic lab, balancing is implemented as a controlled, transparent procedure: compute gradient norms per modality on a validation batch, compute a balancing factor, apply it for a bounded number of steps, and re-evaluate. The decision record references the measured imbalance and the balancing factor used. Stage gates ensure that balancing does not fix dominance at the expense of collapse or spurious alignment.

The broader lesson is that multimodal training is not symmetric by default, even if the loss is symmetric. Dynamics can break symmetry. Therefore, production systems need explicit symmetry controls.

### 3.7.4 Data interventions: sanitation filters and confounder ablations

Some drift types are data problems, not model problems. Pairing drift and confounder drift fall into this category. Interventions must therefore operate on data, not only on training.

**Pairing sanitation.** When pairing drift is suspected, the first principle is: do not fine-tune on corrupted pairs. Instead, we attempt to detect and filter corrupted pairs. In a synthetic lab, we can validate against known corruption. In production, we rely on consistency tests. One such test is similarity-consistency: compute the similarity of each purported pair and compare it to similarity distribution of nearby neighbors. Pairs with unusually low similarity relative to their neighborhood are suspicious. Another test is mutual nearest neighbor consistency: if image $i$'s nearest text is $j$

and text $j$'s nearest image is not $i$, the pair might be inconsistent. These tests are imperfect, but they provide bounded data sanitation.

After filtering, we re-evaluate monitoring metrics on the sanitized set. If metrics improve and instability decreases, that supports the hypothesis of pairing corruption. The decision record should explicitly state that sanitation is a diagnostic and mitigation step, not a proof. Open items should include pipeline validation requirements.

**Confounder ablation.** When confounder drift is suspected, we conduct ablation tests: remove the suspected confounder channel from both modalities and re-evaluate retrieval and probes. In a synthetic lab, confounders are known. In production, suspected confounders might be metadata fields or formatting artifacts. If performance collapses after ablation, that is evidence of shortcut reliance. If factor MI or downstream task performance improves after ablation, that is evidence that shortcuts were harmful.

Ablation is then followed by controlled retraining or fine-tuning on confounder-free data. This is a higher-risk intervention than temperature calibration, so it must be bounded and governed. Stage gates should include not only retrieval but also probe improvements: MI with intended factors should increase relative to confounder MI. If not, the intervention has not achieved its semantic goal.

The key governance principle is: treat data interventions as first-class. Many multimodal failures are data failures. A model-only response can worsen the system by encoding the wrong coupling.

### 3.7.5 Limited fine-tuning with stop rules and rollback discipline

Fine-tuning is the most powerful and most dangerous intervention. It can correct drift by adapting encoders to new distributions, but it can also induce collapse, amplify confounders, and create irreproducible changes. Therefore, Chapter 3 constrains fine-tuning tightly.

We define limited fine-tuning as $K$ steps of gradient-based updates under a fixed learning rate and temperature regime, with strict monitoring at each step. Stop rules are applied immediately if health gates are violated:

- mean cosine exceeds ceiling;
- variance floors drop below floor;
- effective rank drops below floor;
- hubness exceeds ceiling;
- NaNs appear;
- sensitivity proxy exceeds ceiling (indicating instability).

If a stop rule triggers, fine-tuning is aborted and the system rolls back to the last healthy checkpoint. Rollback is not optional. It is the control that prevents intervention-induced drift spirals.

Even when fine-tuning completes $K$ steps, acceptance is not assumed.  We apply stage gates post-intervention. Success requires improvement or recovery in operational metrics *without* violating symmetry and health constraints.  In some cases, fine-tuning can improve retrieval but worsen geometry. In a governance-first posture, that is not success. The system must remain governable.

The decision record for fine-tuning includes: justification (which drift label, which evidence), hyperparameters, stop rules, and outcome. Open items include whether the fine-tuning regime is appropriate for production, whether additional data validation is needed, and whether monitoring thresholds should be revised based on larger-scale evidence. The verification status remains "Not verified."

This constrained fine-tuning framework is pedagogically important. It teaches students that training is a risky intervention, not a default solution. In many organizations, the reflex is "fine-tune it." Chapter 3 replaces that reflex with a controlled policy: prefer low-risk calibrations and data fixes; fine-tune only when evidence supports it; bound fine-tuning; monitor continuously; and roll back if health gates fail.

With intervention policy specified, we are ready to describe the synthetic drift laboratory timeline: how drift episodes are constructed, how orthogonality and interactions are handled, how counterfactual replay is used to isolate causes, and how robustness margins are computed. The next section defines the episode design and the artifact export standards that make lifecycle governance reviewable.

## 3.8   Synthetic Drift Laboratory: Episode Timeline Design

### 3.8.1   Episode construction: 12+ controlled drift events

A lifecycle framework must be tested over time, not only at a single point. The synthetic drift laboratory therefore constructs an explicit timeline of drift episodes. Each episode is a time-indexed configuration of the data-generation pipeline and the evaluation and intervention regime. The intent is to model the kinds of change events that happen in production: new data sources, preprocessing updates, scoring changes, join bugs, and fine-tuning actions. The synthetic nature of the lab allows us to prescribe these events precisely and to know the intended drift type. That ground truth supports pedagogical evaluation of the monitoring and classification framework.

An episode timeline is defined as $\{e = 0, 1, \ldots, E\}$ with $E \geq 12$. Episode $e = 0$ is the baseline. Episodes $e \geq 1$ each introduce one primary drift axis and may include a minor stochastic perturbation. The "one primary axis" rule is deliberate. It allows the student to see clean signatures. In a later variant, mixed drift can be introduced, but the default timeline begins with separable events to teach causal interpretation.

A canonical timeline includes:

- **Modality drift (image)**: progressively increase image noise or blur.
- **Modality drift (text)**: introduce token corruption or positional perturbations.
- **Scoring drift**: rescale similarity scores to simulate temperature change.
- **Confounder introduction**: add a shared shortcut feature.
- **Confounder removal**: remove the shortcut feature.
- **Pairing drift (random)**: swap a fraction of pairs.
- **Pairing drift (systematic)**: apply a fixed offset to create structured join bugs.
- **Preprocessing drift**: change modality feature scaling to simulate normalization changes.
- **Optimization drift**: attempt limited fine-tuning under a risky regime to expose collapse boundaries.
- **Recovery attempt**: apply bounded interventions to restore health after a drift event.

The timeline is designed to include both gradual drifts and step-change events, because real systems experience both. Gradual drifts teach monitoring trends and budgets. Step-change events teach alerting and incident response.

Each episode generates an evaluation dataset using the episode-specific measurement operators. We distinguish between a stable *reference evaluation set* and an episode-specific *current evaluation set*. The reference set is useful for separating model drift from data drift: if the model weights do not change, performance on the reference set should remain stable. If it changes, that indicates that scoring or model changed. The current set captures production reality: the system must operate on current inputs.

This dual-evaluation approach is a core governance lesson. Many teams evaluate only on current data, confounding data drift with model drift. In a governed system, we keep a stable canary set and a current set, and we evaluate on both.

### 3.8.2 Single-axis vs mixed drift: orthogonality and interactions

The "single-axis" design is pedagogically clean, but production drift is rarely orthogonal. Modality drift can coincide with scoring drift; confounders can emerge while pairing integrity degrades; fine-tuning can be triggered precisely when inputs shift. Therefore, the laboratory includes a second layer of episodes that introduce interactions after the student has learned the basic signatures.

We distinguish:

**Single-axis episodes.** One primary drift axis changes, others remain at baseline settings. This supports clean signature learning and classifier calibration.

**Mixed episodes.** Two drift axes change simultaneously, often in a realistic pairing: e.g., modality noise increases and scoring sharpness increases, creating a brittle system; or confounder is introduced while pairing corruption increases, creating contradictory supervision.

Mixed episodes are treated as advanced cases. They test whether the classifier can output "unknown" appropriately and whether intervention policies remain conservative. We do not demand perfect classification under mixed drift; we demand reviewable behavior: the system should recognize that evidence does not match a clean class and should export debug artifacts.

To ensure that mixed episodes are informative, we control interaction magnitudes. For example, if both modality noise and pairing corruption are large, everything fails and interpretation is trivial. Instead, we create moderate interactions where headline retrieval might remain acceptable while health indicators degrade. This teaches the "silent drift" concept: failure can be geometric and mechanistic before it becomes operational.

Orthogonality is also tested via counterfactual replays. If an episode is mixed, we rerun it with one drift axis removed to see which indicators persist. This is a causal probing posture: not "fit a narrative," but "test hypotheses by controlled variation." The notebook implements this by generating alternative episode configs and re-evaluating.

### 3.8.3   Counterfactual replay: removing suspected drift causes

Counterfactual replay is a key method for moving from signatures to mechanisms. If monitoring suggests spurious alignment, we can remove the suspected confounder and re-evaluate. If monitoring suggests scoring drift, we can re-evaluate under baseline temperature. If monitoring suggests modality drift, we can restore preprocessing and see if signatures disappear. If monitoring suggests pairing drift, we can sanitize pairs and see whether stability improves.

In a synthetic lab, counterfactual replay is straightforward because drift causes are explicit. In production, replay can be harder because causes are hidden. But the method still applies: isolate suspected causes by modifying the pipeline or by filtering data and measuring whether signatures change. This is the essence of governed diagnosis: do not guess; test.

Counterfactual replay is implemented as part of the intervention loop. Some interventions are inherently counterfactual tests: confounder ablation, temperature recalibration, sanitation filtering. The notebook records both the test and the outcome. This avoids a common production failure: applying an intervention and attributing improvement to the intervention without testing alternatives. In the lab, we can show that multiple changes could explain improvement, and we train students to record open items rather than claiming certainty.

### 3.8.4   Robustness margins: drift budgets and safe regions

A lifecycle system cannot respond to every change. It must define drift budgets and safe operating regions. The synthetic lab therefore includes budget estimation. For each drift axis, we run a small sweep of drift magnitude and measure how monitoring indicators change. This produces a robustness curve: for example, retrieval@1 vs image noise, effective rank vs temperature scale, MI

with factors vs token corruption.

From these curves, we define safe regions: ranges of drift magnitude within which stage gates still pass. The safe region is multi-dimensional, but in the lab we typically define axis-aligned budgets for simplicity: e.g., image noise $\leq \sigma^\star$, pairing corruption $\leq p^\star$, temperature scale within $[c_{\min}, c_{\max}]$. These budgets are then used in the episode timeline. Episodes that remain within budgets are expected to be classified as benign or low-risk. Episodes that exceed budgets should trigger alerts and potentially interventions.

The pedagogical value of budgets is significant. Students learn that "robustness" is not a moral claim; it is an empirically measured margin. They learn that budgets are chosen intentionally and recorded. They also learn that budgets can be domain-specific: a system used for casual retrieval may tolerate larger drift than a system used for compliance.

In production, budgets correspond to service-level objectives and risk tolerances. The lab's budgeting method teaches how to operationalize these concepts: define metrics, sweep perturbations, measure margins, set thresholds, and record them.

### 3.8.5 Artifact export: per-episode tables and lifecycle reports

The synthetic drift lab produces a large amount of information. Governance requires that this information be exported in structured, reviewable form. Therefore, each episode generates:

- **Monitoring snapshot:** all panel metrics, saved as JSON.
- **Feature vector:** constructed deltas and normalized deltas, saved as JSON.
- **Classification record:** per-class scores, label, and evidence flags, saved as JSON.
- **Decision record:** selected intervention (if any), justification, and stage-gate outcomes, saved as JSON.
- **Plots:** trend plots across episodes, spectrum plots, hubness plots, and retrieval curves.

These per-episode artifacts are aggregated into a lifecycle report. The lifecycle report is a table-like structure: each row is an episode, and columns include drift configuration, key metrics, classification label, intervention action, and success/failure. In the notebook, this table is stored as JSON with explicit keys to avoid ambiguous parsing. The report also includes open items and questions to verify, especially for unknown classifications.

Finally, the entire set of artifacts is packaged into an audit bundle with a run manifest and risk log. This packaging is not cosmetic. It is the mechanism that turns a notebook run into a reviewable experiment. A reviewer can inspect the bundle and reconstruct what happened. The bundle also supports comparison across runs: because the manifest records configurations, we can compare different policy choices or different thresholds.

With the episode timeline design established, we are ready to interpret results and translate them

into production practice. The next section explains how to read drift signatures, which monitors lead and which lag, how to convert monitoring into decision records, and how to map synthetic drift types to real production changes. This is the section where the laboratory becomes a professional lesson rather than a toy demonstration.

## 3.9   Results Interpretation and Production Translation

### 3.9.1   How signatures evolve before failures become visible

A central lesson of lifecycle governance is that failure is often preceded by measurable degradation that does not immediately show up in headline metrics. Students frequently expect failure to look like a sudden drop in retrieval@1. In practice, a multimodal system can drift into a fragile regime while maintaining acceptable retrieval on familiar data. This is why the monitoring panel includes indicators that lead performance.

In the synthetic drift laboratory, we observe several common pre-failure patterns.

First, **geometry degrades before retrieval**. Mean off-diagonal cosine can rise gradually, indicating increasing anisotropy. Effective rank can fall, indicating that the representation is concentrating variance into fewer directions. Hubness can increase, meaning that a small number of items become nearest neighbors for many queries. These changes can occur even when retrieval remains stable because the dataset is still compatible with the learned coupling. However, once drift pushes the system slightly further—e.g., more modality noise, more scoring sharpness, or mild pairing corruption—the brittle geometry causes rapid performance degradation. This is the "cliff" behavior: slow health degradation followed by sudden operational failure.

Second, **symmetry gaps often lead**. Directional retrieval can diverge as one modality becomes noisier or as dominance emerges. A common signature is: one direction remains stable while the other declines steadily. If one reports only an average, the decline is masked. Symmetry gaps therefore function as early warning indicators, particularly for modality drift and dominance.

Third, **probe shifts can reveal silent semantic failure**. In confounder drift episodes, retrieval can improve when a shortcut signal is introduced. However, MI with intended factors declines or stagnates, and MI with the confounder rises. The system is "working" operationally on the evaluation set, but it is not encoding the intended semantics. When the confounder is removed, performance collapses. The probe shift occurs earlier than the collapse; it is a leading indicator.

Fourth, **sensitivity spikes signal instability**. Pairing corruption and optimization drift often manifest as increased sensitivity: small perturbations in one update step can change retrieval substantially. This is a sign that the system is near an unstable boundary. In production, this means that incremental fine-tuning can become dangerous precisely when drift is present. Sensitivity therefore informs intervention choices: when sensitivity is high, avoid aggressive fine-tuning.

These pre-failure signatures teach a practical interpretive posture: do not wait for failure to become obvious. Monitor geometry, symmetry, and probes. Treat rising collapse indicators and probe shifts as warnings that the system's representational policy is changing.

### 3.9.2 Which monitors lead: early warnings vs lagging indicators

Not all monitors are equally early. Some are leading indicators; some are lagging indicators. A professional monitoring system distinguishes between them.

**Leading indicators** include: effective rank, variance floors, mean cosine, hubness, symmetry gap trends, MI proxy shifts, and sensitivity proxy spikes. These indicators often move before retrieval drops. They capture structural health and mechanism changes.

**Lagging indicators** include: headline retrieval@ 1 and retrieval@ 5 on a fixed evaluation slice. These indicators often move after the system has already drifted into a fragile regime. They are still important because they reflect user-visible performance, but they are late.

Spectral measures occupy a mixed role. The eigenvalue spectrum can move early, but interpreting it requires context: some spectrum shifts are benign if data becomes less diverse. Therefore, spectra are early but ambiguous; they become more useful when combined with other signals.

The chapter's governance logic uses this distinction explicitly. Drift detection does not rely solely on lagging indicators. It triggers classification when leading indicators cross budgets, even if retrieval remains acceptable. Intervention decisions also consider lead/lag structure. For example, if retrieval declines but geometry remains healthy, a low-risk recalibration may be sufficient. If retrieval is stable but geometry is degrading, a preemptive intervention or an escalation may be warranted. This prevents the system from "chasing" metrics after the fact.

A key pedagogical objective is to teach students to treat monitoring as a panel, not a scoreboard. A scoreboard invites optimization and persuasion. A panel invites diagnosis and control.

### 3.9.3 Decision records: turning measurements into actions

Measurements do not govern systems. Decisions do. A lifecycle system must translate monitoring into explicit actions, and those actions must be justified in a way that can be reviewed.

A decision record includes:

- **Observed facts:** key monitoring changes (e.g., symmetry gap increased, effective rank fell, MI confounder rose).
- **Classification hypothesis:** which drift class is most supported and why; or "unknown."
- **Action selection:** the chosen intervention from the allowed set, or a decision to take no action.
- **Expected effect:** what metrics should improve if the hypothesis is correct.

- **Stage gates:** acceptance criteria post-action.
- **Rollback plan:** what happens if gates fail.

The decision record enforces a causal posture. It does not merely say "we changed temperature." It says "we detected scoring drift because similarity scale shifted while embedding spectra remained stable; we recalibrated temperature to match baseline logit variance; we expected hubness to decrease and retrieval symmetry to improve; we accepted the change only if health gates were maintained." This is reviewable.

Decision records also enforce honesty about uncertainty. If the classifier is unknown, the record should say: "Evidence is insufficient to identify drift mechanism; we performed a conservative diagnostic intervention (e.g., additional audits, counterfactual evaluation) and escalated for pipeline review." This avoids the trap of pretending that every event is understood.

In the synthetic lab, decision records are part of the lifecycle report and audit bundle. In production, decision records become incident documents and change-control artifacts. Teaching students to write them as part of the technical workflow is a governance-first pedagogical move.

### 3.9.4  Production mapping: what corresponds to each synthetic drift

A synthetic lab is useful only if students can map it to real scenarios. We therefore explicitly map each synthetic drift type to common production events:

**Modality drift** corresponds to: new device types, new compression standards, changes in image preprocessing, OCR degradation, language style shifts, new annotation sources, or changes in tokenization/normalization.

**Scoring drift** corresponds to: downstream service rescaling scores, changes in approximate nearest neighbor indexing, quantization updates, changes in normalization conventions, or changes in ranking thresholds.

**Confounder drift** corresponds to: IDs leaking into both modalities, watermarks, template formatting, metadata injection, dataset source markers, or new logging fields that become correlated with pairing.

**Pairing drift** corresponds to: join key changes, timestamp misalignment, deduplication changes, batch reordering, partial ingestion failures, or ETL bugs that swap records.

**Optimization drift** corresponds to: unreviewed fine-tuning, hyperparameter changes, changes in batch composition and negative sampling, new training data with different confounder structure, or changes in regularization that alter geometry.

These mappings serve a pedagogical function: they show that drift is not an exotic concept. It is the default state of production systems. The lab trains students to anticipate these events, instrument monitors, and respond with bounded controls.

Importantly, the lab also teaches what cannot be mapped directly. MI proxy with true factors is possible in synthetic settings because we know factors. In production, we do not. Therefore, production analogs require proxy labels or curated evaluation sets. This becomes an open item in the report: "How will we approximate factor MI in production? Do we have labeled canary sets? Do we have counterfactual tests?" The notebook's governance artifacts should explicitly include these questions.

### 3.9.5   Minimum control set for real deployments

The chapter culminates in a minimum control set: the smallest set of controls that a professional multimodal deployment should have if it is to be governable under drift. The minimum set is not maximal; it is practical.

A reasonable minimum set includes:

**(1) Versioned manifests.** Every model, pipeline, and scoring configuration must be versioned, with manifests capturing parameters, code hashes, and environment metadata.

**(2) Dual evaluation sets.** Maintain a stable canary set and a current set. Evaluate on both to separate model changes from data/pipeline drift.

**(3) Monitoring panel.** At minimum: directional retrieval@$k$, symmetry gaps, mean cosine, effective rank, hubness, and a sensitivity proxy (or a simpler stability proxy if sensitivity is expensive). These provide lead/lag coverage.

**(4) Shortcut audits.** In production, this means explicit checks for metadata leakage and counterfactual ablations where feasible. Without shortcut audits, spurious alignment will go undetected.

**(5) Stage gates and rollback discipline.** No fine-tuning or scoring changes should be deployed without passing gates on both performance and health metrics. Rollback must be feasible and practiced.

**(6) Decision records and escalation rules.** Unknown cases must trigger escalation; interventions must be documented; thresholds must be pre-specified.

This minimum control set is the practical translation of the synthetic lab. It does not claim to solve all problems. It claims to make systems reviewable and controllable under drift. That is the professional objective.

With interpretation and production translation established, the chapter is prepared to conclude by summarizing the lifecycle governance lesson, emphasizing the shift from point-evaluation to continuous monitoring and bounded intervention, and clarifying limits and the roadmap beyond this chapter. The conclusion will also position the synthetic lab as a pedagogical instrument: it trains the correct mental model and operational habits, while leaving verification to real-world validation and human accountability.

## 3.10 Conclusion

This chapter's core message is structural: multimodal systems do not merely *learn* a shared space, they *operate* a shared space as an internal policy surface. Once images and text are embedded into vectors, the downstream system becomes geometric. Retrieval, ranking, clustering, deduplication, filtering, and even safety gates are all functions of distances, angles, and neighborhoods. That geometry is not static. It changes as the world feeding the system changes, as pipelines evolve, as scoring is recalibrated, and as fine-tuning continues. Therefore, the correct professional framing is not "we trained a multimodal model," but "we are operating a representational policy under change." Lifecycle governance is what turns that policy into something reviewable.

The synthetic drift laboratory exists to make this idea teachable. By construction, it isolates drift causes that are often entangled in production. We treat drift as changes in measurement operators (modality drift), changes in scoring functions (scoring drift), changes in shortcut structure (confounder drift), changes in the coupling relation (pairing drift), and changes induced by our own optimization choices (optimization drift). These are not abstract categories; they correspond to common production events: new device sources, preprocessing updates, downstream service rescaling, metadata leakage, join bugs, and aggressive fine-tuning. The taxonomy gives practitioners a vocabulary that is operational rather than rhetorical. When something goes wrong, the question becomes: which category is most supported by evidence, and what bounded interventions are appropriate?

A second lesson is epistemic: "working at $t_0$" is not evidence of stability. Snapshot evaluation is necessary, but it is not sufficient. The embedding space defines a *relative* neighborhood structure, and neighborhoods change as the distribution of points changes. A model can remain numerically stable while drifting into dominance, spurious alignment, or collapse. Worse, headline metrics can improve while semantics degrade, particularly under confounder drift. This is why lifecycle governance must treat a single metric as a lagging indicator, not as proof. The chapter therefore builds a monitoring panel that includes directional retrieval and symmetry gaps, collapse indicators (mean cosine, variance floors, hubness), spectral measures (eigenvalue spectra, effective rank), and structural probes (MI proxy, CCA proxy, sensitivity proxy). The panel is not a scoreboard; it is an instrument panel. Its purpose is to provide early warnings and mechanism clues before user-visible failures appear.

From monitoring to control, the chapter emphasizes a disciplined decision pipeline. Drift classification is multi-signal and constrained by minimum-evidence rules. The classifier is allowed to output "unknown," and that is treated as a governance feature, not a bug. Unknown outcomes trigger debug artifact exports and escalation rather than forced narratives. This posture prevents one of the most common failures in real deployments: confident misattribution. A system that always "knows" the drift cause is almost always wrong. A governed system represents uncertainty explicitly and uses counterfactual replays and targeted audits to reduce it.

Intervention policy is the third major lesson, and it is where governance becomes real. The bounded action space forces a hierarchy of responses. Low-risk interventions come first: temperature recalibration and scoring normalization for scoring drift; renormalization recalibration for preprocessing drift. Data interventions are treated as first-class: sanitation filters for suspected pairing corruption and confounder ablations for suspected shortcuts. Model interventions exist, but they are constrained: limited fine-tuning with strict stop rules, health gates, and rollback discipline. This matters because many production failures are self-inflicted. Unbounded fine-tuning can push a system across a collapse boundary. Loss can decrease while geometry degrades. The chapter's stop rules—effective rank floors, variance floors, hubness ceilings, and sensitivity ceilings—encode the principle that a model must remain governable, not merely performant on a narrow slice.

The synthetic timeline design reinforces the most important operational habit: evaluate change over episodes and record decisions as artifacts. Every episode produces monitoring snapshots, feature deltas, classification records, decision records, and plots, all packaged into an audit bundle with manifests and risk logs. This is the professional endpoint: reviewability. A reviewer should be able to answer, at any time: what changed, how did it affect the system, what evidence supports our diagnosis, what action did we take, did we pass stage gates, and can we roll back? Without these artifacts, a multimodal deployment becomes a story. With them, it becomes an accountable system.

Two limits must be stated explicitly. First, synthetic evidence is not production verification. The notebook demonstrates mechanisms and failure modes with causal clarity, but it cannot certify real-world robustness. Second, some probes in the synthetic setting rely on knowing true latent factors. In production, factor ground truth is rarely available, so practitioners must build proxy evaluation sets, canary tasks, and counterfactual audits to approximate these probes. The correct governance label therefore remains "Not verified," with open items that specify what must be verified before real deployment claims can be made.

The practical implication is a minimum control set for multimodal deployments. At minimum, organizations should maintain versioned manifests for model, pipeline, and scoring; evaluate on both stable canary sets and current sets; monitor directional retrieval, symmetry gaps, hubness, effective rank, and at least one stability proxy; run shortcut audits and ablations when confounders are plausible; enforce stage gates and rollback discipline for all updates; and produce decision records for interventions and escalations. This minimum set does not eliminate drift. It makes drift governable.

The concluding takeaway for students is deliberately concrete: to use a multimodal embedding model responsibly, you must treat it as a geometric policy that will drift, and you must operate it with the same discipline you would apply to a production control system. The goal is not to win a benchmark. The goal is to keep the system reviewable under change: detect drift early, classify it honestly, intervene conservatively, and leave an artifact trail that can withstand scrutiny. That is the institutional meaning of multimodality when the frontier meets reality.

# Bibliography

[1] H. Hotelling. Relations between two sets of variates. *Biometrika*, 28(3–4):321–377, 1936.

[2] M. Radovanović, A. Nanopoulos, and M. Ivanović. Hubs in space: Popular nearest neighbors in high-dimensional data. *Journal of Machine Learning Research*, 11:2487–2531, 2010.

[3] A. van den Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018.

[4] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020.

[5] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[6] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021.

[7] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y.-H. Sung, Z. Li, and T. Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021.

[8] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2021.

[9] A. Bardes, J. Ponce, and Y. LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations (ICLR)*, 2022.

[10] J. Gama, I. Žliobaitė, A. Bifet, M. Pechenizkiy, and A. Bouchachia. A survey on concept drift adaptation. *ACM Computing Surveys*, 46(4):Article 44, 2014.

# Appendix A

# Companion Colab Notebook Index

Each chapter is paired with a governed Google Colab notebook implementing a synthetic-first laboratory and producing an auditable run bundle.

| Chapter | Notebook (suggested filename) |
| --- | --- |
| Chapter 1 | CH1_MULTIMODAL_GEOMETRY_ALIGNMENT.ipynb |
| Chapter 2 | CH2_MULTIMODAL_FAILURE_MODES.ipynb |
| Chapter 3 | CH3_MULTIMODAL_DRIFT_MONITORING_CONTROL.ipynb |

# Appendix B

# Minimum Governance Standard for All Labs

---

**Artifact (Save This)**

Every companion notebook in this volume must produce:

- Deterministic synthetic generation (seeded) with explicit latent factors.
- Explicit configuration registry (single source of truth) and environment fingerprint.
- Run manifest (`run_id`, timestamp, config hash, library versions).
- Prompts/config log (redacted where needed; hashes always).
- Risk log (capability $\uparrow \Rightarrow$ risk $\uparrow \Rightarrow$ controls $\uparrow$).
- Deliverables folder with plots and strict JSON summaries (facts vs assumptions vs open items).
- Hash ledger for saved artifacts and a zipped audit bundle directory.
- Explicit statement: **Not verified for live deployment; human review required.**

---

# Closing Statement

Multimodality enlarges the surface on which systems can fail. It also enlarges the surface on which systems can be governed.

If you treat multimodality as "more input types," you will get impressive demos and fragile deployments. If you treat it as **shared geometry under drift**, you can build systems that remain reviewable and controllable.

**Governance first. Geometry second. Mechanism always.**